

# High Throughput Sequencing the Multi-Tool of Life Sciences

Lutz Froenicke

DNA Technologies and Expression Analysis  
Cores

UCD Genome Center

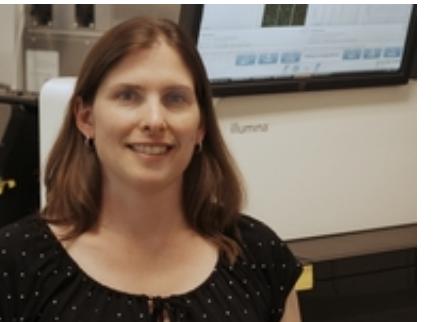
# Outline

- Who are we and what are we doing?
- Overview HTS sequencing technologies
- How does Illumina sequencing work?  
Sequencing library and run QC
- How does RNA-seq work?
- PacBio and Nanopore Sequencing
- Some cutting edge technologies & applications

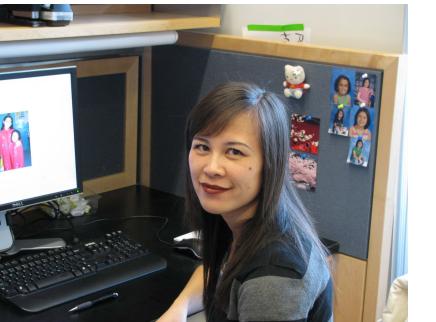
# DNA Technologies & Expression Analysis Cores

- HT Sequencing Illumina
- Long-Read & Linked-Read Sequencing  
PacBio, Oxford Nanopore, 10X Genomics
- HMW DNA isolation
- Illumina microarray (genotyping)
- Single-cell RNA-seq
- Consultations → Experimental Design  
**(Bioinformatics Core & DNA Tech Core)**
- introducing new technologies to the campus
- shared equipment
- teaching (workshops)

# The DNA Tech Core Team



Emily



Oanh



Diana



Siranoosh

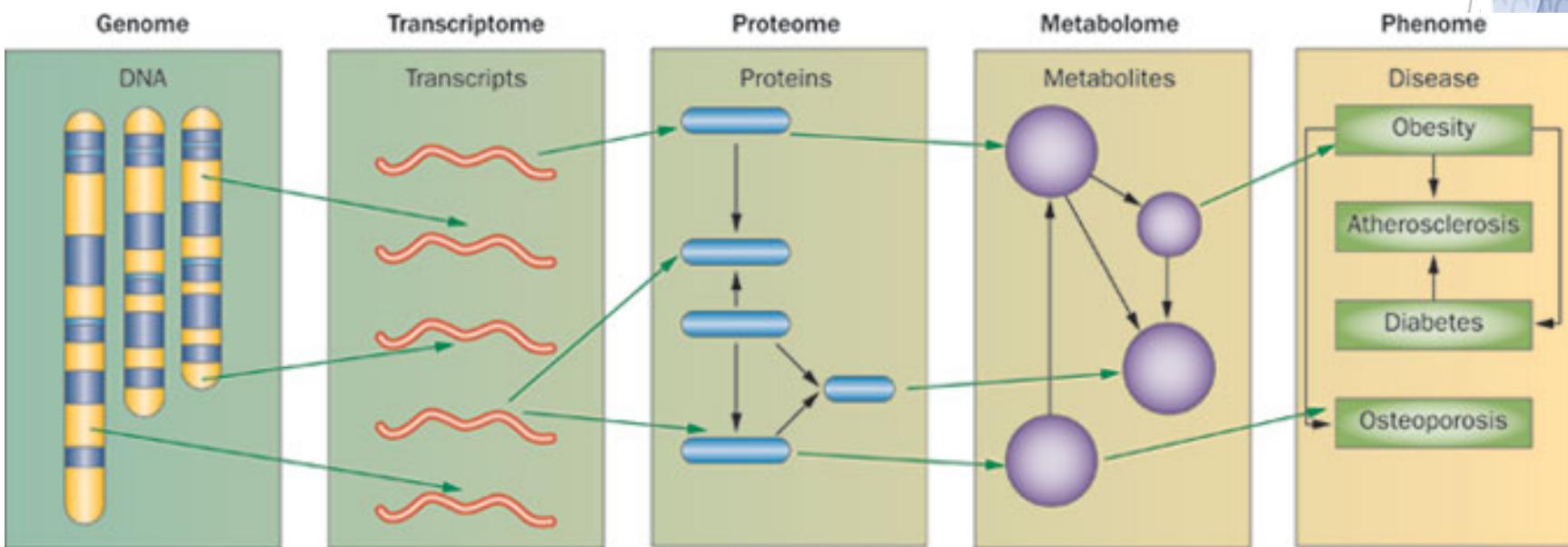


Vanessa



Ruta

# The UCD GENOME CENTER



DNA Tech & Expression Analysis   Proteomics Core   Metabolomics Core

**“DNA makes RNA and RNA makes protein”**

the Central Dogma of Molecular Biology; simplified from Francis Crick  
1958

**nature**  
REVIEWS   **CARDIOLOGY**

MacLellan, W. R. et al. (2012) Systems-based approaches to cardiovascular disease  
*Nat. Rev. Cardiol.* doi:10.1038/nrccardio.2011.208

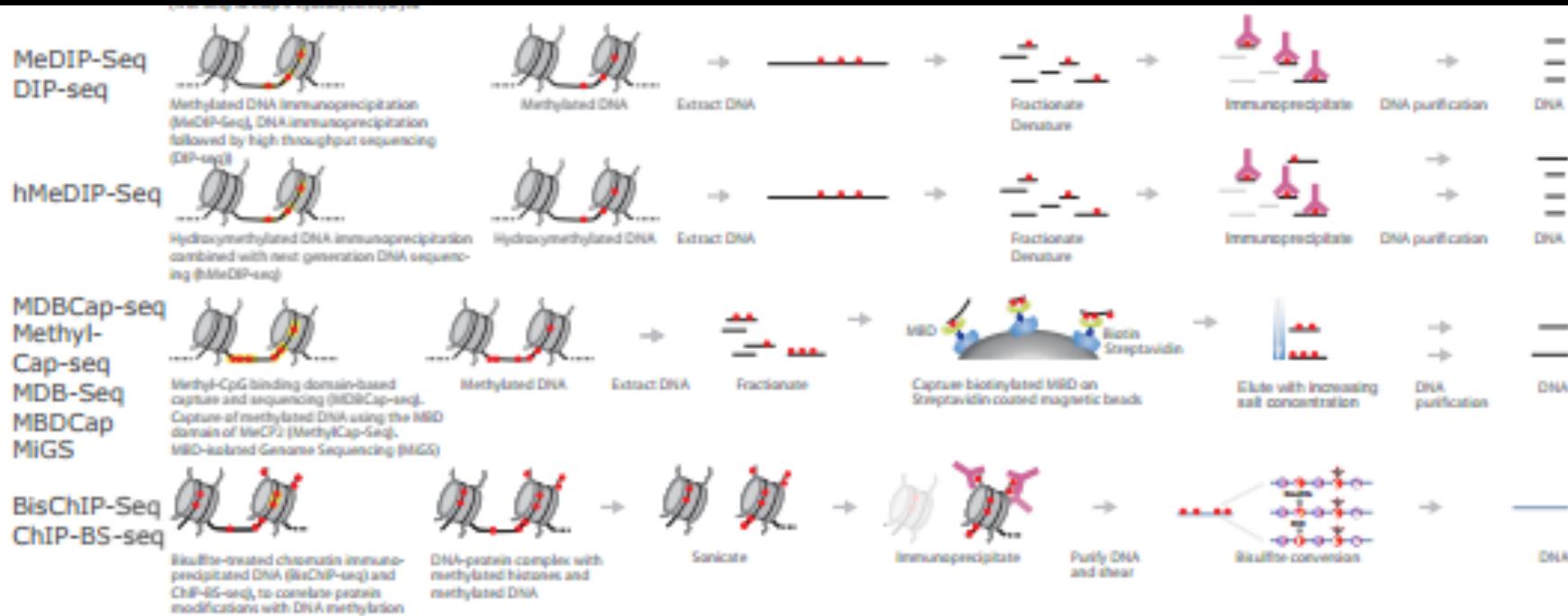
# Complementary Approaches

| Illumina   | PacBio CLR  | PacBio HiFi   | PromethION Nanopore   |
|--|---|---|---|
| Still-imaging of clusters (~1000 clonal molecules) | Video recordings fluorescence of single molecules | Video recordings fluorescence of single molecules                 | Recording of electric current through pores                           |
| Short reads - 2x300 bp<br>MiSeq                    | Up to 70 kb, N50 25 kb                            | Up to 20 kb, N50 18   | Up to 100 kb, N50 30 kb   |
| Repeats are mostly not analyzable                  | spans retro elements                              | accurate enough to assemble through REs                           | spans retro elements  |
| High output - up to 2.4 Tb per lane                | up to 100 Gb per SMRT-cell                        | up to 25 Gb HiFi data per cell                                    | Up to 100 Gb per flowcell   |
| High accuracy (< 0.5 %)                            | Raw data error rate 15 %                          | CCS data < 0.1%   | Raw data error rate 8 %   |
| Considerable base composition bias                 | No base composition bias                          | No base composition bias, but still mononucleotide repeat problem | Some systematic errors,   |
| Very affordable                                    | Costs 3 to 5 times higher                         | Costs 3 to 5 times higher   | Costs 2x higher   |
| De novo assemblies of thousands of scaffolds       | “Near perfect” genome assemblies                  | “Near perfect” genome assemblies; lowest error rate               | “Near perfect” genome assemblies with suppl. data; highest contiguity |

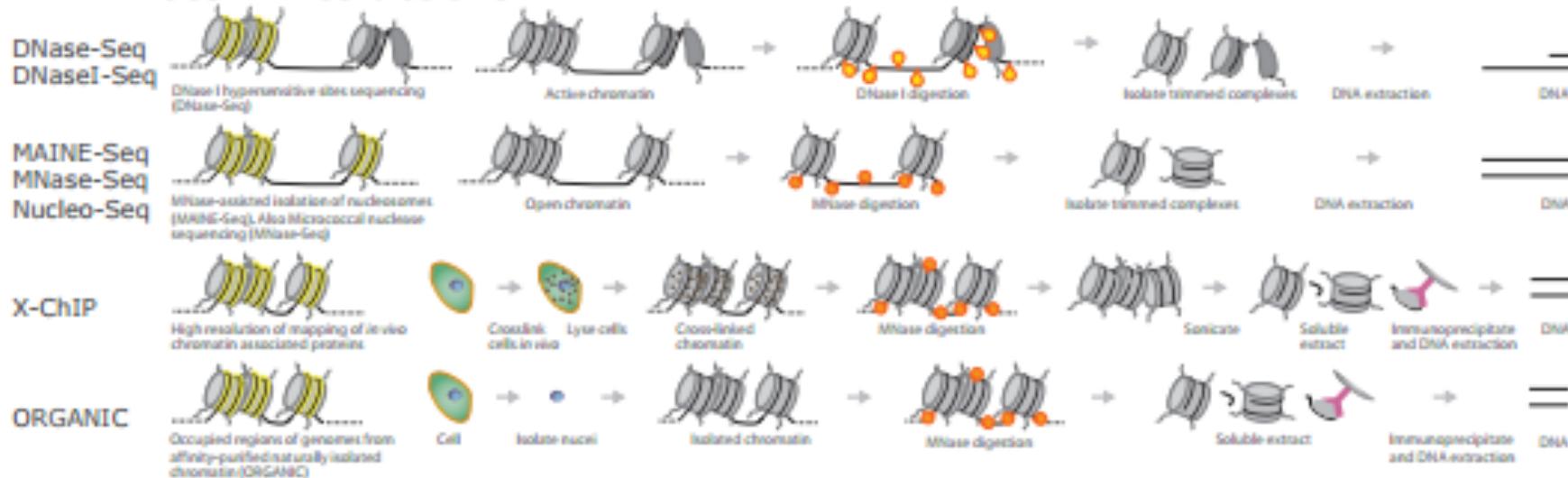
# High Throughput Short Read Sequencing: Illumina



- Whole genome sequencing & Exome sequencing:  
Variant detection (small variants SNPs and indels)  
Copy number variation (CNVs; prenatal diagnostics)
- Genotyping by sequencing
- Genome assemblies: small genomes
- Metagenomics
- RNA-seq: gene expression, transcript expression
- Small RNA-seq
- Single-cell RNA-seq
- Epigenetics: Methyl-Seq:
- ChIP-Seq (detecting molecular interactions)
- 3D Organization of the nucleus (Hi-C)



## DNA-Protein Interactions



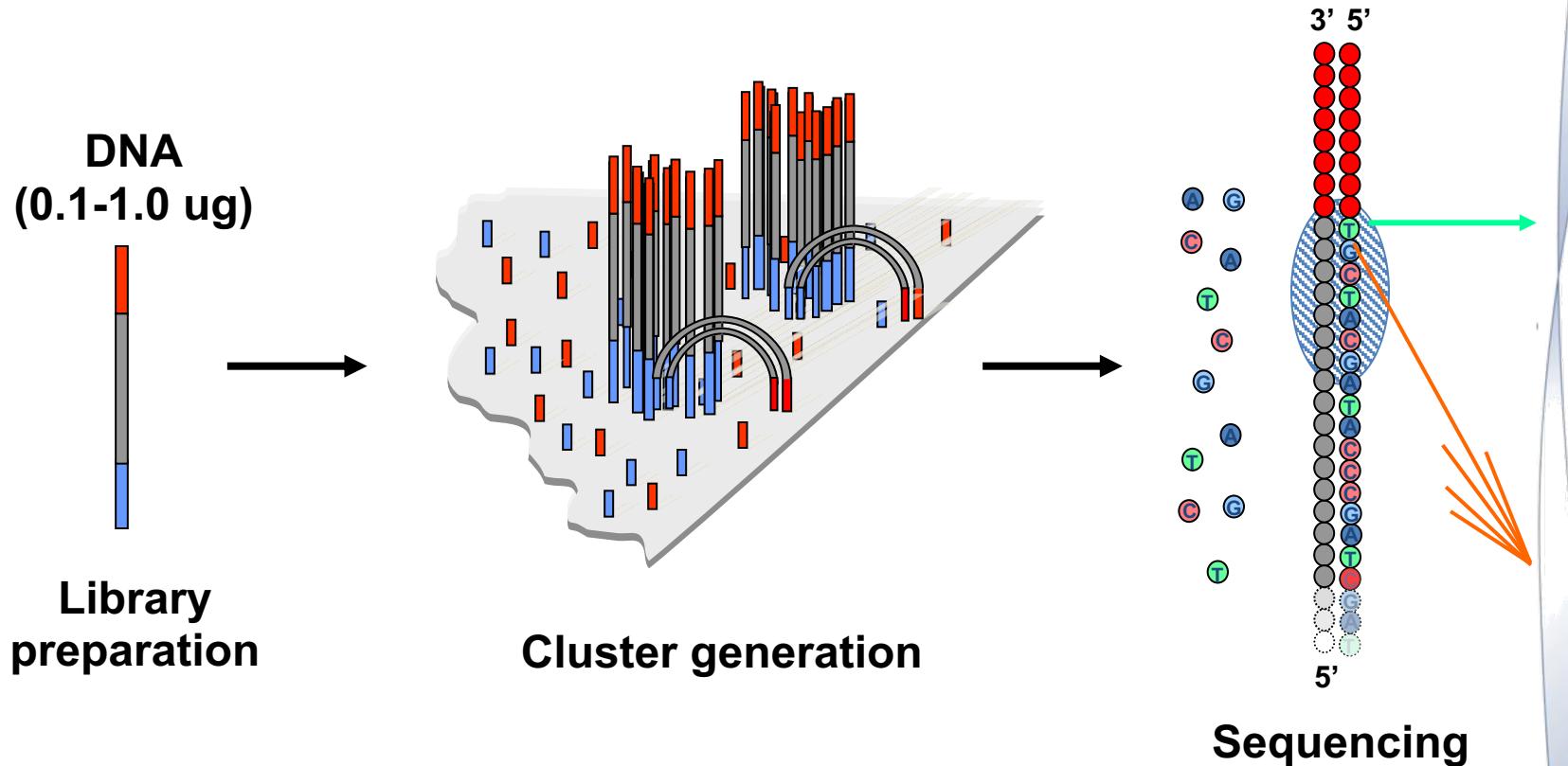
# Long Read Sequencing: PacBio and Nanopore



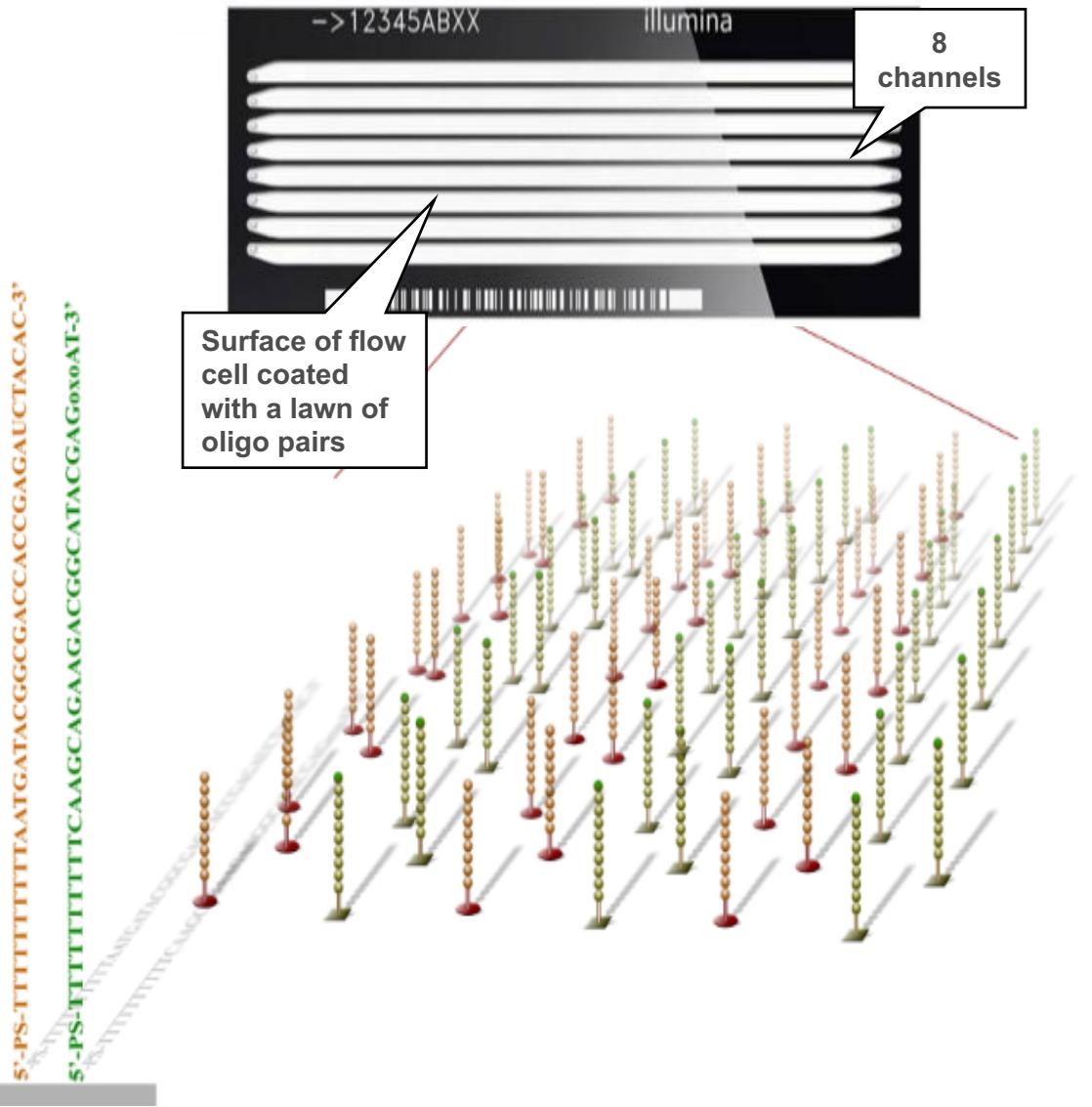
- Whole genome sequencing : Highest quality genome assemblies, Structural variant detection
- RNA-sequencing:
- full transcript data, Iso-form detection and quantification
- Direct RNA-seq identifies base modifications (Nanopore)
- Metagenomics
- Epigenetics (Nanopore: any modified bases, PacBio bacteria)

# Illumina Sequencing Technology

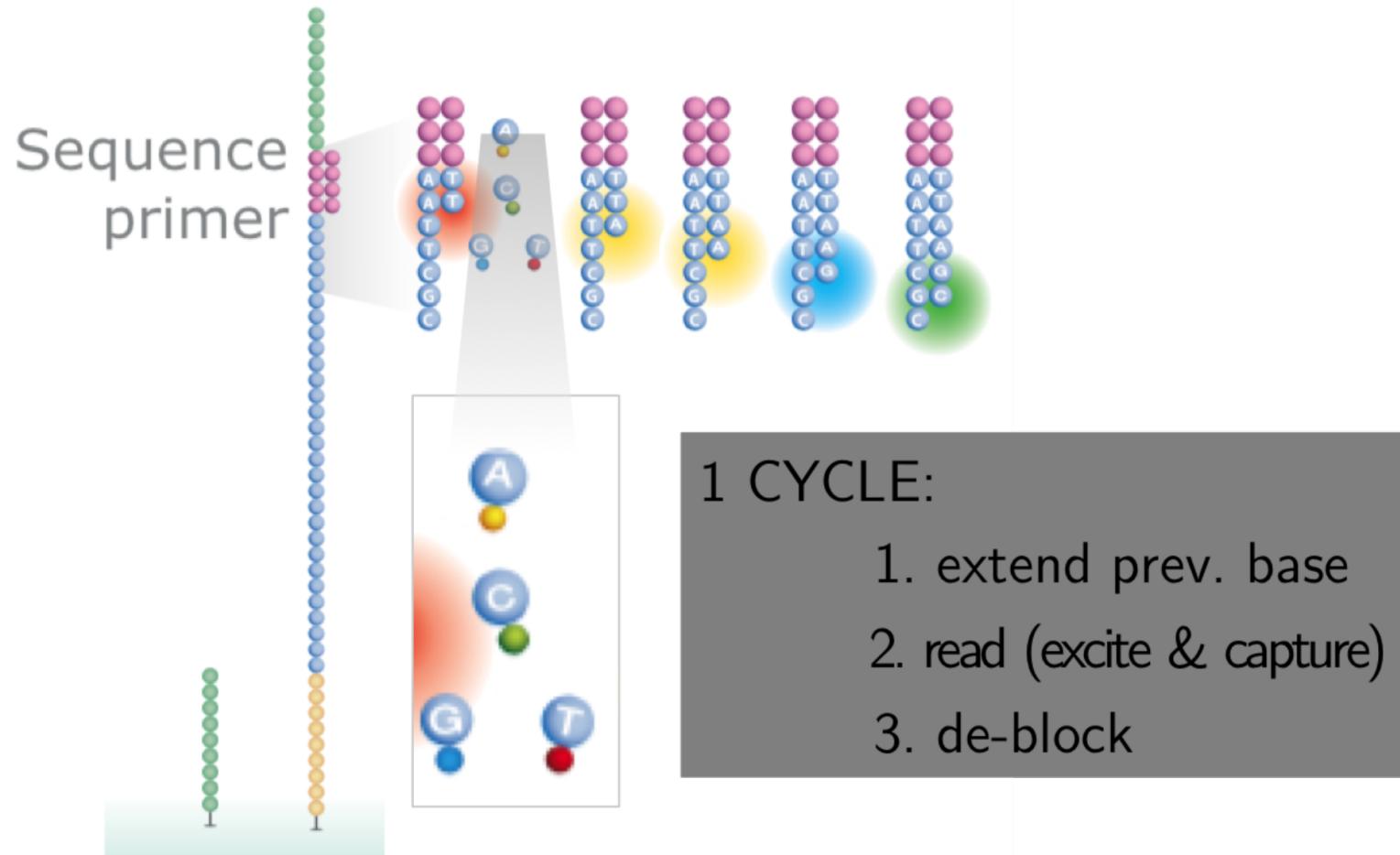
*Sequencing By Synthesis (SBS) Technology*



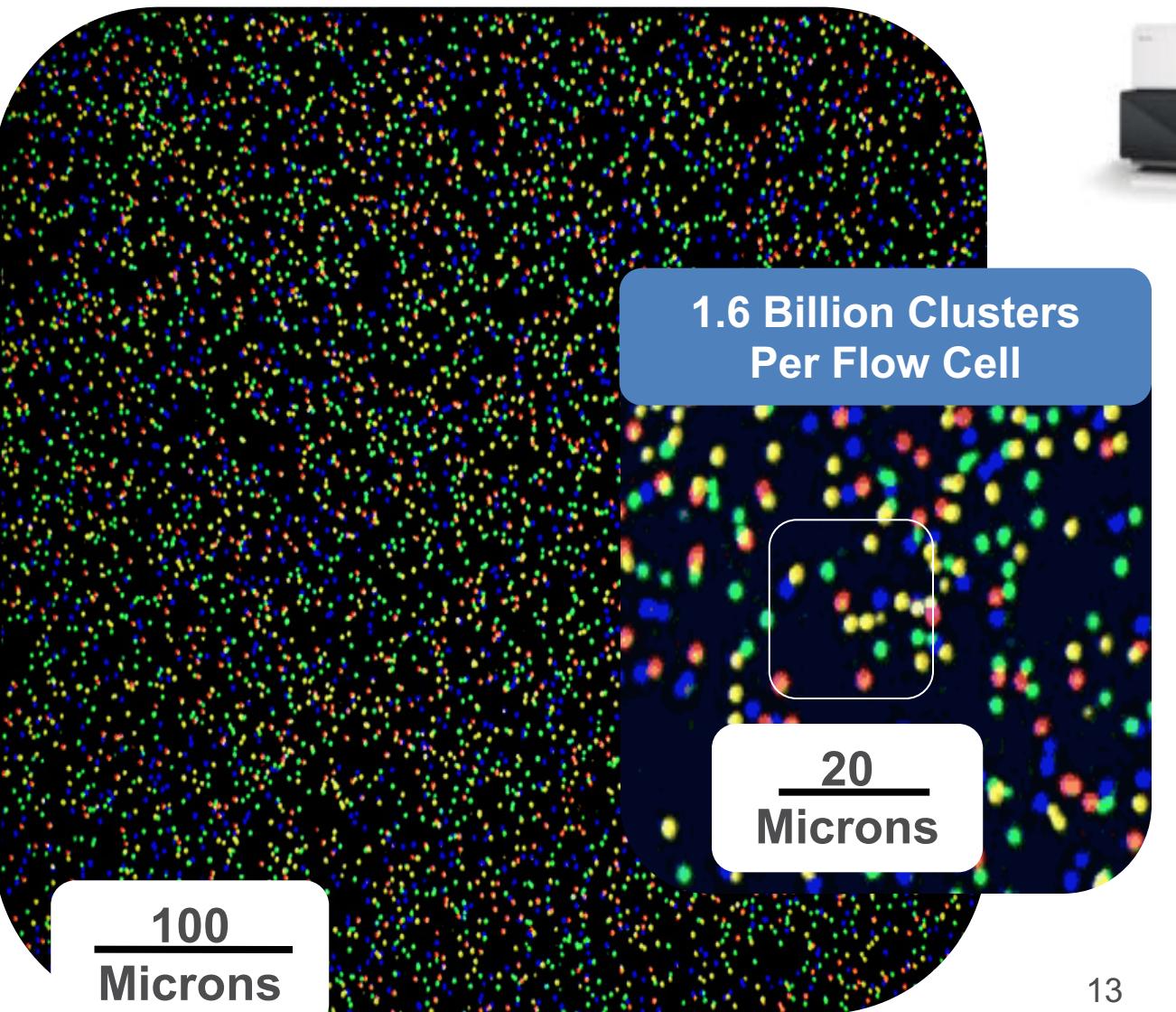
# TruSeq Chemistry: Flow Cell



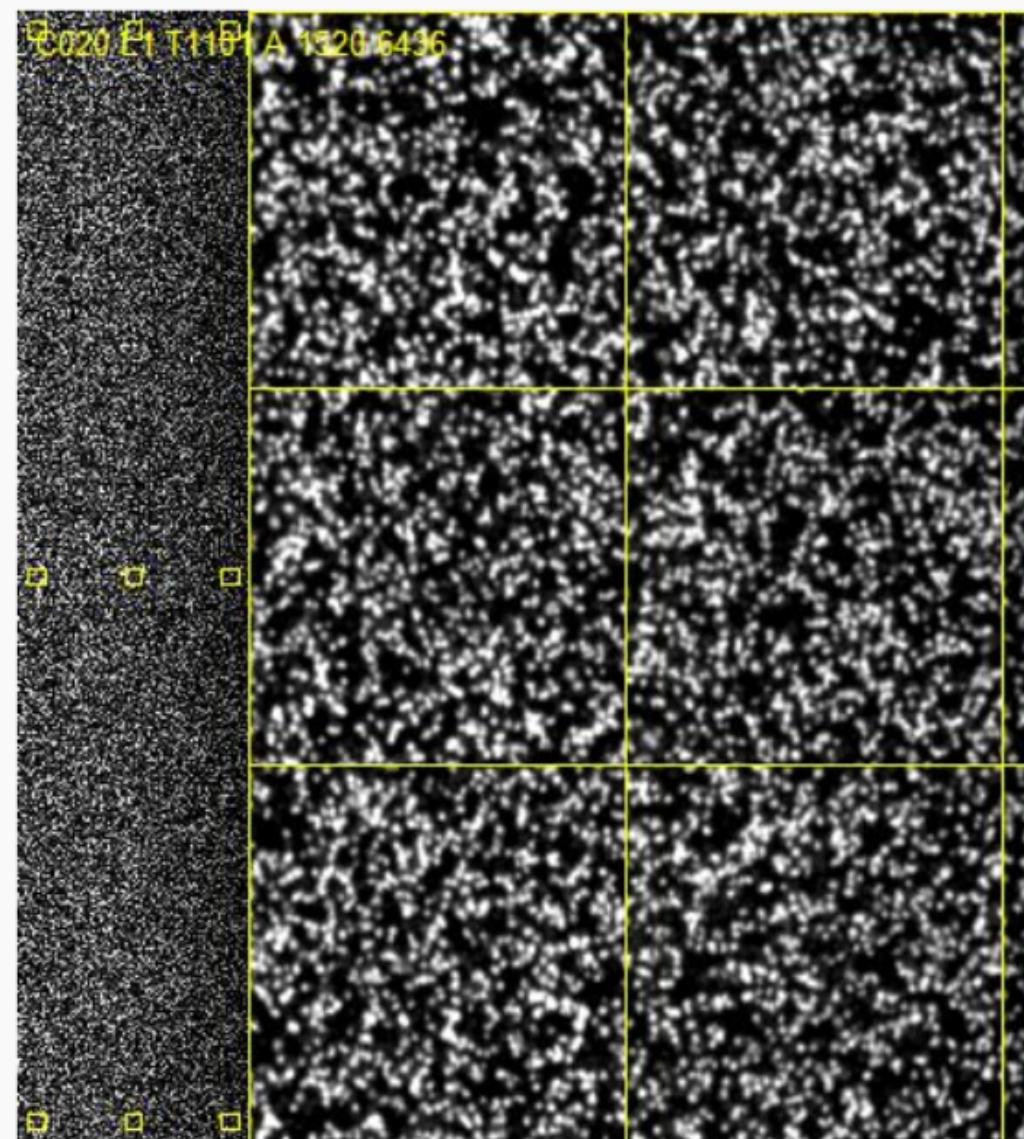
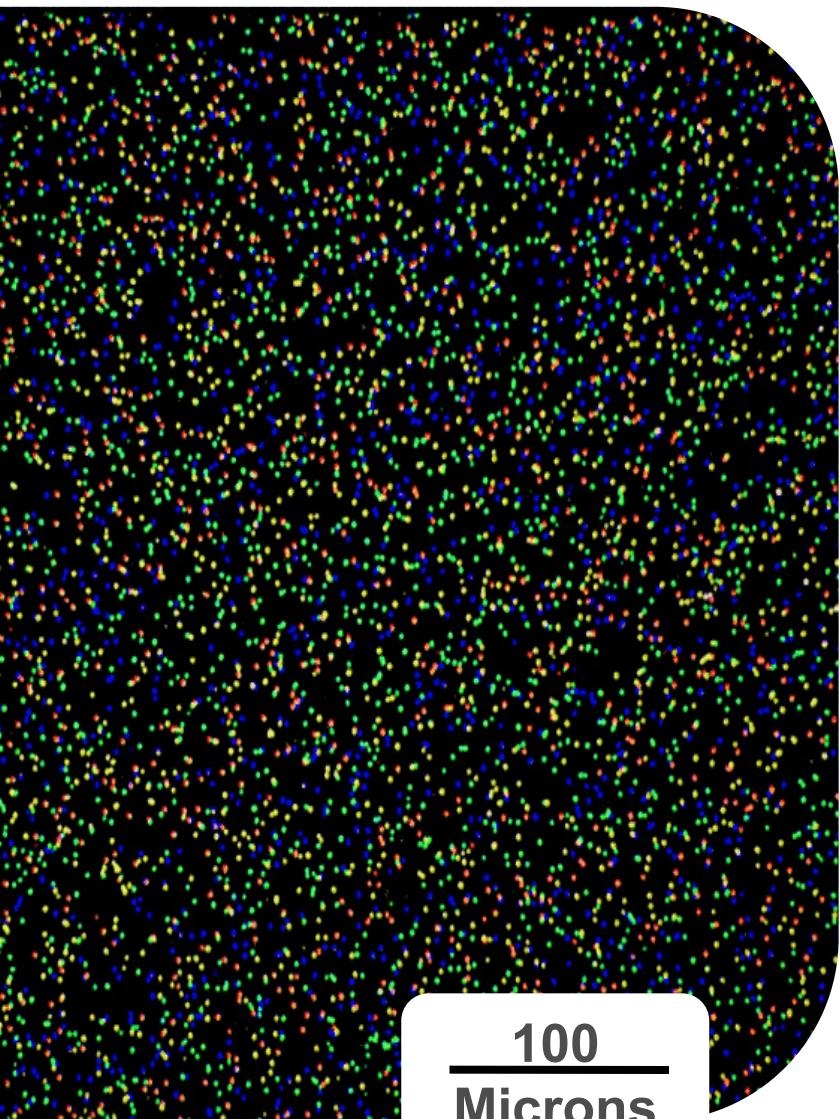
Illumina's sequencing is based on **fluorophore-labelled dNTPs** with **reversible terminator elements** that will become incorporated and excited by a laser one at a time.



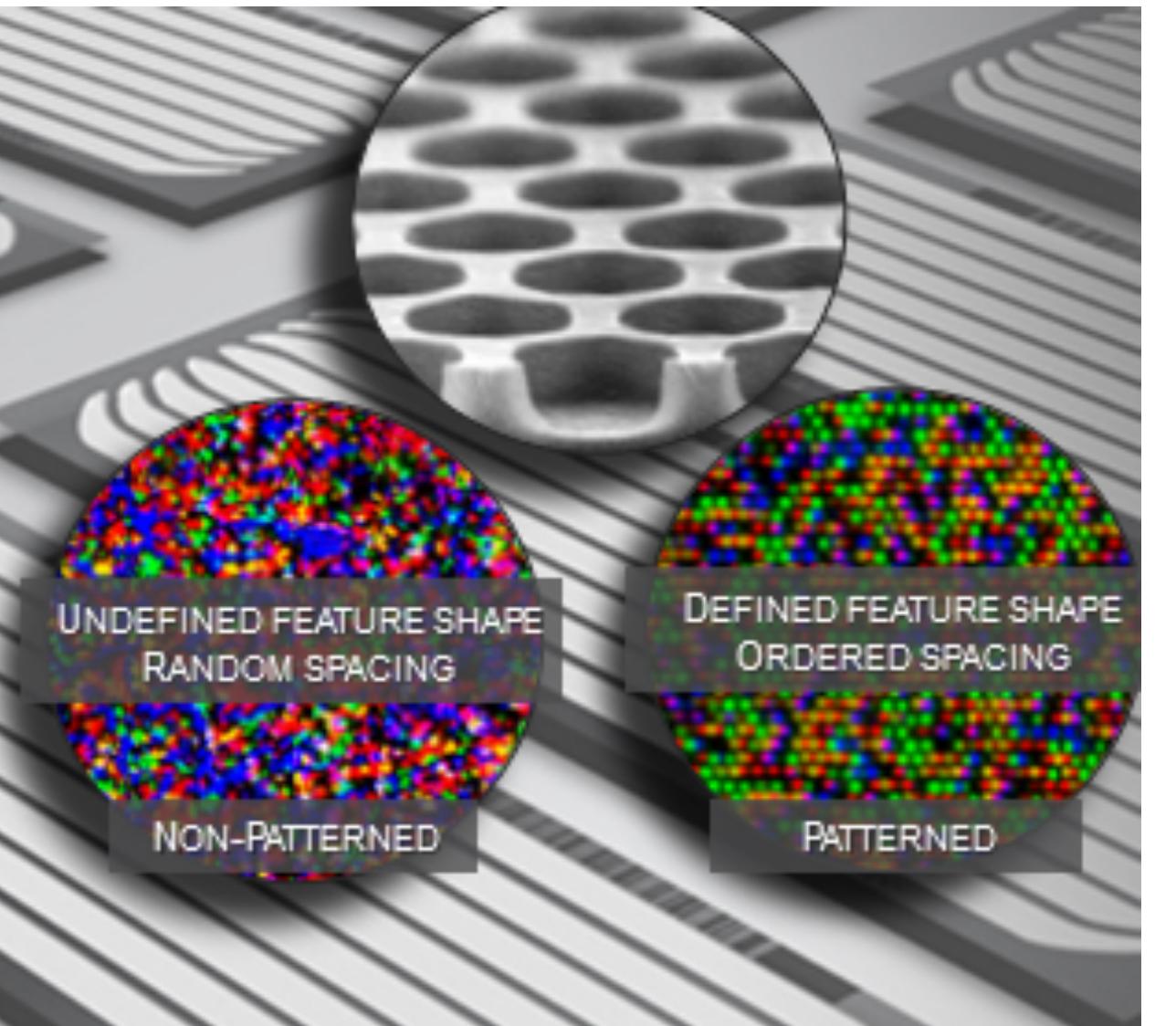
# False colored and merged four channel flowcell images



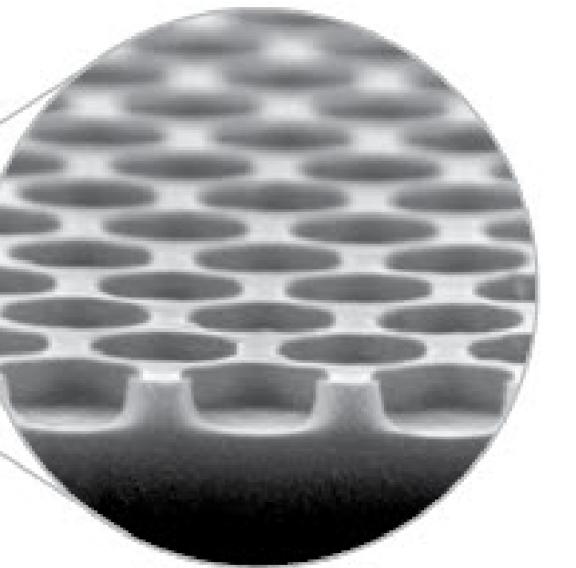
# B&W imaging



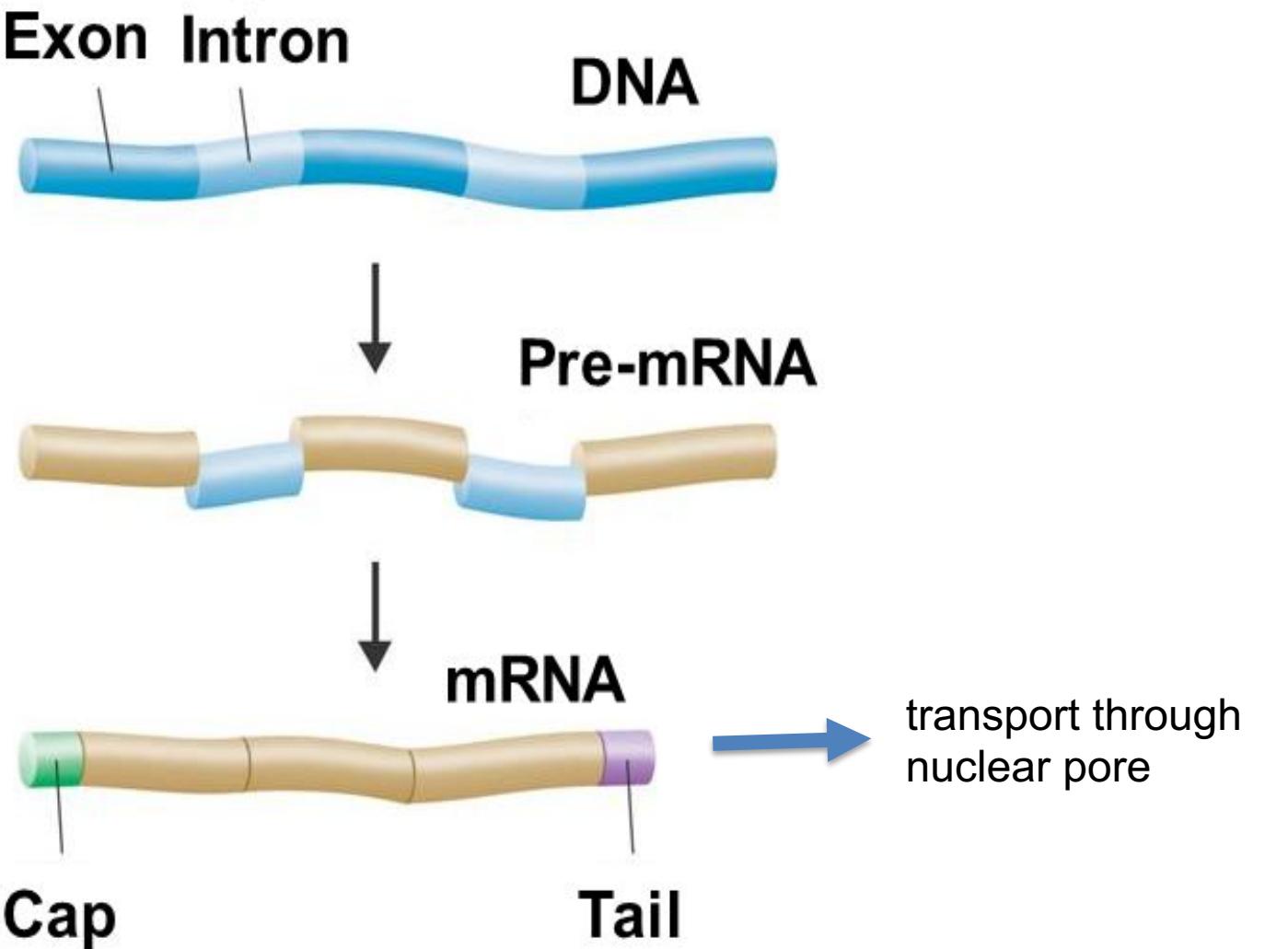
# Patterned Flowcell



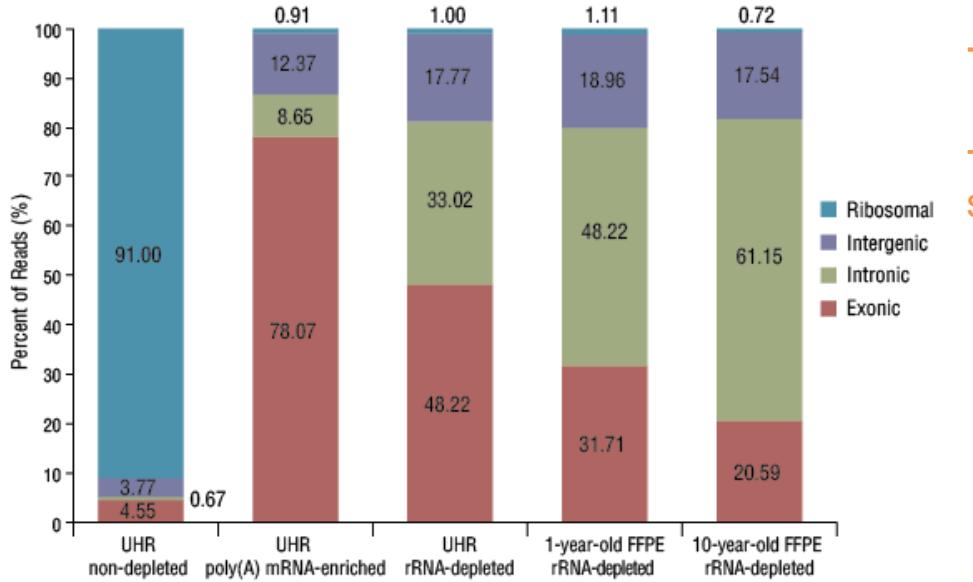
# Hiseq 4000: 478 million nanowells per lane



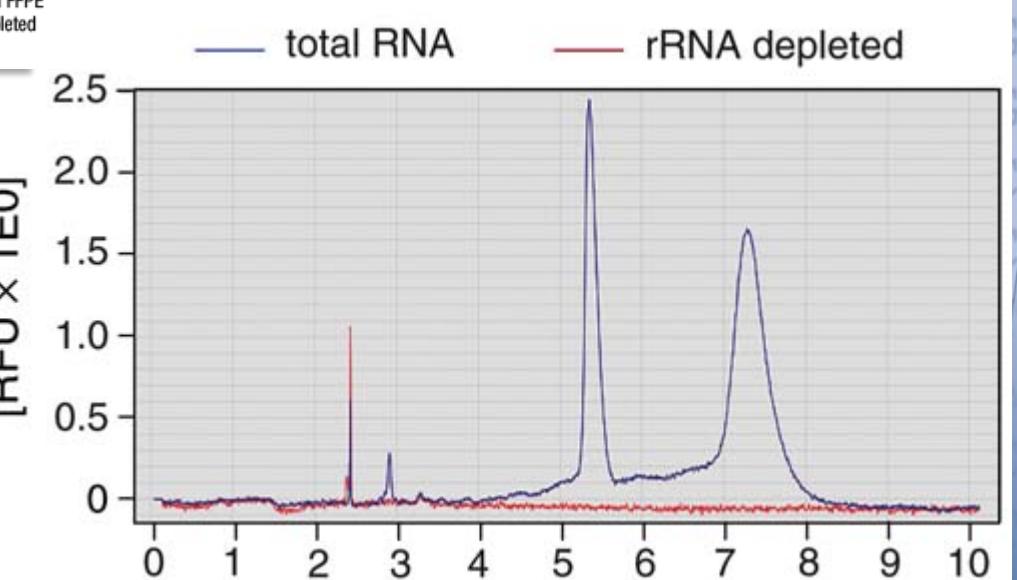
# transcription and processing in nucleus



# mRNA makes up only about 2% of a total RNA sample



- more than 90% rRNA content
- multiple other non-coding RNA species



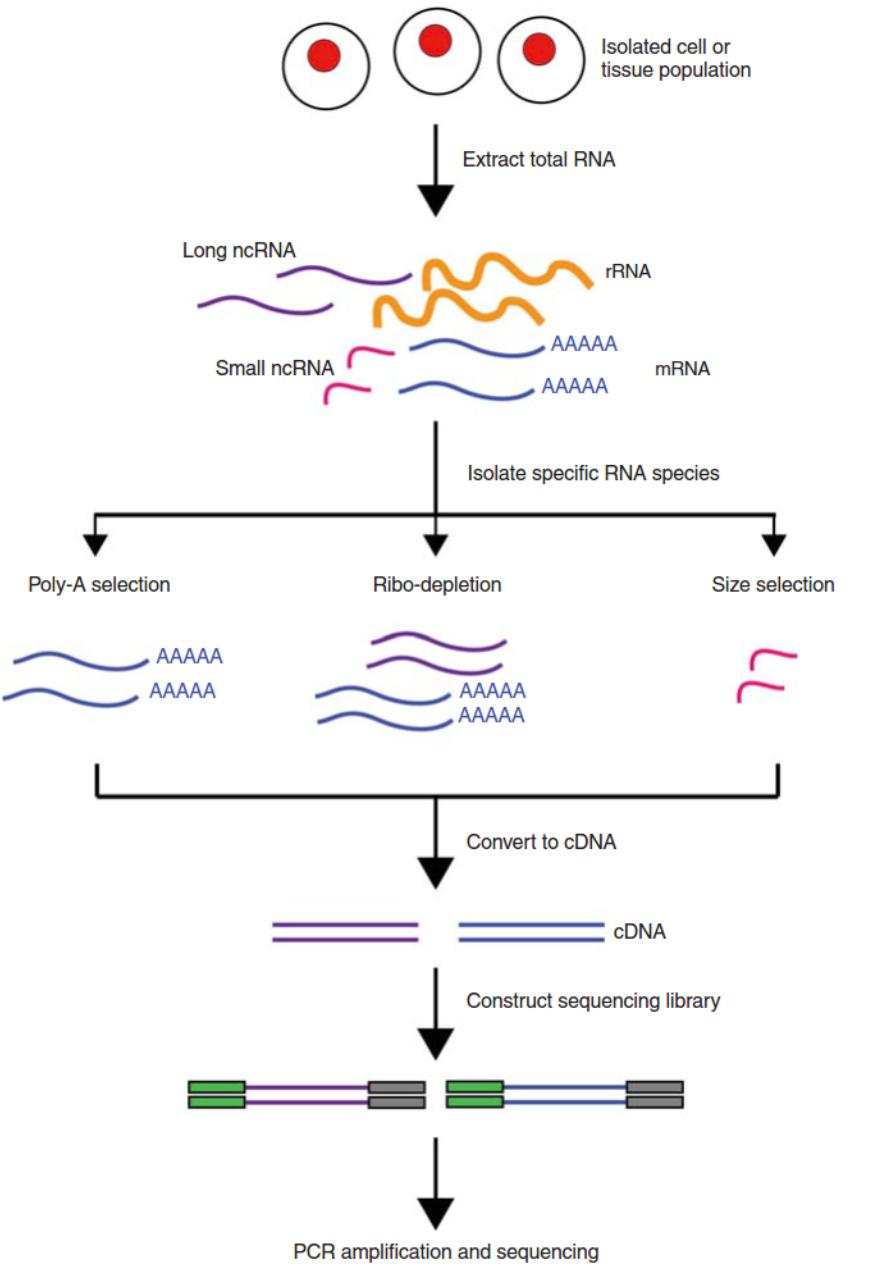
Bioanalyzer trace before and after ribo-depletion

# RNA-Seq library prep procedure

1. RNA-sample QC, quantification, and normalization
2. Removal of ribosomal RNA sequences:  
via positive or negative selection: Poly-A enrichment or ribo-depletion
3. Fragment RNA:  
heating in Mg++ containing buffer – chemical fragmentation has little bias
4. First-strand synthesis:  
random hexamer primed reverse transcription
5. RNase-H digestion:
  - creates nicks in RNA strand; the nicks prime 2nd-strand synthesis
  - dUTP incorporated into 2<sup>nd</sup> strand only
6. A-tailing and adapter ligation exactly as for DNA-Seq libraries
7. PCR amplification of only the first strand to achieve strand-specific libraries - archeal polymerases will not use dUTP containing DNA as template

# Illumina sequencing workflow

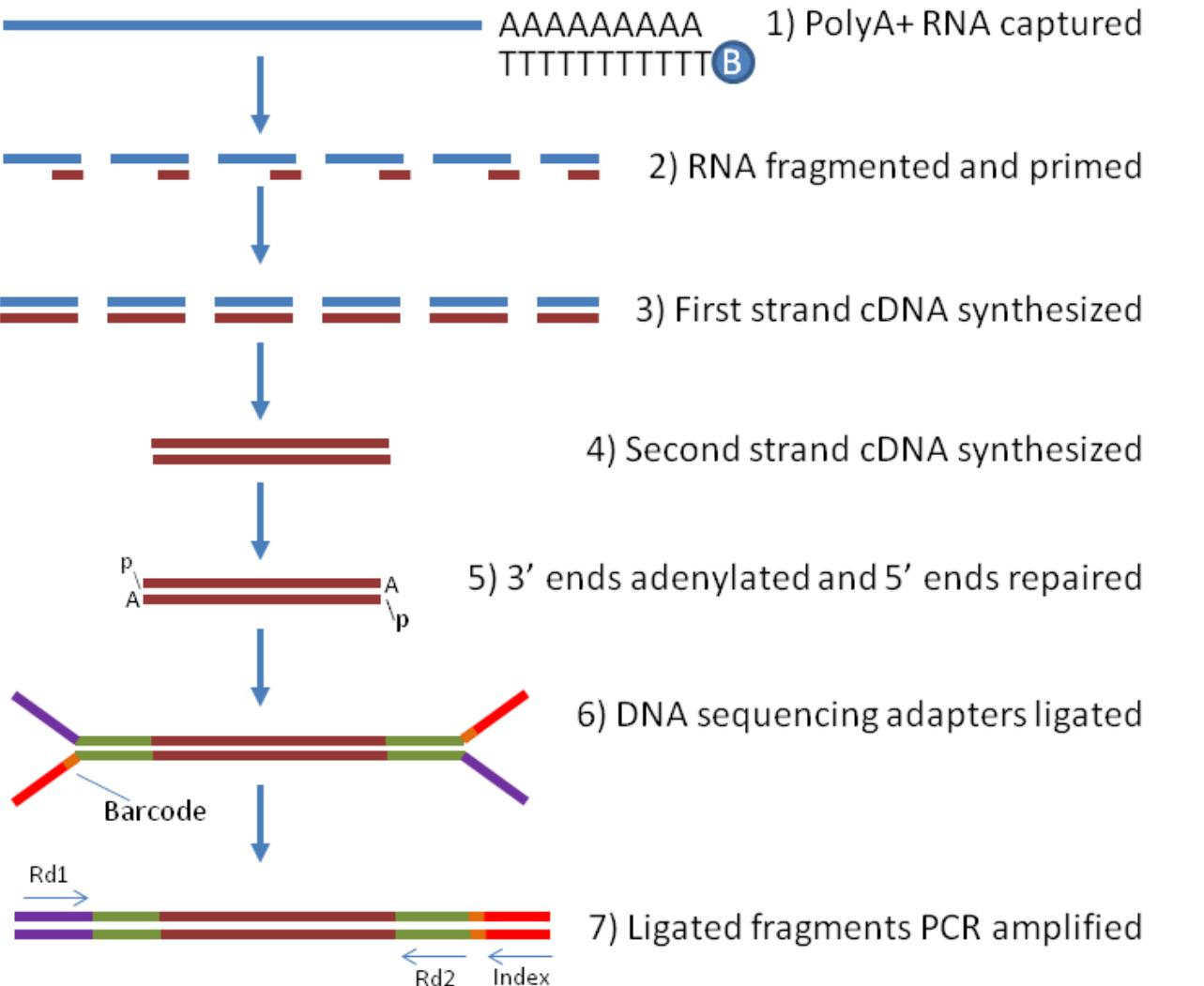
- Library Construction
- Cluster Formation
- Sequencing
- Data Analysis



RNA-seq?

Sorry – Illumina and  
PacBio are only  
sequencing DNA.

# Conventional RNA-Seq library preparation w. Poly-A capture



# What will go wrong ?

- cluster identification
- bubbles
- synthesis errors:

ClusterCluster  
ClustsrCluster  
ClusterCluster  
ClusterCluster  
CllsterCluster

# What will go wrong ?

- synthesis errors:

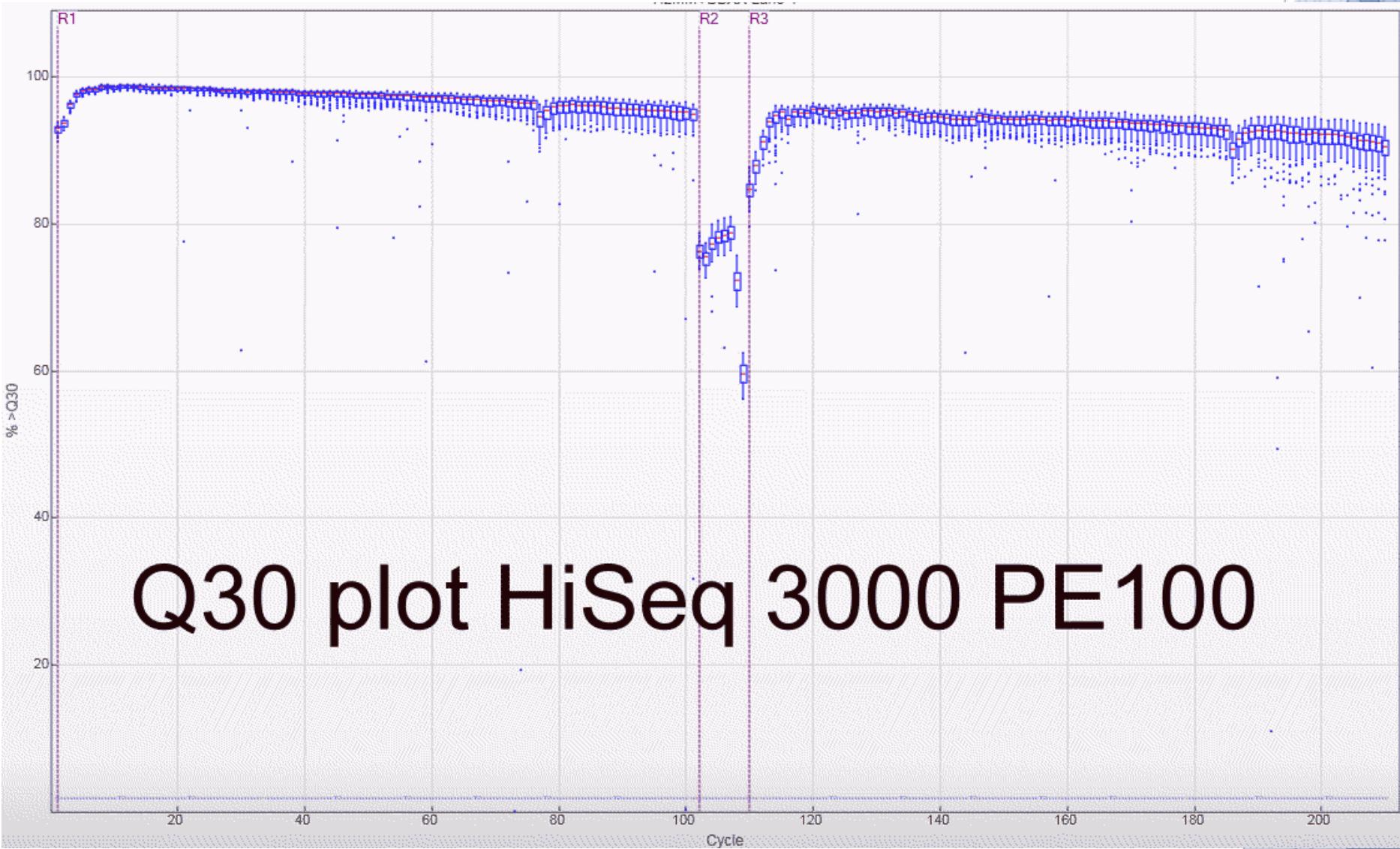
ClusterCluster  
ClustsrCluster  
ClusterCluster  
ClusterCluster  
CllsterCluster

CllsterClusterC  
ClusterCluster  
ClusterCluster  
CllusterCluste  
ClusterCluster

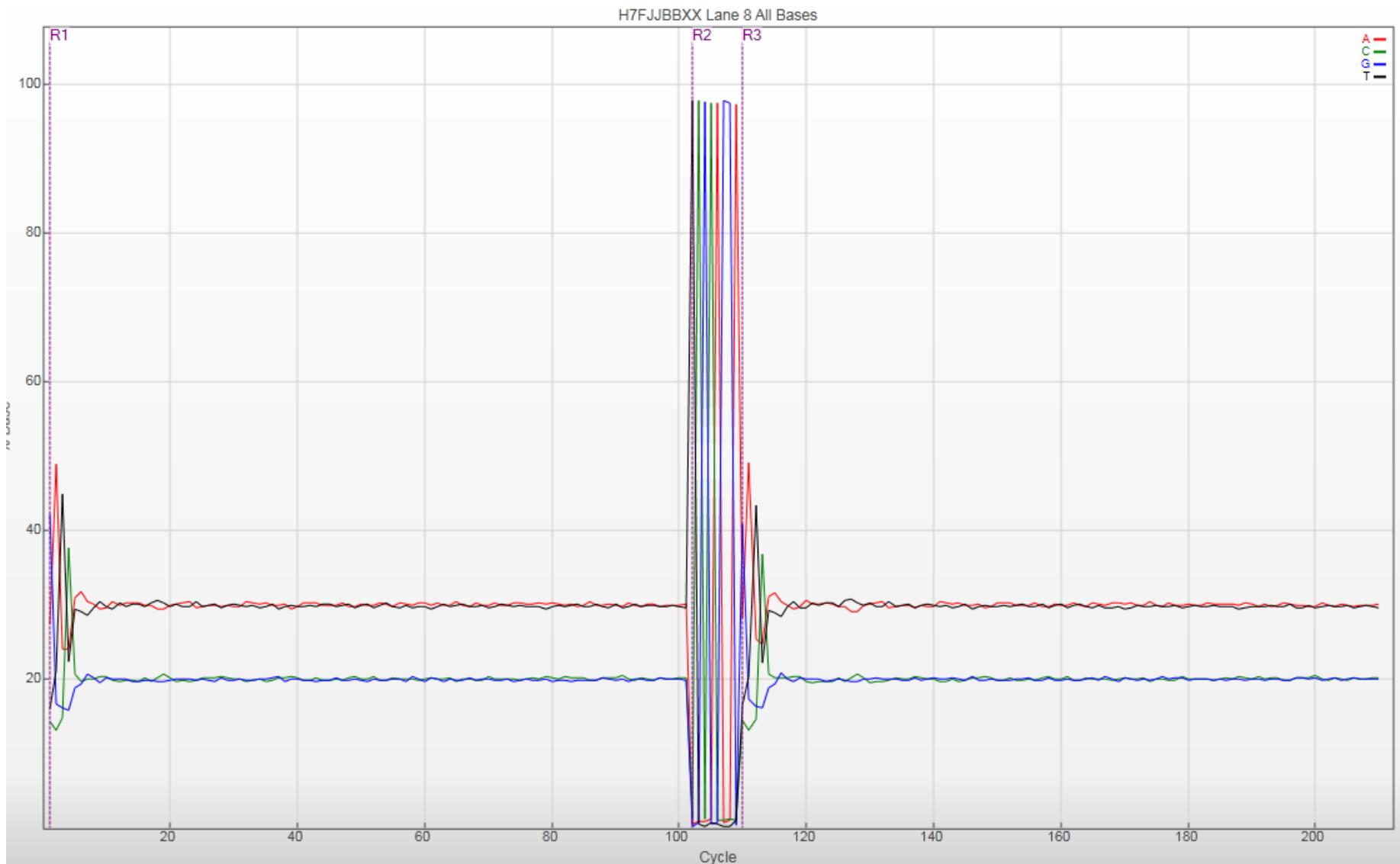
Phasing & Pre-Phasing  
problems

# The first lines of your data

# Illumina SAV viewer

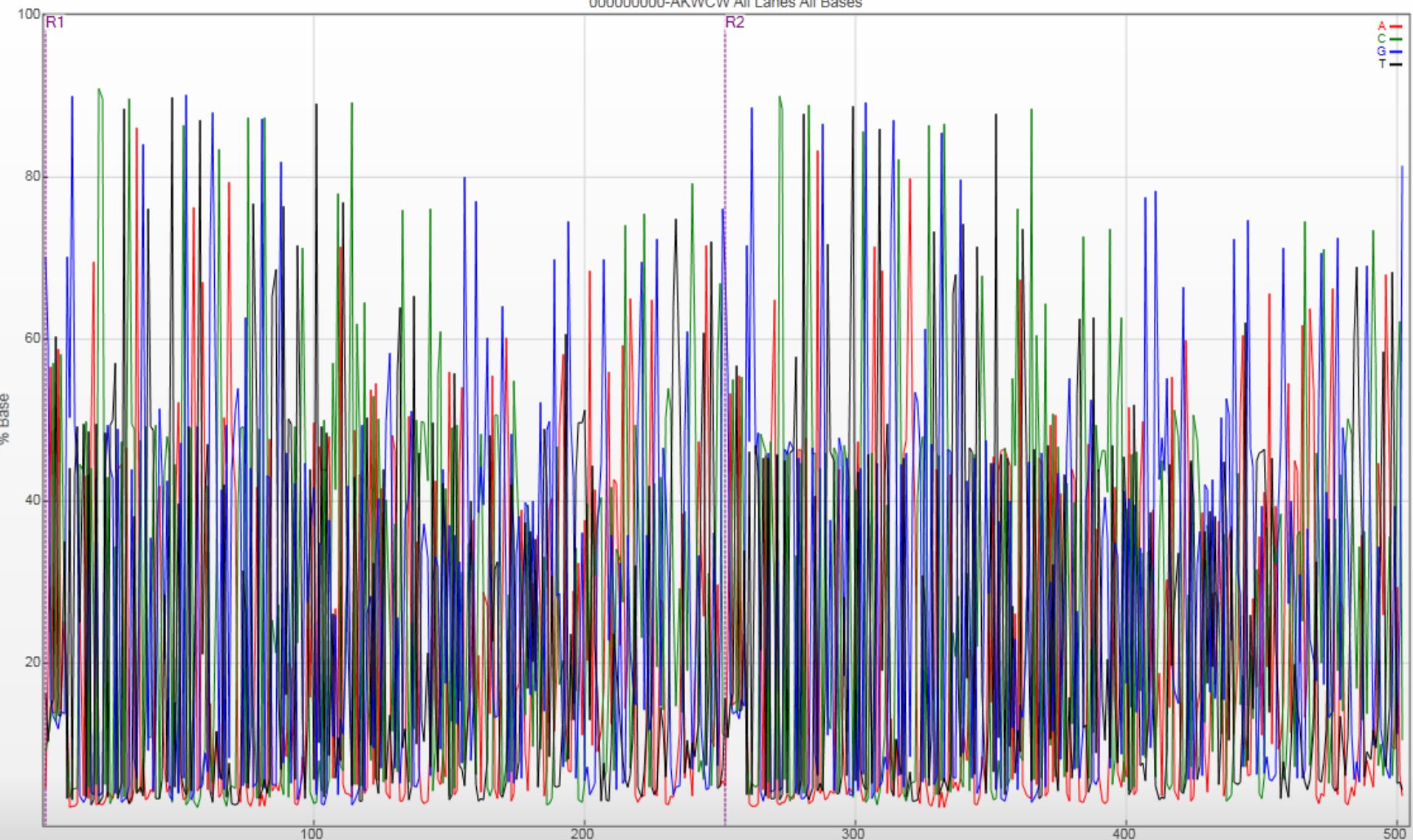


# base composition



# amplicon

000000000-AKWCW All Lanes All Bases

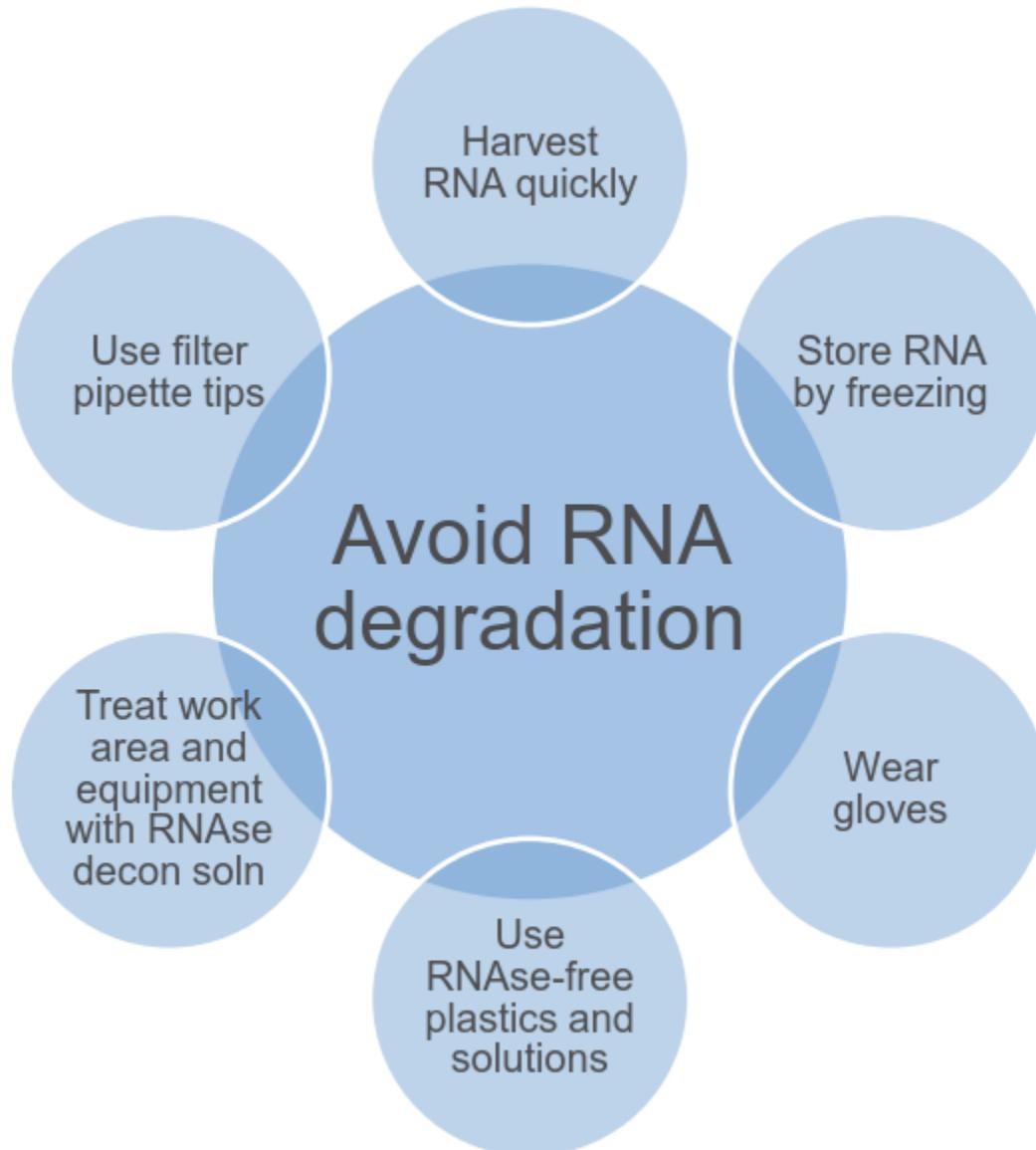


# RNA is not that fragile

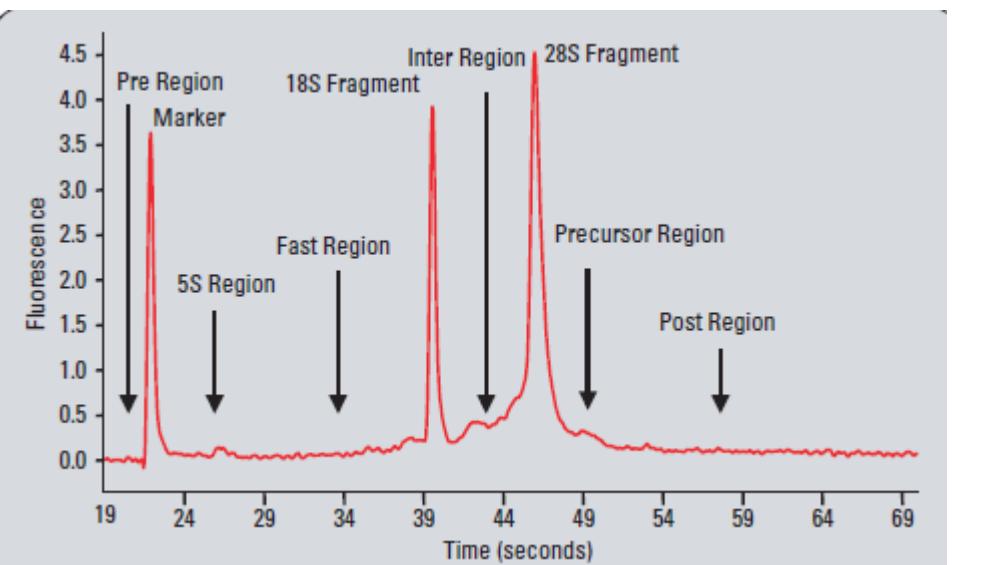


Actually: Avoid DEPC-treated reagents -- remnants can inhibit enzymes

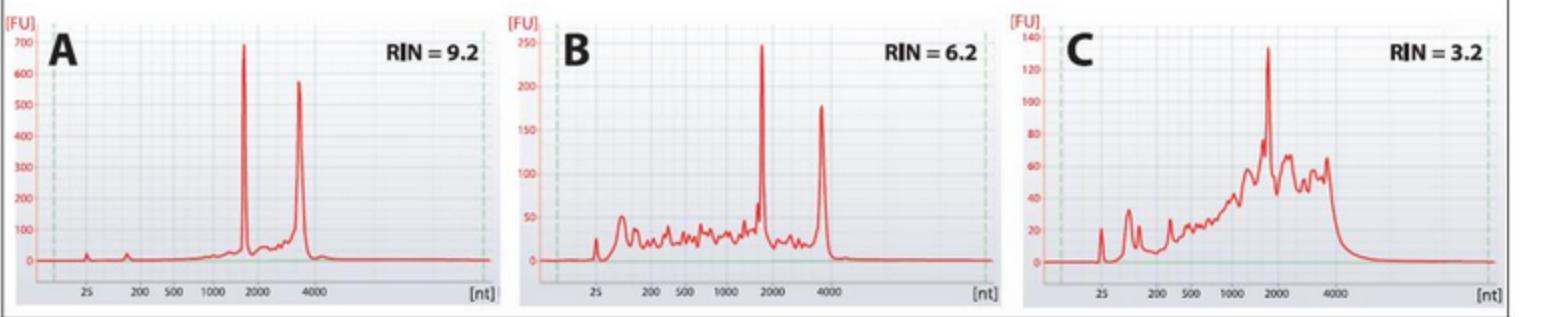
# RNA Handling Best Practices



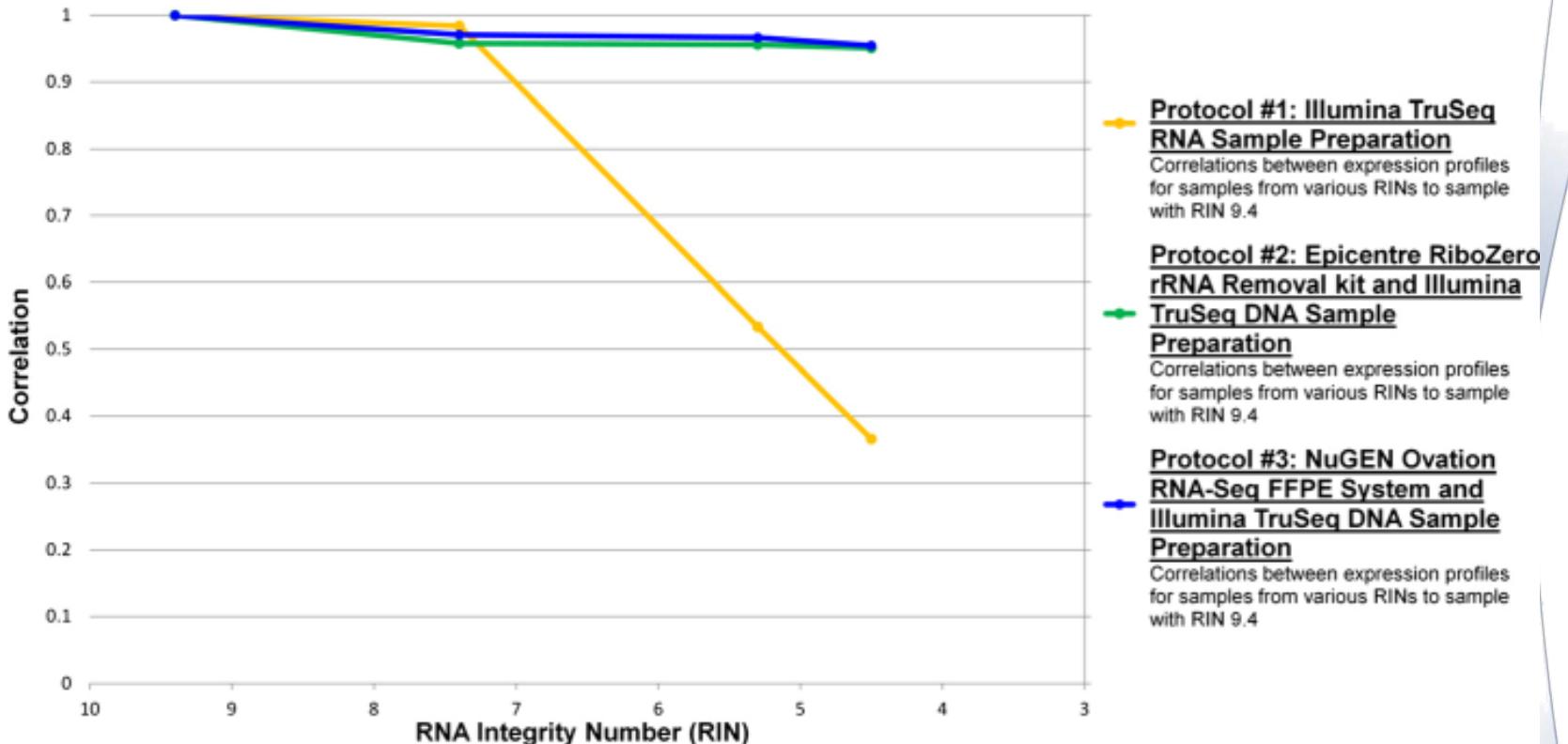
- 18S (2500b) , 28S (4000b)



**Figure 2.1** Example Agilent Bioanalyzer Electropherograms from three different total RNAs of varying integrity. Panel [A] represents a highly intact total RNA (RIN = 9.2), panel [B] represents a moderately intact total RNA (RIN = 6.2), and panel [C] represents a degraded total RNA sample (RIN = 3.2).



# RNA integrity <> reproducibility



Chen et al. 2014

# Quantitation & QC methods

## ➤ Intercalating dye methods (PicoGreen, Qubit, etc.):

Specific to dsDNA, accurate at low levels of DNA

Great for pooling of indexed libraries to be sequenced in one lane

Requires standard curve generation, many accurate pipetting steps

## ➤ Bioanalyzer:

Quantitation is good for rough estimate

Invaluable for library QC

High-sensitivity DNA chip allows quantitation of low DNA levels

## ➤ qPCR

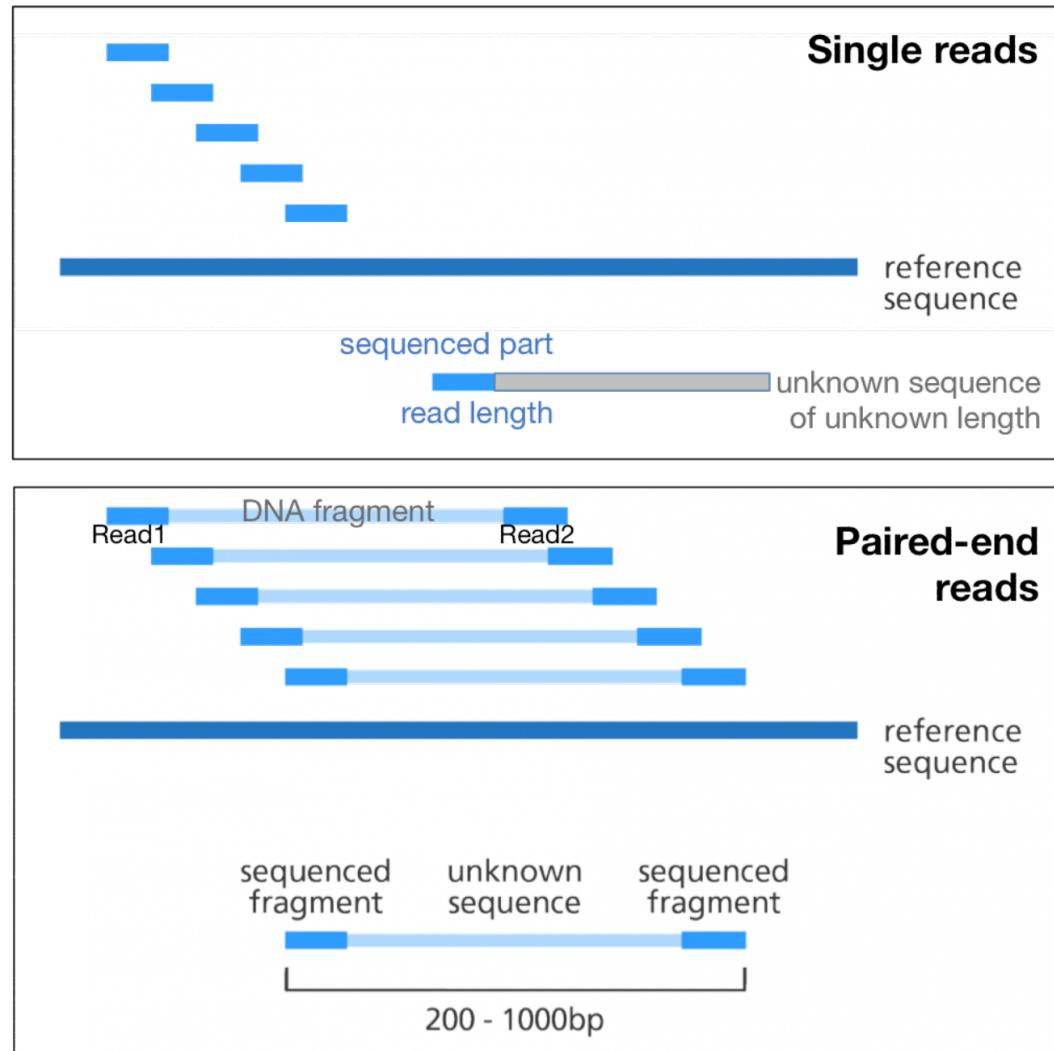
Most accurate quantitation method

More labor-intensive

Must be compared to a control

# Recommended RNA input

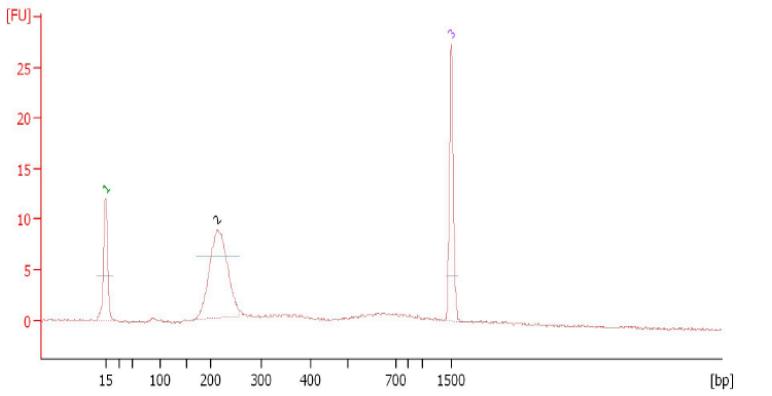
| Library prep kit                                      | Starting material                |
|---|----------------------------------|
| mRNA (TruSeq)   | 100 ng – 4 µg total RNA          |
| Directional mRNA (TruSeq)                             | 1 – 5 µg total RNA or 50 ng mRNA |
| Apollo324 library robot<br>(strand specific)          | 100 ng mRNA                      |
| Small RNA (TruSeq)                                    | 100 ng -1 µg total RNA           |
| Ribo depletion (Epicentre)                            | 500 ng – 5 µg total RNA          |
| SMARTer™ Ultra Low RNA<br>(Clontech)                  | 100 pg – 10 ng                   |
| Ovation RNA seq V2,<br>Single Cell RNA seq<br>(NuGen) | 10 ng – 100 ng                   |



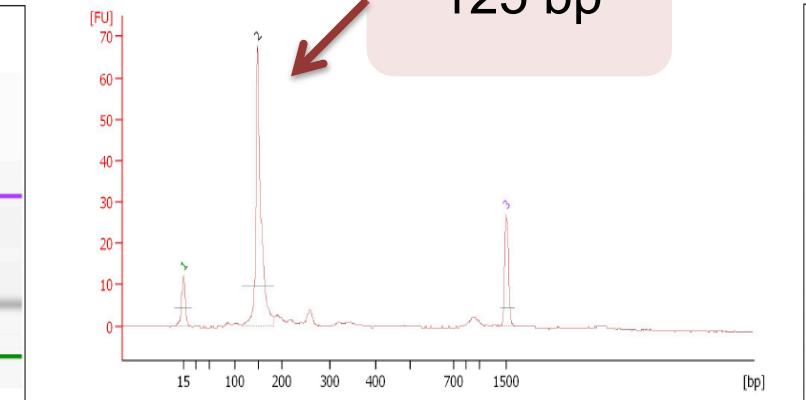
Single reads are the cheaper.  
Paired-end (PE) reads are helpful for:

- **alignment** along repetitive regions
- chromosomal **rearrangements** and gene fusion detection
- *de novo* genome and transcriptome **assembly**
- precise information about the size of the original fragment (**insert size**)
- PCR duplicate identification

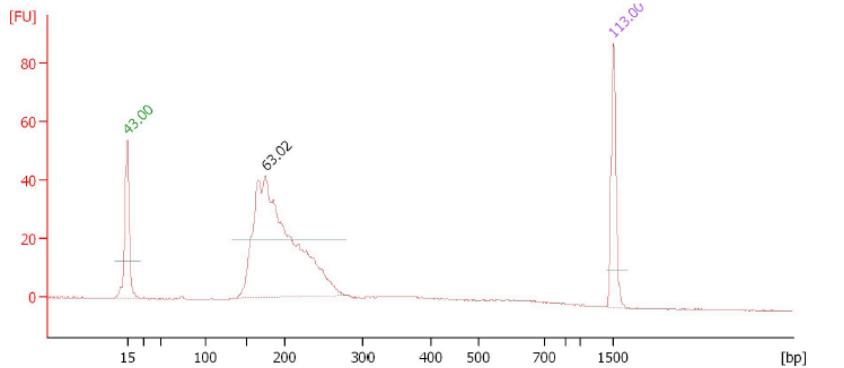
# Library QC by Bioanalyzer



# Beautiful

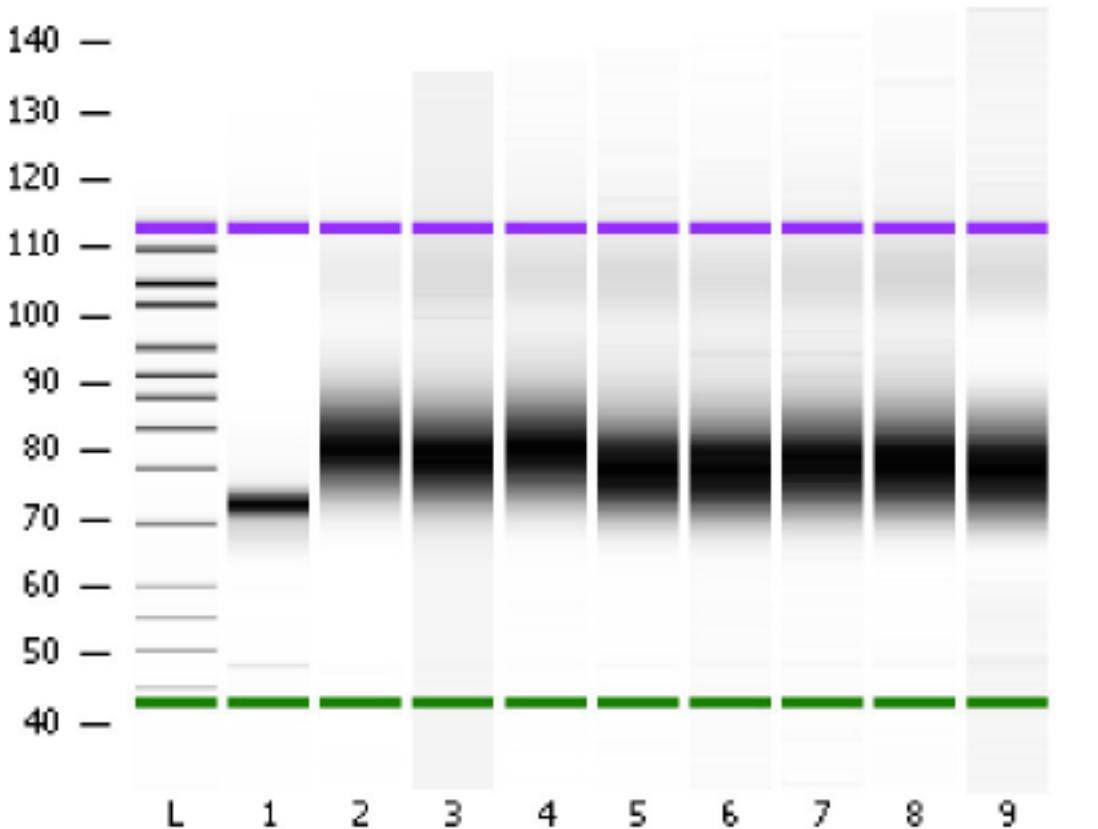


# 00% Adapters

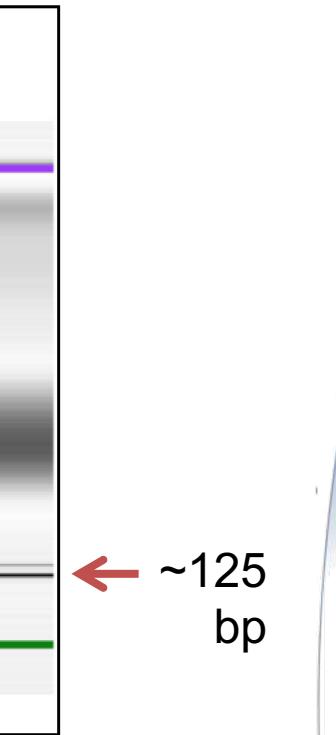


# Beautiful

# Library QC



Examples for successful libraries



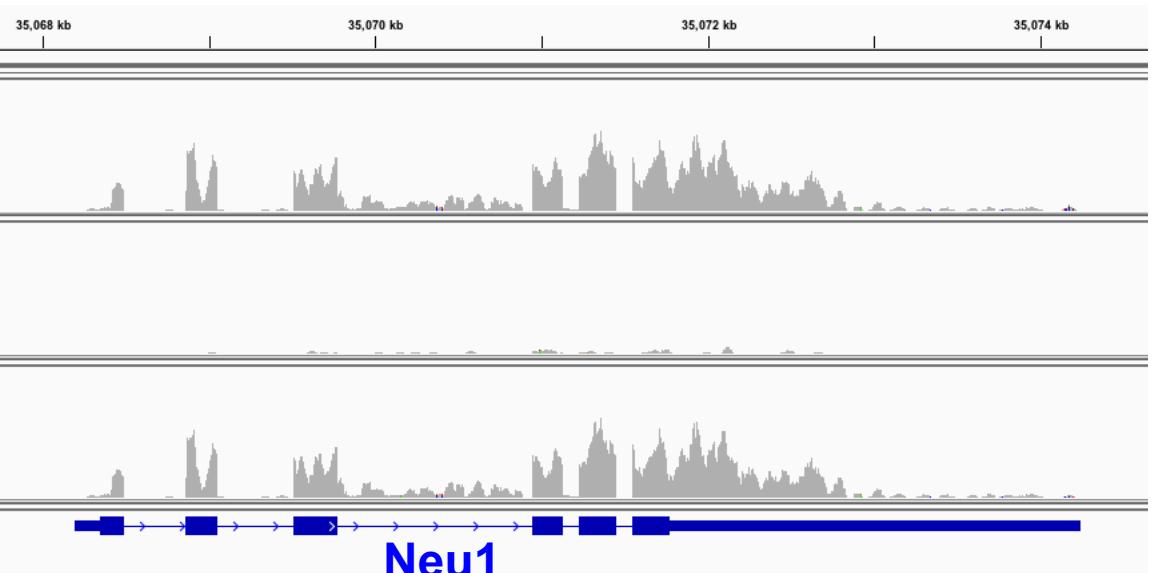
Adapter  
contamination  
at ~125 bp

# Considerations in choosing an RNA-Seq method

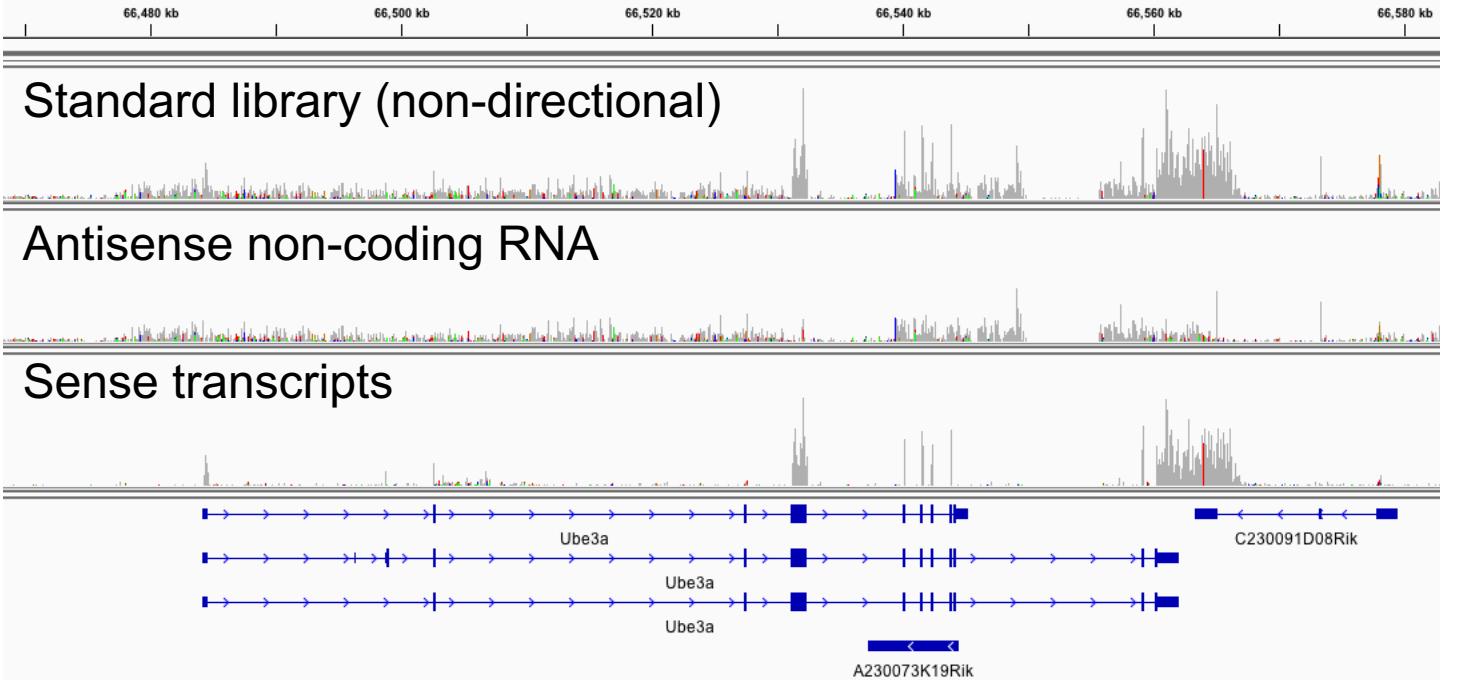
- Transcript type:
  - mRNA, extent of degradation
  - small/micro RNA
- Strandedness:
  - un-directional ds cDNA library
  - directional library
- Input RNA amount:
  - 0.1-4ug original total RNA
  - linear amplification from 0.5-10ng RNA
- Complexity:
  - original abundance
  - cDNA normalization for uniformity
- Boundary of transcripts:
  - identify 5' and/or 3' ends
  - poly-adenylation sites
  - Degradation, cleavage sites

# strand-specific information

Standard library  
(non-directional)



# Strand-specific RNA-seq

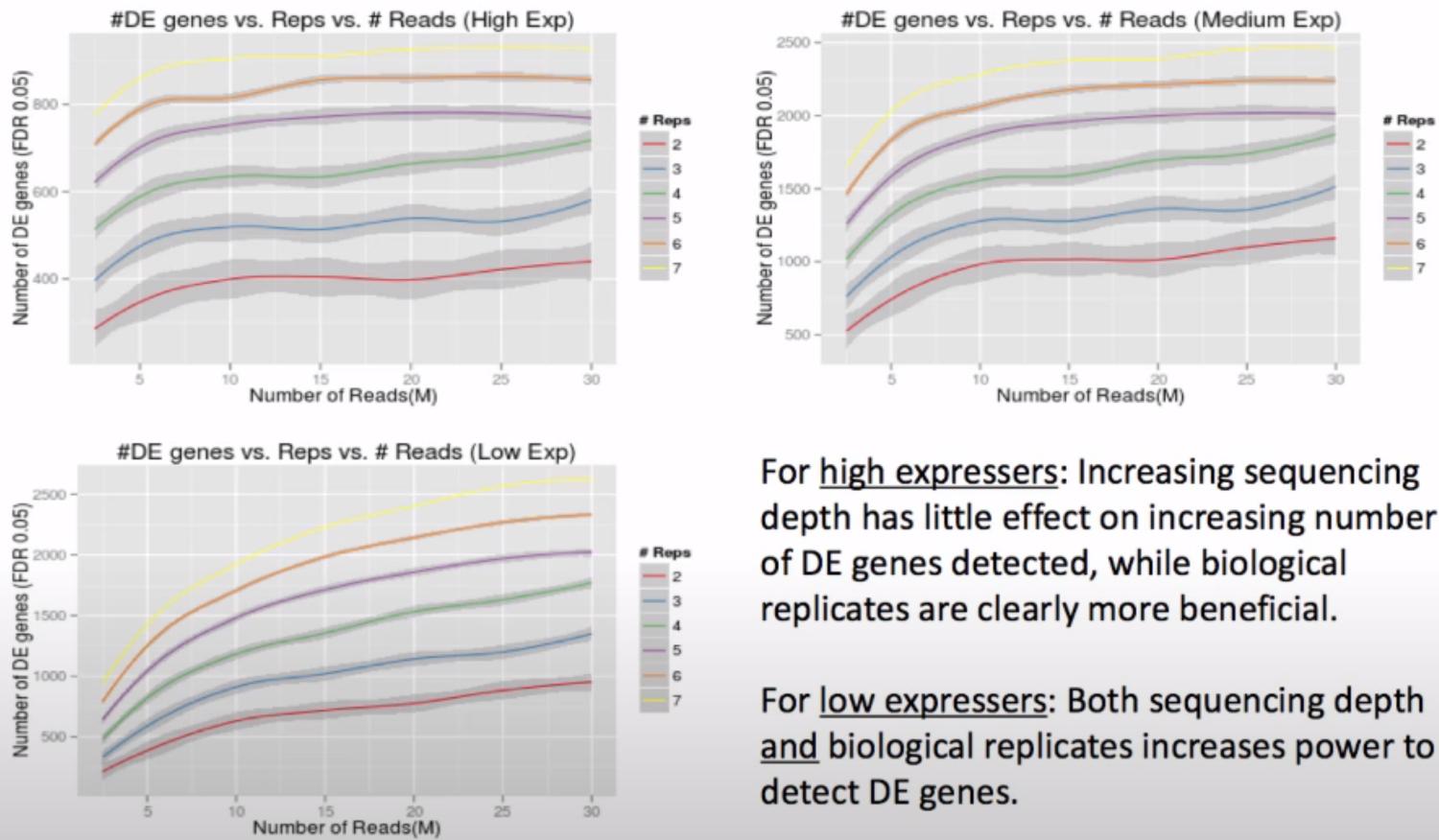


- Informative for non-coding RNAs and antisense transcripts
- Essential when NOT using polyA selection (mRNA)
- No disadvantage to preserving strand specificity

# RNA-seq for DGE

- Differential Gene Expression (DGE)
  - 50 bp single end reads
  - 30 million reads per sample (eukaryotes)
    - 10 mill. reads > 80% of annotated genes
    - 30 mill. . reads > 90% of annotated genes
  - 10 million reads per sample (bacteria)

## Experimental Design



For high expressers: Increasing sequencing depth has little effect on increasing number of DE genes detected, while biological replicates are clearly more beneficial.

For low expressers: Both sequencing depth and biological replicates increases power to detect DE genes.

Liu et al. (2014) RNA-Seq differential expression studies: more sequence or more replication?, Bioinformatics, 30(3):1-4

Image credit: Kevin Knudtson

# RNA-seq reproducibility

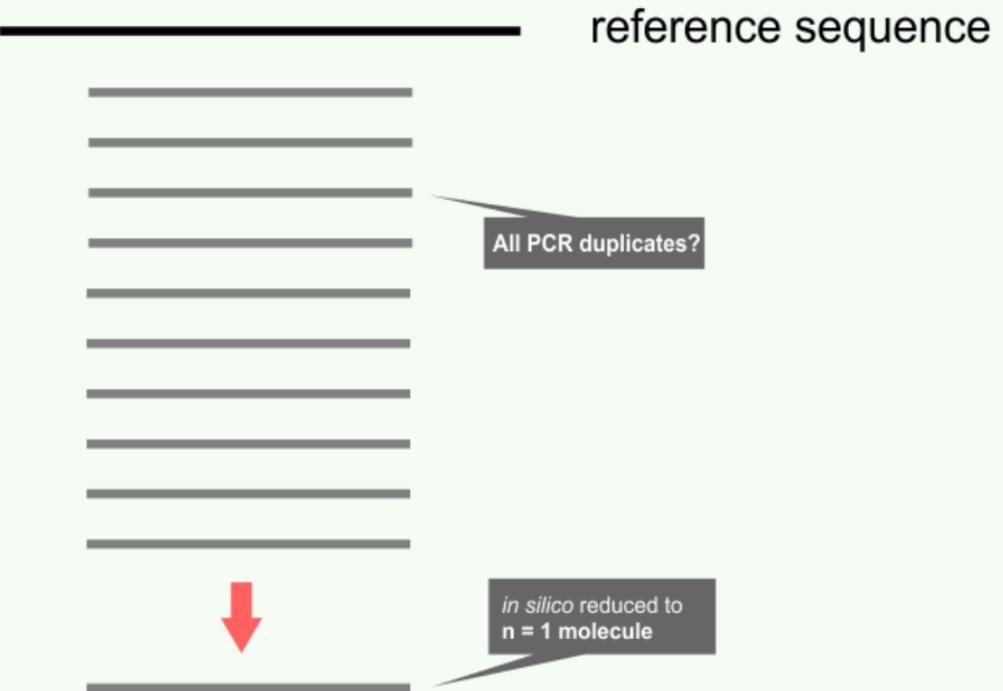
- Two big studies multi-center studies (2014)
- High reproducibility of data given:
  - same library prep kits, same protocols
  - same RNA-samples
  - RNA isolation protocols have to be identical
  - robotic library preps?

# UMIs – Unique Molecular Identifiers

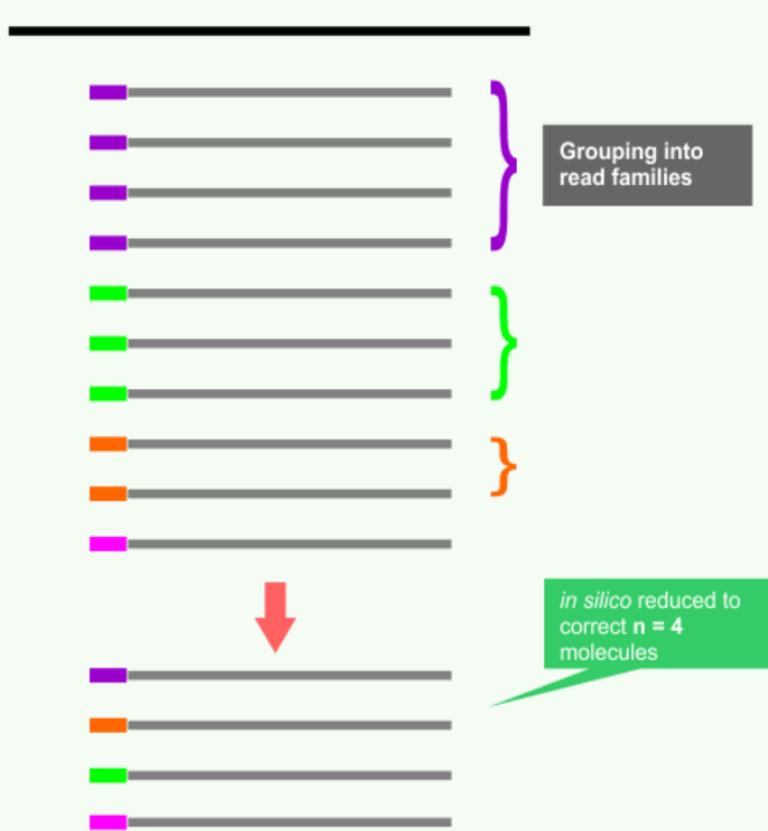
## Molecular indexing for precision counts

UMI application in **quantitative studies** (e.g. RNA-seq, scRNA-seq, miRNA-Seq, ChIP-seq).

PCR duplicate removal without UMIs

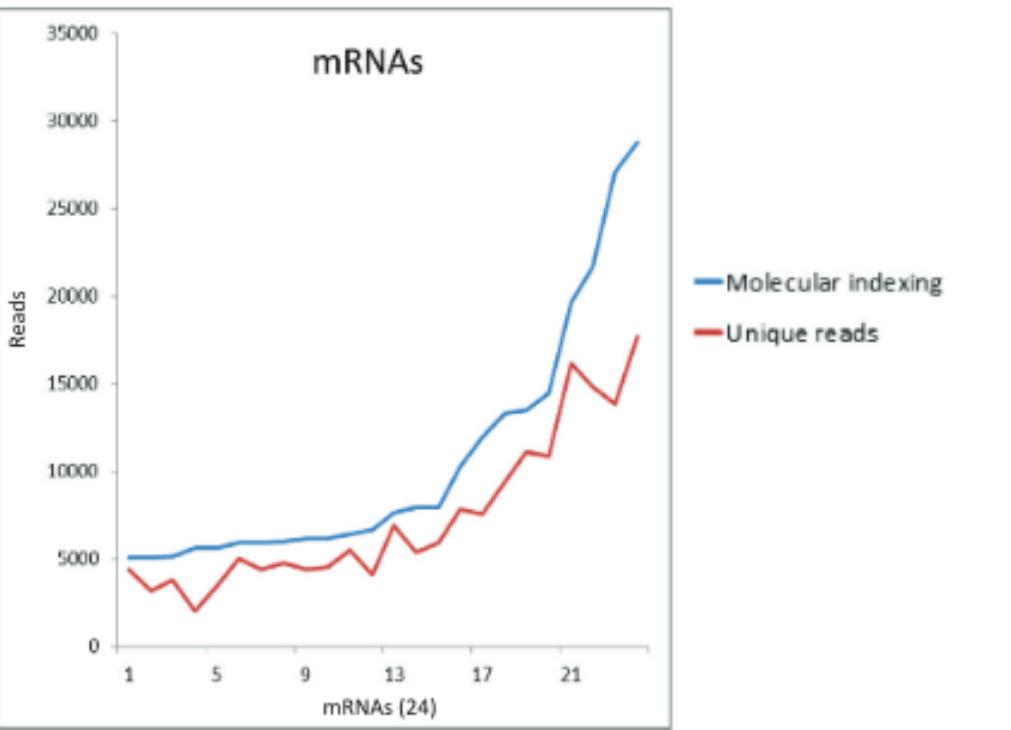


PCR duplicate removal with UMIs



# Molecular indexing – for precision counts

B

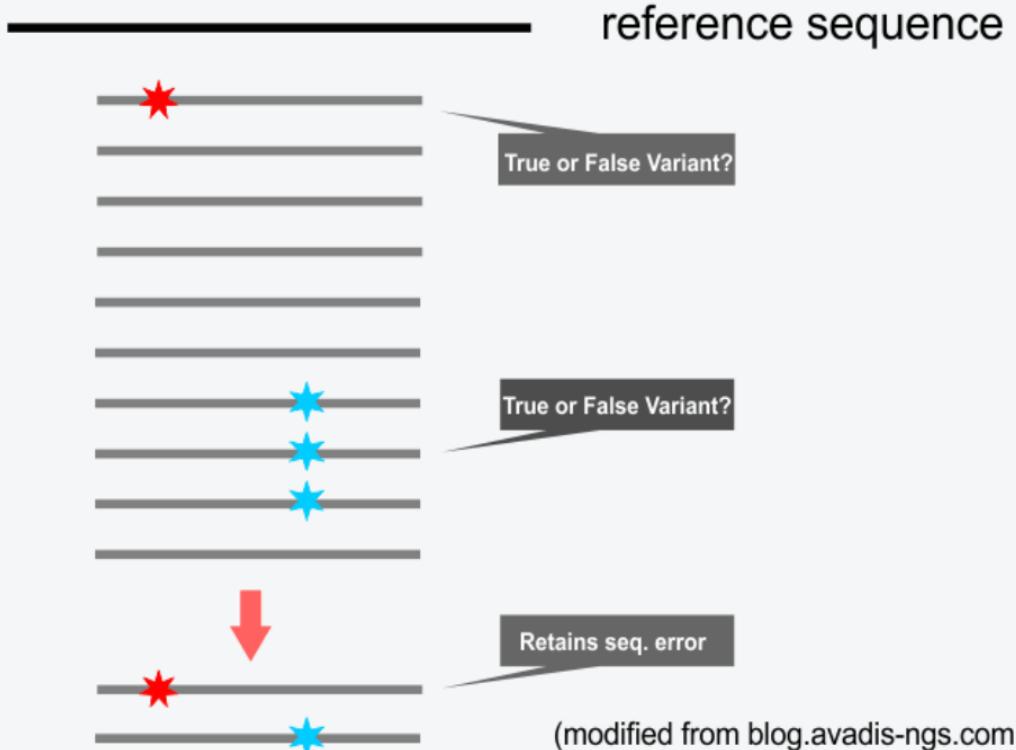


# UMIs – Unique Molecular Identifiers

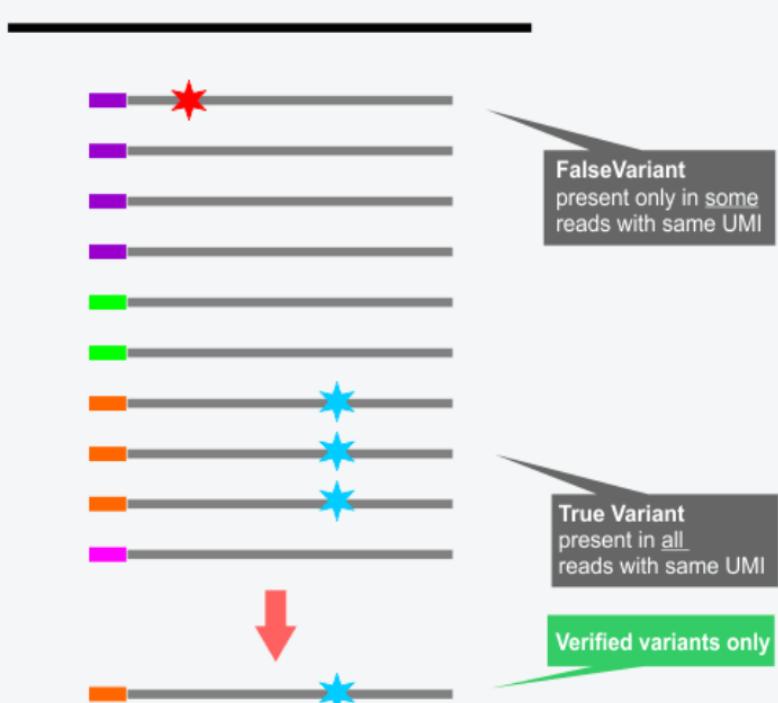
## Molecular indexing for low abundance variants

UMI application in deep sequencing **genomic variation** studies (e.g. WGS, exome capture, cfDNA)

Variant calling without UMIs



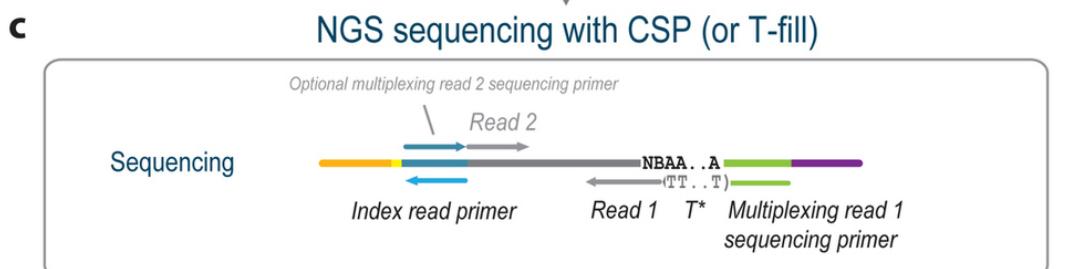
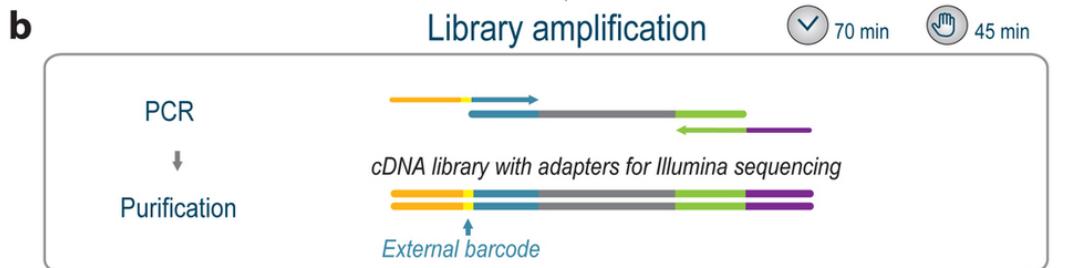
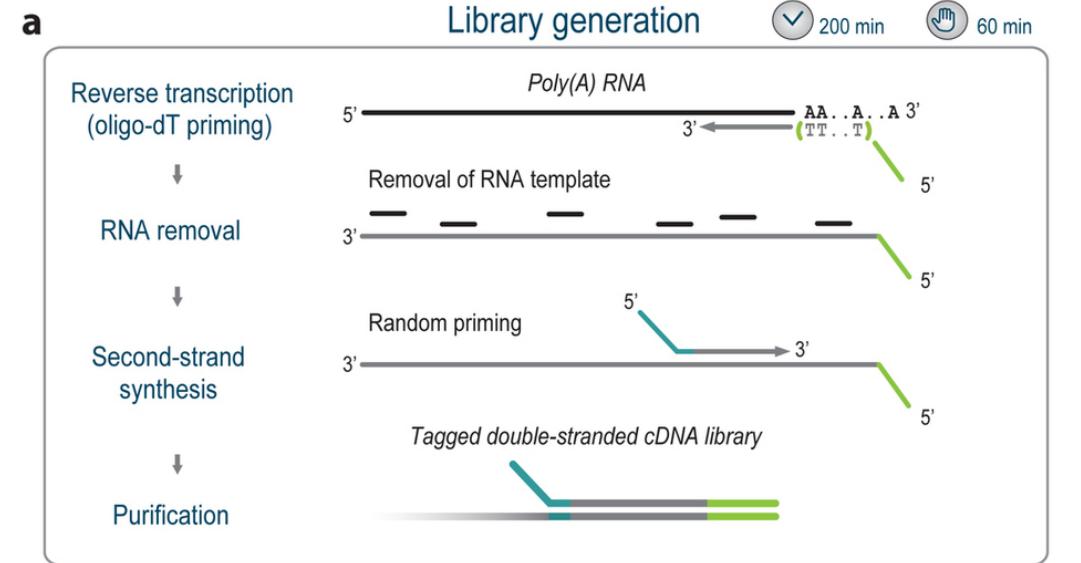
Variant calling with UMIs



# 3'-Tag-Seq

- In contrast to full length RNA-seq
- Sequencing 1/10 for the average transcript
- Less dependent on RNA integrity
- Microarray-like data
- Options:
  - **BRAD-Seq : 3' Digital Gene Expression**
  - **Lexogen Quant-Seq**

# Lexogen Quant-Seq



- we include UMIs

# Other RNA-seq objectives

- Transcriptome assembly:
  - 300 bp paired end **plus**
  - 100 bp paired end
- Long non coding RNA studies:
  - 100 bp paired end
  - 60-100 million reads
- Splice variant studies:
  - 100 bp paired end
  - 60-100 million reads

# RNA-seq targeted sequencing:

- Capture-seq (Mercer et al. 2014)
  - Nimblegen and Illumina
- 
- Low quality DNA (FFPE)
  - Lower read numbers 10 million reads
  - Targeting lowly expressed genes.

# Typical RNA-seq drawbacks

- Very much averaged data:  
Data from mixed cell types & mixed cell cycle stages
- Hundreds of differentially expressed genes  
(which changes started the cascade?)

higher resolution desired

→ beyond steady-state RNA-seq

# mechanisms influencing the mRNA steady-state

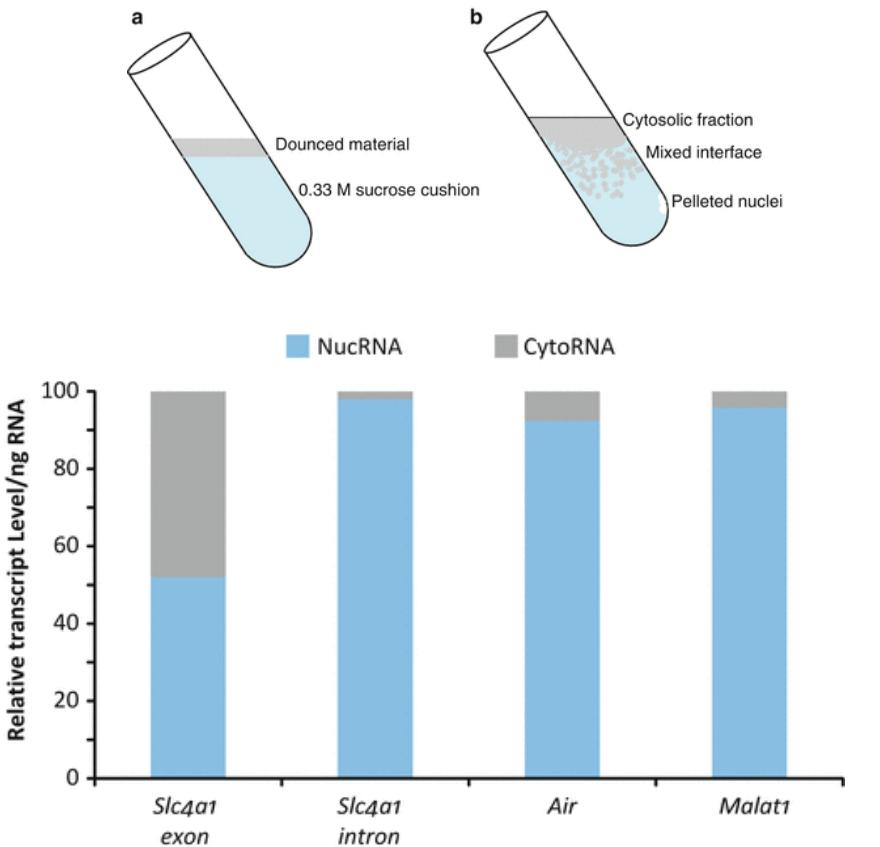
- Transcription rates
- Transport rates
- miRNAs and siRNAs influence both translation and degradation
- RNA modifications (e.g. methylated RNA bases, m<sup>6</sup>A, m<sup>5</sup>C, pseudouridine, ...)
- RNA degradation pathways
- (differential translation into proteins)

# beyond steady-state RNA-seq

- GRO-Seq; PRO-Seq; nuclear RNA-Seq:  
**what is currently transcribed**
- Ribosomal Profiling:  
**what is currently translated**
- Degradome Sequencing:  
**what is ... ?**

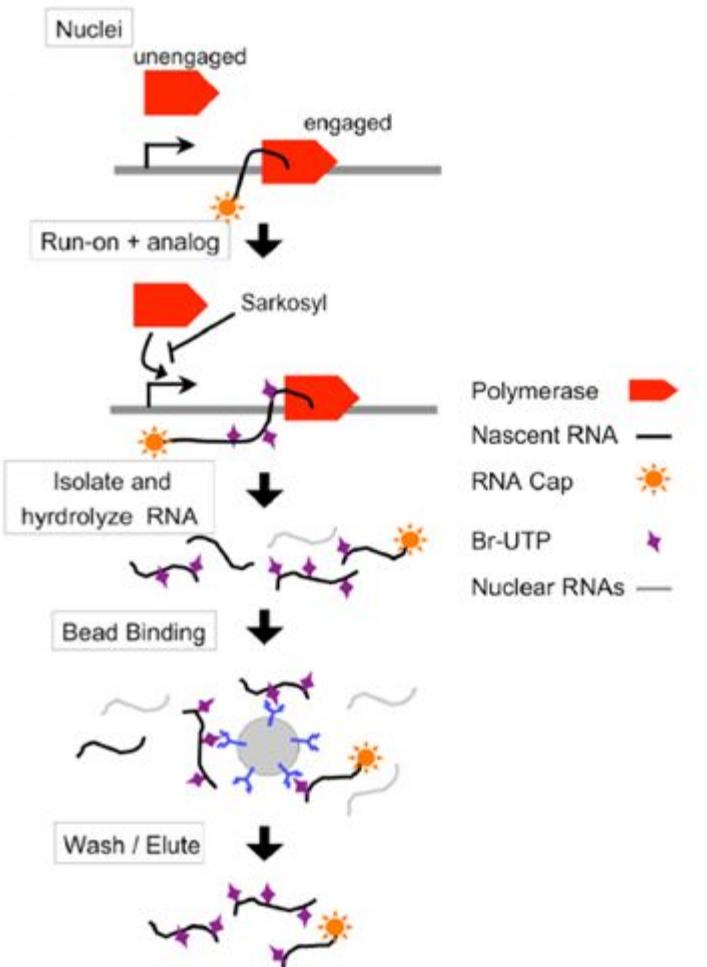
# nucRNA-seq

- Fractioning of nuclei and cytosol
- Studying active transcription



Dhaliwal et al. 2016

# GRO-Seq



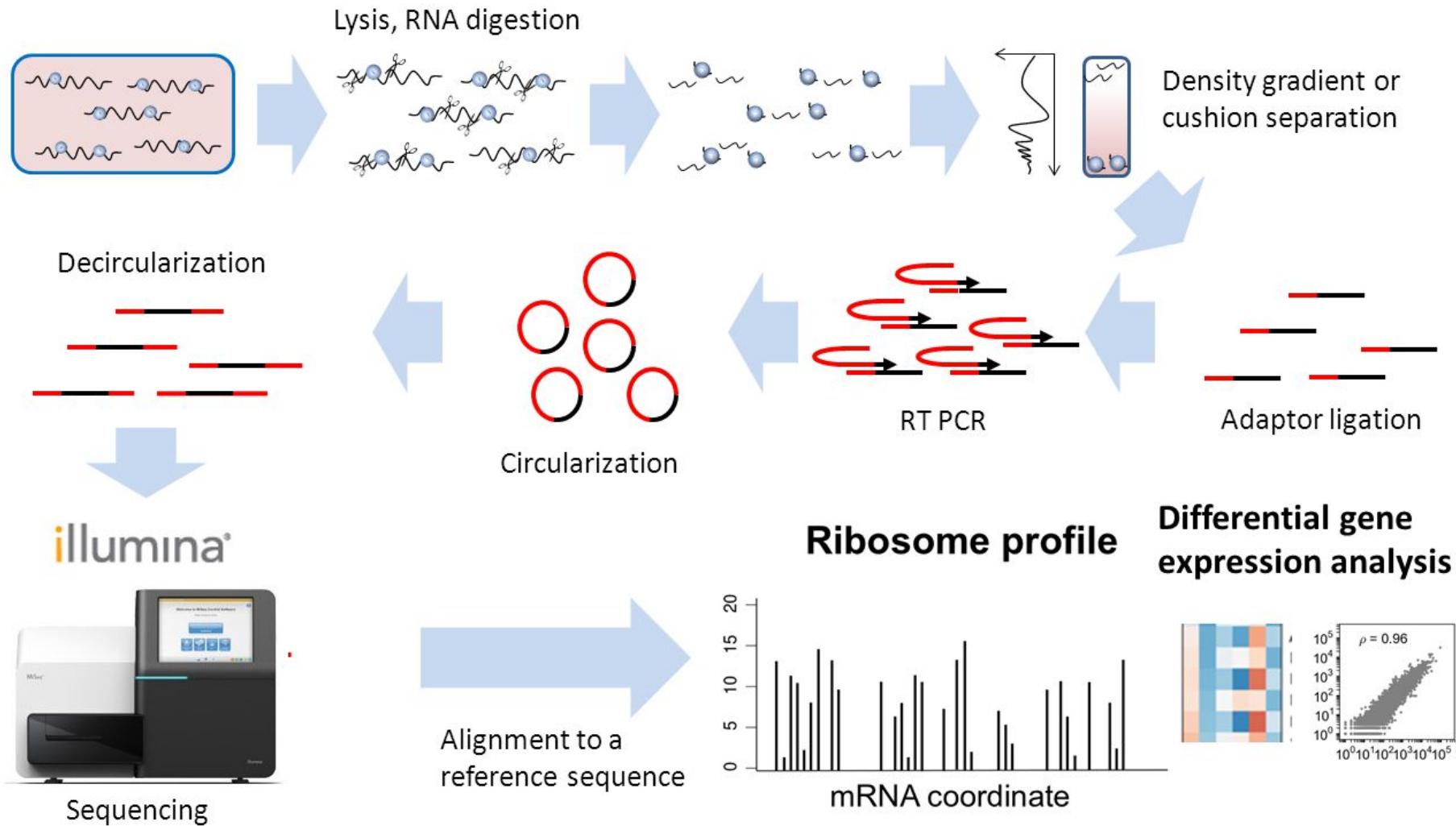
- Global Run-On – sequencing
- pulse-chase experiments (Br-UTP)
- uses isolated nuclei
- sarcosyl prevents binding of polymerase (only transcription in progress will be seq.)
- measures active transcription rather than steady state
- Maps position and orientation
- Earliest changes identify primary targets
- Detection of novel transcripts including non-coding and enhancer RNAs

Core et al, *Science*, 2008

2008: GRO - without the seq

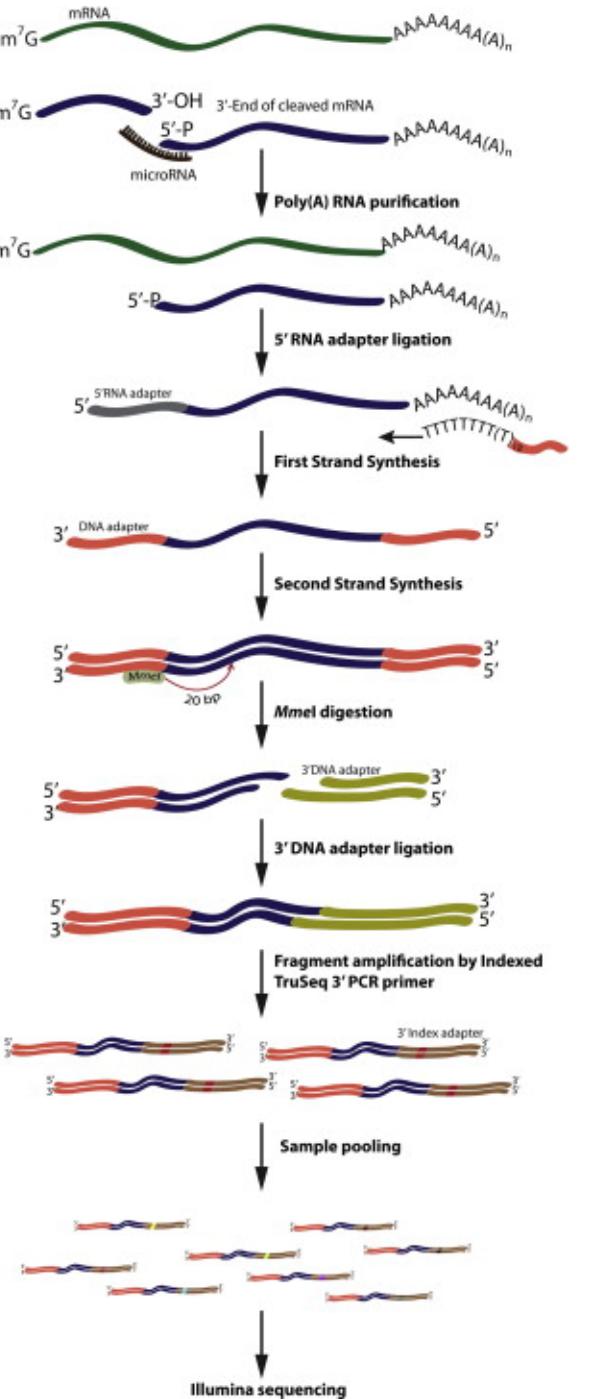
# Ribosomal profiling (ribo-seq)

*Ingolia et al (2009) Science 324: 218-23*



# Degradome Sequencing

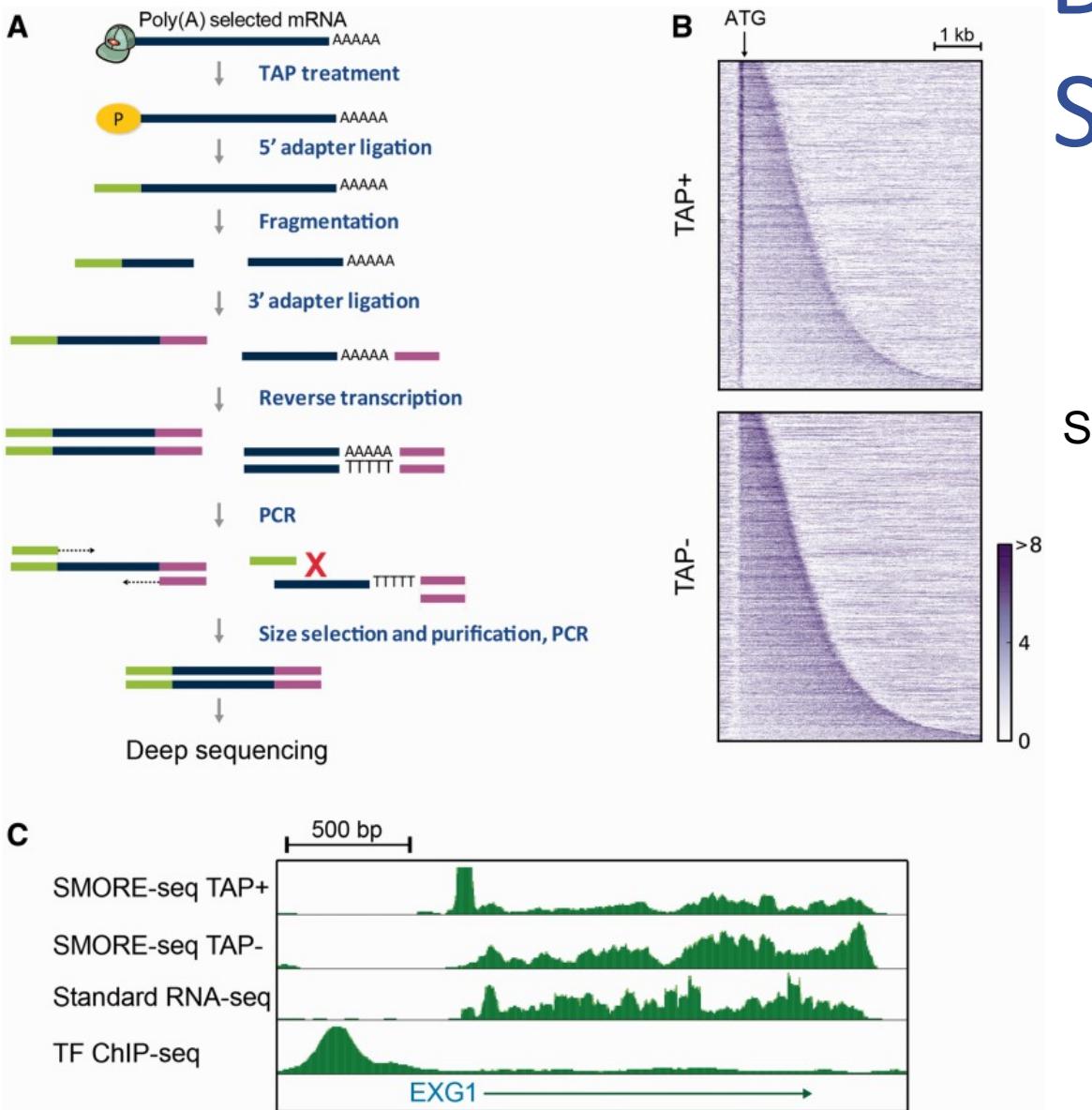
Day 1



## PARE-Seq (Parallel Analysis of RNA Ends)

Zhai et al . 2013

# Degradome Sequencing



Park et al . 2014

# RNA velocity

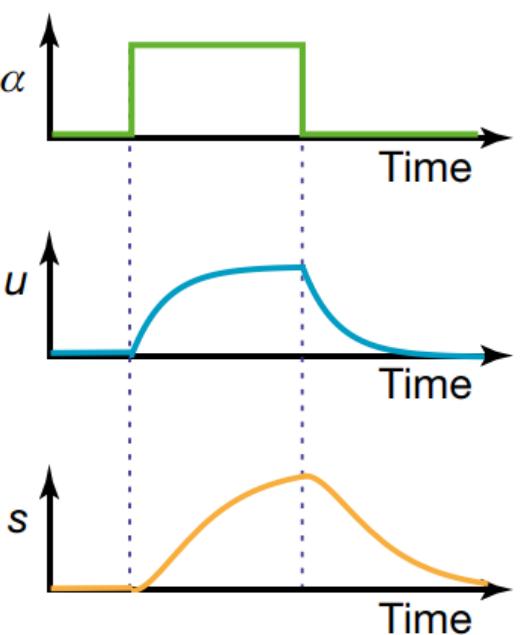
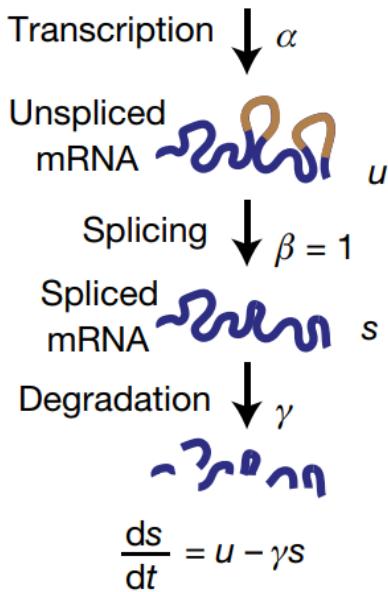
(2018)

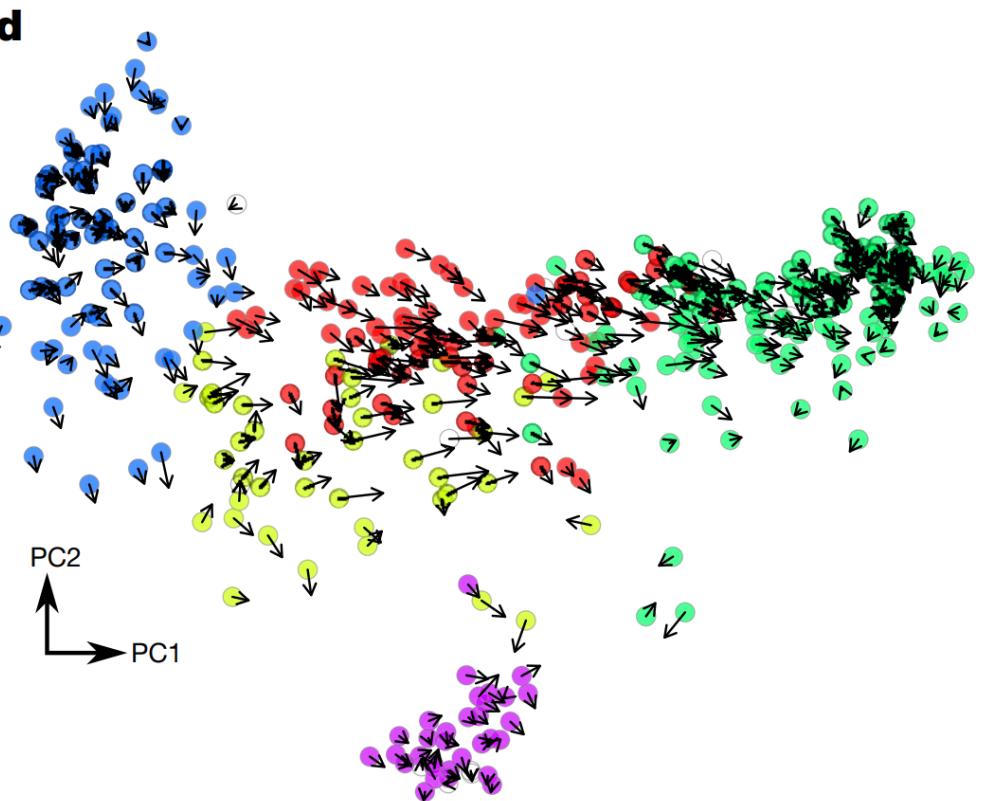
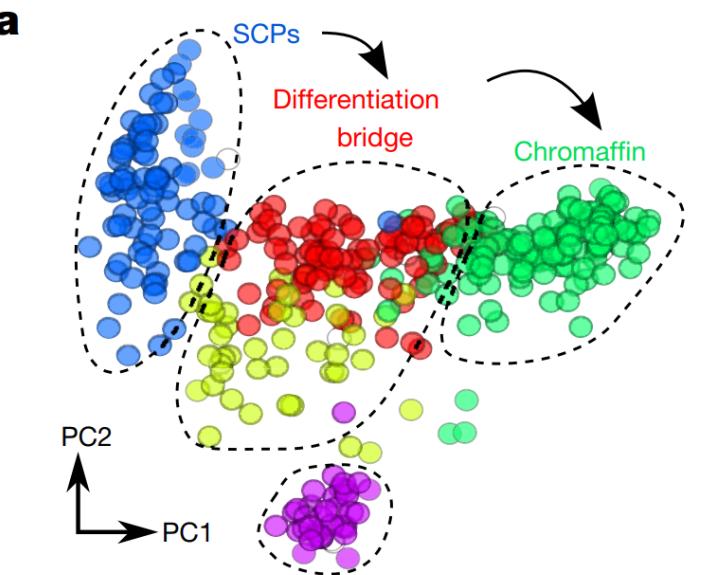
## LETTER

<https://doi.org/10.1038/s41586-018-0414-6>

### RNA velocity of single cells

Gioele La Manno<sup>1,2</sup>, Ruslan Soldatov<sup>3</sup>, Amit Zeisel<sup>1,2</sup>, Emelie Braun<sup>1,2</sup>, Hannah Hochgerner<sup>1,2</sup>, Viktor Petukhov<sup>3,4</sup>, Katja Lidschreiber<sup>5</sup>, Maria E. Kastriti<sup>6</sup>, Peter Lönnerberg<sup>1,2</sup>, Alessandro Furlan<sup>1</sup>, Jean Fan<sup>3</sup>, Lars E. Borm<sup>1,2</sup>, Zehua Liu<sup>3</sup>, David van Bruggen<sup>1</sup>, Jimin Guo<sup>3</sup>, Xiaoling He<sup>7</sup>, Roger Barker<sup>7</sup>, Erik Sundström<sup>8</sup>, Gonçalo Castelo-Branco<sup>1</sup>, Patrick Cramer<sup>5,9</sup>, Igor Adameyko<sup>6</sup>, Sten Linnarsson<sup>1,2\*</sup> & Peter V. Kharchenko<sup>3,10\*</sup>





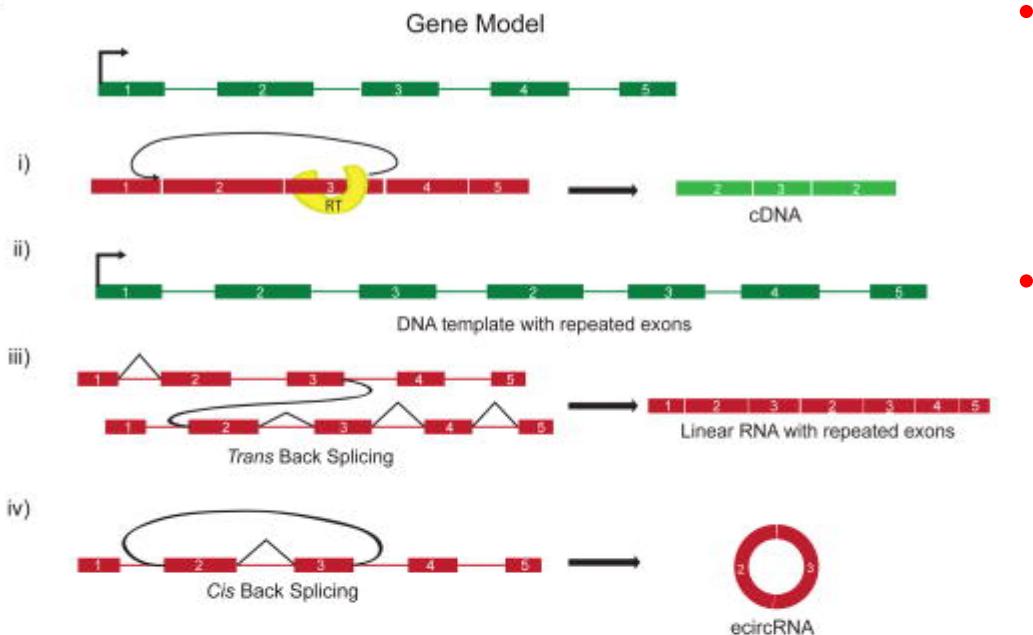
# Circular RNA (circRNA)

- Evolutionary conserved
- Eukaryotes
- Spliced (back-spliced)
- Some tissues contain more circRNA than mRNA
- Sequencing after exonuclease digestion (RNase R)
- Interpretation of ribo-depletion RNA-seq data ????

# Role of circRNAs ?

## Back-splicing and other mechanisms

A



B

Jeck and Sharpless, 2014

- miRNA sponge
- protein expression regulators:  
mRNA traps  
(blocking translation)
- Interactions with RNA binding proteins



PACIFIC  
BIOSCIENCES™

<http://pacificbiosciences.com>

## THIRD GENERATION DNA SEQUENCING

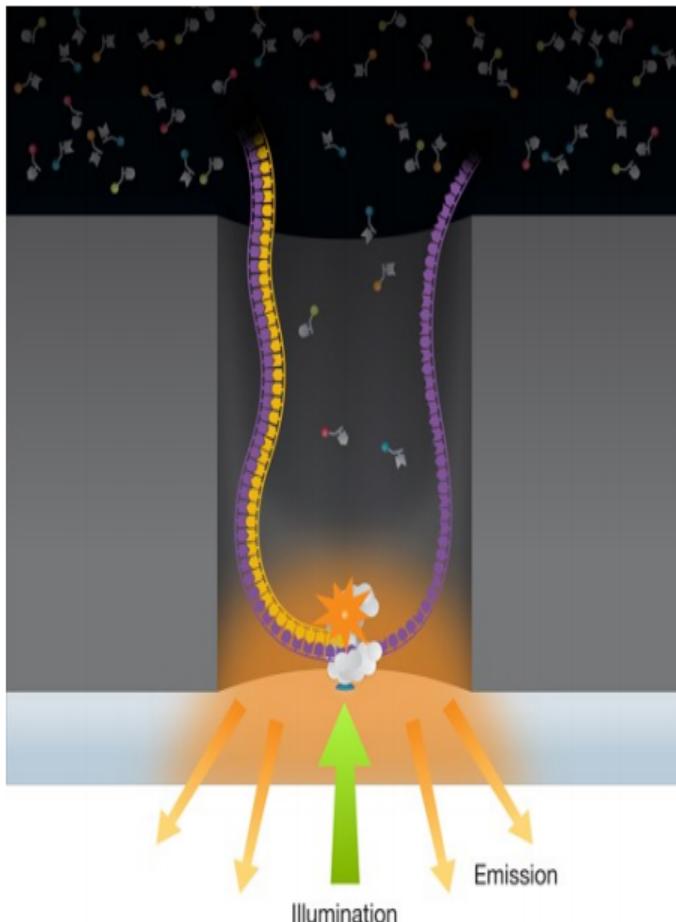


Single Molecule Real Time (SMRT™) sequencing  
Sequencing of single DNA molecule by single  
polymerase  
Very long reads: average reads over 8 kb, up to 30 kb  
High error rate (~13%).  
Complementary to short accurate reads of Illumina

## Third Generation Sequencing : Single Molecule Sequencing

Pacific Biosciences

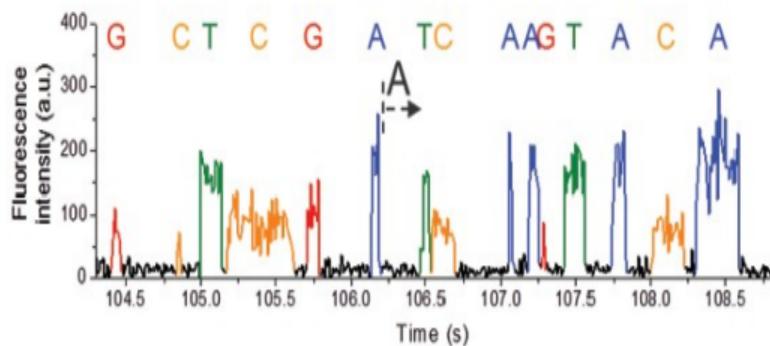
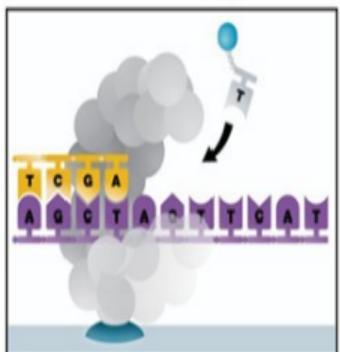
70 nm aperture  
“Zero Mode  
Waveguide”



4 nucleotides with different fluorescent dye simultaneous present

2-3 nucleotides/sec  
2-3 Kb (up to 50) read length  
6 TB data in 30 minutes

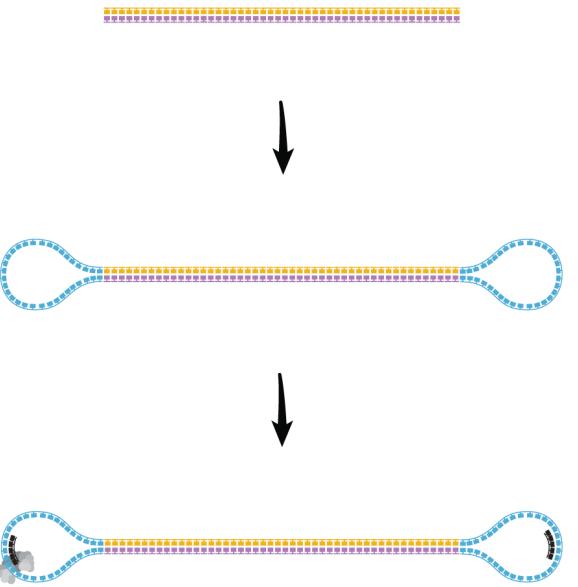
laser damages polymerase



Start with high-quality double stranded DNA

Ligate SMRTbell adapters and size select

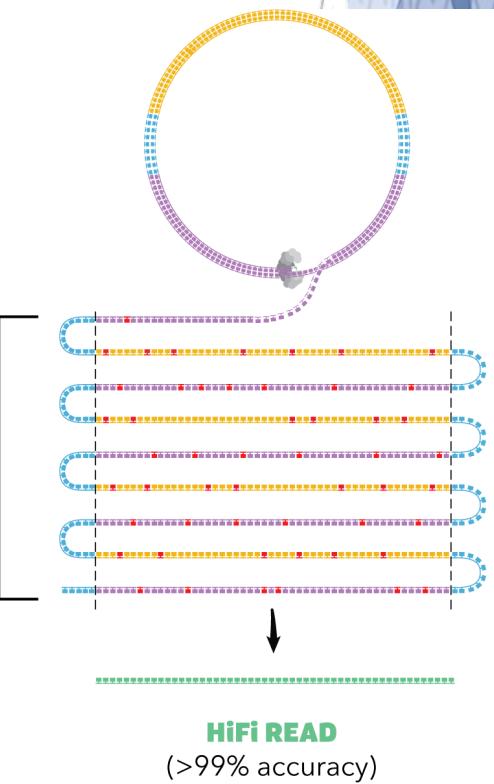
Anneal primers and bind DNA polymerase



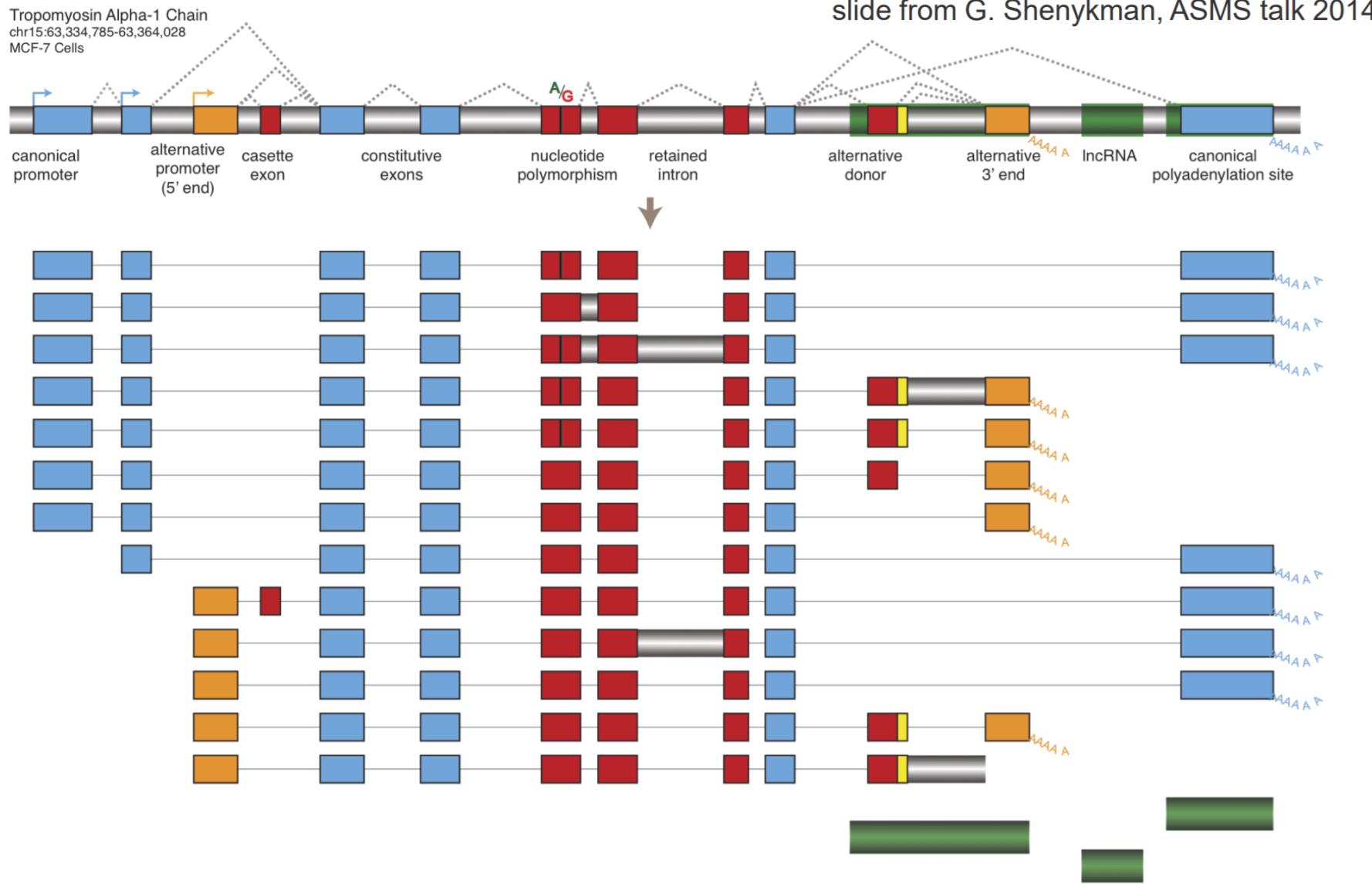
Circularized DNA is sequenced in repeated passes

The polymerase reads are trimmed of adapters to yield subreads

Consensus is called from subreads



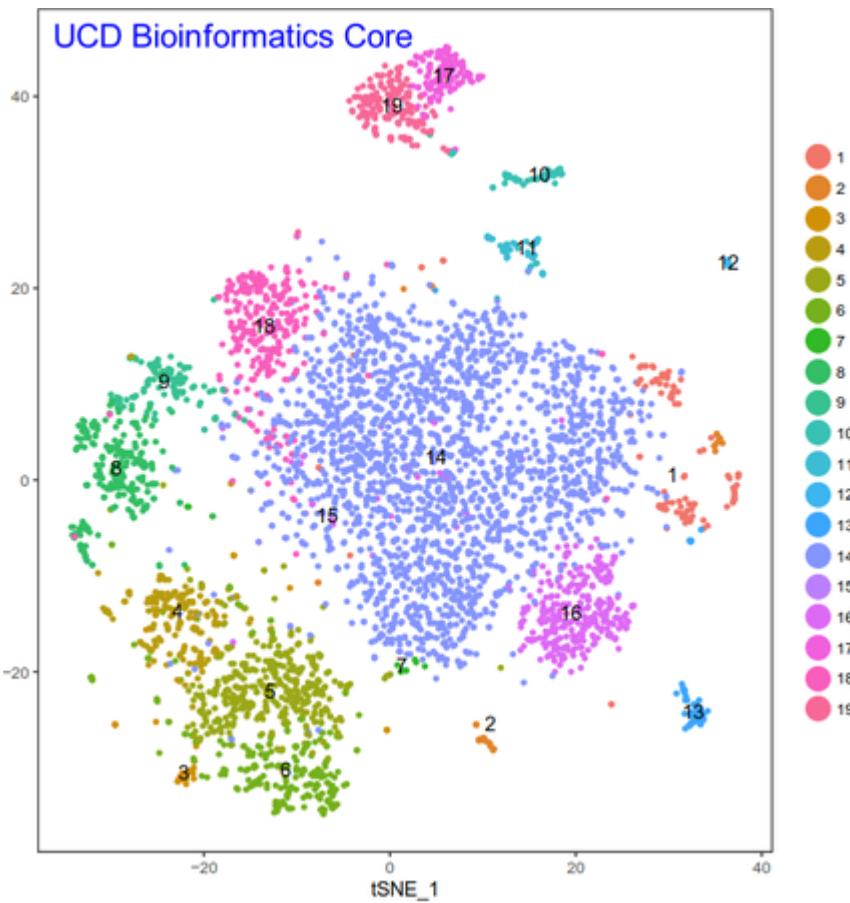
# A Single Gene Locus → Many Transcripts



# Iso-Seq Pacbio

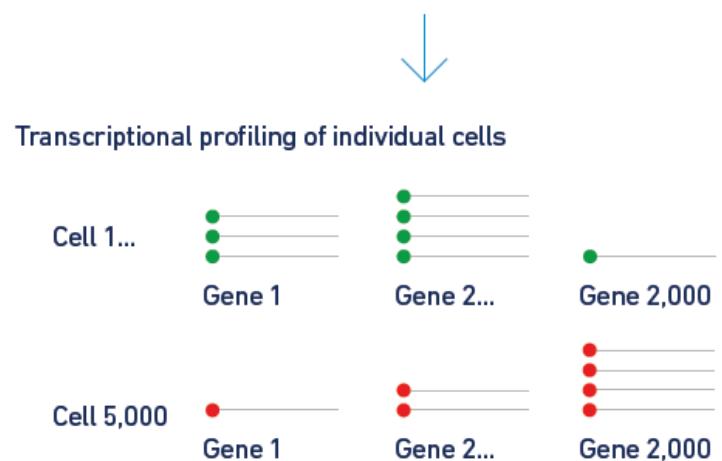
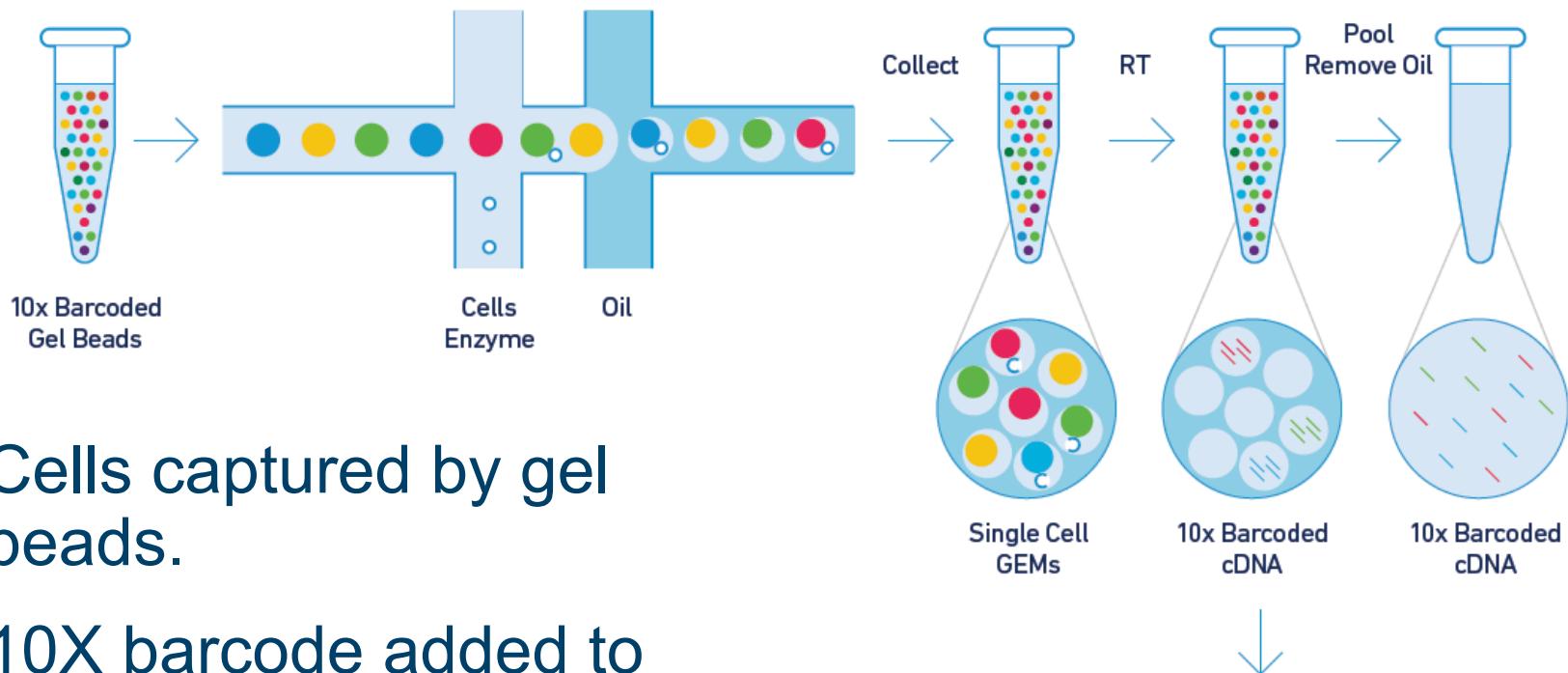
- Sequence full length transcripts  
→ no assembly
- High accuracy (except very long transcripts)
- More than 95% of genes show alternate splicing
- On average more than 5 isoforms/gene
- Precise delineation of transcript isoforms  
( PCR artifacts? chimeras?)

# scRNA-seq (single cells)



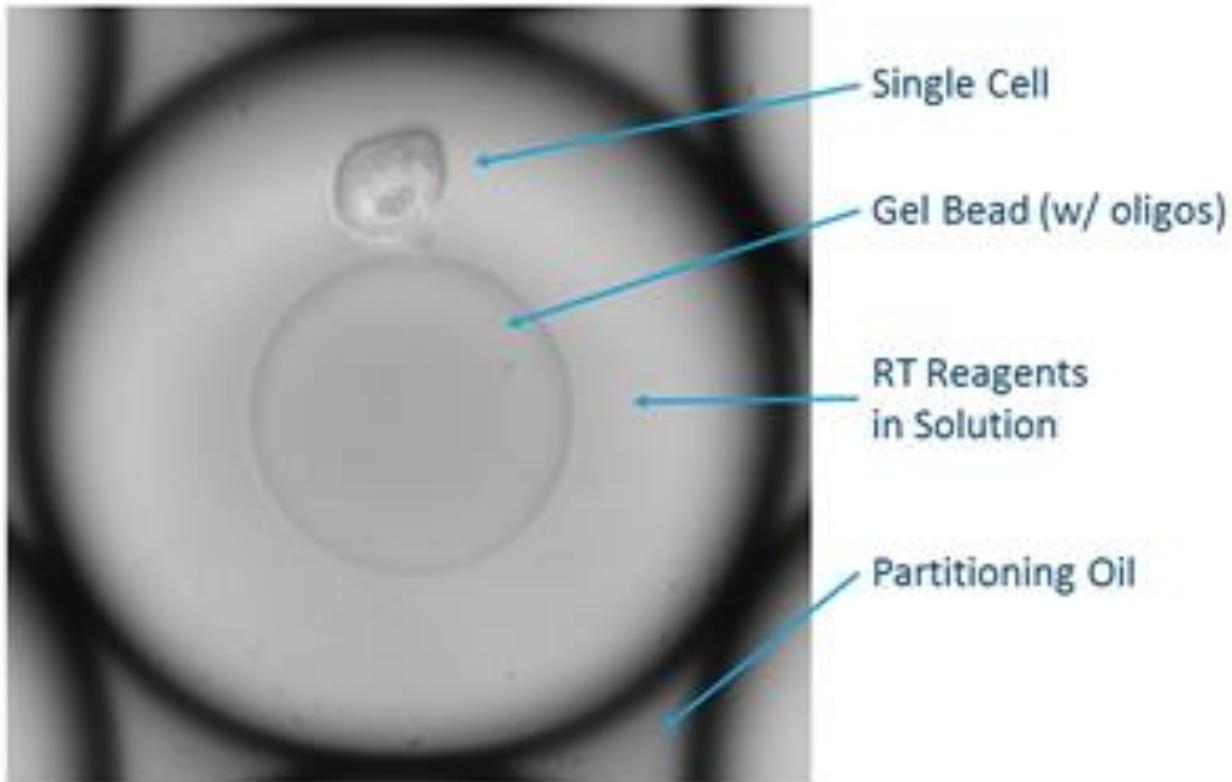
- Gene expression profiling of individual cells.
- Resulting data can distinguish cell types and cell cycle stages - no longer a mix
- Allows the analysis of low abundance cell types

# cDNA preparation



# Cell partitioning into GEMs

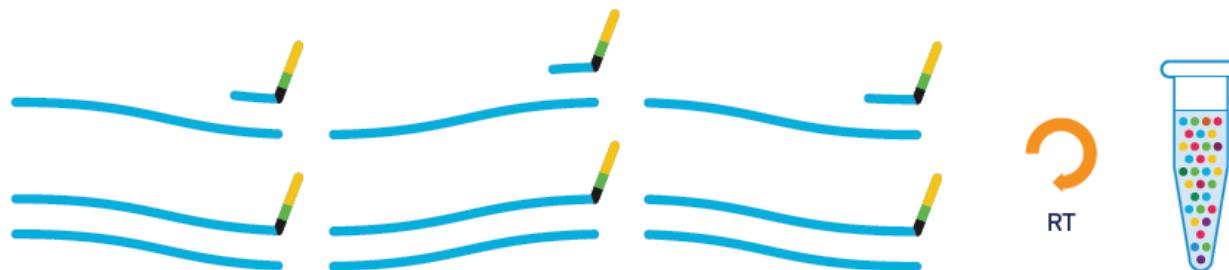
- GEM



Credit: 10X Genomics

# Library preparation

## 1 Molecular Barcoding in GEMs



Credit: 10X Genomics

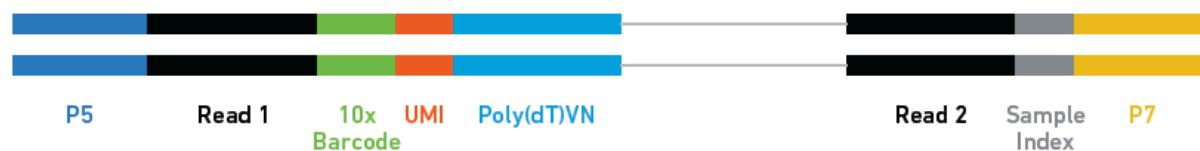
## 2 Pool, Library Prep



## 3 Sequence and Analyze



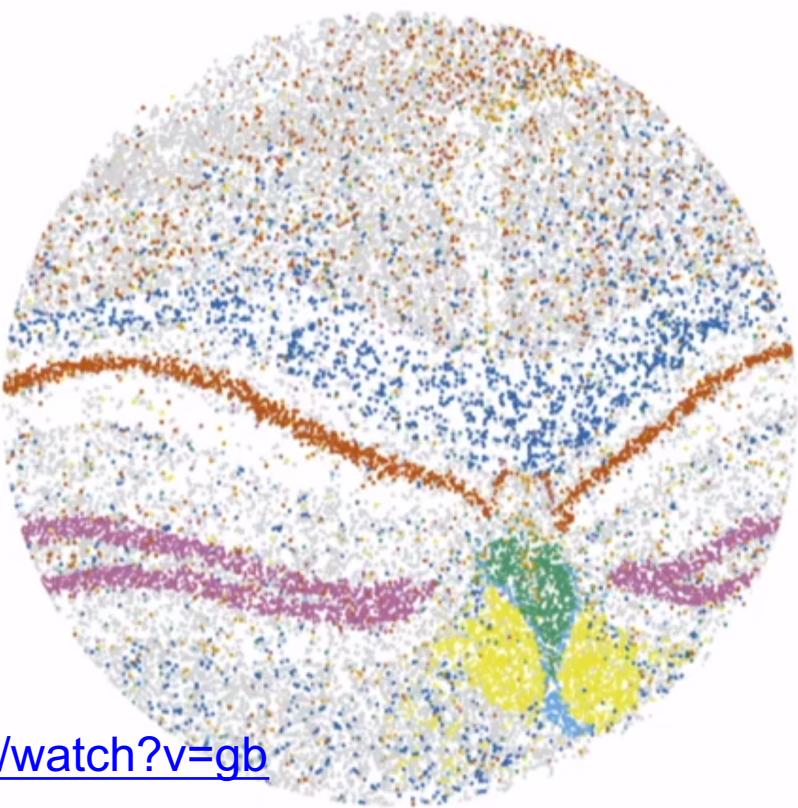
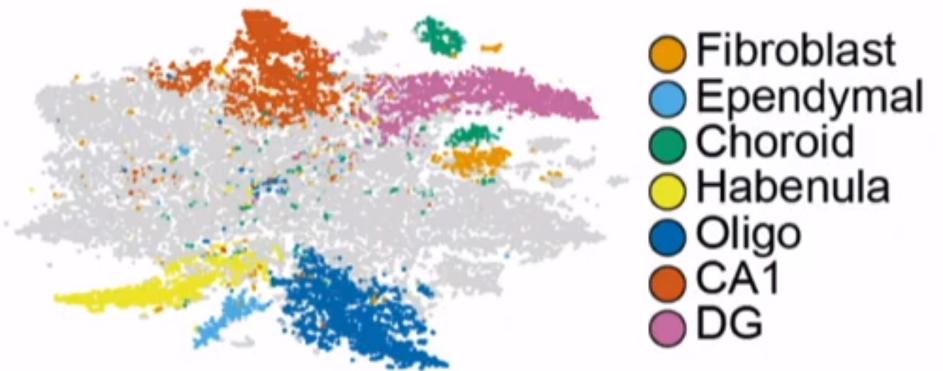
## Final Library Construct



# Spatial Transcriptomics (10XGenomics Visium; Slide-Seq)

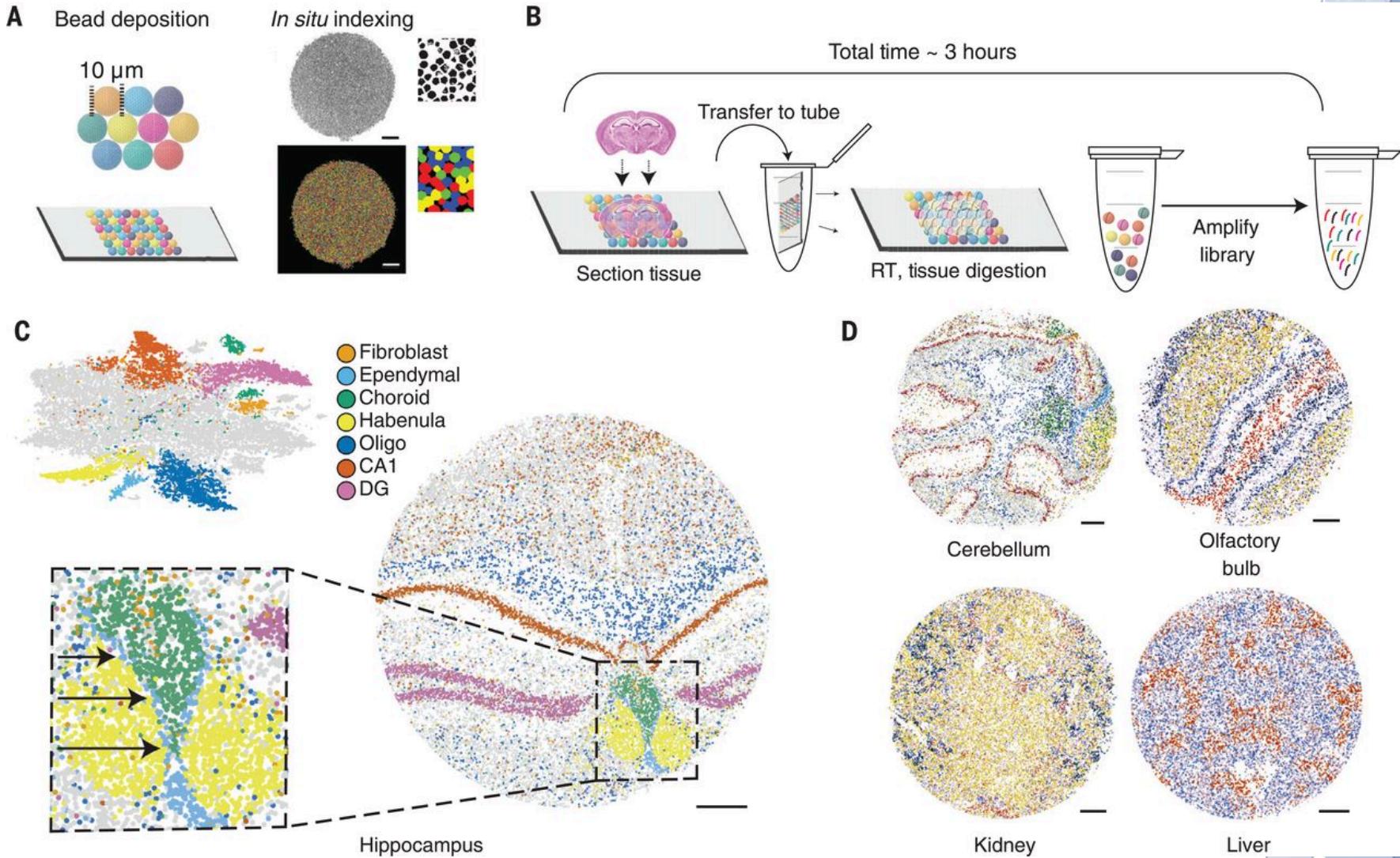
TGAGAT  
TATGAAGG  
TAATCTCT  
TACCCCT  
GCTGAAAC

## 4 Map gene expression into space

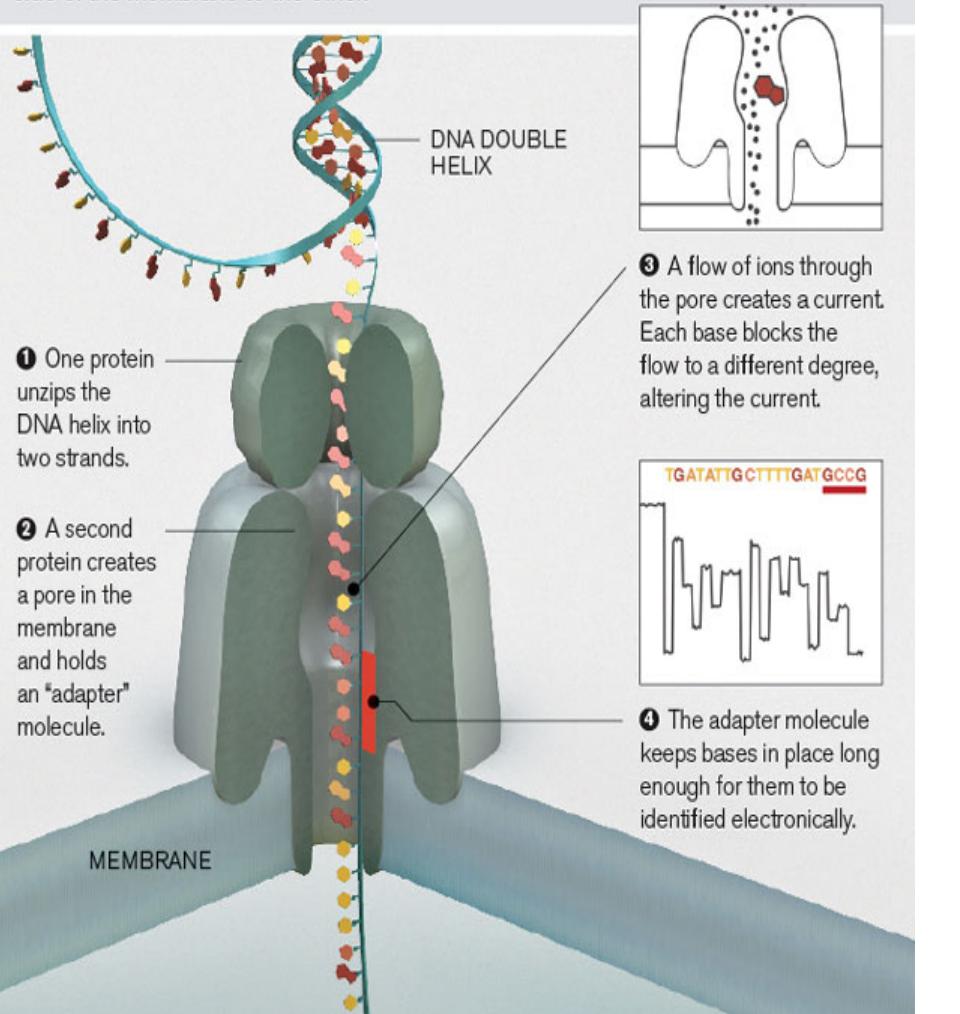


Evan Macosko

<https://www.youtube.com/watch?v=gb0vgwIQPo8&t=2783s>



DNA can be sequenced by threading it through a microscopic pore in a membrane. Bases are identified by the way they affect ions flowing through the pore from one side of the membrane to the other.



# Future's so bright





# Thank you!