

C G T A C G T A  
A C G T A C G T

# The era of long reads

**Sergey Koren**

Staff Scientist, Genome Informatics Section, NHGRI



National Human Genome  
Research Institute

@sergekoren 

—  
The **Forefront**  
of **Genomics**  
—



# Assembly review

# Genome Assembly

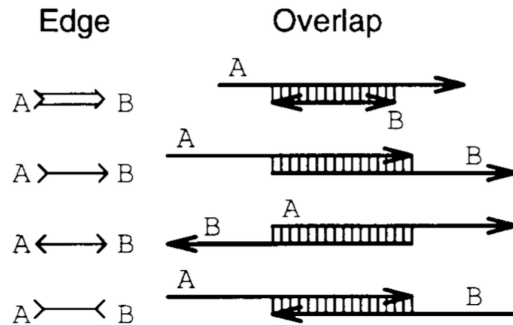
---

- ▶ Assembling a puzzle with a billion pieces



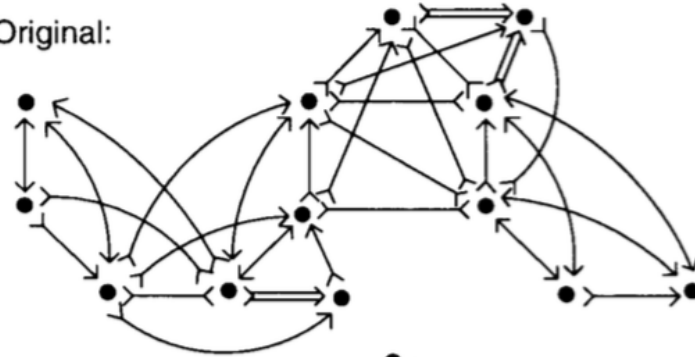
# Assembly the Celera way

- ▶ Step 0:
  - ▶ Find overlaps

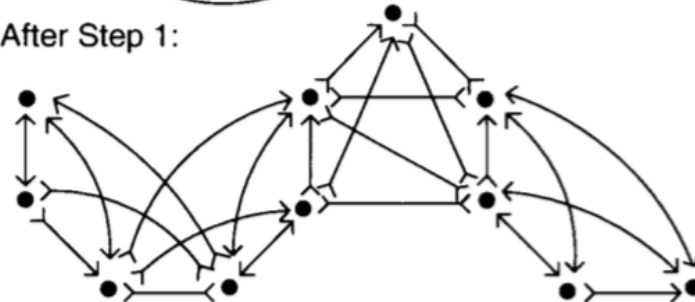


- ▶ Step 1:
  - ▶ Remove contained

Original:



After Step 1:

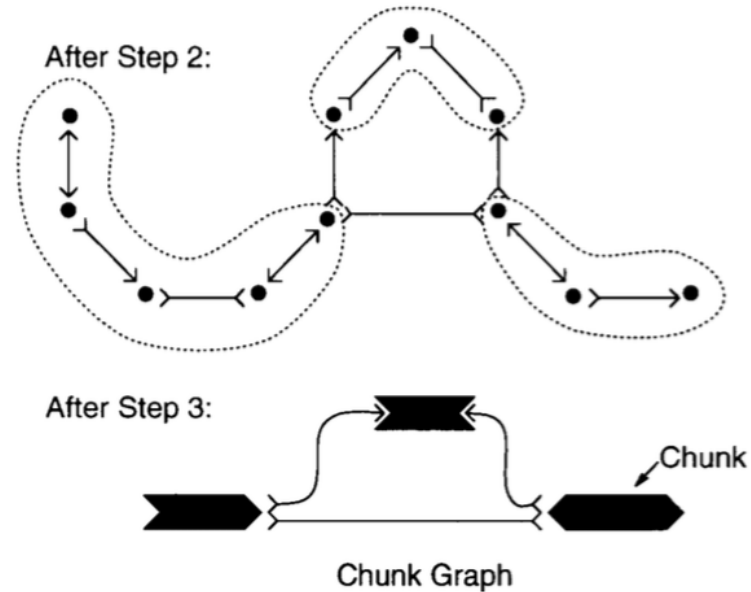




# Assembly the Celera way

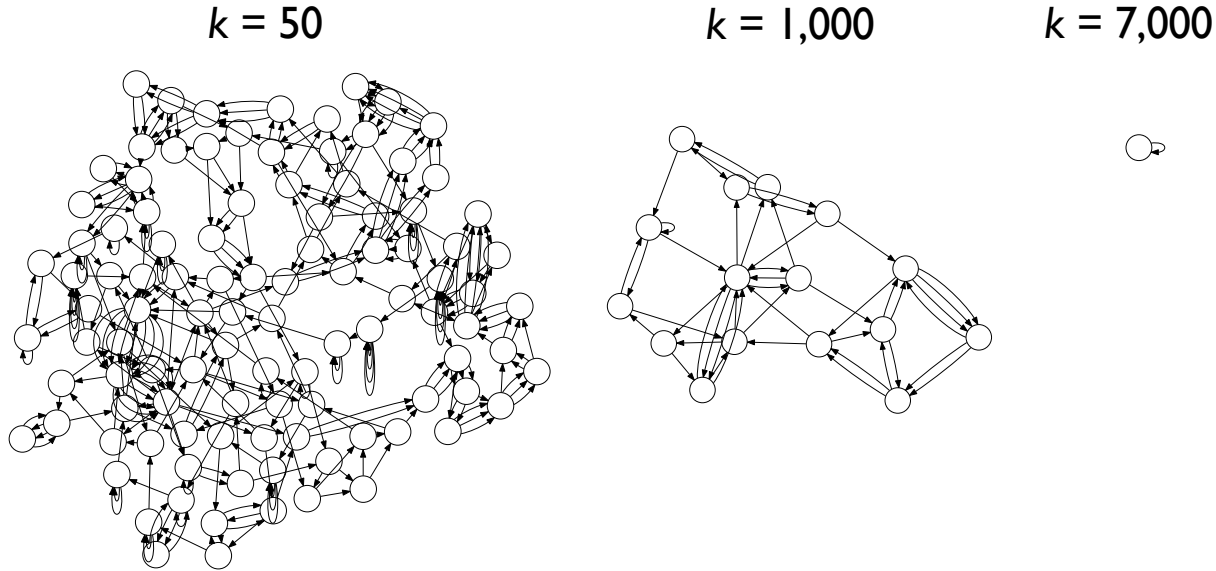
---

- ▶ Step 2:
  - ▶ Transitive reduction
- ▶ Step 3:
  - ▶ Collapse unique
- ▶ Output
  - ▶ “Unitigs”



# Read length matters (*E. coli*)

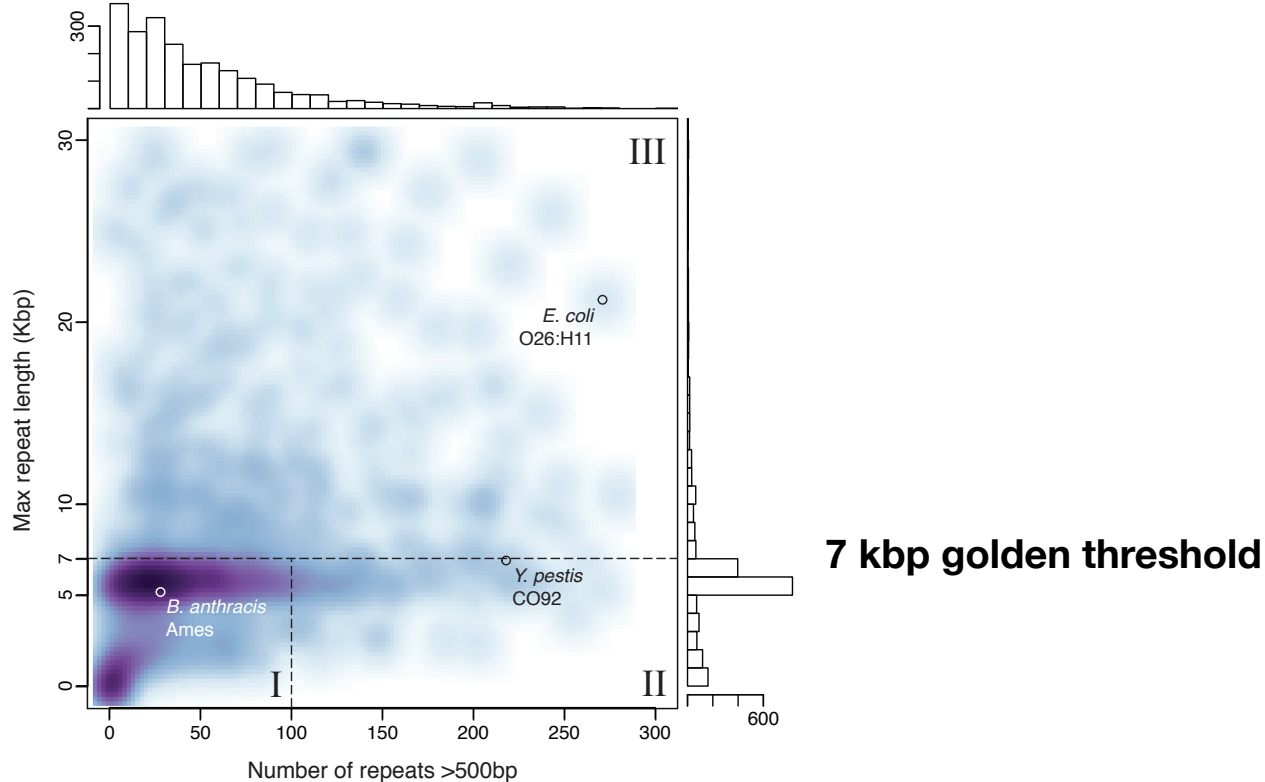
---



\* No errors, perfect coverage, uniform read length

► **One chromosome, one contig: complete microbial genomes from long-read sequencing and assembly.**  
Koren and Phillippy. *Current Opinion Microbiology* (2015)

# How long are microbial repeats?

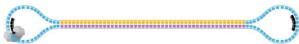


► **Reducing assembly complexity of microbial genomes with single-molecule sequencing.**  
Koren et al. *Genome Biology* (2013)



# A new era of sequencing

# PacBio Sequel II

- Single Molecule sequencer (one DNA strand)
  - Ligate adapters to make a bell 
  - Load molecules onto zero mode waveguides
  - Real-time polymerase sequencing
  - Video analysis
- Capable of sequencing long molecules
  - 10-60 kbp
- High error (85-90% accuracy) but random
  - Can read shorter reads multiple times
  - Converges to near-perfect consensus



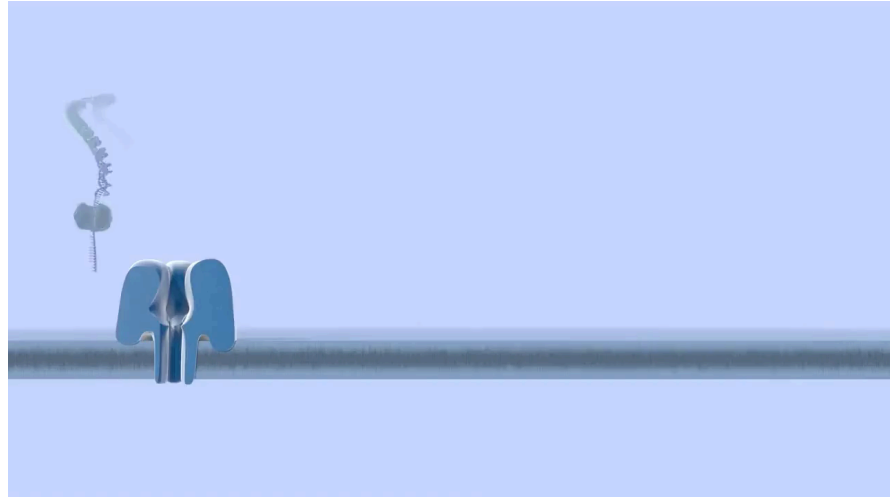
# Oxford Nanopore MinION

---

\$1000 (free) instrument

\$100 / bacterial genome

85–95% read accuracy



Oxford nanopore technologies





# Long read assembly in practice

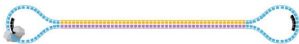
# Real data is messy

---

- ▶ Every technology has its own quirks
- ▶ Tools developed for one don't work on others
- ▶ Best tool may not be the theoretically optimal but best engineered



# Example: PacBio Sequel II

- Single Molecule sequencer (one DNA strand)
  - Ligate adapters to make a bell 
  - Load molecules onto zero mode waveguides
  - Real-time polymerase sequencing
  - Video analysis

## What can go wrong

- More than 1 read loaded into a well
  - Chimeric sequence when basecaller mixes them
- Read goes around adapter
  - Same sequence (forward then complement strand)
- Secondary DNA structure slows down/confuses polymerase



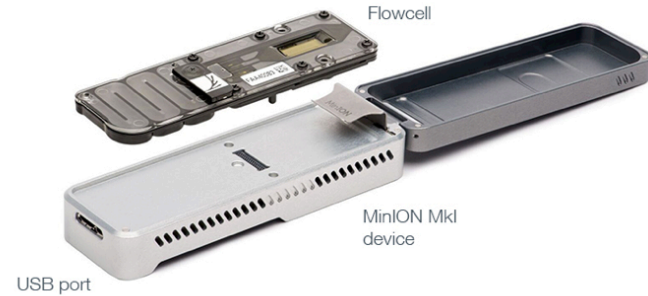
# Example: Oxford Nanopore MinION

---

- Single Molecule sequencer (one DNA strand)
  - Ligation or transposase to add adapter
  - Load molecules onto flowcell guides
  - DNA denatured in real time and passed through pore
  - Signal analysis to identify bases

## What can go wrong

- Two reads pass through same pore quickly
  - Chimeric sequence when not detected
  - Can be same as PacBio chimera (fwd then comp)
- Continuous current mistaken for empty pore
  - Single read split into multiple parts
- DNA structure re-folding on the other side of the pore
  - Can make one strand higher error than the other



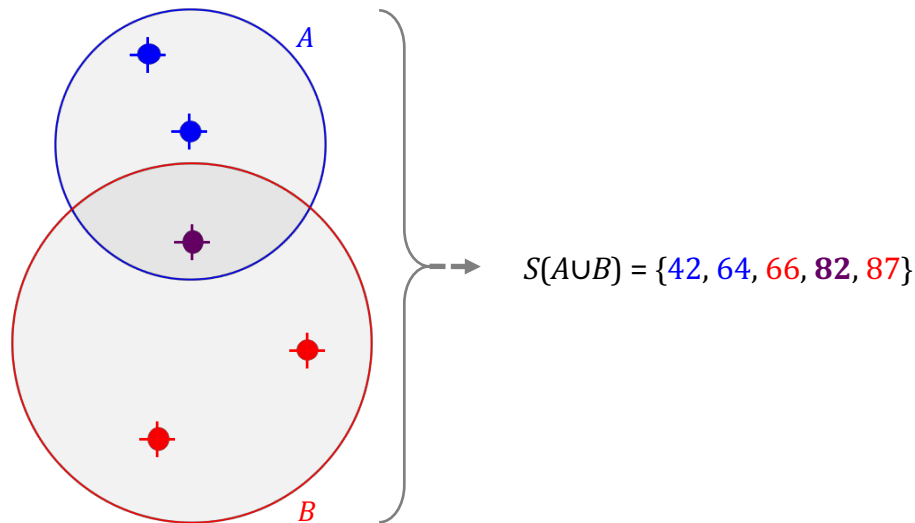
# In summary

---

- ▶ Long-read data is noisy
  - ▶ Base errors
  - ▶ Chimeric reads
  - ▶ *Solution*: read clustering, correction, and trimming
- ▶ Overlaps are long, and graph is big
  - ▶ All-pairs alignment is slow
  - ▶ Full graph is a giant tangle (due to repeats)
  - ▶ *Solution*: MinHash “best” overlap graph
- ▶ *D. melanogaster* results
  - ▶ Celera Assembler v8: **630,000** CPU hours, 15 Mbp NG50
  - ▶ Canu v1: **500** CPU hours, 21 Mbp NG50

# Fast overlapping with MinHash

---



$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \approx \frac{|S(A \cup B) \cap S(A) \cap S(B)|}{|S(A \cup B)|}$$

# *tf-idf* weighted MinHash

---

chief  
detective eat  
elephant followed inspector leave man sir  
thousand

The Stolen White Elephant by Mark Twain

eye head heard heart knew louder  
mannightopen sound

The Tell-Tale Heart by Edgar Allan Poe

away burmans crowd elephant faces people rifle  
seemed shoot shot

Shooting an Elephant by George Orwell

# A few extra details

---

- ▶ Throw hashes in hash table for all-pairs speedup
  - ▶ Only look at reads sharing some minimum number
- ▶ Jaccard based on k-mers, want a base error rate
  - ▶ Estimate from k-mers in the first round of overlapping
  - ▶ Compute exactly in the second round for contigging
- ▶ *tf-idf* weighted MinHash
  - ▶ Common repeats more likely to get larger hash value
  - ▶ Distinctive words more likely to get smaller hash value
  - ▶ Lower memory and runtime *without* k-mer filtering
- ▶ Keep position for each hash
  - ▶ Can be used to approximate the overlap bounds
  - ▶ (See German tank problem)

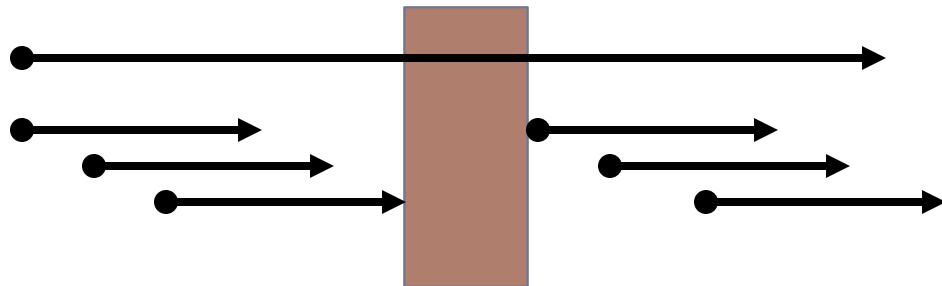


\* And it's written in Java

# Overlap-based correction and trimming

---

- ▶ Every (long) read corrected by its overlaps
  - ▶ Consensus called for covered bases
  - ▶ Missing coverage suggests low-quality or chimeras
  - ▶ Read correction acc: >99% PacBio, <98% Nanopore

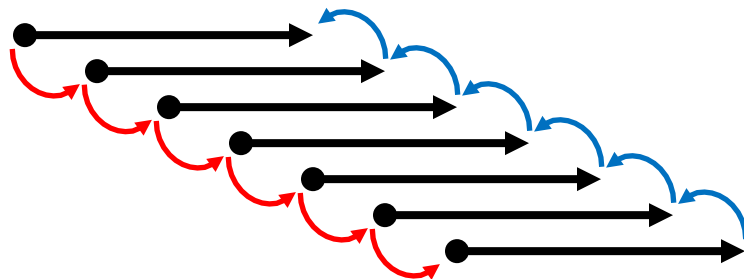


- ▶ Data cleaning is key to assembly
  - ▶ Necessary, not glamorous

# Best overlap graph

---

- ▶ After transitive reduction, only best are left
  - ▶ With enough coverage, nearly a global alignment
  - ▶ Find the “best” 5’ and 3’ overlap for each read
  - ▶ Build a graph from these edges



- ▶ Greedy approach, can be misled by repeats
  - ▶ Works great if given only “true” overlaps



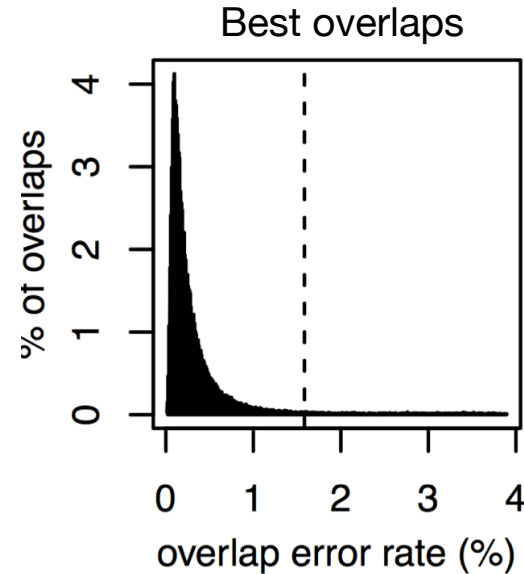
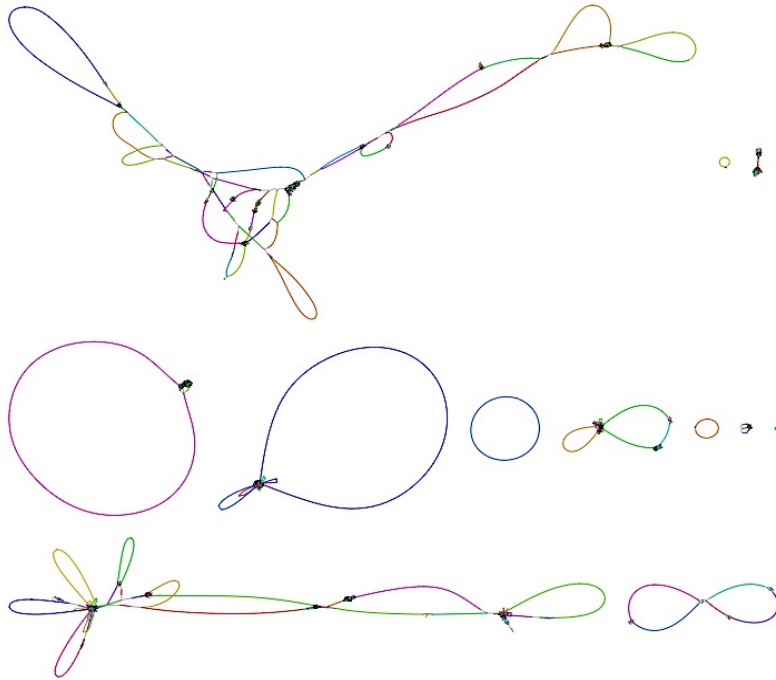
# Check your work

---

- ▶ Overlap filtering + greedy = pretty good
  - ▶ Automatically split divergent repeats and alleles
- ▶ Can still make mistakes, so...
  - ▶ Annotate repeats within contigs using overlaps
  - ▶ Check repeats for spanning reads
  - ▶ Check local error rate across each contig
  - ▶ Break on suspicion of misjoin
- ▶ Complete the graph with non-best overlaps

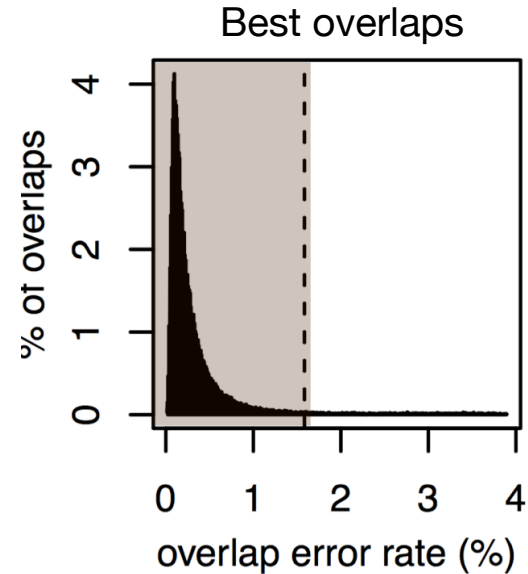
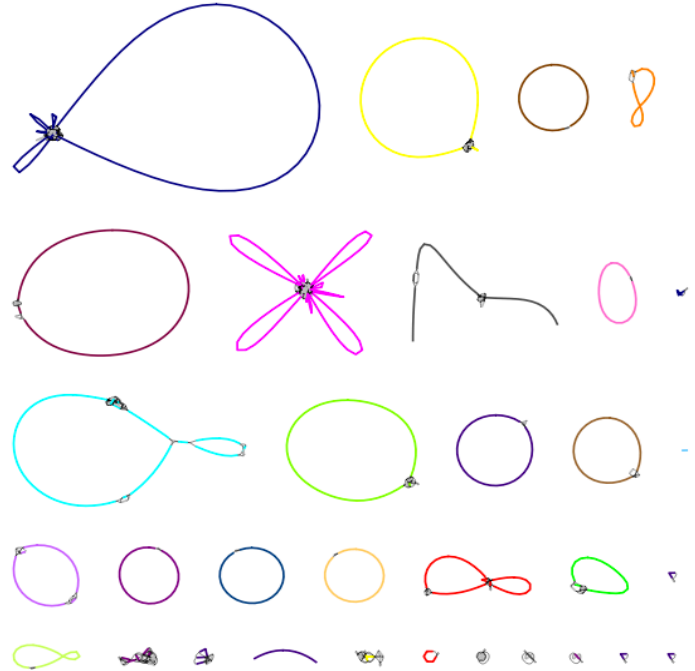
# Repeat and haplotype separation

---



► Don't know the read error rate a priori

# Repeat and haplotype separation



► Differentiate true from false overlaps



# Can long reads solve assembly?

Yes

# How long do reads need to be, for human?

---

## ▶ **How long are the repeats?**

- ▶ 7 kbp LINEs
- ▶ 1 Mbp+ rDNA arrays
- ▶ 1 Mbp+ centromere arrays
- ▶ 10 Mbp+ heterochromatin blocks

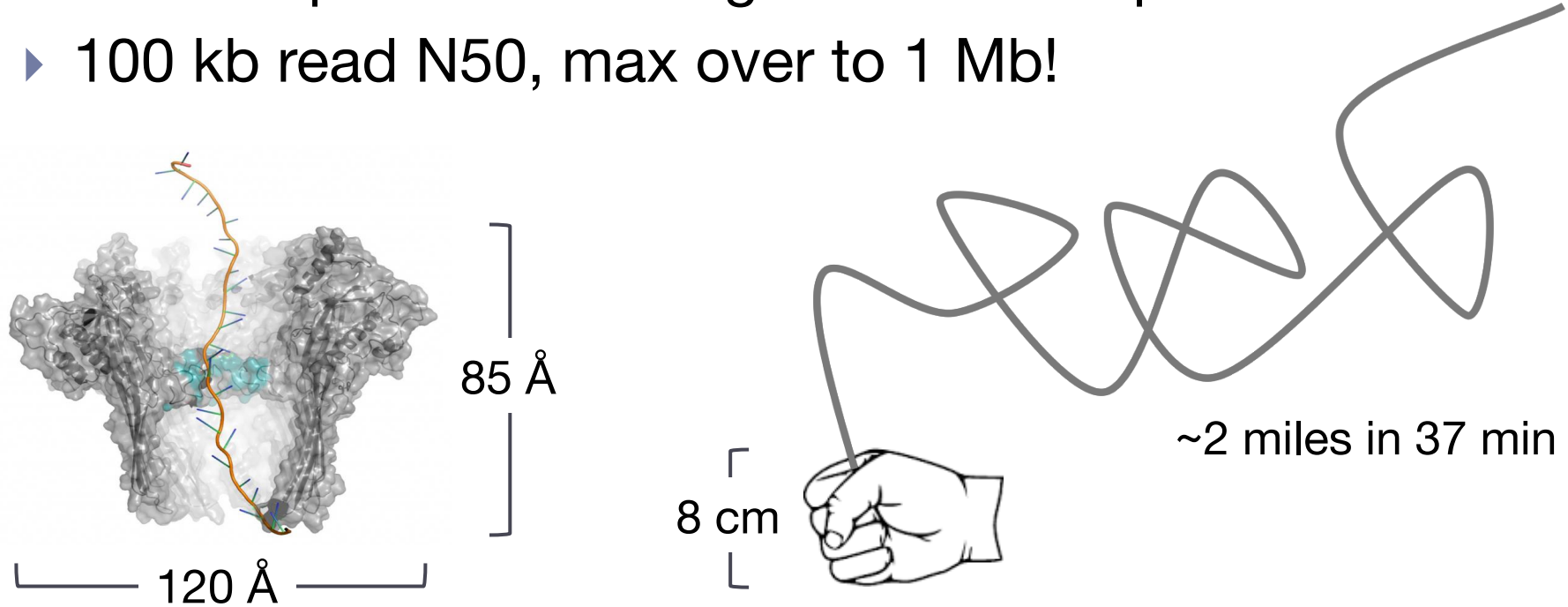
## ▶ **Coverage and accuracy matter too**

- ▶ 1,000X of 100 bp reads at 100% accuracy? **NO**
- ▶ 10X of 10,000,000 bp reads at 100% accuracy, **YES**
- ▶ 100X of 100,000 bp reads at 90% accuracy, **MAYBE?**

# Ultra-long read sequencing



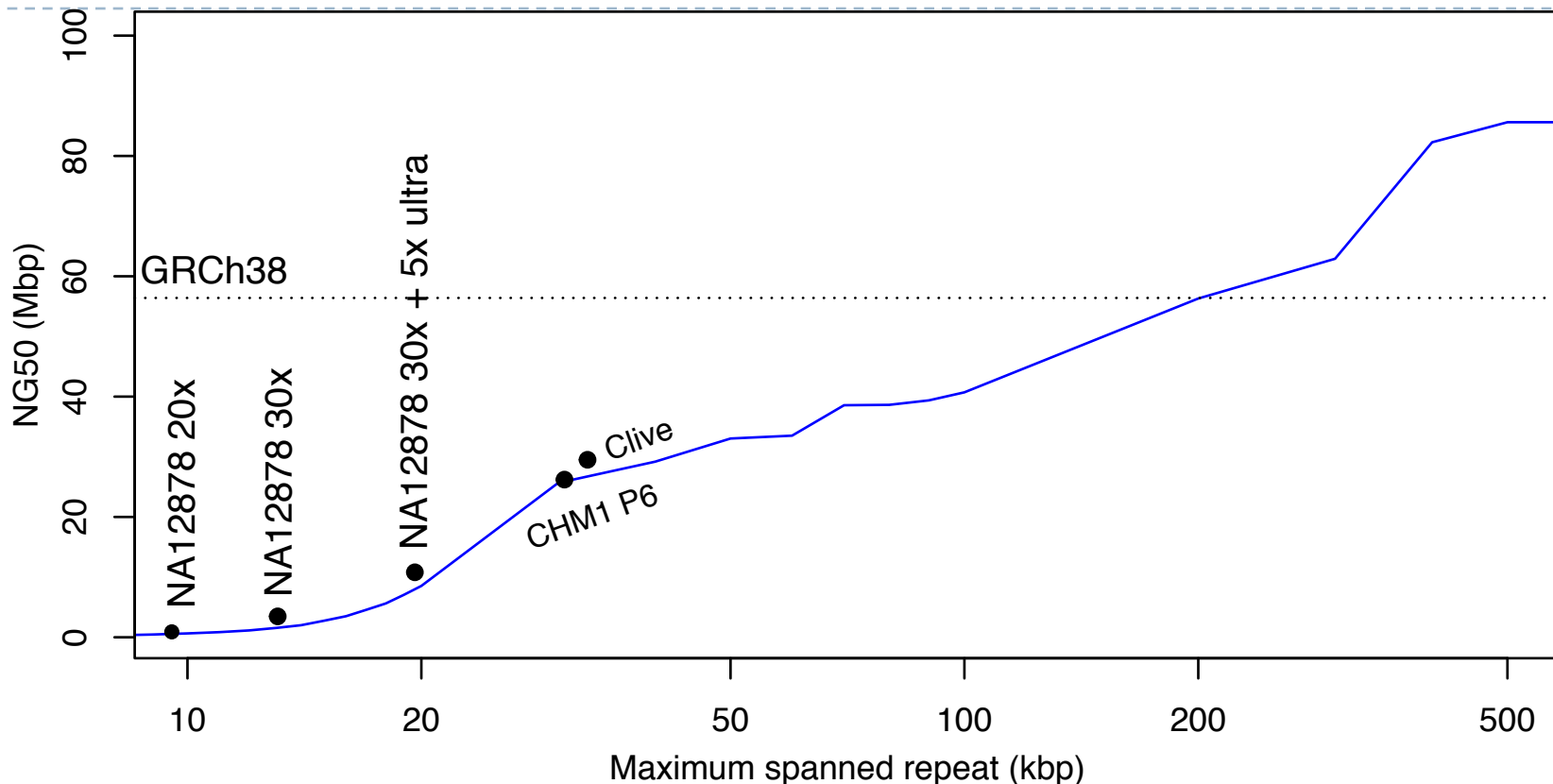
- ▶ ONT R9 pore: *E. coli* CsgG membrane protein
- ▶ 100 kb read N50, max over to 1 Mb!



\*Assuming 3.4 Å per bp, 1 Mbp = 3,400,000 Å (0.34 mm) = 40,000x height of the pore

▶ <http://lab.loman.net/2017/03/09/ultrareads-for-nanopore/> (Josh Quick & Nick Loman, U. Birmingham)

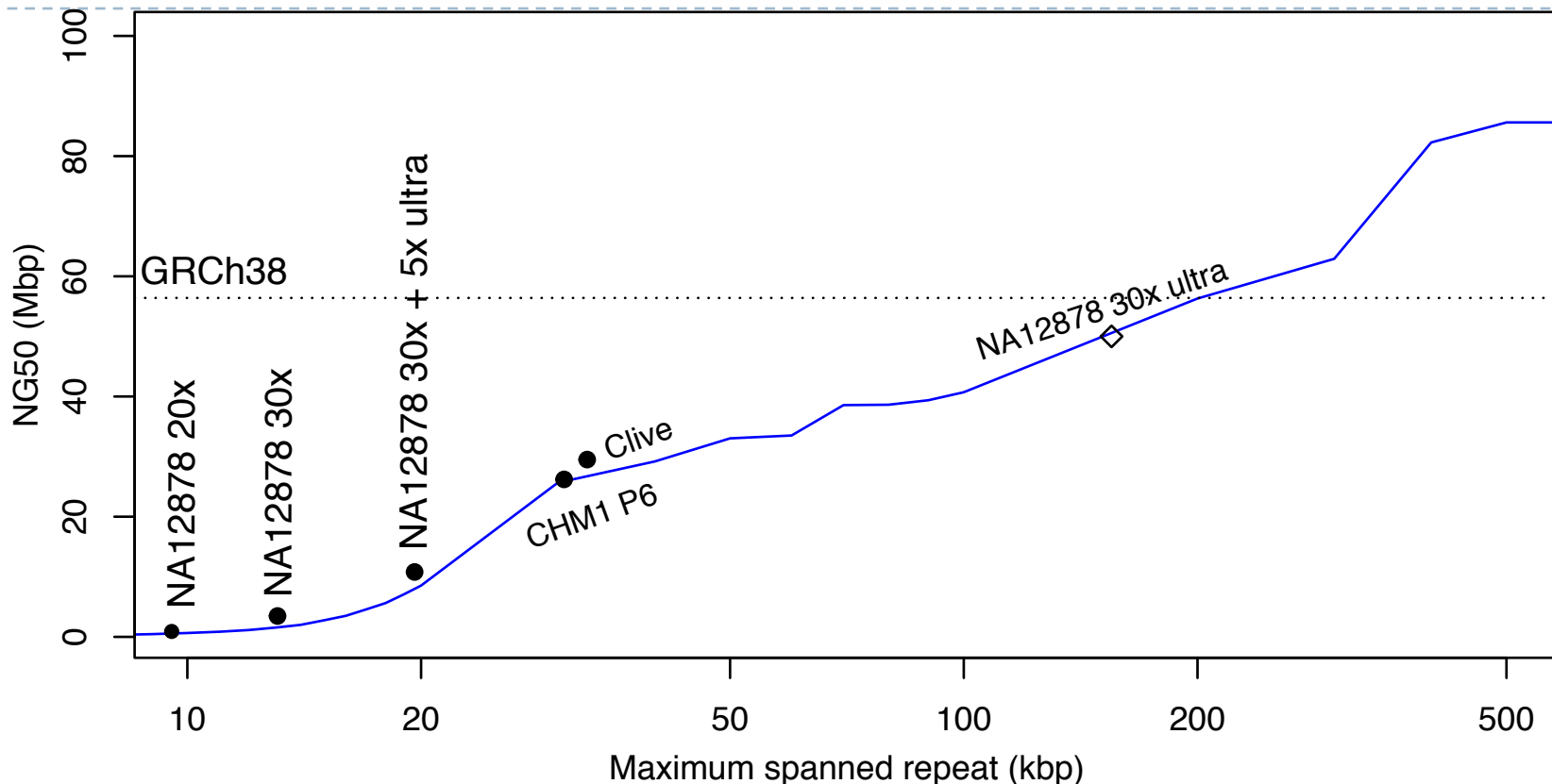
# Ultra-long read benefits



► **Nanopore sequencing and assembly of a human genome with ultra-long reads.**  
Jain, Koren, Miga, Quick, Rand, Sasani, Tyson, et al. *Nature Biotech* (2018)

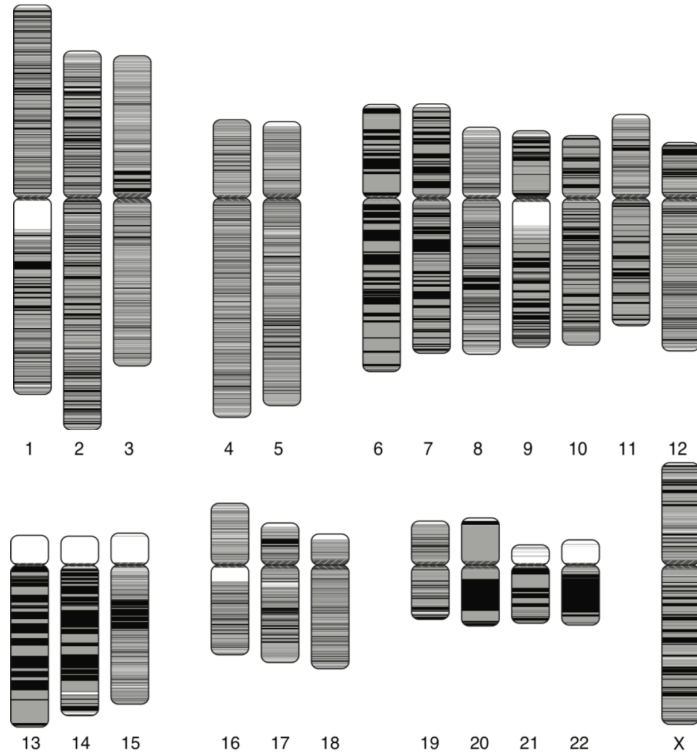


# Ultra-long read benefits



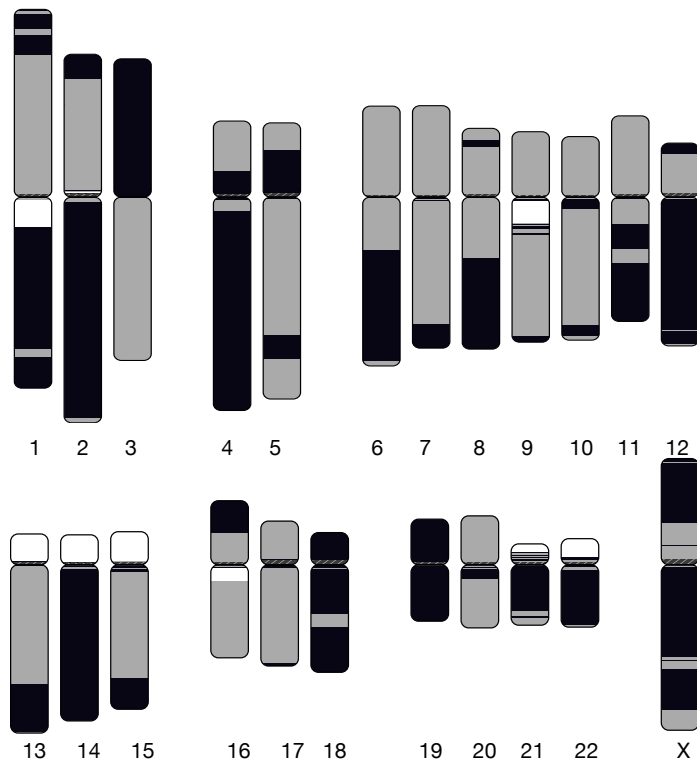
► **Nanopore sequencing and assembly of a human genome with ultra-long reads.**  
Jain, Koren, Miga, Quick, Rand, Sasani, Tyson, et al. *Nature Biotech* (2018)

# Human genome, 2001



ref28 / hg10 : N50 0.5 Mbp

# The human genome, 2017



## GRCh38

The Genome Reference Consortium consists of:



wellcome  
**sanger**  
institute

Wellcome Sanger Institute



MCDONNELL  
GENOME INSTITUTE

The McDonnell Genome Institute at Washington University

EMBL-EBI



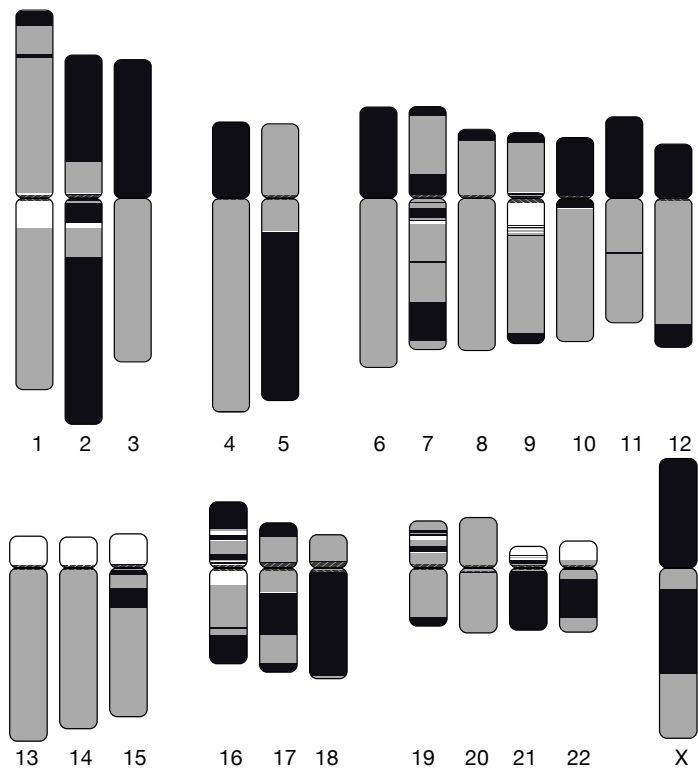
The European Bioinformatics Institute



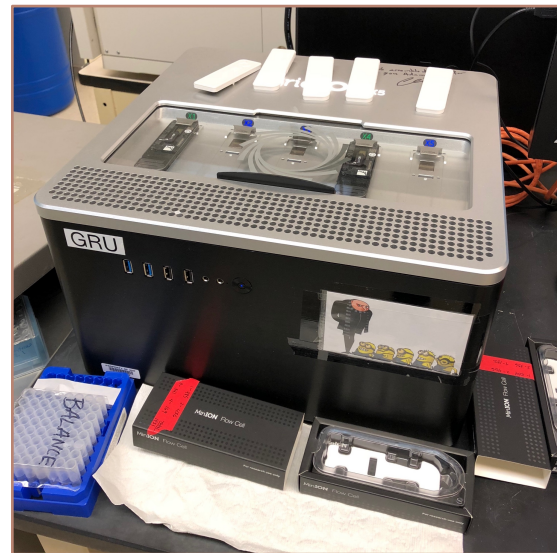
The National Center for Biotechnology Information

GRCh38 NG50 contig 56.4 Mbp

# The human genome, 2018



CHM13 NG50 contig 79.5 Mbp (50x UL ONT)

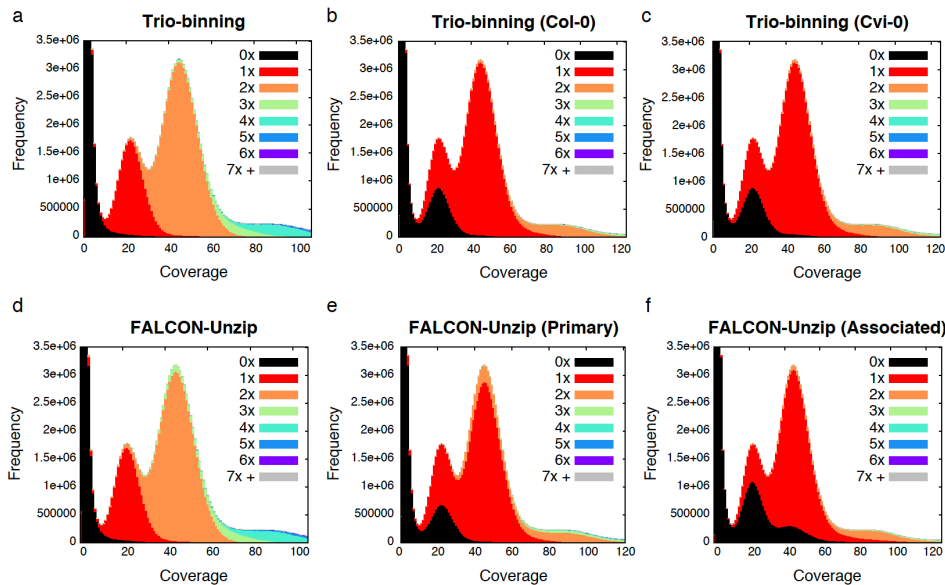




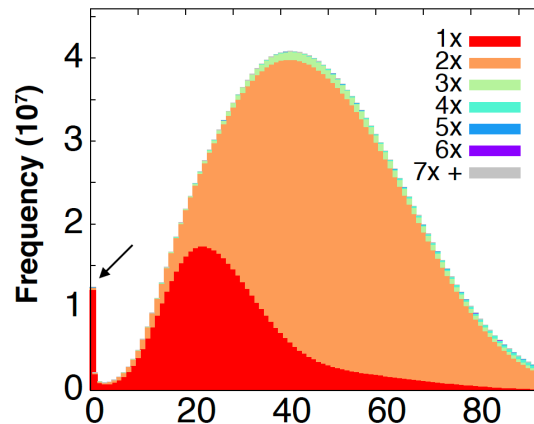
—

# An assembly is a hypothesis

# K-mers as a measure of completeness



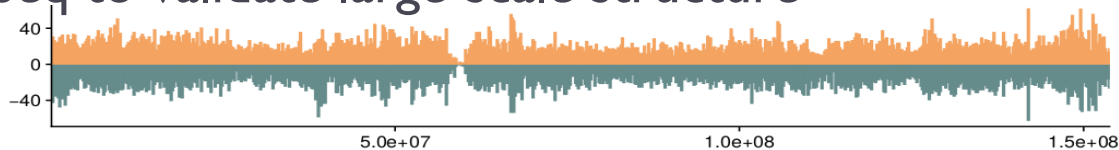
- ▶ K-mers only in assembly (misassembled bps)
- ▶ Haplotype completeness
- ▶ Over-assembled (duplications)
- ▶ Repeat copies ~ exp. copies?



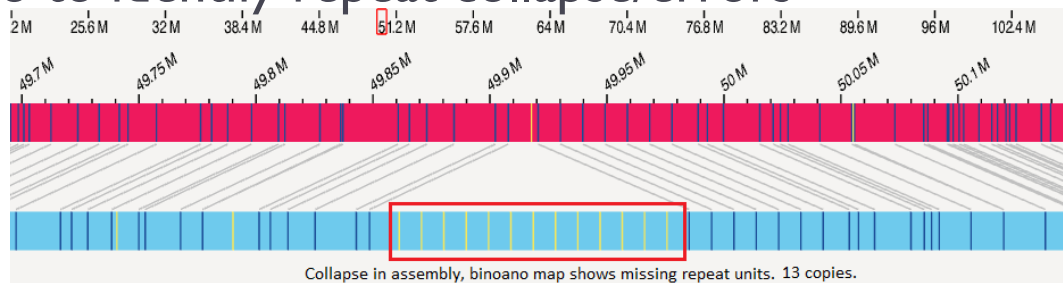
KAT Spectra-cn plots:  
<https://github.com/TGAC/KAT>  
Mapleson *et al.*,  
Bioinformatics (2016)

# Complementary technologies

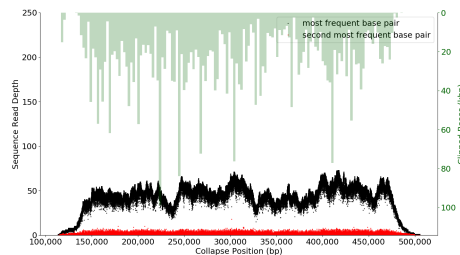
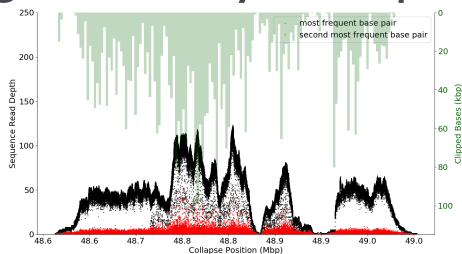
- StrandSeq to validate large-scale structure



- BioNano to identify repeat collapse/errors



- Mapping to identify low-quality regions



# Who said assembly wasn't cool?

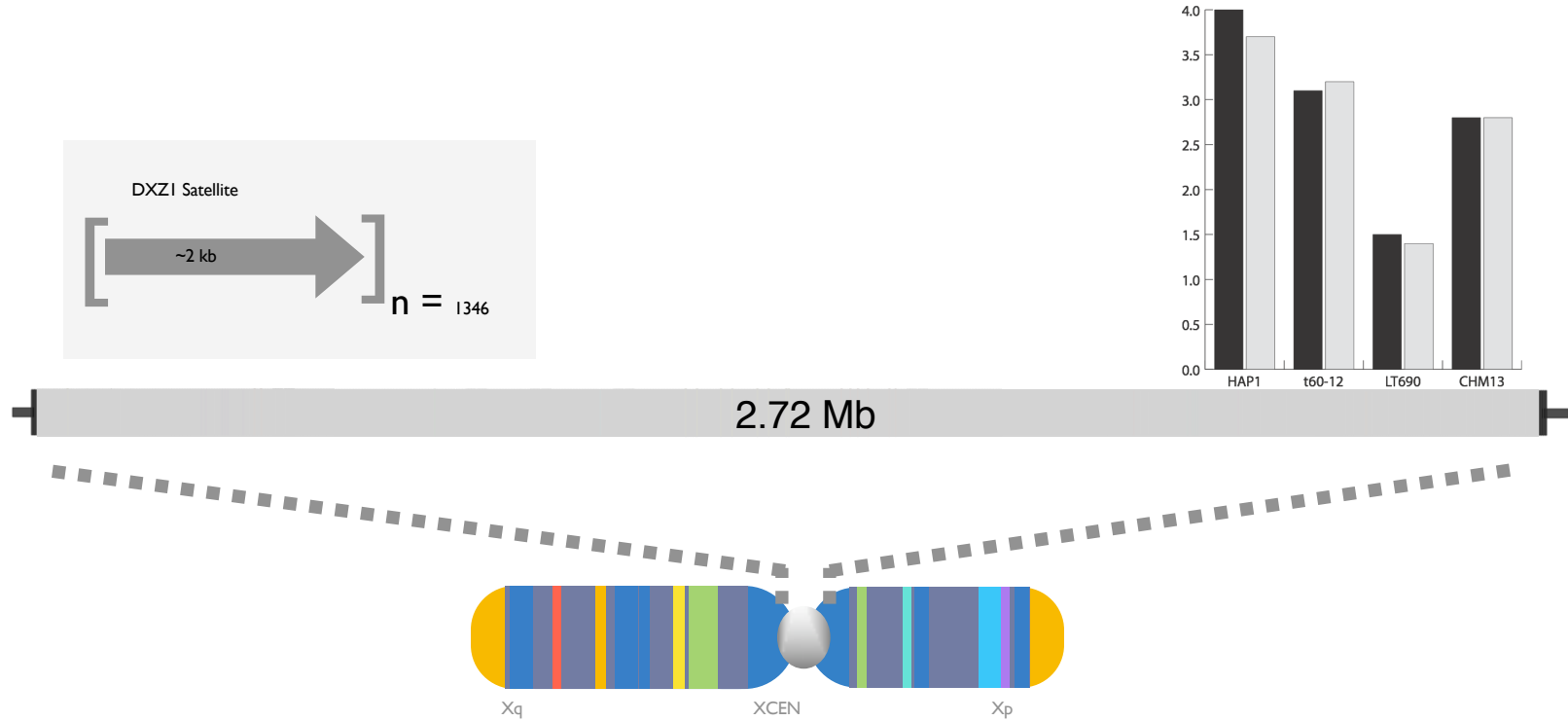






# Assembly is not solved

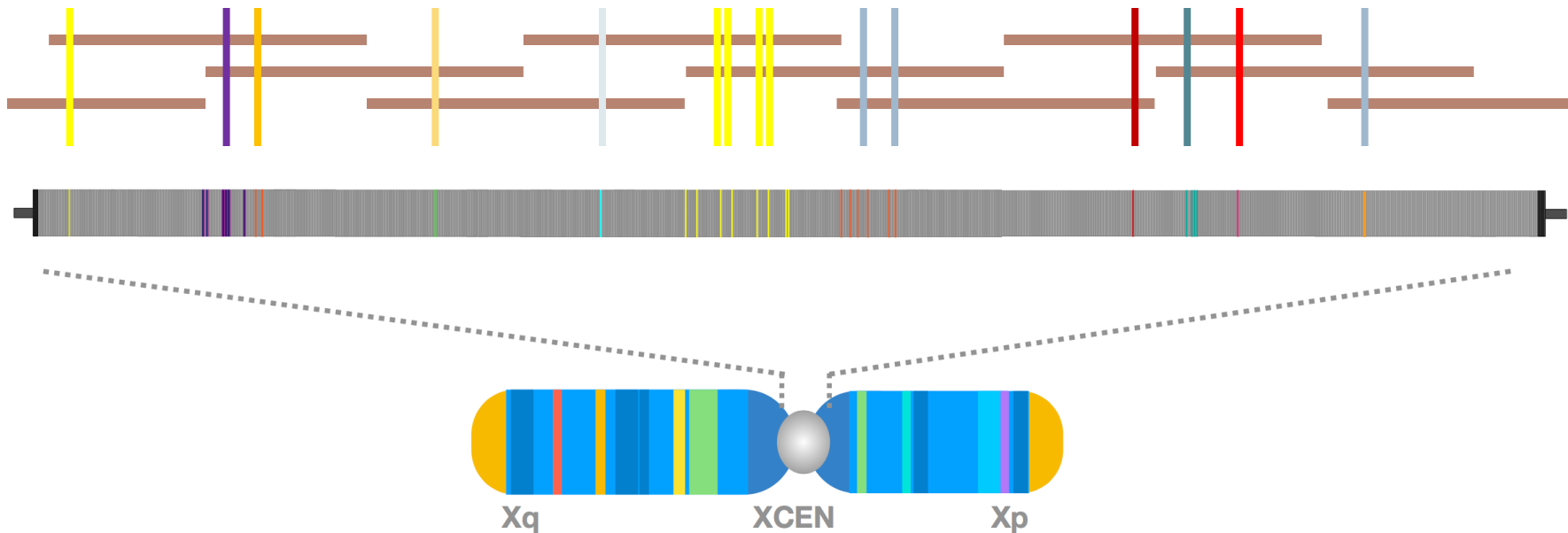
# X Centromere Detail



# Stitching across the X centromere



- ▶ Unique structural variants from PacBio
- ▶ Unique k-mers confirmed by Duplex-Seq

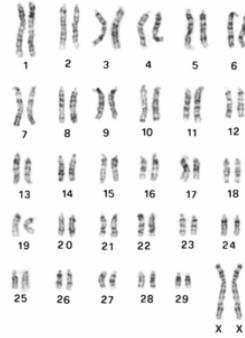




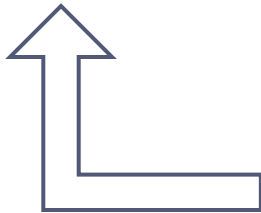
—

# There isn't a single “genome”

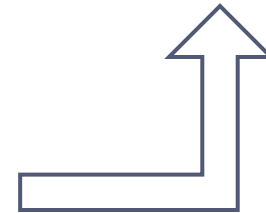
# The *genomes* assembly problem



Duke, highland sire



Molly, yak dam

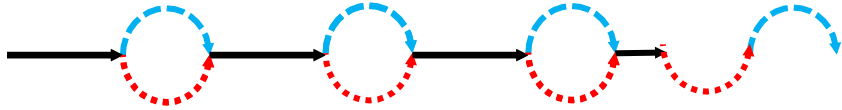


~1% heterozygosity

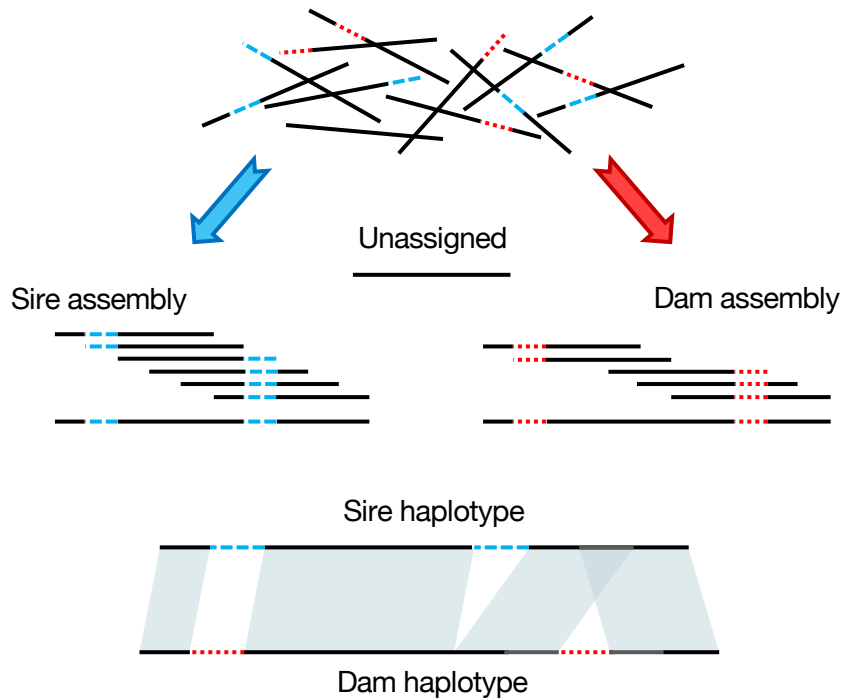
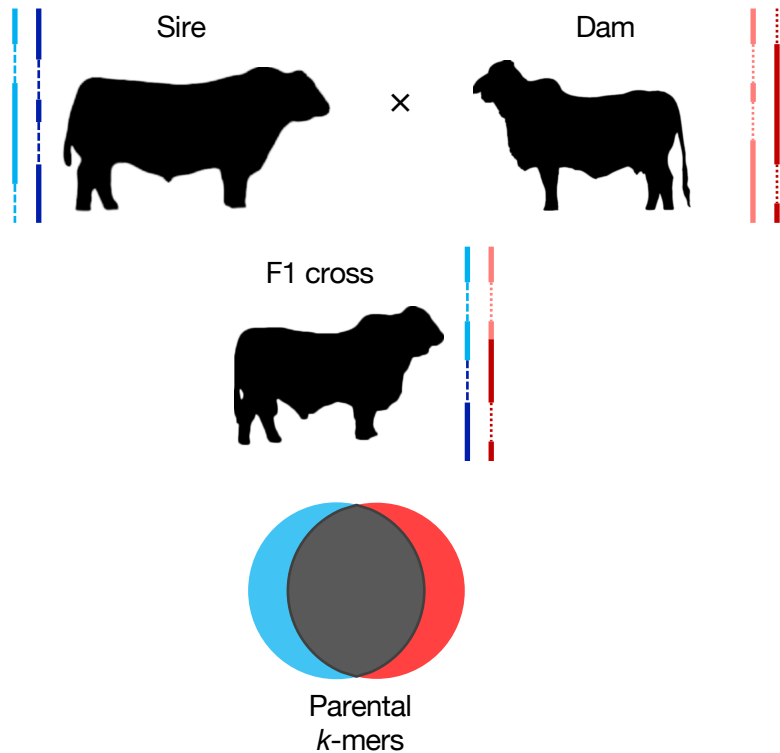


# State of the art: pseudo-haplotype

---



# Trio binning with TrioCanu



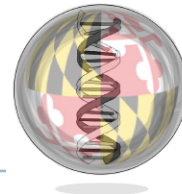
# Esperanza: The nearly perfect diploid



125x PacBio coverage (~60x per haplotype), **TrioCanu** haplotig NG50 70 Mbp, BUSCOs 94%



# Acknowledgements



[genomeinformatics.github.io](https://genomeinformatics.github.io)

- ▶ Sergey Koren
- ▶ Brian Walenz
- ▶ Alexander Dilthey
- ▶ Arang Rhie
- ▶ Brian Ondov
- ▶ Chirag Jain
- ▶ Anna Sappington



AD



CJ



SK



BO



AP



AR



AS



BW

[canu.readthedocs.io](https://canu.readthedocs.io)

- ▶ Sergey Koren
- ▶ Brian Walenz
- ▶ Konstantin Berlin
- ▶ Jason Miller
- ▶ GM12878 collaborators
- ▶ T2T collaborators
- ▶ VGP collaborators
- ▶ Cattle Collaborators