

PRESERVING INTRA-PATIENT VARIANCE IMPROVES PHYLOGENETIC INFERENCE OF HIV TRANSMISSION

AUGUST GUANG

acknowledgments

My CFAR Research Group:

Casey Dunn, Associate Professor of Evolutionary Biology

Rami Kantor, Associate Professor of Medicine

Mia Coetzer, Assistant Professor of Medicine

Mark Howison, Director of Data Science

Colin MacLean, Research Programmer & Charles Lawrence, Professor of Applied Mathematics

My DunnLab Research Group:

Casey Dunn (again)

Zachary Lewis, Postdoc

Catriona Munro, PhD Candidate

Alex Damian Serrano, PhD Candidate

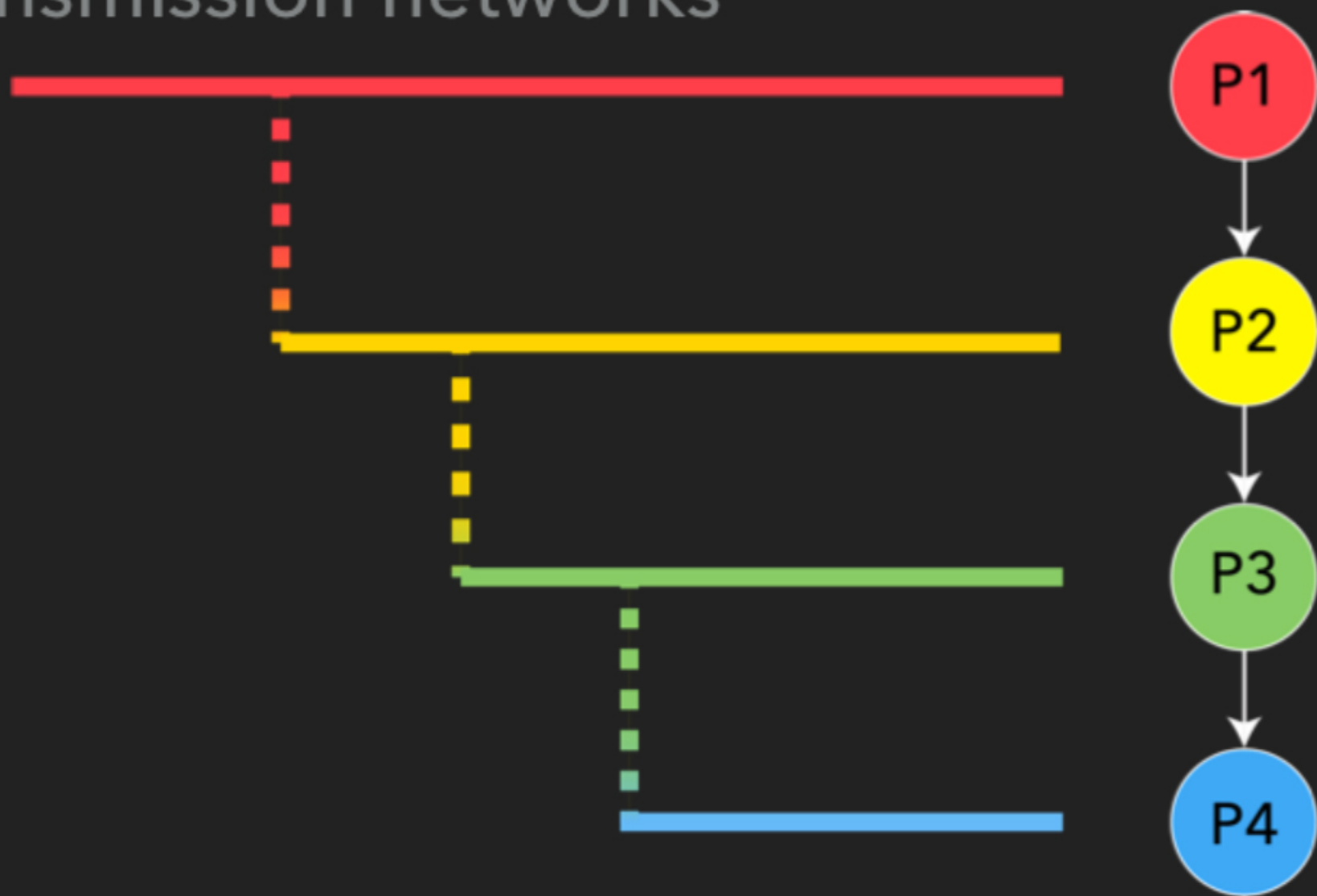


TRANSMISSION NETWORKS AS PHYLOGENIES

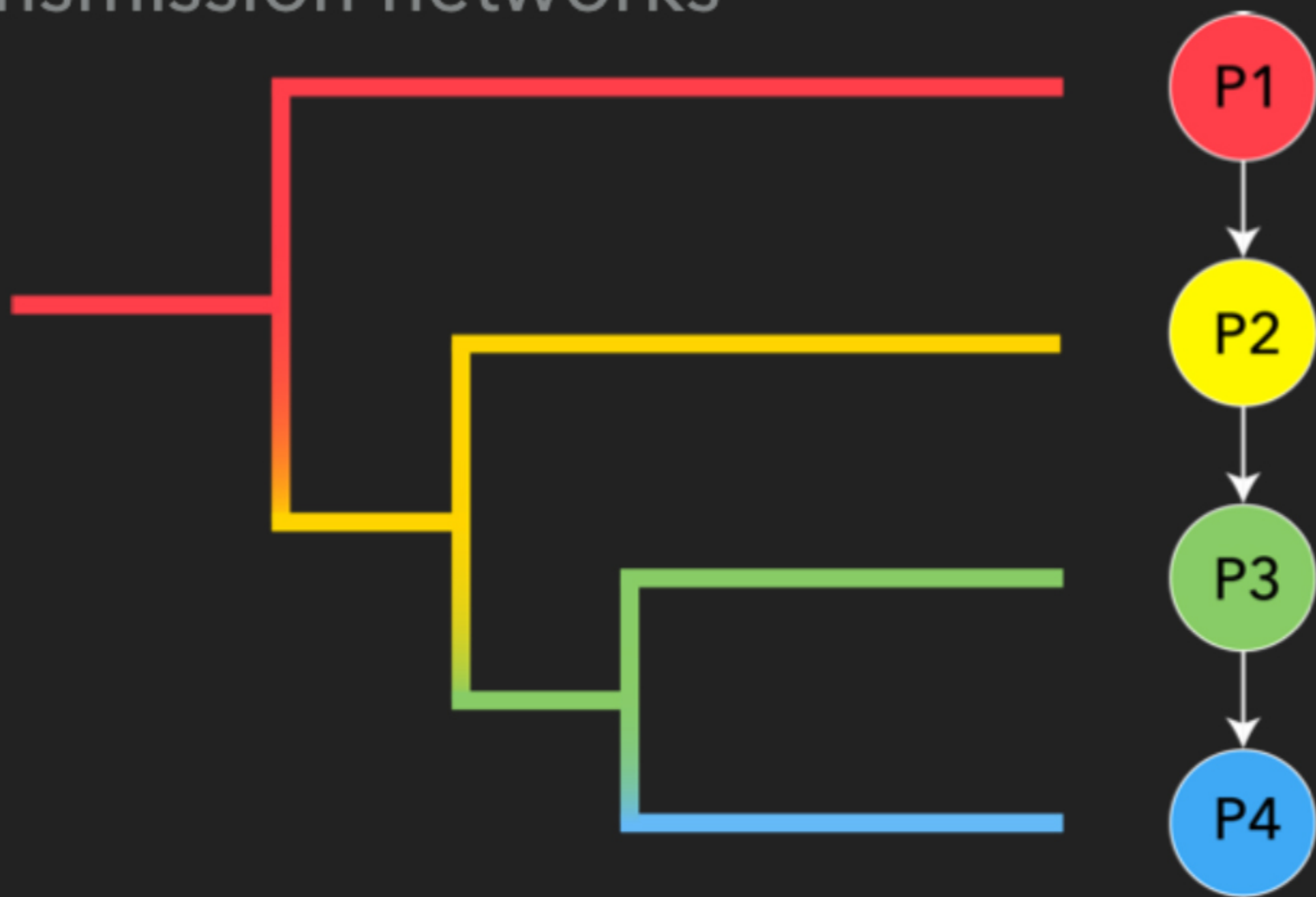


Phylogenetic trees have information about relatedness of organisms at tips

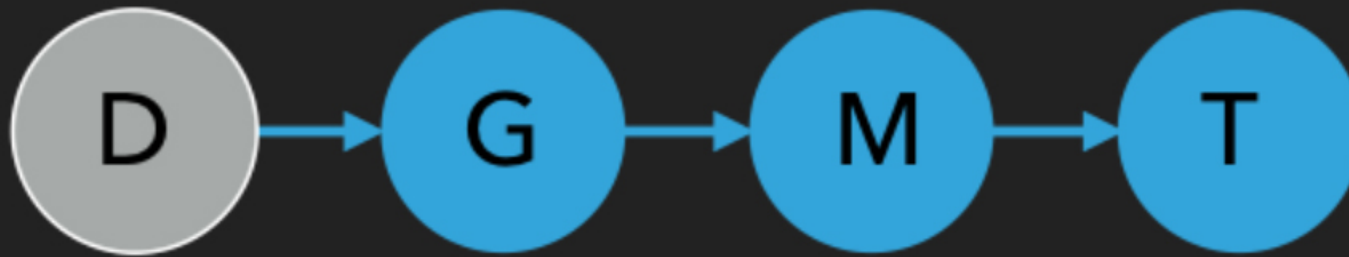
in the absence of reliable patient contact histories, phylogenies can be proxies for transmission networks



in the absence of reliable patient contact histories, phylogenies can be proxies for transmission networks

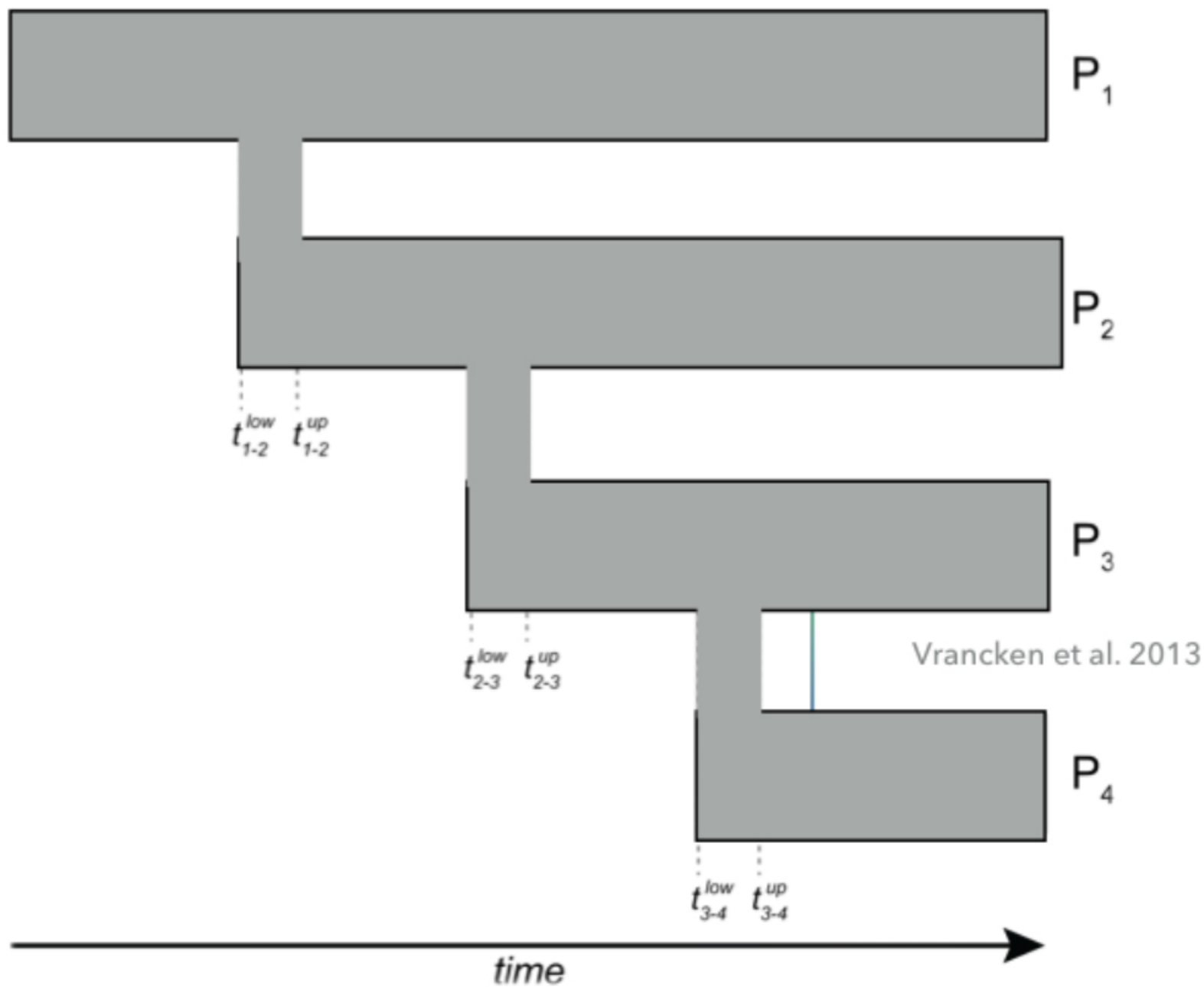


the phylogenetic workflow

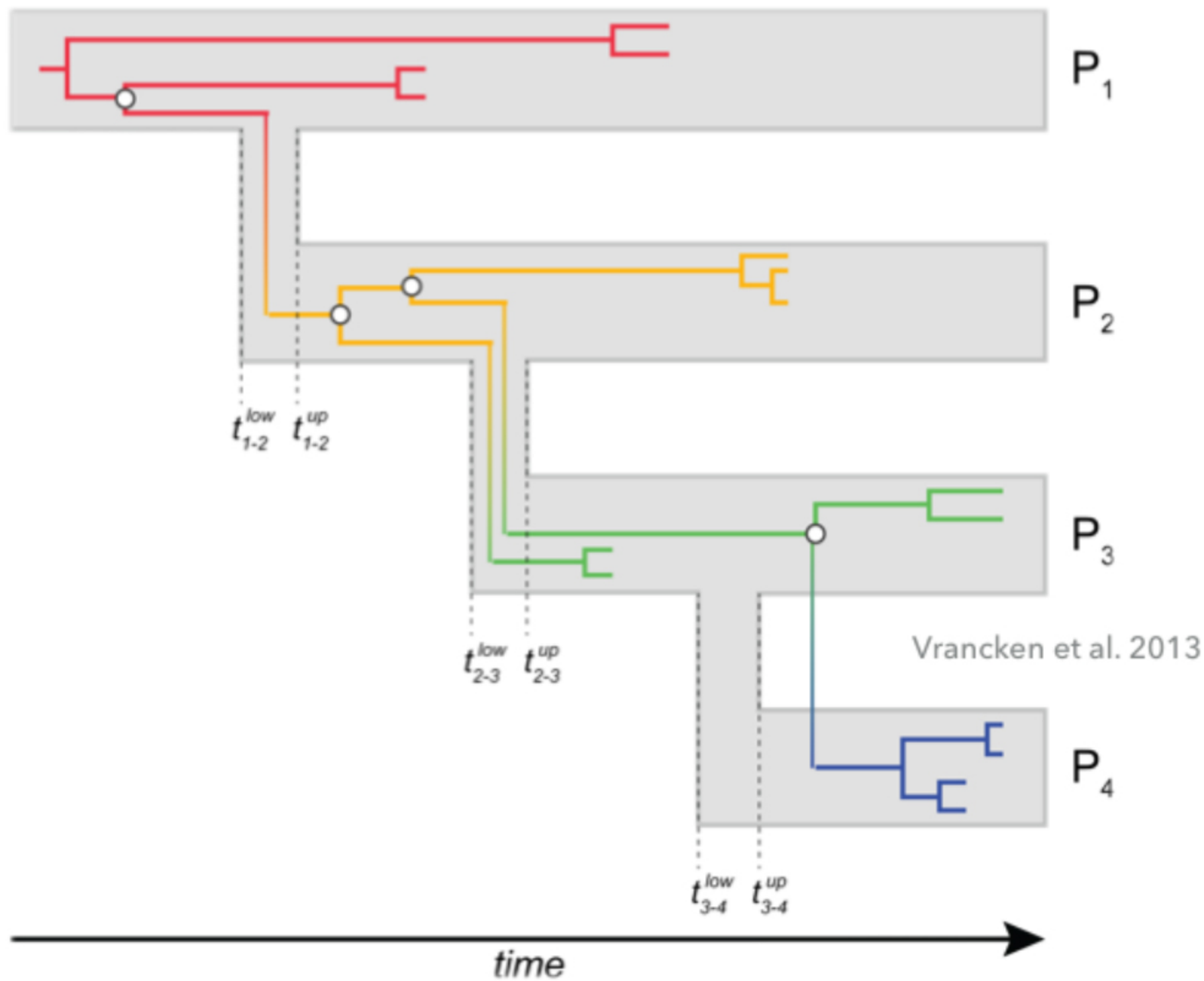


- ▶ D: Data: reads
- ▶ G: Genomes
- ▶ M: Multiple sequence alignment
- ▶ T: Transmission tree

we assume the tips of these trees are single entities



but with rare exceptions they are summaries





this is not an issue if we suspect that the variation can be summarized by the mode, i.e. consensus genome



but is it an issue if it cant?

after all, we ultimately care about the
transmission tree...

**PRESERVING INTRA-
PATIENT VARIATION IS
IMPORTANT**

a simple thought experiment



A

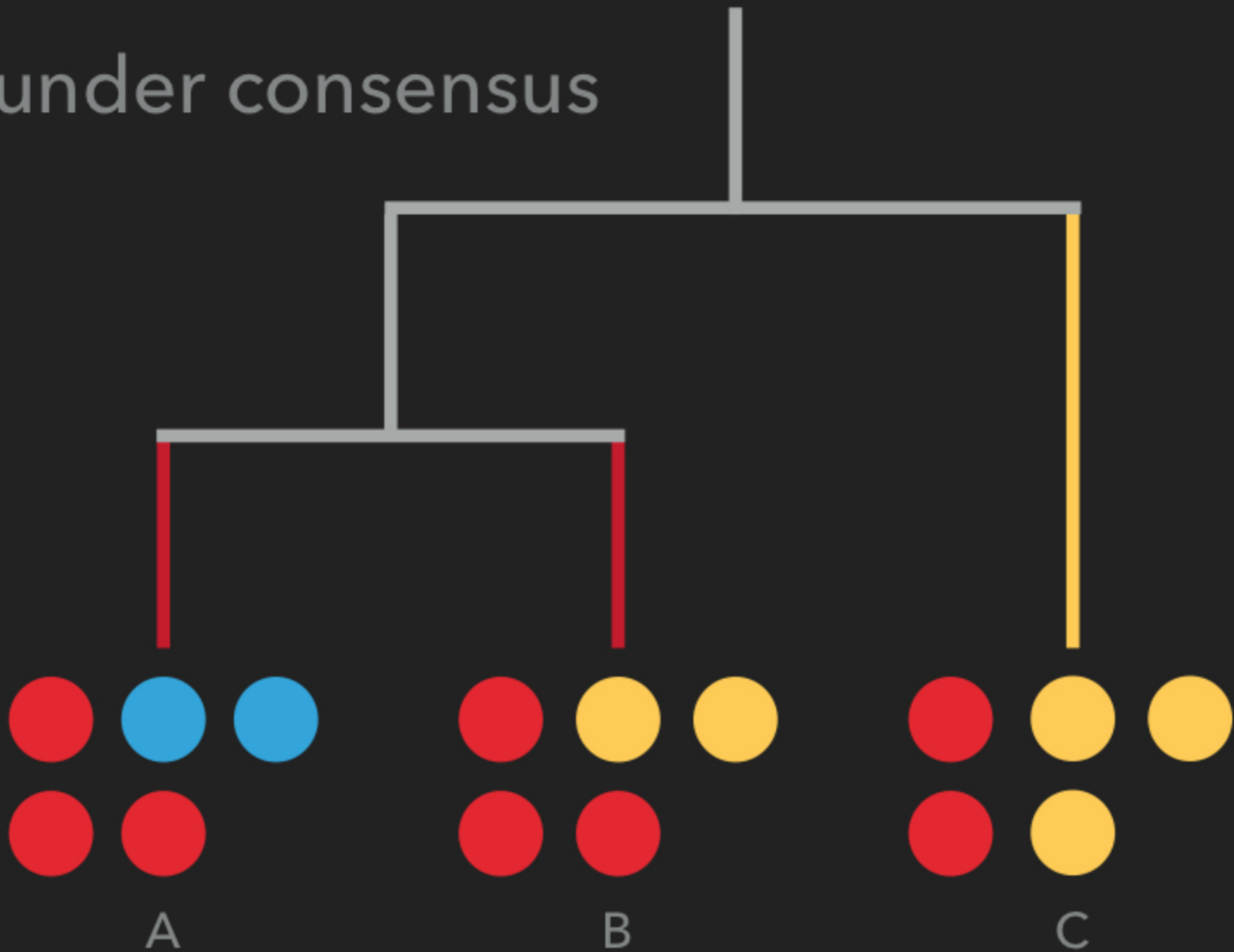


B

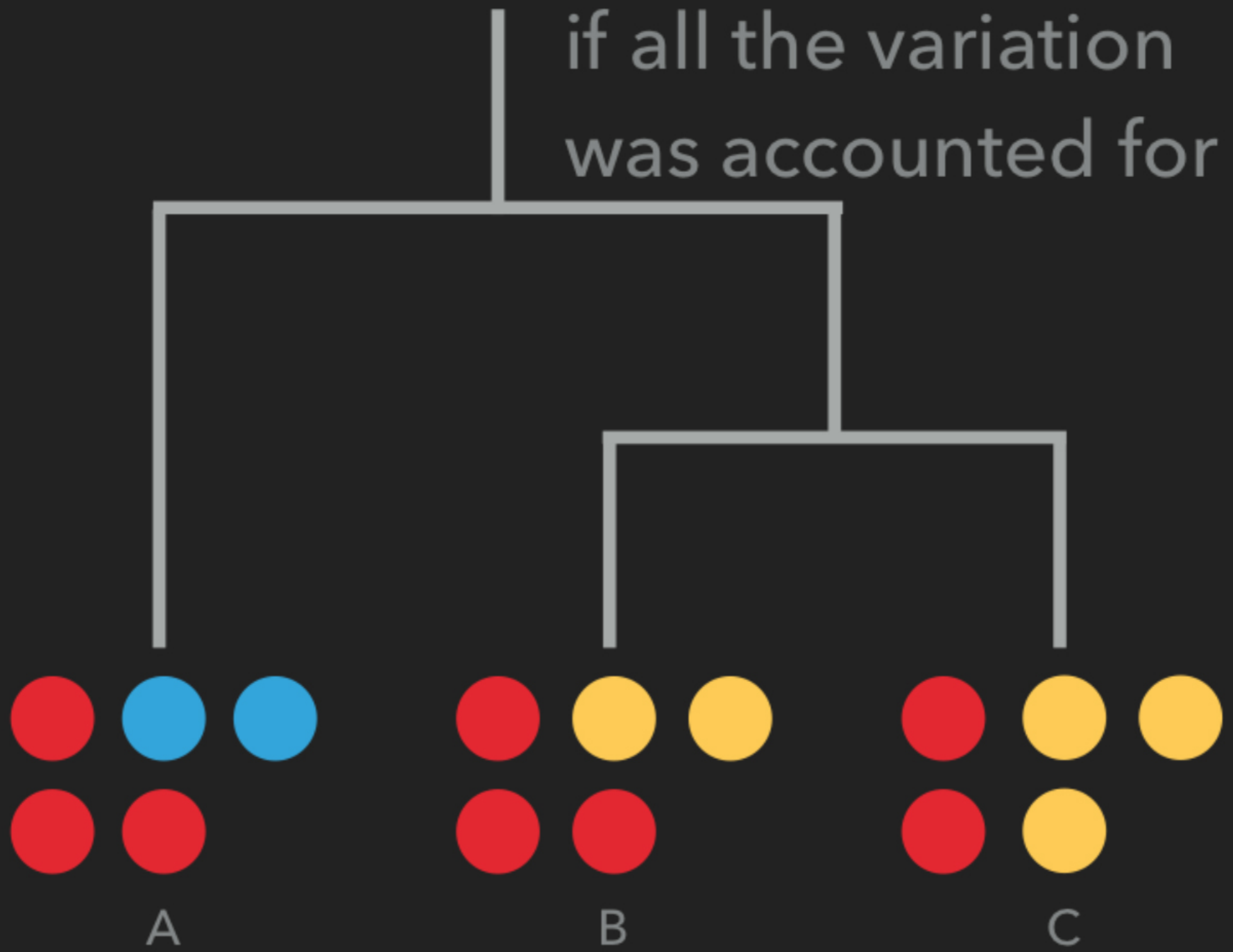


C

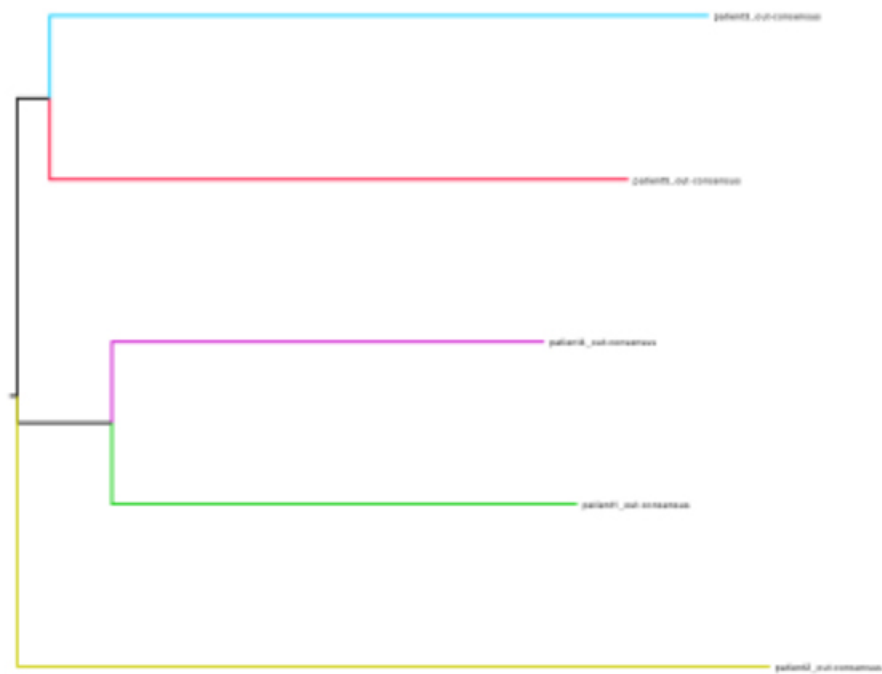
under consensus



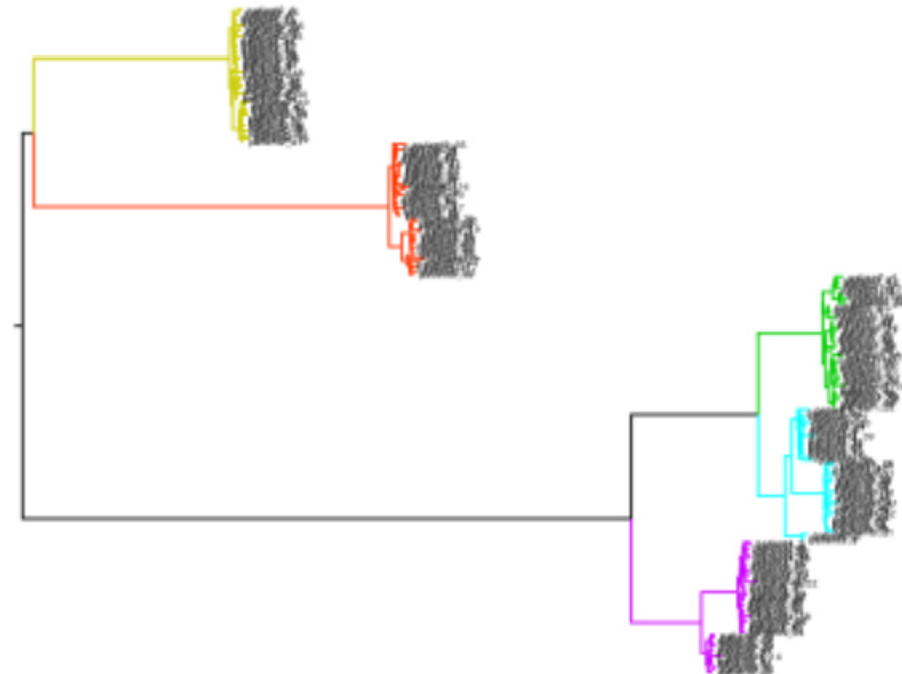
if all the variation
was accounted for



some simple simulations



RAxML tree from consensus genome



true tree

so we don't recover the same tree...

let's try to account for that variation!



A



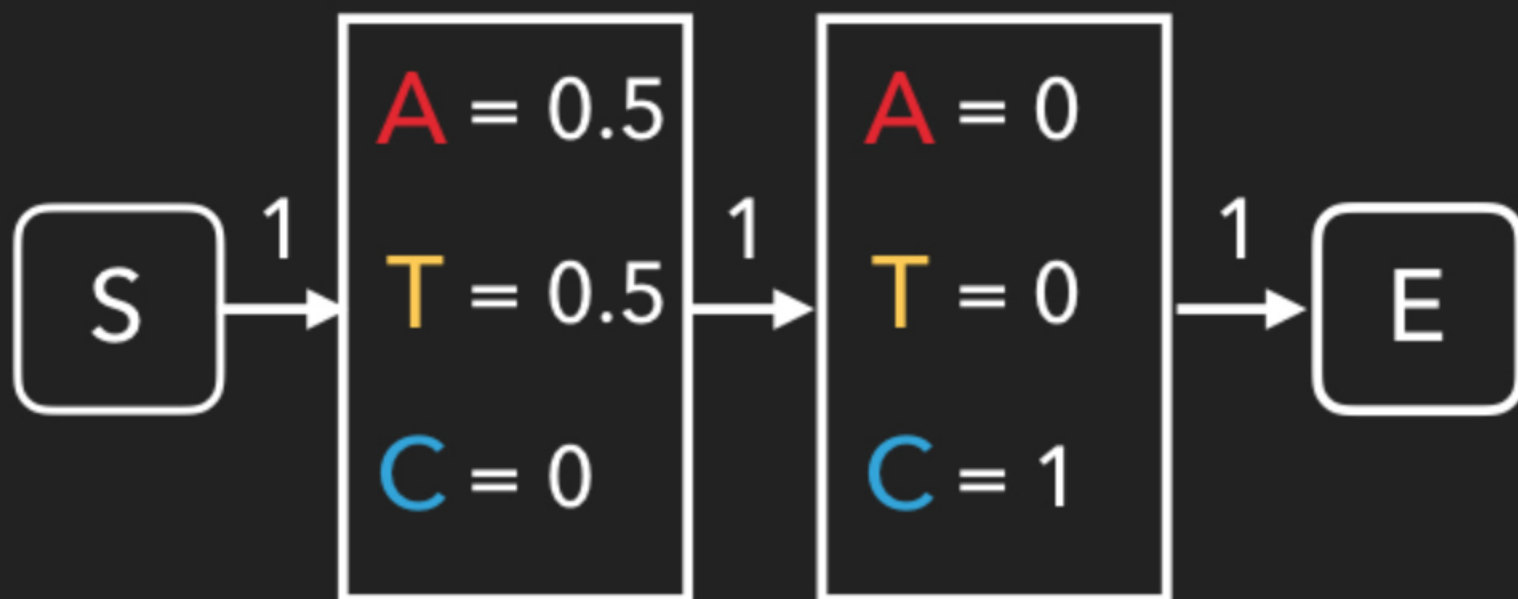
B



C

profile Hidden Markov Models

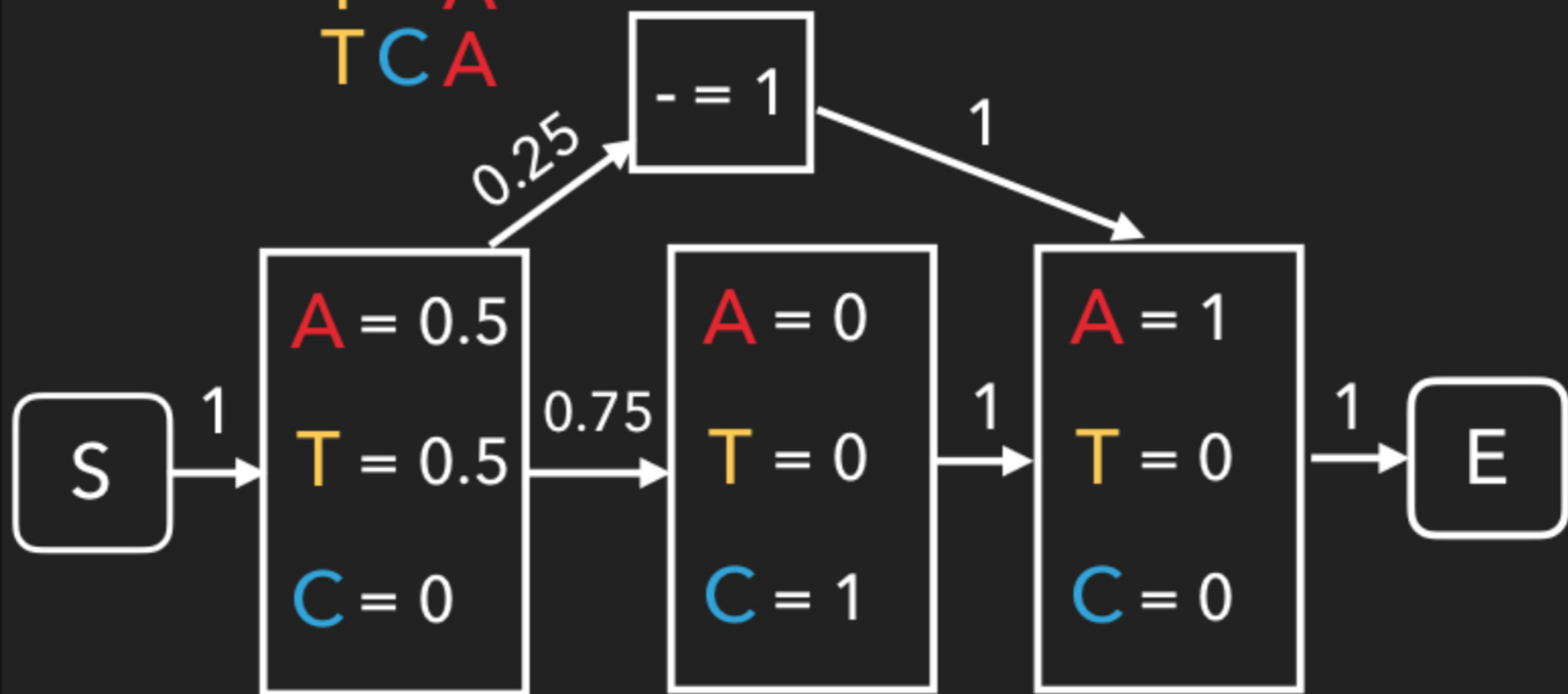
A C
A C
T C
T C



indels are hidden states

ACA
ACA
T - A
TCA

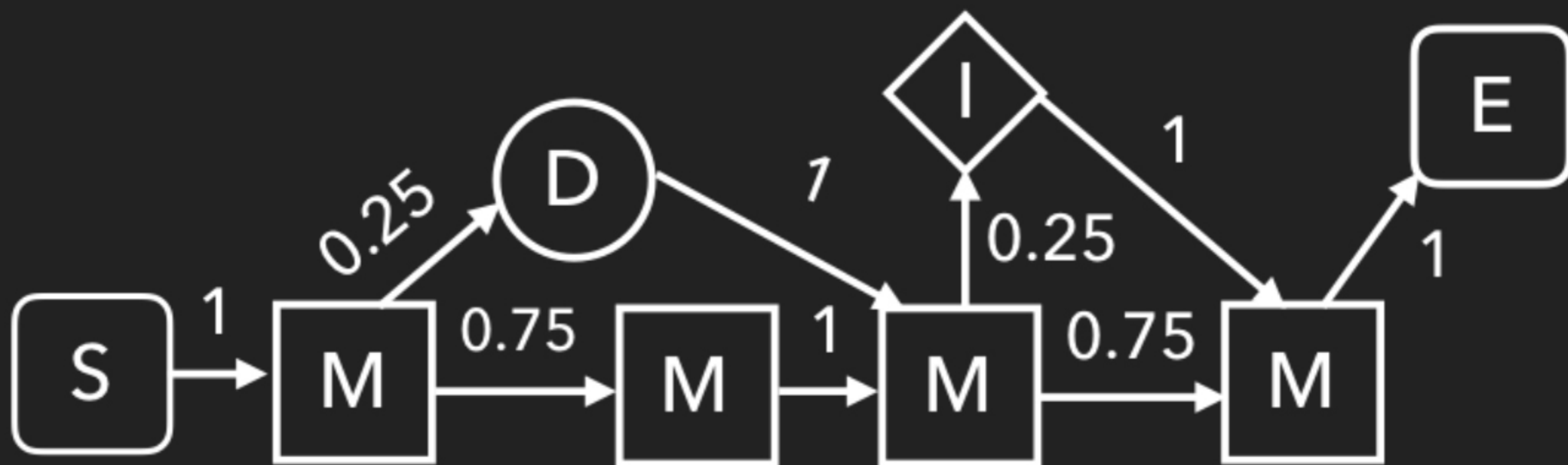
deletion



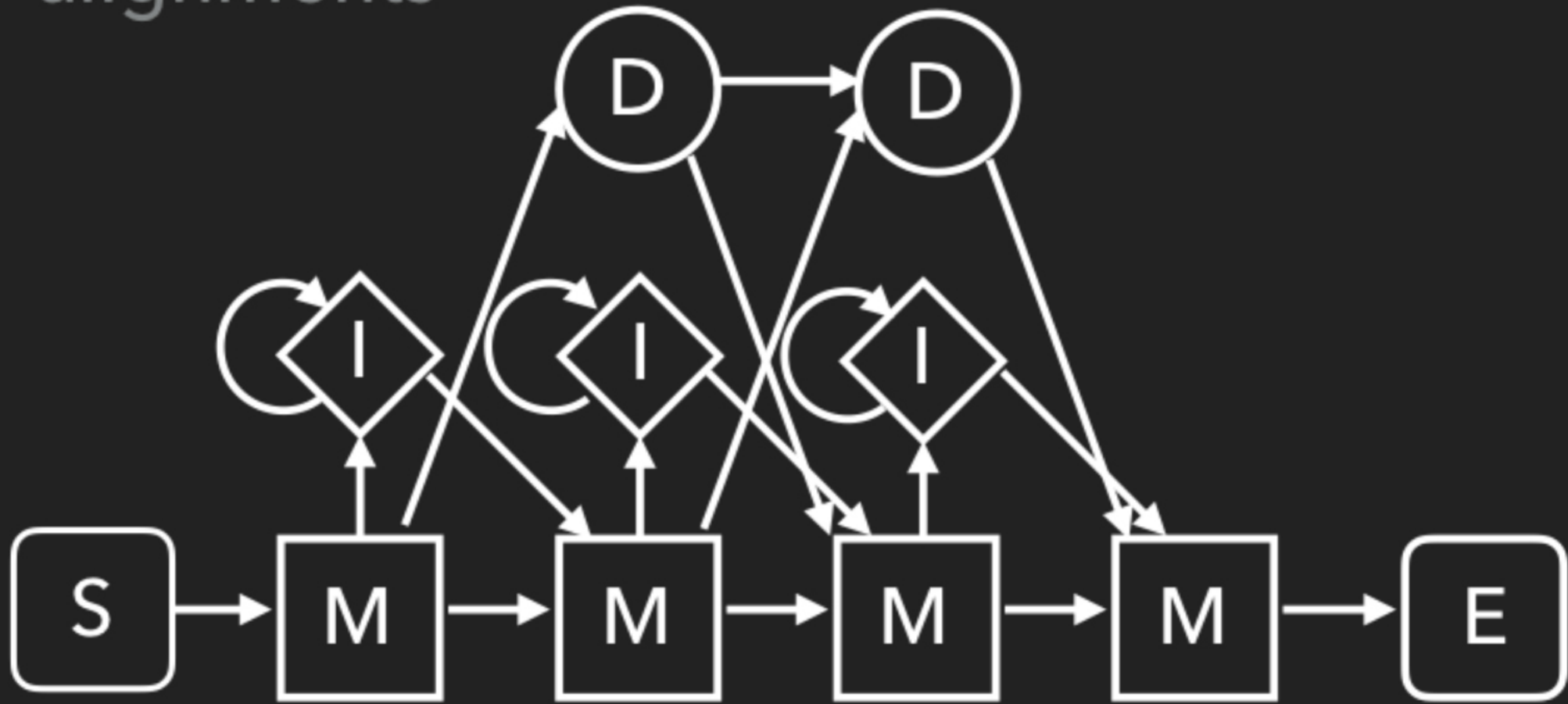
indels are hidden states

A	C	A	A	T
A	C	A		T
T	-	A		T
T	C	A		T

← insertion



we can build read profiles from read alignments



sample genomes from those profiles

A1	A	C	A	A	T	G	A	C	A	A	T	G	G	C	A	A
A2	A	C	A	T	G	A	A	C	T	G	G	C	A			
B1	T	A	T	G	A	A	A	T	G	G	C	A	A			
B2	T	C	A	T	G	A	A	C	T	G	G	C	A			
C1	T	A	T	G	A	C	A	A	T	G	G	C	A	A		
C2	T	C	A	T	G	A	A	C	T	G	G	C	A			

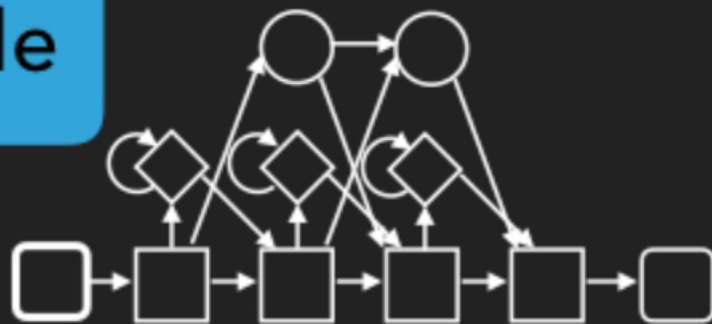
and build alignments and trees
with those samples

A1 A C A A T G A C A A T G G C A A
B1 T A T G A A A T G G C A A
C1 T A T G A C A A T G G C A A

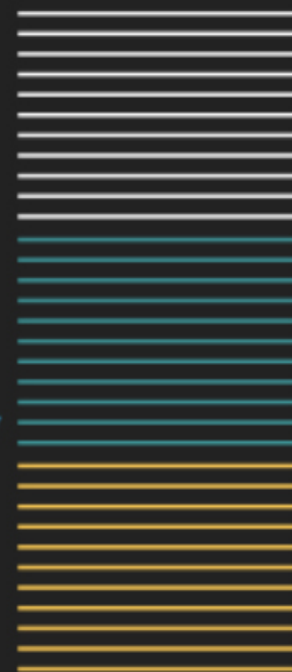
A2 A C A T G A A C T G G C A
B2 T C A T G A A C T G G C A
C2 T C A T G A A C T G G C A

our "synthetic" approach

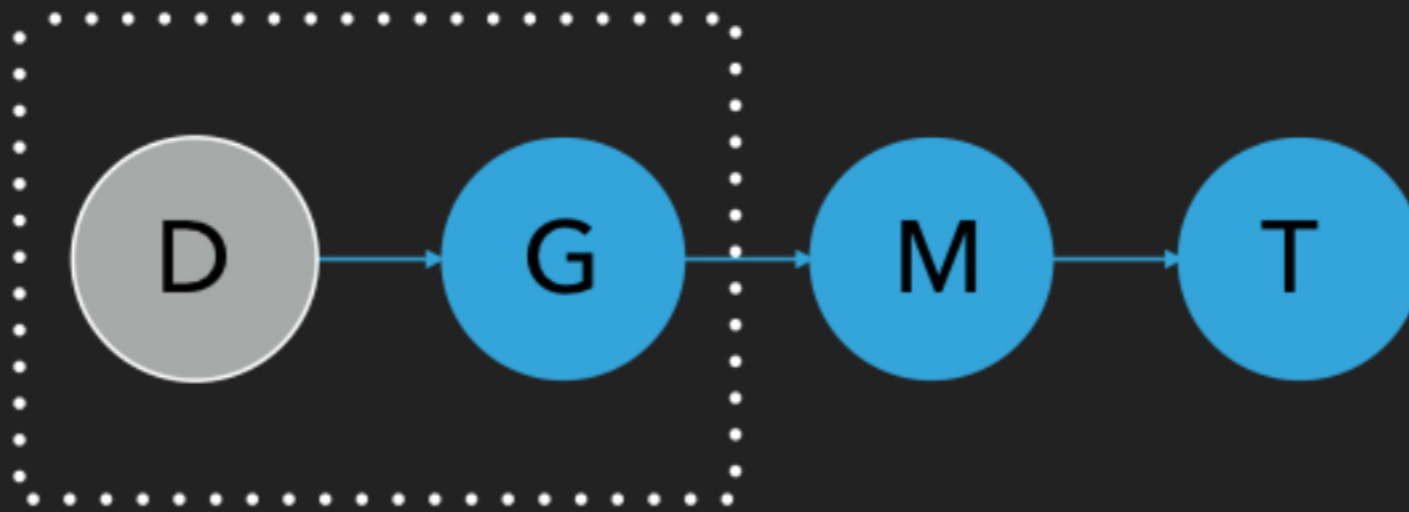
Read Profile



simulate n sequences



each step in this workflow is a high-dimensional inference problem...

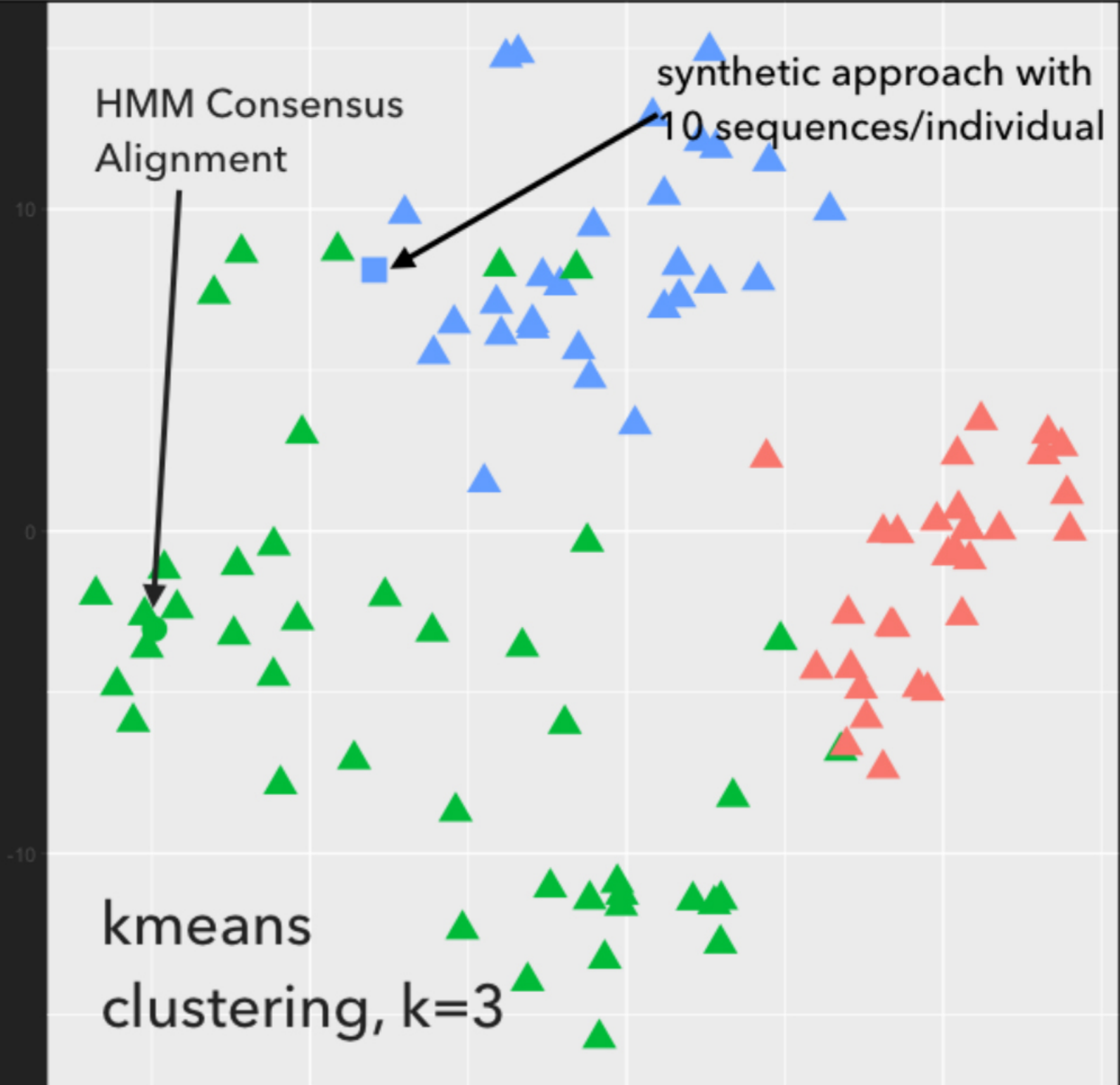


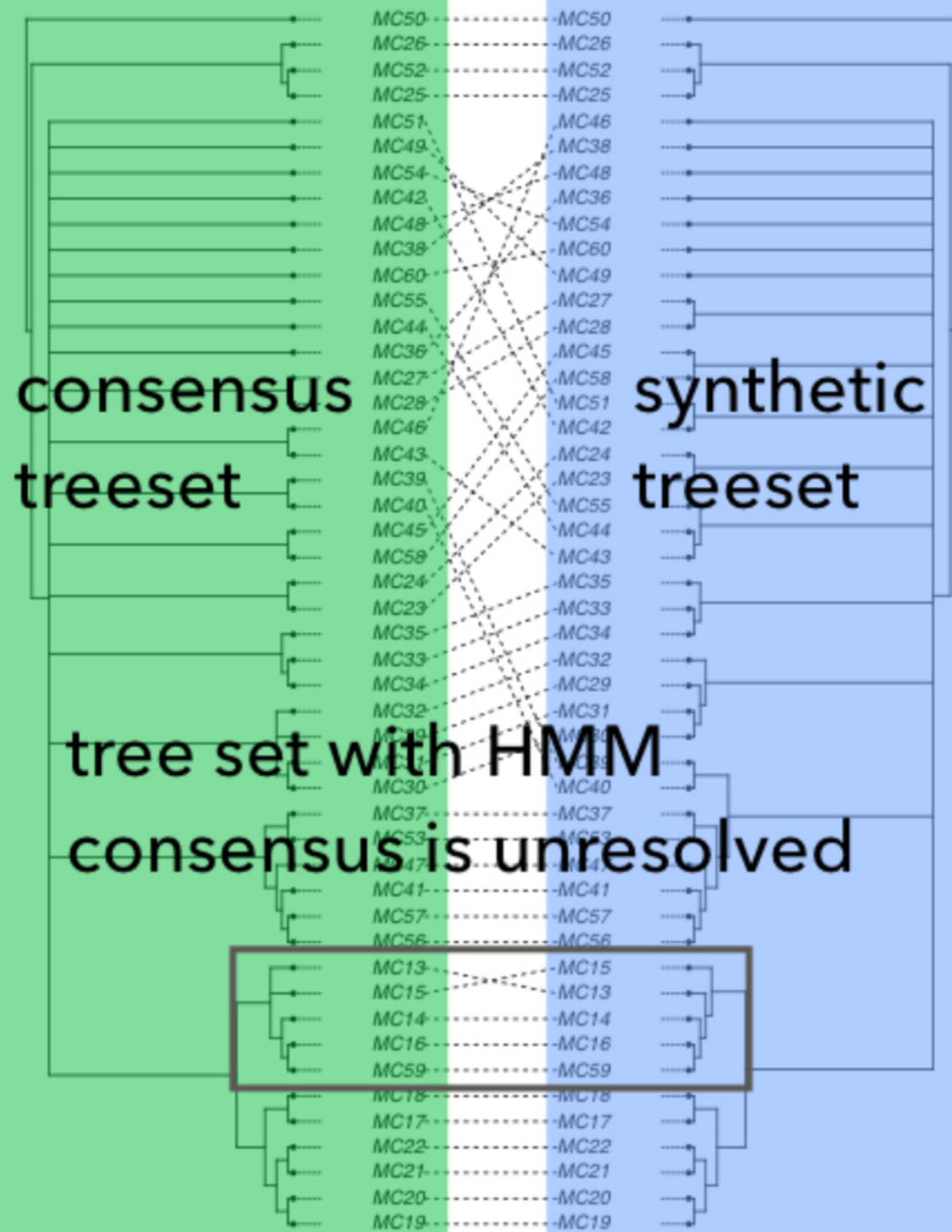
- ▶ D->G: HMMer
- ▶ G->M: mafft
- ▶ M->T: RAxML or MrBayes

BLUE WATERS
SUSTAINED PETASCALE COMPUTING

HIV DATASET

- ▶ Individuals newly-diagnosed with HIV in 2013
- ▶ Knew transmission history for 5 individuals
- ▶ Ran consensus approach; synthetic approach with 10 sequences/individual (collapsed tree); 100 runs of synthetic approach with 1 sequence/individual
- ▶ Computed Robinson-Foulds distance between trees from all approaches and performed Multidimensional Scaling





**TRANSMISSIM (SIMULATED
TRANSMISSION NETWORKS,
PHYLOGENIES, GENOMES, READS)**

a generative model of patient reads from transmission events



- ▶ N: Transmission network
- ▶ T: Transmission tree
- ▶ V: Viral phylogeny
- ▶ G: Genomes
- ▶ D: Data: reads

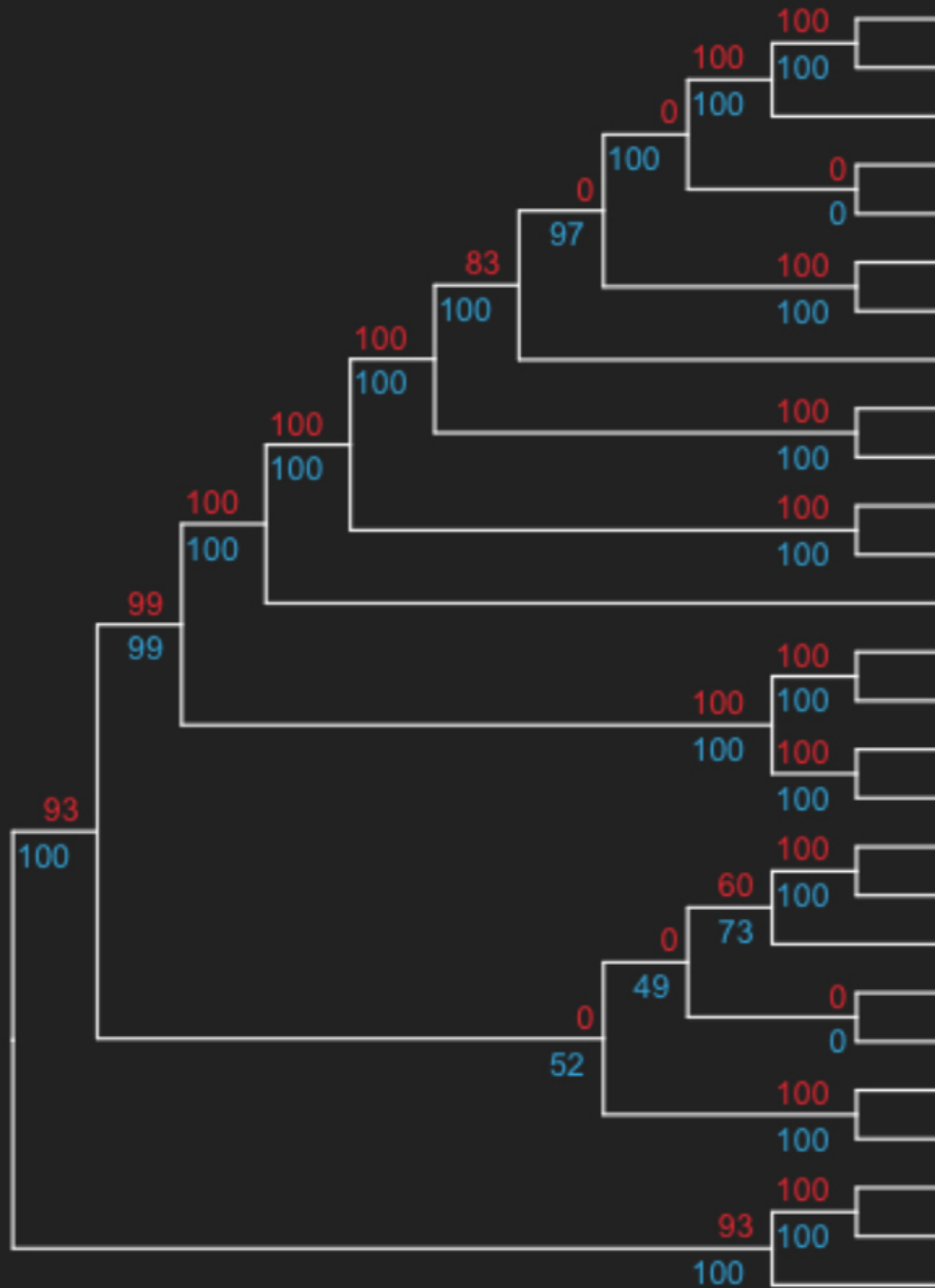
a generative model of patient reads from transmission events



- ▶ N: outbreaker
- ▶ T: binary mapping
- ▶ V: SimPhy
- ▶ G: pyvolve
- ▶ D: ART

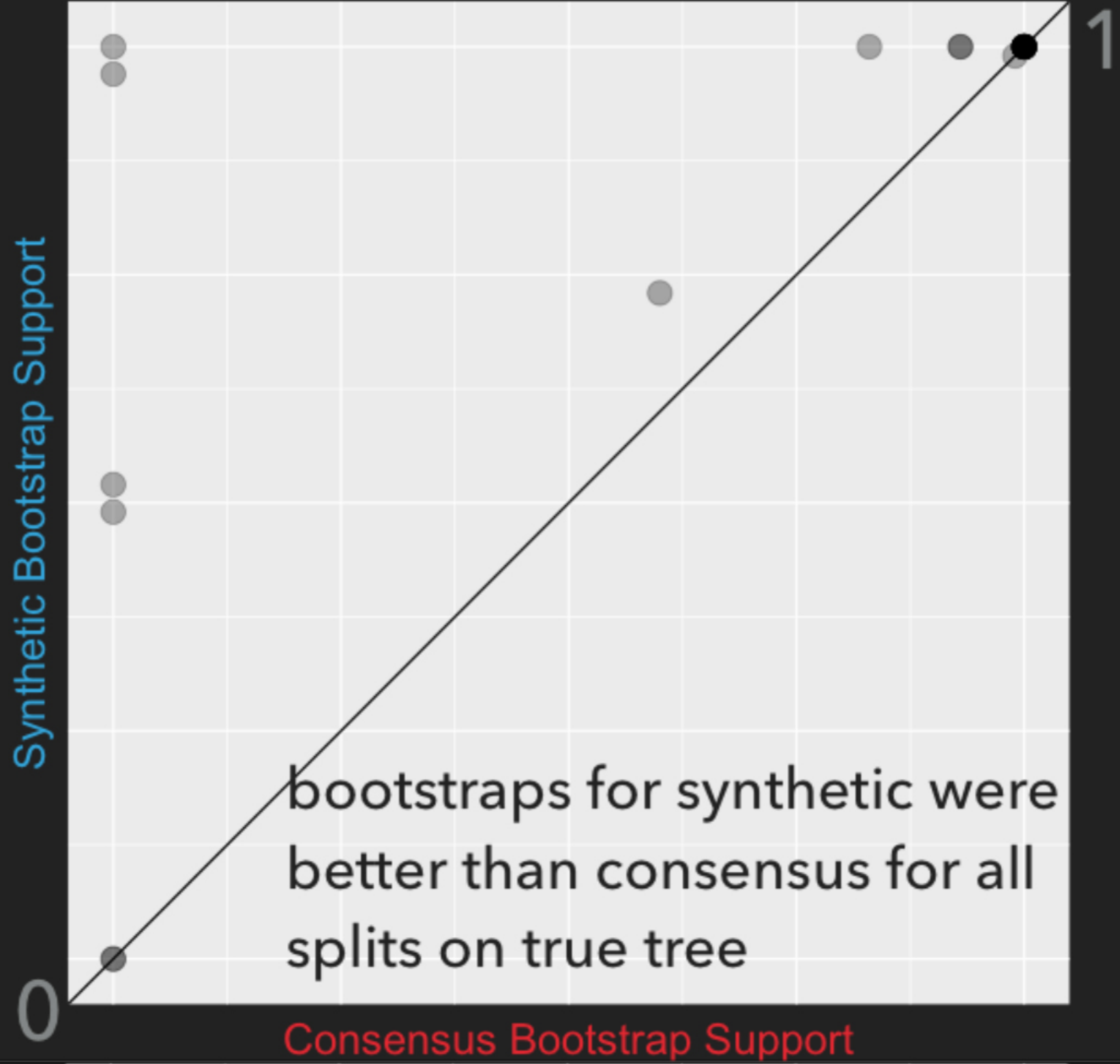
BLUE WATERS
SUSTAINED PETASCALE COMPUTING

bootstrap results on simulations

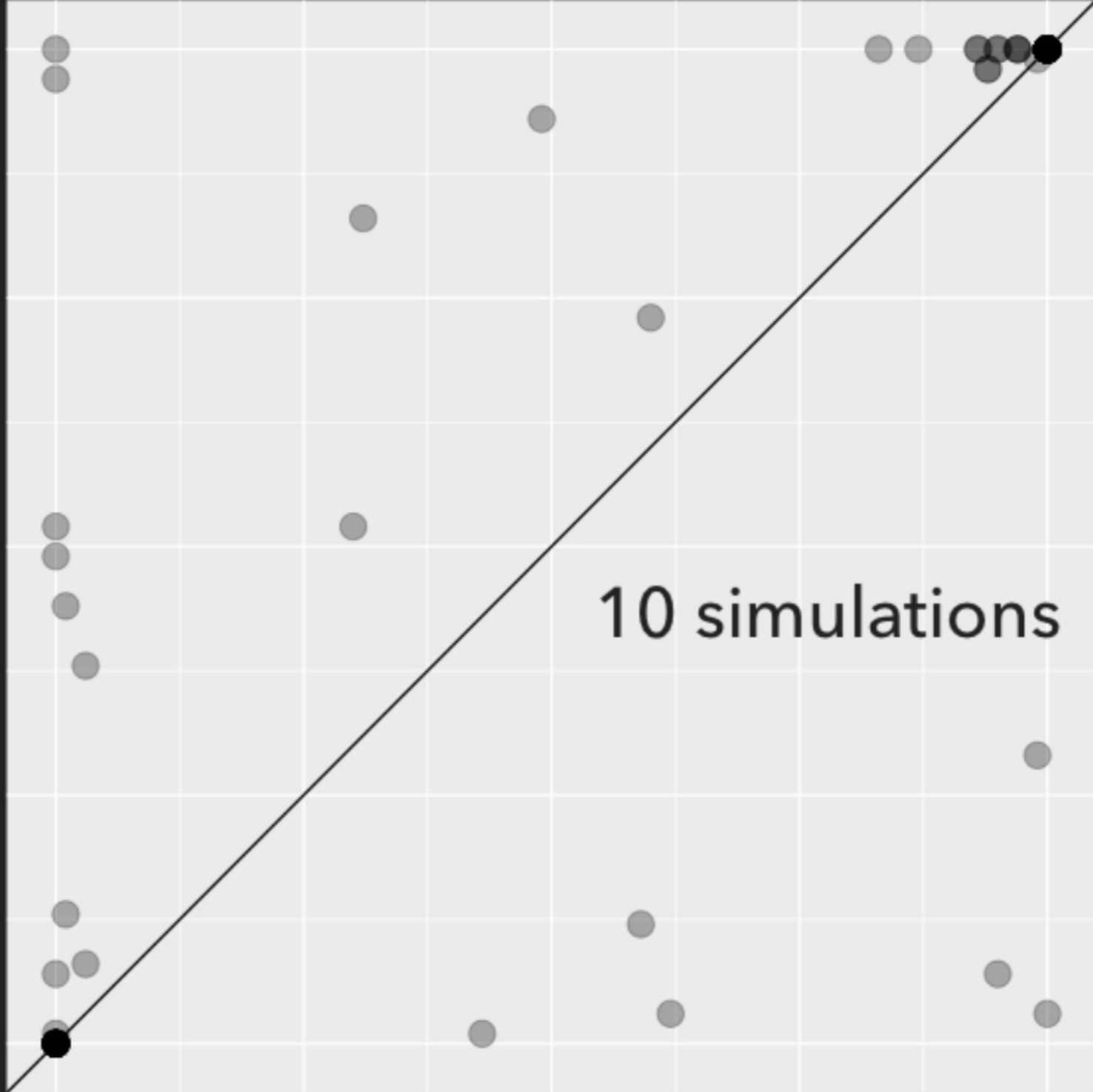


red = consensus

blue = synthetic



Synthetic Bootstrap Support



10 simulations

Consensus Bootstrap Support

WHY BLUE WATERS?

1. 10,000 node hours: Create MDS density plot with 10,000 trees to look for conclusive regions of variation
2. remaining node hours: additional simulations for validation with parameter sweeps

BLUE WATERS PRODUCTS

- ▶ Reproducible manuscript and figures: https://bitbucket.org/aguang/ms_hiv
- ▶ Transmissim: <https://github.com/aguang/transmissim>

acknowledgments

My CFAR Research Group:

Casey Dunn, Associate Professor of Evolutionary Biology

Rami Kantor, Associate Professor of Medicine

Mia Coetzer, Assistant Professor of Medicine

Mark Howison, Director of Data Science

Colin MacLean, Research Programmer & Charles Lawrence, Professor of Applied Mathematics

My DunnLab Research Group:

Casey Dunn (again)

Zachary Lewis, Postdoc

Catriona Munro, PhD Candidate

Alex Damian Serrano, PhD Candidate



a. Whole-genome NGS reads for Patient 1

Read1 ATGGCATATGGAGCATGATGGC
Read2 TGATGCATCGCTGATGCCATAT
Read3 TGGATGCATCGCTGATGGCATA

b. HMMER alignment to reference pHMM

Reference1 TGGATGCATCGCTGATGGCATATGTGATGGCATATT *Los Alamos*
Reference1 TG-ATGCATCGCTGATGGCTTATGGGATGGCAT--- *HIV Sequence Database*

Read1	-----	-----
Read2	TG-ATGCATCGCTGATGCCATAT	-----
Read3	TGGATGCATCGCTGATGGCATA	-----

→ Patient 1 pHMM

c. HMMER re-alignment to Patient 1 pHMM

Read1	-----ATGGCATATGGAGCATGATGGC	<i>Additional sensitivity in re-alignment</i>
Read2	TG-ATGCATCGCTGATGCCATAT	-----
Read3	TGGATGCATCGCTGATGGCATA	-----

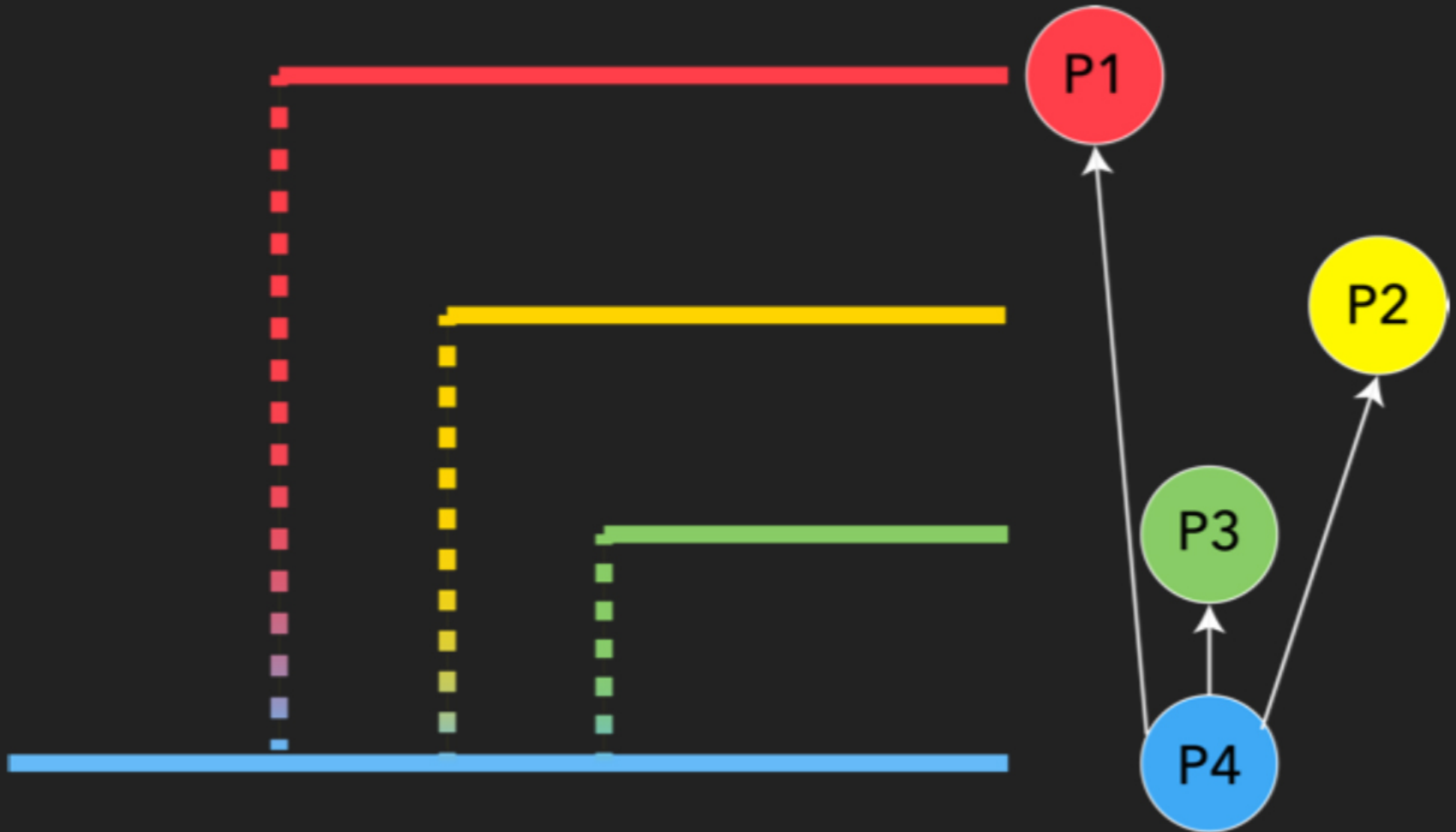
d. Summary sequences for Patient 1 pHMM

Consensus TG-ATGCATCGCTGATGGCATAT----- *Majority-rule*
Synthetic1 T--ATGCATCGCTGATGGCATAT--A-----A---C *Sampled from pHMM*
Synthetic2 -GGATG--TCGCTGATGCCATAT----C-TGA-----

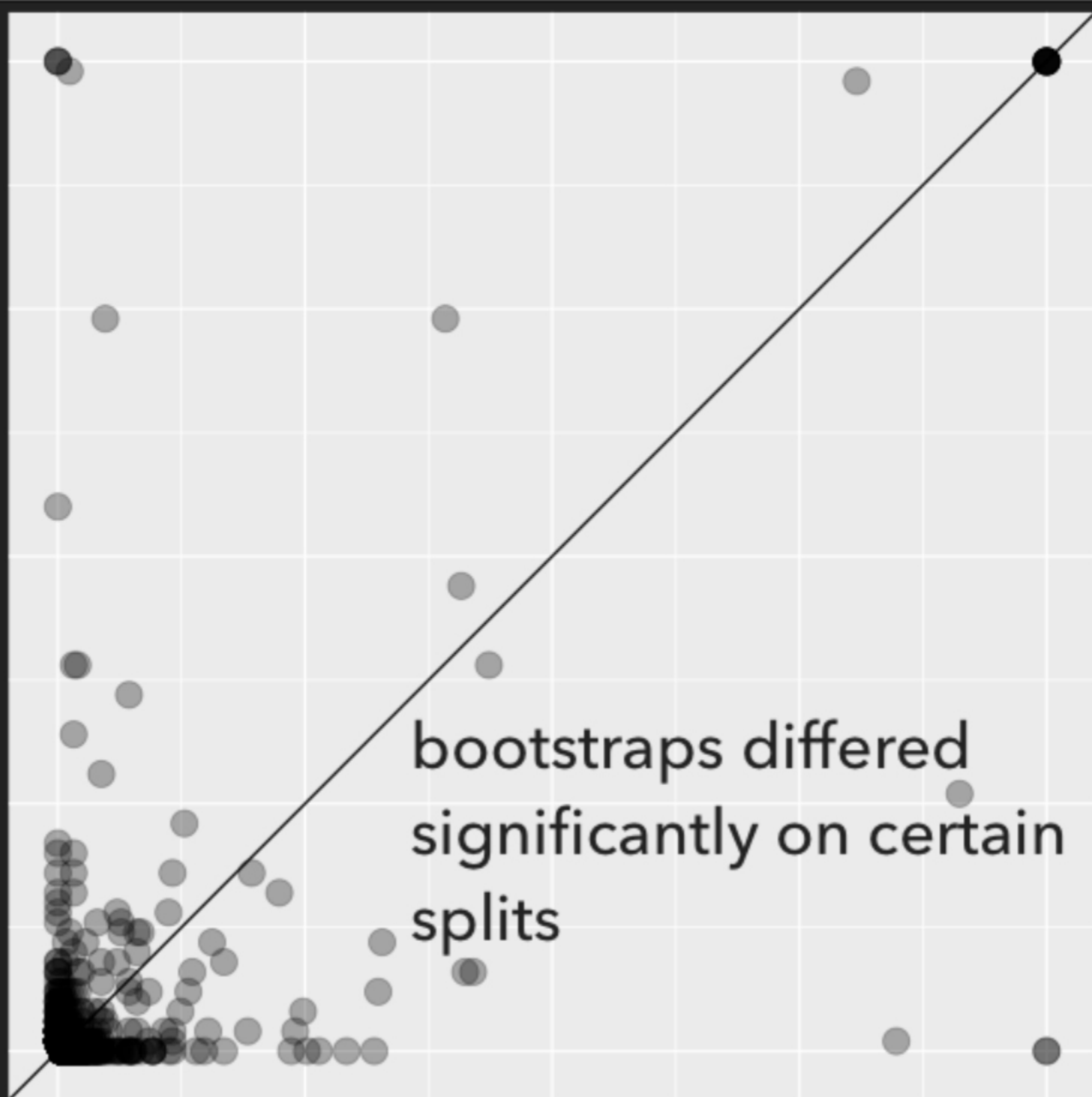
under certain assumptions, transmission networks have a surjective mapping to a phylogeny (transmission tree)



however, the mapping is not injective
and thus not one-to-one



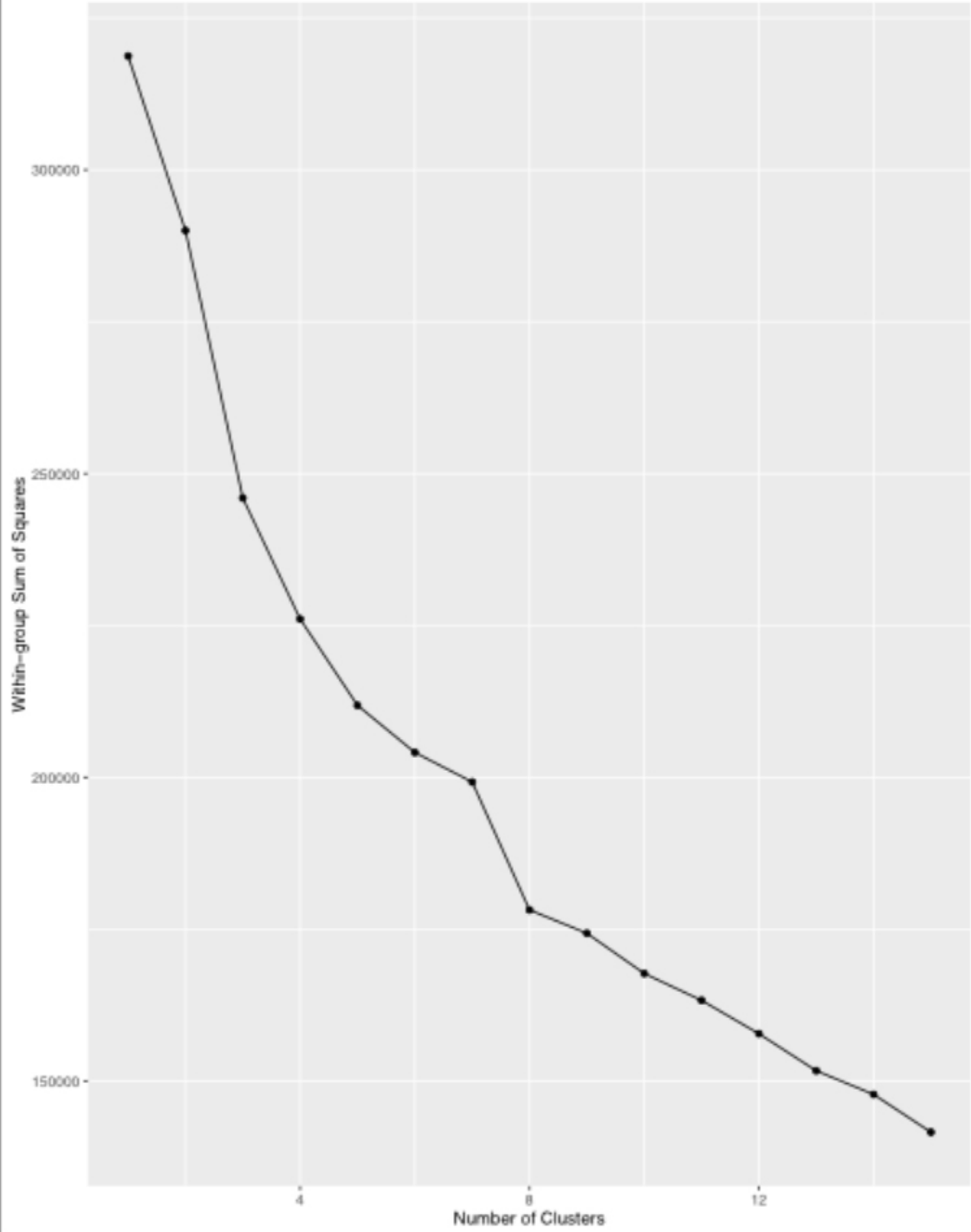
Synthetic Bootstrap Support



bootstraps differed significantly on certain splits

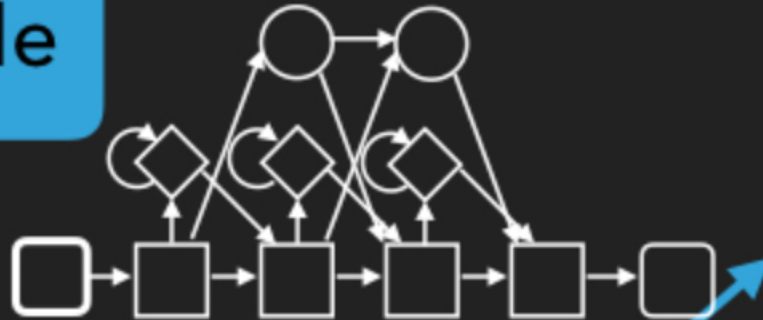
Consensus Bootstrap Support

Elbow plot of K award sampled trees



our "synthetic" approach

Read Profile



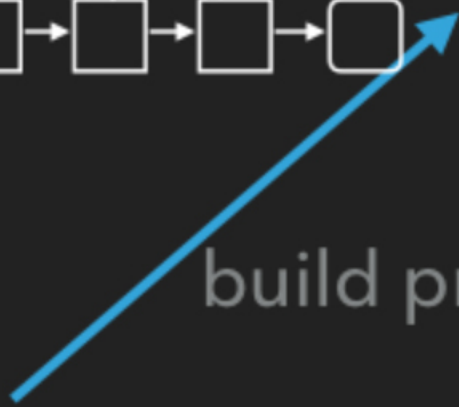
align reads



Reads



build profile



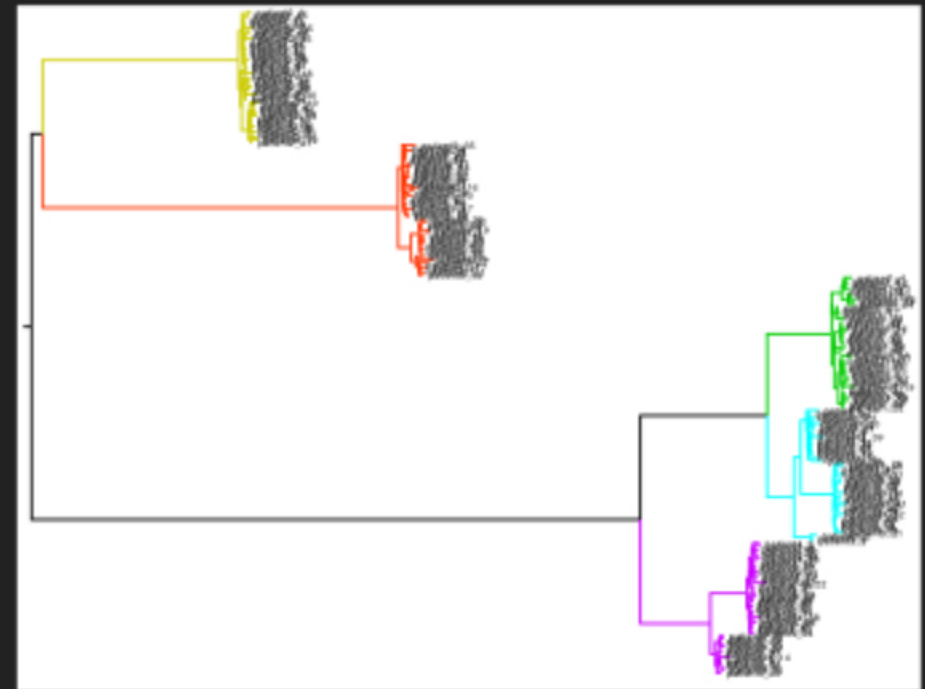
realign reads



some simple simulations



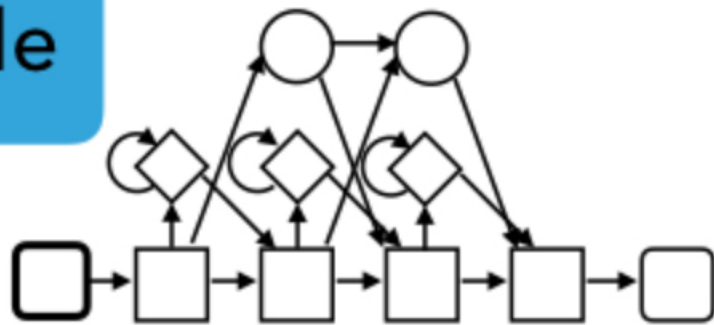
RAxML tree from consensus genome



true tree

so we don't recover the same tree...

Read Profile



simulate n sequences



Kmeans clustering of ML trees from synthetic sequences and ML tree from consensus alignment

