

The human genome

Michael Schatz

Sept 17, 2023

Lecture 6: Applied Comparative Genomics



Assignment 2: Genome Assembly

Due Wednesday Sept 13 by 11:59pm

The screenshot shows a GitHub repository page for 'assignment2'. The left sidebar lists files: main, assignments (assignment1, assignment2), lectures, LICENSE, and README.md. The README.md file is selected. The main content area displays the assignment details:

Assignment 2: Genome Assembly

Assignment Date: Wednesday, September 6, 2023
Due Date: Wednesday, September 13, 2023 @ 11:59pm

Assignment Overview

In this assignment, you are given a set of unassembled reads from a mysterious pathogen that contains a secret message encoded someplace 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 doublecheck 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 chip.

Question 1. Coverage Analysis [20 pts]

Download the reads and reference genome from:
<https://github.com/schatzlab/appliedgenomics2023/blob/main/assignments/assignment2/asm.tgz?raw=true>

Note we have provided both paired-end and mate-pairs reads (see included README for details). Make sure to look at all of the reads for the coverage analysis and kmer analysis, as well as in the assembly.

- Question 1a. How long is the reference genome? [Hint: Try `samtools faidx`]
- Question 1b. How many reads are provided and how long are they? Make sure to measure each file separately [Hint: Try `FastQC`]
- Question 1c. How much coverage do you expect to have? [Hint: A little arithmetic]
- Question 1d. Plot the average quality value across the length of the reads [Hint: Screenshot from `FastQC`]

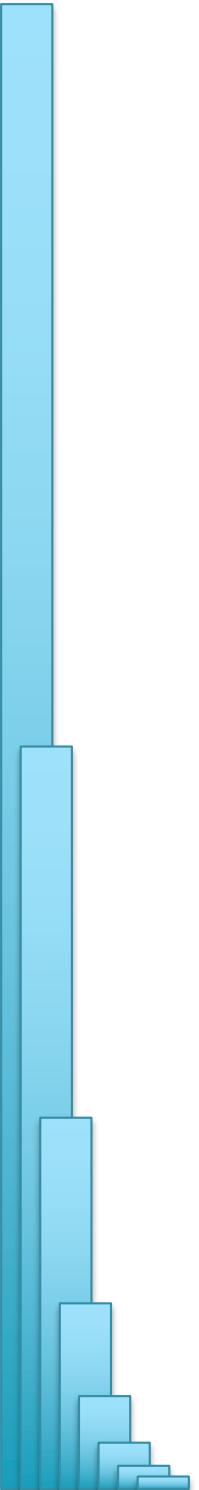
Question 2. Kmer Analysis [20 pts]

Use `Jellyfish` to count the 21-mers in the reads data. Make sure to use the "-C" flag to count canonical kmers, otherwise your analysis will not correctly account for the fact that your reads come from either strand of DNA.

- Question 2a. How many kmers occur exactly 50 times? [Hint: try `jellyfish histo`]

<https://github.com/schatzlab/appliedgenomics2023/tree/main/assignments/assignment2>

Check Piazza for questions!

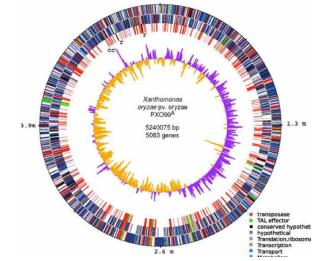


Wednesday's lecture

How many of you know?

- Hash Table
- Suffix Array
- FM Index
- Dynamic Programming
- Edit Distance
- Learned Index Structure

Assembly Summary



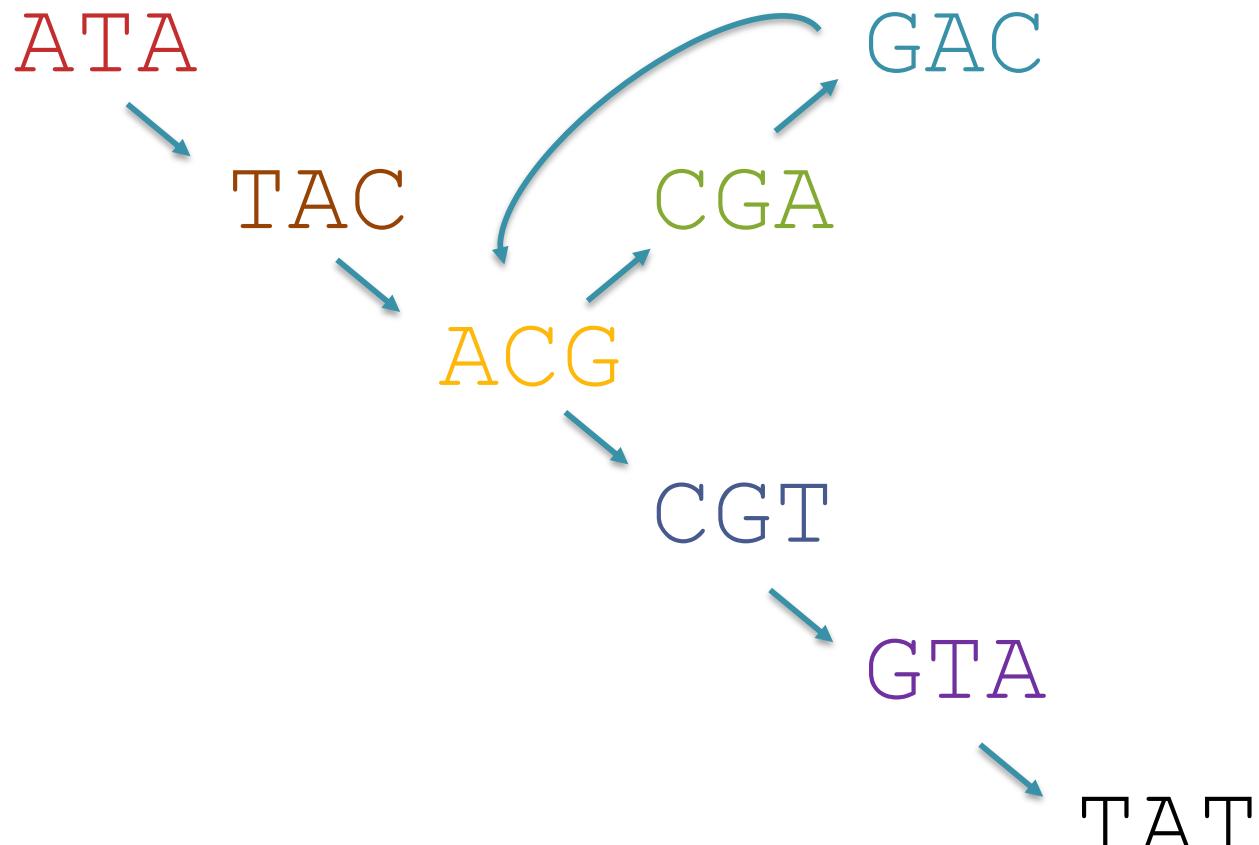
Assembly quality depends on

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

Pop Quiz 2

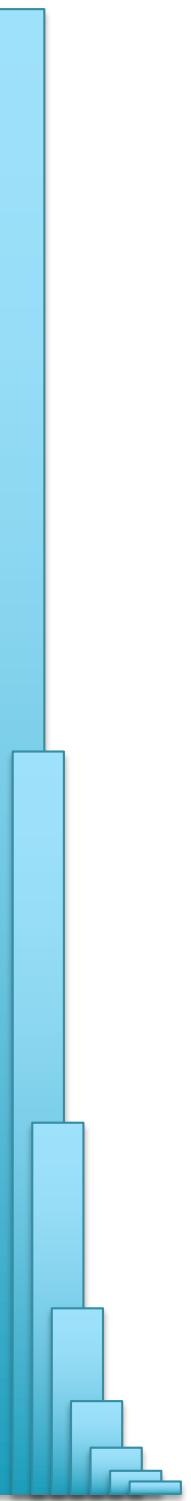
Assemble these reads using a de Bruijn graph approach ($k=3$):

~~-ACGA~~
~~-ACGT~~
~~-ATAC~~
~~-CGAC~~
~~-CGTA~~
~~-GACG~~
~~-GTAT~~
~~-TACG~~



Should we add the edge $TAT \rightarrow ATA$?

ATACGACGTAT



Part 2: 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.
- $40 \text{ trillion} * 4 \text{ meters} = 160 \text{ trillion meters}$.
- $160 \text{ trillion meters} / 1609.34 = 99,750,623,441 \text{ miles}$
- $99,750,623,441 / 92,960,000 = 1,073.05 \text{ trips to the sun.}$

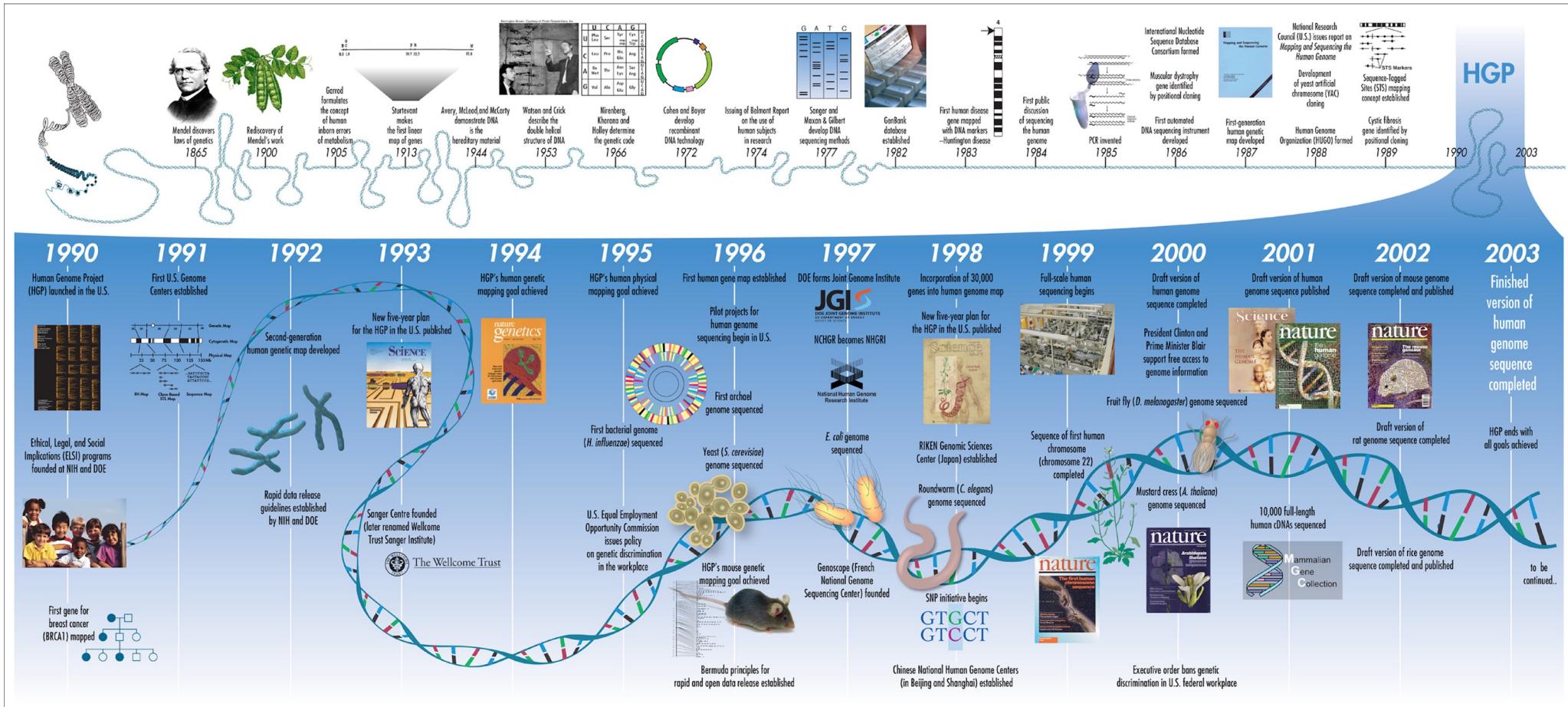
A typical cell replicates about 100 times

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 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?

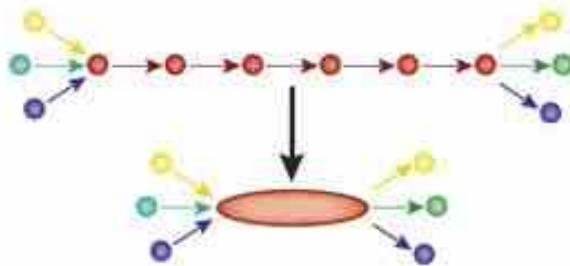
1. Fragment DNA and sequence



2. Find overlaps between reads

...AGCCTAGACCTACA **GGATGCGGCGACAGGT**
 GGATGCGGCGACAGGT CGCATATCCGGT...

3. Assemble overlaps into contigs



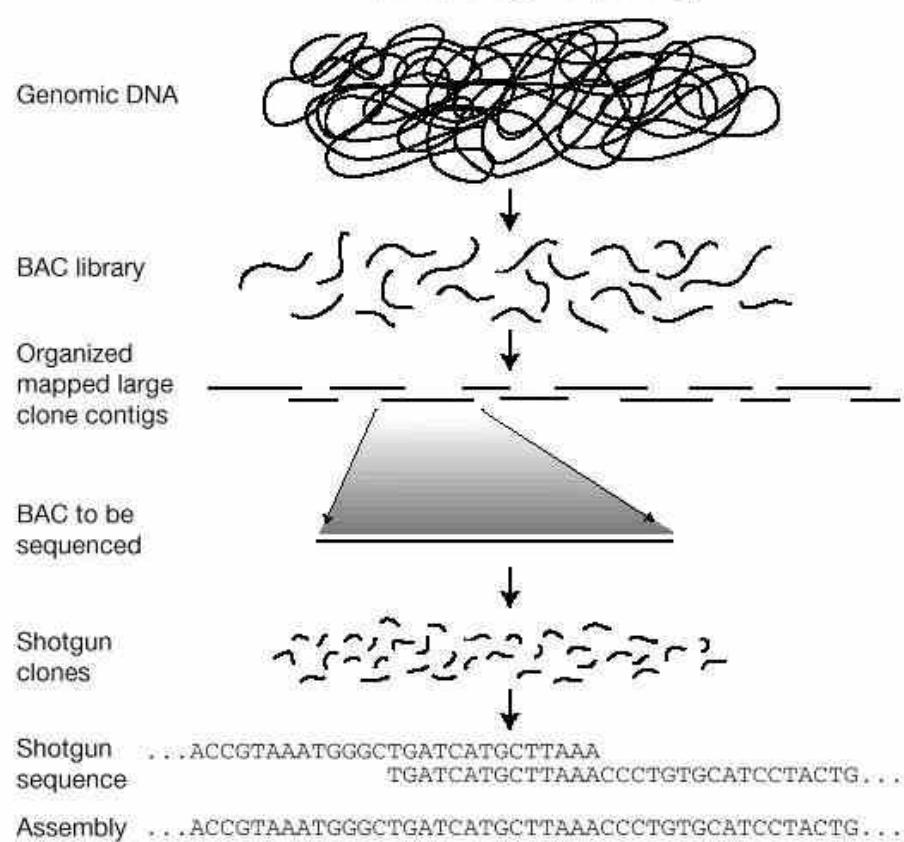
4. Assemble contigs into scaffolds



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."

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

insist on being respectful in our 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 "spirit of the violent militia movement. Non-violence can serve as an anti-government oppression, he added.

"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

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 increase their authority, always at the expense of the people."

ly, government has forgotten that "it's the servant of the people," Parlato added, acting more like it's the master."

to and the Lapps share an abiding non-violent civil disobedience.

insist on being respectful in our 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 "spirit of the violent militia movement. Non-violence can serve as an anti-government oppression, he added.

"Law is unjust or you're given an authority 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.

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 increase their authority, always at expense of the people."

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

insist on being respectful in our of resistance," Barbara Lyn Lapp but if we claim to care about ours, we must protest government in-

violence has to be the watchword, said, calling civil disobedience the of the violent militia movement. Non-violence can serve as an anti-government oppression, he added.

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.

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?

Welcome back: Michael Schatz
 Logout

Search go Advanced search

Journal home > Archive > Editorial > Full Text:

Journal content

- Journal home
- Advance online publication
- Current issue
- Archive**
- Focuses and Supplements
- Methagora blog
- Method of the Year 2016
- Multimedia
- Press releases

Journal Information

- Guide to authors
- Reporting checklist
- Online submission
- Subscribe
 - New Subscription
 - Renew Subscription
 - Paid Subscriptions
 - Change of Address
- Permissions
- For referees
- Contact the journal
- About this site

Nature Research services

- Authors & Referees
- Advertising

EDITORIAL

Nature Methods 7, 331 (2010)
doi:10.1038/nmeth0510-331

E pluribus unum

If the human reference genome is to reflect more of the actual genomic diversity in humans, community participation is needed.

Please visit [methagora](#) to view and post comments on this article.

The human genome is ten years old. We acknowledge its reference assembly as an invaluable resource essential for many purposes such as the assembly of short reads from high-throughput sequencing platforms into chromosome context during resequencing projects. At the same time, we think necessary improvement of the reference genome depends on the willingness of the research community to provide data for the genome's less accessible regions.

First published in 2001, the human reference genome has, since 2007, been in the hands of the Genome Reference Consortium (GRC) a small group of fewer than 20 scientists from the European Bioinformatics Institute, the US National Center for Biotechnology Information, The Sanger Institute and The Genome Center at Washington University in St. Louis, who have committed to the improvement and completion of this reference, with very little financial support.

The reference genome is now in its 19th rendition, and probably the best measure of its improvement over the last ten years is the number of fragments it consists of. The very first version had ~150,000 gaps; the most recent build, GRCh37, has only around 250 gaps.

The only other publicly accessible *de novo* assembly of a human genome that contains chromosome sequences is HuRef. Obtained by traditional capillary sequencing, HuRef is the diploid genome of Craig Venter. It comes in 4,500 pieces and, like any individual genome, it contains many rare alleles.

GRCh37, in contrast, is a mosaic haploid genome derived from about 13 people. It still contains rare alleles, but the GRC recently decided to convert these to common haplotypes. Deciding which alleles are common and which are rare is proving challenging, and the GRC members are collaborating with members of the 1000 Genomes project to collect enough data to make these decisions.

Subscribe to Nature Methods

This issue

- Table of contents
- Next article

Article tools

- Download PDF
- Send to a friend
- CrossRef lists 11 articles citing this article
- Scopus lists 9 articles citing this article
- Export citation
- Rights and permissions

naturejobs

Recruitment of Professors and Associate Professors
School of Materials Science and Engineering, Sun Yat-sen University
Sun Yat-sen University

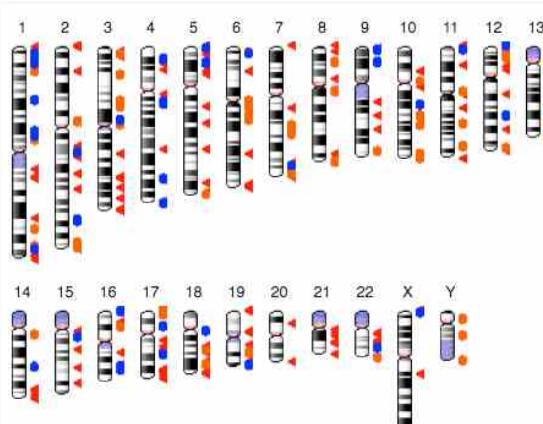
Faculty positions at Institut franco-chinois de l'énergie nucléaire
Institut franco-chinois de l'énergie nucléaire Sun Yat-sen University

More science jobs

Post a job

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 assembly by both generating multiple representations (alternative paths) for each chromosome, represented by a single path. Additionally, we are reassembling the genome to allow users who are interested in a specific locus to do so without affecting users who need chromosome coordinate sets.

Download data:

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

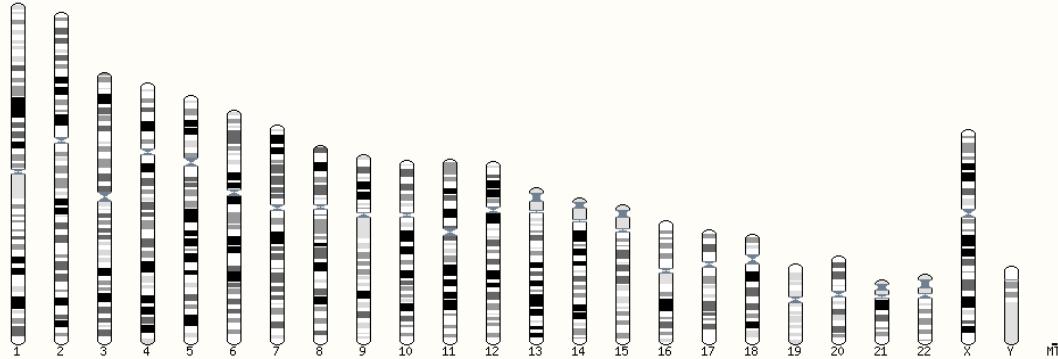
Transitioning to GRCh38? Try the [NCBI Remapper](#) to convert your assembly alignments used by the GRC.

Next assembly update

The next assembly update (GRCh38.p12) will be released in July 2017.



The human genome - basic stats



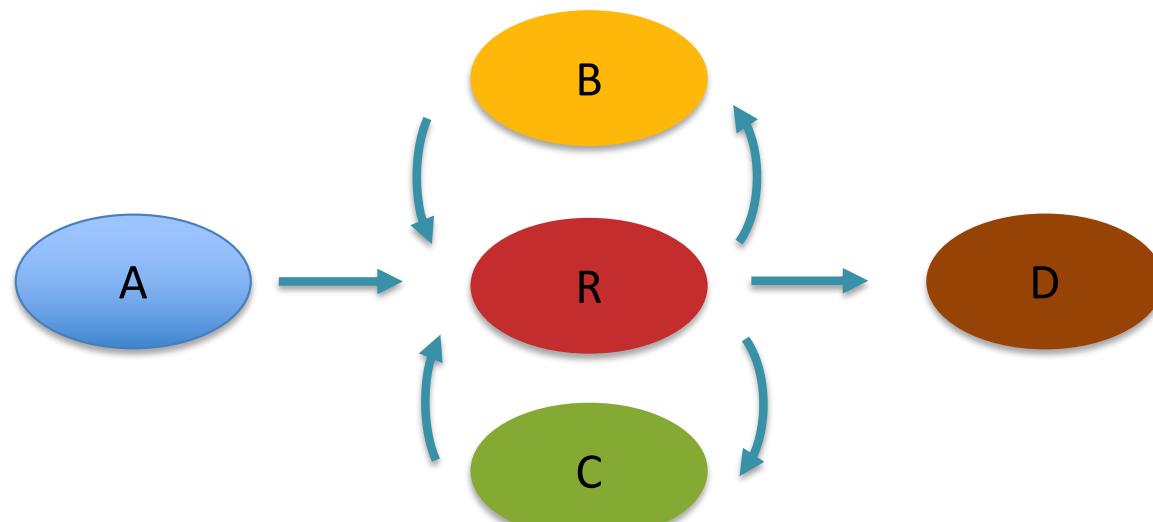
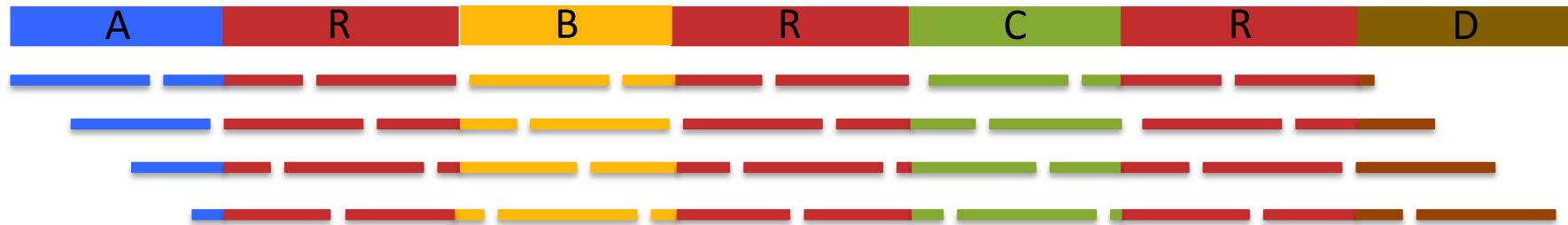
- 3.096 billion base pairs (haploid)
- 20,454 protein coding genes
- 226,950 coding transcripts
(isoforms of a gene that each encode a distinct protein product)

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

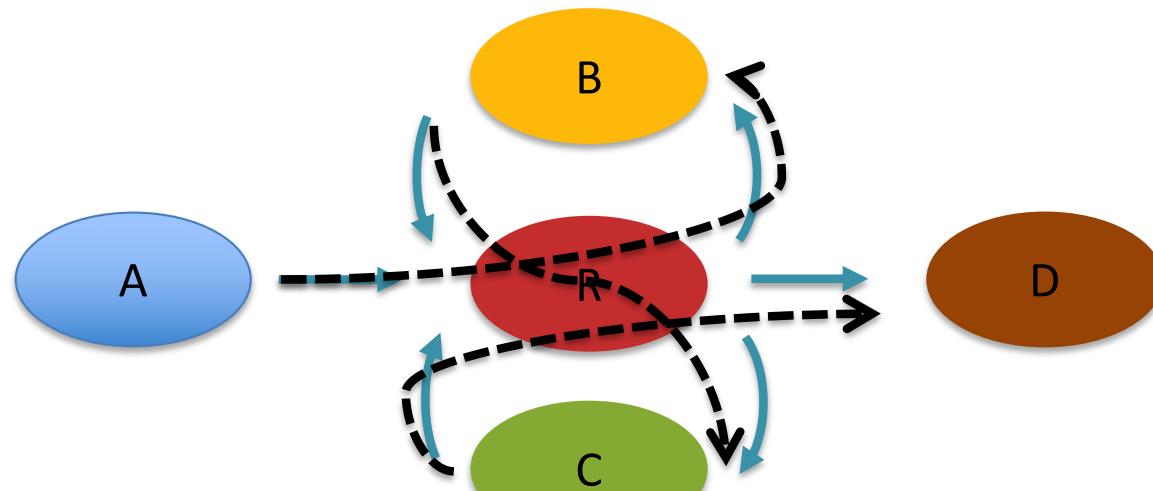
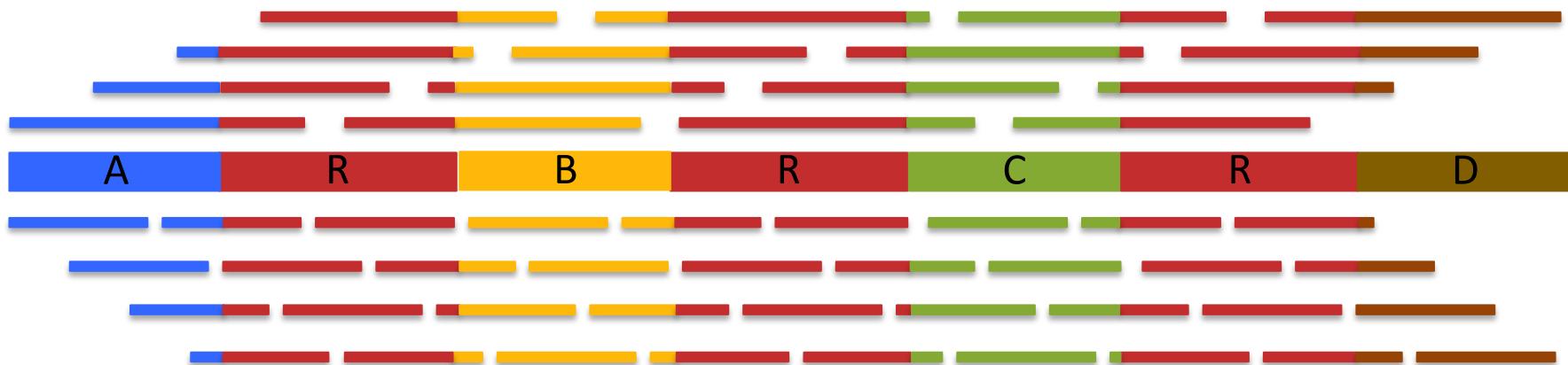
Gene counts (Primary assembly)

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

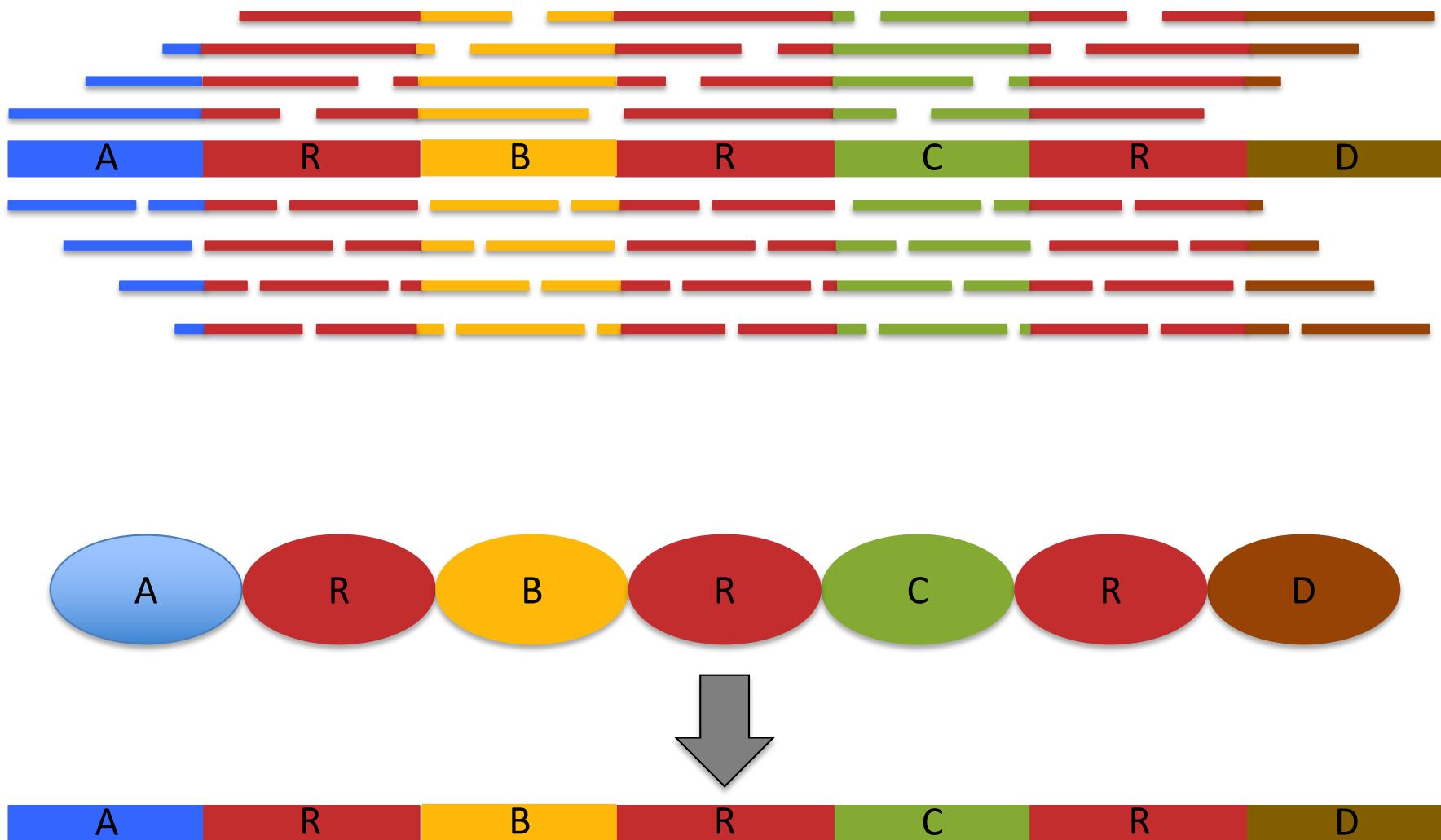
Assembly Complexity



Assembly Complexity



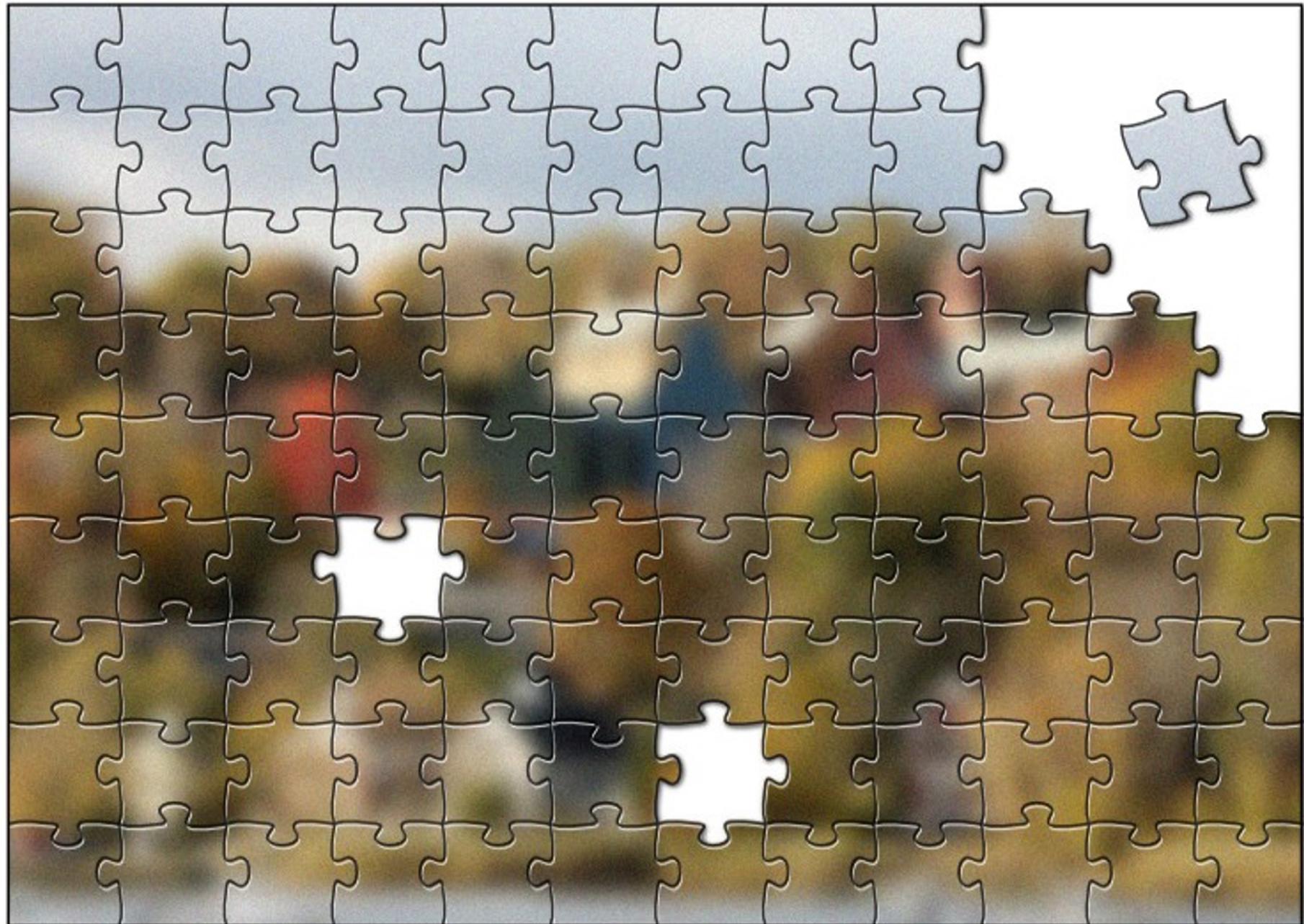
Assembly Complexity



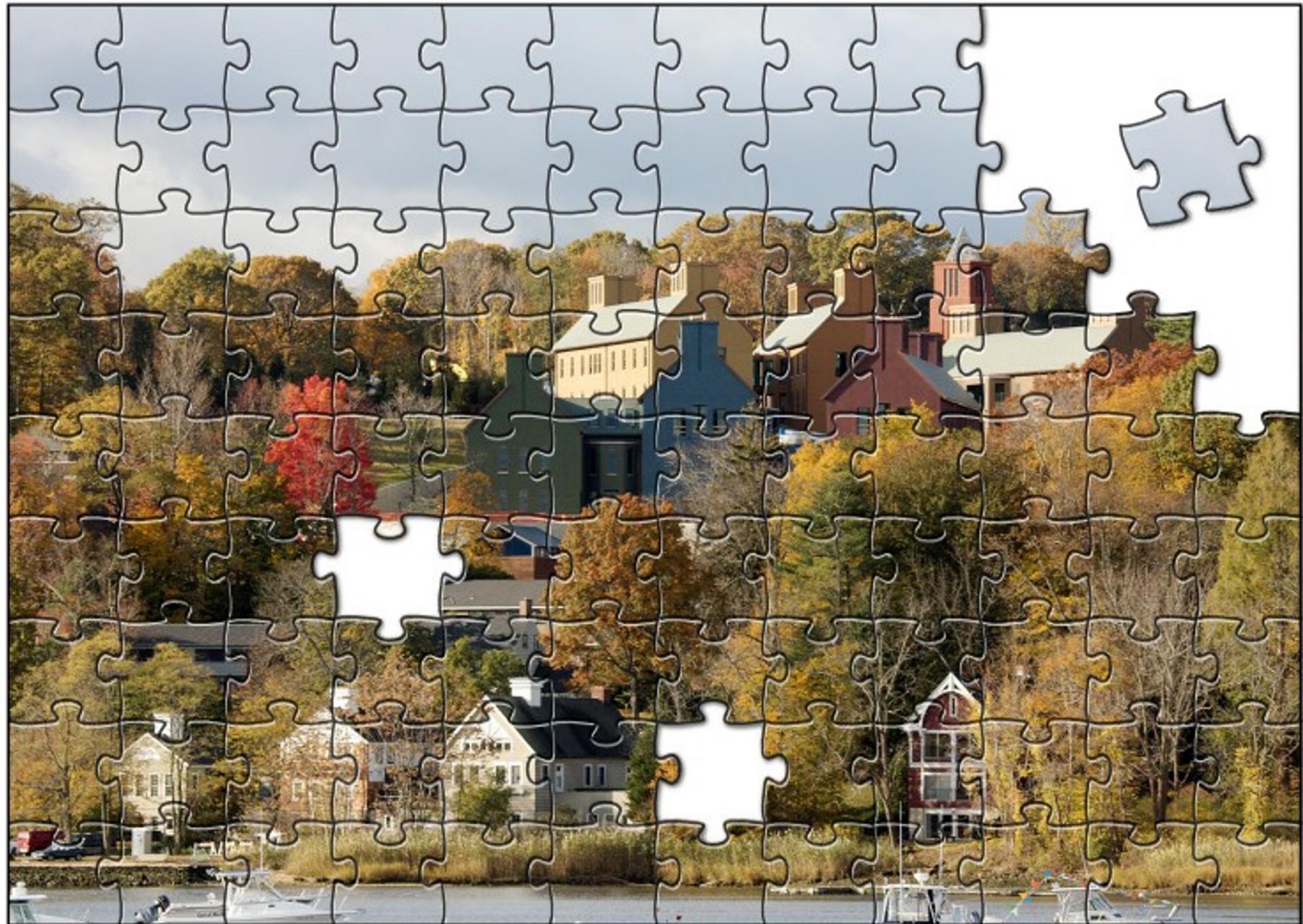
The advantages of SMRT sequencing

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

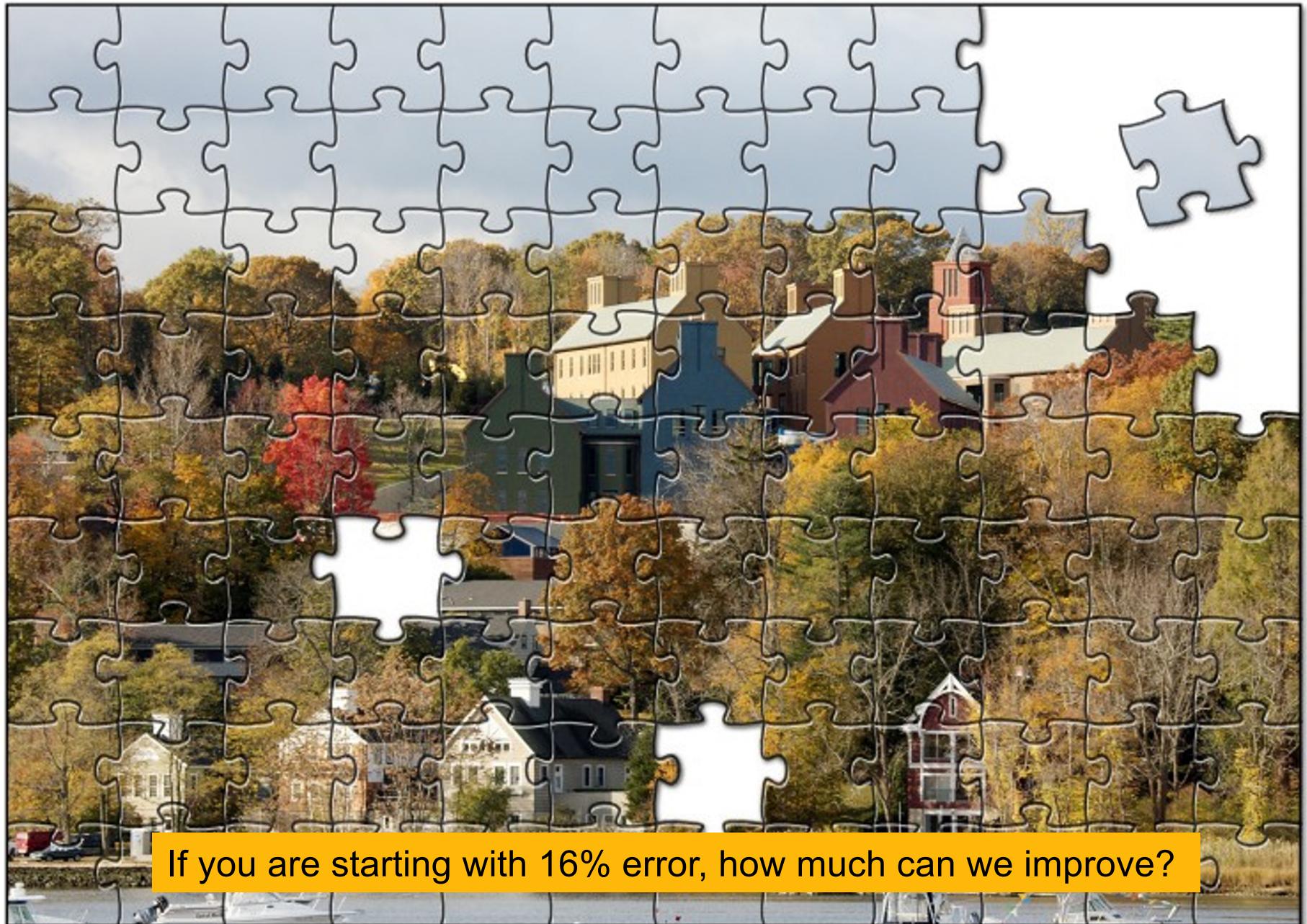
Single Molecule Sequences



“Corrective Lens” for Sequencing

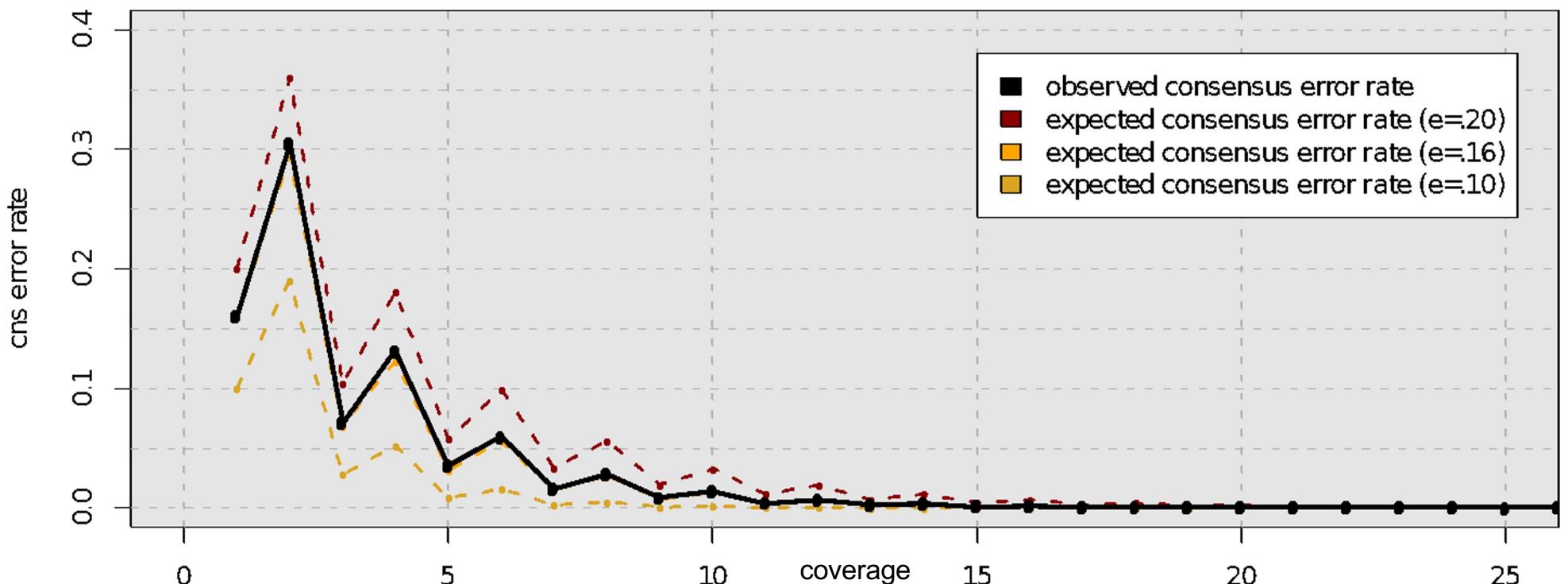


“Corrective Lens” for Sequencing



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

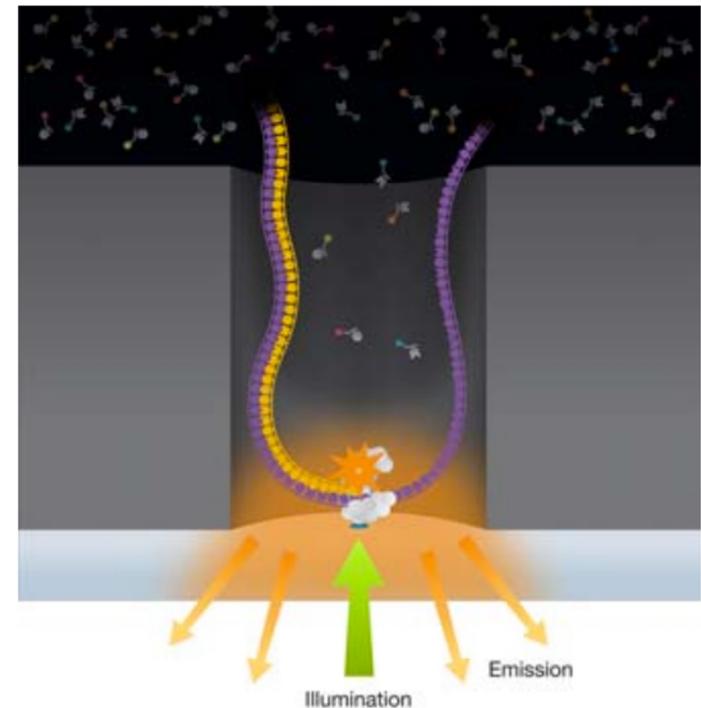
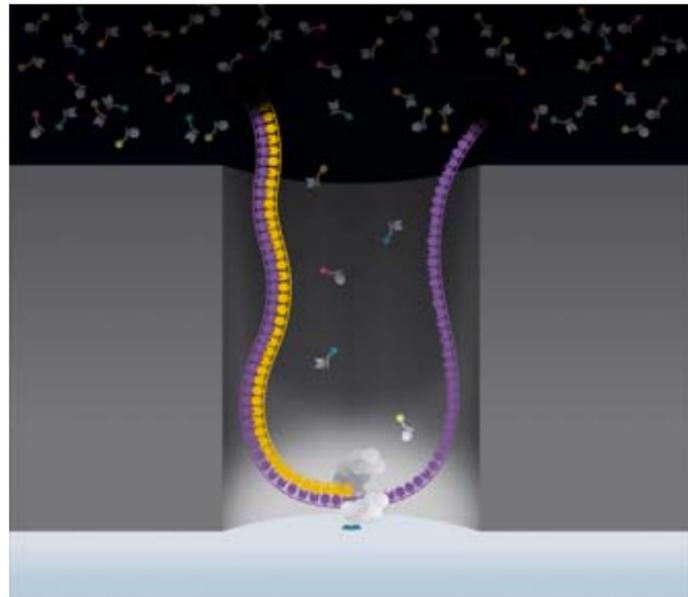
$$CNS\ Error = \sum_{i=\lceil c/2 \rceil}^c \binom{c}{i} (e)^i (1-e)^{n-i}$$



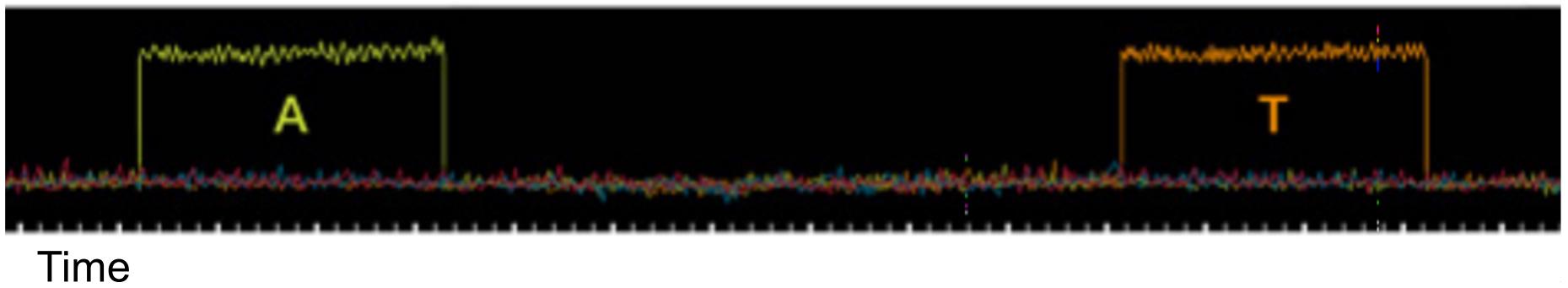
PacBio Single Molecule Real Time Sequencing (SMRT-sequencing)

PacBio: SMRT Sequencing

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



Intensity



[https://www.youtube.com/watch?v= ID8JyAbwEo](https://www.youtube.com/watch?v=ID8JyAbwEo)

“HiFi” Circular Consensus Reads

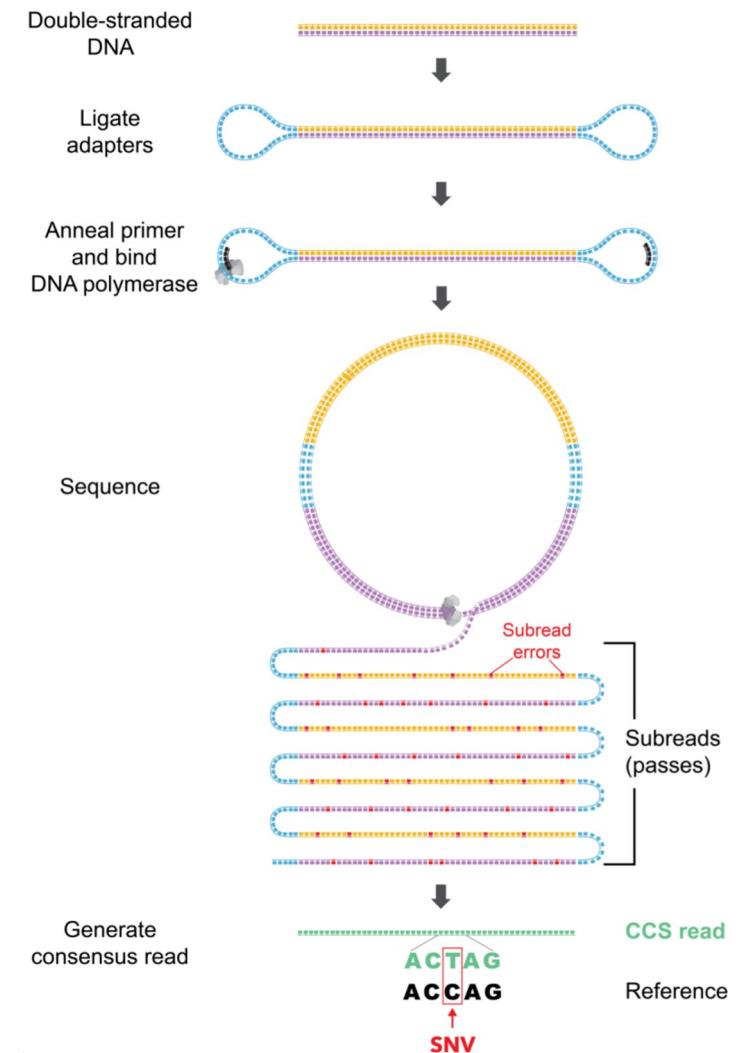
High-quality reads produced by sequencing the same molecule multiple times

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

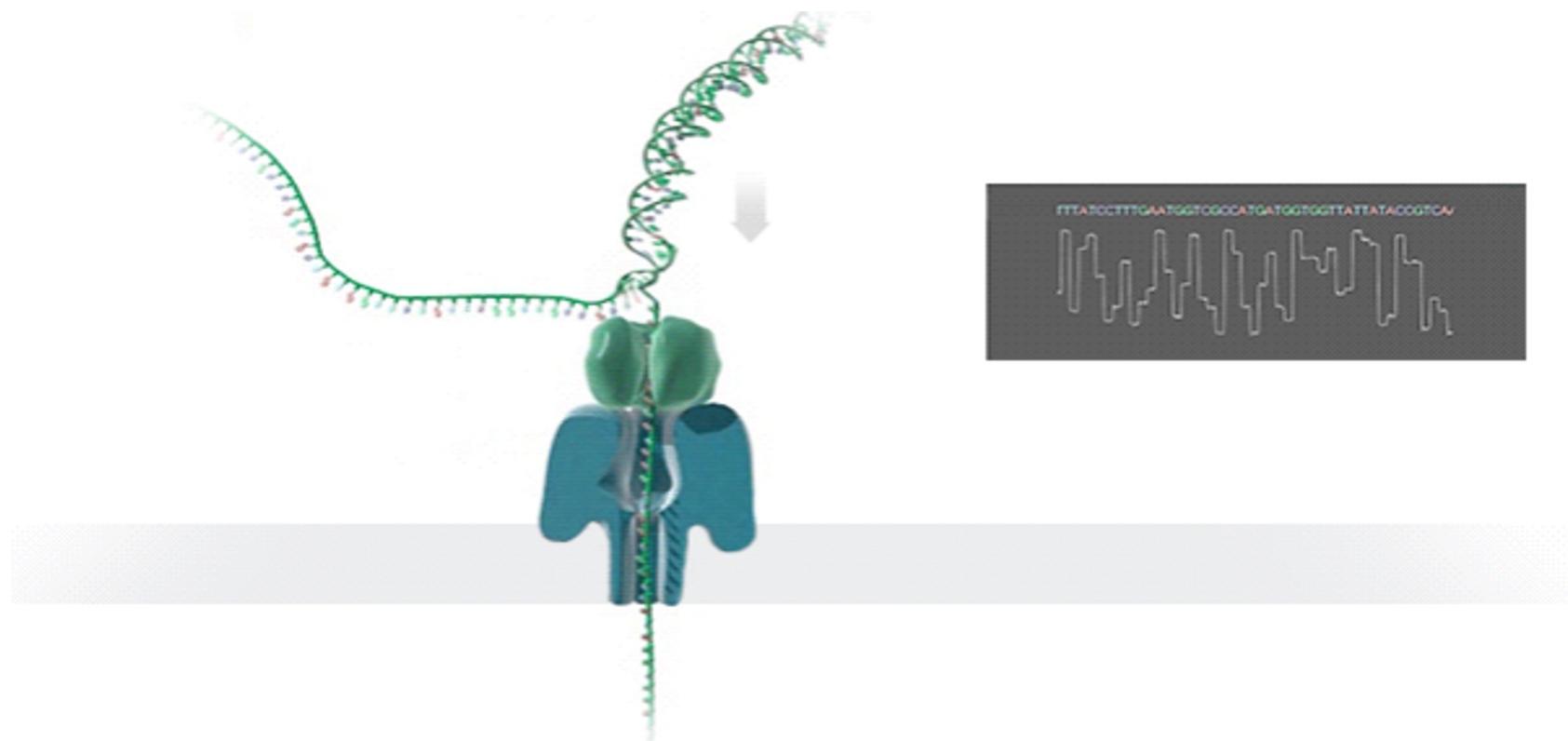


Oxford Nanopore Technologies (ONT)



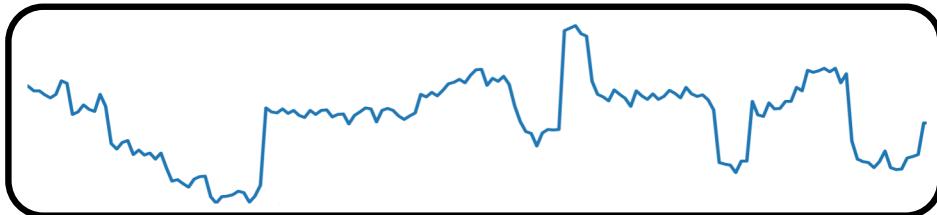
Nanopore Sequencing

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



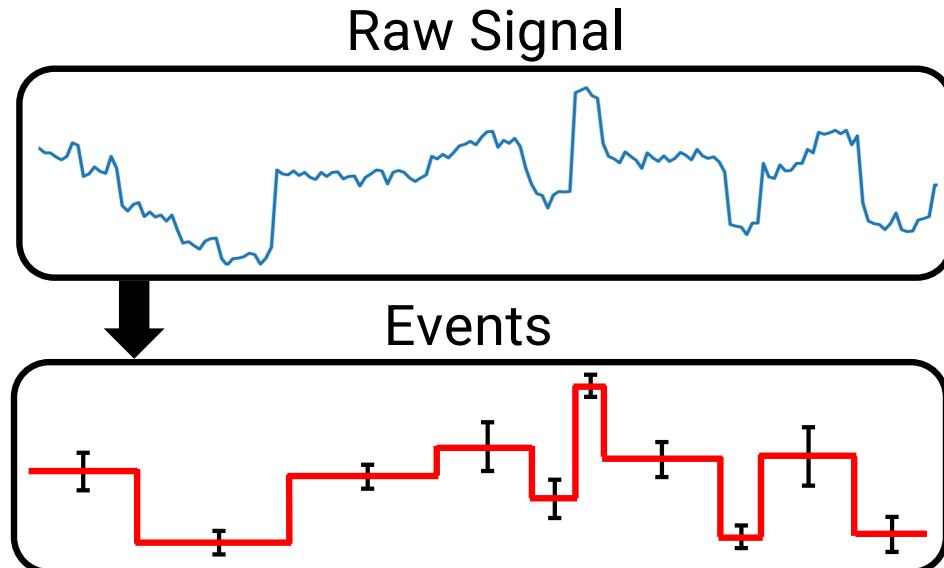
Nanopore Basecalling

Raw Signal



Translation of raw signal
into basepairs

Nanopore Basecalling

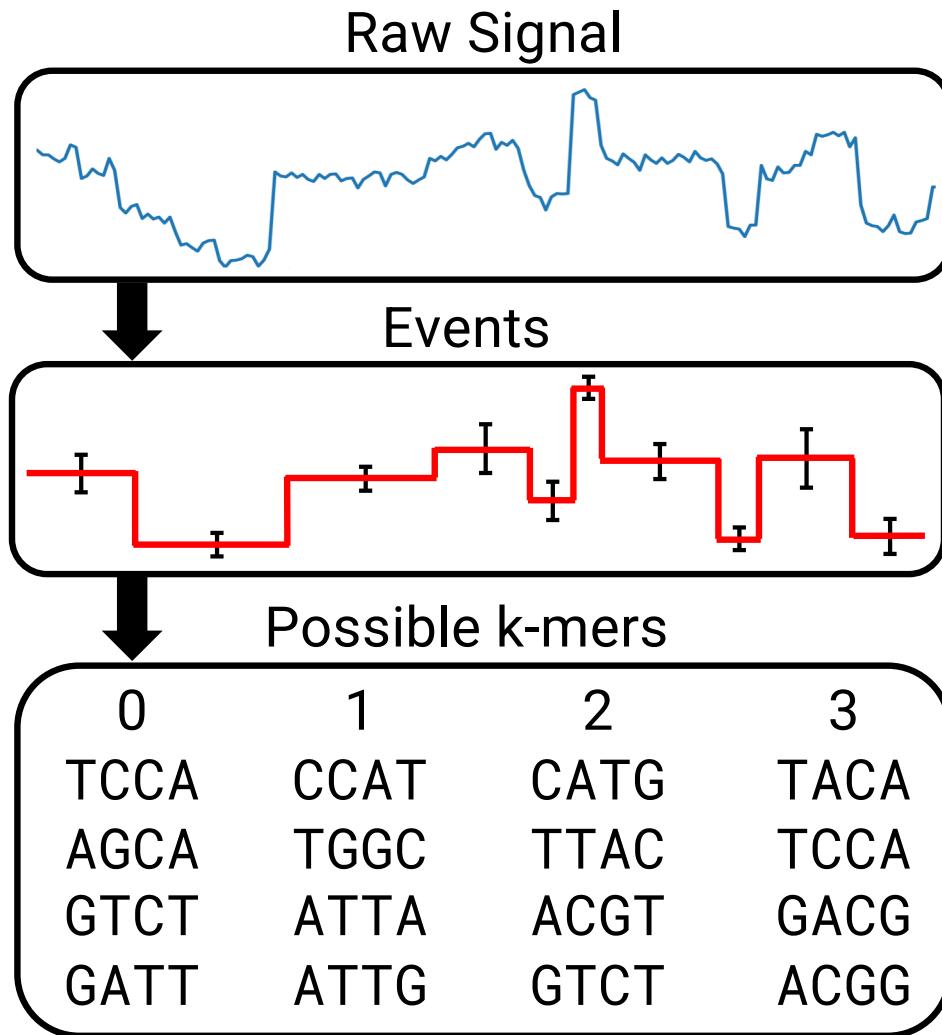


Translation of raw signal
into basepairs

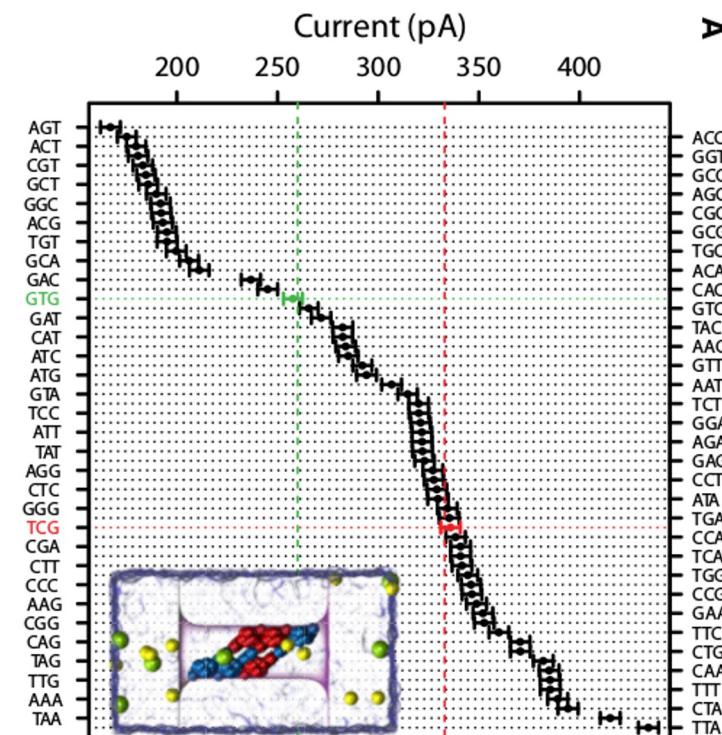
Early basecallers began by
estimating k-mer boundaries
using “events”, which were
then input to an HMM

Modern basecallers use
neural networks directly
on raw signal

Nanopore Basecalling

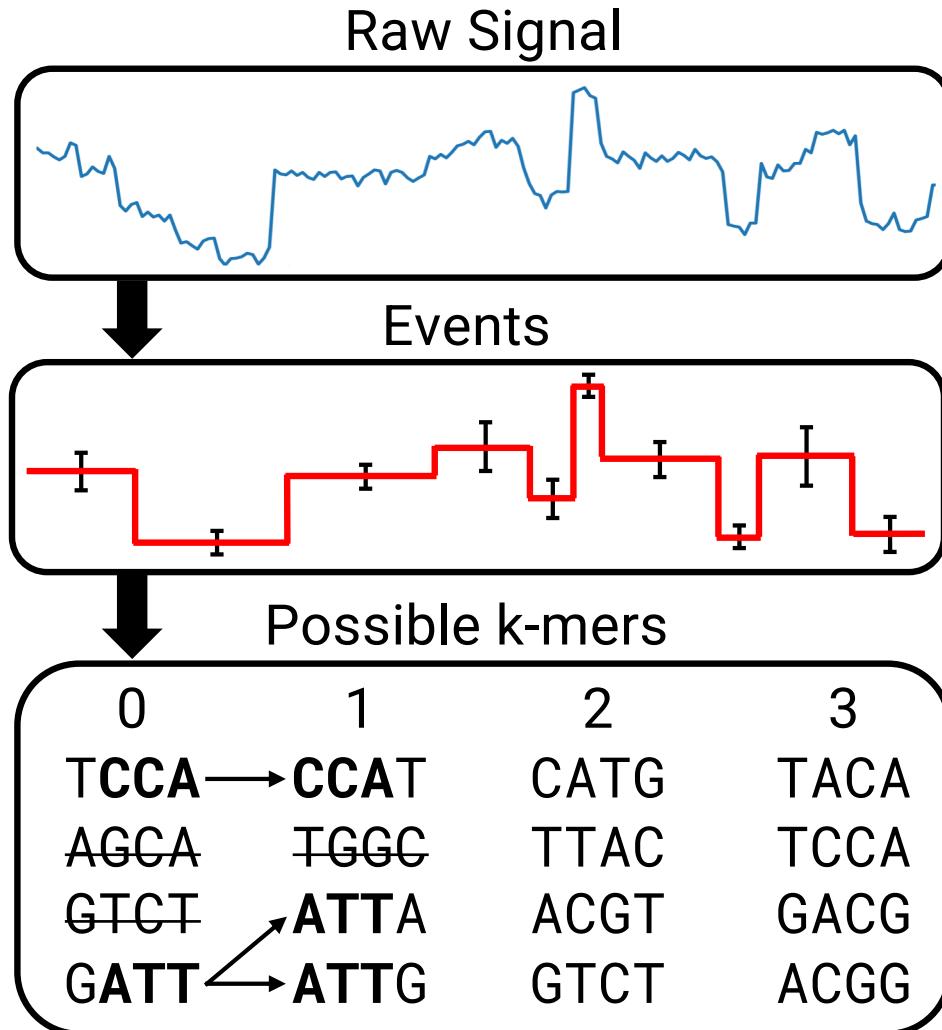


ONT releases k-mer models with expected current distribution of every k-mer

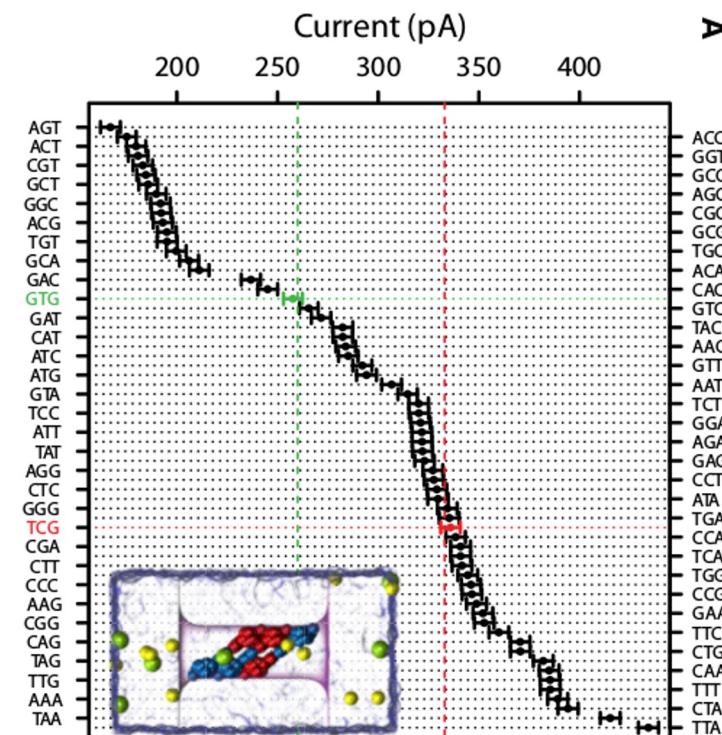


DNA Base-Calling from a Nanopore Using a Viterbi Algorithm
Timp et al. (2012) *Biophysical Journal*

Nanopore Basecalling

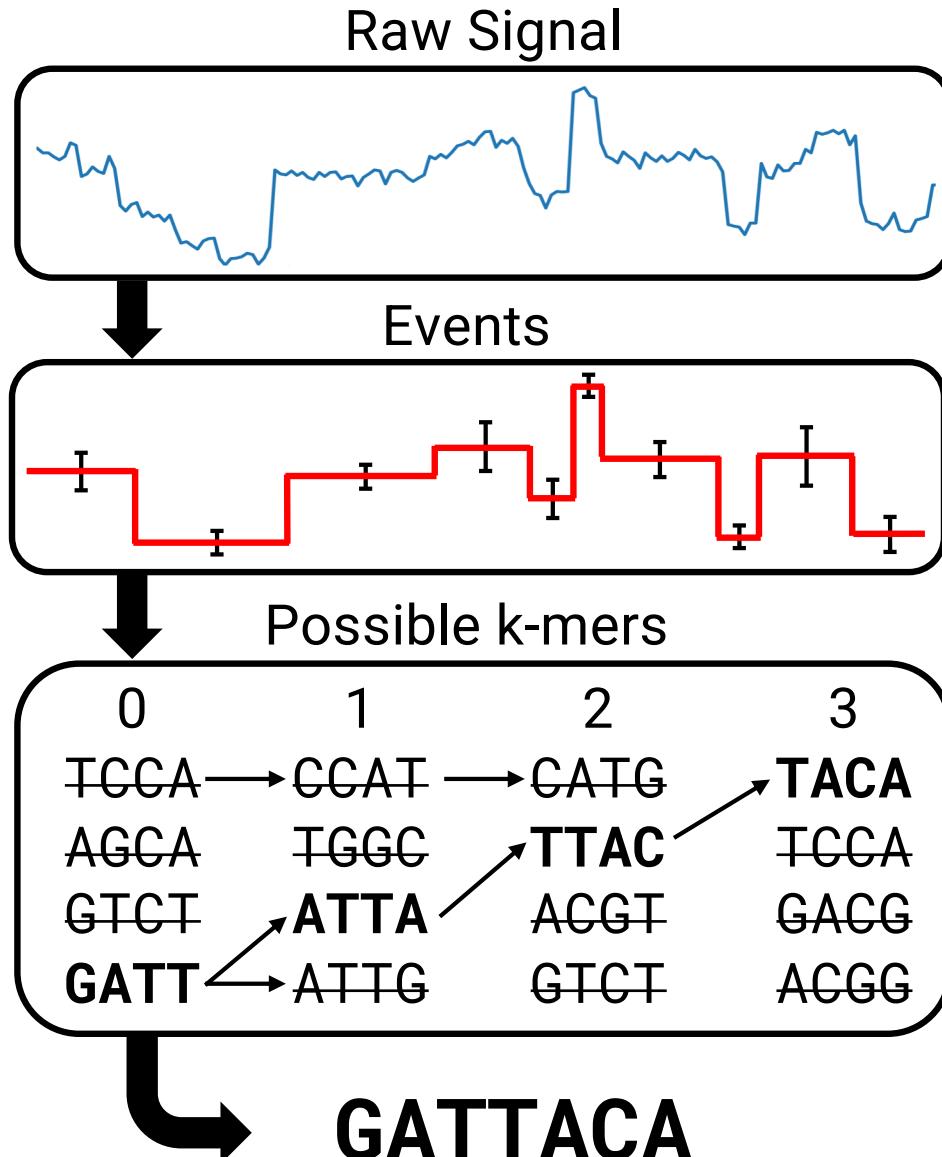


Certain k-mers can be eliminated based on possible transitions

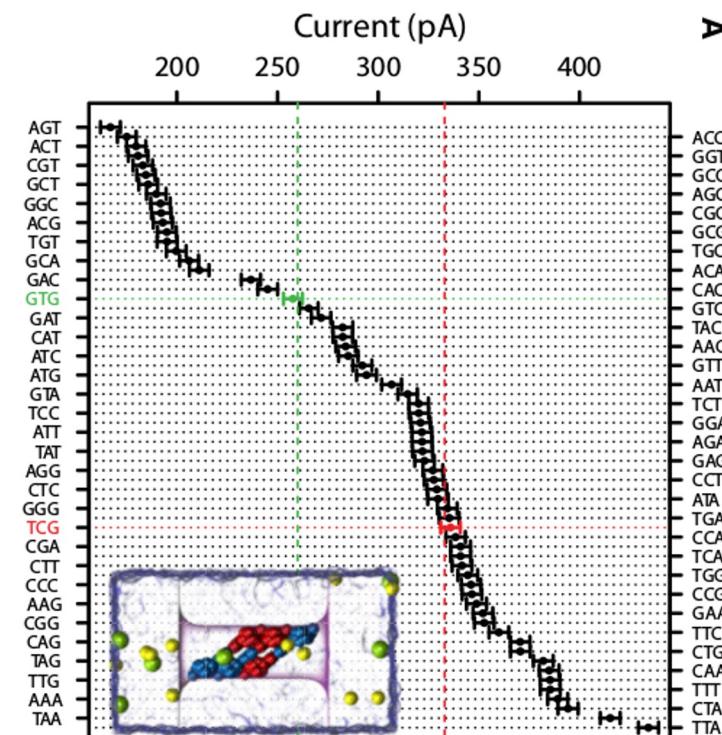


DNA Base-Calling from a Nanopore Using a Viterbi Algorithm
Timp et al. (2012) *Biophysical Journal*

Nanopore Basecalling



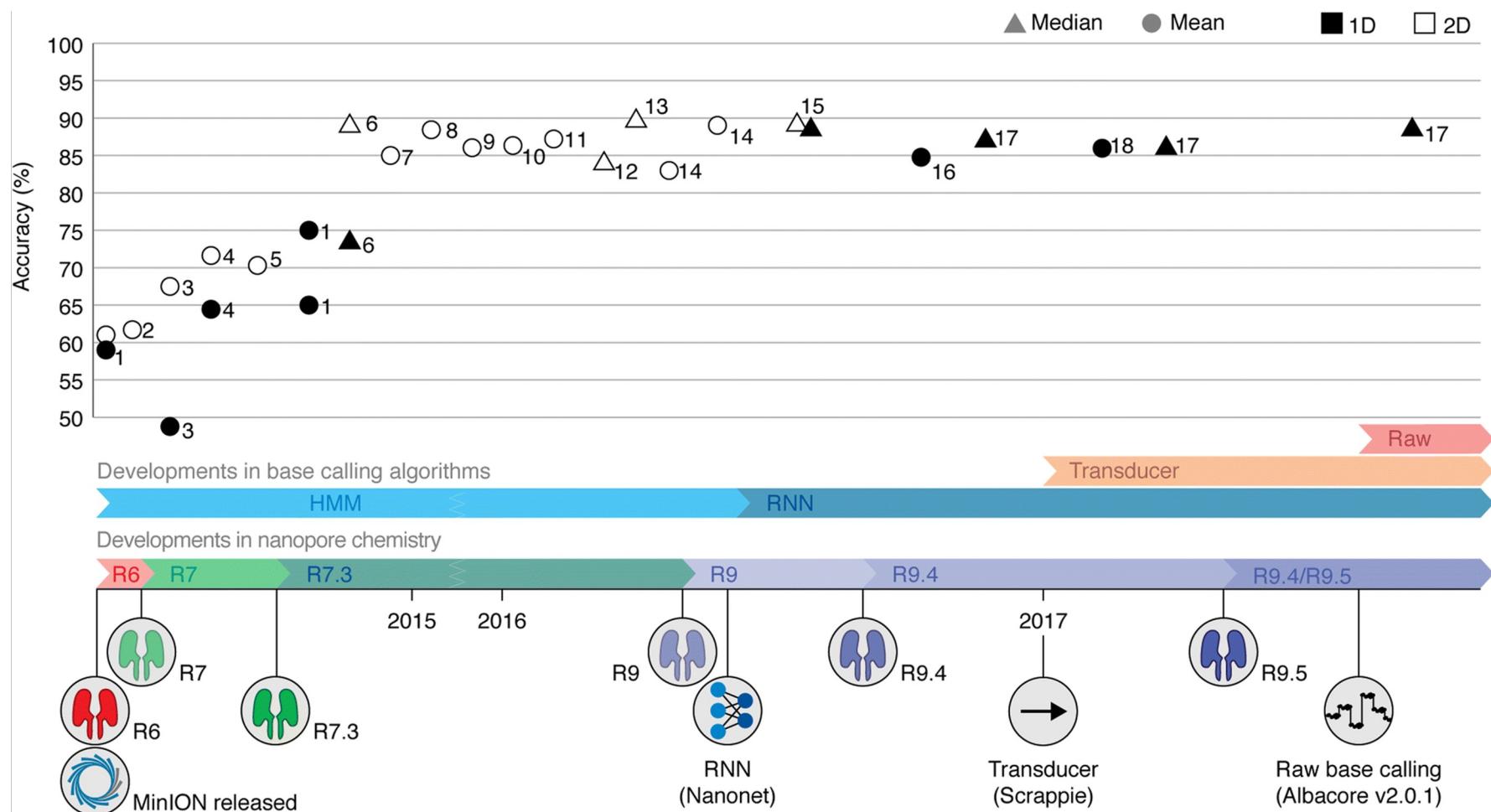
Final sequence determined by most probable k-mers



"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

