

## Lessons from the amylase locus

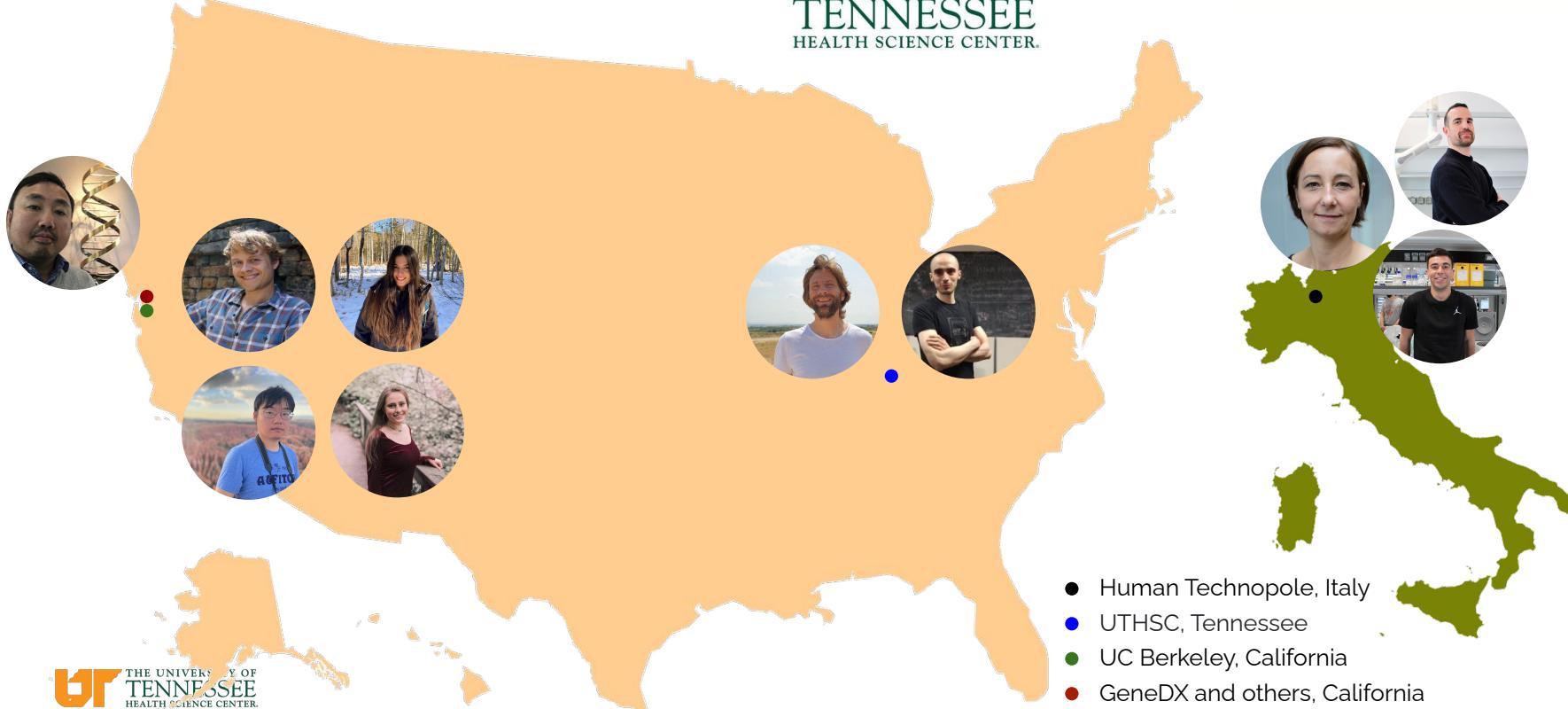
Andrea Guaracino  
Postdoctoral Scholar

 @AndresGuaracino

**HPRC Annual Meeting 2024**  
Seymour Marine Discovery Center, Santa Cruz, USA  
2024/09/05



# Panteam



# Outline

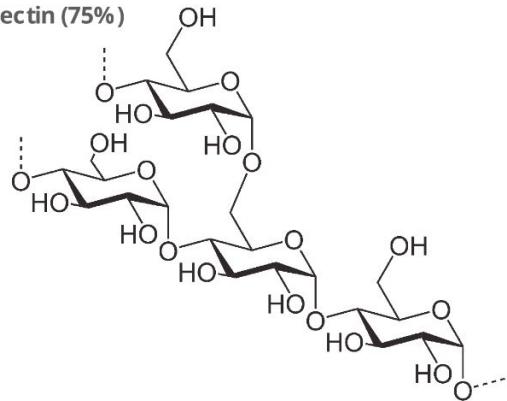
- Introduction
- Method
- Result



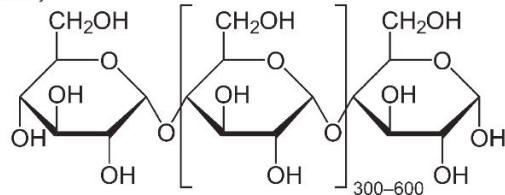
# Amylase breaks down starches to sugars

Starch

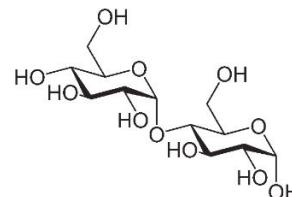
amylopectin (75%)



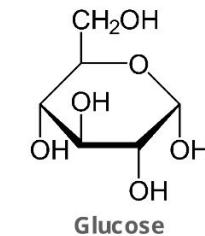
amylose (25%)



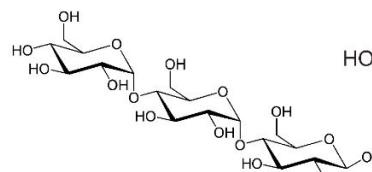
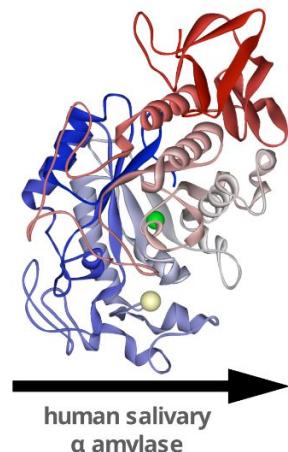
Sugar



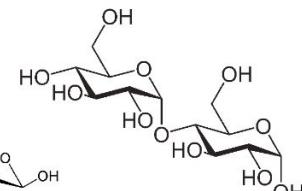
Maltose



Glucose

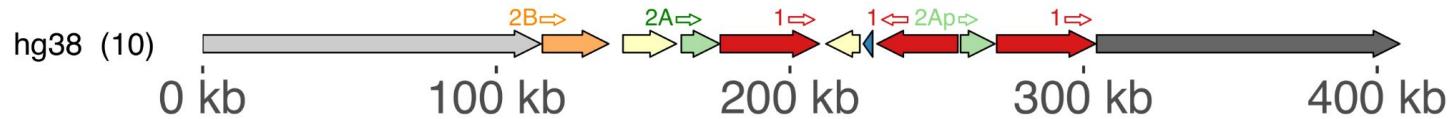


Maltotriose

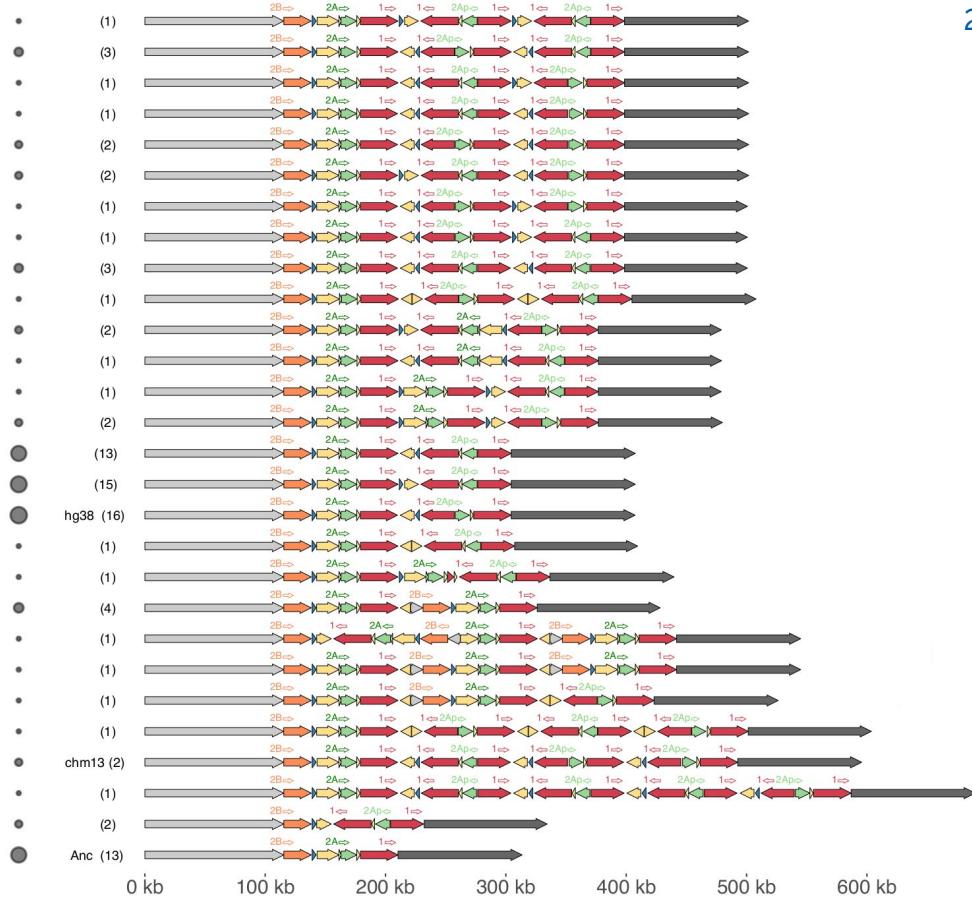


Maltose

## Amylase is a multi-copy gene family

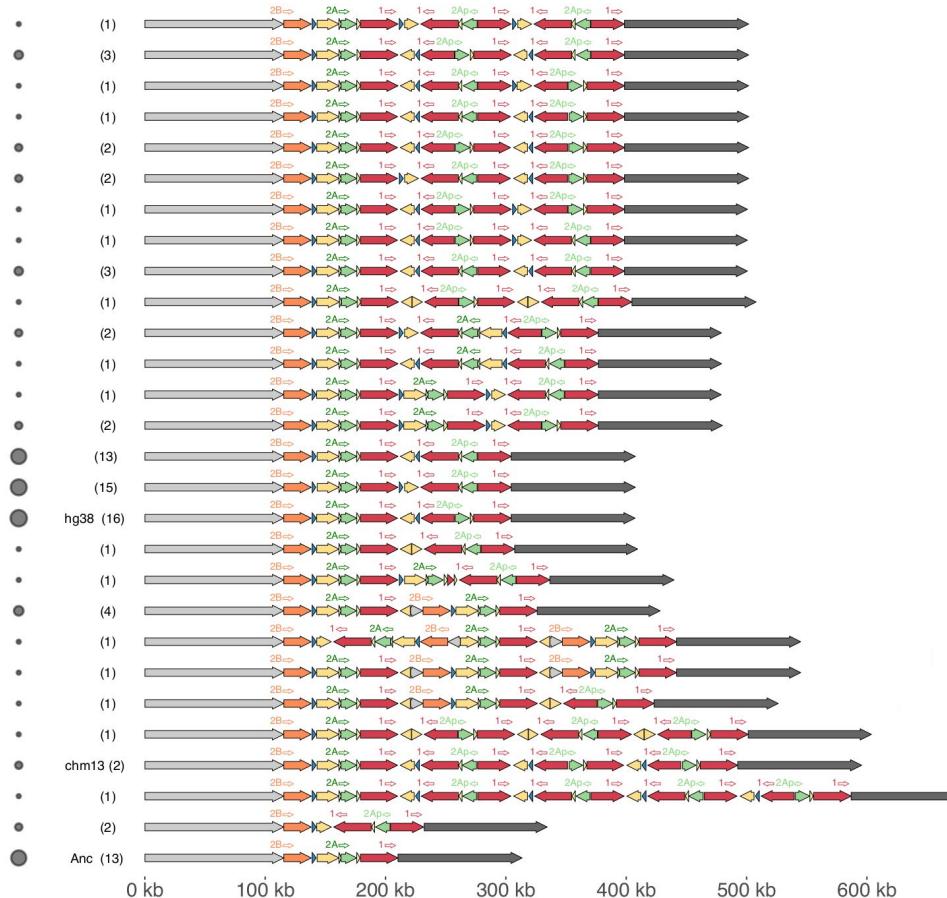


# Human amylase copy number diversity

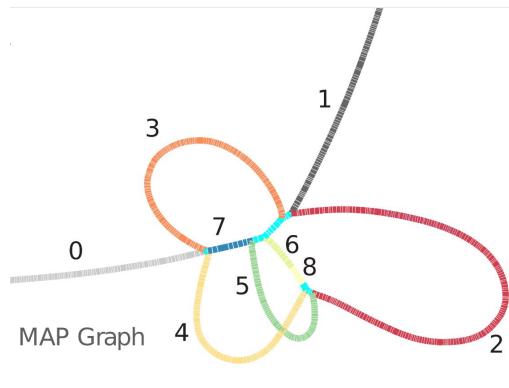


28 amylase structural haplotypes

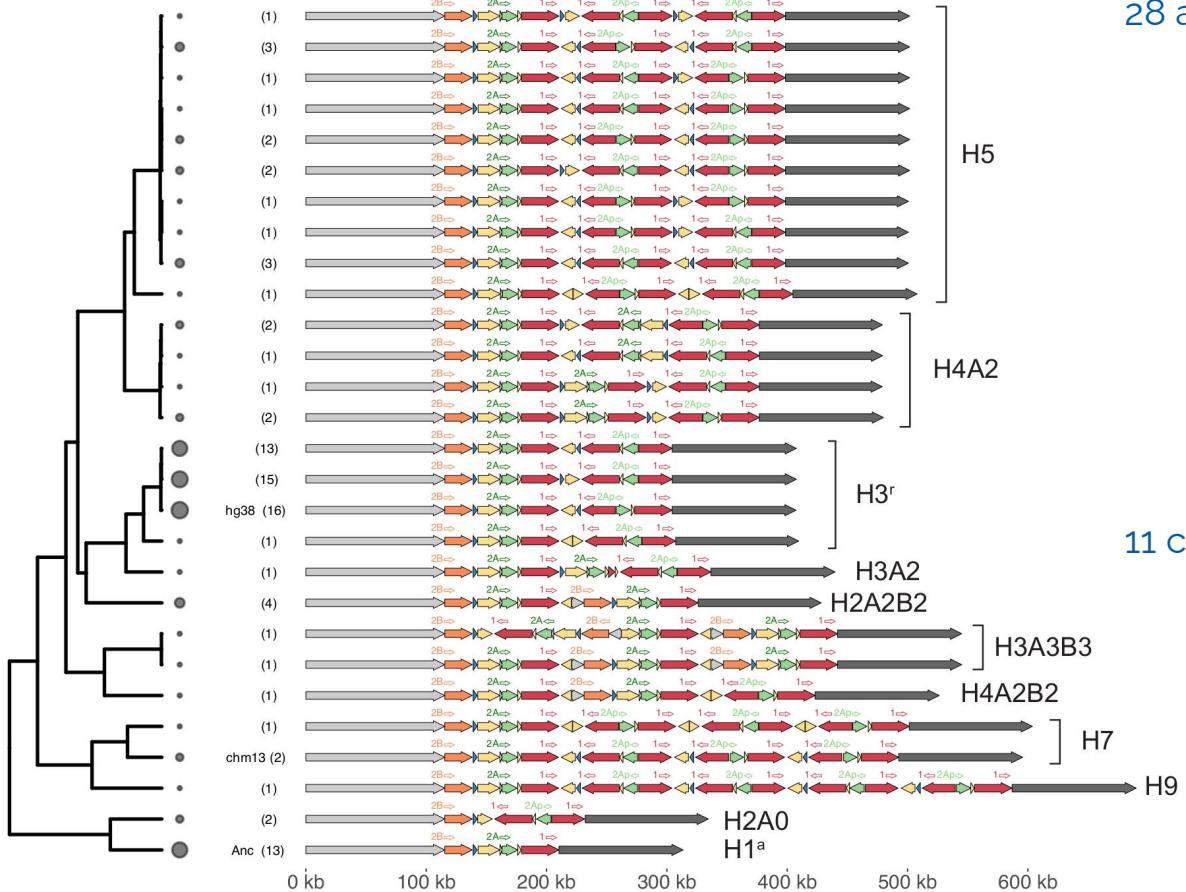
# Human amylase copy number diversity



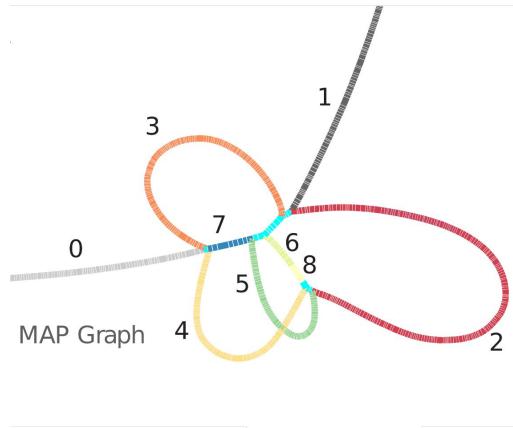
## 28 amylase structural haplotypes



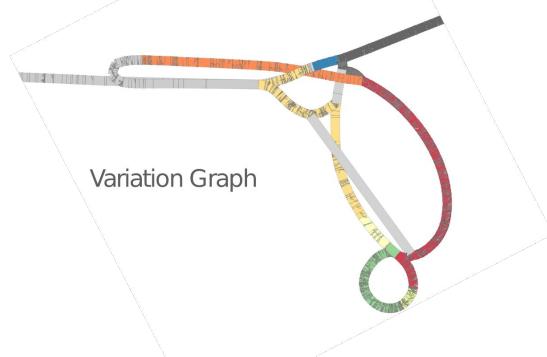
## Human amylase copy number diversity



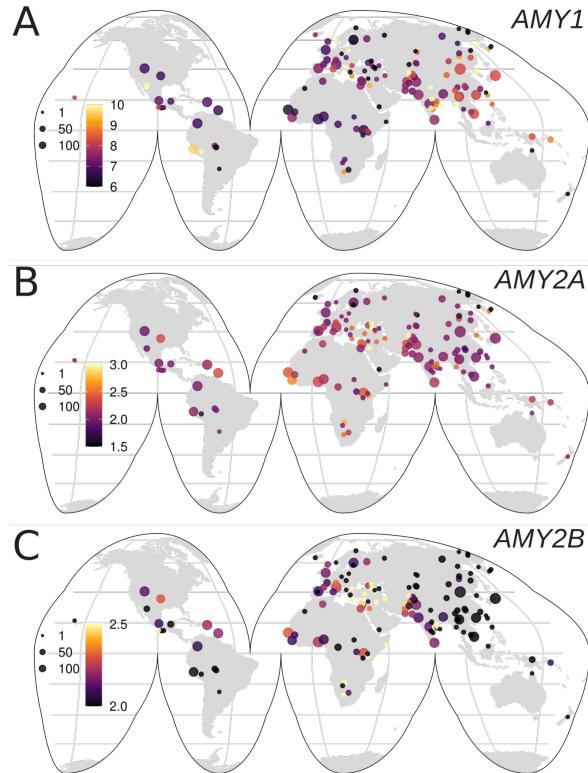
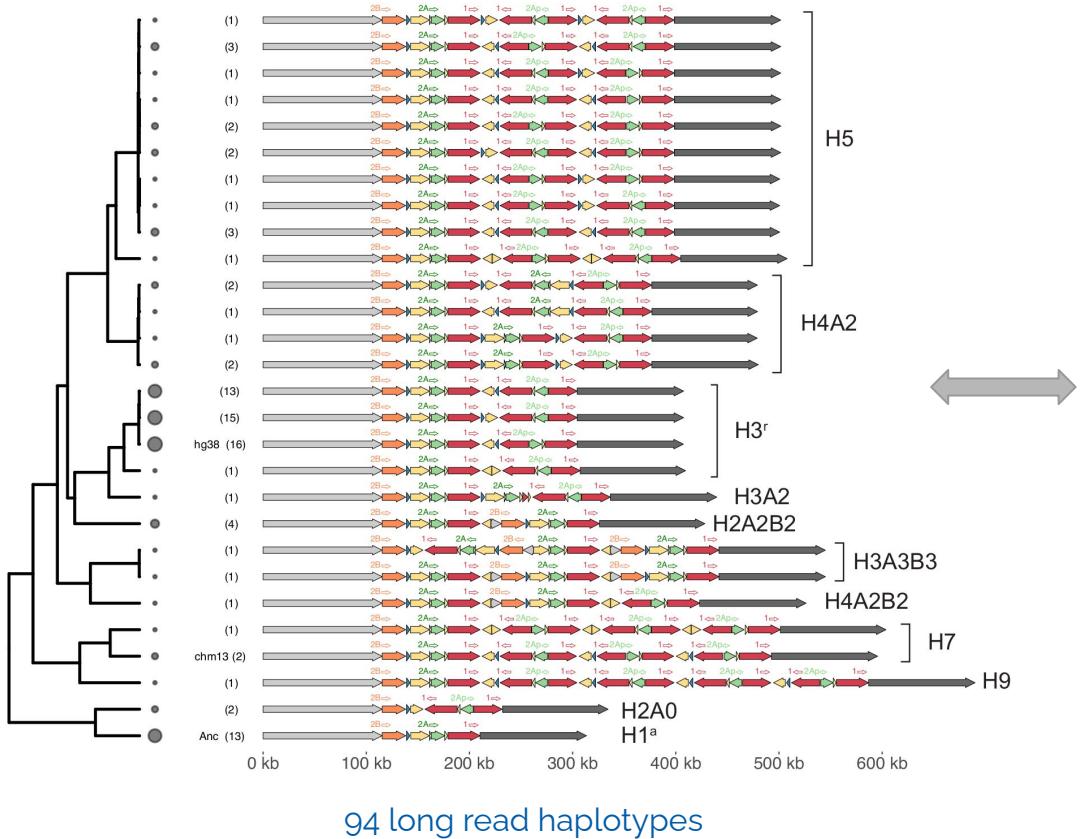
## 28 amylase structural haplotypes



## 11 consensus structures



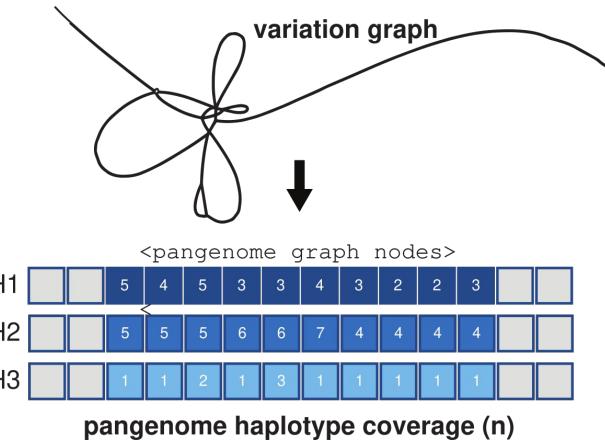
# We want more



~5600 short read haplotypes

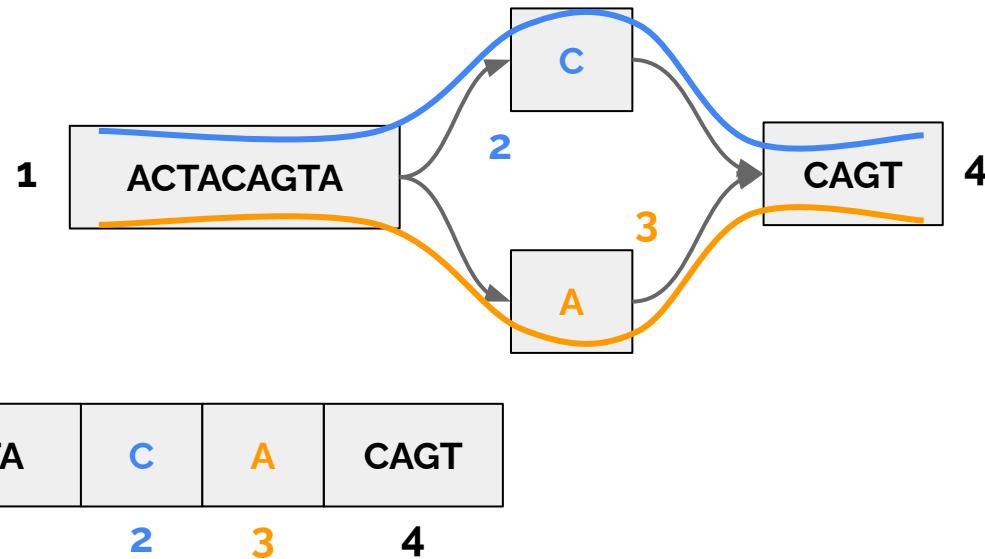
# Outline

- Introduction
- Method
- Result



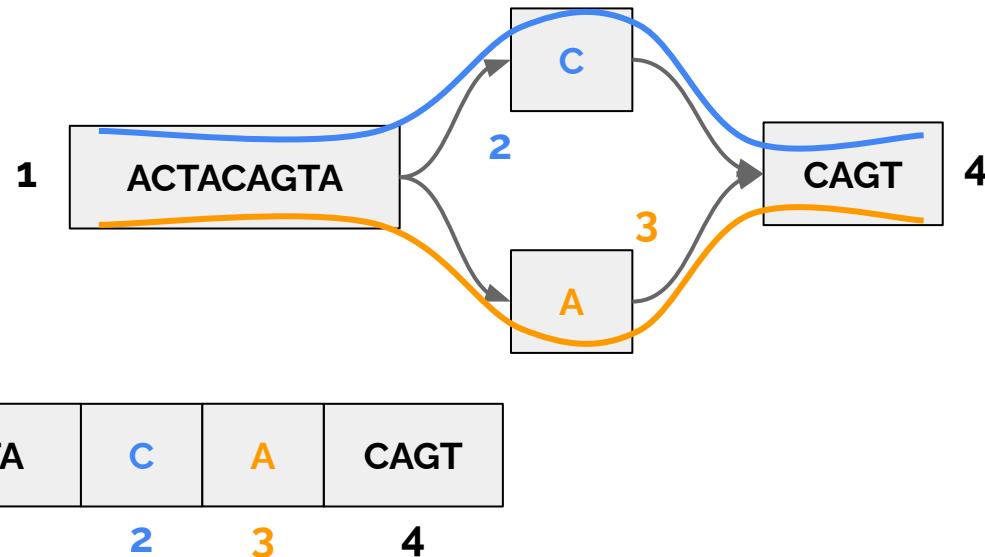
## Node coverage matrix

- Genome 1: ACTACAGT**A**CCAGT
- Genome 2: ACTACAGT**A**A**C**AGT



## Node coverage matrix

- Genome 1: ACTACAGT**A**CCAGT
- Genome 2: ACTACAGT**A**ACAGT

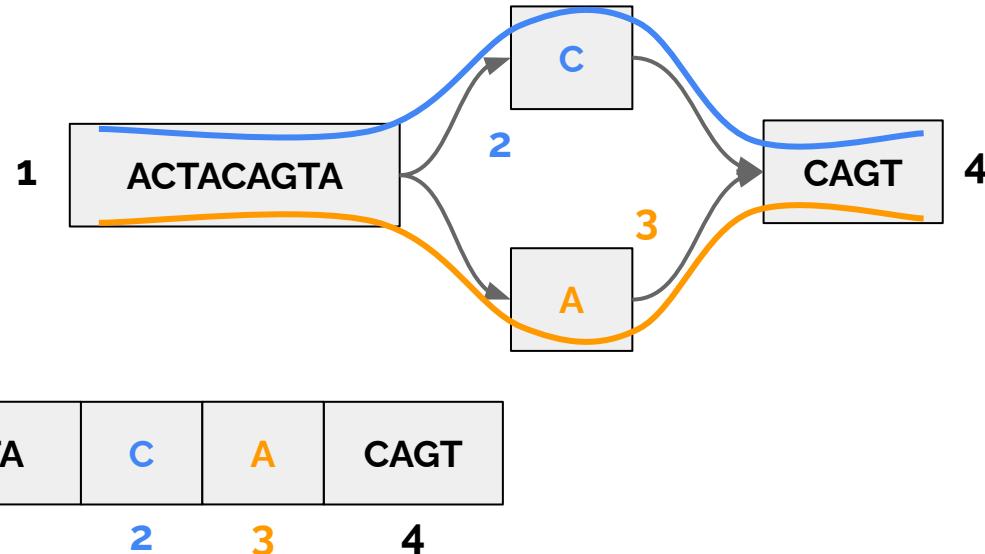


## Path coverage matrix

Genome 1	1	1	0	1
Genome 2	1	0	1	1
	1	2	3	4

## Node coverage matrix

- Genome 1: ACTACAGT**A**CCAGT
- Genome 2: ACTACAGT**A**A**C**AGT



Path coverage matrix

	1	2	3	4
Genome 1	1	1	0	1
Genome 2	1	0	1	1
	1	2	3	4

Sample coverage matrix

	1	2	3	4
Sample 1	29	31	1	27
Sample 2	28	0	30	28
	1	2	3	4

# Graph-based genotyping with COSIGT



<https://github.com/davidebolo1993/cosigt>

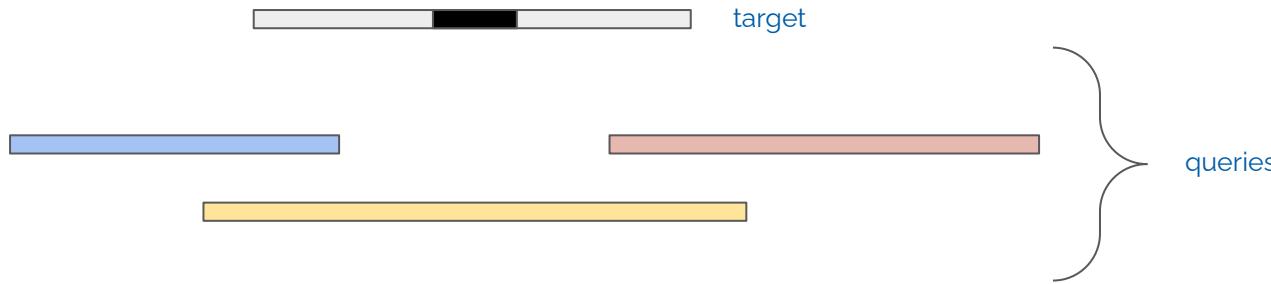
# Graph-based genotyping with COSIGT

build ROI-specific  
graph

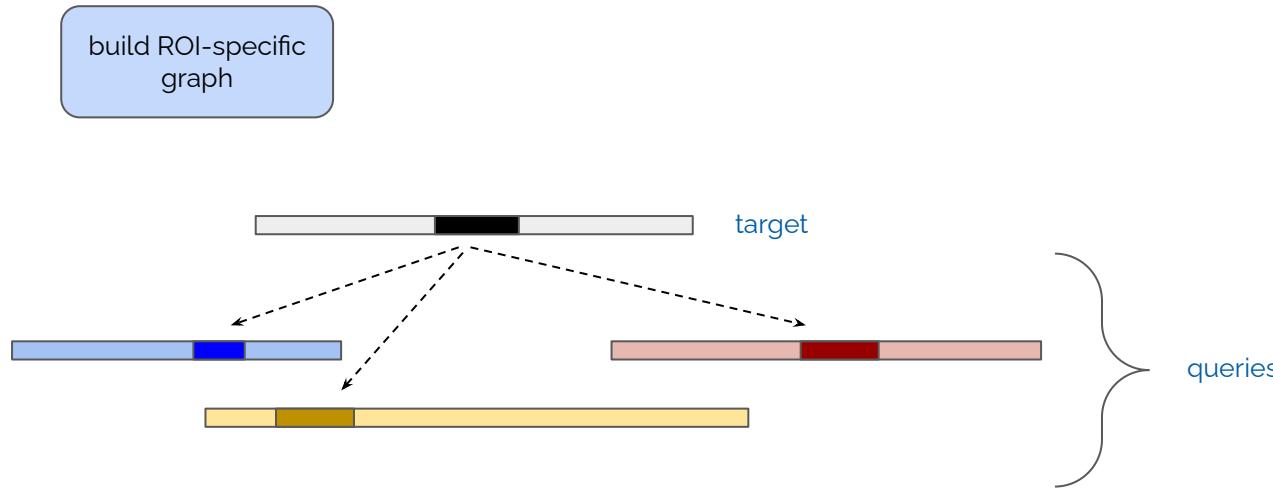


# Graph-based genotyping with COSIGT

build ROI-specific  
graph

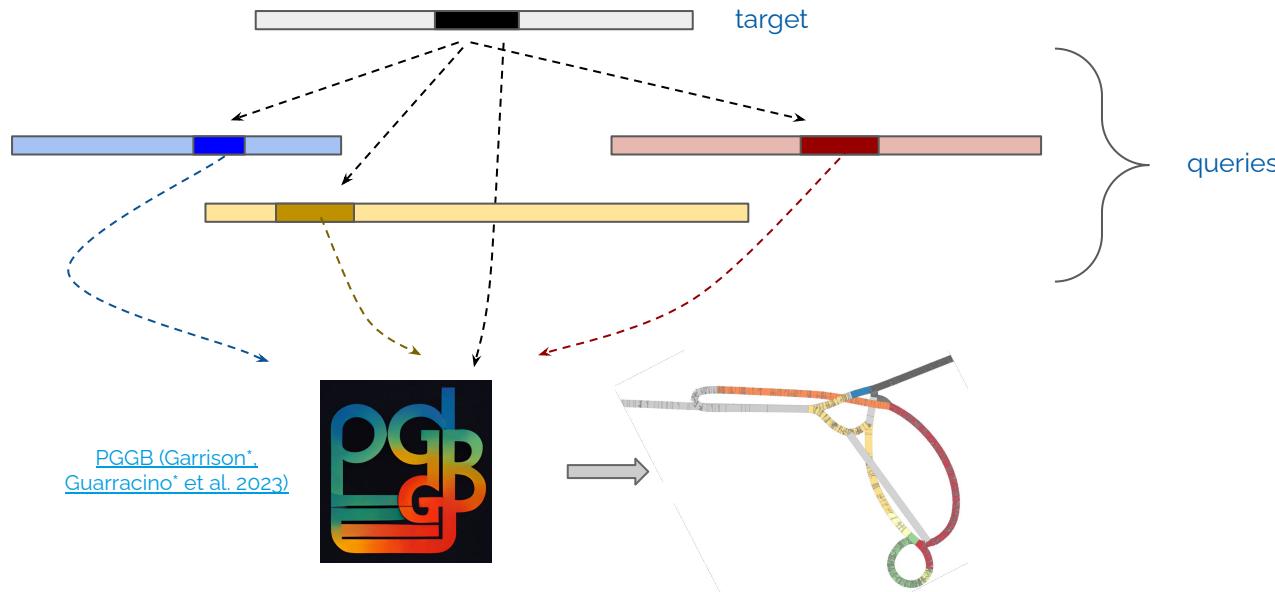


# Graph-based genotyping with COSIGT



# Graph-based genotyping with COSIGT

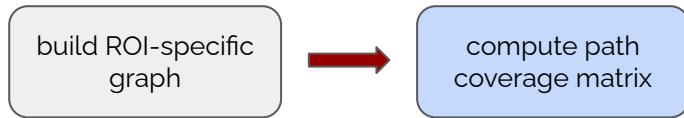
build ROI-specific graph



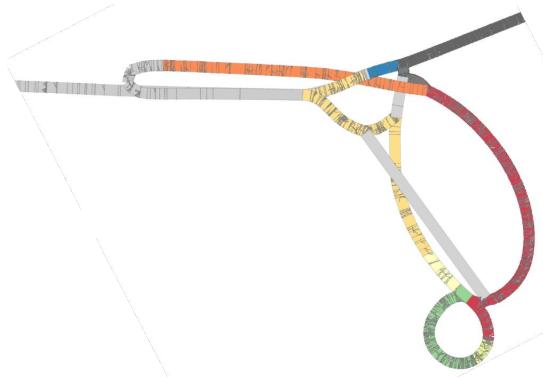
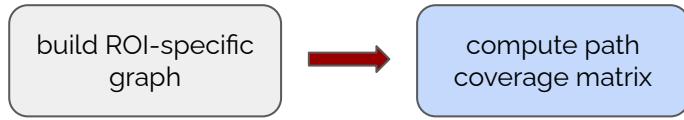
[PGGB \(Garrison\\*,  
Guarracino\\* et al. 2023\)](#)



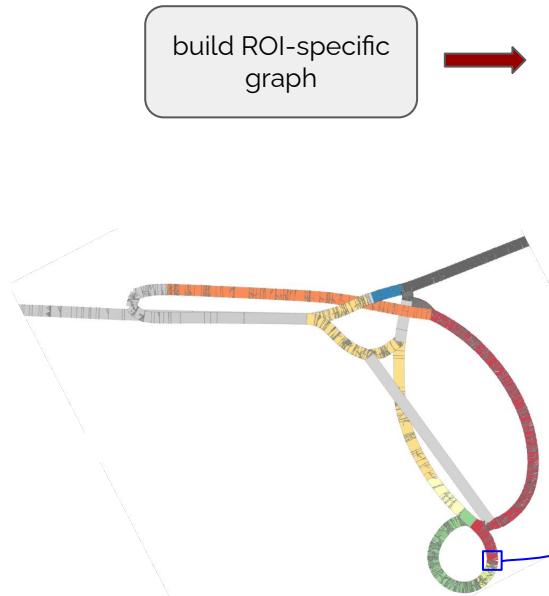
# Graph-based genotyping with COSIGT



# Graph-based genotyping with COSIGT



# Graph-based genotyping with COSIGT

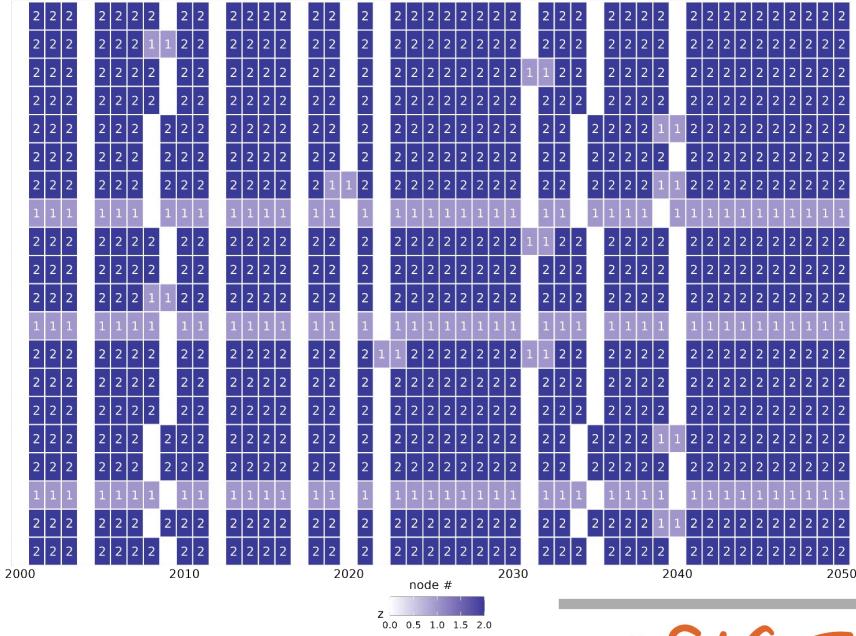


build ROI-specific  
graph



compute path  
coverage matrix

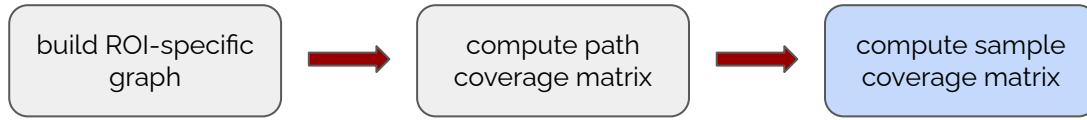
NA21309#1#JHEPC010000026.1.3251471-3328466  
HG03579#2#JAGYVT010000002.1.3306327-3383320  
HG03516#2#JAGYY5010000003.1.32087444-32164440  
HG03492#1#JHEP0100000049.1:6797876-16881236  
HG02886#1#JHAHO010000006.1:23426750-23503734  
HG02717#1#JAHAO5010000073.1:5229589-5306574  
HG02630#2#JHAIAOP010000058.1:24147247-24224230  
HG02572#2#JAHAAOV010000201.1:5173148-5223764  
HG02572#1#JAHAAOW010000052.1:1088220-1165211  
HG02559#2#JAGYVJ010000064.1:31930791-32014151  
HG02486#2#JAGYVL010000026.1:31926549-32003542  
HG02148#2#JHAMF010000034.1:3256183-3308808  
HG02080#2#JAHEDV010000002.1:24389454-24466449  
HG02055#1#JHEPK010000074.1:3306858-3383852  
HG01358#1#JAGYZB010000008.1:7958573-8041931  
HG01258#2#JAGYYU010000011.1:27153790-27230777  
HG01109#2#JAHFOZ010000001.1:28469510-28546498  
HG01106#2#JAHAMBO100000019.1:116915379-116966004  
HG01071#1#JAHBCF010000017.1:706410-783394  
chr13#chr6.31825263-31908622



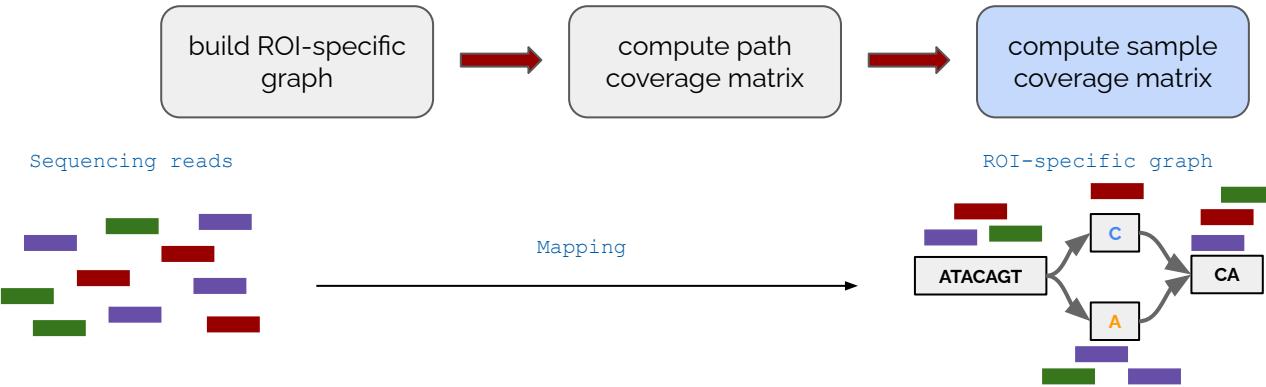
**COSIGT**  
COSINE SIMILARITY-BASED GENOTYPER



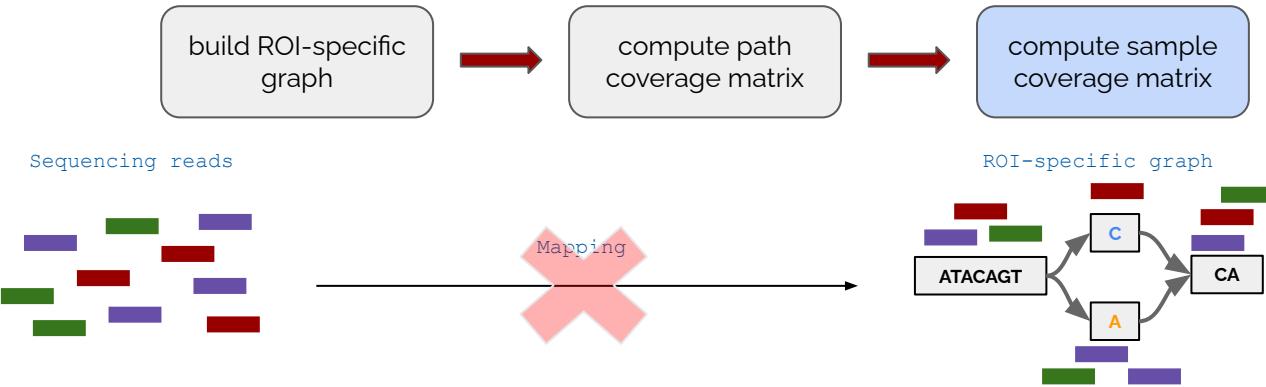
# Graph-based genotyping with COSIGT



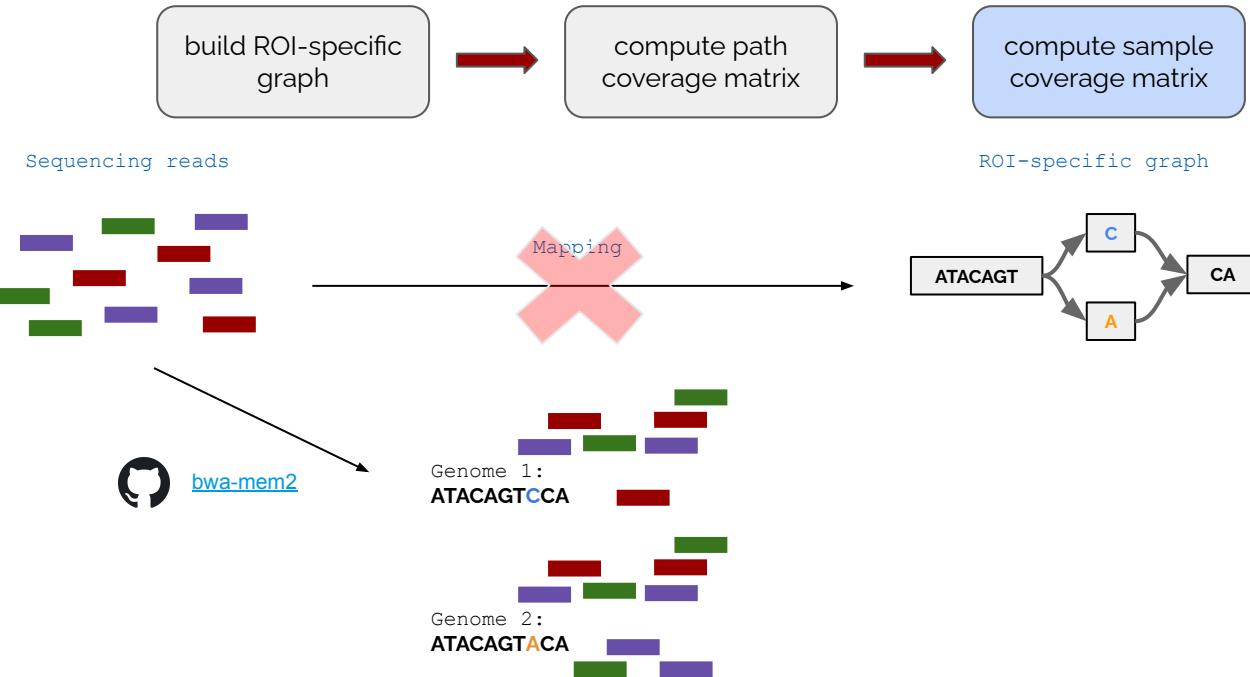
# Graph-based genotyping with COSIGT



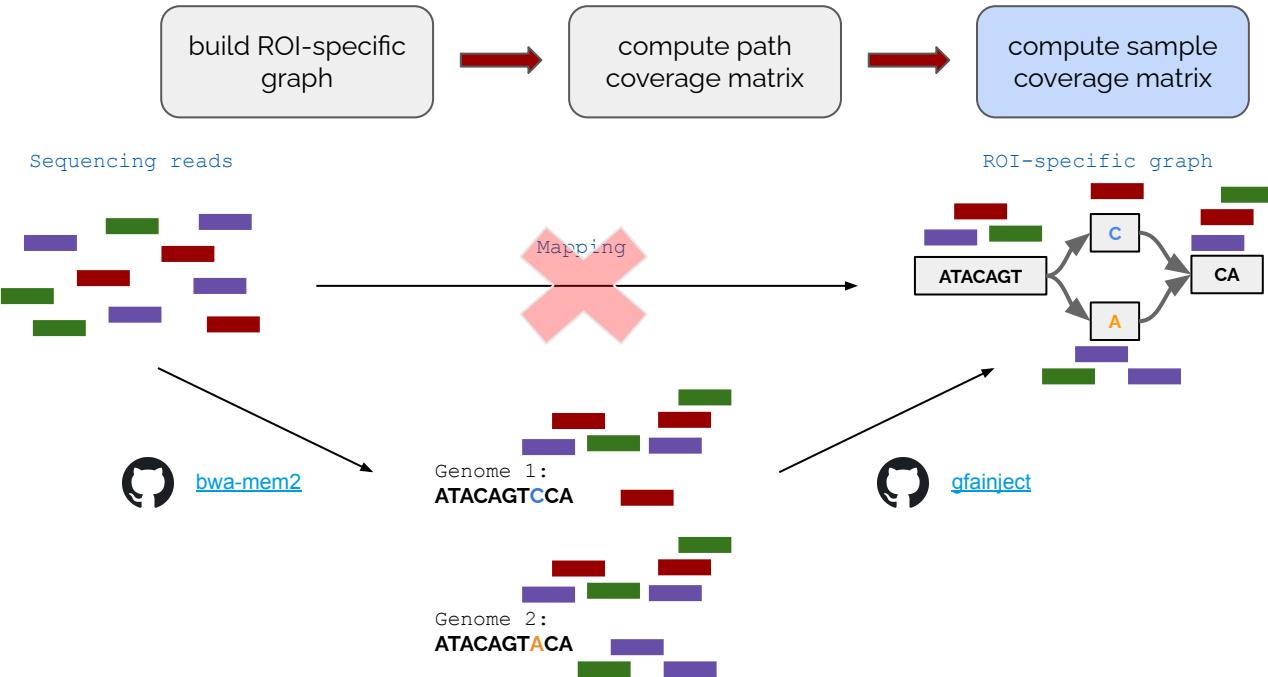
# Graph-based genotyping with COSIGT



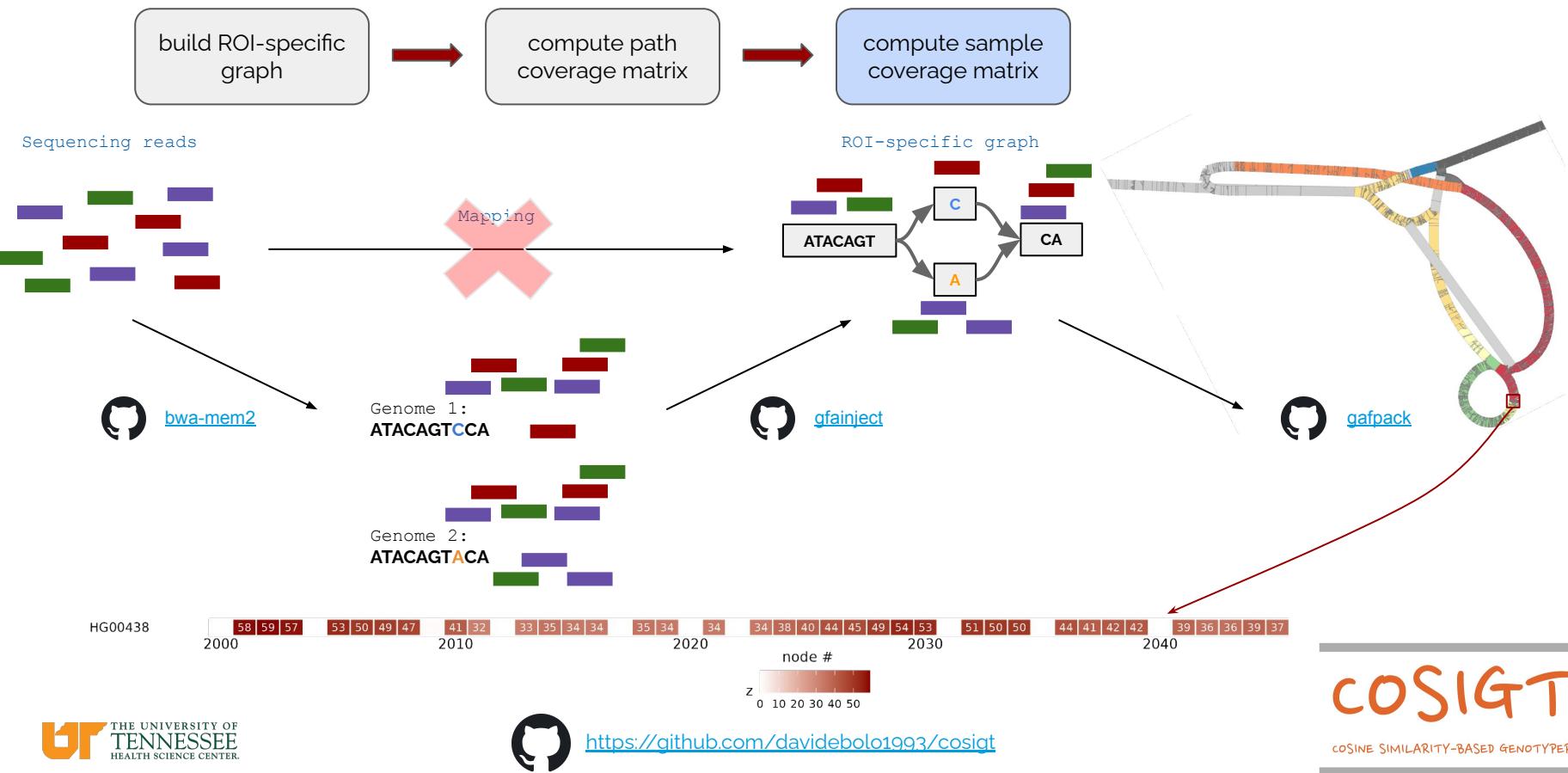
# Graph-based genotyping with COSIGT



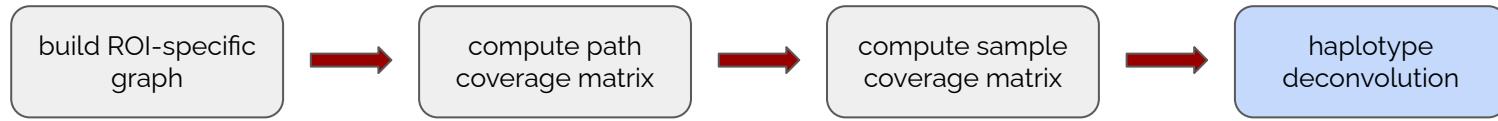
# Graph-based genotyping with COSIGT



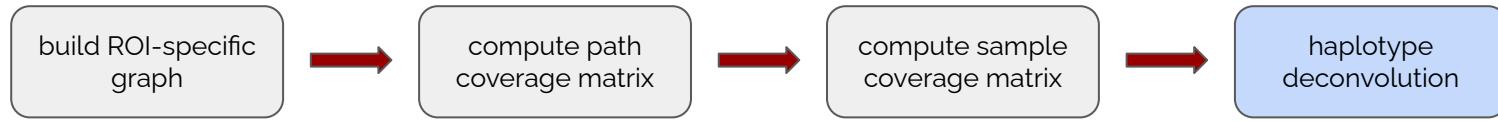
# Graph-based genotyping with COSIGT



# Graph-based genotyping with COSIGT



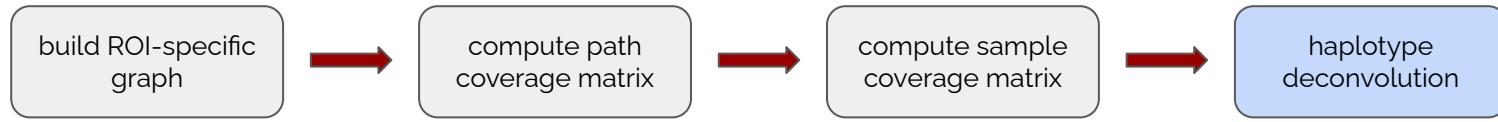
# Graph-based genotyping with COSIGT



**haplotype deconvolution:** find the combination of  $N$  path coverage vectors best representing the trend observed in the sample coverage vector



# Graph-based genotyping with COSIGT



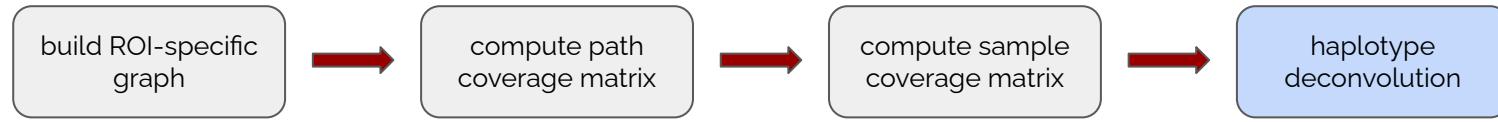
**haplotype deconvolution:** find the combination of  $N$  path coverage vectors best representing the trend observed in the sample coverage vector

path coverage matrix

h#1	[ ]	[ ]	5	4	5	2	2	4	4	2	2	3	[ ]	[ ]
h#2	[ ]	[ ]	5	5	5	6	6	6	7	4	4	4	[ ]	[ ]
h#3	[ ]	[ ]	1	1	2	1	3	1	1	1	1	1	[ ]	[ ]



# Graph-based genotyping with COSIGT



**haplotype deconvolution:** find the combination of  $N$  path coverage vectors best representing the trend observed in the sample coverage vector

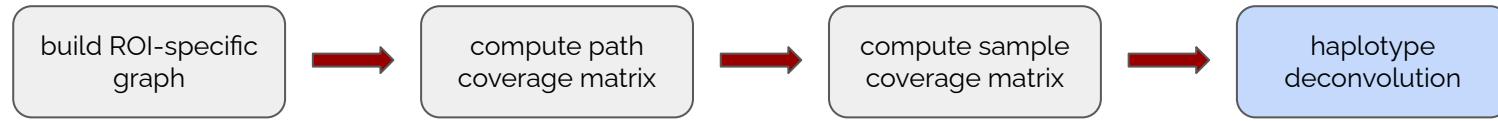
path coverage matrix

h#1	[ ]	[ ]	5	4	5	2	2	4	4	2	2	3	[ ]	[ ]
h#2	[ ]	[ ]	5	5	5	6	6	6	7	4	4	4	[ ]	[ ]
h#3	[ ]	[ ]	1	1	2	1	3	1	1	1	1	1	[ ]	[ ]

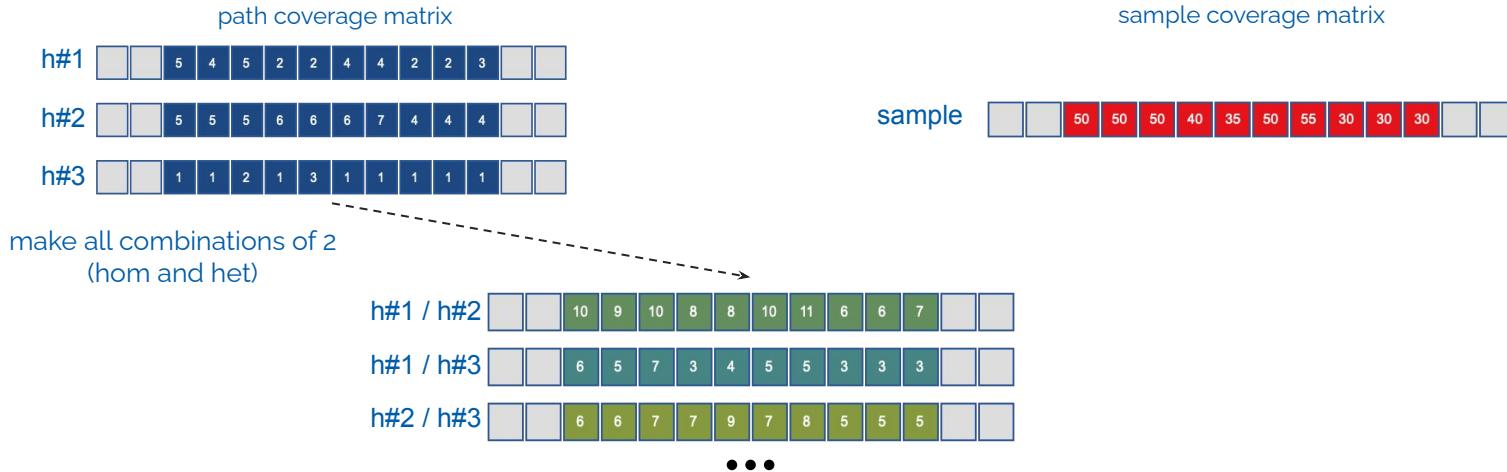
sample coverage matrix

sample	[ ]	[ ]	50	50	50	40	35	50	55	30	30	30	[ ]	[ ]
--------	-----	-----	----	----	----	----	----	----	----	----	----	----	-----	-----

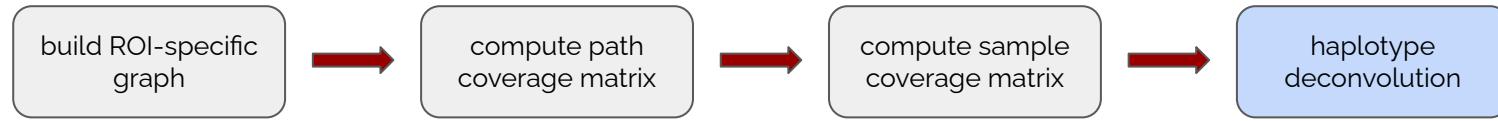
# Graph-based genotyping with COSIGT



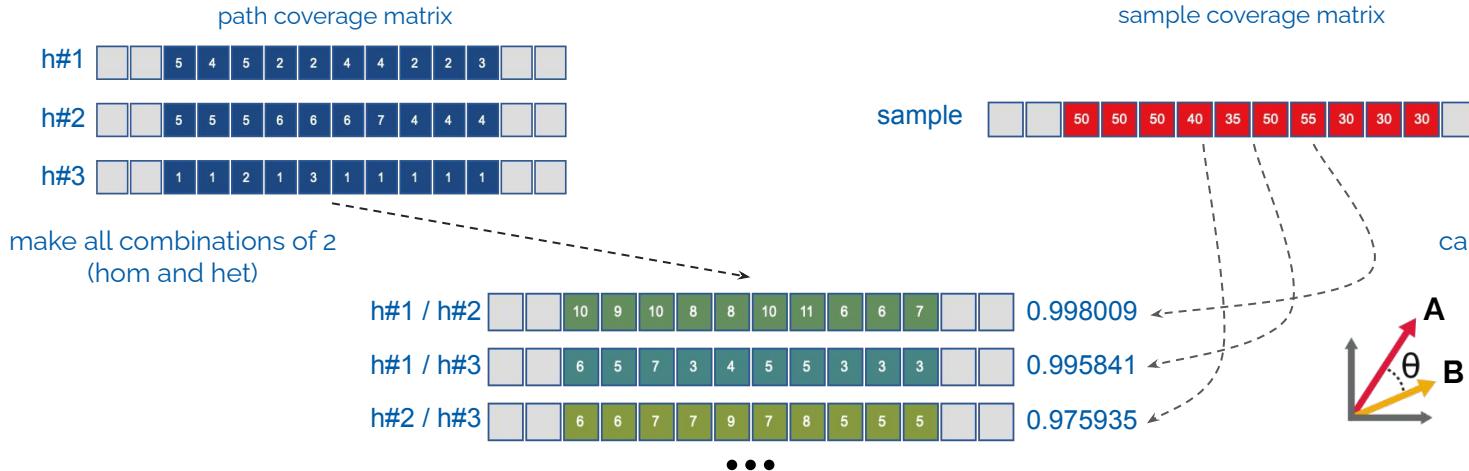
**haplotype deconvolution:** find the combination of  $N$  path coverage vectors best representing the trend observed in the sample coverage vector



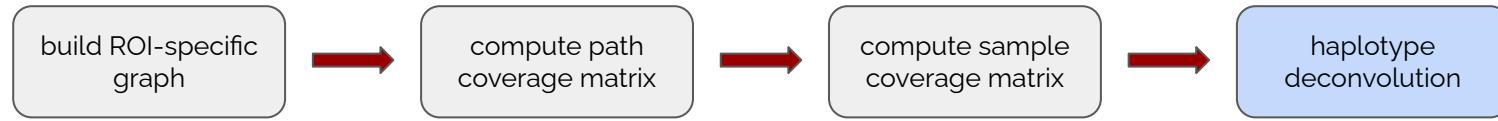
# Graph-based genotyping with COSIGT



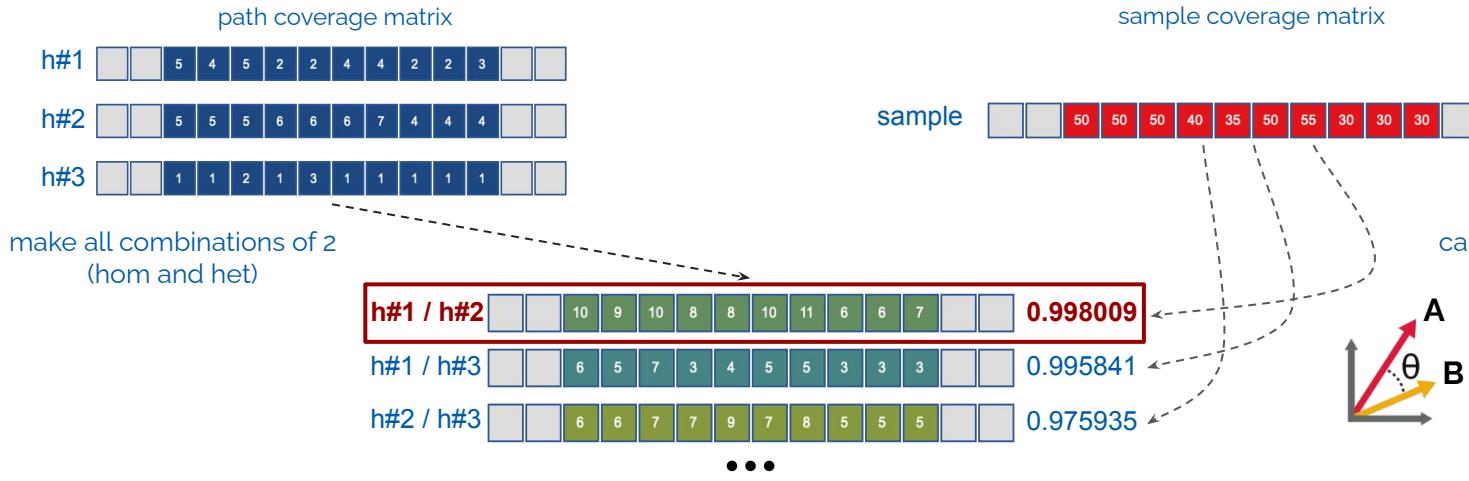
**haplotype deconvolution:** find the combination of  $N$  path coverage vectors best representing the trend observed in the sample coverage vector



# Graph-based genotyping with COSIGT



**haplotype deconvolution:** find the combination of  $N$  path coverage vectors best representing the trend observed in the sample coverage vector

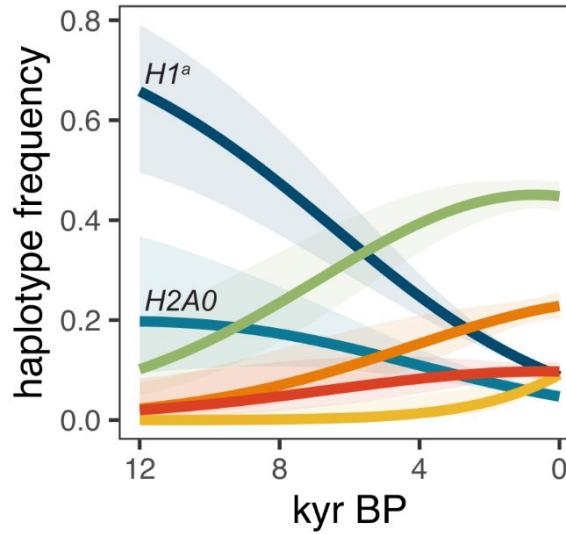


calculate cosine similarity

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

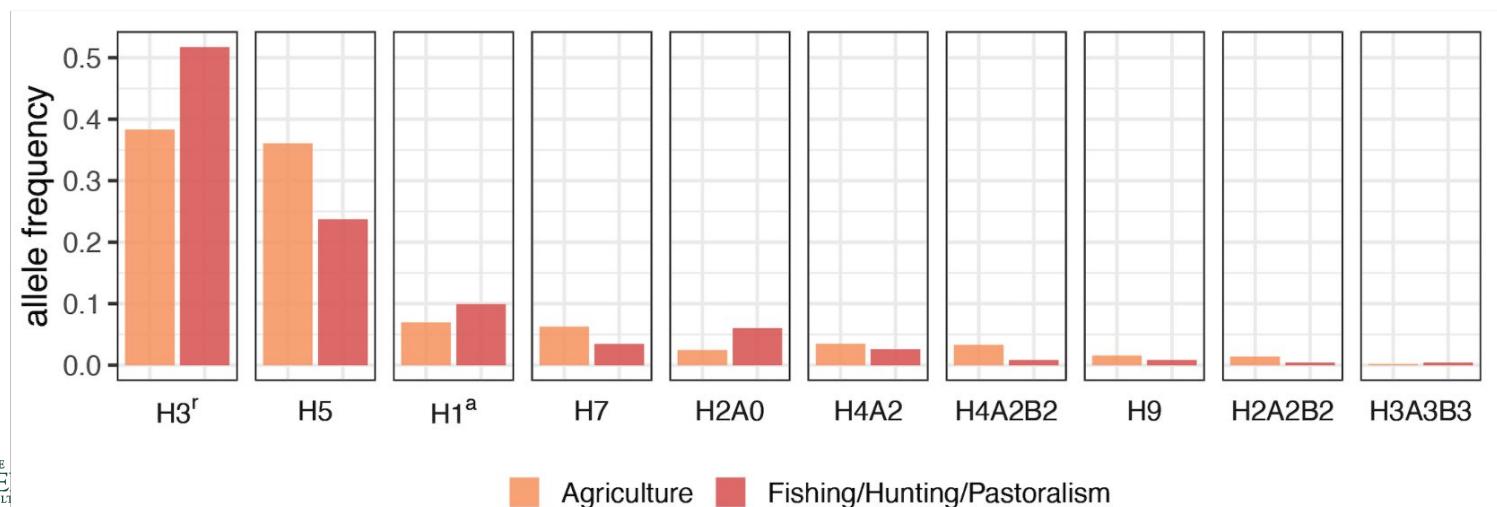
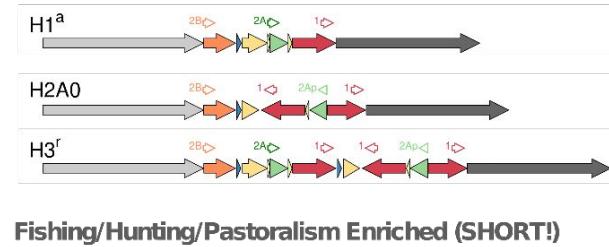
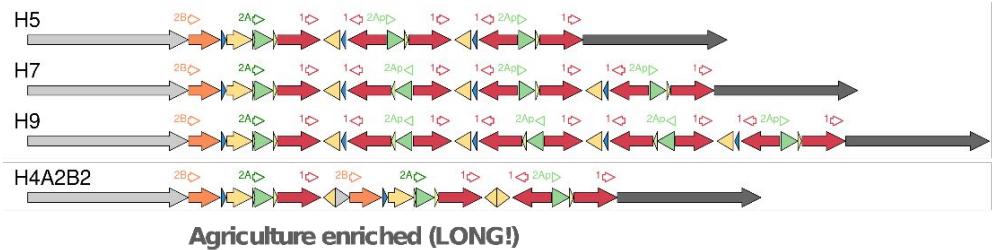
# Outline

- Introduction
- Method
- Result

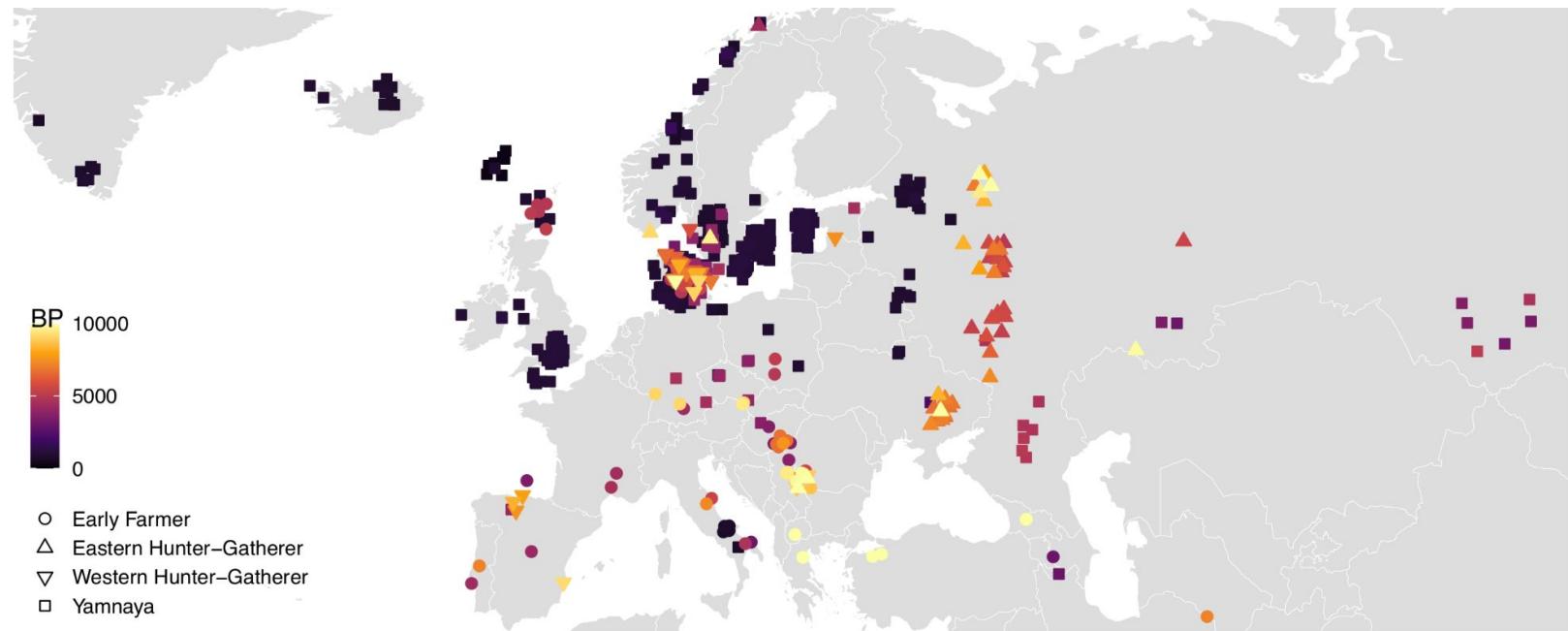


# Allele frequencies of complex structural haplotypes

Pangenome-based haplotype deconvolution for ~5600 short read haplotypes

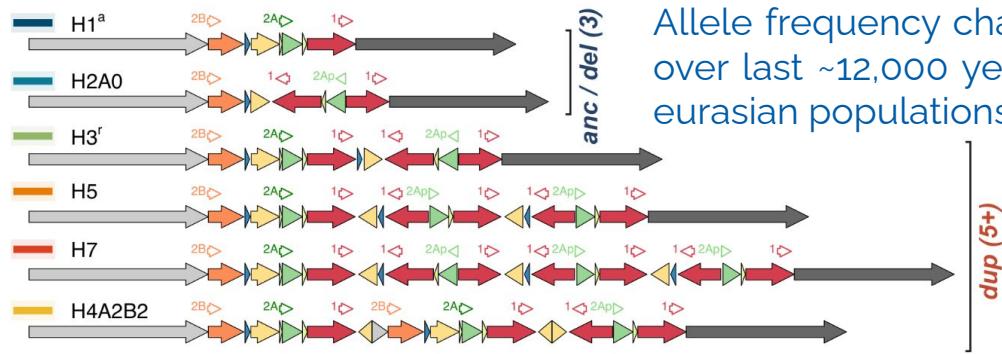
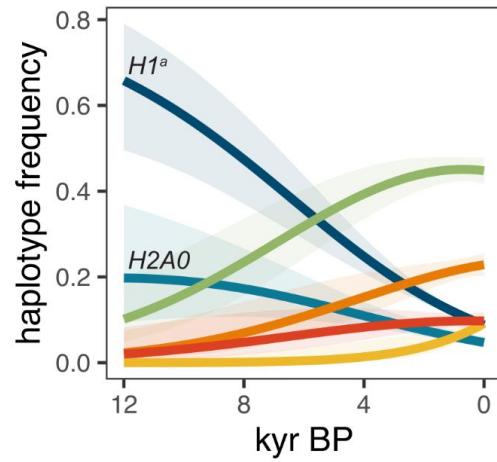


# Recent evolution from 533 ancient European genomes



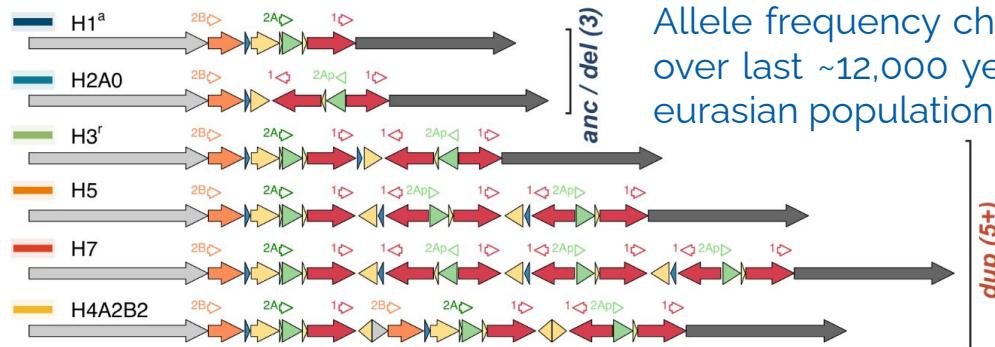
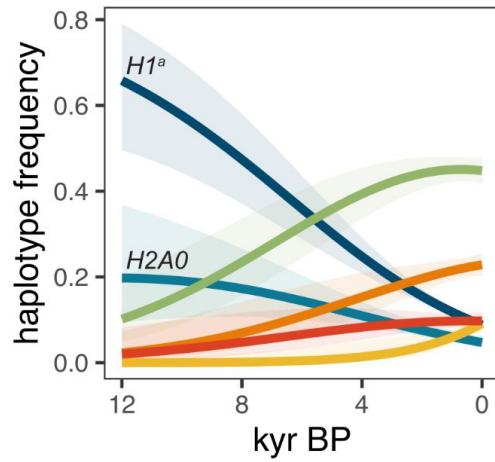
Data from [Allentoft et al. 2024, Nature](#)

# Evidence for selection of high-copy amylase haplotypes

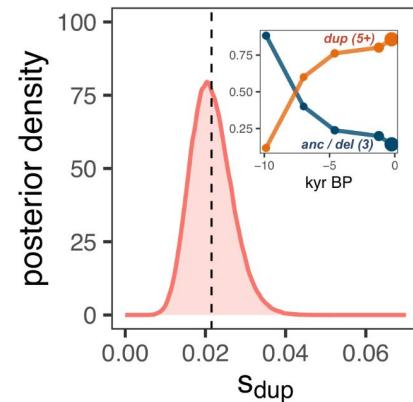


Allele frequency changes over last ~12,000 years in eurasian populations.

# Evidence for selection of high-copy amylose haplotypes



Allele frequency changes over last ~12,000 years in eurasian populations.



Selection coefficient of 0.02 is equivalent to selection at lactase ([Mathieson, 2020](#)).

# Panpublication

## Article

# Recurrent evolution and selection shape structural diversity at the amylase locus

<https://doi.org/10.1038/s41586-024-07911-1>

Received: 29 November 2023

Accepted: 6 August 2024

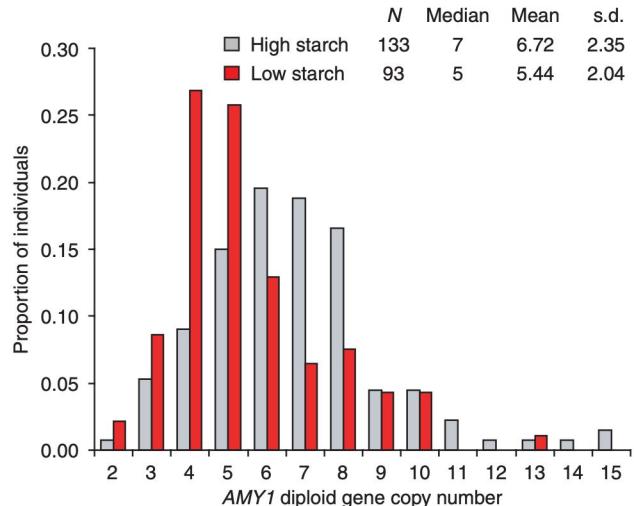
Published online: 04 September 2024

Davide Bolognini<sup>1,10</sup>, Alma Halgren<sup>2,10</sup>, Runyang Nicolas Lou<sup>2,10</sup>, Alessandro Raveane<sup>1,10</sup>, Joana L. Rocha<sup>2,10</sup>, Andrea Guaracino<sup>3</sup>, Nicole Soranzo<sup>1,4,5,6,7</sup>, Chen-Shan Chin<sup>8</sup>, Erik Garrison<sup>3</sup>✉ & Peter H. Sudmant<sup>2,9</sup>✉



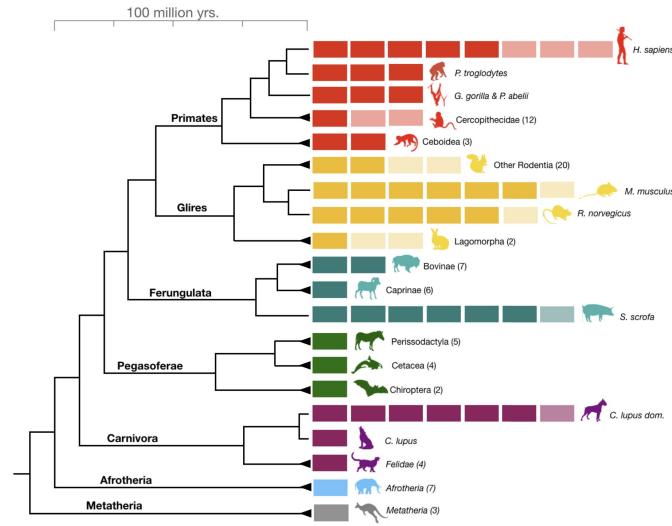


# Diet correlates with amylase copy number



Across humans

[Perry et al., 2007](#)

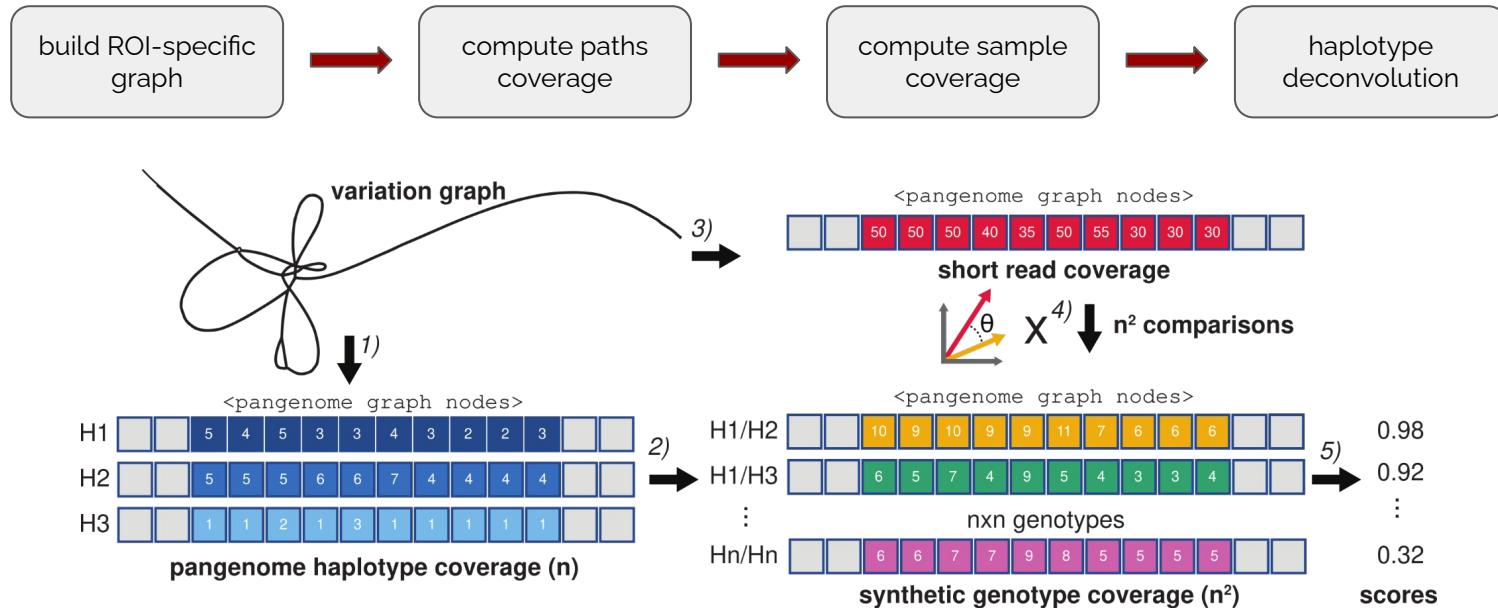


Across mammals

[Pajic et al., 2019](#)

No evidence of recent selection

# Graph-based genotyping with COSIGT works!



- 100% accurate when mapping individuals with both long reads and short reads ( $n=35$ )
- ~95% accurate by mendelian inheritance patterns (561 short read trios)
- 95-99% accurate by copy number prediction compared to read depth ( $n=3102$ )

# Allele frequencies of complex structural haplotypes

Pangenome-based haplotype deconvolution for ~8000 short read haplotypes

