

# Genome assembly

Mark Stenglein, GDW Workshop



# 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



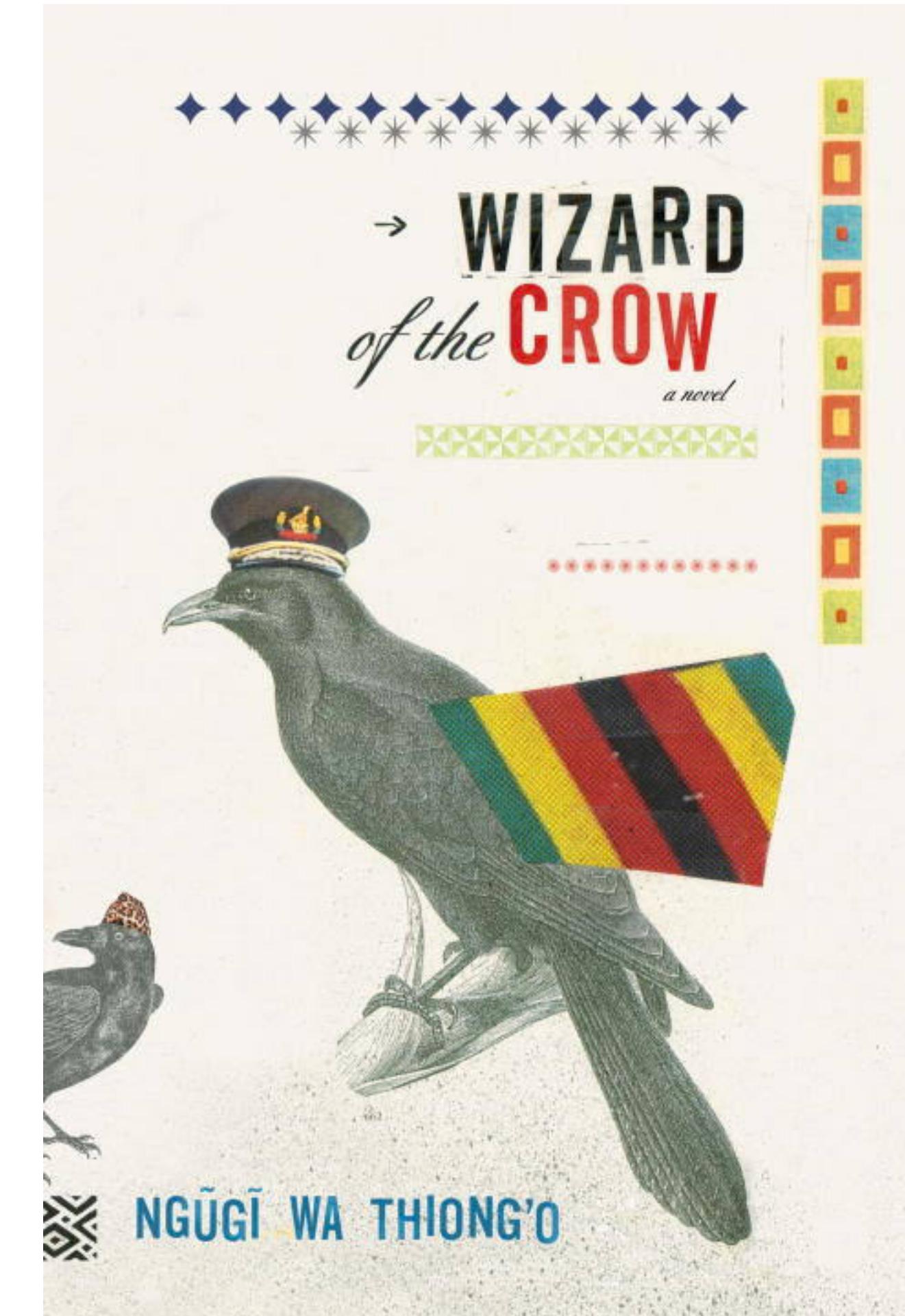
Keith Bradnam, UC Davis

# Genome assembly paper exercise

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

## Rules/info:

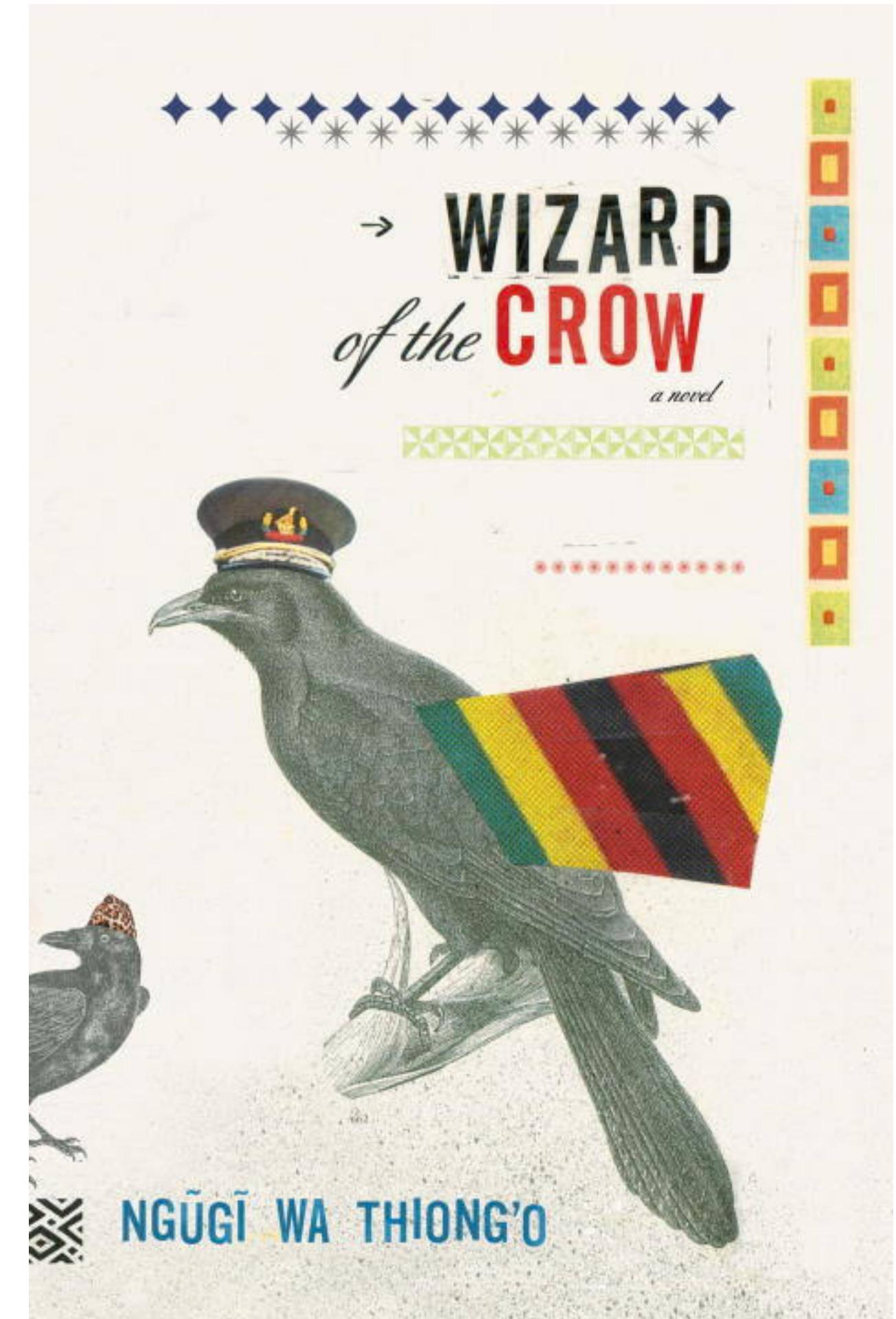
- Like real sequencing data, these reads contain errors.  
The error rate is ~2%
- These are single-end 11-base reads
- The average coverage is ~6x
- You’re not allowed to google the answer
- Also: the answer is in the slides: don’t cheat!
- You can use your computers (i.e. word processors or text editors)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

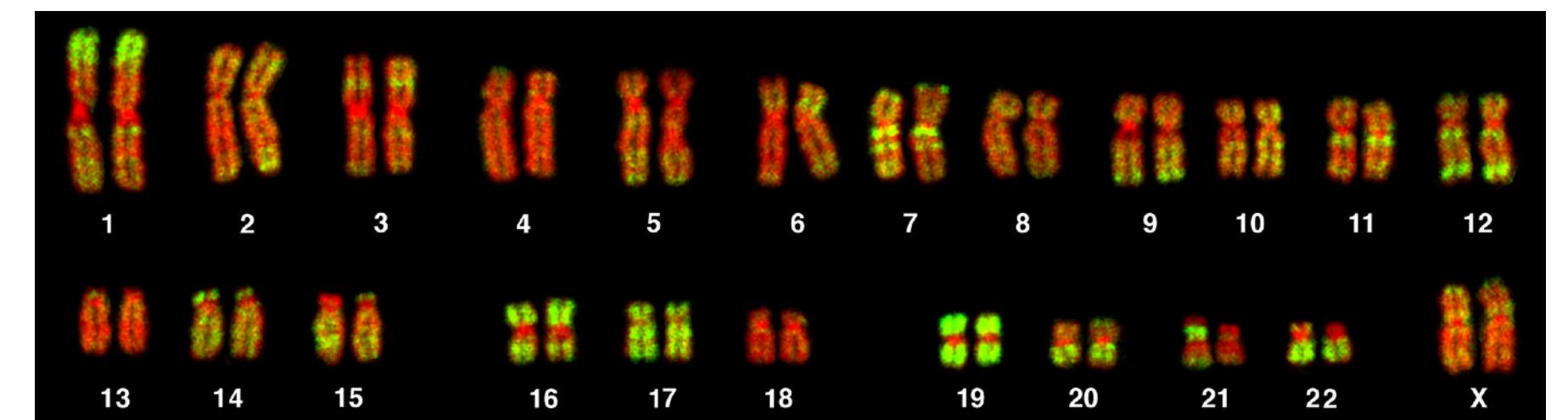


# Conclusion: assembly is not trivial!

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)



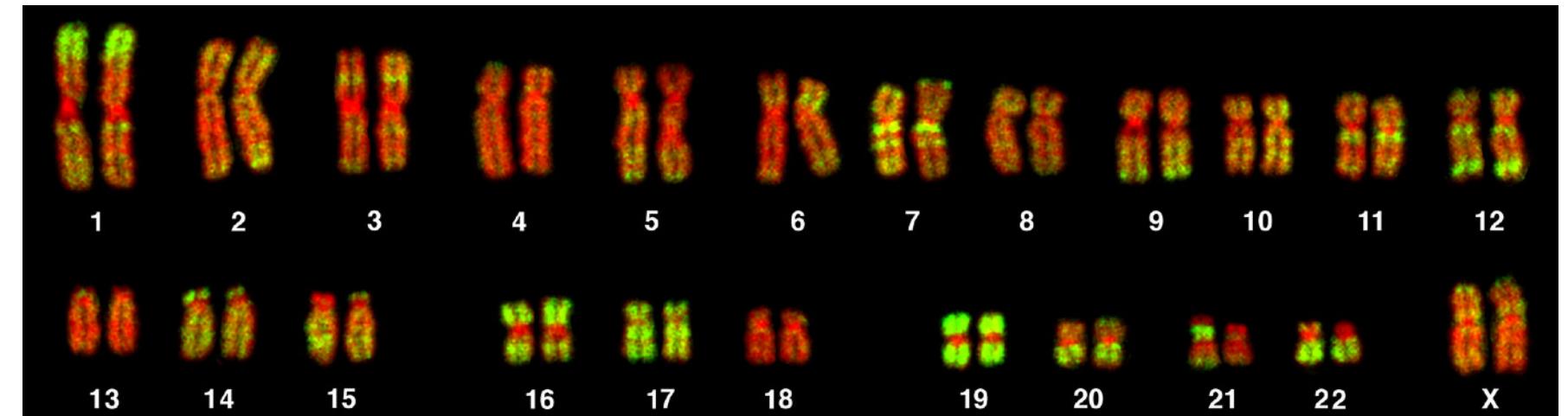
Bolzer et al (2005) PLoS Biol

# 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

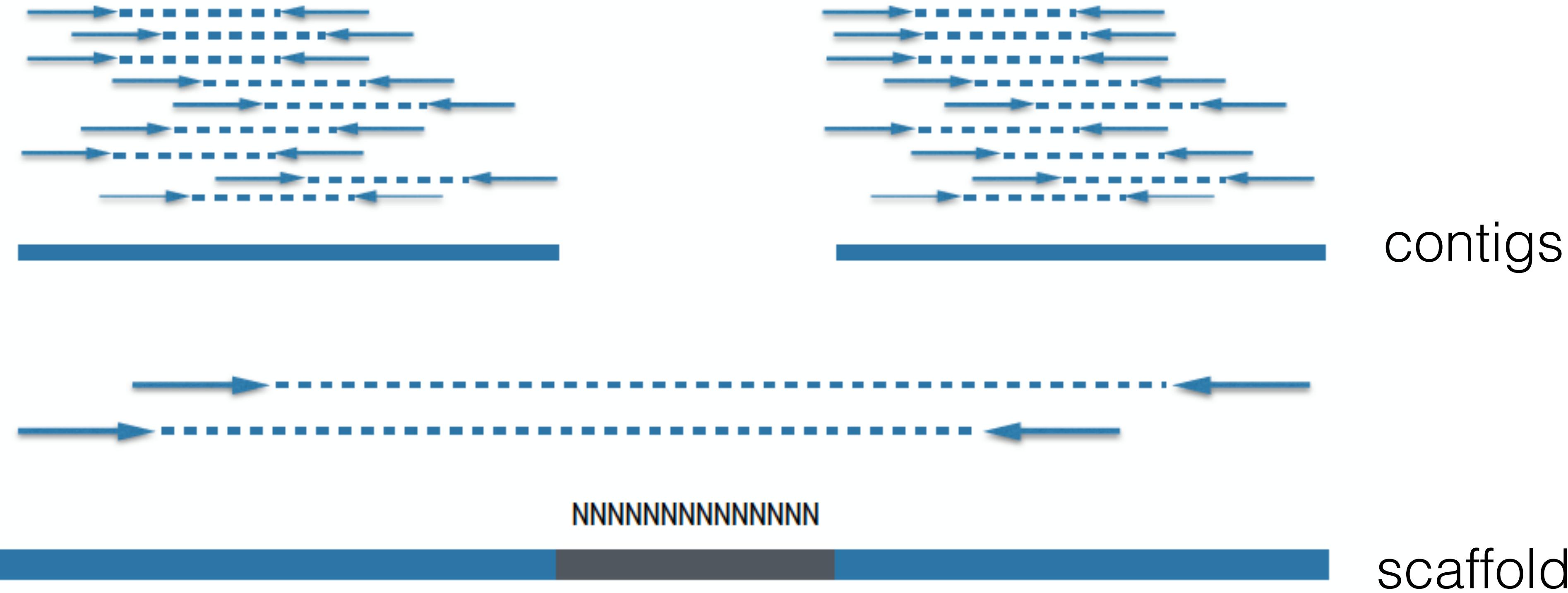
7) Not to mention annotation, which can be as hard as assembly!

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



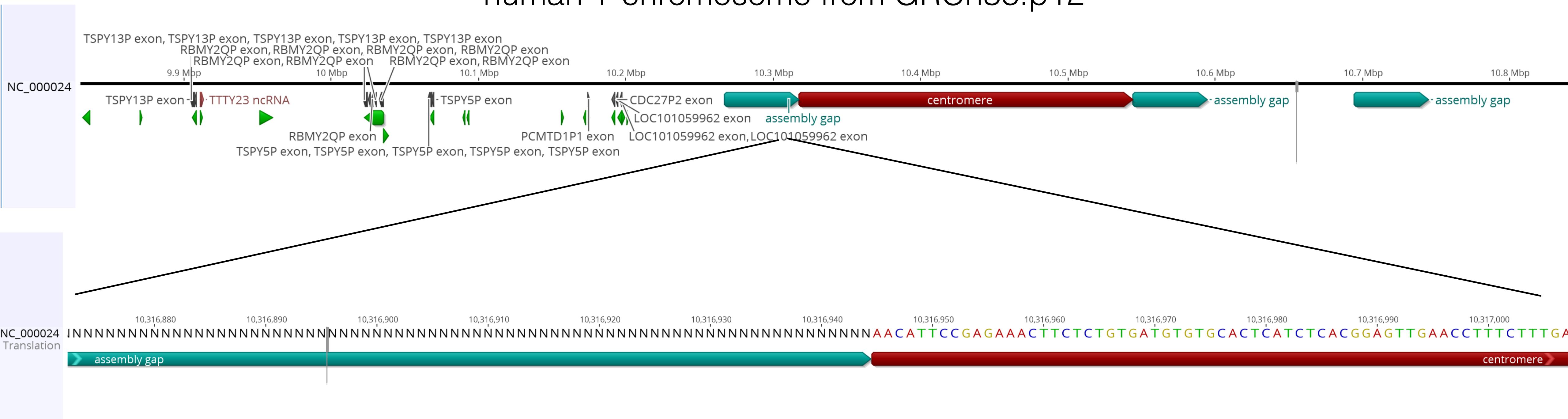


These “contigs” could be scaffolded because we have additional information

(We know about the golden gate bridge)

Sometimes even ‘complete’ assemblies contain gaps

# human Y chromosome from GRCh38.p12



# A truly complete (?) human genome

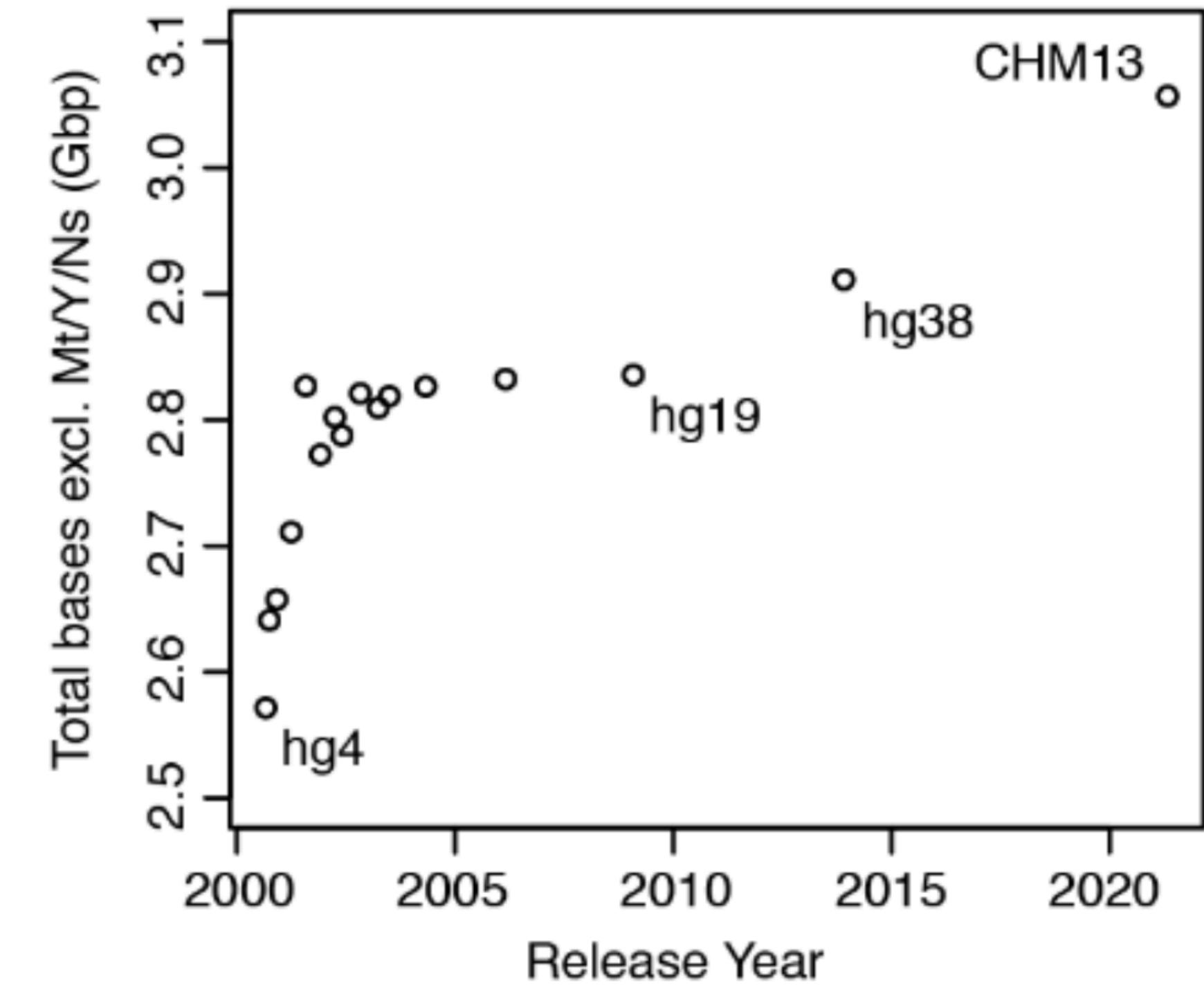
## RESEARCH ARTICLE

### 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'.

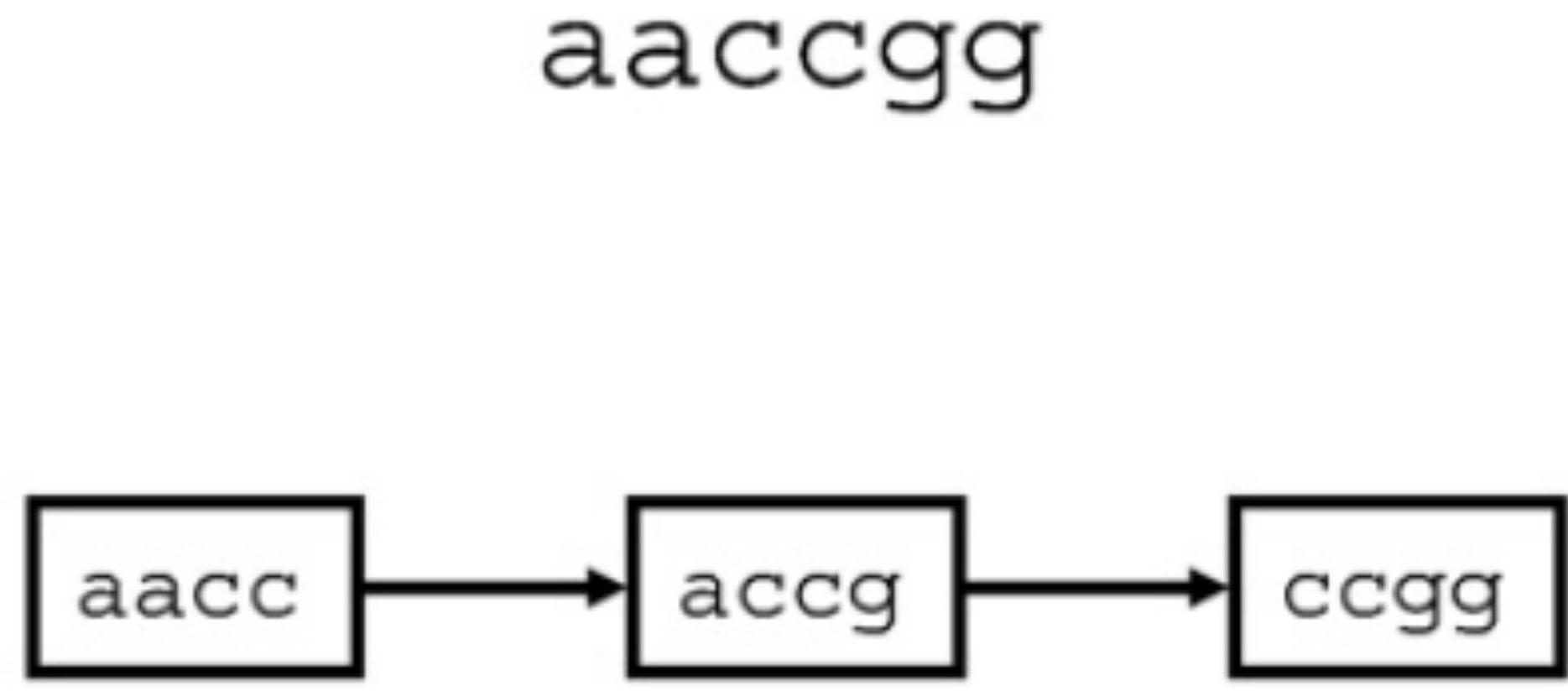
**Introduction**

---

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

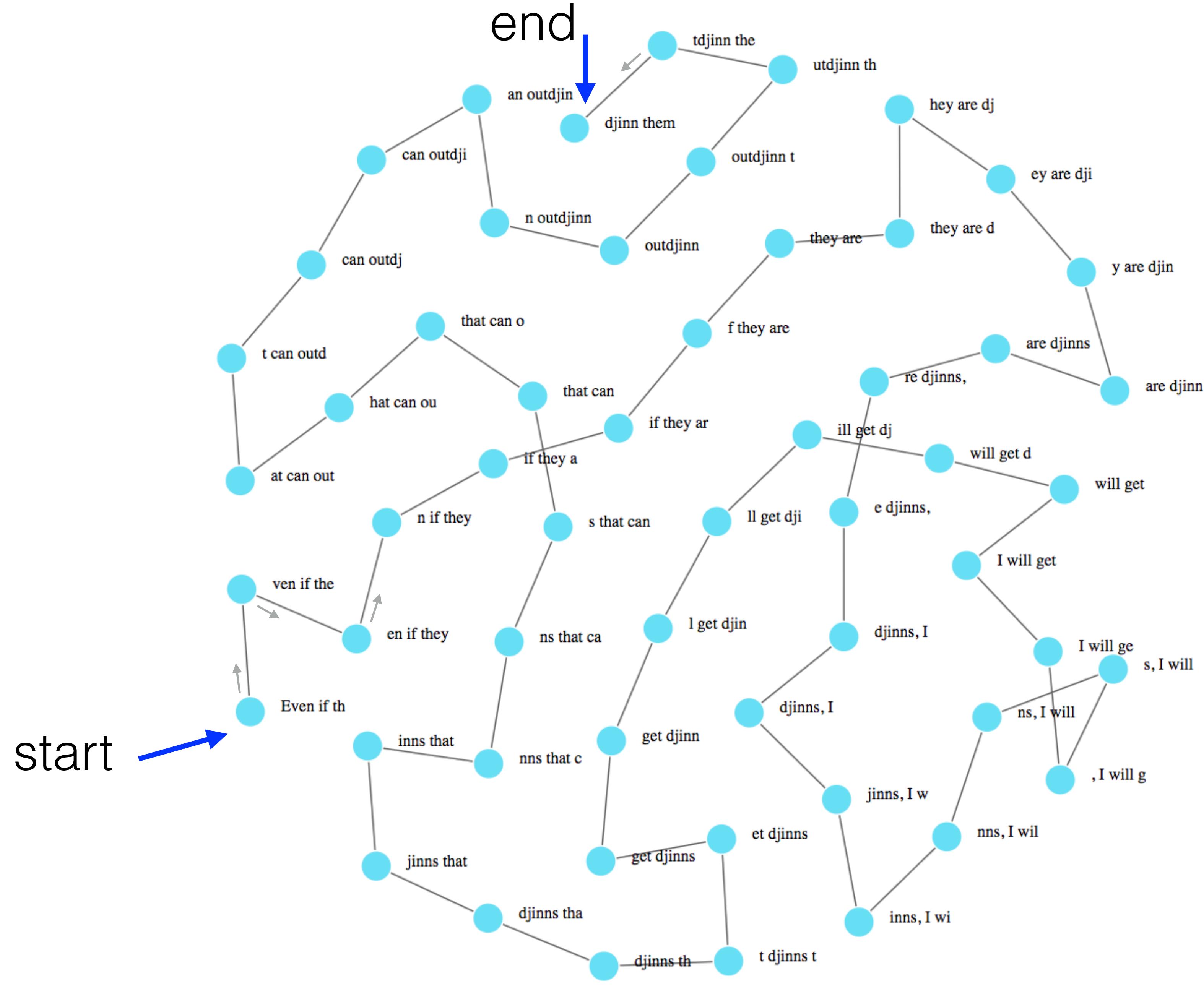
De bruijn graphs are directed graphs with connected nodes of overlapping k-mers



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**

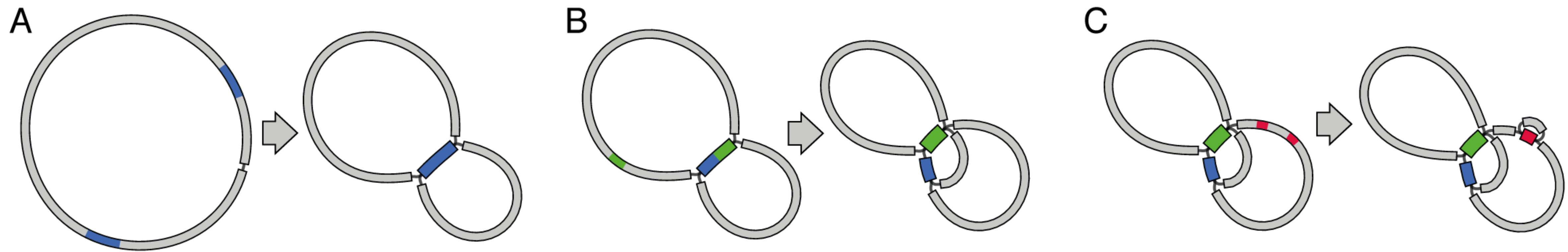
k=10



**k=8**

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

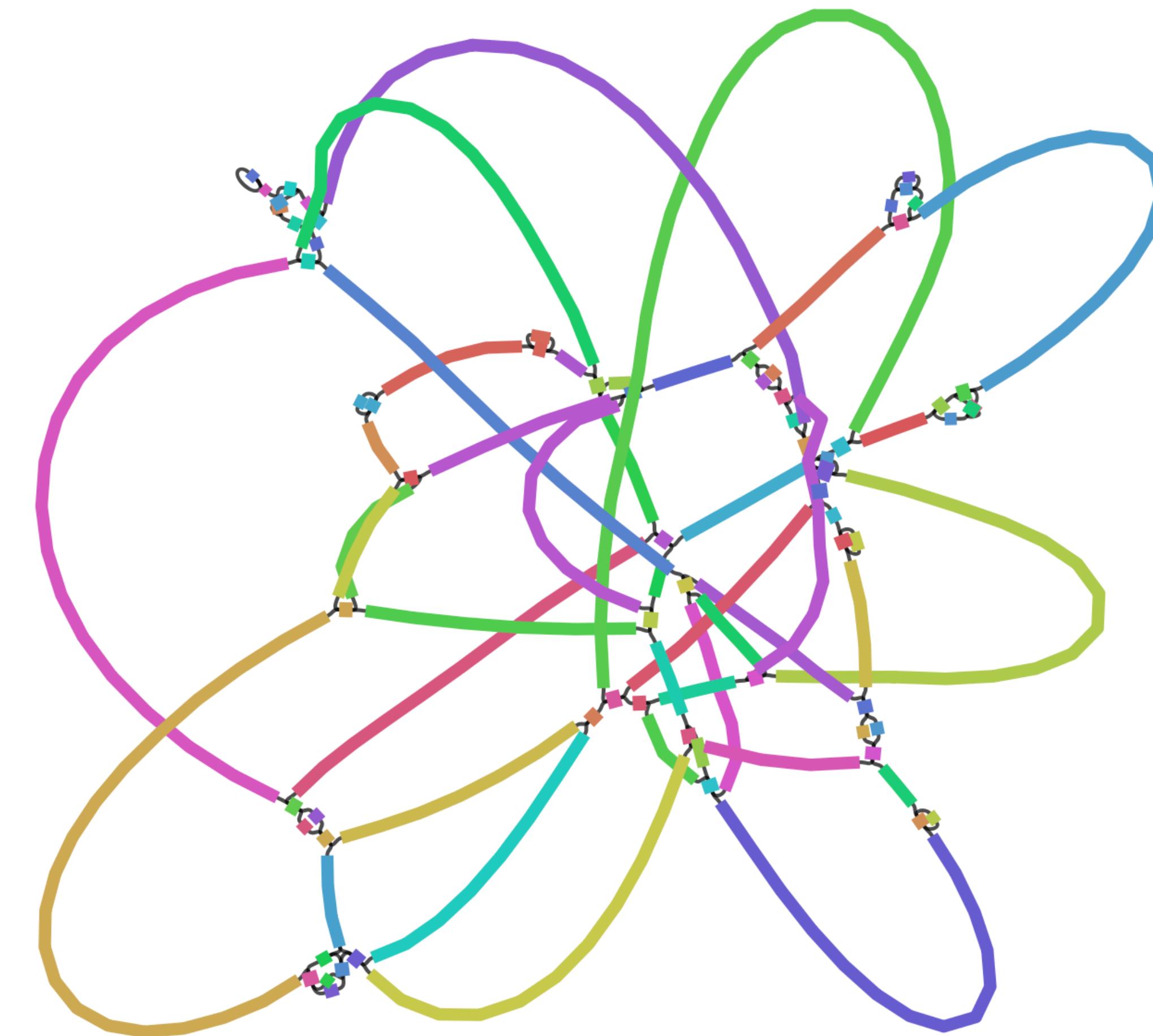
# The impact of additional repeats on graph complexity



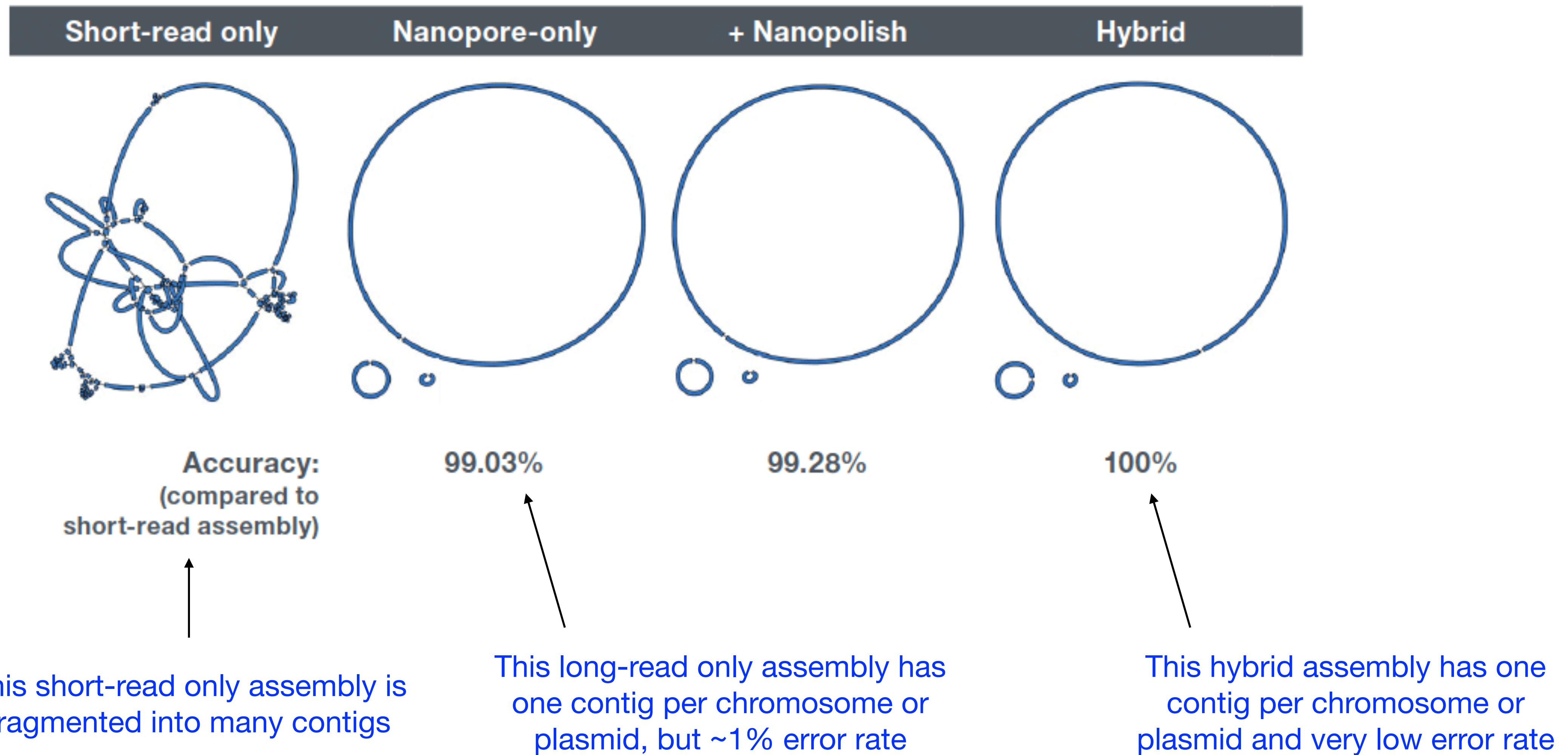
Ryan Wick, Monash University  
<https://github.com/rrwick/Unicycler>

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

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



# 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.

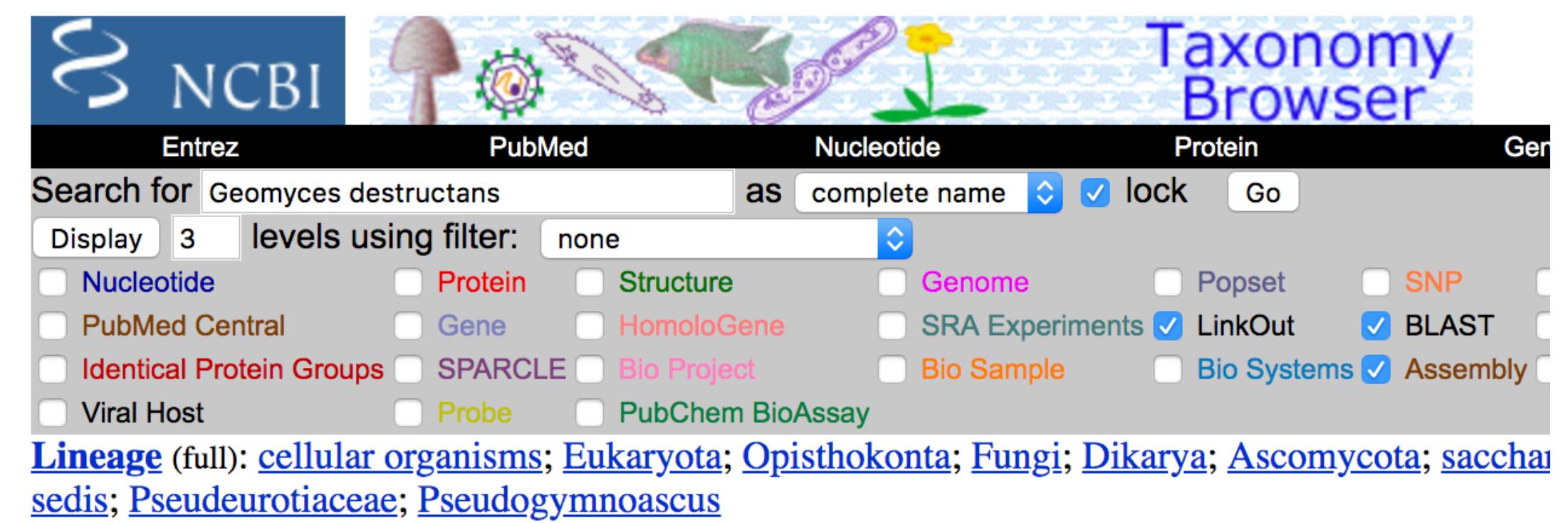
# Assemblies exercise

*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 **Geomycetes 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

I'm painting a bleak picture, but don't be too intimidated:  
genome sequencing and assembly *is* possible.

Not all assembly problems are equally difficult!

tiny ssDNA genome

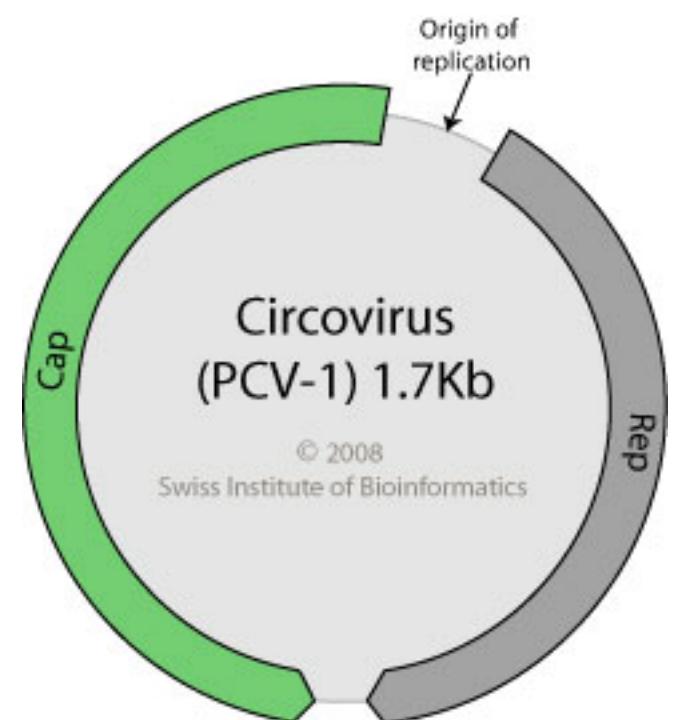
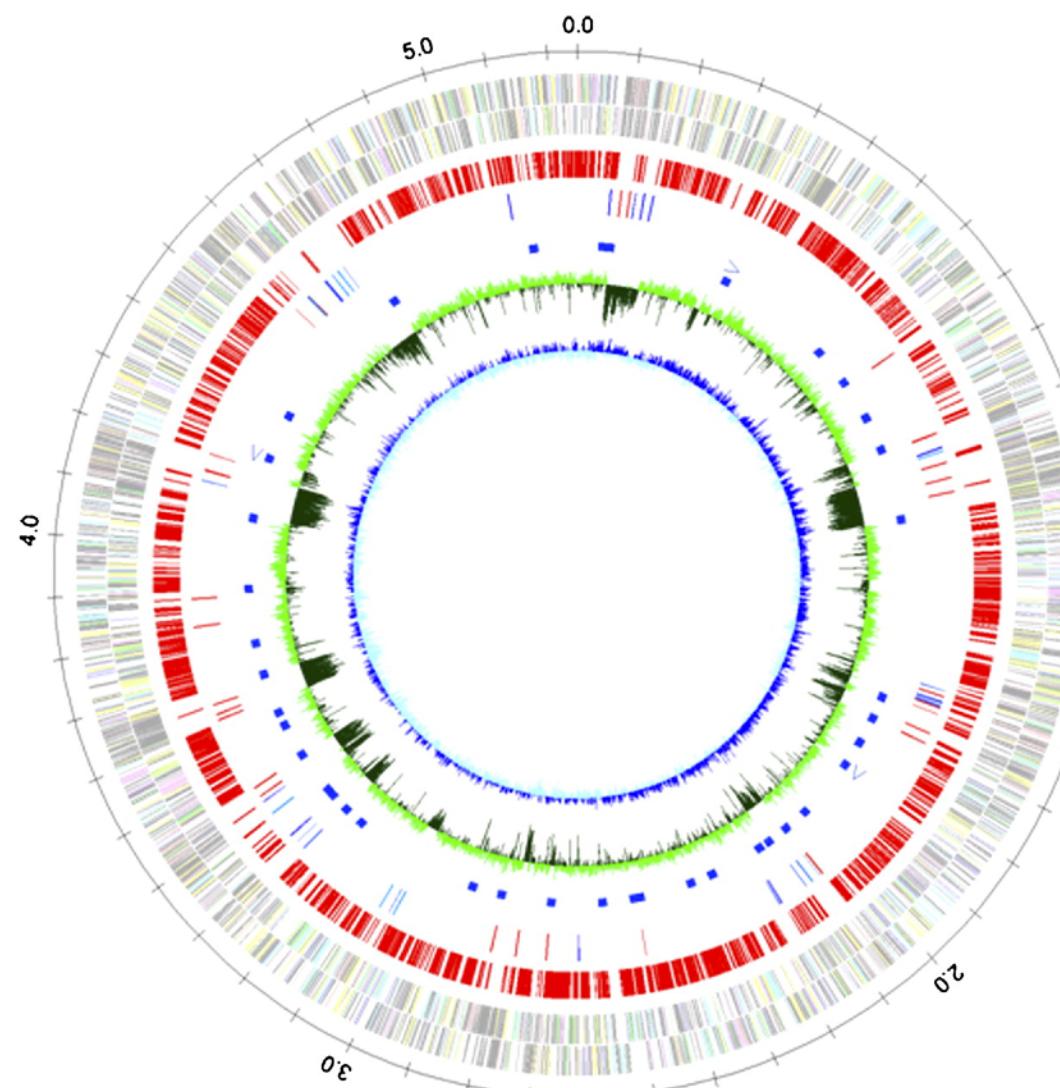


image: viralzone

bacterial genomes ~5 Mbp



Nakazawa et al (2009) Genome Research

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



image: Univ of Alabama

Also: most of the time, you don't need a perfect assembly to answer your biological question

Reading what others have done is a great way to figure out what you could do

MOLECULAR ECOLOGY  
RESOURCES

Molecular Ecology Resources (2016) 16, 314–324

doi: 10.1111/1755-0998.12443

The *de novo* genome assembly and annotation of a female domestic dromedary of North African origin

ROBERT R. FITAK,<sup>\*1</sup> ELMIRA MOHANDESAN,<sup>\*</sup> JUKKA CORANDER<sup>†</sup> and PAMELA A. BURGER<sup>\*</sup>

*\*Institut für Populationsgenetik, Vetmeduni Vienna, Veterinärplatz 1, Vienna 1210, Austria, †Department of Mathematics and Statistics, University of Helsinki, Helsinki, FIN-0014, Finland*

# You could call these ‘bioinformatics protocols’

mina, USA). Preprocessing of the sequence reads included the removal of adapter sequences and removal of reads with >10% uncalled bases and/or >50% of bases with a Phred-scaled quality score <4. After preprocessing, all 100-bp (paired-end) and 50-bp (mate-pair) reads were retained as the set of ‘raw’ reads. We trimmed the 3' end of all raw reads using a modified Mott algorithm in POPOOLATION v1.2.2 (Kofler *et al.* 2011) to a minimum quality score of 20 and a minimum length threshold of 50 bp and 30 bp for the paired-end and mate-pair reads, respectively.

We corrected the trimmed, paired-end reads for substitution sequencing errors using QUAKE v0.3.5 (Kelley *et al.* 2010). Salzberg *et al.* (2012) showed previously that the error correction of sequencing reads can greatly improve the *de novo* assembly of genomes, including genomes assembled using the program ABYSS (Simpson *et al.* 2009). QUAKE uses the distributions of infrequent and abundant *k*-mers to model the nucleotide error rates and subsequently corrects substitution errors. As input to QUAKE and again after error correction, we counted the frequency of 20-mers in the paired-end reads using DSK v1.6066 (Rizk *et al.* 2013). To estimate genome size, we divided the total number of error-free 20-mers by their peak coverage depth.

We assembled the genome using the trimmed and error-corrected paired-end reads with ABYSS v1.3.6. To determine the optimal *k*-mer length, we repeated the assembly using *k* = 40–88 in 8-bp increments. All scaffolding steps were performed using the trimmed mate-pair reads also in ABYSS, and only scaffolds longer than

Read and synthesize a bunch of these like you would ‘wet lab’ protocols

## 4. Sequence assembly

All cleaned sequences were assembled using the Newbler Assembler (25) v2.6 (build version 20110517\_1502) with the following parameters “-scaffold -het -large -cpu 3 -siod -noinfo”. Our decision to use Newbler was influenced by the large proportion of 454 sequences used and the ability for Newbler to handle multiple data, which allowed BACends, Illumina, and 454 data to be combined. Assemblies were run on a 16-processor node with 256 GB of RAM. Our current assembly consists of 43,234 contigs with an average size of 15,456 bp (min= 436 bp; max=287,935 bp), an N50 size of 29,456 bp, and an N50 count of 6,448. Scaffolding by virtue of the cleaned paired-end reads resulted in 5,745 scaffolds, with an average size of 123 kb (min= 1,732 bp; max= 15.98 Mb), an N50 size of 4.93 Mb, and an N50 count of 50. Based on the N90 statistics, 0.00% of our assembled sequence resides within 155 scaffolds, each of which is 1.16 Mb.

*Chamala et al (2016) Science*

# Bioinformatics protocols are analogous to any lab protocol

mina, USA). Preprocessing of the sequence reads included the removal of adapter sequences and removal of reads with >10% uncalled bases and/or >50% of bases with a Phred-scaled quality score <4. After preprocessing, all 100-bp (paired-end) and 50-bp (mate-pair) reads were retained as the set of ‘raw’ reads. We trimmed the 3’ end of all raw reads using a modified Mott algorithm in POPOOLATION v1.2.2 (Kofler *et al.* 2011) to a minimum quality score of 20 and a minimum length threshold of 50 bp and 30 bp for the paired-end and mate-pair reads, respectively.

We corrected the trimmed, paired-end reads for substitution sequencing errors using QUAKE v0.3.5 (Kelley *et al.* 2010). Salzberg *et al.* (2012) showed previously that the error correction of sequencing reads can greatly improve the *de novo* assembly of genomes, including genomes assembled using the program ABYSS (Simpson *et al.* 2009). QUAKE uses the distributions of infrequent and abundant k-mers to model the nucleotide error rates and subsequently corrects substitution errors. As input to QUAKE and again after error correction, we counted the frequency of 20-mers in the paired-end reads using DSK v1.6066 (Rizk *et al.* 2013). To estimate genome size, we divided the total number of error-free 20-mers by their peak coverage depth.

We assembled the genome using the trimmed and error-corrected paired-end reads with ABYSS v1.3.6. To determine the optimal k-mer length, we repeated the assembly using  $k = 40\text{--}88$  in 8-bp increments. All scaffolding steps were performed using the trimmed mate-pair reads also in ABYSS, and only scaffolds longer than

Cells were analyzed using a Cell Lab Quanta SC flow cytometer (Beckman Coulter). CD14-positive cells were stained with CD14-FITC (Miltenyi Biotec). Cells were incubated with propidium iodide to assess cell viability.

**Immunoblotting and antibodies.** Cells were harvested and total protein extracted in a buffer containing 25 mM HEPES (pH 7.4), 10% glycerol, 150 mM NaCl, 0.5% Triton X-100, 1 mM EDTA, 1 mM MgCl<sub>2</sub>, 1 mM ZnCl<sub>2</sub>, and protease inhibitors. The extracts were clarified by centrifugation for 10 minutes at 20,800g at 4°C. The extracted proteins (15 µg) were fractionated by SDS-PAGE, transferred to a polyvinylidene difluoride membrane (Millipore), and probed with an anti-A3A polyclonal antiserum, an anti-GFP monoclonal antibody (Clontech), or an anti-eEF1alpha monoclonal antibody (Upstate). The anti-A3A polyclonal serum was generated by immunizing a rabbit with a peptide corresponding to A3A residues 171-199 (CPFQPWDGLEEHSQALSGRLRAILQNQGN) mixed with TiterMax Gold adjuvant (Sigma). Primary antibodies were detected by incubation with fluorescently labeled secondary antibodies and imaging on an Odyssey imaging device (LI-COR Biosciences).

**DNA cytidine deaminase activity assays.** PBMC or transfected HEK-293T cell lysates were prepared as above for immunoblotting. The deaminase activity in the lysates was determined using a FRET-based assay essentially as described<sup>59</sup>. Briefly, serial dilutions of lysates were incubated for 2 h at 37°C with a DNA oligonucleotide 5’-(6-FAM)-AAA-TTCTAA-TAG-ATA-ATG-TGA-(TAMRA). FRET occurs between the fluorophores, decreasing FAM fluorescence. If

# Questions?

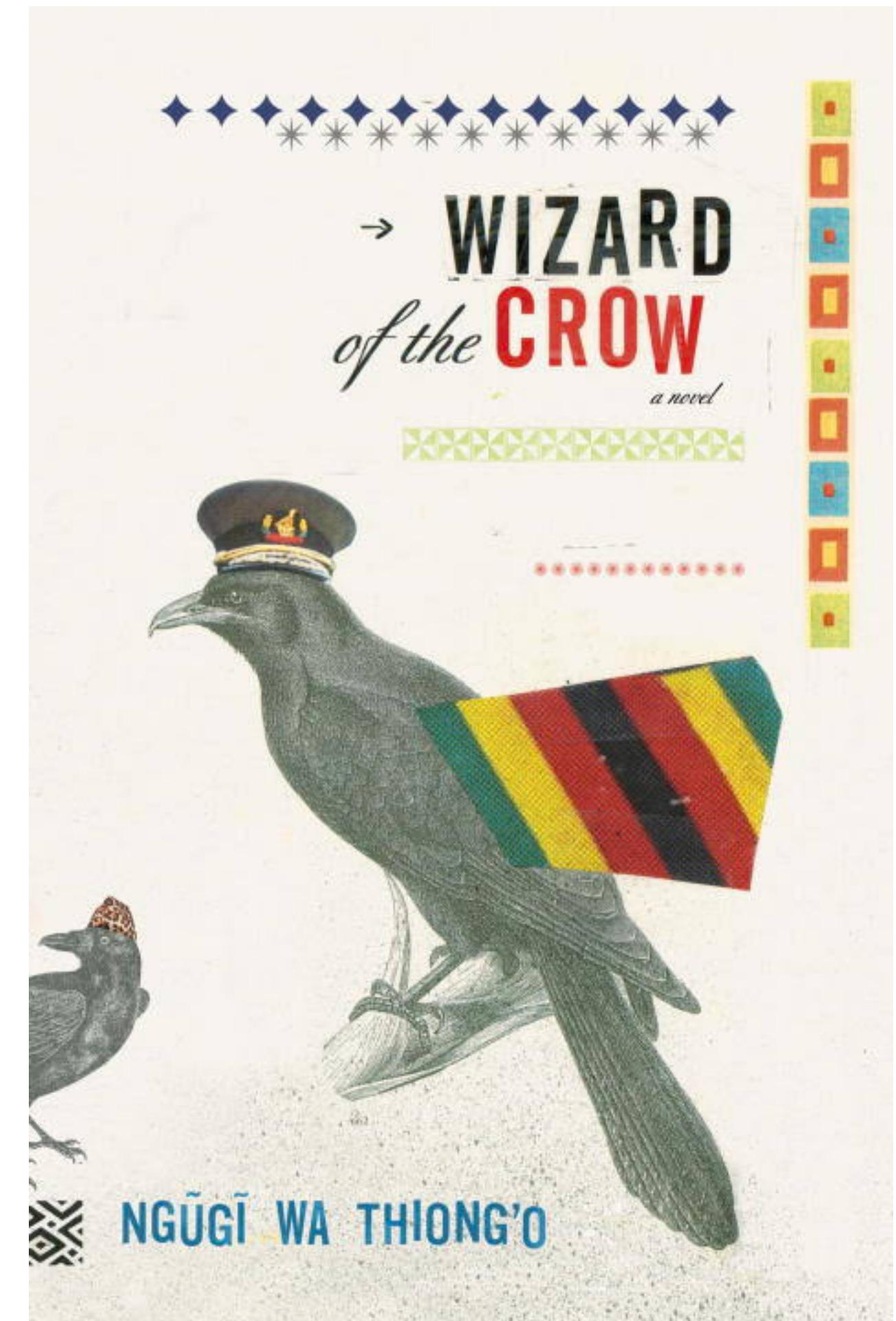


Image: Keith Bradnam, UC Davis