

Gene Architecture & Gene Discovery

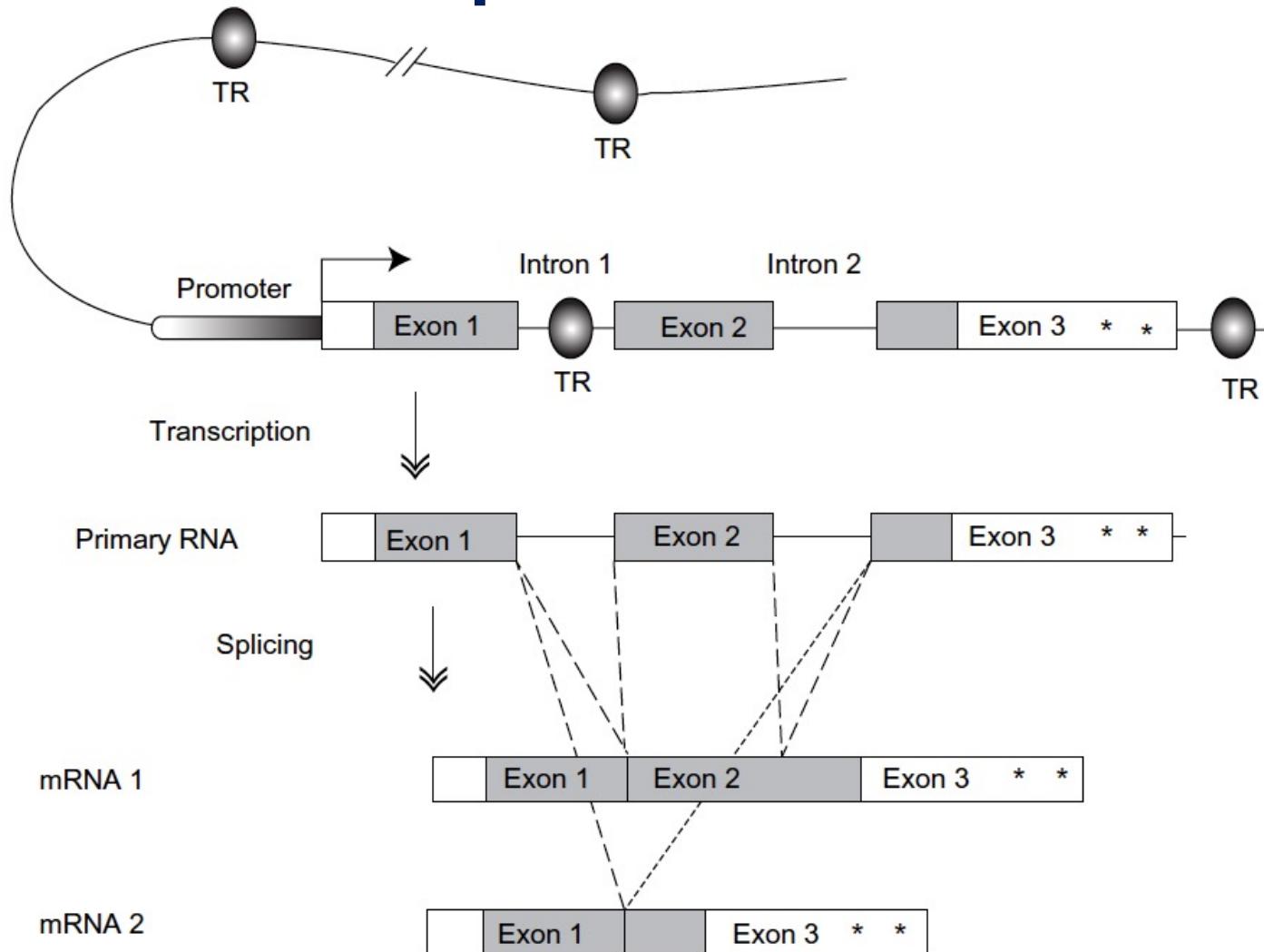


Download Gene Discovery
Files from GitHub!

BIOL 435/535:
Bioinformatics
Feb 1, 2022

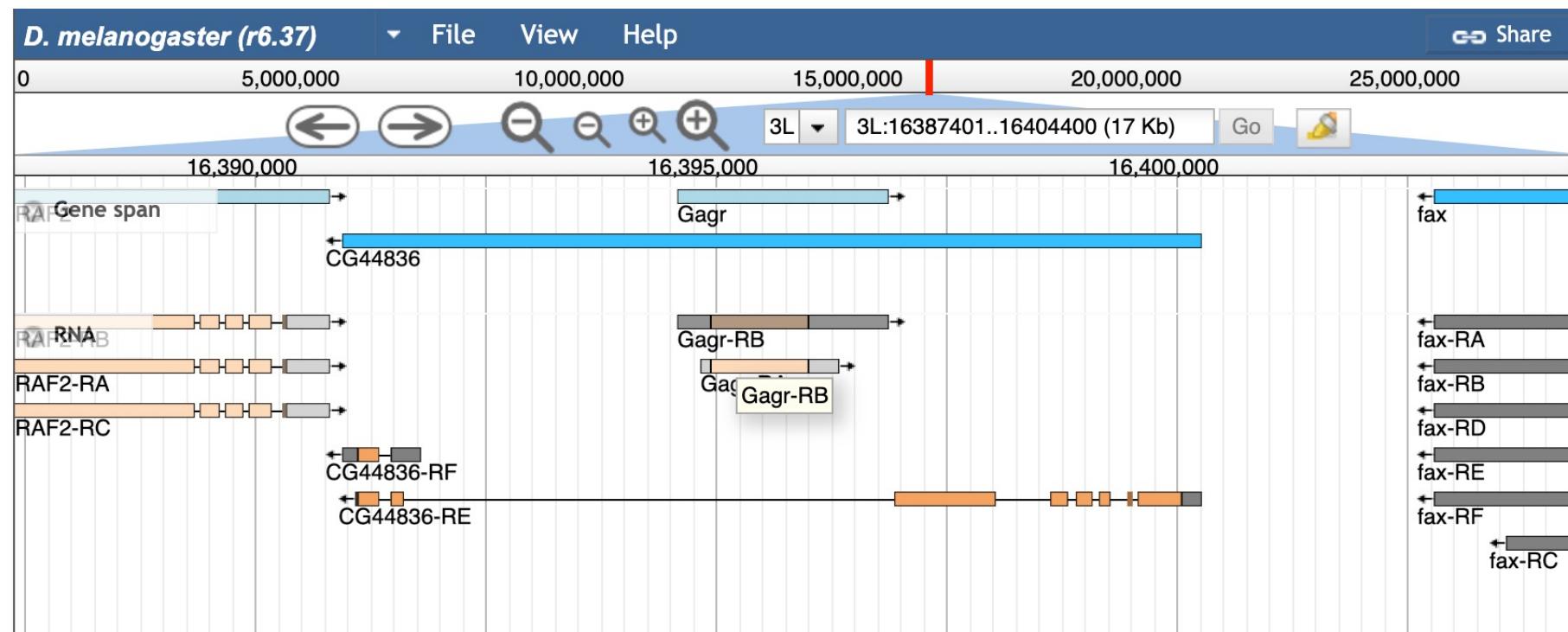
What is a gene?

Kinds of genes – protein coding



Kinds of genes – protein coding

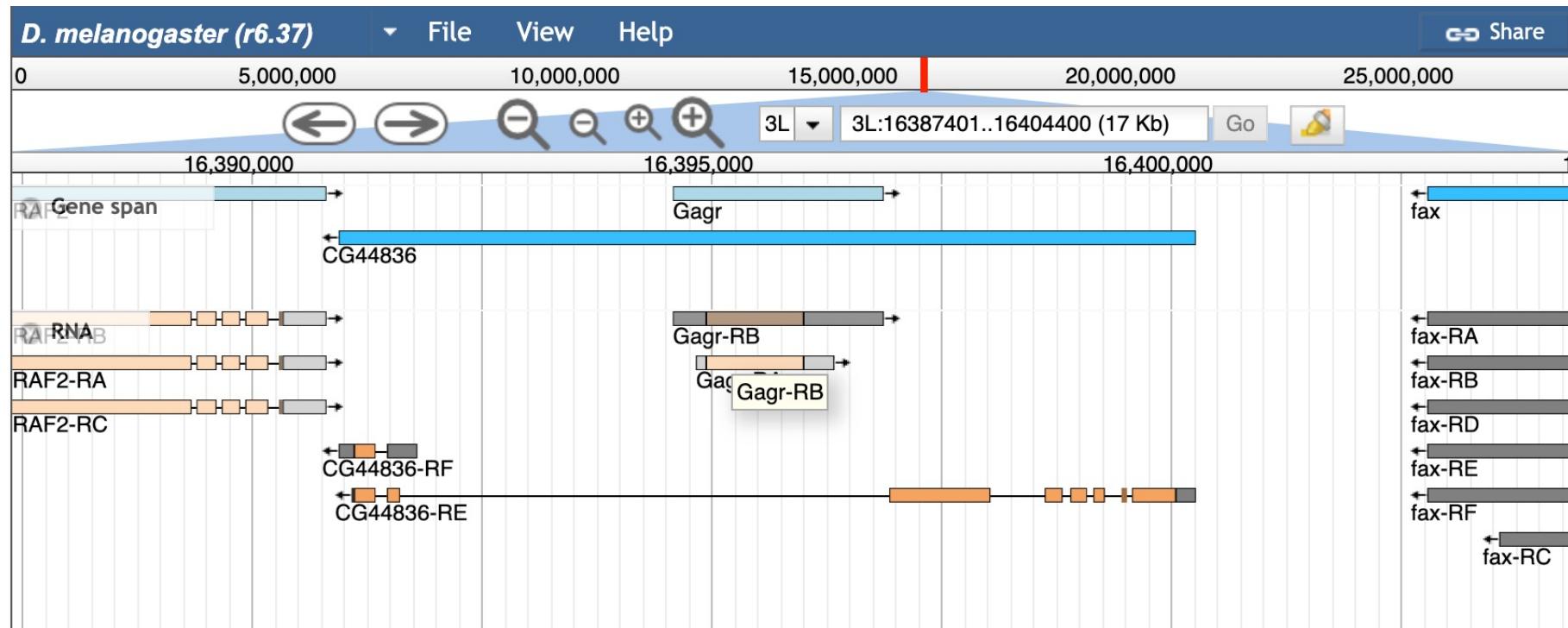
- Genes can be encoded on **both strands** of DNA (i.e., forward and reverse)



Reverse
Complement

Kinds of genes – protein coding

- Genes can be overlapping



Kinds of genes – protein coding

- Redundant, unambiguous genetic code

RNA codon table		2nd position				3rd position	
1st position		U	C	A	G		
U	Phe Phe Leu Leu	Ser Ser Ser Ser	Tyr Tyr stop stop	Cys Cys stop Trp	U C A G		
C	Leu Leu Leu Leu	Pro Pro Pro Pro	His His Gln Gln	Arg Arg Arg Arg	U C A G		
A	Ile Ile Ile Met	Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	U C A G		
G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	U C A G		
Amino Acids							
Ala: Alanine Arg: Arginine Asn: Asparagine Asp: Aspartic acid Cys: Cysteine		Gln: Glutamine Glu: Glutamic acid Gly: Glycine His: Histidine Ile: Isoleucine		Leu: Leucine Lys: Lysine Met: Methionine Phe: Phenylalanine Pro: Proline		Ser: Serine Thr: Threonine Tyr: Tyrosine Val: Valine	

Expasy

Intron/exon boundaries

- Exons often end with **AG**
- Introns often begin with **GT** end with **AG**

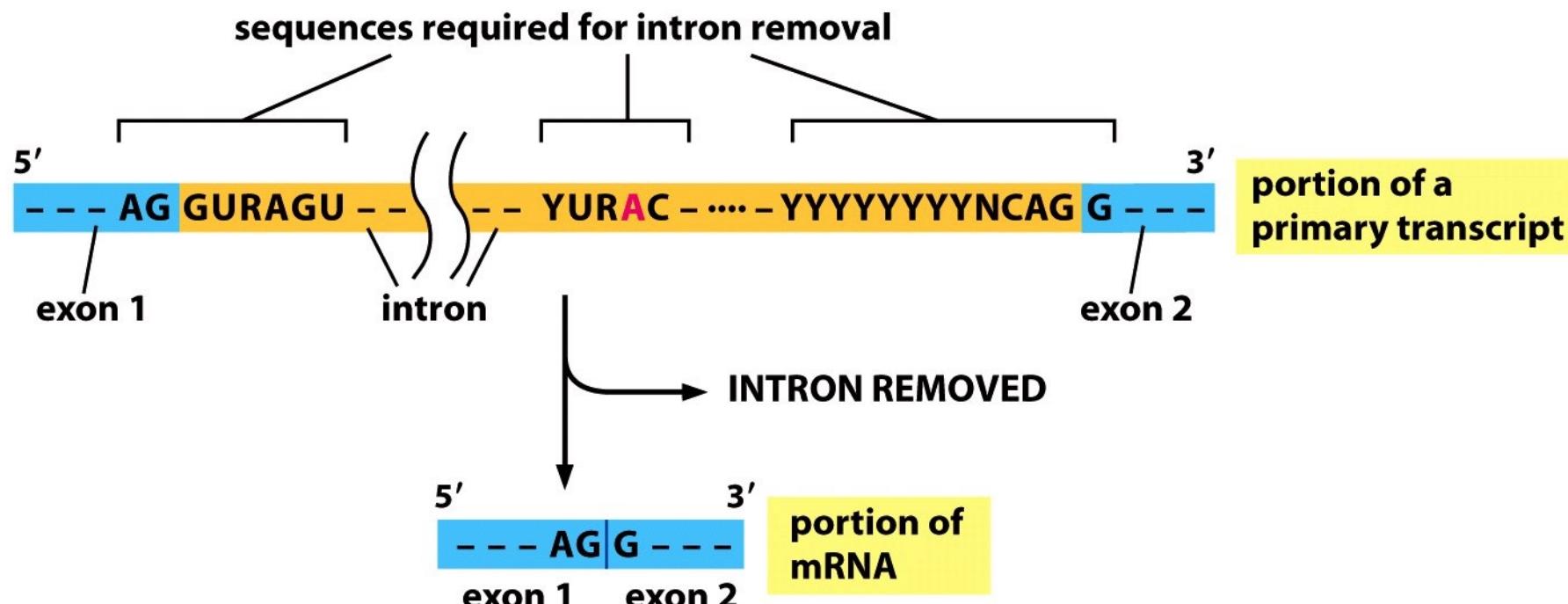
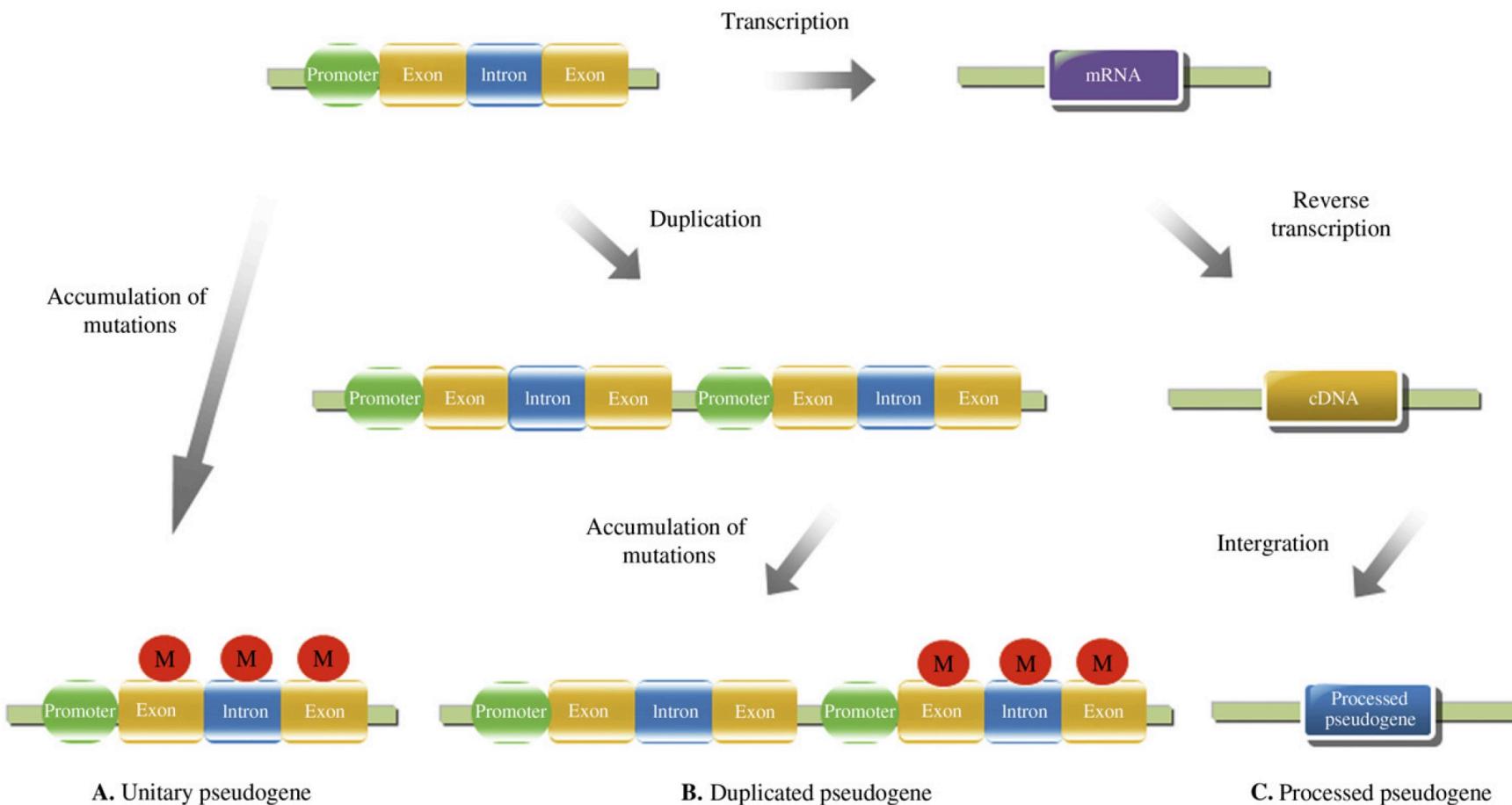


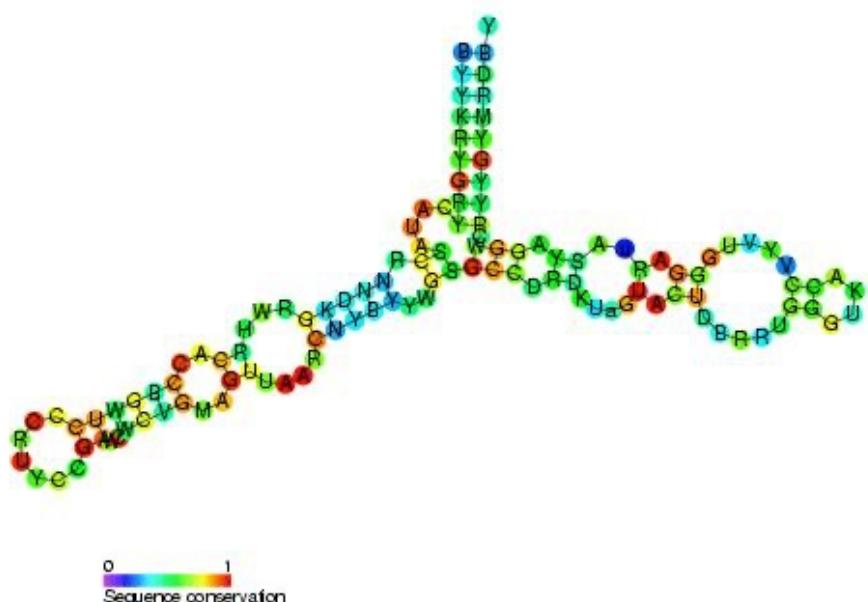
Figure 7-19 Essential Cell Biology 3/e (© Garland Science 2010)

Kinds of genes – pseudogenes

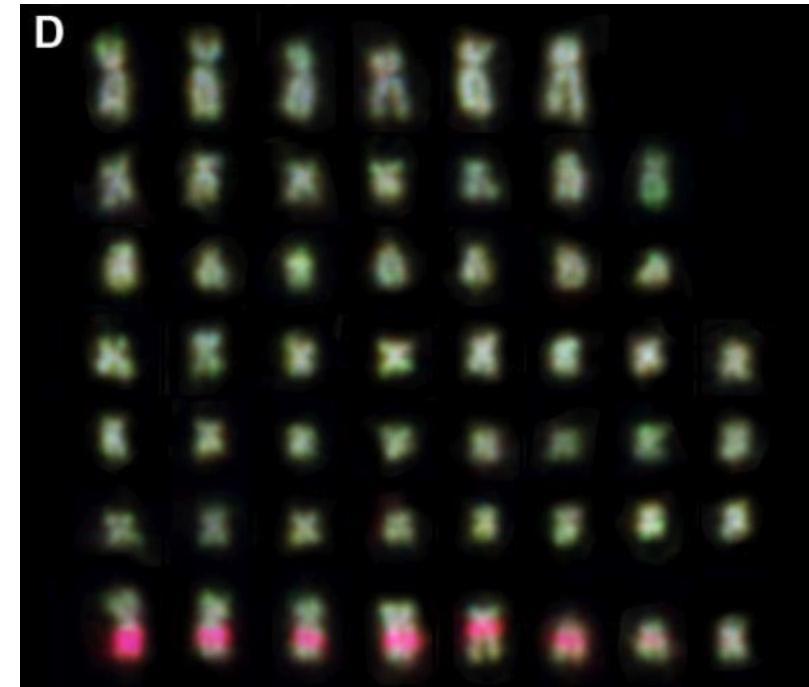


Kinds of genes – rRNAs

- rDNA can represent a substantial fraction of the genome. (many copies)

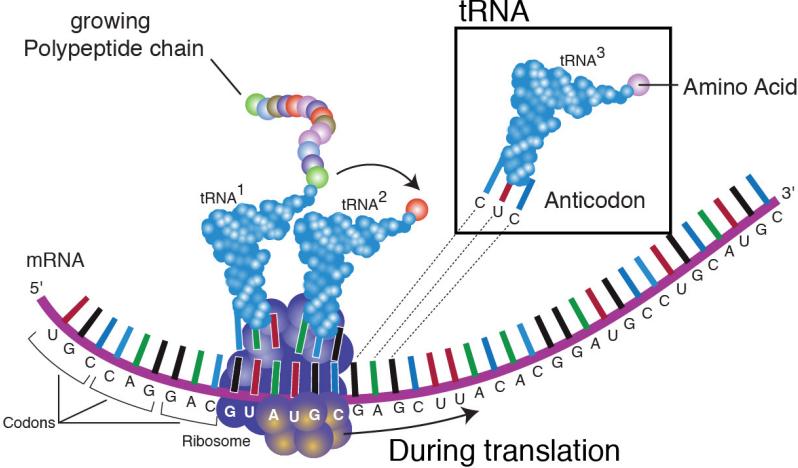


5S Ribosomal subunit
Wikimedia Commons

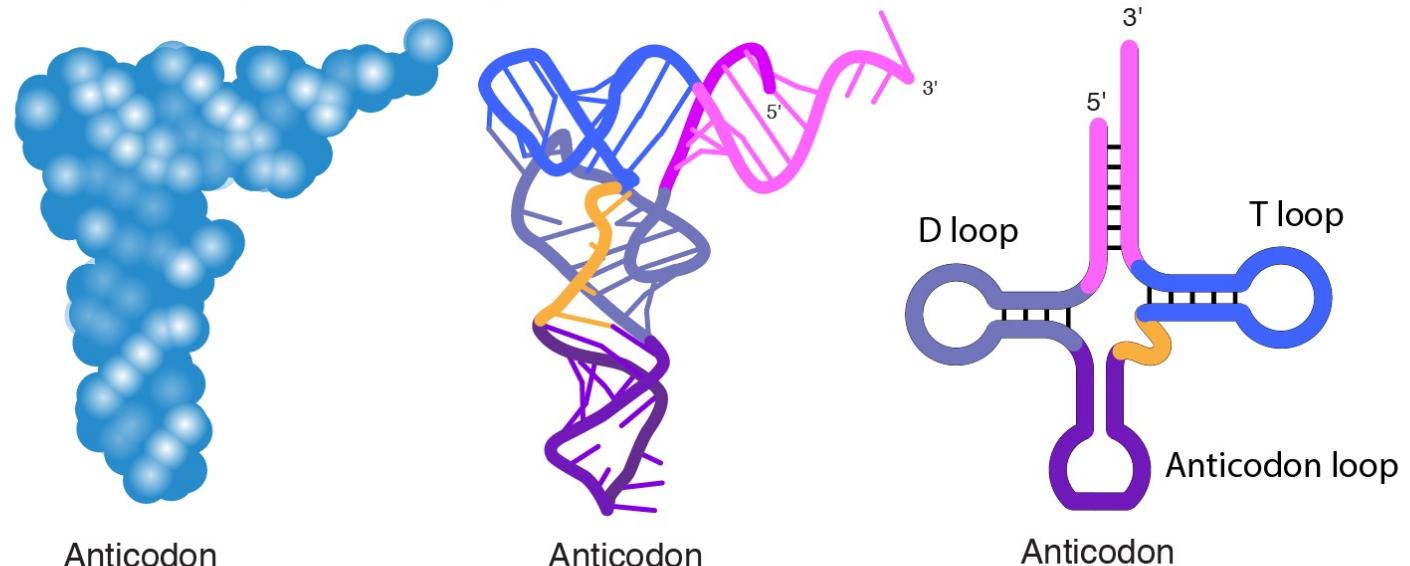


McElroy et al.,

Kinds of genes – tRNAs



Common ways of illustrating tRNA



Kinds of genes – other ncRNAs

Table 1: Classification of ncRNAs.

Type	Abbreviation	Full name	Size
Housekeeping ncRNAs	rRNA	ribosomal RNA	120–4,500 nt
	tRNA	transfer RNA	76–90 nt
	snRNA	small nuclear RNA	100–300 nt
	snoRNA	small nucleolar RNA	60–400
	TERC	telomerase RNA	/
	tRF	tRNA-Derived Fragments	16–28 nt
Regulatory ncRNAs	tiRNA	tRNA halves	29–50 nt
	miRNA	microRNA	21–23 nt
	siRNA	small interfering RNA	20–25 nt
	piRNA	piwi-interacting RNA	26–32 nt
	eRNA	enhancer RNA	50–2,000 nt
	lncRNA	long non-coding RNAs	>200 nt
	circRNA	circular RNA	100–10,000 nt
	Y RNA	Y RNA	/

Genome browsers & Data hubs

- [UCSC Genome Browser](#) – hub browser for lots of model genomes
- [YeastGenome](#) – *S. cerevisiae* genome browser
- [Flybase](#) – *Drosophila* genome browser
- [Wormbase](#) – *C. elegans* genome browser
- [Xenbase](#) – *Xenopus* genome browser
- [TAIR](#) – *Arabidopsis* genome browser
- [SolGenomics](#) – Tomato (and relatives) genome browser
- [Phytozome](#) – Non-model plant genome repository

General Feature Format (GFF) files

1. **seqid** - name of the chromosome or scaffold; chromosome names can be given with or without the 'chr' prefix. **Important note:** the seq ID must be one used within Ensembl, i.e. a standard chromosome name or an Ensembl identifier such as a scaffold ID, without any additional content such as species or assembly. See the example GFF output below.
2. **source** - name of the program that generated this feature, or the data source (database or project name)
3. **type** - type of feature. Must be a term or accession from the SOFA sequence ontology
4. **start** - Start position of the feature, with sequence numbering starting at 1.
5. **end** - End position of the feature, with sequence numbering starting at 1.
6. **score** - A floating point value.
7. **strand** - defined as + (forward) or - (reverse).
8. **phase** - One of '0', '1' or '2'. '0' indicates that the first base of the feature is the first base of a codon, '1' that the second base is the first base of a codon, and so on..
9. **attributes** - A semicolon-separated list of tag-value pairs, providing additional information about each feature. Some of these tags are predefined, e.g. ID, Name, Alias, Parent - see the [GFF documentation](#) for more details.

General Feature Format (GFF) files

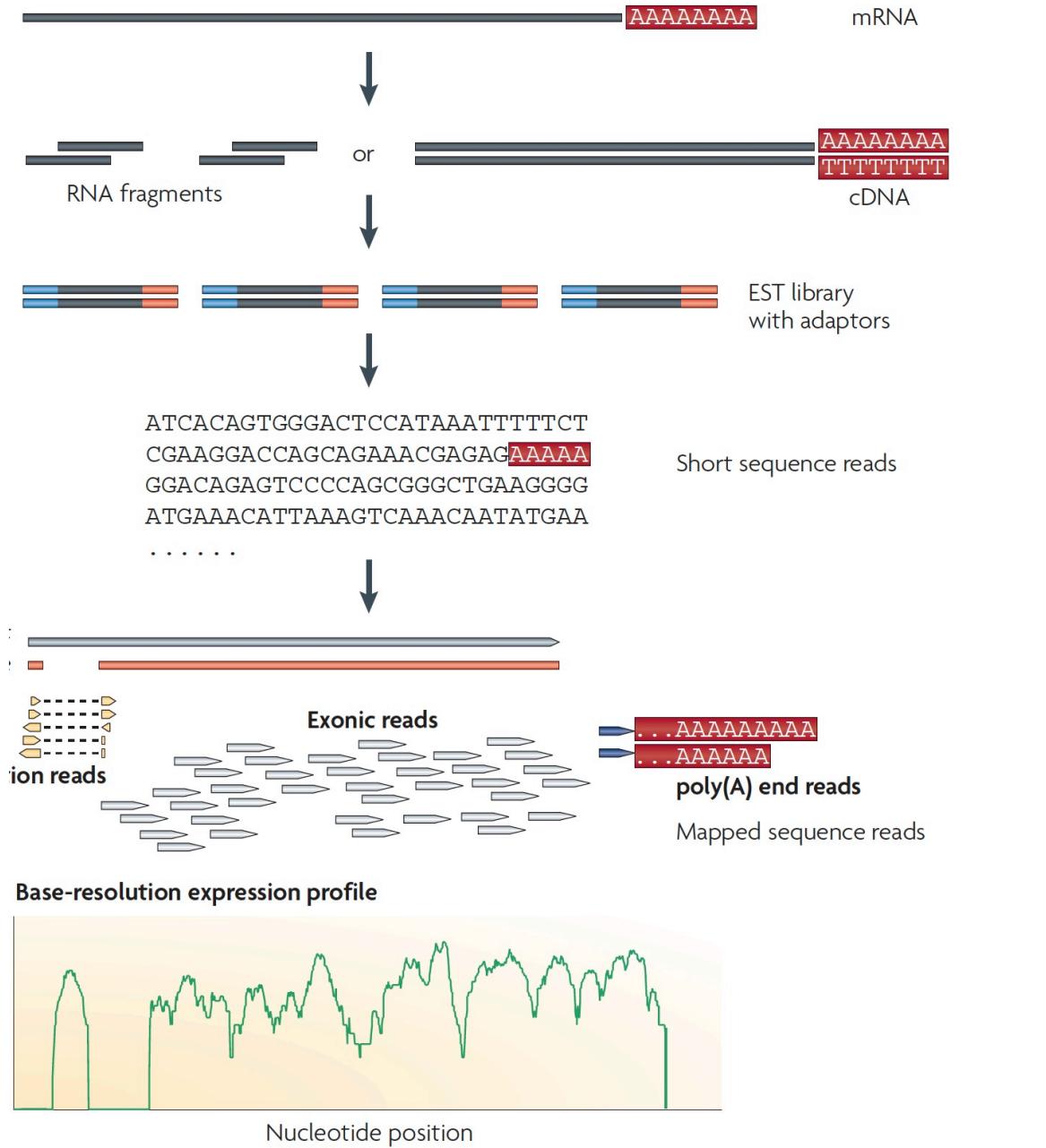
chromosome
gene
mRNA
five_prime_UTR
exon
CDS
three_prime_UTR
ncRNA_gene
lnc_RNA
miRNA
pre_miRNA
tRNA
ncRNA
snoRNA
snRNA
SRP_RNA
rRNA
RNase_MRP_RNA

Arabidopsis_thaliana.TAIR10.41.gff3
(might have to use gunzip to de-compress)

RNA sequencing is the best way to find genes

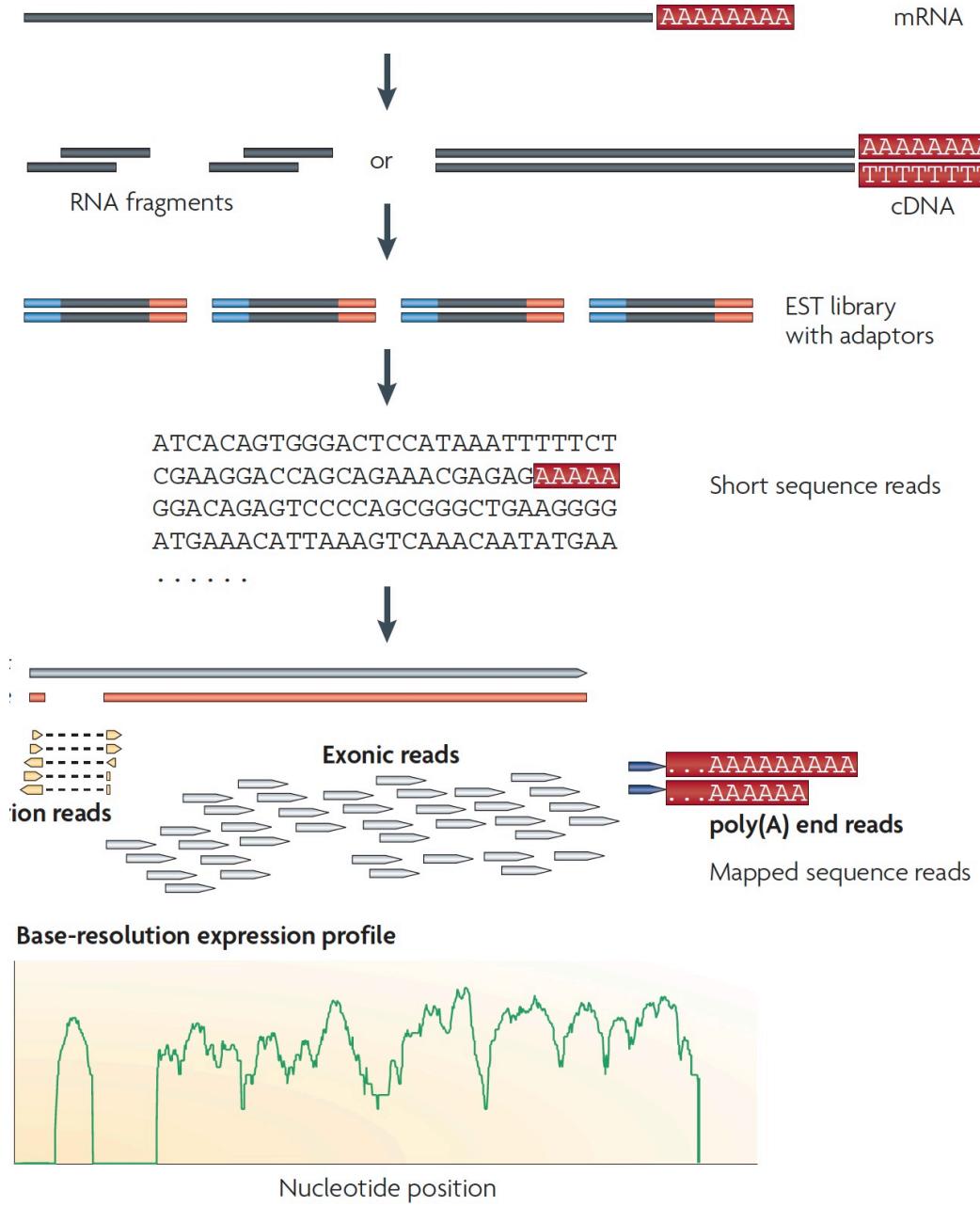
Complementary DNA Sequencing: Expressed Sequence Tags and Human Genome Project

MARK D. ADAMS, JENNY M. KELLEY, JEANNINE D. GOCAYNE, MARK DUBNICK,
MIHAEL H. POLYMEROPoulos, HONG XIAO, CARL R. MERRIL, ANDREW WU,
BJORN OLDE, RUBEN F. MORENO, ANTHONY R. KERLAVAGE,
W. RICHARD McCOMBIE, J. CRAIG VENTER*



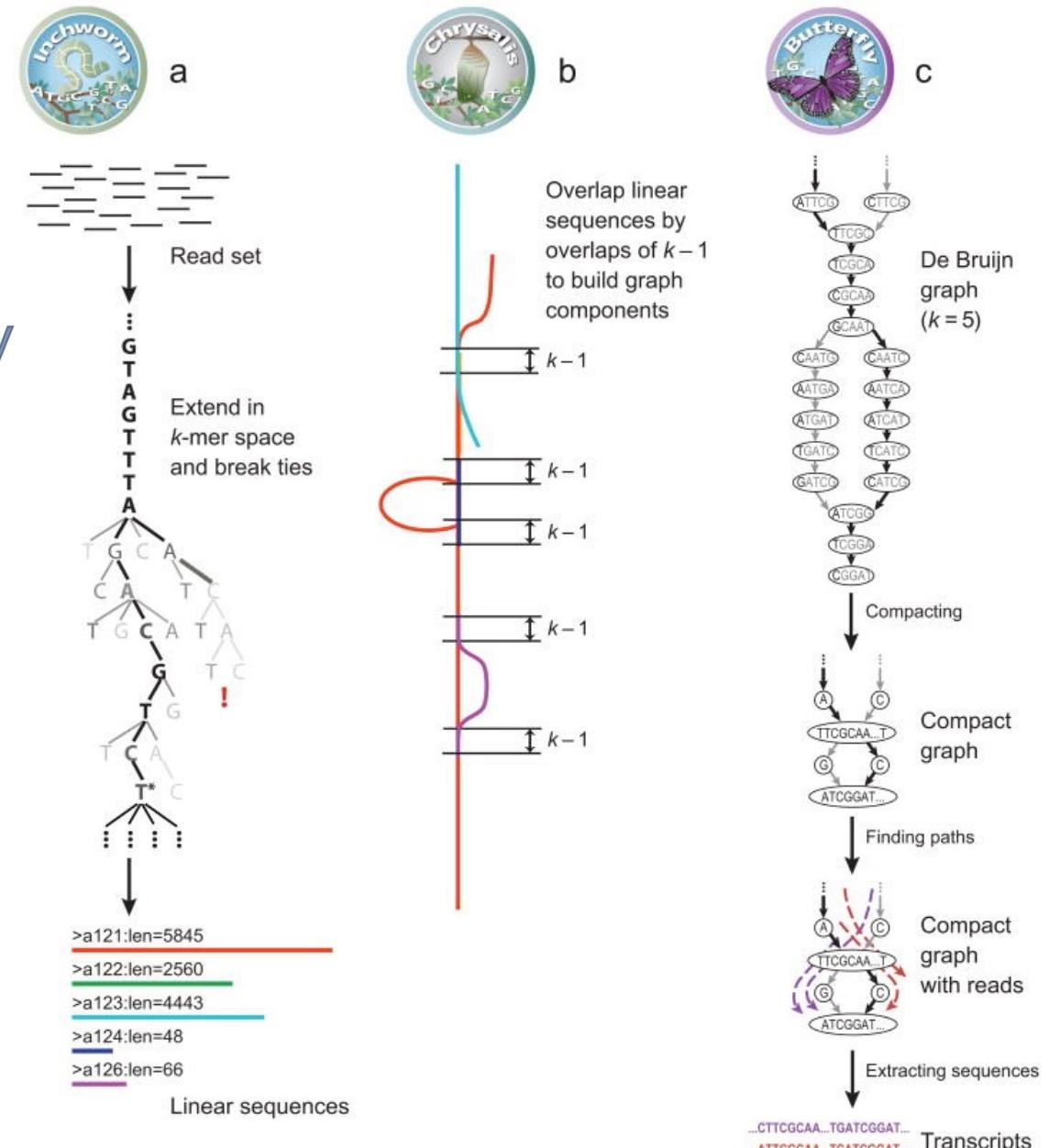
RNA sequencing is the best way to find genes

- Generate/sequence cDNA library



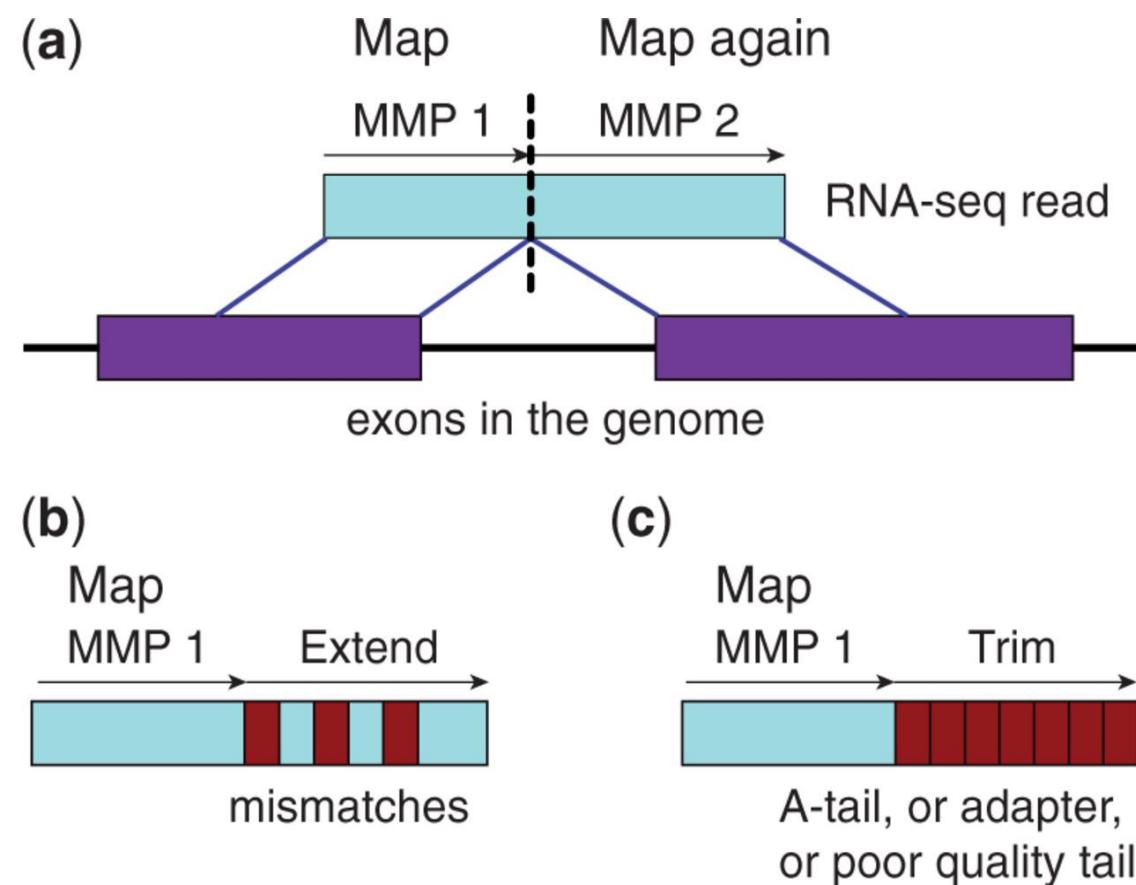
RNA sequencing is the best way to find genes

- Generate/sequence cDNA library
- Assemble RNAseq reads
(optional and genome-free)



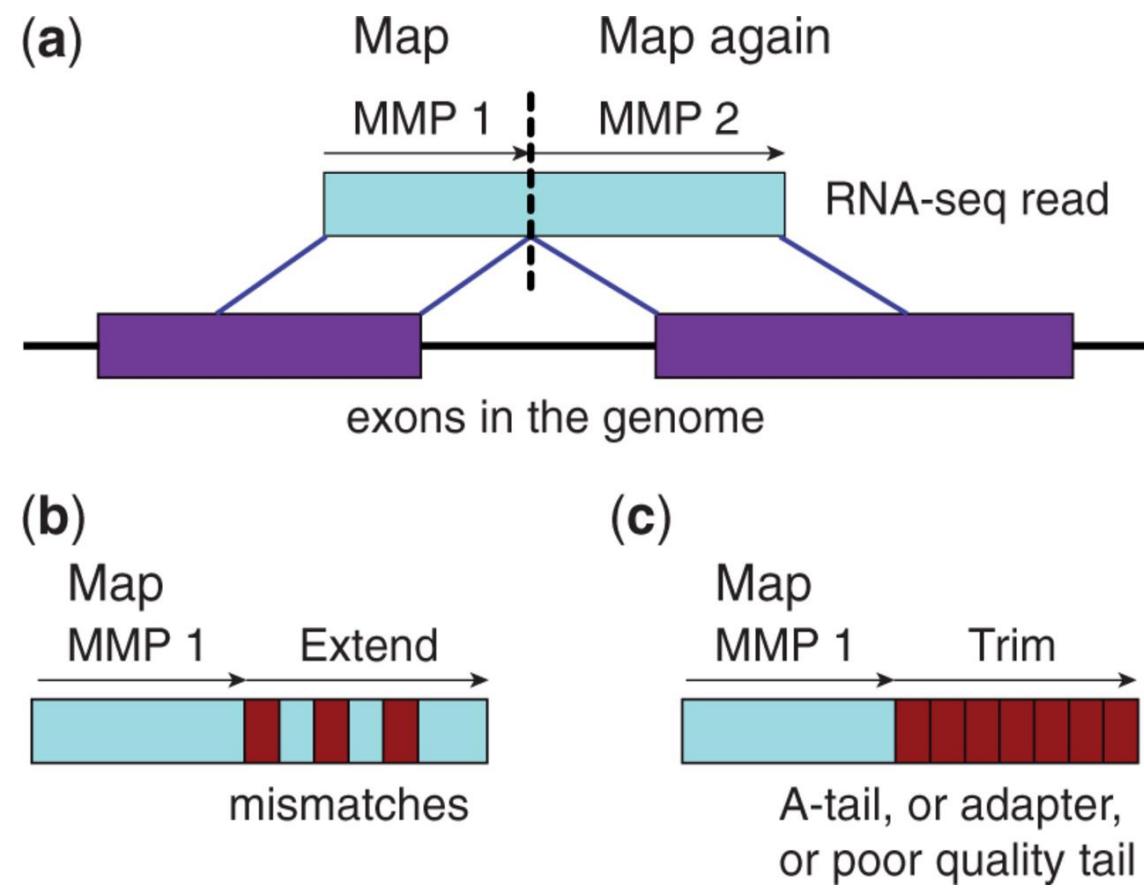
RNA sequencing is the best way to find genes

- Generate/sequence cDNA library
- Assemble RNAseq reads (optional and genome-free)
- Align RNAseq reads to genome sequence

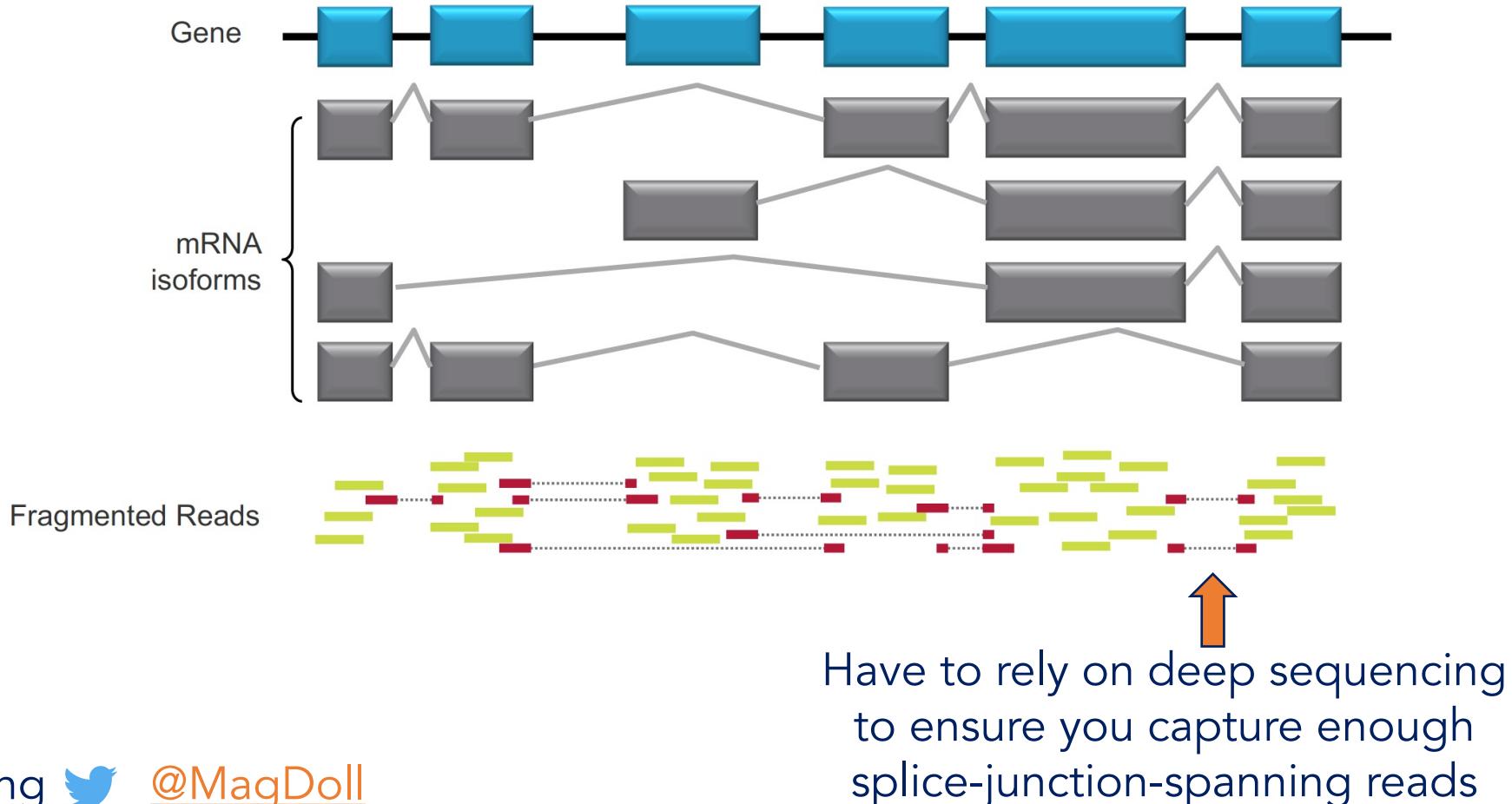


RNA sequencing is the best way to find genes

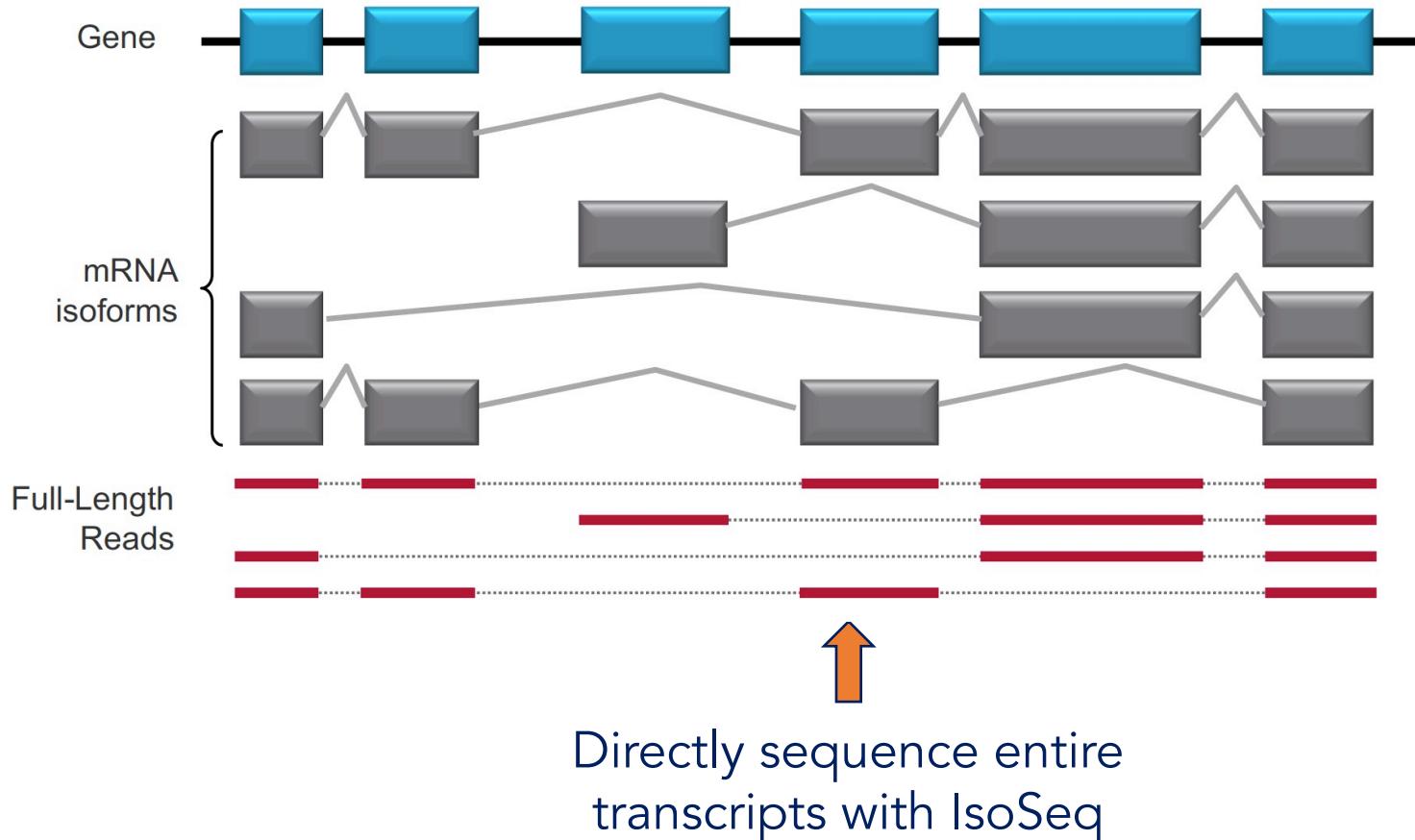
- Generate/sequence cDNA library
- Assemble RNAseq reads (optional and genome-free)
- Align RNAseq reads to genome sequence
- Use gaps to identify intron/exon boundaries, look for ORFs



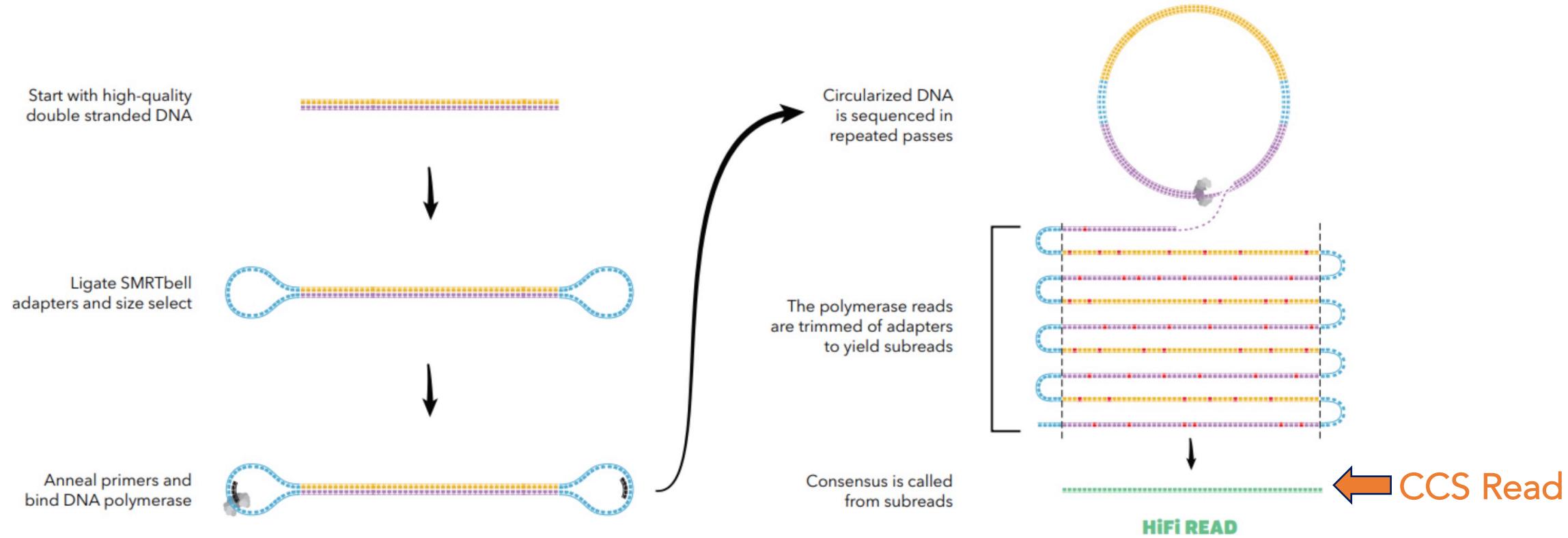
Developing high-quality gene models from short reads is challenging!



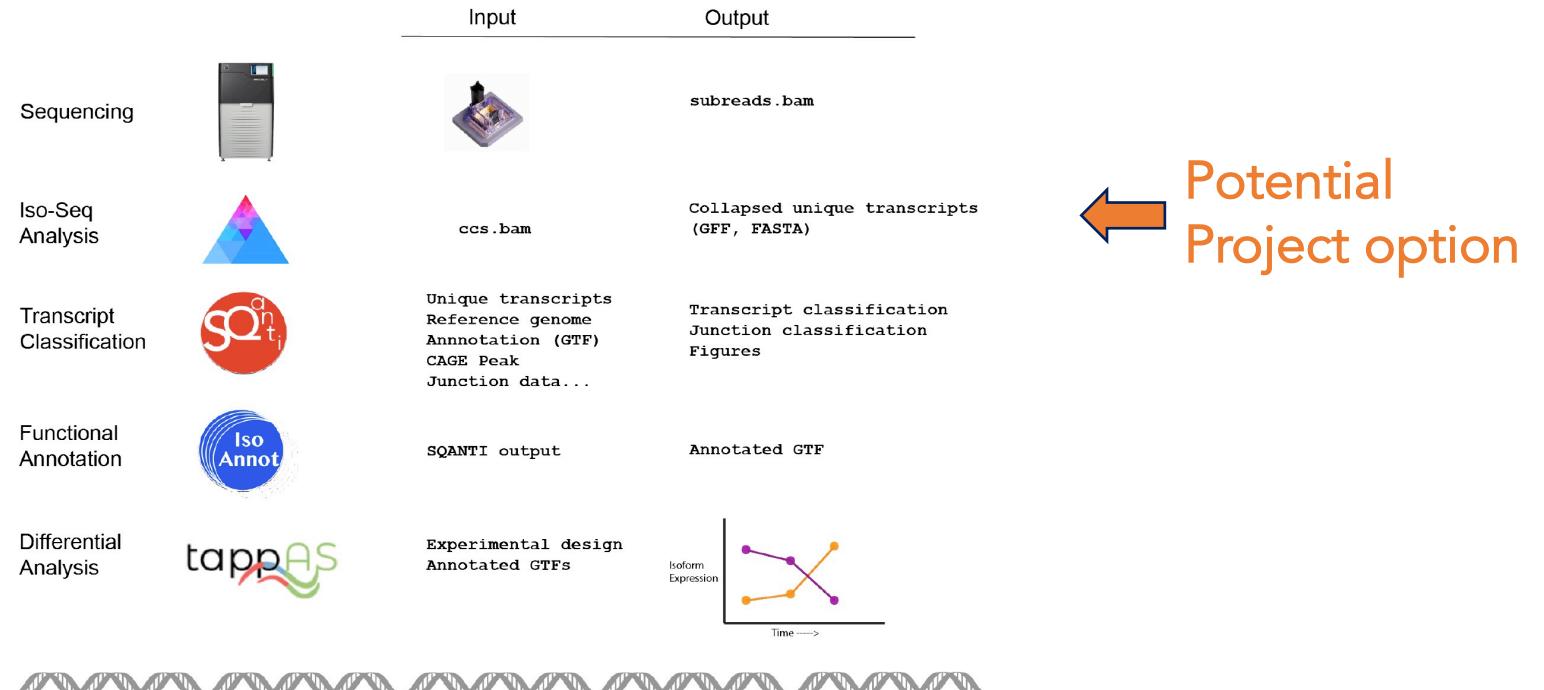
Sequencing cDNAs with long reads mostly alleviates that problem



IsoSeq: Uses single-molecule, real-time sequencing, plus Circular Consensus Sequencing to produce high quality, full-length transcript reads



Bioinformatic tools associated with IsoSeq now coming online



Supporting Tools



- collapse redundant transcripts
- merge multi-sample output
- saturation curve
- file format conversion
- single cell analysis



- gene family finding
- genome reconstruction
- evaluate assembly

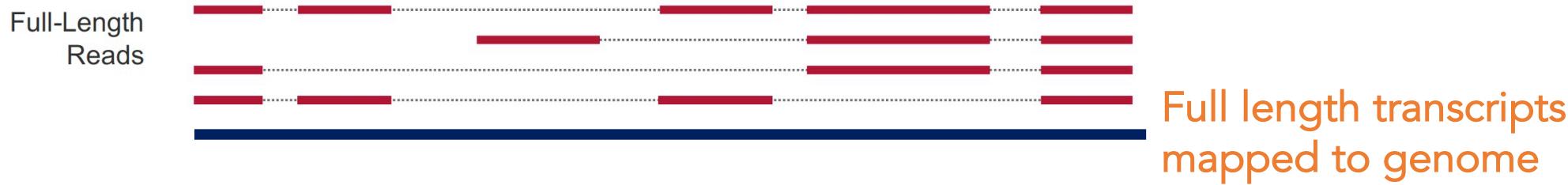


- collapse redundant transcripts
- merge multi-sample output
- NMD/ORF prediction
- transcript filtering

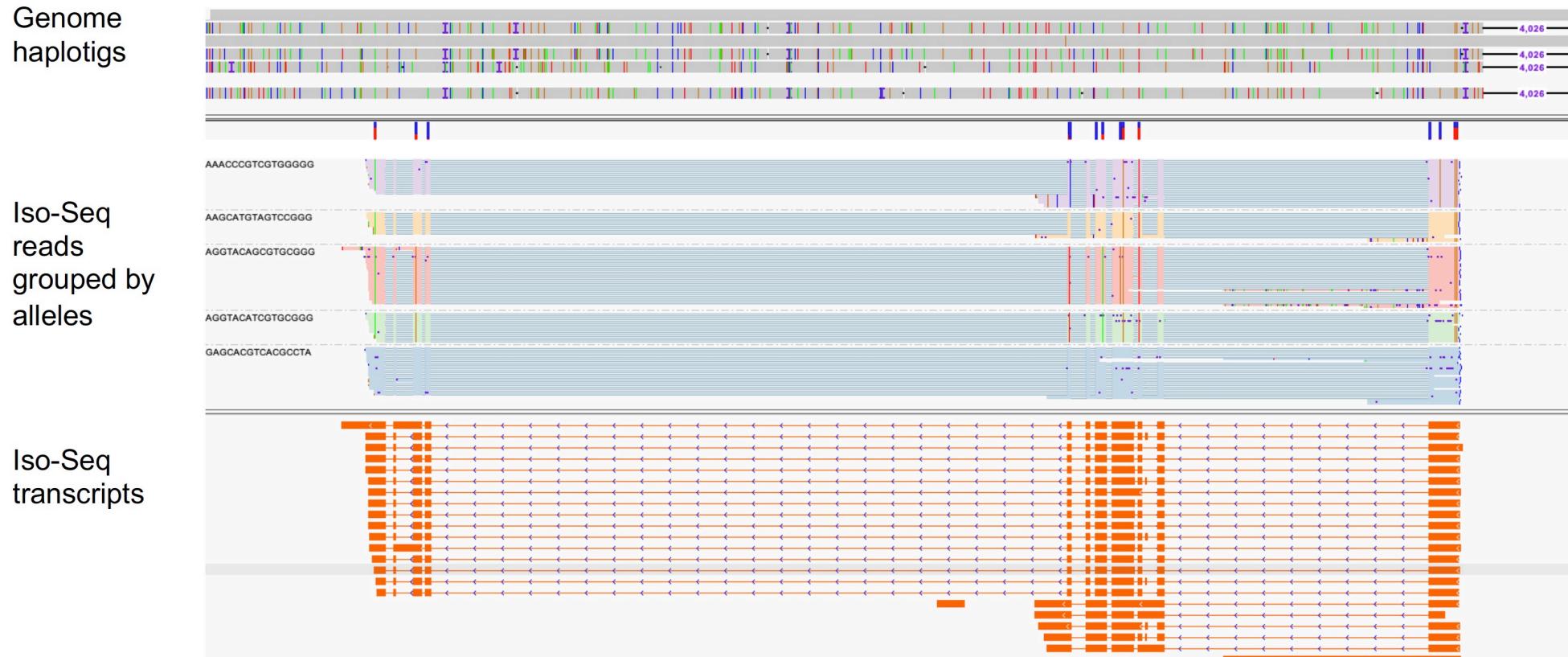


- long read processing & annotation pipeline developed independently by ENCODE4

Gene models are no longer computational inferences, but direct biological observations

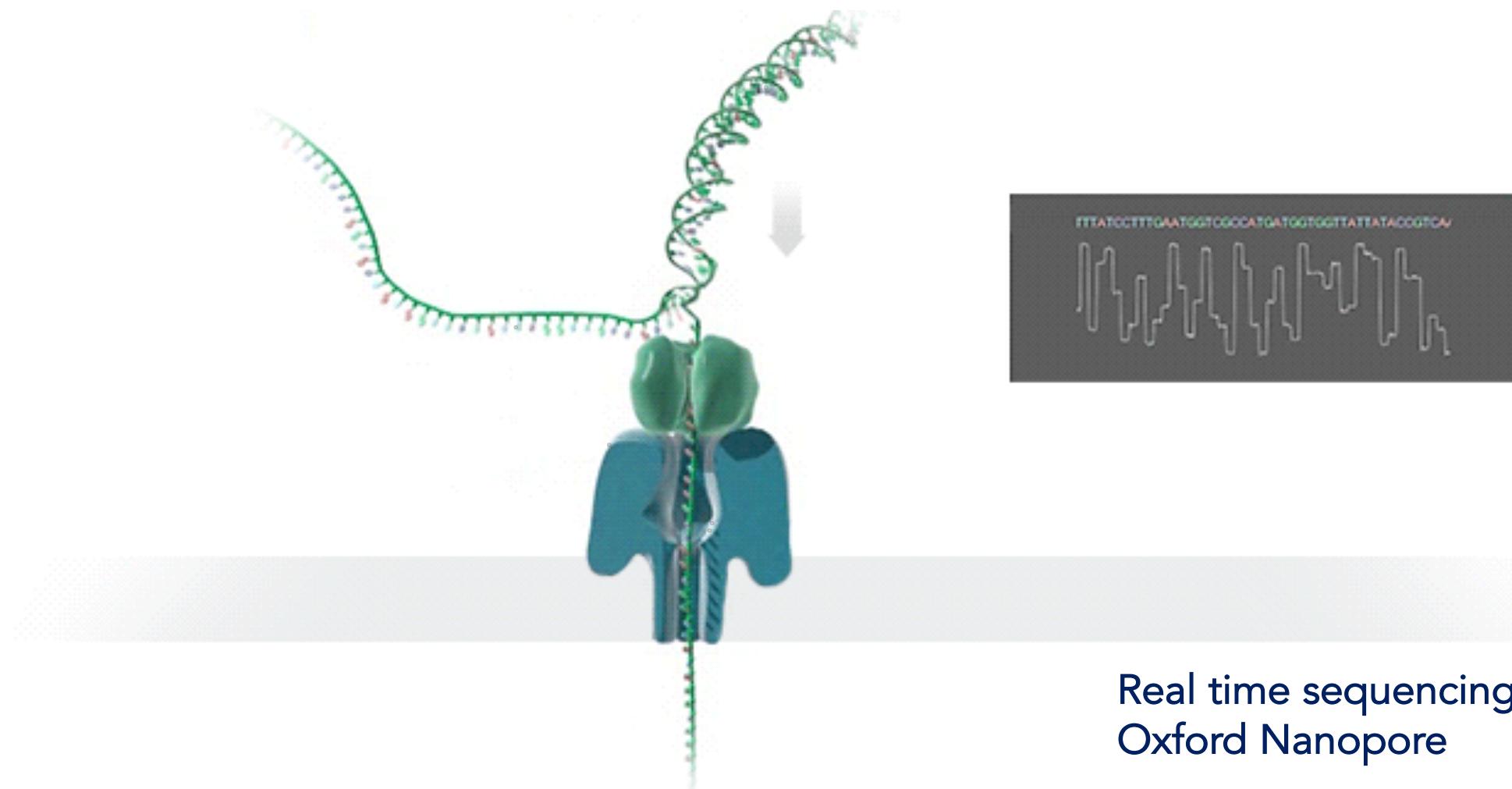


Clustering full-length transcripts with IsoPhase



IsoPhase claims to be able to distinguish IsoForms/Alleles/Paralogs

Other long-read technologies can also be used for RNAseq



IsoSeq is awesome, still doesn't solve everything

What IsoSeq can give you:

- 8 Million reads, 4 Million CCS reads, 500k Full length transcripts
- High quality, direct observations of tens of thousands of expressed genes
- Immediate, assembly-free genome annotation

What IsoSeq cannot give you:

- The full transcriptome
- Quantitative estimates of expression levels
- Analysis of reads produced on sequencing machines using older chemistry

IsoSeq is awesome, still doesn't solve everything

How to differentiate transcripts originating from similar but non-identical genes?

- Paralogs, isoforms, alleles (IsoPhase?)

PacBio-sourced bioinformatic tools are designed to highlight the PacBio technology, not intended for every use

- Need to be aware of the limitations of the technology
- Implement your own workflow, not someone else's!

De novo gene discovery – how to find genes *without* RNAseq

- Identifying genes in the absence of RNAseq data is difficult!
 - RNA can't be reliably sequenced from some organisms
- How might one discover genes without RNA?

De novo gene discovery – how to find genes *without* RNAseq

LilyPad.fasta

Open Reading Frames

GeneMarkS – coding potential

Glimmer

Next up: Primer Design & ENCODE discussion

Please Read: ENCODE 2011

Homework #3 Posted on Canvas/GitHub

