### Genomes and Genomic Tools

**Iestyn Whitehouse** 

### What is a genome?

Telomeres
Centromeres
rDNA
Introns
Exons
Promoters
Repetitive DNA
"Junk DNA"
Transposable elements

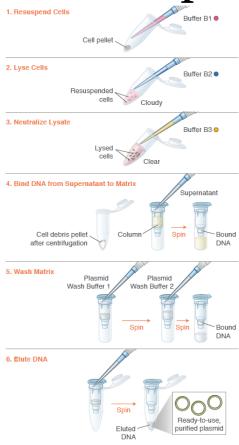
## How do you figure out the composition of a genome?

How much single copy DNA? How much repetitive DNA?

### What does a mini-prep do?

How does a mini-prep work?

### How does a mini-prep work?



### How does a mini-prep work?

#### Buffer P1

50 mM Tris-HCl pH 8.0 10 mM EDTA 100 μg/ml RNaseA

#### **Buffer P2**

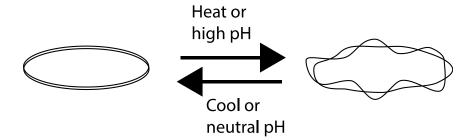
200 mM NaOH 1% SDS

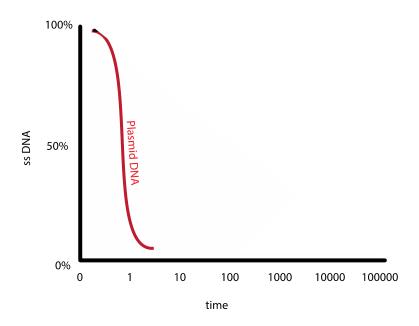
#### Buffer P3

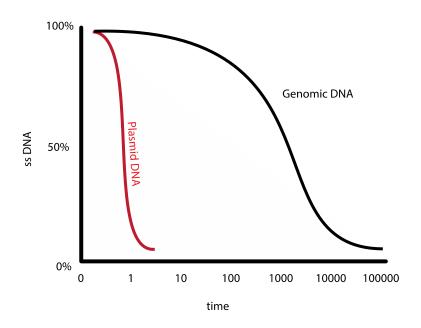
3.0 M potassium acetate pH 5.5

Spin in centrifuge

### mini-prep







### How does a mini-prep work?

Buffer P1

50 mM Tris-HCl pH 8.0 10 mM EDTA 100 μg/ml RNaseA

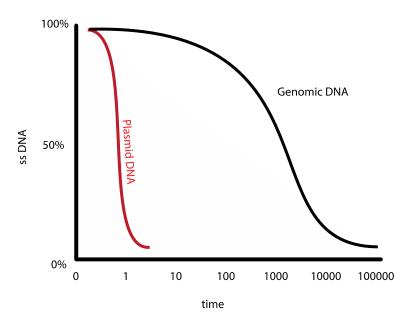
Buffer P2

200 mM NaOH 1% SDS

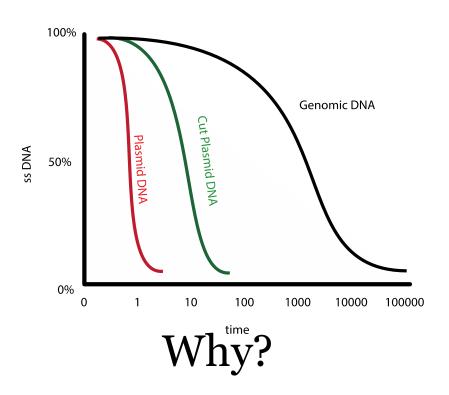
Buffer P3

3.0 M potassium acetate pH 5.5

Spin in centrifuge



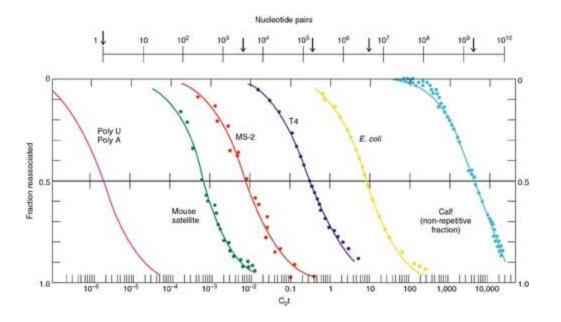
What would happen if you cut the plasmid?



# Annealing is sensitive to concentration and time

Cot analysis

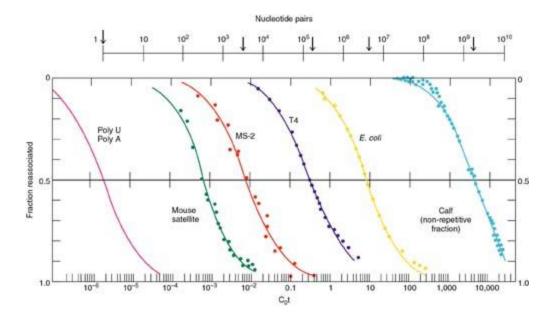
### Reassociation scales according to genome size



Why?

### Bigger genomes can have more combinations of sequence elements

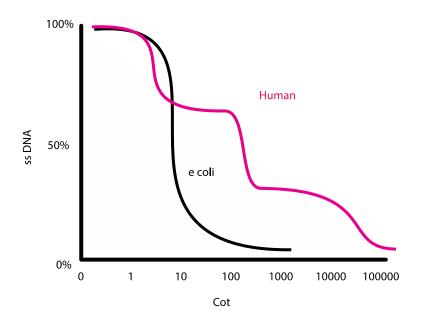
### Not all DNA in the genome is equal!



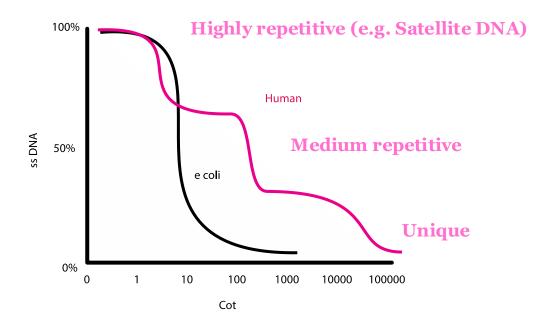
Why does mouse satellite DNA reassociate quickly?

Eukaryotic genomes often contain large quantities of repetitive DNA sequence.

### Composition of eukaryotic genomes

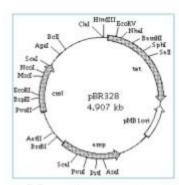


#### Composition of eukaryotic genomes

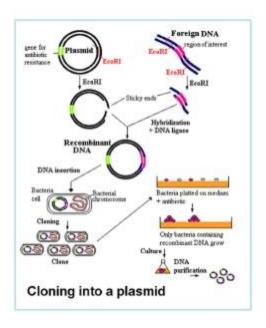


### How do you sequence DNA?

### Preliminaries: Cloning



pBR322: is a vector, an engineered phage. It can reproduce itself inside a bacterial host and do nothing else.





Sub-clone DNA of interest

Transform bacteria

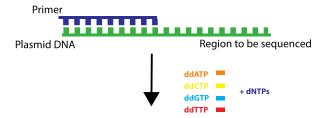
Pick colony

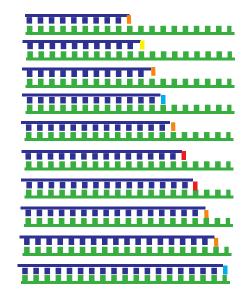
Miniprep

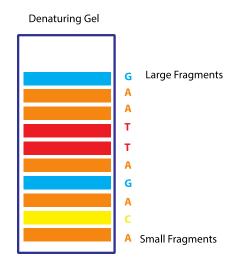
Submit purified plasmid with desired primer for sequencing

Why clone?

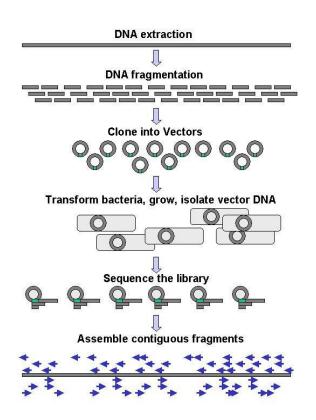
#### Sanger sequencing







#### Sequencing the human genome

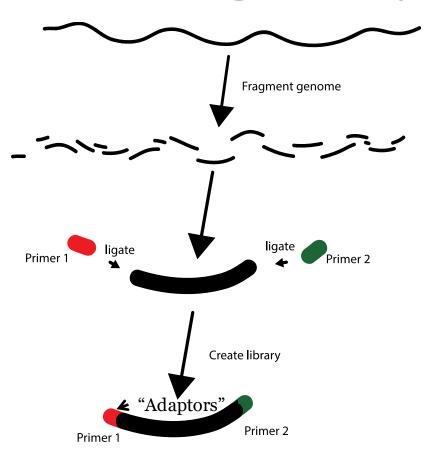


### High throughput "Deep" sequencing

### High throughput "Deep" sequencing

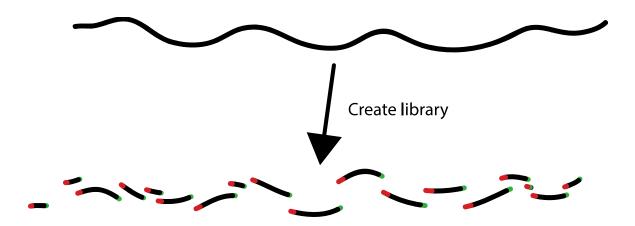
The key to deep sequencing is the generation of billions of individual "clones" without use of microorganisms!

### Generating a library



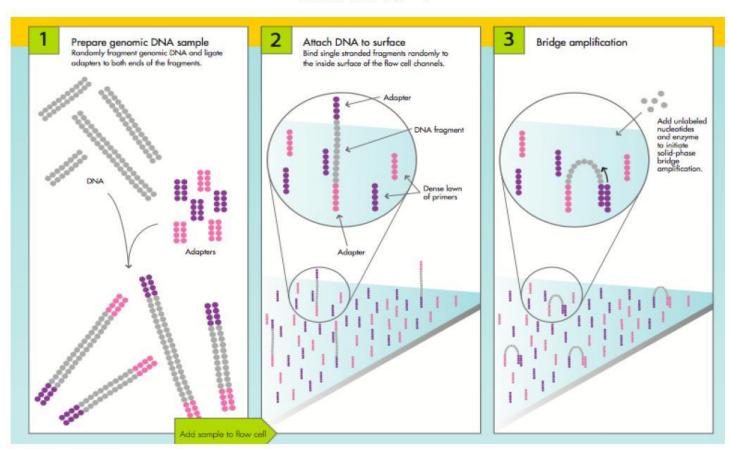
### Generating a library

Generate a random mixture of sequence elements



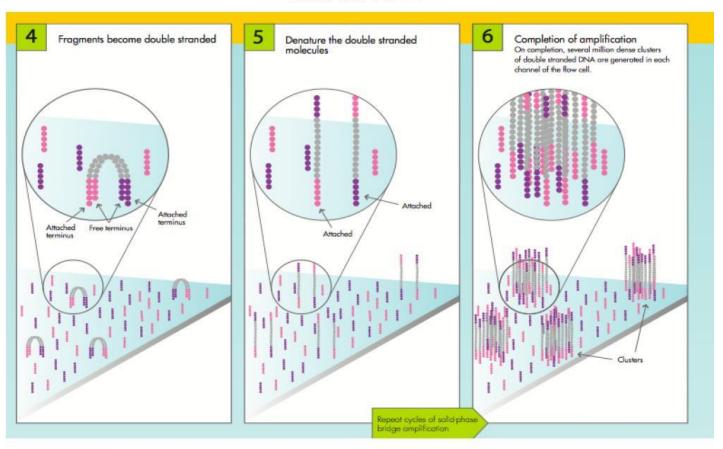
How do you clonally amplify and sequence DNA molecules from the whole population?

#### Illumina 1



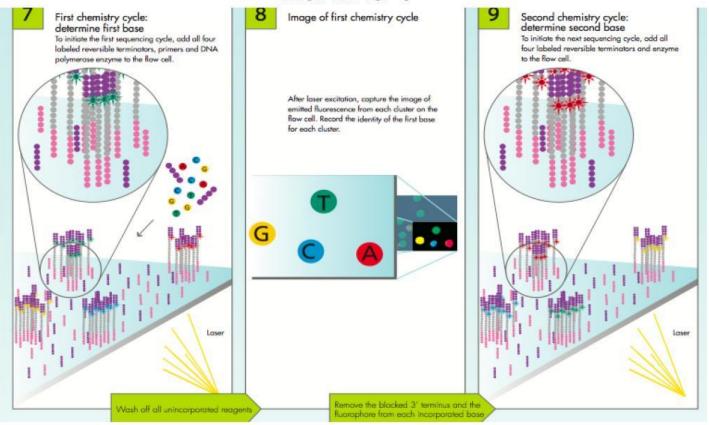
ECE/BioE 416 Lecture 24

#### Illumina 2

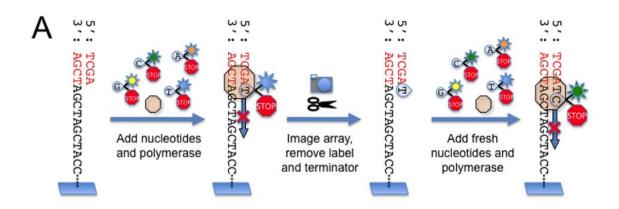


ECE/BioE 416 Lecture 24

#### Illumina 3



### Sequence by synthesis



# Deep sequencing can generate billions of short reads

Use a variety of algorithms to "map" the sequence reads to the genome.

#### FastQ files:

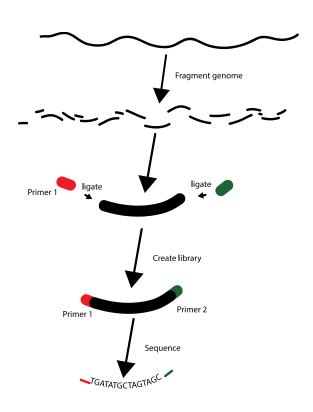
#### \$ head N2-copy-num.R1.fastq

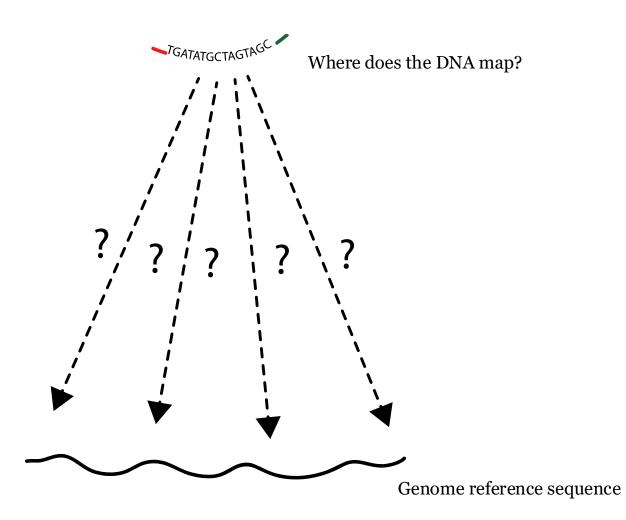
quality value characters in left-to-right increasing order of quality (ASCII):

#### Mapping Data

FastQ file \_\_\_\_\_\_ Processed file for specific purpose

## Read mapping:





Read length matters!

## Can you name a useful restriction enzyme?

### **EcoRI**

GAATTC CTTAAG **EcoRI** 

TaqI

GAATTC CTTAAG TCGA AGCT

## How often will you find an EcoRI site in DNA?

 $4^{6} = 4096$ 

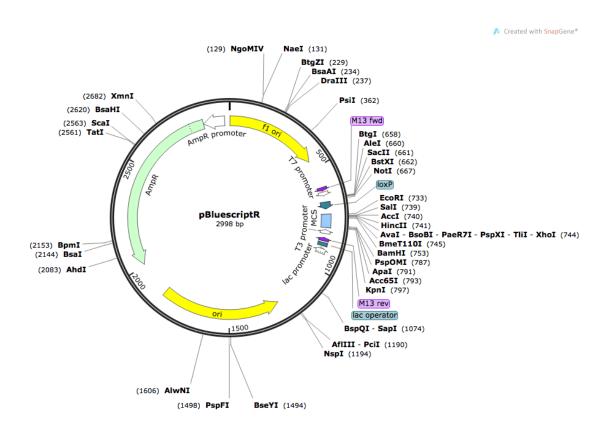
GAATTC CTTAAG

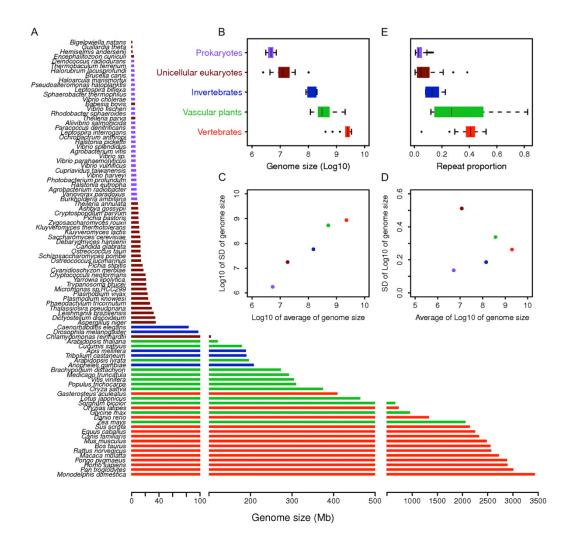
## How often will you find a TaqI site?

TCGA AGCT

 $4^4 = 256$ 

## Why is EcoRI typically more useful than TaqI?



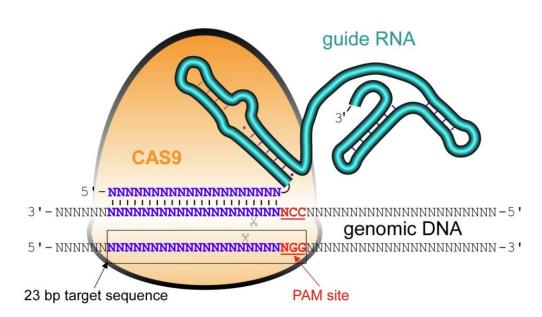


#### Motif length Possible combinations

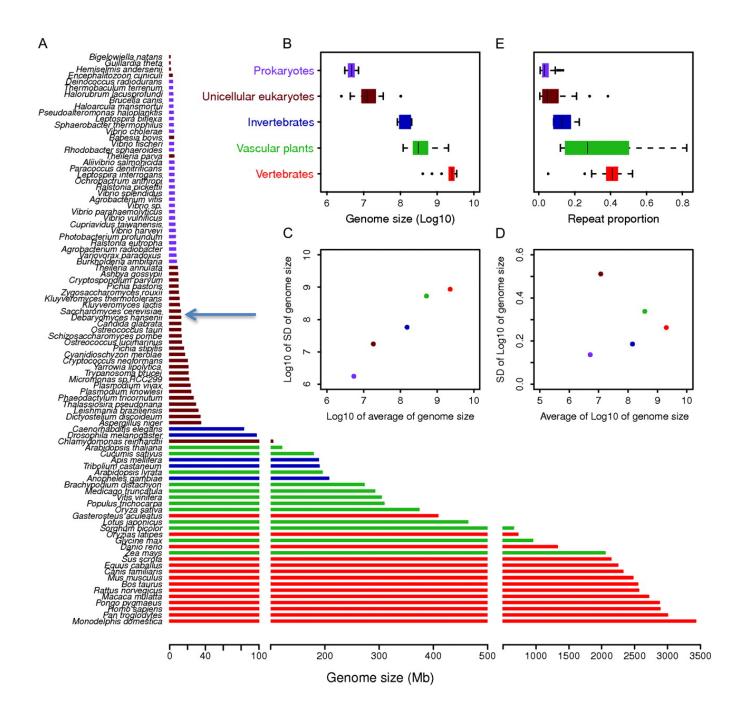
```
4
16
1
2
        64
3
        256
4
5
        1,024
        4,096
6
        16,384
8
        65,536
9
        262,144
10
        1,048,576
       4,194,304
11
        16,777,216
12
        67,108,864
13
        268,435,456
14
        1,073,741,824
15
16
        4,294,967,296
        17,179,869,184
17
18
        68,719,476,736
        274,877,906,944
19
        1,099,511,627,776
20
        4,398,046,511,104
21
        17,592,186,044,416
22
23
        70,368,744,177,664
24
        281,474,976,710,656
        1,125,899,906,842,620
25
26
        4,503,599,627,370,500
        18,014,398,509,482,000
27
28
        72,057,594,037,927,900
        288,230,376,151,712,000
29
        1,152,921,504,606,850,000
30
```

## Why is CrisprCas9 such a useful enzyme?

## Why is CrisprCas9 such a useful enzyme?



How long a "read" would you need to match uniquely to the budding yeast genome?



#### Motif length Possible combinations

```
1
        4
16
2
3
        64
        256
4
5
        1,024
        4,096
6
        16,384
8
        65,536
9
        262,144
10
        1,048,576
11
        4,194,304
                                 12bp!
        16,777,216
12
        67,108,864
13
        268,435,456
14
        1,073,741,824
15
16
        4,294,967,296
        17,179,869,184
17
18
        68,719,476,736
        274,877,906,944
19
        1,099,511,627,776
20
        4,398,046,511,104
21
        17,592,186,044,416
22
23
        70,368,744,177,664
24
        281,474,976,710,656
        1,125,899,906,842,620
25
26
        4,503,599,627,370,500
        18,014,398,509,482,000
27
28
        72,057,594,037,927,900
        288,230,376,151,712,000
29
        1,152,921,504,606,850,000
30
```

Entered nucleotide pattern: GTAGGCAATCTC

Total Hits: 1

Sequences Searched: 17

Dataset: COMPLETE S. cerevisiae Genome DNA

Strand: both strands

Entered nucleotide pattern: GGGAAAATTATC

Total Hits: 6

Sequences Searched: 17

Dataset: COMPLETE S. cerevisiae Genome DNA

Strand: both strands

Entered nucleotide pattern: AAATTTAAATTT

Total Hits: 34

Sequences Searched: 17

Dataset: COMPLETE S. cerevisiae Genome DNA

Strand: both strands

Genomes are not composed of random sequence: some regions are very biased in sequence composition so have a low complexity.

This means that you typically need sequencing reads much longer than what's predicted just from genome size.

Sequence complexity describes the number of possible "word" combinations within a given sequence

GATAGATAGATCATCA

Complex sequence

AAAAAAAAAAA Simple sequence

Sequence complexity describes the number of possible "word" combinations within a given sequence

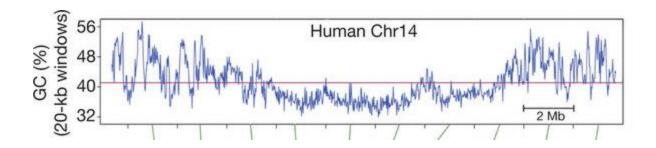
#### 10kb plasmid composed of GATC:

Need a sequence read of **7bp** for a unique match  $(4^7 = 16,384)$ 

#### 10kb plasmid composed of AT only

Need a sequence read of **14bp** for a unique match  $(2^{14} = 16,384)$ 

Different regions of the genome have different base composition



#### "mapability"

Mapability is a measure of both sequence complexity and uniqueness.

Typically (but not always) sequences with low complexity have low mapability.

**Sequence 1: TGATAGATCGATCGATCGATCGA** 

**Sequence 2: AGTCGATTCGAT** 

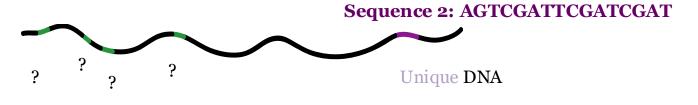
#### Motif length Possible combinations

```
1
       4
       16
2
       64
3
        256
4
5
       1,024
       4,096
6
       16,384
8
       65,536
9
        262,144
10
       1,048,576
       4,194,304
11
       16,777,216
12
       67,108,864
13
       268,435,456
14
       1,073,741,824
15
16
       4,294,967,296
       17,179,869,184
17
18
       68,719,476,736
       274,877,906,944
19
       1,099,511,627,776
20
       4,398,046,511,104
21
       17,592,186,044,416
22
23
       70,368,744,177,664
24
       281,474,976,710,656
       1,125,899,906,842,620
25
26
       4,503,599,627,370,500
       18,014,398,509,482,000
27
28
       72,057,594,037,927,900
       288,230,376,151,712,000
29
30
       1,152,921,504,606,850,000
```

#### Read length and "mappability"

#### **Sequence 1: TGATAGATCGATCGATCGATCGA**

Repetitive DNA

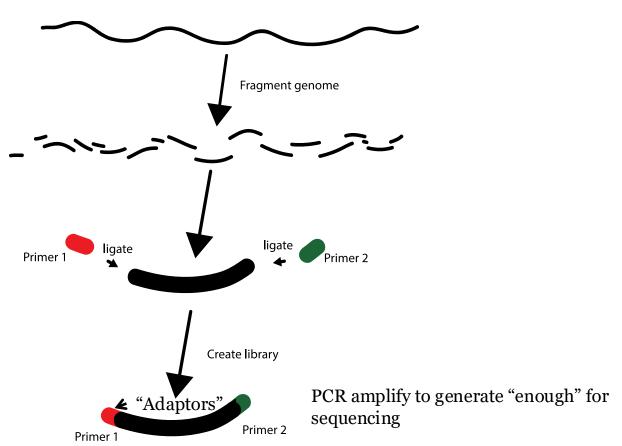


### Genomics approaches

Nearly all approaches create a "library" with adapters on either end.



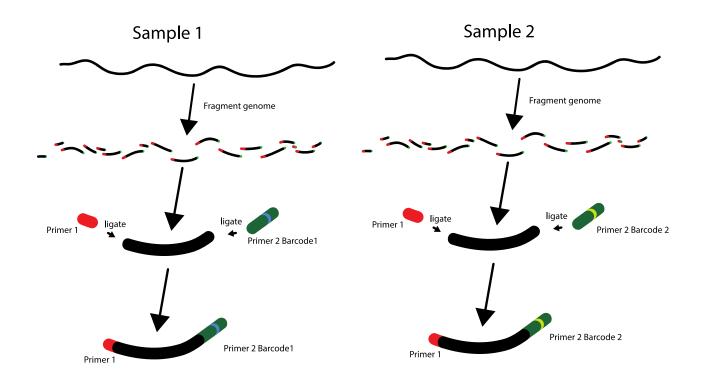
## Generating a library



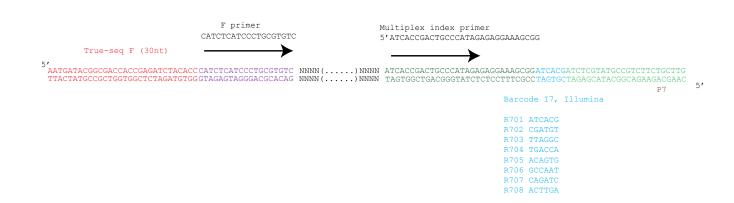
Most sequencing machines produce a fixed number of reads per flow cell.

What do you do if you are working with a small genome and only need a few million reads per sample?



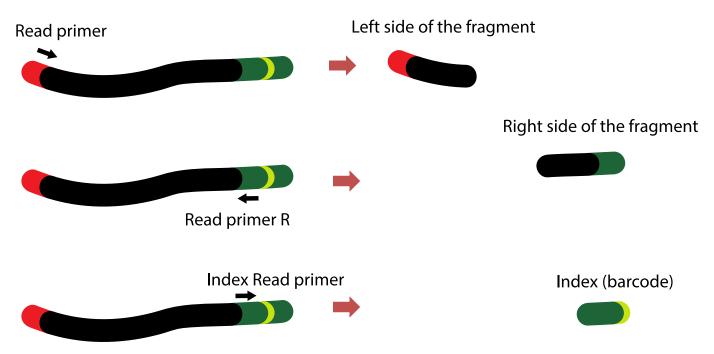






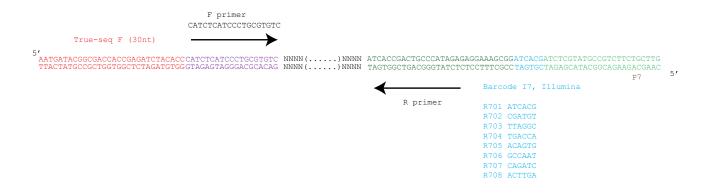
## Paired end sequencing Sequence both ends

# Paired end sequencing Sequence both ends of the molecule



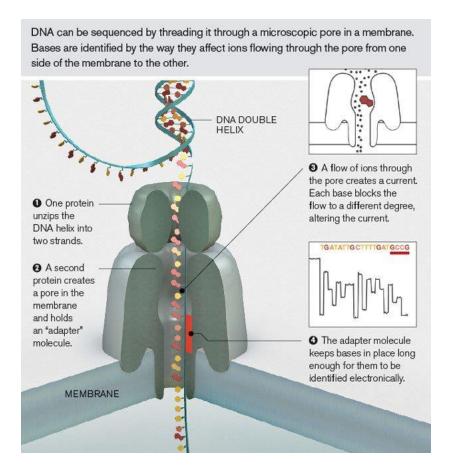
#### Paired end sequencing

#### Sequence both ends of the molecule!



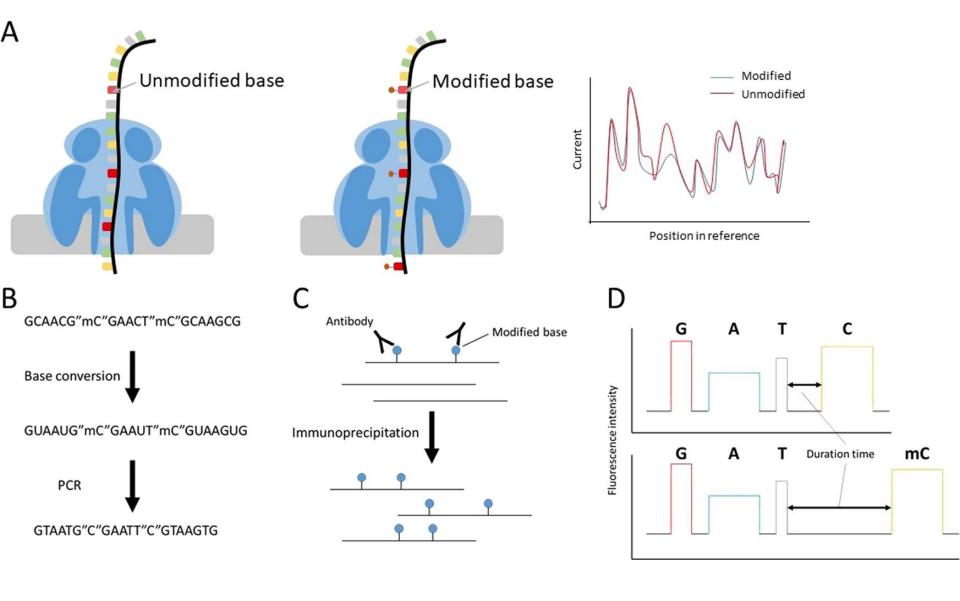
# Long-read sequencing

# Nanopore sequencing



Single molecule, Long reads: up to 3MB Low accuracy, low throughput

#### DNA methylation



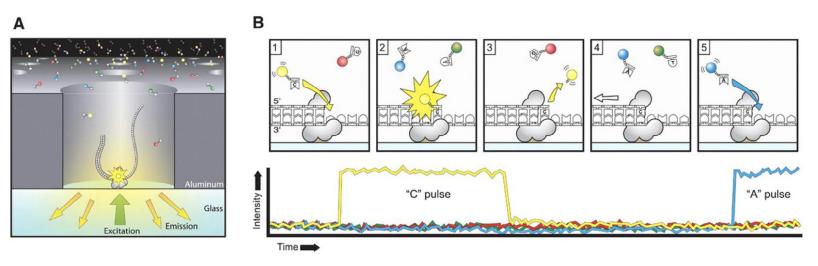
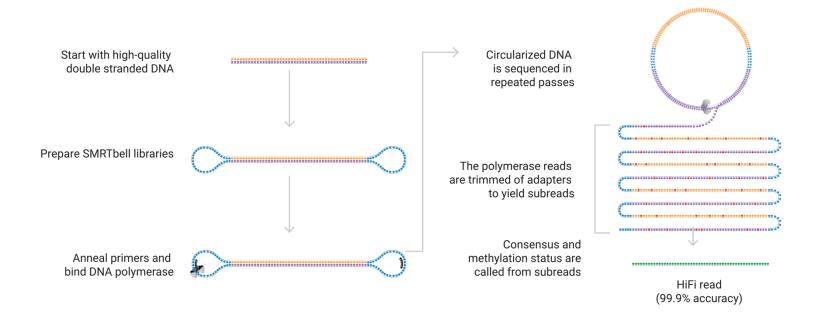


Figure 3 Sequencing via light pulses

A. A SMRTbell (gray) diffuses into a ZMW, and the adaptor binds to a polymerase immobilized at the bottom. **B.** Each of the four nucleotides is labeled with a different fluorescent dye (indicated in red, yellow, green, and blue, respectively for G, C, T, and A) so that they have distinct emission spectrums. As a nucleotide is held in the detection volume by the polymerase, a light pulse is produced that identifies the base. (1) A fluorescently-labeled nucleotide associates with the template in the active site of the polymerase. (2) The fluorescence output of the color corresponding to the incorporated base (yellow for base C as an example here) is elevated. (3) The dye-linker-pyrophosphate product is cleaved from the nucleotide and diffuses out of the ZMW, ending the fluorescence pulse. (4) The polymerase translocates to the next position. (5) The next nucleotide associates with the template in the active site of the polymerase, initiating the next fluorescence pulse, which corresponds to base A here. The figure is adapted from [4] with permission from The American Association for the Advancement of Science.

#### PAC Bio



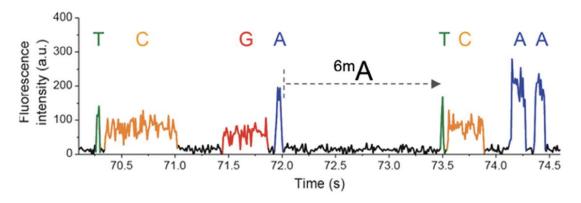
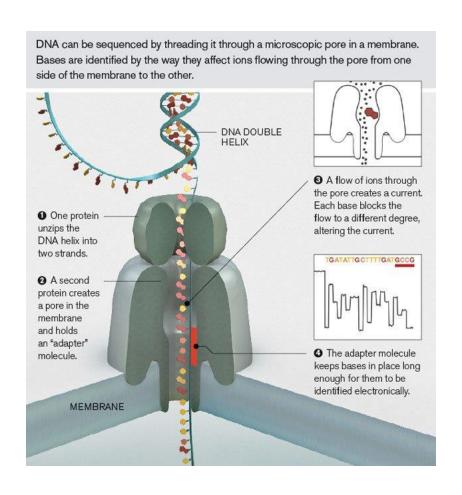


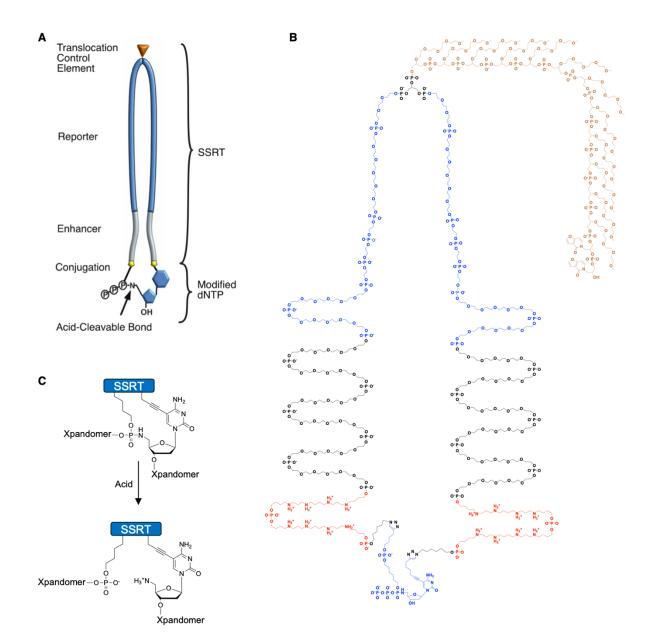
Figure 5 Detection of methylated bases using PacBio sequencing

PacBio sequencing can detect modified bases, including m<sup>6</sup>A (also known as <sup>6m</sup>A), by analyzing variation in the time between base incorporations in the read strand. The figure is adapted with permission from Pacific Biosciences [72]. a.u. stands for arbitrary unit.

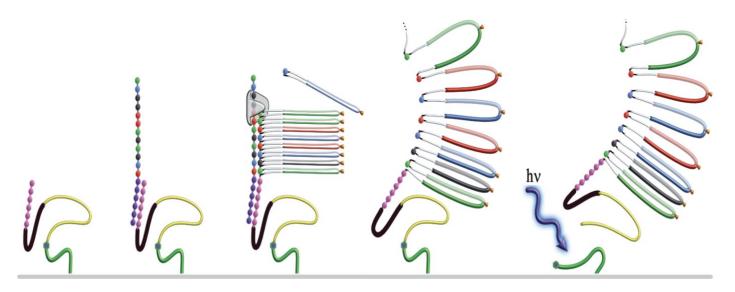
#### Nanopore sequencing



#### SBX Sequencing by Expandomers

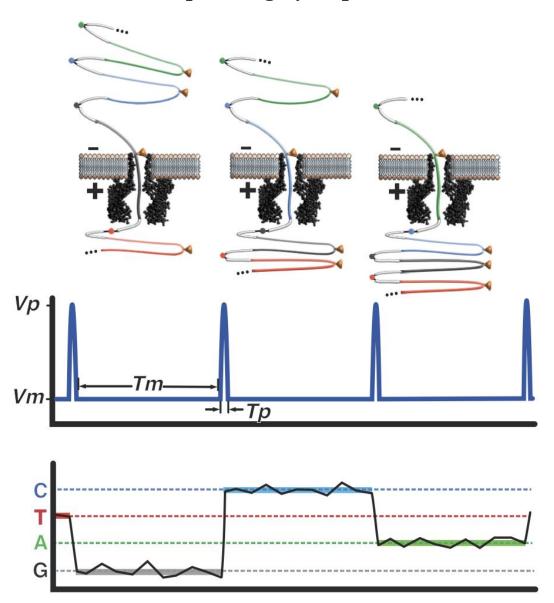


#### SBX Sequencing by Expandomers



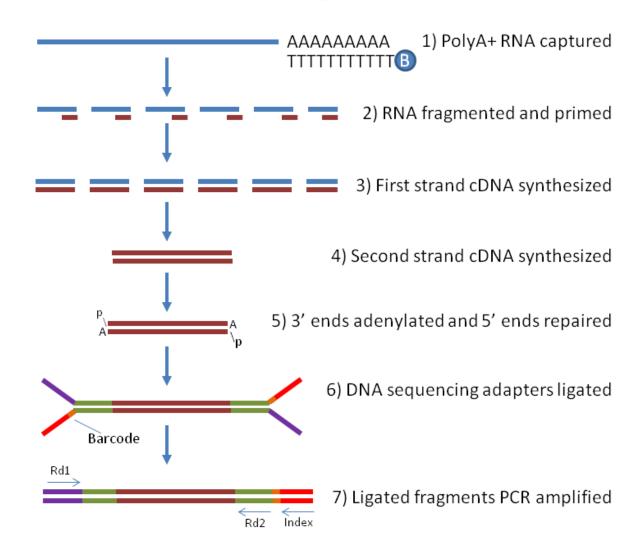
Flow Channel Surface

#### SBX Sequencing by Expandomers



# Some things to consider when analyzing genomics data.

# RNA-seq



PCR is used to amplify the library before sample submission

~20ng of library is usually sufficient for the genomics core facility to work with your sample.

20ng = 0.1 pMol of library (300bp product)

PCR is used to amplify the library before sample submission

~20ng of library is usually sufficient for the genomics core facility to work with your sample.

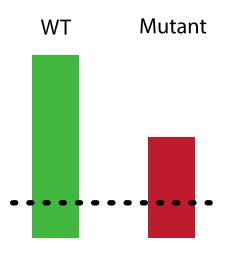
20ng = 0.1 pMol of library (300bp product)

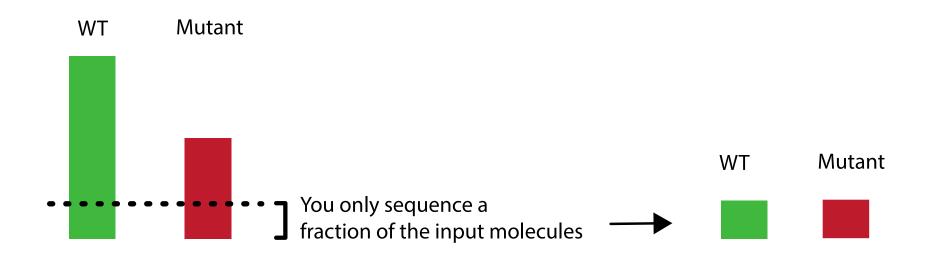
~ 60,000,000,000 molecules!

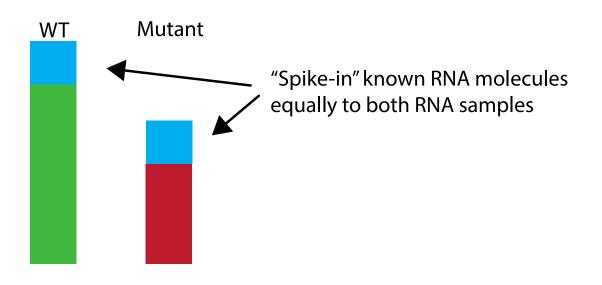
Illumina machines will now produce ~5,000,000,000 reads per flowcell

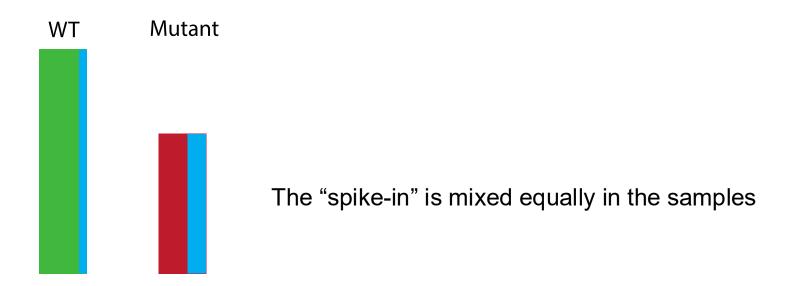
Imagine you conduct an experiment to compare mRNA in WT and a mutant of RNA polymerase II

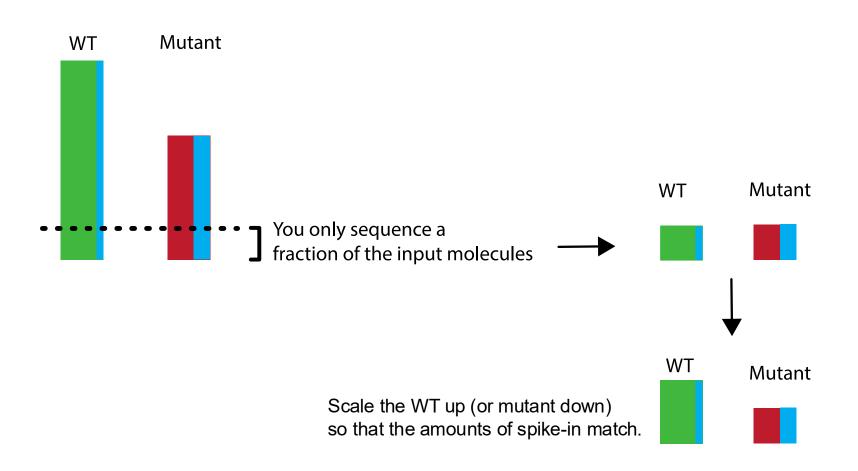
Through other experiments you are fairly convinced that the polymerase mutant reduces all mRNA by about 2-fold.







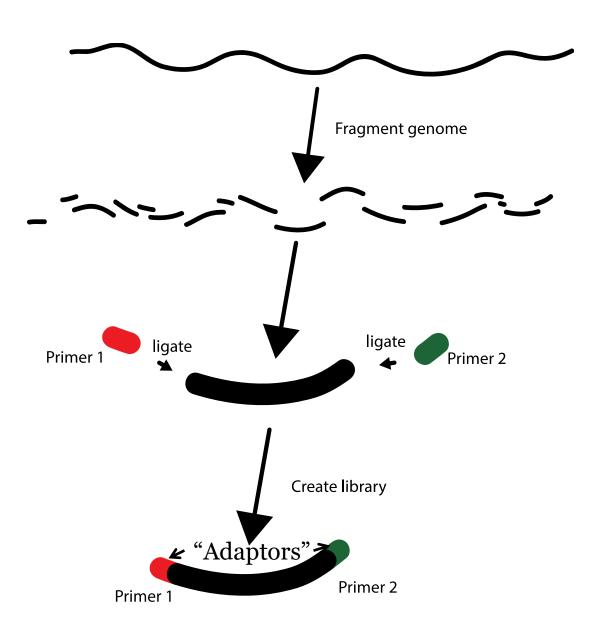




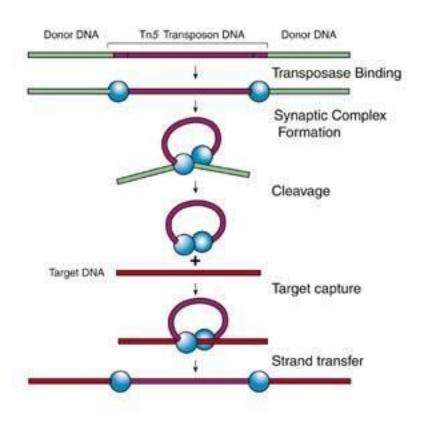
### ATAC seq

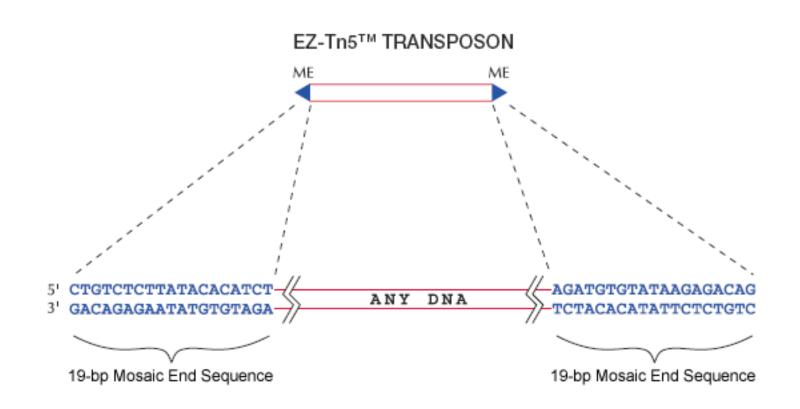
Assay for Transposase-Accessible Chromatin using sequencing

# Generating a library

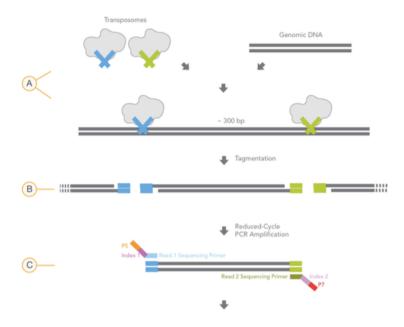


#### Tn5 Transposon





#### ATAC seq

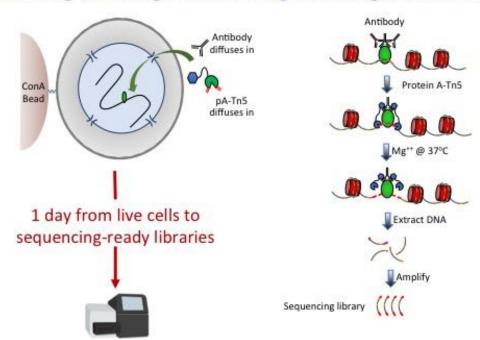


#### ATAC seq



#### Cut and Tag Cut and Run

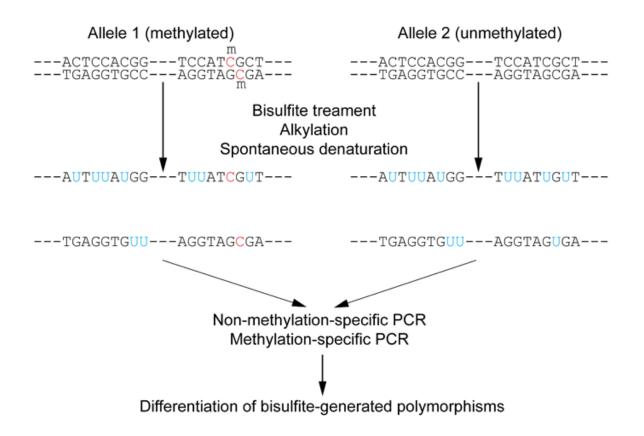
#### CUT&Tag (Cleavage Under Targets & Tagmentation)



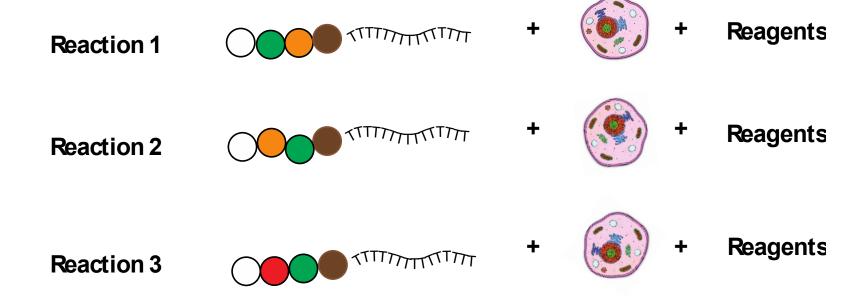
# Bisulphite sequencing

5-methylcytosine

# Bisulphite sequencing



#### Single cell RNA seq



Combinatorial indexing

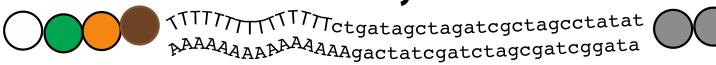
#### **RNA** annealing

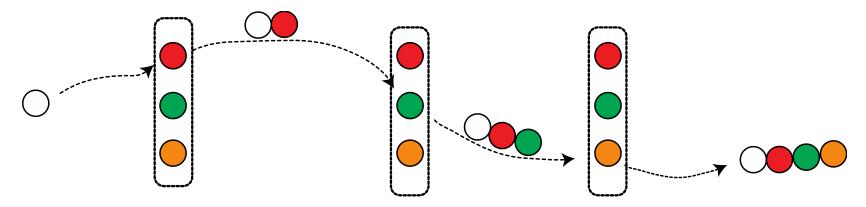


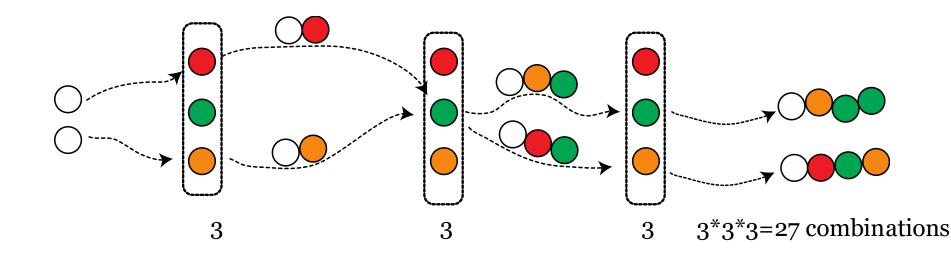
#### **Reverse transcription**

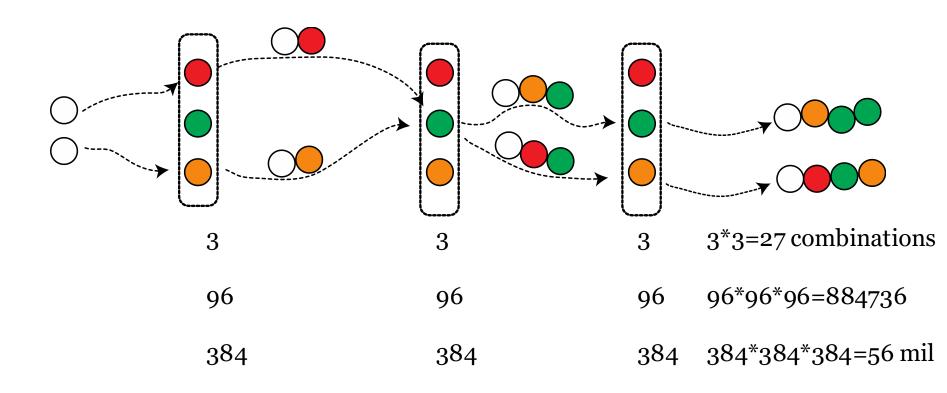


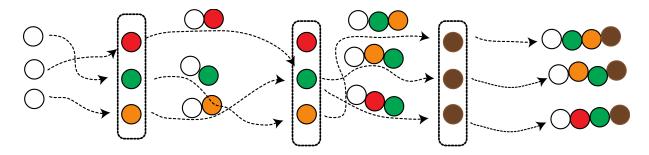
#### cDNA synthesis

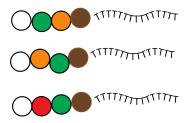




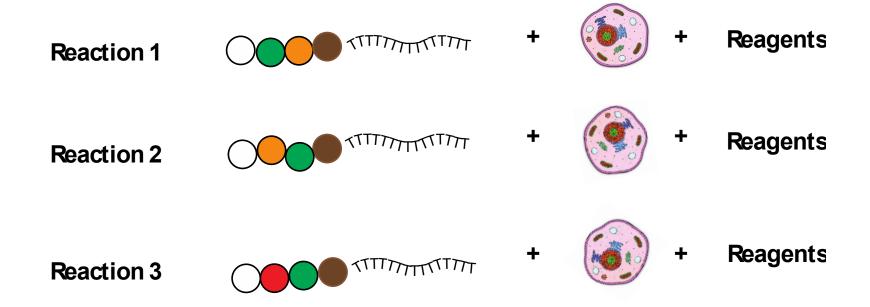


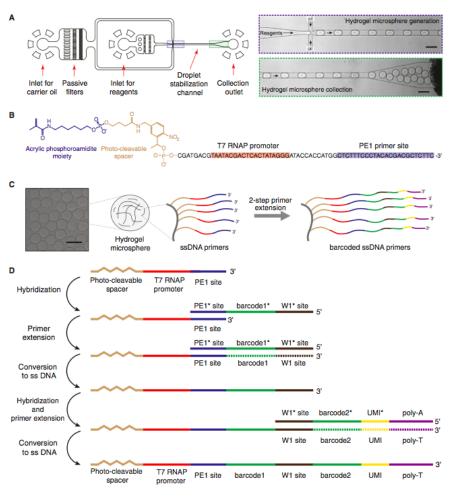


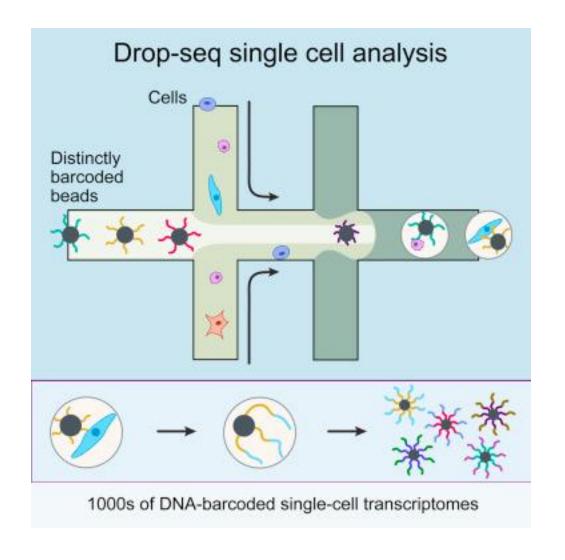




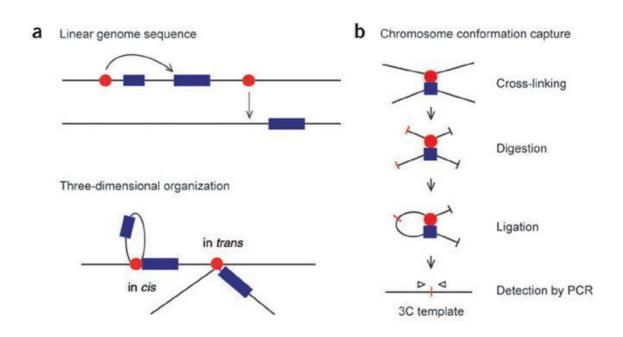
#### Combinatorial indexing

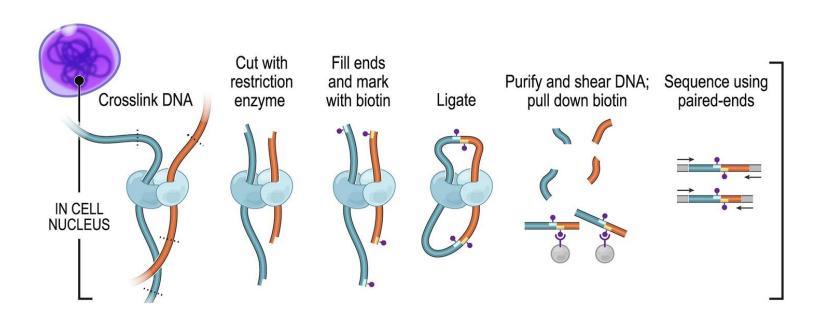


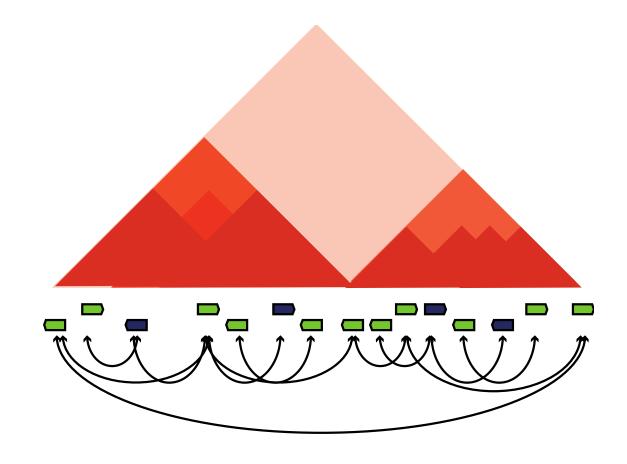


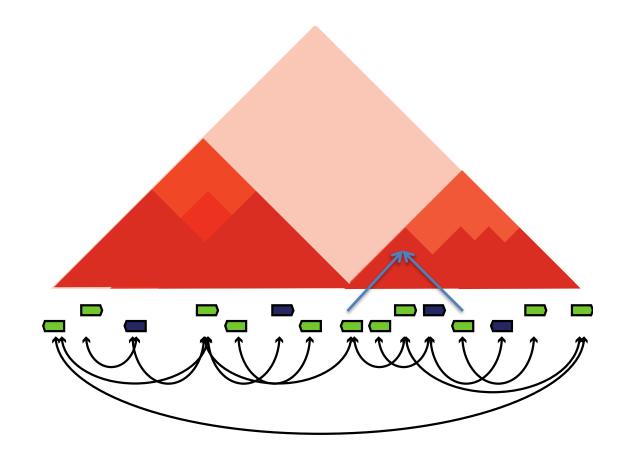


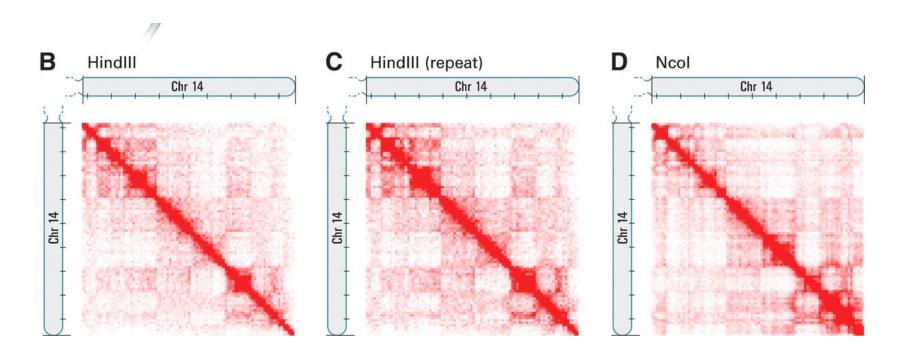
#### Chromosome conformation capture (3C)











# Ribosome-profiling

