



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



Characteristics:

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

Figure 1: HiSeq 2000.

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.

FastX / FastQC.

We use the Trimmomatic / FastX toolkit for data cleaning.

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

FastX / FastQC.

We use the Trimmomatic / FastX toolkit for data cleaning.

- Remove linker sequences.
- Clip low quality reads at the end of the read.
- Judge 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.

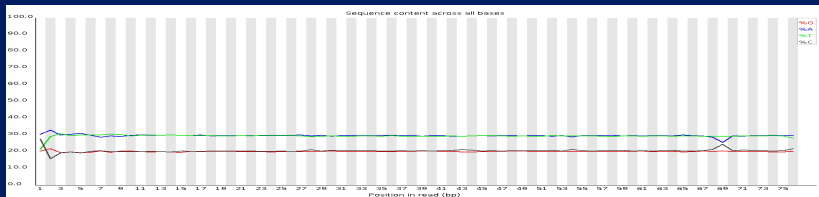


Figure 3: Per base sequence content.

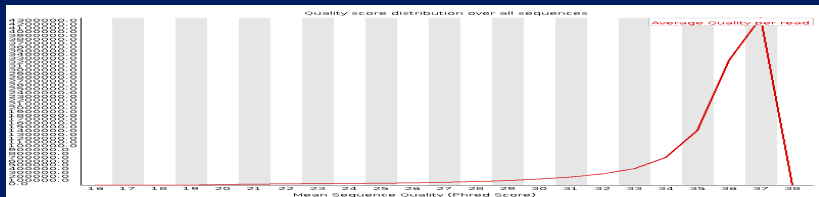


Figure 4: Per sequence quality.

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.
- Motives (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/>

Cufflinks.

Input:

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

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

Cufflinks.

Input:

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

What it can do:

- Assembled transcripts.
- Estimated transcript abundance.

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

Cufflinks.

Input:

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

What it can do:

- Assembled transcripts.
- Estimated 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

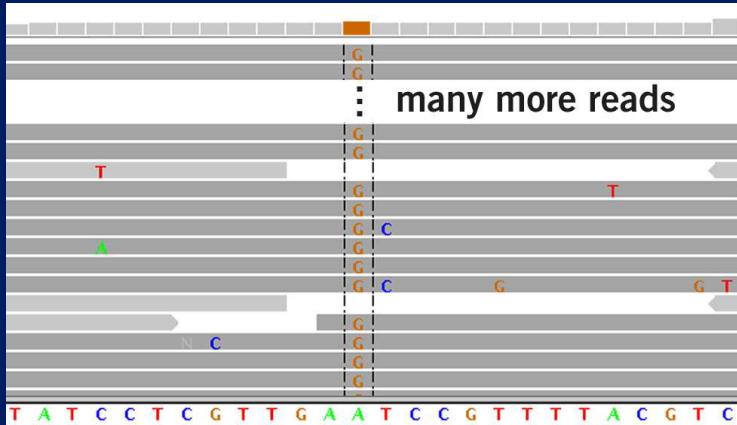


Figure 5: Result of an alignment.

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.

Overview.

Data intensive biology for everyone.

- Open source.
- Web based.
 - No installation required.

<http://galaxy.psu.edu/>

<http://galaxy.nbic.nl/>

Overview.

Data intensive biology for everyone.

- Open source.
- Web based.
 - No installation required.
- Wrapper for command line utilities.
- User friendly.
- Point and click.
- Workflows.
 - Save all the steps you did in your analysis.
 - Rerun the entire analysis on a new dataset.
 - Share your workflow with other people.
 - ...

<http://galaxy.psu.edu/>

<http://galaxy.nbic.nl/>

The Galaxy GUI.

The screenshot displays the Galaxy web interface for the Netherlands Bioinformatics Centre. The main panel shows the configuration for the GMAP (version 2.0.0) tool. The configuration includes a dropdown for 'Use a built-in index', a 'Select a reference genome' dropdown set to 'Human_UCSC_hg19_complete', and various options for 'kmer size', 'Look for splicing involving known sites or known introns', 'Select an mRNA or EST dataset to map', 'Protocol for input quality scores', 'Select the output format', and 'SAM paired reads'. There are also input fields for 'Value to put into read-group id (RG-ID) field', 'Value to put into read-group name (RG-SM) field', 'Value to put into read-group library (RG-LB) field', and 'Value to put into read-group library platform (RG-PL) field'.

On the right side, the 'History' panel shows a list of recent jobs with their names and sizes, such as '12: SAM-to-BAM on data 2: converted BAM (Genome Coverage [testGlaph])' (212.9 Mb) and '11: Cufflinks on data 8 and data 1: assembled transcripts'.

Figure 6: Galaxy panels.




Galaxy icons.

Figure 7: Collapsed history item.

- Eye: view.
- Pencil: edit (rename).
- Cross: delete.




- Click on the title for a more detailed view.

Galaxy icons.

34: Mpileup on data 33:   

Output

~1,100,000 genomic coordinates
format: pileup, database: hgtest

1. Chrom	2. Start	3. Base	4	5	6
chr1	25620470	N	1	^	A =
chr1	25620471	N	1	G	=
chr1	25620472	N	1	A	>
chr1	25620473	N	1	T	>
chr1	25620474	N	1	A	>
chr1	25620475	N	1	T	>

Figure 8: History item.

- Diskette: save.
- Blue looping arrow: rerun.

Outline of the practical

1. Do a typical RNA-seq analysis.
 - Expression.
 - Novel transcripts.
2. Variant calling.
3. Workflows.
 - Rerun the analysis with no effort.

Acknowledgements:

Hailiang Mei
Michiel van Galen
Martijn Vermaat
Johan den Dunnen