

User Guide to



SeqMan NGen

DNASTAR, Inc. 2021

Table of Contents

Welcome to SeqMan NGen	6
SeqMan NGen Tutorials	7
Whole genome reference-guided workflow	8
Whole genome de novo workflow with mate pair data	18
De novo assembly using Sanger data	23
Analysis of a whole genome de novo assembly	28
RNA-Seq de novo transcriptome workflow	32
RNA-Seq reference-guided workflow with analysis in ArrayStar	38
Part A: Setting up the RNA-Seq reference-guided assembly in SeqMan NGen	39
Part B: Analyzing the results in ArrayStar using quick gene sets	43
Part C: Analyzing the results in ArrayStar using advanced filtering	49
ChIP-Seq workflow with analysis in ArrayStar	52
Copy number variation (CNV) workflow with analysis in ArrayStar and GenVision Pro	56
Part A: Setting up the CNV project in SeqMan NGen	57
Part B: Finding a putative duplication in the reference sequence using ArrayStar	59
Part C: Confirming the duplication using GenVision Pro	62
Whole genome reference-guided workflow with analysis in ArrayStar	64
Part A: Setting up the assembly in SeqMan NGen	65
Part B: Analyzing the results in ArrayStar	72
Long-read analysis with accuracy evaluation	78
Part A: Running the assembly in SeqMan NGen and viewing it in SeqMan Ultra	79
Part B (optional): Evaluating assembly accuracy using QUAST	86
Exome workflow with analysis in ArrayStar	89
Wizard screen descriptions	90
Welcome	91
Workflow	92
De novo genome assembling and editing workflows	93
Create a reference-guided assembly to use in the “SNP to Structure” workflow	96
Remove PhiX control reads from Illumina data prior to import	102
Metagenomics workflows	103
RNA-seq/transcriptomics workflows	104
Include DESeq2 or edgeR statistics	106
Variant analysis/resequencing workflows	110
Variant calling accuracy workflow	112
Variant Call Format (VCF) files workflows	113
Combine/Reanalyze Existing Assemblies	114
Analysis Options	115
RNA-seq normalization methods	119
ChIP-seq peak detection methods	122

Assembly Log	123
Assembly Options	124
Assembly Output.....	128
Assembly Summary	131
Cloud Monitor	133
Define Binding Proteins.....	134
Input Assemblies.....	138
Input Host Files	140
Input Sequences	141
Specify read technology.....	143
Specify paired-end data	144
Example regular expressions.....	148
Specify single sample, multi-sample or replicate data	150
Specify RNA-Seq options.....	154
Input Short Read Sequences	155
Input VCF Files	157
Input Viral Genomes	158
Post Assembly Options	159
Preassembly Options	161
Preassembly Options for long-read workflows.....	162
Preassembly Options for all other workflows	164
Reference Sequence / Input Draft Genome	167
Annotate reference sequences prior to import.....	169
Manually specify an isoform prior to import	170
Use RNA-Seq de novo transcriptome output as a reference	171
Specify a VCF, BED or Manifest file.....	172
Make a custom VCF file.....	173
Make a custom BED file	175
Troubleshoot a Manifest file.....	177
Run Assembly	179
Monitor the progress of a Cloud Assembly	181
Set Contaminant	183
Set Up Experiments	184
Set Up Replicate Sets	185
(Short Read) Polishing Options	186
Transcript Annotation Database	188
Add a DNASTAR transcriptome package	190
Create a custom transcript annotation database	191
Use a local copy of RefSeq as a transcript annotation database	192
Annotation Options dialog	193
Options tabs.....	196
Alignment tab.....	197

Alignment tab (Preassembly Options, long read)	198
Alignment tab (Preassembly Options, all others)	200
Alignment tab (Assembly Options)	201
Layout tab	204
Layout tab (Preassembly Options, long read)	205
Layout tab (Assembly or Analysis Options)	207
Peak Detection tab	209
Scans tab	211
Trimming tab	213
Trimming tab (Preassembly Options, all others)	214
Trimming tab (Assembly Options)	216
Variants tab	217
Filter based on "P not Ref"	221
Log in to Cloud Assemblies	222
Use the DNASTAR Cloud Data Drive	224
License and Credential Requirements	225
The DNASTAR Cloud Data Drive User Interface	226
Access the DNASTAR Cloud Data Drive	229
Create a New Cloud Folder	231
Transfer a Folder from a Physical Computer to the Cloud	232
Transfer Files from a Physical Computer to the Cloud	234
Transfer Files or Folders from the Cloud to a Physical Computer	236
Permanently Remove Files and Folders from the Cloud	238
Close the DNASTAR Cloud Data Drive	239
Navigate between wizard screens	240
Add and remove files in the wizard	241
Add sequences from your computer or the cloud	242
Add a genome template from DNASTAR	244
Add a genome template from NCBI	247
Remove a sequence from the list	250
Use editing commands in the wizard	251
Monitor the progress of a cloud assembly	252
Access and understand output files	253
View the Project Report	254
Project Report contents for reference-guided workflows	255
Project Report contents for de novo workflows	257
Reference-guided workflow output	260
Contents of the .assembly package	261
Contents of the -reports folder	265

Contents of the -zinternal folder	266
De novo workflow output	267
RNA-Seq reference-guided workflow output	268
RNA-Seq de novo transcriptome workflow output	269
Appendix	273
SeqMan NGen calculations	274
Calculation of “match percentage”	275
Detection of structural variations	277
Handling of repeats	279
Handling of sex chromosomes	280
How “mer tags” are chosen	282
Run SeqMan NGen through the command line	284
XNG, SNG, and QNG assemblers	285
XNG commands	286
assembleTemplate	287
computeSNP	305
createGenomeTemplate	308
diskPath	309
dumpConsensus	310
dumpSNP	311
execute	312
exportSplits	313
exportVCF	314
extractPairs	315
include	316
loadAssembly	317
loadBAM	318
mergelonTorrentShortReads	321
message	322
pairFilePattern	323
pause	324
quit	325
removeDuplicateSeqs	326
runScript	327
set	329
setDefaultDirectory	330
setMachineMemory	331
setParam	332
SNG commands	333
Project management commands	334
closeProject	335
runScript	336

saveProject	337
saveReport	339
writeUnassembledSeqs	340
File loading commands	341
load454PairedEnd	342
loadConstraint	343
loadContaminant	344
loadLayout	345
loadRepeat	346
loadSeq	347
loadTemplate	349
loadVector	350
openProject	351
setDefaultDirectory	352
Parameter settings commands	353
setContaminantParam	354
setParam	355
setQualityParam	365
setRepeatParam	367
setVectorParam	369
Preprocessing and assembling commands	371
appendToAssembly	372
assemble	373
convertReads	374
extendContigs	375
fixedTrim	376
include	377
makeSeqNamesUnique	378
realignContigs	379
removeSmallContigs	380
set	381
setAssemblyReport	382
setPairSpecifier	383
splitLinkerReads	384
splitMIDSeqs	385
splitPairs	386
splitTemplates	388
trimVector	389
Specifying XNG or SNG/SMNG when running a script	390
Turn off usage logging	391
Non-English keyboards	392
Installed Lasergene file locations	393


Troubleshoot failure to launch	396
Research references	397

Welcome to SeqMan NGen

Lasergene Genomics provides everything you need for assembly and analysis of genomic, metagenomic, exomes/gene panels and transcriptomic sequencing data, and supports all popular [file formats](#). Most workflows will start with sequence assembly in SeqMan NGen.

SeqMan NGen supports both *de novo* and reference-guided (templated) [workflows](#) on all major sequencing platforms ([Sanger](#), [Illumina](#), [Ion Torrent](#), and [Pacific Biosciences](#)). SeqMan NGen supports reference-guided assemblies of billions of sequence reads and *de novo* assemblies of up to 30 million sequence reads (genome sizes up to 50 megabases).

About this User Guide:

- For help **INSTALLING** Lasergene, see our separate [Installation Guide](#).
- To **PRINT** the current page of the User Guide, click the printer icon in the top right corner ().
- To download a **PDF** of the entire User Guide, scroll to the bottom of the table of contents on the left, and press **Download as PDF**.

✱ **Important:** If you decide to print the PDF, note that the last quarter of the PDF (110+ pages) is a “command-line scripting” manual that applies to only a handful of users. You may wish to omit this when printing.

SeqMan NGen Tutorials

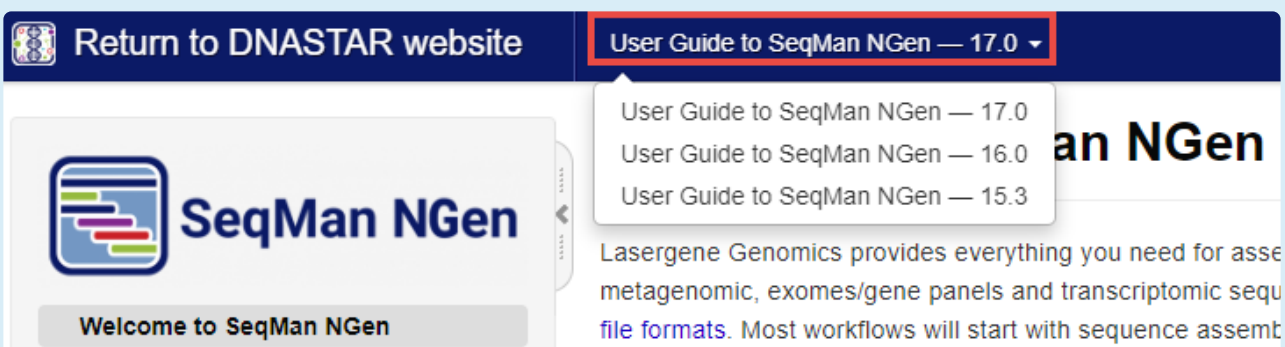
The tutorials listed below all launch from SeqMan Ultra and end with analysis in SeqMan Ultra. The tutorials are available in both this User Guide and the SeqMan Ultra User Guide. Within each tutorial is a link for downloading the corresponding data in archived (.zip) format.

- [Whole genome reference-guided workflow](#)
- [Whole genome *de novo* workflow with mate pair data](#)
- [Analysis of a whole genome *de novo* assembly](#)
- [De novo assembly using Sanger data](#)
- [RNA-Seq *de novo* transcriptome workflow](#)

Tutorials that end with analysis in any other application are listed below, and are found only in this SeqMan NGen User Guide. Within each tutorial is a link for downloading the corresponding data in archived (.zip) format.

- [RNA-Seq reference-guided workflow with analysis in ArrayStar](#)
- [ChIP-Seq workflow with analysis in ArrayStar](#)
- [Copy number variation \(CNV\) workflow with analysis in ArrayStar and GenVision Pro](#)
- [Whole-genome reference-guided workflow with analysis in ArrayStar |](#)
- [Long read analysis with accuracy evaluation](#)

✿ **Note:** SeqMan NGen went through a major update with Lasergene 17.0, which was released in early 2020. Therefore, there are two sets of tutorials and tutorial data. The set included in this version of the User Guide is for Lasergene version 17.0 and later. If you are using Lasergene version 16.0 or earlier, you will want to switch to that version of the User Guide before starting the tutorials. To switch, just use the drop-down menu right above the topic name (“SeqMan NGen Tutorials”) and choose the User Guide that corresponds with your version of Lasergene.

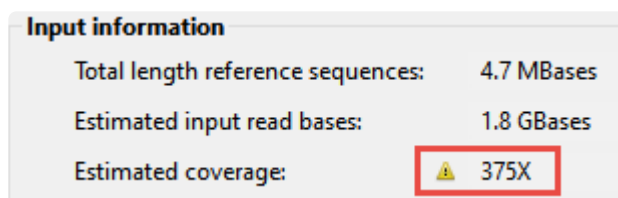


Whole genome reference-guided workflow

In this tutorial, you will create a reference-guided assembly using SeqMan NGen and then analyze the results using SeqMan Ultra. The time required for the assembly component is approximately 2-5 minutes.

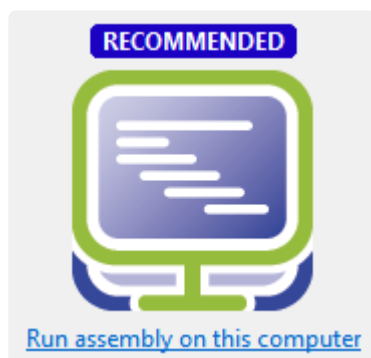
Running a reference-guided assembly in SeqMan NGen:

1. Download [T1_Whole_Genome_Ref.zip](#) (1.4 GB) and extract the contents to any convenient location (e.g., your desktop). The folder contains the following sequences:
 - Reference sequence *DH10B_NC010473.gbk*
 - Paired end sample sequences *SRR1284938_1.fastq* and *SRR1284938_2.fastq*
2. Launch SeqMan Ultra and choose **New Assembly** on the left. On the right, click on the **Genomics** workflow named **Variant analysis and resequencing**. This causes SeqMan NGen to open at the Workflow screen.
3. Choose the **NGS-Based** workflow named **Whole genome**.
4. In the Reference Sequence screen, use the **Add** button to add the sequence *DH10B_NC010473.gbk*. Then click **Next**.
5. In the Input Sequences screen, change the **Experiment setup** to **Single sample**. Press the **Add** button to add the paired reads *SRR1284938_1.fastq* and *SRR1284938_2.fastq*. Then click **Next**.
6. In the Assembly Options dialog, click **Next**.
7. In the Analysis Options screen, change the **Variant detection mode** to **Haploid** (since this is a bacterial organism) and the **SNP filter stringency** to **High**. Then click **Next**.
8. In the Assembly Output screen, type “Templated E coli” into the **Project Name** text box. This name will be assigned to all output files, including the finished assembly. Use the **Browse** button to specify a **Project Folder** for your assembly output files. For local users, an alternative way to select a location is to drag and drop a folder from the file explorer onto the **Project folder** row. The folder you choose must be in a writable location. Then click **Next**.
9. In the **Run Assembly Project** screen, note that the **Estimated coverage** is “375X.”



A coverage of 50-100X is ideal, and additional coverage only serves to increase the assembly time. There are several wizard parameters you can change to reduce coverage; here, you'll reduce the **Maximum total reads**.

10. Click **Assembly Options** in the left margin to return to that screen. Check the **Maximum total reads** box and type in 2000000 (2 followed by six zeros).
11. Click **Run Assembly Project** in the left margin to return to that screen. Note that the **Estimated coverage** is now "64X." Lower down the screen, observe that the recommendation is to perform a local assembly. Press **Run assembly on this computer**.




12. Wait until being informed that assembly has finished (approx. 2-5 minutes), then click **Next**.

The 1 assembleTemplate command succeeded
XNG done

Assembly finished successfully.


13. From the Assembly Summary screen, click **Open assembly** to launch the results in SeqMan Ultra.
14. Close the SeqMan NGen project by clicking the **Finish** button and confirming **Yes**.

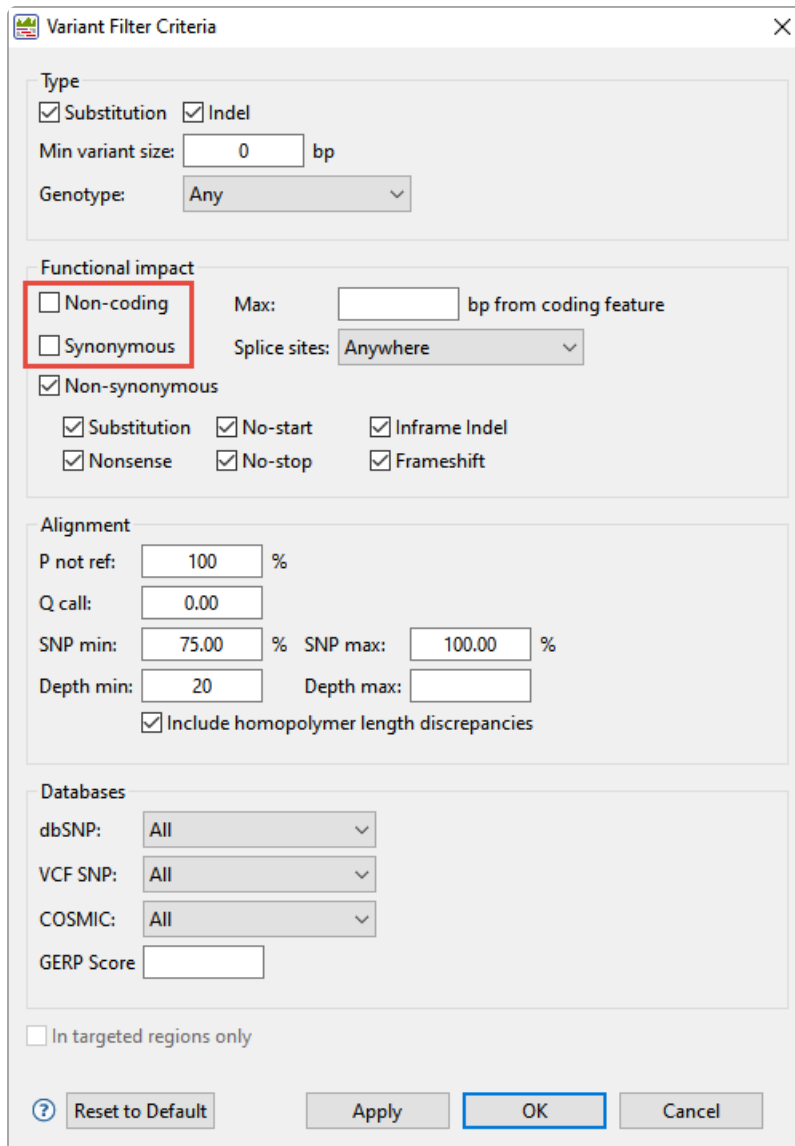
Analyzing variants and structural variations in SeqMan Ultra:

1. SeqMan Ultra's Explorer panel (**Explorer** tab on the upper right of the window) contains a single contig containing the bacterial chromosome. The contig is already selected by default.
2. Press the **Show table of variants** tool () to the right of the Explorer panel to reveal 151 variant positions. This number is displayed in the view header.

SNP	Ref ID	Cons Pos	Ref Pos	Type	Ref Base	Called...	Genotype	Impact	Homopolymer	P Not Ref	Q Call
?	NC_010473.1	63195	63193	SNP	A	C	Variant			99.9%	60.000
?	NC_010473.1	63209	63207	SNP	C	G	Variant			100.0%	60.000
?	NC_010473.1	63249	63247	SNP	G	A	Variant			100.0%	60.000
?	NC_010473.1	63251	63249	SNP	T	G	Variant			100.0%	60.000
?	NC_010473.1	63345	63343	SNP	A	G	Variant			100.0%	60.000
?	NC_010473.1	63347	63345	SNP	A	C	Variant			100.0%	60.000
?	NC_010473.1	63369	63367	SNP	G	A	Variant			100.0%	60.000
?	NC_010473.1	63422	63420	SNP	C	T	Variant			100.0%	60.000
?	NC_010473.1	63427	63425	SNP	G	A	Variant			100.0%	60.000
?	NC_010473.1	63577	63575	SNP	G	T	Variant			100.0%	60.000
?	NC_010473.1	63676	63674	SNP	C	A	Variant			100.0%	60.000
?	NC_010473.1	63805	63803	SNP	C	T	Variant	Synonymous		100.0%	60.000
?	NC_010473.1	63860	63858	SNP	C	T	Variant	Non-synonymous		100.0%	60.000
?	NC_010473.1	64538	64536	SNP	G	A	Variant	Synonymous		100.0%	60.000
?	NC_010473.1	64541	64539	SNP	C	T	Variant	Synonymous		99.9%	57.730

Imagine that you wish to identify those variants which are most likely to alter function of genes in *E. coli*. You will therefore want to focus on variations predicted to have a non-synonymous impact on the *E. coli* genome. In the next step, you will filter out some of the variants and only leave those predicted to cause a coding change.

- To show only variants that cause coding changes, click the **Filter all variant tables** tool () in the upper right of the view. In the Variant Filter Criteria dialog, remove the checkmarks next to **Non-coding** and **Synonymous** and click **OK**.



Variant Filter Criteria

Type

☒ Substitution ☒ Indel

Min variant size: bp

Genotype:

Functional impact

☐ Non-coding Max: bp from coding feature

☐ Synonymous Splice sites:

☒ Non-synonymous

☒ Substitution ☒ No-start ☒ Inframe Indel

☒ Nonsense ☒ No-stop ☒ Frameshift

Alignment

P not ref: %

Q call:

SNP min: % SNP max: %

Depth min: Depth max:

☒ Include homopolymer length discrepancies

Databases

dbSNP:

VCF SNP:

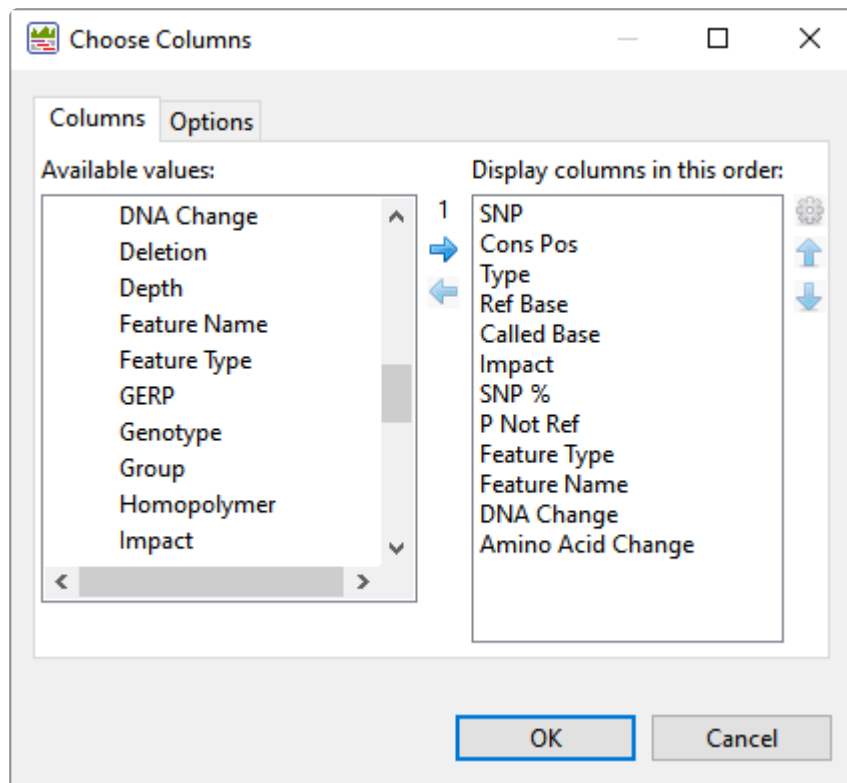
COSMIC:

GERP Score

☐ In targeted regions only

The remaining 22 variants all have a non-synonymous impact or consist of an in-frame deletion.

4. Use the scrollbar on the bottom of the view to scroll right, and observe the many data columns available in this report. While each column is useful in certain circumstances, many will not be needed for this tutorial workflow.
5. To remove unneeded columns from the report, click the **Change alignment options** tool (⚙️). Using the image below as a guide, move all of the wanted columns “up” into the position shown below using the up arrow. Then select the remaining, unwanted, columns and press the left arrow to dismiss them. Press **OK**.



Note that the view is now much more compact and may be easier to view in its entirety, depending on your monitor size.

Now turn your attention to the black arrows in some of the **SNP** column cells. These represent “coalesced variants.” In .assembly projects like this one (but not .sqd projects), variants in adjacent columns are coalesced into a single insertion or deletion if they are of the same type, and if at least 80% of the reads with the called variant in one column have a variant in the adjacent column.

- Click any of the triangles to open the coalesced multiple-base indels and reveal individual variants. After clicking a triangle, information on each position of the insertion/deletion is displayed in a separate row.

NC_010473.1 SNPs Summary 22 variant positions											
0 confirmed 0 rejected 22 putative 0 mixed 158 filtered											
SNP	Cons ...	Type	Ref Base	Called Base	Impact	SNP %	P Not Ref	Feature Type	Feature Name	DNA Change	Amino Acid Change
?	63860	SNP	C	T	Non-synonymous	100.0%	100.0%	CDS	mraZ	c.121C>T	p.H41Y
?	218679	SNP	T	A	Non-synonymous	100.0%	100.0%	CDS	yafJ	c.239T>A	p.L80Q
▼ ?	418879	SNP [2]	AA	GG	Non-synonymous	100.0%	100.0%	CDS	hha	c.13_14delinsCC	p.F5P
?	418879	SNP	A	G	Non-synonymous	100.0%	100.0%	CDS	hha	c.14T>C	p.F5S
?	418880	SNP	A	G	Non-synonymous	100.0%	100.0%	CDS	hha	c.13T>C	p.F5L
?	751890	SNP	A	G	Non-synonymous	100.0%	100.0%	CDS	nagD	c.289T>C	p.Y97H
?	844877	SNP	A	G	Non-synonymous	100.0%	100.0%	CDS	galE	c.368T>C	p.F123S
▶ ?	862407	SNP [2]	AG	CA	Non-synonymous	100.0%	100.0%	CDS	bioA	c.40_41delinsTG	p.L14W
?	971162	SNP	A	G	Non-synonymous	100.0%	100.0%	CDS	ybjD	c.1494A>G	p.I498M

- Look at the **Impact** column for the variant in the first row (Cons Pos = 63860), and note that a non-synonymous mutation is predicted.

8. Double-click on the row to launch the Alignment view with the variant highlighted.

The Alignment View displays the consensus sequence, and the alignment of all sequences making up a contig, at the nucleotide level. The consensus sequence for the selected contig is displayed at the top of the window, immediately beneath the ruler. Constituent sequences are displayed in the bottom section.

9. Click to either side of the highlighted column to temporarily clear the selection. This will allow for a better view of the color-coding used by SeqMan Ultra.

NC_010473.1 Position=17

Ruler

Consensus

Coverage

NC_010473.1(1>4686377) CATTGACATTTCATCACCCTGCTGCT

ID_1698742_f → C

ID_655022_r ← CATT

ID_820868_r → CATT

ID_713240_r ← CATTGACA

ID_1852302_r ← CATTGACATTAT

ID_914064_r ← CATTGACATTATCAC

ID_1636180_f → CATTGACATTATCAC

ID_115384_r → CATTGACATT

ID_1775994_f ← CATTGACATTATCACCC

ID_461262_f → CATTGACATTATCACCCG

ID_1254702_f ← CATTGACATTATCACCCGTGCT

ID_1471690_f ← CATTGACATTATCACCCGTGCTGCT

ID_352602_r → CATTGACATTATCACCCGTGCACT

ID_461714_f → CATTGACATTATCACCCGTGCTGCT

ID_514280_f ← CATTGACATTATCACCCGTGCTGCT

NC_010473.1 x

NC_010473.1 SNPs Summary 1 selected rows; 22 variant positions


0 confirmed 0 rejected 22 putative 0 mixed 158 filtered

SNP	Cons Pos	Type	Ref Base	Called...	Impact	SNP %
?	63860	SNP	C	T	Non-synonymous	100.0%
?	218679	SNP	T	A	Non-synonymous	100.0%

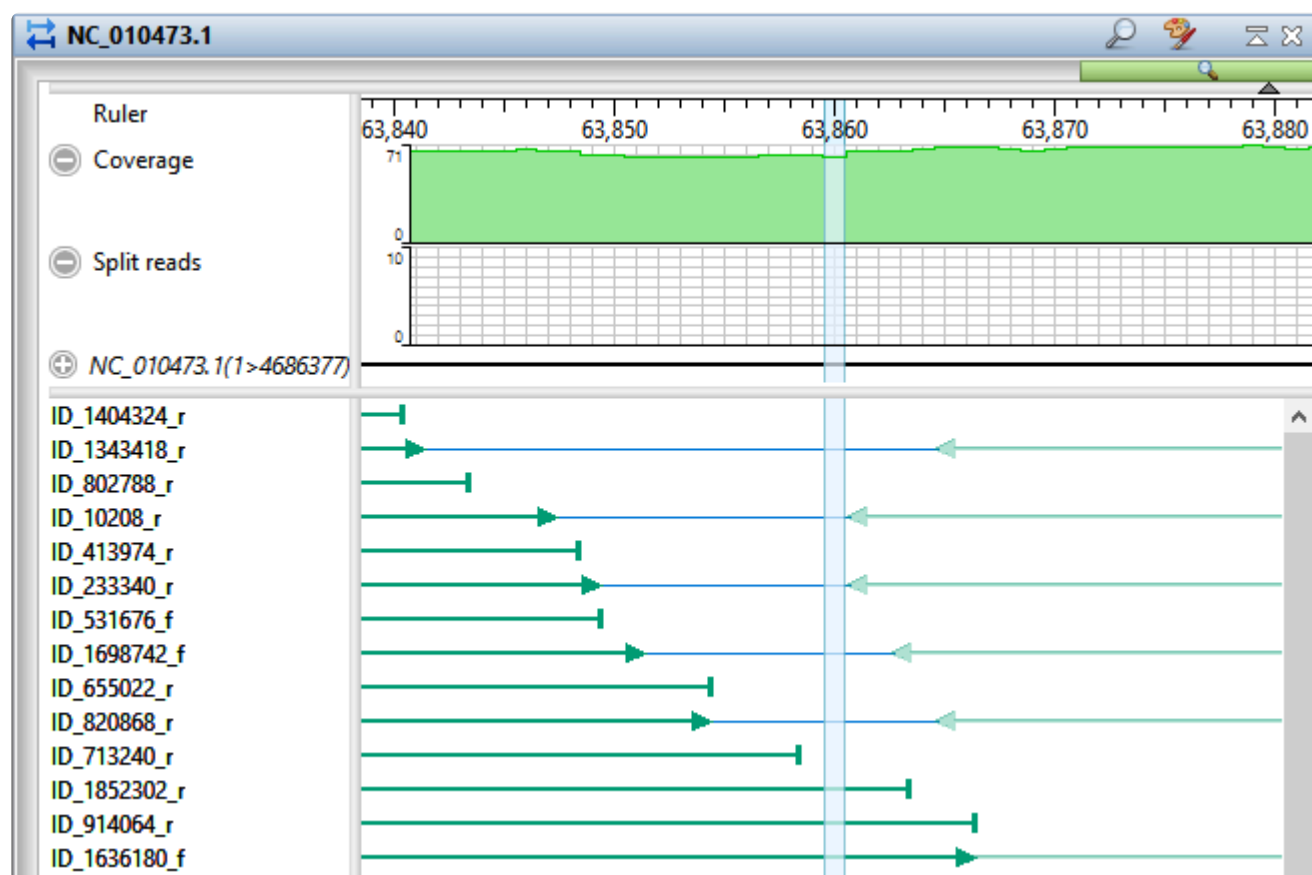
Note that the **Reference** base at the position of interest is **C**, while all of the samples contain ***T** at that position. The yellow highlighting used for variants is specified in the Alignment area of the Style panel and can be changed, if desired.

10. To view the variant in the Strategy view, select **View > Strategy > Show Strategy View**. Because two views are already open, the Strategy view appears in a pop-up window. The view graphically

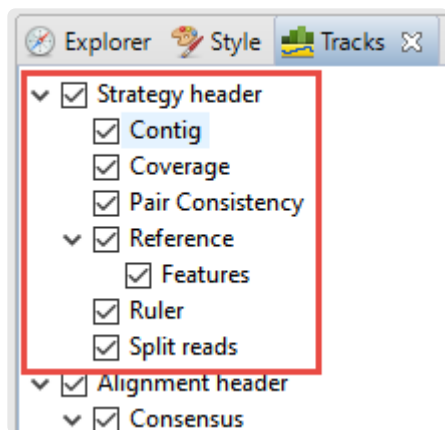
summarizes the position and orientation of every constituent sequence in the selected contig.

11. Zoom in as far as possible by dragging the green Zoom slider () all the way to the right.
12. Click on the top row of the Variants view table to again move the cursor to the variant of interest.

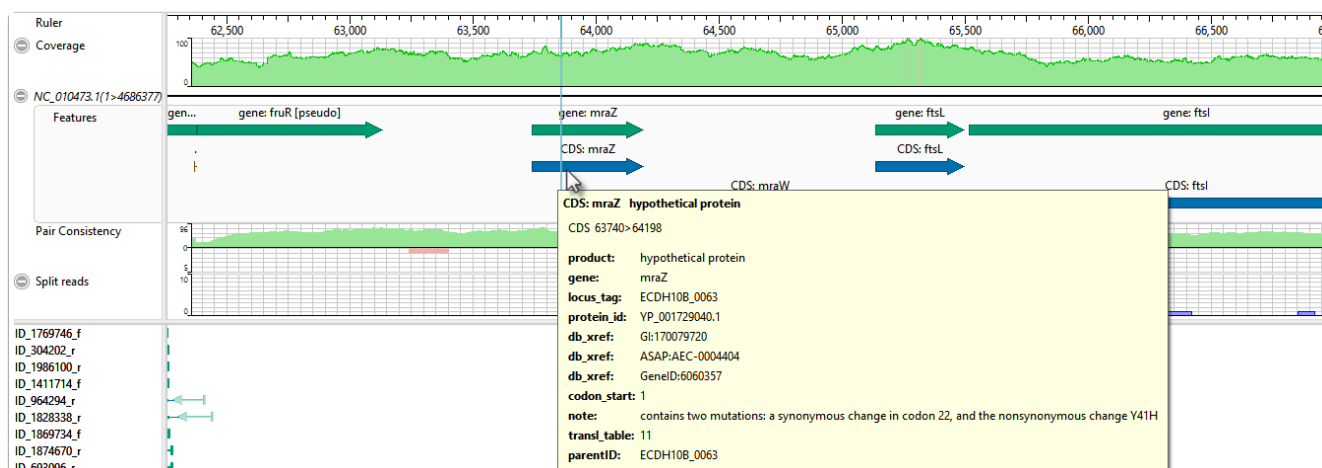
In the lower part of the window, observe the green arrows. Arrows are used to graphically represent the individual sequences in the contig. The point of the arrowhead represents the 3' end and the tail represents the 5' end. The green coloration of arrows in this region denotes consistent, correctly-specified paired reads.



13. To add data tracks to enhance the Strategy view, click on the Tracks tab and check all boxes in the "Strategy header" section.



14. Returning to the floating Strategy view window, click any plus signs appearing next to the header tracks to reveal the “hidden” tracks.
15. Move the horizontal zoom slider (at the top) until you can read the feature names, as shown below. Hover over any part of a histogram to display information about that location.



- **Coverage** – Represents the coverage within the assembly, with green representing above-threshold coverage and red representing areas of low coverage.
 - **Split reads** – Shows the location and depth of split reads.
 - **Features** – The putative variant is located within a labeled CDS called *mraZ*.
 - **Pair Consistency** – Shows consistent reads in green and inconsistent in red. The deeper the histogram, the more recognized pairs of reads there are.
16. Close the floating Strategy view window by clicking the ‘x’ in the top corner. Then direct your view to the Variants view at the bottom of the main window.
 17. Confirm the variant that you just investigated by clicking on the **SNP** (leftmost) column to change the question mark to a checkmark. (If you had found supporting evidence was lacking, you could have

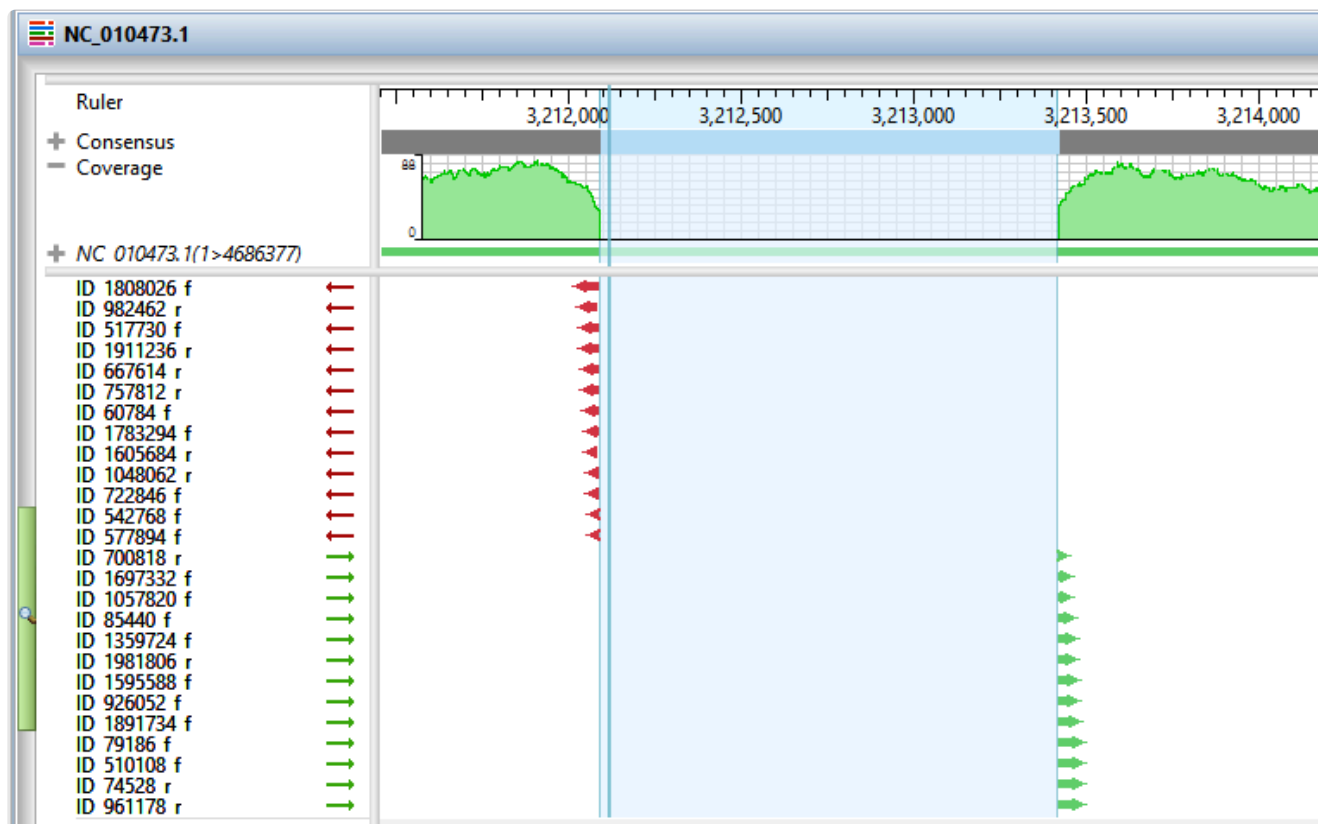
clicked a second time to reject the variant and add an 'x'.)

Short indels are included in the Variants view. However, any indels long enough to inhibit assembly are instead gathered in another report called the Structural Variation Report.

18. To see the longer indels, choose **Contig > NC_010473.1 Structural Variation**. Scroll down to the putative Indel at 3211929.

NC_010473.1 Contig Structural Variation						
Contig ID	Ref Pos	Length	Type	Splits	Cover...	Region Capture
NC_010473.1	3003961	1331	Del	29	34.0	no
NC_010473.1	3019404	56934	Indel	-----	62.0	no
NC_010473.1	3170620	1190	Del	44	44.0	no
NC_010473.1	3170626	1184	Indel	-----	44.0	no
NC_010473.1	3199470	1328	Indel	-----	.0	no
NC_010473.1	3211929	1328	Indel	-----	.0	no
NC_010473.1	3229073	44443	Indel	-----	63.0	no
NC_010473.1	3336619	20201	Indel	-----	62.0	no
NC_010473.1	3381608	11894	Indel	-----	63.0	no

19. Double-click on that table row to select it and center the Alignment view on that position. Use the horizontal zoom slider and vertical scrollbar and note that this area is devoid of reads.



This marks the end of this tutorial.

Whole genome de novo workflow with mate pair data

In this tutorial, you will use SeqMan NGen to de novo assemble an *E. coli* K12 MG155 data set composed of Illumina MiSeq (2×300) paired end reads.

The default stringency setting is designed to produce longer contigs; however, some may contain false joins. You'll learn how to manually examine contigs for false joins and split them into two contigs. Finally, you'll realign all the contigs, including the newly split ones, using SeqMan Ultra.

Assembling the data in SeqMan NGen:

1. Download [T2_Whole_Genome_MatePair.zip](#) (231 MB) and extract it to any convenient location (e.g., your computer's desktop). The folder contains the following paired end reads:

Ecoli_S1_L001_R1_001.fastq

Ecoli_S1_L001_R2_001.fastq

2. Launch SeqMan Ultra and choose **New Assembly** on the left. On the right, click on the **Genomics** workflow named **De novo genome assembly and editing**. This causes SeqMan NGen to open at the Workflow screen.
3. On the right of the screen, under **NGS-Based**, click the workflow named **De novo assembly**.
4. In the Input Sequences screen, press the **Add** button to add the two paired reads whose names begin with *Ecoli_S1_L001_R*. Then click **Next**.
5. In the Preassembly Options screen, change the **Maximum total reads** to 800000 (8 followed by five zeros). Click **Next**.
6. In the Assembly Options dialog, type in an **Estimated genome length** of 4600000 (46 followed by five zeros). Click **Next**.
7. In the Assembly Output screen, type "De novo mate pair" into the **Project Name** text box. This name will be assigned to all output files, including the finished assembly. Use the **Browse** button to specify a **Project Folder** for your assembly output files. For local users, an alternative way to select a location is to drag and drop a folder from the file explorer onto the **Project folder** row. The folder you choose must be in a writable location. Then click **Next**.
8. In the **Run Assembly Project** screen, press **Run assembly on this computer**.
9. Wait until being informed that assembly has finished (approx. 1 hour), then click **Next**.

10. From the Assembly Summary screen, click **Open assembly** to launch the results in SeqMan Ultra.
11. Close the SeqMan NGen project by clicking the **Finish** button and confirming **Yes**.

Correcting misjoined contigs and realigning data in SeqMan Ultra:

After following the steps above, your *.sqd* assembly file is now open in SeqMan Ultra. If analyzing your own data in real life, you would keep the current file open and proceed with the analysis. For purposes of this tutorial, however, you'll close this assembly and open a similar one that is included in your tutorial data. Why? No two assemblies will be completely identical, due to randomization of the order in which contigs are created by the assembler. Because this tutorial will refer to specific contigs, it will be easier to follow along if you are using a pre-made assembly. Therefore, you will close the current assembly in the first step, below.

1. Use **File > Close Editor** to close the current file, *De novo mate pair.sqd*.
2. Open the pre-made assembly by using **File > Open** to open the demo data file *Sample de novo mate pair.sqd*.
3. Observe the information in the Project Overview on your screen. Briefly note the high **Contig N50** value of 149 kb. This means that half of the genome sequence is located in contigs equal to or longer than 149 kb.

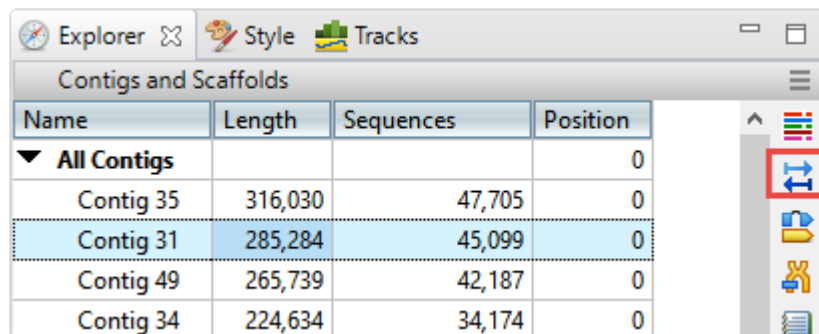
Sample de novo mate pair.sqd

C:\Users\yildizs\Desktop\T2_Whole_Genome_MatePair

Contig N50:	149.0 kb
Largest contig size:	316,030 bp
Number of contigs:	66
Reads assembled:	720,452

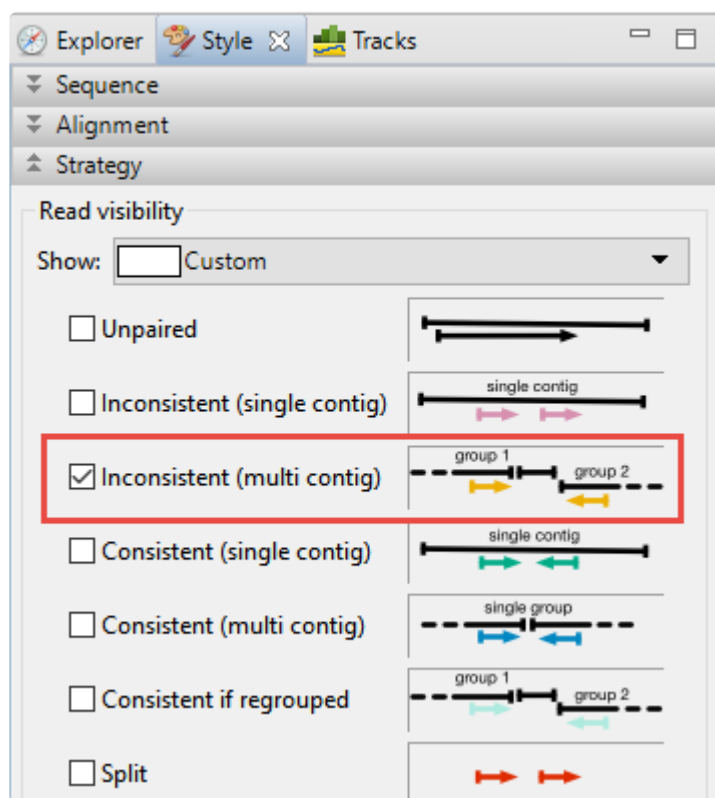
Using a high stringency in SeqMan NGen would likely have resulted in about ten times as many contigs. However, a tradeoff is that some contigs might have been joined when they should not have been. You will address that issue in the next few steps. In “real life,” you would check all of the contigs for false joins (or at least the contigs longer than 100k) and split them where needed. For purposes of this tutorial, you will only examine two contigs: 31 and 38.

4. In the Explorer panel on the right, select Contig 31, which is one of the longer contigs in the project. Then press the **Show strategy view** tool.

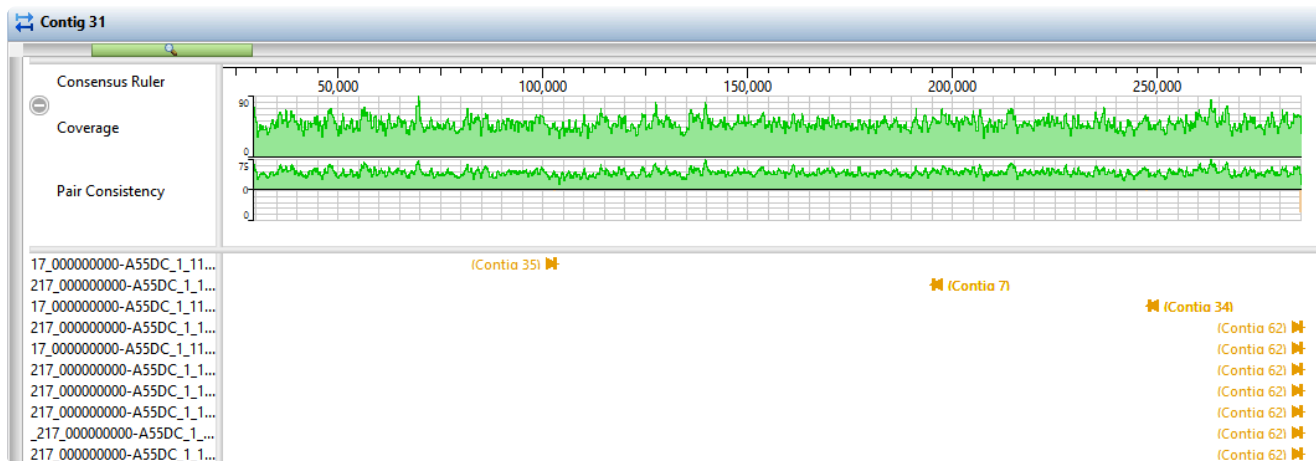


Name	Length	Sequences	Position
▼ All Contigs			0
Contig 35	316,030	47,705	0
Contig 31	285,284	45,099	0
Contig 49	265,739	42,187	0
Contig 34	224,634	34,174	0

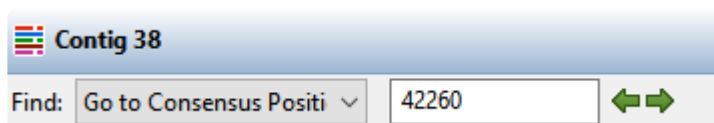
5. You will be looking for locations along the contig that are flagged for inconsistent pairs. To easily see inconsistent pairs, click the **Style** tab on the right. In the **Strategy** section, uncheck all boxes except **Inconsistent (multi-contig)**.



6. In the Strategy view on the left, slide the green horizontal zoom slider to the left so the entire contig is visible. Note that the **Pair Consistency** histogram near the top does not show any areas of orange, meaning that there are no areas of inconsistency. Looking at the orange arrows at the bottom, ignore the ones at the right edge; these indicate that Contig 31 is likely adjacent to Contig 62. Instead, look for orange arrows that are not on the edges. In this case, there are only three orange arrows that represent possible inconsistent pairs. The very small number indicates that these three are likely artifacts. Contig 31 thus exemplifies a correctly-built contig without false joins.



7. Click the **Explorer** tab and select Contig 38. Click the **Show strategy view** tool, then click the **Show Alignment view** tool just above it.
8. In the Strategy view (upper view), slide the green horizontal “zoom slider” all the way to the left so you can see the entire sequence. In the Alignment view (lower view), slide the zoom slider all the way to the right for a detailed view.
9. In the Strategy view (upper view), look at the area around 42,000. Both the **Pair Consistency** histogram and the large number of arrows denote that this is a likely misjoin. The arrows facing left have pairs in Contig 47, while the arrows facing right have pairs in Contig 55. Contig 38 represents a contig with at least one false join. (There is another candidate at the very left of the contig).
10. Near the upper right corner of the Alignment view (lower view), click the **Search alignment** tool (🔍). On the left, use the **Find** menu to select **Go to Consensus Position**. Then type **42260** into the search bar and press either green arrow to navigate to the position.



11. The position 42,260 is now selected, but the cursor needs to be between two nucleotides. In the Alignment view, click just to the right or left of 42,260 to remove the selection
12. Choose **Contig > Split At Insertion**. A new contig is added to the Explorer panel. The new contig is the only one in the table with a **Position** other than zero.

In a real-world project, we recommend going through each contig as above. For this tutorial, however, splitting Contig 38 is sufficient. In the next steps, you will create contig scaffolds.

13. In the Explorer panel, select the uppermost row, **All Contigs**. Right-click on the selection, choose **Order Contigs into Scaffolds** and press **OK**. SeqMan Ultra creates 23 scaffolds numbered 100 to 230.

14. Select the **Scaffold 100** row. Choose **Contig > Align Contigs End to End** and press **OK**. Repeat for all other Scaffold rows. Most scaffolds will collapse into 1-2 contigs. Some scaffolds (e.g., Scaffold 170) will not collapse at all, and will give a message to that effect. In these case, simply move on to the next scaffold on the list.
15. One at a time, drag each contig that remains in a scaffold and drop it anywhere above **Scaffold 100**. Then use **Shift+click** to select all of the now empty scaffolds and press the **Delete** key on your keyboard.
16. Choose **Project > Project Overview** and observe how the steps above have improved the assembly, even though you only split one false join.

Sample de novo mate pair.sqd

C:\Users\yildizs\Desktop

Contig N50: 204.0 kb
Largest contig size: 361,504 bp
Number of contigs: 44
Reads assembled: 720,452




Note that the **Contig N50** has increased from **149kb** to **204kb**, meaning the typical contig is now almost 37% longer than before. The longest contig is over 361,000 bp in length. The number of contigs decreased from 66 to 44. Splitting off additional false joins between Steps 8-9 would likely have improved the assembly even more.

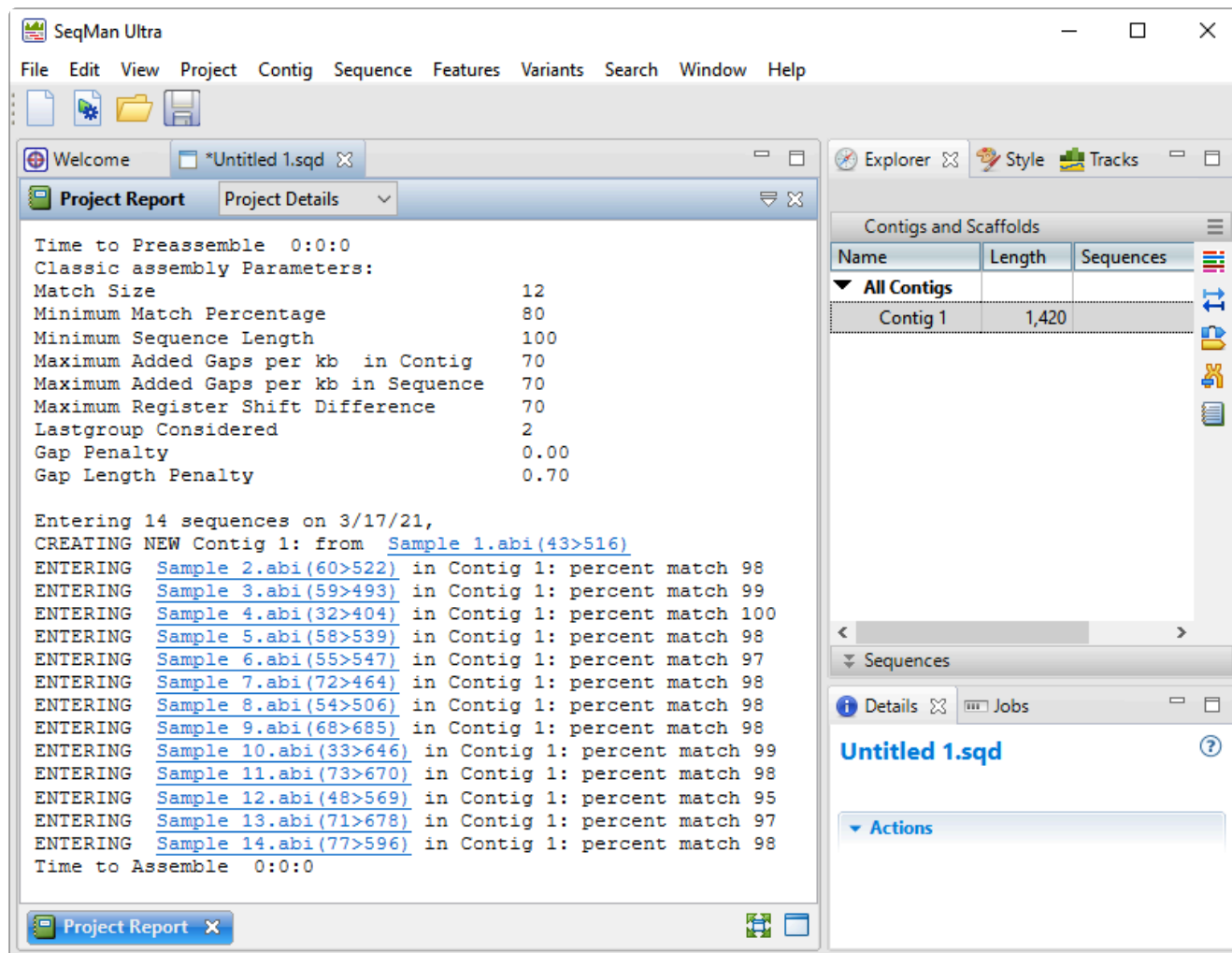
This marks the end of this tutorial.

De novo assembly using Sanger data

In this tutorial, you will *de novo* assemble fourteen short trace sequences from PE Applied Biosystems, Inc. and then analyze the resulting contig in SeqMan Ultra. Scroll to the bottom of this topic for a short video showing a slightly different version of this tutorial using the same data set.

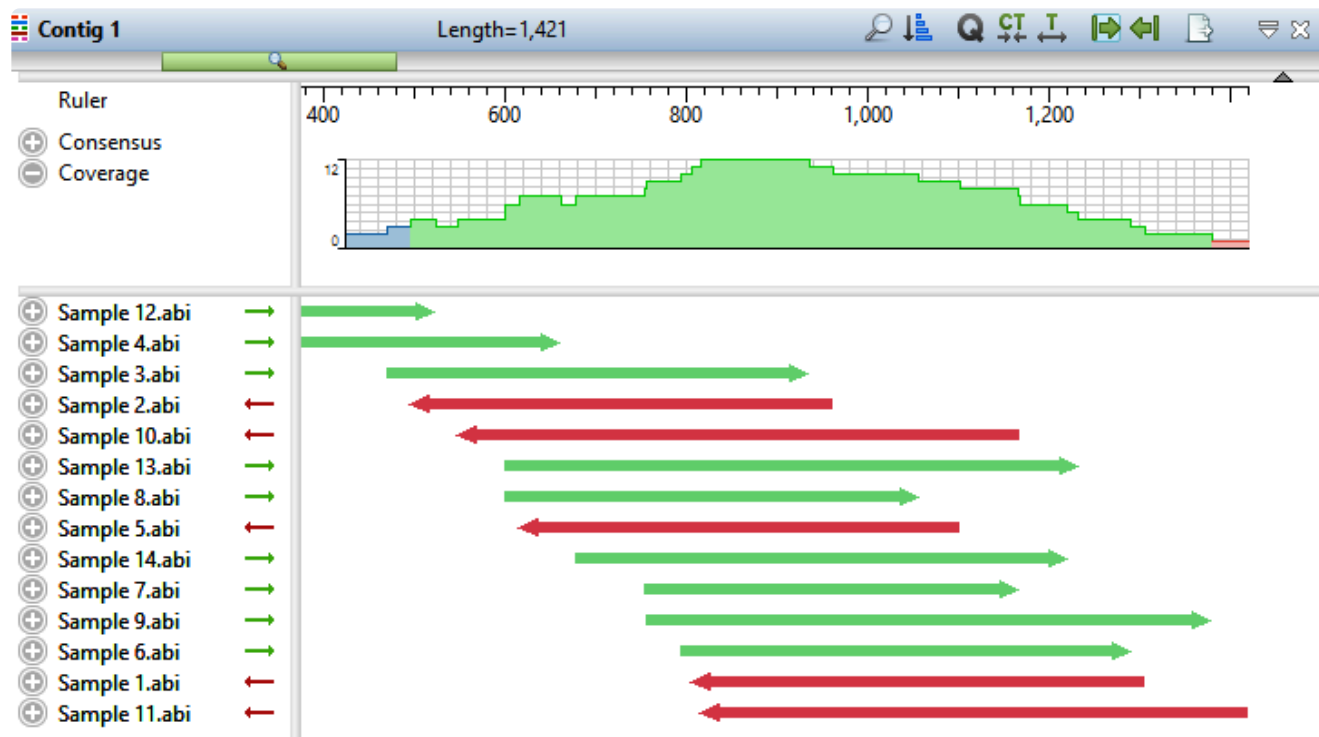
Assembling Sanger trace reads:

1. Download [T4_Sanger_DeNovo.zip](#) (1 MB) and extract it to any convenient location (i.e., your desktop). The data set consists of a folder of *Janus vectors* and fourteen *.abi* sample sequences.
2. Launch SeqMan Ultra. On the left of the window, choose the blue bar named **New Assembly**.
3. Under **Molecular biology**, click on **New Sanger/ABI assembly**. The “Unassembled” window opens.
4. Press the **Add** button near the top left ( **Add**), select the 14 *.abi* sequences and press **Open**.
5. In the header area, leave **Quality Trim** checked. This tells SeqMan Ultra to trim read ends based on trace quality evaluation. To remove Janus vector, add a checkmark next to **Vector Trim** and click on the file path just to its right. Navigate to and select the folder *Janus vectors* then press **Select Folder**.
6. Press the **Trim Now** button ( **Trim Now**). In the **Vector** column in the table, observe that the vector *InvJanus* was discovered and removed.
7. Press the **Assemble** button ( **Assemble**) to begin the assembly. The assembly should take a few seconds, at the most.



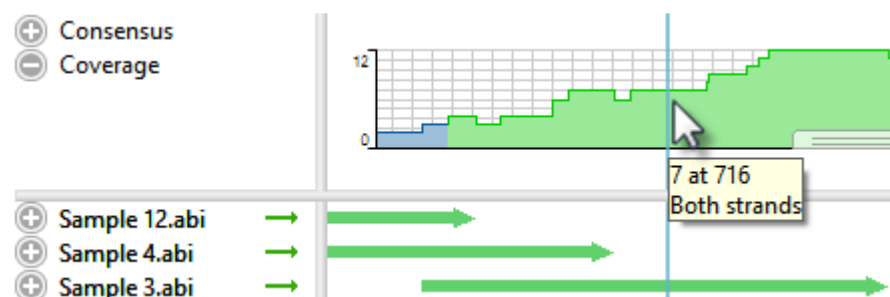
Examining the alignment:



1. In the Explorer tab on the right, note that the assembly resulted in one contig.
2. Double-click on **Contig 1** to open it in the Alignment view. Move the green horizontal zoom slider (near the top of the view) to the left until you can see the entire contig.



The Coverage graph has areas of blue, green and red. Blue indicates single-direction coverage, while red shows single-read coverage. Green denotes coverage on both strands; the height of the histogram corresponds to the depth of coverage.

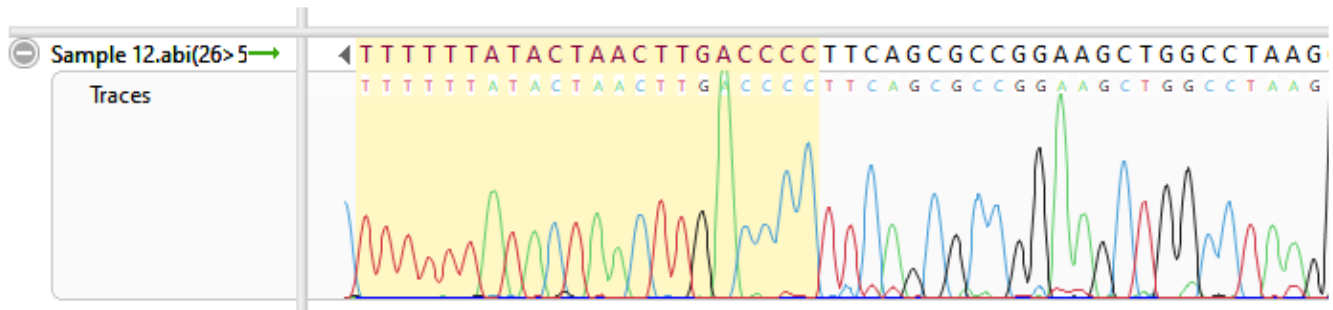
3. Hover the mouse over different parts of this histogram to see tooltips showing the coverage at a given position and whether or not it meets threshold requirements.



4. To zoom in to view details, click the **Restore default zoom** tool () at the top right of the view.
5. To locate conflicts, shown by default with yellow highlighting, click the **Search alignment** tool (). Use the **Find** menu to choose **Conflict**, then press the green arrow keys to navigate from one conflict to the next.
6. To view the trace data, right-click on any sample name on the left and choose **Expand All**.
7. Scroll all the way to the left so that only Sample 12.abi is visible.

During assembly, SeqMan NGen trimmed ends for the constituent sequences based on trace data quality and presence of vector. Although sufficient data remained to assemble the sequences into a single contig, there are cases when restoring some of the trimmed data may allow SeqMan Ultra to join multiple contigs into a single contig. You may also wish to restore data in order to verify the consensus in a low-coverage area.

To reveal the trimmed trace data, grab the black triangle to the left of the sequence and drag it to the left. Trimmed sequence appears with a yellow background by default. Conflicts between the restored data and the consensus are shown via red text.



8. To unmask the trimmed regions at both ends of the contig, use **Contig > Extend Trimmed Ends > Extend 3' and 5' Sequence Ends**.

- Trimmed portions of notably poor quality (e.g., misshapen and overlapping peaks) were removed because the average peak quality fell below the acceptable stringency threshold.
- Data removed due to vector contamination is characterized by normal peak quality in combination with a high number of conflicts. In this example, such regions were removed because they originated from the Janus vector. Regardless of its deceptively high peak quality, it is not recommended that you restore vector data.

If you thought the trimmed ends merited being kept in the alignment (they don't in this case), you could have used the command **Contig > Extend Trimmed Ends > Extend and Align 3' and 5' Sequence Ends** to both extend the ends and reassemble the reads.

The following video illustrates the tutorial workflow using the same data set but slightly different steps.

This marks the end of this tutorial.

Analysis of a whole genome de novo assembly

DNASTAR Lasergene lets you set up a *de novo* assembly with ease. One of the outputs is an editable *.sqd* file that can be opened and edited in SeqMan Ultra. In SeqMan Ultra, you can evaluate the assembled contigs, edit them, organize them into scaffolds, and close any gaps in those scaffolds.

This tutorial uses a *de novo* whole genome assembly that was assembled from two MiSeq 2×300 paired-end read files from *E.coli* K12 MG1655. SeqMan NGen was used to assemble 2.5M reads from this data set, producing an assembly with a large contigN50 of 203Kbases in 51 contigs. Due to the size of the reads (28 GB combined), this tutorial will begin at the stage of downstream analysis in SeqMan Ultra.

1. Download [T3_Whole_Genome_DeNovo.zip](#) (332 MB) and extract it to any convenient location (i.e., your desktop). The data set consists of a single SeqMan Ultra project named *E.coli K12 MG1655 MiSeq de novo.sqd*.
2. Launch SeqMan Ultra and use **File > Open** to open *E.coli K12 MG1655 MiSeq de novo.sqd*. Note that the Project Overview shows that the Contig N50 is 204.0 kb, which is quite large. Also observe that the assembly resulted in 51 contigs.


E.coli K12 MG1655 MiSeq de novo.sqd

C:\Users\yildizs\Desktop\T3_Whole_Genome_DeNovo

Contig N50: 204.0 kb
 Largest contig size: 361,321 bp
 Number of contigs: 51
 Reads assembled: 2,182,744

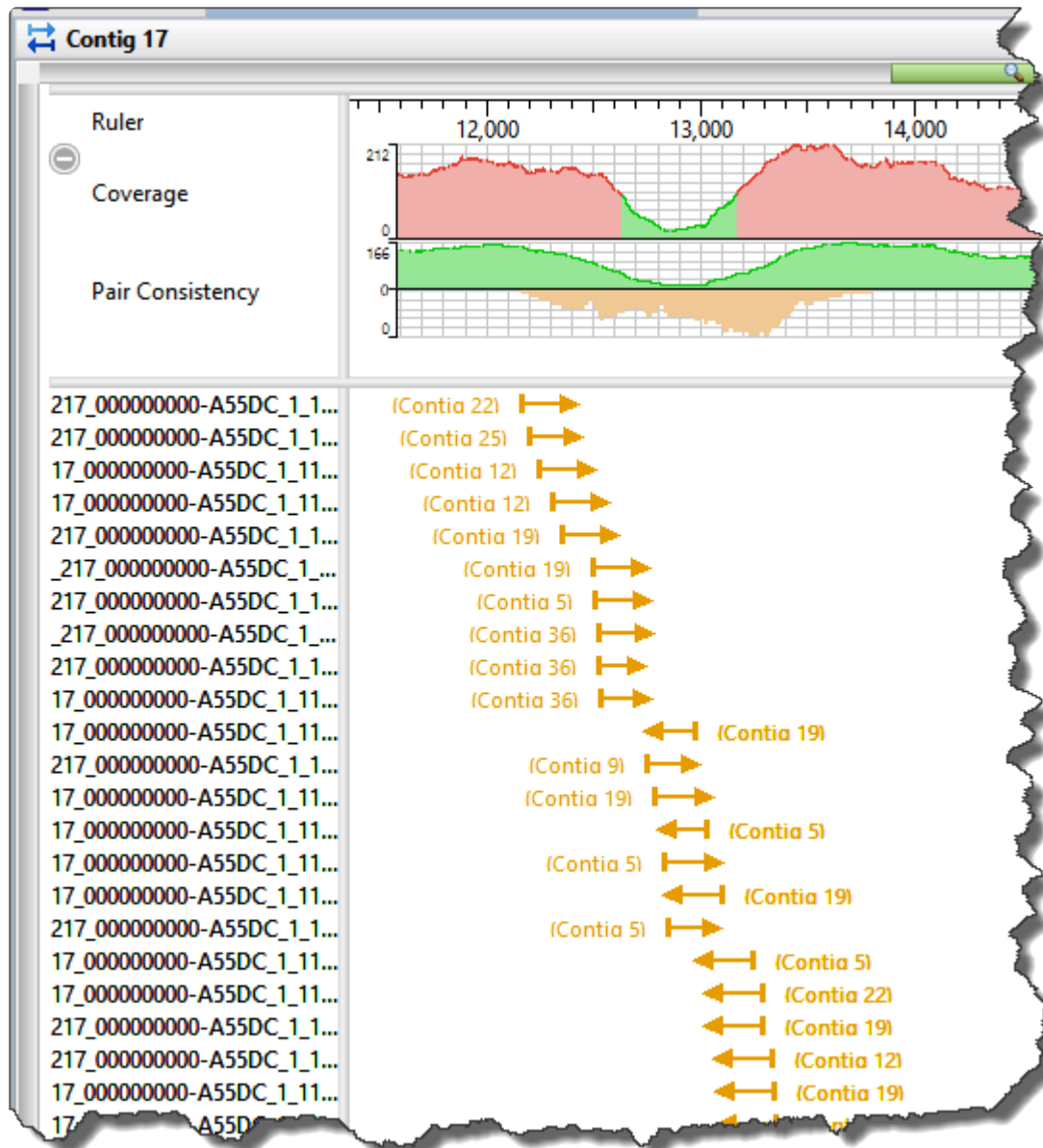
3. The header of this section contains a drop-down menu. Use the menu to select **Project Details**. This section has additional statistics followed by the complete script that was used by SeqMan NGen to assemble the reads. In the upper section of the report, observe that there are an average of 42798 sequences per contig and that average depth of coverage is 135.

In the next few steps, you will check for potential assembly errors and then correct them. Assume that you have already performed these steps for Contigs 1-16 and have found no errors.

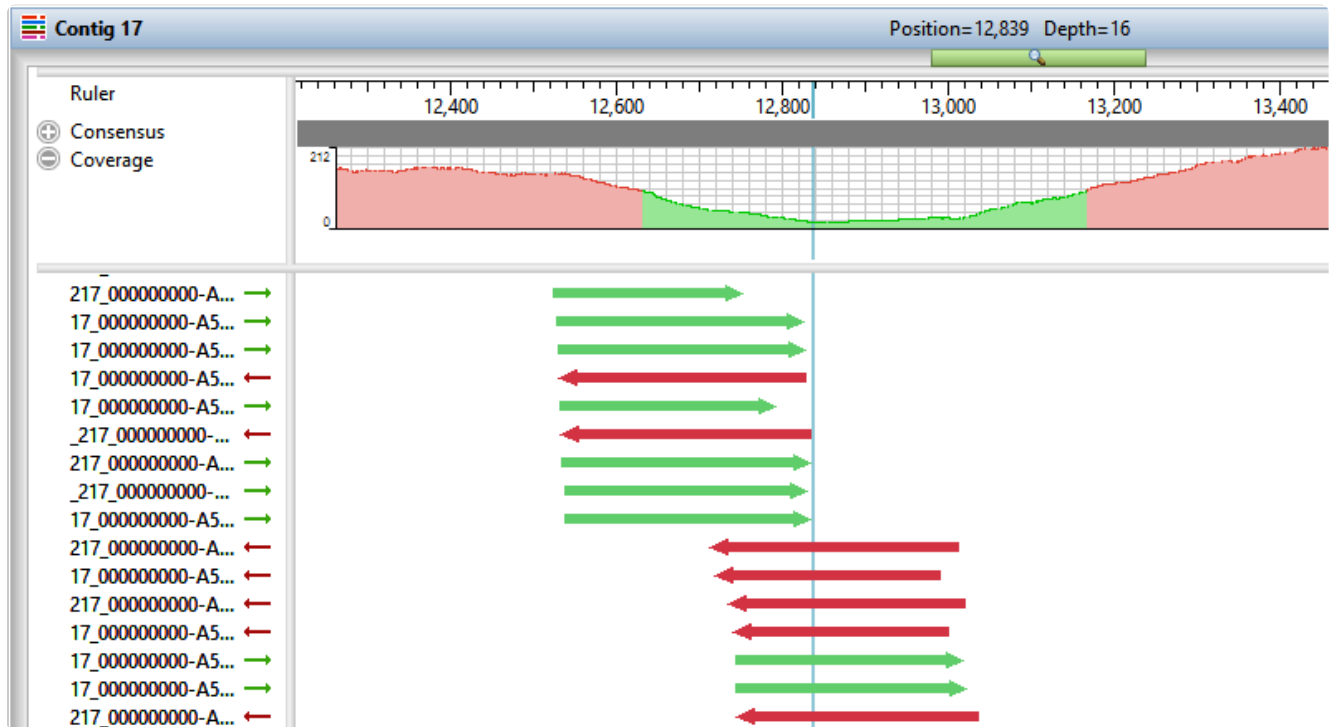
4. In the Explorer panel on the right, select **Contig 17** and press the **Show strategy view** tool () on the right.
5. To limit the display to show only arrows for inconsistent pairs, click the **Style** tab on the right. In the Strategy area, use the **Show** menu to select **Inconsistent**. The orange arrows represent pairs in the middle of contigs that are not matched up with their mate; the numbers associated with these pairs

indicate the contig where their “mate” resides.

- Use the horizontal zoom slider and scroll until you are centered on the area around 13,000. Observe the large number of orange (inconsistent) arrows in this section, as well as orange coloration in the **Pair Consistency** graph in the header. This represents a possible mis-assembled contig. Note that Contigs 5, 19 and 22 are referenced by both left- and right-facing arrows, indicating the region is likely a repeat. This is a good candidate for a contig split.



- Open the Alignment view for this contig by right-clicking in any white portion and choose **Show in Alignment view**. Use the horizontal zoom slider and the horizontal scroll bar to navigate to the same area visible in the Strategy view. Insert the cursor in the low pair-consistency region; the thinnest area of green in the Coverage histogram.



8. Choose **Contig > Split at Insertion** to split the contig into two contigs.
9. Return to the Explorer panel by clicking the **Explorer** tab on the right.
10. Select the uppermost row, **All Contigs** then right-click and choose **Order Contigs into Scaffolds**. When prompted, confirm that you wish to perform scaffolding. SeqMan Ultra uses the pair information at the ends of the contigs to order them into scaffolds.

Fourteen scaffolds are created, numbered from 100-230. The reason that microbial genomes do not assemble into a single contig (when using short reads) is a high prevalence of repetitive elements, like transposons that occur throughout the genome.

11. Now that the contigs have been ordered into scaffolds, you can now merge adjacent overlapping contigs. To do this, select Scaffold 100 and choose **Contig > Align Contigs End to End**. Keep the default settings and press **OK**. Repeat this procedure for all fourteen scaffolds. On occasion, you may get a message saying that the alignment didn't work. In this case, just move on to the next scaffold.
12. Use **Ctrl/Cmd+click** to select each contig that is in a scaffold. Then drag and drop the entire selection anywhere above **Scaffold 100**. Finally, select all of the now-empty scaffolds and press the **Delete** key on your keyboard to remove them.
13. To see the effects of these edits and realignments, use **Project > Project Overview**.

E.coli K12 MG1655 MiSeq de novo.sqd

C:\Users\yildizs\Desktop\T3_Whole_Genome_DeNovo

Contig N50: 290.0 kb

Largest contig size: 496,706 bp

Number of contigs: 36

Reads assembled: 2,182,744

Note that the **Contig N50** has increased from 204 kb to 290 kb, and that the number of contigs has decreased from 51 to 36.

In a real-world situation, you could continue merging contigs and closing gaps via one or more of the following:

- Continue to make contig edits, create scaffolds and merge contigs to close additional gaps, as above.
- [Perform a BLAST search](#) on contig ends to determine genome coordinates, then manually create new scaffolds and attempt additional end-to-end alignments.
- Resolve remaining gaps by adding new data (e.g., Sanger reads) and using the [gap closure workflow](#).

This marks the end of this tutorial.

RNA-Seq de novo transcriptome workflow

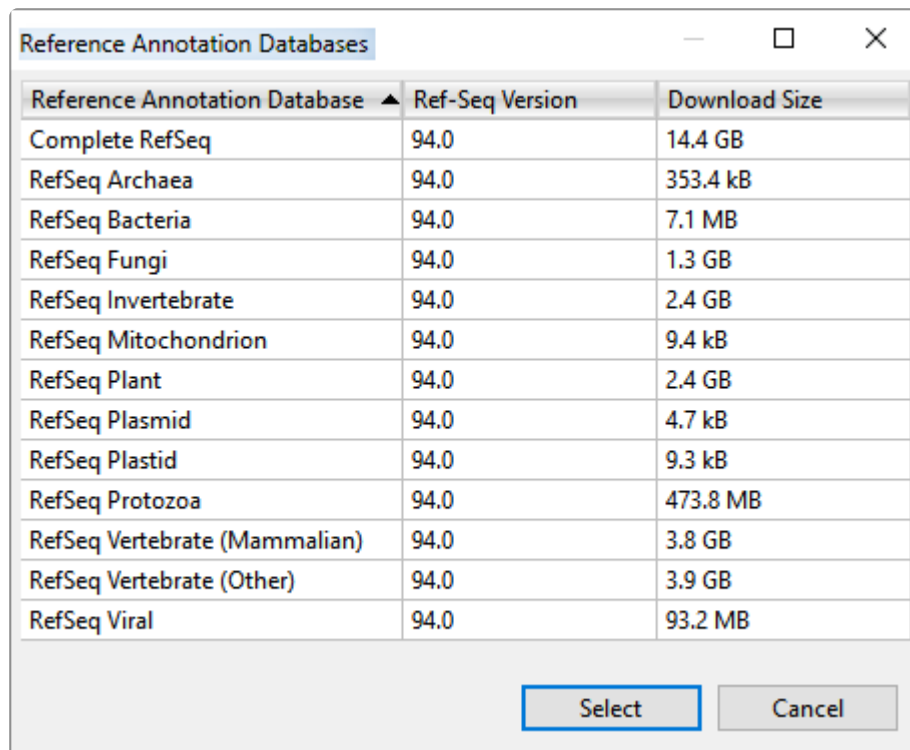
In this tutorial, you will *de novo* assemble an abbreviated set of paired end RNA-Seq sequences from *Saccharomyces cerevisiae* (yeast) from [Nookaew I et al., 2012](#). This workflow uses an abbreviated yeast data set with about 1 million reads per file.

With other applications, *de novo* assembly of RNA-Seq data can potentially result in thousands of unlabeled contigs representing the expressed transcripts. By contrast, SeqMan NGen automatically attempts to group contigs from the same gene, and then name and annotate them based on the best match to a collection of annotated reference sequences (the “Transcript Annotation Database”) extracted from data on NCBI’s [RefSeq](#) website. Results from this workflow are non-quantitative.

Running the transcriptome assembly in SeqMan NGen:

In this part of the tutorial, you will use SeqMan NGen to *de novo* assemble and annotate the RNA-Seq data.

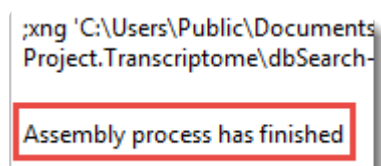
1. Download [T5_RNA-Seq_DeNovo_Transcriptome.zip](#) (147 MB) and extract it to any convenient location (i.e., your desktop). The tutorial data consist of the paired-end reads *Yeast_RNASeq_1Mreads_1.fastq* and *Yeast_RNASeq_1Mreads_2.fastq*.
2. Launch SeqMan NGen and choose **RNA-Seq/Transcriptomics** on the left. On the right, click on **De novo transcriptome**.
3. In the Set Contaminant screen, take the opportunity to verify that you are logged in by looking at the key icon in the bottom left corner. If there is a green check mark, click **Next**. If there is a yellow triangle, click the icon and enter the same login credentials you use for the DNASTAR website. Once you return to the Set Contaminant screen, click **Next**.
4. In the Input Sequences screen, press **Add** and add the *Yeast_RNASeq_1Mreads_1.fastq* and *Yeast_RNASeq_1Mreads_2.fastq* files. Click **Next**.
5. In the Transcript Annotation Database screen, click the **Download Database** button. Choose **RefSeq Fungi** and press **Select**. Then click **Next**.



6. In the Assembly Options screen, click **Next**.
7. In the Assembly Output screen, type *Transcriptome* into the **Project Name** text box, then use the **Browse** button to specify a **Project Folder** for your assembly output files. Click **Next**.
8. In the Run Assembly Project screen, note that:
 - The estimated disk requirement of 2.1 TB is based on the total length of the fungal Transcript Annotation Database, which is 4.2 GB: larger than a human genome. That estimate is based on reference guided genome assemblies that have fixed 50X coverage, not reference guided transcriptome assemblies, which have highly variable coverage. The assembly in this tutorial has extremely low coverage and uses far less disk space than what is estimated here.
 - Cloud Assembly is not offered for the *de novo* transcriptome workflow because most data sets exceed the 48 hour time limit.

Click the link “Run assembly on this computer.” The assembly will take approximately one hour on a standard laptop.

9. Wait until being informed that assembly has finished, then click **Next**.



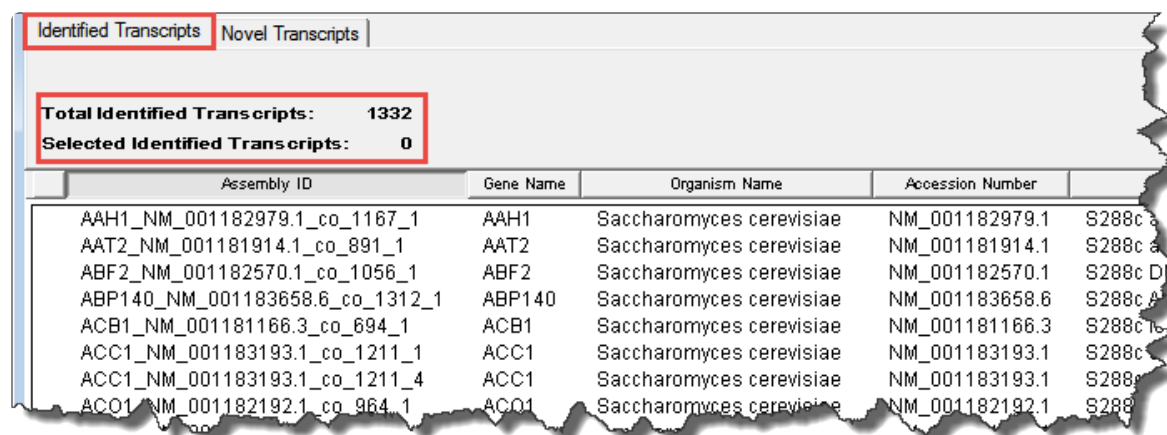
10. In the Assembly Summary screen, note the button **View assembled transcripts**. In the future, this button will allow you to open transcriptome results in SeqMan Ultra. As of version 17, however, *.transcriptome* results can only be analyzed in SeqMan Pro.
11. Click **Finish** to close SeqMan NGen and press **Yes** when prompted.

Viewing transcripts in SeqMan Pro:

During the assembly process, the *de novo* transcriptome assembly output was saved to a package called *Transcript Project.Transcriptome*. Any assembled transcripts with a database match exceeding the specified thresholds were termed “Identified Transcripts,” while assembled transcripts that did not have a database match were called “Novel Transcripts.” This part of the tutorial shows how to load the annotated transcripts into the SeqMan Pro application for downstream analysis.

1. Launch SeqMan Pro and drag and drop the result file *Transcriptome.Transcriptome* from your file explorer onto the SeqMan Pro window.
2. Observe that the ensuing All Transcripts window contains two tabs. Each tab’s heading shows the total number of transcripts in the table, and the number currently selected. The tables in the two tabs support a wide variety of sortable columns which can be displayed or hidden, as desired.

The **Identified Transcripts** tab is active, by default. You should see over 1300 **Total Identified Transcripts**. Since you haven’t yet made any selections, the number of **Selected Identified Transcripts** is zero.



Assembly ID	Gene Name	Organism Name	Accession Number
AAH1_NM_001182979.1_co_1167_1	AAH1	Saccharomyces cerevisiae	NM_001182979.1
AAT2_NM_001181914.1_co_891_1	AAT2	Saccharomyces cerevisiae	NM_001181914.1
ABF2_NM_001182570.1_co_1056_1	ABF2	Saccharomyces cerevisiae	NM_001182570.1
ABP140_NM_001183658.6_co_1312_1	ABP140	Saccharomyces cerevisiae	NM_001183658.6
ACB1_NM_001181166.3_co_694_1	ACB1	Saccharomyces cerevisiae	NM_001181166.3
ACC1_NM_001183193.1_co_1211_1	ACC1	Saccharomyces cerevisiae	NM_001183193.1
ACC1_NM_001183193.1_co_1211_4	ACC1	Saccharomyces cerevisiae	NM_001183193.1
ACO1_NM_001182192.1_co_964_1	ACO1	Saccharomyces cerevisiae	NM_001182192.1

3. Click on the **Novel Transcripts** tab. You should see approximately 50-70 **Total Novel Transcripts**. This table lists the assembled contigs that did not have any match to the Transcript Annotation Database that met the search criteria thresholds and therefore, were not labeled with any match information. Note that this table contains only three columns.

Identified Transcripts		
Novel Transcripts		
Total Novel Transcripts: 56 Selected Novel Transcripts: 0		
Assembly ID	Transcript Length	Assembled Reads
CRH1_NM_001181318.1_co_735_1	1626	627
PPN1_NM_001180760.3_co_614_2	1121	260
RPS22A_NM_001181623.1_co_821_2	534	562
cl_1168_2	262	350
cl_1982_1	963	724
cl_2307_1	311	906
cl_2855_1	354	192
cl_3239_4	537	311
cl_3300_1	453	109
cl_3650_2	326	180
cl_4563_1	285	103

4. Return to the **Identified Transcripts** tab and experiment with the following:
- To show or hide columns - Right-click and choose **Show/Hide Column**, then check or uncheck boxes. Each column is [described in detail](#) in the SeqMan Pro help.

All Transcripts From: Transcriptome_AllTranscripts.searchresults

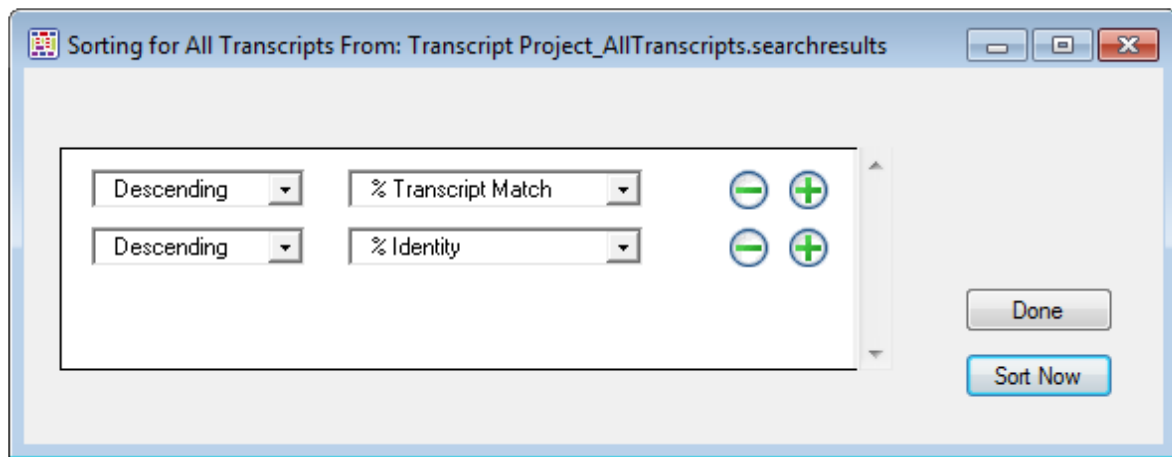
Show Columns:

- ☒
- ☒ Assembly ID
- ☒ Gene Name
- ☒ Organism Name
- ☒ Accession Number
- ☒ Description
- ☒ Database
- ☒ Transcript Length
- ☒ Transcript Start
- ☒ Transcript End
- ☒ % Transcript Match
- ☒ Gene Length
- ☒ % of Full Length
- ☒ Gene Start
- ☒ Gene End
- ☒ % Gene Match
- ☒ % Identity
- ☒ Bit Score
- ☒ eValue
- ☒ Assembled Reads

- To move a column - Use the mouse to drag and drop it in the desired location.
- To sort data in alphabetical or numerical order - Click on the column header that you wish to sort. Note that the resulting groups are also shown in different colors to help visually differentiate between them.

Identified Transcripts		Novel Transcripts			
Total Identified Transcripts:		1322			
Selected Identified Transcripts:		0			
Gene Name	Description	% Transcript Match	% Gene Match	% Identity	
TOS1	S288c Tos1p	78.17 %	100.00 %	100.00 %	
RIB5	S288c riboflavin synthase	92.80 %	70.15 %	100.00 %	
ARO4	S288c 3-deoxy-7-phosphoheptulonate synthase ARO4	87.91 %	100.00 %	100.00 %	
TIM17	S288c protein transporter TIM17	65.98 %	100.00 %	100.00 %	
	S288c hypothetical protein partial mRNA	99.67 %	46.57 %	100.00 %	
APE3	S288c Ape3p	96.30 %	100.00 %	100.00 %	
HYR1	S288c peroxiredoxin HYR1	88.97 %	100.00 %	100.00 %	
YSY6	S288c Ysy6p	67.35 %	100.00 %	100.00 %	
RPL2B	S288c ribosomal 60S subunit protein L2B	63.38 %	23.53 %	100.00 %	
RPA34	S288c DNA-directed RNA polymerase I subunit RPA34	93.34 %	97.86 %	100.00 %	
UMP1	S288c Ump1p	76.80 %	100.00 %	100.00 %	
HIS4	S288c trifunctional histidinol d...sphoribosyl-ATP diphosphatase	92.95 %	100.00 %	100.00 %	
RPS21B	S288c ribosomal 40S subunit protein S21B	55.23 %	100.00 %	100.00 %	
GBP2	S288c Gbp2p	100.00 %	83.64 %	100.00 %	
PGI1	S288c glucose-6-phosphate isomerase	89.56 %	100.00 %	100.00 %	
GDT1	S288c putative ribosome biosynthesis protein GDT1	86.37 %	100.00 %	100.00 %	
RRP7	S288c Rrp7p	92.62 %	47.76 %	100.00 %	
SNL1	S288c Snl1p	97.74 %	90.00 %	100.00 %	
TIF4631	S288c Tif4631p	87.08 %	81.60 %	99.96 %	
CDC48	S288c AAA family ATPase CDC48	95.29 %	100.00 %	99.96 %	
MNN1	S288c alpha-1,3-mannosyltransferase MNN1	90.76 %	100.00 %	99.96 %	
HOM3	S288c aspartate kinase	87.36 %	100.00 %	99.94 %	
	S288c RDN37-2 miscRNA	97.60 %	26.16 %	99.94 %	
ADE17	S288c bifunctional phosphorib...rase/IMP cyclohydrolase ADE17	89.48 %	100.00 %	99.94 %	
PAB1	S288c polyadenylate-binding protein	85.29 %	100.00 %	99.94 %	
SSZ1	S288c Ssz1p	94.12 %	100.00 %	99.94 %	
ATR1	S288c Atr1p	91.31 %	100.00 %	99.94 %	
MET2	S288c homoserine O-acetyltransferase	86.39 %	100.00 %	99.93 %	
COR1	S288c ubiquinol--cytochrome-c reductase subunit COR1	74.86 %	100.00 %	99.93 %	
GAS5	S288c 1,3-beta-glucanosyltransferase	92.03 %	100.00 %	99.93 %	
PFK2	S288c 6-phosphofructokinase subunit beta	92.75 %	100.00 %	99.93 %	
TIM50	S288c protein translocase subunit TIM50	92.81 %	98.32 %	99.93 %	

- To open individual contigs in SeqMan Pro for visualization and editing, double-click on a row of interest to navigate to the corresponding contig assembly. The appropriate .sqd file is loaded with the Alignment view of the selected contig displayed. All the usual visualization and editing tools in SeqMan Pro are available.
- To set stringency thresholds, return to the Identified Transcripts tab and choose **View > Sort**. Set up the dialog to match the image below. To add the second row, click on the plus sign in the original row. Press **Sort Now**.



- To save the sub-set of transcripts that met the stringency thresholds:
 1. In the Identified Transcripts table, click the **%Transcript Match** header to sort the column in decreasing order.
 2. Select all rows where **Transcript Match** is ≥ 99.00 , noting that the **Identity** is also ≥ 99 for those rows. In the header, you will see approximately 130-160 selected transcripts.
 3. Choose **File > Save Selected Transcripts for**. In the Save dialog, designate a name and location for the output file, and then click **Save**. Two files will be created, each with a different extension (*.fas* and *.searchresults*).
 4. (optional) To see what is contained in each of these file types, open the files in any suitable text editor.

This marks the end of this tutorial.

RNA-Seq reference-guided workflow with analysis in ArrayStar

RNA-Seq uses next-gen sequencing to show the presence and quantity of RNA in a genome at a particular moment. DNASTAR's SeqMan NGen application is the starting point for both reference-guided and *de novo* RNA-Seq workflows. Because this tutorial involves a reference-guided workflow, you will then use ArrayStar to analyze the completed RNA-Seq assembly.

In this tutorial, you will compare stationary phase RNA from wild-type *E. coli* cells with that from two different mutant cells, each lacking one of two genes that together encode a transcription factor for flagella and chemotaxis operons. An “operon” is a group of one or more genes that are transcribed as a single RNA unit.

In Part A, you will use SeqMan NGen to create a reference-guided DESeq2-normalized RNA-Seq assembly for two experimental replicates and one wild type control of *E. coli* transcription factors FlhC and FlhD. The data you will assemble is publicly-available single-end Illumina data from *E. coli* ([Fitzgerald et al, 2014](#)). In Parts B & C, you will use ArrayStar to perform downstream analysis and will learn to identify flagella-related operons using two different techniques.

IMPORTANT! There are two options for following this tutorial, one significantly longer than the other.

Item	Full tutorial	Abbreviated tutorial
Description	Assemble in SeqMan NGen, then do downstream analysis in ArrayStar	Do only the downstream analysis in ArrayStar, using a provided project
Download size	4.0 GB (unzips to 16 GB)	4.3 MB
Assembly wait time	1-3 hours	none
Where to begin	Perform the steps in Part A: Setting up the RNA-Seq reference-guided assembly in SeqMan NGen	Read the steps in Part A: Setting up the RNA-Seq reference-guided assembly in SeqMan NGen without performing them

Part A: Setting up the RNA-Seq reference-guided assembly in SeqMan NGen

In this part of the tutorial, you will learn how to set up the project in SeqMan NGen and (optionally) run the assembly. To perform the assembly in Part A, you must download a 4 GB zipped data folder that unpacks to 14.3 GB. We expect most to simply read this section, then proceed to Parts B & C, where you can download a 4 MB data set and perform the downstream analysis part in ArrayStar.

1. (optional) If you plan to run the assembly (most readers will not), download [T1_RNA-Seq.zip](#) (4.0 GB) and extract the contents to any convenient location (e.g., your computer's desktop). The folder contents consist of:
 - The *E. coli* reference sequence: *Escherichia coli str. K-12 substr. MG1655.U00096.gbk*
 - The sample sequences: six folders beginning with *flhC*, *flhD* and *WT*.
2. Launch SeqMan NGen and press **New Assembly**.
3. Select the **RNA-Seq / Transcriptomics** tab on the left and choose **RNA-seq** from the **Quantitative Analysis** section on the right.
4. In the Reference Sequence screen, press the **Add** button and open the *Escherichia coli str. K-12 substr. MG1655.U00096.gbk*. (Alternatively, drag the file from your file explorer and drop it onto the large white space in the middle of the wizard screen.) Click **Next**.
5. In the Input Sequences screen:
 - a. Keep the **Read technology** set at **Illumina**, but uncheck the **paired-end data** box.
 - b. Next to **Experiment setup**, select **Multi-sample with replicates**.
 - c. Use the **Add Folder** button six times, each time adding one of the six sample folders from the tutorial data folder.

Input Sequences

Input sequence files and define experiments or individual replicates

✓ Workflow

✓ Reference Sequence

✓ **Input Sequences**

Set Up Replicate Sets

Input sequences

Read technology: ILLUMINA ☐ Paired-end data

Experiment setup: Multi-sample with replicates ☐ Stranded RNA-Seq reads

Select sequence reads and create replicates

(Replicate sets will be defined on the next page)

Sequence File	Replicate
WT_rep1	WT_rep1
WT_rep2	WT_rep2
flhC_del_rep1	flhC_del_rep1
flhC_del_rep2	flhC_del_rep2
flhD_del_rep1	flhD_del_rep1
flhD_del_rep2	flhD_del_rep2

Add...

Add from Cloud...

Add Folder...

Add Folder from Cloud...

Remove

Create Replicate

Undo Replicate

Auto Name

?

< Back

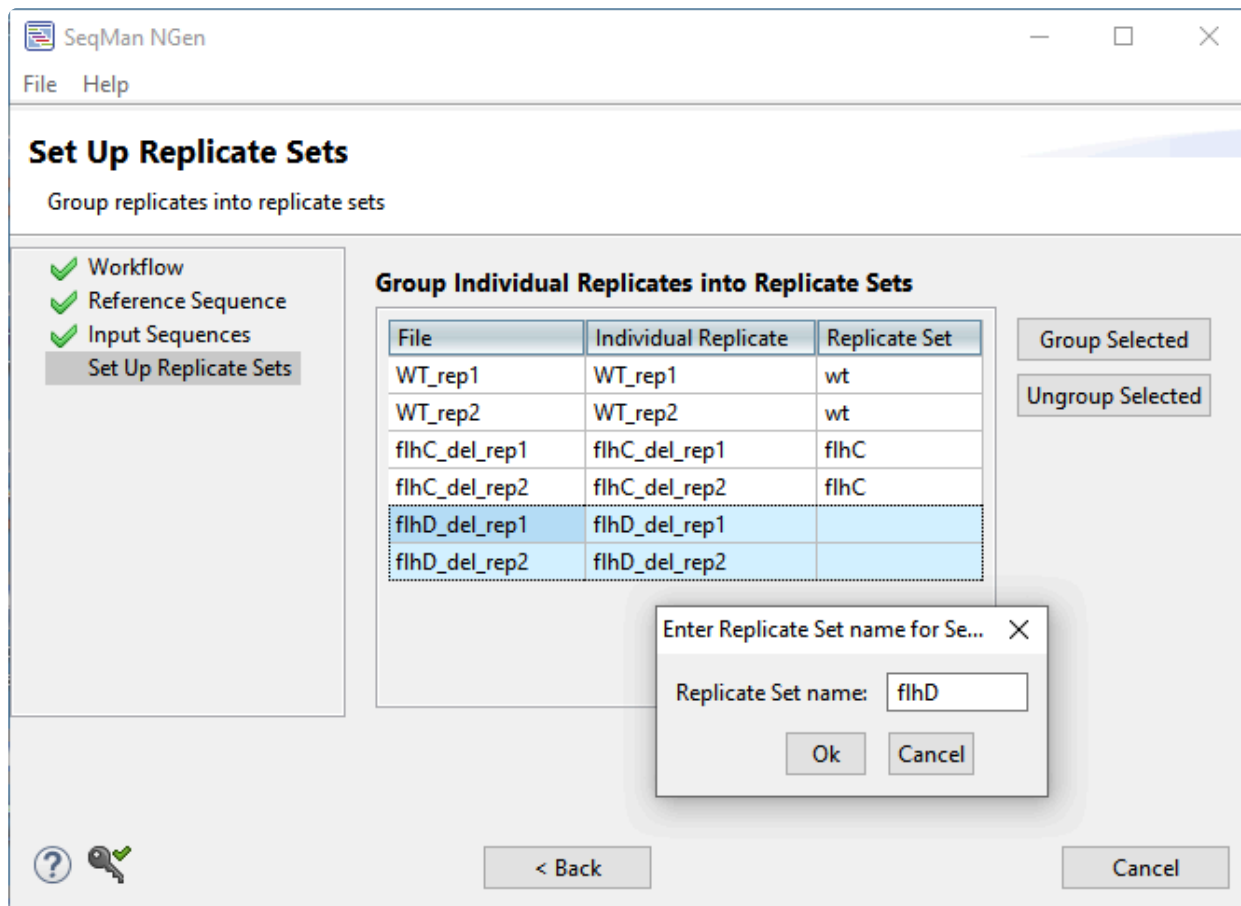
Next >

Cancel

d. Click **Next**.

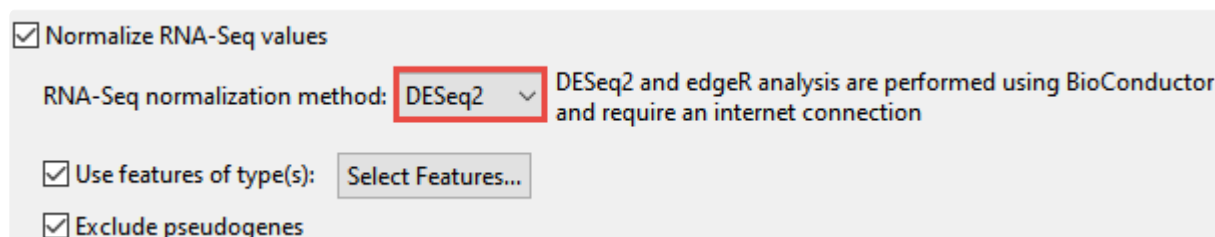
6. In the Set Up Replicate Sets screen:

- Select the two **WT** replicates and click on the **Group Selected** button. In the dialog, name the set "wt" and click **OK**.
- Do the same for the two **flhC** replicates, naming the set "flhC."
- Do the same for the two **flhD** replicates, naming the set "flhD."



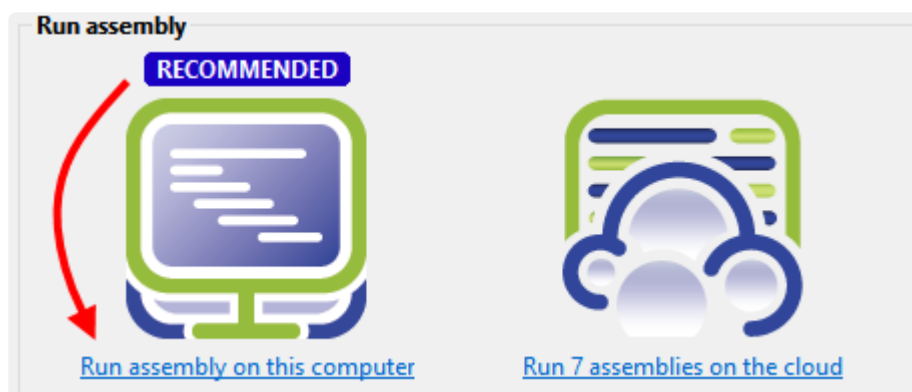
d. Click **Next**.

- In the Set Up Experiments screen, specify the “wt” group as the control. To do this, check the **Is Control** box to the right of “wt,” then click **Next**.
- In the Assembly Options screen, check the box next to **Maximum total reads** and enter 5000000 (5 followed by 6 zeros) to reduce the assembly time. Click **Next**.
- In the Analysis Options screen, make no changes, but observe that the **RNA-Seq normalization method** is **DESeq2**. Click **Next**.



- In the Assembly Output screen, type “Templated RNA-seq” into the **Project Name** text box. This name will be assigned to all output files, including the finished assembly. Use the **Browse** button to specify a **Project Folder** for your assembly output files, then click **Next**.

11. In the Run Assembly Project screen, observe the recommendation for where to run the assembly. Press the link under that recommendation.



12. Wait until being informed that assembly has finished (approximately 10-20 minutes for a local assembly), then click **Next**.
13. From the Assembly Summary screen, click **Analyze differential gene expression**. This launches ArrayStar.
14. (optional) Close SeqMan NGen by clicking the **Finish** button.
15. Within ArrayStar, use **File > Save Project** to save the project as *Templated_RNA-Seq.astar*.

Proceed to [Part B: Analyzing the results in ArrayStar using quick gene sets](#).

Part B: Analyzing the results in ArrayStar using quick gene sets

In [Part A](#) of this tutorial, you set up and ran a templated RNA-Seq assembly using SeqMan NGen. In this part of the tutorial, you will analyze assembly results in ArrayStar using a “quick gene set” and the Gene Table to locate a potential operon structure.

1. If you ran the assembly in Part A, skip to Step 3. Otherwise, download [T1_RNA-Seq \(abbrev\).zip](#) (4.3 MB) and extract it to any convenient location (e.g., your computer’s desktop). The data set consists of a single file, *Templated RNA-Seq.astar*.
2. Double-click on *Templated RNA-Seq.astar* to open the project in ArrayStar.
3. In ArrayStar, select **Graphs > Venn Diagrams** from the menu, and then press the **Quick gene set creation** button.
4. Because this project used DESeq2 normalization, Step 1 is not available. Instead, the dialog opens at **Step 2 - Experiments to compare**. Keep the default (all boxes checked) and press **Move to Step 3 (Comparisons)**.
5. In **Step 3 – Select Comparisons**, set up the filter to find genes in the mutant samples that have a ≥ 5 -fold change, compared to the wild type, and an rlog signal value ≥ 100 . To do this:

- a. Keep the **Signal Threshold (linear)** box checked and change the value from **10.00** to **100.00**.

Note that signal calculation varies depending on which normalization method was selected in SeqMan NGen. In this tutorial, DESeq2 was used, so the Signal is in terms of rlog values. These values are transformed to the linear scale in this dialog and in ArrayStar’s Gene and Isoform tables. If normalization had instead been edgeR or RPKM, the Signal would have been in terms of the TMM or RPKM values, respectively.

- b. Keep the **Fold change (linear)** box checked and change the value from **2.00** to **5.00**.

Quick gene set creation

Step 1 - Choose one Comparison Workflow

Step 2 - 3 experiments selected for comparison.

Step 3 - Select Comparisons

☒ **Signal threshold (linear)** \geq

Compare all genes against this signal threshold to determine if they should be included in the results. When comparing pairs, the higher of the two experiments must exceed this threshold.

☒ **Fold change (linear)** \geq

log₂ fold change (MAP) against wt from DESeq2

A baseline or pairwise workflow is required to use the Fold change comparison.

☒ **P value** \leq

Statistical Test

BH (FDR) adjusted p-values against wt from DESeq2

Calculate the P value for each gene in a pair of replicate set experiments to determine if the gene should be included in the results.

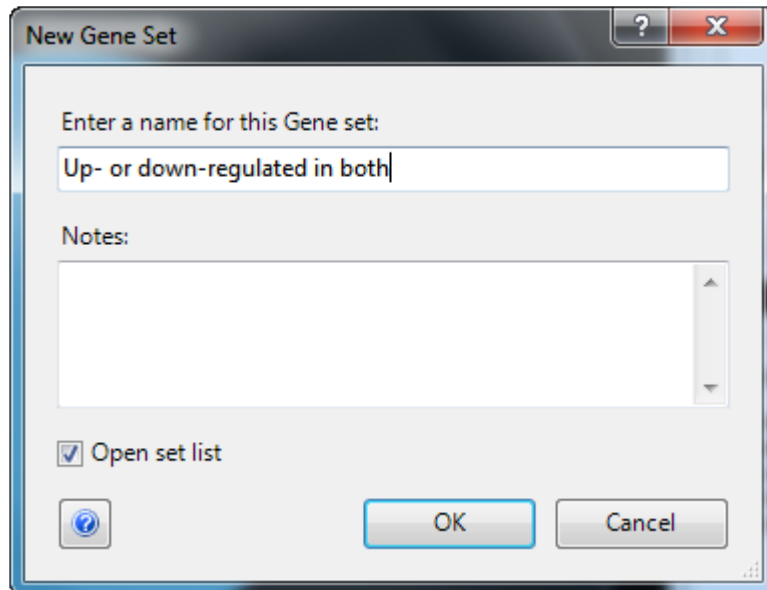
Finish **Cancel**

c. Press **Finish**.

6. In the upper right section of the ArrayStar window, note that two sets have been created. Click the box next to **Global Gene Sets** to select both sets. The Venn Diagram becomes populated with data.

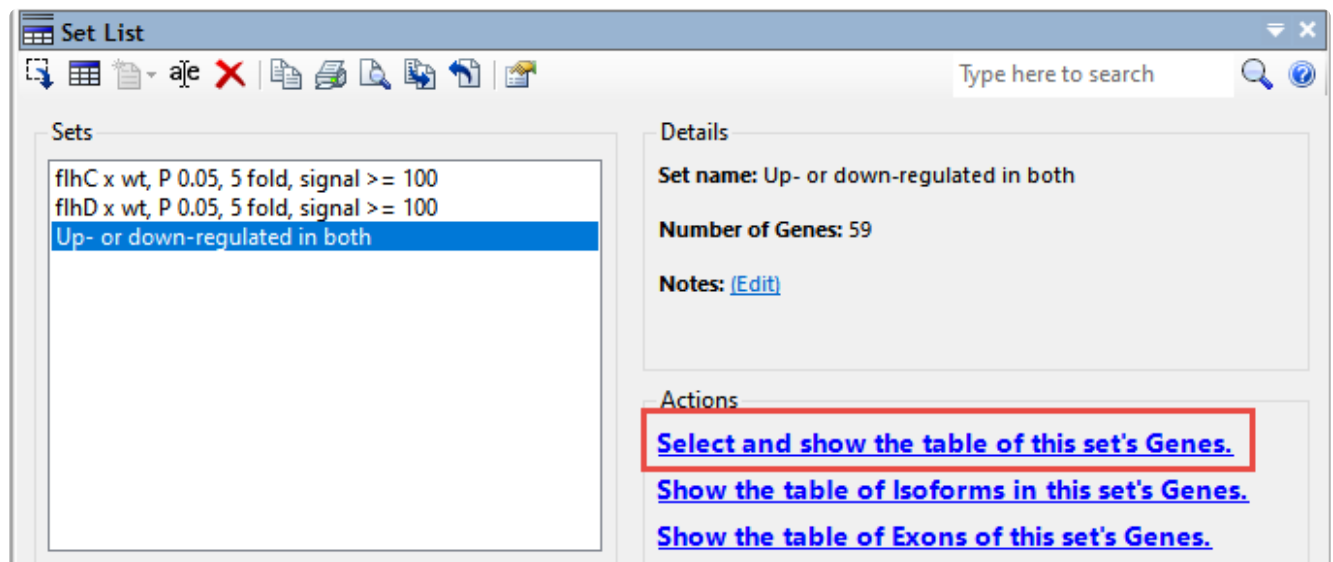
The blue circle labeled “AB” denotes the intersection of the two sets, and represents genes which are up- or down-regulated for both samples: flhC and flhD.

7. Click on **AB** to select it, and then press the link in the bottom right of the window: **Remember the selection as a Gene set**.
8. When prompted to enter a name for the gene set, type “Up- or down-regulated in both” and press **OK**.

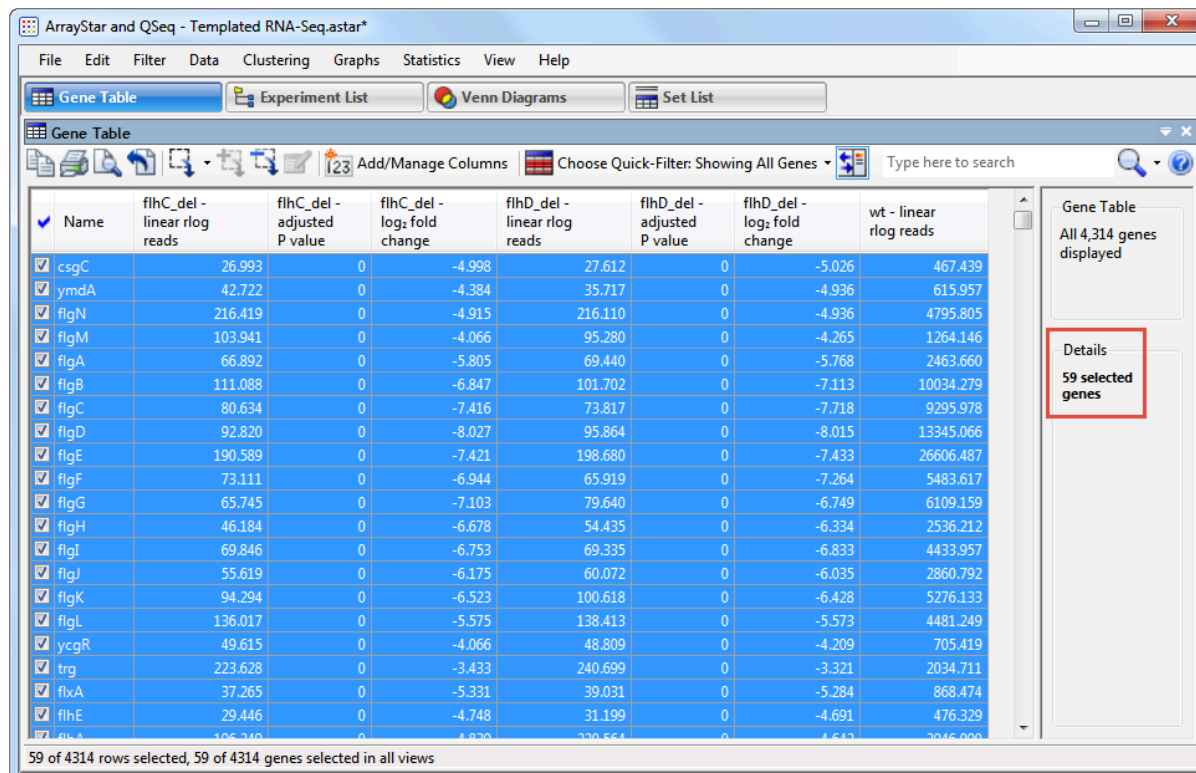




The Set List view opens, with the new set already selected.

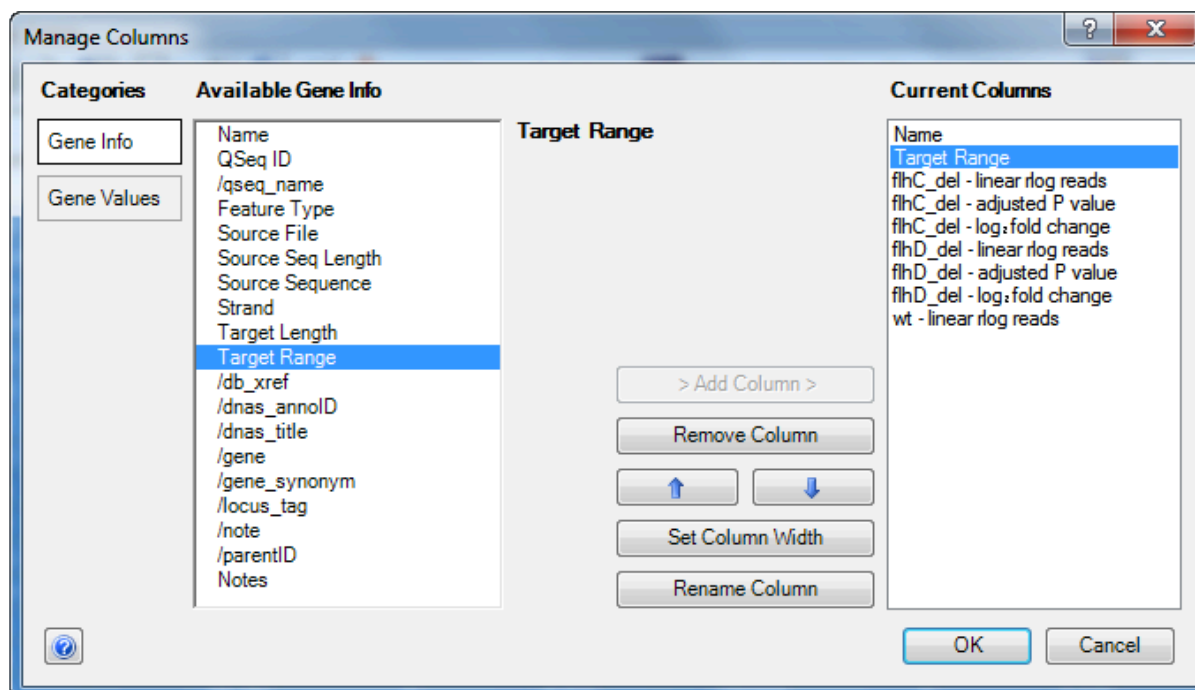
9. Click the link **Select and show the table of this set's Genes.**



The Gene Table opens with 59 genes selected. Do NOT click anywhere in the table for now, as you need to retain the selection (blue highlighting).



10. Click the **Add/Manage Columns** tool ( **Add/Manage Columns**) to open the Manage Columns dialog:
- The **Gene Info** button is active by default. Select **Target Range**, then press the **> Add Column >** button to add the item to the Current Columns list.
 - Press the up arrow button () repeatedly until **Target Range** appears just after **Name**.



c. Click **OK** to close the Manage Columns dialog and return to the Gene Table.


11. In the Gene Table, click once on the column header for **Target Range** to sort all of the genes in the project in ascending order of appearance on the assembly.
12. To identify possible operon structures, scroll down the Gene Table. Genes with blue highlighting represent the downregulated genes in the “quick gene set.” Large blue-highlighted sections with overlapping or consecutive genes that show similar trends in expression levels and fold changes are candidates for operons.

The first operon candidate begins at **Target Range 1129830..1129414**.

Name	Target Range	flhC_del - linear rlog reads	flhC_del - adjusted P value	flhC_del - log ₂ fold change	flhD_del - linear rlog reads	flhD_del - adjusted P value	flhD_del - log ₂ fold change
murJ	1127839..1129374	340.597	0.0198	-0.590	275.798	0.0000719	-0.970
flgN	1129830..1129414	216.419	0	-4.915	216.110	0	-4.936
flgM	1130128..1129835	103.941	0	-4.066	95.280	0	-4.265
flgA	1130863..1130204	66.892	0	-5.805	69.440	0	-5.768
flgB	1131018..1131434	111.088	0	-6.847	101.702	0	-7.113
flgC	1131438..1131842	80.634	0	-7.416	73.817	0	-7.718
flgD	1131854..1132549	92.820	0	-8.027	95.864	0	-8.015
flgE	1132574..1133782	190.589	0	-7.421	198.680	0	-7.433
flgF	1133802..1134557	73.111	0	-6.944	65.919	0	-7.264
flgG	1134729..1135511	65.745	0	-7.103	79.640	0	-6.749
flgH	1135564..1136262	46.184	0	-6.678	54.435	0	-6.334
flgI	1136274..1137371	69.846	0	-6.753	69.335	0	-6.833
flgJ	1137371..1138312	55.619	0	-6.175	60.072	0	-6.035
flgK	1138378..1140021	94.294	0	-6.523	100.618	0	-6.428
flgL	1140033..1140986	136.017	0	-5.575	138.413	0	-5.573
rne	1144367..1141182	13417.958	0.941	-0.017	13218.518	0.853	-0.042
yceQ	1144502..1144822	2990.779	0.999	0.001	2955.350	0.961	-0.017
rluC	1144940..1145899	1283.232	0.774	0.095	1084.979	0.458	-0.190
psrD	1146589..1146757	1239.108	0.724	0.120	1107.895	0.847	-0.070
yceF	1146595..1146011	553.776	0.206	-0.205	571.210	0.208	-0.225

Note that the P-values for both flhC and flhD are 0 in every case, signifying extremely high confidence. Also note that the log₂ fold changes for both flhC and flhD are all between approximately -4 and -8, indicating strong downward regulation for both. By comparison, less than 1.5% of the 4,300+ genes have a log₂ fold change less than -2.

Other candidates for operons begin at: **1962972..1962580**; **2001024..200473**; and **2013014..2012700**. In all cases, all P-values are zero, and log₂ fold changes are similar to those in the first operon.

13. To see only the genes in the ‘quick gene set,’ click on the **Choose Quick-Filter** tool ( **Choose Quick-Filter: Showing All Genes**) above the table, and select **Show Only Gene Set**. In the pop-up, select ‘Up- or down-regulated in both’* and press **OK**.

14. Click once on the **Name** column header to sort alphabetically by gene name.

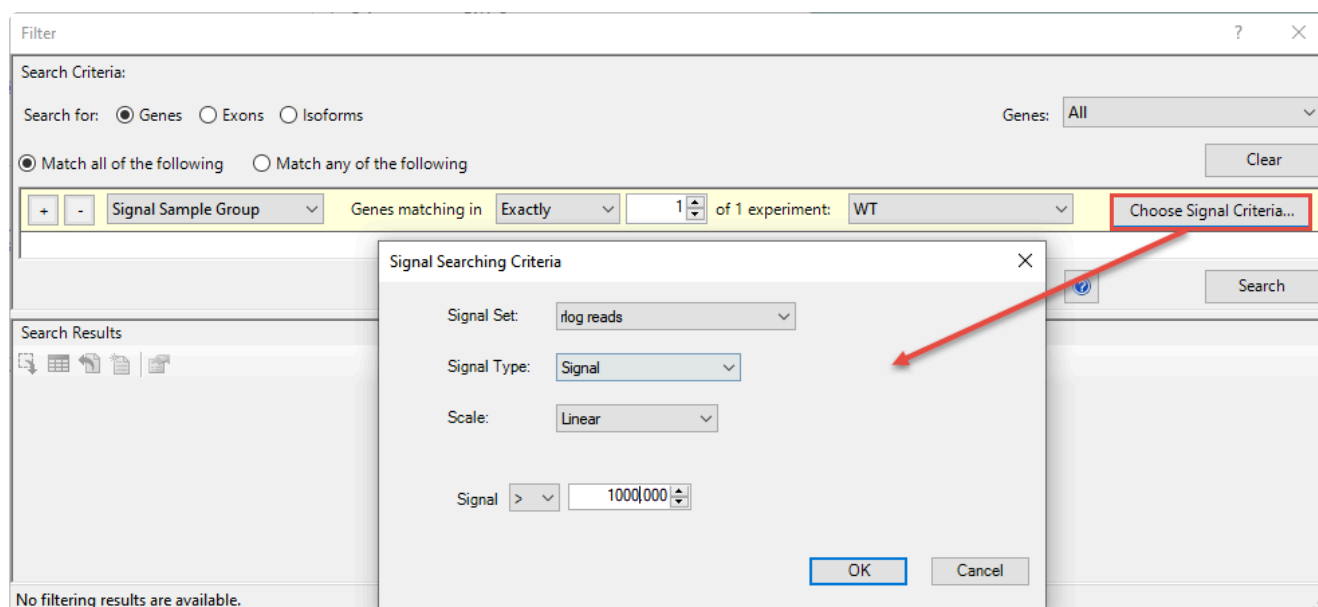
Observe that most of the down-regulated genes have a name prefix of '*fli*', '*flg*' or '*che*.' The first two gene types encode flagella-related proteins while the third type encodes chemotaxis-related proteins. These findings corroborate that these samples are missing coding genes related to flagellar and chemotaxis functions.

Proceed to [Part C: Analyzing the results in ArrayStar using advanced filtering](#).

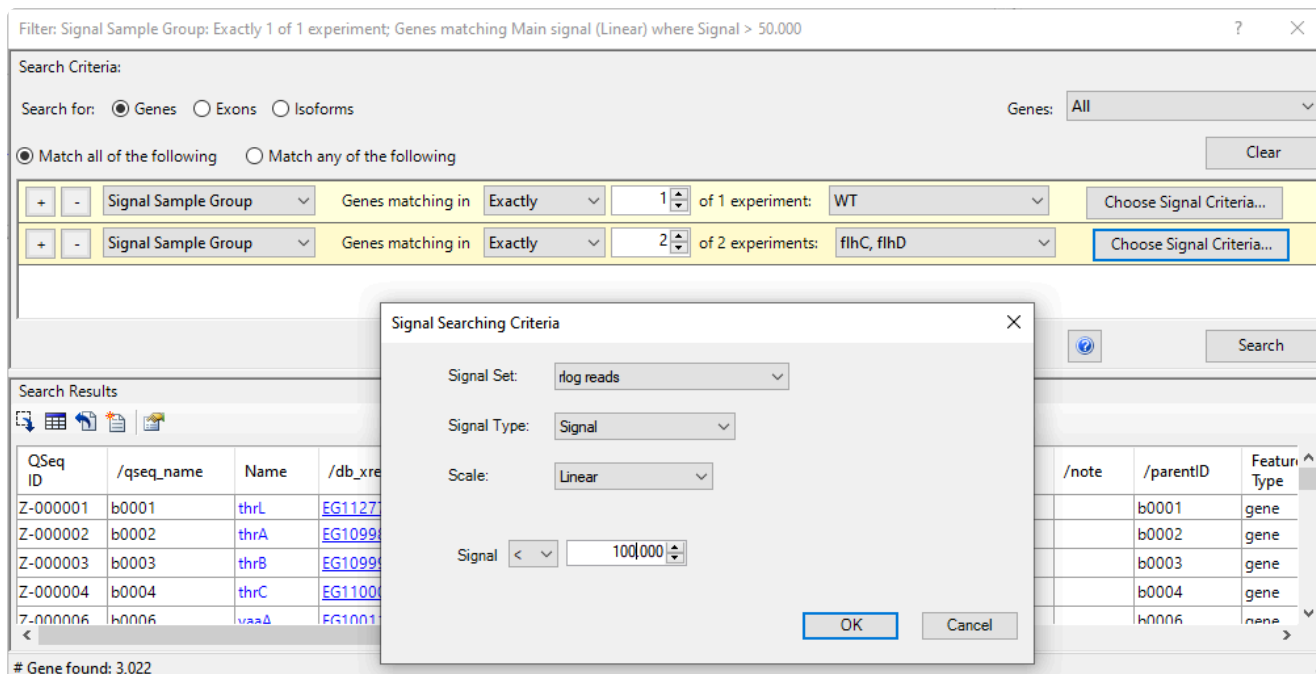
Part C: Analyzing the results in ArrayStar using advanced filtering

In [Part A](#) of this tutorial, you set up and ran a templated RNA-Seq assembly using SeqMan NGen. In this part of the tutorial, you will analyze the assembly results using ArrayStar's advanced filtering functionality to create a set of genes that relate to the flagella structure.

1. Choose **Filter > Filter All** to open the advanced Filter dialog.
2. In the yellow row, select options to match the image below. Note that you must press the **Choose Signal Criteria** button to open the pop-up dialog shown in the image. In the pop-up, change the **Scale** to **Linear** and the **Signal** to **> 1000** (note greater-than sign). Press **OK** to close the pop-up.

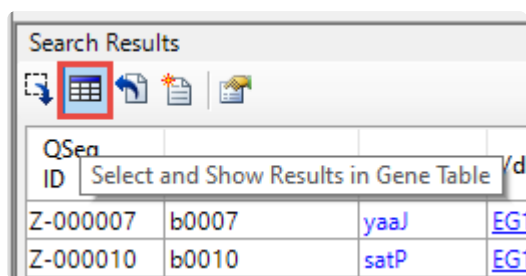


3. Click the plus sign (+) to the left of the yellow row to add a second row. In the new row, set up options to match the image below. In the **Choose Signal Criteria** pop-up, change the **Scale** to **Linear** and the **Signal** to **< 100**. Press **OK** to close the pop-up.



Press **Search**. This two-filter combination results in approximately 25-35 genes.

- Press the **Select and show results in Gene table** tool.



Observe that most of the genes that met the filter thresholds are *fli* and *flg* genes. These genes are involved in producing flagella in *E. coli* and are known to be regulated by the transcription factor encoded by the *flhC* and *flhD* genes. The *che*, *tap* and *tar* genes are involved in chemotaxis and are also known to be regulated by *FlhCD*.

- Use the **Quick filter** tool in the Gene Table header to **Show all genes**.
- Sort by the adjusted p-values from smallest to largest by clicking either of the p-value column headers.

✓	Name	flhC - linear rlog reads	flhC - adjusted P value	flhC - log ₂ fold change	flhD - linear rlog reads	flhD - adjusted P value	flhD - log ₂ fold change	WT - linear rlog reads
---	------	--------------------------	-------------------------	-------------------------------------	--------------------------	-------------------------	-------------------------------------	------------------------

Note that there are approximately 40-60 genes where the adjusted p-value is zero for both *flhC* and *flhD*. These include the genes found by filtering (still highlighted in blue), as well as additional flagella

and chemotaxis related genes.

123 Add/Manage Columns Choose Quick-Filter: Showing All Genes								
<input checked="" type="checkbox"/>	Name	flhC - linear rlog reads	flhC - adjusted P value	flhC - log ₂ fold change	flhD - linear rlog reads	flhD - adjusted P value	flhD - log ₂ fold change	WT - linear rlog reads
<input checked="" type="checkbox"/>	fliZ	39.924	0	-6.078	42.368	0	-6.005	1643.292
<input checked="" type="checkbox"/>	tar	33.124	0	-6.365	35.279	0	-6.302	1580.940
<input checked="" type="checkbox"/>	flgG	35.011	0	-6.936	40.163	0	-6.630	2399.959
<input type="checkbox"/>	flgH	27.004	0	-6.030	29.191	0	-5.904	960.120
<input checked="" type="checkbox"/>	cheW	55.532	0	-5.159	39.145	0	-6.075	1410.894
<input checked="" type="checkbox"/>	flgI	35.957	0	-6.337	36.607	0	-6.398	1742.687
<input type="checkbox"/>	cheA	107.454	0	-5.490	72.177	0	-6.360	3547.725
<input checked="" type="checkbox"/>	flgJ	32.227	0	-5.856	33.422	0	-5.822	1112.089
<input checked="" type="checkbox"/>	flgK	45.033	0	-6.202	47.428	0	-6.144	2024.721
<input checked="" type="checkbox"/>	motB	52.468	0	-5.166	40.574	0	-5.884	1380.328
<input checked="" type="checkbox"/>	flgE	79.665	0	-7.618	85.667	0	-7.544	10458.105
<input type="checkbox"/>	motA	50.079	0	-4.463	32.257	0	-5.686	916.399
<input type="checkbox"/>	trg	87.620	0	-3.287	102.989	0	-3.319	766.896
<input checked="" type="checkbox"/>	flgD	45.129	0	-7.723	47.125	0	-7.696	5270.712
<input checked="" type="checkbox"/>	flgF	33.210	0	-6.822	32.441	0	-7.001	2104.961
<input checked="" type="checkbox"/>	flgC	36.982	0	-7.483	34.534	0	-7.798	3651.014
<input checked="" type="checkbox"/>	fliI	26.213	0	-6.355	25.460	0	-6.573	1158.409
<input type="checkbox"/>	fliJ	24.147	0	-5.539	25.480	0	-5.477	615.607
<input checked="" type="checkbox"/>	fliK	43.963	0	-6.129	50.579	0	-5.857	1841.078
<input checked="" type="checkbox"/>	fliL	34.795	0	-6.364	37.487	0	-6.248	1626.622
<input checked="" type="checkbox"/>	fliF	63.726	0	-6.824	59.129	0	-7.060	4541.831
<input checked="" type="checkbox"/>	fliM	41.475	0	-6.230	41.429	0	-6.310	1910.774
<input checked="" type="checkbox"/>	flhA	84.942	0	-4.659	96.517	0	-4.432	1508.108

This marks the end of this tutorial.

ChIP-Seq workflow with analysis in ArrayStar

Lasergene Genomics' ChIP-Seq analysis workflow enables you to locate the binding sites of DNA-associated proteins and determine how these proteins interact with the DNA to affect expression in nearby genes.

In this tutorial, you will use SeqMan NGen to create a ChIP-Seq assembly for two experimental replicates and one wild type control of *E. coli* transcription factor FlhD. Once the assembly is complete, you will import the data into ArrayStar's ChIP-Seq workflow and filter the results for IP fragments that match in both replicates, and that meet a specified signal threshold. Finally, you will look for genes that the fragments intersect, and see whether they share any common characteristics.

Setting up and running the assembly in SeqMan NGen:

1. Download [T2_ChIP-Seq.zip](#) (2.4 GB) and extract the contents to any convenient location (e.g., your computer's desktop). The folder contents consist of:
 - The *E. coli* reference sequence: *Escherichia coli str. K-12 substr. MG1655.U00096.gbk*
 - The sample sequences: six folders beginning with *flhC*, *flhD* and *WT*.
2. Launch SeqMan NGen and choose **New Assembly**.
3. Choose **RNA-Seq / Transcriptomics** on the left, then choose **ChIP-Seq** on the right.
4. In the Reference Sequence screen, click **Add** and add the reference sequence *MG1655_U00096.3.gb*. Press **Next**.
5. In the Input Sequences screen, uncheck the box next to **Paired-end reads**. Press the **Add Folder** button to add the folder starting with *WT*. Do the same for each of the next two sample folders. Click **Next**.
6. In the Set Up Experiments screen, check the box to the right of **WT_rep1**, indicating that it is the control. Press **Next**.
7. In the Assembly Options screen, check the box next to **Maximum total reads** and enter 5000000 (5 followed by 6 zeros) to reduce the assembly time. Click **Next**.
8. In the Analysis Options screen, leave the **ChIP-Seq detection method** at **MACS** and press **Next**.
9. In the Define Binding Proteins screen, use the **Known binding site motif** drop-down menu to select **Transcription Factor Database**. Use the **Organism** drop-down menu to select **Bacteria**. Press **Select** and choose **FlhCD_CS / FlhCD**, then press **OK**.

Define Binding Proteins

Known binding site motif: Transcription Factor Database ▼

Organism: Bacteria ▼

Site name/Factor name: FlhCD_CS / FlhCD Select...

PubMed ID: [12181488](#)

Binding Protein Label: FlhCD

Press **Next**.

10. In the Assembly Output screen, type “ChIP-seq” into the **Project Name** text box. Use the **Browse** button to specify a **Project Folder** for your assembly output files, then press **Next**.
11. In the Run Assembly Project screen, check whether the recommendation is to run the assembly locally or on the cloud. Press the corresponding link to begin assembly. Typical assembly times are 20 minutes (cloud) or 30-45 minutes (local).
12. Wait until being informed that assembly has finished, then click **Next**.
13. From the Assembly Summary screen, click **Analyze peaks** to open the results in ArrayStar.

Analyzing the results in ArrayStar:

1. In ArrayStar, click on the Fragment Table tab and ensure that approximately 40 IP fragments were found.

1 of 40 rows selected, 0 of 40 IP fragments selected in all views

2. Similarly, click on the Peak Table tab and ensure that about 50-60 peaks were found.
3. To filter for fragments found in both replicate samples, choose **Filter > Filter All** and make selections to match the image below.

Filter

Search Criteria:

☐ Genes
 ☒ IP Fragments
 ☐ IP Peaks

Search for:
 ☐ Exons
 ☐ Isoforms

☒ Match all of the following
 ☐ Match any of the following

+

-

Signal Sample Group ▼

IP Fragments matching in

Exactly ▼

2 ▼

of 2 experiments:

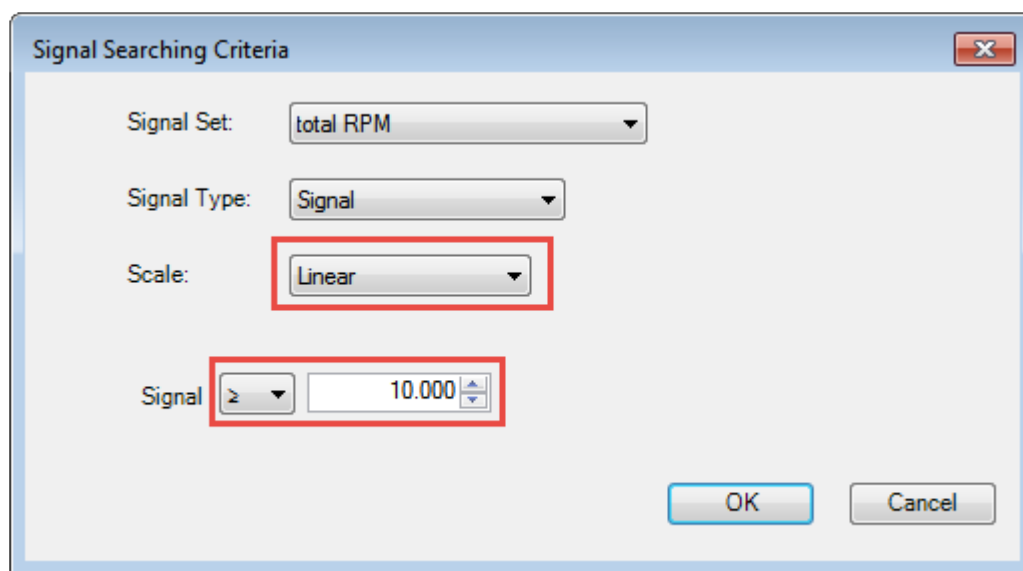
FlhD_3xFLAG_rep1, FlhD_3x ▼


Choose Signal Criteria...

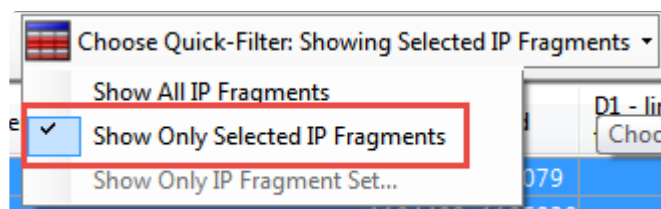
IP Fragments: All ▼

Clear

- Press the **Choose Signal Criteria** button and make selections to match the following image. Press **OK**.



- Press **Search**. The search should yield about 17 IP fragments.
- Press the **Select and Show Results in IP Fragment Table** tool () , located above the result table.
- In the Fragment Table, use the **Quick-Filter** tool to choose **Show Only Selected IP Fragments**.



- Click on **Add/Manage Columns**. In the IP Fragment Info section, choose **Intersecting Genes**, then **>Add Column>**, then **OK**.
- Look at the newly-added column and note that many of the fragments intersect genes beginning with "flg", "fli" and "flh".

Intersecting Genes

ampH, sbmA
csgA, csgC, ymdA
flgA, flgB
yciK, sohB
flhB, cheZ
azuC, yecR, ftnA
fliA, fliC
fliE, fliF
fliK, fliL
ygdH, sdaC
ygeR
ppiA, tsgA
gntK, gntR
gntR, yhhW
mdtF
yhjV, dppF
gltP, yjcO

These are the major genes for producing flagella in *E. coli* and are known to be regulated by the transcription factor encoded by the *flhC* and *flhD* genes. The *che* gene is involved in chemotaxis and are also known to be regulated by *FlhCD*.

This marks the end of this tutorial.

Copy number variation (CNV) workflow with analysis in ArrayStar and GenVision Pro

Copy-number variation (CNV) is defined as genomic regions that have been repeated one or more times and these variations play an important role in normal genetic variation and in some diseases. DNASTAR's CNV workflow is used to analyze genomic variation by considering changes over a region, as indicated by deleted or duplicated gene copies.

In this tutorial, you will run a set of paired-end Illumina sequence files from MG1655 *E. coli* against the DH10B reference genome. MG1655 has several deleted regions compared to the DH10B reference sequence. In addition, the reference sequence contains a large duplication.

Begin with [Part A: Setting up the CNV project in SeqMan NGen](#).

Part A: Setting up the CNV project in SeqMan NGen

In this part of the tutorial, you will use the SeqMan NGen wizard to import data and run the assembly. You will then press a button to open the results in ArrayStar.

1. Download [T3_CNV.zip](#) (1.4 GB) and extract the contents to any convenient location (e.g., your computer's desktop). The folder contains the following sequences:
 - Reference sequence *DH10B_NC010473.gbk*
 - Paired end sample sequences *SRR1284938_1.fastq* and *SRR1284938_2.fastq*
2. Launch SeqMan NGen and choose **New Assembly**.
3. Choose **Variant Analysis / Resequencing** on the left, then **Whole genome** under **NGS-Based** on the right.
4. In the Reference Sequence screen, press **Add** and add the file *DH10B_NC010473.gbk*. Click **Next**.
5. In the Input Sequences screen, press **Add** and add the paired reads *SRR1284938_1.fastq* and *SRR1284938_2.fastq*. Use the **Experimental setup** menu to choose **Single sample**, then click **Next**.
6. In the Assembly Options screen, click **Next**.
7. In the Analysis Options screen, choose **Haploid**, since this is a bacterial genome. Check the box next to **Detect CNVs** and keep the default method of **RPK_CN**. Click **Next**.
8. In the Assembly Output screen:
 - a. Type "CNV" into the **Project Name** text box. Use the **Browse** button to specify a **Project Folder** for your assembly output files. Click **Next**.
9. In the Run Assembly Project screen, note that the **Estimated coverage** is 375X. A coverage of 50-100X is adequate and additional coverage simply slows down the assembly.
10. To reduce coverage, click **Assembly Options** on the left. Check the box next to **Maximum total reads** and type in 1600000 (16 followed by 5 zeros). Then click **Run Assembly Project** to return to that screen. Note that the **Estimated coverage** is now near 50X.
11. Note the recommendation to assemble on your local computer or on the cloud and click the corresponding link to begin assembly. Local assembly should take approximately 30 minutes.

12. After being informed that assembly has finished, click **Next**.
13. From the Project Report screen, click **Analyze and compare variants** to open the project in ArrayStar.
14. After ArrayStar opens, Use **File > Save Project** to save the project as *CNV.astar*.
15. Close the SeqMan NGen project by clicking the **Finish** button and confirming you would like to close the application.

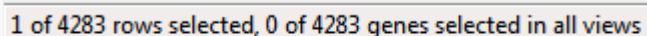
Proceed to [Part B: Finding a putative duplication in the reference sequence using ArrayStar](#).

Part B: Finding a putative duplication in the reference sequence using ArrayStar


In [Part A](#) of this tutorial, you ran an assembly and launched the results in ArrayStar. In this part, you will use the ArrayStar Gene Table to locate potential duplications in the reference sequence.

Imagine that you would like to find a region that is repeated in the reference, but present only once in the MG1655 sample. In the Gene Table these regions will have a linear weighted RPKM-CN of approximately 0.5.

1. Click on the Gene Table tab near the top of the ArrayStar window. The footer shows that the Gene Table contains 4,283 genes.




1 of 4283 rows selected, 0 of 4283 genes selected in all views

2. Click on the header **CNV (2) – linear weighted RPKM-CN** to sort the column from small to large values
3. Drag the mouse to select all table rows with a linear weighted RPKM-CN from 0.4-0.6, inclusive.
4. Right-click anywhere on the highlighted area and choose **Select and Remember as a Gene Set**.
5. Type in the name **Duplicated in reference**, and then press **OK**.
6. Return to the Gene Table by clicking its tab near the top of the ArrayStar window.
7. Without disturbing the selection, click on the **Add/Manage Columns** tool ( **Add/Manage Columns**). Select **Target Range** and press the **>Add Column>** button, then click **OK**.
8. Order the genes by ascending location by clicking on the **Target Range** column header.
9. Use the vertical scrollbar to scroll down to the large block of selected rows beginning around Target Range **515000**.

<input checked="" type="checkbox"/>	Name	CNV (2) - linear weighted RPKM-CN	Target Range ▲
<input type="checkbox"/>	ybcN	1.133	511021..511476
<input type="checkbox"/>	ninE	0.760	511476..511646
<input type="checkbox"/>	ybcO	0.779	511639..511929
<input type="checkbox"/>	rusA	0.882	511926..512288
<input type="checkbox"/>	ylcG	0.686	512285..512425
<input type="checkbox"/>	ybcQ	1.207	512511..512894
<input type="checkbox"/>	insH-2	0.715	514308..513292
<input checked="" type="checkbox"/>	essD	0.509	515953..516168
<input checked="" type="checkbox"/>	ybcS	0.542	516168..516665
<input checked="" type="checkbox"/>	rzpD	0.527	516662..517123
<input type="checkbox"/>	rzoD	0.114	516882..517064
<input checked="" type="checkbox"/>	borD	0.533	517448..517155
<input checked="" type="checkbox"/>	ybcV	0.462	518149..517739
<input type="checkbox"/>	ybcW	0.322	518435..518641
<input checked="" type="checkbox"/>	nohB	0.464	519389..519934
<input type="checkbox"/>	appY	0.382	522236..522985
<input checked="" type="checkbox"/>	ompT	0.573	524188..523235
<input checked="" type="checkbox"/>	envY	0.481	525463..524702
<input checked="" type="checkbox"/>	ybcH	0.415	526536..525646
<input checked="" type="checkbox"/>	nfrA	0.538	529509..526537
<input checked="" type="checkbox"/>	nfrB	0.554	531733..529496
<input checked="" type="checkbox"/>	cusS	0.522	533325..531883
<input checked="" type="checkbox"/>	cusR	0.414	533998..533315
<input checked="" type="checkbox"/>	cusC	0.442	534155..535528
<input checked="" type="checkbox"/>	cusF	0.440	535686..536018
<input checked="" type="checkbox"/>	cusB	0.474	536034..537257
<input checked="" type="checkbox"/>	cusA	0.553	537269..540412
<input checked="" type="checkbox"/>	pheP	0.583	540514..541890
<input checked="" type="checkbox"/>	ybdG	0.532	543218..541971
<input checked="" type="checkbox"/>	nfnB	0.596	543979..543326
<input checked="" type="checkbox"/>	ybdF	0.497	544441..544073

This area marks a possible duplication in the DH10B reference sequence.

- From the **Choose Quick Filter** menu ( Choose Quick-Filter: Showing All Genes) choose **Show Only Gene Set**. Then select **Duplicated in reference** and press **OK**.

The Gene Table now contains only the ~230 putative duplicated genes.

- Use **File > Save Project** to save updates to the project.

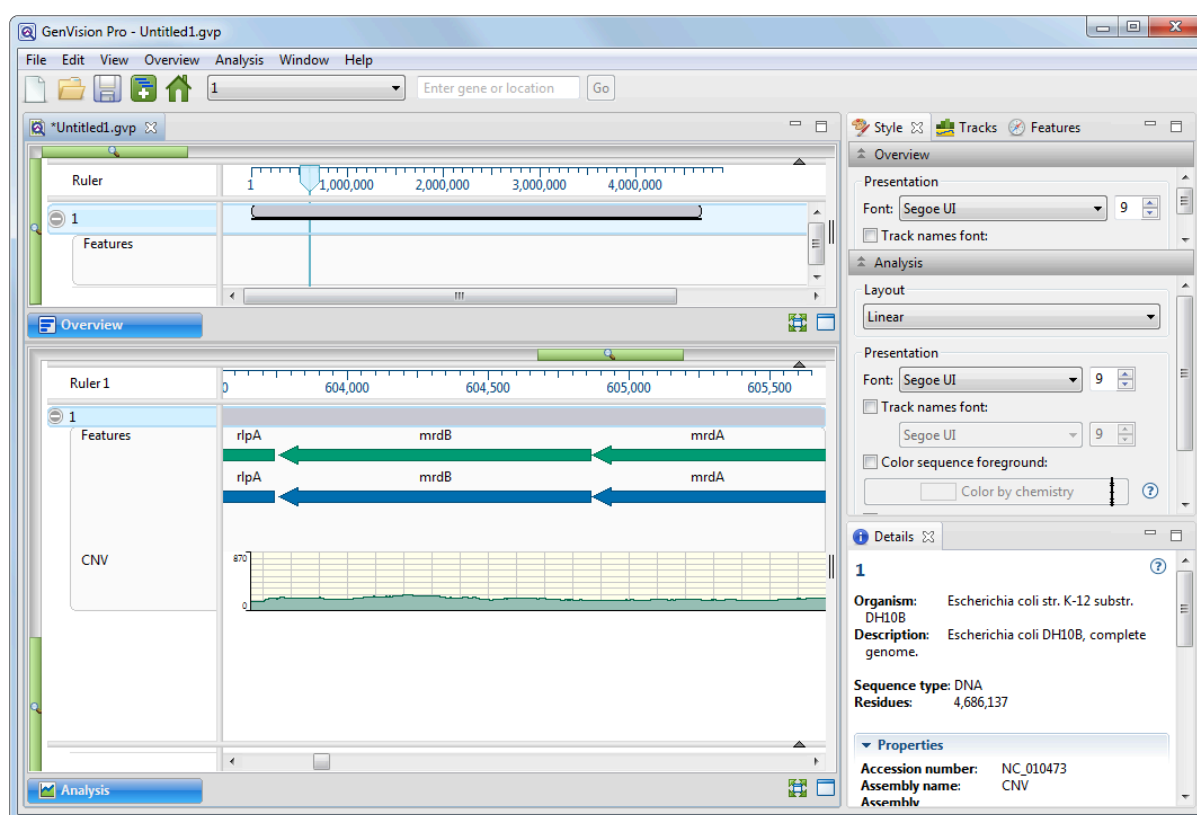
Proceed to [Part C: Confirming the duplication using GenVision Pro](#).

Note: The following brief video is not part of this tutorial, but also explores copy-number variation by loading SeqMan NGen assemblies into ArrayStar:

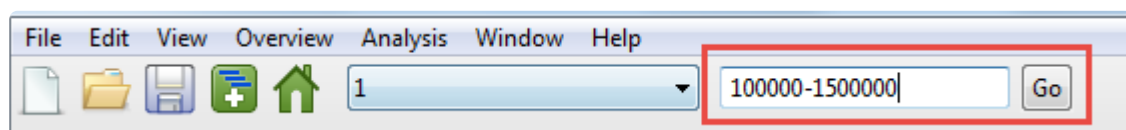
Part C: Confirming the duplication using GenVision Pro

In [Part B](#) of this tutorial, you located a putative duplication in the reference sequence using ArrayStar. Deletions or duplications can be confirmed graphically by sending them to SeqMan Pro or GenVision Pro. In this section, you will use GenVision Pro to view a graphical representation confirming the putative duplication found in the previous section.

1. In the Gene Table, click on a random row to select it. Then right-click on the row and choose **Send Selection to GenVision Pro**. GenVision Pro launches with the display centered on the target range you selected in Part B.



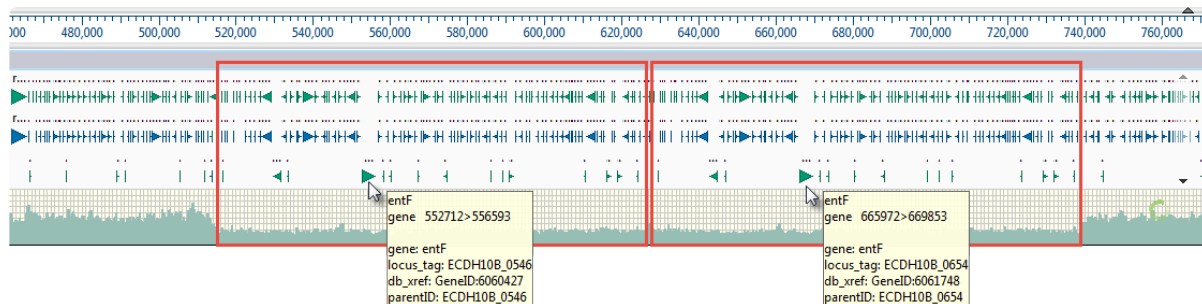
2. In the navigation area near the top of the GenVision Pro window, enter the range **100000-1500000** and press **Go**.



3. Slide the Analysis View's green horizontal slider to the left until you can see the range from approximately 100,000 to 1,500,000.



4. Look at the green coverage graph (“CNV”) at the bottom of the view. Note how the coverage drops in the area from 500,000-700,000. This indicates a likely duplication in the reference sequence. The sample reads have been divided between the two copies, resulting in half the expected coverage in this range. (To the right, you can also see a deletion.)
5. Zoom in on the putative duplication until it fills the entire window. Even though the features are very small at this zoom level, compare the feature sizes and relative positions. It appears there are two duplicate sets of features, side-by-side.
6. Hover over “sibling” features with the mouse, one at a time, to confirm that the feature names match. This further corroborates that the region has been duplicated in the reference sequence.



This marks the end of this tutorial.

Whole genome reference-guided workflow with analysis in ArrayStar

In this tutorial, you will look for deleterious genes in the Caucasian/Utah/Mormon father-mother-daughter trio data from the [NIST Genome in a Bottle](#) project. This is commonly referred to as the “CEPH Trio.”

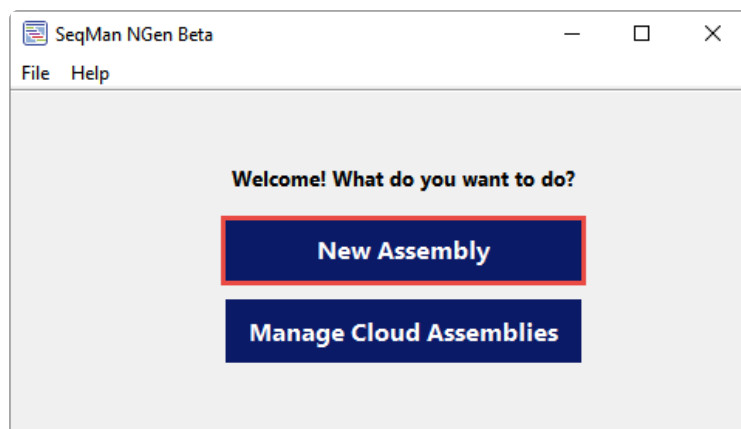
No data are provided for Part A of the tutorial. The samples consist of 2 × 150 bp paired end Illumina reads for NA12891 (father), NA12892 (mother) and NA12878 (daughter). The data set used for assembly is 160 GB in size and requires cloud assembly to perform. Because of this, you will simply read through the steps used to create the whole genome reference-guided assembly in SeqMan NGen. In Part B, you will download the much-smaller finished assembly and follow along with the downstream analysis in ArrayStar.

Begin with [Part A: Setting up the assembly in SeqMan NGen](#).

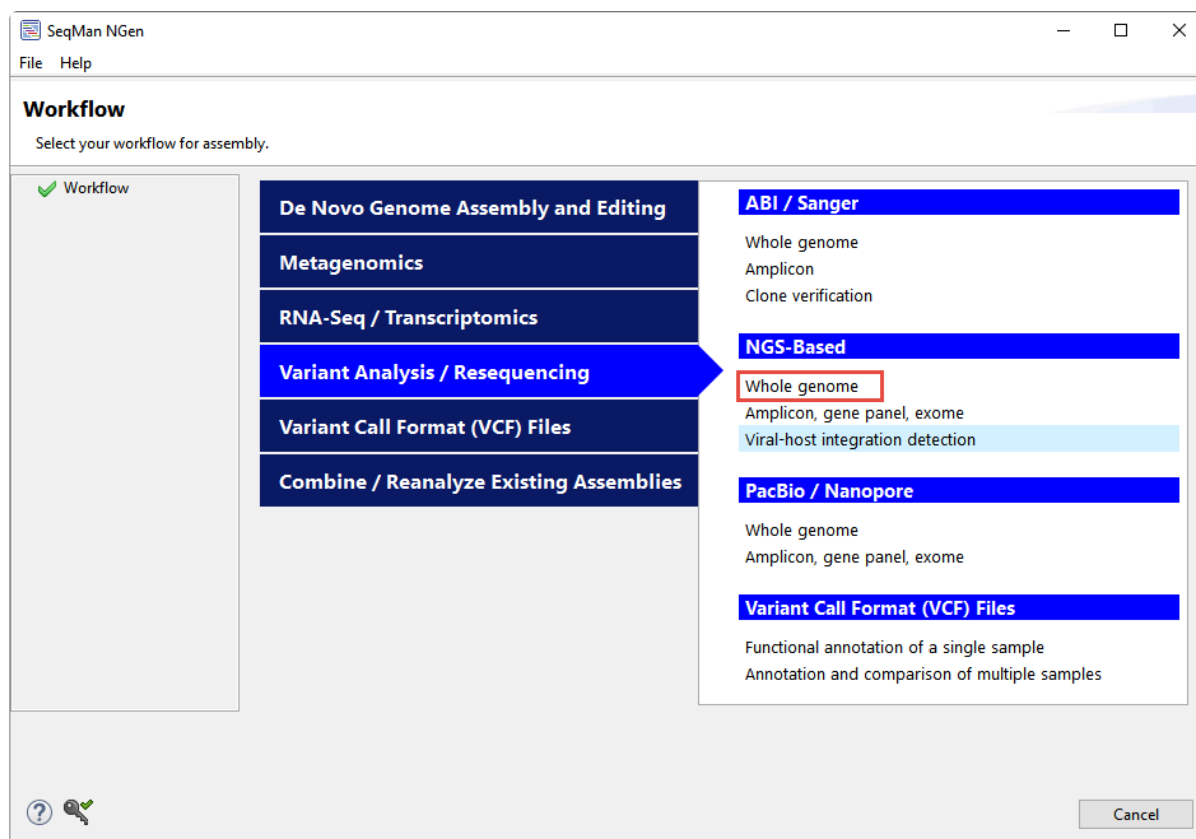
Part A: Setting up the assembly in SeqMan NGen

In this part of the tutorial, you will simply read the steps but won't follow them yourself. No data are provided, as the data set used is 16 GB in size.

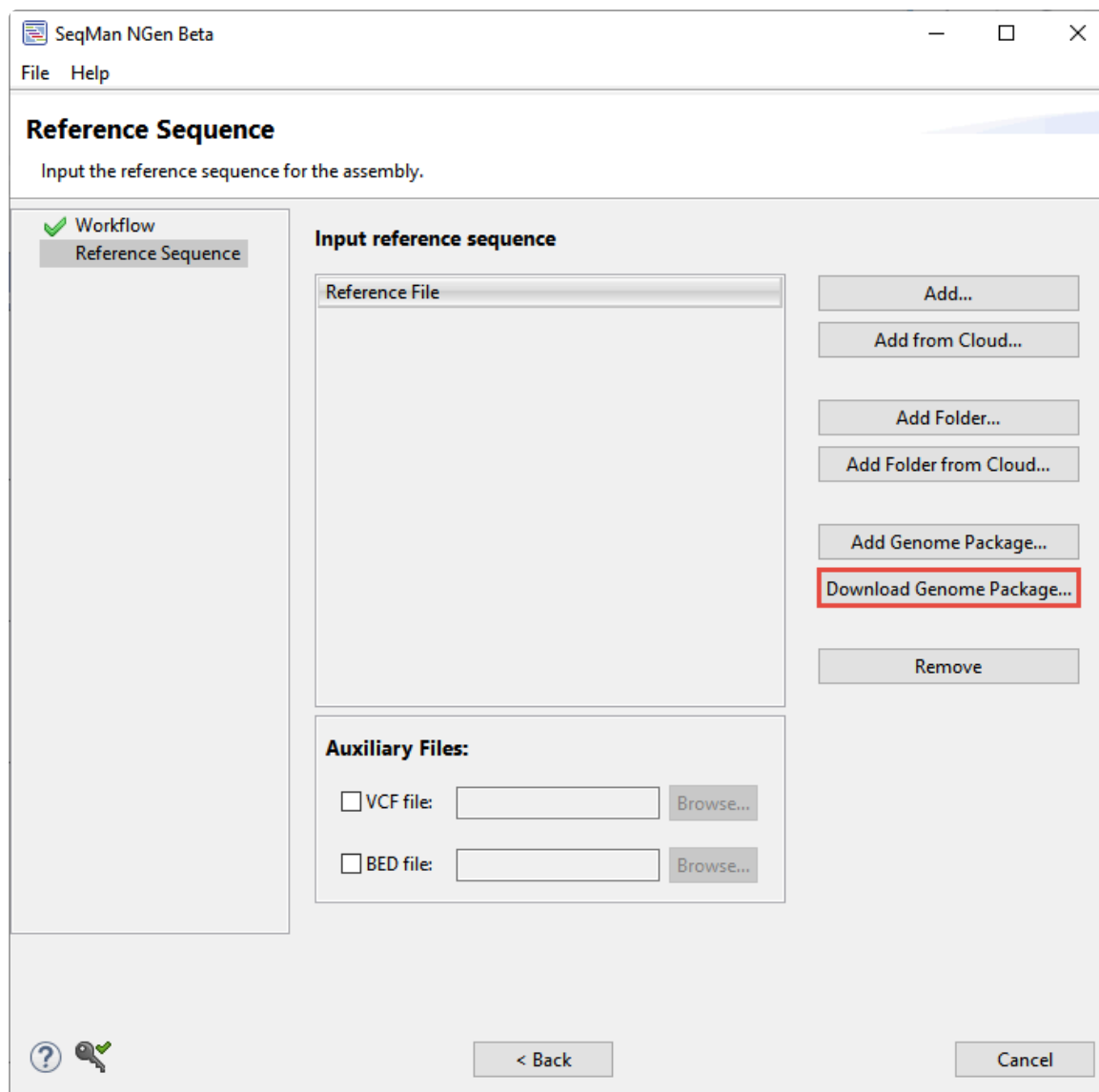
1. Launch SeqMan NGen and choose **New Assembly**.



2. In the Workflow screen, choose **Variant Analysis / Resequencing** on the left. On the right, choose the **NGS-Based** workflow **Whole genome**.



3. In the Reference Sequence screen, click **Download Genome Package**.



Select **Homo sapiens** build **GRCh38.p7** and press **Select**.

Organism ▲	dbSNP Build	Assembly	Download Size
Apis mellifera	149	Amel_4.5	129.9 MB
Arabidopsis thaliana	138	TAIR10	87.1 MB
Bos taurus	148	UMD_3.1.1	1.8 GB
Canis lupus familiaris	138	UMD_3.1.1	2.8 GB
Canis lupus familiaris	138	UMD_3.1.1	2.8 GB
Escherichia coli K12 MG1655	-	ASM584v2	3.1 MB
Gallus gallus	147	5.0	625.6 MB
Homo sapiens	150	GRCh37.p13	4.9 GB
Homo sapiens	150	GRCh38.p7	3.5 GB
Homo sapiens - Ensembl w/PDB	150	GRCh38.p10 + Ensembl 91	3.5 GB
Homo sapiens mitochondrion	150	GRCh38.p7	33.3 kB
Homo sapiens, CEU reference	142	GRCh37.p10	3.8 GB
Homo sapiens, CHBJPT reference	142	GRCh37.p10	3.8 GB
Homo sapiens, YRI reference	142	GRCh37.p10	3.8 GB
Microbial Genome database	-	-	1.5 GB

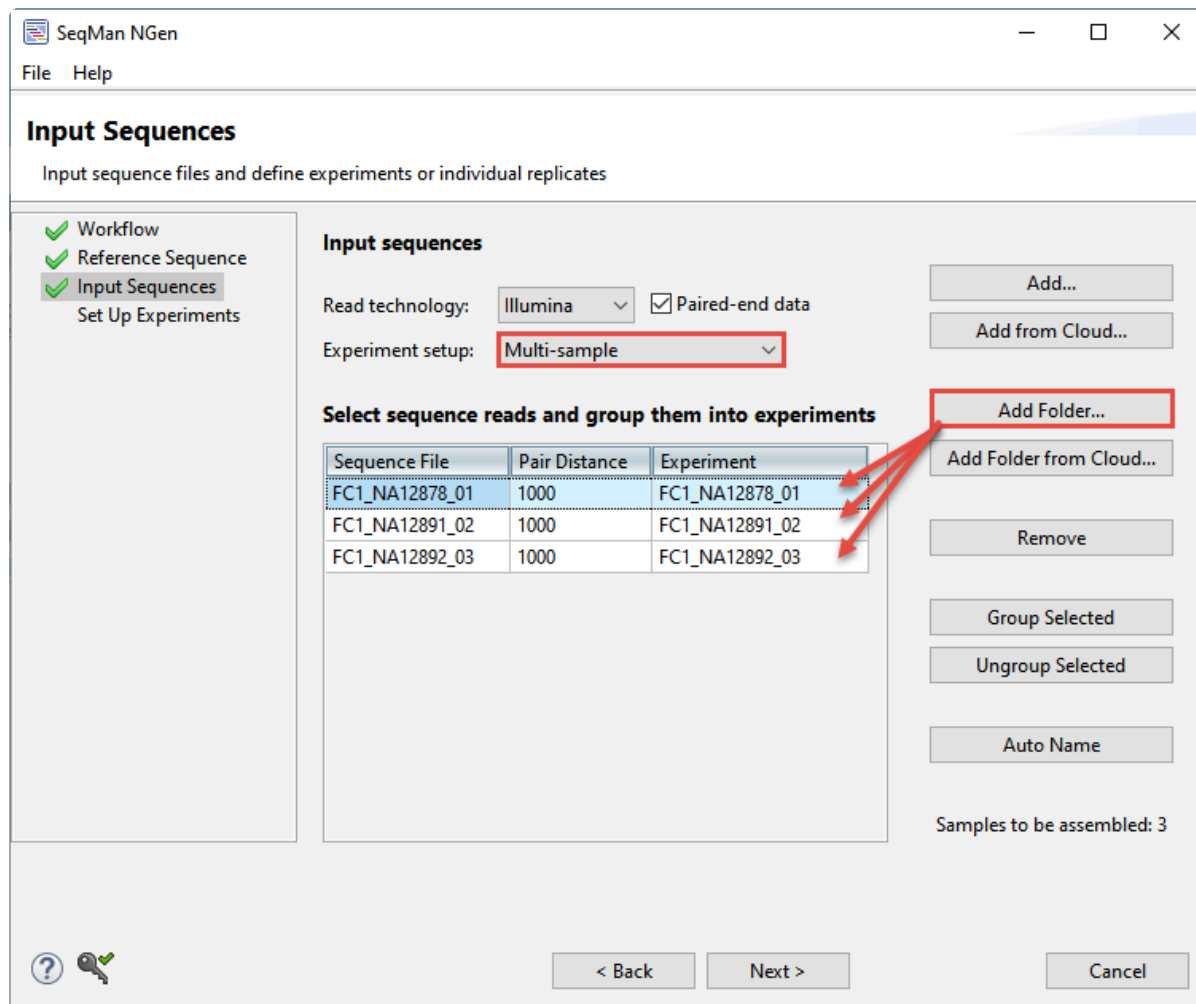
4. Check the box next to **BED file** and then use **Browse** to upload *NexteraRapidCapture_Exome_TargetedRegions_v1.2Used.bed*. Click **Next**.

Auxiliary Files:

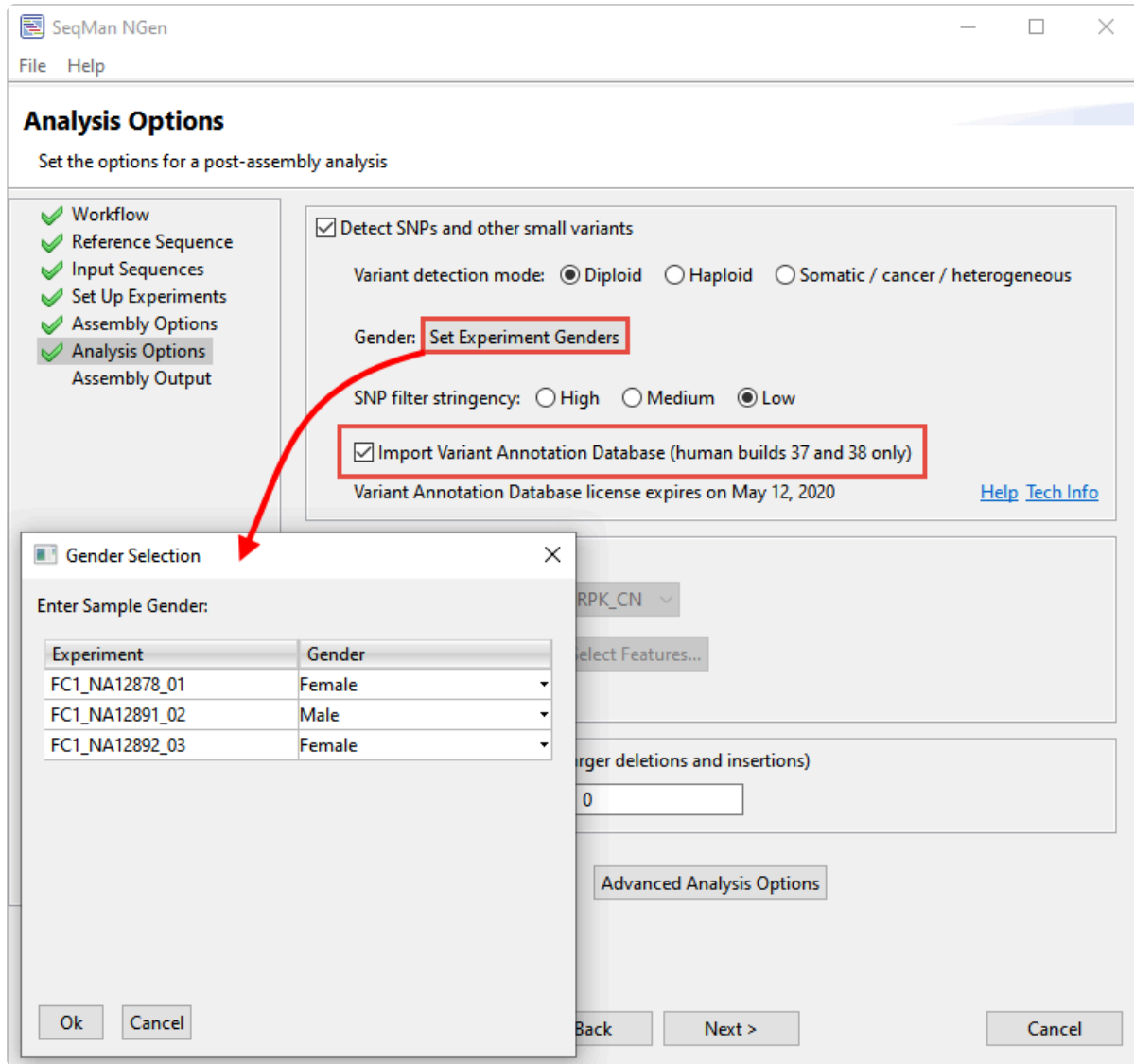
☐ VCF file:

☒ BED file:

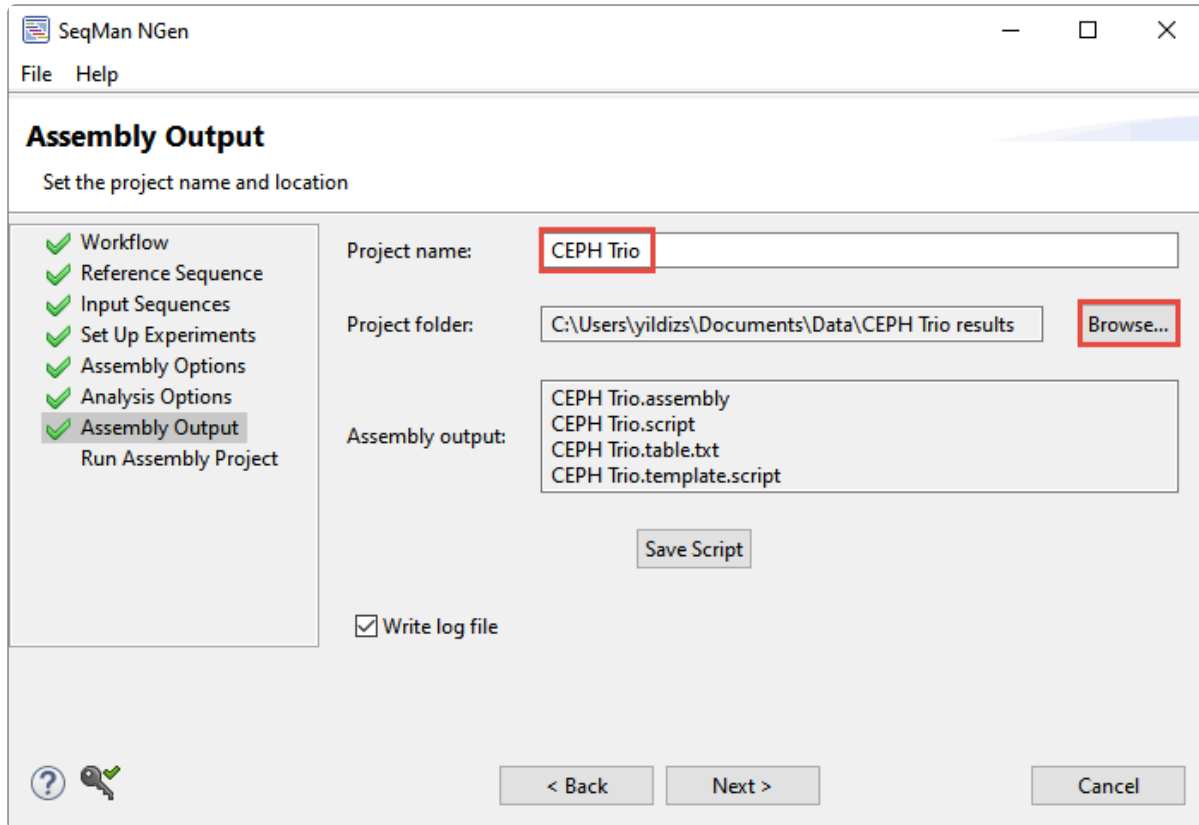
5. In the Input Sequences screen, change the **Experiment setup** to **Multi-sample**. Click **Add Folder** and add *FC1_NA12878_01* (daughter). Repeat twice to add *FC1_NA12891_02* (father) and *FC1_NA12892_03* (mother). Click **Next**.



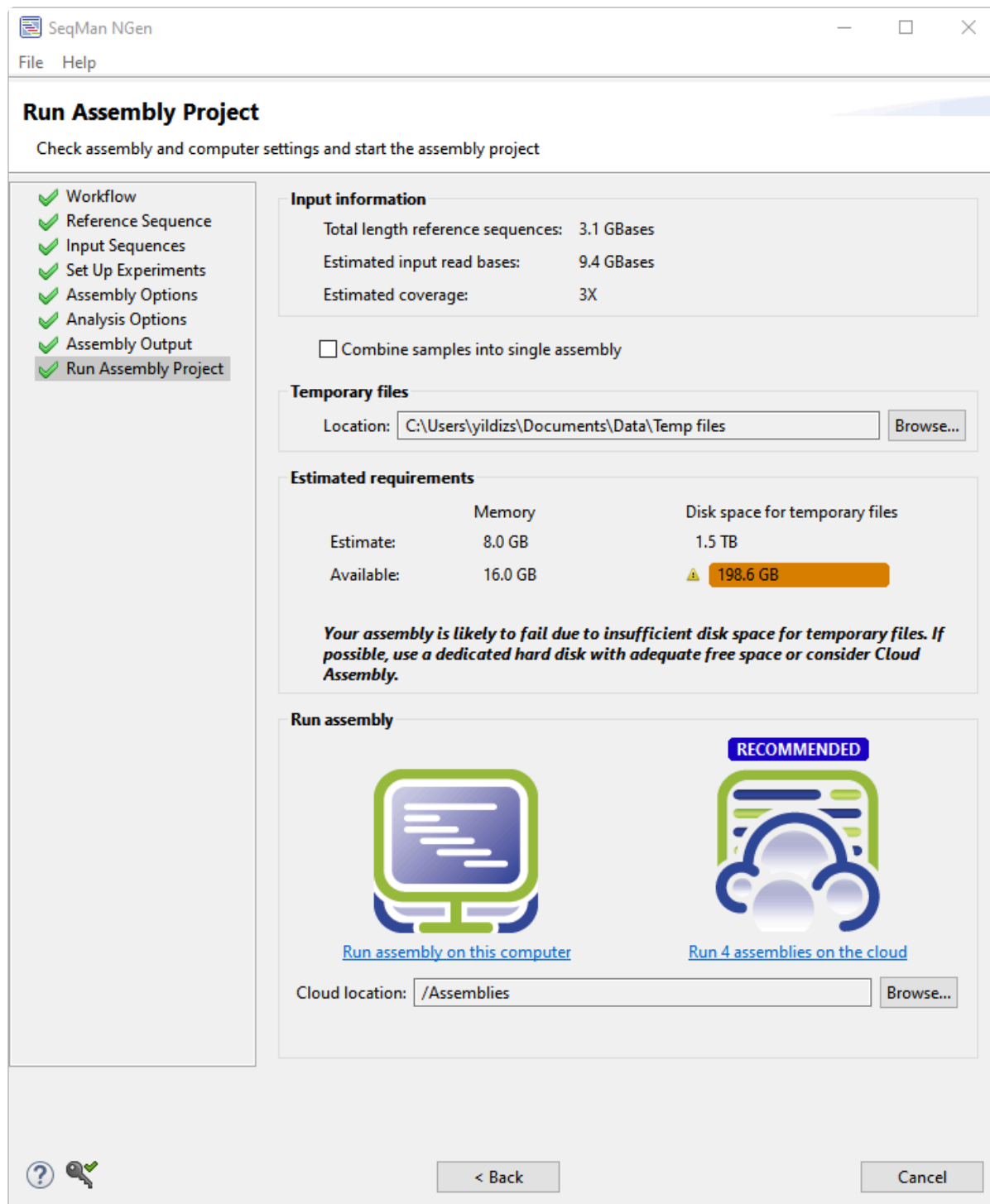
6. In the Set Up Experiments screen, do not check any of the boxes. Click **Next**.
7. In the Assembly Options screen, click **Next**.
8. In the Analysis Options screen, click **Set Experiment Genders**. Select **Female** for the mother, **Male** for the father and **Female** for the daughter. Click **OK**. Check the box next to **Import Variant Annotation Database**. This will provide a highly enriched source of annotations for genes and other items of interest. Click **Next**.



- In the Assembly Output screen, type in a **Project name** of **CEPH Trio**. Use **Browse** to select a writable location for the results. Click **Next**.



10. In the Run Assembly Project screen, note that the assembly requires about 1.5TB of disk space for temporary files. The recommendation is to run the assembly on the cloud.



Click the link “Run 4 assemblies on the cloud.” (The number of cloud assemblies needed is usually the number of samples plus one). The assemblies take about 3.5 hours to run on the cloud.

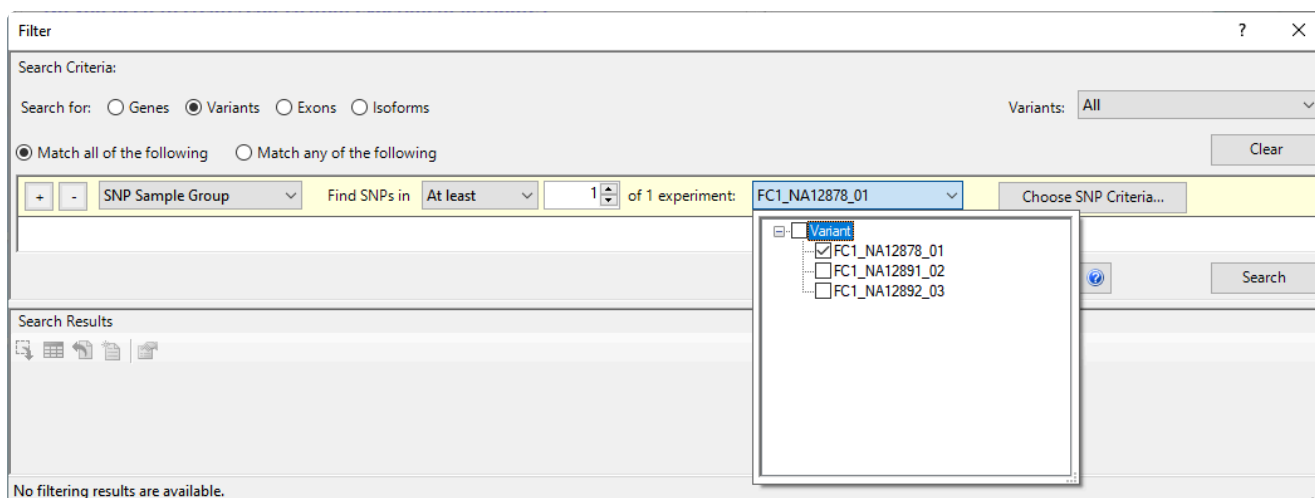
Proceed to [Part B. Analyzing the results in ArrayStar](#).

Part B: Analyzing the results in ArrayStar

In [Part A](#) of this tutorial, you read about how the SeqMan NGen assembly project was set up. In this part, you will download the *CEPH Trio.aster* assembly results file and perform downstream analysis in ArrayStar.

ArrayStar is a discovery tool that provides many different ways to analyze data. The following workflow shows just a sampling of the ways in which you could analyze the CEPH Trio data set. In this case, you will be using advanced filtering to search for potentially deleterious non-synonymous changes in the daughter, and will use the SNP Table to find the source of the SNPs.

1. Download [T4_Whole_Genome.zip](#) (40 MB) and extract it to any convenient location (i.e. your desktop). The data set consists of a single file, *CEPH Trio.aster*.
2. Double-click on *CEPH Trio.aster* to open it in ArrayStar.
3. Choose **Filter > Filter All**. In the gray heading, keep the default settings **Search for = Variants** and **Variants = All**. In the yellow filter row, keep the default **SNP Sample Group** and find SNPs in **At least 1** of 1 experiments. To change the second “3” to a “1,” click the drop-down menu that begins with **FC1_NA12878** and uncheck all of the boxes except for **FC1_NA12878_01** (the daughter).



Deleterious changes are unlikely to occur as homozygotes in germline. To reflect this:

4. Click **Choose SNP Criteria**.
 - a. In the starting tab (**General**), change **SNP Genotype** to **Heterozygous w/Reference**. In the **Translation** section, check **Non-synonymous**. This will add checkmarks to all subclasses as well.

SNP Searching Criteria

SNP Selection Criteria Summary: Non-synonymous, Heterozygous w/Reference

General Statistics Population Genetics Functional Prediction Evolutionary Conservation Pathogenicity

SNP Type: Any type

SNP Genotype: Heterozygous w/Reference

SNP Class

Called as change: ☒ Yes ☐ No ☐ Don't check

Located in region: ☐ Intergenic ☐ Genic ☒ Coding ☐ Non-coding RNA
☐ Don't check

In splice site: ☐ Yes ☐ No ☒ Don't check

Translation:

- ☐ Synonymous
- ☒ Non-synonymous
 - ☒ Substitution
 - ☒ No-start
 - ☒ No-stop
 - ☒ Nonsense
 - ☒ Frameshift
 - ☒ Inframe indel

Within targeted regions: ☒ Don't check ☐ Yes ☐ No

Use Defaults Set Defaults OK Cancel

- b. Select the **Statistics** tab. Check the box next to **Filter minimum P not ref** and change the number to **90.00** (90%). Check the **Filter minimum depth** box and change the number to **20**.

SNP Searching Criteria

SNP Selection Criteria Summary: Non-synonymous, Heterozygous w/Reference, P not ref $\geq 90.00\%$, Depth ≥ 20

General Statistics Population Genetics Functional Prediction Evolutionary Conservation Pathogenicity

Statistics

☒ Filter minimum P not ref % ☐ Require
☐ Filter minimum Q call ☐ Require
☐ Filter SNP % minimum % maximum % ☐ Require
☒ Filter minimum depth ☐ Require

dbSNP: ☒ Don't check ☐ Present ☐ Absent
 User VCF: ☒ Don't check ☐ Present ☐ Absent
 (Provided in SeqMan NGen assembly)

Use Defaults Set Defaults OK Cancel

Press **OK**.

- Back in the Filter dialog, press **Search**. As shown in the bottom left corner of the dialog, 616 matching variants have been found.

Search Results

Ref ID	Ref Pos	Ref Seq
NC_000001.11	69428	T
NC_000001.11	981169	A

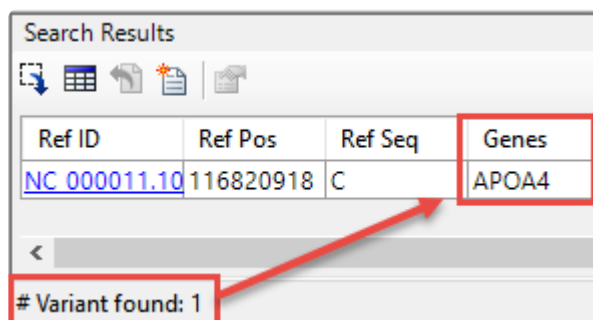
<


Variant found: 616

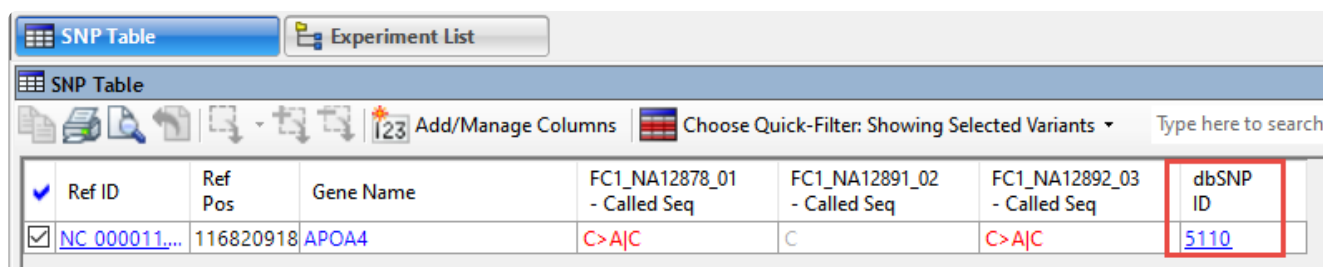
To cut down the number of matches, you will next restrict matches to those variants predicted to be pathogenic in NCBI's [ClinVar](#) database.

- Again press **Choose SNP Criteria**.

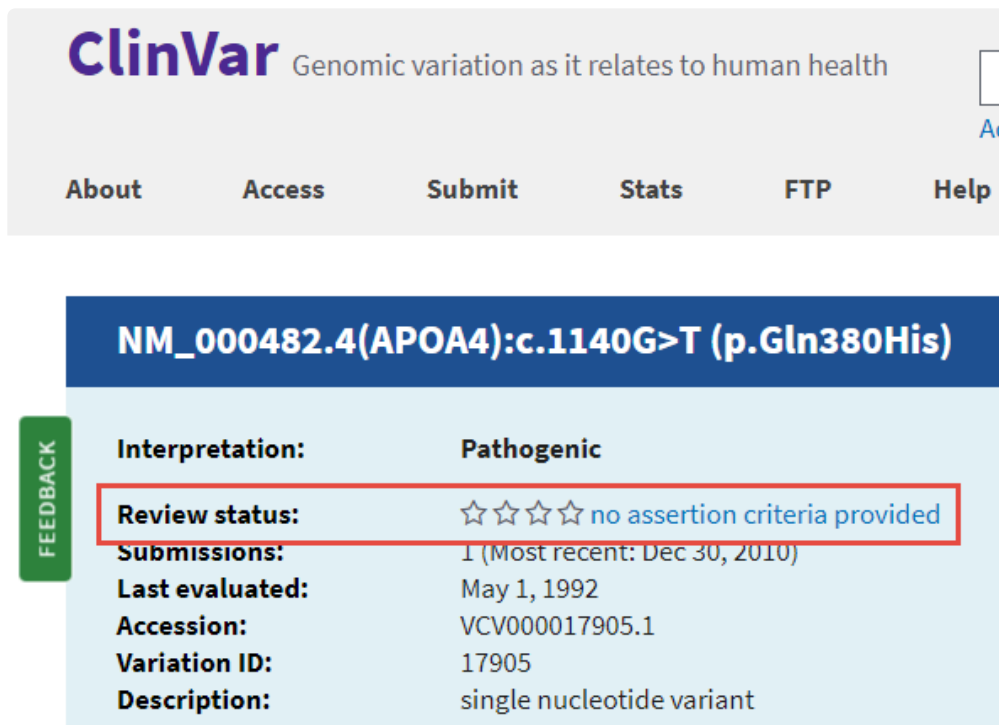
- a. Select the **Pathogenicity** tab. On the left, click on **clinvar_clnsig**. On the right, add checkmarks next to **Likely pathogenic** and **Pathogenic**. Press **Add Filter to Set** and then press **OK**.
7. Back in the Filter dialog, press **Search**. Now, a single variant is found that matches all criteria. This variant occurs in the *APOA4* gene.



8. Press the **Select and Show Results in Variant Table** tool ().
9. In the ensuing SNP Table, click on the link in the **dbSNP ID** column to go to the corresponding web page.



10. On the web page, click the **ClinVar** link to open a page where you can see submitted interpretations, citations and other information about the putative deleterious variant. Note that the Review Status near the top indicates that no assertion criteria have been provided that the variant is pathogenic. This result is therefore inconclusive.



ClinVar Genomic variation as it relates to human health

About Access Submit Stats FTP Help

NM_000482.4(APOA4):c.1140G>T (p.Gln380His)

Interpretation: Pathogenic

Review status: ☆☆☆☆ no assertion criteria provided

Submissions: 1 (Most recent: Dec 30, 2010)

Last evaluated: May 1, 1992

Accession: VCV000017905.1

Variation ID: 17905

Description: single nucleotide variant

You will next try a different line of inquiry. You will change SNP filtering criteria to filter to those variants predicted to be deleterious rather than pathogenic.

11. Once again, press **Choose SNP Criteria**.

- a. Select the **Pathogenicity** tab and click the **Remove** button to remove the ClinVar criterion.
- b. Select the **Functional Prediction** tab.
 - i. On the left, click **LRT_pred**. On the right, check the box next to **Deleterious (D)**. Click **Add Filter to Set**.
 - ii. On the left, choose **MutationTaster_pred**. On the right, check both **Disease causing – automatic (A)** and ****Disease causing (D)***, then click **Add Filter to Set**.
 - iii. On the left, choose **SIFT_pred**. On the right, check **Damaging**, then click **Add Filter to Set**.

Deleterious mutations are likely to be rare in a given population. So in the next step, you will add criteria to filter to variants that are rare in the European population (MAF <5%).

- c. Select the **Population Genetics** tab. Click on the plus sign next to **1000Gp3_MAF** to expand it, then choose **1000Gp3_EUR_MAF** (the European group). Click **Add Filter to Set**.
- d. Click **OK**.

12. Back in the Filter dialog, press **Search**. This search yields 4 variants.

In the final step, you will identify the parent from which each variant was inherited.

13. Click the **Select and Show Results in the Variant Table** tool ()

The SNP Table contains three columns for **Called Seq**. The column ending in 01 pertains to the daughter; 02 to the father; and 03 to the mother. Look at the red text in the **Called Seq** columns. This indicates that two of the four variants were inherited from the mother; one from the father; and one arose *de novo* in the daughter.

SNP Table

Experiment List

SNP Table

<

The variant in this gene...	... is found in:
PTK7	mother and daughter
SMARCA2	daughter only
POMT1	mother and daughter
JMJD8	father and daughter

This marks the end of this tutorial.

Long-read analysis with accuracy evaluation

The following tutorial shows how to do long read assembly in SeqMan NGen, and contains optional steps for assembly validation in [QUAST](#) (Quality Assessment Tool for Genome Assemblies). The sequence data consists of an Oxford Nanopore Technologies (ONT) MAP006-1 *.fastq* file from *E. coli* strain MG1655.

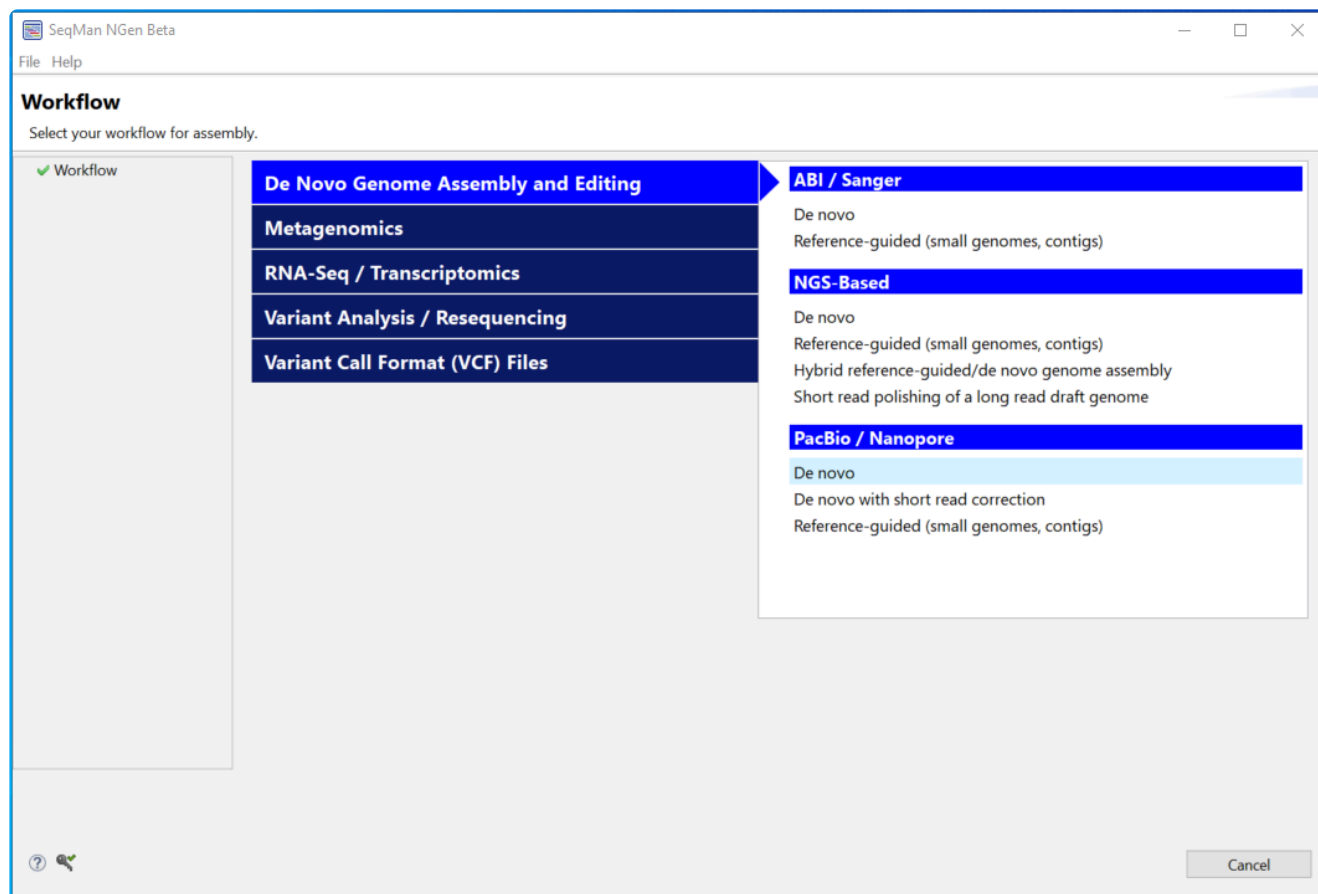
This tutorial includes a step in which you will add a reference sequence in order to assist SeqMan NGen in creating a scaffold. Reference-guided scaffolding is optional, but is commonly called for when: 1) the assembly was broken up into multiple contigs, 2) the genome involves multiple chromosomes (e.g. yeast) and/or plasmids or 3) both of the above. Since long reads contain no pair information, reference-guided scaffolding is a way to provisionally order and orient the contigs. For an organism like the yeast *S. cerevisiae* that has sixteen chromosomes and two plasmids, mitochondrion and two micron, the scaffolding automatically associates each contig with its canonical chromosome rather than have eighteen or more contigs in random order to sort through. It's important to note that, since structural variation is common between species and strains, the ordering and orientations are provisional. As always, the closer the reference sequence is to the sequence strain, the better. More distantly related sequences will likely yield inconsistent/inaccurate scaffolding.

Begin with [Part A: Running the assembly in SeqMan NGen and viewing it in SeqMan Ultra](#).

Part A: Running the assembly in SeqMan NGen and viewing it in SeqMan Ultra

In Part A, you will use SeqMan NGen to run the assembly, then launch the results in SeqMan Ultra.

1. Download [T5_Long_Read.zip](#) (232 MB) and extract it to any convenient location (i.e., your desktop). Part A of the tutorial will use the files *MAP006-1_2D_pass.fastq* (sample) and *U00096.3.gbk* (ref).
2. Launch the workflow in the SeqMan NGen wizard using either of the methods below:
 - Launch the DNASTAR Navigator, click on the **Genomics** tab on the left. On the right, click **De novo genome assembly and editing**. This launches SeqMan NGen at the Workflow screen. On the right, under **PacBio/Nanopore** options, click **De novo**.
 - Launch SeqMan NGen and choose **New Assembly**. In the Workflow screen, choose **De Novo Genome Assembly and Editing** on the left. On the right, under **PacBio/Nanopore** options, click **De novo**.



3. In the Input Sequences screen, choose a **Read Technology** of **ONT**. Click the **Add** button and add the sequence *MAP006-1_2D_pass.fastq*.

SeqMan NGen Beta

File Help

Input Sequences

Input long read sequence files

- ✓ Workflow
- ✓ **Input Sequences**
- ✓ Preassembly Options
- ✓ Post Assembly Options
- ✓ Assembly Output
- ✓ Run Assembly Project
- ✓ Assembly Log
- ✓ Assembly Summary

Input long read sequences

Read technology: ONT

Sequence File

MAP006-1_2D_pass.fastq

Add...
Add from Cloud...
Add Folder...
Add Folder from Cloud...
Remove

< Back Next > Finish

Click **Next**.

- In the Preassembly Options screen, enter an **Expected genome length** of 4600000 (46 followed by 5 zeroes). Leave the **Desired depth of coverage of final assembly** at the default setting of 100. Leave **Use longest reads in data set to achieve depth** selected.

SeqMan NGen Beta

File Help

Preassembly Options

Set the most important pre-assembly options for a successful assembly

- ✓ Workflow
- ✓ Input Sequences
- ✓ **Preassembly Options**
- ✓ Post Assembly Options
- ✓ Assembly Output
- ✓ Run Assembly Project
- ✓ Assembly Log
- ✓ Assembly Summary

Preassembly Options

For optimal results, use a subset of reads to achieve 100X depth of coverage

☐ Assemble all reads
☒ Use subset of reads

Expected genome length: bp

Desired depth of coverage of final assembly:

☒ Use the longest reads in data set to achieve depth
☐ Use the first n reads in data set to achieve depth

Correct the longest, higher-quality reads, then assemble only these corrected reads

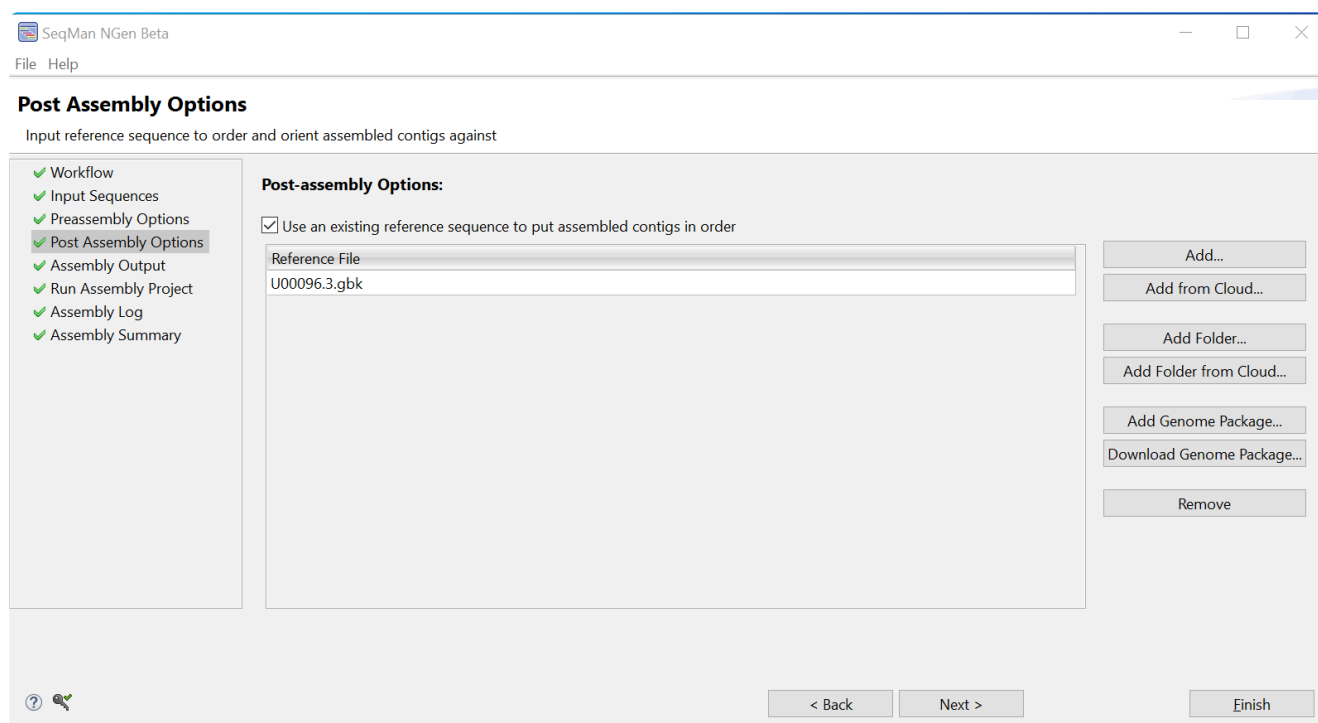
☐ Run a first-pass correction assembly

Advanced Assembly Options

< Back Next > Finish

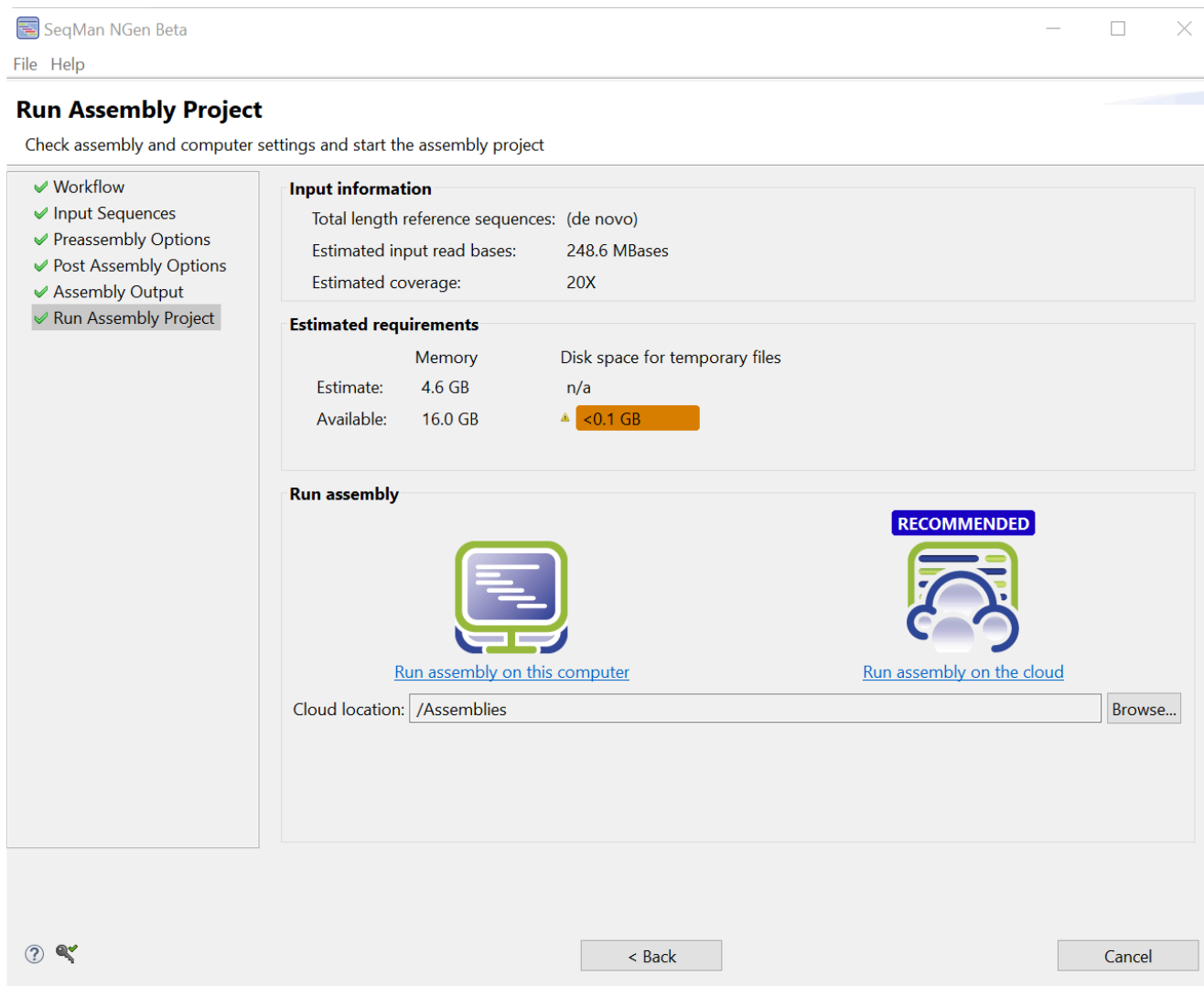
Click **Next**.

5. In the Post Assembly Options screen, keep the box checked and use the **Add** button to add the reference genome *U00096.3.gbk*.




Click **Next**.

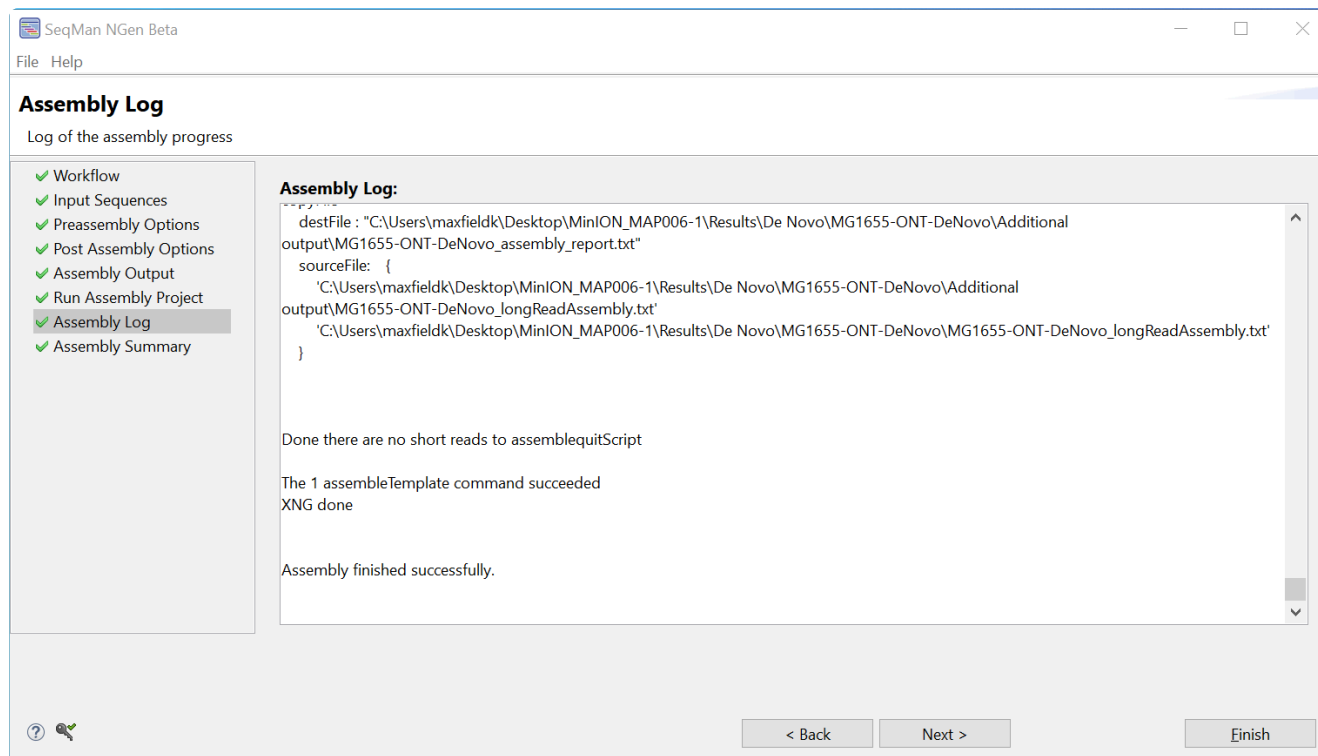
6. In the Assembly Output screen, type in a **Project Name**. Then use the **Browse** button to select a save location for the project. Click **Next**.
7. In the Run Assembly Project screen, look at the “Recommended” method under **Run assembly**.



Regardless of the recommendation, it is probably safe to choose **Run assembly on this computer** (see Note below). Typically, local assembly takes 30-60 minutes.

 **Note:** SeqMan NGen's long-read workflow is still considered a 'beta' version and has a known issue where available disk space may be calculated incorrectly for *de novo* assemblies without short read correction. In this situation, SeqMan NGen may erroneously encourage you to run assembly on the cloud. For this tutorial, most users can simply ignore this error. At worst, the assembly will begin but not finish.

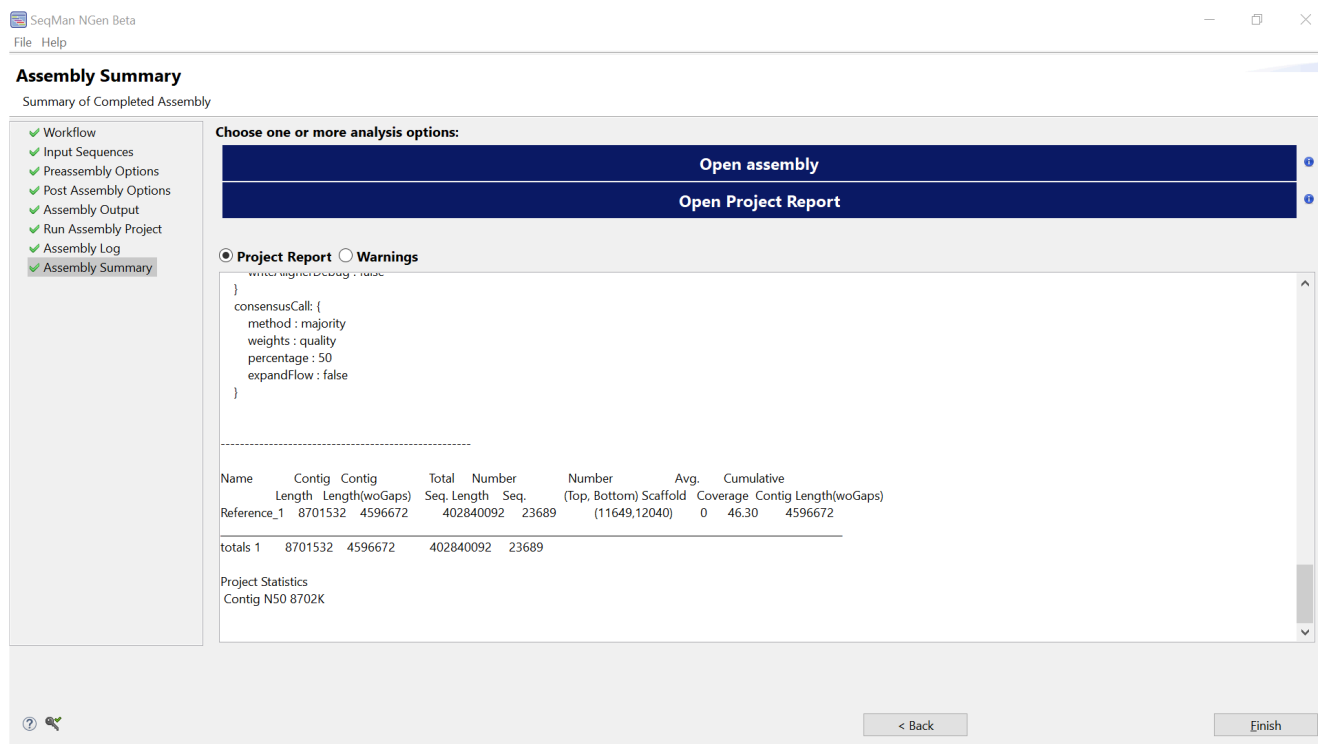
8. A message should appear at the bottom of the Assembly Log screen indicating that the assembly has finished successfully.



Click **Next**.

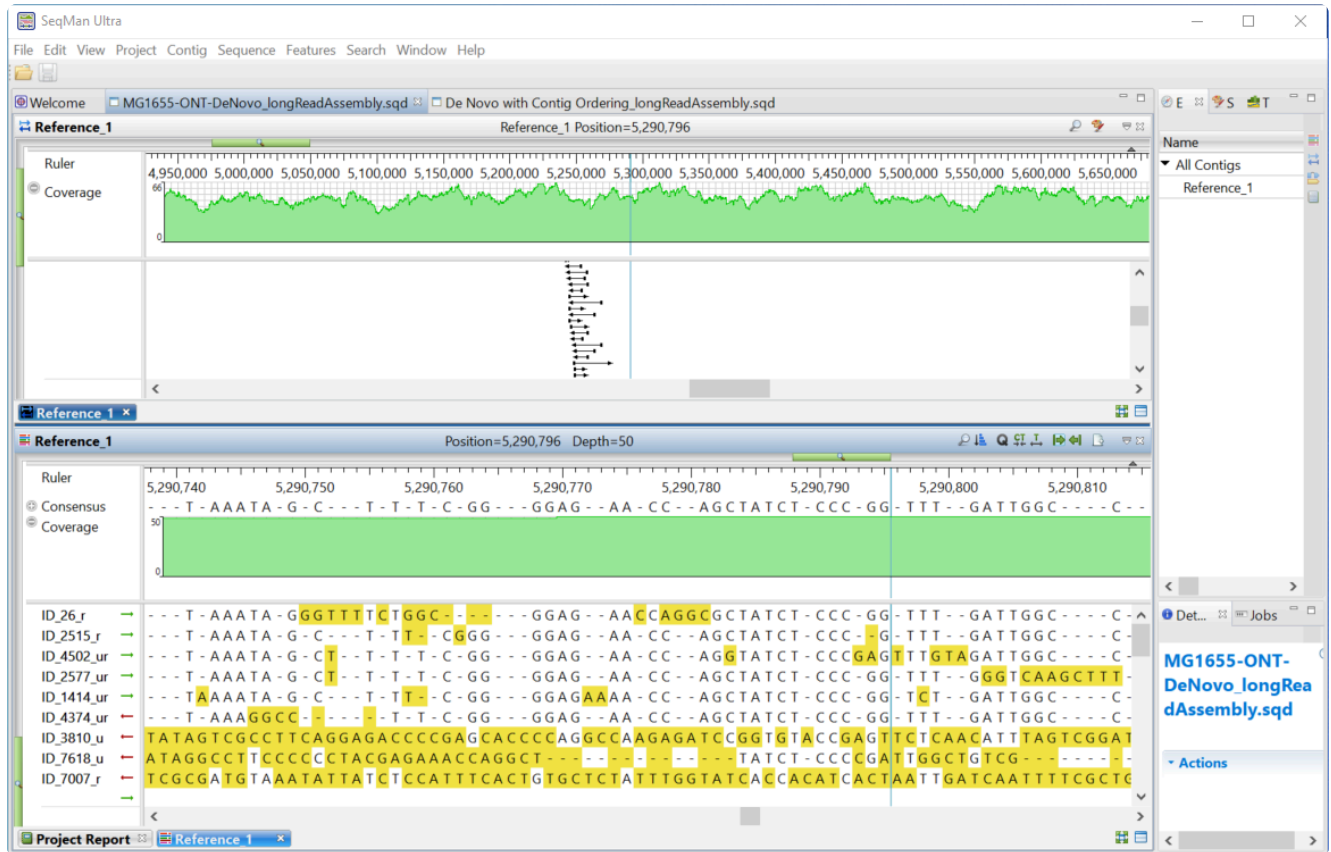
9. By scrolling through the bottom part of the Assembly Summary screen, you will see it consists of three sections:

- **Run statistics** is a high-level summary that includes the number of input and assembled reads and the assembly time.
- **Script** displays a portion of the run script and contains useful debugging information.
- **Contig/Scaffold statistical summary** includes information on the number of contigs assembled, their length and average depth of coverage. Note that the contig N50 value is includes consensus gaps.



As shown in the image above, the single contig produced by this assembly may be somewhat shorter (e.g., ~4.60Mb) than the complete_E. coli_genome (~4.64Mb) . This is due to a number of bases being erroneously “deleted”.

10. Click the blue **Open assembly** button to launch assembly results in SeqMan Ultra.
11. In SeqMan Ultra, open the Alignment and Strategy views and scroll to inspect the assembly. You will notice a significant number of gaps in the alignment caused by the relatively high error rate of the reads.



If desired, proceed to [Part B: Evaluating assembly accuracy using QUAST](#).

Part B (optional): Evaluating assembly accuracy using QUASt

After running the long read assembly in [Part A](#), the following steps can be used to evaluate the base level accuracy of the assembled sequence. The [Quality Assessment Tool for Genome Assemblies](#) (QUASt) is utilized to compare the assembled consensus sequence to the *E. coli* reference sequence.

Before you begin, you will need to decide whether you want to use the [QUASt Web Interface](#) or [download](#) and install the local version of QUASt. In order to use the local version of QUASt, Python must be installed on your computer before you begin.

Both the reference sequence and the assembly output must be in *.fas* or *.fasta* format. The sequences in this case are:

- Assembled sequence: This is the exported consensus from the completed assembly from Part A. It can be found in the Project Folder you specified during assembly setup in the **Additional output** subfolder. The file name will end in *longReadAssembly_consensusSeqs.fas*.
- *E. coli* reference sequence: *U00096.3.fas*.

To perform the evaluation using the QUASt web interface:

1. Go to the [QUASt Web Interface](#).
2. Under **Assemblies**, upload the sequence file ending in *longReadAssembly_consensusSeqs.fas*.
3. Under **Genome**, click **Another genome** and upload *U00096.3.fas* as the **Reference**. Enter a **Name** for the genome.
4. Click **Evaluate**.

When analysis is complete, the results will appear in the browser window. QUASt provides an option to enter your email address so that you can return to the page later.

To perform the evaluation using a locally-installed version of QUASt:

1. [Download](#) and install the local version of QUASt (you must already have Python installed).
2. Launch the Command Line (Win) or Terminal (Mac).
3. Type the following text: `>[path]\quast.py [path]\[project name]_longReadAssembly_consensusSeqs.fas -R [path]\U00096.3.fas`. Press **Enter**.

When analysis is complete, the results are saved by default in the *Users* directory in a *quast_results* folder.

Interpreting QUAST results:

QUAST produces a number of useful output files for evaluation.

QUAST	
Quality Assessment Tool for Genome Assemblies by CAB	
12 May 2020, Tuesday, 20:04:03	
View in Icarus contig browser	
All statistics are based on contigs of size ≥ 500 bp, unless otherwise noted (e.g., "# contigs (≥ 0 bp)" and "Total length (≥ 0 bp)" include all contigs).	
Aligned to "U00096" 4 641 652 bp 1 fragment 50.79 % G+C	
Genome statistics	MG1655_ONT_DeNovo_longReadAss...
Genome fraction (%)	99.484
Duplication ratio	0.994
Largest alignment	1 062 699
Total aligned length	4 588 966
NGA50	628 951
LGA50	3
Misassemblies	
# misassemblies	11
Misassembled contigs length	4 596 672
Mismatches	
# mismatches per 100 kbp	78.29
# indels per 100 kbp	808.15
# N's per 100 kbp	24
Statistics without reference	
# contigs	1
Largest contig	4 596 672
Total length	4 596 672
Total length (≥ 1000 bp)	4 596 672
Total length (≥ 10000 bp)	4 596 672
Total length (≥ 50000 bp)	4 596 672
Extended report	

To calculate a first approximation base-level percent accuracy value for the consensus sequence:

1. Open the *report.pdf* file.
2. On the Report page find the **Total aligned length**.
3. On the Misassemblies Report page, find the entries **# Mismatches** and **Indels length**.
4. **Accuracy** = $100 * (\text{Total aligned length} - \text{\# Mismatches} - \text{Indels length}) / \text{Total aligned length}$

The initial SeqMan NGen beta version of the long-read assembler typically produces a single full length contig of the MG1655 genome with an approximate base-level percent accuracy of 98.65%. The total length of the consensus sequence may be shorter than the expected reference sequence value of 4,641,652 largely due to the number of "deletions" caused by errors in data as well as outstanding issues in the aligner of this beta release.

This marks the end of this tutorial.

Exome workflow with analysis in ArrayStar

The following video shows how to set up an exome assembly and illustrates downstream analysis in ArrayStar. No data files are provided for this tutorial.

Wizard screen descriptions

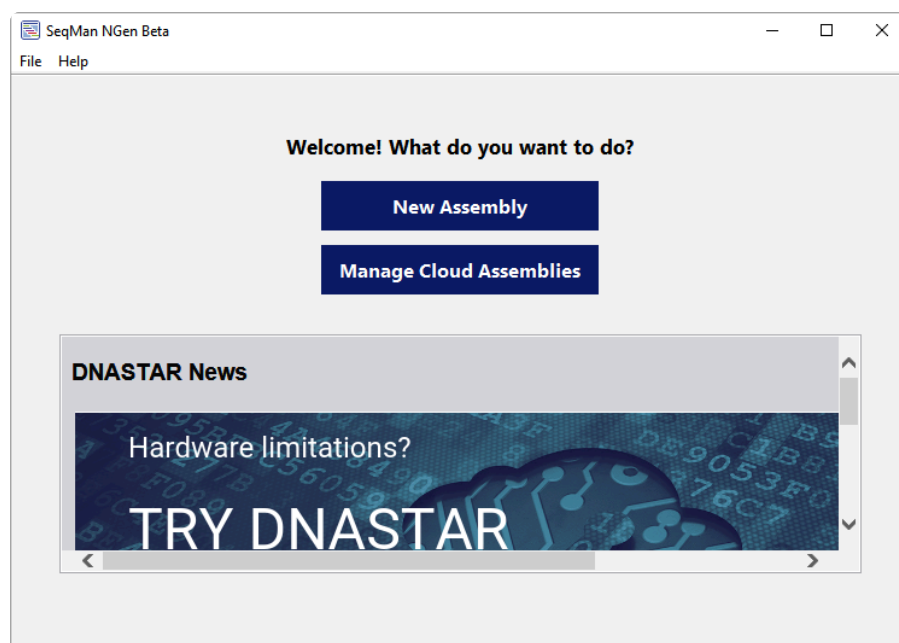
SeqMan NGen's project setup wizard lets you upload files and optimize parameters for your assembly. After choosing your desired workflow in the first ("Workflow") screen, subsequent wizard screens will vary according to the workflow chosen. Most workflows include 5-8 wizard screens. Each screen is described in a separate topic:

- [Welcome screen](#) (appears only when you first open SeqMan NGen)
- [Workflow](#) (used in all workflows)
- [Analysis Options](#)
- [Assembly Log](#)
- [Assembly Options](#)
- [Assembly Output](#) (used in all workflows)
- [Assembly Summary](#)
- [Cloud Monitor](#)
- [Define Binding Proteins](#)
- [Input Assemblies](#)
- [Input Host Files](#)
- [Input Reference](#) (Sequence, Genome, for Scaffolding, etc.)
- [Input Sequences](#) (including Long Read, Short Read)
- [Input VCF Files](#)
- [Input Viral Genomes](#)
- [Post Assembly Options](#)
- [Preassembly Options](#)
- [Run Assembly Project](#)
- [Set Contaminant](#)
- [Set Up Experiments](#)
- [Set Up Replicate Sets](#)
- [Polishing Options](#)
- [Transcript Annotation Database](#)
- [Options tabs](#)

Additional information pertaining to the wizard can be found in [Add and remove files in the wizard](#) and [Use editing commands in the wizard](#).

Welcome

Each time you launch SeqMan NGen, the first screen that appears is the “Welcome” screen.



The top of the screen has two option buttons.

- Press **New Assembly** to set up and run an assembly project using the SeqMan NGen wizard. You will start at the [Workflow](#) screen. Subsequent screens depend on which option you choose in the Workflow screen.
- Press **Manage Cloud Assemblies** to go to the Cloud Monitor screen, where you can monitor current or past cloud assemblies. If you are not currently logged in, you will be prompted to log in using the same email and password you use to access your information on the DNASTAR website.

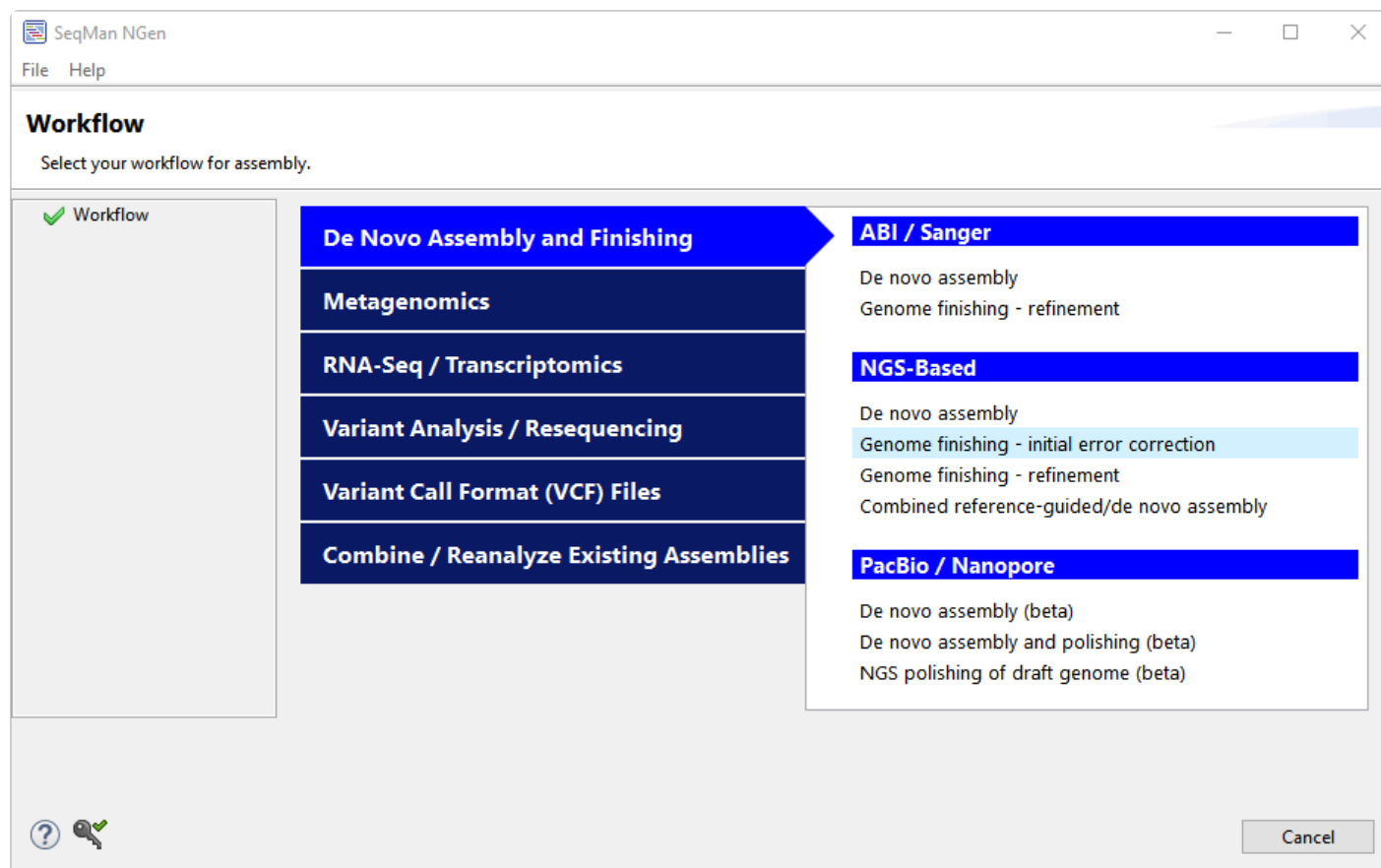
Login to the DNASTAR cloud to view cloud assembly jobs:

Email: Password: ☒ Save password

The bottom of the screen contains news, upcoming events, or new feature announcements from DNASTAR.

Workflow

“Workflow” is the first wizard screen in SeqMan NGen and is where you select the assembly workflow. Each group of workflows is accessed by clicking its dark blue “bar-shaped” tab on the left.



Click a link below for descriptions of each workflow in that group.

- [De novo genome assembly and editing](#)
- [Metagenomics](#)
- [RNA-seq/transcriptomics](#)
- [Variant analysis/resequencing](#)
- [Variant Call Format \(VCF\) files](#)
- [Combine/Reanalyze Existing Assemblies](#)

De novo genome assembling and editing workflows

The following table describes each of the workflows available in the **De novo genome assembly and editing** tab of the [Workflow](#) screen.

Group	Workflow	Description
ABI / Sanger	De novo assembly	Fast, accurate trimming and assembly of Sanger trace data, creating a project file that can be edited in SeqMan Ultra or SeqMan Pro. A non templated assembly of up to 30 million sequence reads and up to a 50 Mbase total length for all contigs combined. The capacity is determined by the amount of available RAM. When assembling a data set <i>de novo</i> , we recommend using paired end data, if available.
	Genome finishing – refinement	Align Sanger data to a draft sequence for further refinement of small errors. (Note: Use Variant Analysis/Resequencing workflow if your primary intent is SNP analysis). This workflow is most frequently used for extending off the ends of saved contig consensus sequences and correcting small errors within the contigs. This type of assembly can include up to 10 million reads and up to a 100 Mbase genome. It can be edited at a later time using a utility like SeqMan Ultra or SeqMan Pro.
NGS-based	De novo assembly	Assembly of Sanger, Illumina and Long-read sequencing data that produces a file that can be edited with SeqMan Ultra or SeqMan Pro. A non templated assembly of up to 30 million sequence reads and up to a 50 Mbase total length for all contigs combined. The capacity is determined by the amount of available RAM. When assembling a data set <i>de novo</i> , we recommend using paired end data, if available.
	Genome finishing – initial error correction	<p>Align NGS data to a draft genome or contigs to correct large misalignments and smaller errors. This option utilizes both reference-guided and <i>de novo</i> assembly steps to resolve both single nucleotide and small multibase replacements (indels) as well as three types of larger structural variation (SV): insertions, deletions and large indels with minimal user intervention. In this workflow, your data should be from a haploid genome with at least one mate pair data set with read lengths of 100 bases or greater. Your total number of reads should be 10 million or less. If you use a larger data set, only the first 10 million reads will be used. For mate pair data, equal numbers of matching forward and reverse reads are processed. The SQD-formatted assembly can be edited at a later time using SeqMan Ultra or SeqMan Pro. When opened in either application, contigs will already be organized into scaffolds in the Explorer panel.</p> <p>This workflow replaces the “gap closure workflow” from Lasergene 16.0 and before. This newer version features an additional “refinement” stage before the “gap closure” stage and some additional “finishing” steps after the gap closure portion takes place. During assembly, data is processed in several stages:</p>

		<ol style="list-style-type: none"> 1. Data is mapped and aligned to a user-defined set of consensus sequences from which a new consensus sequence is determined. Five rounds of this consensus refinement process are performed to remove the majority of single nucleotide and small multibase errors. 2. Data is mapped and aligned to refined consensus sequence(s) from stage 1 and then analyzed for characteristic SV motifs. 3. The reference sequence is split at the detected SV sites, forming a series of ordered contigs. 4. Mate pair and split reads from each SV event are collected in site-specific pools and assembled <i>de novo</i>. Deletions are detected using three types of data: split reads, spanning paired-end reads, and sequence coverage information. For insertions and replacements, mate pair reads corresponding to the new sequence are collected from the unassembled read pool. Only reads anchored by mates flanking the SV in the main assembly are used at this stage. 5. The <i>de novo</i> assembled contigs are then brought into the main assembly and positioned consistently with the mate pair information. 6. For SVs where the gap is not completely covered by the <i>de novo</i> assembled contigs (e.g. insertions longer than twice the size of the insert library), additional reads from the unassembled read pool matching and extending the ends of the joining contigs are added in an attempt to “walk” across the gap. This walk is terminated when either no new reads are found or when a repeated element is encountered. 7. Note that the final <i>de novo</i> assembly performed in stage 6 typically results in additional contigs added to the final assembly project. These are often small contigs with redundant sequences of chromosomal segments. However, they can also represent plasmids, for example, that were not present in the input consensus sequences. <p>Click here to see benchmarks for SeqMan NGen vs. three open source tools.</p>
	Genome finishing – refinement	Align NGS data to a draft genome for further refinement of small errors and closing small gaps between contigs. This workflow is most frequently used for extending off the ends of saved contig consensus sequences and correcting small errors within the contigs. This type of assembly, which uses mate-pair data, can include up to 10 million reads and up to a 100 Mbase genome. It can be edited at a later time using a

PacBio/ Nanopore		utility like SeqMan Ultra or SeqMan Pro.
	Combined reference-guided/ <i>de novo</i> assembly	This workflow aligns paired end NGS data from a new strain/isolate to a closely-related reference genome (>90% identity) to replace SNVs and small indels as well as larger structural variants in the reference with the sequences of the new organism. This workflow is analogous to the Genome finishing – initial error correction workflow above and uses the same series of stages to construct the new sequence from the starting reference.
	<i>De novo</i> assembly (beta)	<i>De novo</i> assembly of long-read-only data sets with an option to first “correct” a genome spanning set of overlapping read prior to assemblies. The “simple” version is described in the tutorial De novo assembly using Sanger data . A second option known as “correct first” is initiated by choosing the same <i>de novo</i> workflow and then proceeding to the Preassembly Options page. On that page, select the Run a first-pass correction assembly option. The “correct first” mode consists of two stages. First, the set of primary overlapping reads covering each contig from end to end are identified and combined with their overlapping and containment reads in a series of mini assemblies, the consensus sequences of which represent “corrected” sequence reads. Second, the corrected read sequences are <i>de novo</i> assembled into a final assembly from which new consensus sequences are determined. In this beta version, note that this workflow typically produces more contigs than the standard single stage <i>de novo</i> assembly, but consensus sequences are usually of somewhat higher base-level accuracy.
	<i>De novo</i> assembly and polishing (beta)	<i>De novo</i> assembly of long-read-only data sets followed by NGS polishing to correct assembly errors. Choosing this workflow will first <i>de novo</i> assemble a long read data set and then automatically run the Genome finishing – initial error correction workflow starting from the <i>de novo</i> assembled consensus sequence(s). Polishing and Genome finishing are largely synonymous terms with the former coined in connection with long read assemblies. You will be asked to specify both the long read and short read data on separate pages in the wizard.
	NGS polishing of draft genome (beta)	This option takes an existing set of long read assembled contig consensus sequences together with a NGS paired end data set from the same organism and runs the Genome finishing – initial error correction workflow . Polishing and Genome finishing are largely synonymous terms with the former coined in connection with long read assemblies. You will be asked to specify both the long read and short read data on separate pages in the wizard.

Create a reference-guided assembly to use in the “SNP to Structure” workflow

If you are working with reference-guided human assemblies, Lasergene’s “SNP to Structure” workflow lets you combine genomic sequencing and variant level data with structure files from the [RCSB Protein Data Bank](#) (PDB) to model point mutations on the protein structure and assess the effect on protein stability. By combining structural bioinformatics with sequencing technologies, this integrated workflow can guide genomic and molecular biology researchers to create structure-based hypotheses and to investigate possibilities not evident by sequences alone.

This workflow requires that you be licensed to use several Lasergene applications: SeqMan NGen, SeqMan Pro / SeqMan Ultra and Protean 3D (all required), and ArrayStar (optional).

Only Part A of the workflow involves SeqMan NGen. However, all parts of the workflow are described below.

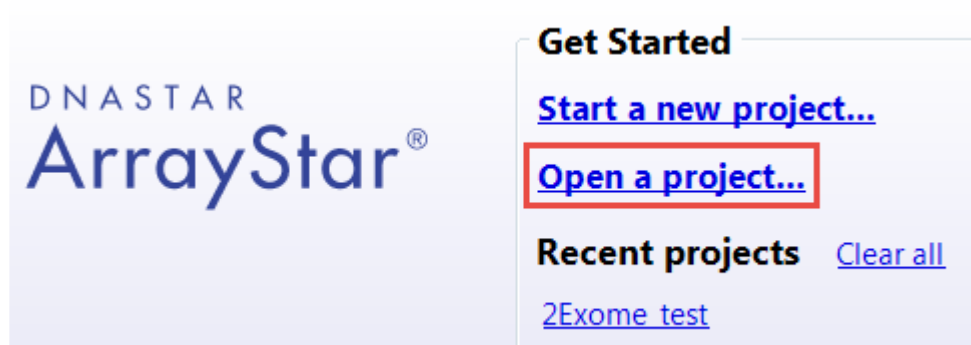
Part A: Create a reference-guided assembly in SeqMan NGen:

1. In SeqMan NGen’s [Workflow](#) screen, choose a reference-guided workflow.
2. In the [Reference Sequence](#) screen, [add the DNASTAR genome template](#) Homo_sapiens-GRCh38-Ensembl-dbsNP150.zip_. This template contains the mapping of the sequences of PDB structures to the human genomic coordinates, and will later allow SeqMan Pro to communicate with Protean 3D. SeqMan NGen outputs an *.astar* and an *.assembly* package to use in later steps.
3. Follow the rest of the wizard steps and create the assembly.

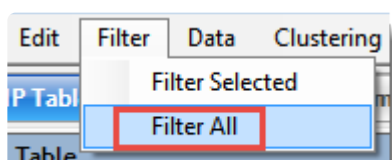
Part B (optional): Filter for variants of interest in ArrayStar:

If your assembly has a large number of variants, you can use ArrayStar to filter them down to a smaller group of interest before sending them to SeqMan Pro for viewing. This step is highly recommended for all but very small assemblies.

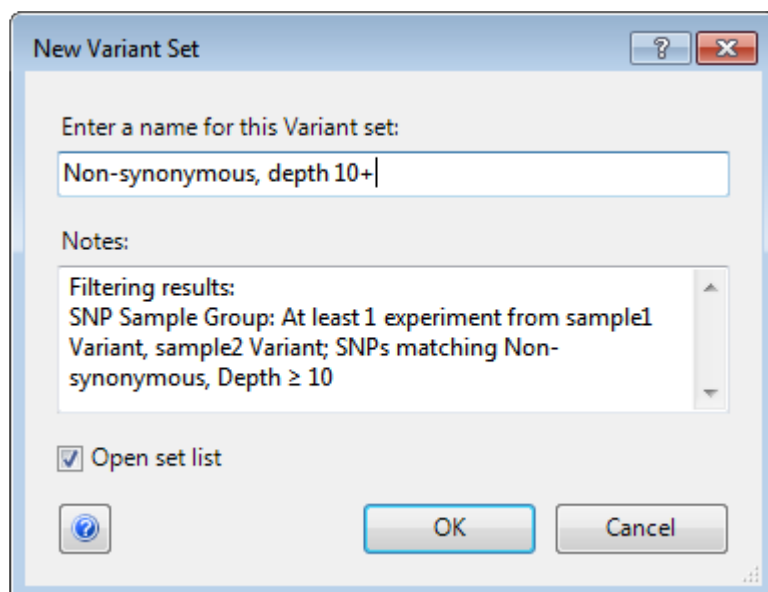
1. Launch ArrayStar and choose **Open a project**.



2. Navigate to and select the *.astar* file output by SeqMan NGen.
3. When the file has loaded, click on the **SNP Table** tab.
4. Use **Filter > Filter All** to perform any desired filtering.

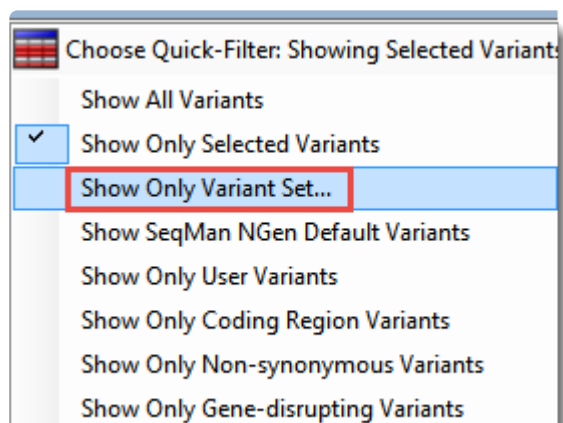


5. At the conclusion of filtering, click the **Remember Results as a Variant Set** tool () above the Search Results table. Type in any desired name and press **OK**.

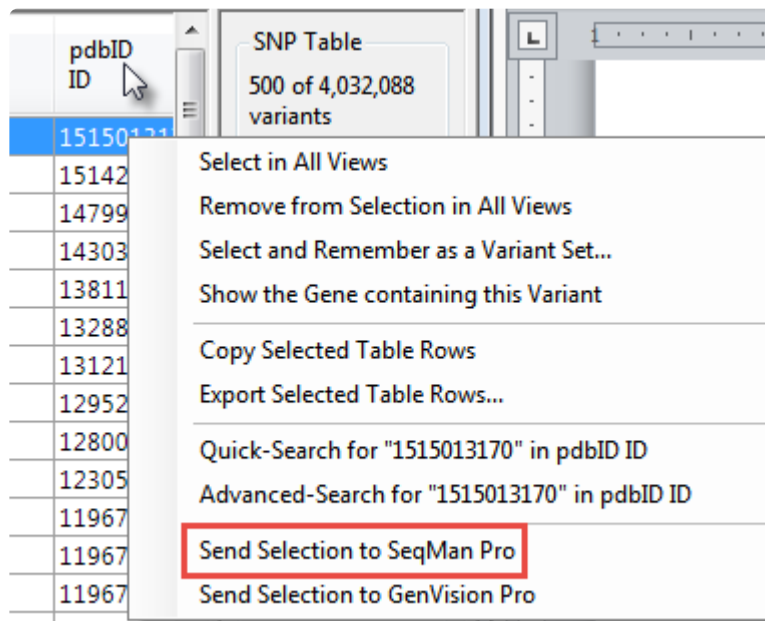


6. In the Action section in the center-right of the window, click the link **Select and show the table of this set's Variants**.
7. Use the **Choose Quick-Filter** drop-down to select **Show Only Variant Set**. In the ensuing dialog box,

select the named set and click **OK**.



8. Click on **Add/Manage Columns**. Select **SeqMan NGen Assembly Variants**, then **pdbID**. Click **Add Column** and **OK**.
9. In the SNP Table, click on the **pdbID** column header to sort the column and locate rows with PDB entries.
10. Within the subset of rows with PDB entries, select any row of interest. Then right-click on it and choose **Send Selection to SeqMan Pro**.

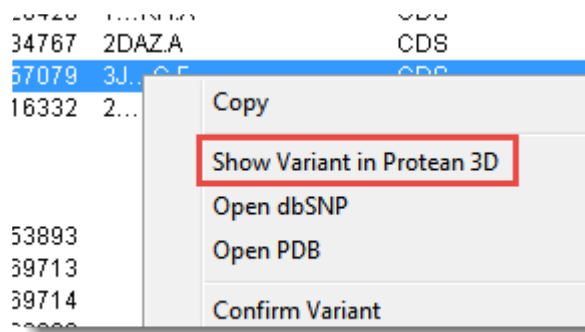


11. If prompted, select an individual sample of interest and press **OK**.

Part C: View variants in SeqMan Pro or SeqMan Ultra:

Instructions below pertain to SeqMan Pro, but are similar in SeqMan Ultra.

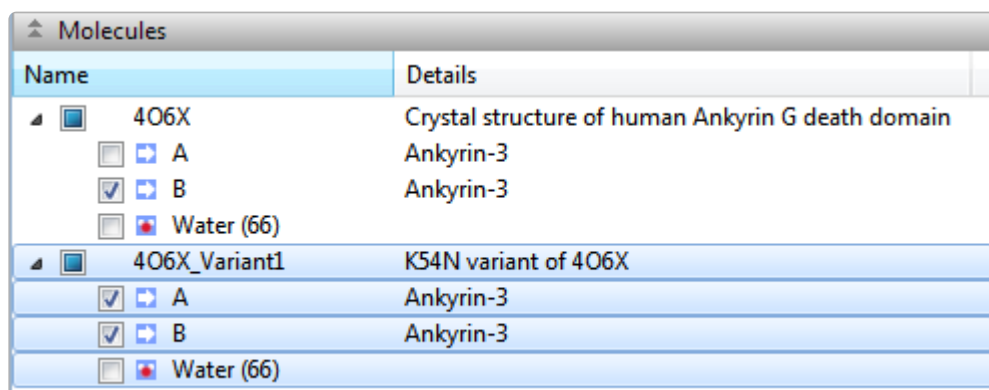
1. If you are coming from Part B, above, the sample is automatically selected and its Alignment view is opened. Continue directly to Step 2. If you are coming from Part A, launch SeqMan Pro and use **File > Open** to open the .assembly file.
2. Choose **Variant > Variant Report**. (If desired variants are being filtered out, you may need to click on the **Show All** button.)
3. Click on the **PDB ID** column header to sort items with PDB IDs to the top. Select one or more rows, then right-click within the selection and choose **Show Variant in Protean 3D** (or use **Variant > Show Variant in Protean 3D**).



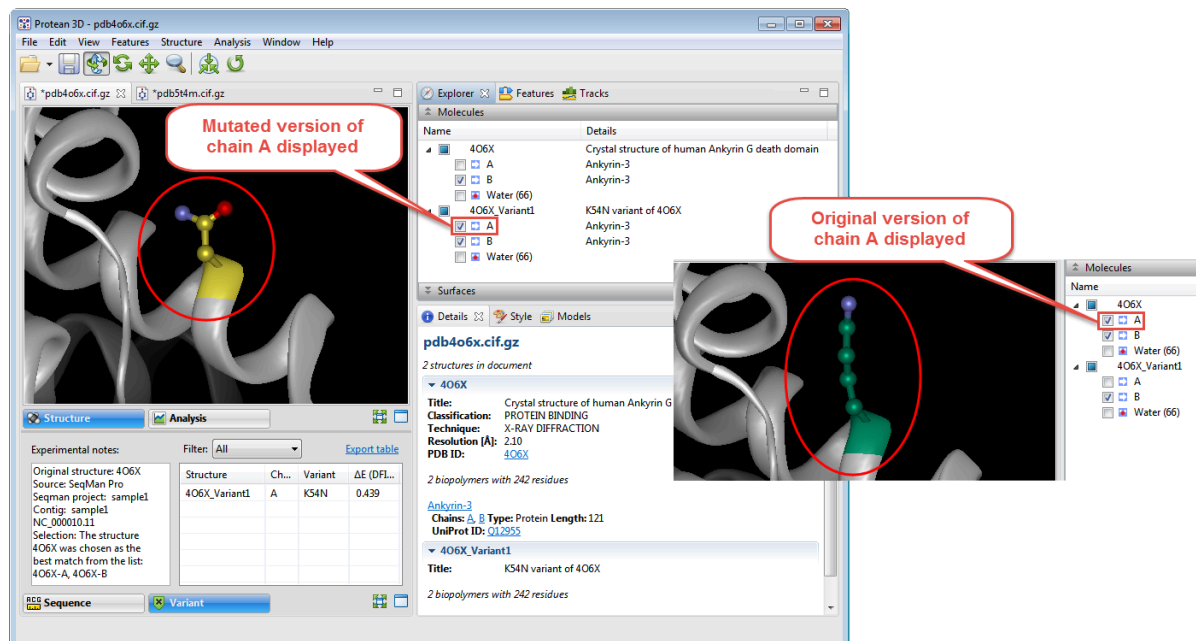
Part D: View the protein structure in Protean 3D:

After finishing Part C, the protein structure with the variant of interest opens in Protean 3D.

In the Molecules area, two near-identical copies of the structure appear. The upper structure is the original structure from the Protein Data Bank. The lower structure is the variant version calculated by Protean 3D.



- Use the Structure view to observe the mutated side chain along the backbone. To show/hide the different versions of each chain, check/uncheck boxes in the Molecules area.



- To see notes about the structure chosen as the best match, look in the Experimental notes box at the bottom of the Variant view. If the Variant view is not visible, click on the **Variant** tab at the bottom of the Protean 3D window.


Original structure: 4WKA
 Source: SeqMan Pro
 Seqman project: sample1
 Contig: sample1 NC_000001.11
 Selection: The structure 4WKA was chosen as the best match from the list:
 1GUV-A, 1HKI-A, 1HKJ-A, 1HKK-A, 1HKM-A,
 1LG1-A, 1LG2-A, 1LQ0-A, 1WAW-A,
 1WB0-A, 4WJX-A, 4WK9-A, 4WKA-A,
 4WKF-A, 4WKH-A

✿ **Note:** If the residue has no atoms, you will see an error message in this box. In this case, return to SeqMan Pro and choose a different row in the Variants table.

Protean 3D uses several metrics to determine the “best” PDB file to display when a variant is located in a CDS that is associated with more than one PDB file. Quality is the first consideration, with high-resolution crystal structures > NMR structures > low-resolution crystal structures > other techniques. This ordering is refined by alignment to the corresponding Uniprot sequences. This refinement considers the percent match and the number of gaps before the variant position. If two structures are still tied as the “best,” the largest structure is chosen.

- To predict whether the mutation is stabilizing or destabilizing to the protein structure, use the table on the right of the Variants view. The **delta-E (DFIRE-A)** column displays the change in energy value based on the DFIRE calculation ([reference](#)). This number can be used to predict whether the mutation is stabilizing or destabilizing to the protein structure. A positive number is considered destabilizing to

the structure when compared to the original amino acid; a negative number is considered non-destabilizing.

 **Note:** If multiple entries appear in the table, you can sort them by clicking on their column headers, or filter them using the **Filter** drop-down menu. If desired, you can save the table by clicking the **Export table** link.

- To further explore potential impact of mutation on protein stability and function, apply a solvent-accessible surface.
- Analyze secondary structure characteristics to interrogate the mutation's effect on protein flexibility, amphiphilicity, charge density, hydrophathy, and more.

Remove PhiX control reads from Illumina data prior to import

During *de novo* assembly, contamination of Illumina data with PhiX control sequence may result in the generation of spurious contigs. For background information, see [Mukherjee et al., 2015](#).

Note that:

- Not all Illumina data are contaminated with PhiX.
- PhiX contamination is not a major concern for reference-guided assemblies.
- In informal tests at DNASTAR, the amount of contamination in most data sets was so low that the spurious contigs were automatically discarded for being “under the minimum coverage” per the SeqMan NGen defaults. And when coverage was higher, contaminated contigs contained *only* PhiX174 reads and could be readily recognized post-assembly.

If you are following a *de novo* workflow, you can easily remove PhiX contamination prior to assembly by following these steps:

1. Download the sequence [NC_001422.1](#) (*Enterobacteria phage phiX174 sensu lato, complete genome*) from the NCBI website.
2. Launch SeqMan NGen.
3. In the [Workflow](#) screen, choose a *de novo* assembly option.
4. In the [Assembly Options](#) or [Preassembly Options](#) screen, check the **Contaminant Scan** box. Use the associated **Add** button to add the PhiX174 sequence downloaded in the previous step.
5. Proceed through the rest of the SeqMan NGen wizard screens and assemble as usual.

Any PhiX174 sequence will be removed prior to assembly.

Metagenomics workflows

The following table describes each of the workflows available in the **Metagenomics** tab of the [Workflow](#) screen.

Group	Workflow	Description
ABI / Sanger	De novo	Cluster metagenomic and other “mixed” sequence samples into contigs. This workflow allows optional removal of specified host DNA. The default parameters for this workflow have been optimized to take into account the short read lengths and presence of repetitive DNA sequences common to metagenomic and 16S rRNA data.
	Reference-guided	Align/sort metagenomic and other mixed sequence samples to a set of reference sequences, including microbial genome database. This workflow allows optional removal of host DNA before assembling/aligning the remaining reads onto one or more reference sequences/templates. The default parameters for this workflow have been optimized to take into account the short read lengths and presence of repetitive DNA sequences common to metagenomic and 16S rRNA data.
NGS-based	De novo	Cluster metagenomic and other “mixed” sequence samples into contigs. This workflow allows optional removal of host DNA. The default parameters for this workflow have been optimized to take into account the short read lengths and presence of repetitive DNA sequences common to metagenomic and 16S rRNA data.
	Reference-guided	Align/sort metagenomic and other mixed sequence samples to a set of reference sequences, including microbial genome database. This workflow allows optional removal of host DNA before assembling/aligning the remaining reads onto one or more reference sequences/templates.. The default parameters for this workflow have been optimized to take into account the short read lengths and presence of repetitive DNA sequences common to metagenomic and 16S rRNA data.

RNA-seq/transcriptomics workflows

The following table describes each of the workflows available in the **RNA-seq/transcriptomics** tab of the [Workflow](#) screen.


Group	Workflow	Description
Quantitative analysis	RNA-seq	<p>RPKM gene expression quantification and differential gene expression using DESEQ2 and EdgeR from BioConductor. First-pass assembly is done using the usual XNG assembler. Second-pass assembly utilizes the QNG analysis module to determine expression level statistics. For each sample in a project, two new files for each contig/chromosome are put into the .assembly package. These contain the QNG calculated expression values for each gene and its isoforms: -[contig number].genes-features and -[contig number].isoforms-features.</p> <p>To learn how to use the output of an RNA-Seq <i>de novo</i> transcriptome assembly as input for the RNA-Seq reference-guided workflow, see Use RNA-Seq de novo transcriptome output as a reference.</p> <p>After performing an RNA-Seq reference-guided assembly, you can view the results in any of three applications:</p> <ul style="list-style-type: none"> • SeqMan Pro – Use SeqMan Pro’s 3-tabbed Feature Table for downstream analysis. Each tab (All Features, Gene Features, CDS Features) displays expression values in a column entitled “RPKM”. If the sample is part of a replicate set, a second column entitled “RPKM – Replicate” displays the expression value for the feature determined from the replicate set. • GenVision Pro – Display a Sashimi plot for the assembly. Sashimi plots are designed to display data indicative of mRNA splicing, and are generated automatically during RNA-Seq assembly. See RNA-Seq reference-guided workflow output for a list of output files resulting from this type of assembly. • ArrayStar – Use ArrayStar’s Gene and Isoform tables to filter for differentially-expressed genes of interest. ArrayStar tables can also display any DESeq2 or edgeR statistics included in the assembly.
	ChIP-seq	Choose from several different normalization and peak detection methods, including ERANGE and MACS.
	miRNA	miRNA gene expression quantitation and discovery of new miRNAs.
	De novo	Large capacity assembly of transcriptome sequence data with auto-annotation of

Assembly	transcriptome	assembled transcripts. In the past, <i>de novo</i> assembly of RNA-Seq data could result in thousands of contigs representing the expressed transcripts, without any context or labels. For Lasergene 13.0 and later, SeqMan NGen automatically attempts to group contigs from the same gene, and then name and annotate them based on the best match to a collection of annotated reference sequences. Two different SeqMan NGen assembly engines are used to optimize your results. Note that results from this workflow are non-quantitative. Result files for this workflow are described in detail in RNA-Seq <i>de novo</i> transcriptome workflow output . (Also called the “StarBlast” workflow)
	De novo miRNA	<i>De novo</i> assemble novel miRNAs.

Include DESeq2 or edgeR statistics

To view statistics from DESeq2 or edgeR in ArrayStar, you first need to create one or more SeqMan NGen assemblies in which one of these statistical packages was specified as the [normalization method](#). That specification also triggers the use of the package for differential expression analysis.

1. In SeqMan NGen:
 - a. Choose a reference-guided RNA-Seq [workflow](#).
 - b. In the [Analysis Options](#) screen, specify that **DESeq2** or **edgeR** be used as the **Normalization method**.
 - c. Run the assembly.
 - d. If you performed the assembly on a Windows machine and will be doing the ArrayStar analysis on the same machine, click the **Compare variants / differential gene expression between samples** button in the [Assembly Summary](#) screen. Then skip ahead to Step 3. Otherwise, continue to Step 2.
2. If you did not press the **Compare variants / differential gene expression between samples** button in SeqMan NGen, you can use an alternative method to open the assembly in ArrayStar:
 - a. Launch ArrayStar and begin an RNA-Seq project.



DNASTAR
ArrayStar®

Experiment Type

☐ Microarray gene expression analysis
☐ Variant comparison
☐ Copy Number Variation (CNV)
☒ RNA-Seq gene expression analysis
☐ ChIP-Seq analysis
☐ miRNA analysis

- b. Use ArrayStar's Project Setup wizard to import one or more *.assembly* packages created in Step 1, above.
- c. In the Set Up Preprocessing page, choose **DESeq2**, or **DESeq2-Local**, **edgeR**, or **edgeR-Local**.

RNA-Seq Project Setup

Set Up Preprocessing

Preprocessing Method

Select the desired processing: QSeq

Quantifies gene expression values for RNA-Seq, Copy Number expression levels by mapping DNA sequences to a reference genome for this purpose.

Preprocessing Parameters

Normalization method: DESeq2

Please specify the sequencing method: flhC_del flhC_del 1-0.1

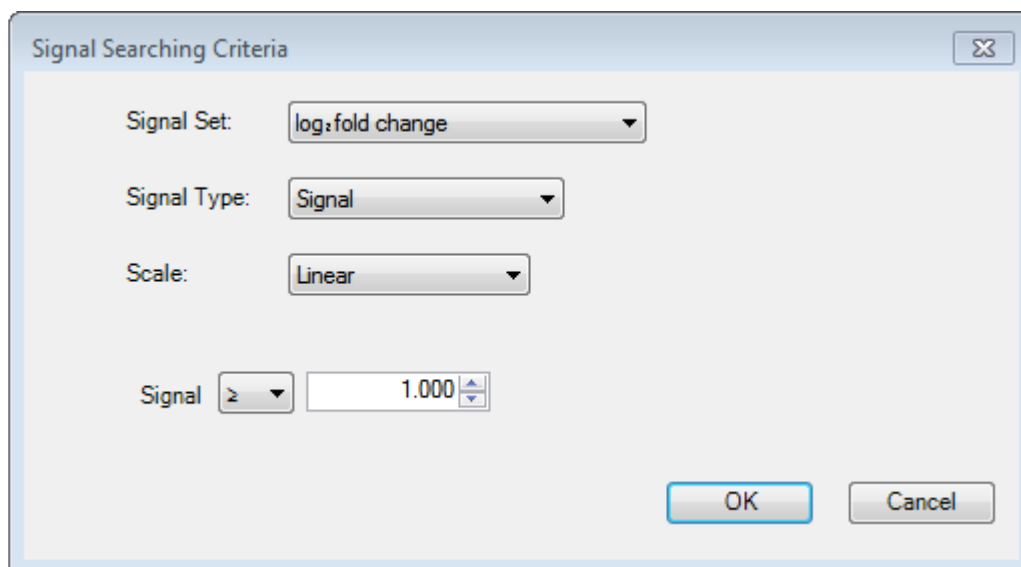
None
 RPM
 RPKM
 RPK
 DESeq2
 DESeq2-Local
 edgeR
 edgeR-Local

The Genomics Cloud-based options (without “**-Local**”) are suitable for most users. Windows users who have already installed R, Bioconductor and DESeq2 can elect to run locally by selecting the version ending in **-Local**.

3. Access DESeq2 or edgeR statistics in ArrayStar using either of these methods:

- Open the Gene or Isoform tables and use the Add/Manage Columns tool to add DESeq2-related columns from the **Gene Values** or **Isoform Values** tabs.
- Use **Filter > Filter All** to open the Advanced Filtering dialog. In the header, elect to search for **Genes** or **Isoforms**. Set up one or more filter rows. In each row, use the left-most drop-down menu to select **Signal Sample Group**, and use the middle section of each row to set up a comparison of interest. Finally, click the **Choose Signal Criteria** button in each row and choose the desired settings.

For example, in the first row you might search for reads with a **log₂ fold change** greater than or equal to **1**.

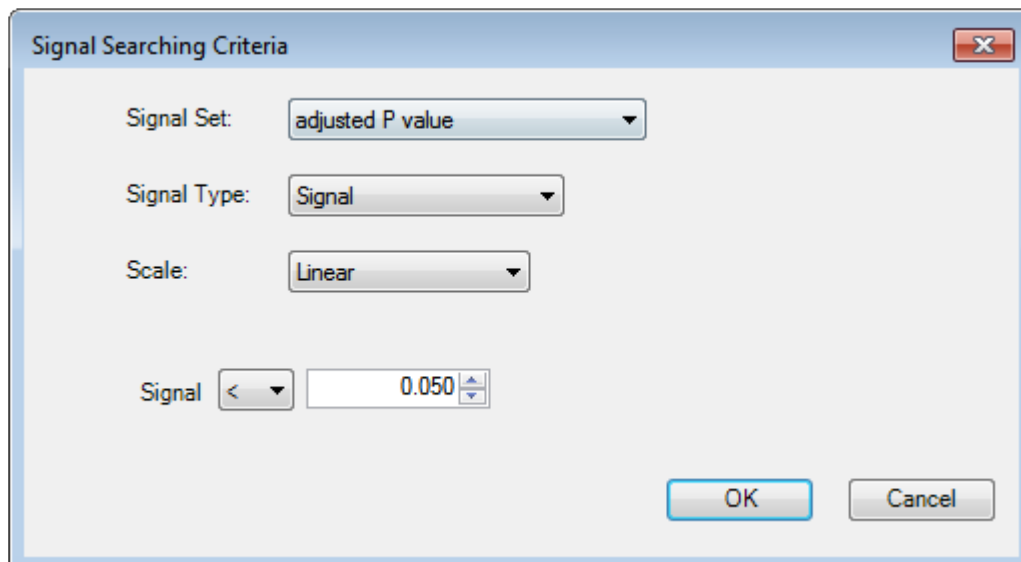


The image shows a dialog box titled "Signal Searching Criteria". It contains four rows of settings:

- Signal Set:** A dropdown menu with "log₂ fold change" selected.
- Signal Type:** A dropdown menu with "Signal" selected.
- Scale:** A dropdown menu with "Linear" selected.
- Signal:** A dropdown menu with "≥" selected, followed by a text input field containing "1.000".

At the bottom right, there are two buttons: "OK" and "Cancel".

In the second row, you could search for items with an **adjusted P-value** less than **0.05**.



Use the **Search** button and then select all the filtered isoforms. Click on the **Remember Results as an Isoform Set** tool, choose a name and press **OK**. In the Set List, click the link **Show the** table of Genes containing this set's Isoforms. In the Gene table, use the Add Fold Change tool to add the **Fold Change** column.

4. (optional) To view or copy a Bioconductor script/log which you can use for QC/QA purposes or include in publications, open the Experiment List, select an experiment, and use the Show log links in the Info Pane on the right.

Variant analysis/resequencing workflows

The following table describes each of the workflows available in the **Variant analysis/resequencing** tab of the [Workflow](#) screen.

Group	Workflow	Description
ABI / Sanger	Whole genome	Align Sanger trace data from one or multiple samples to a genomic reference or genome template package for accurate SNP/Indel analysis. This type of assembly can include billions of reads and large eukaryotic genomes. After assembly, compare results in ArrayStar using the SNP Report.
	Amplicon	Align Sanger trace data from one or multiple samples to targeted genes or genomic regions for accurate SNP/Indel analysis. Assembles a region of interest produced by PCR amplification.
	Clone verification	Align reads to confirm clone integrity and insert orientation. (Note: For a dedicated clone verification workflow, see this topic in the SeqBuilder Pro User Guide).
NGS-based	Whole genome	Align NGS sequence data from one or multiple samples to a genomic reference or genome template package for accurate SNP/Indel analysis. This type of assembly can include billions of reads and large eukaryotic genomes.
	Amplicon, gene panel, exome	Align NGS sequence data from one or multiple samples to targeted genes or genomic regions for accurate SNP/Indel analysis. Gene panels look at specific gene regions, usually those corresponding to known defects. Exome assembly saves assembly time and resources by specifically targeting only exons and coding regions, but do require you to have the corresponding <i>.bed</i> file from the capture kit. For instance, if you used Human Genome build 38 as the reference, for example, the corresponding <i>.bed</i> file might be called <i>Human genome build38.bed</i> . If using this workflow with cancer samples, check the box next to Somatic/Cancer/Heterogeneous in the Analysis Options screen. In most cases, downstream analysis of these finished assembly will place in ArrayStar.
	Viral-host integration detection	Locate prophage and retro-viral insertion sites in host genome. Available in Lasergene 17.1 and later. Used to locate putative viral insertion sites or to predict the location of other inserted sequences, such as transposable elements. When you select this workflow, SeqMan NGen automatically sets up a templated assembly that is optimized for locating viral insertion sites. Since chimeric reads (sequences consisting of both host and viral DNA) usually indicate viral insertion sites, SeqMan NGen looks for chimeric reads in a multi-step process. First, the viral genome is used as the initial assembly template. Next, the sub-set of reads that mapped to the viral genome is then re-assembled against the host template. During both reference-guided assembly steps, SeqMan NGen “masks” (trims) whichever half of the chimeric read does not match the template for that step. The host template assembly results are output in BAM file format. SeqMan Ultra is used to explore

		possible viral insertion sites post-assembly. Launch SeqMan Ultra and use Contig > Contig Coverage to view tabular data for the individual contigs. Navigate to positions with multiple reads, as evidenced in the depth column. The reads at these positions should be trimmed to the same base indicating the insertion site. You may “untrim” the reads to verify that they also contain viral sequence.
Variant Call Format (VCF) files	Functional annotation of a single sample	Annotates the variant positions with functional information from a database, including affected genes and impact on protein encoding regions and/or splice sites.
	Annotation and comparison of multiple samples	Allows multiple samples in VCF format to be annotated and then compared to identify genes and/or variants of interest in ArrayStar. This workflow is designed to use with assemblies created outside SeqMan NGen (e.g., using BWA + GATK). Such assemblies often have .vcf files as their only output.

Variant calling accuracy workflow

While SeqMan NGen 17 no longer has a workflow named “Variant calling accuracy” (AKA “validation control accuracy,” “reference SNP accuracy”), you can still perform this workflow as follows:

1. In the [Workflow](#) screen, select the **Variant analysis/Resequencing tab** and choose the “NGS-based” **Amplicon, gene panel, exome** workflow.
2. In the [Reference Sequence](#) screen, [add the relevant DNASTAR genome package](#); in most cases, this will be “human”. Check the **VCF file** and browse to a VCF file of the true variants. If you have a *.bed* file, you can also check the **BED file** box (optional) and browse to that file.
3. In the [Input Sequences](#) screen, add read data for the Validation Control.
4. In [Assembly Options](#), check the box next to the desired **Variant detection mode**.
5. Follow the rest of the wizard screens and initiate the assembly.
6. Variant calling accuracy is evaluated within ArrayStar. Once assembly is complete, see the ArrayStar help topic [Validation Control Accuracy](#) for further instructions.

Variant Call Format (VCF) files workflows

The following table describes each of the workflows available in the **Variant Call Format (VCF) files** tab of the [Workflow](#) screen. These workflows are used if you have assembled your data and done SNP calling in another application and have VCF files with a *.vcf* or *.abi* file extension that you wish to compare to a database or to one another. These workflows both output an *.astar* file that can be opened in ArrayStar for downstream analysis.

Group	Workflow	Description
VCF File Analysis	Functional annotation of a single sample	Annotates the variant positions with functional information from a database, including affected genes and impact on protein encoding regions and/or splice sites.
	Annotation and comparison of multiple samples	Allows multiple samples in VCF format to be annotated and then compared to identify genes and/or variants of interest in ArrayStar. This workflow is designed to use with assemblies created outside SeqMan NGen (e.g., using BWA + GATK). Such assemblies often have <i>.vcf</i> files as their only output.

Combine/Reanalyze Existing Assemblies

The following table describes each of the workflows available in the **Combine/Reanalyze Existing Assemblies** tab of the [Workflow](#) screen.

Group	Workflow	Description
Combine/ Reanalyze	Combine existing assemblies	Once you have prepared individual assemblies for projects of the same type (e.g., two or more RNA-Seq assemblies) using the same reference sequence, this workflow lets you merge them into a single ArrayStar (. <i>astar</i>) project. You can then perform downstream analysis in ArrayStar and/or send data from ArrayStar to DNASTAR's genome visualizer, GenVision Pro.
	Rerun analysis of existing assemblies	<p>This workflow, available to local users only, reanalyzes existing .<i>assembly</i> packages from one or more workflows and puts them into a single ArrayStar-readable project (.<i>astar</i>). First use SeqMan NGen to create templated assemblies for various projects, using the same reference sequence for each. Then create the merged SeqMan NGen assembly by following this workflow. Perform downstream analysis in ArrayStar and/or send data from ArrayStar to DNASTAR's genome visualizer, GenVision Pro.</p> <p>Examples of when reanalysis is indicated include:</p> <ul style="list-style-type: none"> • After adding a new replicate • After changing the normalization method • After adding VAD information.

Analysis Options

The Analysis Options wizard screen allows you to specify the analysis parameters to use for your assembly. This screen comes in several variations, one of which is shown below:

SeqMan NGen

File Help

Analysis Options

Set the options for a post-assembly analysis

- ✓ Workflow
- ✓ Reference Sequence
- ✓ Input Sequences
- ✓ Assembly Options
- Analysis Options**

☒ Detect SNPs and other small variants

Variant detection mode: ☒ Diploid ☐ Haploid ☐ Somatic / cancer / heterogeneous

Gender: Select-

SNP filter stringency: ☐ High ☐ Medium ☒ Low

☐ Annotate with the Variant Annotation Database (VAD)

Includes: Mastermind literature citations from Genomenon
1000 Genomes Project and Exome Variant Server allele and genotype frequencies
dbNFSPv4.1 (Clinvar, SIFT, GERP++, more...)

[Variant Annotation Database FAQ](#)

☐ Detect CNVs

CNV normalization method: RPK_CN

☒ Use features of type(s): Select Features...

☐ Exclude pseudogenes

☒ Detect structural variations (larger deletions and insertions)

Deletions - minimum depth:

Advanced Analysis Options

?

< Back Cancel

Depending on the [workflow](#), only a subset of the following options will be available:

Category	Options and Descriptions
----------	--------------------------

SNP detection

To enable SNP detection, start by checking **Detect SNPs and other small variants**. If this box is checked, you can later open the assembly in ArrayStar or SeqMan Pro to view the SNPs.

Use **Variant detection mode** to specify genome ploidy for SNP detection purposes. Choosing **Haploid** or **Diploid** establishes the statistical model SeqMan NGen will use in estimating the probability that a given called variant is real (i.e., that the sequence really differs from the reference). Selecting **Somatic/cancer/heterogeneous** (e.g. for a polyploid genome, cancer panel, etc.) prevents SeqMan from calculating probabilities.

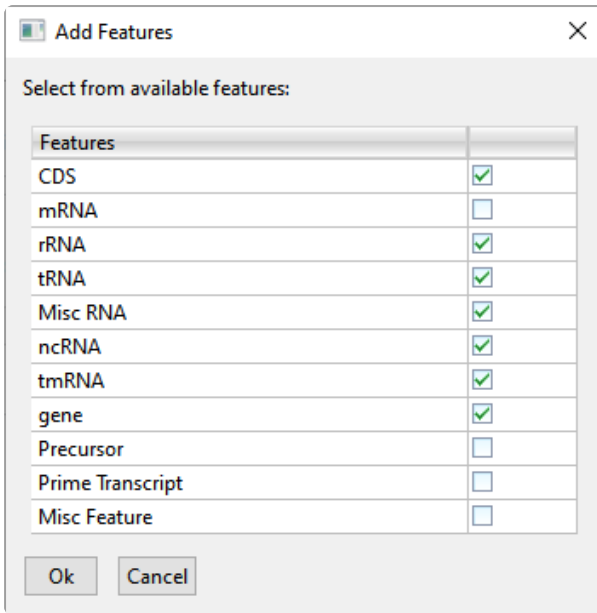
If the **Gender** checkbox is present, specify the gender of the subject (**Male/Female**), if known. Otherwise, select **Unknown**. This checkbox appears only if you are using a DNASTAR genome template package and have chosen a genome ploidy other than **Haploid**.

SNP filter stringency specifies two key settings for placing a read in the layout. When building an assembly, SeqMan NGen uses a three stage strategy: overlap, layout, and alignment. In the overlap stage of a reference-guided assembly, for example, each read and the reference are broken up into an overlapping set of substrings or “mers” of a specified length (“mer length” or “mer size”). Identical mer matches are an indication that the read matches the reference at that position. The more overlapping mers between two sequences, the stronger the indication that the match is real. The layout stage uses that overlap information and attempts to place each read in its true position on the reference. The final layout of all the reads is then sent to an aligner that produces the final fully gapped alignment. Layout stringency settings can be used to adjust the extent of overlap data required to include a putative match in the final layout.

The radio buttons specify stringency levels for “soft” filtering of SNPs. Soft filtering means that SNPs of the least interest to you will be automatically hidden when SNP reports/tables are viewed in SeqMan Pro, SeqMan Ultra or ArrayStar. Your selection in this screen controls the three assembly parameters shown in the table below. For more information on PnotRef, see [Filter based on P not Ref](#).

- **High** has Min SNP%=15, PnotRef%=99.9, and Depth=20. This option has a lower false discovery rate (FDR) for SNPs and is recommended for whole genome workflows.
- **Medium** (where available) has Min SNP%=15, PnotRef%=99, and Depth=20.
- **Low** has Min SNP%=15, PnotRef%=90, and Depth=20. This option has a higher true positive rate (TPR) for SNPs and is recommended for all workflows other than whole genome.

Note 1: If a BED, manifest and/or VCF file was specified during project setup and a SNP table is opened in SeqMan Pro or ArrayStar, then only the variants in the targeted regions and at the positions specified in the VCF within those targeted regions will be shown by default.

	<p>Note 2: These “soft filtered” SNPs are <i>not</i> removed from the assembly, and can be made visible again by changing the SNP filtering parameters in SeqMan Pro, SeqMan Ultra or ArrayStar. This is in contrast to “hard filtering” of SNPs, which is done through the Layout tab (see information below this table).</p>
Variant Annotation Database	<p>Check Annotate with Variant Annotation Database (VAD) if you are working with human samples and would like to import variant annotations from a specific portion of the NCBI RefSeq database maintained on the DNASTAR website. This checkbox is only available for human samples assembled against builds 37 or 38. This is an extremely powerful analysis tool that we highly recommend using any type you are analyzing human samples. To learn more about the Variant Annotation Database, click on the link in this section to read our Frequently Asked Questions (FAQ).</p>
CNV detection	<p>Check Detect CNVs if you wish to calculate copy number variants (CNV) as part of the assembly. If the box is checked, you may choose between two CNV Normalization method options: RPK-CN and None (i.e., no data normalization). If you also select a Variant detection mode other than Do not calculate variants, then CNVs, SNPs and small indels will be calculated from the assembly. After assembly, you can then use ArrayStar to view all three types, or SeqMan Pro or SeqMan Ultra to view only the SNPs and small indels.</p> <p>Check Use features of type(s) to only report results when a specific type of feature is used as the target for mapping reads. Note that mapping occurs regardless of the type of feature annotation. However, when you check this option, the mapping results for unwanted feature types will not be reported. Put checkmarks next to the feature types you wish to use, then press OK.</p>  <p>If you leave Use features of type(s) unchecked or if the reference sequence has no feature annotations, each individual sequence in the reference set will be used as a separate transcript</p>

	<p>(i.e., a single gene feature).</p> <p>Check Exclude pseudogenes to <i>not</i> report mapping results for features with /pseudo in their annotations. As with the previous option, mapping occurs regardless of the type of feature annotation.</p>
RNA-Seq normalization	<p>Check Normalize RNA-Seq values if you want to apply a normalization method to the data on a per-isoform basis. If you check this option, use the drop-down menu to choose the desired RNA-Seq normalization method. In order to enable this option, some workflows require you to check the Calculate Copy Number Variation option. See CNV detection, above, to learn about the Use features of type(s) and Exclude pseudogenes options.</p>
ChIP-Seq normalization	<p>Check Set ChIP-seq peak detection method if you want to apply a normalization method to the data. If you check this option, use the drop-down menu to choose the desired ChIP-seq peak detection method. See CNV detection, above, to learn about the Use features of type(s) and Exclude pseudogenes options.</p>
Structural variations	<p>If you want SeqMan NGen to flag structural variations for viewing during downstream analysis, check the box next to Detect structural variations (larger deletions and insertions). If desired, type in a minimum depth to be considered a deletion.</p>

To access additional options, click the **Advanced Analysis Options** button to open a multi-tabbed dialog. Each tab has changeable parameters for different parts of the analysis process. Different workflows have different subsets of tabs. In addition, tabs with the same names may contain different options depending on the workflow. For details on each tab, see [Peak Detection tab](#), [Variants tab](#) and [Layout tab](#).

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

RNA-seq normalization methods

“Normalization” refers to the standardization of sequencing data on the basis of sequencing depth and gene length. Some versions of the Assembly Options”:#assembly-options screen allow you to specify a data normalization method, or to select **None**, in which case data will not be normalized. Some methods are described in the table below, while **DESeq2** and **edgeR** are discussed below the table.

Normalization method	Description
zRPKM	zRPKM (Krumm et al., 2012) is available for the CNV workflow and is calculated: $zRPKM = (RPKM_{\text{exon, sample}} - \text{Median}_{\text{exon}}) / \text{StdDev}_{\text{exon}}$. This is the optimal normalization method for CNV projects with at least three groups/projects. Otherwise, we recommend using the RPK-CN normalization method.
RPK-CN	RPKM-CN (Krumm et al., 2012) is available for CNV experiments and is calculated as: $RPKM-CN = RPKM / \text{median of the exon's RPKMs}$; where $RPKM > 1$. RPKM-CN calculates the copy number by taking the ratio of the RPKM of an exon versus the median RPKM of any exon in the experiment. The final number is a ratio (or log ratio) indicating a relative copy number with no units, since the units are cancelled out in the ratio. The variable M is a constant: the number of millions of mapped reads in the experiments. The ultimate meaning of the ratio comes from the different reads “R” and length “K” of each exon and the median. The constant, M, drops out of the equation and only affects scaling for initial filtering-out of low-coverage exons. We only recommend using RPKM-CN if you don’t have enough samples to provide a good standard deviation for each exon when using the zRPKM normalization method. Otherwise, zRPKM is the preferred method for the CNV workflow.
Quantile	Quantile is available for RNA-Seq workflows only. Quantile normalization adjusts all of the values in your project so that the distribution is the same across all of the experiments.
RMP	RPM (reads assigned per million mapped reads) is available for RNA-Seq, ChIP-Seq, and miRNA experiments, and is the only normalization method available for ChIP-Seq and miRNA experiments. When RPM is selected the signal values for each experiment will be divided by the total number of mapped reads divided by one million.
RPKM	RPKM (reads assigned per kilobase of target per million mapped reads) is available for RNA-Seq data. When RPKM is selected, the signal values for each experiment will be divided by the total bases of target sequence divided by one thousand; and the resulting number divided by the total number of mapped reads divided by one million.

DESeq2 and edgeR:

[DESeq2](#) ([Love et al. 2014](#)) and [edgeR](#) ([Robinson et al. 2010](#)) are statistical packages in [Bioconductor](#) used to assess differential expression in RNA-Seq experiments.

DESeq2 or edgeR statistics for an assembly can be analyzed by opening the assembly in ArrayStar. For information about setting up an assembly suitable for analyzing DESeq2 or edgeR statistics in ArrayStar, see [Create an assembly using DESeq2 or edgeR statistics](#).

Both methods require a control group to be specified, and both require replicate samples for each experimental condition and for the control. Note that when multiple experimental conditions are being considered, the same control group is used for multiple tests. The original P-values from the statistical tests are then adjusted using the [Benjamini-Hochberg](#) (1995) procedure.

Differences between DESeq2 and edgeR are shown in the table below:

Calculation	DESeq2	edgeR
Normalization method	<p>Uses a median of ratios method to normalize read counts to account for sequencing depth and RNA composition. Provides two methods: regularized logarithm (<i>rlog</i>) and Variance Stabilizing Transformations (VST).</p> <p>DESeq2 does not attempt to account for transcript length since it is comparing counts between samples for the same gene and assumes the length does not change. This assumption holds true except in rare cases where the dominant transcript length changes between samples due to alternative splicing for example.</p>	<p>Uses "trimmed mean of M-values" (<i>TMM</i>) (Robinson & Oshlack, 2010 topic=Research References). The TMM normalized read count can be viewed in the ArrayStar tables, where counts are represented as $\log_2(\text{counts-per-million-reads})$.</p> <p>Normalized counts generated by a different method, <i>RLE</i>, are also available within ArrayStar but these values are not used for the actual statistical tests. RLE is similar to the RLOG normalization method used by DESeq2.</p>
Statistical tests for differential expression	<p>DESeq2 uses raw counts, rather than normalized count data, and models the normalization to fit the counts within a Generalized Linear Model (GLM) of the negative binomial family with a logarithmic link. Statistical tests are then performed to assess differential expression, if any.</p>	<p>Data are normalized to account for sample size differences and variance among samples. The normalized count data are used to estimate per-gene fold changes and to perform statistical tests of whether each gene is likely to be differentially expressed.</p> <p>EdgeR uses an exact test under a negative binomial distribution (Robinson and Smyth, 2008 topic=Research References). The statistical test is related to Fisher's exact test, though Fisher uses a different distribution.</p>
Data reporting	<p>In ArrayStar, the <i>rlog</i> values are used by default in the scatter plot and for clustering. VST values are</p>	<p>In ArrayStar, the $\log_2(\text{CPM})$ values calculated using TMM are used by</p>

method	displayed as Gene Table data columns.	default in the scatter plot. In the Gene Table, values for fold change compared to the control are represented as $\log(\text{fold change})$.
---------------	---------------------------------------	--

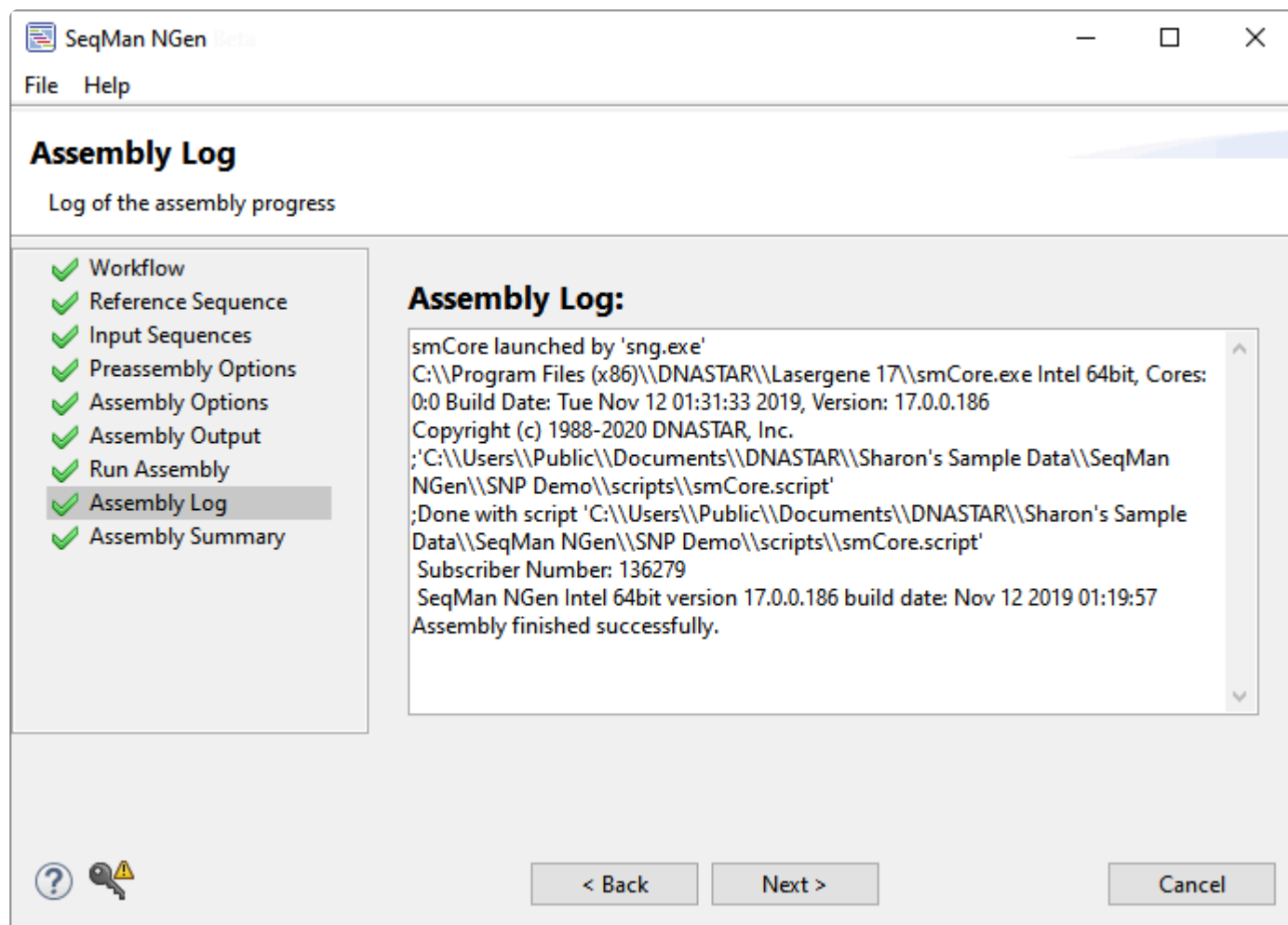
ChIP-seq peak detection methods

If you are following the ChIP-seq [workflow](#), the [Analysis Options](#) screen allows you to specify a peak detection method. Available methods are described in the table below:

Name	Description
MACS	<p>The MACS Peak Finder is based on the peak detection algorithm (Zhang, et al., 2008). This is a model-based algorithm that expects there to be paired peaks of reads on either side of a true binding site. Paired flanking peaks reflect the fact that only the 5' ends of immuno-precipitated fragments are usually sequenced. Because of this, the majority of the locations of sequences associated with peaks don't correspond to the location of the binding site.</p> <p>The MACS algorithm attempts to build a model of the distance between these peaks and takes this distance into account to shift reads forwards or backwards, resulting in a peak centered over the true binding site.</p> <p>The MACS Peak Detection algorithm reports the number of reads within a peak as the signal value for that peak. It also calculates a P-value based on the distribution of reads near a peak region to try to compensate for uneven background noise across the genome. When present, control data are used to filter the peaks that are called and to assign each peak an FDR score which is the false discovery rate likelihood that the peak is not valid.</p>
ERANGE2	<p>The ERANGE2 Peak Finder (Johnson et al., 2007) is a simple "sliding window" peak detection algorithm that looks for a specified number of reads within a window of a specified length. If a peak is found, it is extended as long as there are reads within the window width. If control data are present, it can be used to disqualify any peaks that do not have a minimum fold enrichment over the same region in the control data.</p>
ERANGE3	<p>The ERANGE3 Peak Finder is based on the ERANGE 3.1 Algorithm for ChIP-Seq and RNA-Seq Analysis (see Mortazavi et al., 2008). This peak detection algorithm calculates peaks in a normalized reads-per-million space. Features of this algorithm include simple read shifting and repeat read handling. This algorithm also considers the directionality of reads when calling peaks.</p>

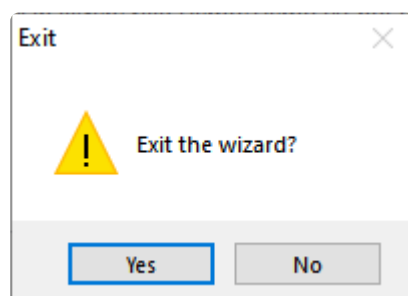
Assembly Log

After pressing the **Run assembly on this computer** link from the [Run Assembly](#) screen, the Assembly Log opens to show the status of the assembly.



Once the assembly runs to completion without failure, the following text will be displayed: “Assembly finished successfully.”

Click **Next >** to proceed to the [Assembly Summary](#) screen or **< Finish** to exit from SeqMan NGen. In the latter case, the following confirmation popup will appear:



Assembly Options

The Assembly Options wizard screen allows you to specify assembly parameters. This screen comes in several variations, two of which are shown below (click on either to see full-size):

SeqMan NGen
File Help

Assembly Options

Set the most important options for a successful assembly

- ✓ Workflow
- ✓ Input Sequences
- ✓ Preassembly Options
- Assembly Options**

Assembly options

Coverage Calculation

☒ Estimated contig length bp

☐ Estimated coverage X

Assembly Options

Mer size: ☒ Automatic ☐ Custom nt

Minimum match %: ☒ Automatic ☐ Custom %

☐ Allow large inserts

Post Assembly Options

☒ Minimum contig size requirements

Minimum sequences: Minimum length: bp

[Advanced Options](#)

[? < Back](#) [Cancel](#)

SeqMan NGen
File Help

Assembly Options

Set the most important options for a successful assembly

- ✓ Workflow
- ✓ Reference Sequence
- ✓ Input Sequences
- Assembly Options**
- Assembly Output

Mer size: ☒ Automatic ☐ Custom nt

Minimum match %: ☒ Automatic ☐ Custom %

☐ Maximum total reads:

☐ Vector / adapter scan [Add File...](#) [Add Folder...](#)

☐ Remove Host DNA / Contaminant scan [Add File...](#) [Add Folder...](#)

[Advanced Options](#)

[? < Back](#) [Next >](#) [Cancel](#)

Depending on the [workflow](#), only a subset of the following options will be available:

Category	Options and Descriptions
Coverage calculation	<p>If you know the approximate length of the genome or fragment being assembled, select Estimated genome/contig length and specify a length. SeqMan NGen will then calculate the expected average coverage empirically from the amount of data. This, in turn, allows repeat regions to be identified and handled more accurately, resulting in a better assembly. If the approximate genome length is not known, use the Expected coverage option.</p> <p>If you do not know the length of the genome/fragment, select Estimated coverage and provide an estimate of the depth of the sequencing. The default value for this field is 20, and the maximum allowable value is 65,535. If you enter a value larger than the maximum, you may receive an error message and be prevented from continuing until you choose a value less than or equal to the maximum. Use caution when estimating the value for Expected coverage. If the value you use is significantly lower than the actual depth, the assembly may take a much longer time to complete and may have too many mers flagged as repeats. We recommend using Expected genome length whenever possible.</p>
Mer size	<p>The minimum length of a mer (overlapping region of a fragment read), in bases, required to be considered a match when arranging reads into contigs. Mer size information is used to identify matches during the assembly layout phase. The default mer size is determined by the selected read technology and is shown in the window. For more information, see How mer tags are chosen.</p> <ul style="list-style-type: none"> • Automatic – Select this button to automatically set the size based on assembly type and sequencing technology. • Custom – Select this button to choose the size yourself. You must enter the desired number of base pairs in the field at right. Lowering the mer size increases the sensitivity of finding matches, but also increases the likelihood of finding spurious matches in addition to the correct match. Lowering the mer size can also greatly increase the requirements for storing intermediate and temporary files with large projects.
Minimum match %	<p>Specifies the minimum percentage of matches in an overlap required to join two sequences in the same contig. SeqMan NGen determines the percentage to use based on the sequencing technology you specified in the Assembly Options dialog. For more information, see Calculation of.</p> <ul style="list-style-type: none"> • Automatic – Select this button to automatically set the percentage based on assembly type and sequencing technology. • Custom – Select this button to designate the percentage yourself. You must enter a number in the field at right.

Post assembly options	<p>Check Minimum contig size requirements and type values for one or both of the following to remove assembled, untemplated contigs that do not meet minimum thresholds. This can lead to a desirable decrease in project size.</p> <ul style="list-style-type: none"> • Minimum sequences – Disassembles any untemplated contigs with fewer than the specified number of sequences. This option affects only untemplated contigs. No templated contigs are removed. • Minimum length – Using this option disassembles any untemplated contigs shorter than the specified length. This option affects only untemplated contigs. No templated contigs are removed.
Maximum total reads	<p>Check the box and enter a value if you wish to limit the read depth. Utilizing this option can make the assembly proceed faster.</p> <p>TIP: To see the effect of changes to this parameter, look at the Estimated coverage in the Run Assembly screen. In general, a coverage (AKA “depth”) of 50-100 is ideal and additional depth does nothing but slow the assembly. If your depth differs from the ideal, return to this screen and change Maximum total reads until the Estimated coverage is satisfactory. Note that this value is calculated on a per-assembly basis. If you set up 5 assemblies and set this value to 10 million reads, the cap would be 10 million reads for each assembly, not 2 million for each assembly.</p>
Vector/ adapter scan	<p>Adapter sequences are added to the ends of fragments during sequencing library preparations, and can interfere with downstream processing, if not removed. Check Vector/ adapter scan to add either one single- or multiple-sequence <i>.fasta</i> file or one folder of <i>.fasta</i> files containing known or suspected adapter sequences. The file(s) must be in <i>.fasta</i> format. During assembly, sequencing reads will be scanned for the presence of each of the specified adapters and when detected, trimmed off of that read. The trimmed read will then be used in any downstream processing. There is no specific header formatting. There is a minimum exact match length of 11 bases, and a minimum overall match of 15 bases, that allows for some mismatching. Both ends are searched within a specified range (default = 130), and all bases from an identified match to that end of the read are trimmed off.</p>
Remove Host DNA / Contaminant scan	<p>If you want SeqMan NGen to ignore host or contaminant DNA during assembly, check the box and navigate to the file or folder containing the host DNA or contaminant sequence(s). See Remove PhiX control reads from Illumina data.</p>
De novo PacBio/ Nanopore options	<p>Choose Assemble all reads or Use subset of reads. If you choose the latter, you may (or must) enter:</p> <ul style="list-style-type: none"> • Expected genome length – If you know the approximate length of the genome/ fragment being assembled, select this button and specify a length. SeqMan NGen will then calculate the expected average coverage empirically from the amount of

data. This, in turn, allows repeat regions to be identified and handled more accurately, resulting in a better assembly.

- **Desired depth of coverage of final assembly** – For optimal results, use a subset of reads to achieve 100x depth of coverage.
- **Use the longest reads in data set to achieve depth**
- **Use the first n reads in data set to achieve depth**

To set additional assembly options, press the button named **Advanced Options** or **Advanced Assembly Options** to open a multi-tabbed dialog. For details on the settings in each tab, see “Alignment tab”:#alignment-tab, “Layout tab”:#layout-tab-assembly-or-analysis-options,” Trimming tab”:#trimming-tab-assembly-options and “Scans tab”:#scans-tab.

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Assembly Output

When the Assembly Output screen appears in the wizard, you must select a name and location for your project before proceeding further.

- **Project name** – Enter a name for all output files, including the finished assembly. The finished assembly will be saved in BAM format.
- **Project folder** – Use the **Browse** button to select a location for your assembly output files.
 - For non-Cloud assemblies, **Browse** launches your file explorer. Navigate to the desired location and then click **Open** to exit. The required disk space may range from 1 GB to 5 TB, depending on a variety of factors. See our [technical requirements](#) page for more information.

IMPORTANT: Never save the assembly output files directly to the desktop, as the many intermediate files and folders created during assembly may hamper or prevent further computer operations. However, files may be saved to a folder on the desktop.

IMPORTANT: You must save to a writable location. For example, the **User** folder on Windows

is not writeable.

- For Cloud assemblies, **Browse** opens the DNASTAR Cloud Data Drive and displays your files on the DNASTAR Cloud. Navigate to the desired location and highlight the target folder, then



click the green check mark () to exit from the DNASTAR Cloud Data Drive.

- **Assembly output** displays the output file names and extensions based on current workflow and other selections. This is for information only, and cannot be edited.
- **Save Script** – Press this button if you wish to save your project and convert your wizard choices into a SeqMan NGen assembly script (*.script*) prior to assembly. This button is not available for Cloud assemblies. The resulting assembly script is an editable text file that can be modified and re-run if desired. Note that if you use **Save Script** after having checked the **Run as separate assemblies** box in the [Input Sequences](#) screen, a set of three separate scripts will be saved for the project. If you save one or more of these scripts to a location other than the main project folder, any attempt to run the assemblies from the SeqMan NGen project script will fail. Moving the projects back to the main project folder will allow assembly to proceed.
- **Write log file** – To create and save a text-formatted log file that can be used for troubleshooting any issues with the assembly. The file is saved in the same project folder as the *.assembly* and is assigned the name *.log*. The log includes the SeqMan NGen script for the project, followed by a list of steps that were performed and their outcomes. This log is especially useful in troubleshooting an assembly that will not complete.

```
Preprocessing complete.

Mapping data

Mapping data to genes in project...

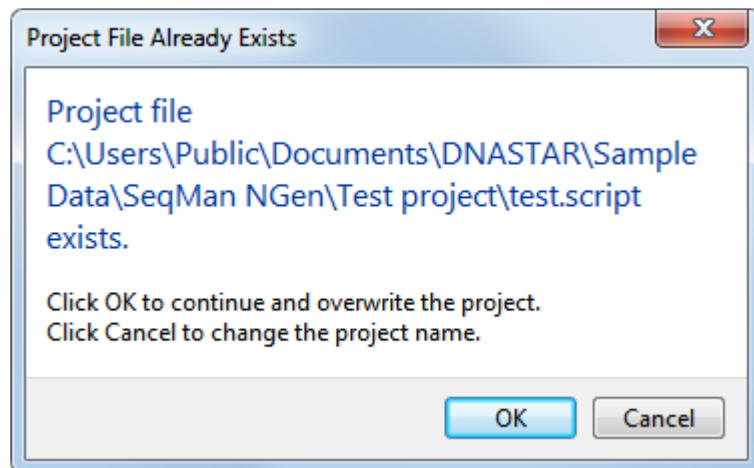
Loading data names and annotations

Loading signal data

Loading exon repeat_distrib_percent values for Ecoli templated project

Finishing
```

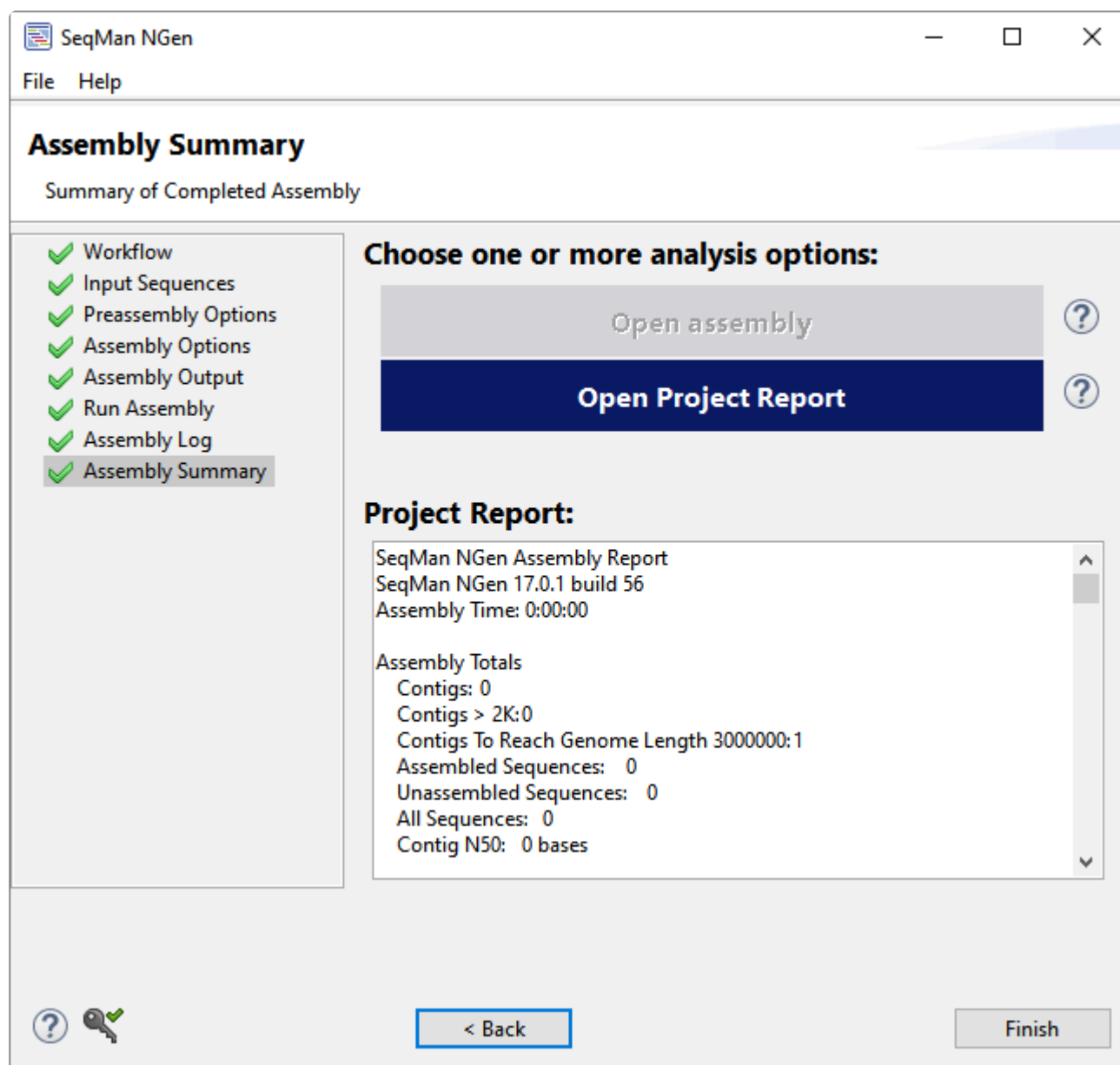
Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen. If you choose a name that already exists in the chosen location, clicking **Next >** will cause the following warning.



Click **OK** to continue and over-write the earlier project; or **Cancel** to return to the wizard screen, where you may change the project name and/or location.

Assembly Summary

After a local assembly has finished in the [Assembly Log](#), clicking **Next >** takes you to the Assembly Summary screen.



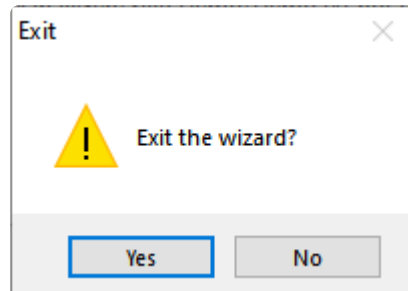
If assembly failed, the dialog displays the message "Assembly failed. No report available." Otherwise, you will see the [Project Report](#) information in the body of the screen.

Under "Choose one or more analysis options," select the button of interest. The availability of a button depends on the workflow and/or operating system and will include a subset of the following:

Button name	Result
-------------	--------

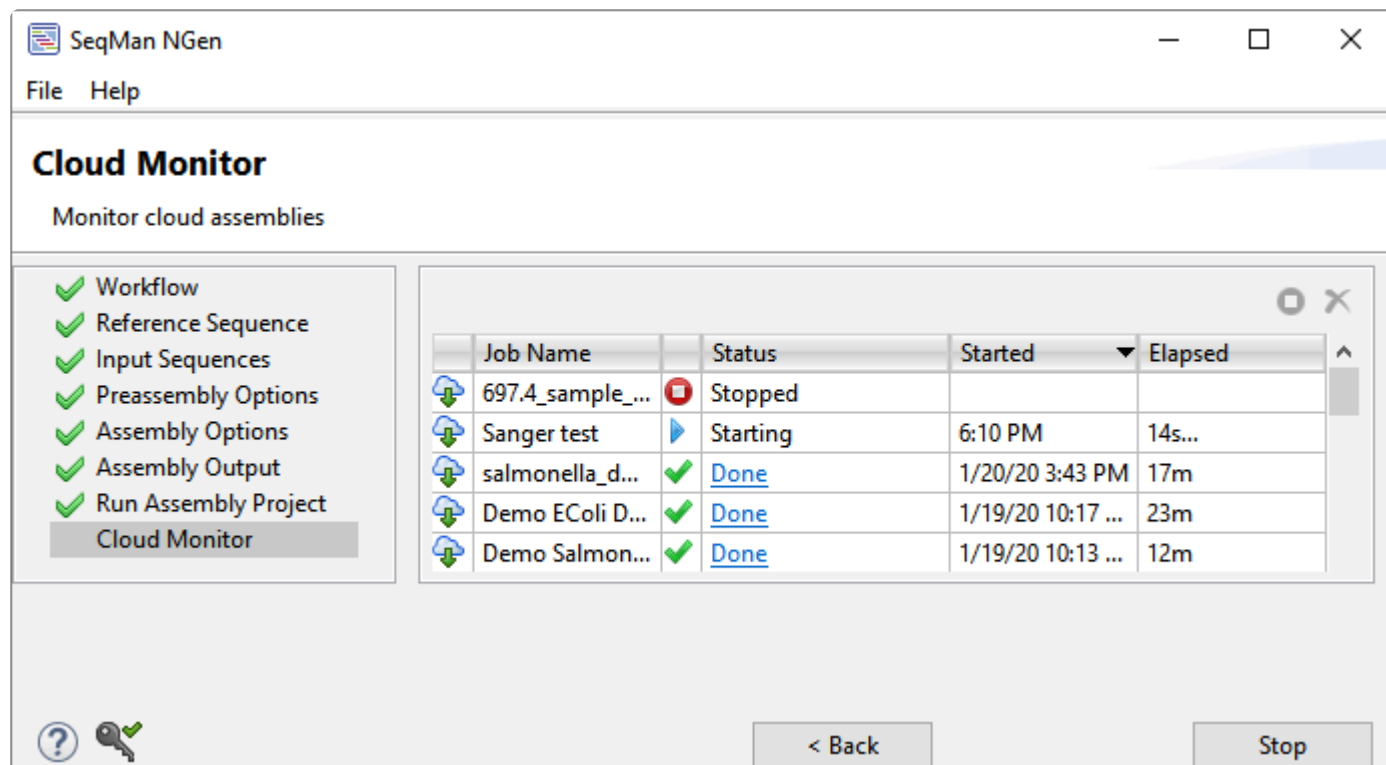
Open assembly	Launches the results in SeqMan Ultra. If multiple <i>.assembly</i> projects were created, you will be prompted to choose the names of those you wish to open. If you are following a variant workflow, we recommend downstream analysis using SeqMan Ultra's Variant view or searching for Variants in the search bar of the Analysis view.
Compare variants / differential gene expression between samples	Launches the results in ArrayStar. If you are using Macintosh, a warning message will appear asking you to open the assembly on Windows. Move the completed assembly to a Windows computer (or Macintosh running Parallels) in order to view the assembly in ArrayStar.
View peaks	This option is available for ChIP-seq only and opens the results in GenVision Pro.
Open Project Report	Opens the Project Report in your default text editor. For other ways to view the report and to learn about the report contents, see View the Project Report .
View and Compare Sashimi Plots	This option is available for RNA-seq workflows only and opens the completed assembly in GenVision Pro.

Click **< Back** to return to the previous screen or **Finish** to exit from SeqMan NGen. In most cases, you will choose **Finish**. The following confirmation popup will appear:

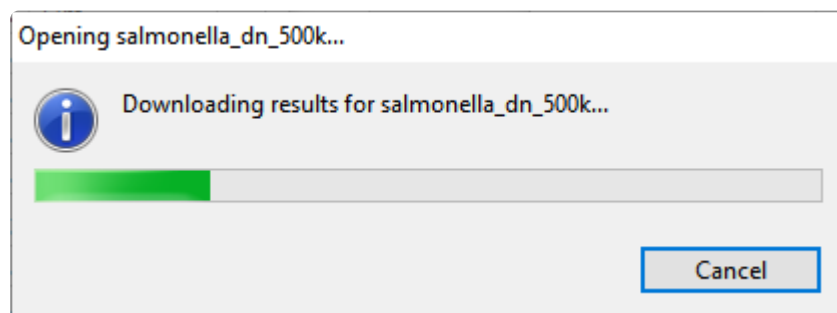


Cloud Monitor

After pressing the **Run assembly on the Cloud** link in the [Run Assembly](#) screen, you will be taken to the Cloud Monitor screen. This is where you monitor in-progress and completed Cloud Assemblies. You can also get to this screen by launching SeqMan NGen and selecting **Monitor cloud assembly** from the [Welcome](#) screen.



Once an assembly has finished successfully, the **Status** for that run is shown with the hyperlinked word "Done." Click on "Done" to download the result files. A progress bar may appear:



Once the download is complete, the results files will be displayed in your file explorer. You can then open these in a suitable application for downstream analysis, such as SeqMan Ultra, GenVision Pro or ArrayStar.

Define Binding Proteins

This wizard screen only appears for reference-guided ChIP-seq [workflows](#). The Define Binding Protein screen allows you to define binding sites for your experiment.

Start by making a selection from the **Known binding site motif** drop-down menu. The Define Binding Proteins screen will change and offer different options based on your selection.

Choose this option to use the whole genome to find peaks. Type a name into the **Binding Protein Label** box.

Type-in pattern

SeqMan NGen

File Help

Define Binding Proteins

Set the options for a post-assembly analysis

- Workflow
- Reference Sequence
- Input Sequences
- Assembly Options
- Analysis Options
- Define Binding Proteins**
- Assembly Output



Known binding site motif: Type-in Pattern

Binding Site Sequence:

Binding Protein Label:

Nucleic Acid Code	Base(s)
R	A or G
Y	C or T
M	C or A
K	T or G
W	T or A
S	C or G
B	C T or G (not A)
D	A T or G (not C)
H	A T or C (not G)
V	A C or G (not T)
N	A C T or G (any base)
X	A C T or G (any base)

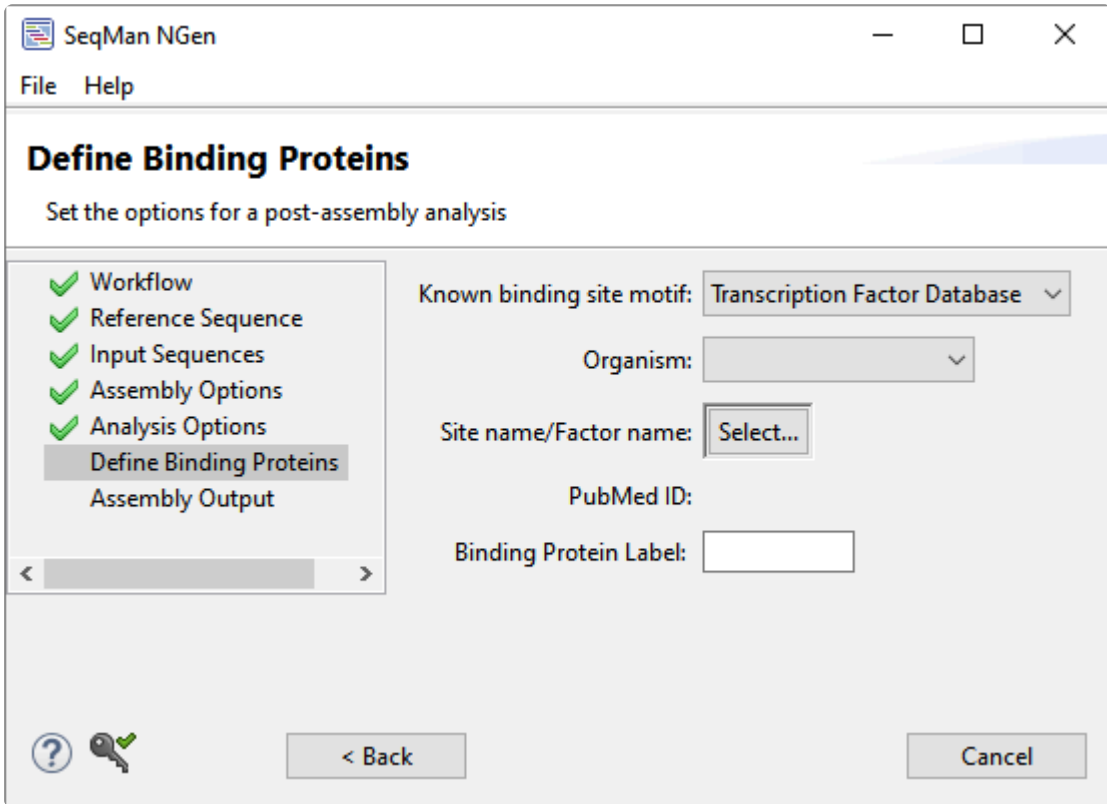
Expression Type	Syntax	Example
Inclusive	[]	[CG] means C or G is acceptable
Exclusive	{}	{AT} means neither A nor T is acceptable
Repeated	()	(AG(3)t searches for the string AGGT
Gapped	(.)	AN(1.3)G would recognize ANG ANNG ANNNG

< Back Cancel

Select this option to specify your own binding site pattern for your binding protein. Type a pattern into the **Binding Site Sequence** text box and a name into the **Binding Protein Label** text box. SeqMan NGen recognizes IUPAC nucleic acid and regular expression syntax for these patterns. A key to the syntax is provided within the wizard screen.

**Transcription
Factor
database**



Choose this option to select a binding site pattern from DNASTAR's transcription factor database, which is for prokaryotic organisms only. Select the **Organism** from the drop-down menu. Click the **Select** button to choose the site/factor name from a list. Make a selection and choose **OK**. The **PubMed ID** field will be filled in automatically. Finally, type a name into the **Binding Protein Label** text box.

JASPAR (PWM)

Select this option if you want to use the JASPAR position weight matrix to locate binding sites for eukaryotic organisms. If you choose this option, SeqMan NGen will calculate the log-odds for each sequence given the selected matrix. The score for a single character at a particular position in the matrix is equal to the \log_2 of the likelihood of seeing that character at that position in the data used to generate the matrix divided by the background likelihood of seeing that character at that position.

For example, if the matrix is derived from 80 sequences and in 70 of those sequences there is an "A" in position 1, the log odds score of seeing the character "A" in position 1 is $\log_2((70/80)/(20/80)) = 1.80$. If a "C" occurs 1 time in position 1 of the training sequences, the log odds score of seeing the character "C" in position 1 is $\log_2((1/80)/(20/80)) = -4.32$. To get the log odds score for the whole sequence, SeqMan NGen sums the log odds scores of each character in the sequence.

A sequence is considered to "match" the matrix if its score is greater than or equal to the specified Threshold. By default the threshold value is half of the average of the log-odds scores of sequences that were used to generate the pattern. You can increase the threshold for more stringency or decrease it for more matches. Once you [initiate sequence assembly](#), all detected peaks will be scanned for the presence of sites that pass the JASPAR scoring threshold.

Once you have chosen this option, select the **Organism** from the drop-down menu. Click the **Select** button to choose the site/factor name from a list. Make a selection and choose **OK**. The remaining fields will be filled in automatically. If you wish to view online entries for the selected site/factor, click on the corresponding links.

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Input Assemblies

If you are following a [Combine/Reanalyze Existing Assemblies](#) workflow, you must input at least two assemblies into the Input Assemblies screen before proceeding to the next wizard screen.

SeqMan NGen

File Help

Input Assemblies

Input assemblies and define individual replicates

✓ Workflow
Input Assemblies

Input assemblies

(Note: Two or more .assembly files are required for combine / reanalyze)

☐ Samples have replicates

Sequence File	Workflow	Experiment
---------------	----------	------------

Add Assembly...

Add Assembly from Cloud...

Add Folder...

Add Folder from Cloud...

Remove

Group Selected

Ungroup Selected

Auto Name

? Key

< Back

Cancel

- Add at least two assemblies to the list. All assemblies must have been created using the same reference sequence.
 - To add a single assembly, use the **Add Assembly** or **Add Assembly from Cloud** buttons.
 - To add all assemblies from a folder, use the **Add Folder** or **Add Folder from Cloud** buttons.

- To remove an assembly from the list, select it and click **Remove**.
- If your project involves replicate samples, check the box next to **Samples have replicates**.

If you specified that your samples have replicates:

You will see an **Experiment** column on the right of the table. Before proceeding, each blank **Experiment** cell must be assigned a name. All rows with identical **Experiment** names will be treated as a group.

- To assign experiment names manually, click on each individual cell and type in a name.
- To assign experiment names automatically based on their file names minus the file extensions, select one or more rows using click, **Ctrl+click**, **Cmd+click** or **Shift+click** and then press the **Auto Name** button.
- To group and name experiments, select one or more rows and press the **Group Selected** button. A popup will appear prompting you to type in a name for the selected file(s). Any names assigned to the row before using **Group Selected** will be overwritten.
- To remove experiments from a group, select one or more rows and press the **Ungroup Selected** button. In the selected rows, the **Experiment** column will return to its original (blank) state.

After entering **Experiment** names, click **Next** to proceed to the [Set Up Experiments](#) screen.

If you did **not** specify that your samples have replicates:

You will see an **Replicate** column on the right of the table. Before proceeding, each blank **Replicate** cell must be assigned a name. Data files that share identical names in this column will be treated as replicates.

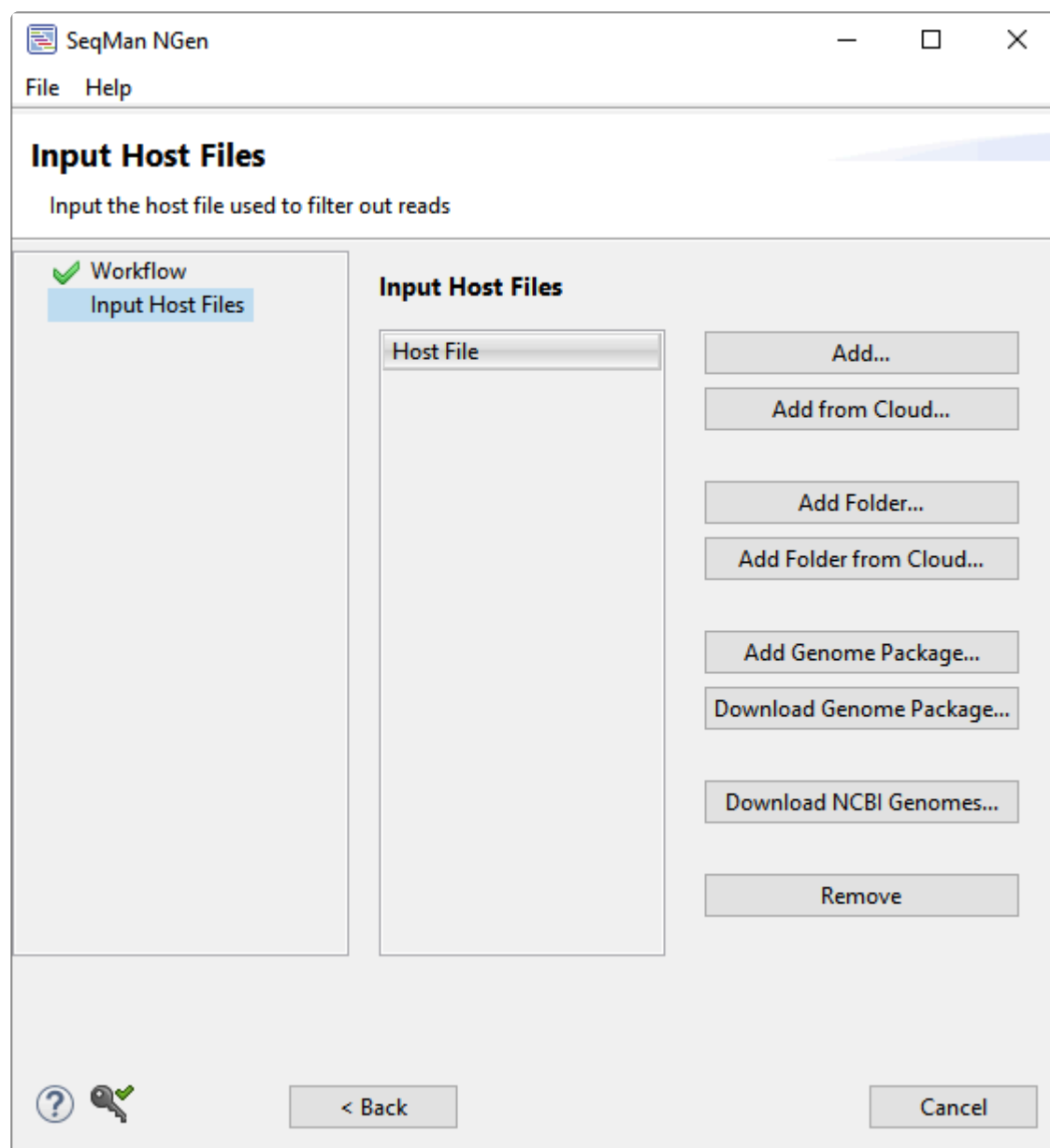
To assign replicate names manually, click on each individual cell and type in a name.

- To assign replicate names automatically based on their file names minus the file extensions, select one or more rows using click, **Ctrl+click**, **Cmd+click** or **Shift+click** and then press the **Auto Name** button.
- To group and name replicates, select one or more rows and press the **Create Replicate** button. A popup will appear prompting you to type in a name for the selected replicates. Any names assigned to the row before using **Create Replicate** will be overwritten.
- To remove items from a replicate group, select one or more rows and press the **Undo Replicate** button. In the selected rows, the **Replicate** column will return to its original (blank) state.

After entering **Replicate** names, click **Next** to proceed to the [Set Up Replicate Sets](#) screen.

Input Host Files

If you select the Viral-host integration detection workflow from the [Workflow](#) screen, Input Host Files will be the next screen to appear.



Before proceeding, you must input the host file used to filter out reads. To learn how to add files from a variety of sources, see [Add and remove files in the wizard](#).

Once you have added one or more host files, click **Next >** to proceed to the [Input Viral Genomes](#) screen or **< Back** to return to the [Welcome](#) screen.

Input Sequences

If the Input Sequences screen appears, you must specify the required options and enter one or more read files in this dialog before proceeding to the next screen. If you are following a long-read workflow, this screen will direct you to “Input long read sequences.” Depending on your workflow, only a subset of the options discussed below may be included on this screen.

The upper part of the dialog consists of a drop-down menu for specifying the read technology of your sequence files. Other options in this dialog will vary depending on the workflow and the selection made in the **Read technology** drop-down menu.

See the following topics to learn how to specify options and add data in this screen:

[Pre-import steps \(optional\):](#)

- [Check paired-end data naming specifications](#)

Specify data options:

- [Specify read technology](#)
- [Specify paired-end data](#)
- [Specify RNA-Seq options](#)
- [Specify single sample, multi-sample or replicate data](#)

Add and remove sequences:

- [Add sequences from your computer or the Cloud](#)
- [Remove a sequence from the list](#)

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Specify read technology

To specify read technology in the [Input Sequences](#) screen, make a selection from the **Read technology** drop-down menu. Default values for parameters and other assembly options in subsequent panels will be based on this selection.

Considerations when choosing a read technology:

- If you are following a long-read workflow, the only options are **PacBio** (Pacific Biosciences) or **ONT** (Oxford Nanopore Technologies).
- If you choose **Illumina**, SeqMan NGen assumes that you have paired-end data > 50 bp in length, and with a 500 bp insert distance. For all other technologies, SeqMan NGen presumes single-end data. If the read length is shorter than 50 bp, you may wish to specify a shorter **Mer size** in the [Assembly Options](#) screen. When using very short reads, you may also consider optimizing the **Minimum aligned length** and **Maximum gap size** in the [Alignment tab](#).
- For de novo workflows, if you select **Illumina** and enter an insert size of 150 bp or less in the **Set Pair Information** dialog, the assembler will assume the reads overlap and will attempt to create a single “super-read” from each pair. Read pairs that cannot be merged, either because they do not overlap or have numerous errors in the overlapping region, will not be included in the assembly. See [Remove PhiX control reads from Illumina data prior to import](#) for a description of how to use **Contaminant scan** to remove PhiX174 control sequence from Illumina data prior to assembly.
- Because the technology does not support paired reads, **PacBio** is not available when you select the [Hybrid reference-guided/de novo genome assembly workflow](#), as the gap closure step is dependent on the presence of paired reads.
- Both types of **Ion Torrent** paired reads—“mate pairs” and “paired ends”—are supported.
- If using both Sanger and Illumina data, choose **Illumina** for all reads.

Specify paired-end data

Depending on the workflow and the read technology selected, the [Input Sequences](#) screen may allow you to specify paired reads.

To specify paired reads, check the **Paired-end data** box. This causes the Pair Distance dialog to pop up. Type in the pair distance and press **OK**. The default pair distance of **500** bp is suitable for most projects.

Preparing paired-end reads:

Paired end reads are typically in two files with the forward reads in one file and the reverse reads in the other. SeqMan NGen assumes the pair will be from opposite ends of the same DNA fragment, and sequenced from the end of the fragment inwards.

To enable SeqMan NGen to identify pairs, a sequence naming convention must systematically distinguish between different pair reads while specifying which pair reads are associated. Forward and reverse sequences must have identical names except for the unique portion that determines the direction of the clone. Expressions for these naming conventions are created using a subset of *regular expressions*, which utilize elements of the Grep language. The following rules apply:

- Two parallel files must use standard naming convention (e.g. s_7_1_sequence and s_7_2_sequence).
- “Forward” and “reverse” reads must be in *exactly* the same order in the two files.
- Both forward and reverse reads must be present for every pair, including pairs where one of the reads failed or is of very low quality.

As an example, forward and reverse Sanger pair files are named as follows: 01f.abi and 01r.abi, where “01” distinguishes that they are members of the same pair. The “f” and “r” at the end of each sequence name distinguishes the orientation.

In Grep, the naming convention would be written as follows:

- Forward convention: `(.*)f\..*$`
- Reverse convention: `(.*)r\..*$`

For more information on Grep name patterns, see [Example regular expressions](#).

SeqMan NGen considers paired-end reads whose forward and reverse reads start at the same position in two reads to be clonal. In these cases, the reads with highest scores are retained, while the other reads are ignored.

Conventions for Sanger pairs:

Paired end Sanger reads are typically all in multiple files with the forward pairs having an “f” or “forward” in the name and the reverse pairs having “r” or “reverse” in the name.

Conventions for Illumina pairs:

Paired end Illumina reads are typically in two files, or a small number of files if they are from multiple runs or lanes. These pairs are specified by a naming convention used in the *.fasta* file comment line.

For *de novo* assemblies with paired end reads, SeqMan NGen automatically adds the following information to the script:

setPairSpecifier pairs:

```
{ {
  forward: "(.*)/1"
  reverse: "(.*)/2"
  min: 0
  max: 750
  key: Illumina
} }
```

If reads do not match one of the pair specifiers, or if the forward and reverse specifiers are represented by empty strings (“”), the assembler will attempt to match using the whole name of the sequence. If exactly two reads have the same name, they will be considered a match.

For reference-guided assemblies, SeqMan NGen adds the following information:

```
{
  is Pair: true
  file: "****"
  SeqTech: "Illumina"
  minDist: 0
  maxDist: 750
}
```

For reference-guided assemblies with paired-end reads, SeqMan NGen recognizes the pairs by their file names. The following examples demonstrate some of the filename formats that SeqMan NGen supports for reference-guided pairs. Large-bold text in the examples is used to highlight the region of each filename that specifies the forward and reverse reads:

“R_2011_11_21_11_06_08_user_C29-100_PE_DH10B_11_Auto_C29-100_PE_DH10B_11_4120_reverse_**pe2**.fastq

“R_2011_11_21_11_06_08_user_C29-100_PE_DH10B_11_Auto_C29-100_PE_DH10B_11_4120_forward_**pe1**.fastq

```

“Strain1234_L7_*R1*_ATCACG_Index1.fastq”,
“Strain1234_L7_*R2*_ATCACG_Index1.fastq”,

“K12-1-B_TGACCA_L006_R1.fastq”,
“K12-1-B_TGACCA_L006_R2.fastq”,

“GBBC920_GGCTAC_L008_R1.filt.50bp.fastq”,
“GBBC920_GGCTAC_L008_R2.filt.50bp.fastq”

“tiny*_1*.txt”,
“tiny*_2*.txt”,

“tiny*_1*_sequence.txt”,
“tiny*_2*_sequence.txt”,

“tiny1._qseq”,
“tiny2._qseq”,

“s_1*_1*_sequence.txt”
“s_1*_2*_sequence.txt”

“C29-129_forward_pe1.fastq”
“C29-129_forward_pe2.fastq”

The Grep used to match the pairFileNames is shown below:

“(‘name’.*)_R1_(?‘ext’.*)\\.fastq”,
“(‘name’.*)_R2_(?‘ext’.*)\\.fastq”,

“(‘name’.*)_R1\\. (?‘ext’.*)\\.fastq”,
“(‘name’.*)_R2\\. (?‘ext’.*)\\.fastq”,

“(‘name’.*)_forward_pe1(?‘ext_p’\\.fastq”,
“(‘name’.*)_reverse_pe2(?‘ext_p’\\.fastq”,

“(‘name’.*)_ {0,1}1\\.fastq”,
“(‘name’.*)_ {0,1}2\\.fastq”,

“(‘name’.*)1\\.fastq”,
“(‘name’.*)2\\.fastq”,

“(‘name’.*)_1_sequence\\.txt”,
“(‘name’.*)_2_sequence\\.txt”,

```

```
"(? 'name'.*?)1\\.txt",
```

```
"(? 'name'.*?)2\\.txt",
```

```
"(? 'name'.*?)1\\.qseq",
```

```
"(? 'name'.*?)2\\.qseq",
```

```
"(? 'name'.*?)1\\.fq",
```

```
"(? 'name'.*?)2\\.fq",
```

The following script command can be used to add support for a new filename format. The command must be executed before assembly. The pattern will be used for all subsequent **assembleTemplate** commands for that run of the reference-guided assembler.

```
pairFilePattern forward: "(? 'name'.*?)_R1_(? 'ext'.*)\\.fastq" reverse: "(? 'name'.*?)_R2_(? 'ext'.*)\\.fastq"
```

Example regular expressions

Examples of expressions you may find useful regarding paired end naming specifications follow. Please note this is not a complete list of regular expressions, and the definitions of the terms used are limited to their application to SeqMan NGen paired end naming specifications.

Special Characters	
[]	<p>Character class used to enclose a list of alternatives. For example:</p> <p>[Aa]bc matches abc and Abc.</p> <p>If the first character is a carat (^), it means anything but the characters on the list. Thus: [^a]bc matches xbc but not abc.</p>
\	A switch that makes special characters literal and literal characters special.
()	<p>Grouping--used to delimit a string comprising a "phrase." Phrases are necessary in paired end specification so you can match a pair of forward and reverse reads while still distinguishing their orientation. In SeqMan NGen, phrases in parentheses must match for two reads to qualify as a pair; phrases outside the parentheses are used to distinguish members of the same pair.</p>
\d	Any digit (0-9)
\D	Any non-digit character.
\w	Any alphanumeric "word" character (including "_")
.	Any character
	Alternate--either the term before " " or after " "
^	Match at the beginning of the line only.
\$	Match at the end of the line only.
Numerical Modifiers	
*	0 or more
+	1 or more
?	1 or 0
{n}	Exactly n
{n,}	At least n
{n,m}	At least n but not more than m
Example Expressions and Their Meanings	

d	Literally the letter d
\d	Any digit (0-9)
\d*	Zero or more digits
\d+	One or more digits
(\d+)	A phrase comprising one or more digits--same as “\d+”, but causes SeqMan NGen to match the names from the string inside the phrase when other characters in the name may not match.
\.	Literally the period symbol (.)
.	Any character
.+	One or more of any characters
.*	Zero or more of any characters
a b	a OR b
ab[i1]	abi or ab1
abi\$	Ends with abi
[\.\d]	A period OR a digit
[abc]	a OR b OR c
[abc]+	One or more characters from the set a, b, c
.*f	Any number of any characters followed by the letter “f”
(.*)f*	A phrase comprising any number of any characters, followed by the letter “f”--same as “.*f”, but causes SeqMan NGen to match the phrase in parentheses without matching the “f” in a read name
(\D+)r(\d+)	One or more non-digit characters followed by “r” followed by one or more digits.
(\d{2,4})f(\.abi)	Two, three or four digits followed by “f” followed by “.abi”

Specify single sample, multi-sample or replicate data

Depending on the workflow, the [Input Sequences](#) screen may include an **Experiment setup** drop-down menu. You must make a selection from this menu before proceeding to the next screen.

To run each sample individually:

1. Specify the **Read technology** using the drop-down menu.
2. Choose **Single sample** from the **Experiment setup** drop-down menu.
3. Check or uncheck the **Paired-end data** box as applicable to your data and enter a pair distance when prompted.
4. [Add the sequence files](#). In most cases, this will be done using the **Add** or **Add Folder** buttons on the right.
5. Press **Next**.

To run multi-sample data without replicates as one or multiple assemblies:

1. Specify the **Read technology** using the drop-down menu.
2. Choose **Multi-sample** from the **Experiment setup** drop-down menu. If you choose this option, an **Experiment** column is added to the sequence file table.
3. Check or uncheck the **Paired-end data** box as applicable to your data and enter a pair distance when prompted.
4. Specify whether you want to create a separate assembly for each sample or to run all samples as a single assembly.
 - [To create a separate assembly for each sample](#), uncheck the box next to **Run samples in single assembly**. When multi-sample assemblies are run separately, each data set is run against the reference sequence independently and an .assembly package is created for each sample.
 - [To create a single assembly](#), check the box next to **Run samples in single assembly**. This is a common workflow to use with Sanger reads (e.g., .abi and .ab1 files) and is commonly referred to as “assemble in groups”.

5. [Add the sequence files](#). In most cases, this will be done using the **Add** or **Add Folder** buttons on the right.
6. Replace each blank **Experiment** cell with a name; this must be done before you can proceed to the next screen. Data files that share identical names in this column will be assembled together.
- [To assign experiment names automatically based on a portion of the name \(recommended\)](#), use the **Group Selection Tool** at the bottom of the screen. In the box, select only the *unique* portion of the sample name, i.e., the part that *changes* from one sample to the next. This causes the **Experiment** column to change accordingly and to become simplified.

Use group selection tool to automatically group Sanger reads using the unique portions of their name

Sequence File	Pair Distance	Experiment
C-gig-2-SOSi16_M13F.ab1	1000	2
C-gig-2-SOSi16_M13R.ab1	1000	2
C-gig-6-SOSi16_M13F.ab1	1000	6
C-gig-6-SOSi16_M13R.ab1	1000	6
C-gig-8-SOSi16_M13F.ab1	1000	8
C-gig-8-SOSi16_M13R.ab1	1000	8

Group Selection Tool

Select a unique portion of the sample sequence to define the group

C-gig-**2**-SOSi16_M13F.ab1

- [To assign experiment names manually](#), click on each individual cell and type in a name. All rows with identical **Experiment** names will be treated as a group.
 - [To assign experiment names automatically based on their file names minus the file extensions](#), select one or more rows using click, **Ctrl+click**, **Cmd+click** or **Shift+click** and then press the **Auto Name** button.
 - [To group and name experiments](#), select one or more rows and press the **Group Selected** button. A popup will appear prompting you to type in a name for the selected file(s). Any names assigned to the row before using **Group Selected** will be overwritten.
 - [To remove experiments from a group](#), select one or more rows and press the **Ungroup Selected** button. In the selected rows, the **Experiment** column will return to its original (blank) state.
7. Verify that the expected number of samples and assemblies are shown in the bottom right corner of the screen.

Samples to be assembled: 6
Assemblies created: 1

8. Press **Next**. Note that a [Set Up Experiments](#) screen will be included among subsequence wizard screens.

To run multi-sample data with replicates as a single assembly:

When multi-sample data is run as a single assembly process, the data from all the samples are processed together and a single .assembly package is produced. In that case, there is a single alignment view with the data separated into the corresponding sample groups, each with a pseudo-consensus. Note that in assembling multi-sample data as a single assembly, SeqMan NGen considers all samples together. This can affect the final gapped alignment and therefore potentially yield slightly different results than assembling each sample individually.

1. Specify the **Read technology** using the drop-down menu.
2. Choose **Multi-sample with replicates** from the **Experiment setup** drop-down menu. When you choose this option, a **Replicate** column is added to the sequence file table.
3. Check or uncheck the **Paired-end data** box as applicable to your data and enter a pair distance when prompted.
4. [Add the sequence files](#). In most cases, this will be done using the **Add** or **Add Folder** buttons on the right.
5. Replace each blank **Replicate** cell with a name; this must be done before you can proceed to the next screen. Data files that share identical names in this column will be treated as replicates.
 - [To assign replicate names automatically based on a portion of the name \(recommended\)](#), use the **Replicate Selection Tool** at the bottom of the screen. In the box, select only the *unique* portion of the sample name, i.e., the part that *changes* from one sample to the next. This causes the **Replicate** column to change accordingly and to become simplified.

Input Sequences
Input sequence files and define experiments or individual replicates

Workflow
Reference Sequence
Input Sequences
Set Up Replicate Sets

Input sequences

Read technology: Sanger ☐ Paired-end data
Experiment setup: Multi-sample with replicates

Select sequence reads and create replicates
(Replicate sets will be defined on the next page)

Sequence File	Pair Distance	Replicate
C-gig-2-SOSi16_M13F.ab1	1000	2
C-gig-2-SOSi16_M13R.ab1	1000	2
C-gig-6-SOSi16_M13F.ab1	1000	6
C-gig-6-SOSi16_M13R.ab1	1000	6
C-gig-8-SOSi16_M13F.ab1	1000	8
C-gig-8-SOSi16_M13R.ab1	1000	8

Replicate Selection Tool

Select a unique portion of the sample sequence to define the replicate

C-gig-8-SOSi16_M13R.ab1

Buttons: Add..., Add from Cloud..., Add Folder..., Add Folder from Cloud..., Remove, Create Replicate, Undo Replicate, Auto Name

Samples to be assembled: 6
Assemblies created: 6

< Back Next > Cancel

- To assign replicate names manually, click on each individual cell and type in a name.
 - To assign replicate names automatically based on their file names minus the file extensions, select one or more rows using click, **Ctrl+click**, **Cmd+click** or **Shift+click** and then press the **Auto Name** button.
 - To group and name replicates, select one or more rows and press the **Create Replicate** button. A popup will appear prompting you to type in a name for the selected replicates. Any names assigned to the row before using **Create Replicate** will be overwritten.
 - To remove items from a replicate group, select one or more rows and press the **Undo Replicate** button. In the selected rows, the **Replicate** column will return to its original (blank) state.
6. Verify that the expected number of samples and assemblies are shown in the bottom right corner of the screen.
- Samples to be assembled: 6
Assemblies created: 6
7. Press **Next**. Note that a [Set Up Replicate Sets](#) screen will be included among subsequence wizard screens.

Specify RNA-Seq options

When following the reference-guided RNA-Seq workflow, the [Input Sequences](#) screen has an additional option: **Stranded RNA-Seq reads**.

Some library preparation methods preserve the directionality of reads, i.e., reverse reads always point 5' to 3' in the direction of transcription, while forward reads point from 3' to 5. If you selected the RNA-seq/ Transcriptomics > Quantitative analysis > RNA-seq [workflow](#), the Input Sequences screen provides a **Stranded RNA-seq data** checkbox. Check the box if you want SeqMan NGen to determine whether reads come from the top or bottom strand of the genome. Strandedness information is used in two ways:

- To disambiguate transcription from overlapping genes on opposite strands.
- To allow the assembler to properly allocate reads from overlapping genes.

If you check the box, you can open the finished assembly in SeqMan Pro or SeqMan Ultra and view results in the form of color-coded reads.

Input Short Read Sequences

If the Input Short Read Sequences screen appears, you must specify the required options and enter one or more read files in this dialog before proceeding to the next screen.

See the following topics to learn how to specify options and add data in this screen:

Pre-import steps (optional):

- [Check paired-end data naming specifications](#)
- [Remove PhiX control reads from Illumina data](#)

Specify data options:

- [Specify read technology](#)
- [Specify paired-end data](#)
- [Specify RNA-Seq options](#)
- [Specify single sample, multi-sample or replicate data](#)

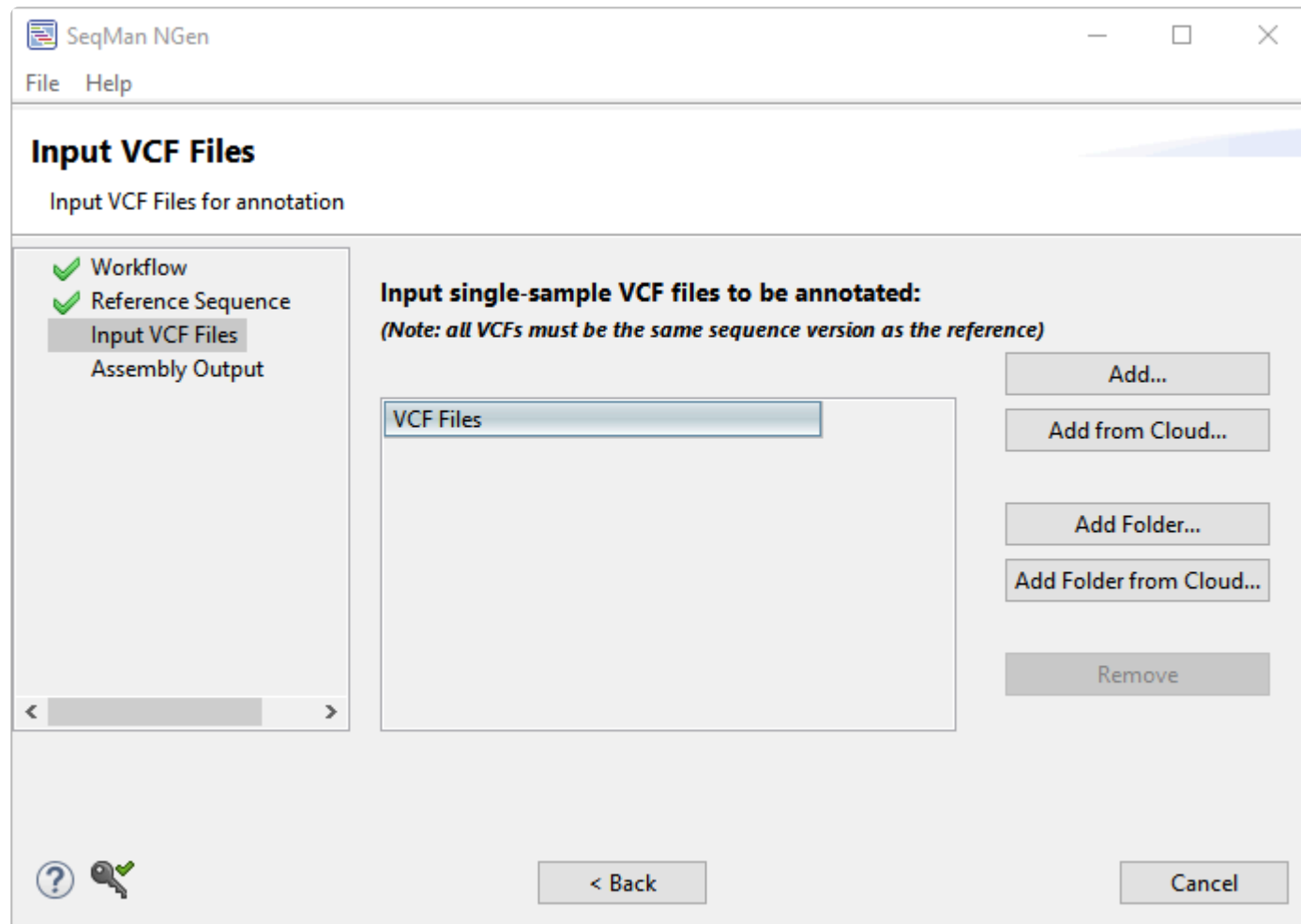
Add and remove sequences:

- [Add sequences from your computer or the Cloud](#)
- [Remove a sequence from the list](#)

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Input VCF Files

In certain workflows, the Input VCF Files screen prompts you to input the VCF files to be annotated. All VCF files must be the same sequence version as the reference sequence added in the [Reference Sequence](#) screen.



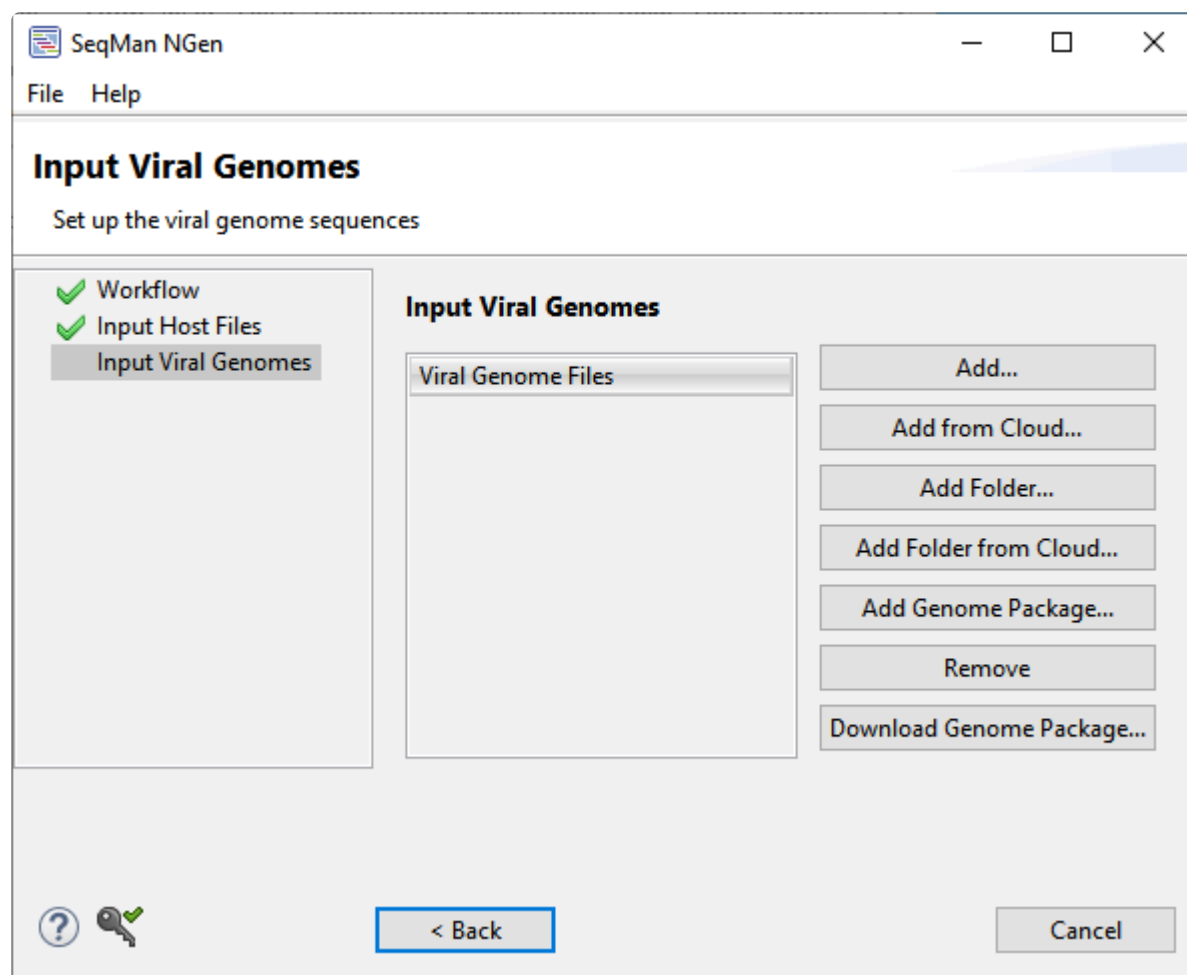
See the following topics to learn how to:

- [Add sequences from your computer or the Cloud](#)
- [Remove a sequence from the list](#)
- [Make a custom VCF file](#)

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Input Viral Genomes

If you select the Viral-host integration detection workflow from the [Workflow](#) screen, the wizard will include the Input Viral Genomes screen.



Before proceeding, you must input one or more viral genome sequences. To learn how to add sequences from a variety of sources, see [Add and remove files in the wizard](#).

Once you have added one or more viral genome sequences, click **Next >** to proceed to the [Input Sequences](#) screen or **< Back** to return to the [Input Host Files](#) screen.

Post Assembly Options

The Post Assembly Options screen lets you enter a reference sequence against which to order and orient assembled contigs. You can also set requirements for minimum contig size, if you would like to do so. Different subsets of options are available depending on the workflow.

Post Assembly Options
Input reference sequence to order and orient assembled contigs against

Workflow
Input Sequences
Assembly Options
Input Short Read Sequences
Short Read Refinement Options
Post Assembly Options
Assembly Output

Post-assembly Options:

☒ Minimum contig size requirements
Min sequences: 30 Min length: 250 bp

☐ Use an existing reference sequence to put assembled contigs in order

Reference File

Add...
Add Folder...
Add Genome Package...
Remove
Add Features...

? < Back Next > Cancel

- Check the box if you wish to specify **Minimum contig size requirements**. This option is not available for long-read workflows. If you check the box, you will need to enter the desired values for:
 - **Min(imum) sequences** – The minimum allowed number of sequences that make up a valid contig.
 - **Min(imum) length** -The size of the minimum allowed contig length for a contig to be included in the assembly.
- Check the box if you wish to **Use an existing reference sequence to put assembled contigs in order**. If you check this box, you will need to add at least one reference sequence using the various

Add buttons in the lower half of the screen. To learn how to do this, see:

- [Add a reference sequence from your computer or the Cloud](#)
- [Add a genome template from DNASTAR](#)
- [Remove a sequence from the list](#)


Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Preassembly Options

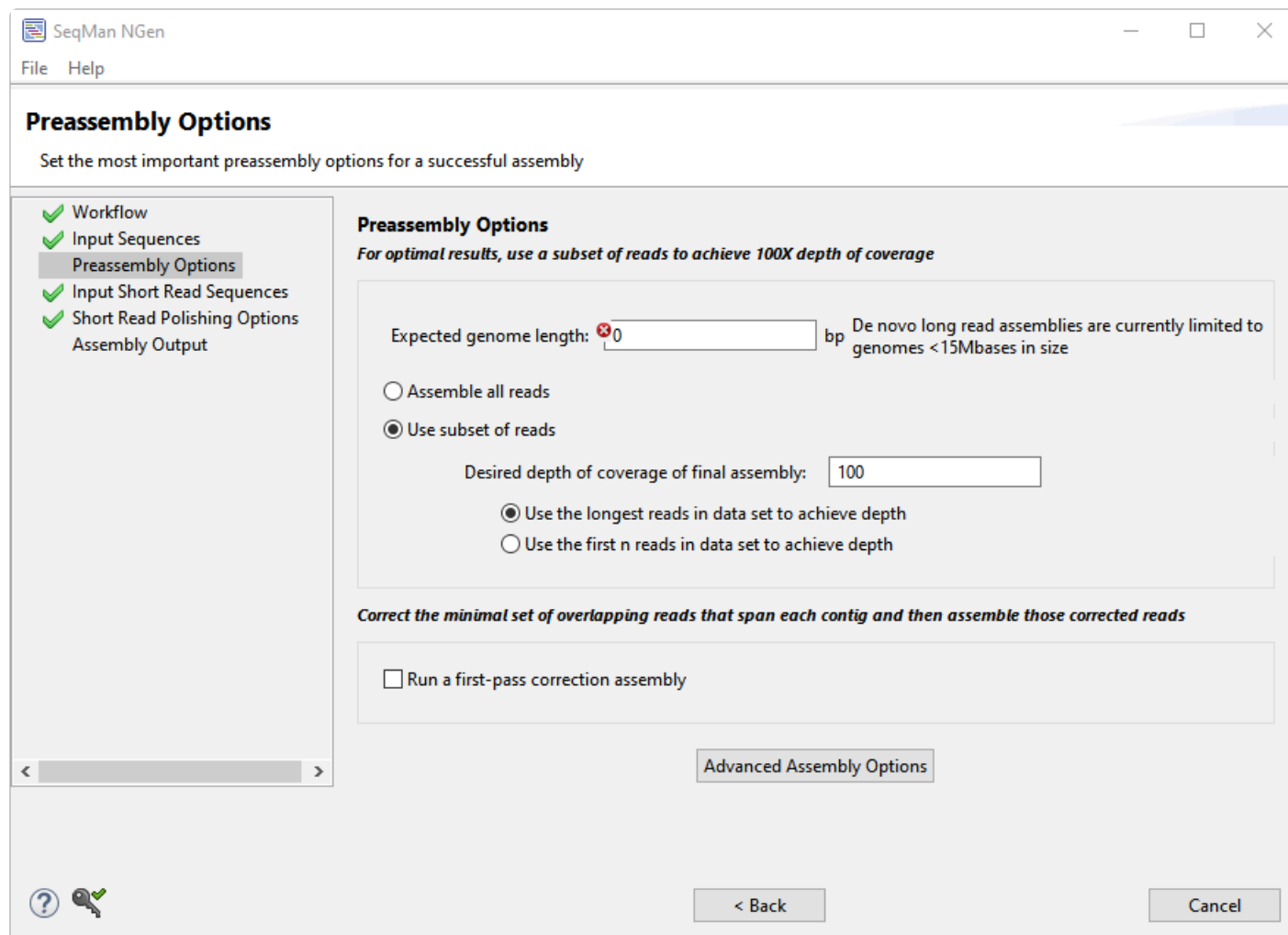
If your workflow includes the Preassembly Options screen, you can adjust the parameters used for running pre-assembly scans. There are two versions of this screen depending on the workflow:

- [Preassembly Options for long-read workflows](#)
- [Preassembly Options for all other workflows](#)

Preassembly Options for long-read workflows

 **Note:** Does your wizard screen look different from the image below? For the standard version of this wizard screen, see [Preassembly Options for all other workflows](#).

If your workflow includes the long-read version of the Preassembly Options screen, you can adjust the parameters used for running pre-assembly scans.




- **Expected genome length** – De novo long read assemblies are currently limited to genomes <15 mb in size.
- Choose between **Assemble all reads** and **Use subset of reads** (recommended). If you choose the latter, you must also enter the **Desired depth of coverage of final assembly**. We recommend keeping the default of **100** times depth of coverage for optimal results. You will also need to specify whether you want SeqMan NGen to **Use the longest reads in data set to achieve depth** or **Use the first n reads in data set to achieve depth**.

- If you wish to correct the longest, higher-quality reads and assembly only those corrected reads, check the box next to **Run a first-pass correction assembly**.

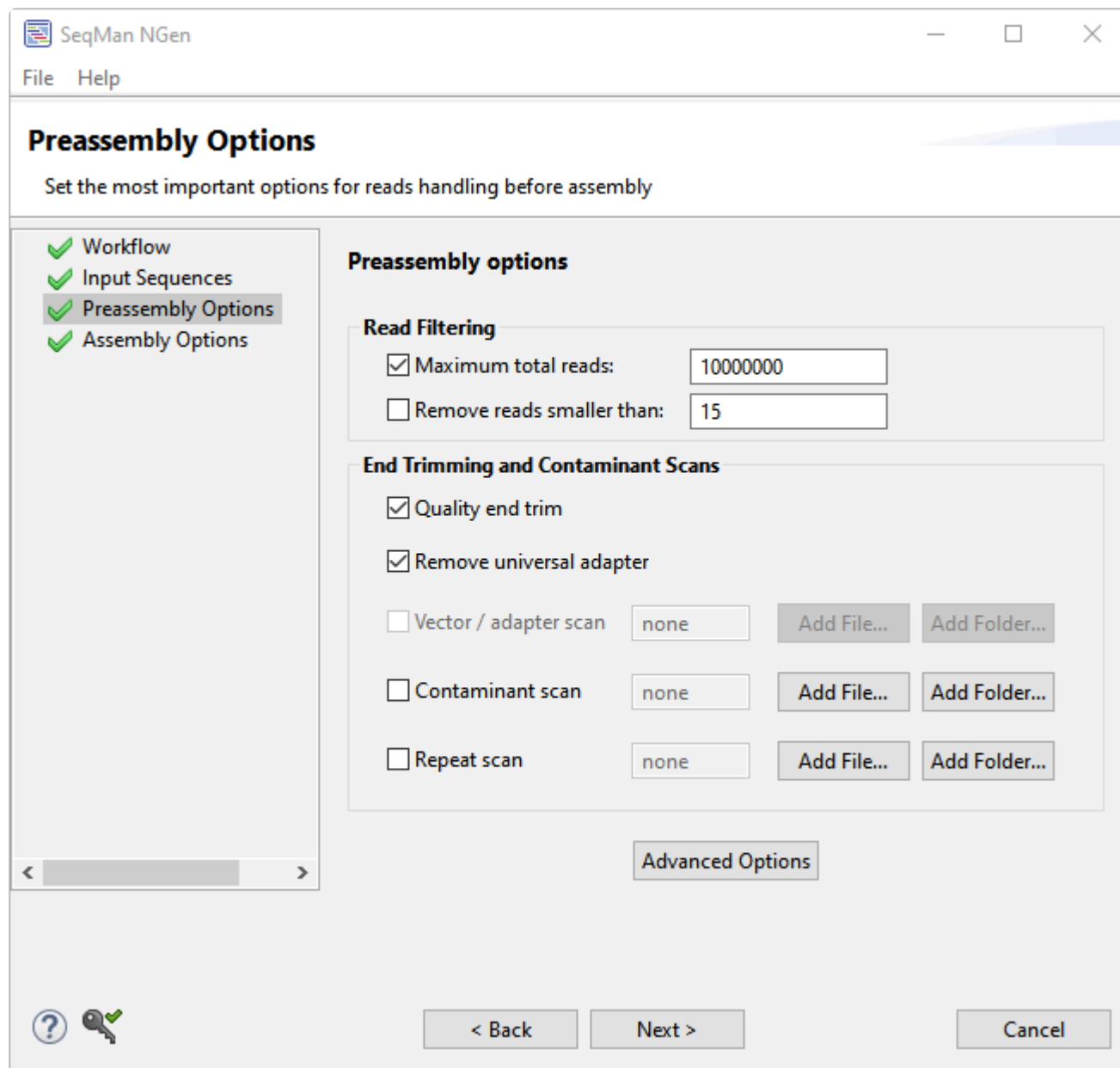
If desired, press the **Advanced Assembly Options** button to access additional options. A two-tabbed dialog opens. See [Alignment tab](#) and [Layout tab](#) for descriptions.

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Preassembly Options for all other workflows

 **Note:** Does your wizard screen look different from the image below? For the standard version of this wizard screen, see [Preassembly Options for long-read workflows](#).

If your workflow includes the standard Preassembly Options screen, you can adjust the parameters used for running pre-assembly scans.



SeqMan NGen

File Help

Preassembly Options

Set the most important options for reads handling before assembly

- ✓ Workflow
- ✓ Input Sequences
- ✓ Preassembly Options
- ✓ Assembly Options

Preassembly options

Read Filtering

☒ Maximum total reads: 10000000

☐ Remove reads smaller than: 15

End Trimming and Contaminant Scans

☒ Quality end trim

☒ Remove universal adapter

☐ Vector / adapter scan none Add File... Add Folder...

☐ Contaminant scan none Add File... Add Folder...

☐ Repeat scan none Add File... Add Folder...

Advanced Options

< Back Next > Cancel

In the Read Filtering section.

- To set a threshold for **Maximum total reads**, check the box and enter the threshold value. The default is for this option to be checked and the threshold to be 10000000. When using Illumina

technology, we recommend checking this box and specifying a value to limit the number of reads used in the assembly. For 454 technologies, we suggest unchecking the box and leaving the field blank. If you check this option, be sure to add individual read files rather than folders in the [Input Sequences](#) screen. Adding files individually causes SeqMan NGen to use an equal number of reads from each file. If you instead add a folder, SeqMan NGen may potentially use reads from only the first file(s).

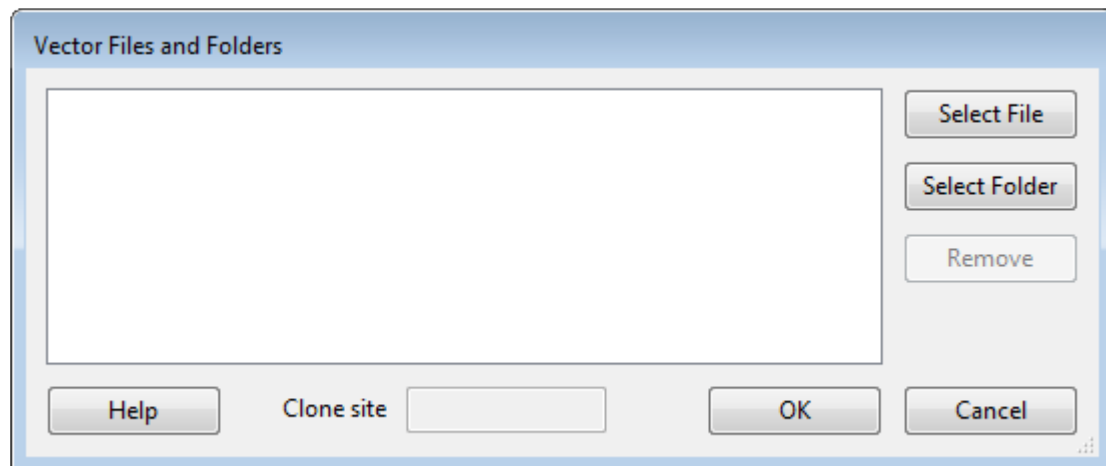
TIP: To see the effect of changes to this parameter, look at the **Estimated coverage** in the [Run Assembly](#) screen. In general, a coverage (AKA “depth”) of 50-100 is ideal and additional depth does nothing but slow the assembly. If your depth differs from the ideal, return to this screen and change **Maximum total reads** until the **Estimated coverage** is satisfactory.

- To **Remove reads smaller than** a specified threshold, check the box and enter the threshold value. The default is for this option to be unchecked.

In the End Trimming and Contaminant Scans section, specify whether you would like SeqMan NGen to perform any of the following pre-assembly tasks:

- To automatically trim reads prior to assembly based on quality scores and specified quality end trimming parameters, check **Quality end trim**.
- To use specified vector/adaptor scan parameters to scan and trim reads for the vector or adaptor, check **Vector/adaptor scan**.
- To use specified contaminant scan parameters to scan and remove reads that contain contaminant sequences, check **Contaminant scan**. If desired, learn how to [use Contaminant scan to remove PhiX174 control sequence from Illumina data prior to assembly](#).
- To use specified repeat scan parameters to scan reads for known repetitive sequences, check **Repeat scan**. If the option is checked, all sequences identified as repeats will be added to the assembly last, after all non-repeats have been assembled.

If you check any of the last three options, click the **Add** button to its right to select the desired vector, contaminant, or repetitive sequence(s). This causes a Files and Folders dialog to pop up.



- **Select File** – Click to navigate to and select one or more individual sequences.
- **Select Folder** – Click to navigate to and select an entire folder of sequences.
- **Clone site** (Vector dialog only) – Enter the position of the cloning site where insertion occurs.
- **Remove** – Click to remove a selected (highlighted) file from the list.

Once you are finished, click **OK** to save your changes and return to the Preassembly Options dialog.

If desired, click the **Advanced Options** button to access additional settings. This opens a multi-tabbed dialog. For details, see [Trimming tab](#), [Scans tab](#) and [Alignment tab](#).

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Reference Sequence / Input Draft Genome

If the Reference Sequence screen appears, you must input some type of reference sequence or genome template package before proceeding further in the wizard.

If you are following one of the [Genome finishing workflows](#), this screen is called **Input Draft Genome or Contigs** or **Input High Quality Draft Genome**. These screens each contain a subset of the options in the **Reference Sequence** screen.

SeqMan NGen

File Help

Reference Sequence

Input the reference sequence for the assembly.

✓ Workflow
Reference Sequence

Input reference sequence

Reference File

Add...

Add from Cloud...

Add Folder...

Add Folder from Cloud...

Add Genome Package...

Download Genome Package...

Download NCBI Genomes...

Remove

Auxiliary Files:

☐ VCF file: Browse...

☐ BED file: Browse...

? < Back Cancel

Optional pre-import steps for the Reference Sequence screen:

- [Annotate reference sequences prior to import](#)
- [Manually specify an isoform prior to import](#)
- [Make a custom VCF file](#)
- [Make a custom BED file](#)
- [Troubleshoot a Manifest file](#)

Add and remove reference sequence files or draft genomes:

SeqMan NGen can read and produce output using a variety of common chromosome naming conventions, including “chr1” and “ch1,” as well as Arabic and Roman numerals. See our [Supported File Types](#) page for allowable file extensions.

- [Add a reference sequence or draft genome from your computer or the Cloud](#)
- [Add a genome template from DNASTAR](#) – (RECOMMENDED) Genome template packages include dbSNP information that is very useful for downstream analysis.
- [Add a genome template from NCBI](#)
- [Use RNA-Seq de novo transcriptome output as a reference](#)
- [Remove a sequence from the list](#)

Options pertaining only to the Reference Sequence screen:


- [Specify a VCF, BED or Manifest file](#)

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Annotate reference sequences prior to import

Using annotated reference sequences in SeqMan NGen may enable you to better analyze the identified putative SNPs when viewing your assembled project in SeqMan Pro or SeqMan Ultra. If desired, annotate your reference sequence in [SeqBuilder Pro](#) (the Lasergene application for sequence editing and visualization) prior to adding it to the [Reference Sequence](#) dialog.

1. Launch SeqBuilder Pro.
2. Go to **File > Open** and select the reference sequence.
3. Select the range of sequence where a feature will be added. (Use **Edit > Go to Position** to navigate quickly up and down your sequence.)
4. Go to **Features > New Feature**. A new “misc_feature” will be added to your sequence and displayed in the Feature List.
5. Click on “misc_feature” from within the Feature List and select the appropriate feature type from the list provided*. * For example:
 - For SNPs, choose **Variation > variation**.
 - For exons, choose **Gene > exon**.
 - For CDS features, choose **Transcript > CDS**.
 - For origin of replication, choose **Structure > rep_origin**.

 **Note:** The next feature you create will automatically be of the same feature type you just selected, enabling you to create all the features of one type more quickly.

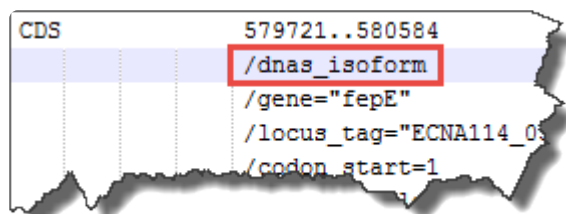
6. Repeat steps 3-5 until all of your features have been added. Then go to **File > Save As** and save your sequence in *.sbd*, *.seq* or *.gbk* format. Your annotated reference sequence is now ready for assembly in SeqMan NGen and subsequent analysis in SeqMan Pro.

Manually specify an isoform prior to import

By default, SeqMan NGen chooses the longest CDS as the isoform for SNP calling. If desired, you may override the automated choice by specifying the preferred isoform manually in the reference sequence.

To do so, follow these steps prior to importing the reference sequence into SeqMan NGen's [Reference Sequence](#) screen:

1. Open the reference sequence in a text editor.
2. Locate the feature of interest. Just below its location coordinates, type in **/dnas_isoform**.

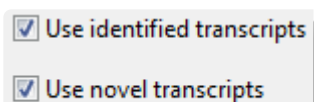


3. Save the edited template sequence.
4. Input the saved version of the template sequence into SeqMan NGen's [Reference Sequence](#) screen.


Use RNA-Seq de novo transcriptome output as a reference

You may use the contigs output from an RNA-seq *de novo* workflow as reference sequences in the templated RNA-Seq workflow. Doing this may allow you to quantify the relative abundances of transcripts using ArrayStar. Note that this use case assumes the same reads are used for both rounds of assembly.

1. Follow the RNA-Seq de novo transcriptome workflow. After assembly is complete, close the SeqMan NGen wizard.
2. Launch SeqMan NGen again and follow the RNA-seq templated workflow.
3. In the [Reference Sequence](#) screen do one of the following:
 - Press **Add Folder**, navigate to the Transcriptome package that was output in Step 1, and press **OK**. Then choose which types of transcripts to input by checking one or both boxes at the bottom of the screen.



- If you exported a subset of transcripts from SeqMan Ultra and wish to use those, rather than the full Transcriptome package, press **Add** and add a single *.fas* file.

 **Note:** SeqMan NGen does not support the addition of more than one *.fas* file from the *.Transcriptome* folder.

4. In the [Input Sequences](#), load the desired reads. These can be the same reads that were used as input in the original transcript annotation workflow.
5. Set other options, as desired, and run the assembly.


Specify a VCF, BED or Manifest file

Certain workflows allow or require you to import a BED file, Manifest file, or a custom VCF SNP file with data from one or more assemblies. These options, if available, will be offered at the bottom of the [Reference Sequence](#) screen.

To add a BED, Manifest, or targeted regions file:

These files can be used in many types of assemblies, and should always be uploaded when doing exome assemblies.

To add a file of these types, check the box next to **BED file** and then use the corresponding **Browse** button to navigate to the file. If you are doing a local assembly, select the file and click **Open**. If you are doing a


Cloud assembly, select the file and click the green check mark (). Note that:

- You must upload the *.bed* file from the capture kit, not any random *.bed* file. For instance, if you used Human Genome build 38 as the reference, for example, the corresponding *.bed* file might be called *Human genome build38.bed*. Sequencing services do not always send the *.bed* file to customers, but can often provide it on request.
- BED files must have the extension *.bed*. For information on making a BED file, see [Make a custom BED file](#).
- Manifest files are typically used to represent coordinates of regions that were captured in procedures, such as exon capture performed prior to sequencing. They can have various extensions (usually *.txt*), but must be in the [correct format](#).

To add a VCF file:

Most users will not have a VCF file to upload, but are looking forward to producing a VCF file as part of the assembly output. However, researchers working with standard reference data sets may have a VCF file to upload to the wizard at this stage.

To upload a VCF file, check the box next to **VCF file** and then use the corresponding **Browse** button to navigate to the file. If you are doing a local assembly, select the file and click **Open**. If you are doing a

Cloud assembly, select the file and click the green check mark (). Positions within the VCF file will be given a VCF SNP ID during the assembly process. After assembly, information about each position can be viewed in the SeqMan Pro SNP Report (**SNP > SNP Report**) or the SeqMan Ultra Variants view. SeqMan NGen only supports one VCF file per assembly project. If you have multiple VCF files (e.g., one per chromosome), you must merge the information into a single VCF file before browsing to the file. For more information, see [Make a custom VCF file](#).

Make a custom VCF file

Variant Call Format (VCF) files have multiple uses. For instance, they can provide a way to flag previously known SNPs and to filter them in SNP tables. In DNASTAR's SeqMan NGen, these SNPs are called "annotated SNPs"; in ArrayStar, they are referred to as "user variants." VCF files can also be used to keep track of previously identified variants so that they can be verified in a new assembly or experiment.

The following brief video is an overview of how VCF files can be used in an assembly and in downstream analysis:

VCF files can be custom-made or automatically generated by sequencing software. For instance, you can create a VCF file using software such as SeqMan Pro, SeqMan Ultra or ArrayStar. Certain SeqMan NGen assemblies also output a VCF file called [assembly_name].sample.vcf. VCF files are also available from other sources, such as the [UCSC Genome Browser](#), and the [Genome In a Bottle](#) project. For a description of various VCF version specifications, see the Sourceforge [VCF Specification page](#).

These two columns are REQUIRED , and must be in the order shown. All cells in these columns must be filled.		These four columns are OPTIONAL . If optional columns are present, the assembler will check the length of the string and compare against the length of the called variant. The base identities will not be checked.				Columns 7 and beyond are allowed, but will be ignored.
#CHROM	POS	ID	REF	ALT	INFO	(Misc.)
Chromosome identifier. Numbers are preferred, but chr or ch prefixes	Position in the reference sequence.	For known dbSNP entries, the rs ID. The valid format is rs followed by a series of	The reference base(s). For	The variant base(s). For	User ID and source assembly information.	These columns may contain data, but they will be

are allowed. All cells in this column must be filled.	All cells in this column must be filled.	digits. For unknown or nonexistent IDs, a period (.)	unknown bases, a period (.)	unknown bases, a period (.)	For unknown bases, a period (.)	ignored by the SeqMan NGen assembler.
---	--	--	-----------------------------	-----------------------------	---------------------------------	---------------------------------------

- The table portion of the file must be sorted numerically, first by #CHROM, and then by POS. Make sure to sort the columns numerically (1, 2, 3...) and not alphabetically (1, 11, 12...). If you attempt to run the assembly after loading an improperly-sorted VCF file, multiple red error messages will be displayed during the assembly.
- When you try to open extremely large VCF files in a spreadsheet program or text editor for sorting purposes, you may receive an “insufficient memory” warning. If you need to sort a VCF file that is too big to open on your machine, we recommend using Sourceforge’s [VCFTools](#).
- If quotation marks are used anywhere in the VCF file, they must be straight quotes, not curly or “smart” quotes. In addition, quotation marks should not be used in lines beginning with ##contig, ##UnifiedGenotyper, or ##INFO. If these rules are not followed, an error message will appear during assembly stating that “the VCF file has an incorrect or missing header.” Though the assembly will continue, the VCF SNP file that is output will be empty.
- Chromosome names are captured from genome template packages and used to assign contig IDs to entries from BED, VCF and manifest files. SeqMan NGen can read and produce output using common naming conventions (i.e., “chr” and “ch”) and Arabic numerals. It understands that chr1, ch1, or 1 can all be used to represent “the first template in the index,” and so on. In addition, Genome Template Packages sometimes internally define “short names” for particular chromosomes. For example, the *C. elegans* template package names its chromosomes using the standard convention for that organism: “I”, “II”, “III”, “IV”, “V”, “X”, “M.” SeqMan NGen does not normally recognize Roman numerals, but can in this case, because the numbers are “short names” that have been mapped to specific chromosomes.

Make a custom BED file

Certain workflows allow or require you to import a targeted regions file, also known as a BED file, within the [Reference Sequence](#) screen.

BED files are used to define capture regions in the assembly, and can be generated by the sequence provider or made by hand. These files are basically tab-separated text files whose extension has been changed to *.bed*. See the UCSC Genome Bioinformatics [BED file](#) page for detailed information.

The following brief video shows how BED files can be incorporated into an assembly for later downstream analysis:

SeqMan NGen can read and produce output using a variety of common chromosome naming conventions, including “chr1” and “ch1,” as well as Arabic and Roman numerals. Chromosome names are captured from genome template packages and used to assign contig IDs to entries from BED, VCF and manifest files.

The BED file can consist of multiple sections, each with a different track name. Text is allowed between the tables without restriction.

These three columns are REQUIRED , and must be in the order shown. All cells in these columns must be filled.			Columns 4 and beyond are allowed, but will be ignored.
chrom	chromStart	chromEnd	(Misc.)
The name of the	Starting position for the feature. The coordinates for	Ending	Data in these

chromosome or scaffold. Numbers are preferred, but chr or ch prefixes are allowed.	BED intervals are in 0-based coordinates as follows: (1-0) .. (100+1-1). Therefore, base 1-100 of a chromosome is represented in a BED file as 0-100.	position for the feature.	columns are ignored.
--	---	---------------------------	----------------------

- A header row is optional and can contain any text; text need not match that shown in the table header row above.
- **IMPORTANT:** Each table in the file must be primarily sorted by the first column, and secondarily sorted by the second column. The columns must be sorted numerically (1, 2, 3...) and not alphabetically (1, 11, 12...). If only chromosome 1 (and possibly 11) appears in SeqMan Pro's "Coverage of Targeted Regions" report (**Project > Show coverage of target regions**), this is indicative of incorrect sorting.

Troubleshoot a Manifest file

In certain workflows, you may be able to specify a Manifest file in the [Reference Sequence](#) screen.

Manifest files are tab-separated files used to define the chromosomal coordinates of gene targets in the assembly, and are normally generated automatically, e.g., by Illumina. (See Illumina's [manifest file PDF](#) for a description.) Manifest files can have various file extensions, though *.txt* is commonly used.

SeqMan NGen can read and produce output using a variety of common chromosome naming conventions, including “chr1” and “ch1,” as well as Arabic and Roman numerals. Chromosome names are captured from genome template packages and used to assign contig IDs to entries from BED, VCF and manifest files.

The following examples show the most basic required columns for the manifest file, as well as two formats that are used by Illumina. These examples are provided in case you need to trouble-shoot problems with existing manifest files. If you want to create your own targeted regions file, we recommend making a BED file rather than a manifest file. See [Make a custom BED file](#) for detailed instructions.

The columns below can appear in any order. Columns with **blue** headers are required, while columns with **orange** headers are optional.

User-made file (most basic version):

[Regions]			
Name	Chromosome	Start	End
28324371	chrM	577	647

Illumina manifest file – format 1:

[Targets]						
TargetA	TargetB	Target Number	Chromosome	Start Position	End Position	Probe Strand
chr1.43815008.43815009	chr1.43815008.43815009	1	chr1	43814982	43815163	+
chr1.115256528.115256531	chr1.115256528.115256531	1	chr1	115256500	115256680	+

Illumina manifest file – format 2:

[Regions]					
Name	Chromosome	Start	End	Upstream Probe Length	Downstream Probe Length
28324371	chrM	577	647	0	0

Run Assembly

The Run Assembly screen follows the [Assembly Output](#) screen and prompts you to specify a name for the project and a location in which to save temporary files. You can also review system memory information that will help you make an educated decision whether to assemble on your local computer or on the cloud.

SeqMan NGen

File Help

Run Assembly

Check assembly and computer settings and start the assembly

- ✓ Workflow
- ✓ Reference Sequence
- ✓ Input Sequences
- ✓ Preassembly Options
- ✓ Assembly Options
- ✓ Assembly Output
- ✓ **Run Assembly**

Input information

Total length reference sequences: 538 Bases

Estimated input read bases: 4.5 KBases

Estimated coverage: 8X

Temporary files


Location: C:\Users\Public\Documents\DNASTAR\Sample Data\SeqMan Browse...

Estimated requirements


	Required:	This computer and temporary files drive:
Memory:	4.0 GB	16 GB
Free disk space:	<0.1 GB	351.5 GB

Run assembly


RECOMMENDED



[Run assembly on this computer](#)



[Run assembly on the cloud](#)

?  < Back Cancel

- The “Input information” section shows a quick overview regarding the size of the planned assembly, including the total length of the reference sequences, estimated size of the read files, and estimate coverage area.
- Under “Temporary files,” the **Location** box shows any location previously specified for temporary files

created during the run. If you have not previously specified a location, the box will be blank. To specify or change the location, use the **Browse** button or open your computer's file explorer and drag and drop the directory in the **Temporary file location** box. SeqMan NGen will remember and use the temporary file location for future assemblies. Here are some tips regarding temporary files:

- If possible, use an external hard drive as the temporary file location.
 - Never save the assembly output files or temporary files directly to the desktop, as the many intermediate files and folders created during assembly may hamper or prevent further computer operations. However, files may be saved to a folder on the desktop.
 - By default, most temporary files are deleted when the assembly is complete. Other files (e.g., [template_name].FasInfo.sqlite and [template_name].mer) may remain in the temporary file location in order to facilitate efficient reassembly of data in the future.
 - You do not need to specify a temporary file location when following a *de novo* or special reference-guided workflow.
- The “Estimated requirements” section consists of two columns. The **Estimate** column shows the estimated **Memory** and **Free disk space** required for the current assembly to run without failure. The **This computer and temporary files** column shows the amount of memory and free disk space on your local computer. If either value on the left is larger than its equivalent on the right, we strongly recommend you run the assembly on the cloud. Unless you have an extremely powerful computer, most non-bacterial assemblies should be run on the cloud.
 - In the “Run assembly*” section, click the link to **Run assembly on this computer** or **Run assembly on the cloud**. Base your decision on the “Estimated requirements” section above. If you choose to run on the cloud and are not already logged in, you will be prompted to [log in to Cloud Assemblies](#).
 - If you choose **Run assembly on this computer**, you must wait for one assembly to finish before beginning another. Note that *de novo* projects require ample amounts of RAM, while most reference-guided assemblies require large amounts of CPU and free hard drive space. Unless you have an unusually powerful computer, local assembly is best reserved for small, reference-guided assembly projects.
 - If you choose **Run assembly on the cloud**, you can perform any number of SeqMan NGen assemblies simultaneously, without using up your computer's resources. The assembly takes place on a powerful Amazon cloud computer, and output data are stored securely on the cloud. To try Cloud Assembly for free*, [contact your DNASTAR representative](#) to request free trial access and a secure data storage folder.

If you want to make changes to a screen before proceeding, click the screen name on the left of the wizard or choose **< Back** to return to the previous screen.

Monitor the progress of a Cloud Assembly

To monitor the progress of a Cloud Assembly:

Once you press the **Start Assembly** button in the SeqMan NGen wizard, the assembly begins. The Cloud Assembly tutorials in this User Guide each take 30 minutes to several hours to complete. You can monitor the progress using either of two methods:

- From the wizard: – After pressing **Start Assembly**, wait for the **Next** button to become active. Press **Next** to go to the Cloud Assemblies page, where you can monitor the progress of the assembly.

Cloud Assemblies			
Job	Status	Start Time	End Time
Arabidopsis strains-qng	Queued	03/11/19 15:22:22	
Mutant 3.7z	Starting	03/11/19 15:22:23	
Mutant 2.7z	Starting	03/11/19 15:22:23	
Mutant 1.7z	Starting	03/11/19 15:21:58	

Close the wizard at any time by pressing **Finish**. After pressing **Finish**, you will need to use the other monitoring method (next bullet point) if you wish to continue monitoring the assembly.

- From the DNASTAR website: – [Open the Web Monitor](#), entering your DNASTAR login credentials if prompted to do so. See the “Web Monitor help”https://www.dnastar.com/Web_Monitor_Help/#!/Documents/usethewebmonitor.htm to learn about features in this browser window.

DNASTAR Cloud Assemblies

Refresh

Legend

Web Monitoring

- NGen Cloud service - Number of assemblies remaining on DNASTAR Cloud Assemblies license: 16
- 9 total assemblies (4 done, 0 failed, 1 stopped, 1 queued, 3 in-process) - Updated 3/11/2019 3:36:30 PM
- For help, please use [this link](#).

Show latest assemblies ▼

Limit to 10 assemblies

Name	Status	Started	Ended	Message	Files
Arabidopsis strains-qng	Queued				
Mutant 3.7z	Starting	3/11/2019 3:22:23 PM		Starting	
Mutant 2.7z	Starting	3/11/2019 3:22:23 PM		Starting	
Mutant 1.7z	Starting	3/11/2019 3:21:58 PM		Starting	

Once an assembly is **Done**, you can download the assembly results for analysis with SeqMan Pro, ArrayStar, GenVision Pro, etc.

Set Contaminant

If you are following certain RNA-Seq workflows, the wizard includes the Set Contaminant screen. The screen prompts you to remove rRNA sequences if they are present. We highly encourage doing so, as these sequences can represent up to 80% of the reads in a de novo transcriptome data set and can swamp the assembly due to their large numbers.

Check **Scan for rRNA contamination** to enable the three add/remove sequence buttons on the right. You will then need to specify rRNA sequences here, as described below. Typically, you will add reference sequence(s) from a 16S rRNA database (e.g., [Silva](#), [Greengenes Ribosomal Database Project](#), etc.

See the following topics to learn how to:

- [Add rRNA sequence files or folders of files from your computer or the Cloud](#)
- [Remove a biome genome from the list](#)

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Set Up Experiments

If, in the [Input Sequences](#) screen, you choose the **Experiment setup** option **Multi-sample**, Set Up Experiments will appear as the next wizard screen.

SeqMan NGen
File Help

Set Up Experiments

Set the most important options for a successful assembly

- ✓ Workflow
- ✓ Reference Sequence
- ✓ Input Sequences
- ✓ **Set Up Experiments**
- ✓ Assembly Options

Set control experiment:

Experiment	Is Control?
wildtype	<input checked="" type="checkbox"/>
red eyes	<input type="checkbox"/>

< Back Next > Cancel

The **Experiment** column is pre-loaded with the “Experiment” names specified in the [Input Sequences](#) screen. If you wish to edit an experiment name, you must use the **< Back** button and edit them in one of those previous screens.

In the **Is Control** column, check one or more boxes to designate which experiments should be used as the baseline control for variant or CNV analysis.

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Set Up Replicate Sets

If, in the [Input Sequences](#) screen, you choose the **Experiment setup** option **Multi-sample with replicates**, Set Up Replicate Sets will appear as the next wizard screen. Replicates are defined in Input Sequences, while replicate sets are specified in the Set Up Replicate Sets screen. The image below shows an example screen after the information has been filled in.

Set Up Replicate Sets
Group replicates into replicate sets

Workflow
Reference Sequence
Input Sequences
Set Up Replicate Sets
Set Up Experiments
Assembly Options

Group Individual Replicates into Replicate Sets

File	Individual Replicate	Replicate Set
Sample1_f.scf	Asian smokers	smokers
Sample1_r.scf	Asian smokers	smokers
Sample2_f.scf	African smokers	smokers
Sample2_r.scf	African smokers	smokers
Sample3_f.scf	Asian non-smokers	non-smokers
Sample3_r.scf	Asian non-smokers	non-smokers

Group Selected
Ungroup Selected

< Back Next > Cancel

- To group and name replicate sets automatically, select one or more replicates using click, **Ctrl+click**, **Cmd+click** or **Shift+click** and then press the **Group Selected** button to group the replicates and give them a single shared name. Type the name when prompted and click **OK**. All selected rows will now share the same “Replicate Set” name.
- To enter a single replicate set name manually, click on an individual “Replicate Set” cell and type in a name.
- To remove replicate sets, select one or more rows using click, **Ctrl+click**, **Cmd+click** or **Shift+click** and then press the **Ungroup Selected** button. In the selected rows, the “Replicate Set” column will return to its original (blank) state.

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

(Short Read) Polishing Options

This screen appears in long read workflows and pertains to the short read sequences being used to refine a long-read assembly. Some workflows call this the Polishing Options screen, while others call it the Short Read Polishing Options screen, but both screens are identical.

SeqMan NGen

File Help

Short Read Polishing Options

Set options for the reads being used to polish the reference sequence

- ✓ Workflow
- ✓ Input Sequences
- ✓ Preassembly Options
- ✓ Input Short Read Sequences
- ✓ **Short Read Polishing Options**
- Assembly Output

Polishing Options

Pre-assembly Options:

☒ Maximum total reads: 5000000

Assembly Options:

Mer size: ☒ Automatic ☐ Custom 21 bp

Minimum match %: ☒ Automatic ☐ Custom 93 %

Post-assembly Options:

Minimum contig size requirements:

☒ Default ☐ Custom

Min sequences: 50 Min length: 250 bp

? < Back Next > Cancel

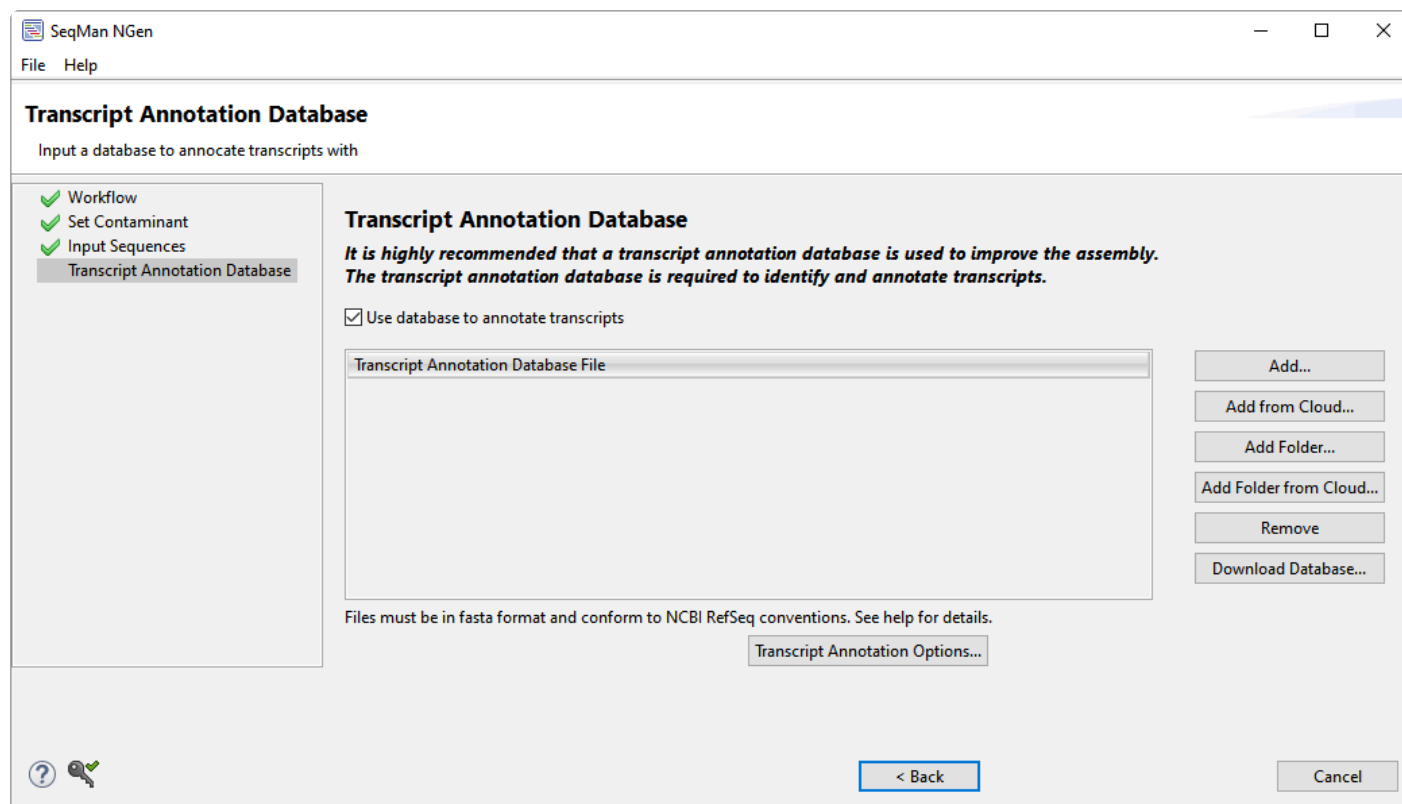
Category	Options and Descriptions
Pre-assembly Options	Check Maximum total reads and enter a value if you wish to limit the read depth. Utilizing this option can make the assembly proceed faster.
Assembly Options	Mer size is the minimum length of a mer (overlapping region of a fragment read), in bases, required to be considered a match when arranging reads into contigs. Mer size information is used to identify matches during the assembly layout phase. The default mer size is determined by the selected read technology and is shown in the window. For more information, see How mer tags are

	<p>chosen.</p> <ul style="list-style-type: none"> • Automatic – Select this button to automatically set the size based on assembly type and sequencing technology. • Custom – Select this button to choose the size yourself. You must enter the desired number of base pairs in the field at right. Lowering the mer size increases the sensitivity of finding matches, but also increases the likelihood of finding spurious matches in addition to the correct match. Lowering the mer size can also greatly increase the requirements for storing intermediate and temporary files with large projects. <p>Minimum match % specifies the minimum percentage of matches in an overlap required to join two sequences in the same contig. SeqMan NGen determines the percentage to use based on the sequencing technology you specified in the Assembly Options dialog. For more information, see Calculation of match percentage.</p> <ul style="list-style-type: none"> • Automatic – Select this button to automatically set the percentage based on assembly type and sequencing technology. • Custom – Select this button to designate the percentage yourself. You must enter a number in the field at right.
Post-assembly Options	<p>The Minimum contig size requirements options direct SeqMan NGen to remove assembled, untemplated contigs that do not meet minimum thresholds. This can lead to a desirable decrease in project size. Choose whether to use the Default values or enter Custom values. If you choose the latter, enter values for the following:</p> <ul style="list-style-type: none"> • Min sequences – Disassembles any untemplated contigs with fewer than the specified number of sequences. This option affects only untemplated contigs. No templated contigs are removed. • Min length – Using this option disassembles any untemplated contigs shorter than the specified length. This option affects only untemplated contigs. No templated contigs are removed.

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Transcript Annotation Database

Selecting the RNA-seq/Transcriptome workflow **De novo transcriptome** on the [Workflow](#) screen, causes the wizard to include the Transcript Annotation Database screen.



We highly recommend that you use a transcript annotation database when doing this workflow.

If you want to annotate transcripts during assembly using information from a specified database, leave **Use database to annotate transcripts** checked. You will need to specify a transcript annotation database, as described below. If you do **not** want to annotate transcripts during assembly, uncheck the box. If you do this, note that all other buttons in the dialog will be disabled.

To add or remove transcript annotation databases:

If you check the box as described above, you must add a transcript annotation database using any of the methods below.

- [Add a DNASTAR transcript package](#)
- [Create a custom transcript annotation database](#) and [add it to the list](#). Any files added must conform to strict formatting specifications and must be in *.fasta* format.
- [Use a local copy of RefSeq as a database](#)
- [Remove a sequence from the list](#)

To access and edit advanced options:

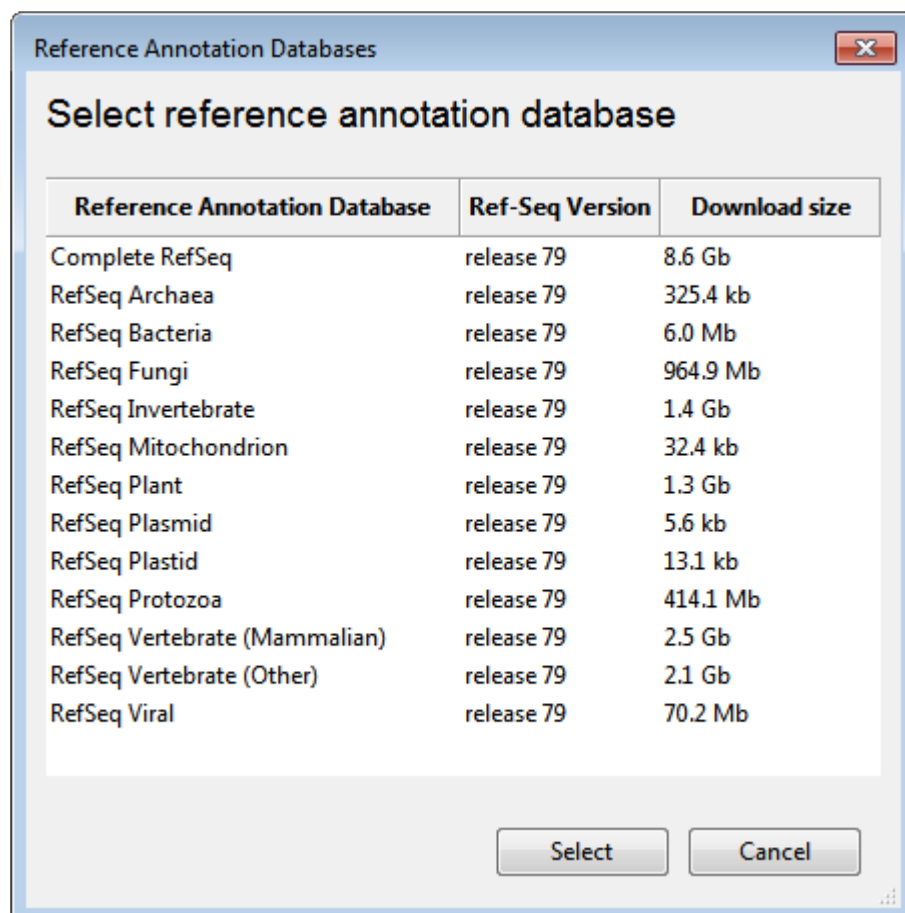
Press the **Transcript Annotation Options** button to launch the [Annotation Options dialog](#).

Click **Next >** to proceed to the next wizard screen or **< Back** to return to the previous screen.

Add a DNASTAR transcriptome package

The [Transcript Annotation Database](#) screen provides the option for licensed users to add a DNASTAR database of transcript annotations extracted from data on NCBI's [RefSeq](#) website.

To do this, click the **Download Database** button, select one or more items from the list and then press **Select**.



Create a custom transcript annotation database

The Transcript Annotation Database screen allows you upload a *.fasta_-formatted database for use in _de novo* RNA-seq [workflows](#).

A custom database must meet the same formatting specifications as NCBI [RefSeq](#) files. They must:

- Be in *.fasta* format (either single or multi-sequence files are supported)
- Use the field delimiter '|' (without quotes) between fields
- Have a header line for each entry, written in the format:

ref | [Accession] | [Organism Name] [Description] ([Gene Name])

... where:

- - **Accession** – All characters between third and fourth field delimiters
 - **Organism Name** – The first two words after fourth field delimiter
 - **Description** – All words after **Organism Name** up to the end of the line, or up to a comma or parentheses, if the gene name exists
 - **Gene Name** – All characters in parentheses after **Description**

Example:

ref | XM_005842486.1 | *Chlorella variabilis* hypothetical protein (CHLNCDRAFT_144668) mRNA, partial cds

Use a local copy of RefSeq as a transcript annotation database

To use a local copy of RefSeq as a transcript annotation database:

1. Download the latest version of the RefSeq package from this [NCBI ftp site](#).
2. Choose a de novo transcriptome RNA-seq [workflow](#).
3. On the [Transcript Annotation Database screen](#), select **Add** or **Add Folder** and select the file or folder for the RefSeq package that you downloaded in Step 1. Ignore the **Transcript Annotation Options** button on this screen. Since you are still using a RefSeq package, the default grep string and naming convention will be used automatically.
4. Finish setting up and running the assembly as usual for this workflow type.

Annotation Options dialog

To open the Annotation Options dialog, press the **Transcript Annotation Options** button in the [Transcript Annotation Database](#) screen. The Annotation Options dialog can be used to change the default naming convention used for [RefSeq](#) packages, or to specify a custom GREP and naming convention for non-RefSeq packages.

Annotation Options

Annotation Naming Convention:

☒ Use default naming convention

Primary name	Secondary name	Tertiary name
accession	none	none
description	accession	accession
geneName	description	description
organismName	geneName	geneName
	organismName	organismName

Example: geneName accession

☐ Use custom annotation name matching

Custom grep match:

Annotation Match Quality Settings:

Minimum percent of reference sequence matching the transcript: %

Minimum percent of transcript matching the reference sequence: %

Reset to default Help OK Cancel

Annotation Naming Convention:

Check the **Use default naming convention** box if you want to keep or to return to the default naming convention (geneName accession). Uncheck the box to enable dialog options allowing you to customize the naming convention.

If you uncheck the **Use default naming convention** box, the dialog provides two ways to customize names.

- Manual selection of naming components – Using the **Primary**, **Secondary** and **Tertiary** name boxes, you can select up to three “compound name” components in any order. Click in each box you want to use to enable it, then make a selection from the list: **accession**, **description**, **geneName**, **organismName** or **none**. The Example below the boxes shows your current selection(s).

Annotation Naming Convention:

☐ Use default naming convention

Primary name	Secondary name	Tertiary name
accession	none	none
description	accession	accession
geneName	description	description
organismName	geneName	geneName
	organismName	organismName

Example: geneName organismName accession

If you made the selections as shown in the preceding image, for example, then a sample *E. coli* transcript would automatically be assigned a name like: *thrB - Escherichia coli str. K-12 substr. DH10B - 6058639*

- Annotation name matching using grep syntax – If you uploaded a transcript annotation database, it automatically describes a GREP string that uses the FASTA headers to define rules for the naming convention. The same string will define the “default naming convention”. If you want to edit this GREP string, thus changing the rule regarding the extraction of contig name fragments from the FASTA headers, check the **Use custom annotation name matching** check box. After checking the box, you may edit the [regular expression/GREP](#) by typing in the **Custom grep match** textbox. An example is provided in the text box: `gi\|.*\|ref\|(?‘accession’.*)\|[]*(PREDICTED:)*\s*(?‘organismName’\w*\w*)\s*(?‘description’.*)([^(]*)(\|(?‘geneName’\w*)\|)[^([*,.]*)(\|(?‘description’.*),.)*(?‘description’.*))`

Annotation Match Quality Settings:

This section of the dialog supplies two metrics for specifying what constitutes a valid match between the reference sequence (the “query”) and the database entry.

- **Minimum percent of reference sequence matching the transcript** – Enter the minimum percentage of the query that must match the database entry. The default is 80%.
- **Minimum percent of transcript matching the reference sequence** - Enter the minimum percentage of the database entry that must match the query. The default is 50%.

To save the current selections and close the dialog, click **Save**. To return to the default settings, press **Reset to Default**. To close the dialog without saving changes, click **Cancel**.

Options tabs

Multi-tabbed advanced options dialogs are available from the:

- [Preassembly Options](#) screen for long-read workflows, by pressing the **Advanced Assembly Options** button. The following tabs are available: [Alignment](#), [Layout](#).
- [Preassembly Options](#) screen for all other workflows, by pressing the **Advanced Options** button. The following tabs are available: [Trimming](#), [Scans](#), [Alignment](#).
- [Assembly Options](#) screen, by pressing the **Advanced Options** button. The following tabs are available: [Trimming](#), [Scans](#), [Alignment](#), [Layout](#).
- [Analysis Options](#) screen, by pressing the **Advanced Analysis Options** button. The following tabs are available: [Variants](#), [Layout](#), [Peak Detection](#).

Alignment tab

This is the disambiguation page for “Alignment tab.” In SeqMan NGen, there are three advanced options dialogs that contain an Alignment tab. See the topics below for information about each version of this tab.

If you reached the Alignment tab this way...	See this User Guide topic
From the Preassembly Options screen for long read workflows, by pressing the Advanced Assembly Options button.	Click here
From the Preassembly Options screen for all other workflows, by pressing the Advanced Options button.	Click here
From the Assembly Options screen, by pressing the Advanced Options button.	Click here

Alignment tab (Preassembly Options, long read)

The Alignment tab of the Preassembly Options dialog is used to set parameters for the alignment phase of the assembly. To access the tab from the [Preassembly Options](#) screen, click the **Advanced Assembly Options** button then click on the **Alignment** tab.

Advanced Options

Alignment Layout

Advanced Alignment Options

Set parameters for both passes of the alignment phase of the long read assembly

	First pass	<input checked="" type="checkbox"/> Second Pass
Mer size: <input checked="" type="radio"/> Automatic <input type="radio"/> Custom	13 bp	8 bp
Mer search range:	300 nt	50 nt
Minimum depth percentage to confirm:	0.2	0.3

OK Cancel

The table below shows editable options. Default parameters vary according to the sequencing technology and project type specified elsewhere in the wizard, and values seldom need to be changed.

Parameter	Description
Mer size	<p>The minimum length of a mer (overlapping region of a fragment read), in bases, required to be considered a match when arranging reads into contigs. Mer size information is used to identify matches during the assembly layout phase. The default mer size is determined by the selected read technology and is shown in the window. For more information, see How mer tags are chosen.</p> <ul style="list-style-type: none"> Automatic – Select this button to automatically set the size based on assembly type and sequencing technology. Custom – Select this button to choose the size yourself. You must enter the desired number of base pairs in the field at right. Lowering the mer size increases the

	<p>sensitivity of finding matches, but also increases the likelihood of finding spurious matches in addition to the correct match. Lowering the mer size can also greatly increase the requirements for storing intermediate and temporary files with large projects. If you choose Custom, you need to enter a threshold for First pass. If you wish to do a second pass, check Second pass and enter a threshold for that pass as well.</p>
Mer search range	<p>The maximum distance upstream that the long read aligner will look for exact mer matches of the specified length. Distance is calculated in gapped coordinates. Enter thresholds for both the First pass and Second pass of the alignment phase.</p>
Minimum depth percentage to confirm	<p>Enter thresholds for both the First pass and Second pass of the alignment phase.</p>

Once you are finished, click the [Layout tab](#) to change settings there, or click **OK** to save changes and return to the [Pressembly Options](#) screen.

Alignment tab (Preassembly Options, all others)

The Alignment tab is used to set parameters for the alignment phase of the assembly. To access the tab from the [Preassembly Options](#) screen, click the **Advanced Options** button then click on the **Alignment** tab.

Advanced Options

Trimming Scans Alignment

Advanced Alignment Options

Set parameters for the alignment phase of the assembly

Gap penalty: Mismatch penalty:

Maximum gaps: Match window:

?

OK Cancel

Default parameters vary according to the sequencing technology and project type specified elsewhere in the wizard, and values seldom need to be changed.

Parameter	Description
Gap penalty	The penalty for opening or extending a gap during an alignment. This penalty is deducted from the pairwise score used to calculate match percentage . A high gap penalty suppresses gapping, while a low value promotes gapping.
Maximum gaps	The theoretical maximum length of a gap that could be inserted. In practice, the maximum gap size will usually be about half of this value. The maximum allowable value is 99.
Mismatch penalty	The penalty for a base mismatch during an alignment. This penalty is deducted from the pairwise score used to calculate match percentage .
Match window	The size of the window used to calculate match percentage .

Once you are finished, click the [Trimming](#) or [Scans](#) tabs, or click **OK** to save changes and return to the [Preassembly Options](#) screen.

Alignment tab (Assembly Options)

The Alignment tab of the Assembly Options dialog is used to set parameters for the alignment phase of the assembly. To access the tab from the [Assembly Options](#) screen, click the **Advanced Options** button then click on the **Alignment** tab. The options available in this tab vary depending on the [workflow](#).

Advanced Options

Alignment Layout

Advanced Alignment Options

Set parameters for the alignment phase of the assembly

Minimum aligned length:

Maximum gap size:

Read Handling

☒ Auto trim reads (alignment-based trimming)

☐ Trim to targeted regions

☒ Combine duplicate reads

☐ Remove clonal reads

Repeat Handling

☐ Place repeat reads All ▾

Deep Coverage Handling

☐ Do not limit

☒ Only in mitochondria and chloroplasts

☐ Limit all deep regions

OK Cancel

The table below shows editable options in alphabetical order; each workflow includes a subset of these options. Default parameters vary according to the sequencing technology and project type specified elsewhere in the wizard, and values seldom need to be changed.

Parameter	Description
Auto trim reads (alignment-based trimming)	If this box is checked, the ends of reads are trimmed to best match alignment to the reference. SeqMan NGen will mark the portion of the read that aligns well to the reference and will set the trimming to skip any of the poorly aligning parts of the read. Checking this option optimizes the end trimming of reads to maintain as much of the read as possible, while still meeting the minimum match percentage threshold. However, checking the box can also lead to the removal of true variant bases located near the ends of reads. The box is checked by default.
Combine duplicate reads	Duplicate reads are those which share the same starting position and the same sequence. Check this box if you wish to combine the reads and only enter one of them into the alignment. Any duplicates will be scored but not aligned. Combining duplicate reads collapses reads with identical sequences with the same start and stops and replaces them with a single entry with a suffix “[dup #]” where # is the number of collapsed reads. However, this option does not take the location of a paired end read into consideration. It is used primarily to reduce issues with alignment and visualization of very deep sequence regions, typical of RNA-Seq data for highly expression genes.
Deep coverage handling	<p>This section of the Alignment tab dialog lets you specify whether and when to filter deep coverage regions. The default selection is made automatically by SeqMan NGen based on the current workflow. However, you may make any desired selection. The way in which very deep coverage is handled can greatly affect assembly time. For example, if unlimited deep coverage is allowed, it could take upwards of eight hours to align human NA24385 sequences (Genome in a Bottle “Ashkenazim Trio Son”) to the mitochondrial (MT) genome. That is because most of that human sample has a coverage depth of 35,000. In such a case, limiting deep coverage regions can allow the assembly to proceed much more quickly. Choose between:</p> <ul style="list-style-type: none"> • Do not limit deep coverage regions – Use all reads, regardless of depth. With the exceptions described below, this is the default for most reference-guided workflows. • Only limit deep coverage regions for Mitochondria and Chloroplasts – Limit the depth of coverage only in the specified areas. The default for all assemblies using references that include mitochondrial and chloroplast genomes. • Limit all deep coverage regions – Limit the depth of coverage along the entire length of the sequence. The default for all miRNA and microbial genome workflows.
Layout align	In cases where a read has an identical, or nearly identical, overlap score to more than one location on the reference, indicative of a repeated sequence, the read can be evaluated by attempting a fully gapped alignment to each potential mapping position and selecting the position with the best score. In case of ties, the read is placed in one of the locations at random. Checking this box will further lower the false discovery rate (FDR), but may substantially increase the assembly time. The default is for this box to be unchecked.
Layout stringency	To specify the non-permanent “soft” filters for SNP data. SNPs that do not meet thresholds specified in this section are removed from certain displays (e.g., tables) but are still retained in

	the final project and may be displayed in downstream analysis, if desired. Specify Low , High or Custom stringency. Choosing Custom enables additional options.
Maximum gap size	The theoretical maximum length of a gap that could be inserted. In practice, the maximum gap size will usually be about half of this value. The maximum allowable value is 99.
Minimum aligned length	The minimum length of at least one aligned segment of a read after trimming. The default value varies depending on the read technology you selected. Allowed values are 0-999.
Place repeat reads	<p>Choose to place repeated reads Once, All or Never. The default is All for metagenomics workflow and Once for all other workflows.</p> <p>The way in which very deep coverage is handled can greatly affect assembly time. For example, if unlimited deep coverage is allowed, it could take upwards of eight hours to align human NA24385 sequences (Genome in a Bottle “Ashkenazim Trio Son”) to the mitochondrial (MT) genome. That is because most of that human sample has a coverage depth of 35,000. In such a case, limiting deep coverage regions can allow the assembly to proceed much more quickly.</p>
Remove clonal reads	<p>Clonal reads, where the sequence and endpoints of both reads in a pair match those in another pair, are usually the result of PCR artifacts. Check this box if you wish to retain one of the pairs in the assembly, but completely remove the clones (duplicate pairs) after the alignment phase of assembly. If the box is checked, cloned reads will not be scored, and will not be included in SNP calculation or gene quantification. This option can be useful in genome/exome/gene panel sequencing workflows where clonal reads can skew variant calculations. Checking the box may add substantially to the time required for assembly.</p> <p>Checking this option does not remove a pair if its two reads are duplicates of different pairs. It only removes duplicate pairs if the entire pair is completely identical to another pair. For example, SeqMan NGen will not remove a pair whose forward read is a duplicate of a read from pair A, but whose reverse read is a duplicate of a read from pair B.</p> <p>Note: Do not check both Combine duplicate reads and Remove clonal reads, as this will lead to unpredictable results due to the order in which SeqMan NGen removes clones and combines duplicates.</p>
Trim to targeted regions	This box is only enabled for workflows that offer the ability to add a <i>.bed</i> file, and where a <i>.bed</i> file was specified in the Reference Sequence screen. If this box is checked, reads extending beyond the 5' or 3' end of a targeted region will be trimmed to the target boundary. The box is unchecked by default.

Once you are finished, click the [Layout](#) or [Trimming](#) or [Scans](#) tabs to change those settings. Or click **OK** to save all changes and return to the [Assembly Options](#) screen.

Layout tab

This is the disambiguation page for “Layout tab.” In SeqMan NGen, there are three advanced options dialogs that contain a Layout tab. See the topics below for information about each version of this tab.

If you reached the Layout tab this way...	See this User Guide topic
From the Preassembly Options screen for long read workflows, by pressing the Advanced Assembly Options button.	Click here
From the Assembly Options and Analysis Options screens, by pressing the Advanced Options button and then the Layout tab.	Click here

Layout tab (Preassembly Options, long read)

The Layout tab of the Preassembly Options dialog is used to set parameters for the layout phase of the assembly. To access the tab from the [Preassembly Options](#) screen, click the **Advanced Assembly Options** button then click on the **Layout** tab.

Advanced Options [X]

Alignment | **Layout**

Advanced Layout Options

Set parameters for the layout phase of the long read assembly

☒ Repeat handling:

- ☒ Auto-calculate mer count thresholds
- ☐ Set mer count thresholds: to

Mer size: ☒ Automatic ☐ Custom

First pass: bp

☒ Second pass: bp

Max overhang: Max overhang ratio:

Min match length: Max gap between mer matches:

Min match length ratio: Max mismatch percentage:

Max offset between mer matches:

☒ Use chimeric read detection

[?]

OK Cancel

Default parameters vary according to the sequencing technology and project type specified elsewhere in the wizard, and values seldom need to be changed.

Parameter	Description
-----------	-------------

Repeat handing	Check this option to enable the choice between automatic or manually-specified thresholds for dealing with repeated sequence. If you check this box, you will need to specify either Auto-calculate mer count thresholds or Set mer count thresholds . If you choose the latter, you need to enter minimum and maximum values in the boxes on the right.
Mer size	Choose whether to use Automatic or Custom mer size thresholds for the first pass, and optionally, the second pass. If you choose Custom , you need to enter a threshold for First pass . If you wish to do a second pass, check Second pass and enter a threshold for that pass as well.
Max overhang	The maximum allowed length of the sum of the two overhang regions.
Min match length	The minimum length of at least one aligned segment of a read after trimming. The default value varies depending on the read technology you selected. Allowed values are 0-999.
Min match length ratio	The minimum allowable overlap to read length ratio.
Max offset between mer matches	The maximum offset in bases from the main matching diagonal for a given mer match.
Use chimeric read detection	Checking the box enables a method to identify artifactual chimeric reads in the data set.
Max overhang ratio	The maximum allowable overhang to overlap length ratio.
Max gap between mer matches	The theoretical maximum length of a gap that could be inserted between mer matches. In practice, the maximum gap size will usually be about half of this value. The maximum allowable value is 99
Max mismatch percentage	The maximum value of the estimated percentage of mismatching bases in the overlap based on the estimated sequencing error rate.

Once you are finished, click the [Alignment tab](#) to change settings there, or click **OK** to save changes and return to the [Pressembly Options](#) screen.

Layout tab (Assembly or Analysis Options)

The Layout tab in the Assembly Options and Analysis Options dialogs is used to set parameters for the layout phase of the assembly. To access the tab from the [Assembly Options](#) or [Analysis Options](#) screens, click the **Advanced Options** button then click on the **Layout** tab. The options available in this tab may vary depending on the [workflow](#).

Advanced Options

Alignment Layout

Advanced Layout Options

Set parameters for the layout phase of the assembly

☐ Alignment-based layout refinement

Layout Stringency: ☐ High - Recommended for whole genomes
☐ Low - Recommended for all other workflows
☒ Custom Minimum layout length: nt

OK Cancel

Parameter	Description
Alignment-based layout refinement	If you are following a whole-genome workflow, checking this button may improve the accuracy of the assembly.
Layout Stringency	<p>Choose from the following:</p> <ul style="list-style-type: none"> • High – Recommended for whole genomes. • Low – Recommended for all other workflows. • Custom – If you wish to enter the Minimum layout length manually. This is the minimum number of identical matching bases (from the mer analysis only) for a read to be included in the layout. It is specified by an integer, with a default of 50 nucleotides. For reads shorter than 100 bases, the setting is automatically adjusted to the mer size. Increasing this number decreases the false discovery rate (FDR) and true positive rate (TPR). Allowed values are 0-999. The default is -1 nt (nucleotides).

Once you are finished, click the [Alignment](#) or [Trimming](#) or [Scans](#) tabs to change those settings. Or click **OK** to save all changes and return to the [Assembly Options](#) screen.

Peak Detection tab

The Peak Detection tab is used to set parameters for MACS peak detection. To access the tab from the [Analysis Options](#) screen, click the **Advanced Analysis Options** button then click on the **Peak Detection** tab. The options available in this tab may vary depending on the [workflow](#).

Advanced Options

Peak Detection Variants Layout

MACS Options

Set advanced parameters for MACS peak detection

Minimum Fold Enrichment over control:

P-Value cutoff:

☒ Use local lambda for every peak region

☐ Build shifting model

Shift size:

Bandwidth:

☒ Automatically calculate tag size

Tag Size:

OK Cancel

Parameter	Description
Minimum Fold Enrichment over control	This parameter controls how enriched a peak must be, compared to the background read distribution, in order to be considered in building the peak model. If a previous assembly attempt returned the error message “too few paired peaks to build a model,” we recommend using a lower number in this box.
P-Value cutoff	Local read distribution is compared to the threshold value entered here in order to calculate whether a peak should be counted as “real.” SeqMan NGen calculates the likelihood that a detected peak is actually a peak based on the local read distribution and only returns peaks with values below the P-Value cutoff . The default P-value cutoff is 0.0001. This doesn’t mean that values greater than 10^{-5} are filtered out, but rather that they are not included in the count

	of “real” peaks.
Use local lambda for every peak region	Lambda is a parameter used to define a Poisson distribution which MACS uses to determine the expected number of reads in a given region. When this option is checked, the Poisson distribution is calculated for the peak region and for three regions surrounding the peak. MACS can use this information to determine the expected number of reads in a given region. If the option is unchecked, then the local distributions are not calculated. Instead the expected distribution is based on the total number of reads and the effective size of the genome.
Build shifting model	<p>Check this option to have MACS build a model based on the data to determine the width of and the distance between the “paired peaks.” Alternatively, leave this option unchecked to set Shift Size and Bandwidth values manually. The Shift Size is the distance each of the paired peaks will be shifted to try to center them over the actual binding site. The Bandwidth value defines the expected width of peaks. SeqMan NGen will search for peaks using a window twice as long as the bandwidth.</p> <p>Note: The creators of the MACS algorithm advise disabling model building when dealing with broad peaks (i.e. binding sites); for example, when studying histone binding.</p>
Shift size	This parameter is described in Zhang et al., 2008 .
Bandwidth	This parameter is described in Zhang et al., 2008 .
Automatically calculate tag size / Tag Size	MACS treats all reads as though they have equal length. To explicitly specify that length, check Automatically calculate tag size and enter a length in the Tag Size text box.

Scans tab

The Scans tab is used to set scanning parameters. The tab can be accessed in either of two ways:

- From the [Assembly Options](#) screen, click the **Advanced Options** button then click on the **Scans** tab.
- From the [Preassembly Options](#) screen, click the **Advanced Options** button then click on the **Scans** tab.

Advanced Options

Alignment Trimming **Scans**

Advanced Scan Options

Set scanning parameters

Vector / Adaptor scan

Mer length: Minimum matches:

Trim length: Trim to end:

Repeat scan

Mer length: Minimum matches:

Flag length:

Contaminant scan

Mer length: Minimum matches:

OK Cancel

Default parameters vary according to the sequencing technology and project type specified elsewhere in the wizard, and values seldom need to be changed.

Section	Parameter	Description
Vector/Adaptor scan section	Mer length	The minimum length of a mer required to be considered an exact match when searching for vector.
	Trim length	The minimum length required for a mer to be considered as a match for

		vector trimming.
	Minimum matches	The minimum number of matching mers required to start an alignment.
	Trim to end	The distance to the endpoint where trimming will go all the way to the end of the sequence.
Repeat scan section	Mer length	The minimum length of a mer required to be considered an exact match when scanning for repeats.
	Flag length	The minimum length required for a mer to be flagged as a repeat.
	Minimum matches	The minimum number of matching mers required to be considered a repeat.
Contaminant scansection	Mer length	The minimum length of a mer required to be considered an exact match when scanning for contaminants.
	Minimum matches	The minimum number of matching mers required to mark the sequence as a contaminant.

Once you are finished, click the [Alignment](#), [Layout](#) or [Trimming](#) tabs (if you came via [Assembly Options](#)) or the [Trimming](#) or [Alignment](#) tabs (if you came through [Preassembly Options](#)) to change those settings. Or click **OK** to save all changes and return to the [Assembly Options](#) or [Preassembly Options](#) screens.

Trimming tab

This is the disambiguation page for “Trimming tab.” In SeqMan NGen, there are two advanced options dialogs that contain a Trimming tab. See the topics below for information about each version of this tab.

If you reached the Trimming tab this way...	See this User Guide topic
From the Preassembly Options screen for non-long-read workflows, by pressing the Advanced Assembly Options button.	Click here
From the Assembly Options screen, by pressing the Advanced Options button.	Click here

Trimming tab (Preassembly Options, all others)

The Trimming tab is used to set parameters for the trimming phase of the assembly. To access the tab from the [Preassembly Options](#) screen, click the **Advanced Options** button then click on the **Trimming** tab.

Advanced Options

Trimming Scans Alignment

Advanced Trimming Options
Set parameters for the trimming phase of the assembly

Quality-based end trimming

Trim quality:

Window: bp

Fixed-end trimming

☐ 5' end bp ☐ 3' end bp

☐ Measure from 5' end

?

OK Cancel

Default parameters vary according to the sequencing technology and project type specified elsewhere in the wizard, and values seldom need to be changed.

Parameter	Description
Trim quality	The minimum averaged quality score of the evaluated window that is required in order to be considered low-quality.
Window	The length of the window to be used for averaging quality scores.
Fixed-end	If desired, check the 5' and/or 3' boxes and enter the number of bases to trim from the 5' and/or 3'

trimming

end of each read. Alternatively, check the **Measure from 5' end** box and enter a specific coordinate to which reads should be trimmed.

Once you are finished, click the [Scans](#) or [Alignment](#) tabs, or click **OK** to save changes and return to the [Pressembly Options](#) screen.

Trimming tab (Assembly Options)

The Trimming tab of the Assembly Options dialog is used to set parameters for the trimming phase of the assembly. To access the tab from the [Assembly Options](#) screen, click the **Advanced Options** button then click on the **Trimming** tab.

Advanced Options

Alignment Trimming Scans

Advanced Trimming Options

Set parameters for the trimming phase of the assembly

Quality-based end trimming

NGS

Minimum quality: Window: bp

OK Cancel

Default parameters vary according to the sequencing technology and project type specified elsewhere in the wizard, and values seldom need to be changed.

Parameter	Description
Minimum quality	The minimum averaged quality score of the evaluated window that is required in order to be considered low-quality.
Window	The length of the window to be used for averaging quality scores.

Once you are finished, click the [Alignment Layout](#) or [Scans](#) tabs to change those settings. Or click **OK** to save all changes and return to the [Assembly Options](#) screen.

Variants tab

The Variants tab is used to set parameters for the variant analysis phase of the assembly. To access the tab from the [Analysis Options](#) screen, click the **Advanced Analysis Options** button then click on the **Variants** tab. The options available in this tab may vary depending on the [workflow](#).

Advanced Options

Peak Detection Variants Layout

Advanced Variant Options

Set parameters for the variant analysis phase of the assembly

Editable Variant Filters

Filter stringency: Low

Minimum variant percentage: 15

P not ref: 90.0

Depth: 20

Fixed Variant Filters

Minimum variant percentage: 5.0

P not ref: 10.0

Minimum variant count: 2

Minimum base quality score: 5

Minimum strand coverage: 0

Maximum strand bias:

Bases to mask at ends of reads: 0

☐ Bayesian-based removal of heterozygous indels

OK Cancel

The table below shows editable options in alphabetical order; each workflow includes a subset of these options. Default parameters vary according to the sequencing technology and project type specified elsewhere in the wizard, and values seldom need to be changed.

Parameter	Description
Heterozygous peak threshold	This option is only available if you are using Sanger trace data. It is designed to identify positions in a read that contain two different bases that are both real. This can occur, for example, when you sequence a PCR product from a diploid genome at sites that are heterozygous. The percentage threshold is the minimum height of the secondary peak relative to the primary peak that's required to call the second base. Increasing the percentage increases the stringency at the cost of potentially increasing false negatives; decreasing the percentage calls more positions at the cost of potentially increasing false positives.
Editable Variant Filters section	
This section lets you specify the non-permanent “soft” filters for SNP data. SNPs that do not meet thresholds specified in this section are removed from certain displays (e.g., tables) but are still retained in the final project and may be displayed in downstream analysis, if desired.	
Filter stringency	Use the down menu to specify low , medium , high or custom stringency. Choosing custom enables the next three options in the dialog, Otherwise, these options are disabled and instead populated with unchangeable default values based on your stringency selection.
Minimum variant percentage	The minimum percent of non-reference bases required to call a SNP. When it performs SNP passes, SeqMan NGen will include regions in an assembly that have coverage less than or equal to the specified value. The default value is 5. A non-zero value is recommended when using Ion Torrent data, or working with larger genomes or doing population studies. Very low values will lead to larger files, but do not necessarily result in better SNP calls. This is only enabled when Custom is chosen as the Filter stringency .
P not ref	The minimum SNP quality score (Q_{call}) required to include a position as a putative SNP. For more information on the several ways to set P not Ref, see the topic Filter based on P not ref . This is only enabled when Custom is chosen as the Filter stringency .
Depth	The minimum depth of coverage required to include a position as a putative SNP. This is only enabled when Custom is chosen as the Filter stringency .
Fixed Variant Filters section	
This section lets you specify permanent “hard” filters for SNP data. SNPs that do not meet thresholds specified in this section are permanently deleted without saving and will not be displayed at any point downstream.	
Minimum variant percentage	The minimum percent of non-reference bases required to call a SNP. When it performs SNP passes, SeqMan NGen will include regions in an assembly that have coverage less than or equal to the specified value. The default value is 5. A non-zero value is recommended when using Ion Torrent data, or working with larger genomes or doing population studies. Very low values will lead to larger files, but do not necessarily result in better SNP calls. Minimum variant percentage and Minimum variant count can be used in tandem to control the number of reportable SNPs, and by extension, the size of the SNP table.
P not ref	The minimum SNP quality score (Q_{call}) required to include a position as a putative SNP. For more information on the several ways to set P not Ref, see Filter based on.

Minimum variant count	The minimum number of non-reference bases required to call a SNP. When it performs SNP passes, SeqMan NGen will include regions in an assembly that have coverage less than or equal to the specified value. Minimum SNP percentage* and Minimum variant count can be used in tandem to control the number of reportable SNPs, and by extension, the size of the SNP table.
Minimum base quality score	The minimum quality score below which a base will not be considered.
Minimum strand coverage	The minimum number of reads from each strand required to call a variant at a given position.
Maximum strand bias	<p>Strand Bias (SB) for a SNP is the bias for the SNP appearing on one strand versus the other. It is measured relative to the strand bias in the assembly at the location of the SNP. For example, in a column with 60 forward reads and 40 backward reads, 6 SNP bases on the forward strands, and 4 on the reverse strands would be unbiased.</p> <p>SB is given by the formula: $SB = SNP\%_f - SNP\%_r / \text{Total SNP}\%$</p> <p>...where $SNP\%_f$ and $SNP\%_r$ are the percentage of reads containing the variant on the forward (top) and reverse (bottom) strands, respectively; and $SNP\%$ is the total percentage of reads containing the variant. SB is calculated based on an “absolute value,” and will therefore be a positive number.</p> <ul style="list-style-type: none"> • 0 – Perfectly balanced (unbiased) strands. Reads with variants are present on both strands, and variants appear equally on both stands. . • 0-1, not inclusive – As the number ‘1’ is approached, more variants are called with unbalanced variants containing reads at that position. . • 1 – All variant-containing reads are on a single strand. <p>Note: If Maximum strand bias is blank or absent in the wizard, this indicates that the corresponding scripting parameter has been turned off in the script. For more information and an example, search this help document for the scripting parameter snp_maxStrandBias.</p>
Bases to mask at ends of reads	The specified number of bases from both the 5’ and 3’ ends of each read will be masked from the SNP caller and will not be considered during variant calling.
Bayesian-based removal of	Check this box to turn on <i>H-factor</i> , a Bayesian-based model that excludes heterozygous calls. If you want to view the MID column in the ArrayStar SNP Report, you must check this box. By default, the box is unchecked.

heterozygous indels	
--------------------------------	--

Once you are finished, click **OK** to save changes and return to the Assembly Options screen, or **Cancel** to return without saving changes.

Filter based on “P not Ref”

In reference-guided [workflows](#), “P not Ref” is the probability that the base does not match the reference. The P not Ref cutoff can be set using “hard” and/or “soft” filters. The following table describes the wizard parameters and their corresponding scripting commands that relate to P not Ref filtering.

To specify a “hard” filter:

In a “hard” filter, data not matching the criterion are permanently removed from the assembly. To set a hard filter, use the **P not ref** parameter located in the [Variants tab](#), accessed by pressing the **Advanced Analysis Options** button in the [Analysis Options](#) screen.

To specify a “soft” filter:

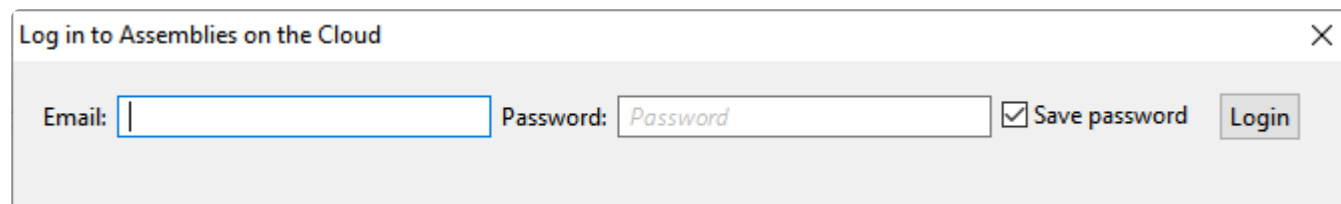
In a “soft” filter, data not matching the criterion are removed from the default display of the SeqMan SNP table. To set a soft filter:



- In the wizard – Use the **SNP filter stringency** radio buttons in the [Analysis Options](#) screen.
- In a [manual script](#) – Use the [computeSNP](#) > snpFilter > pNotRefMinVal parameter.

Log in to Cloud Assemblies

To log in to Cloud Assemblies or to access the Cloud Data Drive:

If you are not logged into your DNASTAR account and press a button associated with Cloud Assemblies or the Cloud Data Drive (e.g., the **Add from Cloud** button in the Reference Sequence screen), a popup dialog will prompt you to enter your **Email** and **Password**.

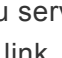


You can also access this popup dialog by pressing the “key” icon in the lower left corner of the SeqMan NGen wizard. The key has different appearances depending whether you are currently logged in () or not ().

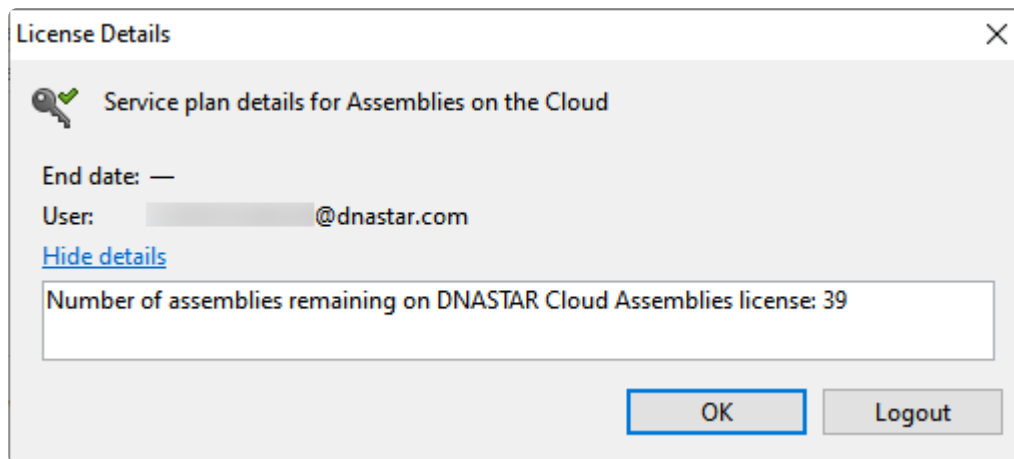
To log in, enter the information that you use when logging in to the DNASTAR website. Check the **Save password** box if you would like SeqMan NGen to save this information so you do not need to retype it in the future. Press **Log in**.

If you forgot your login credentials, click the **Forgot password** link.

To view your Cloud Assembly service plan details:

Once you are logged in as above, click the “key” icon () to display the License Details popup. Your service plan **End date** and your **User** name are shown. To see additional details, click the **Show details** link.

- If you have assemblies remaining on a current service plan, the popup expands to show the number of assemblies remaining on your license. Click **OK** to close the message or **Logout** to exit from Cloud Assemblies.



- If no assemblies remain or your service plan has expired, the message will instead ask you to purchase a Cloud license through DNASTAR Support (support@dnastar.com).

Use the DNASTAR Cloud Data Drive

The DNASTAR Cloud Data Drive works as both a Cloud file browser and a mechanism for transferring data between your desktop/laptop computer and the DNASTAR Cloud. Your data and results are stored in a private, encrypted Amazon Web Services (AWS) account visible only to you.

As a file browser:

The DNASTAR Cloud Data Drive lets you view files and folders, and create/delete folders and subfolders.

As a data transfer tool:

The DNASTAR Cloud Data Drive lets you upload data files and folders to the DNASTAR Cloud, or download them from the Cloud to your own computer. The status and progress of uploads and downloads are continuously displayed in the lower pane of the dialog.

The following brief video shows how easy it is to upload your data and then perform a SeqMan NGen Cloud Assembly.

License and Credential Requirements

Requirements for utilizing the DNASTAR Cloud Data Drive:


You must have Internet access and an active license for Assemblies on the Cloud, NovaFold, or the DNASTAR Cloud Desktop. To purchase these items, please [request an online quote](#).

Licensing options for the DNASTAR Cloud:

There are two licensing options:

- DNASTAR provides you with credentials for Assemblies on the Cloud, DNASTAR Cloud Desktop, and/or NovaFold. These are accessed through your [DNASTAR account](#).
- You set up an account with Amazon Web Services (AWS) and create an Access Key for the DNASTAR Cloud Desktop or NovaFold. To sign up for AWS, follow the instructions on [this AWS page](#).

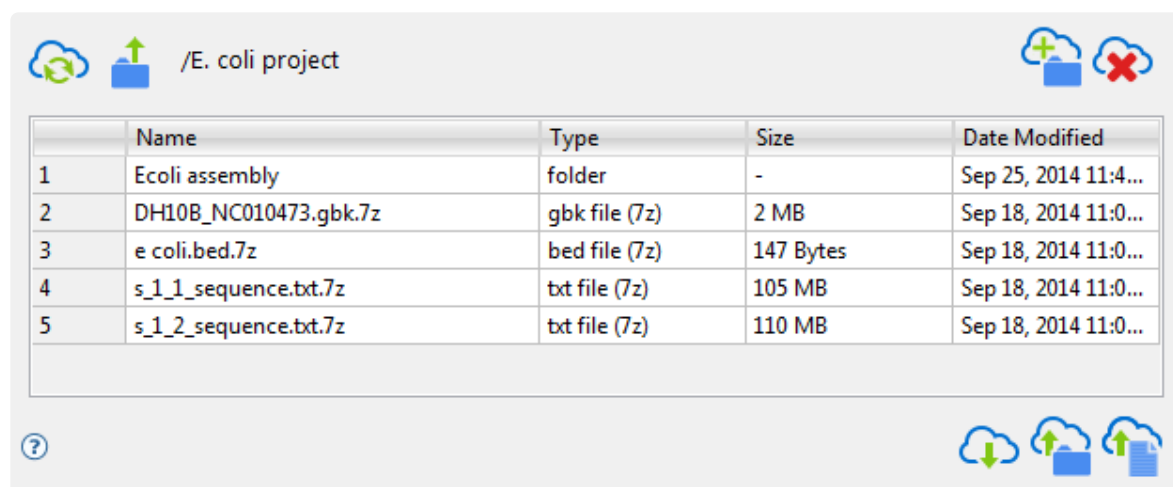
You may need to enter these credentials the first time you attempt to access the DNASTAR Cloud Data Drive.






 **Note:** During the time that your DNASTAR Cloud license is in effect, we recommend that you not change the access key. For more information about managing access keys, please refer to [this AWS page](#).




The DNASTAR Cloud Data Drive User Interface

The DNASTAR Cloud Data Drive consists of an upper section with colorful tool icons and a lower section with buttons.

The table in the top half of the dialog shows files and folders already being stored in the Cloud.



Task	How to...
To refresh the current Cloud folder	Click this icon:  . Use this option if a file/folder that you recently uploaded is not yet visible in the top table of the Cloud Data Drive, or if a file/folder you recently deleted continues to be visible there.
To create a new folder on the Cloud	Click this icon:  . See Create a New Cloud Folder for detailed instructions.
To open a folder	Double-click on the folder. Note that before uploading data (see second table in this help topic), you must open the folder in which the data will be contained.
To move up one folder level	Click this icon:  .
To delete the selected items from the Cloud	Click this icon:  . See Permanently Remove Files and Folders from the Cloud for detailed instructions.
To upload a folder and its contents from your physical computer to the Cloud	Click this icon:  . Once uploaded, folders cannot be moved into or out of another folder. This means that before you upload a folder, you need to open the parent folder (if any) that will contain it. For details on uploading, see Transfer a Folder from a Physical Computer to the Cloud .

To upload files from your physical computer to the Cloud	Click this icon:  . Once uploaded, files cannot be moved into or out of another folder. This means that before you upload a file, you need to open the parent folder (if any) that will contain it. For details on uploading, see Transfer Files from a Physical Computer to the Cloud .
To download files from the Cloud to your physical computer	Click this icon:  . For details, see Transfer Files or Folders from the Cloud to a Physical Computer .
To view the online help	Click the help icon:  .

The table in the bottom half of the dialog displays files or folders that are in the process of being transferred between the Cloud and a desktop/laptop computer.

Recent Uploads/Downloads (4)

	Name	Status	Elapsed	Progress	Remove
1	DH10B_NC010473.gbk	Upload	00:00:18	Complete	No
2	e coli.bed	Upload	00:00:01	Complete	No
3	s_1_1_sequence.txt	Uploading	00:00:10	48%	No
4	s_1_2_sequence.txt	Compressing	00:00:10	41 MB	No

The following table contains descriptions for each of the columns:

Column	Description
Name	The file name and extension.
Status	<p>The stage of the upload or download process that is currently in progress.</p> <ul style="list-style-type: none"> For an upload, the steps are: Compressing, Compression, Uploading, Upload (shown both when Progress = Pending and Complete). For a download, the steps are: Decompressing, Decompress, Download (shown both when Progress = Pending and Complete).
Elapsed	Time elapsed since initiating the upload/download in Hours:Minutes:Seconds .
Progress	During compression/decompression of files, this column displays the size of the file. During upload/download, it instead displays information about the percent of the job completed so far: Pending (i.e., 0%), n% , Complete (i.e., 100%).
Remove	Press the Remove button if you would like to remove the file from the queue and thereby cancel the



upload/download.

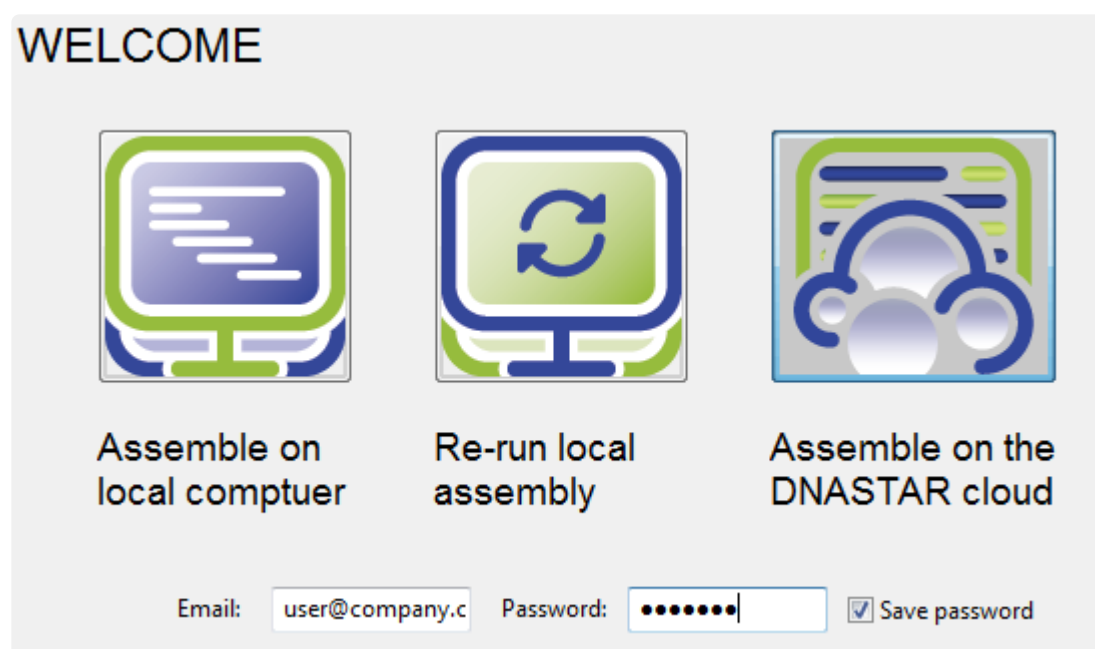
The buttons at the bottom right of the Cloud Data Drive are used to [exit from the application](#).

Access the DNASTAR Cloud Data Drive


Once you have met the [license and credential requirements](#), the Data Drive can be accessed in several ways:

Through the SeqMan NGen wizard:

1. In the Welcome screen, select **Assemble on the DNASTAR Cloud** (currently available on Windows and Macintosh only) if you are licensed for this option and wish to perform one or more assemblies on the Cloud, rather than using your own physical machine.
2. Enter your **Email** and **Password**. Enter the information that you use when logging in to the DNASTAR website. Check the **Save password** box if you would like SeqMan NGen to save this information so you do not need to retype it in the future.



3. Click **Next**. If you are licensed, a message will display the number of assemblies remaining on your license. Click **OK** to close the message and continue to the next screen.

 **Note:** If you are *not* licensed, a message will appear, asking you to purchase a Cloud license through DNASTAR Support (support@dnastar.com).

4. In the Cloud Assembly page:
 - If you intend to start a new assembly project, click the **Upload Data** button to access the DNASTAR Cloud Data Drive. From within the Cloud Data Drive, upload any [files](#) or [folders](#) that

will be needed for the assembly. As you progress through the SeqMan NGen wizard, you can access the files on the Cloud Data Drive by using the **Add**, **Add Folder** or **Browse** buttons.

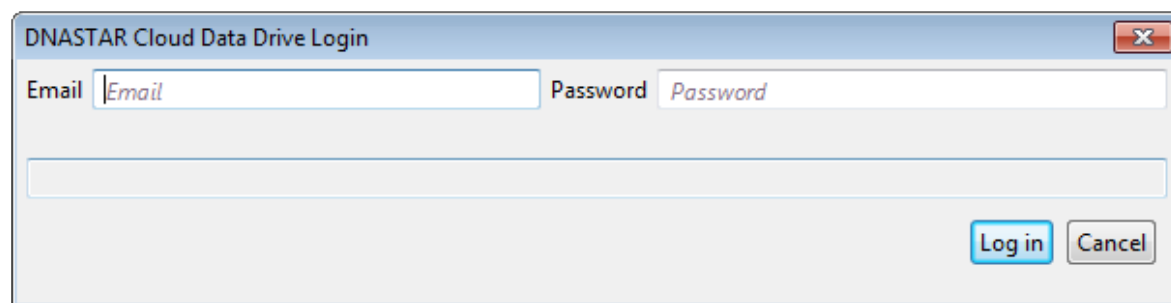
or...

- If you want to monitor an existing assembly, click **Manage/monitor existing cloud projects** and click **Next** to open the Cloud Assemblies screen. From this screen, make a selection in the upper table and then choose either **Download Assembly** or **View Data** to open the DNASTAR Cloud Data Drive.

As a standalone application:

From your computer's hard drive or from the DNASTAR Cloud Desktop, launch the DNASTAR Navigator and then click on **DNASTAR Cloud Data Drive**.


If you are not currently logged in to the DNASTAR Cloud Data Drive, you will first be prompted to log in.

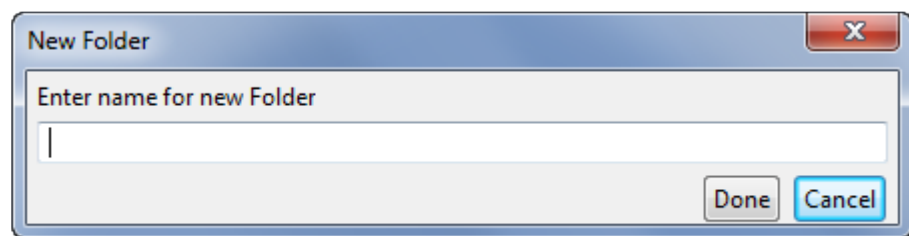
A screenshot of a Windows-style dialog box titled "DNASTAR Cloud Data Drive Login". The dialog has a light gray background and a blue title bar. It contains two input fields: "Email" with a placeholder text "Email" and "Password" with a placeholder text "Password". Below these fields is a large empty rectangular box. At the bottom right, there are two buttons: "Log in" (highlighted in blue) and "Cancel".

Enter the User ID and Password from your DNASTAR account, then click **Log in**.

Create a New Cloud Folder

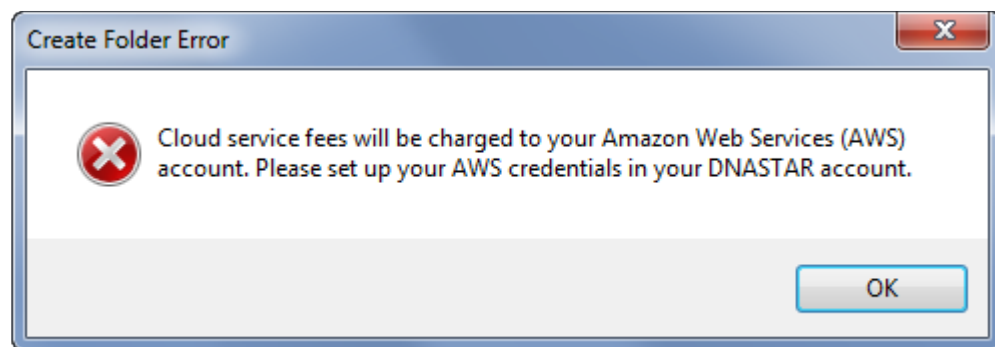
To create a new folder in the DNASTAR Cloud:

Click the **Create a new folder on the Cloud** tool () and type in a name for the new folder.



Press **Done**; or choose **Cancel** to exit without creating a folder.

If your AWS credentials have not been set up, you will receive the following error message:




Click **OK** to exit from the message. To address the credential issue, refer to [License and Credential Requirements](#) or contact DNASTAR Support at support@dnastar.com.

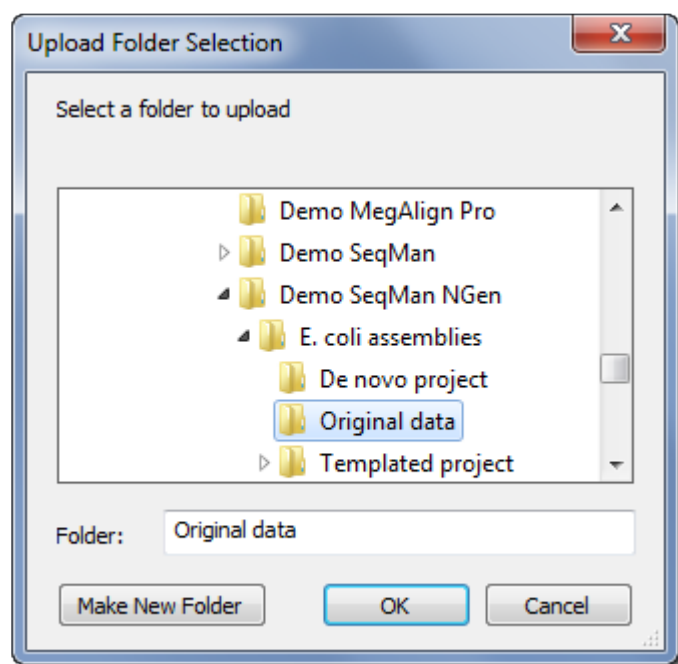
Transfer a Folder from a Physical Computer to the Cloud

To transfer a folder and its contents from your desktop or laptop computer to the DNASTAR Cloud:

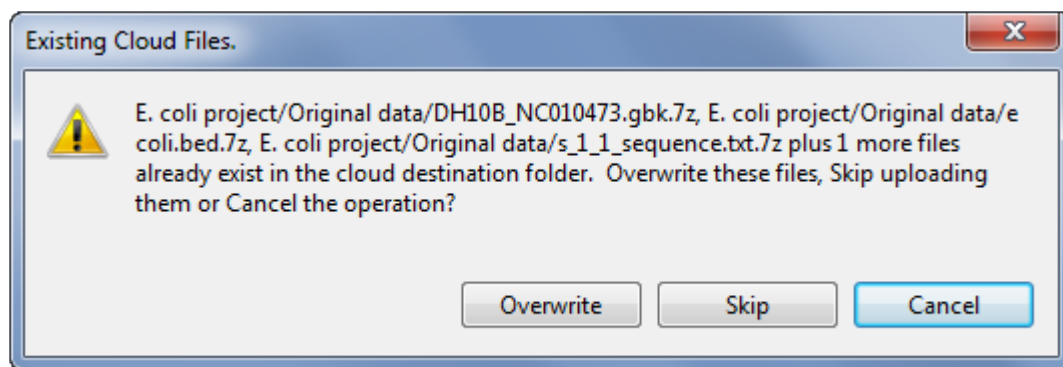
1. Within the Cloud Data Drive, open the folder (if any) that will contain the folder you are transferring.

✿ **Note:** Once a folder has been transferred to the DNASTAR Cloud Data Drive, you cannot move it into a different “containing folder.”

2. Click the **Upload a folder and all its contents to the Cloud** tool ).
3. Navigate to the file or folder location and select it.

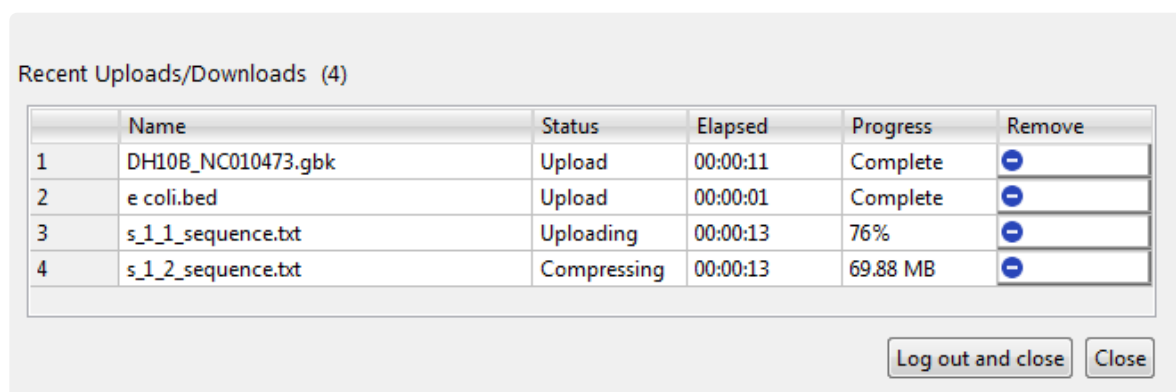


4. Press **OK**.
 - If the same folder has been uploaded previously to the same location, the following message will appear.



Choose whether to **Overwrite** the existing folder, **Skip** uploading the folder or **Cancel** the entire upload operation. In the case of a folder, the final two options are identical.

- If you have not previously uploaded the folder and its files, or if you have chosen to **Overwrite** existing files, the files within the selected folder will begin uploading in the lower half of the Cloud Data Drive dialog. If the folder was not previously uploaded, the folder will be created in the upper half of the Cloud Data Drive dialog and the files placed within.





At the completion of the upload, the folder will be displayed in the upper table of the dialog. Double-clicking on the folder name will reveal the folder contents.

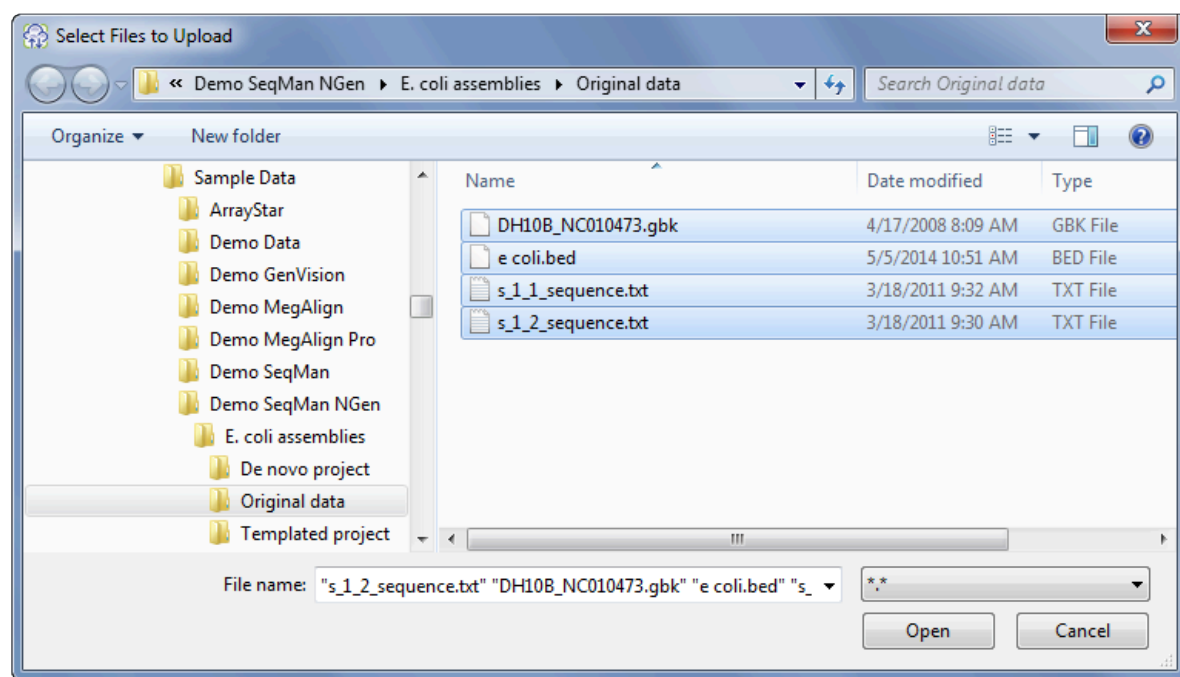
Transfer Files from a Physical Computer to the Cloud

To transfer one or more files from your desktop or laptop computer to the DNASTAR Cloud:

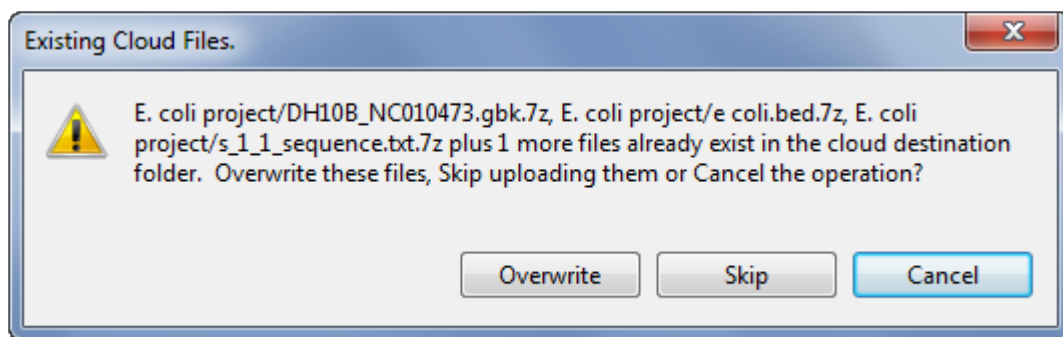
1. Within the Cloud Data Drive, open the folder (if any) that will contain the files you are transferring.

 **Note:** Once a file has been transferred to the DNASTAR Cloud Data Drive, you cannot move it into a different “containing folder.”

2. Click the **Upload files to the Cloud** tool ().
3. Navigate to the file location and select the desired file(s).

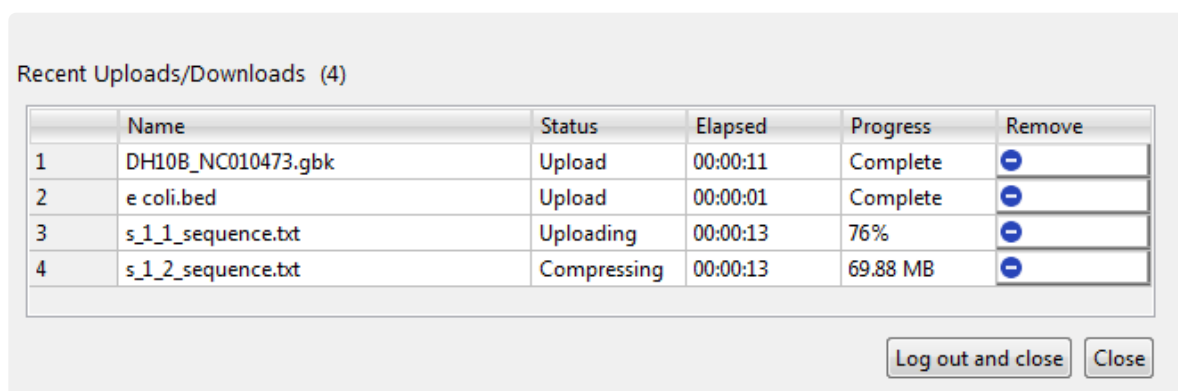


4. Press **Open**.
 - If one or more files have been uploaded to the same location previously, a warning message will appear.



Choose whether to **Overwrite** the file(s), **Skip** uploading the duplicate file(s) or **Cancel** the entire upload operation.

- If you have not previously uploaded the files, or if you have chosen to **Overwrite** existing files, the files will begin uploading in the lower half of the Cloud Data Drive dialog. If the files were not previously uploaded, they will be added to the upper half of the Cloud Data Drive dialog.




At the completion of the upload, the files will be displayed in the upper table of the dialog.

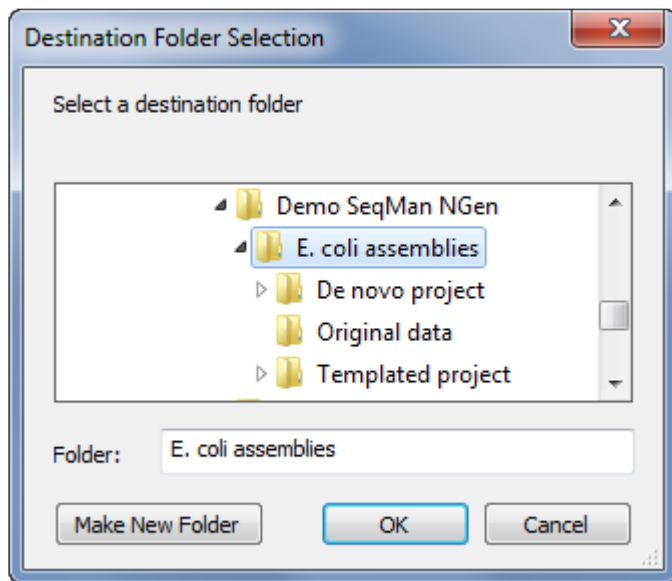
Transfer Files or Folders from the Cloud to a Physical Computer

To download files or folders from the Cloud to a physical computer:

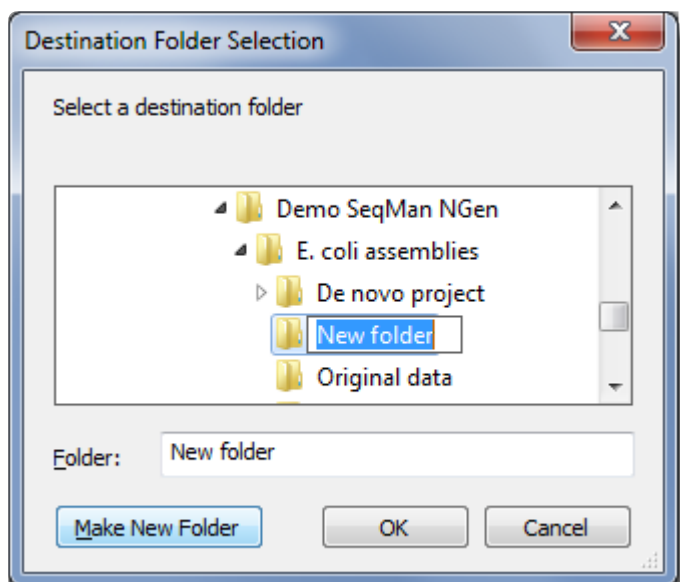
1. Select the files and/or folders from the upper table by clicking, **Shift+clicking** or **Ctrl/Cmd+clicking**.

✿ **Note:** If you do not make a selection, you will receive a warning message after Step 2.

2. Select the **Download** tool ()
3. In the ensuing dialog, navigate to the destination folder.



4. (optional) If desired, create a new sub-folder in the selected directory as follows:
 - a. Select the parent folder.
 - b. Click **Make New Folder**. A new folder is created, ready to be renamed.



c. Type a name for the folder.

5. Press **OK**.


The progress of the download(s) will be shown in the bottom half of the Cloud Data Drive. After downloading is complete, the files/folders can be found in the destination directory chosen in Steps 3 or 4.

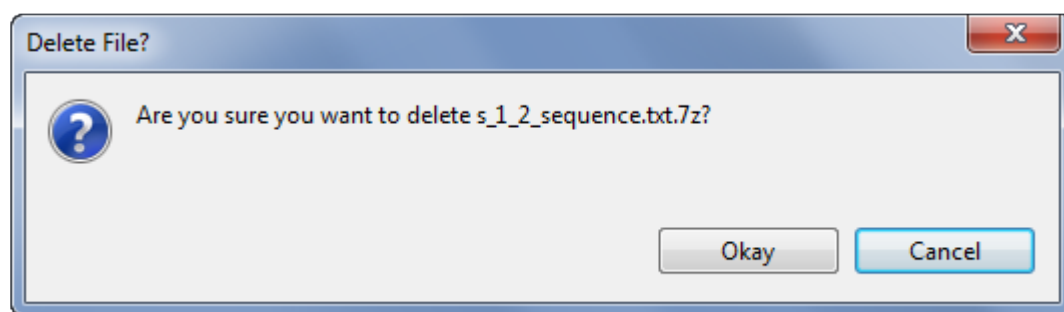
Permanently Remove Files and Folders from the Cloud

To permanently remove files and/or folders from the DNASTAR Cloud:

1. Select their rows from the upper table by clicking, **Shift+clicking** or **Ctrl/Cmd+clicking**.

✿ **Note:** If nothing is selected, you will receive a warning message after Step 2.

2. Press the **Delete** tool ().
3. In the ensuing confirmation dialog(s), confirm that you wish to delete the file(s)/folder(s) by pressing **Okay**. Alternatively, you may abort the deletion process for a given file by clicking **Cancel**.



Close the DNASTAR Cloud Data Drive

To close the DNASTAR Cloud Data Drive application:

- Click **Close** to exit.

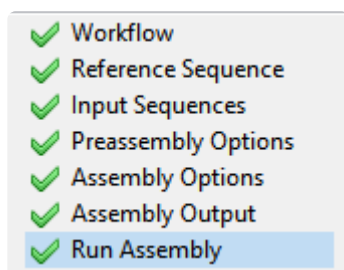
or...

- Click **Log out and close** to both close the Cloud Data Drive and log off from your DNASTAR account. After using this button, the SeqMan NGen wizard will “forget” your saved credentials, and will prompt you for your **Email** and **Password** the next time you request access to Assemblies on the Cloud.

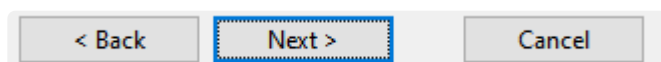
Navigate between wizard screens

There are two different ways to navigate between screens in the SeqMan NGen wizard.

- If you have already visited a particular screen, you can return to it by clicking its name from the menu at the left of the SeqMan NGen wizard. The menu contains a list of wizard screens that you have visited so far. If you return to a previous screen and make a selection that changes the workflow type, subsequent screens that are not part of the new workflow will be removed from the list



- Whether or not you have yet been to a screen, you can navigate using the buttons at the bottom of each dialog.



- Click the **Help** button (Win) or the question mark icon (Mac) to launch the user's guide topic for the current panel.
- Click **< Back** and **Next >** to navigate to the previous or next panel.
- Click **Quit** to exit SeqMan NGen.

Add and remove files in the wizard

Many [wizard screens](#) allow or require you to specify files that will be used in the assembly. There are several ways to add these files, each discussed in a separate topic:

- [Add sequences from your computer or the cloud](#)
- [Add a genome template from DNASTAR](#)
- [Add a genome template from NCBI](#)


You can also [remove a sequence from the list](#) if you change your mind.

Add sequences from your computer or the cloud

Some wizard screens require you to add at least one sequence or genome template before proceeding. In these cases, the “Next” button only appears after the item has been added.

- To add one or more sequences from your computer (i.e. “local” sequences), press **Add** or **Add Folder**. Navigate to and select the desired file(s) or folder of files, then click **Open**.

✿ **Note:** If your reference sequence is a GFF feature file, with or without a separate FASTA file, we recommend adding the GFF file first. If the FASTA is built into the GFF file, you don’t need to do anything further. If the FASTA file is separate, SeqMan NGen will automatically prompt you to add the accompanying FASTA file.

- To add one or more sequences from the DNASTAR Cloud Data Drive, press **Add from Cloud** or **Add Folder from Cloud**. Navigate to and select the desired file(s) or folder of files, then click the green check mark ()

The first time you press a button associated with Cloud Assemblies (or the Cloud Data Drive), you may be prompted to enter your **Email** and **Password**. Enter the information that you use when logging in to the DNASTAR website. Check the **Save password** box if you would like SeqMan NGen to save this information so you do not need to retype it in the future. If you are a licensed Cloud user, a message will display the number of assemblies remaining on your license. Click **OK** to close the message and continue to the next screen. If you are *not* licensed, a different message will appear, asking you to purchase a Cloud license through DNASTAR Support (support@dnastar.com).

The following brief video shows how easy it is to upload your data and then perform an assembly on the cloud:

Add a genome template from DNASTAR

If your workflow includes the [Reference Sequence](#) or [Set Contaminant](#) screen, you must add at least one reference sequence, biome genome, or genome template before proceeding.

Curated and up-to-date DNASTAR genome template packages are available for common model organisms. Each template package contains template sequence, annotations, and database linking information. If you wish to use DNASTAR's database association features (e.g., dbSNP, GERP, and COSMIC), you must input one of these genome packages in the appropriate screen for your workflow.

To add a DNASTAR genome template from your local computer:

If you have not previously downloaded the genome package to your local computer, start with step 1. Otherwise, start at step 5.

1. Press **Download Genome Package**. If the button is disabled, see the first and last notes below.
2. Select a package from the list.
3. Click **Select**. and choose a save location. The package is saved with the extension *.genometemplate*.
4. When the download finishes, click **OK**.
5. Press **Add Genome Package**.
6. Navigate to the location where you saved the automatically downloaded (or manually downloaded & extracted) package, and click **Open**.

To add a DNASTAR genome template from the Cloud:

1. Press **Add Genome Package**.
2. Select a package from the list and click **Select**.

After creating an assembly with a genome template, you can access dbSNP information as described in the following brief video:

Notes:

- The **Add Genome Package** and **Download Genome Package** buttons are disabled if you have already added files using the **Add** or **Add Folder** buttons.
- SeqMan NGen can read and produce output using a variety of common chromosome naming conventions, including “chr1” and “ch1,” as well as Arabic and Roman numerals. Chromosome names are captured from genome template packages and used to assign contig IDs to entries from BED, VCF and manifest files.
- Chromosome names are captured from genome template packages and used to assign contig IDs to entries from [BED and Manifest files](#).
- In the human genome template packages provided by DNASTAR, the “unlocated contig” is actually a concatenated, multi-sequence contig containing the alternate loci sequences. These loci are used for large regions where the human population contains variation so divergent that it cannot be adequately described by simple substitutions and small indels. Examples of these regions include the LRC/KIR complex on chromosome 19 and the MHC on chromosome 6.
- If you are performing a local assembly, there are some circumstances in which it is necessary to download and extract the package manually prior to using the wizard. To do this, go to DNASTAR’s [Genome Template Packages web page](#) and download a template package with the genome of interest. Downloaded genome packages are saved on your computer as ZIP files, and must be extracted prior to use. On Macintosh, double-click on the ZIP file. The files will be automatically extracted via the Archive Utility. On Windows, use any archive utility to extract the files. One method is to double-click on the ZIP file. In the ensuing Explorer window, click **Extract all files** from the top left. Choose a location for the files and select **Extract**.

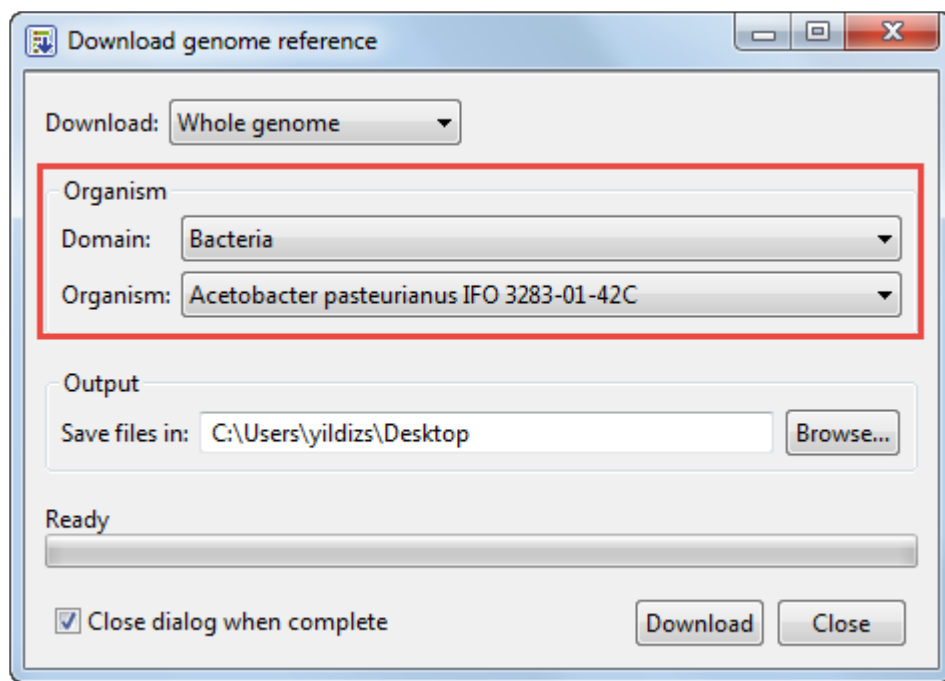
- The Human Genome Package includes Chr 1-22, X, Y, MT and “Unlocated.” DNASTAR's “Unlocated” contigs match GenBank's “Unplaced sequence,” i.e., items specifically annotated as “Homo sapiens concatenated unplaced genomic contigs.”

Add a genome template from NCBI

If your workflow includes the [Reference Sequence](#) or [Set Contaminant](#) screen, you must add at least one reference sequence, biome genome, or genome template before proceeding.

If you are doing a local (i.e. non-Cloud) assembly, you may download and/or add genomes directly from the NCBI database in either GenBank or FASTA formats.

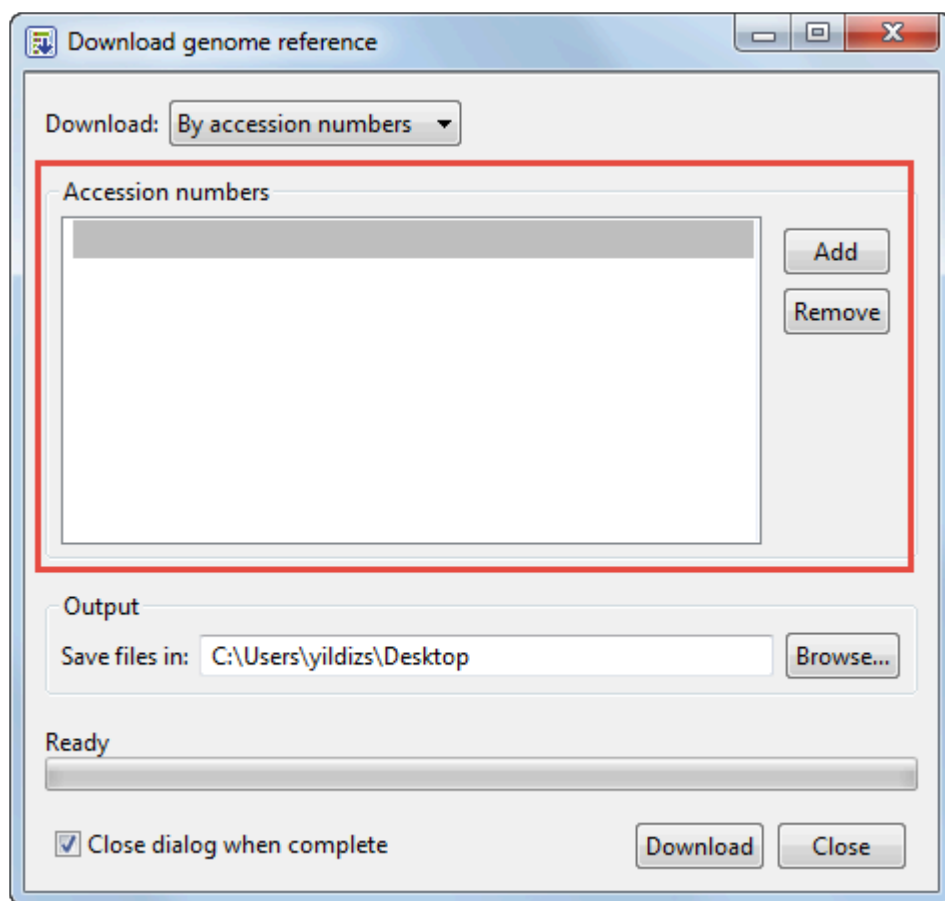
1. Press the **Download NCBI Genomes** button. This launches the Download genome reference dialog.
2. Use the **Download** dropdown menu to choose between downloading a **Whole genome** or **By accession numbers**.
 - If you select **Whole genome**, SeqMan NGen will retrieve the most recent build of the selected genome. Use the next two drop-down menus to select the **Organism type** and **Organism**.



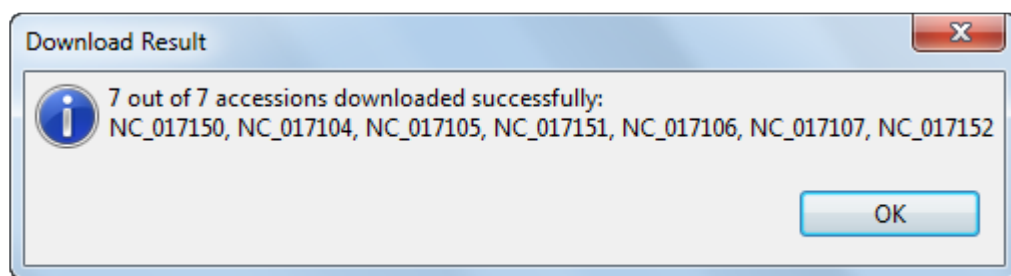
SeqMan NGen will download all the reference sequences from the NCBI Entrez Genome Project database for the selected genome. These downloads may include auxiliary genomes such as mitochondria and chloroplasts. They may also include some contigs which have not yet been placed by the genome finishing process.

- If you select **By accession numbers**, the Organism section disappears and is replaced by an Accession numbers entry area. Type an accession number or paste it from your clipboard, then press **Add** to add a number to the list. You can type accession numbers with or without explicit version numbers (e.g., NC_000913.3 or NC_000913). If the version number is omitted, then NCBI's latest version is returned.

Continue adding numbers, as desired. Multiple accession numbers should be separated using a space, comma, semi-colon or line break. To remove an accession number from the list, select it and click **Remove**.



3. Click **Browse** to select a name and location in which to save the downloaded genome files.
4. If you do not need to download additional genomes, you can check **Close dialog when complete**. Otherwise, leave the box unchecked to keep the dialog open after initiating the current download.
5. Press **Download**. Once the download is complete, a message like the one below will appear.



6. Click **OK** to close the dialog and add the accessions to the Reference Sequence screen. If you

checked **Close dialog when complete**, the Download genome reference dialog will also close. Otherwise, it will remain open so that you can download additional NCBI genomes.

Remove a sequence from the list

To remove a sequence that has been added to the list on the left of the [Input Sequences](#) screen, select the file from the list and press **Remove**.

Use editing commands in the wizard

SeqMan NGen's basic editing commands are similar to those found in Microsoft Windows and other text editing programs. These commands may be available through context (right-click) menu options in text boxes throughout the SeqMan NGen wizard.

Task	Keyboard shortcut
To remove the selected portion of text and place it on the clipboard	Ctrl/Cmd+X
To copy the selected portion of text to the clipboard	Ctrl/Cmd+C
To paste an item in the clipboard at the cursor insert point	Ctrl/Cmd+V
To delete a portion of text without placing it on the clipboard	Delete key
To select all text	Ctrl+A

Monitor the progress of a cloud assembly

If you choose to perform an assembly on the cloud, you can monitor assembly progress using the online Web Monitor. Click [this link](#), and enter your DNASTAR login credentials if prompted to do so. See the [Web Monitor help](#) to learn about features in this browser window.

DNASTAR Cloud Assemblies

[Refresh](#)
[Legend](#)

Web Monitoring

- NGen Cloud service - Number of assemblies remaining on DNASTAR Cloud Assemblies license: 16
- 9 total assemblies (4 done, 0 failed, 1 stopped, 1 queued, 3 in-process) - Updated 3/11/2019 3:36:30 PM
- For help, please use [this link](#).

Show latest assemblies

Limit to 10 assemblies

Name	Status	Started	Ended	Message	Files
Arabidopsis strains-qng	Queued				
Mutant 3.7z	Starting	3/11/2019 3:22:23 PM		Starting	
Mutant 2.7z	Starting	3/11/2019 3:22:23 PM		Starting	
Mutant 1.7z	Starting	3/11/2019 3:21:58 PM		Starting	

Access and understand output files

The output file structure for a SeqMan NGen assembly varies depending upon the workflow. For a description of output files, see [Reference-guided workflow output](#) or [De novo workflow output](#).

Note that FASTQ files “created with the SeqMan NGen wizard will have a *.fastq* extension, while those [created via a command line script](#) will have a *.fas* extension.

View the Project Report

The Project Report summarizes the assembly statistics, including the parameters used, the number of assembled/unassembled sequences and contigs in your project, and the average quality scores.

Opening the Project Report:

Use any of the methods below:

- Click the **Project Report** button in the [Assembly Summary](#) screen.
- Open the output file *Report.txt* in any suitable text editor.
- Open the assembly in SeqMan Pro and choose **Project > Report**.
- Open the assembly in SeqMan Ultra and choose **Project > Project Details**.

Project Report contents:

- Click these links for a list of the Project Report contents for [reference-guided](#) and [de novo](#) workflows.
- The contigs in the project are named as follows:
 - If you performed a reference-guided [workflow](#), the resulting contig will take the name of the reference sequence name.
 - If you scanned your assembly for [known repeats](#), then the contigs in your project containing sequences flagged as possible repeats will be named: Repeat-00001, Repeat-00002, Repeat-00003, etc.
 - If none of the above applies, the contigs in your project will be named Contig 00001, Contig 00002, Contig 00003, etc.
- If you've opted to create both an *.sqd* and an *.assembly* output, you may notice that the files do not exactly match. That's because *.sqd* files, unlike *.assembly* files, allow sequences to extend beyond either end of the reference sequence.

Project Report contents for reference-guided workflows

The [Project Report](#) for reference-guided assemblies will contain a subset of the following results:

Run Statistics	
Reference Seq Cnt	The total number of sequences in the reference (template).
Sequence Cnt	The total number of reads in the sample.
Total Reads Assembled	
Pair Seqs Cnt	The number of paired sequences included in the assembly.
Single Pair Seq Cnt	The number of paired sequences of which only one pair was included in the assembly.
Split Seq Cnt	The number of sequences that were split in the assembly.
Bad Split Seq Cnt	The number of sequences that were split, and of which only one portion was included in the assembly.
Single Seq Cnt	The number of single (unpaired) sequences in the assembly.
Consistent Pair Cnt	The number of paired sequences that met pair constraints. One “pair” in this statistics represents two sequences.
Inconsistent Pair Cnt	The number of putative paired sequences that did not meet pair constraints.
Seqs score < 80%	Percentage of reads that exactly matched the template (i.e. “alignment score”).
Seqs score < 90%	
Seqs score < 100%	
Seqs score 100%	
Unassembled Sequences	
Unaligned Cnt	Total number of reads not included in the finished assembly.
LayoutMiss Cnt	<p>The number of reads that didn’t match the template at all. In other words, the number of sequences that contained no mer which matched a mer on a template sequence. This number is affected by the assembly parameters merSize and merSkip.</p> <p>Example: A sequence that has no 21-mer in common with the template but does have a matching 17-mer would be included in LayoutMiss Cnt at a mer size of 21, but not at a mer size of 17.</p>
LayoutPoor Cnt	The number of reads with an insufficient number of mer matches to be included in the

	assembly. This number is affected by the assembly parameter merLayoutMin.
Bad Seq Cnt	The number of reads with $\geq 25\%$ ambiguous Ns in the sequence. Filtered Illumina data is sometimes included in this count, as well.
Excluded Seq Cnt	The number of contaminated reads.
ExcessiveCov. Seq Cnt	The number of reads unused due to excessive coverage.
SNP Info	
Found SNP Cnt (incl. indel lengths)	The number of SNP positions plus the total number of coalesced bases, minus the number of multi-base indel entries.
Found User SNP	The number of SNPs found that match those in the user-supplied VCF SNP file.
Missing User SNP coverage	The number of SNPs from the user-supplied VCF SNP file that were <i>not</i> found, even though the area had coverage.
Missing User SNP zero coverage	The number of SNPs from the user-supplied VCF SNP file that were not found because the area had no coverage.
Assembly Parameters	
merSize	The values specified in the SeqMan NGen wizard prior to assembly.
merSkip	
merSkipQuery	
merLayoutMin	
templateHitCntThresh	

Project Report contents for de novo workflows

The [Project Report](#) for *de novo* assemblies will contain a subset of the following results:

Assembly Totals	
Contigs	Total number of contigs assembled.
Contigs > 2K	Total number of assembled contigs that are more than 2000 base pairs in length.
Contigs to Reach Genome Length 'x'	The number of contigs needed to cover the genome length specified in the Workflow pane.
Contigs removed due to small size	The number of contigs removed due to being smaller than the threshold value.
Assembled Sequences	The number of sequences utilized in the assembly.
Unassembled Sequences	The number of sequences excluded from the assembly. These may be further categorized as: 1) Sequences not assembled due to complete trimming, and 2) Sequences removed due to small contig size.
All Sequences	Total number of sequences in the project.
Contig N50	Contig size at which 50% of the sequence data are represented. Note: In a typical microbial genome assembly, Contig N50 values exceed 80K base pairs and genome coverage is attained in less than 100 contigs. In many assemblies, contig N50 exceeds 100K with genome coverage attained in 25 contigs. If paired-end Roche 454 Life Sciences data are used, contigs can be ordered into a handful of large scaffolds to attain genome coverage that greatly facilitates gap closure and completion of the genome assembly.
Average Coverage	Average depth of coverage in the assembly.
Average Totals	
Sequences Per Contig	Average number of sequences used for each contig.
Average Lengths	

Contigs	Average contig length.
Assembled Sequences	Average length of sequences used in the assembly.
Unassembled Sequences	Average length of sequences excluded from the assembly.
All Sequences	Average length of all sequences in the project.
Average Quality	
Assembled Sequences	Average quality score of sequences used in the assembly.
Unassembled Sequences	Average quality score of sequences excluded from the assembly.
All Sequences	Average quality score of all sequences in the project.
Assembled Pair Statistics	
Read Pairs	Total number of paired reads in the project.
Assembled Pairs	The number of paired reads included in the assembly.
Pairs Consistent Within a Contig	The number of paired sequences within a single contig that met pair constraints. One “pair” in this statistics represents two sequences.
Pairs Inconsistent Within a Contig	The number of putative paired sequences within a contig that did not meet pair constraints.
Split Pair Statistics (Ion Torrent paired reads and 454 data only)	
Reads Split into Pairs	The number of reads that were split into pairs at the linker.
Unsplit Reads with Pair Linker(s)	The number of reads that were not split into reads because the linker was too far to one side.
Unsplit Reads without Pair	The number of reads that were not associated with a linker.

Linker(s)	
Assembly Parameters	
Match Size	The values specified in the SeqMan NGen wizard prior to assembly.
Match Spacing	
Minimum Match Percentage	
Match Score	
Mismatch Penalty	
Gap Penalty	
Max Gap	
Genome Length	
Expected Coverage	

Reference-guided workflow output

Reference-guided workflows vary in the number and contents of output files and folders. Only a subset of items in the table below may appear for a particular workflow.

In the file names below, the project name should be understood to precede any hyphen (-) or period (.) used at the beginning of file and folder names.

Single assemblies:

The project folder has the name specified in the [Assembly Output screen](#) and contains:

- **.script** file
- **.assembly** package (called **.transcriptome** for the transcript annotation workflow)
- **-noSplit.assembly** package (Reference-guided assembly with gap closure workflows only)
- **-Reports** folder
 - **-zinternal** folder
 - **info** folder

All Exome and Gene Panel projects and multiple sample assemblies run as separate assemblies:


The project folder contains the name specified in the [Assembly Output screen](#) followed by the suffix **_assemblies** (called **_RNA-Seq** in the reference-guided RNA-seq workflows). This folder contains:

- **.script** file (if saved in the [Run Assembly](#) screen).
- **Results.txt** file - Overview information and statistics for each assembly.
- **.table.txt** file
- **.template.script** file
- **_arstar.script** file – A script to load all assemblies as a SNP project in ArrayStar.
- **_arstarValidation.script** file – (only if validation control was present) A script to load the validation control assembly and associated VCF file as a SNP project in ArrayStar and to automatically calculate the accuracy statistics.
- **.assembly** packages (one per sample)
- **-Reports** folders (one per sample)
 - **-zinternal** folder
 - **info** folder

Contents of the .assembly package

The **.assembly** package is part of the output for [XNG](#) workflows. (The contents of the **-noSplit.assembly** package are similar to those of the **.assembly** package.)

In the file names below, the project name should be understood to precede any hyphen (-) or period (.) used at the beginning of file and folder names.

 **Note for Windows users:** To open text reports with the correct formatting displayed, we recommend using Wordpad, Notepad++, or Microsoft Excel®, and not the default Windows text editor, Notepad.

File Suffix	Description
It is intended that the entire .assembly package be opened in SeqMan Pro or SeqMan Ultra for viewing and analysis of the assembly. However, the following individual files also contain useful information.	
.vcf	<p>A VCF file (.vcf) is automatically created for all assemblies with variants. The file is modified in three ways adhering to the Variant Call File (VCF) v. 4.2 specification:</p> <ul style="list-style-type: none"> * In the FILTER field, each row is marked with one of three qualifiers to show whether or not a position was covered: <ul style="list-style-type: none"> ** "PASS" for positions where a call could be determined based on the sequence read data. ** "NC" for positions with no sequence read coverage (this will be denoted at the top of the file under ##FILTER.) ** "." for positions when data for a call is missing or a call could not be made. These changes to the FILTER field apply to both single-sample and multi-sample VCFs, but not to VCFs lacking any sample information. * In the QUAL field, a Phred-scaled quality score is provided for the assertion made in the ALT column. The score is calculated as $-10 \log_{10} \text{prob}(\text{call in ALT is wrong})$. <ul style="list-style-type: none"> ** In rows where the ALT column contains '.' (i.e. no variant was called), the column contains $-10 \log_{10} \text{prob}(\text{variant})$. ** In rows where the ALT does not contain '.' (i.e. a variant call), the column contains $-10 \log_{10} \text{prob}(\text{no variant})$. ** A missing value is specified as "." * The PA field contains the Pnotref value. Note that the QUAL scale is reversed relative to Pnotref when ALT is "."; that is, when a position is in the

	reference. However, in one direction or the other, it will scale logarithmically with Pnotref. This does mean that it will be closer to Qcall (or "GQ") in cases where there isn't "homozygous vs. heterozygous" call ambiguity. However, when the ambiguity is present, it will diverge.
.bed, .txt, etc.	The target region file (.bed or manifest) for the assembly, if one was specified.
.templateInfo	Contains general information for each contig in the assembly.
.enrichment_Summary.txt	Contains the textual information for the Project > Show coverage of target regions option in SeqMan Pro and SeqMan Ultra.
.sqd	This file is only created when the <i>.assembly</i> is first opened in SeqMan Pro or SeqMan Ultra. It contains saved display specific information such as SNP filtering criteria. Double-clicking on this file will open the <i>.assembly</i> package in SeqMan Pro or SeqMan Ultra.
-Transcriptome table folder containing the file .table.txt	This folder and its file, showing the putative gene identity for each transcript, are created for the de novo transcriptome RNA-seq workflow only.
There is normally no reason to open the following files.	
-0.assemblyInfo	Contains information about assembly parameters which can be used for combining multiple assemblies. This file is not present for SeqMan NGen assemblies made prior to version 14.0. In 14.0 and later, it is present in templated miRNA, ChIP-Seq, and RNA-Seq workflows.
[project name]Transcriptome.table.txt	This file is present in RNA-seq workflows that used a <i>.Transcriptome</i> package as a template. It is equivalent to the <i>.table.txt</i> file in the Transcriptome table folder of the <i>.Transcriptome</i> package.
.auxPair	(internal use only)
.bam	The BAM formatted alignment file.
.bam.bai	The BAM index file.
.capture.userSNP.vcf	(internal use only)
.combined.snpExt	(internal use only)
.coverage	Contains information at each position along the contig where the coverage changes.
.coverage2	Contains information for the maximum coverage of 100 base pair intervals across the contig.
.coverage4	Contains information for the maximum coverage of 10,000 base pair intervals across the contig.
.coverage.missingSNP	Contains information about positions in dbSNP that had coverage and were called the reference base in the assembly.

.exomeCapture-features	(internal use only)
.info	Contains files used by SeqMan Pro and SeqMan Ultra.
.midinfo	(internal use only)
missing.fas	A fasta file of reads with no mers matching the reference.
missing.fas.qual	A base quality file of reads with no mers matching the reference.
.nocoverage.missingSNP	Contains information about positions in dbSNP that had no coverage in the assembly.
outofOrder.txt	A text file of sequence reads not included in the final assembly due to excessive trimming during the alignment phase.
.pair	(internal use only)
.pairDist	Contains information about the position and distance between paired end reads.
pairSpecifiers.txt	(internal use only)
poor.fas	A fasta file of reads rejected at the layout phase due to match scores below the threshold.
poor.fas.qual	A base quality file of reads rejected at the layout phase due to match scores below the threshold.
.quant	Reprises information in the .coverage4 .coverage2 and /or .coverage files.
.region_capture.bed	(internal use only)
report.txt	Contains the textual information for the Project > Report option in SeqMan Pro. See View the Project Report for information about the report contents for XNG and SNG workflows.
.snp	Contains all the information for SNPs called using the “Simple” method.
.snpExt	Contains all the information for SNPs called using either the “Diploid” or “Haploid” method.
SNPs.log	An optional text form of the .snpExt table that contains information on how each was calculated. If you encounter a problem, this file is useful for DNASTAR Support to help you with trouble-shooting.
.splitExt	(internal use only)
.template-comment	Contains the comment information for that contig.
.template-features	Contains the feature information for that contig.
.template-features2	(internal use only)

.template.fof	A file-of-files containing the path and file names of the reference sequences.
.template-gapped-seq	A .seq file of the template containing gaps.
.template-gaps	A binary file of the template gap information.
.template-seq	A .seq file of the template without gaps.
unaligned.fas.qual	A base quality file of reads rejected at the alignment phase.

Contents of the -reports folder

The *-reports* folder is part of the XNG [.assembly package](#).

In the table below (and in the sentence above), it should be understood that the project name precedes any hyphen (-) or period (.) used at the beginning of file and folder names.

File Suffix or Extension	Description
-zinternal	See Contents of the -zinternal folder for details.
-enrichment_Summary.txt	(internal use only)
-perAssemblyResults.txt	<p>The table consists of one row of overview information and assembly statistics per assembly. Information from subsequent assemblies is appended to the existing table.</p> <p>Here are some statistics of interest from perAssemblyResults.txt:</p> <ul style="list-style-type: none"> * NumSeqs – The total number of sequences. * Export_Split_Cnt – The number of exported split reads. * Export_Aligned – The number of exported reads. <p>Note: Since split reads are counted in each location where they align, it is possible for Export_Aligned to exceed NumSeqs.</p>
-perTemplateResults.txt	Overview information and assembly statistics per contig.
-projectReport.txt	Overview information of the assembly. The same report can be viewed within SeqMan Pro using the Project > Report menu command.
-unassembled.fastq	<p>The unassembled reads from the assembly in Fastq format. If production of this file is not specified in the script, three files are created instead:</p> <ul style="list-style-type: none"> * missing.fastq – Unassembled reads with no hits to any template. * poor.fastq – Unassembled reads with scores too low to include in the layout. * unaligned.fastq – Unassembled reads included in the layout, but rejected by the aligner.

Contents of the -zinternal folder

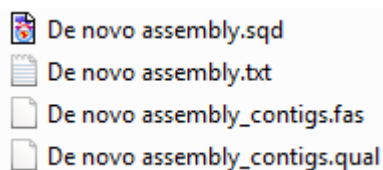
The *-zinternal* folder is located in the [-reports folder](#), which in turn is part of the XNG [.assembly package](#).

In the table below (and the nomenclature used in the sentence above), it should be understood that the project name precedes any hyphen (-) or period (.) used at the beginning of file and folder names.

File Suffix or Extension	Description
info	The files -[templateID].insertion2 and -[templateID].sV_Edges.txt both contain structural variation information used in SeqMan Pro reports.
bamToSQD.script	For converting the assembly to .sqd format.
pairScheme.info	(internal use only)
results.txt	<p>Here are some statistics of interest from the results.txt file:</p> <p>* templateCoverage% – The percentage of the length of the template that is covered by one or more reads.</p> <p>* medianCoverage – The average depth for regions that with nonzero coverage.</p>
The following files instruct SeqMan NGen to convert unassembled reads into a separate SQD project:	
batchUnassembled	A UNIX executable file for <i>de novo</i> assembly of unassembled reads with v. and vi.
batchUnassembled.table.txt	A table of values for running an SNG assembly of the missing .fas reads against the template sequence.
batchUnassembled.template.script	The SNG script containing variables that are specified by the batchMissing.table.txt.

De novo workflow output

De novo workflows output a results folder containing the following files:



File Suffix	Description
.sqd	The main assembly output. To view and analyze the assembly, open this file with SeqMan Pro or SeqMan Ultra.
.txt	(internal use only)
-contigs.fas	Created when contigs are saved in FASTA format.
-contigs.qual	Created when contigs are saved in FASTA format. The values in the file are the sum of the base qualities at each position in the contig, up to a maximum of 90.

✿ **Note for Windows users:** To open text reports with the correct formatting displayed, we recommend using Wordpad, Notepad++, or Microsoft Excel®, and not the default Windows text editor, Notepad.

RNA-Seq reference-guided workflow output

If you are following a reference-guided RNA-Seq [workflow](#), output results are saved in an .assembly package folder labeled with the user-specified [project name](#) and the suffix **_RNA-Seq**. This folder contains the following files:

Subfolder Name	File Name	Description
	[project name].astar	An ArrayStar project file containing data from all of the samples.
	[project name].script	The main script for a reference-guided assembly project.
	[project name].table.txt	The table containing project descriptions and corresponding files; used by the main script.
	[project name].template.script	The reference-guided workflow script, which is applied to all projects listed in .table.txt .
	[project name].astar.script	The script file executed by ArrayStar when the Launch in ArrayStar button is pressed in the Project Report.
	query'n'.script	The script containing query sequences. This script is generated only if a multiple project assembly includes groups that contain non-equal numbers of files.
	_tableScript.template.script	The script template, executed by the XNG assembler in order to run the QNG assembler. This script is not present before the run and is generated by XNG.
	[project name]-qng.script	After SeqMan NGen finishes running the XNG part of the assembly, QNG assembly is performed using this script. In multi-sample projects, this file references .table.txt .
	[project name].Results.txt	A tab-separated table showing the assembly results.
Sample'n'.assembly	sample'n'-0.isoforms-features	Each .assembly project contains a set of files relating to a single contig.
	sample'n'-0.genes-features	
	sample'n'-0.bam	
	sample'n'-0.coverage	
Sample'n'-Reports	'n'-0.report.txt	A summary of the run; also shown in the Project Report.

RNA-Seq de novo transcriptome workflow output

If you are following the *de novo* transcriptome RNA-Seq [workflow](#), output results are saved in a folder called *[project name] De Novo Transcriptome Assembly*. This folder contains the following subfolders and files:

Subfolder	File/Folder Name	Description
	[project name]_rnaAssemble.script	Input script used to create the assembly results. This file can be opened in SeqMan Pro in order to examine isoforms using the Feature Table .
Assemblies	[project name]_novel_transcripts.sqd	SQD assembly of all contigs that did not have a database match.
	[project name]_unassembled.fastq	Multi-sequence FASTQ file with all unclustered and unassembled sequences.
	sub_0 (folder)k	Folder containing sub-folders (sub_0, sub_1, etc.) with a separate .sqd document for each final assembly. If available, gene and organism names are used to create the file names.
Intermediate Assembly Results	cluster (folder)	Intermediate results are deleted by default at the end of the assembly, but can be retained by designing the input script such that the assembleTemplate command's deleteIntermediates parameter is set to false .
	combine (folder)	
	intermediateFiles (folder)	
Reports	[project name].AllTranscripts.SearchResults.txt	<p>Excel file containing summary information for each of the final assembled contigs. The table automatically opens for viewing when you open a .Transcriptome package in SeqMan Pro. The table, known in SeqMan Pro as the “All Transcripts” table, contains the following columns:</p> <ul style="list-style-type: none"> * Assembly ID – Name assigned to the assembled sequence, using the criteria specified in the wizard. * Gene name, Custom column #1* – Best matching gene meeting criteria defined in the wizard. * Organism name, Custom column #2* – Organism from which the best matching gene came.

* **Accession number, Custom column #3*** – Accession number of the best match.

* **Description, Custom column #4*** – Description of the best match.

***Custom columns:** The four “custom columns,” above, use default names (e.g., **Gene name**, **Organism name**) if one of the default RefSeq databases was used in the SeqMan NGen assembly. However, if you used a custom GREP expression or a custom database that did not include these fields, these columns may have different names or be absent from the table.

* **Database** – Database (e.g. [RefSeq](#), Custom, etc.) from which the best matching gene came.

* **Transcript length** – Length of the assembled sequence, in bases.

* **Transcript start** – Position in the assembled sequence where the match begins.

* **Transcript end *** – Position in the assembled sequence where the match ends.

* **% Transcript match** – Length of the matching segment in the transcript x 100, divided by the total length of the transcript.

* **Gene length** – Length of the database entry, in bases.

* **% of Full length** – Length of the assembled sequence x 100, divided by the length of the corresponding database entry. Values greater than 100% indicate that the assembled sequence is longer than the database entry.

* **Gene start** – Position in the database entry where the match begins.

		<p>* Gene end – Position in the database entry where the match ends.</p> <p>* % Gene match – Length of the matching segment in the database entry x 100, divided by the total length of the database entry.</p> <p>* % Identity – Total number of identical bases in the matching region x 100, divided by the total number of bases in the matching region.</p> <p>* Bit score – Normalized value calculated from the raw score and expressed in units of “bits,” a common measure in information theory.</p> <p>* eValue – “Expectation value,” an estimate of the probability of obtaining the observed alignment score with two random sequences. Expectation values are less sensitive to length than Bit scores and are therefore generally a better measure of alignment quality.</p> <p>* Assembled reads – Total number of assembled reads for that sequence.</p>
	<p>[project name].AllTranscripts.Table.txt</p>	<p>Excel file containing summary information for each of the final assembled contigs. The table contains the following columns:</p> <p>* Assembly ID – Name assigned to the assembled sequence, using the criteria specified in the wizard.</p> <p>* Type – Type of matching gene (e.g., mRNA, tRNA, rRNA, etc.)</p> <p>* Gene length – Length of the database entry, in bases.</p> <p>* % of Full length – Length of the assembled sequence x 100, divided by the length of the corresponding database entry. Values greater than 100% indicate that the assembled sequence is longer than the database entry.</p>

		<p>* Assembled reads – Total number of assembled reads for that sequence.</p> <p>* Depth – Average depth of coverage.</p>
Transcripts	[project name]_identified_transcripts.fas	Multi-sequence <i>.fasta</i> file containing the consensus sequences from all the assembled contigs that had a database match. Header lines for each entry contain the name and sequence length.
	[project name]_novel_transcripts.fas	Multi-sequence <i>.fasta</i> file containing the consensus sequences from all the assembled contigs that did not have a database match. Header lines for each entry contain the name and sequence length.

Appendix

The Appendix contains the following topics:

- [Non-English keyboards](#)
- [SeqMan NGen calculations](#)
- [Access and understand output files](#)
- [Turn off usage logging](#)
- [Installed Lasergene file locations](#)
- [Research references](#)
- [Run SeqMan NGen through the command line](#)

SeqMan NGen calculations

The following topics describe how SeqMan NGen handles various situations or makes calculations:

- [Calculation of match percentage](#)
- [Detection of structural variations](#)
- [Handling of repeats](#)
- [Handling of sex chromosomes](#)
- [How mer tags are chosen](#)

Calculation of “match percentage”

By default, SeqMan NGen uses a local match percentage which requires that the match percentage threshold be met in each overlapping window of 50 bases. The size of this window can be adjusted by specifying a different value for the **match window** parameter. To access this parameter from the [Preassembly Options](#) screen for non-long-read workflows, press the **Advanced Options** button, then the [Alignment tab](#).

An example containing a repeated region follows.

A genome fragment has repeated regions labeled A and A', and two unique regions labeled B and C.



When the fragment is sequenced, one of the sequences contains parts of regions A and B, and another contains parts of regions A' and C:

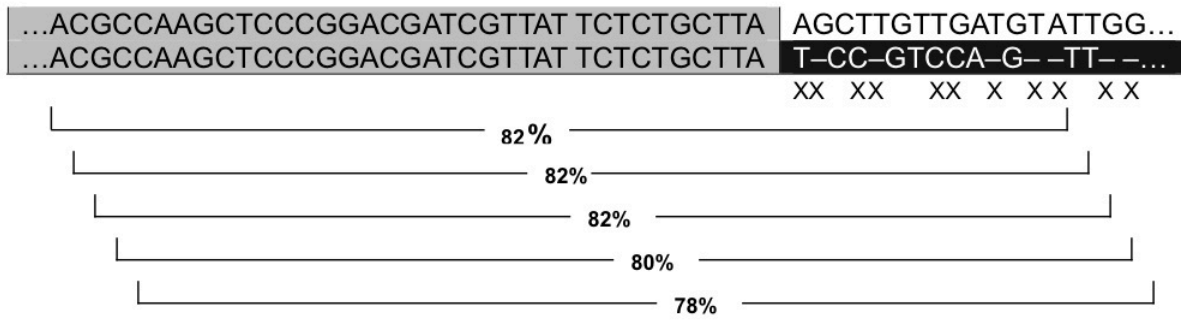
Sequence 1: 

Sequence 2: 

In this example, a **minimum match percentage** of 80% is used. When the two sequences are aligned, the 400 bases in the overlapping A and A' regions match 100%. The 200 bases in the overlapping B and C regions match 42%. Over the entire alignment, 484 out of 600 bases match, yielding a global match percentage of 81%.

However, SeqMan NGen checks the match percentage for every alignment of 50 bases. The alignment below shows the last 36 overlapping bases of A and A' and the first 18 overlapping bases of B and C. Each mismatch in the overlap is marked by an X below the alignment. In the first 50 bases shown, there are 41 matches, and the match percentage is 82%. This is above the threshold of 80%, so the match percentage of the next 50 bases is checked and is also found to be 82%.

Each fifty bases are checked along the overlap as long as the match percentage is at or above the threshold. In this case, the alignment fails once it gets far enough into the overlap of the unique regions, B and C, that the match percentage drops to 78%. The sequences will not be assembled together into a contig, which is correct for this data set.



Detection of structural variations

In addition to SNPs and small insertions and deletions, genetic variation can also involve large scale rearrangements. These rearrangements may include large insertions and deletions, inversions, and translocations — collectively known as *structural variations* (SV's). To view a tabular report with structural variation findings for an assembly, open the assembly in SeqMan Pro and select **Contig > Structural Variation Report**.

The rest of this topic describes how these structural variations were detected during assembly with SeqMan NGen.

During a templated data assembly, SeqMan NGen automatically detects insertions or deletions (*indels*) greater than 10 bp based on a combination of two data types:

Coverage – read depth can be suggestive of larger indels, including duplications or the collapse of a duplication.

Split reads – reads spanning a deletion in the new genome relative to the reference genome can be “split” into two segments based on matches to discontinuous regions on the reference. For example, the following split read alignment indicates there is a 35bp deletion in the new genome:

Ref: AGGCTGACCTCCTGGGCCTAAGACACTGAGTGCCCCAATATGACTCGACTAGCA

Read: AGGCTGACCTCGACTAGCA

Ref: AGGCTGACCTCCTGGGCCTAAGACACTGAGTGCCCCAATATGACTCGACTAGCA

Split: AGGCTGACCTC GACTAGCA

The SeqMan NGen algorithm requires that four criteria be met for splitting a read:

1. At least 20 bases on each “half” of the split must match the reference. This means that reads must be at least 40bp long, though in practice they should be > 60 bp.
2. The first mer match must be within 10 bases of the start of the read, and the final mer match must be within 10 bases of the end of the read. This increases the likelihood that the entire read will align after splitting.
3. The distance between the two closest mers on either side of the split must be within 20% of the total read length. For example, in a 100 base read where bases 5-30 make up the mer match on the 5’ “half” of the read, then the first mer match on the 3’ half must start between bases 31 and 50 ($30 + (100 \times 0.2) = 50$) of the read. This relatively simple requirement allows for SNPs or sequencing errors near the actual split to be tolerated and resolved during alignment.

4. The two “halves” must be aligned in the same orientation.

In practice, two copies of the read are given to the aligner: one seeded with the 5' mer match and the other with the 3' mer match. The aligner then extends the alignment on both sides of each copy, and then trims each copy to maximize the final alignment score. It is the final trimmed internal position for each copy that is reported in SeqMan Pro's Structural Variation Report.

✿ **Note:** Because of the trimming and the flexibility in the location where the nearest internal mer match must begin (criteria #3 above), it is possible that base substitutions are present in the split region that will not be displayed in SeqMan Pro's Alignment View. Thus, while split reads have far greater resolution power than coverage, the breakpoints identified must be considered as provisional.

Handling of repeats

Repeat handling parameters compute a threshold for deciding the number of identical subsequences of bases (mers) used to indicate a putative repeat. Mers that are common to two or more fragment reads are aligned to determine the overall layout of reads. For additional information, see [How mer tags are chosen](#).

Repeat handling is controlled via the **Place repeat reads** drop-down menu in the [Alignment tab](#). This tab is accessed from the [Assembly Options](#) screen by pressing the **Advanced Options** button, then the **Alignment** tab.

The repeat threshold can also be computed by multiplying the **Match repeat percent** parameter value in the Advanced Assembly Options dialog by the **Expected coverage**. Any mer that occurs more frequently than the computed threshold is not considered for use as a mer tag in determining overlaps. Coverage can be determined in two ways:

- By using the length of the genome/fragment being sequenced, as specified by the **Expected genome length** parameter in the Assembly Options dialog. If this option is used, SeqMan NGen calculates expected coverage by dividing the **Expected genome length** value by the total length of all sequences in the project.
- By using a fixed number for expected coverage, as specified by the **Expected coverage** parameter in the [Assembly Options](#) dialog.

Handling of sex chromosomes

When using the SeqMan NGen wizard, certain workflows allow you to specify the subject's **Gender** in the [Assembly Options](#) screen.

✿ **Note for command-line users:** When using the [command-line version of XNG](#), the sex of the subject is instead specified using the [assembleTemplate](#) command by setting the **query** parameter value to **male**, **female**, or **unknown**.

SeqMan NGen treats all non-template package chromosomes as diploid, unless otherwise specified. However, SeqMan NGen's [XNG assembler](#) does recognize sex chromosomes in DNASTAR genome template packages. Since some regions of the X and Y chromosomes are homologous, read placement may be more accurate in females when reads are not falsely assigned to the Y.

Special handling of sex chromosomes occurs in two circumstances:

During the placement/layout of reads:

All samples in an assembly must be of the same sex; different sexes cannot be specified for individual MID samples. Chromosomes are recognized by GTP shortname (X Y W Z). One of the names must match exactly. Humans and many other animals have X/Y, while the chicken (*Gallus gallus*) has a W/Z. For some GTPs such as the cow (*Bos taurus*), Y exists biologically, but has not been sequenced and is not provided in the genome template package.

Sex influences read placement as follows:

- If sex is left **Unknown**, all chromosomes will be available.
- If sex is set to **Female**, no reads are placed against the Y chromosome.
- If sex is set to **Male**, the haploid variant caller is used for both X and Y chromosomes; no reads are placed against the W chromosome.

+When calling SNPs and using the Bayesian SNP caller: +

Bayesian SNP calling is also modified by taking the sex chromosome into account. The XNG script normally controls whether calling is done as diploid or haploid. Whenever haploid is specified as a SNP method all templates are haploid. However even when diploid is chosen:

- chromosomes Y or Z are always haploid

- chromosome X is haploid for males, diploid for female or unknown
- chromosome W is haploid for females, diploid for male or unknown

The mitochondrion is considered diploid, though it would usually be polyploid. If mitochondria are of interest, a separate assembly should be done using the simple SNP method, since the Bayesian caller depends on knowing that there are one or two chromosomes.

Special handling of sex chromosomes can be disabled in an XNG script by setting `noSexChromosomes:true`. Changing the shortname of a chromosome in a GTP will cause it to revert to normal autosomal behavior in all scripts using the GTP.

How “mer tags” are chosen

The SeqMan NGen layout algorithm relies on unique subsequences of bases, or *mers*, which occur in overlapping regions of fragment reads. Mers that are common to two or more fragment reads are aligned to determine the overall layout of reads. Overlapping reads have many mers in common, but only a few mers per overlapping region are needed to identify the overlap. These mers are called *mer tags*. The use of mers to tag fragments and identify overlaps is illustrated in the following figure:

Original DNA Sequence:

CGAATGTCATATGGCAGTACACGGCGTACGTTAGGTTTCTGAGGGATTTTCGAG

Fragment Reads:

1. CGAATGTCATATGGCAGTA
2. TATGGCAGTACACGGCGTACGT
3. GCGGTACGTTAGGTTT
4. TTAGGTTTCTGAGGGATT
5. AGGTTTCTGAGGGATTTTCGAG

Fragment Read Layout:

1. CGAATGTCATATGGCAGTA
2. TATGGCAGTACACGGCGTACGT
3. GCGGTACGTTAGGTTT
4. TTAGGTTTCTGAGGGATT
5. AGGTTTCTGAGGGATTTTCGAG

As shown in the above figure, a 54bp original DNA sequence is covered by five overlapping fragment reads. The 6-mer tags for each fragment read are underlined. Matching mer tags are aligned to determine the layout of the reads.

The power of using mer tags relies on the ability of SeqMan NGen to choose mers that are most likely to occur only once in the original DNA sequence. It is important to avoid choosing mers that occur in repeated regions since the result may be fragment reads that are incorrectly aligned together.

Three parameters are involved in choosing mer tags: **Match Size**, **Repeat Handling**, and **Match Spacing**. All of these parameters can be adjusted in the multi-tabbed Advanced Assembly Options dialog. To access this dialog from the [Assembly Options](#) screen, press the **Advanced Options** button.

The **Match Size** and **Repeat Handling** parameters help to choose tags that are most likely to be unique in the original DNA sequence. **Match Size** sets the length of the mers. The longer the mer, the higher the probability that it is unique. **Repeat Handling** parameters help to identify which mers are not likely to be unique. If a mer occurs more often than expected in the dataset, the mer may be part of a repeated region.

Match Spacing specifies the preferred distance between mer tags. The smaller the **Match Spacing** parameter value, the more memory and more time the assembly will take. If a fragment read is shorter than the **Match Spacing** value, multiple mer tags are still chosen for the read.

✿ **Note:** During assembly, any given read will only be assigned to one contig, even if it matches the hit criteria for more than one contig. If there is no information linking the read to a specific contig (e.g. a unique SNP or a paired-end constraint), SeqMan NGen will assign the sequence randomly to one of the contigs for which it meets the criteria.

Run SeqMan NGen through the command line

To create a script:

Either:

- Use the SeqMan NGen [wizard](#) to create the script. Rather than pressing the **Assemble** button at the end, simply save the script.
- Create a text file from scratch using the commands and parameters for [XNG](#) or [SNG](#). Save the script with the file extension *.script*.

To run the script using the command line:

1. Open the command shell (Win) or Terminal (Mac).
2. Type one of the following commands:
 - `sng [path to script]`
 - `xng [path to script]`



Note: If SeqMan NGen does not run immediately, it is possible that the application has not been added to your environment variable PATH. After installing SeqMan NGen on a computer for the first time, command-line users may want to restart the computer before using the application. Restarting will add the application to your environment variable PATH. If you are unable to restart your computer, the following two options are available:

- * Use the SeqMan NGen wizard to run the script.
- * At the command-line, specify the path to the application, followed by the path to the script.

XNG, SNG, and QNG assemblers

SeqMan NGen uses three powerful assemblers: XNG, SNG and QNG.

The XNG assembler:

The XNG assembler (patent pending) is used for all reference-guided assemblies. This assembler features an algorithm for fast, accurate assembly of extremely large genomes and is capable of assembling data sets of any size, given [sufficient disk resources and modest RAM requirements](#). The XNG assembler uses multiple cores, but the exact number varies over the course of the assembly.

The primary output is an assembly folder (Win) or package (Mac) containing a binary *.bam* file for each reference sequence to which reads could be mapped. The assembly folder also contains accessory files. Note that BAM files cannot be edited. Reference-guided workflows use the XNG assembler. The reference-guided assembly with gap closure uses both the XNG and SNG assemblers, but the output files are most similar to XNG outputs.

SNG/SMNG assembler:

The SNG assembler generates finished assemblies in SeqMan Pro / SeqMan Ultra (.sqd) or BAM (.assembly) formats. The *.sqd* files are editable within SeqMan Pro, but the number of data reads is limited to 10 million or fewer. BAM files of any size can be created, but may not be edited. The SNG/SMNG assembler uses one core during assembly. SNG/SMNG workflows include:

- All *de novo* assemblies. SNG/SMNG generates finished assemblies in editable SeqMan Pro / SeqMan Ultra format.
- Reference-guided assemblies for small genome (less than 30MB) reconstruction projects where editing is required. In order to perform editing, an output format of SeqMan Pro / SeqMan Ultra must be selected.

QNG assembler:

In this type of assembly, XNG is used to make a first-pass assembly. Then, a version of ArrayStar's QSeq algorithm is applied to the results to create the finished assembly.

For a list of the output files for a given workflow, see [Access and Understand Output Files](#) and its subtopics. In addition, [View the Project Report](#) describes how to access the most commonly-viewed output file.

XNG commands

The following commands can be used in the script for an XNG assembly. Click on the name of a command in the table below (shown in alphabetical order) to see a description and example, and to see the parameters associated with that command. All XNG commands and parameters are assumed to be optional unless the description states that it is required.

assembleTemplate	exportVCF	pause
computeSNP	extractPairs	quit
createGenomeTemplate	include	removeDuplicateSeqs
diskPath	loadAssembly	runScript
dumpConsensus	loadBAM	set
dumpSNP	mergelonTorrentShortReads	setDefaultDirectory
execute	message	setMachineMemory
exportSplits	pairFilePattern	setParam

assembleTemplate



Note: All parameters are assumed to be optional unless the description is prefaced by “required.”

assembleTemplate is a required command, and Initiates the assembly of the loaded sequences using the specified template as a reference.

Example:

XNG script used in the “clustering” step of the de novo transcriptome RNA-seq [workflow](#):

```
merSize: 25
minNewClusterSize: 5
minSingleMergeClusterSize: 7
minMultiMergeClusterSize: 7
minMultiMergeIgnoreFactor: (currently not used by default)
minClusterSizeToOutput: 100
```

Parameter	Description	Allowed values (default value underlined)
alignmentCutoff	Used in the “clustering” step of the de novo transcriptome RNA-seq workflow .	[number] Default = <u>20</u>
assemble	Specifies whether to use the part of the query that matches the contaminant sequence(s), the part that doesn’t match, or both.	[<u>matchContaminant</u> / <u>noMatchContaminant</u> / <u>all</u>]
assemblyInfo	Contains information about the assembly.	[text string]
assemblyInfoAlt	Contains pairs of keys and values which will be written to the -0.assemblyInfo file.	
autoTrim	Specifies whether mismatching ends of reads should automatically be trimmed.	[<u>true</u> / false]
autoTrim	Specifies whether mismatching ends of reads should automatically be trimmed.	[<u>true</u> / false]
boneyardAssembly	Specifies whether sequences not used in the original or incremental XNG assemblies should be added to the assembly project by the SNG assembler. This command pertains only to reference-guided assemblies with gap closure. By default, during this type of assembly, the XNG	[<u>true</u> / false]

	assembler first finds structural variations (SVs) then splits the contig after each SV. Elements of this process can be modified using this command. (Note: “Boneyard” is a term for sequences that were not assigned to any contig).	
combineDuplicateSeqs	Specifies whether the duplicate reads will be clustered.	[true / false]
contaminant	<p>Use of this parameter partitions the query data by running an additional mer-match (layout) against the specified contaminant sequence(s). A full assembly is then run using the part of the query that either matches or does not match the contaminant sequence(s). This parameter can be used for removing reads originating from an organism(s) that may have also been present in the query data set (e.g., reads from human DNA present in a metagenomic sample from the human gut).</p> <p>file: [directory/filename enclosed in quotes] the file with contaminant sequences.</p> <p>assembleContam: [matchContam/noMatchContam/all]</p> <p>merLayoutMin: [number]</p> <p>unassembled: [directory/filename enclosed in quotes] the file containing no contaminant reads.</p>	[directory/filename enclosed in quotes]
dbSNPTable	(Intended for internal use only).	[directory/filename enclosed in quotes]
delayAlignInserts	Use of this flag turns the delay reads that cause inserts on or off. ‘True’ means that gap causing reads will be delayed. Reads will be added such that reads causing the lowest number of inserts (length of inserts is not considered) will be added before those causing more inserts.	[true / false] Defaults: true for ‘Illumina’ and ‘Solexa’ technologies, false for ‘Other’ read technologies.
deleteIntermediates	Specifies whether intermediate files are saved or deleted. These files can be large with large-scale projects.	[true / false / all / notTemplate]
directoryMer	Specifies the path and directory where both the template and query data mer files will be stored. Alternatively, separate directories for the template and query mer files can be specified using the parameters below. If no directory is specified, the mer file will be created in the directory containing the sequence data.	[directory/filename enclosed in quotes]
directoryQueryMer	(required) Specifies the path and directory where the query mer file will be stored.	[directory/filename enclosed in quotes]

directoryTemplateMer	(required) Specifies the path and directory where the template mer file will be stored.	[directory/ filename end in quotes]
filterDeepLayout	(optional) Specifies that XNG remove superfluous sequences in areas of deep coverage. Wizard equivalent: Using 'true' is equivalent to selecting the Limit all deep coverage regions radio button from the Alignment tab . This tab is accessed from the Assembly Options screen by pressing the Advanced Options button.	[true / false Set to 'false', default, exce projects invo miRNA or microbial genomes, w it is set to 'tr
filterDeepLayoutOrganelle	(optional) Specifies that XNG remove superfluous sequences in areas of deep coverage. Wizard equivalent: Using 'true' is equivalent to selecting Advanced Assembly Options > Alignment tab > Only limit deep coverage regions for Mitochondria and Chloroplasts radio button	[true / false Set to 'false', default, exce projects invo a mitochond chloroplast template (i.e those with a name of 'MT' or 'CHL' or 'chloro'), wh is set to 'true
forceFullForwardAlign	Start the alignment at the 5' end of the sequence.	[true / false
forceMake	Specifies whether new intermediate mer files will be created. A value of 'false' means that existing valid intermediate files will be used.	[true / false query / hit / layout]
format	Specifies the format of the alignment output file. If 'none' is entered, the assembly is run to include the alignment phase, but no alignment output is generated. This parameter can be used to remove reads from a contaminant source.	[BAM / SQD NONE / NONE_align Aux_align]
gap5Prime	Put the gap on the 5' side of the sequence.	[true / false
gapPenalty	The penalty for opening or extending a gap during an alignment. This penalty is deducted from the pairwise score used to calculate match percentage. A high gap penalty suppresses gapping, while a low value promotes gapping.	[number] Default = <u>30</u> most workflo <u>50</u> for the de transcriptom

		RNA-seq workflow .
gapExtensionPenalty	Used in the “clustering” step of the de novo transcriptome RNA-seq workflow .	[number] Default = <u>5</u>
geneticCode	This parameter specifies the genetic code to use with a reference sequence.	[filepath/star Lasergene genetic code name]
hits	(required) Specifies the path and name of the hit file. Incomplete paths will be appended to the default directory.	[directory/ filename end in quotes]
increaseRunGapPen	This parameter is a flag to increase the gap open penalty in HP runs.	[true / <u>false</u>]
layout	(required) Specifies the path and name of the layout file. Incomplete paths will be appended to the default directory.	[directory/ filename end in quotes]
layoutAlign	Specifies that a pairwise alignment should be performed at the payout phase in order to pick the best position for a given read.	[true / <u>false</u>]
layoutMaxTemplateGap	The maximal number of gaps introduced into the alignment used during layout.	[number]
layoutRSRange	The maximal Register Shift difference used while building the layout.	[number]
layoutType	Specifies how reads are to be laid out.	[unique / <u>on</u> multiple / multipleAll]
matchScore	The score for a base match during an alignment. This score contributes to the pairwise score used to calculate match percentage. Increasing the matchScore value allows for longer or more frequent gaps, thus forcing bases that match to be assembled together.	[number] Default = <u>10</u>
MaxGap	The theoretical maximum length of a gap that could be inserted. In practice, the maximum gap size will usually be about half of this value.	[number from 0-99] Default = <u>6</u> for most workflo <u>30</u> for the de transcriptom RNA-seq workflow
maxMergeSize	When linking clusters into a scaffold, only link them together if the overall	

	number of reads in the scaffold would not exceed this threshold. Used in the “clustering” step of the de novo transcriptome RNA-seq workflow .	
maxNCnt	(optional) This parameter removes sequential reads of the IUPAC ambiguity code ‘N’ that are greater than or equal to the number specified. Use of this parameter may help in assemblies whose reads contain large clusters of spurious N’s.	[integer]
maxSecondaryTrimLength	During alignment, a read can be trimmed from both ends. This parameter defines the longest allowable length for the smaller of the two trimmed ends.	[number]
maxSeqs	Specifies the maximum number of query sequences to add to an assembly. Use of this command can speed up assembly.	[number]
merCntThresh	Minimum number of mers needed in order to be recorded in the mer file.	[number]
merLayoutMin	Specifies the minimum length (in bases) of at least one stretch of matching mers used to identify matches between the reference and query data. The minimum value is equal to the mer. The maximum value is the read length, which would require the entire read be an exact match. For example, with a merSize of 19 and a merLayoutMin of 21, at least one stretch of three consecutive mers in a read would have to match for the read in order to be included in the layout.	[number from 11-1000] Default = <u>25</u>
merMinimizer	(Intended for internal use only)	[number]
merSize, merLength or matchSize	(required) Specifies the length (in bases) of mers used to identify matches between the reference and query data.	[number]
merSkip	(Intended for internal use only) Specifies the number of positions to ignore or “skip” when creating the template mer file. Normally, mers are only skipped in the query (see merSkipQuery , below). The first and last mer of every read are always included. Increasing the value reduces the size of the intermediate files as well as the overall assembly time. However, larger values can also reduce the number of reads included in the assembly, especially with short read data. 0 = do not skip 2 = skip every second base 3 = skip every third base etc.	[number] Default = <u>0</u>
merSkipQuery	Specifies the number of positions to ignore or “skip” when creating the query mer file. The first and last mer of every read are always included. Increasing the value reduces the size of the intermediate files as well as the overall assembly time. However, larger values can also reduce the number of reads included in the assembly, especially with short read	[number] Default = <u>0</u>

	<p>data.</p> <p>0 = do not skip 2 = skip every second base 3 = skip every third base etc.</p>	
method	<p>Defines how to handle splits in the assembly:</p> <ul style="list-style-type: none"> * normal – normal assembly method * splitOnly – only reads which have been split will be included in the assembly * noSplit – no reads will be split 	[normal/split/noSplit]
minAlignedLength	<p>Specifies the minimum number of bases that must align after trimming for a read to be included in the assembly.</p>	<p>[number from 11-999]</p> <p>Default = <u>25</u> most workflow <u>50</u> for the de transcriptome RNA-seq workflow.</p>
minClusterSizeToOutput	<p>Threshold for the number of reads that a cluster must contain in order for the cluster to be passed along to SNG for assembly in the next step of the program. Used in the “clustering” step of the de novo transcriptome RNA-seq workflow.</p> <p>Note that this command is present only for the clusterParam block of the rnaAssemble command.</p>	[number]
minMatchPercent	<p>The minimum percentage of matches in an overlap required to join two sequences in the same contig.</p>	<p>[number]</p> <p>Default = <u>93</u> most workflow <u>60</u> for the de transcriptome RNA-seq workflow.</p>
minMultiMergeClusterSize	<p>When two or more clusters overlap the same k-mer, the minimum number of reads (depth) required at that k-mer for a cluster to consider that cluster significant.</p> <p>If three or more clusters exceed this threshold, the k-mer is considered</p>	[number]

	<p>“noisy” and a potential false join, and will not be merged. This is reported as a “multi-cluster link that was not merged”.</p> <p>If two significant clusters overlap and have similar enough depth, the clusters are considered linked and are scaffolded together. Otherwise, if only one cluster is significant, all reads at that k-mer which have no assigned cluster are merged directly into it as described for the minSingleMergeClusterSize option. This parameter is used in the “clustering” step of the de novo transcriptome RNA-seq workflow.</p> <p>Note that this command is present only for the clusterParam block of the rnaAssemble command.</p>	
minMultiMergeIgnoreFactor	<p>When two or more clusters overlap the same k-mer and may be linked, they must be within this ratio of one other. Used in the “clustering” step of the de novo transcriptome RNA-seq workflow.</p> <p>Note that this command is present only for the clusterParam block of the rnaAssemble command.</p>	[number]
minSeqsPerTemplate	Minimum number of sequences sufficient to build the layout or alignment.	[number]
minSingleMergeClusterSize	<p>The minimum number of reads (depth) matching an existing cluster at a single k-mer required to extend that cluster by immediately adding all new reads for that k-mer to the cluster. Used in the “clustering” step of the de novo transcriptome RNA-seq workflow.</p> <p>Note that this command is present only for the clusterParam block of the rnaAssemble command.</p>	[number]
minNewClusterSize	<p>Minimum number of matching reads at a single k-mer (i.e., “depth”) required to create a new cluster. Used in the “clustering” step of the de novo transcriptome RNA-seq workflow.</p> <p>Note that this command is present only for the clusterParam block of the rnaAssemble command.</p>	[number]
mismatchPenalty	The penalty for a base mismatch during an alignment. This penalty is deducted from the pairwise score used to calculate match percentage.	[number] Default = 20
noSexChromosomes	Disables special handling of sex chromosomes.	[true / false]
noSVPairSort	Specifies whether to turn off the calculation of pairs for structural variations. This may potentially reduce XNG assembly time.	[true / <u>false</u>]
onePackage	Specifies whether an assembly containing multiple reference sequences should be bundled into a single .assembly package. If ‘false’ is entered,	[<u>true</u> / false]

	one .assembly package is created per contig.	
openInSeqman	(optional) Specifies whether the completed assembly should immediately be launched in SeqMan.	[true / <u>false</u>]
output	(required) Specifies the path and directory of the output files. Incomplete paths are appended to the default directory.	[directory/ filename enclosed in quotes]
pairDist	(Intended for internal use only)	[true/ <u>false</u>]
pickTemplate	Defines the number of templates from which to choose, and finds the template that is the best match for the input sequence.	[number]
placeHit	(Intended for internal use only)	[true / <u>false</u>]
probe	(Intended for internal use only)	[number]
query	<p>(required) Specifies the directory and file name(s) of the query data to be assembled. A folder with one or data files can also be used in place of individual file names.</p> <p>Properties for query:</p> <p>file: [directory/filename enclosed in quotes] Specifies the directory and file/folder.</p> <p>isPair: [true/false] Specifies whether the query files contain paired end data.</p> <p>minDist: [number] (required if isPair is 'true') Specifies the minimum expected distance in bases between paired end reads. Default is <u>0</u>.</p> <p>maxDist: [number] (required if isPair is 'true') Specifies the maximum expected distance in bases between paired end reads. Defaults are <u>750</u> for Illumina; <u>4500</u> for 454 and Sanger, <u>7500</u> for Other, and <u>user-defined</u> for Ion Torrent</p> <p>seqTech: [unknown IonTorrent IlluminaLongReads 454 PacBio normalScore Other] Specifies the offset to be used when converting compressed quality scores into numerical values. These are the offsets used for the technology specified:</p>	[directory/ filename enclosed in quotes]

Data Type	Value	Offset
IonTorrent	IonTorrent	33
Illumina	IlluminaLongReads	33
Roche 454	454	33
Other types	normalScore	33

Note 1: For 454, quality scores for homopolymeric runs of ≥ 2 are oriented from 5' to 3' on the top strand.

Note 2: If possible, the data type of unknown data is determined automatically based on the first data file.

pairTech : [unknown|LucigenRsaI|LucigenBfaI|RsaI|BfaI|Custom]

pairLinker: [string]

groupName: [string] The name of a group this file belongs to. Used for running multiple samples in one file.

sex: [unknown|female|male]

trim: [true / false] Specifies whether vector trimming needs to be applied to the reads.

sngTrim: contains parameters for fast vector trimming (See the SNG command [trimVector](#))

scan: [true / false] Specifies whether reads needs to be scanned for contaminants

contaminantScan: Contains the assembleTemplate command with contaminant file used as a template and parameters: directoryTemplateMer, hits, layout, output, unassembled, results, format, mersize, ignorePolyMers and deleteIntermediates. The format parameter has valuenone_ALIGN.

Example:

```
query: {{file: "/data/home/proj/Illumina_s_5_1.txt"}
      {file: "/data/home/proj/Illumina_s_5_2.txt "}
isPair: true
minDist: 400
maxDist: 700
```

	seqTech: Illumina}	
recordSplitsOnly	Functional only when used in the same program as splitTemplateContigs or recordStructVariations (both described below). Specifies whether or not to turn off contig splitting while still recording SVs for later inclusion in the Structural Variation Report.	[true / <u>false</u>]
recordStructVariations	Specifies under which circumstances structural variations (SVs) should be calculated and recorded. 0 false = Don't calculate SVs 1 true = Calculate SVs at zero coverage 2 = Calculate SVs at insertions and deletions 3 = Calculate SVs at zero coverage and at insertions	[integer bet 0-3 / true / fa Default = <u>2</u>
removeDuplicateSeqs	Completely removes clonal reads after the alignment phase of assembly. Clonal reads, where the endpoints of both reads in a pair match those in another pair, are usually the result of PCR artifacts. If 'true,' the reads will not be scored, and will not be included in SNP calculations. Marking this parameter to 'true' may substantially increase the time needed for assembly.	[true / false]
removeUniqueInserts	Removes reads that cause an insert which no other read would create. This parameter is only enabled when delayAlignInserts (described under the assembleTemplate command) is true.	[true / false] Defaults: <u>true</u> for Illumina and Torrent read technologies <u>false</u> for all o types.
repeatPenaltyScale	Indicates the quality penalty (using the Phred scale) to use for a read which places in two locations identically. Higher repeat counts are further penalized relative to this on a log ₂ scale such that repeats placing in four locations have a double penalty, in eight locations have a triple penalty, and so on. This penalty is applied to a ceiling of Phred score 30 if the other methods are disabled or have a higher score.	[number] Default = <u>8</u>
repeatThreshMax	Specifies the maximum number of occurrences of a mer in the reference sequence(s) for it to be considered repeated. Mers exceeding this number will not be used for identifying matches.	[number from 1-10000] Default = <u>10</u>
repeatThreshMin	Specifies the minimum number of occurrences of a mer in the reference sequence(s) for it to be considered repeated. Mers less than this number will not be used for identifying matches.	[number]
reportFiles	Defines the kind of report file to be generated.	

	<p>perProject: [true / false] Generate a per project report.</p> <p>perTemplate: [true / false] Generate a per template report.</p> <p>removeInteral: [true / false] Remove intermediate reports.</p>	
repeatmermax	Threshold number of occurrences in a data set for a mer to be considered “repeated.” Used in the “clustering” step of the de novo transcriptome RNA-seq workflow .	
results	Specifies the path and name of the result summary file. This file contains a compilation of assembly statistics and uses the extension fileSize.txt. Incomplete paths will be appended to the default directory.	[directory/ filename end in quotes]
saveUnSplitAssembly	Specifies whether XNG should save both the normal assembly output, [filename].assembly, and the unsplit intermediate assembly, [filename]-noSplit.assembly. The latter file contains SVs but no SNPs, and can be used to validate splits in the final assembly.	[true / false]
sex	Specifies the sex of the subject, used for read placement and SNP calling. See How sex chromosomes are handled for details.	[male / female / unknown]
showCDSVariant	Specifies whether or not XNG should show all variants of a CDS feature contacted by a SNP. The version number for the CDS variant will then appear in brackets when viewed in the SNP report in SeqMan Pro.	[true / false]
sngConvertOptions	(Intended for internal use only)	[text string]
snp	Specifies whether or not a SNP detection pass of the gapped alignment should be made during the assembly.	[true / false]
snp_checkStrandedness	Specifies whether or not the strand that each read comes from is considered in the SNP calculation. This is ignored by the simple SNP calling method (used when genome ploidy is “Heterogeneous”).	[true / false]
snp_combineSubs	This parameter is used to coalesce adjacent substitutions.	[true / false]
snp_excludeBases3p	(internal use only) This parameter causes the specified number of bases from the 3’ end of each read to not be considered during variant calling.	[integer]
snp_excludeBases5p	(internal use only) This parameter causes the specified number of bases from the 5’ end of each read to not be considered during variant calling.	[integer]
snp_excludeBasesEdge	This parameter causes the specified number of bases from both the 5’ and 3’ ends of each read to not be considered during variant calling.	[integer] For the simple SNP calling method (used when genome

		ploidy is “Heterogeneous” the default is 0. For the Bayesian SNP calling methods (used when genome ploidy is Diploid or Haploid), the default is 0.
snp_limitEndPos	Specifies the 3' most coordinate of the specified template from which to stop calculating SNPs.	[number between 1 and the length of the template]
snp_limitStartPos	Specifies the 5' most coordinate of the specified template from which to begin calculating SNPs. A value between 1 and the length of the template must be entered.	[number] Default = 1
snp_limitTemplateID	Specifies a single template ID for which to calculate SNPs.	[number] Default = 0
snp_logEndPos	Specifies the 3' most coordinate of the specified template from which to stop storing a detailed log of SNP information. A value between 1 and the length of the template must be entered.	[number] Default = 1
snp_logLevel	Specifies the level of detailed logging to store in the “shared” project directory as “SNP.log.” Level 0 specifies that no log will be stored. Level 1 stores detailed info on the SNPs which were called, level 2 also logs columns where the preliminary filtered passed but the final filtering failed, and level 3 logs all columns. This is ignored by the simple SNP calling method (used when genome ploidy is “Heterogeneous”).	[whole number from 0-3] Default = 0
snp_logStartPos	Specifies the 5' most coordinate of the specified template from which to begin storing a detailed log of SNP information. A value between 1 and the length of the template must be entered.	[number] Default = 1
snp_logTemplateID	Specifies a single template from which to store a detailed log of SNP information.	[number] Default = 0
snp_maxRun	Specifies the maximum length of a homopolymeric run for an indel to be considered during variant calling. For example, a snp_maxRun of '5' will allow a portion of sequence up to 5 bases in length to be called as a SNP.	[integer] Defaults are 454 and Ion Torrent read technologies

		for all others
snp_maxStrandBias	<p>Strand Bias (SB) for a SNP is the bias for the SNP appearing on one strand versus the other. It is measured relative to the strand bias in the assembly at the location of the SNP. For example, in a column with 60 forward reads and 40 backward reads, 6 SNP bases on the forward strands, and 4 on the reverse strands would be unbiased. SB is given by the formula:</p> $SB = SNP\%_f - SNP\%_r / \text{Total SNP}\%$ <p>...where $SNP\%_f$ and $SNP\%_r$ are the percentage of reads containing the variant on the forward (top) and reverse (bottom) strands, respectively; and $SNP\%$ is the total percentage of reads containing the variant. SB is calculated based on an “absolute value,” and will therefore be a positive number.</p> <p>The effect of different SB thresholds is shown below:</p> <p>-1 – A negative number cannot normally be generated by the equation above. However, you may use ‘-1’ in the script to turn off the snp_maxStrandBias parameter. In the wizard, SeqMan NGen indicates the parameter is turned off by making Maximum strand bias (see Variants tab) either blank or absent.</p> <p>0 – Perfectly balanced (unbiased) strands. Reads with variants are present on both strands, and variants appear equally on both strands.</p> <p>Between 0-1, not inclusive – As the number ‘1’ is approached, more variants are called with unbalanced variants containing reads at that position</p> <p>1 – All variant-containing reads are on a single strand.</p> <p>Note: In cases where all the reads covering a base are on one strand only, the $SNP\%$ of the other strand cannot be calculated (due to a “division by zero” error). These positions will not be removed by the snp_maxStrandBias filter. To remove these variants, instead set snp_minStrandCov to ≥ 1.</p> <p>Example:</p> <p>In a homozygous case ($SNP\% = 100$) with a depth of 100, where 75 variant containing reads are on the top strand (75%) and 25 variant</p>	<p>[integer]</p> <p>Defaults for Bayesian SNP calling method (used when genome ploidy is “Diploid” or “Haploid”) are 454 and Torrent read technologies shown (blank) for all others.</p> <p>Defaults for simple SNP calling method (used when genome ploidy is “Heterogeneous”) are 0.25 for read technologies</p>

	containing reads are on the bottom strand (25%), the strand bias would equal: $(75 - 25)/100 = 0.5$.	
snp_minHomopolDelDepth	Specifies the minimum read depth required to call a deletion in a homopolymeric run.	[integer] Default = 0
snp_minHomopolDelFrac	Specifies the minimum fraction of reads required to call a deletion in a homopolymeric run.	[integer] Default = 0
snp_minHomopolInsDepth	Specifies the minimum read depth required to call an insertion in a homopolymeric run.	[integer] Default = 0
snp_minHomopolInsFrac	Specifies the minimum fraction of reads required to call an insertion in a homopolymeric run.	[integer] Default = 0
snp_minPctToScore	Specifies minimum percentage of reads in a column which must differ from the reference in order to score the column. For the simple SNP calling method (used when genome ploidy is “Heterogeneous”), this is the only criteria used to call a SNP. For the Bayesian SNP calling methods (used when genome ploidy is “Diploid” or “Haploid”), this is a filter applied before the other parameters.	[number from 0 to 1] Default = 0.0
snp_minProbNonrefToCall	Specifies the minimum probability of a SNP column which is required to call a SNP, expressed as a number from 0 and 1. The probabilities of all genotypes other than Homozygous Reference are totaled and checked against this number. This is the final filter applied during the Bayesian SNP calling methods (used when genome ploidy is “Diploid” or “Haploid”) and is ignored by the simple SNP calling method (used when genome ploidy is “Heterogeneous”).	[number from 0 to 1] Default = 0.1, requiring a minimum 10% change.
snp_minStrandCov	Specifies the minimum number of reads from each strand required to call a variant at a given position.	[integer] In the Bayesian SNP calling methods (used when genome ploidy is “Diploid” or “Haploid”), default is 0. In the simple SNP calling method (used when genome ploidy is “Heterogeneous”), default is 1.

		“Heterogeneous” the default is
snp_minVariantDepthToScore	(required if “snp” is true) Specifies the minimum depth required for a specific base (or deletion) in a column before it is considered usable for SNP calling. This is the second filter applied during the Bayesian SNP calling methods (used when genome ploidy is “Diploid” or “Haploid”) and is ignored by the simple SNP calling method (used when genome ploidy is “Heterogeneous”).	[number from 0-100] Default = 2
snp_minWeight	Called “Minimum base quality score” in the SeqMan NGen wizard, this parameter specifies the minimum quality score for a base to be considered in the SNP calculation.	[number] In the simple calling meth (used when genome plo “Heterogene the default is In the Bayes SNP calling methods (us when genom ploidy is “Dip or “Haploid” default is 5.
snp_reportUserMissing	Specifies what kind of positions to put in the missingUser file, including one or more of the following: dbSNP = dbSNP Pos user = in user VCF SNP file zeroCoverage = include zero coverage regions cosmic = in COSMIC database allcaptured = include all positions in capture regions captured = include only positions in capture regions Example: snp_reportUserMissing: [user allcaptured captured] [kParamTypeStrFixedVocab]	
snp_runVar	Uses a Bayesian probabilistic model to exclude heterozygous insertions and deletions in homopolymeric runs. Intended for use with Ion Torrent data.	[true / false] Defaults: tru 454 and Ion Torrent read

		technologies <u>false</u> for all others.
snp_showAllFeatures	Specifies whether XNG should count SNPs multiple times if the SNP contacts different versions (variants) of a CDS feature.	[<u>true</u> / false]
snp_writeExtended	Specifies whether the additional values produced by the Haploid or Diploid SNP calculation methods are included in the SNP table. Wizard equivalent: Advanced Options > Alignment tab > Trim to targeted regions	[<u>true</u> / false]
snpMethod	Specifies the SNP detection method to use. Simple produces a count of each type of base in the column and calculates the percent of non-reference bases. Haploid uses a Bayesian statistical model to calculate a probability score that the position contains a polymorphism and give a quality score for the base called at that position. Diploid uses a Bayesian statistical model to calculate a probability score that the position contains a polymorphism and give a quality score for the base(s) called at that position. Based on the scores, it also calls the genotype at each position.	[simple / haploid / <u>diploid</u>]
splitTemplateContigs	Specifies under which circumstances contigs should be cut after a templated assembly. Any split contigs will be grouped into scaffolds with a defined position to allow for easy sorting when the project is viewed in SeqMan Pro. This command pertains only to reference-guided assemblies with gap closure. By default, during this type of assembly, the XNG assembler first finds structural variations (SVs) then splits the contig after each SV. Elements of this process can be modified using this command. 0 false = Don't split 1 true = Split at locations with zero coverage 2 = Split at insertions and deletions 3 = Split at zero coverage and at insertions	[integer between 0-3 / true / false] Default = 2
template	(required) Specifies the directory and file name of the reference sequence file. A folder with one or more reference sequence files can also be used in place of individual file names. Each entry must also be enclosed by brackets. If more than template entry is used, the list must also be enclosed by an additional set of brackets. Properties for template: file: [directory/filename enclosed in quotes]	[directory/filename enclosed in quotes]

	<p>Specifies the directory and file/folder.</p> <p>feature: [directory/filename enclosed in quotes] (optional) Specifies the directory and file name for annotated features when the reference sequence and feature annotations are in separate files.</p> <p>transcriptKind: [both identified novel] if the <i>.Transcriptome</i> package is used as a template, defines which transcripts will be used as a template.</p> <p>userSNP: [directory/filename enclosed in quotes]</p> <p>exomeCapture: file: [directory/filename enclosed in quotes] The BED file name.</p> <p>track: [string] the region of interest (Optional)</p> <p>merMask: [true / false] Specifies if mers from outside of the capture region should be excluded from assembly.</p> <p>Examples for template:</p> <p>Sequence and annotation in one file:</p> <p>AssembleTemplate</p> <p>template: {{file: "/data/home/proj/MG1655.gbk"} {file: "/data/home/proj/W3110.gbk"}}}</p> <p>Sequence and annotation in separate files:</p> <p>AssembleTemplate</p> <p>template: {file: "/Library/ABC_proj/references/MG1655.fas" feature: "/Library/ABC_proj/references/MG1655.gff"}</p>	
templateHitCntThresh	(Intended for internal use only)	[number]
trimToTargetRegions	<p>Controls whether reads are trimmed, by default, to the boundaries of the targeted regions, as defined by the .bed or manifest file. The default of <u>true</u> indicates that the reads are trimmed to the stated boundaries. If conditions are not met, the SeqMan NGen wizard does not change this parameter to 'false,' but instead omits it from the script. The parameter status is only shown in the script for control workflows.</p>	[<u>true</u> / false]

	Wizard equivalent: Trim to targeted regions in the Alignment tab . This tab is accessed from the Assembly Options screen by pressing the Advanced Options button.	
unassembled		[directory/ filename end in quotes]
verify		[true / false]

computeSNP

Sets parameters for the SNP computation phase of the assembly. The command is designed for use with existing BAM files that have not been analyzed for SNPs, or to re-analyze an existing file with different parameters.

Most of the parameters for **computeSNP** are identical to parameters for [assembleTemplate](#) and are discussed in that topic:

showCDSVariant
snp_logLevel
snp_minProbNonrefToCall
snp_checkStrandedness
snp_logStartPos
snp_minStrandCov
snp_combineSubs
snp_logTemplateID
snp_minVariantDepthToScore
snp_excludeBases3p
snp_maxRun
snp_minWeight

All other parameters are described in the table below:

Parameter	Description	Allowed values (defaults are underlined)
calcJunctionSeqs	In the structural variation workflow, specifying 'false' prevents junction sequences from being calculated.	[<u>true</u> / false]
concurrentAligns	(Intended for internal use only)	[number]
file	(required) Specifies the path and name of one or more .assembly projects from which to compute SNPs.	[directory/ filename enclosed in quotes]
snp_writeMissingDBSnps	In a SNP assembly, specifying 'false' causes missing SNPs not to be recorded, saving time and file space.	[<u>true</u> / false]
snpFilter	Specifies whether SNP filtering is turned on or off. Properties for snpFilter:	

capture: [true / false]

Specifies whether there is an exome capture file. If an exon capture file is added in the SeqMan NGen wizard or through a script, this value is set to 'true.' In the absence of an exome capture file, the SeqMan NGen wizard automatically sets this property to 'false.'

pNotRefMinVal: [number]

In the unusual case that the hard filter is missing, this property is used to set the minimum value that can be displayed in the SeqMan SNP table. Otherwise, this property is ignored. Default is 10.

userOnly: [true / false / All]

Specifies whether there is a VCF SNP file. The SeqMan NGen wizard always calls this as 'true' (or 'yes') but ignores the property if no VCF SNP file has been loaded.

pNotRef: [number]

This parameter is equivalent to Wizard option **Assembly Options > SNP Filter Stringency > pNotRef**. This parameter is a "soft" filter used to specify a PnotRef threshold. Data not matching the criterion are removed from the default display of the SeqMan Pro SNP table. This option is only available for the Bayesian SNP calling methods (used when genome ploidy is "Diploid" or "Haploid"). Wizard values include Low (90%), Medium (99%) and High (99.9%).

minSnpFilter: [number]

This parameter does not relate to any setting in the SeqMan NGen wizard, but corresponds to "SNP%" in SeqMan Pro and "minSNPFilter" in ArrayStar. In the simple SNP calling method (used when genome ploidy is "Heterogeneous"), the default is 5% for 454 and Ion Torrent read technologies; 1% for all others. In Bayesian SNP calling methods (used when genome ploidy is "Diploid" or "Haploid"), the default depends on stringency and ploidy rather than the read technology. The default for Diploid is 15% for all stringency levels. The default for Haploid is 25% for low stringency, 50% for medium and 75% for high.

minDepth: [number]

(optional) Specifies a minimum sequence depth threshold. This parameter does not relate to any setting in the SeqMan NGen wizard, but corresponds to **Depth** in SeqMan Pro and **minDepth** in ArrayStar. In the simple SNP calling method (used when genome ploidy is "Heterogeneous"), the default is 50. In Bayesian SNP

	<p>calling methods (used when genome ploidy is Diploid or Haploid), the default is <u>20</u>.</p> <p>The following set of SNP filters are used by ArrayStar and SeqMan Pro:</p> <p>codonOnly: [Coding / CodingChange / Nonsense / All] maxDepth: [number] maxCodingFeatureDistance: [number] minSnpFilter: [number] qCall: [number] synonymousCodingChange: [true/false] substitutionCodingChange: [true/false] noStartCodingChange: [true/false] noStopCodingChange: [true/false] nonsenseCodingChange: [true/false] frameshiftCodingChange: [true/false] notCodingCodingChange: [true/false] inFrameIndelCodingChange: [true/false] refOnly: [Reference/Unique/All] cosmicOnly : [Yes/No/All] minIndelSize: [number] gerpScore: [number] substitution: [true/false] showIndels: [true/false]</p>	
userSNP	Specifies a location for storing the VCF SNP table.	[directory/ filename enclosed in quotes]

createGenomeTemplate

The command **createGenomeTemplate** is intended for internal use only.

Parameter	Description	Allowed values (defaults in bold/ underline)
file	Specifies the directory and file/folder of the input file.	[directory/filename enclosed in quotes]
output	The path and name of the output file.	[directory/filename enclosed in quotes]

diskPath

The **diskPath** command is required, and defines the default directory where temporary intermediate files from the assembly will be stored. The files can be large with large scale projects. Visit our website to view [space requirements](#) for a range of representative projects.

Parameter	Description	Allowed values (defaults underlined)
clean	Specifies whether or not to clean the merge disk. When automated scripts are being run simultaneously or sequentially, this command can be useful for emptying the merge disk between assemblies.	[true false]
pathMac	Specifies the default path and file name for Macintosh.	[directory/ filename enclosed in quotes]
pathWin	Specifies the default path and file name for Windows.	[directory/ filename enclosed in quotes]
path	(required) Specifies the default path and file name. Example: diskPath path: "/data/proj/"	[directory/ filename enclosed in quotes]

dumpConsensus

The command **dumpConsensus** is intended for internal use only, and is used to convert the binary consensus file created during assembly into a text file.

Parameter	Description	Allowed values (defaults underlined)
file	Specifies the directory and file/folder.	[directory/filename enclosed in quotes]

dumpSNP

The command **dumpSNP** is intended for internal use only, and creates a tab delimited text file from one or more SNP containing binary files generated during assembly. SNP binary files include those with the .snpExt suffix contained in an .assembly package as well as those with either the .coverage.missingSNP or .nocoverage.missingSNP suffix contained in the _shared folder. To convert all the .snpExt files in a package simply use the .assembly name.

Parameter	Description	Allowed values (defaults underlined)
file	(required) Specifies the path and name of .assembly package (all SNP files will be included), one or more individual .snpExt files or either/both of the missingSNP files.	[directory/filename enclosed in quotes]
output	(required) Specifies the path and name of the output file.	[directory/filename enclosed in quotes]
refPos_end	To export SNPs with positions lower than this value.	[number]
refPos_start	To export SNPs with positions higher than this value.	[number]
snp_maxProbNonrefToCall	Lower limit for probability scores for exported SNPs.	[number]
snp_minProbNonrefToCall	Lower limit for probability scores for exported SNPs.	[number]
snp_type	Specifies which SNP file from the .assembly to use as an input.	[simple SNP missing user stats userIDOnly]
templateID	Defines the template for which the SNP will be exported.	[number]
onefile	Defines whether all SNPs should be placed into one file.	[true false]

execute

The command **execute** executes any shell script command.

Parameter	Description	Allowed values
command	Text for any shell script command.	[text string]

exportSplits

The command **exportSplits** is intended for internal use only, and is used to convert the binary splits file created during assembly into a text file.

Parameter	Description	Allowed values
file	Specifies the directory and file/folder.	[directory/filename enclosed in quotes]
output	The path and name of the output file.	[directory/filename enclosed in quotes]

exportVCF

The command **exportVCF** is used to accept the exome capture file and VCF file and builds another VCF file containing SNPs only in the capture regions.

Parameter	Description	Allowed values
userSNP	User SNP file.	[directory/filename enclosed in quotes]
exomeCapture	file: [directory/filename enclosed in quotes] Exome capture file. track: [text string] The name of the region of interest.	
output	The output VCF file.	[directory/filename enclosed in quotes]

extractPairs

The command **extractPairs** is used to create a tab delimited table of pair end information.

Parameter	Description	Allowed values
file	The path and name of any pair distance file (<i>.pairdist</i> file) from within a project's shared folder.	[directory/filename enclosed in quotes]
output	The path and name of the output file.	[directory/filename enclosed in quotes]

include

The command **include** is used to call up additional lines of script previously stored in a text file. In this way, a group of commands can be shared between two or more scripts.

Parameter	Description	Allowed values
file	Specifies the directory and file/folder.	[directory/filename enclosed in quotes]

loadAssembly

The command **loadAssembly** is intended for internal use only.

Parameter	Description	Allowed values
file	Specifies the directory and file/folder.	[directory/filename enclosed in quotes]

loadBAM

The command **loadBAM** is used to set parameters for analyzing existing BAM files. It allows ungapped BAM files to be converted into a fully gapped assembly file or to re-gap an existing file with different parameters. The command also permits SNPs to be calculated or re-calculated with different parameters starting with an existing BAM file. The associated parameters are also available for full assemblies and are described under the [assembleTemplate](#) command.

Parameter	Allowed values
align	
delayAlignInserts	[true false]
format	
gapPenalty	
increaseRunGapPen	[true false]
layout	
matchScore	
minAlignedLength	[0-999]
minMatchPercent	
mismatchPenalty	
output	
removeUniqueInserts	[true false]
snp	
snp_checkStrandedness	
snp_clusteredPosFilterMinDev	[number]
snp_clusteredPosFilterMinFromEdge	[number]
snp_hetKnownThresh	[number]
snp_hetThresh	[number]
snp_limitEndPos	
snp_limitStartPos	
snp_limitTemplateID	
snp_logEndPos	

snp_logLevel	
snp_logStartPos	
snp_logTemplateID	
snp_maxRun	[number]
snp_maxStrandBias	[number]
snp_minHomopolDelDepth	[number]
snp_minHomopolDelFrac	[number]
snp_minHomopolInsDepth	[number]
snp_minHomopolInsFrac	[number]
snp_minPctToScore	
snp_minProbNonrefToCall	
snp_minStrandCov	[number]
snp_minVariantDepthToScore	
snp_minWeight	
snp_nIMutationRate	<p>The chance that any single base is different from the reference. The default value of 0.0013 is equivalent to ~4 million variations in a Human sample against the reference or several thousand in a bacterial genome.</p> <p>[number]</p>
snp_observedInControlFilterMaxCount	[number]
snp_observedInControlFilterMaxFrac	[number]
snp_proximalGapFilterMaxDel	[number]
snp_proximalGapFilterMaxIns	[number]
snp_proximalGapFilterWindowSize	[number]
snp_reportUserMissing	[dbSNP user zeroCoverage cosmic allcaptured captured]
snp_runVar	[true false]
snp_showAllFeatures	[true false]
snp_writeExtended	
snp_writeMissingDBSnps	[true false]
snpMethod	

snpRefAsm	[quoted file name]
template	

mergelonTorrentShortReads

When using Ion Torrent data, the **mergelonTorrentShortReads** command causes overlapping short reads to merge into mini-contigs.

Parameter	Description	Allowed values
output	(required) Specifies the path and directory of the output files.	[directory/ filename enclosed in quotes]
query	(required) Specifies the directory and file name(s) of the query data to be assembled. A folder with one or data files can also be used in place of individual file names.	[directory/ filename enclosed in quotes]

message

The **message** command is used to write out the string to the standard output.

Parameter	Description	Allowed values
str	Specifies the string to be written to the standard output.	[text string]

pairFilePattern

The **pairFilePattern** command allows you to specify the pattern for pair files using the GREP language.

Example:

pairFilePattern

forward: "(? 'name'.*)_R1_(? 'ext'.*)\fastq

reverse: "(? 'name'.*)_R2_(? 'ext'.*)\fastq

Parameter	Description	Allowed values
forward	A naming pattern to match forward clones.	[text string enclosed in quotes]
reverse	A naming pattern to match reverse clones.	[text string enclosed in quotes]

pause

The command **pause** is used to create a pause. It can be used when running table scripts to stop at any point.

Example:

pause

prompt: "Table script paused. Press Enter to continue."

Parameter	Description	Allowed values
prompt	Text that should appear in the console. The pause is terminated by hitting the Enter key.	[text string enclosed in quotes]

quit

The **quit** command is used to terminate a script. This command does not have any parameters.

removeDuplicateSeqs

The **removeDuplicateSeqs** command is used to coalesce multiple identical reads at the same position into a single read, provided the reads match the reference exactly. If this feature is active, at the end of assembly, XNG will print the message: "Coalesced \$Ild identical reads that matched the template exactly." This command does not have any associated parameters.

Allowable values are [true|false] The default is false.

runScript

The **runScript** command allows you to batch multiple projects of the same type (e.g. assembly, computeSNPs). There are required three file: 1) a runScript file with variables, 2) a file with a table of values for the variables, and 3) a script file specifying the action to be carried out.

Example (runScript file):

```
setDefaultDirectory directory: "."
set $force: false
set $DataDisk: "/Volumes/Raid/DataDisk"
set $ResultDisk: "/Volumes/ResultDisk"
set $MergeDisk: "/Volumes/MergeDisk0"
set $snp:true
set $snpMethod:"Diploid"
set $repCnt:100
set $merLayoutMin:19
diskPath path: {"${MergeDisk}/mergeSort Data"}}
runScript table: "testAssembly.txt" script: "testAssembly.template.script"
```

Example (table file):

```
defaultDir template query isPair seqTech project merSize snp snpMethod
"${ResultDisk}/rice" "${DataDisk}/rice.genome" "${DataDisk}/rice" FALSE Illumina rice 21 TRUE Diploid
"${ResultDisk}/ecoli" "${DataDisk}/Ecoli.gbk" "${DataDisk}/ecoli" TRUE Illumina Ecoli 21 TRUE Diploid
"${ResultDisk}/Exome" "${DataDisk}/GRCh37.gbk" "${DataDisk}/Sample1" FALSE 454 HuEx 19 TRUE Diploid
```

Example (script file):

```
; "assembly.template.script"
setMachineMemory memory:32
setDefaultDirectory directory: $defaultDir
compareSeqs template: $template
query: {file: $query
isPair: $isPair
seqTech: $seqTech}
directoryMer: "intermediateFiles"
; directoryQueryMer: "intermediateFiles"
hits: "intermediateFiles/${project}.hits"
layout: "intermediateFiles/${project}.layout"
output: "results_${mersize}_${merSkipQuery}/${project}"
; results per project results: "${project}.results.txt"
```

```

; aggregate all results
results: "${ResultDisk}/assembly.results.txt"
merSize: $mersize
merSkipQuery: $merSkipQuery
repeatCnt: $repCnt
merLayoutMin: $merLayoutMin
layoutType: once
maxGap: 6
format: BAM
onePackage: true
snp: $snp
snpMethod: $snpMethod
; snp_writeExtended: true
forceMake: $force

```

Parameters for this command are described below:

Parameter	Description	Allowed values
script	The filename and location of the script.	[directory/filename enclosed in quotes]
table	The filename and location of the file containing text strings and numbers values for each variable.	[directory/filename enclosed in quotes]
inline	Executes the list of commands and parameters.	

set

The command **set** is used to set variables. It does not have any associated parameters. See the example below and those under the [runScript](#) command.

Example:

```
set $snp:true  
set $snpMethod:"Diploid"
```

setDefaultDirectory

The **setDefaultDirectory** command (required) defines the default directory for the project. When a default directory is specified (see table below), files located in that directory only need to be identified by their subfolder and/or file name in subsequent commands.

Parameter	Description	Allowed values
directory or defaultDirectory	(required) Specifies the default directory. Previously called defaultDirectory.	[directory/filename enclosed in quotes]
directoryMac or defaultMacDirectory	Specifies the default directory for Macintosh. Previously called defaultMacDirectory.	[directory/filename enclosed in quotes]
directoryWin or defaultWinDirectory	Specifies the default directory for Windows. Previously called defaultWinDirectory.	[directory/filename enclosed in quotes]

Example:

```
setDefaultDirectory
  directory: "/data/home/proj/"
```

setMachineMemory

The **setMachineMemory** command defines the amount of random access memory (RAM) that the program will use. Limiting the amount of RAM available to the assembler allows you to use the computer for other purposes while an assembly is running. However, this will likely slow down the assemblies and is not recommended for large projects.

Parameter	Description	Allowed values
memory	(required) Amount of RAM (in GB) to be used, entered in multiples of four. Entering a value greater than the available RAM causes all RAM to be used.	[number that is a multiple of 4]

Example:

```
setMachineMemory  
  memory: 32
```

setParam

The **setParam** command adjusts the stringency of one or more of the assembling parameters for the project. SeqMan NGen will use the default values for any parameter that is not specified within the script.

All of the parameters for setParam are identical to the same parameters described in the [assembleTemplate](#) command topic:

delayAlignInserts
gapPenalty
increaseRunGapPen
matchScore
minAlignedLength
minMatchPercent
mismatchPenalty
removeUniqueInserts

SNG commands

The following commands can be used in the script for an SNG assembly. Click on the name of a command in the table below (shown in alphabetical order) to see a description and example, and to see the parameters associated with that command. All SNG commands and parameters are assumed to be optional unless the description states that it is required.

Project management	File loading	Parameter settings	Preprocessing and assembling
closeProject	load454PairedEnd	setContaminantParam	appendToAssembly
runScript	loadConstraint	setContaminantParam	assemble
saveProject	loadContaminant	setParam	convertReads
saveReport	loadLayout	setQualityParam	extendContigs
writeUnassembledSeqs	loadRepeat	setRepeatParam	fixedTrim
	loadSeq	setVectorParam	include
	loadTemplate		makeSeqNamesUnique
	loadVector		realignContigs
	openProject		removeSmallContigs
	setDefaultDirectory		set
			setAssemblyReport
			setPairSpecifier
			splitLinkerReads
			splitMIDSeqs
			splitPairs
			splitTemplates
			trimVector

Project management commands

SNG “project management” commands include:

- [closeProject](#)
- [runScript](#)
- [saveProject](#)
- [saveReport](#)
- [writeUnassembledSeqs](#)

closeProject

The command **closeProject** closes the current project and frees the memory in use so that the system is ready for additional assemblies. This can be useful if you want to run multiple assemblies in one script.

runScript

The **runScript** command allows you to run a table script within the current script. A table script references variable values for specified parameters and other elements in a script. This enables you to run multiple projects from the same script, substituting new parameter values and other variables each time. SeqMan NGen will run the table script repeatedly, using the variable values from one row of the table for each iteration of the script until all of the rows have been used.

Parameter	Description	Allowed values
file	Specifies the directory and file/folder.	[directory/filename enclosed in quotes]
script	(required) Specifies the directory and file name of the table script you wish to run.	[directory/filename enclosed in quotes]
table	(required) Specifies the delimited text file containing the variable values.	[directory/filename enclosed in quotes]

Example:

runScript

script: "/Library/abc_Project/abc_script.script"

table: "/Library/abc_Project/table.txt"

saveProject

The **saveProject** command saves the assembly to a project file. By default, the SeqMan Pro / SeqMan Ultra project file format (.sqd) is used. Phrap (.ace) and FASTA (.fas) formats may also be specified by using the format parameter, and specifying the desired file extension using the file parameter.

* **Note:** As a command-line tool, SeqMan NGen will not prompt you if you try to save a new project file with the same name as an existing file in the same location. When you run a script multiple times, be sure to change the file name of the project to be saved each time to prevent existing project files from being overwritten.

Parameter	Description	Allowed values
file	(required) Specifies the directory and file name of the project file to be saved.	[directory/filename enclosed in quotes]
format	<p>Specifies the output file format.</p> <ul style="list-style-type: none"> SeqMan - Saves a 64-bit SeqMan Pro project file (.sqd) that is compatible with SeqMan Ultra or SeqMan Pro version 8.1 and higher (default). SeqMan8 - Saves a 32-bit SeqMan Pro project file (.sqd) that is compatible with SeqMan Ultra or SeqMan Pro version 8.0 and higher. SeqMan7 - Saves a 32-bit SeqMan Pro project file (.sqd) that is compatible with SeqMan Ultra or SeqMan Pro version 7.2 and higher. Note that this project file will be much bigger than the same project created in either of the SeqMan formats listed above. 	[SeqMan SeqMan8 SeqMan7 Phrap Fasta BAM SAM]

	<ul style="list-style-type: none"> • Phrap - Saves an <i>.ace</i> file. • Fasta - Saves <i>.fas</i> and <i>.qual</i> files of the consensus sequence for each contig. • BAM - Saves a BAM file (SNG/SMNG reference-guided assemblies only). • SAM - Saves a SAM file (SNG/SMNG reference-guided assemblies only). 	
onePackage	Specifies whether an assembly containing multiple reference sequences should be bundled into a single <i>.assembly</i> package. If 'false' is entered, one <i>.assembly</i> package is created per contig.	[true false]
openInSeqMan	Specifies whether to automatically launch SeqMan Pro / SeqMan Ultra and open the completed assembly once the script has completed.	[true false]

Example:

SaveProject

file: "/Library/My projects/ABC_project.sqd"

format:seqman

openInSeqMan:true

saveReport

The **saveReport** command exports a report as a text file that summarizes assembly statistics, including the parameters used, the number of assembled/unassembled sequences and contigs, average quality scores, and the number of sequences excluded from the assembly due to exceeding the **maxAssemblyCoverage** parameter of the [*setParam*](#) command.

The same information contained within this report is also saved within the SeqMan Pro project file (.sqd) regardless of whether you choose to export the report by setting this parameter. The report can be viewed in SeqMan Pro / SeqMan Ultra using the **Project > Report** command.

Parameter	Description	Allowed values
file	(required) Specifies the directory and file name of the report to be saved.	[directory/filename enclosed in quotes]

Example:

saveReport

file: "/Library/abc_Project/abc_report.txt"

writeUnassembledSeqs

The **writeUnassembledSeqs** command saves all sequences that were not assembled in the project as *.fas* and *.qual* files.

Parameter	Description	Allowed values
file	(required) Specifies the directory and file name of the unassembled sequences to be saved.	[directory/filename enclosed in quotes]
saveTrimmed	Specifies whether to save only the trimmed portion of the unassembled sequences.	[true / <u>false</u>]

File loading commands

SNG “file loading” commands include:

- [load454PairedEnd](#)
- [loadConstraint](#)
- [loadContaminant](#)
- [loadLayout](#)
- [loadRepeat](#)
- [loadSeq](#)
- [loadTemplate](#)
- [loadVector](#)
- [openProject](#)
- [setDefaultDirectory](#)

load454PairedEnd

The **load454PairedEnd** command loads a file of Roche 454 sequences and checks for the presence of a linker defining the paired end sequences. If the linker is found, the linker is removed and the remaining portion is split into two sequences linked with a paired end constraint.

Parameter	Description	Allowed values
DiscardLinkerless	Specifies whether to discard any read where no portion of the mate pair linker was found. In this way, reads that do not have a linker sequence will be discarded from the assembly.	[true false]
file	The directory and file name of the .fas, .fna, or .sff file containing the 454 sequences.	[directory/ filename enclosed in quotes]
linker	The directory and file name of the .fas, fna, or .sff file containing the 454 linker sequences. If not specified, SeqMan NGen will use its default 454 linker sequence: GTTGGAACCGAAAGGGTTTGAATTCAAACCCTTTCGGTTCCAAC.	[directory/ filename enclosed in quotes]
max, maxDistance	The maximum distance for the paired end constraint.	[number] Default = <u>10000</u>
min, minDistance	The minimum distance for the paired end constraint.	[number] Default = <u>0</u>

Example:

load454PairedEnd

file: "/Library/454 data/123_Pairedend.fas"

linker: "/Library/454 data/123_linkerseqs.fas"

min: 0

max: 10000

DiscardLinkerless: false

loadConstraint

The **loadConstraint** command loads a constraint file. The file can be in the NCBI ancillary file format, or in the CAP3 constraint file format. SeqMan NGen uses constraint files to identify paired end reads, similar to using the [setPairSpecifier](#) command. Constraint files in the NCBI ancillary file format also contain trimming information, which SeqMan NGen will load and use. SeqMan NGen will create a CAP3 file when saving a Phrap project (.ace) that used paired end constraints.

Parameter	Description	Allowed values
file	The directory and file name of the constraint sequence file.	[directory/filename enclosed in quotes]

Example:

loadConstraint

file: "/Library/constraints/123_xyz.con"

loadContaminant

The **loadContaminant** command loads a contaminant sequence file to be used to identify known contaminants, such as primers, in the assembly. Sequences that contain at least 12 matching 17-mers are flagged as contaminant sequences and will be removed from the assembly.

Parameter	Description	Allowed values
file	The directory and file name of the contaminant sequence file. A folder may also be specified, in which case all of the sequence files within that folder will be loaded and used for contaminant screening.	[directory/ filename enclosed in quotes]

Example:

loadContaminant

file: "/Library/contaminants/123_abc.seq"

loadLayout

The **loadLayout** command loads a layout file to be used for an assembly. The format may be either a SOLiD General Feature Format file (.gff) or a File of Filenames file (.fof). When this command is used, SeqMan NGen still aligns each read from the file to the reference, but uses the information contained within the specified file to determine the overall layout of reads.

Parameter	Description	Allowed values
layoutFile	(required) Specifies the directory and file name of the layout file. Both .gff and .fof formats are accepted	[directory/filename enclosed in quotes]
templateFile	(required) Specifies the directory and file name of the reference sequence file.	[directory/filename enclosed in quotes]

Example:

loadLayout

templateFile: "/Library/123_project/template.seq"

layoutFile: "/Library/123_project/layoutfile.gff"

loadRepeat

The **loadRepeat** command loads a sequence file to be used to identify repeat sequences in the assembly. All sequences identified as repeats will be added to the assembly last, after all non-repeats have been assembled. See our website for a list of supported file types.

Parameter	Description	Allowed values
file	(required) Specifies the directory and file name of the repeat sequence file. A folder may also be specified, in which case all of the sequence files within that folder will be loaded and used as repetitive sequences.	[directory/ filename enclosed in quotes]

Example:

loadRepeat

file: "/Library/repetitive_seqs/123_repeat.seq"

loadSeq

The **loadSeq** command loads a sequence file or files for assembly. See our website for a list of [supported file types](#).

Parameter	Description	Allowed values
blockContig	Used in the reference-guided workflow.	[text string]
blockContigID	Used in the reference-guided workflow.	[number]
blockName	Used in the reference-guided workflow.	[text string]
blockPos	Used in the reference-guided workflow.	[number]
DiscardLinkerless	Specifies whether reads that do not have a linker sequence should be discarded from the assembly.	[true / false]
file	(required) Specifies the directory and file name of the sequence file(s) to be loaded. A folder may also be specified, in which case all of the sequence files within that folder will be loaded.	[directory/filename enclosed i
groupName	Used to identify the multi-sample group name for a read file.	[text string]
isPair	Specifies whether the query files contain paired end data.	[true / false]
linker	The directory and file name of the .fas, fna, or .sff file containing the 454 linker sequences. If not specified, SeqMan NGen will use its default 454 linker sequence: GTTGGAACCGAAAGGGTTTGAATTCAAACCCTTTTCGGTTCCAAC.	[directory/filename enclosed i
max	The maximum distance for the paired end constraint.	[number] Default = <u>10000</u>
maxSeqs	Specifies the maximum number of reads to load from a file.	[number]
mergePairs	Specifies whether the reads are paired end data that overlap and should therefore be merged.	[true / false]
min	The minimum distance for the paired end constraint.	[number] Default = <u>0</u>
minSeqLen	Minimum length of a sequence required to include it in the assembly.	[number]
multiplex	Specifies whether reads are from a multi-sample run.	[true / false]
seqTech	Specifies the offset to be used when converting compressed quality scores into numerical values. These are the offsets used for the technology specified:	[IonTorrent SOLiD Illumina 45

	<p>Data Type / Value / Offset</p> <p>IonTorrent / IonTorrent / 33</p> <p>Applied Biosystems SOLiD / SOLiD / 33</p> <p>Illumina / Illumina / 64</p> <p>Roche 454 / 454 / 33</p> <p>Other types / normalScore / 33</p> <p>Note 1: For 454, quality scores for homopolymeric runs of ≥ 2 are oriented from 5' to 3' on the top strand.</p> <p>Note 2: If possible, the data type of unknown data is determined automatically based on the first data file.</p>	
templateFragment	Used in reference-guided assemblies with gap closure.	[number]

Example:

loadSeq

file: "/Library/ABC_project/ABC_sequences.fas"

loadTemplate

The **loadTemplate** command loads a sequence file to be used as a reference for all other sequences to be assembled to. The sequence will be displayed as a “reference” sequence in SeqMan Pro for SNP analysis.

Parameter	Descriptions	Allowed values
file	(required) Specifies the directory and file name of the reference sequence file to be loaded. A folder may also be specified, in which case all of the sequence files within that folder will be loaded and treated as reference sequences.	[directory/ filename enclosed in quotes]

Example:

loadTemplate

file: “/Library/abc_Project/abc_template.seq”

loadVector

The **loadVector** command loads a vector sequence file to be used for vector trimming.

Parameter	Description	Allowed values
cloneSite	This parameter specifies the position of the cloning site on the vector where insertion occurs.	[number]
file	(required) Specifies the directory and file name of the vector sequence file to be used for vector trimming.	[directory/filename enclosed in quotes]

Example:

loadVector

file: "/Library/vectors/123_vector.seq"

cloneSite:826

openProject

The **openProject** command loads an existing assembly project into memory.

Parameter	Description	Allowed values
file	(required) Specifies the directory and file name of the project file to be loaded.	[directory/filename enclosed in quotes]

setDefaultDirectory

The **setDefaultDirectory** command is required, and defines the default directory for the project. When a default directory is specified, files located in that directory only need to be identified by their subfolder and/or file name in subsequent commands.

Parameter	Description	Allowed values
directory	(required) Specifies the default directory.	[directory/filename enclosed in quotes]
defaultMacDirectory	Specifies the default directory for Macintosh.	[directory/filename enclosed in quotes]
defaultWinDirectory	Specifies the default directory for Windows.	[directory/filename enclosed in quotes]

Example:

```
setDefaultDirectory: "/Library/ABC_proj/"
```

Once you have set a default directory, you may use two periods before a file name to specify that the file you wish to use is located in the parent folder of the default directory you specified. For example, the following line specifies that the vector file, 123Vector.fas, is located in the ABC Data folder, the parent folder of the default directory.

```
loadVector file: "../123Vector.fas"
```

Parameter settings commands

SNG “parameter setting” commands include:

- [setContaminantParam](#)
- [setParam](#)
- [setQualityParam](#)
- [setRepeatParam](#)
- [setVectorParam](#)

setContaminantParam

The **setContaminantParam** command allows you to adjust the parameters used for scanning for contaminant sequences. In order to be applied, this command must appear in the script before the [loadContaminant](#) command, and the **contamScan** parameter for the [assemble](#) command must be set to **true**.

Parameter	Description	Allowed values
MerLength	The minimum length of a mer required to be considered an exact match when scanning for contaminants. Wizard equivalent: Advanced Trim/Scan Options > Mer length .	[number from 5-50] Default = <u>17</u>
MinMerMatch	The minimum number of matching mers required to mark the sequence as a contaminant. Wizard equivalent: Advanced Trim/Scan Options > Minimum matches .	[number from 1-50] Default = <u>12</u>

Example:

```
setContaminantParam MerLength:17
setContaminantParam MinMerMatch:12
```

setParam

The **setParam** command allows you to adjust the stringency of one or more of the assembling parameters for the project. SeqMan NGen will use the default values for any parameter that is not specified within the script.

Parameter	Description	Allowed values
AllowConstraintBased	Specifies whether the assembler should use constraints during assembly.	[<u>true</u> / false]
AssembleBoneyard	Specifies whether, after a reference-guided assembly has been completed, the unassembled sequences remaining should be assembled into contigs. If the reference has been split, SeqMan NGen will attempt to join the split contigs together in new arrangements. (Note: “Boneyard” is a term for sequences that were not assigned to any contig). Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > De novo assemble unassembled reads .	[true / <u>false</u>]
CoverageType	Specifies the type of coverage to be used for repeat handling. ‘Genome’ uses the length of the genome being assembled to calculate the expected coverage. ‘Fixed’ uses a fixed value as the expected coverage. If you know the length of the genome/fragment being assembled, we recommend using ‘genome’ for this parameter and then specifying the length using the genomeLength parameter. If you do not know the genome/fragment length, use ‘fixed’ and provide the most accurate estimate of expected coverage for the FixedCoverage value.	[<u>genome</u> / fixed]
DefaultQuality	The value used for the base quality of sequences without quality scores. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only):	[number from 5-100] Default = <u>15</u>

	Assembly Options > Default quality.	
FixedCoverage	The estimated depth of the sequencing, which can be used instead of the genome length for repeat handling. Use caution when estimating the value for fixedCoverage. If the value you use is significantly lower than the actual depth, the assembly may take a much longer time to complete and may have too many mers flagged as repeats.	[number from 1-65535] Default = <u>20</u>
GapPenalty	The penalty for opening or extending a gap during an alignment. This penalty is deducted from the pairwise score used to calculate match percentage. A high gap penalty suppresses gapping, while a low value promotes gapping. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Gap penalty.	[number from 0-1000] Default = <u>30</u> for most workflows; <u>50</u> for the <i>de novo</i> transcriptome RNA-seq workflow .
GenomeLength	Specifies the length of the genome or fragment being assembled. This is used to calculate expected coverage in determining repeat handling. (Note: this parameter was called “setGenomeParam” prior to SeqMan NGen 2.0.).	[number from 0-1015 ULL] Default = <u>0</u>
HaploidSNP	Specifies whether to use the second most common base at a position when performing SNP passes. (See the snpPasses parameter). Using this parameter will increase the SNP percentage for SNPs occurring on one allele of a diploid genome in a reference-guided assembly. When haploidSNP is set to ‘true,’ the lowCoverageThreshold parameter value should be greater than zero.	[true / <u>false</u>]
HaploidThreshold	The minimum number of times that the second most common base must occur at a position in order for it to be used to find SNPs during haploid SNP passes.	[number from 0-100] Default = <u>0</u>

	(See the haploidSNP parameter above).	
LowCoverageThreshold	The minimum coverage required in an assembly to be excluded from SNP passes. SeqMan NGen will include regions in an assembly that have coverage less than the value specified as well as regions with zero coverage when it performs SNP passes. (See the snpPasses parameter). Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > SNP low cover cutoff .	[number from 0-10000] Default = <u>0</u>
MatchRepeatPercent	The percent frequency a mer occurs compared to its expected frequency. Mers exceeding this value are flagged as repeated and not used as mer tags in determining overlaps. (Note: this parameter was called “ maxCoverageRatio ” prior to SeqMan NGen 2.0.). Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Match repeat percent .	[number from 100-1000] Default = <u>150</u>
MatchScore	The score for a base match during an alignment. This score contributes to the pairwise score used to calculate match percentage. Increasing the matchScore value will allow for longer or more frequent gaps, thus forcing bases that match to be assembled together. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Match score .	[number from 1-1000] Default = <u>10</u>
MatchSize	The minimum number of matching consecutive bases required to determine the overlap of sequence reads. If an even number is entered, SeqMan NGen will automatically increase the value to the next odd number. (Note: this parameter was	[odd whole number] Default = <u>21</u>

	called setParamMerLength prior to SeqMan NGen 2.0.). Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Mer size.	
MatchSpacing	The length of the window of a sequence read where at least one mer tag will be chosen. (Note: this parameter was called “merTagWindow” prior to SeqMan NGen 2.0.). Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Match spacing.	[number from 1- 1000000] Default = <u>50</u>
MatchWindowLength	The size of the window used to calculate the match percentage. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Match window.	[number from 10-1000] Default = <u>50</u>
MaxAssemblyCoverage	The maximum depth of coverage allowed in the reference-guided assembly. SeqMan NGen will not exceed the coverage specified by this threshold. This parameter is only available for reference-guided assemblies, and should be used with caution as it will limit the number of sequences included in the assembly. A value of 0 indicates unlimited coverage. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Maximum coverage.	[number from 0-65535] Default = <u>0</u>
MaxContigs	The maximum number of contigs to write to an .assembly project. This command is not generally needed due to SeqMan’s capacity to handle a very large number of contigs.	[number]
MaxGap	The theoretical maximum length of a gap that could be inserted. In practice, the maximum gap size will usually be about half of this value. Wizard	[number from 0-99] Default = <u>6</u>

	equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Max gap.	
MaxUsableCount	Any mers occurring more frequently than FixedCoverage multiplied by MaxUsableCount are disregarded as mer tags from the assembly. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Max usable.	[number from 1-65535] Default = <u>25</u>
MinContigSeqs	The minimum number of sequences in a contig. After an assembly has been completed, any contigs without a reference sequence will be disassembled if they contain fewer sequences than the number specified. The use of this parameter is recommended when performing <i>de novo</i> assemblies using data from Next Generation sequencing technologies, such as Illumina, as these types of assemblies can produce tens of thousands of very small contigs.	[number from 0-10000] Default = <u>0</u>
Minimizer	(Intended for internal use only). An experimental way of choosing mer tags that may save time and memory. The accuracy of this parameter has not been verified by DNASTAR.	[number]
MinMatchPercent	The minimum percentage of matches in an overlap required to join two sequences in the same contig. (Note: this parameter was called “minMatchPercentage” prior to SeqMan NGen 2.0.). Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Minimum match percentage.	[number from 0-100] Default = *+93+
MismatchPenalty	The penalty for a base mismatch during an alignment. This penalty is deducted from the pairwise score used to calculate match percentage. Wizard	[number from 0-1000] Default = <u>20</u>

	equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Mismatch penalty .	
SkipRealign	This parameter only affects <i>de novo</i> assemblies, and specifies whether to skip the realignment step of the assembly. The realignment step will then analyze each sequence at the nucleotide level to determine the exact position of each sequence in the alignment. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Realign reads after assembly .	[true / <u>false</u>]
SNP	Specifies whether a SNP detection pass of the gapped alignment is made during the assembly.	[<u>true</u> / false]
snp_checkStrandedness	Specifies whether the strand that each read comes from is considered in the SNP calculation. This is ignored by the simple SNP calling method (used when genome ploidy is “Heterogeneous”).	[true / <u>false</u>]
snp_minPctToScore	Specifies minimum percentage of reads in a column which must differ from the reference in order to score the column. For the simple SNP calling method (used when genome ploidy is “Heterogeneous”), this is the only criteria used to call a SNP. For the Bayesian SNP calling methods (used when genome ploidy is “Diploid” or “Haploid”), this is a filter applied before the other parameters.	[number from 0-1] Default = <u>0.05</u>
snp_minProbNonrefToCall	Specifies the minimum probability of a SNP column which is required to call a SNP, expressed as a number from 0 and 1. The probabilities of all genotypes other than Homozygous Reference are totaled and checked against this number. This is the final filter applied	[number from 0-1] Default = <u>0.1</u> , requiring a minimum 10% change

	during the Bayesian SNP calling methods (used when genome ploidy is “Diploid” or “Haploid”). This is ignored by the simple SNP calling method (used when genome ploidy is “Heterogeneous”).	
snp_minVariantDepthToScore	(required if “snp” is true) Specifies the minimum depth required for a specific base (or deletion) in a column before it is considered usable for SNP calling. This is the second filter applied during the Bayesian SNP calling methods (used when genome ploidy is “Diploid” or “Haploid”). This is ignored by the simple SNP calling method (used when genome ploidy is “Heterogeneous”).	[number from 0-100] Default = <u>2</u>
snp_minWeight	Called “Minimum base quality score” in the SeqMan NGen wizard, this parameter specifies the minimum quality score for a base to be considered in the SNP calculation.	[number]
SNPMatchPercentage	The minimum match percentage required during passes to fill in SNP regions. See the snpPasses parameter. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > SNP match percent .	[number from 0-100] Default = <u>90</u>
snpMethod	Specifies the SNP detection method to use. Simple produces a count of each type of base in the column and calculates the percent of non-reference bases. Haploid uses a Bayesian statistical model to calculate a probability score that the position contains a polymorphism and give a quality score for the base called at that position. Diploid uses a Bayesian statistical model to calculate a probability score that the position contains a polymorphism and give a quality score for the base(s) called at	[simple haploid *diploid* population]

	that position. Based on the scores, it also calls the genotype at each position.	
SNPPasses	The number of times SeqMan NGen will cycle through a reference-guided assembly, attempting to fill in regions with low coverage or no coverage due to SNPs. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > SNP passes .	[number from 0-10] Default = <u>2</u>
SplitFalseJoins	Specifies whether the assembler should identify and splits false joins based on the set of false join parameters indicated.	[true / <u>false</u>]
SplitTemplateContigs	Specifies whether, after a reference-guided assembly has been completed, the template should be split into contigs at areas where there is zero coverage. Split contigs will be grouped into scaffolds with a defined position to allow for easy sorting when the project is viewed in SeqMan Pro. Annotations on the reference sequence will also be split, and any /codon_start qualifiers will be adjusted to stay in frame.	[true / <u>false</u>]
TemplateDefaultQuality	The value used for the base quality of template sequences without quality scores. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Default template quality .	[number from 5-50000] Default = <u>500</u>
TrimToMer	Specifies whether to trim the reads to the matching mer tags within the read. For each read, SeqMan NGen looks for mers that exist in the template (for templated assemblies) or in any other read in the assembly (for <i>de novo</i> assemblies). It then sets the trimming for the read to the start of the first mer found and the end of the last mer	[true / <u>false</u>]

	found. Trimming to mer may be useful when assembling data without accurate quality scores, data with very short linkers, or when assembling SOLiD data.	
UseRepeatHandling	Specifies whether to use the repeat probabilities to determine if a mer occurs too frequently to use. This parameter should only be used for <i>de novo</i> assemblies, unless the assembleBoneyard parameter is set to 'true' for the templated assembly. Wizard equivalent (<i>de novo</i> or special reference-guided workflows only): Assembly Options > Repeat handling.	[<u>true</u> / false]

Example A:

```

setParam SNP: true
setParam snp_minVariantDepthToScore: 2
setParam snp_minWeight: 5
setParam snp_combineSubs: true
setParam snp_excludeBasesEdge: 0
setParam snp_maxRun: -1
setParam snp_maxStrandBias: -1
setParam snp_minHomopolDelDepth: 0
setParam snp_minHomopolDelFrac: 0
setParam snp_minHomopolInsDepth: 0
setParam snp_minHomopolInsFrac: 0
setParam snp_minSoftDepth: -1
setParam snp_minSoftPnotRefPct: -1
setParam snp_minSoftSnpPct: -1
setParam snp_minStrandCov: 0
setParam snp_runVar: false
setParam snp_checkStrandedness: false
setParam snp_minProbNonrefToCall: 0.1
setParam SNPmethod: diploid
setParam snp_minPctToScore: 0.05

```

Example B:

In the de novo transcriptome RNA-seq [workflow](#), reads clustered with XNG are reassembled using SNG. In order to minimize mis-joins, the initial phase of the assembly is done at high stringency using the following parameters:

```
setParam  
merLength: 21  
minMatchPercent: 97  
useRepeatHandling: false  
minContigSeqs: 101
```

Two assembly passes are performed for each read cluster. During the first pass, contigs are assembled from the reads after which those with less than 101 reads are dis-assembled and added to the unassembled sequences pool for that cluster. During a second pass SNG attempts to merge the assembled contigs and add any of the unassembled sequence reads from the first pass. To facilitate merging, **minMatchPercent** is lowered to 85 for this pass.

```
setParam  
minMatchPercent: 85
```

setQualityParam

The **setQualityParam** command allows you to adjust the parameters used for quality trimming. In order to be applied, the **trimEnds** parameter for the **assemble** command must be set to 'true.'

Parameter	Description	Allowed values
EndRegion	The number of bases at the end of a sequence considered to be the “end region” which is used by other quality parameters.	[number from 1-100] Default = <u>5</u>
MaxN	The maximum number of “N” bases permitted in the window used for N-based quality trimming.	[number from 1-100] Default = <u>2</u>
MaxNHiQual	The maximum number of “N” bases permitted in the window used for N-based quality trimming to meet the high-quality threshold.	[number from 0-100] Default = <u>1</u>
MinAveHiQual	The minimum averaged quality score of the evaluated window required to be considered high-quality.	[number from 10-40] Default = <u>22</u>
MinAveLowQual	The minimum averaged quality score of the evaluated window required to be considered low-quality. Wizard equivalent: Advanced Trim/Scan Options > Minimum quality.	[number from 5-40] Default = <u>20</u>
MinEndBaseQual	The minimum quality base score required in the specified end region.	[number from 5-40]

		Default = <u>15</u>
NTrimWinLength	The length of the window used for “N-based” quality trimming. N-based quality trimming trims bases that are called “N” and is used only when quality scores are not available.	[number from 5-100] Default = <u>7</u>
WinLength	The length of the window used for averaging quality scores. Wizard equivalent: Advanced Trim/Scan Options > Window.	[number from 2-100] Default = <u>5</u>

Example:

```

setQualityParam winLength:30
setQualityParam minAveLowQual:14
setQualityParam minAveHiQual:18
setQualityParam minEndBaseQual:15
setQualityParam endRegion:15
setQualityParam nTrimWinLength:50
setQualityParam maxN:2
setQualityParam maxNHiQual:1

```


setRepeatParam

The **setRepeatParam** command allows you to adjust the parameters used for scanning for repetitive sequences. In order to be applied, this command must appear in the script before the **loadRepeat** command, and the **repeatScan** parameter for the **assemble** command must be set to 'true.'

Parameter	Description	Allowed values
AlignCutoff	The minimum acceptable alignment score. When the alignment score drops below the specified value, this indicates that the end of the alignment between the read and the repeat has been reached, and the alignment will stop.	[number from 10-1000000] Default = <u>100</u>
MaxMerGap	The maximum distance between two mers required to be considered a matching pair.	[number from 0-50] Default = <u>10</u>
MerLength	The minimum length of a mer required to be considered an exact match when scanning for repeats. Wizard equivalent: Advanced Trim/Scan Options > Mer length .	[number from 5-50] Default = <u>17</u>
MinEndFlagLen	The minimum length required for a mer to be flagged as a repeat if the segment is bound by the end of the read.	[number from 5-1000000] Default = <u>25</u>
MinFlagLength	The minimum length required for a mer to be flagged as a repeat. Wizard equivalent: Advanced Trim/Scan Options > Flag length .	[number from 5-1000000] Default = <u>50</u>
MinMerMatch	The minimum number of matching mers required to start an alignment. Wizard equivalent: Advanced Trim/Scan Options > Minimum matches .	[number from 2-25] Default = <u>2</u>

Example:

```
setRepeatParam merLength:17
```

```
setRepeatParam minMerMatch:2  
setRepeatParam maxMerGap:10  
setRepeatParam minFlagLength:50  
setRepeatParam alignCutoff:100  
setRepeatParam minEndFlagLength:25
```

setVectorParam

The **setVectorParam** command allows you to adjust the parameters used for vector trimming. In order to be applied, this command must appear in the script before the **loadVector** or **TrimVector** command, and the **vectScan** parameter for the **assemble** command must be set to 'true.'

Parameter	Description	Allowed values (defaults in bold)
AlignCutoff	The minimum acceptable alignment score. When the alignment score drops below the specified value, this indicates that the end of the alignment between the read and the vector has been reached, and the alignment will stop.	[number from 10-1000000] Default = <u>100</u>
EndCutOff	The distance to the endpoint where trimming will go all the way to the end of the sequence. Wizard equivalent: Advanced Trim/Scan Options > Trim to end.	[number from 0-1000000] Default = <u>25</u>
EndMerMatch	The minimum number of mer matches required to start an alignment in the specified end region.	[number from 1-25] Default = <u>1</u>
EndRegion	The number of bases at the end of a sequence where a lower stringency for matching and trimming is used.	[number from 0-1000000] Default = <u>15</u>
MaxMerGap	The maximum distance between two mers required to be considered a matching pair.	[number from 0-50] Default = <u>5</u>
MergeTrimGap	Maximum distance between two trim segments that will cause the segments to be merged. MergeTrimGap limits trimming to the ends of sequence reads, while EndCutOff doesn't. Controls how sensitive trimming should be in areas where some portions of the sequence match a vector and other portions don't. The higher the number the more likely the vector trimmer will find all the vector sequence in a region of poor	[number from 0-1000000] Default = <u>7</u> , which is

	quality. The smaller the number, the more confidence there is that the bases trimmed are actually vector and not a spurious match.	suitable for trimming linkers from the ends of sequences.
MerLength	The minimum length of a mer required to be considered an exact match when searching for vector. Wizard equivalent: Advanced Trim/Scan Options > Mer length .	[number from 5-25] Default = <u>9</u>
MinEndTrimLength	The minimum length to be trimmed when a vector matches the end of a read. This parameter can be useful in preventing small spurious matches from being trimmed, which may be significant with short read technologies.	[number from 5-1000000] Default = <u>5</u>
MinMerMatch	The minimum number of matching mers required to start an alignment. Wizard equivalent: Advanced Trim/Scan Options > Minimum matches .	[number from 1-25] Default = <u>3</u>
MinTrimLength	The minimum length required for a mer to be considered as a match for vector trimming. Wizard equivalent: Advanced Trim/Scan Options > Trim length .	[number from 5-1000000] Default = <u>30</u>

Example:

```
setVectorParam merLength:9
setVectorParam minMerMatch:3
setVectorParam MerGap:5
setVectorParam minTrimLength:30
setVectorParam minEndTrimLength:5
setVectorParam alignCutoff:100
setVectorParam endRegion:15
setVectorParam endCutoff:25
setVectorParam endMerMatch:1
```

Preprocessing and assembling commands

SNG “preprocessing and assembling commands” include:

- [appendToAssembly](#)
- [assemble](#)
- [convertReads](#)
- [extendContigs](#)
- [fixedTrim](#)
- [include](#)
- [makeSeqNamesUnique](#)
- [realignContigs](#)
- [removeSmallContigs](#)
- [set](#)
- [setAssemblyReport](#)
- [setPairSpecifier](#)
- [splitLinkerReads](#)
- [splitMIDSeqs](#)
- [splitPairs](#)
- [splitTemplates](#)
- [trimVector](#)

appendToAssembly

The **appendToAssembly** command is for the reference-guided workflow and is intended for internal use only.

assemble

The **assemble** command is required and reprocesses and assembles the sequences that have been loaded. Preprocessing may include quality trimming, and scanning for vector, repetitive, and contaminant sequences.

Parameter	Description	Allowed values
assembleBlocks	Specifies whether the assembly is a reference guided assembly.	[true / <u>false</u>]
contamScan	If true, sequences will be scanned for the specified contaminant sequences before assembling. Also see loadContaminant	[true / <u>false</u>]
doAssemble	If false, only the preprocessing will be done, and the sequences will not be assembled.	[<u>true</u> / false]
repeatScan	If true, sequences will be scanned for the specified known repetitive sequences before assembling. Also see loadRepeat	[true / <u>false</u>]
trimEnds	If true, the sequences will be trimmed based on quality scores before assembling. Wizard equivalent: Read options > Quality end trim .	[true / <u>false</u>]
vectScan	If true, the sequences will be scanned and trimmed for vector before assembling. Also see loadVector .	[true / <u>false</u>]

Example:

```
assemble
  trimEnds:false
  vectScan:false
  repeatScan:false
  contamScan:false
  doAssemble:true
```

convertReads

The **convertReads** command converts a sequence from one file format to another. This command is particularly useful for converting SOLiD .csfasta files into .fastq files that can be used by the XNG assembler.

Parameter	Description	Allowed values
destination	The location and filename for the output.	[directory/filename enclosed in quotes]
file, reads	The input file containing the reads.	[directory/filename enclosed in quotes]
format	Specifies the format of the output file. If 'genbank' is entered, the output will be in .gbk format. If fastq is entered, the output will be in .fastq format.	[genbank / fastq]

extendContigs

The **extendContigs** command is intended for internal use only.

Parameter	Description	Allowed values
extendPasses		[number]
mergeContigsInScaffold		[true/false]

fixedTrim

The **fixedTrim** command trims reads prior to assembly using fixed values. Based on the parameter settings for this command, SeqMan NGen will trim reads either by a specified number of bases from each end, or to a specified range.

Parameter	Description	Allowed values
end3	If trimRelative (see below) is set to 'true,' then this value indicates the number of bases for SeqMan NGen to trim from the 3' end of each read. If trimRelative is set to 'false,' then this value indicates the specific 3' coordinate to which reads should be trimmed. Wizard equivalent: Advanced Trim/Scan Options > 3' trim .	[number from 0-1000000] Default = <u>0</u>
end5	If trimRelative (see below) is set to 'true,' then this value indicates the number of bases for SeqMan NGen to trim from the 5' end of each read. If trimRelative is set to 'false,' then this value indicates the specific 5' coordinate to which reads should be trimmed. Wizard equivalent: Advanced Trim/Scan Options > 5' trim .	[number from 0-1000000] Default = <u>0</u>
trimRelative	Specifies whether the value for the end3 and end5 parameters should indicate the number of bases for SeqMan NGen to trim from the 3' or 5' end of each read. When 'false,' the value specified for the end3 or end5 parameter indicates the specific coordinate to which reads should be trimmed.	[<u>true</u> / false]

Example:

```
fixedTrim
  end5:10
  end3:20
  trimRelative:true
```

include

When building a script, the **include** command can be used to call up additional lines of script previously stored in a text file. In this way, a group of commands can be shared between two or more scripts.

Parameter	Description	Allowed values
file	Specifies a directory and name for the file.	[directory/filename enclosed in quotes]

makeSeqNamesUnique

The **makeSeqNamesUnique** command is intended for internal use only.

realignContigs

The **realignContigs** command causes SeqMan NGen to perform another pass through a reference-guided assembly once the initial assembly is complete, and realigns contigs as needed. (This step occurs automatically for *de novo* assemblies.) Using this command may improve the accuracy of the final assembly by correcting occasional misalignments that can occur in gapped regions, however note that this step may significantly increase the time to assemble. This command must appear in the script after the [assemble](#) command.

removeSmallContigs

The **removeSmallContigs** command disassembles any contigs without reference sequences that have fewer than the specified number of sequences.

Parameter	Description	Allowed values
minLength	Specifies the minimum length of a contig to prevent it from being disassembled. Wizard equivalent (<i>de novo</i> , special reference-guided workflows only): Assembly Options > Minimum length .	[number] Default = <u>0</u>
minSeqs	(required) Specifies the minimum number of sequences necessary in a contig to prevent it from being disassembled. Wizard equivalent (<i>de novo</i> , special reference-guided workflows only): Assembly Options > Minimum sequences .	[number] Default = <u>100</u>

set

The **set** command is used to set variables. See the example below and those under the [runScript](#) command.

Example:

```
set $snp:true  
set $snpMethod:"Diploid"
```

setAssemblyReport

The **setAssemblyReport** command is intended for internal use only. It is used to designate a file for a tab delineated report, similar to a report that XNG generates. This is useful during development to test how code changes impact results.

Parameter	Description	Allowed values
file name	Specifies the folder and file name.	[directory/filename enclosed in quotes]

setPairSpecifier

The **setPairSpecifier** command defines the paired end pair specifier for the paired Sanger and Illumina sequences in the assembly. This command must appear in the script before the assemble command, but after sequences have been loaded using the [loadSeq](#) command. For more information on assembling 454 paired end data, see the [load454PairedEnd](#) command. Pair specifiers define the naming convention for sequence pairs, as well as requirements for a minimum and maximum distance between the opposite ends of the inserts. Expressions for forward and reverse naming conventions should be created using the paired end specification language. Forward and reverse sequences must have identical names except for the unique portion that determines the direction of the clone.

Parameter	Description	Allowed values
pairs	This parameter lists the paired end constraints, specified by the following four values. Each value should be separated by a space and the list of values enclosed in double brackets {}. An additional set of brackets is required around all of the paired end constraints, regardless of whether one or multiple pair constraints are specified.	[forward reverse min max]
forward	A naming pattern to match forward clones.	[text string enclosed in quotes]
max	The maximum distance for the paired end sequences to be separated.	[number]
min	The minimum distance for the paired end sequences to be separated.	[number]
reverse	A naming pattern to match reverse clones.	[text string enclosed in quotes]

Example:

(defines 2 pair specifiers each with different size ranges)

setPairSpecifier

```
pairs:{{forward:"(*) (2kb) (*)-FP.*$" reverse:"(*) (2kb) (*)-RP.*$" min: 1500 max: 2500}
{forward:"(*) (8kb) (*)-FP.*$" reverse:"(*) (8kb) (*)-RP.*$" min: 7000 max: 9000}}
```

splitLinkerReads

The **splitLinkerReads** command splits specified reads based on their match to given linker sequences. Reads that align to the linker and include the linker site (as specified by the linkerSite parameter or by the cloneSite option in an .fof file) will be split into two reads. The two newly split reads will be designated by _A and _B appended to the name.

Parameter	Description	Allowed values
linkerFile	The directory and file name of the linker file.	[directory/filename enclosed in quotes]
linkerSite	The position indicating where reads should be split.	[number]
seqFile	The directory and file name of the sequence reads.	[directory/filename enclosed in quotes]

Example:

splitLinkerReads

seqFile: "/Library/123_project/reads.fas"

linkerFile: "/Library/123_project/linker.fas"

linkerSite:30

splitMIDSeqs

The **splitMIDSeqs** command is used to split 454 MID reads into individual files with one file per MID tag.

Parameter	Description	Allowed values
destination	The location and filename for the output.	[directory/filename enclosed in quotes]
file, reads	The input file containing the reads.	[directory/filename enclosed in quotes]

splitPairs

The **splitPairs** command is used to split 454 or ion torrent mate pair files into forward and reverse (and singleton) files.

Parameter	Description	Allowed value
destination	The location and filename for the output.	[directory/filename enclosed in quotes]
DiscardLinkerless	Specifies that reads without a linker sequence should be discarded from the assembly.	[true / <u>false</u>]
file, reads	The location and filename for the input.	[directory/filename enclosed in quotes]
seqTech	<p>Specifies the offset to be used when converting compressed quality scores into numerical values. These are the offsets used for the technology specified:</p> <p>Data Type / Value / Offset IonTorrent / IonTorrent / 33 Applied Biosystems SOLiD / SOLiD / 33 Illumina / Illumina / 64 Roche 454 / 454 / 33 Other types / normalScore / 33</p> <p>Wizard equivalent: Input Sequence Files > Read technology.</p> <p>Note 1: For 454, quality scores for homopolymeric runs of ≥ 2 are oriented from 5' to 3' on the top strand.</p> <p>Note 2: If possible, the data type of unknown data is determined automatically based on the first data file.</p>	[IonTorrent SOLiD Illumina 454 normalScore Other]

Example:

SplitPairs

destination:"c:data\splitReads\"

reads: {

{ file:"C:data\reads\file1.fas" format: IonTorrent }

{ file: "C:data\reads\file2.fas" format:454 discardLinkerless: true}

}

splitTemplates

The **splitTemplates** command splits reference contigs into multiple contigs in areas where there is zero coverage. Split contigs will be grouped into scaffolds with a defined position to allow for easy sorting when the project is viewed in SeqMan Pro / SeqMan Ultra. Annotations on the reference sequence will also be split, and any /codon_start qualifiers will be adjusted to stay in frame.

trimVector

The **trimVector** command is used for fast trimming vector sequence. Each read file is processed and the trimmed file is saved to the destination folder. If the file with the same name exists, the number will be appended to the file name. The file is saved in .fastq format, including trimming statistics.

Parameter	Description	Allowed values
file, reads	The location and filename for the input.	[directory/filename]
LinkerFile	The location of file or folder with vector sequence.	[directory/filename]
destination	The location of output folder.	[directory]

Example:

```
setVectorParam
  EndCutOff: 130
  MatchSize: 11
  MinTrimLength: 15
```

```
TrimVector
  reads: {
    file: "C:\data\input.fastq"
  }
  LinkerFile: "c:\data\adapter.fas"
  destination: "c:\data\Out\"
```

Specifying XNG or SNG/SMNG when running a script

SeqMan NGen utilizes [several assemblers](#) with different capabilities and scripting languages. Therefore, it is essential to match the correct assembler with the type of assembly project to be done.

To specify which assembler to use to run your script, type **xng** or **sng** followed by the path and script file name after the command prompt. Alternatively, add either the **#!/usr/bin/xng** or **#!/usr/bin/sng** command as first line of the script and execute through the command line.

Turn off usage logging

By default, [usage logging](#) is enabled in Lasergene version 11 and later. To opt out of usage logging, launch the DNASTAR Navigator and go to **View > Preferences > Lasergene** (Win) or **DNASTAR Navigator > Preferences > Lasergene** (Mac). Uncheck the box and click **Apply** and then **OK**.

Non-English keyboards

SeqMan NGen recognizes only standard English-keyboard characters as input. If you are using a non-English keyboard, we recommend that you switch to a “virtual” English keyboard. Click a link for instructions: [Windows 7 & 8](#), [Macintosh 10.11](#).

Installed Lasergene file locations

The following file names use 'x' to represent the version number.

File Category	Application	Path
Application ²	ArrayStar	Windows: C:\Program Files (x86)\DNASTAR\Lasergene x\ ArrayStar
	SeqNinja (command line)	Windows: C:\Program Files (x86)\DNASTAR\Lasergene x\ SeqNinjaCL
	All others	Windows: C:\Program Files (x86)\DNASTAR\Lasergene x Macintosh: /Applications/ DNASTAR/Lasergene x
Data Manager ² (<i>DMx</i> , <i>DMx.exe</i>)	SeqBuilder Pro, Protean 3D, GeneQuest, MegAlign Pro	Windows: C:\Program Files (x86) \DNASTAR\Lasergene x Macintosh: /Applications/ DNASTAR/Lasergene x
Data Manager State File	SeqBuilder Pro, Protean 3D, GeneQuest, MegAlign Pro	Windows: C:\Program Data\ DNASTAR\DataManager Macintosh: ~/Library/ Application Support/ DNASTAR/DataManager, ~/Library/Preferences/ DNASTAR/DataManager
License Manager ²	All	Windows: C:\Program Files (x86)\DNASTAR\License Manager Macintosh: /Applications/ DNASTAR/License Manager
Commuter License Manager	ArrayStar	Windows: C:\Program Files (x86)\DNASTAR\ArrayStar 'x'

		Macintosh: /Applications/ DNASTAR/ArrayStar 'x'
	SeqMan NGen	Windows: C:\Program Files (x86)\DNASTAR\SeqMan NGen 'x' Macintosh: /Applications/ DNASTAR/SeqMan NGen 'x'
	All others	Windows: C:\Program Files (x86)\DNASTAR\Lasergene 'x' Macintosh: /Applications/ DNASTAR/Lasergene 'x'
Server License File (<i>/servrc</i>), Server License Manager ² , <i>Server Executables</i> (<i>_lserv</i> , <i>lservnt.exe</i>)	All	Windows: C:\Program Files (x86)\DNASTAR- LicenseServer\Server Macintosh: ~/Library/ DNASTAR-LicenseServer
Standalone & Trial Licenses (<i>*.license</i>), License Server Client License (<i>*.lshost</i>), Key Server Client License (<i>*.keyhost</i>)	All	Windows: C:\Program Data\ DNASTAR\Licenses Macintosh: ~/Library/ Application Support/ DNASTAR/Licenses
Preferences	Protean 3D, Navigator, SeqNinja (DNA*), GenVision Pro	Windows: C:\Users\<user>\DNASTAR Macintosh: ~/Library/ DNASTAR
	ArrayStar	Windows: C:\Users\<User>\AppData\ Roaming\DNASTAR\ArrayStar
	All others	Windows: C:\Users\<user>\AppData\ Local\DNASTAR\ Macintosh: ~/Library/ Preferences



AppData is a hidden folder in Windows. To unhide the folder, go to **Organize > Folder and Search Options > View > Show Hidden files and folders**.

Troubleshoot failure to launch

If you attempt to launch SeqMan NGen, without an updated .NET Framework Service Pack 1 (Windows) or Mono package (Mac) installed, you may receive an error message, and SeqMan NGen may fail to open.

The best way to ensure the correct software is installed is to perform all recommended updates when prompted by your operating system.

To instead resolve this issue manually:

- Windows – Download and install .NET Framework 4.7.2 from www.microsoft.com/net/download/windows/run.
- Macintosh – Download and install the Mono 5.10.1 stable package from www.mono-project.com/download/stable/#download-mac.

Research references

Benjamini Y and Hochberg Y (1995). "Controlling the false discovery rate: a practical and powerful approach to multiple testing." *Journal of the Royal Statistical Society. Series B (Methodological)*, Vol. 57, No. 1 (1995), pp. 289-300. ([See online.](#))

Dong C, Wei P, Jian X, Gibbs R, Boerwinkle E, Wang K and Liu X (2015). "Comparison and integration of deleteriousness prediction methods for nonsynonymous SNVs in whole exome sequencing studies." *Human Molecular Genetics* 24(8):2125-2137. (Pertains to dbNSFP; [see online.](#))

Fitzgerald DM, Bonocora RP, Wade JT (2014). "Comprehensive Mapping of the Escherichia coli Flagellar Regulatory Network." *PLoS Genet* 10(10): e1004649. doi:10.1371/journal.pgen.1004649.

Johnson et al. (2007). "Genome-Wide Mapping of in Vivo Protein-DNA Interactions." *Science* June 8, 2007; 316(5830):1497-502.

Krumm N, Sudmant PH, Ko A, et al. (2012). "Copy number variation detection and genotyping from exome sequence data." *Genome Res.* published online May 14, 2012. (Pertains to the zRPKM and RPKM-CN normalization methods.)

Li H, Ruan J, and Durbin R (2008). "Mapping short DNA sequencing reads and calling variants using mapping quality scores." *Genome Res.* 2008 Nov;18(11):1851-8. doi.org/10.1101/gr.078212.108. Epub 2008 Aug 19.

Liu X, Jian X, and Boerwinkle E (2011). "dbNSFP: a lightweight database of human non-synonymous SNPs and their functional predictions." *Human Mutation*. 32:894-899. (Pertains to dbNSFP; [see online.](#))

Liu X, Wu C, Li C and Boerwinkle E (2016). "dbNSFP v3.0: A One-Stop Database of Functional Predictions and Annotations for Human Non-synonymous and Splice Site SNVs." *Human Mutation*. 37:235-241. (Pertains to dbNSFP; [see online.](#))

Love MI, Huber W and Anders S (2014) "Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2." *Genome Biology* 15:550. doi.org/10.1186/s13059-014-0550-8.

Mortazavi A, Williams BA, McCue K, Schaeffer L and Wold B (2008). "Mapping and quantifying mammalian transcriptomes by RNA-Seq." *Nature Methods*, 5, 621-628. (Pertains to the ChIP-Seq Peak Finder algorithm.)

Mukherjee, Supratim & Huntemann, Marcel & Ivanova, Natalia & Kyripides, Nikos & Pati, Amrita. (2015). Large-scale contamination of microbial isolate genomes by Illumina PhiX control. *Standards in genomic sciences*. 10. 18. 10.1186/1944-3277-10-18. ([See online.](#))

Nookaew I et al. (2012). "A comprehensive comparison of RNA-Seq-based transcriptome analysis from reads to differential gene expression and cross-comparison with microarrays: a case study in

Saccharomyces cerevisiae.” Nucleic Acids Res. 2012 Nov 1;40(20):10084-97. doi: 10.1093/nar/gks804 ([View on PubMed.](#))

Robinson MD, McCarthy DJ, Smyth GK (2010). “edgeR: a Bioconductor package for differential expression analysis of digital gene expression data.” Bioinformatics, 26(1), 139-140.

Robinson, MD, and Oshlack, A (2010). “A scaling normalization method for differential expression analysis of RNA-seq data.” Genome Biology 11, R25

Robinson, MD, and Smyth, GK (2008). “Small sample estimation of negative binomial dispersion, with applications to SAGE data.” Biostatistics 9, 321 – 332.

Zhang Y, Liu T, Meyer CA et al. (2008) Model-based Analysis of ChIP-Seq (MACS). Genome Biol 9, R137. ([See online.](#))