

High Throughput Sequencing the Multi-Tool of Life Sciences

Lutz Froenicke

DNA Technologies and Expression Analysis
Cores

UCD Genome Center



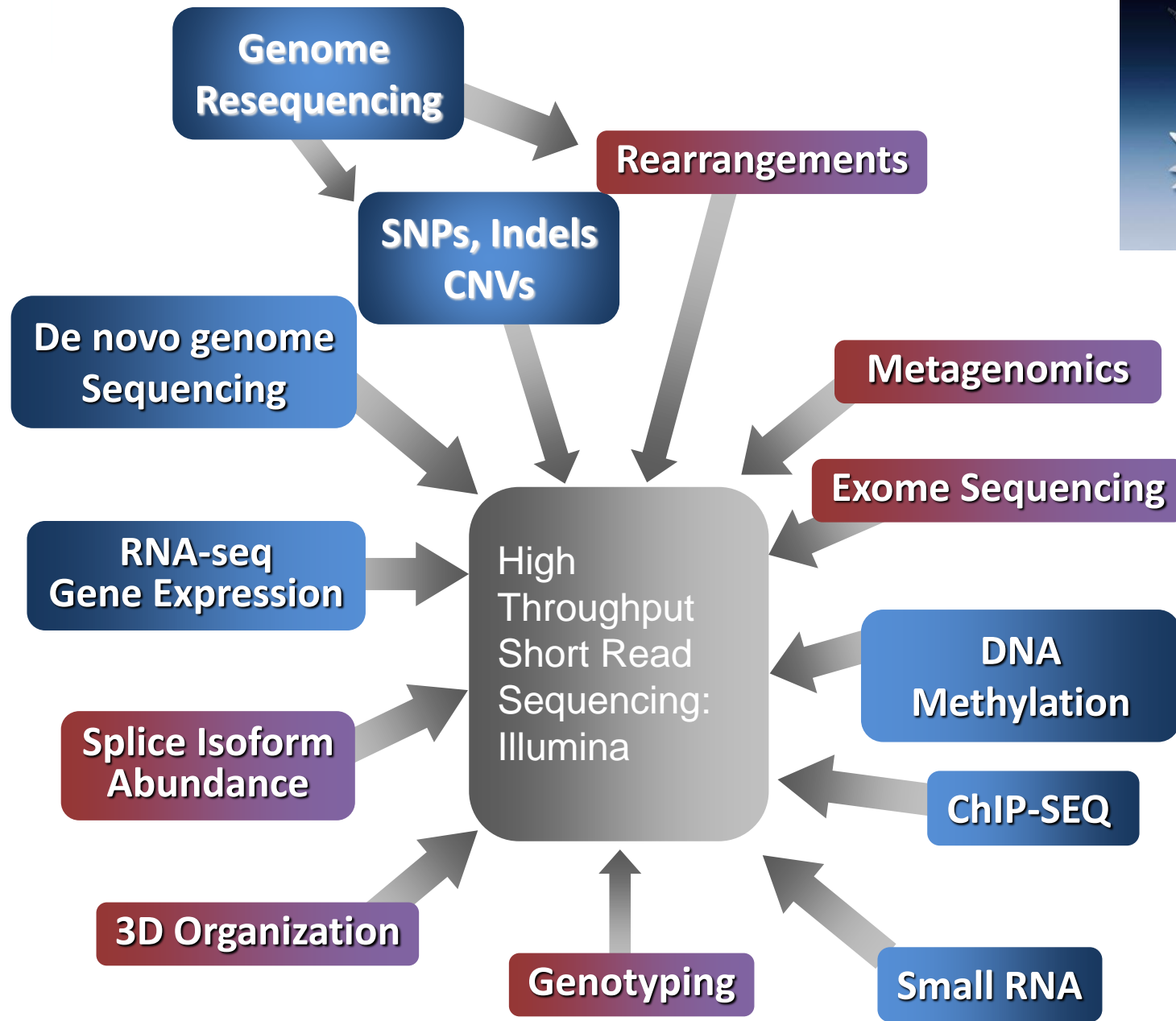
DNA Technologies & Expression Analysis Cores

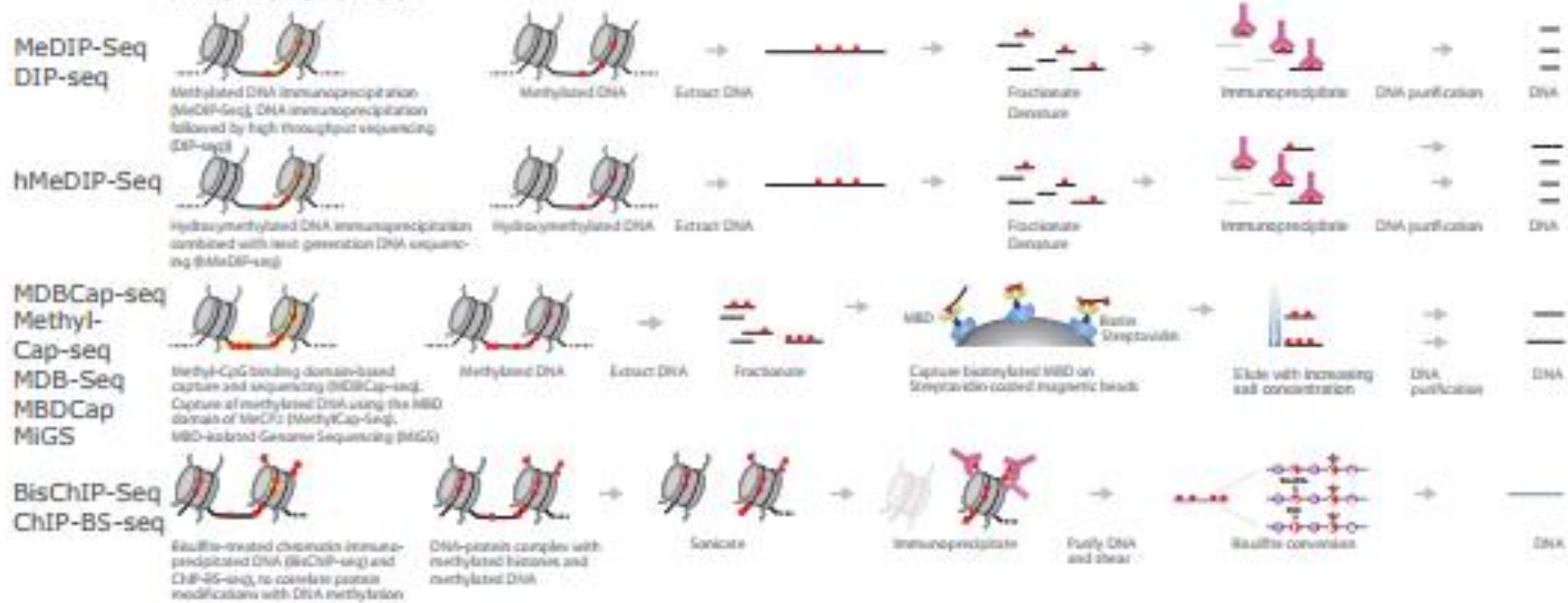
- HT Sequencing (Illumina & PacBio)
- Illumina microarray (for genotyping – Illumina has discontinued expression analysis)
- consultations
- introducing new technologies to campus
- shared equipment (accessible after training)
- teaching (workshops)



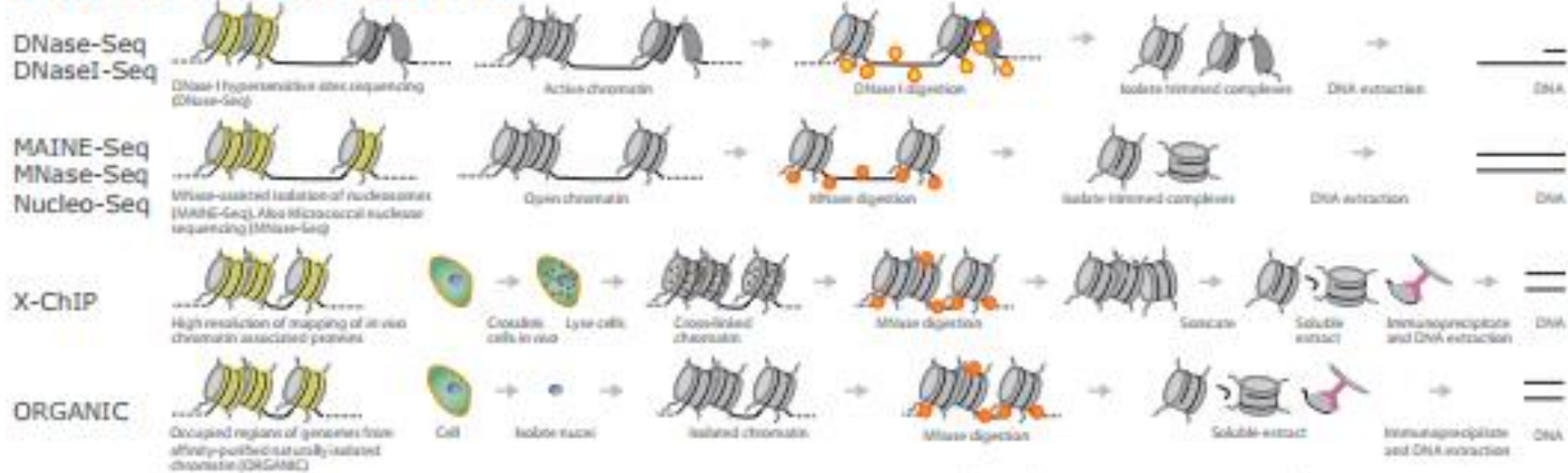
Complementary Approaches

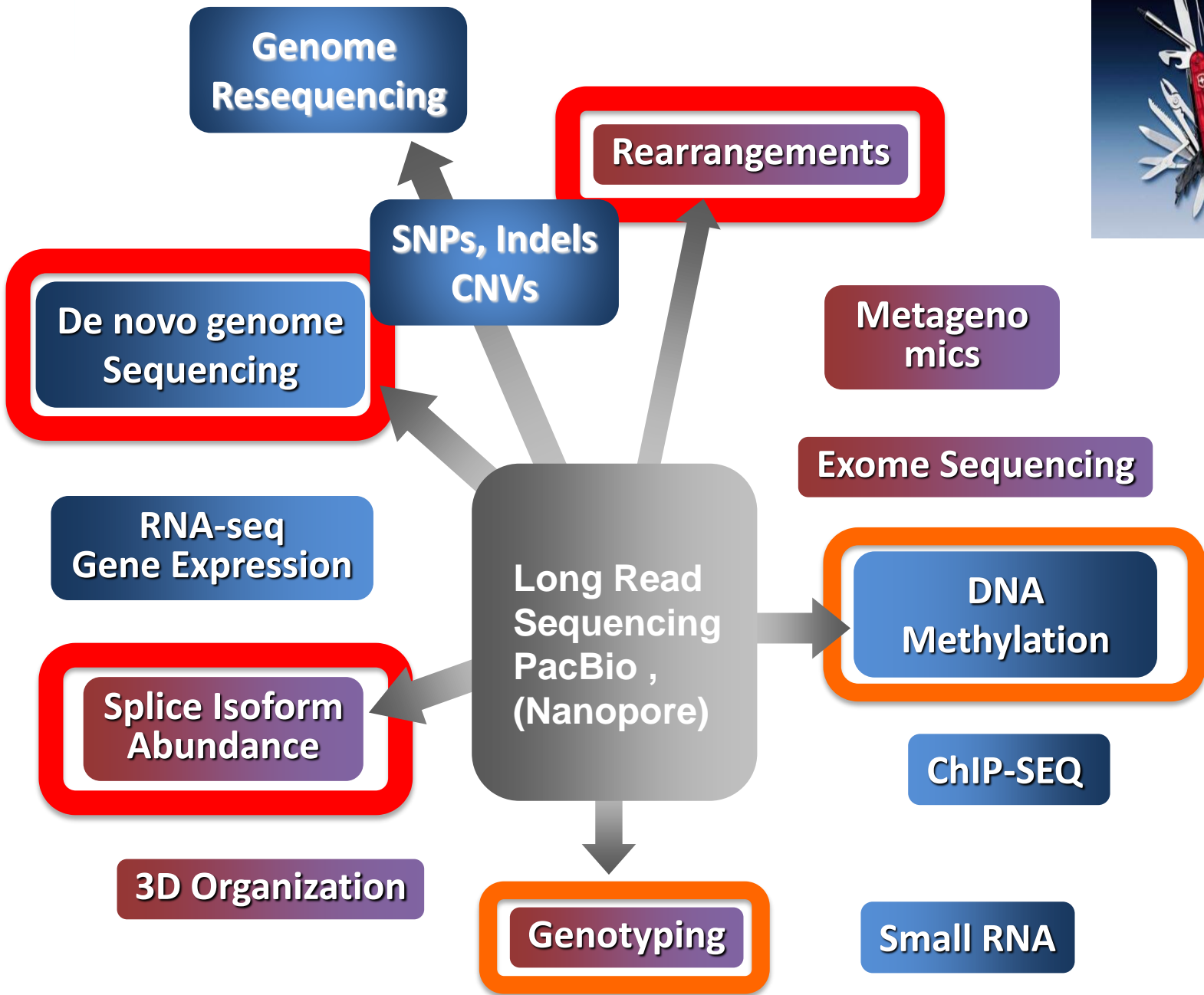
| ILLUMINA | PACBIO |
|--|--|
| Still-imaging of clusters (~1000 clonal molecules) | Movie recordings fluorescence of single molecules |
| Short reads - 2x300 bp Miseq | Up to 60 kb, N50 23 kb |
| Repeats are mostly not analyzable | spans retro elements |
| High output - up to 100 Gb per lane | up to 1,3 Gb and 5 Gb per SMRT-cell |
| High accuracy (< 0.5 %) | Error rate 15 % |
| Considerable base composition bias | No base composition bias |
| Very affordable | Costs 5 to 10 times higher |
| <i>De novo</i> assemblies of thousands of scaffolds | “Near perfect” genome assemblies |





DNA-Protein Interactions





Illumina sequencing workflow

- **Library Construction**
- Cluster Formation
- Sequencing
- Data Analysis



Fragmentation

- Mechanical shearing:

- BioRuptor
- Covaris

DNA, RNA

- Enzymatic:

- Fragmentase, RNAse3

DNA, RNA

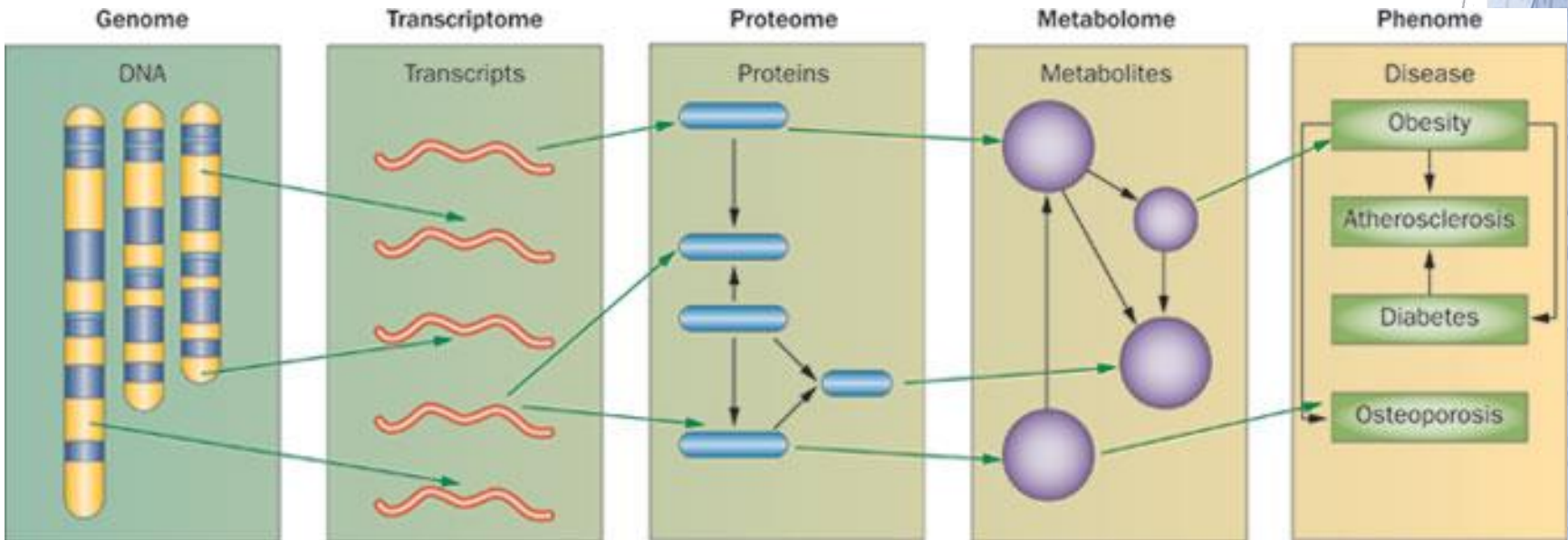
- Chemical: Mg^{2+} , Zn^{2+}

→ RNA



“DNA makes RNA and RNA makes protein”

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

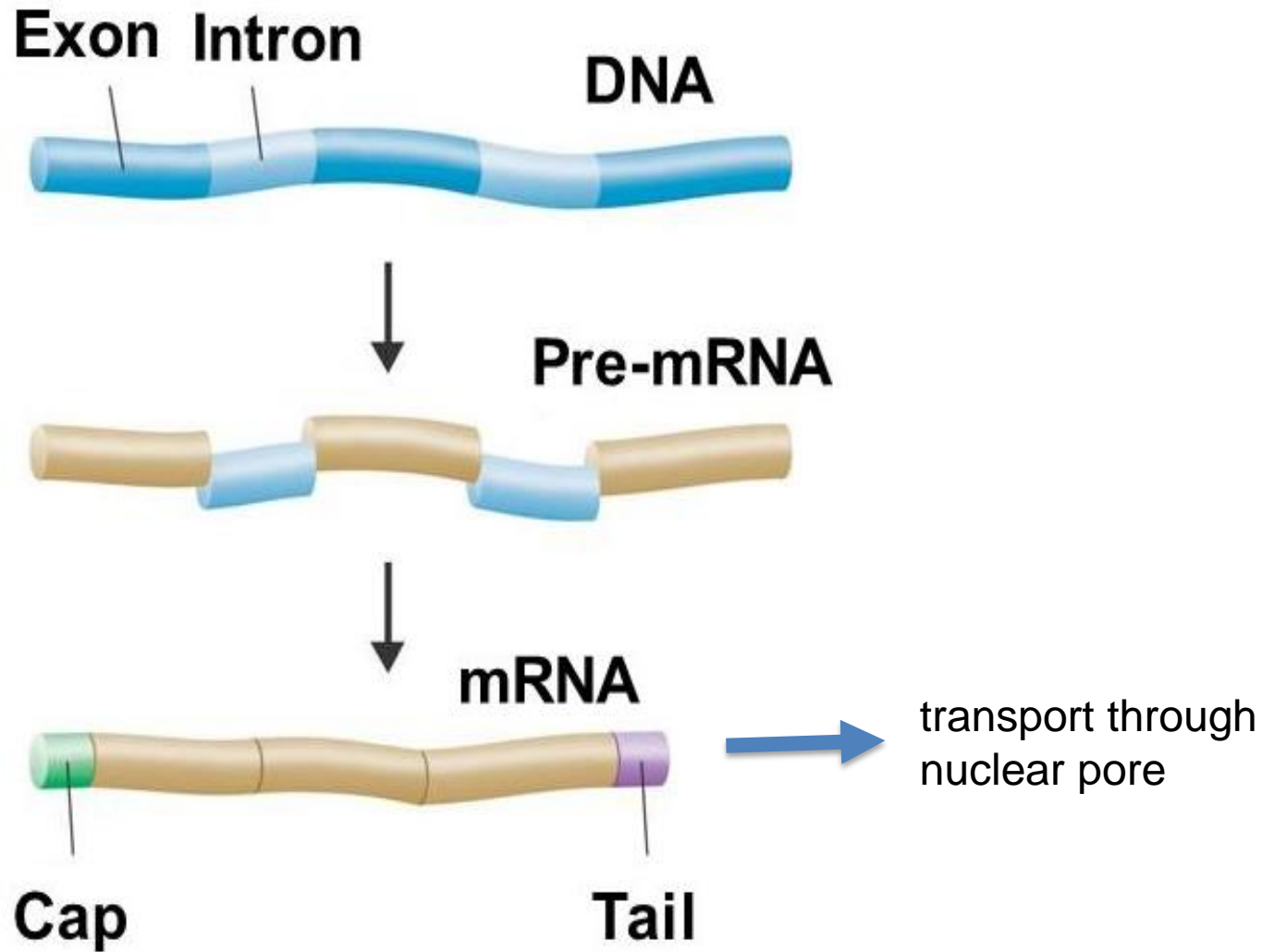


DNA Tech & Expression Analysis Proteomics Core Metabolomics Core

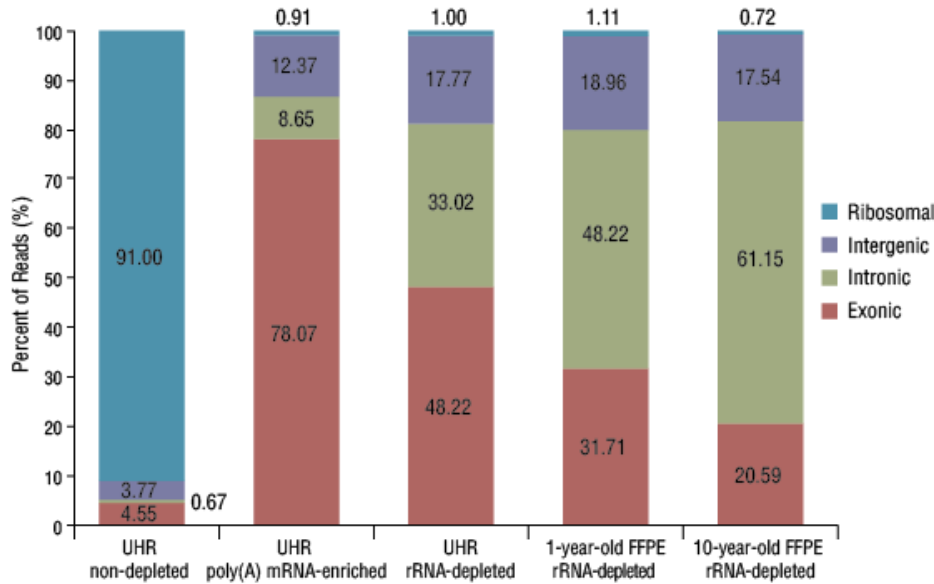
UCD Genome Center

nature
REVIEWS **CARDIOLOGY**

transcription and processing in nucleus



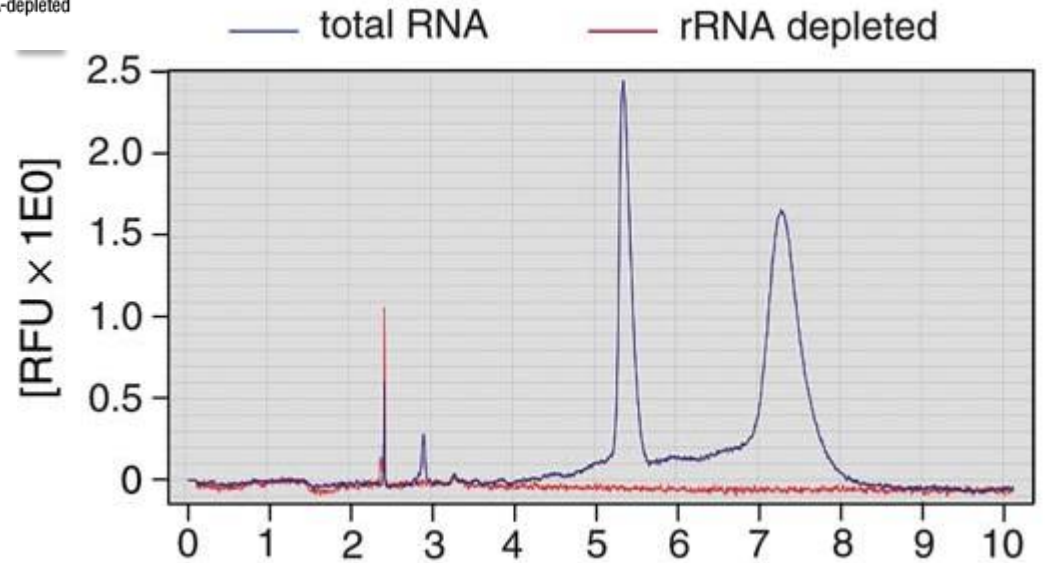
mRNA makes up only about 2% of a total RNA sample



- more than 90% rRNA content

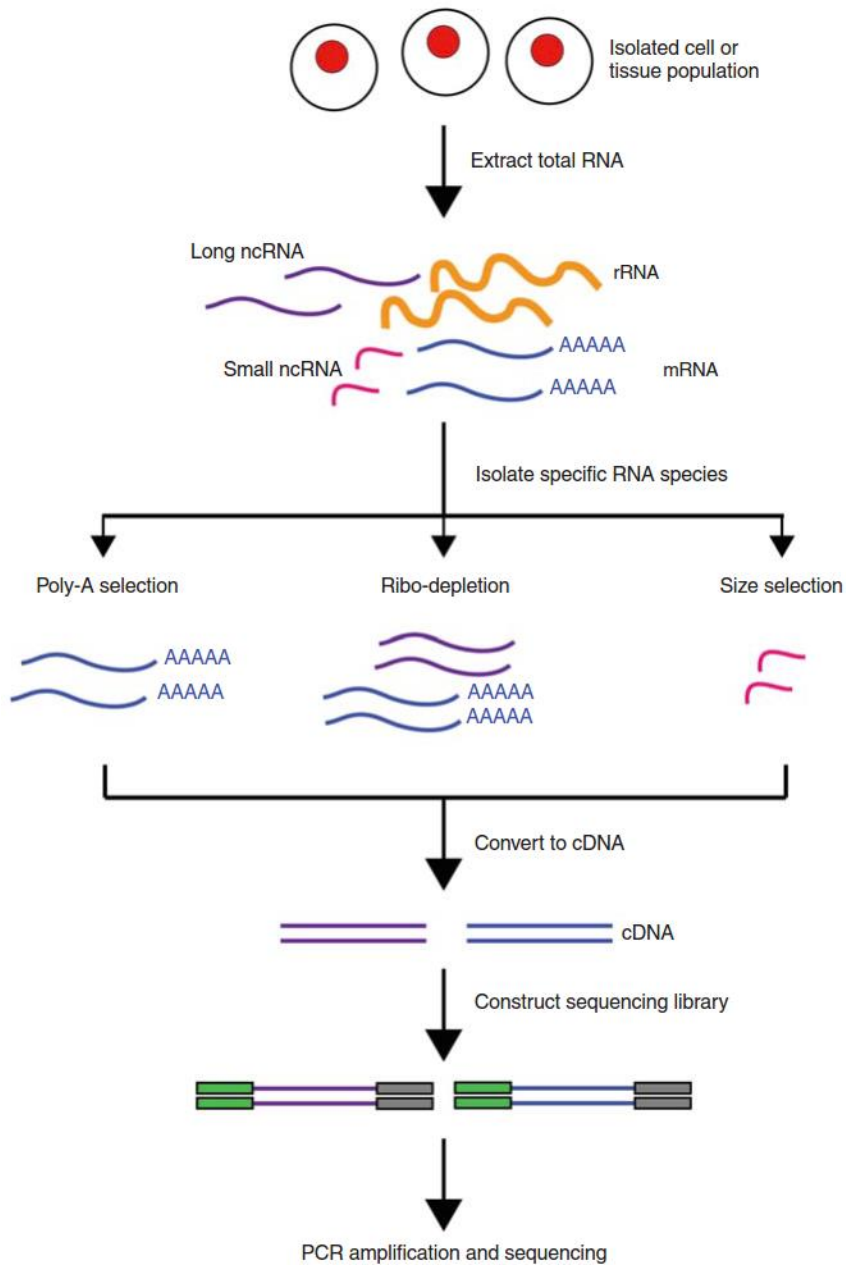
- multiple other non-coding RNA species

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

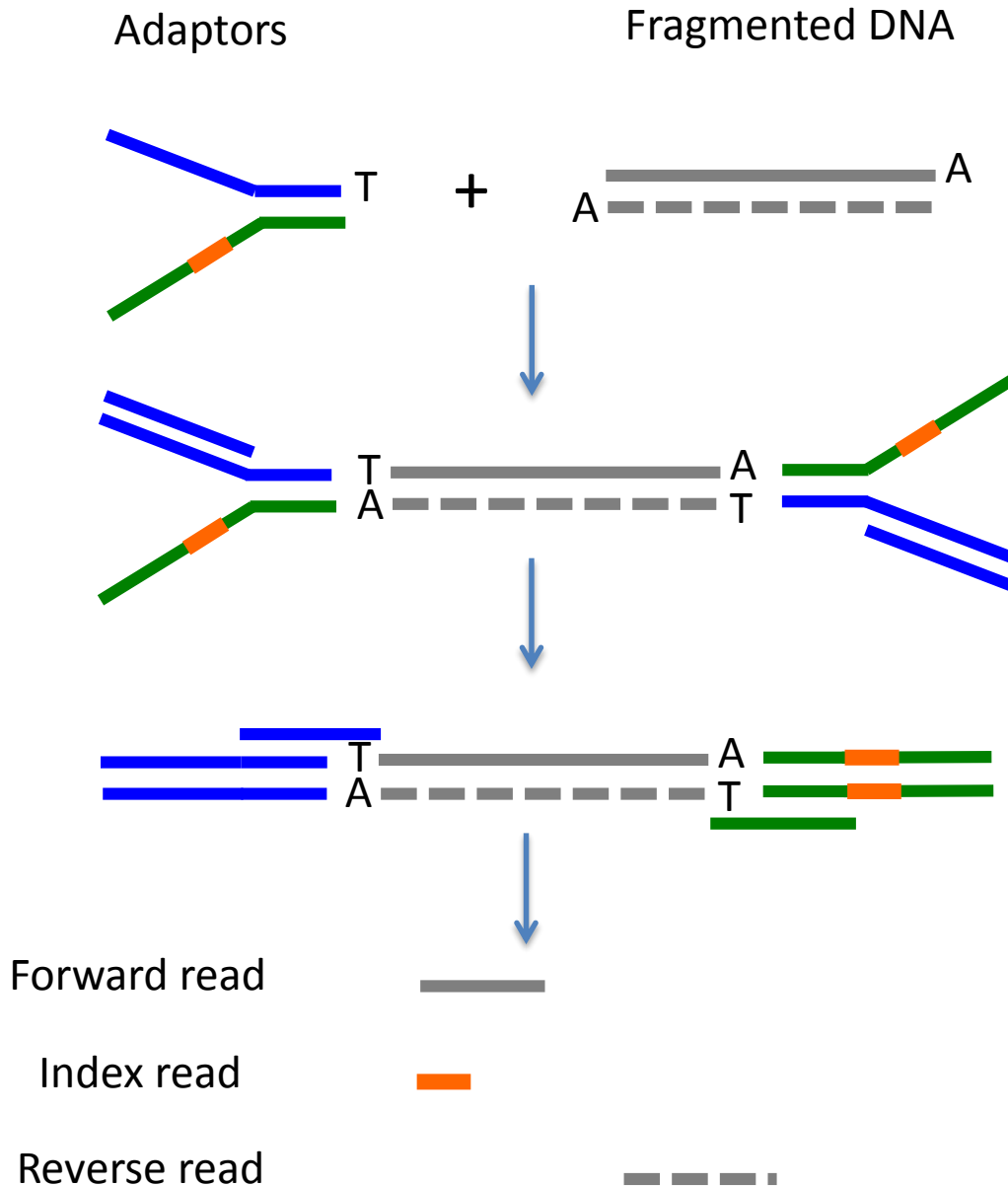


RNA-seq?

Sorry - we are only sequencing DNA.



“Truseq –style” indexed adaptors



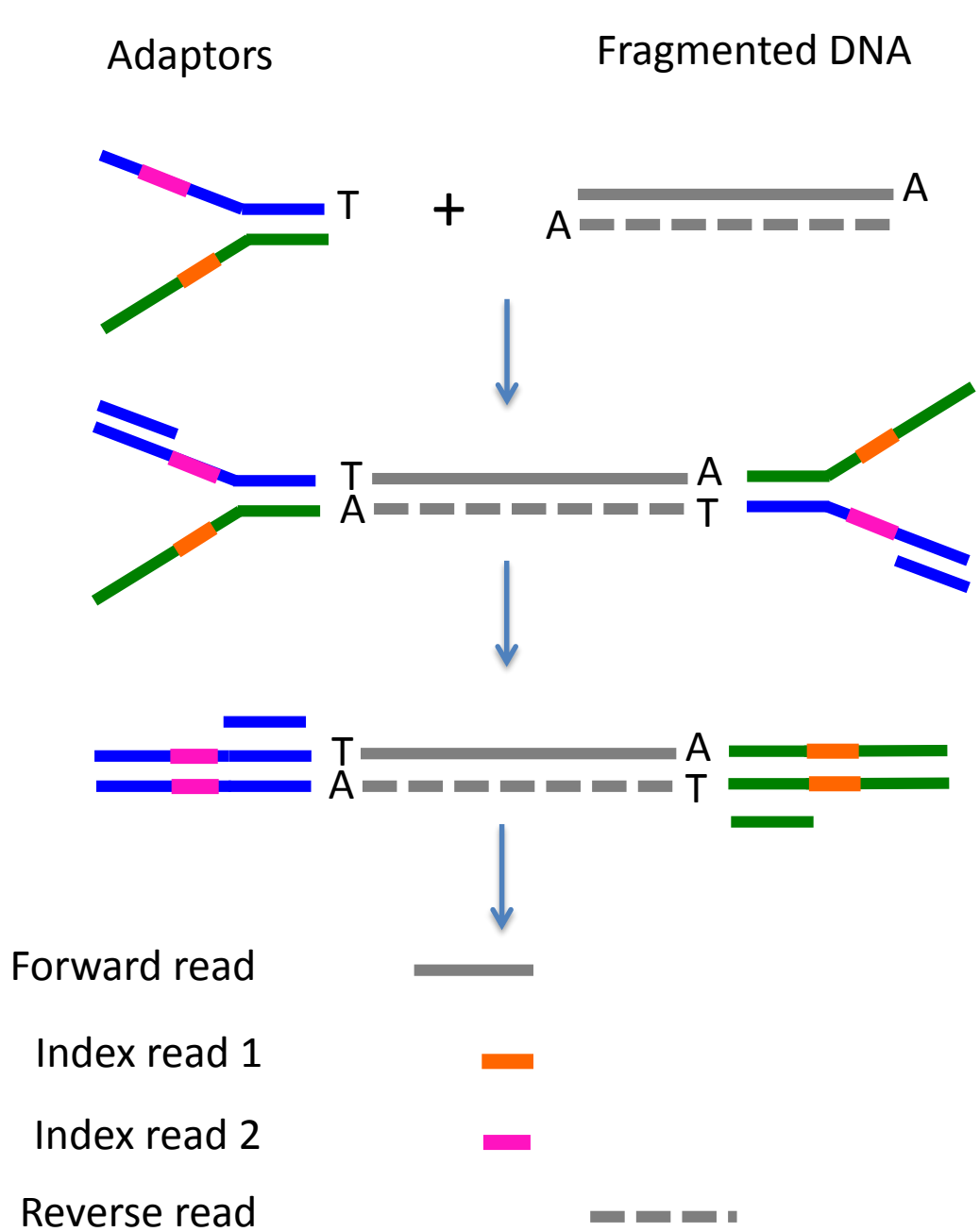
Advantages:

- Index independent of read
- > more data
- > no more clustering problems

Problems:

- Need more reagents
- Index only on one side





“Dual indexed” adaptors

For 96 reactions

Simple index:

96 B adaptors

1 A adaptor

Dual index:

12 A adaptors

8 B adaptors

Advantages:

Cheaper

Indexing information on both sides

Problems:

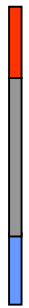
TBA...



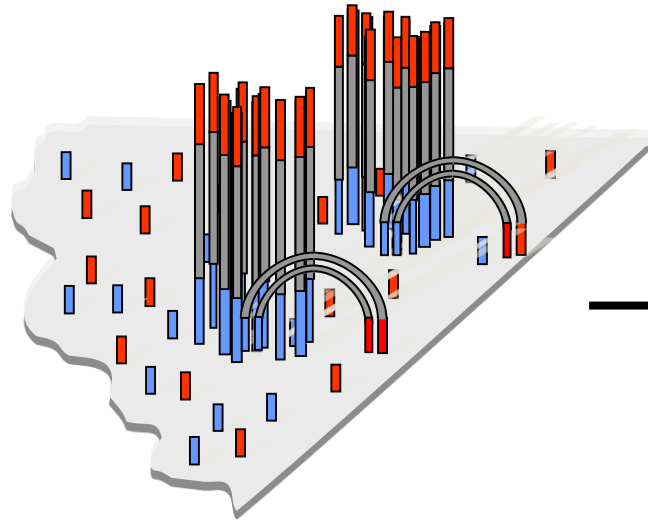
Illumina Sequencing Technology

Sequencing By Synthesis (SBS) Technology

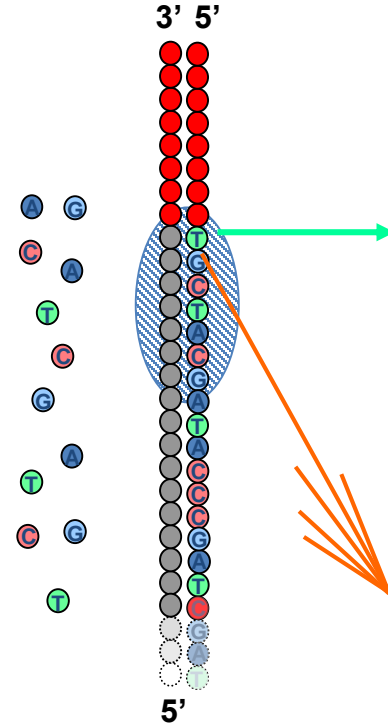
DNA
(0.1-1.0 ug)



Library
preparation



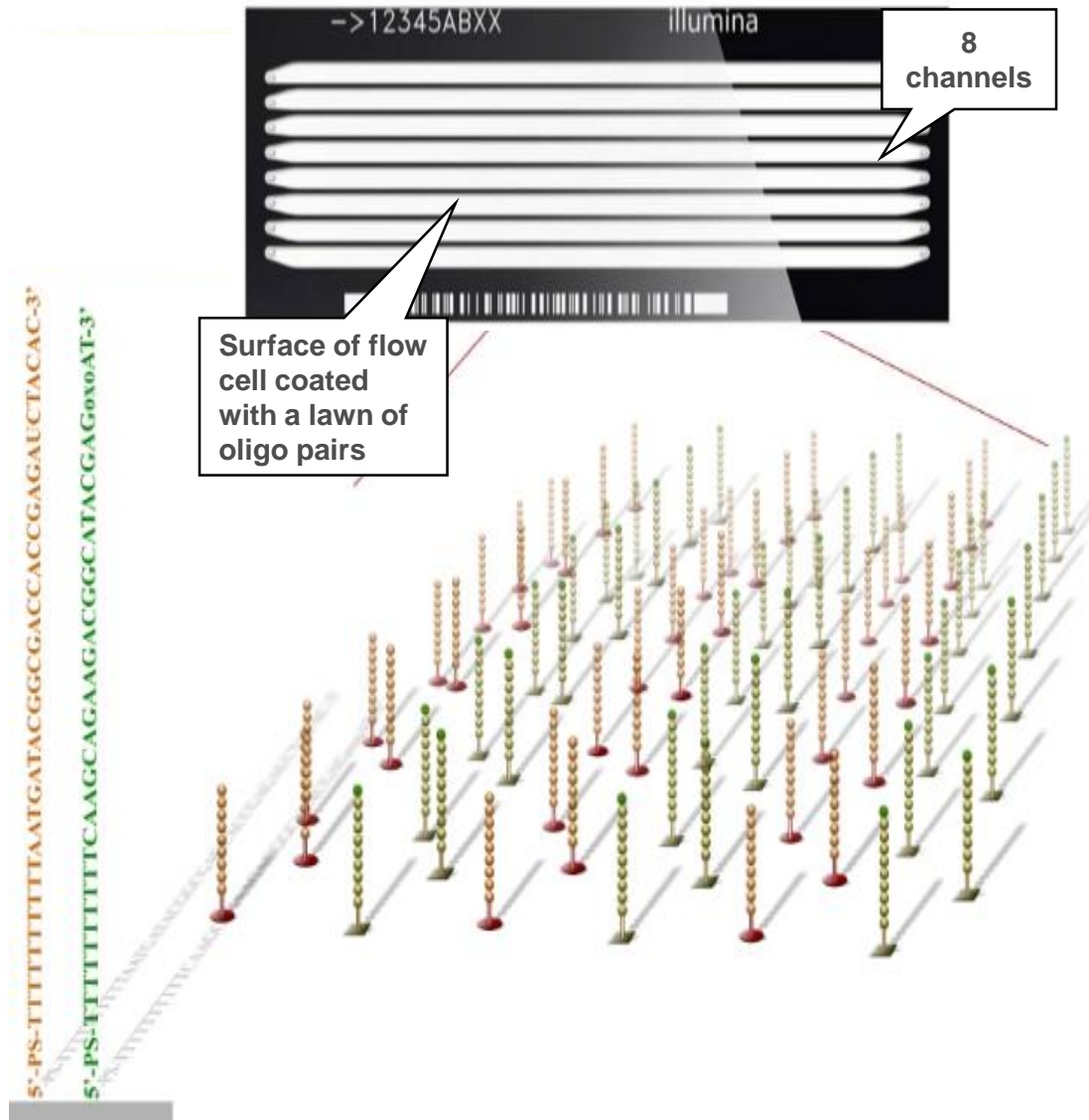
Cluster generation



Sequencing

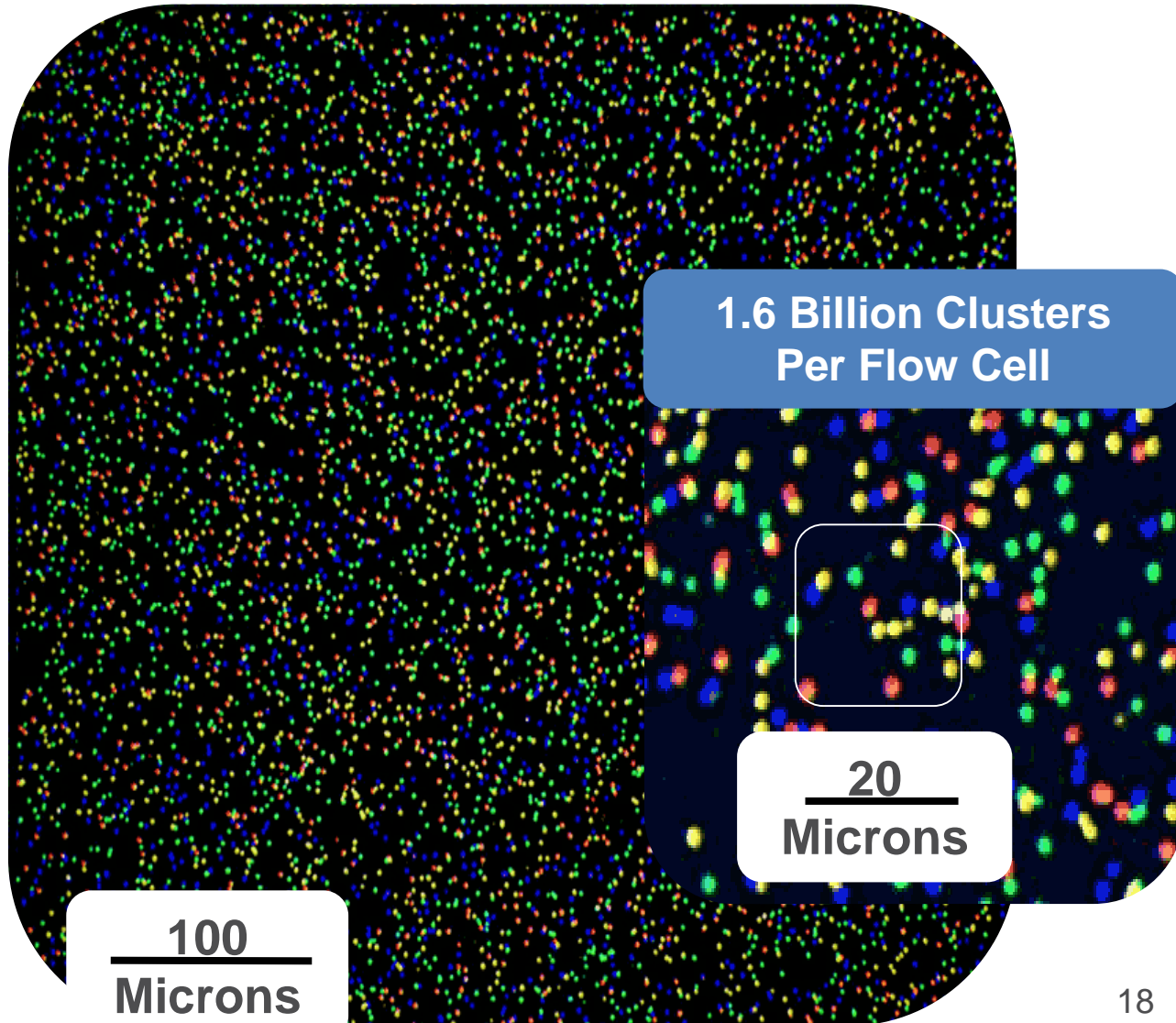


TruSeq Chemistry: Flow Cell

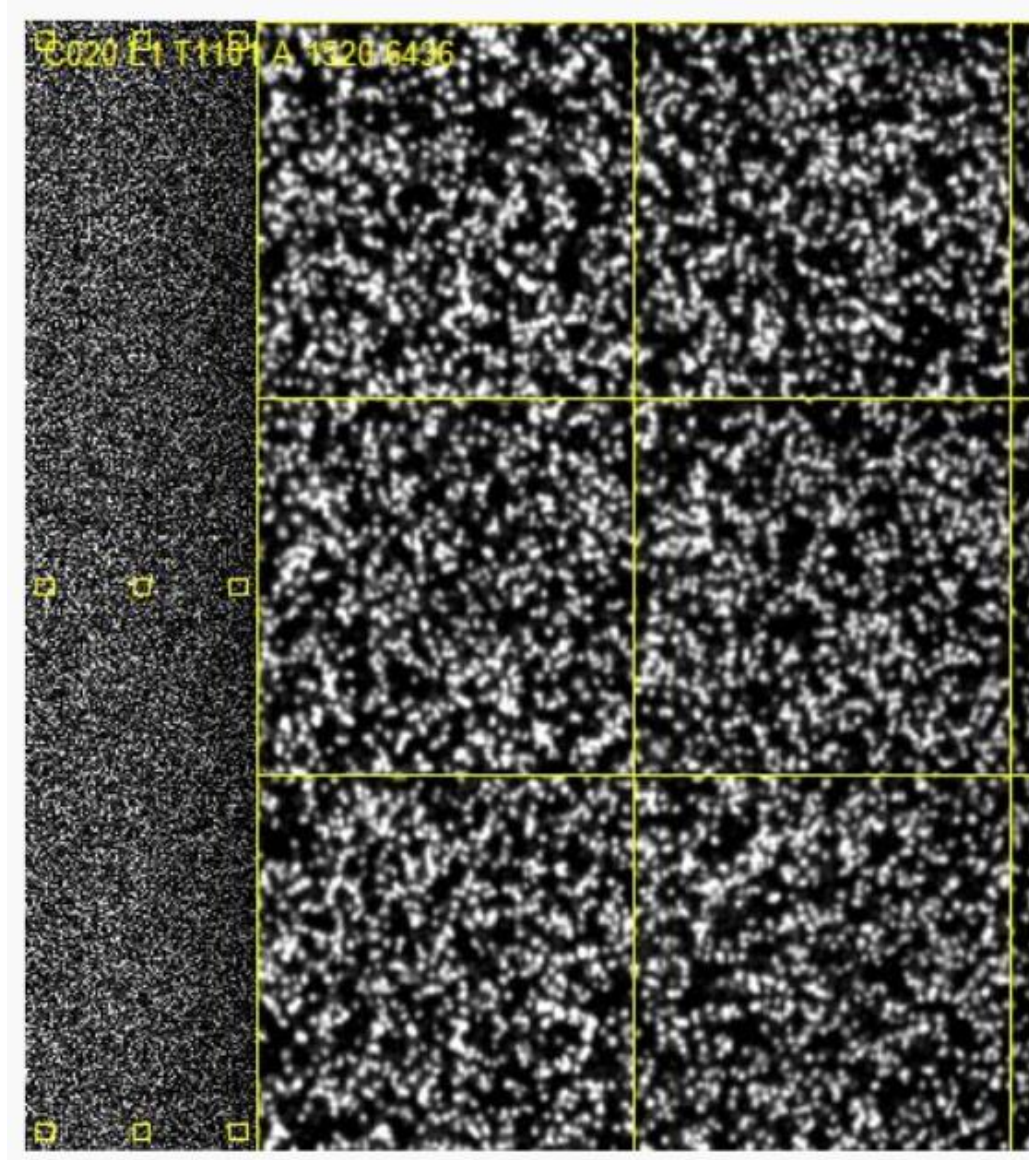
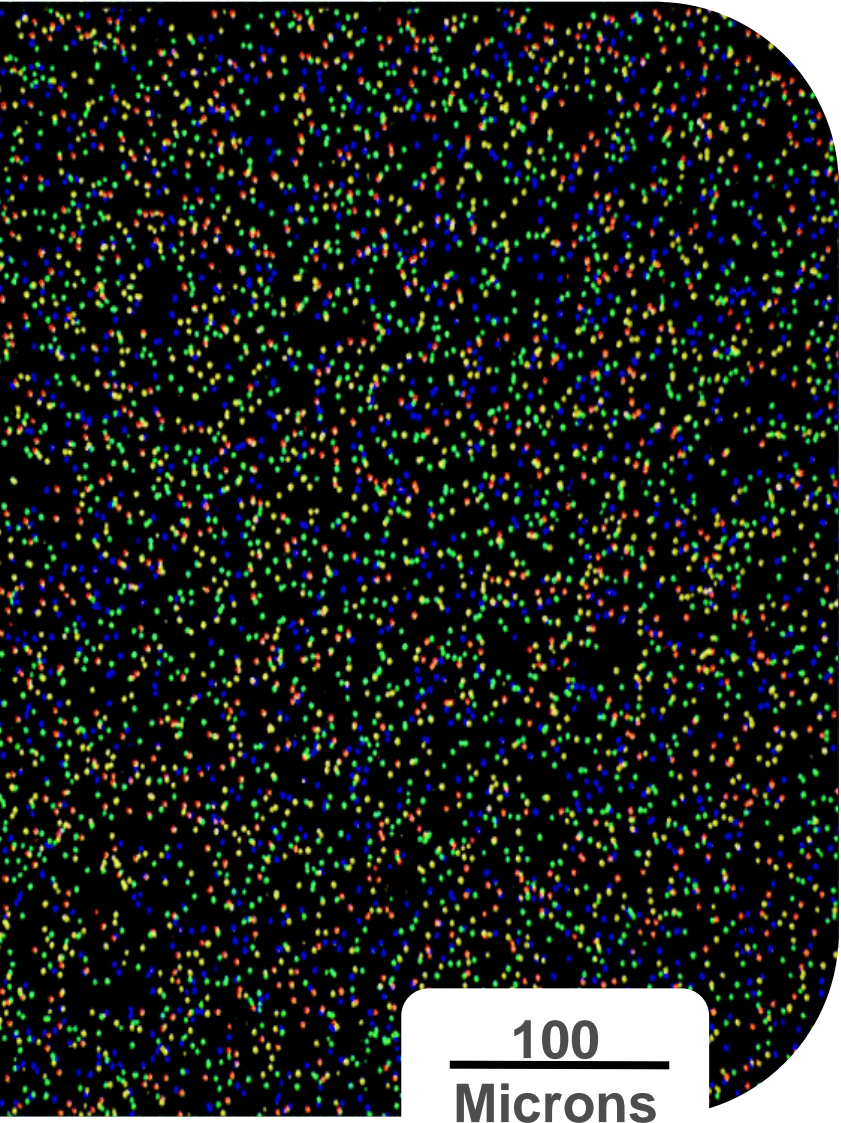


GCTGAAAC
ATCCCT
TCTGGGAA
GAAATT
TGTTGA
AAGGAG
TTTGGG
CGCCAG
TCCCAG
AATTGC
TCTCCA
AAGGCT
AATTGA
GCACAA
ATACCA
GCTTTT
TTZATA

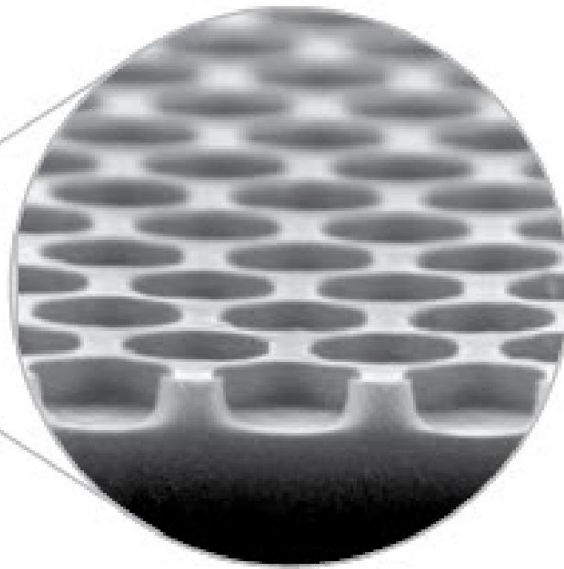
Sequencing



Sequencing



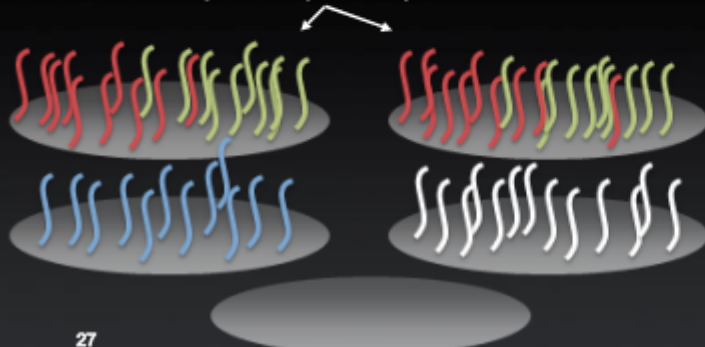
Hiseq 3000: 478 million nanowells per lane



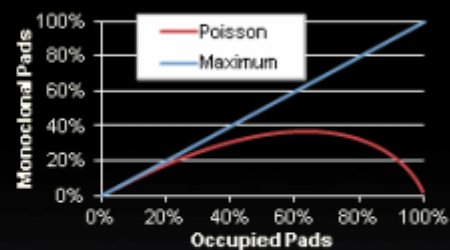
CONCEPTUAL CHALLENGE— BEATING POISSON

Amplification Phase

Polyclonal (non-PF) Pads



Maximizing Well Occupancy and Monoclonality

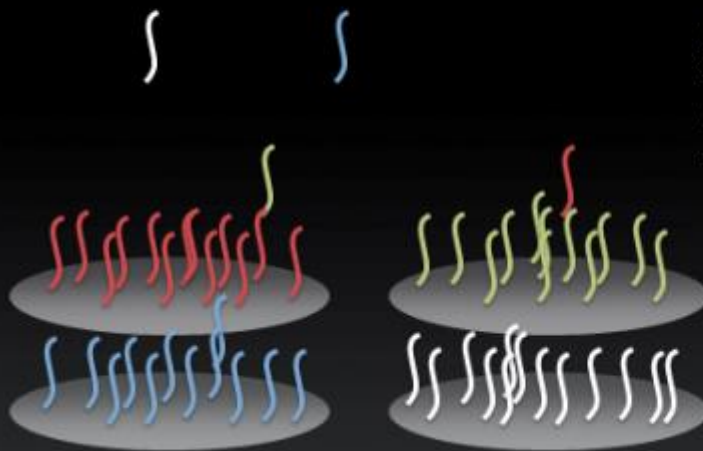


Poisson statistics limit max
monoclonal occupancy < 40%

Polyclonality rises as occupancy
increases

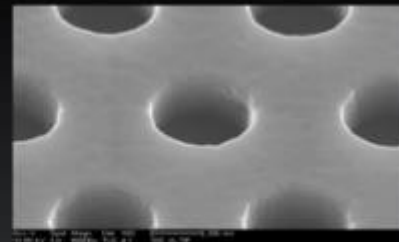
SIMULTANEOUS SEEDING AND AMPLIFICATION

Maximizing Well Occupancy and Monoclonality



Amplification occurs at rate \gg faster
than seeding rate

Templates excluded from occupied
wells



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

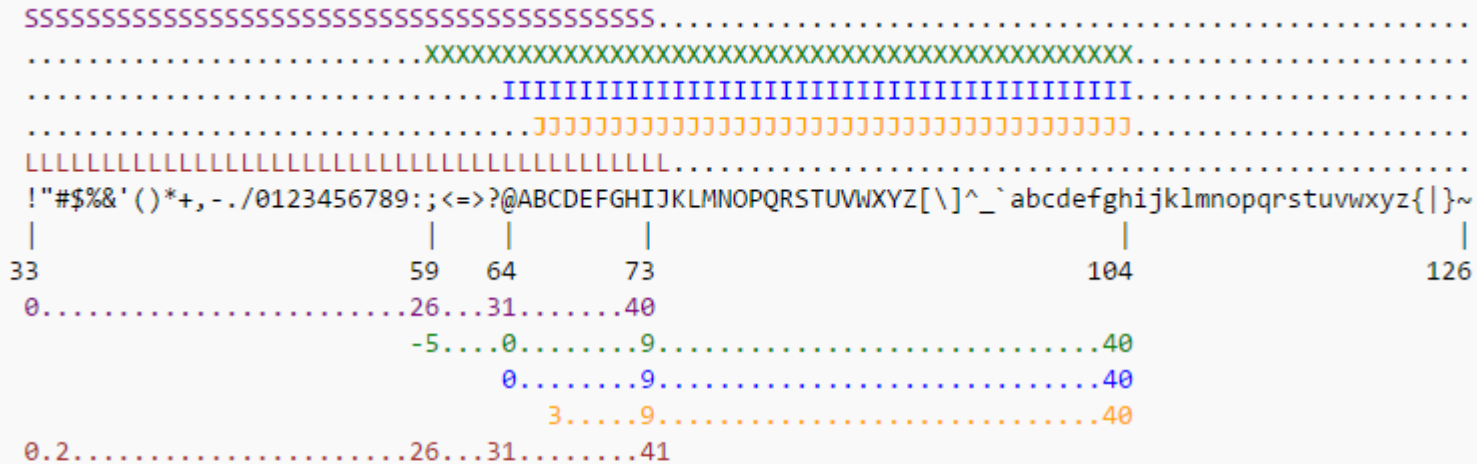
Phasing & Pre-Phasing
problems




```

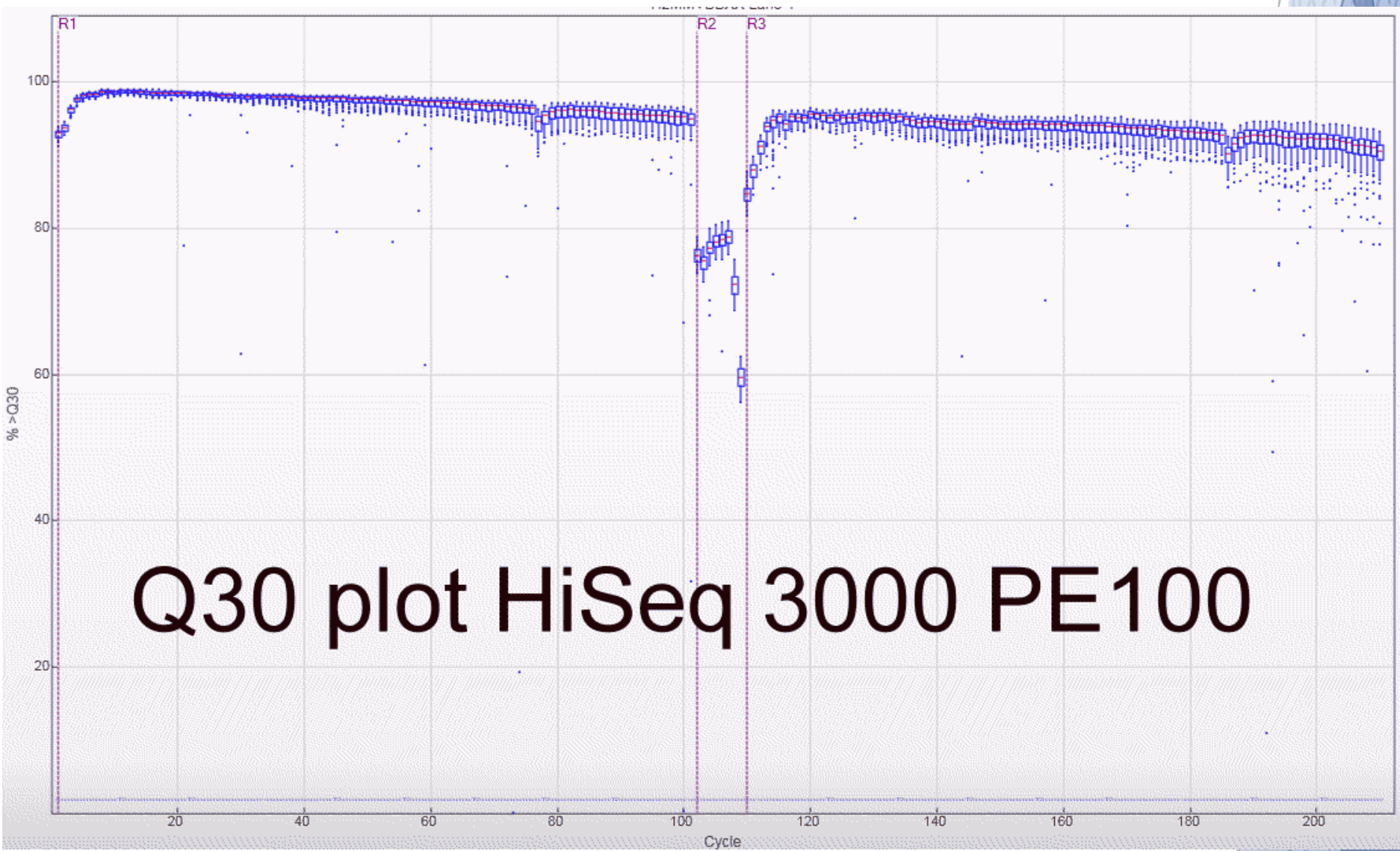
@M02034:265:000000000-AN3L2:1:2102:8707:16197 2:N:0:85
GATGAACATAATAAGCAATGACGGCAGCAATAAACTCAACAGGAGCAGGA
+
AAAAAFFFFFFFFFGFFGFFBE5GEAAAEDCFDFAEG5CFGHFGGFEGHHHG

```



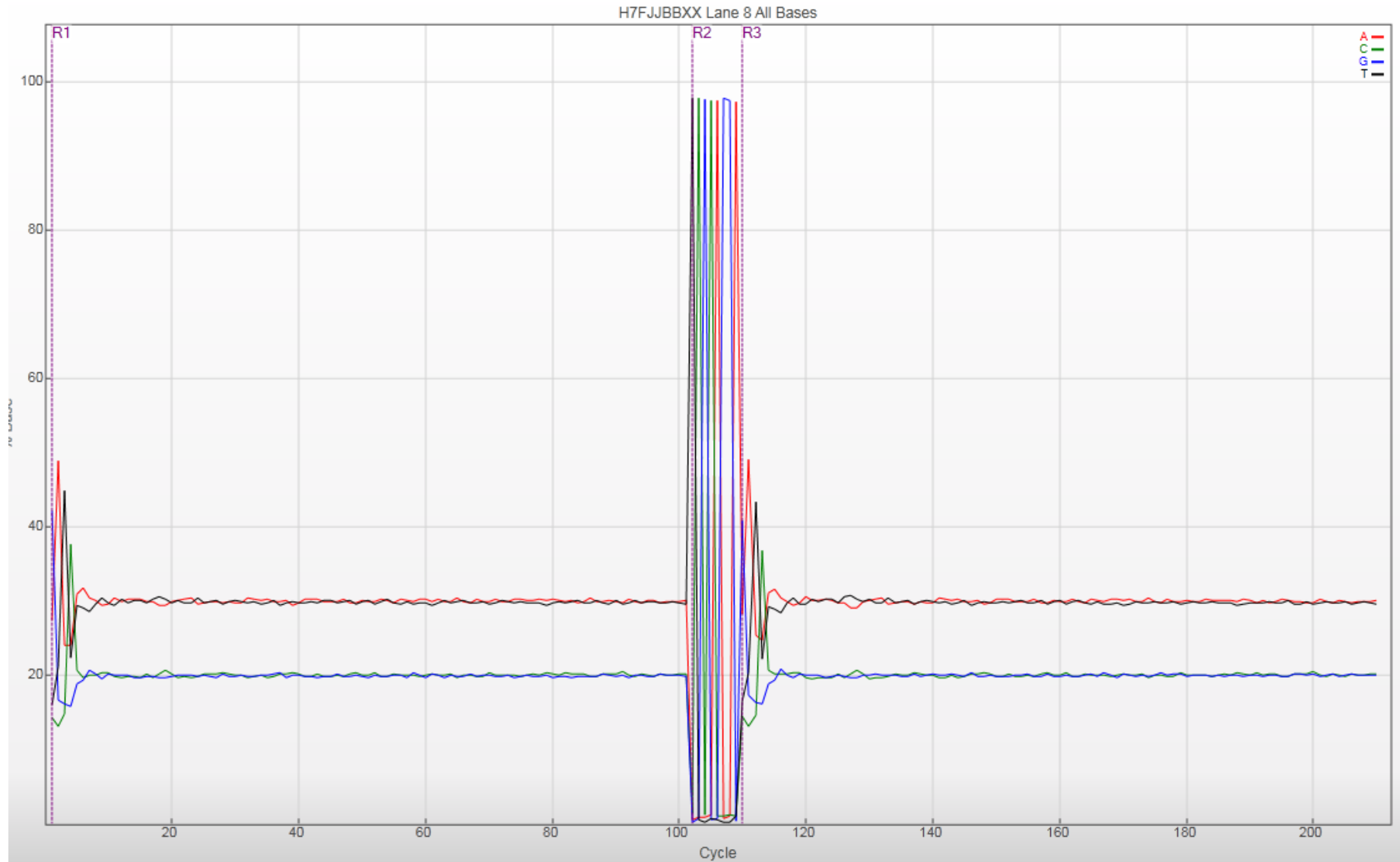
- S - Sanger Phred+33, raw reads typically (0, 40)
- X - Solexa Solexa+64, raw reads typically (-5, 40)
- I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)
- J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)
with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
(Note: See discussion above).
- L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)

Illumina SAV viewer

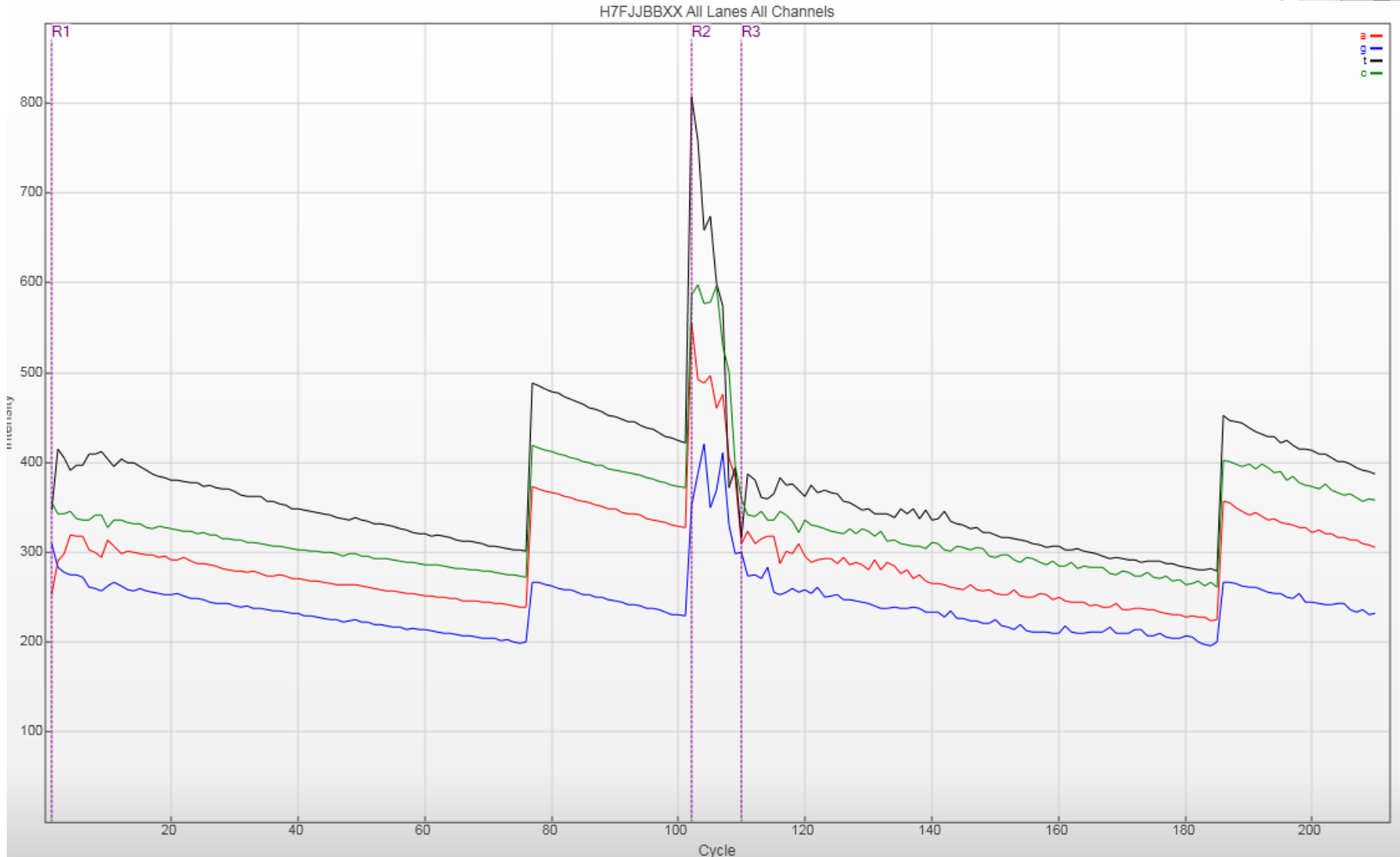


Q30 plot HiSeq 3000 PE100

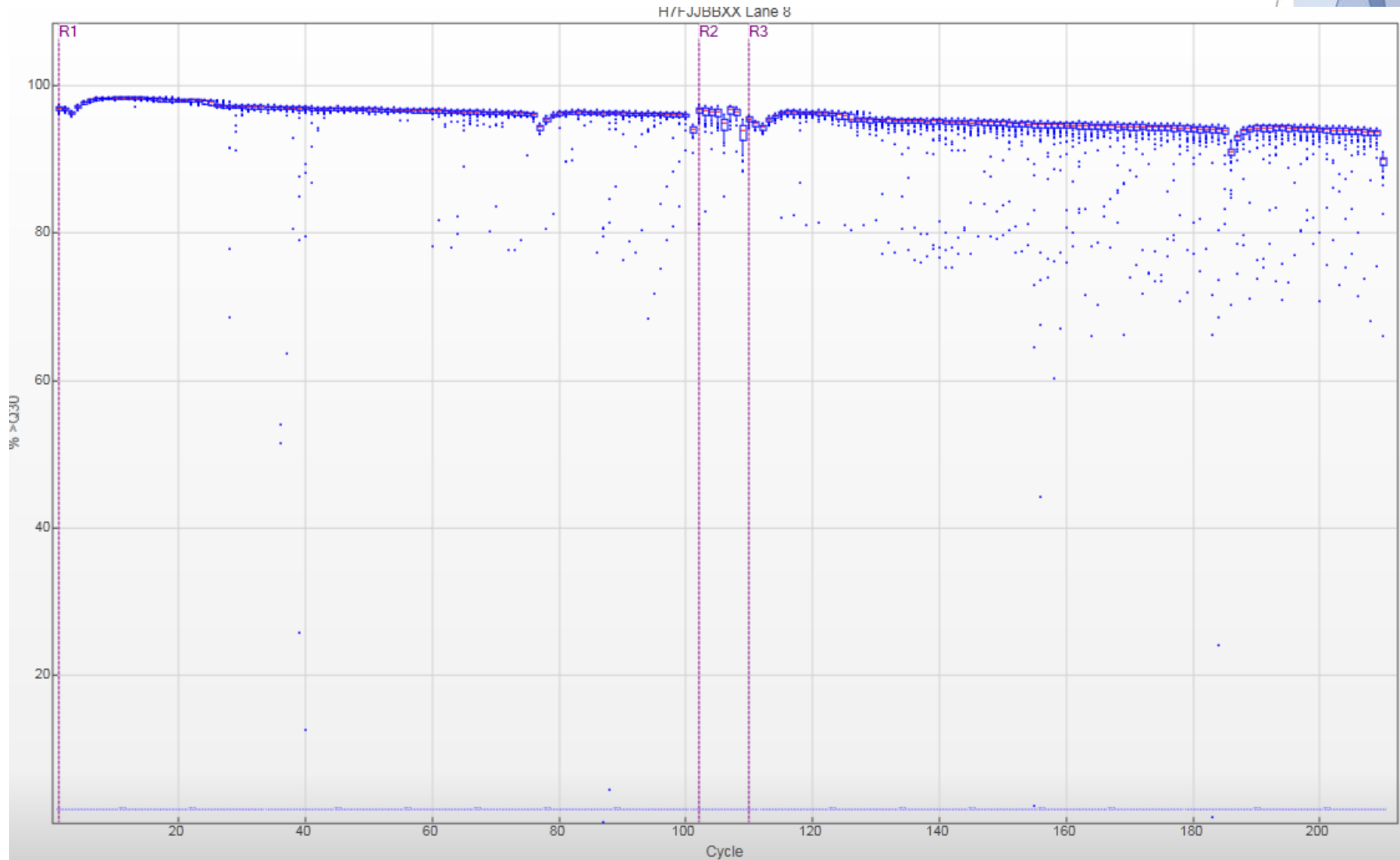
base composition



fluorescence intensity

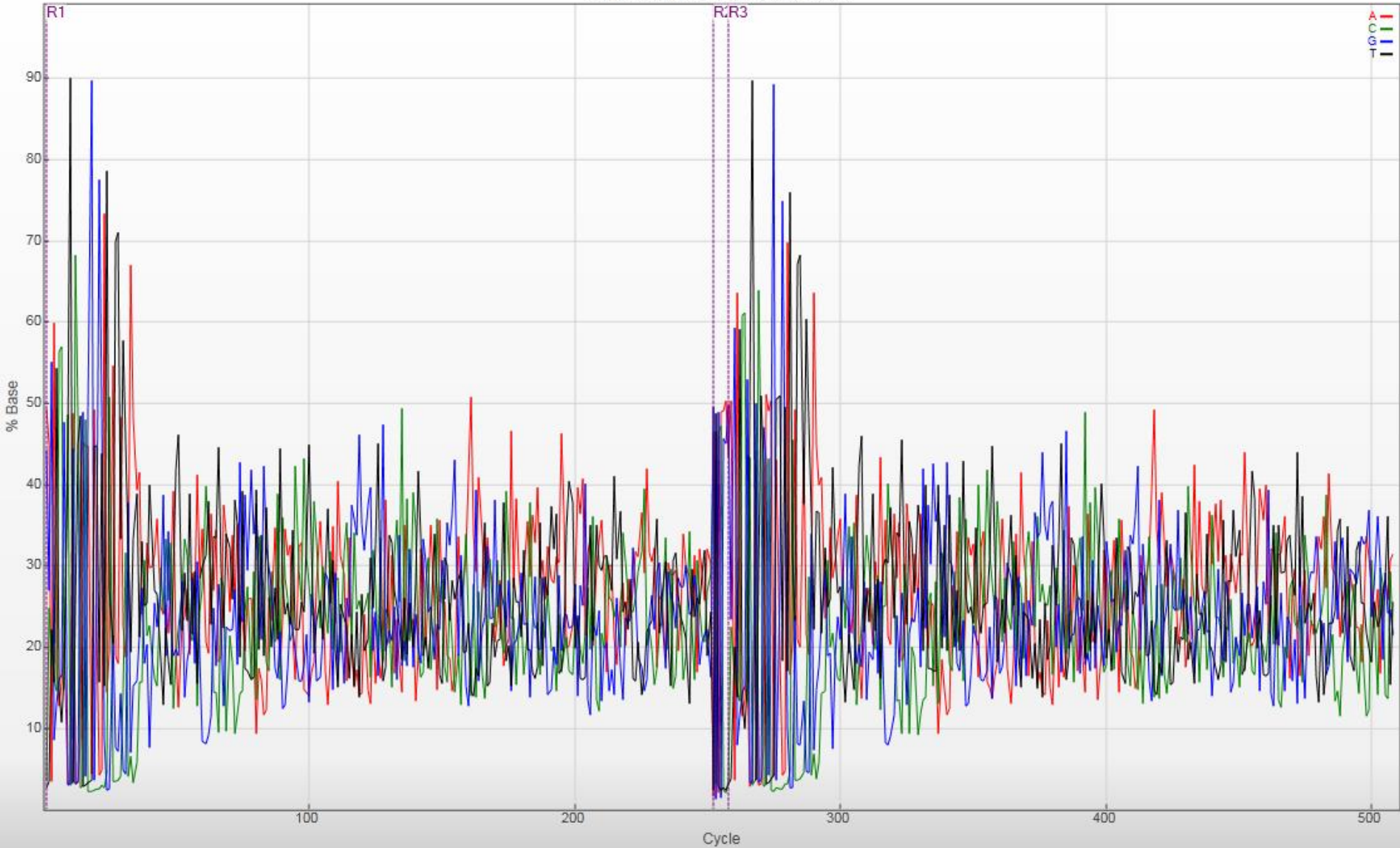


fluorescence intensity

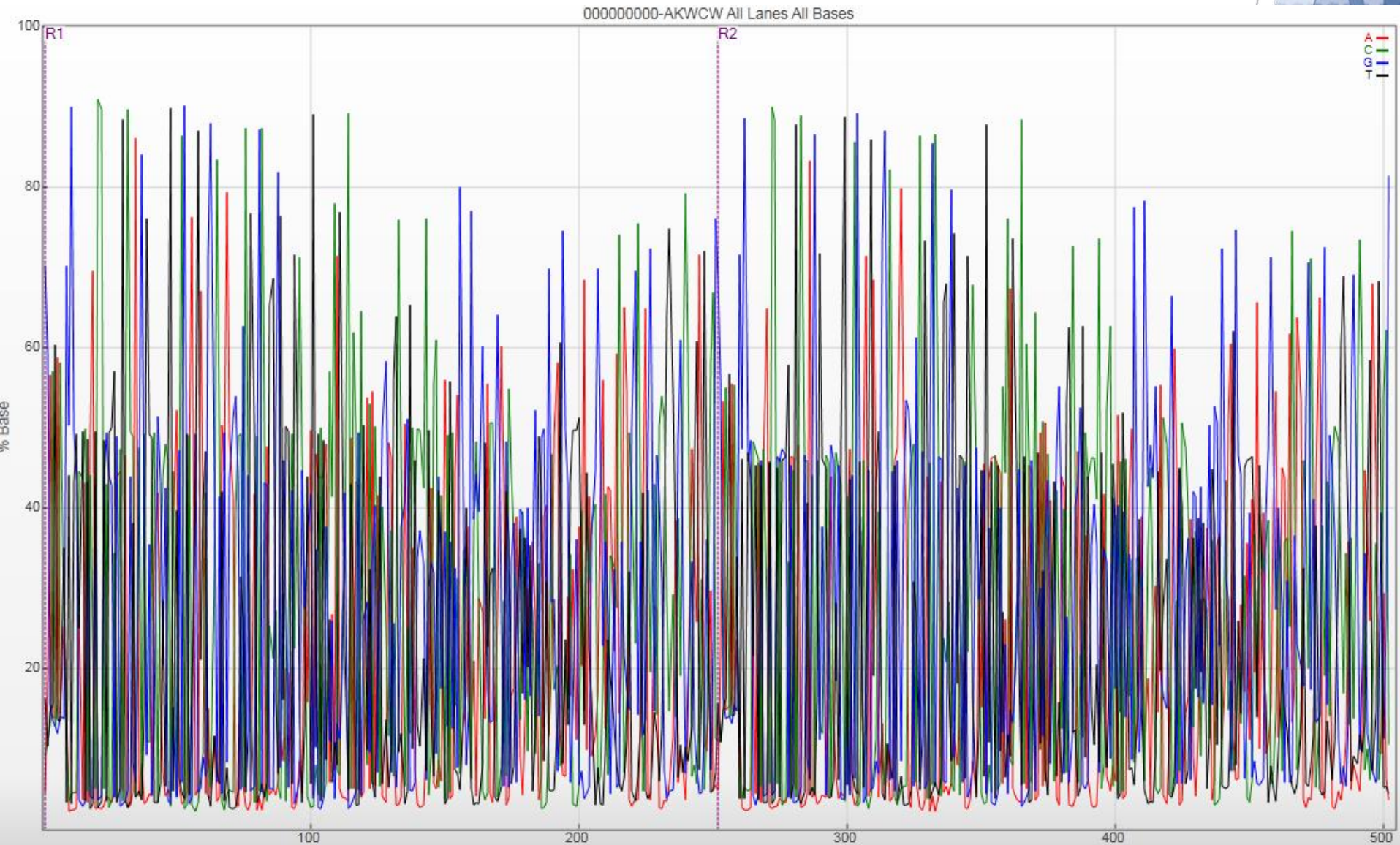


amplicon mix

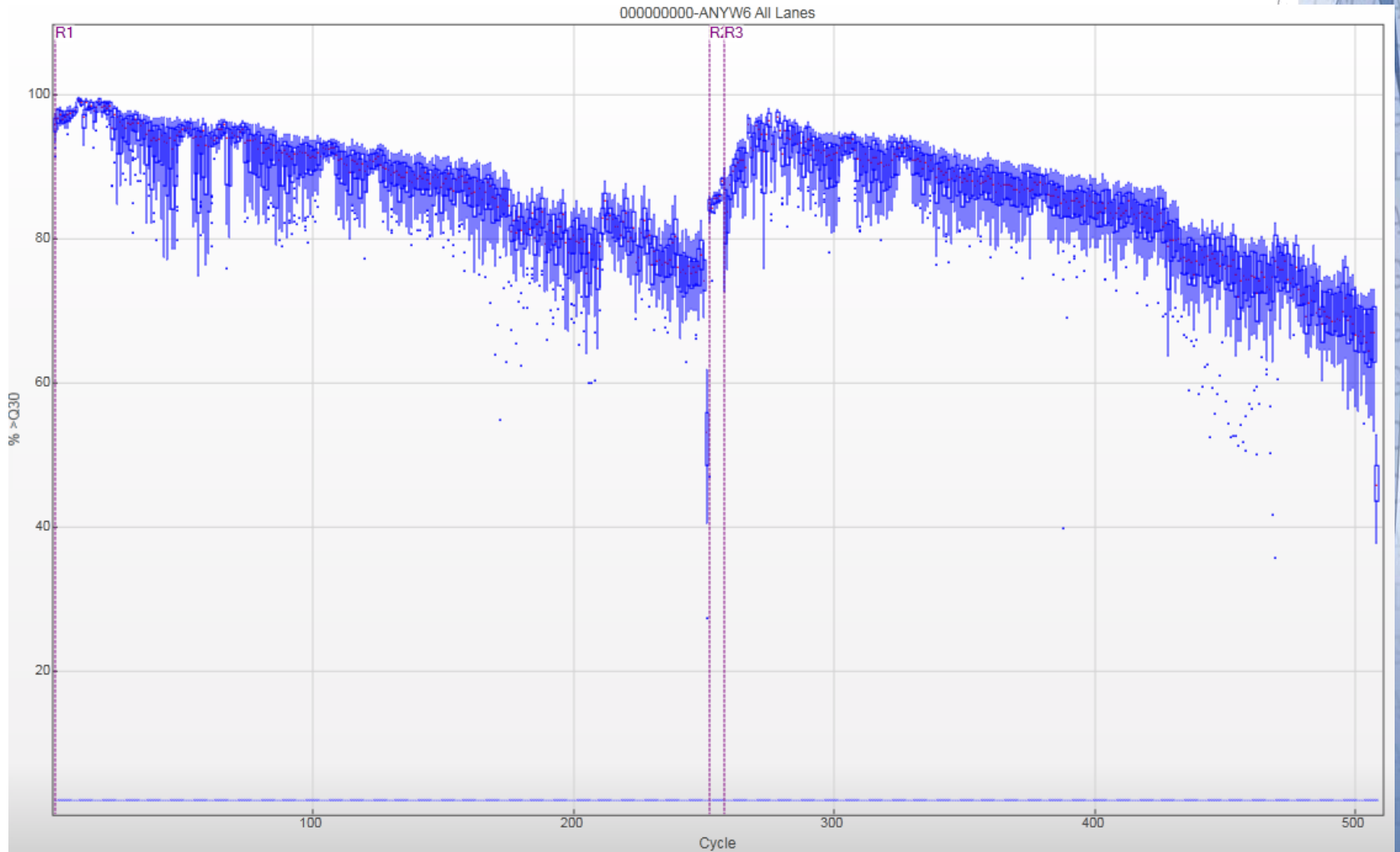
00000000-ANYW6 All Lanes All Bases



amplicon



amplicon mix Q30



FASTQC

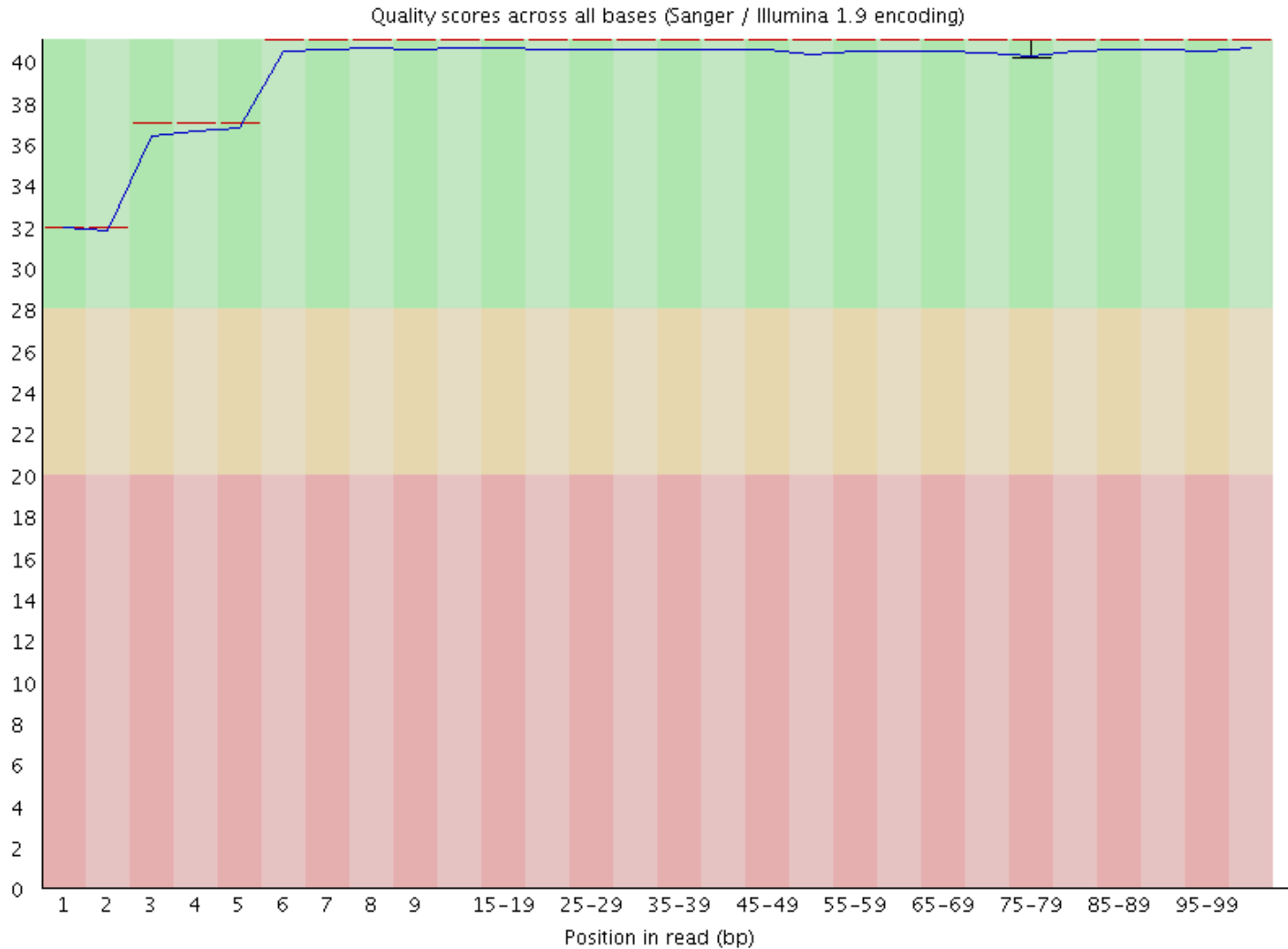


Basic Statistics

| 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 |

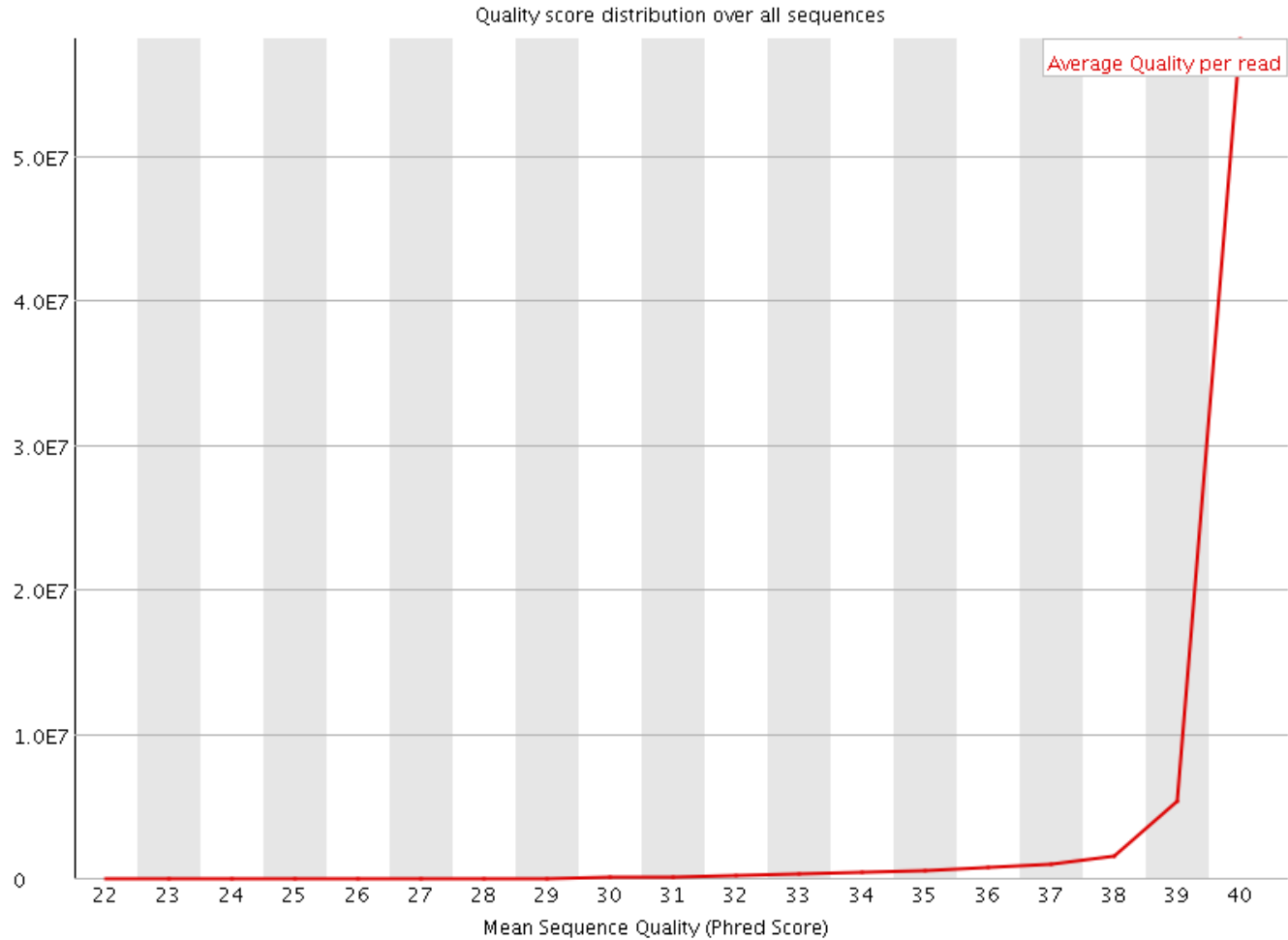
FASTQC

✔ Per base sequence quality



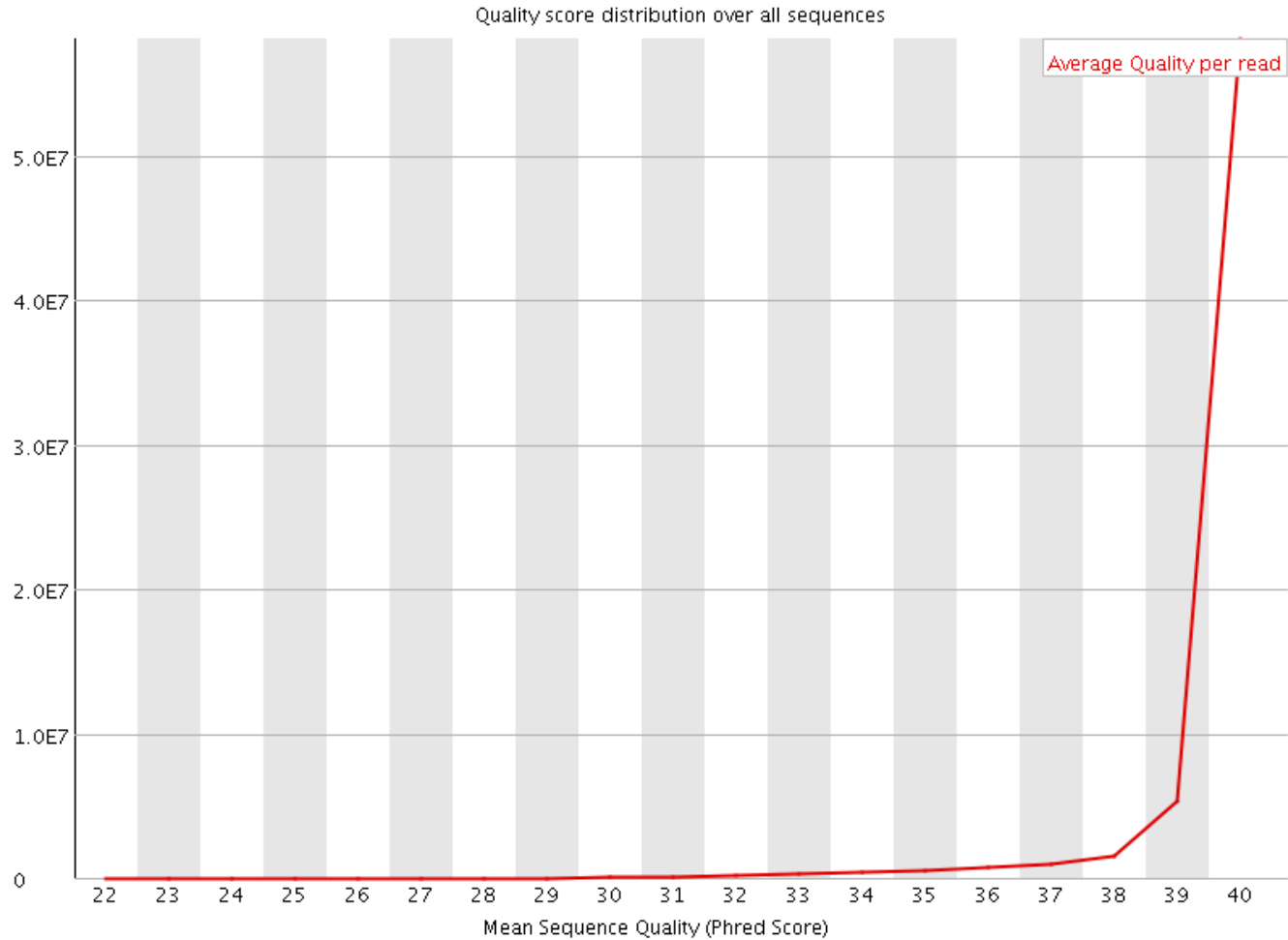
FASTQC

✔ Per sequence quality scores



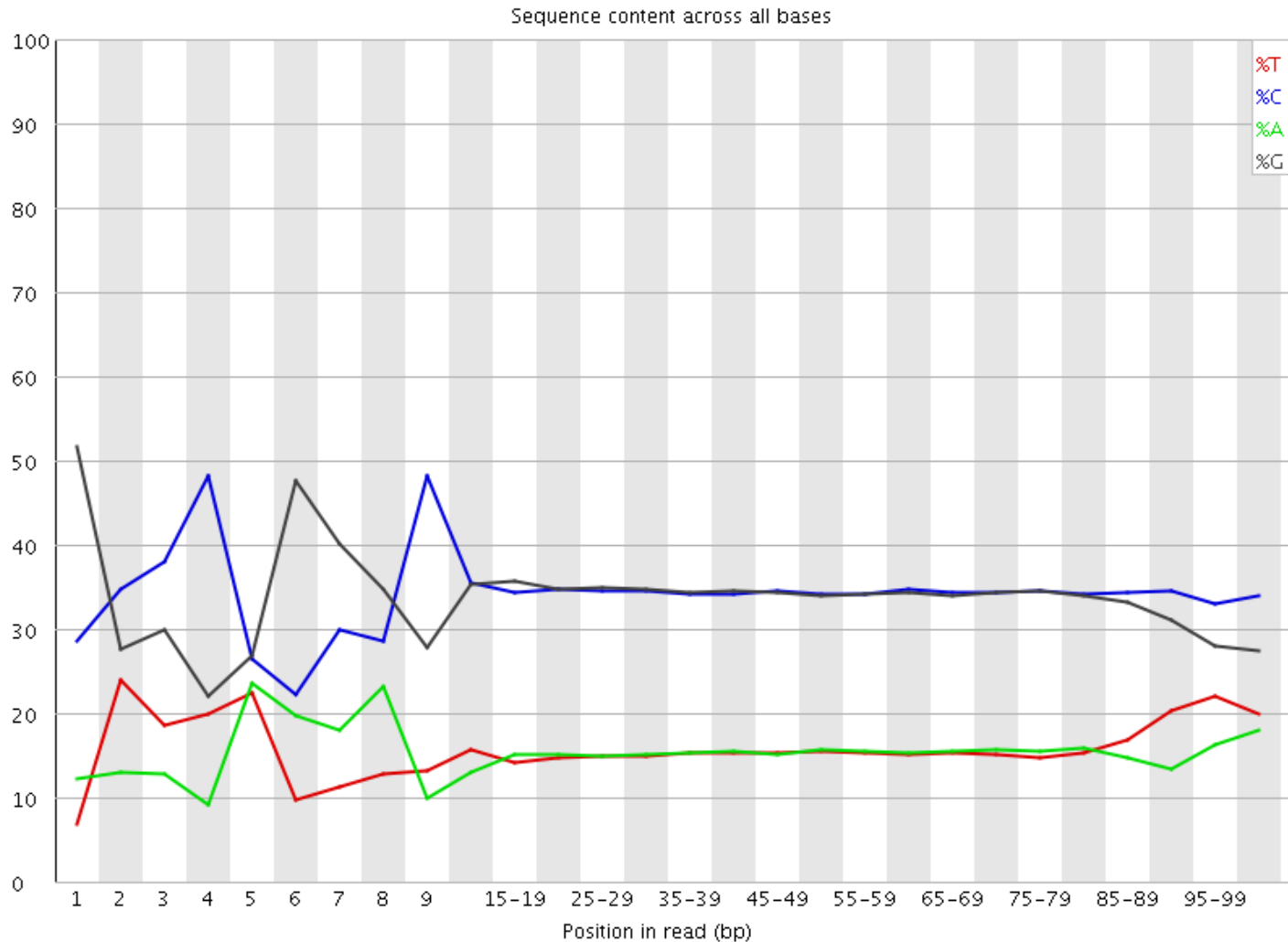
FASTQC

✔ Per sequence quality scores



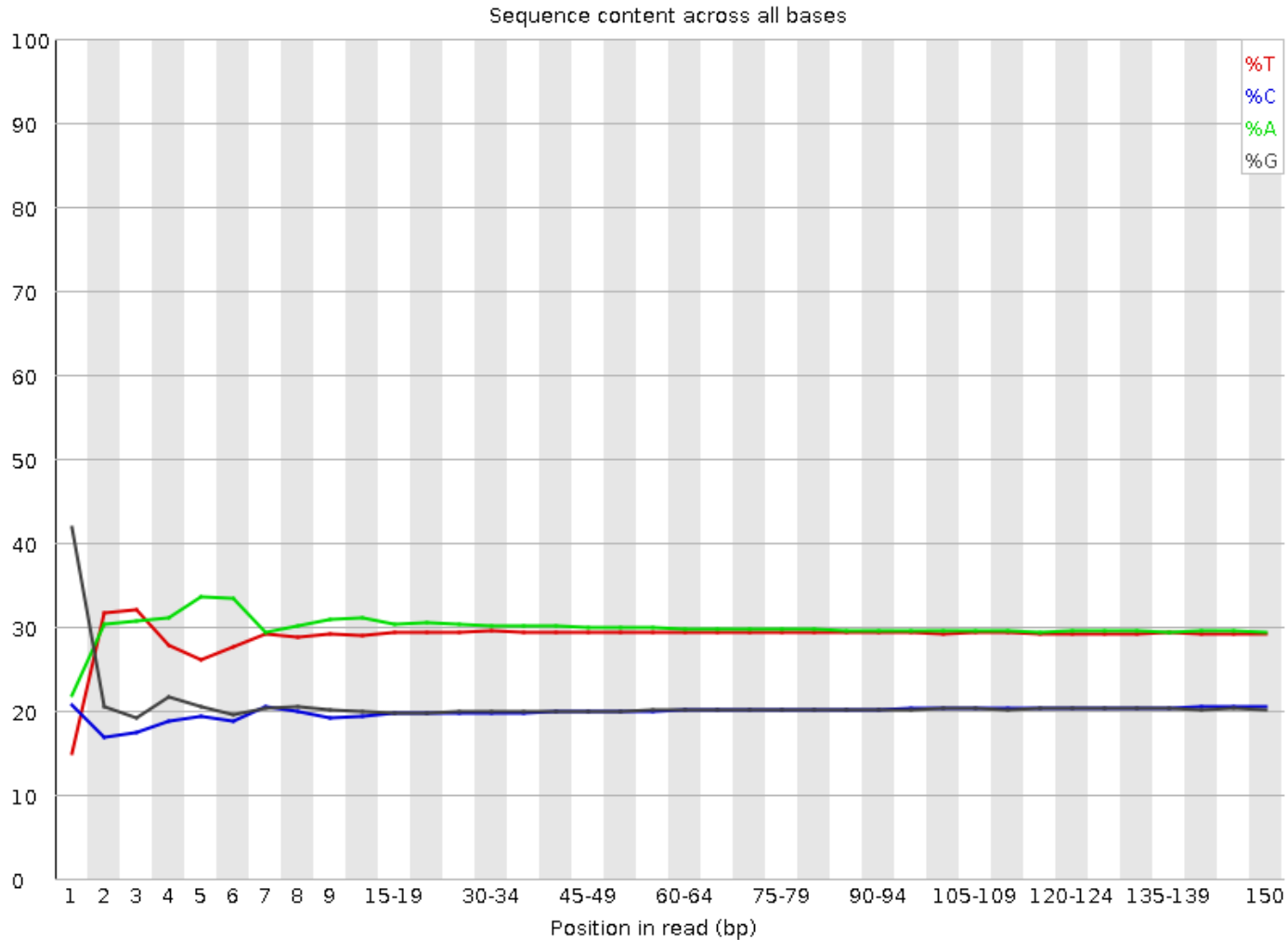
FASTQC - Nextera

✖ Per base sequence content



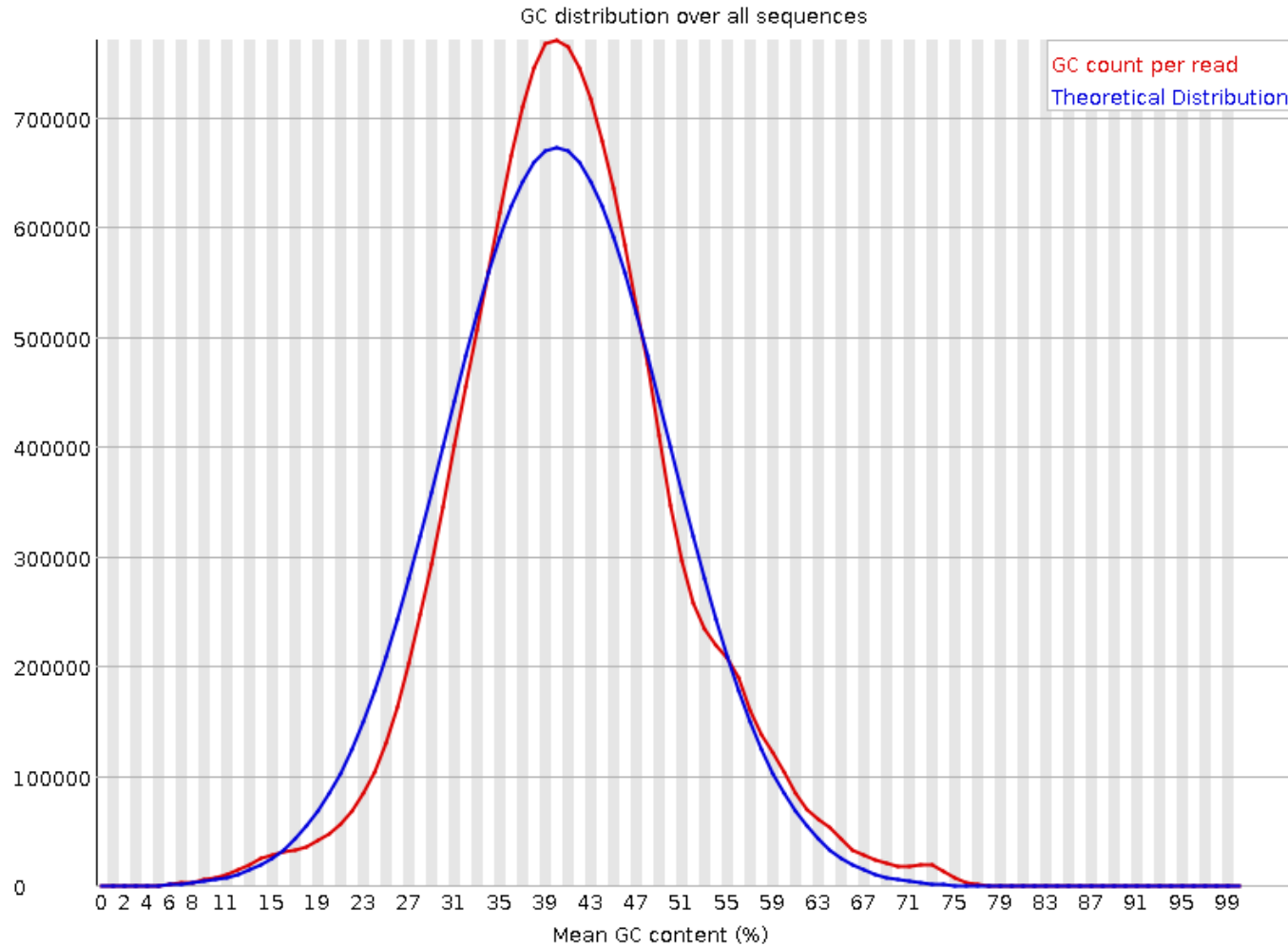
FASTQC

✖ Per base sequence content



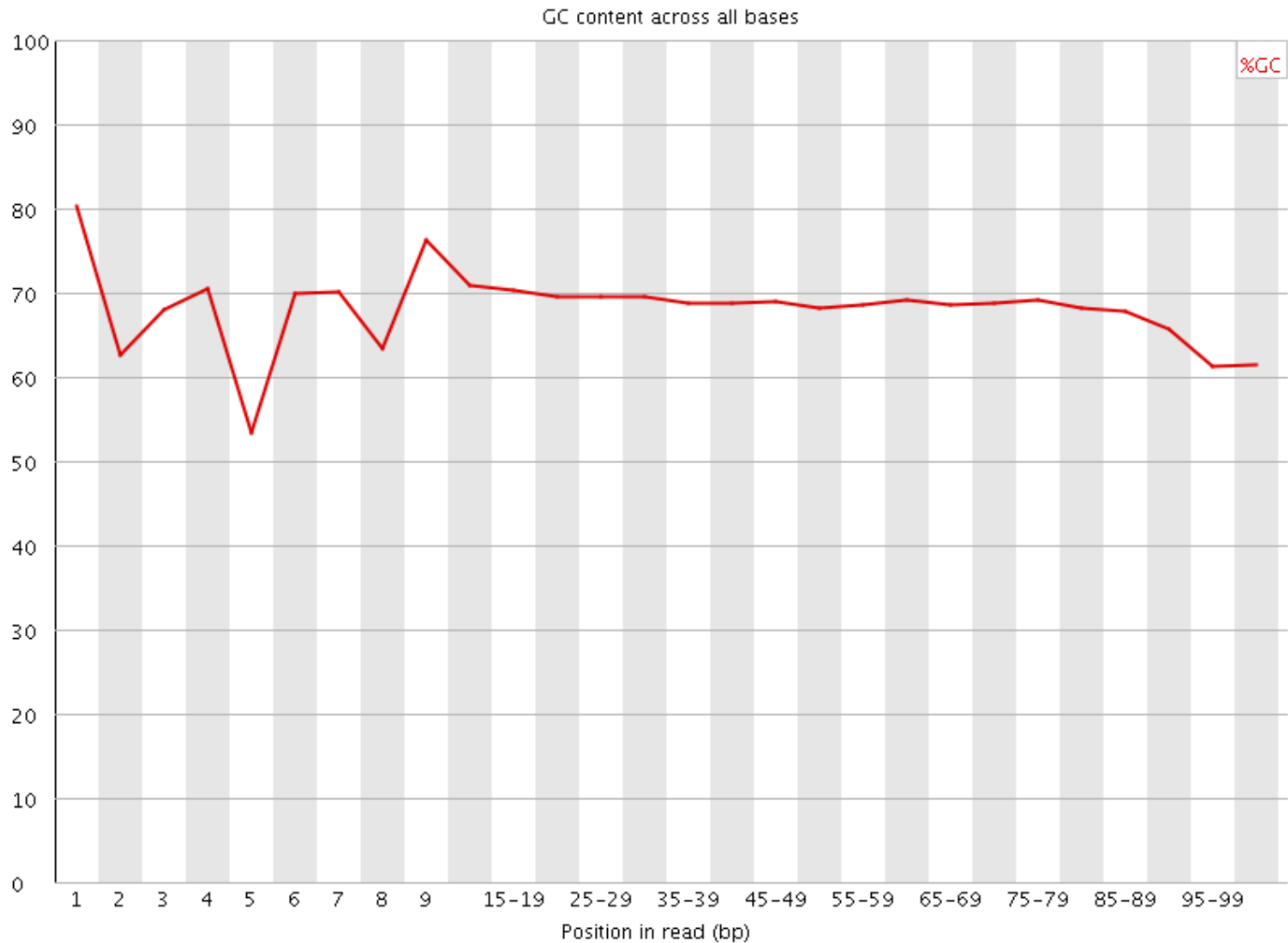
FASTQC

✔ Per sequence GC content



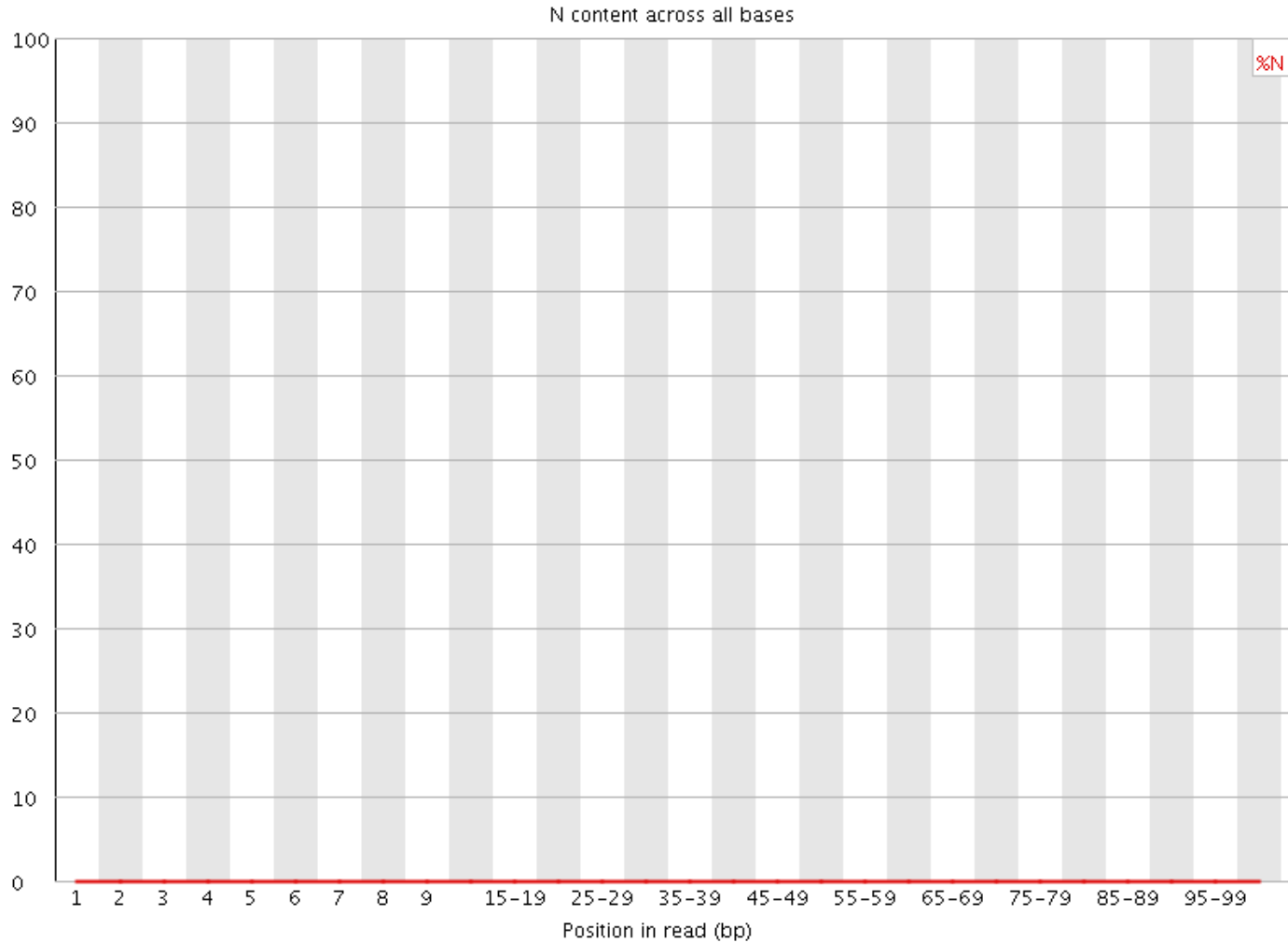
FASTQC

❌ Per base GC content



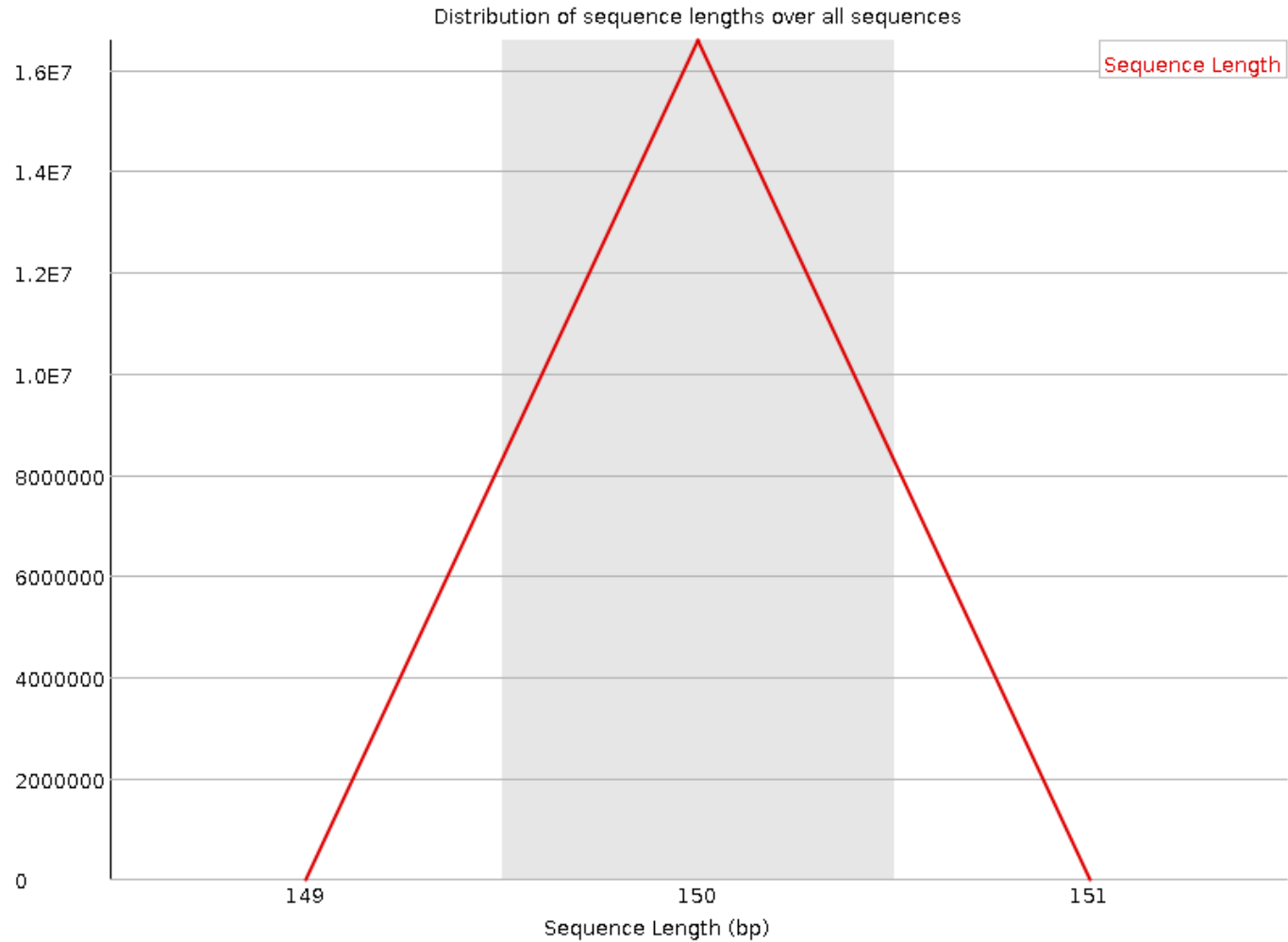
FASTQC

✔ **Per base N content**



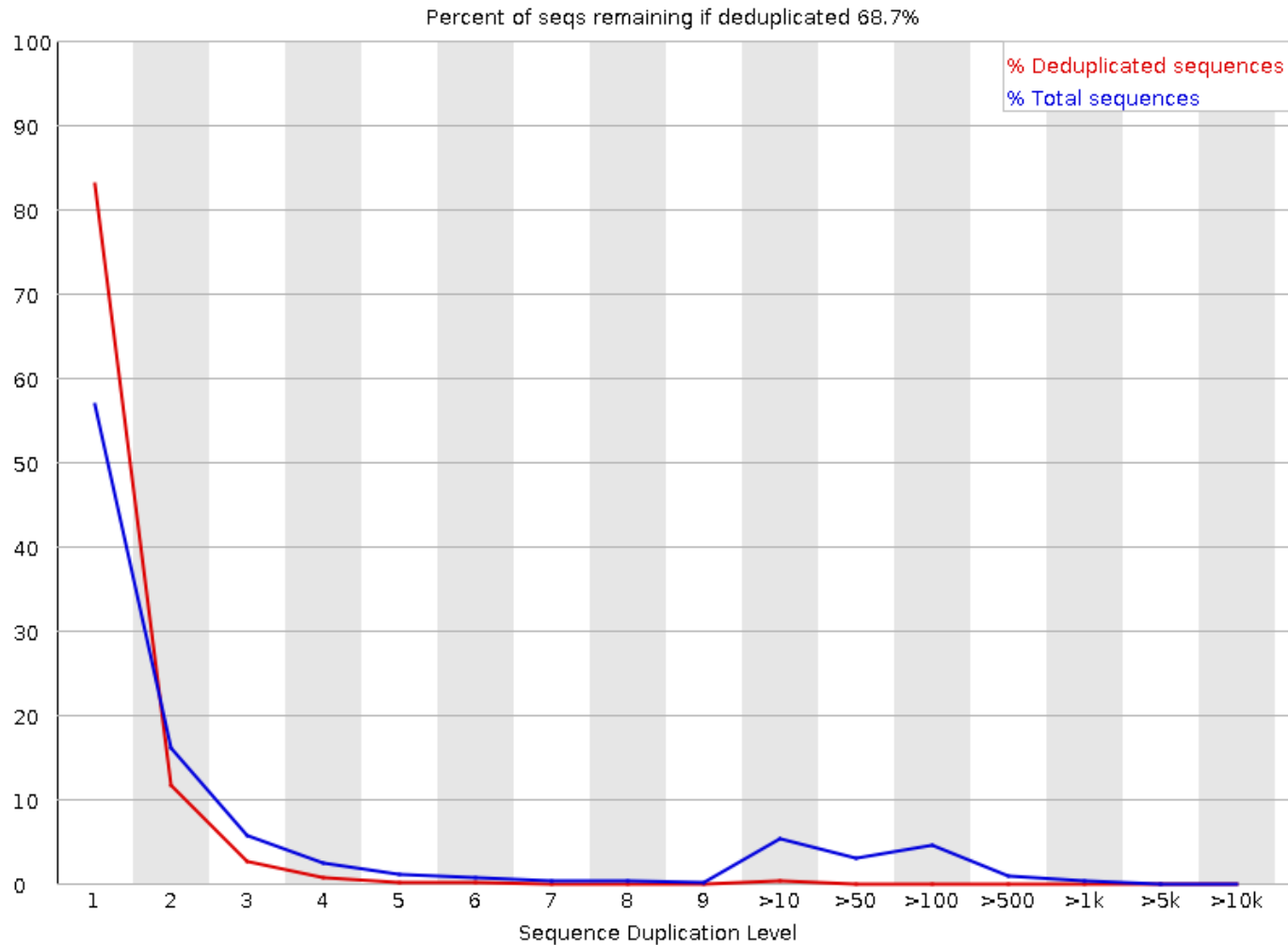
FASTQC

✔ Sequence Length Distribution



FASTQC

⚠ Sequence Duplication Levels

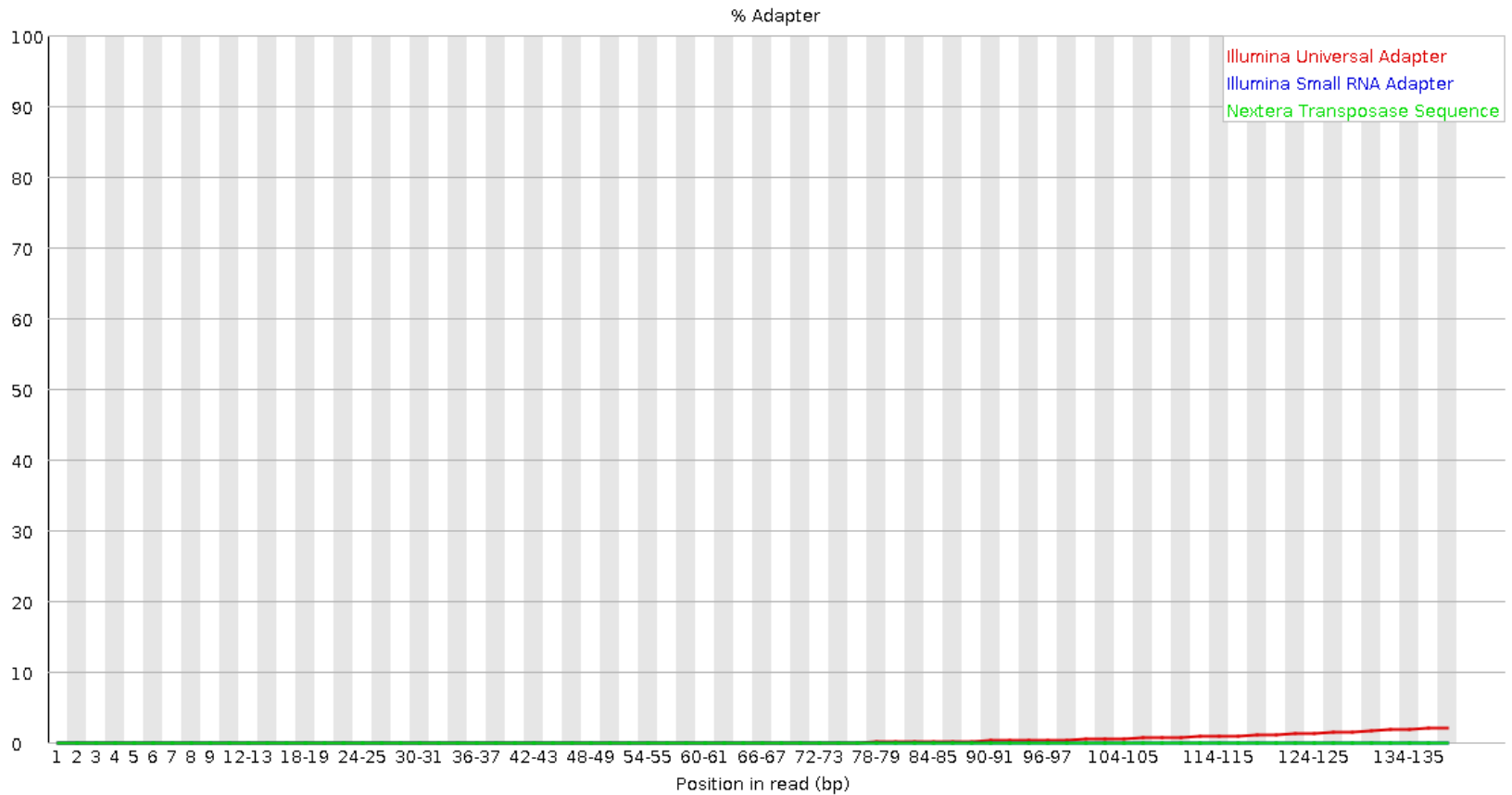


FASTQC

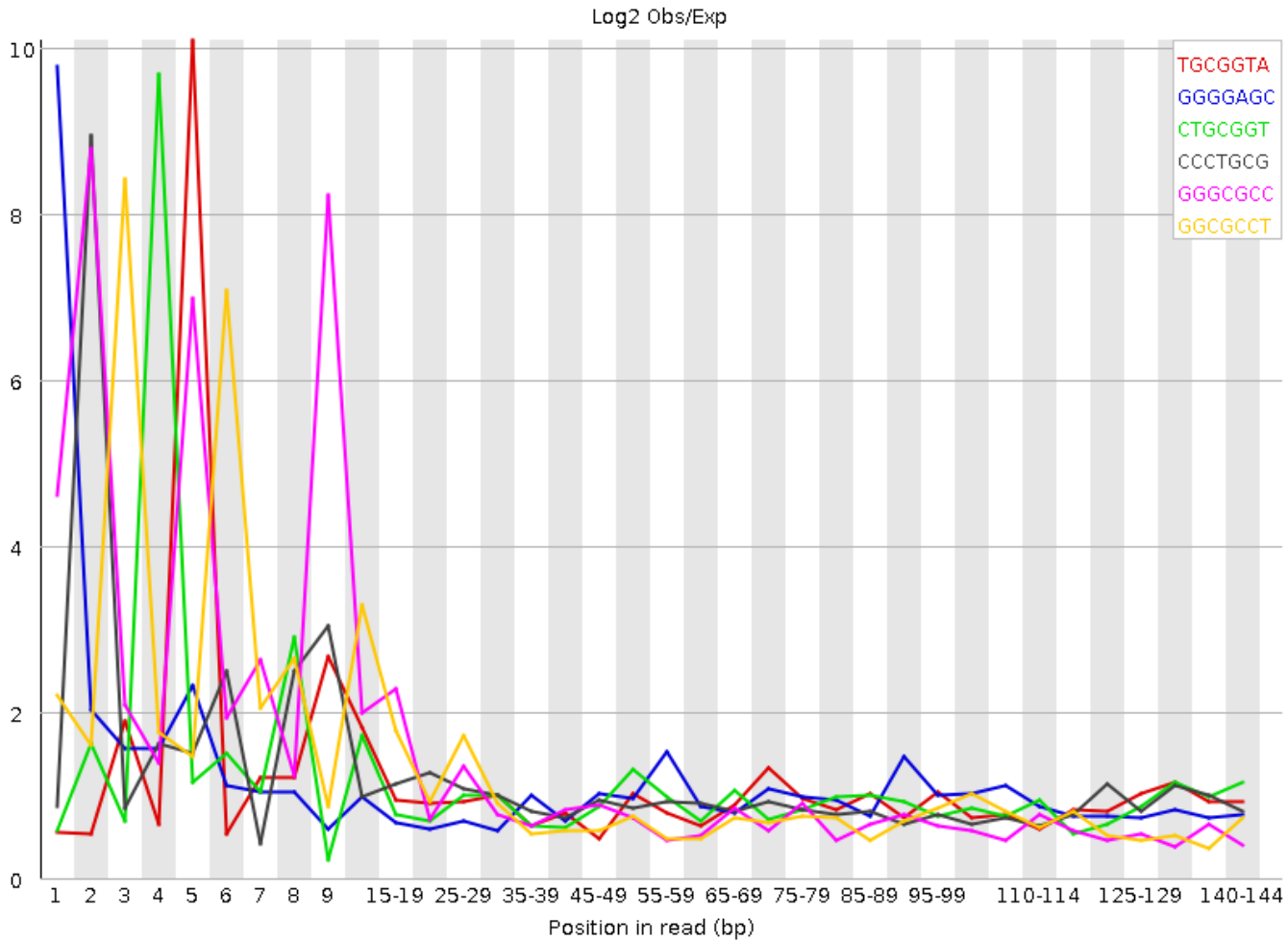
✔ Overrepresented sequences

No overrepresented sequences

✔ Adapter Content

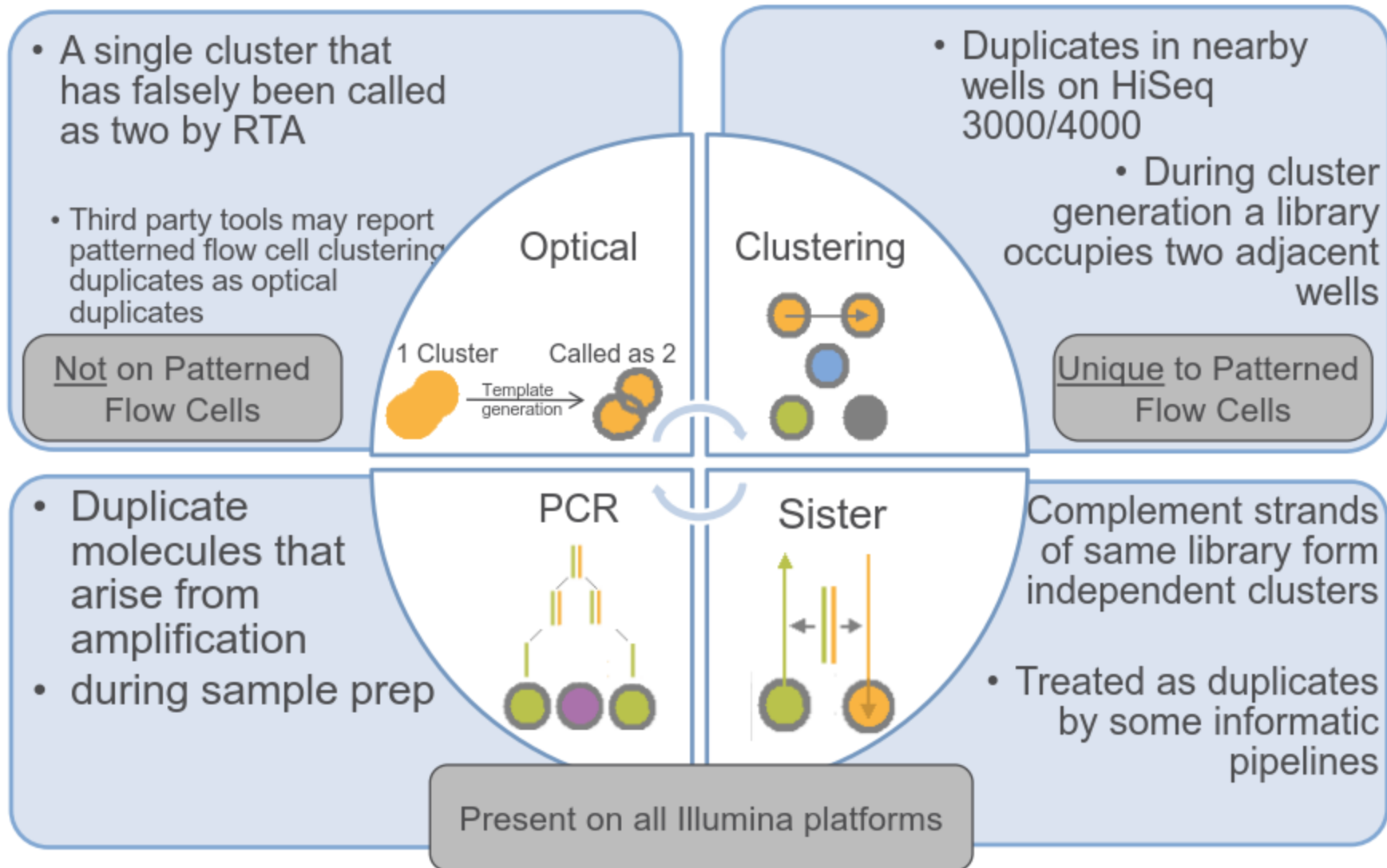


Kmer Content



| 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 |

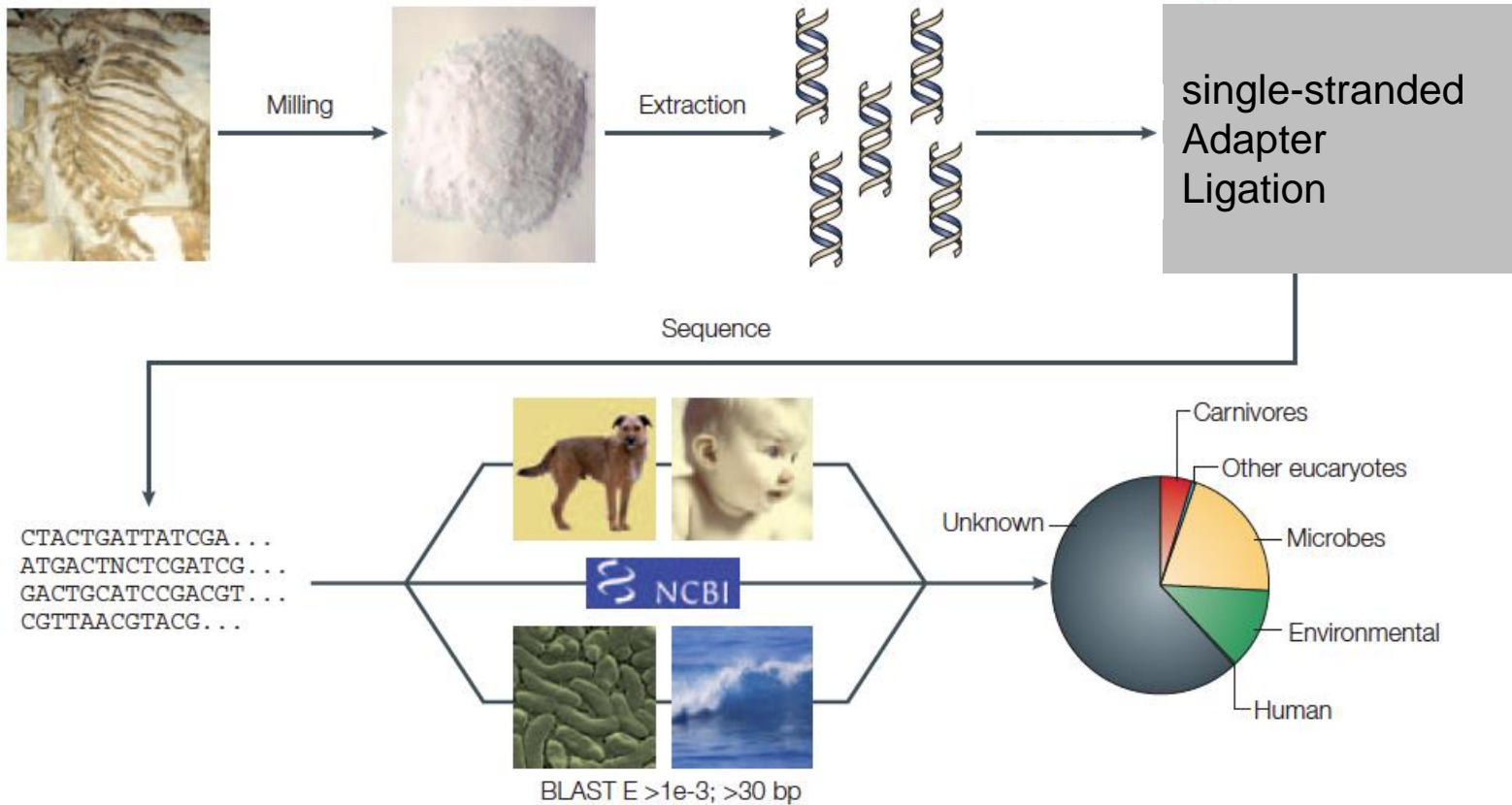
A Review of Sequencing Duplicate Types



**“If you can put adapters on it,
we can sequence it!”**



Know your sample



No need to be scared of HTS

UC Davis Center for Plant Diversity/Herbarium

- The Herbarium archives contain over 300,000 dried specimens.
- Search for **Grapevine Red Blotch-Associated Virus**
- Virus traces found by PCR



Maher Al Rwahnih
UCD Plant Foundation
Plant Services

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

Optional: PCR-free libraries

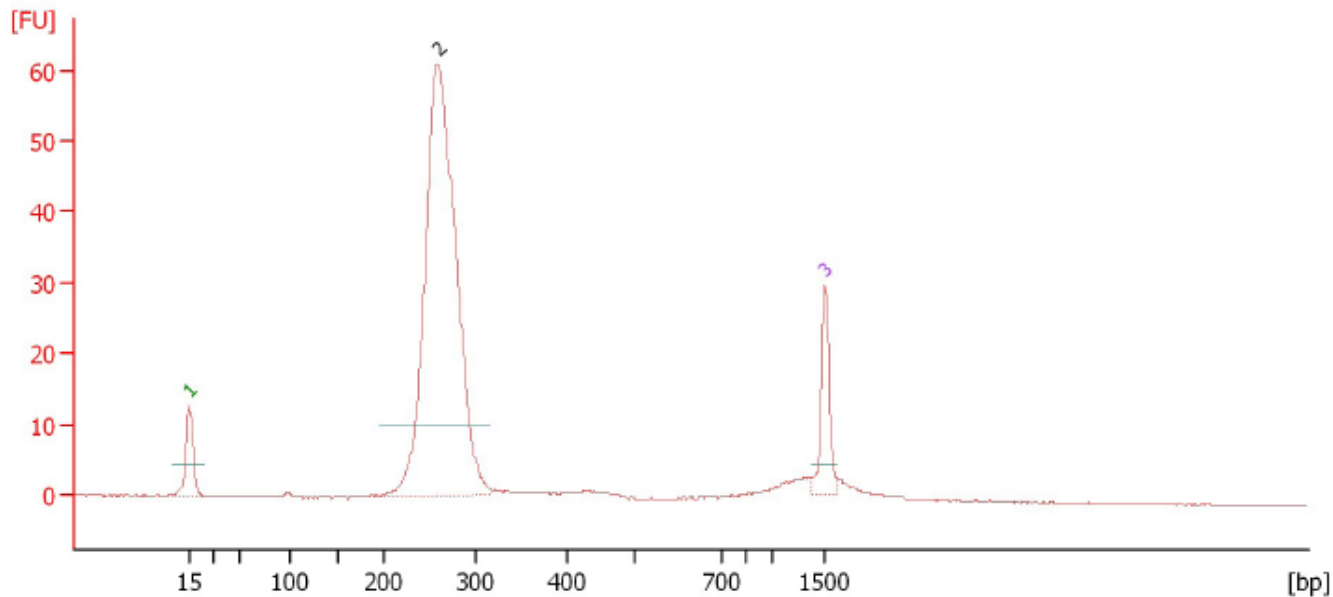
- PCR-free library:
 - if concentration allows
 - Reduction of PCR bias against e.g. GC rich or AT rich regions, especially for metagenomic samples

OR

- Library enrichment by PCR:
 - Ideal combination: high input and low cycle number; low-bias polymerase



Library QC by Bioanalyzer

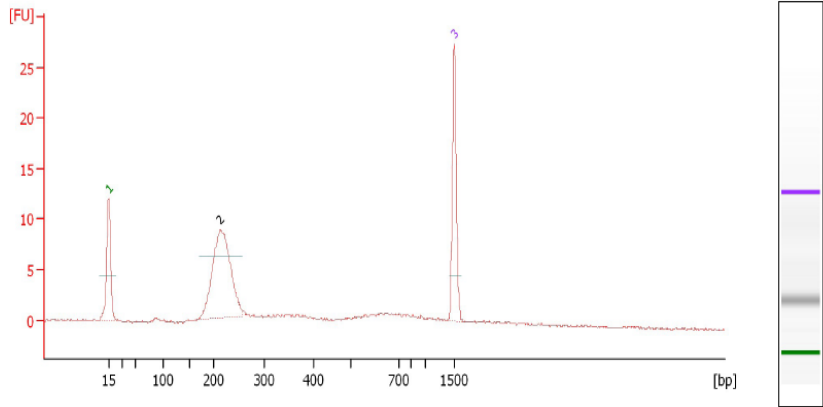


Predominant species of appropriate MW

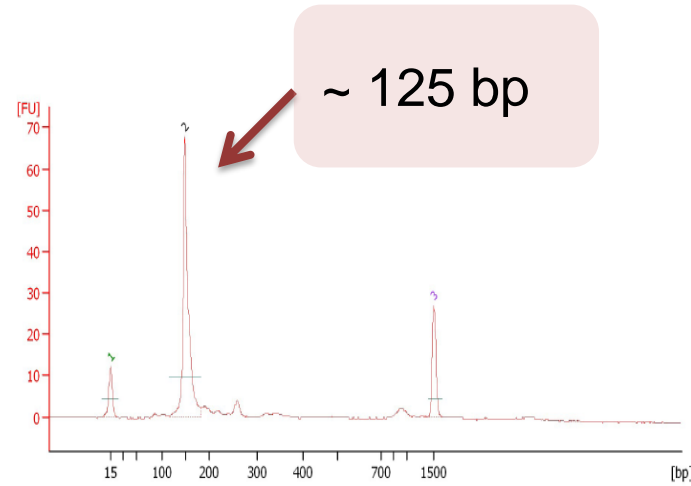
Minimal primer dimer or adapter dimers

Minimal higher MW material

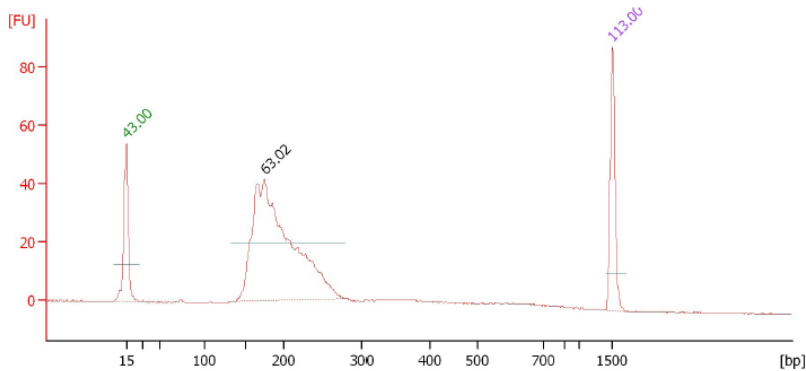
Library QC by Bioanalyzer



Beautiful

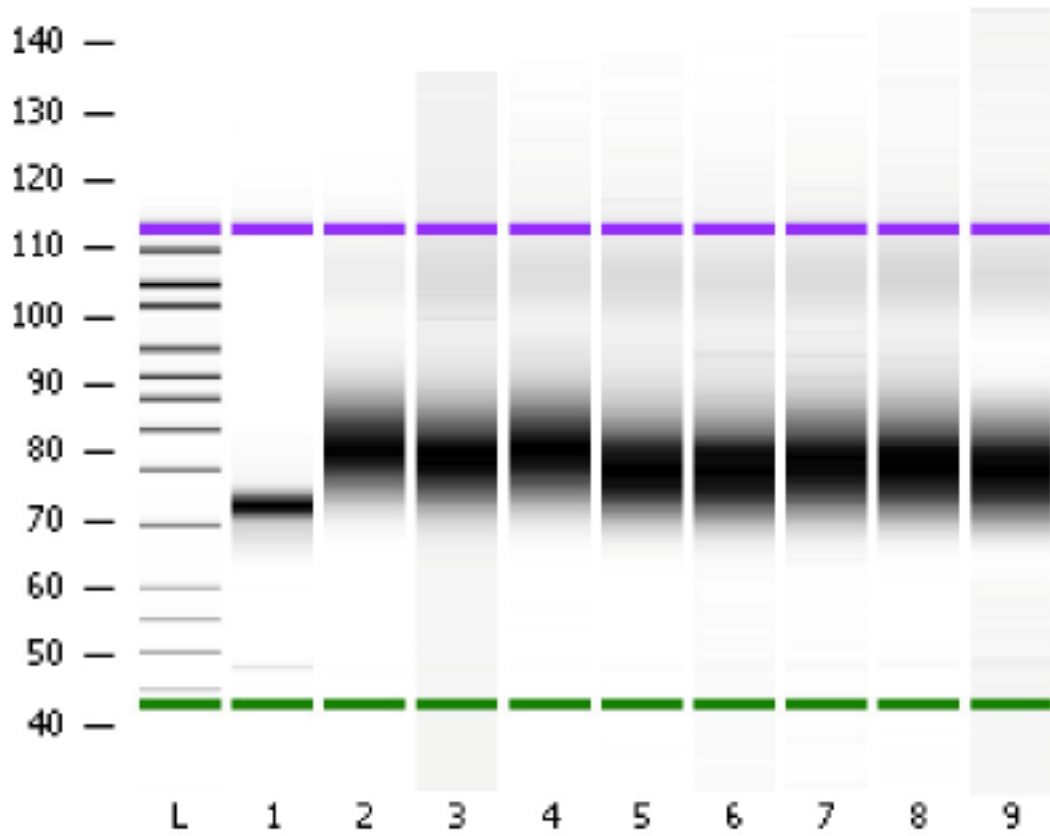


100% Adapters



Beautiful

Library QC



Examples for successful libraries



← ~125 bp

Adapter contamination at ~125 bp



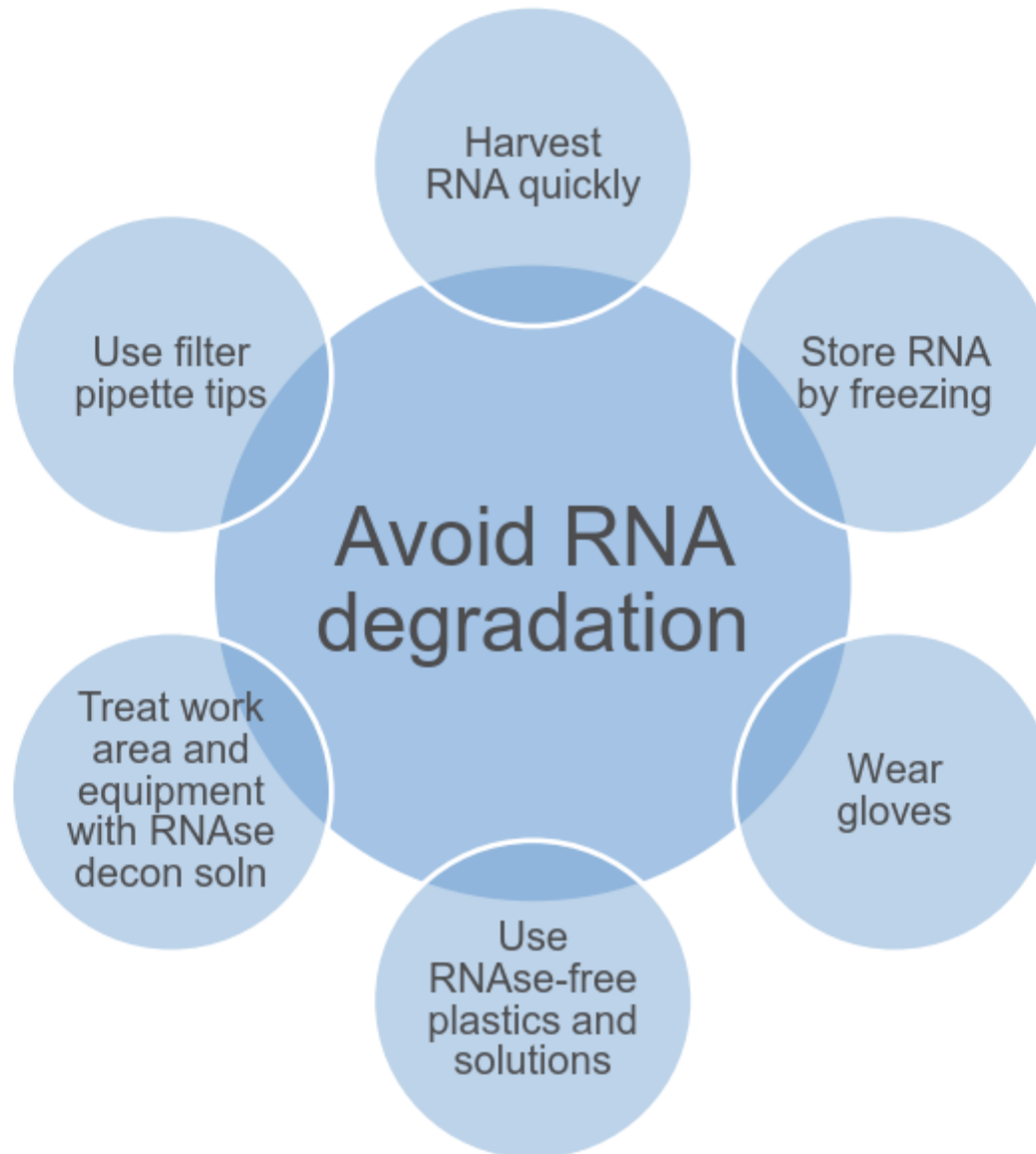
RNA is not that fragile



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



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 |

Standard RNA-Seq library protocol

- QC of total RNA to assess integrity
- Removal of rRNA (most common)
 - mRNA isolation
 - rRNA depletion
- Fragmentation of RNA
- Reverse transcription and second-strand cDNA synthesis
- Ligation of adapters
- PCR Amplification
- Purify, QC and Quantify



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

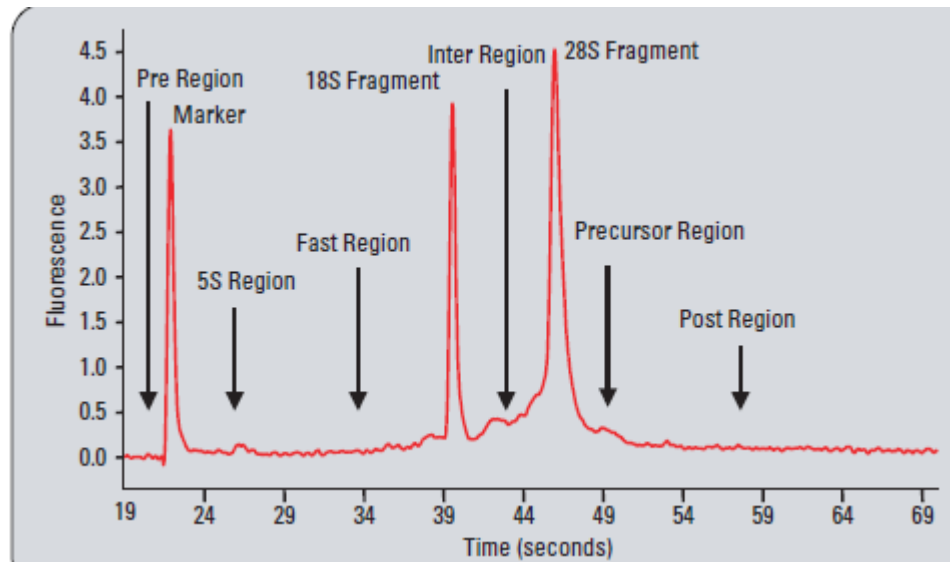
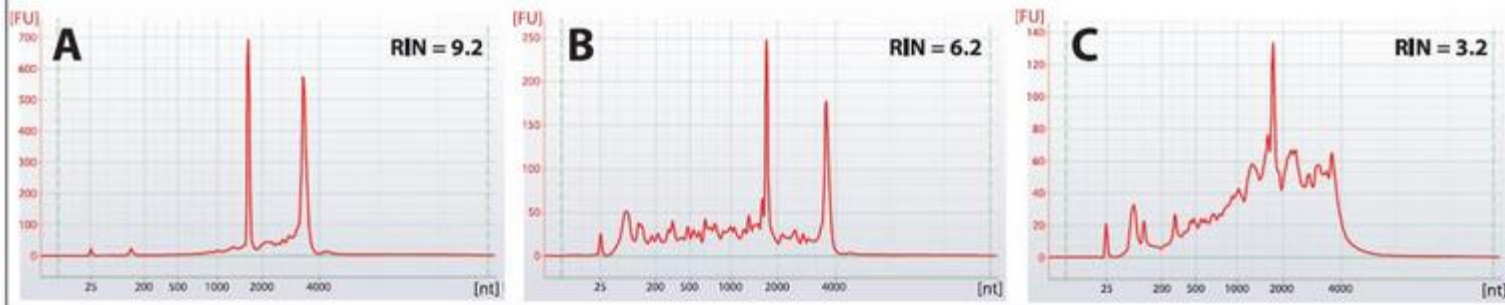
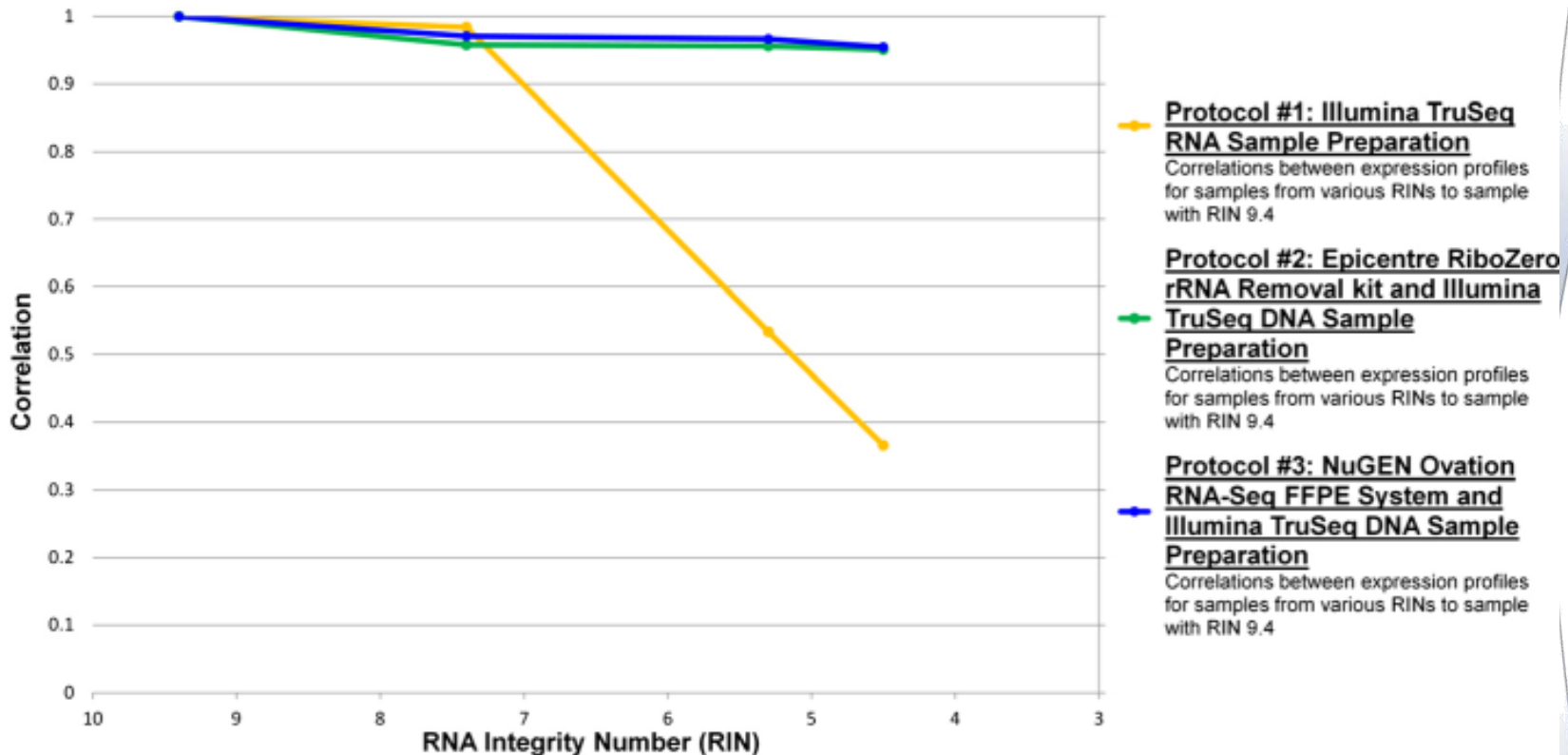


Figure 2.1 Example Agilent Bioanalyzer Electropherograms from three different total RNAs of varying integrity. Panel [A] represents a highly intact total RNA (RIN = 9.2), panel [B] represents a moderately intact total RNA (RIN = 6.2), and panel [C] represents a degraded total RNA sample (RIN = 3.2).



RNA integrity <> reproducibility



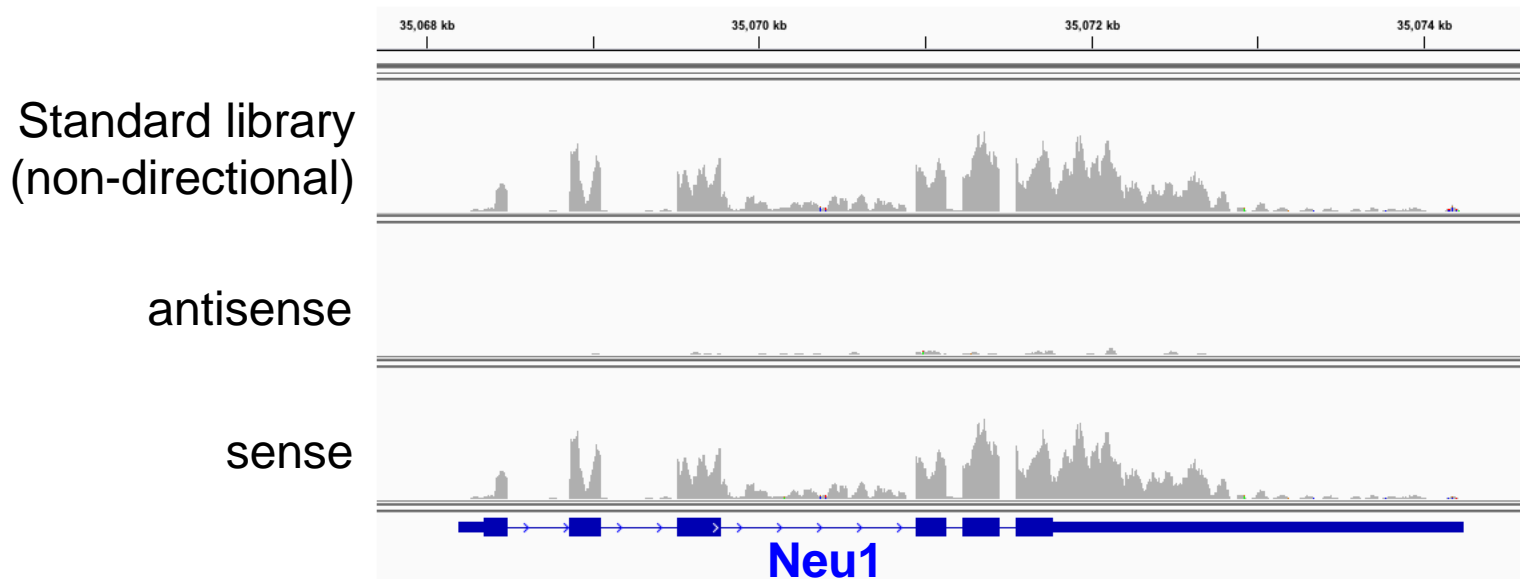
Chen et al. 2014

Considerations in choosing an RNA-Seq method

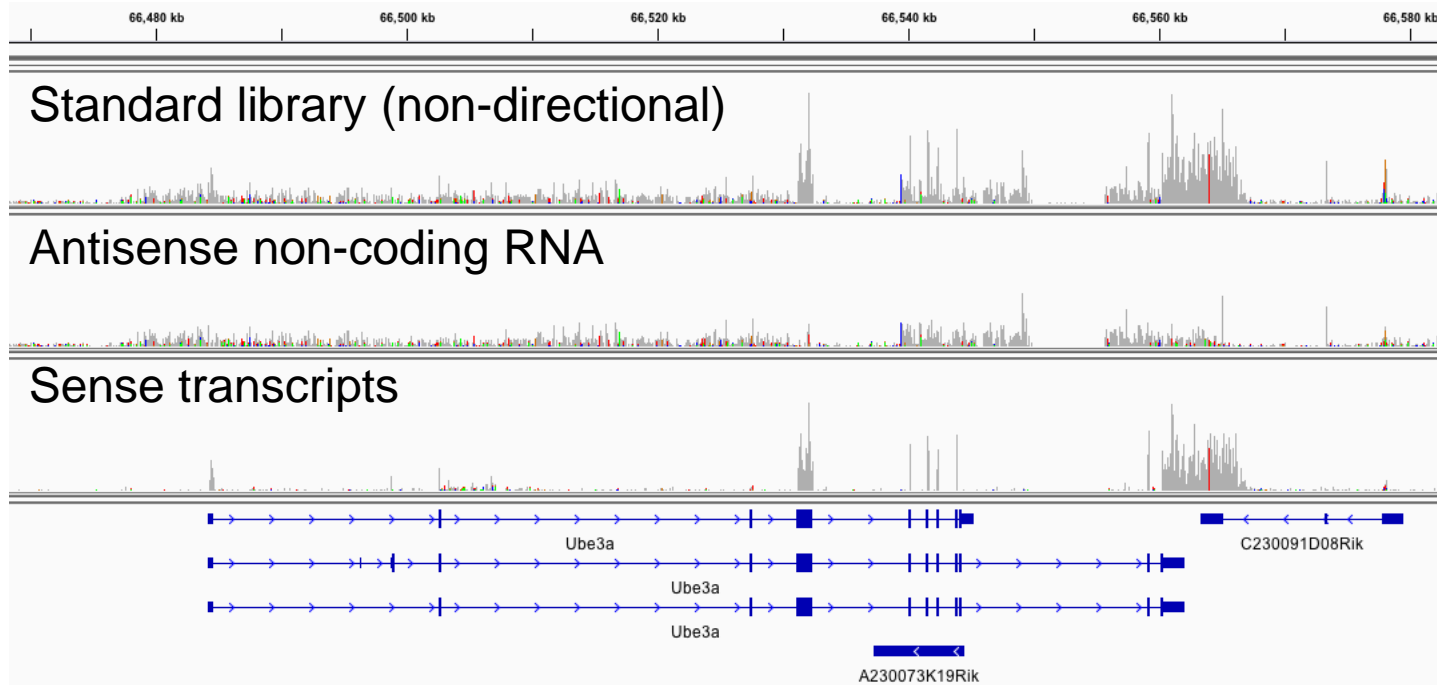
- Transcript type:
 - mRNA, extent of degradation
 - small/micro RNA
- Strandedness:
 - un-directional ds cDNA library
 - directional library
- Input RNA amount:
 - 0.1-4ug original total RNA
 - linear amplification from 0.5-10ng RNA
- Complexity:
 - original abundance
 - cDNA normalization for uniformity
- 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

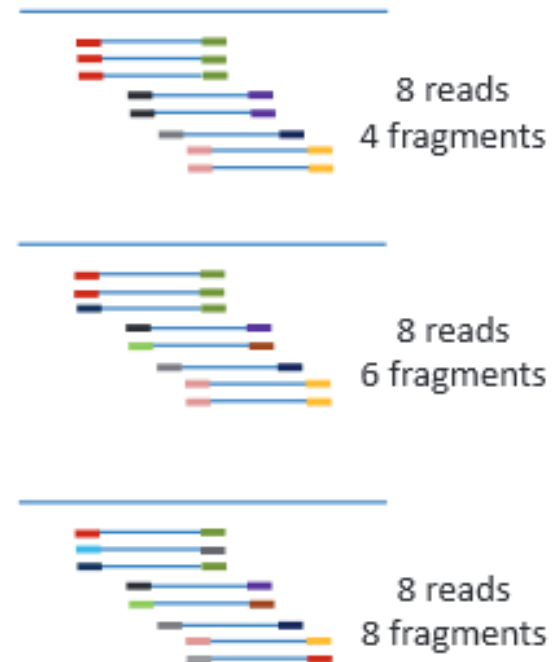
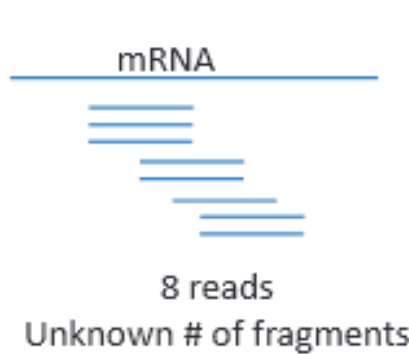
RNA-seq for DGE

- 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)



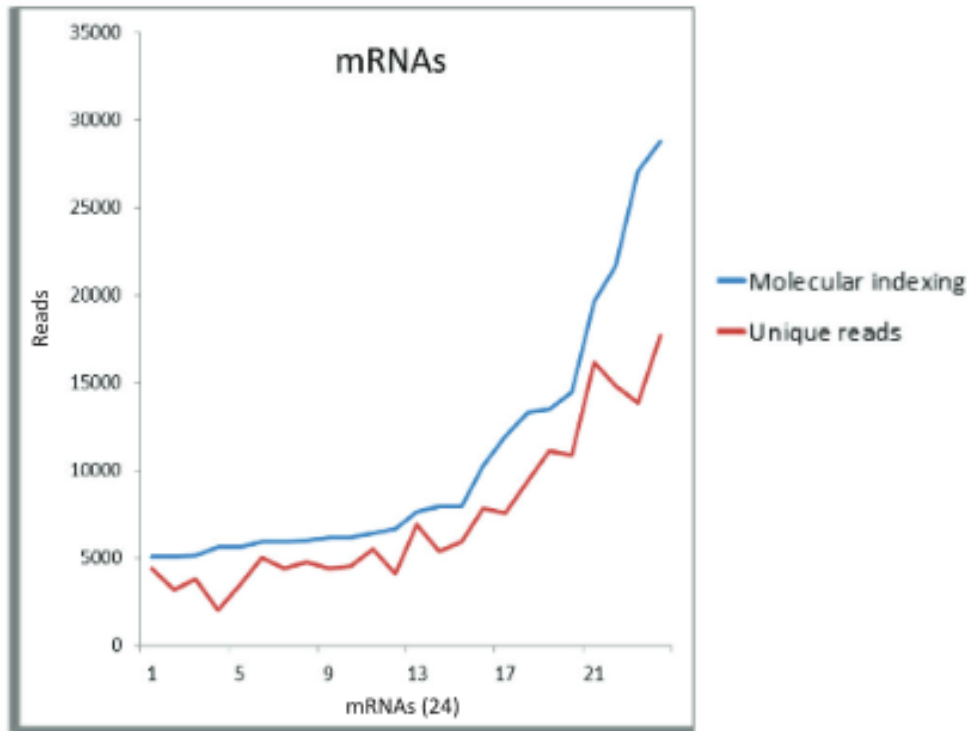
Molecular indexing – for precision counts

Conventional RNA-Seq
Without Molecular Indexing



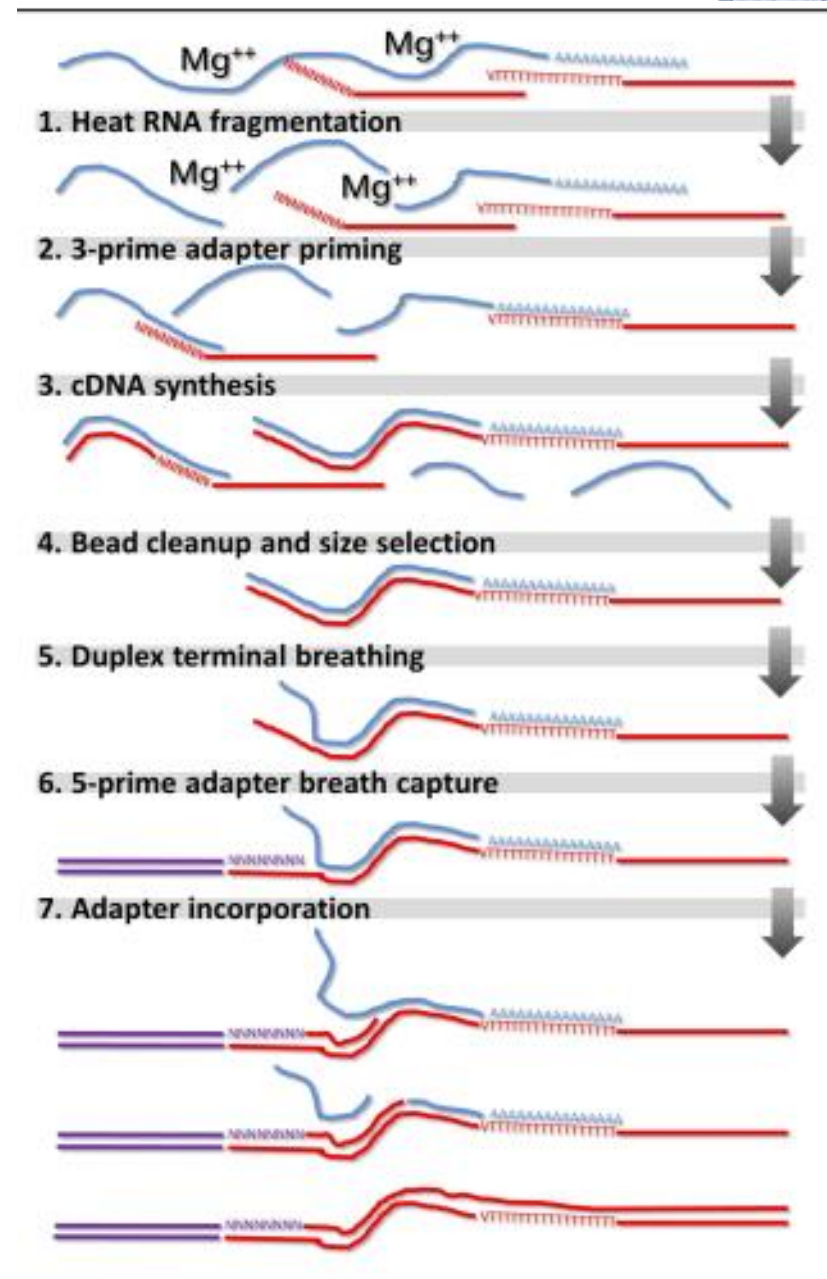
Molecular indexing – for precision counts

B



RNA-seq: cheap and dirty

- 3' Tag-sequencing
- Micro-array-like data
- Quant-Seq
- Brad-Seq (Townesley 2015)



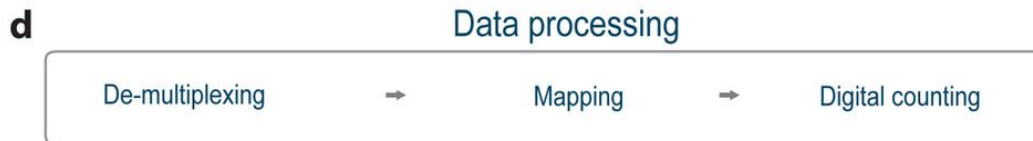
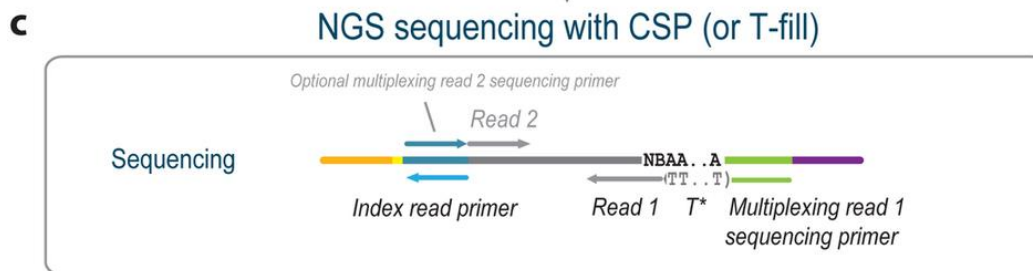
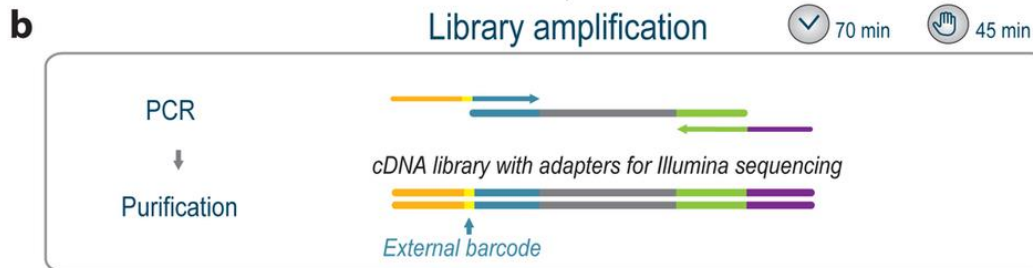
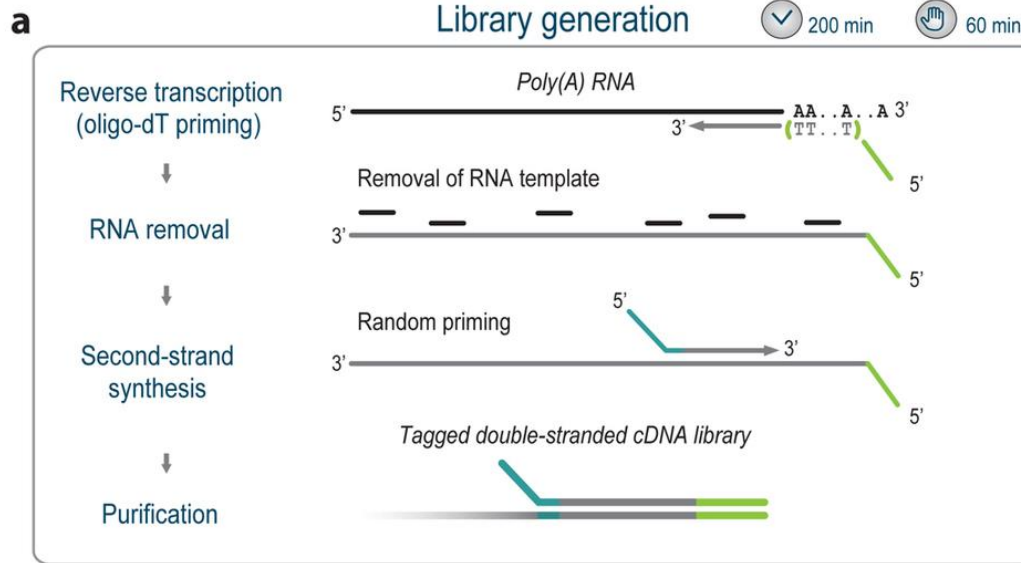
3'-Tag-Seq

- 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



DGE protocols

| | Ribo-depletion | Poly-A enriched | 3-Tag-Seq | Single-cell RNA-seq |
|--|---|---------------------------------------|---------------------------------|------------------------------|
| | all non rRNA transcripts full length | protein encoding genes full length | one "tag" per poly-A transcript | no averaging over cell types |
| | Immature transcripts lnc-RNAs circular RNAs tRNAs, etc... | mature mRNAs | mature mRNAs | mature mRNAs |
| | high noise | medium noise | low noise | high noise |
| | | | | |



Other RNA-seq

- 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



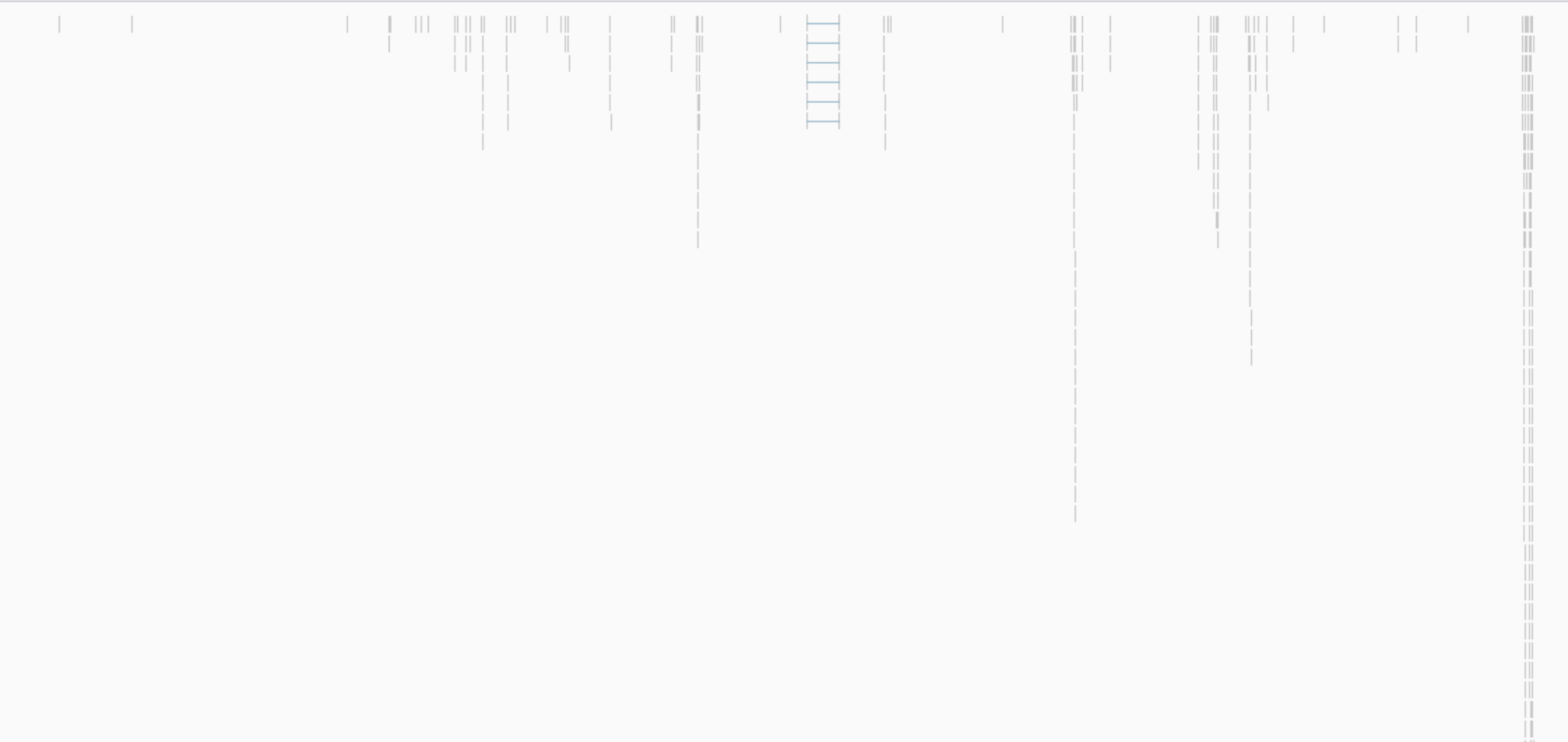
RNA-seq targeted sequencing:

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



Biology

intronic reads ???



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)



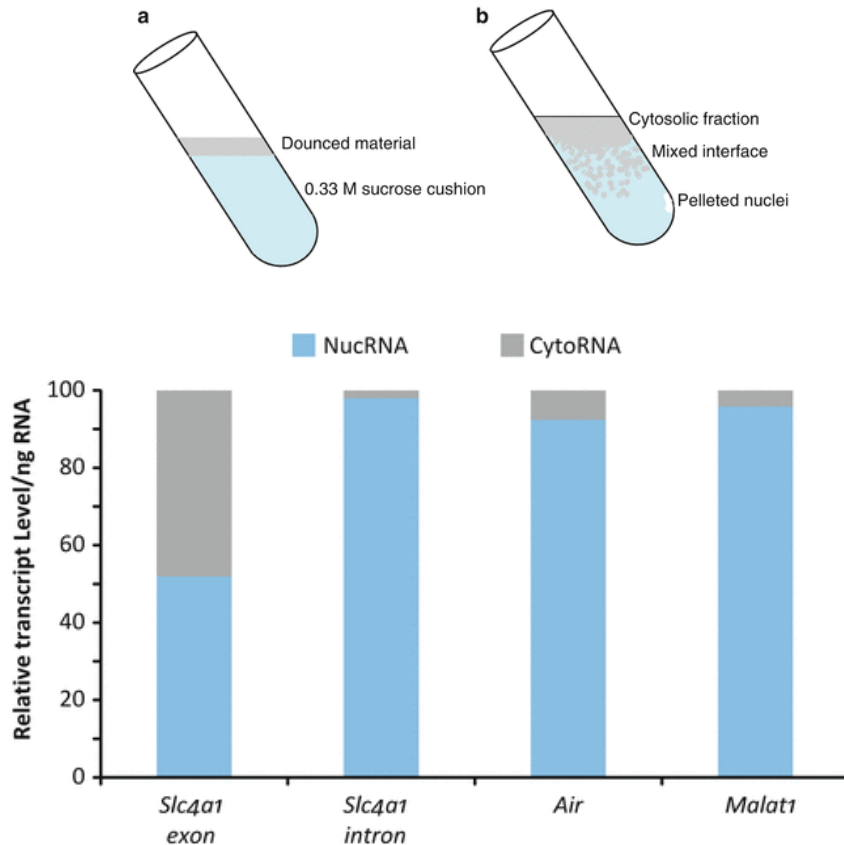
beyond steady-state RNA-seq

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



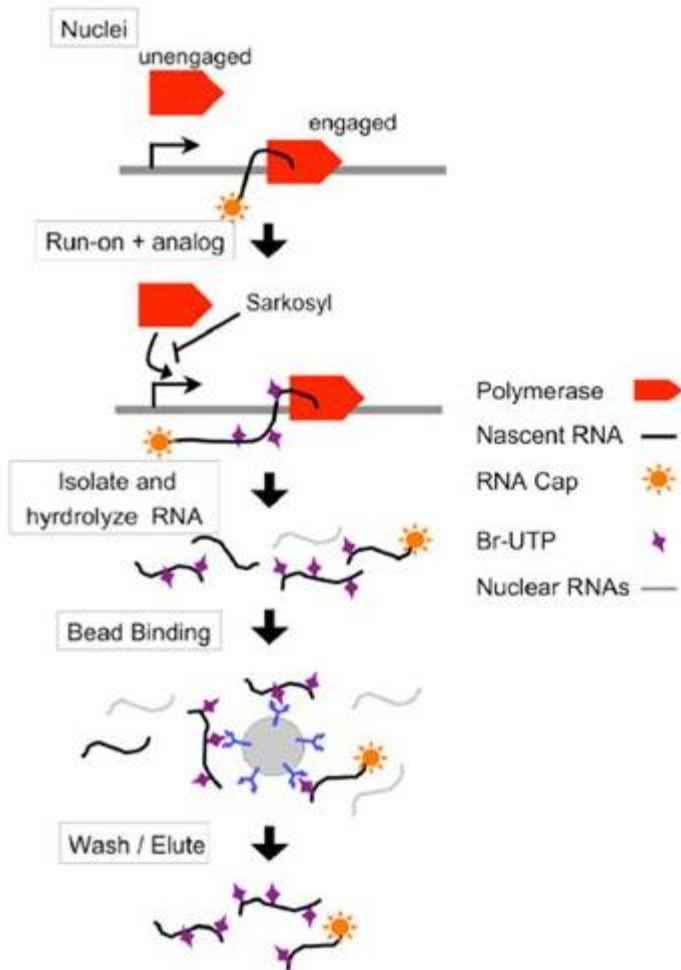
nucRNA-seq

- Fractioning of nuclei and cytosol
- Studying active transcription



Dhaliwal et al. 2016

GRO-Seq



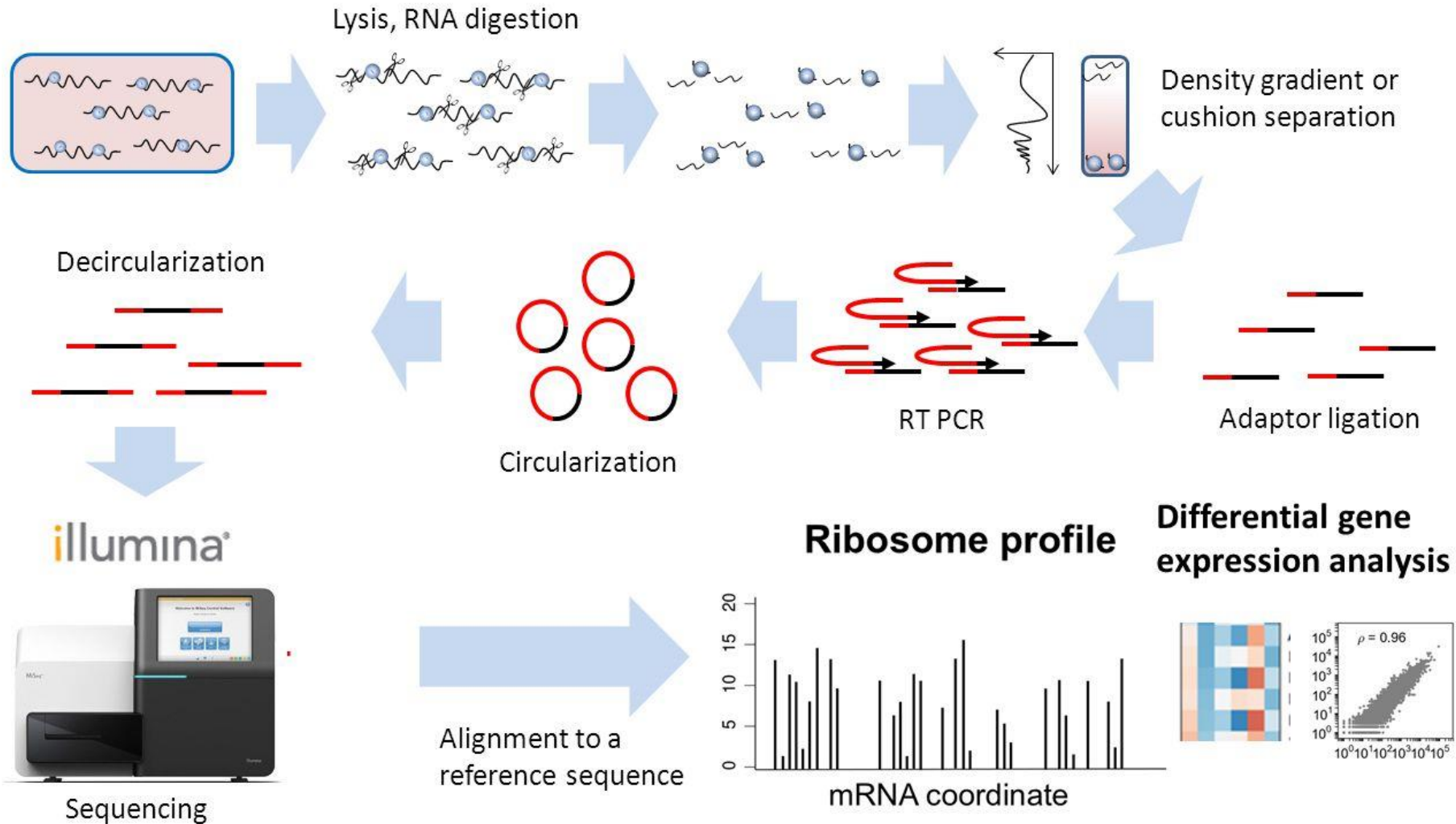
- Global Run-On – sequencing
- pulse-chase experiments (Br-UTP)
- uses isolated nuclei
- sarcosyl prevents binding of polymerase (only transcription in progress will be seq.)
- measures active transcription rather than steady state
- Maps position and orientation
- Earliest changes identify primary targets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



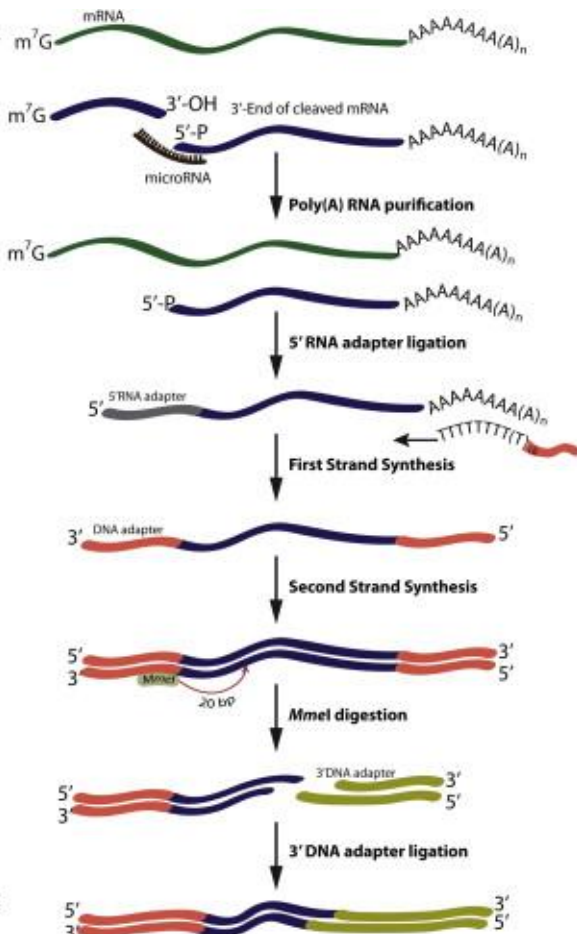
Degradome Sequencing

PARE-Seq

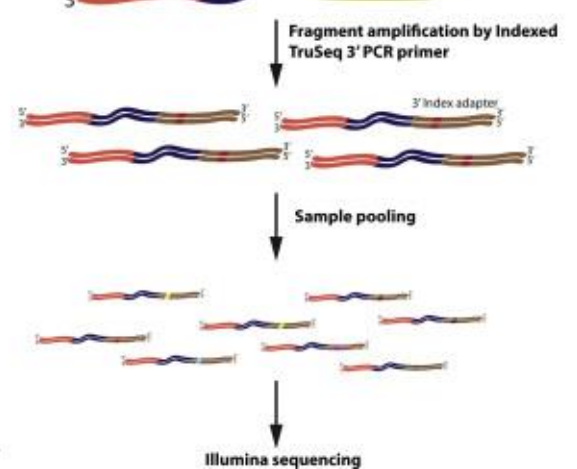
(Parallel Analysis of RNA Ends)

Zhai et al . 2013

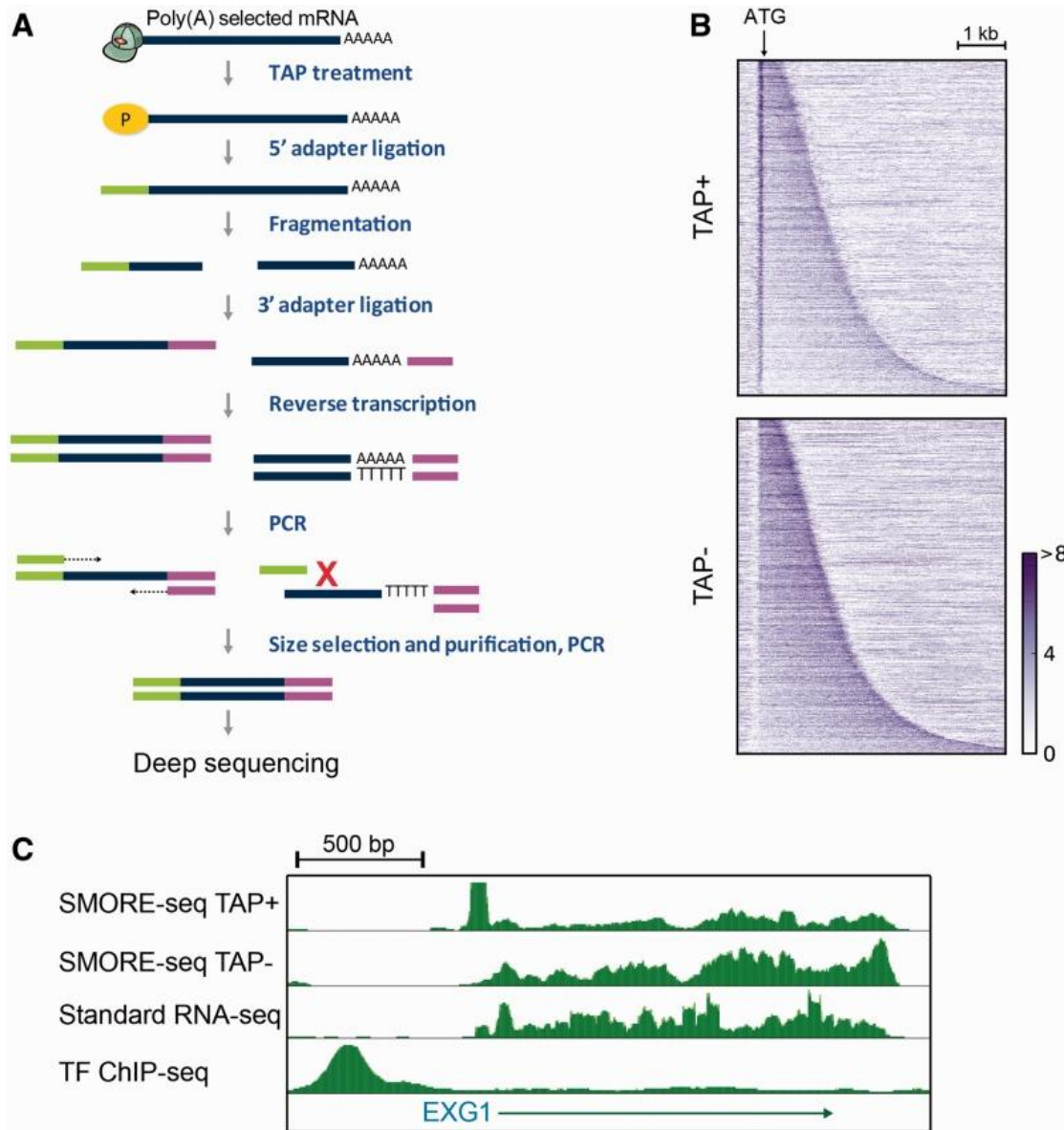
Day 1



Day 2-3



Degradome Sequencing



SMORE-Seq

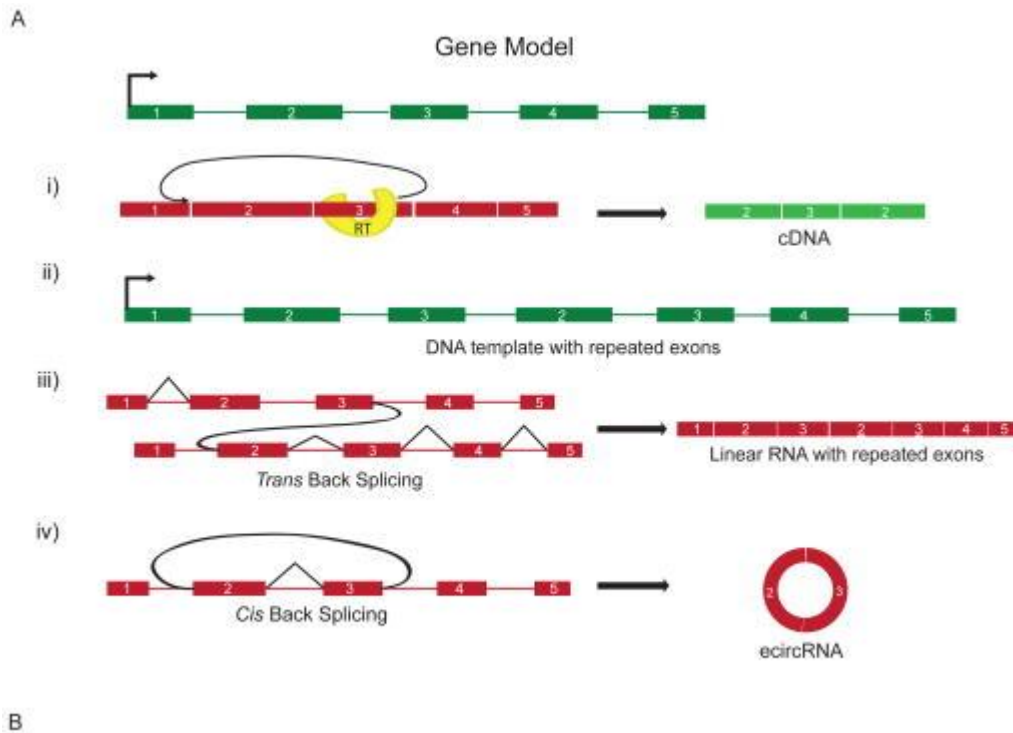
Circular RNA (circRNA)

- 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 ????



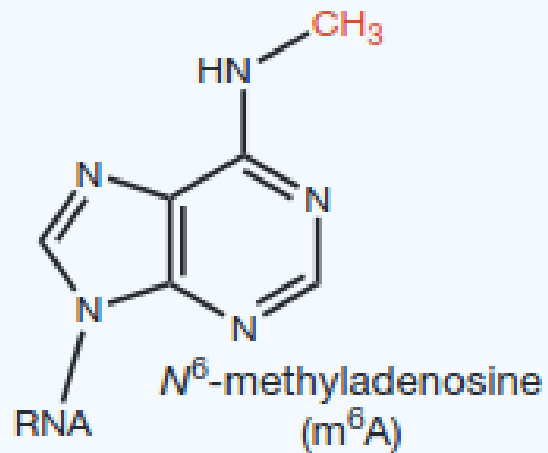
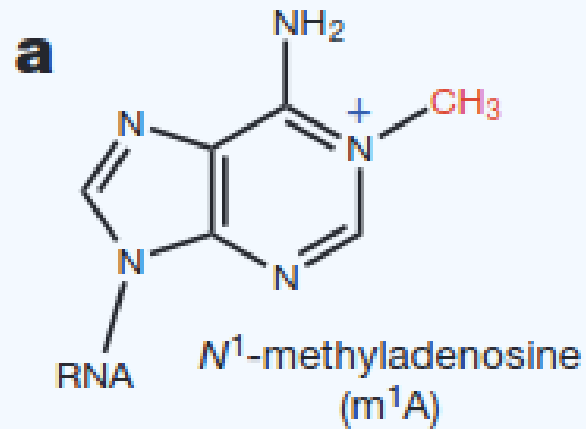
Role of circRNAs ?

Back-splicing and other mechanisms

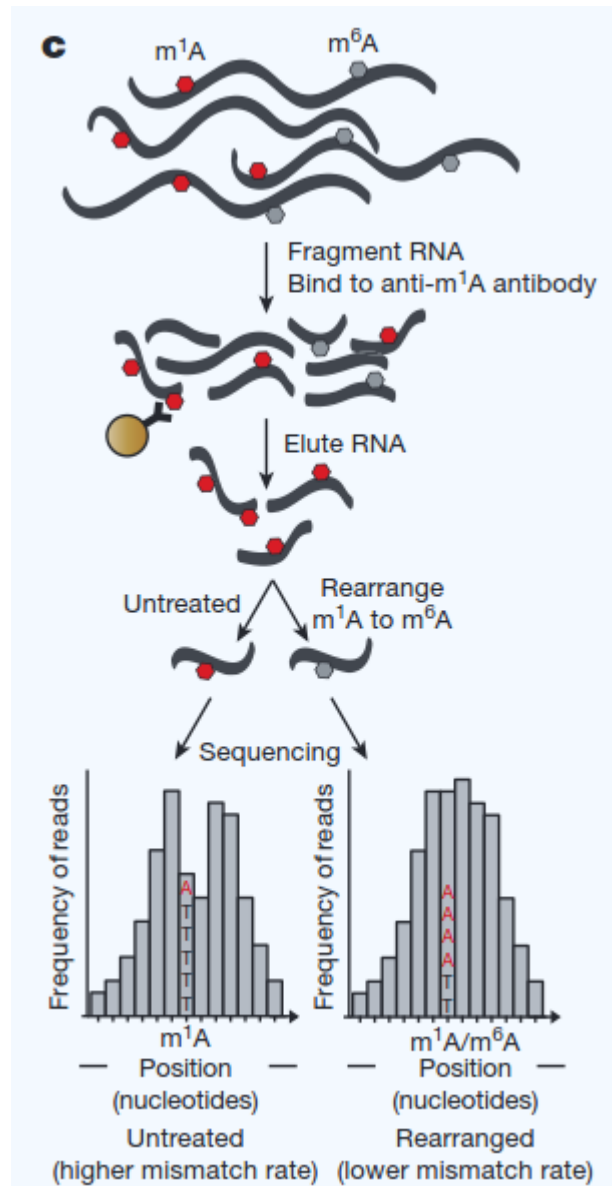


- miRNA sponge
- protein expression regulators:
mRNA traps
(blocking translation)
- Interactions with RNA binding proteins

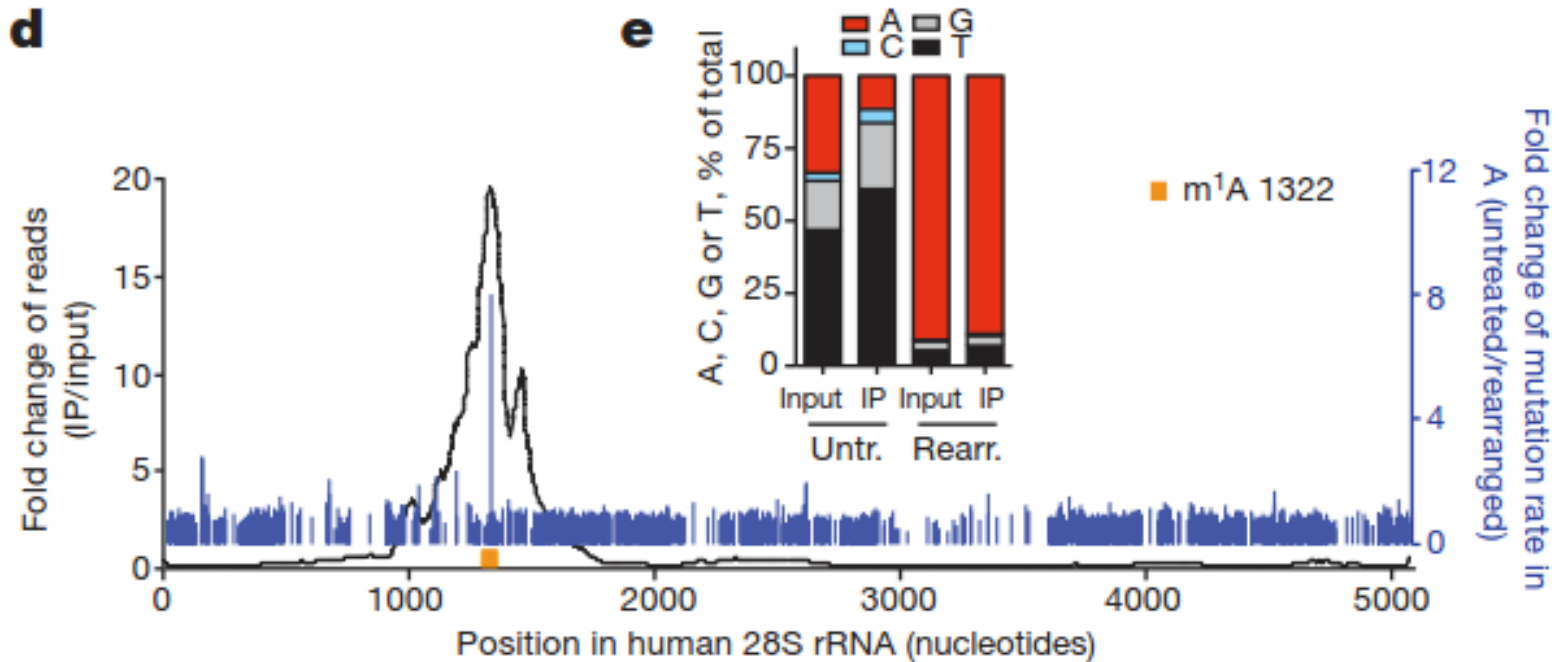
Methylated mRNAs



Methylated mRNAs

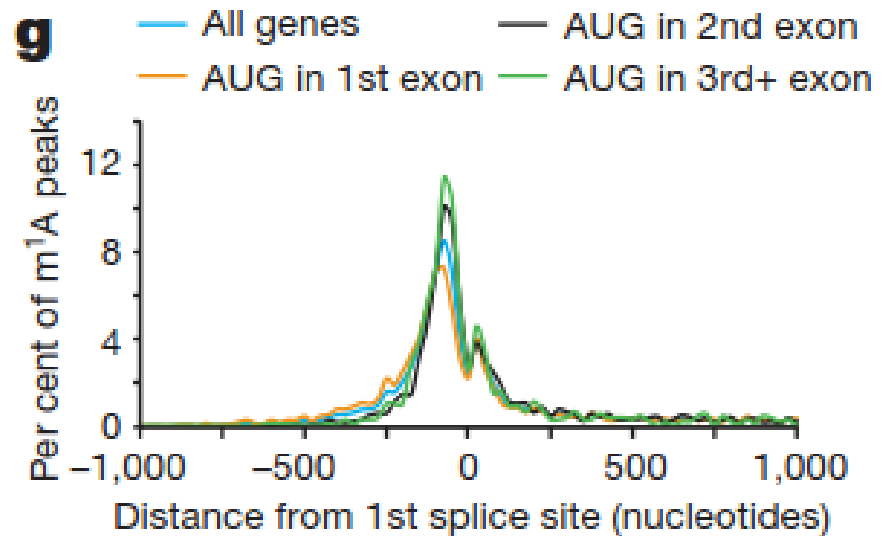
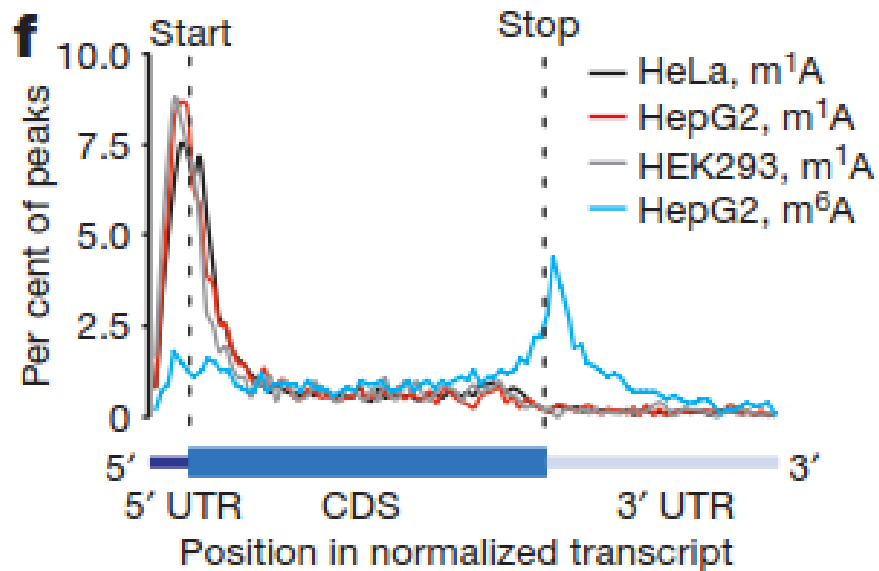


Methylated noncoding RNAs



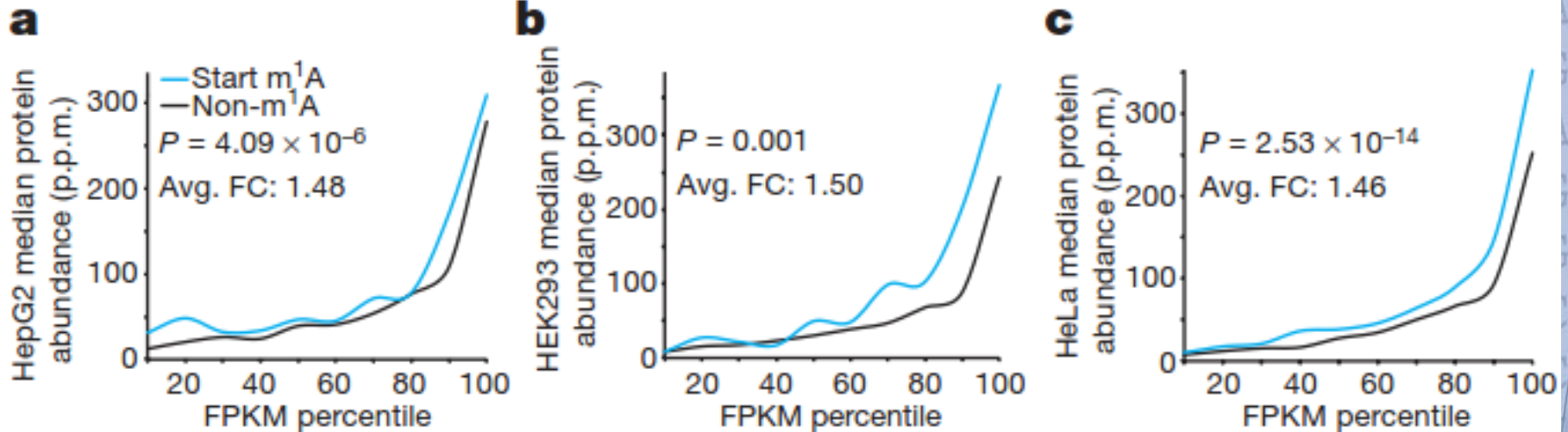
Methylated mRNAs

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

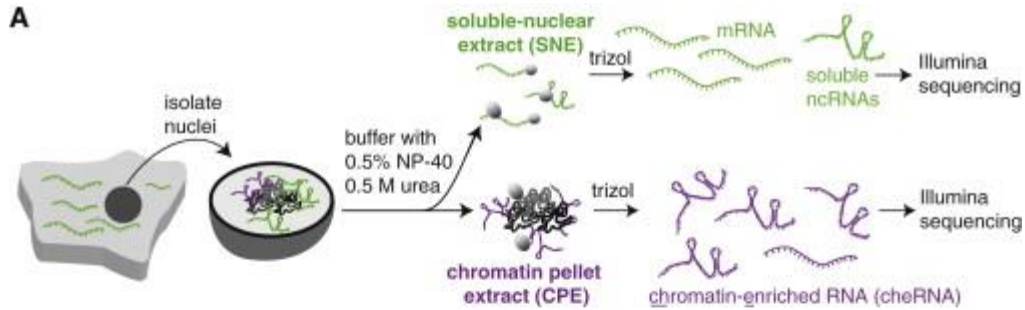


Methylated mRNAs

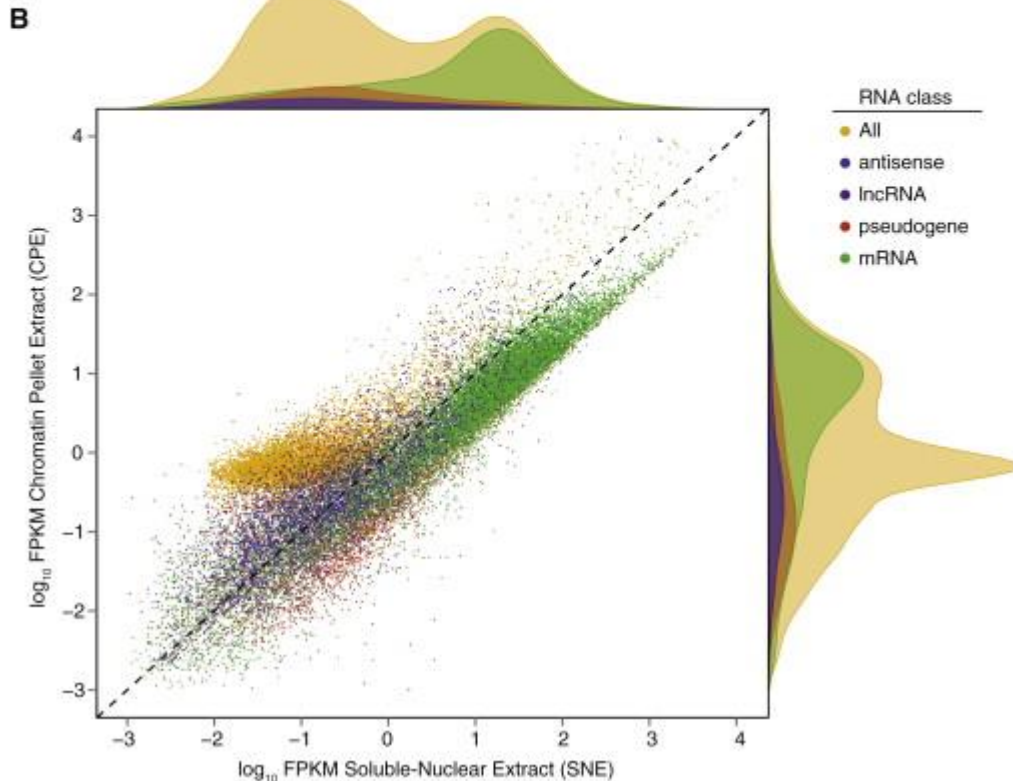
- m¹A around the start codon correlates with higher protein



chromatin-enriched RNAs



- Soluble vs. chromatin bound lncRNAs



Werner et al. 2015

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?



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%).

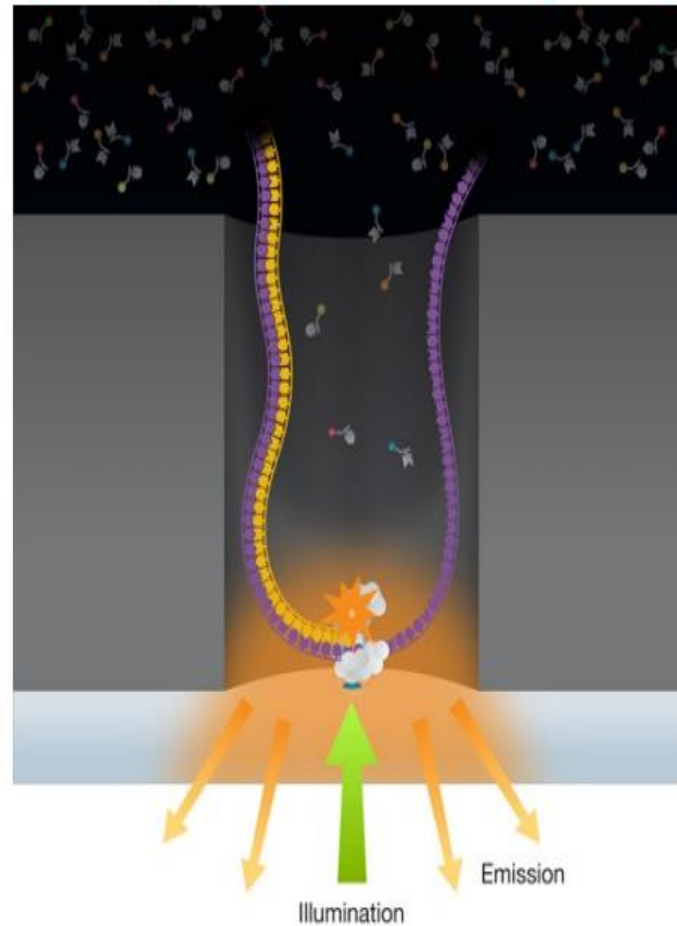
Complementary to short accurate reads of Illumina

TGAGAT
TATGAGC
TAAATCTC
TACCCCT
GCTGAAAC
ATTCCT
TCTGGGA
GAAATT
TGTTGA
AAGGAG
TTTGGG
CGCCAG
TCCCAG
AATTGC
TCTCCA
AAGGCT
AATTGA
GCACAA
ATACCA
GCTTTT
TTTATA

Third Generation Sequencing : Single Molecule Sequencing

Pacific Biosciences

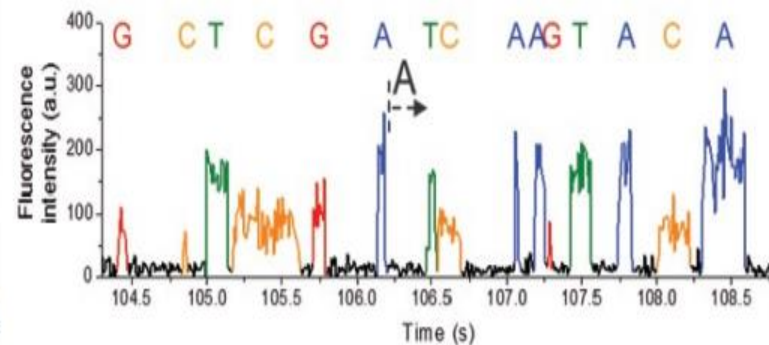
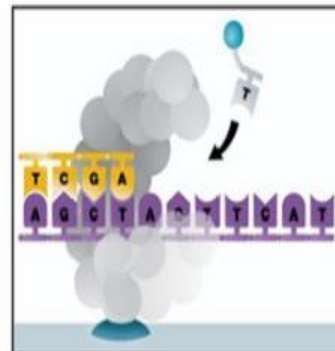
70 nm aperture
“Zero Mode Waveguide”

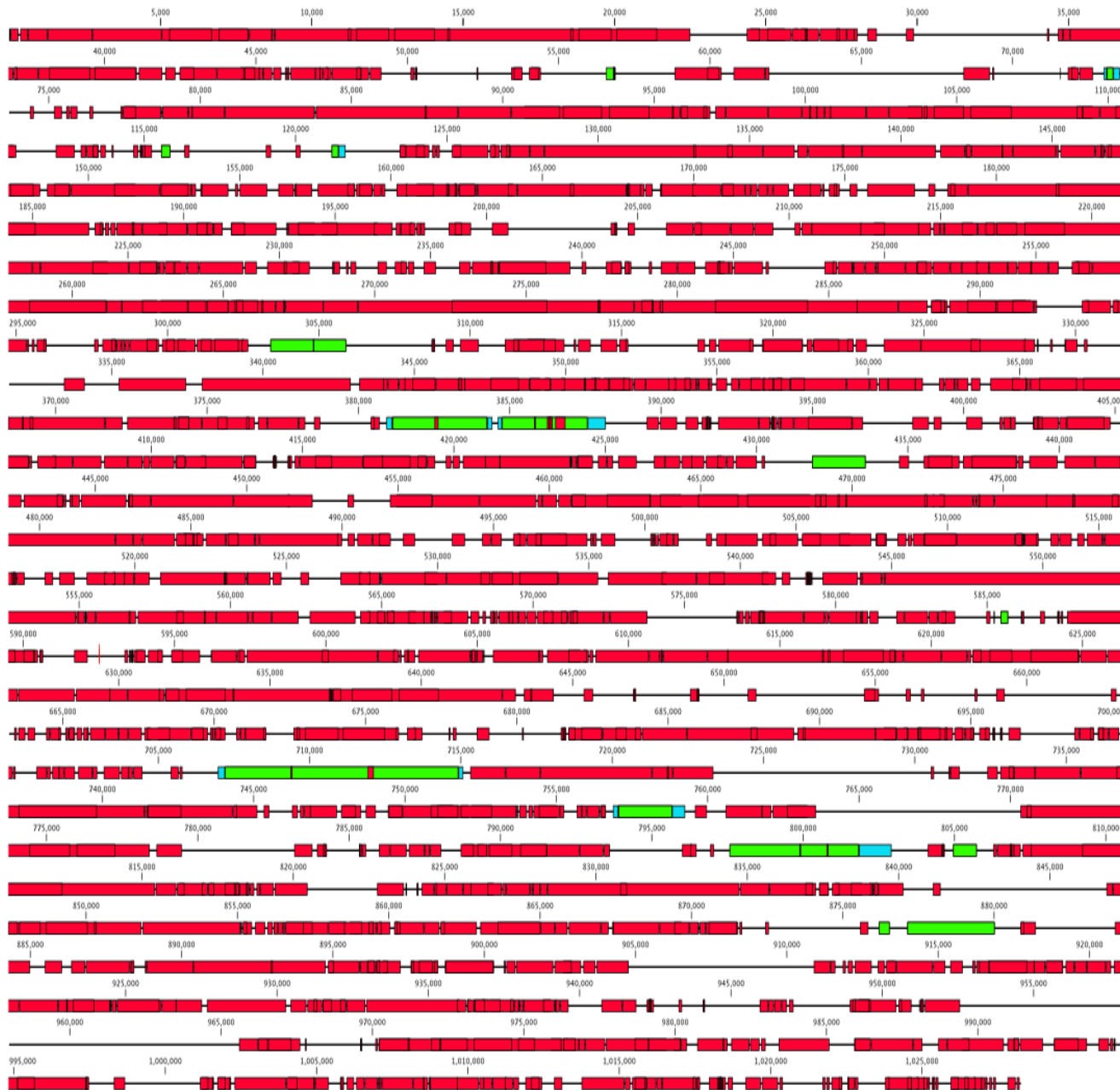


4 nucleotides with different fluorescent dye simultaneous present

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

laser damages polymerase





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

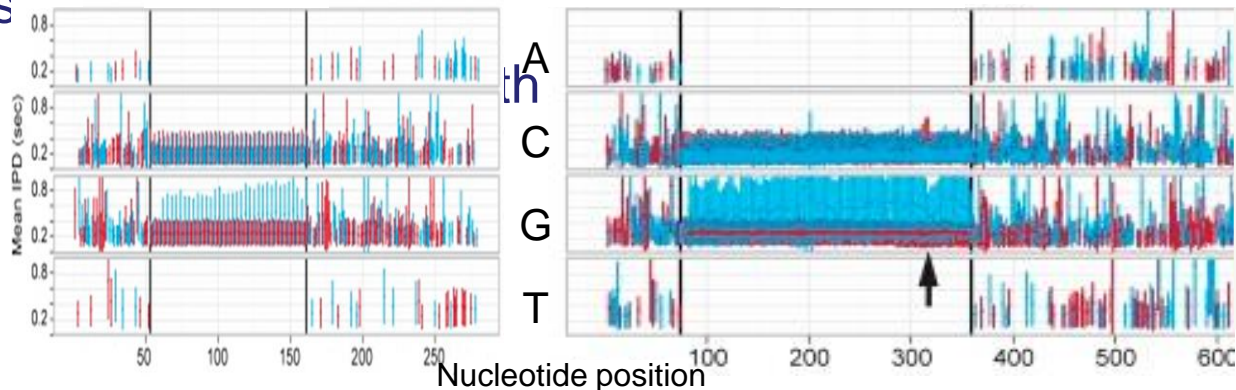


Paul Hagerman, Biochemistry and Molecular Medicine, SOM.

- 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 disorder

CGG₃₆ CGG₉₅

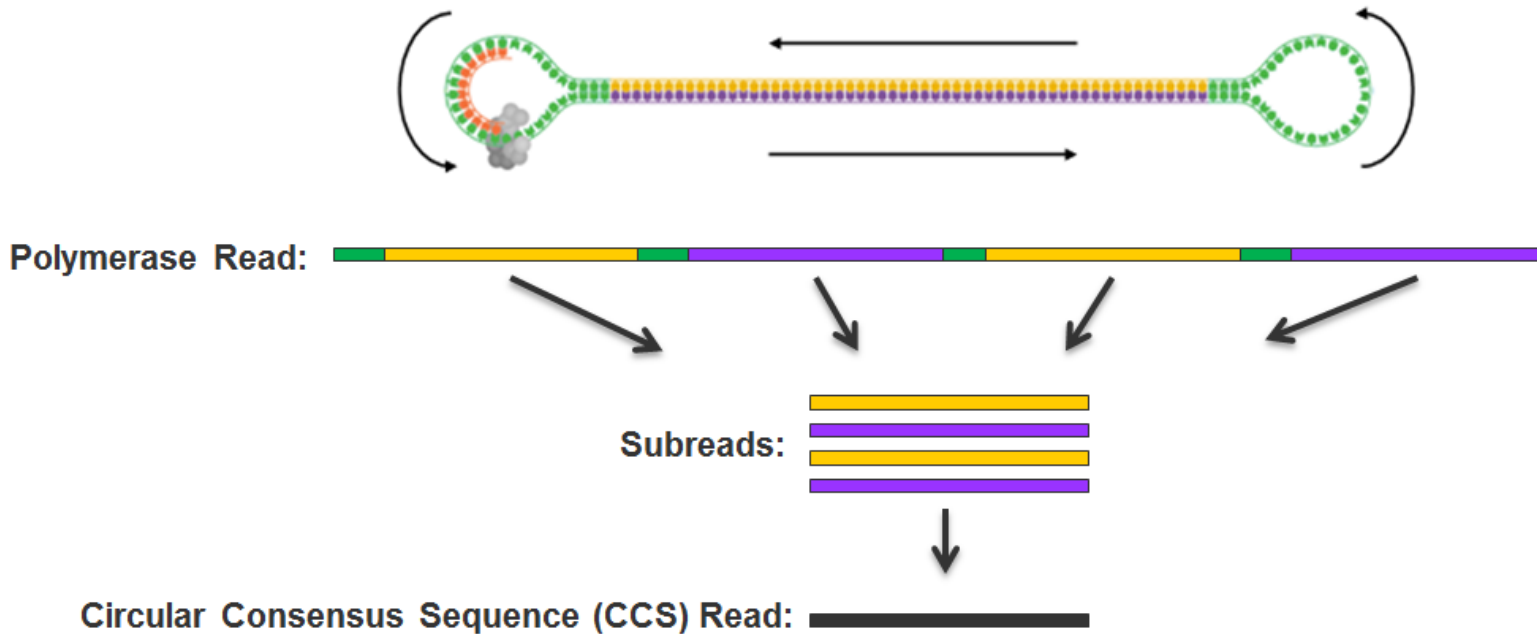


Iso-Seq Pacbio

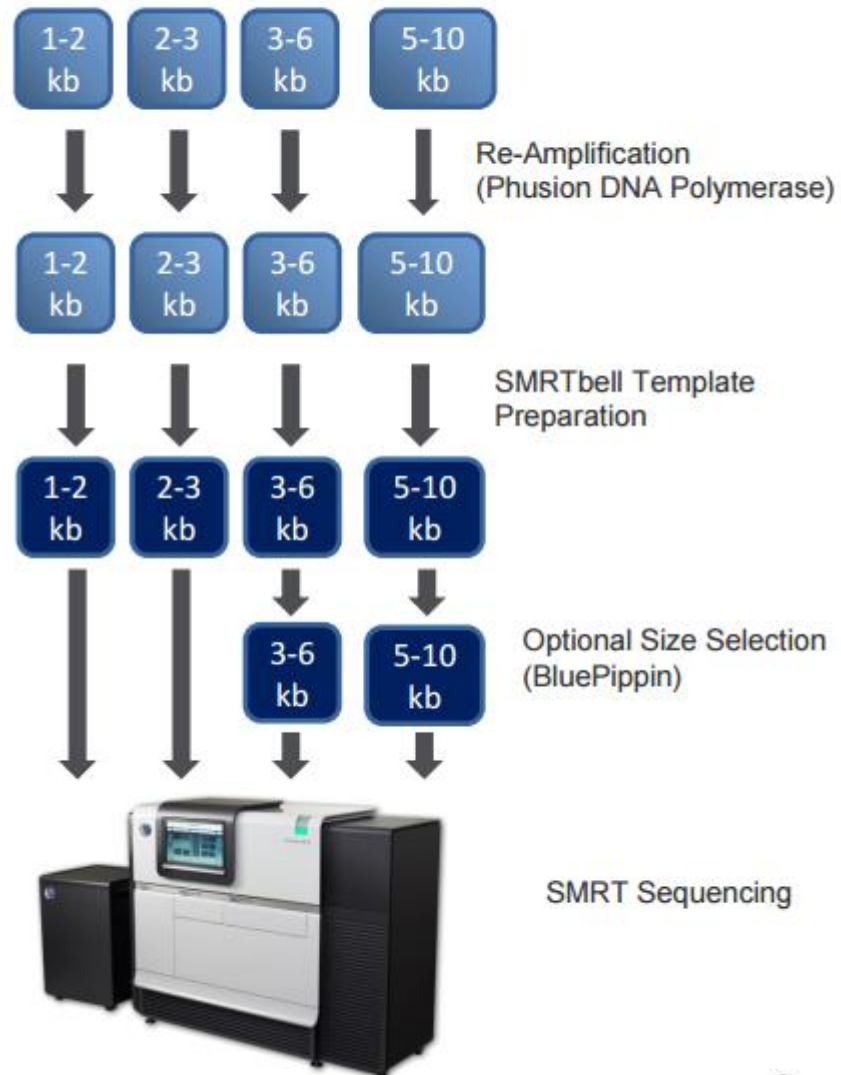
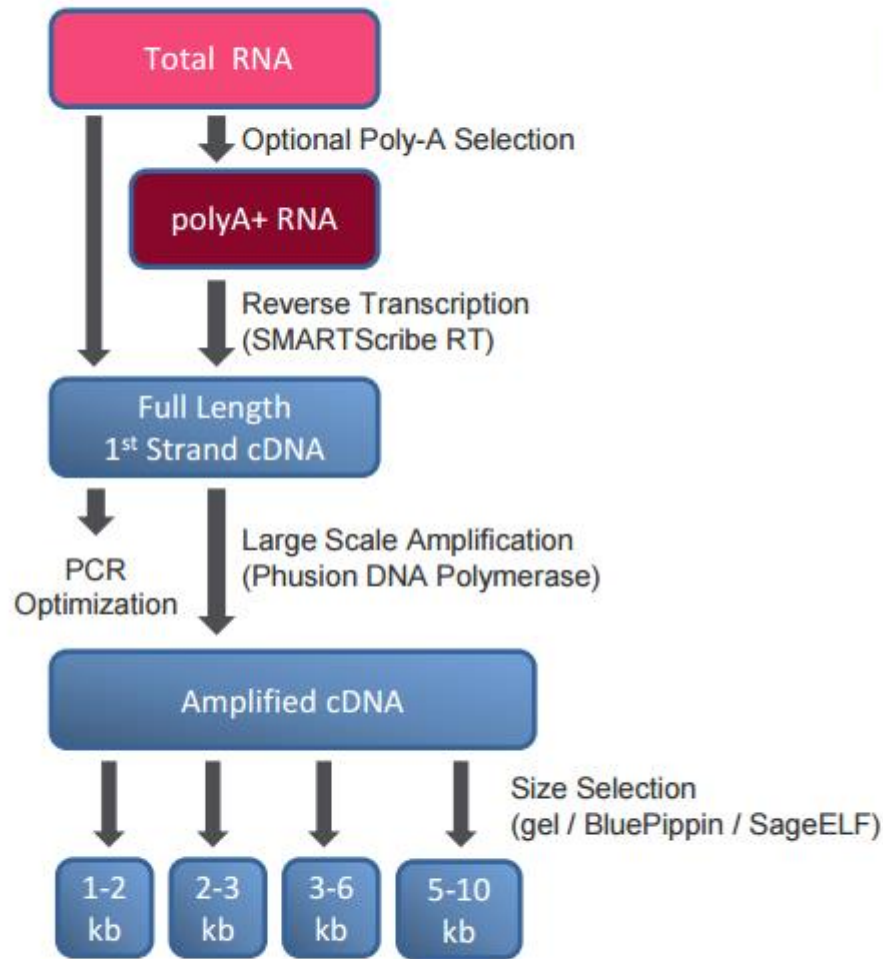
- Sequence full length transcripts
→ no assembly
- High accuracy (except very long transcripts)
- More than 95% of genes show alternate splicing
- On average more than 5 isoforms/gene
- Precise delineation of transcript isoforms
(PCR artifacts? chimeras?)



SMRT-bell adapters circular sequencing

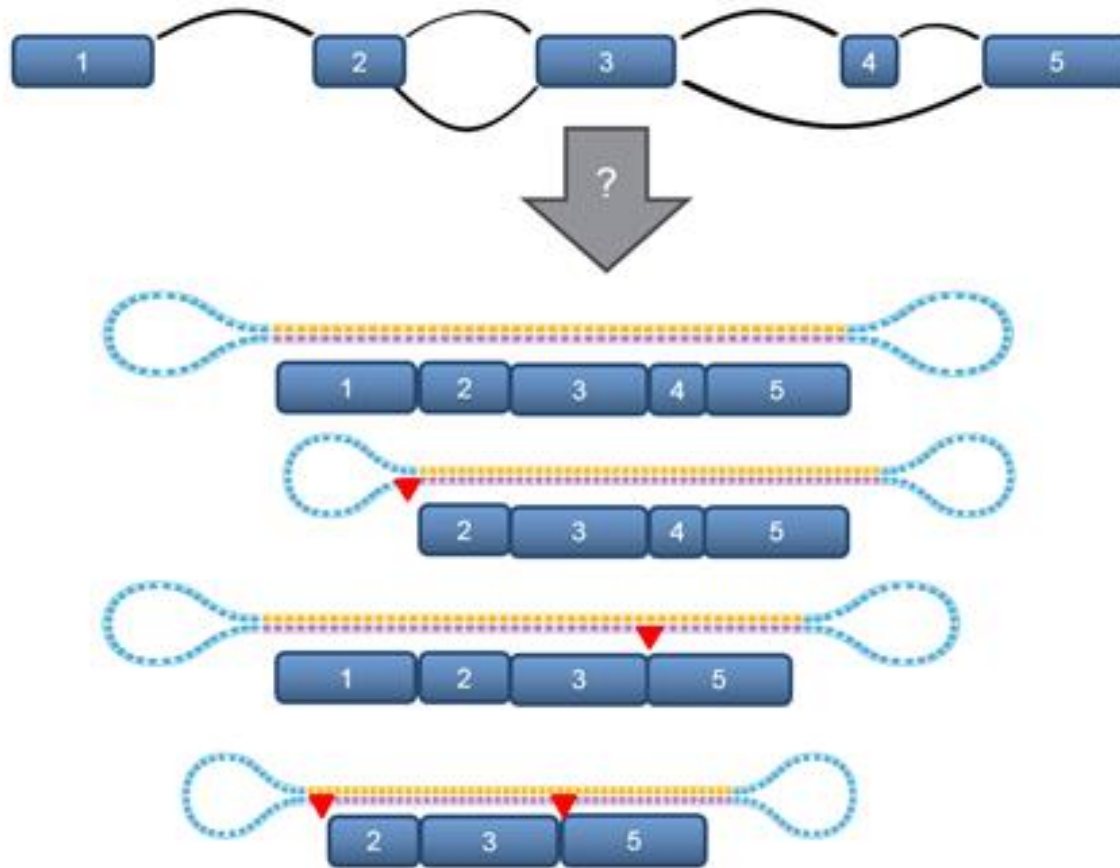


Iso-Seq Library Workflow



Size cuts can be arbitrary
Current max FL transcript seen: 15 kb

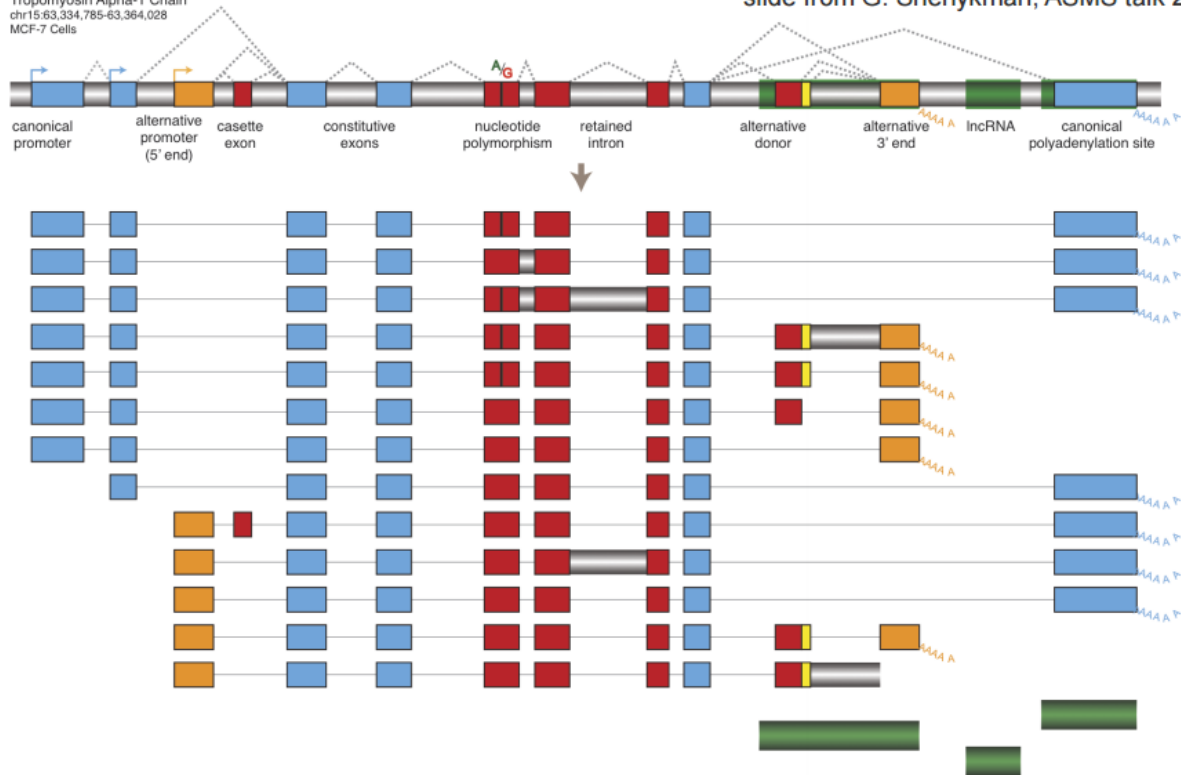
Alternative Splicing



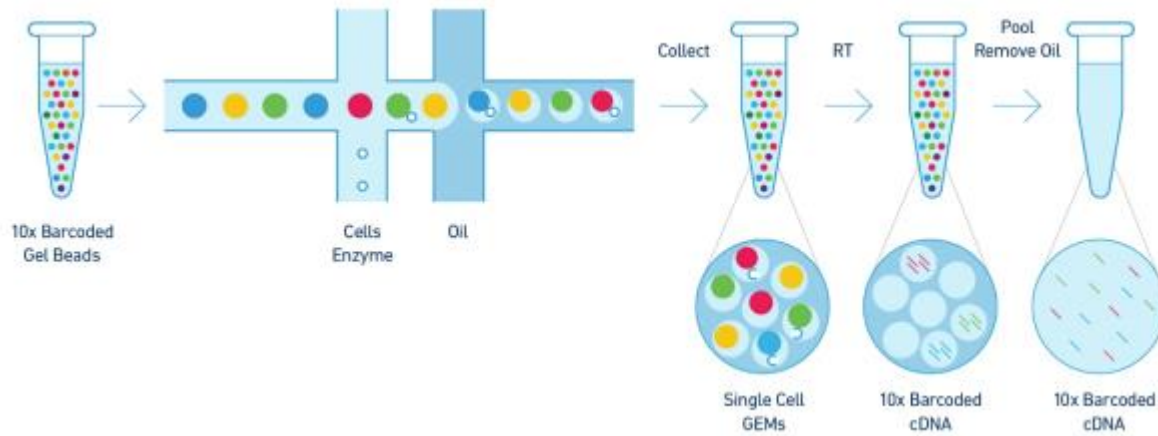
A Single Gene Locus → Many Transcripts

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

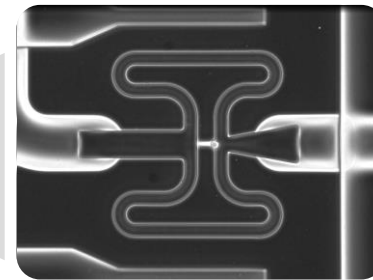
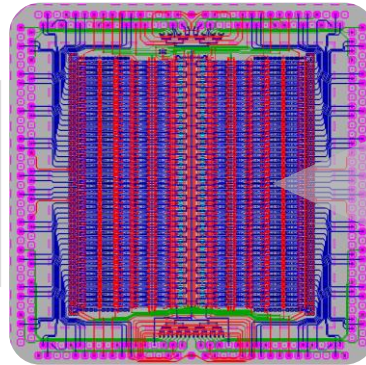
slide from G. Shenykman, ASMS talk 2014



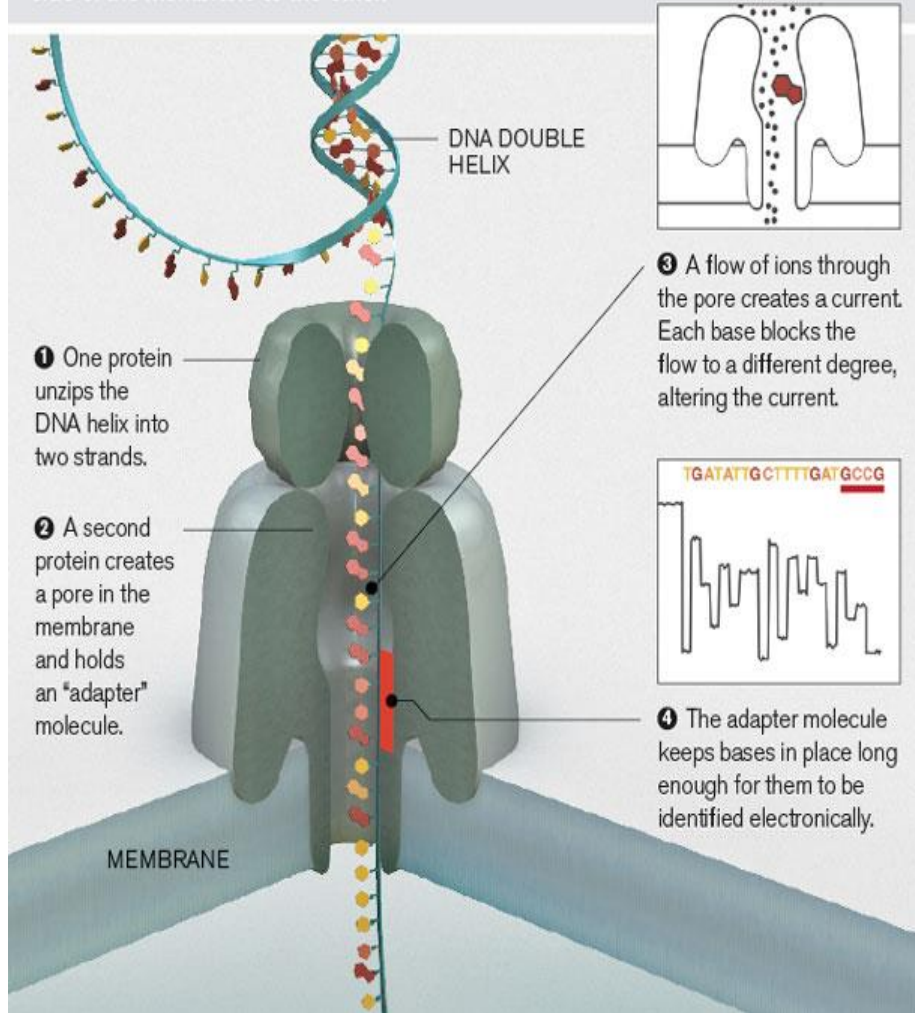
10X Genomics single-cell Drop-Seq



C₁ Single cell capture



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



TYGAGATT
IATGAGG
TAAATCTC
IACCACCT
GCTGAAAG
ATTCCCT
TCTGGGAA
GAAATTAT
TGTTGA
AAGGAG
TTTGGG
CGCCAG
TCCCAG
AATTGC
TCTCCA
AAGGCT
AATTGA
GCACAA
ATACCA
GCTTTT
TTTATA



Thank you!

