

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



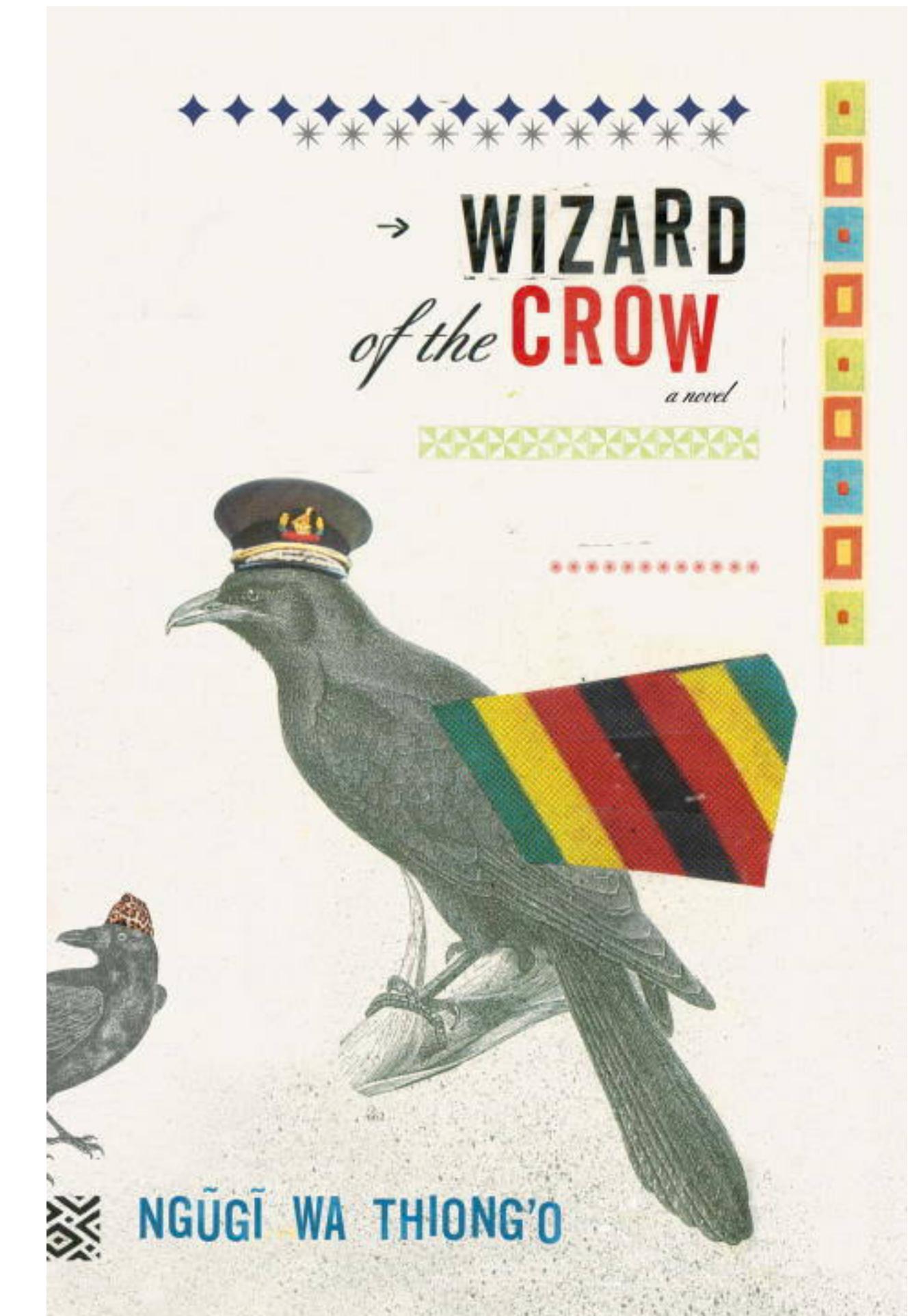
Keith Bradnam, UC Davis

Genome assembly paper exercise

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

Rules/info:

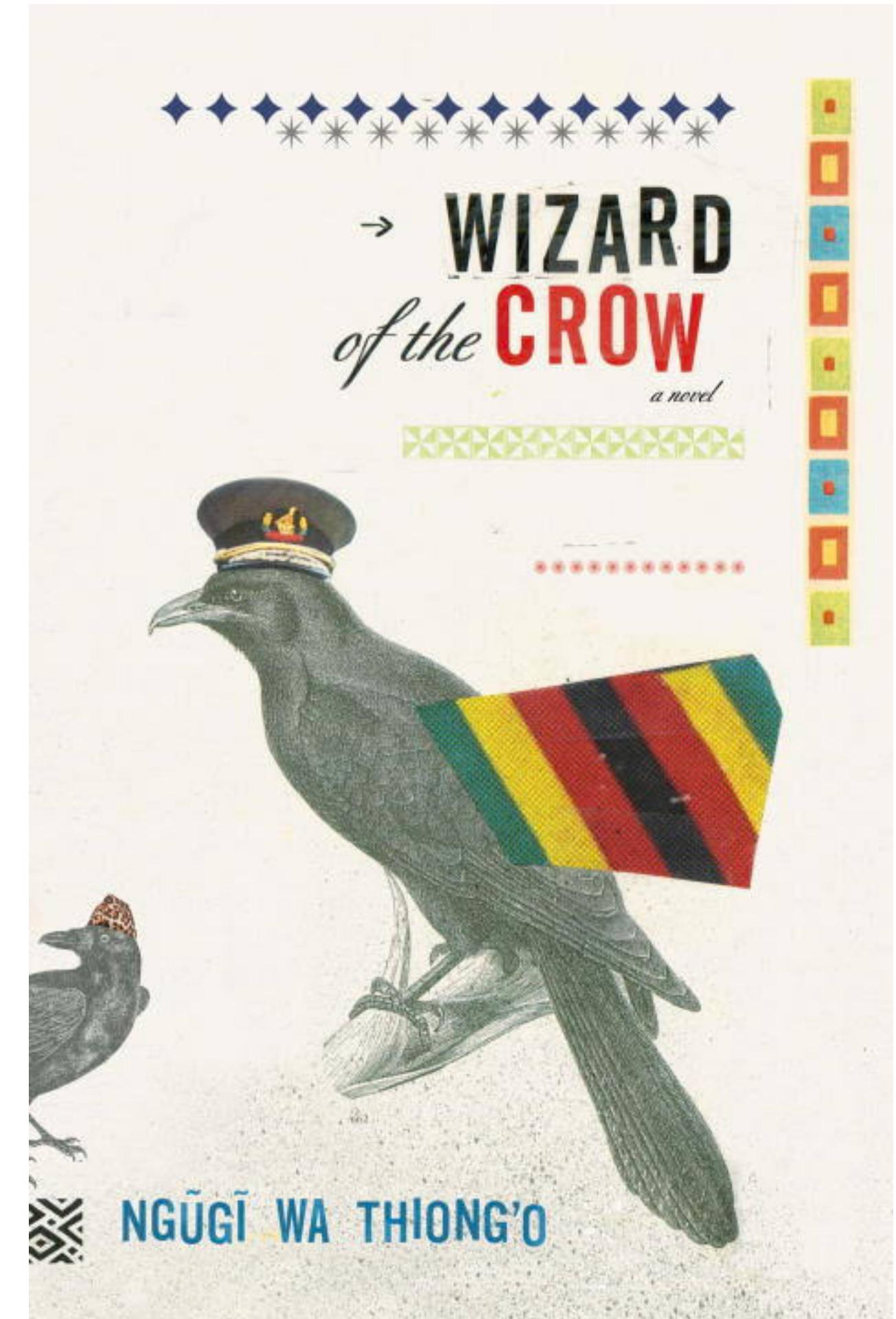
- Like real sequencing data, these reads contain errors.
The error rate is ~2%
- These are single-end 11-base reads
- The average coverage is ~6x
- You’re not allowed to google the answer
- Also: the answer is in the slides: don’t cheat!
- You can use your computers (i.e. word processors or text editors)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

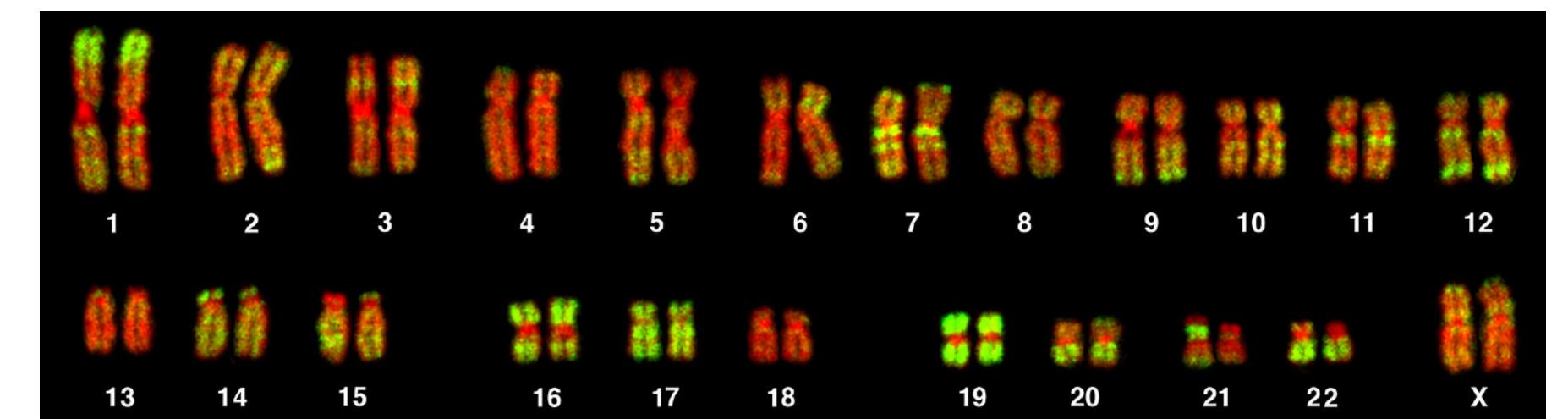


Conclusion: assembly is not trivial!

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)



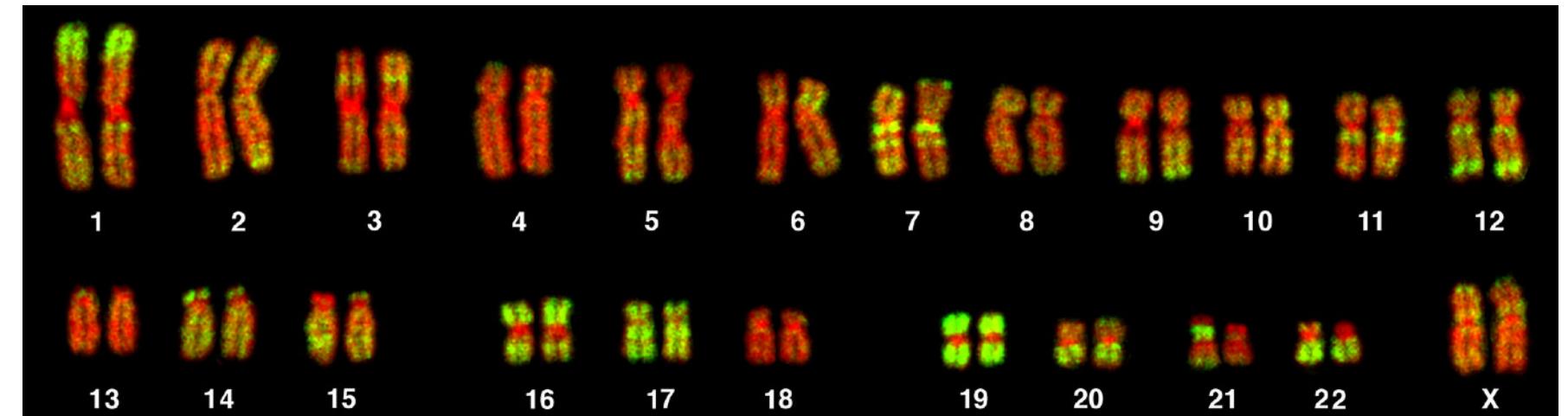
Bolzer et al (2005) PLoS Biol

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

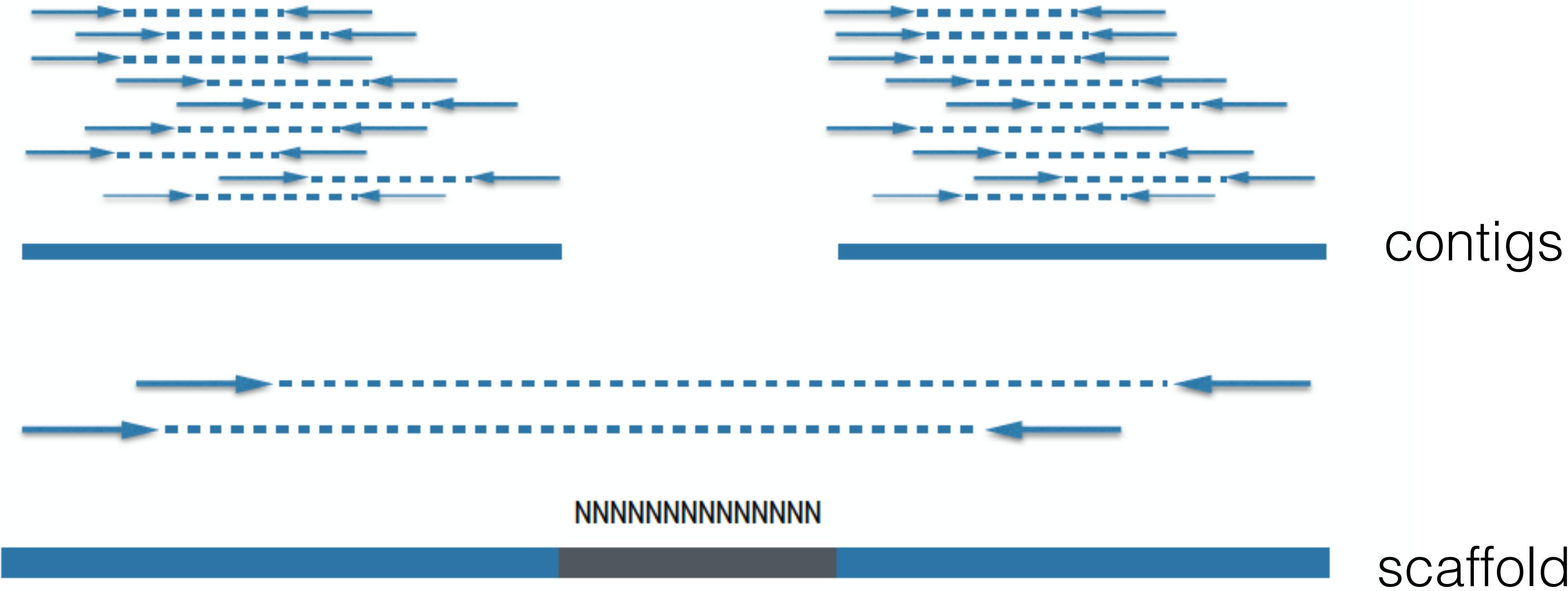
7) Not to mention annotation, which can be as hard as assembly!

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

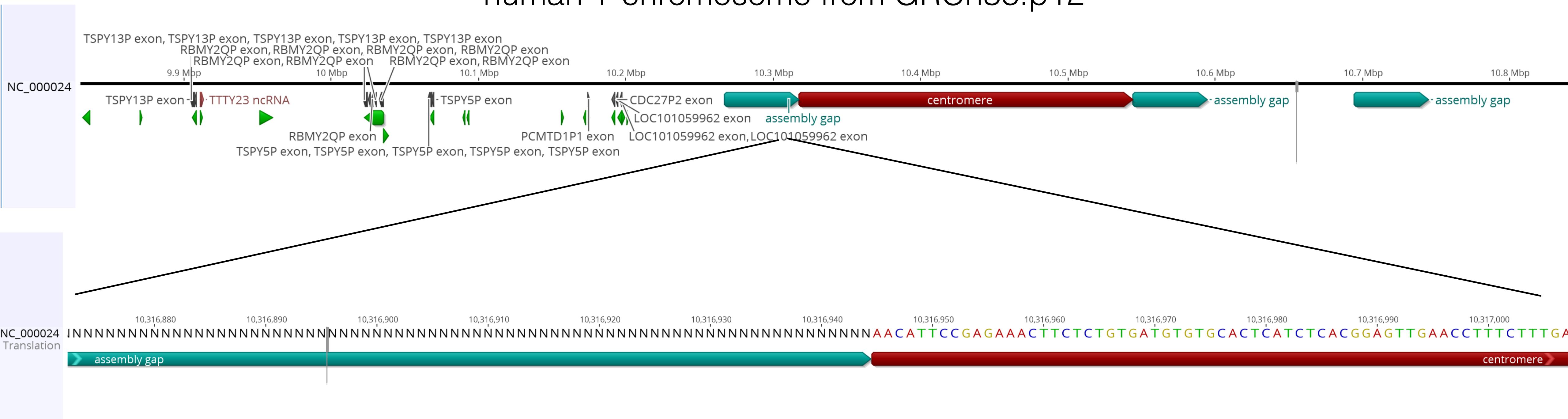




These “contigs” could be scaffolded because we have additional information

Sometimes even ‘complete’ assemblies contain gaps

human Y chromosome from GRCh38.p12



A truly complete (?) human genome

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.26.445798>; this version posted May 27, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

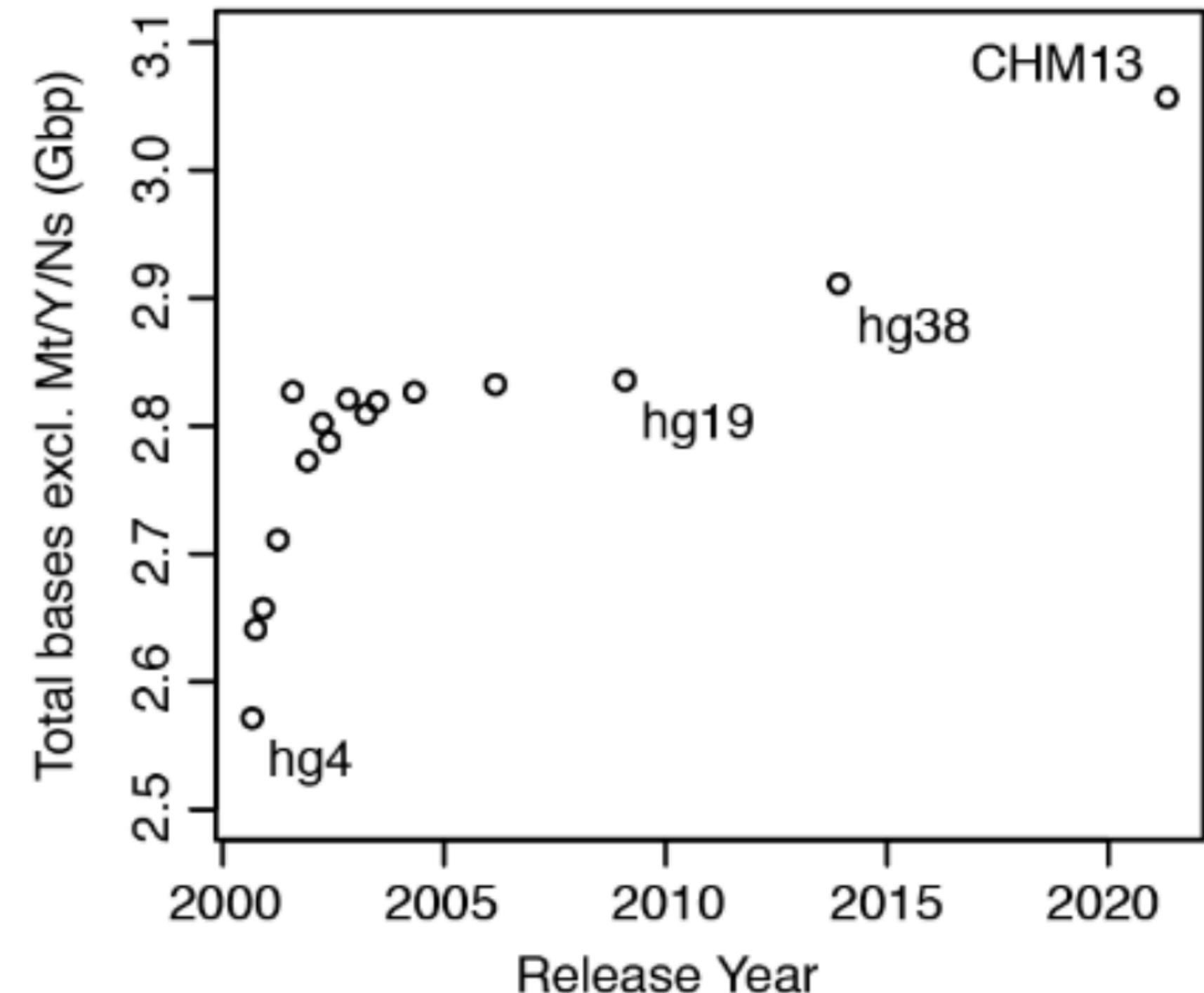
The complete sequence of a human genome

Sergey Nurk^{1,*}, Sergey Koren^{1,*}, Arang Rhee^{1,*}, Mikko Rautiainen^{1,*}, Andrey V. Bzikadze², Alla Mikheenko³, Mitchell R. Vollger⁴, Nicolas Altemose⁵, Lev Uralsky^{6,7}, Ariel Gershman⁸, Sergey Aganezov⁹, Savannah J. Hoyt¹⁰, Mark Diekhans¹¹, Glennis A. Logsdon⁴, Michael Alonge⁹, Stylianos E. Antonarakis¹², Matthew Borchers¹³, Gerard G. Bouffard¹⁴, Shelise Y. Brooks¹⁴, Gina V. Caldas¹⁵, Haoyu Cheng^{16,17}, Chen-Shan Chin¹⁸, William Chow¹⁹, Leonardo G. de Lima¹³, Philip C. Dishuck⁴, Richard Durbin²¹, Tatiana Dvorkina³, Ian T. Fiddes²², Giulio Formenti^{23,24}, Robert S. Fulton²⁵, Arkarachai Fungtammasan¹⁸, Erik Garrison^{11,26}, Patrick G.S. Grady¹⁰, Tina A. Graves-Lindsay²⁷, Ira M. Hall²⁸, Nancy F. Hansen²⁹, Gabrielle A. Hartley¹⁰, Marina Haukness¹¹, Kerstin Howe¹⁹, Michael W. Hunkapiller³⁰, Chirag Jain^{1,31}, Miten Jain¹¹, Erich D. Jarvis^{23,24}, Peter Kerpeljiev³², Melanie Kirsche⁹, Mikhail Kolmogorov³³, Jonas Korlach³⁰, Milinn Kremitzki²⁷, Heng Li^{16,17}, Valerie V. Maduro³⁴, Tobias Marschall³⁵, Ann M. McCartney¹, Jennifer McDaniel³⁶, Danny E. Miller^{4,37}, James C. Mullikin^{14,29}, Eugene W. Myers³⁸, Nathan D. Olson³⁶, Benedict Paten¹¹, Paul Peluso³⁰, Pavel A. Pevzner³³, David Porubsky⁴, Tamara Potapova¹³, Evgeny I. Rogaev^{6,7,39,40}, Jeffrey A. Rosenfeld⁴¹, Steven L. Salzberg^{9,42}, Valerie A. Schneider⁴³, Fritz J. Sedlazeck⁴⁴, Kishwar Shafin¹¹, Colin J. Shew²⁰, Alaina Shumate⁴², Yumi Sims¹⁹, Arian F. A. Smit⁴⁵, Daniela C. Soto²⁰, Ivan Sovic^{30,46}, Jessica M. Storer⁴⁵, Aaron Streets^{5,47}, Beth A. Sullivan⁴⁸, Françoise Thibaud-Nissen⁴³, James Torrance¹⁹, Justin Wagner³⁶, Brian P. Walenz¹, Aaron Wenger³⁰, Jonathan M. D. Wood¹⁹, Chunlin Xiao⁴³, Stephanie M. Yan⁴⁹, Alice C. Young¹⁴, Samantha Zarate⁹, Urvashi Surti⁵⁰, Rajiv C. McCoy⁴⁹, Megan Y. Dennis²⁰, Ivan A. Alexandrov^{3,7,51}, Jennifer L. Gerton¹³, Rachel J. O'Neill¹⁰, Winston Timp^{8,42}, Justin M. Zook³⁶, Michael C. Schatz^{9,49}, Evan E. Eichler^{4,24,†}, Karen H. Miga^{11,†}, Adam M. Phillippy^{1,†}

¹⁻⁵¹ Affiliations are listed at the end

* Equal contribution

† Corresponding authors: Evan E. Eichler (eee@gs.washington.edu); Karen H. Miga (khmiga@ucsc.edu); Adam M. Phillippy (adam.phillippy@nih.gov)



A truly complete (?) human genome

Required lots of data!

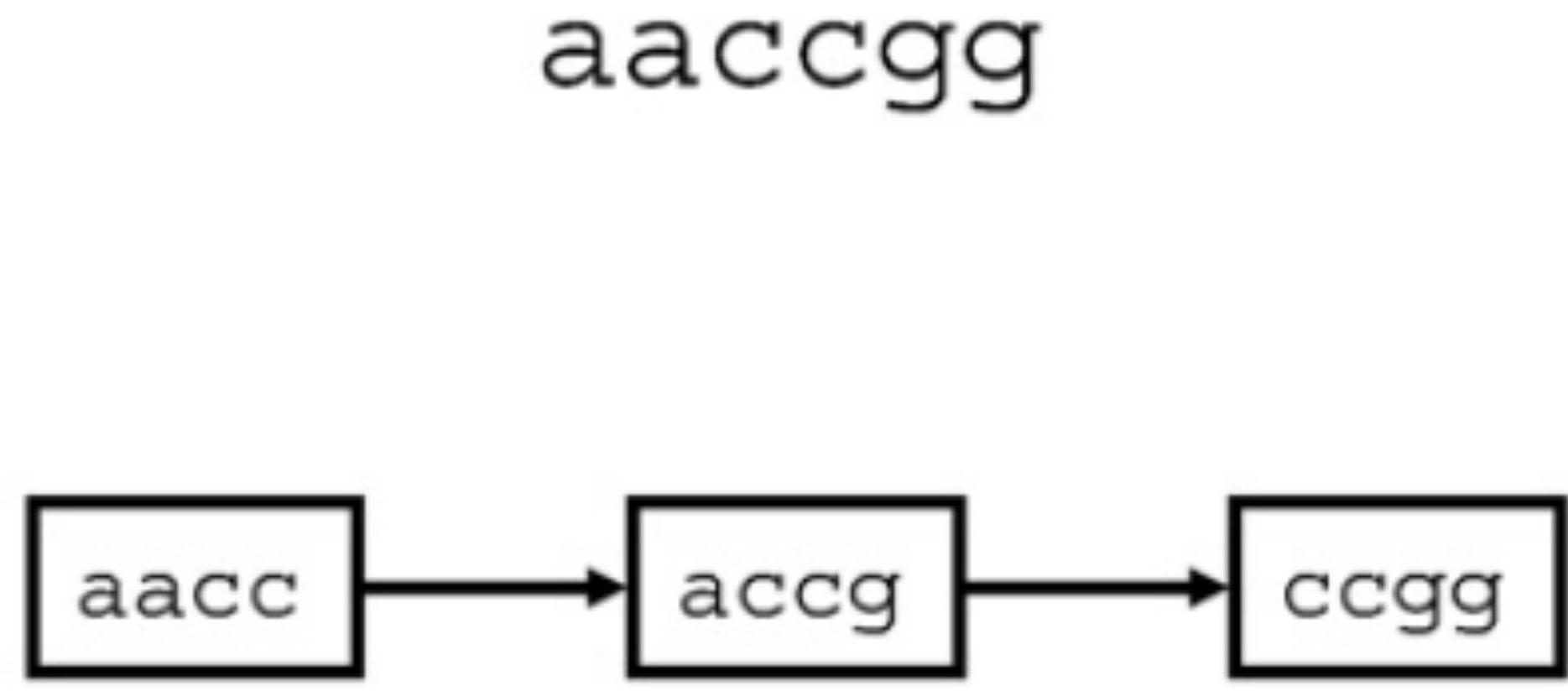
The screenshot shows a GitHub repository page for 'marbl/CHM13'. The URL in the address bar is <https://github.com/marbl/CHM13>. The page displays the contents of the 'README.md' file. The first section is titled 'Introduction'.

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:

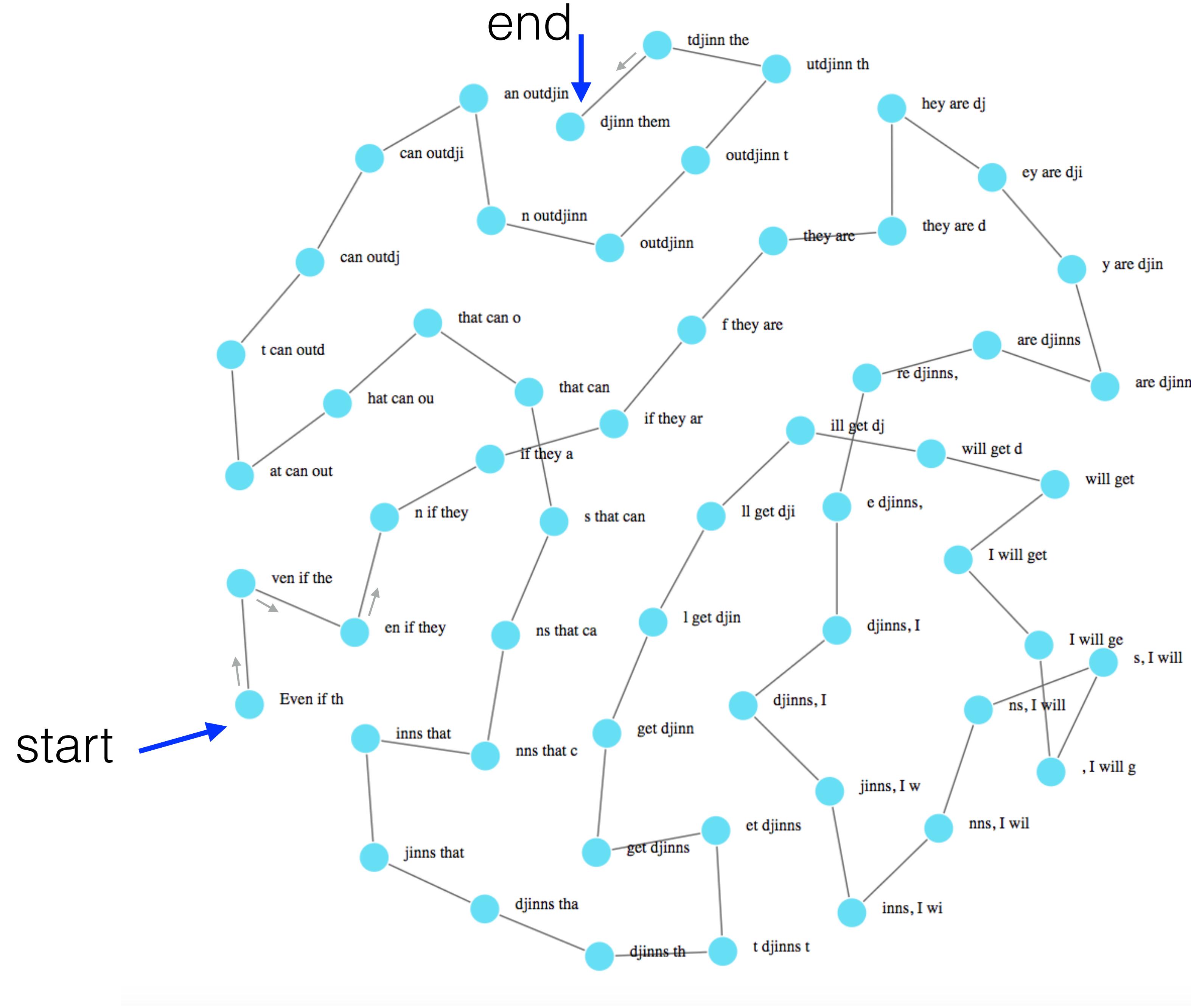
What do you call a tetranucleotide that's had a base added or removed?

A former 4-mer.

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

19 Retweets 109 Likes

k=10

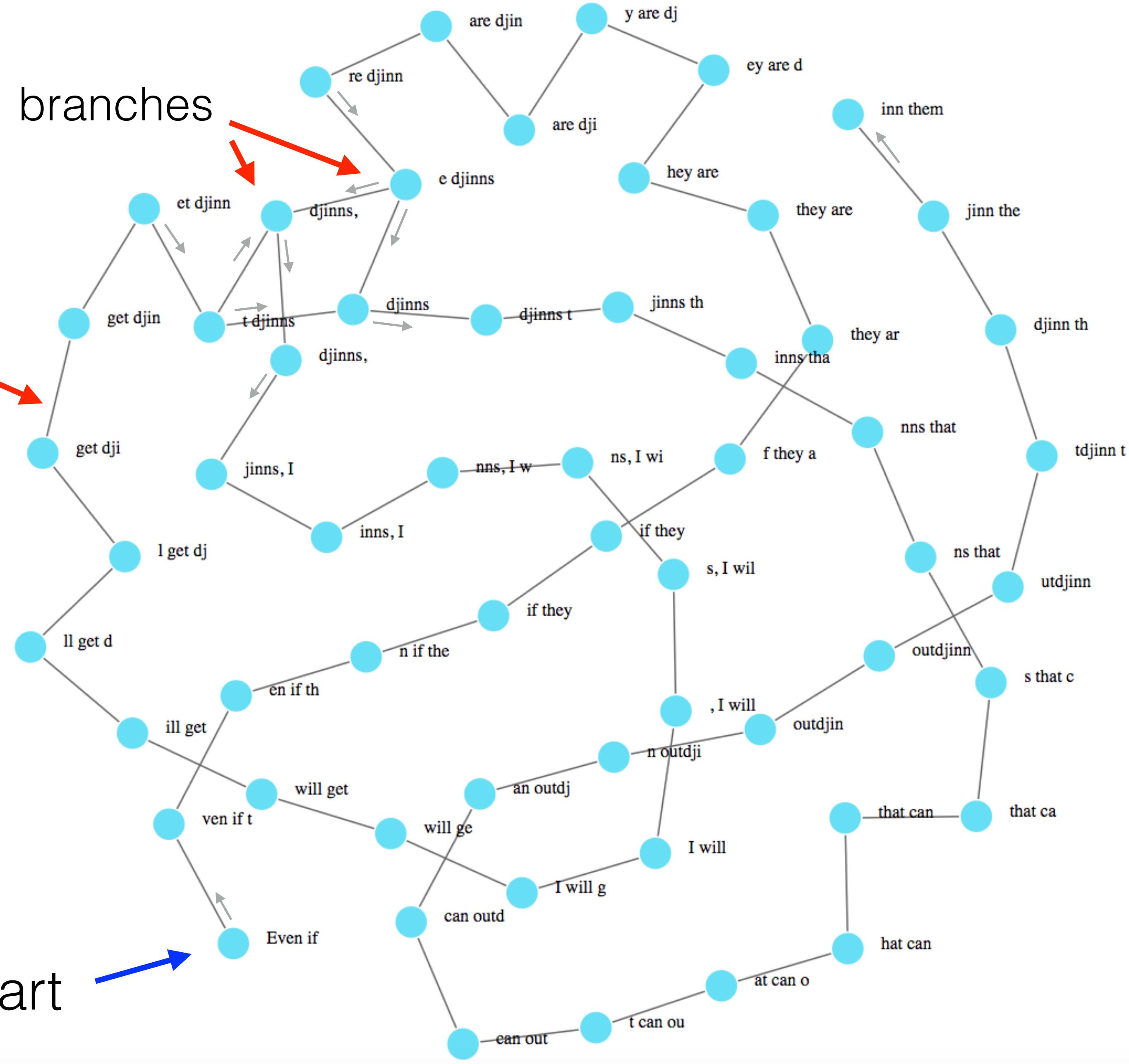


k=8

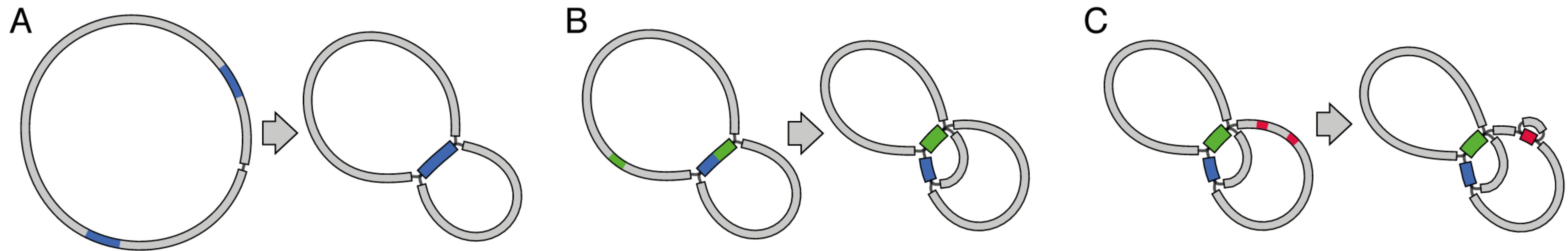
bubble
(circular path)

art

branches



The impact of additional repeats on graph complexity



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

Assemblers use a variety of strategies to try to resolve graph complexity

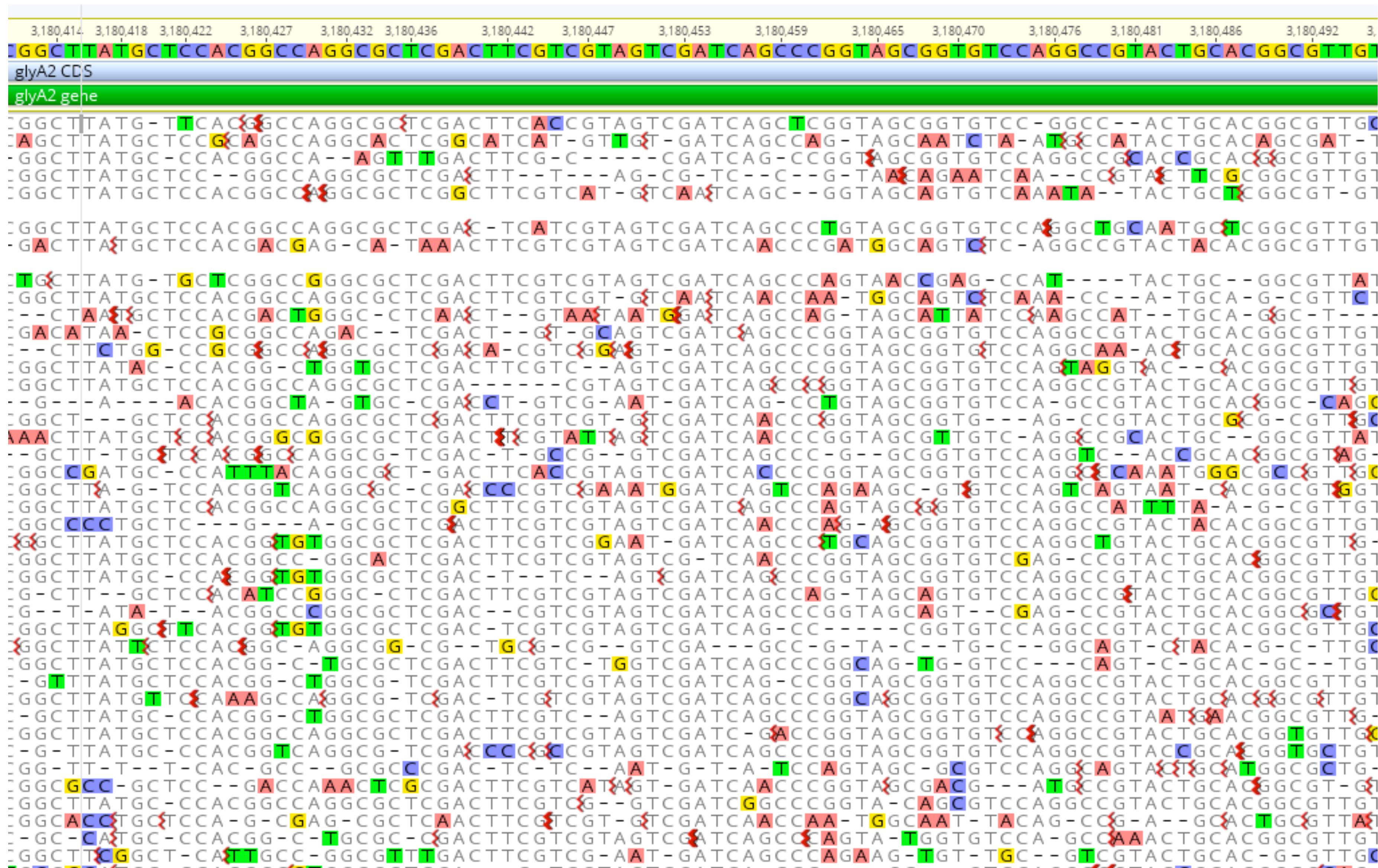
To read more about these strategies:

- Miller JR, Koren S, Sutton G. Assembly algorithms for next generation sequencing data. *Genomics* 2010;95:315–27.
- Compeau PE, Pevzner PA, Tesler G. How to apply de Bruijn graphs to genome assembly. *Nat Biotechnol* 2011;29:987–91.
- Nagarajan N, Pop M. Sequence assembly demystified. *Nat Rev Genet* 2013;14:157–67.
- Sohn JI, Nam JW. The present and future of de novo whole-genome assembly. *Brief Bioinform*. 2016 Oct 14. pii: bbw096.

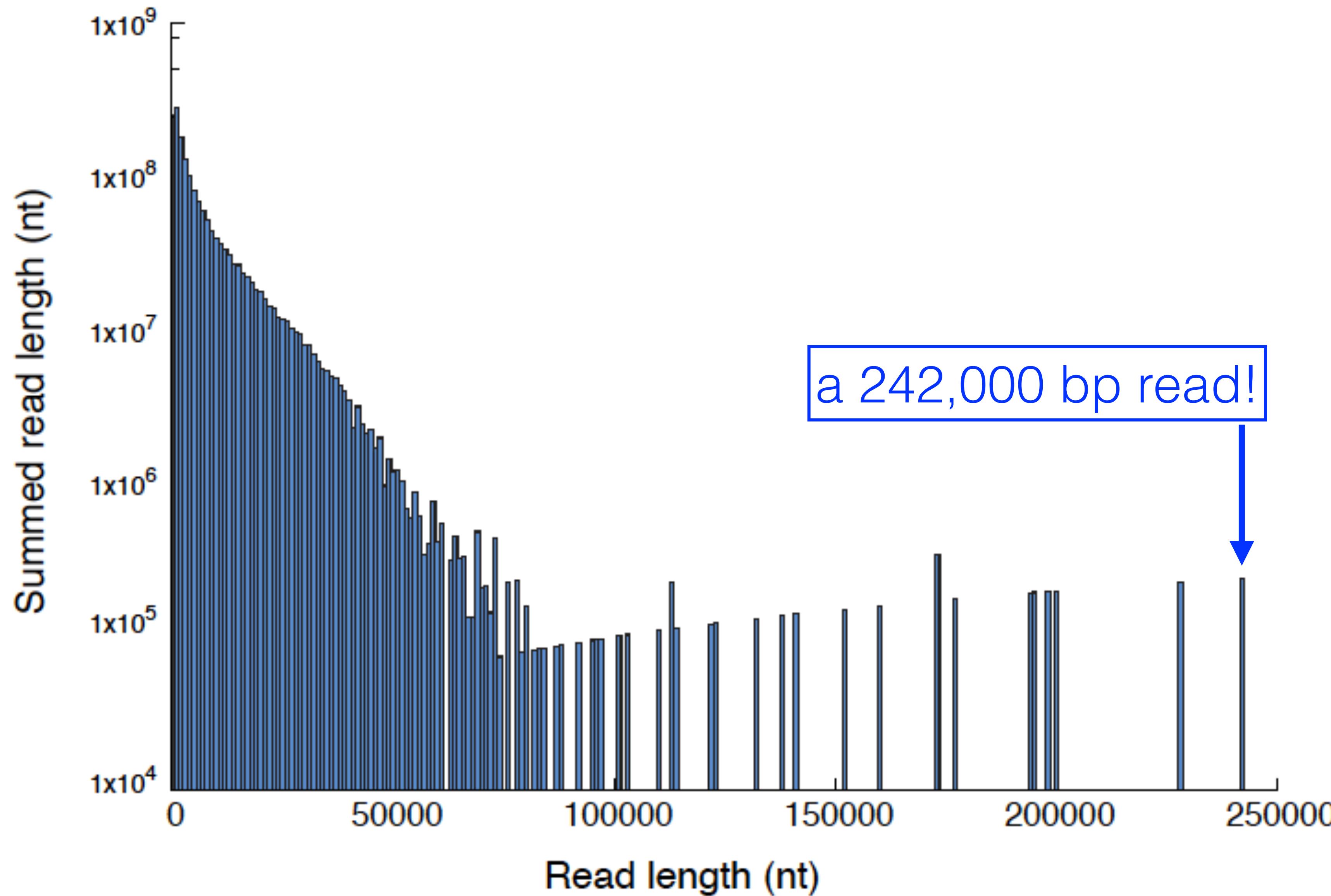
To some extent, the issues attributable to short-read only assembly are becoming moot as long read technology becomes cheaper and better.

Assemblies that mix long and short reads are called ‘hybrid’ assemblies, and they are increasingly the norm.

Long reads have high error rates - but their consensus has a lower error rate



Long reads can also be combined with short reads to create high quality ‘hybrid’ assemblies.



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?
- What fraction of reads map to 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 *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](#)

- [**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 somewhat bleak picture, but don't be too intimidated:
genome sequencing and assembly *is* possible.

Not all assembly problems are equally difficult!

tiny ssDNA genome

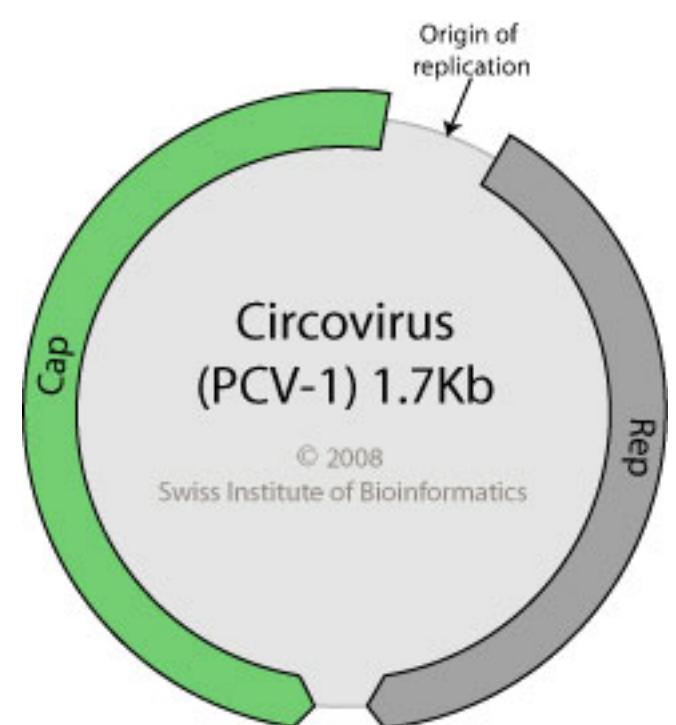
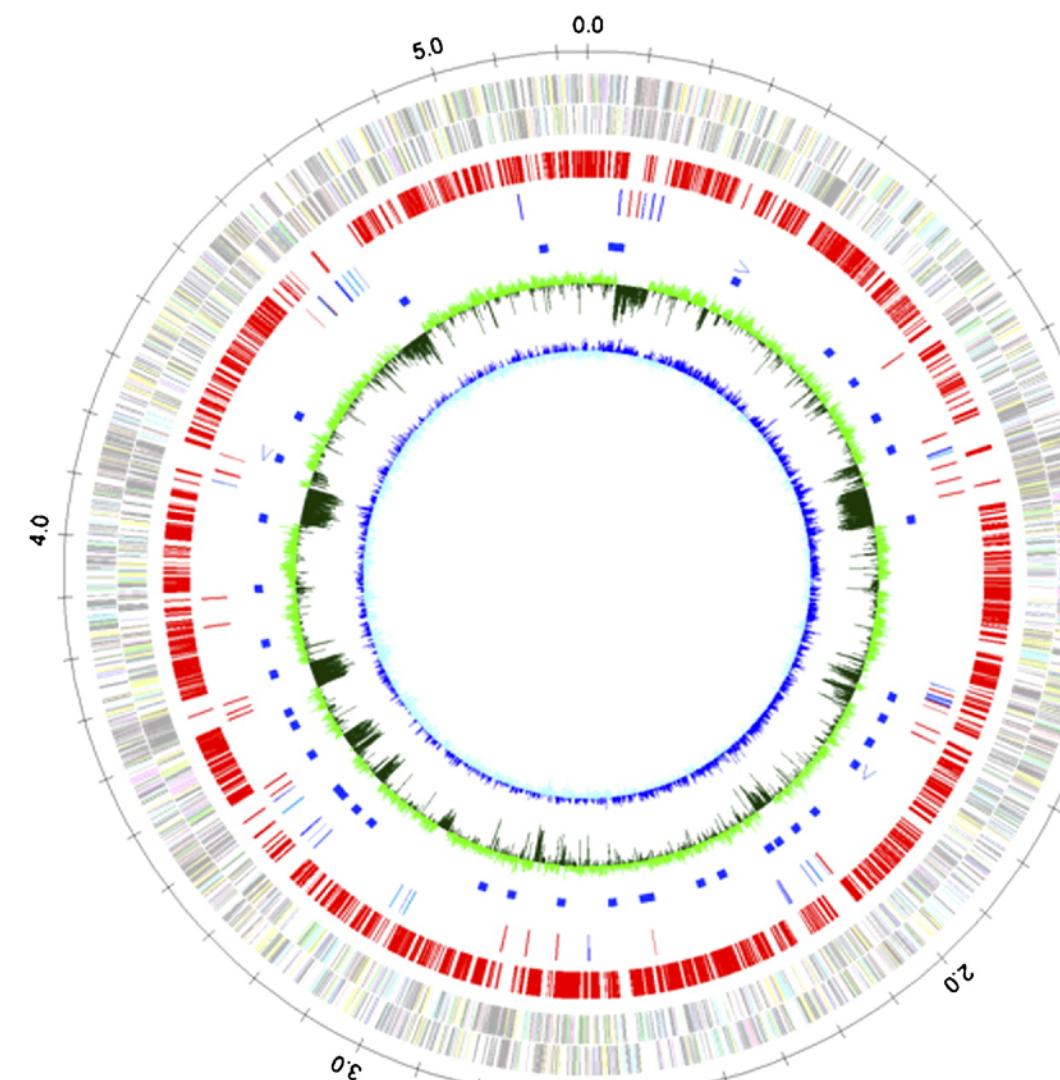


image: viralzone

bacterial genomes ~5 Mbp



Nakazawa et al (2009) Genome Research

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



image: Univ of Alabama

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,^{*1} ELMIRA MOHANDESAN,^{*} JUKKA CORANDER[†] and PAMELA A. BURGER^{*}

**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 MgCl₂, 1 mM ZnCl₂, 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 described⁵⁹. 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?

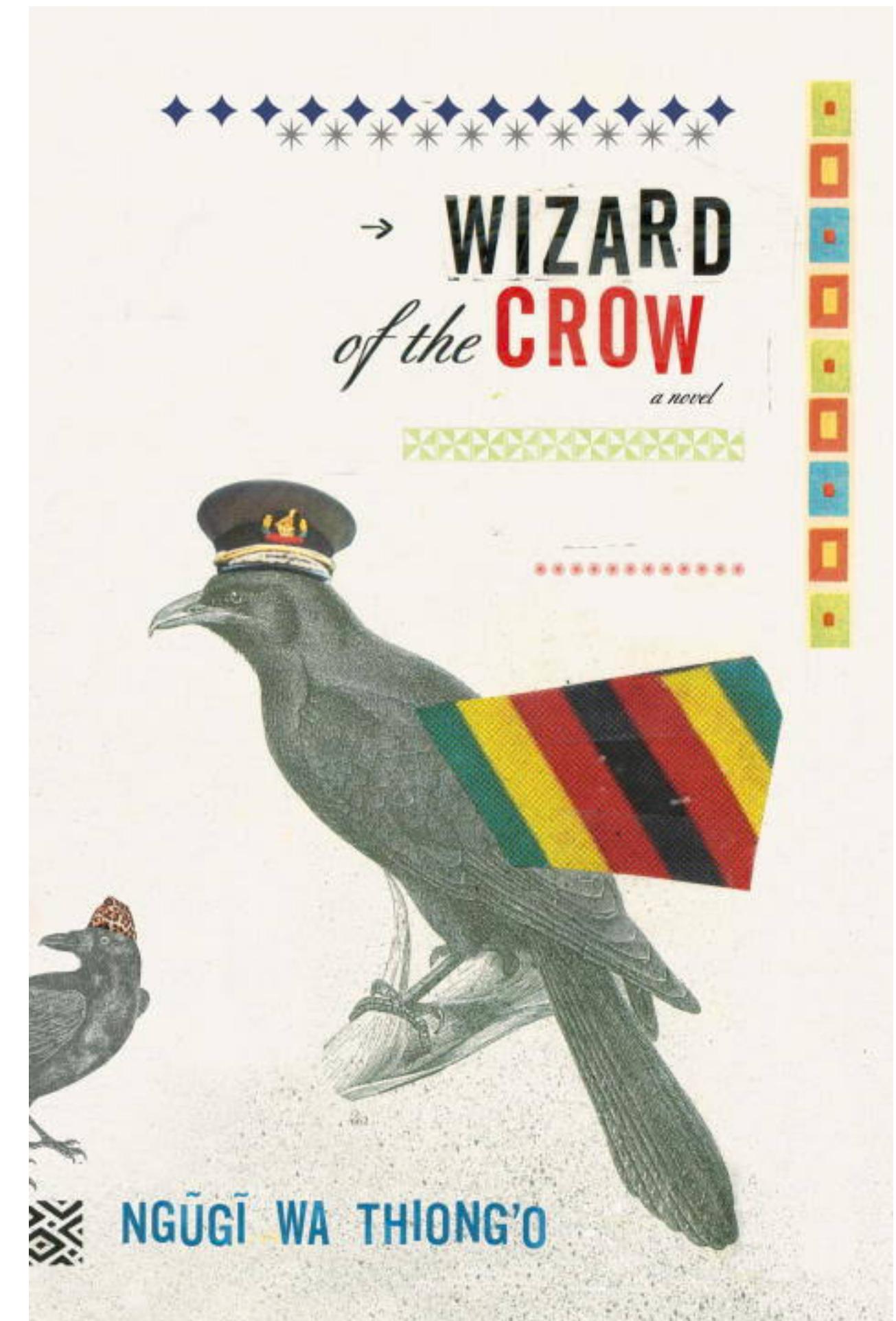


Image: Keith Bradnam, UC Davis