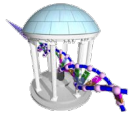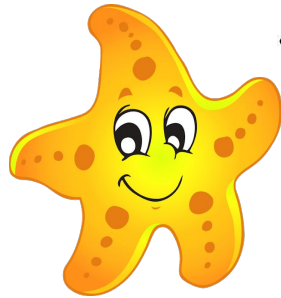# BCB 716 - Sequence Analysis

The fishy business behind RNAseq analysis

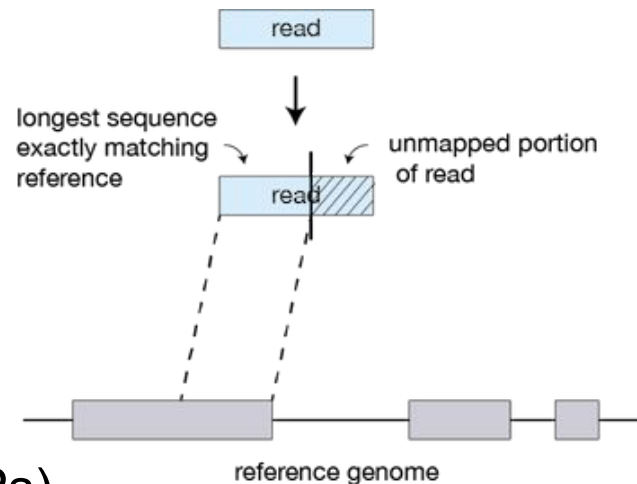- Regarding the problem sets... at least you are all in the same boat

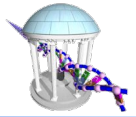RNA Analysis Pipelines

# Alignment based RNAseq tools

Splice-aware genome aligners (STAR, HiSat2, MapSplice2, GSNAP, etc)

- STAR (Spliced Transcripts Alignment to a Reference)
- Fast two-stage alignment
  - Seed search
  - Clustering, stitching, and scoring
- STAR searches for the longest exact match to one or more genomic locations. These longest matching sequences are called the Maximal Mappable Prefixes (MMPs)
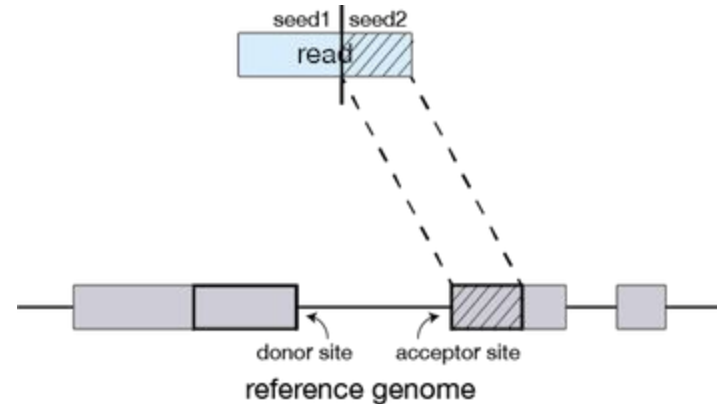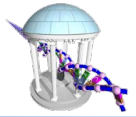
# Seed Mapping

Different parts of the read, called 'seeds', are mapped separately. The first MMP that is mapped to the genome is called seed1.

STAR next searches for the unmapped portion of the read to find the next MMP, which will be seed2.

This sequential "greedy" search for unmapped portions of reads underlies the efficiency of the STAR algorithm.

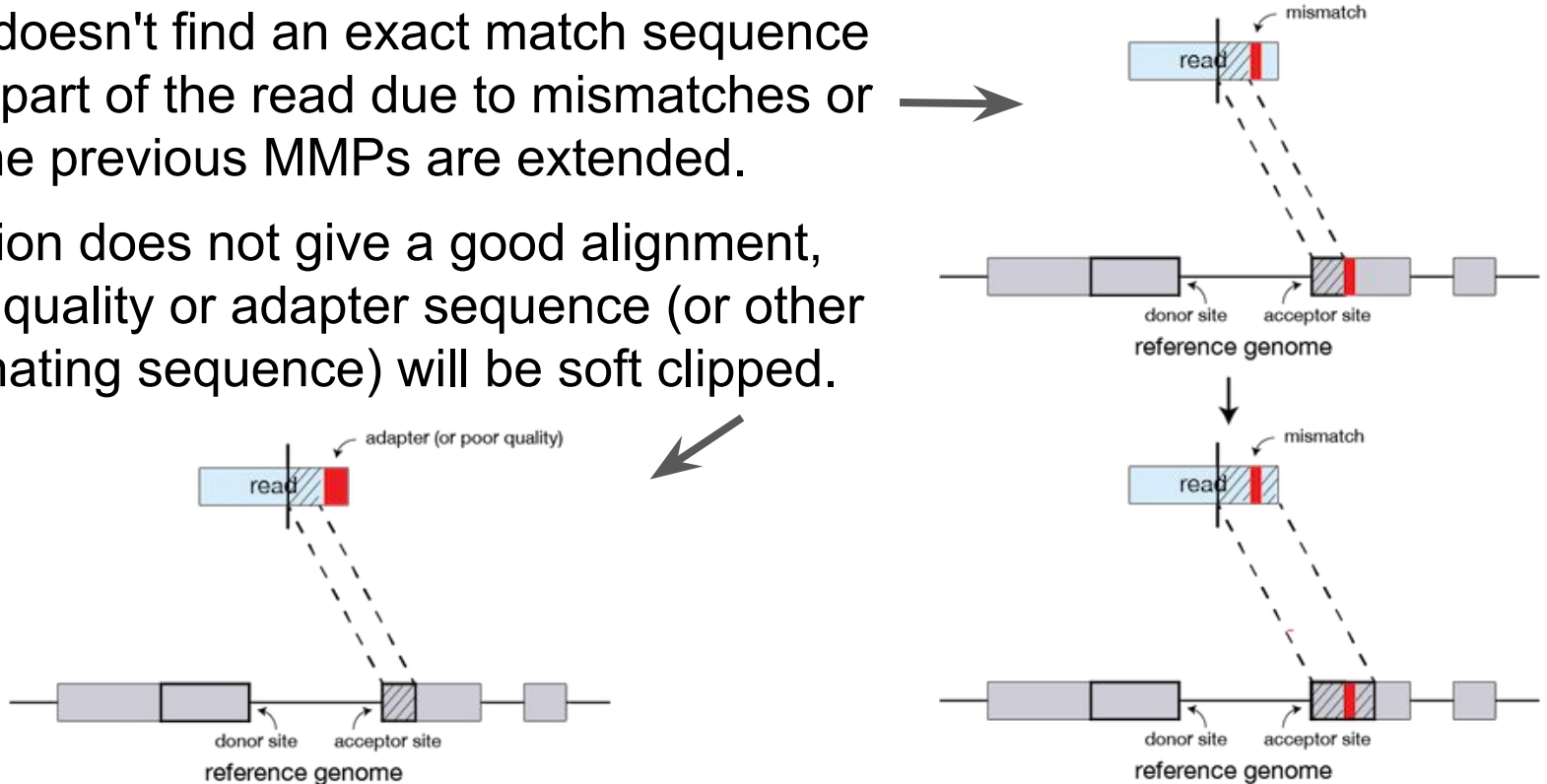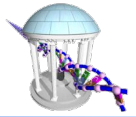STAR uses an uncompressed suffix array (SA) to search for MMPs.

# Handling Mismatches

If STAR doesn't find an exact match sequence for each part of the read due to mismatches or indels, the previous MMPs are extended.

If extension does not give a good alignment, the poor quality or adapter sequence (or other contaminating sequence) will be soft clipped.
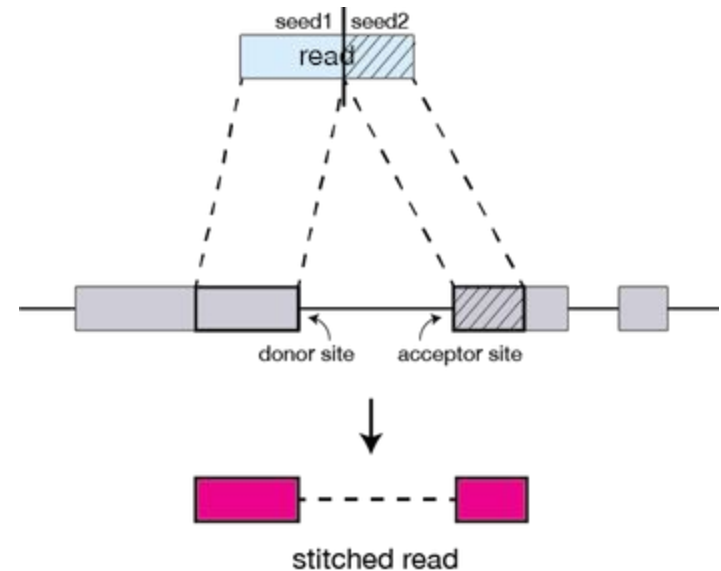
# Clustering, Stitching, and Scoring

Separate seeds are stitched together to create a complete read by first clustering the seeds together based on proximity to a set of 'anchor' seeds, or seeds that are not multi-mapping.

Then the seeds are stitched together based on the best alignment for the read (scoring based on mismatches, indels, gaps, etc.).

# Star Command-line Options

The basic options to **generate genome indices** using STAR are as follows:

- --genomeDir: /path/to/store/genome_indices
- --readFilesIn: /path/to/FASTQ_file
- --outFileNamePrefix: prefix for all output files
- --runThreadN: number of threads
- --outSAMtype: output filetype (SAM default)
- --outSAMunmapped: what to do with unmapped reads
- --sjdbOverhang: readlength -1

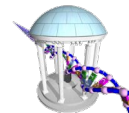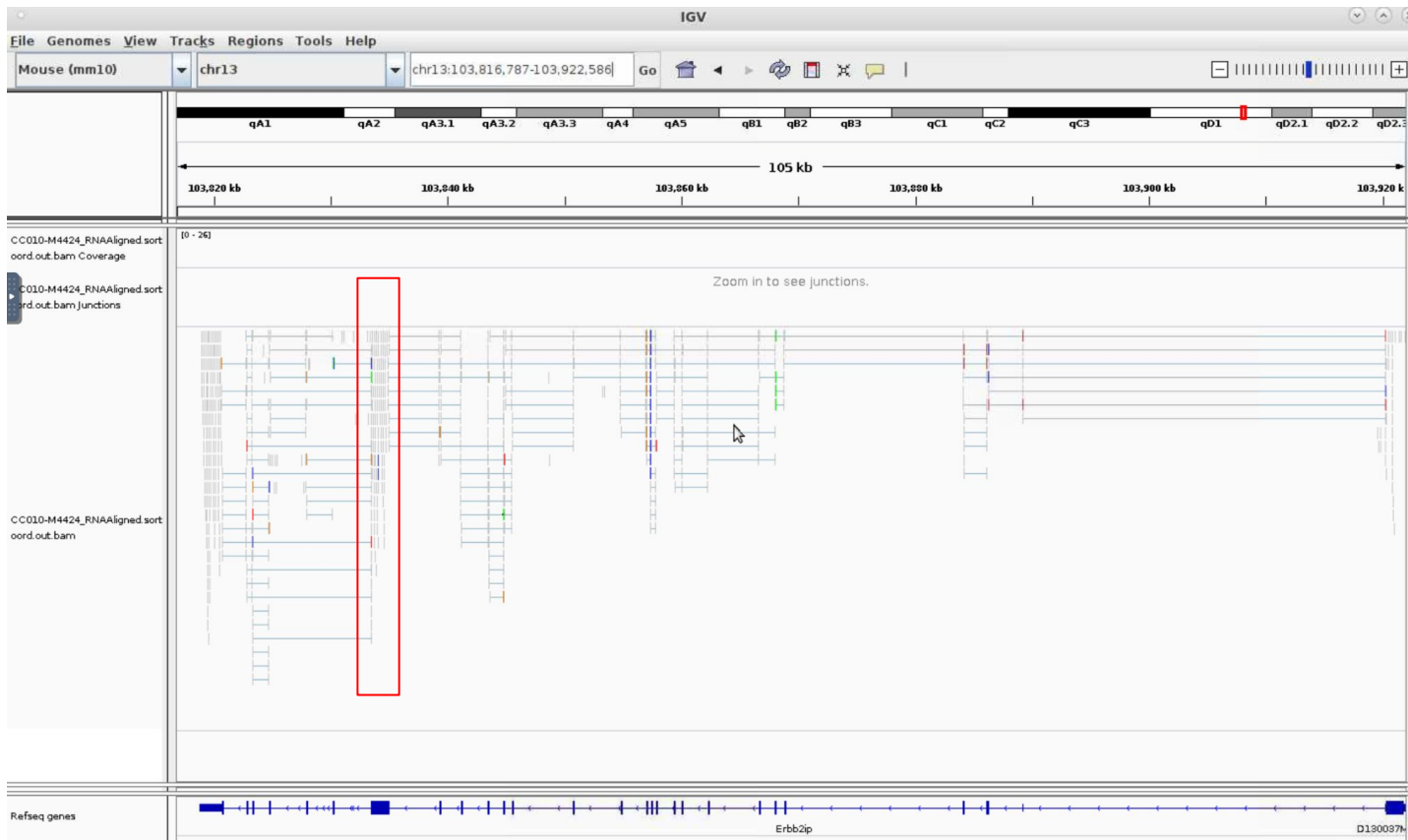# Star Sbatch Script

```
longleaf-login4$ cat bin/runStar
#!/bin/bash
#SBATCH --ntasks=6
#SBATCH --time=4:00:00
#SBATCH --mem=32G
reffile=$1
outfile="${!#}"
fqfiles="${@:2:$#-2}"
echo "ref genome:" $reffile
echo "input files:" $fqfiles
echo "output file:" $outfile
star --genomeDir $reffile --runThreadN 6 --readFilesIn $fqfiles --outFileNamePrefix
$outfile --outSAMtype BAM SortedByCoordinate --outSAMunmapped Within --outSAMattributes
Standard
longleaf-login4$ sbatch runStar /proj/seq/data/STAR_genomes_v277/GRCm38_p6_GENCODE_primary
/proj/mcmillanlab/BCB716F22/RNAseq/CC001_F5694_L001_R1.fastq
/proj/mcmillanlab/BCB716F22/RNAseq/CC001_F5694_L001_R2.fastq
/pine/scr/m/c/mcmillan/alignments/CC001_F5694_RNA
Submitted batch job 33793629
longleaf-login4$ samtools index
/pine/scr/m/c/mcmillan/alignments/CC001_F5694_RNAAligned.sortedByCoord.out.bam
```
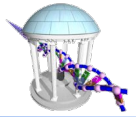
**Runs typically take less than 1 hour**
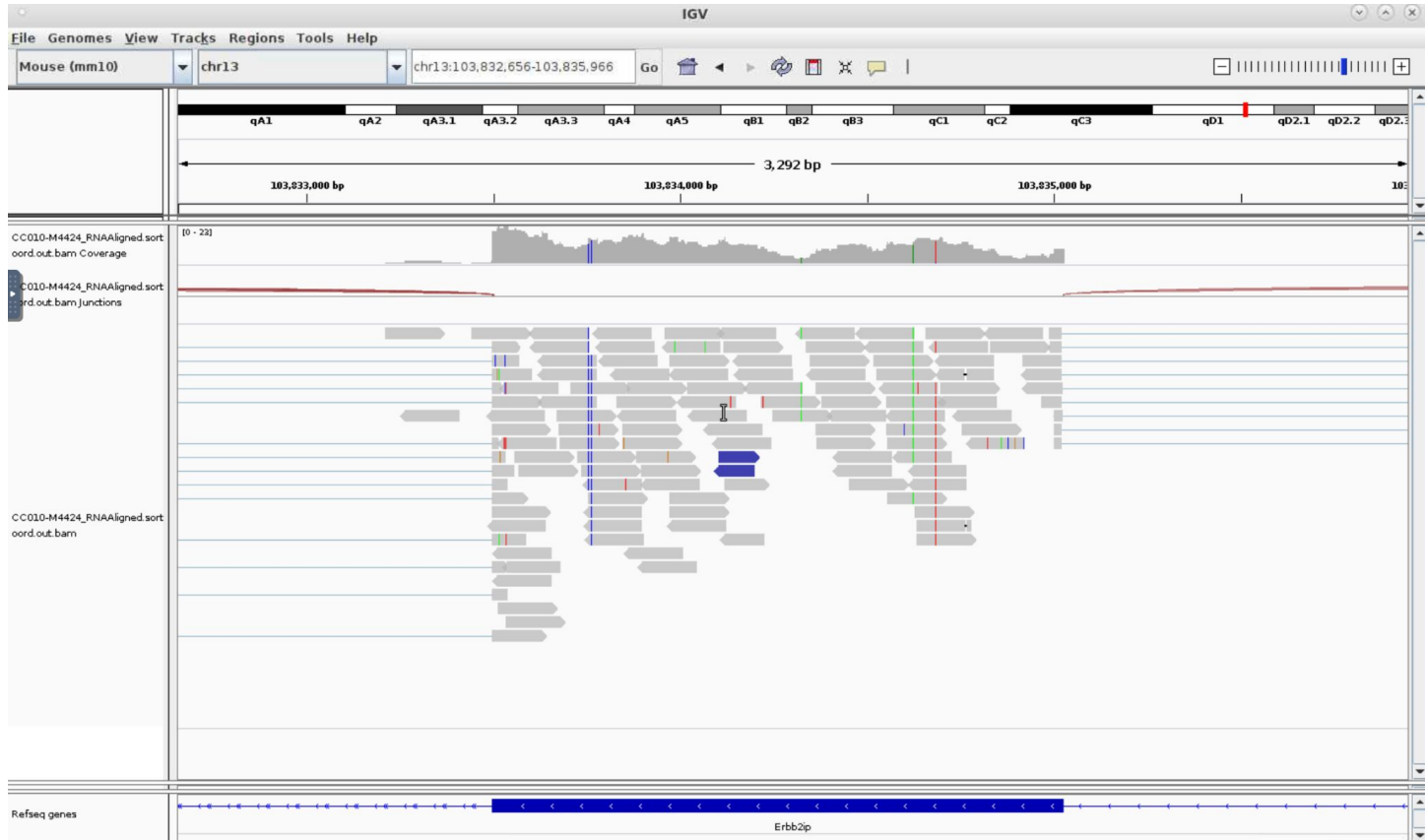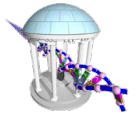
# IGV an RNAseq dataset
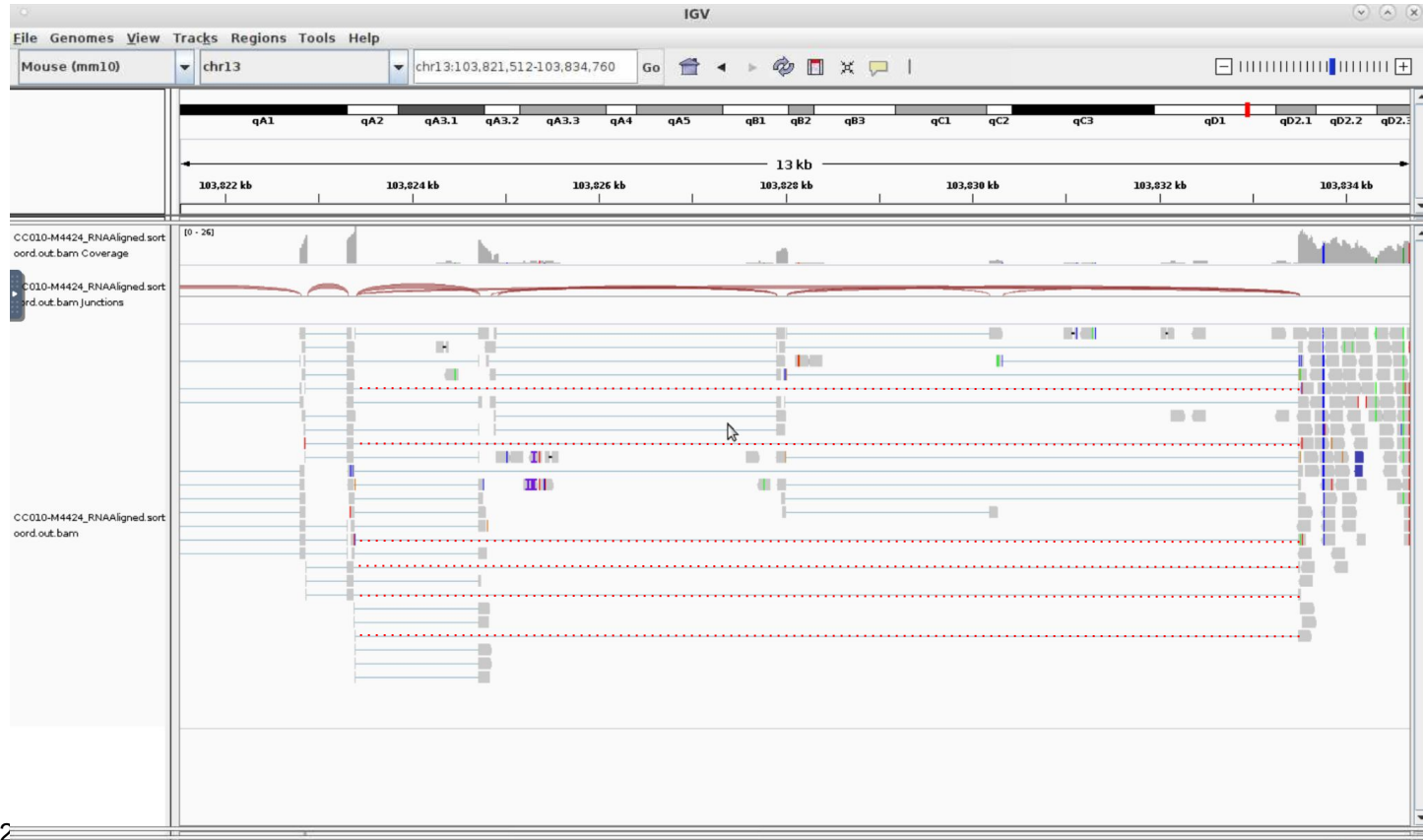
**An expressed gene**

# IGV an RNAseq dataset
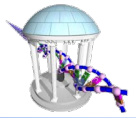
**4 variants
in an Exon**

# IGV an RNAseq dataset
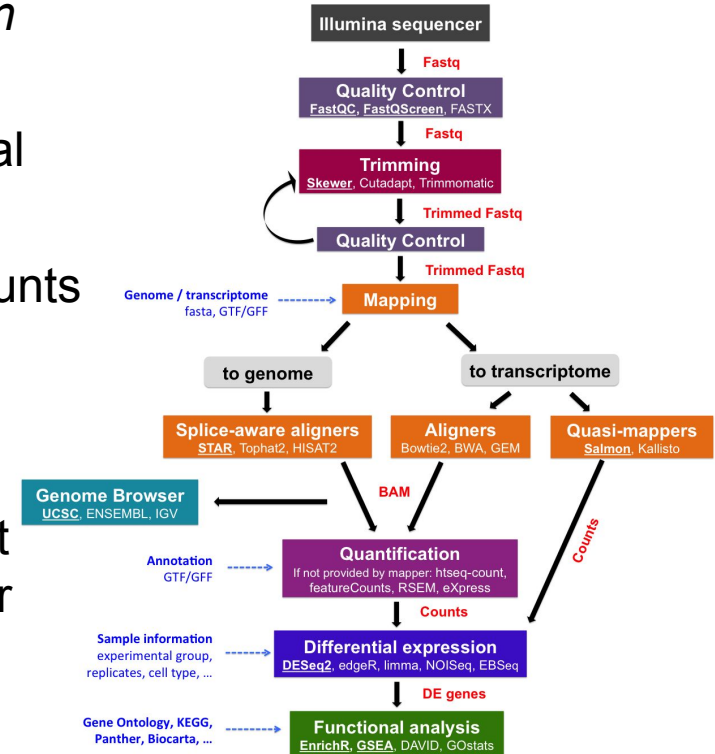
**Alternate splice versions**

# RNAseq quantification

RNA-Seq quantification estimates *gene expression levels* from RNAseq data.

Used for cross-sample comparisons and differential expression analysis

Results are reported in two forms raw mapping counts (reads that map to a gene) and normalized TPM (Transcripts per Million mapped reads) values.

Quantification can be done using alignment-based methods to either the genome or transcriptome but the trend is to quantify using the much faster k-mer based methods.
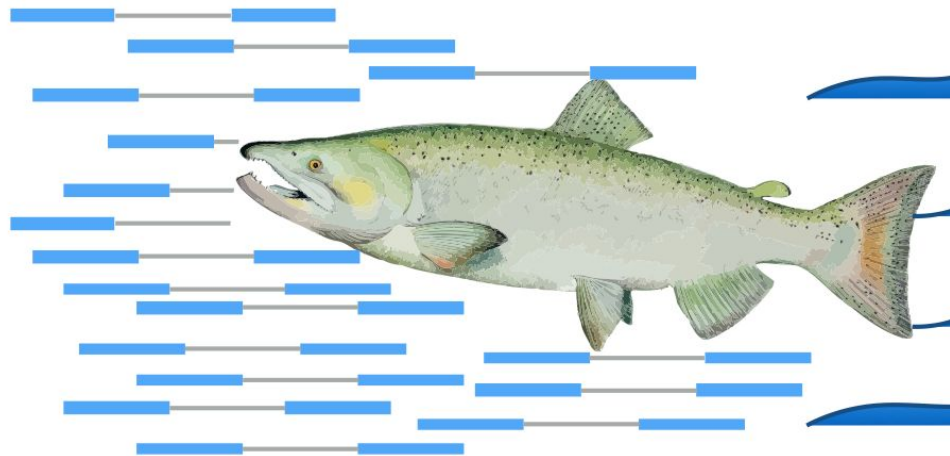
# Salmon

Salmon is a tool for quantifying the expression of transcripts using RNA-seq data.

It is a quasi-mapper as it doesn't produce the read alignments (and doesn't output BAM/SAM files).

Salmon "quasi-maps" reads to the *transcriptome* rather than to the *genome* as STAR does.

Salmon can make use of genome alignments (in the form of a SAM/BAM file) to the transcripts rather than using raw reads in the FASTQ format, this is often a good sanity check.

# Building a Salmon index

**To make an index for Salmon, we need transcript sequences in the FASTA format. Salmon does not need any decompression of the input, so we can index by using this command:**
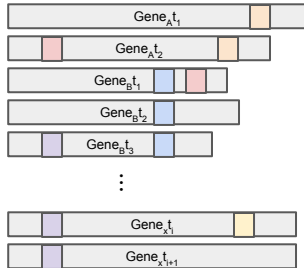
```
$ salmon index -t annotations/gencode.v29.transcripts.fa.gz \
        -i indexes/transcripts --gencode

Version Info: This is the most recent version of salmon.
index ["transcripts"] did not previously exist  . . . creating
it
[2019-04-30 18:12:59.272] [jLog] [info] building index
[2019-04-30 18:12:59.275] [jointLog] [info] [Step 1 of 4] :
counting k-mers

[....]
[2019-04-30 18:18:07.251] [jLog] [info] done building index
```

# Salmon's Index

Salmon requires an index for quazi-mapping that is derived from a FASTA file of annotated transcripts. This index tracks the association of *informative* k-mers to transcripts.
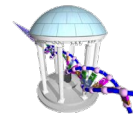
**Transcriptome**



**Index**



The index is used to associate reads with transcripts. Coordinates/offsets aren't important, thus "quasi-mapping"

# Salmon processing

1. Reads are scanned from left to right until a k-mer appears in the hash table
2. The k-mer is looked up in the hash table and the involved transcripts are added to a list
3. Scanning continues after the found k-mer
4. This process is repeated until the end of the read.
5. The final set of mappings is determined by a consensus mechanism, i.e. which of the gene/transcript pairs have the most support considering the transcript's orientation
6. A weighted count is added to a table for all potential transcripts consistent with the read.
7. Salmon corrects abundance estimates for any sample-specific biases/factors.
    a. Fragment CG bias
    b. Transcript-length corrections
    c. Fragment-length distributions
    d. Sequence specific biases
8. The weights are then redistributed using an Expectation Maximization (EM) approach.
   Recall if k-mer appeared in multiple transcripts it's count was distributed across all possibilities.
   However, if other k-mers are one of those transcripts are not seen, we can redistribute the partial counts to transcripts that are supported. This is done iteratively.
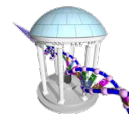
# Running Salmon

A slurm script for running Salmon, although it typically takes less than 10 mins

```
longleaf-login4$ cat bin/runSalmon
#!/bin/bash
#SBATCH --ntasks=6
#SBATCH --time=4:00:00
#SBATCH --mem=32G
reffile=$1
outfile="${!#}"
fqfiles="${@:2:$#-2}"
echo "ref genome:" $reffile
echo "input files:" $fqfiles
echo "output file:" $outfile
salmon quant -i $reffile --libType A -r $fqfiles -p 16 --validateMappings -o $outfile
longleaf-login4$ sbatch runSalmon /proj/mcmillanlab/BCB716F21/genome/Mouse/salmon
/proj/mcmillanlab/BCB716F22/RNAseq/CC011_M4692_L001_R*
/pine/scr/m/c/mcmillan/alignments/CC011_M4692_Salmon
Submitted batch job 33828974
```

# Salmon Results

# Salmon Verification

# Next Time

## Reference-free Sequence Analysis Approaches