

RNA Sequencing & Gene Expression

02-25 I

Slides by Carl Kingsford

The Central Dogma

DNA



RNA Polymerase
(transcription)

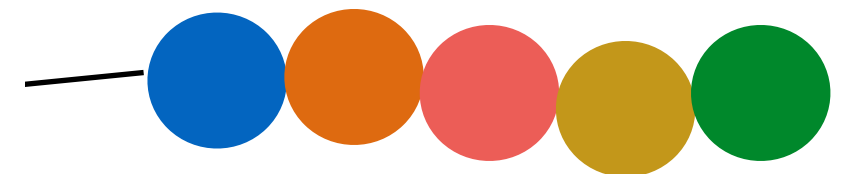
Changes in mRNA expression can
have vast effects on cellular function
“downstream”

RNA

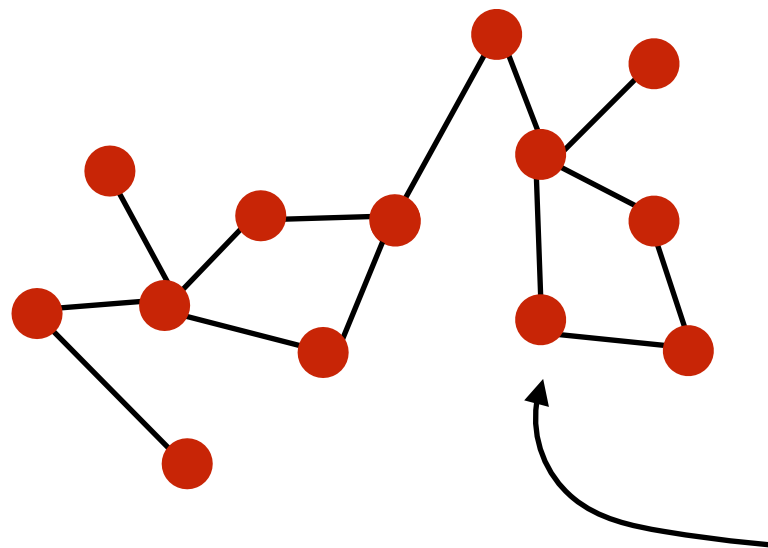


Ribosomes
(translation)

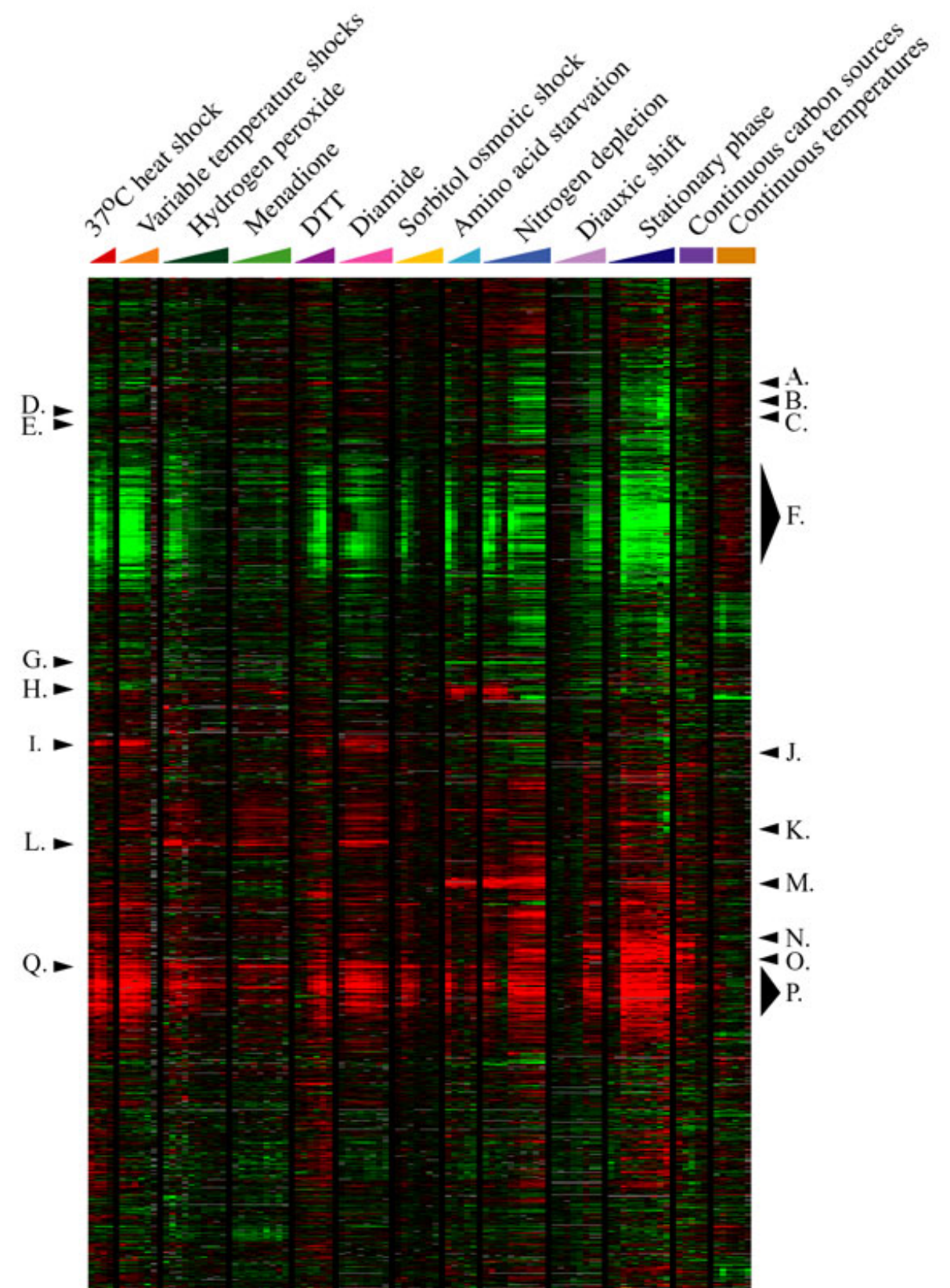
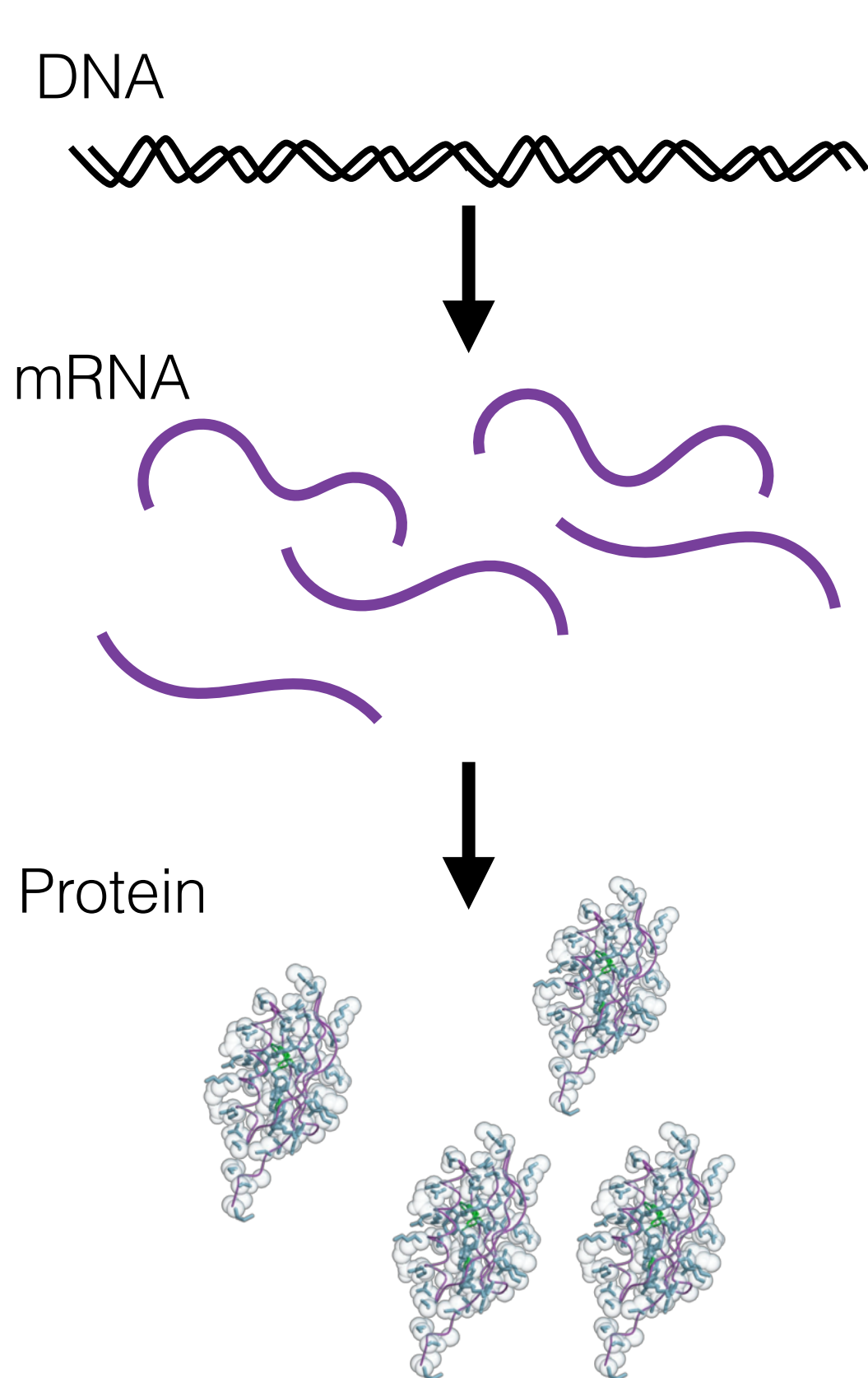
Protein



Form networks &
pathways; perform a
vast set of cellular
functions



Gene Expression Varies By Condition



The Genetic Code

- There are 20 different amino acids & 64 different codons.
- Lots of different ways to encode for each amino acid.
- The 3rd base is typically less important for determining the amino acid
- Three different “stop” codons that signal the end of the gene
- Start codons differ depending on the organisms, but AUG is often used.

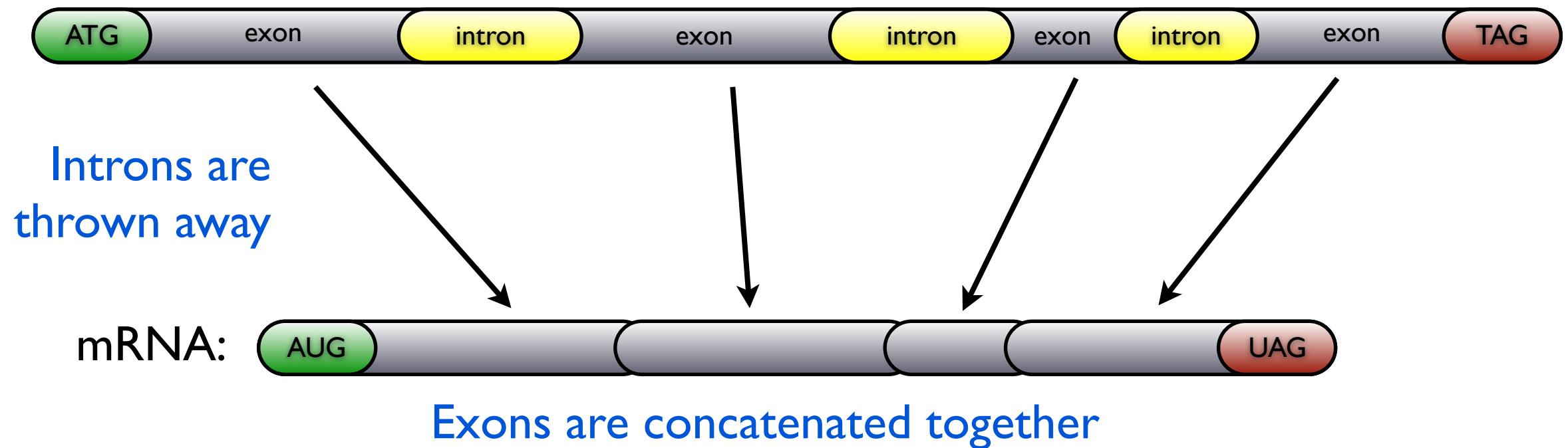
		2nd base			
		U	C	A	G
1st base	U	UUU (Phe/F) Phenylalanine	UCU (Ser/S) Serine	UAU (Tyr/Y) Tyrosine	UGU (Cys/C) Cysteine
		UUC (Phe/F) Phenylalanine	UCC (Ser/S) Serine	UAC (Tyr/Y) Tyrosine	UGC (Cys/C) Cysteine
		UUA (Leu/L) Leucine	UCA (Ser/S) Serine	UAA Ochre Stop	UGA Opal Stop
		UUG (Leu/L) Leucine	UCG (Ser/S) Serine	UAG Amber Stop	UGG (Trp/W) Tryptophan
	C	CUU (Leu/L) Leucine	CCU (Pro/P) Proline	CAU (His/H) Histidine	CGU (Arg/R) Arginine
		CUC (Leu/L) Leucine	CCC (Pro/P) Proline	CAC (His/H) Histidine	CGC (Arg/R) Arginine
		CUA (Leu/L) Leucine	CCA (Pro/P) Proline	CAA (Gln/Q) Glutamine	CGA (Arg/R) Arginine
		CUG (Leu/L) Leucine	CCG (Pro/P) Proline	CAG (Gln/Q) Glutamine	CGG (Arg/R) Arginine
	A	AUU (Ile/I) Isoleucine	ACU (Thr/T) Threonine	AAU (Asn/N) Asparagine	AGU (Ser/S) Serine
		AUC (Ile/I) Isoleucine	ACC (Thr/T) Threonine	AAC (Asn/N) Asparagine	AGC (Ser/S) Serine
		AUA (Ile/I) Isoleucine	ACA (Thr/T) Threonine	AAA (Lys/K) Lysine	AGA (Arg/R) Arginine
		AUG [A] (Met/M) Methionine	ACG (Thr/T) Threonine	AAG (Lys/K) Lysine	AGG (Arg/R) Arginine
	G	GUU (Val/V) Valine	GCU (Ala/A) Alanine	GAU (Asp/D) Aspartic acid	GGU (Gly/G) Glycine
		GUC (Val/V) Valine	GCC (Ala/A) Alanine	GAC (Asp/D) Aspartic acid	GGC (Gly/G) Glycine
		GUA (Val/V) Valine	GCA (Ala/A) Alanine	GAA (Glu/E) Glutamic acid	GGA (Gly/G) Glycine
		GUG (Val/V) Valine	GCG (Ala/A) Alanine	GAG (Glu/E) Glutamic acid	GGG (Gly/G) Glycine

Eukaryotic Genes & Exon Splicing

Prokaryotic (bacterial) genes look like this:

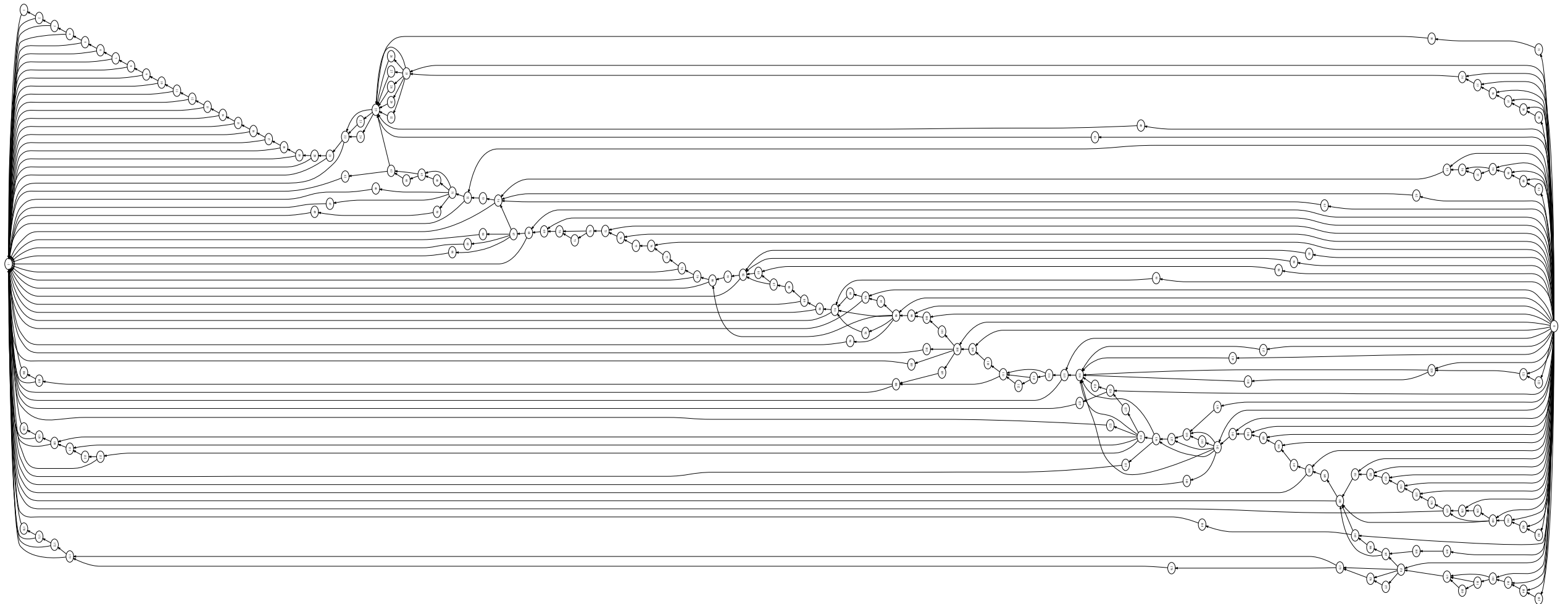


Eukaryotic genes usually look like this:



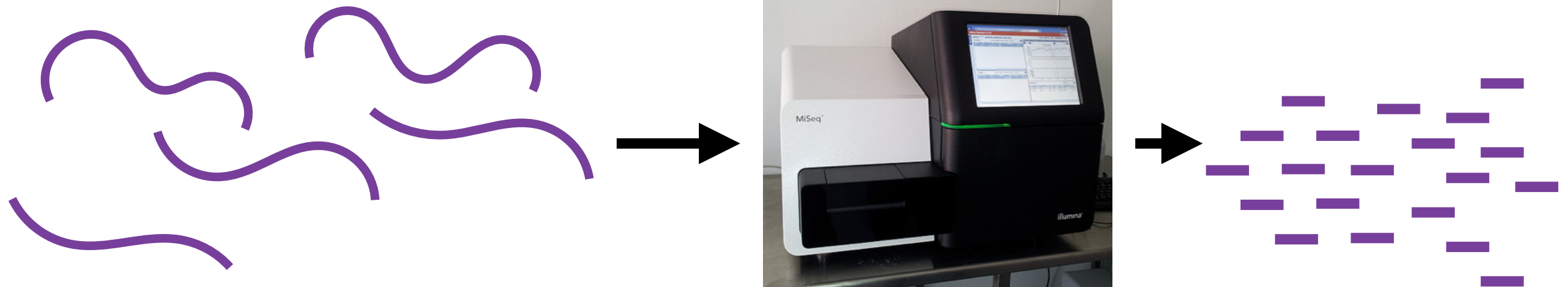
This spliced RNA is what is translated into a protein.

Alternative Splicing: It can get pretty complex...



- Splice graph of a nurexin, which is a presynaptic protein that helps to connect neurons at the synapse.
 - Node = exon (or part of an exon)
 - Edge (a,b) = sequence b can follow sequence a in some transcript

Transcriptome Sequencing



mRNA in a cell under
a given condition

Given:

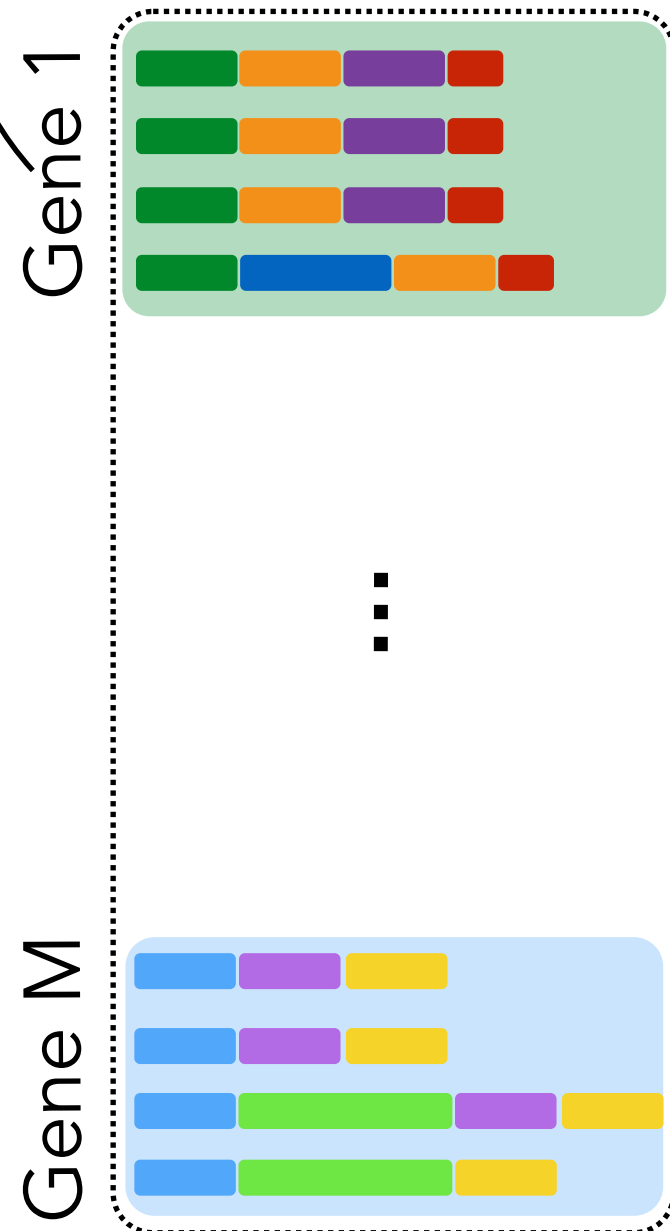
- Collection of short reads
- A set of known transcript sequences

Estimate:

- The relative abundance of each transcript

Transcript Quantification: An Overview

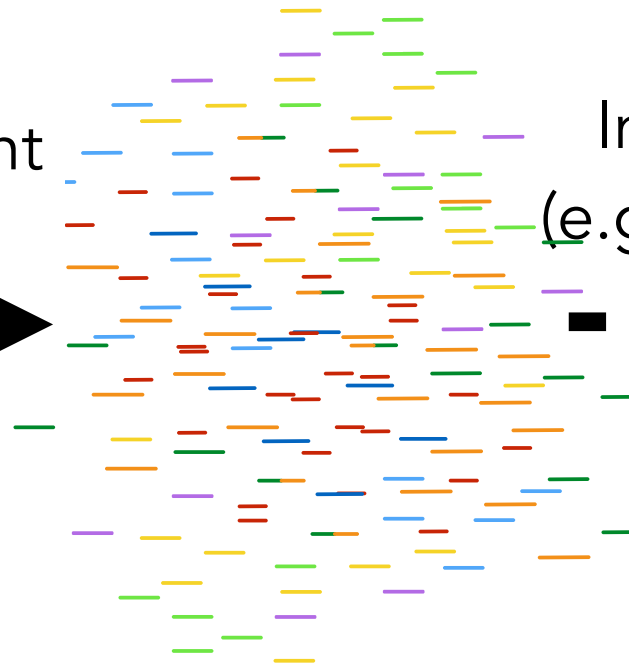
1 gene \Rightarrow many variants (isoforms)



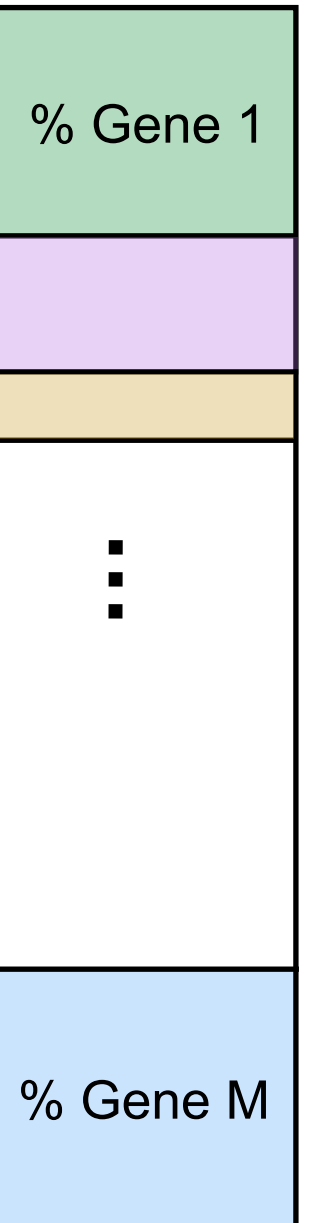
Sample

Measurement
(RNA-seq)

10s-100s of millions of
short (35-300 character) "fragments"



Inference
(e.g. Salmon)



isoform A

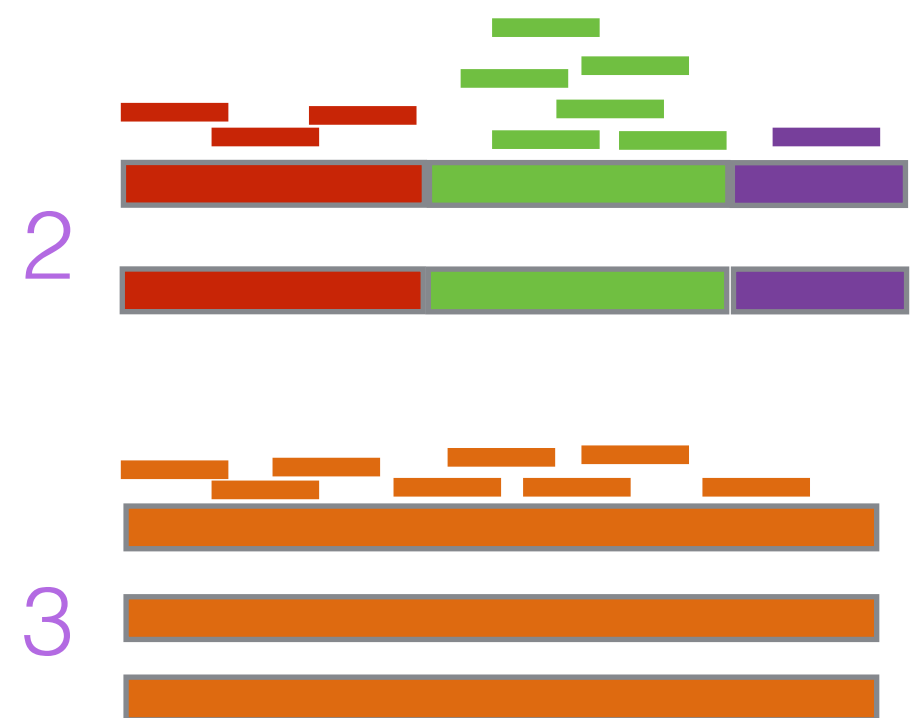
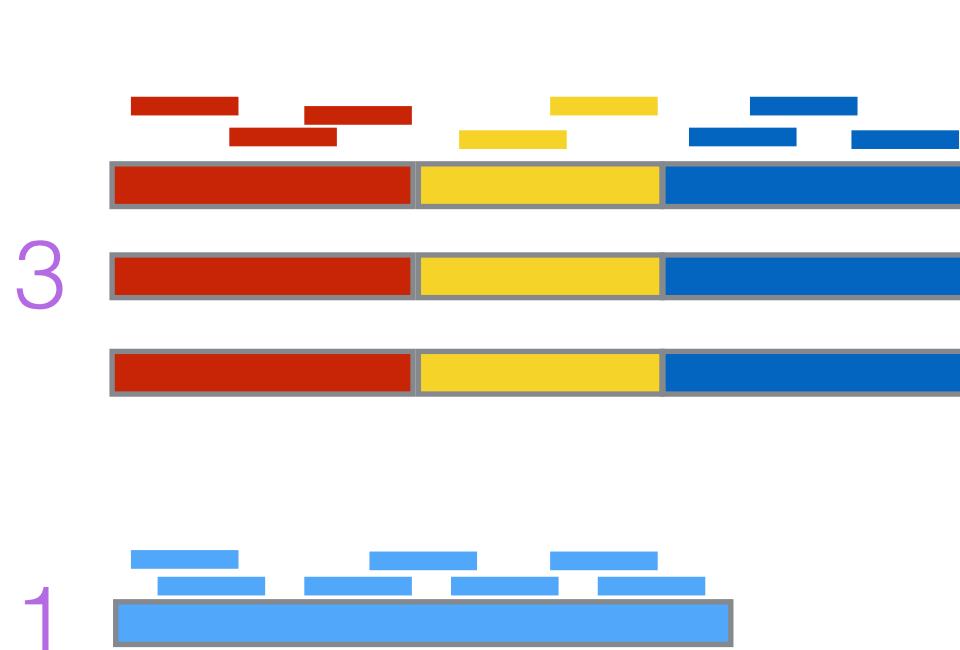
isoform B

isoform C

Abundance Estimates

Alternative Splicing is the Main Challenge

Goal: estimate the **abundance** of each kind of transcript given short reads sampled from the expressed transcripts.

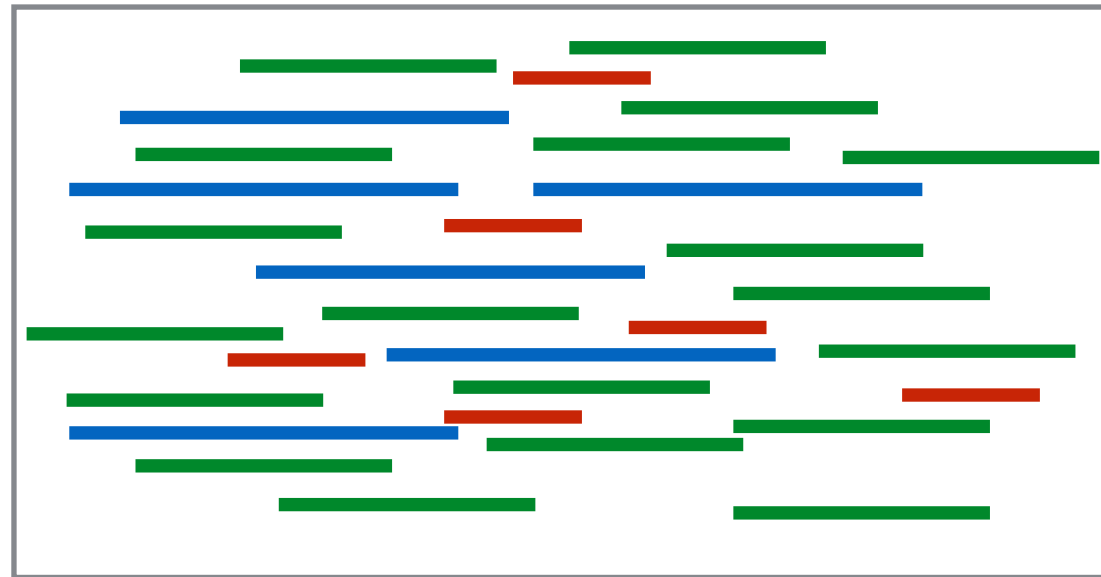



Challenges:

- hundreds of millions of short reads per experiment
- finding locations of reads (mapping) is traditionally slow
- **alternative splicing** creates ambiguity about where reads came from
- **sampling of reads is not uniform**


Inference Problem

Experimental
mixture:



length() = 100 x 6 copies = 600 nt ~ 30% blue

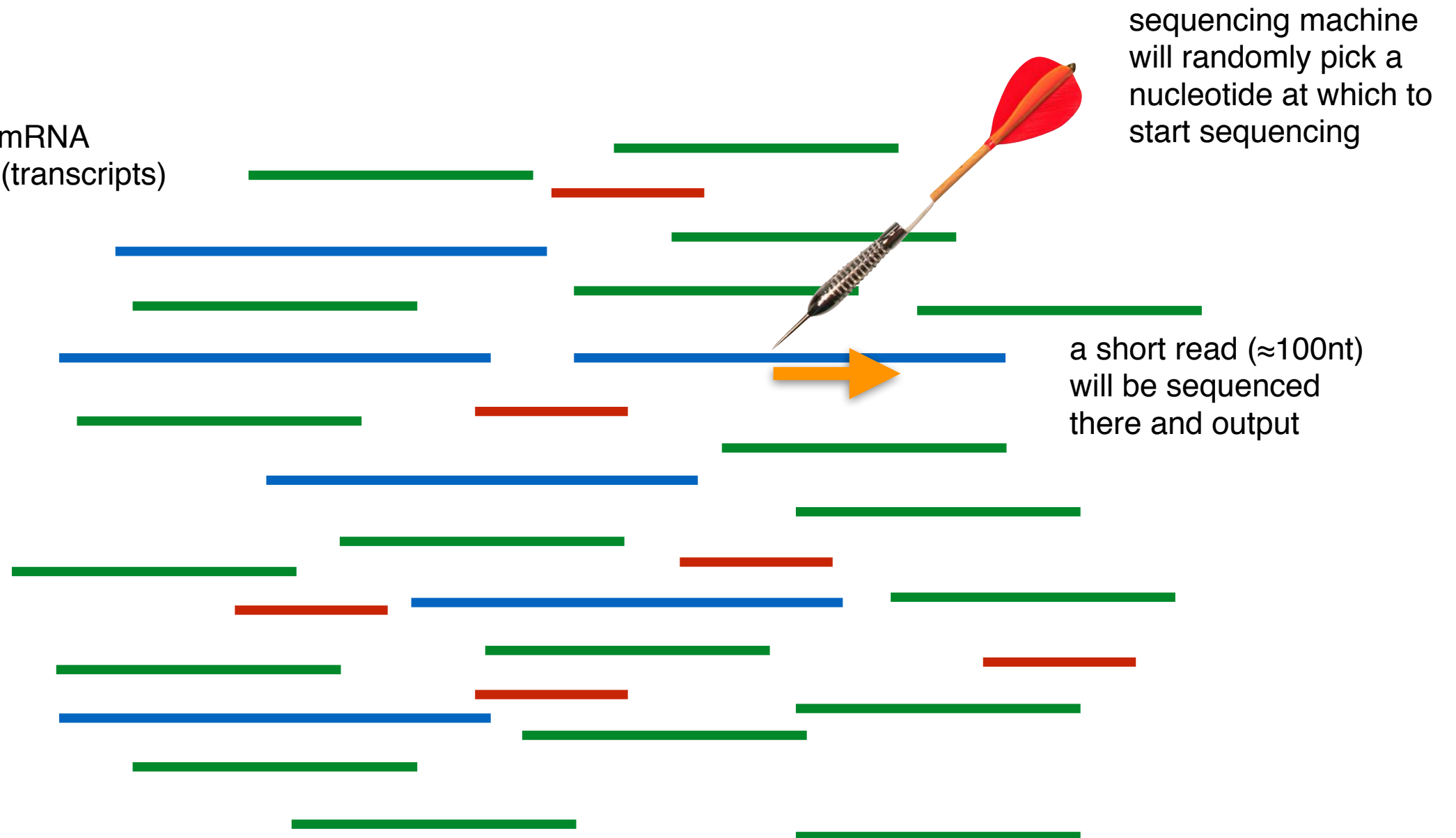
length() = 66 x 19 copies = 1254 nt ~ 60% green

length() = 33 x 6 copies = 198 nt ~ 10% red

These values $\eta = [0.3, 0.6, 0.1]$ are the *nucleotide fractions*;
they are the quantities we want to infer

Model of Sequencing

Mixture of mRNA
molecules (transcripts)
in the cell



Motif Finding \leftrightarrow Gene Expression

Gene Expression

- **Have:** reads
- **Want:** abundance vector (model)
- **Hidden data:** which transcript each read came from

$x, \text{"data"}, f_i$

η or M

z

Motif Finding

- **Have:** sequences
- **Want:** sequence profile matrix (model)
- **Hidden data:** where the motifs start in each sequence

\Rightarrow perhaps a similar EM approach will work

To Apply EM

- Recall the main EM equation:

$$\underbrace{\mathbb{E}_z \log \Pr(\text{data} \mid M_{\text{new}})}_{\substack{\uparrow \\ Q(M_{\text{new}} \mid M_{\text{old}})}} = \sum_z \underbrace{\Pr(z \mid \text{data}, M_{\text{old}})}_{\substack{\text{in expectation over choices of hidden variables (drawn} \\ \text{according to the old model)}}} \underbrace{\log \Pr(\text{data}, z \mid M_{\text{new}})}_{\text{quality of new model}}$$

data	sequenced fragments
M (model)	transcript abundances (η)
z	which transcript each fragment came from

- Let's write $\Pr(\text{data} \mid \text{model})$ in terms of gene expression notation (next slide)

Gene Expression Inference

- **Want:** expression $\eta(t)$ for each transcript t .
- **Observed:** Sequence fragments f_i sampled from molecules in the cell.

- **Hidden variables:**

$$z_{ti} = \begin{cases} 1 & \text{if fragment } i \text{ came from transcript } t \\ 0 & \text{otherwise} \end{cases}$$

indicator variable: if we knew this, answer would be easy to compute

- **Goal:**

$$\underset{\eta}{\text{maximize}} \quad g(\eta) = \Pr(\{f_i\} \mid \eta) = \sum_z \Pr(\{f_i\}, z \mid \eta)$$

observed data

unknown model

hidden variables

green box is our focus

Why do we introduce these hidden variables?

- Computing $\Pr(\{f_i\} \mid \eta)$ directly is complicated
- Once we introduce z , we have, by conditional probability:

green box
from last
slide

$$\Pr(\{f_i\}, z \mid \eta) = \Pr(z \mid \eta) \Pr(\{f_i\} \mid z, \eta)$$

- and we only need to compute:

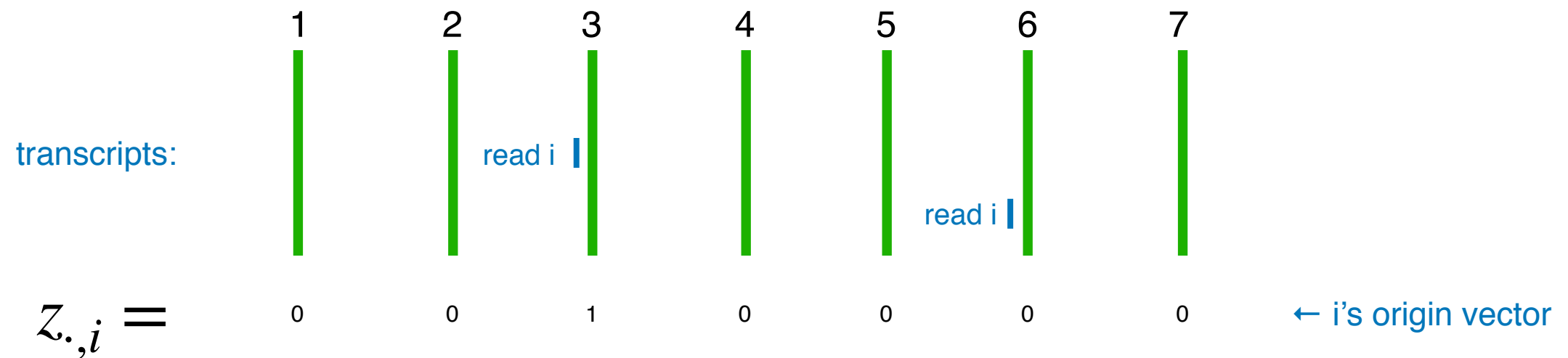
$$\Pr(z \mid \eta)$$

Probability of picking transcript, given abundances

$$\Pr(\{f_i\} \mid z, \eta)$$

Probability of generating a fragment, given
where it came from

Treat Each Sequence Read Independently



$$\Pr(z \mid \eta) = \prod_i \Pr(z_{.,i} \mid \eta)$$

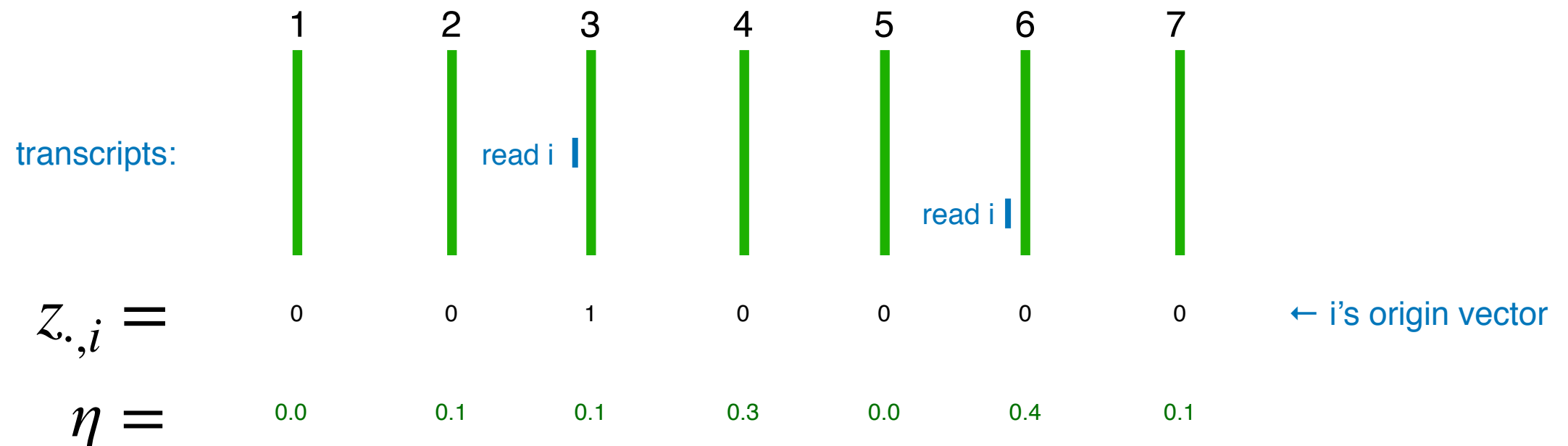
← probability of fragment i's origin vector

$$= \prod_i \sum_t \Pr(z_{ti} = 1 \mid \eta)$$

← because each read has only 1 origin, we can just ask which z_{ti} is 1 in its origin vector.

↑ Probability that read i came from transcript t given the abundances

But this is "easy":



$$\Pr(z_{ti} = 1 \mid \eta) = \frac{\eta(t)}{\sum_q \eta(q)}$$

← q sums over the transcripts where the read mapped

In example above:

$$\Pr(z_{3i} = 1 \mid \eta) = 0.1 / (0.1 + 0.4) = 0.2$$

$$\Pr(z_{6i} = 1 \mid \eta) = 0.4 / (0.1 + 0.4) = 0.8$$

$\Pr(z_{2i} = 1 \mid \eta)$ because i didn't map to transcript 2

What about the other probability?

Again, treat each fragment independently:

$$\Pr(\{f_i\} \mid z, \eta) = \prod_i \Pr(f_i \mid z_{\cdot,i}, \eta)$$

When looking at fragment i ,
only i 's origin vector matters

Since the “ z ” is given (after the \mid in the probability), we can assume we know which transcript f_i came from.

Need to compute only:

$$\Pr(f_i \mid z_{ti} = 1, \eta) = \begin{cases} 0 & \text{if } f_i \text{ doesn't map to transcript } t \\ \text{probability of generating read } f_i \text{ from the sequence of transcript } t & \text{(see next slide)} \end{cases}$$

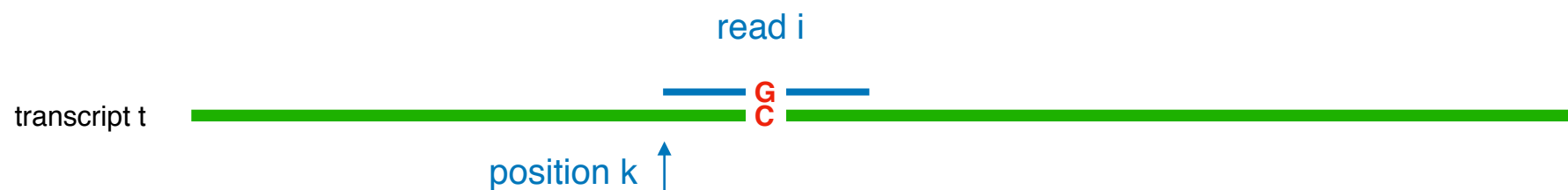
where you can encode
“intelligence” about the
experimental process

Fragment Generation Probability

Fragment generation probability is independent of abundances *once you know which fragment the transcript came from*:

$$\Pr(f_i \mid z_{ti} = 1, \eta) = \Pr(f_i \mid z_{ti} = 1)$$

Can estimate the above probability using various terms that model the experiment:



$$\Pr(f_i \mid z_{ti} = 1) \approx$$

Prob. of seeing
a fragment of
length $\text{len}(f_i)$

x

Prob. of
sequencing a
fragment at
position k

x

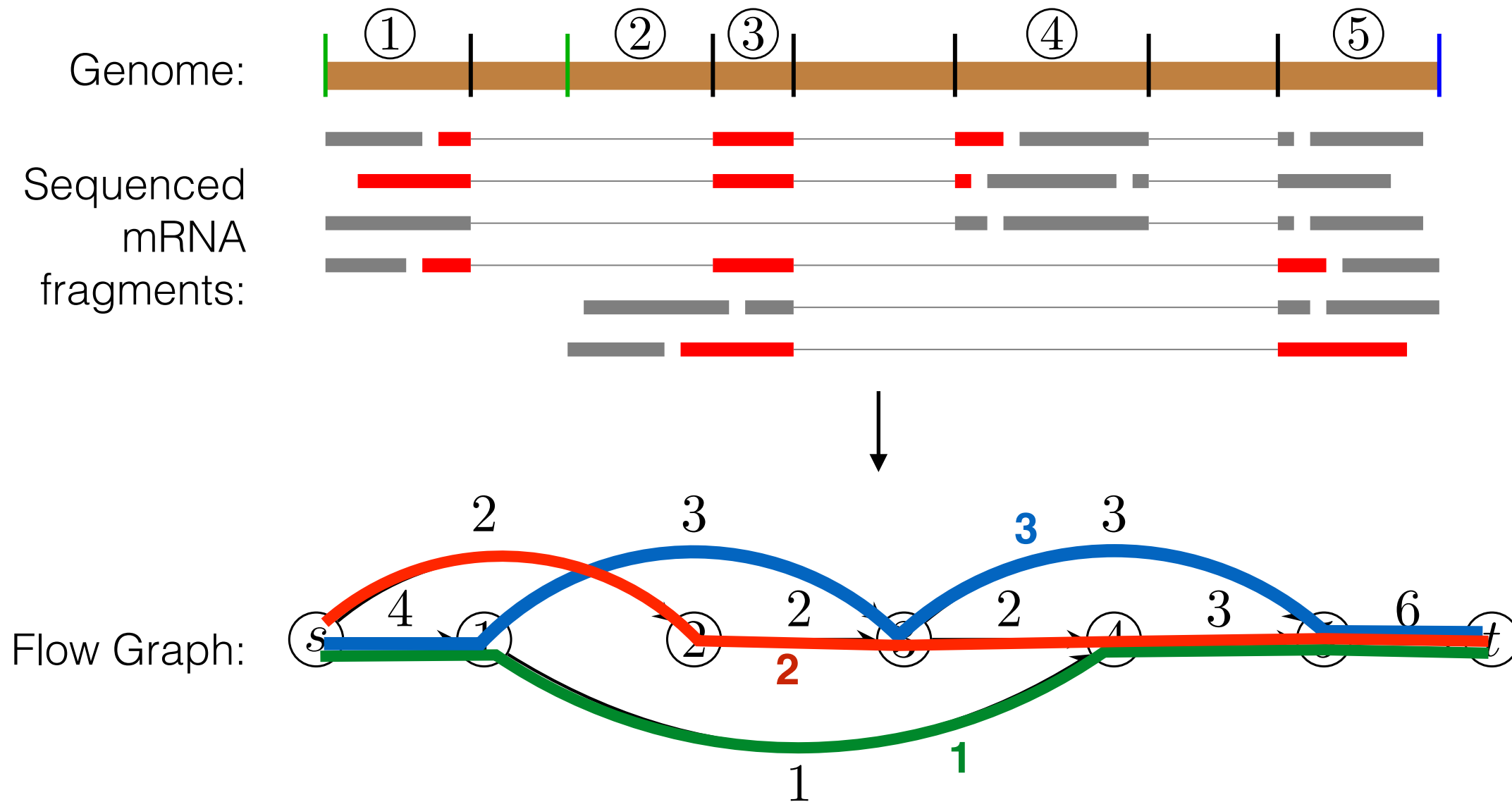
Prob. of seeing so
many differences
between read and
transcript

Expression Quantification EM Summary

- Want: expression vector; missing knowledge of where sequenced fragments came from
- Solution: add “hidden” (aka latent) variables z , and estimate $\Pr(\text{fragments}, z \mid \text{model})$.
- Do this by breaking it into 2 parts:
 $\Pr(z \mid \text{model}, \text{fragments})$
 $\Pr(\text{fragments} \mid z, \text{model})$
each of which is easier to estimate
- Once we can compute the above probabilities, we can apply EM.

Transcript Assembly

Sequencing Isoforms



Mingfu Shao and Carl Kingsford. [Scallop Enables Accurate Assembly Of Transcripts Through Phasing-Preserving Graph Decomposition](#). *Nature Biotechnology* (2017)