

Next Generation Sequencing Hands-on Session

BMI7830

Gulcin Ozer, PhD

gulcin.ozer@osumc.edu

Department of Biomedical Informatics
The Ohio State University

November 3rd, 2015



Wexner Medical Center

UCSC Genome Browser Home

genome.ucsc.edu/index.html Reader +

Apple iCloud Facebook Twitter Wikipedia SSRSearch.a...stTime=true Yahoo! News Popular url

UCSC Genome Bioinformatics

Genomes - Blat - Tables - Gene Sorter - PCR - VisiGene - Session - FAQ - Help

About the UCSC Genome Bioinformatics Site

Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also provides portals to [ENCODE](#) data at UCSC (2003 to 2012) and to the [Neandertal](#) project. Download or purchase the Genome Browser source code, or the Genome Browser in a Box ([GBiB](#)) at our [online store](#).

We encourage you to explore these sequences with our tools. The [Genome Browser](#) zooms and scrolls over chromosomes, showing the work of annotators worldwide. The [Gene Sorter](#) shows expression, homology and other information on groups of genes that can be related in many ways. [Blat](#) quickly maps your sequence to the genome. The [Table Browser](#) provides convenient access to the underlying database. [VisiGene](#) lets you browse through a large collection of *in situ* mouse and frog images to examine expression patterns. [Genome Graphs](#) allows you to upload and display genome-wide data sets.

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the [UC Santa Cruz Genomics Institute](#) and the Center for Biomolecular Science and Engineering ([CBSE](#)) at the University of California Santa Cruz ([UCSC](#)). If you have feedback or questions concerning the tools or data on this website, feel free to contact us on our [public mailing list](#).

The Genome Browser project team relies on public funding to support our work. Donations are welcome – we have many more ideas than our funding supports! If you have ideas, drop a comment in our [suggestion box](#).

[DONATE NOW](#)

News [News Archives ▶](#)

To receive announcements of new genome assembly releases, new software features, updates and training seminars by email, subscribe to the [genome-announce](#) mailing list. Please see our [blog](#) for posts about Genome Browser tools, features, projects and more.

2 February 2015 – Host a Genome Browser Workshop

New timeslots are now available to host a Genome Browser workshop at your institution. Thanks to the funding support of NHGRI, we offer hands-on Genome Browser training onsite at your institution, tailored to your audience's level of expertise.

For more information or to submit a request to host a workshop, please visit [our signup](#).

23 January 2015 – Genome Browser YouTube Channel

We are pleased to announce the release of the UCSC Genome Browser [YouTube channel](#). The channel contains short videos showing how to use the Genome Browser and associated tools to solve selected problems. The videos were produced by Robert Kuhn and Pauline Fujita, with

UCSC Genome Browser [Help](#)

genome.ucsc.edu

UCSC Genome Bioinformatics

Genomes - Blat - Tables - Gene Sorter - PCR - VisiGene - Session - FAQ - Help

Genome Browser

Ebola

Blat

Table Browser

Gene Sorter

In Silico PCR

Genome Graphs

Galaxy

VisiGene

Utilities

Downloads

Release Log

Custom Tracks

Cancer Browser

About the UCSC Genome Bioinformatics Site

Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also provides portals to [ENCODE](#) data at UCSC (2003 to 2012) and to the [Neandertal](#) project. Download or purchase the Genome Browser source code, or the Genome Browser in a Box ([GBiB](#)) at our [online store](#).

We encourage you to explore these sequences with our tools. The [Genome Browser](#) zooms and scrolls over chromosomes, showing the work of annotators worldwide. The [Gene Sorter](#) shows expression, homology and other information on groups of genes that can be related in many ways. [Blat](#) quickly maps your sequence to the genome. The [Table Browser](#) provides convenient access to the underlying database. [VisiGene](#) lets you browse through a large collection of *in situ* mouse and frog images to examine expression patterns. [Genome Graphs](#) allows you to upload and display genome-wide data sets.

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the [UC Santa Cruz Genomics Institute](#) and the Center for Biomolecular Science and Engineering ([CBSE](#)) at the University of California Santa Cruz ([UCSC](#)). If you have feedback or questions concerning the tools or data on this website, feel free to contact us on our [public mailing list](#).

The Genome Browser project team relies on public funding to support our work. Donations are welcome -- we have many more ideas than our funding supports! If you have ideas, drop a comment in our [suggestion box](#). [DONATE NOW](#)

News

[News Archives ▶](#)

To receive announcements of new genome assembly releases, new software features, updates and training seminars by email, subscribe to the [genome-announce](#) mailing list. Please see our [blog](#) for posts about Genome Browser tools, features, projects and more.

2 February 2015 - Host a Genome Browser Workshop

genome.ucsc.edu/cgi-bin/hgUserSuggestion [e to host a Genome Browser workshop at your institution. Thanks to the funding support of](#)

Human (Homo sapiens) [G](#) genome.ucsc.edu/cgi-bin/hgGateway

Genomes Genome Browser Tools Mirrors Downloads My Data Help About Us

Human (*Homo sapiens*) Genome Browser Gateway

The UCSC Genome Browser was created by the [Genome Bioinformatics Group of UC Santa Cruz](#). Software Copyright (c) The Regents of the University of California. All rights reserved.

group	genome	assembly	position	search term
Mammal	Human	Feb. 2009 (GRCh37/hg19)	chr21:33,031,597-33,041,570	enter position, gene symbol or search terms <input type="button" value="submit"/>

[Click here to reset](#) the browser user interface settings to their defaults.
[track search](#) [add custom tracks](#) [track hubs](#) [configure tracks and display](#)

Human Genome Browser – hg19 assembly (sequences)

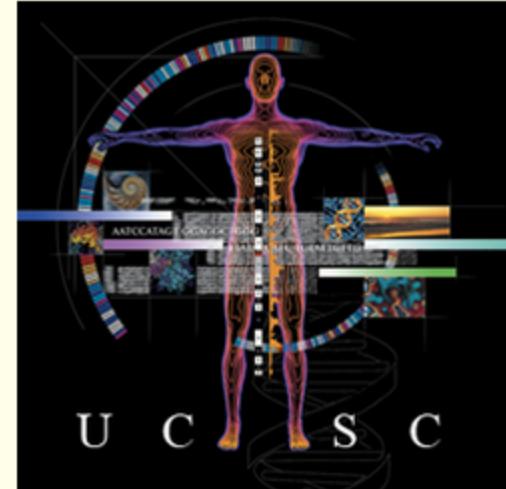
The February 2009 human reference sequence (GRCh37) was produced by the [Genome Reference Consortium](#). For more information about this assembly, see [GRCh37](#) in the NCBI Assembly database.

Sample position queries

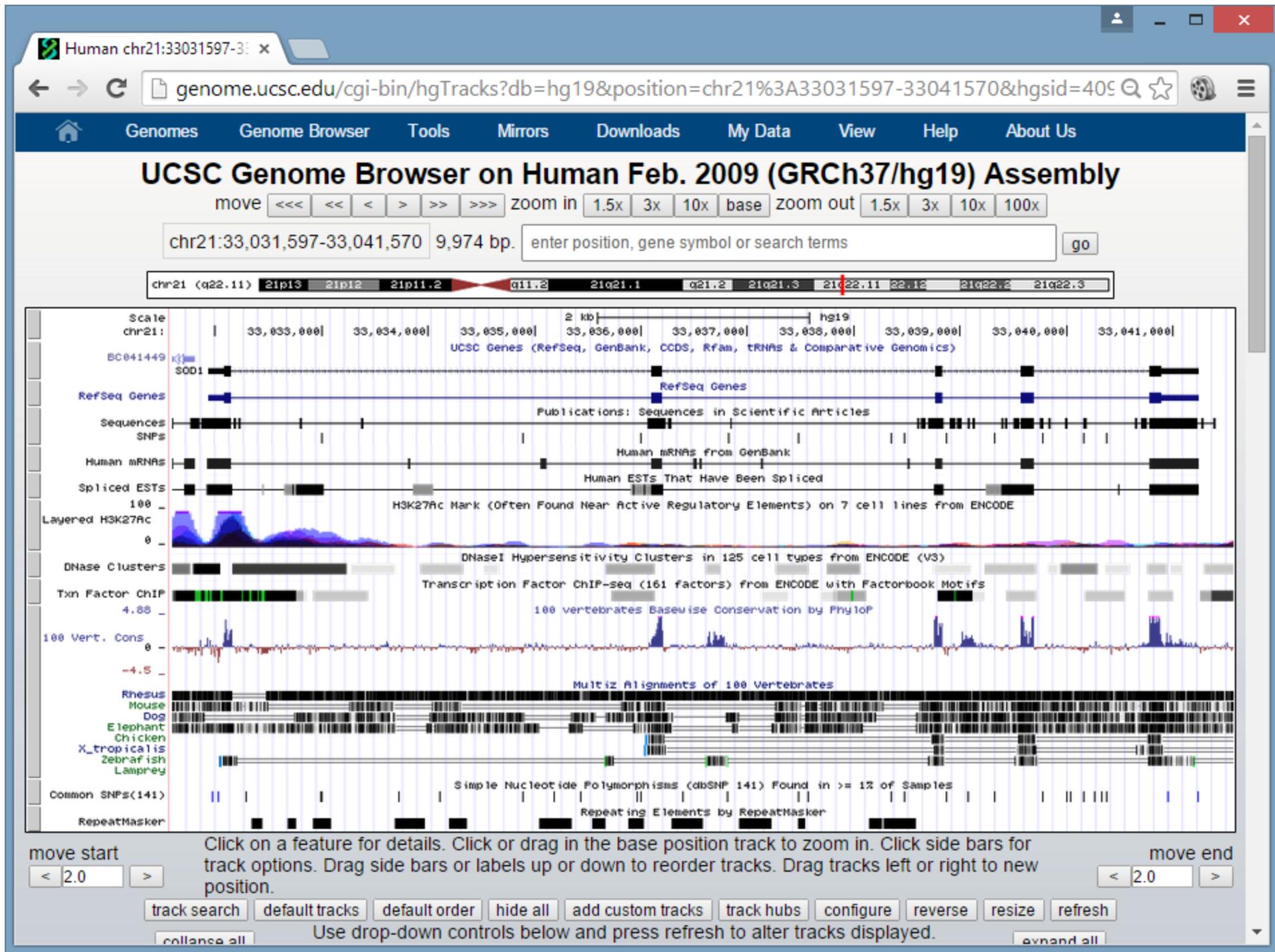
A genome position can be specified by the accession number of a sequenced genomic clone, an mRNA or EST or STS marker, a chromosomal coordinate range, or keywords from the GenBank description of an mRNA. The following list shows examples of valid position queries for the human genome. See the [User's Guide](#) for more information.

Request:	Genome Browser Response:
chr7	Displays all of chromosome 7
chrUn_g1000212	Displays all of the unplaced contig g1000212
20p13	Displays region for band p13 on chr 20
chr3:1-1000000	Displays first million bases of chr 3, counting from p-arm telomere
chr3:1000000+2000	Displays a region of chr3 that spans 2000 bases, starting with position 1000000

DH10001-DH00175 Displays region between genome landmarks, such as the STS markers



(Graphic courtesy of [CBSE](#))



Human chr21:33031597-33041570

genome.ucsc.edu/cgi-bin/hgTracks?db=hg19&position=chr21%3A33031597-33041570&hgsid=409

Tracks with lots of items will automatically be displayed in more compact modes.

Mapping and Sequencing

Genes and Gene Predictions

UCSC Genes	RefSeq Genes	AceView Genes	CCDS	Ensembl Genes	EvoFold
pack ▾	dense ▾	hide ▾	hide ▾	hide ▾	hide ▾
Exoniphy	GENCODE...	Geneid Genes	Genscan Genes	H-Inv 7.0	IKMC Genes
hide ▾	hide ▾	hide ▾	hide ▾	hide ▾	Mapped
lncRNAs...	LRG Transcripts	MGC Genes	N-SCAN	Old UCSC Genes	ORFeome Clones
hide ▾	hide ▾	hide ▾	hide ▾	hide ▾	hide ▾
Other RefSeq	Pfam in UCSC Gene	Retroposed Genes	SGP Genes	SIB Genes	sno/miRNA
hide ▾	hide ▾	hide ▾	hide ▾	hide ▾	hide ▾
TransMap...	tRNA Genes	UCSC Alt Events	UniProt	Vega Genes	Yale Pseudo60
hide ▾	hide ▾	hide ▾	hide ▾	hide ▾	hide ▾

Phenotype and Literature

Publications	ClinVar Variants	Coriell CNVs	COSMIC	DECIPHER	GAD View
dense ▾	hide ▾	hide ▾	hide ▾	hide ▾	hide ▾
GeneReviews	GWAS Catalog	HGMD Variants	ISCA	LOVD Variants	MGI Mouse QTL
hide ▾	hide ▾	hide ▾	hide ▾	hide ▾	hide ▾
OMIM AV SNPs	OMIM Genes	OMIM Pheno Loci	RGD Human QTL	RGD Rat QTL	UniProt Variants
hide ▾	hide ▾	hide ▾	hide ▾	hide ▾	hide ▾

mRNA and EST

Human mRNAs	Spliced ESTs	CGAP SAGE	Gene Bounds	H-Inv	Human ESTs
dense ▾	dense ▾	hide ▾	hide ▾	hide ▾	hide ▾
Human RNA Editing	Other ESTs	Other mRNAs	Poly(A)	PolyA-Seq	SIB Alt-Splicing
hide ▾	hide ▾	hide ▾	hide ▾	hide ▾	hide ▾

Human chr21:33031597-33041570

genome.ucsc.edu/cgi-bin/hgTracks?db=hg19&position=chr21%3A33031597-33041570&hgsid=409974595

Expression

- Affy Exon Array [hide](#)
- Affy GNF1H [hide](#)
- Affy RNA Loc [hide](#)
- Affy U133 [hide](#)
- Affy U133Plus2 [hide](#)
- Affy U95 [hide](#)
- Allen Brain [hide](#)
- Burge RNA-seq [hide](#)
- CSHL Small RNA-seq [hide](#)
- ENC Exon Array... [hide](#)
- ENC ProtGeno... [hide](#)
- GIS RNA PET [hide](#)
- GNF Atlas 2 [hide](#)
- Illumina WG-6 [hide](#)
- PeptideAtlas... [hide](#)
- gPCR Primers [hide](#)
- RIKEN CAGE Loc [hide](#)
- Sestan Brain [hide](#)

Regulation

- ENCODE Regulation... [show](#) [hide](#)
- ENC RNA Binding... [hide](#)
- CD34 DnaseI [hide](#)
- CpG Islands... [hide](#)
- ENC Chromatin... [hide](#)
- ENC DNA Methyl... [hide](#)
- DNase/FAIRE... [hide](#)
- ENKE Histone... [hide](#)
- ENKE Nucleosome [hide](#)
- ENKE TF Binding... [hide](#)
- FSU Repli-chip [hide](#)
- Genome Segments [hide](#)
- NKI Nuc Lamina... [hide](#)
- ORegAnno [hide](#)
- Stanf Nucleosome [hide](#)
- SUNY SwitchGear [hide](#)
- SwitchGear TSS [hide](#)
- TFBS Conserved [hide](#)
- TS miRNA sites [hide](#)
- UCSF Brain Methyl [hide](#)
- UMMS Brain Hist [hide](#)
- UW Repli-seq [hide](#)
- Vista Enhancers [hide](#)

Comparative Genomics

- Conservation [full](#) [hide](#)
- Cons 46-Way [hide](#)
- Cons Indels MmCf [hide](#)
- Evo Cpg [hide](#)
- GERP [hide](#)
- phastBias gBGC [hide](#)
- Primate Chain/Net [hide](#)
- Placental Chain/Net [hide](#)
- Vertebrate Chain/Net [hide](#)

Indel-based Conservation for Human hg19, Mouse mm8 and Dog canFam2

Neandertal Assembly and Analysis

genome.ucsc.edu/cgi-bin/hgTrackUi?hgsid=409974595

UCSC Genome Browser [Help](#)

genome.ucsc.edu

UCSC Genome Bioinformatics

Genomes - Blat - Tables - Gene Sorter - PCR - VisiGene - Session - FAQ - Help

About the UCSC Genome Bioinformatics Site

Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also provides portals to [ENCODE](#) data at UCSC (2003 to 2012) and to the [Neandertal](#) project. Download or purchase the Genome Browser source code, or the Genome Browser in a Box ([GBiB](#)) at our [online store](#).

We encourage you to explore these sequences with our tools. The [Genome Browser](#) zooms and scrolls over chromosomes, showing the work of annotators worldwide. The [Gene Sorter](#) shows expression, homology and other information on groups of genes that can be related in many ways. [Blat](#) quickly maps your sequence to the genome. The [Table Browser](#) provides convenient access to the underlying database. [VisiGene](#) lets you browse through a large collection of *in situ* mouse and frog images to examine expression patterns. [Genome Graphs](#) allows you to upload and display genome-wide data sets.

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the [UC Santa Cruz Genomics Institute](#) and the Center for Biomolecular Science and Engineering ([CBSE](#)) at the University of California Santa Cruz ([UCSC](#)). If you have feedback or questions concerning the tools or data on this website, feel free to contact us on our [public mailing list](#).

The Genome Browser project team relies on public funding to support our work. Donations are welcome -- we have many more ideas than our funding supports! If you have ideas, drop a comment in our [suggestion box](#). [DONATE NOW](#)

News [News Archives ▶](#)

To receive announcements of new genome assembly releases, new software features, updates and training seminars by email, subscribe to the [genome-announce](#) mailing list. Please see our [blog](#) for posts about Genome Browser tools, features, projects and more.

2 February 2015 - Host a Genome Browser Workshop

genome.ucsc.edu/cgi-bin/hgUserSuggestion [e to host a Genome Browser workshop at your institution. Thanks to the funding support of](#)

Downloads

Genome Browser
Ebola
Blat
Table Browser
Gene Sorter
In Silico PCR
Genome Graphs
Galaxy
VisiGene
Utilities
Release Log
Custom Tracks
Cancer Browser

UCSC Genome Browser: D x

hgdownload.soe.ucsc.edu/downloads.html

UCSC Genome Bioinformatics

Home - Genomes - Blat - Tables - Gene Sorter - PCR - FAQ - Help

Sequence and Annotation Downloads

This page contains links to sequence and annotation data downloads for the genome assemblies featured in the UCSC Genome Browser. Table downloads are also available via the Genome Browser [FTP server](#). For quick access to the most recent assembly of each genome, see the [current genomes](#) directory. This directory may be useful to individuals with automated scripts that must always reference the most recent assembly.

To view the current descriptions and formats of the tables in the annotation database, use the "describe table schema" button in the Table Browser. The [Description of the annotation database](#) page (no longer maintained) also provides descriptions of selected tables in the database.

All tables in the Genome Browser are freely usable for any purpose except as indicated in the README.txt files in the download directories. To view restrictions specific to a particular data set, click on the corresponding download link and review the README text. These data were contributed by many researchers, as listed on the Genome Browser [credits](#) page. Please acknowledge the contributor(s) of the data you use.

VERTEBRATES - Complete annotation sets

Human	Hedgehog	Platypus
Alpaca	Horse	Rabbit
American alligator	Kangaroo rat	Rat
Armadillo	Lamprey	Rhesus
Atlantic cod	Lizard	Rock hyrax
Baboon	Manatee	Sheep
Budgerigar	Marmoset	Shrew
Bushbaby	Medaka	Sloth
Cat	Medium ground finch	Squirrel
Chicken	Megabat	Squirrel monkey
Chimpanzee	Microbat	Stickleback

UCSC Genome Browser: D x

hgdownload.soe.ucsc.edu/downloads.html#human

Human Genome

Dec. 2013 (hg38, GRCh38)

- [Full data set](#)
- [Data set by chromosome](#)
- [Annotation database](#)
- [Protein database for hg38](#)
- [SNP141-masked FASTA files](#)
- [LiftOver files](#)
- Pairwise Alignments
 - [Human/Chimp \(panTro4\)](#)
 - [Human/Rhesus \(rheMac3\)](#)
 - [Human/Mouse \(mm10\)](#)
 - [Human/Rat \(rn5\)](#)
 - [Human/Dog \(canFam3\)](#)
 - [Human/Opossum \(monDom5\)](#)
- Multiple Alignments
 - [Multiple alignments of 7 vertebrate genomes with Human](#)
 - [Conservation scores for alignments of 7 vertebrate genomes with Human](#)
 - [Basewise conservation scores \(phyloP\) of 7 vertebrate genomes with Human](#)
 - [FASTA alignments of 7 vertebrate genomes with Human for CDS regions](#)

Feb. 2009 (hg19, GRCh37)

Index of /goldenPath/hg38

hgdownload.soe.ucsc.edu/goldenPath/hg38/bigZips/

```
gunzip <file>.fa.gz
```

GenBank Data Usage

The GenBank database is designed to provide and encourage access within the scientific community to the most up to date and comprehensive DNA sequence information. Therefore, NCBI places no restrictions on the use or distribution of the GenBank data. However, some submitters may claim patent, copyright, or other intellectual property rights in all or a portion of the data they have submitted. NCBI is not in a position to assess the validity of such claims, and therefore cannot provide comment or unrestricted permission concerning the use, copying, or distribution of the information contained in GenBank.

<u>Name</u>	<u>Last modified</u>	<u>Size</u>	<u>Description</u>
Parent Directory		-	
analysisSet/	14-Apr-2014 09:51	-	
est.fa.gz	02-Feb-2015 15:01	1.5G	
est.fa.gz.md5	02-Feb-2015 15:01	44	
hg38.2bit	09-Jan-2014 14:14	797M	
hg38.agp.gz	15-Jan-2014 20:55	842K	
hg38.chromFa.tar.gz	23-Jan-2014 17:18	938M	
hg38.chromFaMasked.tar.gz	23-Jan-2014 17:10	487M	
hg38.fa.align.gz	08-Jan-2014 23:43	2.4G	
hg38.fa.gz	15-Jan-2014 21:14	938M	
hg38.fa.masked.gz	15-Jan-2014 21:24	487M	
hg38.fa.out.gz	15-Jan-2014 20:56	172M	
hg38.trf.bed.gz	15-Jan-2014 20:56	7.9M	
md5sum.txt	03-Mar-2014 10:31	451	
mrna.fa.gz	02-Feb-2015 14:34	173M	
mrna.fa.gz.md5	02-Feb-2015 14:34	45	
refMrna.fa.gz	02-Feb-2015 15:02	49M	
refMrna.fa.gz.md5	02-Feb-2015 15:02	48	
xenoMrna.fa.gz	02-Feb-2015 14:44	4.1G	
xenoMrna.fa.gz.md5	02-Feb-2015 14:44	49	
xenoRefMrna.fa.gz	02-Feb-2015 15:01	218M	
xenoRefMrna.fa.gz.md5	02-Feb-2015 15:01	52	

hgdownload.soe.ucsc.edu/goldenPath/hg38/bigZips/

UCSC Genome Browser [Help](#)

genome.ucsc.edu

UCSC Genome Bioinformatics

Genomes - Blat - Tables - Gene Sorter - PCR - VisiGene - Session - FAQ - Help

Genome Browser

Ebola

Blat

Table Browser

Gene Sorter

In Silico PCR

Genome Graphs

Galaxy

VisiGene

Utilities

Downloads

Release Log

Custom Tracks

Cancer Browser

About the UCSC Genome Bioinformatics Site

Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also provides portals to [ENCODE](#) data at UCSC (2003 to 2012) and to the [Neandertal](#) project. Download or purchase the Genome Browser source code, or the Genome Browser in a Box ([GBiB](#)) at our [online store](#).

We encourage you to explore these sequences with our tools. The [Genome Browser](#) zooms and scrolls over chromosomes, showing the work of annotators worldwide. The [Gene Sorter](#) shows expression, homology and other information on groups of genes that can be related in many ways. [Blat](#) quickly maps your sequence to the genome. The [Table Browser](#) provides convenient access to the underlying database. [VisiGene](#) lets you browse through a large collection of *in situ* mouse and frog images to examine expression patterns. [Genome Graphs](#) allows you to upload and display genome-wide data sets.

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the [UC Santa Cruz Genomics Institute](#) and the Center for Biomolecular Science and Engineering ([CBSE](#)) at the University of California Santa Cruz ([UCSC](#)). If you have feedback or questions concerning the tools or data on this website, feel free to contact us on our [public mailing list](#).

The Genome Browser project team relies on public funding to support our work. Donations are welcome -- we have many more ideas than our funding supports! If you have ideas, drop a comment in our [suggestion box](#).

[DONATE NOW](#)

News

[News Archives ►](#)

To receive announcements of new genome assembly releases, new software features, updates and training seminars by email, subscribe to the [genome-announce](#) mailing list. Please see our [blog](#) for posts about Genome Browser tools, features, projects and more.

2 February 2015 - Host a Genome Browser Workshop

genome.ucsc.edu/cgi-bin/hgUserSuggestion [e to host a Genome Browser workshop at your institution. Thanks to the funding support of](#)

Table Browser genome.ucsc.edu/cgi-bin/hgTables

Genomes Genome Browser Tools Mirrors Downloads My Data Help About Us

Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). To examine the biological function of your set through annotation enrichments, send the data to [GREAT](#). Send data to [GenomeSpace](#) for use with diverse computational tools. Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data. All tables can be downloaded in their entirety from the [Sequence and Annotation Downloads](#) page.

clade: Mammal **genome:** Human **assembly:** Feb. 2009 (GRCh37/hg19)

group: Genes and Gene Predictions **track:** UCSC Genes

table: knownGene **describe:**

region: genome ENCODE Pilot regions

identifiers (names/accessions): [paste list](#) [upload](#)

filter: [create](#)

intersection: [create](#)

correlation: [create](#)

output format: all fields from selected table

output file: (leave empty to output to browser)

file type returned: plain text gzip compressed

[get output](#) [summary/statistics](#)

To reset **all** user cart settings (including custom tracks)

Using the Table Browser

This section provides brief line-by-line descriptions of the controls in the Table Browser. For more information on using this program, see the [Table Browser User's Guide](#).

- **clade:** Specifies which clade the organism is in
- **genome:** Specifies which organism data to use

UCSC Genes
RefSeq Genes
GENCODE Genes V19
GENCODE Genes V17
GENCODE Genes V14
GENCODE Genes V7
TransMap UCSC
TransMap RefGene
TransMap mRNA
TransMap ESTs
AceView Genes
CCDS
Ensembl Genes
EvoFold
Exoniphy

41570 [lookup](#) [define regions](#)

Galaxy GREAT GenomeSpace

Browser)

Common Data Formats

- Fastq (reads with quality information)
- SAM/BAM (detailed alignment information)
- BED
- WIG
- GTF
- VCF



| Wexner Medical Center

Sequence Reads with Quality Information

FASTQ is the most commonly used format



Wexner Medical Center

Sequence Alignment/Map (SAM) Format

- Standard format for short read alignment

- SAM is plain text
 - BAM is binary and indexed
 - Details of the format

<http://samtools.github.io/hts-specs/SAMv1.pdf>

- Two sections
 - Header
 - Alignment



The SAM

- Each type
- File formats used

record
ble
ams

Tag	Description
<code>@HD</code>	The header line. The first line if present.
<code>VN*</code>	Format version. Accepted format: <code>/^ [0-9]+ \. [0-9]+ \$/</code> .
<code>SO</code>	Sorting order of alignments. Valid values: <code>unknown</code> (default), <code>unsorted</code> , <code>queryname</code> and <code>coordinate</code> . For coordinate sort, the major sort key is the RNAME field, with order defined by the order of <code>@SQ</code> lines in the header. The minor sort key is the POS field. For alignments with equal RNAME and POS, order is arbitrary. All alignments with '*' in RNAME field follow alignments with some other value but otherwise are in arbitrary order.
<code>@SQ</code>	Reference sequence dictionary. The order of <code>@SQ</code> lines defines the alignment sorting order.
<code>SN*</code>	Reference sequence name. Each <code>@SQ</code> line must have a unique SN tag. The value of this field is used in the alignment records in RNAME and PNEXT fields. Regular expression: <code>[!-]+-<>-[!-]*</code>
<code>LN*</code>	Reference sequence length. Range: <code>[1,2³¹-1]</code>
<code>AS</code>	Genome assembly identifier.
<code>M5</code>	MD5 checksum of the sequence in the uppercase, excluding spaces but including pads (as '*'s).
<code>SP</code>	Species.
<code>UR</code>	URI of the sequence. This value may start with one of the standard protocols, e.g http: or ftp:. If it does not start with one of these protocols, it is assumed to be a file-system path.
<code>@RG</code>	Read group. Unordered multiple <code>@RG</code> lines are allowed.
<code>ID*</code>	Read group identifier. Each <code>@RG</code> line must have a unique ID. The value of ID is used in the RG tags of alignment records. Must be unique among all read groups in header section. Read group IDs may be modified when merging SAM files in order to handle collisions.
<code>CN</code>	Name of sequencing center producing the read.
<code>DS</code>	Description.
<code>DT</code>	Date the run was produced (ISO8601 date or date/time).
<code>FO</code>	Flow order. The array of nucleotide bases that correspond to the nucleotides used for each flow of each read. Multi-base flows are encoded in IUPAC format, and non-nucleotide flows by various other characters. Format: <code>/* [ACMGRSVTWYHKDBN]+\/</code>
<code>KS</code>	The array of nucleotide bases that correspond to the key sequence of each read.
<code>LB</code>	Library.
<code>PG</code>	Programs used for processing the read group.
<code>PI</code>	Predicted median insert size.
<code>PL</code>	Platform/technology used to produce the reads. Valid values: CAPILLARY, LS454, ILLUMINA, SOLID, HELICOS, IONTORRENT and PACBIO.
<code>PU</code>	Platform unit (e.g. flowcell-barcode.lane for Illumina or slide for SOLiD). Unique identifier.
<code>SM</code>	Sample. Use pool name where a pool is being sequenced.
<code>@PG</code>	Program.
<code>ID*</code>	Program record identifier. Each <code>@PG</code> line must have a unique ID. The value of ID is used in the alignment PG tag and PP tags of other <code>@PG</code> lines. PG IDs may be modified when merging SAM files in order to handle collisions.
<code>PN</code>	Program name
<code>CL</code>	Command line
<code>PP</code>	Previous <code>@PG-ID</code> . Must match another <code>@PG</code> header's ID tag. <code>@PG</code> records may be chained using PP tag, with the last record in the chain having no PP tag. This chain defines the order of programs that have been applied to the alignment. PP values may be modified when merging SAM files in order to handle collisions of PG IDs. The first PG record in a chain (i.e. the one referred to

The SAM/BAM Format - Alignment

Each line has 11 mandatory fields

Col	Field	Type	Regexp/Range	Brief description
1	QNAME	String	[!-?A-~]{1,255}	Query template NAME
2	FLAG	Int	[0,2 ¹⁶ -1]	bitwise FLAG
3	RNAME	String	* [!-()+-<>-~] [!-~]*	Reference sequence NAME
4	POS	Int	[0,2 ³¹ -1]	1-based leftmost mapping POSition
5	MAPQ	Int	[0,2 ⁸ -1]	MAPping Quality
6	CIGAR	String	* ([0-9]+[MIDNSHPX=])+	CIGAR string
7	RNEXT	String	* = [!-()+-<>-~] [!-~]*	Ref. name of the mate/next read
8	PNEXT	Int	[0,2 ³¹ -1]	Position of the mate/next read
9	TLEN	Int	[-2 ³¹ +1,2 ³¹ -1]	observed Template LENgth
10	SEQ	String	* [A-Za-z.=.]+	segment SEQuence
11	QUAL	String	[!-~]+	ASCII of Phred-scaled base QUALity+33



The SAM/BAM Format - Alignment

CIGAR string - 9M32N8M

REF.: AGCTAGCATCGTGTAAACCGGTCTAGCAACGCTAGTCAGCTAGTCAGACTAGTCGATCGATGTG

READ: GTGTAACCC.....TCAGAATA



Wexner Medical Center

BED file format

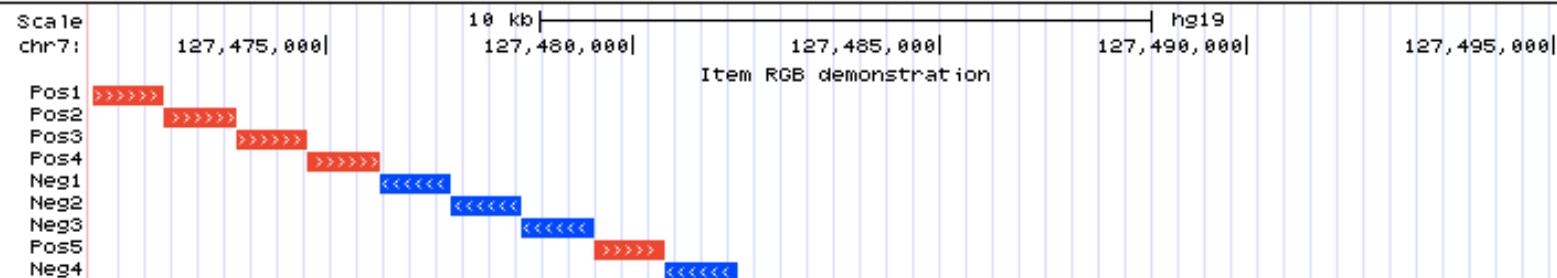
UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly

move <<< << < > >> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x 100x

chr7:127,471,196-127,495,720 24,525 bp.

enter position, gene symbol or search terms

go



```
browser position chr7:127471196-127495720
browser hide all
track name="ItemRGBDemo" description="Item RGB demonstration" visibility=2
itemRgb="On"
chr7    127471196    127472363    Pos1    0    +    127471196    127472363    255,0,0
chr7    127472363    127473530    Pos2    0    +    127472363    127473530    255,0,0
chr7    127473530    127474697    Pos3    0    +    127473530    127474697    255,0,0
chr7    127474697    127475864    Pos4    0    +    127474697    127475864    255,0,0
chr7    127475864    127477031    Neg1    0    -    127475864    127477031    0,0,255
chr7    127477031    127478198    Neg2    0    -    127477031    127478198    0,0,255
chr7    127478198    127479365    Neg3    0    -    127478198    127479365    0,0,255
chr7    127479365    127480532    Pos5    0    +    127479365    127480532    255,0,0
chr7    127480532    127481699    Neg4    0    -    127480532    127481699    0,0,255
```

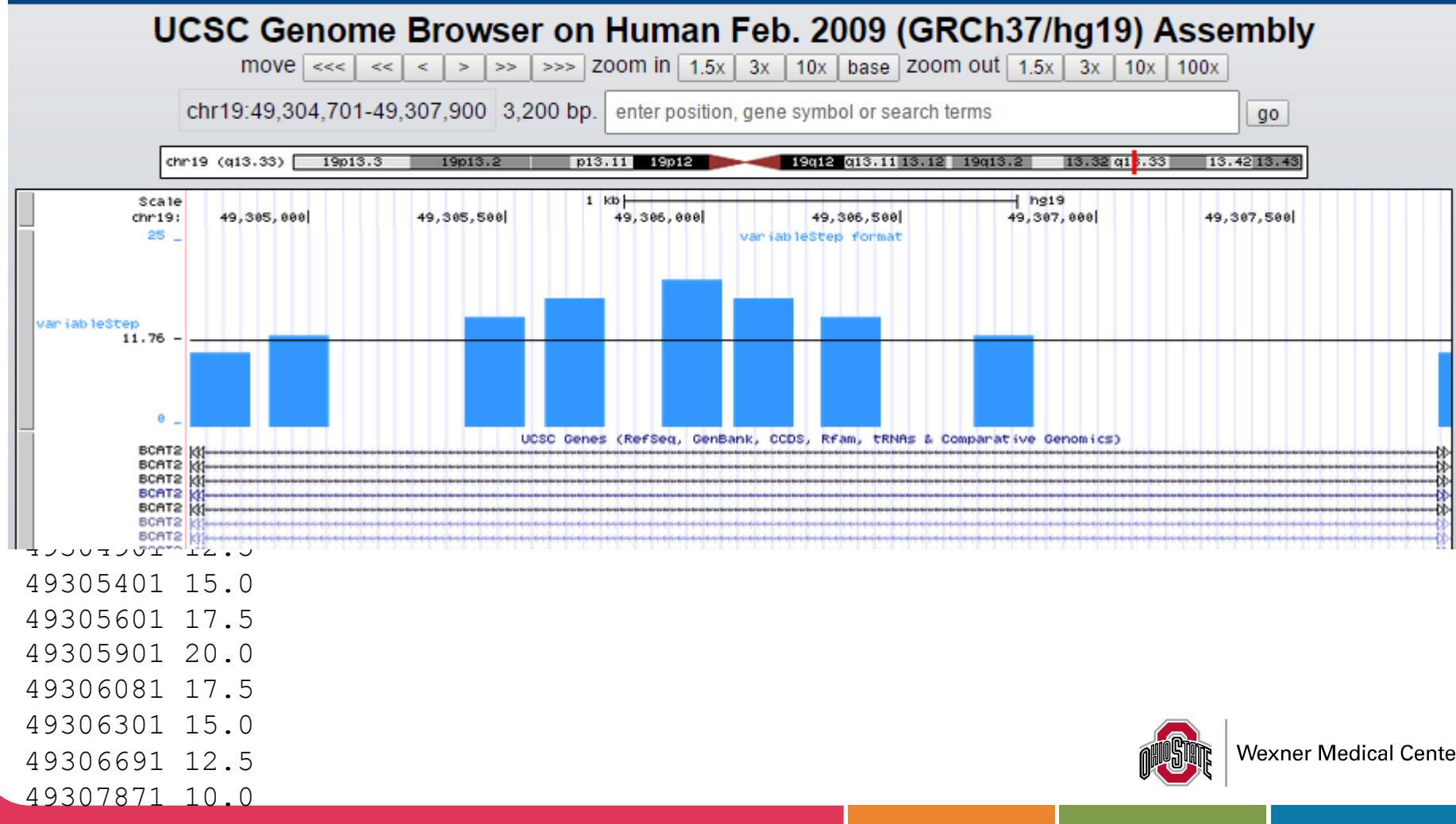


Wexner Medical Center

WIG File Format

- Wiggle format (WIG) allows the display of continuous-valued data in a track format.

The **bigWig** format is for display of dense, continuous data that will be



VCF File Format

- Variant Call Format (VCF) is a flexible and extendable format for variation data

##fileformat=VCFv4.1
##fileDate=20090805
##tcgaversion=1.1
##vcfProcessLog=<InputVCF=<file1.vcf>, InputVCFVer=<1.0>, InputVCFParam=<a1,b>, InputVCFgeneAnno=<anno1.gaf>>
##reference=<ftp://ftp.ncbi.nih.gov/genbank/genomes/Eukaryotes/vertebrates_mammals/Homo_sapiens/GRCh37/special_requests/GRCh37-lite.i>
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">

##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">

##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">

##SAMPLE=<ID=NORMAL,Individual=TCGA-01-1000,File=TCGA-01-1000-1.bam,Platform=Illumina,Source=dbGAP,Accession=1234>
##SAMPLE=<ID=TUMOR,Individual=TCGA-01-1000,File=TCGA-01-1000-2.bam,Platform=Illumina,Source=dbGAP,Accession=4567>
##PEDIGREE=<Name_0=TUMOR,Name_1=NORMAL>

HEADER

BODY

Fixed fields

CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;DB
20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T
20	1234567	microsat1	GTC	G,GTCTC	50	PASS	NS=3;DP=9;AA=G

Optional: FORMAT field specifying data type + Per-sample genotype data

FORMAT	NORMAL	TUMOR
GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51
GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3
GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2
GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51
GT:GQ:DP	0/1:35:4	0/2:17:2

[https://wiki.nci.nih.gov/display/TCGA/TCGA+Variant+Call+Format+\(VCF\)+1.1+Specification](https://wiki.nci.nih.gov/display/TCGA/TCGA+Variant+Call+Format+(VCF)+1.1+Specification)



Wexner Medical Center

GTF File Format

- Gene Transfer Format
- Structure

```
<seqname> <source> <feature> <start> <end> <score>  
<strand> <frame> [attributes] [comments]
```

- Simple example with 3 translated exons. Order of rows is not important.

```
AB000381 Twinscan CDS      380    401    .    +    0    gene_id "001"; transcript_id "001.1";  
AB000381 Twinscan CDS      501    650    .    +    2    gene_id "001"; transcript_id "001.1";  
AB000381 Twinscan CDS      700    707    .    +    2    gene_id "001"; transcript_id "001.1";  
AB000381 Twinscan start_codon 380    382    .    +    0    gene_id "001"; transcript_id "001.1";  
AB000381 Twinscan stop_codon 708    710    .    +    0    gene_id "001"; transcript_id "001.1";
```



Wexner Medical Center

UCSC Genome Bioinformatics: FAQ

genome.ucsc.edu/FAQ/FAQformat.html

Reader

Apple iCloud Facebook Twitter Wikipedia SSRSearch.a...stTime=true Yahoo! News Popular url

UCSC Genome Bioinformatics

Home - Genomes - Blat - Tables - Gene Sorter - PCR - VisiGene - Session - Help

Frequently Asked Questions: Data File Formats

General formats:

- [Axt format](#)
- [BAM format](#)
- [BED format](#)
- [BED detail format](#)
- [bedGraph format](#)
- [bigBed format](#)
- [bigWig format](#)
- [Chain format](#)
- [GenePred table format](#)
- [GFF format](#)
- [GTF format](#)
- [HAL format](#)
- [MAF format](#)
- [Microarray format](#)
- [Net format](#)
- [Personal Genome SNP format](#)
- [PSL format](#)
- [VCF format](#)
- [WIG format](#)

ENCODE-specific formats:

- [ENCODE broadPeak format](#)
- [ENCODE gappedPeak format](#)
- [ENCODE narrowPeak format](#)
- [ENCODE pairedTagAlign format](#)
- [ENCODE peptideMapping format](#)

<http://genome.ucsc.edu/FAQ/FAQformat.html>





Upload a fastq file

Quality control (FastQC)

Mapping (BWA)

Metrics

Visualization (Trackster)

Visualization (IGV)

Variant Calling (VarScan2)

Variant Annotation (Annovar)

Variant Annotation (VEP)

<https://osu.box.com/exome-seq-example>

Create a Galaxy Account

<https://usegalaxy.org/>

Galaxy

Tools

search tools

Contacts

Send Data

Lift-Over

Text Manipulation

Convert Formats

Filter and Sort

Join, Subtract and Group

NGS: QC and manipulation

NGS: Mapping

NGS: RNA-seq

NGS: SAMtools

NGS: BAM Tools

NGS: Picard

NGS: VCF Manipulation

Extract Features

Fetch Sequences

Fetch Alignments

Get Genomic Scores

Operate on Genomic Intervals

Statistics

Graph/Display Data

Phenotype Association

snpEff

BEDTools

Genome Diversity

EMBOSS

Galaxy start here or consult our [help resources](#).

The central panel features a large blue banner with the text "Try Galaxy on the Cloud" and a smaller text below it stating "Now you can have a personal Galaxy within the infinite Universe". Below the banner is a series of small circular icons. To the right is a "History" window pane showing a single history named "workshop-exome-seq" which contains an unnamed history entry with 0 bytes and no matching datasets found.

History window

- ✓ All analysis steps are saved
- ✓ Data is not overwritten
- ✓ Output can be visualized and downloaded
- ✓ Can create a workflow to repeat an analysis

- ✓ Public instance
- ✓ No need of programming/cluster computing experience
- ✓ Integrates many bioinformatics tools within one interface
- ✓ Keeps track of all the steps performed in an analysis

Galaxy

Analyze Data Workflow Shared Data Visualization Cloud Help User

Tools

search tools

Get Data

Upload File from your computer

UCSC Main table browser

UCSC Archaea table browser

EBI SRA ENA SRA

BioMart Central server

GrameneMart Central server

Flymine server

modENCODE fly server

modENCODE modMine server

MouseMine server

Ratmine server

YeastMine server

modENCODE worm server

WormBase server

ZebrafishMine server

EuPathDB server

GenomeSpace import from browser

Download data directly from web or upload files from your disk

You added 1 file(s) to the queue. Add more files or click 'Start' to proceed.

Name	Size	Type	Genome	Settings	Status
100326_FC6107FAAXX-chr22_35-38M.fastq.gz	1.7 MB	fastqsanger	Human Feb. 2009 (...)		

Type (set all): Auto-detect Genome (set all): ----- Additional Species ...

Choose local file Choose FTP file Paste/Fetch data Start Pause Reset Close

Example dataset: <https://osu.app.box.com/exome-seq-example>

The screenshot shows the Galaxy web interface with the following details:

- Left Sidebar (Tools):** A list of tools categorized under "NGS: QC and manipulation". The "FastQC Read Quality reports" tool is highlighted with a red circle.
- Tool Panel (Top):** The "FastQC Read Quality reports (Galaxy Tool Version 0.63)" panel is active. It includes tabs for "Short read data from your current history" and "Contaminant list". The "Short read data from your current history" tab has a file input field containing "1: 100326_FC6107FAAXX-chr22_35-38M.fastq.gz", which is also circled in red.
- Tool Panel (Bottom):** Includes sections for "Submodule and Limit specifying file" and "Execute" button.
- Right Panel (History):** Shows a history named "workshop-exome-seq" with an "Unnamed history" entry containing 0 bytes and a message "No matching datasets found".



Galaxy

Analyze Data Workflow Shared Data Visualization Cloud Help User

Using 0%

Tools

[Send Data](#)

[Lift-Over](#)

[Text Manipulation](#)

[Convert Formats](#)

[Filter and Sort](#)

[Join, Subtract and Group](#)

[NGS: QC and manipulation](#)

[NGS: Mapping](#)

Bowtie2 - map reads against reference genome

BWA - map short reads (< 100 bp) against reference genome

BWA-MEM - map medium and long reads (> 100 bp) against reference genome

Parse blast XML output

Megablast compare short reads against htgs, nt, and wgs databases

Map with BWA for Illumina

Map with Bowtie for Illumina

Lastz map short reads against reference sequence

[NGS: RNA-seq](#)

[NGS: SAMtools](#)

[NGS: BAM Tools](#)

[NGS: Picard](#)

[NGS: VCF Manipulation](#)

BWA - map short reads (< 100 bp) against reference genome (Galaxy Tool Version 0.1)

Load reference genome from Local cache

Using reference genome Human (Homo sapiens) (b37): hg19

Select genome from the list

Select input type Single fastq

Select between fastq and bam datasets and between paired and single end data

Select fastq dataset 1: 100326_FC6107FAAXX-chr22_35-38M.fastq

Specify dataset with single reads

Set advanced single end options? Do not set

Provides additional controls

Set readgroups information? Do not set

Specifying readgroup information can greatly simplify your downstream analyses by allowing combining multiple datasets. See help below for more details

Select analysis mode 1.Simple Illumina mode

Job Resource Parameters Use default job resource parameters

Execute

History

search datasets

Unnamed history 3 shown 9.4 MB

3: FastQC on data 1: Raw Data

2: FastQC on data 1: Web page

1: 100326_FC6107FAAXX-chr22_35-38M.fastq

Galaxy

- Analyze Data
- Workflow
- Shared Data
- Visualization
- Cloud
- Help
- User

Using 0%

1 job has been successfully added to the queue – resulting in the following datasets:

4: BWA on data 1 (mapped reads in BAM format)

You can check the status of queued jobs and view the resulting data by refreshing the History pane. When the job has been run the status will change from 'running' to 'finished' if completed successfully or 'error' if problems were encountered.

Tools

- [Send Data](#)
- [Lift-Over](#)
- [Text Manipulation](#)
- [Convert Formats](#)
- [Filter and Sort](#)
- [Join, Subtract and Group](#)
- [NGS: QC and manipulation](#)
- [NGS: Mapping](#)
 - [Bowtie2 – map reads against reference genome](#)
 - [BWA – map short reads \(< 100 bp\) against reference genome](#)
 - [BWA-MEM – map medium and long reads \(> 100 bp\) against reference genome](#)
 - [Parse blast XML output](#)
 - [Megablast compare short reads against htgs, nt, and wgs databases](#)
 - [Map with BWA for Illumina](#)
 - [Map with Bowtie for Illumina](#)
 - [Lastz map short reads against reference sequence](#)
- [NGS: RNA-seq](#)
- [NGS: SAMtools](#)
- [NGS: BAM Tools](#)

History

search datasets

Unnamed history

4 shown

11.4 MB

4: BWA on data 1 (mapped reads in BAM format)

1.9 MB

format: **bam**, database: **hg19**

```
[bwa_aln] 17bp reads: max_diff = 2
[bwa_aln] 38bp reads: max_diff = 3
[bwa_aln] 64bp reads: max_diff = 4
[bwa_aln] 93bp reads: max_diff = 5
[bwa_aln] 124bp reads: max_diff = 6
[bwa_aln] 157bp reads: max_diff = 7
[bwa_aln] 190bp reads: max_diff = 8
[bwa_a]
```

display at UCSC main

display at Ensembl Current

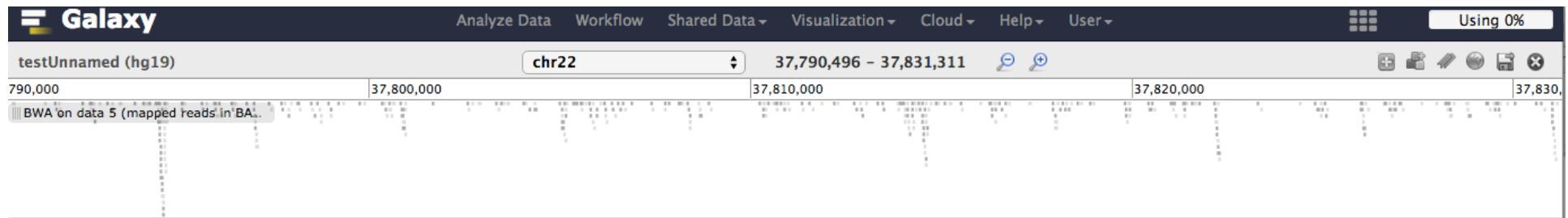
display with IGV local Human hg19

display in IGD View

Binary bam alignments file

3: FastQC on data 1: Raw Data

2: FastQC on data 1: Web page



Tools

[EMBOSS](#)[Regional Variation](#)[FASTA manipulation](#)[Evolution](#)[Multiple Alignments](#)[Metagenomic analyses](#)[Motif Tools](#)

NGS TOOLBOX BETA

[NGS: Peak Calling](#)[NGS: Variant Analysis](#)[NGS: GATK Tools \(beta\)](#)[NGS: Picard \(beta\)](#)

[FASTQ to BAM](#) creates an unaligned BAM file

[SAM to FASTQ](#) creates a FASTQ file

[BAM Index Statistics](#)[SAM/BAM Alignment Summary Metrics](#)[SAM/BAM GC Bias Metrics](#)[Estimate Library Complexity](#)

[Insertion size metrics for PAIRED data](#)

[SAM/BAM Hybrid Selection Metrics for targeted resequencing data](#)

[Add or Replace Groups](#)[Reorder SAM/BAM](#)[Replace SAM/BAM Header](#) **SAM/BAM Alignment Summary Metrics (Galaxy Tool Version 1.56.0)**

▼ Options

SAM/BAM dataset to generate statistics for 6: BWA on data 1 (mapped reads in BAM format)
If empty, upload or import a SAM/BAM dataset.**Title for the output file**

Use this remind you what the job was for.

Select Reference Genome**Check the assigned reference genome**

Galaxy thinks that the reads in your dataset were aligned against this reference. If this is not correct, use the 'Select a build-in reference genome' option of the 'Select Reference Genome' dropdown to select appropriate Reference.

Assume the input file is already sorted Yes No**Input file contains Bisulphite sequenced reads** Yes No**Adapter sequences**
One per line if multiple**Larger paired end reads and inter-chromosomal pairs considered chimeric**

Tools

[EMBOSS](#)[Regional Variation](#)[FASTA manipulation](#)[Evolution](#)[Multiple Alignments](#)[Metagenomic analyses](#)[Motif Tools](#)[NGS TOOLBOX BETA](#)[NGS: Peak Calling](#)[NGS: Variant Analysis](#)[NGS: GATK Tools \(beta\)](#)[NGS: Picard \(beta\)](#)

[FASTQ to BAM](#) creates an unaligned BAM file

[SAM to FASTQ](#) creates a FASTQ file

[BAM Index Statistics](#)

[SAM/BAM Alignment Summary](#)

[Metrics](#)

[SAM/BAM GC Bias Metrics](#)

[Estimate Library Complexity](#)

[Insertion size metrics for PAIRED data](#)

[SAM/BAM Hybrid Selection Metrics](#)
for targeted resequencing data

[Add or Replace Groups](#)

[Reorder SAM/BAM](#)

[Replace SAM/BAM Header](#)

Option	Description
INPUT=File	SAM or BAM file Required.
OUTPUT=File	File to write insert size metrics to Required.
REFERENCE_SEQUENCE=File	Reference sequence file Required.
ASSUME_SORTED=Boolean	If true (default), unsorted SAM/BAM files will be considered coordinate sorted
MAX_INSERT_SIZE=Integer	Paired end reads above this insert size will be considered chimeric along with inter-chromosomal pairs. Default value: 100000.
ADAPTER_SEQUENCE=String	This option may be specified 0 or more times.
IS_BISULFITE_SEQUENCED=Boolean	Whether the SAM or BAM file consists of bisulfite sequenced reads. Default value: false.
CREATE_MD5_FILE=Boolean	Whether to create an MD5 digest for any BAM files created.

The output produced by the tool has the following columns:

1. CATEGORY: One of either UNPAIRED (for a fragment run), FIRST_OF_PAIR when metrics are for only the first read
2. TOTAL_READS: The total number of reads including all PF and non-PF reads. When CATEGORY equals PAIR this value
3. PF_READS: The number of PF reads where PF is defined as passing Illumina's filter.
4. PCT_PF_READS: The percentage of reads that are PF (PF_READS / TOTAL_READS)
5. PF_NOISE_READS: The number of PF reads that are marked as noise reads. A noise read is one which is composed
6. PF_READS_ALIGNED: The number of PF reads that were aligned to the reference sequence. This includes reads tha
7. PCT_PF_READS_ALIGNED: The percentage of PF reads that aligned to the reference sequence. PF_READS_ALIGNED / F
8. PF_HQ_ALIGNED_READS: The number of PF reads that were aligned to the reference sequence with a mapping qualit
9. PF_HQ_ALIGNED_BASES: The number of bases aligned to the reference sequence in reads that were mapped at high
10. PF_HQ_ALIGNED_Q20_BASES: The subset of PF_HQ_ALIGNED_BASES where the base call quality was Q20 or higher.
11. PF_HQ_MEDIAN_MISMATCHES: The median number of mismatches versus the reference sequence in reads that were ali
12. PF_HQ_ERROR_RATE: The percentage of bases that mismatch the reference in PF HQ aligned reads.
13. MEAN_READ_LENGTH: The mean read length of the set of reads examined. When looking at the data for a single la
14. READS_ALIGNED_IN_PAIRS: The number of aligned reads who's mate pair was also aligned to the reference.
15. PCT_READS_ALIGNED_IN_PAIRS: The percentage of reads who's mate pair was also aligned to the reference. READS_
16. BAD_CYCLES: The number of instrument cycles in which 80% or more of base calls were no-calls.
17. STRAND_BALANCE: The number of PF reads aligned to the positive strand of the genome divided by the number of
18. PCT_CHIMERAS: The percentage of reads that map outside of a maximum insert size (usually 100kb) or that have
19. PCT_ADAPTER: The percentage of PF reads that are unaligned and match to a known adapter sequence right from t

Warning on SAM/BAM quality

Many SAM/BAM files produced externally and uploaded to Galaxy do not fully conform to SAM/BAM specifications. Galaxy deals with this by using the LENIENT flag when it runs Picard, which allows reads to be discarded if they're empty or don't map. This appears to be the only way to deal with SAM/BAM that cannot be parsed.

Tools		
EMBOSS		## METRICS CLASS net.sf.picard.analysis.AlignmentSummaryMetrics
Regional Variation	CATEGORY	UNPAIRED
FASTA manipulation	TOTAL_READS	511766
Evolution	PF_READS	511766
Multiple Alignments	PCT_PF_READS	1
Metagenomic analyses	PF_NOISE_READS	1
Motif Tools	PF_READS_ALIGNED	511102
NGS TOOLBOX BETA	PCT_PF_READS_ALIGNED	0.998703
NGS: Peak Calling	PF_ALIGNED_BASES	35776127
NGS: Variant Analysis	PF_HQ_ALIGNED_READS	491401
NGS: GATK Tools (beta)	PF_HQ_ALIGNED_BASES	34397052
NGS: Picard (beta)	PF_HQ_ALIGNED_Q20_BASES	34194818
FASTQ to BAM creates an unaligned BAM file	PF_HQ_MEDIAN_MISMATCHES	0
SAM to FASTQ creates a FASTQ file	PF_MISMATCH_RATE	0.001478
BAM Index Statistics	PF_HQ_ERROR_RATE	0.001458
SAM/BAM Alignment Summary Metrics	PF_INDEL_RATE	0.00005
SAM/BAM GC Bias Metrics	MEAN_READ_LENGTH	70
Estimate Library Complexity	READS_ALIGNED_IN_PAIRS	0
Insertion size metrics for PAIRED data	PCT_READS_ALIGNED_IN_PAIRS	0
SAM/BAM Hybrid Selection Metrics for targeted resequencing data	BAD_CYCLES	0
Add or Replace Groups	STRAND_BALANCE	0.497341
Reorder SAM/BAM	PCT_CHIMERAS	0
Replace SAM/BAM Header	PCT_ADAPTER	0.000006
	SAMPLE	
	LIBRARY	
	READ_GROUP	

Tools

search tools

[Get Data](#)

[Send Data](#)

[Lift-Over](#)

[Text Manipulation](#)

[Convert Formats](#)

[Filter and Sort](#)

[Join, Subtract and Group](#)

[NGS: QC and manipulation](#)

[NGS: Mapping](#)

[NGS: RNA-seq](#)

[NGS: SAMtools](#)

MPileup call variants

[bcftools view](#) Converts BCF format to VCF format

[Reheader](#) copy SAM/BAM header between datasets

[Split BAM](#) dataset on readgroups

[Stats](#) generate statistics for BAM dataset

[BAM-to-SAM](#) convert BAM to SAM

[Sort BAM](#) dataset

[CalMD](#) recalculate MD/NM tags

[BedCov](#) calculate read depth for a set of genomic intervals

[IdxStats](#) tabulate mapping statistics for BAM dataset

MPileup call variants (Galaxy Tool Version 2.1)

Choose the source for the reference genome

Use a built-in genome

BAM file

1: BAM file

BAM file

6: BWA on data 1 (mapped reads in BAM format)

+ Insert BAM file

Using reference genome

Human (Homo sapiens): hg19

Genotype Likelihood Computation

Do not perform genotype likelihood computation (output pileup)

Output base positions on reads

Yes No

--output-BP

Output mapping quality

Yes No

--output-MQ

Set advanced options

Basic

Execute

What it does

NGS: Variant Analysis

Galaxy

Analyze Data Workflow Shared Data Visualization Cloud Help User

History Using 0%

Tools

- CloudMap: Variant detection
- Mapping with WGS data Map a mutation by plotting recombination frequencies resulting from crossing to a highly polymorphic strain
- CloudMap: Variant Discovery Mapping with WGS data Map a mutation using in silico bulk segregant linkage analysis using variants that are already present in the mutant strain of interest (rather than those introduced by a cross to a polymorphic strain).
- CloudMap: in silico complementation Perform in silico complementation analysis on multiple tabular snpEff output files
- Variant Annotator process variant counts
- Phylorelatives Relatedness of minor allele sequences in NJ tree
- MAF boxplot Minor Allele Frequency Boxplot
- FASTA from allele counts Generate major and minor allele sequences from alleles table
- VarScan for variant detection**
- ANNOVAR Annotate VCF with functional information using ANNOVAR

VarScan for variant detection (Galaxy Tool Version 0.1) Options

Pileup dataset: 5: MPileup on data 4 (circled)

Analysis type: single nucleotide variation

Minimum read depth: 5 (circled)

Minimum supporting reads: 2

Minimum base quality at a position to count a read: 15

Minimum variant allele frequency threshold: 0.01

Minimum frequency to call homozygote: 0.75

p-value threshold for calling variants: 0.99

Ignore variants with >90% support on one strand: no

sample names:

History

- search datasets
- workshop-exome-seq 6 shown
- 32.8 MB
- 6: MPileup on data 4 (log) (circled)
- 5: MPileup on data 4
- 4: BWA on data 1 (mapped reads in BAM format)
- 1.9 MB format: bam, database: hg19
- [bwa_aln] 17bp reads: max_diff = 2
- [bwa_aln] 38bp reads: max_diff = 3
- [bwa_aln] 64bp reads: max_diff = 4
- [bwa_aln] 93bp reads: max_diff = 5
- [bwa_aln] 124bp reads: max_diff = 6
- [bwa_aln] 157bp reads: max_diff = 7
- [bwa_aln] 190bp reads: max_diff = 8
- [bwa_a]

display at UCSC main
display at Ensembl Current
display with IGV local Human hg19
display in IGB View

Binary bam alignments file

3: FastQC on data 1: Raw Data

Galaxy

Analyze Data Workflow Shared Data Visualization Cloud Help User Using 0%

Tools

- [CloudMap: Variant Discovery](#)
- [Mapping with WGS data](#)
- [Map a mutation using in silico bulk segregant linkage analysis using variants that are already present in the mutant strain of interest \(rather than those introduced by a cross to a polymorphic strain\).](#)
- [CloudMap: in silico complementation](#)
- [Perform in silico complementation analysis on multiple tabular snpEff output files](#)
- [Variant Annotator](#)
- [process variant counts](#)
- [Phylorelatives](#)
- [Relatedness of minor allele sequences in NJ tree](#)
- [MAF boxplot](#)
- [Minor Allele Frequency Boxplot](#)
- [FASTA from allele counts](#)
- [Generate major and minor allele sequences from alleles table](#)
- [Varscan for variant detection](#)
- [ANNOVAR Annotate VCF with functional information using ANNOVAR](#)
- [Annotate a VCF dataset with custom filters](#)
- [Slice VCF to get data from selected regions](#)

ANNOVAR Annotate VCF with functional information using ANNOVAR (Galaxy Tool Version 0.1)

Options

Variants

7: Varscan on data 5

Gene Annotations

Select/Unselect all

refGene

Annotation Regions

Select/Unselect all

Annotation Databases

Select/Unselect all

1000g2012apr_all snp137 cosmic67

Execute

What it does

This tool will annotate variants using specified gene annotations, regions, and filtering databases. Input is a VCF dataset, and output is a table of annotations for each variant in the VCF dataset.

ANNOVAR Website and Documentation

Website: <http://www.openbioinformatics.org/annovar/>

Paper: <http://nar.oxfordjournals.org/content/38/16/e164>

Important Usage Note

ANNOVAR is open-source and free for non-profit use. If you use it for commercial purposes, please contact BIOBASE (info@biobase-international.com) directly for license related issues. Also see http://www.openbioinformatics.org/annovar/annovar_faq.html#license

History

search datasets

workshop-exome-seq

7 shown

32.8 MB

7: Varscan on data 5

6: MPileup on data 4 (log)

5: MPileup on data 4

4: BWA on data 1 (mapped reads in BAM format)

1.9 MB

format: bam, database: hg19

[bwa_align] 17bp reads: max_diff = 2
[bwa_align] 38bp reads: max_diff = 3
[bwa_align] 64bp reads: max_diff = 4
[bwa_align] 93bp reads: max_diff = 5
[bwa_align] 124bp reads: max_diff = 6
[bwa_align] 157bp reads: max_diff = 7
[bwa_align] 190bp reads: max_diff = 8
[bwa_align]

display at UCSC main

display at Ensembl Current

display with IGV local Human hg19

display in IGB View

Binary bam alignments file

Galaxy

Analyze Data Workflow Shared Data Visualization Cloud Help User

Tools

resulting from crossing to a highly polymorphic strain

[CloudMap: Variant Discovery Mapping with WGS data](#) Map a mutation using *in silico* bulk segregant linkage analysis using variants that are already present in the mutant strain of interest (rather than those introduced by a cross to a polymorphic strain).

[CloudMap: in silico complementation](#) Perform *in silico* complementation analysis on multiple tabular snpEff output files

[Variant Annotator](#) process variant counts

[Phylorelatives](#) Relatedness of minor allele sequences in NJ tree

[MAF boxplot](#) Minor Allele Frequency Boxplot

[FASTA from allele counts](#) Generate major and minor allele sequences from alleles table

[Varscan for variant detection](#)

[ANNOVAR Annotate VCF](#) with functional information using ANNOVAR

[Annotate a VCF dataset with custom filters](#)

[Slice VCF](#) to get data from selected regions

Chromosome	Start	End	Type	Ref	Alt	Effect	Description	Score	Format
chr22	35193534	35193534	C	G	intergenic	LARGE(dist=877118),ISX(dist=268596)	NA	NA	History
chr22	35590706	35590706	G	A	intergenic	ISX(dist=107326),HMGXB4(dist=62739)	NA	NA	
chr22	35591035	35591035	G	A	intergenic	ISX(dist=107655),HMGXB4(dist=62410)	NA	NA	
chr22	36099882	36099882	A	G	intergenic	APOL6(dist=35426),APOL5(dist=14037)	NA	NA	
chr22	36100410	36100410	C	A	intergenic	APOL6(dist=35954),APOL5(dist=13509)	NA	NA	
chr22	36294050	36294050	C	A	intronic	RBFOX2	NA	NA	
chr22	36449831	36449831	C	T	intergenic	RBFOX2(dist=25246),APOL3(dist=86540)	NA	NA	
chr22	36449862	36449862	C	T	intergenic	RBFOX2(dist=25277),APOL3(dist=86509)	NA	NA	
chr22	36478729	36478729	G	A	intergenic	RBFOX2(dist=54144),APOL3(dist=57642)	NA	NA	
chr22	36593197	36593197	G	C	intronic	APOL4	NA	NA	
chr22	36848062	36848062	T	C	intergenic	MYH9(dist=63999),TXN2(dist=15031)	NA	NA	
chr22	36863160	36863160	C	T	UTR3	TXN2	NA	NA	
chr22	36863247	36863247	A	C	UTR3	TXN2	NA	NA	
chr22	36863491	36863491	G	A	UTR3	TXN2	NA	NA	
chr22	36863705	36863705	C	T	UTR3	TXN2	NA	NA	
chr22	36864063	36864063	G	T	intronic	TXN2	NA	NA	
chr22	36872750	36872750	T	G	intronic	TXN2	NA	NA	
chr22	36933668	36933668	G	A	intergenic	EIF3D(dist=8391),CACNG2(dist=23248)	NA	NA	
chr22	36991133	36991133	G	A	intronic	CACNG2	NA	NA	
chr22	37154444	37154444	C	A	exonic	IFT27	stopgain SNV	IFT27	
chr22	37154445	37154445	C	A	exonic	IFT27	nonsynonymous SNV	IFT27	
chr22	37163329	37163329	C	T	intronic	IFT27	NA	NA	
chr22	37268868	37268868	G	A	intronic	NCF4	NA	NA	
chr22	37268898	37268898	G	A	intronic	NCF4	NA	NA	
chr22	37318426	37318426	C	T	intronic	CSF2RB	NA	NA	
chr22	37318446	37318446	G	A	intronic	CSF2RB	NA	NA	
chr22	37318515	37318515	T	C	intronic	CSF2RB	NA	NA	
chr22	37319425	37319425	G	T	intronic	CSF2RB	NA	NA	
chr22	37325527	37325527	C	T	exonic	CSF2RB	nonsynonymous SNV	CSF2RB	
chr22	37325833	37325833	C	T	exonic	CSF2RB	synonymous SNV	CSF2RB	
chr22	37326504	37326504	C	T	exonic	CSF2RB	nonsynonymous SNV	CSF2RB	
chr22	37328751	37328751	A	G	intronic	CSF2RB	NA	NA	
chr22	37328966	37328966	G	C	intronic	CSF2RB	NA	NA	

Using 0%

History

search datasets

workshop-exome-seq
8 shown

32.8 MB

8: ANNOVAR Annotate VCF on data 7

7: VarScan on data 5

6: MPileup on data 4 (log)

5: MPileup on data 4

4: BWA on data 1 (mapped reads in BAM format)

1.9 MB

format: bam, database: hg19

[bwa_align] 17bp reads: max_diff = 2
[bwa_align] 38bp reads: max_diff = 3
[bwa_align] 64bp reads: max_diff = 4
[bwa_align] 93bp reads: max_diff = 5
[bwa_align] 124bp reads: max_diff = 6
[bwa_align] 157bp reads: max_diff = 7
[bwa_align] 190bp reads: max_diff = 8
[bwa_align]

display at UCSC main
display at Ensembl Current
display with IGV local Human hg19
display in IGB View
Binary bam alignments file

In this section**Web interface**[Input form](#)
[Results](#)**VEP script**[Tutorial](#)
[Download and install](#)
[Running the script](#)
[Caches and databases](#)
[Filters](#)
[Custom](#)
[Plugins](#)
[Examples](#)
[Other](#)
[Data files](#)
[FAQ](#)**Upload your
VCF file**[Search documentation](#)**Go****Make sure to select
correct genome
version GRCh38 or
GRCh37 (hg19)**

Variant Effect Predictor

The VEP determines the effect of your variants (SNPs, insertions, deletions, CNVs or structural variants) on genes, transcripts, and protein sequence, as well as regulatory regions. Simply input the coordinates of your variants and the nucleotide changes to find out the:

- **genes and transcripts** affected by the variants
- **location** of the variants (e.g. upstream of a transcript, in coding sequence, in non-coding RNA, in regulatory regions)
- **consequence** of your variants on the protein sequence (e.g. stop gained, missense, stop lost, frameshift)
- **known variants** that match yours, and associated minor allele frequencies from the **1000 Genomes Project**
- **SIFT** and **PolyPhen** scores for changes to protein sequence



IGV

Integrative Genomics Viewer



Wexner Medical Center

<http://www.java.com/en/>



java web start launcher



Apple

Yahoo!

Google Maps

YouTube

Wikipedia

News (13)

Popular

Galaxy

Downloads | Integrat...

java.com: Java + You



Download Help

Search



JAVA + YOU, DOWNLOAD TODAY!

[Free Java Download](#)



Feedback

» [What is Java?](#) » [Do I have Java?](#) » [Need Help?](#)

About Java



Java Developer



Java + Alice



Java + Greenfoot



Java + BlueJ



Oracle Academy



Java Magazine

Java Downloads for All Operating Systems

http://www.java.com/en/download/manual.jsp

java web start launcher

Apple Yahoo! Google Maps YouTube Wikipedia News (13) Popular

Galaxy Downloads | Integrat... Java Downloads for ...

Search

Java™

Download Help

Available Operating Systems

- » Windows
- » Mac
- » Linux
- » Solaris

Help Resources

- » Troubleshoot Java

Java 7

- » Where can I get Java 7?

JDK

- » Looking for the JDK?

Java Downloads for All Operating Systems

Recommended Version 8 Update 31

Select the file according to your operating system from the list below to get the latest Java for your computer.

> [Remove Older Versions](#) > [What is Java?](#)

By downloading Java you acknowledge that you have read and accepted the terms of the [end user license agreement](#)

 Windows  [Which should I choose?](#)

 Windows Online (32-bit) filesize: 624 KB	Instructions	After installing Java, you may need to restart your browser in order to enable Java in your browser.
 Windows Offline (32-bit) filesize: 29.0 MB	Instructions	
 Windows Offline (64-bit) filesize: 89.1 MB	Instructions	

If you use 32-bit and 64-bit browsers interchangeably, you will need to install both 32-bit and 64-bit Java in order to have the Java plug-in for both browsers. » [FAQ about 64-bit Java for Windows](#)

Feedback

www.broadinstitute.org/igv/

igv
Integrative Genomics Viewer

- Home
- Downloads
- Documents
 - Hosted Genomes
 - FAQ
 - IGV User Guide
 - File Formats
 - Release Notes
 - IGV for iPad
 - Credits
- Contact

Search website

[Broad Home](#) [Cancer Program](#)

BROAD INSTITUTE
© 2013 Broad Institute

Home

Integrative Genomics Viewer



What's New

NEWS September 2014. The IGV iPad app can now be installed from the Apple App Store. *IGV for iPad* is a lightweight genomic data viewer that provides some of the functionality available in our regular desktop IGV. See the [IGV for iPad documentation](#) for details.

Overview

The Integrative Genomics Viewer (IGV) is a high-performance visualization tool for interactive exploration of large, integrated genomic datasets. It supports a wide variety of data types, including array-based and next-generation sequence data, and genomic annotations.

Downloads

Please [register](#) to download IGV. After registering, you can log in at any time using your email address. Permission to use IGV is granted under the [GNU LGPL license](#).

Citing IGV

To cite your use of IGV in your publication:

James T. Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander, Gad Getz, Jill P. Mesirov, **Integrative Genomics Viewer**. *Nature Biotechnology* 29, 24–26 (2011)

Helga Thorvaldsdóttir, James T. Robinson, Jill P. Mesirov. **Integrative Genomics Viewer (IGV): high-performance genomics data visualization and exploration**. *Briefings in Bioinformatics* 14, 178–192 (2013).

Funding

Development of IGV is made possible by funding from the [National Cancer Institute](#), the [National Institute of General Medical Sciences](#) of the [National Institutes of Health](#), and the [Starr Cancer Consortium](#).

IGV participates in the [GenomeSpace](#) initiative, which is funded by the [National Human Genome Research Institute](#).

Integrative Genomics Viewer (IGV) (Version 2.3)

Install IGV

Options for installing and running IGV:

1. (Mac only) Download and run the Mac application; or
2. (All systems) Use the Java Web Start buttons (Mac users: see below for limitations); or
3. (All systems) Download the binary distribution and run IGV from the command line.

1. Mac Application

Download and unzip the Mac App archive, then double-click the IGV application to run it. The application can be moved to the "Applications" folder, or anywhere else. **Note: This requires Java 7. Mac users with Java 6 (JRE 1.6) should use the binary distribution archive or the Java Web Start buttons below.**

[Download Mac App](#)



2. Java Web Start

The buttons below use Java Web Start (JWS) to install and launch IGV directly from our web site.

***Mac Users:** The Java Web Start option does not work for some users due to security settings. The recommended solution is to use the bundled Mac App from the link above. Alternatively you can try to work around this by right-clicking on the buttons and saving the "jnlp" file, then right-click on the "jnlp" file and select "Open With Java Web Start".

Chrome: Chrome does not automatically launch the Java Webstart files by default. Instead, the launch buttons below will download a "jnlp" file. This should appear in the lower left corner of the browser. Double-click the downloaded file to run, or if on a Mac right-click and select "Open With Java Web Start".

Windows users: To run with more than 1.2 GB of memory you must install 64-bit Java. **Most Windows installs do not include 64-bit Java by default, even if the operating system is 64-bit.** Attempting to use the 2GB or greater launch options with 32-bit Java will result in the error "could not create virtual machine".



 Launch Launch with 750 MB	 Launch Launch with 1.2 GB Maximum usable memory for Windows OS with 32-bit Java.	 Launch Launch with 2 GB Maximum usable memory for 32-bit MacOS.	 Launch Launch with 10 GB For large memory machines with 64-bit Java.
--	--	---	--

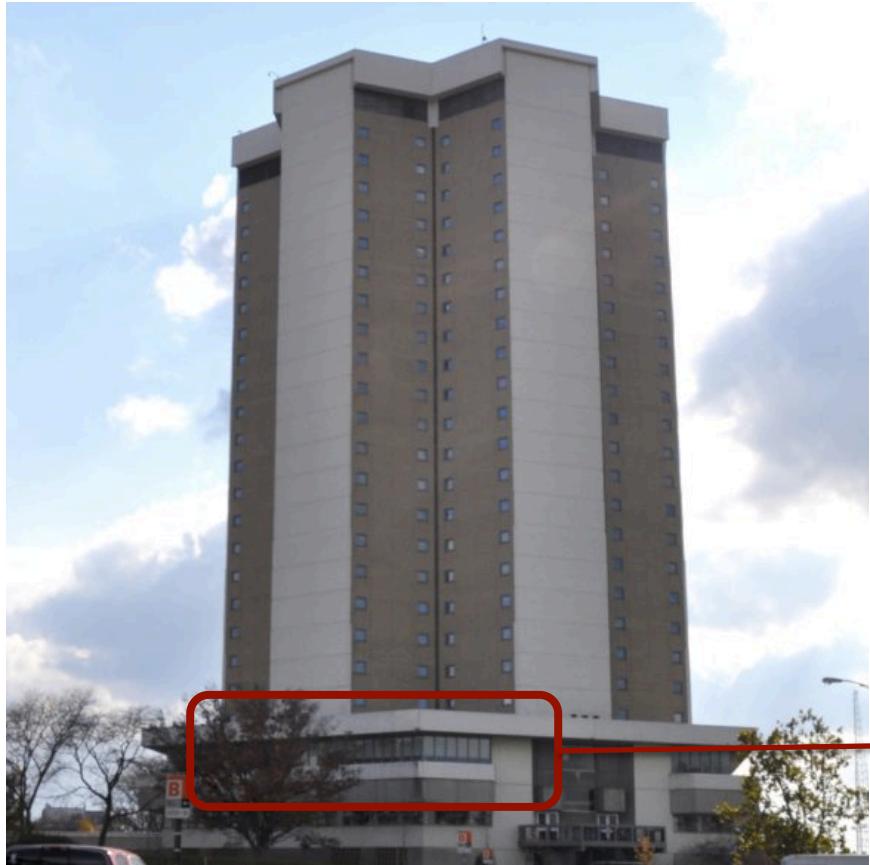
Acknowledgement



Selen Yilmaz, MS



| Wexner Medical Center



Gulcin Ozer, PhD
340C Lincoln Tower
1800 Cannon Dr.

Gulcin.Ozer@osumc.edu



Wexner
Medical
Center