# Analysis workflow



**Reads**
fastq

**Quality control and cleaning**
fastq                FastQC
                     trim_galore

**Reference genome**
fasta

**Genome indexation**
fasta,               Samtools faidx
fai,                 STAR genomeGenerate
star*

**Transcriptome de référence**
gtf / gff

**Transcriptome indexation**
rsem*                *rsem-prepare-reference*

**Alignment on reference**
                     STAR
bam, bai             samtools sort
                     samtools index

**Transcripts discovery**
                     cufflinks
                     cuffmerge
gtf                  samtools merge

**Quantification**
                     RSEM (rsem-
                     calculate-expression)
csv / txt            featureCount
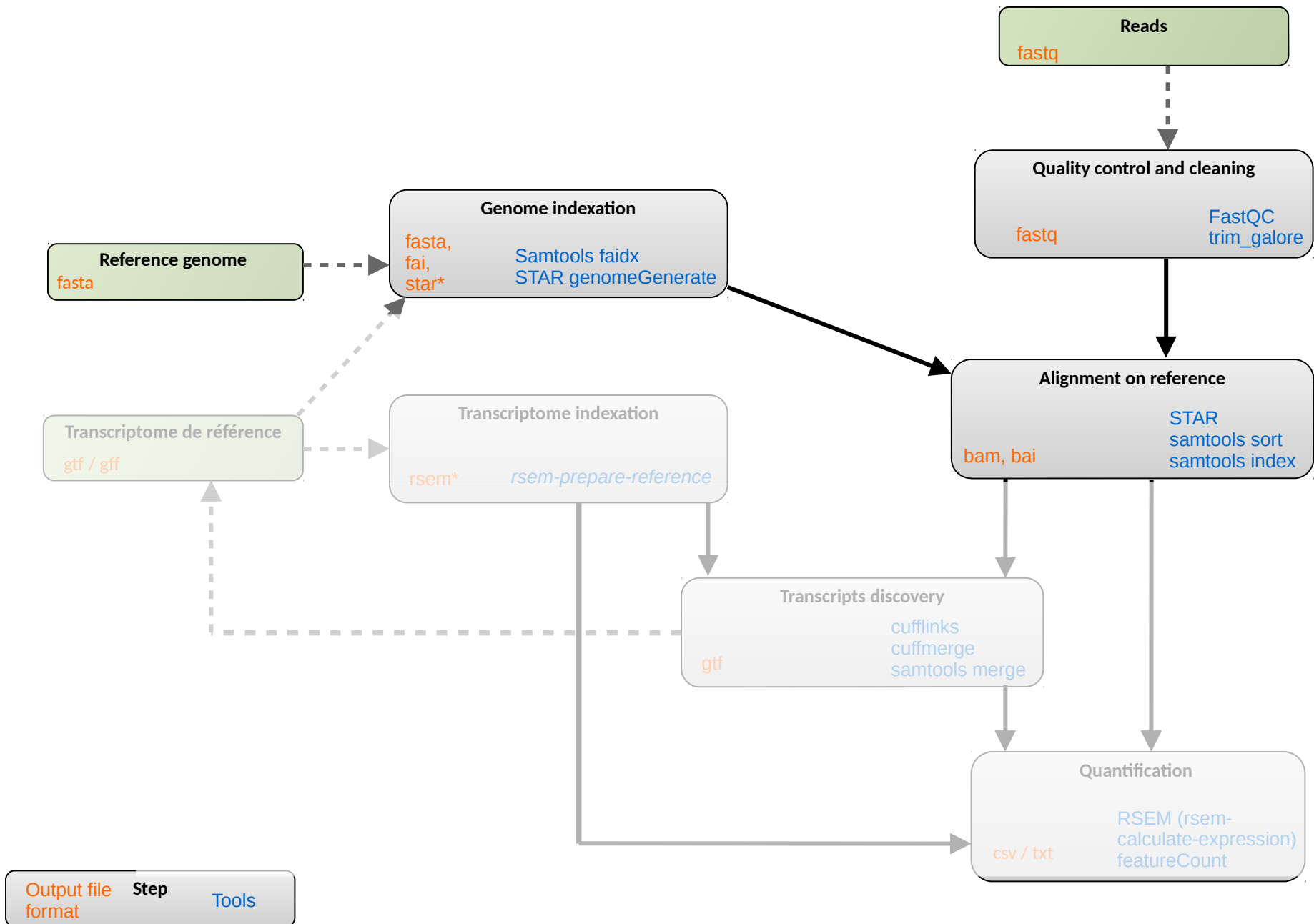
Output file     **Step**
format                      Tools

# Summary -

## Spliced read mapping & Visualisation

1. What is a spliced aligner?

2. Reference genome & transcriptome files formats

3. STAR principle and usage
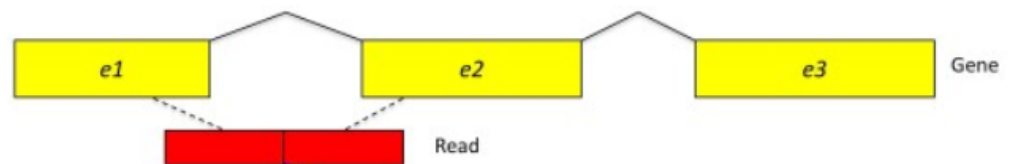
4. BAM & Bed files formats

5. Visualisation with IGV

# Aim -

## Spliced read mapping & Visualisation

**Aim**: Discover the true location (origin) of each read on the reference.

**Problems**:
- Some features (repetitive regions, assembly errors, missing information) make it impossible for some reads.
- Reads may be split by potentially thousands of bases of intronic sequence.



**And**:
Do it in/with reasonable time/resources.

# Splice sites

– Canonical splice site:
– which accounts for more than 99% of splicing
– GT and AG for donor and acceptor sites



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

– Non-canonical site:
– GC-AG splice site pairs, AT-AC pairs

Nucleic Acids Res. 2000 Nov 1;28(21):4364-75.

**Analysis of canonical and non-canonical splice sites in mammalian genomes.**

Burset M, Seledtsov IA, Solovyev VV.

– Trans-splicing:
splicing that joins two exons that are not within the same RNA transcript

# **Hard case**

- Lot of variations (sequencing errors, mutations)
- Repeats
- Reads spanning 3+ exons
- Gene or pseudogene



- Small end "anchor"

Kim et al, Genome Biology, 2013



- Unknown junction inside poorly rarely expressed gene

Mapping

# Most used tools

Tools for splice-mapping:

- ~~Tophat:~~
- HISAT

- STAR:

# **Benchmarks**

Run time:

Mapping

# Tuning parameters



on the human-T3-data base-level statistics.

« Therefore, an algorithm that is robust to parameter settings and exhibits good performance using defaults is desirable »

« most reliable general-purpose aligners appear to be CLC, Novoalign, GSNAP, and STAR. »

# rnaSTAR

## STAR: ultrafast universal RNA-seq aligner

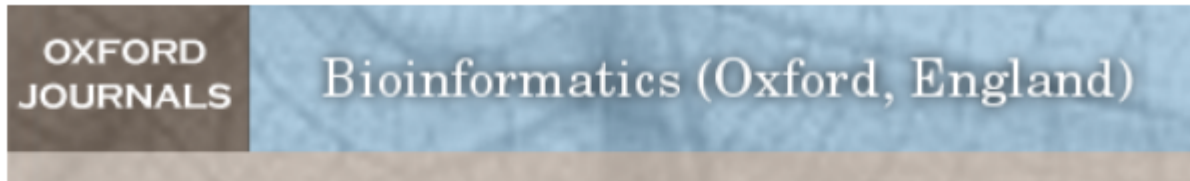Alexander Dobin,[1,*] Carrie A. Davis,[1] Felix Schlesinger,[1] Jorg Drenkow,[1] Chris Zaleski,[1] Sonali Jha,[1] Philippe Batut,[1] Mark Chaisson,[2] and Thomas R. Gingeras[1]

- Spliced Transcripts Alignment to a Reference
- Outperforms other aligners by more than a factor of 50 in mapping speed

## rnaSTAR strategy

- search for a MMP from the 1st base

- MMP search repeated for the unmapped portion next to the junction

- do it in both fwd and rev directions

- cluster seeds from the mates of paired-end RNA-seq reads

Soft-clipping is the main difference between Tophat and STAR



Dobin *et al*, Bioinformatics, 2011

# STAR : two passes strategy



« Improved ability to align reads by short spanning lengths is sufficient to explain the quantification benefit of two-pass alignment »

Veeneman et al, Bioinformatics, 2016

```
module load bioinfo/starXXX
```
**STAR --runMode genomeGenerate --genomeDir**
*genome_dir* **--genomeFastaFiles** *genome.fasta*

To use *N* CPUs, add: `--runThreadN` *N*
With an annotation: **--sjdbGTFfile** *annot.gtf*

Some precomputed indices are already available:
http://labshare.cshl.edu/shares/gingeraslab/www-data
/dobin/STAR/STARgenomes
or on your preferred platform: /bank/STARdb

# Where to find a reference genome?

Retrieving the genome file (fasta):

- The Genome Reference Consortium
http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/

- NCBI chromosome naming with « | » not well supported by mapping software

- Prefer EMBL
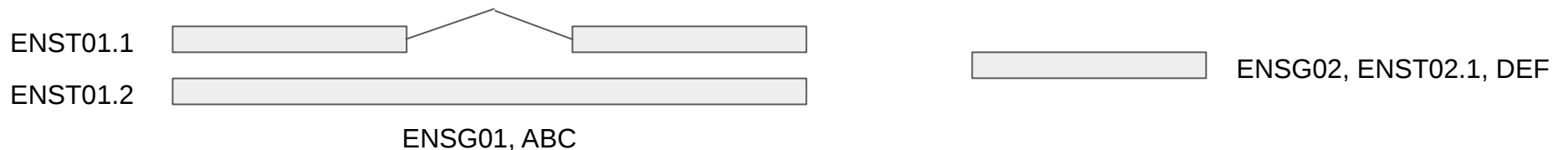http://www.ensembl.org/info/data/ftp/index.html

# Reference transcriptome file

What is a **GTF** file ?

- An annotation file: loci of coding genes (transcripts, CDS, UTRs), non-coding genes, etc.
- Gene Transfer Format (derived from GFF)

```
chr source   feature start end   score strand frame [attributes]
1   ENSEMBL exon    1000  2000 .     +      .      gene_id "ENSG01"; transcript_id "ENST01.1"; gene_name "ABC";
1   ENSEMBL exon    3000  4000 .     +      .      gene_id "ENSG01"; transcript_id "ENST01.1"; gene_name "ABC";
1   ENSEMBL exon    1000  4000 .     +      .      gene_id "ENSG01"; transcript_id "ENST01.2"; gene_name "ABC";
1   ENSEMBL exon    5000  6000 .     +      .      gene_id "ENSG02"; transcript_id "ENST02.1"; gene_name "DEF";
```

ENST01.1

ENST01.2

ENSG01, ABC

ENSG02, ENST02.1, DEF

- `gene_id` *value* : unique identifier for the gene.
- `transcript_id` *value* : unique identifier for the transcript.

**The chromosome names MUST be the same in the gtf file and fasta files (e.g. `chr1` vs `Chr1` vs `1`).**

http://genome.ucsc.edu/FAQ/FAQformat.html#format4

Mapping

**Exercise n°3**

Create a directory for the genome and annotation files.

Get the FASTA and GTF files from:
http://genoweb.toulouse.inra.fr/~formation/19_Rnaseq_Cli/data/reference/

Create the STAR index.
Tip: you can allocate *N* CPUs with the `sbatch -c 8`

```
module load bioinfo/starXXX
STAR --genomeDir genome_dir
--readFilesIn read1.fastq.gz read2.fastq.gz
--readFilesCommand zcat
--sjdbGTFfile transcriptome.gtf
--alignIntronMin 20 --alignIntronMax 500000
--outSAMtype BAM SortedByCoordinate    → sort
--outSAMstrandField intronMotif    → for cufflinks
--alignSoftClipAtReferenceEnds No    → for cufflinks
--outSAMattrIHstart 0    → for cufflinks or StringTie
--outFilterType BySJout    → filter by splice site
--outFilterIntronMotifs RemoveNoncanonical    → filter
--quantMode TranscriptomeSAM GeneCounts    → for RSEM
--outSAMattributes All    → more information
--outFileNamePrefix sampleName
--runThreadN 4
```

Intron size

```
 --alignIntronMin  20
 --alignIntronMax 500000
```

Allow soft-clipping past the end of chr (for cufflinks No)

`--alignSoftClipAtReferenceEnds No` [default Yes]

Output format:

`--outSAMtype BAM SortedByCoordinate [SAM]`

Output SAM/BAM alignments to transcriptome into a separate file (for RSEM)

```
--quantMode TranscriptomeSAM
    → need :--sjdbGTFfile annot.gtf
```

Output read unmapped

`--outReadsUnmapped Fastx`

# STAR options

Add more tags:

`--outSAMattributes All`

Default output file name: `Aligned.bam` Modify prefix:

`--outFileNamePrefix prefix`

Infer strand using intron motifs (for Cufflinks)

`--outSAMstrandField intronMotif [None]`

Start IH at `--outSAMattrIHstart 0 [1]` (for Cufflinks)

# STAR options

Remove reads that did not pass the junction filter:

`--outFilterType BySJOut [Normal]`

**Filter out alignments with non-canonical intron motifs**

`--outFilterIntronMotifs RemoveNoncanonical`

Mismatches :

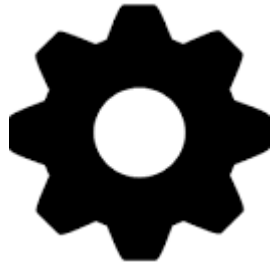`--outFilterMismatchNmax [default: 10]`

Limit multimap outputed:

`--outFilterMultimapNmax [Default: 10]`
> Flag of secondary alignment 0x100

Too short alignemnt

`--outFilterMatchNminOverLread 0.66`
`--outFilterScoreMinOverLread 0.66`

# STAR - two passes mode

- First pass: discover new junctions.

- Second pass: run again with knowing the new junctions.
(most useful for poorly annotated genomes.)

```
--twopassMode [None|Basic]
```

Defines the number of reads to be mapped in the 1st pass :
```
--twopass1readsN [-1]
```

# STAR Output files

Outputs (w/o specific options except `BAM SortedByCoordinate`):

- `Aligned.sortedByCoord.out.bam`: list of read alignments in SAM format compressed

- `Log.out`: main log file with a lot of detailed information about the run (for troubleshooting)

- `Log.progress.out`: reports job progress statistics

- `Log.final.out`: summary mapping statistics after mapping job is complete, very useful for quality control.

- `SJ.out.tab`: contains high confidence collapsed splice junctions in tab-delimited format

(chr, intron start, end, strand, intron motif, in database, # uniquely mapping reads, # multi, max. overhang)

Mapping

# STAR technical issues

- Temporary disk space:
  - Indexing the mouse genome requires 128GB and 1 hour on 6 slots.
  - Mapping a 16M paired-end reads requires 110GB and 4 mins on 6 slots.

- Available cluster:
  - New : 48 nodes with 32 cores and 256 GB of ram per node
  - Old : 68 nodes with 20 cores and 256 GB of ram per node

Mapping

**Exercise n°3**

Map the 2 FASTQ files.

*Do not forget to provide a different output file name for each set.*

Index the output BAM files with:

`samtools index file.bam`

→ Then BAM format presentation.

# SAM / BAM formats

Sequence Alignment/Map format:
- Each line stores an alignment/map

```
Coor     12345678901234  567890123456789012345678901234 5
ref      AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1        TTAGATAAAGGATA*CTG
+r002        aaaAGATAA*GGATA
+r003      gcctaAGCTAA
+r004                   ATAGCT..............TCAGC
-r003                            ttagctTAGGC
-r001/2                                     CAGCGGCAT

name flag chr start mapQ  cigar     nNext sNext tlen      seq           qual    tags
r001   99 ref   7    30 8M2I4M1D3M =     37    39 TTAGATAAAGGATACTG *
r002    0 ref   9    30 3S6M1P1I4M *      0     0 AAAAGATAAGGATA     *
r003    0 ref   9    30        5S6M *     0     0 GCCTAAGCTAA        * SA:Z:ref,29,-,6H5M,17,0;
r004    0 ref  16    30      6M14N5M *    0     0 ATAGCTTCAGC        *
r003 2064 ref  29    17        6H5M *     0     0 TAGGC              * SA:Z:ref,9,+,5S6M,30,1;
r001  147 ref  37    30          9M =     7   -39 CAGCGGCAT          * NM:i:1
```

- Header stores genome information
```
@HD VN:1.5 SO:coordinate
@SQ SN:ref LN:45
```

24                                                                   Mapping

```
Coor       12345678901234   5678901234567890123456789012345
ref        AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1         TTAGATAAAGGATA*CTG
+r002           aaaAGATAA*GGATA
+r003       gcctaAGCTAA
+r004                       ATAGCT..............TCAGC
-r003                             ttagctTAGGC
-r001/2                                       CAGCGGCAT

name flag chr start mapQ  cigar    nNext sNext tlen      seq           qual    tags
r001    99 ref    7     30 8M2I4M1D3M =      37       39 TTAGATAAAGGATACTG *
r002     0 ref    9     30 3S6M1P1I4M *       0        0 AAAAGATAAGGATA    *
r003     0 ref    9     30        5S6M *       0        0 GCCTAAGCTAA       * SA:Z:ref,29,-,6H5M,17,0;
r004     0 ref   16     30      6M14N5M *       0        0 ATAGCTTCAGC       *
r003  2064 ref   29     17        6H5M *       0        0 TAGGC             * SA:Z:ref,9,+,5S6M,30,1;
r001   147 ref   37     30          9M =       7      -39 CAGCGGCAT         * NM:i:1
```

- Flags: https://broadinstitute.github.io/picard/explain-flags.html
- MapQ: similar to a phred score
- nNext: =  means same chr
- In general, * means NA

# CIGAR

```
Coor        12345678901234  567890123456789012345678901234567
ref         AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1        TTAGATAAAGGATA*CTG
+r002        aaaAGATAA*GGATA
+r003     gcctaAGCTAA
+r004                     ATAGCT..............TCAGC
-r003                           ttagctTAGGC
-r001/2                                  CAGCGGCAT

name  flag chr start mapQ  cigar     nNext sNext tlen     seq            qual    tags
r001    99 ref    7    30 8M2I4M1D3M =      37    39 TTAGATAAAGGATACTG *
r002     0 ref    9    30 3S6M1P1I4M *       0     0 AAAAGATAAGGATA    *
r003     0 ref    9    30      5S6M *        0     0 GCCTAAGCTAA       * SA:Z:ref,29,-,6H5M,17,0;
r004     0 ref   16    30    6M14N5M *       0     0 ATAGCTTCAGC       *
r003  2064 ref   29    17      6H5M *        0     0 TAGGC             * SA:Z:ref,9,+,5S6M,30,1;
r001   147 ref   37    30        9M =        7   -39 CAGCGGCAT         * NM:i:1
```

- 30M means 30 matches or mismatches
- I and D : insertion/deletion
- S and H : soft/hard clipping

```
Coor        12345678901234   567890123456789012345678901234567890
ref         AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1        TTAGATAAAGGATA*CTG
+r002         aaaAGATAA*GGATA
+r003       gcctaAGCTAA
+r004                      ATAGCT..............TCAGC
-r003                            ttagctTAGGC
-r001/2                                   CAGCGGCAT

name  flag chr start mapQ  cigar     nNext sNext tlen    seq          qual    tags
r001    99 ref   7    30 8M2I4M1D3M =      37    39 TTAGATAAAGGATACTG *
r002     0 ref   9    30 3S6M1P1I4M *       0     0 AAAAGATAAGGATA    *
r003     0 ref   9    30       5S6M *       0     0 GCCTAAGCTAA        * SA:Z:ref,29,-,6H5M,17,0;
r004     0 ref  16    30     6M14N5M *      0     0 ATAGCTTCAGC        *
r003  2064 ref  29    17       6H5M *       0     0 TAGGC             * SA:Z:ref,9,+,5S6M,30,1;
r001   147 ref  37    30        9M =        7   -39 CAGCGGCAT         * NM:i:1
```

- Format: *2-letter name*:*format*:*value* (many different)
- NM: # mismatches
- SA: chimeric reads
- NH, HI: # hits for this sequence, hit index
- AS: alignment score
- nM: # mismatches per fragment

27    

BAM (Binary Alignment/Map) format:
- Compressed binary representation of SAM
- Greatly reduces storage space requirements to about 27% of original SAM
- samtools: reading, writing, and manipulating BAM files
- Most tools require a sorted and indexed BAM file.

- To be viewed a bam file must be indexed : `samtools index`

# samtools

```
Program: samtools (Tools for alignments in the SAM format)
Version: 1.8 (using htslib 1.8)

Usage:   samtools <command> [options]

Commands:
  -- Indexing
     dict           create a sequence dictionary file
     faidx          index/extract FASTA
     index          index alignment

  -- Editing
     calmd          recalculate MD/NM tags and '=' bases
     fixmate        fix mate information
     reheader       replace BAM header
     targetcut      cut fosmid regions (for fosmid pool only)
     addreplacerg   adds or replaces RG tags
     markdup        mark duplicates

  -- File operations
     collate        shuffle and group alignments by name
     cat            concatenate BAMs
     merge          merge sorted alignments
```

`module load bioinfo/samtools-1.8`

Bam → sam
`samtools view in.bam`
Sam → bam
`samtools view in.sam > out.bam`

Sort
`samtools sort -o out.bam in.bam`

Index
`samtools sort in.bam`

Global options nb threads:
`-@ 4`

# Visualizing alignments on IGV



http://www.broadinstitute.org/igv/home

Integrative genomics viewer

James T Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S Lander, Gad Getz & Jill P Mesirov

Affiliations | Corresponding authors

# Visualizing alignments on IGV

- High-performance visualization tool
- Interactive exploration of large, integrated datasets
- Supports a wide variety of data types
- Documentations
- Developed at the Broad Institute of MIT and Harvard

**File Formats**

- File Extension Identifies Format
- Recommended File Formats
- BAM
- BED
- CBS
- CN
- Cytoband
- FASTA
- GCT
- genePred
- GFF
- GISTIC
- HDF5
- IGV
- LOH
- Birdsuite Files
- MUT
- RES
- SAM
- Sample Information
- SEG
- SNP
- TAB
- TDF
- Track Line
- Type Line
- WIG

Mapping

# Visualizing alignments on IGV



Mapping

# IGV : Load reference genome



**Select a fasta file, the index .fai must exists in the same directory**

# IGV : Load annotation



**Go to** position or **gene** (enter gene name)

Load **GTF or GFF**, to get annotation track

Mapping

# IGV : Load alignment



Select a bam file, the index **.bai must exists** in the same directory

Mapping

Mapping

# Find library orientation

Color alignment by > first-of-pair strand

# Exercices 5