

The human genome

Michael Schatz

Sept 11, 2024

Lecture 5: Applied Comparative Genomics



Assignment 2: Genome Assembly

Due Monday Sept 16 by 11:59pm

The screenshot shows a GitHub repository page for 'appliedgenomics2024'. The repository is public and contains several files. The 'assignment2' folder is expanded, showing 'README.md' as the selected file. The README file content is as follows:

```
Assignment 2: Genome Assembly

Assignment Date: Monday, September 9, 2024
Due Date: Monday, September 16, 2023 @ 11:59pm

Assignment Overview

In this assignment, you will explore the steps for de novo genome assembly. This will start with constructing and analyzing the de Bruijn graph of reads using a short python/R script. Next we will evaluate the expected and observed coverage in a set of reads. These reads come from a mysterious pathogen that contains a secret message encoded somewhere in the genome. The secret message will be recognizable as a novel insertion of sequence not found in the reference. Your task is to assess the quality of the reads, assemble the genome, identify, and decode the secret message. If all goes well the secret message should decode into a recognizable english text, otherwise double check your coordinates and try again. As a reminder, any questions about the assignment should be posted to Piazza.

For this assignment, we recommend you install and run the tools using bioconda. There are some tips below in the Resources section. Note on Mac, we highly recommend you install the x86_64 package even if you are using an M1/M2 chip.

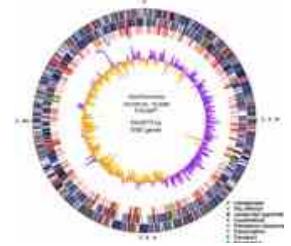
Question 1. de Bruijn Graph construction [10 pts]

• Q1a. Write a script (in python, R, C++, etc) to draw the de Bruijn graph for the following reads using k=3 (assume all reads are from the forward strand, no sequencing errors, complete coverage of the genome). You may find graphviz to be helpful (see below).
```

<https://github.com/schatzlab/appliedgenomics2024/tree/main/assignments/assignment2>

Check Piazza for questions!

Assembly Summary



Assembly quality depends on

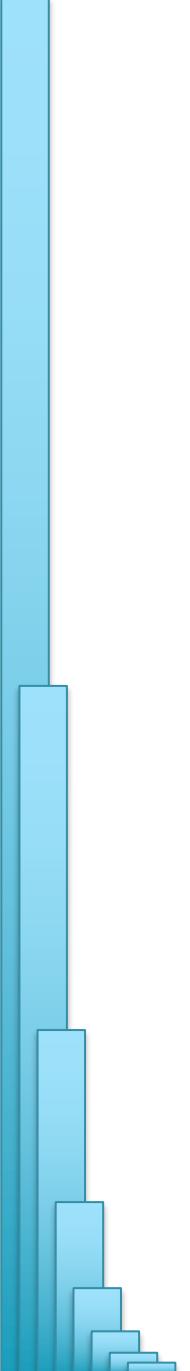
1. **Coverage:** low coverage is mathematically hopeless
 2. **Repeat composition:** high repeat content is challenging
 3. **Read length:** longer reads help resolve repeats
 4. **Error rate:** errors reduce coverage, obscure true overlaps
-
- Assembly is a hierarchical, starting from individual reads, build high confidence contigs/unitigs, incorporate the mates to build scaffolds
 - Extensive error correction is the key to getting the best assembly possible from a given data set
 - Watch out for collapsed repeats & other misassemblies
 - Globally/Locally reassemble data from scratch with better parameters & stitch the 2 assemblies together

SV Types

Insertion into Reference R: AIB Q: AB	Insertion into Query R: AB Q: AIB
Collapse Query R: ARRB Q: ARB	Collapse Reference R: ARB Q: ARRB
Collapse Query w/ Insertion R: ARIRB Q: ARB Exact tandem alignment if I=R	Collapse Reference w/ Insertion R: ARB Q: ARIRB Exact tandem alignment if I=R
Collapse Query R: ARRRB Q: ARRB	Collapse Reference R: ARRB Q: ARRRB
Inversion R: ABC Q: AB'C	Rearrangement w/ Disagreement R: ABCDE Q: AFCBE

- Different structural variation types / misassemblies will be apparent by their pattern of breakpoints
- Most breakpoints will be at or near repeats
- Things quickly get complicated in real genomes

[http://mummer.sf.net/manual/
AlignmentTypes.pdf](http://mummer.sf.net/manual/AlignmentTypes.pdf)



The human genome

The scale of DNA in our body is staggering.

- A typical human is comprised of roughly 40 trillion human cells
(excluding trillions of bacterial cells in our gut)
- If stretched out, each haploid genome would be roughly 2 meters.
- So, each cell has 4 meters of DNA.
- $40 \text{ trillion} * 4 \text{ meters} = 160 \text{ trillion meters}$.
- $160 \text{ trillion meters} / 1609.34 = 99,750,623,441 \text{ miles}$
- $99,750,623,441 / 92,960,000 = 1,073.05 \text{ trips to the sun.}$

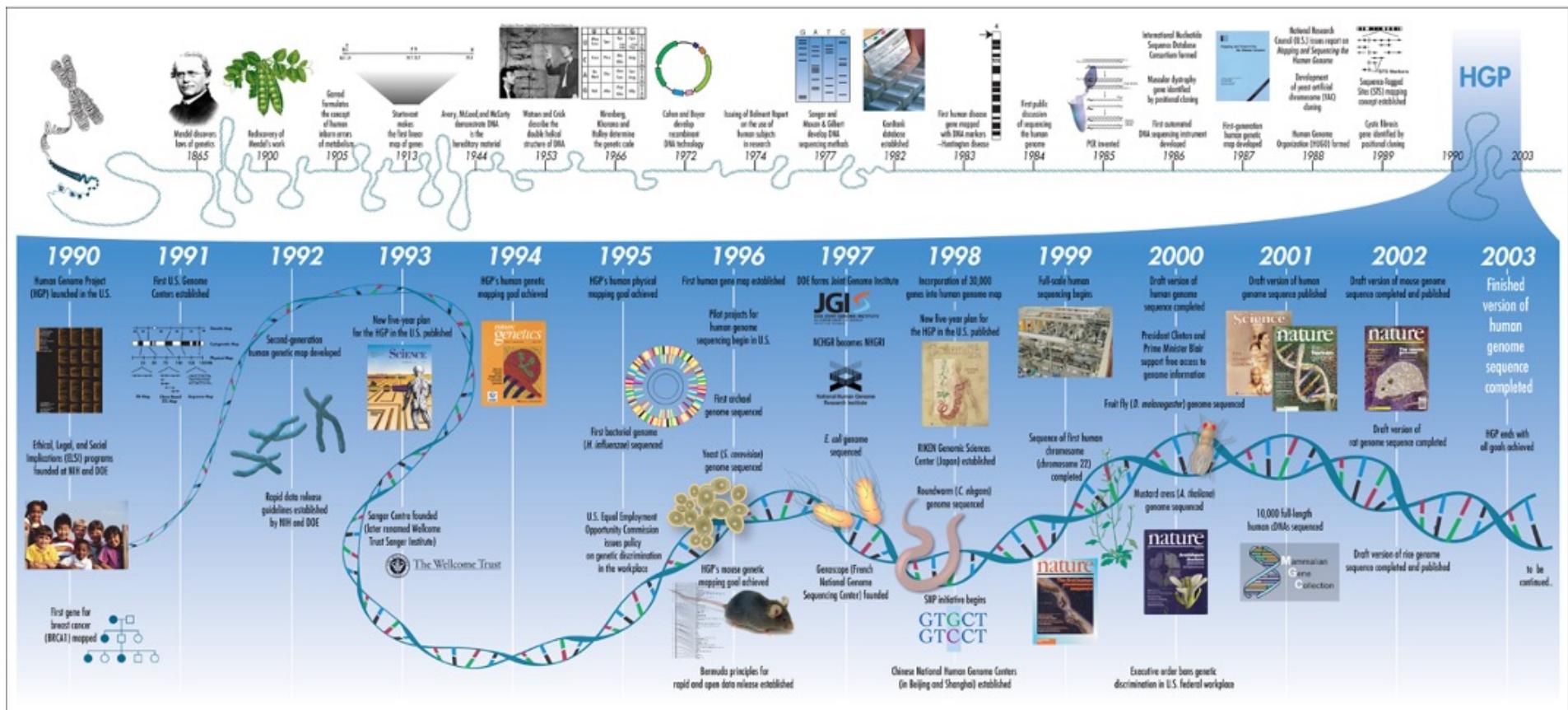
A typical cell replicates about 100 times

160 trillion meters x 100 =

1.69123746 light years

More info

History of the Human Genome Project



The reference human genome



“Without a doubt, this is the most important, most wondrous map ever produced by humankind.”

*Bill Clinton
June 26, 2000*

The reference human genome



“Without a doubt, this is the most important, most wondrous map ever produced by humankind.”

*Bill Clinton
June 26, 2000*



The Sequence of the Human Genome

Venter et al.

Science 291, pp 1304-1351 (2001)



Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium

Nature 409, pp 860–921 (2001)

Two Human Genomes?

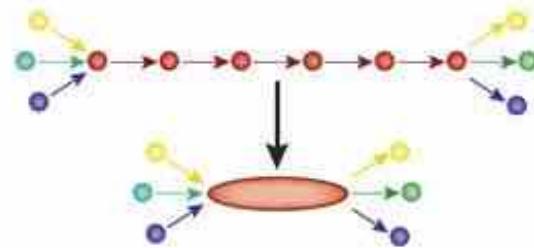
1. Fragment DNA and sequence



2. Find overlaps between reads

...AGCCTAGACCTACA **GGATCGGCGACAGT**
 GGATCGGCGACAGT CGCATATCCGGT...

3. Assemble overlaps into contigs



4. Assemble contigs into scaffolds



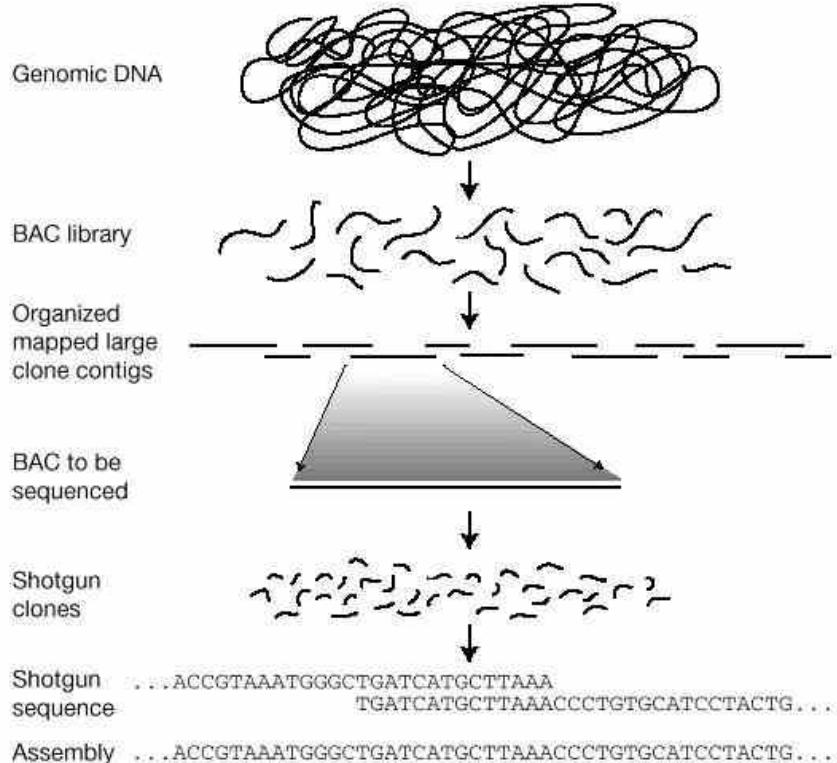
The Sequence of the Human Genome

Venter et al.

Science 291, pp 1304-1351 (2001)

(Figure from Baker (2012) Nature Methods)

Hierarchical shotgun sequencing



Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium

Nature 409, pp 860–921 (2001)

Who is the reference human?

The Buffalo News/Sunday, March 23, 1997

ment abuse, civil disobedience

ople. But the very nature of government creates a mind set that inspires increase their authority, always at the expense of the people."

"government has forgotten that it's the people," Parlato added, "acting more like it's the master." To and the Lapps share an abiding non-violent civil disobedience.

"We insist on being respectful in our form of resistance," Barbara Lyn Lapp said. "But if we claim to care about our rights, we must protest government instead of violence has to be the watchword," said, calling civil disobedience the "is of the violent militia movement." Non-violence can serve as an anti-government oppression, he added. "A law is unjust or you're given an authority without moral or legal authority,

Rachel Lapp says she believes government can be good, when it controls the aggressors in society. Instead, it too often comes down on the side of the aggressors, who enforce child-protection laws, compulsory education, disclosure rules on tax forms and seat belt laws.

"We want people to see the correlation between what happened to us and what can happen to anyone when government gets out of hand," Rachel Lapp said.

The Lapps and Parlato will be joined by Samuel Radford III, a critic of public education who was arrested and pleaded guilty to reduced charges following a 1993 disturbance at the City Campus of Erie Community College.

WANTED
20 Volunteers
to participate in the
Human Genome Project
a very large international scientific research effort.

The goal is to decode the human hereditary information (*human blueprint*) that determines all individual traits inherited from parents. The outcome of the project will have tremendous impact on future progress of medical science and lead to improved diagnosis and treatment of hereditary diseases.

Volunteers will receive information about the project from the Clinical Genetics Service at Roswell Park, and sign a consent form before participating.

No personal information will be maintained or transferred.

Individuals must be at least 18 years of age.
Persons who have undergone chemotherapy are not eligible.

For more information please contact the
Clinical Genetics Service
845-5720 (9:00 am - 3:00 pm)
March 24 - 26, 1997

ROSWELL PARK
CANCER INSTITUTE

WANTED

20 Volunteers

to participate in the

Human Genome Project

a very large international scientific research effort.

The goal is to decode the human hereditary information (*human blueprint*) that determines all individual traits inherited from parents. The outcome of the project will have tremendous impact on future progress of medical science and lead to improved diagnosis and treatment of hereditary diseases.

Volunteers will receive information about the project from the Clinical Genetics Service at Roswell Park, and sign a consent form before participating.

No personal information will be maintained or transferred.

Volunteers will provide a one-time donation of a small blood specimen. A small monetary reimbursement will be provided to the participants for their time and effort.

Individuals must be at least 18 years of age.
Persons who have undergone chemotherapy are not eligible.

For more information please contact the
Clinical Genetics Service
845-5720 (9:00 am - 3:00 pm)
March 24 - 26, 1997

ROSWELL PARK
CANCER INSTITUTE

Pieter de Jong, RPCI

Who is the reference human?

The Buffalo News/Sunday, March 23, 1997

ment abuse, civil disobedience

ople. But the very nature of government creates a mind set that inspires them to increase their authority, always at the expense of the people."

"Government has forgotten that it's the people," Parlato said, adding more like it's the master." "The Lapps share an abiding non-violent civil disobedience. We insist on being respectful in our form of resistance," Barbara Lyn Lapp said. "But if we claim to care about our rights, we must protest government instead of violence."

Violence has to be the watchword, said, calling civil disobedience the is of the violent militia movement. Non-violence can serve as an anti-government oppression, he added. "A law is unjust or you're given an without moral or legal authority,

You should refuse it," Parlato said. "And, if need be, you have to be brave enough to accept the consequences."

Rachel Lapp says she believes government can be good, when it controls the aggressors in society. Instead, it too often comes down on the side of the aggressors, who enforce child-protection laws, compulsory education, disclosure rules on tax forms and seat belt laws.

"We want people to see the correlation between what happened to us and what can happen to anyone when government gets out of hand," Rachel Lapp said.

The Lapps and Parlato will be joined by Samuel Radford III, a critic of public education who was arrested and pleaded guilty to reduced charges following a 1993 disturbance at the City Campus of Erie Community College.

WANTED
20 Volunteers
to participate in the
Human Genome Project
a very large international scientific research effort.

The goal is to decode the human hereditary information (*human blueprint*) that determines all individual traits inherited from parents. The outcome of the project will have tremendous impact on future progress of medical science and lead to improved diagnosis and treatment of hereditary diseases.

Volunteers will receive information about the project from the Clinical Genetics Service at Roswell Park, and sign a consent form before participating.

No personal information will be maintained or transferred.

Volunteers will provide a one-time donation of a small blood specimen. A small monetary reimbursement will be provided to the participants for their time and effort.

Individuals must be at least 18 years of age.
Persons who have undergone chemotherapy are not eligible.

ROSWELL PARK CANCER INSTITUTE

For more information please contact the
Clinical Genetics Service
845-5720 (9:00 am - 3:00 pm)
March 24 - 26, 1997



Pieter de Jong, RPCI

Who is the reference human?

The Buffalo News/Sunday, March 23, 1997

ment abuse, civil disobedience

ople. But the very nature of government creates a mind set that inspires increase their authority, always at the expense of the people."

"Government has forgotten that it's the people," Parlato added, "acting more like it's the master." Tom and the Lapps share an abiding non-violent civil disobedience.

"We insist on being respectful in our form of resistance," Barbara Lyn Lapp said, "but if we claim to care about our rights, we must protest government instead of violence."

Violence has to be the watchword, said, calling civil disobedience the way of the violent militia movement. Non-violence can serve as an anti-government oppression, he added. "A law is unjust or you're given an without moral or legal authority,

Rachel Lapp says she believes government can be good, when it controls the aggressors in society. Instead, it too often comes down on the side of the aggressors, who enforce child-protection laws, compulsory education, disclosure rules on tax forms and seat belt laws.

"We want people to see the correlation between what happened to us and what can happen to anyone when government gets out of hand," Rachel Lapp said.

The Lapps and Parlato will be joined by Samuel Radford III, a critic of public education who was arrested and pleaded guilty to reduced charges following a 1993 disturbance at the City Campus of Erie Community College.

WANTED
20 Volunteers
to participate in the
Human Genome Project
a very large international scientific research effort.

The goal is to decode the human hereditary information (*human blueprint*) that determines all individual traits inherited from parents. The outcome of the project will have tremendous impact on future progress of medical science and lead to improved diagnosis and treatment of hereditary diseases.

Volunteers will receive information about the project from the Clinical Genetics Service at Roswell Park, and sign a consent form before participating.

No personal information will be maintained or transferred.

Volunteers will provide a one-time donation of a small blood specimen. A small monetary reimbursement will be provided to the participants for their time and effort.

Individuals must be at least 18 years of age. Persons who have undergone chemotherapy are not eligible.

For more information please contact the Clinical Genetics Service 845-5720 (9:00 am - 3:00 pm) March 24 - 26, 1997

ROSWELL PARK CANCER INSTITUTE

Pieter de Jong, RPCI

Appendix: Identifying the ancestry of segments of the human genome reference sequence

To compare Neandertal to present-day human haplotypes for the purpose of population genetic analysis, we needed to have long haploid sequences from present-day humans that were of known ancestry. To identify such segments, we took advantage of the fact that the human reference sequence is haploid over scales of tens of kilobases, because it is comprised of a tiling-path of Bacterial Artificial Chromosomes (BACs) or other clone types that are of typical size 50-150 kb (S92). We do not know of any other substantial source of high quality human haploid sequences of the requisite size.

Determining the ancestries of the libraries in the human genome reference sequence using HAPMIX

It is crucial to know the 'ancestry' of a clone to use it in a meaningful population genetic analysis. In what follows, we define 'ancestry' as the geographic region in which a clone's ancestor lived 1,000 years ago, inferred based on its genetic proximity to other individuals from that region today. This definition allows us to classify clones from Chinese Americans as "East Asian," from European Americans as "European", and from African Americans as either "West African" or "European".

To identify the ancestries of the libraries comprising most of the human genome reference sequence, we used a list of 26,558 clones tiling the great majority of the genome, most of which we were able to assign to a library of origin. Restricting to the autosomes, we identified 21,156 clones that seemed to fall into 9 libraries based on the naming scheme: CTA (n=199), CTB (n=356), CTC (n=452), CTD (n=1,426), RPCI-1 (n=740), RPCI-3 (n=456), RPCI-4 (n=716), RPCI-5 (n=802) and RPCI-11 (n=16,009). (In a subsequent re-examination, we identified additional clones that we likely could have classified into libraries, including 953 from RPCI-11, 632 from RPCI-1, and 490 from another library RPCI-13.) The median span of the 21,156 clones we analyzed was 112 kb, and 80% are >50kb in size. About 2/3 came from a single library, RPCI-11.

1. **RPCI-11 is an African American:** RPCI-11, the individual who contributed most of the human genome reference sequence, is consistent with having African American ancestry, with 42% of the clones of confident West African ancestry and 42% of the clones of confident European ancestry, and the ancestry of the remaining clones less confidently inferred. The finding of likely African American ancestry for RPCI-11 was previously reported in a study of the ancestry of RPCI-11 clones spanning the Duffy blood group locus (S93), and here we confirm this finding, and also expand the inference to the whole genome.
2. **CTD is an East Asian:** The majority of clones from CTD, the second largest library in its contribution to the human genome sequence, is likely an East Asian. In a HAPMIX analysis with CEU (European) – CHB+JPT (East Asian) as the proposed ancestral populations, the majority of clones are of confident East Asian origin, and there is no secondary mode of confident European ancestry, as might be expected from a Latino or South Asian individual.
3. **The remaining 7 libraries are European:** The remaining libraries (CTA, CTB, CTC, RPCI-1, RPCI-3, RPCI-4 and RPCI-5) are inferred to be of European ancestry, since they all have consistent distributions of inferred clone ancestries, with the majority of clones of confident European ancestry in both our HAPMIX analyses and no secondary modes.

A Draft Sequence of the Neandertal Genome

Green et al (2010) Science. DOI: 10.1126/science.1188021
Supplemental Note 16 (pg 145-146)

Who is the reference human?

The screenshot shows the homepage of the **nature methods** journal. At the top, there's a dark banner with the journal title and a sub-headline: "Techniques for life scientists and chemists". On the right side of the banner, it says "Welcome back: Michael Schatz" with links for "Logout" and "Cart". Below the banner, there's a search bar with "Search" and "go" buttons, and a link to "Advanced search".

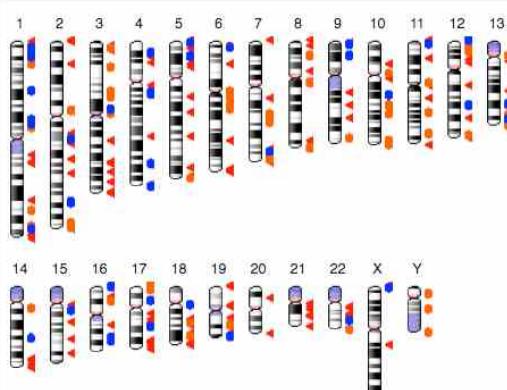
The main content area has a sidebar on the left with sections like "Journal content", "Archive", "Focuses and Supplements", "Methagora blog", "Method of the Year 2016", "Multimedia", and "Press releases". Another sidebar on the right lists "Article tools" such as "Download PDF", "Send to a friend", and "Export citation".

The central article is titled "**E pluribus unum**". It discusses the need for the human reference genome to reflect more of the actual genomic diversity in humans through community participation. It mentions the Human Reference Consortium (GRC) and the completion of GRCh37. It also compares GRCh37 to HuRef, another *de novo* assembly.

At the bottom of the article, there's a note about GRCh37 being a mosaic haploid genome derived from about 13 people, containing rare alleles, and the GRC's decision to convert these to common haplotypes.

Human Genome Overview

Information about the continuing improvement of the human genome



- ◀ Region containing alternate loci
- Region containing fix patches
- Region containing novel patches

Ideogram of the latest human assembly, GRCh38.p11

[GRCh38.p11](#)

[GRCh37.p13](#)

[GRCh37](#)

GRCh38.p11

Release date: June 14, 2017

Release type: minor

Release notes: GRCh38.p11 is the eleventh patch release for the GRCh38 reference assembly. No chromosome coordinate changes were made. The total number of patch scaffolds is now: 64 FIX and 59 NOVEL.

Assembly accessions: GenBank: [GCA_000001405.26](#), RefSeq: [GCF_000001405.37](#)

Pseudoautosomal regions

Name	Chr	Start	Stop
PAR#1	X	10,001	2,781,479
PAR#2	X	155,701,383	156,030,895
PAR#1	Y	10,001	2,781,479
PAR#2	Y	56,887,903	57,217,415

The GRC is working hard to provide the best possible by both generating multiple representations (alternatively represented by a single path. Additionally, we are re allows users who are interested in a specific locus to affecting users who need chromosome coordinate s.

Download data:

- [GRCh38.p11 \(latest minor release\) FTP](#)
- [GRCh38 \(latest major release\) FTP](#)
- [Genomic regions under review FTP](#)
- [Current Tiling Path Files \(TPFs\)](#)

Transitioning to GRCh38? Try the [NCBI Remapping](#) assembly alignments used by the GRC.

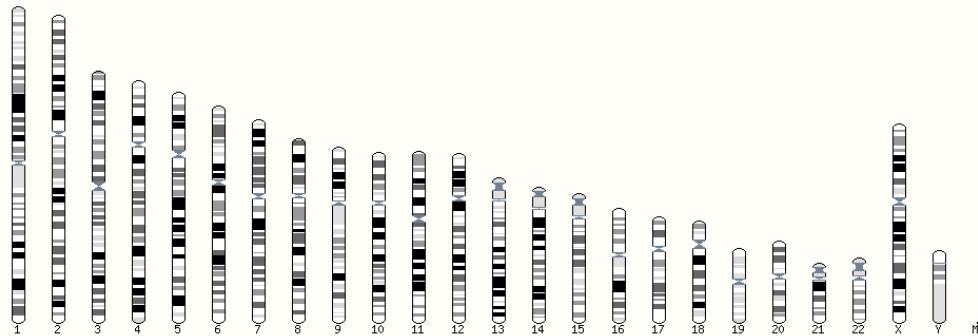
Next assembly update

The next assembly update (GRCh38.p12) will be



©HumanWorlds.com

The human genome - basic stats

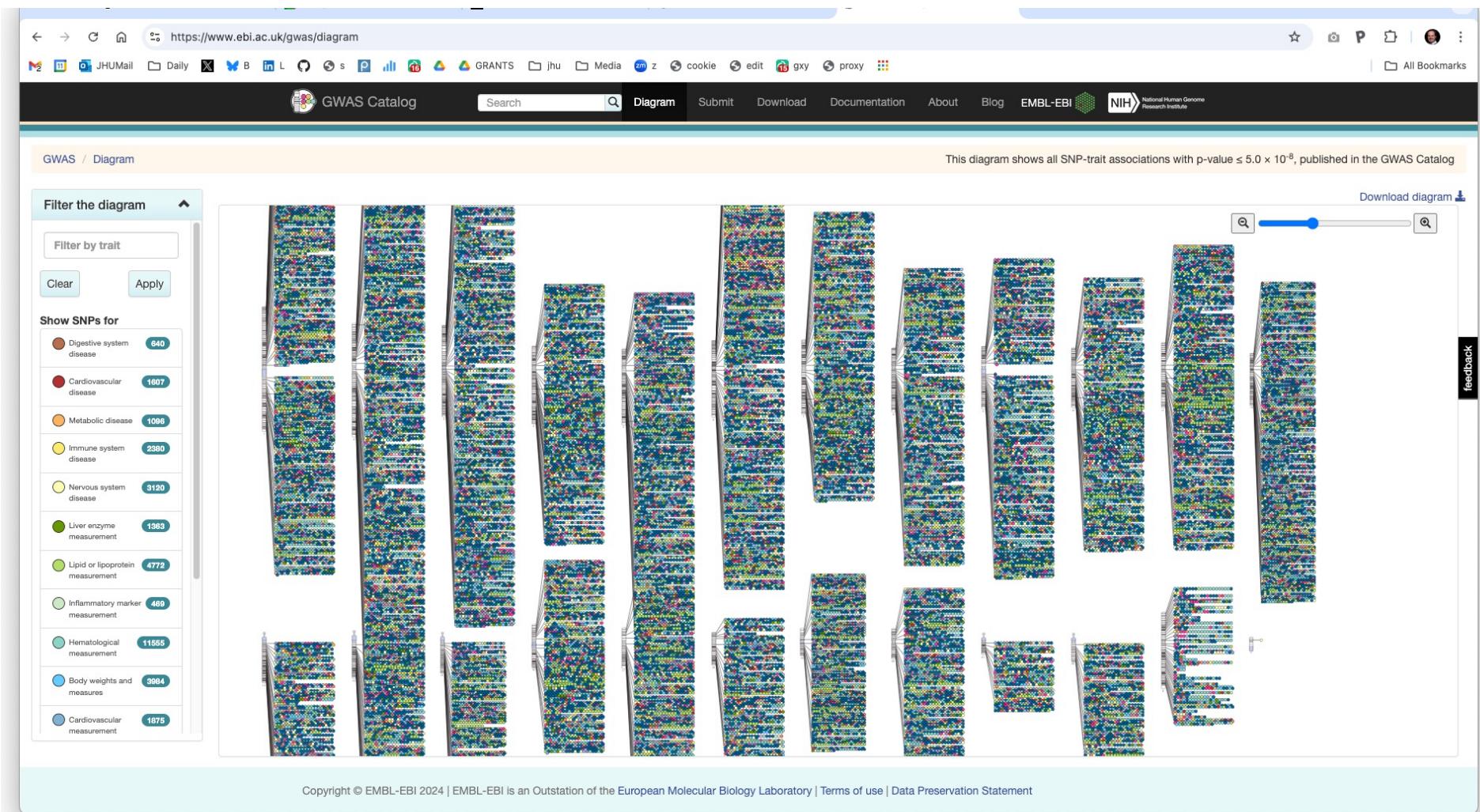


- 3.096 billion base pairs (haploid)
- 20,454 protein coding genes
- 226,950 coding transcripts
(isoforms of a gene that each encode a distinct protein product)

Assembly	GRCh38.p12 (Genome Reference Consortium Human Build 38), INSDC Assembly GCA_000001405.27 , Dec 2013
Base Pairs	3,609,003,417
Golden Path Length	3,096,649,726
Annotation provider	Ensembl
Annotation method	Full genebuild
Genebuild started	Jan 2014
Genebuild released	Jul 2014
Genebuild last updated/patched	Mar 2019
Database version	97.38
Gencode version	GENCODE 31

Gene counts (Primary assembly)

Coding genes	20,454 (incl 660 readthrough)
Non coding genes	23,940
Small non coding genes	4,871
Long non coding genes	16,848 (incl 302 readthrough)
Misc non coding genes	2,221
Pseudogenes	15,204 (incl 8 readthrough)
Gene transcripts	226,950



- As of 2024-08-12, the GWAS Catalog contains 6960 publications, 668514 top associations and 97045 full summary statistics.
- GWAS Catalog data is currently mapped to Genome Assembly GRCh38.p14 and dbSNP Build 156.

SPECIAL REPORT

The untold story of the Human Genome Project: How one man's DNA became a pillar of genetics

By Ashley Smart — Undark July 9, 2024



Pieter De Jong, who led the Human Genome Project work at Roswell Park Cancer Institute, stores DNA samples from the anonymous donor known as RP11 (or RPCI-11) at minus 80 degrees Celsius in a freezer at his home in Redmond, Wash.

JOVELLE TAMAYO FOR STAT



Reprints

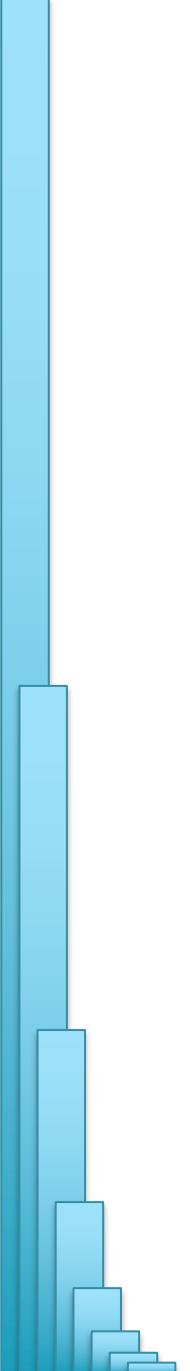
STAT is co-publishing this [investigation by Undark](#).

hey numbered 20 in all — 10 men and 10 women who came to a sprawling medical campus in downtown Buffalo, N.Y., to volunteer for

1/3 FREE STORIES

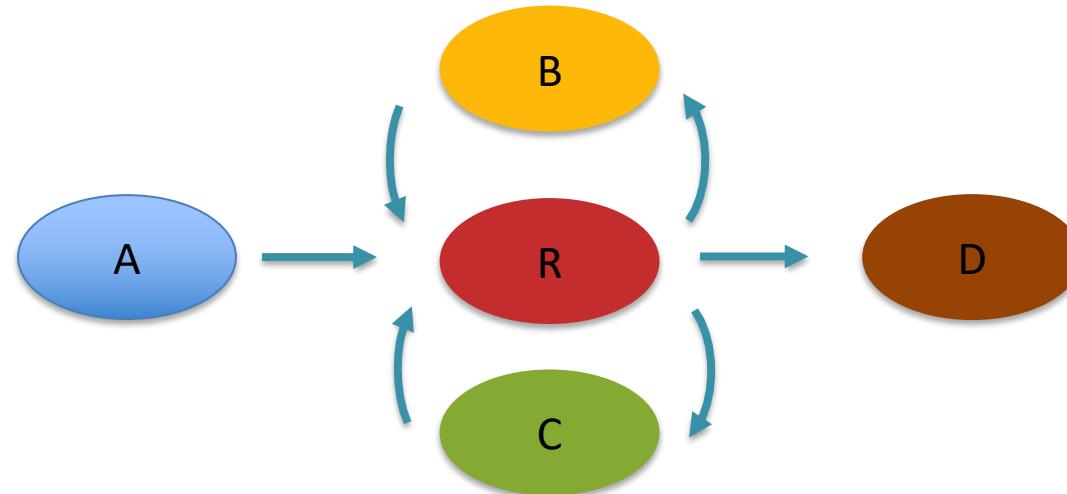
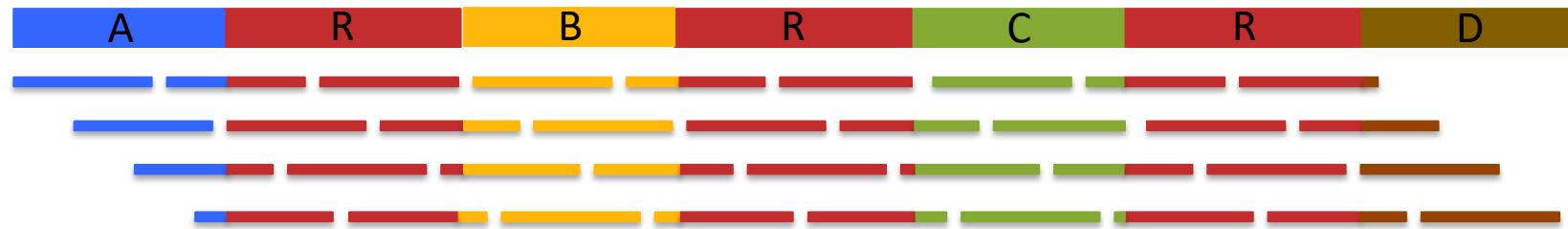
Already have an account? [Log in](#)

<https://www.statnews.com/2024/07/09/human-genome-project-untold-story-how-single-volunteer-became-genetics-foundation/>

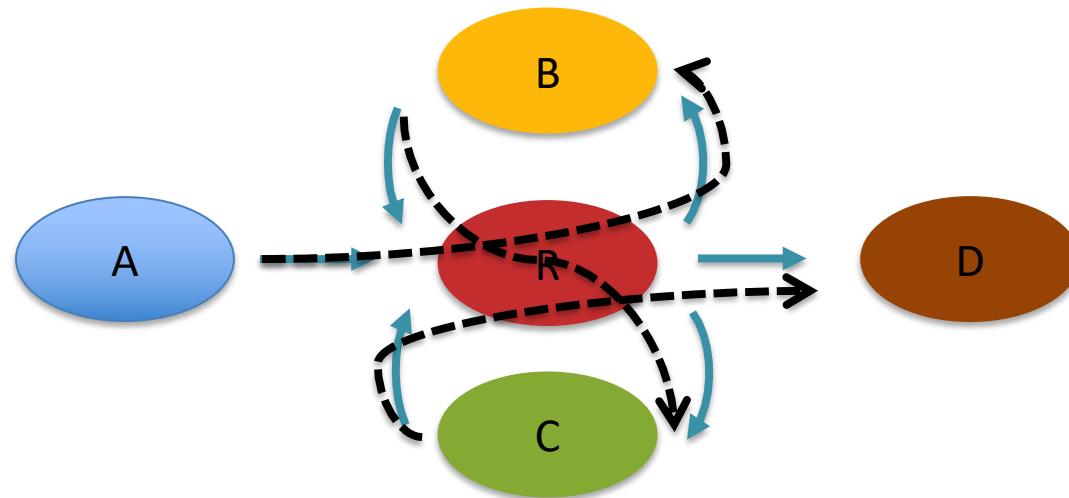


Part II: Long Read Assembly

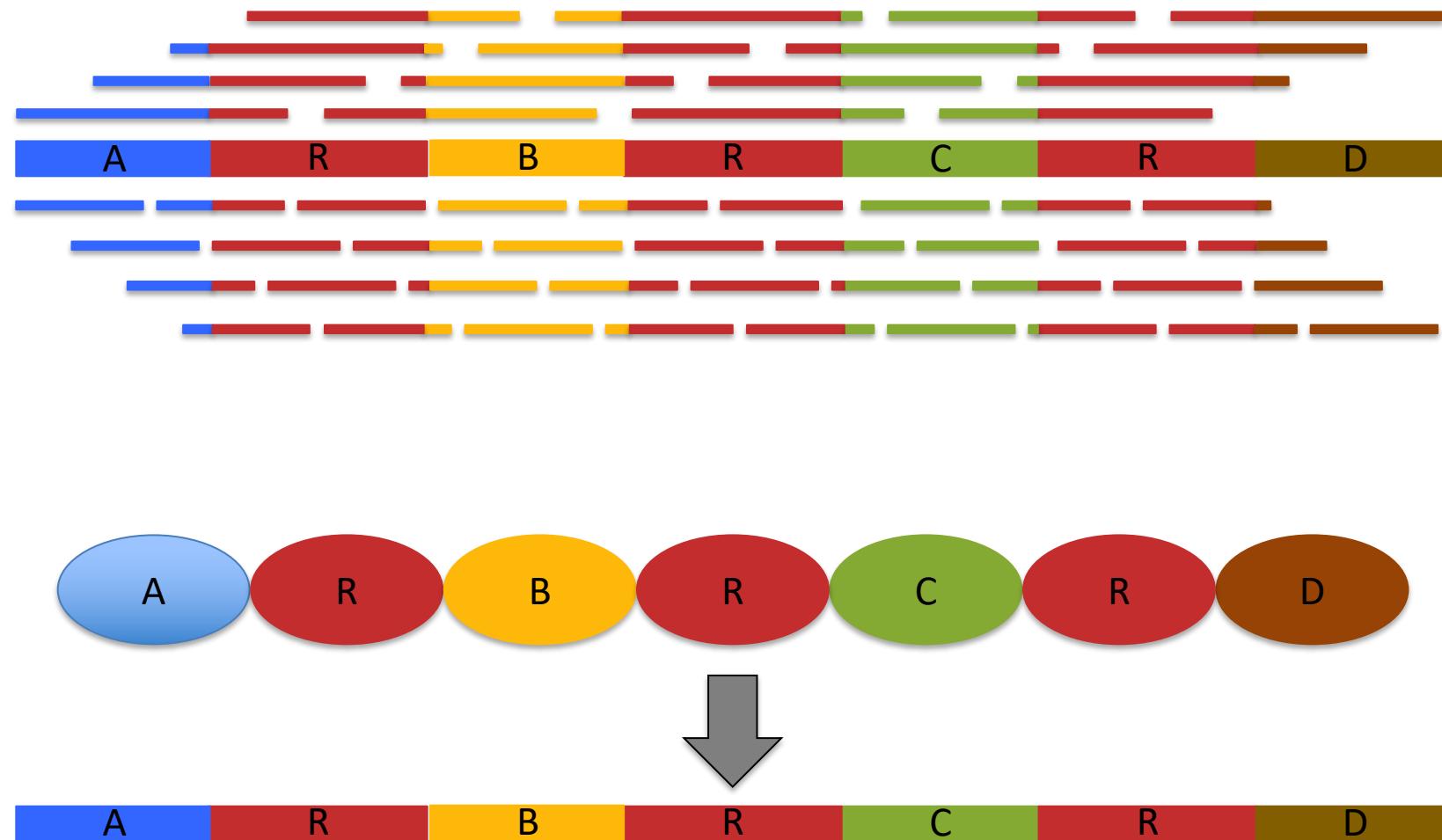
Assembly Complexity



Assembly Complexity



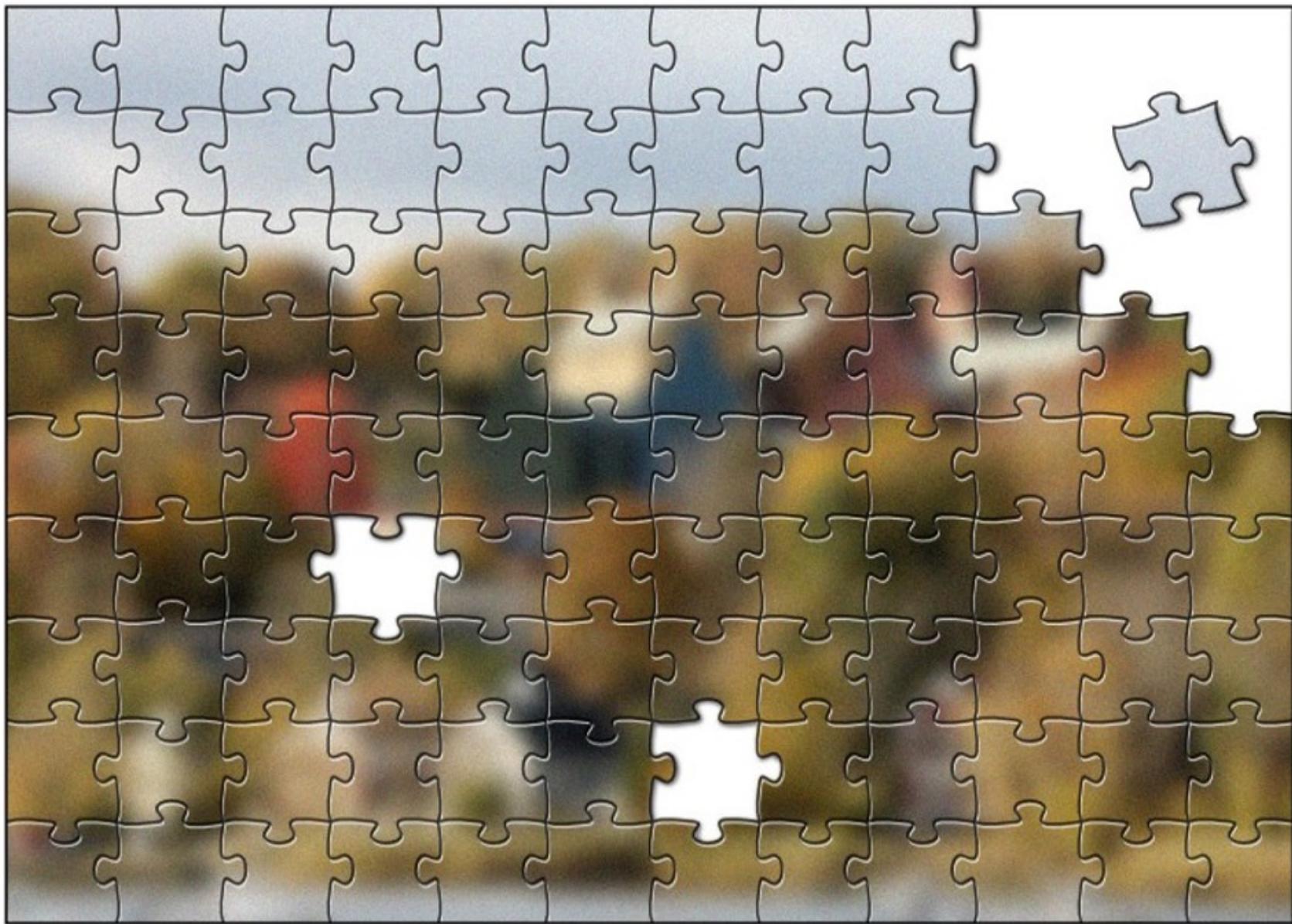
Assembly Complexity



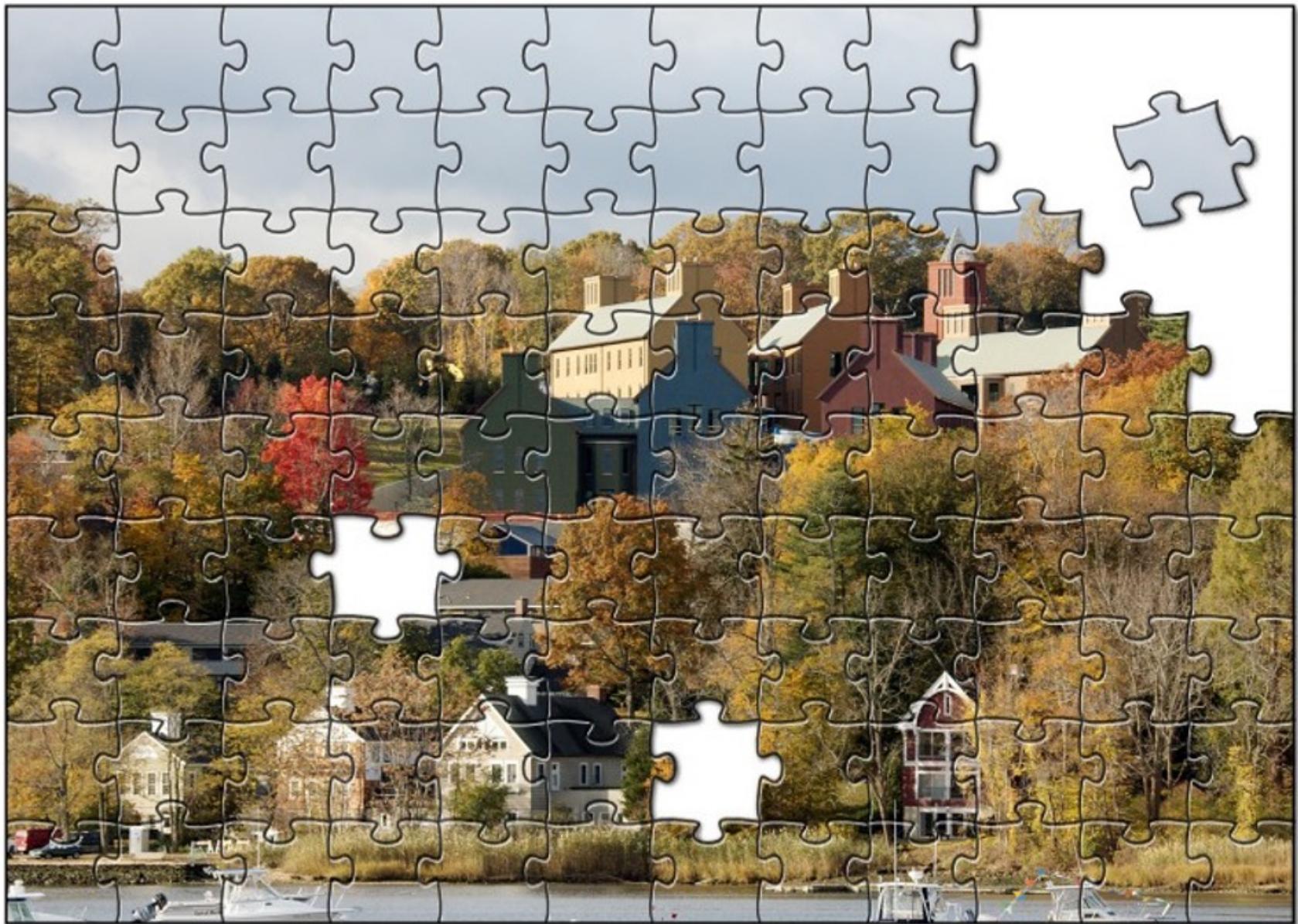
The advantages of SMRT sequencing

Roberts, RJ, Carneiro, MO, Schatz, MC (2013) *Genome Biology*. 14:405

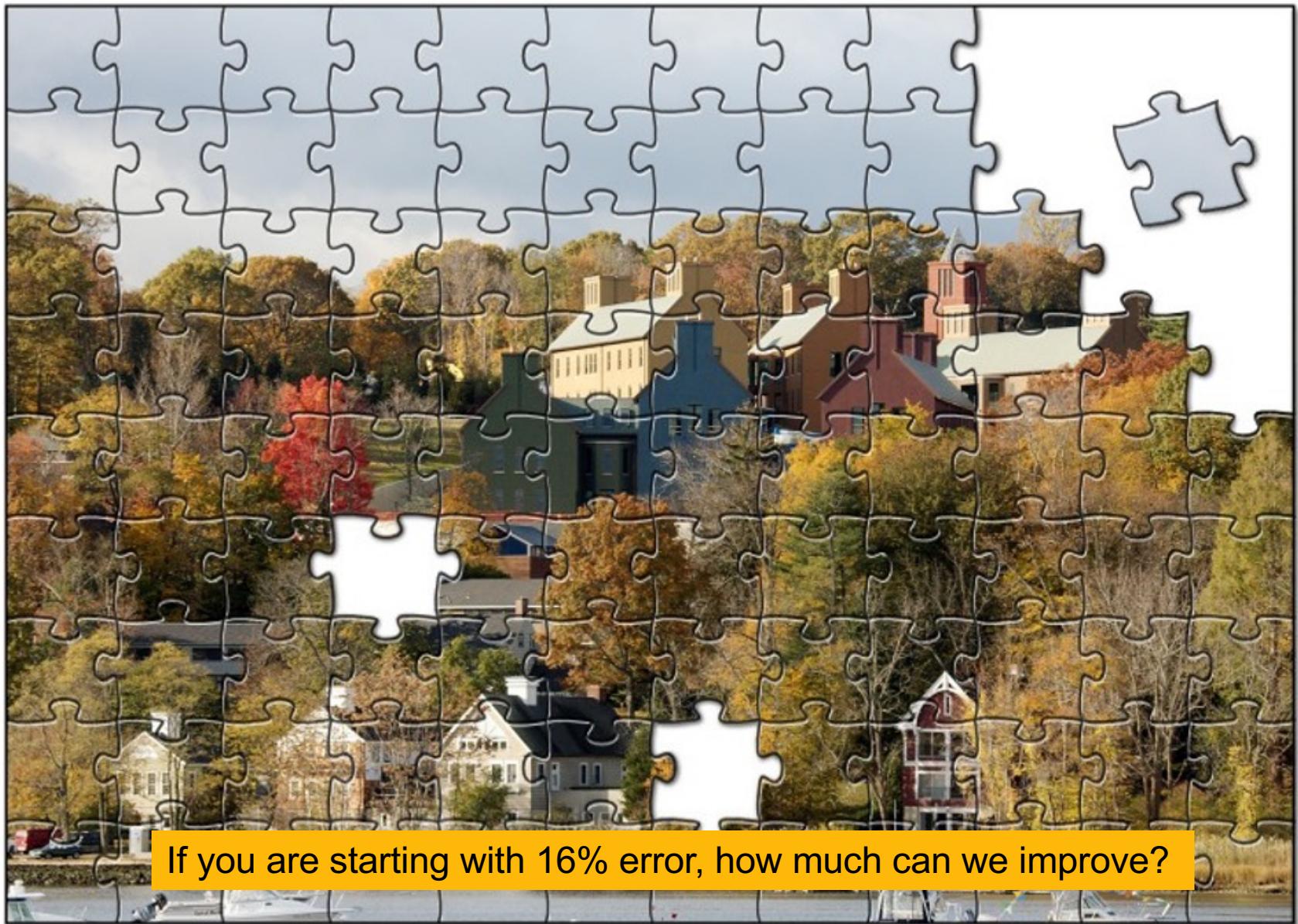
Single Molecule Sequences



“Corrective Lens” for Sequencing

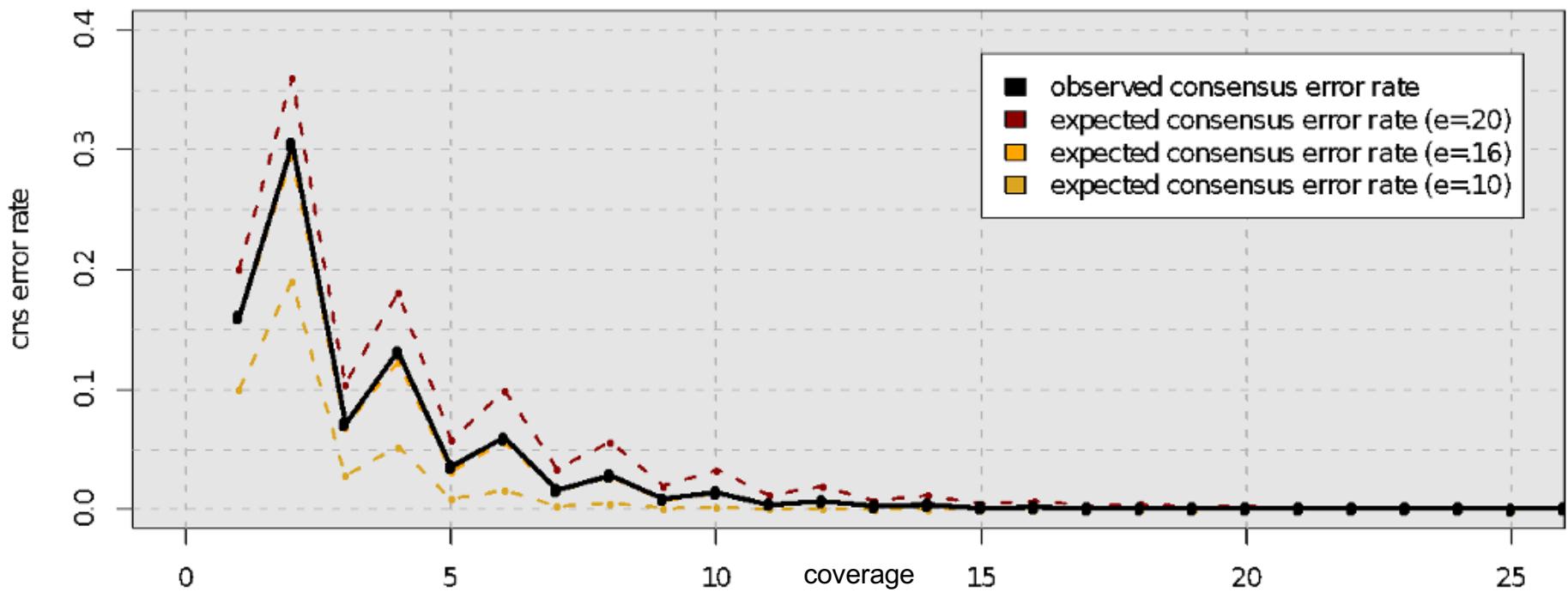


“Corrective Lens” for Sequencing



If you are starting with 16% error, how much can we improve?

Consensus Accuracy and Coverage



Coverage can overcome random errors

- Dashed: error model from binomial sampling; solid: observed accuracy
- For same reason, CCS is extremely accurate when using 5+ subreads

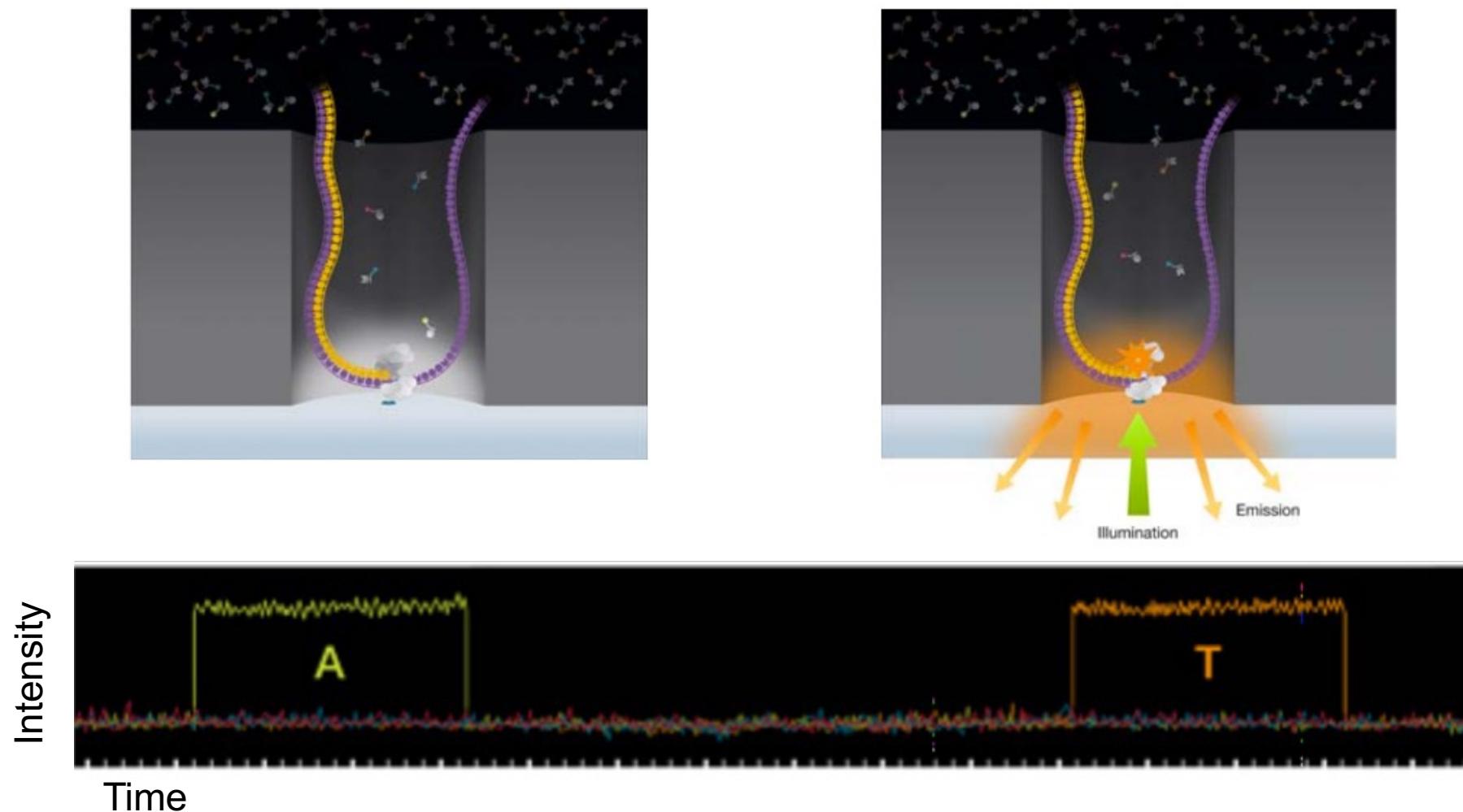
$$CNSError = \sum_{i=\lceil c/2 \rceil}^c \binom{c}{i} (e)^i (1-e)^{n-i}$$



PacBio Single
Molecule Real Time
Sequencing
(SMRT-sequencing)

PacBio: SMRT Sequencing

Imaging of fluorescent phospholinked labeled nucleotides as they are incorporated by a polymerase anchored to a Zero-Mode Waveguide (ZMW).



https://www.youtube.com/watch?v=_ID8JyAbwEo

“HiFi” Circular Consensus Reads

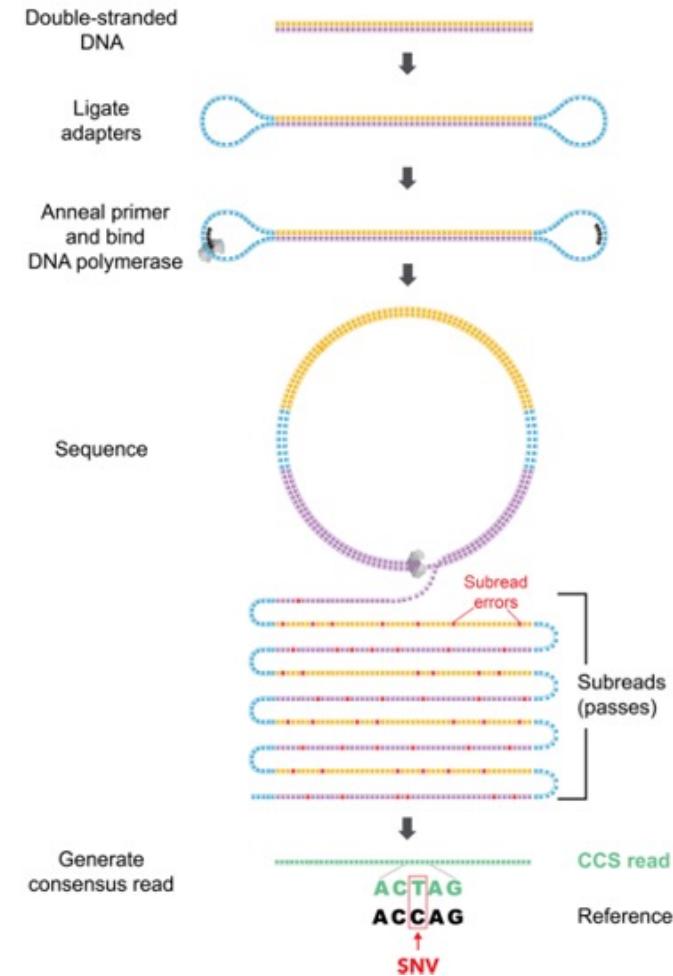
High-quality reads produced by sequencing the same molecule multiple times

Higher accuracy for low-coverage sequences like somatic variants or lowly expressed transcripts in RNA-seq, more interpretable alignments, better & faster assembly

Limits read length, used to be very expensive but more manageable now

Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome

Wenger et al (2019) Nature Biotechnology doi:10.1038/s41587-019-0217-9

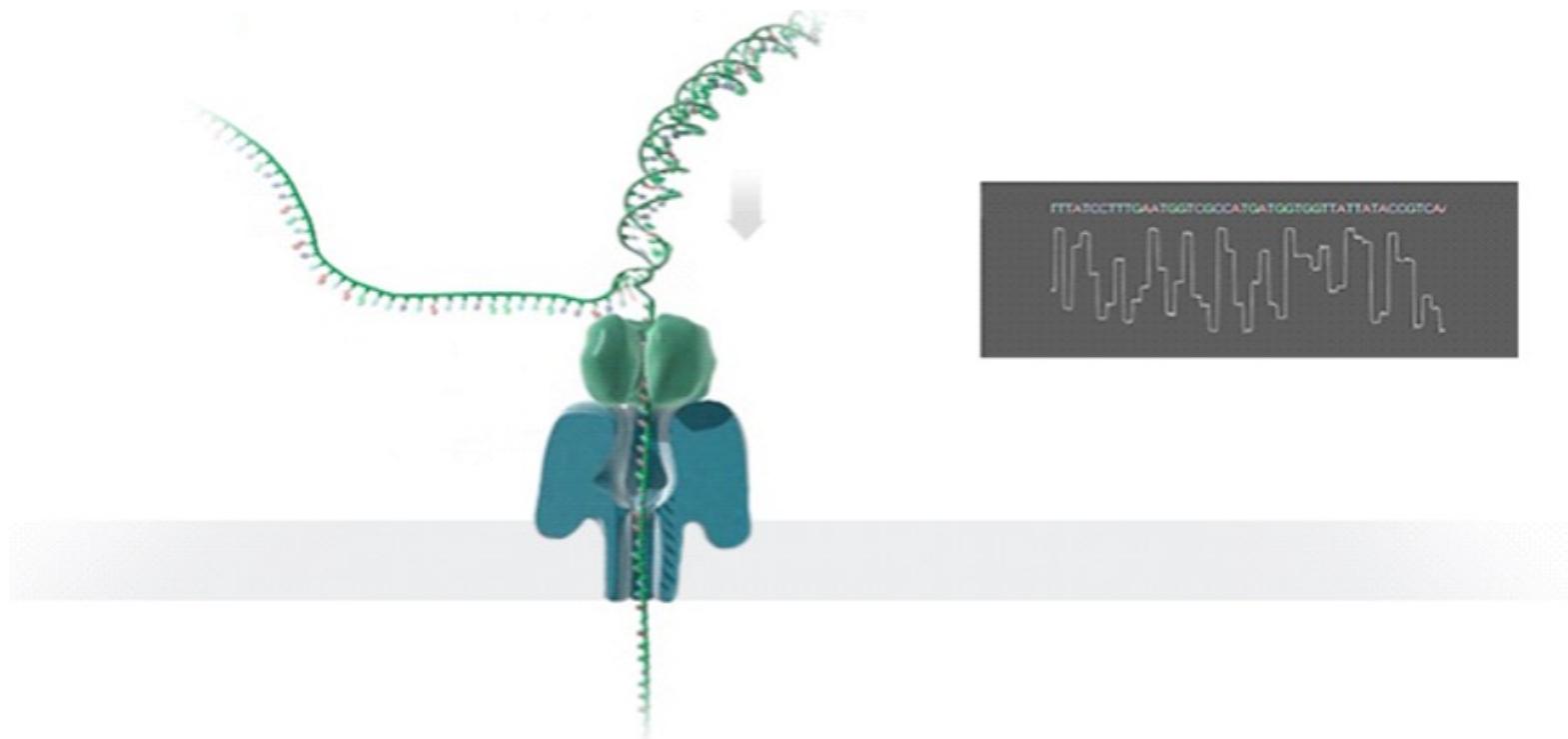


Oxford Nanopore Technologies (ONT)



Nanopore Sequencing

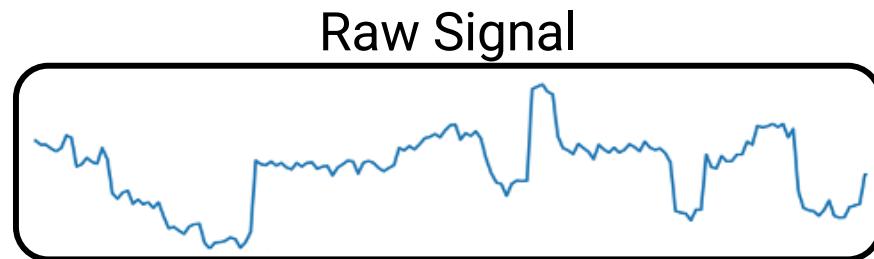
Sequences DNA/RNA by measuring changes in ionic current as nucleotide strand passes through a pore



<https://www.youtube.com/watch?v=RcP85JHLmnI>

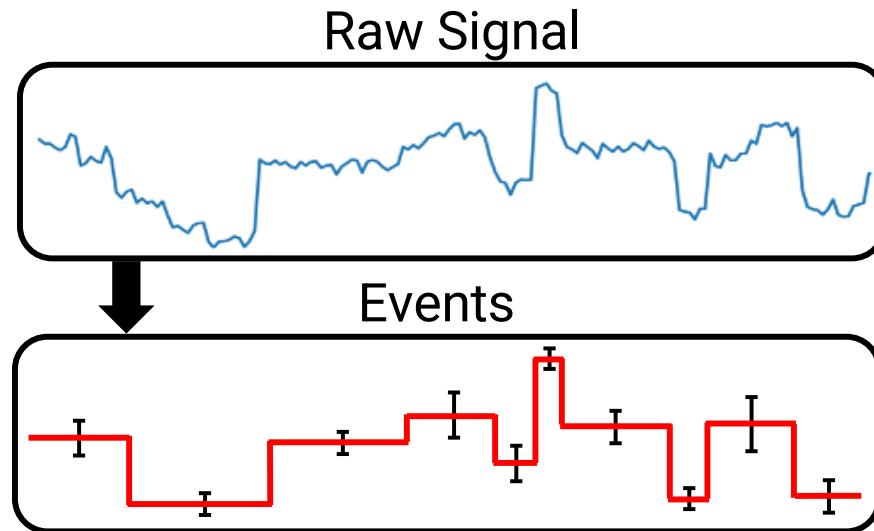
nanoporetech.com/applications/dna-nanopore-sequencing

Nanopore Basecalling



Translation of raw signal
into basepairs

Nanopore Basecalling

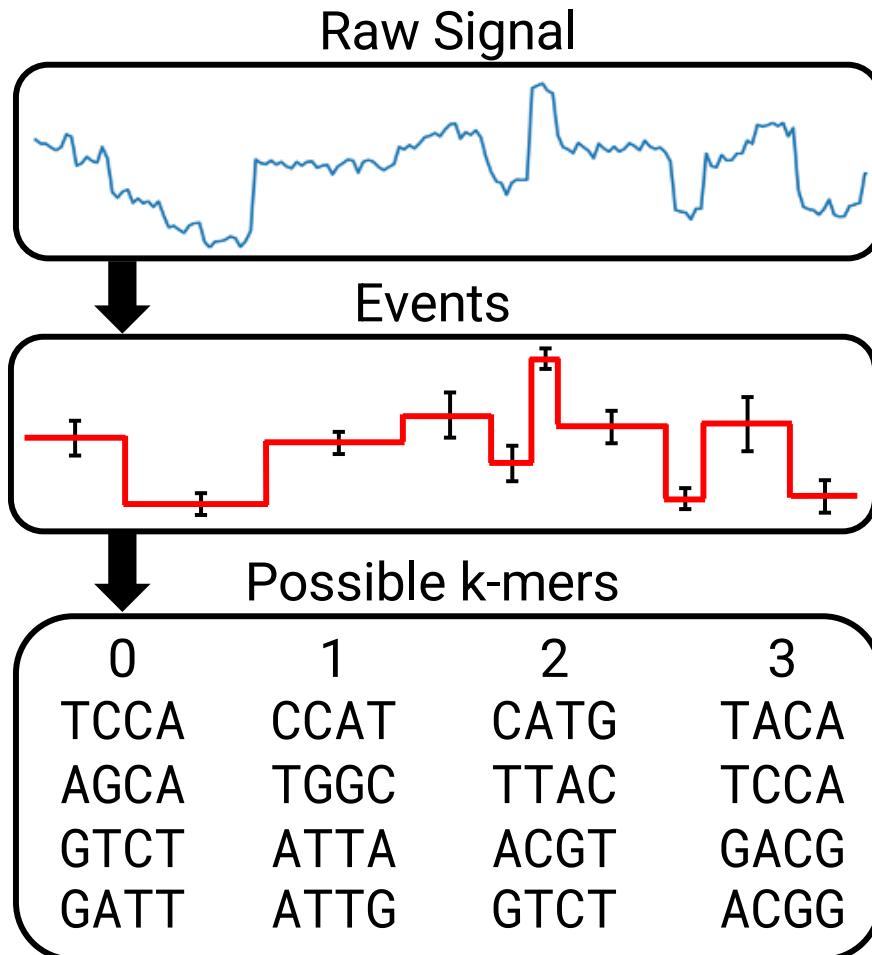


Translation of raw signal
into basepairs

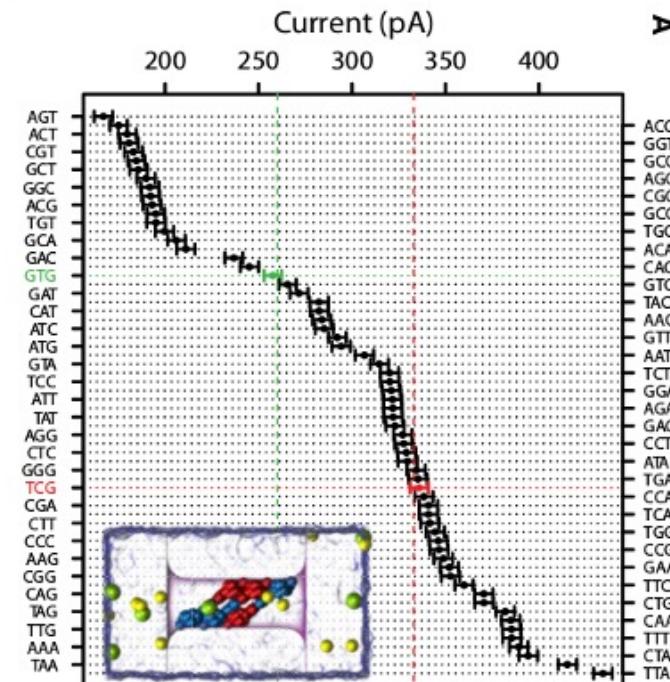
Early basecallers began by
estimating k-mer boundaries
using “events”, which were
then input to an HMM

Modern basecallers use
neural networks directly
on raw signal

Nanopore Basecalling

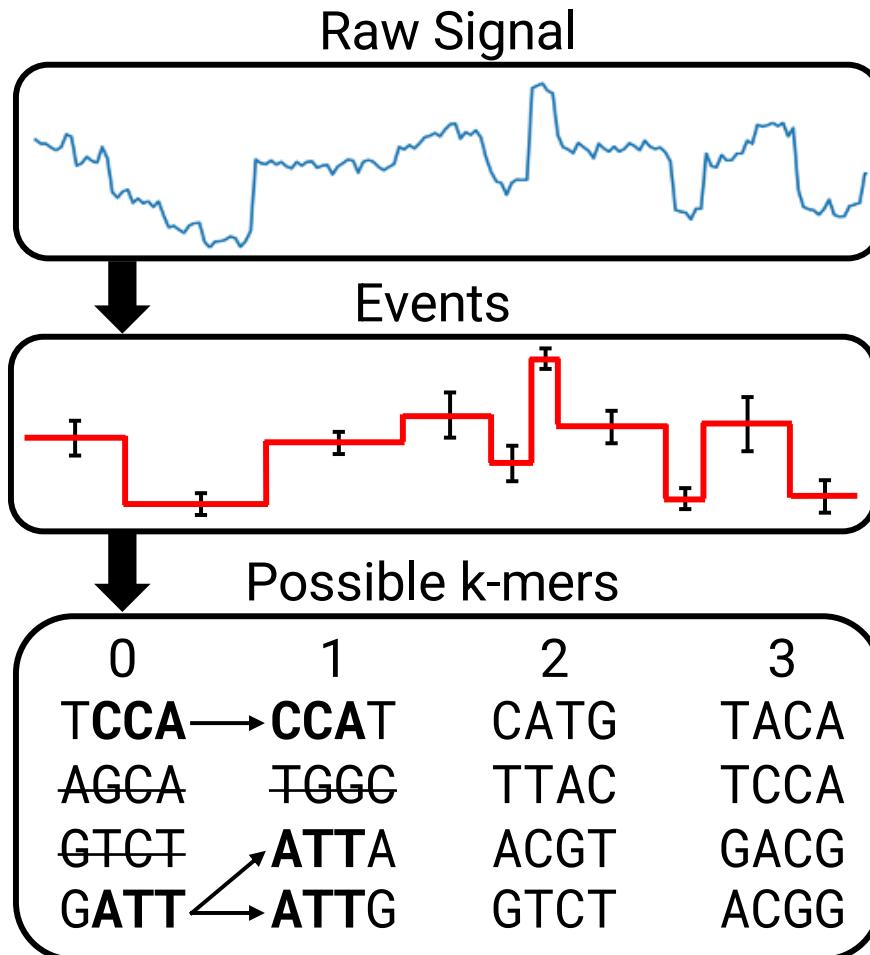


ONT releases k-mer models with expected current distribution of every k-mer

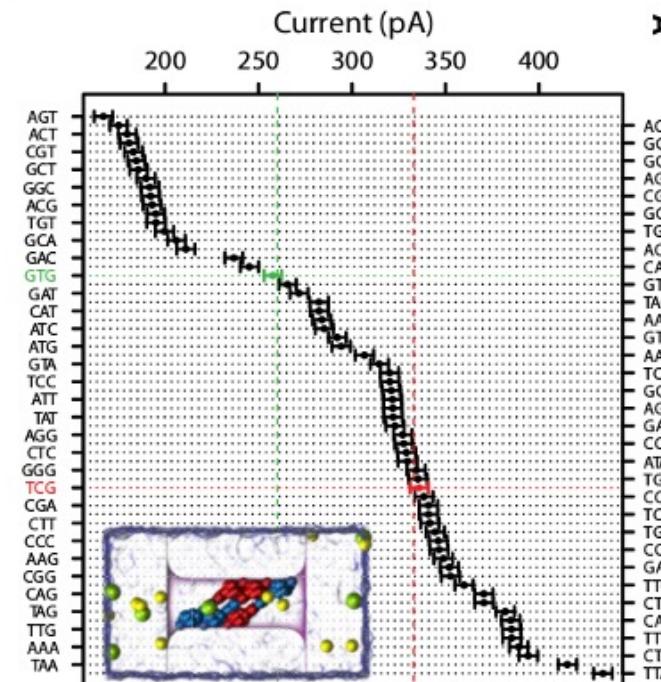


DNA Base-Calling from a Nanopore Using a Viterbi Algorithm
Timp et al. (2012) *Biophysical Journal*

Nanopore Basecalling

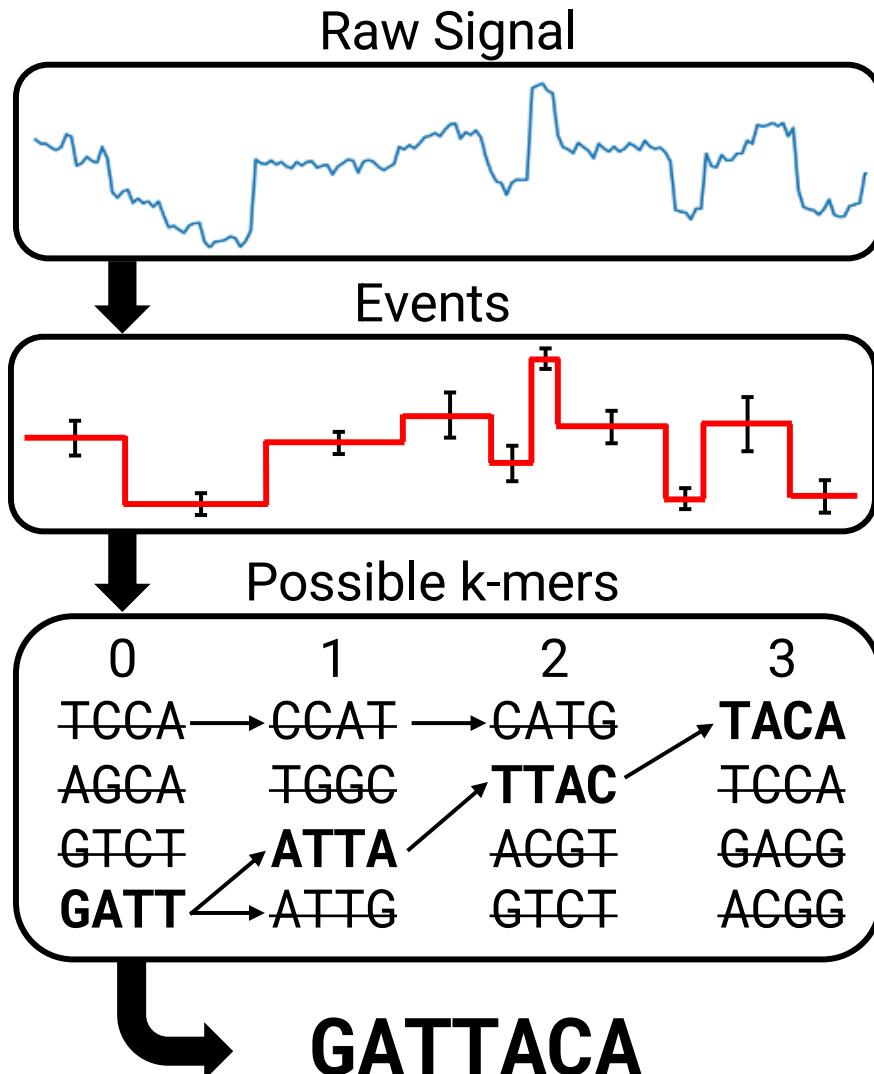


Certain k-mers can be eliminated based on possible transitions

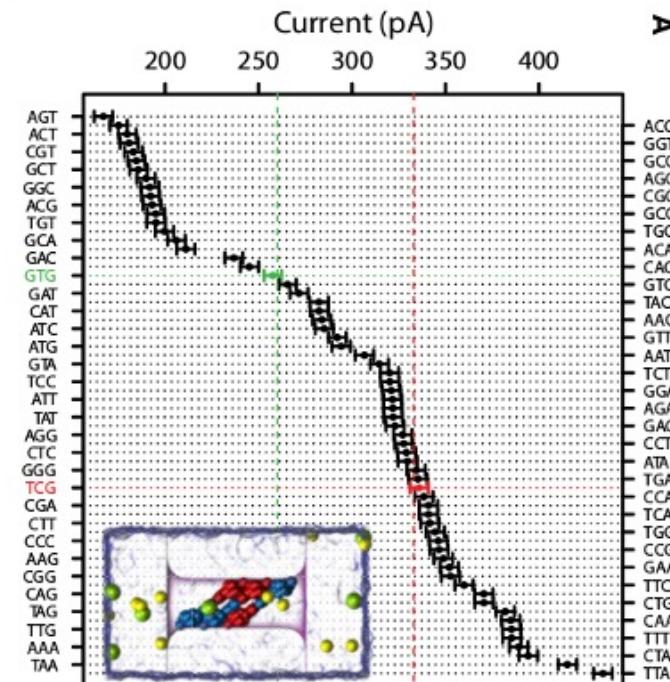


DNA Base-Calling from a Nanopore Using a Viterbi Algorithm
Timp et al. (2012) *Biophysical Journal*

Nanopore Basecalling



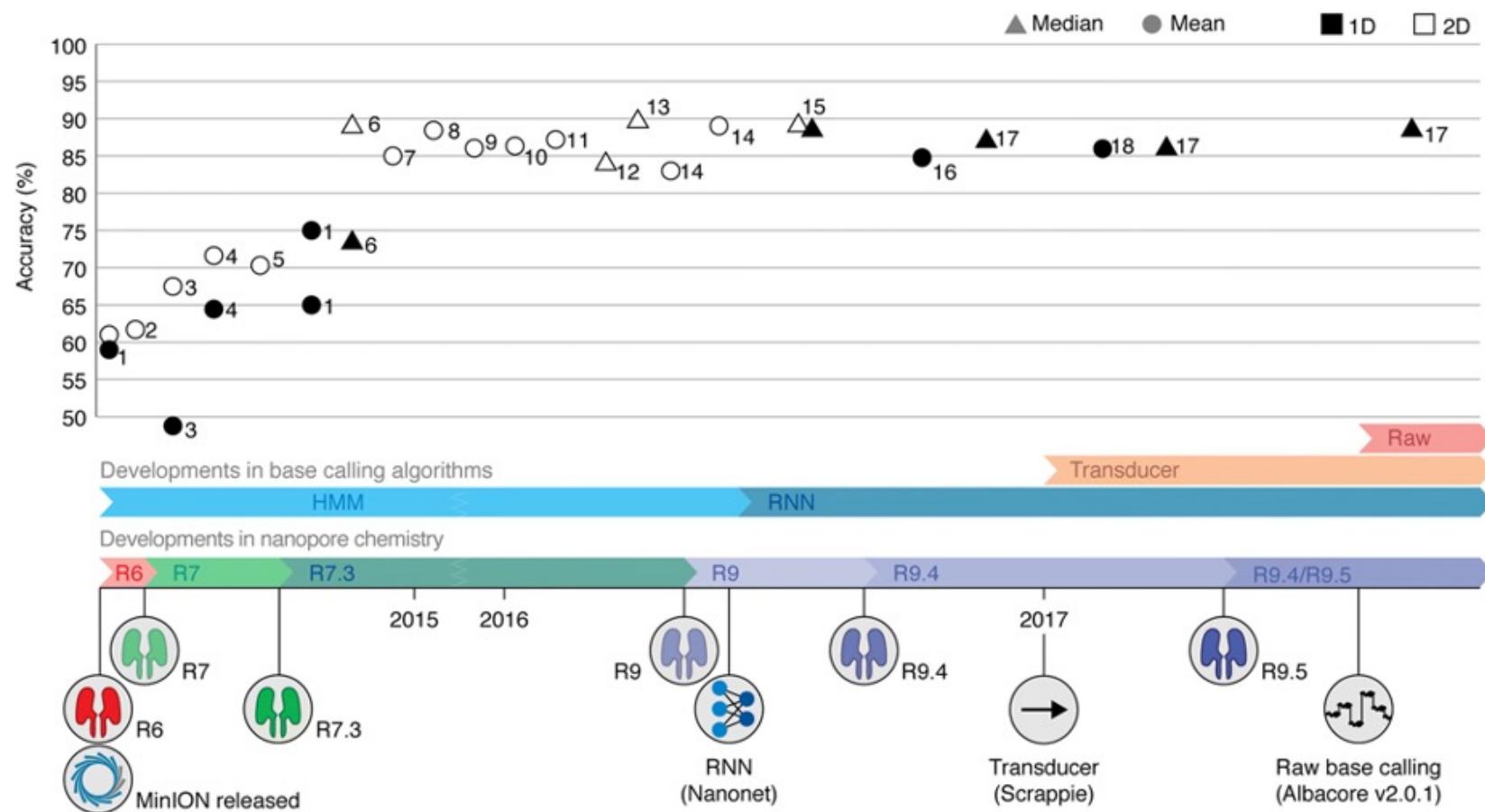
Final sequence determined by
most probable k-mers



"DNA Base-Calling from a Nanopore Using a Viterbi Algorithm"
Timp et al. (2012) *Biophysical Journal*

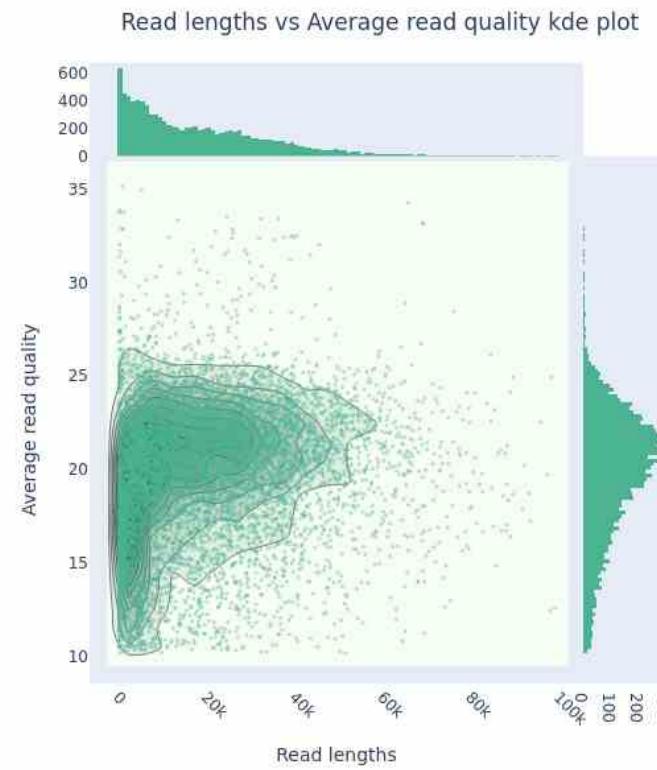
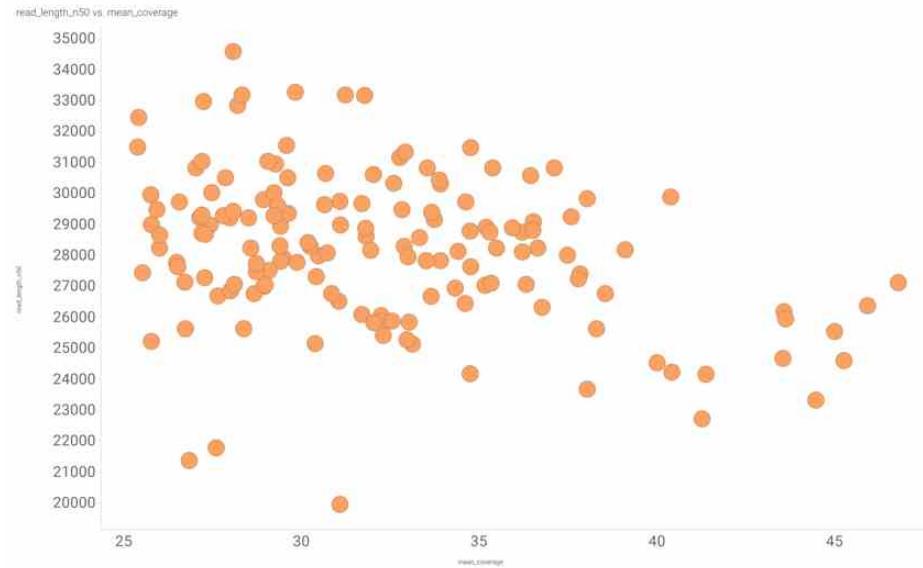
Basecaller/Pore Timeline

Development of both pore chemistry and basecalling algorithms is responsible for improvement in accuracy



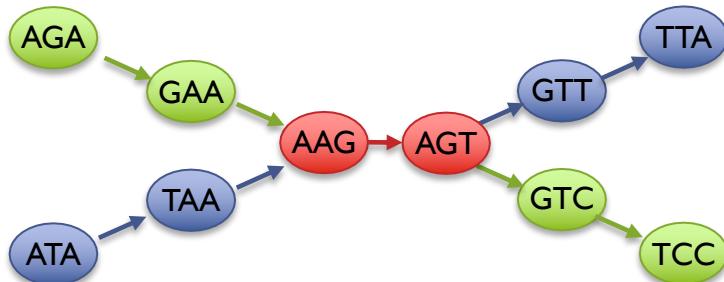
From squiggle to basepair: computational approaches for improving nanopore sequencing read accuracy
Rang et al (2018) *Genome Biology*. <https://doi.org/10.1186/s13059-018-1462-9>

ONT at JHU



Two Paradigms for Assembly

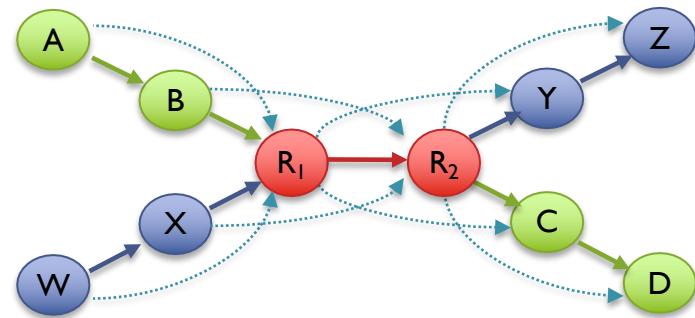
de Bruijn Graph



Short read assemblers

- Repeats depends on word length
- Read coherency, placements lost
- Robust to high coverage

Overlap Graph

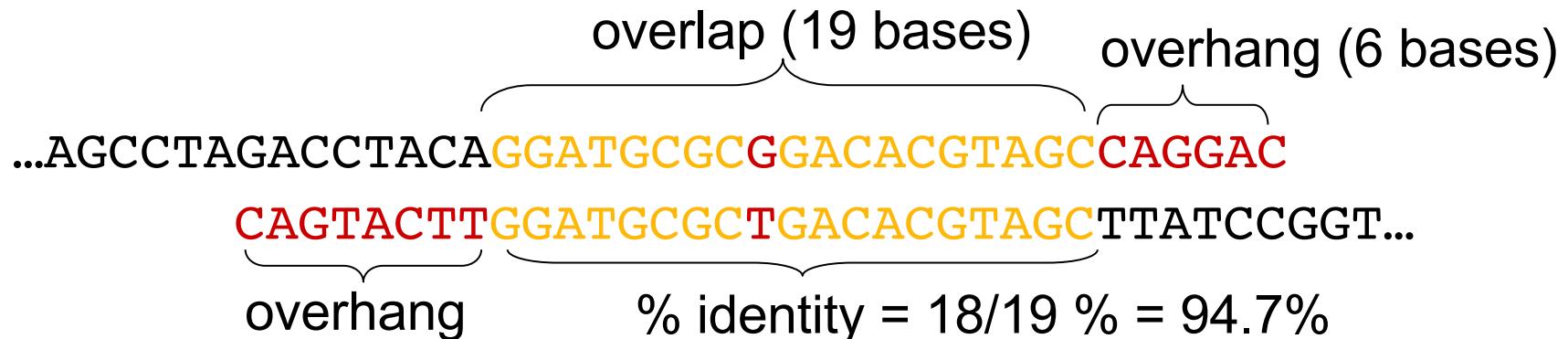


Long read assemblers

- Repeats depends on read length
- Read coherency, placements kept
- Tangled by high coverage

Assembly of Large Genomes using Second Generation Sequencing
Schatz MC, Delcher AL, Salzberg SL (2010) Genome Research. 20:1165-1173.

Overlap between two sequences



overlap - region of similarity between regions

overhang - un-aligned ends of the sequences

The assembler screens merges based on:

- length of overlap
- % identity in overlap region
- maximum overhang size.

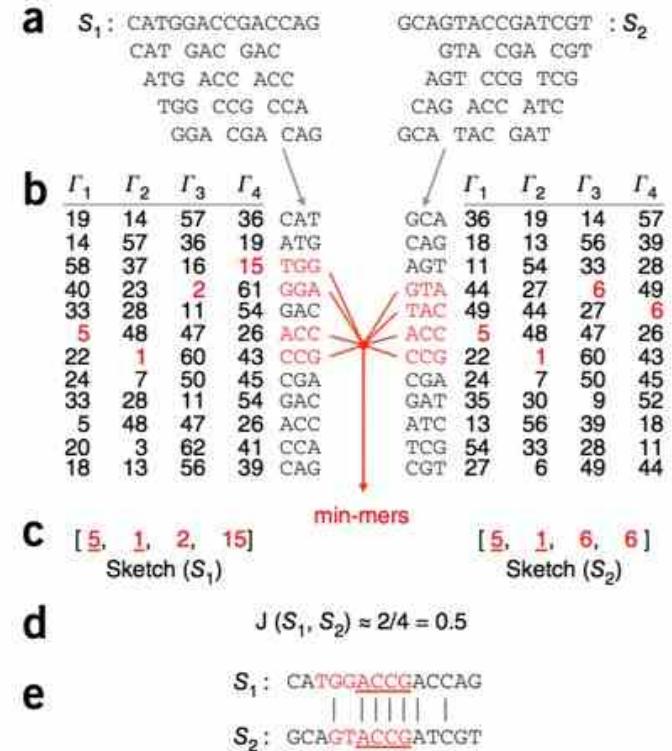
[How do we compute the overlap?]

[Do we really want to do all-vs-all?]

Very fast approximate overlapping

Maybe we don't need to compute the exact identity of the overlap region, just approximate it

- If two reads overlap, they should share many of the same kmers: Their Jaccard coefficient should be high: $|\text{intersection}| / |\text{union}|$
- But tracking all of the kmers for a read is a lot of overhead
- Instead, compare the “sketch” of the reads: a small fraction of kmers carefully chosen
- LSH: Find the sketch by applying N hash functions to the kmers, and keeping the minimum hash values reported from each (N=4 in example)
- This forms a nice “random” sample of the reads, and the Jaccard coefficient is a good approximation of the sequence similarity



Assembling large genomes with single-molecule sequencing and locality-sensitive hashing
Berlin et al (2015) *Nature Biotechnology*

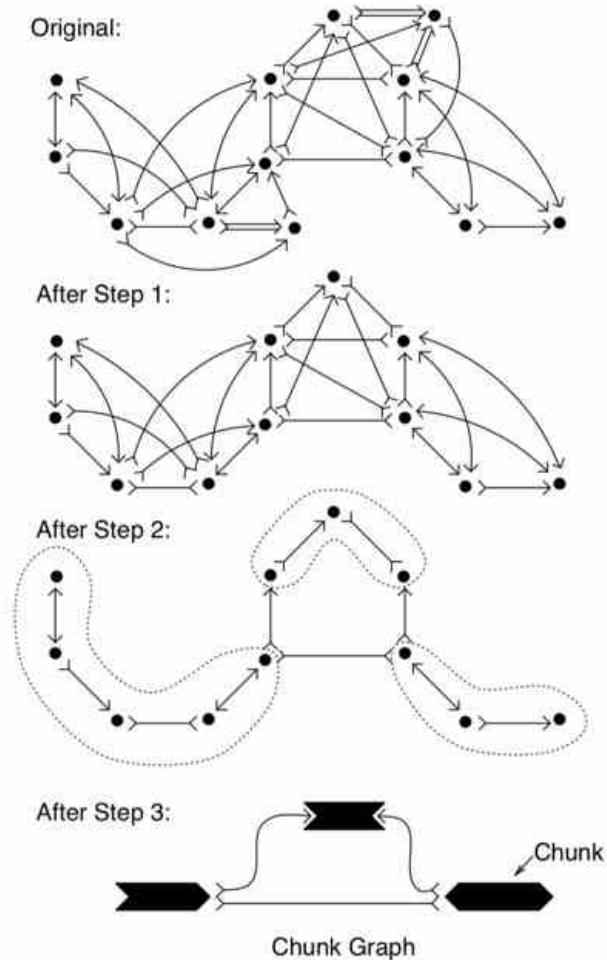
Unitigging: Pruning the Overlap Graph

The overlap graph has many redundant edges:

- If the average coverage is D, we should expect D overlaps at the beginning of the read, and D at the end

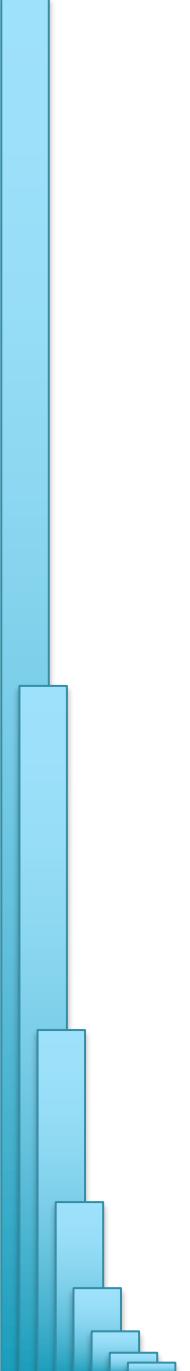
Transform the graph to simplify the assembly problem (without changing the valid solutions):

1. **Contained reads removal:** Short reads that are substrings of longer reads don't advance the assembly, remove those nodes and all of the edges
2. **Transitive edge removal:** If A \rightarrow B, and B \rightarrow C, remove the transitive edge A \rightarrow C
3. **“Chunkification”:** Linear subgraphs define uniquely assemblable segments: “unitigs”



Towards Simplifying and Accurately Formulating Fragment Assembly

Myers (1995) J Comput Biol. Summer;2(2):275-90.

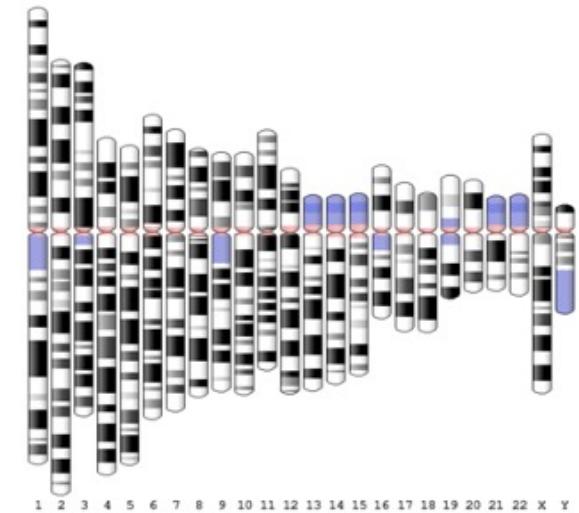


Part II: The T2T human genome

Finishing the human genome

238Mbp is missing or incorrect

- Centromeres and telomeres
- Segmentally duplicated genes
- Tandem gene arrays (e.g. rDNAs)
- And an unknown number of errors...



Why does it matter?

Variation in these regions is unexplored

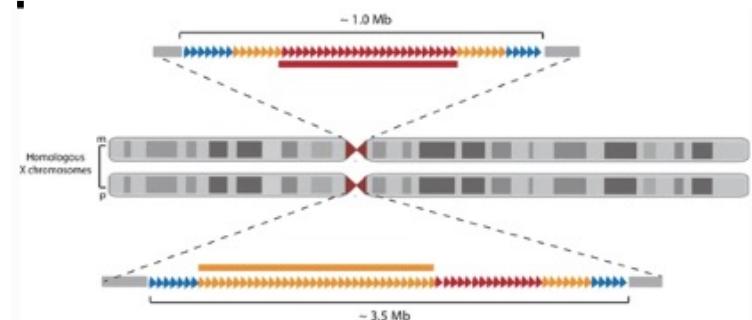
Functional studies need sequence

Reference gaps lead to artifacts

We don't know what we don't

Why has it taken so long?

Repeats, repeats, repeats...



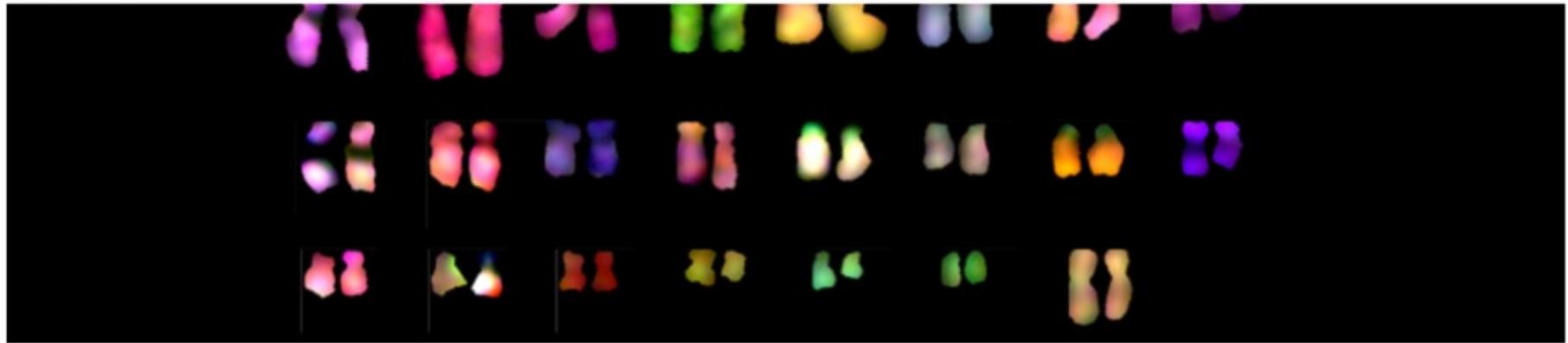
Miga 2015

Let's finish a human genome



T2T Working Group

[Home](#) · [Technology](#) · [Data](#) · [CHM13 Cell Line](#) · [Remaining Challenges](#) ▾ · [Who We Are](#) · [Join Us](#) 



The Telomere-to-Telomere (T2T) consortium is an open, community-based effort to generate the first complete assembly of a human genome.

CHM13 homozygous 46,XX cell line from Urvashi Surti, Pitt; SKY karyotype from Jennifer Gerton, Stowers



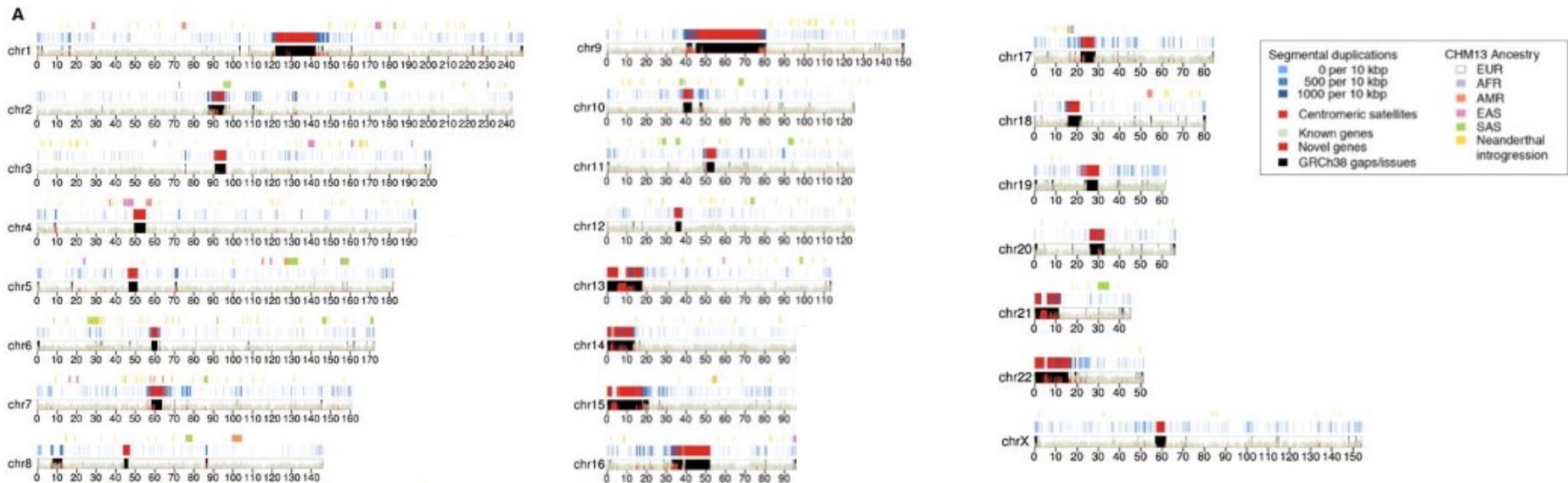
CHM13 assembly graph



HiCanu: accurate assembly of segmental duplications, satellites, and allelic variants from high-fidelity long reads.

Nurk et al. *Genome Research* (2020)

The complete sequence of a human genome



CHM13v1.1 genome size is **3.057 Gbp with zero Ns**
Every chromosome is telomere-to-telomere, quality estimated >Q70
~190 Mbp (~8%) of new sequence vs. GRCh38, fixes thousands of errors

(Nurk et al. Science, 2022)



NEWS CAREERS COMMENTARY JOURNALS ▾

Science

HOME > COLLECTIONS > COMPLETING THE HUMAN GENOME

COMPLETING THE HUMAN GENOME

A fully sequenced human genome was announced more than 20 years ago. However, owing to technological limitations, some genomic regions remained unresolved. Here, *Science* and other journals present research by the Telomere-to-Telomere (T2T) Consortium, reporting on the endeavor to complete a comprehensive human reference genome.

FILTERS

6 RESULTS FOUND

SPECIAL ISSUE RESEARCH ARTICLE

Segmental duplications and their variation in a complete human genome

BY MITCHELL R. VOLLSER, XAVI GUTIÁR, PHILIP C. DISHICK, LUDOVICA MEROLI, WILLIAM T. HARVEY, ARIEL GERSHMAN, MARK DICKHANS, ARVIS SULONARI, KATHERINE M. MUNDON, ALEXANDRA P. LEWIS, [...] EVAN E. EICHLER

SCIENCE • VOL. 376, NO. 6588 • 81 APR 2022

SPECIAL ISSUE RESEARCH ARTICLE

Complete genomic and epigenetic maps of human centromeres

BY NICOLAS ALTEMOSSE, GLENNIS A. LOGGISON, ANDREY V. BZKADZE, PRADYA SIDHWANI, SASHA A. LANGLEY, GINA V. CALDAS, SAVANNAH J. HOYT, LEV URALSKY, FEDOR D. RYABIKO, COLIN J. SHEK, [...] KAREN H. MIGA

SCIENCE • VOL. 376, NO. 6588 • 81 APR 2022

SPECIAL ISSUE RESEARCH ARTICLE

From telomere to telomere: The transcriptional and epigenetic state of human repeat elements

BY SAVANNAH J. HOYT, JESSICA M. STOREY, GABRIELLE A. MARTLEY, PATRICK G. S. GRADY, ARIEL GERSHMAN, LEONARDO G. DE UMA, CHARLES LIMOUSE, REZA HILARIOU, LIKE MOJENSKI, MATIAS RODRIGUEZ, [...] RACHEL J. COOPER, +16 authors

SCIENCE • VOL. 376, NO. 6588 • 81 APR 2022

SPECIAL ISSUE RESEARCH ARTICLE

A complete reference genome improves analysis of human genetic variation

BY SERGEY AGAFONOV, STEPHANIE M. YAN, DANIELA C. SOTO, MÉLANIE KIRSCH, SAMANTHA ZAKRZE, PAVEL AVDEYEV, DYLAN J. TAYLOR, KISHWAR SHAFIN, ALAINA SHUMATE, CHUNLIN XIAO, [...] MICHAEL C. SCHATZ

SCIENCE • VOL. 376, NO. 6588 • 81 APR 2022

SPECIAL ISSUE RESEARCH ARTICLE

Epigenetic patterns in a complete human genome

BY ARIEL GERSHMAN, MICHAEL E. G. SAURIA, XAVI GUTIÁR, MITCHELL R. VOLLSER, PAUL W. HOOK, SAVANNAH J. HOYT, MITEN JAIN, ALAINA SHUMATE, ROHAM RAZAGHI, SERGEY KOREN, [...] WINSTON TIMP, +9 authors

SCIENCE • VOL. 376, NO. 6588 • 81 APR 2022

SPECIAL ISSUE RESEARCH ARTICLE

The complete sequence of a human genome

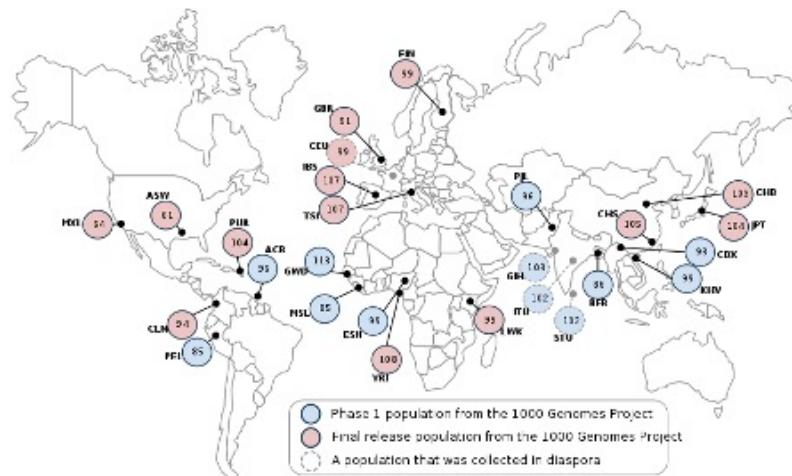
BY SERGEY NURK, SERGEY KOREN, ARANS RHEE, MIKOŁO RAUTIAINEN, ANDREY V. BZKADZE, ALLA MIKHAELOVA, MITCHELL R. VOLLSER, NICOLAS ALTEMOSSE, LEV URALSKY, ARIEL GERSHMAN, [...] ADAM M. PHILLIPPI

SCIENCE • VOL. 376, NO. 6588 • 31 MAR 2022: 44–53

T2T-chrY:Variation across the *entire* genome

1000 Genomes Project (1KGP)

3,202 samples from 26 populations



(Byrska-Bishop et al., Cell, 2022)

Simons Genome Diversity Project (SGDP)

279 open access samples from 130 populations



(Mallick et al., Nature, 2016)

The complete sequence of a human Y chromosome

Rhee et al. (2023) *Nature*

<https://doi.org/10.1038/s41586-023-06457-y>



Stephen Hwang



Dylan Taylor

T2T Team @ Santa Cruz, August 2022

