

Statistical/computational phylogenetics

Alexandre Bouchard-Côté

Department of Statistics,
University of British Columbia

Organization

- Background in phylogenetics
- Non-standard application areas:
 - Clonal evolution inside an individual cancer tumour
 - Reconstructing ancient languages
- Current research on Divide-and-Conquer Sequential Monte Carlo methods

Background on phylogenetics

Notation/background on trees

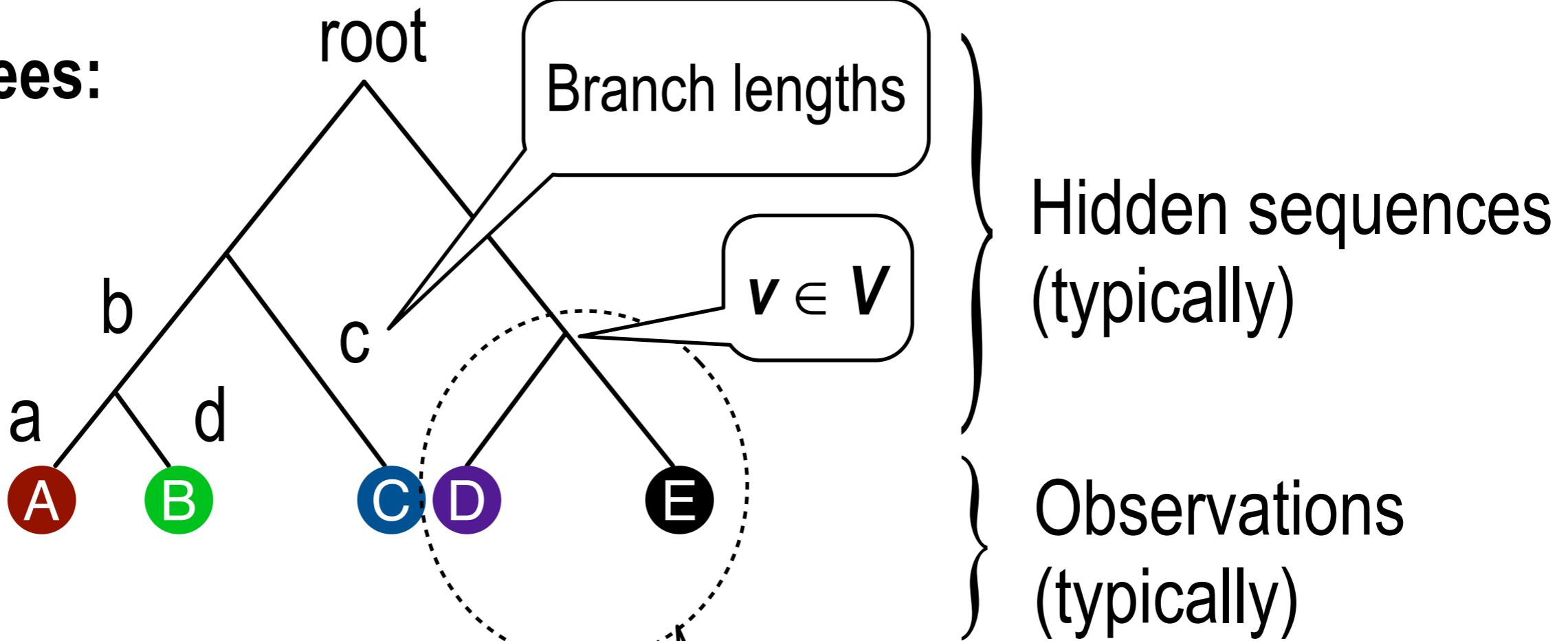
Clock trees:

$$a+b = c$$

$$d+b = c$$

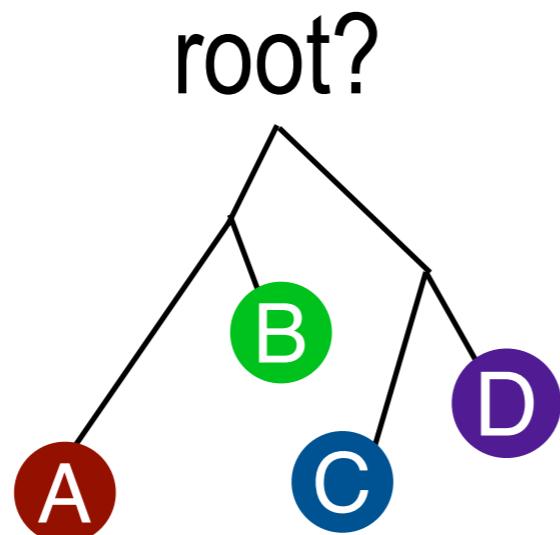
$$a = d$$

...



Non-clock tree:

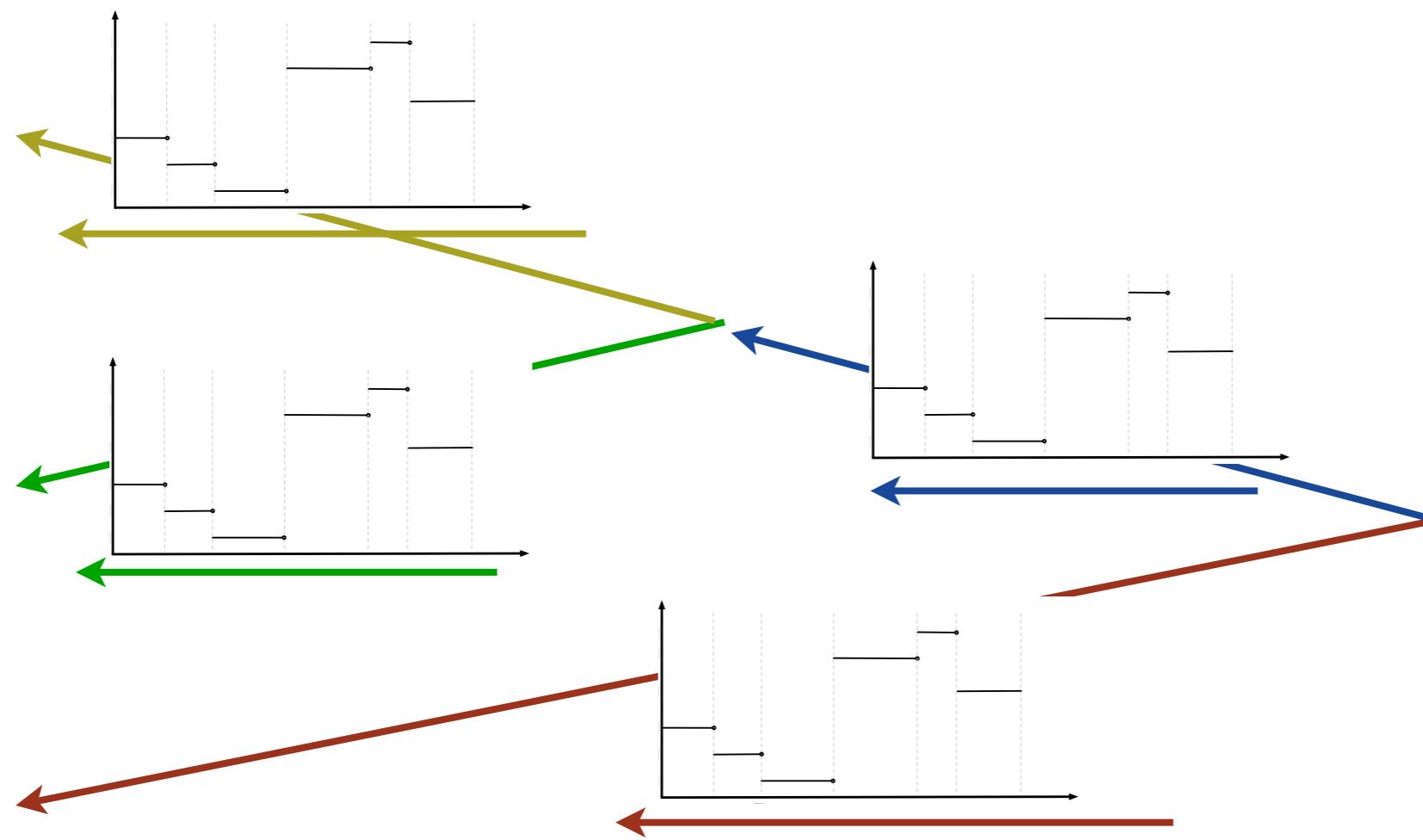
remove additivity
restrictions on
branch lengths



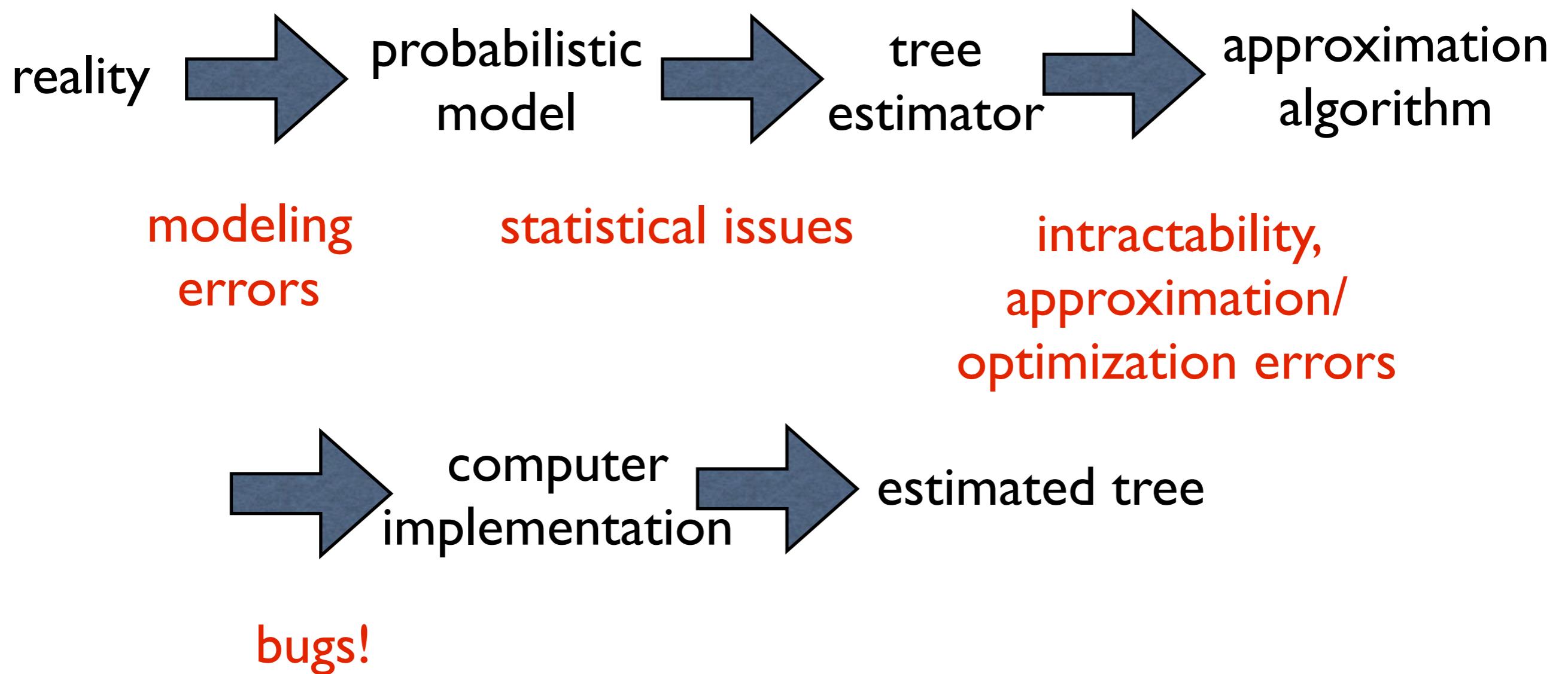
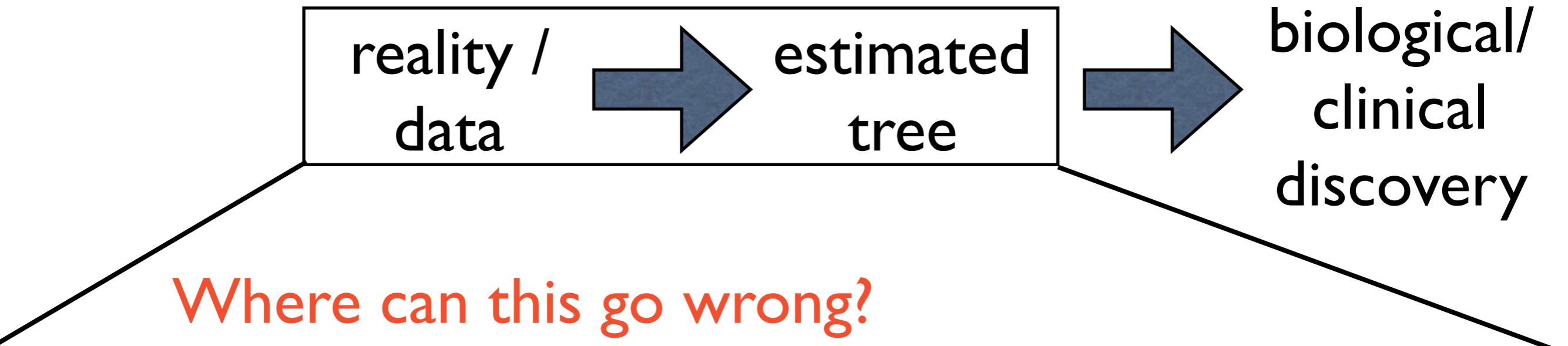
Clade: group containing a taxa and all its descendants. (I will restrict them to leaves, e.g. {D,E})

Stochastic process on a tree

- Evolution modeled by a stochastic process
 - Discrete data: Continuous time Markov chains (CTMC)
 - Continuous data: Diffusions

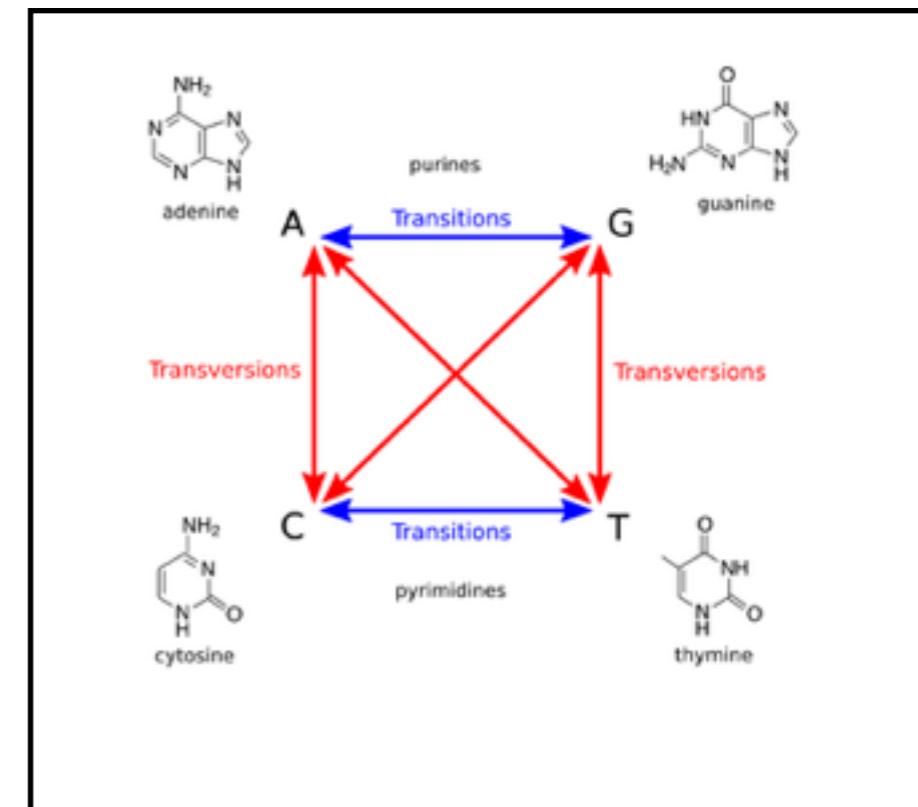


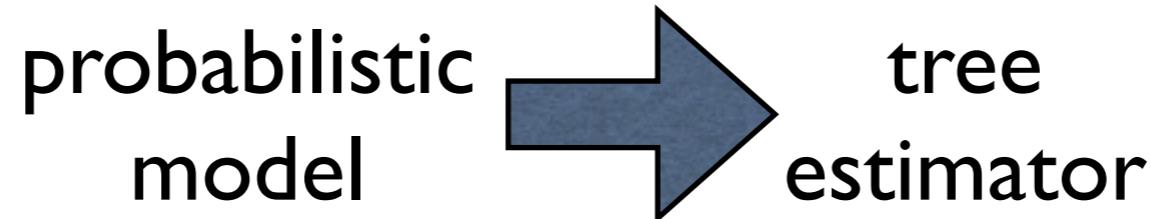
Phylogenetics: challenges



reality → probabilistic model

- Common assumption: we have replicates (loci/sites) of the evolutionary stochastic process (often iid given tree)

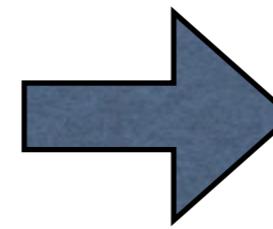




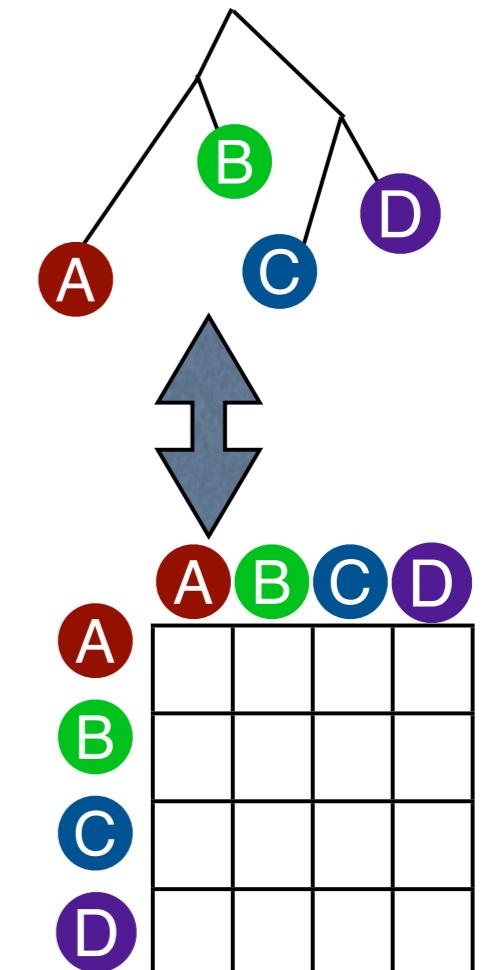
statistical issues

- Identifiability/consistency: can we find the tree given ‘infinite data’
 - What does infinite data mean?
 - # loci/sites
 - # leaves
 - sequencing depth
- Efficiency: how fast will the error decay as we add more data? (cf. computational efficiency)

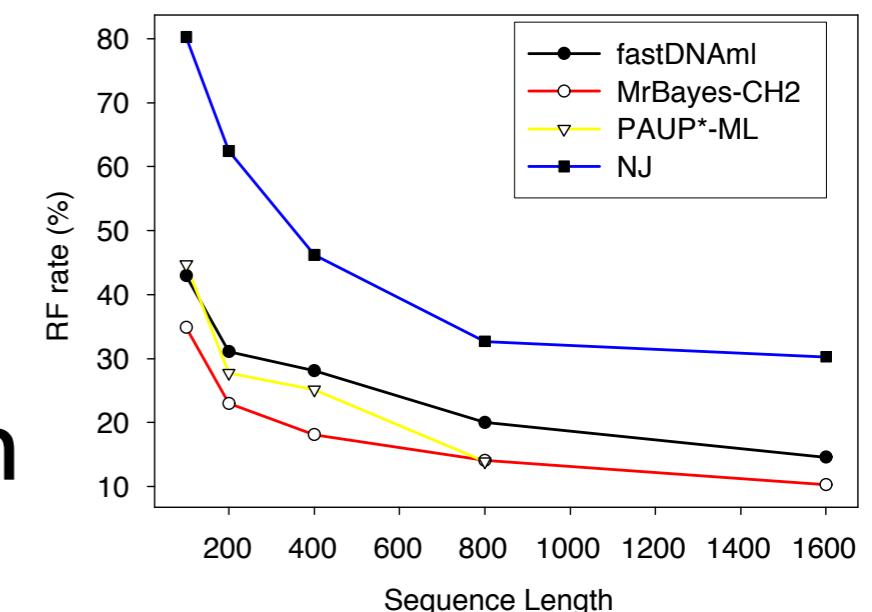
probabilistic
model



tree
estimator

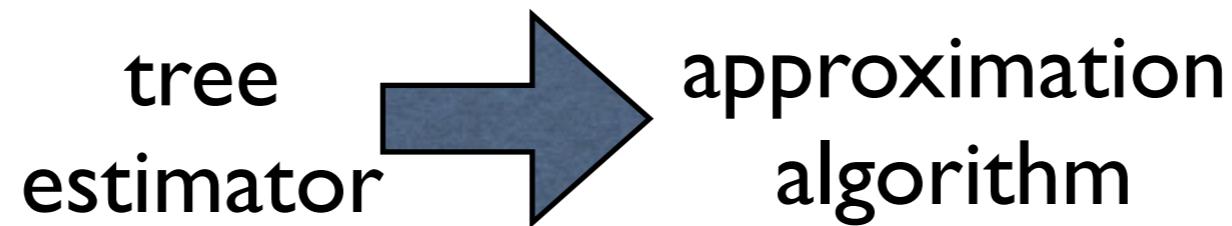


- Earlier methods: parsimony, distance-based
- State-of-the-art: likelihood-based (maximum likelihood and Bayesian method)
- key idea: score many hypothetical trees, marginalizing over unknown ancestral states



Taxa = 40, Height = 1.0

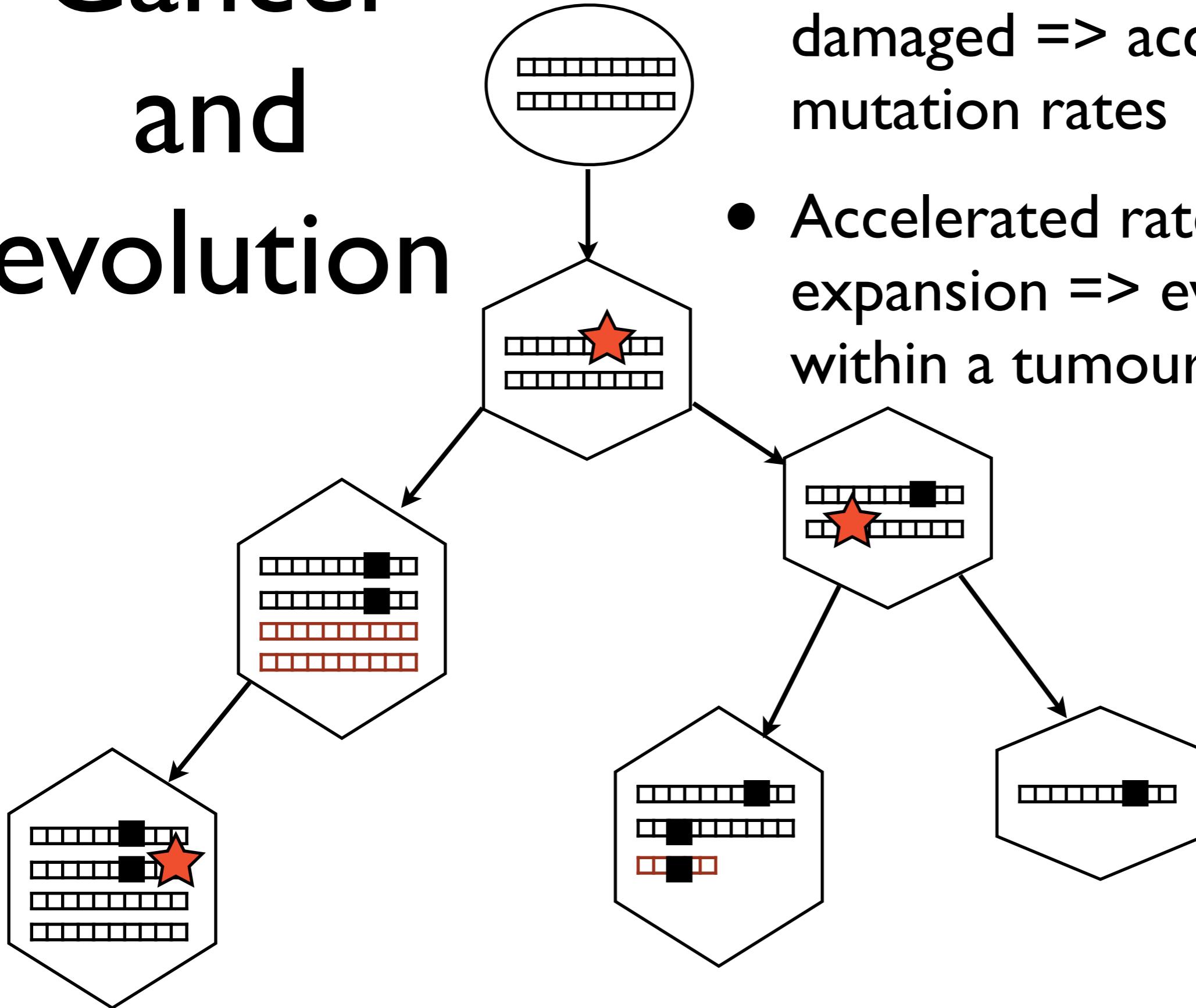
Williams and
Moret, 03



- Statistically efficient tree reconstruction methods are generally computationally inefficient
- Parallelization, distribution
- Approximation methods:
 - MCMC
 - Sequential Monte Carlo
 - Annealing

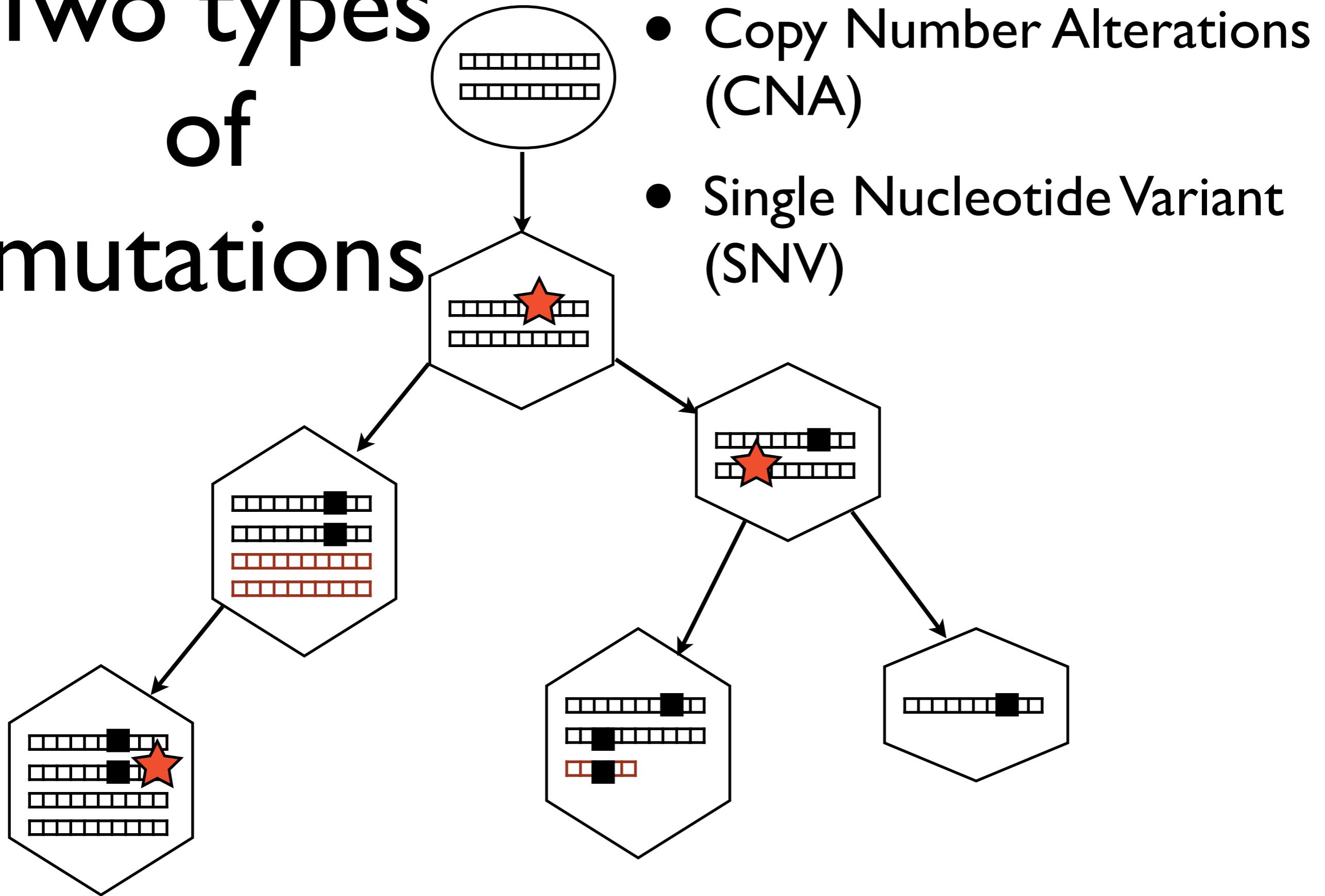
Application: Inferring the Phylogeny of Cancer Cells

Cancer and evolution



- DNA repair pathways damaged => accelerated mutation rates
- Accelerated rate + clonal expansion => evolution within a tumour

Two types of mutations

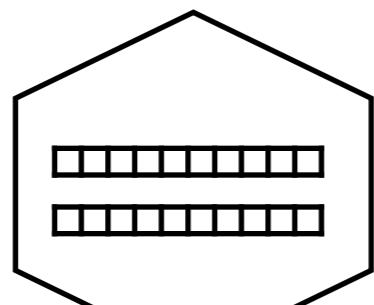


- Copy Number Alterations (CNA)
- Single Nucleotide Variant (SNV)

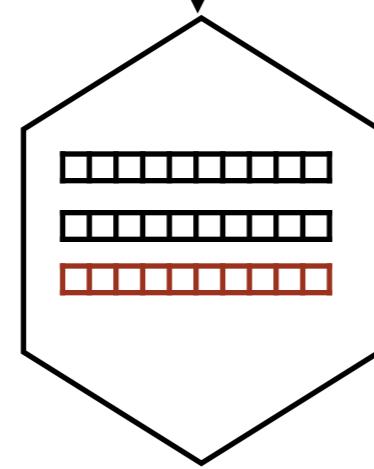
Copy number alterations

Whole
chromosome

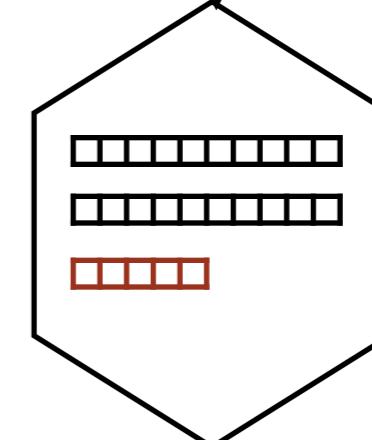
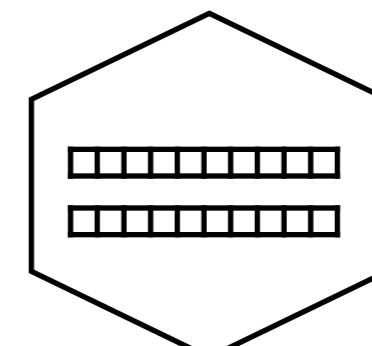
Parent
cancer
cell



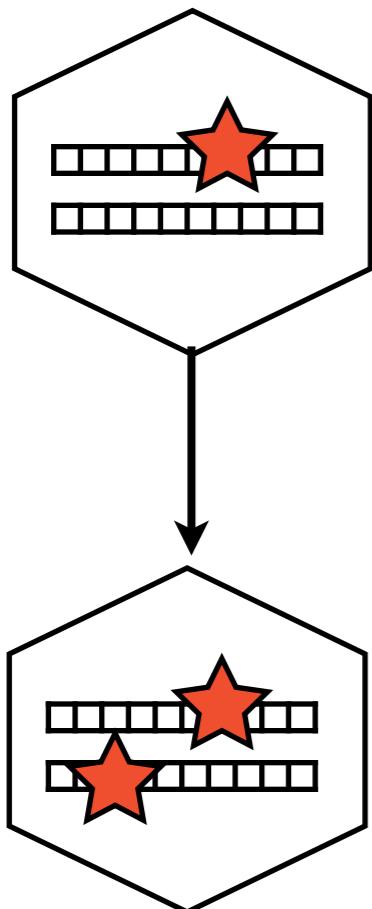
Child
cancer
cell



Local aberration

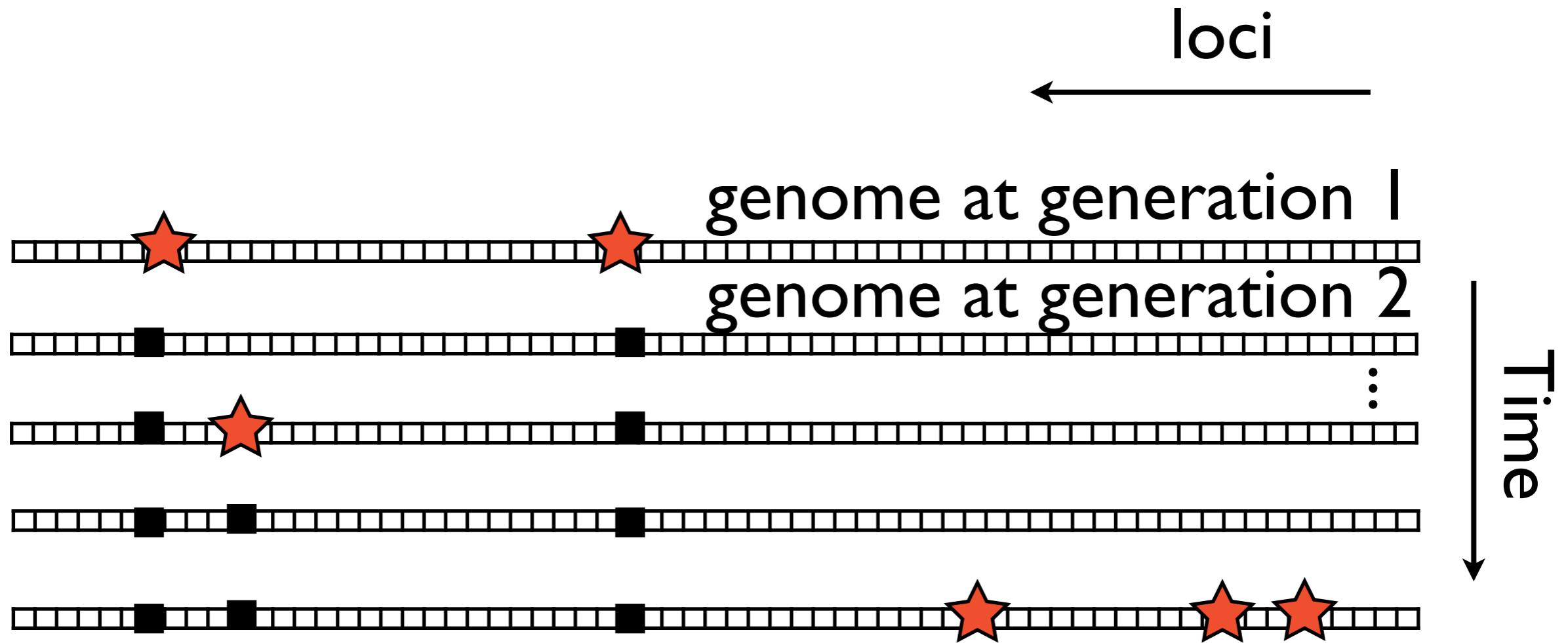


Single nucleotide variants (SNVs)



- Like SNPs, but arise in tumour instead of already in healthy cells (somatic instead of germline)

Accumulation of SNVs



Mathematical simplification: at most one single nucleotide hit per locus ('infinite site model', $p^2 \ll p$ when p is small)

=> binary state for SNVs

Motivations for studying this process

- Clinical:
 - Evolution => Adaptation => Relapse
- Scientific:
 - Evolution in the fitness landscape far from a local optimum

Bayesian non-parametric clonal inference

Trying to identify the ‘nodes’ of the tree

Simplification (pedagogical)

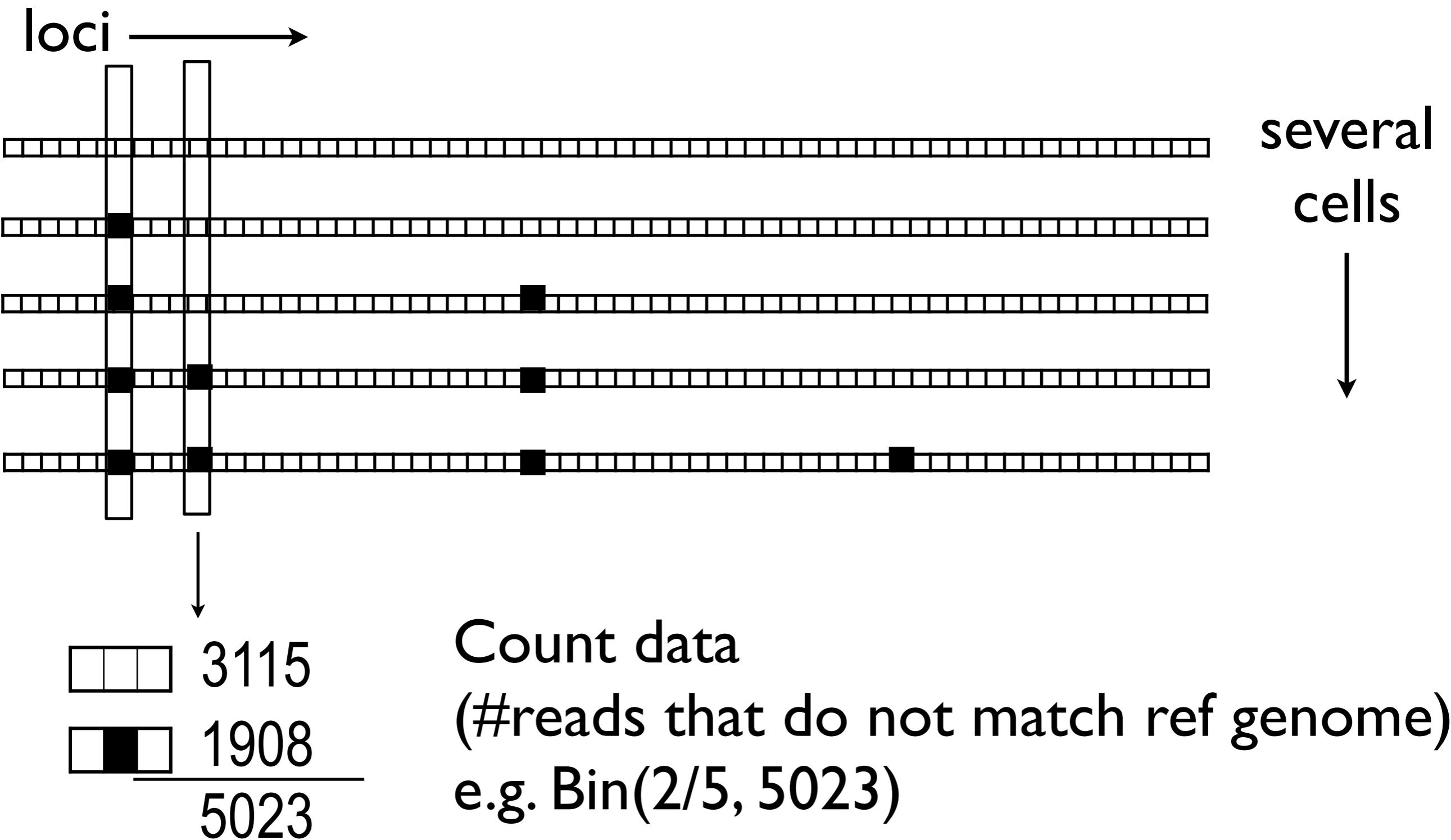
Simplification: for now, 1 chromosome, no CNA

one
cell

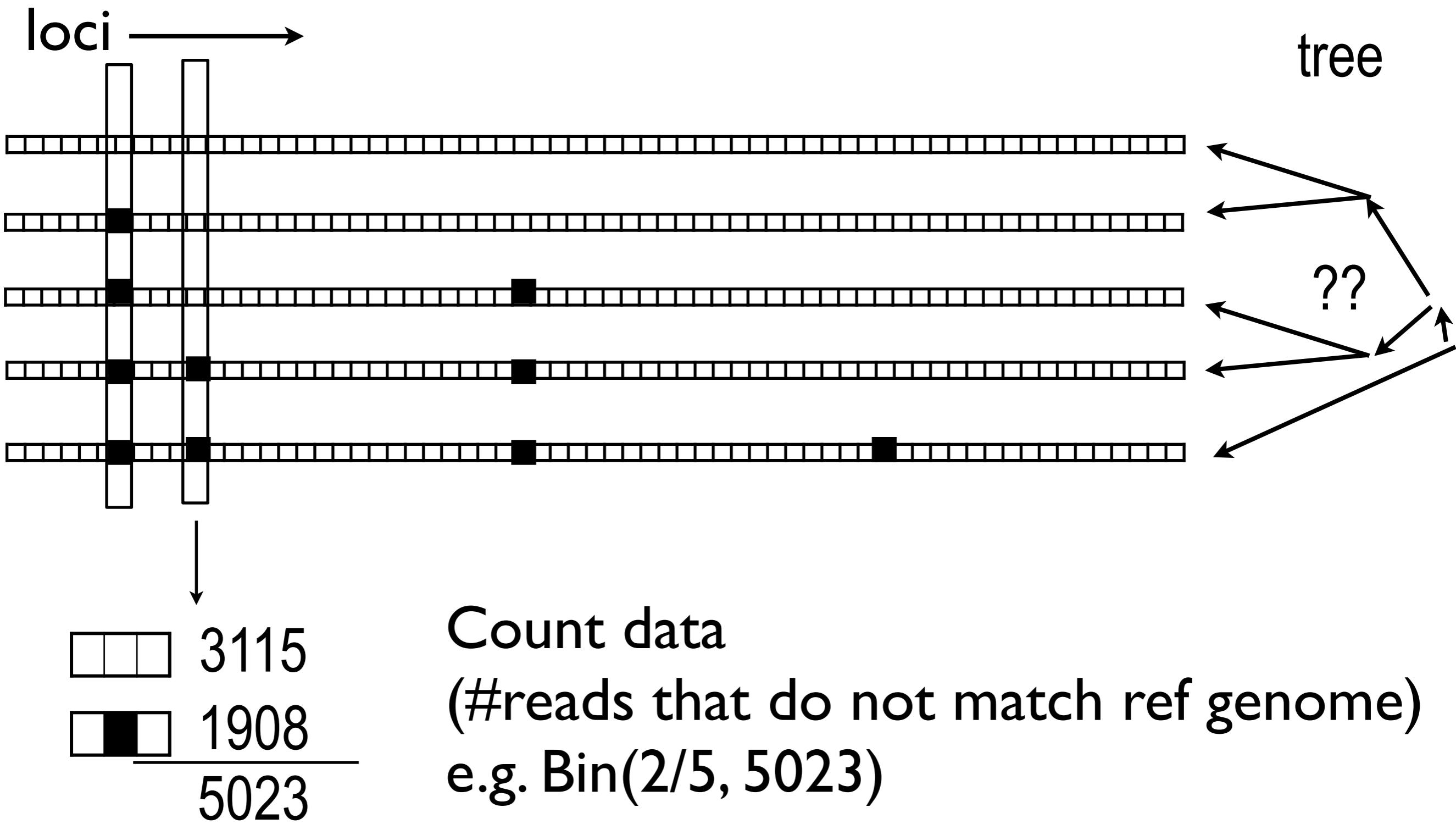


loci —→

Observation: ‘bulk data’



Observation: ‘bulk data’

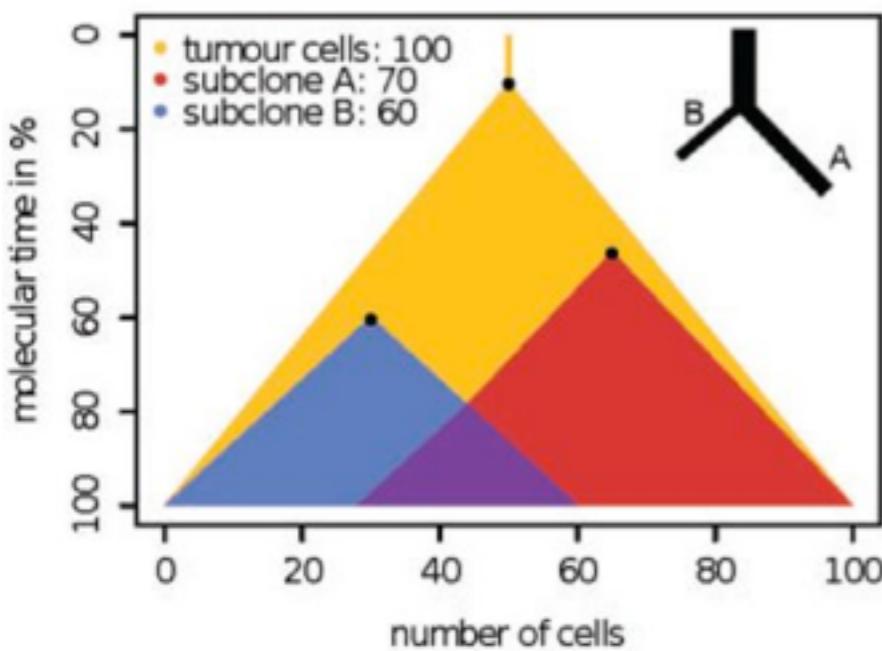


Mutation (partial) order from clonal prevalence

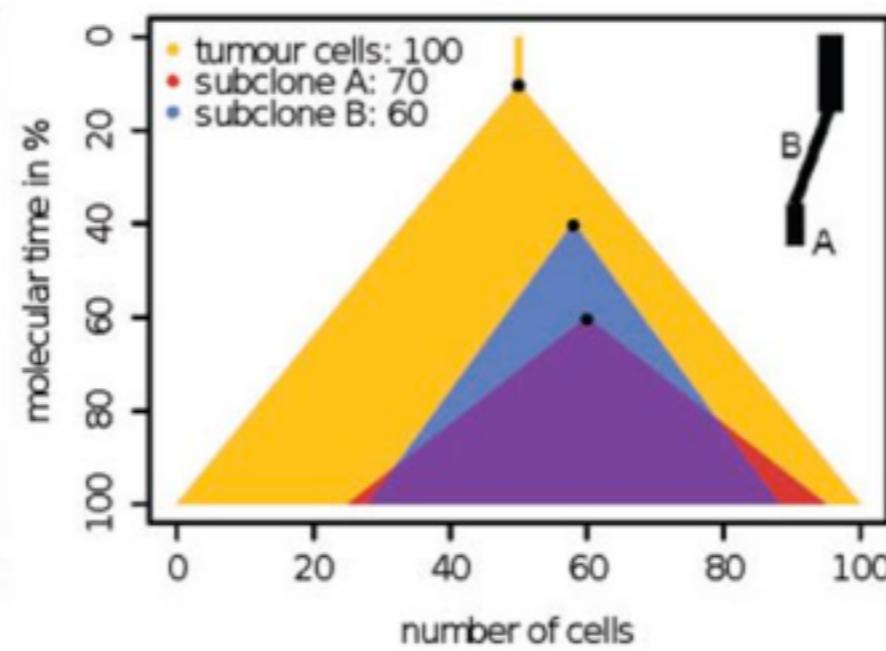
Example: one mutation in 70% of cells, another in 60%

Consequence of infinite site assumption & pigeon hole principle:

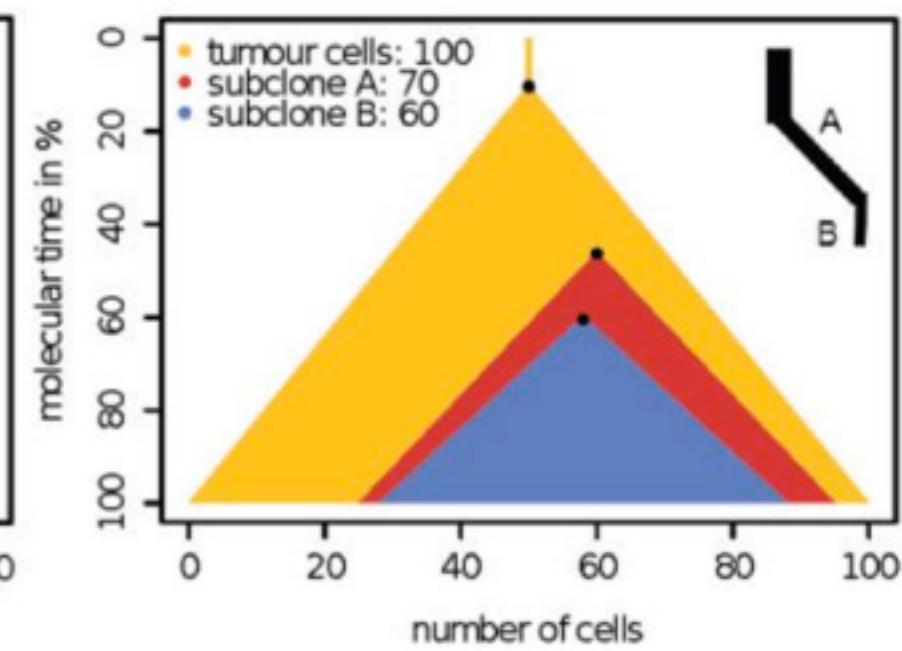
a) non feasible



b) non feasible

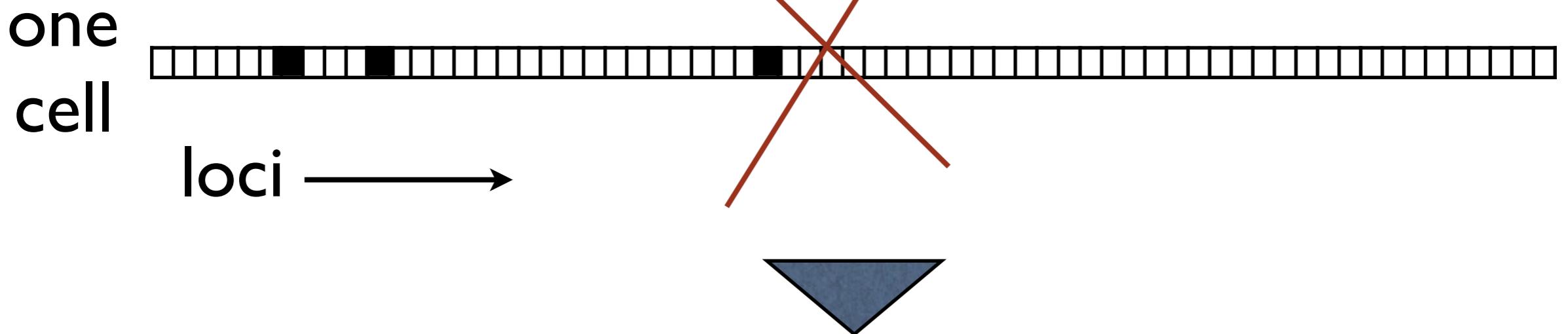


c) feasible

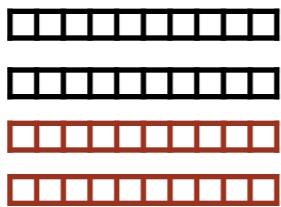


Simplification (pedagogical)

~~Simplification: for now, 1 chromosome, no CNA~~



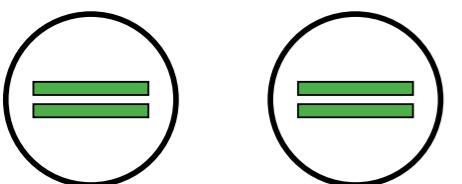
Complex, aberrant copy number profiles



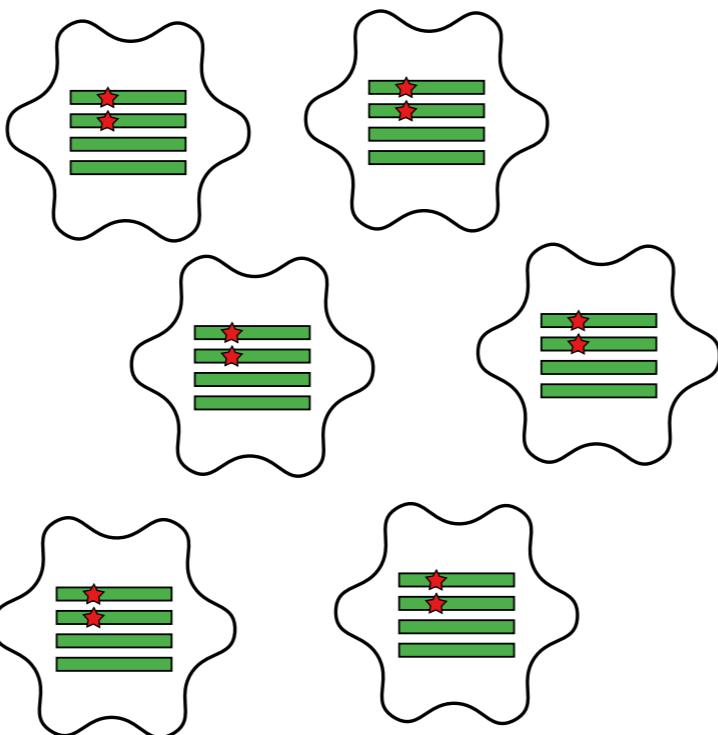
etc

Deconvolution of bulk data

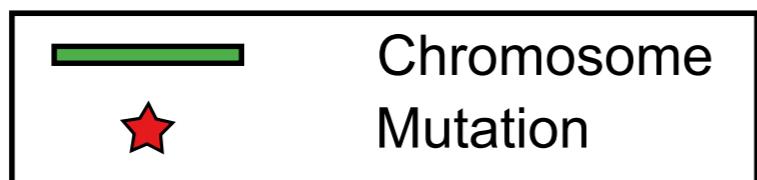
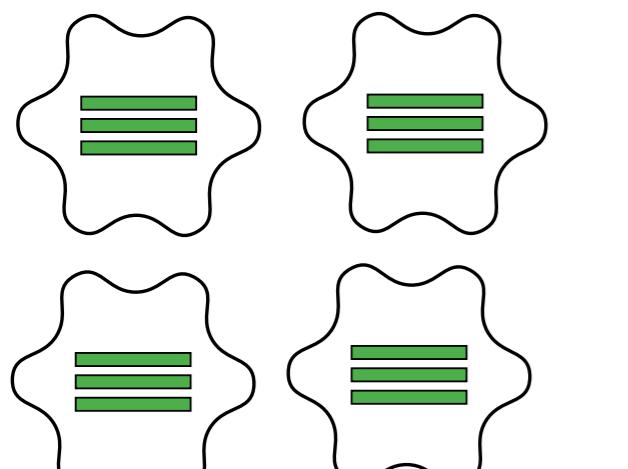
Normal Population



Variant Population



Reference Population



Variant allelic prevalence

= # variant strands
/ # strands

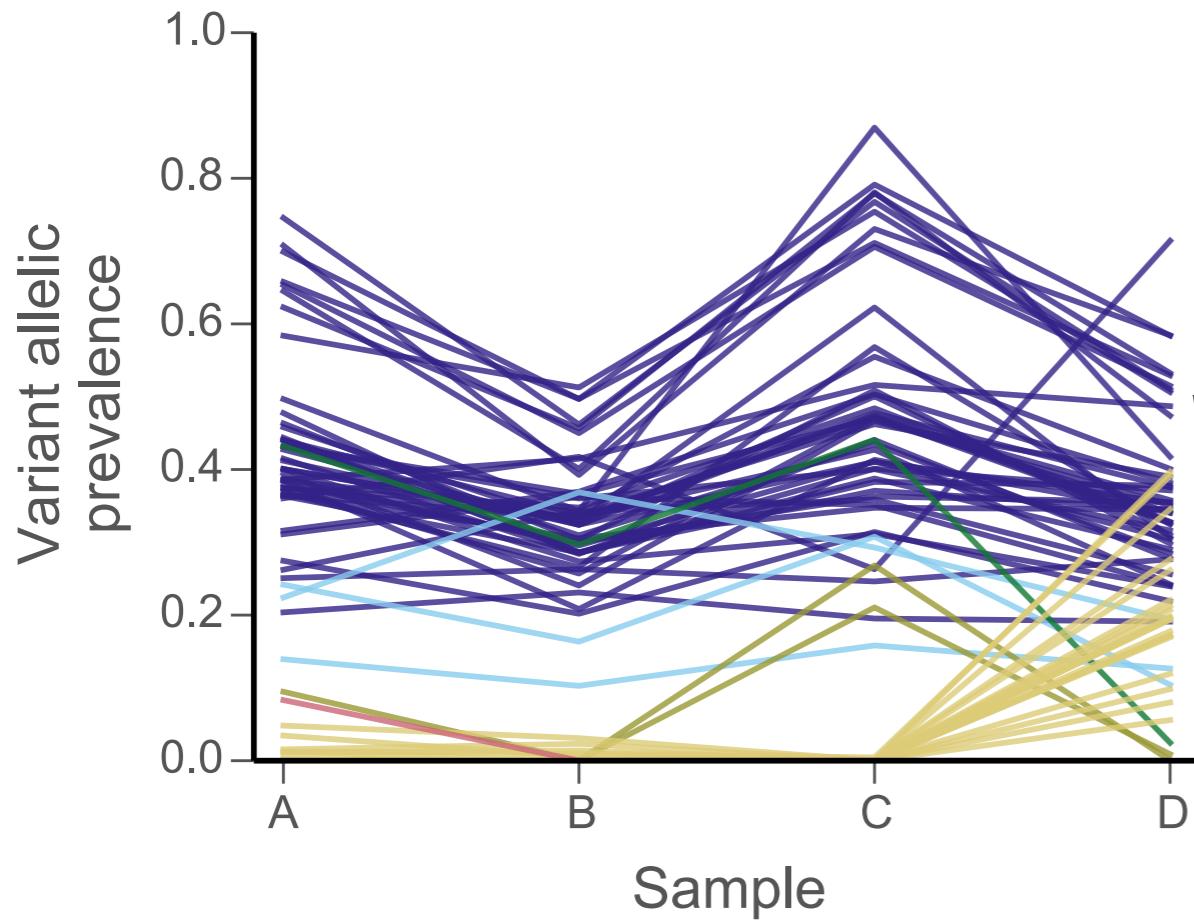
= 12/40 ('observed')

Variant cellular prevalence

= # variant cells / # cells

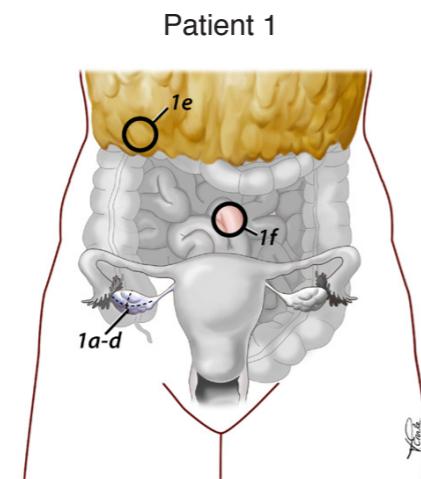
= 6/12 (inferred)

Example of data

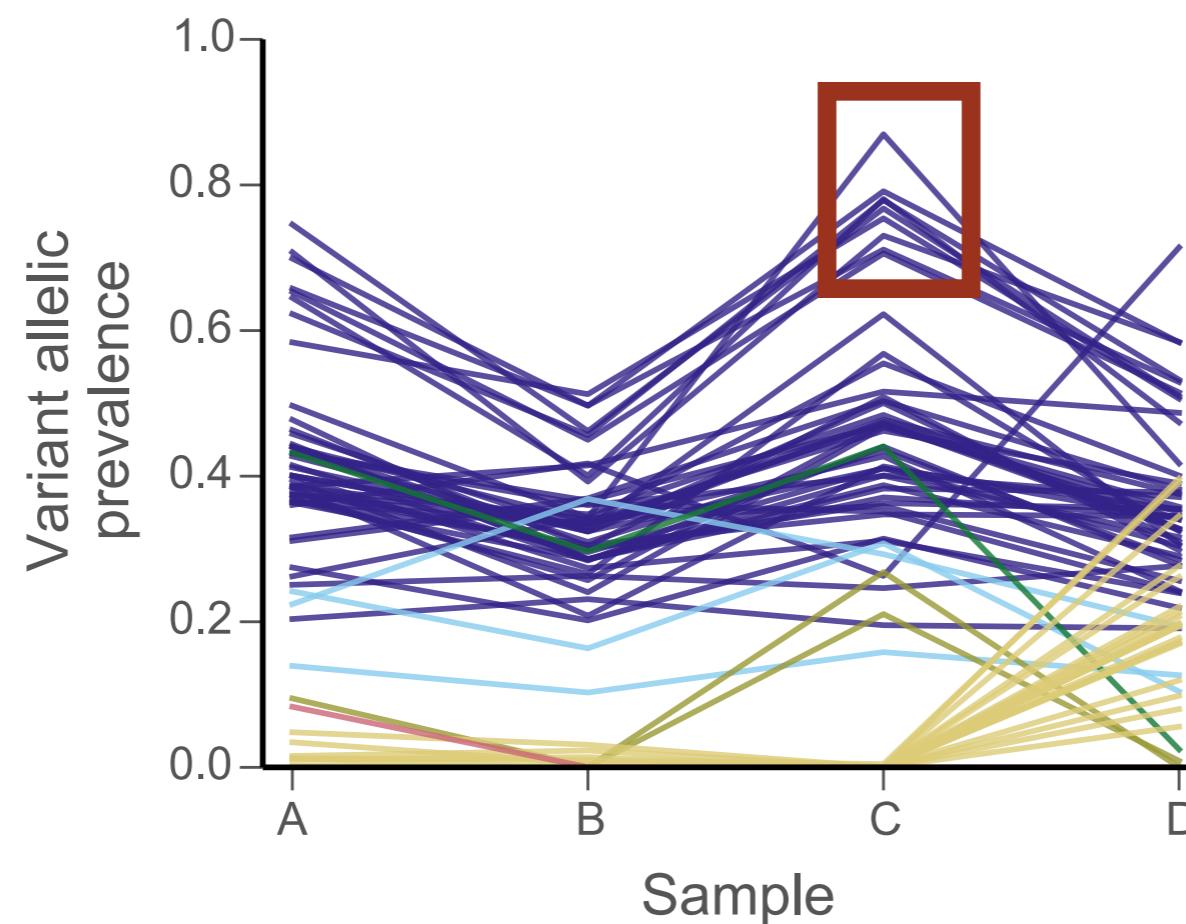


Each curve shows the allele prevalence of a mutation
(depth not shown)

4 biopsies (can be temporal or spatial)

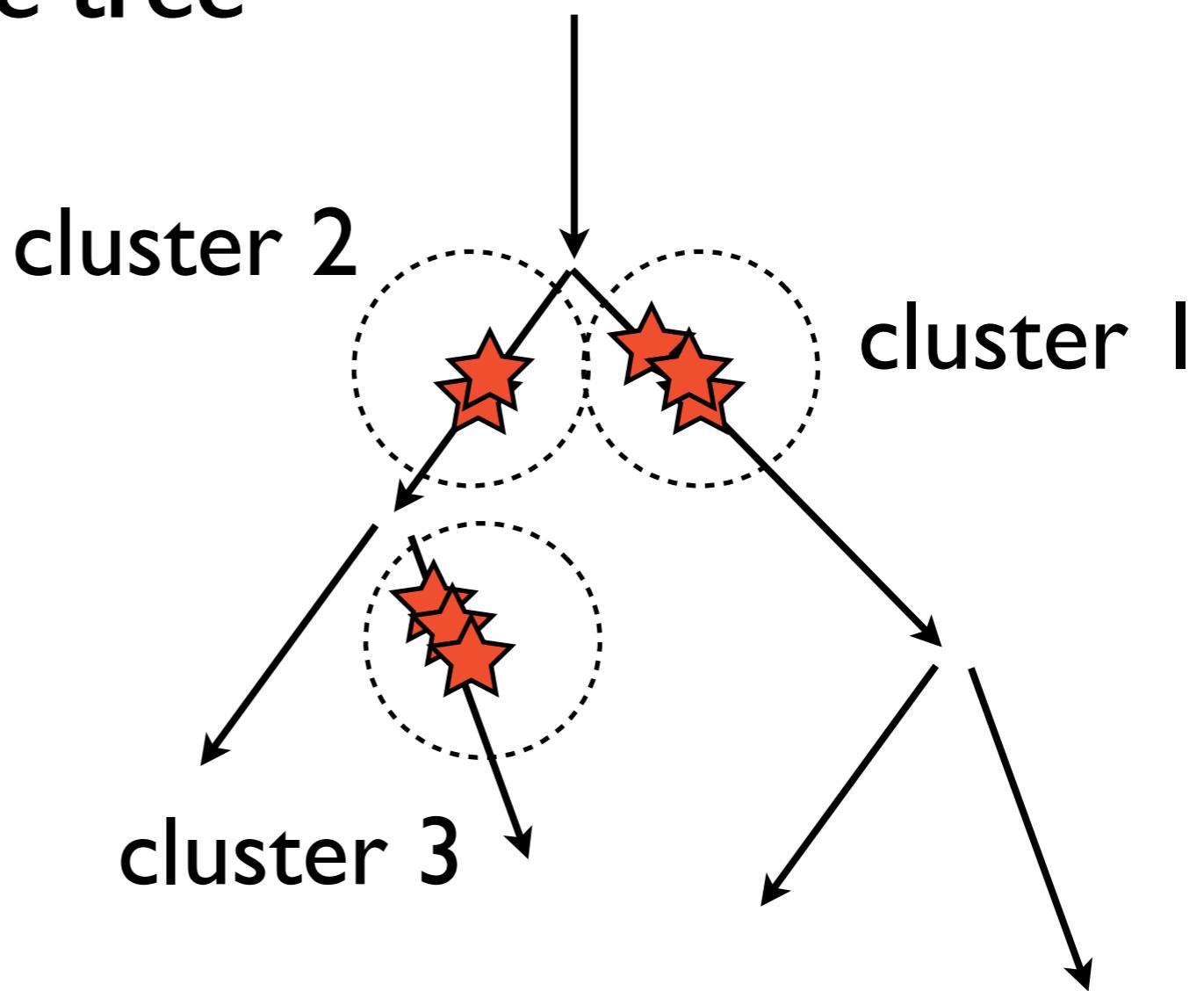


Clustering: motivation



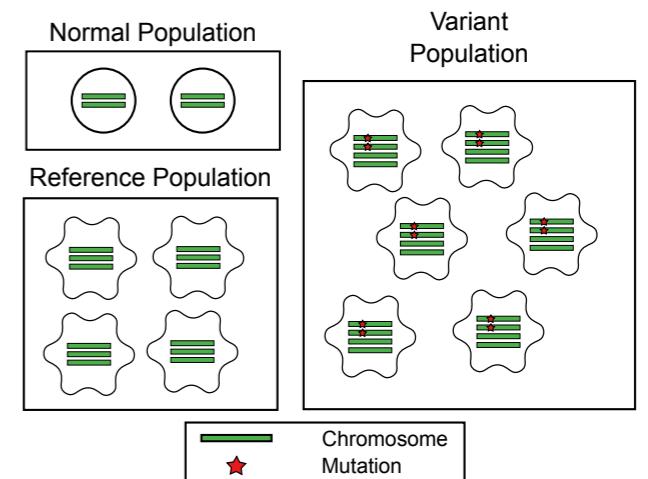
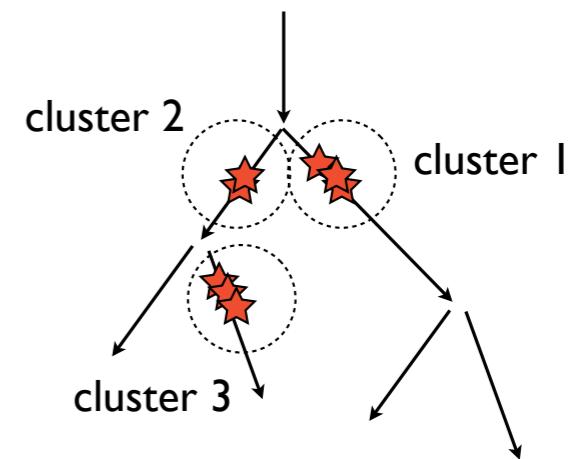
Clustering: motivation

Cluster ≈ edge of the tree



Inferential questions

- Cluster mutations
- Cellular prevalence of each sub-population.

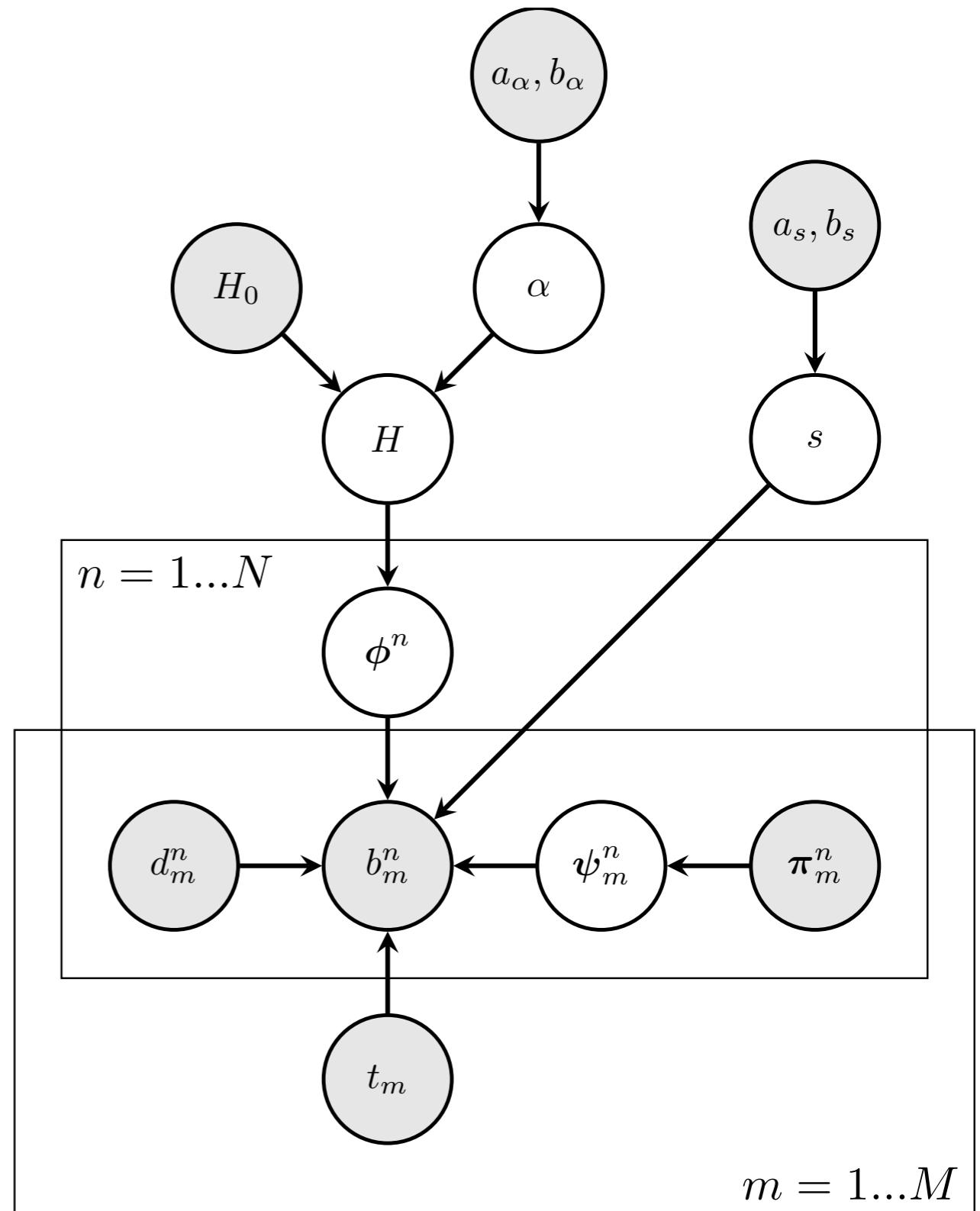


Challenges

- Noisy data
 - High level of uncertainty, partial identifiability
 - Many data types and sources of prior information
- => Good fit for Bayesian non-parametrics

PyClone

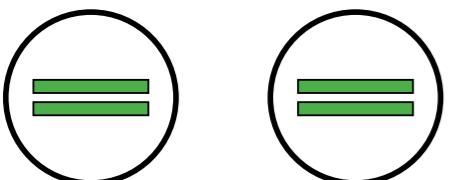
- Input: bulk ultra-deep sequencing data
[+ some prior info]
- Output: posterior distribution over clustering and cellular prevalences
- Based on the Dirichlet Process mixture model



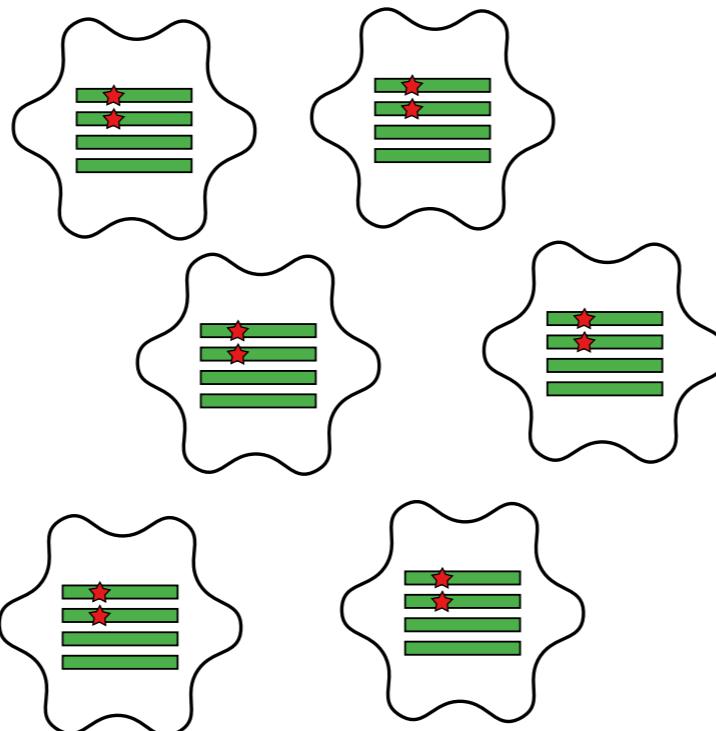
Roth et al. (2014) Nature Methods

PyClone: main parameter

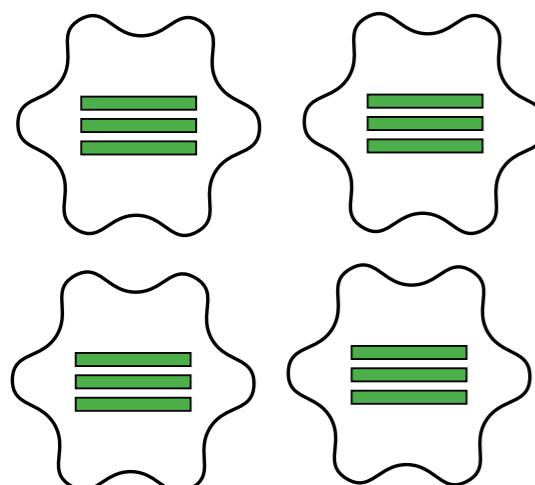
Normal Population



Variant Population



Reference Population



Chromosome
Mutation

Variant allelic prevalence
= # variant strands
/ # strands
= 12/40 ('observed')

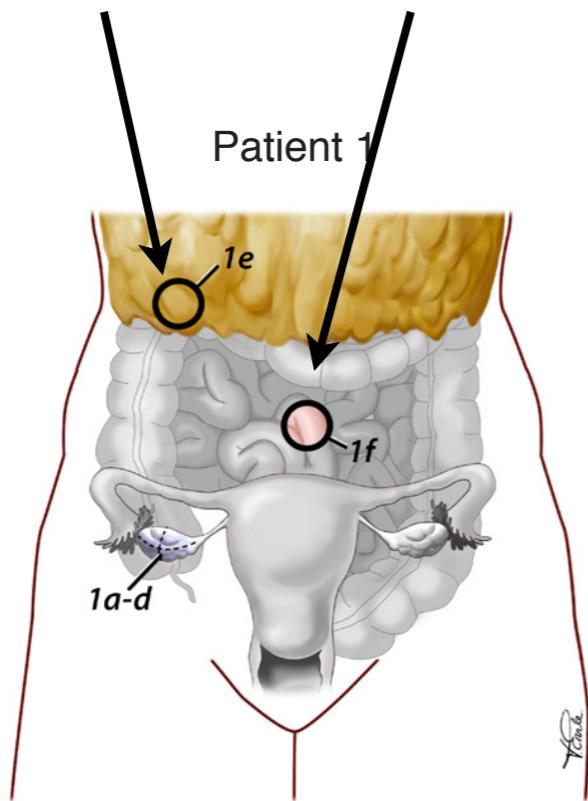
Variant cellular prevalence
= # variant cells / # cells
= 6/12 (inferred)

ϕ^n : cellular
prevalence vector
of mutation n

PyClone: main parameter

Note: the cellular prevalence varies across anatomical samples

$$\phi^n = (\phi_1^n, \dots, \phi_M^n)$$

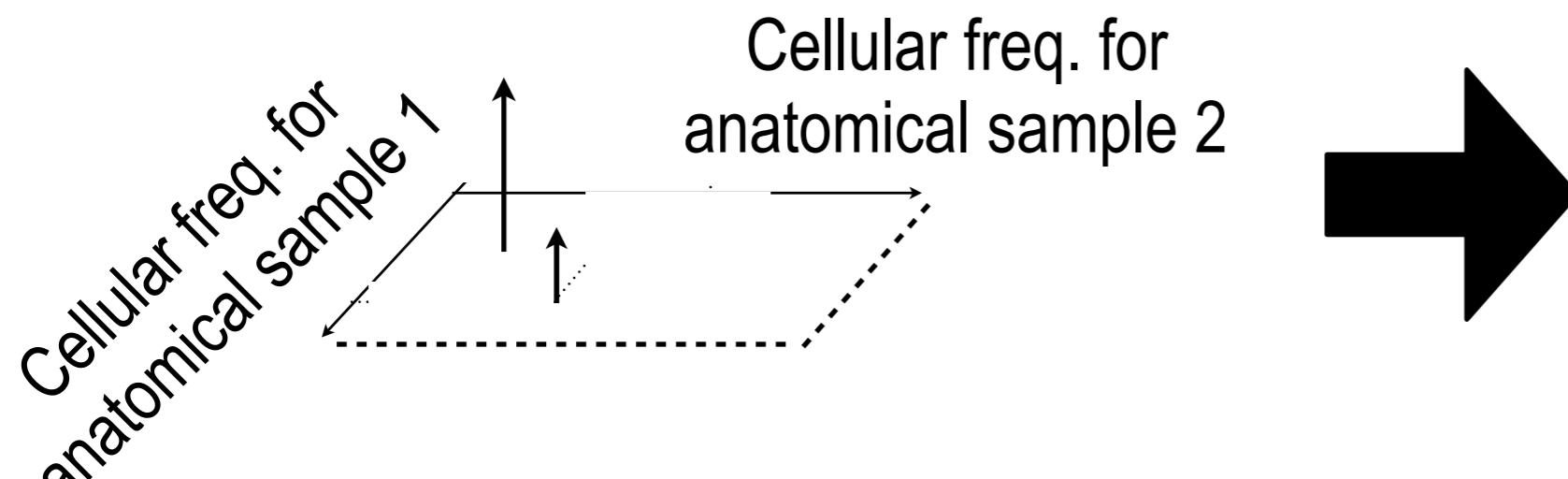


Variant cellular prevalence
= # variant cells / # cells
= 6/12 (inferred)

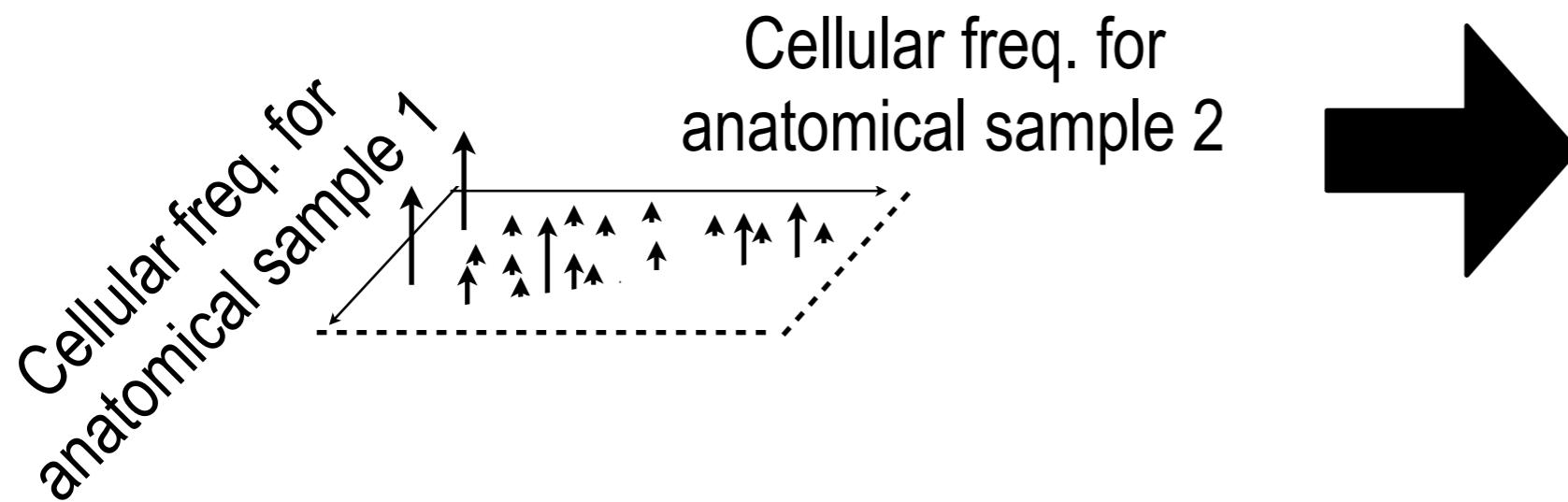
ϕ^n : cellular prevalence vector of mutation n

DP mixture model

Parametric mixture model:



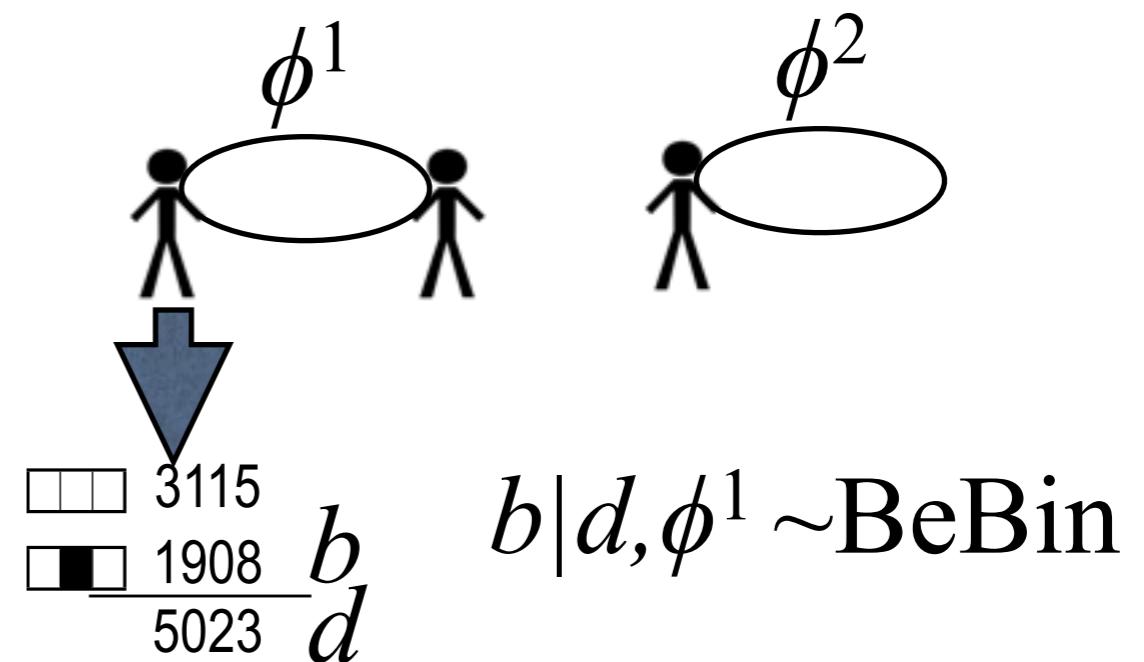
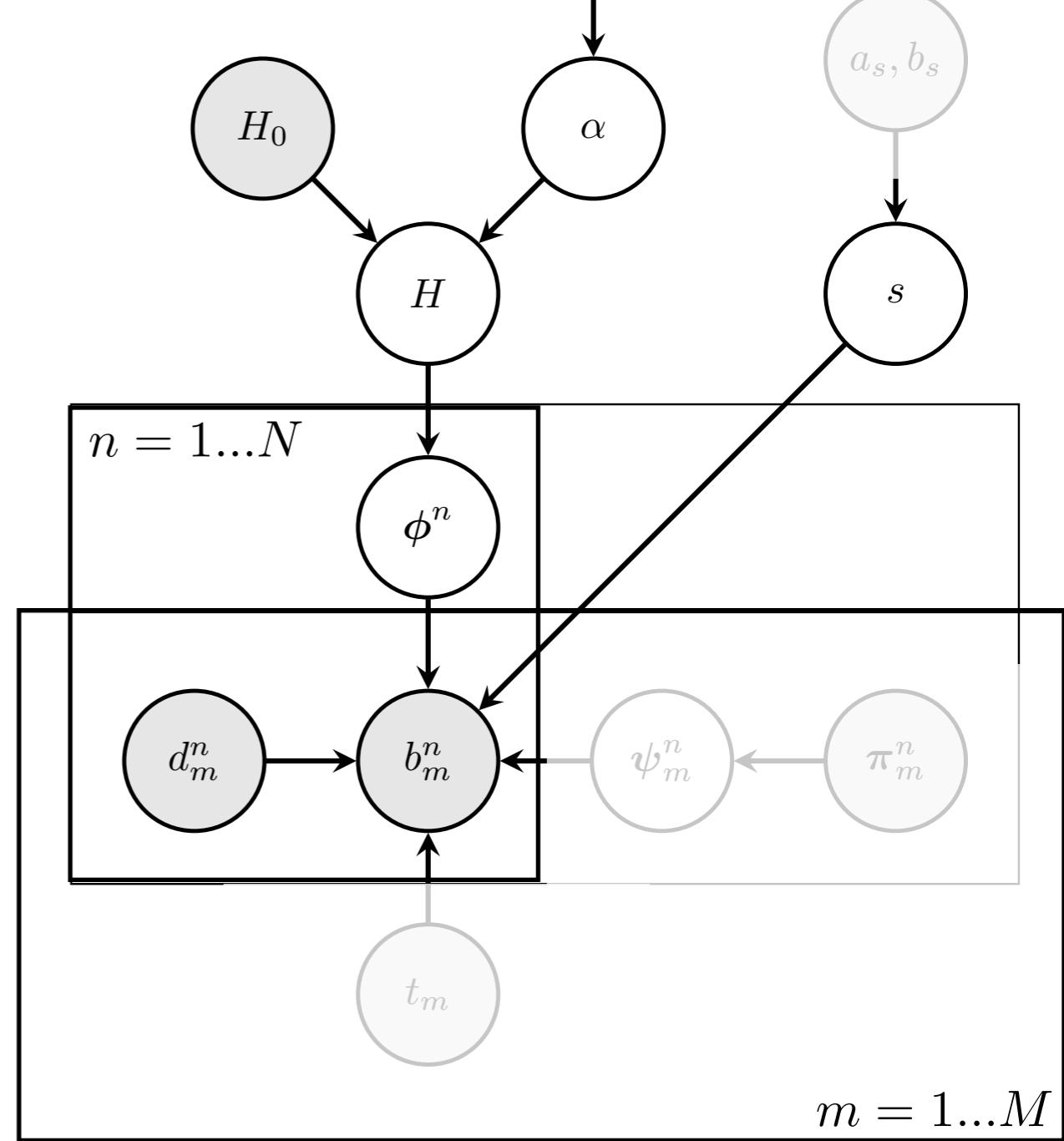
DP mixture model:



same

DP mixture

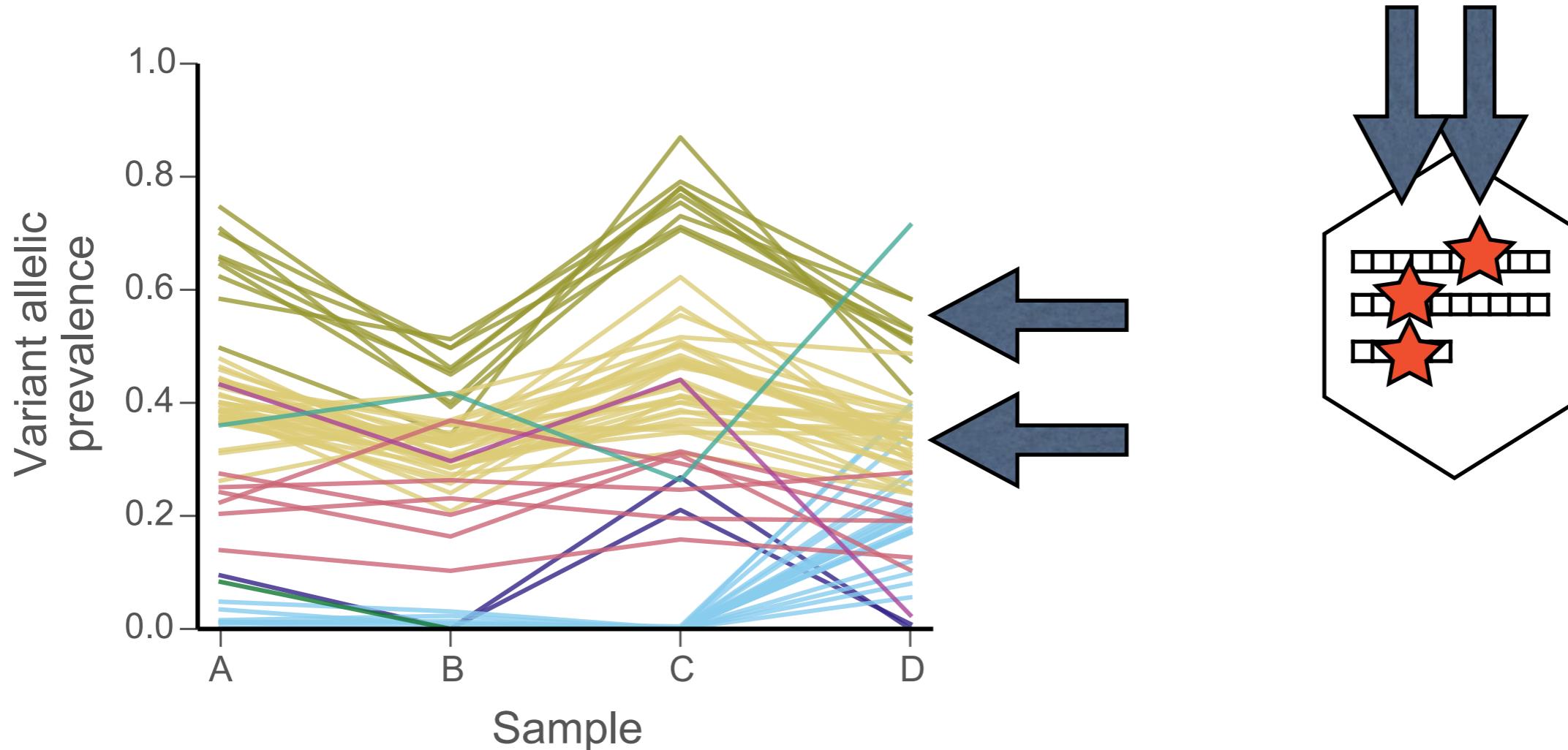
- n : index of the SNV loci
- m : index of sample (biopsy)
- ϕ^n : cellular prevalence vector of mutation n
- d : depth (total # of reads)
- b : number of mutant



Posterior inference

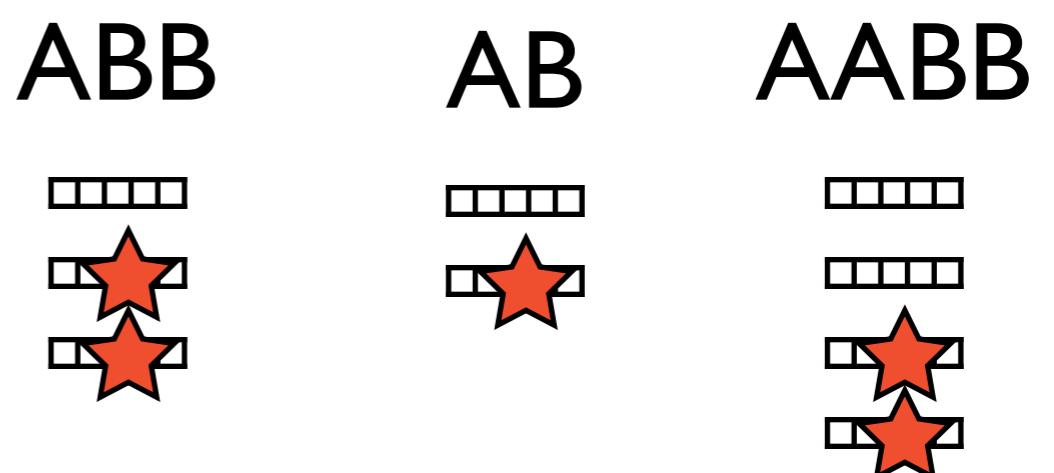
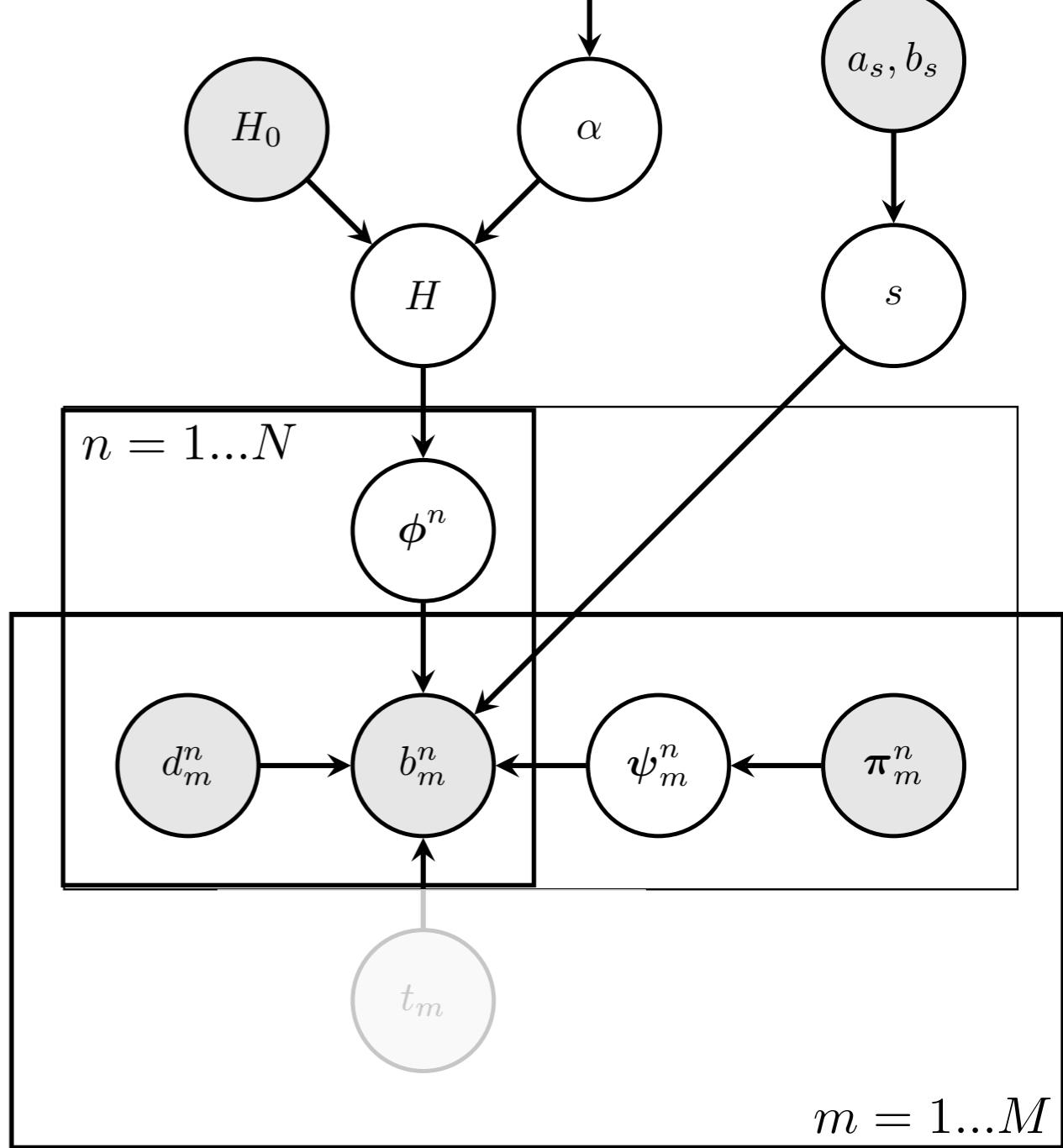
- MCMC: standard auxiliary variable method for non-conjugate models ('algorithm 8')
- Issue: scaling to large # of mutations
- Latest work with Andy Roth and Arnaud Doucet: a particle MCMC split merge algorithm

Example of clustering with the model so far

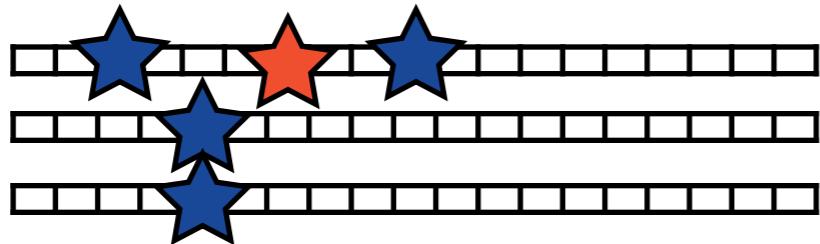


Genotypes

- n : index of the SNV loci
- m : index of sample (biopsy)
- ϕ^n : cellular prevalence vector of mutation n
- d : depth (total # of reads)
- b : number of mutant alleles
- s : overdispersion parameter
- ψ : genotype ($\in \{\text{AB}, \text{AAB}, \text{AABB}, \dots\}$)



Informed priors over genotypes

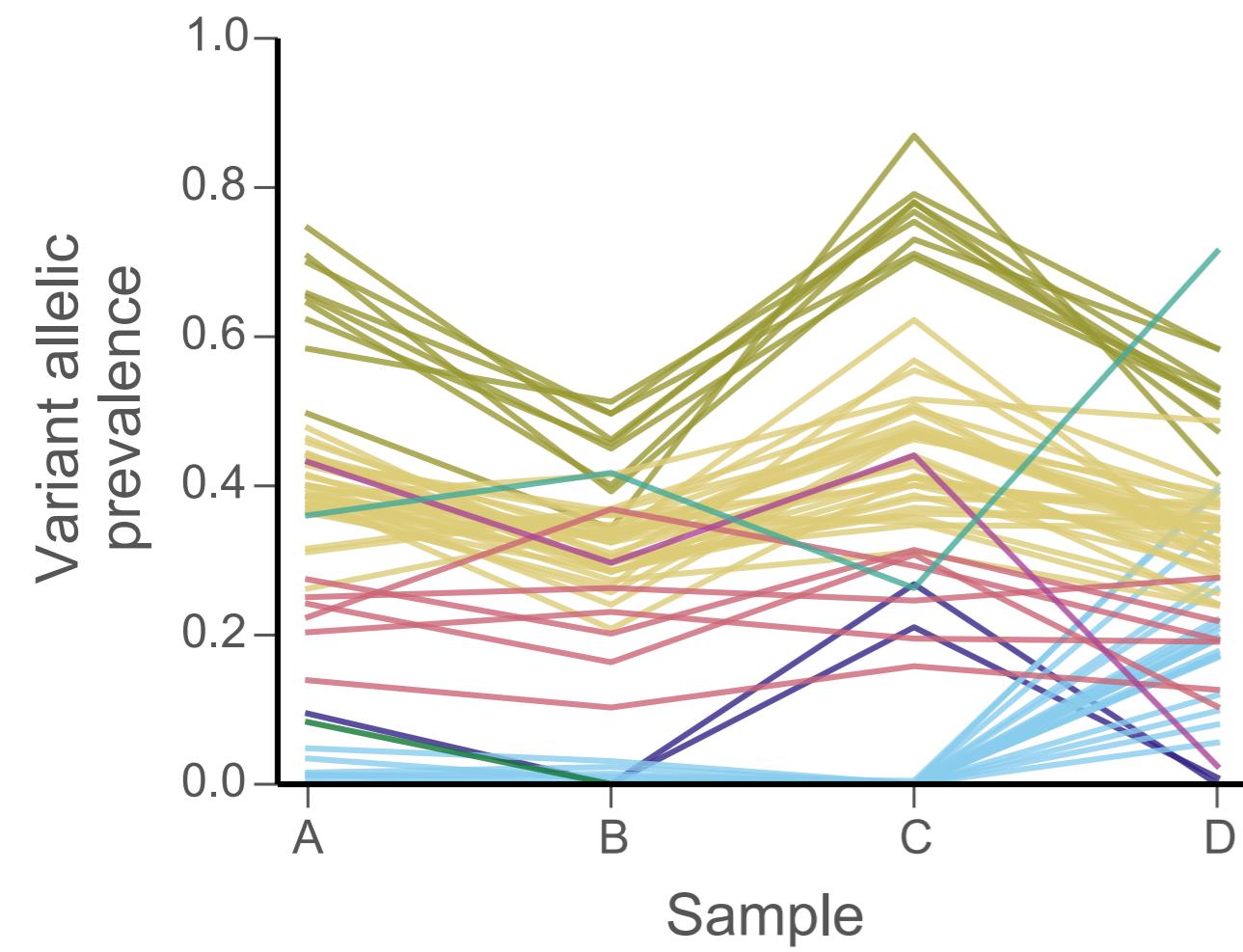


- ★ SNV (cancer)
- ★ SNP (from your parents)

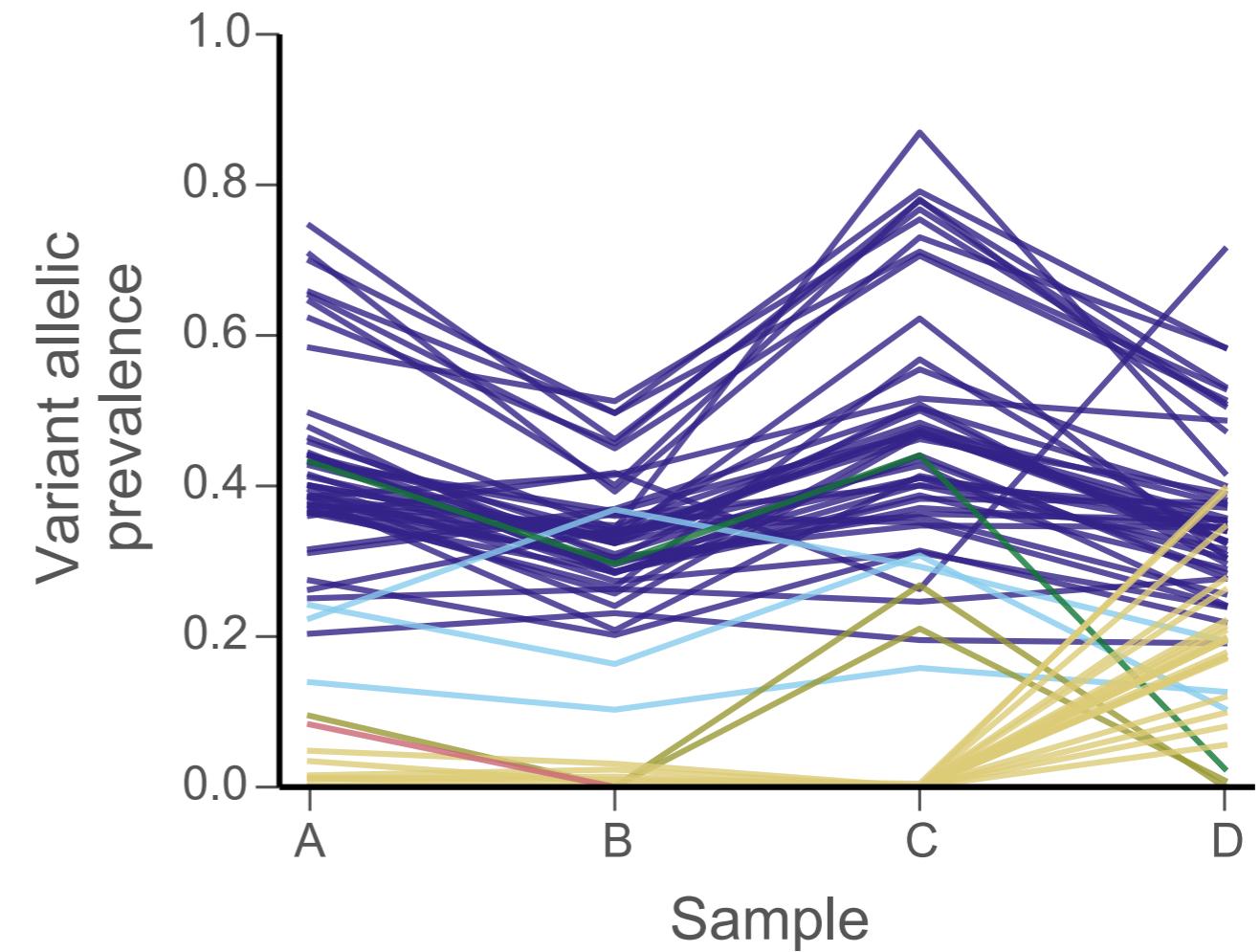
- Goal: reduce the space of possible genotypes
- Key idea: we can also use germline SNPs to get extra information on the copy number at a loci
- Hijack standard micro-array platforms to get log intensity in the neighborhood of the SNV of interest

Clustering of mutations

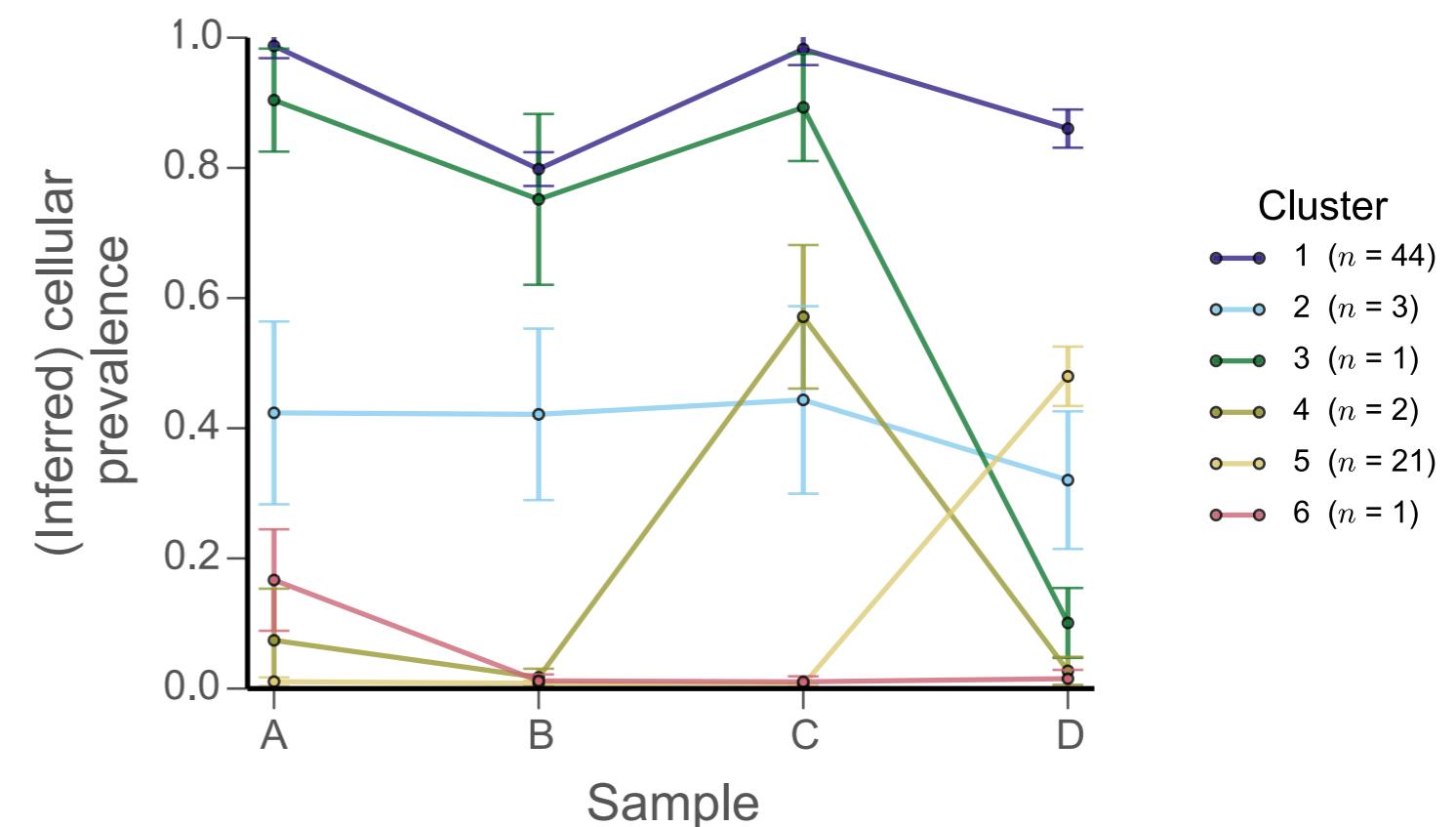
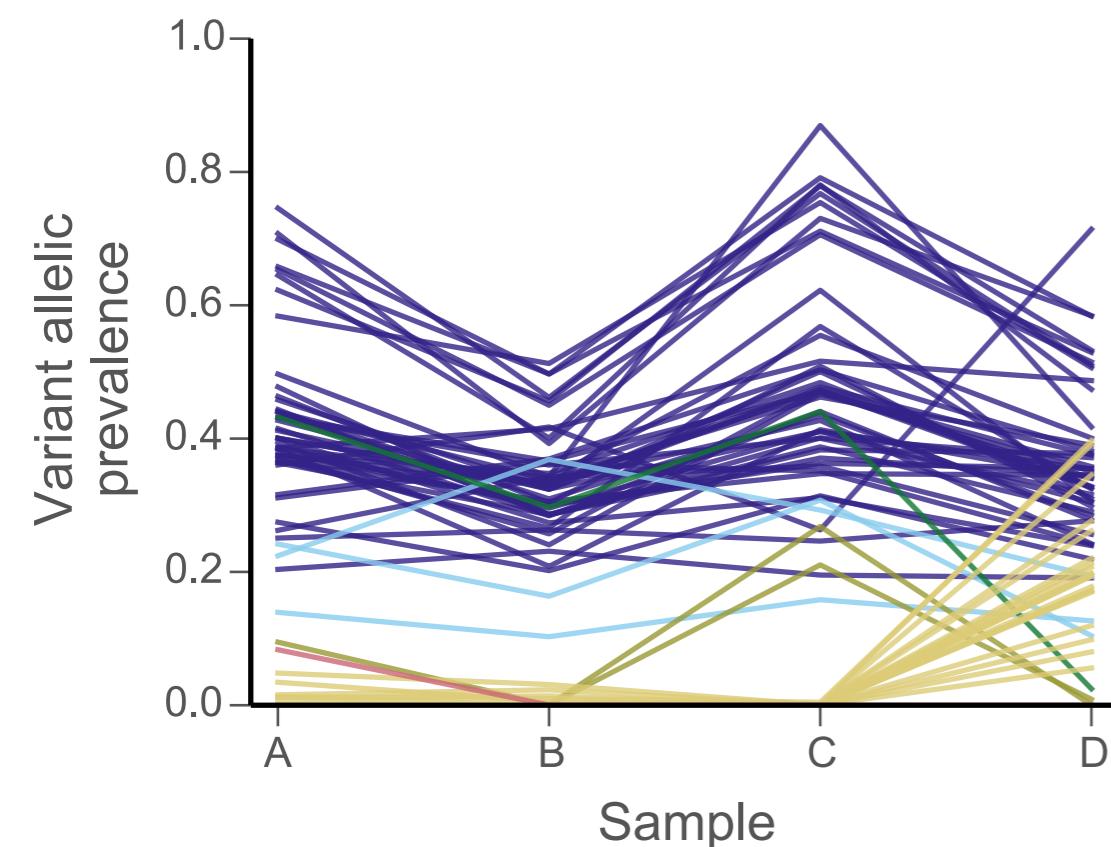
Without genotypes



With genotypes

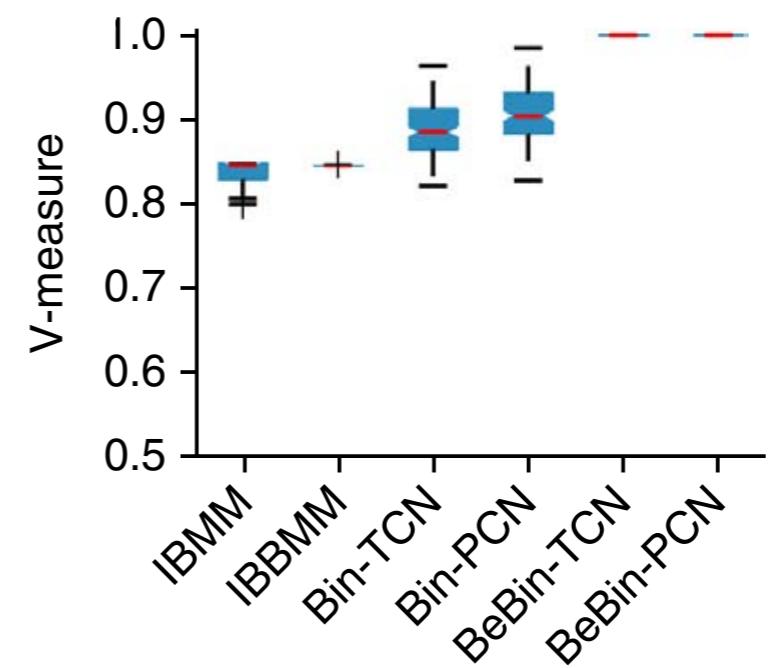
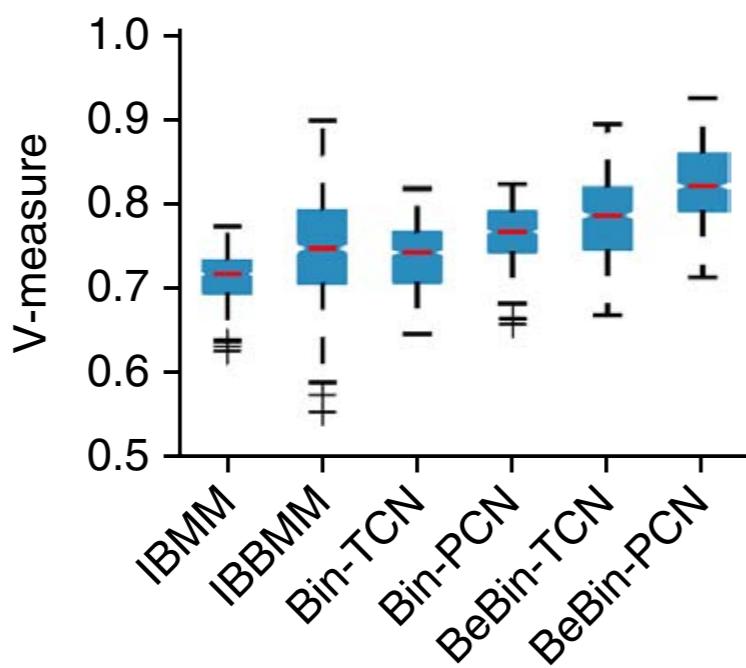


Full output (consensus)

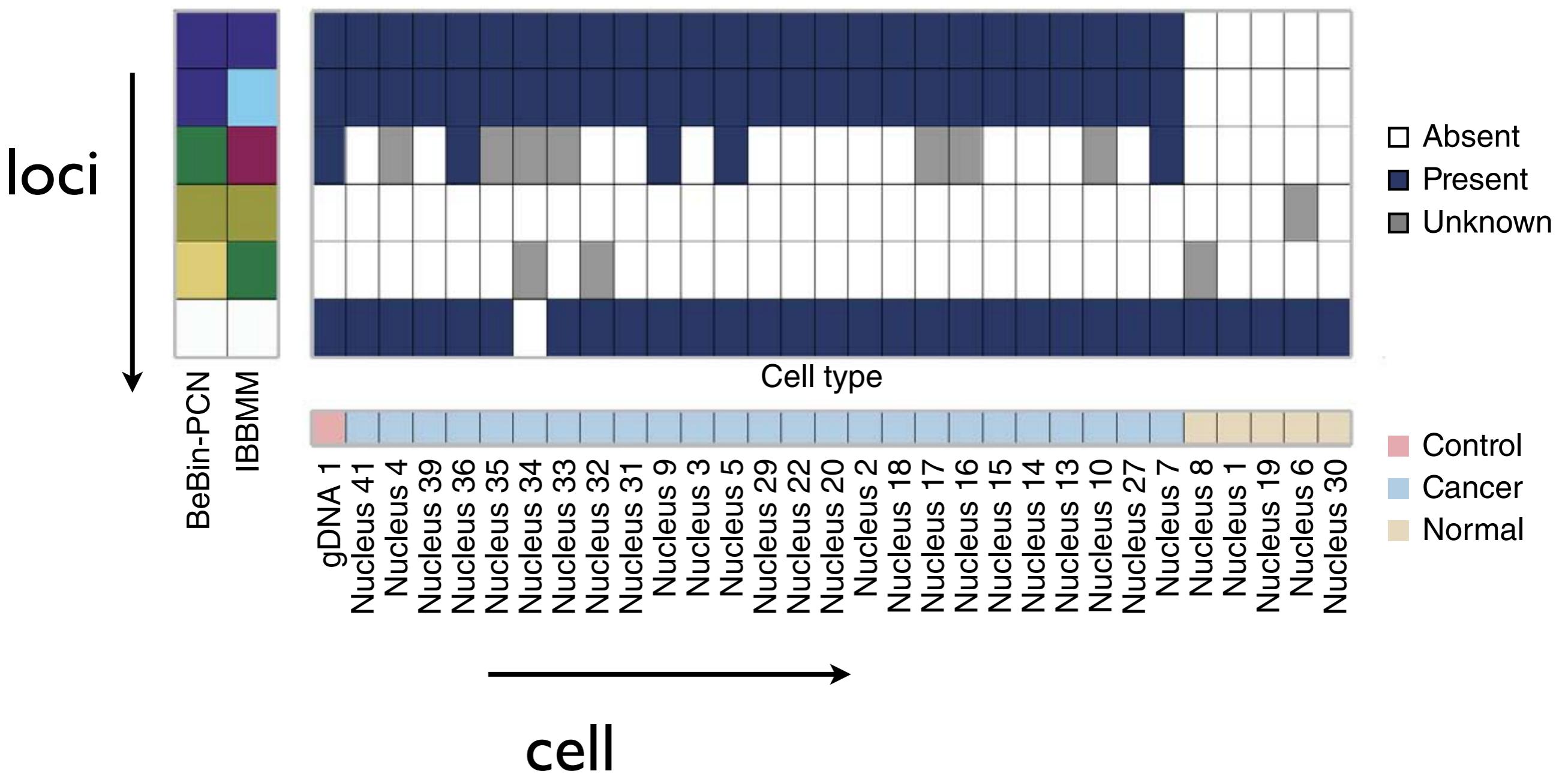
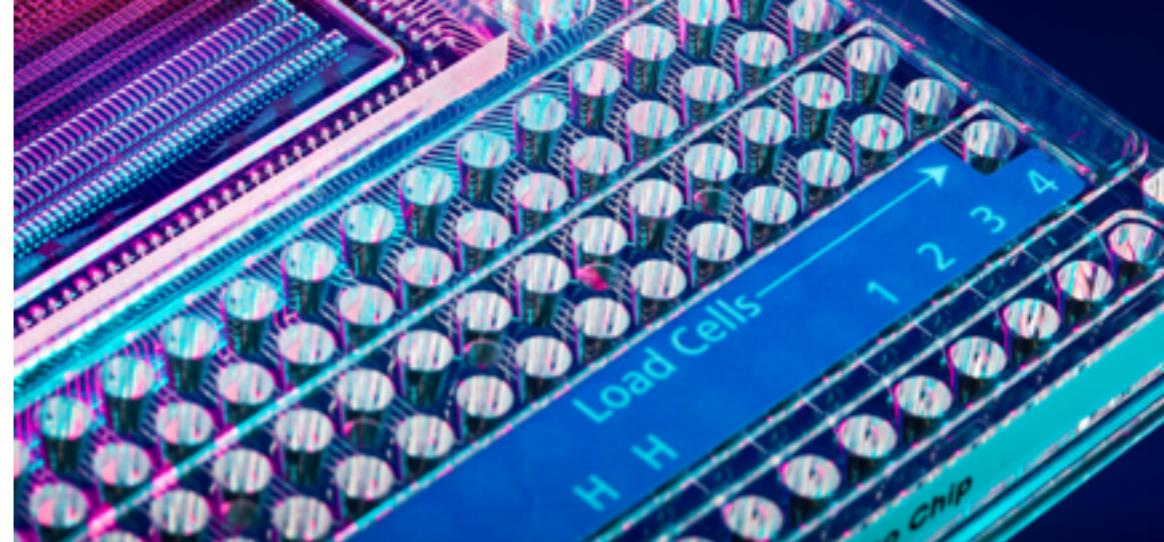


How to validate?

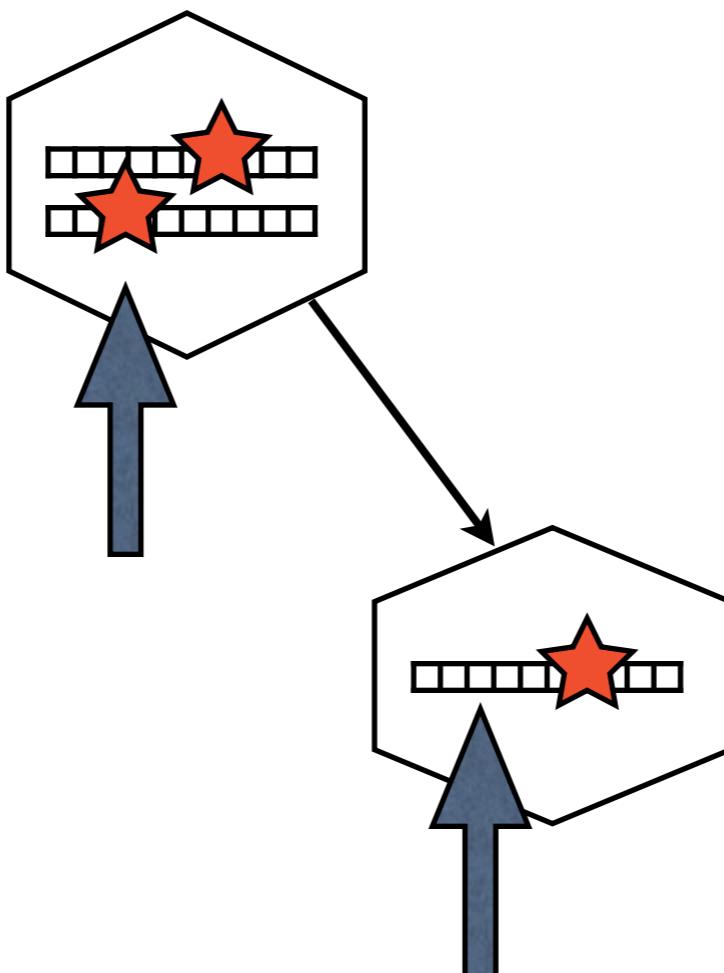
- Cross validation
- Synthetic data (in silico and in vitro)



Single-cell sequencing



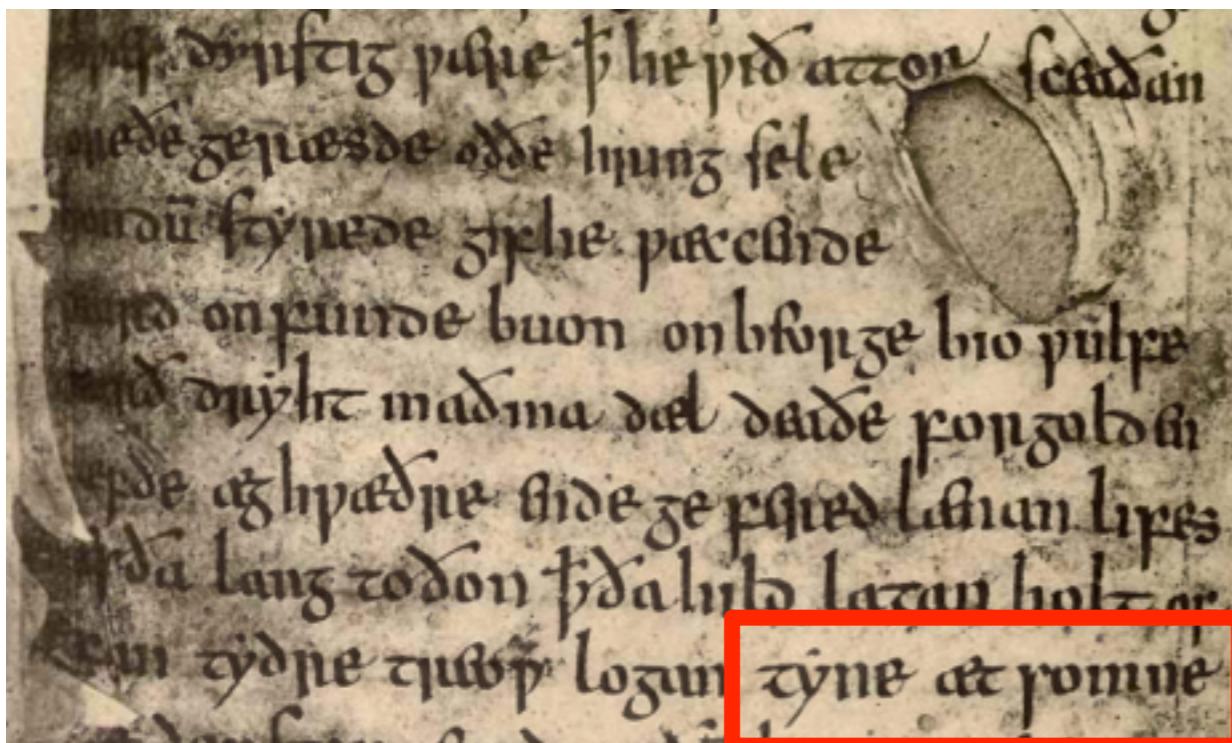
Issue: loss of SNVs



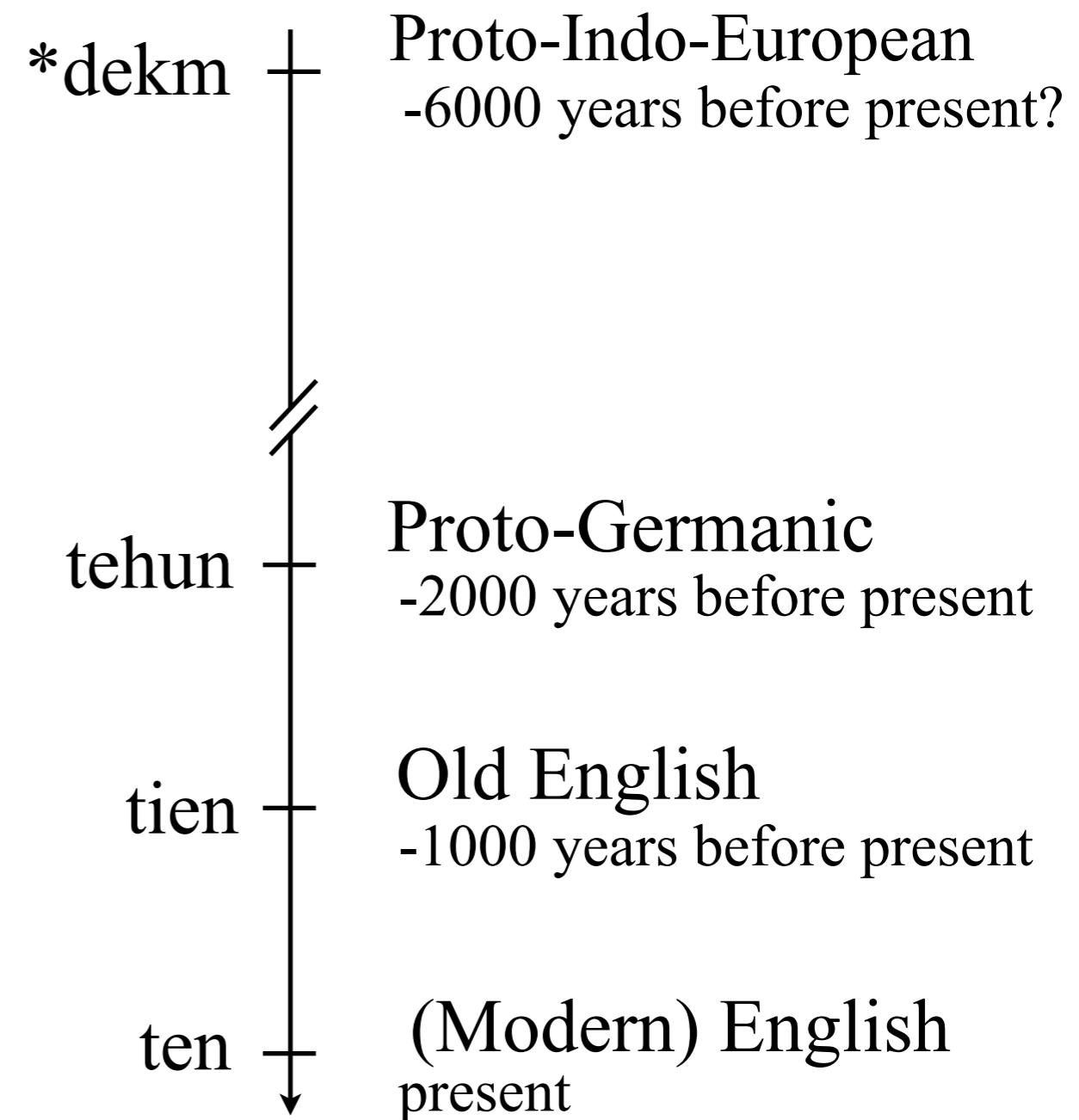
**Application: reconstructing
ancient languages**

Examples of language change

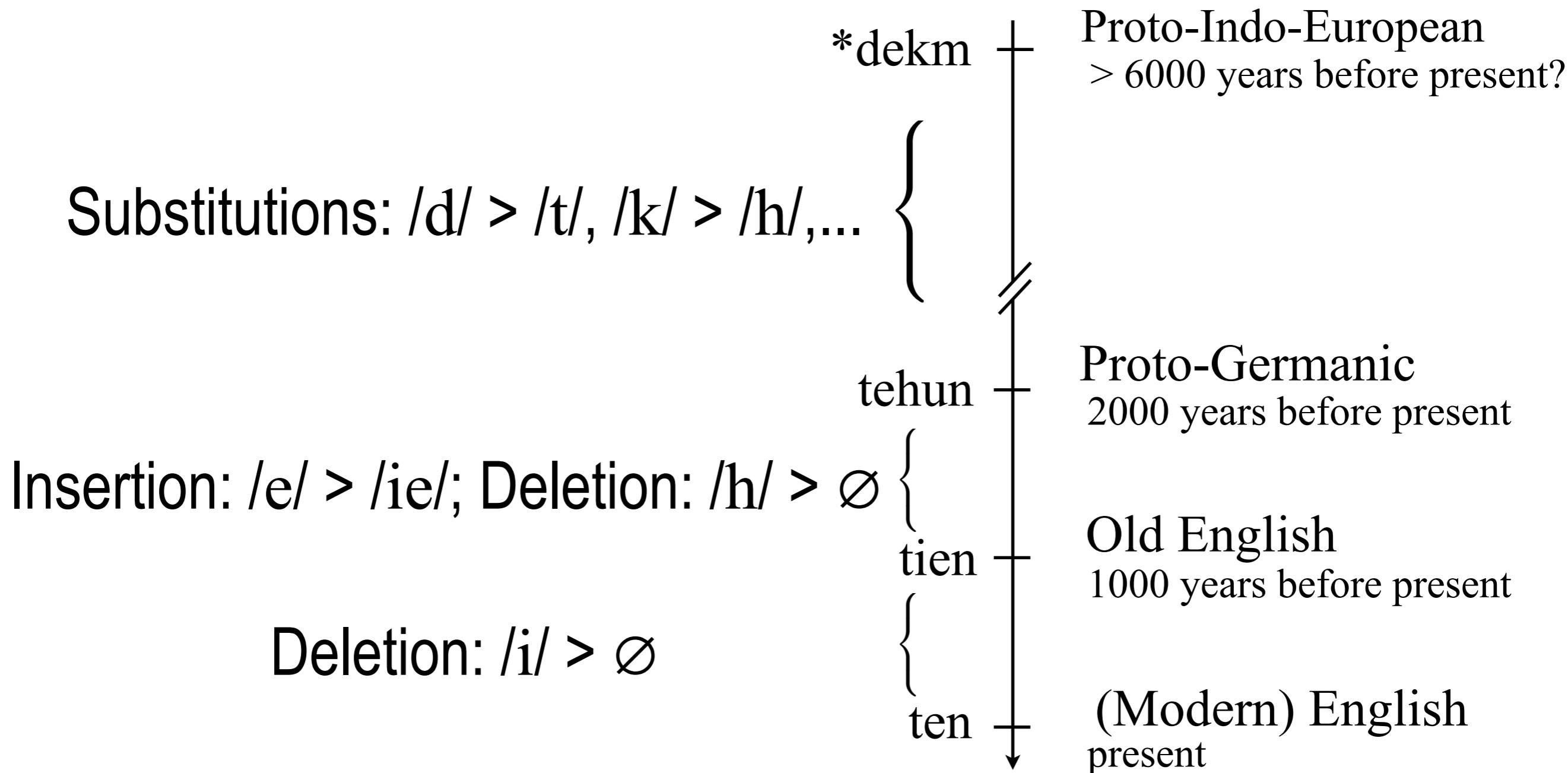
Beowulf (~ 8th - 11th century)



tyne ætsomne
'ten altogether'



Examples of language change



Example of data

Austronesian dataset:



	'fish'	'fear'
Hawaiian	i?a	maka?u
Samoan	i?a	mata?u
Tongan	ika	
Proto-Oceanic	*ika	*mataku

:

Size: 706 languages, 150k word forms in IPA

[Greenhill et al, '08]

Oceanic language



Task: comparative reconstruction



	'fish'
Hawaiian	i?a
Samoan	i?a
Tongan	ika
Maori	ika

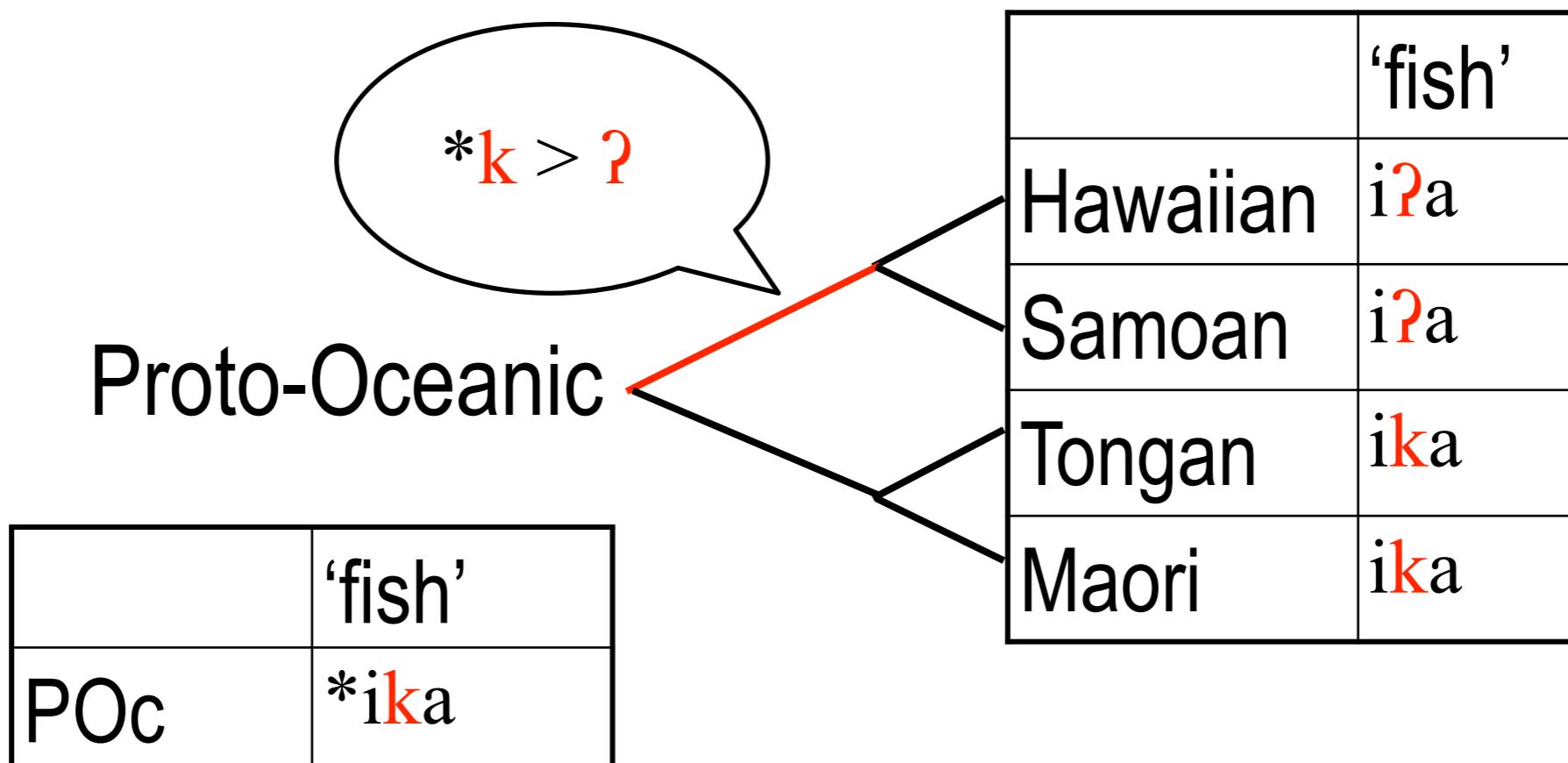
Task: comparative reconstruction

	'fish'
Hawaiian	i? ² a
Samoan	i? ² a
Tongan	ika
Maori	ika

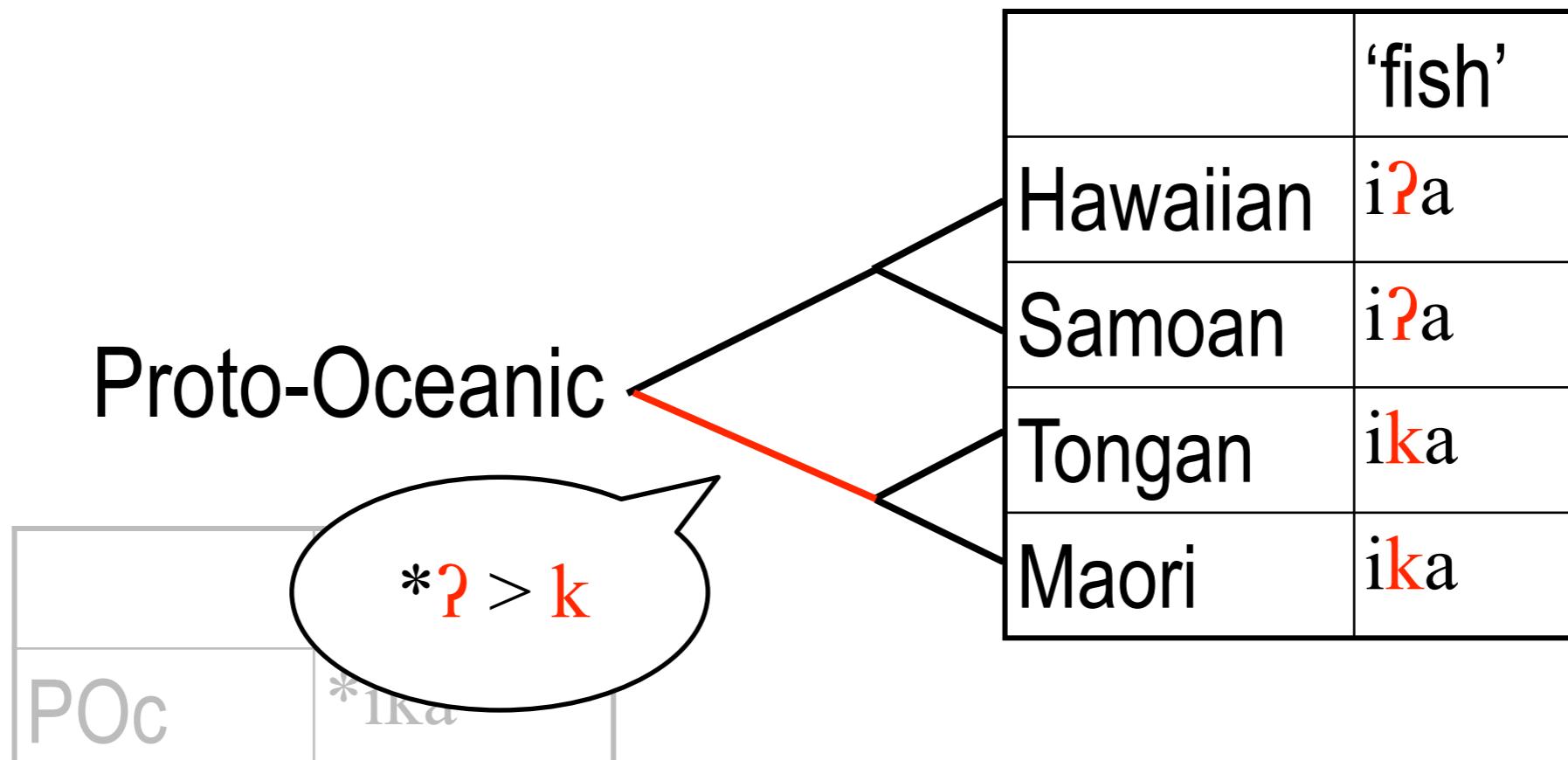
Proto-Oceanic

```
graph LR; Root[Proto-Oceanic] --- Hawaiian[Hawaiian]; Root --- Samoan[Samoan]; Root --- Tongan[Tongan]; Root --- Maori[Maori]
```

Task: comparative reconstruction



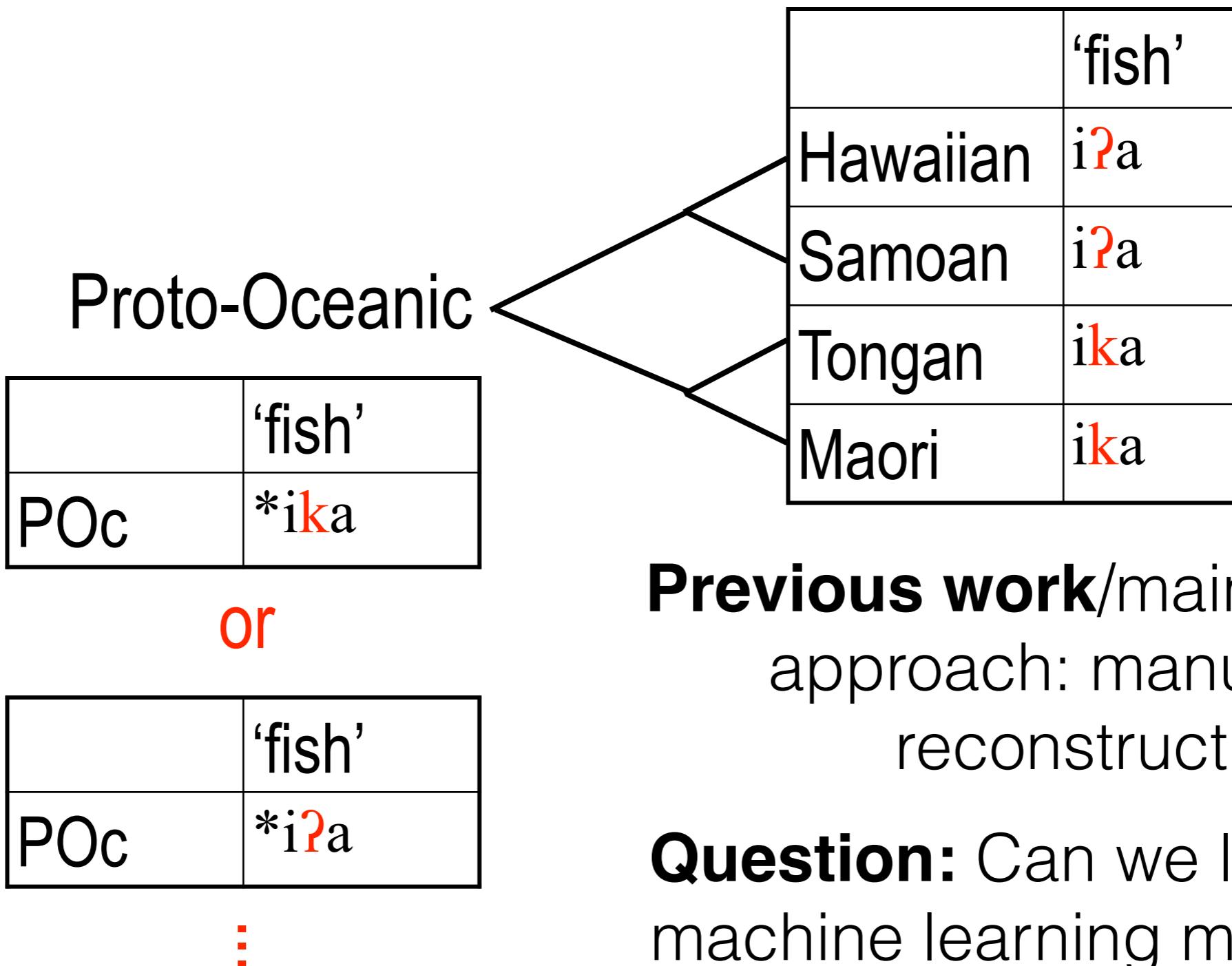
Task: comparative reconstruction



or

	'fish'
POc	$*\text{i}\text{?}$ a

Task: comparative reconstruction



Can we harness more languages?



	'fish'
Hawaiian	i [?] a
Samoan	i [?] a
Tongan	ika
Maori	ika
Geser	ikan
Rapanui	ika
Nukuoro	iga
Niue	ika

Task 2: cognate inference

Fact: New words are invented & old words fall out of usage

Consequence: not all words forms for a given concept will have cousins in related languages.

	'fish'	'fear'	'fear'
Hawaiian (G)	i?a	maka?u	
Samoan (S)	i?a	mata?u	
Tongan (T)	ika		manavahē
Maori (M)	ika	mataku	

Column:
'cognate
set' (think as a
cluster)

Task 2: cognate inference

Hawaiian	i?a, maka?u, ...
Samoan	i?a, maka?u, ...
Tongan	ika, manavahē, ...
Maori	ika, mataku, ...



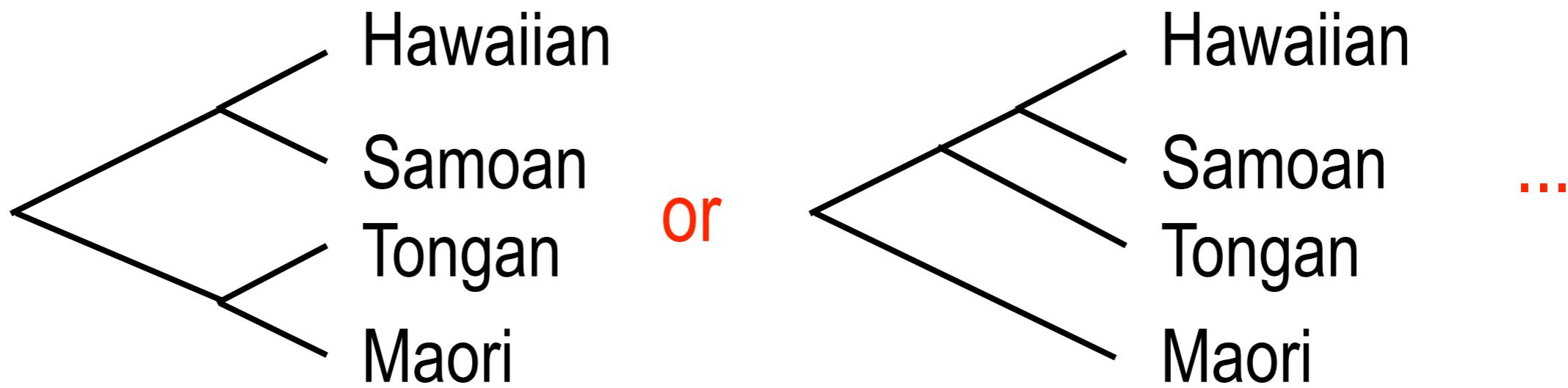
Hawaiian	i?a	maka?u
Samoan	i?a	mata?u
Tongan	ika	manavahē
Maori	ika	mataku

or

Hawaiian	i?a	maka?u	
Samoan	i?a	mata?u	
Tongan	ika		manavahē
Maori	ika	mataku	

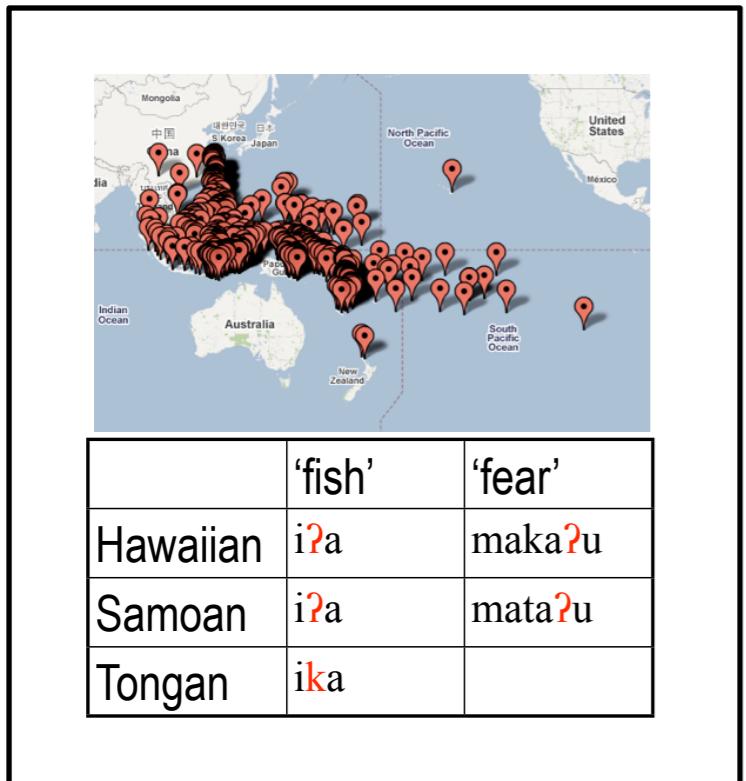
...

Task 3: tree inference

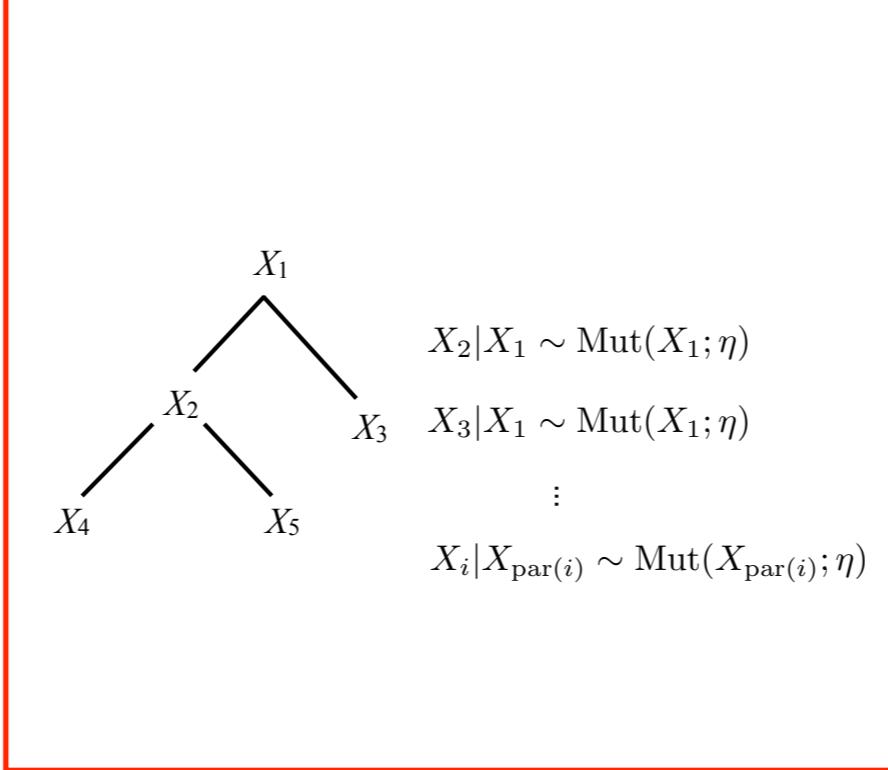


Model-based approach

Input: modern words



Probabilistic
models of change



Reconstruction

	'fish'
POc	*i?ka

or

	'fish'
POc	*i?ka

⋮

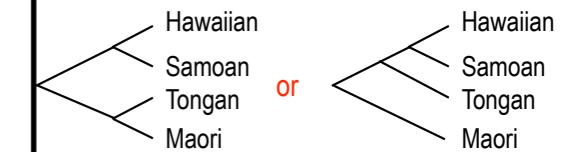
Cognates

	'fish' (1)	'fish' (2)
Hawaiian	i?a	
Samoan	i?a	
Tongan	ika	
Marshallese	yapil	

or

	'fish' (1)	'fish' (2)
Hawaiian	i?a	
Samoan	i?a	
Tongan	ika	
Marshallese	yapil	

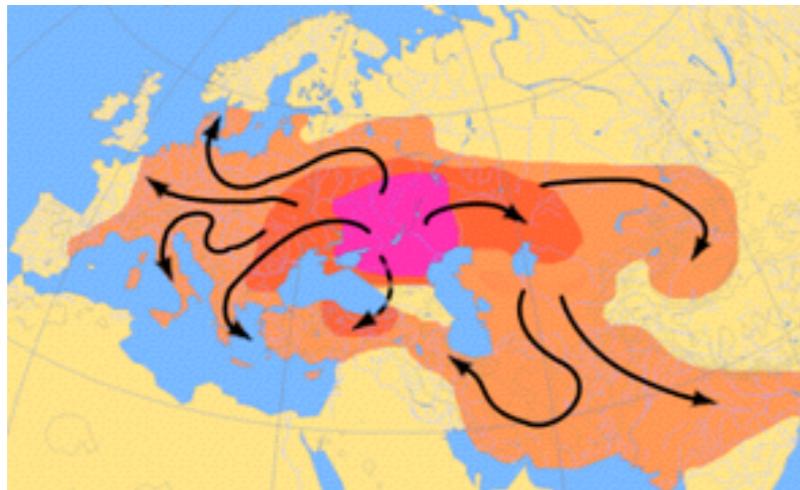
Trees



Motivation

Why reconstruct?

- Can answer a large number of questions about our past
 - Learn about ancient populations' migrations
 - Decipherment of ancient scripts



Motivation: typology

- What are the statistical regularities of sound change?
 - Role of functional load in sound change
- Useful for synchronic typology as well
 - Is an observed pattern the result of structural constraints, or the result of common descent?

Motivation: machine translation (MT)

- Long-term challenge: usable machine translation between *all* pairs of languages
 - Existing MT rely on large corpora in both languages in the pair,
 - in particular, *bitexts* (e.g. europarl)
- Concrete step:
large scale cognate detection

Hawaiian	i?a, maka?u, ...
Samoan	i?a, maka?u, ...
Tongan	ika, manavahē, ...
Maori	ika, mataku, ...



Hawaiian	i?a	maka?u
Samoan	i?a	mata?u
Tongan	ika	manavahē
Maori	ika	mataku

or

Hawaiian	i?a	maka?u	
Samoan	i?a	mata?u	
Tongan	ika		manavahē
Maori	ika	mataku	

Background

International Phonetic Alphabet

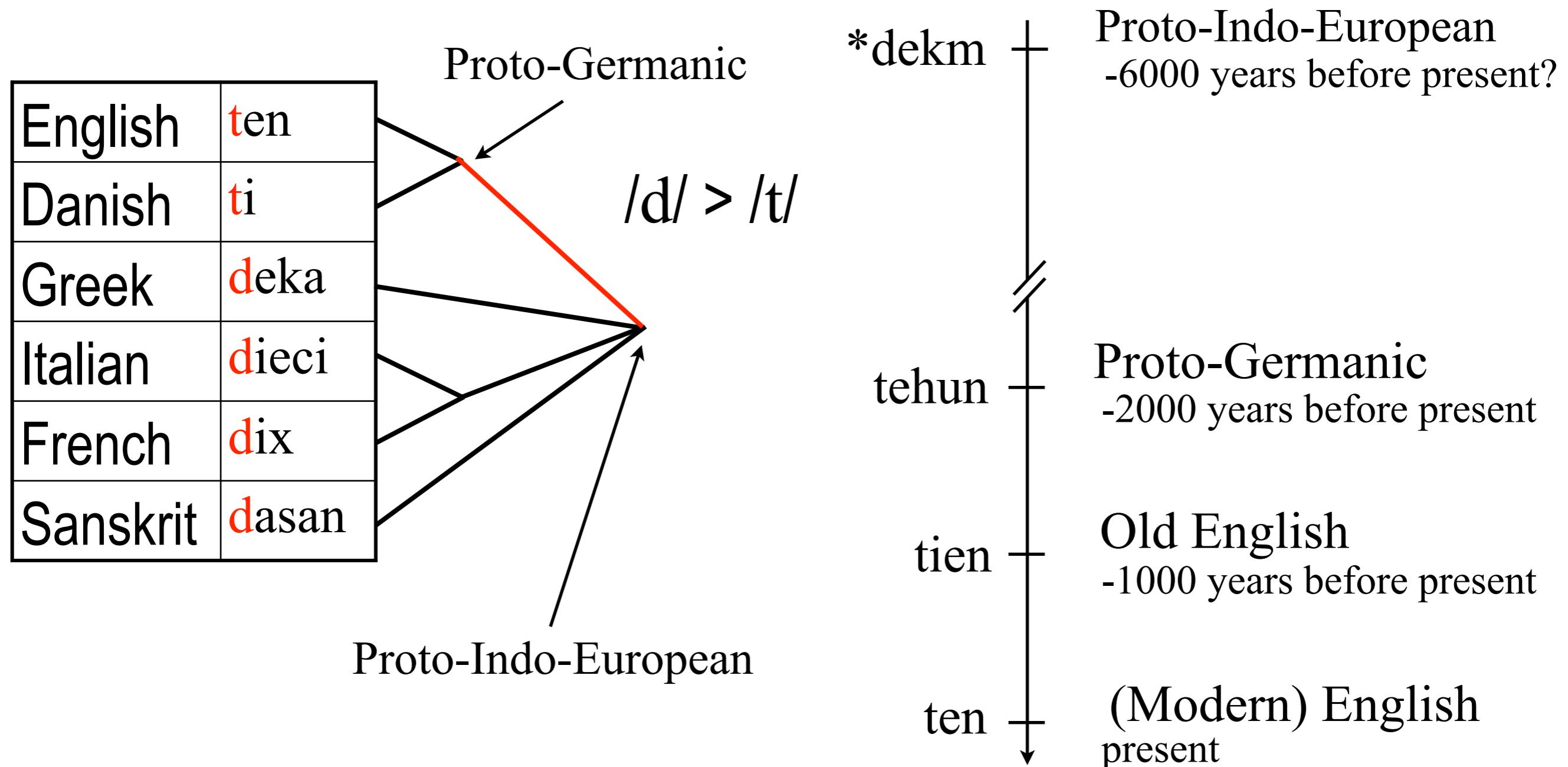
Instead of {A,C,G,T}, here the alphabet is a set of basic sound units (*phonemes*):

m	mj	n	ŋ	N
p b		t d	k g	q G ʔ
ɸ β	f v	s z	x γ	χ h ɦ

i y	i u	ɯ u
e ø	ɛ e	ɤ o
ɛ œ	ɜ ə	ʌ ɔ
a ə		ɑ ɒ

	'fish'	'fear'
Hawaiian	iɸa	makaɸu
	iɸa	mataɸu
An	ika	
to-Oceanic	ika	mataku

Sound changes



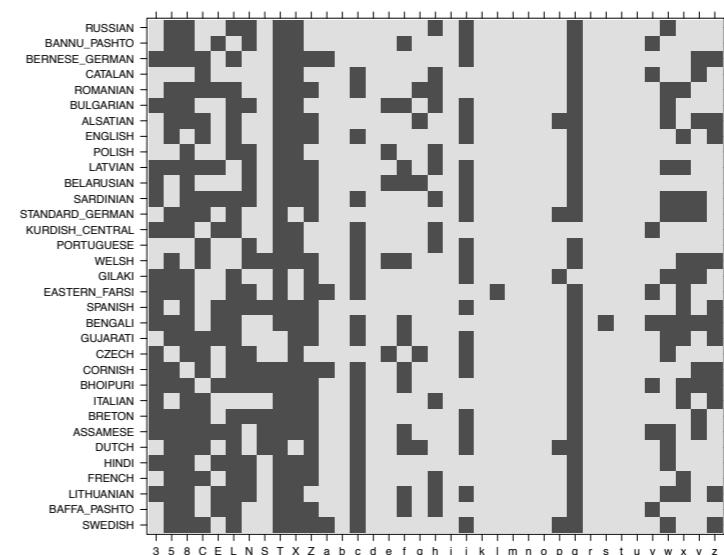
Regularity in sound change

English	ten	tooth
Danish	ti	tand
Greek	deka	donti
Italian	dieci	dente
French	dix	dent
Sanskrit	dasan	danta

/d/ > /t/ in word initial position

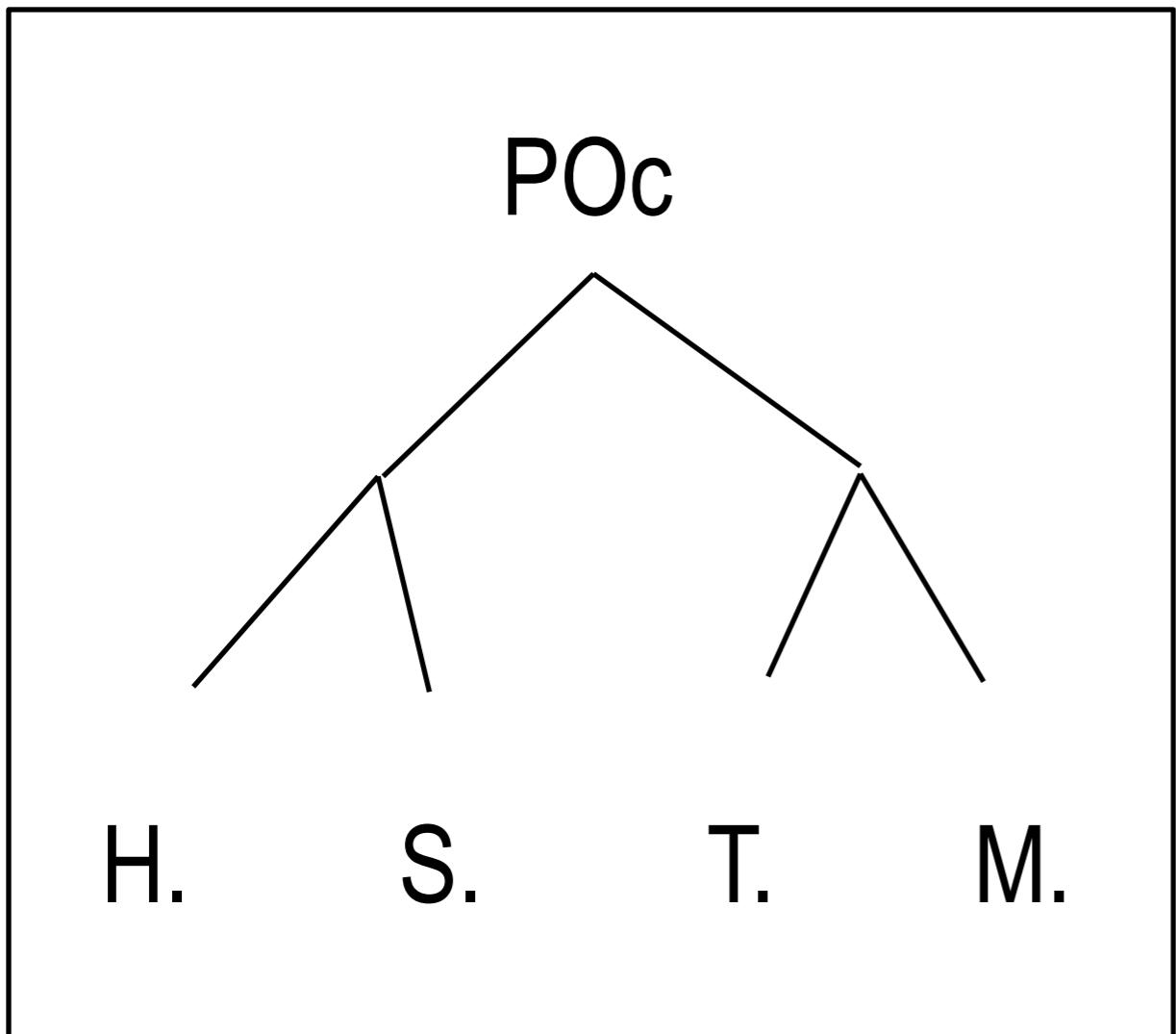
Proto-Indo-European

- Probabilistic models of inventory changes?
 - Open computational problem



Modelling phonological change

Simplifying assumptions

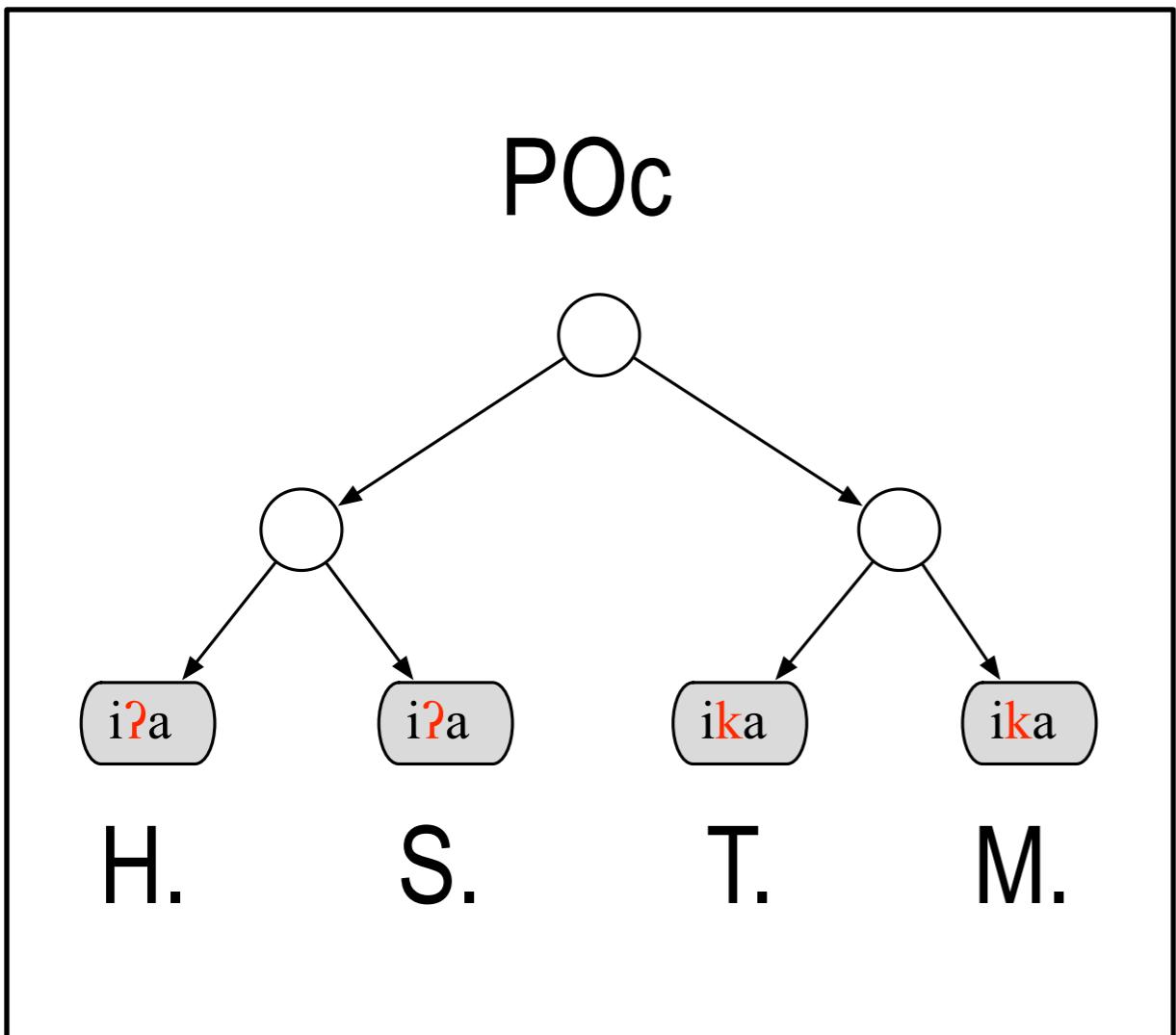


Fixed tree

	'fish'	'fear'
Hawaiian	i?a	maka?u
Samoan	i?a	mata?u
Tongan	ika	
Maori	ika	mataku

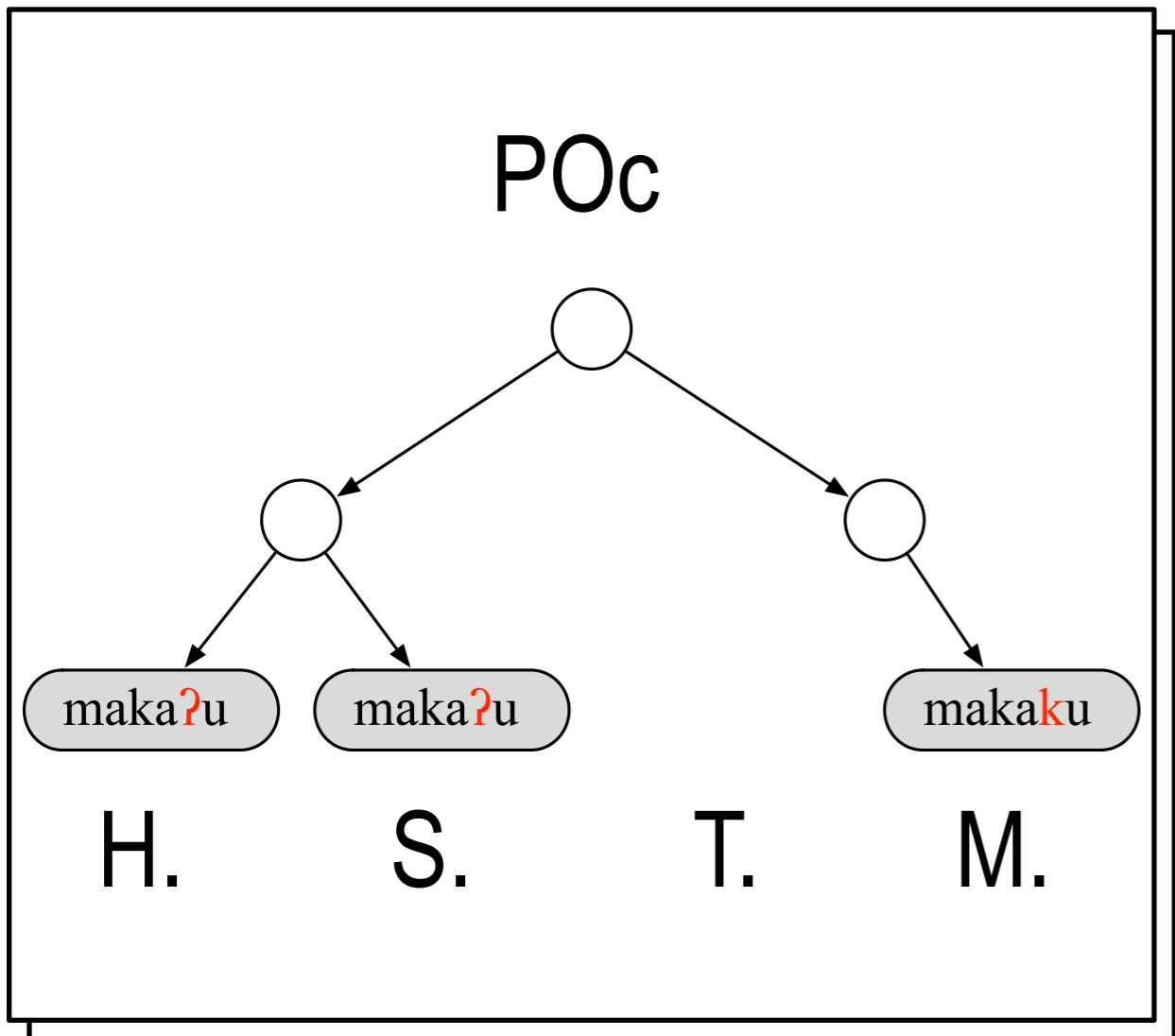
Fixed cognate sets

Graphical model



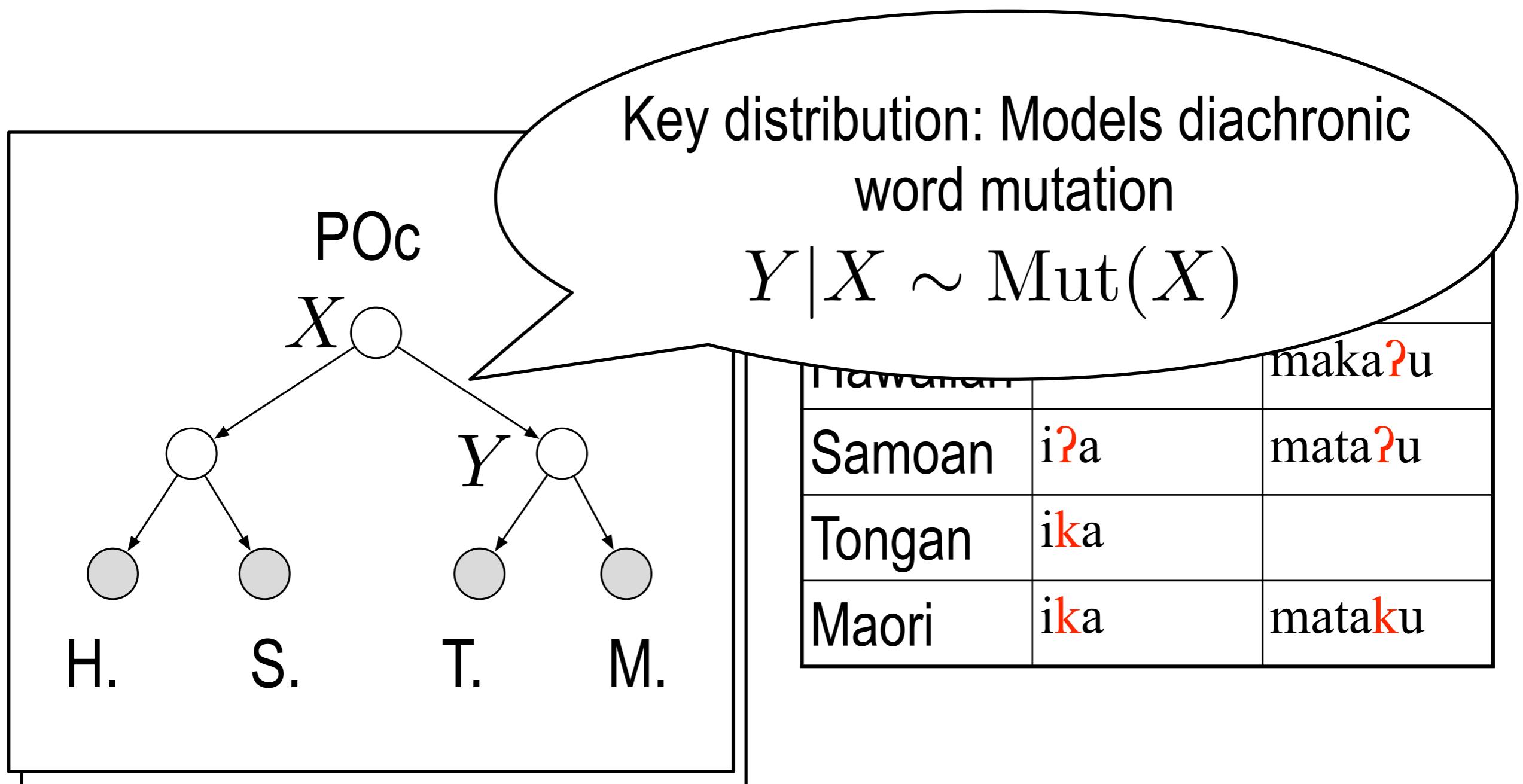
	'fish'	'fear'
Hawaiian	i?a	maka?u
Samoan	i?a	mata?u
Tongan	ika	
Maori	ika	mataku

Graphical model



	'fish'	'fear'
Hawaiian	i?a	maka?u
Samoan	i?a	mata?u
Tongan	ika	
Maori	ika	mataku

Graphical model

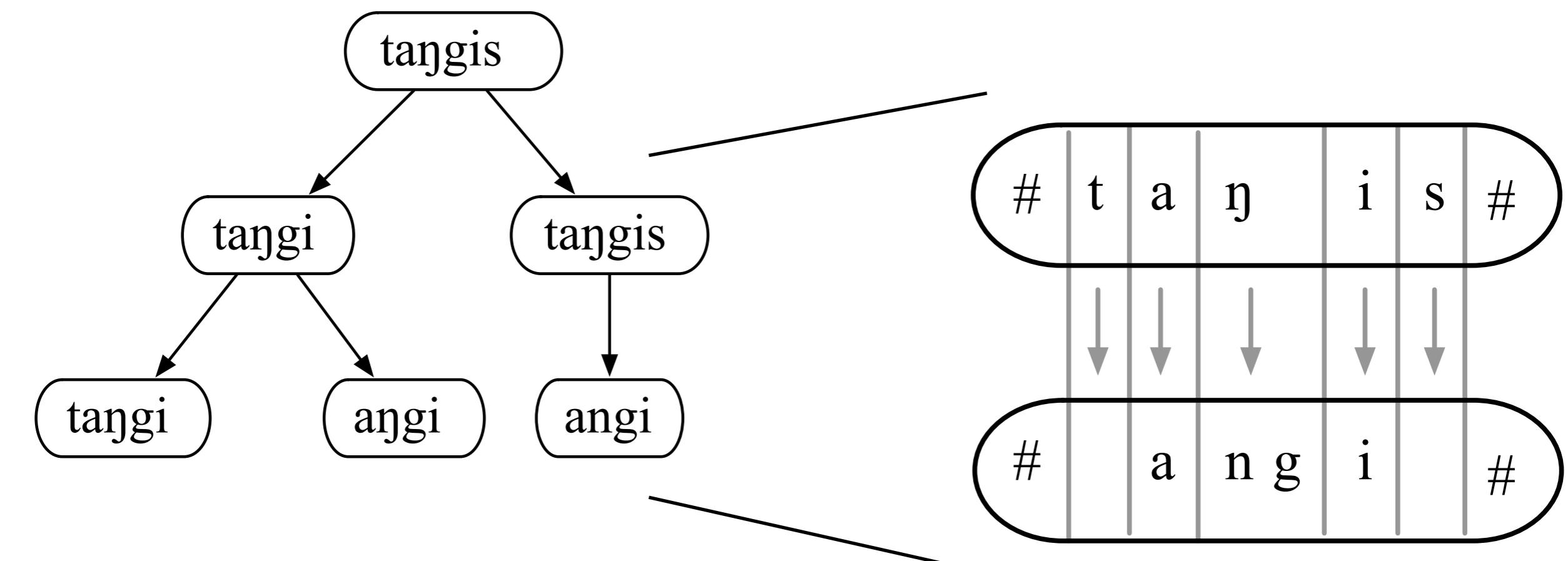


Modeling string mutation

- What kind of string mutations need to be captured?
 - Substitution
 - $*k > ?$
 - Deletion (& insertions)
 - $*? > \emptyset$
 - Contextualized change
 - $*? > \emptyset / V_V$

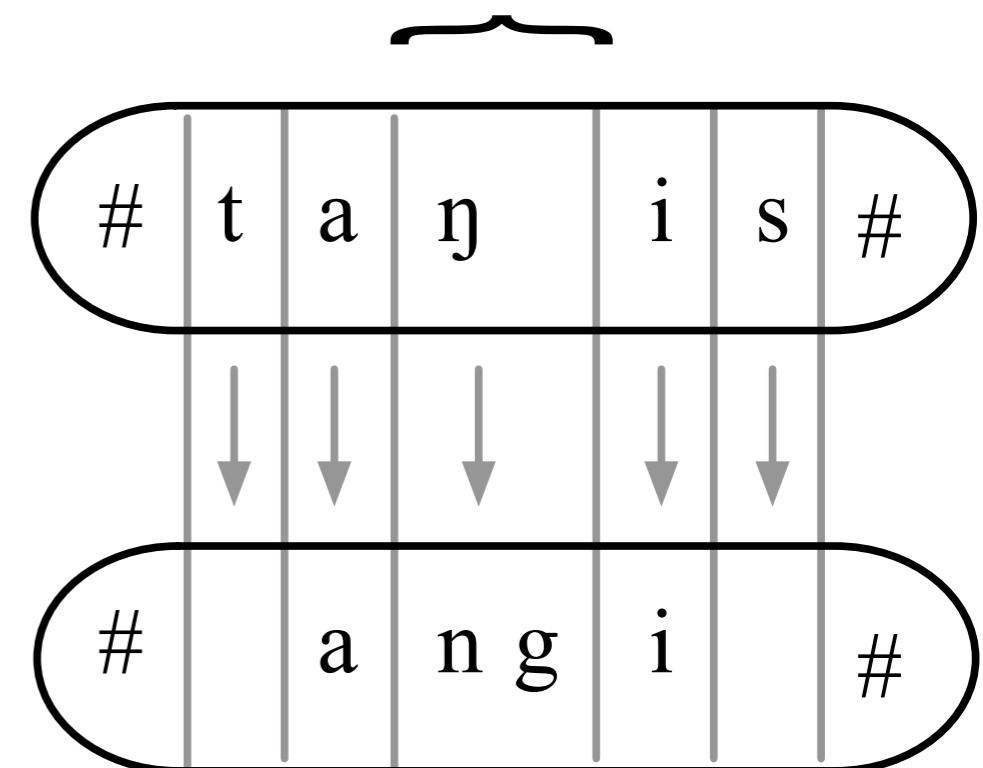
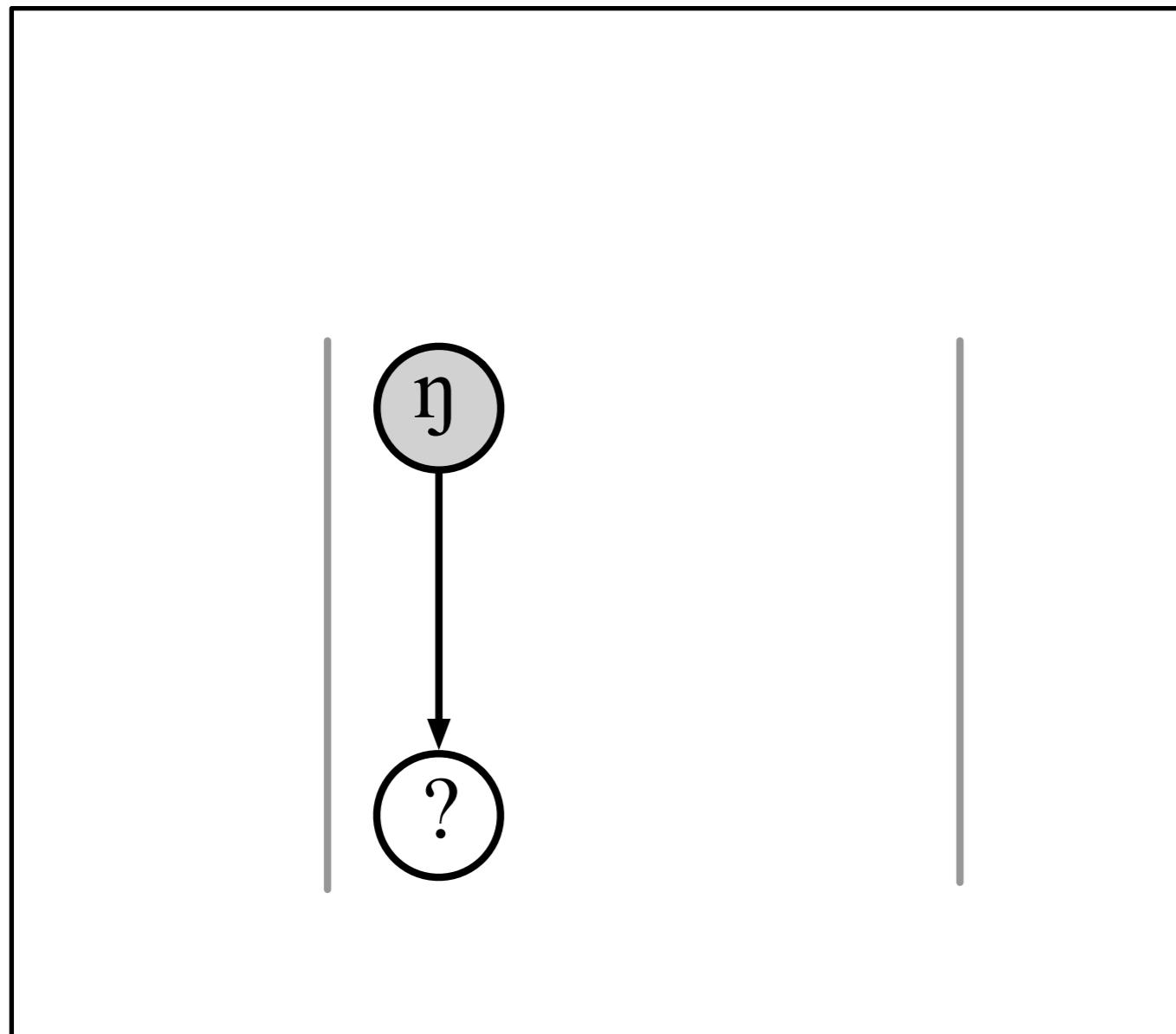
	'fish'	'voice'
Hawaiian	i?a	leo
Samoan	i?a	leo
Tongan	ika	le?o
Maori	ika	reo

String transducer



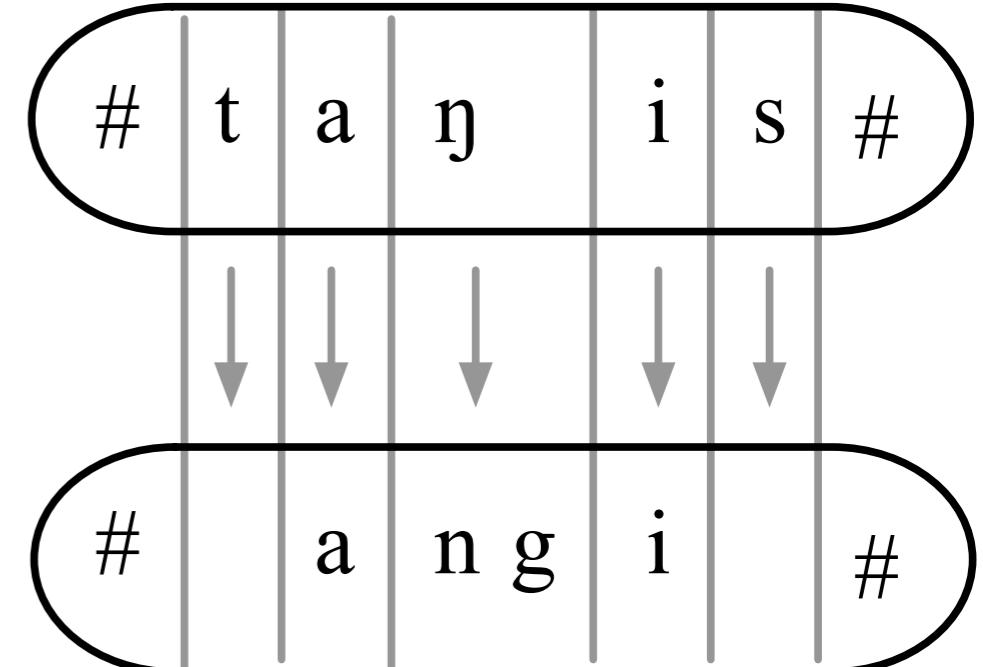
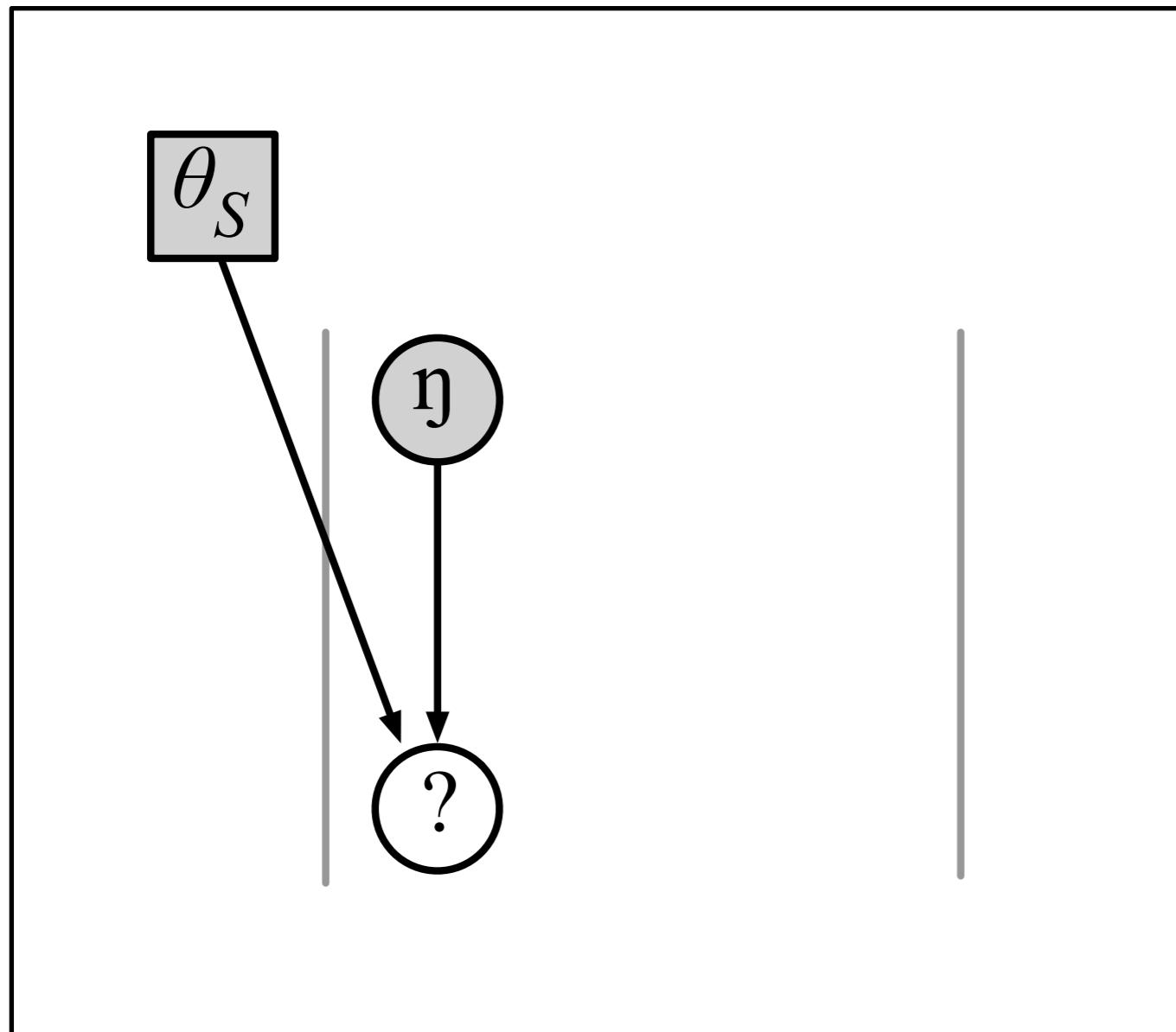
'to cry'

String transducer



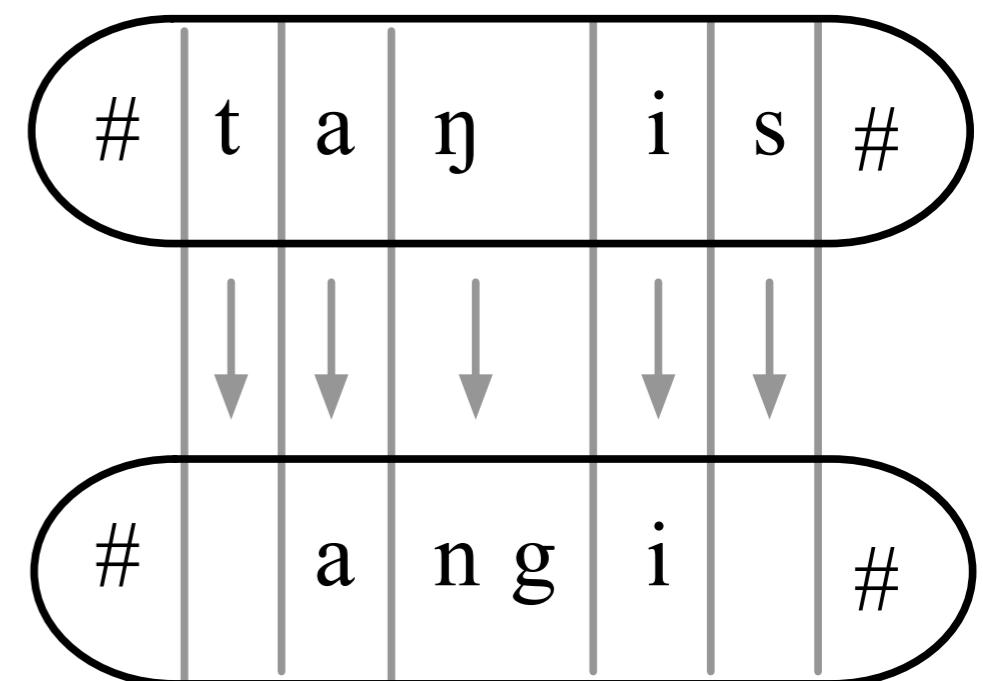
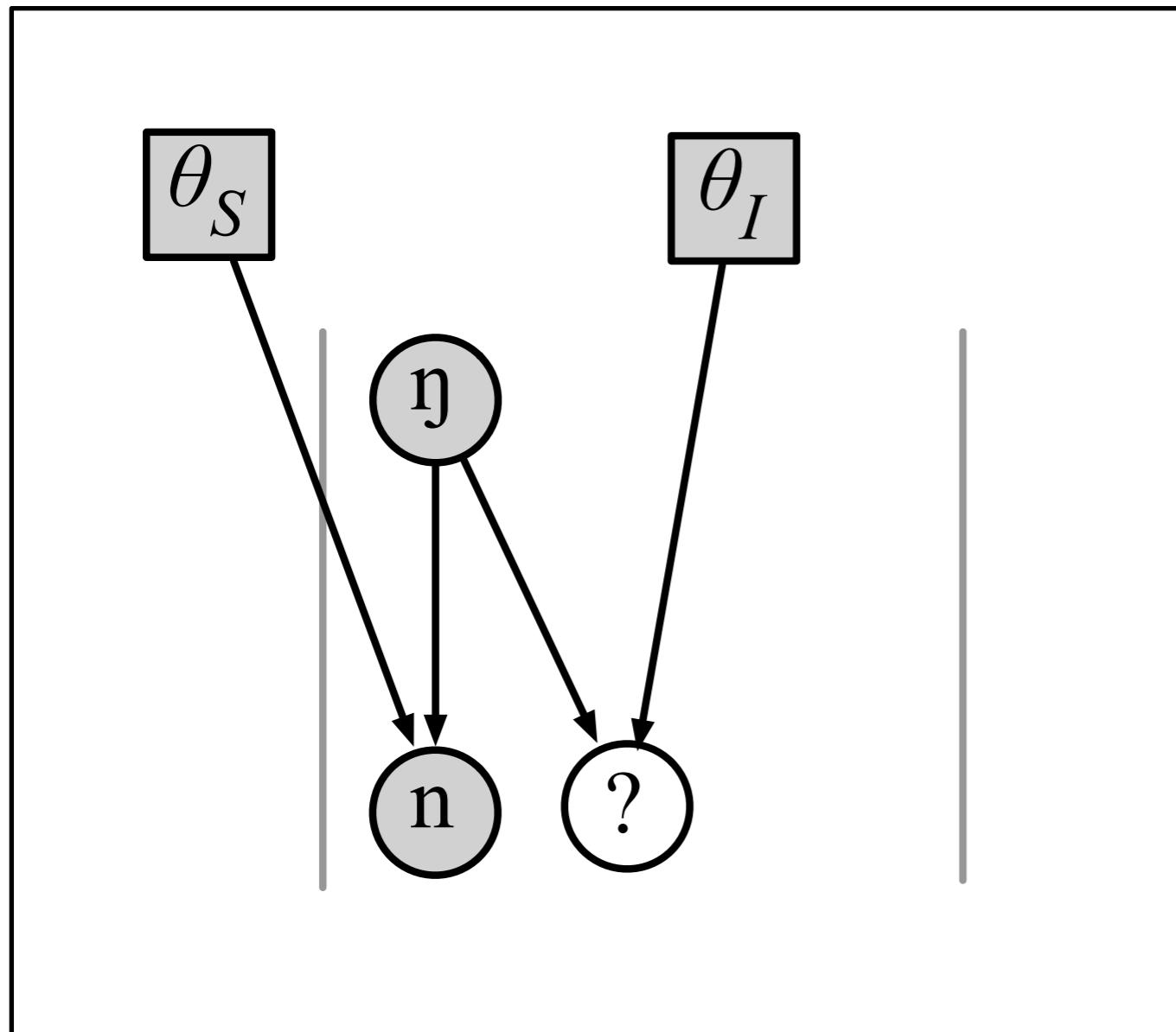
String transducer

θ_S : Substitution/Deletion
Parameters

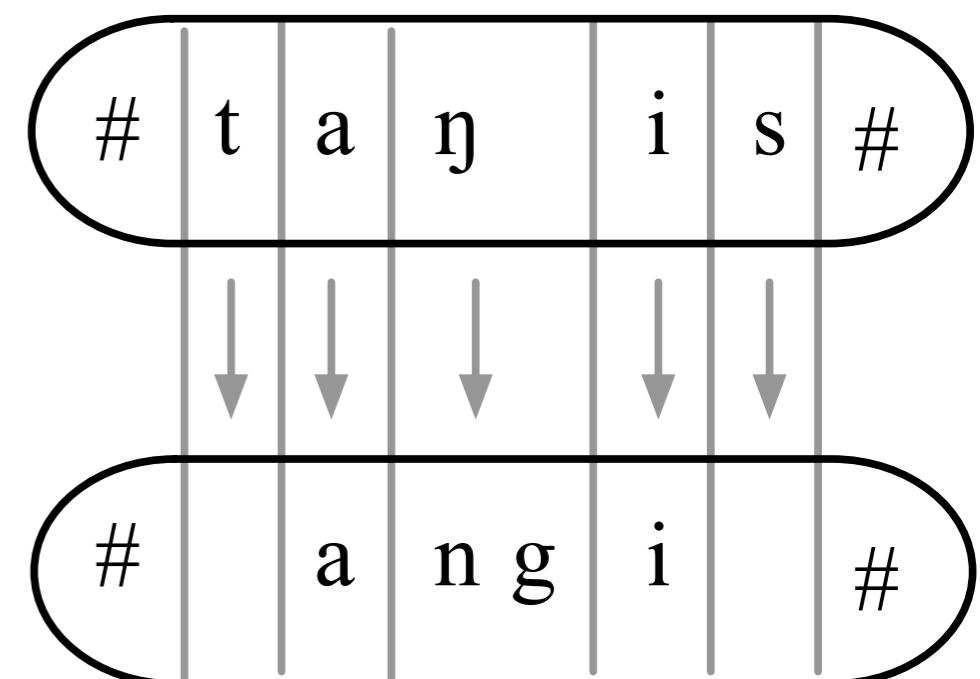
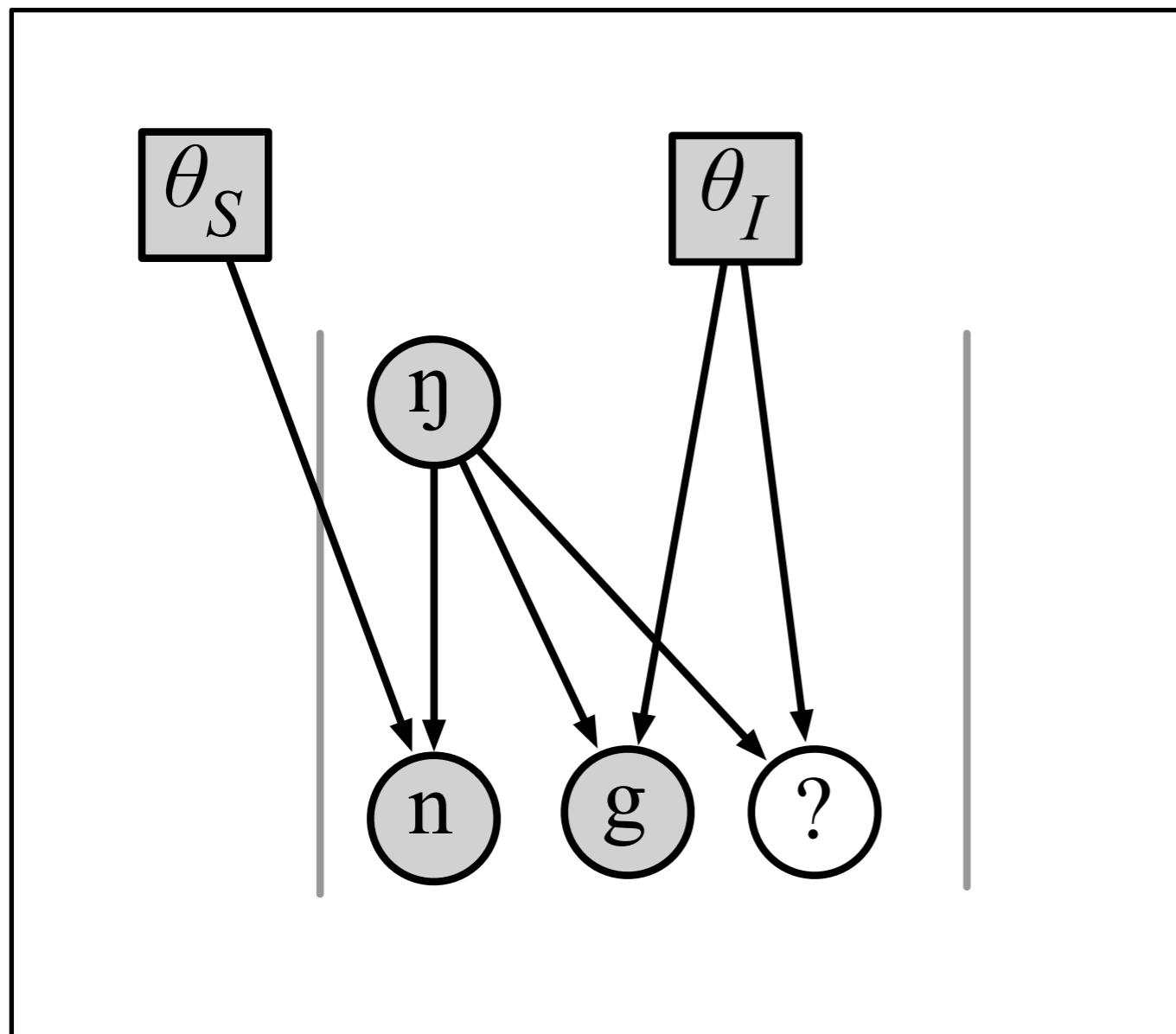


String transducer

θ_I : Insertion Parameters



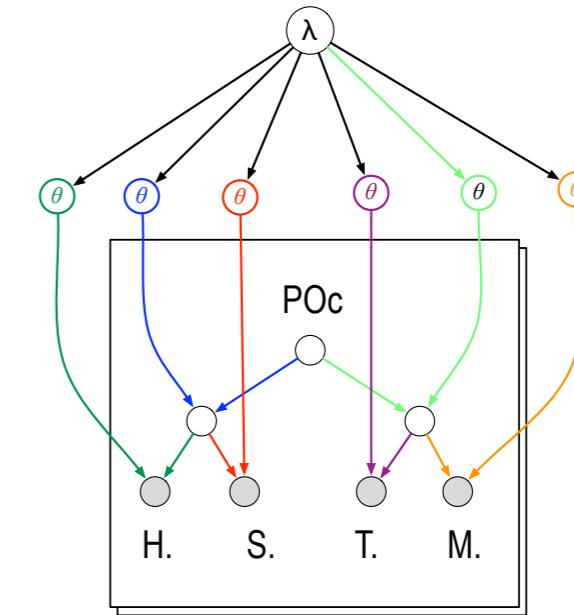
String transducer



Bayesian inference

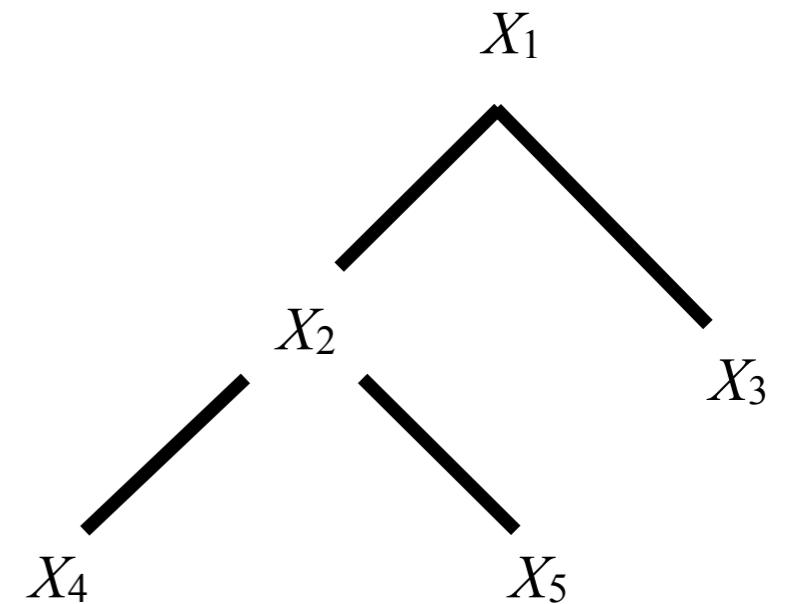
Computational bottleneck

Parameters fitting: can be approached using standard tools



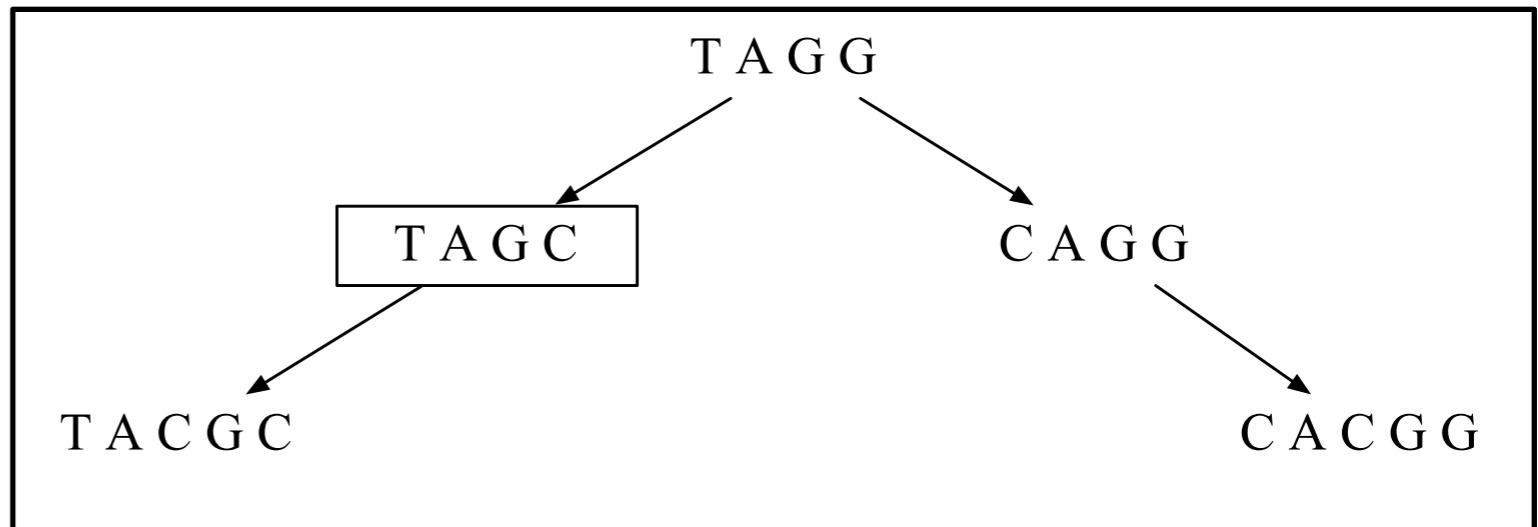
More difficult: string-valued expectations

$$X_1, X_2 | X_3, X_4, X_5$$

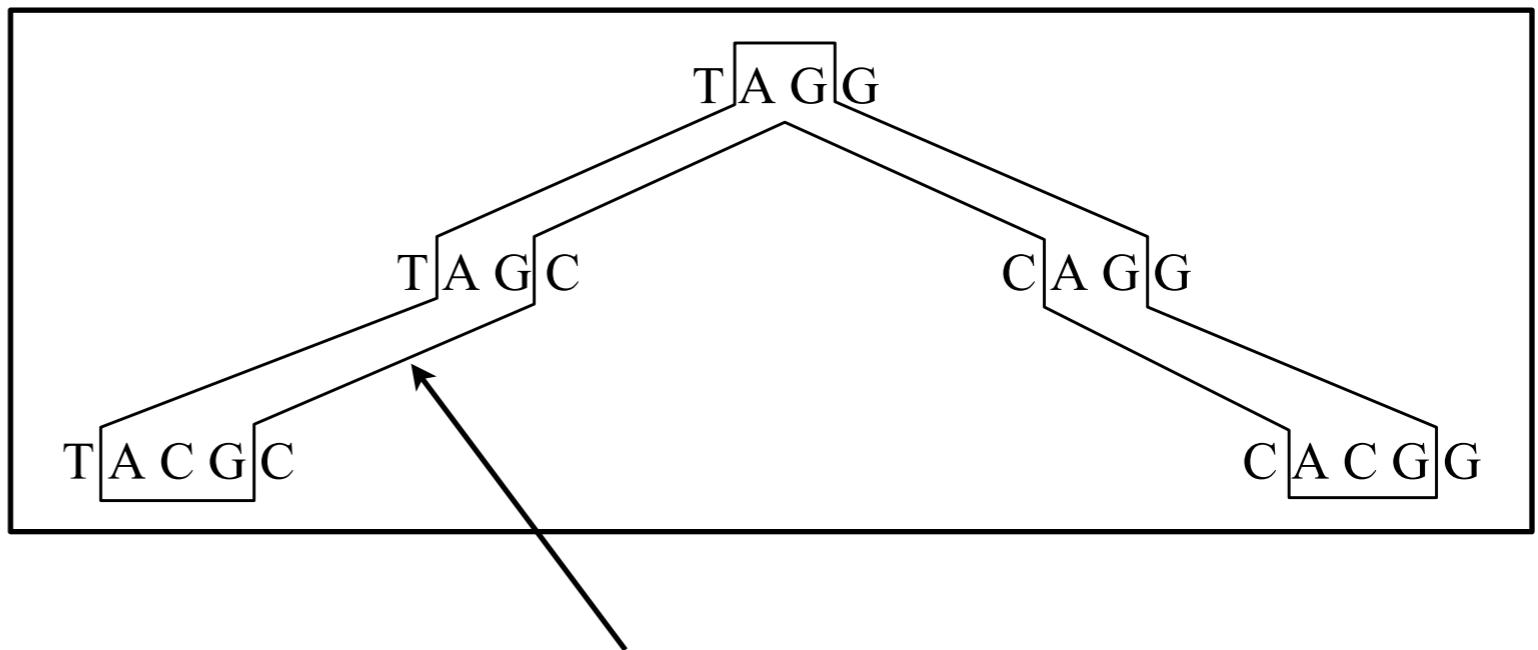


Various MCMC algorithms

Gibbs sampling



Ancestry resampling



Thin vertical section

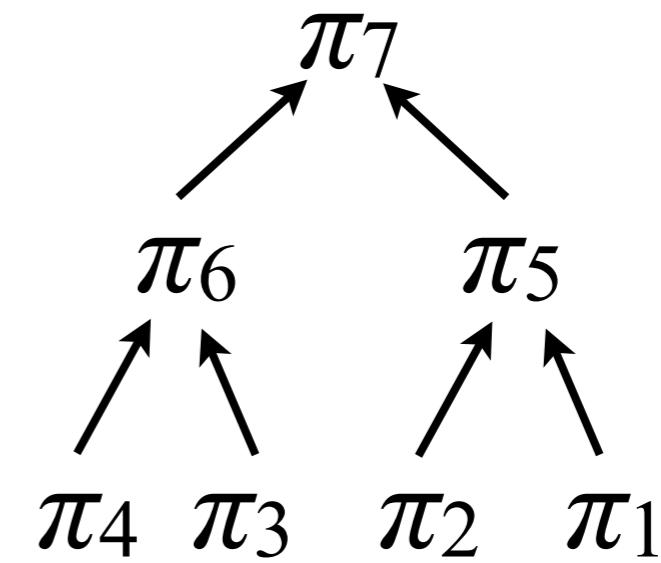
Sequential Monte Carlo

- SMC: An alternative to MCMC
 - Approximates posterior with particles and weights
- D&C SMC: generalization more suitable to phylogenies
 - tinyurl.com/divide-and-conquer-smc
 - Fredrik Lindsten, Adam M. Johansen, Christian A. Naesseth, Bonnie Kirkpatrick, Thomas B. Schön, John Aston, and Alexandre Bouchard-Côté. (2016) Divide-and-Conquer with Sequential Monte Carlo. JCSG, In Press.

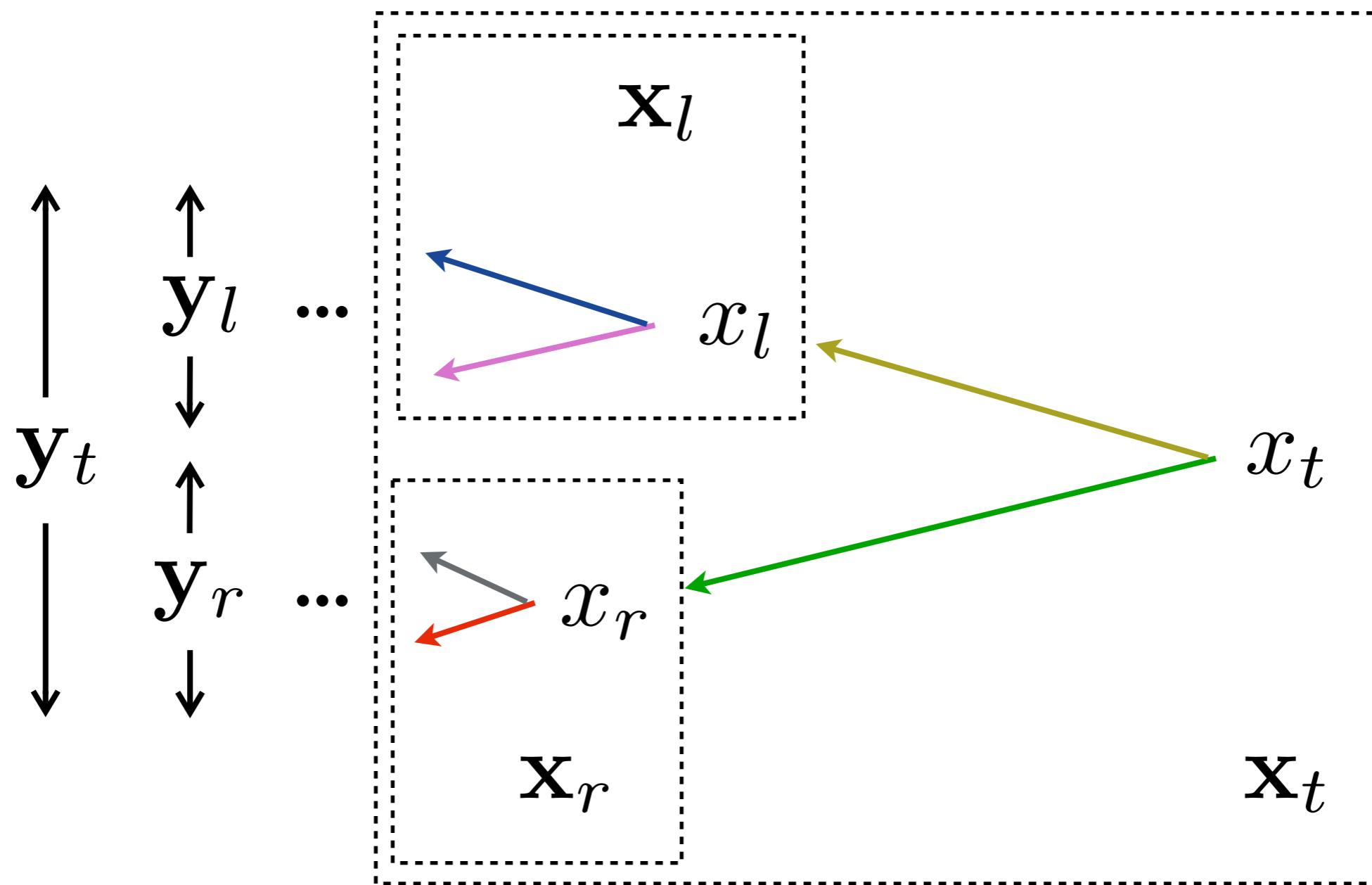
Standard SMC samplers

$$\pi_1 \longrightarrow \pi_2 \longrightarrow \pi_3 \longrightarrow \dots$$

Divide and Conquer
(D&C) SMC samplers

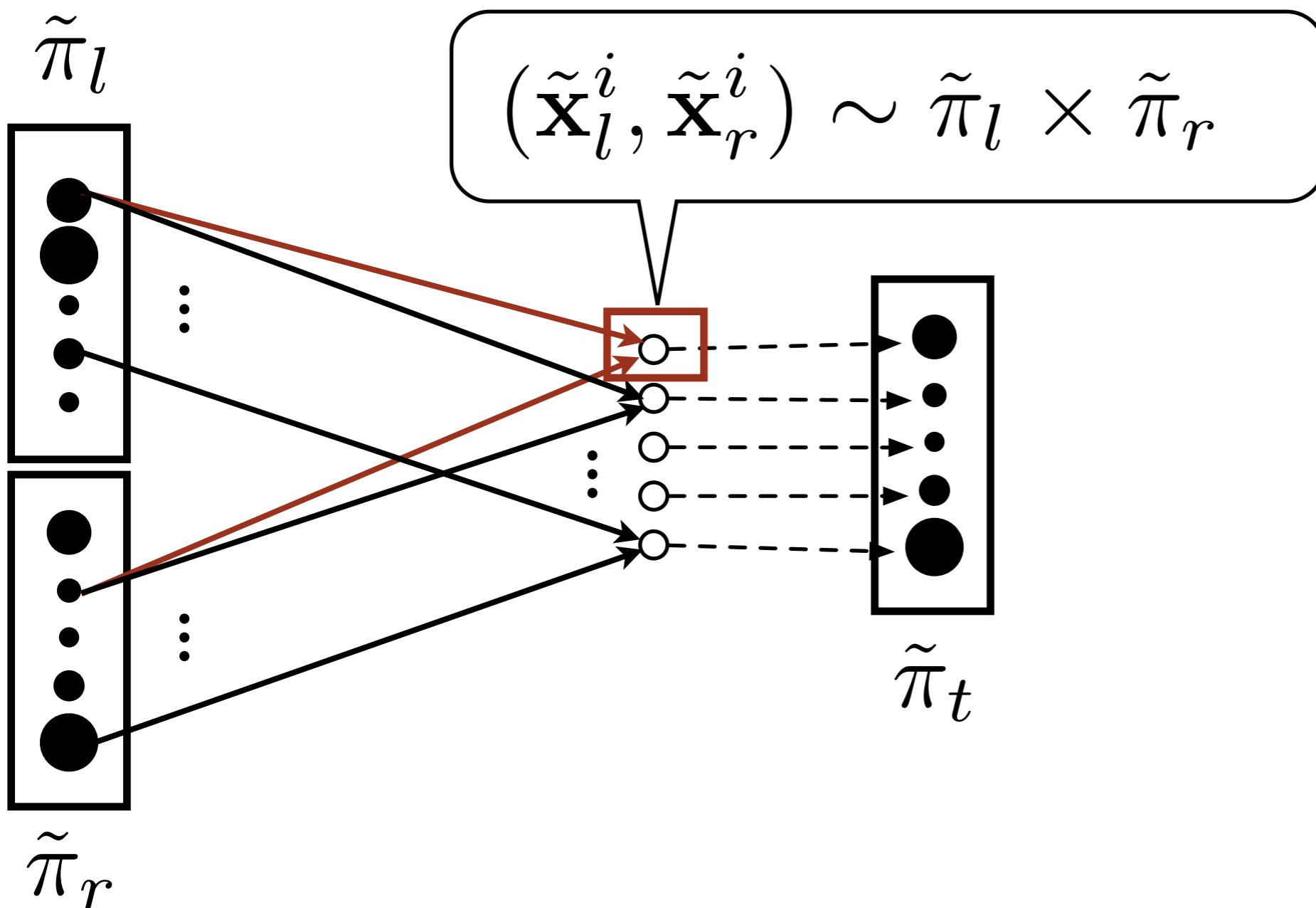


D&C SMC: notation



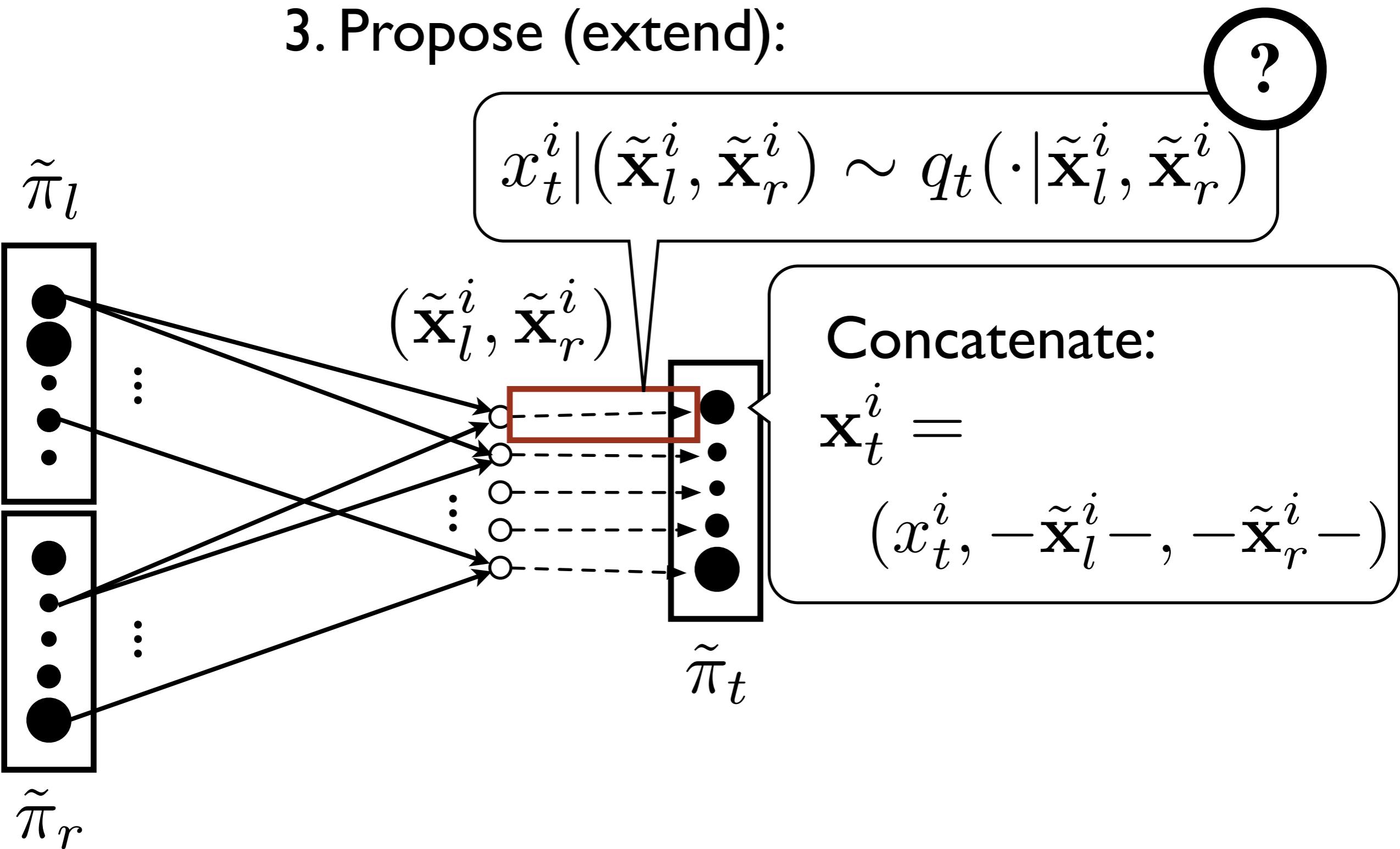
D&C (SIR):

1. Assume inductively we have $\tilde{\pi}_l$ & $\tilde{\pi}_r$
2. Sample from the product measure:
 - amounts to two independent multinomial sampling steps



D&C (SIR):

1. Assume inductively...
2. Sample from the product measure:
3. Propose (extend):

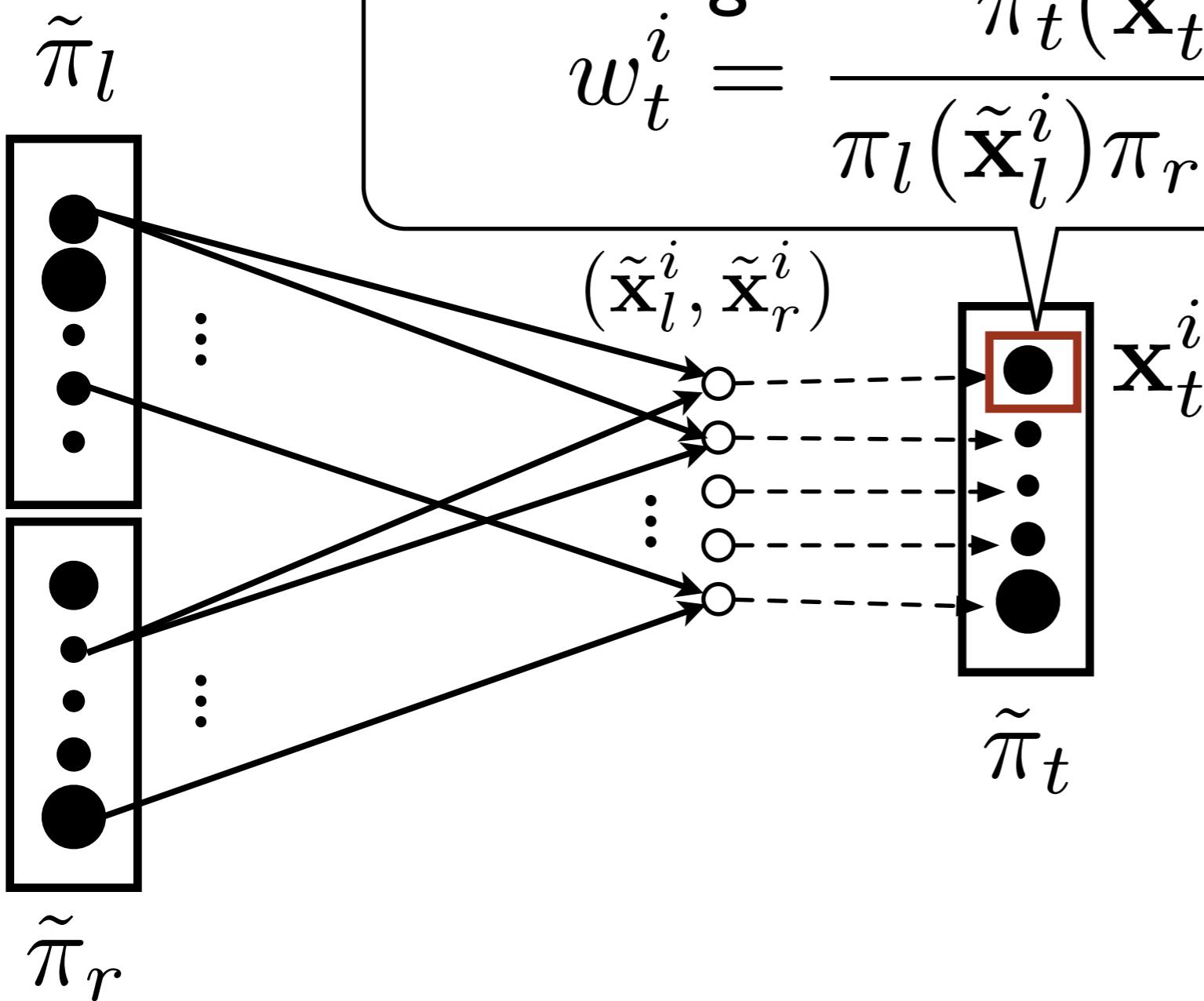


D&C (SIR):

1. Assume inductively...
2. Sample from the product measure
3. Propose (extend)

4. Reweigh:

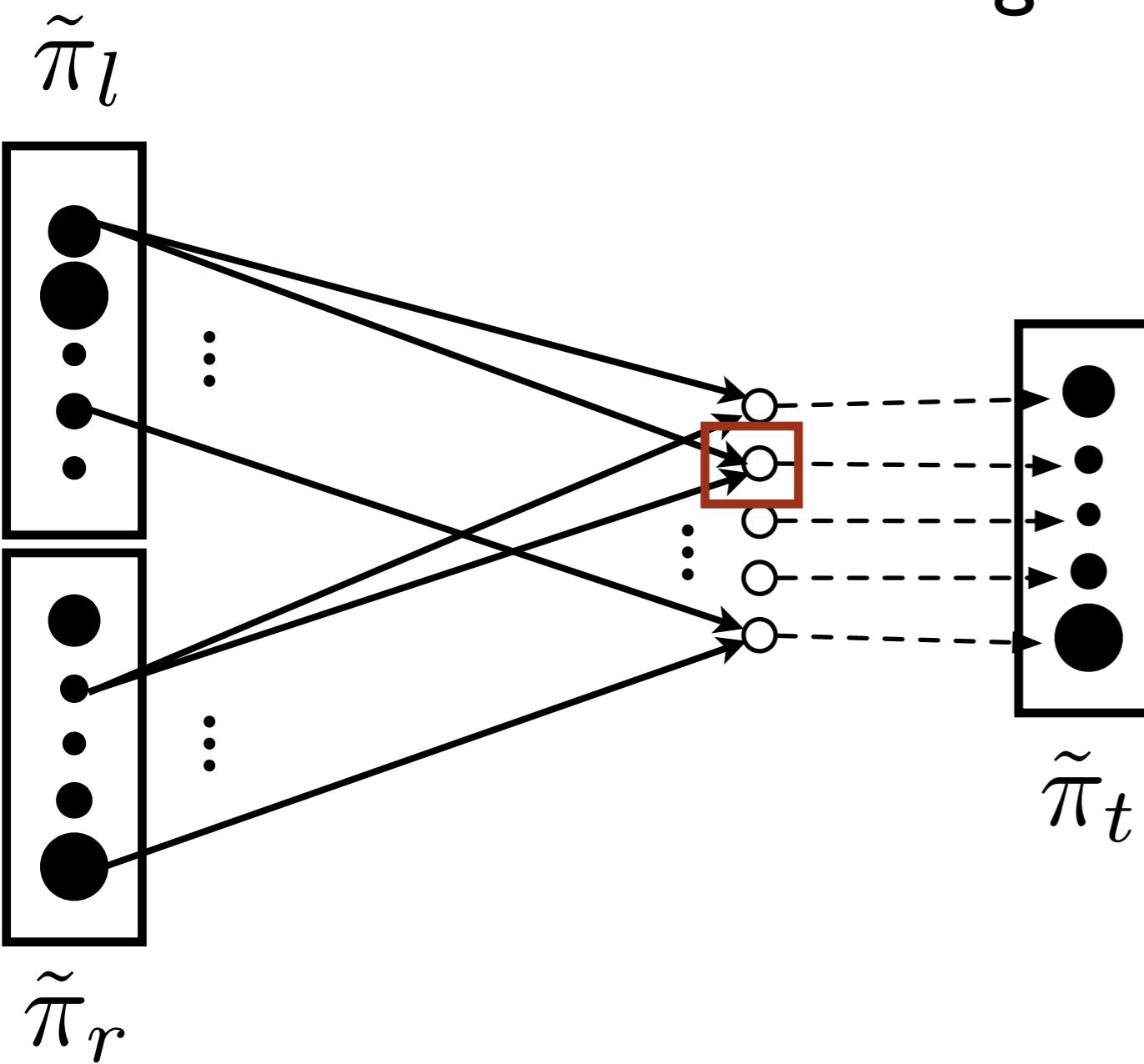
$$w_t^i = \frac{\pi_t(\mathbf{x}_t^i)}{\pi_l(\tilde{\mathbf{x}}_l^i)\pi_r(\tilde{\mathbf{x}}_r^i)} \frac{1}{q_t(x_t^i | \tilde{\mathbf{x}}_l^i, \tilde{\mathbf{x}}_r^i)}$$



D&C (SIR):

Repeat for
each particle

1. Assume inductively...
2. Sample from the product measure
3. Propose (extend)
4. Reweigh



Potential applications

- Inference over non-local/catastrophic events (explicit sound change operators)
- Joint cognate inference
- Statistical multiple sequence alignment
- Relaxed clock models
- Many other models that would otherwise have to be handled with very expensive methods such as Approximate Bayesian Computation

Numerical examples

Validation methodology

Holding-out the manual reconstructions

	'fish'	'fear'
Hawaiian	i?a	maka?u
Samoan	i?a	mata?u
Tongan	ika	
Proto-Oceanic	*ika	*mataku



	'fish'	'fear'
Hawaiian	i?a	maka?u
Samoan	i?a	mata?u
Tongan	ika	
Proto-Oceanic	???	???

Evaluation criterion: Edit distance

Smallest number of substitutions, insertions and deletions needed to go from one string to the other

$$\text{e.g.: } d(/ika/ , /ga/) = 2$$

$/ika/ \rightarrow /iga/ \rightarrow /ga/$

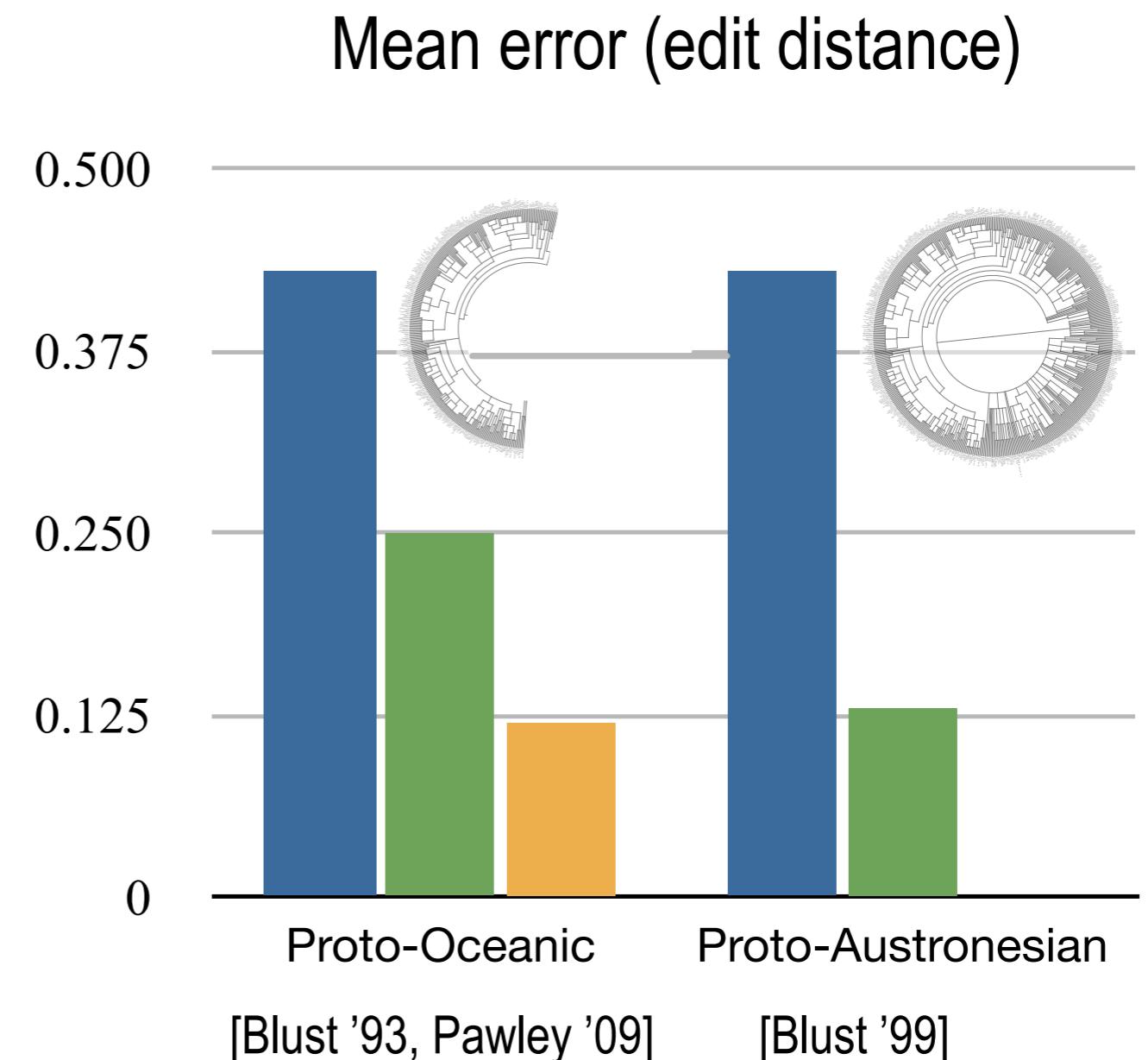
Datasets used for reconstruction

- Austronesian languages dataset [Greenhill et al. 08]
 - 706 languages, 150k words
 - 17 annotated reconstructions
 - Ideal for testing large scale reconstruction
 - Has continued to grow since the time experiments were performed

Comparisons in large phylogenies

Comparisons on Proto-Oceanic and Proto-Austronesian

- Baseline: sampling from the observed words
- This work
- Distance between two linguists' reconstructions



Examples of reconstruction

Gloss	Known Modern Languages				Reconstructed Ancestors		
	Fijian	Pazeh	Melanau	Inabaknon	Manual	Automated	Δ
star	kalokalo	mintol	biten	bitu'on	*bituqen	*bituqen	0
to hold	taura	ma.ra?	magem	kumkom	*gemgem	*gemgem	
house	vale	xuma?	lebu?	ruma	*Rumaq	*Rumaq	
bird	manumanu	aiam	manuk	manok	*qayam	*qayam	
to cut, hack	tata	ta:tatak	tutek	hadhad	*taraq	*taraq	
at	e	N/A	ga?	N/A	*i	*i	
what?	cava	?axai	ua? inew	ay	*nanu	*anu	1
this	oqo	?imini	itew	yayto	*ini	*ani	
wind	cagi	varə	panjay	bariyo	*bali	*beliu	2

Colors here
indicate
cognate sets

Stratified
sampling by edit
distance