

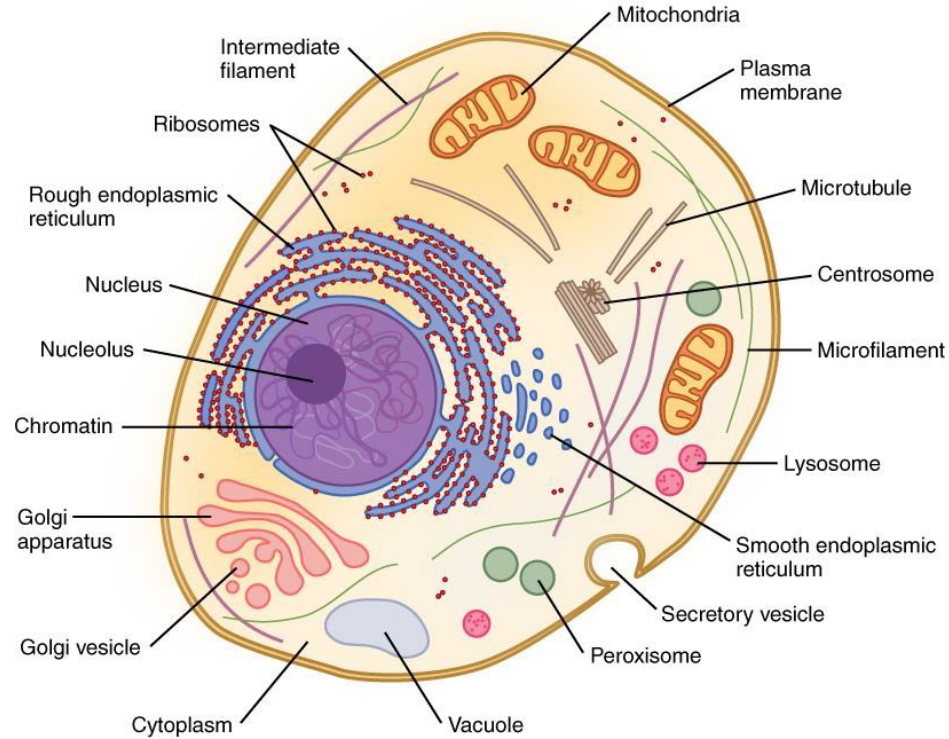
Sequencing for Statisticians



Meng Wang
Senior Bioinformatic Scientist
Stanton Lab, Center for Childhood Cancer Resarch
Nationwide Children's Hospital

Outline

- Background
- Sequencing Platforms
- Different Types of Sequencing



Slides courtesy of Ryan Roberts



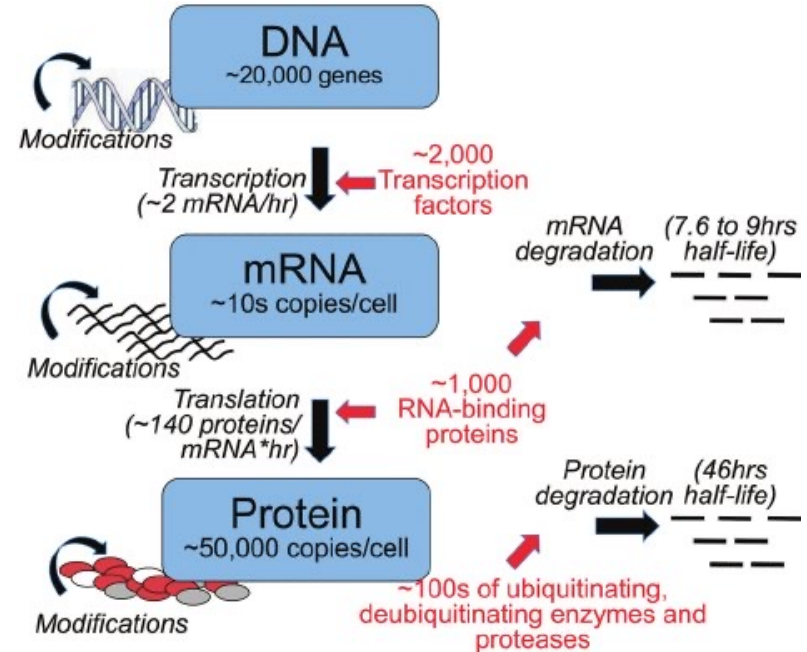
NATIONWIDE CHILDREN'S
When your child needs a hospital, everything matters.

Why sequencing?

- Knowledge of our genome
- Transcription differences
- Transcriptional regulations
- Structural Variations
- etc, etc

What can we quantify?

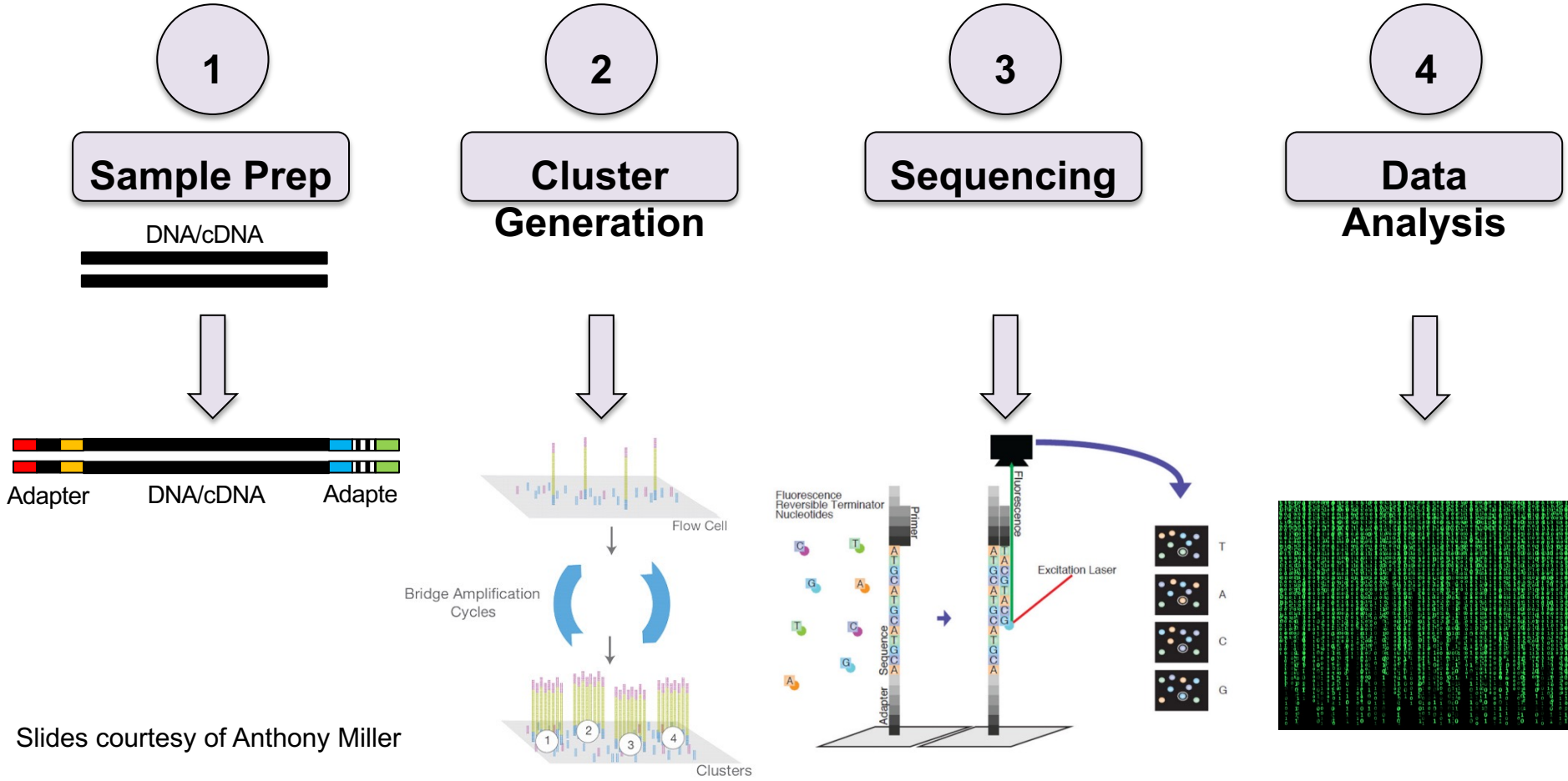
- mRNA (RNA-seq)
- DNA (WGS, ChIP-seq, ATAC-seq, Hi-C)
- Protein (proteomics)



Sequencing Platforms

- Sequencing consists of two steps:
 - Library prep (generating sequencable fragments)
 - Sequencing
 - Illumina (dominant platform)
 - Pacbio/Nanopore (long read sequencing)
-

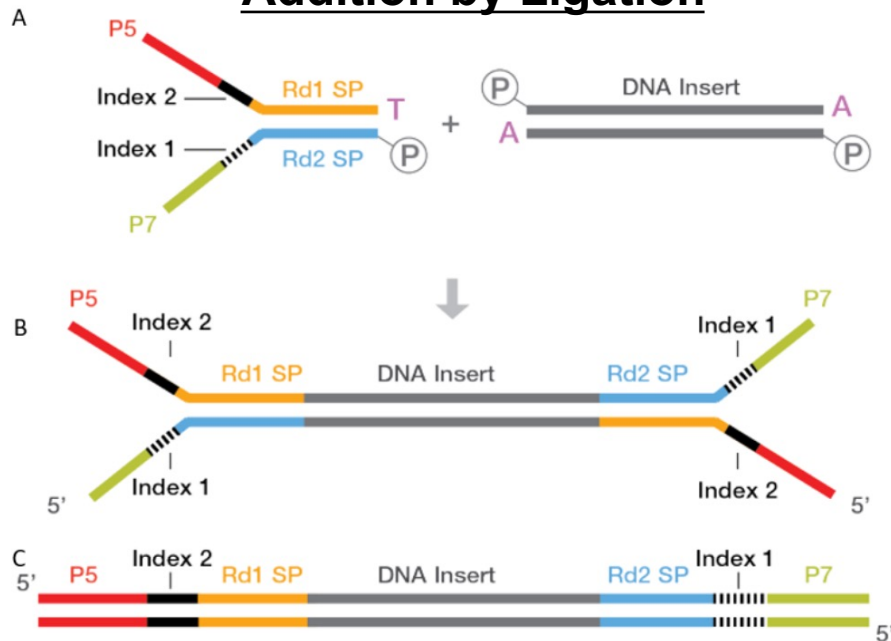
Overview of Illumina Workflow



Slides courtesy of Anthony Miller

Illumina Adapter – Foundation for Cluster Generation

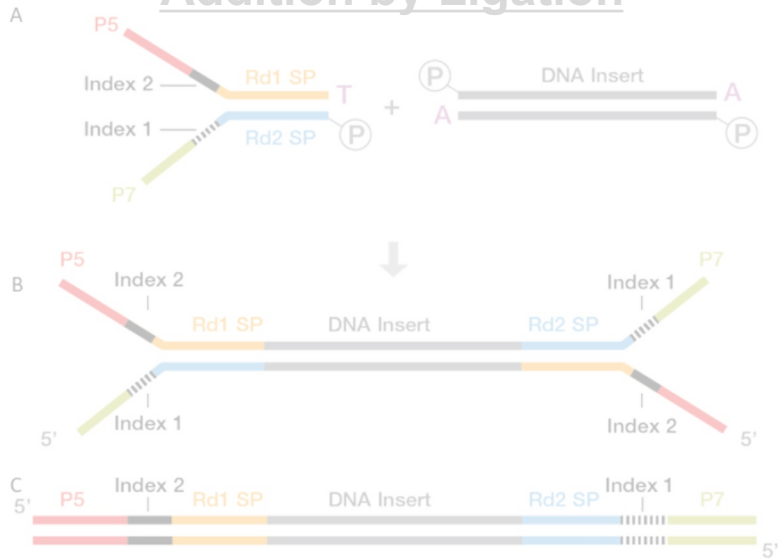
Addition by Ligation



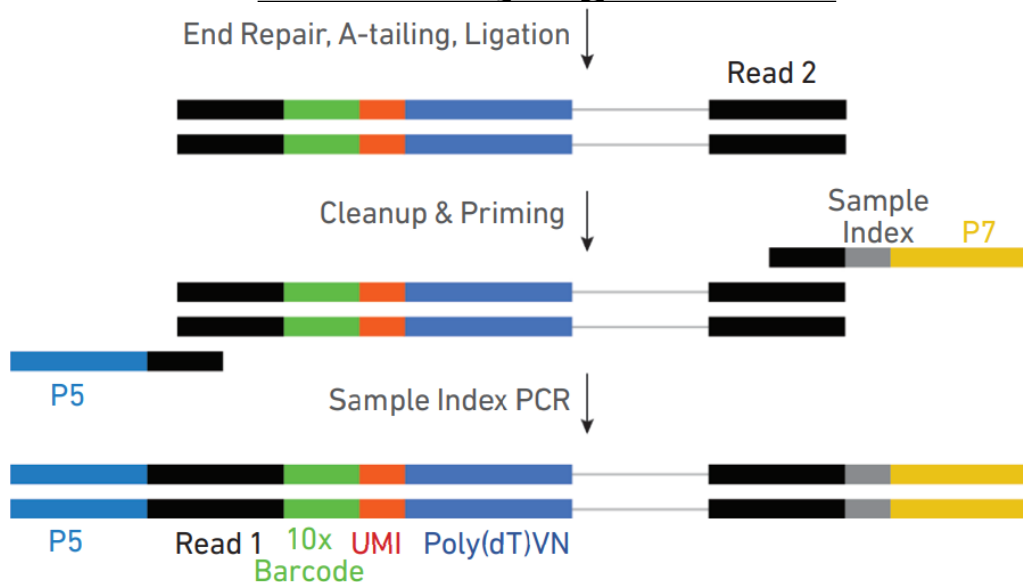
Illumina Adapter – Foundation for Cluster Generation

Sample Prep

Addition by Ligation

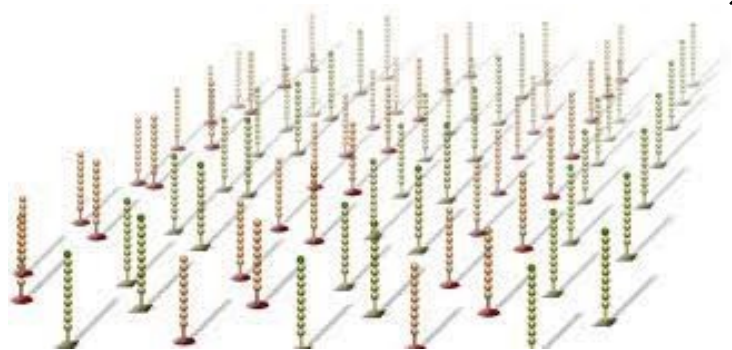


Addition by Ligation/PCR



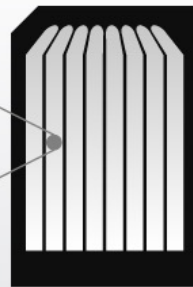
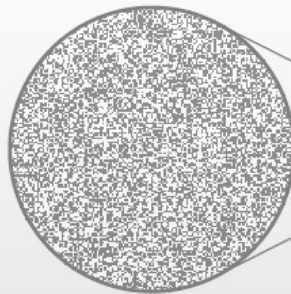
Illumina Flowcell – Cluster Generation

Cluster Generation

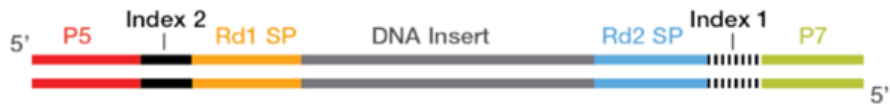
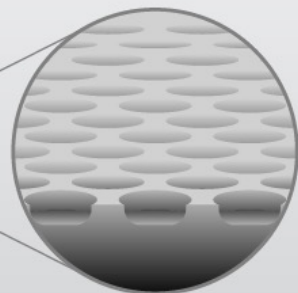
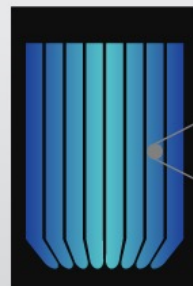


Flow cell contains lawn of oligos complementary to the P5/P7 sequences

Random Flow Cell



Patterned flow cells



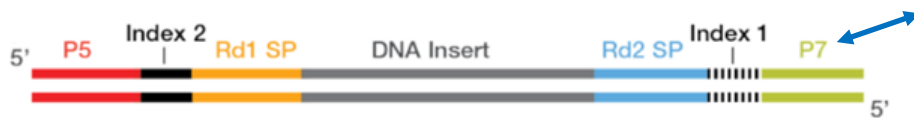
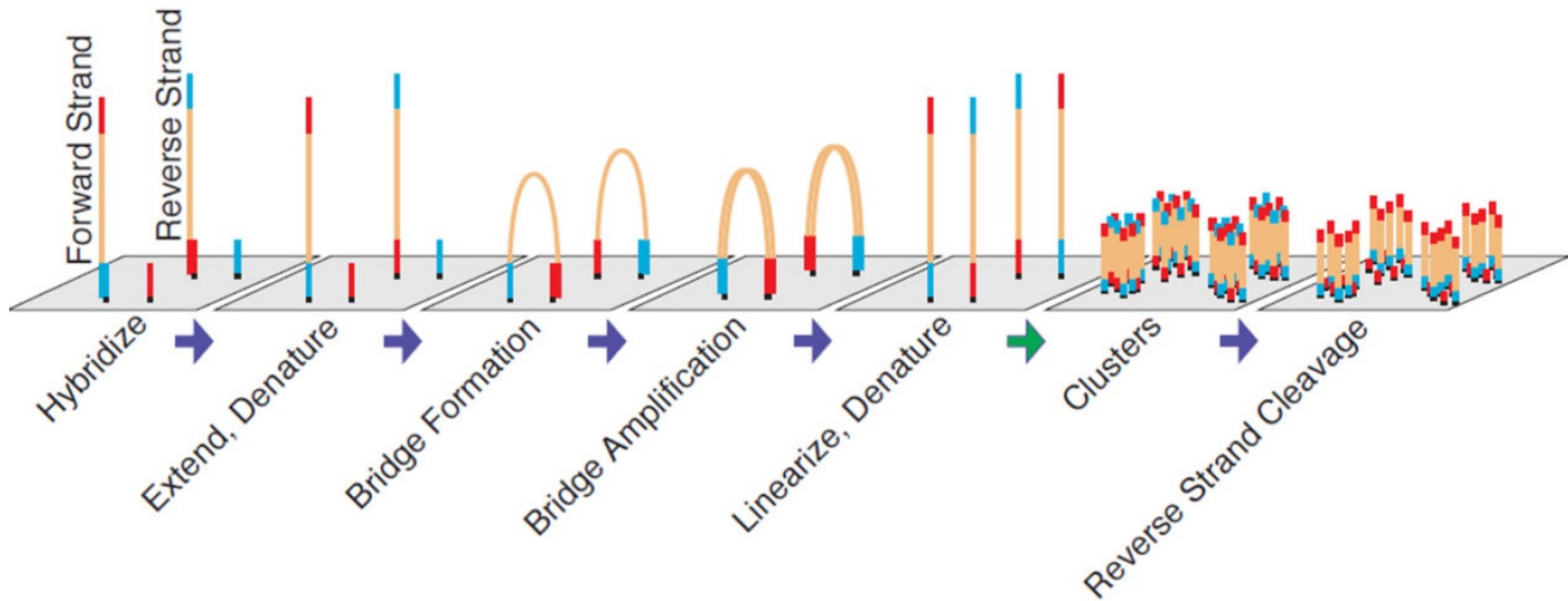
2

3

Illumina Flowcell – Cluster Generation

Cluster Generation

Sequencing



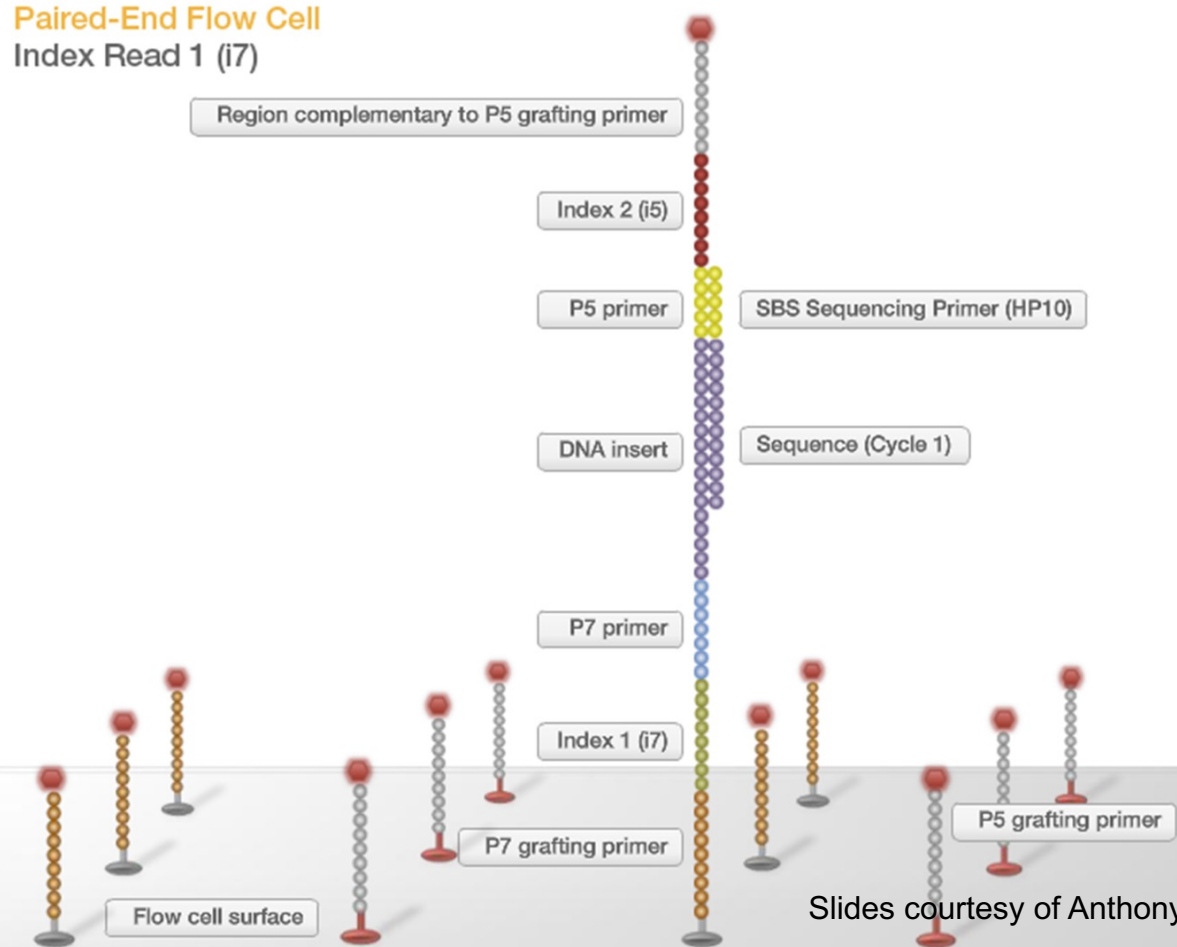
P7 is blue in above diagram

Slides courtesy of Anthony Miller

Illumina Sequencing By Synthesis (SBS) – Read 1

Sequencing

- Initiated by HP10 primer (Rd1 SP)
- Fluorescently labeled and reversibly terminated nucleotides
- Clusters are excited by light, fluorescent signal emitted
- Terminator remove for next round of nucleotide addition

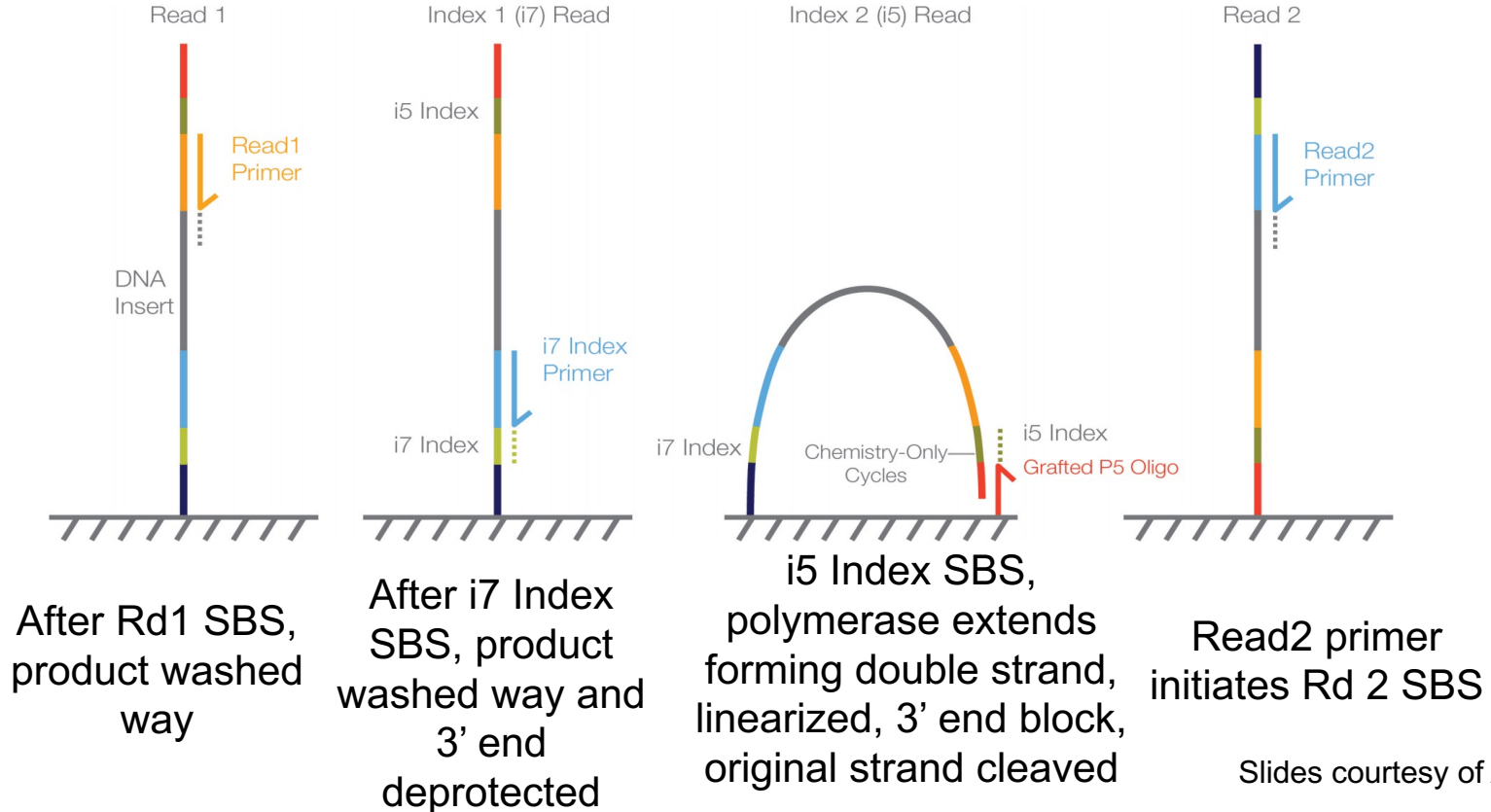


Slides courtesy of Anthony Miller

Illumina Sequencing By Synthesis (SBS) – Index(s) and Read 2

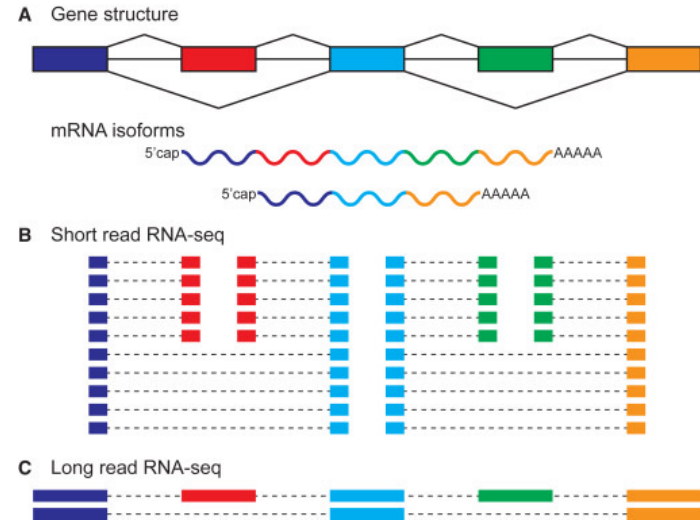
NovaSeq 6000 w/ v1 chem, MiniSeq w/ rapid chem, MiSeq, HiSeq 2000/2500

Figure 2 Dual-Indexed Sequencing on a Paired-End Flow Cell (Workflow A)



Alternative Splicing

- Genes are not continuous coding.
 - Exons: coding regions
 - Introns: silent.
- Different combinations of exons yields **isoforms**.
- This phenomenon is called alternative splicing.

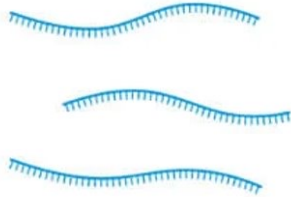


Types of Sequencing

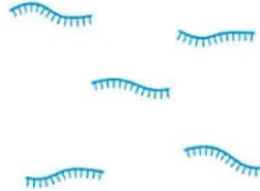
- RNA: quantifying transcription
- WGS: sequencing the genome
- Epigenetics (sequencing DNA):
 - ChIP-seq: transcription factors and histone modification
 - ATAC-seq: chromatin accessibility
 - Hi-C: 3D chromatin structure.
 - bisulfite sequencing (DNA methylation).

RNA Sequencing

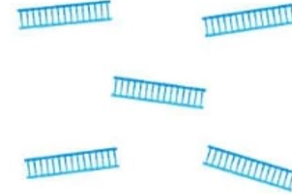
① Isolate RNA from samples



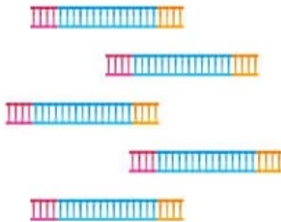
② Fragment RNA into short segments



③ Convert RNA fragments into cDNA



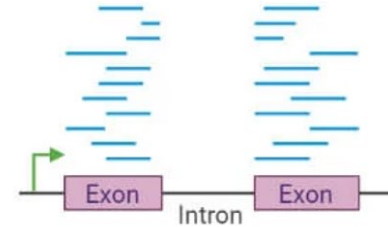
④ Ligate sequencing adapters and amplify



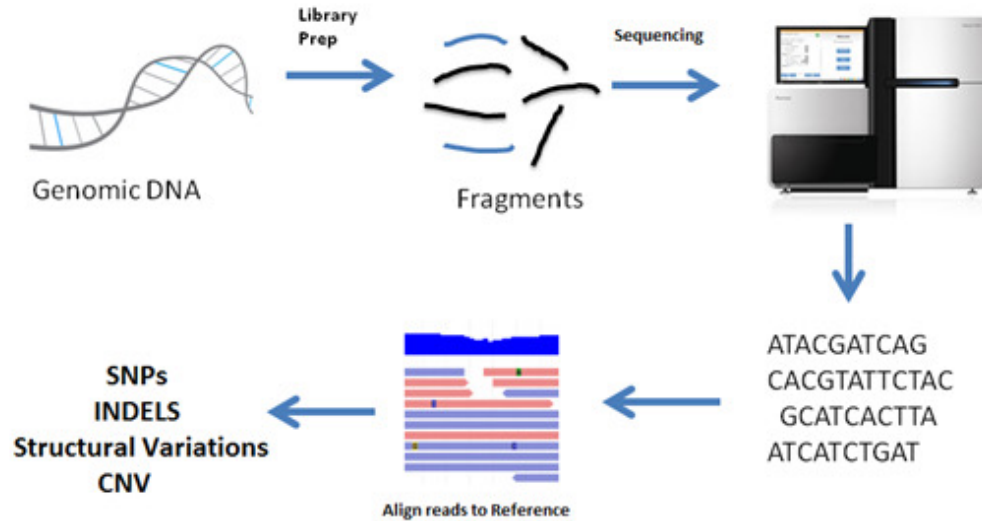
⑤ Perform NGS sequencing



⑥ Map sequencing reads to the transcriptome/genome



Whole Genome Sequencing

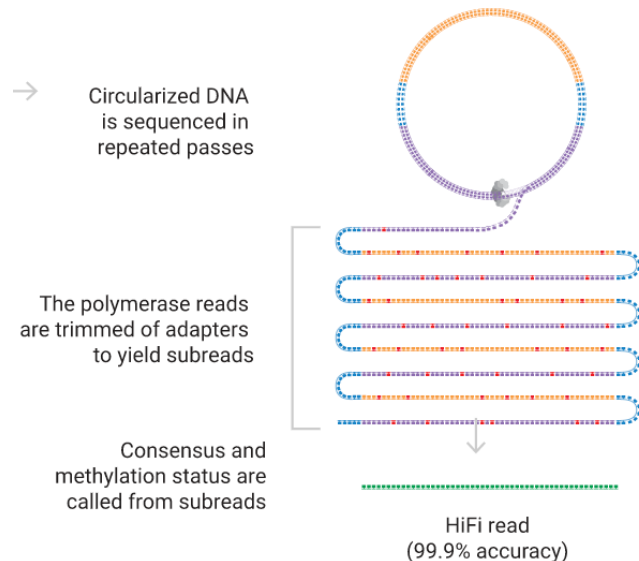


PacBio HiFi Sequencing

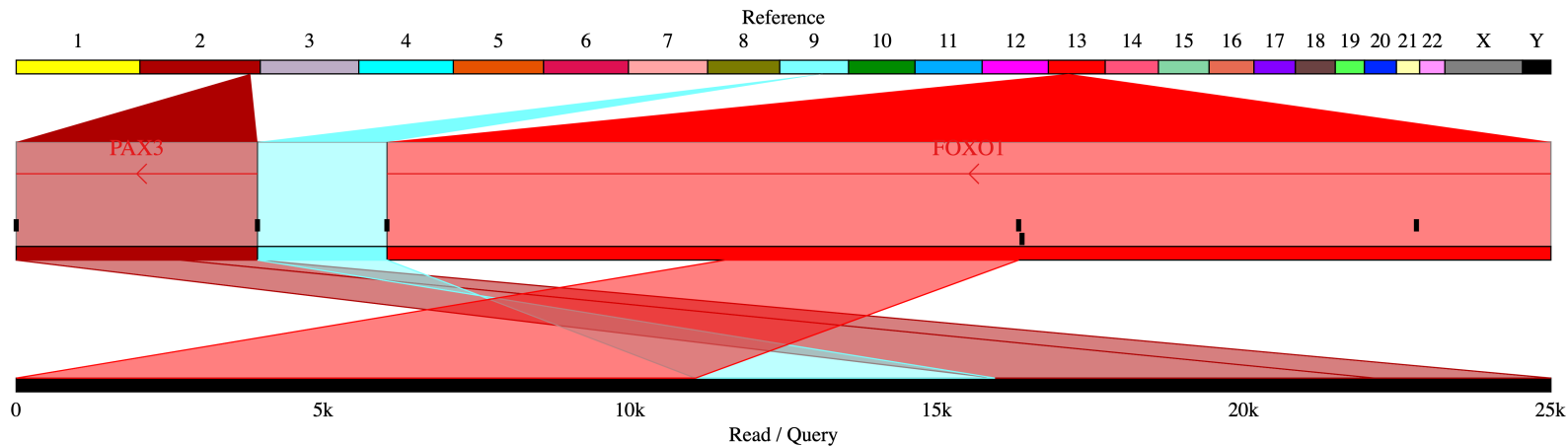
- Illumina platforms are restricted to 300bp reads, and 500bp fragments.
- Great overall, but not that good for isoform or structure variation detection.

PacBio Long Read

- 10k-30kb read length.
- Amazing for:
 - genome assembly,
 - CNV and SV calls,
 - isoform analysis (RNA),
 - Microsatellite repeats,
 - Haplotyping!!!
- Single cell version is being developed.
- Terrible error rate one pass (10%)
- Consensus (10 times or more) reads are very accurate.

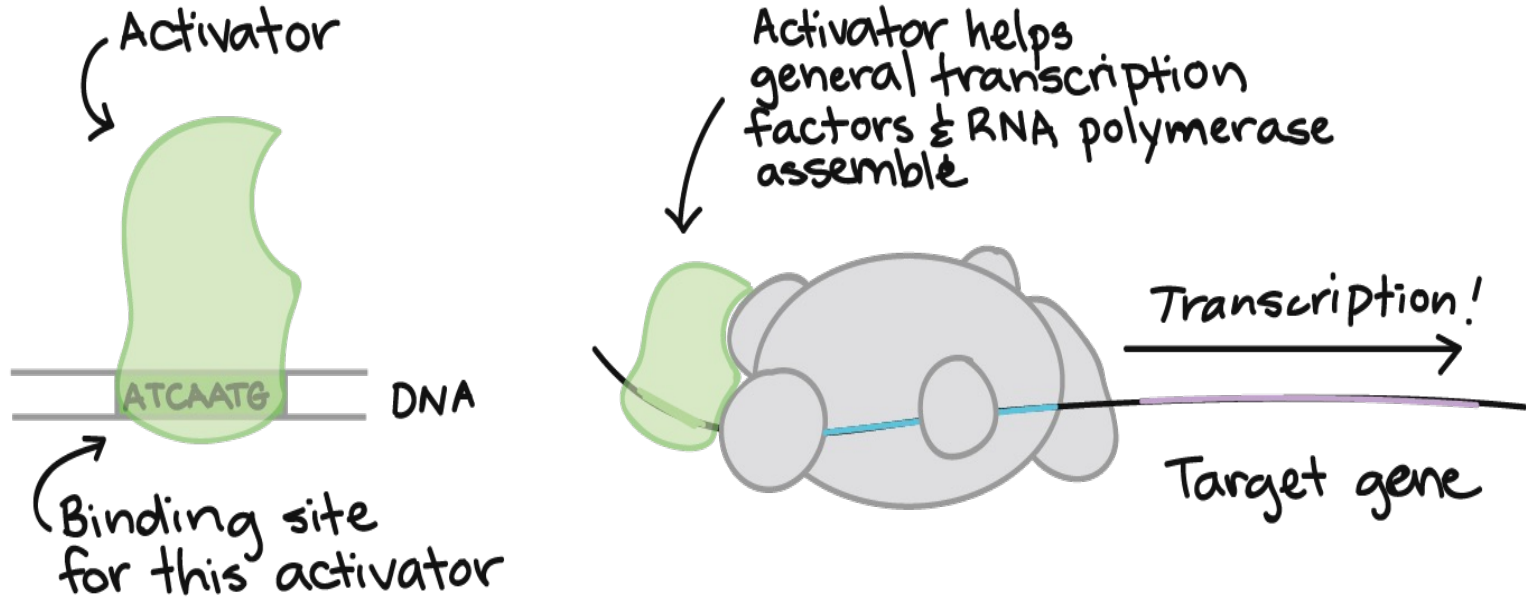


PacBio Mapping of a cancer fusion gene



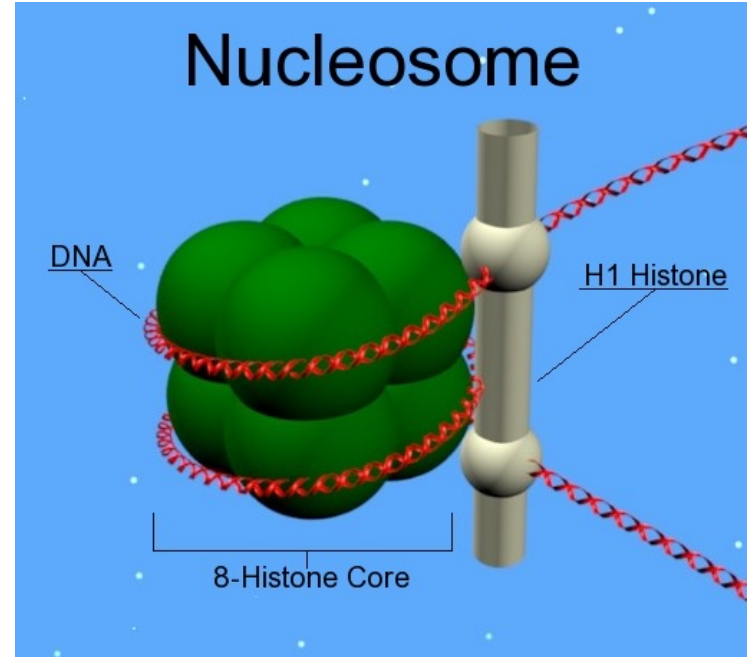
Transcription Factor

Proteins that modulate gene transcription



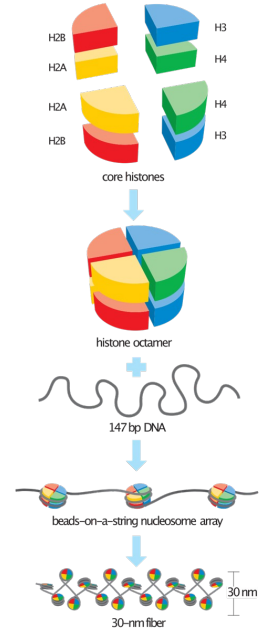
Nucleosomes

- Nucleosome: basis of chromatin: 147bp of DNA wrapped around a nucleosome (8 histone).
- Heterochromatin = tightly packed nucleosomes + DNA wrapped around it, usually repressed.
- Euchromatin = "free" chromatin, usually transcriptionally active.

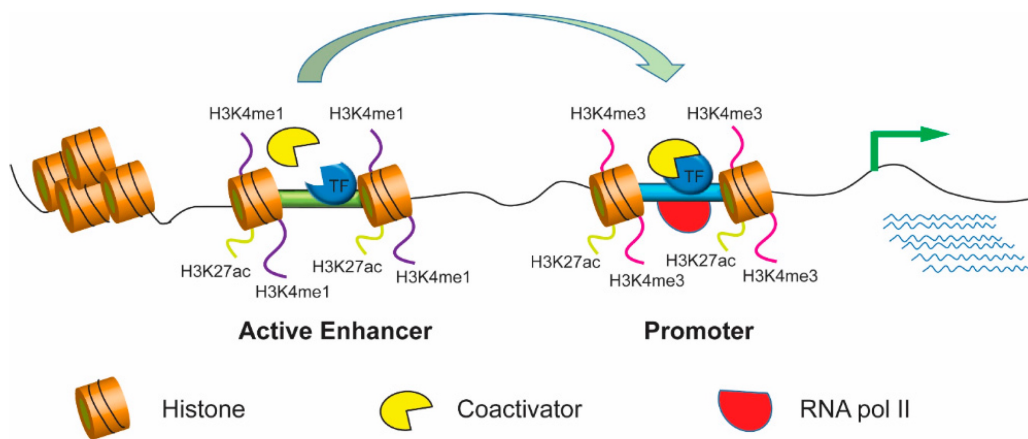
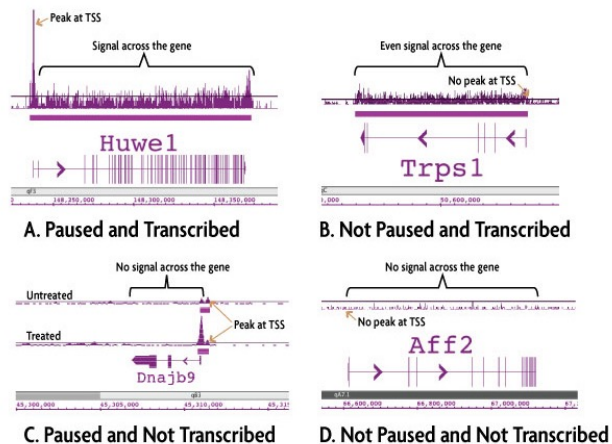


Nucleosomes

- Histones can be modified chemically. Those modifications activate or repress expression.
- For example:
 - H3K9me3 = repressive (3rd histone, 9th lysine, methylated 3 times)
 - H3K27ac = active.
- Histone Modification are protein and therefore can be assessed with ChIP-Seq.



Transcription Regulation

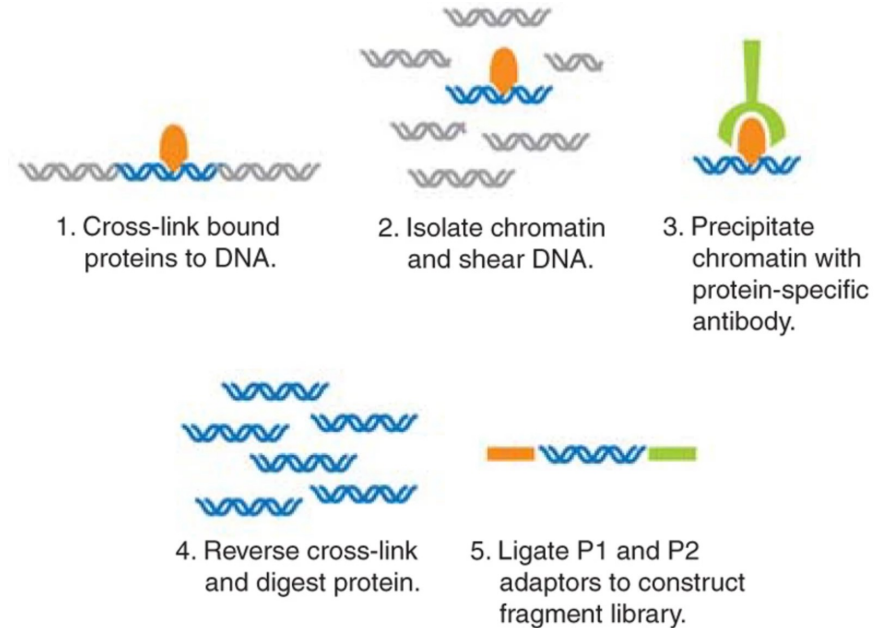


RNA Pol II ChIP-Seq

ChIP-Seq

- Chromatin ImmunoPrecipitation:

1. cross link
2. sonication
3. ChIP
4. Remove protein



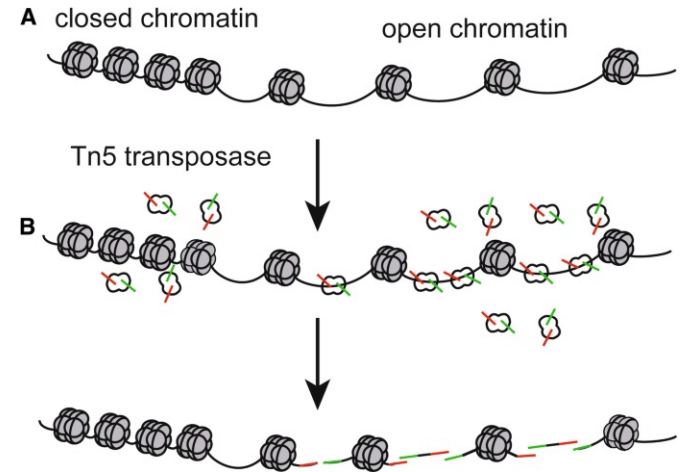
Cut-and-Tag

- Alternative to ChIP-seq
- Uses Tn5 enzyme instead of sonication
 - Cut where there isn't a nucleosome.

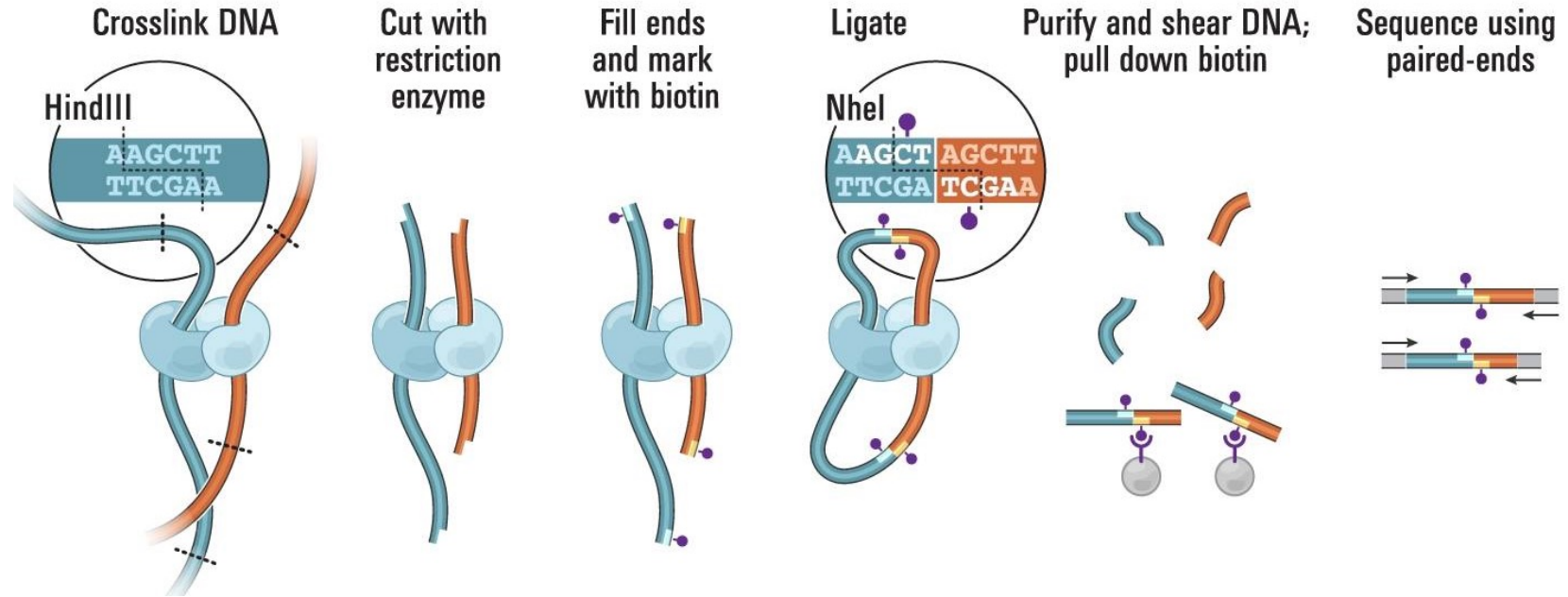


ATAC-seq

- Assay for Transposase Accessible Chromatin
- Sub-nucleosome (<150bp), mono-nucleosome (~150bp), di-nucleosome(300bp) reads.
- Nucleosome positioning, eviction etc is very important in epigenetics.
- cut-and-run: MNase instead of Tn5.
 - MNase: cut into a single nucleotide and start digesting them, until stopped by a nucleosome.



Hi-C



Variation

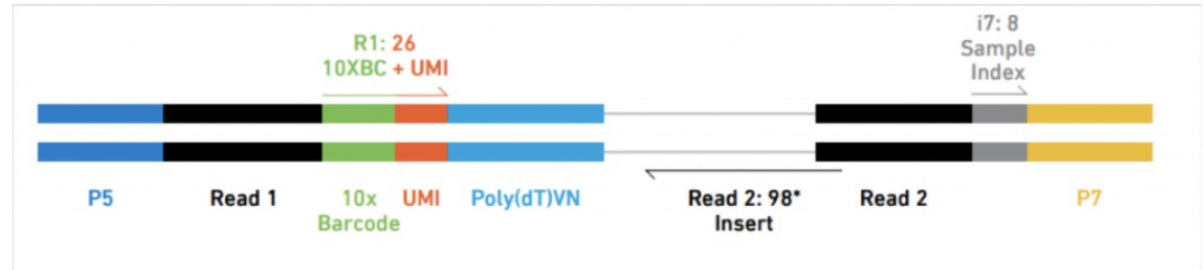
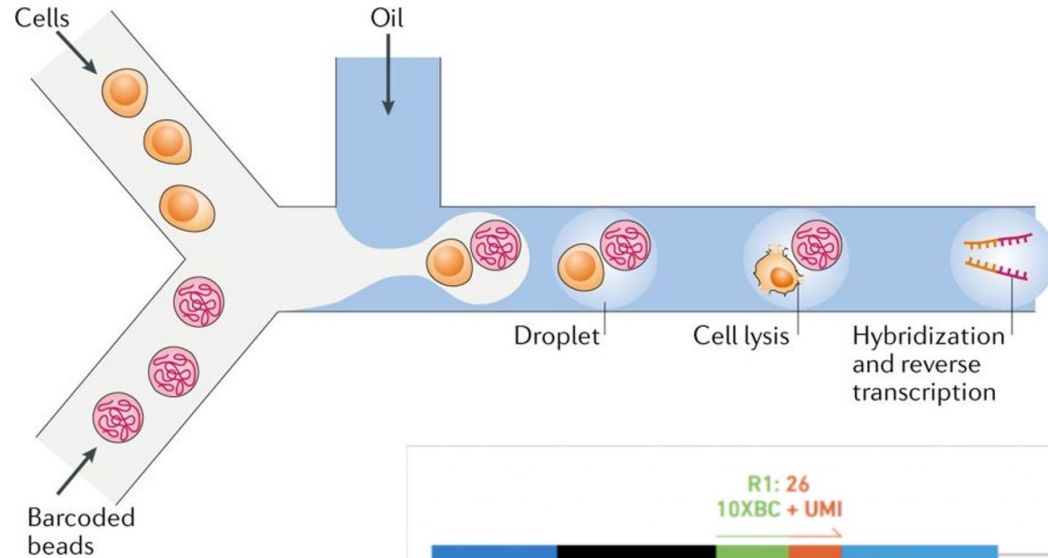
- Combine Hi-C + ChIP = Hi-ChIP
 - long-range interactions/loops
- Cut&Tag, Cut&Run
 - Alternative to ChIP and ATAC-seq

And many more...

Single Cell

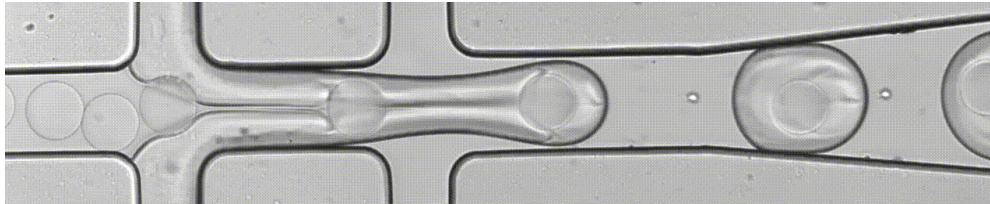
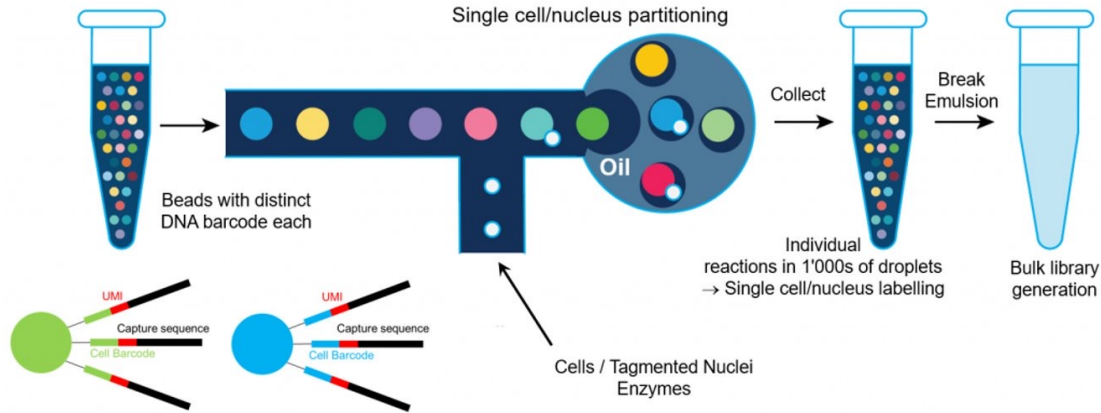
- sc-RNA-seq
- sc-ChIP (not really)
- sc-ATAC-seq
- sc-Hi-C
- spatial transcriptomics

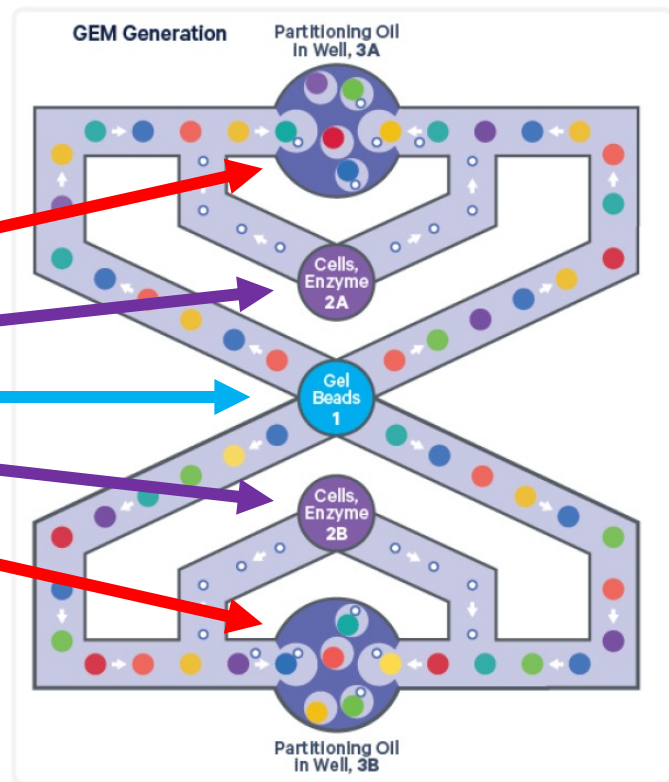
sc-RNA-seq (Droplet based)



10X GENOMICS

sc-RNA-seq (Droplet based)





Chromium controller

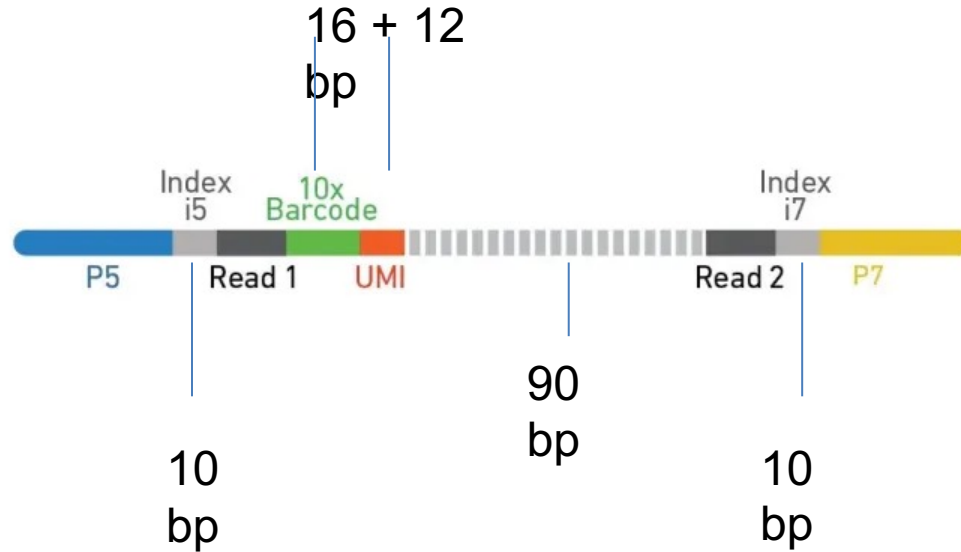


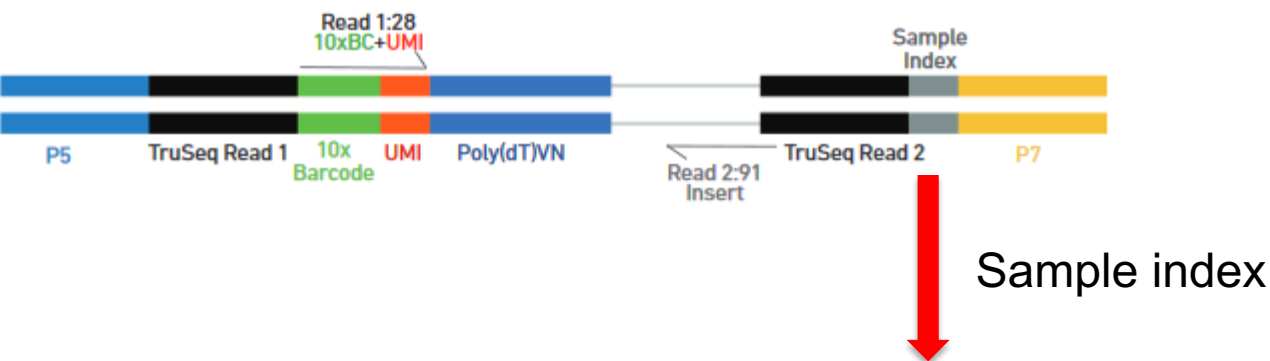
Slides courtesy of Matt Cannon

Sequence data

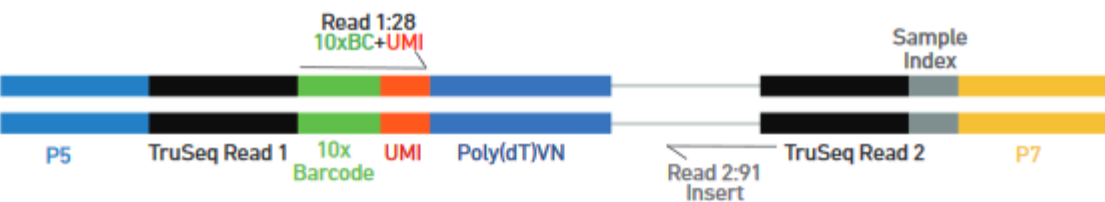
“Run recipe” = 28 + 10 + 10 + 90

This gets sequenced on standard Illumina sequencer (about 150bp per read)





```
@A00498:356:H53CMDRXY:1:2101:1371:1016 1:N:0:CATGCGAT
CATGCGAT
+
FF:FFFFF
@A00498:356:H53CMDRXY:1:2101:1443:1016 1:N:0:CATGCGAT
CATGCGAT
+
FFFFFFFFF
@A00498:356:H53CMDRXY:1:2278:31349:37059 1:N:0:ACCCGACG
ACCCGACG
+
FFFFFFFFF
```



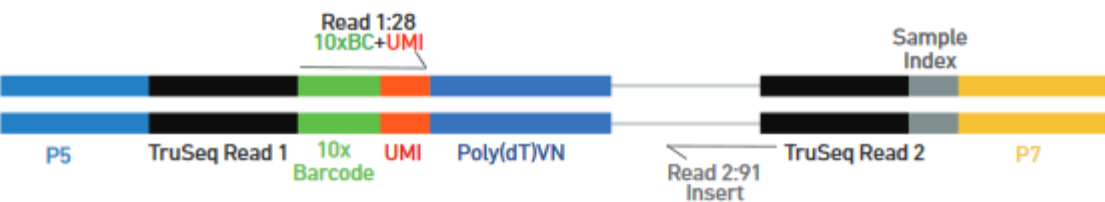
Read 1

```
@A00498:356:H53CMDRX:1:2101:4390:1266 1:N:0:CATGCGAT
AATAGAGAGTCTGTACTTTGACAACCGT
+
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
@A00498:356:H53CMDRX:1:2101:15130:1141 1:N:0:CATGCGAT
ATCTCTATCCAACTGACAGTTAACTGGT
+
:FFFFFFF,FFFFFFFFFFFFFFFFFFFFFF
@A00498:356:H53CMDRX:1:2101:15167:1141 1:N:0:CATGCGAT
GATGGAGTCGTGGACCACTCGGGCAGCC
+
FFFFFFFFFFFFFFFF:FFFF:FFFFFFFFFFFF
@A00498:356:H53CMDRX:1:2101:15528:1141 1:N:0:CATGCGAT
TCTACATAGTCTACCATCCAAACGAAC
+
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
```

10X barcode (cell)

TCTACATAGTCTACCA

Unique molecular identifier (UMI)

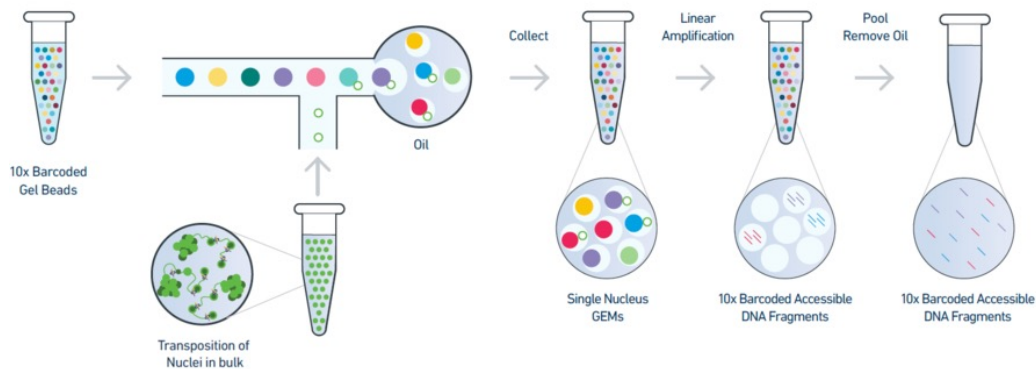


Read 2 = insert = transcript of interest

```
@A00498:356:H53CMDRX:1:2101:28221:1204 2:N:0:CATGCGAT
AGATGATCTGTTCAAGCGTAGGATGTTGAAGTCCCCAACTGTTATTGTGTTGGAGTCTATTTGTCTCTTTAGGTTTAATAATATT
+
FFF:FFFFFFFF:FFFFF:FFF,F,FFFF,F,:FFFFFFFFFFFFFFFF:FFF,FFFFFFFFFFFFF:,FFFF,FF:FFFFFFFFFFFFFFFF,FFFF
@A00498:356:H53CMDRX:1:2101:28583:1204 2:N:0:CATGCGAT
ATGCCCTAGCCCACTTCTTACCACAAGGCACACCTACACCCATTATCCCCATACTACTTATAATCGAAACCATCAGACTACACATTCAACC
+
F::FFF,F,FFFFFFFF,,FF:FFFFFF,FFF:FFF,FFFF,:FF::,F:FF,:F,,F,,F:FFFFFF:FFFFFF,,FFFF,:F,FF:FF,
@A00498:356:H53CMDRX:1:2101:28673:1204 2:N:0:CATGCGAT
GTGAAGAGGATCTTGAATTTCTTAATGCATCCAAAGCCTTCTGGCAAACCATTGCCGAAATCTAAAAAACTTGTAGCAAAGGCAGTAAA
+
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF:FFFF,FFFFFFFFFFFFFFFF:FFF:FFFFFFFFFFFFFFFF:FFFFFFFFFFFFFFFF
@A00498:356:H53CMDRX:1:2101:28709:1204 2:N:0:CATGCGAT
CGCGAGGTGGGGGCGTCGTGTAAGCAGCGGAGGATGGGGGGGCGGTGCACGTGGGTGGGCGTGGCTGAGATCTAAGTGTCTCTGCAGCTGTG
+
FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
```

scATAC-seq

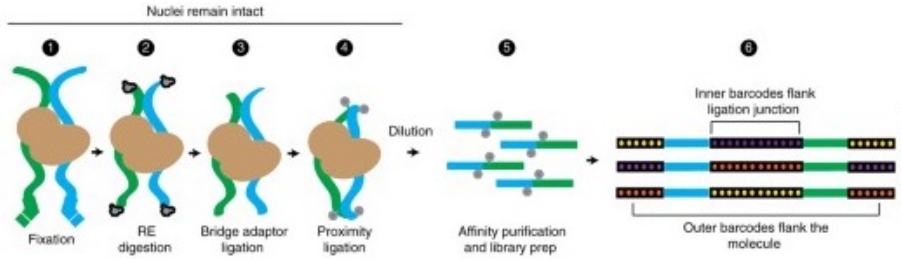
- Droplet based.
- tn5 in oil.



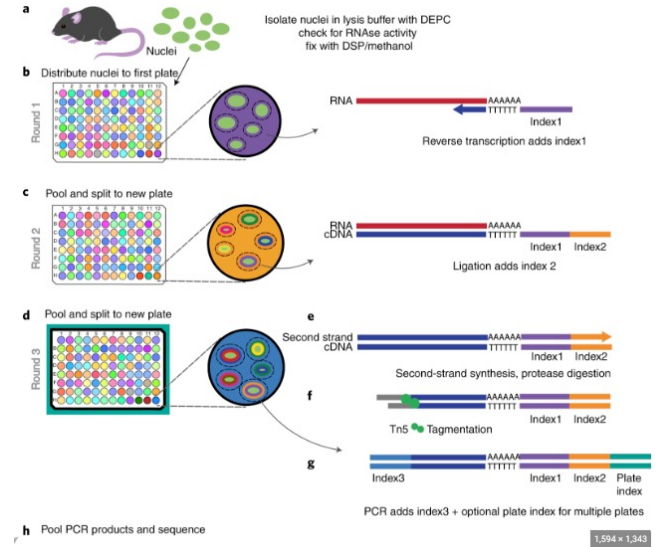
scChIP-seq

- Droplet based
- Uses cut-and-tag (oil contain tn5+antibody)
- Work somewhat decent for K27ac (easiest to ChIP)

sc-Hi-C



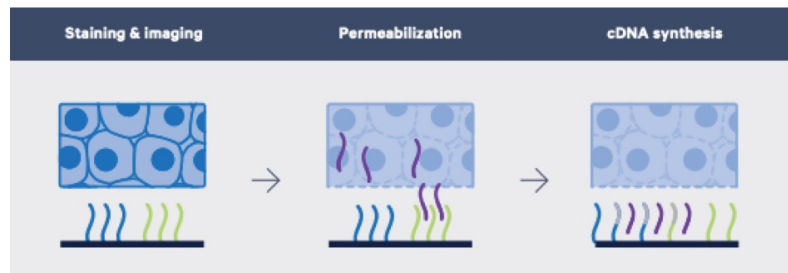
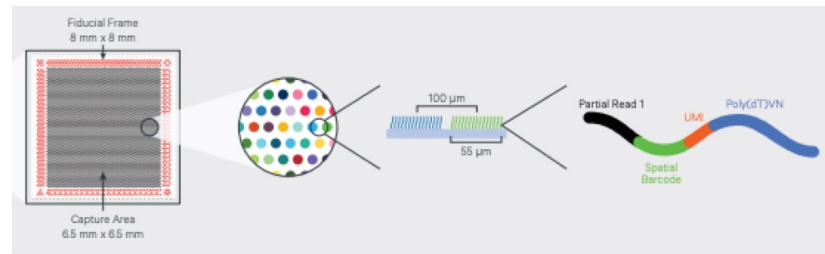
Proximity ligation with nuclei intact



Combinatorial Indexing

Spatial Transcriptomics

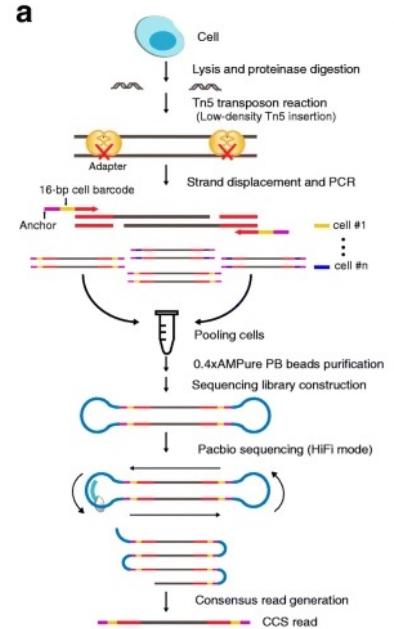
- Each microscope slide includes ~5000 “spots”.
- Each spot is about ~55µm (1-10 cells).
- mRNAs released from cell,
- mRNAs bind to spatially barcoded oligos,
- RT to produce cDNA.
- Not exactly single cell.



sc-PacBio

sc-WGS and sc-RNA have both been done and reported on PacBio HiFi platform.

Fig. 1



Costs

- Bulk: ~\$200-\$500 per sample (including reagents but not time cost)
- Hi-C is more expensive due to sequencing depth.
- Pacbio \$3000 per sample for bulk.
- sc-RNA-seq could cost over \$5000 per sample.

In the next 5 years...according to Katie

- Freedom from <expensive> hardware (e.g., Parse and Fluent Biosciences)
- Decreasing sequencing costs
 - ...therefore increase number of cells to be analyzed
- Enhanced/more-sensitive mRNA capture rates
- More user-friendly computational pipelines
- Multi-omics approaches
 - Proteomics
 - CRISPR
 - Non-coding RNA
 - Epigenomics
 - Long-reads
- Clinical applications

Special Thanks to...



Ben Stanton
PI, CCCR



Ryan Roberts,
PI, CCCR



Matt Cannon
bioinformatician



Katherine Miller
PI, IGM



Anthony Miller
Director of TechDev, IGM