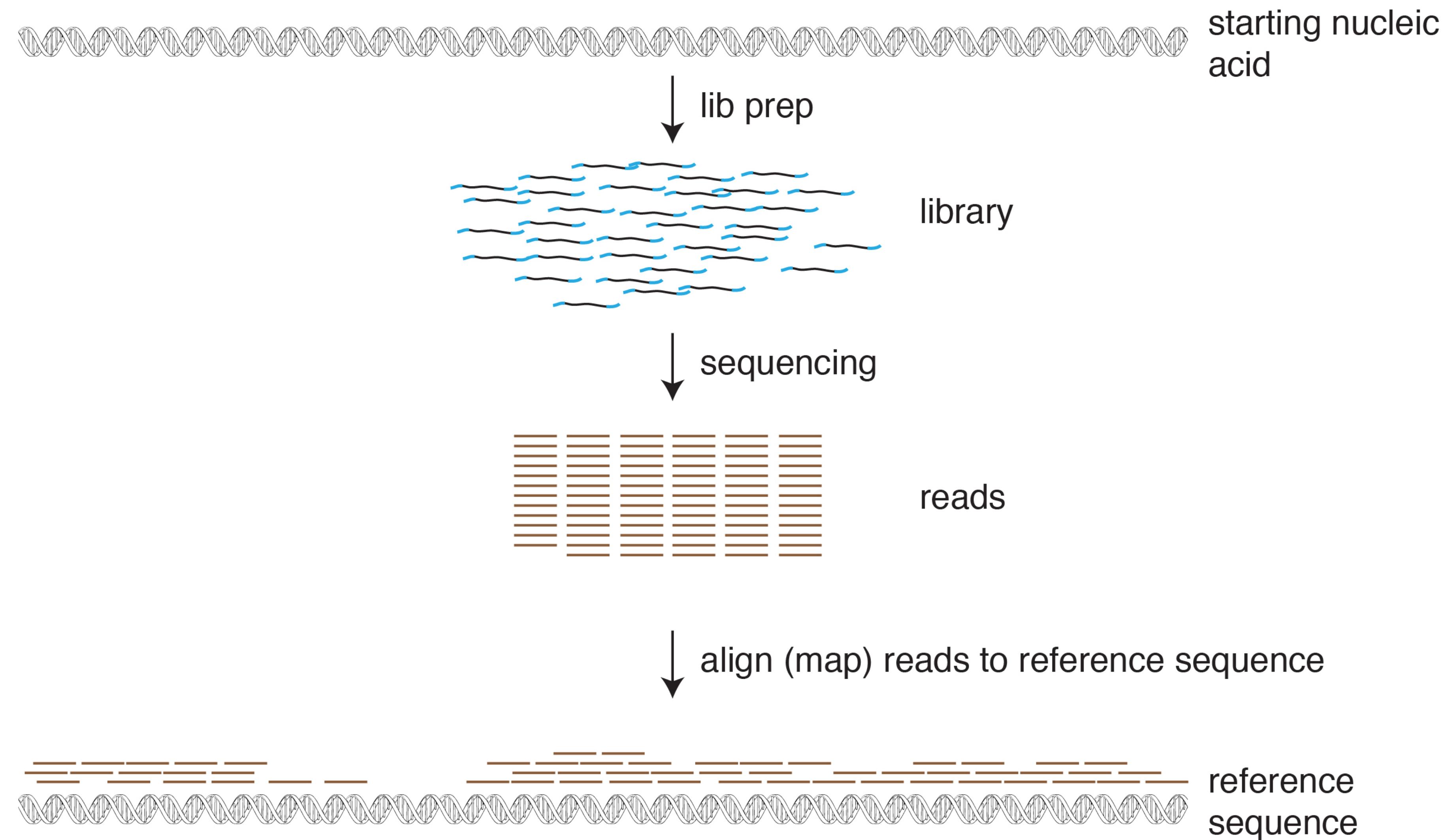


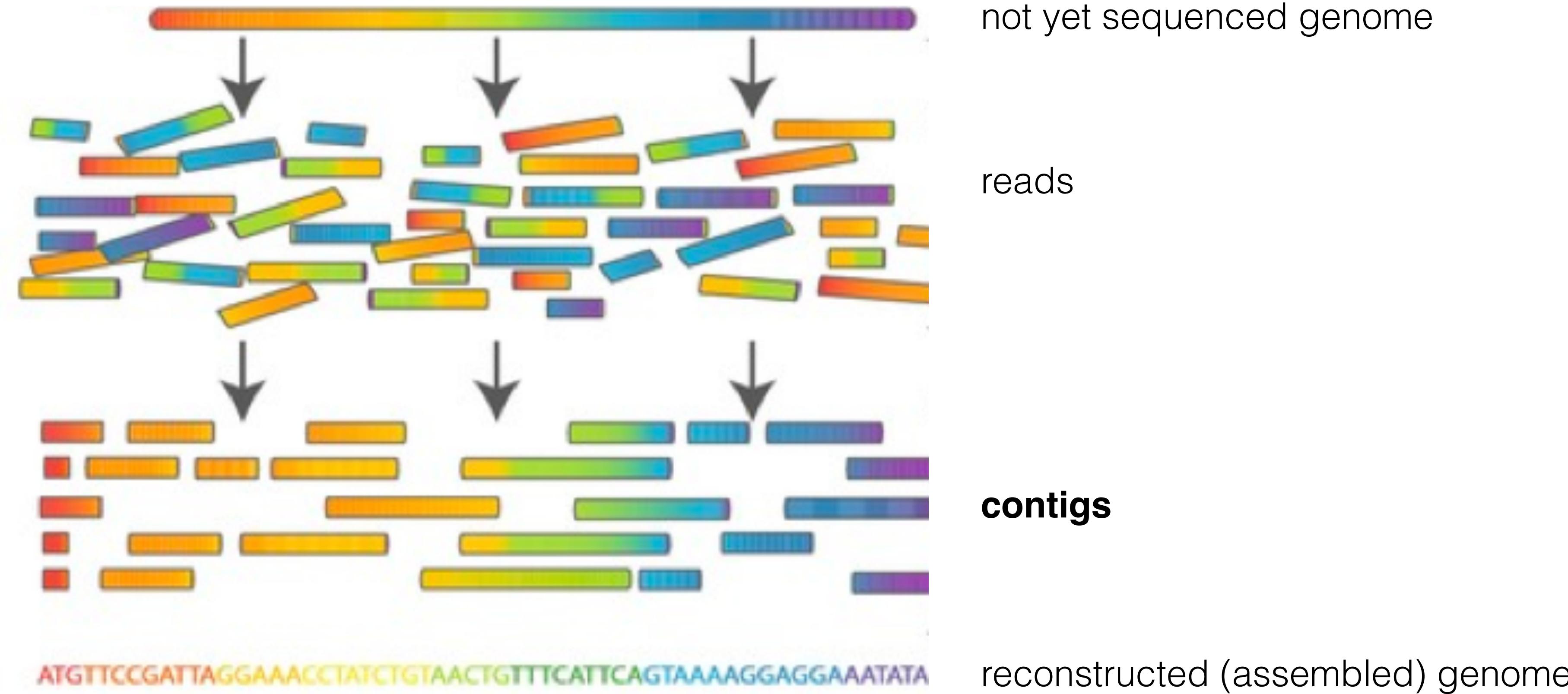
# Genome Assembly

Mark Stenglein, MIP 280A4

**Mapping** is the process by which sequencing reads are aligned to the region of a genome from which they derive.

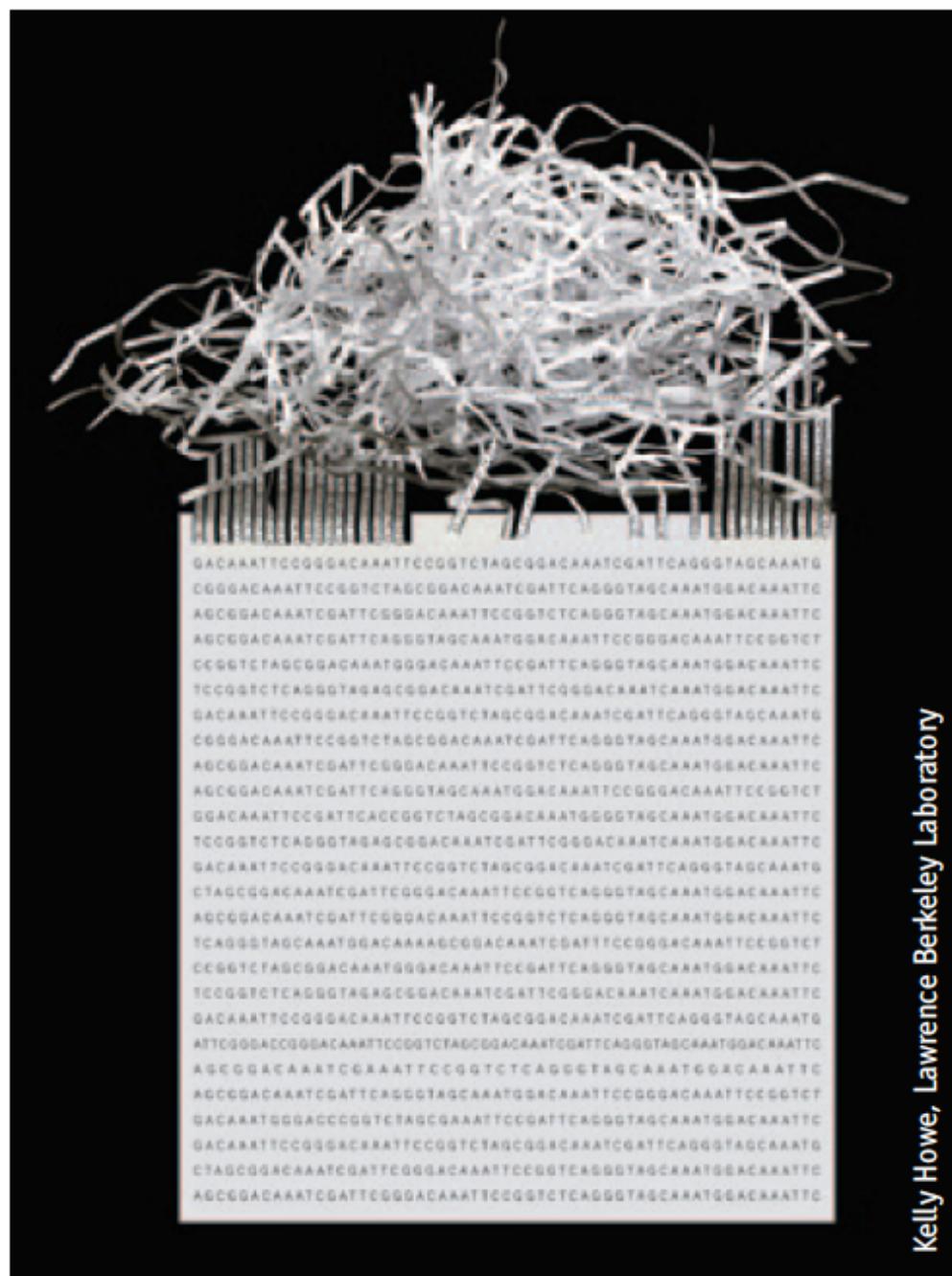


Genome **assembly** is the process of trying to reconstruct a genome sequence from reads (making a new reference sequence)



# Genome assembly is the process of *attempting* to reconstruct a genome sequence

An assembly is only a “putative reconstruction” of the genome sequence [Miller, Koren, Sutton (2010)]



Kelly Howe, Lawrence Berkeley Laboratory

Baker M (2012) Nat Methods



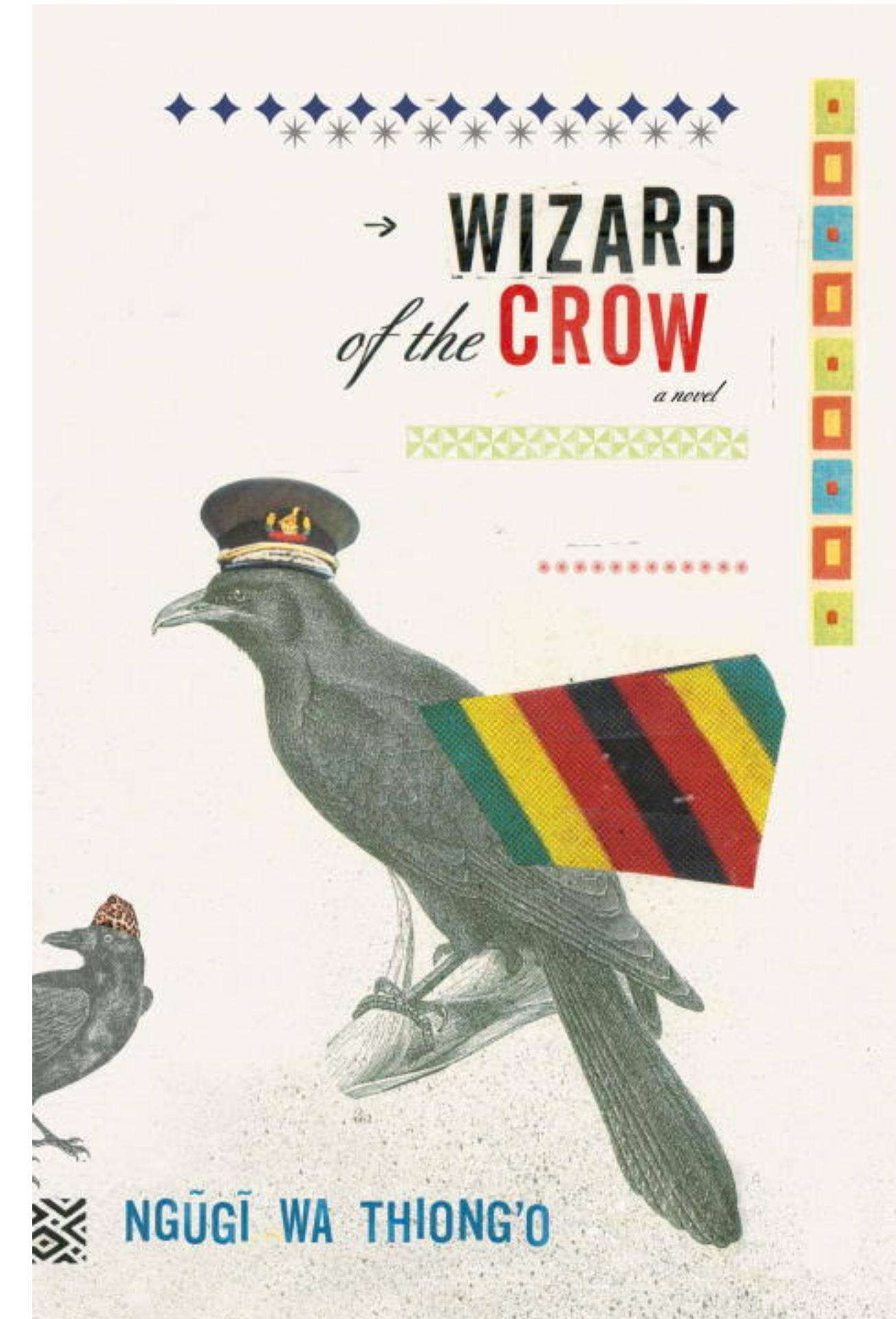
Keith Bradnam, UC Davis

# Genome assembly exercise

Your job is to assemble the ‘genome’ from which the ‘reads’ you’ve been given derive.

## Rules/info:

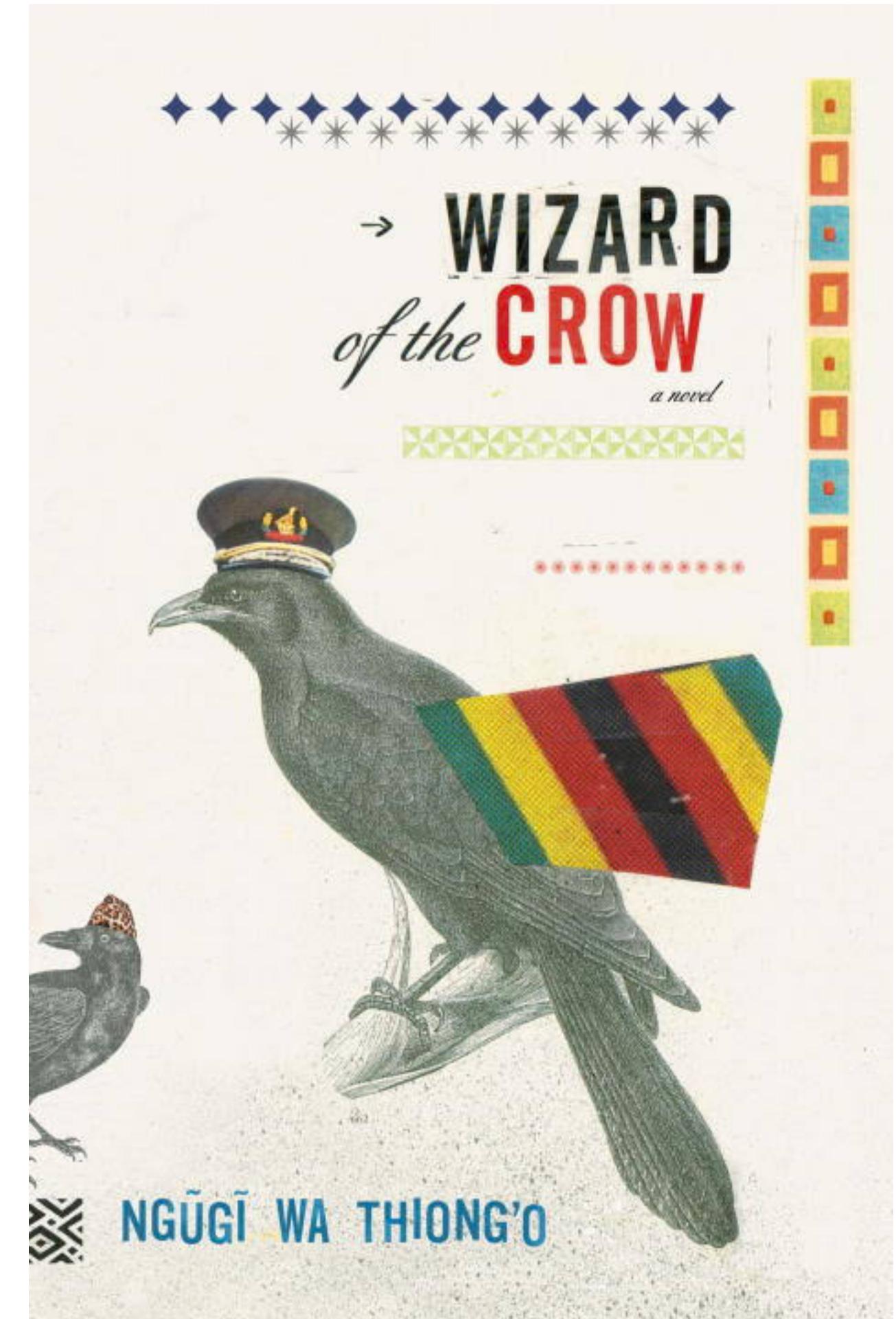
- Like real sequencing data, these reads contain errors.  
The error rate is ~2%
- These are single-end 11-base reads
- The average coverage is ~6x
- You’re not allowed to google the answer
- Also: the answer is in the slides: don’t cheat!
- You can use your computers (i.e. word processors or text editors) or paper and whatever strategy you want to do the assembly...



# Genome assembly paper exercise

“Jinn (Arabic), also romanized as djinn … are supernatural creatures in early Arabian and later Islamic mythology and theology.”

<https://en.wikipedia.org/wiki/Jinn>

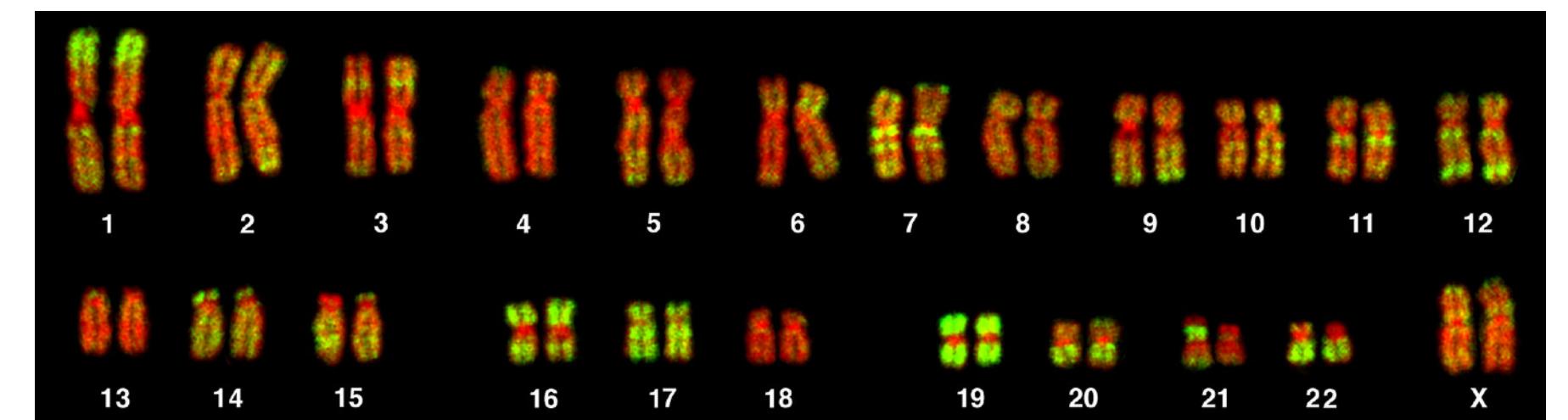


# Conclusion: assembly is not trivial!

In this exercise, the ‘genome’ was only 65 positions long, and its alphabet contained 26 ‘bases’ (more information rich)

the human *haploid* genome is 3 Gb

Eukaryotic genomes can have billions of bases and there are only 4 bases (less information)



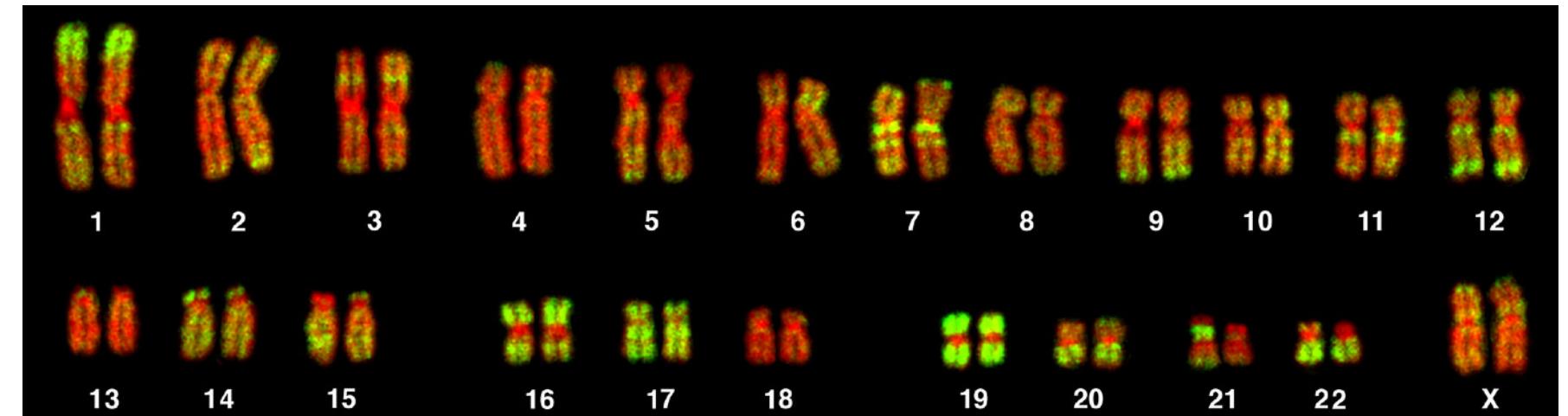
Bolzer et al (2005) PLoS Biol

# Some of the reasons that assembly is difficult

1) Genomes are full of repetitive sequences

Alu sequences in the human genome  
1 million copies, ~10% of the mass

2) Reads contain errors



Bolzer et al (2005) PLoS Biol

\_gew\_kjinns

get\_djinns\_

l\_get\_djinn

3) Uneven coverage, including possibly no coverage for particular regions (e.g. GC-rich regions)

4) Even with fast computers, it's still computationally difficult

5) Since you don't know what the 'answer' is, it can be difficult to assess whether your assembly is 'good' or not

6) Polyploidy means you are effectively assembling >1 closely related, but not identical, genome

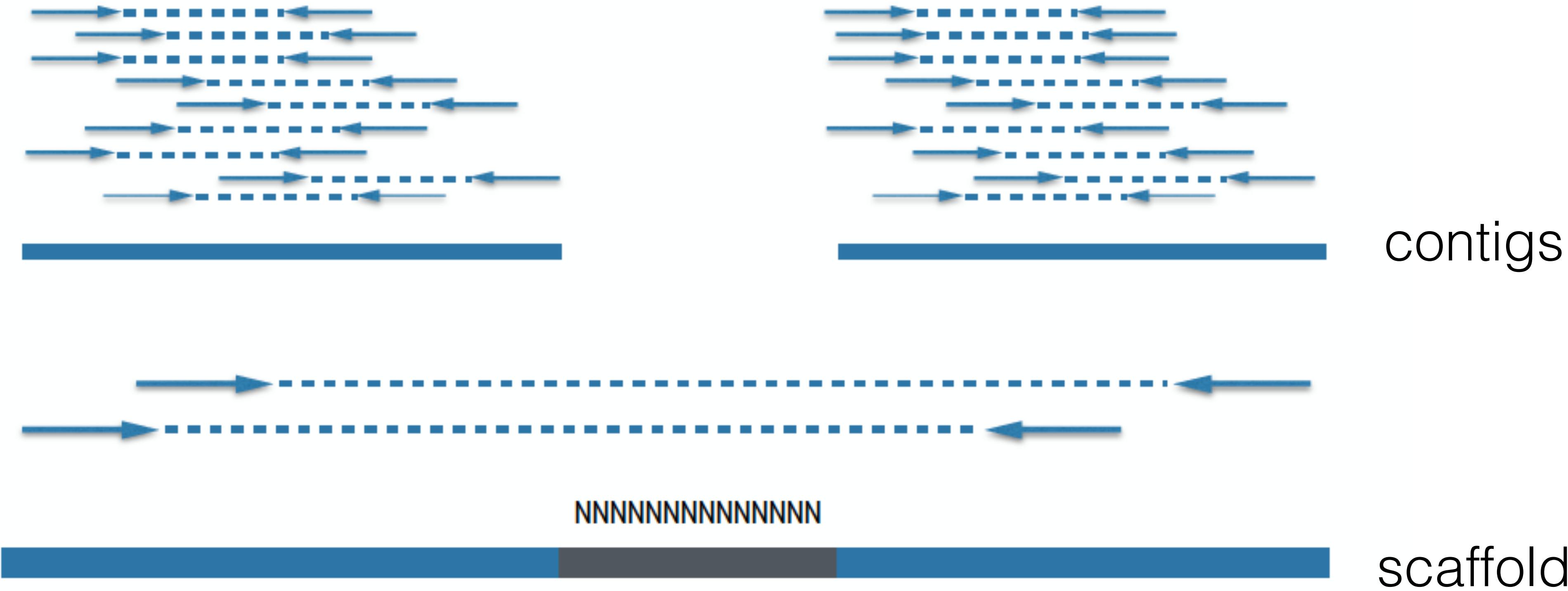
7) Not to mention annotation, which can be as hard as assembly!

De novo assembly is like doing a jigsaw puzzle without the picture on the box



Images, metaphor: *Keith Bradnam, UC Davis*

Reads are assembled into **contigs**, contigs into **scaffolds**,  
and scaffolds into chromosomes or genomes



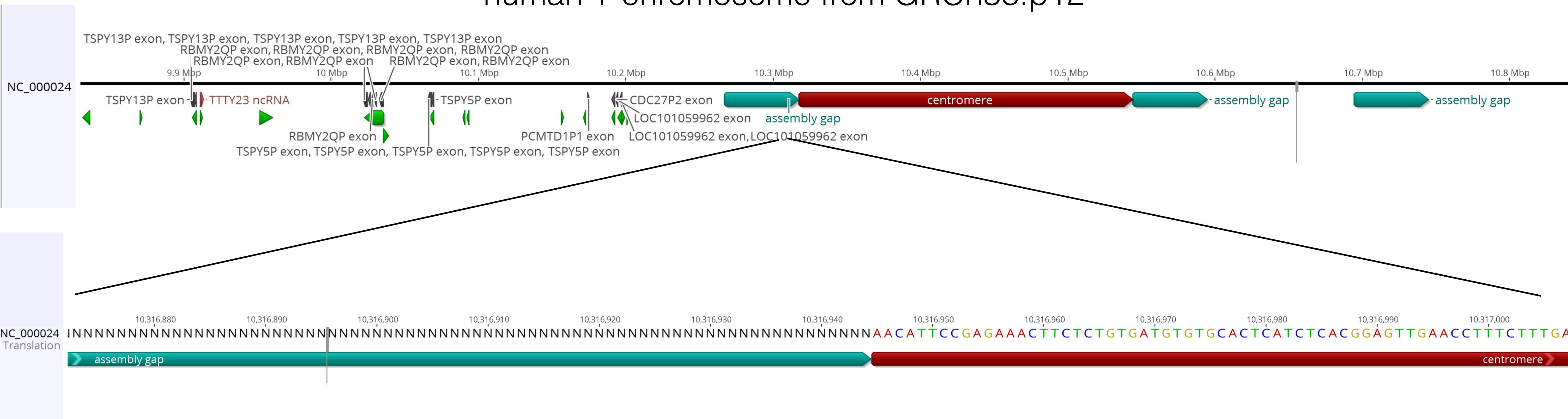


These “contigs” could be scaffolded because we have additional information.

(We know that the two halves of the golden gate bridge should go together)

Sometimes even ‘complete’ assemblies contain gaps

# human Y chromosome from GRCh38.p12



# A truly complete (?) human genome

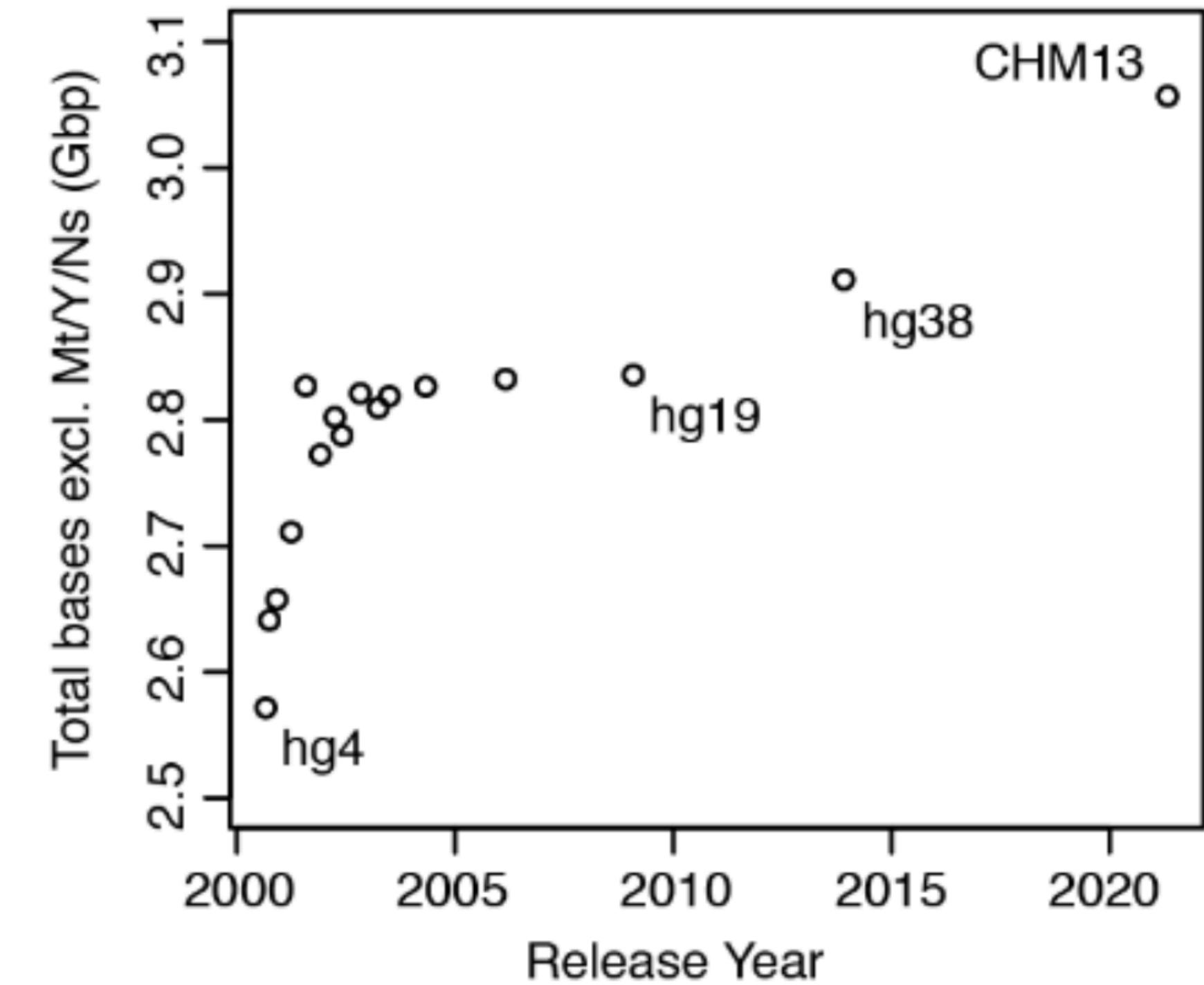
## RESEARCH ARTICLE

### HUMAN GENOMICS

## The complete sequence of a human genome

Sergey Nurk<sup>1†</sup>, Sergey Koren<sup>1†</sup>, Arang Rhie<sup>1†</sup>, Mikko Rautiainen<sup>1†</sup>, Andrey V. Bzikadze<sup>2</sup>, Alla Mikheenko<sup>3</sup>, Mitchell R. Vollger<sup>4</sup>, Nicolas Altemose<sup>5</sup>, Lev Uralsky<sup>6,7</sup>, Ariel Gershman<sup>8</sup>, Sergey Aganezov<sup>9‡</sup>, Savannah J. Hoyt<sup>10</sup>, Mark Diekhans<sup>11</sup>, Glennis A. Logsdon<sup>4</sup>, Michael Alonge<sup>9</sup>, Stylianos E. Antonarakis<sup>12</sup>, Matthew Borchers<sup>13</sup>, Gerard G. Bouffard<sup>14</sup>, Shelise Y. Brooks<sup>14</sup>, Gina V. Caldas<sup>15</sup>, Nae-Chyun Chen<sup>9</sup>, Haoyu Cheng<sup>16,17</sup>, Chen-Shan Chin<sup>18</sup>, William Chow<sup>19</sup>, Leonardo G. de Lima<sup>13</sup>, Philip C. Dishuck<sup>4</sup>, Richard Durbin<sup>19,20</sup>, Tatiana Dvorkina<sup>3</sup>, Ian T. Fiddes<sup>21</sup>, Giulio Formenti<sup>22,23</sup>, Robert S. Fulton<sup>24</sup>, Arkarachai Fungtammasan<sup>18</sup>, Erik Garrison<sup>11,25</sup>, Patrick G. S. Grady<sup>10</sup>, Tina A. Graves-Lindsay<sup>26</sup>, Ira M. Hall<sup>27</sup>, Nancy F. Hansen<sup>28</sup>, Gabrielle A. Hartley<sup>10</sup>, Marina Haukness<sup>11</sup>, Kerstin Howe<sup>19</sup>, Michael W. Hunkapiller<sup>29</sup>, Chirag Jain<sup>1,30</sup>, Miten Jain<sup>11</sup>, Erich D. Jarvis<sup>22,23</sup>, Peter Kerpeljiev<sup>31</sup>, Melanie Kirsche<sup>9</sup>, Mikhail Kolmogorov<sup>32</sup>, Jonas Korlach<sup>29</sup>, Milinn Kremitzki<sup>26</sup>, Heng Li<sup>16,17</sup>, Valerie V. Maduro<sup>33</sup>, Tobias Marschall<sup>34</sup>, Ann M. McCartney<sup>1</sup>, Jennifer McDaniel<sup>35</sup>, Danny E. Miller<sup>4,36</sup>, James C. Mullikin<sup>14,28</sup>, Eugene W. Myers<sup>37</sup>, Nathan D. Olson<sup>35</sup>, Benedict Paten<sup>11</sup>, Paul Peluso<sup>29</sup>, Pavel A. Pevzner<sup>32</sup>, David Porubsky<sup>4</sup>, Tamara Potapova<sup>13</sup>, Evgeny I. Rogaev<sup>6,7,38,39</sup>, Jeffrey A. Rosenfeld<sup>40</sup>, Steven L. Salzberg<sup>9,41</sup>, Valerie A. Schneider<sup>42</sup>, Fritz J. Sedlazeck<sup>43</sup>, Kishwar Shafin<sup>11</sup>, Colin J. Shew<sup>44</sup>, Alaina Shumate<sup>41</sup>, Ying Sims<sup>19</sup>, Arian F. A. Smit<sup>45</sup>, Daniela C. Soto<sup>44</sup>, Ivan Sovic<sup>29,46</sup>, Jessica M. Storer<sup>45</sup>, Aaron Streets<sup>5,47</sup>, Beth A. Sullivan<sup>48</sup>, Françoise Thibaud-Nissen<sup>42</sup>, James Torrance<sup>19</sup>, Justin Wagner<sup>35</sup>, Brian P. Walenz<sup>1</sup>, Aaron Wenger<sup>29</sup>, Jonathan M. D. Wood<sup>19</sup>, Chunlin Xiao<sup>42</sup>, Stephanie M. Yan<sup>49</sup>, Alice C. Young<sup>14</sup>, Samantha Zarate<sup>9</sup>, Urvashi Surti<sup>50</sup>, Rajiv C. McCoy<sup>49</sup>, Megan Y. Dennis<sup>44</sup>, Ivan A. Alexandrov<sup>3,7,51</sup>, Jennifer L. Gerton<sup>13,52</sup>, Rachel J. O'Neill<sup>10</sup>, Winston Timp<sup>8,41</sup>, Justin M. Zook<sup>35</sup>, Michael C. Schatz<sup>9,49</sup>, Evan E. Eichler<sup>4,53\*</sup>, Karen H. Miga<sup>11,54\*</sup>, Adam M. Phillippy<sup>1\*</sup>

Since its initial release in 2000, the human reference genome has covered only the euchromatic fraction of the genome, leaving important heterochromatic regions unfinished. Addressing the remaining 8% of the genome, the Telomere-to-Telomere (T2T) Consortium presents a complete 3.055 billion-base pair sequence



This “telomere-to-telomere” assembly required a ton of data and many different types of data

None of this was standard Illumina short read data

The screenshot shows a GitHub repository page for 'marbl/CHM13'. The URL in the address bar is <https://github.com/marbl/CHM13>. The page displays the contents of the 'README.md' file. The first section is titled 'Introduction'.

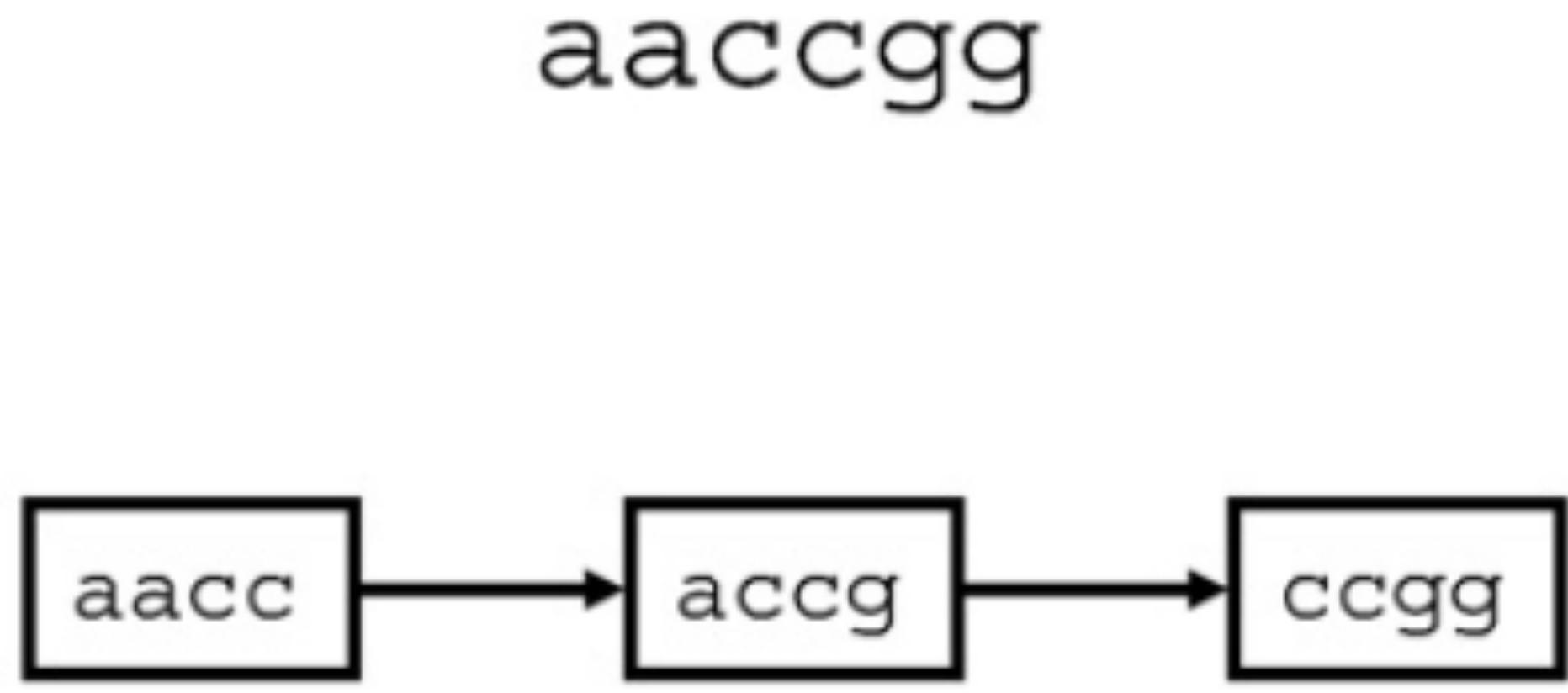
**Introduction**

---

We have sequenced the CHM13hTERT human cell line with a number of technologies. Human genomic DNA was extracted from the cultured cell line. As the DNA is native, modified bases will be preserved. The data includes 30x PacBio HiFi, 120x coverage of Oxford Nanopore, 70x PacBio CLR, 50x 10X Genomics, as well as BioNano DLS and Arima Genomics HiC. Most raw data is available from this site, with the exception of the PacBio data which was generated by the University of Washington/PacBio and is available from NCBI SRA.

Nearly all short read assemblers use a de Bruijn graph-based algorithm

De bruijn graphs are directed graphs with connected nodes of overlapping k-mers



Generic simplified strategy:

- Attempted error correction
- Break reads into overlapping k-mers (here  $k = 4$ )
- Construct de Bruijn graph of k-mers
- Trace path through graph:  
**Tada! Genome sequence**

kmers are just sequences of length k



Ryan Wick  
@rrwick



Bioinformatics joke:

What do you call a tetranucleotide that's had a base added or removed?

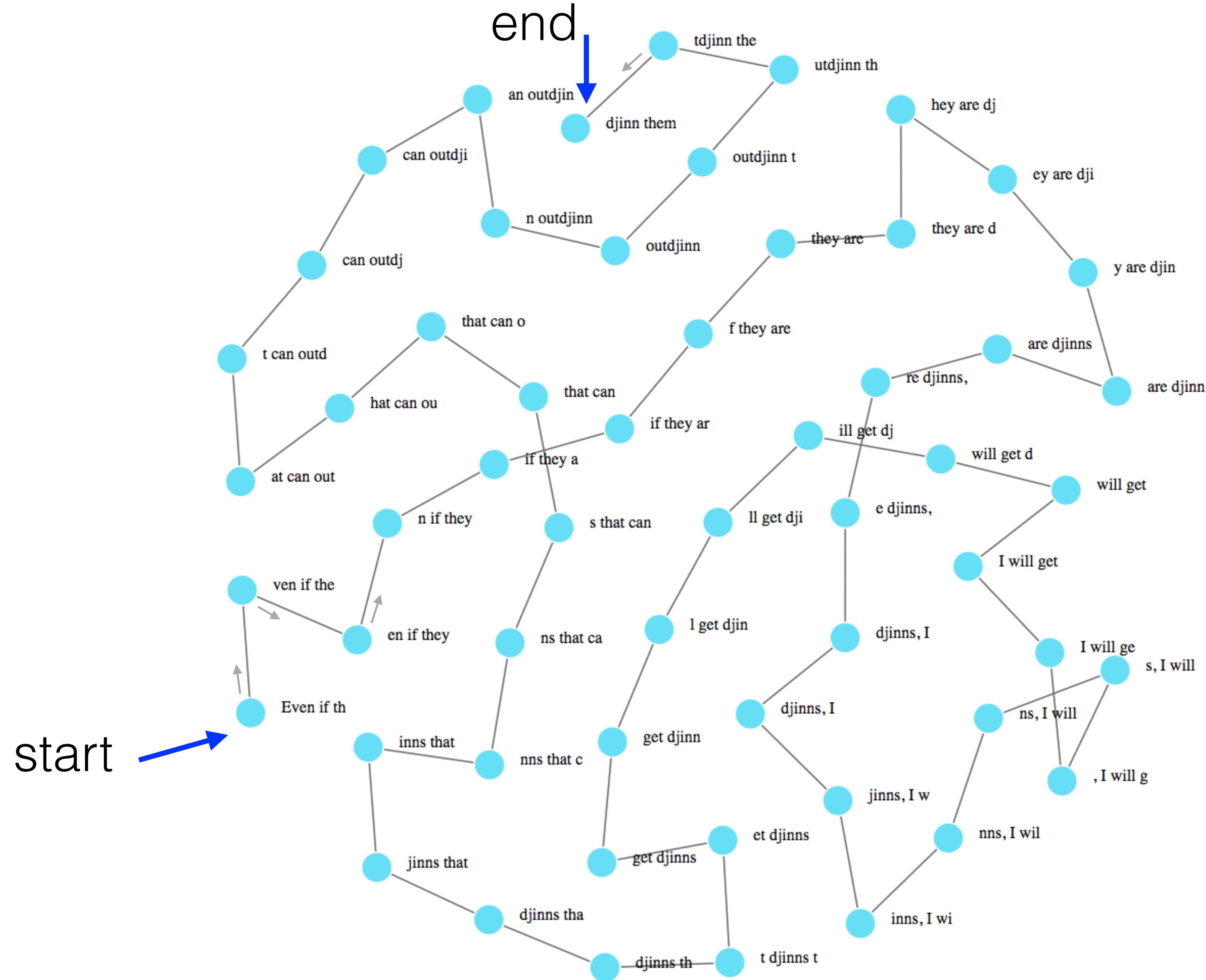
A former 4-mer.

10:52 PM · Jun 24, 2020 · [Twitter Web App](#)

---

19 Retweets 109 Likes

k=10

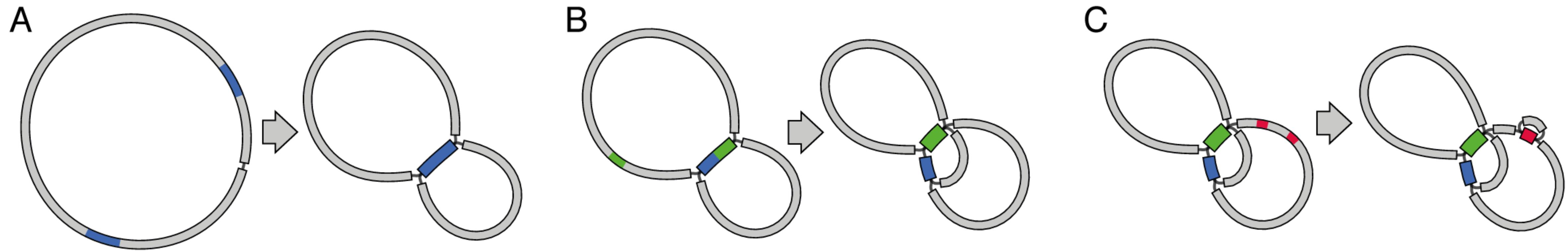


**k=8**

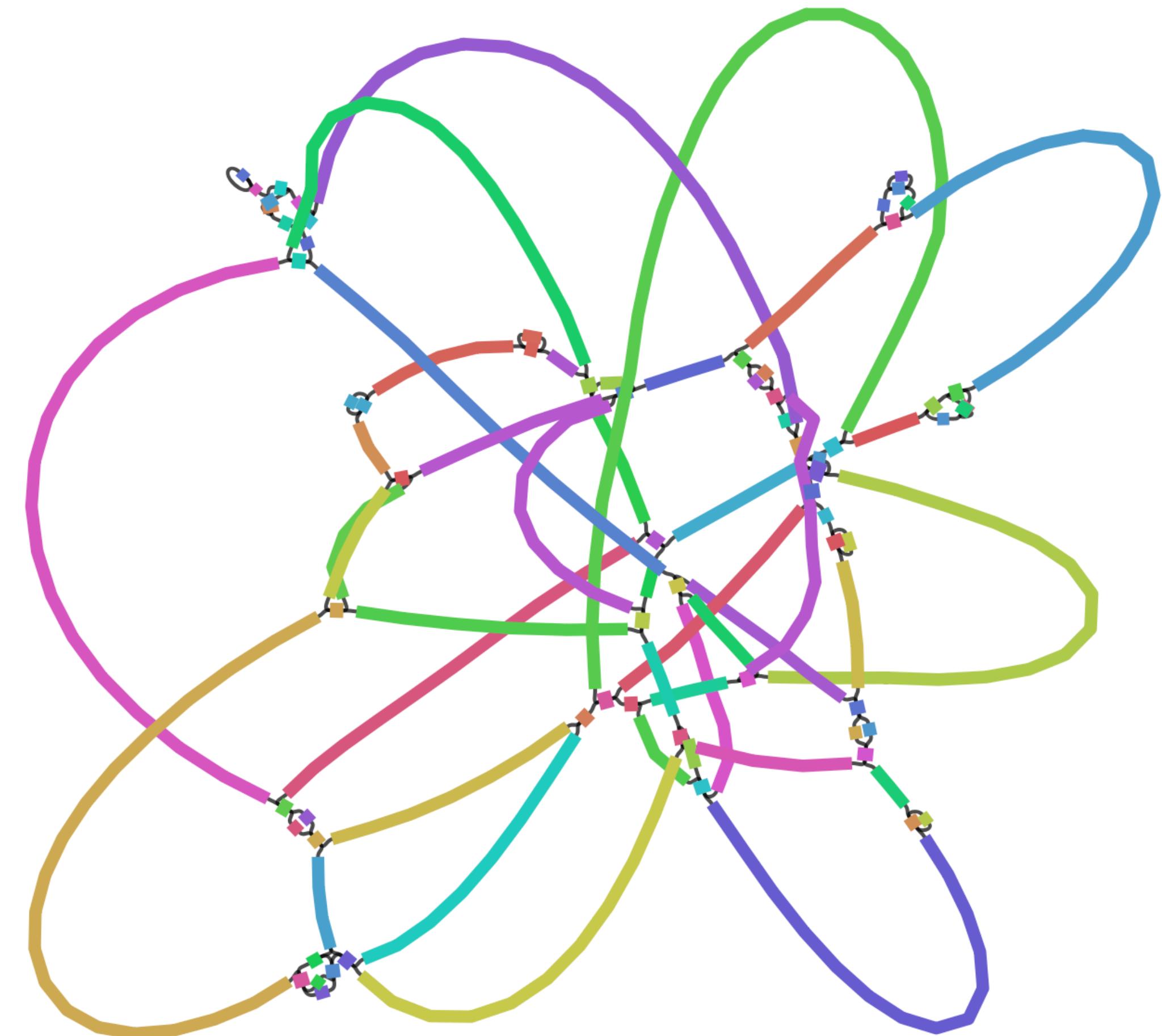
<http://debruijn.herokuapp.com/graph>

# The impact of additional repeats on graph complexity

Even bacterial genomes have repeated sequences (e.g. rRNA loci, duplicated genes)

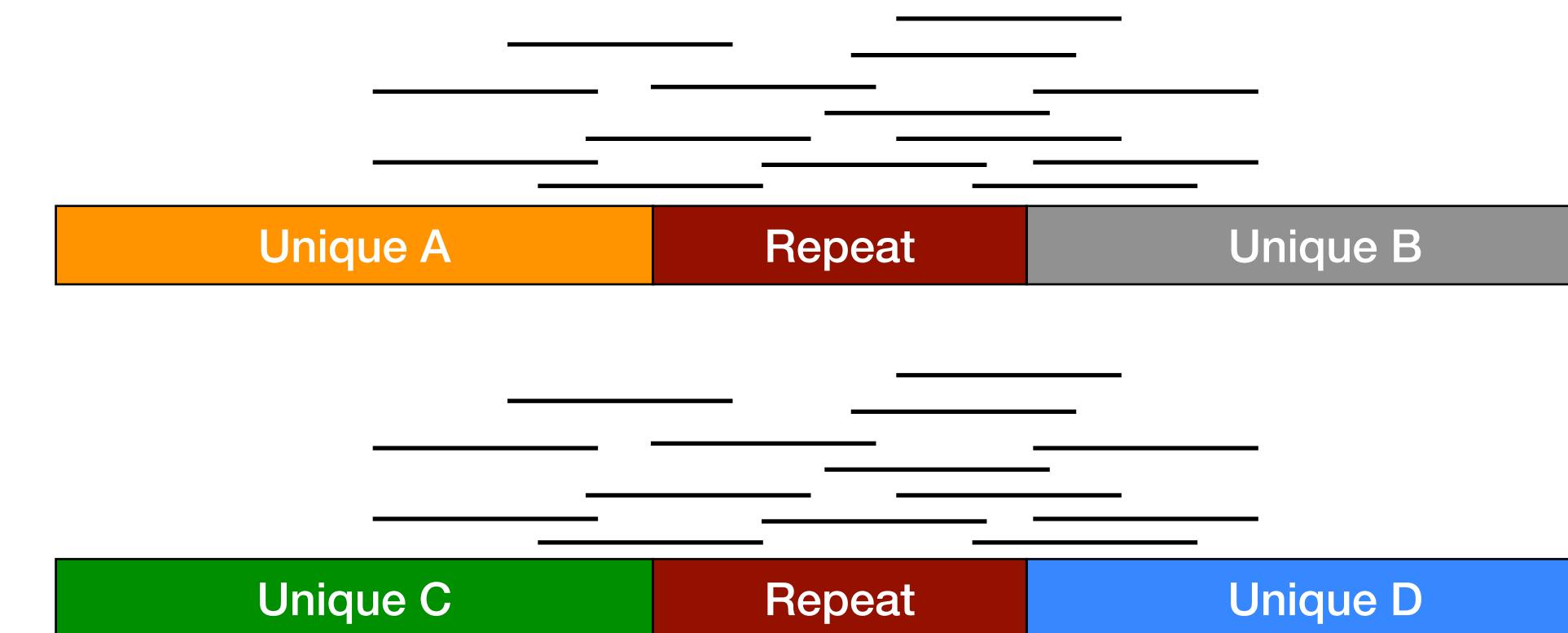


A real life graph from Illumina only data: A *Pseudomonas aeruginosa* isolate  
Illumina 2x250 paired end data with ~100x average depth of coverage

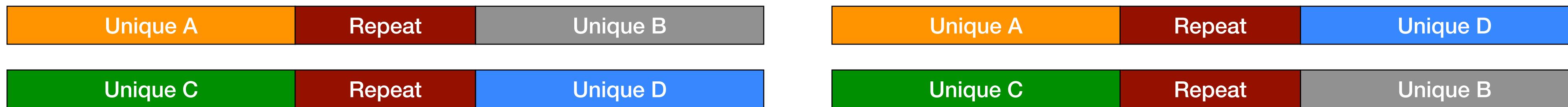


# Short reads alone do not contain enough information to resolve repeats

No read with read length < repeat length can bridge the repeat to link unique sequence on each side

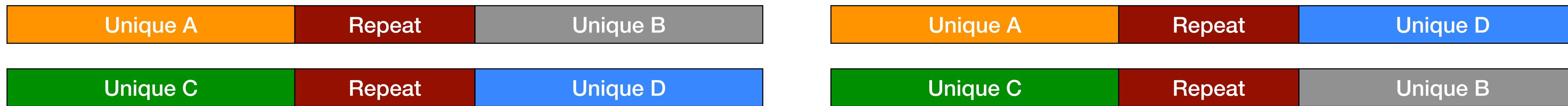


2 equally plausible combinations given short read data

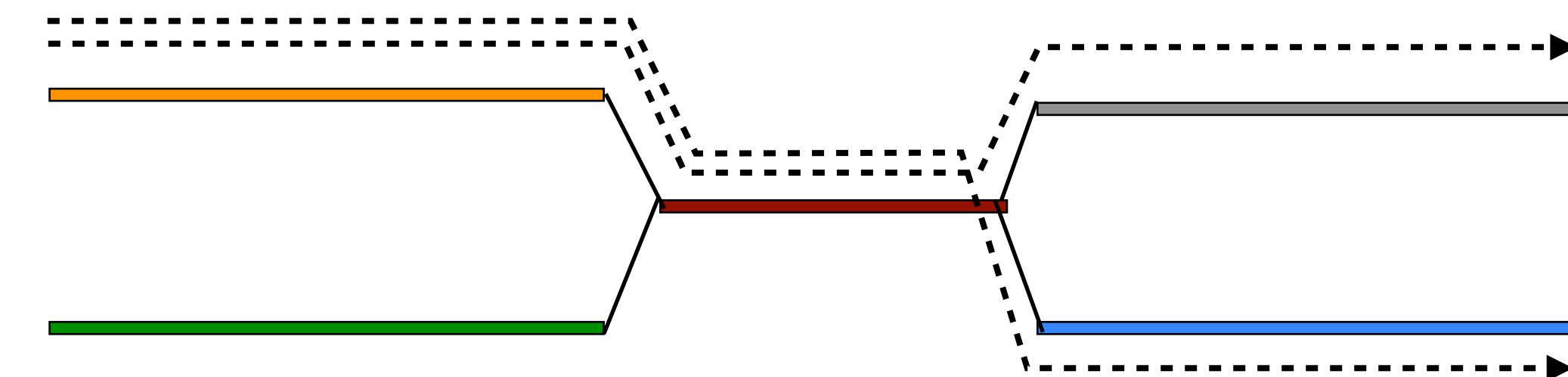


# Short reads alone do not contain enough information to resolve repeats

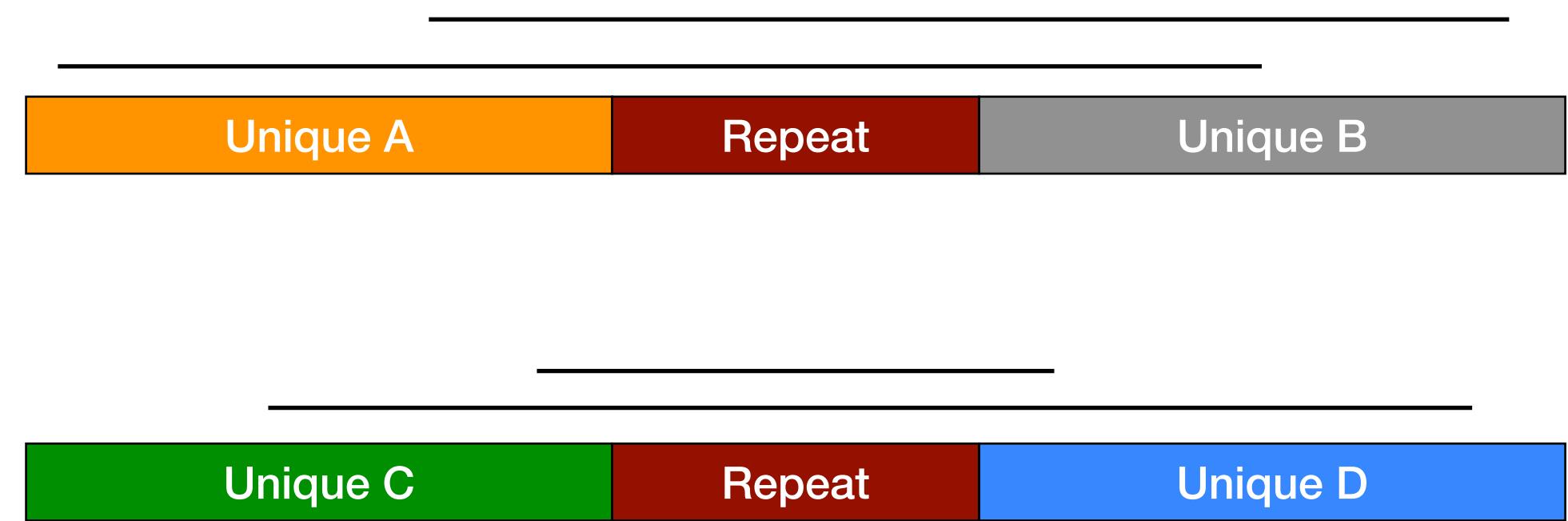
2 equally plausible combinations given short read data



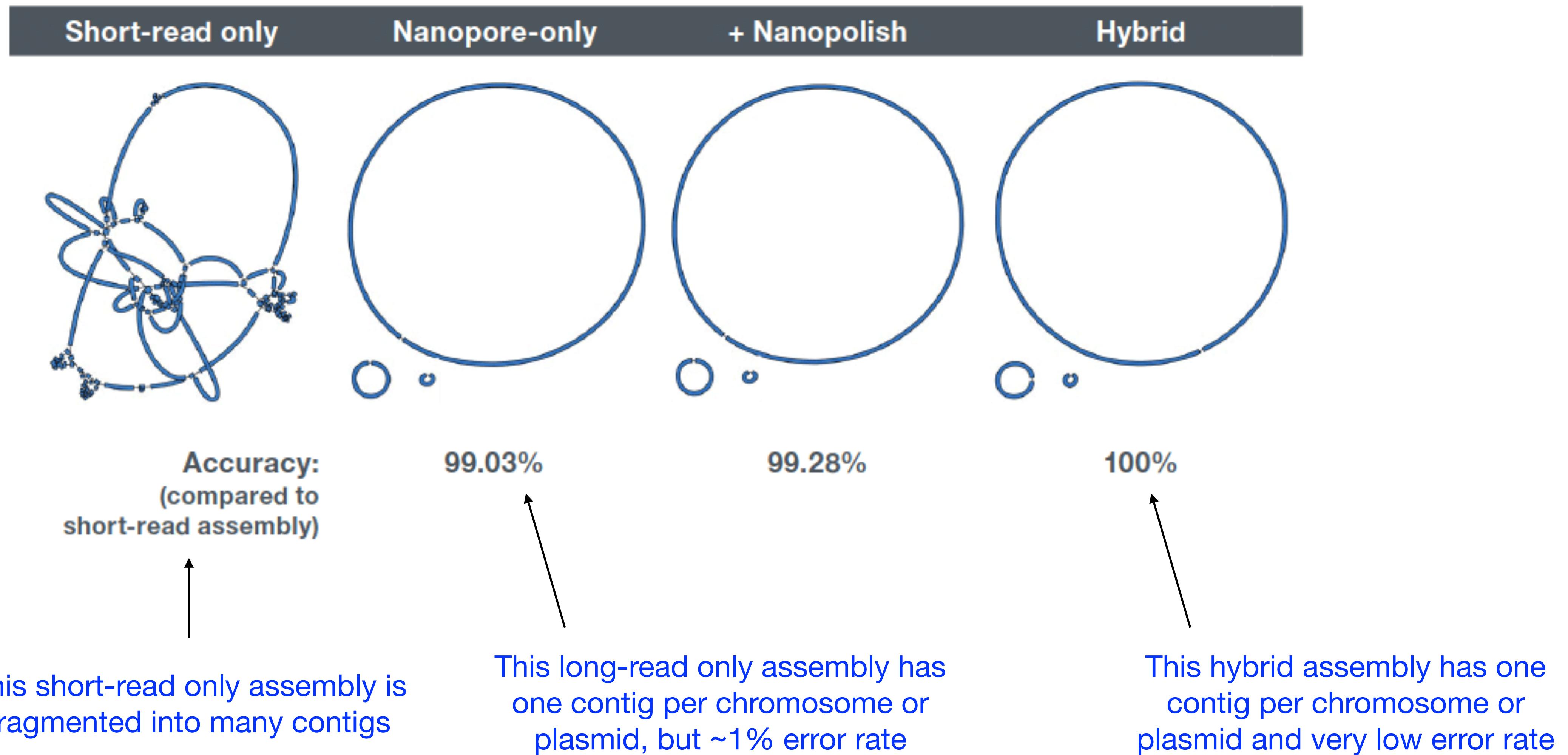
This corresponds to 2 equally plausible paths through the graph



Long reads span repeats to link unique sequence on each side



# Combining long + short read produces assemblies with high contiguity and low error rate



## How do you know if your assembly is good?

- Size of the assembly: does it match estimates from other means?
- Size of the contigs/scaffolds: are they reasonably long?
- Are the expected ‘core genes’ present in the assembly?
- Does the assembly contain sequences of contaminating organisms?
- Is the assembly consistent with independently derived data? (optical mapping, transcriptome sequencing, genomes of related organisms?)

For what purpose do you need the assembly?

These questions apply to assemblies in databases too.

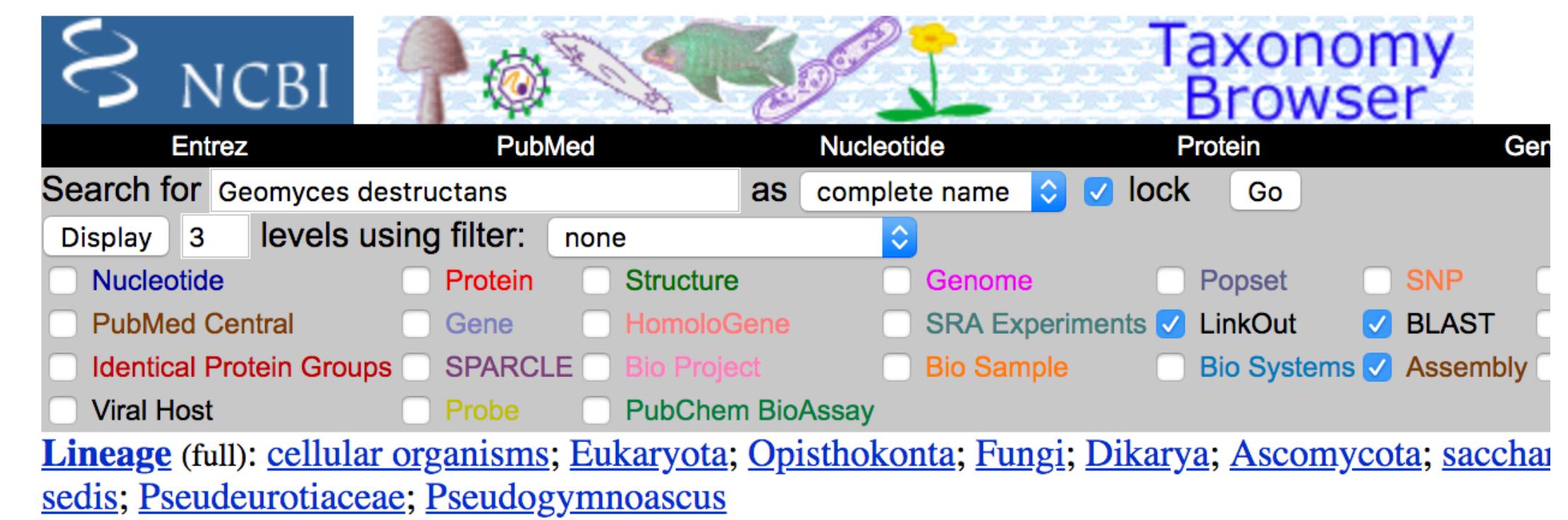
# Mini exercise

*Pseudogymnoascus destructans*  
cause of white nose syndrome



image: Marvin Moriarty/USFWS

Visit the pages for the 3 assemblies.  
How were they made? What type of data?  
Is one obviously better? Which would you use?



NCBI Taxonomy Browser

Search for Geomyces destructans as complete name lock Go

Display 3 levels using filter: none

Nucleotide Protein Structure Genome Popset SNP  
PubMed Central Gene HomoloGene SRA Experiments LinkOut BLAST  
Identical Protein Groups SPARCLE Bio Project Bio Sample Bio Systems Assembly  
Viral Host Probe PubChem BioAssay

[Lineage](#) (full): [cellular organisms](#); [Eukaryota](#); [Opisthokonta](#); [Fungi](#); [Dikarya](#); [Ascomycota](#); [saccharis](#); [Pseudeurotiaceae](#); [Pseudogymnoascus](#)

- o [Pseudogymnoascus destructans](#) 3 [LinkOut](#) [BLAST page](#) Click on organism name to get more information
  - [Pseudogymnoascus destructans 20631-21](#) 1 [LinkOut](#)
  - [Pseudogymnoascus destructans M1379](#) 1 [LinkOut](#)

a common assembly metric:

**N50**: a measure of the average size of contigs & scaffolds

# Not all assembly problems are equally difficult!

tiny ssDNA genome

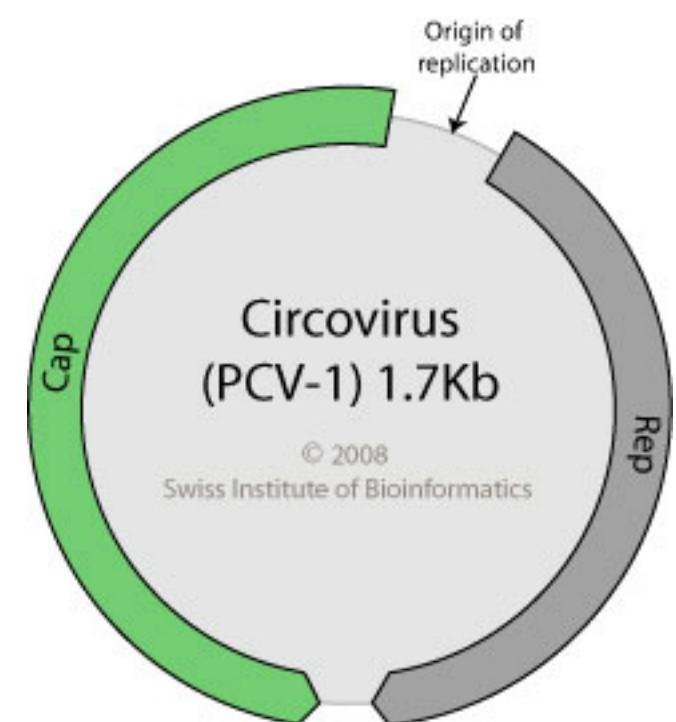
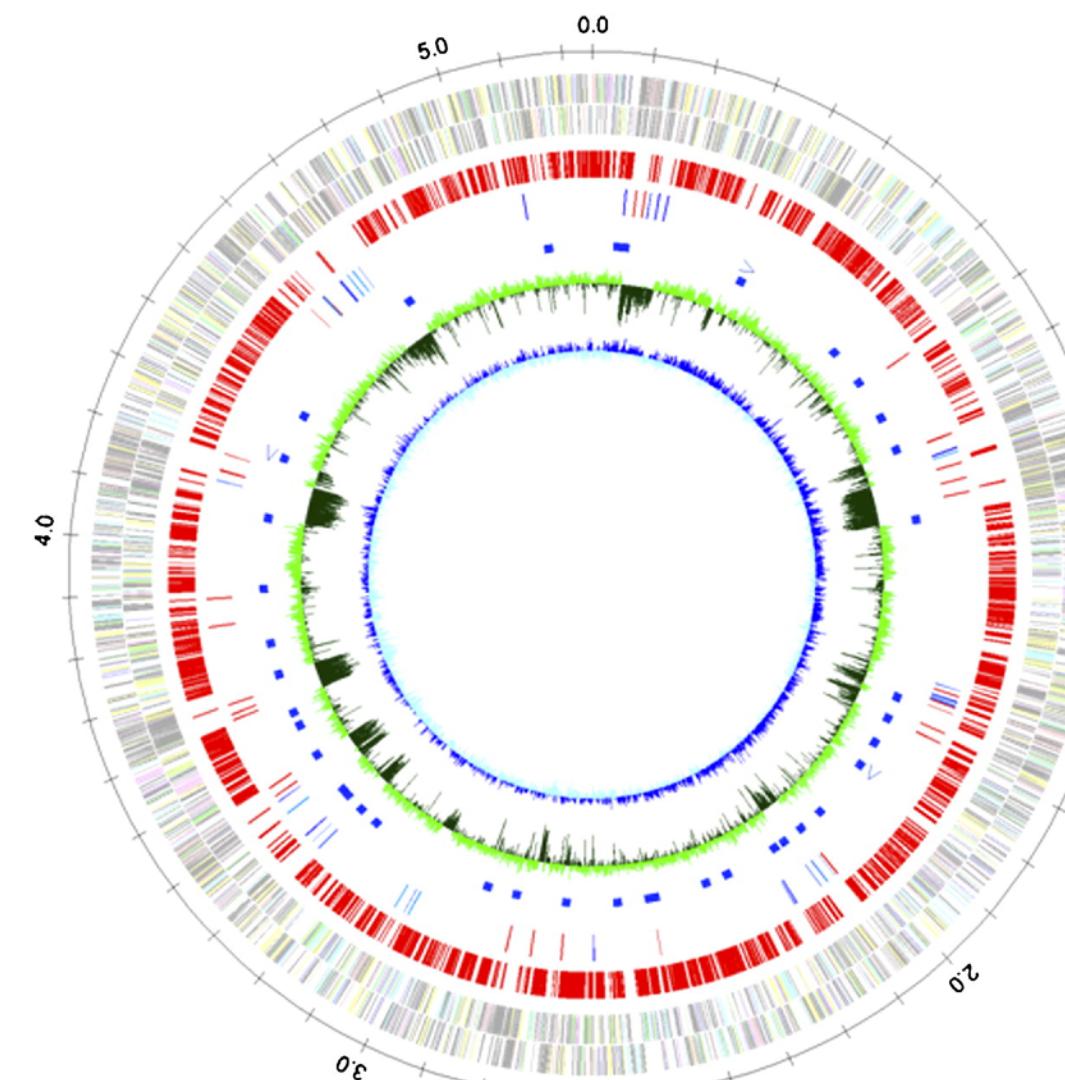


image: viralzone

bacterial genomes ~5 Mbp



Nakazawa et al (2009) Genome Research

Loblolly pine (*Pinus taeda*)  
22 Gbp genome!



image: Univ of Alabama