# Genomic Resources

Michael Schatz
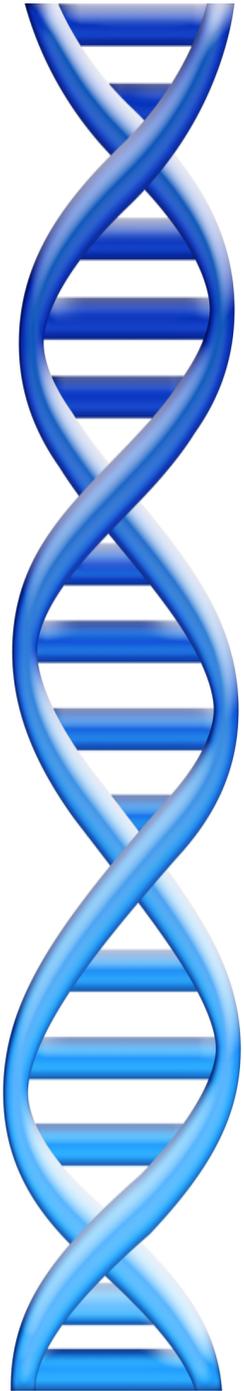
CSH

# Outline

Part 1: Overview & Fundamentals

Part 2: Sequence Analysis Theory

**Part 3: Genome Resources**
- **Public: NCBI, UCSC**
- **CSHL: Intranet, Meetings, Galaxy**

Part 4: Unix Scripting

Part 5: Example Analysis

# NCBI

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

# UCSC Genome Browser

## http://genome.ucsc.edu/

# Intranet

http://intranet.cshl.edu/IT-HPCC/blacknblue.html



**BlackNBlue**

BlackNBlue is an institutionally shared compute cluster introduced in 2012. The cluster is intended to support the full spectrum of CSHL research computing efforts and accommodates both standard batch processing and calculations implemented in the Hadoop framework.

BlackNBlue is a 1,696-core IBM System x solution based on the M4 server line with Intel Xeon E5 (Sandy Bridge-EP) processors. The cluster was designed from 106 servers, configured as development, compute, and management nodes, using 10 Gigabit per second Ethernet networking.

Two development nodes provide the sole point of user access to the cluster and allow for interactive development work as well as submission of batch and Hadoop jobs to the compute nodes. The cluster is administered from a pair of management nodes running UGE (formerly SGE), a "fair share" resource management system for the equitable allocation of compute resources. The management nodes are configured for failover protection that ensures uninterrupted execution of batch jobs in the event that the primary management node becomes unavailable.

The development nodes and 100 compute nodes have Xeon E5-2665 processors running at 2.40 GHz. The development nodes have 64GB of memory, the compute nodes 128GB. Each node has two sockets with 8 cores per socket, for a total of 16 physical cores. Hyperthreading doubles the number of physical cores, resulting in 32 virtual cores per node, which provides a total of 3,200 UGE job slots over the 128GB compute nodes.

In addition to the standard compute nodes, the cluster has two high memory nodes, each with 1.5TB of memory. The high-memory nodes have Xeon E5-4650 processors running at 2.70 GHz. Each node has 4 sockets with 8 cores per socket, for a total of 32 physical cores. With hyperthreading, users see 64 virtual cores, or UGE job slots, for each high-memory node.

The BlueArc, Isilon, and IBM SONAS storage systems are connected to all nodes via NFS.

Last Updated (Wednesday, 05 December 2012 14:27)

# Conferences and Journals

## CSHL Yearly Conferences

| | | |
|---|---|---|
| Biology of Genomes | May | Latest advances in biology, genomics, and medicine |
| Symposium | May/June | Latest advances with yearly themes |
| Genome Informatics | Sept/Nov | Computational Biology |
| Personal Genomes | Sept/Nov | Computational Biology |
| In-house Symposium | Nov | 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 |

# Galaxy

http://usegalaxy.org  http://genomics.cshl.edu

# Genotyping

Heterozygous variant?                                    Homozygous variant

GGTATAC…

Subject
```
…CCATAG      TGTGCGCCC      CGGAAATTT   CGGTATAC
…CCAT     CTATGTGCG        TCGGAAATT    CGGTATAC
…CCAT  GGCTATGTG        CTATCGGAAA     GCGGTATA
…CCA  AGGCTATAT        CCTATCGGA      TTGCGGTA   C…
…CCA  AGGCTATAT     GCCCTATCG        TTTGCGGT     C…
…CC   AGGCTATAT     GCCCTATCG  AAATTTGC      ATAC…
…CC  TAGGCTATA  GCGCCCTA      AAATTTGC  GTATAC…
```

Reference    …CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC…

- Sequencing instruments make mistakes
  - Quality of read decreases over the read length

- A single read differing from the reference is probably just an error, but it becomes more likely to be real as we see it multiple times
  - Often framed as a Bayesian problem of more likely to be a real variant or chance occurrence of N errors
  - Accuracy improves with deeper coverage

# Illumina Quality

| QV | p_error |
|----|---------|
| 40 | 1/10000 |
| 30 | 1/1000 |
| 20 | 1/100 |
| 10 | 1/10 |

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



```
  SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS...............................................
  ..........................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....................
  ...............................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII..................
  .................................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ...................
  LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.................................................
  !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
  |                        |    |        |                                    |               |
 33                       59   64       73                                  104             126

S - Sanger        Phred+33,  raw reads typically (0, 40)
X - Solexa        Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64,  raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64,  raw reads typically (3, 40)
   with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
   (Note: See discussion above).
L - Illumina 1.8+ Phred+33,  raw reads typically (0, 41)
```

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

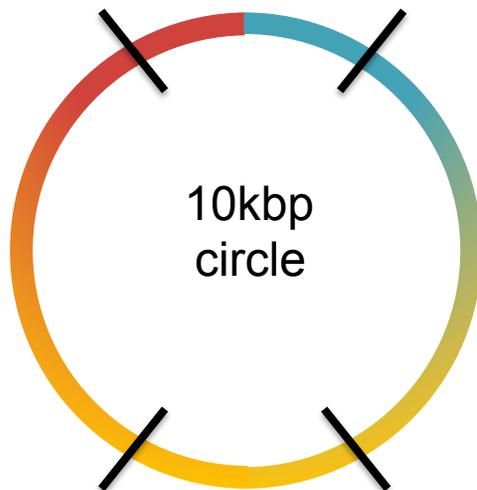# Paired-end and Mate-pairs

**_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

300bp

**_Mate-pair sequencing_**

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

10kbp

10kbp circle

2x100 @ ~10kbp (outies)

2x100 @ 300bp (innies)

# Galaxy Exercise

1. Download data:
   – 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