# Preprocessing genomic data

HackBio
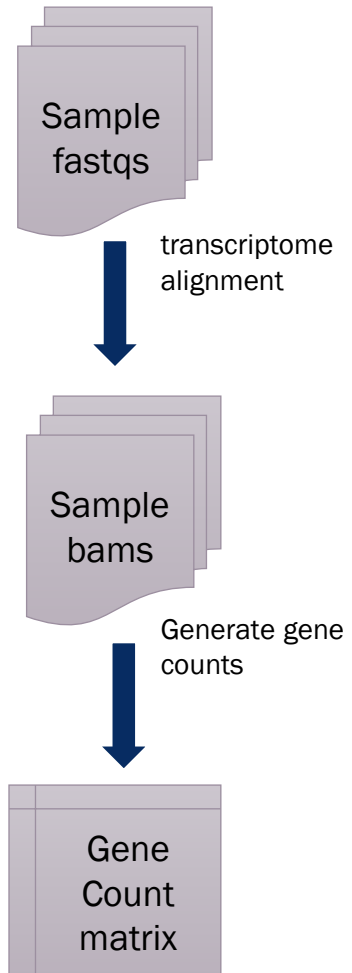
-----------------------------------------------------------------------

Melyssa Minto

West Lab, Duke Neurobiology

Computational Biology and Bioinformatics

HB

# Transcriptomics pre-proccesing workflow

Preprocessing



Sample fastqs

transcriptome alignment

Sample bams

Generate gene counts

Gene Count matrix

# Transcriptomics pre-proccesing workflow

**Preprocessing**

Sample fastqs

transcriptome alignment

Sample bams

Generate gene counts

Gene Count matrix
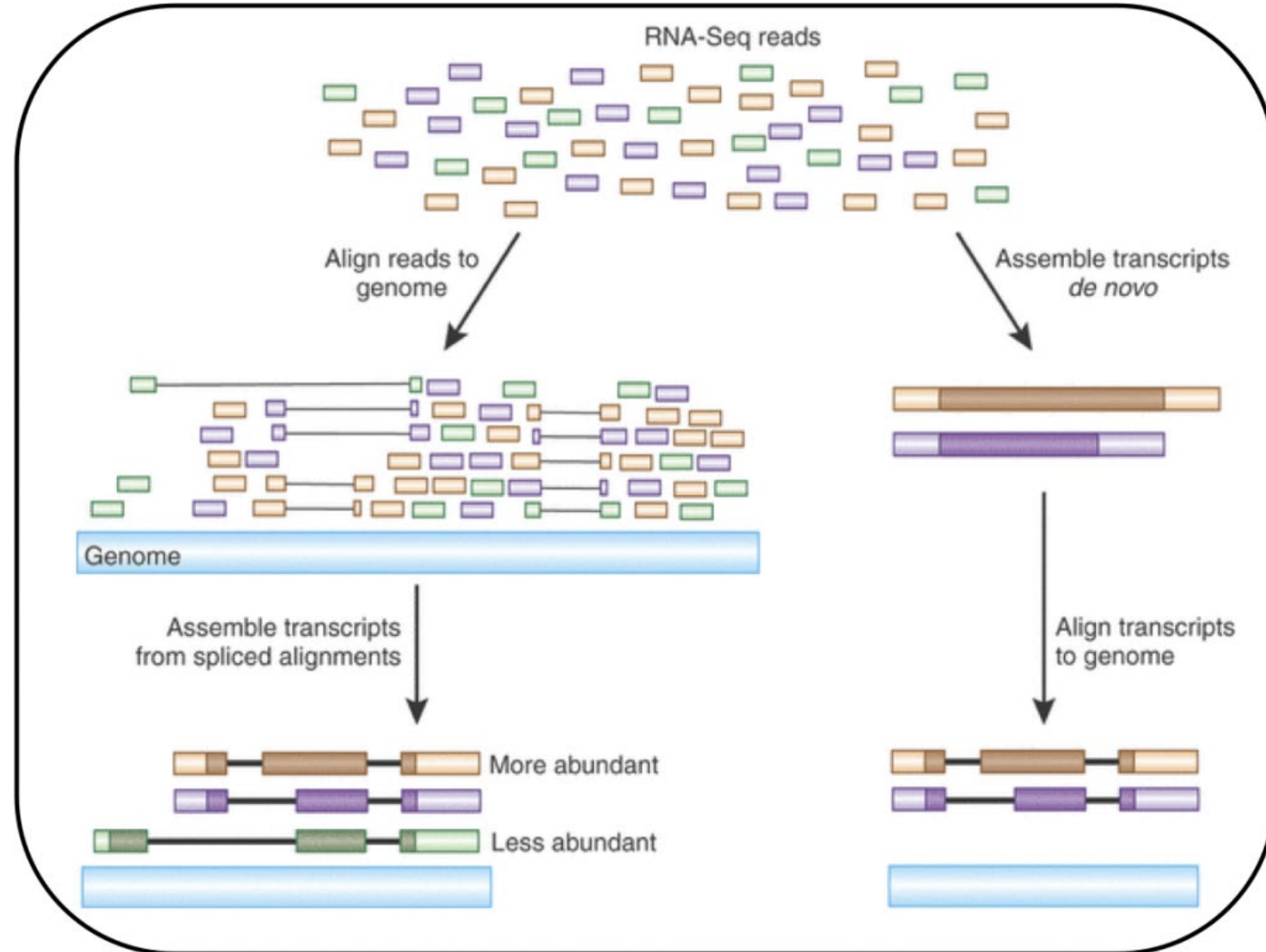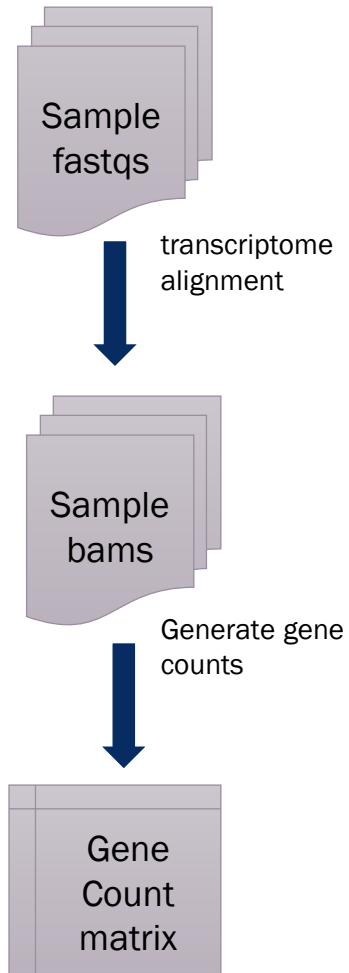
```
@A00257:355:HK7CTDRXX:1:2101:3522:1204 1:N:0:GACTACGA
@A00257:355:HK7CTDRXX:1:2101:3522:1204 1:N:0:GACTACGA
CNCTTGAATGCTGAGATTACAGATGTGCTCATAGACAACAGTAGCCACATC
@A00257:355:HK7CTDRXX:1:2101:3522:1204 1:N:0:GACTACGA
CNCTTGAATGCTGAGATTACAGATGTGCTCATAGACAACAGTAGCCACATC
+
F#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
@A00257:355:HK7CTDRXX:1:2101:3577:1204 1:N:0:GACTACGA
CNGGGAGAACCAGGTTAAAATTGAAGGTAGAAAACACTATAAGATGGAGGA
+
F#FFFFFFFFFFFFFF:FFFFFFFFFFFFFFF:FFFFFFFFFFFFFFFFFFF
@A00257:355:HK7CTDRXX:1:2101:3703:1204 1:N:0:GACTACGA
CNTATCCATATAAGAATTCAACAGAGAAACGGCAGGAAGACCCTTACCACT
+
F#FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
```
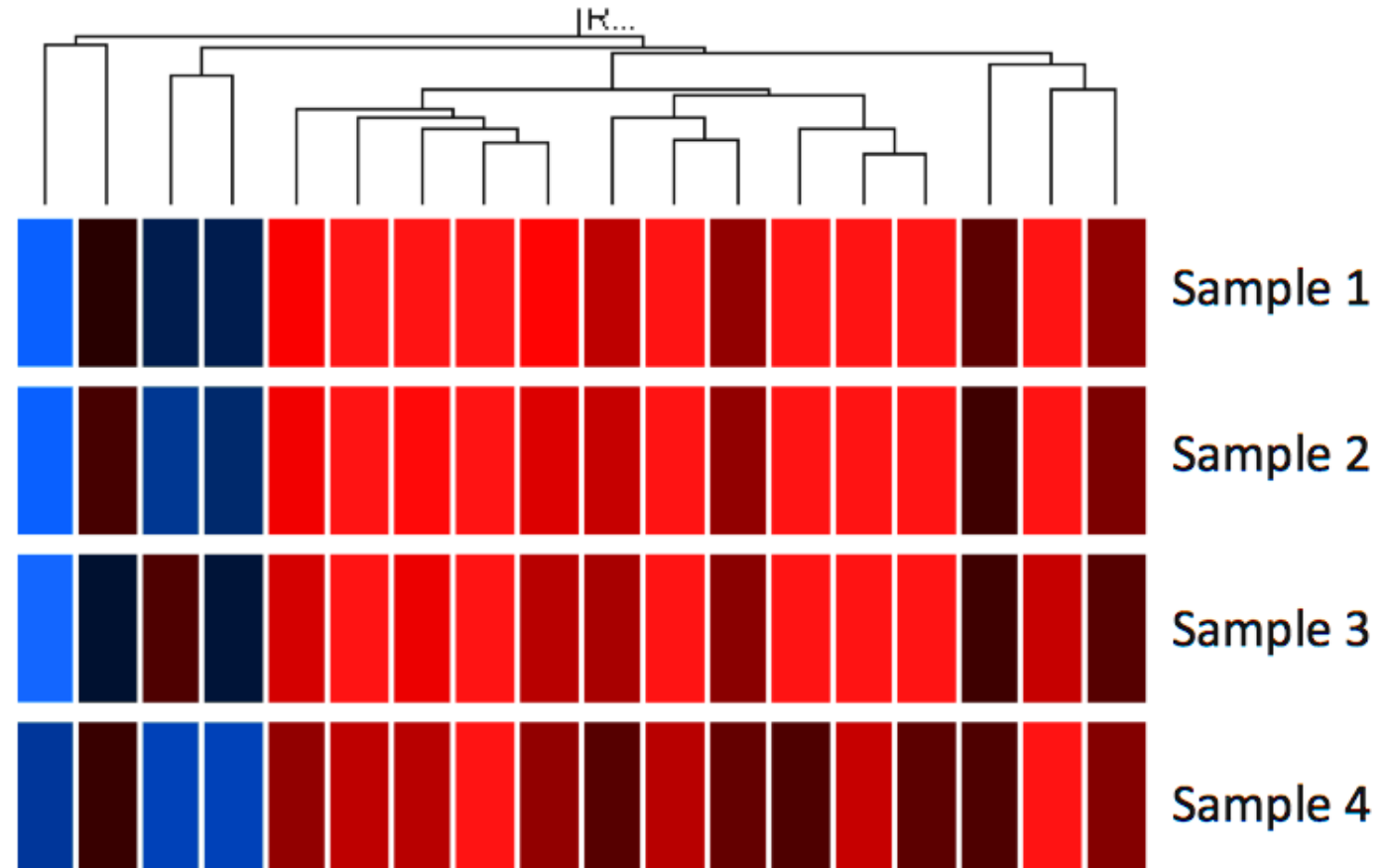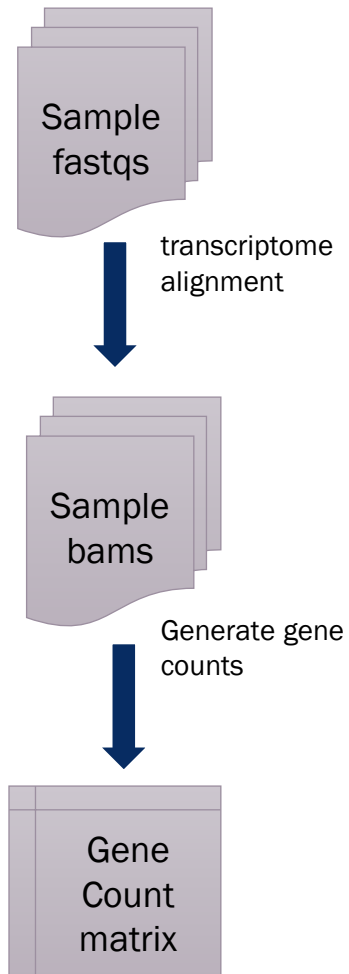
HB

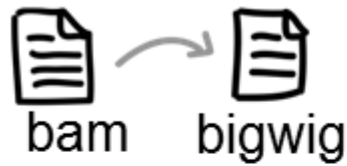# Transcriptomics pre-proccesing workflow

**Preprocessing**

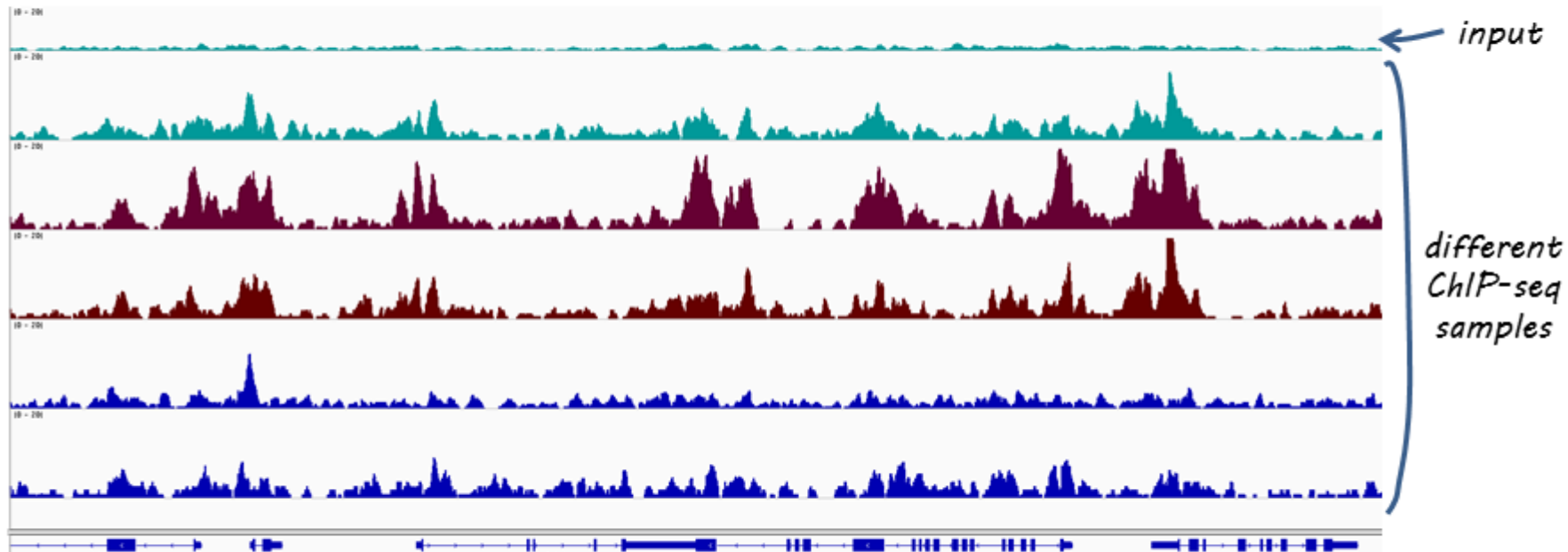# Transcriptomics pre-proccesing workflow

**Preprocessing**



Sample fastqs

transcriptome alignment

Sample bams

Generate gene counts

Gene Count matrix

Sample 1

Sample 2

Sample 3

Sample 4

# Visualizing aligned genomic data



bam → bigwig

for visualizing continuous data, e.g. in the UCSC Genome Browser or IGV, bigWig files come in really handy

input

different ChIP-seq samples

remember that there are 2 deepTools for bam → bigWig conversion:
- ❖ **bamCoverage**: for individual files (like those shown here)
- ❖ **bamCompare**: to normalize two files to each other

# Quality Control of Alignments

- Mapping logs
- samtools flagstat <bam>
- Read depth

# Extracting reads – what is a gtf file?

# Extracting reads – what is a gtf file?

1. **Sequence Name**
2. Source of Annotation
3. Feature
4. Start
5. End
6. Score
7. Strand
8. Frame
9. Attribute

# Extracting reads – what is a gtf file?

1. Sequence Name
2. **Source of Annotation**
3. Feature
4. Start
5. End
6. Score
7. Strand
8. Frame
9. Attribute

# Extracting reads – what is a gtf file?

1. Sequence Name
2. Source of Annotation
3. **Feature**
4. Start
5. End
6. Score
7. Strand
8. Frame
9. Attribute

```
##description: evidence-based annotation of the mouse genome (GRCm38), version M21 (Ensembl 96)
##provider: GENCODE
##contact: gencode-help@ebi.ac.uk
##format: gtf
##date: 2019-03-27
chr1    HAVANA  gene    3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; gene_type "TEC"; gene
chr1    HAVANA  transcript      3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; transcript_id
3401J01Rik-201"; level 2; transcript_support_level "NA"; tag "basic"; havana_gene "OTTMUSG00000049935.1"; havana_tran
chr1    HAVANA  exon    3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; transcript_id "ENSMUS
ik-201"; exon_number 1; exon_id "ENSMUSE00001343744.1"; level 2; transcript_support_level "NA"; tag "basic"; havana_g
chr1    ENSEMBL gene    3102016 3102125 .       +       .       gene_id "ENSMUSG00000064842.1"; gene_type "snRNA"; ge
chr1    ENSEMBL transcript      3102016 3102125 .       +       .       gene_id "ENSMUSG00000064842.1"; transcript_id
```

HB

# Extracting reads – what is a gtf file?

1. Sequence Name
2. Source of Annotation
3. Feature
4. **Start**
5. End
6. Score
7. Strand
8. Frame
9. Attribute

```
##description: evidence-based annotation of the mouse genome (GRCm38), version M21 (Ensembl 96)
##provider: GENCODE
##contact: gencode-help@ebi.ac.uk
##format: gtf
##date: 2019-03-27
chr1    HAVANA  gene    3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; gene_type "TEC"; gene
chr1    HAVANA  transcript      3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; transcript_id
3401J01Rik-201"; level 2; transcript_support_level "NA"; tag "basic"; havana_gene "OTTMUSG00000049935.1"; havana_tran
chr1    HAVANA  exon    3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; transcript_id "ENSMUS
ik-201"; exon_number 1; exon_id "ENSMUSE00001343744.1"; level 2; transcript_support_level "NA"; tag "basic"; havana_g
chr1    ENSEMBL gene    3102016 3102125 .       +       .       gene_id "ENSMUSG00000064842.1"; gene_type "snRNA"; ge
chr1    ENSEMBL transcript      3102016 3102125 .       +       .       gene_id "ENSMUSG00000064842.1"; transcript_id
```

# Extracting reads – what is a gtf file?

1. Sequence Name
2. Source of Annotation
3. Feature
4. Start
5. **End**
6. Score
7. Strand
8. Frame
9. Attribute

# Extracting reads – what is a gtf file?

1. Sequence Name
2. Source of Annotation
3. Feature
4. Start
5. End
6. **Score**
7. Strand
8. Frame
9. Attribute

```
##description: evidence-based annotation of the mouse genome (GRCm38), version M21 (Ensembl 96)
##provider: GENCODE
##contact: gencode-help@ebi.ac.uk
##format: gtf
##date: 2019-03-27
chr1    HAVANA  gene    3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; gene_type "TEC"; gene
chr1    HAVANA  transcript      3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; transcript_id
3401J01Rik-201"; level 2; transcript_support_level "NA"; tag "basic"; havana_gene "OTTMUSG00000049935.1"; havana_tran
chr1    HAVANA  exon    3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; transcript_id "ENSMUS
ik-201"; exon_number 1; exon_id "ENSMUSE00001343744.1"; level 2; transcript_support_level "NA"; tag "basic"; havana_g
chr1    ENSEMBL gene    3102016 3102125 .       +       .       gene_id "ENSMUSG00000064842.1"; gene_type "snRNA"; ge
chr1    ENSEMBL transcript      3102016 3102125 .       +       .       gene_id "ENSMUSG00000064842.1"; transcript_id
```

# Extracting reads – what is a gtf file?

1. Sequence Name
2. Source of Annotation
3. Feature
4. Start
5. End
6. Score
7. **Strand**
8. Frame
9. Attribute

# Extracting reads – what is a gtf file?

1. Sequence Name
2. Source of Annotation
3. Feature
4. Start
5. End
6. Score
7. Strand
8. **Frame**
9. Attribute

# Extracting reads – what is a gtf file?

1. Sequence Name
2. Source of Annotation
3. Feature
4. Start
5. End
6. Score
7. Strand
8. Frame
9. **Attribute**

```
##description: evidence-based annotation of the mouse genome (GRCm38), version M21 (Ensembl 96)
##provider: GENCODE
##contact: gencode-help@ebi.ac.uk
##format: gtf
##date: 2019-03-27
chr1    HAVANA  gene    3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; gene_type "TEC"; gene
chr1    HAVANA  transcript      3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; transcript_id
3401J01Rik-201"; level 2; transcript_support_level "NA"; tag "basic"; havana_gene "OTTMUSG00000049935.1"; havana_tran
chr1    HAVANA  exon    3073253 3074322 .       +       .       gene_id "ENSMUSG00000102693.1"; transcript_id "ENSMUS
ik-201"; exon_number 1; exon_id "ENSMUSE00001343744.1"; level 2; transcript_support_level "NA"; tag "basic"; havana_g
chr1    ENSEMBL gene    3102016 3102125 .       +       .       gene_id "ENSMUSG00000064842.1"; gene_type "snRNA"; ge
chr1    ENSEMBL transcript      3102016 3102125 .       +       .       gene_id "ENSMUSG00000064842.1"; transcript_id
```

# Extracting reads – what is a gtf file?

# Extracting reads – different tools

- [Subreads feature count](#)

- [Htseq](#)

- [RSEM](#)

Summarize a BAM format dataset:

*featureCounts -t exon -g gene_id -a annotation.gtf -o counts.txt mapping_results_SE.bam*

Summarize multiple datasets at the same time:

*featureCounts -t exon -g gene_id -a annotation.gtf -o counts.txt library1.bam library2.bam library3.bam*

Perform strand-specific read counting (use '-s 2' if reversely stranded):

*featureCounts -s 1 -t exon -g gene_id -a annotation.gtf -o counts.txt mapping_results_SE.bam*

Summarize paired-end reads and count fragments (instead of reads):

*featureCounts -p -t exon -g gene_id -a annotation.gtf -o counts.txt mapping_results_PE.bam*

Summarize multiple paired-end datasets:

*featureCounts -p -t exon -g gene_id -a annotation.gtf -o counts.txt library1.bam library2.bam library3.bam*

# Extracting reads – different tools

- [Subreads feature count](#)
- [Htseq](#)

```
htseq-count [options] <alignment_files> <gff_file>
```

- [RSEM](#)

**-f** `<format>`, **--format**=`<format>`
    Format of the input data. Possible values are `sam` (for text SAM files) and `bam` (for binary BAM files). Default is `sam`.

**-r** `<order>`, **--order**=`<order>`
    For paired-end data, the alignment have to be sorted either by read name or by alignment position. If your data is not sorted, use the `samtools sort` function of `samtools` to sort it. Use this option, with `name` or `pos` for `<order>` to indicate how the input data has been sorted. The default is `name`.

    If `name` is indicated, `htseq-count` expects all the alignments for the reads of a given read pair to appear in adjacent records in the input data. For `pos`, this is not expected; rather, read alignments whose mate alignment have not yet been seen are kept in a buffer in memory until the mate is found. While, strictly speaking, the latter will also work with unsorted data, sorting ensures that most alignment mates appear close to each other in the data and hence the buffer is much less likely to overflow.

**--max-reads-in-buffer**=`<number>`
    When <alignment_file> is paired end sorted by position, allow only so many reads to stay in memory until the mates are found (raising this number will use more memory). Has no effect for single end or paired end sorted by name. (default: `30000000`)

**-s** `<yes/no/reverse>`, **--stranded**=`<yes/no/reverse>`
    whether the data is from a strand-specific assay (default: `yes`)

**-m** `<mode>`, **--mode**=`<mode>`
    Mode to handle reads overlapping more than one feature. Possible values for *<mode>* are `union`, `intersection-strict` and `intersection-nonempty` (default: `union`)

# Extracting reads – different tools

- [Subreads feature count](#)
- [Htseq](#)
- [RSEM](#)

`htseq-count [options] <alignment_files> <gff_file>`

-**f** <format>, --**format**=<format>
    Format of the input data. Possible values are `sam` (for text SAM files) and `bam` (for binary BAM files). Default is `sam`.

-**r** <order>, --**order**=<order>
    For paired-end data, the alignment have to be sorted either by read name or by alignment position. If your data is not sorted, use the `samtools sort` function of `samtools` to sort it. Use this option, with `name` or `pos` for <order> to indicate how the input data has been sorted. The default is `name`.

    If `name` is indicated, `htseq-count` expects all the alignments for the reads of a given read pair to appear in adjacent records in the input data. For `pos`, this is not expected; rather, read alignments whose mate alignment have not yet been seen are kept in a buffer in memory until the mate is found. While, strictly speaking, the latter will also work with unsorted data, sorting ensures that most alignment mates appear close to each other in the data and hence the buffer is much less likely to overflow.

--**max-reads-in-buffer**=<number>
    When <alignment_file> is paired end sorted by position, allow only so many reads to stay in memory until the mates are found (raising this number will use more memory). Has no effect for single end or paired end sorted by name. (default: `30000000`)

-**s** <yes/no/reverse>, --**stranded**=<yes/no/reverse>
    whether the data is from a strand-specific assay (default: `yes`)

-**m** <mode>, --**mode**=<mode>
    Mode to handle reads overlapping more than one feature. Possible values for <mode> are `union`, `intersection-strict` and `intersection-nonempty` (default: `union`)

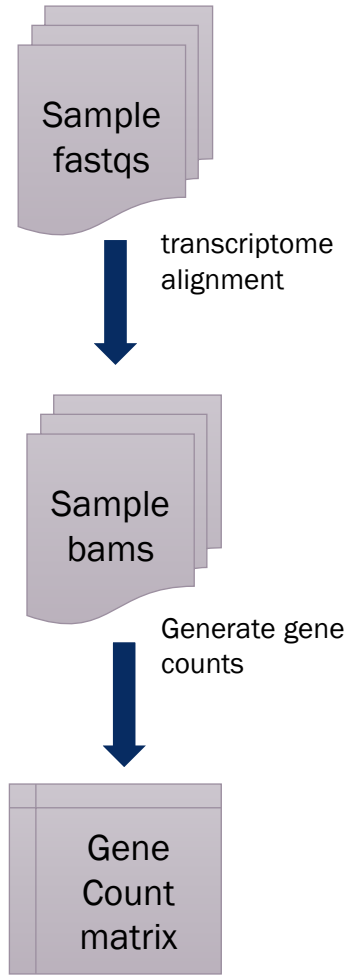| | union | intersection_strict | intersection_nonempty |
|---|---|---|---|
| read / gene_A | gene_A | gene_A | gene_A |
| read / gene_A | gene_A | no_feature | gene_A |
| read / gene_A gene_A | gene_A | no_feature | gene_A |
| read read / gene_A gene_A | gene_A | gene_A | gene_A |
| read / gene_A / gene_B | gene_A | gene_A | gene_A |
| read / gene_A / gene_B | ambiguous (both genes with --nonunique all) | gene_A | gene_A |
| read / gene_A / gene_B | ambiguous (both genes with --nonunique all) | | |
| read / gene_A gene_B | alignment_not_unique (both genes with --nonunique all) | | |

# Extracting reads – different tools

- Subreads feature count
- Htseq
- RSEM

```
software/RSEM-1.2.25/rsem-calculate-expression -p 8 --paired-end \
                                --bam \
                                --estimate-rspd \
                                --append-names \
                                --output-genome-bam \
                                exp/LPS_6h.bam \
                                ref/mouse_ref exp/LPS_6h
```
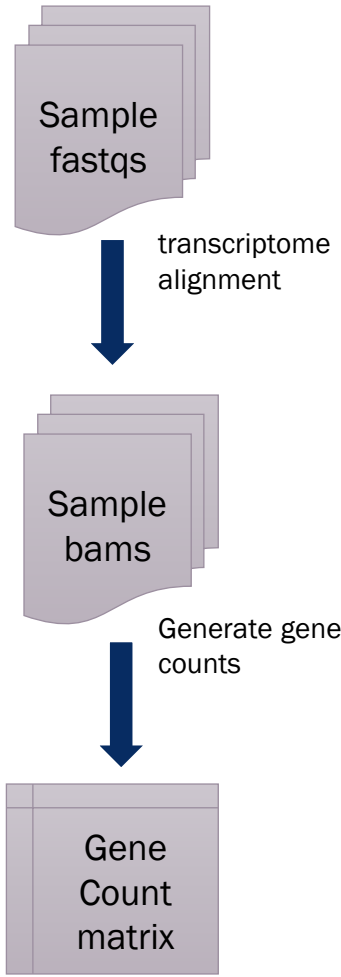
# Transcriptomics pipeline/workflow

Preprocessing

Sample
fastqs

transcriptome
alignment

Sample
bams

Generate gene
counts

Gene
Count
matrix

HB

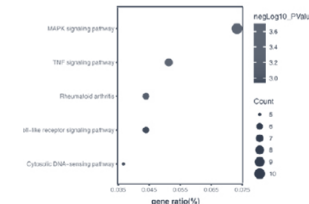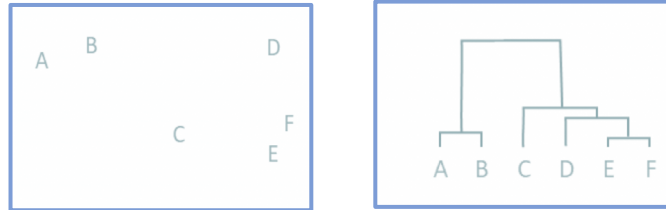# Transcriptomics pipeline/workflow

**Preprocessing**

**Analyses**

Sample fastqs

*transcriptome alignment*

Sample bams

*Generate gene counts*

Gene Count matrix

Clustering

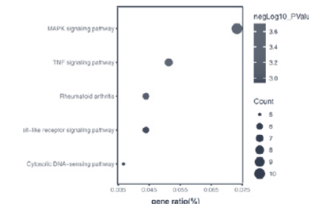**HB**

# Transcriptomics pipeline/workflow

**Preprocessing**

**Analyses**

# Transcriptomics pipeline/workflow

# Transcriptomics pipeline/workflow