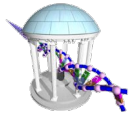# BCB 716 - Sequence Analysis



- A combined problem set 1 and 2 will go out on Thursday

- Your course logins should work now
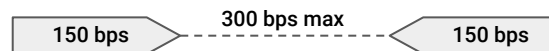
- Password is your PID

DNA Variant Calling and Analysis
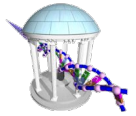
# From last time

- **Aligners generate SAM files**
    - An attempt is made to find the closest match for a given read, or read-pair to a reference
    - Alignments are performed independently and in parallel
    - SAM files include
        - the original sequence and quality string from the FASTQ file
        - Initially read pairs are considered together
            - Alignment tolerances
            - Opposite strands
            - Must satisfy a maximum gap distance
        - A placement of the first base that is "normalized" to reference orientation
        - An alignment represented as a CIGAR string
        - Various alignment scores (edit distances, etc.)
- **SAM files are a lot to interpret**
    - Statistic provide a rough idea
    - Localized analysis provides more insights

150 bps    300 bps max    150 bps

# SAM to BAM

- **SAM files tend to be large and difficult to index and manipulate**
- **Converted into Binary Alignment Maps (BAM files)**
- **This is done using a toolset called SAMtools**
- **First to convert a SAM file to a BAM file**

```
$ samtools view -S -b CC053.sam -o CC053.bam
$ ls -l CC053.*
-rw-rw-r-- 1 mcmillan its_faculty_psx 5.1G Nov  8 14:57 CC053.bam
-rw-rw-r-- 1 mcmillan its_faculty_psx  24G Nov  8 14:22 CC053.sam
```

- **BAM files are smaller, and not simply text, making them easier to search**

```
$ samtools view CC053.bam | head -1
A00434:231:H2K7FDSX2:1:1101:10529:1157  99       14       55067154        42
100M     =       55067503        449
GGCTGGAGATGGGGCTGGAGAAGGCGGCTGATCAGGGCTTTCTGAGGGCTCCCTGGAGCCCTCGACTGGCGCCAGGGAAGG
CTCAAGAGGAGGATCTGGG
FFFFFFF:FFFFFFF:FFFFFFFFFF:FFFFFFFFFFF:FFFFF:FFFF:FFFFFFFFFFFFFFFFFFFFFFFFFF:FF
FFFFFFFFFFFFFFFFFFFF AS:i:-5 XN:i:0  XM:i:1  XO:i:0       XG:i:0  NM:i:1
MD:Z:77G22       YS:i:-4 YT:Z:CP
```
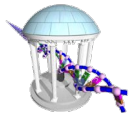
# Sorted and Indexed BAMs

- **The reads in a BAM file are roughly in the order they can out of the sequencer**
- **SAM tools provides a tool to sort the reads genomically**

```
$ samtools sort CC053.bam -o CC053.sorted.bam
$ ls -l CC053.*
-rw-rw-r-- 1 mcmillan its_faculty_psx 5.1G Nov  8 14:57 CC053.bam
-rw-rw-r-- 1 mcmillan its_faculty_psx  24G Nov  8 14:22 CC053.sam
-rw-rw-r-- 1 mcmillan its_faculty_psx 3.0G Nov  8 15:11 CC053.sorted.bam
```

- **BAM files are even smaller, nearby sequences overlap and compress better**
- **Last of all we build an index so that the BAM file is easier to search/load**

```
$ samtools index CC053.sorted.bam
$ ls -l CC053.*
-rw-rw-r-- 1 mcmillan its_faculty_psx 5.1G Nov  8 14:57 CC053.bam
-rw-rw-r-- 1 mcmillan its_faculty_psx  24G Nov  8 14:22 CC053.sam
-rw-rw-r-- 1 mcmillan its_faculty_psx 3.0G Nov  8 15:11 CC053.sorted.bam
-rw-rw-r-- 1 mcmillan its_faculty_psx 3.0M Nov  8 15:18 CC053.sorted.bam.bai
```

# Exercise

**Go to the following website:**

**https://ondemand.rc.unc.edu**

**You will need to authenticate with your ONYEN**

**Eventually you will get here:**

Click here and pick:

# Wait here for a few seconds



Wait for this button to appear.
Then press it

# Eventually you'll get here

# Now type a few commands at the command line

- **Install an initial set of bioinformatic modules:**

```
$ cp /proj/mcmillanlab/BCB716F21/loadModules .
$ loadModules
$ module list



Currently Loaded Modules:
1) samtools/1.9     3) bowtie2/2.4.1     5) minimap2/2.17
2) bwa-mem2/2.2.1   4) igv/2.8.7
```
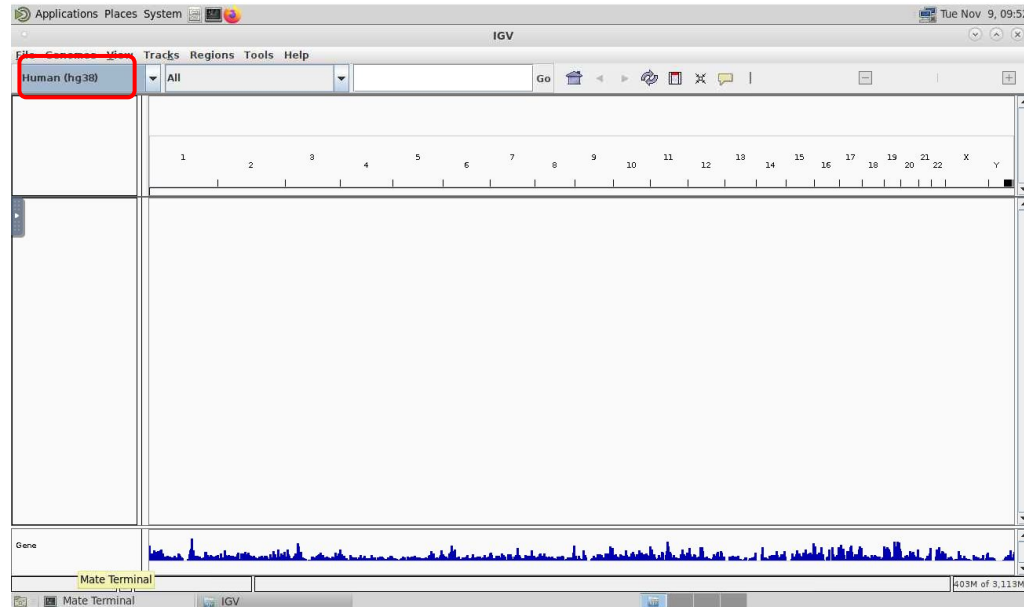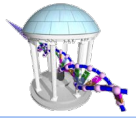
- **Today we'll discuss IGV**

# Integrative Genomics Viewer (IGV)

- **I typed:**

  `$ igv &     # starts the viewer as a background process`

- **After some machinations, and maximizing**

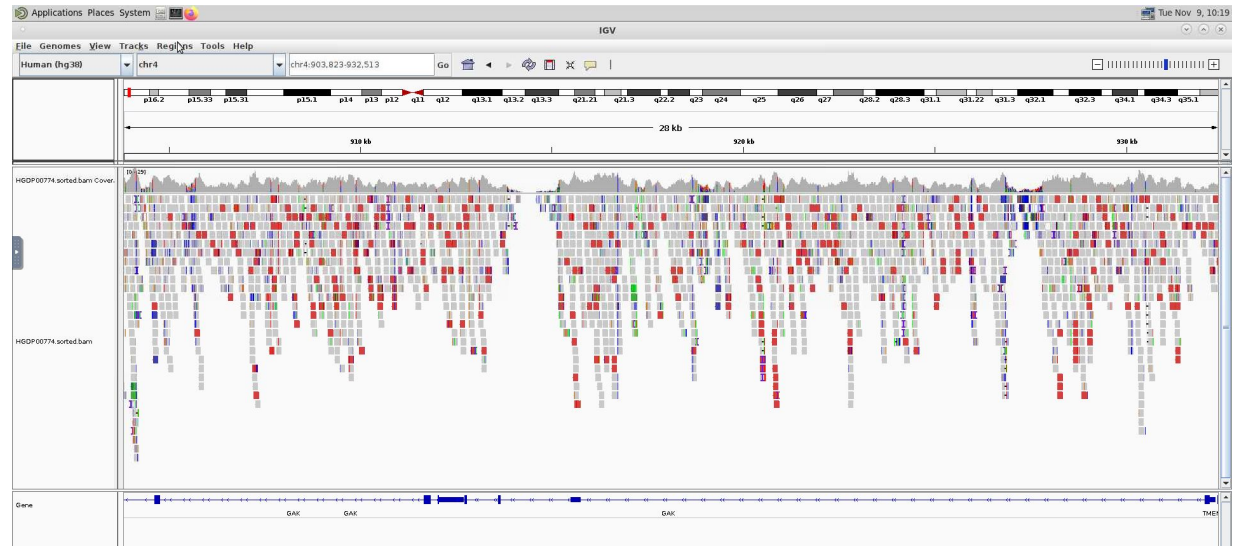First, you'll need to make sure you are using the correct genome.
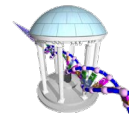
I'll use Human (hg38)

# Visualizing BAM files

- **The Interactive Genome Viewer (IGV) is a standard tool for visualizing sorted BAM files with index files**

- **You won't see any reads until you get to a window smaller than 30 kb (configurable, but)**
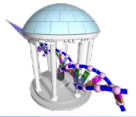
- **Coverage above**

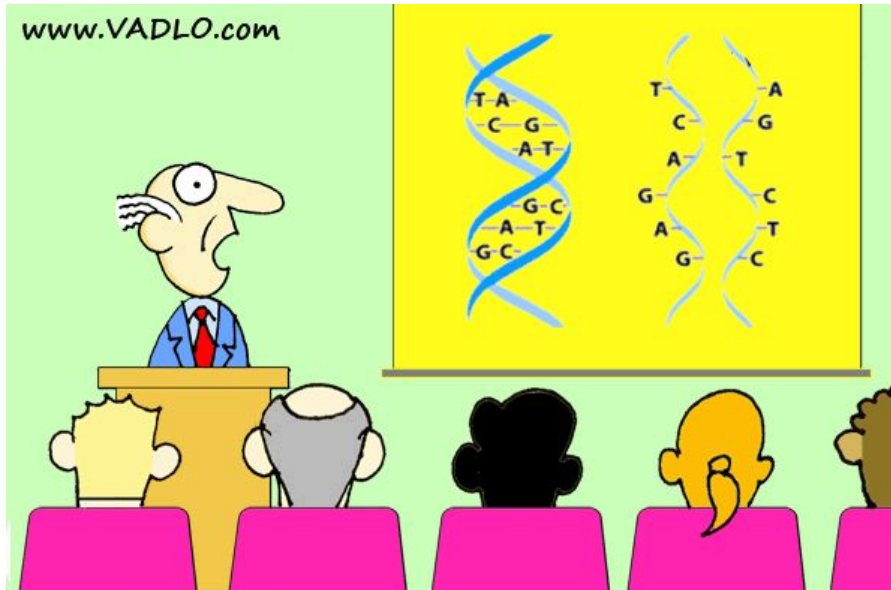- **Alignments below**

# Visualizing BAM files

- **The reads are labelled with variants and INDELS that differ from the reference**

- **Red reads are separated from mates by a larger gap than expected**

# Next Time

## Visualizing, Interpreting, and Analyzing Alignment outputs