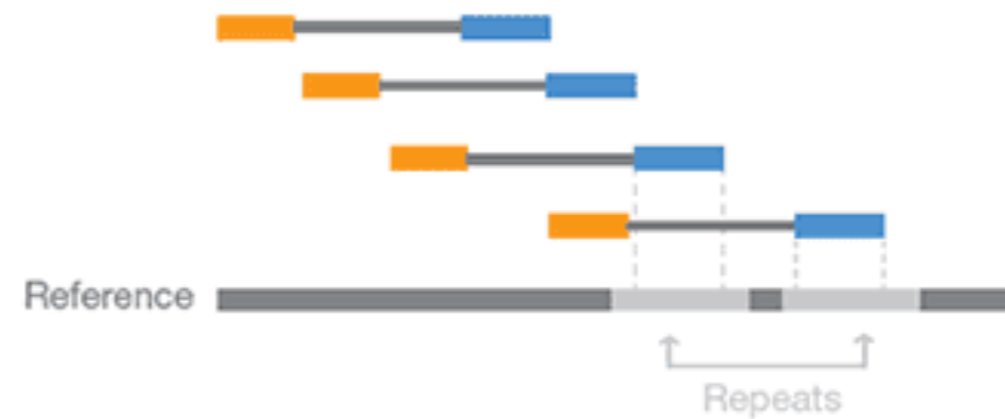# High-Throughput Sequencing: from Raw Reads to Variants

Brice A. J. Sarver
University of Montana
Division of Biological Sciences
ConGen 2015

Paired-end sequencing enables both ends of the DNA fragment to be sequenced. Because the distance between each paired read is known, alignment algorithms can use this information to map the reads over repetitive regions more precisely. This results in much better alignment of the reads, especially across difficult-to-sequence, repetitive regions of the genome.

HTS data analysis is ultimately about placing reads and accounting for uncertainty

Illumina Corporation

# Sources of Error

- (Library Preparation)

- Sequencing

- Read Cleaning

- Assembly

- Mapping

- Post-mapping Processing (e.g., indel realignment, etc.)

- Variant Calling

- Post-processing

# Goals

- Introduce a (general) workflow for data analysis

- Describe the data structures of common files

- Perform analyses using the Unix command line

- After: Tiago will (among other things) analyze the same dataset using Galaxy

- Please ask me about your own datasets!

# Recommendations

If you are going to be analyzing large datasets or lots of libraries, using the command line or custom scripts may be the best way to go

- Unix (+ a shell language like bash)

- A scripting language (R or Python, possibly Perl)

- Application development: a compiled language (C, C++)

Lots of help available in several great communities

# Exercises

- All commands in the text

- Starting from *cleaned reads*

  - You need to do some processing beforehand!

- We'll stop periodically and examine some of the files that we've been generating
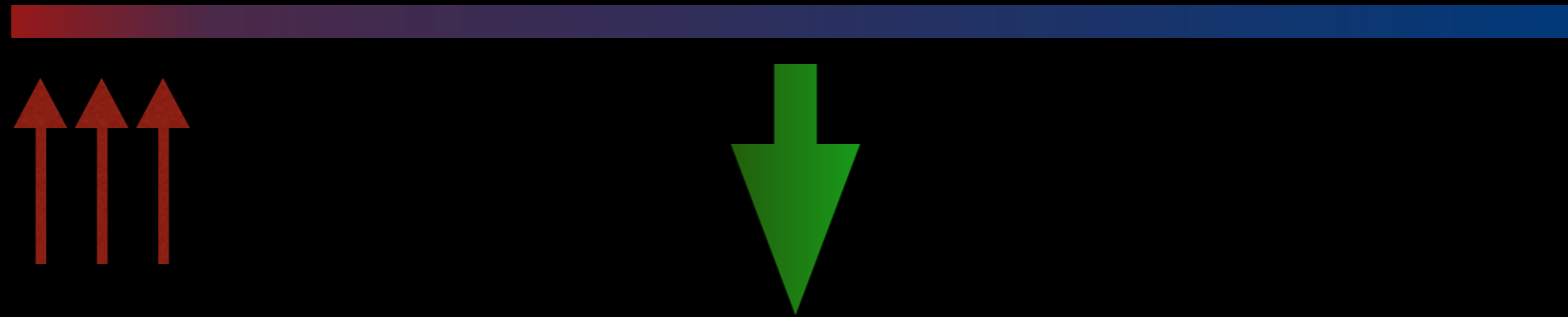
# Mouse exome capture example



*Mus spretus* (Algerian mouse)

- 55 Mbp capture
  - Genome is ~2.8 Gbp (~2%)
- Carry over annotation information from the reference
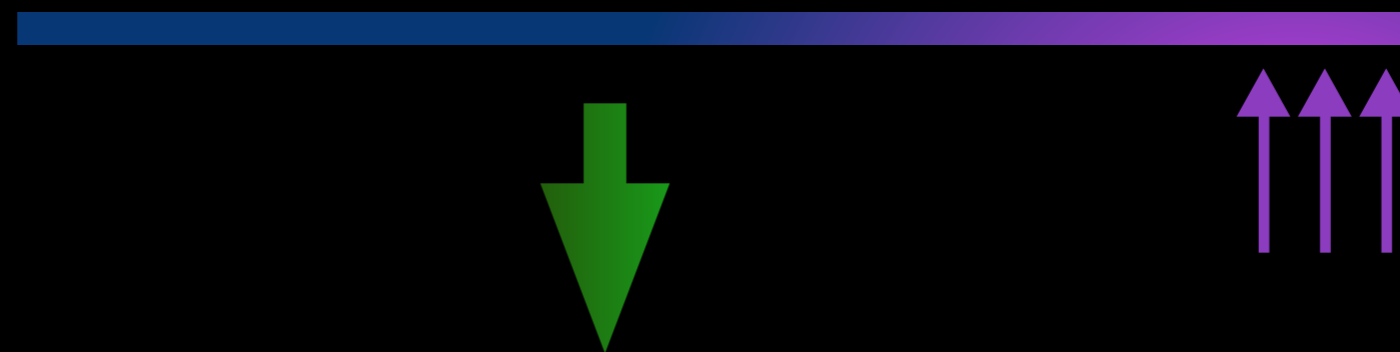- In place of the whole genome, we will be working with a 20 Mbp region of chromosome 1 containing approximately 50 transcripts
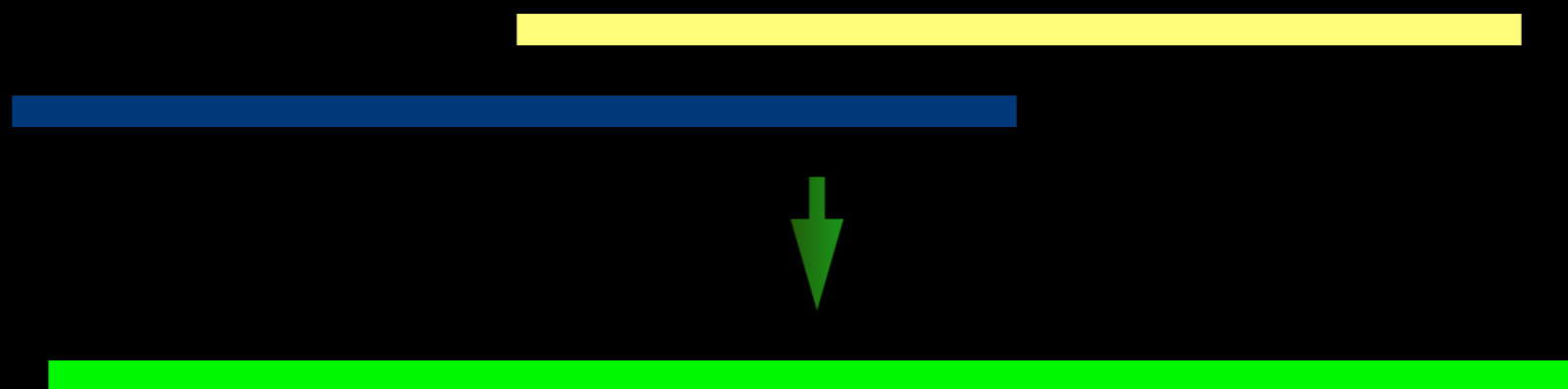
# Pre-processing of Reads

Removal of low-quality base calls

Removal of residual adapter sequence

Merge reads that overlap

# The FASTQ format

```
congen@congen-VirtualBox:~/data/brice/test$ zcat 10252.final.fq.gz | head -n 12
@HS3:309:D2385ACXX:5:2210:7285:54270
TATCCAGCCAGCCTGGCTTAGATGGTGAGTGAGCGCCAGGCCAATGAGGAACCTGTGCCATGGACCGGGCCTAGTCAGCTCCCCTCAATTCGTGGGAATC
+
BBBFFFFFFFFFFFIIIIIIIIIIIIIIFIIIFIIIIIIIIIIIIIIIIIIIIIIIIIFIIIIFFIIIFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFB
@HWI-ST665R:136:C1G9MACXX:2:1302:14922:92414
TATCCAGCCAGCCTGGCTTAGATGGTGAGTGAGCGCCAGGCCAATGAGGAACCTGTGCCATGGACCGGGCCTAGTCAGCTCCCCTCAATTCGTGGGAATC
+
CCCFFFFFHGHHHIIFIGIIJGIJ@ABFHHGHEBGAGHBHBFEGGHHGIJHIGGGHIJIJGEHB?CB>B@@BCCCCDDCACCBDDBBBDDDDB>8<>?ABC
@DQNZZQ1:722:C2J9EACXX:7:2106:18279:13220
AAATGCCACGGACTGGGCTCAGTAGGCCCCCCTCAATCCATGGGAATCAGGGTTTCGGACAGATGGGCACAGAGTCGGTGAAAATAGGGTGACAAACAGACAGGACATAAGGAAGTGTGCTGAATCTGAATGT
+
<@@DDDDEDADAFGIBBHHGFG@9CGBHIIHIEDHB)=BBCFHGGA?CDGEHB?ECE?=A>CD5DA<GCF@8F:?=HGGED?BCGE@HG?FBGEHEBDBEIJIIFHH>IIIIIGIIGGGEHHFHDDDDDDDB@@
```

## Four lines per read:

1. Description - starts with @
   - @[Machine Identifier]:[Run Identifier]:[Flowcell ID]:[Flowcell Lane]:[Tile]:[x-coordinate]:[y-coordinate]
2. The base calls
3. +
4. Quality scores (ASCII)

# The SAM format



**Mandatory Fields (tab-delimited):**

QNAME: Query template/pair NAME
FLAG: bitwise FLAG
RNAME: Reference sequence NAME
POS: 1-based leftmost POSition/coordinate of clipped sequence
MAPQ: MAPping Quality (Phred-scaled)
CIGAR: extended CIGAR string
MRNM: Mate Reference sequence NaMe (`=' if same as RNAME)
MPOS: 1-based Mate POSistion
TLEN: inferred Template LENgth (insert size)
SEQ: query SEQuence on the same strand as the reference
QUAL: query QUALity (ASCII-33 gives the Phred base quality)
OPT: variable OPTional fields in the format TAG:VTYPE:VALUE

# The Variant Call Format



GT: genotype
AD: allelic depth
DP: depth of coverage
GQ: genotype quality
PL: Phred-scaled genotype likelihoods