

A complete reference genome improves analysis of human genetic variation

Michael Schatz

November 18, 2021
CSHL Advanced Sequencing Course



 @mike_schatz

Schatzlab Overview



Human Disease Genetics

Aganezov *et al.* (2020)
Feigin *et al.* (2017)



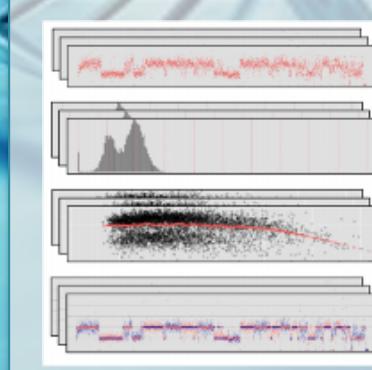
Agricultural Genomics

Alonge *et. al.* (2020)
Chen *et al.* (2019)



Algorithmics & Systems Research

Kirsche *et al.* (2020)
Schatz (2009)

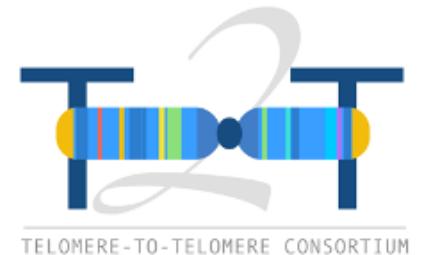


Single Cell & Single Molecule

Kovaka *et al.* (2020)
Sedlazeck *et al.* (2018)

Outline

- Assembly by analogy
- T2T-CHM13 Assembly
 - How the T2T-CHM13 reference improves the analysis of genetic variation
- Beyond T2T-CHM13



Shredded Book Reconstruction

- Dickens accidentally shreds the first printing of A Tale of Two Cities
 - Text printed on 5 long spools

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

- How can he reconstruct the text?

- 5 copies x 138,656 words / 5 words per fragment = 138k fragments
- The short fragments from every copy are mixed together
- Some fragments are identical

Greedy Reconstruction

It was the best of
age of wisdom, it was
best of times, it was
it was the age of
it was the age of
it was the worst of
of times, it was the
of times, it was the
of wisdom, it was the
the age of wisdom, it
the best of times, it
the worst of times, it
times, it was the age
times, it was the worst
was the age of wisdom,
was the age of foolishness,
was the best of times,
was the worst of times,
wisdom, it was the age
worst of times, it was

It was the best of
was the best of times,
the best of times, it
best of times, it was
of times, it was the
of times, it was the
times, it was the worst
times, it was the age

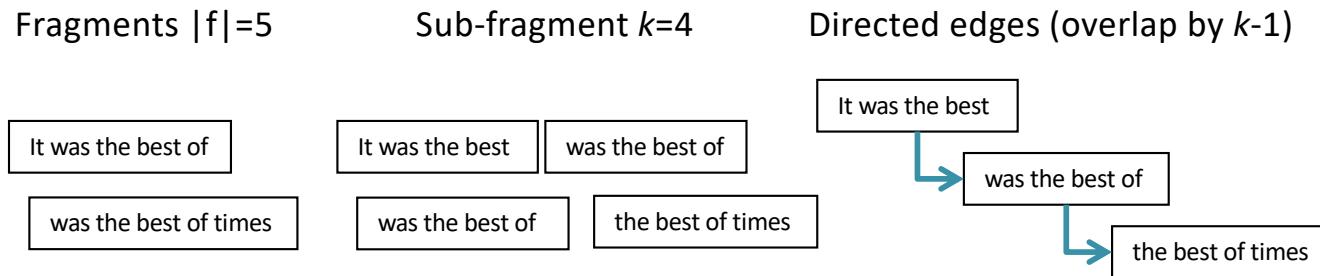
The repeated sequence make the correct reconstruction ambiguous

- It was the best of times, it was the [worst/age]

Model the assembly problem as a graph problem

de Bruijn Graph Construction

- $G_k = (V, E)$
 - V = Length- k sub-fragments
 - E = Directed edges between consecutive sub-fragments
 - Sub-fragments overlap by $k-1$ words



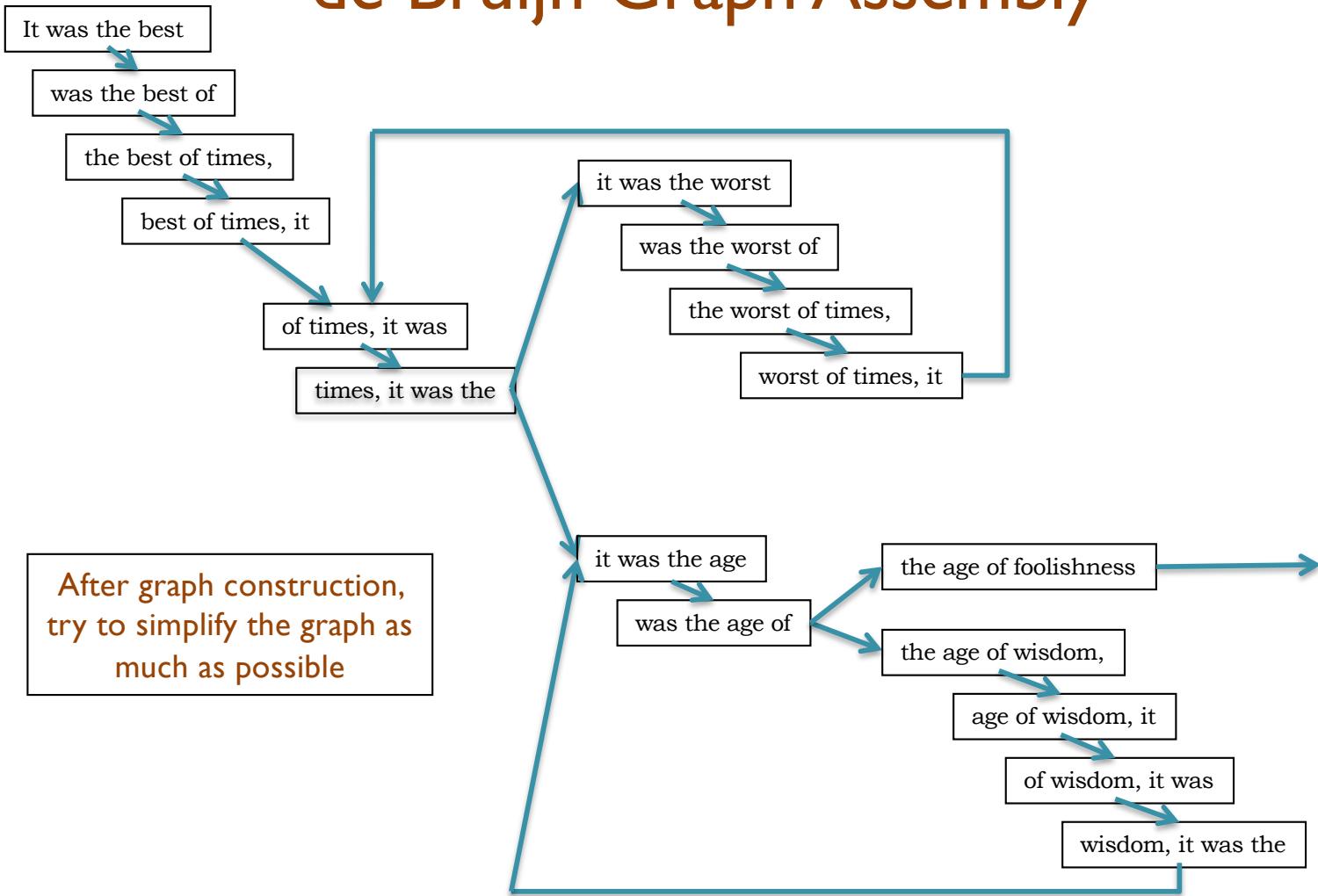
- Overlaps between fragments are implicitly computed

How to pronounce:

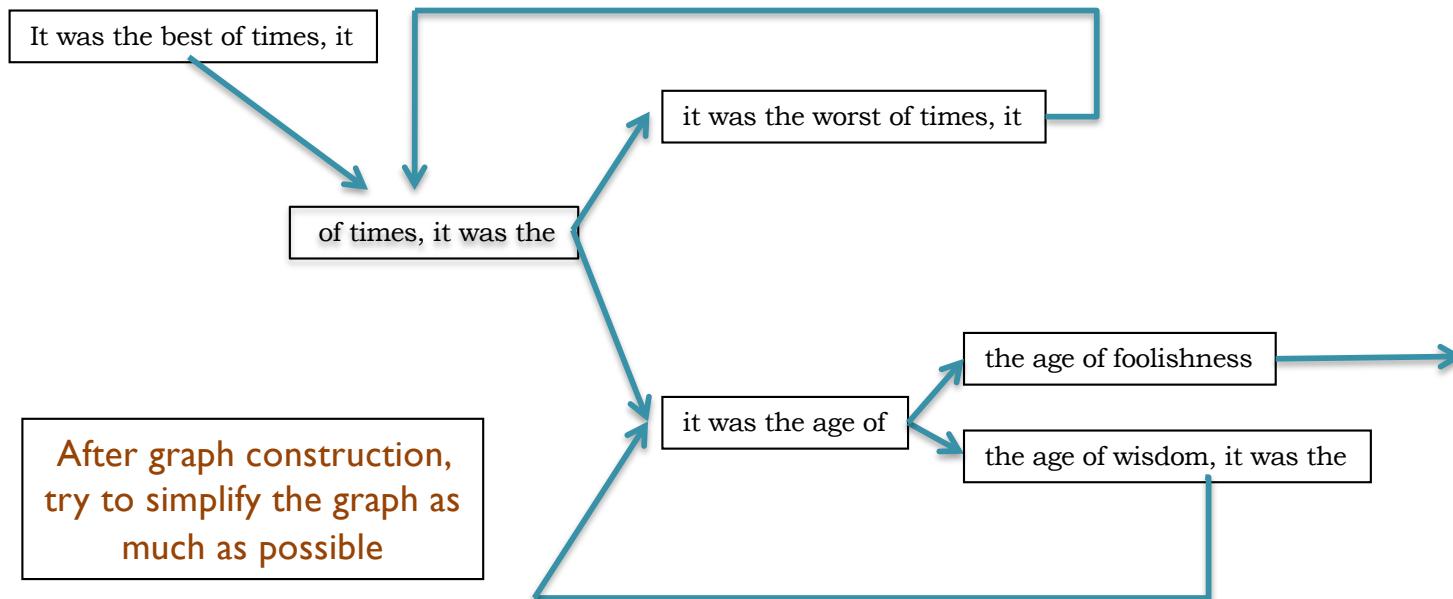
https://forvo.com/word/de_briuin/

de Bruijn, 1946
Idury et al., 1995
Pevzner et al., 2001

de Bruijn Graph Assembly



de Bruijn Graph Assembly



The full tale

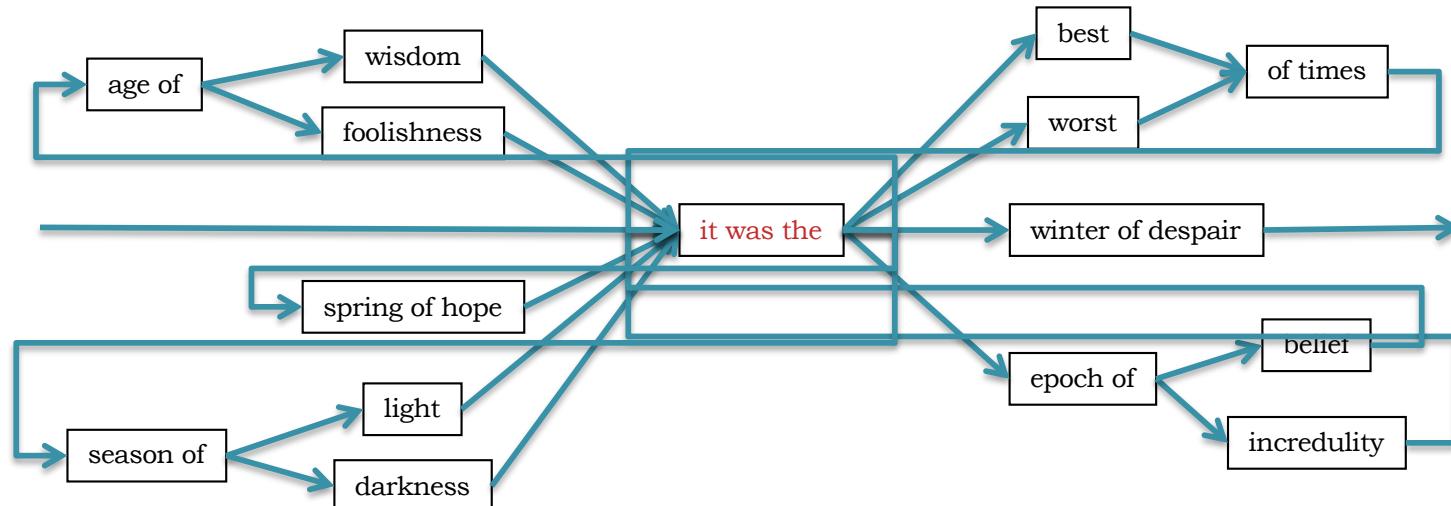
... it was the best of times it was the worst of times ...

... it was the age of wisdom it was the age of foolishness ...

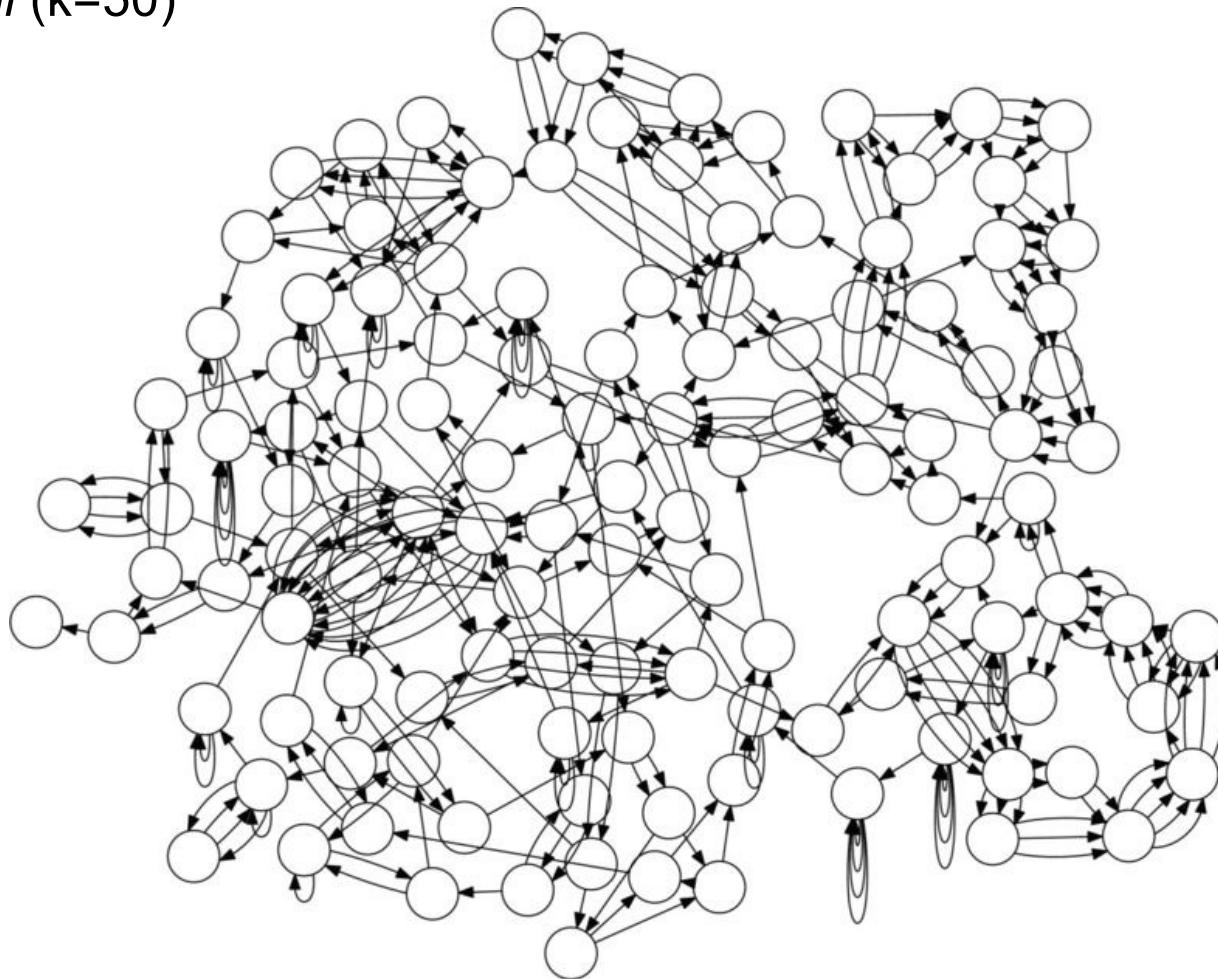
... it was the epoch of belief it was the epoch of incredulity ...

... it was the season of light it was the season of darkness ...

... it was the spring of hope it was the winter of despair ...

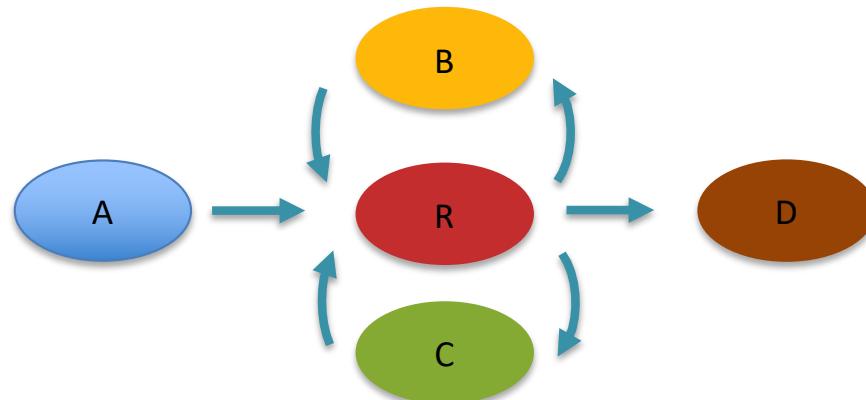
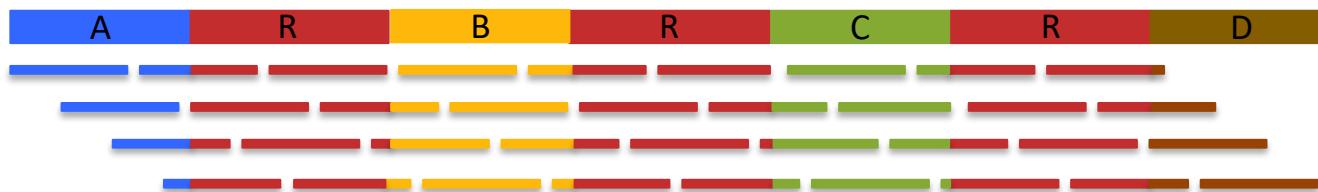


E. coli ($k=50$)

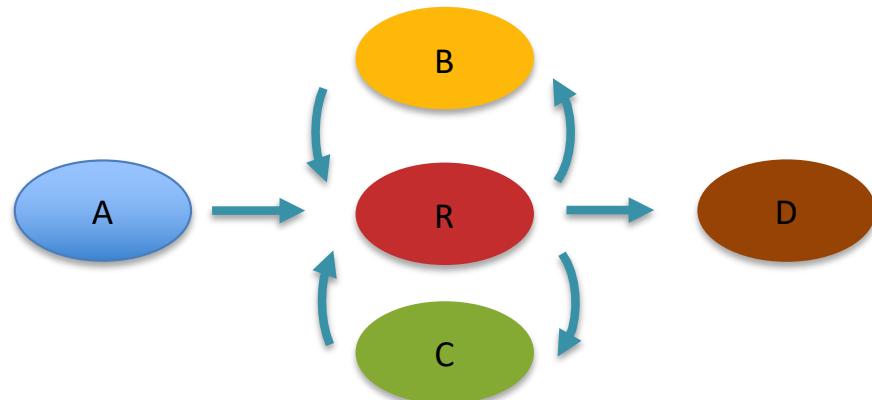
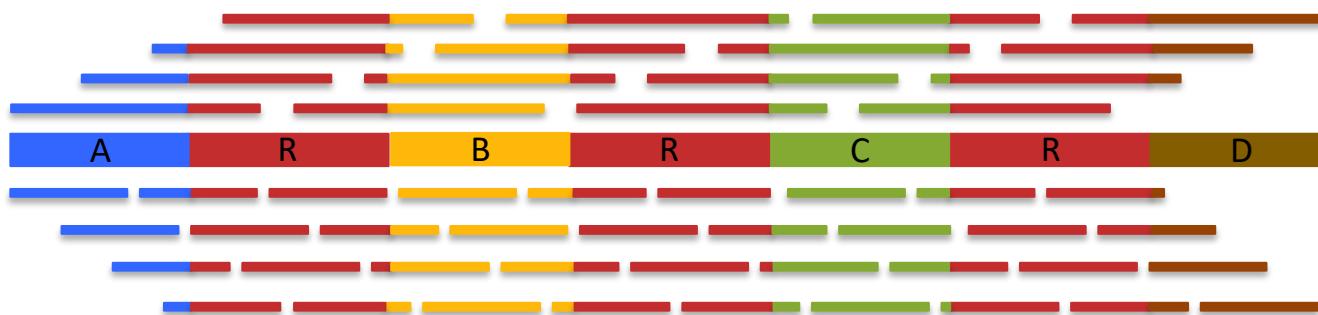


Reducing assembly complexity of microbial genomes with single-molecule sequencing
Koren et al (2013) Genome Biology. **14**:R101 <https://doi.org/10.1186/gb-2013-14-9-r101>

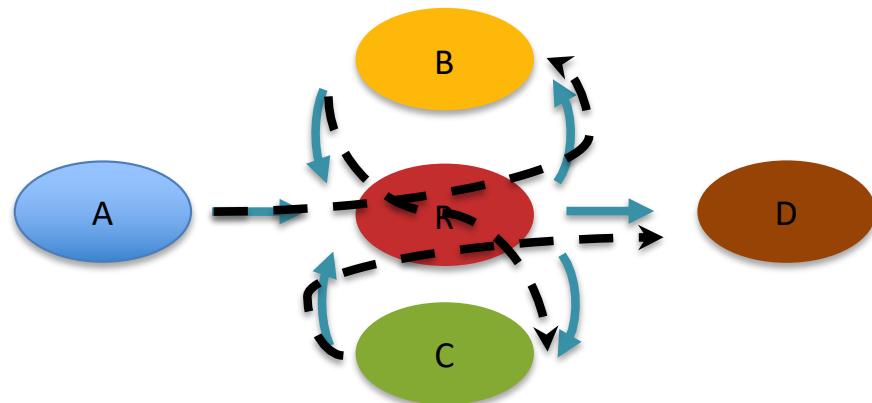
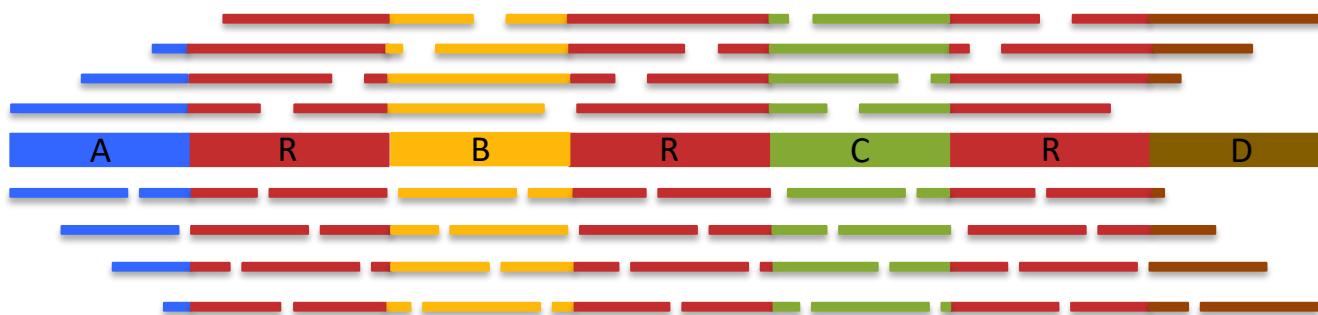
Assembly Complexity



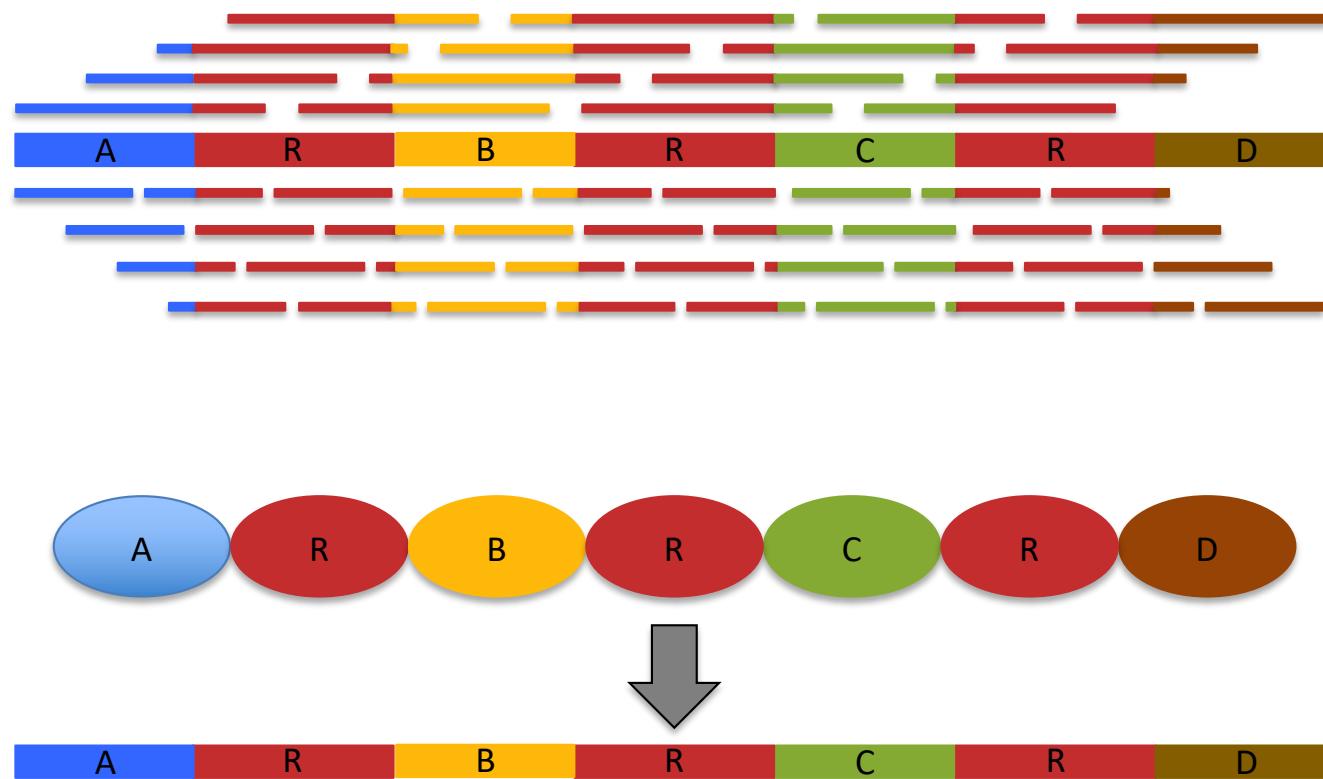
Assembly Complexity



Assembly Complexity



Assembly Complexity

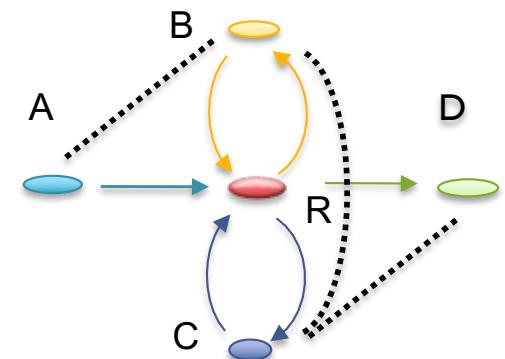


The advantages of SMRT sequencing

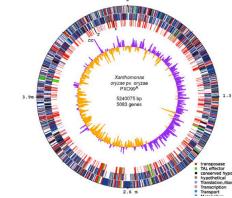
Roberts, RJ, Carneiro, MO, Schatz, MC (2013) *Genome Biology*. 14:405

Scaffolding

- Initial contigs (aka unipaths, unitigs) terminate at
 - Coverage gaps: especially extreme GC
 - Conflicts: errors, repeat boundaries, heterozygosity
- Use mate-pairs or other mapping data to resolve correct order through assembly graph
 - Place sequence to satisfy the mate constraints
 - Mates through repeat nodes are tangled
- Final scaffold may have internal gaps called “sequencing gaps”
 - We know the order, orientation, and spacing, but just not the bases.
Fill with Ns instead



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

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

ATTA

GATT

TACA

TTAC

Pop Quiz

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

ATTA: ATT → TTA

GATT: GAT → ATT

TACA: TAC → ACA

TTAC: TTA → TAC

Pop Quiz

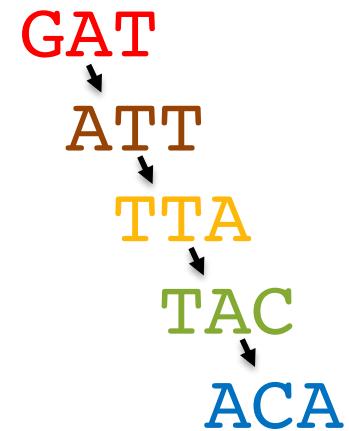
Assemble these reads using a de Bruijn graph approach (k=3):

ATTA: ATT → TTA

GATT: GAT → ATT

TACA: TAC → ACA

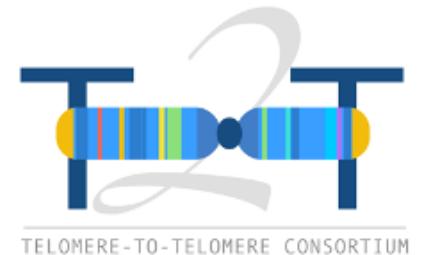
TTAC: TTA → TAC



GATTACA

Outline

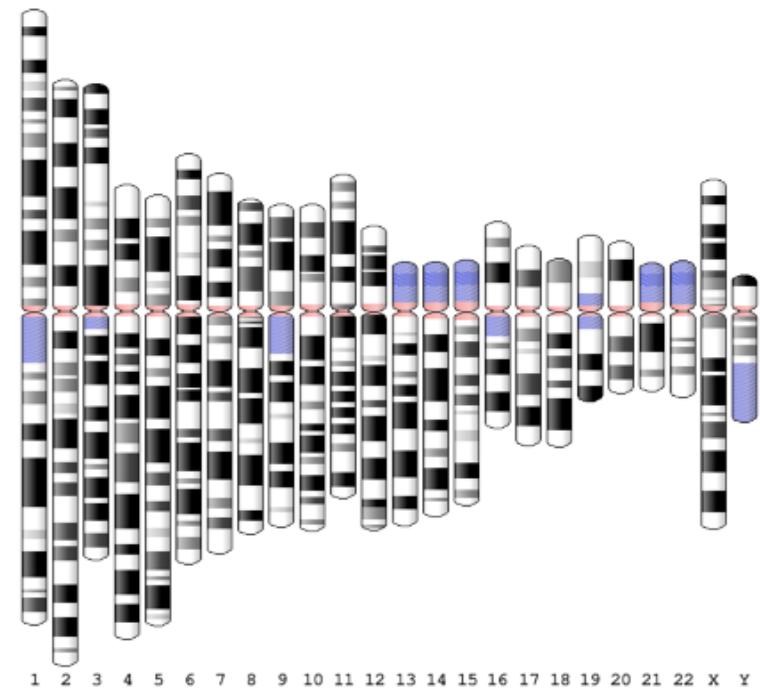
- Assembly by analogy
- T2T-CHM13 Assembly
 - How the T2T-CHM13 reference improves the analysis of genetic variation
- Beyond T2T-CHM13



The Human Genome Today

About 8% is missing or incorrect

- Over 100M “Ns” in the reference
- Centromeres and telomeres
- Segmentally duplicated genes
- Tandem gene arrays (e.g. rDNAs)
- Thousands of haplotype switches
- And an unknown number of errors...



Finishing the human genome

Why does it matter?

Variation in these regions is unexplored

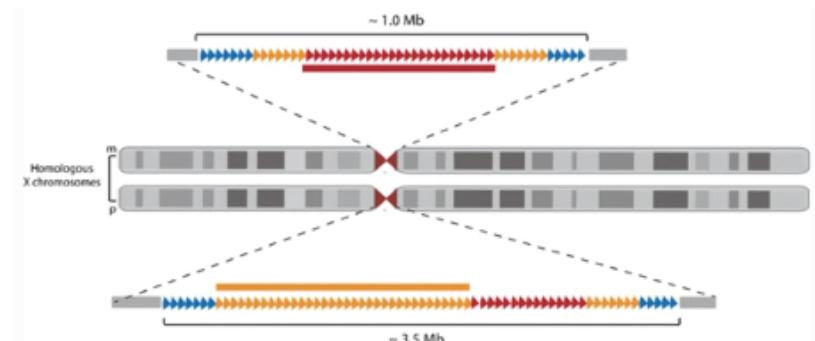
Functional studies need sequence

Reference gaps lead to artifacts

We don't know what we don't know...

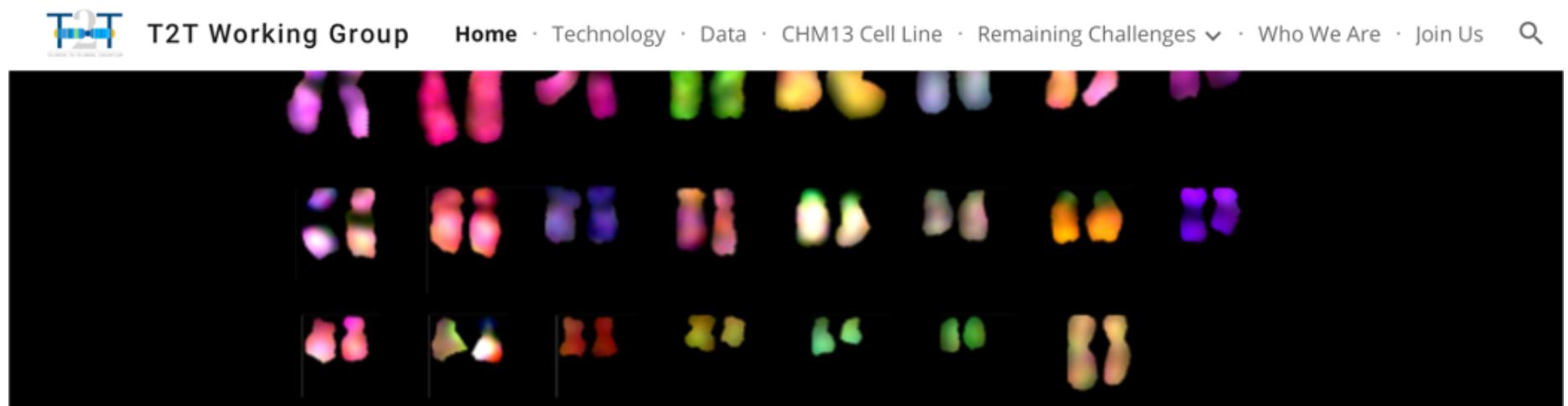
Why has it taken so long?

Repeats, repeats, repeats...



Miga
2015

Let's finish a human genome!



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

Single Molecule Long Read Sequencing

PacBio
Sequel II



Oxford Nanopore
PromethION





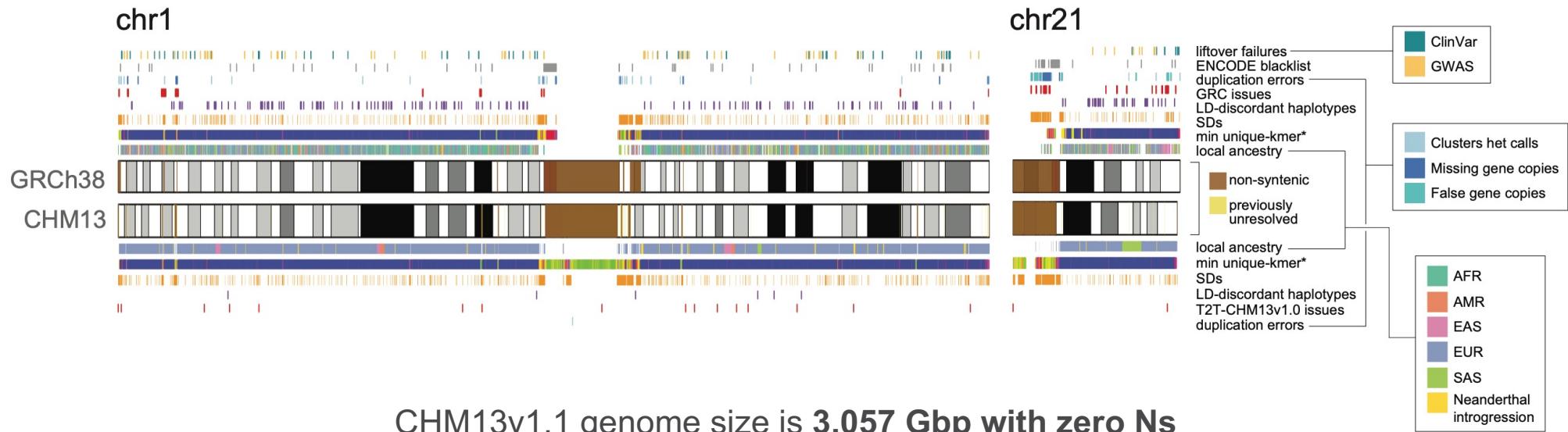
CHM13 HiFi assembly graph



The complete sequence of a human genome

Nurk et al (2021) *bioRxiv*. doi: <https://doi.org/10.1101/2021.05.26.445798>

The complete sequence of a human genome



The complete sequence of a human genome

Nurk et al (2021) *bioRxiv*. doi: <https://doi.org/10.1101/2021.05.26.445798>

T2T-Variants Analysis Summary



- Large-scale cloud analysis of short-read data
- Short-read alignment and variant calling
- Long-read alignment and variant calling
- Variants in previously unresolved regions
- Clinical implications of variants



Sergey Aganezov



Stephanie Yan



Daniela Soto



Melanie Kirsche



Samantha Zarate

A complete reference genome improves analysis of human genetic variation

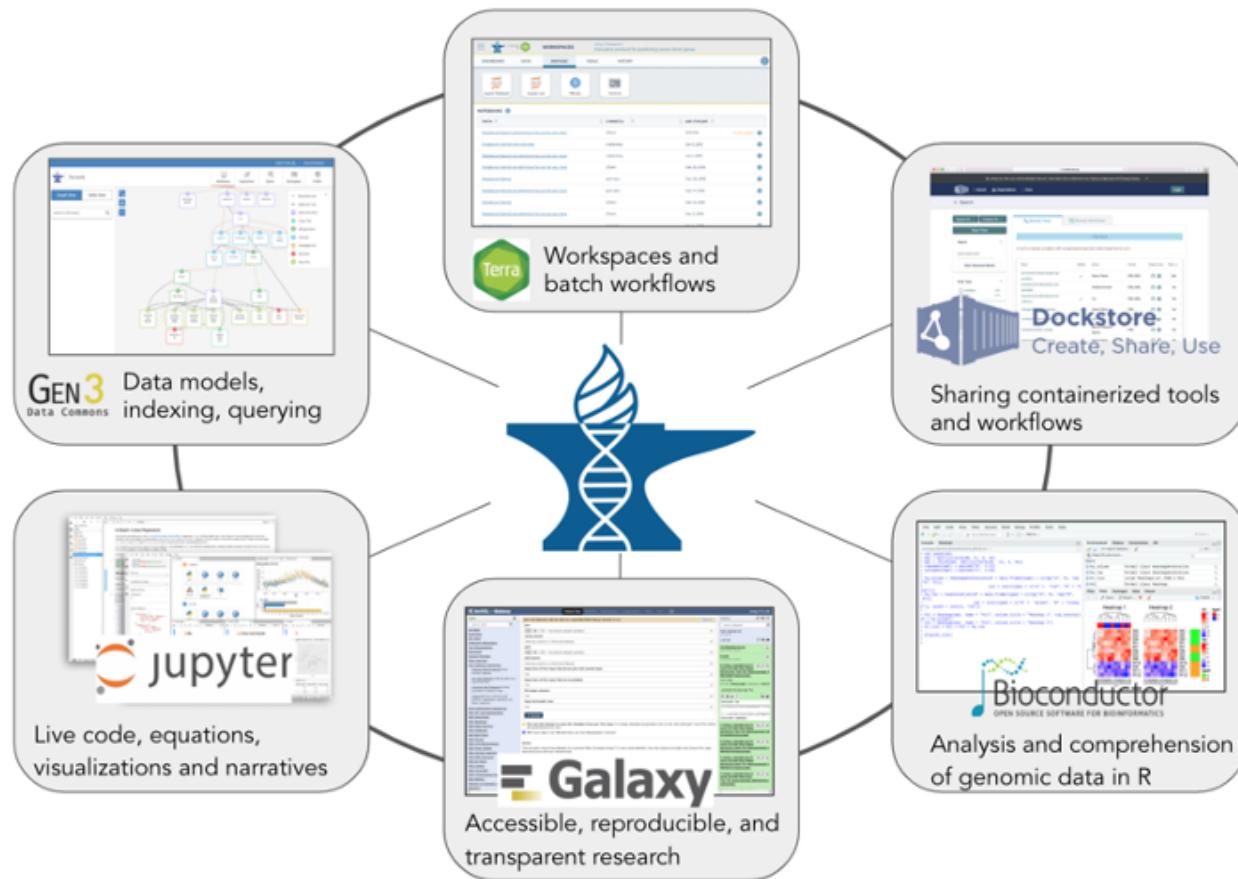
Aganezov, S*, Yan, SM*, Soto, DC*, Kirsche, M*, Zarate, S*, et al. (2021) *bioRxiv*. doi: <https://doi.org/10.1101/2021.07.12.452063>

T2T Variants: 1000 Genomes Project

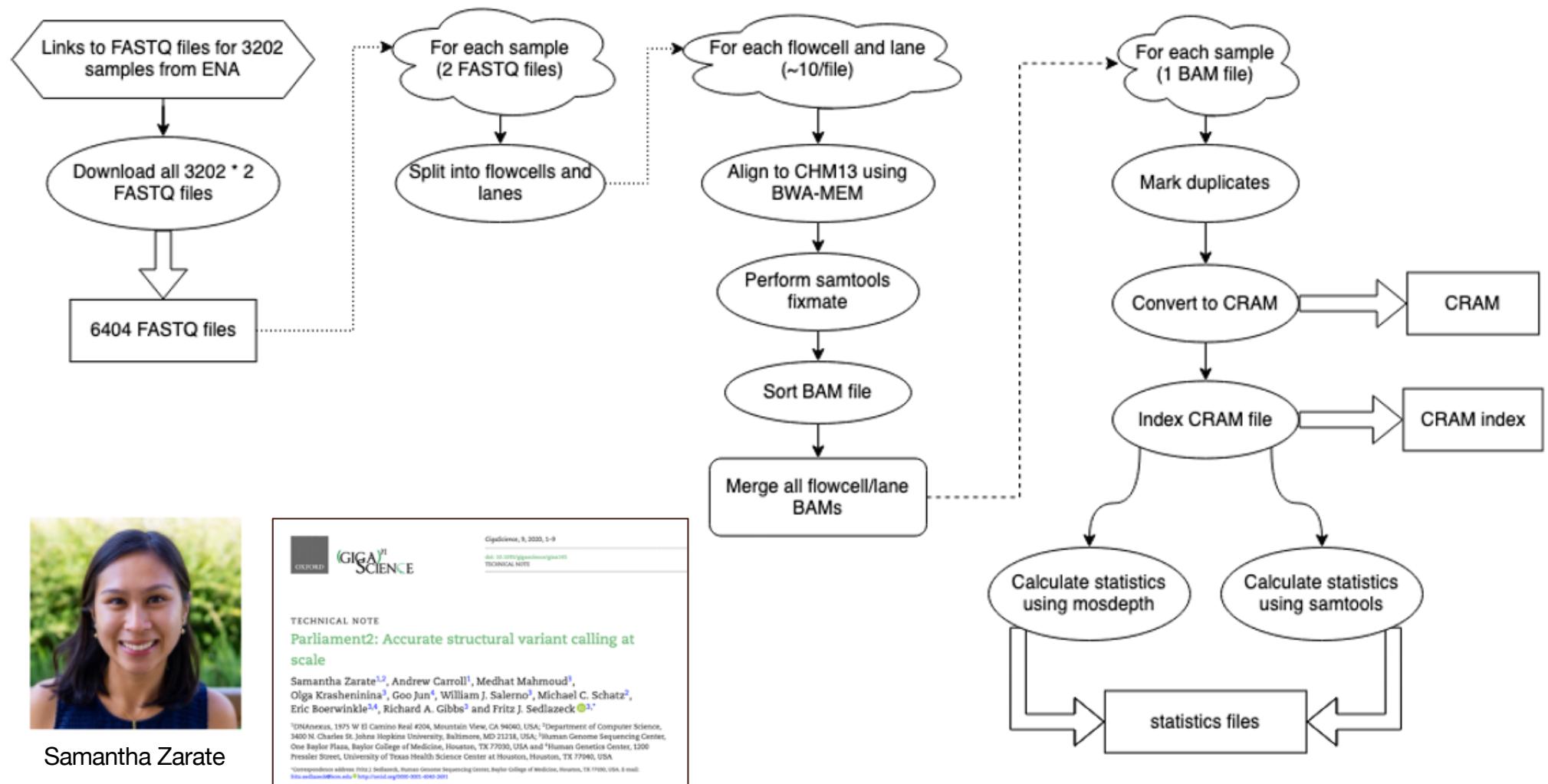


26 populations from 5 superpopulations (continental regions)
3202 samples (2504 core genomes + 698 offspring)

3202 samples x 30Gb = 96Tb input data | >5Pb of intermediate data | >>1M core hours



**Inverting the model of genomics data sharing with the
NHGRI Genomic Data Science Analysis, Visualization, and Informatics Lab-space (AnVIL)**
Schatz, Philippakis et al. (2021) *bioRxiv* doi: <https://doi.org/10.1101/2021.04.22.436044>



Core usage over 24 hours



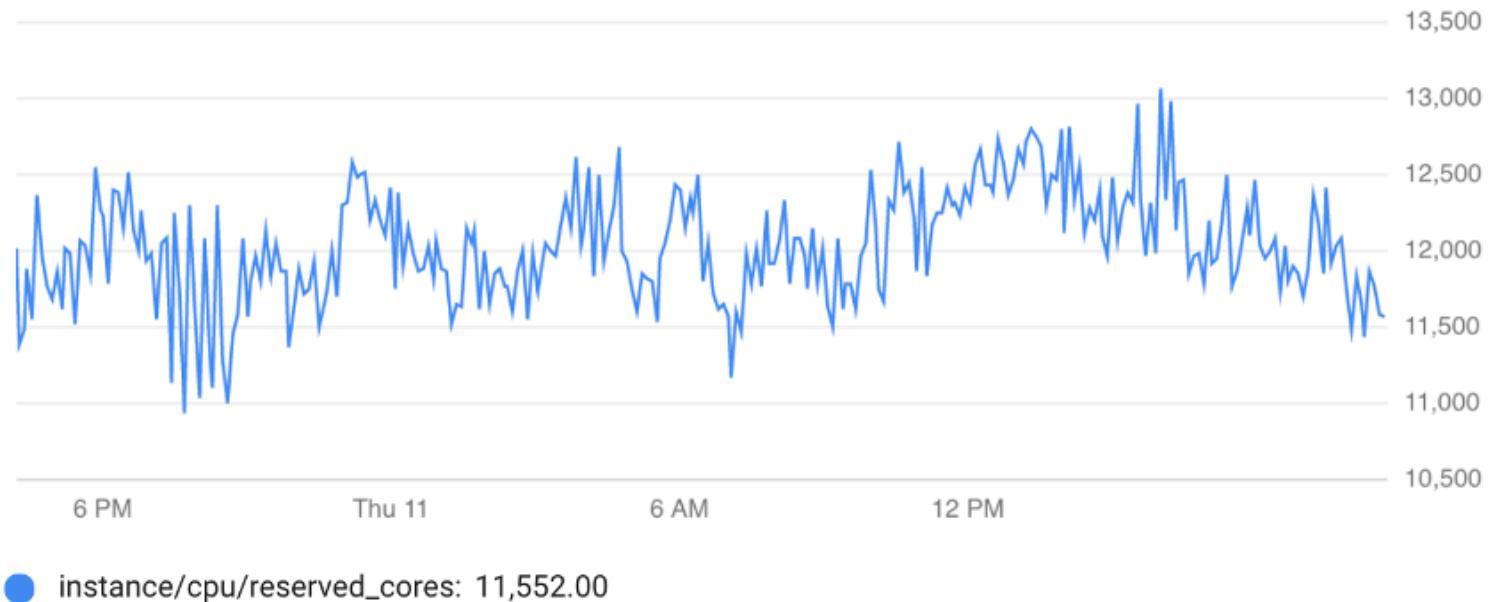
Samantha Zarate

Preview

1 hour

4 hours

1 day



AnVIL: RStudio in the Cloud

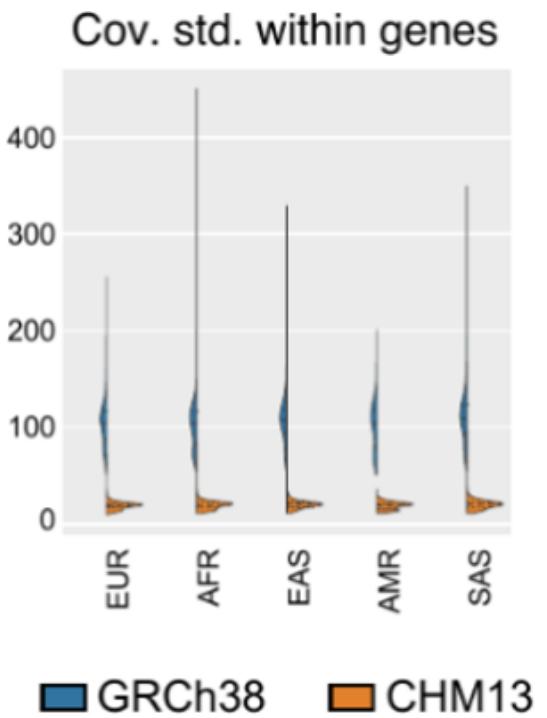
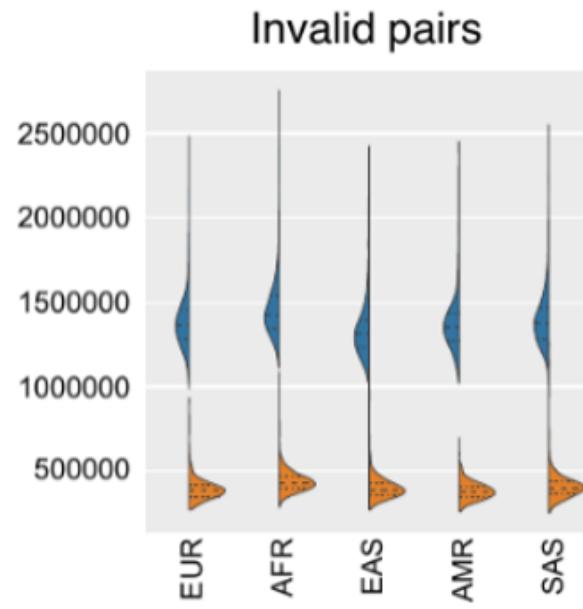
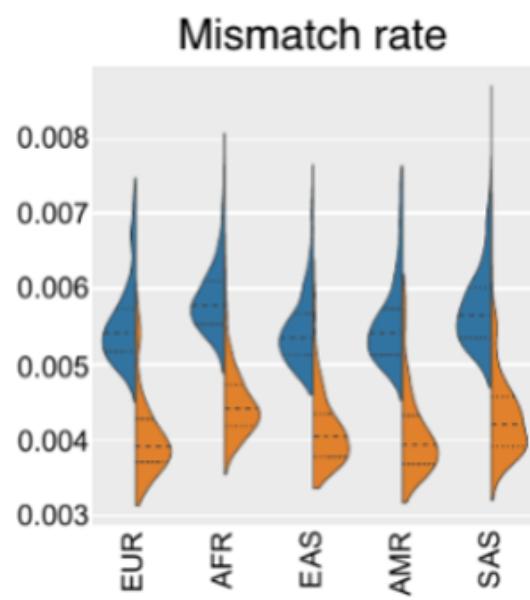
The screenshot shows the AnVIL RStudio interface running in a web browser. The top navigation bar includes links for Home, Dashboards, Workspaces, and Applications. The main window has tabs for Dashboard, Data, Notebooks, Workflows, and Job History. The current tab is 'Terminal'.

The terminal window displays R code being run. The code involves filtering data by 'superpop' (which is set to 'pop') and creating ggplot objects for different population groups ('AFR', 'AMR', 'EAS', 'EUR', 'SAS', 'SAM'). These plots show the distribution of 'percentage_of_properly_paired_reads' for 'female' (pink) and 'male' (teal) samples. The plots are faceted by 'superpop' and have legends indicating the population and sex.

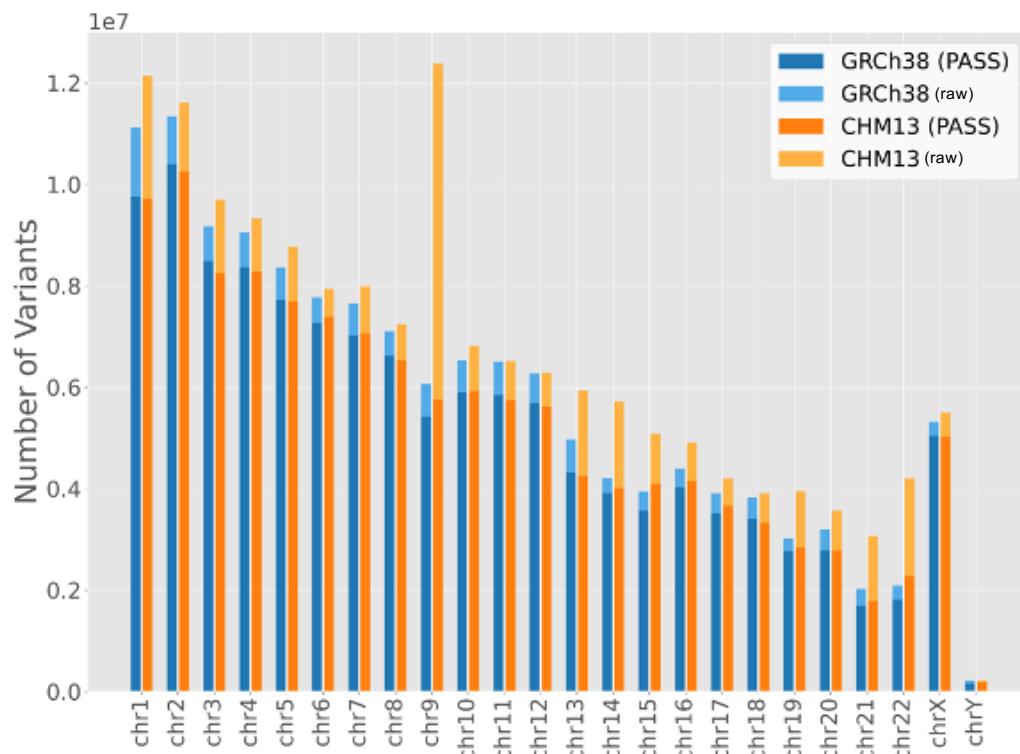
```
R
File Edit Code View Plots Session Build Debug Profile Tools Help
Source Terminal
> print(pop)
[1] "pop"
> superpop = filter(xstats, Superpopulation_code==pop)
> ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=Superpopulation_code))
#facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0,100)
+   geom_histogram(aes(x= percentage_of_properly_paired_reads_00), fill=pop) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code)
+   print("done")
}
ggplot(xstats, aes(x=percentage_of_properly_paired_reads_00 , fill=Superpopulation_code)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code)
{
  print("done")
}
print(pop)
superpop = filter(xstats, Superpopulation_code==pop)
#ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=Superpopulation_code)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=pop)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
readline("Press [enter] to continue")
}
[1] "AFR"; superpop = filter(xstats, Superpopulation_code==pop)
#ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=Superpopulation_code)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=pop)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
readline("Press [enter] to continue")
}
[1] "AMR"; superpop = filter(xstats, Superpopulation_code==pop)
#ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=Superpopulation_code)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=pop)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
readline("Press [enter] to continue")
}
[1] "EAS"; superpop = filter(xstats, Superpopulation_code==pop)
#ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=Superpopulation_code)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=pop)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
readline("Press [enter] to continue")
}
[1] "EUR"; superpop = filter(xstats, Superpopulation_code==pop)
#ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=Superpopulation_code)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=pop)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
readline("Press [enter] to continue")
}
[1] "SAS"; superpop = filter(xstats, Superpopulation_code==pop)
#ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=Superpopulation_code)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=pop)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
readline("Press [enter] to continue")
}
[1] "SAM"; superpop = filter(xstats, Superpopulation_code==pop)
#ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=Superpopulation_code)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
ggplot(superpop, aes(x=percentage_of_properly_paired_reads_00 , fill=pop)) + geom_density(alpha=0.4) + themeLegend(position="top") + facet_grid(Superpopulation_code) + ggtitle(paste("Superpopulation: ", pop)) + xlim(0, 100)
readline("Press [enter] to continue")
}
[1] "AFR"; Press [enter] to continue
[1] "AMR"; Press [enter] to continue
[1] "EAS"; Press [enter] to continue
[1] "EUR"; Press [enter] to continue
[1] "SAS"; Press [enter] to continue
[1] "SAM"; Press [enter] to continue
[1] "SAM"; Press [enter] to continue
Removed 1 rows containing non-finite values (stat_density).
> |
```

The viewer area displays four density plots for the 'AFR', 'AMR', 'EAS', and 'EUR' populations. Each plot shows two overlapping bell-shaped curves: a pink one for females and a teal one for males. The x-axis is labeled 'percentage_of_properly_paired_reads_(%)' and ranges from 90.0 to 100.0. The y-axis is labeled 'density'. A legend at the top right indicates 'Sex' with 'female' in pink and 'male' in teal.

I000G Mapping on T2T-CHM13

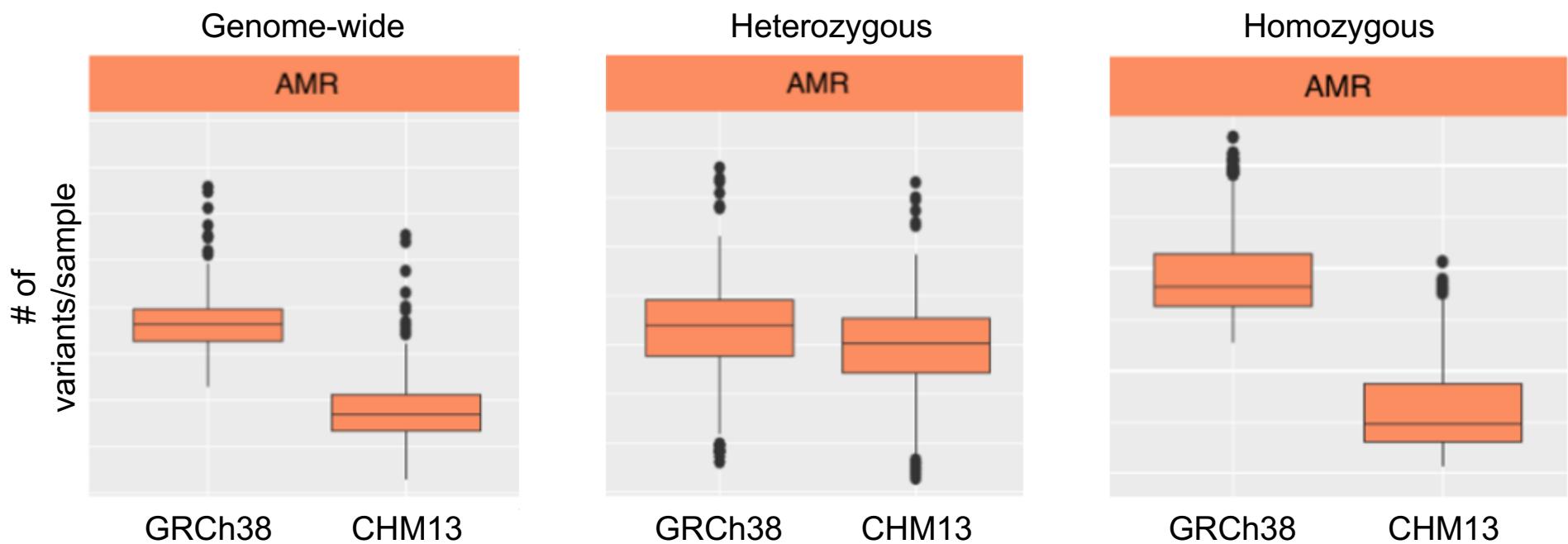


Across 1000G, Many More Variants Found Using CHM13

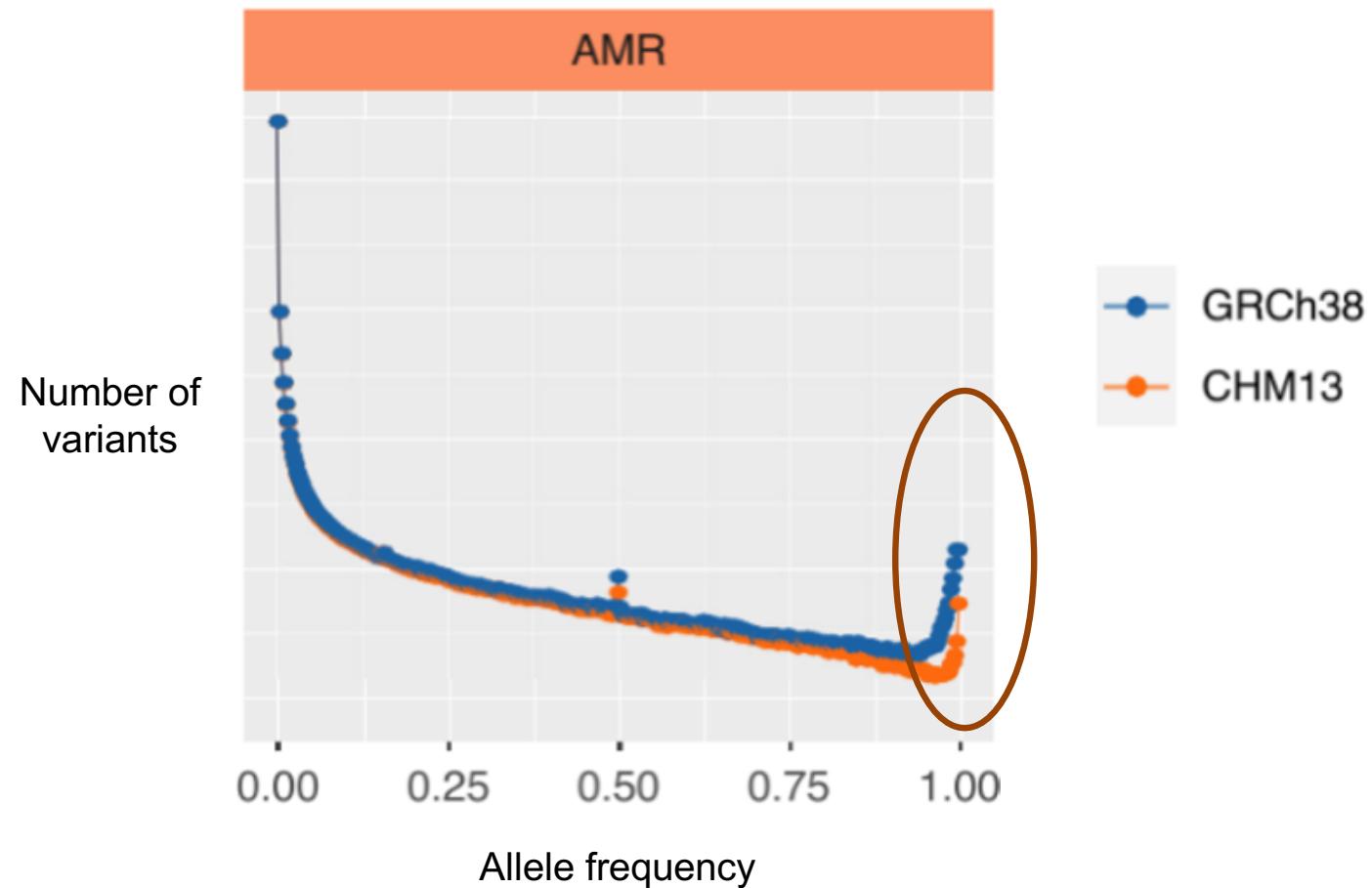


	GRCh38	CHM13
# PASS variants	125,484,020	126,591,489

1000G Per-Sample Variant Counts on T2T-CHM13



Explaining Decreased Per-Sample Count

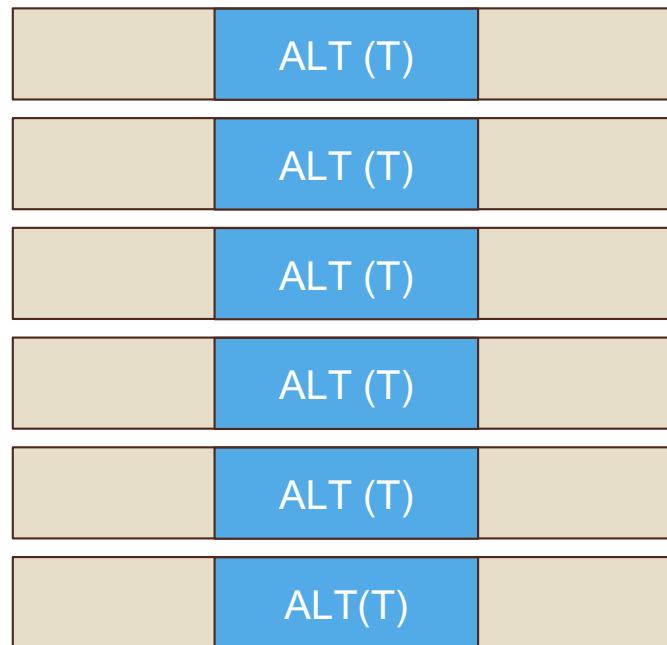


Allele Frequency = I: Reference error / private variant

GRCh38



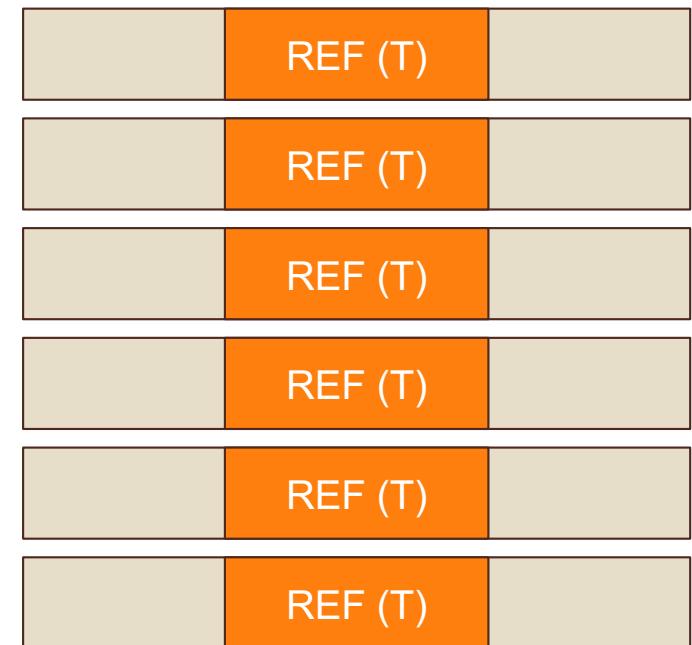
Samples



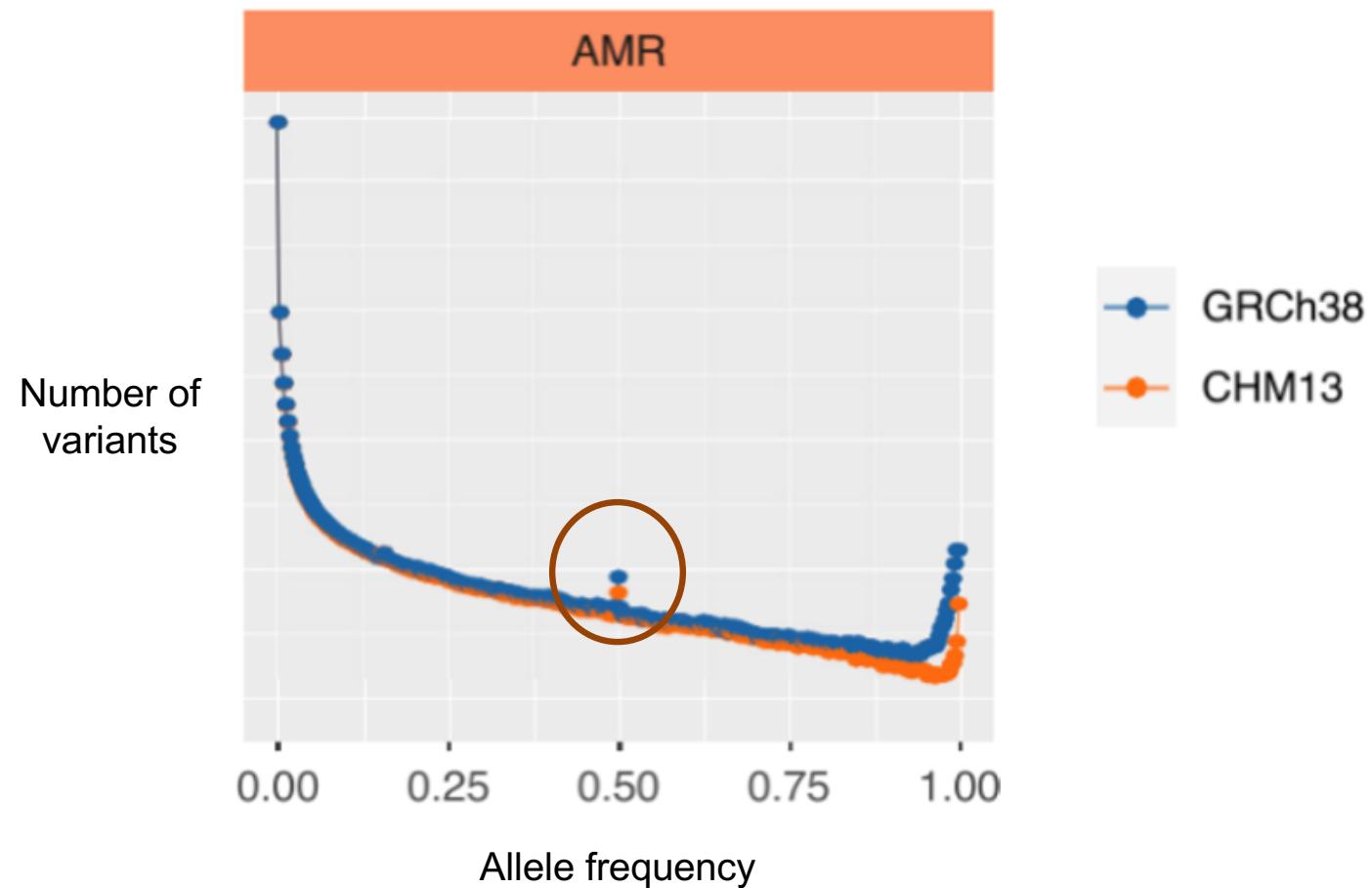
CHM13



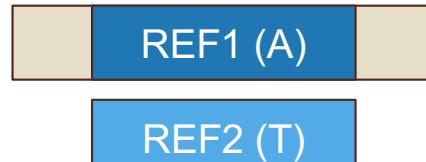
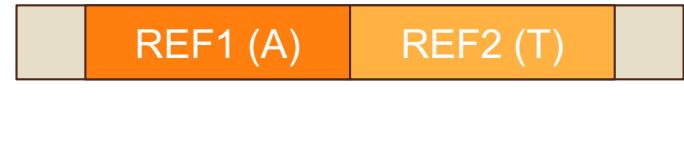
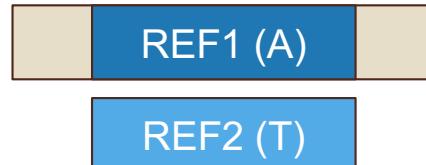
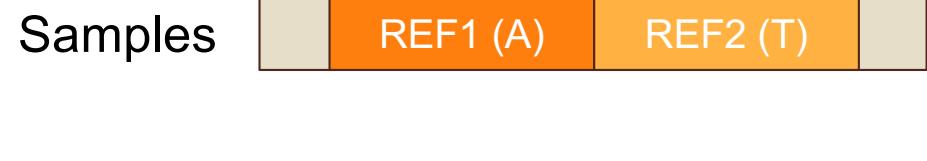
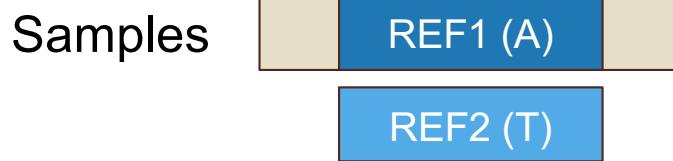
Samples



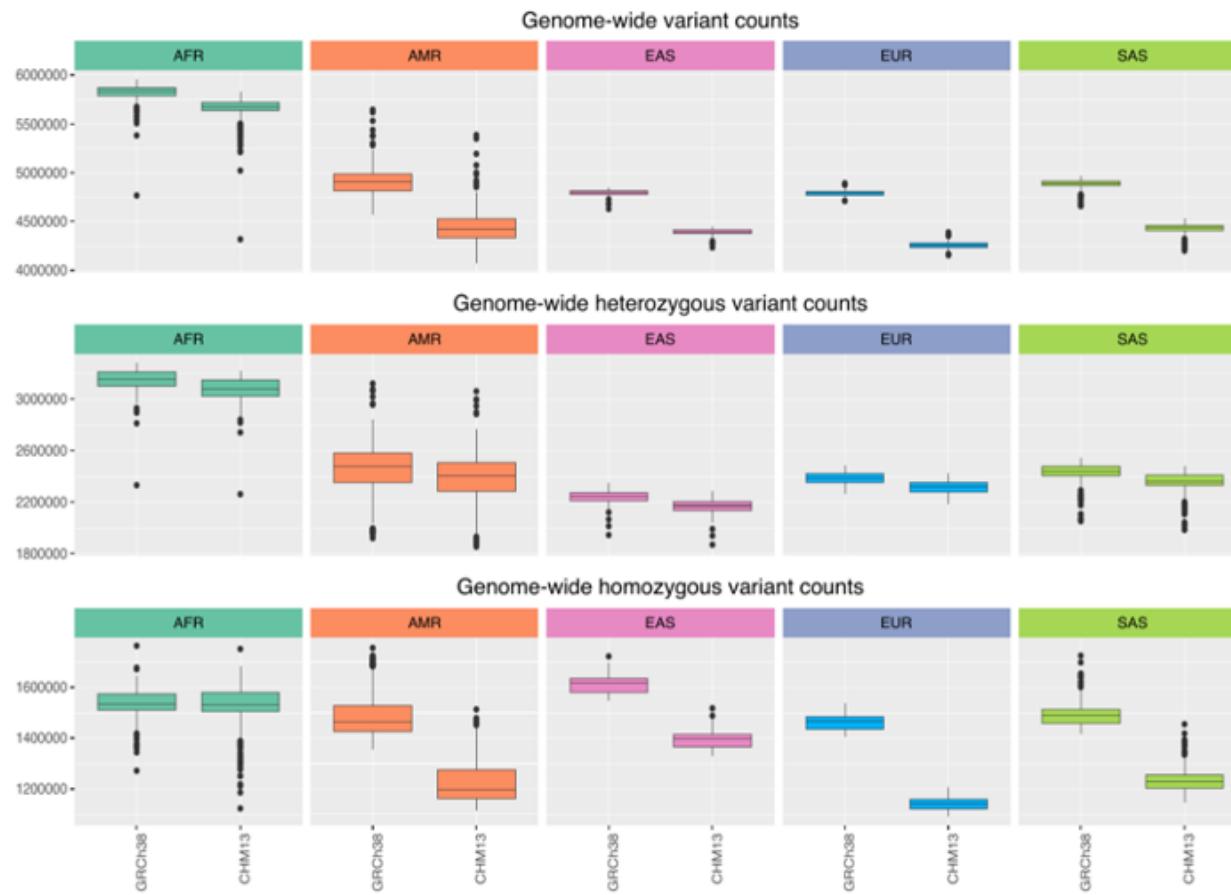
Explaining Decreased Per-Sample Count



Allele Frequency \approx 0.5: Collapsed duplication



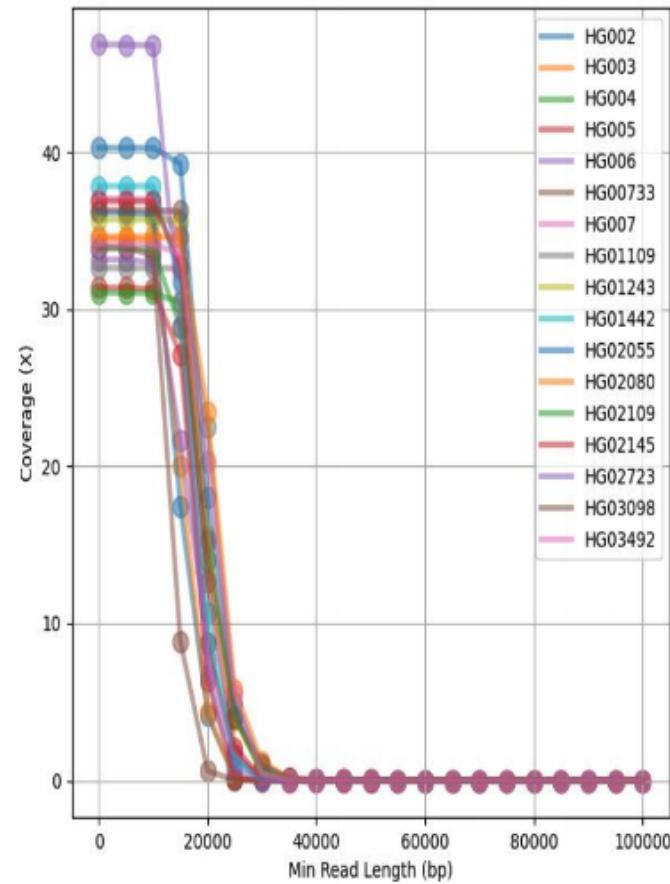
1000G Per-Sample Variant Counts on T2T-CHM13



Long Read Analysis

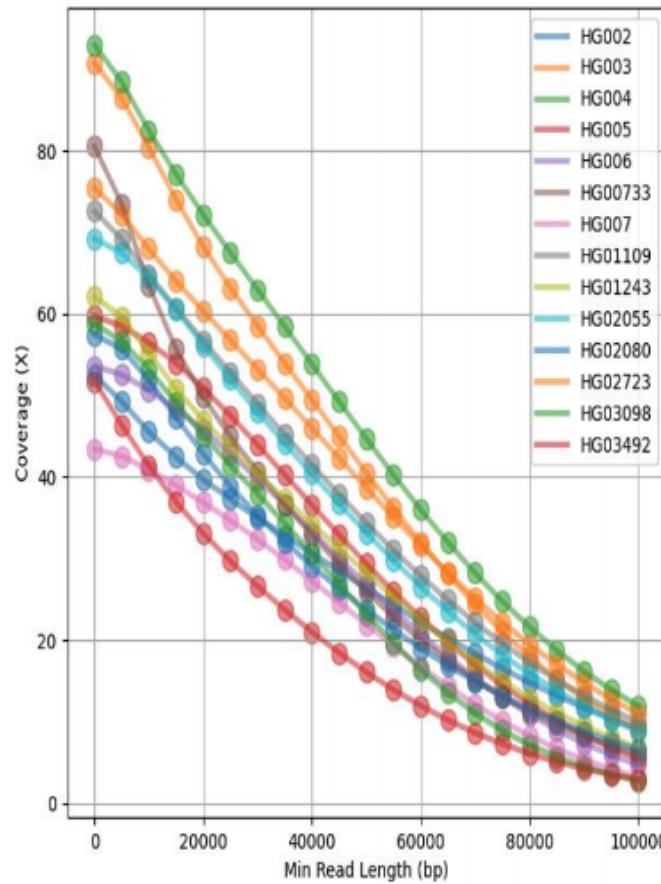
B.

HiFi



C.

ONT



Melanie Kirsche



Sergey Aganezov

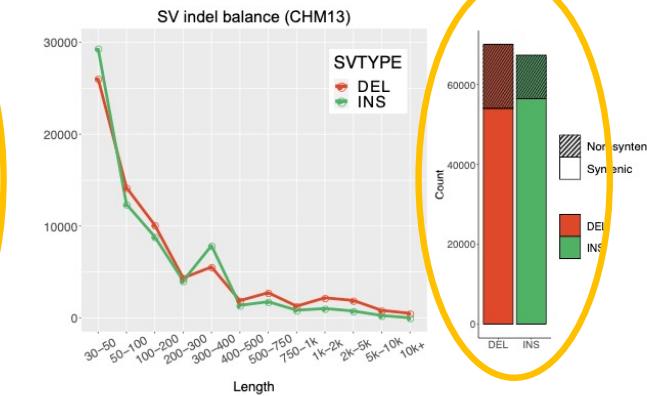
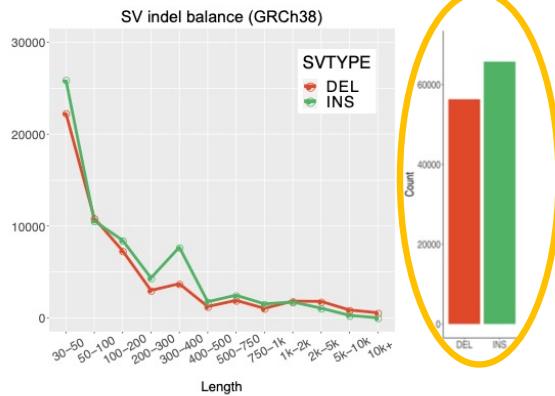
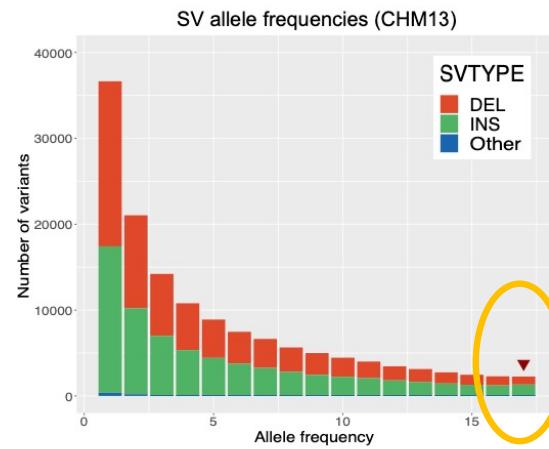
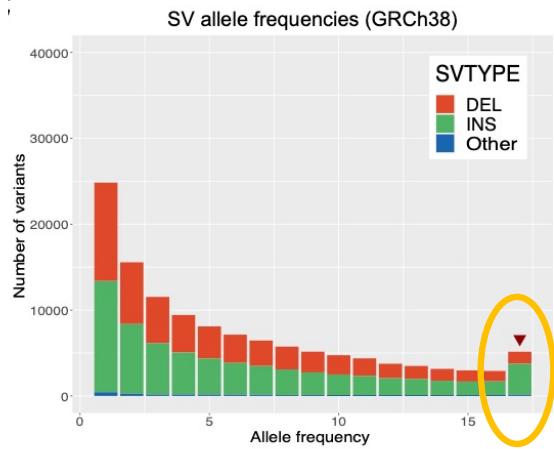
Long-Read Analysis with T2T-CHM13



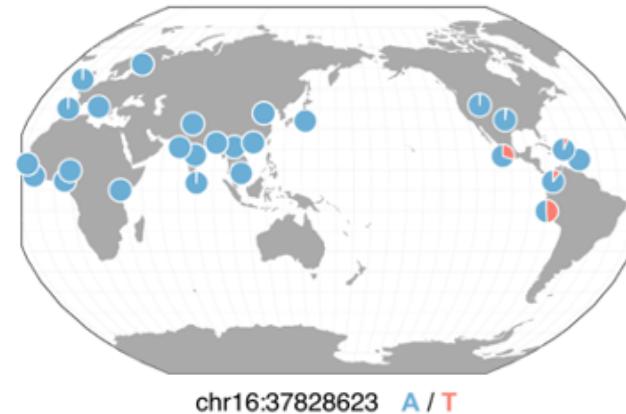
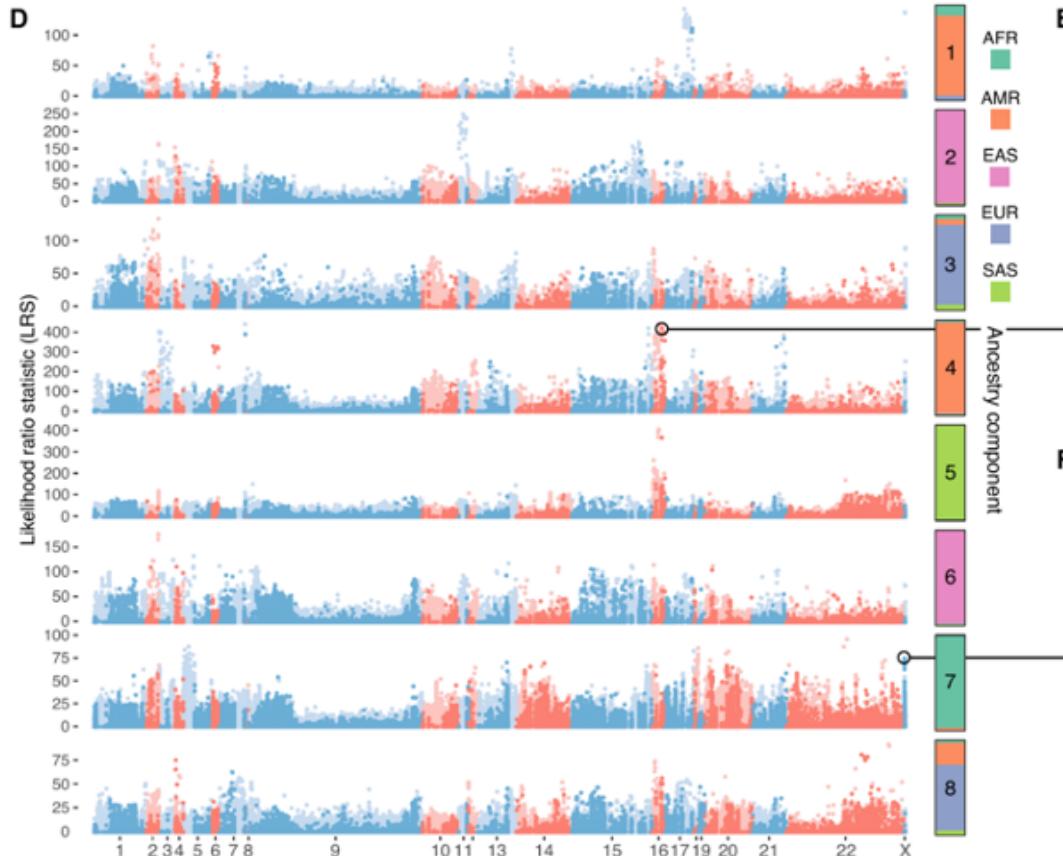
Melanie
Kirsche



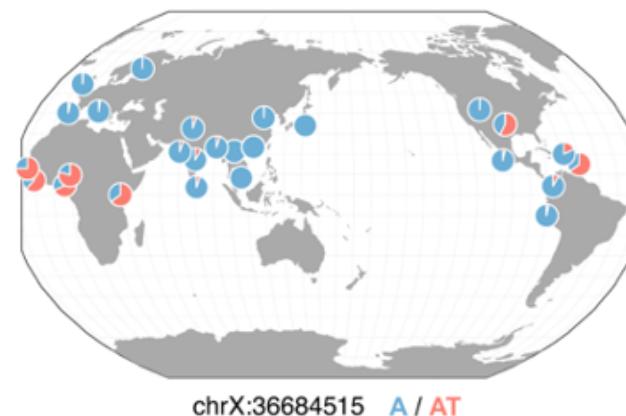
Sergey
Aganezov



Variants in Newly Resolved Regions

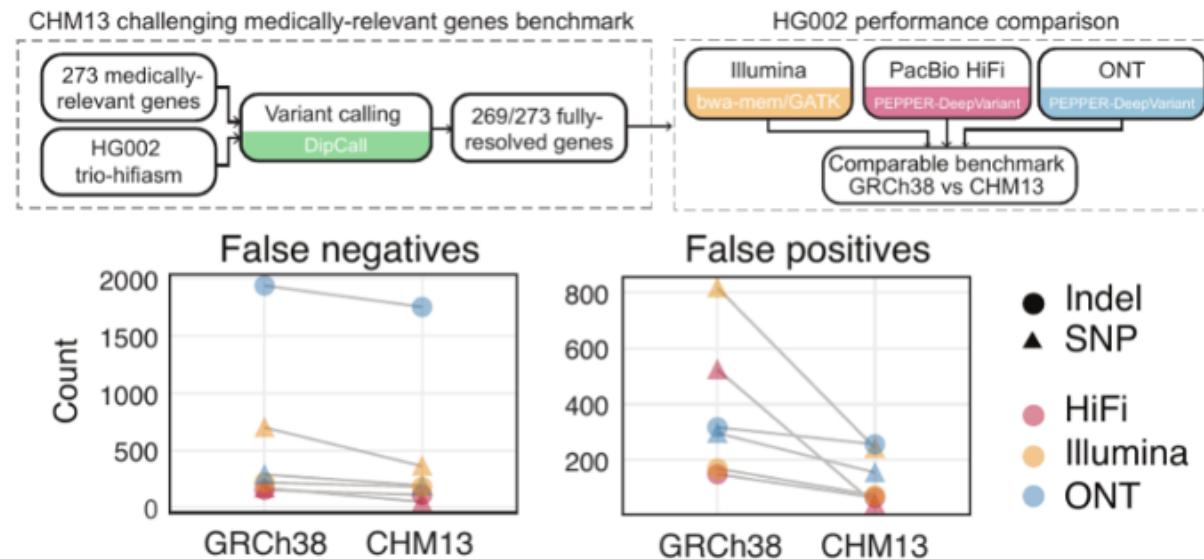


Stephanie Yan



Rajiv McCoy

T2T-CHM13 Improves Clinical Genomics Variant Calling



Danny Miller



Daniela Soto



Megan Dennis



Justin Zook

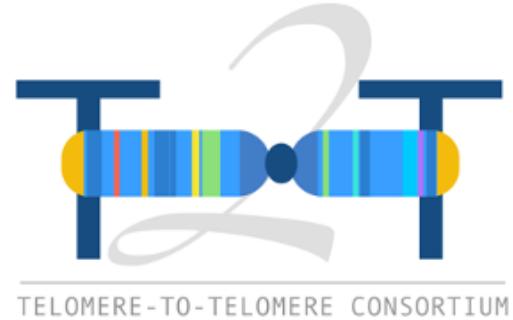
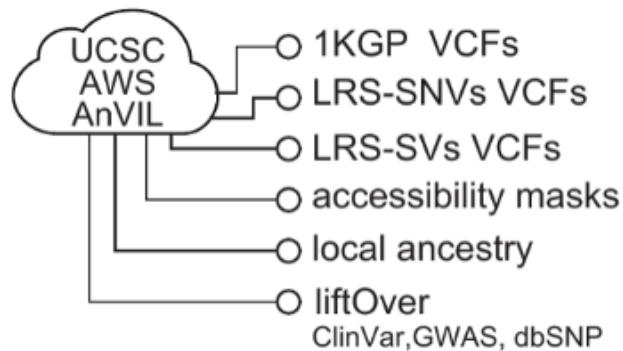


Fritz Sedlazeck

For more information, see: (Wagner, J et al., bioRxiv, 2021)

Summary of Results

- Large-scale cloud analysis of short-read data
- Short-read alignment and variant calling
- Long-read alignment and variant calling
- Variants in previously unresolved regions
- Clinical implications of variants





THE PREPRINT SERVER FOR BIOLOGY

The complete sequence of a human genome

Sergey Nurk, Sergey Koren, Arang Rhee, Mikko Rautiainen, Andrey V. Bzikadze, Alla Mikheenko, Mitchell R. Vollger, Nicolas Altemose, Lev Uralsky, Ariel Gershman, Sergey Aganezov, Savannah J. Hoyt, Mark Diekhans, Glennis A. Logsdon, Michael Alonge, Stylianos E. Antonarakis, Matthew Borchers, Gerard G. Bouffard, Shelise Y. Brooks, Gina V. Caldas, Haoyu Cheng, Chen-Shan Chin, William Chow, Leonardo G. de Lima, Philip C. Dishuck, Richard Durbin, Tatiana Dvorkina, Ian T. Fiddes, Giulio Formenti, Robert S. Fulton, Arkarachai Fungtammasan, Erik Garrison, Patrick G.S. Grady, Tina A. Graves-Lindsay, Ira M. Hall, Nancy F. Hansen, Gabrielle A. Hartley, Marina Haukness, Kerstin Howe, Michael W. Hunkapiller, Chirag Jain, Miten Jain, Erich D. Jarvis, Peter Kerdjiev, Melanie Kirsche, Mikhail Kolmogorov, Jonas Korlach, Milinn Kremitzki, Heng Li, Valerie V. Maduro, Tobias Marschall, Ann M. McCartney, Jennifer McDaniel, Danny E. Miller, James C. Mullikin, Eugene W. Myers, Nathan D. Olson, Benedict Paten, Paul Peluso, Pavel A. Pevzner, David Porubsky, Tamara Potapova, Evgeny I. Rozaev, Jeffrey A. Rosenfeld, Steven L. Salzberg, Valerie A. Schneider, Fritz J. Sedlazeck, Kishwar Shafin, Colin J. Shew, Alaina Shumate, Yumi Sims, Arian F.A. Smit, Daniela C. Soto, Ivan Sović, Jessica M. Storer, Aaron Streets, Beth A. Sullivan, Françoise Thibaud-Nissen, James Torrance, Justin Wagner, Brian P. Walenz, Aaron Wenger, Jonathan M. D. Wood, Chunlin Xiao, Stephanie M. Yan, Alice C. Young, Samantha Zarate, Urvashi Surti, Rajiv C. McCoy, Megan Y. Dennis, Ivan A. Alexandrov, Jennifer L. Gerton, Rachel J. O'Neill, Winston Timp, Justin M. Zook, Michael C. Schatz, Evan E. Eichler, Karen H. Miga, Adam M. Phillippy

doi: <https://doi.org/10.1101/2021.05.26.445798>

A complete reference genome improves analysis of human genetic variation

Sergey Aganezov, Stephanie M. Yan, Daniela C. Soto, Melanie Kirsche, Samantha Zarate, Pavel Avdeyev, Dylan J. Taylor, Kishwar Shafin, Alaina Shumate, Chunlin Xiao, Justin Wagner, Jennifer McDaniel, Nathan D. Olson, Michael E.G. Sauria, Mitchell R. Vollger, Arang Rhee, Melissa Meredith, Skylar Martin, Joyce Lee, Sergey Koren, Jeffrey A. Rosenfeld, Benedict Paten, Ryan Layer, Chen-Shan Chin, Fritz J. Sedlazeck, Nancy F. Hansen, Danny E. Miller, Adam M. Phillippy, Karen H. Miga, Rajiv C. McCoy, Megan Y. Dennis, Justin M. Zook, Michael C. Schatz

doi: <https://doi.org/10.1101/2021.07.12.452063>

Complete genomic and epigenetic maps of human centromeres

Nicolas Altemose, Glennis A. Logsdon, Andrey V. Bzikadze, Pragy Sidhwani, Sasha A. Langley, Gina V. Caldas, Savannah J. Hoyt, Lev Uralsky, Fedor D. Ryabov, Colin J. Shew, Michael E.G. Sauria, Matthew Borchers, Ariel Gershman, Alla Mikheenko, Valery A. Shepelev, Tatiana Dvorkina, Olga Kunyavskaya, Mitchell R. Vollger, Arang Rhee, Ann M. McCartney, Mobin Asri, Ryan Lorig-Roach, Kishwar Shafin, Sergey Aganezov, Daniel Olson, Leonardo Gomes de Lima, Tamara Potapova, Gabrielle A. Hartley, Marina Haukness, Peter Kerdjiev, Fedor Gusev, Kristof Tigray, Shelise Brooks, Alice Young, Sergey Nurk, Sergey Koren, Sofie R. Salama, Benedict Paten, Evgeny I. Rozaev, Aaron Streets, Gary H. Karpen, Abby F. Dernburg, Beth A. Sullivan, Aaron F. Straight, Travis J. Wheeler, Jennifer L. Gerton, Evan E. Eichler, Adam M. Phillippy, Winston Timp, Megan Y. Dennis, Rachel J. O'Neill, Justin M. Zook, Michael C. Schatz, Pavel A. Pevzner, Mark Diekhans, Charles H. Langley, Ivan A. Alexandrov, Karen H. Miga

doi: <https://doi.org/10.1101/2021.07.12.452052>

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

Savannah J. Hoyt, Jessica M. Storer, Gabrielle A. Hartley, Patrick G. S. Grady, Ariel Gershman, Leonardo G. de Lima, Charles Limouse, Reza Halabian, Luke Wojenski, Matias Rodriguez, Nicolas Altemose, Leighton J. Core, Jennifer L. Gerton, Wojciech Makalowski, Daniel Olson, Jeb Rosen, Arian F.A. Smit, Aaron F. Straight, Mitchell R. Vollger, Travis J. Wheeler, Michael C. Schatz, Evan E. Eichler, Adam M. Phillippy, Winston Timp, Karen H. Miga, Rachel J. O'Neill

doi: <https://doi.org/10.1101/2021.07.12.451456>

Epigenetic Patterns in a Complete Human Genome

Ariel Gershman, Michael E.G. Sauria, Paul W. Hook, Savannah J. Hoyt, Roham Razaghi, Sergey Koren, Nicolas Altemose, Gina V. Caldas, Mitchell R. Vollger, Glennis A. Logsdon, Arang Rhee, Evan E. Eichler, Michael C. Schatz, Rachel J. O'Neill, Adam M. Phillippy, Karen H. Miga, Winston Timp

doi: <https://doi.org/10.1101/2021.05.26.443420>

Segmental duplications and their variation in a complete human genome

Mitchell R. Vollger, Xavi Guitart, Philip C. Dishuck, Ludovica Mercuri, William T. Harvey, Ariel Gershman, Mark Diekhans, Arvis Sulovari, Katherine M. Munson, Alexandra M. Lewis, Kendra Hoekzema, David Porubsky, Ruiyang Li, Sergey Nurk, Sergey Koren, Karen H. Miga, Adam M. Phillippy, Winston Timp, Mario Ventura, Evan E. Eichler

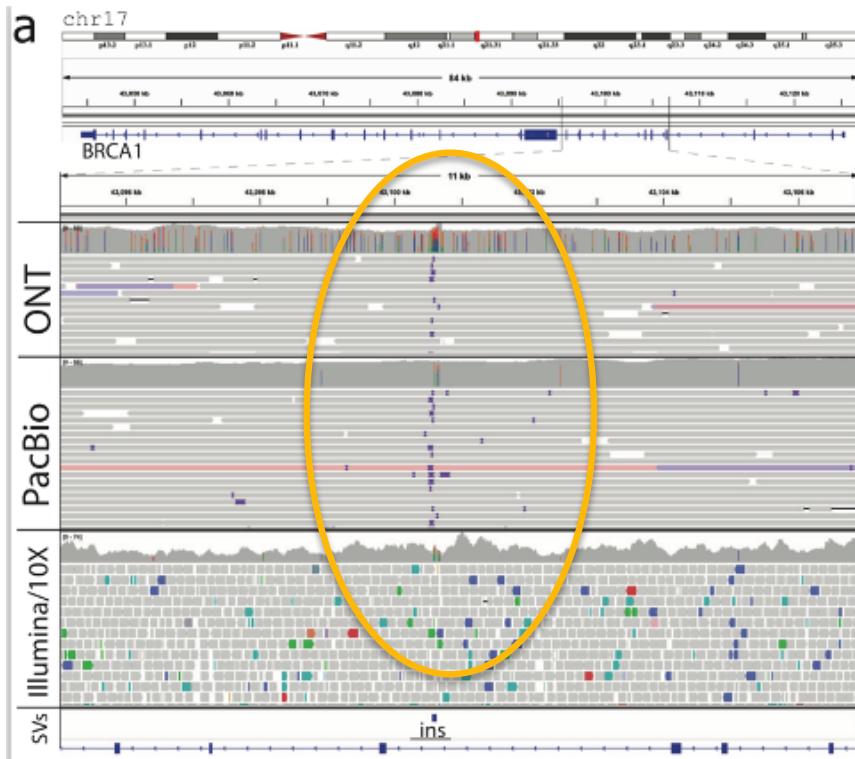
doi: <https://doi.org/10.1101/2021.05.26.445678>

Outline

- Assembly by analogy
- T2T-CHM13 Assembly
 - How the T2T-CHM13 reference improves the analysis of genetic variation
- Beyond T2T-CHM13



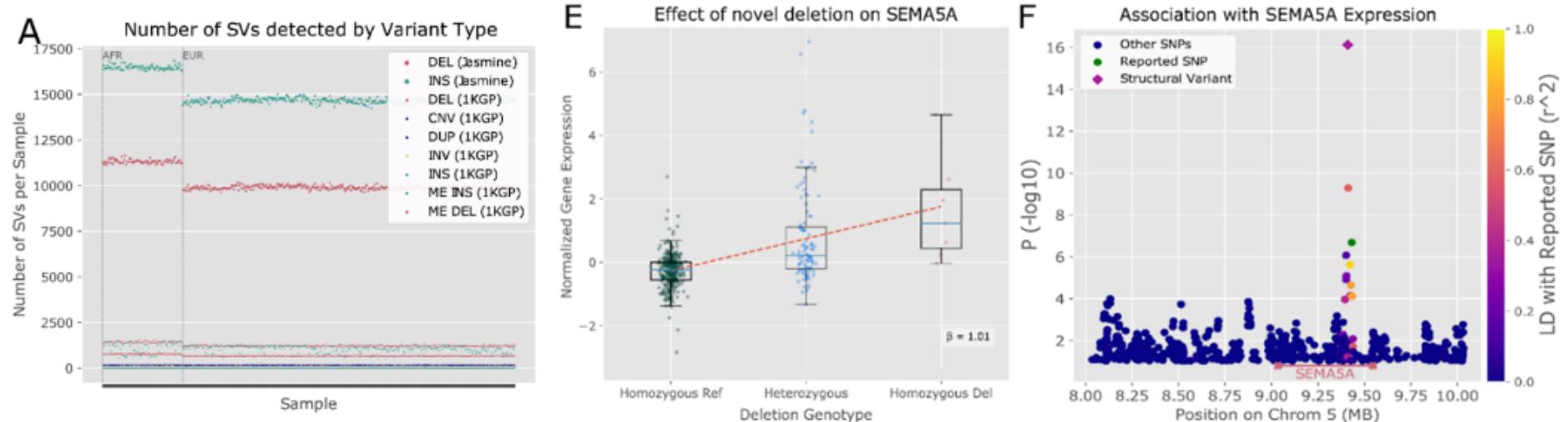
Hidden Variants in Breast Cancer Genes



62bp repeat expansion in BRCA1 detected in normal tissue that is undetectable using a cancer panel or short read sequencing

Comprehensive analysis of structural variants in breast cancer genomes using single molecule sequencing
Aganezov et al. (2020) Genome Research. doi:10.1101/gr.260497.119

Novel association of 37bp deletion with SEMA5A expression



- Genotyped ~200k SVs from our long-read cohort in 444 samples from the 1000 Genomes Project using ParaGRAPH (Chen et al, 2019) which also had gene expression data available
- The deletion is in strong LD with the reported SNP-eQTL, but has a much higher effect size on SEMA5A

Kirsche, et al (2021) bioRxiv doi: <https://doi.org/10.1101/2021.05.27.445886>

T2T Genomes and Epigenomes Across the Tree of Life

Cell

Article

Major Impacts of Widespread Structural Variation on Gene Expression and Crop Improvement in Tomato

Graphical Abstract

The graphical abstract illustrates the workflow: 100 diverse tomato varieties are sequenced using long-read sequencing to identify over 200k structural variants (SVs). These SVs are then analyzed for their impact on gene expression, trait variation (flavor, fruit size, productivity), and breeding applications like genome editing with Cas9.

Authors

Michael Alonte, Xingang Wang,
Matthias Benoit, ..., Esther van der Knaap,
Michael C. Schatz, Zachary B. Lippman

Correspondence

m.schatz@cs.jhu.edu (M.C.S.),
lippman@cshl.edu (Z.B.L.)

In Brief

Comprehensive structural variant identification in tomato genomes allows insight into the evolution and domestication of tomato and serves as a resource for phenotype-directed breeding.

RESEARCH

RESEARCH ARTICLE SUMMARY

PLANT SCIENCE

The genetic and epigenetic landscape of the *Arabidopsis* centromeres

Matthew Naish,¹ Michael Alonte,¹ Piotr Włodzimierz,¹ Andrew J. Tock,¹ Bradley W. Abramson,¹ Anna Schmücker,¹ Terezie Mandáková,² Bhagyashree Jange,¹ Christophe Lambing,¹ Pallas Kuø,¹ Natasha Yelina,¹ Nolan Hartwick,¹ Kelly Coft,¹ Lisa M. Smith,¹ Jurriaan Ton,¹ Tetsuji Kakutani,¹ Robert A. Martienssen,¹ Korbinian Schneeberger,¹ Martin A. Lysak,¹ Frédéric Berger,¹ Alexandros Bousios,¹ Todd P. Michael,¹ Michael C. Schatz*,¹ Ian R. Henderson*

INTRODUCTION: The centromeres of eukaryotic chromosomes assemble the multiprotein kinetochore complex and thereby position attachment to the spindle microtubules, allowing chromosome segregation during cell division. The key function of the centromere is to load nucleosomes containing the CENTROMERE SPECIFIC HISTONE H3 (CENH3) histone variant [also known as centromere protein A (CENPA)], which directs kinetochore formation. Despite their conserved function during chromosome segregation, centromeres show radically diverse organization between species at the sequence level, ranging from single nucleosomes to megabase-scale satellite repeat arrays, which is termed the centromere paradox. Centromeric satellite repeats are variable in sequence composition and length when compared between species and show a high capacity for evolutionary change, both at the levels of primary sequence and array position along the chromosome. However, the genetic and epigenetic features that contribute to centromere function and evolution are incompletely understood, in part because of the challenges of centromere sequence assembly and functional genomics of highly repetitive sequences. New long-read DNA sequencing technologies can now resolve these complex repeat arrays, revealing insights into centromere architecture and chromatin organization.

RATIONALE: *Arabidopsis thaliana* is a model plant species; its genome was first sequenced in 2000, yet the centromeres, telomeres, and ribosomal DNA repeats have remained unassembled, owing to their high repetition and similarity. Genomic repeats are difficult to assemble from fragmented sequencing reads, with longer, high-identity repeats being the most challenging to correctly assemble. As sequencing reads have become longer and more accurate, eukaryotic de novo genome assemblies have captured an increasingly complete picture of the repetitive component of the genome, including the centromeres. For example, Oxford Nanopore Technologies (ONT) reads have

become longer and more accurate, now reaching >100 kilo-base pairs (kbp) in length with 95 to 99% modal accuracy. PacBio high-fidelity (HF) reads, although shorter (<15 kbp), are ~99% accurate. Using ONT and HF reads, it is possible to bridge across interspersed unique marker sequences and accurately assemble centromere sequences. In this study, we used long-read DNA sequencing to generate a genome assembly of the *A. thaliana* accession Columbia (Col-0) that resolves all five centromeres. We use the Col-CEN assembly to derive insights into the chromatin and recombination landscapes within the *Arabidopsis* centromeres and how these regions evolve.

RESULTS: The Col-CEN assembly reveals that the *Arabidopsis* centromeres consist of mega-base-scale tandemly repeated satellite arrays,

The figure shows the structure of Arabidopsis centromere 1. Top: Fluorescence in situ hybridization (FISH) showing CENH3 binding. Middle: ChIP-seq signal for CENH3 across the centromere. Bottom: Heatmap of sequence identity across nonoverlapping 5-kbp windows. The heatmap shows a gradient of color from red (0%) to green (100%), indicating sequence identity across the centromere.

Assembly of the *Arabidopsis* centromeres. The structure of *Arabidopsis* centromere 1 is shown by fluorescence in situ hybridization (top) [upper-arm bacterial artificial chromosomes (BACs) (green), *ATHILA* (purple), CENH3 (blue), the telomeric repeat (green), and bottom-arm BACs (yellow)] and a long-read genome assembly (middle). The density of centromeric histone CENH3 binding measured by ChIP-seq is shown (black), alongside the frequency of CENH3 centromere satellite repeats. Red and blue represent forward- and reverse-strand satellites, respectively. The heatmap (bottom) shows patterns of sequence identity across the centromere on nonoverlapping 5-kbp windows. Chr, chromosome 1.

Downloaded from https://www.science.org/10.1126/science.abi7489 by Johns Hopkins University on November 15, 2021

The list of author affiliations is available in the full article online.

*Corresponding author. Email: m.schatz@cs.jhu.edu

†These authors contributed equally to this work.

Cite this article as M. Naish et al., Science 374, eab7489 (2021). DOI: 10.1126/science.abi7489

S READ THE FULL ARTICLE AT
<https://doi.org/10.1126/science.abi7489>

1 of 1

Highlights

- Long-read sequencing of 100 tomato genomes uncovered 238,490 structural variants
- Transposons underlie many SVs, and SV hotspots revealed large introgressions
- SVs associated with genes are predictive of population-scale changes in expression
- New genome assemblies resolved complex breeding QTLs caused by SVs

Alonte et al., 2020, Cell 182, 145–161
July 9, 2020 © 2020 Elsevier Inc.
<https://doi.org/10.1016/j.cell.2020.05.021>

CellPress

Naish et al., Science 374, 840 (2021) 12 November 2021

Our genomics future



We are on the cusp of a biotechnology revolution: “tricorders for all”

- We all carry professional digital cameras in our pocket with quality unimaginable 20 years ago
- Endless applications to learn more about ourselves, improve our health, and study the world around us

Genomics 2041

- Interpreting genomics results requires very deep knowledge that is currently out of reach for many
- We must immediately take action to ensure genomics equality for all
- In 1000+ years, we will still be learning new meaning in the genome

Acknowledgements

Schatz Lab

Enis Afghan	Stephen Mosher
Mike Alonge	Bohan Ni
Arun Das	Alex Ostrovsky
Dannon Baker	Cokie Parker
David Clements	Gautam Prabhu
Sergey Golitsynskiy	Srividya
Sam Guerler	Ramakrishnan
Katie Jenike	Michael Sauria
Melanie Kirsche	Vikram Shivakumar
Sam Kovaka	Margaret Starostik
Natalie Kucher	Jenn Vessio
Alex Mahmoud	Samantha Zarate

AnVIL, Galaxy, & T2T Teams

Miga, Phillippy, Nekrutenko, Goecks, Tan, Leek, Morgan, Carey, Philippakis et al.

JHU

Battle Lab
Langmead Lab
Leek Lab
McCoy Lab
Salzberg Lab
Taylor Lab
Timp Lab

CSHL

McCombie Lab
Lippman Lab

Baylor

Sedlazeck Lab

U. Cambridge

Henderson Lab





Thank you!

@mike_schatz
<http://schatz-lab.org>