

# Ancient sequencing technology – Sanger Vectors



DNA



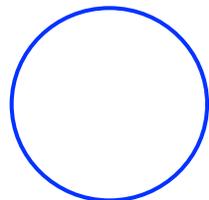
Shake



DNA fragments



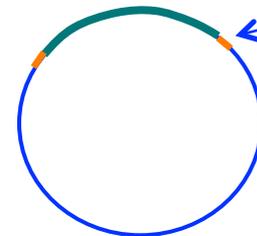
Vector  
Circular genome  
(bacterium, plasmid)



+



=



Known location

(restriction site)





# Fluorescent Sanger sequencing trace

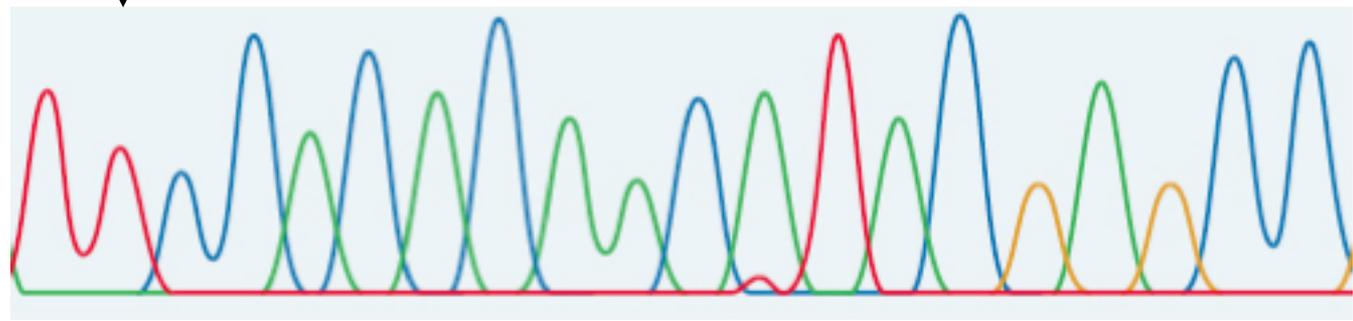
Lane signal



(Real fluorescent signals from a lane/capillary are much uglier than this).

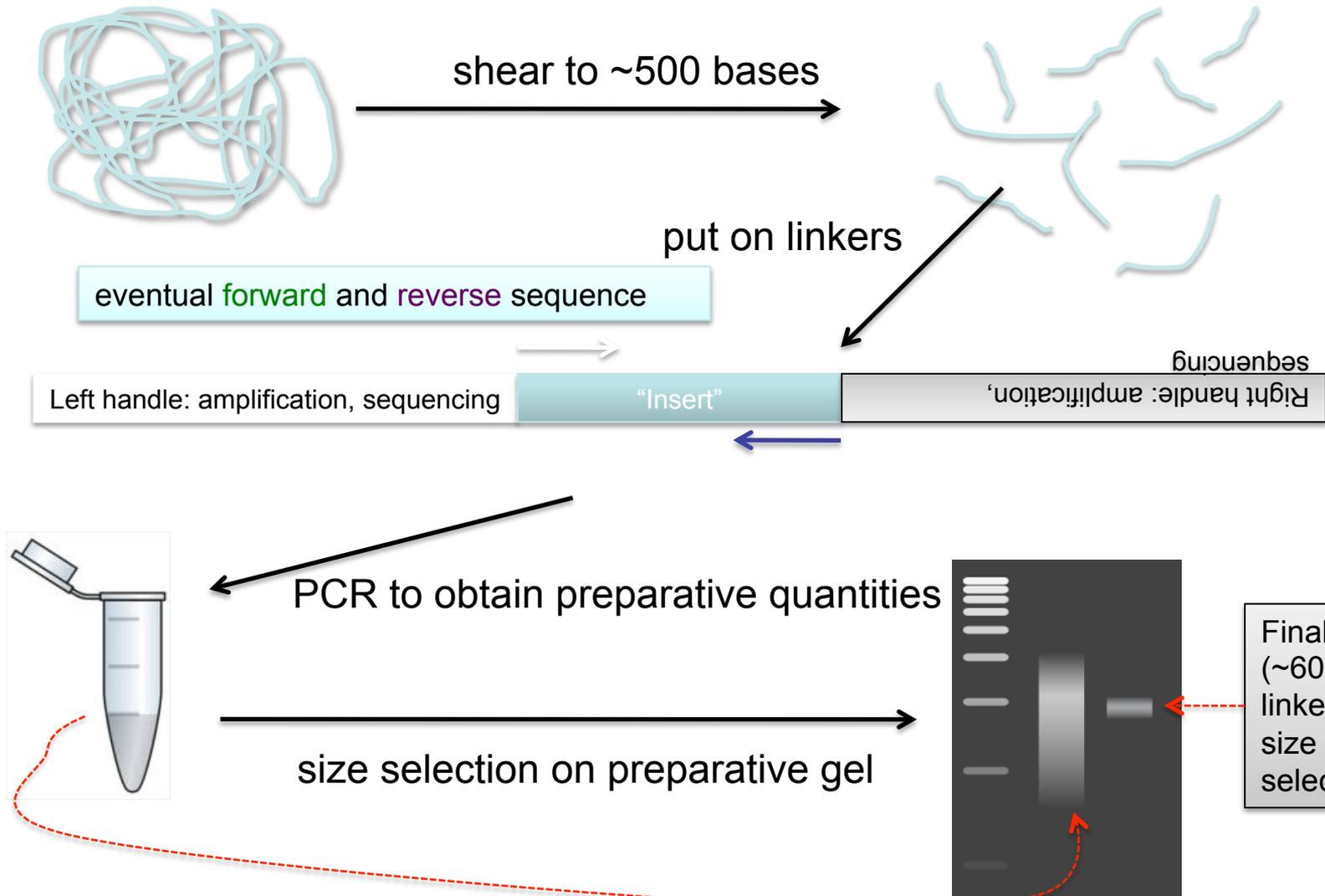
A bunch of magic to boost signal/noise, correct for dye-effects, mobility differences, etc, generates the 'final' trace (for each capillary of the run)

Trace





# Making a Library (present)



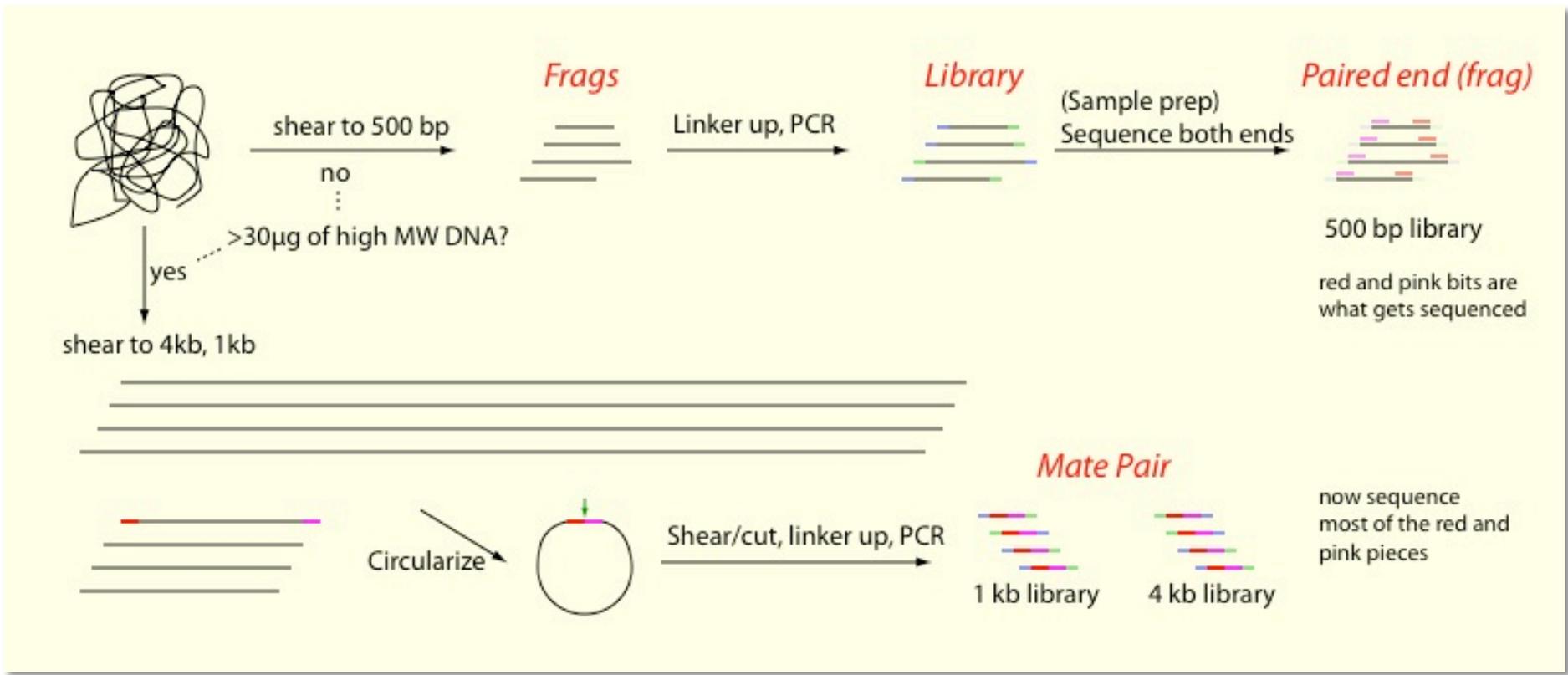
# Library



- Library is a massively complex mix of -initially- individual, unique fragments
- Library amplification mildly amplifies each fragment to retain the complexity of the mix while obtaining preparative amounts
  - (how many-fold do 10 cycles of PCR amplify the sample?)



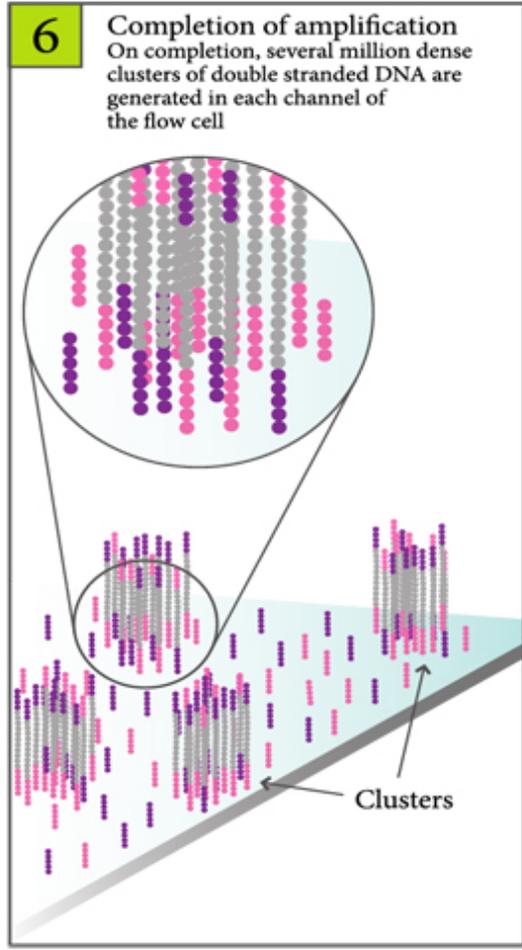
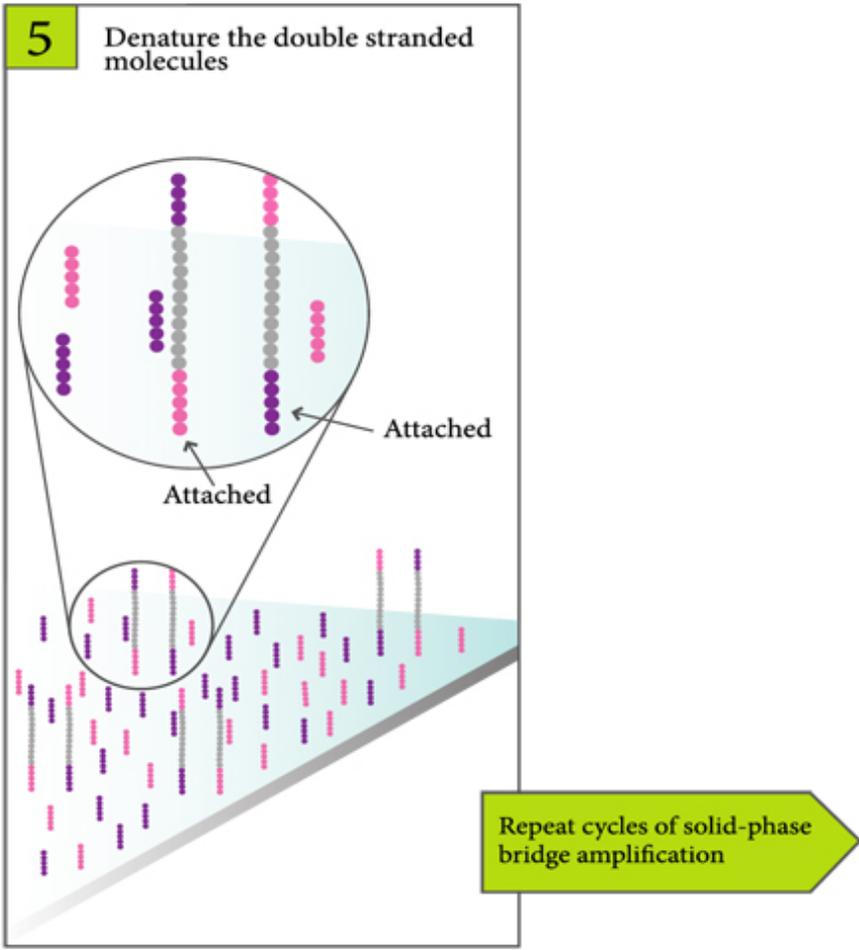
# Fragment vs Mate pair ('jumping')



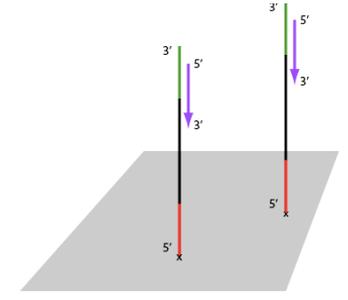
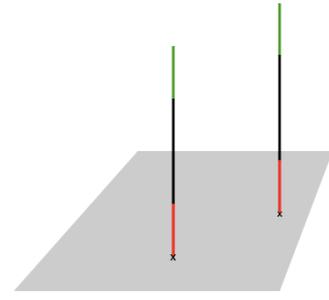
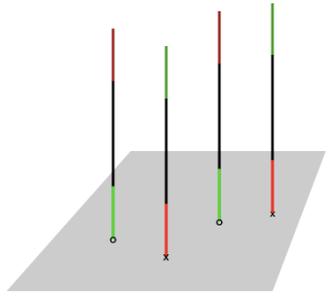
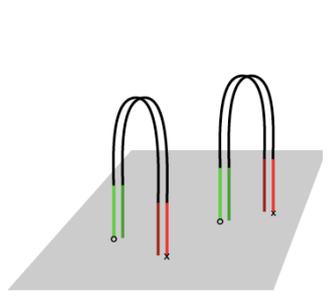
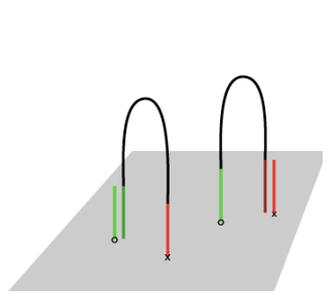
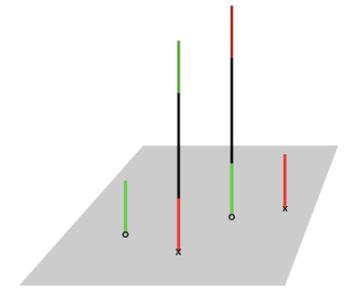
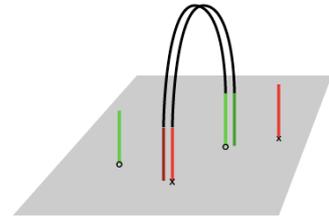
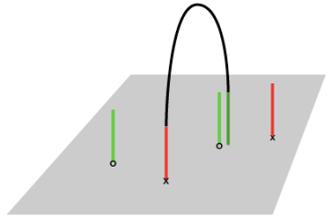
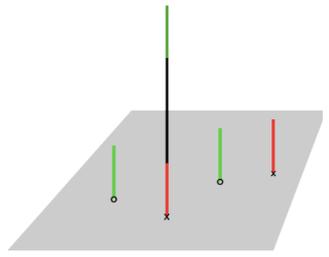
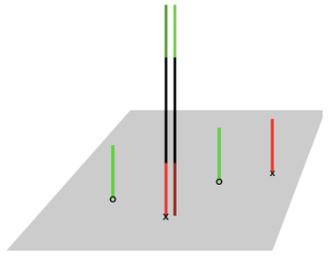
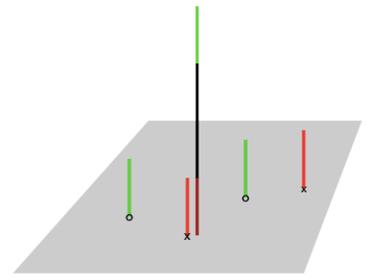
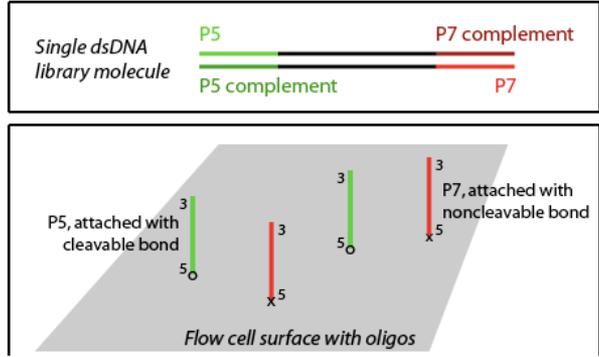
(Illumina has new kits/methods with which mate pair libraries can be built with less material)



# Illumina cluster concept

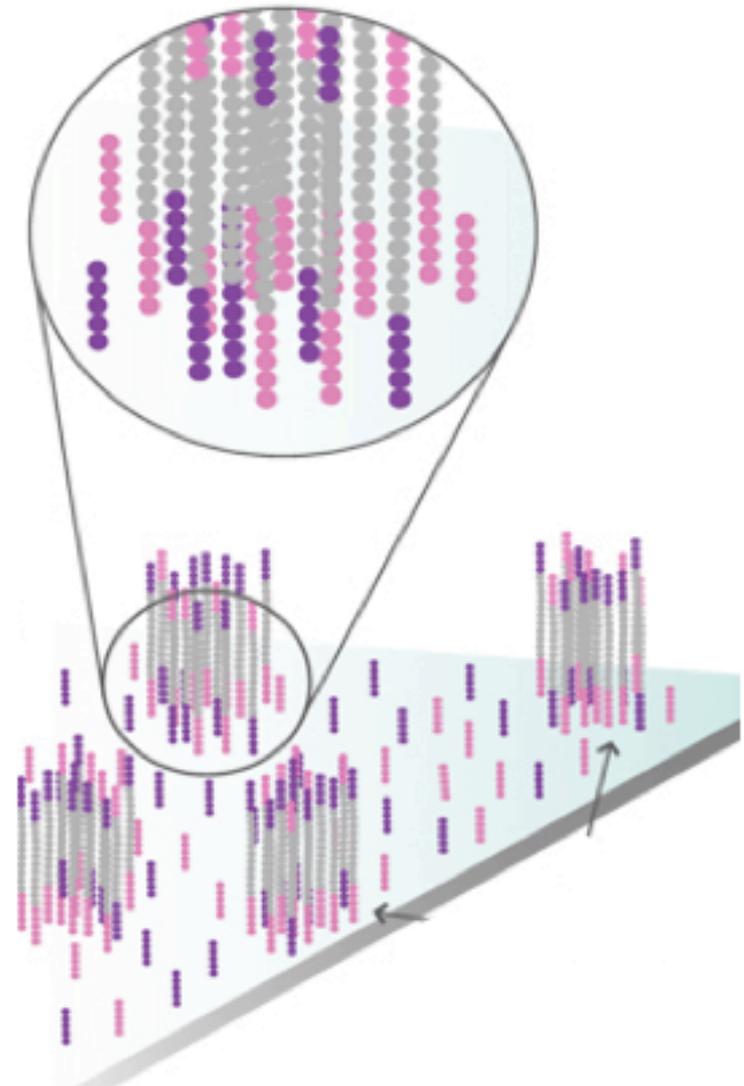
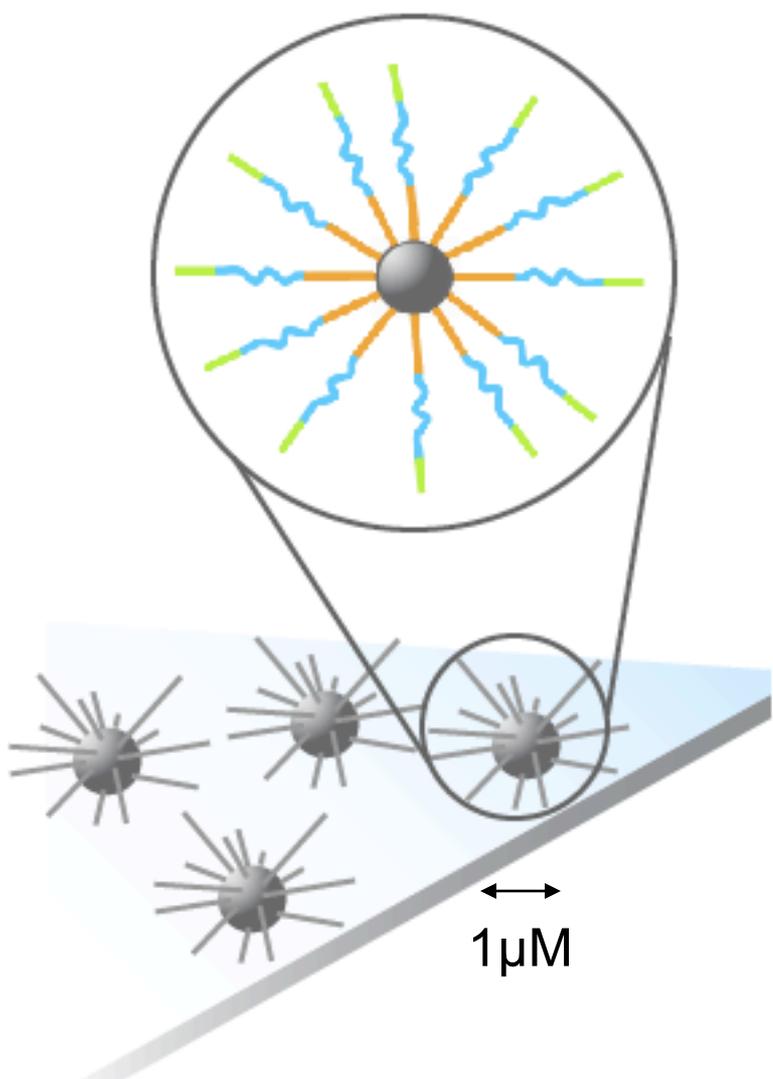


# Cluster generation ('bridge amplification')



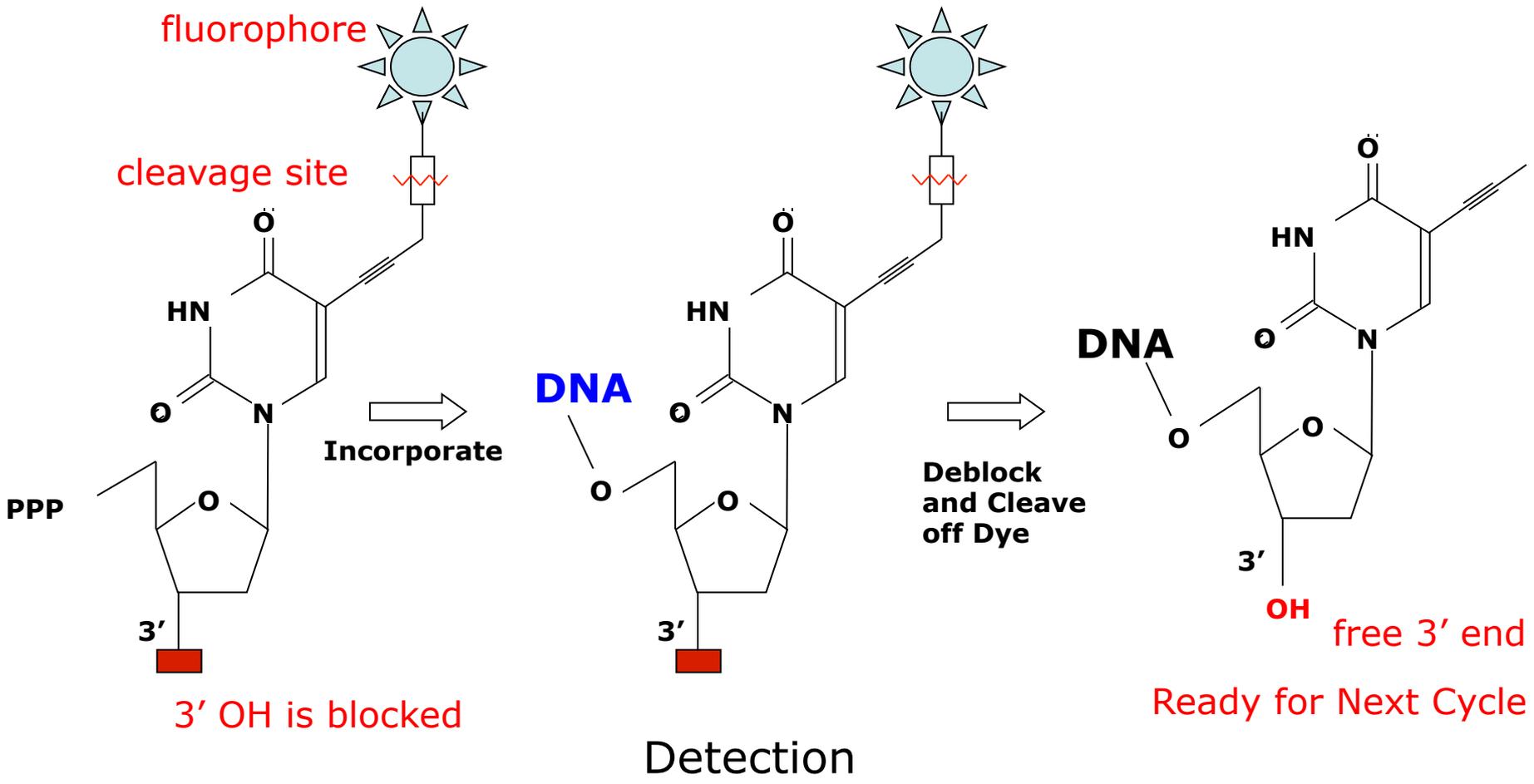


# Clonally Amplified Molecules on Flow Cell



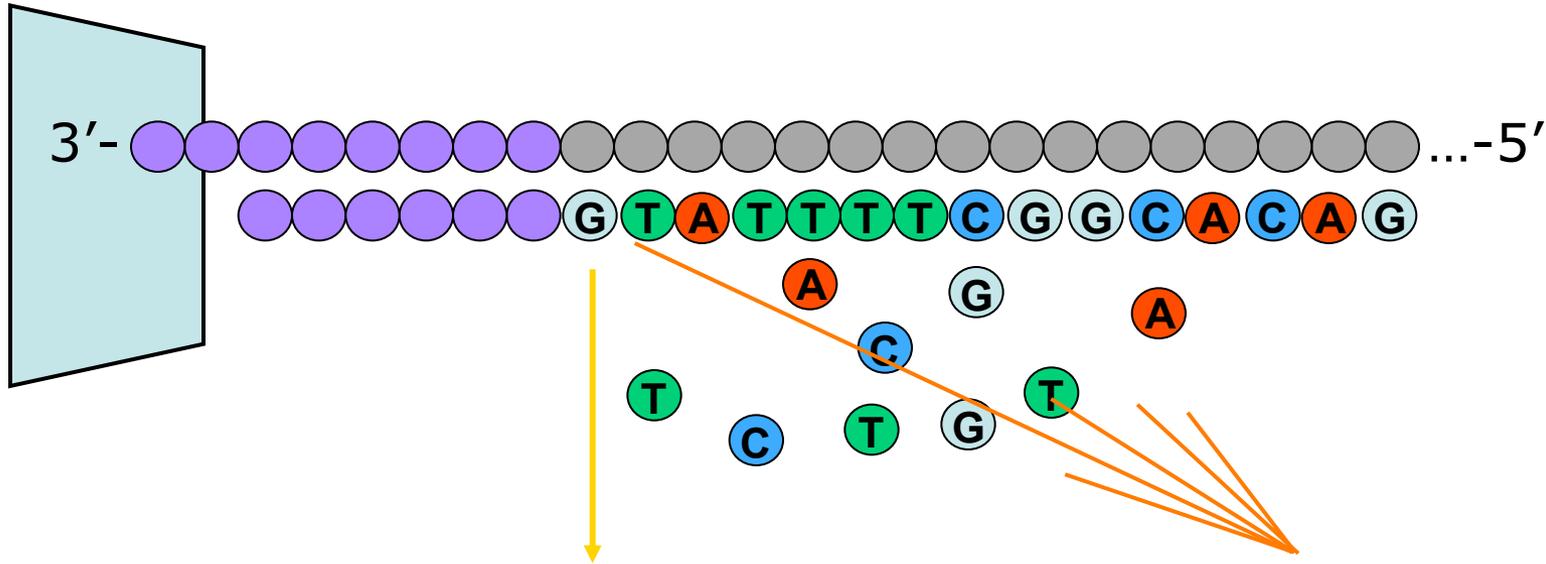


# Reversible Terminators





# Sequencing by Synthesis, One Base at a Time



Cycle 1:            Add sequencing reagents  
                      First base incorporated  
                      Remove unincorporated bases  
                      Detect signal

Cycle 2-n:            Add sequencing reagents and repeat



# Sequencing power for every scale.

Find the sequencing system that's right for your lab.

Compare key specifications across the whole portfolio of Illumina sequencing systems. Understand the differences between the MiniSeq, MiSeq, NextSeq, HiSeq, and HiSeq X Series.



	 MiniSeq System	 MiSeq Series	 NextSeq Series	 HiSeq Series	 HiSeq X Series*
<b>Key Methods</b>	Amplicon, targeted RNA, small RNA, and targeted gene panel sequencing.	Small genome, amplicon, and targeted gene panel sequencing.	Everyday exome, transcriptome, and targeted resequencing.	Production-scale genome, exome, transcriptome sequencing, and more.	Population- and production-scale whole-genome sequencing.
<b>Maximum Output</b>	7.5 Gb	15 Gb	120 Gb	1500 Gb	1800 Gb
<b>Maximum Reads per Run</b>	25 million	25 million <sup>†</sup>	400 million	5 billion	6 billion
<b>Maximum Read Length</b>	2 × 150 bp	2 × 300 bp	2 × 150 bp	2 × 150 bp	2 × 150 bp
<b>Run Time</b>	4–24 hours	4–55 hours	12–30 hours	<1–3.5 days (HiSeq 3000/HiSeq 4000) 7 hours–6 days (HiSeq 2500)	<3 days
<b>Benchtop Sequencer</b>	Yes	Yes	Yes	No	No
<b>System Versions</b>	<ul style="list-style-type: none"> <li>MiniSeq System for low-throughput targeted DNA and RNA sequencing</li> </ul>	<ul style="list-style-type: none"> <li>MiSeq System for targeted and small genome sequencing</li> <li>MiSeq FGx System for forensic genomics</li> <li>MiSeqDx System for molecular diagnostics</li> </ul>	<ul style="list-style-type: none"> <li>NextSeq 500 System for everyday genomics</li> <li>NextSeq 550 System for both sequencing and cytogenomic arrays</li> </ul>	<ul style="list-style-type: none"> <li>HiSeq 3000/HiSeq 4000 Systems for production-scale genomics</li> <li>HiSeq 2500 Systems for large-scale genomics</li> </ul>	<ul style="list-style-type: none"> <li>HiSeq X Five System for production-scale whole-genome sequencing</li> <li>HiSeq X Ten System for population-scale whole-genome sequencing</li> </ul>



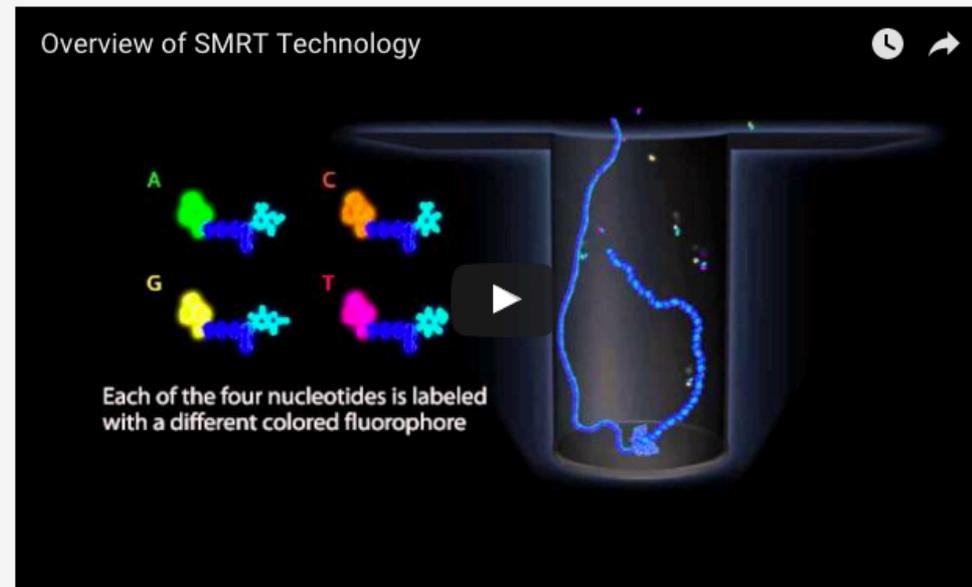
# Pacific Biosciences SMRT technology

## The SMRT Sequencing advantage

SMRT Sequencing is ideal for a variety of research applications and offers many benefits, including:

- [Longest average read lengths](#)
- [Highest consensus accuracy](#)
- [Uniform coverage](#)
- [Simultaneous epigenetic characterization](#)
- [Single-molecule resolution](#)

## An overview of SMRT Sequencing



# Oxford Nanopore



## MinION

Portable, real-time biological analyses



MinION is a portable device for molecular analyses that is driven by nanopore technology. It is adaptable for the analysis of DNA, RNA, proteins or small molecules with a straightforward workflow. The MinION product specification is available here.

[More about sequencing with MinION](#) ▾

[Explore all publications](#) >

[Start using MinION](#) >

### Simple workflows



Sample preparation — MinION — PromethION — Analyses

- 

Simple sample preparation  
(Coming soon: automated sample preparation from Voltrax)

[Learn about Voltrax](#) >
- 

Pocket-sized MinION for analysis anywhere

[Learn about MinION](#) >
- 

Desktop PromethION for high throughput analysis

[Learn about PromethION](#) >
- 

Real time analysis solutions from Metrichor

[Learn about Metrichor](#) >



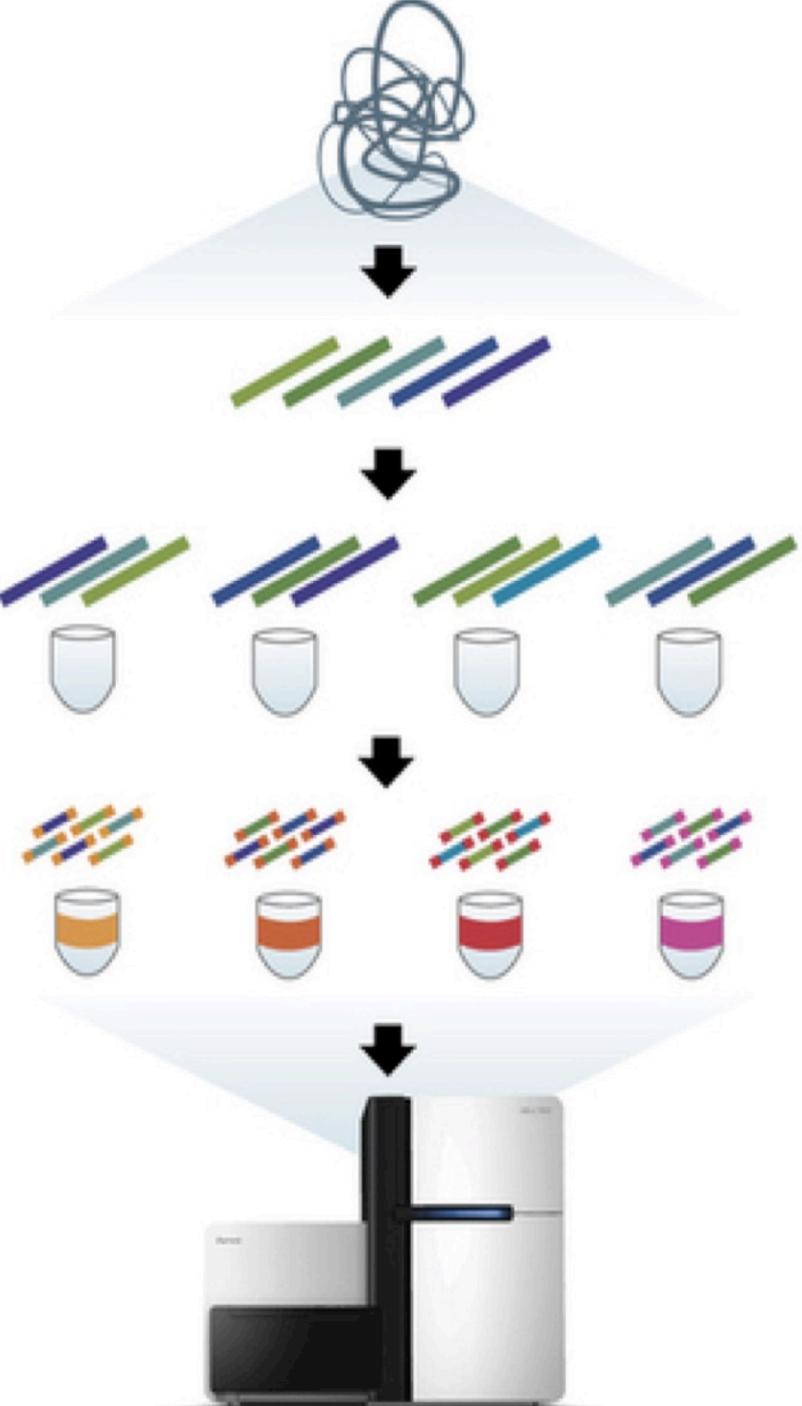
# Moleculo Overview

1. Sample DNA is sheared into fragments of about 10 kbp

2. Fragments are diluted and placed into 384 wells

3. Fragments are amplified through long-range PCR, cut into short fragments and barcoded

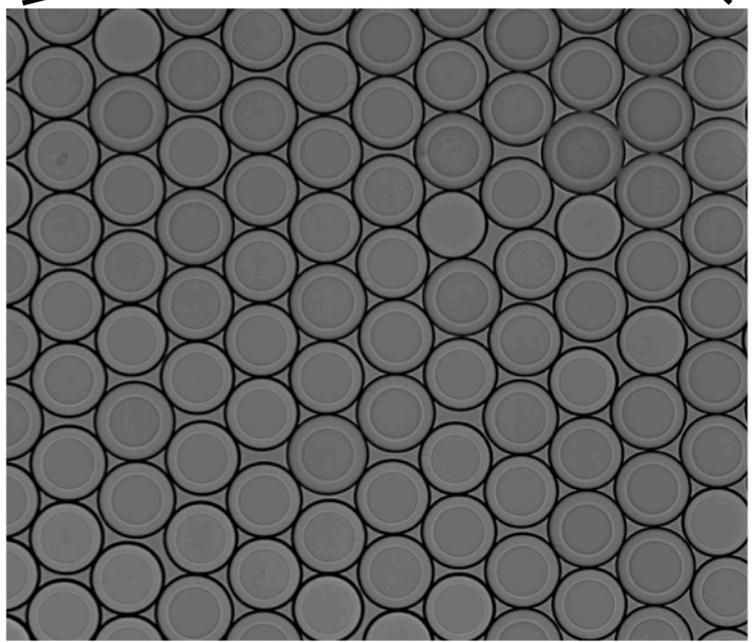
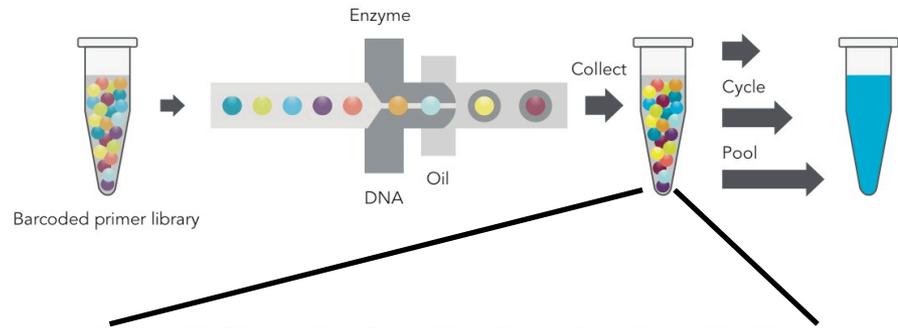
4. Short fragments are pooled together and sequenced





# 10x System

## Massively Parallel Partitioning



X 200,000+

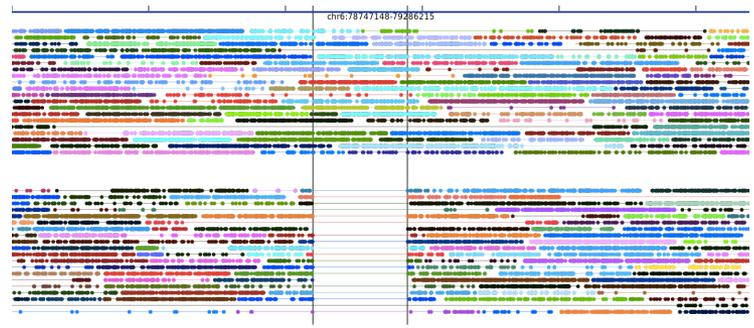
## 10X Instrument & Reagents



## Read Clouds ("linked reads")

Hap1

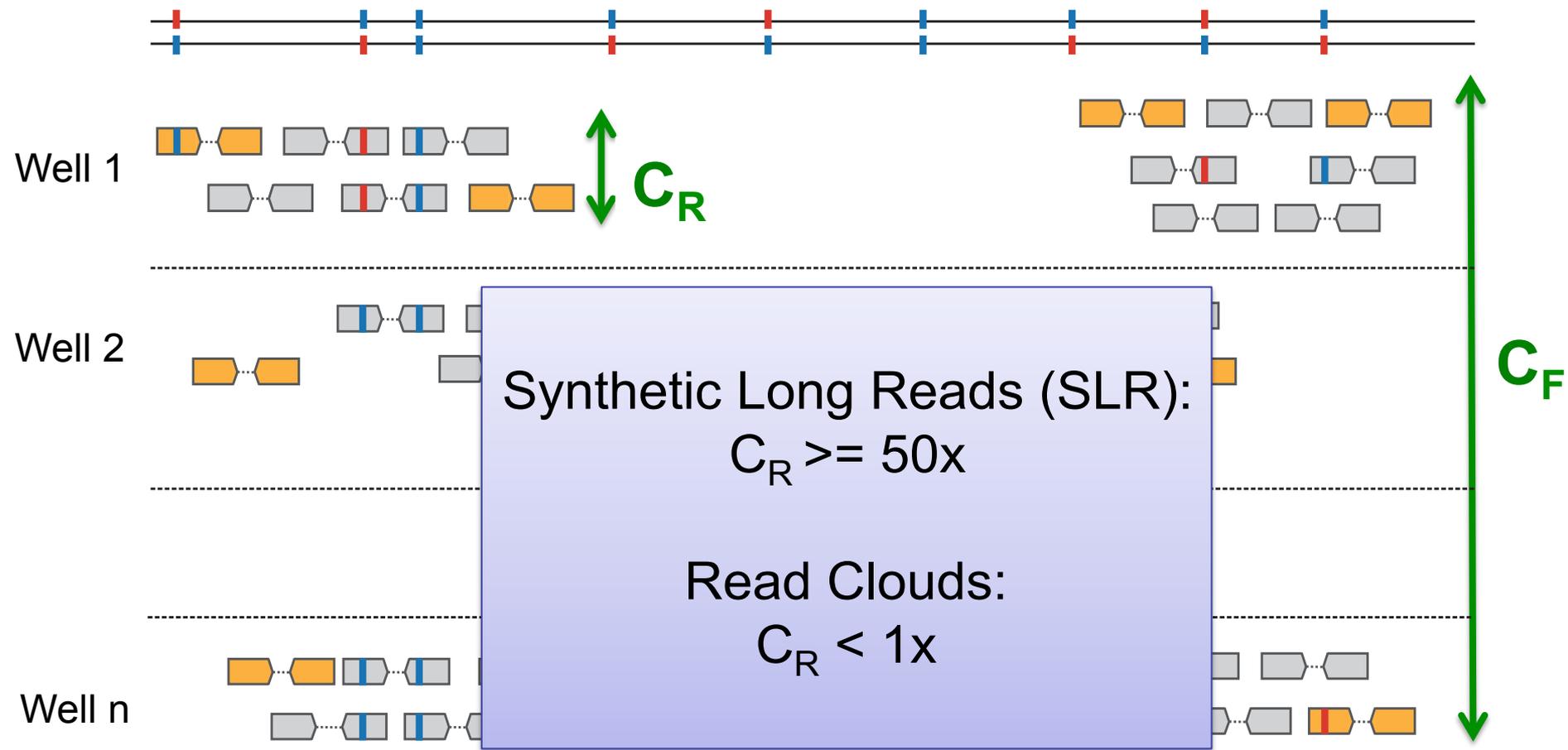
Hap2



Phased 60Kb deletion



# Read Clouds

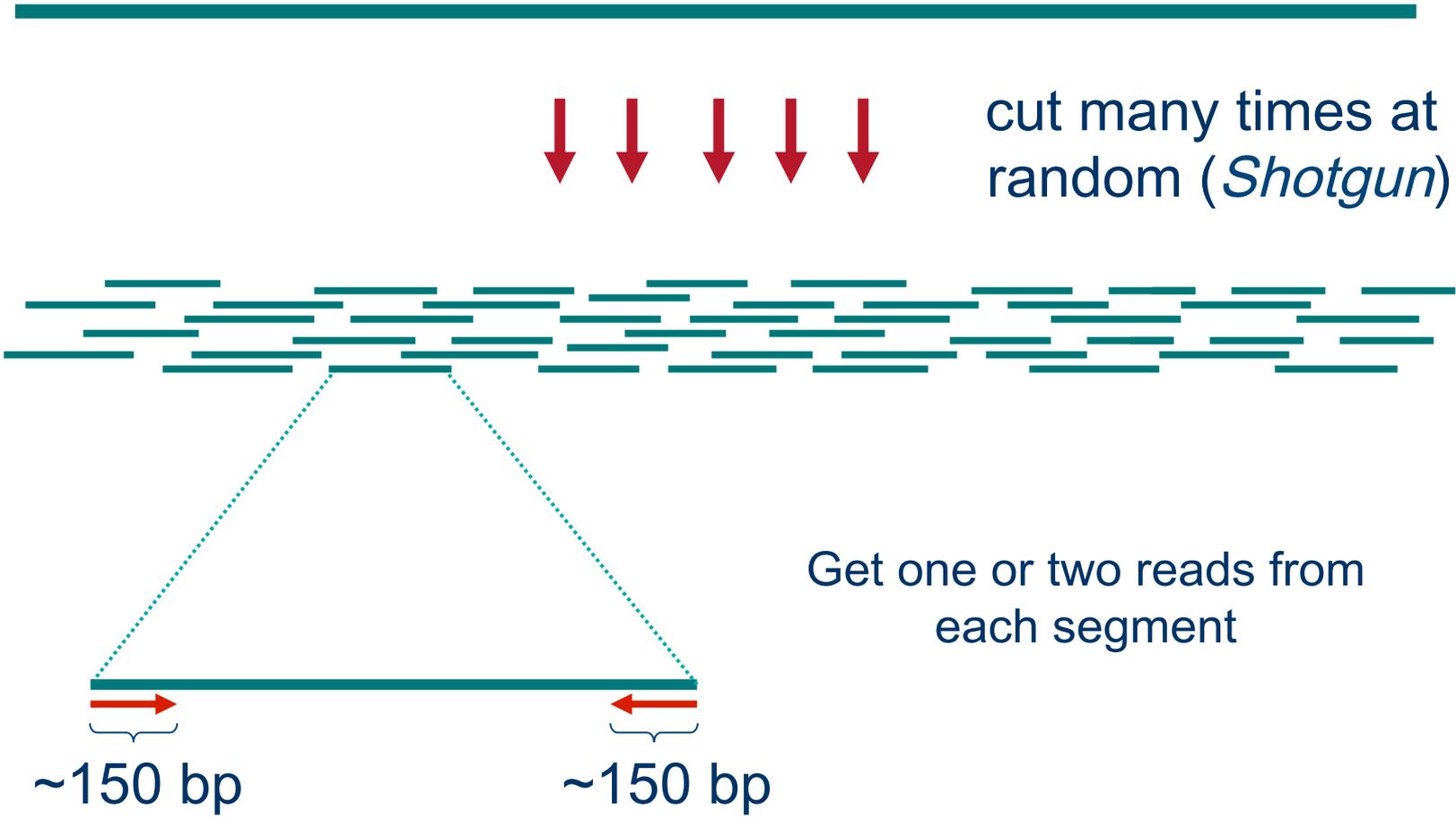


Coverage =  $C_F C_R$



# Shotgun Sequencing

genomic segment





# Two main assembly problems

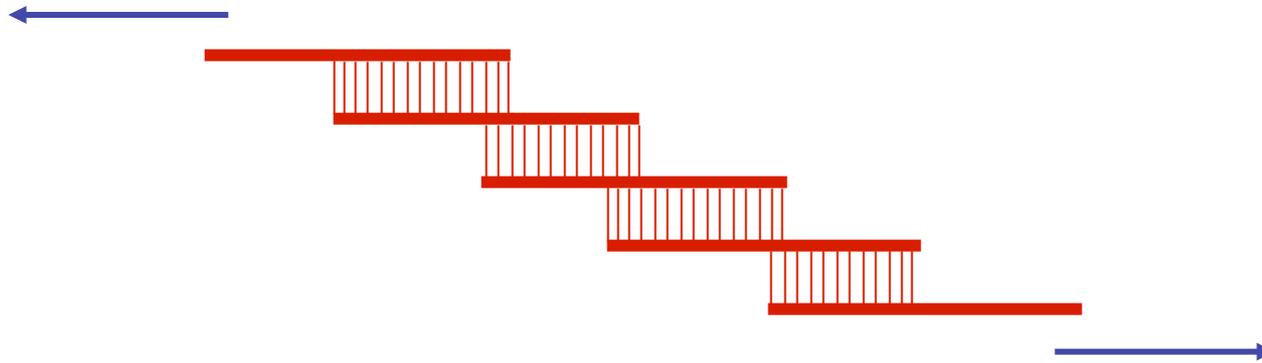
- De Novo Assembly



- Resequencing



# Reconstructing the Sequence (De Novo Assembly)



Cover region with high redundancy

Overlap & extend reads to reconstruct the original genomic region



# Definition of Coverage

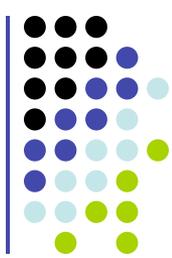


Length of genomic segment: **G**  
Number of reads: **N**  
Length of each read: **L**

**Definition:** Coverage  **$C = N L / G$**

How much coverage is enough?

**Lander-Waterman model:** **Prob[ not covered bp ] =  $e^{-C}$**   
Assuming uniform distribution of reads, **C=10** results in 1 gapped region /1,000,000 nucleotides



# Repeats

Bacterial genomes: 5%  
Mammals: 50%

## Repeat types:

- **Low-Complexity DNA** (e.g. ATATATATACATA...)
- **Microsatellite repeats**  $(a_1 \dots a_k)^N$  where  $k \sim 3-6$   
(e.g. CAGCAGTAGCAGCACCAG)
- **Transposons**
  - **SINE** (Short Interspersed Nuclear Elements)  
e.g., ALU: ~300-long,  $10^6$  copies
  - **LINE** (Long Interspersed Nuclear Elements)  
~4000-long, 200,000 copies
  - **LTR retroposons** (Long Terminal Repeats (~700 bp) at each end)  
cousins of HIV
- **Gene Families** genes duplicate & then diverge (paralogs)
- **Recent duplications** >100,000-long, very similar copies



# Sequencing and Fragment Assembly



$3 \times 10^9$  nucleotides



50% of human DNA is composed



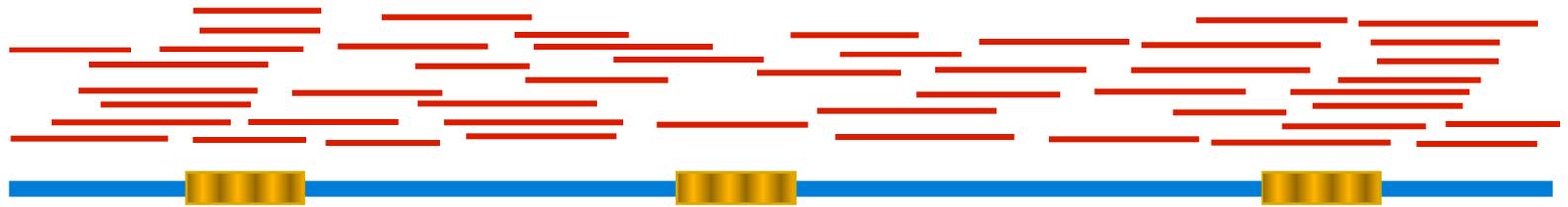
Error!  
Glued together two distant regions



# What can we do about repeats?

Two main approaches:

- Cluster the reads



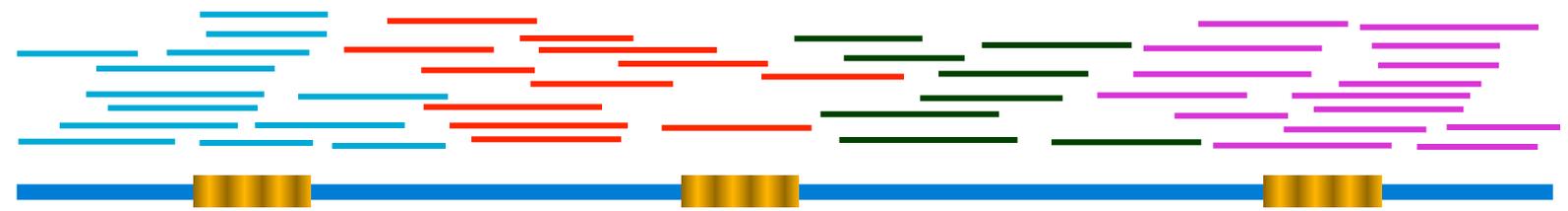
- Link the reads



# What can we do about repeats?

Two main approaches:

- Cluster the reads



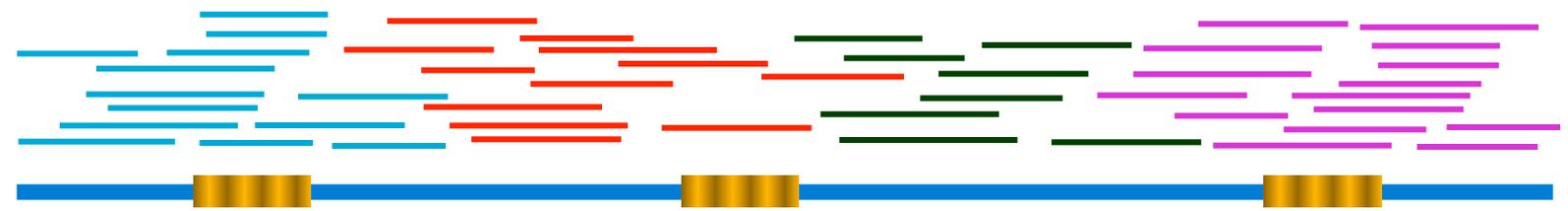
- Link the reads



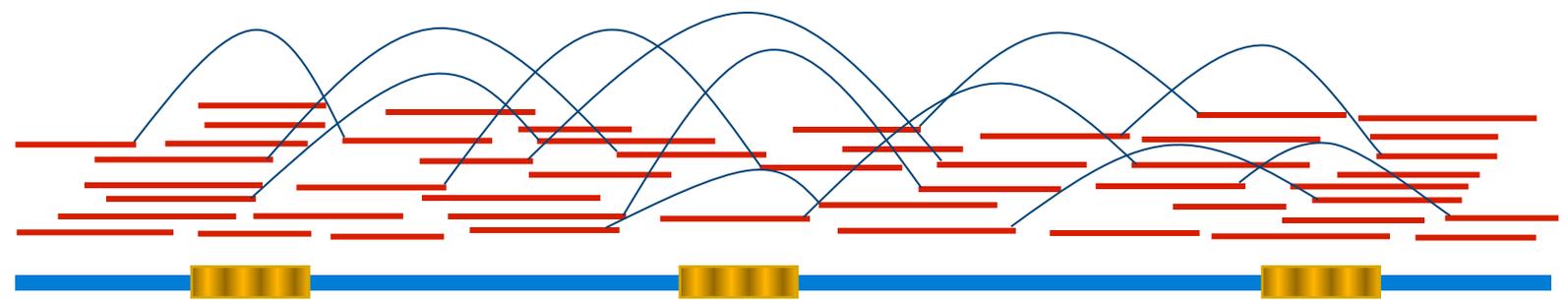
# What can we do about repeats?

Two main approaches:

- Cluster the reads



- Link the reads

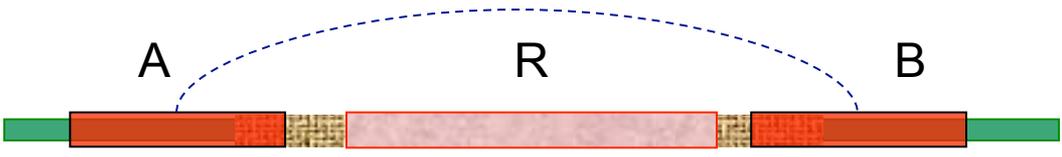




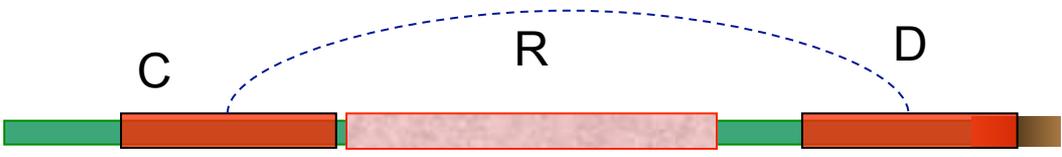
# Sequencing and Fragment Assembly



$3 \times 10^9$  nucleotides



ARB, CRD



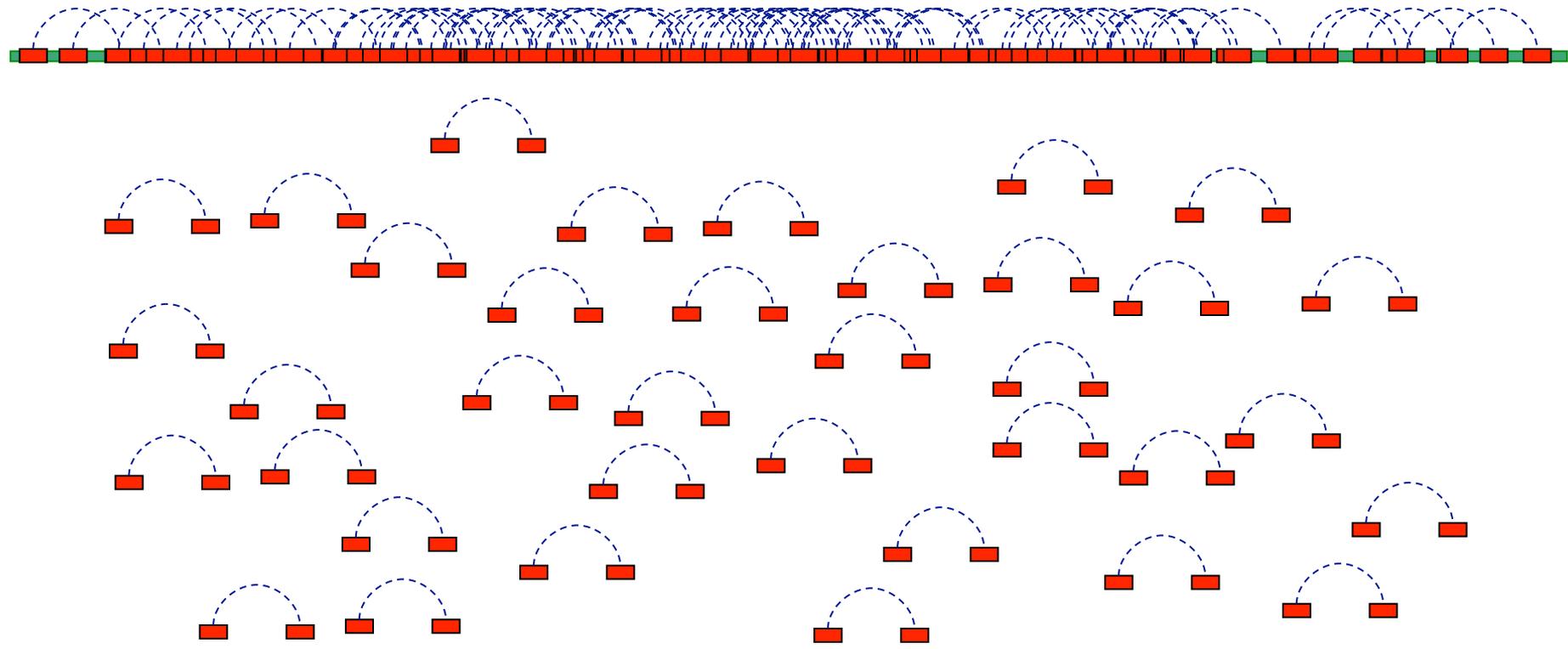
or  
~~ARD, CRB ?~~



# Sequencing and Fragment Assembly



$3 \times 10^9$  nucleotides





# Fragment Assembly (in whole-genome shotgun sequencing)





# Fragment Assembly



I THINK I FOUND A CORNER PIECE.

**Given N reads...  
Where N ~ 30 million...  
We need to use a linear-time algorithm**

SHEPHERD the Ledger

# Steps to Assemble a Genome



## Some Terminology

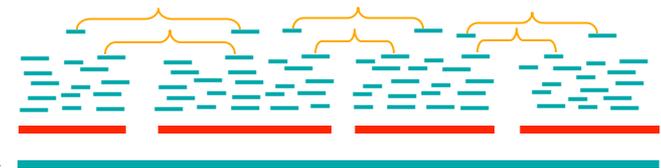
**read** a 500-900 long word that comes out of sequencer

**mate pair** a pair of reads from two ends of the same insert fragment

**contig** a contiguous sequence formed by several overlapping reads with no gaps

**supercontig (scaffold)** an ordered and oriented set of contigs, usually by mate pairs

**consensus sequence** sequence derived from the multiple alignment of reads in a contig



..ACGATTACAATAGGTT..



# 1. Find Overlapping Reads

aaactgcagtacggatct  
aaactgcag  
  aactgcagt  
...  
          gtacggatct  
          tacggatct  
gggcccaactgcagtac  
gggcccaa  
  ggcccaaac  
...  
          actgcagta  
          ctgcagtac  
gtacggatctactacaca  
gtacggatc  
  tacggatct  
...  
          ctactacac  
          tactacaca

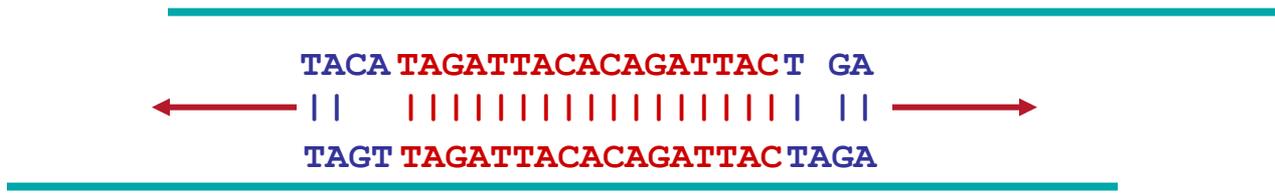
(read, pos., word, orient.)  
aaactgcag  
aactgcagt  
actgcagta  
...  
gtacggatc  
tacggatct  
gggcccaa  
ggcccaaac  
gcccaaact  
...  
actgcagta  
ctgcagtac  
gtacggatc  
tacggatct  
acggatcta  
...  
ctactacac  
tactacaca

(word, read, orient., pos.)  
aaactgcag  
aactgcagt  
**acggatcta**  
actgcagta  
actgcagta  
cccaaactg  
**cggatctac**  
**ctactacac**  
ctgcagtac  
ctgcagtac  
gcccaaact  
ggcccaaac  
gggcccaa  
gtacggatc  
gtacggatc  
tacggatct  
tacggatct  
tactacaca



# 1. Find Overlapping Reads

- Find pairs of reads sharing a k-mer,  $k \sim 24$
- Extend to full alignment – throw away if not  $>98\%$  similar

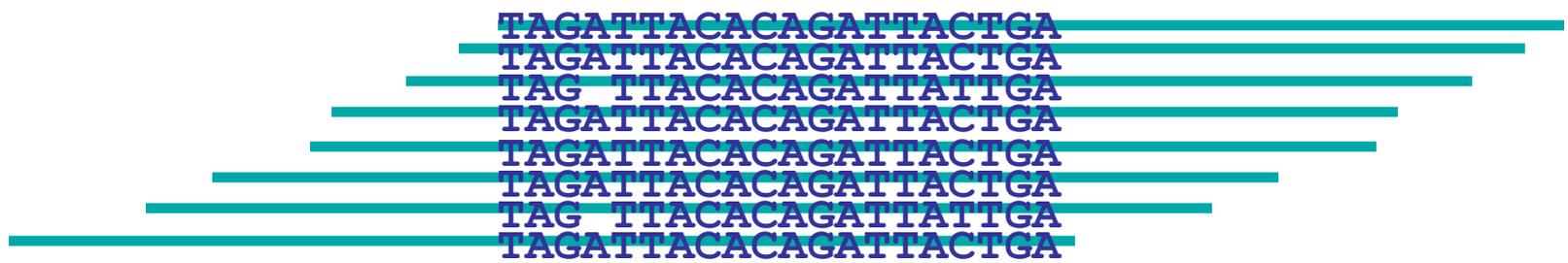


- Caveat: repeats
  - A k-mer that occurs  $N$  times, causes  $O(N^2)$  read/read comparisons
  - ALU k-mers could cause up to  $1,000,000^2$  comparisons
- Solution:
  - Discard all k-mers that occur “too often”
    - Set cutoff to balance sensitivity/speed tradeoff, according to genome at hand and computing resources available



# 1. Find Overlapping Reads

Create local multiple alignments from the overlapping reads





# 1. Find Overlapping Reads

- Correct errors using multiple alignment

```
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
```

insert A

replace T with C

```
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
```

correlated errors—  
probably caused by repeats  
⇒ disentangle overlaps

```
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
TAGATTACACAGATTACTGA
```

In practice, error correction removes up to 98% of the errors

```
TAG-TTACACAGATTACTGA
TAG-TTACACAGATTACTGA
```

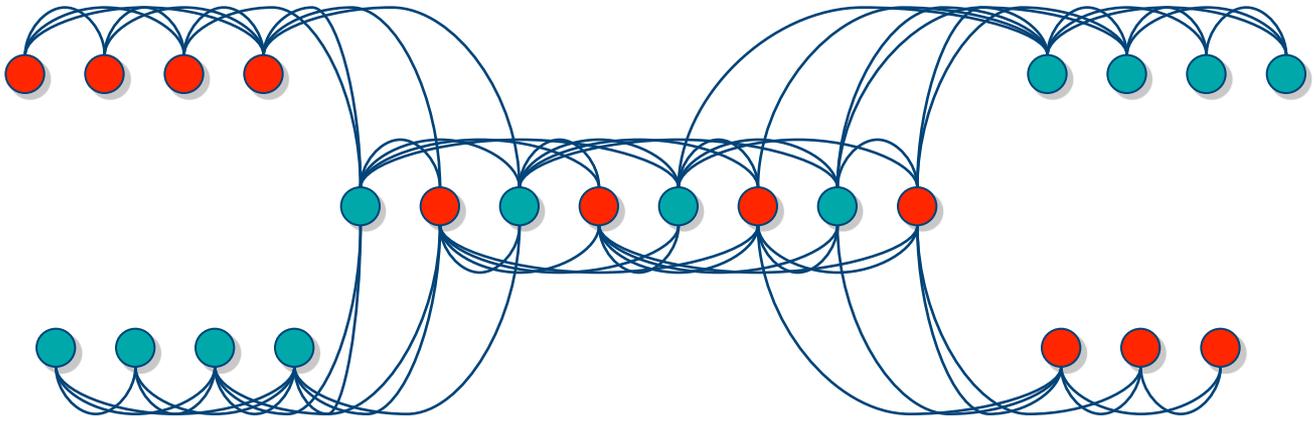


# 2. Merge Reads into Contigs

- Overlap graph:
  - Nodes: reads  $r_1 \dots r_n$
  - Edges: overlaps  $(r_i, r_j, \text{shift}, \text{orientation}, \text{score})$



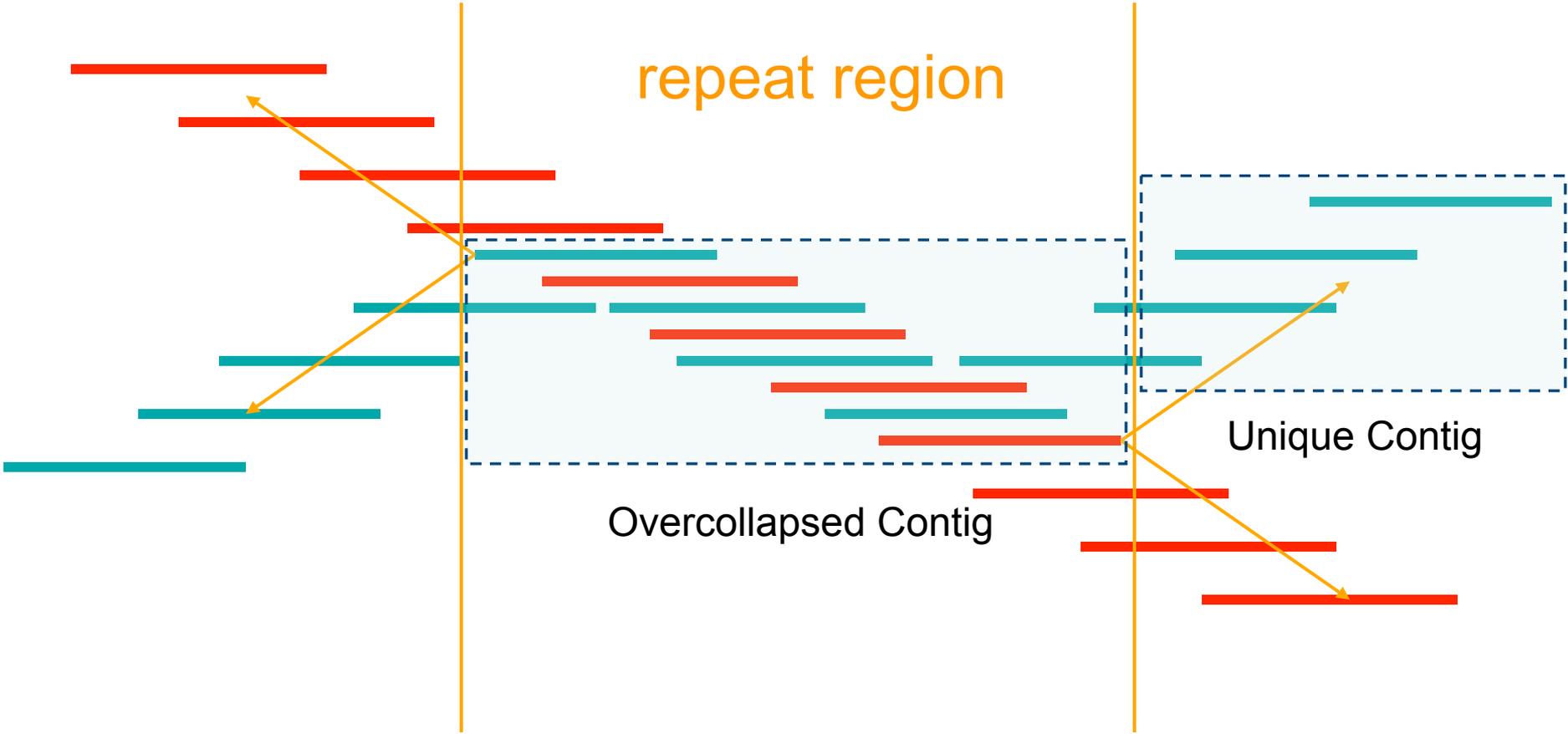
Reads that come from two regions of the genome (blue and red) that contain the same repeat



Note:  
of course, we don't know the "color" of these nodes



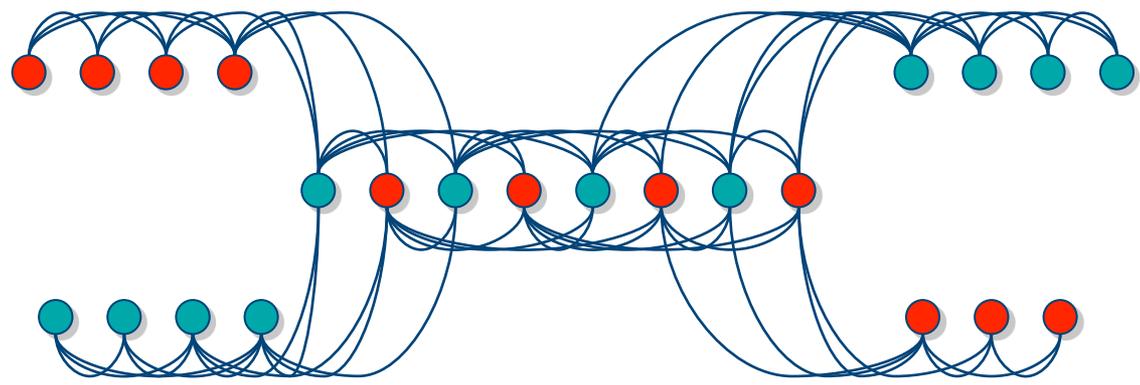
# 2. Merge Reads into Contigs



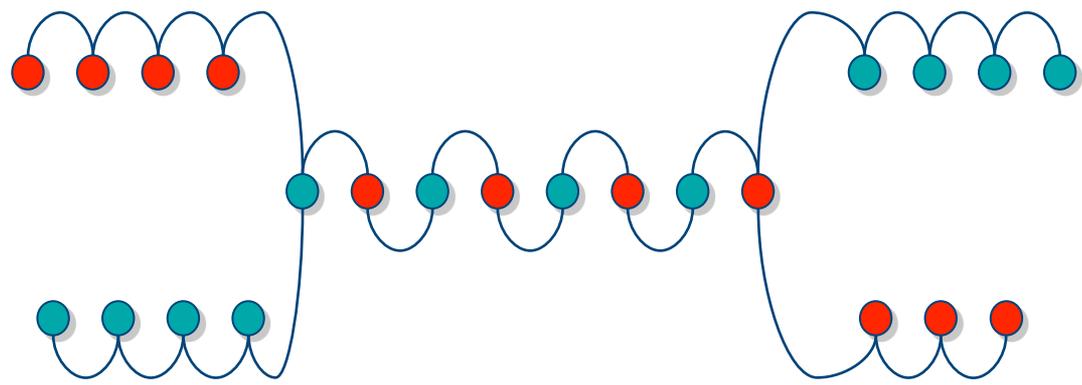
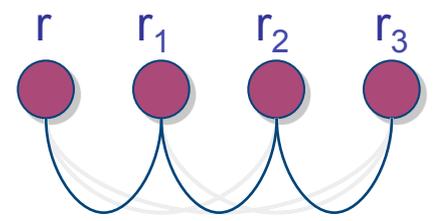
We want to merge reads up to potential repeat boundaries



# 2. Merge Reads into Contigs

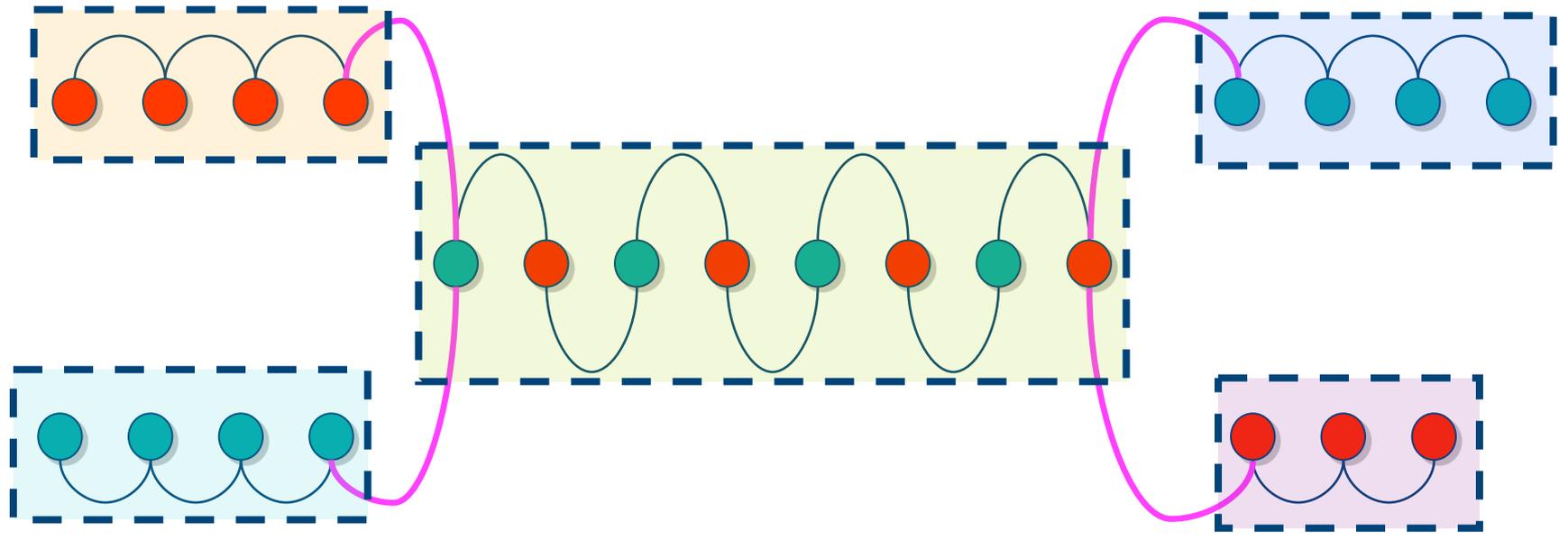


- Remove transitively inferable overlaps
  - If read  $r$  overlaps to the right reads  $r_1, r_2$ , and  $r_1$  overlaps  $r_2$ , then  $(r, r_2)$  can be inferred by  $(r, r_1)$  and  $(r_1, r_2)$





# 2. Merge Reads into Contigs





# Repeats, errors, and contig lengths

- Repeats shorter than read length are easily resolved
  - Read that spans across a repeat disambiguates order of flanking regions
- Repeats with more base pair diffs than sequencing error rate are OK
  - We throw overlaps between two reads in different copies of the repeat
- To make the genome **appear** less repetitive, try to:
  - Increase read length
  - Decrease sequencing error rate

## Role of error correction:

Discards up to 98% of single-letter sequencing errors  
decreases error rate  
⇒ decreases effective repeat content  
⇒ increases contig length



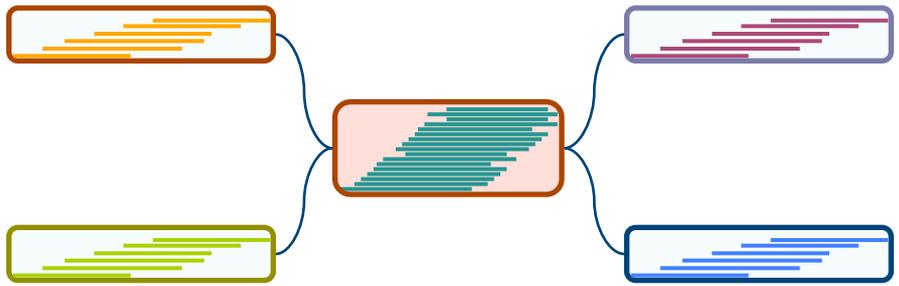
# 3. Link Contigs into Supercontigs



Normal density



Too dense  
⇒ Overcollapsed



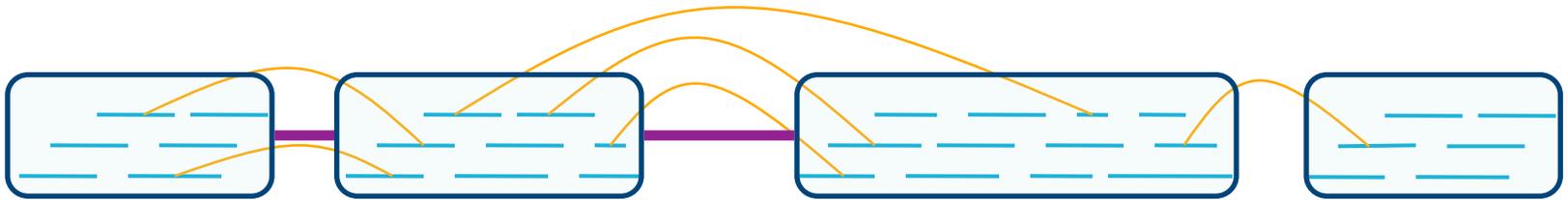
Inconsistent links  
⇒ Overcollapsed?



# 3. Link Contigs into Supercontigs

Find all links between unique contigs

Connect contigs incrementally, if  $\geq 2$  forward-reverse links



*supercontig*  
(aka scaffold)

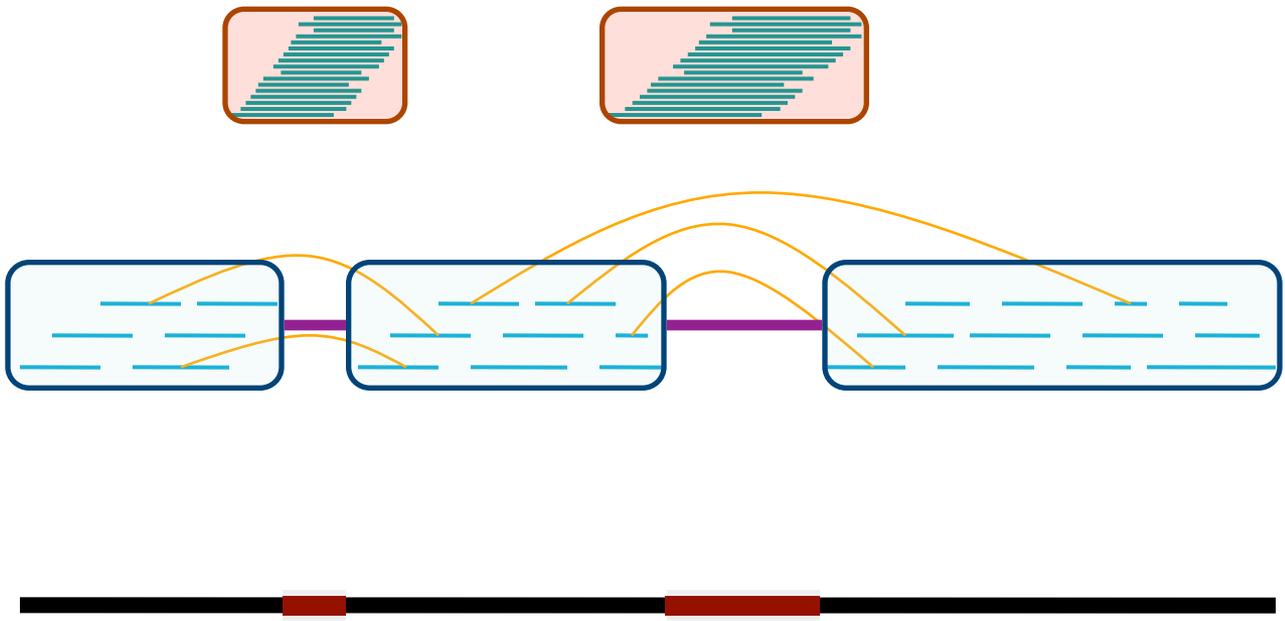


# 3. Link Contigs into Supercontigs

Fill gaps in supercontigs with paths of repeat contigs

Complex algorithmic step

- Exponential number of paths
- Forward-reverse links





# 4. Derive Consensus Sequence



Derive **multiple alignment** from pairwise read alignments

Derive each consensus base by weighted voting

(Alternative: take maximum-quality letter)



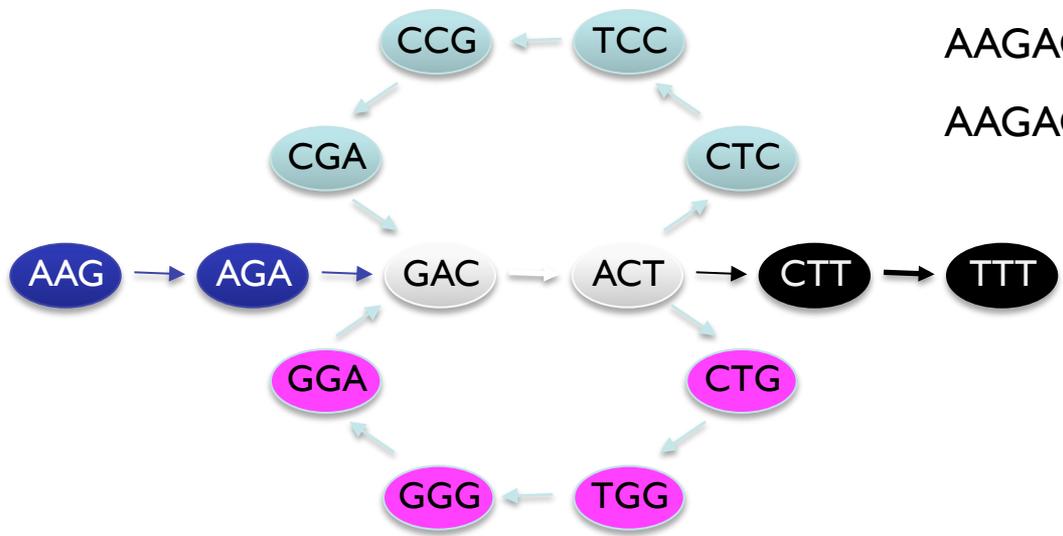
# De Bruijn Graph formulation

- Given sequence  $x_1 \dots x_N$ , k-mer length k,  
Graph of  $4^k$  vertices,  
Edges between words with (k-1)-long overlap

Reads

AAGA  
ACTT  
ACTC  
ACTG  
AGAG  
CCGA  
CGAC  
CTCC  
CTGG  
CTTT  
...

de Bruijn Graph

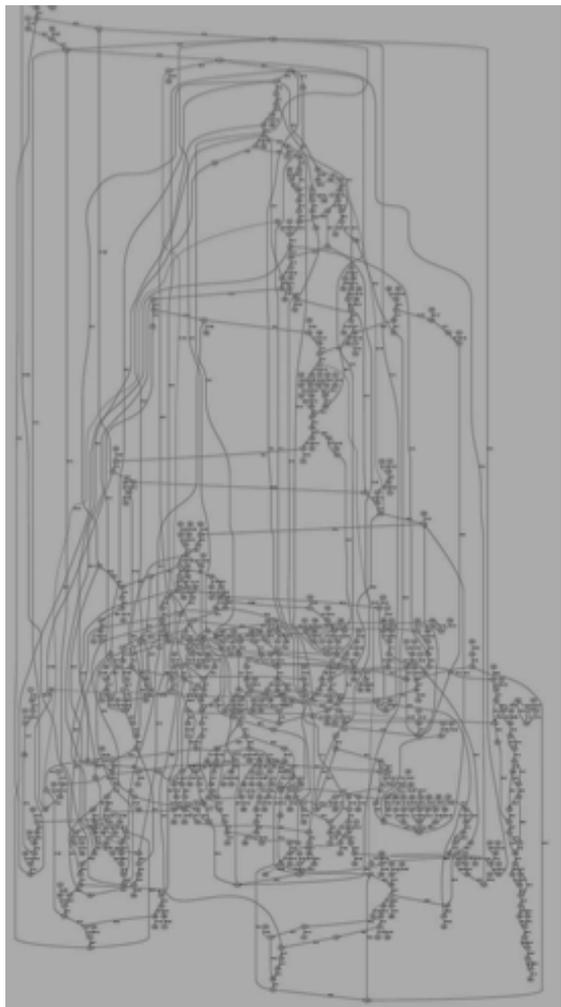


Potential Genomes

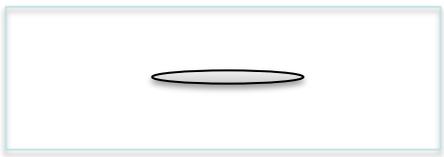
AAGACTCCGACTGGGACTTT  
AAGACTGGGACTCCGACTTT



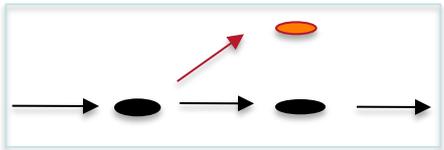
# Node Types



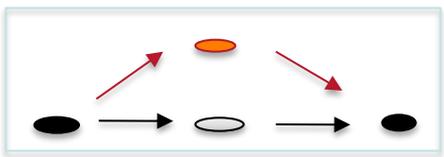
(Chaisson, 2009)



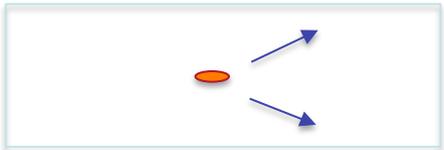
Isolated nodes (10%)



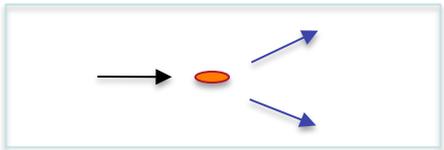
Tips (46%)



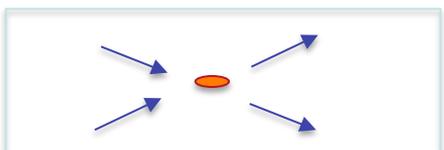
Bubbles/Non-branch (9%)



Dead Ends (.2%)



Half Branch (25%)

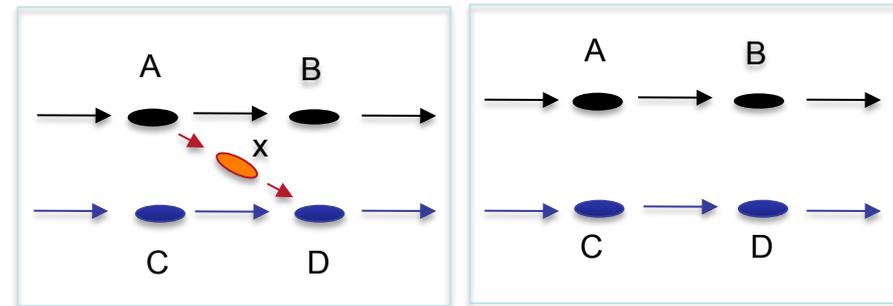
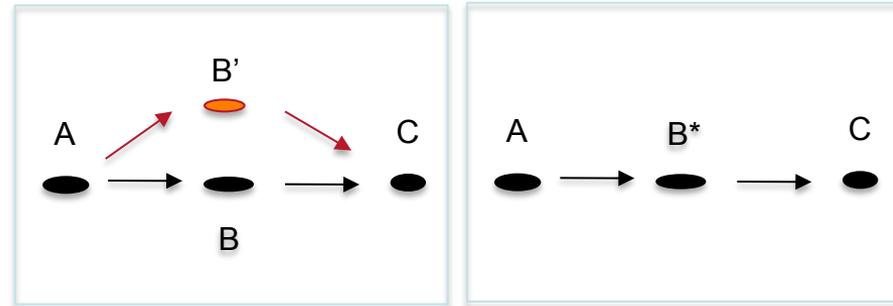
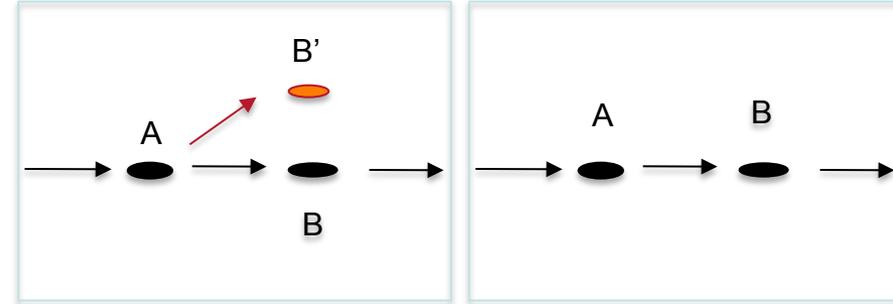


Full Branch (10%)



# Error Correction

- Errors at end of read
  - Trim off 'dead-end' tips
- Errors in middle of read
  - Pop Bubbles
- Chimeric Edges
  - Clip short, low coverage nodes





# De Bruijn Graph formulation

