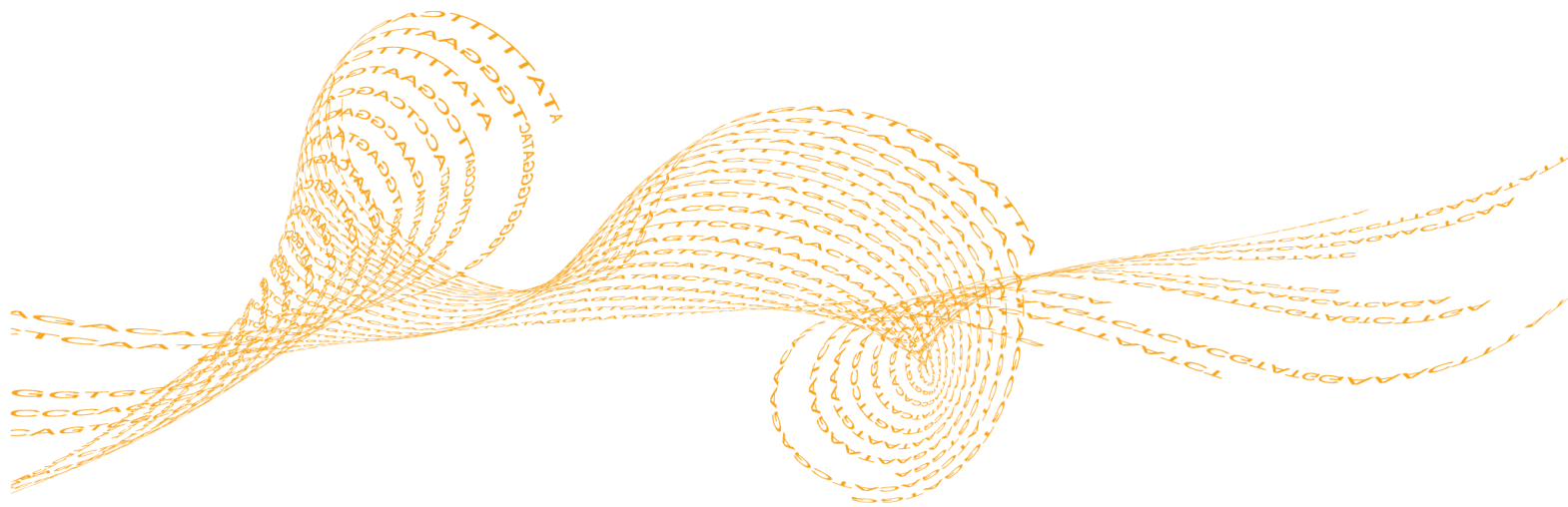


Isaac Enrichment v2.0 App

Introduction	3
Running Isaac Enrichment v2.0	5
Isaac Enrichment v2.0 Output	7
Isaac Enrichment v2.0 Methods	31
Technical Assistance	



This document and its contents are proprietary to Illumina, Inc. and its affiliates ("Illumina"), and are intended solely for the contractual use of its customer in connection with the use of the product(s) described herein and for no other purpose. This document and its contents shall not be used or distributed for any other purpose and/or otherwise communicated, disclosed, or reproduced in any way whatsoever without the prior written consent of Illumina. Illumina does not convey any license under its patent, trademark, copyright, or common-law rights nor similar rights of any third parties by this document.

The instructions in this document must be strictly and explicitly followed by qualified and properly trained personnel in order to ensure the proper and safe use of the product(s) described herein. All of the contents of this document must be fully read and understood prior to using such product(s).

FAILURE TO COMPLETELY READ AND EXPLICITLY FOLLOW ALL OF THE INSTRUCTIONS CONTAINED HEREIN MAY RESULT IN DAMAGE TO THE PRODUCT(S), INJURY TO PERSONS, INCLUDING TO USERS OR OTHERS, AND DAMAGE TO OTHER PROPERTY.

ILLUMINA DOES NOT ASSUME ANY LIABILITY ARISING OUT OF THE IMPROPER USE OF THE PRODUCT(S) DESCRIBED HEREIN (INCLUDING PARTS THEREOF OR SOFTWARE) OR ANY USE OF SUCH PRODUCT(S) OUTSIDE THE SCOPE OF THE EXPRESS WRITTEN LICENSES OR PERMISSIONS GRANTED BY ILLUMINA IN CONNECTION WITH CUSTOMER'S ACQUISITION OF SUCH PRODUCT(S).

FOR RESEARCH USE ONLY

© 2014 Illumina, Inc. All rights reserved.

Illumina, 24sure, BaseSpace, BeadArray, BlueFish, BlueFuse, BlueGnome, cBot, CSPro, CytoChip, DesignStudio, Epicentre, GAllx, Genetic Energy, Genome Analyzer, GenomeStudio, GoldenGate, HiScan, HiSeq, HiSeq X, Infinium, iScan, iSelect, ForenSeq, MiSeq, MiSeqDx, MiSeq FGx, NeoPrep, Nextera, NextBio, NextSeq, Powered by Illumina, SeqMonitor, SureMDA, TruGenome, TruSeq, TruSight, Understand Your Genome, UYG, VeraCode, verifi, VeriSeq, the pumpkin orange color, and the streaming bases design are trademarks of Illumina, Inc. and/or its affiliate(s) in the U.S. and/or other countries. All other names, logos, and other trademarks are the property of their respective owners.

Introduction

The BaseSpace® app Isaac Enrichment v2.0 analyzes DNA that has been enriched for particular target sequences using Nextera® Rapid Capture. Illumina offers fixed panels, add-on, and full custom library prep kits using the Nextera Rapid Capture technology. Alignment is performed with Isaac Aligner, and variant calling with Isaac Variant Caller. Variant analysis is performed for just the target regions. Statistics reporting accumulates coverage and enrichment-specific statistics for each target as well as overall metrics.

The main output files generated by the Isaac Enrichment v2.0 app are:

- ▶ BAM files, containing the reads after alignment.
- ▶ VCF files, containing the variant calls like indels and SNVs, Single Nucleotide Variants.
- ▶ Genome VCF (.genome.vcf) files, describing the calls for all variant and non-variant sites in the genome.



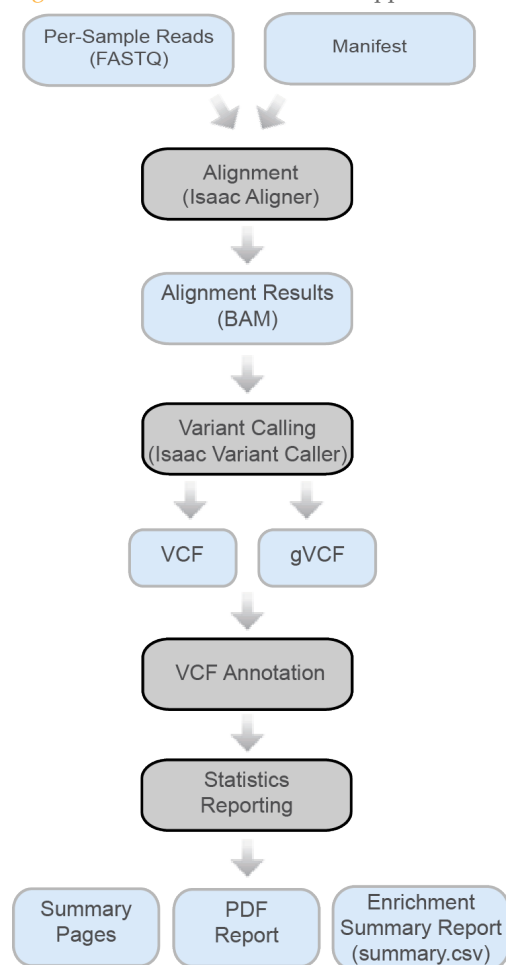
NOTE

VCF and Genome VCF files can be loaded into VariantStudio for viewing; see www.illumina.com/clinical/clinical_informatics/illumina-variantstudio.ilmn.

In addition, there are summary pages and PDF reports.

See *Isaac Enrichment v2.0 Methods* on page 31 and *Isaac Enrichment v2.0 Output* on page 7 for more information.

Figure 1 Isaac Enrichment v2.0 App Workflow



Versions

The following module versions are used in the Isaac Enrichment v2.0 app:

- ▶ Isaac: 02.14.09.26
- ▶ Isaac Variant Caller: 2.0.17
- ▶ Picard: 1.79
- ▶ SAMtools: 0.1.19-isis-1.0.1
- ▶ Tabix: 0.2.5 (r1005)
- ▶ IAS (Annotation Service): v3

Current Limitations

Before running the Isaac Enrichment v2.0 app, be aware of the following limitations:

- ▶ hg19 reference only
- ▶ Read length of at least 32 bp
- ▶ Data set size fewer than 200 gigabytes
- ▶ No minimum number of reads, but use a reasonable input size to get your required coverage
- ▶ Cannot mix single-end and paired-end samples in a single analysis

Running Isaac Enrichment v2.0

- 1 Navigate to the project or sample that you want to analyze.
- 2 Click the **Apps** button and select **Isaac Enrichment v2.0**.
- 3 Fill out the required fields in the Isaac Enrichment input form:
 - a **Analysis Name:** Provide the analysis name. Default name is the app name with the date and time the analysis was started.
 - b **Save Results To:** Select the project that stores the app results.
 - c **Sample(s):** Browse to the sample you want to analyze, and select the checkbox. You can analyze multiple samples.
 - d **Reference Genome:** Select the reference genome. Currently, you can only use hg19.
 - e **Targeted Regions:** Select the targeted region of your enrichment.
 - f **Custom Targeted Manifest:** Select the custom targeted manifest file for analysis. This option is only available when **Custom Manifest** is selected in the **Targeted Regions** drop-down list.
 If you need to upload a custom manifest, perform the following:
 - Navigate to your project in BaseSpace.
 - Click **Import**.
 - Drag and drop your Targeted Region and Probes txt files to the Import window.
 These files will now be added to the project and available under Custom Targeted Manifest.
 - g **Base Padding:** Select the padding you want. Padding defines the amount of sequence immediately upstream and downstream of the targeted regions that is also used in enrichment analysis.
 - h **Annotation:** Choose which gene and transcript annotation reference database to use.
- 4 If desired, fill out the advanced fields in the Isaac Enrichment input form:
 - a **Trim Nextera Rapid Capture Adapters:** If selected, Nextera Rapid Capture adapters are trimmed. Use this setting only if not already applied as a sample sheet setting.
 - b **Flag PCR Duplicates:** If selected, PCR duplicates are flagged in the BAM files and not used for variant calling. PCR duplicates are defined as two clusters from a paired-end run where both clusters have the exact same alignment positions for each read. Optical duplicates are already filtered out during RTA processing. Not applicable for single-end samples.
 - c **Generate Picard HS Metrics and Per Target Coverage Information:** If selected, Picard HS metrics are generated.
 See *Picard Metrics* on page 33 for more information.
 - d **Custom Probes Manifest:** Select the custom probes manifest file for analysis. Required when both **Custom Targeted Manifest** and **Generate Picard HS Metrics and per target coverage information** are selected.

Figure 2 Isaac Enrichment v2.0 Input Form

Analysis Name:	Isaac Enrichment v2.0 12/15/2014 3:19::: ⓘ
Save Results To:	Select Project(s): ⓘ
Sample(s):	Select Sample(s): ⓘ
Reference Genome:	Human (UCSC hg19) ⓘ
Targeted Regions:	Nextera Rapid Capture Exome v1.2 ⓘ
Custom Targeted Manifest:	Select File(s): ⓘ
Base Padding:	150 ⓘ
Annotation:	<input checked="" type="radio"/> RefSeq <input type="radio"/> Ensembl ⓘ
▼ Advanced Options	
Trim Nextera Rapid Capture Adapters:	<input type="checkbox"/> ⓘ
Flag PCR Duplicates:	<input checked="" type="checkbox"/> ⓘ
Generate Picard HS Metrics and Per Target Coverage Information:	<input type="checkbox"/> ⓘ
Custom Probes Manifest:	Select File(s): ⓘ

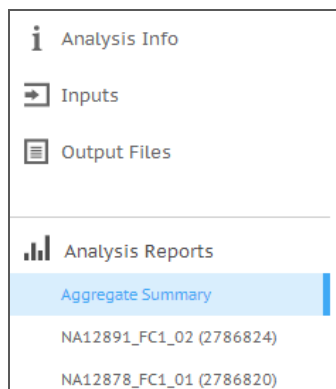
5 Click **Continue**.

The Isaac Enrichment v2.0 app now starts analyzing your sample. When completed, the status of the app session is automatically updated, and you receive an email.

Isaac Enrichment v2.0 Output

This chapter describes the Isaac Enrichment v2.0 app output. To go to the results, click the **Projects** button, then the project, then the analysis.

Figure 3 Isaac Enrichment v2.0 Output Navigation Bar



When the analysis is completed, you can access your output through the left navigation bar, which provides the following:

- ▶ **Analysis Info:** an overview of the app session settings. For more information, see *Analysis Info* on page 20
- ▶ **Inputs:** overview of input settings. For more information, see *Inputs* on page 21
- ▶ **Output Files:** access to the output files, organized by sample and app session. For more information, see *Isaac Enrichment v2.0 Output Files* on page 21.
- ▶ **Aggregate Summary:** access to analysis metrics for the aggregate results. The Aggregate Summary is only displayed if multiple samples are analyzed. For more information, see *Aggregate Summary Page* on page 7.
- ▶ **Sample Pages:** access to analysis reports for each sample. For more information, see *Sample Summary Page* on page 13.

Aggregate Summary Page

The Isaac Enrichment v2.0 app provides an overview of metrics for all samples combined on the Aggregate Summary page. You can view the numerical data and histograms with metrics graphed by sample. You can also download the Enrichment Sequencing Report as PDF (see also *Enrichment Sequencing Report* on page 10).



NOTE

PCR duplicate reads are not removed from statistics. Results are not directly comparable to Picard HsMetrics.

Enrichment Summary

Aligned bases are plotted against sample. The metrics displayed are explained here:

- ▶ Read Level

Statistic	Definition
Total Aligned Reads	The total number of reads passing filter present in the data set that aligned to the reference genome.

Statistic	Definition
Percent Aligned Reads	The percentage of reads passing filter that aligned to the reference genome.
Target Aligned Reads	Number of reads that aligned to the target.
Read Enrichment	$100 * (\text{Target aligned reads} / \text{Total aligned reads})$.
Padded Target Aligned Reads	Number of reads that aligned to the padded target.
Padded Read Enrichment	$100 * (\text{Padded target aligned reads} / \text{Total aligned reads})$.

► Base Level

Statistic	Definition
Total Aligned Bases	The total number of bases present in the data set that aligned to the reference genome.
Target Aligned Bases	Total aligned bases in the target region.
Base Enrichment	$100 * (\text{Total Aligned Bases in Targeted Regions} / \text{Total Aligned Bases})$.
Padded Target Aligned Bases	Total aligned bases in the padded target region.
Padded Base Enrichment	$100 * (\text{Total Aligned Bases in Padded Targeted Regions} / \text{Total Aligned Bases})$.

Variant Summary

► SNVs

Number of SNVs passing are plotted against sample. The metrics displayed are explained here:

Statistic	Definition
Total Passing	Total number of Single Nucleotide Variants present in the data set passing the quality filters.
Het/Hom Ratio	Number of heterozygous SNVs/Number of homozygous SNVs.
Ts/Tv Ratio	The number of Transition SNVs that pass the quality filters divided by the number of Transversion SNVs that pass the quality filters. Transitions are interchanges of purines (A, G) or of pyrimidines (C, T). Transversions are interchanges of purine and pyrimidine bases (for example, A to T).
Percent Found in dbSNP	$100 * (\text{Number of SNVs in dbSNP} / \text{Number of SNVs})$. The SNVs that were found in the dbSNP are annotated accordingly.

► Indels

Number of indels passing are plotted against sample. The metrics displayed are explained here:

Statistic	Definition
Total Passing	Total number of indels present in the data set passing the quality filters.
Het/Hom Ratio	Number of heterozygous indels/Number of homozygous indels.
Percent Found in dbSNP	100*(Number of Indels in dbSNP/Number of Indels).

Coverage Summary

► Mean Region Coverage Depth

The mean region coverage depth is plotted against sample. The metrics displayed are explained in the next table.

► Uniformity of Coverage

The uniformity of coverage is plotted against sample. The metrics displayed are explained in the next table.

► Depth of Coverage in Targeted Regions

The depth of coverage in targeted regions is plotted against sample. Coverage values for the targeted bases covered can be found in the downloadable Export (CSV) file. The metrics displayed are explained in the next table.

Statistic	Definition
Mean Region Coverage Depth	The total number of aligned bases to the targeted region divided by the targeted region size.
Uniformity of Coverage (Pct > 0.2*mean):	The percentage of targeted base positions in which the read depth is greater than 0.2 times the mean region target coverage depth.
Target Coverage at 1X	Percentage targets with coverage greater than 1X.
Target Coverage at 10X	Percentage targets with coverage greater than 10X.
Target Coverage at 20X	Percentage targets with coverage greater than 20X.
Target Coverage at 50X	Percentage targets with coverage greater than 50X.

Fragment Length Summary

► Fragment Length Medians

Fragment length medians are plotted against sample. The metrics displayed are explained here:

Statistic	Definition
Fragment Length Median	Median length of the sequenced fragment. The fragment length is calculated based on the locations at which a read pair aligns to the reference. The read mapping information is parsed from the BAM files.

Statistic	Definition
Minimum	Minimum length of the sequenced fragment.
Maximum	Maximum length of the sequenced fragment.
Standard Deviation	Standard deviation of the sequenced fragment length.

Duplicates Summary

► Percent Duplicates Paired Reads

Percent duplicate paired reads are plotted against sample. The metric displayed is explained here:

Statistic	Definition
Percent Duplicate Paired Reads	Percentage of paired reads that have duplicates.

Enrichment Sequencing Report

The Isaac Enrichment v2.0 app provides an aggregate summary PDF report for all samples combined on the Summary page.

This section contains a description of the report.



NOTE

PCR duplicate reads are not removed from statistics. Results are not directly comparable to Picard HsMetrics.

Manifest Information

Provides the target manifest name, and the following metrics:

Statistic	Definition
Total length	The total length of the sequenced bases in the target region.
Padding Size	The length of sequence immediately upstream and downstream of the enrichment targets that is included for a padded target.

Sample Information

Defines the sample numbers and names in the report.

Enrichment Summary

► Read Level Enrichment

Statistic	Definition
Total Aligned Reads	The total number of reads passing filter present in the data set that aligned to the reference genome.

Statistic	Definition
Targeted Aligned Reads	Number of reads that aligned to the target.
Read Enrichment	$100 * (\text{Target aligned reads} / \text{Total aligned reads})$.
Padded Target Aligned Reads	Number of reads that aligned to the padded target.
Padded Read Enrichment	$100 * (\text{Padded target aligned reads} / \text{Total aligned reads})$.

► Base Level Enrichment

Statistic	Definition
Total Aligned Bases	The total number of bases present in the data set that aligned to the reference genome.
Targeted Aligned Bases	Total aligned bases in the target region.
Base Enrichment	$100 * (\text{Total Aligned Bases in Targeted Regions} / \text{Total Aligned Bases})$.
Padded Target Aligned Bases	Total aligned bases in the padded target region.
Padded Base Enrichment	$100 * (\text{Total Aligned Bases in Padded Targeted Regions} / \text{Total Aligned Bases})$.

The Base Enrichment histogram graphs the total aligned bases, padded targeted aligned bases, and padded base enrichment by sample.

SNV Summary

Statistic	Definition
SNVs	Total number of Single Nucleotide Variants present in the data set passing the quality filters.
SNVs (Percent Found in dbSNP)	$100 * (\text{Number of SNVs in dbSNP} / \text{Number of SNVs})$. The SNVs that were found in the dbSNP are annotated accordingly.
SNV Ts/Tv Ratio	The number of Transition SNVs that pass the quality filters divided by the number of Transversion SNVs that pass the quality filters. Transitions are interchanges of purines (A, G) or of pyrimidines (C, T). Transversions are interchanges of purine and pyrimidine bases (for example, A to T).
SNV Het/Hom Ratio	Number of heterozygous SNVs/Number of homozygous SNVs.

The SNVs histogram graphs the number of SNVs passing by sample.

Indel Summary

Statistic	Definition
Indels	Total number of indels present in the data set passing the quality filters.
Indels (Percent Found in dbSNP)	$100 * (\text{Number of Indels in dbSNP} / \text{Number of Indels})$.
Indel Het/Hom Ratio	Number of heterozygous indels/Number of homozygous indels.

The Indels histogram graphs the number of indels passing by sample.

Coverage Summary

Statistic	Definition
Mean Region Coverage Depth	The total number of aligned bases to the targeted region divided by the targeted region size.
Uniformity of Coverage (Pct > 0.2*mean):	The percentage of targeted base positions in which the read depth is greater than 0.2 times the mean region target coverage depth.
Target Coverage at 1X	Percentage targets with coverage greater than 1X.
Target Coverage at 10X	Percentage targets with coverage greater than 10X.
Target Coverage at 20X	Percentage targets with coverage greater than 20X.
Target Coverage at 50X	Percentage targets with coverage greater than 50X.

In addition, the app provides two graphs:

- ▶ A Coverage Mean and Uniformity histogram that plots the mean region coverage depth and uniformity of coverage by sample.
- ▶ A Depth of Coverage in Targeted Regions histogram that plots the number of targeted sequences by the depth of coverage.

Statistic	Definition
Depth of Sequencing Coverage	The coverage depth of a position in the genome refers to the number of sequenced bases that align to that position.
Number of Targeted Bases Covered	Number of targeted bases that have at least the indicated depth of coverage.

Fragment Length Summary

Statistic	Definition
Fragment Length Median	Median length of the sequenced fragment. The fragment length is calculated based on the locations at which a read pair aligns to the reference. The read mapping information is parsed from the BAM files.
Minimum	Minimum length of the sequenced fragment.
Maximum	Maximum length of the sequenced fragment.
Standard Deviation	Standard deviation of the sequenced fragment length.

The Fragment Length Medians histogram graphs the fragment length median by sample.

Duplicates Summary

Statistic	Definition
Percent Duplicate Paired Reads	Percentage of paired reads that have duplicates.

The Percent Duplicate Paired Reads histogram graphs the percent duplicate paired reads by sample.

Sample Summary Page

The Isaac Enrichment v2.0 app provides an overview of statistics per sample on the sample pages. You can also download each sample's Enrichment Sequencing Report as PDF.



NOTE

PCR duplicate reads are not removed from statistics. Results are not directly comparable to Picard HsMetrics.

Sample Information

Statistic	Definition
Total PF Reads	The number of reads passing filter for the sample.
Percent Q30	The percentage of bases with a quality score of 30 or higher.

Enrichment Summary

Provides the target manifest name, and the following metrics:

Statistic	Definition
Total Length of Targeted Reference	The total length of the sequenced bases in the target region.

Statistic	Definition
Padding Size	The length of sequence immediately upstream and downstream of the enrichment targets that is included for a padded target.

► Read Level Enrichment

Statistic	Definition
Total Aligned Reads	The total number of reads passing filter present in the data set that aligned to the reference genome.
Percent Aligned Reads	The percentage of reads passing filter that aligned to the reference genome.
Target Aligned Reads	Number of reads that aligned to the target.
Read Enrichment	$100 * (\text{Target aligned reads} / \text{Total aligned reads})$.
Padded Target Aligned Reads	Number of reads that aligned to the padded target.
Padded Read Enrichment	$100 * (\text{Padded target aligned reads} / \text{Total aligned reads})$.

► Base Level Enrichment

Statistic	Definition
Total Aligned Bases	The total number of bases present in the data set that aligned to the reference genome.
Target Aligned Bases	Total aligned bases in the target region.
Bases Enrichment	$100 * (\text{Total Aligned Bases in Targeted Regions} / \text{Total Aligned Bases})$.
Padded Target Aligned Bases	Total aligned bases in the padded target region.
Padded Base Enrichment	$100 * (\text{Total Aligned Bases in Padded Targeted Regions} / \text{Total Aligned Bases})$.

Small Variants Summary

This table provides metrics about the number of SNVs, insertions, and deletions.

Statistic	Definition
Total Passing	The total number of variants present in the data set that passed the variant quality filters.
Percent Found in dbSNP	$100 * (\text{Number of variants in dbSNP} / \text{Number of variants})$.

Statistic	Definition
Het/Hom Ratio	Number of heterozygous variants/Number of homozygous variants.
Ts/Tv Ratio	Transition rate of SNVs that pass the quality filters divided by transversion rate of SNVs that pass the quality filters. Transitions are interchanges of purines (A, G) or of pyrimidines (C, T). Transversions are interchanges between purine and pyrimidine bases (for example, A to T).

Variants by Sequence Context

Statistic	Definition
Number in Genes	The number of variants that fall into a gene.
Number in Exons	The number of variants that fall into an exon.
Number in Coding Regions	The number of variants that fall into a coding region.
Number in UTR Regions	The number of variants that fall into an untranslated region (UTR).
Number in Splice Site Regions	The number of variants that fall into a splice site region.
Number in Mature microRNA	The number of variants that fall into a mature microRNA.

Variants by Consequence

Statistic	Definition
Frameshifts	The number of variants that cause a frameshift.
Non-synonymous	The number of variants that cause an amino acid change in a coding region.
Synonymous	The number of variants that are within a coding region, but do not cause an amino acid change.
Stop Gained	The number of variants that cause an additional stop codon.
Stop Lost	The number of variants that cause the loss of a stop codon.

Coverage Summary

Statistic	Definition
Mean Coverage	The total number of aligned bases to the targeted region divided by the targeted region size.
Uniformity of Coverage (Pct > 0.2*mean):	The percentage of targeted base positions in which the read depth is greater than 0.2 times the mean region target coverage depth.

Statistic	Definition
Target Coverage at 1X	Percentage targets with coverage greater than 1X.
Target Coverage at 10X	Percentage targets with coverage greater than 10X.
Target Coverage at 20X	Percentage targets with coverage greater than 20X.
Target Coverage at 50X	Percentage targets with coverage greater than 50X.

Fragment Length Summary

Statistic	Definition
Fragment Length Median	Median length of the sequenced fragment. The fragment length is calculated based on the locations at which a read pair aligns to the reference. The read mapping information is parsed from the BAM files.
Minimum	Minimum length of the sequenced fragment.
Maximum	Maximum length of the sequenced fragment.
Standard Deviation	Standard deviation of the sequenced fragment length.

Duplicate Information

Statistic	Definition
Percent Duplicate Paired Reads	Percentage of paired reads that have duplicates.

Enrichment Sequencing Reports by Sample

The Isaac Enrichment v2.0 app provides an enrichment statistics PDF report for each sample.



NOTE

PCR duplicate reads are not removed from statistics. Results are not directly comparable to Picard HsMetrics.

Sample Information

Provides the sample ID and name, and the following metrics:

Statistic	Definition
Total PF Reads	The number of reads passing filter for the sample.
Percent Q30	The percentage of bases with a quality score of 30 or higher.

Statistic	Definition
Median Read Length	Median length of the sequenced fragment. The fragment length is calculated based on the locations at which a read pair aligns to the reference. The read mapping information is parsed from the BAM files.
Adapters Trimmed	Whether adapter trimming was used.

Enrichment Summary

Provides the target manifest name, and the following metrics:

Statistic	Definition
Total Length of Targeted Reference	The total length of the sequenced bases in the target region.
Padding Size	The length of sequence immediately upstream and downstream of the enrichment targets that is included for a padded target.

Read Level Enrichment

Statistic	Definition
Total Aligned Reads	The total number of reads passing filter present in the data set that aligned to the reference genome.
Percent Aligned Reads	The percentage of reads passing filter that aligned.
Targeted Aligned Reads	Number of reads that aligned to the target.
Read Enrichment	$100 * (\text{Target aligned reads} / \text{Total aligned reads})$.
Padded Target Aligned Reads	Number of reads that aligned to the padded target.
Padded Read Enrichment	$100 * (\text{Padded target aligned reads} / \text{Total aligned reads})$.

Base Level Enrichment

Statistic	Definition
Total Aligned Bases	The total number of bases present in the data set that aligned to the reference genome.
Percent Aligned Bases	The percentage of bases that aligned to the reference genome.
Targeted Aligned Bases	Total aligned bases in the target region.

Statistic	Definition
Base Enrichment (not padded)	$100 * (\text{Total Aligned Bases in Targeted Regions} / \text{Total Aligned Bases})$.
Padded Target Aligned Bases	Total aligned bases in the padded target region.
Padded Base Enrichment	$100 * (\text{Total Aligned Bases in Padded Targeted Regions} / \text{Total Aligned Bases})$.

Small Variants Summary

This table provides metrics about the number of SNVs, insertions, and deletions.

Statistic	Definition
Total Passing	The total number of variants present in the data set that passed the variant quality filters.
Percent Found in dbSNP	$100 * (\text{Number of variants in dbSNP} / \text{Number of variants})$.
Het/Hom Ratio	Number of heterozygous variants/Number of homozygous variants.
Ts/Tv Ratio	Transition rate of SNVs that pass the quality filters divided by transversion rate of SNVs that pass the quality filters. Transitions are interchanges of purines (A, G) or of pyrimidines (C, T). Transversions are interchanges between purine and pyrimidine bases (for example, A to T).

Variants by Sequence Context

Statistic	Definition
In Genes	The number of variants that fall into a gene.
In Exons	The number of variants that fall into an exon.
In Coding Regions	The number of variants that fall into a coding region.
In UTR Regions	The number of variants that fall into an untranslated region (UTR).
In Splice Site Regions	The number of variants that fall into a splice site region.
In Mature microRNA	The number of variants that fall into a mature microRNA.

Variants by Consequence

Statistic	Definition
Frameshifts	The number of variants that cause a frameshift.
Non-synonymous	The number of variants that cause an amino acid change in a coding region.

Statistic	Definition
Synonymous	The number of variants that are within a coding region, but do not cause an amino acid change.
Stop Gained	The number of variants that cause an additional stop codon.
Stop Lost	The number of variants that cause the loss of a stop codon.

Coverage Summary

Statistic	Definition
Mean region Coverage depth	The total number of aligned bases to the targeted region divided by the targeted region size.
Uniformity of Coverage (Pct > 0.2*mean):	The percentage of targeted base positions in which the read depth is greater than 0.2 times the mean region target coverage depth.
Target Coverage at 1X	Percentage targets with coverage greater than 1X.
Target Coverage at 10X	Percentage targets with coverage greater than 10X.
Target Coverage at 20X	Percentage targets with coverage greater than 20X.
Target Coverage at 50X	Percentage targets with coverage greater than 50X.

In addition, the app provides two graphs:

- ▶ A Mean Coverage by Targeted Region graph that plots the mean coverage by the targeted region.
- ▶ A Depth of Coverage in Targeted Regions graph that plots the number of targeted sequences by the depth of coverage.

Statistic	Definition
Depth of Sequencing Coverage	The coverage depth of a position in the genome refers to the number of sequenced bases that align to that position.
Number of Targeted Bases Covered at Depth	Number of targeted bases that have at least the indicated depth of coverage.
Total Targeted Bases Covered	Total bases aligning to the target regions that have at least the indicated depth of coverage.
Target Coverage	Percent of targeted bases that reach the indicated depth of coverage.

Fragment Length Summary

Statistic	Definition
Fragment Length Median	Median length of the sequenced fragment. The fragment length is calculated based on the locations at which a read pair aligns to the reference. The read mapping information is parsed from the BAM files.
Minimum	Minimum length of the sequenced fragment.
Maximum	Maximum length of the sequenced fragment.
Standard Deviation	Standard deviation of the sequenced fragment length.

Gaps Summary

The app also provides a Targeted Regions Gap Length Distribution graph that plots the number of gaps on a log scale by the length of the gap in bases.

Duplicates Information

Statistic	Definition
Percent Duplicate Paired Reads	Percentage of paired reads that have duplicates.

Analysis Details

Provides the analysis settings, software versions, and data collections used.

Analysis Info

This app provides an overview of the analysis on the Analysis Info page.

A brief description of the metrics is below.

Table 1 Analysis Info

Row	Definition
Name	Name of the app session.
Application	App that generated this analysis.
Date Started	Date and time the app session started.
Date Completed	Date and time the app session completed.
Duration	Duration of analysis.
Session Type	The number of nodes used.
Size	Total size of all output files.
Status	Status of the app session.

Log Files

Clicking the **Log Files** link on the Analysis Info page provides access to the app log files.

The key log files to help follow data processing and debugging are the following:

- ▶ **CompletedJobInfo.xml**: Contains information about the completed job.
- ▶ **EnrichmentStatistics.xml**: Contains statistics about the completed job.
- ▶ **Logging.zip**: Contains all detailed workflow log files for each step of the workflow.
- ▶ **output-timestamp.log**: Shows the raw console output from the Enrichment app.
- ▶ **SampleSheet.csv**: Sample sheet.
- ▶ **SampleSheetUsed.csv**: A copy of the sample sheet, generated at the end of a run.
- ▶ **spacedock-{timestamp}.log**: Shows console output from the SpaceDock and BaseSpace communication and input/output file staging.
- ▶ **spacedock-infrastructure-{timestamp}.log**: Log file used for debugging.
- ▶ **WorkflowError.txt**: Workflow standard error output (contains error messages created while running the workflow).
- ▶ **WorkflowLog.txt**: Workflow standard output (contains details about workflow steps, command line calls with parameters, timing, and progress).

Isaac Enrichment v2.0 Status

For single samples, the status of the Isaac Enrichment v2.0 app session can have the following values:

- 1 Launching Isis
- 2 Alignment
- 3 Variant analysis
- 4 Calculate Picard HS Metrics (if selected in *Running Isaac Enrichment v2.0* on page 5)
- 5 Statistics evaluation
- 6 Report generation
- 7 Aggregate generation (for multiple samples)

Inputs

The Inputs page provides an overview of the input samples and settings that were specified when the Isaac Enrichment v2.0 project was set up.

Isaac Enrichment v2.0 Output Files

The Output Files page provides access to the output files. See the following pages for descriptions:

- ▶ *Enrichment Sequencing Report* on page 10
- ▶ *Enrichment Sequencing Reports by Sample* on page 16
- ▶ *BAM Files* on page 22
- ▶ *VCF Files* on page 22
- ▶ *gVCF Files* on page 23
- ▶ *Enrichment Summary Report (*.summary.csv)* on page 27
- ▶ *Manifest Output Files* on page 30

BAM Files

The Sequence Alignment/Map (SAM) format is a generic alignment format for storing read alignments against reference sequences, supporting short and long reads (up to 128 Mb) produced by different sequencing platforms. SAM is a text format file that is human-readable. The Binary Alignment/Map (BAM) keeps the same information as SAM, but in a compressed, binary format that is only machine readable.

If you use an app in BaseSpace that uses BAM files as input, the app locates the file when launched. If using BAM files in other tools, download the file to use it in the external tool.

Go to samtools.sourceforge.net/SAM1.pdf to see the exact SAM specification.

VCF Files

VCF is a text file format that contains information about variants found at specific positions in a reference genome. The file format consists of meta-information lines, a header line, and then data lines. Each data line contains information about a single variant.

If you use an app in BaseSpace that uses VCF files as input, the app locates the file when launched. If using VCF files in other tools, download the file to use it in the external tool.

A detailed description of the VCF format is provided in the *BaseSpace User Guide*.

Additional entries are described in the section *Isaac Enrichment v2.0 VCF Entries* on page 22.

Isaac Enrichment v2.0 VCF Entries

The VCF files for Isaac Enrichment v2.0 can have the following entries in the FILTER, FORMAT, and INFO fields:

Table 2 VCF FILTER Entries

Entry	Description
IndelConflict	Locus is in region with conflicting indel calls
SiteConflict	Site genotype conflicts with proximal indel call, typically a heterozygous SNV call made inside of a heterozygous deletion
LowGQX	Locus GQX is less than 30 or not present
HighDPFRatio	The fraction of base calls filtered out at a site is greater than 0.4
HighSNVSB	SNV strand bias value (SNVSB) exceeds 10
HighDepth	Locus depth is greater than 3x the mean chromosome depth
OffTarget	Variant is not on target

Table 3 VCF FORMAT Entries

Entry	Description
GQX	Minimum of {Genotype quality assuming variant position, Genotype quality assuming non-variant position}
GT	Genotype

Entry	Description
GQ	Genotype Quality
DP	Filtered base call depth used for site genotyping
DPF	Base calls filtered from input before site genotyping
AD	Allelic depths for the ref and alt alleles in the order listed. For indels, this value only includes reads that confidently support each allele (posterior probability 0.999 or higher that read contains indicated allele vs all other intersecting indel alleles)
DPI	Read depth associated with indel, taken from the position preceding the indel.

Table 4 VCF INFO Entries

Entry	Description
SNVSB	SNV site strand bias
SNVHPOL	SNV contextual homopolymer length
CIGAR	CIGAR alignment for each alternate indel allele
RU	Smallest repeating sequence unit extended or contracted in the indel allele relative to the reference. RUs longer than 20 bases are not reported.
REFREP	Number of times RU is repeated in reference.
IDREP	Number of times RU is repeated in indel allele.
END	End position of the region described in this record
BLOCKAVG_min30p3a	Non-variant site block. All sites in a block are constrained to be non-variant, have the same filter value, and have all sample values in range [x,y], $y \leq \max(x+3, (x*1.3))$. All printed site block sample values are the minimum observed in the region spanned by the block

gVCF Files

This application also produces the Genome Variant Call Format file (gVCF). gVCF was developed to store sequencing information for both variant and non-variant positions, which is required for human clinical applications. gVCF is a set of conventions applied to the standard variant call format (VCF) 4.1 as documented by the 1000 Genomes Project. These conventions allow representation of genotype, annotation, and other information across all sites in the genome in a compact format. Typical human whole-genome sequencing results expressed in gVCF with annotation are less than 1 Gbyte, or about 1/100 the size of the BAM file used for variant calling. If you are performing targeted sequencing, gVCF is also an appropriate choice to represent and compress the results.

gVCF is a text file format, stored as a gzip compressed file (*.genome.vcf.gz). Compression is further achieved by joining contiguous non-variant regions with similar properties into single 'block' VCF records. To maximize the utility of gVCF, especially for high stringency applications, the properties of the compressed blocks are conservative. Block properties like depth and genotype quality reflect the minimum of any site in the block. The gVCF file can be indexed (creating a *.tbi file) and used with existing VCF

tools such as tabix and IGV, making it convenient both for direct interpretation and as a starting point for further analysis.

For more information, see sites.google.com/site/gvcftools/home/about-gvcf.

The following conventions are used in the variant caller gVCF files.

Samples per File

There is only one sample per gVCF file.

Non-Variant Blocks Using END Key

Contiguous non-variant segments of the genome can be represented as single records in gVCF. These records use the standard 'END' INFO key to indicate the extent of the record. Even though the record can span multiple bases, only the first base is provided in the REF field to reduce file size.

The following is a simplified segment of a gVCF file, describing a segment of non-variant calls (starting with an A) on chromosome 1 from position 51845 to 51862.

```
##INFO=<ID=END,Number=1,Type=Integer,Description="End position  
of the variant described in this record">#CHROM POS ID REF  
ALT QUAL FILTER INFO FORMAT NA19238chr1 51845 . A . . PASS  
END=51862
```

Any field provided for a block of sites, such as read depth (using the DP key), shows the minimum value that is observed among all sites encompassed by the block. Each sample value shown for the block, such as the depth (DP), is restricted to a range where the maximum value is within 30% or 3 of the minimum. For example, for sample value range [x,y], $y \leq x + \max(3, x * 0.3)$. This range restriction applies to each of the sample values printed in the final block record.

Indel Regions

Sites that are "filled in" inside of deletions have additional changes:

All deletions:

- ▶ Sites inside of any deletion are marked with the deletion filters, in addition to any filters that have already been applied to the site.
- ▶ Sites inside of deletions cannot have a genotype or alternate allele quality score higher than the corresponding value from the enclosing indel.

Heterozygous deletions:

- ▶ Sites inside of heterozygous deletions are altered to have haploid genotype entries (e.g. "0" instead of "0/0", "1" instead of "1/1").
- ▶ Heterozygous SNV calls inside of heterozygous deletions are marked with the "SiteConflict" filter and their genotype is unchanged.

Homozygous deletions:

- ▶ Homozygous reference and no-call sites inside of homozygous deletions have genotype "."
- ▶ Sites inside of homozygous deletions that have a non-reference genotype are marked with a "SiteConflict" filter, and their genotype is unchanged.
- ▶ Site and genotype quality are set to "."

The described modifications reflect the notion that the site confidence is bound within the enclosing indel confidence.

On occasion, the variant caller produces multiple overlapping indel calls that cannot be resolved into two haplotypes. If this case, all indels and sites in the region of the overlap are marked with the *IndelConflict* filter.

Genotype Quality for Variant and Non-variant Sites

The gVCF file uses an adapted version of genotype quality for variant and non-variant site filtration. This value is associated with the key GQX. The GQX value is intended to represent the minimum of {Phred genotype quality assuming the site is variant, Phred genotype quality assuming the site is non-variant}. The reason for using this value is to allow a single value to be used as the primary quality filter for both variant and non-variant sites. Filtering on this value corresponds to a conservative assumption appropriate for applications where reference genotype calls must be determined at the same stringency as variant genotypes, ie:

- ▶ An assertion that a site is homozygous reference at $GQX \geq 30$ is made assuming the site is variant.
- ▶ An assertion that a site is a non-reference genotype at $GQX \geq 30$ is made assuming the site is non-variant.

Section Descriptions

The gVCF file contains the following sections:

- ▶ Meta-information lines start with ## and contain metadata, config information, and define the values that the INFO, FILTER, and FORMAT fields can have.
- ▶ The header line starts with # and names the fields that the data lines use. These fields are #CHROM, POS, ID, REF, ALT, QUAL, FILTER, INFO, FORMAT, followed by one or more sample columns.
- ▶ Data lines that contain information about one or more positions in the genome.

If you extract the variant lines from a gVCF file, you produce a conventional variant VCF file.

Field Descriptions

The fixed fields #CHROM, POS, ID, REF, ALT, QUAL are defined in the VCF 4.1 standard provided by the 1000 Genomes Project. The fields ID, INFO, FORMAT, and sample are described in the meta-information.

- ▶ **CHROM:** Chromosome: an identifier from the reference genome or an angle-bracketed ID String ("**<ID>**") pointing to a contig.
- ▶ **POS:** Position: The reference position, with the first base having position 1. Positions are sorted numerically, in increasing order, within each reference sequence CHROM. There can be multiple records with the same POS. Telomeres are indicated by using positions 0 or N+1, where N is the length of the corresponding chromosome or contig.
- ▶ **ID:** Semi-colon separated list of unique identifiers where available. If this ID is a dbSNP variant, it is encouraged to use the rs number. No identifier is present in more than one data record. If there is no identifier available, then the missing value is used.
- ▶ **REF:** Reference bases: A,C,G,T,N; there can be multiple bases. The value in the POS field refers to the position of the first base in the string. For simple insertions and deletions in which either the REF or one of the ALT alleles would otherwise be null/empty, the REF and ALT strings include the base before the event. This modification is reflected in the POS field. The exception is when the event occurs at

position 1 on the contig, in which case they include the base after the event. If any of the ALT alleles is a symbolic allele (an angle-bracketed ID String "<ID>"), the padding base is required. In that case, POS denotes the coordinate of the base preceding the polymorphism.

- ▶ **ALT:** Comma-separated list of alternate non-reference alleles called on at least one of the samples. Options are:

- Base strings made up of the bases A,C,G,T,N
- Angle-bracketed ID String ("<ID>")
- Break-end replacement string as described in the section on break-ends.

If there are no alternative alleles, then the missing value is used.

- ▶ **QUAL:** Phred-scaled quality score for the assertion made in ALT. ie $-10\log_{10}$ probability (call in ALT is wrong). If ALT is "." (no variant), this score is $-10\log_{10}$ p (variant). If ALT is not ".", this score is $-10\log_{10}$ p(no variant). High QUAL scores indicate high confidence calls. Although traditionally people use integer Phred scores, this field is permitted to be a floating point to enable higher resolution for low confidence calls if desired. If unknown, the missing value is specified. (Numeric)
- ▶ **FILTER:** PASS if this position has passed all filters, ie a call is made at this position. Otherwise, if the site has not passed all filters, a semicolon-separated list of codes for filters that fail. gVCF files use the following values:

- *PASS:* position has passed all filters.
- *IndelConflict:* Locus is in region with conflicting indel calls.
- *SiteConflict:* Site genotype conflicts with proximal indel call, which is typically a heterozygous SNV call made inside of a heterozygous deletion.
- *LowGQX:* Locus GQX (minimum of {Genotype quality assuming variant position, Genotype quality assuming non-variant position}) is less than 30 or not present.
- *HighDPFRatio:* The fraction of base calls filtered out at a site is greater than 0.3.
- *HighSNVSB:* SNV strand bias value (SNVSB) exceeds 10. High strand bias indicates a potential high false-positive rate for SNVs.
- *HighSNVHPOL:* SNV contextual homopolymer length (SNVHPOL) exceeds 6.
- *HighREFREP:* Indel contains an allele that occurs in a homopolymer or dinucleotide track with a reference repeat greater than 8.
- *HighDepth:* Locus depth is greater than 3x the mean chromosome depth.

- ▶ **INFO:** Additional information. INFO fields are encoded as a semicolon-separated series of short keys with optional values in the format: <key>=<data>[,<data>]. gVCF files use the following values:

- *END:* End position of the region described in this record.
- *BLOCKAVG_min30p3a:* Non-variant site block. All sites in a block are constrained to be non-variant, have the same filter value, and have all sample values in range [x,y], $y \leq \max(x+3, x*1.3)$. All printed site block sample values are the minimum observed in the region spanned by the block.
- *SNVSB:* SNV site strand bias.
- *SNVHPOL:* SNV contextual homopolymer length.
- *CIGAR:* CIGAR alignment for each alternate indel allele.
- *RU:* Smallest repeating sequence unit extended or contracted in the indel allele relative to the reference. If longer than 20 bases, RUs are not reported.
- *REFREP:* Number of times RU is repeated in reference.
- *IDREP:* Number of times RU is repeated in indel allele.

- ▶ **FORMAT:** Format of the sample field. FORMAT specifies the data types and order of the subfields. gVCF files use the following values:

- *GT:* Genotype.
- *GQ:* Genotype Quality.

- *GQX*: Minimum of {Genotype quality assuming variant position, Genotype quality assuming non-variant position}.
 - *DP*: Filtered base call depth used for site genotyping.
 - *DPF*: Base calls filtered from input before site genotyping.
 - *AD*: Allelic depths for the ref and alt alleles in the order listed. For indels, this value only includes reads that confidently support each allele (posterior probability 0.999 or higher that read contains indicated allele vs all other intersecting indel alleles).
 - *DPI*: Read depth associated with indel, taken from the site preceding the indel.
- **SAMPLE**: Sample fields as defined by the header.

Enrichment Summary Report (*.summary.csv)

The Isaac Enrichment v2.0 app produces an Enrichment Summary Report and the aggregate result in a comma-separated values (CSV) format: *.summary.csv. This report is an overview of statistics for each sample. These files are located in the results folder for each sample and the aggregate summary folder.

A brief description of the metrics is below.



NOTE

PCR duplicate reads are not removed from statistics. Results are not directly comparable to Picard HsMetrics.

Statistic	Definition
Sample ID	IDs of samples reported on in the file.
Sample Name	Names of samples reported on in the file.
Run Folder	Run folders for samples reported on in the file.
Reference Genome	Reference genome selected.
Target Manifest	The target manifest file used for analysis. This file specifies the targeted regions for the aligner and variant caller.
Total Length of Targeted Reference	The total length of the sequenced bases in the target region.
Padding Size	The length of sequence immediately upstream and downstream of the enrichment targets that is included for a padded target.
Total PF Reads	The number of reads passing filter for the sample.
Total Aligned Reads	The total number of reads passing filter present in the data set that aligned to the reference genome.
Percent Aligned Reads	The percentage of reads passing filter that aligned to the reference genome.
Percent Duplicate Paired Reads	Percentage of paired reads that have duplicates.

Statistic	Definition
Targeted Aligned Reads	Number of reads that aligned to the target.
Padded Target Aligned Reads	Number of reads that aligned to the padded target.
Read Enrichment	$100 * (\text{Target aligned reads} / \text{Total aligned reads})$.
Padded Read Enrichment	$100 * (\text{Padded target aligned reads} / \text{Total aligned reads})$.
Total PF Bases	The number of bases passing filter for the sample.
Percent Q30	The percentage of bases with a quality score of 30 or higher.
Percent Q30 Aligned	Percent of bases with a quality score of 30 or higher that aligned to the reference genome.
Total Aligned bases	The total number of bases present in the data set that aligned to the reference genome.
Percent Aligned bases	Percent aligned bases in the target region.
Targeted Aligned bases	Total aligned bases in the target region.
Padded Target Aligned bases	Total aligned bases in the padded target region.
Base Enrichment	$100 * (\text{Total Aligned Bases in Targeted Regions} / \text{Total Aligned Bases})$.
Padded Base Enrichment	$100 * (\text{Total Aligned Bases in Padded Targeted Regions} / \text{Total Aligned Bases})$.
Mean Region Coverage Depth	The total number of aligned bases to the targeted region divided by the targeted region size.
Uniformity of Coverage (Pct > 0.2*mean):	The percentage of targeted base positions in which the read depth is greater than 0.2 times the mean region target coverage depth.
Target Coverage at 1X	Percentage targets with coverage greater than 1X.
Target Coverage at 10X	Percentage targets with coverage greater than 10X.
Target Coverage at 20X	Percentage targets with coverage greater than 20X.

Statistic	Definition
Target Coverage at 50X	Percentage targets with coverage greater than 50X.
Fragment Length Median	Median length of the sequenced fragment. The fragment length is calculated based on the locations at which a read pair aligns to the reference. The read mapping information is parsed from the BAM files.
Fragment Length Min	Minimum length of the sequenced fragment.
Fragment Length Max	Maximum length of the sequenced fragment.
Fragment Length SD	Standard deviation of the sequenced fragment length.
SNVs, Indels, Insertions, Deletions	Total number of variants present in the data set that pass the quality filters.
SNVs (All), Indels (All), Insertions (All), Deletions (All)	Total number of predicted variants in the data set.
SNVs, Indels, Insertions, Deletions (Percent Found in dbSNP)	100*(Number of variants in dbSNP/Number of variants).
SNV Ts/Tv ratio	The number of Transition SNVs that pass the quality filters divided by the number of Transversion SNVs that pass the quality filters. Transitions are interchanges of purines (A, G) or of pyrimidines (C, T). Transversions are interchanges of purine and pyrimidine bases (for example, A to T).
SNVs, Indels, Insertions, Deletions Het/Hom ratio	Number of heterozygous variants/Number of homozygous variants.
SNVs, Insertions, Deletions in Genes	The number of variants that fall into a gene.
SNVs, Insertions, Deletions in Exons	The number of variants that fall into an exon.
SNVs, Insertions, Deletions in Coding Regions	The number of variants that fall into a coding region.

Statistic	Definition
SNVs, Insertions, Deletions in Mature miRNA	The number of variants that fall into a mature microRNA.
SNVs, Insertions, Deletions in UTR Region	The number of variants that fall into an untranslated region (UTR).
SNVs, Insertions, Deletions in Splice Site Region	The number of variants that fall into a splice site region.
Stop Gained SNVs, Insertions, Deletions	The number of variants that cause an additional stop codon.
Stop Lost SNVs, Insertions, Deletions	The number of variants that cause the loss of a stop codon.
Frameshift Insertions, Deletions	The number of variants that cause a frameshift.
Non- synonymous SNVs, Insertions, Deletions	The number of variants that cause an amino acid change in a coding region.
Synonymous SNVs	The number of variants that are within a coding region, but do not cause an amino acid change.

Manifest Output Files

The Isaac Enrichment v2.0 app produces BED and TXT manifest output files that specify the regions that were used in the analysis. If there were any duplicates or overlapping regions, those files contain the corrected version. The BED file can be used in VariantStudio or IGV to highlight the targeted regions.

The two output files are in the aggregate output directory for multi-sample inputs and in the single sample directory for single-sample inputs.

Isaac Enrichment v2.0 Methods

This chapter describes the methods that are used in the Isaac Enrichment v2.0 app.

Isaac Aligner

The Isaac aligner aligns DNA sequencing data, single or paired-end, with read lengths and low error rates using the following steps:

- ▶ **Candidate mapping positions**—Identifies the complete set of relevant candidate mapping positions using a 32-mer seed-based search.
 - ▶ **Mapping selection**—Selects the best mapping among all candidates.
 - ▶ **Alignment score**—Determines alignment scores for the selected candidates based on a Bayesian model.
 - ▶ **Alignment output**—Generates final output in a sorted duplicate-marked BAM file and summary file.
- 1 Come Raczky, Roman Petrovski, Christopher T. Saunders, Ilya Chorny, Semyon Kruglyak, Elliott H. Margulies, Han-Yu Chuang, Morten Källberg, Swathi A. Kumar, Arnold Liao, Kristina M. Little, Michael P. Strömberg and Stephen W. Tanner (2013) Isaac: Ultra-fast whole genome secondary analysis on Illumina sequencing platforms. *Bioinformatics* 29(16):2041-3
bioinformatics.oxfordjournals.org/content/29/16/2041

Candidate Mapping

To align reads, the Isaac aligner first identifies a small but complete set of relevant candidate mapping positions. The Isaac aligner begins with a seed-based search using 32-mers from the extremities of the read as seeds. Isaac performs another search using different seeds for only those reads that were not mapped unambiguously with the first pass seeds.

Mapping Selection

Following a seed-based search, the Isaac aligner selects the best mapping among all the candidates. For paired-end data sets, all mappings where only one end is aligned (called orphan mappings) trigger a local search to find additional mapping candidates. These candidates (called shadow mappings) are defined through the expected minimum and maximum insert size. After optional trimming of low quality 3' ends and adapter sequences, the possible mapping positions of each fragment are compared. This step takes into account pair-end information (when available), possible gaps using a banded Smith-Waterman gap aligner, and possible shadows. The selection is based on the Smith-Waterman score and on the log-probability of each mapping.

Alignment Scores

The alignment scores of each read pair are based on a Bayesian model, where the probability of each mapping is inferred from the base qualities and the positions of the mismatches. The final mapping quality is the alignment score, truncated to 60 for scores above 60, and possibly corrected to known ambiguities in the reference as flagged in the seeds. Following alignment, reads are sorted. Further analysis is performed to identify duplicates and optionally to realign indels.

The alignment scores of each read pair are based on a Bayesian model, where the probability of each mapping is inferred from the base qualities and the positions of the

mismatches. The final mapping quality is the alignment score, truncated to 60 for scores above 60. Following alignment, reads are sorted. Further analysis is performed to identify duplicates and optionally to realign indels.

Alignment Output

After sorting the reads, the Isaac aligner generates compressed binary alignment output files, called BAM (*.bam) files, using the following process:

- ▶ **Marking duplicates**—Detection of duplicates is based on the location and observed length of each fragment. The Isaac aligner identifies and marks duplicates even when they appear on oversized fragments or chimeric fragments. Optical duplicates are already filtered out during RTA processing.
- ▶ **Realigning indels**—The Isaac aligner tracks previously detected indels, over a window large enough for the current read length, and applies the known indels to all reads with mismatches.
- ▶ **Generating BAM files**—The first step in BAM file generation is creation of the BAM record, which contains all required information except the name of the read. The Isaac aligner reads data from base call (BCL) files that were written during base calling on the sequencer to generate the read names. Data are then compressed into blocks of 64 kb or less to create the BAM file.

Isaac Variant Caller

The Isaac Variant Caller identifies single nucleotide polymorphisms (SNPs) and small indels using the following steps:

- ▶ **Read filtering**—Filters out reads failing quality checks.
- ▶ **Indel calling**—Identifies a set of possible indel candidates and realigns all reads overlapping the candidates using a multiple sequence aligner.
- ▶ **SNP calling**—Computes the probability of each possible genotype given the aligned read data and a prior distribution of variation in the genome.
- ▶ **Indel genotypes**—Calls indel genotypes and assigns probabilities.
- ▶ **Variant call output**—Generates output in a VCF file and a compressed genome variant call (gVCF) file. See *VCF Files* on page 22 and *gVCF Files* on page 23 for details.

Indel Candidates

Input reads are filtered by removing any of the following:

- ▶ Reads that failed base calling quality checks.
- ▶ Reads marked as PCR duplicates.
- ▶ Paired-end reads not marked as a proper pair.
- ▶ Reads with a mapping quality less than 20.

Indel Calling

The variant caller proceeds with candidate indel discovery and generates alternate read alignments based on the candidate indels. As part of the realignment process, the variant caller selects a representative alignment to be used for site genotype calling and depth summarization by the SNP caller.

SNP Calling

The variant caller runs a series of filters on the set of filtered and realigned reads for SNP calling without affecting indel calls. First, any contiguous trailing sequence of N base calls is trimmed from the ends of reads. Using a mismatch density filter, reads having an unexpectedly high number of disagreements with the reference are masked, as follows:

- ▶ The variant caller treats each insertion or deletion as a single mismatch.
- ▶ Base calls with more than two mismatches to the reference sequence within 20 bases of the call are ignored.
- ▶ If the call occurs within the first or last 20 bases of a read, the mismatch limit is applied to a 41-base window at the corresponding end of the read.
- ▶ The mismatch limit is applied to the entire read when the read length is 41 or shorter.

Indel Genotypes

The variant caller filters out all bases marked by the mismatch density filter and any N base calls that remain after the end-trimming step. These filtered base calls are not used for site-genotyping but appear in the filtered base call counts in the variant caller output for each site.

All remaining base calls are used for site-genotyping. The genotyping method heuristically adjusts the joint error probability that is calculated from multiple observations of the same allele on each strand of the genome. This correction accounts for the possibility of error dependencies.

This method treats the highest-quality base call from each allele and strand as an independent observation and leaves the associated base call quality scores unmodified. Quality scores for subsequent base calls for each allele and strand are then adjusted. This adjustment is done to increase the joint error probability of the given allele above the error expected from independent base call observations.

Variant Call Output

After the site and indel genotyping methods are complete, the variant caller applies a final set of heuristic filters to produce the final set of non-filtered calls in the output.

The output in the genome variant call (gVCF) file captures the genotype at each position and the probability that the consensus call differs from reference. This score is expressed as a Phred-scaled quality score.

Picard Metrics

Picard is a suite of tools in Java that work with next-generation sequencing data in BAM format. Isaac Enrichment uses the CalculateHsMetrics tool in Picard to compute a set of Hybrid Selection specific metrics from an aligned SAM or BAM file. If a reference sequence is provided, AT/GC dropout metrics are calculated. GC and mean coverage information for every target can also be computed.

For more information, see: picard.sourceforge.net/command-line-overview.shtml

Notes

Technical Assistance

For technical assistance, contact Illumina Technical Support.

Table 5 Illumina General Contact Information

Website	www.illumina.com
Email	techsupport@illumina.com

Table 6 Illumina Customer Support Telephone Numbers

Region	Contact Number	Region	Contact Number
North America	1.800.809.4566	Italy	800.874909
Australia	1.800.775.688	Netherlands	0800.0223859
Austria	0800.296575	New Zealand	0800.451.650
Belgium	0800.81102	Norway	800.16836
Denmark	80882346	Spain	900.812168
Finland	0800.918363	Sweden	020790181
France	0800.911850	Switzerland	0800.563118
Germany	0800.180.8994	United Kingdom	0800.917.0041
Ireland	1.800.812949	Other countries	+44.1799.534000

Safety Data Sheets

Safety data sheets (SDSs) are available on the Illumina website at support.illumina.com/sds.html.

Product Documentation

Product documentation in PDF is available for download from the Illumina website. Go to support.illumina.com, select a product, then click **Documentation & Literature**.



Illumina

San Diego, California 92122 U.S.A.

+1.800.809.ILMN (4566)

+1.858.202.4566 (outside North America)

techsupport@illumina.com

www.illumina.com