# SevenBridges

---

# Intro to RNA-seq

Milena Stanojevic
milena.stanojevic@sbgenomics.com

April 2021

# Beyond DNA

Recap

# Central dogma of molecular biology



Topics covered so far:

- Sequencing technologies

- DNA assembly

- DNA alignment

- DNA variants and variant calling
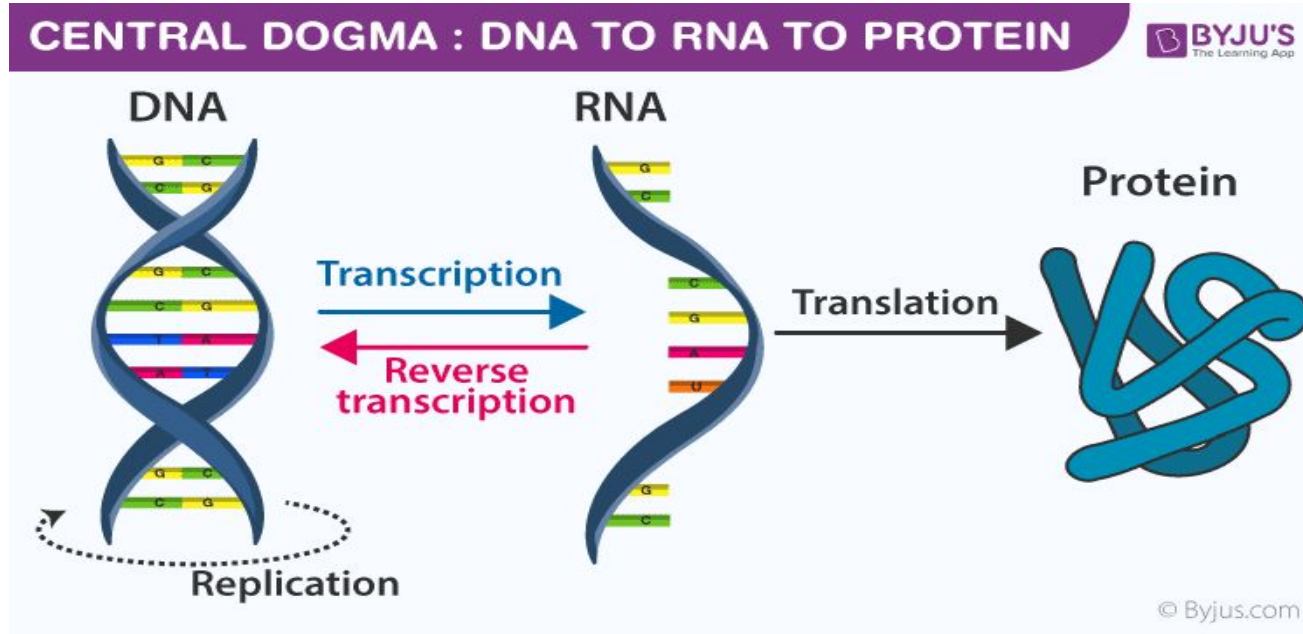
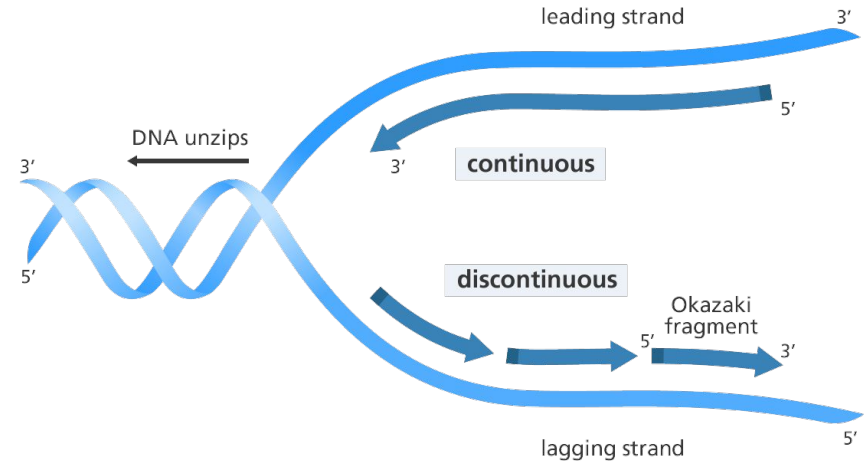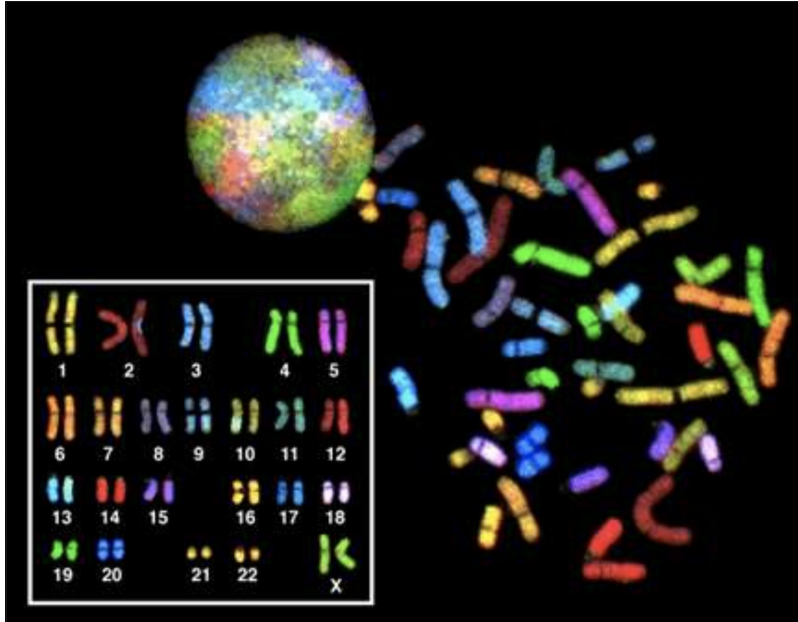# Central dogma of molecular biology



Following topics:

- **RNAs: transcription and translation, types of RNA**

- **mRNA: splicing, transcripts/isoforms**

- **RNA seq and alignment**

- **RNA quantification**

- **Differential expression**

+   Cancer bioinformatics

# Central dogma of molecular biology

# DNA replication



leading strand

DNA unzips

continuous

discontinuous

Okazaki fragment

lagging strand

Results in **all** cells in a body having same DNA
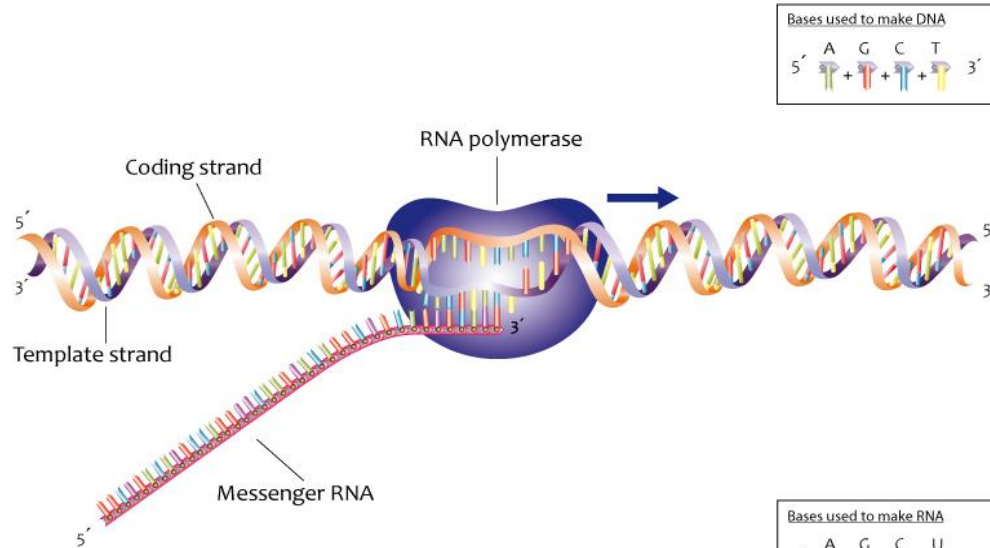
How is that we have different cells in a body?

# Transcriptomics

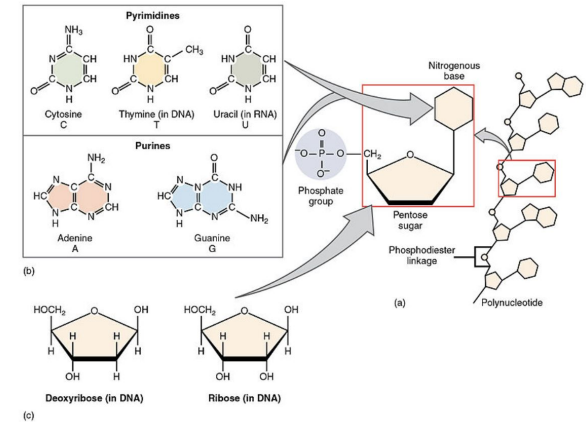Lots of RNAs, splicing, GTF, translation

Synthesis of RNAs from some relatively small genomic regions of DNA



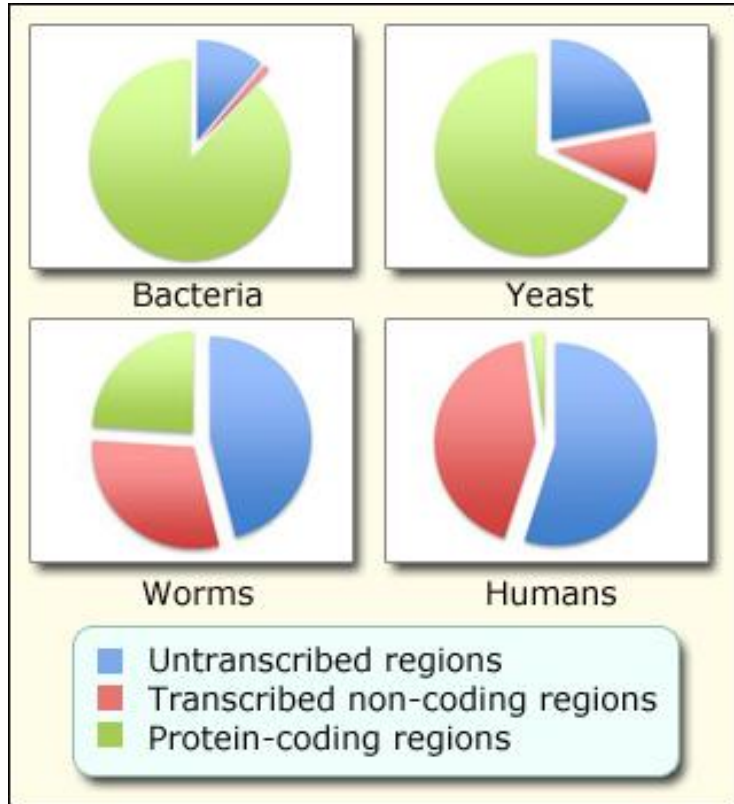**RNA transcript**
single-stranded nucleic acid

RNA vs DNA - difference?

# Transcription - how much of DNA is transcribed?



Bacteria · Yeast · Worms · Humans

Untranscribed regions
Transcribed non-coding regions
Protein-coding regions

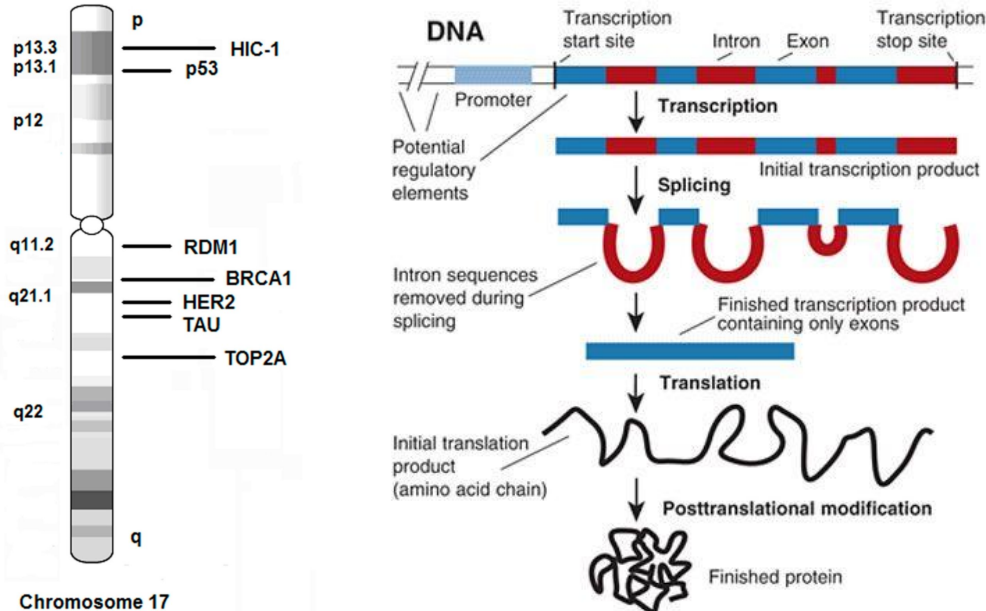**Gene** - segment of DNA which is transcribed into RNA which then has a function in cell

If RNA codes for protein that RNA is called **mRNA** and the region of genome from which it is transcribed is called **protein-coding gene** (green)

Genes which code for RNA with different functions other than protein coding - structural, regulatory, transport etc. - **non-coding genes** (red)

Some regions of DNA (most of it) are not transcribed at all (blue)

# Transcription

**Process**: Synthesis of RNAs from some relatively small genomic regions
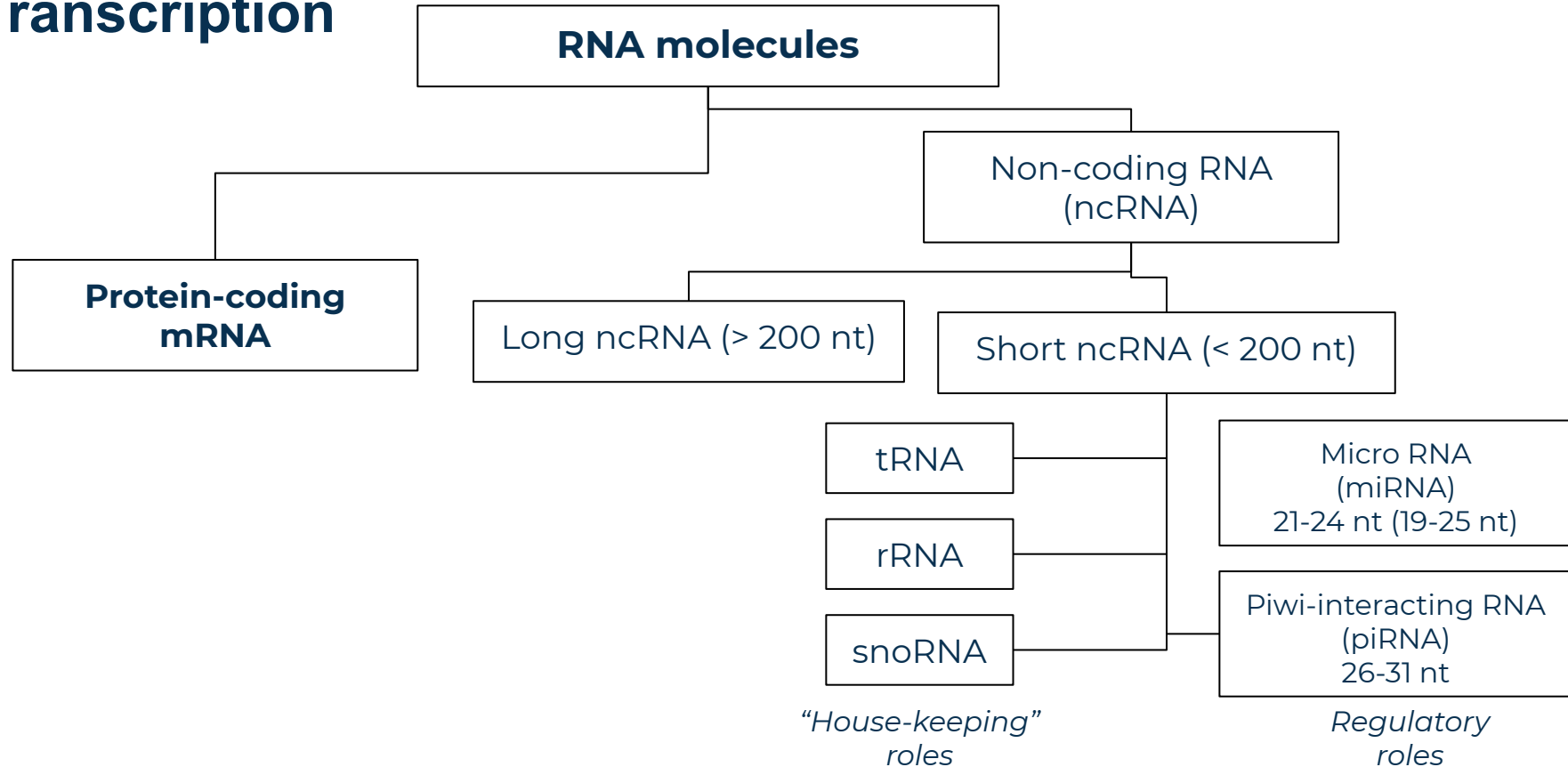


Chr 17 - protein coding genes example

Eucariotic gene structure

Human cells use splicing and other processes to make multiple proteins from the instructions encoded in a single gene
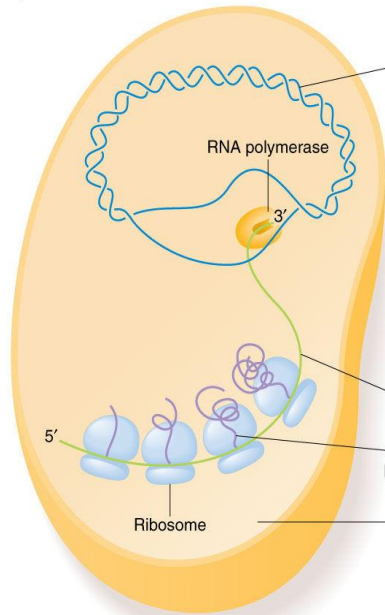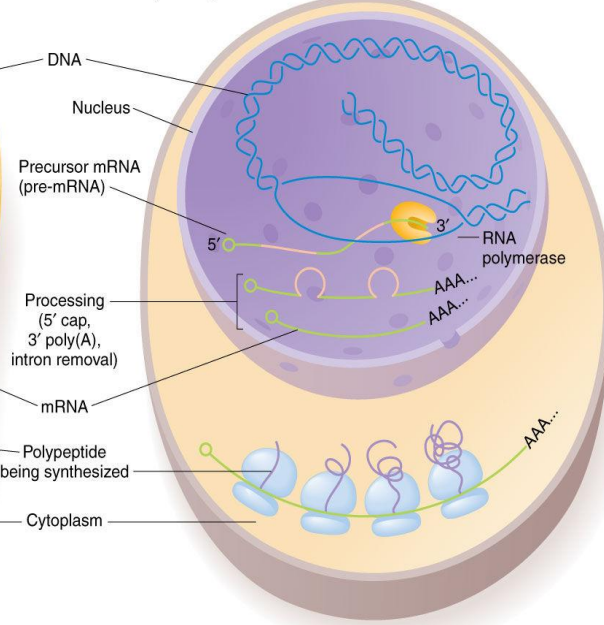
sevenbridges.com

# Transcription

```
                          ┌─────────────────────────┐
                          │      RNA molecules      │
                          └─────────────────────────┘
```

**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*

# mRNAs

## Synthesis and maturation in nucleus

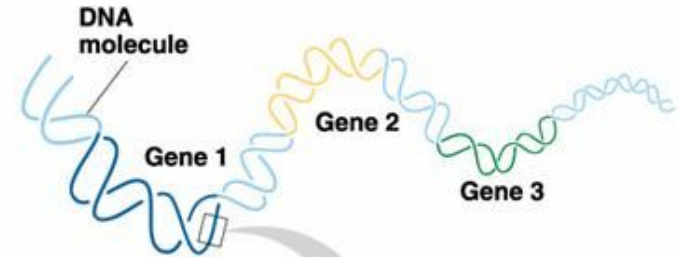# mRNAs: maturation

mature mRNA = spliced transcript + poly-A tail

# mRNAs: alternative splicing

transcripts from same gene: *isoforms*

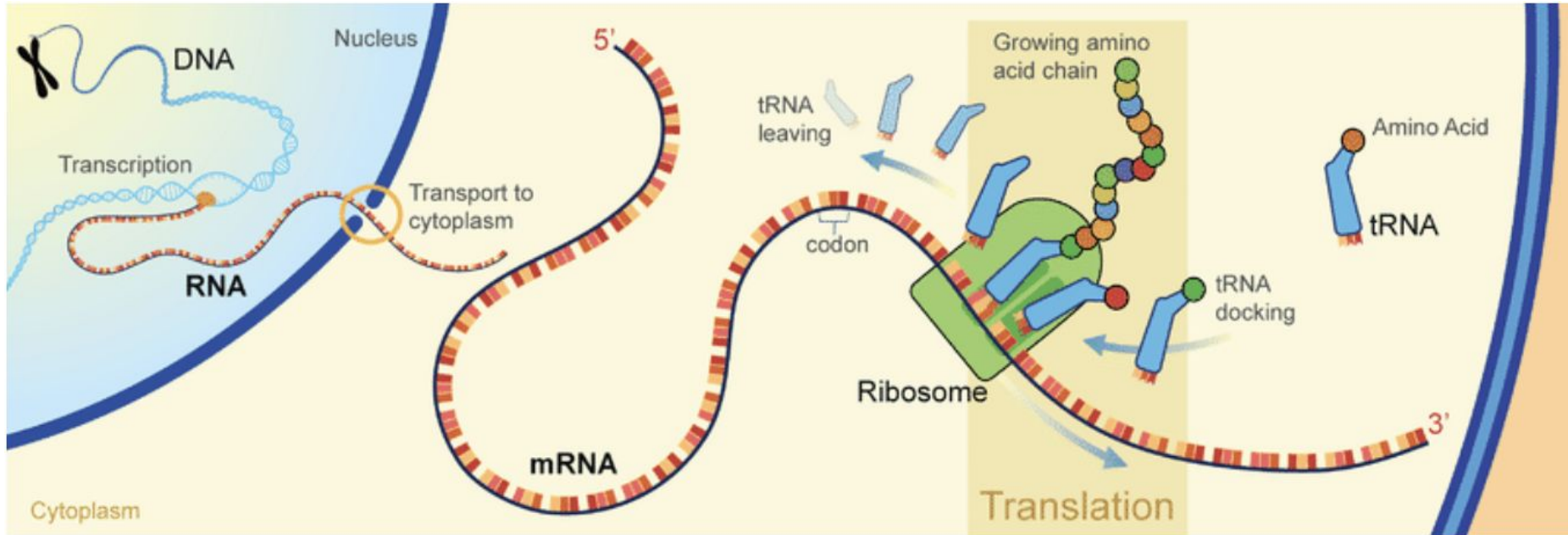# mRNAs: translation to proteins

# mRNAs: translation to proteins
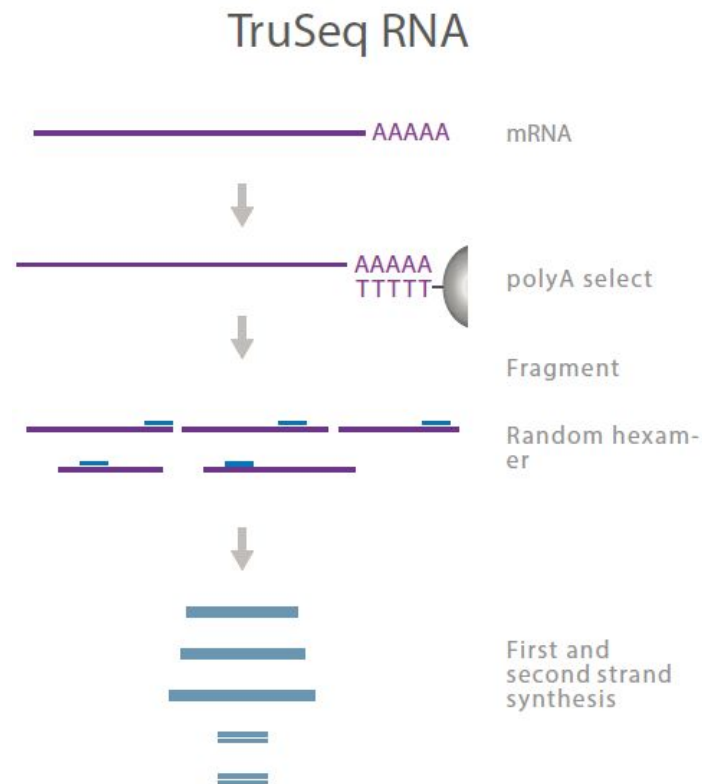
# GTF (gene transfer format)
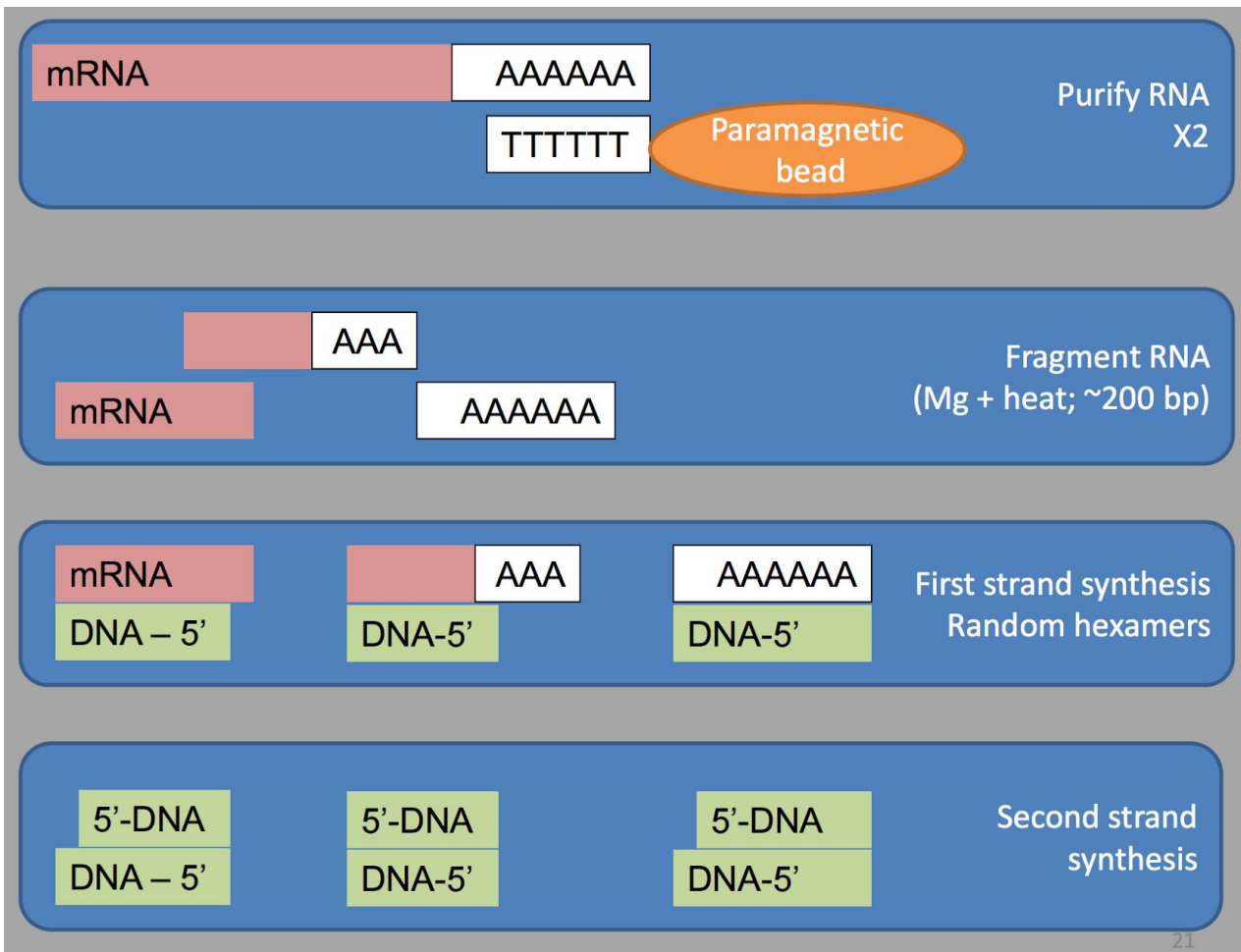
# RNA-seq

Library preparation

# RNA-seq library prep (1/2)

Protocols differ in :

- what types of RNA they target
  (total RNA, mRNA)
- on fragment sizes
- strand specificity
- bulk or single-cell



TruSeq RNA

mRNA

polyA select

Fragment

Random hexam-er

First and second strand synthesis

End repair

5'-DNA
DNA – 5'

5'-DNA
DNA-5'

5'-DNA
DNA-5'

'A' overhang

5'-DNA     A
A   DNA – 5'
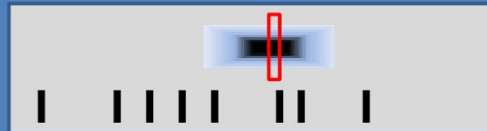
5'-DNA     A
A   DNA-5'

5'-DNA     A
A         DNA-5'

Ligate adapters

T   5'-DNA   A              Adapter
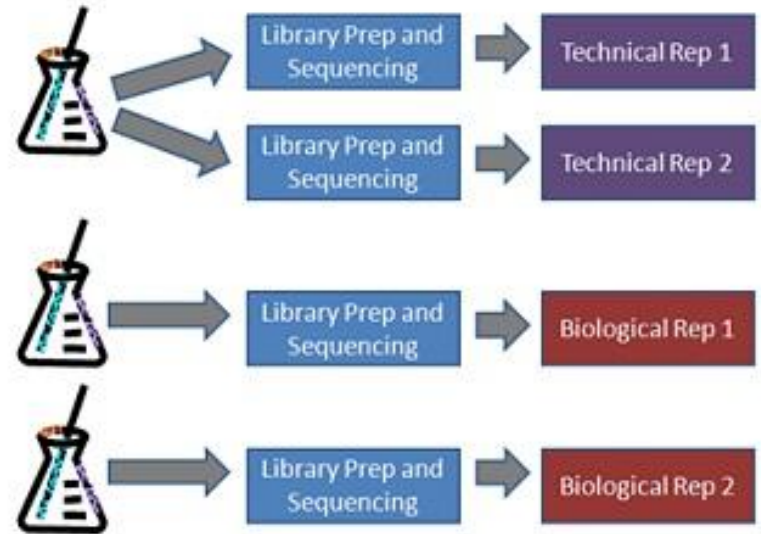Adapter      A   DNA-5'   T

Size select electrophoresis

PCR Enrichment

22

# RNA-seq library prep (2/2)

- two sources of variation in RNA-seq

- technical or biological replicates

- Important to estimate # replicates
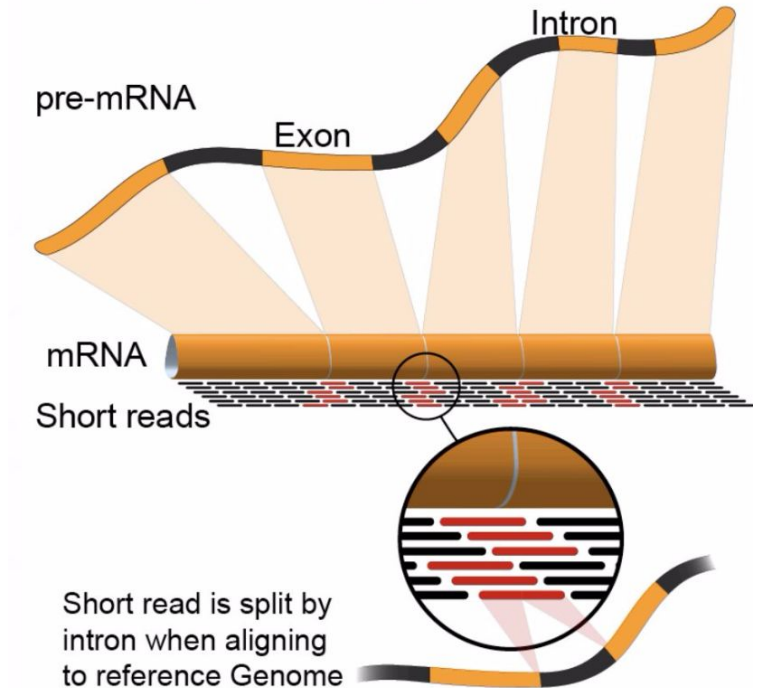  (power analysis)



Source: http://hdl.handle.net/2345/3145
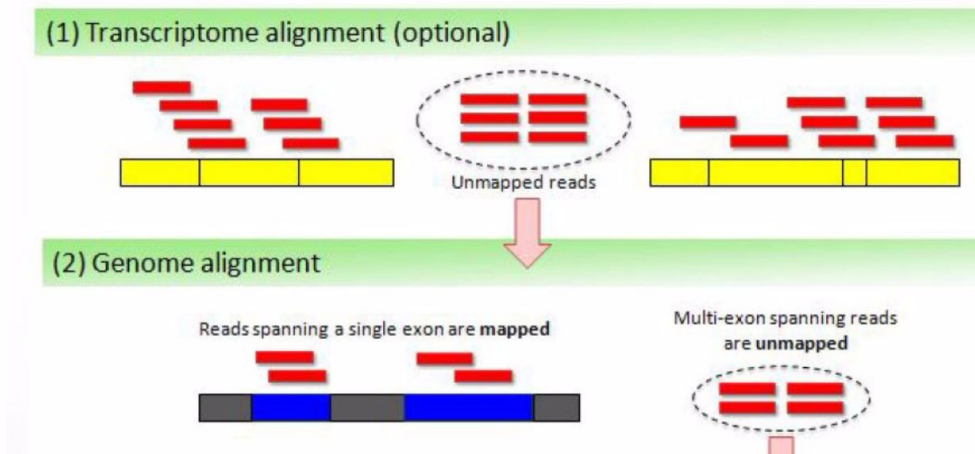
# RNA-seq

Splice-aware alignment

# Splice-aware alignment

- Average gene size ~ 10-15 kbp
- Average length of mRNA ~ 2200bp
- Average exon  ~ 230bp
- Average number of exons ~ 9.5
- For 100bp reads ~ 35% of reads would span exons
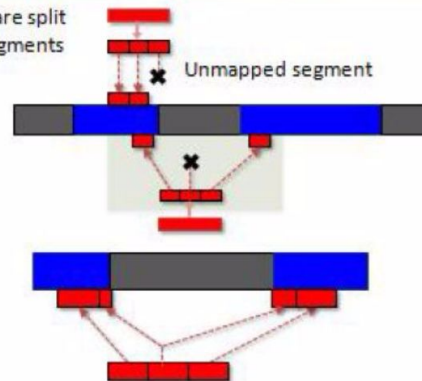
sevenbridges.com

# Splice-aware alignment

# Splice-aware alignment

sevenbridges.com

# Questions?