

# **RNA-Seq Data Analysis**

**Dr. Jung Soh  
657.000 SS 2017**

# What is RNA-seq?

- ▶ An experimental protocol that uses **next-generation sequencing** technologies to sequence the **RNA** molecules within a biological sample in an effort to determine the **primary sequence** and **relative abundance** of each RNA
  - Martin JA, Wang Z (2011) Next-generation transcriptome assembly. *Nat Rev Genet.* 12(10):671-682
- ▶ Also known as “Whole Transcriptome Shotgun Sequencing” (WTSS)

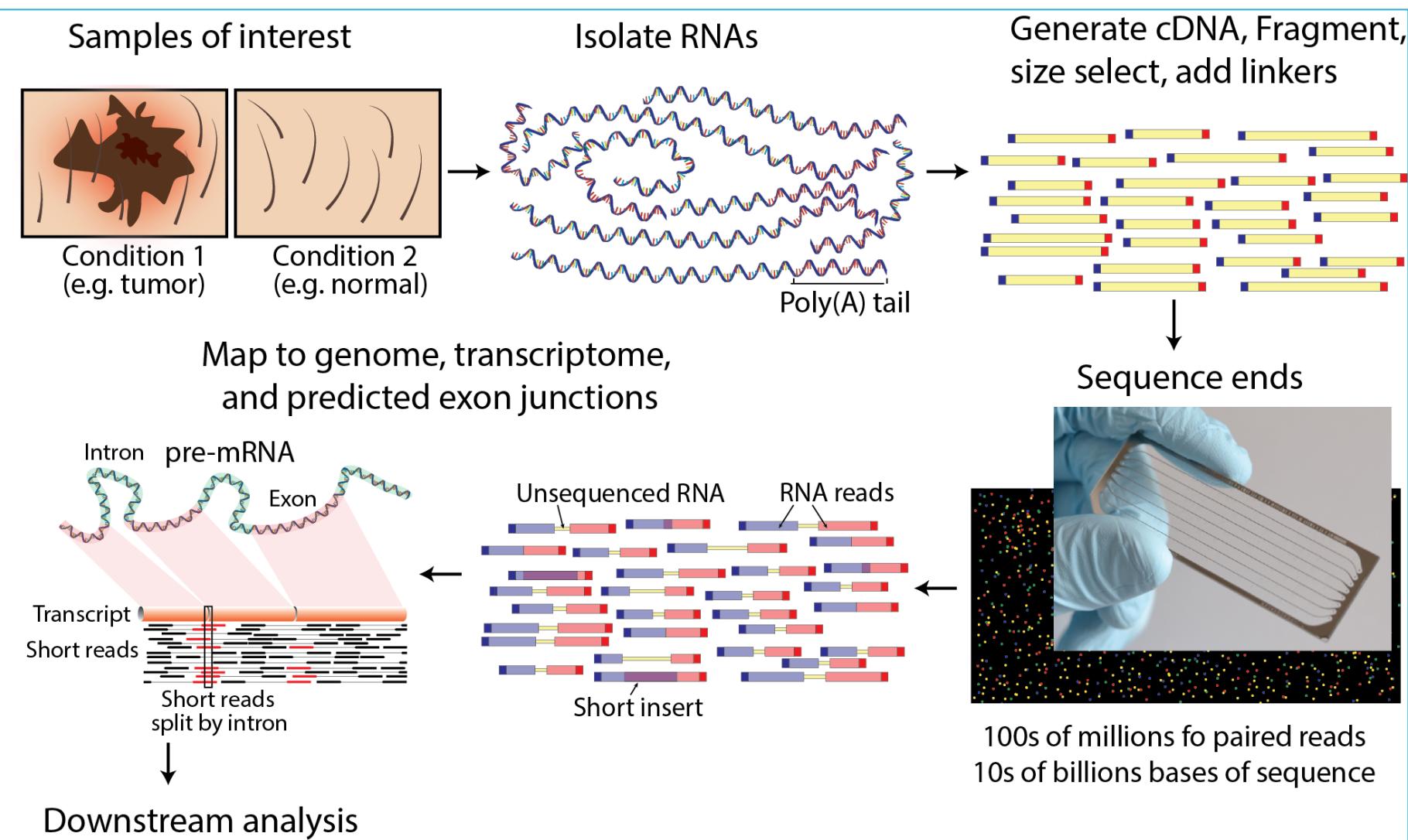
# Why RNA-seq?

- ▶ Some studies are not possible with DNA sequences
- ▶ To study functions based on gene expression changes
  - Drug treatment vs. no treatment
  - Patients vs. healthy people
  - Wild type vs. knock-out
- ▶ Some features available only at RNA level
  - Alternative isoforms, RNA editing, transcript fusion

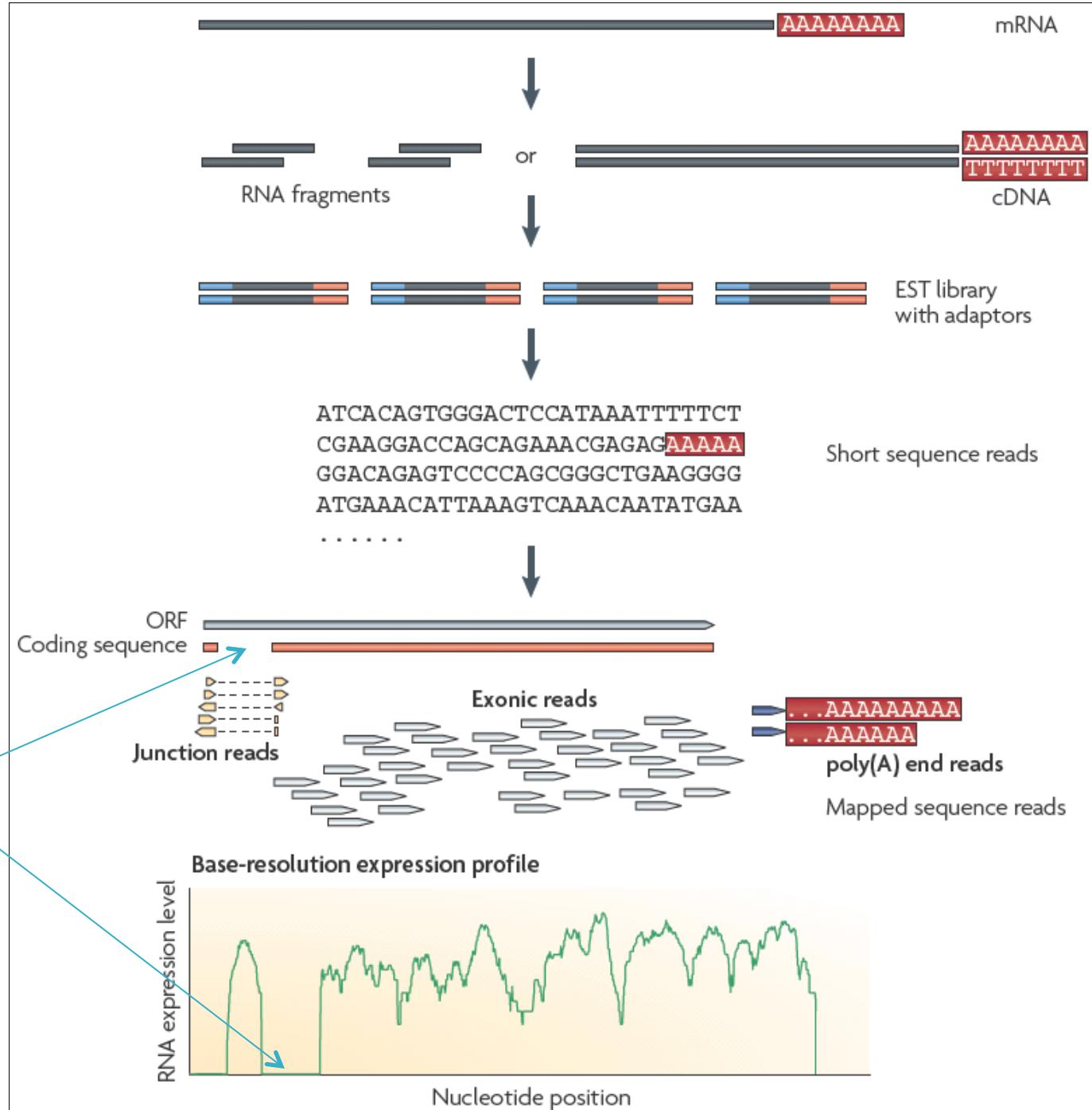
# Goals of RNA-seq experiments

- ▶ Differential expression (DE) analysis
- ▶ Alternative expression analysis
- ▶ Novel transcript discovery
- ▶ Allele-specific expression analysis
  - SNPs, mutations
- ▶ Gene fusion detection

# RNA sequencing



# RNA-Seq workflow



Wang Z, Gerstein M, Snyder M  
(2009) RNA-Seq: a revolutionary  
tool for transcriptomics. Nat Rev  
Genet. 10(1):57-63.

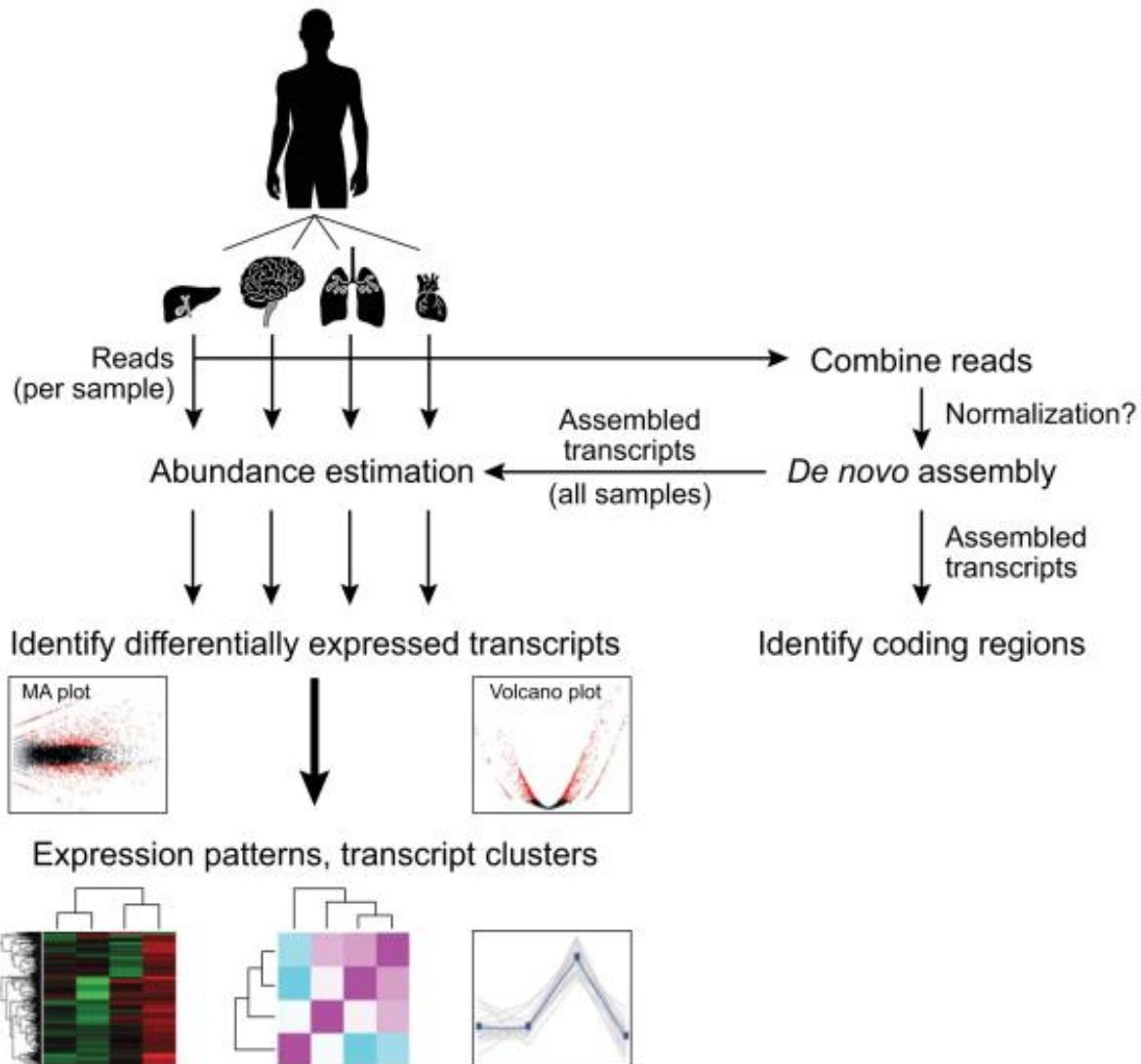
# RNA-seq poses more challenges

- ▶ Difficulty with sampling
  - More fragile than DNA
  - Different sizes of RNA
- ▶ Relative abundance of RNAs hard to control
  - Can vary by orders of magnitude
  - Uneven coverage
    - Many reads from a small number of highly expressed genes
- ▶ Computational analysis challenges
  - Assembly, (spliced) alignment, quantification, normalization, visualization, ...

# RNA-Seq vs microarray

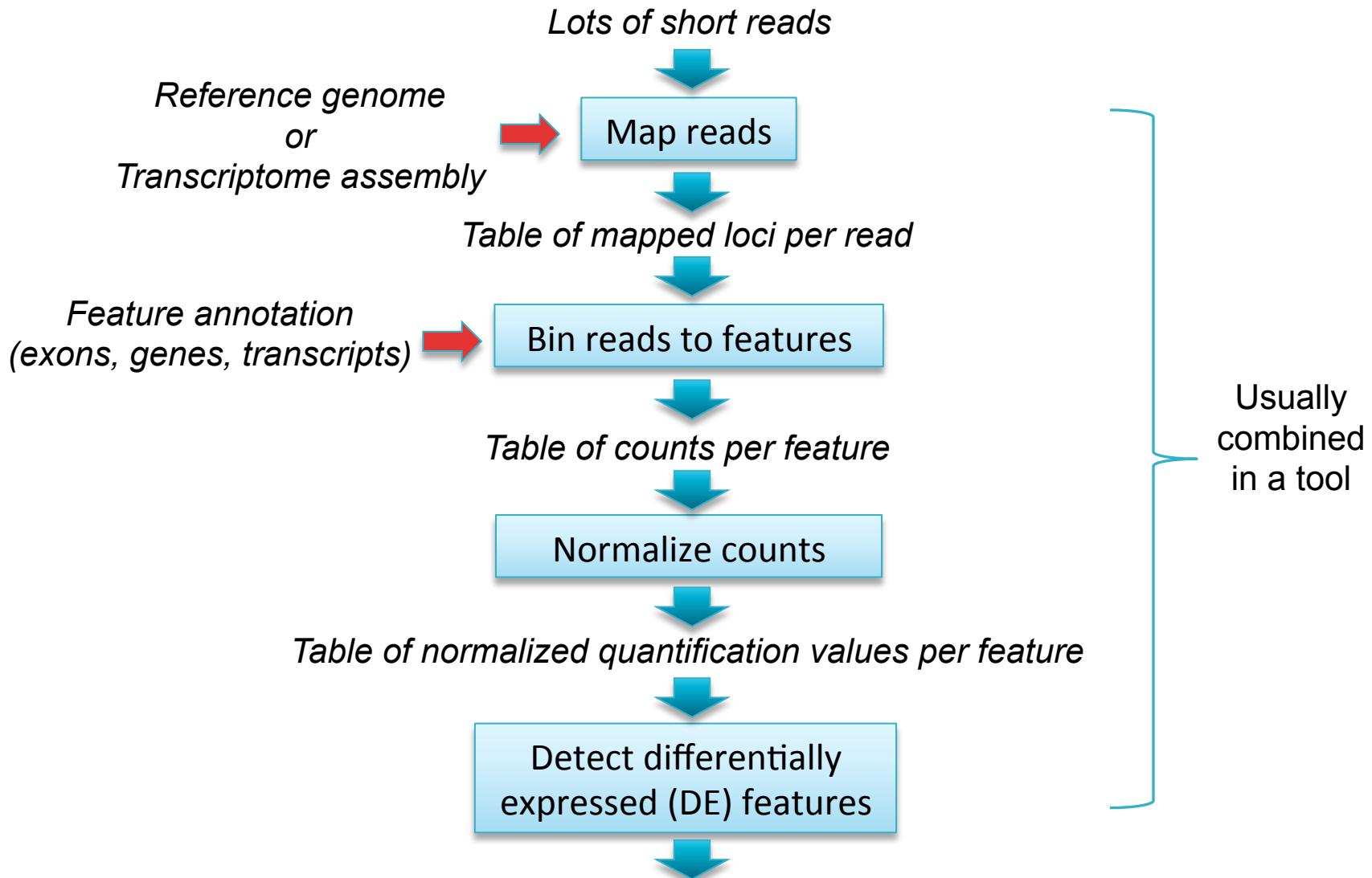
Characteristics	RNA-Seq	Microarray
Which transcripts?	All in a sample	Only those for which probes are designed
Transcript sequence generation	Yes	No
Low-abundance transcript detection	Yes	Limited
Abundance info source	Count (of the reads aligned to transcript)	Fluorescence level (of the probe spot for gene)
Resolution	Base	Probe sequence
Background noise	Low	High
Additional info	Alternative splicing, transcriptome-level variation	Limited

# *De novo* transcriptome assembly and analysis



Haas, B. J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P. D., Bowden, J., ... Regev, A. (2013). *De novo* transcript sequence reconstruction from RNA-Seq: reference generation and analysis with Trinity. *Nature Protocols*, 8(8)

# RNA-Seq DE analysis data flow



# Mapping reads

- ▶ Need a reference genome
  - Or *de novo* assembled transcriptome
- ▶ Issues
  - Reads spanning across exon junction
  - Alternative splicing
  - Reads mapping to multiple locations in the genome
  - Huge amounts of data
- ▶ Most common mapping results format
  - SAM: sequence alignment/map
  - BAM: binary format of SAM
- ▶ Many tools
  - Bowtie, SOAP, BWA, SHRiMP, mrFAST, mrsFAST, ZOOM, SSAHA2, Mosaik

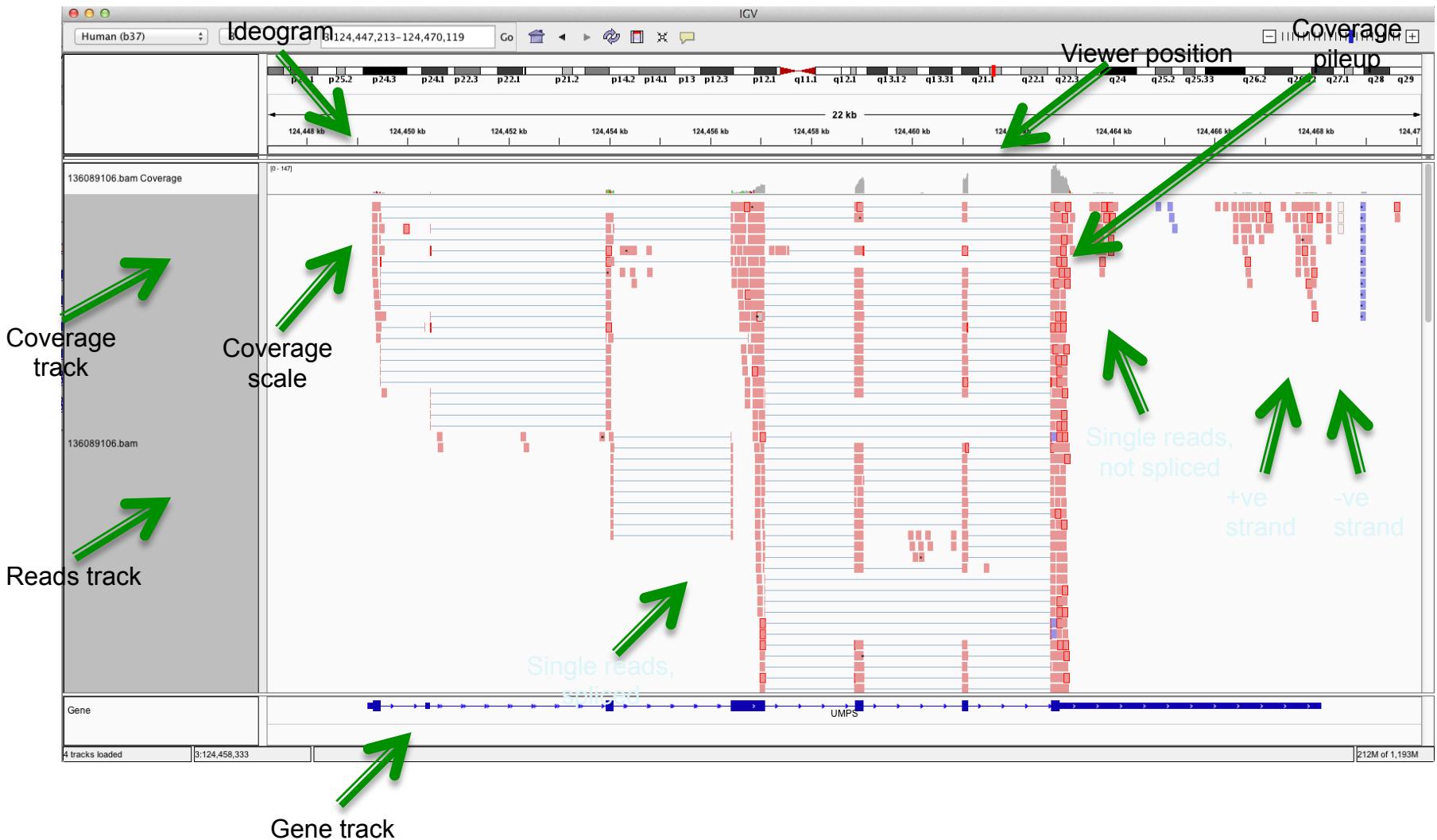
# SAM/BAM format

## Example SAM/BAM header section (abbreviated)

```
mgriffit@linus270 ~$ samtools view -H /gscmnt/gc13001/info/model_data/2891632684/build136494552/alignments/136080019.bam | grep -P "SN:22|HD|RG|PG"
@HD VN:1.4 SO:coordinate
@SQ SN:22 LN:51304566 UR:ftp://ftp.ncbi.nih.gov/genbank/genomes/Eukaryotes/vertebrates_mammals/Homo_sapiens/GRCh37/special_requests/GRCh37-lite.fa.gz AS:GRCh37-lite M5:a718acaa6135fdca8357d5bfe9
4211dd SP:Homo sapiens
@RG ID:2888721359 PL:illumina PU:D1BA4ACXX.3 LB:H_KA-452198-0817007-cDNA-3-lb1 PI:365 DS:paired end DT:2012-10-03T19:00:00-0500 SM:H_KA-452198-0817007 CN:WUGSC
@PG ID:2888721359 VN:2.0.8 CL:tophat --library-type fr-secondstrand --bowtie-version=2.1.0
@PG ID:MarkDuplicates PN:MarkDuplicates PP:2888721359 VN:1.85(exported) CL:net.sf.picard.sam.MarkDuplicates INPUT=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/scratch-ILg6Y/H_KA-452198-0817007-cDNA-3-lb1-2888360300.bam OUTPUT=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/staging-liuJS/H_KA-452198-0817007-cDNA-3-lb1-2888360300.metrics REMOVE_DUPLICATES=false ASSUME_SORTED=true MAX_FILE_HANDLES_FOR_READ_ENDS_MAP=9500 TMP_DIR=/gscmnt/gc13001/info/build_merged_alignments/merged-alignment-blade10-2-5.gsc.wustl.edu-jwalker-15434-136080019/scratch-ILg6Y VALIDATION_STRINGENCY=SILENT MAX_RECORDS_IN_RAM=500000 PROGRAM_RECORD_ID=MarkDuplicates PROGRAM_GROUP_NAME=MarkDuplicates MAX_SEQUENCE_LENGTHS_FOR_DISK_READ_ENDS_MAP=50000 SORTING_COLLECTION_SIZE_RATIO=0.25 READ_NAME_REGEX=[a-zA-Z0-9]+:[0-9]+:[0-9]+:[0-9]+.* OPTICAL_DUPLICATE_PIXEL_DISTANCE=100 VERBOSITY=INFO QUIET=false COMPRESSION_LEVEL=5 CREATE_INDEX=false CREATE_MD5_FILE=false
mgriffit@linus270 ~$
```

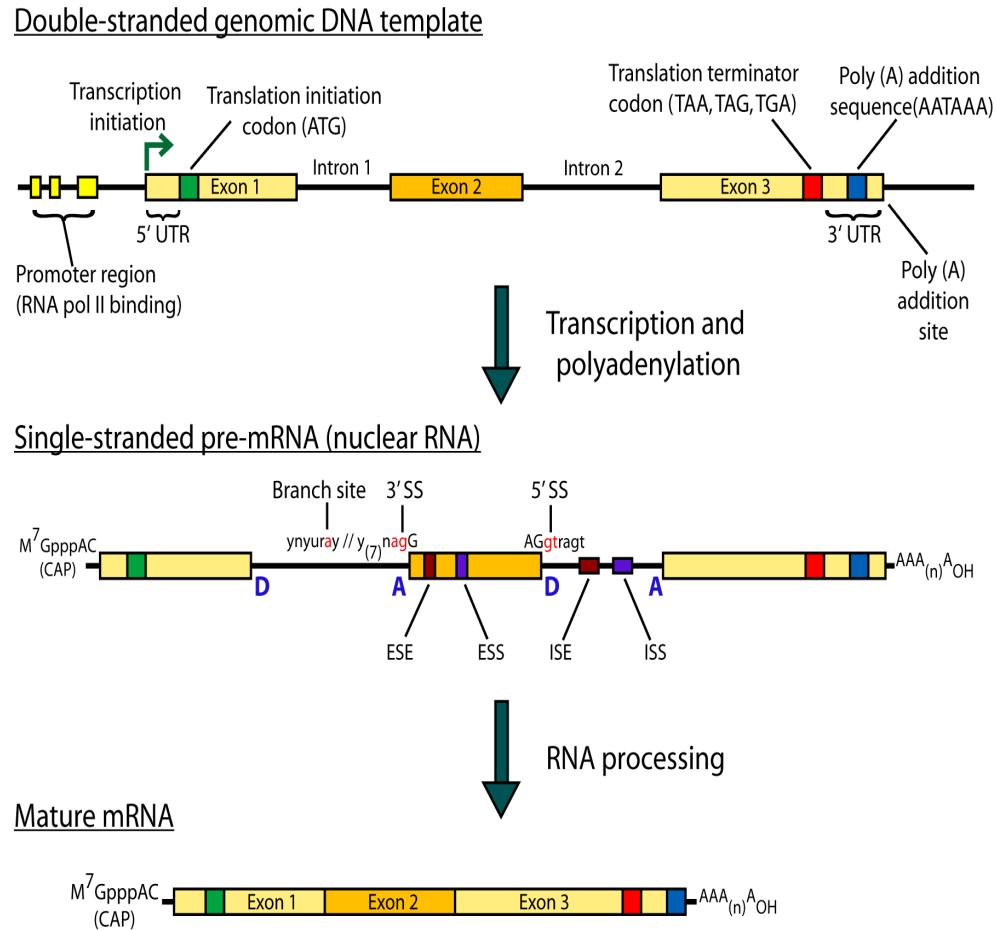
## Example SAM/BAM alignment section (only 10 alignments shown)

# Alignment visualization with IGV



# Splice-aware alignment

- ▶ RNA-seq reads represent mRNA with introns removed
  - RNA-seq reads may span large introns
- ▶ But we usually align these reads back to reference genome
- ▶ Unless your reads are short (<50bp) you should use a splice-aware aligner
  - TopHat, STAR, MapSplice, etc.



# Downstream analysis

- ▶ Transcript abundance estimation
- ▶ Sample/replicate quality check
- ▶ Transcript functional annotation (*de novo* assembly)
- ▶ Differential expression analysis
- ▶ Coding region identification
- ▶ Alternative expression analysis

# Transcript abundance estimation

- ▶ For each transcript, count total number of reads mapped
  - Also called “binning” the reads
  - Counts are not directly comparable across features or samples (yet)
  - Followed by normalization into expression values

# Normalizing counts

- ▶ Why normalize?
  - Longer features (naturally) can have more reads mapped
  - Deeper sequencing produces more reads
- ▶ RPKM (or FPKM) most commonly used
  - Reads (Fragments) per Kilobase per Million reads
  - Defined as  $C/(LN)$ 
    - $C$  = number of reads mapped to a feature
    - $L$  = length of the feature (in kilobases)
    - $N$  = total number of reads from the sample (in millions)

# *RPKM* examples

Gene A 600 bases

$$\text{RPKM} = 12/(0.6*6) = 3.33$$

C=12

Gene B 1100 bases

$$\text{RPKM} = 24/(1.1*6) = 3.64$$

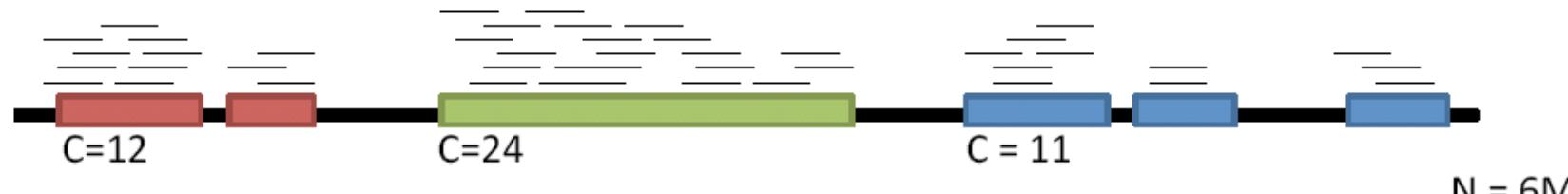
C=24

Gene C 1400 bases

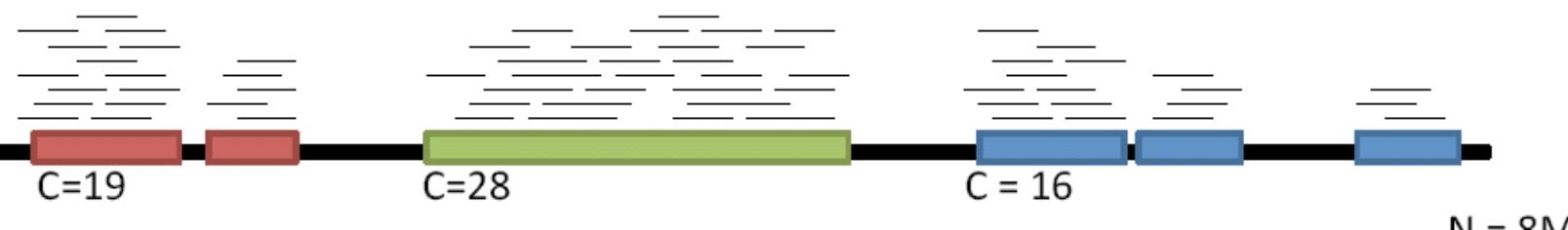
$$\text{RPKM} = 11/(1.4*6) = 1.31$$

C = 11

Sample 1



Sample 2



$$\text{RPKM} = 19/(0.6*8) = 3.96$$

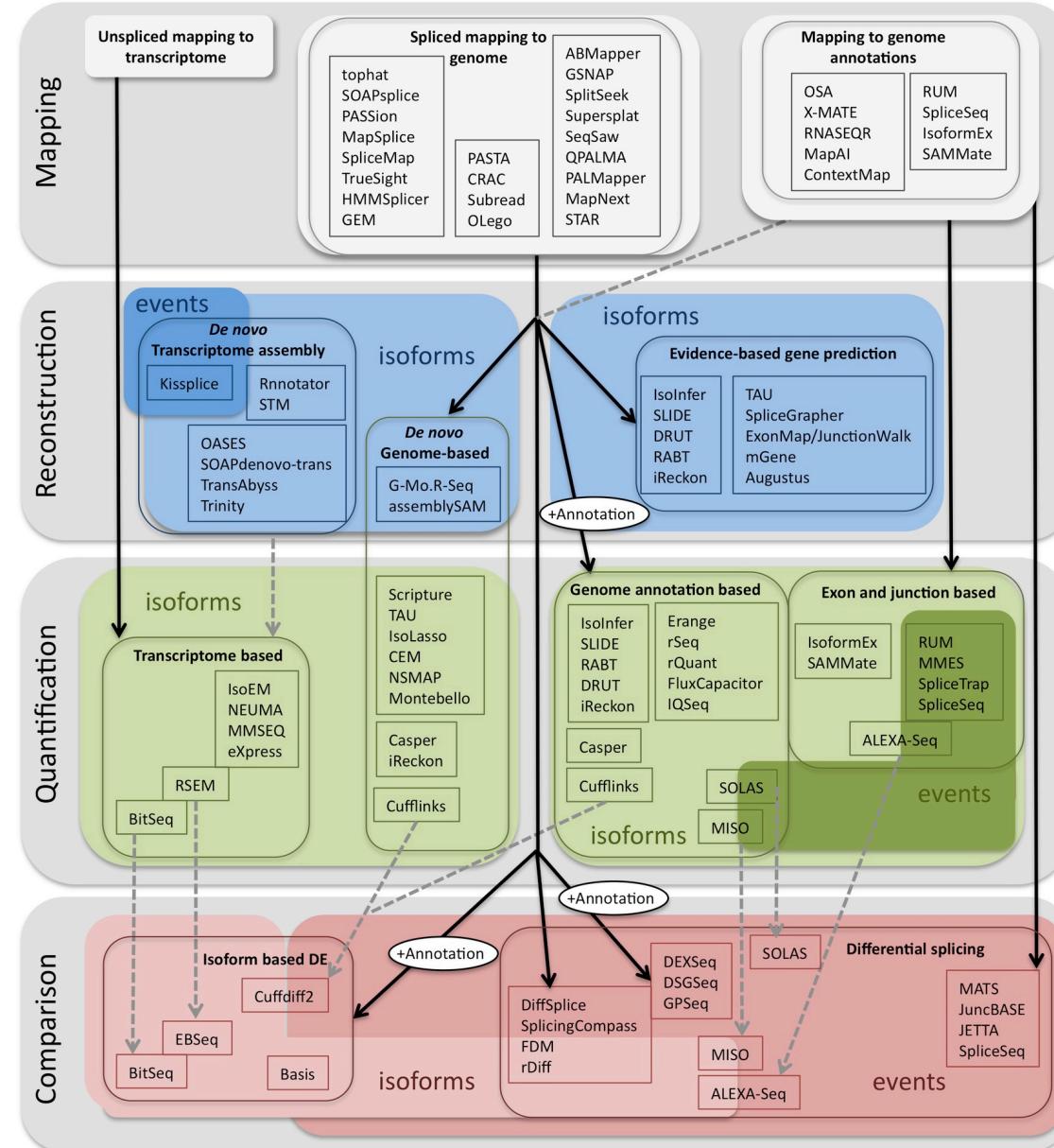
$$\text{RPKM} = 28/(1.1*8) = 1.94$$

$$\text{RPKM} = 16/(1.4*8) = 1.43$$

# Differential expression analysis

- ▶ Compare quantification values across samples or across features
  - The goal is to find differentially expressed (DE) features
- ▶ Many tools summarize/normalize counts and suggest DE features
  - Cufflinks/Cuffdiff, R packages (DESeq, edgeR, baySeq, TSPM), Samtools
- ▶ Determination of DE features depends on
  - Fold changes (FC)
  - Statistical significance of FC (FDR, p-value)

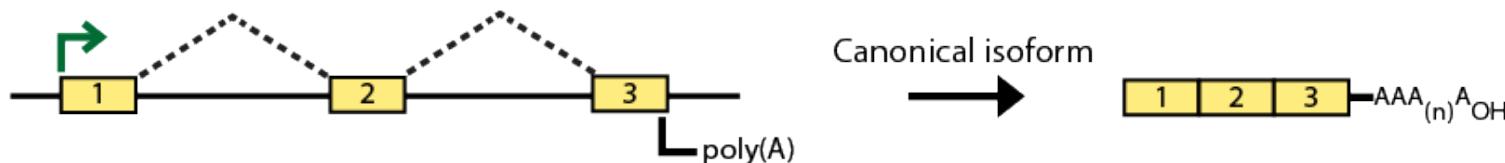
# Studying splicing from RNA-seq



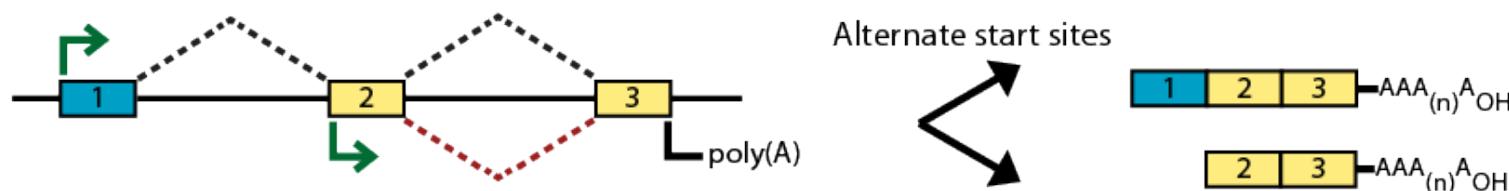
<http://www.rna-seqblog.com/data-analysis/splicing-junction/methods-to-study-splicing-from-rna-seq/>

# Type of alternative expression

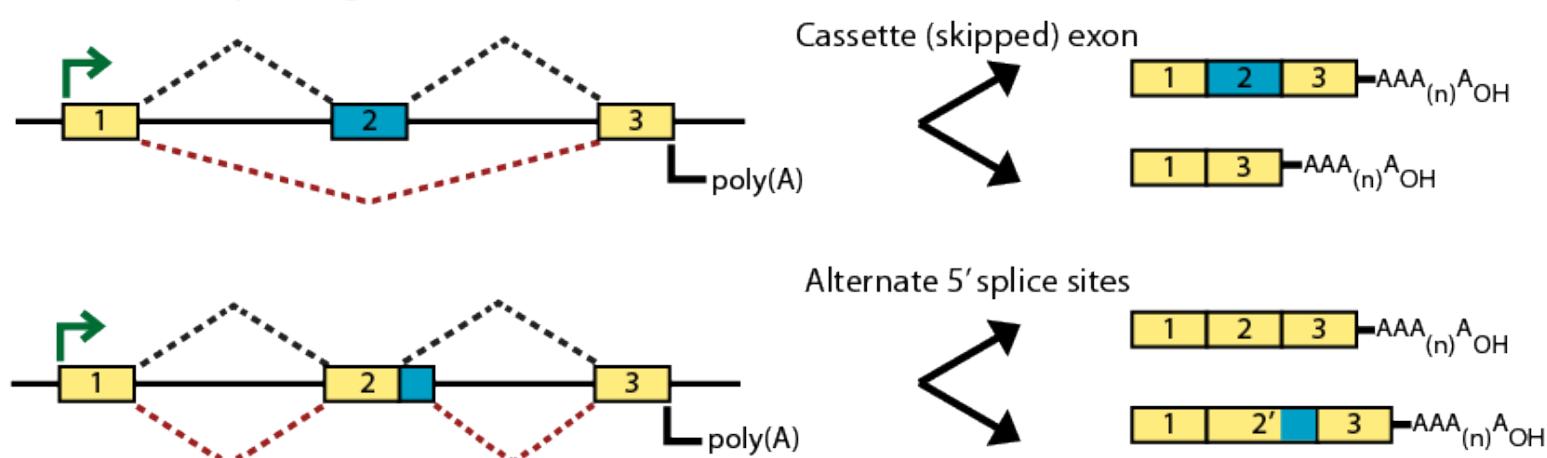
## Simple transcription



## Alternative transcript initiation



## Alternative splicing



# Types of alternative expression

