

Third-Party Analysis Software and Utilities

Several applications are available for analysis and visualization of Illumina sequencing data.

Highlights

- Several commonly-used feature-rich analysis packages support Illumina workflows
- Some third-party tools are a good starting point for next-generation sequencing (NGS) users with limited informatics support
- Third-party visualization and process quality control tools are useful in both research and production environments

Why Use a Third-Party Analysis Tool?

The number of available software packages and third-party utilities have continued to grow along with the types of kits and applications available for Illumina sequencing systems. Some of these packages have become commonly-used tools by next-generation sequencing users, and represents a good starting point for those who might have limited analysis experience or informatics support.

This document describes the following third-party analysis tools that support Illumina sequencing data. Most of the tools listed are open source tools, unless otherwise noted¹.

Alignment	BWA Bowtie MAQ Novoalign ¹
De novo Assembly	ABYSS ALLPATHS-LG SSAKE Velvet
SNP/Indel Discovery	GATK SAM Tools SOAPsnp
Genome Annotation Browser	Anno-J Integrative Genomics Viewer UCSC Genome Browser
ChIP-Seq Applications	CisGenome MACS PeakSeq
Transcriptome Analysis	Partek Genomic Suite ¹ RNA-Star TopHat
Metagenomics	MEGAN Qiime
Integrated Solutions	Avadis NGS ¹ CLCBio Genomics Workbench ¹ Galaxy
Utilities and QC Tools	FastQC Picard ¹

¹ Commercially-licensed tool

Alignment

Bowtie

bowtie-bio.sourceforge.net/index.shtml

Bowtie is an ultra-fast, memory-efficient short read aligner able to align large sets of short DNA sequences (reads) to large genomes. Bowtie is capable of aligning 35 bp reads to the human genome at a rate of 25 million reads per hour on a typical workstation. Bowtie indexes the genome with a Burrows-Wheeler index to keep its memory footprint small. For example, for the human genome, the index is typically about 2.2 GB (for unpaired alignment) or 2.9 GB (for paired-end alignment).

- Supports multiples processors
- Command line for Windows, Mac OS X, Linux, and Solaris
- Developed with C++
- Available under academic license from Johns Hopkins University and listed authors

BWA

bio-bwa.sourceforge.net

BWA maps low-divergent sequences against a large reference genome, such as the human genome. It consists of three algorithms: BWA-backtrack, BWA-SW, and BWA-MEM. The first algorithm is designed for Illumina sequence reads up to 100 bp, while the remaining two algorithms are better suited for longer sequences ranging from 70 bp to 1 Mbp.

- Supports gapped alignment and paired-end mapping
- Developed with C
- Available under GNU General Public License version 3.0 (GPLv3)

MAQ

maq.sourceforge.net

MAQ builds assemblies by mapping short reads to reference sequences generated by next-generation sequencing instruments. MAQ first aligns reads to reference sequences and then calls the consensus. At the mapping stage, MAQ performs ungapped alignment. For single-end reads, MAQ is able to find all hits with up to two or three mismatches, depending on a command-line option. For paired-end reads, it finds all paired hits with one of the two reads containing up to one mismatch. At the assembling stage, MAQ calls the consensus based on a statistical model.

- Supports single reads and paired-end reads
- Available under GNU General Public License version 3.0 (GPLv3)

Novoalign

www.novocraft.com/main/index.php

Novocraft Technologies is a bioinformatics company specializing in the development of fast and accurate tools for NGS analysis. The Novoalign package is an aligner for single-ended and paired-end reads from Illumina sequencers, and finds global optimum alignments using the full Needleman-Wunsch algorithm, including affine gap penalties.

- Uses base qualities at all steps in the alignment for greater accuracy
- Generates native SAM/BAM alignment output format
- Available under free, evaluation, and commercial licensing terms from Novocraft

De novo Assembly

ABYSS

www.bcgsc.ca/platform/bioinfo/software/abyss

Assembly By Short Sequences (ABYSS) is a *de novo* sequence assembler that is designed for very short reads. The single-processor version is useful for assembling genomes up to 40–50 Mbases in size. The parallel version is implemented using MPI and capable of assembling larger genomes.

- Merges overlapping read pairs
- Uses sequence overlap graph to layout and merge contigs
- Attempts to fill scaffold gap with consensus of all paths between contigs
- Developed with C++

ALLPATHS-LG

www.broadinstitute.org/software/allpaths-lg/blog

ALLPATHS is a short read genome assembler from the Computational Research and Development group at the Broad Institute. ALLPATHS-LG is a whole-genome shotgun assembler that generates high-quality genome assemblies using short reads (~100 bp) such as those produced by the new generation of sequencers. The significant difference between ALLPATHS and traditional assemblers such as Arachne is that ALLPATHS assemblies are not necessarily linear, but instead are presented in the form of a graph. This graph representation retains ambiguities, such as those arising from polymorphism, uncorrected read errors, and unresolved repeats; thereby, providing information that has been absent from previous genome assemblies.

- Supports Illumina-based sequencing output, optimized for, but not necessarily limited to, reads of length 100 bases
- Available under open access license from the Broad Institute

SSAKE

www.bcgsc.ca/platform/bioinfo/software/ssake

The Short Sequence Assembly by K-mer search and 3' read Extension (SSAKE) is a genomics application for aggressively assembling millions of short nucleotide sequences by progressively searching for perfect 3'-most k-mers using a DNA prefix tree. SSAKE is designed to help leverage information from short sequence reads by stringently clustering them into contigs for use in characterizing novel sequencing targets.

- SSAKE is written in PERL and runs on Linux
- Available under the terms of the GNU General Public License

Velvet

www.ebi.ac.uk/~zerbino/velvet

Velvet is a *de novo* genomic assembler specifically designed for short read sequencing technologies. Velvet currently takes short read sequences and removes errors, and then produces high quality unique contigs. Using paired-end reads and long read information, when available, Velvet retrieves the repeated areas between contigs.

- Tested on Linux 64-bit, Mac OS X, and Cygwin; designed for 64-bit Linux-compatible environments
- Developed by Daniel Zerbino and Ewan Birney at the European Bioinformatics Institute (EMBL-EBI) near Cambridge in the United Kingdom
- Velvet source code is freely available under the GPL agreement

SNP/Indel Discovery

GATK

www.broadinstitute.org/gatk

The Genome Analysis Toolkit (GATK) analyzes next-generation resequencing data. GATK offers a wide variety of tools with a primary focus on variant discovery and genotyping, as well as strong emphasis on data quality assurance. Its robust architecture, powerful processing engine, and high-performance computing features make it capable of taking on projects of any size.

- Developed at the Broad Institute
- Designed to run on Linux and other POSIX-compatible platforms (i.e., Mac OS X, Cygwin)
- Supports variant discovery, genotyping, and filtering workflows and protocols
- Full-featured GATK is available for non-commercial use, or licensed through Appistry

MACS

liulab.dfci.harvard.edu/MACS

Model-based Analysis of ChIP-Seq (MACS) analyzes short reads. MACS empirically models the length of the sequenced ChIP fragments, which tend to be shorter than sonication or library construction size estimates, and uses it to improve the spatial resolution of predicted binding sites. MACS also uses a dynamic Poisson distribution to effectively capture local biases in the genome sequence, which allows for more sensitive and robust prediction. MACS compares favorably to existing ChIP-Seq peak-finding algorithms. It is a publicly available open source used for ChIP-Seq, with or without control samples.

- Uses a dynamic Poisson distribution to effectively capture local biases in the genome sequence
- Distributed under the terms of Artistic License

PeakSeq

info.gersteinlab.org/PeakSeq

PeakSeq is a program for identifying and ranking peak regions in ChIP-Seq experiments. As input, PeakSeq takes mapped reads from a ChIP-Seq experiment and mapped reads from a control experiment, and outputs a file with peak regions ranked with increasing Q-values.

- Supports multiple input read formats, such as SAM, ELAND, default Bowtie format, and tagAlign
- Supports BAM by piping SAM output from SamTools
- Developed with C/Perl

Transcriptome Analysis

Partek Genomic Suite

www.partek.com

Partek Genomic Suite (GS) empowers biologists to analyze RNA-Seq, ChIP-Seq, Methyl-Seq and DNA-Seq easily by following dedicated workflows such as Data import, Quality Control, Statistical Analysis, Clustering, Gene Ontology, and Pathway analysis, including highly interactive and intuitive visualizations, dedicated genome browser, and data analysis reporting features. At the core of Partek GS, annotations are processed efficiently to enable the correlation of various data formats and save time.

- Peak representation of protein binding sites of interest from aligned reads
- Supports differential expression and alternative splicing based on known mRNA annotation
- Enables the differential expression of known microRNAs (miRNA) between different samples
- Commercial grade package with free 14-day trial license available

RNA-Star

code.google.com/p/rna-star

STAR aligns RNA-seq reads to a reference genome using uncompressed suffix arrays. Accurate alignment of high-throughput RNA-seq data is a challenging and yet unsolved problem due to the non-contiguous transcript structure, relatively short read lengths, and constantly increasing throughput of sequencing technologies. Currently available RNA-seq aligners suffer from high mapping error rates, low mapping speed, read length limitation, and mapping biases.

- Discover non-canonical splices and chimeric (fusion) transcripts
- Capable of mapping full-length RNA sequences
- Implemented as a standalone C++ code
- Free open source software distributed under GPLv3 license

TopHat

tophat.cbcb.umd.edu

TopHat 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 Institute of Genetic Medicine at Johns Hopkins University, the Department of Mathematics and Department of Molecular and Cell Biology at the University of California, Berkeley, and the Department of Stem Cell and Regenerative Biology at Harvard University.

- Finds splice junctions without a reference annotation
- Identifies potential exons, since many RNA-Seq reads will contiguously align to the genome
- Distributed under the terms of Artistic License

Metagenomics

MEGAN

ab.inf.uni-tuebingen.de/software/megan/welcome.html

MEGAN is a tool for studying the taxonomical content of a set of DNA reads, typically collected in a metagenomics project. In a pre-processing step, a sequence comparison of all reads with a suitable database of reference DNA or protein sequences must be performed to produce an input file for the program. MEGAN is suitable for DNA reads (metagenome data), RNA reads (metatranscriptome data), peptide sequences (metaproteomics data), and 16S rRNA data (amplicon sequencing) using a suitable synonyms file that maps SILVA IDs to taxon IDs.

- Interactively explore the data set
- Interactively inspect the assignment of reads to a specific node
- Free academic licensing is available under certain conditions

