



LEIDEN UNIVERSITY MEDICAL CENTER

RNA-seq using Galaxy

Jeroen F. J. Laros

Leiden Genome Technology Center

Department of Human Genetics

Center for Human and Clinical Genetics



Sequencers: HiSeq



Figure 1 : HiSeq 2000.

Characteristics:

- High throughput.
- Paired end.
- High accuracy.
- Read length $2 \times 150\text{bp}$.
- Relatively long run time.
- Relatively expensive.

Sequencers: Ion Torrent



Figure 2 : Ion torrent.

Characteristics:

- Moderate throughput.
- Single end (for now).
- High accuracy.
- Read length ± 200 bp.
- Short run time.
- Cheap runs.

General layout of an RNA-seq pipeline.

1. Pre-alignment.
 - QC.
 - Data cleaning.

General layout of an RNA-seq pipeline.

1. Pre-alignment.
 - QC.
 - Data cleaning.
2. Alignment.
 - Use a specialised (RNA) aligner.

General layout of an RNA-seq pipeline.

1. Pre-alignment.
 - QC.
 - Data cleaning.
2. Alignment.
 - Use a specialised (RNA) aligner.
3. Expression (gene, transcripts) analysis.
 - Known transcripts.

General layout of an RNA-seq pipeline.

1. Pre-alignment.
 - QC.
 - Data cleaning.
2. Alignment.
 - Use a specialised (RNA) aligner.
3. Expression (gene, transcripts) analysis.
 - Known transcripts.
4. Transcript assembly.
 - New transcripts, alternative splicing, etc.

Sickle / FastQC.

We use the Sickle for data cleaning.

For adapter clipping, we use Trimmomatic or the FastX toolkit (not in this practical).

- Remove linker sequences.
- Clip low quality reads at the end of the read.
- Evaluate the part of the read that is left.

Sickle / FastQC.

We use the Sickle for data cleaning.

For adapter clipping, we use Trimmomatic or the FastX toolkit (not in this practical).

- Remove linker sequences.
- Clip low quality reads at the end of the read.
- Evaluate the part of the read that is left.

The FastQC tool kit is used for quality control (both before and after the data cleaning step).

- GC content.
- GC distribution.
- Quality scores distribution.
- ...

FastQC report.

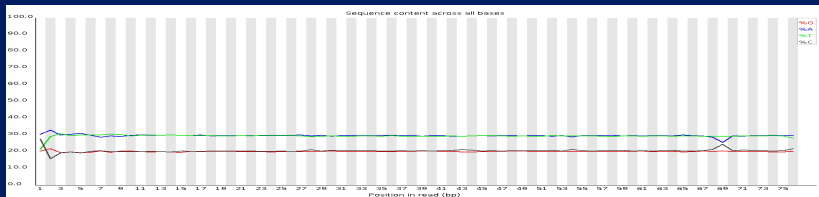


Figure 3 : Per base sequence content.

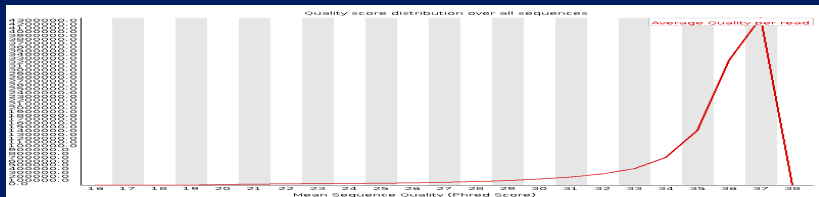


Figure 4 : Per sequence quality.

RNA aligners.

Difference with DNA:

- Splicing.

RNA aligners.

Difference with DNA:

- Splicing.

This affects:

- Insert sizes.
- Mapping of reads that cover an exon-exon boundary.

RNA aligners.

Difference with DNA:

- Splicing.

This affects:

- Insert sizes.
- Mapping of reads that cover an exon-exon boundary.

Available tools:

- Tophat.
- Gmap / Gsnap.
- PASSion.
- MapSplice.
- HMMSplicer.
- ...

Choose your aligner carefully.

If you work with pre-mRNA, the options are limited.

- Some tools find exons first, then use this to break up reads.
- Some tools prefer splitting reads over mapping them in an intron.

Choose your aligner carefully.

If you work with pre-mRNA, the options are limited.

- Some tools find exons first, then use this to break up reads.
- Some tools prefer splitting reads over mapping them in an intron.

Some tools heavily rely on annotation.

- A list of known splice sites.
- Motifs (canonical splice sites).

Gmap.

Gmap: A Genomic Mapping and Alignment Program for mRNA and EST Sequences.

Gsnap: Genomic Short-read Nucleotide Alignment Program.

<http://research-pub.gene.com/gmap/>

Gmap.

Gmap: A Genomic Mapping and Alignment Program for mRNA and EST Sequences.

Gsnap: Genomic Short-read Nucleotide Alignment Program.

Some features:

- Split read alignment.
 - Split both ends.
 - Split a read into many pieces.
- Fast.
- Memory efficient.
 - No limit on intron size.

<http://research-pub.gene.com/gmap/>

Tophat

Tophat: A spliced read mapper for RNA-Seq.

<http://tophat.cbcb.umd.edu>

Tophat

Tophat: A spliced read mapper for RNA-Seq.

Approach:

- Identify potential exons.
 - Split all reads into small segments, align independently.
- Make a database of splice junctions.
- Map the reads to confirm the splice junctions.

<http://tophat.cbcb.umd.edu>

Tophat

Tophat: A spliced read mapper for RNA-Seq.

Approach:

- Identify potential exons.
 - Split all reads into small segments, align independently.
- Make a database of splice junctions.
- Map the reads to confirm the splice junctions.

Some considerations:

- The software is optimized for reads 75bp or longer.
- Mixing paired- and single- end reads together is not supported.

<http://tophat.cbcb.umd.edu>

Cufflinks.

Input:

- Aligned reads.
 - Gmap / Gsnap.
 - Tophat.

<http://cufflinks.ccb.umd.edu/>

Cufflinks.

Input:

- Aligned reads.
 - Gmap / Gsnap.
 - Tophat.

What it can do:

- Assemble transcripts.
- Estimate transcript abundance.

<http://cufflinks.cbc.umd.edu/>

Cufflinks.

Input:

- Aligned reads.
 - Gmap / Gsnap.
 - Tophat.

What it can do:

- Assemble transcripts.
- Estimate transcript abundance.

Differential expression and regulation (**Cuffcompare**).

<http://cufflinks.cbc.umd.edu/>

Cufflinks.

Modes of operation:

- Use predefined transcripts.
- Assemble transcripts assisted by known transcripts.
- Assemble transcripts with no prior knowledge.

<http://cufflinks.ccb.umd.edu/>

Cufflinks.

Modes of operation:

- Use predefined transcripts.
- Assemble transcripts assisted by known transcripts.
- Assemble transcripts with no prior knowledge.

When to use:

- Only interested in expression.
- Alternative splicing.

<http://cufflinks.ccb.umd.edu/>

Principle of variant calling

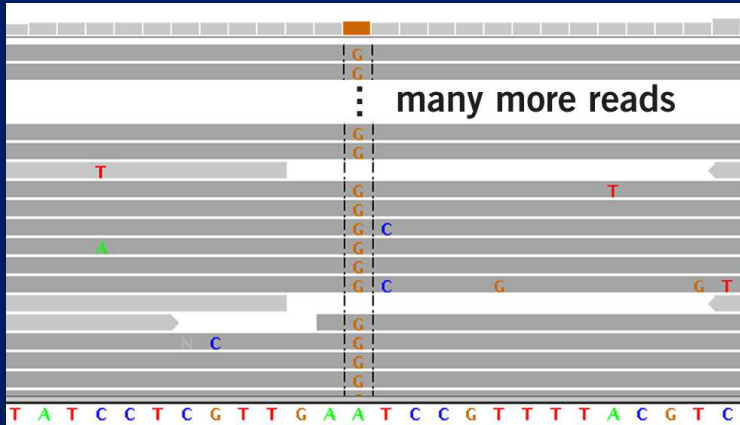


Figure 5 : Result of an alignment.

Principle of variant calling

In principle, we call a variant when we are confident we have seen one.

Principle of variant calling

In principle, we call a variant when we are confident we have seen one.

But when are we confident?

- More than x times?
- In more than y percent of the reads covering the variant?

Principle of variant calling

In principle, we call a variant when we are confident we have seen one.

But when are we confident?

- More than x times?
- In more than y percent of the reads covering the variant?

Variant callers can use:

- Fixed settings.
- Statistical models.

Some considerations

Things a variant caller might take into account:

- Strand balance.
- Base quality.
- Mapping quality.
 - Distribution within the reads.
- Ploidity of the organism in question.

Some considerations

Things a variant caller might take into account:

- Strand balance.
- Base quality.
- Mapping quality.
 - Distribution within the reads.
- Ploidity of the organism in question.

Some complications when analysing RNA:

- Allele specific expression.
 - Heterozygosity may not be detected.

Some considerations

Things a variant caller might take into account:

- Strand balance.
- Base quality.
- Mapping quality.
 - Distribution within the reads.
- Ploidity of the organism in question.

Some complications when analysing RNA:

- Allele specific expression.
 - Heterozygosity may not be detected.
- Tissue specific expression.
 - Some variants will be missed completely.

Some considerations

Things a variant caller might take into account:

- Strand balance.
- Base quality.
- Mapping quality.
 - Distribution within the reads.
- Ploidity of the organism in question.

Some complications when analysing RNA:

- Allele specific expression.
 - Heterozygosity may not be detected.
- Tissue specific expression.
 - Some variants will be missed completely.
- RNA editing.
 - Some variants will not be present on DNA.

Some considerations

Things a variant caller might take into account:

- Strand balance.
- Base quality.
- Mapping quality.
 - Distribution within the reads.
- Ploidity of the organism in question.

Some complications when analysing RNA:

- Allele specific expression.
 - Heterozygosity may not be detected.
- Tissue specific expression.
 - Some variants will be missed completely.
- RNA editing.
 - Some variants will not be present on DNA.
- Strand specific sampleprep.

Combining tools in a pipeline.

```
1 bwa aln -t 8 $reference $i > $i.sai
2 bwa samse $reference $i.sai $i > $i.sam
3 samtools view -bt $reference -o $i.bam $i.sam
```

Listing 1 : Shell script.

Combining tools in a pipeline.

```
1 bwa aln -t 8 $reference $i > $i.sai
2 bwa samse $reference $i.sai $i > $i.sam
3 samtools view -bt $reference -o $i.bam $i.sam
```

Listing 1 : Shell script.

```
1 %.sai: %.fq
2 $(BWA) aln -t $(THREADS) $(call MKREF, $@) $< > $@
3
4 %.sam: %.sai %.fq
5 $(BWA) samse $(call MKREF, $@) $^ > $@
6
7 %.bam: %.sam
8 $(SAMTOOLS) view -bt $(call MKREF, $@) -o $@ $<
```

Listing 2 : Makefile.

Outline of the practical

1. Do a typical RNA-seq analysis.
 - Expression.
2. Workflows.
 - Rerun the analysis with no effort.
3. Differential expression analysis.

Acknowledgements:

Wibowo Arindrarto

Wai Yi Leung

Hailiang Mei

Irina Pulyakhina

Peter-Bram 't Hoen

Johan den Dunnen