

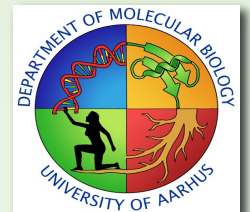
Next generation sequencing

RNA-seq

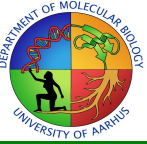
Stig Uggerhøj Andersen, PhD

Department of Molecular Biology

University of Aarhus



Lecture overview



- Applications and data scale
- Sample preparation
- Strand specificity
- Single-cell sequencing
- Spliced alignments
- Quantification and normalization
- Long read RNA-seq

Counting or Profiling

- 10 million total reads of 50 bp length from poly-A selected RNA will give performance better than any microarray

Studying Alternative Splicing or quantifying cSNPs for most transcripts

- Deeper profiling of 50 to 100 million reads, with read lengths of 50 to 100 bps, from poly-A selected RNA using mRNA-Seq assay

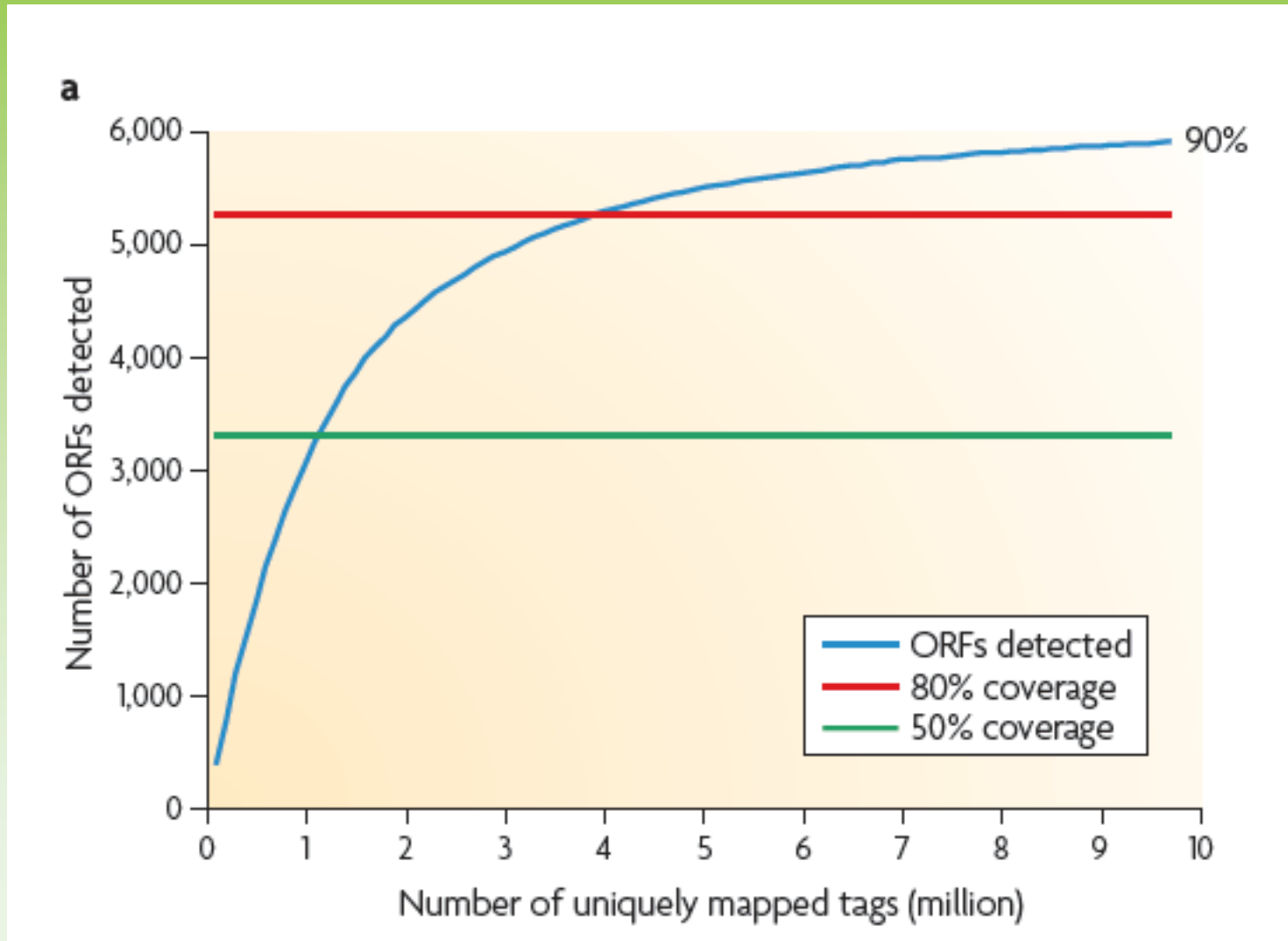
Complete Annotation of an entirely New Transcriptome

- ~500 Million reads of 100 bp read length from multiple tissues
- Normalized stranded mRNA-Seq
- Normalized stranded Total RNA-Seq for looking at ncRNAs
- Small RNA-Seq for microRNAs

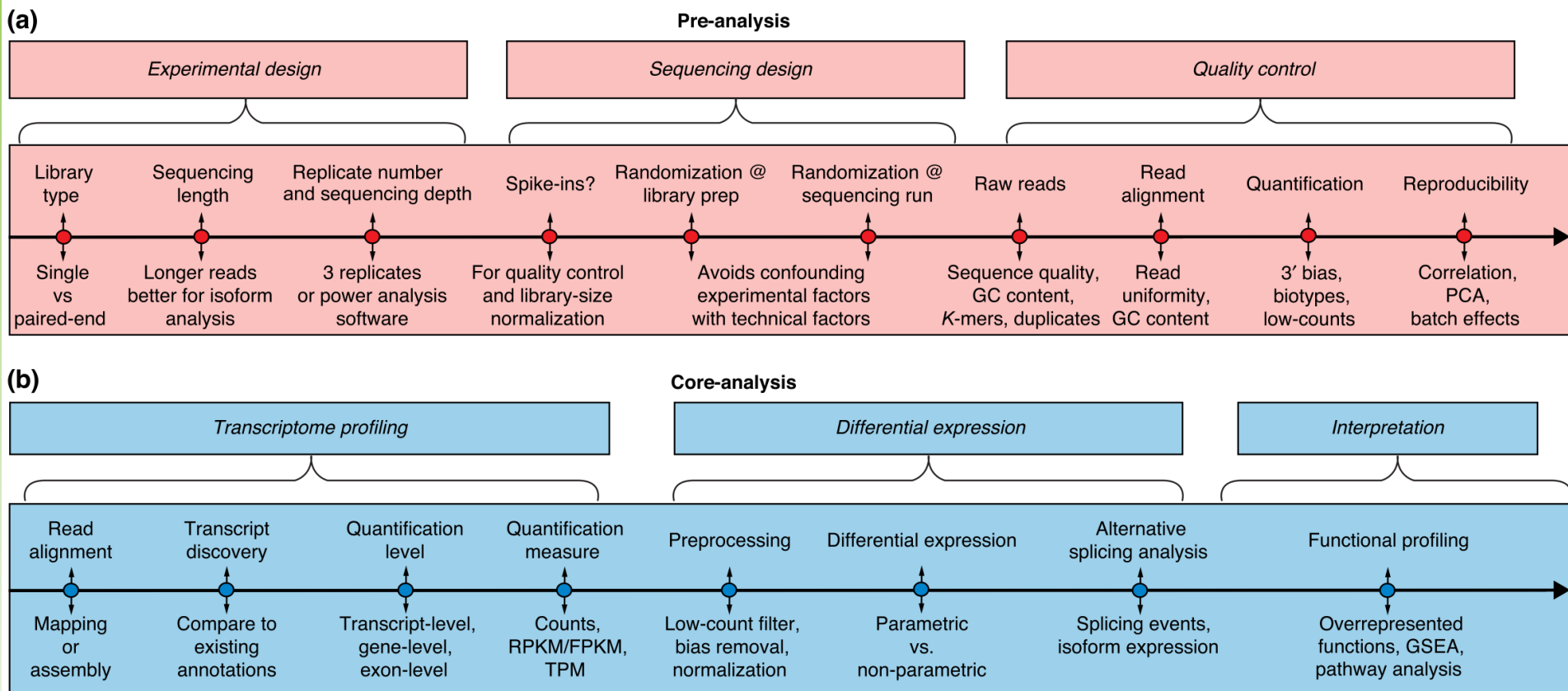
Consider Hifi reads for annotation of a new transcriptome

You should profile ~15 million full length transcripts at average 2kb length

Applications and data scale



RNA-seq overview



Experimental design

- Beware of batch effects
 - circadian rhythm (time of day)
 - sampling method
 - sequencing strategy (adapters, lanes)
- Record all information and check for co-variation
- Three replicates to keep statistical analysis options open

Sample preparation - isolating mRNA

Classes of RNA Molecules in Human Cells

- ▶ **Ribosomal RNA - rRNA**
 - 28 S
 - 18 S
 - 5.8 S
- ▶ **Non-coding RNA - ncRNA**
 - tRNA
 - snoRNA
 - lincRNA
 - miRNA
 - Many, many others...
- ▶ **Mitochondrial RNA - mtRNA**
- ▶ **Messenger RNA - mRNA**
 - Highly expressed transcripts (>10,000 copies per cell)
 - Rarely expressed transcripts (~1 copy per cell)

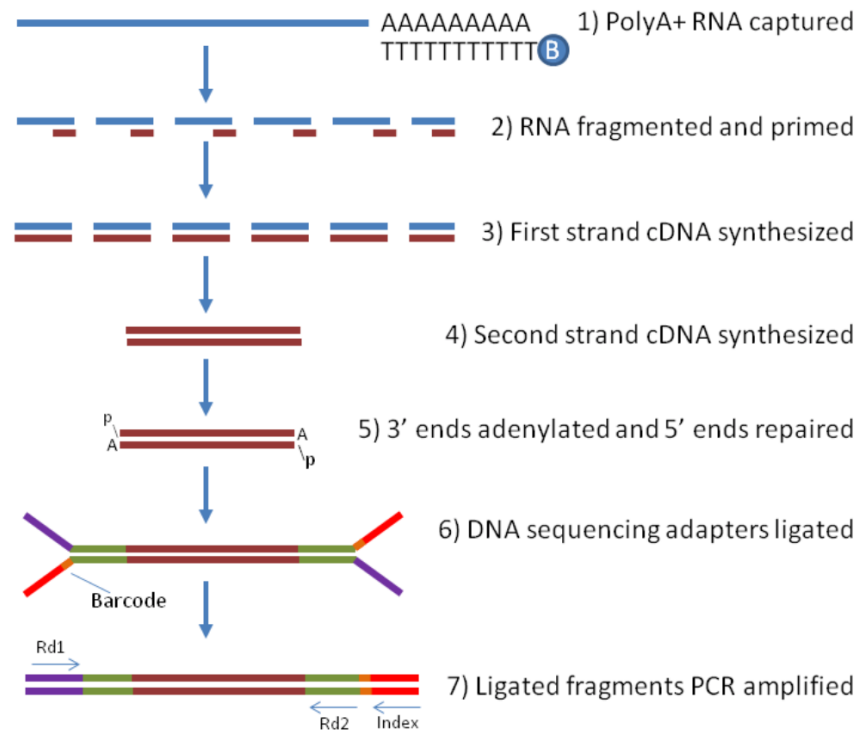
Enrichment of RNA or cDNA samples by selection or counterselection

Hybridization based

- ▶ oligo-dT
 - *selects pA+ RNA*
- ▶ Ribominus (Plus)
 - *counterselection of ribosomal RNAs (Plus includes mitochondrial rRNAs)*
- ▶ Capture arrays/bead capture
 - *for selection of an array of 'hand-picked' targets*
- ▶ Custom-made selection or counterselection strategies
 - *e.g. through hybridization to biotinylated RNA or DNA molecules*
- ▶ DSN-normalization
 - *elimination of cDNAs arising from highly abundant RNAs in total RNA samples by a differential hybridization approach coupled with duplex-specific nuclease*

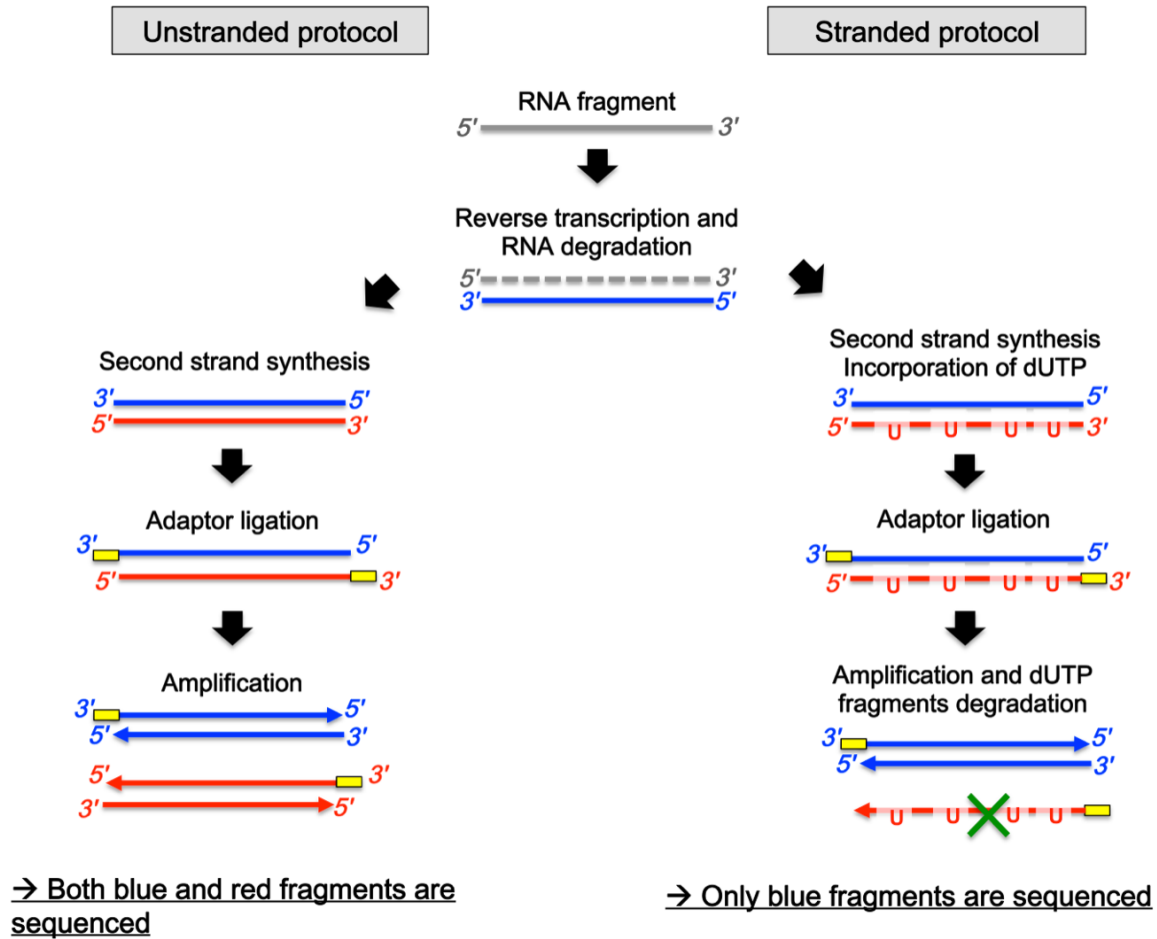
Sample preparation – Illumina sequencing

Illumina TruSeq stranded poly(A) library preparation steps

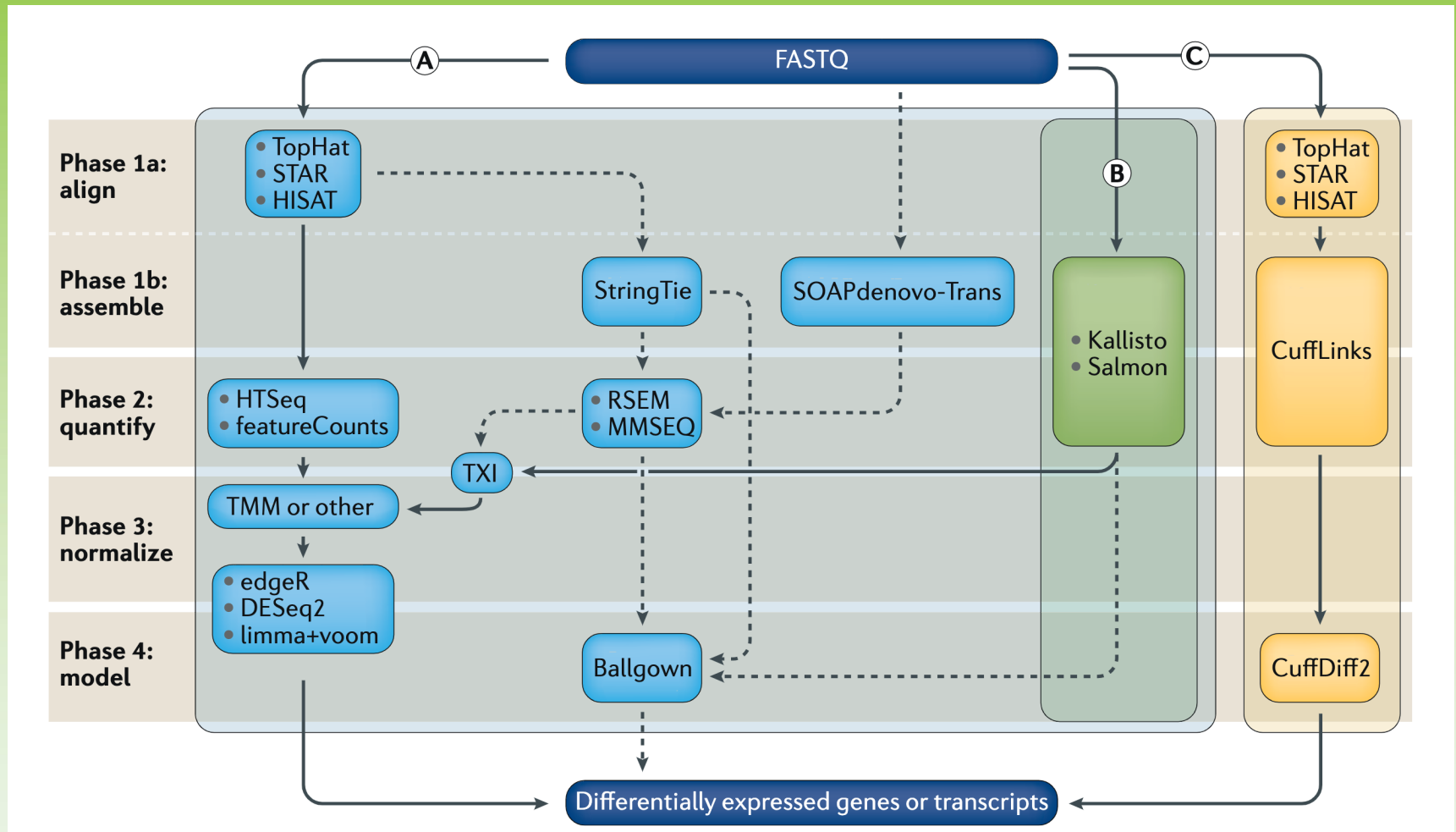


from <https://www.labome.com/method/RNA-seq.html>(<https://www.labome.com/method/RNA-seq.html>)

Stranded RNA-seq

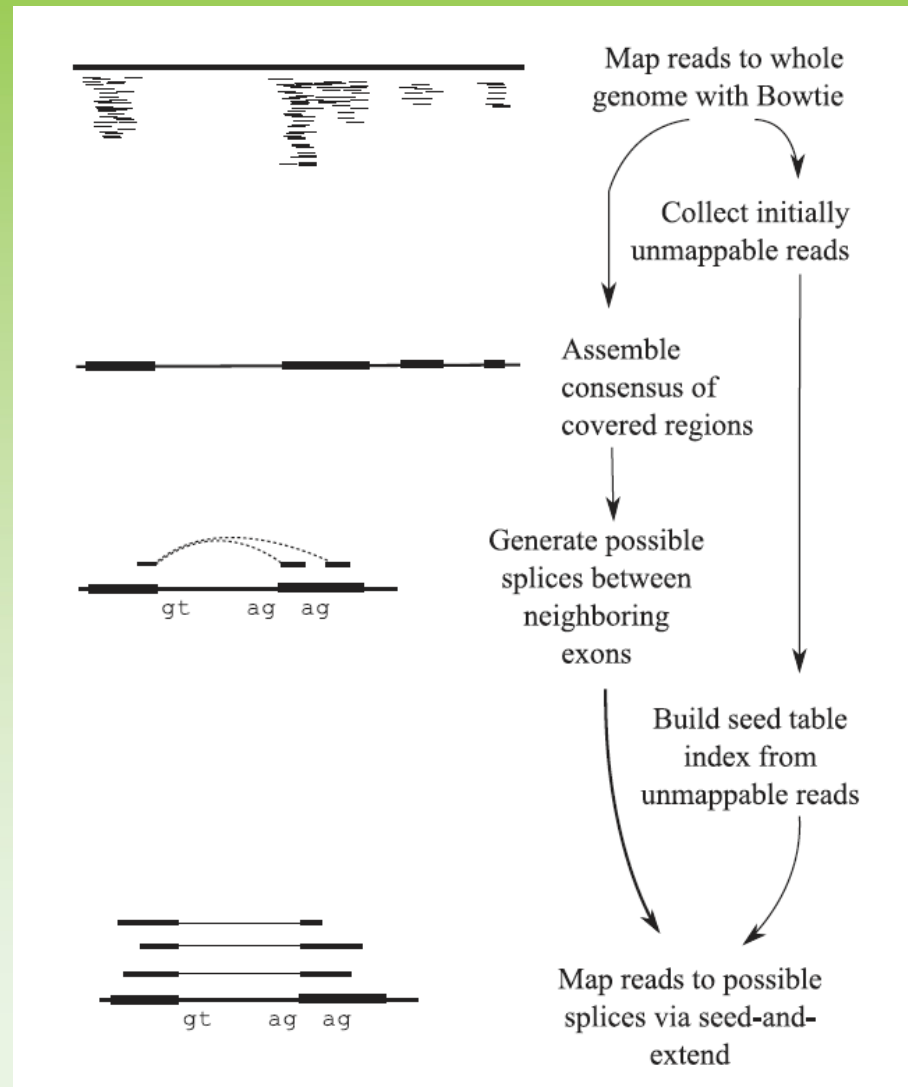


Differential gene expression analysis



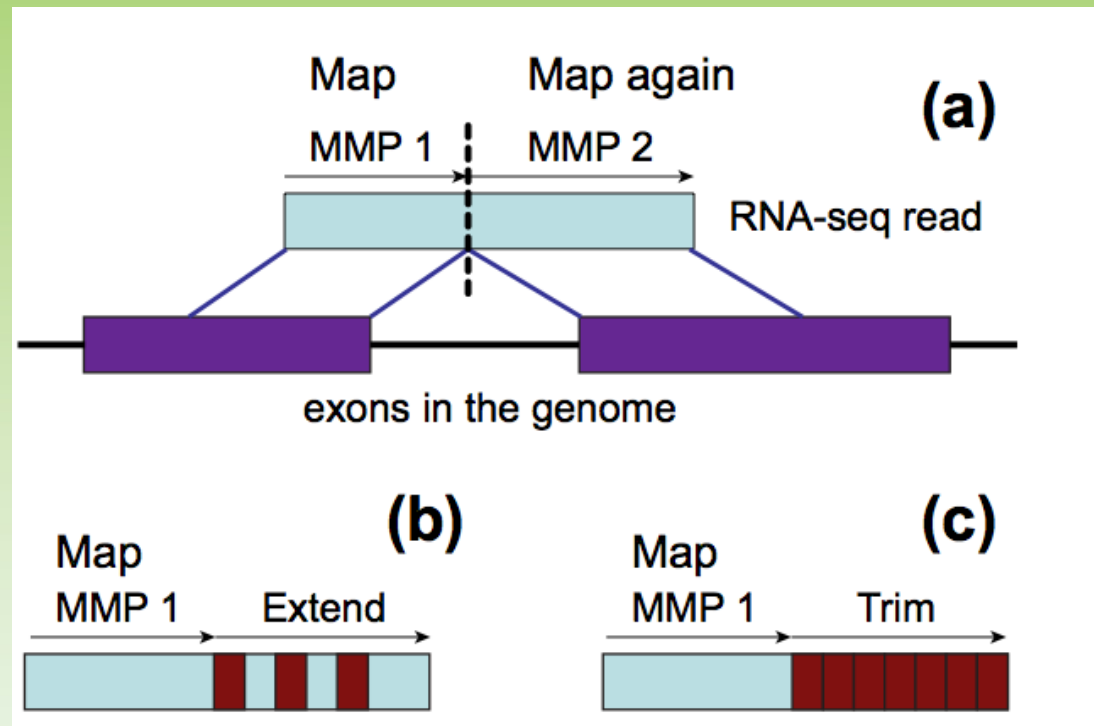
Spliced read alignment problem

Tophat
> 3000 citations

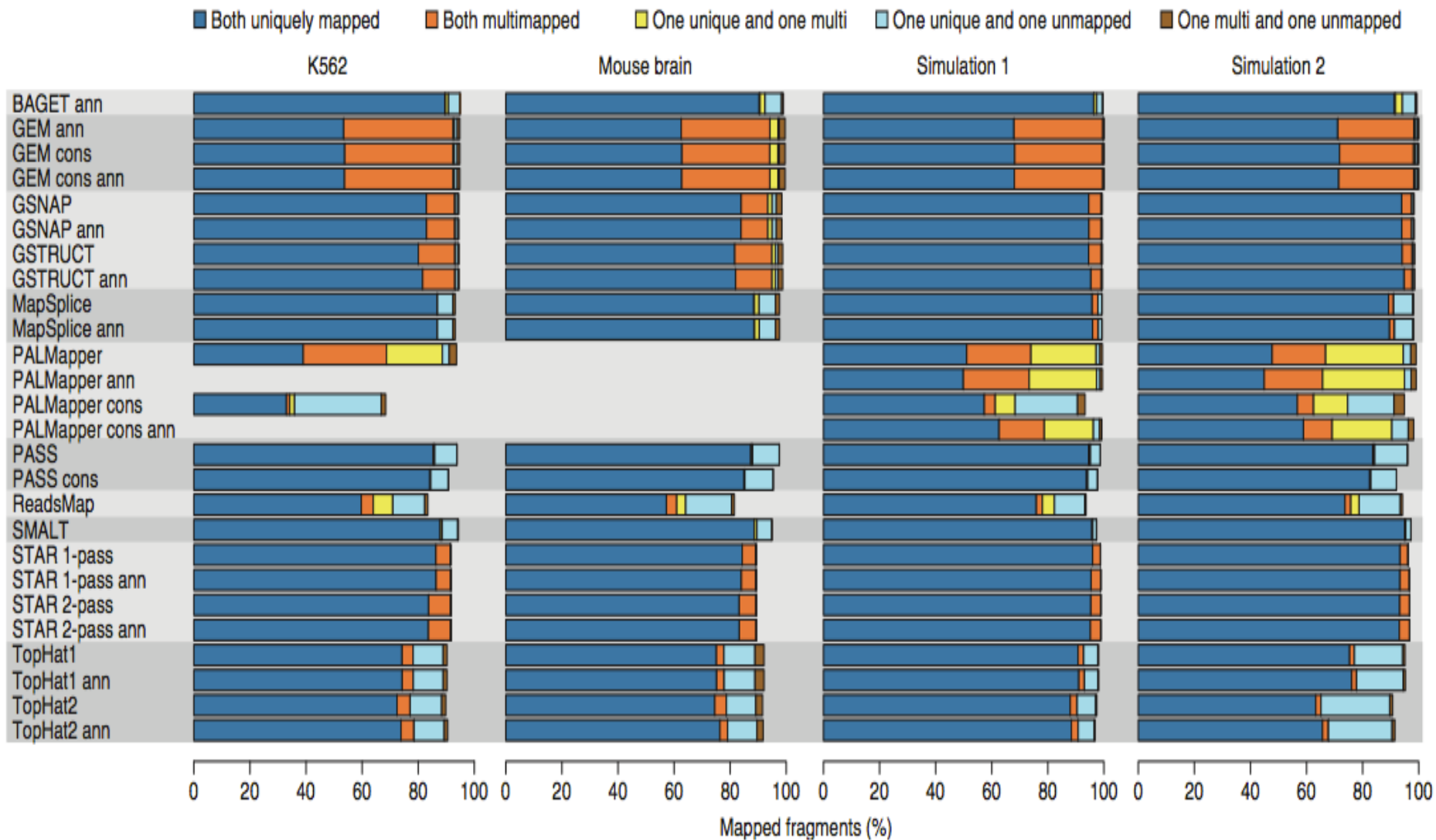


Spliced read alignment problem

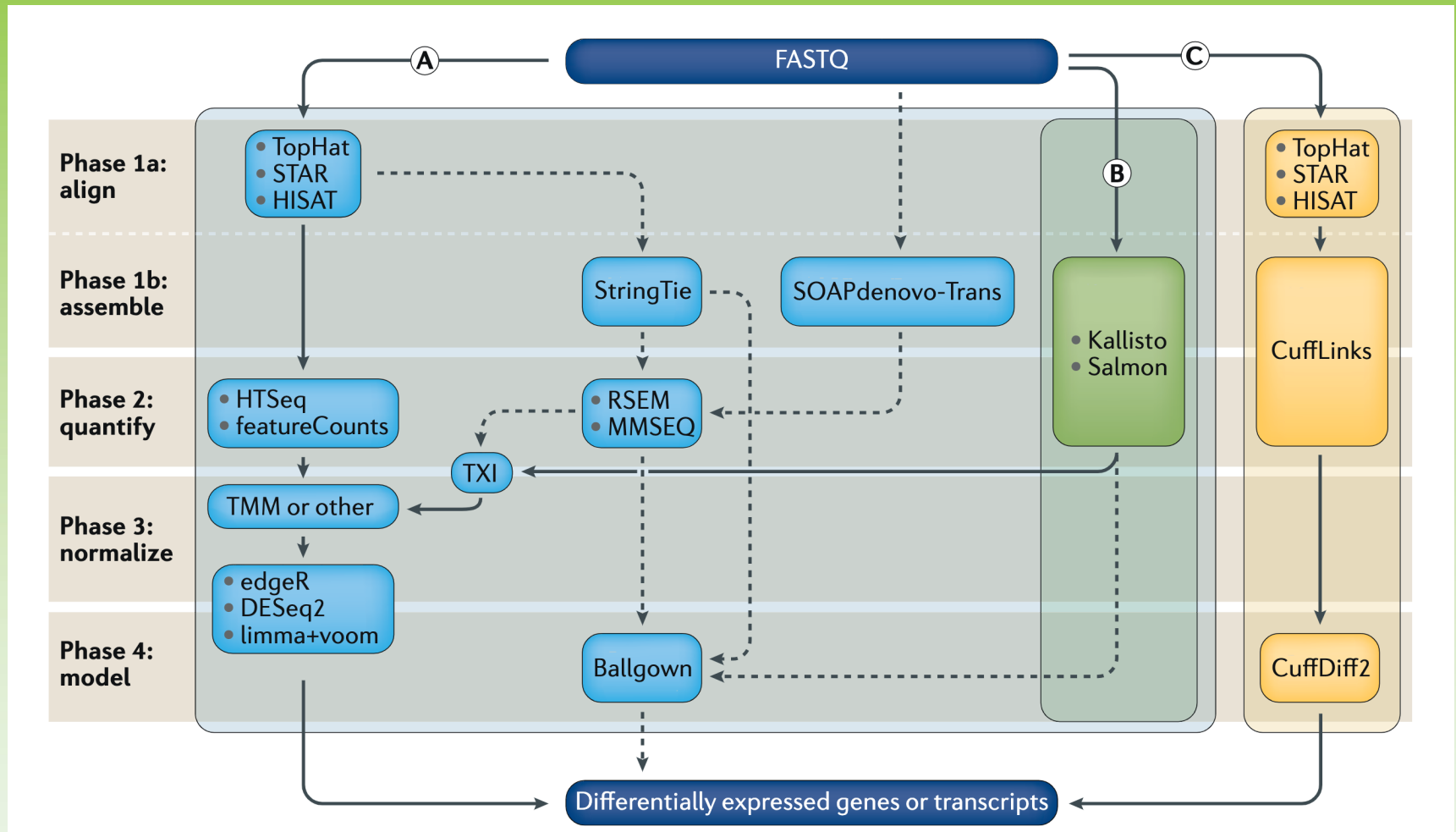
STAR >2000 citations



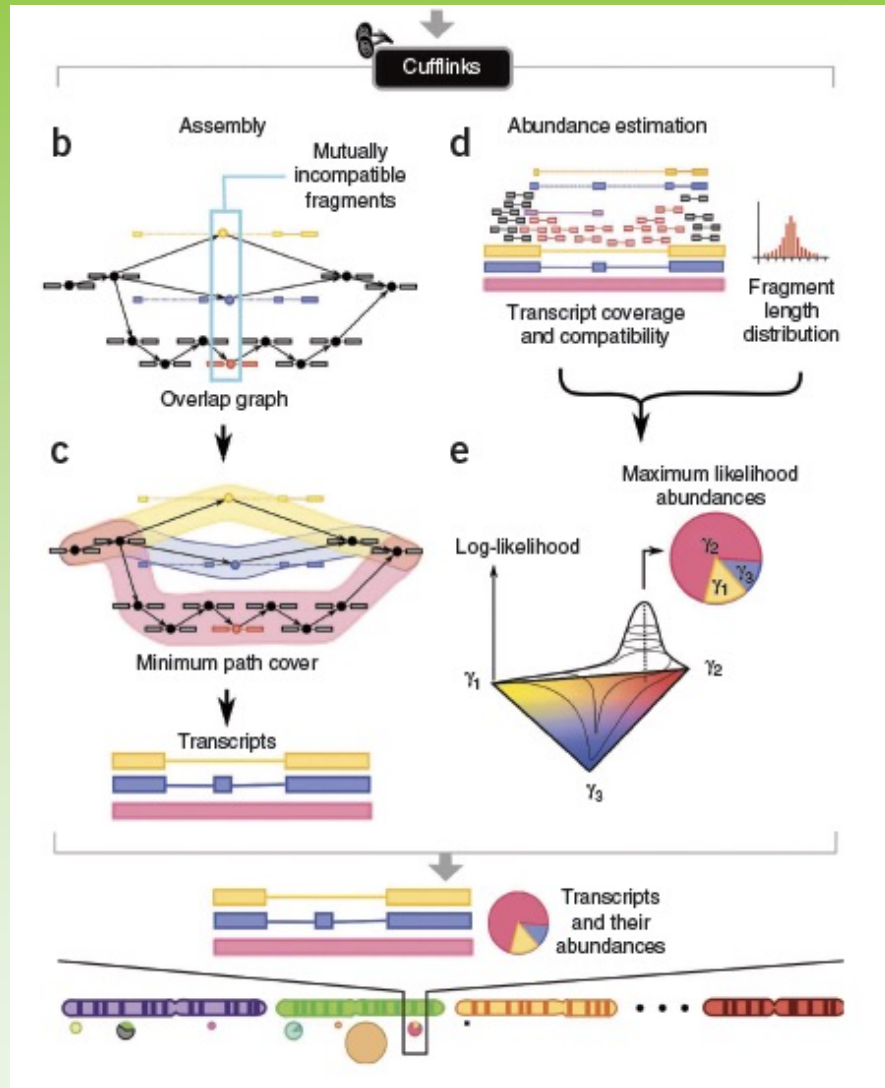
Spliced read aligners



Differential gene expression analysis



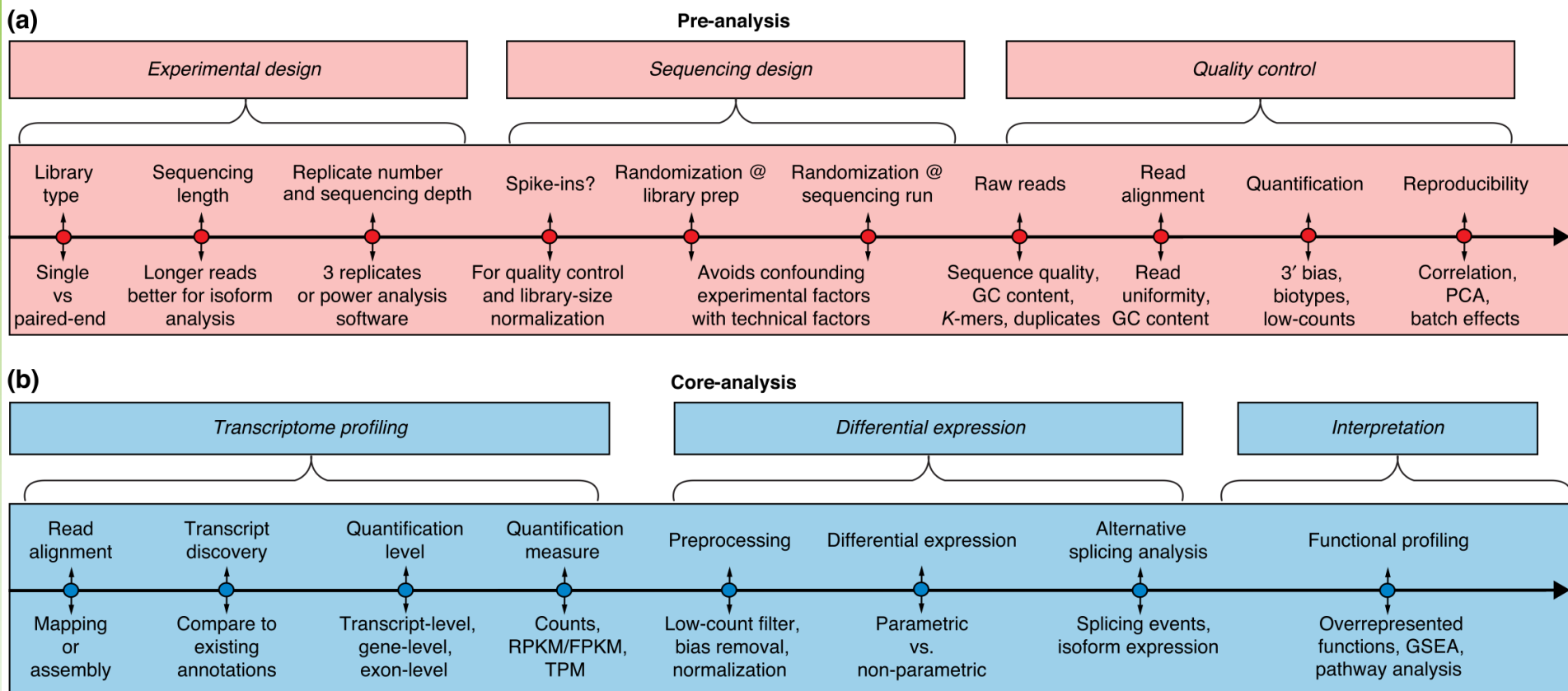
Cufflinks isoform quantification



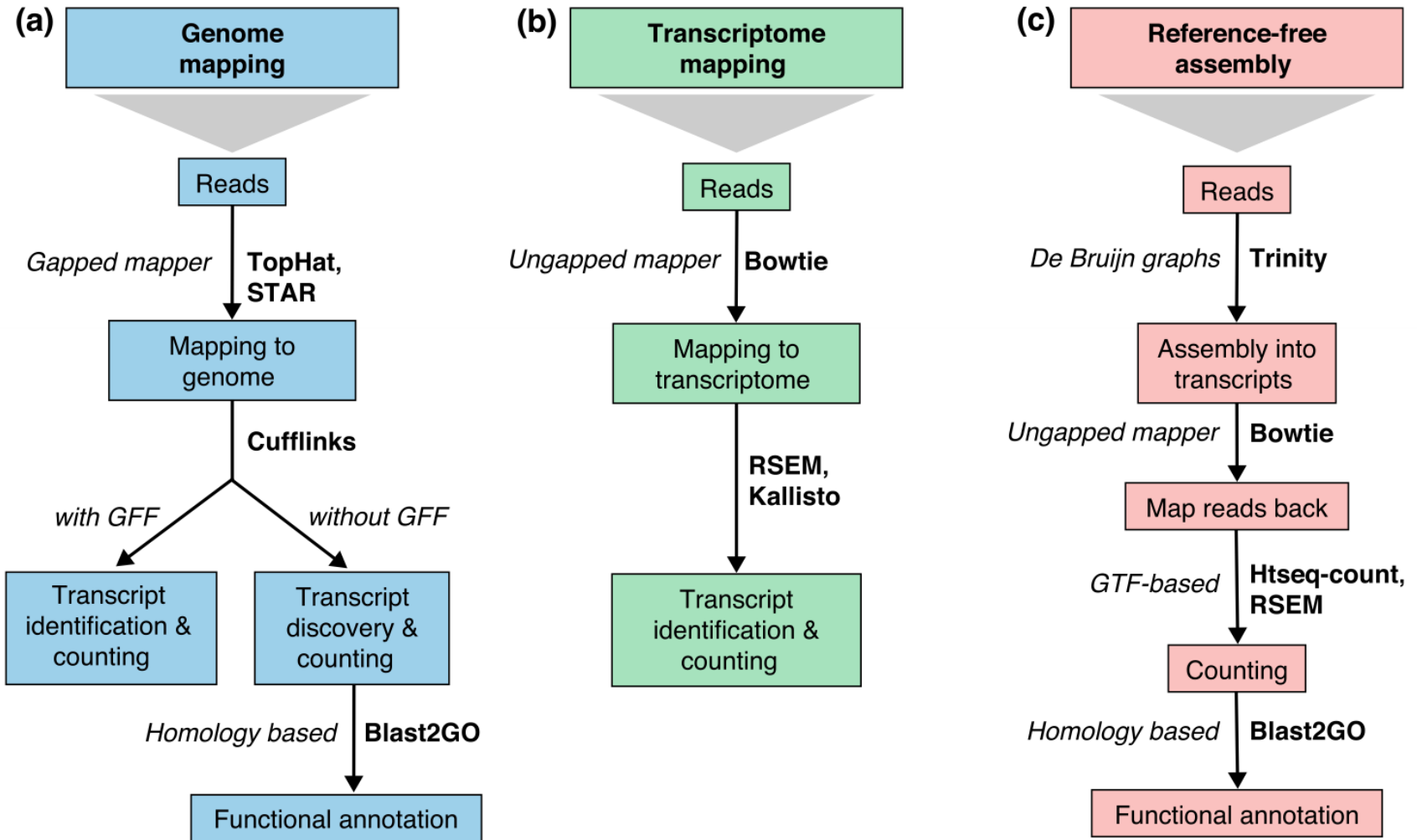
The *de novo* way

- Trans-ABYSS (PMID: 20935650)
- Trinity (PMID: 21572440)
- Mapping reads more sensitive
- *de novo* has no problems with intron/exon structure
- With long reads, *de novo* assembly of transcripts is becoming obsolete (most transcripts are significantly shorter than HiFi reads).

RNA-seq overview



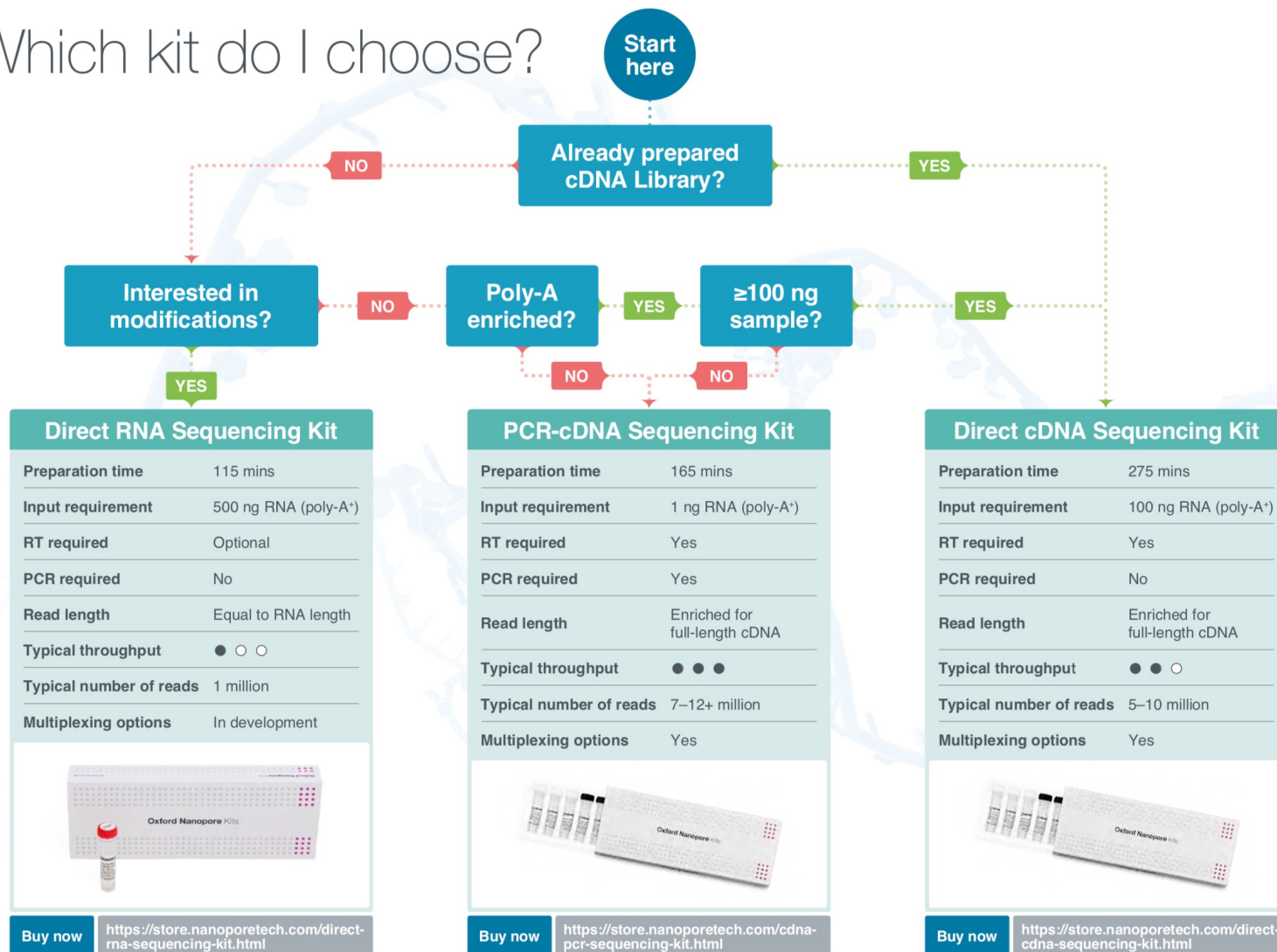
RNA-seq overview



Nanopore RNA-seq

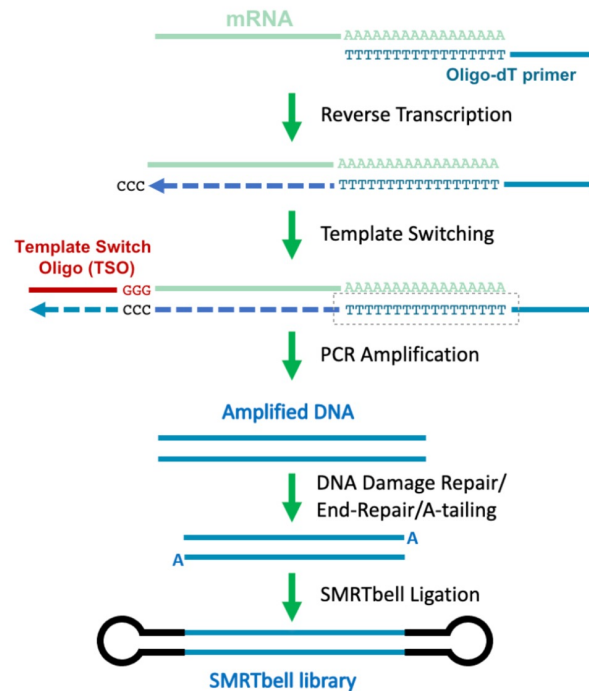
Which kit do I choose?

Start here



Iso-Seq Express Kit [2]

- Input 60-300 ng total RNA
- Full-length cDNA
- Multiplexing support



**LIBRARY
PREP
1 DAY**









Sequel II System

- 1 SMRT Cell 8M for whole transcriptome
- Up to 4 million full-length reads
- Accuracy 99-99.9%

	Unique Genes	Unique Transcripts	Unique ORFs
Single Cell, Human Brain Organoid	14,737	60,815	34,697
Single Cell, Human Cell Line	17,767	237,951	89,399
Bulk, UHRR [3]	16,328	183,689	60,649
Bulk, Alzheimer Brain [4]	17,670	162,290	80,539

Main Bioinformatics Tools

		Input	Output
Sequencing			subreads.bam
Iso-Seq Analysis		subreads.bam or ccs.bam	Collapsed unique transcripts (GFF, FASTA)
Transcript Classification		Unique transcripts Reference genome Annotation (GTF) CAGE Peak Junction data...	Transcript classification Junction classification Figures
Functional Annotation		SQANTI output	Annotated GTF
Differential Analysis		Experimental design Annotated GTFs	