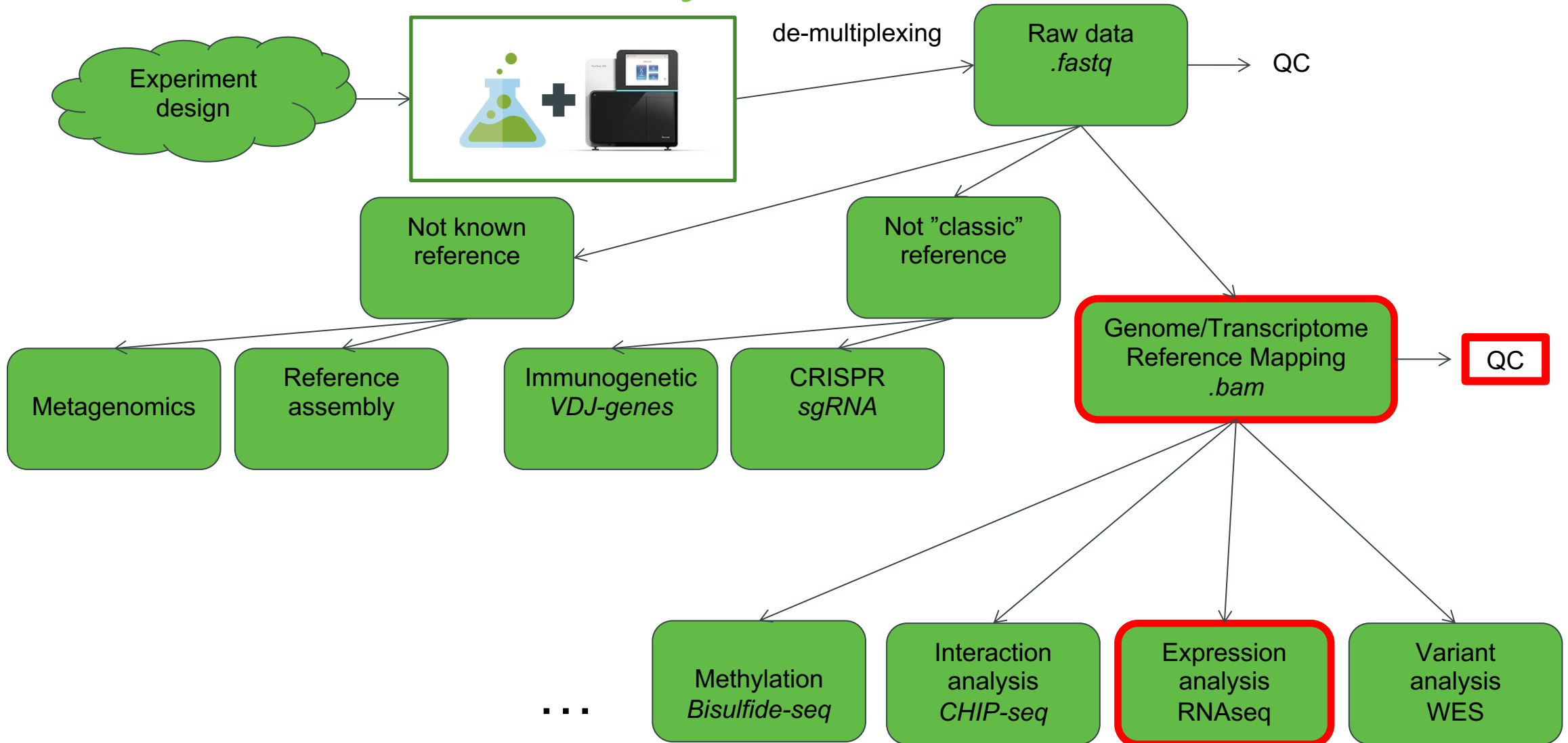**Modern methods for genome analysis (PřF:Bi7420)**

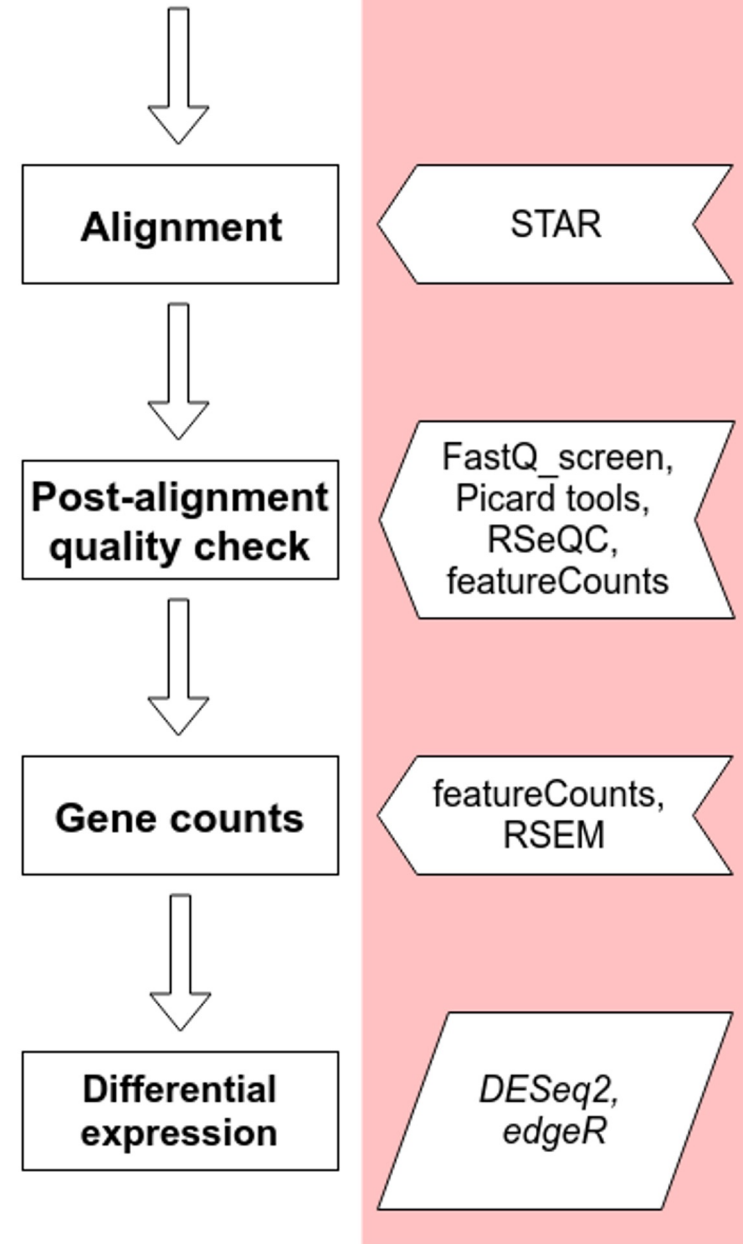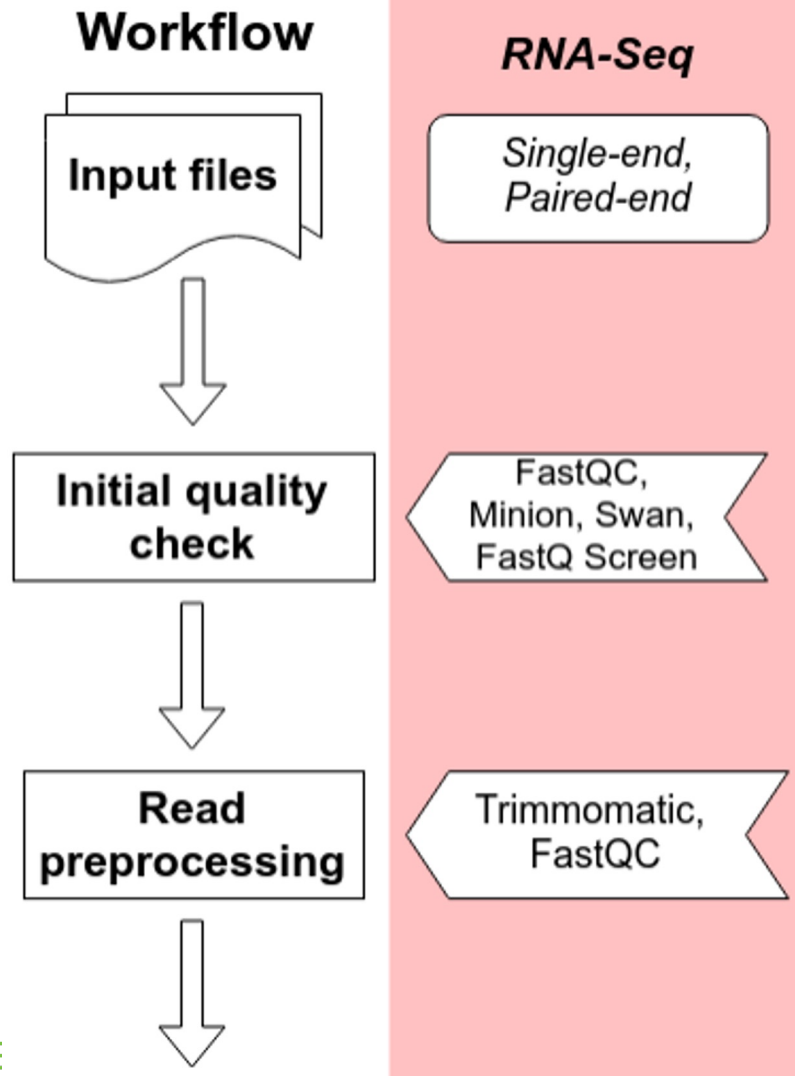# Lecture 5 : RNA-seq primary analysis

Vojta Bystry
vojtech.bystry@ceitec.muni.cz

# NGS data analysis

# The RNA-Seq workflow



**Workflow** — **RNA-Seq**

Input files — *Single-end, Paired-end*

Initial quality check — FastQC, Minion, Swan, FastQ Screen

Read preprocessing — Trimmomatic, FastQC

Alignment — STAR

Post-alignment quality check — FastQ_screen, Picard tools, RSeQC, featureCounts

Gene counts — featureCounts, RSEM

Differential expression — *DESeq2, edgeR*

CEITE

# Alignment

- Mapping to genome or transcriptome?
- Genome
  - Requires spliced alignment
  - Can find novel genes/isoforms/exons
  - Information about whole genome/transcriptome
- Transcriptome
  - No spliced alignments necessary
  - Many reads will map to multiple transcripts (shared exons)
  - Cannot find anything new
  - Difficult to determine origin of reads (multiple copies of transcripts)

CEITEC

# Alignment

- Our choice is the `STAR` aligner

- It performs genome alignment

- Offers a lot of settings to support splicing, soft-clipping, chimeric alignments, ...

- Other techniques (`Salmon` or `Kallisto`) do not use alignment per se and can give you the gene count information right away

  - They use only transcriptome as a reference and are very quick

  - Drawback is you see only what's in the transcriptome and nothing else

CEITEC

# Duplication removal - UMI

- PCR duplicates

- Optical duplicates


- How the tools recognize duplicates
  - Maps to the exact same place
- Problem is it could be identical fragment not PCR duplicate
- UMI helps
  - Maps to the exact same place
  - AND have identical UMI sequence

# Post-alignment QC

- Number of mapped reads - unique + multi mapped
- Mapped locations – intron, exon, intergenic
- Duplication rates
- Library strand specificity
- Captured biotypes
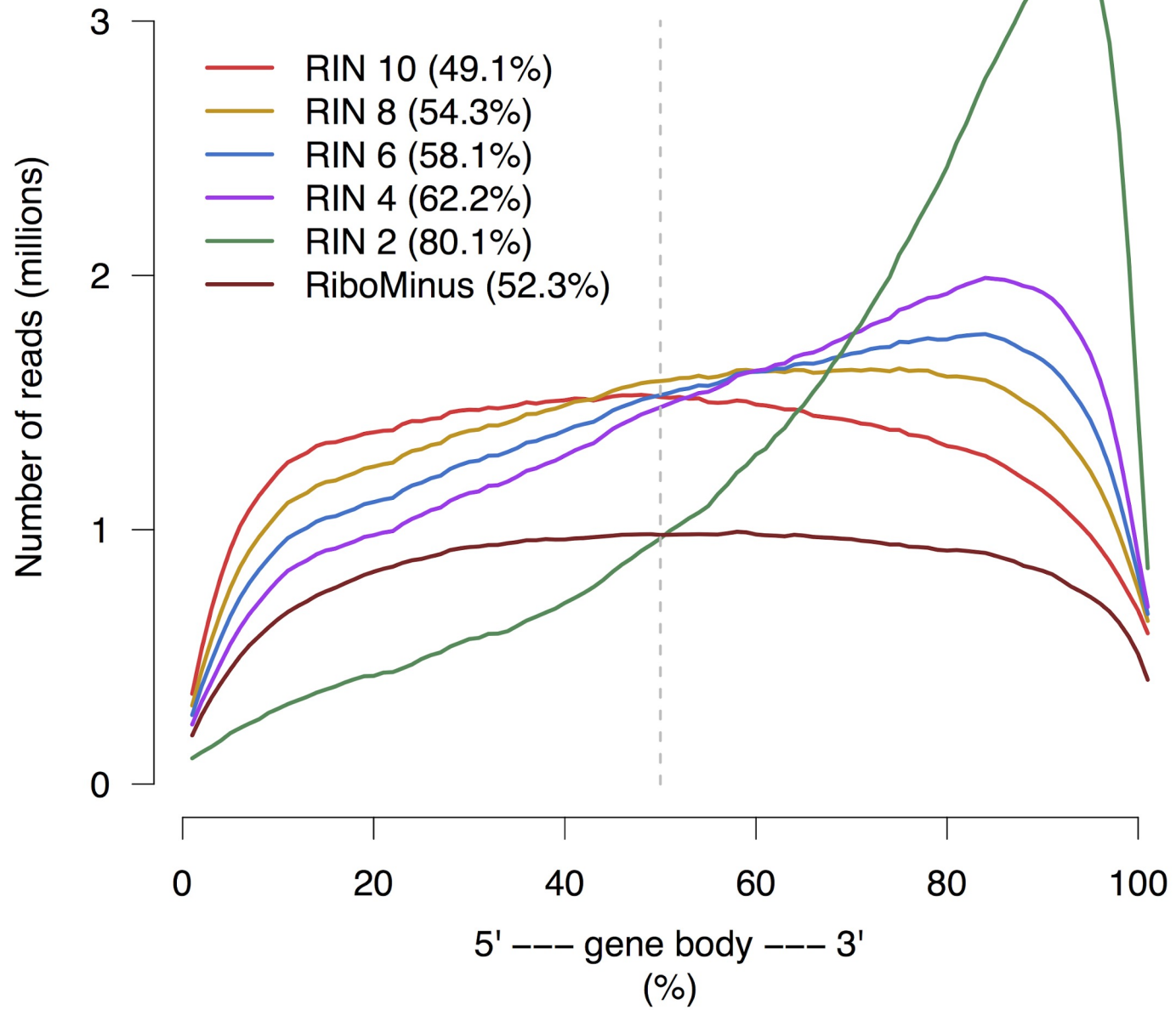- Contamination (rRNA, non-self)
- 5' to 3' end coverage bias

CEITEC

# Post-alignment QC - Tools

- Aligner report
  - `STAR` – most direct assesment

- General QC tools
  - `RSeQC`
  - `Picard`
  - `Qualimap`

- Feature counting tools
  - `featureCounts`
  - `RSEM`

- Non-aligment tools
  - `FastQ screen`
  - `Biobloom`

# Note: Gene body coverage

- Often, libraries with high fragmentation (and low RIN numbers) combined with polyA selection might have strong 3' end bias

  - This is a result of polyA "pulled" fragments

- Some kits, however, target only the polyA tail or sequences close to it

  - An example is Lexogen QuantSeq which sequences only one read per mRNA molecule close to polyA tail

## Gene body coverage



Legend:
- RIN 10 (49.1%)
- RIN 8 (54.3%)
- RIN 6 (58.1%)
- RIN 4 (62.2%)
- RIN 2 (80.1%)
- RiboMinus (52.3%)

Number of reads (millions)

5' ––– gene body ––– 3'
(%)

Coverage

Percent through gene

Created with MultiQC

# Feature counting

- Now, when we know our alignments are solid we need to get the number of reads mapped to a gene (or other feature)

  - From there, we can calculate the differential expression

- The question is, how do we summarize the counts

  - Do we want only uniquely mapped reads

  - Do we want also multi mapped? And how do we assign them? All? One random? Somehow else?

  - And what if we have multiple genes which overlap each other?

# Strand specific library

- We can basically have three strand specificities
  - **Non stranded/Unstranded** - not very common anymore
    - Direction of the read mapping is completely random (50/50)
  - **Forward (sense) stranded** - common for target kits and "bacterial kits"
    - Direction of the read mapping is the **same** as the gene it originates from
  - **Reverse (antisense) stranded** - "default" for Illumina and NEB kits
    - Direction of the read mapping is the **opposite** as the gene it originates from

- In case of paired-end sequencing it's measure by the first (R1) read orientation (FR, RF)

# Feature counting

- The regular settings are - summarize reads mapping to exons (-t exon) and sum them up to gene id (-g gene_id)

- Other possibilities:
  - Count per exons
  - Include introns
  - …

# Gene counts - Tools

- `featureCounts` is build around the "classic" read to gene assignment
  - By default, assigns only uniquely mapped reads an only reads uniquely assignable to a single gene (but both can be changed)
  - Gives you **raw read counts** per **gene**
- `RSEM` is efficient in counting also multi mapped reads and can estimate expression of individual gene isoforms
  - Tries to "weight" the probability a mapped position of a multi mapped read and assign it correctly to the real source
  - Gives you **estimated counts** per **gene** as well as per **isoform** and normalized **TPM** = **Transcripts per million transcripts**

- But, there is a **big differences** in the **minimal required** "good" aligned reads

CEITEC

# Minimal number of reads and expression I

- `RSEM` is less precise in low read counts (<40-50M reads) and for low expressed RNAs (difficult to estimate)

- For lower read counts it's safer to go for `featureCounts`

- Our best practices for a minimal read count for each tools:
  - Less than **40-50M aligned reads** (to the good stuff) -> `featureCounts`
  - More than **40-50M aligned reads** (to the good stuff) -> `RSEM`

- But if you want isoforms!!! -> `RSEM`

CEITEC

# Feature count results

# Post-alignment QC - example

CEITEC

@CEITEC_Brno

Thank you for your attention!

Vojta Bystry
vojtech.bystry@ceitec.muni.cz

www.ceitec.eu

18