

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](#))

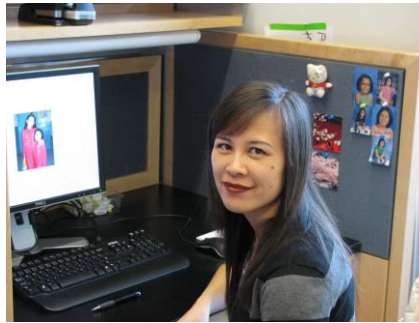
- introducing new technologies to the campus
- shared equipment
- teaching (workshops)



The DNA Tech Core Team



Emily



Oanh



Diana



Siranoosh



Vanessa



Ruta





Complementary Approaches

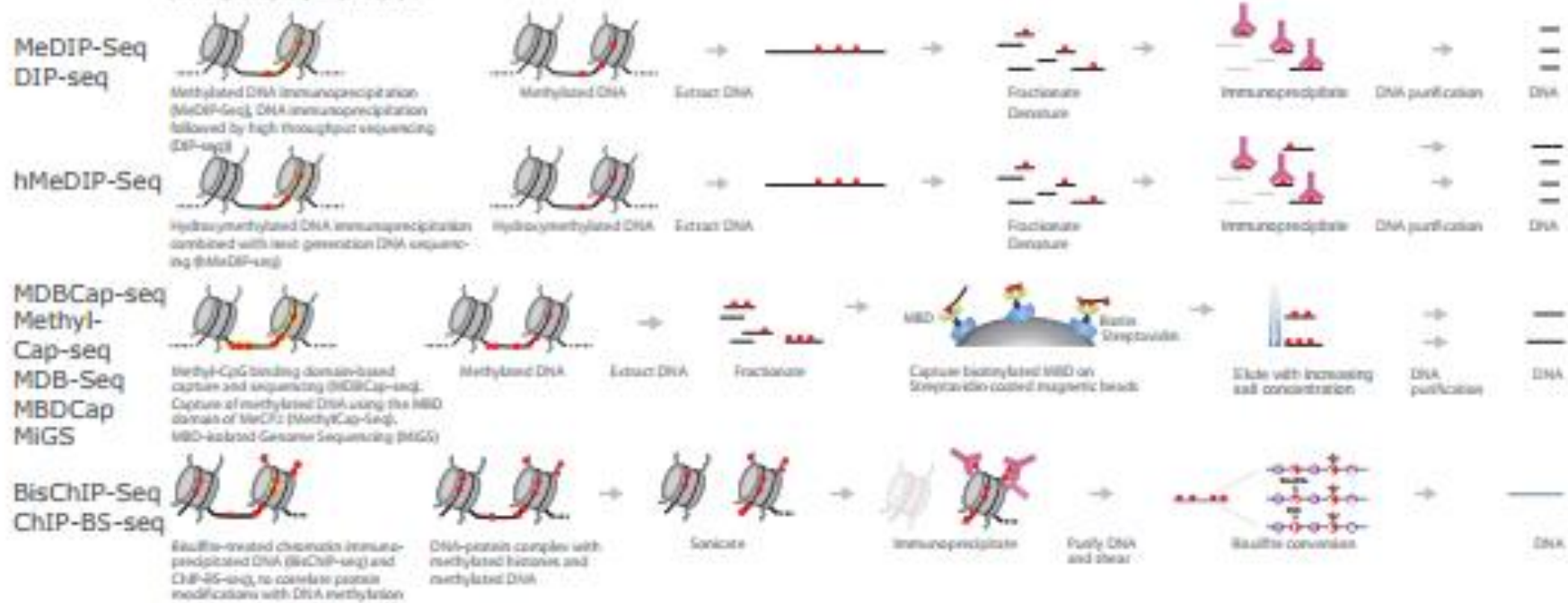
ILLUMINA	PACBIO	PROMETHION NANOPORE
Still-imaging of clusters (~1000 clonal molecules)	Movie recordings fluorescence of single molecules	Recording of electric current through a pore
Short reads - 2x300 bp Miseq	Up to 70 kb, N50 25 kb	Up to 70 kb, N50 25 kb
Repeats are mostly not analyzable	spans retro elements	spans retro elements
High output - up to 2.4 Tb per lane	up to 100 Gb per SMRT-cell, up to 20 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 8-10 %
Considerable base composition bias	No base composition bias	Some systematic errors
Very affordable	Costs 5 to 10 times higher	Costs same or 2x higher
<i>De novo</i> assemblies of thousands of scaffolds	“Near perfect” genome assemblies; lowest error rate	“Near perfect” genome assemblies with suppl. data; highest contiguity

High Throughput Short Read Sequencing: Illumina

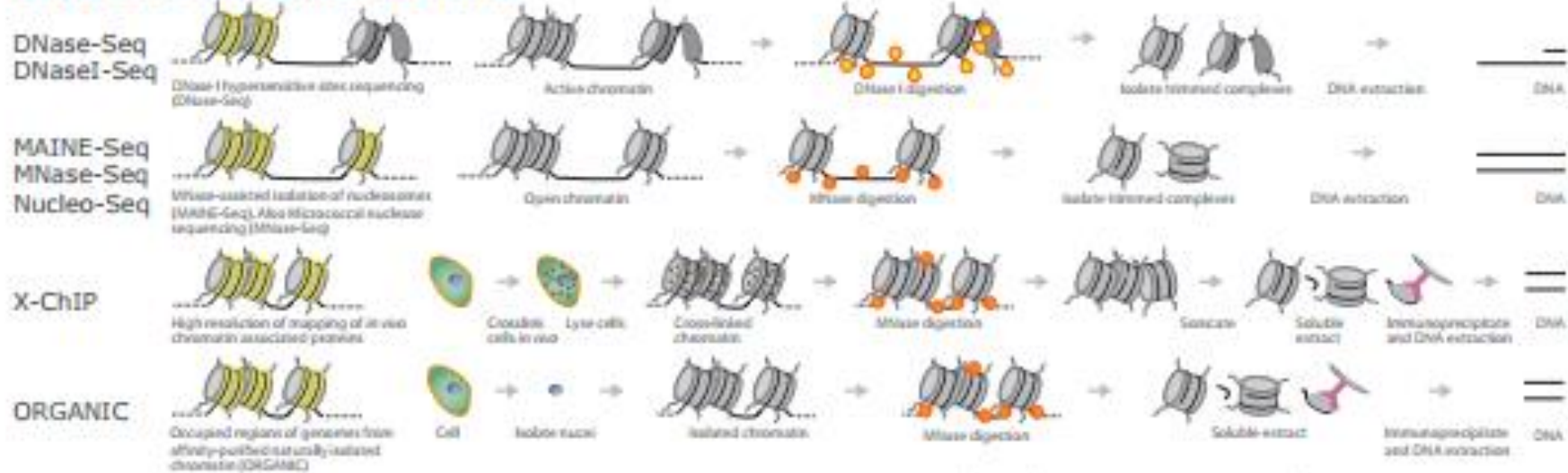


- Whole genome sequencing & Exome sequencing:
Variant detection (small variants SNPs and indels)
Copy number variation (CNVs; prenatal diagnostics)
- Genotyping by sequencing
- Genome assemblies: small genomes
- Metagenomics
- RNA-seq: gene expression, transcript expression
- Small RNA-seq
- Single-cell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

TCTGGGA
GAAATT
TGTTGA
AAGGAG
TTTGGG
CGCCAG
TCCCAG
AATTGC
TCTCCA
AAGGCT
AATTTG
GCACAA
ATACCA
GCTTTT
TTTATC



DNA-Protein Interactions



Long Read Sequencing: PacBio and Nanopore

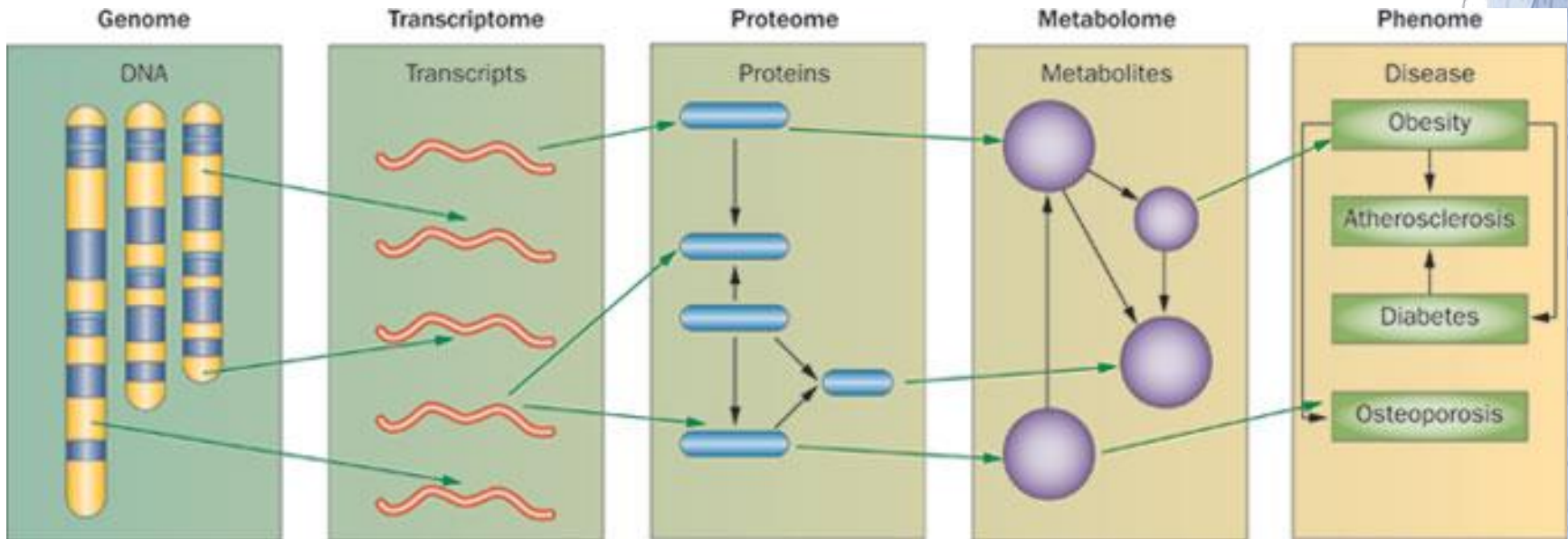


- Whole genome sequencing : Highest quality genome assemblies, Structural variant detection
- RNA-sequencing:
 - full transcript data, Iso-form detection and quantification
 - Direct RNA-seq
- Metagenomics
- Epigenetics (Nanopore: modified DNA and RNA bases)

CTGGGA
GAAATT
TGTTGA
AAGGAG
TTTGGG
CGCCAG
TCCCAG
AATTGC
TCTCCA
AAGGCT
AATTGA
GCACAA
ATACCA
GCTTTT
TTTATC

“DNA makes RNA and RNA makes protein”

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

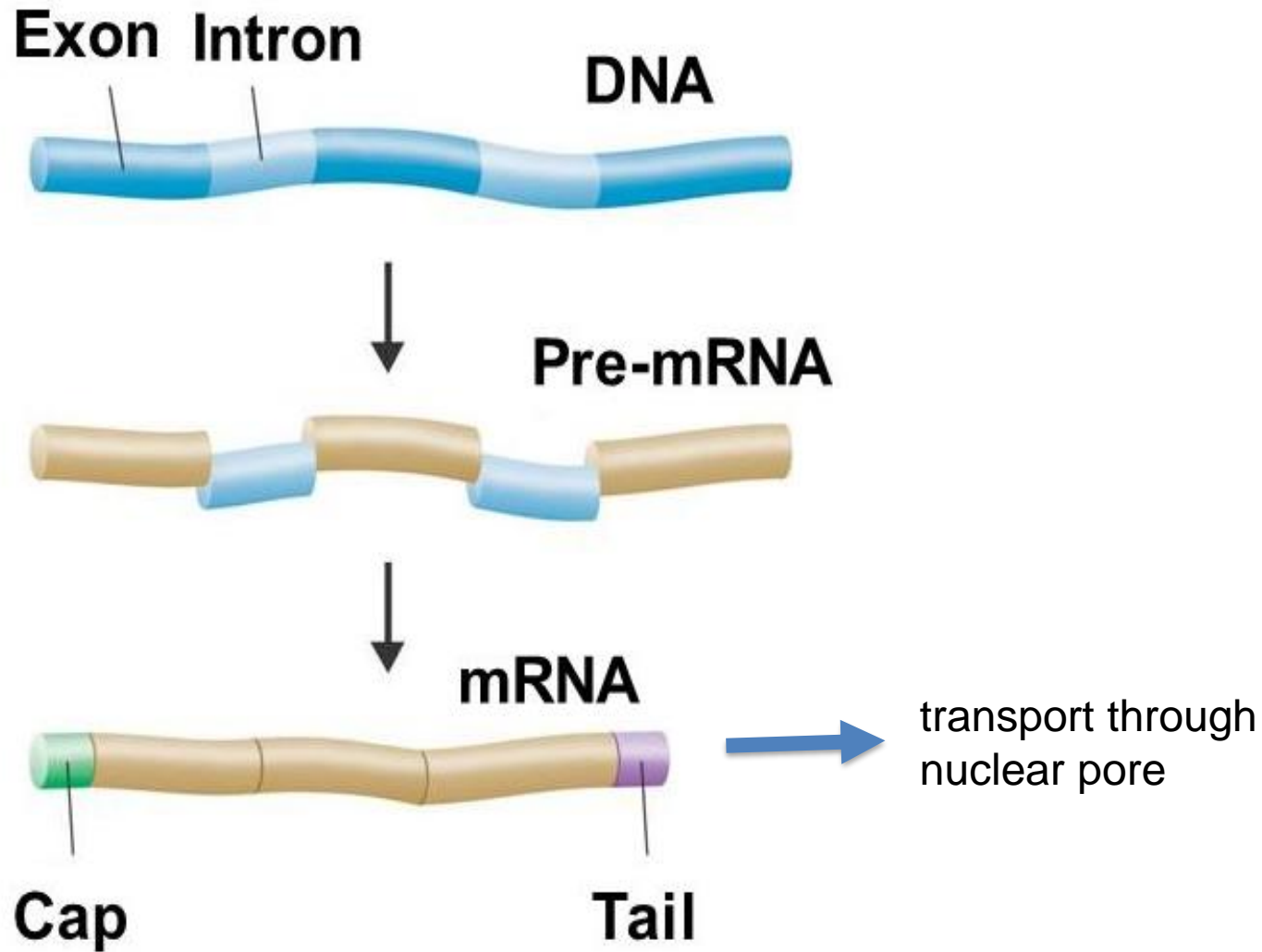


DNA Tech & Expression Analysis Proteomics Core Metabolomics Core

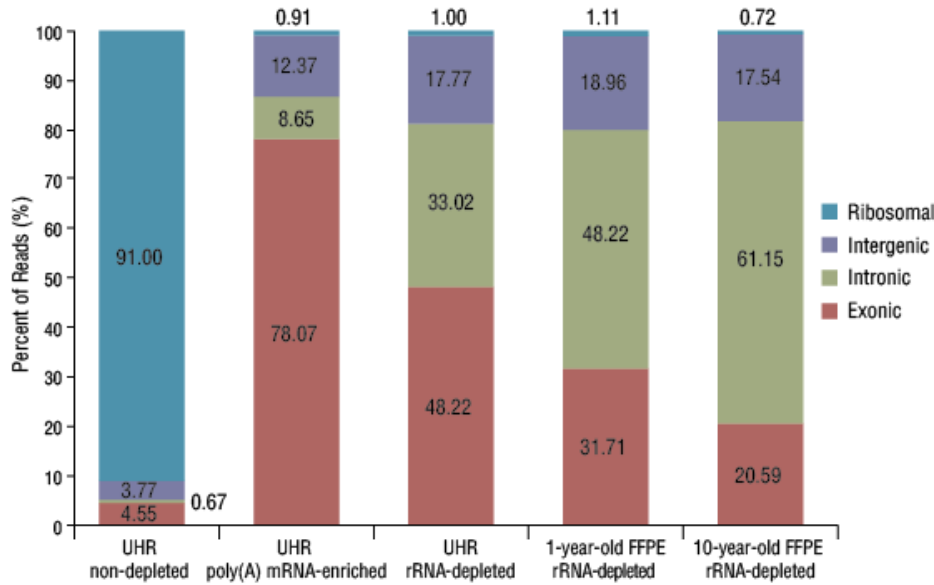
UCD Genome Center

nature
REVIEWS CARDIOLOGY

transcription and processing in nucleus

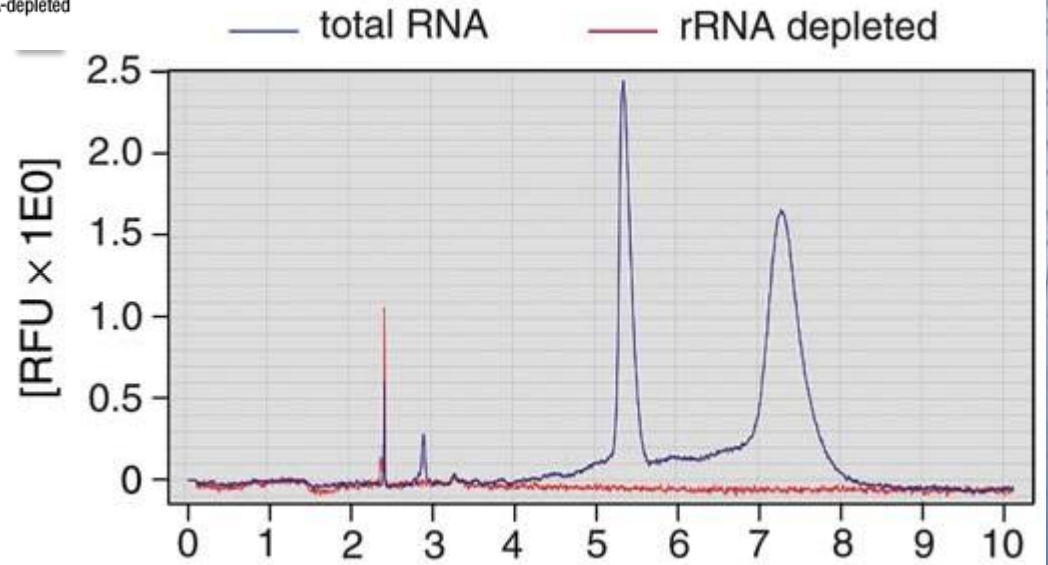


mRNA makes up only about 2% of a total RNA sample



- 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 2nd 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



Fragmentation

- Mechanical shearing:

- BioRuptor
- Covaris

DNA, RNA

- Enzymatic:

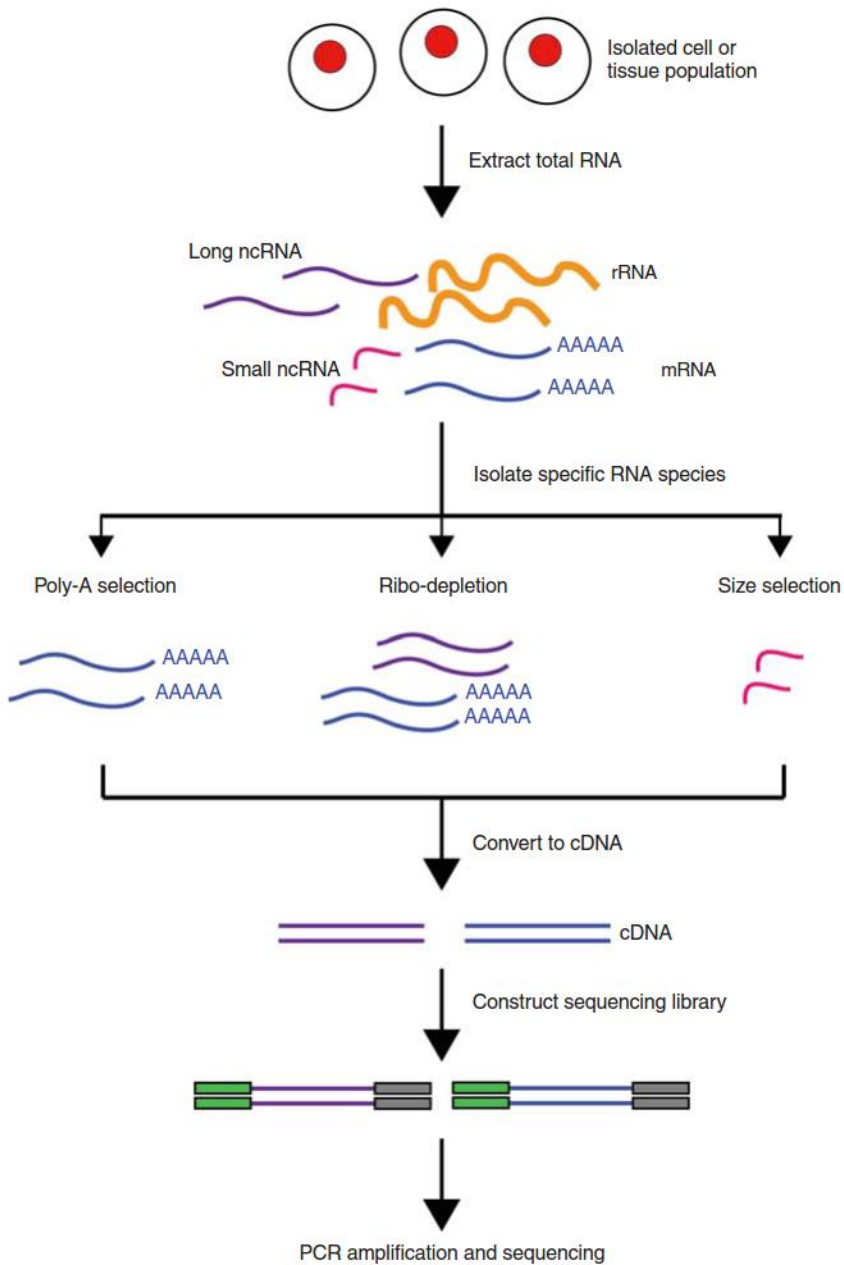
- Fragmentase, RNAse3

DNA, RNA

- Chemical: Mg^{2+} , Zn^{2+}

→ RNA



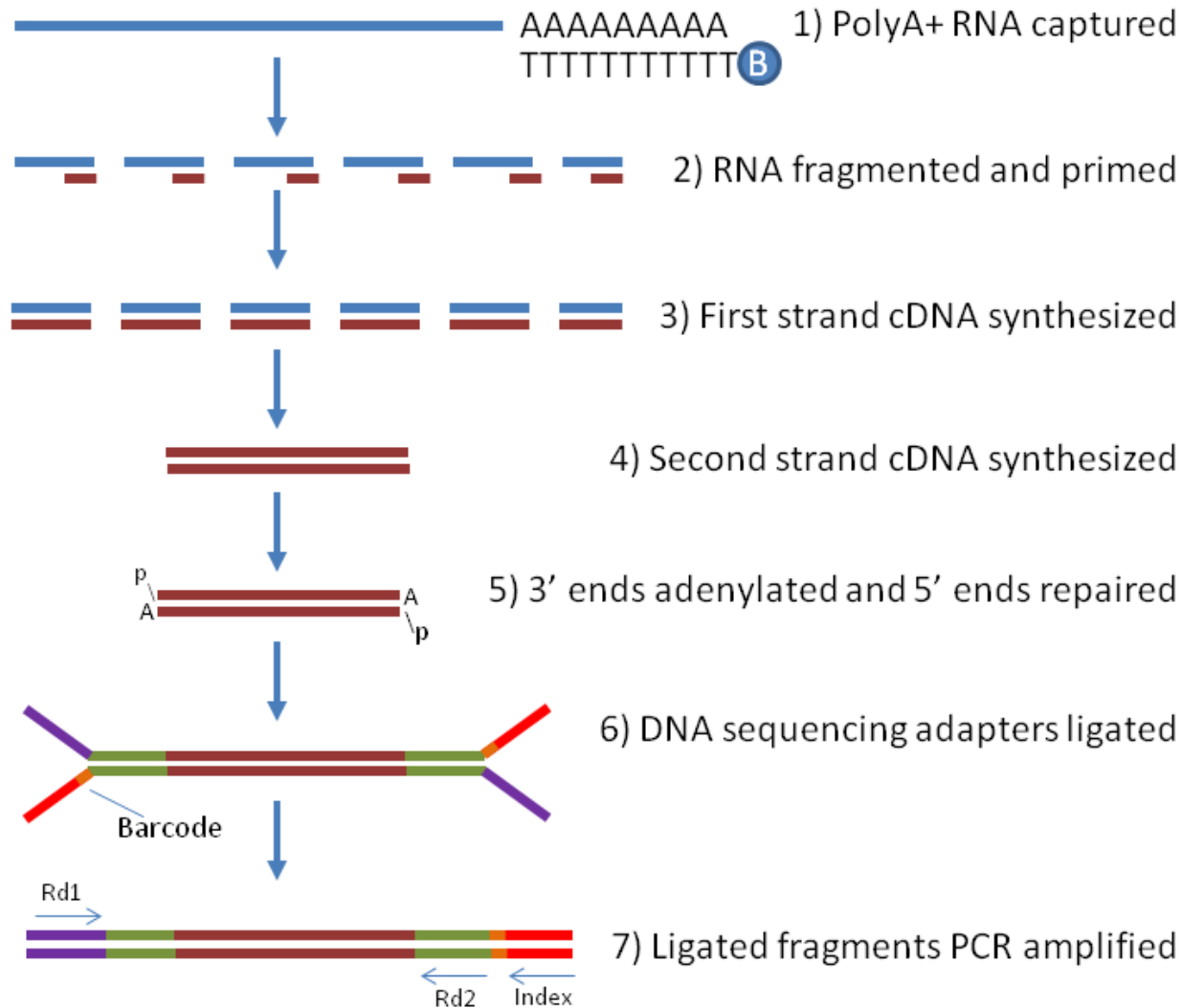


RNA-seq?

Sorry – Illumina and PacBio are only sequencing DNA.



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



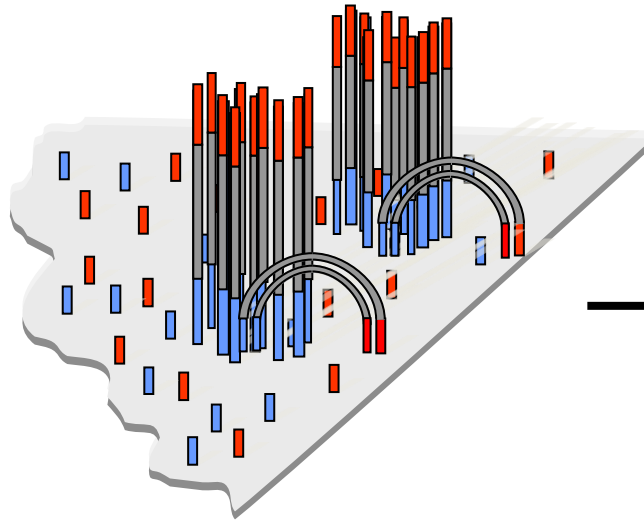
Illumina Sequencing Technology

Sequencing By Synthesis (SBS) Technology

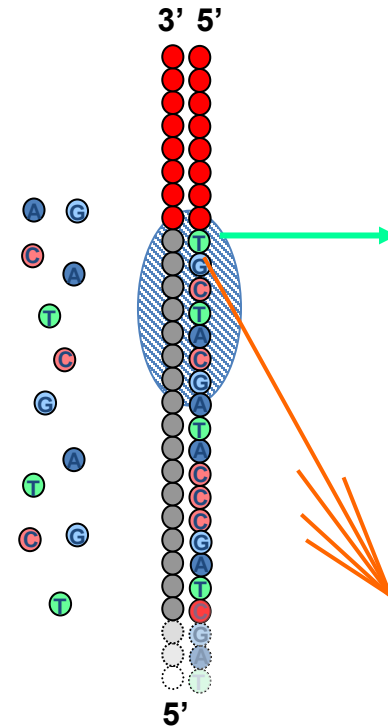
DNA
(0.1-1.0 ug)



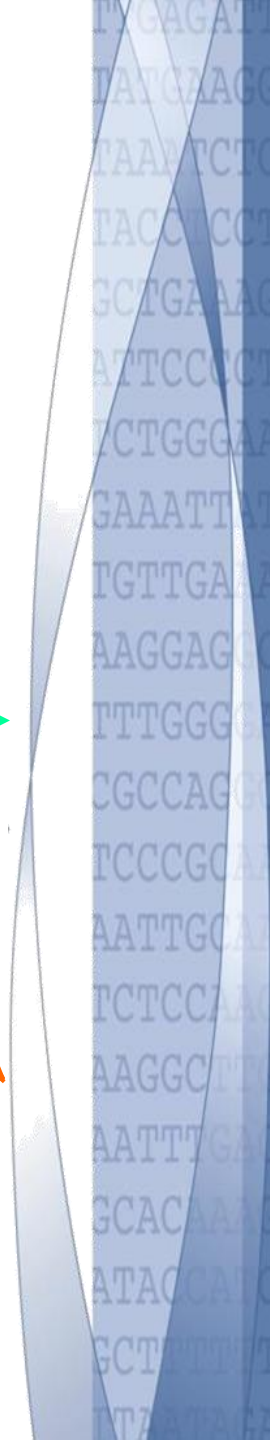
Library
preparation



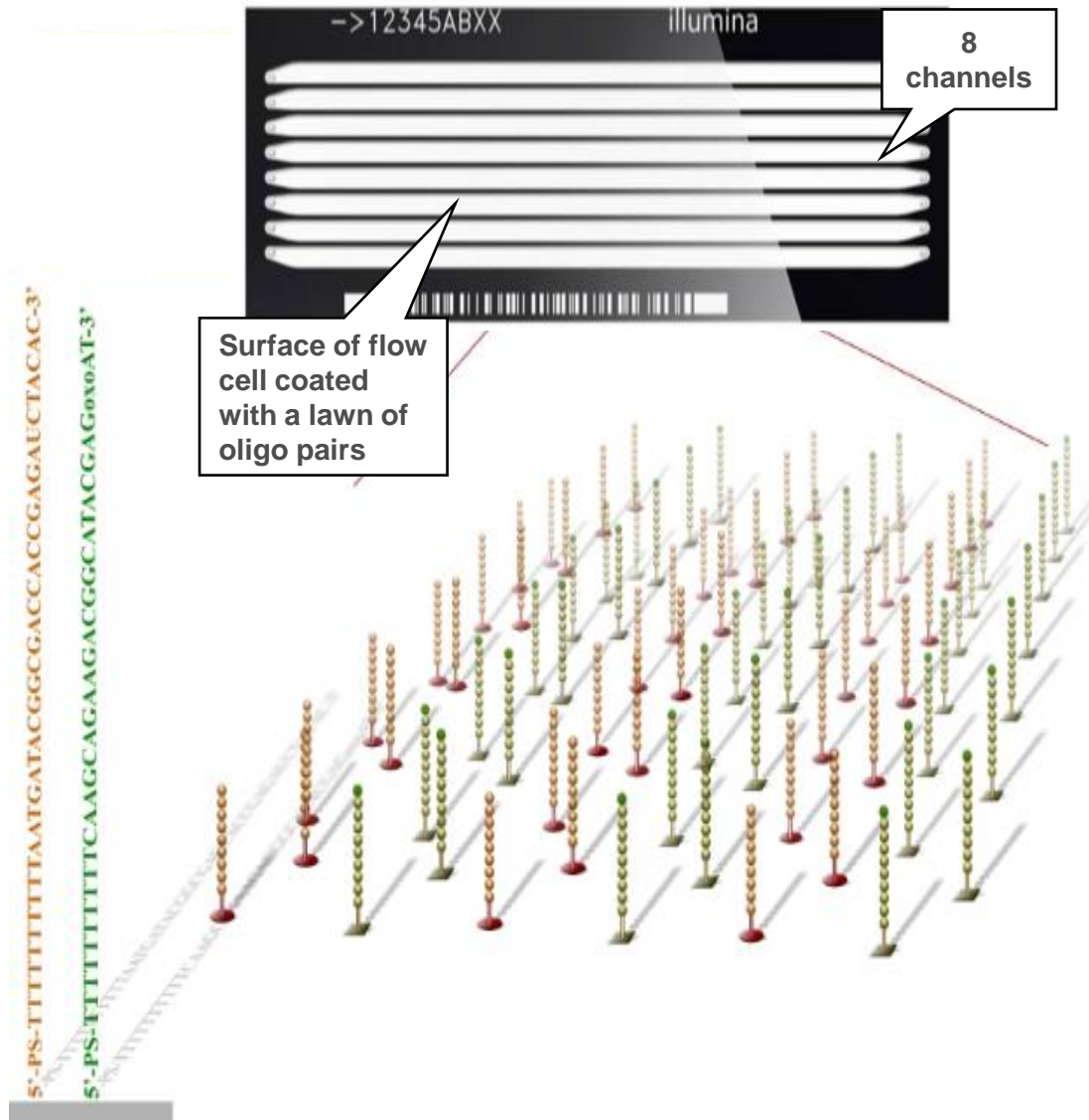
Cluster generation



Sequencing

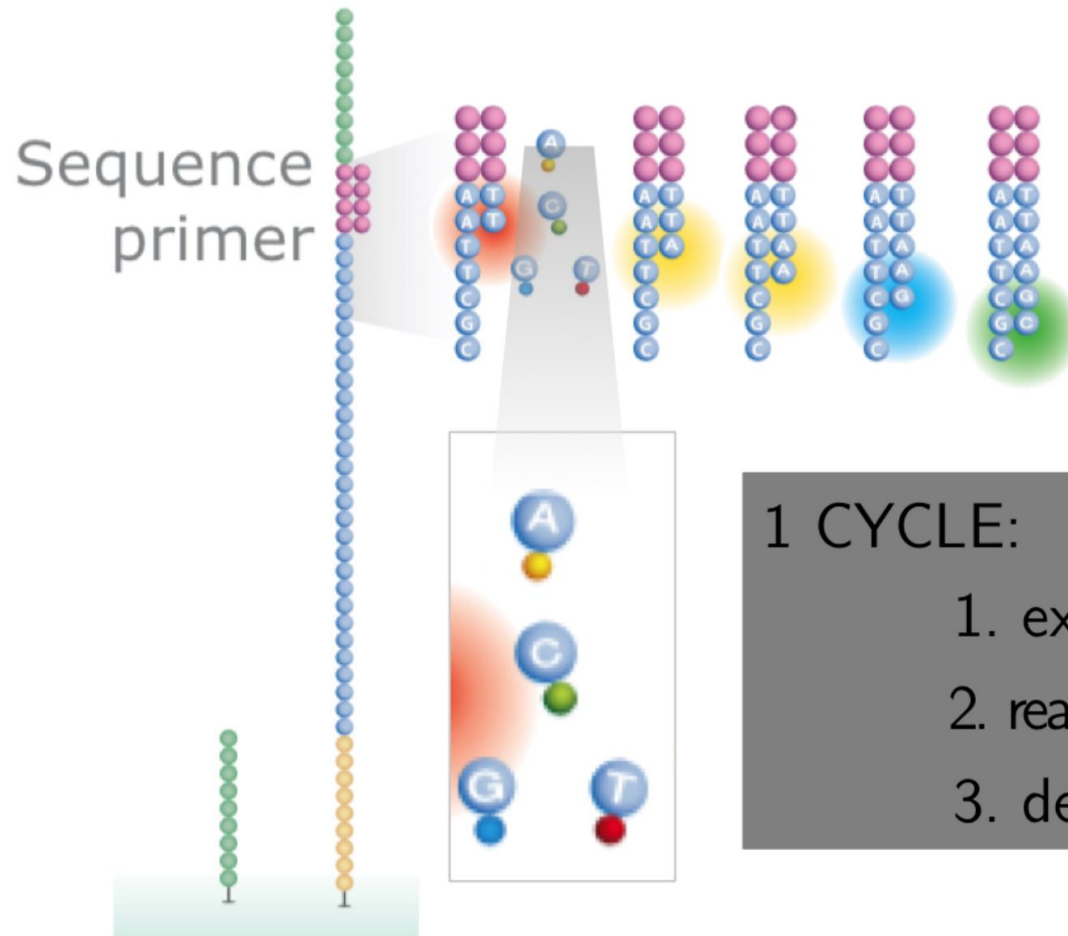


TruSeq Chemistry: Flow Cell



GCTGAAAC
AATCCCT
TCTGGGA
GAAATT
TGTTGA
AAGGAG
TTTGGG
CGCCAG
TCCCAG
AATTGC
TCTCCA
AAGGCT
AATTGA
GCACAA
ATACCA
GCTTTT
TTTATA

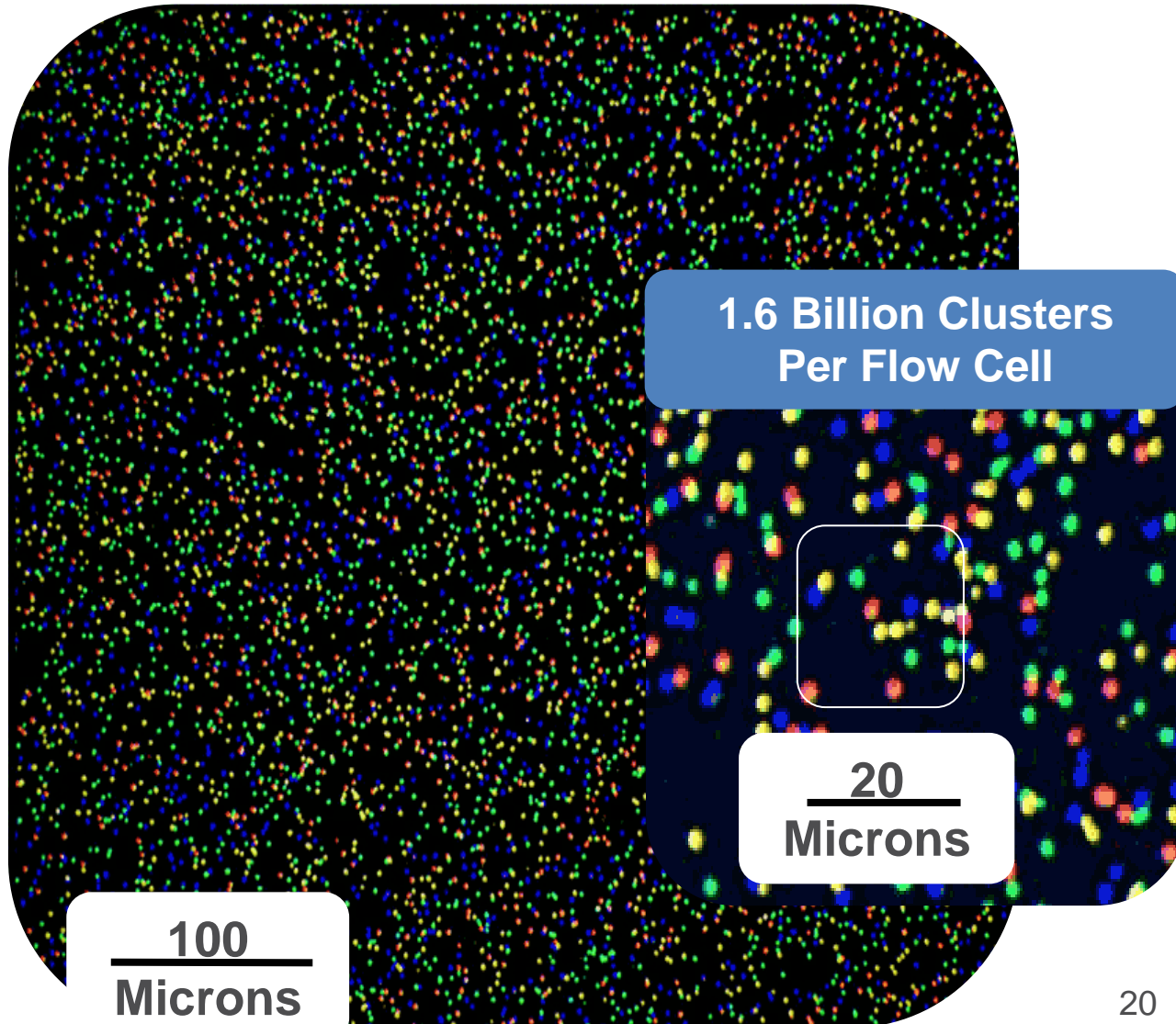
Illumina's sequencing is based on **fluorophore-labelled dNTPs** with **reversible** terminator elements that will become incorporated and excited by a laser one at a time.



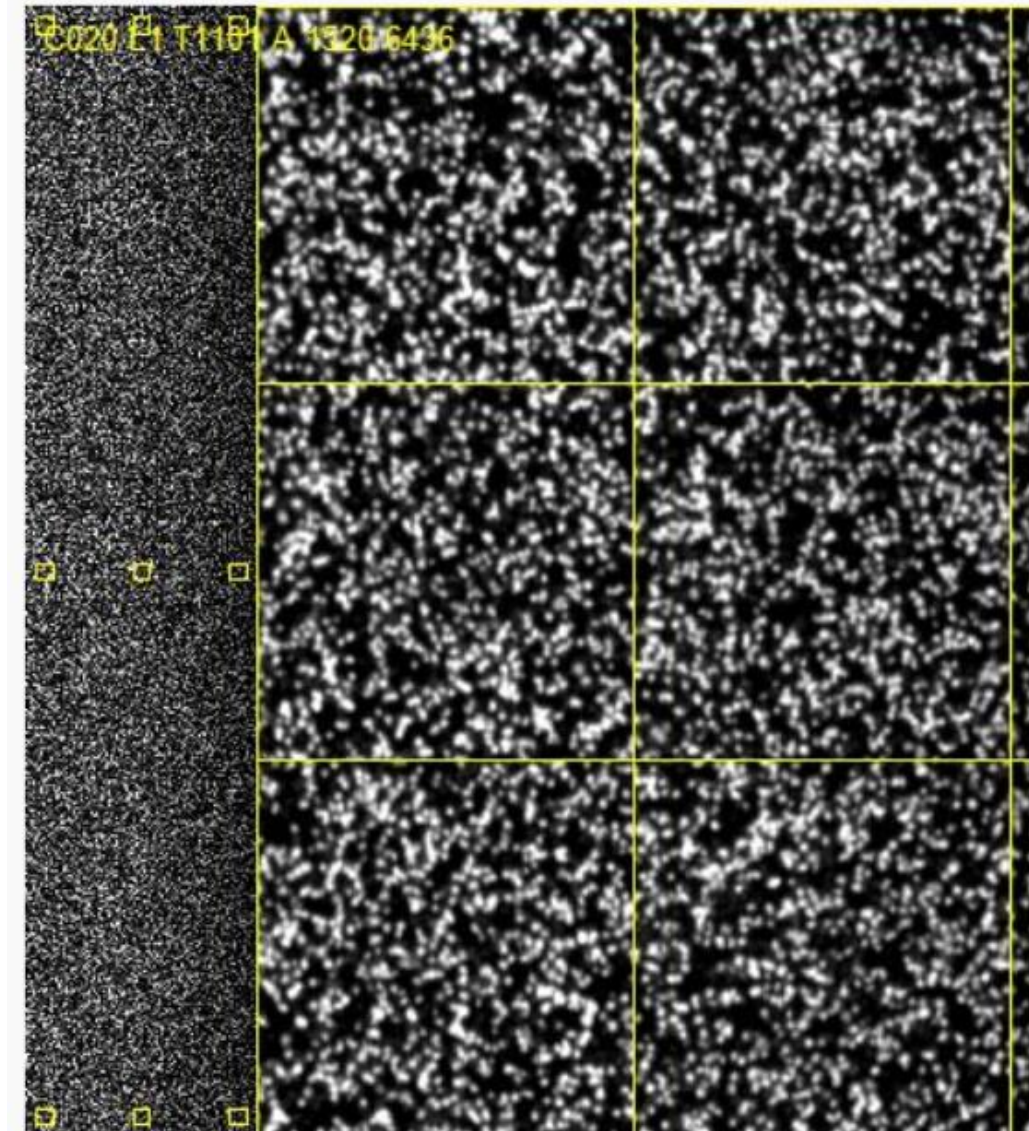
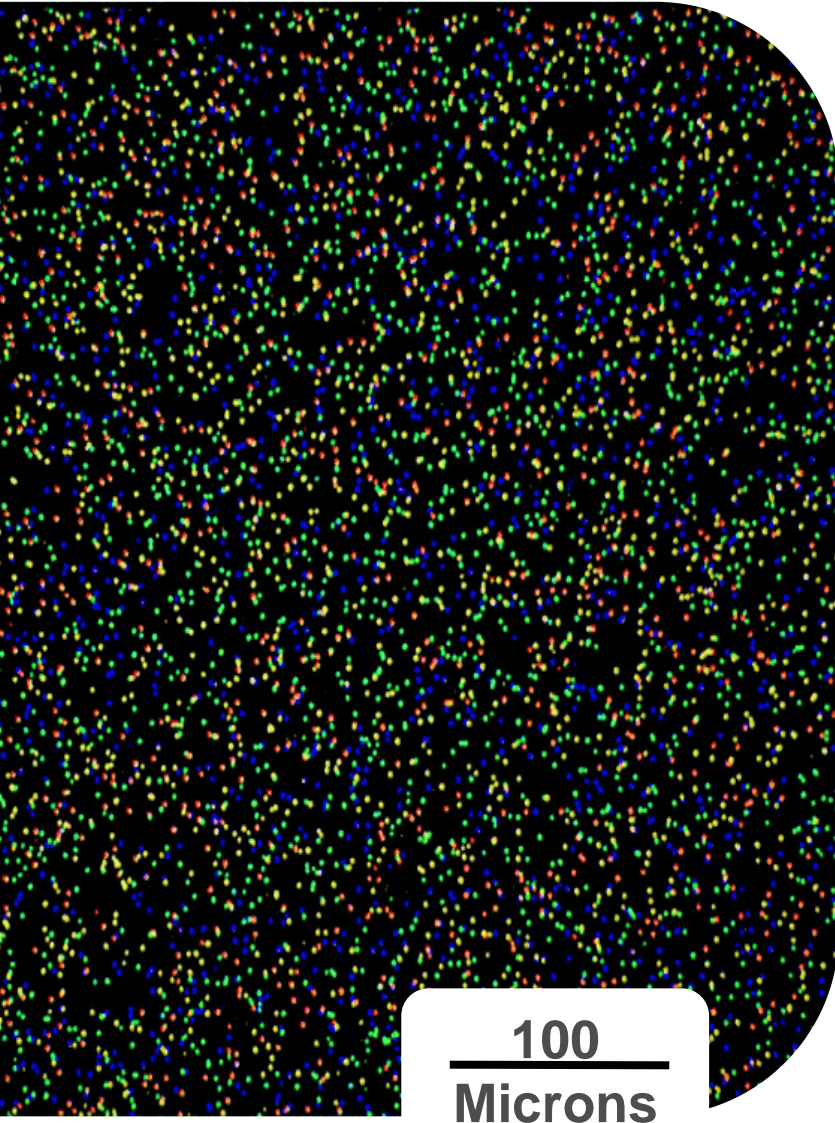
1 CYCLE:

1. extend prev. base
2. read (excite & capture)
3. de-block

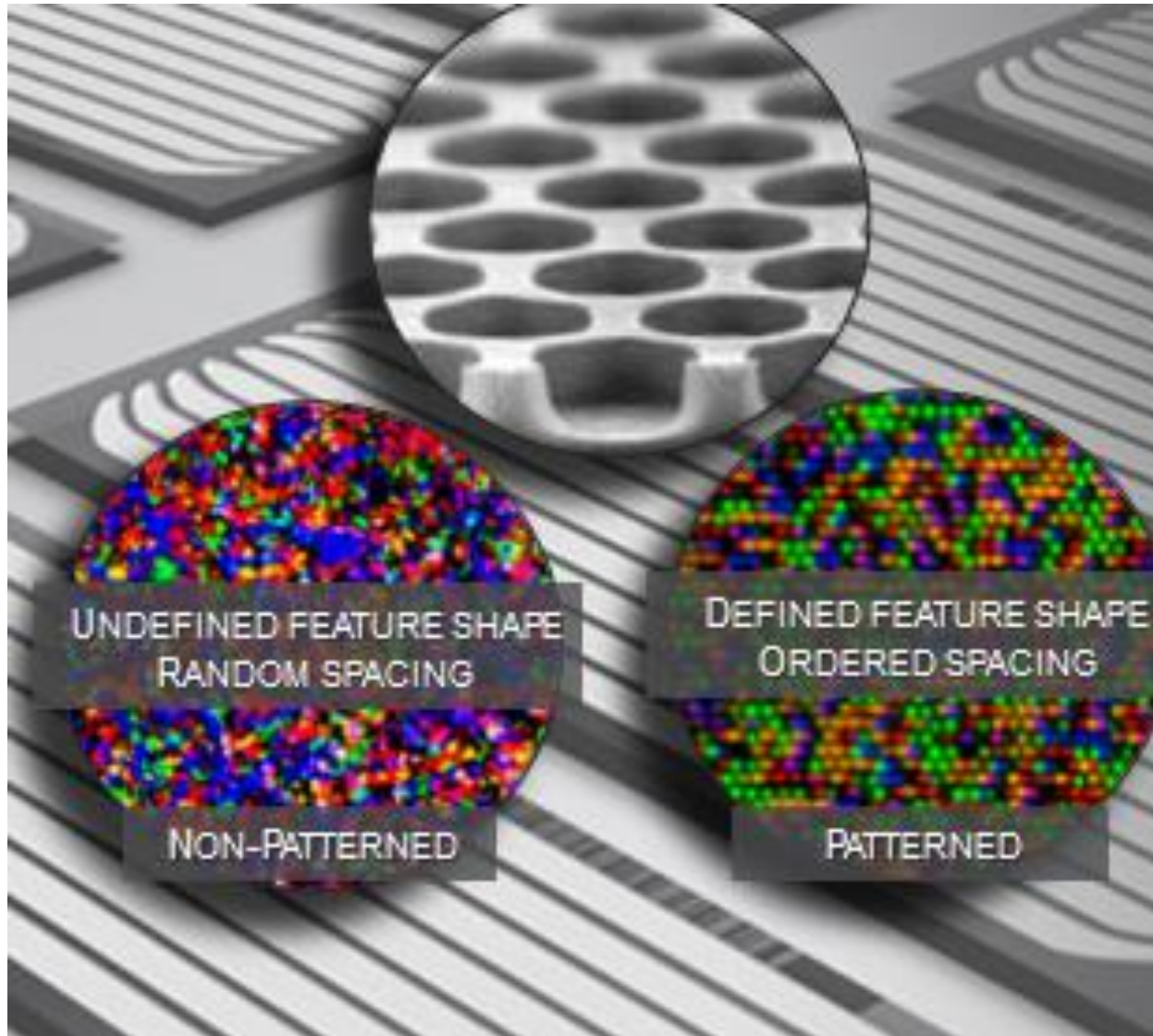
Sequencing



Sequencing

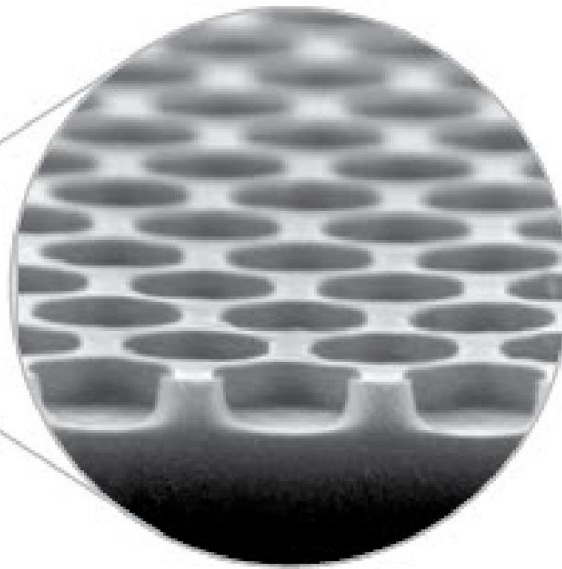


Patterned Flowcell



TTGAGAT
TATGAGC
TAAATCTC
TACCACCT
GCTGAAAC
ATTCCCT
TCTGGGAA
GAAATTAT
TGTGAA
AAGGAG
TTTGGG
CGCCAG
TCCCAG
AATTGCA
TCTCCA
AAGGCT
AATTGA
GCACAA
ATACCA
GCTTTT
TTTATC

Hiseq 3000: 478 million nanowells per lane



What will go wrong ?

- cluster identification
- bubbles
- synthesis errors:

ClusterCluster
Clusts^rCluster
ClusterCluster
ClusterCluster
Cl^lsterCluster



What will go wrong ?

➤ synthesis errors:

ClusterCluster
Clusts^rCluster
ClusterCluster
ClusterCluster
Cl^lsterCluster

Cl^lsterClusterC
ClusterCluster
ClusterCluster
Cl^lusterCluste
ClusterCluster

Phasing & Pre-Phasing
problems



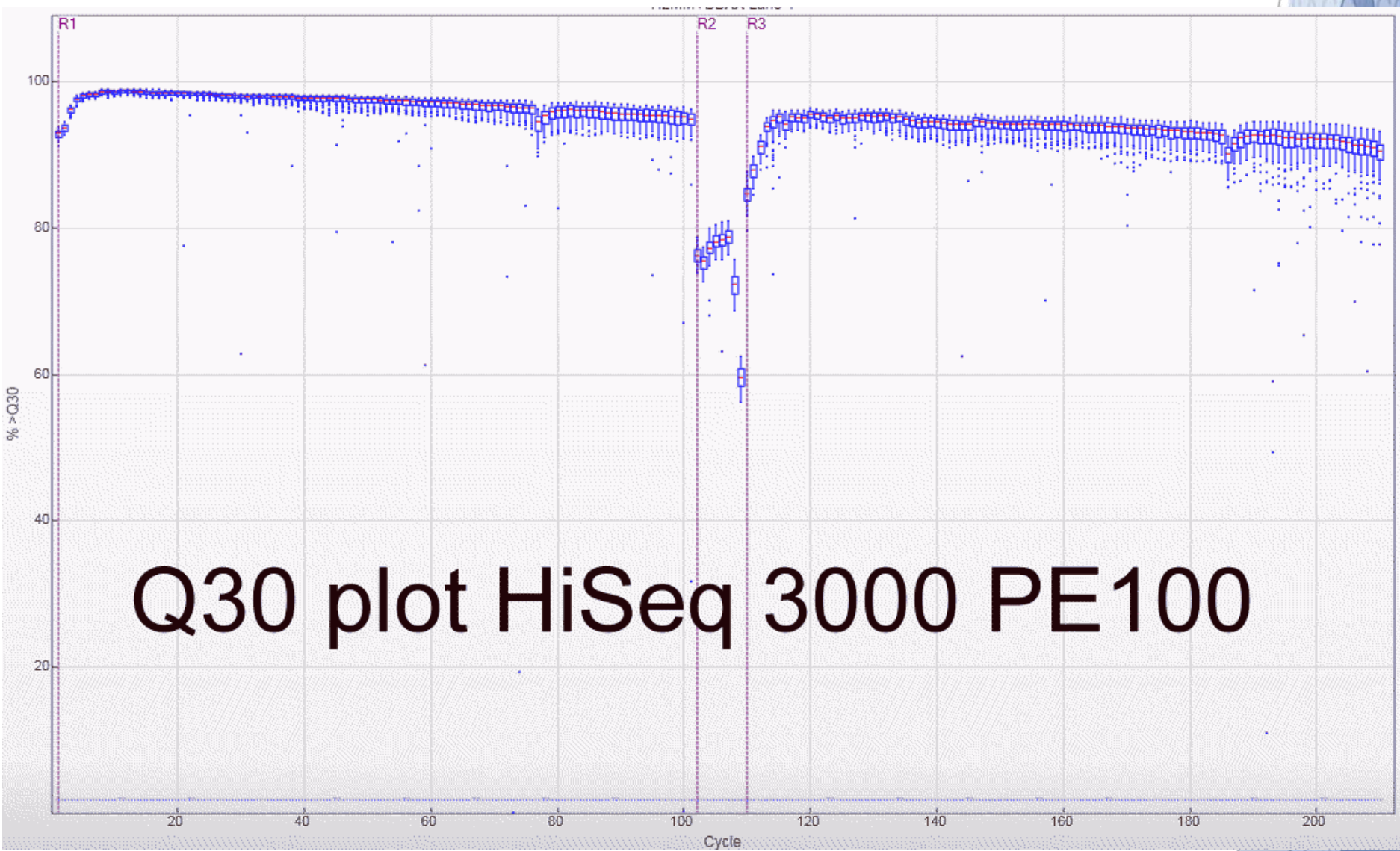
The first lines of your data

@700593F:586:HTWJJBCXX:1:1107:2237:10031 2:N:0:GGAGAACA
NNNNNNNNNNNNNNNNNNNNNNNAGGCCAGCCATAGAACGCTCCCGGCTTCACGGACGT
CATATAGTCAGGCACGAGGTCGGCGCCGAGTTCGTACGCTCGTTGACGACCGCCAT
ACCGCTTGATTTGCGGGGTTGATCGCTAGCGCGGTCCGATTGCGAATGCCCGAGGCAT
ACGTCCGATGGGCCCGCTGGCGCGGTCCACTTCCCATACAACCGCGCGTTCCTCTTC
CAGCGCCATGCCGCGTT

+

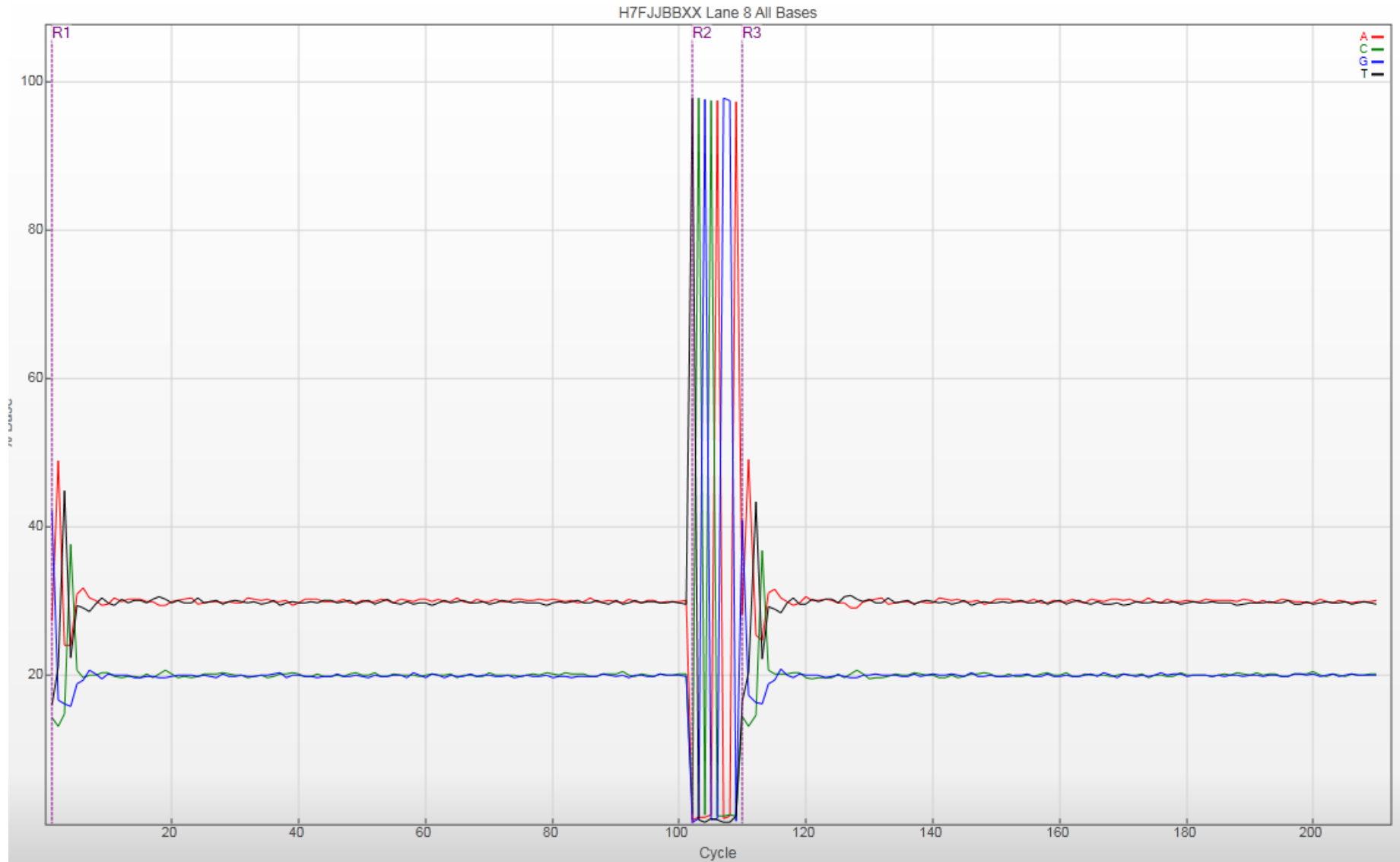
@@@@@@@@@@@@@@@@@@@@@IHHFIIIIIIHHIIIIIIIIIIEHHIIIIIIHHFHIIIIII
IIGHHHHHIHHIIIIIEHHIIHHIHHIGIIHHDHIIIIIIHHHHEEHHIHHCG/EH@GHHII
GIHHIHHIIIIICHGHHIHHIHHIIIIIHHIHHIIHIGIIIGIEH?H?HHHIGHHIIIGHI
GHIDHH@AHHIHHI=HI=GHGHHGD

Illumina SAV viewer

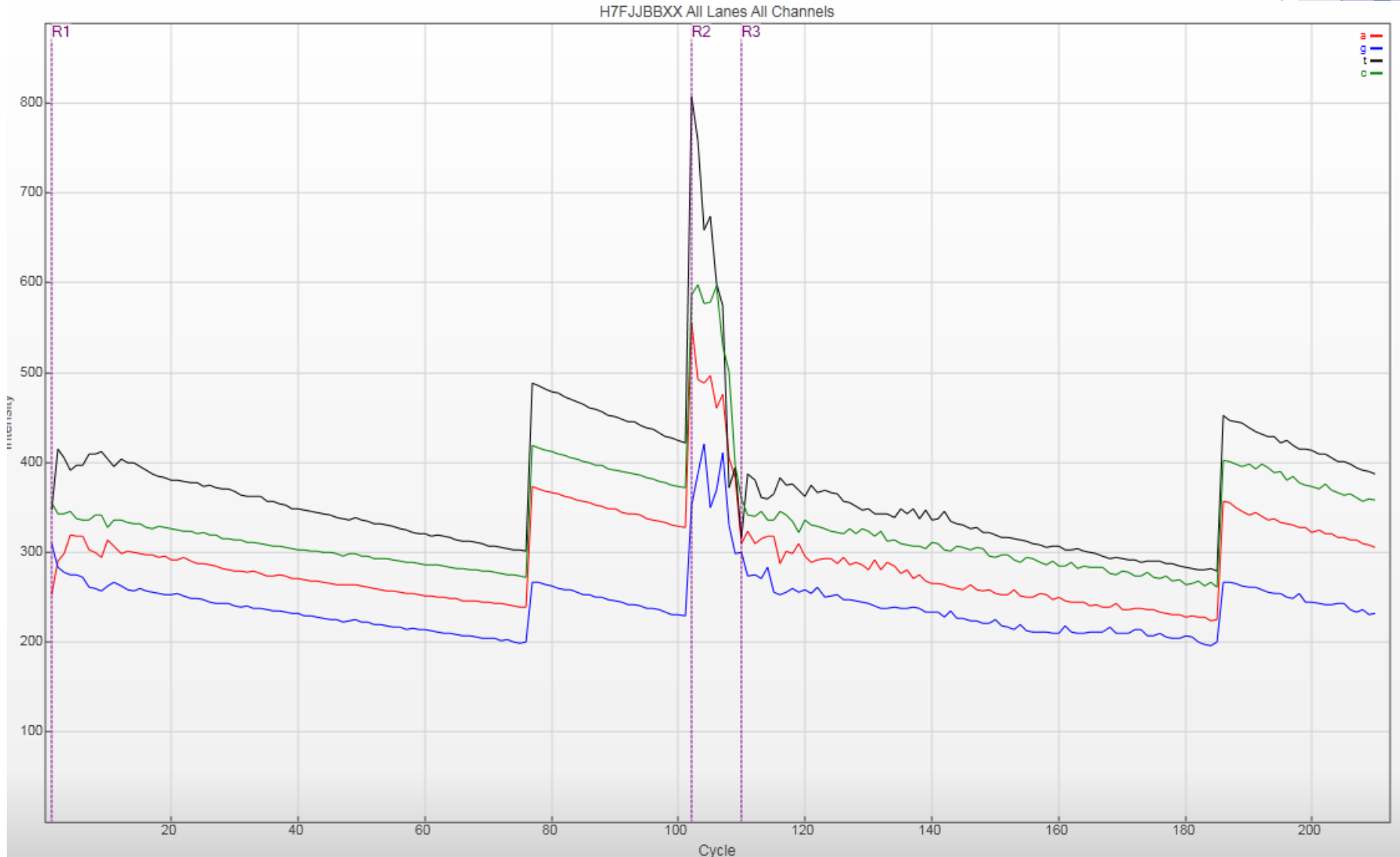


Q30 plot HiSeq 3000 PE100

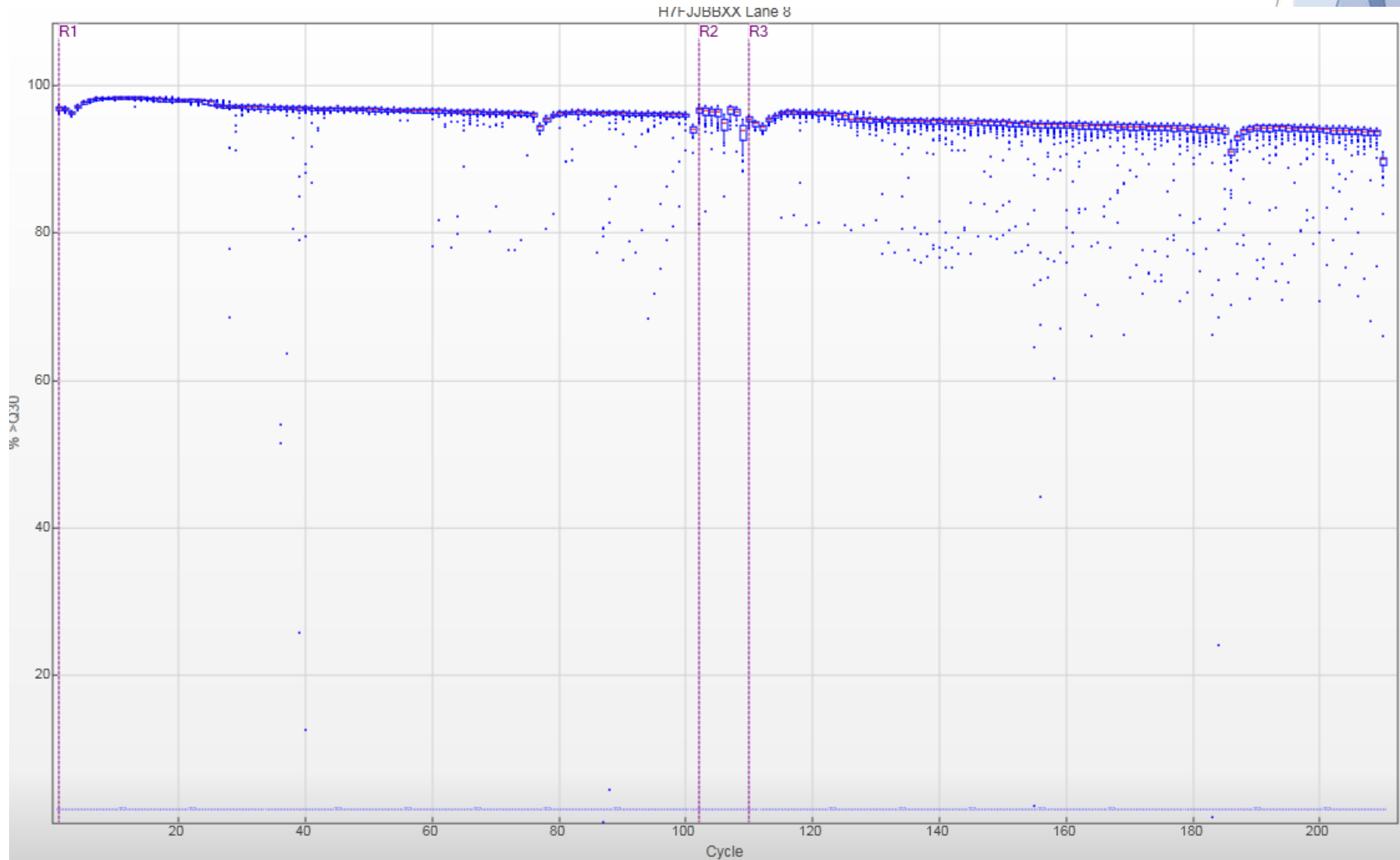
base composition



fluorescence intensity

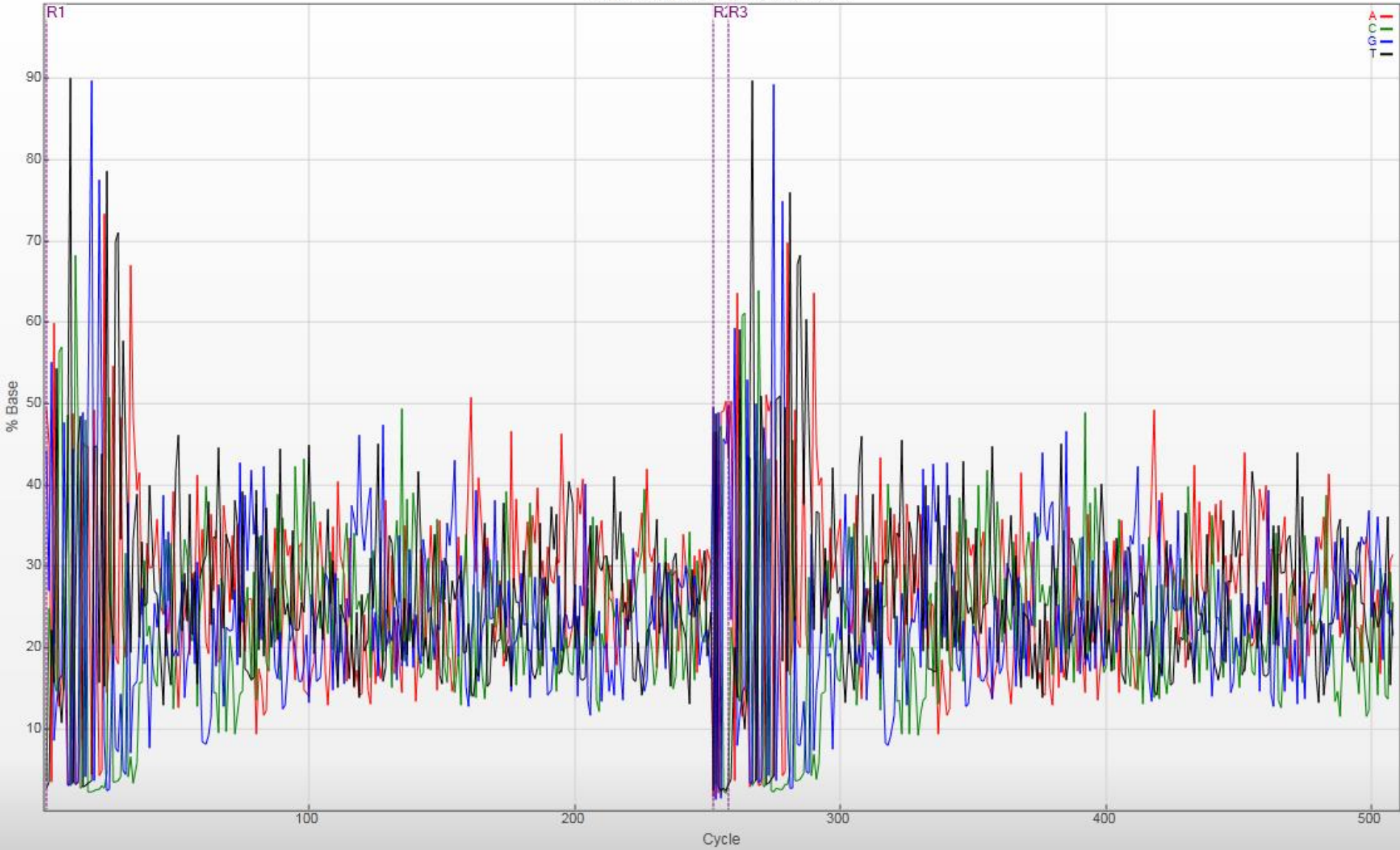


fluorescence intensity

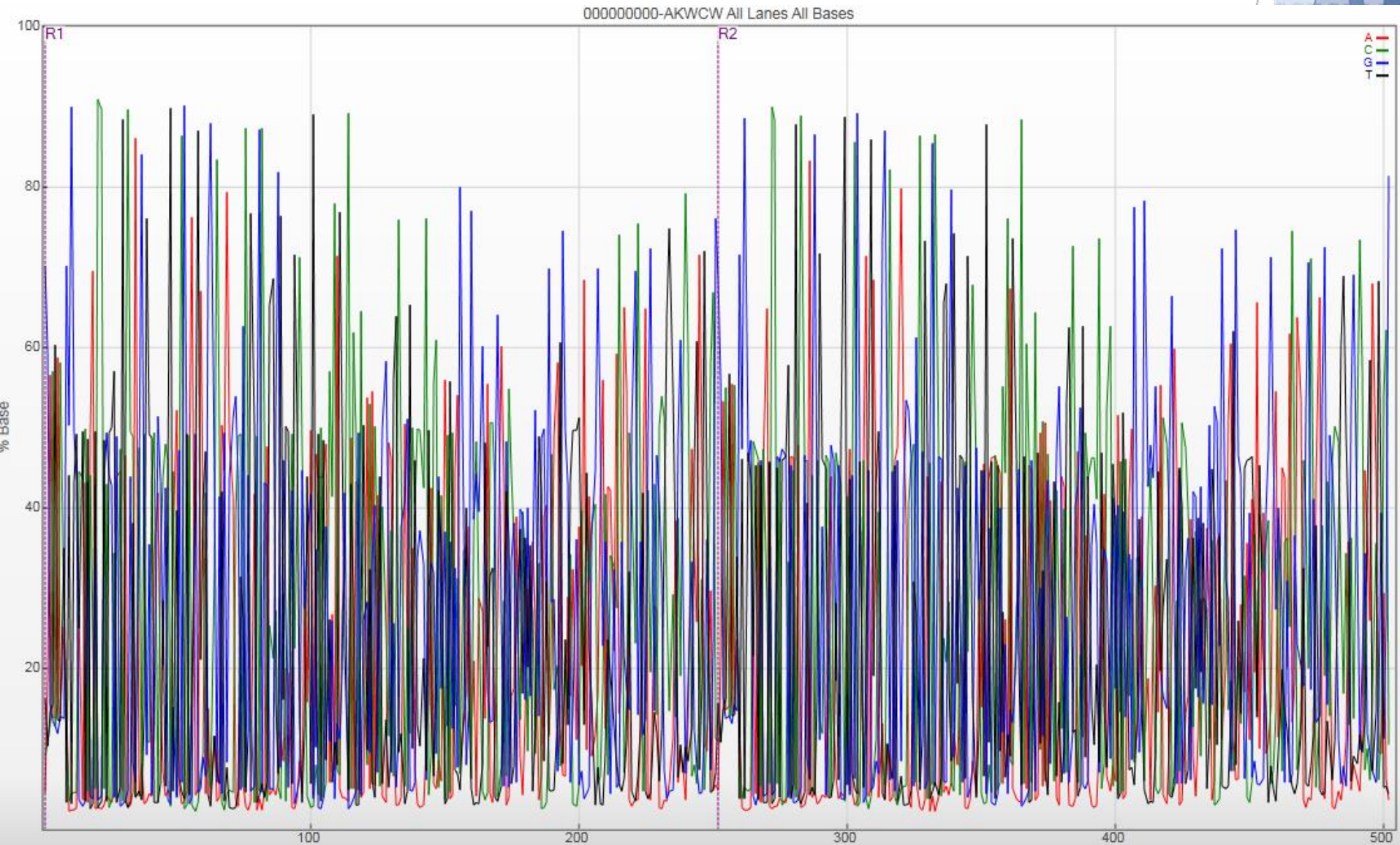
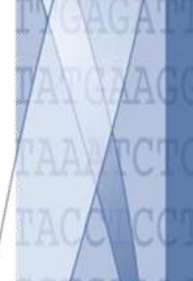


amplicon mix

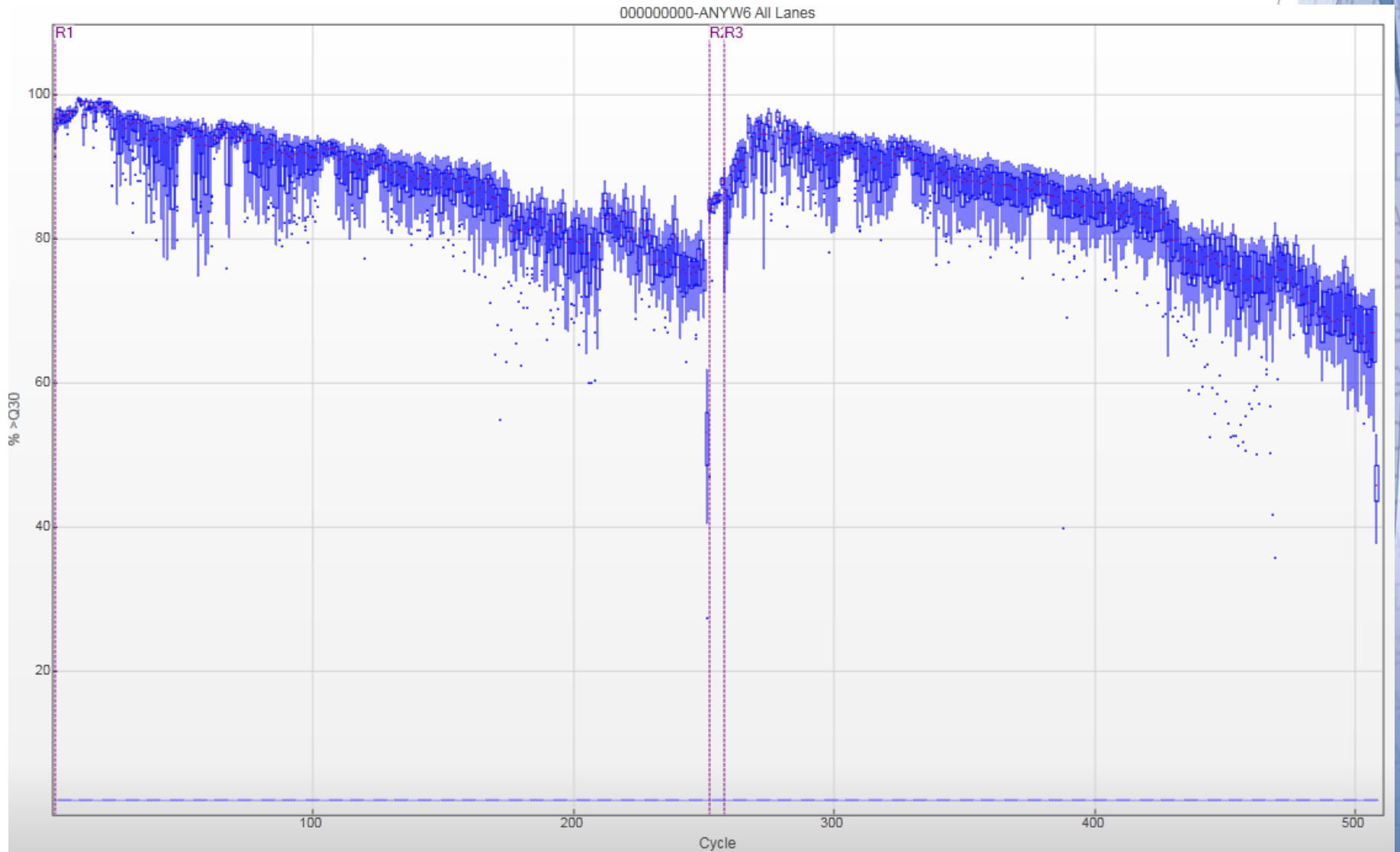
00000000-ANYW6 All Lanes All Bases



amplicon



amplicon mix Q30



FASTQC

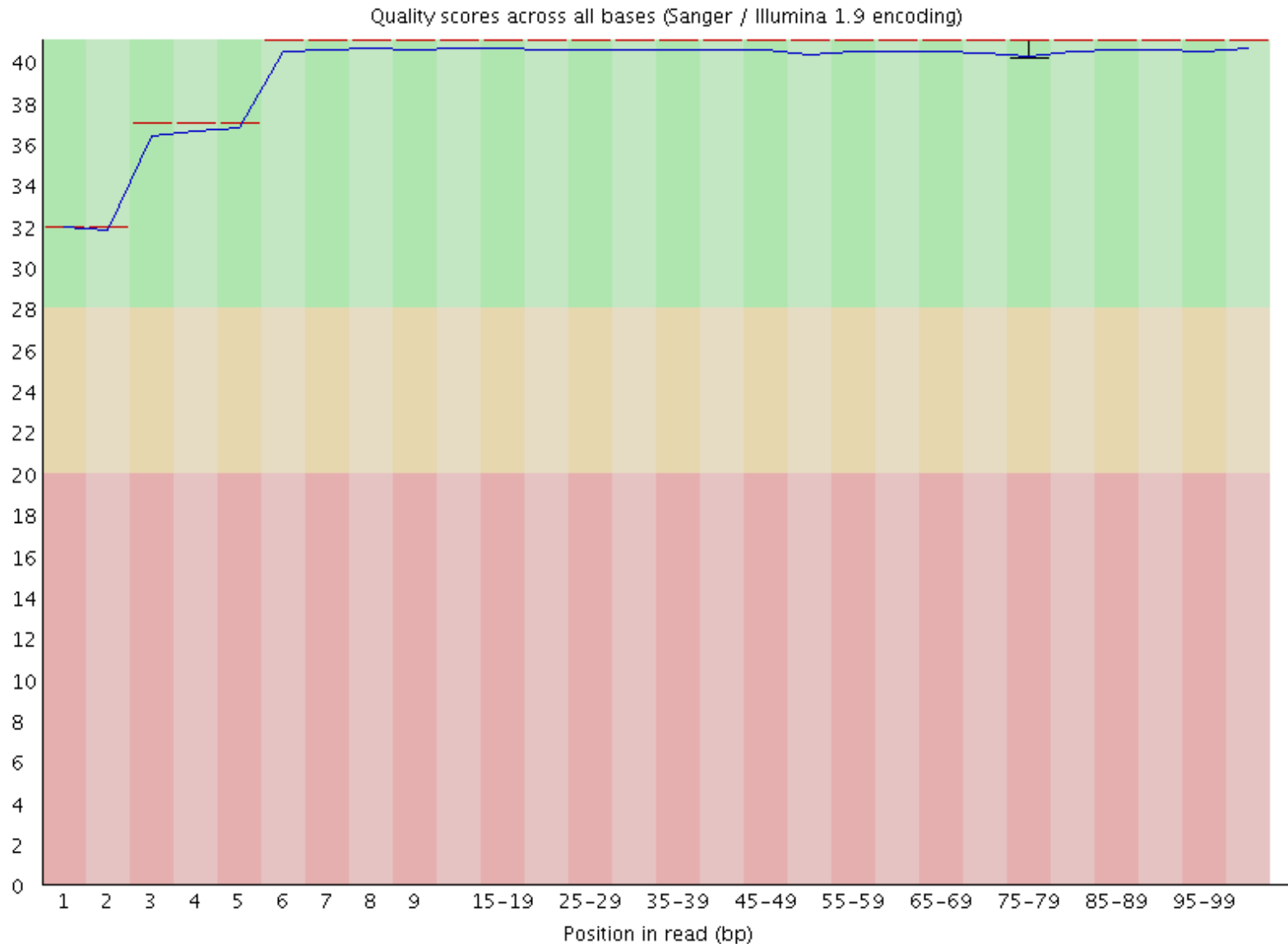


Basic Statistics

Measure	Value
Filename	3_S16_L008_R1_001.fastq.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	16574908
Sequences flagged as poor quality	0
Sequence length	150
%GC	40

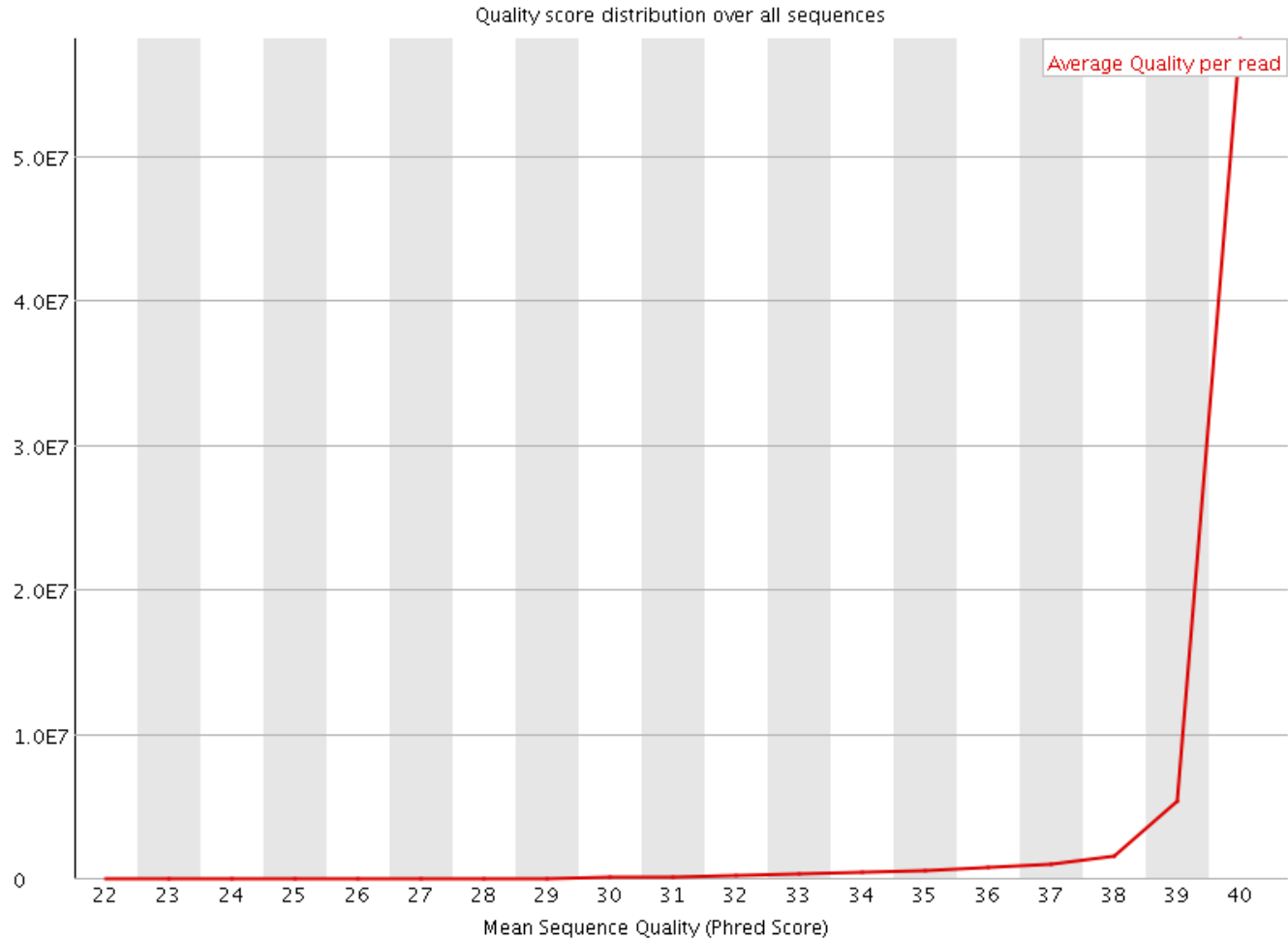


Per base sequence quality



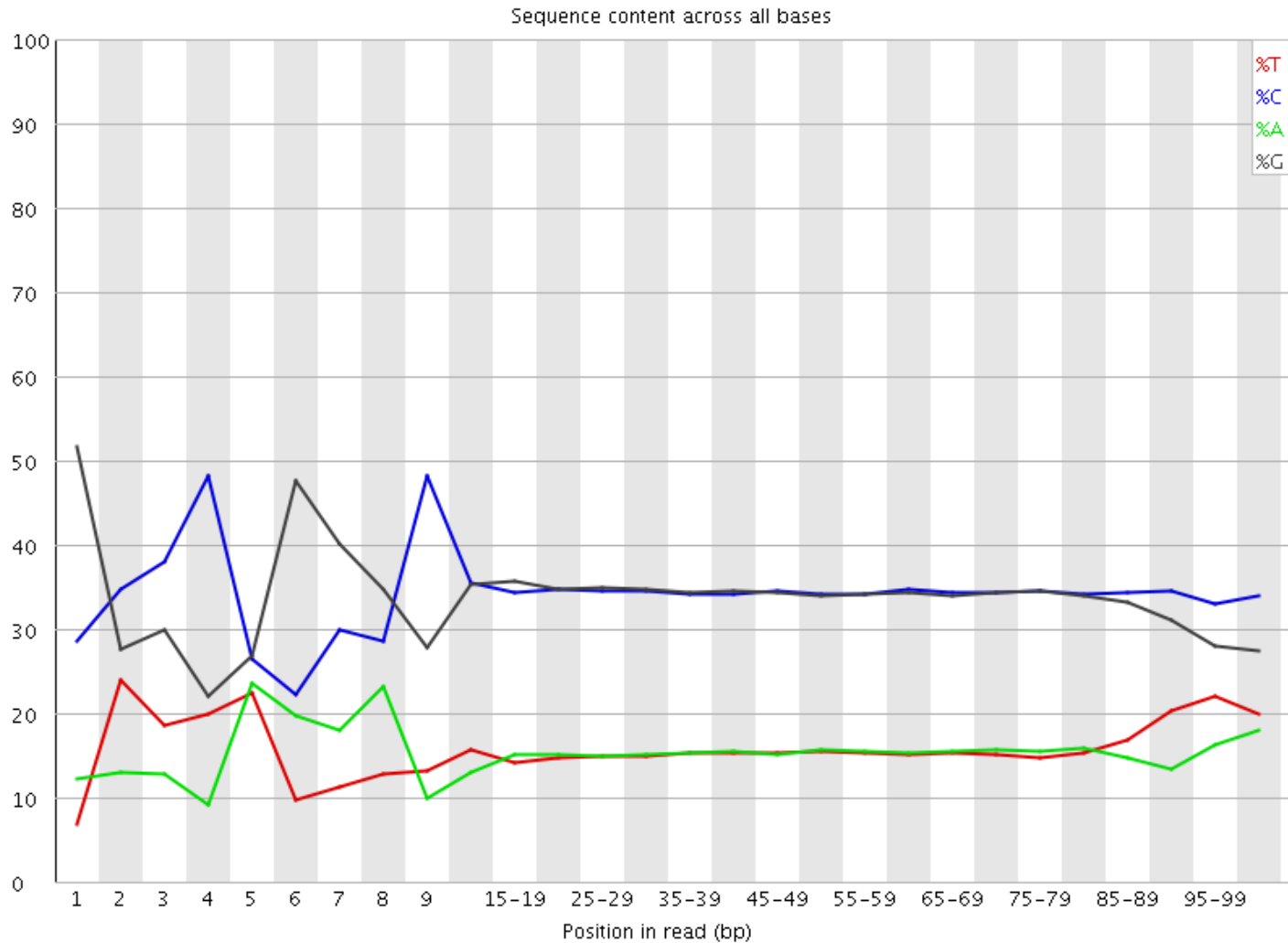
FASTQC

✔ Per sequence quality scores



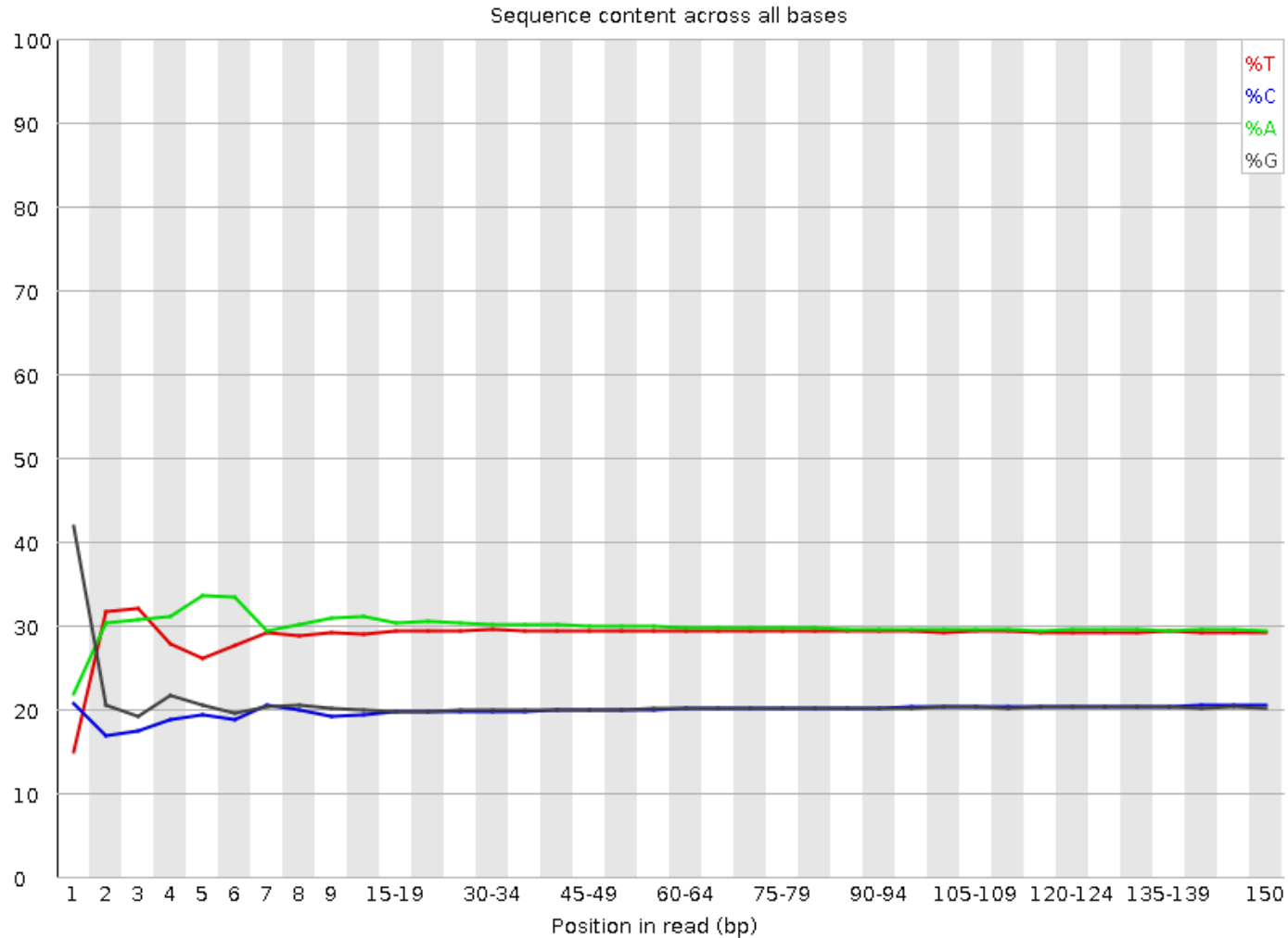
FASTQC - Nextera

✖ Per base sequence content



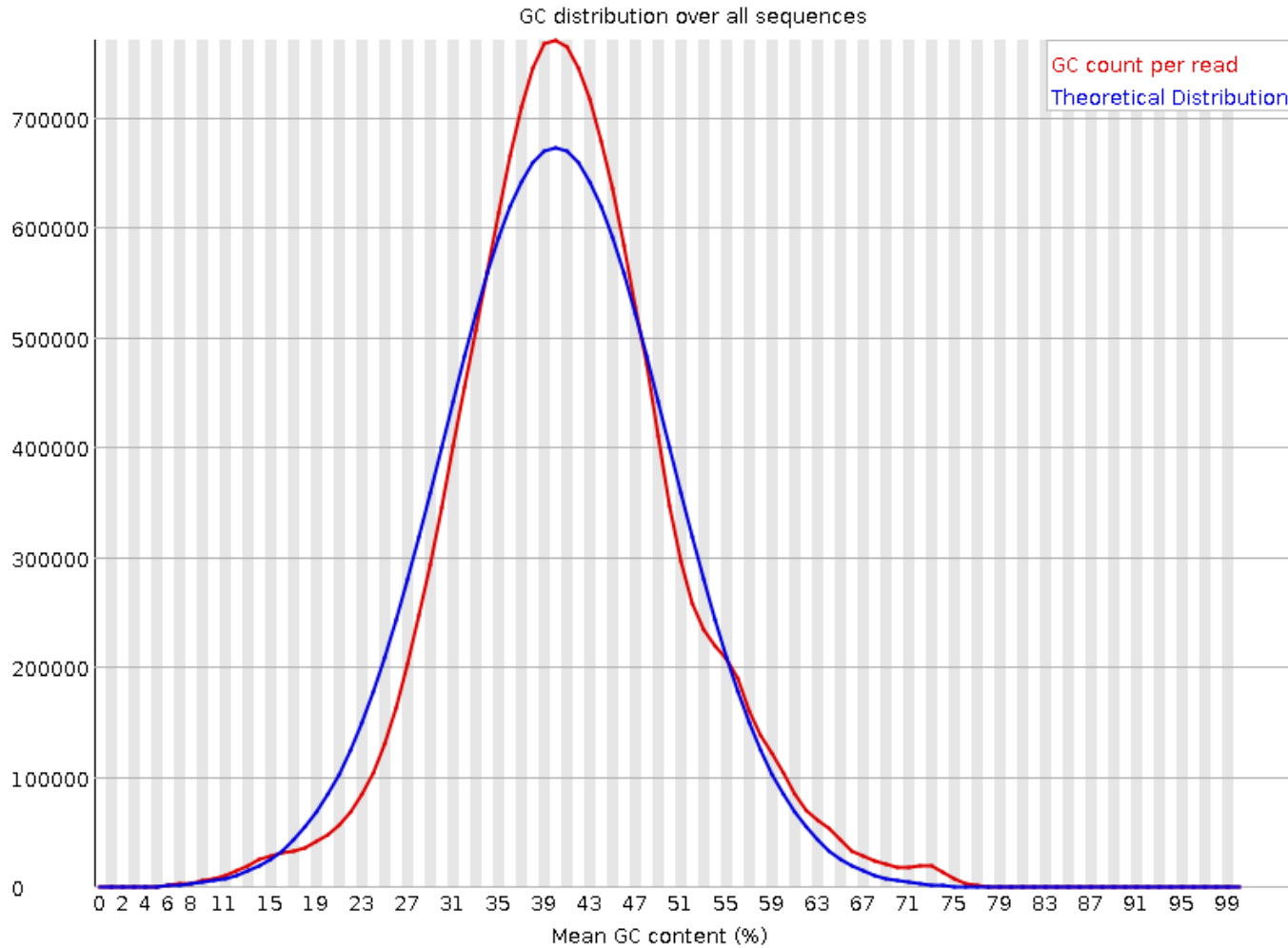
FASTQC

✖ Per base sequence content



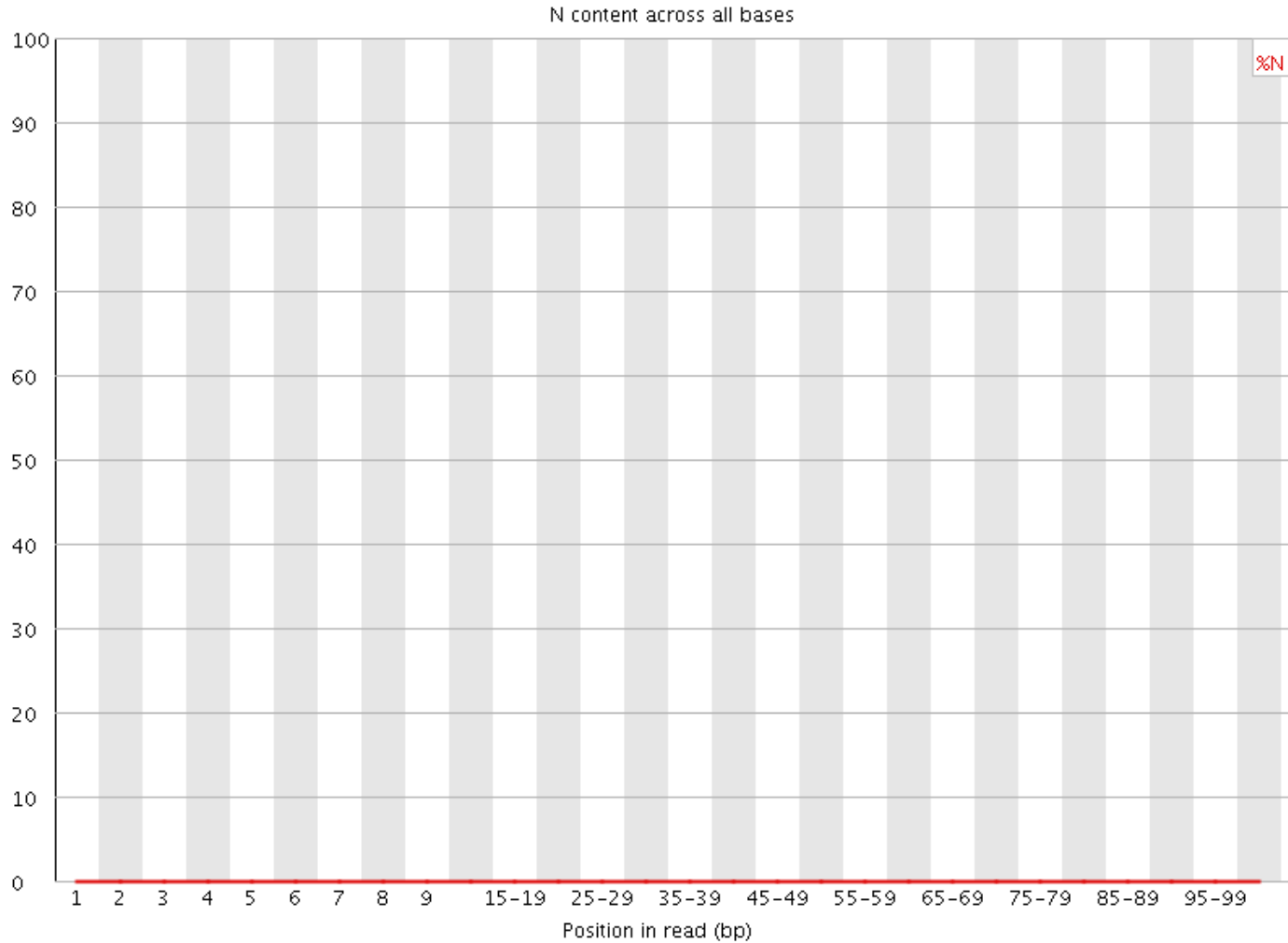
FASTQC

✔ Per sequence GC content



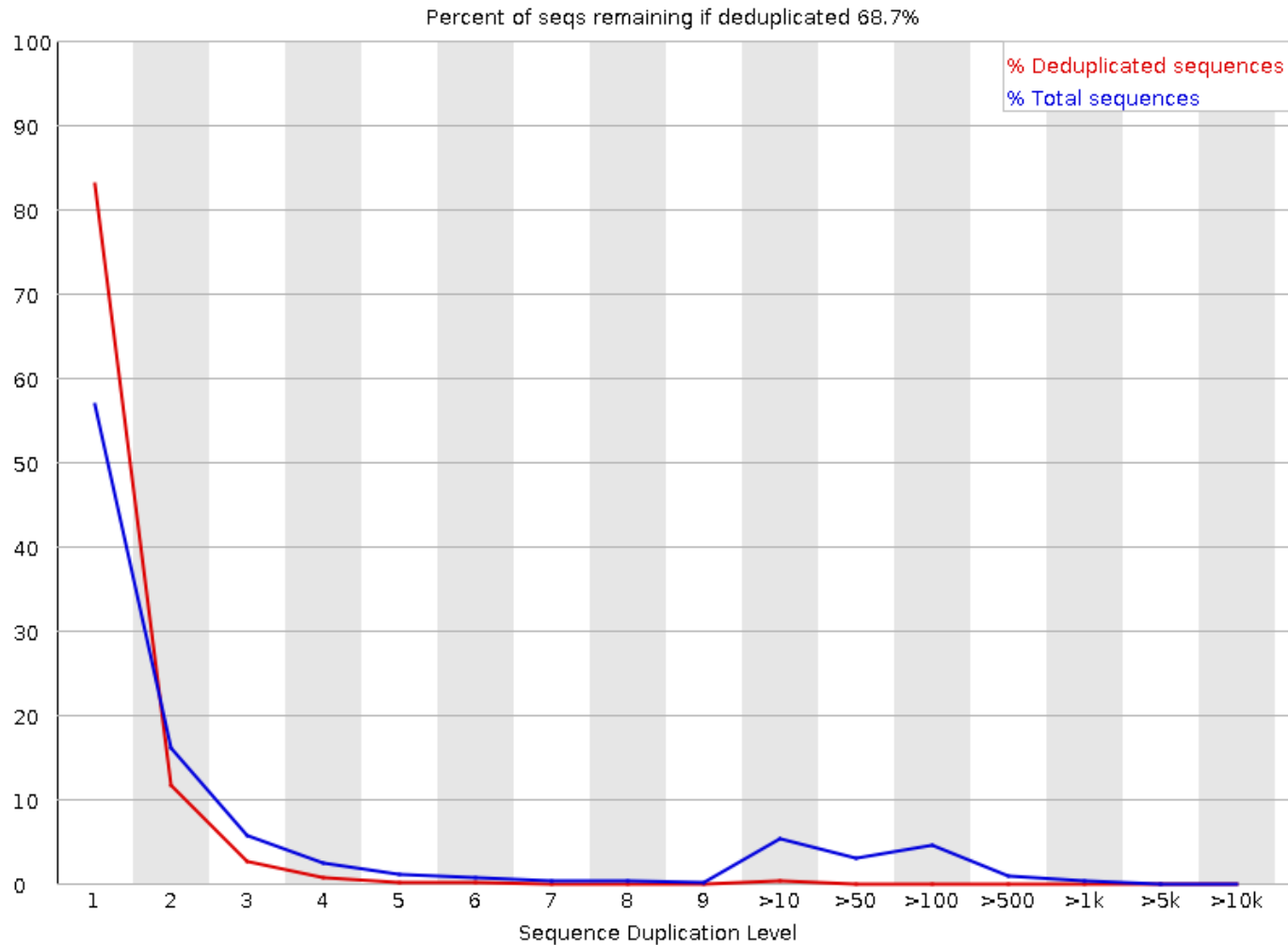
FASTQC

✔ Per base N content



FASTQC

⚠ Sequence Duplication Levels

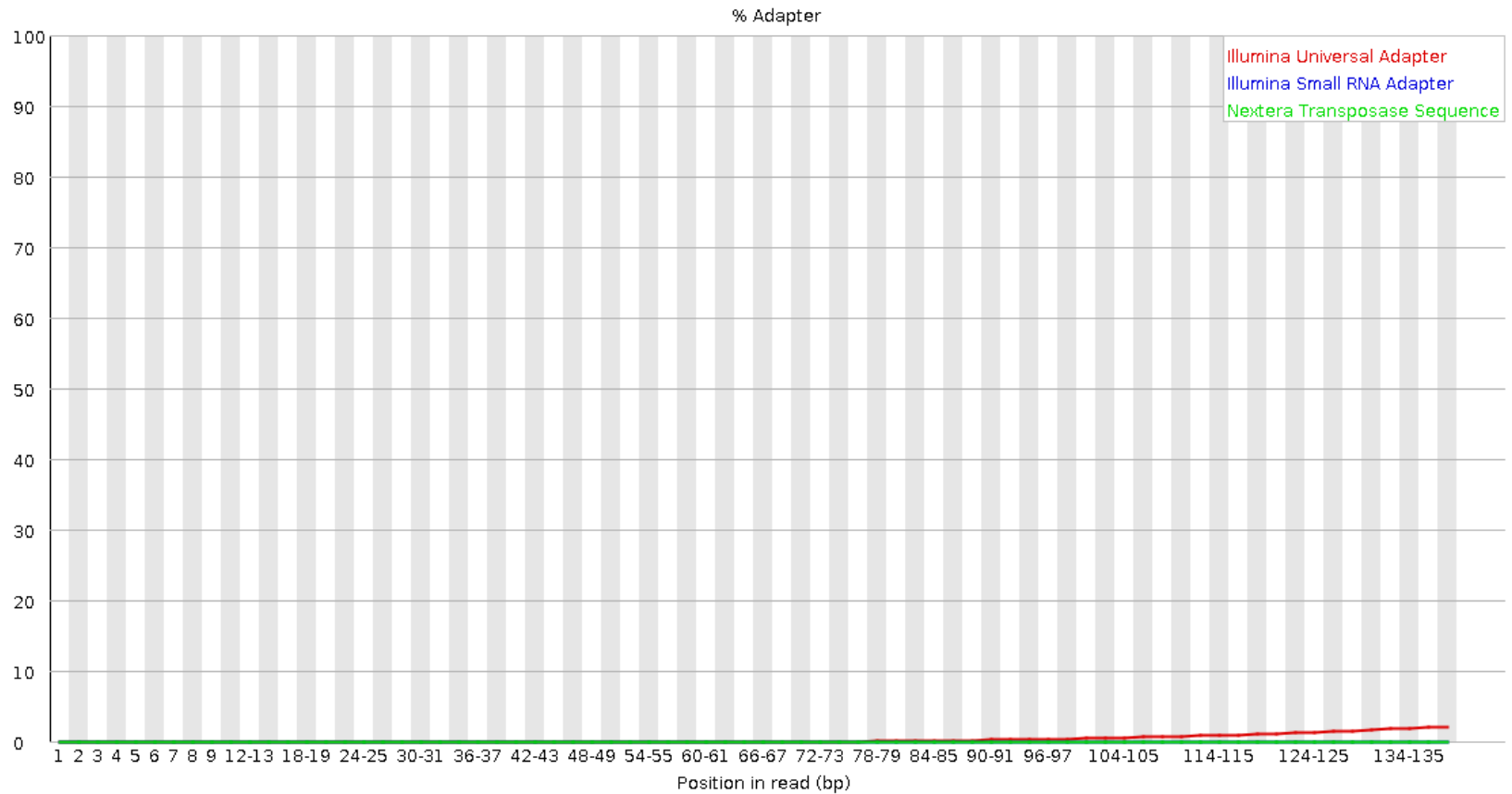


FASTQC

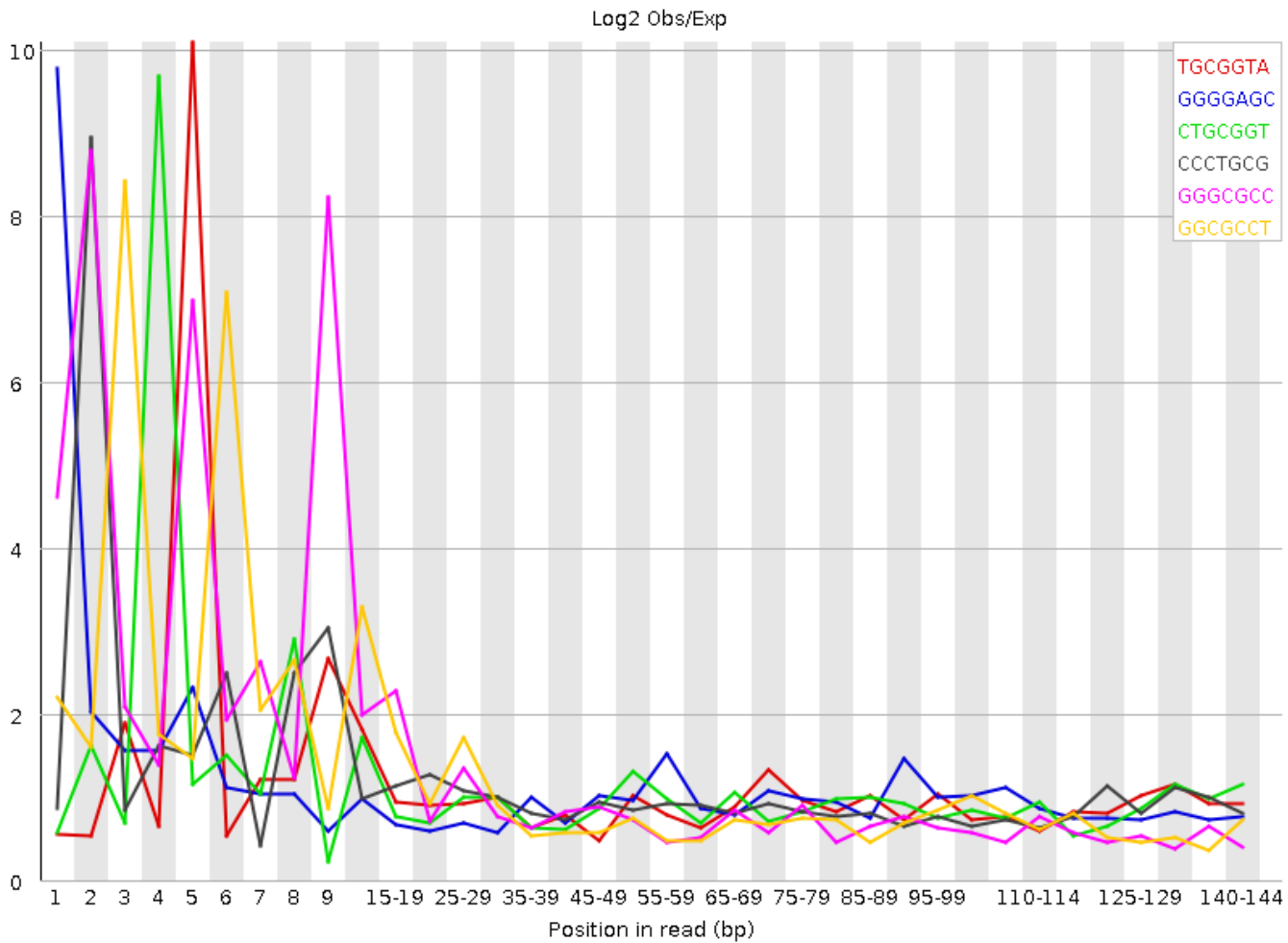
✔ Overrepresented sequences

No overrepresented sequences

✔ Adapter Content

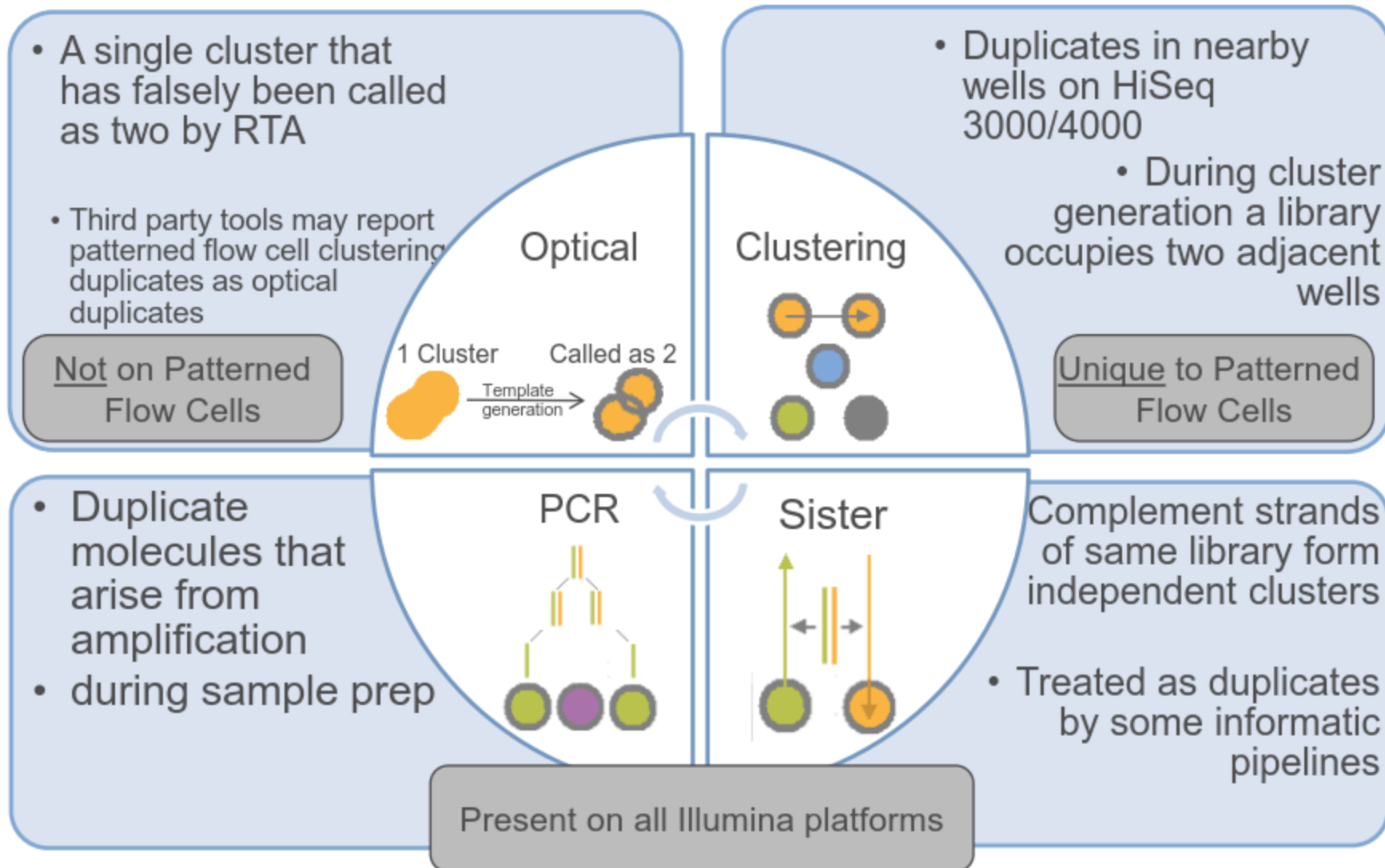


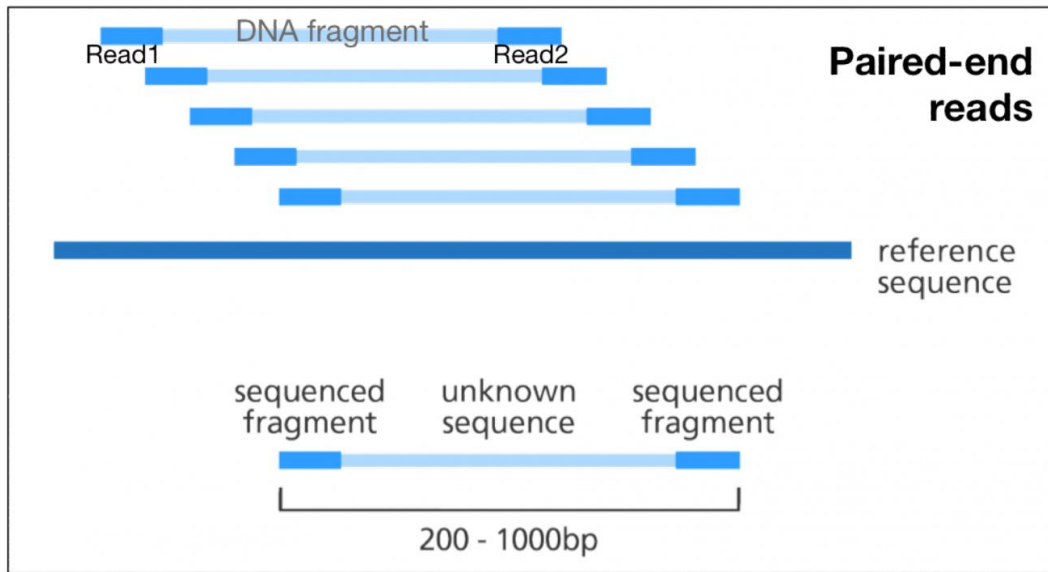
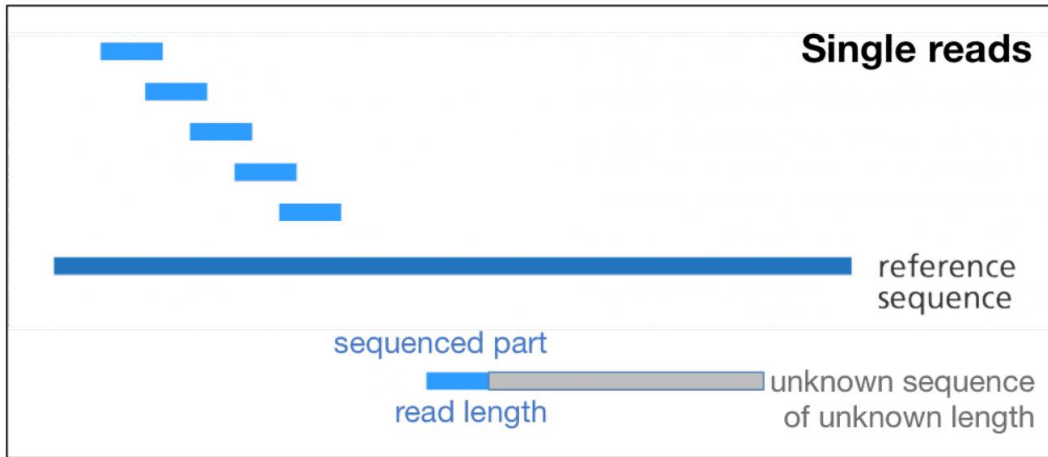
Kmer Content



Sequence	Count	PValue	Obs/Exp Max	Max Obs/Exp Position
TGGCGTA	6425	0.0	10.080686	5
GGGGAGC	9540	0.0	9.778594	1
CTGCGGT	6170	0.0	9.680999	4
CCCTGCG	6605	0.0	8.939233	2
GGGCGCC	5155	0.0	8.799765	2

A Review of Sequencing Duplicate Types





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

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

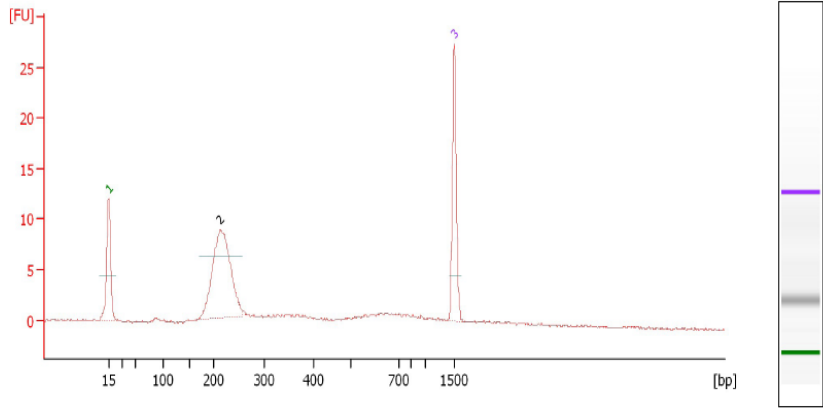
➤ qPCR

Most accurate quantitation method

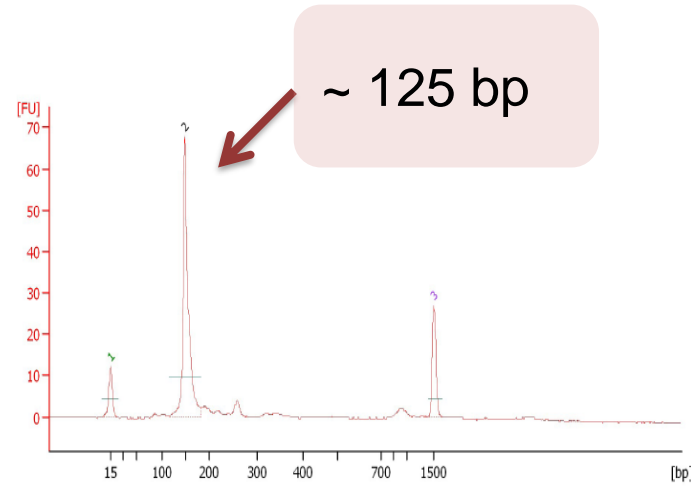
More labor-intensive

Must be compared to a control

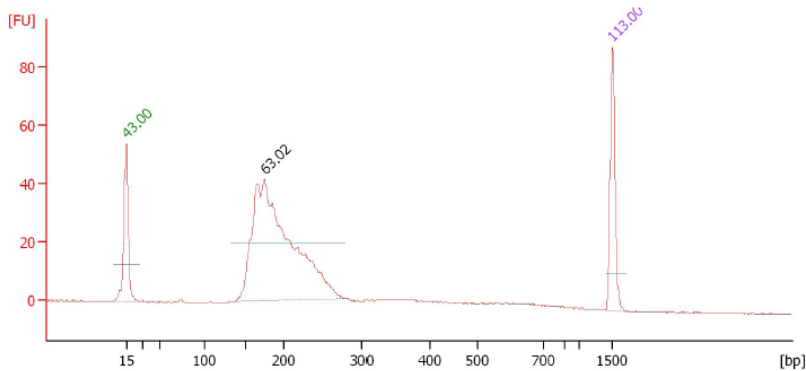
Library QC by Bioanalyzer



Beautiful

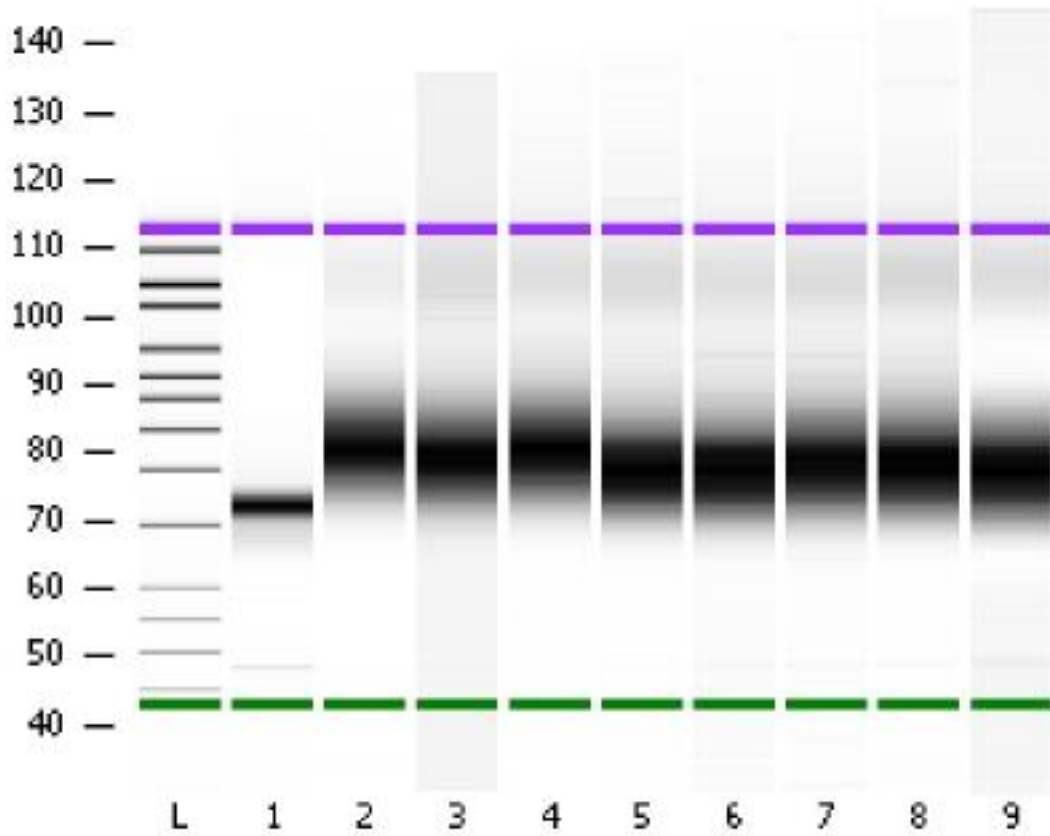


100% Adapters



Beautiful

Library QC



Examples for successful libraries



Adapter contamination at ~125 bp



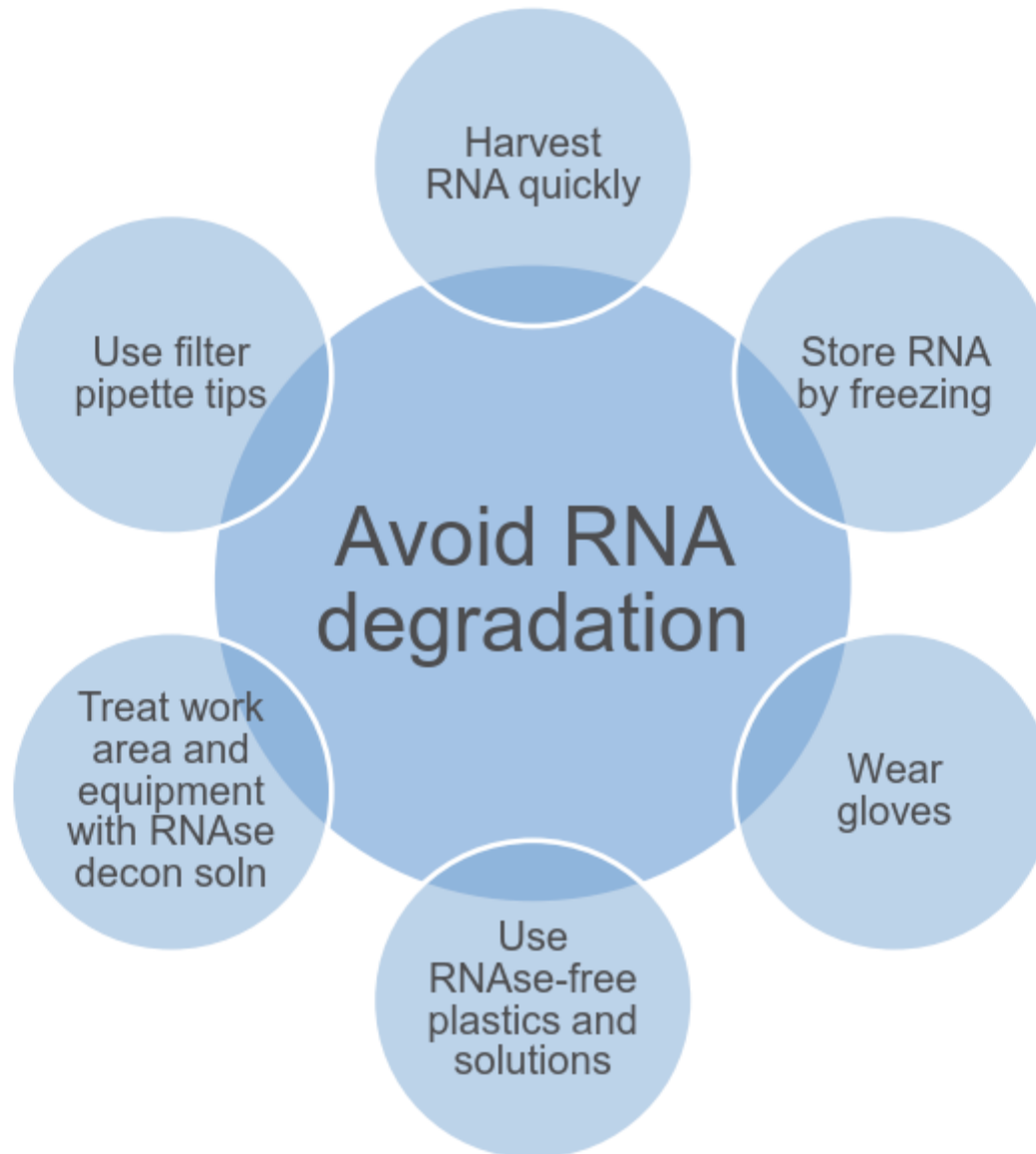
RNA is not that fragile



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

TTGAGAT
TATGAGC
TAAATCT
TACCACCT
GCTGAAAC
ATTCCCT
TCTGGGAA
GAAATTAT
TGTTGAA
AAGGAGC
TTTGGG
CGCCAGC
TCCCAGC
AATTGCAT
TCTCCAA
AAGGCTT
AATTGGA
GCACAA
ATACCA
GCTTTT
TTTATC

RNA Handling Best Practices



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 (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 (Clontech)	100 pg – 10 ng
Ovation RNA seq V2, Single Cell RNA seq (NuGen)	10 ng – 100 ng

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

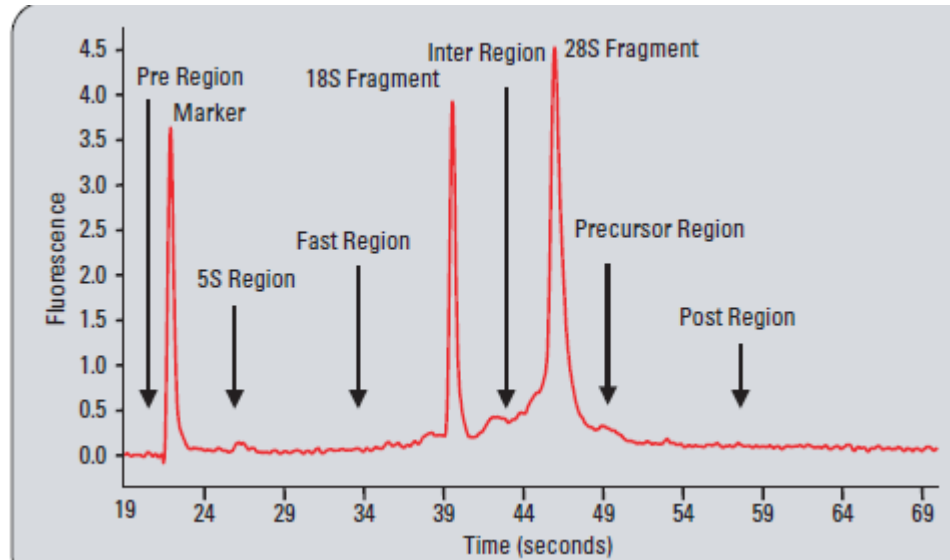
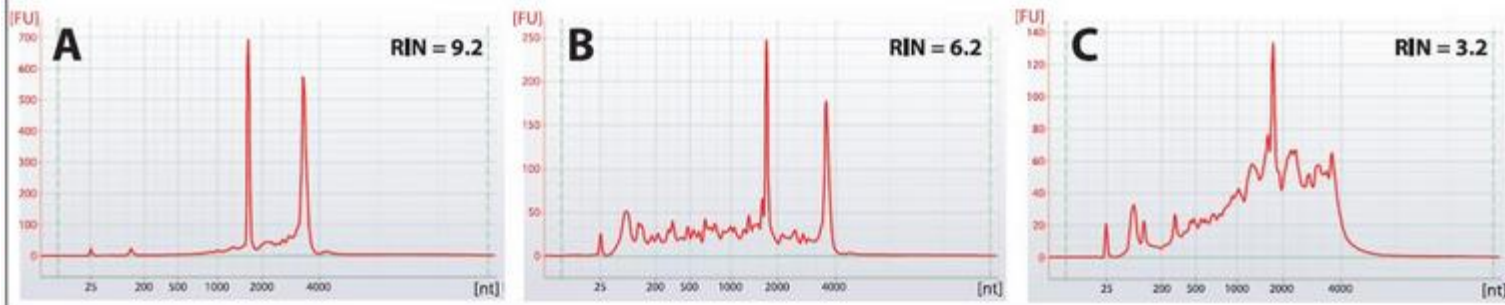
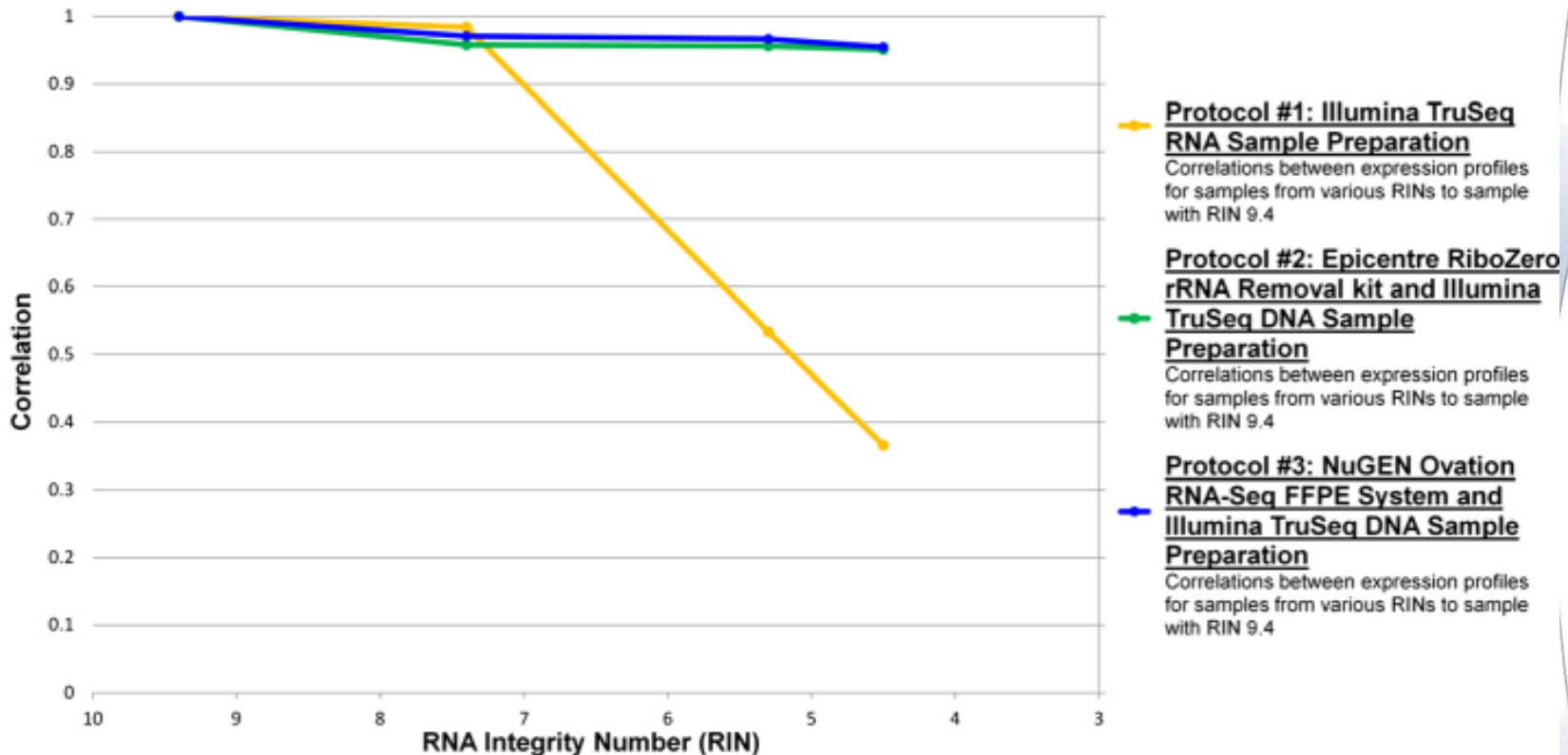


Figure 2.1 Example Agilent Bioanalyzer Electropherograms from three different total RNAs of varying integrity. Panel [A] represents a highly intact total RNA (RIN = 9.2), panel [B] represents a moderately intact total RNA (RIN = 6.2), and panel [C] represents a degraded total RNA sample (RIN = 3.2).



RNA integrity <> reproducibility



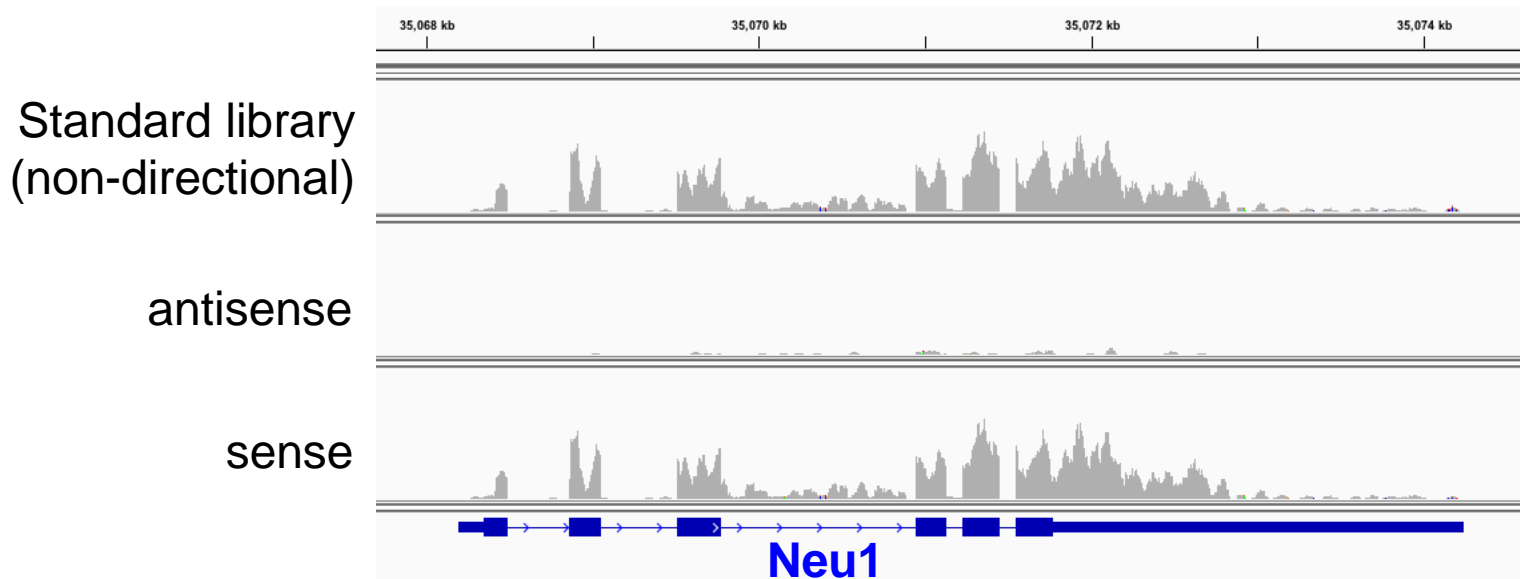
Chen et al. 2014

Considerations in choosing an RNA-Seq method

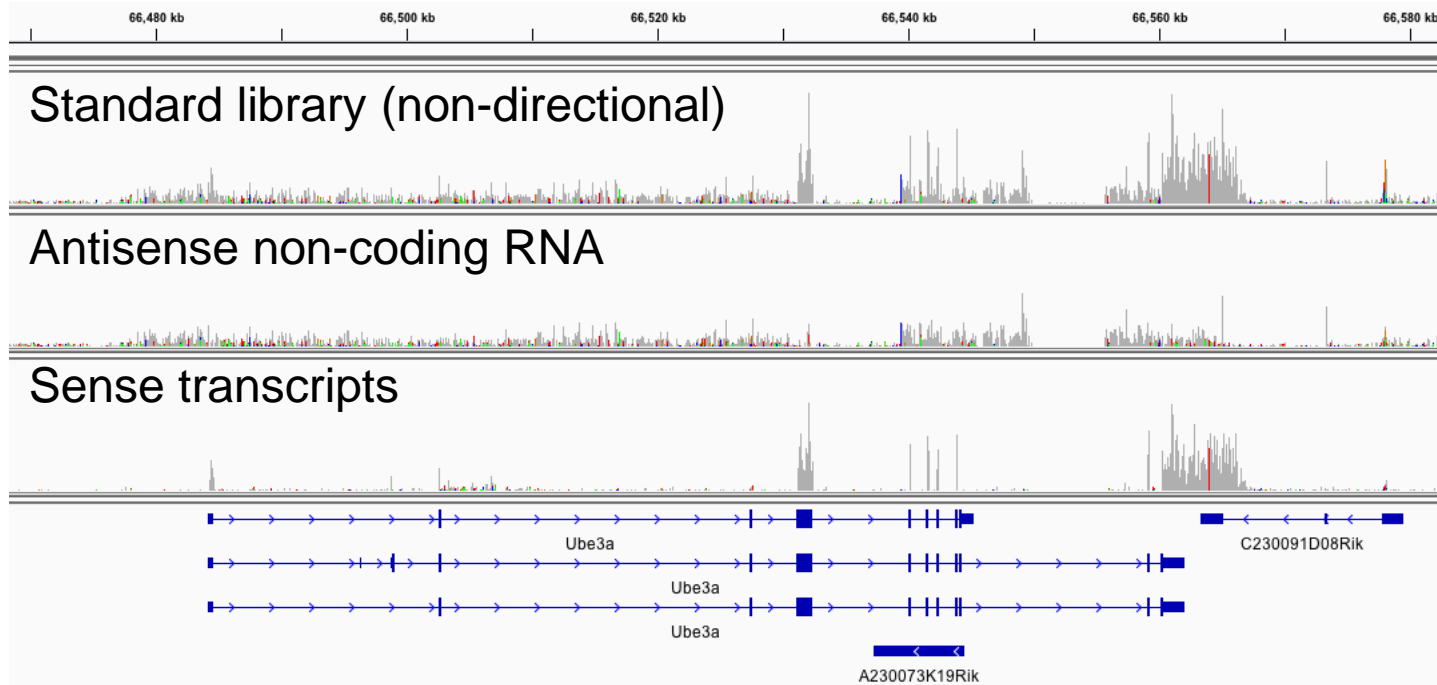
- Transcript type:
 - mRNA, extent of degradation
 - small/micro RNA
- Strandedness:
 - un-directional ds cDNA library
 - directional library
- Input RNA amount:
 - 0.1-4ug original total RNA
 - linear amplification from 0.5-10ng RNA
- Complexity:
 - original abundance
 - cDNA normalization for uniformity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



Is strand-specific information important?



Strand-specific RNA-seq



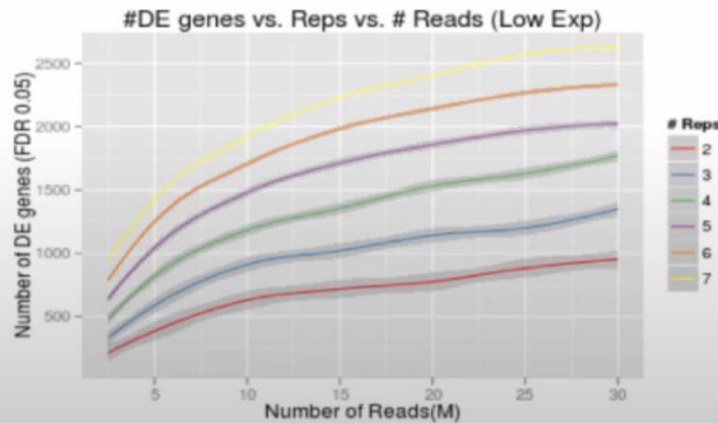
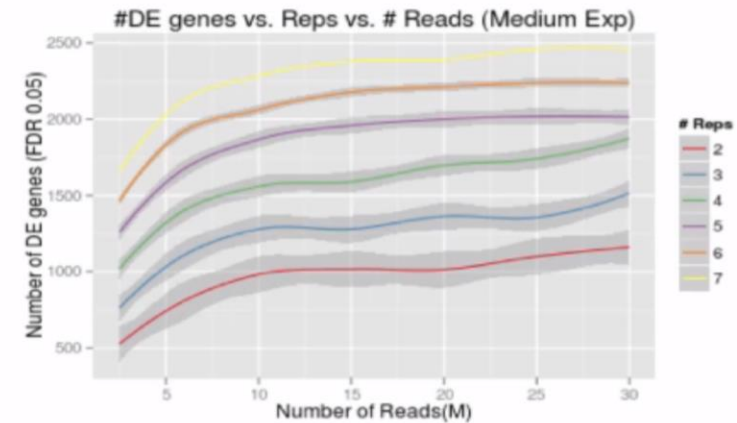
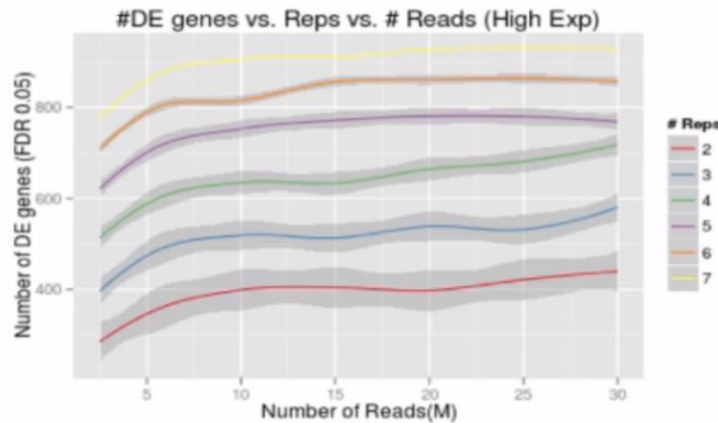
- Informative for non-coding RNAs and antisense transcripts
- Essential when NOT using polyA selection (mRNA)
- No disadvantage to preserving strand specificity

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.

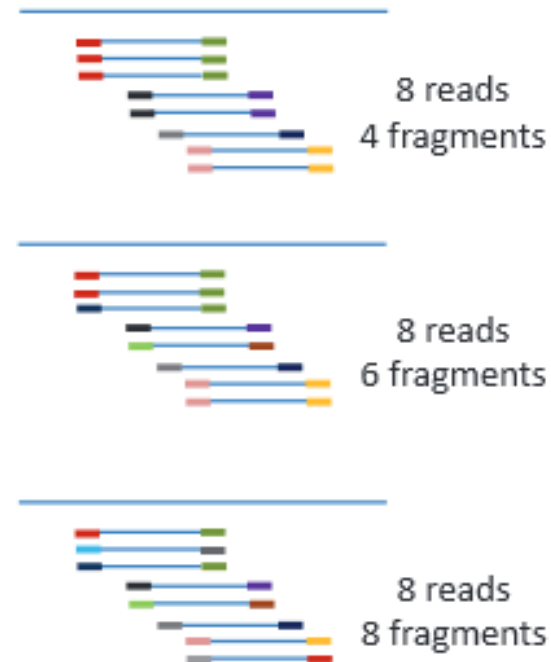
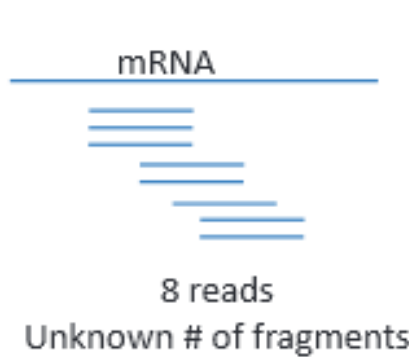
Liu et al. (2014) RNA-Seq differential expression studies: more sequence or more replication?, Bioinformatics, 30(3):1-4

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?

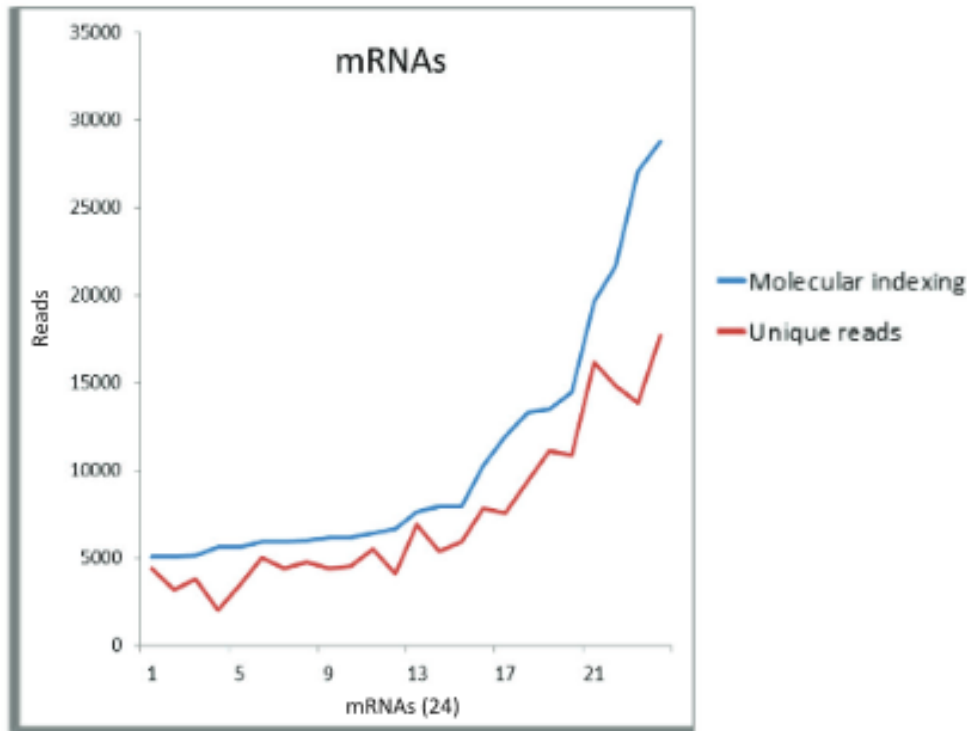
Molecular indexing – for precision counts

Conventional RNA-Seq
Without Molecular Indexing



Molecular indexing – for precision counts

B



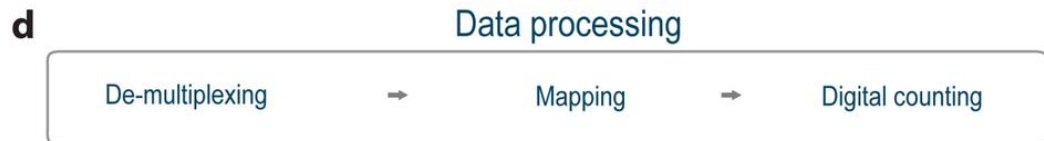
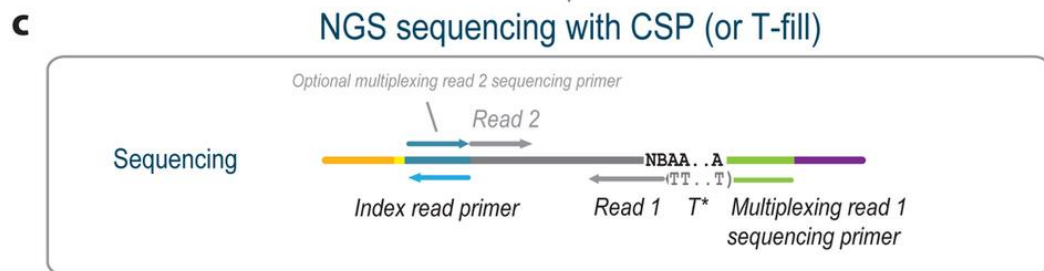
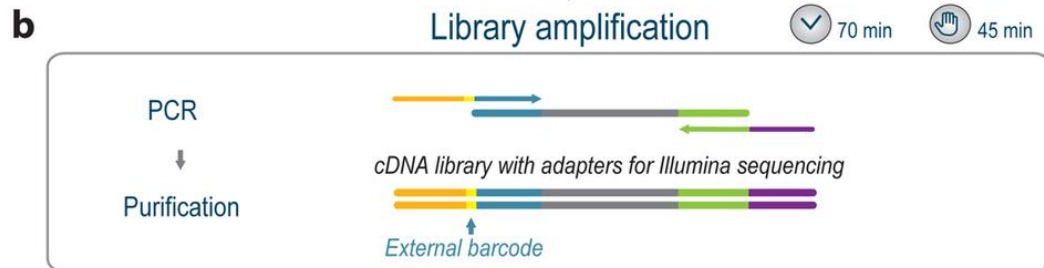
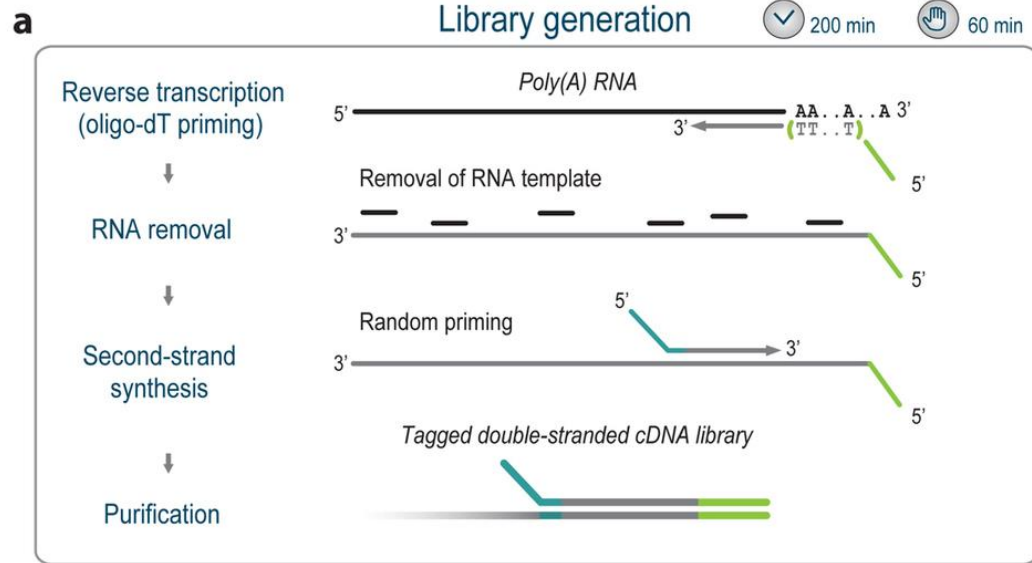
3'-Tag-Seq

- 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, m⁶A, m⁵C, pseudouridine, ...)
- RNA degradation pathways
- (differential translation into proteins)



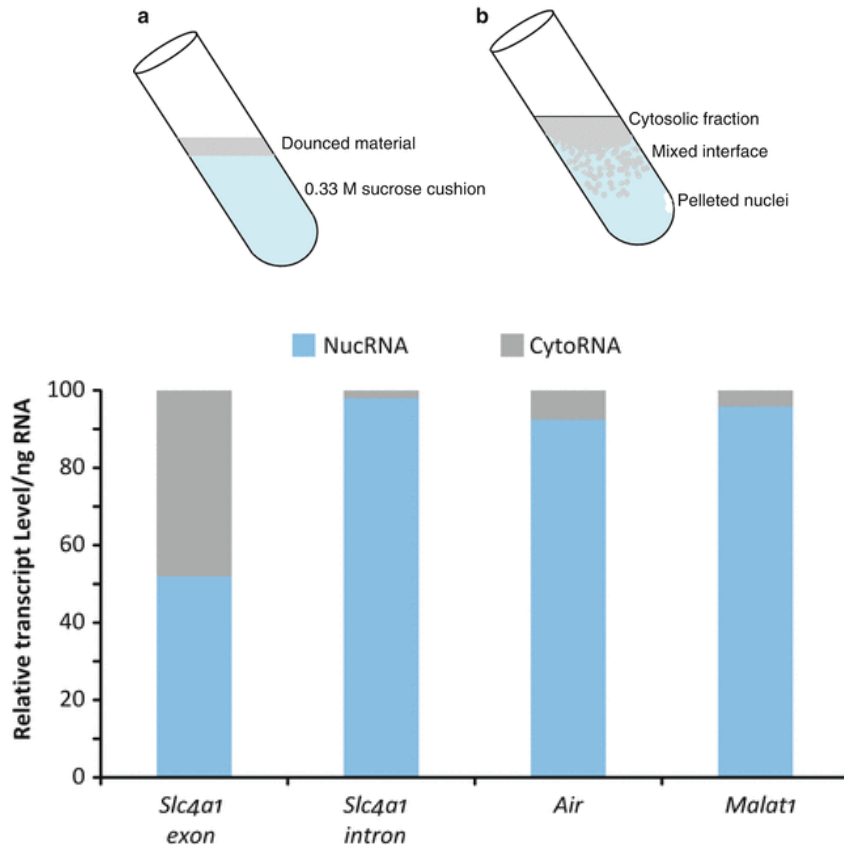
beyond steady-state RNA-seq

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



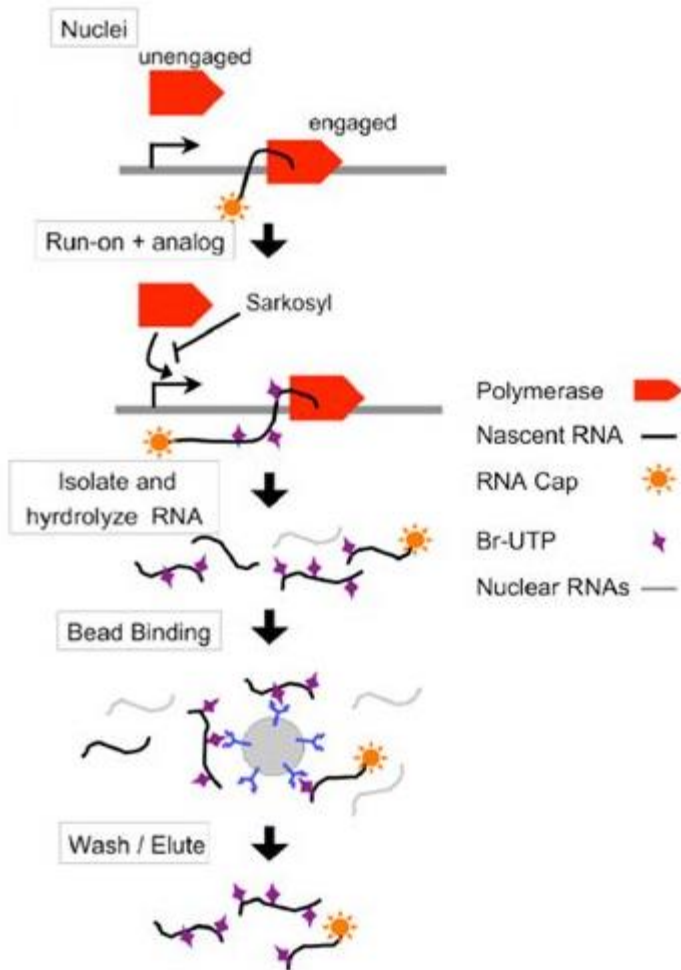
nucRNA-seq

- Fractioning of nuclei and cytosol
- Studying active transcription



Dhaliwal et al. 2016

GRO-Seq



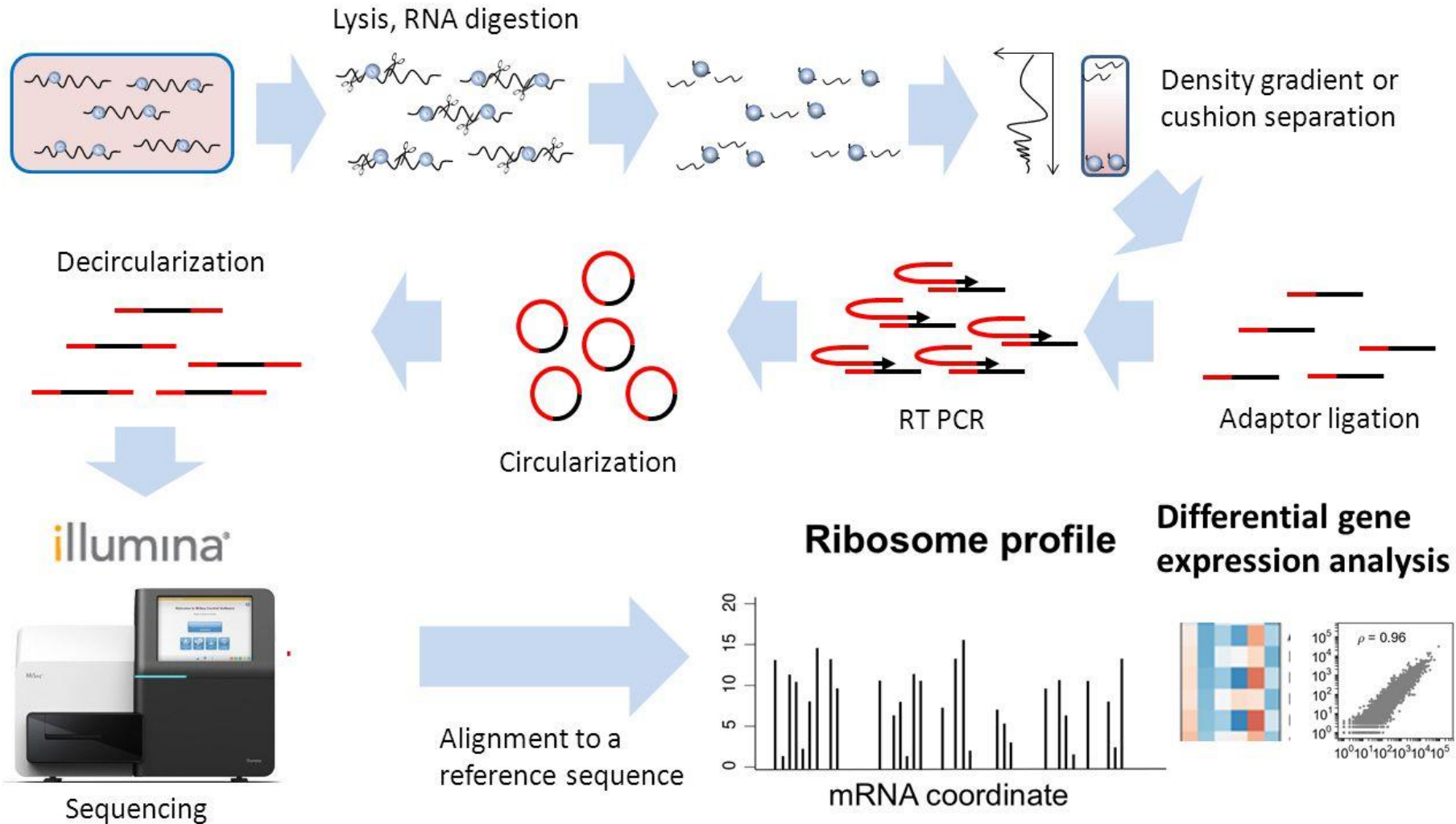
- Global Run-On – sequencing
- pulse-chase experiments (Br-UTP)
- uses isolated nuclei
- sarcosyl prevents binding of polymerase (only transcription in progress will be seq.)
- measures active transcription rather than steady state
- Maps position and orientation
- Earliest changes identify primary targets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



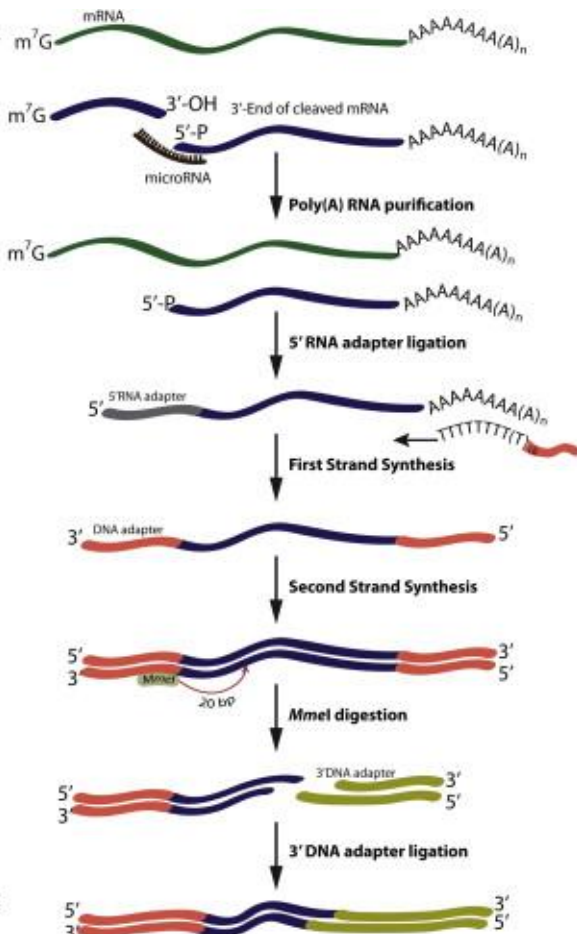
Degradome Sequencing

PARE-Seq

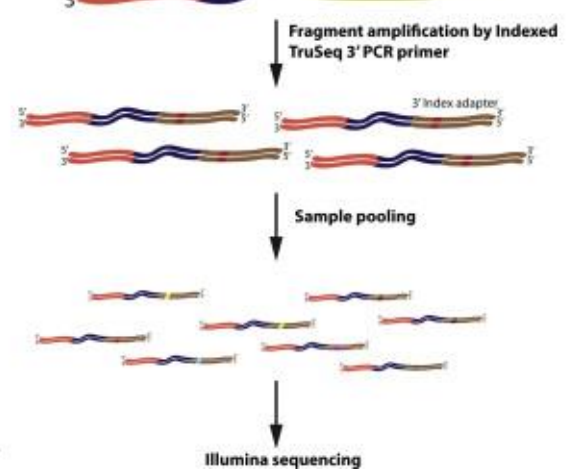
(Parallel Analysis of RNA Ends)

Zhai et al . 2013

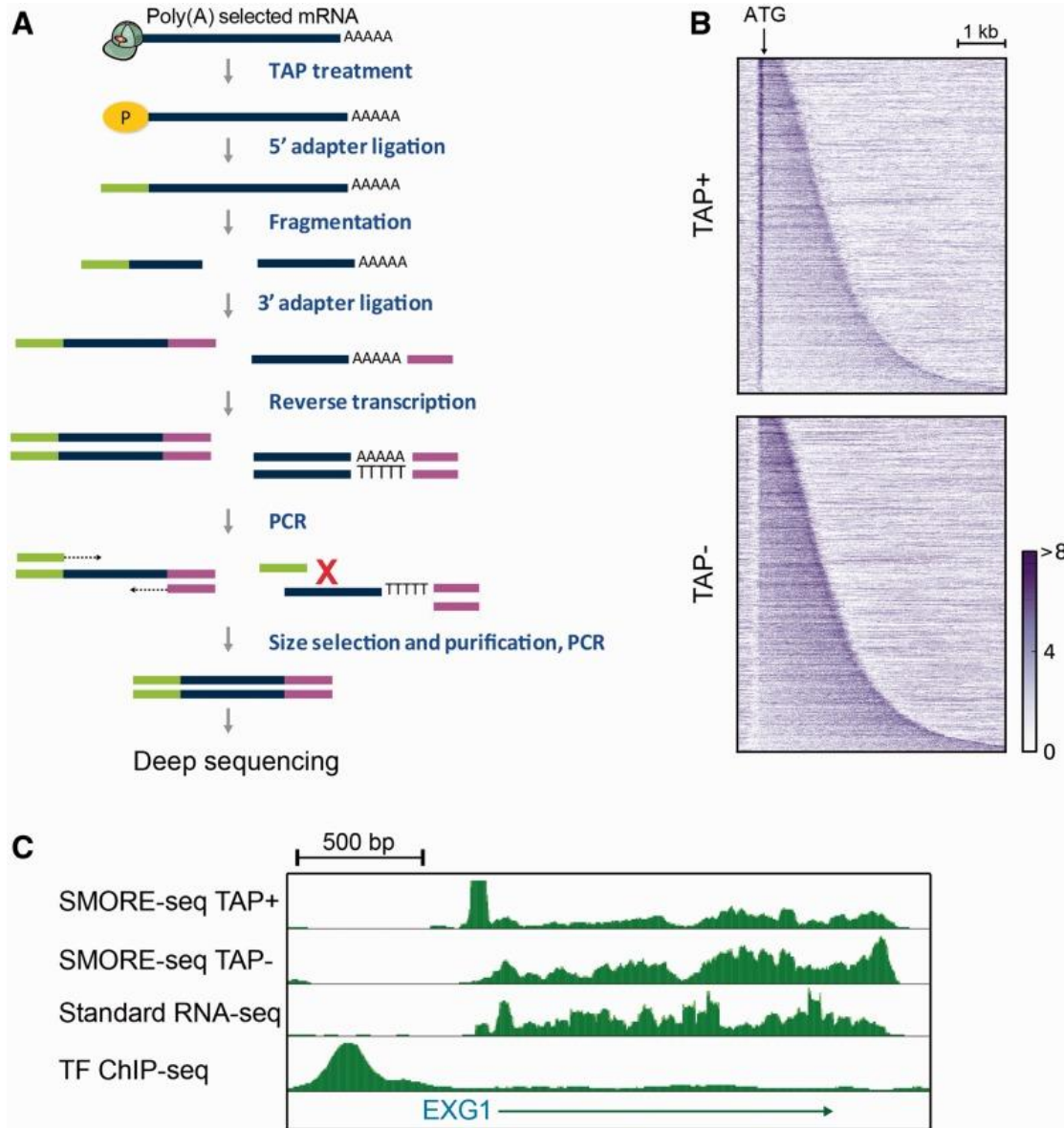
Day 1



Day 2-3



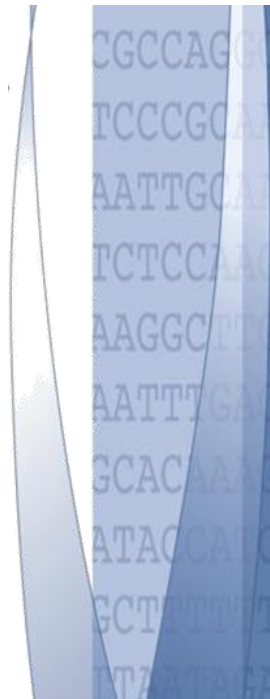
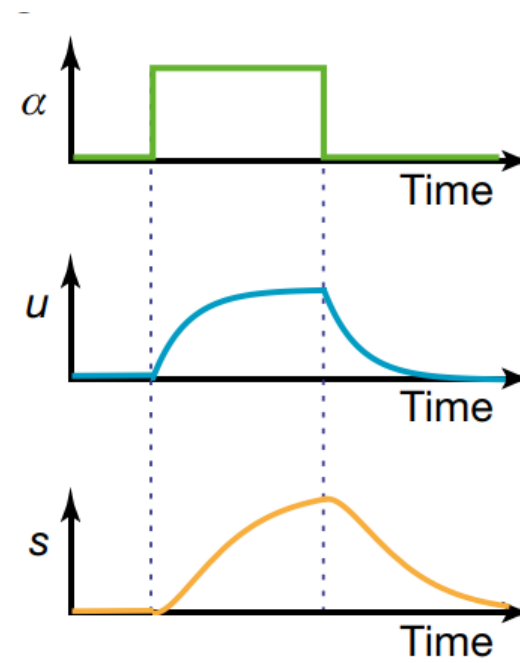
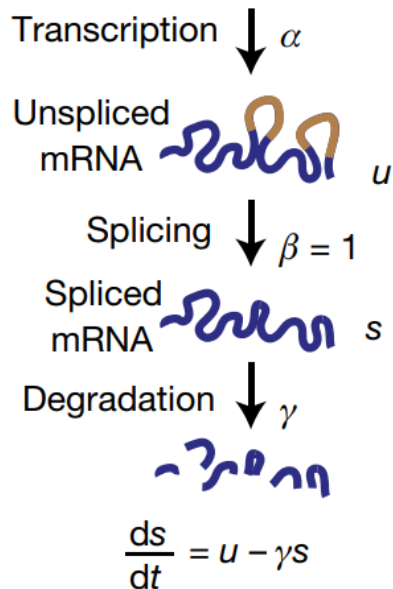
Degradome Sequencing

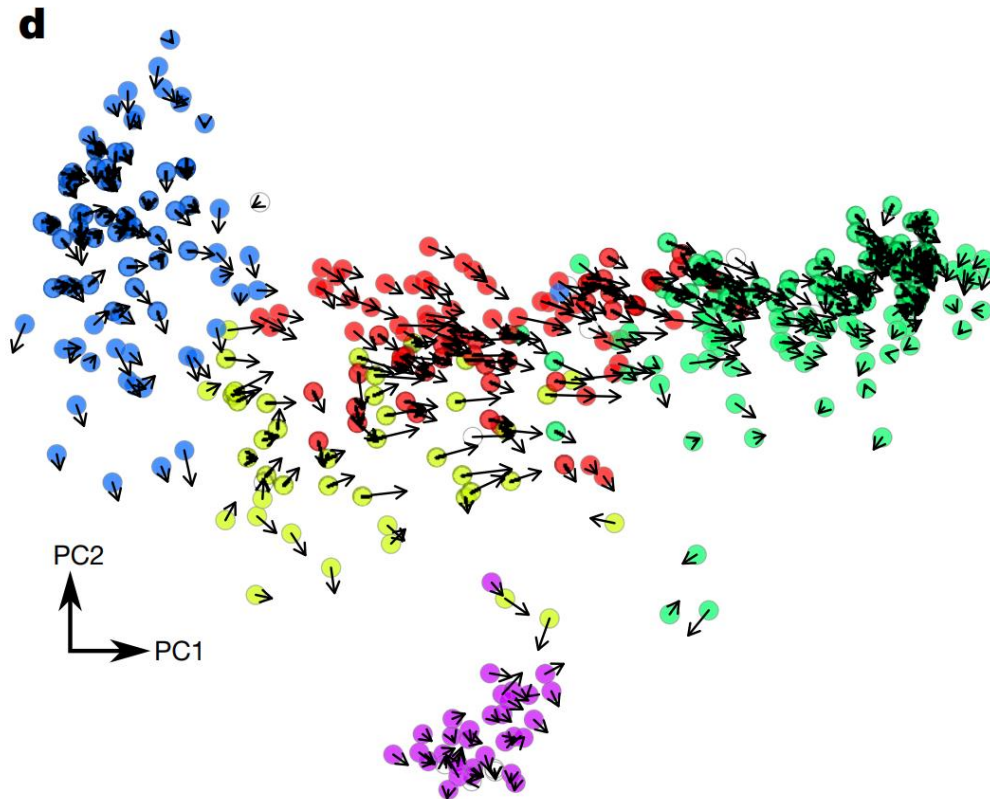
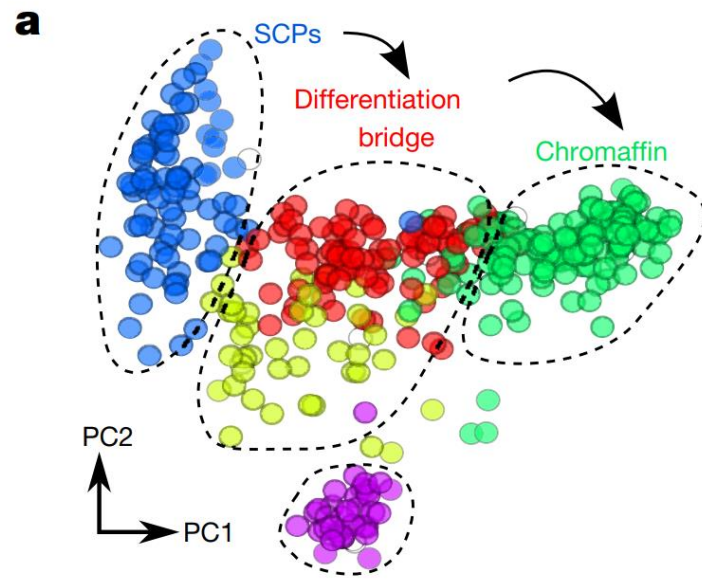


SMORE-Seq

RNA velocity of single cells

Gioele La Manno^{1,2}, Ruslan Soldatov³, Amit Zeisel^{1,2}, Emelie Braun^{1,2}, Hannah Hochgerner^{1,2}, Viktor Petukhov^{3,4}, Katja Lidschreiber⁵, Maria E. Kastriiti⁶, Peter Lönnerberg^{1,2}, Alessandro Furlan¹, Jean Fan³, Lars E. Borm^{1,2}, Zehua Liu³, David van Bruggen¹, Jimin Guo³, Xiaoling He⁷, Roger Barker⁷, Erik Sundström⁸, Gonçalo Castelo-Branco¹, Patrick Cramer^{5,9}, Igor Adameyko⁶, Sten Linnarsson^{1,2*} & Peter V. Kharchenko^{3,10*}





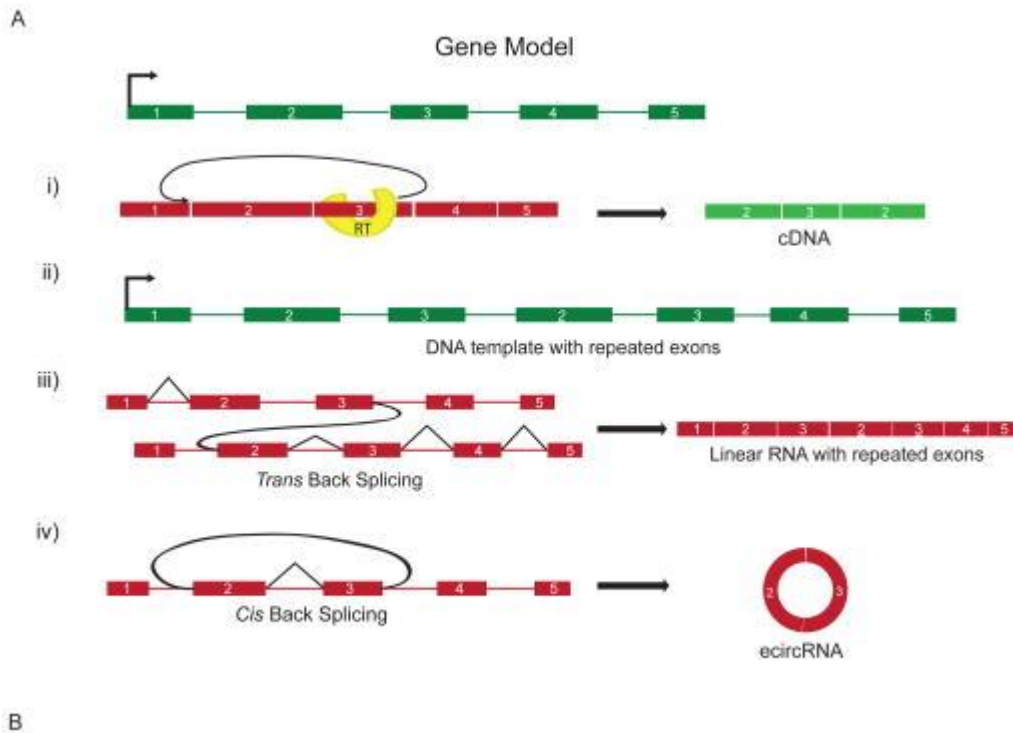
Circular RNA (circRNA)

- 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 ????



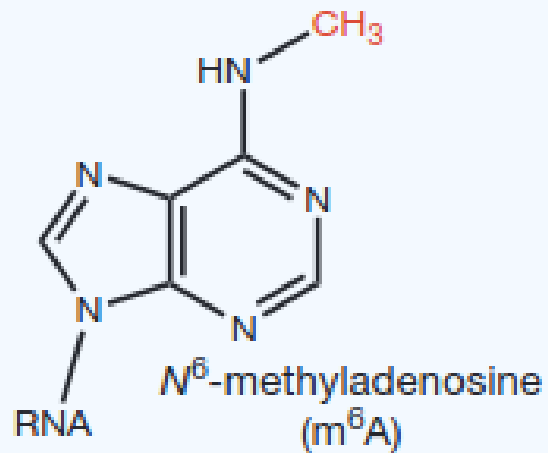
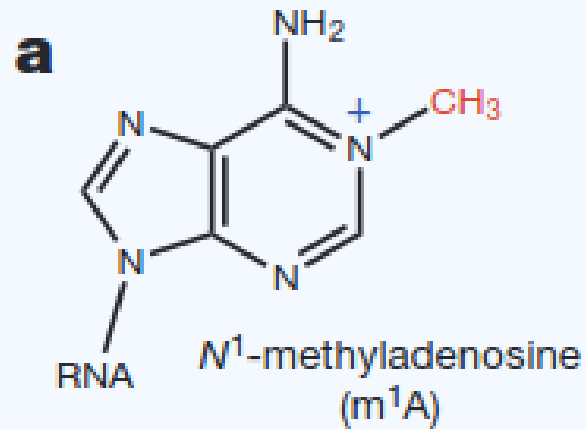
Role of circRNAs ?

Back-splicing and other mechanisms

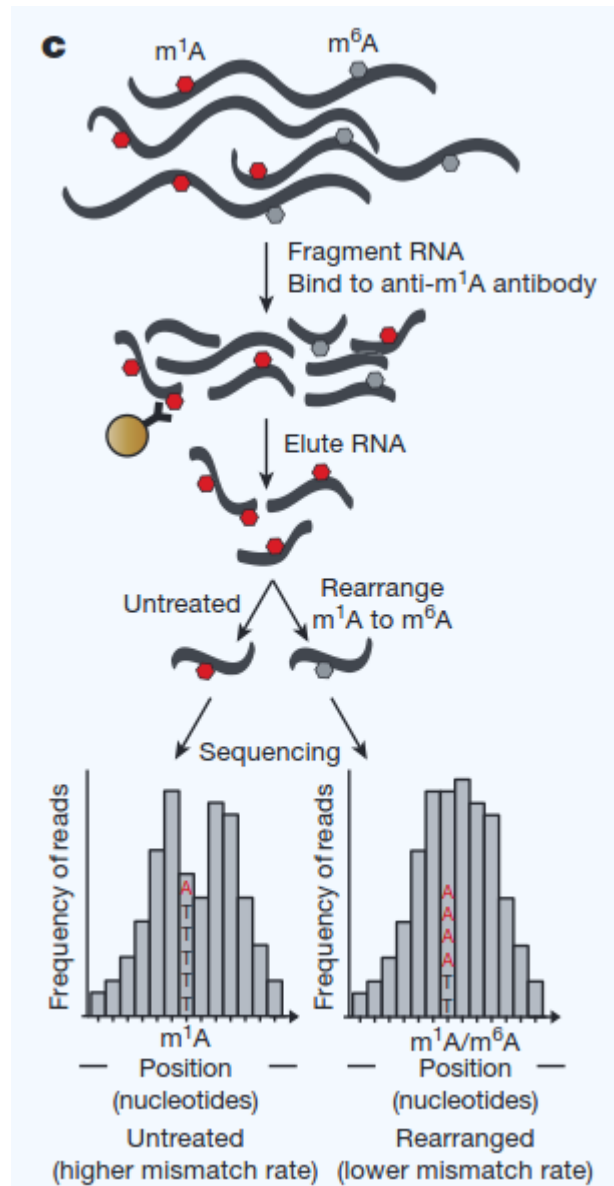


- miRNA sponge
- protein expression regulators:
mRNA traps
(blocking translation)
- Interactions with RNA binding proteins

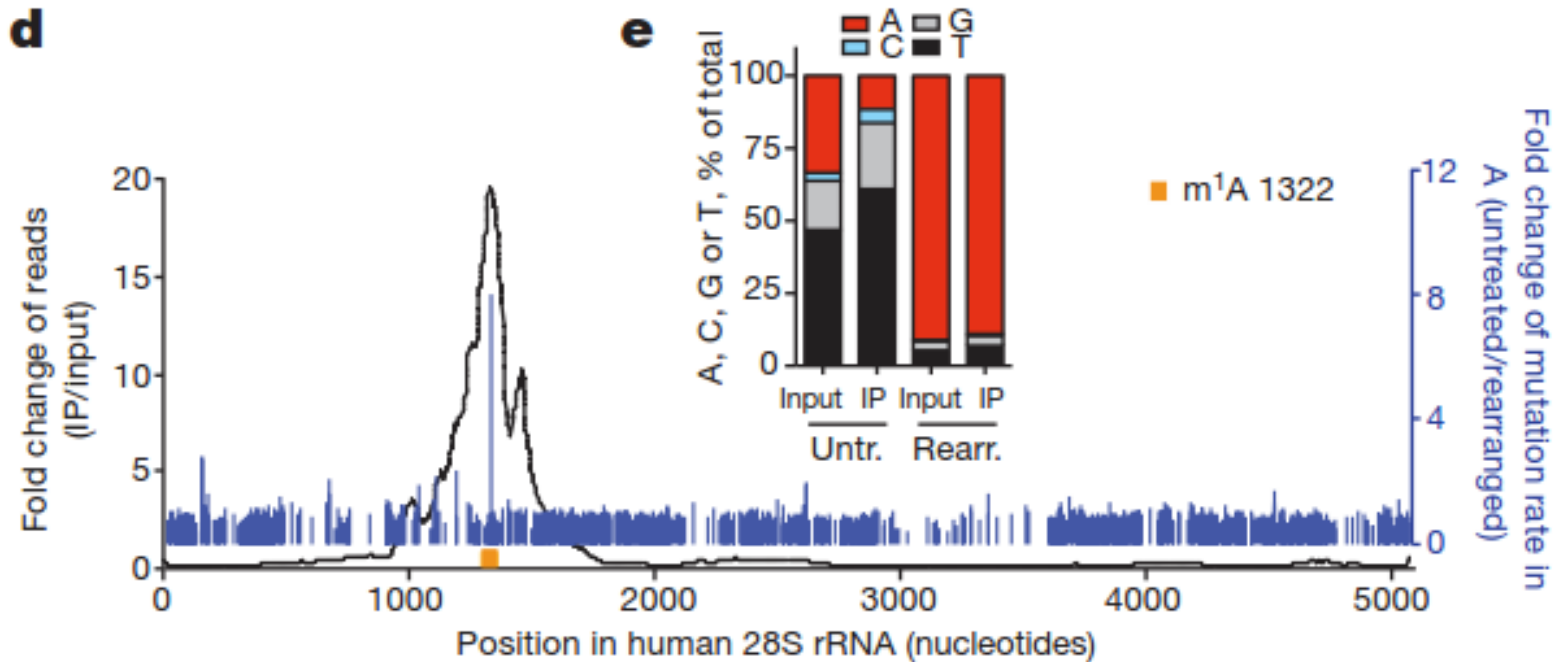
Methylated mRNAs



Methylated mRNAs

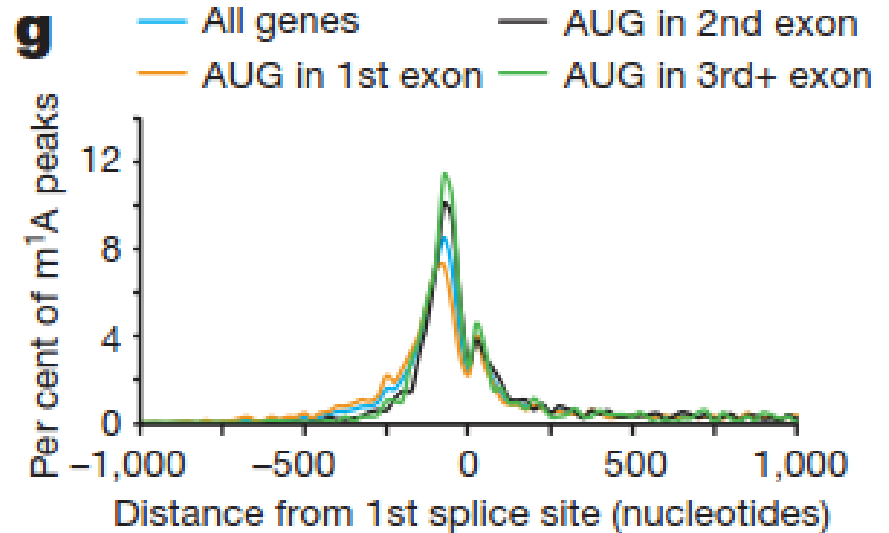
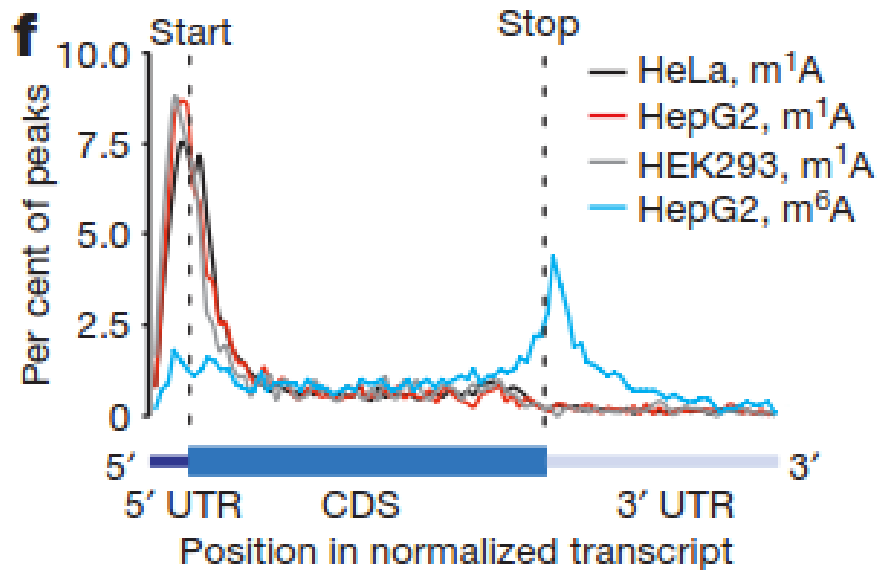


Methylated noncoding RNAs



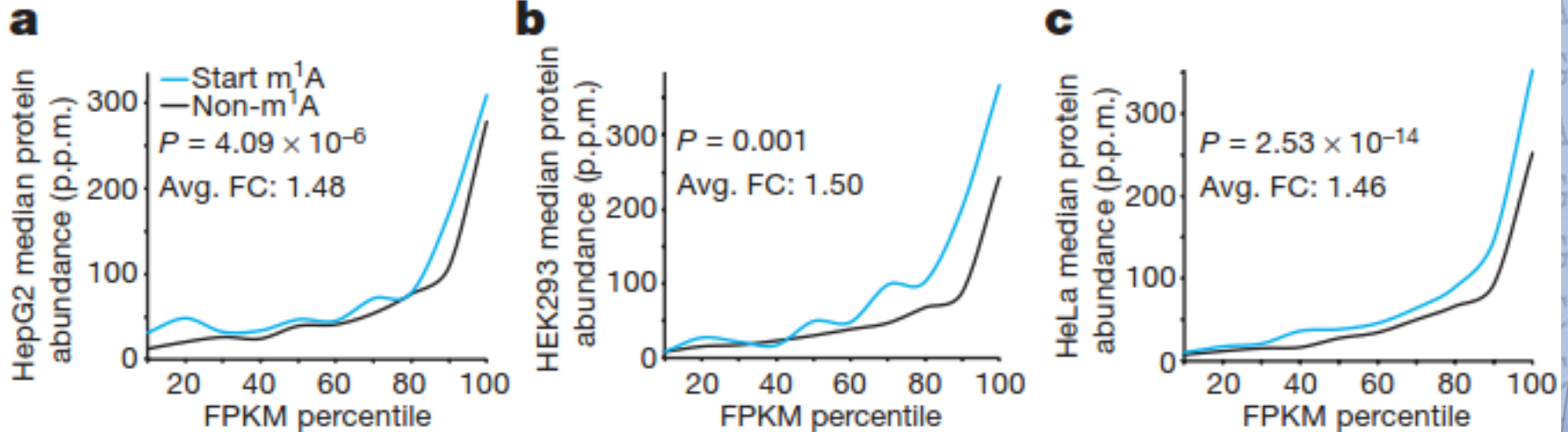
Methylated mRNAs

- Associated with translation starts and stops
- Correlated to splice sites

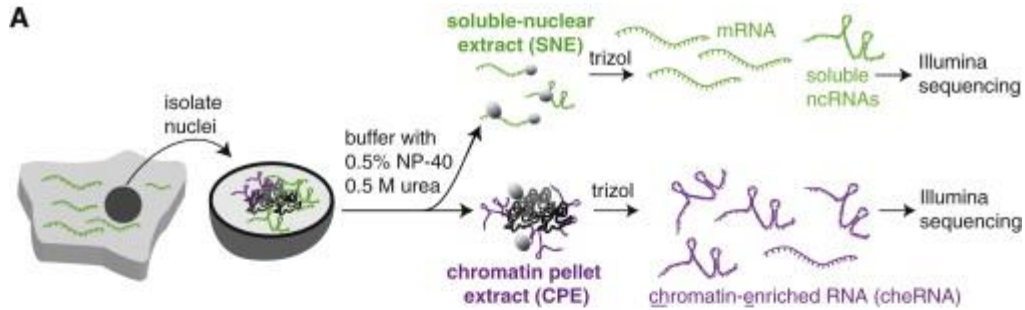


Methylated mRNAs

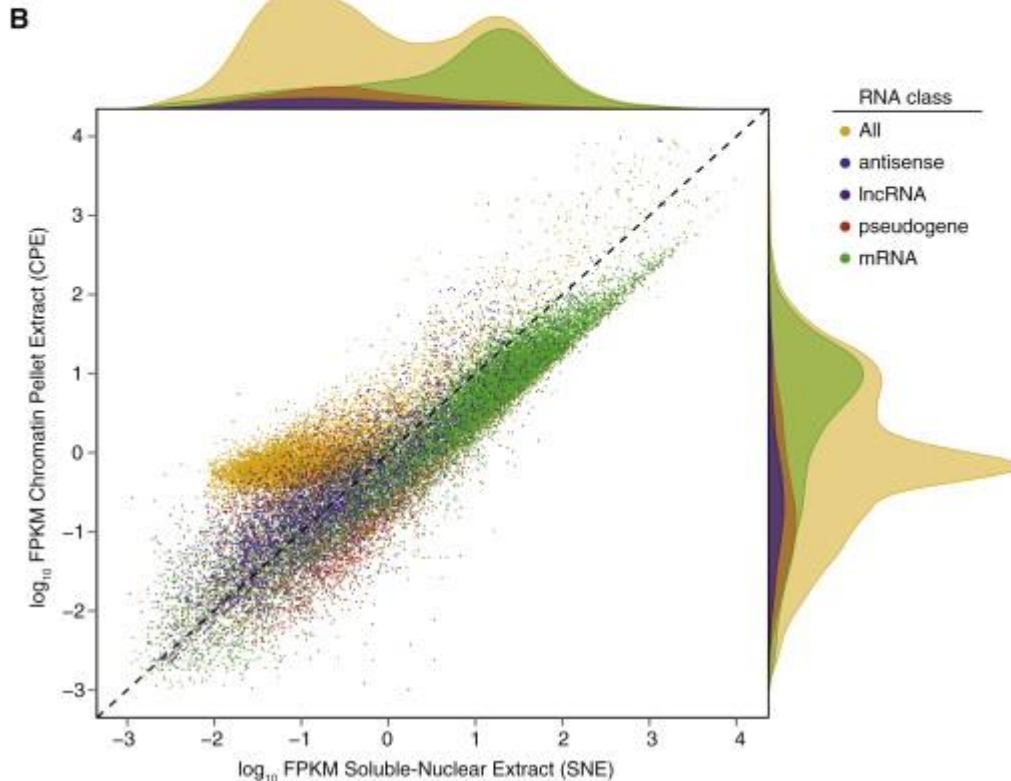
- m¹A around the start codon correlates with higher protein



chromatin-enriched RNAs



- Soluble vs. chromatin bound lncRNAs



Werner et al. 2015



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%).

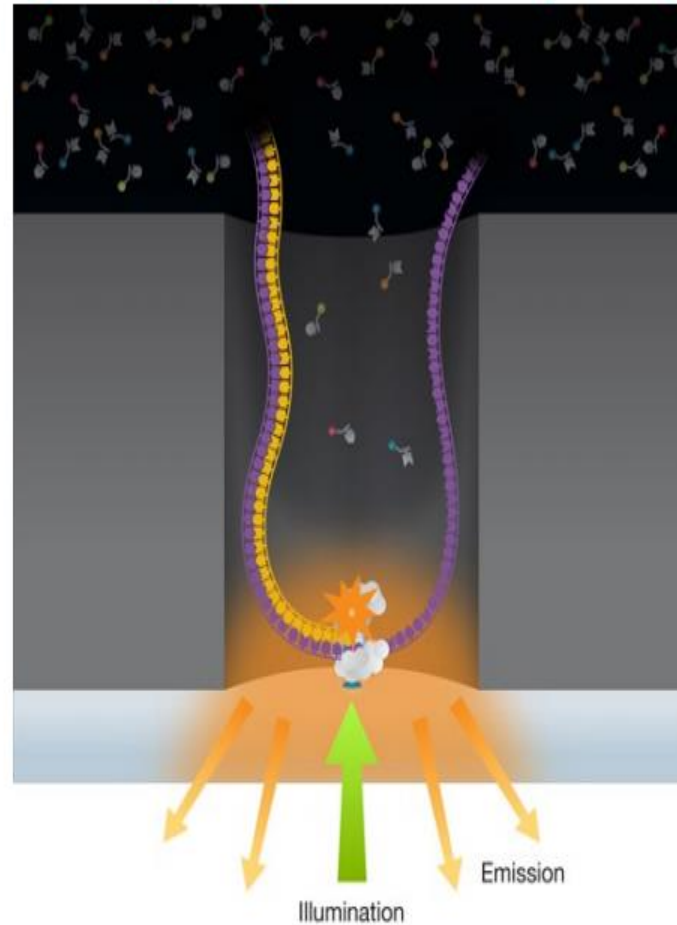
Complementary to short accurate reads of Illumina

TGAGAT
TATGAGC
TAAATCTC
TACCCCT
GCTGAA
ATTCCCT
TCTGGGA
GAAATT
TGTGAA
AAGGAG
TTTGGG
CGCCAG
TCCCAG
AATTGC
TCTCCA
AAGGCT
AATTGA
GCACAA
ATACCA
GCTTTT
TTTATC

Third Generation Sequencing : Single Molecule Sequencing

Pacific Biosciences

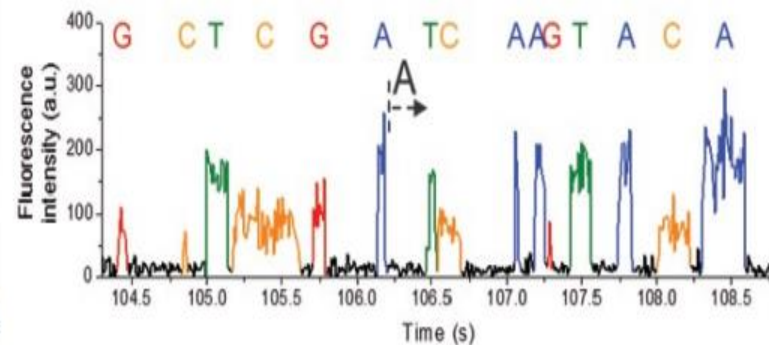
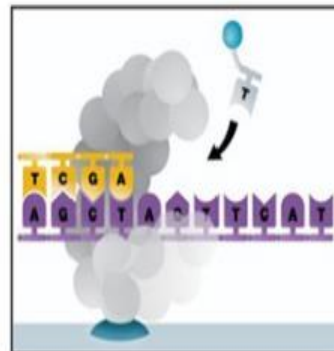
70 nm aperture
“Zero Mode Waveguide”



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



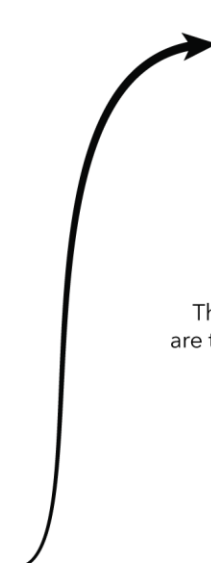
Start with high-quality double stranded DNA



Ligate SMRTbell adapters and size select



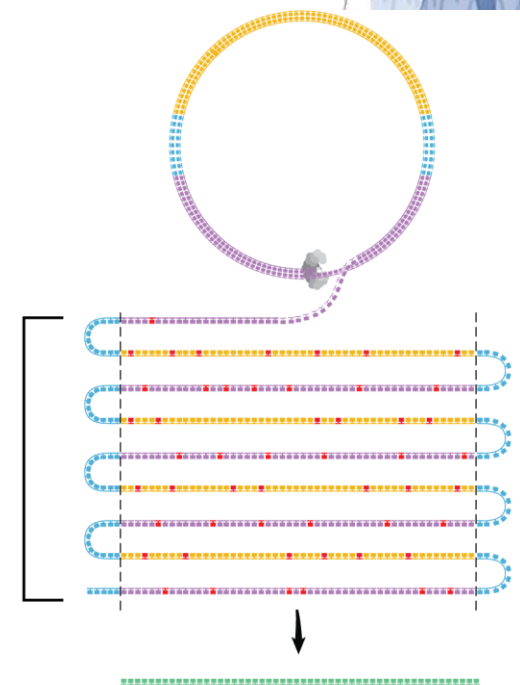
Anneal primers and bind DNA polymerase



Circularized DNA is sequenced in repeated passes

The polymerase reads are trimmed of adapters to yield subreads

Consensus is called from subreads



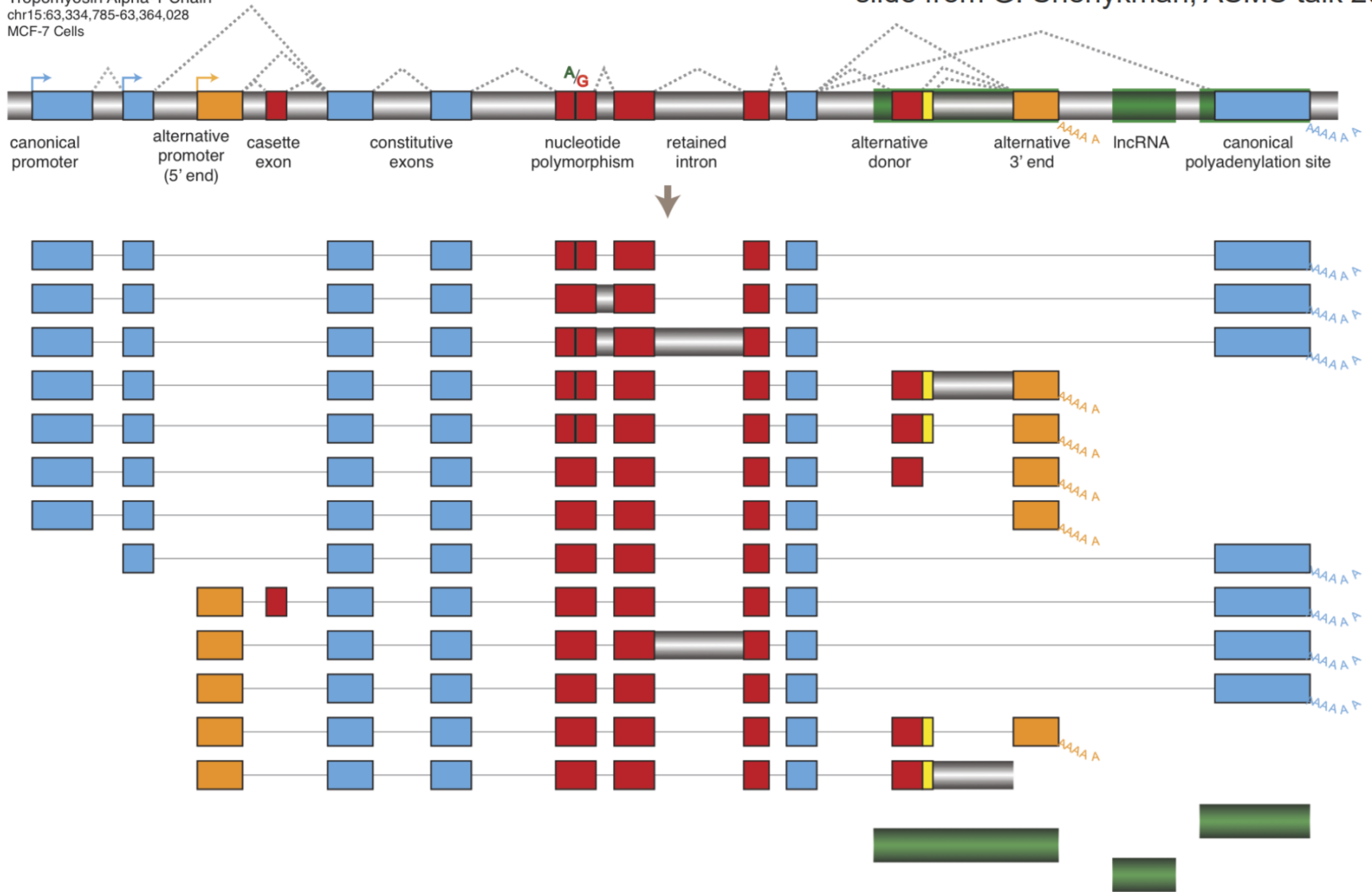
HiFi READ
(>99% accuracy)



A Single Gene Locus → Many Transcripts

Tropomyosin Alpha-1 Chain
chr15:63,334,785-63,364,028
MCF-7 Cells

slide from G. Shenykman, ASMS talk 2014

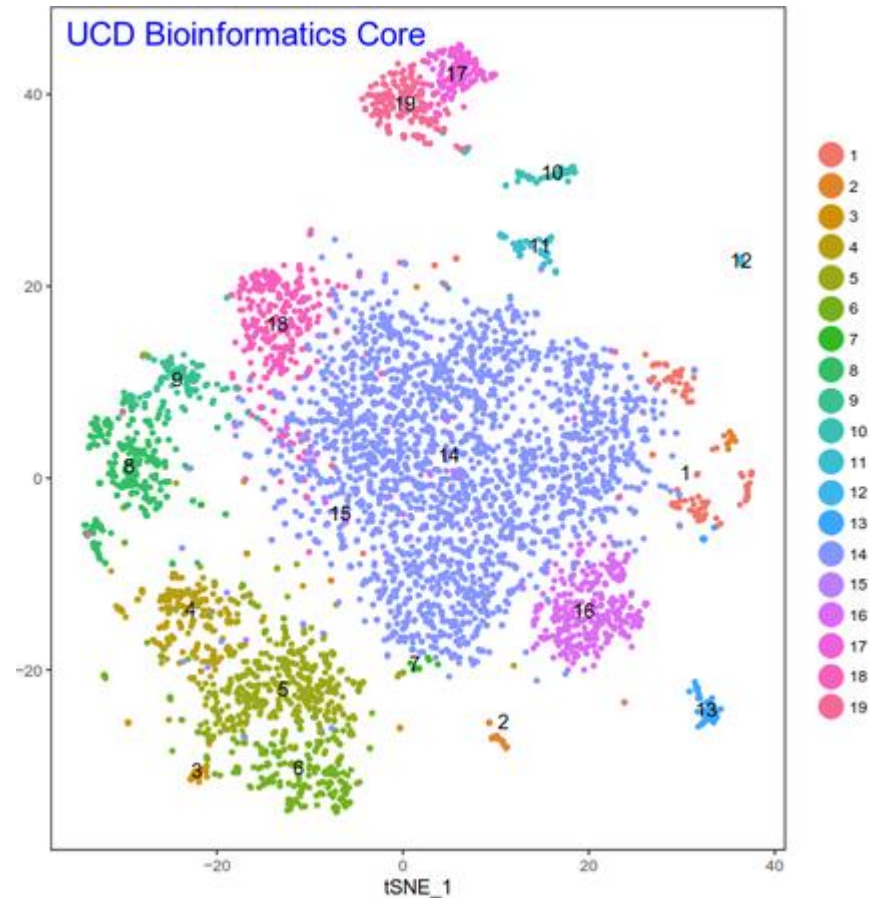


Iso-Seq Pacbio

- 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?)

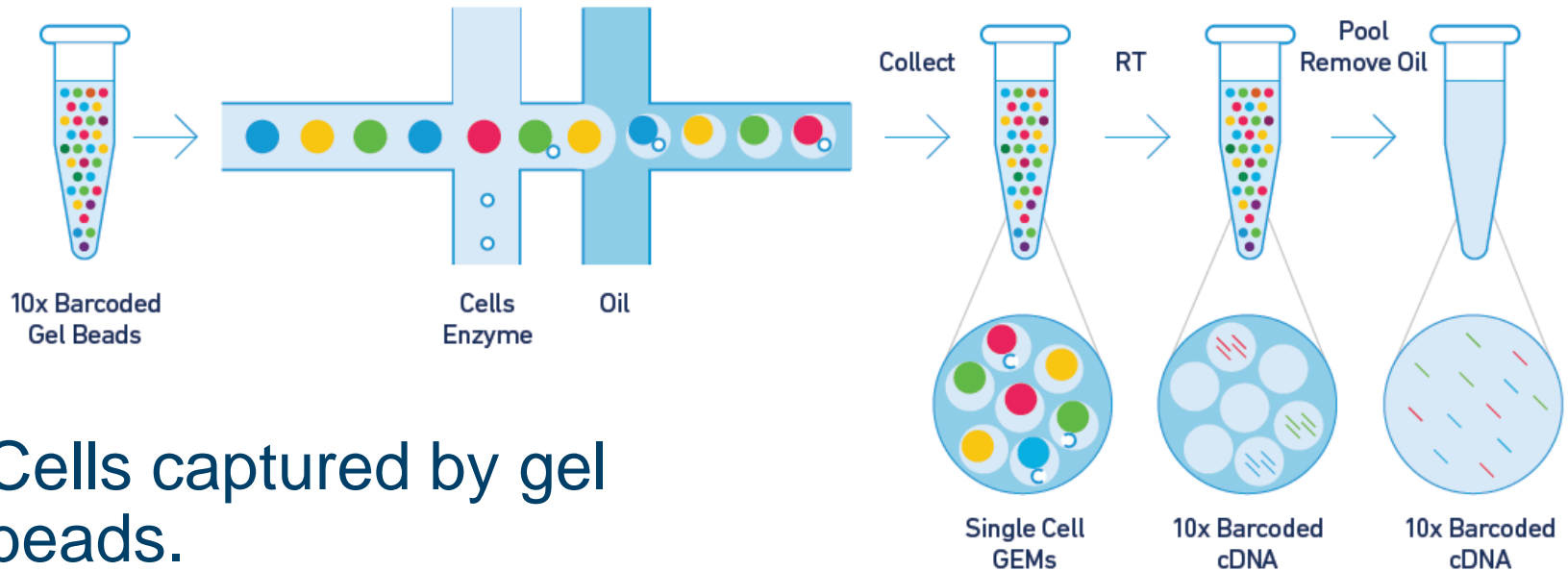


scRNA-seq (single cells)



- 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



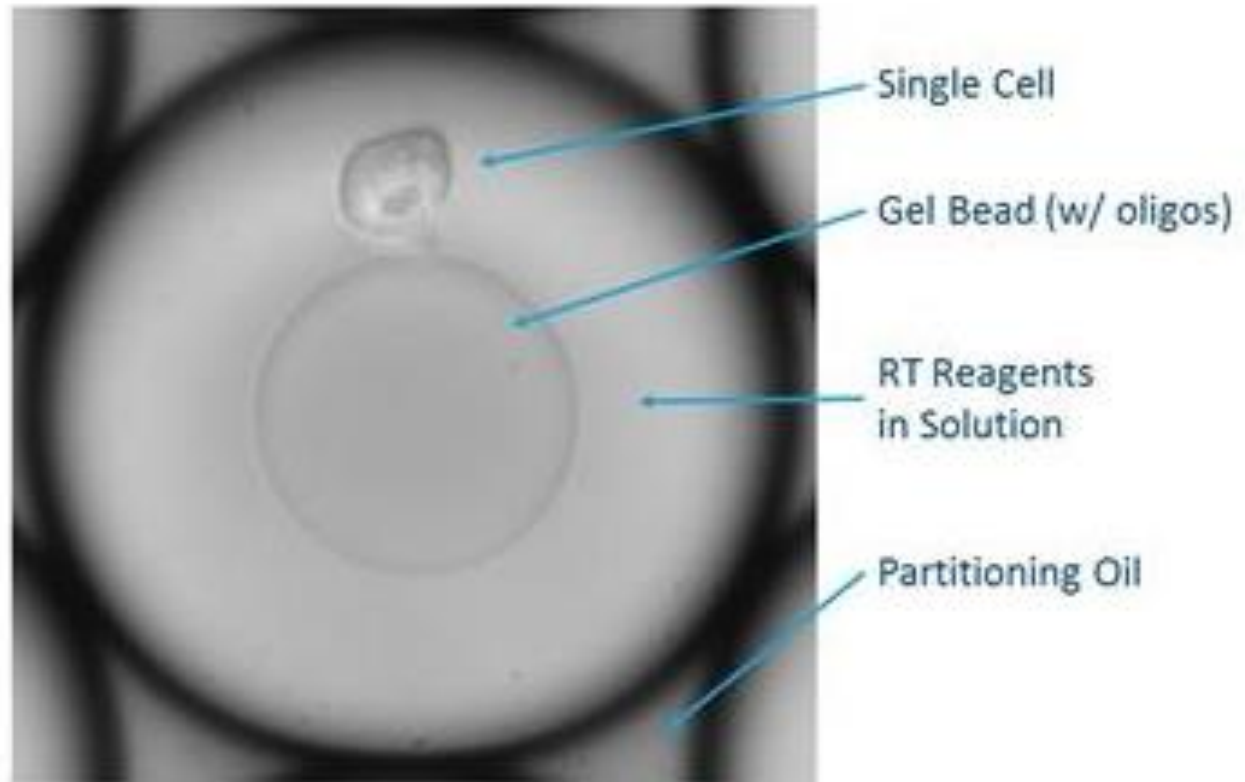
- Cells captured by gel beads.
- 10X barcode added to transcript.
- cDNA amplification.
- Transcriptional profiling of individual cells due to unique barcodes / UMIs.

Transcriptional profiling of individual cells



Cell partitioning into GEMs

- GEM

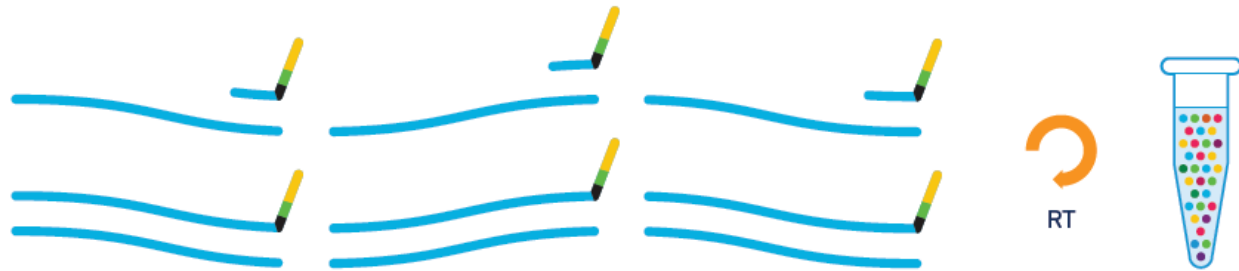


Credit: 10X Genomics

Library preparation

1 Molecular Barcoding in GEMs

Credit: 10X Genomics



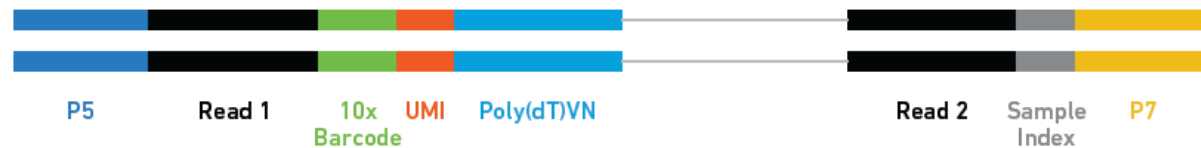
2 Pool, Library Prep



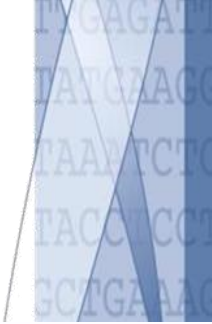
3 Sequence and Analyze



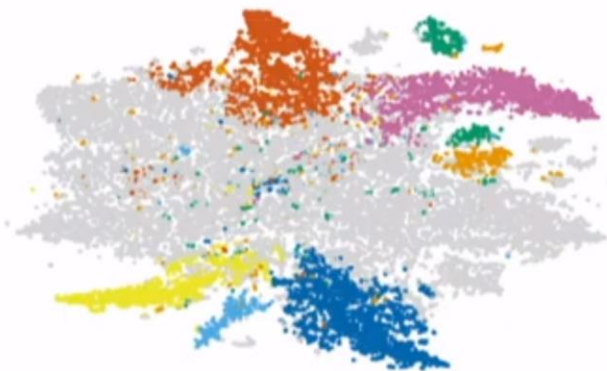
Final Library Construct



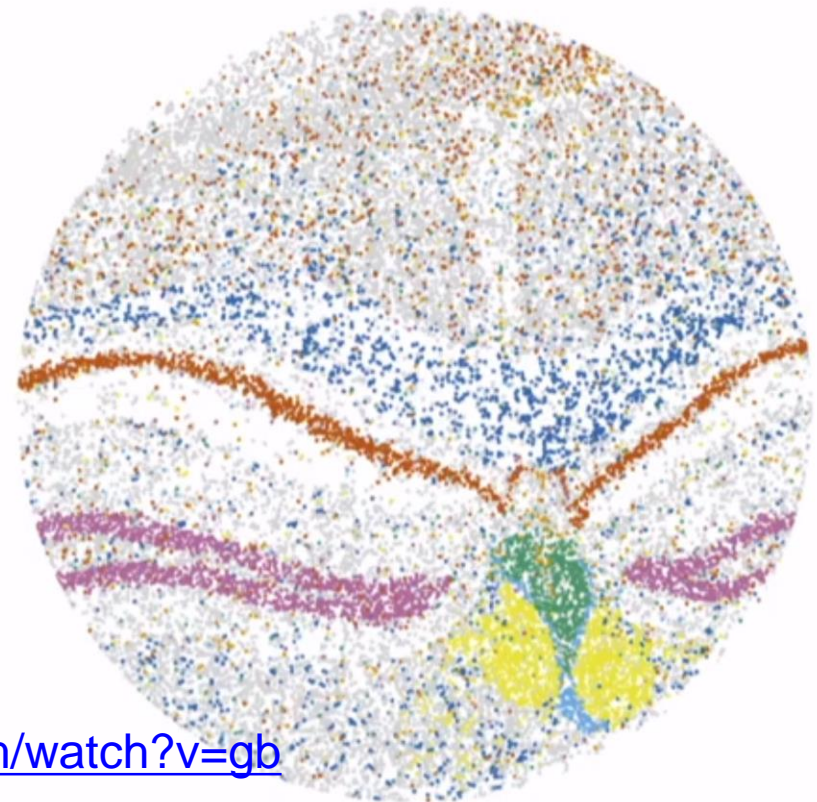
Spatial Transcriptomics (10XGenomics Visium; Slide-Seq)



4 Map gene expression into space



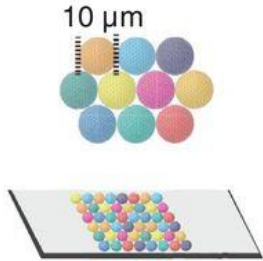
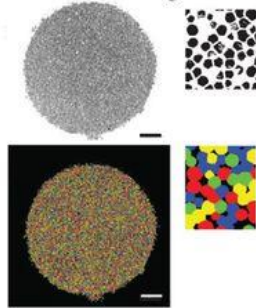
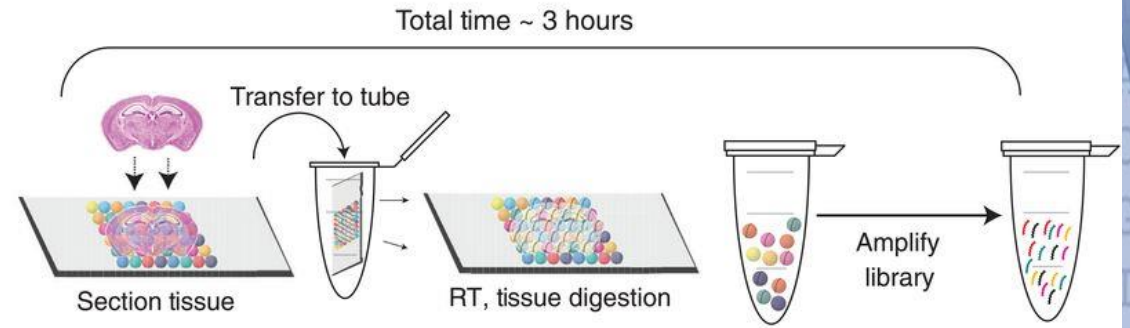
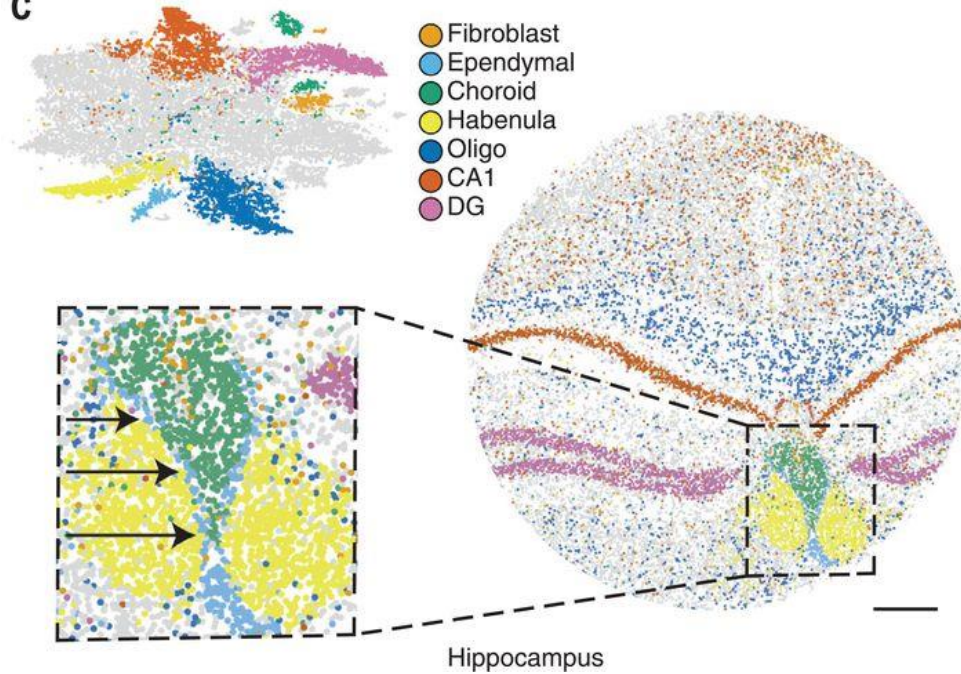
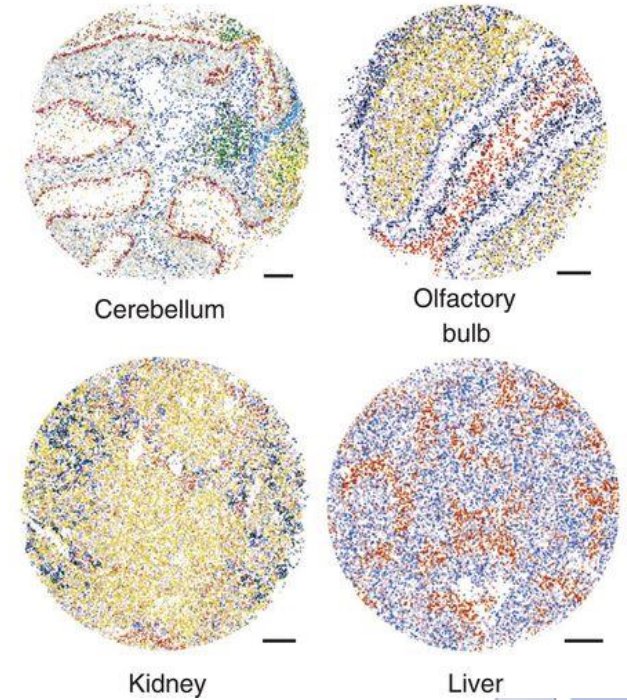
- Fibroblast
- Ependymal
- Choroid
- Habenula
- Oligo
- CA1
- DG



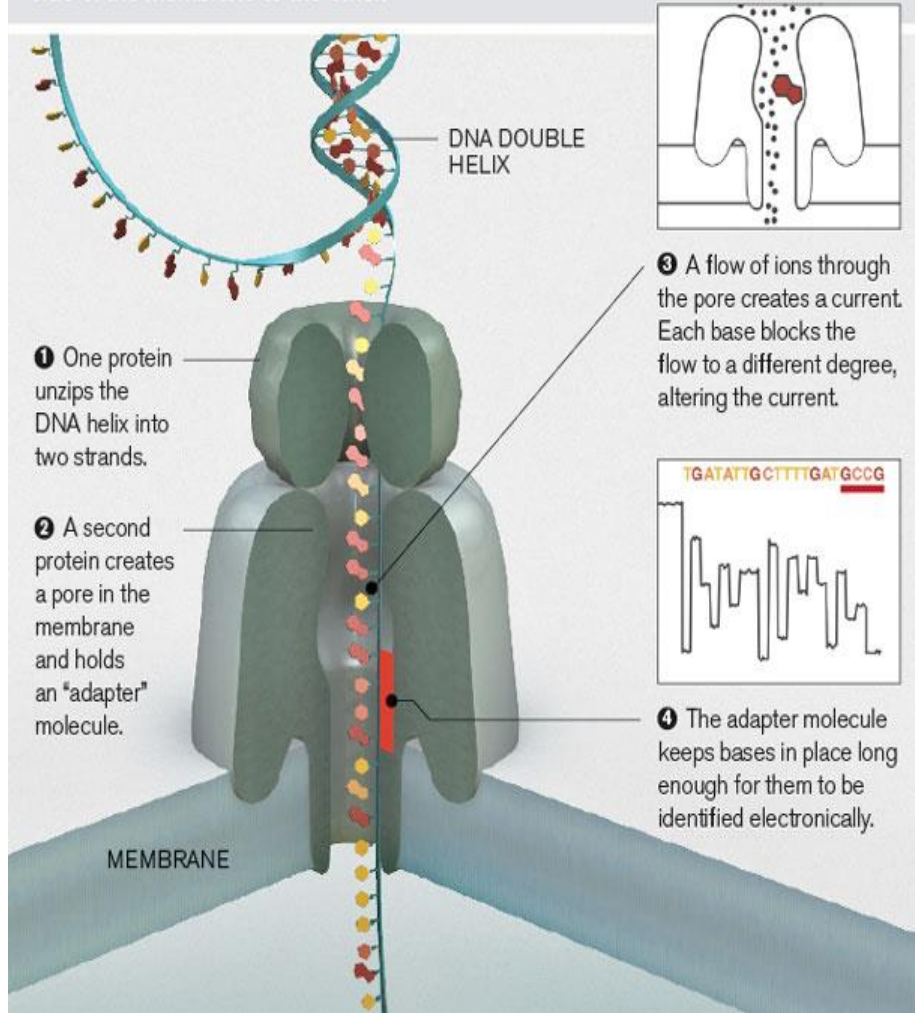
Evan Macosko

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



A Bead deposition*In situ* indexing**B****C****D**

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.



TGAGAT
IATGAGC
TAAATCTC
TACCCTCC
GCTGAAAC
ATTCCCTC
TCTGGGAA
GAAATTAT
TGTTGAA
AAGGAGC
TTTGGG
CGCCAGC
TCCCAGC
AATTGCA
TCTCCAA
AAGGCTT
AATTGAA
GCACAA
ATACCA
GCTTTT
TTTATC

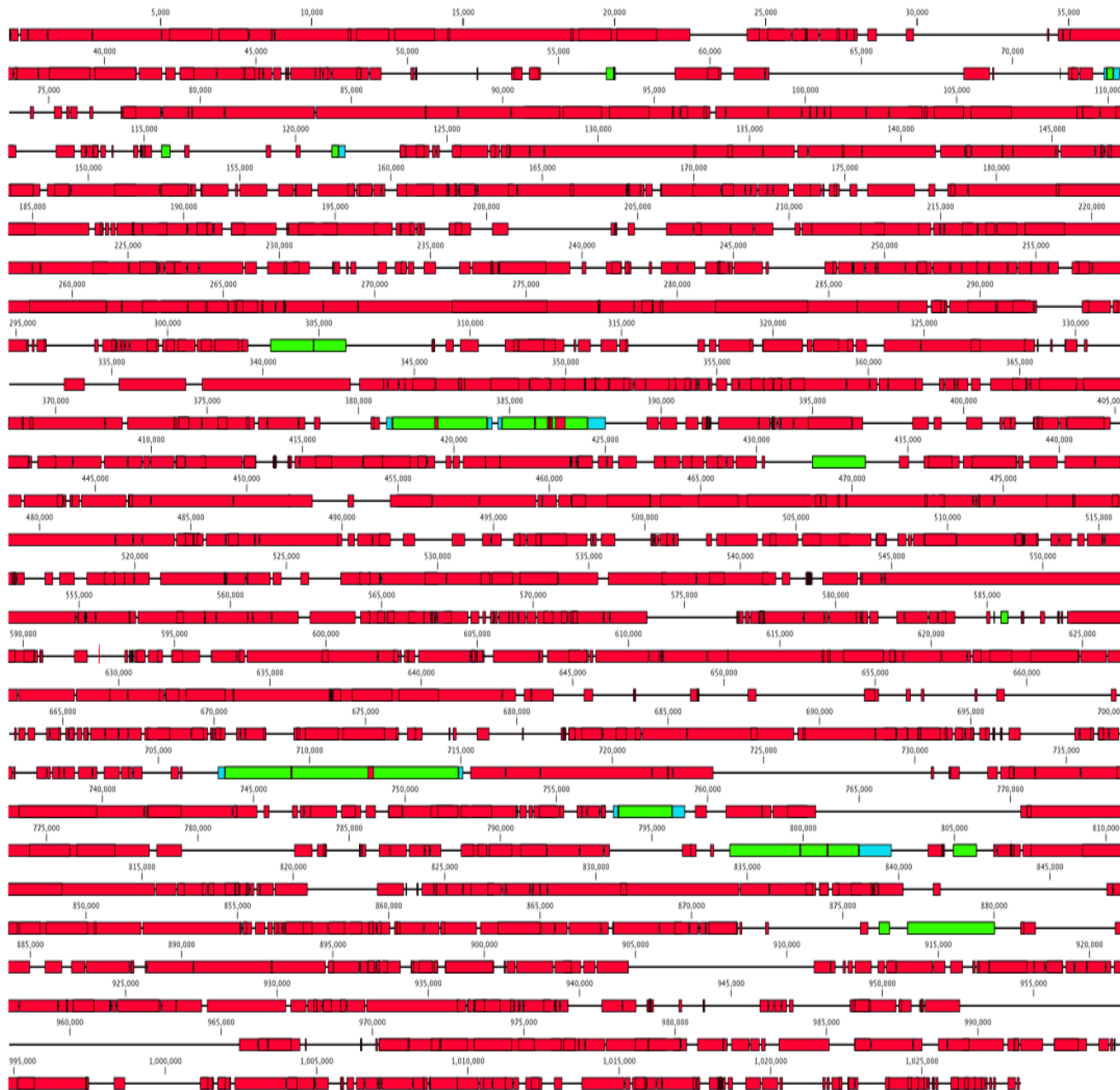
Future's so bright





Thank you!





First Sequencing of CGG-repeat Alleles in Human Fragile X Syndrome using PacBio RS Sequencer



Paul Hagerman, Biochemistry and Molecular Medicine, SOM.

- Single-molecule sequencing of pure CGG array,
 - first for disease-relevant allele. Loomis *et al.* (2012) *Genome Research*.
 - applicable to many other tandem repeat disorders.
- Direct genomic DNA sequencing of methyl groups,
 - direct epigenetic sequencing (paper under review).
- Discovered 100% bias toward methylation of 20 CGG-repeat allele in female,
 - first methylated DNA sequence in human

dis

