



Introduction to RNAseq Analysis

Milena Kraus
Apr 18, 2016

Agenda

What is RNA sequencing used for?

1. Biological background
2. From wet lab sample to transcriptome
 - a. Experimental procedure
 - b. Raw data
 - c. Processing pipeline(s)
 - d. Downstream analysis

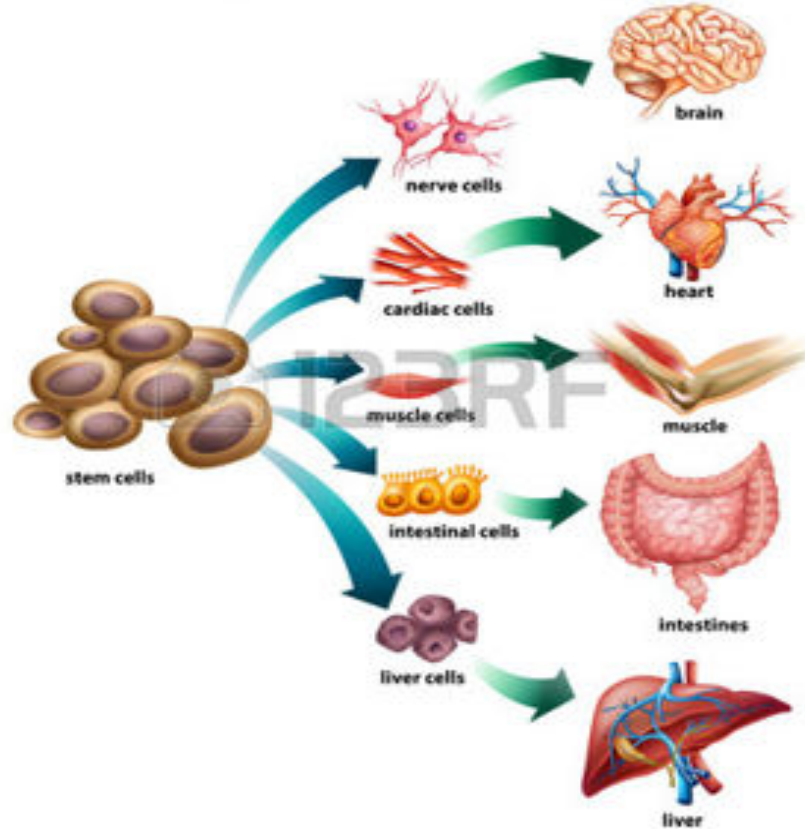
RNAseq Intro

Milena Kraus, Apr
19, 2016

Chart 2

How is a muscle cell different from a liver cell?

- Every cell in your body contains the same DNA as every other cell
- The DNA codes for every process in the cell

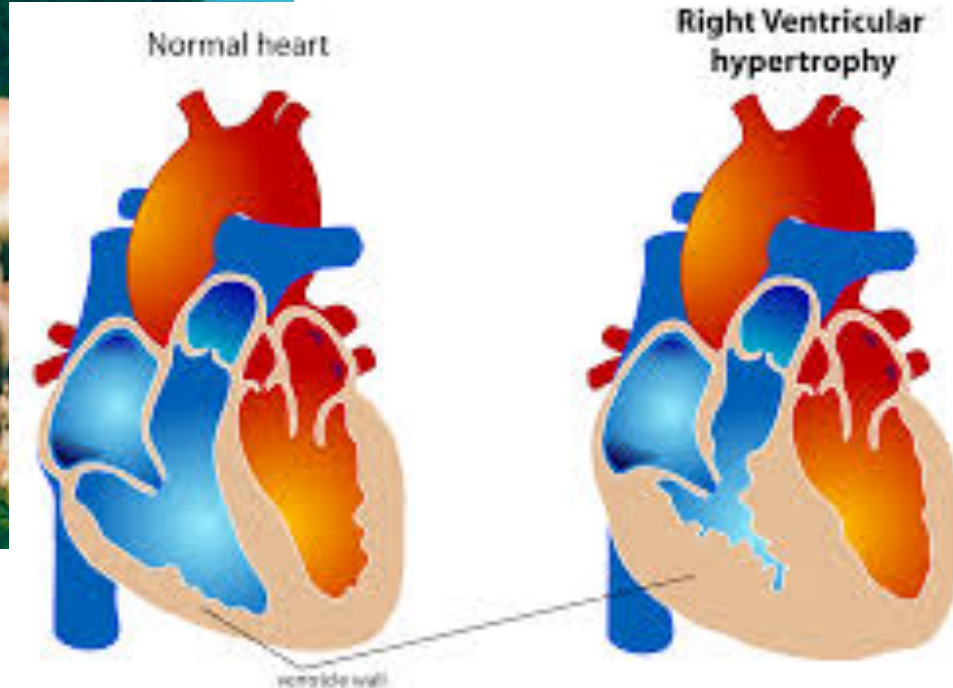


RNAseq Intro

Milena Kraus, Apr 19, 2016

Chart 3

What is the difference between a healthy heart and a sick heart?

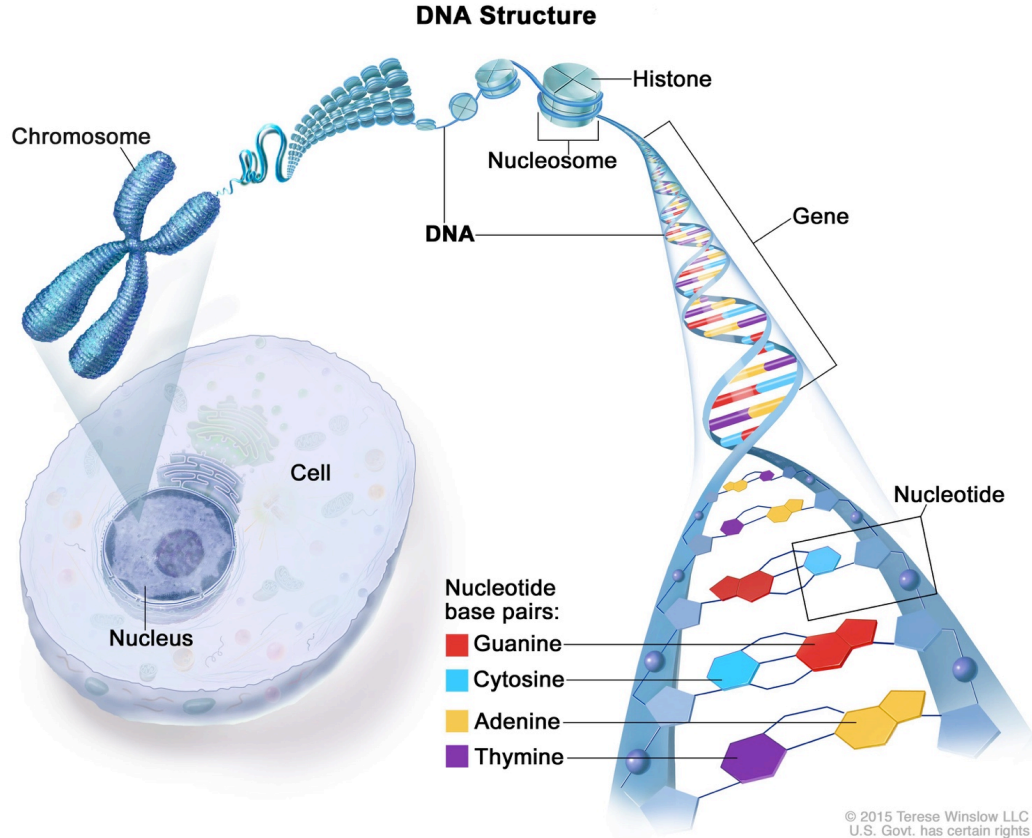


RNAseq Intro

Milena Kraus, Apr 19, 2016

Chart 4

Biological Background



RNAseq Intro

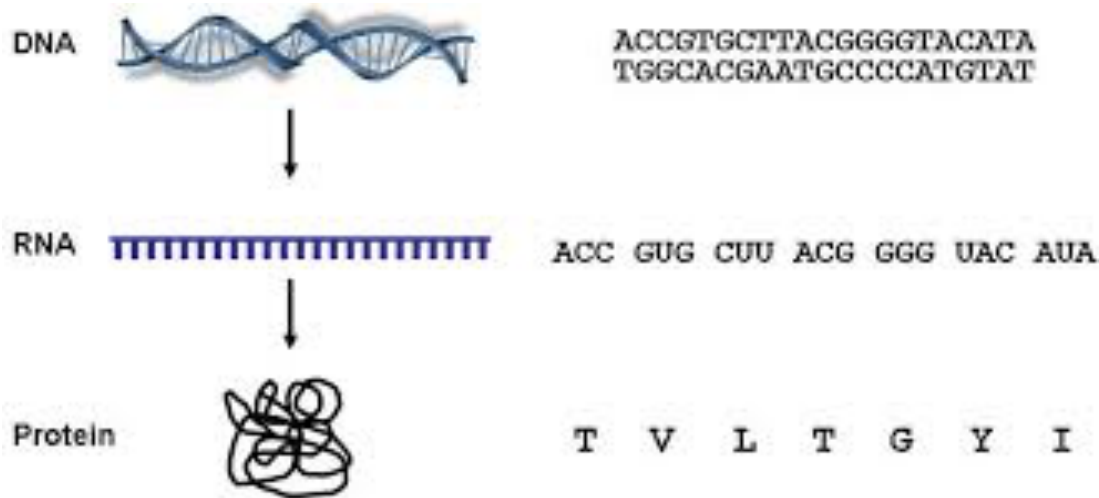
Milena Kraus, Apr 19, 2016

Chart 5

Central Dogma of Molecular Biology

From DNA to RNA to protein

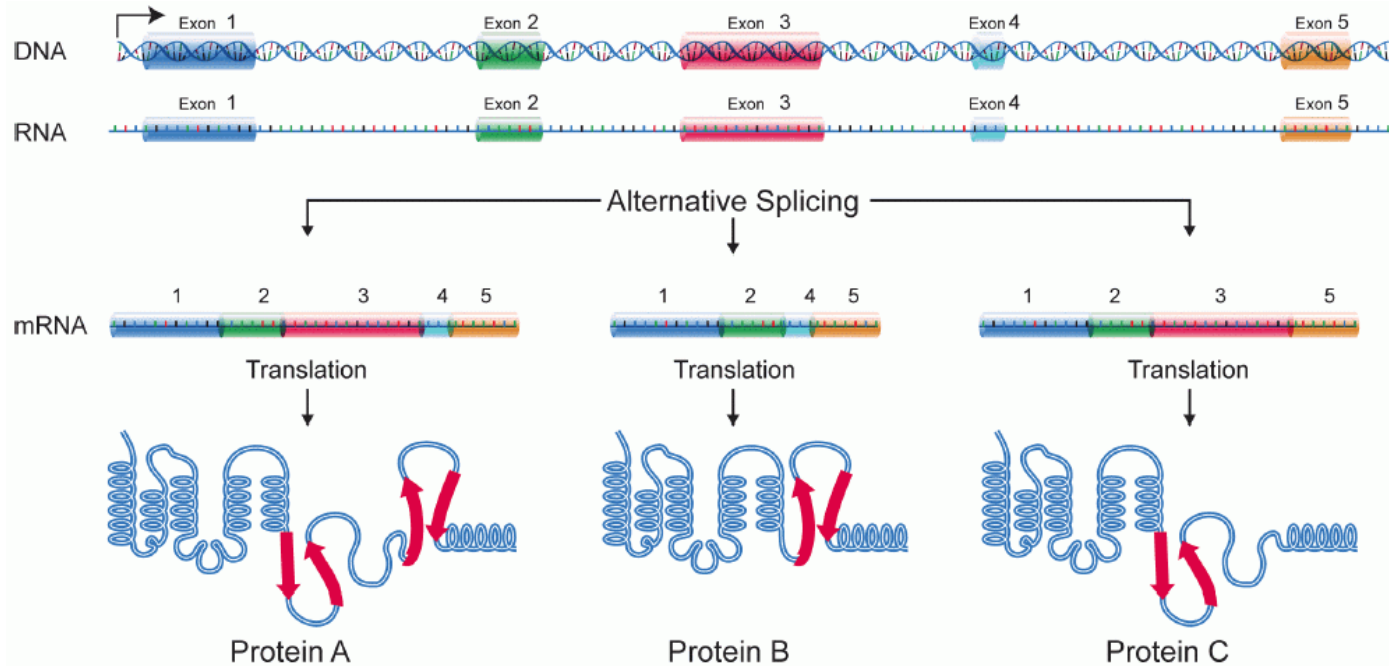
- Interesting information from RNA:
 - Sequence
 - Quantity



RNAseq Intro

Milena Kraus, Apr 19, 2016

One more bio fact before we start: Alternative Splicing



- Von National Human Genome Research Institute - http://www.genome.gov/Images/EdKit/bio2j_large.gif, Gemeinfrei, <https://commons.wikimedia.org/w/index.php?curid=2132737>

RNAseq Intro

Milena Kraus, Apr
19, 2016

Chart 7

From Wet Lab Experiment to Transcriptome

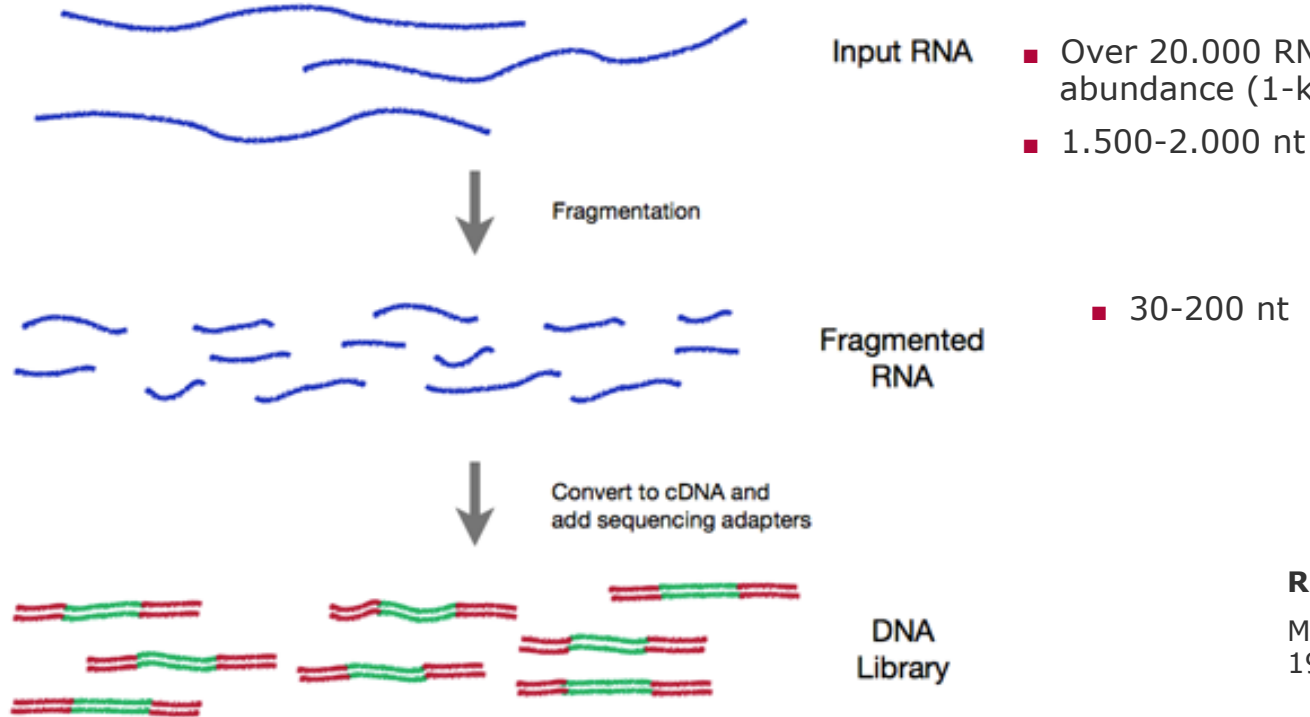


RNAseq Intro

Milena Kraus, Apr
19, 2016

Chart **8**

Experimental Procedure



RNAseq Intro

Milena Kraus, Apr 19, 2016

Sequencer



SOLID 5500



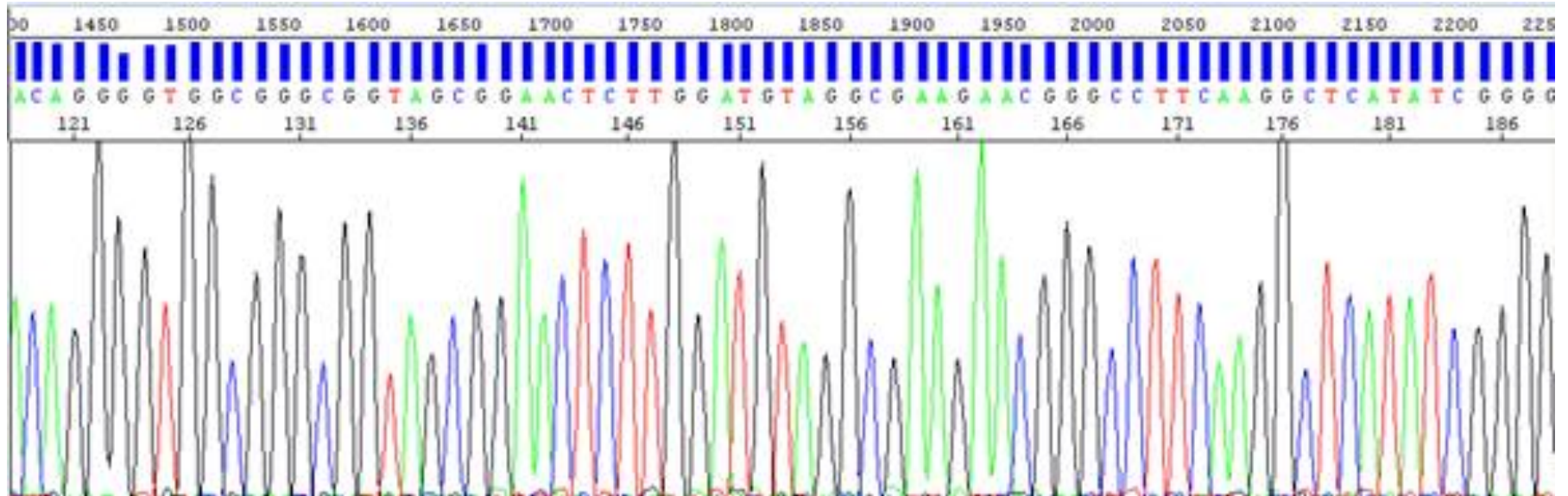
Illumina HiSeq2000

RNAseq Intro

Milena Kraus,
19.04.2016

Chart 10

Sequencing Signal



RNAseq Intro

Milena Kraus, Apr 19, 2016

Chart 11

Raw data FASTQ files

```
@SRR831012.1 HWI-ST155_0742:7:1101:1284:1981/1
NGAGATGAAGCACTGTAGCTTGGAATTCTCGGGTGCCAAGGAACTCCAGT
+
%1=DDDFHHHGGFIHHIIIIIIIIIIIIIEHIIIIIIIFIIIIII
```

```
@SRR831012.2 HWI-ST155_0742:7:1101:2777:1998/1
NGAGATGAAGCACTGTAGCTCTTTGGAATTCTCGGGTGCCAAGGAACTCC
+
%1=DFFFHHHHHHIIIIIIIIIIIIIIIIIIIIIGIIIIIIIIIIIIIG
```

Quality score (increasing from worst to best):

```
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNPOQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
```

@SampleID.ReadNr

Experimental Setup

In our setting:

- ~1.4 GB per file
- ~8 Mio reads per file
- 80 files

RNAseq Intro

Milena Kraus, Apr
19, 2016

Chart **12**

Raw data

Reference genome

- FASTA-file

>Sequenz 1

;comment A

```
ACAAGATGCCATTGTCCCCCGGCCTCCTGCTGCTGCTGCTCTCCGGGGCCACGGC  
CACCGCTGCCCTGCCCCTGGAGGGTGGCCCCACCGGCCGAGACAGCGAGCATAT  
GCAGGAAGCGGCAGGAATAAGGAAAAGCAGCCTCCTGACTTTCTCGCTTGGTGG  
TTTGAGTGGACCTCCCAGGCCAGTGCCGGGCCCCCTCATAGGAGAGGAAGCTCGG  
GAGGTGGCCAGGCGGCAGGAAGGCGCACCCCCCAGCAATCCGCGCGCCGGGA  
CAGAATGCCCTGCAGGAAGTTCTTCTGGAAGACCTTCTCCTCCTGCAAATAAACCT  
CACCCATGAATGCTCACGCAAGTTTAATTACAGACCTGAA...
```

In our setting:

- Indexed HG19 (Humane Genome V19)
- HG consists of approx 3.2B nucleotides

RNAseq Intro

Milena Kraus, Apr
19, 2016

Chart **13**

Raw data

Gene library

- 20k-25k protein coding genes representing small part of the genome
- Using the annotation to speed up processing
- If the discovery of new genes in a sample is expected, a custom annotation can be calculated from the reads

<u>Col 1</u>	<u>Col 2</u>	<u>Col 3</u>	<u>Col 4</u>	<u>Col 5</u>	<u>Col 6</u>	<u>Col 7</u>	<u>Col 8</u>	<u>Col 9</u>
chr21	HAVANA	transcript	10862622	10863067	.	+	.	gene_id "ENSG00000169..
chr21	HAVANA	exon	10862622	10862667	.	+	.	gene_id "ENSG00000169..
chr21	HAVANA	CDS	10862622	10862667	.	+	0	gene_id "ENSG00000169..
chr21	HAVANA	start_codon	10862622	10862624	.	+	0	gene_id "ENSG00000169..
chr21	HAVANA	exon	10862751	10863067	.	+	.	gene_id "ENSG00000169..
chr21	HAVANA	CDS	10862751	10863064	.	+	2	gene_id "ENSG00000169..
chr21	HAVANA	stop_codon	10863065	10863067	.	+	0	gene_id "ENSG00000169..
chr21	HAVANA	UTR	10863065	10863067	.	+	.	gene_id "ENSG00000169..

- In our setting: geneshg19.gtf

RNAseq Intro

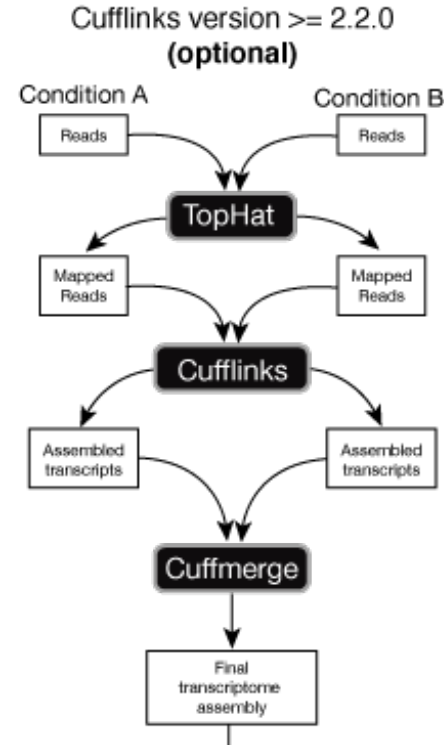
Milena Kraus, Apr
19, 2016

Chart 14

Processing pipeline

Gold standard – tophat/cufflinks

- **TopHat** aligns RNA-Seq reads to mammalian-sized genomes using the ultra high-throughput short read aligner BOWTIE, and then analyzes the mapping results to identify splice junctions between exons.
- **Cufflinks** assembles mapped RNA-Seq reads into transcripts.
- **Cuffmerge** creates an assembly of all transcripts to build the transcriptome (occurrence transcripts).

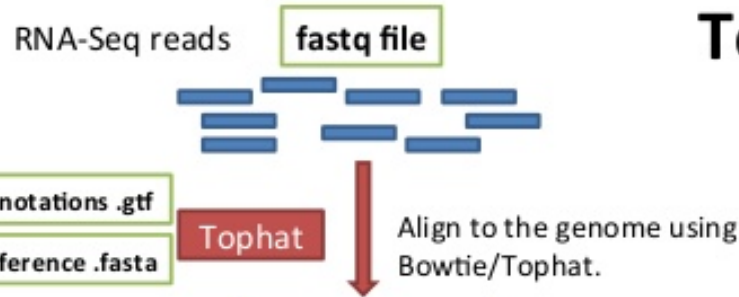


RNAseq Intro

Milena Kraus, Apr 19, 2016

Chart 15

Tophat/Cufflinks Workflow



SAM/BAM file



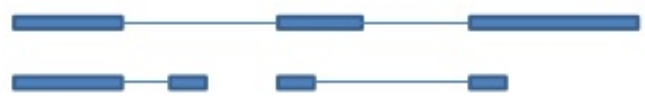
Genomic mapped reads may identify novel isoforms.

Spliced Fragments align to known exon-exon junctions.

Gene annotations .gtf
Genome reference .fasta

Cufflinks

Transcript isoforms **Gene/transcript quantification**



Cufflinks identifies mutually exclusive exons. Graph-based analysis uses a shortest-path algorithm to determine

Cuffmerge

- Input: Transcript library condition-wise (.gtf)
- Algorithm: Counts/Assembles all transcripts found in the different conditions
- Output: Library of all transcripts over all conditions (.gtf)

- The transcriptome ...
 - Serves as a reference for further analysis,
 - Contains all found transcripts over all conditions, and
 - Resembles a rough profile of the studied tissue.

Processing Pipeline

New approach: DESeq and DEXseq

- Preprocessing to generate count tables from .bam files with htseq-count

DESeq

- Input: count table including all conditions
- Algorithm: Estimates variance-mean dependence in count data using a negative binomial distribution instead of maximum likelihood.
- Output: table containing gene identifiers and their normalized counts

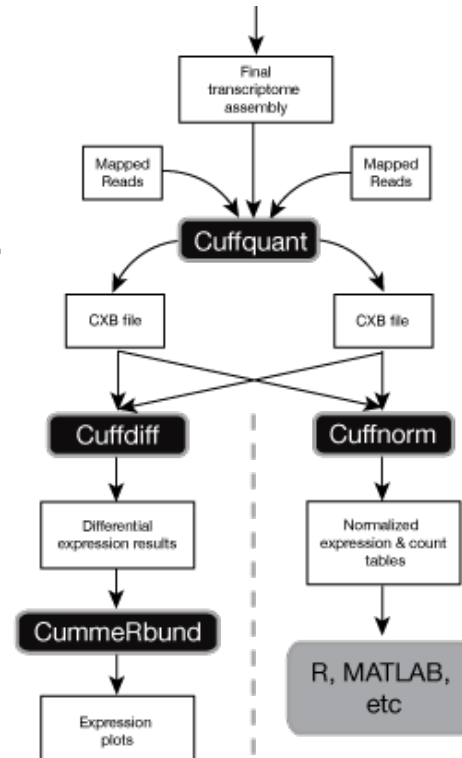
DEXseq

- Same stat. method as DESeq but the output shows differentially expressed exons

Downstream Processing

Statistical analysis and visualization

- **Cuffquant** is an intermediate step that helps to serialize and parallelize analysis.
- **Cuffdiff** compares expression levels of transcripts and shows differentials spliced genes and isoforms.
- **Cuffnorm** normalizes expression levels for exact comparison (usually optional).
- **CummeRbund** is an R package that provides various methods to visualize the data.



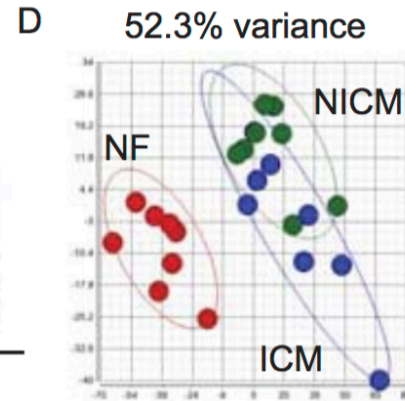
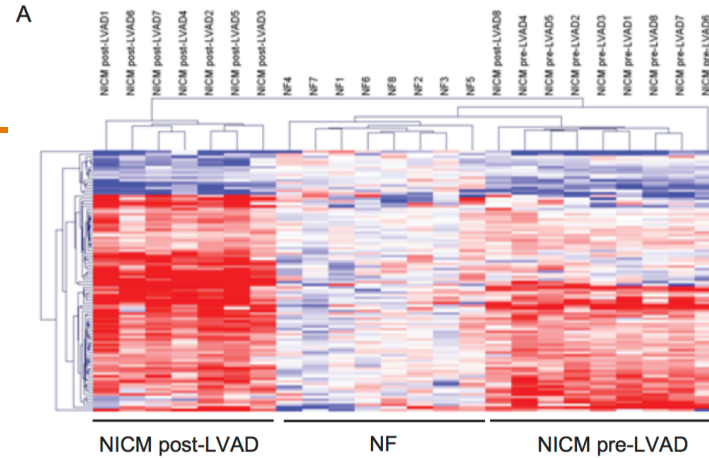
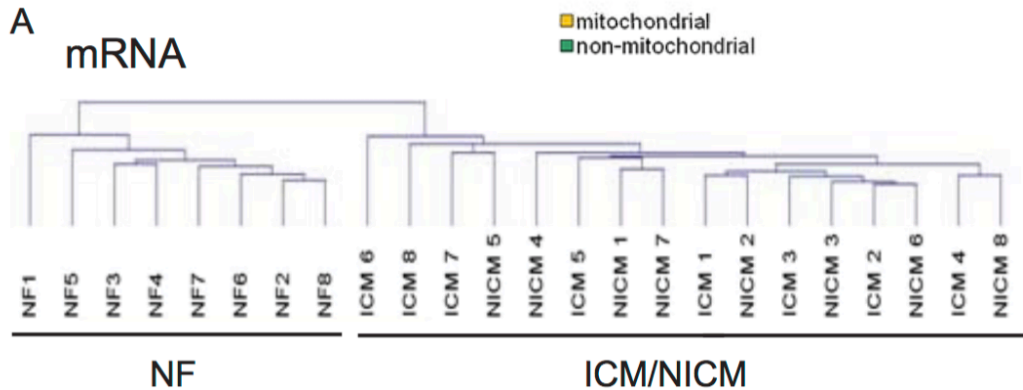
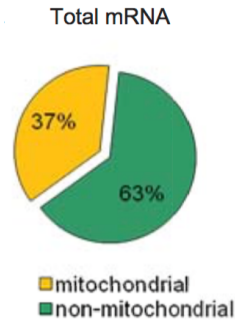
RNAseq Intro

Milena Kraus, Apr 19, 2016

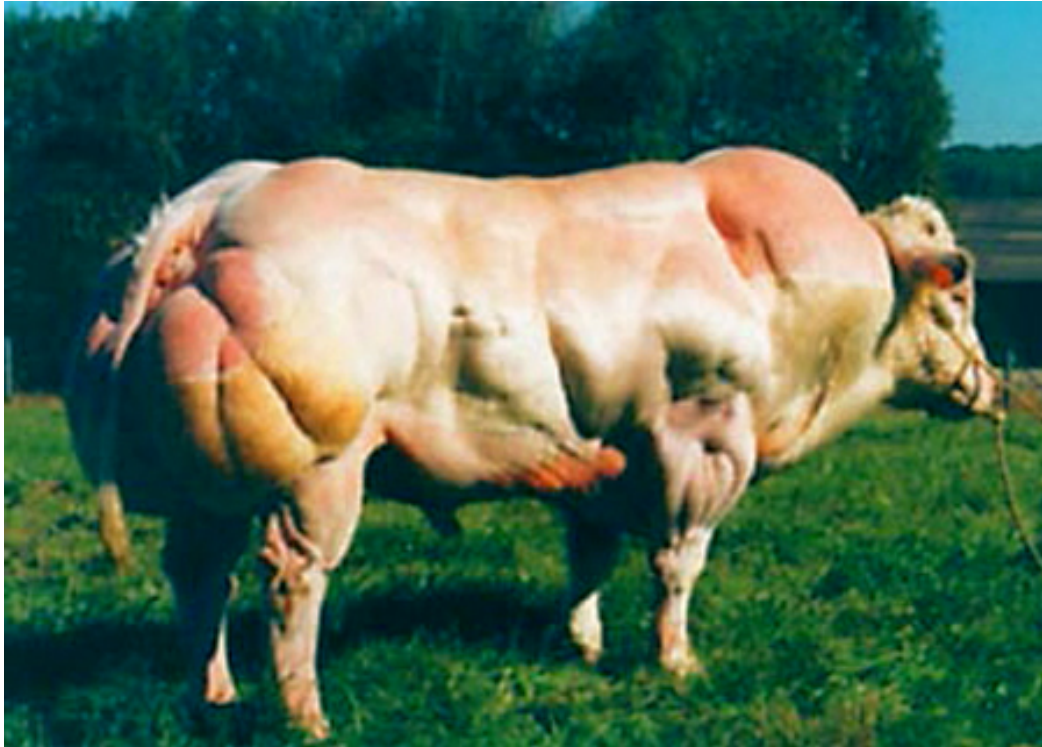
Chart 19

CummeRbund

- R package with common methods for
 - Statistical analysis
 - Visualization



Variant calling from RNAseq data



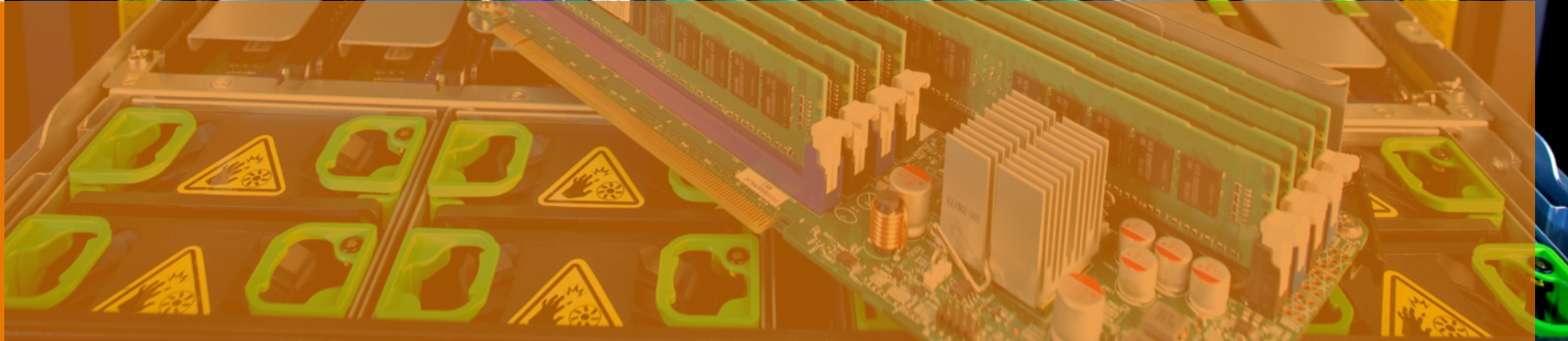
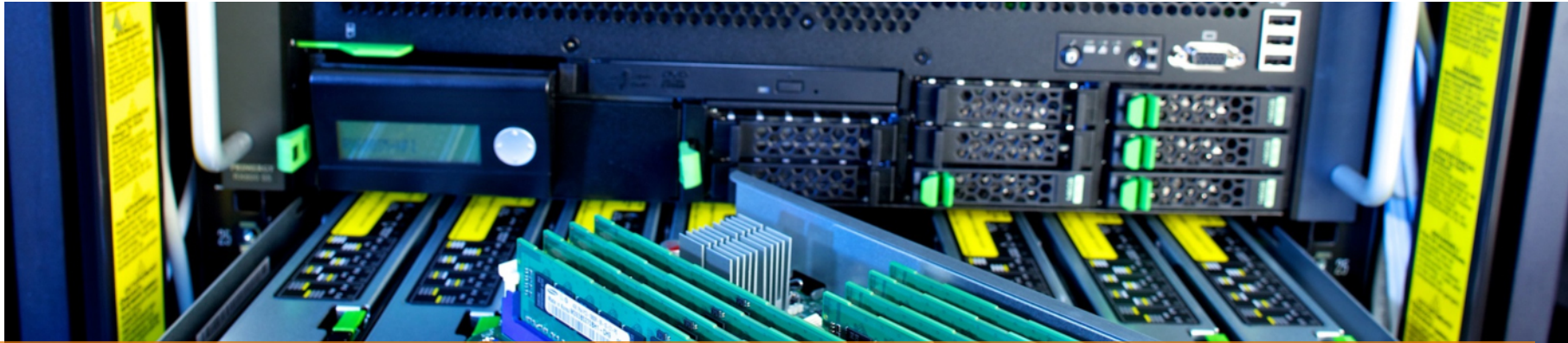
RNAseq Intro

Milena Kraus, Apr
19, 2016

Chart **22**

Thank you
for your attention!

Speaker
Job Description
Institute







Explanations

- Text layers

First text layer for running text.

- Second level for bullet points
 - Third level for bullet points
 - Fourth level for bullet points
- 1. Fifth level for numberings
 - a) Sixth level for listings

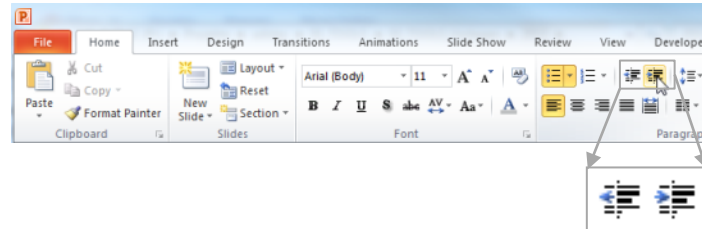
**SEVENTH TEXT LAYER
FOR CORE MESSAGES**

In this template, we pre-formatted different text layers (as you can see on the right side).

You don't have to generate bullet points manually.

By the way: Please avoid this!

To change from one text layer to the next, use the Increase / Decrease List Level buttons:



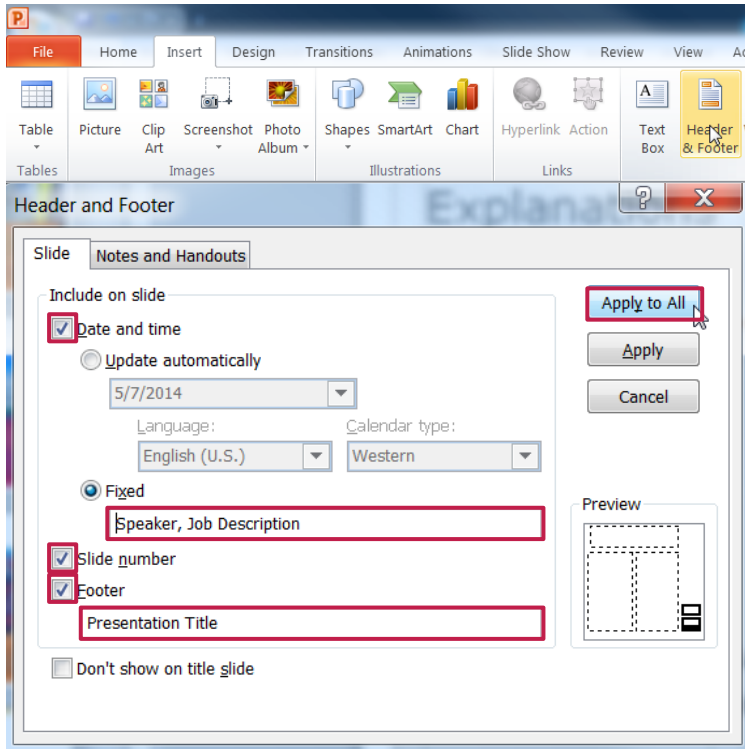
RNAseq Intro

Milena Kraus, Apr 19, 2016

Chart 27

Explanations

- Footer



You can insert or change your presentation's footer. Click on the Insert-tab | Header and Footer | After filling in your descriptions click on **Apply to All**.

Descriptions:

- Activate date and time and write in: *Speaker, Job Description*
- Activate the slide number.
- Activate the footer and write in: *Presentation Title*

Don't use the template without the complete footer.

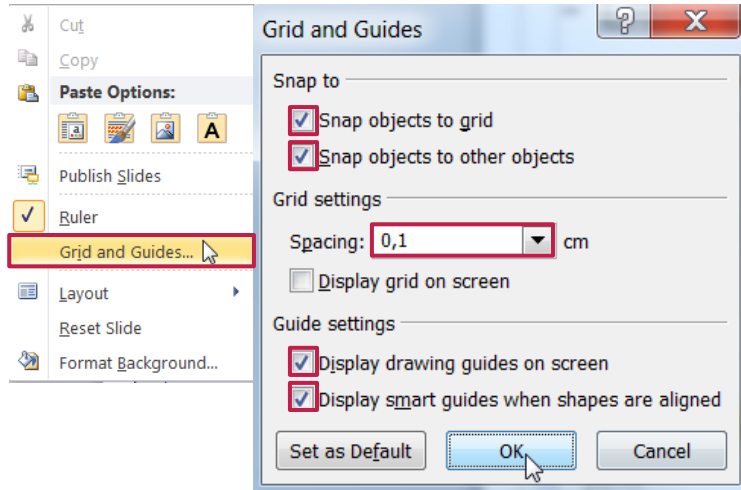
RNaseq Intro

Milena Kraus, Apr 19, 2016

Chart 28

Explanations

- Drawing guides

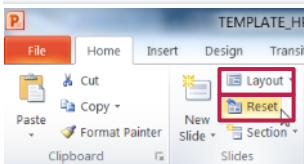
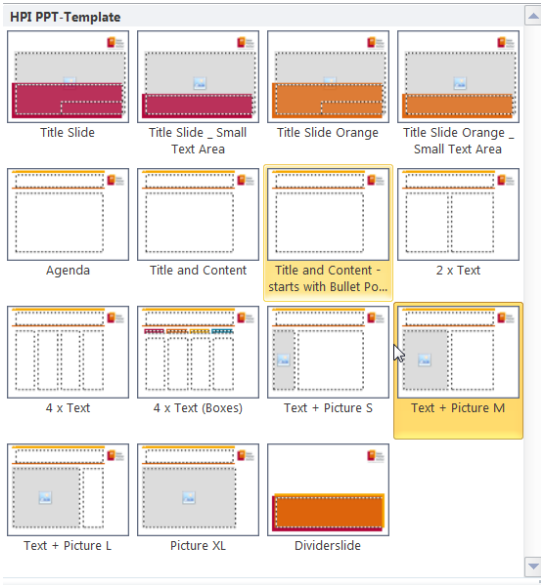


You can enable your guide-lines to align objects on the slide (**View | Show | Select the option „Guides“**)

Or hit the right mouse button outside the slide and go at „Grid and Guides...“

Explanations

- Slide layouts



You can choose between different slide layouts.

These pre-defined layouts gives you the oportunity to use text and visualisations just the right way.

To use these layouts:

Click on the Home-tab | New Slide or Layout | and choose one out of the layouts

Click „Reset“ to reset to the predefined slide layout.

RNaseq Intro

Milena Kraus, Apr 19, 2016

Chart **30**