

NextGENe®

Next Generation Sequencing Software for Biologists



Application Modules for:

Variant Analysis, Targeted, WES, WGS

- SNPs
- Indels
- Structural Variations
- Somatic Mutation Mining

Copy Number Variation (CNV) Analysis

- CNV Tool
- MLPA-like Analysis with Batch CNV Tool
- Sensitive Aneuploidy Detection (SAD) Tool

HLA Analysis

Patient Comparison

- Groups
- Families
- Trios

Functional Prediction Calling

RNA-Seq/Transcriptome Analysis

- Digital Gene Expression
- miRNA Discovery and Quantification

Assembly

- de novo*
- Paired end

SOFTGENETICS®

Software PowerTools for Genetic Analysis

NextGENe[®]

Benefits

Instant Knowledge...

Single Annotated data review screen

Increased Accuracy...

Unique technologies increase accuracy by

- ✓ Removing sequencing errors
- ✓ Platform specific technologies

Compatible with all major sequencing systems

Track Manager

- ✓ Functional Prediction
- ✓ Disease Association
- ✓ Conservation Scores
- ✓ Population Frequencies

Automated format conversion tool

Low-Cost Hardware Requirements

Biologist Friendly Windows[®] interface...

- ✓ Application driven
- ✓ Automated inspection of input files to set analysis parameters
- ✓ Requires no scripting
- ✓ Reduces bioinformatics requirements
- ✓ Unattended batch processing capabilities

Annotated results in single easy-to-navigate view...

- ✓ Automated pipeline tool speeds analysis
- ✓ Multiple integrated, exportable reports
- ✓ Analysis filters

Automated Analysis Pipeline

- ✓ Automated linkage to Geneticist Assistant™ NGS Interpretive Workbench
- ✓ Custom Templates for streamlined setup of batch processing

NextGENe Software is a complete, "free-standing" analysis package designed for use by biologists in the analysis of data from Next Generation Sequencing systems. The icon driven, easy-to-use Windows[®] interface significantly reduces bioinformatics requirements, provides annotated analysis review, while reducing sequencing errors to improve analysis accuracy and speed.



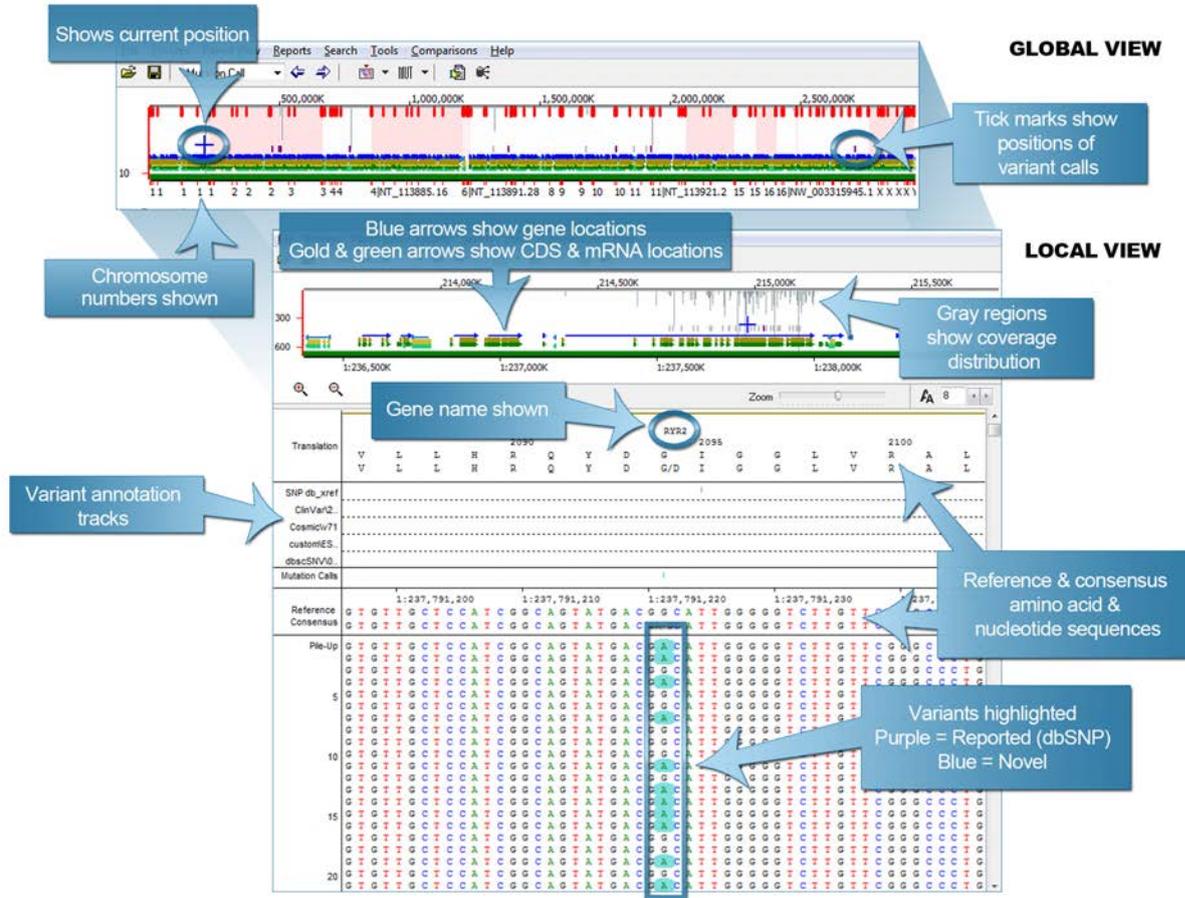
**Bred to
Track!**

Features

Instant Knowledge

- Annotation
- Easy Navigation
- Exportable

NextGENe's analysis browser provides a highly interactive review of annotated analysis results in a single view. Navigation is as simple as drawing "boxes" to zoom in or out, graphics and text reports are hyperlinked to speed data review and "hot" keys ease navigation.



Example of annotated Whole Human Genome data review with NextGENe browser. Navigation is simple either using Hot Keys or by dragging mouse over screen to move across the genome or zoom in on selected areas. Text Reports are linked to browser for quick, easy data review.

Mutation Report is hyperlinked to graphical NextGENe browser and databases including dbSNP and COSMIC. Several filtering options are available to speed and ease analysis review:

Index	Chromosom	Gene	CDS	Chr	Referer	Covered	PhyloP	PolyPh	SIFT	Mutatio	LRT	1000Ge	Score	A Retri	C Ret	G Ret	T Ret	Ins Ret	Del Ret	SNP db_xref	Mutation Call	Amino Acid Change	
1	2337277	PEX10	1	C		318							16.5	0.00	50.31	0.00	49.69	0.00	0.00	rs11586985	C>CT		
2	2340073	PEX10	3	1	C	134	N	B	T	N	N	NA	16.9	0.00	52.24	47.76	0.00	0.00	0.00	rs76530653	C>CG	140G>GR	
3	20972048	PINK1	1	G		255							18.7	99.61	0.39	0.00	0.00	0.00	0.00	rs3131713	G>A		
4	21904131	ALPL	11	1	T	253	N	NA	D	N	N	0.1083	8.0	0.00	35.97	0.40	51.78	0.00	11.86	rs34605986	T>CT	522V>AV	
5	41296828	KCNQ4	10	1	T	433	N	B	T	D	N	0.1042	18.5	0.23	0.00	46.19	53.58	0.00	0.00	rs34287852	T>GT	455H>HQ	
6	94512565	ABCA4	19	1	C	446	C	D	T	D	N	0.0417	18.7	0.00	48.65	0.00	51.35	0.00	0.00	rs1801581	C>CT	943R>RQ	
7	100672060	DBT	9	1	T	397	C	NA	T	P	N	0.8659	20.5	0.25	99.75	0.00	0.00	0.00	rs12021720	T>C	384S>G		
8	103354138	COL11A1	62	1	A	455	C	NA	B	T	P	N	0.8156	16.0	52.31	0.00	45.49	0.00	0.22	2.20	rs1676486	A>AG	1547S>SP
9	116247826	CASQ2	9	1	T	494	C	P	D	D	NA	21.0	0.00	49.60	0.00	50.20	0.00	0.20	rs72703607	T>CT	309D>GD		
10	116310967	CASQ2	1	1	T	360	N	B	T	P	N	0.4274	17.5	0.28	48.33	0.00	51.39	0.56	0.00	rs4074536	T>CT	66T>AT	
11	171076966	FMO3	3	1	G	472	C	B	T	P	N	0.3464	20.2	48.09	0.00	50.64	0.00	0.85	1.27	rs2266782	G>AG	158E>KE	
12	201331088	TNNT2	12	1	A	285	C	D	D	D	NA	15.9	55.79	0.00	44.21	0.00	0.35	0.00	rs45520032	A>AG	228I>IT		
13	215844373	USH2A	63	1	C	747	C	D	T	N	U	NA	21.5	0.00	48.06	0.00	51.94	0.00	0.00	rs45549044	C>CT	4692G>GR	
14	215914826	USH2A	59	1	T	282	N	B	T	N	N	0.1285	19.3	0.00	47.87	0.00	52.13	0.00	0.00	rs35309576	T>CT	386M>VM	

Optimized Analysis Performance

NextGENe has been optimized to perform on any Windows® 64 bit operating system from Vista through Windows 10, or Windows based server 2008R2 and forward. NextGENe also performs well on Mac® hardware when utilizing VMware, boot camp or similar boot managers along with Windows OS. NextGENe is multi-threaded, utilizing all available processors.

Performance examples of alignment and variant detection on various fasta data sets.

Computer Type	Data Set Type	Number of Cores Used	Total Analysis Time	Million Reads per Hour
Intel® Core™ i7 2.8 GHz 8 core, 16GB RAM	10 GB Illumina, single end 100 bp reads	6	1 hour	61
Intel® Core™ i7 2.8 GHz 8 core, 16GB RAM	1GB Illumina single end 100bp reads	6	5 Min	119
Intel® Core™ i7 2.8 GHz 8 core, 16GB RAM	5GB x2 Illumina paired end 100bp reads	6	1.4 hours	48
Intel® Xeon™ 2.3 GHz 16 core, 60GB RAM	10GB x2 Illumina paired end 100bp reads	14	2.6 hours	46
Intel® Xeon™ 2.3 GHz 16 core, 60GB RAM	0.8 GB Ion PGM AmpliSeq™, 117bp reads	4	2 Min	19
Intel® Xeon™ 2.3 GHz 16 core, 60GB RAM	4.5GB Ion Torrent Proton WES 126bp reads	14	1.3 hours	26
Intel® Xeon™ 2.3 GHz 16 core, 60GB RAM	4.5 GB Ion Torrent Proton WES 126 bp reads	4	2.7 hours	13

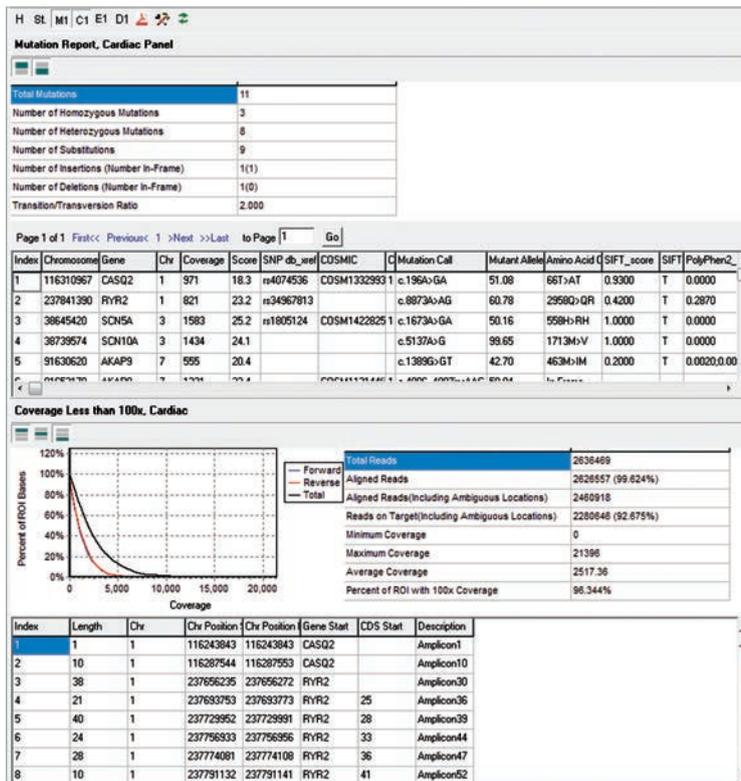
Reporting of Results and Quality Control Metrics

- Customizable
- Various Formats including VCF & SIFT

NextGENe produces several reports, including Mutation Report, Coverage Curve Report, Distribution Report, Expression Report, Paired Read Reports, and more. Each report can be saved in various formats. For example, the Mutation Report can be saved in VCF, SIFT or a tab delimited text format with your choice of columns of information. NextGENe Viewer's main display contains several panes of information. The reporting pane can display several different reports, one of which is the Summary Report. The Summary Report can show several statistics and reports in a single view, wrapping up the results of a single project to show quality, coverage, mutation calls and more in one report. The Summary Report is tied in to the Post Processing step of the Project Wizard, allowing you to set up each report's configuration prior to the analysis. The Summary Report can be customized after the project is completed also, allowing you to tailor the report to each project's specific needs.

NextGENe software's custom report builder allows users to create reports that concentrate only on the information required by their application or institution. In the above example the report contains total found mutations, type, chromosome position, gene, coverage levels by position, confidence score, dbSNP ID, COSMIC reference, mutation call, mutant allele %, amino acid change, SIFT and PolyPhen2 scores, as well as run quality data including read balance, total and total aligned reads, reads on target, percent of ROI covered and more.

Customized Clinical Research Report



Mutation Confidence Scoring

- Overall mutation confidence score provided for every mutation
- Any penalty score can be disabled
- Quickly view the distribution of scores in a project
- Filter based on the overall score or on penalty sub-scores

NextGENe software includes a proprietary mutation confidence scoring system designed to make it easier to find the called mutations that are most likely to represent true variations. The overall score is the product of the coverage score and several penalty sub-scores:

- **Coverage score** – starts at 0 and has no upper limit, but is rarely higher than 32. It is calculated as $8 * \log_{10}$ (adjusted coverage). The adjusted coverage gives greater weight to the higher quality 5' end of reads and less weight to the lower quality 3' end.

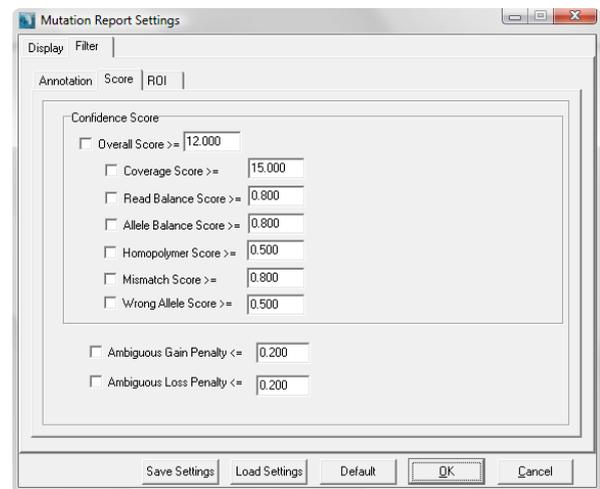
- **Read Balance Score (0 to 1)** – A score of 1 indicates perfect or near perfect balance between the number of forward and reverse reads. Unbalanced data may indicate misalignment or may allow basecalling biases to cause false positives. If the data is expected to be biased (as it is for some targeted sequencing applications) then this score should be disabled.

- **Allele Balance Score (0 to 1)** – Measures the major and minor allele balance and compares it to the read balance. The calculation is similar to a chi-square test. If the allele balance is different from the read balance, then there is strong evidence that the mutation may be an error.

- **Homopolymer Score (0 to 1)** – Penalizes indels occurring in homopolymer regions for data that has that error profile.

- **Mismatch Score (0 to 1)** – Penalizes mutations when many mismatches occur in a small area. This usually indicates untrimmed adapter sequence or misalignment.

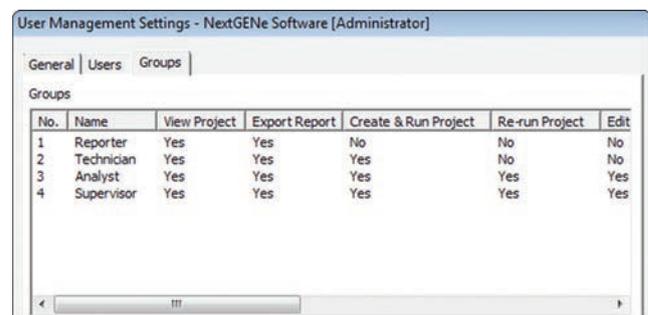
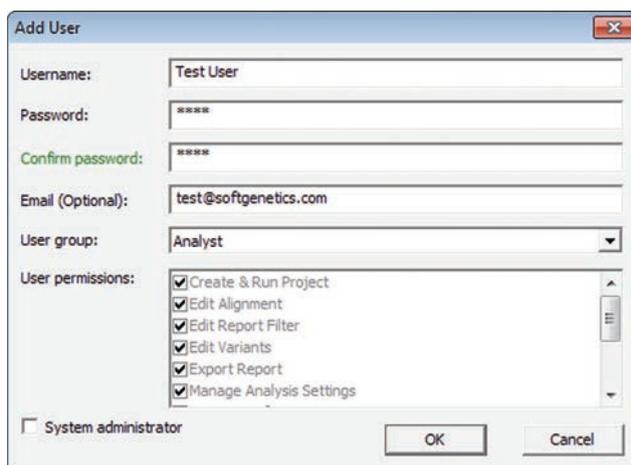
- **Wrong Allele Score (0 to 1)** – Penalizes mutations when a third allele is found which makes more than one possible mutation call possible (such as 60% A, 20%C, 20%T). This score is especially helpful for targeted capture data.



User Management

- Enable to require user login
- Create and manage permissions for users and/or groups
- Audit trail records username of who created or modified each project and/or report

NextGENe's new User Management function provides the ability to control access for different users or groups of users and to track the user who created or modified any project or report. When enabled, user login is required to open NextGENe.



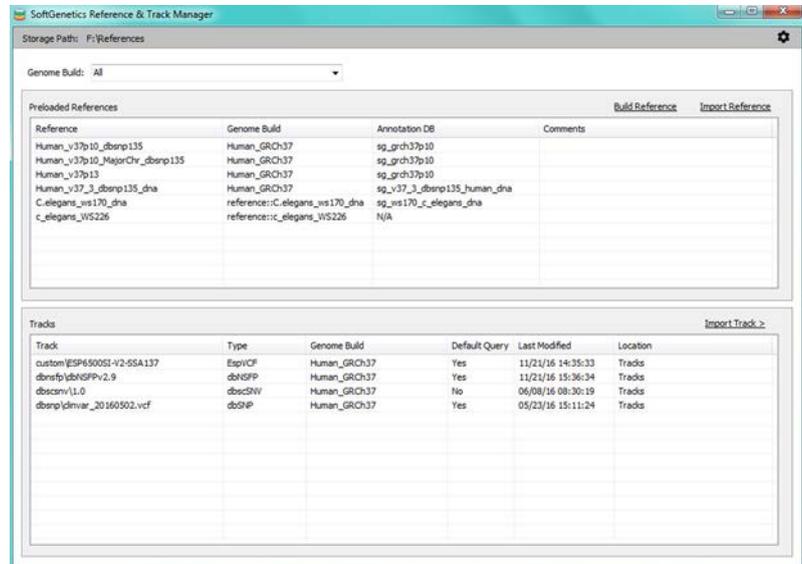
User accounts can be created for each user, with permissions assigned on an individual basis, or groups can be utilized to assign the same privileges to multiple users.

Variant Annotation Tracks and Pathogenicity Calling

- **Functional Prediction information**
SIFT, PolyPhen-2, LRT, MutationTaster, FATHMM, CADD & MutationAssessor
- **Disease association**
ClinVar & COSMIC
- **Conservation scores**
phyloP, GERP++, phastCons & SiPhy
- **Population frequencies**
1000 Genomes and Exome Variant Server
- **COSMIC - Catalogue of Somatic Mutations in Cancer from Sanger Institute**
- **dbSNP from NCBI**
- **Custom Tracks**

NextGENe's Reference and Track Manager Tool can be used to import variant databases as tracks for annotation. Information from imported tracks can be displayed and used for filtering in the Mutation Report. The Reference and Track Manager includes specialized support for 3 popular databases: the dbNSFP database, which contains pre-calculated functional prediction and conservation scores along with population frequencies from the 1000 Genomes and Exome Variant Server projects, the COSMIC database of cancer variants from the Sanger Institute and the dbSNP database from NCBI. There is also a custom track import feature that can be used to import other public or proprietary databases.

NextGENe's Reference and Track Manager can be used to easily import dbNSFP, COSMIC, dbSNP as well as custom databases as annotation tracks.

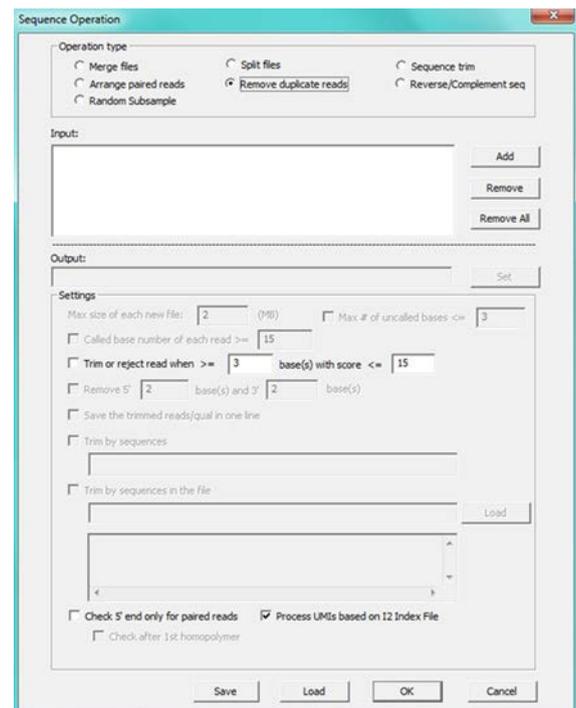


Using Unique Molecular Identifiers (UMIs) to Remove PCR Duplicates

The removal of PCR duplicates provides increased accuracy in determining allele frequencies for improved variant calling. NextGENe Software's Sequence Operation Tool includes the ability to remove PCR duplicates using unique molecular identifiers (UMIs) such as those used NEBNext or HaloPlex chemistries. The UMIs are random sequences of bases that are used to tag each molecule (fragment) of DNA prior to library amplification, which aids in the identification of PCR duplicates. Illumina instruments generate an I2 Index File for paired-end runs, which can store these UMIs.

NextGENe's Sequence Operation Tool uses UMIs within the Illumina I2 files to identify the PCR duplicates. NextGENe identifies all the read pairs that share the same UMI and retains only the pair that has the highest total quality to be processed along with unique paired reads, while duplicate reads are removed from further processing.

NextGENe's Sequence Operation Tool can be used to remove PCR duplicates using UMIs from an I2 index file.

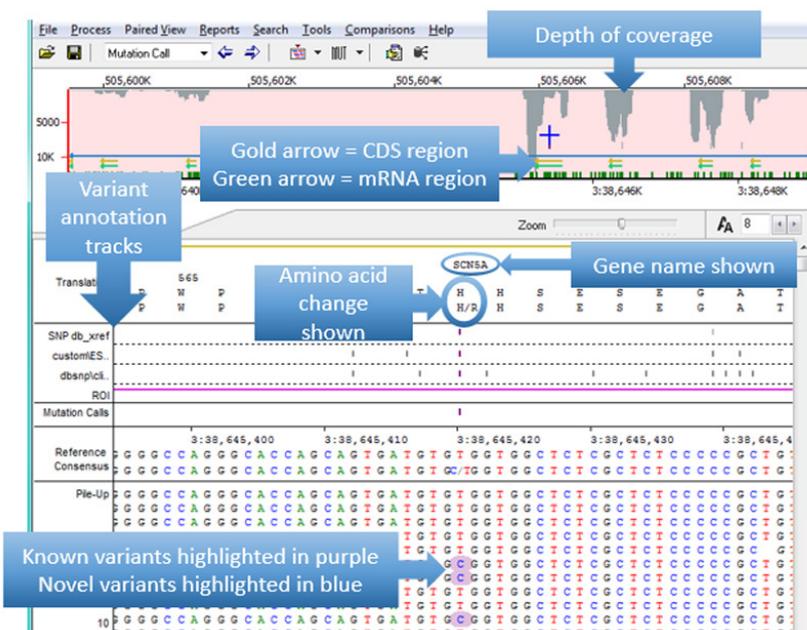


Applications

SNP/INDEL Detection

- SNP Detection
- Indel Detection (up to the read length)
- Low False Positive Rate
- Biologist friendly reporting
- Export results in spreadsheet form or VCF format to database or LIMS system
- Scoring of Variants

SNP's and small Indels, up to the read length, can be detected in sequencing data from all sequencing technologies. Use of the Condensation® Tool elongates short reads, increasing read uniqueness probability in the genome, while polishing the data to remove chemistry and instrumental errors. NextGENE software automatically calculates a confidence score for each found variant.



In the region of aligned sequence reads, novel variant calls are highlighted in blue, previously reported in purple.

The Whole Genome Pane is located at the top of the display – coverage is indicated by gray regions, blue tick marks identify the location of novel SNPs, previously reported SNPs are indicated in purple.

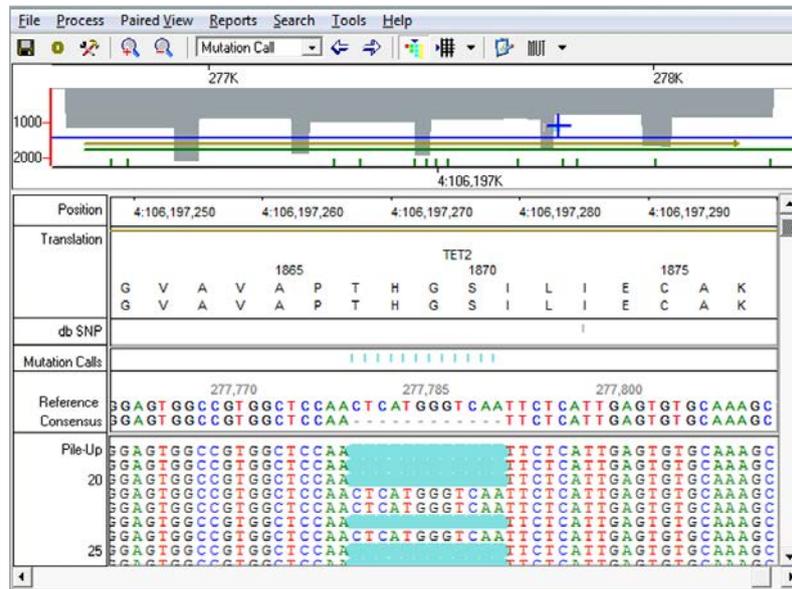
Index	Chromosom Position	Gene	CDS	Chr	Referenc Nucleo	Coverage	PhyloP Classif	PolyPh Classif	SIFT Class	Mutatio Class	LRT Clas	1000Ge	Score	A Rati	C Rati	G Rati	T Rati	Ins Rati	Del Rati	SNP db_xref	Mutation Call	Amino Acid Change
1	2337277	PEX10		1	C	318							16.5	0.00	50.31	0.00	49.69	0.00	0.00	rs11586985	C>CT	
2	2340073	PEX10	3	1	C	134	N	B	T	N	N	NA	16.9	0.00	52.24	47.76	0.00	0.00	0.00	rs76530653	C>CG	140G>GR
3	20972048	PINK1		1	G	255							18.7	99.61	0.39	0.00	0.00	0.00	0.00	rs3131713	G>A	
4	21904131	ALPL	11	1	T	253	N	NA	D	N	N	0.1083	8.0	0.00	35.97	0.40	51.78	0.00	11.86	rs34605986	T>CT	522V>AV
5	41296828	KCNQ4	10	1	T	433	N	B	T	D	N	0.1042	18.5	0.23	0.00	46.19	53.58	0.00	0.00	rs34287852	T>GT	455H>HQ
6	94512565	ABCA4	19	1	C	446	C	D	T	D	N	0.0417	18.7	0.00	48.65	0.00	51.35	0.00	0.00	rs1801581	C>CT	943R>RQ
7	100672060	DBT		9	T	397	C	NA	T	P	N	0.8659	20.5	0.25	99.75	0.00	0.00	0.00	rs12021720	T>C	384S>G	
8	103354138	COL11A1	62	1	A	455	C	B	T	P	N	0.8156	16.0	52.31	0.00	45.49	0.00	0.22	2.20	rs1676486	A>AG	1547S>SP
9	116247826	CASQ2		9	T	494	C	P	D	D	D	NA	21.0	0.00	49.60	0.00	50.20	0.00	0.20	rs72703607	T>CT	309D>GD
10	116310967	CASQ2		1	T	360	N	B	T	P	N	0.4274	17.5	0.28	48.33	0.00	51.39	0.56	0.00	rs4074536	T>CT	66T>AT
11	171076966	FMO3		3	G	472	C	B	T	P	N	0.3464	20.2	48.09	0.00	50.64	0.00	0.85	1.27	rs2266782	G>AG	158E>KE
12	201331068	TNNT2		12	A	285	C	D	D	D	D	NA	15.9	55.79	0.00	44.21	0.00	0.35	0.00	rs45520032	A>AG	228I>IT
13	215844373	USH2A		63	C	747	C	D	T	N	U	NA	21.5	0.00	48.06	0.00	51.94	0.00	0.00	rs45549044	C>CT	4692G>GR
14	215914826	USH2A		59	T	282	N	B	T	N	N	0.1285	19.3	0.00	47.87	0.00	52.13	0.00	0.00	rs35309576	T>CT	3868M>VM

A Mutation Report was generated for the run, showing a list of all variant calls. Calls can be manually reviewed, and can be edited, deleted or added. Options are available for customizing the view of this information, in addition to further filtering. The calls within this report are organized by position within the reference, and each line contains the position within reference, the reference nucleotide, coverage, causative prediction by several databases, 1000 genomes frequency percentages for each allele found, amino acid changes, gene and/or chromosome location and dbSNP identification.

Analysis of Targeted Sequencing Panels

NextGENe software is able to very rapidly process samples from Ion Torrent™, and Illumina platforms in order to find potentially important mutations in AmpliSeq™, HaloPlex™, Multiplicom™, NEBNext™, RainDance™, Roche GS G Type Assay and all other targeted sequencing panels. NextGENe software includes sorting and trimming tools, alignment, and mutation calling on low cost Windows PCs. NextGENe software is also able to annotate the mutations that were found using dbSNP, the dbNSFP database, the COSMIC database, or custom databases. Alignments and variant calling is typically accomplished in minutes, with pre-alignment processing such as barcode sorting taking only a few minutes. All of these steps can be fully automated in order to make processing samples even faster and easier.

Trio results are shown in NextGENe's comprehensive viewer which also allows comparison of up to 20 individual samples and prediction filtering:



A 12 bp deletion in TET2 detected at approximately 41% frequency from Roche GS G Type TET2/CBL/KRAS Panel. Data courtesy of 454 Life Sciences.

Trio Analysis with NextGENe Viewer



Trio results are shown in NextGENe's comprehensive viewer which also allows comparison of up to 20 individual samples and prediction filtering.

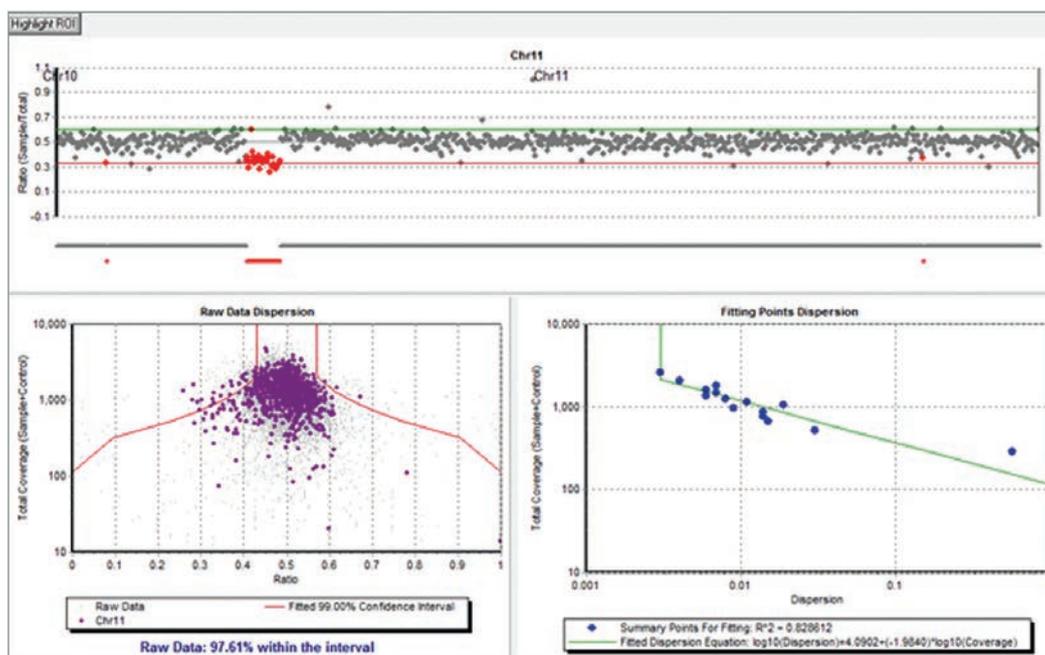
Copy Number Variation Analysis Tools

CNV Tool

- Two powerful CNV algorithms available
- Provides confidence scores for CNV calls
- Classifies regions as potential deletions, duplications, or normal copy number
- Processes individual CNV changes (per exon, or per amplicon) into multi-locus calls that can be made more confidently.

NextGENe includes a new Copy Number Variation (CNV) analysis tool designed to make variant calls on a case-control basis. Two powerful CNV detection algorithms are included, using proprietary technology. The “SNP-Based” method uses a global coverage normalization and uses SNP locations to determine the median coverage for each region in both the sample and control. A log₂ ratio is used to compare the median coverage values to make the CNV calls. The “Dispersion and Hidden Markov Model (HMM)” method utilizes coverage ratios between the sample and control, models the dispersion (noise) at varying coverage levels, and uses this information in conjunction with a Hidden Markov Model to make the CNV calls.

The CNV Tool works well with targeted sequencing data such as Ion AmpliSeq™ panels or the HaloPlex™ Target Enrichment System from Agilent Technologies or consistent depth of coverage whole Exome Sequencing data. Ideally the two samples used in a comparison will be as close as possible in experimental conditions. No special processing is needed to use the tool - any aligned NextGENe projects can be utilized for CNV analysis.



Specialized graphics are displayed for CNV analysis. This figure shows the results view with chromosome 11 selected with a deletion indicated in red

Sample	E.pjt										
Control	F.pjt										
Index	Description	Chr	Chr Start	Chr End	Gene	CDS	Length	Log2 Ratio	Score	Original Coverage (Sample:Control)	Normalized Coverage (Sample:Control)
1	Amplicon206	chr7	91623915	91624109	AKAP9	• 6	195	0.90	3.72	100.55	100.54
2	Amplicon255	chr7	150645512	150645650	KCNH2	- 11	139	-1.03	8.26	197.413	197.402
3	Amplicon358	chr19	35521705	35521783	SCN1B	• 1	79	-2.73	11.00	38.258	38.251

Figure indicates the detection of a known deletion in the KCNH2 gene using HaloPlex™ Cardiac Panel data. The log₂ ratio (-1.03) was very close to the expected value of -1.0.

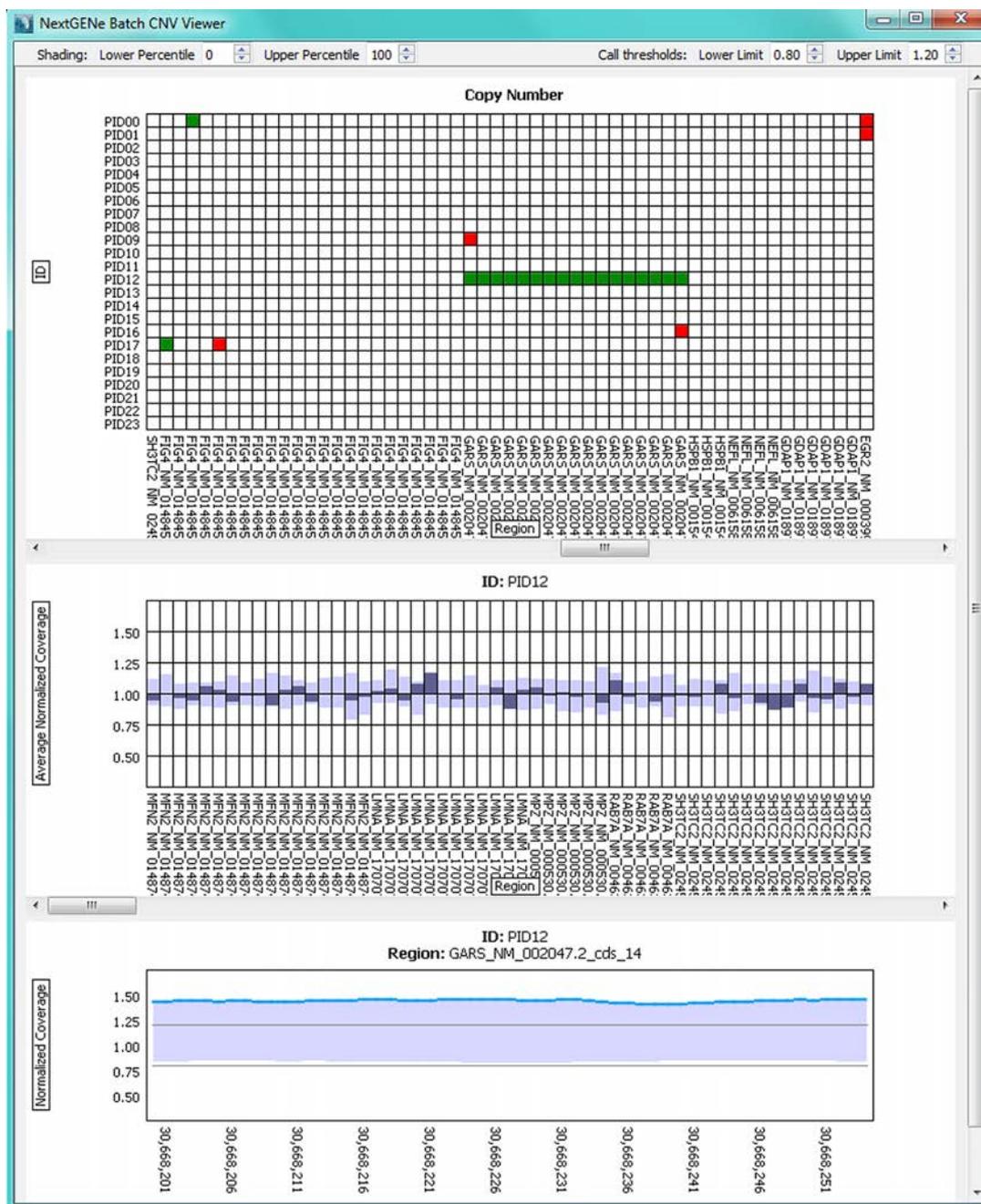
Batch CNV Tool for MLPA-like CNV Detection

The use of NGS and NextGENe software's Batch CNV Tool for screening dramatically increases laboratory throughput while significantly reducing the cost of MLPA Analysis. NextGENe software includes a Batch CNV tool that provides a means to replace Multiplex Ligation-dependent Probe Amplification and similar techniques, typically performed by capillary electrophoresis (CE) fragment analysis, with High-Throughput or Next Generation Sequencing. HTS approach to CNV does not use the problematic steps of the traditional MLPA reaction, which are challenging to optimize.

The tool calculates the pileup coverage at each position in the targeted panel, normalizes it by a constant value in each sample, and reports the normalized values.

The normalization procedure is:

1. Get the total coverage in each sample (adding the coverage at each position in the ROI bed file).
2. Divide these values by the average across all samples to get the normalization factor for each sample.
3. Divide the coverage at each position by the normalization factor for that sample.
4. When reporting a value for a given region, reports the average of all bases within that region.



Global plot with samples in rows and regions in columns. Red squares indicate a deletion call while green squares indicate a duplication call

Sample-specific plot shows coverage distribution of an individual sample in dark gray. Light gray shows the distribution of coverage values for the percentile range defined above

Region-specific plot shows the sample's coverage value across the selected region (blue line) and the distribution of coverage values for the percentile range defined above (light gray region)

Sensitive Aneuploidy Detection (SAD) Tool

The Sensitive Aneuploidy Detection (SAD) tool is designed for analysis of low coverage whole genome sequencing samples in order to detect whole-chromosome aneuploidies. The tool is streamlined for very simple operation and the normalization process allows for detection of aneuploidies even with a high background of non-aneuploid data. A set of at least 30 normal samples (“modeling controls”) is used to set up an analysis model which automatically accounts for biases in coverage across the chromosomes. The model is based on comparing “ratio” values for each tested chromosome to a normal distribution of these “ratio” values in the controls, resulting in normalized chromosome values (NCVs) similar to z-scores. Cutoffs are used to make calls based on the NCVs. The model can be tested for sensitivity and specificity using samples with known genotypes.

One possible application of this tool is for non-invasive prenatal testing (NIPT), based on sequencing of maternal plasma samples to detect fetal aneuploidies.

The tool is used in three stages:

1. Optimize the normalization technique and set up baseline values using the negative controls.
2. Optionally test positive and negative controls to calculate sensitivity and specificity. Use this information to fine-tune cutoff values for making calls based on the Normalized Chromosome Values (NCVs).
3. Test unknown samples one batch (run) at a time. Every sample will have results reported for all tested chromosomes including inconclusive calls if the results were in between the cutoffs for a normal call and an aneuploid call.



Sensitive Aneuploidy Detection (SAD) Tool



Somatic Mutation Comparison Tool

A new “Somatic Mutation Comparison Tool” has been added to the NextGENe Viewer tools menu. This specialized version of the variant comparison tool allows the user to load three exome projects- a tumor sample, a matched normal sample (such as a blood sample), and a pooled sample consisting of normal data from multiple sources. The tumor and normal projects are used to find somatic variants and the pooled project is used to eliminate artifacts due to library preparation and alignment.

Available settings include:

- **“Maximum Contamination”** – somatic variants may have up to this allele frequency in the normal project. This is to account for contamination of the normal sample with tumor sequence. Default setting is 5%. Any allele ratio greater than 5% in normal will be considered as germline mutation.
- **“Number of Pooled Samples”** – The number of samples used in the pooled project. The maximum contamination setting is divided by the square root of this number and used the same way, but for the pooled project instead of the normal project.
- **“Somatic Allele Count”** – The tumor sample must have at least this number of reads containing the variant. Default is 5 counts.
- **“Directional Balance Ratio (T/N)”** – The forward/reverse ratio in the tumor project divided by the same ratio in the normal project must not exceed this value. The default value is between(1/7x, 7x).
- **“Somatic Allele Frequency Ratio (T/N)”** – The ratio of the mutant allele percentage between the tumor and normal projects must be greater than this value. Default is set to 3.
- **“Pooled Allele Count Ratio (T/P)”** – This optional filter requires that the ratio of the number of reads with the mutant allele between the tumor and the pool projects must be greater than this value. Default is 3.

The screenshot displays the Somatic Mutation Comparison Tool interface. At the top, there are three alignment tracks for different samples: T_Unique_wgs, N_L003_Unique_wgs, and T1352L_Unique_wgs. Each track shows sequence alignment with a reference consensus and pile-up views. Below the tracks, a table lists identified variants with columns for ID, Chr, Position, Gene, CDS, Chr, Reference, Category, Change, SNP db, Coverage, and various ratios (A, C, G, T, In, De, A, C, G, T, In, De, Read Balance, Mutation Call, Mutant Allele, and Annotation). The table is currently showing page 1 of 1.

ID	Chr	Position	Gene	CDS	Chr	Reference	Category	Change	SNP db	Coverage	A Ratio	C Ratio	G Ratio	T Ratio	In Ratio	De Ratio	Read Balance	Mutation Call	Mutant Allele	Annotation					
1	5	52200830	SNORA	26	12	C	T	21.4	234	1.0	91.92	0.0	28.22	0.0	0.0	0.43	78.21	0.00	21.37	0.00	0.95	C>T	21.37	p.R1854R	
2	2	23543959	ZNF91	4	12	T	T	18.6	229	0.0	16.46	0.0	57.11	0.0	0.00	27.07	0.00	72.93	0.00	0.47	T>C	27.07	p.H641R		
3	6	6914102	DEFAS	1	9	G	T	28.0	rs142230452	100	15.14	0.0	37.34	0.0	0.0	0.0	29.00	0.00	71.00	0.00	0.92	G>A	29.00	p.G40Q	
4	12	12575498	ZNF783	4	13	G	T	14.7	267	27.30	0.0	88.122	0.0	0.0	0.0	21.35	0.00	78.65	0.00	0.00	0.76	G>A	21.35	p.T413I	
5	16	16079466	LV9	8	1	C	T	37.1	88	0.0	37.18	0.0	15.14	0.0	0.0	0.00	62.50	0.00	37.50	0.00	0.97	C>T	37.50	p.G570G	
6	3	9260230	LOC10028741	4	6	T	T	82.1	rs75345769	23	0.0	0.0	0.4	0.19	0.0	0.0	0.00	17.39	82.61	0.00	0.00	G>T	82.61	p.T127T	
7	5	5395105	ZNF813	3	19	A	A	19.1	153	57.57	0.0	17.22	0.0	0.0	0.0	74.51	0.00	25.49	0.00	0.00	0.94	A>G	25.49	p.H540R	
8	3	30919047	PMPFYVE	19	2	A	A	23.0	112	37.47	0.0	11.17	0.0	0.0	0.0	75.00	0.00	25.00	0.00	0.75	A>G	25.00	p.I1138V		
9	7	7318832	NLGN2	6	17	C	T	19.0	146	20.13	63.50	0.0	0.0	0.0	0.0	22.60	77.40	0.00	0.00	0.00	0.76	C>A	22.60	p.L380I	
10	15	153997562	DAC1	9	X	C	C	17.0	171	0.0	39.103	0.0	10.19	0.0	0.0	0.00	83.04	0.00	16.96	0.00	0.40	C>T	16.96	p.R298V	
11	14	14622226	LOC10506025	1	T	T	T	54.7	125	40.29	0.0	5.5	41.5	0.0	0.0	55.20	0.00	36.80	0.00	0.45	T>A	55.20	p.R1089R		
12	12	12501347	ZNF739	4	19	C	T	15.3	rs2862768	395	0.0	134.65	0.0	52.14	0.0	0.0	0.00	75.09	0.00	24.91	0.00	0.42	C>T	24.91	p.G523E
13	4	470596	TBRV1-1	2	7	NW_003574	T	100.0	13	0.0	0.0	10.3	0.0	0.0	0.0	0.00	100.00	0.00	0.00	0.30	A>G	100.00	p.G30R		

Tests using 5 pooled controls have shown a reduction of up to 70% in the number of false positive calls. Most of these false positives were caused by common variants and systematic errors that appear in the tumor sample but were not called in the control sample due to low coverage.

Structural Variant Detection

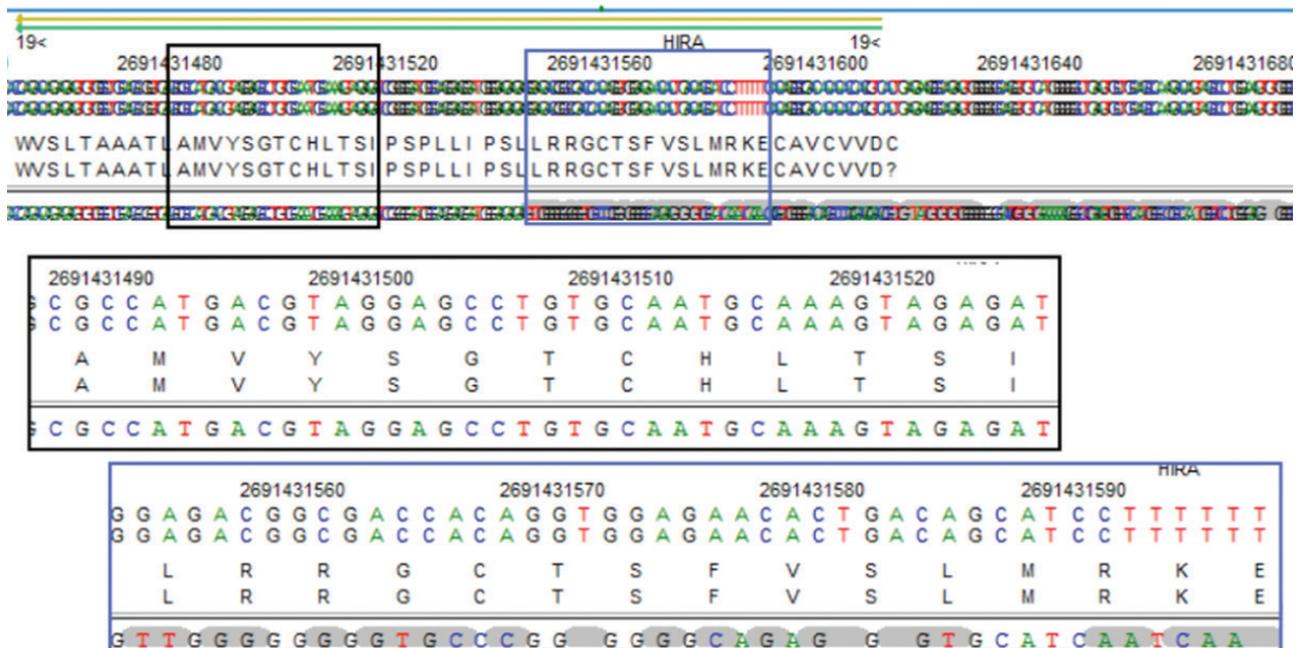
- **Discovery of large Insertion & Deletion; Translocation; Gene Fusion**
- **Flexible Alignment technology allows for large mismatches**
- **Creation of “Pseudo Pairs” allows for detection and mapping of structural variations**

Structural variants (SVs) include insertions, deletions, inversions, and gene fusions that frequently occur across the human genome with over 1000 segmental deletions and over 200 copy number polymorphisms having been reported. These genomic structures have been shown to be important in a number of diseases, usually referred to as genomic disorders.

There are several ways to detect and map SVs, but there are limitations. Microarrays are useful for detecting differences in copy number but are unable to detect smaller SVs and cannot map boundaries. It is possible to detect SVs smaller than 1 kb with sequencing. Paired-end read mapping (PEM) has been used to detect shorter deletions and to hone in on breakage sites but is unable to detect structure variations larger than the library size. NextGENe makes it easy to both find and map structural variants. Targeted sequencing methods such as exon capture with Roche Nimblegen or Agilent SureSelect™ assays can be used to significantly reduce the cost and time requirements of these experiments.

NextGENe allows large mismatches when aligning to the genome in order to find SVs and display information about those regions in a structural variation report. The structural variation report uses a specialized algorithm to list regions with high variant frequency. Interference from false positives caused by sequencing errors is rarely detected in this report since multiple errors are unlikely to occur in a local region.

NextGENe then generates pseudo-paired reads for the sequences aligned to these regions by breaking the original reads into pairs. These can be aligned to the reference genome in order to map the SVs, as seen in figure. Detailed information on where these reads align is available in the Paired Read Reports.



A detailed view shows how NextGENe highlights the mismatched portion of a read. This is an example of fusion gene discovery using NextGENe software. In this example Pseudo-paired reads would be created and split with each half aligning to their appropriate gene. The paired read report identifies the location of each half.



Variant Comparison Tool, with functional prediction

NextGENe includes an advanced mutation comparison tool. When searching for causative mutations of rare diseases, it can be used to narrow down the list of mutations from tens of thousands to a few dozen that can then be examined for possible clinical significance. The tool can also be used to compare unrelated samples.

Excellent for comparing:

- Trios
- Families
- Groups

Advanced Comparison between multiple projects - 5 options

- Manually set expected SNP types
 - Homozygous, Heterozygous, Present, Negative, Undetermined, etc
- Load an inheritance template (Autosomal recessive, X-linked dominant, etc)
- Compound Heterozygous
 - Includes a report listing all valid pairs of mutations
- Shows shared or differences between all individuals
- Gene Association

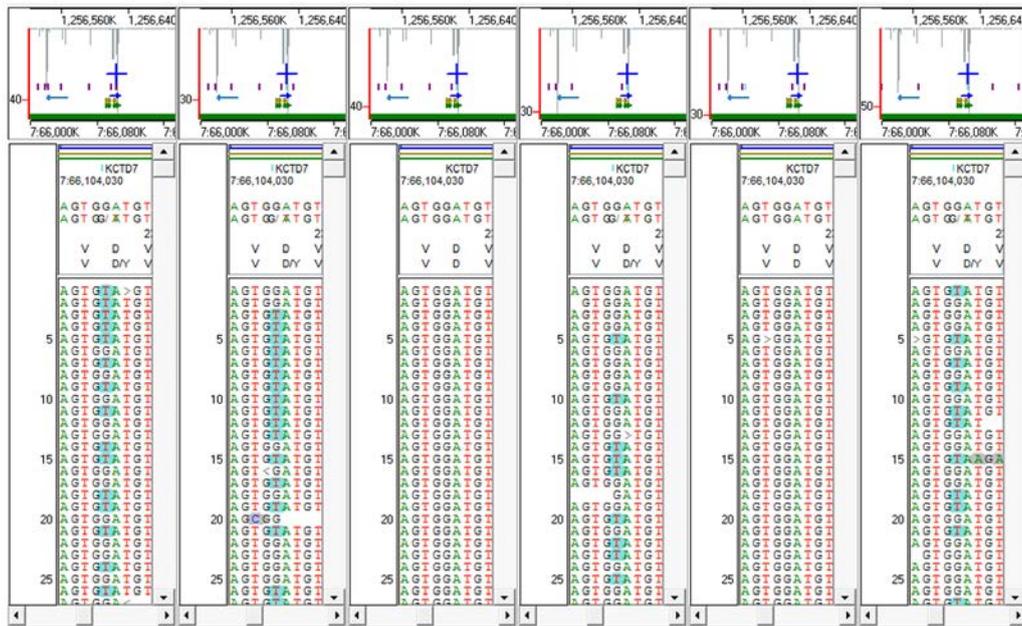
Many advanced filtering options

- Annotation (mRNA, CDS, Splice sites)
- Mutation Confidence Score
- dbSNP mutations
- Substitutions (Non-coding, Silent, Mis-sense, Nonsense) or Indels
- Advanced ROI filtering
- Imported Tracks such as Prediction Information
- Concordance

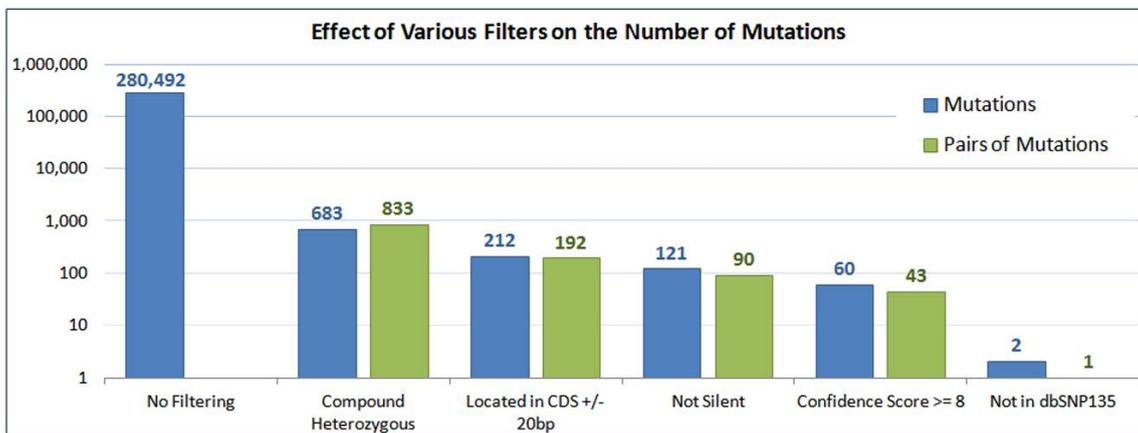
Prediction database integration

- dbNSFP, includes 1000 genomes frequency, phyloP, PolyPhen-2, MutationTaster, SIFT
- COSMIC
- Custom, allows import of proprietary or other public databases

View multiple projects side-by-side



One of two causative mutations found as compared across all six projects. The projects (left to right) are the two affected children, the mother, the father, and the two unaffected children. The second unaffected child has this mutation but does not have the other mutation in the same gene.



Number of Mutations Left After Each Filtering Step Within NextGENe software.

MHC/HLA Analysis

Typing of the Major Histocompatibility Complex (MHC/HLA) in humans and non-human primates using next-generation sequencing allows for phasing of multiple alleles. When using Sanger sequencing, multiple alleles are combined into one signal and minor alleles (due to PCR amplification bias or indel frame shift) are harder to detect. NGS provides a more distinct signal - each read is separate from the others. This allows for detection even in cases of strong PCR bias and insertion and deletion. It is possible to determine HLA genotypes across all exons of an HLA gene using whole gene amplification via long-range PCR. Currently most next-generation sequencers are capable of producing reads that are several hundred bases long, which allows for coverage of exon 2 and 3 from MHC Class I and Class II genes without issue.

NextGENe software processes HLA data by first aligning reads to a reference of selected HLA genes. For each exon the software generates the two most common consensus sequences and sorts aligned reads into three categories based on them- allele one, allele two, and "other". The consensus sequences are then compared to an HLA dictionary to determine to most likely HLA alleles according to amino acid differences, synonymous changes, and intronic changes.

The screenshot displays the NextGENe Viewer interface for HLA analysis. Key components include:

- Whole Genome View:** A genomic map at the top left showing the location of HLA genes.
- Reference Dictionary:** A table showing the reference nucleotide sequence and a consensus sequence for the HLA gene.
- Best Matched Alleles:** A section showing the top matching alleles and their corresponding read alignments.
- Reads matching Allele 1 and Allele 2:** Two panes showing the pile-up of reads grouped by allele.
- HLA Summary Report:** A table listing the top allele picks for each gene, including columns for Index, Sample, Allele 1, Allele 2, Score, Coverage, Poor Coverage, and Mismatches in CDS.
- Allele Matching Report:** A table showing discrepancies between the allele call and the sequence within the reads.
- Allele Coverage Report:** A table showing coverage for various positions across the HLA gene.

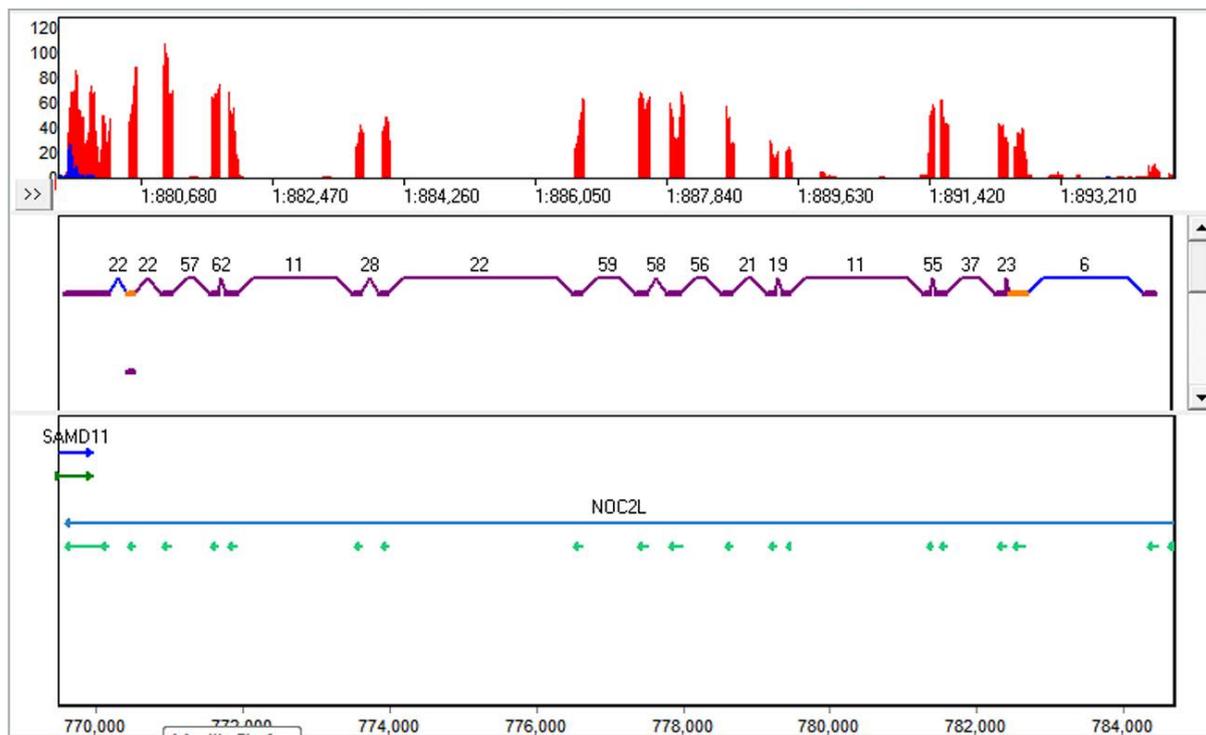
NextGENe Viewer includes several views and reports within a single screen. On the left side, the top pane shows a whole genome view of the HLA genes, the next pane shows the reference nucleotide sequence and a consensus sequence of all alleles, the middle pane shows the top matching alleles and the following two panes show the pile-up of reads grouped by alleles. On the right side, the top pane is the HLA Report listing the top allele picks for each gene and some summary statistics, the next pane down shows the Matching Report identifying the locations where there is a discrepancy between the allele call and the sequence within the reads and the last pane shows the coverage report listing all positions with inadequate coverage.

RNA-Seq Analysis

- Detect multiple transcripts- Insertions, Deletions, Fusions, Un-annotated Transcripts including new genes
- Expression Analysis – Actual coverage (min, max, avg) or normalized (Reads Per Thousand, RPKM, FPKM)
- SNP detection including RNA editing
- Use GBK files or provided whole-genome references.
- Strand-specific analysis

Due to reference sequence difficulties associated with alternative splicing and fusion genes, alignment of RNA-seq data is more challenging than alignment of DNA sequences. Short reads - especially those that fall within large exons - are able to align normally since they will generally match the reference with very few mismatches. Reads that span an exon-exon junction are more difficult because they must be split at the correct position and each part of the read must align correctly. Fusion genes provide even more of a challenge because the partial reads can align almost anywhere in the genome.

Different solutions to these challenges have been implemented in various software packages. Q-PALMA uses a machine learning algorithm and training datasets in order to identify splice junctions. SuperSplat divides sequence reads at multiple positions and tries to find mapping sites where the sub-reads are separated by an intron in a certain size range. TopHat is a software package that first finds potential exons based on coverage and then finds splice sites and links using canonical splice site sequence information. NextGENe uses a novel algorithm to correctly align reads belonging to annotated and novel transcripts while providing the added benefit of a highly graphical interface that doesn't require use of scripting or the command line. Analysis can be performed on a desktop PC in just a few hours without any training datasets or pre-filtering of the reads.



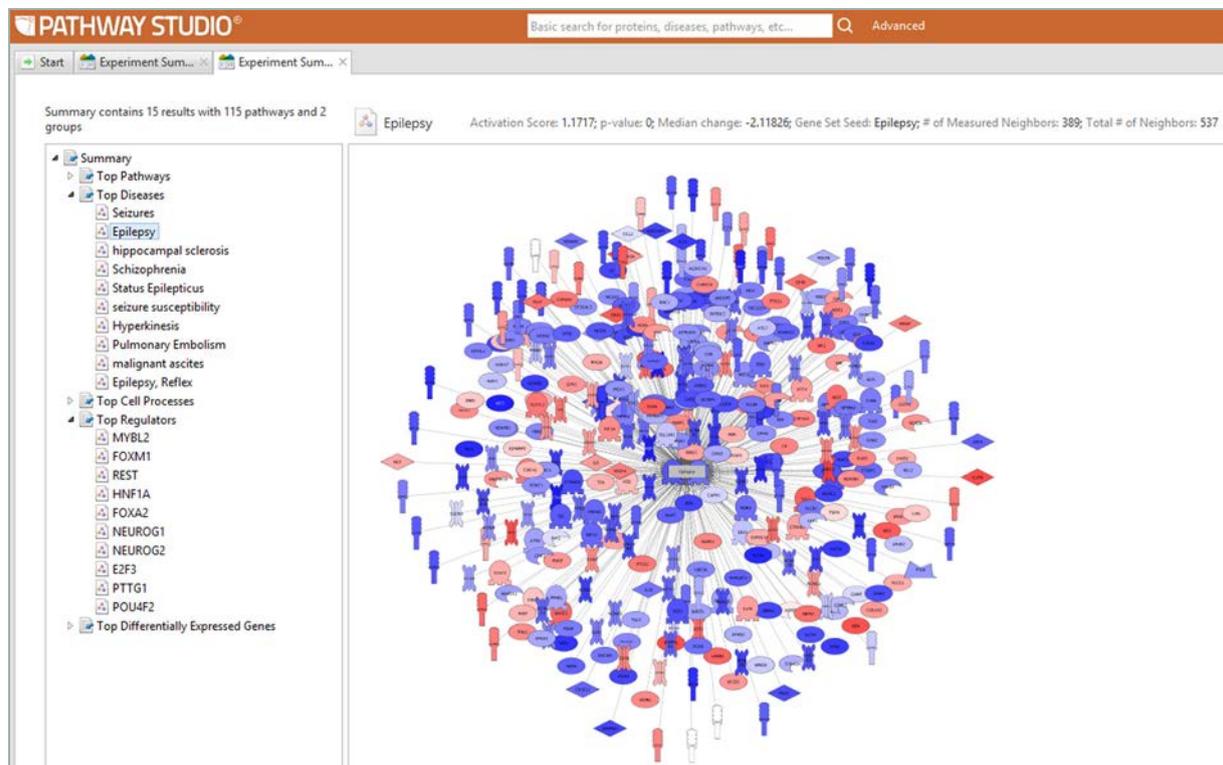
The new RNA-Seq Transcript View. The purple links and exons are known, the blue links are novel, and the orange exons are alternative splicing. The figure of Ion Torrent 318 chip RNA-Seq data set, available on Ion Torrent Development community website.

Integration with Pathway Studio for Differential Gene Expression and Biological Pathway Analysis

Identifying the gene products that are differentially expressed is a key area for understanding essential mechanisms causing diseases such as cancer and autoimmune disorders. Understanding the biological pathways of the genes involved is also critical for the study and characterization of diseases in order to develop effective treatments.

Pathway Studio®, a product of Elsevier R&D Solutions, utilizes a knowledgebase of biological relationships developed from millions of articles and abstracts as well as clinical trials to assist researchers in this important endeavor.

SoftGenetics' NextGENe software now provides a direct link into Elsevier's Pathway Studio for Pathway Studio subscribers, allowing seamless mining of biological systems. NextGENe software's Expression Comparison tool can be used to identify the differentially expressed genes and with a simple click of a button, NextGENe software's expression analysis results can be uploaded to Pathway Studio, enabling the visualization of the disease mechanisms, gene expression and more. This automated linkage can improve laboratory efficiency by allowing researchers to quickly access a broad knowledge base without requiring tedious and error prone manual searches.



de novo Assembly

NextGENe software includes several options to assist in creating fast and accurate assemblies. These include traditional assemblers and several new technologies detailed below.

Floton™ Assembler

- Exclusively for Ion Torrent and Roche technology
- Reduces homopolymer errors to substitutions, greatly improving assembly capability
- Fast, Accurate assemblies

Ion Torrent systems are fundamentally different from most other sequencing systems. The Ion PGM and Proton systems use a “post-light” technology because it doesn’t depend on detection of light emission from nucleotide incorporation. Instead, it uses a silicon chip containing millions of individual pH meters. Its flow-based approach detects pH changes caused by release of hydrogen ion during incorporation of unmodified nucleotides in DNA replication. Because of this different approach to sequencing, the instrument and reagents are much less expensive and it has a unique error profile- most errors are indels rather than substitutions, especially in homopolymer regions.

These errors are more problematic for assembly than substitution errors because of the increased complexity of gapped comparison. However, NextGENe software now includes the newly developed Floton assembler which is able to treat these homopolymer errors as substitution errors. In doing so, it is possible to correct the errors during assembly. This method condenses the sequence into flow calls of individual bases and the number of bases in each flow (see figure). By converting the sequence data into this format, the indels are essentially converted into substitution errors (different base count numbers), allowing for faster computation time and correction of most homopolymer errors. When adjusting the index size, it is important to note that the index size is based on a number of flows, rather than a number of bp.

GGTCCGAAAAACGCCG
↓
GTCGACGCCG
2 1 2 1 6 1 1 2 1

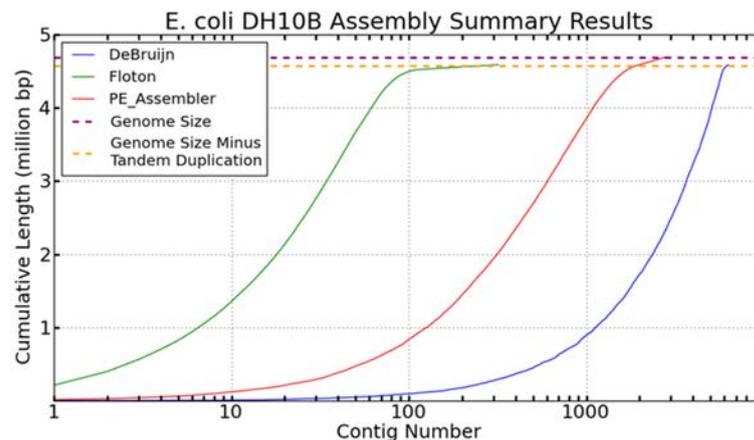
Conversion of base calls into flow calls. In this example a flow size of 9 corresponds to a 17 bp sequence.

Assembly of STO-409 (316 chip) E. coli DH10B

Quality-filtered data:

- Average Coverage = 18x
- Reads = 2,068,174
- Length (Average) = 113.8
- Length (Range) = 25 to 275

	Floton	PE Assembler	DeBruijn
Contigs	320	2,825	6,309
Max Contig Length	212,035	14,571	9,188
Average Contig Length	14,339	1,662	727
N50	64,441	3,750	1,230



Comparison of 3 different assemblers in NextGENe software

The NextGENe software’s Floton Assembler is specially designed to handle data from flow-based sequencing technologies that tend to have indel errors rather than substitution errors. When paired with the Ion Torrent or Roche systems, it provides a very fast and inexpensive method for *de novo* sequencing of bacterial and other small genomes. It is designed to run on relatively inexpensive hardware- these projects were run on a typical laptop - and to be easy to use. Often the default settings do not need to be adjusted at all, and in this case only the index size and minimum contig length were adjusted.

Assembly of Illumina MiSeq™ data

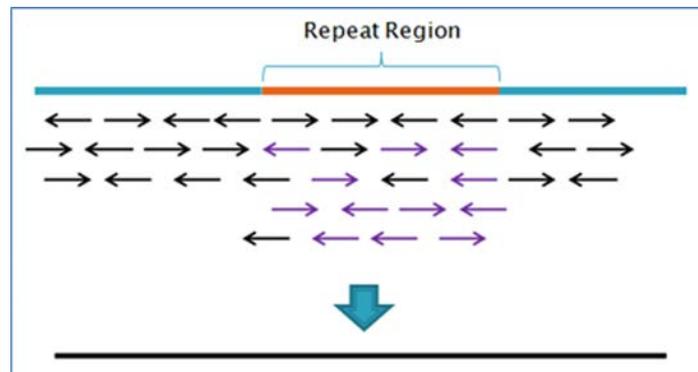
• NextGENe software Stepwise Assembler

Sequences that repeat throughout the genome can pose a problem for the assembly of short reads. Normally the repeat reads would assemble within the incorrect contig. This causes two problems - the repeat regions elsewhere in the genome have reduced coverage and contigs terminate at repeat regions due to the formation of multiple ambiguous contigs.

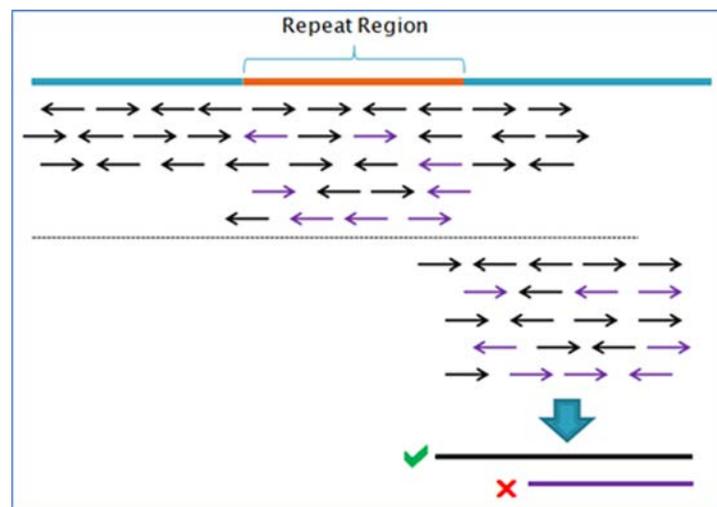
NextGENe software solves this problem with a stepwise paired-end assembly. The end of a contig produced by assembly may indicate a repeat region. The software first calls in the reads paired to those assembled (using overlaps) at the end of the contig (up to 1 ½ times the library size from the end). These reads are then assembled and the software chooses the most complete assembly to continue the contig. Shorter assemblies are not used because they are formed by repetitive sequences erroneously assembled together due to their similarity.

Processing Time	1 hour, 48 min
Total Reads	7,267,665
Reads Used	7,101,508
% Reads Used	97.7%
Number of Contigs	65
Average Length	70,747 bp
Minimum Length	363
Maximum Length	494,530 bp
N50	212,730 bp

E. coli data courtesy Illumina Inc.



Incorrect assembly of reads in a repeat region. Purple reads are from the same repeat sequence elsewhere in the genome.



NextGENe calls in the paired reads and assembles them. Reads paired to those that do not belong in the original assembly will form a separate contig (purple). Only the longest contig is used in the assembly.

Assembly Condensation® Tool

NextGENe now includes a tool to assist in rapid and accurate assemblies of data from all NGS systems. This condensation based tool first scans the short reads removing the excess coverage including repeat regions, retaining longer high-quality reads. This is followed by trimming low quality bases from 3' sequences and finally removes low frequency reads prior to commencing with assembly operation.

Benefits:

Removes excess coverage (including in repeat regions)

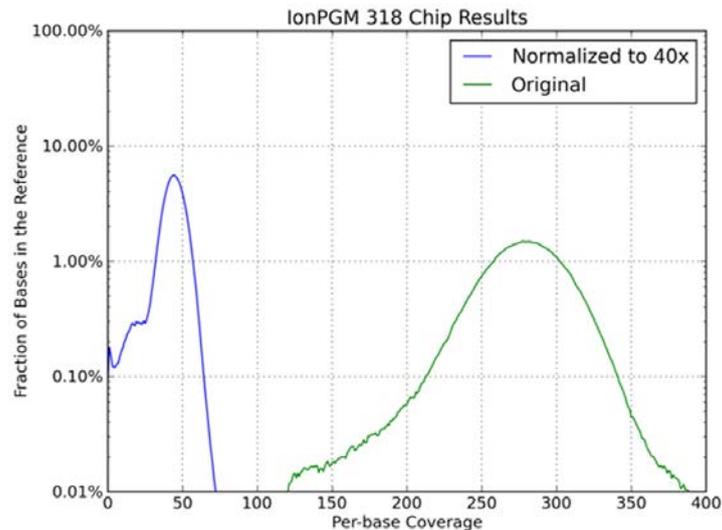
- Longer, higher-quality reads are retained
- Assembly speed and accuracy are improved
- Less RAM is required

Operation:

1. Calculate “flow-mer” frequencies (32 flows each)
2. Remove reads from high-coverage regions. Reads are kept randomly based on the product of three fractions:
 - Normalization Factor (Desired Coverage / Estimated Coverage)
 - Quality adjustment
 - Read Length adjustment
3. Trim low frequency (< 1/10 of desired coverage) 3' flow-mers and remove low frequency reads

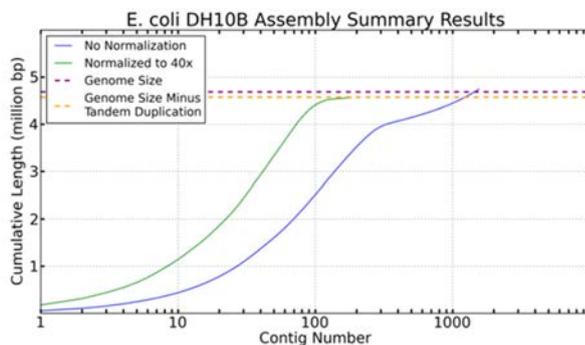
Results:

E coli DH10B sequence run on a 318 chip (C23-140)



	Before	After	Difference
Reads	5,973,581	912,010	-84.7%
Bases	1,323,802,002	199,676,721	-84.9%

Coverage distribution before and after.



	No Condensation	With Condensation	Difference
Contigs	1,555	182	-88.3%
Maximum Length	54,303	177,056	3.3x
Average Length	3,051	25,084	8.2x
N50	16,290	57,186	3.5x
Process Time	69 min, 22 seconds	27 min, 26 seconds	-60%

Assembly Results.

Data courtesy of Ion Community

Please open disc to review your applications of interest and to install a **30-day** trial of NextGENe software.

Disc Contains:

Application Notes

30-day free NextGENe software Trial

User Manual

Minimum Computer Recommendations:

Desktop PC

64 bit, Windows® Vista, 7, 8 or 10 Operating System

Windows® Server 2008 R2 or 2012

Processor: Dual Quad Core Processors

RAM: 12GB

2TB Hard Drive

Intel Powered Macintosh

OS: 10.4.6, with Parallels desktop for MAC or

Apple Boot Camp

64 bit, Windows® Vista, 7, 8 or 10 Operating System

Windows® Server 2008 R2 or 2012

Processor: Dual Quad Core

RAM: 12GB

2TB Hard Drive

Note: Some applications will require additional memory



Bred to Track!

If trial disc is not present please email info@softgenetics.com
for a free 30-day trial

SOFTGENETICS®
Software PowerTools for Genetic Analysis

www.softgenetics.com

SOFTGENETICS®

Software PowerTools for Genetic Analysis

SoftGenetics
Oakwood Centre
100 Oakwood Avenue
Suite 350
State College PA 16803 USA
info@softgenetics.com
www.softgenetics.com

© 2011-2017 Registered Trademarks are property of their respective owners.

For Research Use Only