

The human genome

Michael Schatz

Sept 14, 2022

Lecture 5: Applied Comparative Genomics



Assignment 2: Genome Assembly

Due Monday Sept 19 by 11:59pm

The screenshot shows a GitHub repository page for 'appliedgenomics2022'. The repository is public and has 2 forks and 11 stars. The README.md file is displayed, containing the assignment details. The assignment is titled 'Assignment 2: Genome Assembly' and specifies the assignment date as Monday, September 12, 2022, and the due date as Monday, September 19, 2022, at 11:59pm. It includes an 'Assignment Overview' section describing the task of assembling a genome from unassembled reads to find a secret message. It also recommends using bioconda or Docker for tool installation.

Assignment 2: Genome Assembly

Assignment Date: Monday, September 12, 2022
Due Date: Monday, September 19, 2022 @ 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](#).

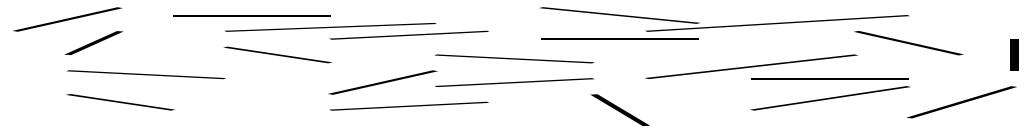
For this assignment, we recommend you install and run the tools using [bioconda](#). There are some tips below in the Resources section.
Alternatively, you can try running the tools using [Docker](#). Docker is a powerful containerization tool to make software easier to distribute. This will

<https://github.com/schatzlab/appliedgenomics2022/tree/main/assignments/assignment2>
Check Piazza for questions!

Part I: Recap

Assembling a Genome

I. Shear & Sequence DNA

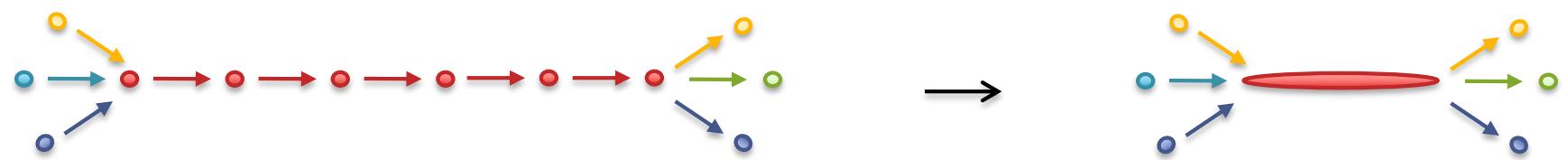


2. Construct assembly graph from reads (de Bruijn / overlap graph)

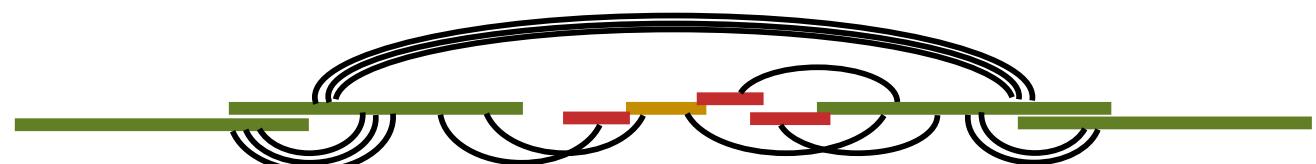
...AGCCTAG**GGATGCGCGACACGT**

GGATGCGCGACACGTCGCATATCCGGTTTGGT**CAACCTCGGACGGAC**
CAACCTCGGACGGACCTCAGCGAA...

3. Simplify assembly graph



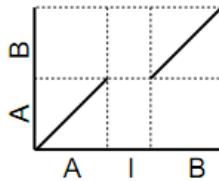
4. Detangle graph with long reads, mates, and other links



SV Types

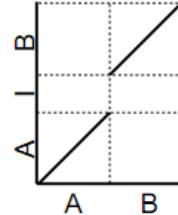
Insertion into Reference

R: AIB
Q: AB



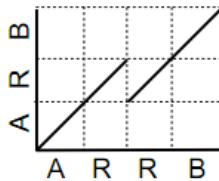
Insertion into Query

R: AB
Q: AIB



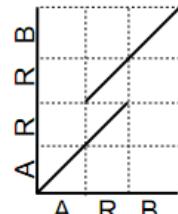
Collapse Query

R: ARRB
Q: ARB



Collapse Reference

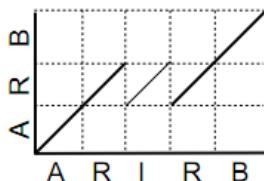
R: ARB
Q: ARRB



Collapse Query w/ Insertion

R: ARIRB
Q: ARB

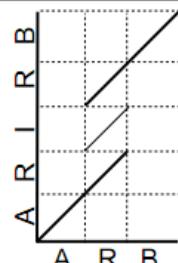
Exact tandem alignment if I=R



Collapse Reference w/ Insertion

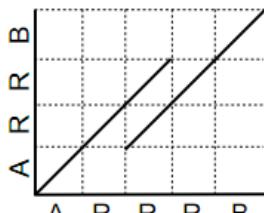
R: ARB
Q: ARIRB

Exact tandem alignment if I=R



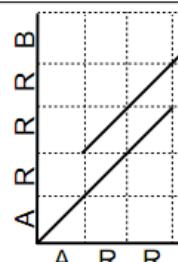
Collapse Query

R: ARRRB
Q: ARRB



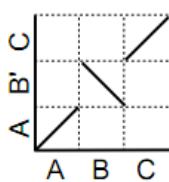
Collapse Reference

R: ARRB
Q: ARRRB



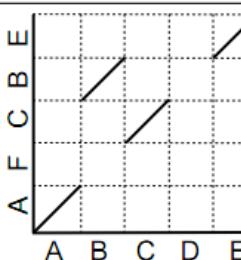
Inversion

R: ABC
Q: AB'C



Rearrangement w/ Disagreement

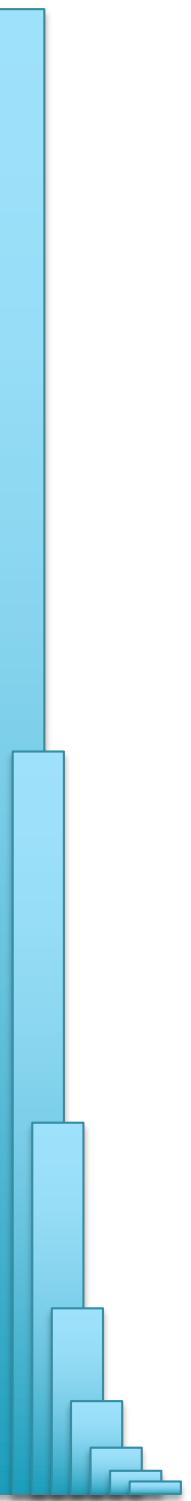
R: ABCDE
Q: AFCBE



- 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



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

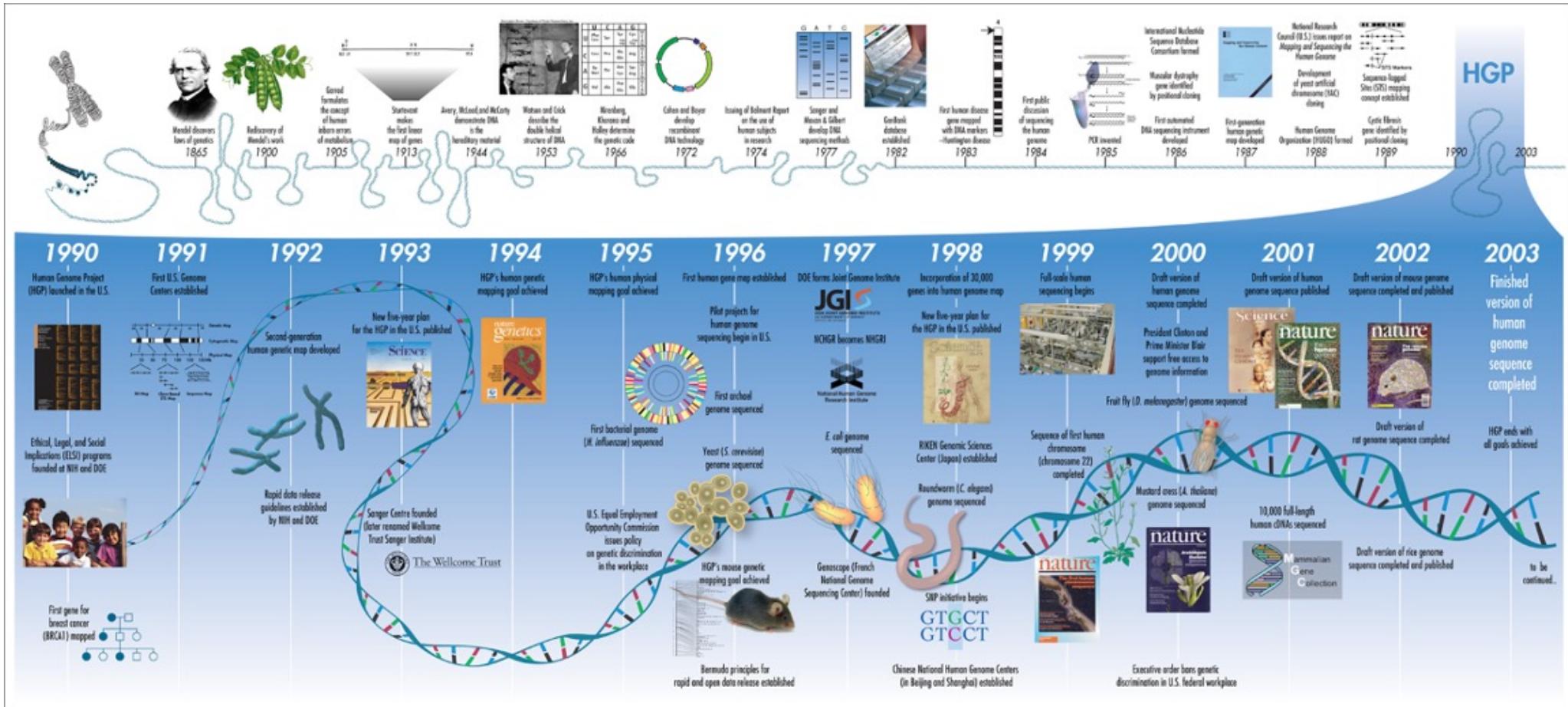
A typical cell replicates about 100 times

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 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?

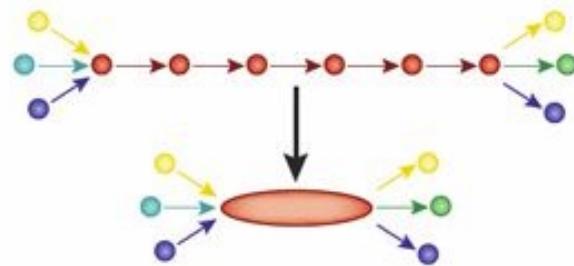
1. Fragment DNA and sequence



2. Find overlaps between reads

...AGCCTAGACCTACA**GGATGCGCGACACGT**
 GGATGCGCGACACGTCGCATATCCGGT...

3. Assemble overlaps into contigs



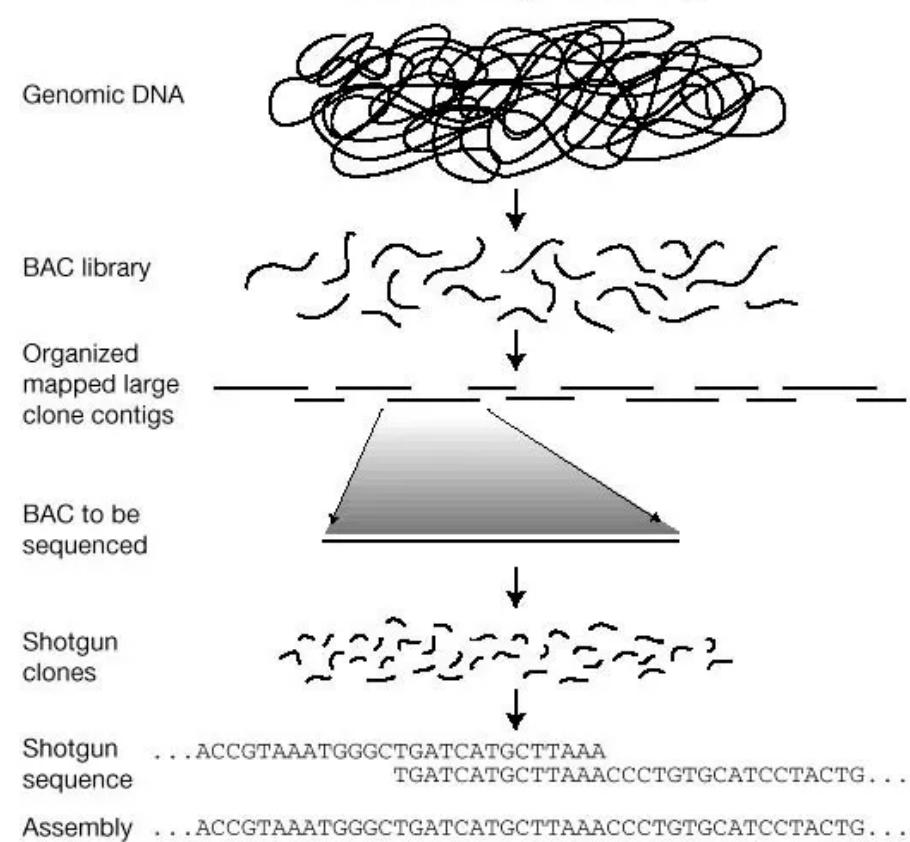
4. Assemble contigs into scaffolds



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 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.

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,

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

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."

"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 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 heart 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

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 resistance," Barbara Lyn Lapp said. "If we claim to care about our rights, we must protest government instead of the watchword, 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](#) [Cart](#)

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

[Subscribe](#)

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

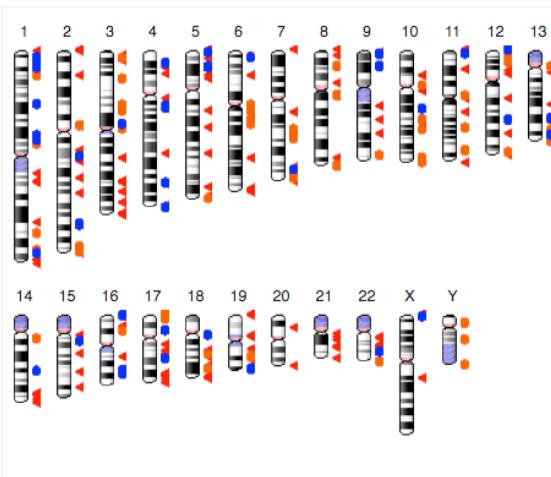
[Faculty positions at Institut franco-chinois de l'énergie nucléaire](#)
Institut franco-chinois de l'énergie nucléaire Sun Yat-sen University

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

The GRC is working hard to provide the best possible assembly by both generating multiple representations (alternatives) for each chromosome, each represented by a single path. Additionally, we are releasing multiple reference assemblies, which allows users who are interested in a specific locus to choose the assembly that is most useful for them. This also allows users who need chromosome coordinate systems to choose the one that is most appropriate for their needs.

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 Remap tool to find the new coordinates for your assembly alignments used by the GRC.

Next assembly update

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



GRCh38.p11

Release date: June 14, 2017

Release type: minor

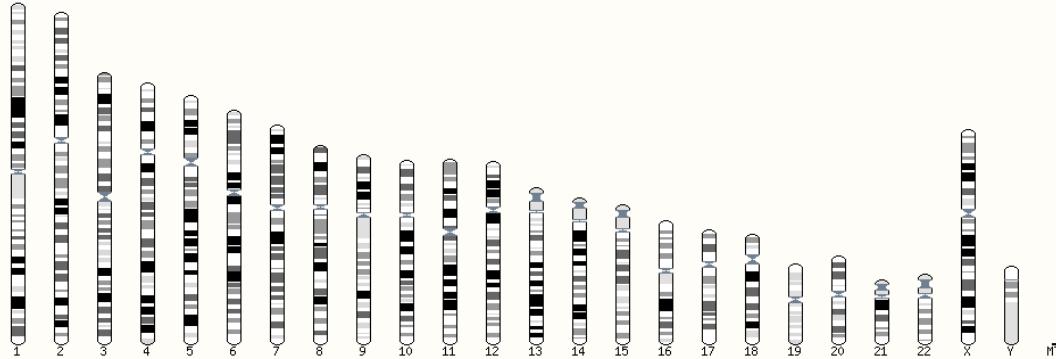
Release notes: GRCh38.p11 is the eleventh patch release for the GRCh38 reference assembly. No chromosome coordinates have changed. The total number of patch scaffolds is now: 64 FIX and 59 NOVEL.

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 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

Genomics Arsenal in the Year 2022

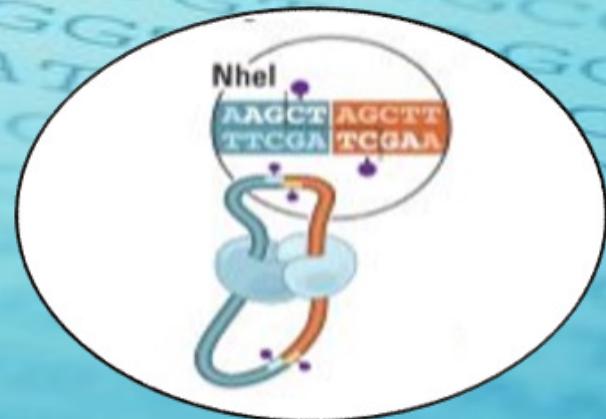
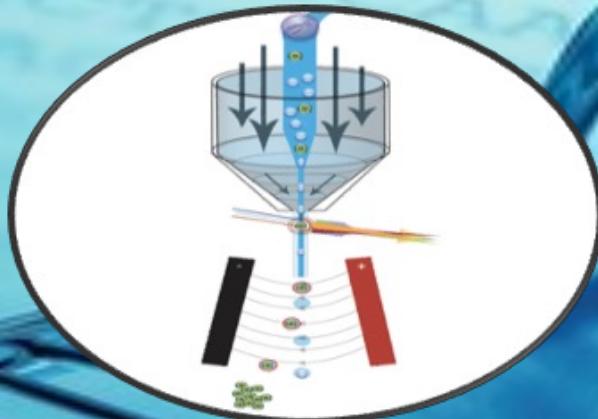
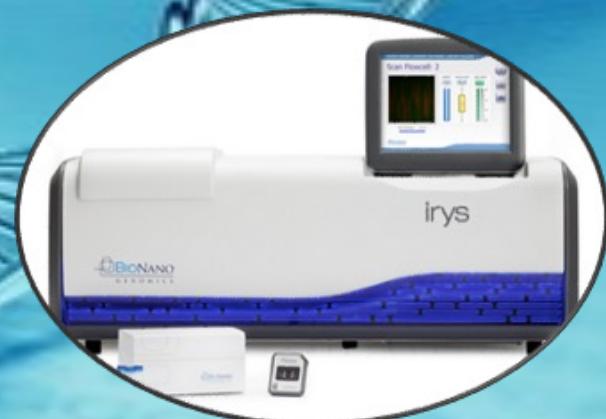
Sample Preparation



Sequencing



Chromosome Mapping

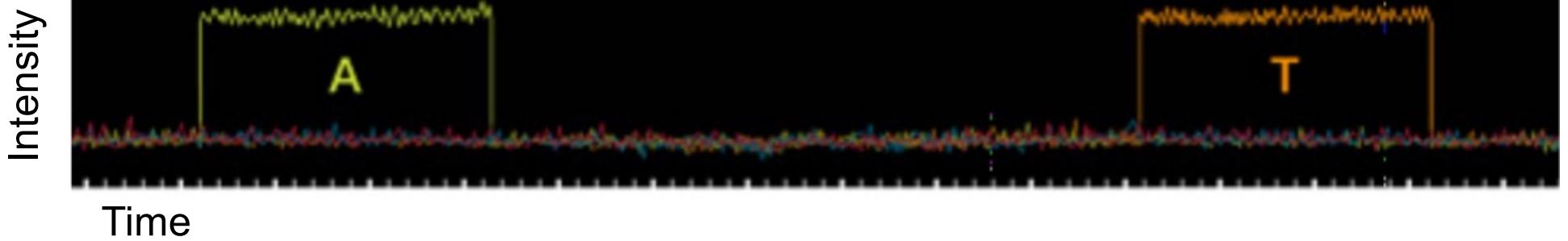
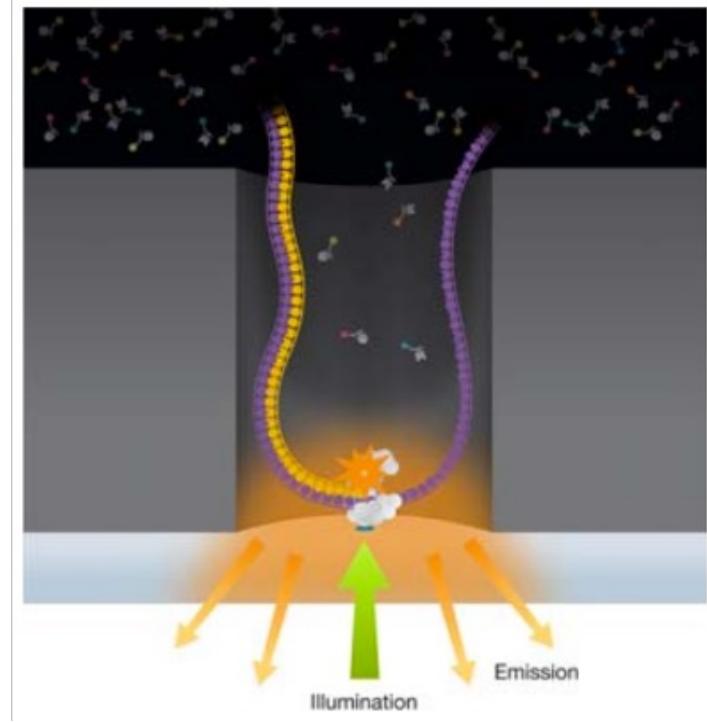
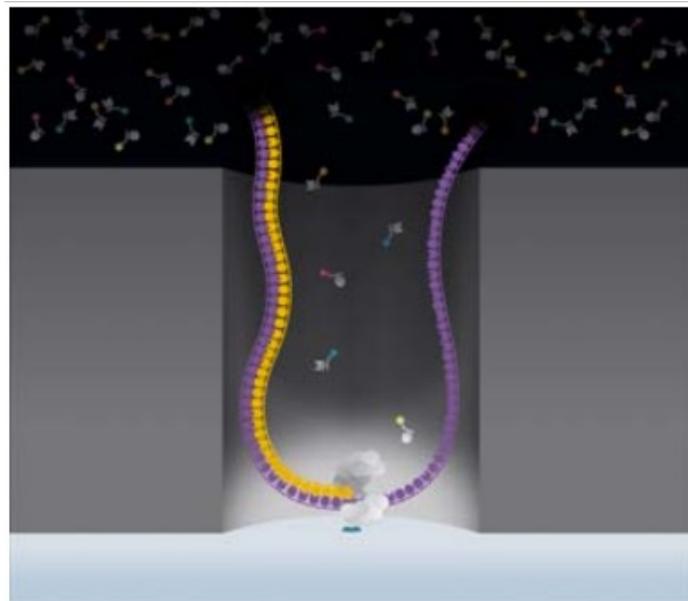


PacBio Single Molecule Real Time Sequencing (SMRT-sequencing)

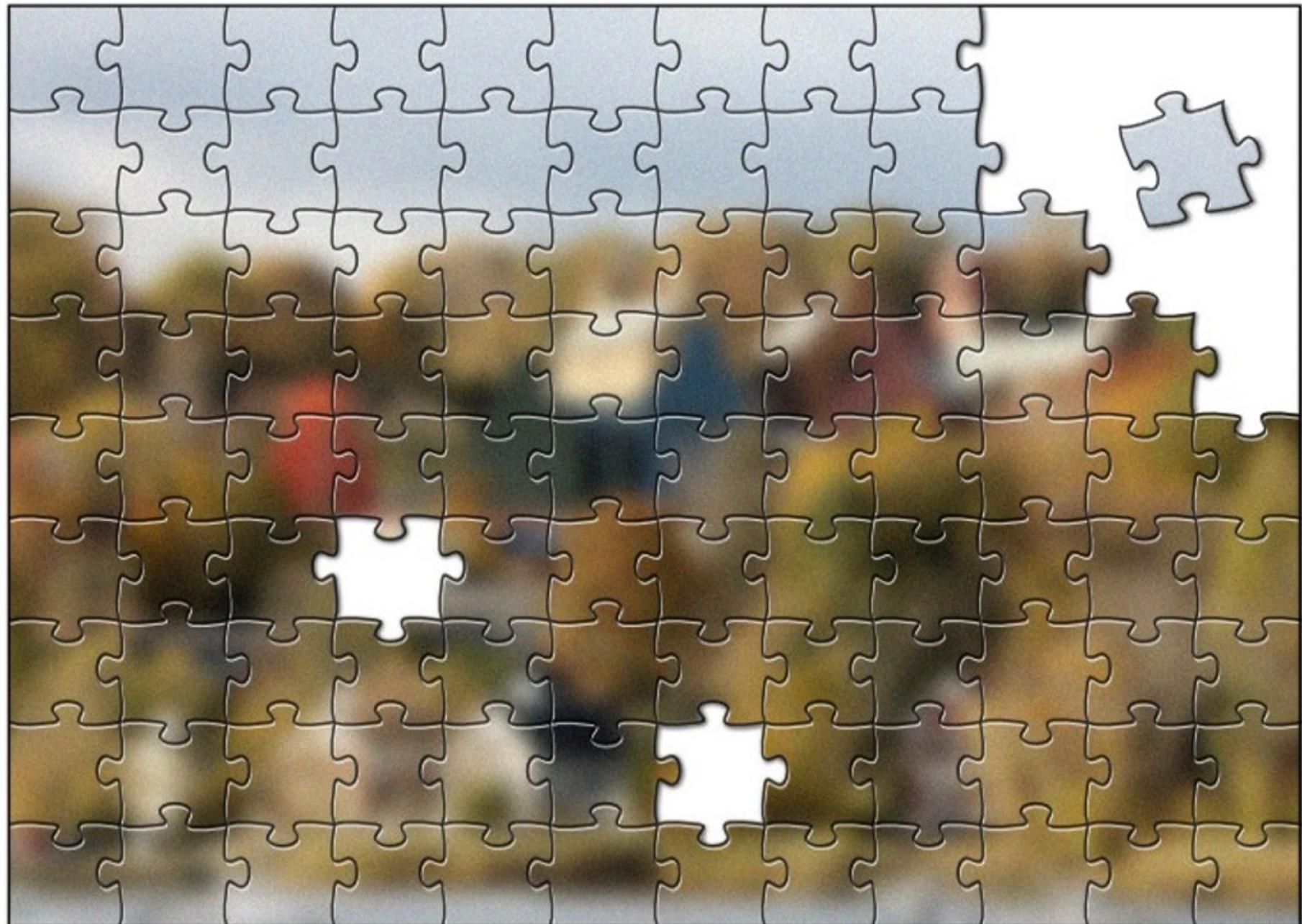


PacBio: SMRT Sequencing

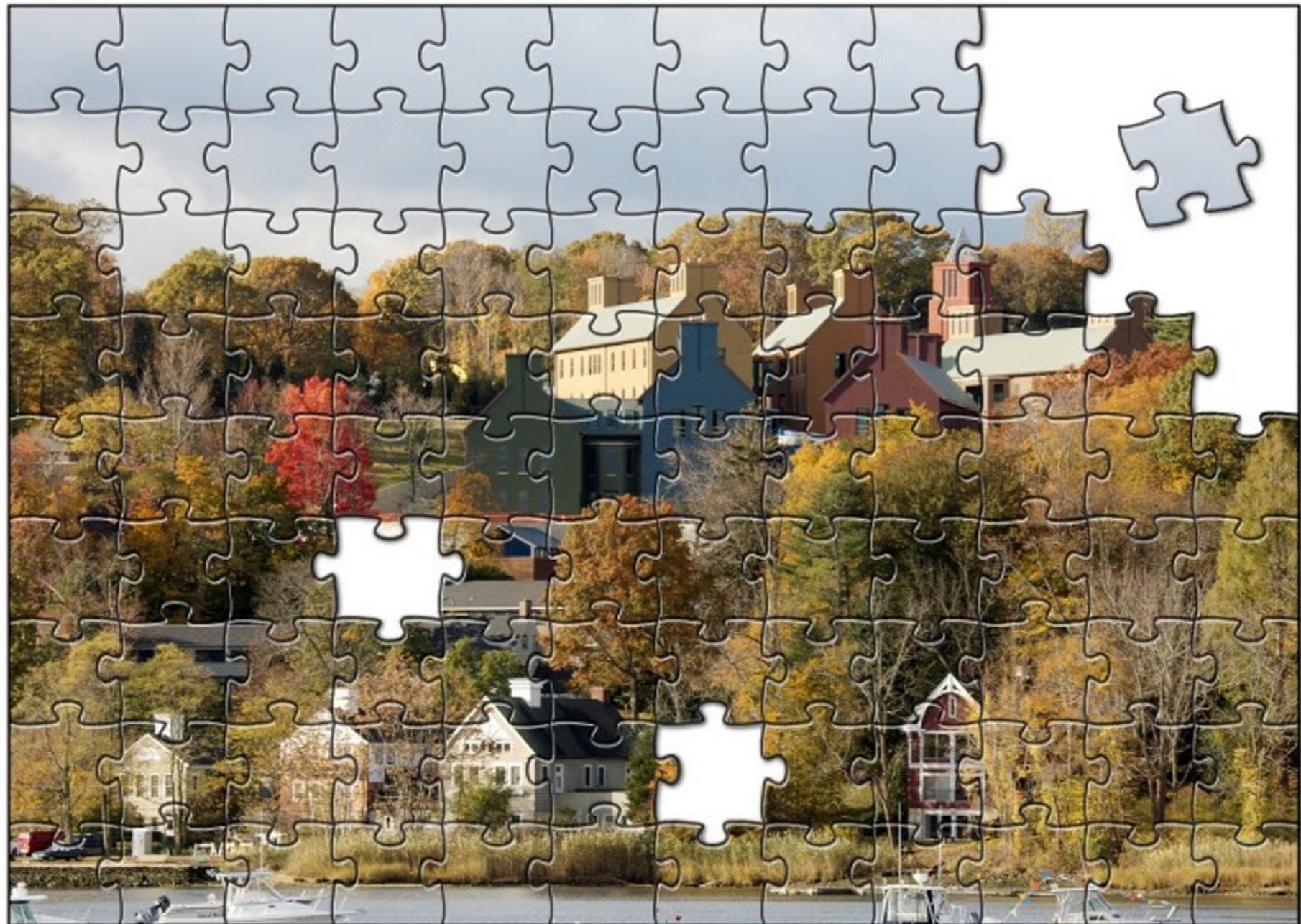
Imaging of fluorescent phospholinked labeled nucleotides as they are incorporated by a polymerase anchored to a Zero-Mode Waveguide (ZMW).



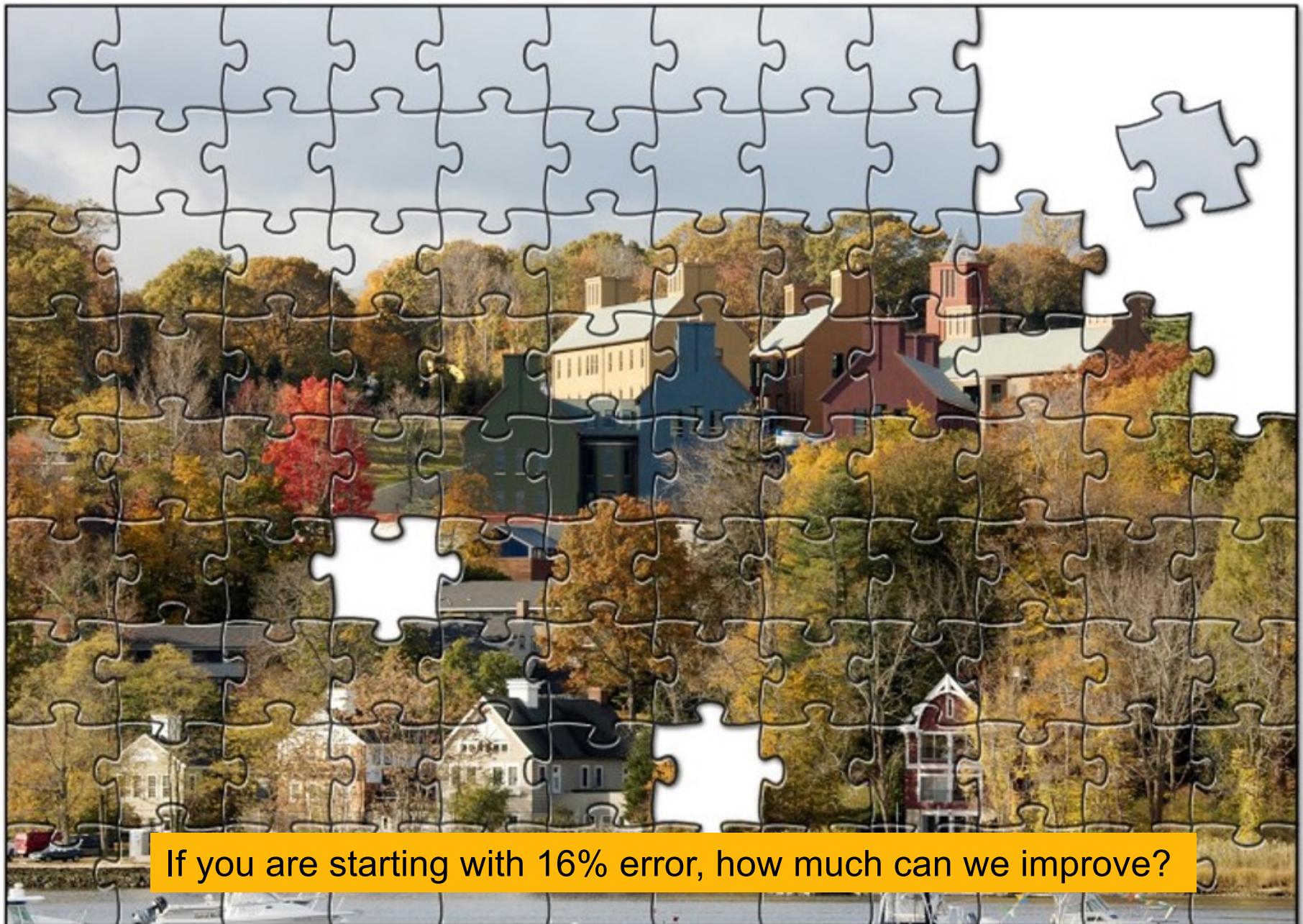
Single Molecule Sequences



“Corrective Lens” for Sequencing

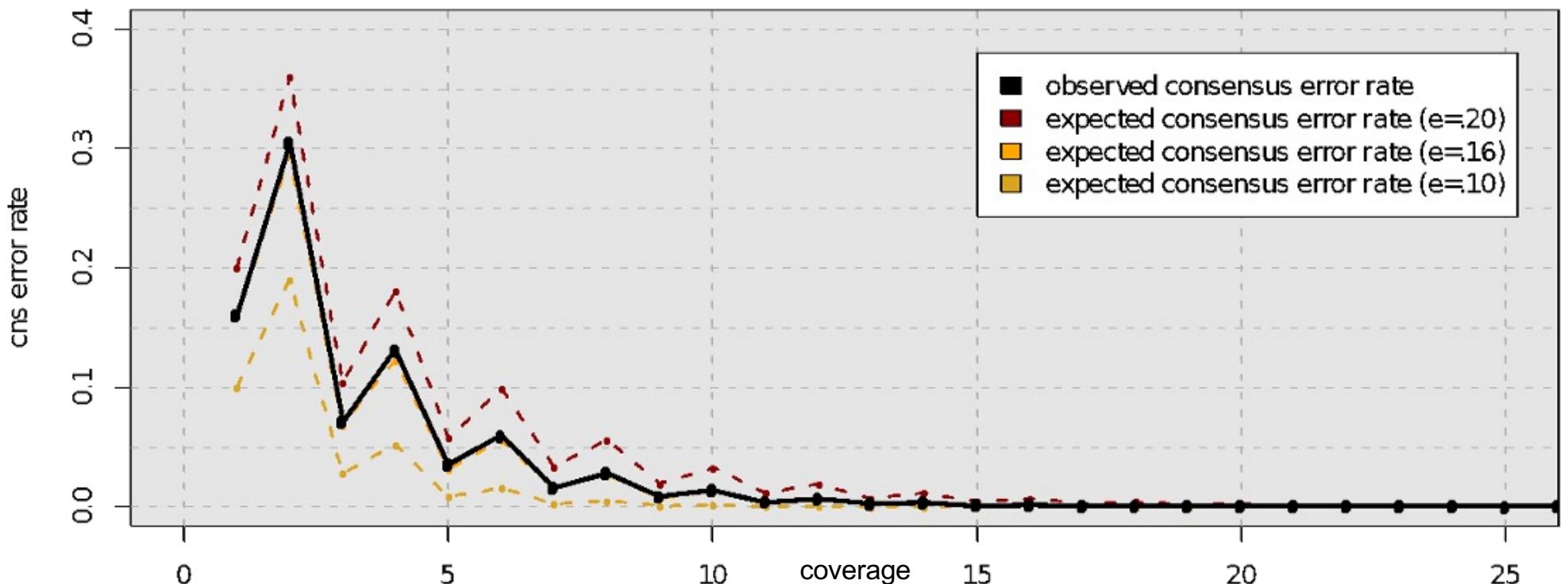


“Corrective Lens” for Sequencing



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

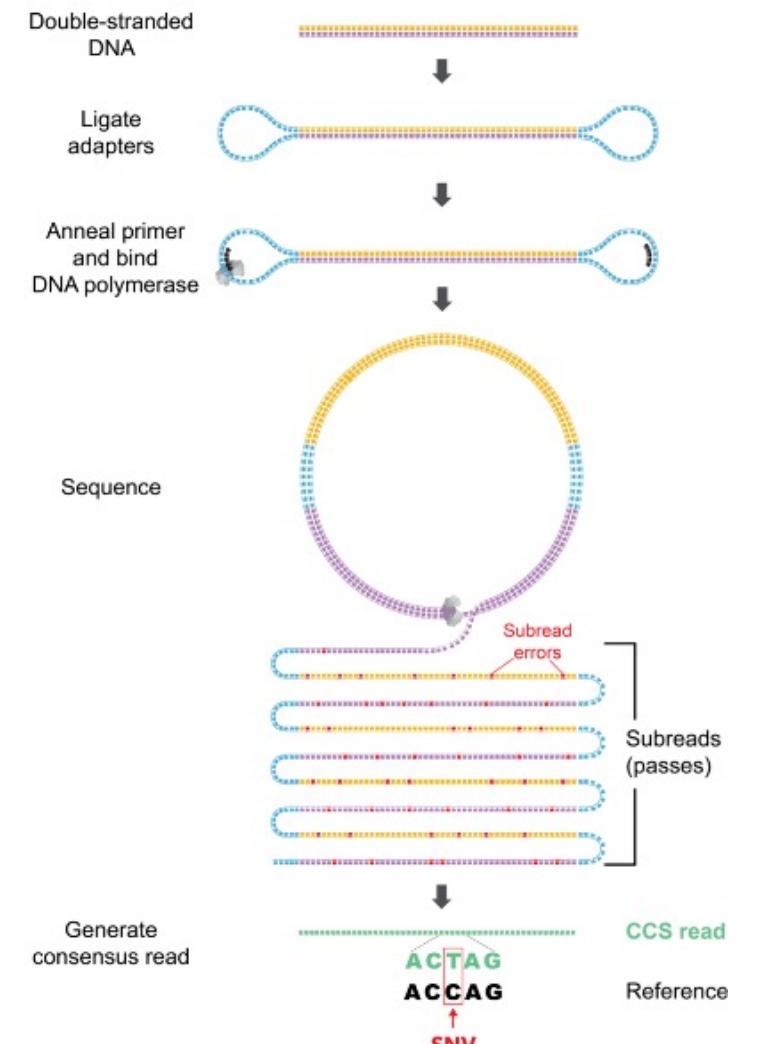
$$CNS Error = \sum_{i=\lceil c/2 \rceil}^c \binom{c}{i} (e)^i (1-e)^{n-i}$$

“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



Pacific Biosciences of California Inc

NASDAQ: PACB

Overview

News

Compare

Financials

Market Summary > Pacific Biosciences of California Inc

6.95 USD

-9.85 (-58.63%) ↓ all time

Closed: Sep 13, 5:18 AM EDT • Disclaimer

Pre-market 6.95 0.00 (0.00%)

1D 5D 1M 6M YTD 1Y 5Y Max



Open	6.75	Mkt cap	1.56B	52-wk high	31.10
High	7.06	P/E ratio	-	52-wk low	3.85
Low	6.48	Div yield	-		

Feedback

More about Pacific Biosciences o... →

About

pacb.com

Pacific Biosciences of California, Inc. is an American biotechnology company founded in 2004 that develops and manufactures systems for gene sequencing and some novel real time biological observation. [Wikipedia](#)

Founded: 2004

Headquarters: Menlo Park, CA

Number of employees: 728 (2021)

Revenue: 130.5 million USD (2021)

Predictions

Ceo

[Disclaimer](#)

Financials



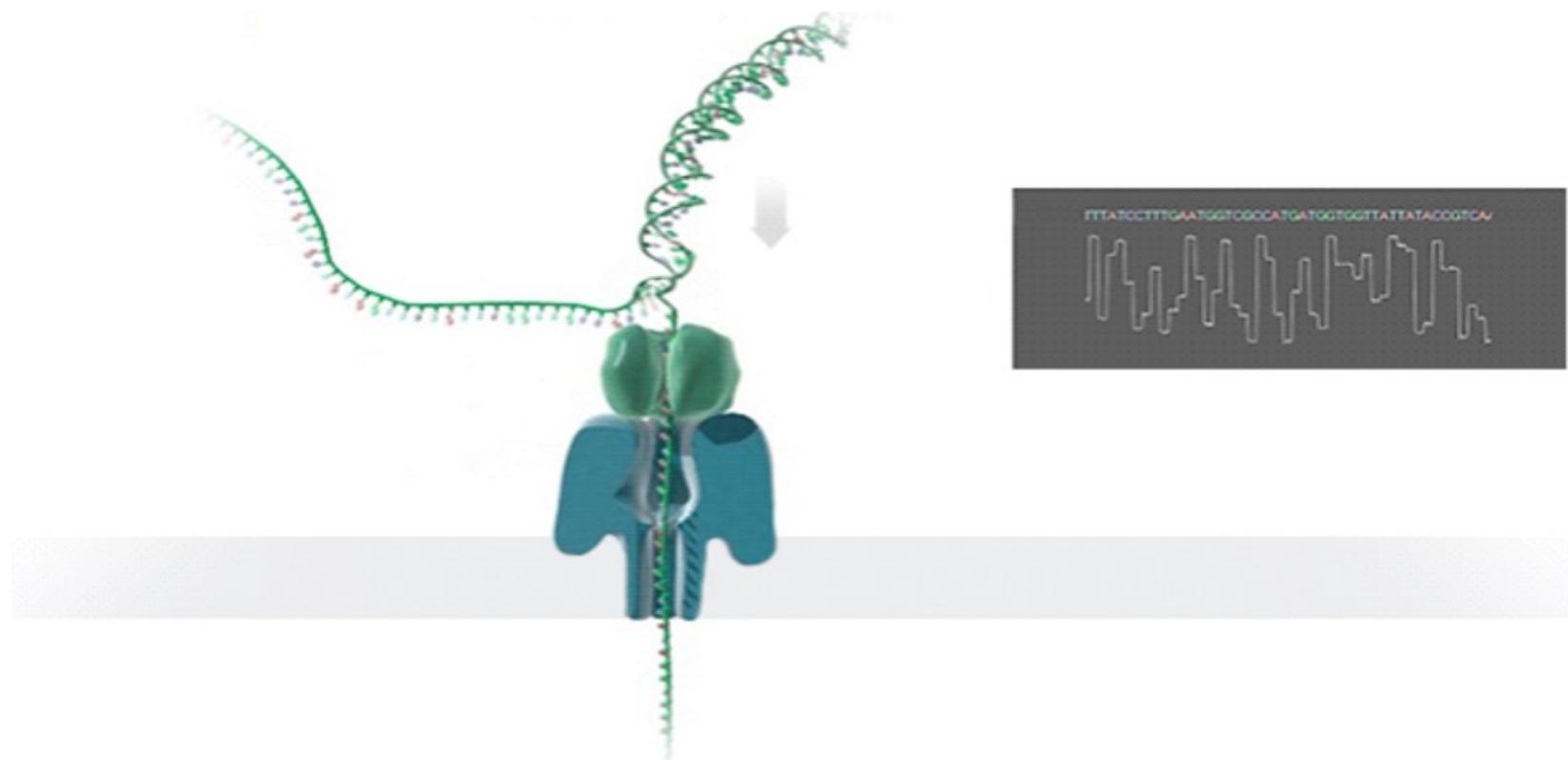
Quarterly financials

Oxford Nanopore Technologies (ONT)



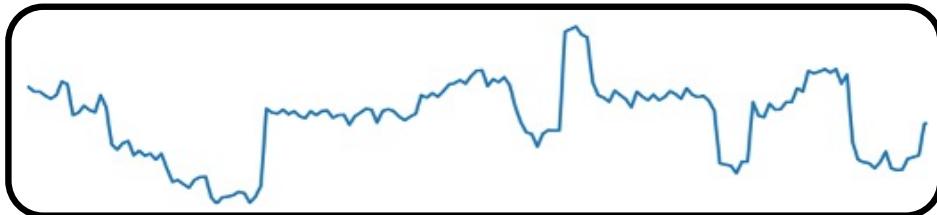
Nanopore Sequencing

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



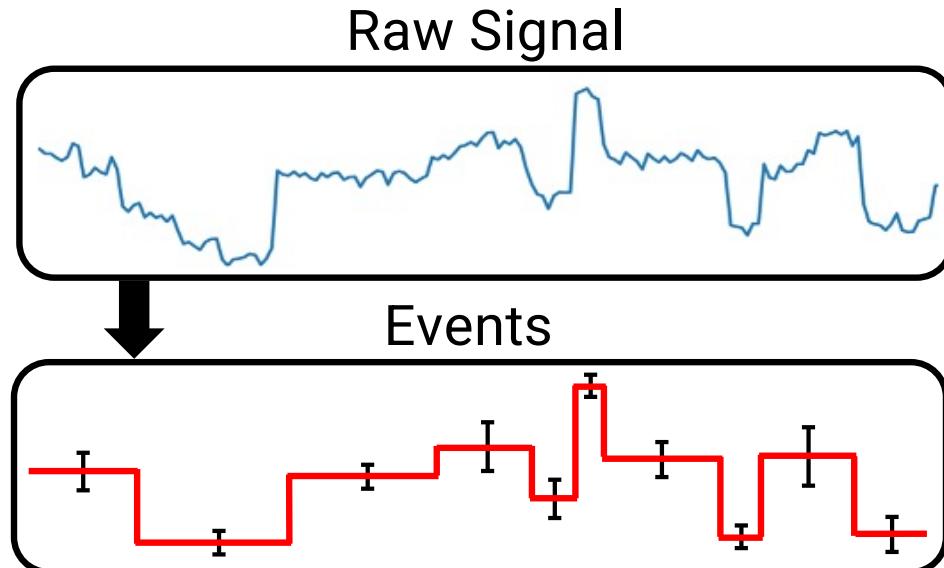
Nanopore Basecalling

Raw Signal



Translation of raw signal
into basepairs

Nanopore Basecalling

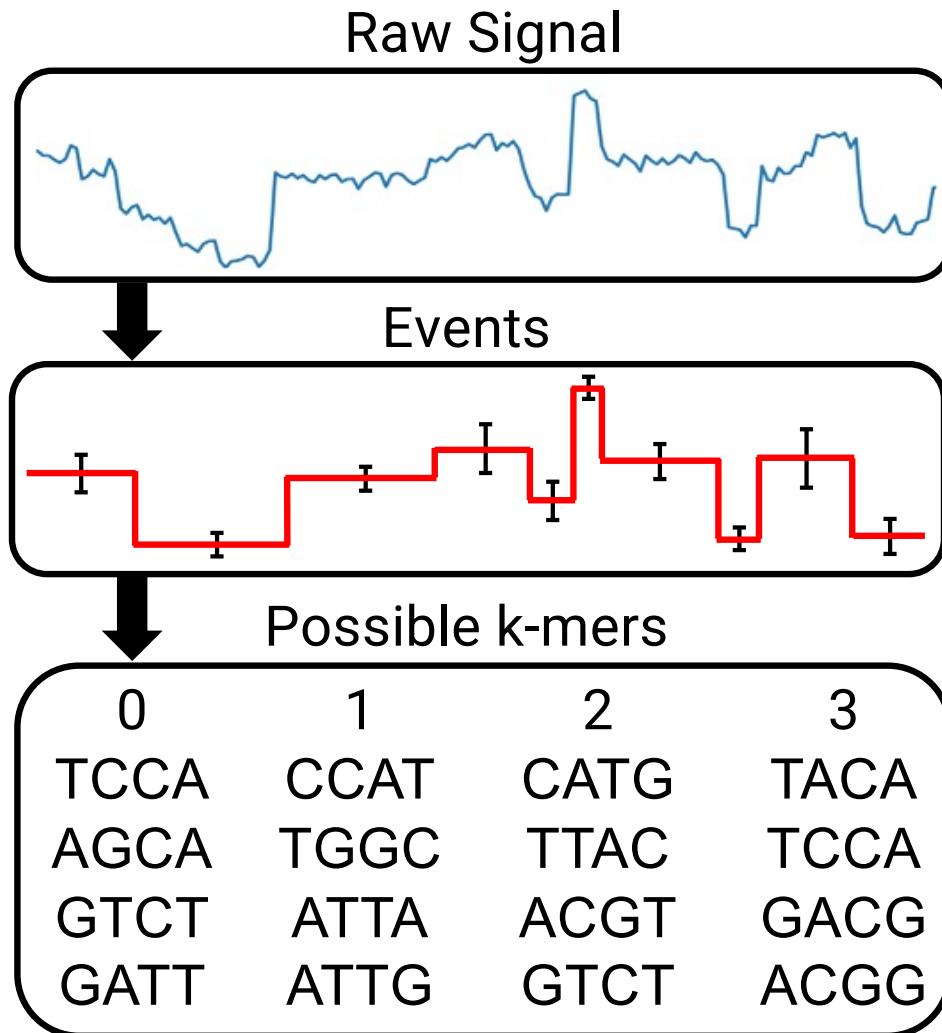


Translation of raw signal
into basepairs

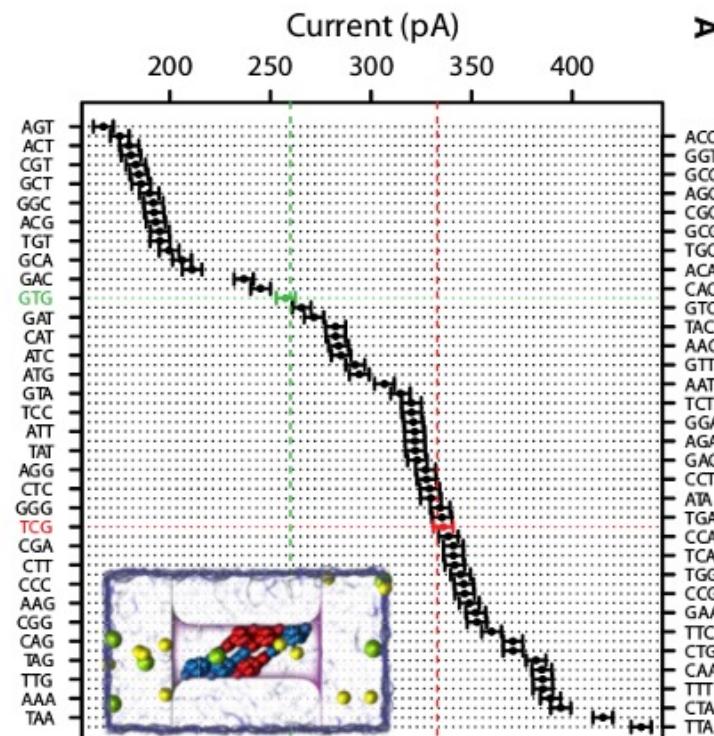
Early basecallers began by
estimating k-mer boundaries
using “events”, which were
then input to an HMM

Modern basecallers use
neural networks directly
on raw signal

Nanopore Basecalling

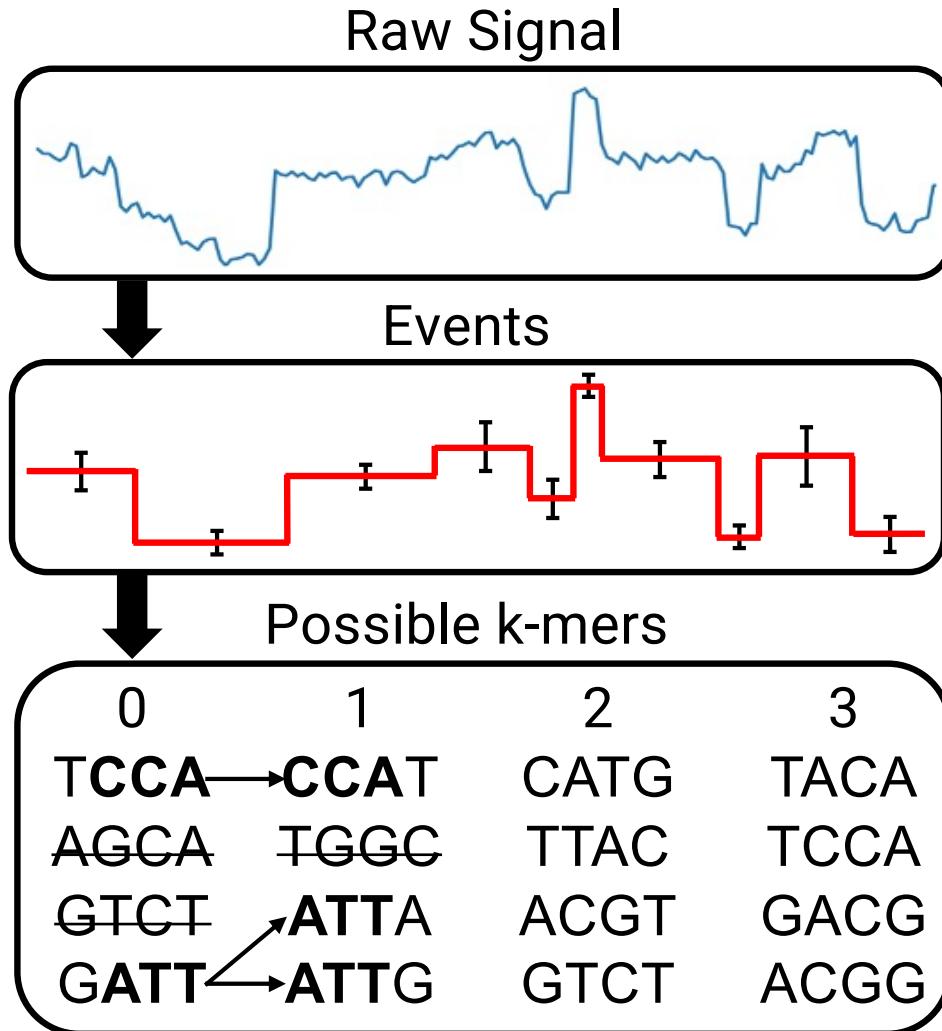


ONT releases k-mer models with expected current distribution of every k-mer

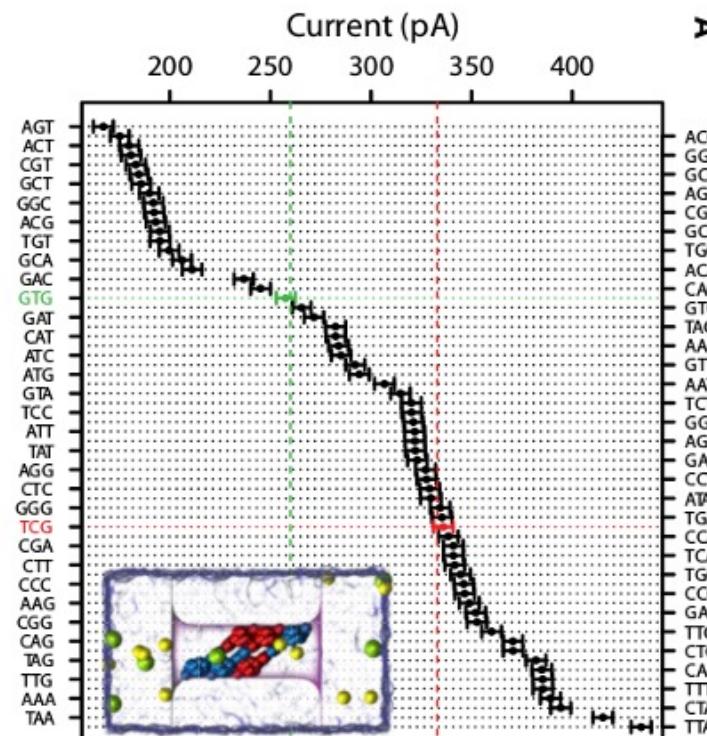


DNA Base-Calling from a Nanopore Using a Viterbi Algorithm
Timp et al. (2012) *Biophysical Journal*

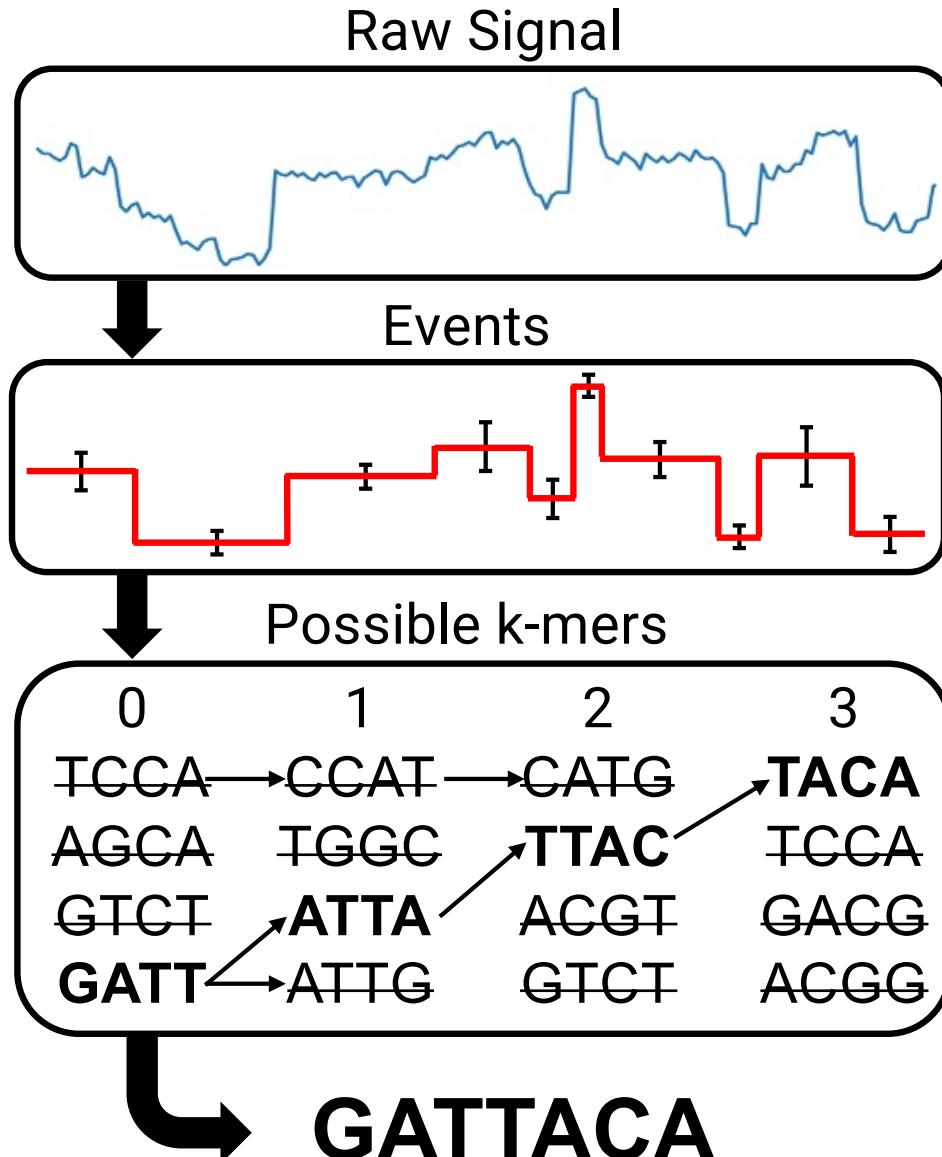
Nanopore Basecalling



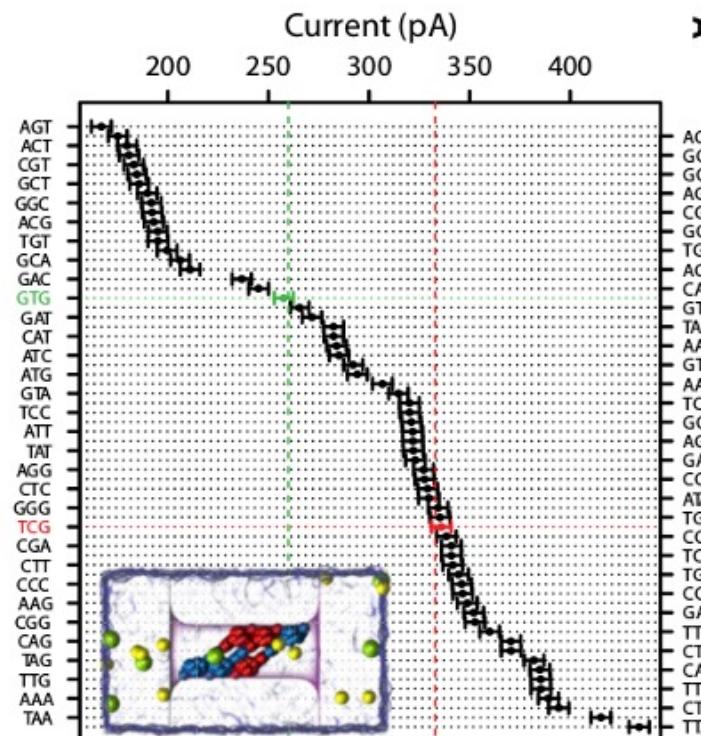
Certain k-mers can be eliminated based on possible transitions



Nanopore Basecalling



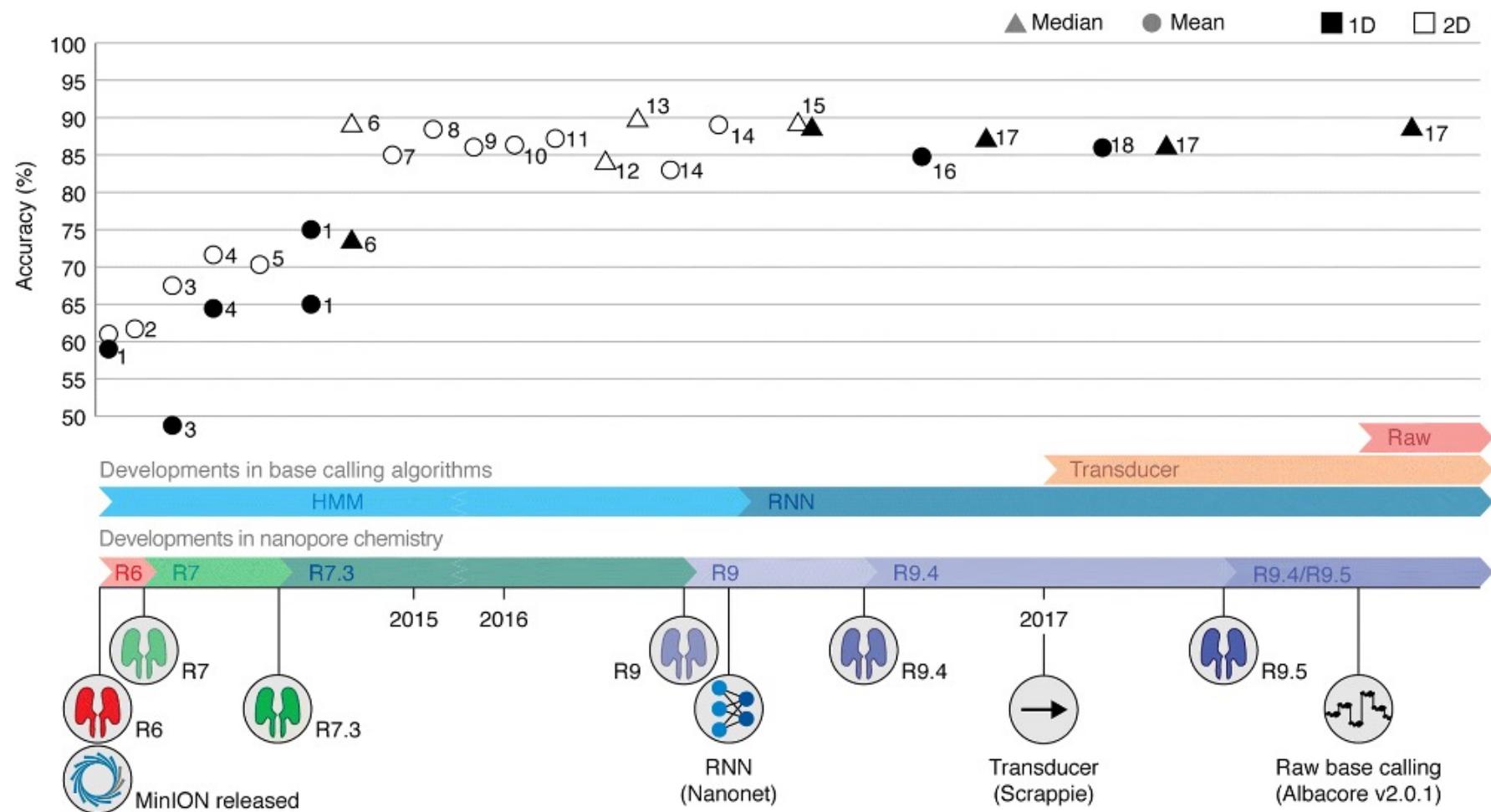
Final sequence determined by most probable k-mers



"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>



Oxford Nanopore Technologies PLC

OTCMKTS: ONTTF :

Overview

News

Compare

Financials

Market Summary > Oxford Nanopore Technologies PLC

3.32 USD

+ Follow

-5.25 (-61.26%) ↓ all time

Sep 12, 4:00 PM EDT • Disclaimer

1D 5D 1M 6M YTD 1Y 5Y Max



Open	3.36	Mkt cap	2.46B GBP	52-wk high	10.20
High	3.36	P/E ratio	-	52-wk low	2.95
Low	3.30	Div yield	-		

More about Oxford Nanopore Tec... →

Feedback

News >

About

nanoporetech.com

Oxford Nanopore Technologies Limited is a UK-based company which is developing and selling nanopore sequencing products for the direct, electronic analysis of single molecules. [Wikipedia](#)

CEO: Gordon Sanghera (May 2005–)

Founded: 2005

Founder: Hagan Bayley

Number of employees: 803 (2021)

Subsidiary: Oxford Nanopore Technologies, Inc.

Headquarters: United Kingdom, Oxford, United Kingdom

Disclaimer

Financials

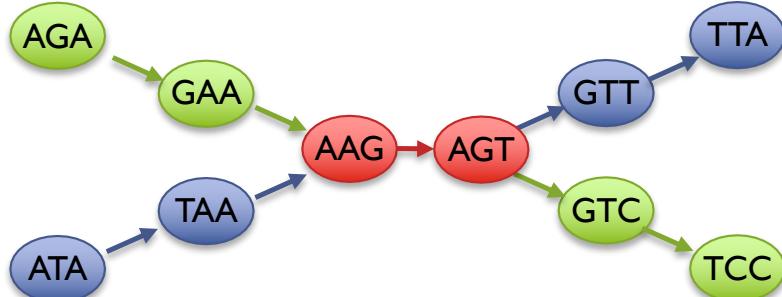


Quarterly financials

(GBP)	Dec 2021	Y/Y
Revenue	37.36M	13.97% ↑

Two Paradigms for Assembly

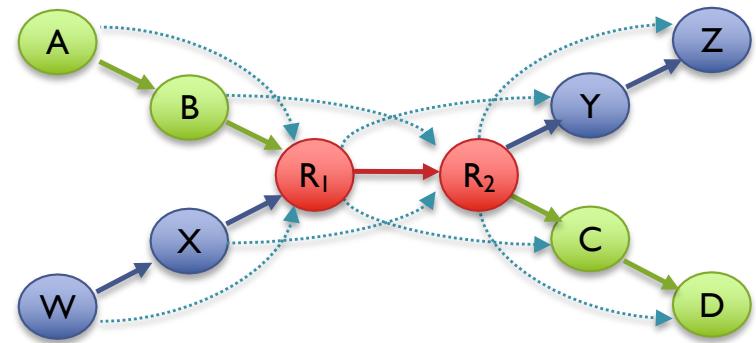
de Bruijn Graph



Short read assemblers

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

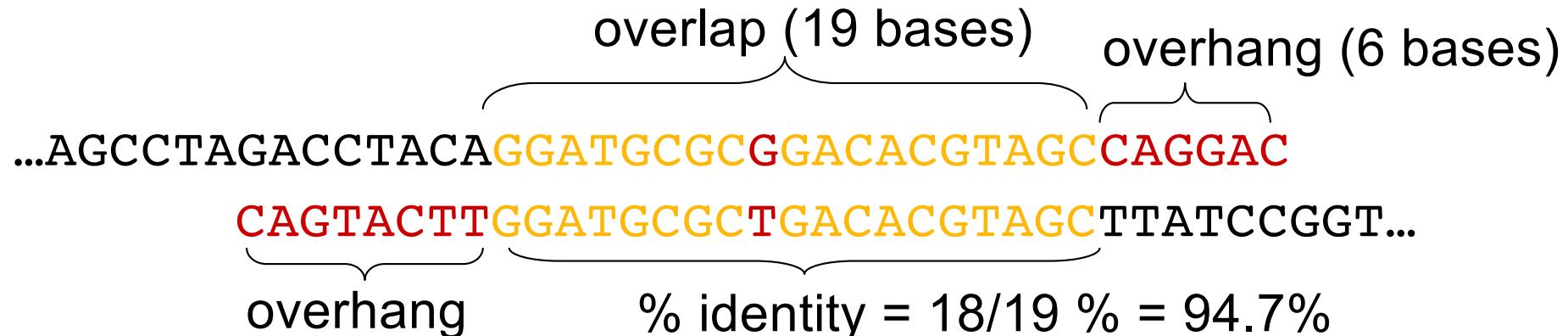
Overlap Graph



Long read assemblers

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

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

