

Genomic Resources

Michael Schatz

Aug 29, 2012

QB Bootcamp Lecture 3





Outline

Part 1: Overview & Fundamentals

Part 2: Sequence Analysis Theory

Part 3: Genome Resources

- Public: NCBI, UCSC
- CSHL: Intranet, Meetings, Galaxy

Part 4: Example Analysis

NCBI

<http://www.ncbi.nlm.nih.gov/>

NCBI

Resources ▾ How To ▾

mschatr@cshl.edu My NCBI Sign Out

NCBI

National Center for Biotechnology Information

All Databases ▾

Search

NCBI Home

Resource List (A-Z)

All Resources

Chemicals & Bioassays

Data & Software

DNA & RNA

Domains & Structures

Genes & Expression

Genetics & Medicine

Genomes & Maps

Homology

Literature

Proteins

Sequence Analysis

Taxonomy

Training & Tutorials

Variation

Welcome to NCBI

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[About the NCBI](#) | [Mission](#) | [Organization](#) | [Research](#) | [RSS Feeds](#)

Get Started

- [Tools](#): Analyze data using NCBI software
- [Downloads](#): Get NCBI data or software
- [How-To's](#): Learn how to accomplish specific tasks at NCBI
- [Submissions](#): Submit data to GenBank or other NCBI databases

Genomic Structural Variation

dbVar archives large scale genomic variation data and associates defined variants with phenotypic information.

1 2 3 4 5 6 7 8

Popular Resources

[PubMed](#)

[Bookshelf](#)

[PubMed Central](#)

[PubMed Health](#)

[BLAST](#)

[Nucleotide](#)

[Genome](#)

[SNP](#)

[Gene](#)

[Protein](#)

[PubChem](#)

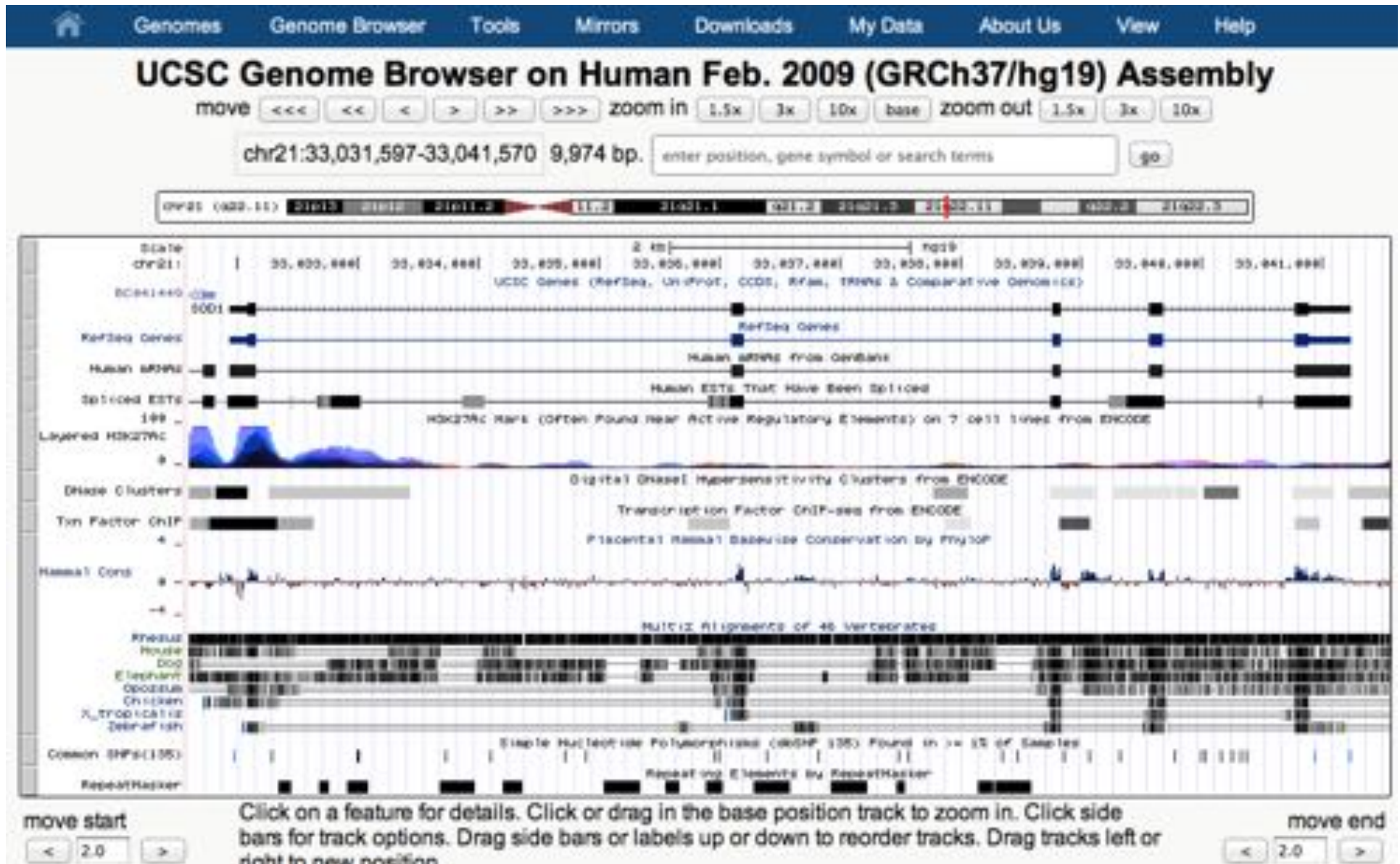
NCBI Announcements

NCBI's July Newsletter is on the Bookshelf

13 Aug 2012

Introduction to the 1000 Genomes Browser, PubMed's Citation Manager and

<http://genome.ucsc.edu/>



Intranet

<http://intranet.cshl.edu/it/bluehelix/>

**COLD SPRING
HARBOR
LABORATORY**

Intranet

Tuesday, August 28, 2012

Search **FACES**

Home • Most Popular • Administration • Education • General Info • Lab Websites • Science • Services •

CSHL Information Technology

Home • About IT • Applications • **BlueHelix** • FAQ • Helpdesk • Stats

 **BlueHelix (HPCC)**

- BlueHelix Home
- HPCC Talks
- Architecture
- Applications
- Databases
- Request an Account
- Login & File Transfer
- Using BlueHelix
- Dev Node Status
- SGE Job Status
- Maintenance Downtime
- Comments or Questions

BlueHelix, the High Performance Computing Cluster (HPCC) is installed at the main IT DataCenter in the Hilltop complex. The system is an IBM E1350 blade based cluster running Linux. It is comprised of 10 racks with a total of 504 compute nodes, 2 management servers, and a number of storage servers. The system has a Gigabit Ethernet communication network.



- Each of the 504 compute nodes contains 2 dual core Opteron 64-bit processors running at 2.0GHz (2016 cores total), and a local 73GB SFF SCSI disk. There are 56 nodes with 2GB of memory, 364 nodes with 4GB of memory, and 84 nodes with 8GB of memory.
- Six of the nodes are development nodes, used for editing, compiling, and debugging user codes, jobs, and so on.

Conferences and Journals

CSHL Yearly Conferences

Biology of Genomes Symposium	May	Latest advances in biology, genomics, and medicine
Genome Informatics	May/June	Latest advances with yearly themes
Personal Genomes	Sept/Nov	Computational Biology
In-house Symposium	Nov	Computational Biology
		Updates from the faculty (Just before Thanksgiving)

You are welcome to attend all meetings at CSHL free of charge:

<http://meetings.cshl.edu/meetings.html>

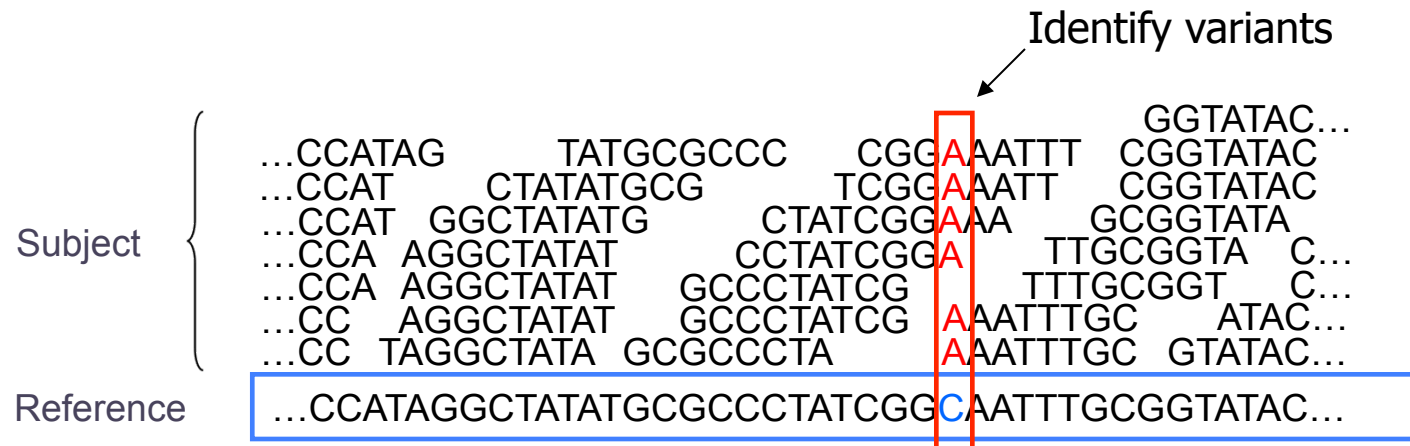
Journals (RSS feeds and eTOC available)

Bioinformatics	Genome Biology	Genome Research
Nature	Nature Biotechnology	Nature Methods
PNAS	PLoS Biology	Science

<http://genomics.cshl.edu>



Short Read Mapping



- Given a reference and many subject reads, report one or more “good” end-to-end alignments per alignable read
 - Fundamental computation to genotyping and many assays
 - RNA-seq Methyl-seq FAIRE-seq
 - ChIP-seq Dnase-seq Hi-C-seq
- Desperate need for scalable solutions
 - Single human requires >1,000 CPU hours / genome
 - **1000 hours * 1000 genomes = 1M CPU hours / project**

Paired-end and Mate-pairs

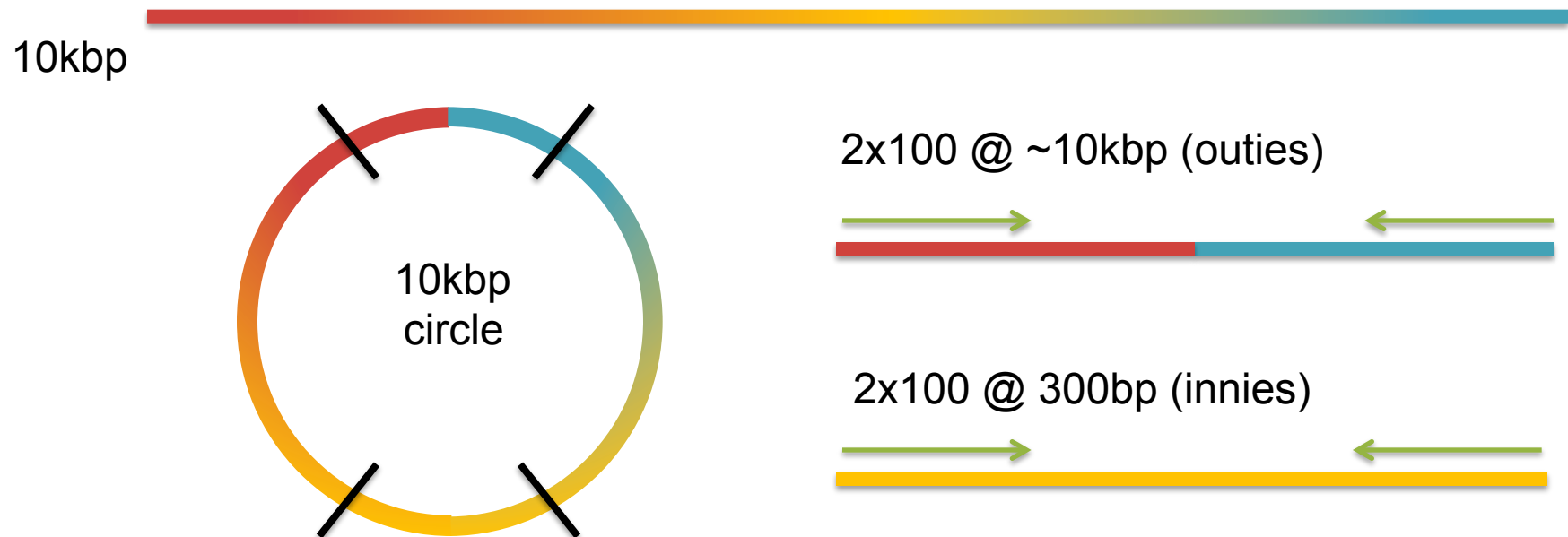
Paired-end sequencing

- Read one end of the molecule, flip, and read the other end
- Generate pair of reads separated by up to 500bp with inward orientation



Mate-pair sequencing

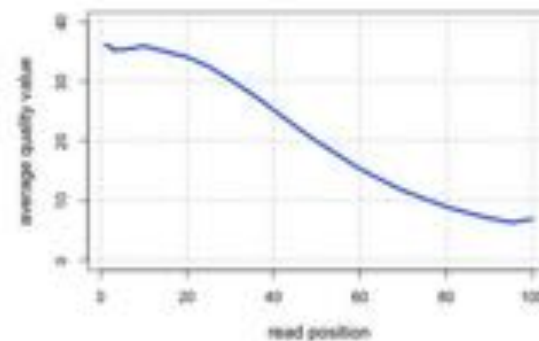
- Circularize long molecules (1-10kbp), shear into fragments, & sequence
- Mate failures create short paired-end reads



Illumina Quality

QV	P _{error}
40	1/10000
30	1/1000
20	1/100
10	1/10

$$Q_{\text{sanger}} = -10 \log_{10} p$$

[illegible]

http://en.wikipedia.org/wiki/FASTQ_format

Galaxy Exercise

1. Download data:
 - `curl -O http://schatzlab.cshl.edu/teaching/exercises/mapping/mapping.tgz`
2. Unpack and upload to Galaxy
 - Set fastq type to fastqillumina of reads
3. Map with Bowtie for Illumina
 - Aligns the reads to the reference genome
4. SAM-to-BAM
 - Converts from ASCII text file to interval representation
5. Coverage Plot of BAM
 - Mapping Statistics
6. Call variants with FreeBayes
 - Print Stats (search vcf)

Other Resources

Resource	URL	Description
Google	http://www.google.com	Internet Search
Google Scholar	http://scholar.google.com/	Literature Searches
SeqAnswers	http://seqanswers.com/	Bioinformatics Forum
Wikipedia	http://www.wikipedia.org/	Overview on anything
Circos	http://circos.ca/	Circular Genome Plots
GraphViz	http://www.graphviz.org/	Graph Visualization
EndNote	http://endnote.com/	Citation Manager
R	http://www.r-project.org/	Stats & Visualizations
Weka	http://www.cs.waikato.ac.nz/ml/weka/	Data Mining
IGV	http://www.broadinstitute.org/igv/	Read Mapping Viz
Schatz Lab	http://schatzlab.cshl.edu/teaching/	Exercises and Lectures

Questions?

<http://schatzlab.cshl.edu>