# RNA-seq quantification

Dusan Randjelovic
Seven Bridges Genomics, Inc.
Belgrade: 2016/2017

SevenBridges

# Recap (1/3)

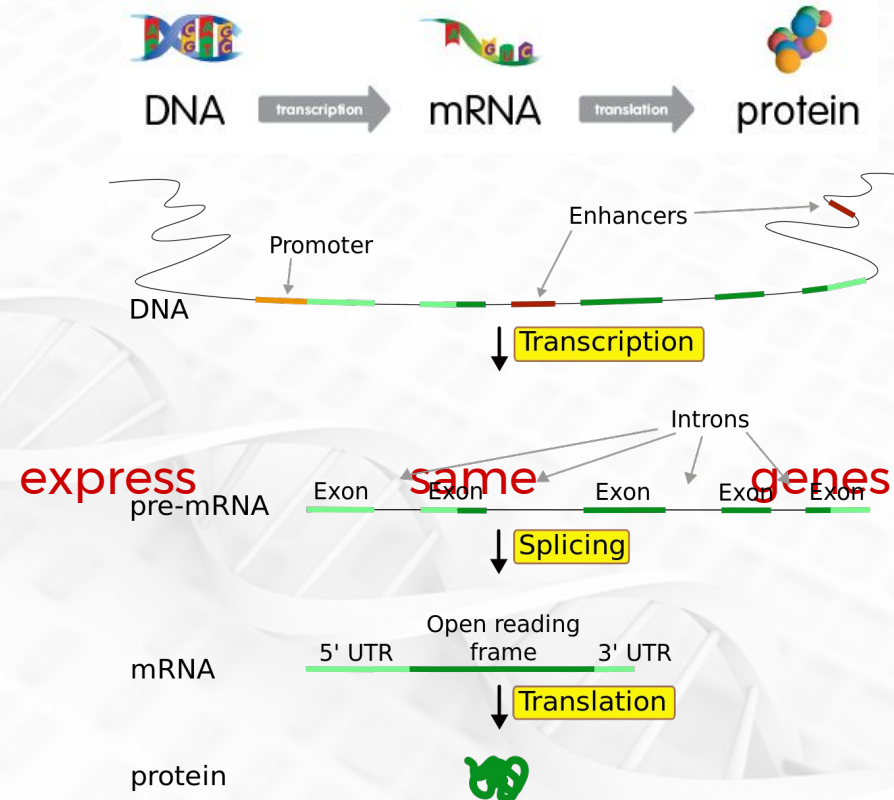- DNA -> RNA -> Protein

- + DNA replication

  (all cells in body have same DNA)

- NOT all cells express same genes

  (due to cell type or cell cycle)
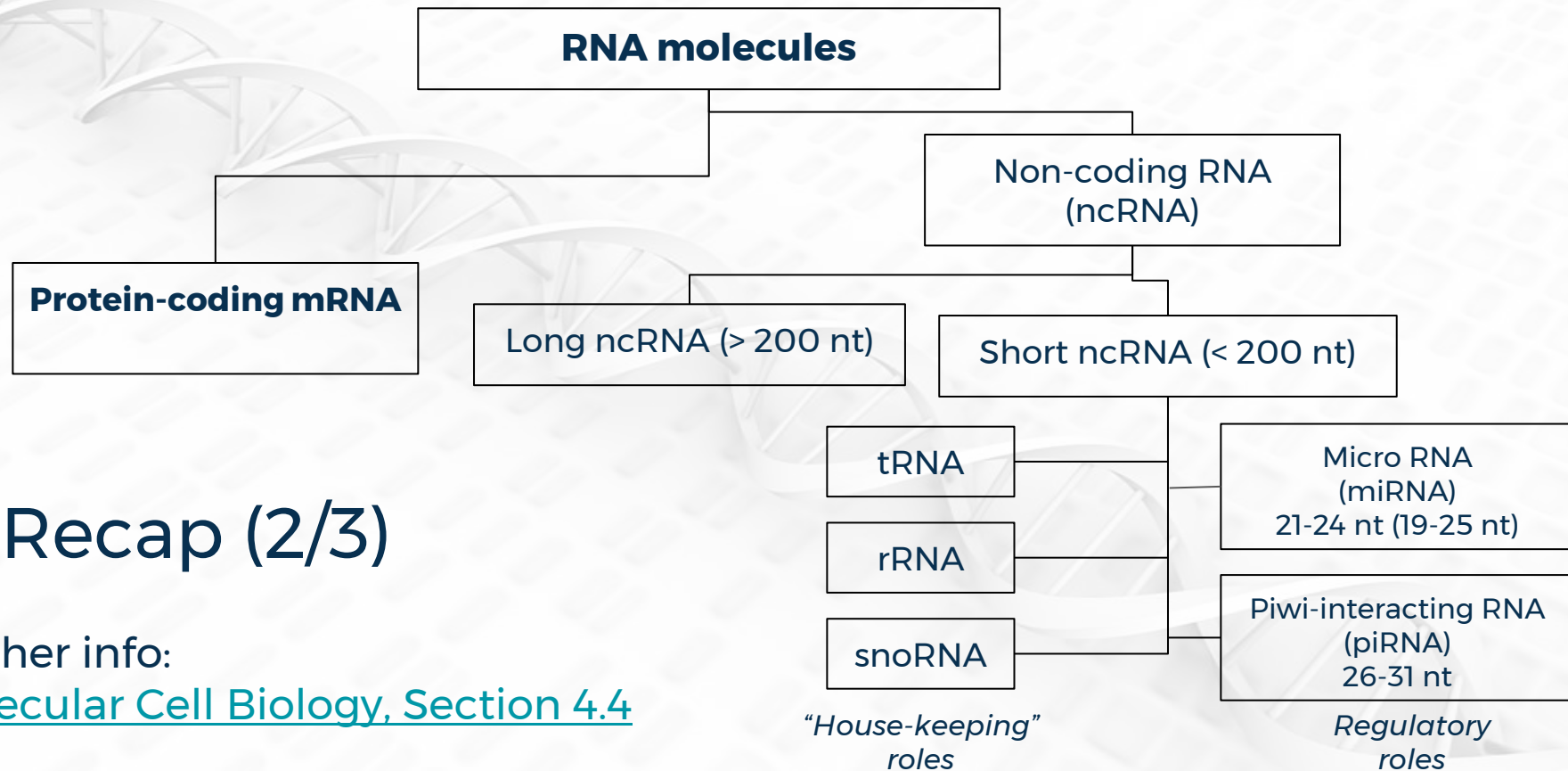
Further info:
animation: DNA from the beginning
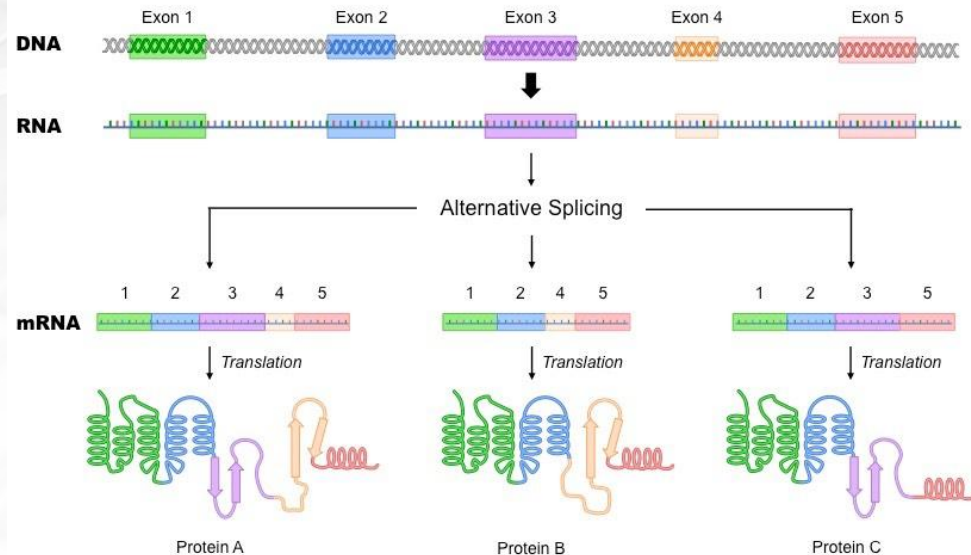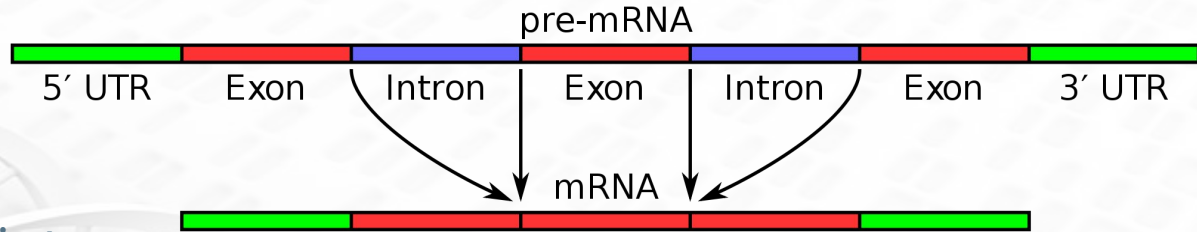Nature Scitable: Gene expression

Central Dogma of Molecular Biology

RNA molecules

Non-coding RNA (ncRNA)

**Protein-coding mRNA**

Long ncRNA (> 200 nt)

Short ncRNA (< 200 nt)

tRNA

rRNA

snoRNA

Micro RNA (miRNA) 21-24 nt (19-25 nt)

Piwi-interacting RNA (piRNA) 26-31 nt

*"House-keeping" roles*

*Regulatory roles*

# Recap (2/3)

Further info:
Molecular Cell Biology, Section 4.4

SevenBridges

# Recap (3/3)

- mRNA = spliced transcript + poly-A tail

- Alternative splicing (isoforms of genes)

Further info:
[DNA Learning Center](#)

# Motivation for RNA quantification

- We (usually) want to check if there is **change in transcription** between conditions (healthy/sick, treated/untreated, different tissues, etc..)
- Typical studies:
  - DNA-seq -> alignment -> variant calling
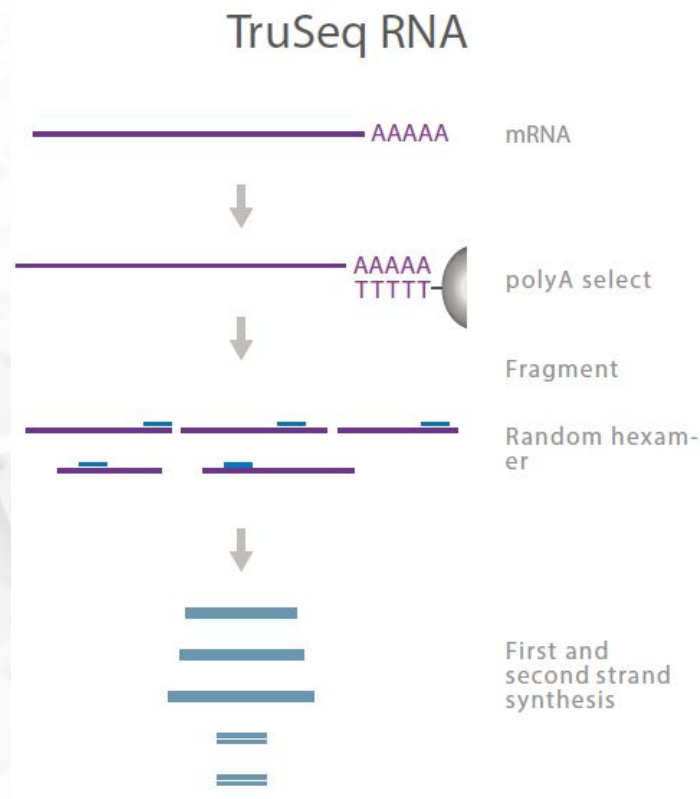  - RNA-seq -> (alignment) -> quantification -> differential expression
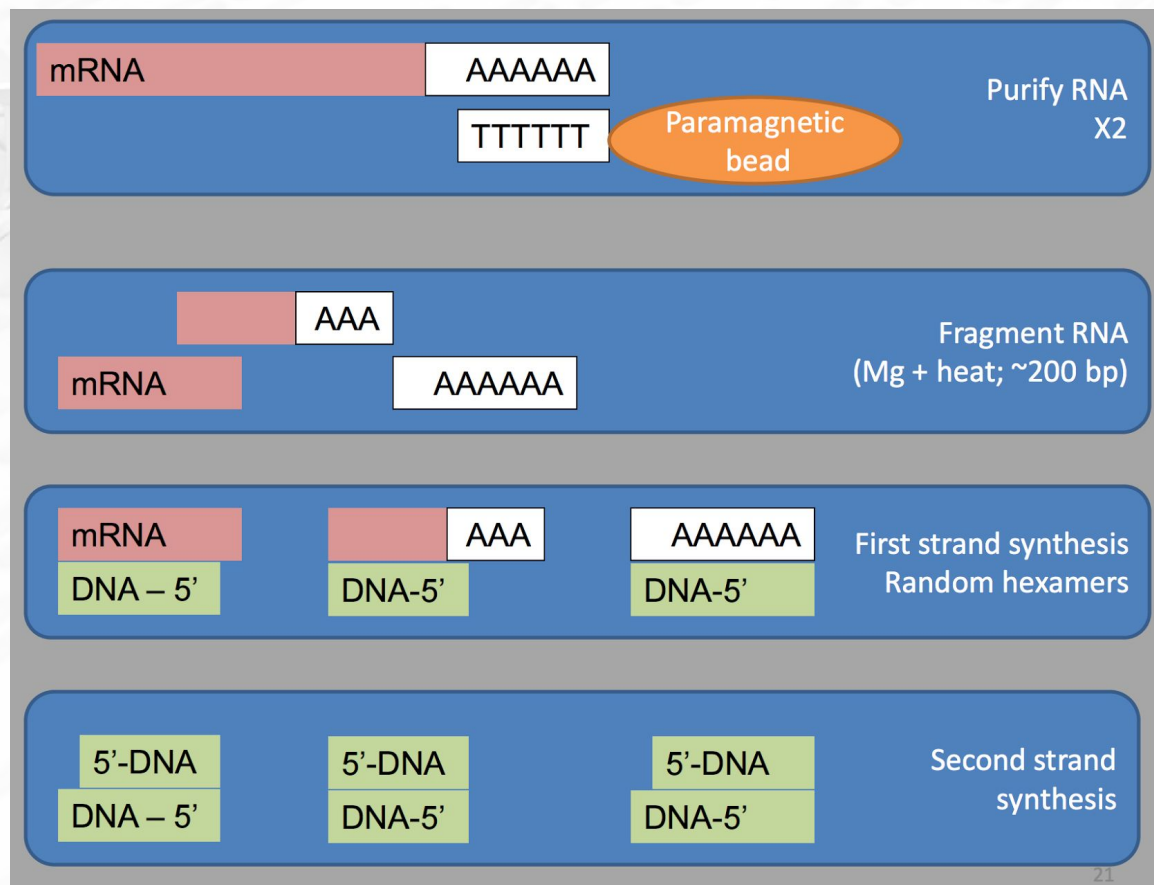
Further info:
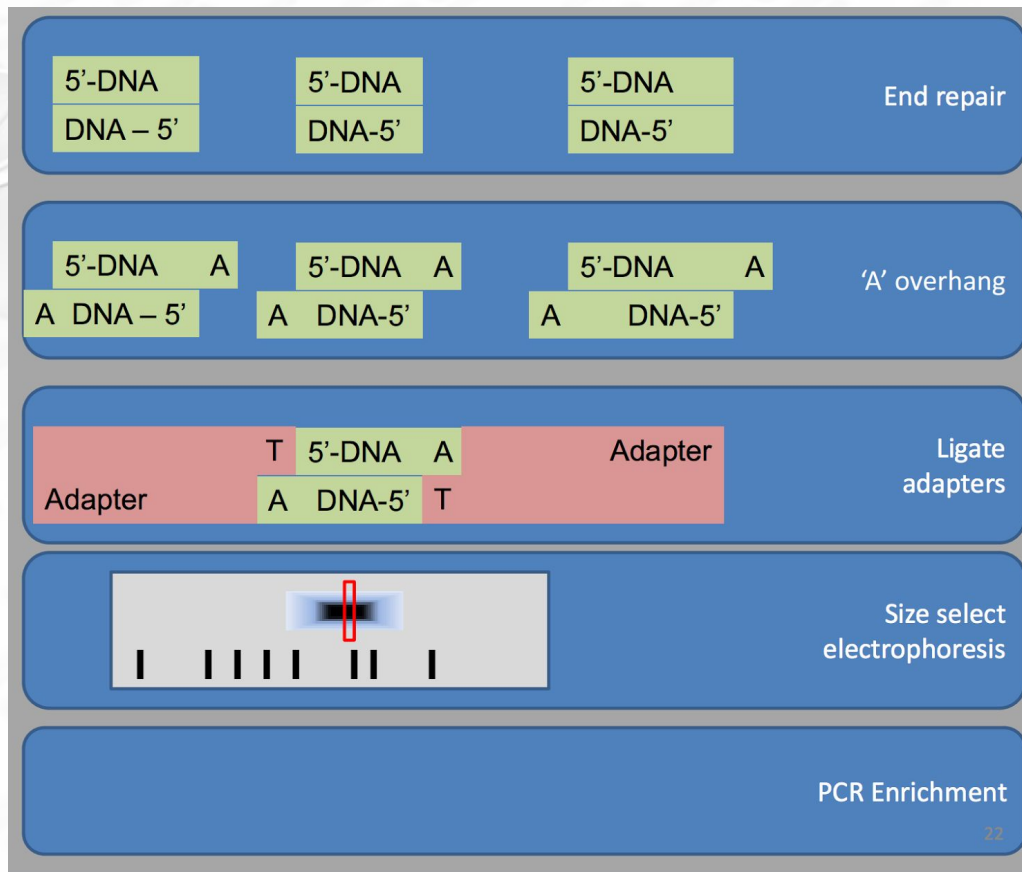Number of mRNAs/cell

**SevenBridges**

# RNA-seq library prep (1/2)

Protocols differ in :

- what types of RNA they target (total RNA, mRNA)

- on fragment sizes

- strand specificity (more info)

- bulk or single-cell

- ligation, priming...



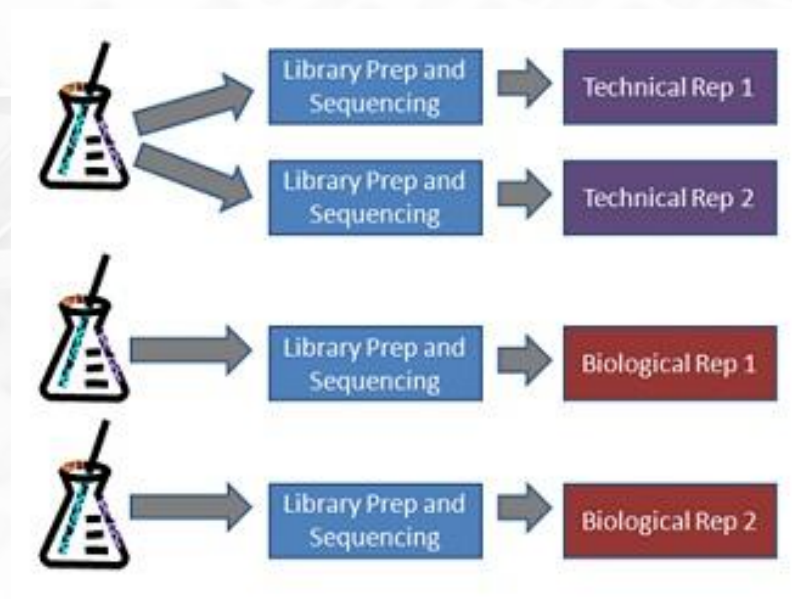TruSeq RNA

# RNA-seq library prep (2/2)

- technical or biological replicates
- Important to estimate # replicates (power analysis)



Source: http://hdl.handle.net/2345/3145

SevenBridges

# RNA-seq data analysis

- alignment, assembly, **relative abundance**, differential expression, functional enrichment analysis

- RNA-seq quantification

  vs.

  microarrays or qRT-PCR

- Latest approach -> single-cell RNA-seq (preserve information of cell origin, usually by priming all transcripts from same cell identically)

**SevenBridges**
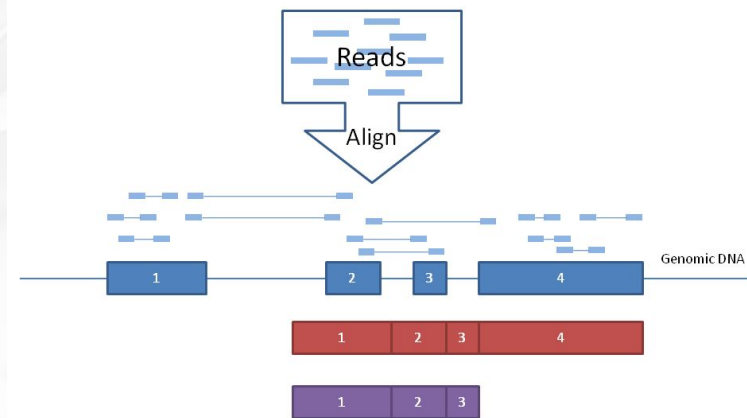
# RNA-seq data (1)
(variance and ambiguity)

- sampling variance (biological replicates)

- technical or *biological* variance

  (both technical and biological rep.)

alternative splicing:

- mapping ambiguity

  (multiple mapping)

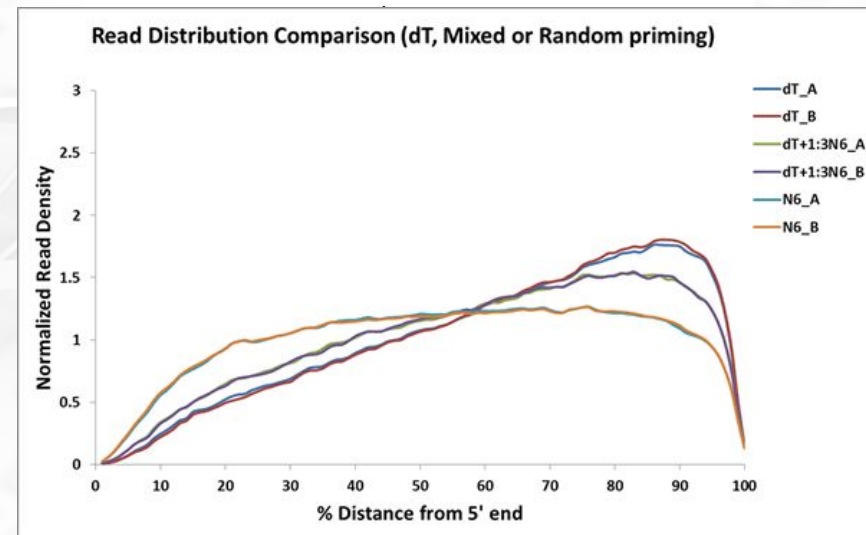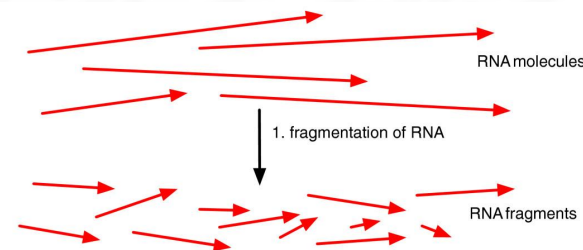Source: http://dx.doi.org/10.13070/mm.en.3.203

# RNA-seq data (2)
(biases)

- fragment length distribution

  (small fragments -> more ambiguity)

- positional and sequence-specific

  (due to priming or fragmentation)

- sequencing errors

  (error model - mismatches & indels)

Source: doi:10.1186/gb-2011-12-3-r22



RNA molecules

1. fragmentation of RNA

RNA fragments

**Read Distribution Comparison (dT, Mixed or Random priming)**

— dT_A
— dT_B
— dT+1:3N6_A
— dT+1:3N6_B
— N6_A
— N6_B

Normalized Read Density

% Distance from 5' end

paired-end read    6. mapping    sense    RNA sequence    anti-sense

SevenBridges

# RNA-seq data (3)
### (normalization)

- **within sample** normalization (transcript length + sequencing depth):
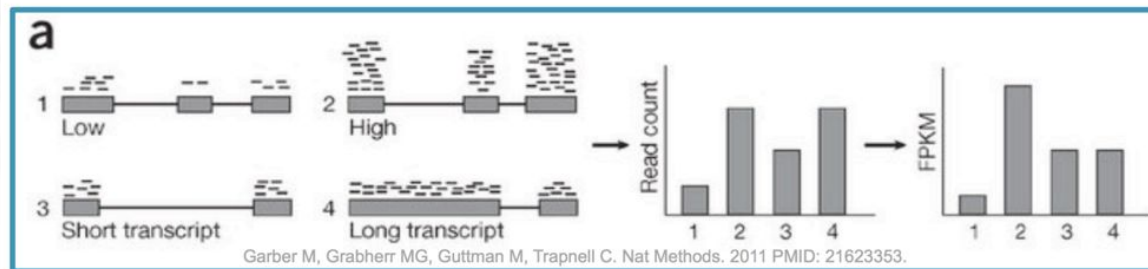
Let $X_i$ be number of reads aligned to $i$th transcript
and $l_i$ is length of $i$th transcript

$$\sum_i X_i \neq \text{expression of a gene}$$

## Adjust for (effective) length



Garber M, Grabherr MG, Guttman M, Trapnell C. Nat Methods. 2011 PMID: 21623353.

Let $N$ = total number of mapped reads
$\tilde{l}_i = l_i - \mu_{FLD} + 1$, $\mu$ is mean of fragment length dist.

$$\text{effCounts}_i = X_i \cdot \frac{l_i}{\tilde{l}_i}$$

SevenBridges

# RNA-seq data (4)

(normalization)

$$\mathbf{TPM}_i = \frac{X_i}{\widetilde{l}_i} \cdot \left( \frac{1}{\sum_j \frac{X_j}{\widetilde{l}_j}} \right) \cdot 10^6, \ \mathbf{FPKM}_i = \frac{X_i}{\left( \frac{\widetilde{l}_i}{10^3} \right) \cdot \left( \frac{N}{10^6} \right)}$$

Some approaches:

- Raw counts: Quantile normalization

- Raw counts: Thinning to the minimum library size

- Transcripts per million

- Fragments per kilobase of exon per million reads

**Further info: What the FPKM**

SevenBridges

# RNA-seq data (5)
(features for quantification)

- Exons, **transcripts** or genes
- Raw counting (aligned reads)

  vs.

  probabilistic

  aligned reads to features

  vs.

  probabilistic

  unaligned/mapped kmers of reads

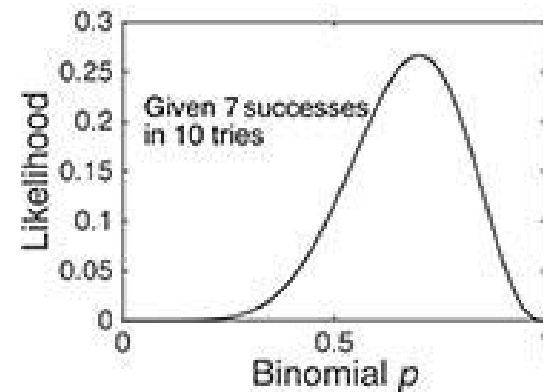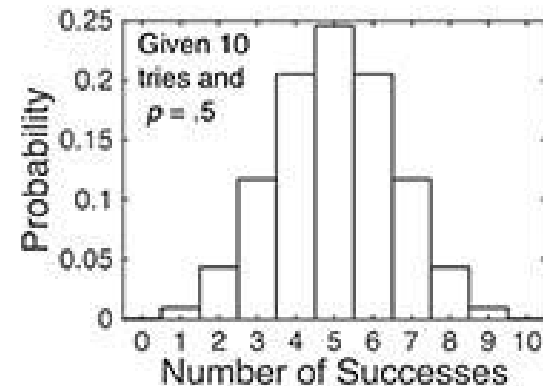HTseq counting model



of

of

SevenBridges

# Statistical background (1)

- type of problem -> counting transcripts
  (Poisson, binomial/multinomial)

- Likelihood vs. probability

If *probability* is a function of data given some params, then

*likelihood* is function of those params given the data

$$L(p; x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x}$$
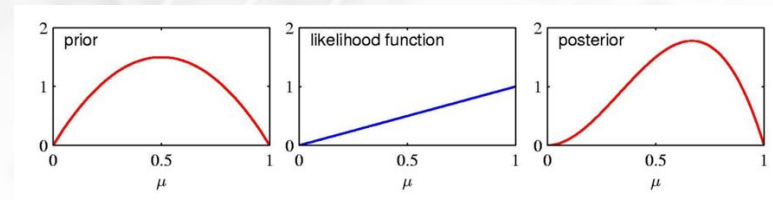
- Maximum Likelihood Estimation (MLE)

SevenBridges

# Statistical background (2)

- MLE: analytically or numerically

- Frequentist: data -> model fit-> validate model params (MLE)

- Bayesian: prior model + likelihood (given data) -> posterior model

- Bayesian interpretation closer to algorithmic approach

  (treats additional data, hidden params)

  Further info: cruncher notebook



**SevenBridges**

# MLE example (RNA)

$i = 5$ single-end, equal-length reads (a,b,c,d,e)

$k = 3$ transcripts (blue, green, red)

$\rho = (\rho_{blue}, \rho_{green}, \rho_{red})$ relative abundances of transcripts
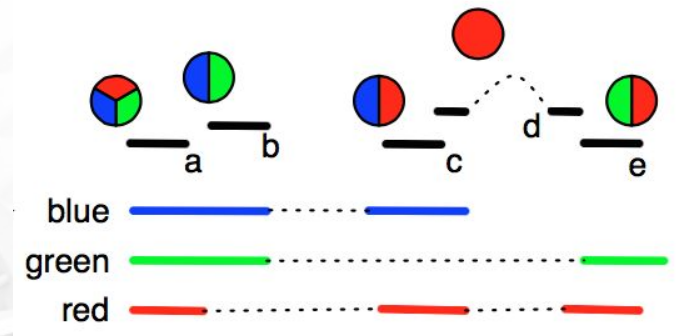
$\sum_k \rho_k = 1$, multinomial distribution

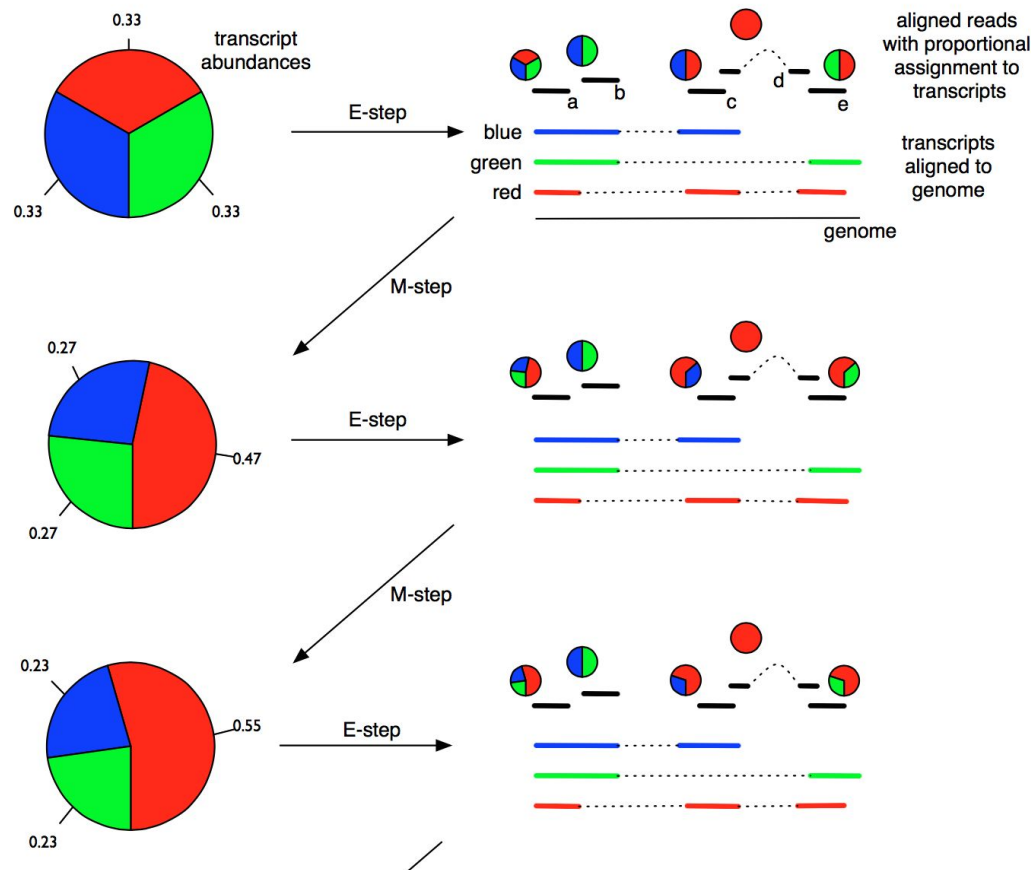$P_i = \sum_k y_{i,k} \cdot \rho_k$, probability of detecting $i$-th read

where $y_{i,k} = 1$ if $i$-th read aligns to $k$-th transcript, otherwise 0

$$L(\rho) = \prod_i \sum_k y_{i,k} \cdot \rho_k$$

Analytical solution $\rho = (0.18, 0.18, 0.64)$



Adapted from: Lior Pachter 2011, arxiv: 1104.3889v2

SevenBridges

# RNA example EM

$$(\rho_{blue}, \rho_{green}, \rho_{red}) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}), \text{ uniform prior}$$

**E1 step:** Proportional assignment
$p_a = (1/3, 1/3, 1/3), \ p_b = (1/2, 1/2, 0),$
$p_c = (1/2, 0, 1/2), \ p_d = (0, 0, 1), \ p_e = (0, 1/2, 1/2)$
**M1 step:** recalculate abundances
$\rho_{blue} = (1/3 + 1/2 + 1/2 + 0 + 0)/5 = 0.27$

**E2 step:** prior $= (0.27, 0.27, 0.46)$
$p_a = (0.27, 0.27, 0.46), \ p_b = (1/2, 1/2, 0),$
$p_c = (\frac{0.27}{0.46 + 0.27}, 0, \frac{0.46}{0.46 + 0.27}), \ p_d = (0, 0, 1), ...$
**M2 step:**
$\rho_{blue} = (0.27 + 1/2 + 0.37 + 0 + 0)/5 = 0.23$

Iterative convergance $\rho_{blue} = 0.33, 0.27, 0.23, ..., 0.18$

SevenBridges

# Raw counts implementations
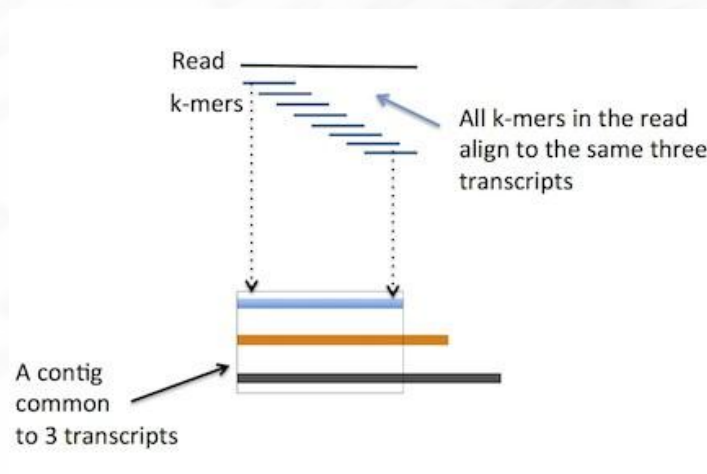
- HTseq count

- featureCounts

# EM implementations

- RSEM

- eXpress

- … EM in other areas of genomics

SevenBridges

# Pseudo-alignment methods

- Find set of transcripts that a read belongs to, BUT not exact position

- Counting kmers instead of reads

- RESULT: no local alignment, much faster EM