# RNA Sequencing with TopHat Alignment v1.0 and Cufflinks Assembly & DE v1.1
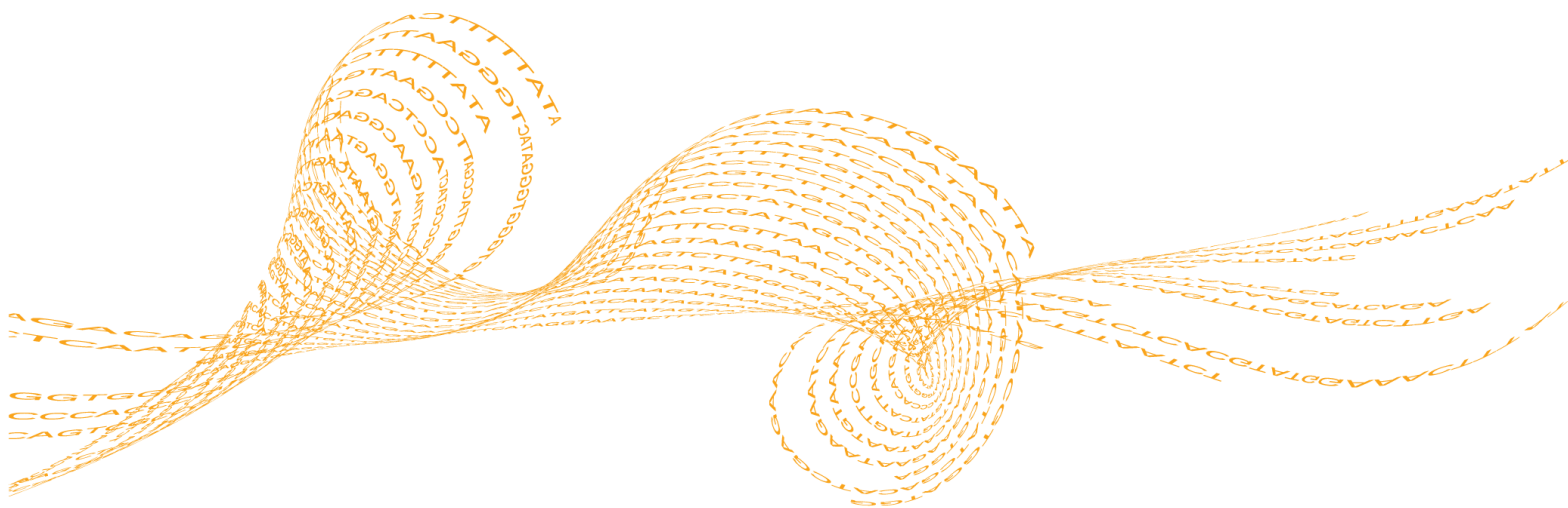# App Guide

For Research Use Only. Not for use in diagnostic procedures.

illumına®

# Introduction

The BaseSpace® Apps *TopHat Alignment v1.0* and *Cufflinks Assembly & DE (Differential Expression) v1.1* are developed for the analysis of RNA sequencing data. The main components of the TopHat and Cufflinks apps are derived from the Tuxedo suite. The 2 apps are designed to perform the following steps:

▸ TopHat Alignment v1.0 runs Bowtie, TopHat, the Isaac Variant Caller, and Cufflinks. The app produces aligned reads, variant calls and FPKM abundance estimates of reference genes and transcripts.

   ▸ See *TopHat App Workflow* on page 24 for more information.

▸ Cufflinks Assembly & DE v1.1 uses previous alignment results produced by the TopHat app as input, and performs novel transcript assembly and/or differential expression analysis.

   ▸ See *Cufflinks App Workflow* on page 25 for more information.

## Versions

The following module versions are used in the TopHat and Cufflinks Apps:

▸ TopHat2 v2.0.7
▸ Bowtie 0.12.9
▸ Cufflinks 2.1.1
▸ Isaac Variant Caller 2.0.5
▸ Picard tools 1.72

## Workflow Requirements

▸ The minimum read length is 35 bp.
▸ The maximum read length is 125 bp.
▸ Each sample should have between 100,000 and 400 million reads.
▸ Supports UCSC hg19 (human), UCSC mm10 (mouse), and UCSC rn5 (rat).

# Set Analysis Parameters—TopHat

1 In BaseSpace, click the **Apps** tab.

2 Click **TopHat Alignment**.

3 From the drop-down list, select **version 1.0**, and then click **Launch** to open the app.

4 Complete the required fields in the TopHat Alignment v1.0 input form.

   ▸ **App Session Name**—Provide the app (analysis) session name. By default, the app session (analysis) name includes the app name, followed by the date and time the app session starts.

   ▸ **Save Results To**—Select the project that stores the app results.

   ▸ **Samples**—Browse to the sample you want to analyze, and select the check box. You can analyze multiple samples.

   ▸ Identify the samples prepared with a stranded sample kit by selecting the **Stranded** checkbox. Samples prepared with the TruSeq Stranded Total RNA Sample Prep Kit and TruSeq Stranded mRNA Sample Prep Kit are stranded.

   ▸ **Reference Genome**—Select the reference genome.

5 Fill out the additional options.

   ▸ **Call Fusions**—When selected, TopHat-fusion is used to detect gene fusions.

   ▸ **Trim TruSeq Adapters**—When selected, the application attempts to trim TruSeq adapters from the FASTQ sequence. Typically, this trimming will be unnecessary as adapter trimming is performed as part of demultiplexing during sample upload. However, if the user did not specify adapter sequences in the sample sheet during upload, this provides a second opportunity to trim the adapters.

6 Click **Continue**.

The TopHat App starts analyzing your sample. When completed, the status of the app session is automatically updated, and you receive an email.

> **NOTE**
> Depending on the size and number of samples, the complete analysis may take between a few hours and several days.

4

# Analysis Output

To view the results, click the **Projects** tab, then the project name, and then the analysis.

After analysis is complete, access the output through the left navigation bar.

▸   **Analysis Info**—Information about the analysis session, including log files.

▸   **Inputs**—Overview of input settings.

▸   **Output Files**—Output files for the sample.

▸   **Sample Analysis Reports**—Analysis reports for each sample.

## Analysis Reports—Summary

The TopHat App provides an overview for all samples on the Summary page.

### Alignment Statistics

| Statistic | Definition |
|---|---|
| Reads | Number and length of reads. |
| Number of Reads | Total number of reads passing filter for this sample. |
| % Total Aligned | Percentage of reads passing filter that aligned to the reference, including abundant reads. |
| % Abundant | Percentage of reads that aligns to abundant transcripts, such as mitochondrial and ribosomal sequences. |
| % Unaligned | Percentage of passed filter reads that do not align to the reference. |
| Median CV Coverage Uniformity | The median coefficient of variation of coverage of the 1000 most highly expressed transcripts, as reported by the CollectRnaSeqMetrics utility from Picard tools. Ideal value = 0. |
| % Stranded | Percentage of reads that align to the correct strand, as reported by the CollectRnaSeqMetrics utility from Picard tools. |

### Insert Length Distribution

The Insert Length Distribution graphs the frequency of each insert length. This histogram is only available for paired-end reads.

Figure 2   Insert Length Distribution



## Alignment Distribution

The Alignment Distribution histogram graphs the number of bases by genomic region (coding, UTR, intron, or intergenic).

Figure 3   Alignment Distribution Histogram



## Transcript Coverage

The Transcript Coverage graphs the coverage of a transcript position by that transcript position as reported by the CollectRnaSeqMetrics utility from Picard tools. The coverage is normalized for the coverage of that transcript. You can get the series names from the tooltips by hovering over the graph.

Transcript Coverage



## Analysis Reports—Sample Pages

The TopHat app provides an overview of statistics per sample on the sample pages.

### Primary Analysis Information

| Statistic | Definition |
|---|---|
| Read Length | Length of reads. |
| Number of Reads | Total number of reads passing filter for this sample. |
| Bases (GB) | The total number of bases for this sample. |
| Q30 Bases (GB) | The total number of bases with a quality score of 30 or higher. |

### Insert Information

Only available for paired-end reads.

| Statistic | Definition |
|---|---|
| Insert 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. |
| Insert Length S.D. | Standard deviation of the sequenced fragment length. |
| Duplicates (% Reads) | Percentage of paired reads that have duplicates., from a sub-sampled set of 4 M reads. |

## Alignment Quality

| Statistic | Definition |
|---|---|
| Total Aligned Reads (% Reads) | The percentage of reads passing filter that aligned to the reference genome. |
| Abundant Reads (% Reads) | Percentage of reads that aligns to abundant transcripts, such as mitochondrial and ribosomal sequences. |
| Unaligned Reads (% Reads) | Percentage of reads passing filter that are not aligned to the reference. |
| Reads with spliced alignment (% Aligned Reads) | Percentage of aligned reads with a spliced alignment. |
| Reads aligned at multiple loci (% Aligned Reads) | Percentage of aligned reads to multiple loci. |

## Coverage Uniformity

| Statistic | Definition |
|---|---|
| Median CV | The median coefficient of variation of coverage of the 1000 most highly expressed transcripts, as reported by the CollectRnaSeqMetrics utility from Picard tools. Ideal value = 0. |
| Median 3' | The median uniformity of coverage of the 1000 most highly expressed transcripts at the 3' end, as reported by the CollectRnaSeqMetrics utility from Picard tools. 3' bias is calculated per transcript as: mean coverage of the 3' most 100 bases divided by the mean coverage of the whole transcript. |
| Median 5' | The median uniformity of coverage of the 1000 most highly expressed transcripts at the 5' end, as reported by the CollectRnaSeqMetrics utility from Picard tools. 5' bias is calculated per transcript as: mean coverage of the 5' most 100 bases divided by the mean coverage of the whole transcript. |
| Reads aligned to correct strand | Percentage of reads that align to the correct strand, as reported by the CollectRnaSeqMetrics utility from Picard tools. |

## Alignment Information

This table presents metrics for the categories in two columns:

▸ Fold Coverage: The total number of bases in the category divided by the size of the entire category.

▸ % Bases: The total number of bases aligned to this region relative to the total number of aligned bases.

These metrics are reported for the following categories, derived from RefSeq:

| Statistic | Definition |
|---|---|
| Coding | Metrics based on coding bases. |

| Statistic | Definition |
|---|---|
| UTR | Metrics based on bases in untranslated regions (UTR). |
| Intron | Metrics based on bases in introns. |
| Intergenic | Metrics based on bases in intergenic regions. |

## Variant Calls

| Statistic | Definition |
|---|---|
| Homozygous Reference | Number of homozygous reference calls. |
| Heterozygous | Number of heterozygous variant calls. |
| Homozygous variant | Number of homozygous variant calls. |
| SNV | Total number of SNVs detected for the sample. |
| Indel | The number of indels detected for the sample. |
| $T_n/T_v$ | 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). |

### Histograms

The sample pages also contain the Insert Length Distribution, Alignment Distribution, and Transcript Coverage histograms, which are described in *Analysis Reports—Summary* on page 5

### Important Files

The Important Files section provides direct access to the most important output files.

For detailed descriptions of these files, see:
- *BAM Files* on page 11
- *VCF Files* on page 11
- *gVCF Files* on page 12
- *FPKM Files* on page 22
- *Candidate Fusion List* on page 16

## Analysis Info

The Analysis Info page displays the analysis settings and execution details.

| Row Heading | Definition |
|---|---|
| Name | Name of the analysis session. |
| Application | App that generated this analysis. |

| Row Heading | Definition |
|---|---|
| Date Started | Date and time the analysis session started. |
| Date Completed | Date and time the analysis session completed. |
| Duration | Duration of the analysis. |
| Session Type | Number of nodes used. |
| Status | Status of the analysis session. The status shows either Running or Complete. |

## TopHat Status

The status of the TopHat app session can have the following values:
- Preparing data
- Aligning
- Variant calling
- Estimating reference abundances
- Calculating metrics
- Finalizing results

Depending on the size and number of samples, the complete analysis may take between a few hours and several days.

## TopHat Inputs Overview

The TopHat App provides an overview of the input samples and settings on the Inputs page.

| Statistic | Definition |
|---|---|
| App Session Name | Name of the app session. |
| Save Results To | The project that stores the app results. |
| Samples | Samples selected for this analysis. Checkbox indicates if the samples are stranded. |
| Reference Genome | Reference genome and gene model selected. |
| Call Fusions | If selected, TopHat-fusion is used to detect gene fusions. |
| Trim TruSeq Adapters | If selected, the application trims TruSeq adapters from the FASTQ sequence. |

## TopHat Output Files

This section describes the main output files of the TopHat App:
- *BAM Files* on page 11
- *VCF Files* on page 11
- *gVCF Files* on page 12
- *FPKM Files* on page 22
- *Coverage.BedGraph.gz Files* on page 15

The output files are located in the following locations:

```
<Analysis Result>|filtered
```
▸ <sample name>.abundant_alignments.bam – Alignments of the reads against abundant sequences.
```
<Analysis Result>|alignments
```
▸ <sample name>.alignments.bam – Alignments of filtered reads against the genome (and transcriptome).
▸ <sample. name>.coverage.bedGraph.gz – Genome coverage from aligned RNA-seq reads
▸ <sample name>.junctions.bed – BED file describing junctions discovered in this data set in TopHat format.
```
<Analysis Result>|cufflinks|reference
```
▸ <sample name>.isoforms.fpkm_tracking – Abundance estimates for individual transcripts defined in genome reference.
▸ <sample name>.genes.fpkm_tracking – Abundance estimates for genes defined in genome reference.
```
<Analysis Result>|metrics:
```
▸ Contains CSV and json metrics files used to generate the summary and individual reports.
```
<Analysis Result>|tophat_fusion
```
▸ result.html – HTML summary page describing the discovered fusions.
```
<Analysis Result>|variants
```
▸ <sample_id>.genome.vcf.gz – Variant calls from starling in gVCF format.

## 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*.

## gVCF Files

This application also produces the Genome Variant Call Format file (gVCF). gVCF was developed to store sequencing information for both variant and nonvariant 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 nonvariant 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.

### Nonvariant Blocks Using END Key

Contiguous nonvariant 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 nonvariant 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], $y <= x+max(3,x*0.3)$. This range restriction applies to each of the sample values printed in the final block record.

### Indel Regions

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

All deletions:

▸ Sites inside any deletion are marked with the deletion filters, in addition to any filters that have already been applied to the site.

- Sites inside deletions cannot have a genotype or alternate allele quality score higher than the corresponding value from the enclosing indel.

Heterozygous deletions:

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

Homozygous deletions:

- Homozygous reference and no-call sites inside homozygous deletions have genotype ".".
- Sites inside homozygous deletions that have a nonreference 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 2 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 Nonvariant Sites

The gVCF file uses an adapted version of genotype quality for variant and nonvariant 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 nonvariant}. The reason for using this value is to allow a single value to be used as the primary quality filter for both variant and nonvariant 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 >= 30 is made assuming the site is variant.
- An assertion that a site is a nonreference genotype at GQX >= 30 is made assuming the site is nonvariant.

## Section Descriptions

The gVCF file contains the following sections:

- Metainformation 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 metainformation.

- **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**: Semicolon 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 1 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 1 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 nonreference alleles called on at least 1 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 a heterozygous deletion.
  - *LowGQX*: Locus GQX (minimum of {Genotype quality assuming variant position,Genotype quality assuming nonvariant 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*: nonvariant site block. All sites in a block are constrained to be nonvariant, have the same filter value, and have all sample values in range [x,y], y <= 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 nonvariant 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.

## FPKM Files

Fragments Per Kilobase of sequence per Million mapped reads (FPKM) normalizes the number of aligned reads by the size of the sequence feature and the total number of mapped reads.

In each output directory, this app creates the following output files:

▸ *genes.fpkm_tracking*—Quantifies the expression of genes specified in the GTF annotation file.

▸ *isoforms.fpkm_tracking*—Quantifies the expression of transcripts specified in the GTF annotation file.

## Coverage.BedGraph.gz Files

The *.coverage.BedGraph.gz file can summarize coverage plotting in browsers (IGV, UCSC genome browser); it is much smaller than the BAM file.

For more information, see http://genome.ucsc.edu/FAQ/FAQformat.html#format1.8.

## Junctions.bed Files

The *.junctions.bed file is a bed file with junctions and the coverage of these for visualization in a browser.

For more information, see http://genome.ucsc.edu/FAQ/FAQformat.html#format1

## Candidate Fusion List

The candidate fusion list is an html page that provides an overview of all candidate fusions. The top tables provide the genes that are fused, the chromosomes and positions involved, as well as the number of reads that provide the evidence. This is explained in the table description. In addition, genomic and nucleotide blast searches are provided with the left and right sides of the fusion.
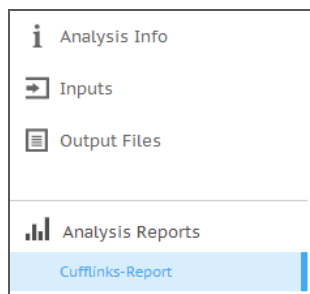
1 In BaseSpace, click the **Apps** tab.

2 Click **Cufflinks Assembly & DE**.

3 From the drop-down list, select **version 1.1**, and then click **Launch** to open the app.

4 Complete the required fields in the Cufflinks Assembly & DE v1.1 input form.
   ▶ **App Session Name**—Provide the app (analysis) session name. By default, the app session (analysis) name includes the app name, followed by the date and time the app session starts.
   ▶ **Save Results To**—Select the project that stores the app results.

5 Fill out the TopHat alignments selection criteria. These settings are used to filter the results available in the app result choosers below. This ensures that only compatible app results are selected for analysis.
   ▶ **Reference Genome**—Select the reference genome.
   ▶ **Stranded**—Indicate if samples were stranded.

6 Fill out the additional options.
   ▶ **Novel Transcript Assembly**—Select if you want Cufflinks to perform novel transcript assembly.

7 Fill out the control group information.
   ▶ **Group Label**—Provide the control group label. Default name is control.
   ▶ **Select App Result**—Browse to the app result from the TopHat App that you want to use as control, and select the checkbox. You can use multiple app results as control.
   ▶ **Adjust transcript assembly for samples without polyA selection**—Select if performing transcript assembly and your samples were generated without polyA selection. This will enable additional filtering of intronic alignments during transcript assembly to reduce the false identification of transcripts.

8 Fill out the comparison group information.
   ▶ **Group Label**—Provide the comparison group label. Default name is comparison.
   ▶ **Select App Result**—Browse to the app result from the TopHat app that you want to analyze against the control, and select the check box. You can analyze multiple app results.
   ▶ **Adjust transcript assembly for samples without polyA selection**—Select if performing transcript assembly and your samples were generated without polyA selection.

9 Click **Continue**.

The Cufflinks App starts analyzing your sample. When completed, the status of the app session is automatically updated, and you receive an email.

# Cufflinks App Output

This chapter describes the output that is produced by the Cufflinks App. To view the results, click the **Projects** button, then the project, and then the analysis.

Figure 5   Cufflinks App Output Navigation Bar



When the analysis is complete, you can access your output through the left navigation bar.

▸ **Analysis Info**—Overview of the analysis. See *Analysis Info* on page 20.

▸ **Inputs**—Overview of the input samples and settings. See *Cufflinks Inputs Overview* on page 21.

▸ **Output Files**—Access to the output files, organized by sample and app session. See *Cufflinks Output Files* on page 21.

▸ **Cufflinks Report**—Access to analysis metrics for the aggregate results. See *Cufflinks Report* on page 18.

## Cufflinks Report

The Cufflinks App Result Page provides an overview of statistics of the Cufflinks App Session. A brief description of the metrics is below.

### Overview

Provides links to the control and comparison samples, as well as to the FPKM tables for genes and transcripts (see *FPKM Files* on page 22).

### Assembly

The Assembly table provides metrics derived from the assembly of the control and comparison samples, as well as the two samples merged.

| Statistic | Definition |
|---|---|
| Gene Count | Number of identified genes per sample. |
| Transcript Count | Number of identified transcipts per sample. |
| Link to gene models | Link to the GTF results (see *GTF Files* on page 22). |
| Relation to reference transcripts | Link to reference information. |
| Equal (=) | Complete match of intron chain. |

| Statistic | Definition |
|---|---|
| Potentially novel (j) | Potentially novel isoform (fragment): at least one splice junction is shared with a reference transcript. |
| Unknown, intergenic (u) | Unknown, intergenic transcript. |
| Overlap with opposite-strand exon (x) | Exonic overlap with reference on the opposite strand. |
| Other | Other types of transcripts. |

## Differential Expression

| Statistic | Definition |
|---|---|
| Gene Count | Number of identified genes assessed in differential expression calculation. |
| ΔGene Count | Number of genes for which differential expression is significant. |
| Transcript Count | Number of identified transcripts assessed in differential expression calculation. |
| ΔTranscript Count | Number of transcripts for which differential expression is significant. |
| Selected CuffDiff results | Links to selected Cuffdiff results. |

### Sample Correlation

Sample correlation shows the similarity of the samples, based on the correlation of expression levels.

This is determined the following way:
▸ Any gene with a reported confidence interval for the FPKM estimate that is greater than 25% of the estimated FPKM value is filtered out.
▸ Any gene with an FPKM value less than 1.0 is filtered out.
▸ For each sample pair, we calculate a correlation over the shared genes that were not filtered out in the above step from the natural log of the estimated FPKM values.

The sample correlation section includes a heat plot and clustering dendrogram.

### Differential Expression Gene Browser

The Differential Expression Gene Browser shows an interactive scatter plot of the log2 (FPKM) counts of genes for two groups of samples.

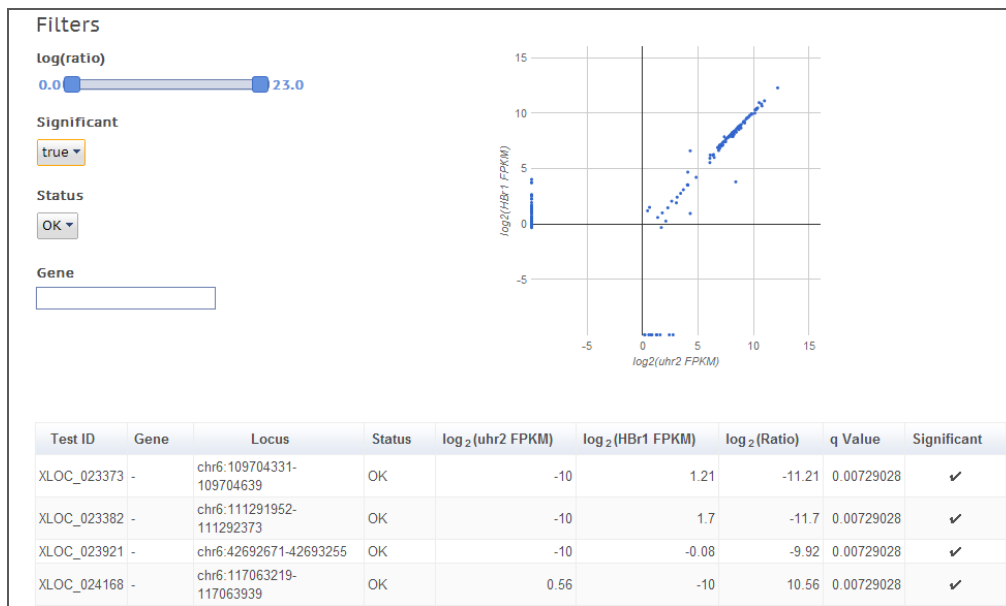You can filter the results by the following metrics:
▸ The absolute value of log2(ratio) of expression levels between two sample groups.
▸ Significance of the differential expression.
▸ Test status
  ▸ OK—Test successful.
  ▸ NOTEST—Not enough alignments for testing.
  ▸ FAIL—When an ill-conditioned covariance matrix or other numerical exception prevents testing.
  ▸ HIDATA—Too many fragments in locus.

- LOWDATA—The region being tested was either too complex or too shallowly sequenced to support a reliable calculation of abundance.
  - Gene—Allows you to search for a particular gene in the plot and the gene table below the scatter plot.

The gene table below the scatter plot shows the individual gene results. If you click a gene, the corresponding dot is circled in the scatter plot. Vice versa, if you click any of the dots in the scatter plot, the gene is highlighted in the gene table. You can also sort the data by clicking a column header.

Figure 6 Differential Expression Gene Browser



| Test ID | Gene | Locus | Status | log$_2$(uhr2 FPKM) | log$_2$(HBr1 FPKM) | log$_2$(Ratio) | q Value | Significant |
|---|---|---|---|---|---|---|---|---|
| XLOC_023373 | - | chr6:109704331-109704639 | OK | -10 | 1.21 | -11.21 | 0.00729028 | ✔ |
| XLOC_023382 | - | chr6:111291952-111292373 | OK | -10 | 1.7 | -11.7 | 0.00729028 | ✔ |
| XLOC_023921 | - | chr6:42692671-42693255 | OK | -10 | -0.08 | -9.92 | 0.00729028 | ✔ |
| XLOC_024168 | - | chr6:117063219-117063939 | OK | 0.56 | -10 | 10.56 | 0.00729028 | ✔ |

## Analysis Info

The Analysis Info page displays the analysis settings and execution details.

| Row Heading | Definition |
|---|---|
| Name | Name of the analysis session. |
| Application | App that generated this analysis. |
| Date Started | Date and time the analysis session started. |
| Date Completed | Date and time the analysis session completed. |
| Duration | Duration of the analysis. |
| Session Type | Number of nodes used. |
| Status | Status of the analysis session. The status shows either Running or Complete. |

## Cufflinks Status

The status of the Cufflinks app session can have the following values:
- Preparing data
- Assembling transcripts

20

- Merging samples
- Calculating differential expression
- Calculating metrics
- Finalizing results

Depending on the size and number of samples, the complete analysis may take between a few hours and several days.

## Cufflinks Inputs Overview

The Cufflinks App provides an overview of the input app results and settings on the Inputs page.

| Statistic | Definition |
| --- | --- |
| App Session Name | Name of the app session. |
| Save Results To | The project that stores the app results. |
| Reference Genome | Reference genome selected. |
| Stranded | Indicates if samples were stranded. |
| Novel Transcript Assembly | Selected when the Cufflinks App is set to perform novel transcript assembly. |
| Group Label | The group label. |
| TopHat Alignment App Result(s) | App results selected for that group. |
| Adjust transcript assembly for samples without polyA selection | Selected if the samples for that group were generated without polyA selection. |

## Cufflinks Output Files

This section describes the main output files of the Cufflinks App:

- *FPKM Files* on page 22
- *GTF Files* on page 22
- *DIFF Files* on page 22

The output files for both apps are located in the following locations:

`<Analysis Result>/control` (if novel transcript assembly requested)

- cuffmerge/control.merged.gtf—Assembled and merged transcripts from designated control samples
- cuffmerge/quantification/<sample name>/genes.fpkm_tracking—Per-gene FPKM estimates for sample against newly assembled transcripts in control.merged.gtf
- cuffmerge/quantification/<sample name>/isoforms.fpkm_tracking—Per-transcript FPKM estimates for sample against newly assembled transcripts in control.merged.gtf
- metrics - Metrics on assembled control transcripts

`<Analysis Result>/comparison` (if novel transcript assembly requested)

- cuffmerge/comparison.merged.gtf—Assembled and merged transcripts from designated comparison samples
- cuffmerge/quantification/<sample name>/genes.fpkm_tracking—Per-gene FPKM estimates for sample against newly assembled transcripts in comparison.merged.gtf

- cuffmerge/quantification/<sample name>/isoforms.fpkm_tracking—Per-transcript FPKM estimates for sample against newly assembled transcripts in comparison.merged.gtf
- metrics—Metrics on assembled comparison transcripts

`<Analysis Result>/differential/cuffmerge` (if novel transcript assembly requested)

- control_vs_comparison.merged.gtf—Assembled and merged transcript for use in differential expression

`<Analysis Result>/differential/cuffdiff`—CuffDiff output from comparison of control and comparison samples

- control_vs_comparison.gene_exp.diff—Differential expression results for genes
- control_vs_comparison.isoform_exp.diff—Differential expression results for transcripts

## FPKM Files

Fragments Per Kilobase of sequence per Million mapped reads (FPKM) normalizes the number of aligned reads by the size of the sequence feature and the total number of mapped reads.

In each output directory, this app creates the following output files:

- *genes.fpkm_tracking*—Quantifies the expression of genes specified in the GTF annotation file.
- *isoforms.fpkm_tracking*—Quantifies the expression of transcripts specified in the GTF annotation file.

## GTF Files

The *merged.gtf* file provides the set of merged transcripts. Each line contains an annotation field ("class_code") that describes the nature of the overlap of this transcript with transcripts from the reference annotation. The table below, taken from the cufflinks manual, (cole-trapnell-lab.github.io/cufflinks/manual/), provides a description of the possible class codes.

| = | Match |
|---|---|
| j | New isoform |
| e | A single exon transcript overlapping a reference exon and at least 10 bp of a reference intron, indicating a possible pre-mRNA fragment |
| i | A single exon transcript falling entirely with a reference intron |
| r | Repeat, currently determined by looking at the reference sequence and applied to transcripts where at least 50% of the bases are lower case |
| p | Possible polymerase run-on fragment |
| u | Unknown, intergenic transcript |
| o | Unknown, generic overlap with reference |
| . | Tracking file only, indicates multiple classifications |

Transcripts annotated with the 'i', 'j', 'u', or 'o' class codes represent novel transcripts of potential interest.

## DIFF Files

The Cufflinks App creates several DIFF files that describe the differential expression. This tab delimited file lists the results of differential expression testing between samples for spliced transcripts, primary transcripts, genes, and coding sequences.

- *isoform_exp.diff*—Transcript differential FPKM.
- *gene_exp.diff*—Gene differential FPKM. Tests differences in the summed FPKM of transcripts sharing each `gene_id`.
- *tss_group_exp.diff*—Primary transcript differential FPKM. Tests differences in the summed FPKM of transcripts sharing each `tss_id`.
- *cds_exp.diff*—Coding sequence differential FPKM. Tests differences in the summed FPKM of transcripts sharing each `p_id` independent of `tss_id`.

The DIFF file has the following format.

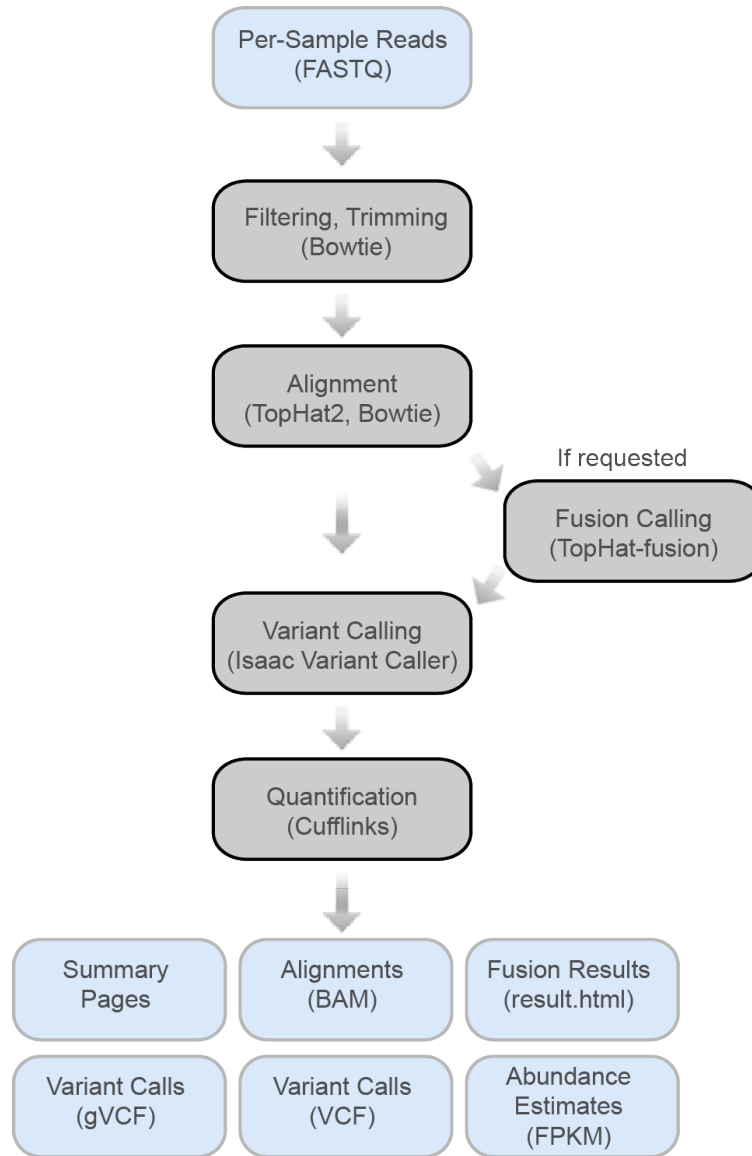| Column number | Column name | Example | Description |
|---|---|---|---|
| 1 | Tested id | XLOC_000001 | A unique identifier describing the transcipt, gene, primary transcript, or CDS being tested |
| 2 | gene | Lypla1 | The gene_name(s) or gene_id(s) being tested |
| 3 | locus | chr1:4797771-4835363 | Genomic coordinates for easy browsing to the genes or transcripts being tested. |
| 4 | sample 1 | Liver | Label (or number if no labels provided) of the first sample being tested |
| 5 | sample 2 | Brain | Label (or number if no labels provided) of the second sample being tested |
| 6 | Test status | NOTEST | Can be one of the following: OK (test successful) NOTEST (not enough alignments for testing) LOWDATA (too complex or shallowly sequenced) HIDATA (too many fragments in locus) FAIL, when an ill-conditioned covariance matrix or other numerical exception prevents testing. |
| 7 | FPKMx | 8.01089 | FPKM of the gene in sample x |
| 8 | FPKMy | 8.551545 | FPKM of the gene in sample y |
| 9 | log2 (FPKMy/FPKMx) | 0.06531 | The (base 2) log of the fold change y/x |
| 10 | test stat | 0.860902 | The value of the test statistic used to compute significance of the observed change in FPKM |
| 11 | p value | 0.389292 | The uncorrected p-value of the test statistic |
| 12 | q value | 0.985216 | The FDR-adjusted p-value of the test statistic |
| 13 | significant | no | Can be either "yes" or "no", depending on whether p is greater then the FDR after Benjamini-Hochberg correction for multiple-testing |

For more information, see the Cufflinks manual at cole-trapnell-lab.github.io/cufflinks/manual/.

# RNAseq Methods

## TopHat App Workflow

1 **Filtering**—The first step of the workflow is to filter the input reads against abundant sequences, such as mitochondrial or ribosomal sequences, as defined by iGenomes (support.illumina.com/sequencing/sequencing_software/igenome.html). The workflow has an internal library of such sequences that it uses for this purpose, including mitochondrial and ribosomal sequences. Bowtie (see *Bowtie* on page 26) performs this alignment. Only sequences that do not align against abundant sequences are passed through to the next phase of the analysis. Read pairs are filtered if at least 1 read aligns to an abundant sequence. This filtering step also performs trimming of 2 bases from the 5' end of the read. This is due to a consistent observation of a high mismatch rate from these two bases in RNA-seq libraries.

2 **Alignment**—TopHat2 (see *TopHat* on page 26), using the Bowtie 1 aligner, performs a spliced alignment of the filtered reads against the genome. Based on the user-specified genome, TopHat uses a list of known transcripts to align reads against known transcripts and splice junctions.

3 **Fusion Calling**—If requested, Tophat-fusion detects gene fusions. Fusion calling occurs in 2 steps. First, TopHat2 is run in a mode to allow the detection of fused alignments. Then, a post-alignment analysis script identifies candidate fusion genes from these fusion alignments. This step does not work with STAR alignments.

4 **Variant Calling**—The Isaac Variant Caller (see *Isaac Variant Caller* on page 26) performs variant calling, producing gVCF output. For stranded library preps, the strand bias filter will be disabled. In addition, the workflow makes use of the - bsnp-diploid-het-bias parameter to expand the allowable range for the heterozygous variant call, in order to account for allele-specific expression.

5 **Quantification**—Cufflinks (see *Cufflinks* on page 28) is used for quantification of reference genes and transcripts.
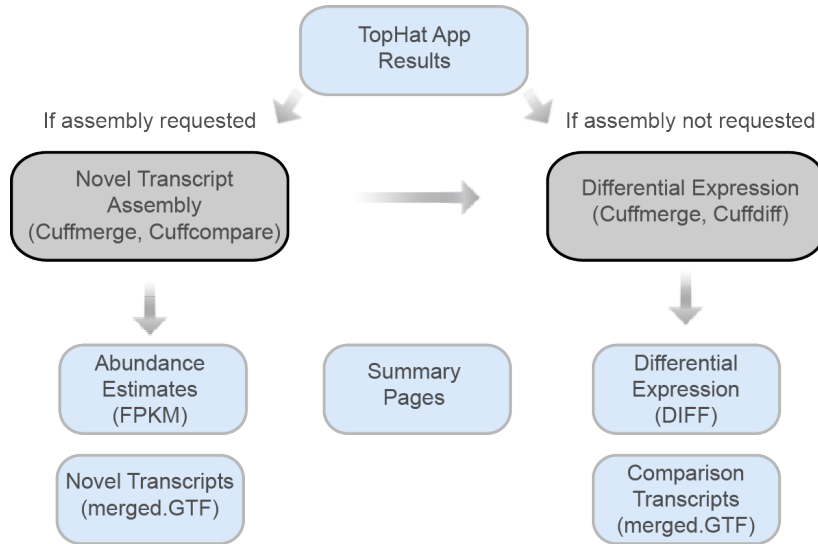
TopHat App Workflow



## Cufflinks App Workflow

1   **Novel transcript assembly**—If novel transcript assembly is requested, transcripts are first assembled and quantified independently for each sample. Cuffmerge merges the assemblies between the samples of each group and combines them with the known gene models from annotation. Cuffcompare then provides annotations of novel transcripts, for example, if they overlap with a known transcript. Lastly abundances will be re-estimated for each transcript, known or novel.

2   **Differential expression**—Cuffdiff calculates differential expression between these two sample groups, estimating variance from the individual samples supplied for each sample group. Cuffdiff calculates a variety of information regarding differential expression, which are summarized in the on-line manual (cole-trapnell-lab.github.io/cufflinks/cuffdiff/index.html#cuffdiff-output-files)) for that tool. Cuffdiff needs a consensus set of transcripts that is compared between different sample groups. If no novel transcript detection is performed, this consensus set is

simply the reference annotation. If novel transcript assembly is requested, then the set of transcripts from each sample must first be merged together. The Cuffmerge tool does this, while also merging in the known gene models from annotation.

Figure 8  Cufflinks App Workflow



# Bowtie

Bowtie[1] is an ultrafast, memory-efficient aligner designed to quickly align large sets of short reads to large genomes. Bowtie indexes the genome to keep its memory footprint small: for the human genome, the index is typically about 2.2 GB for single-read alignment or 2.9 GB for paired-end alignment. Multiple processors can be used simultaneously to achieve greater alignment speed.

Bowtie forms the basis for other tools like TopHat, a fast splice junction mapper for RNA-seq reads, and Cufflinks, a tool for transcriptome assembly and isoform quantitation from RNA-seq reads.

# TopHat

TopHat[2] is a fast splice junction mapper for RNA-Seq reads. It aligns RNA-Seq reads to mammalian-sized genomes using the ultra high-throughput short read aligner Bowtie, and then analyzes the mapping results to identify splice junctions between exons.

TopHat is a collaborative effort between the University of Maryland Center for Bioinformatics and Computational Biology and the University of California, Berkeley Departments of Mathematics and Molecular and Cell Biology.

# Isaac Variant Caller

Strelka identifies single nucleotide variants (SNVs) 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.
- **SNV 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.

26

&#9657; **Variant call output**—Generates output in a vcf and a compressed genome variant call (gVCF) file. See *VCF Files* on page 11 and *gVCF Files* on page 12 for details.

## Indel Candidates

Input reads are filtered by removing any of the following reads:
&#9657; Reads that failed base calling quality checks.
&#9657; Reads marked as PCR duplicates.
&#9657; Paired-end reads not marked as a proper pair.
&#9657; Reads with a mapping quality < 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 SNV caller.

## SNV Calling

The variant caller runs a series of filters on the set of filtered and realigned reads for SNV 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:

&#9657; The variant caller treats each insertion or deletion as a single mismatch.
&#9657; Base calls with more than 2 mismatches to the reference sequence within 20 bases of the call are ignored.
&#9657; 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.
&#9657; 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 SNV and indel genotyping methods are complete, the variant caller applies a final set of heuristic filters to produce the final set of 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.

## Cufflinks

Cufflinks[3] assembles transcripts, estimates their abundances, and tests for differential expression and regulation in RNA-Seq samples. It accepts aligned RNA-Seq reads and assembles the alignments into a parsimonious set of transcripts. Cufflinks then estimates the relative abundances of these transcripts based on how many reads support each one.

Cufflinks is a collaborative effort between the Laboratory for Mathematical and Computational Biology, led by Lior Pachter at UC Berkeley, Steven Salzberg's group at the University of Maryland Center for Bioinformatics and Computational Biology, and Barbara Wold's lab at Caltech.

## References

1   Langmead B, Trapnell C, Pop M, Salzberg SL. (2009) Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biol 10(3):R25.

2   Trapnell C, Pachter L, Salzberg SL.(2009) TopHat: discovering splice junctions with RNA-Seq. Trapnell C, Pachter L, Salzberg SL. Bioinformatics 25(9):1105-11.

3   Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ, Pachter L. (2010) Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. Nat Biotechnol 28(5):511-5.

# Revision History

| Document | Date | Description of Change |
|---|---|---|
| Document # 15050962 v01 | January 2016 | • Update cover with App version information.<br>• Revise instructions in Set Analysis Parameter sections to reflect updated BaseSpace user interface. |

Notes

# Technical Assistance

For technical assistance, contact Illumina Technical Support.

Table 1   Illumina General Contact Information

| | |
|---|---|
| **Website** | www.illumina.com |
| **Email** | techsupport@illumina.com |

Table 2   Illumina Customer Support Telephone Numbers

| Region | Contact Number | Region | Contact Number |
|---|---|---|---|
| North America | 1.800.809.4566 | Japan | 0800.111.5011 |
| Australia | 1.800.775.688 | Netherlands | 0800.0223859 |
| Austria | 0800.296575 | New Zealand | 0800.451.650 |
| Belgium | 0800.81102 | Norway | 800.16836 |
| China | 400.635.9898 | Singapore | 1.800.579.2745 |
| Denmark | 80882346 | Spain | 900.812168 |
| Finland | 0800.918363 | Sweden | 020790181 |
| France | 0800.911850 | Switzerland | 0800.563118 |
| Germany | 0800.180.8994 | Taiwan | 00806651752 |
| Hong Kong | 800960230 | United Kingdom | 0800.917.0041 |
| Ireland | 1.800.812949 | Other countries | +44.1799.534000 |
| Italy | 800.874909 | | |

**Safety data sheets (SDSs)**—Available on the Illumina website at support.illumina.com/sds.html.

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