# Scalable likelihood-based methods to infer lineages and estimate selection in B cell repertoires

Frederick "Erick" Matsen
Fred Hutchinson Cancer Research Center
http://matsen.fredhutch.org/
@ematsen

with *Trevor Bedford (FH), Vladimir Minin (UW),* **Duncan Ralph** *(FH) and David Shaw (FH)*
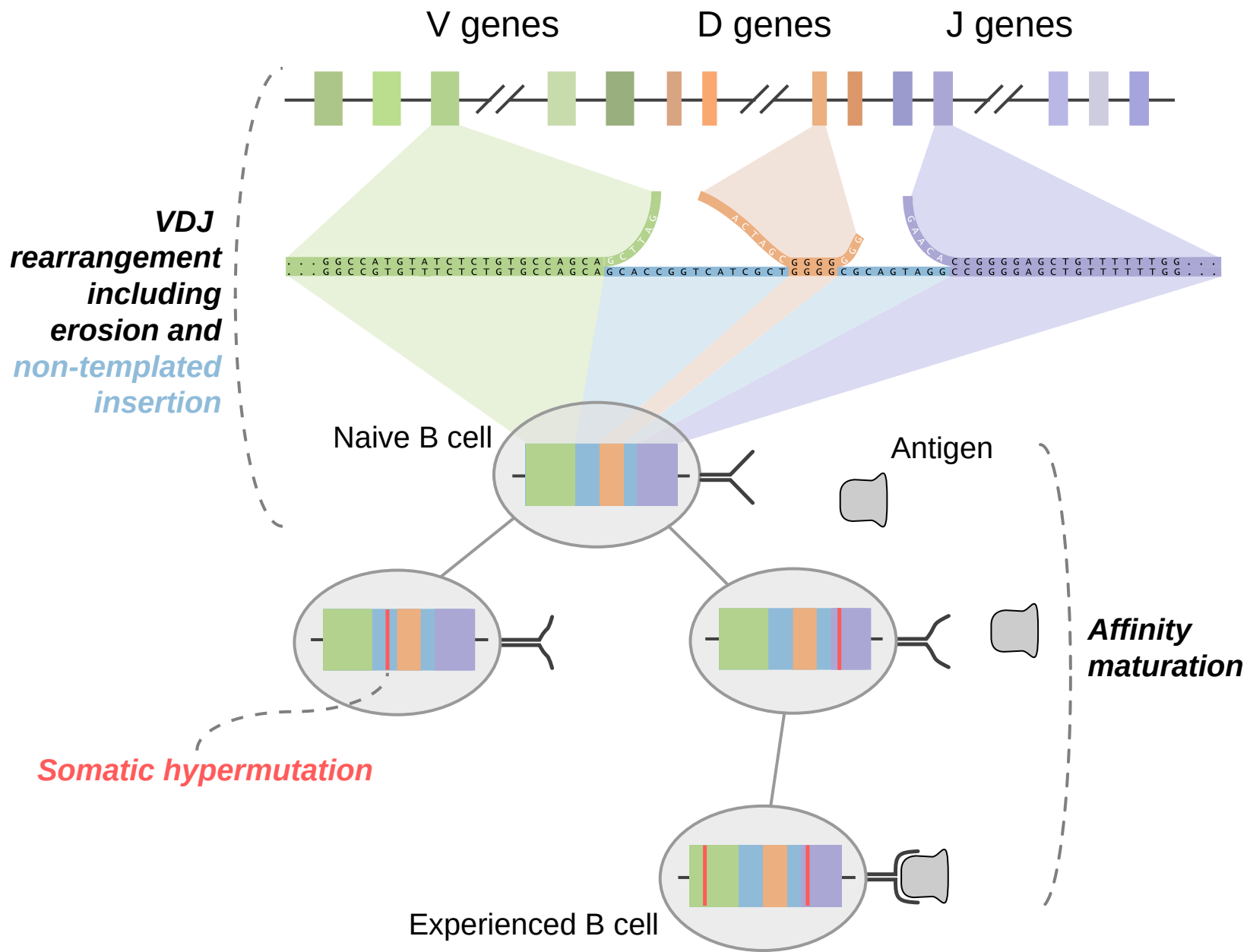
# Philosophy of talk

- Model immune cells (in this case B cells) probabilistically
- Infer parameters describing process via likelihoods
- Use these parameters to improve sequence analysis.
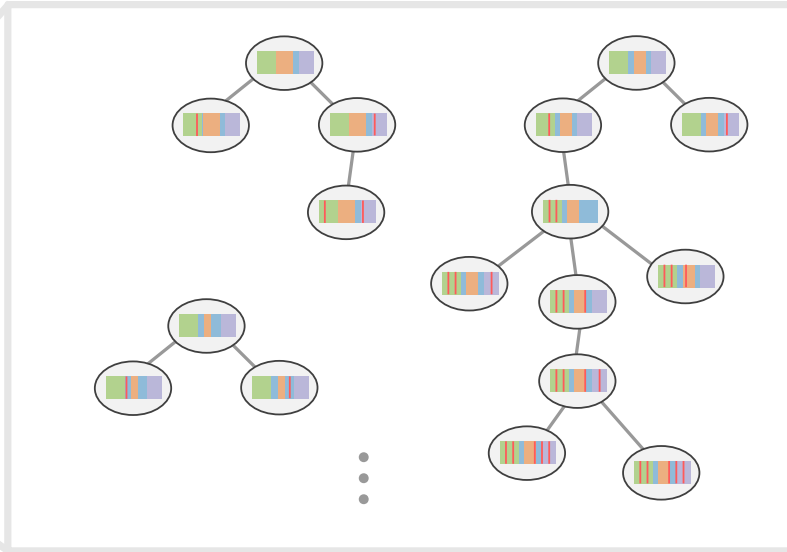
**Work in progress.**

# Statistical phylogenetics

- Develop probabilistic model for sequence evolution
- Write down likelihood function
- Search for the maximum likelihood tree, including optimization heuristics
- Or integrate over trees using MCMC.

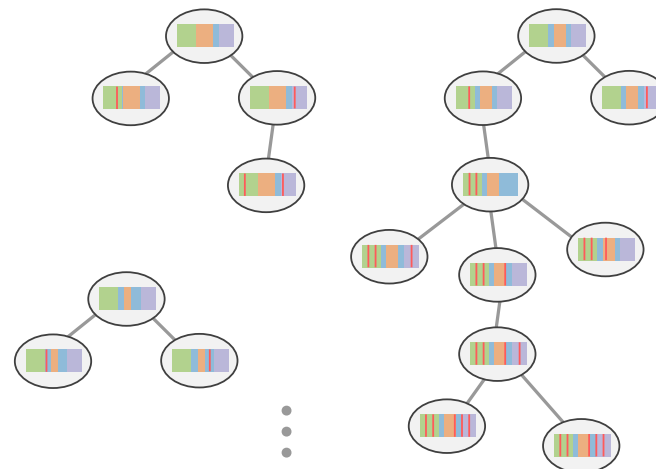Casting phylogenetics as a statistical inference problem provides a solid foundation.

V genes    D genes    J genes

**VDJ rearrangement including erosion and** *non-templated insertion*

Naive B cell

Antigen

*Somatic hypermutation*

*Affinity maturation*

Experienced B cell

reality

inference

ACATGGCTC...
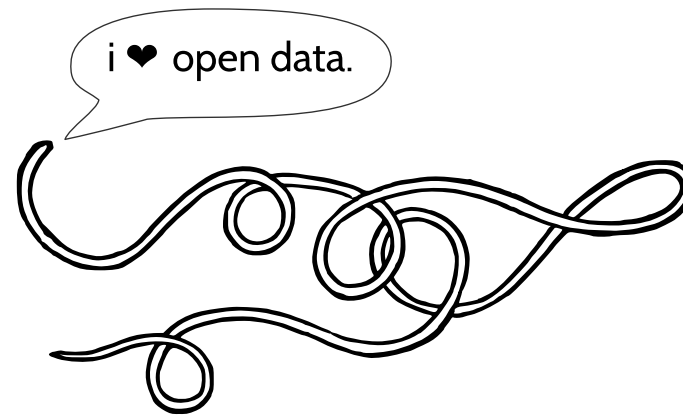ATACGTTCC...
TTACGGTTC...
ATCCGGTAC...
ATACAGTCT...

# To-do list

0. [Generate high-quality data. Hard!]
1. Annotate BCR sequences
2. Find clonal families
3. Reconstruct BCR phylogenetic trees
4. Infer BCR ancestral sequences
5. Evolutionary selection inference for BCRs
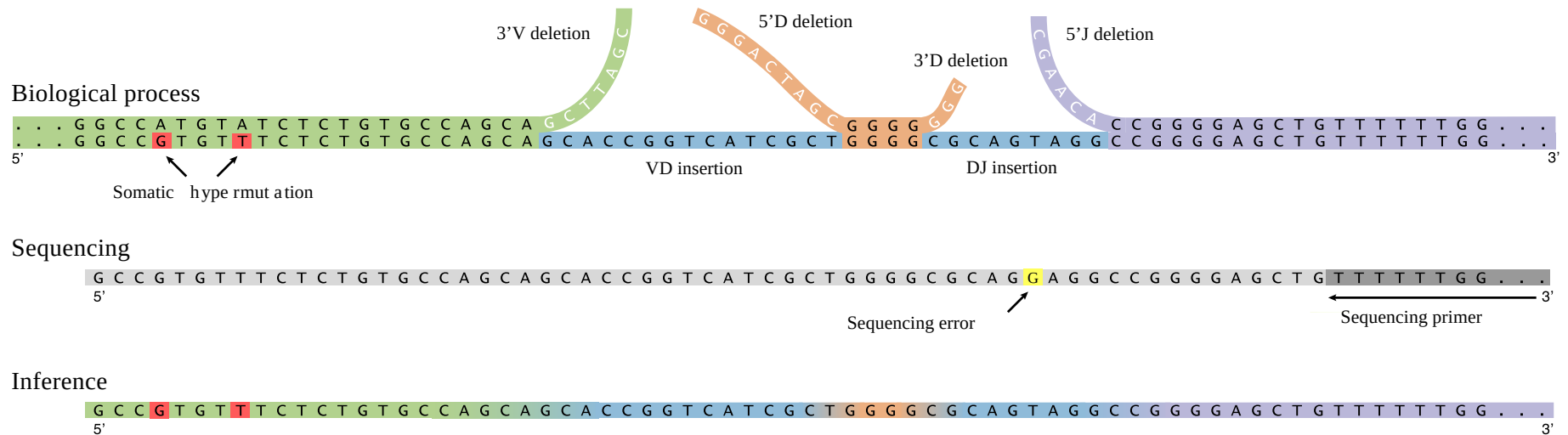
… in a probabilistic framework.

# 0. Gather data



- Data from Adaptive Biotechnologies: 3 healthy individuals, naive/memory sorted, replicate immunosequencing with 188 wells and ~50K cells/well http://adaptivebiotech.com/link/mat2015
- Stern, Yaari, Heiden … O'Connor (2014). B cells populating the multiple sclerosis brain mature in the draining cervical lymph nodes. *Science Translational Medicine*, 6(248), 248ra107.
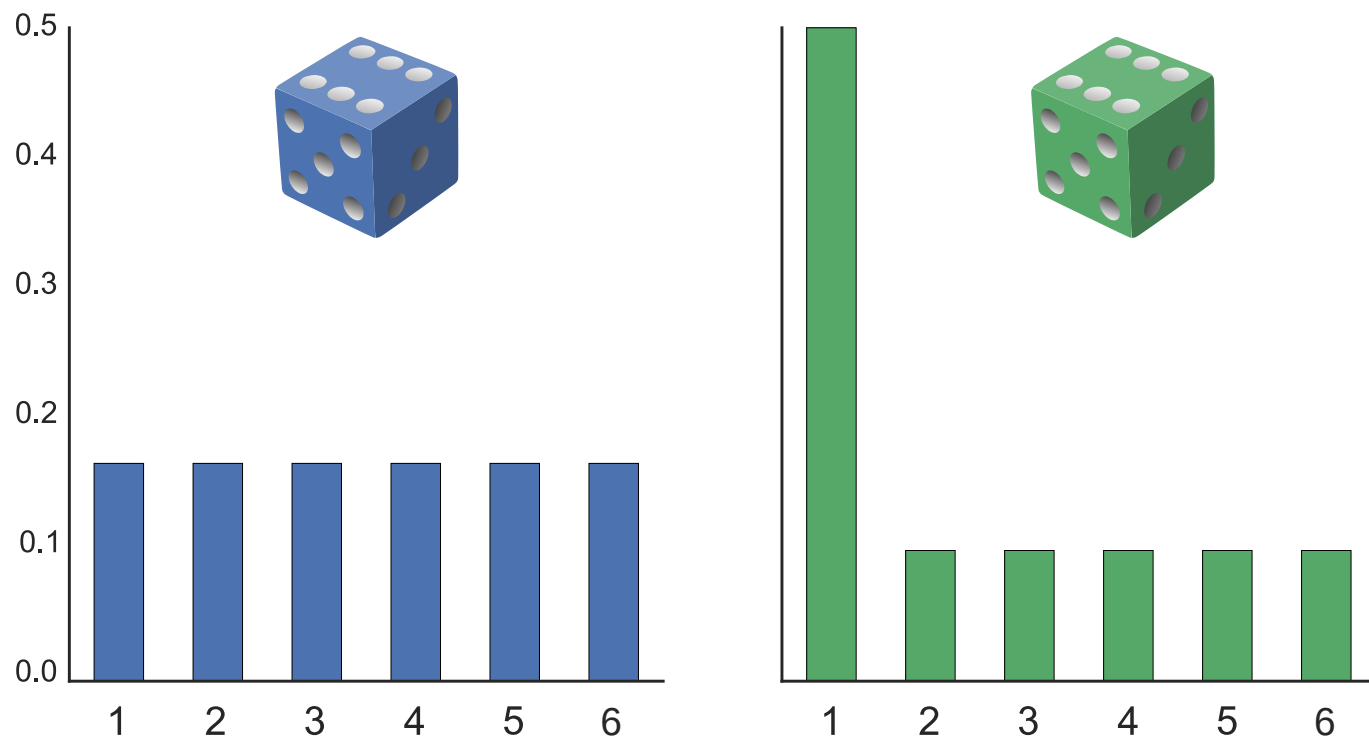
# 1. *Annotate BCR sequences*
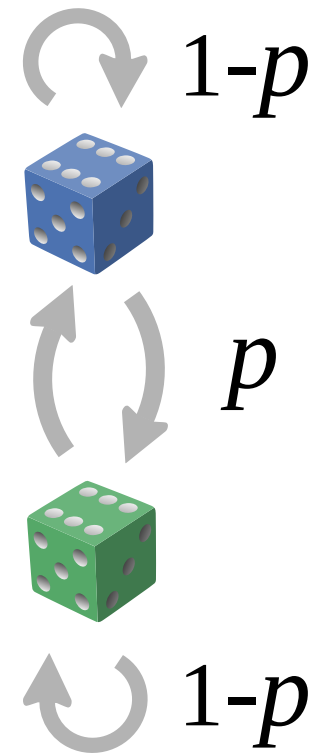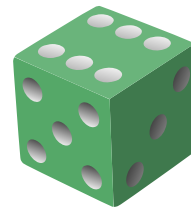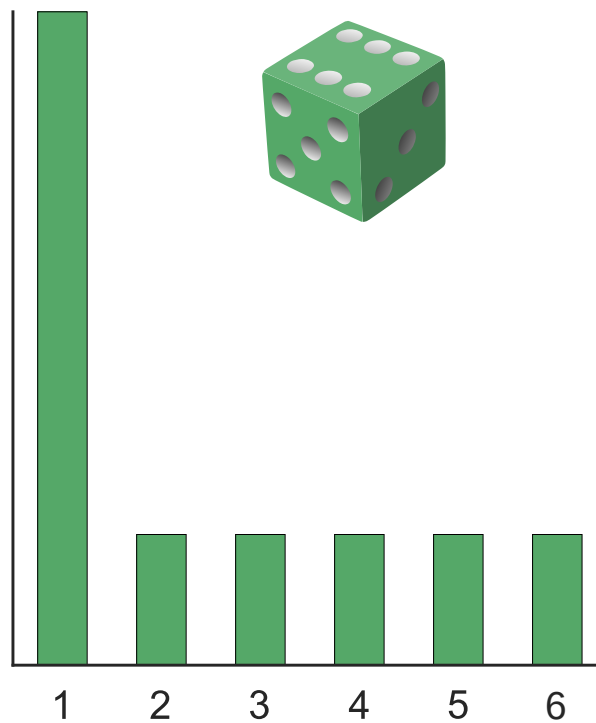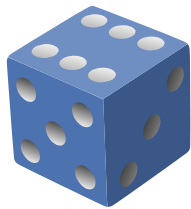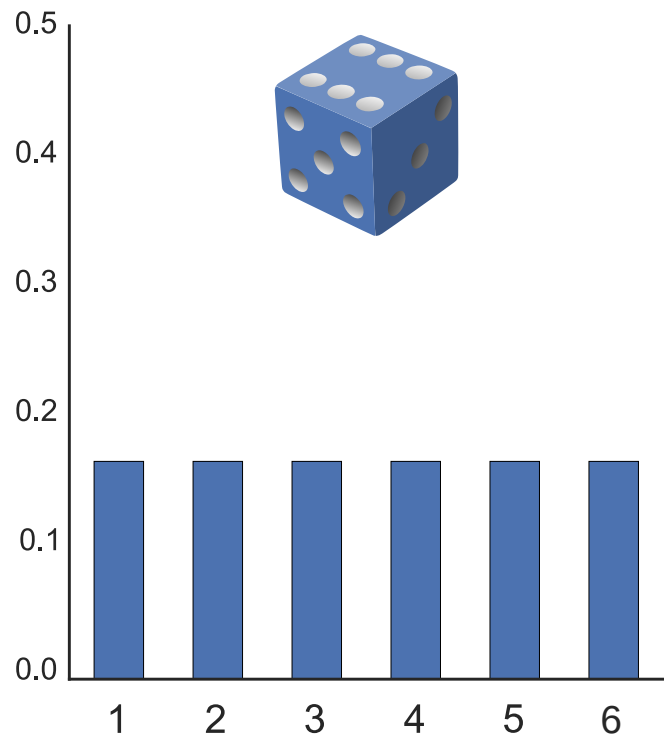# from where did each nucleotide come?



This is a key first step in BCR sequence analysis.
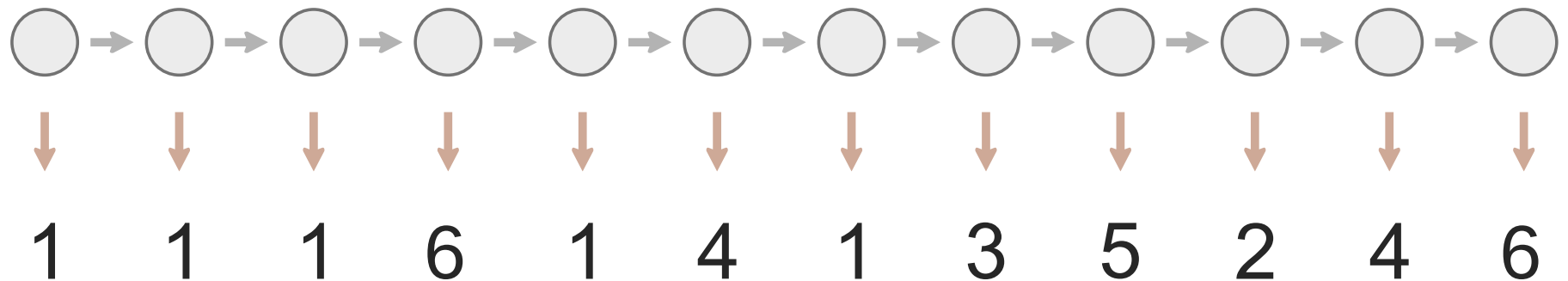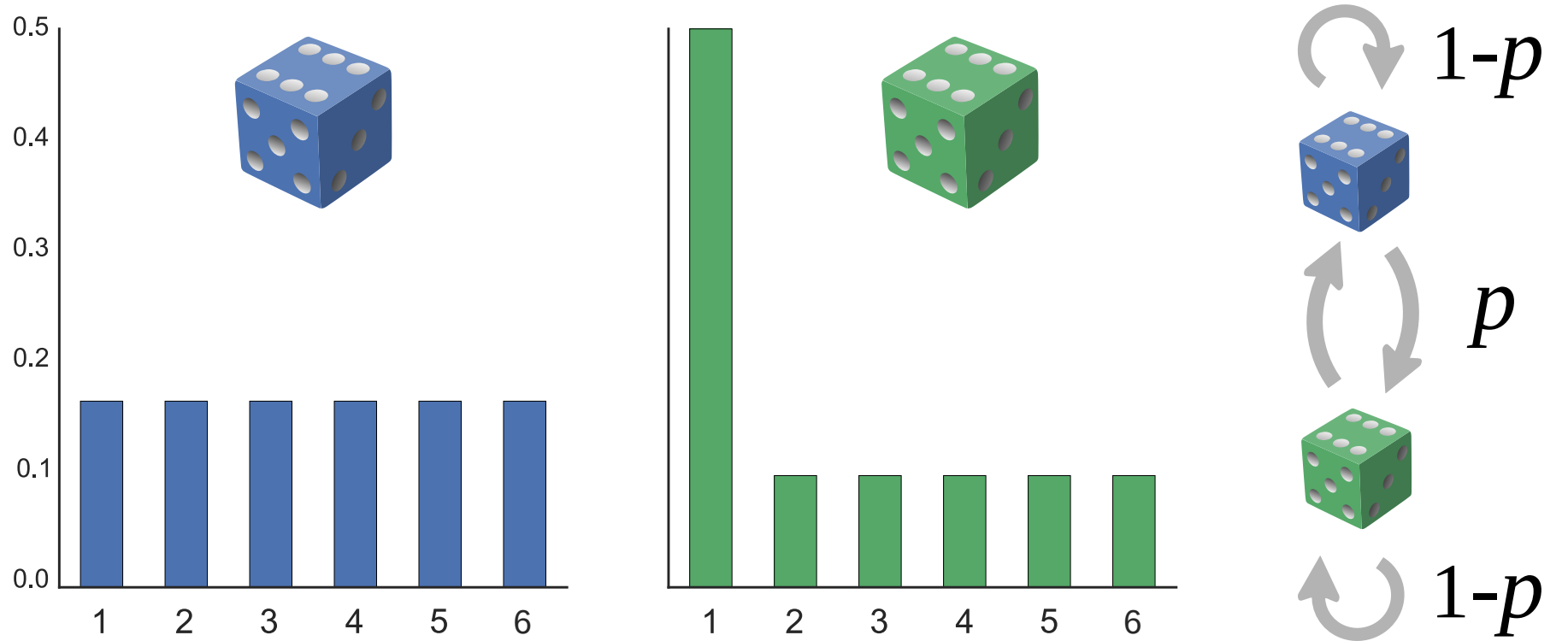
# HMM intro: dishonest casino

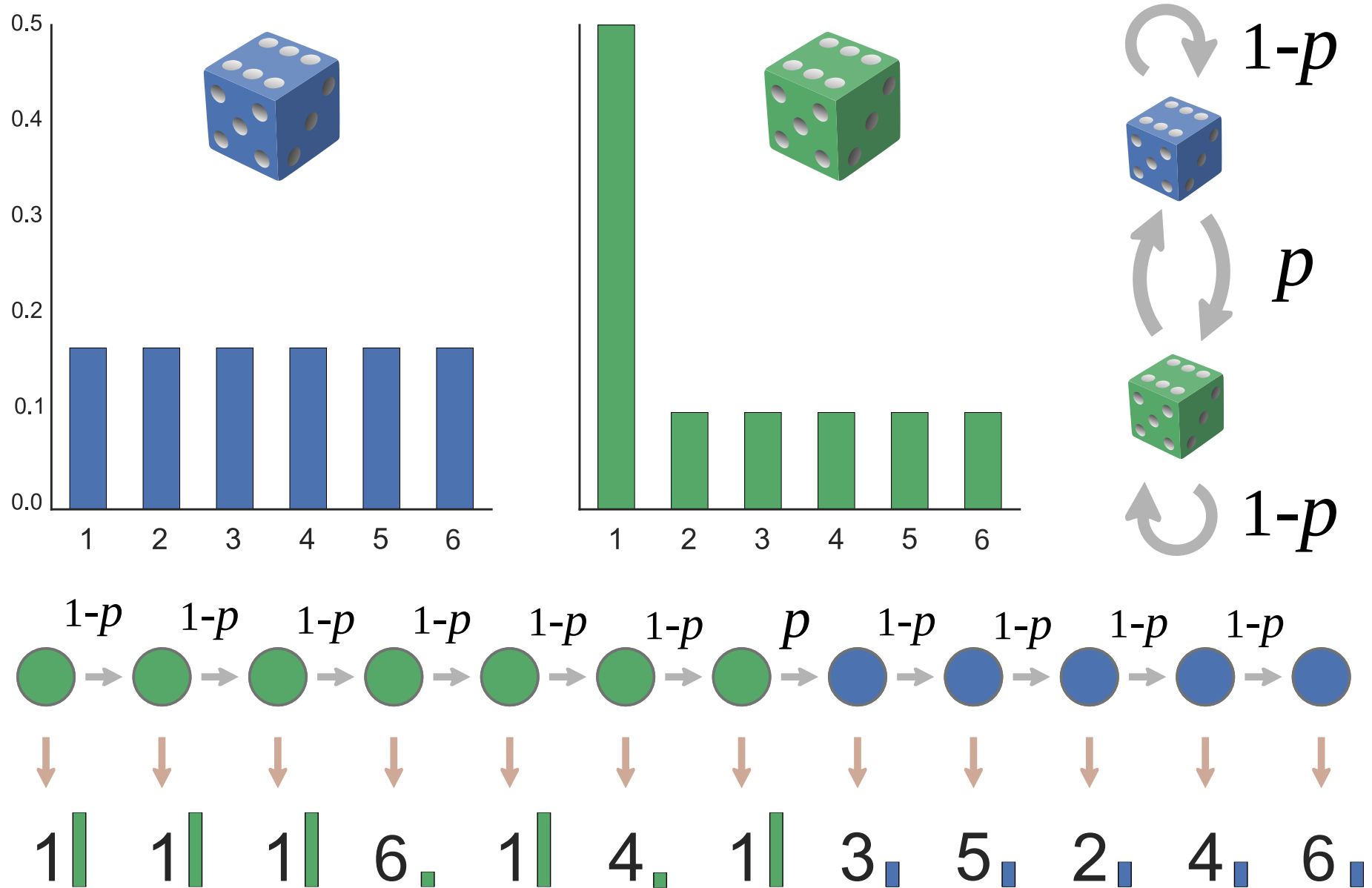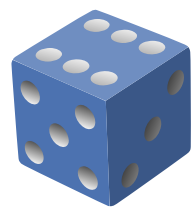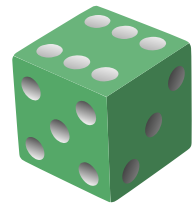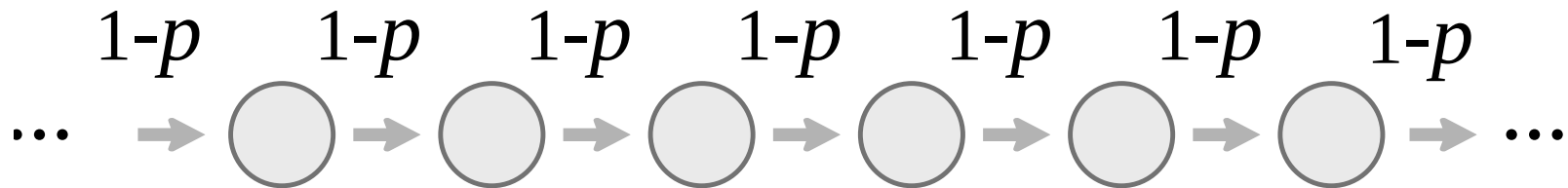# HMM intro: dishonest casino

# HMM intro: dishonest casino

HMM intro: dishonest casino

# "Thread" reads onto structure

# Distributions are reproducibly weird!

# Murugan, Mora, Walczak, Callan (2012)

Incorporating model complexity leads to better inferences

Legend: partis (k=5), partis (k=1), ighutil, iHMMunealign, igblast, imgt

frequency vs. Hamming to unmutated ancestor

# HMMs for BCR annotation

- SoDA: Volpe, Cowell, & Kepler (2005)
- iHMMune-align: Gaëta, Malming, … & Collins (2007)
- SoDA2: Munshaw & Kepler (2010)

*These implementations use standard probability distributions for parameters (e.g. deletion lengths).*

Fit parameter-rich HMMs that are able to capture underlying complexity of the process.

- partis: Ralph & M. (2016)
- repgenHMM: Elhanati, Marcou, Mora & Walczak (2016)

# IgSCUEAL: Frost, Murrell, ... K. Pond (2015)

# 2. Find clonal families



reality

inference of clonal families

ACATGGCTC...
ATACGTTCC...
TTACGGTTC...
ATCCGGTAC...
ATACAGTCT...

# Say we are given *two* sequences

Double roll
of a single die
per turn

vs.

Two independent
die rolling games



$$2 \times \quad 1\text{-}p$$

$$p$$

$$2 \times \quad 1\text{-}p$$

$$1\text{-}p \qquad 1\text{-}p$$

$$p \quad + \quad p$$

$$1\text{-}p \qquad 1\text{-}p$$

# Double roll ↔ Pair HMM

# Two sequences from a single *(unknown)* path?

The forward algorithm for HMMs gives probability of generating observed sequence $x$ from a given HMM:

$$\mathbb{P}(x) = \sum_{\text{paths } \sigma} \mathbb{P}(x; \sigma),$$

$$\mathbb{P}(x, y) = \sum_{\text{paths } \sigma} \mathbb{P}(x, y; \sigma),$$

probability of generating two sequences $x$ and $y$ from the same path through the HMM (i.e. from the same rearrangement event).

This is obtained by *efficiently* summing across paths.

V genes    D genes    J genes

...GGCCATGTATCTCTGTGCCAGCAGCTTAGGCACCGGTCATCGCTGGGGGCGCAGTAGGCCGGGGAGCTGTTTTTTGG...
...GGCCGTGTTTCTCTGTGCCAGCAGCACCGGTCATCGCTGGGGGCGCAGTAGGCCGGGGAGCTGTTTTTTGG...

# Do sets of sequences come from a single rearrangement event?

$$\frac{\mathbb{P}(A \cup B)}{\mathbb{P}(A)\mathbb{P}(B)} = \frac{\mathbb{P}(A \cup B \mid \text{single rearrangement})}{\mathbb{P}(A, B \mid \text{independent rearrangements})}$$

Use this for agglomerative clustering:

# Goal: maximum likelihood clustering

Find the maximum of

$$L(\{C_i\}_{i=1,\ldots,k}) = \prod_i \mathbb{P}(C_i)$$

across clusterings $\{C_i\}_{i=1,\ldots,k}$ of our sequences.

# HMM-based clustering works under simulation

# Likelihood-based clustering of clonal families

Phylogenetic empirical Bayes method for inferring unmutated common ancestor (perhaps?):

- Clonalyst: Kepler (2014-2015)

Use forward algorithm with parameter-rich HMMs for efficient evaluation of marginal probability.

- partis: Ralph & M. (2016) *in prep.*

# 3. Reconstruct BCR phylogenetic trees

reality

inference

ACATGGCTC...
ATACGTTCC...
TTACGGTTC...
ATCCGGTAC...
ATACAGTCT...

**c** Structural development of CAP256-VRC26 lineage



| Heavy chain longitudinal phylogenetic tree (condensed) | Mutations from UCA (aa) | Structure | CDR H3 | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | Disulphide | Tyrosine sulphation | | Electrostatics | | |
| | | | | Residue | Score | Charge | Surface representation | |
| UCA | H 0 / L 0 | | N100a / G100q | Y100h / Y100i | 1.8 / 1.6 | –10 | | |
| VRC26.01 | H 30 / L 9 | | S100a / I100q | Y100h / Y100i | 1.4 / 2.1 | –8 | | |
| VRC26.10 | H 33 / L 11 | | C100a / C100q | Y100h / Y100i | 1.4 / 1.7 | –6 | | |
| VRC26.11 / VRC26.06 | H 34 / L 14 | | C100a / C100r | Y100g / Y100i / Y100j | *1.1 / 1.4 / 0.8 | –7 | | |
| VRC26.05 / VRC26.02 | H 30 / L 16 | | C100a / C100q | Y100h / Y100i | 1.1 / 1.9 | –5 | | |
| VRC26.04 / VRC26.03 / VRC26.07 | H 29 / L 16 | | C100a / C100q | Y100h / Y100i | 1.1 / 1.9 | –6 | | |
| VRC26.12 | H 33 / L 17 | | C100a / C100q | Y100h / Y100i | 1.1 / 1.8 | –4 | | |

No cysteine ↑ / Disulphide ↓

Week
38
48
59
119
176
206

0.02
Evolutionary

VRC26.09

# Likelihood-based phylogenetics

Mutations appear at some rate $\lambda$:

ancestor •————$z$————$z$————• descendant

Mutations change bases according to substitution matrix:

$$\begin{pmatrix} p_{AA} & p_{AG} & p_{AC} & p_{AT} \\ p_{GA} & p_{GG} & p_{GC} & p_{GT} \\ p_{CA} & p_{CG} & p_{CC} & p_{CT} \\ p_{TA} & p_{TG} & p_{TC} & p_{TT} \end{pmatrix}$$

# Traditional phylogenetic approaches assume that the same evolutionary process is happening at each site.



**IGHV3-23D*01**

This does not hold for B cell receptor sequences.

# Context sensitive substitutions



Elhanati et al, 2015

# Context sensitive likelihoods are hard



$P(\text{GACGTG})$

GAAGTG    GACCTG    GACGTG

# Context sensitive likelihoods are hard

- Siepel & Haussler (MBE 2003); Saunders & Green (MBE 2007): context-sensitive likelihoods via along-sequence Markov cond'n
- Lunter & Hein (Bioinformatics 2004): MCMC approach to estimating likelihoods
- Christensen, Hobolth & Jensen (J Comp Biol 2005): pseudo-likelihood analysis using parsimony-ish inference on flanking bases
- Baele, Van de Peer, & Vansteelandt, (Sys Bio 2008): pseudo-likelihood analysis using context-insensitive likelihood inference on flanking bases
- Bérard & Guéguen (Sys Bio 2012): specific context dependent model enabling independence assumption in many cases
- Peter Ralph (unpublished): approximations using interacting particle systems

# Special sauce: per site models

Each site $s$ of every germline gene gets its own substitution rate $\lambda_s$ and mutation rate matrix:
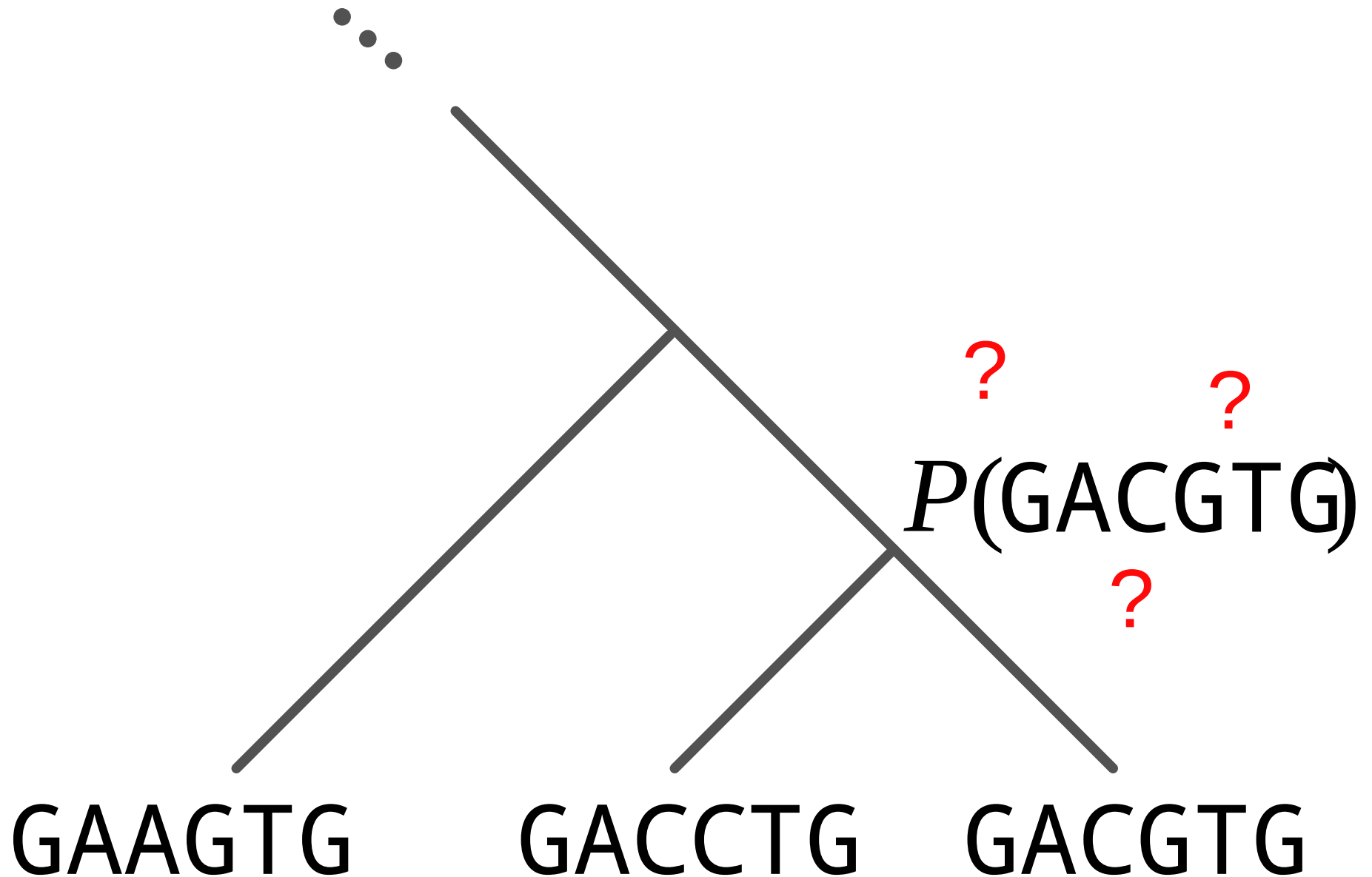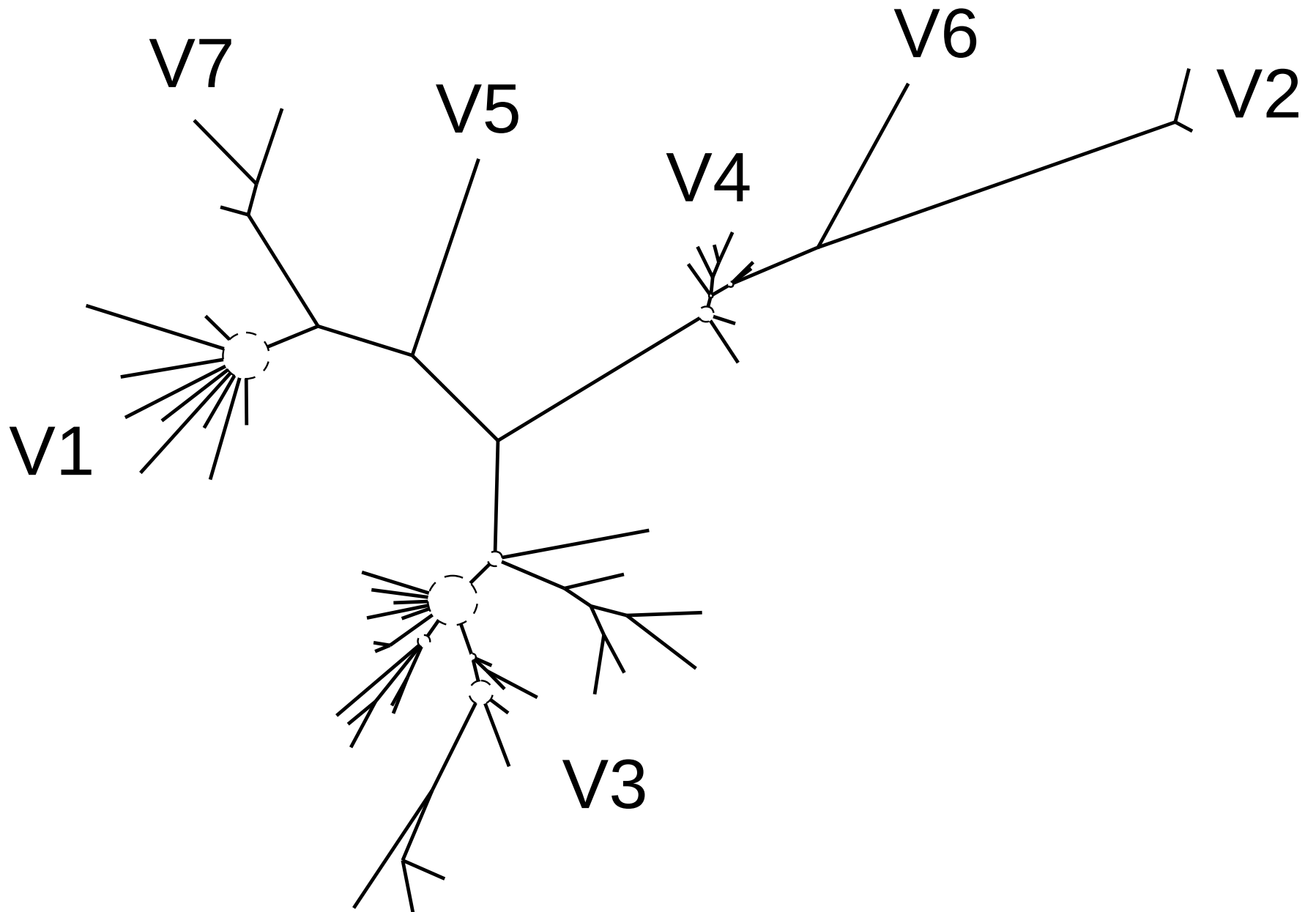
$$
\begin{pmatrix}
p_{\cdot A} & p_{\cdot G} & p_{\cdot C} & p_{\cdot T} \\
p_{\cdot A} & p_{\cdot G} & p_{\cdot C} & p_{\cdot T} \\
p_{\cdot A} & p_{\cdot G} & p_{\cdot C} & p_{\cdot T} \\
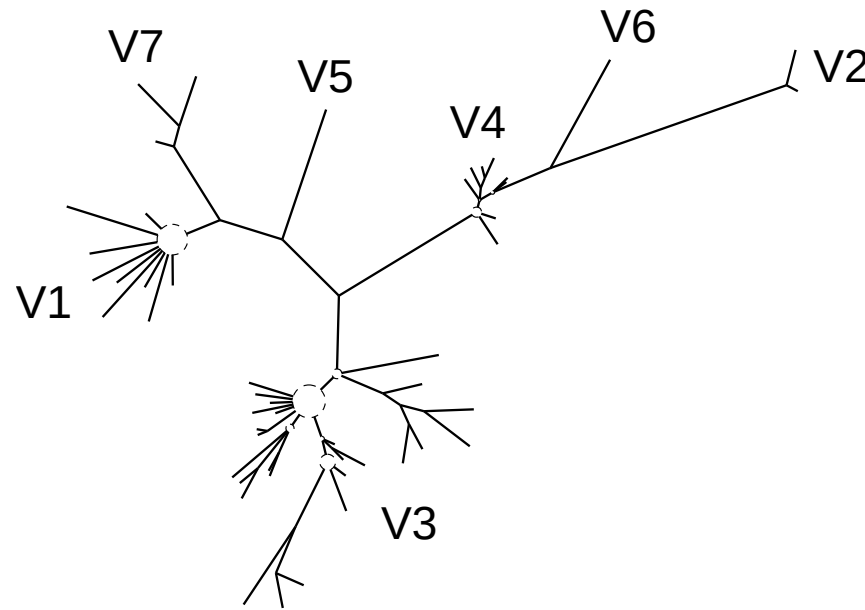p_{\cdot A} & p_{\cdot G} & p_{\cdot C} & p_{\cdot T}
\end{pmatrix}
$$

$$\approx 350 \times 5 \times 300 = 525,000 \text{ parameters}$$

Ouch! Need to be careful.
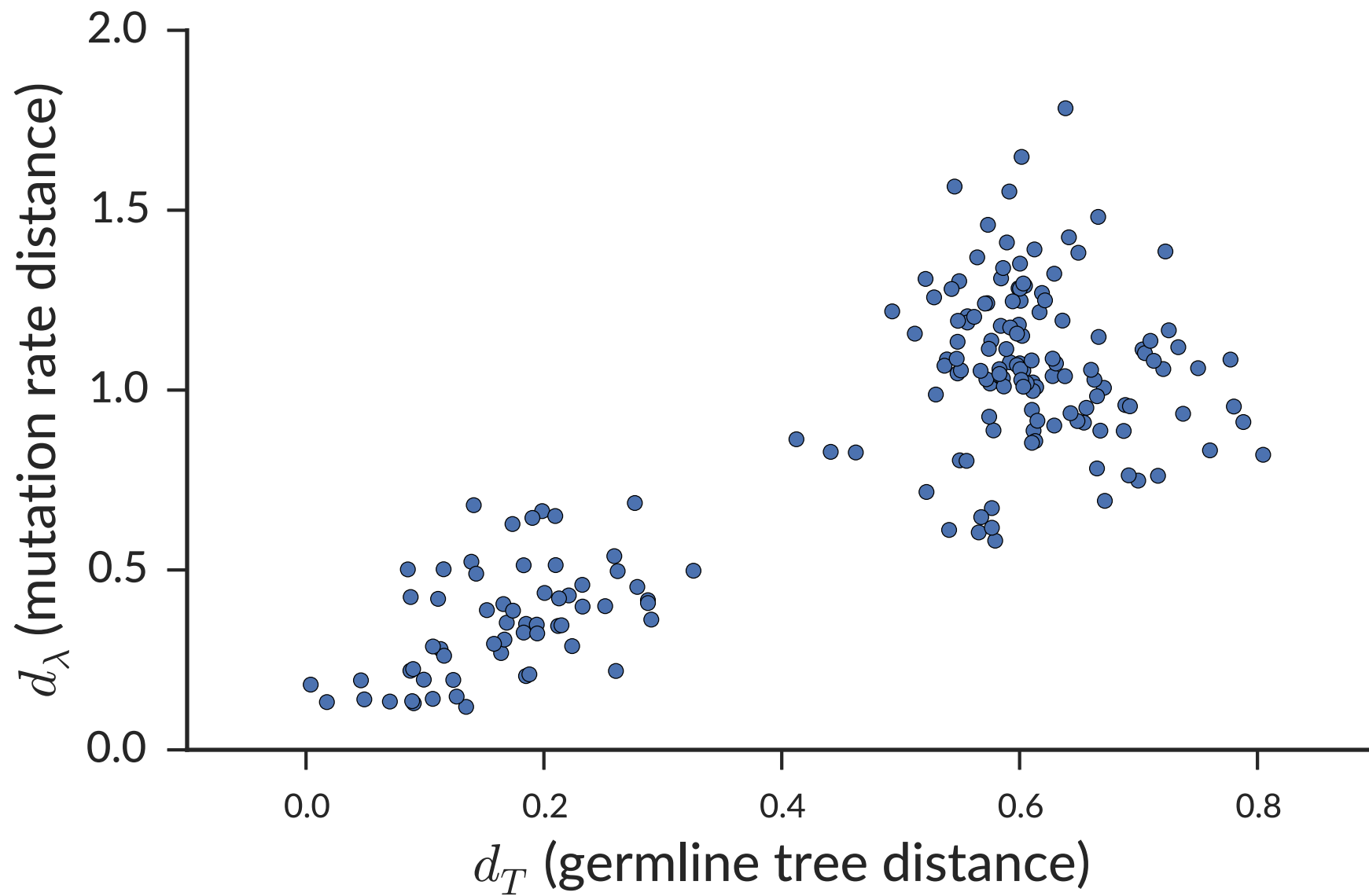
B cell germline gene phylogeny

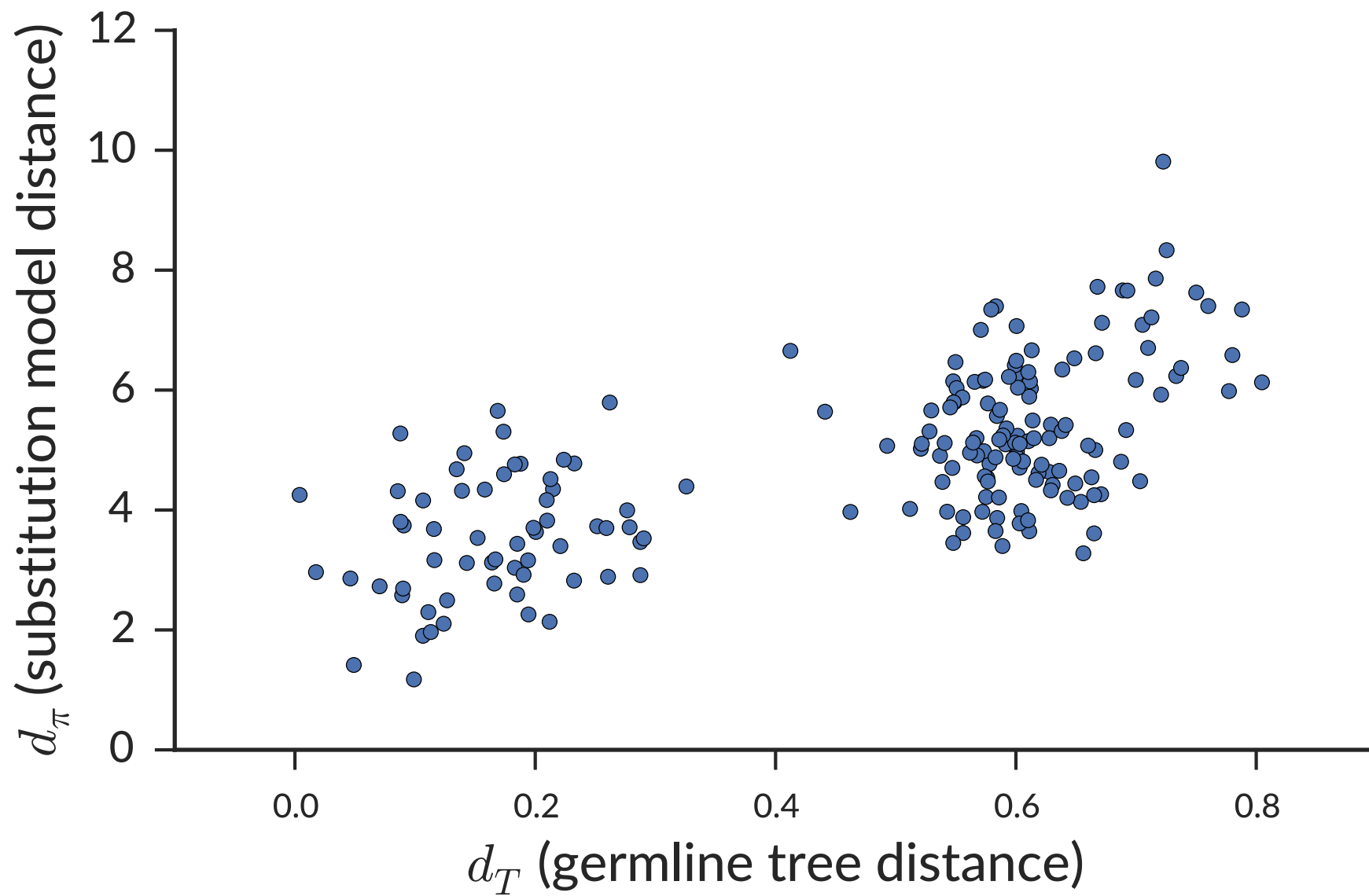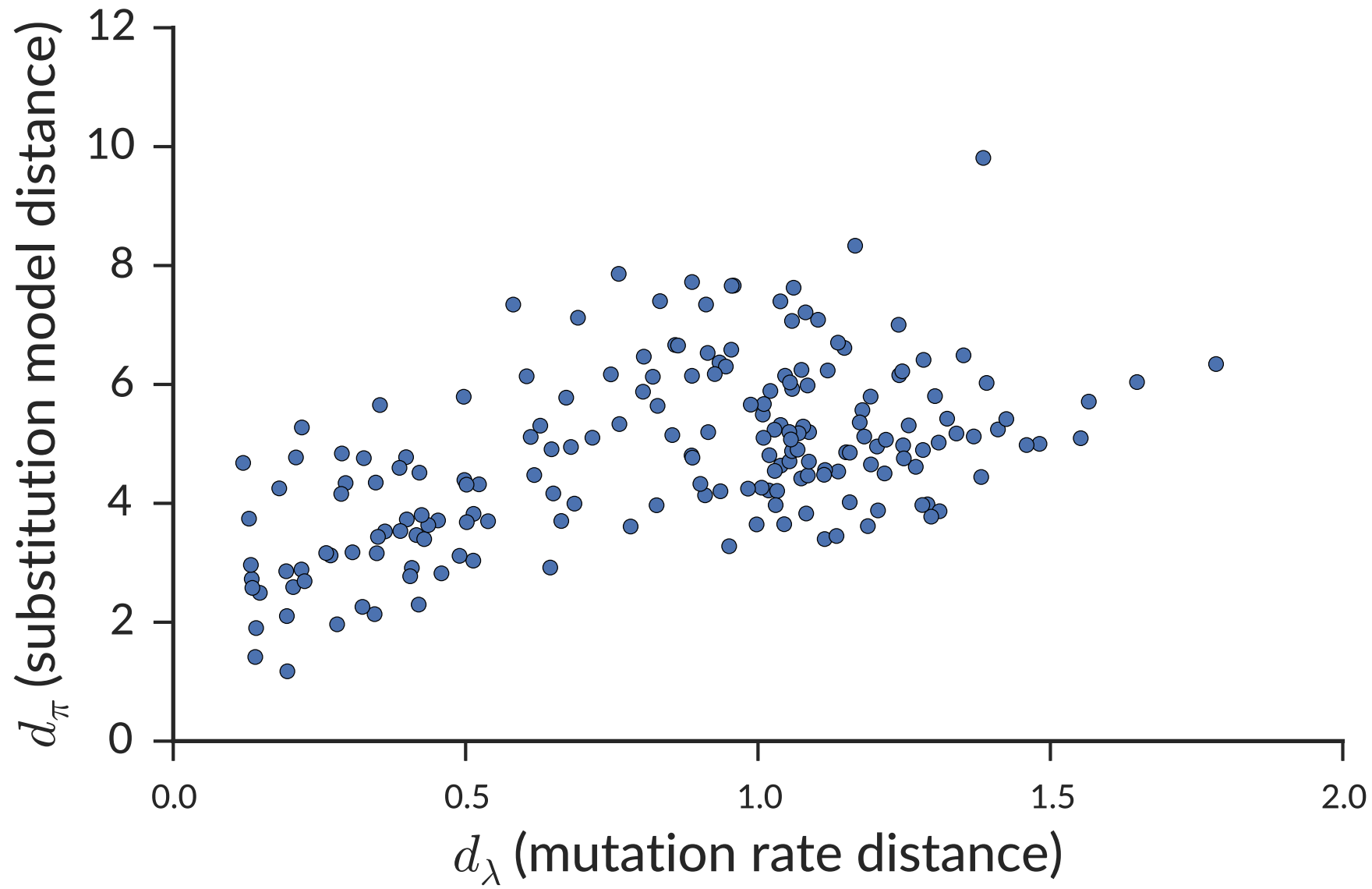# Q: do closely related genes evolve similarly?



- Fit models for the 20 genes for which we have the most data (these are good estimates)
- Compare parameter fits between these genes
- Compute evolutionary distance between these genes

# Top 20 genes

```
IGHV1-18*04      34507
IGHV1-2*04       19432
IGHV1-46*02      18453
IGHV1-69D*01     34218
IGHV3-15*07      18789
IGHV3-23D*01     58627
IGHV3-53*02      16552
IGHV3-64*04      38324
IGHV3-69-1*02    22445
IGHV3-7*01       78868
IGHV3-7*02       17992
IGHV3-74*03      18015
IGHV3-9*02       24010
IGHV3-NL1*01     51790
IGHV4-30-4*06    17419
IGHV4-34*13      14089
IGHV4-4*07       20816
IGHV4-61*02      18944
IGHV4-61*08      18644
```
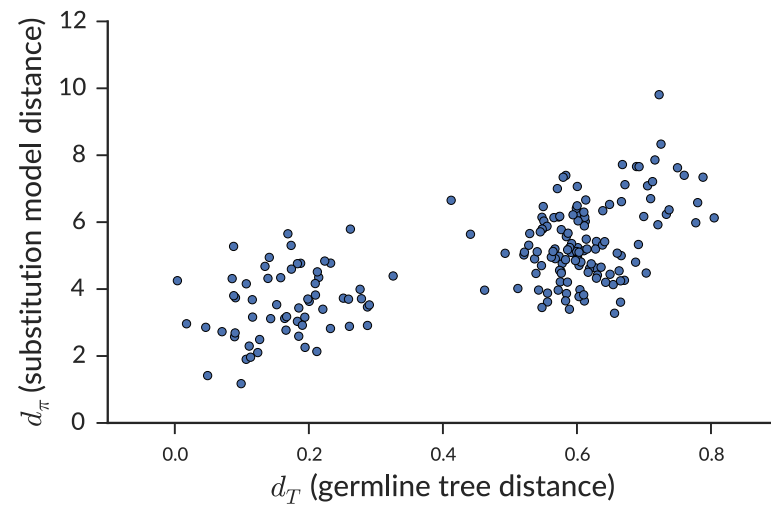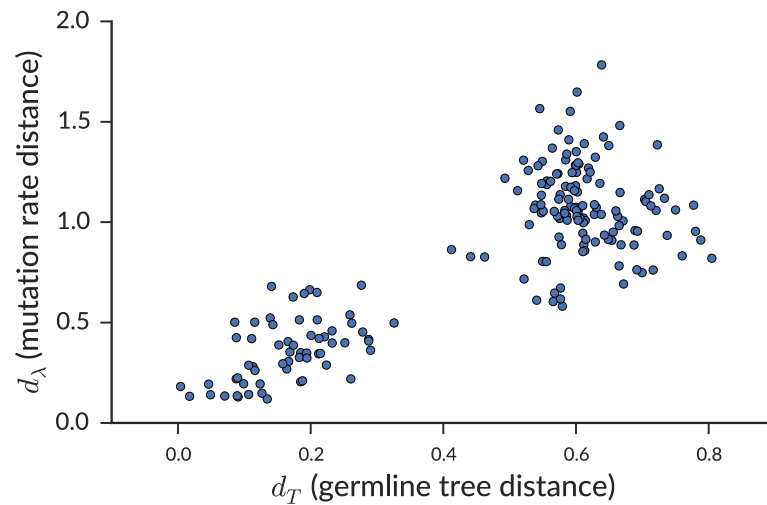
# We can use this for model fitting

Use homologous sites to regularize our rate parameters

# Multiple sequence alignment of V3 germline



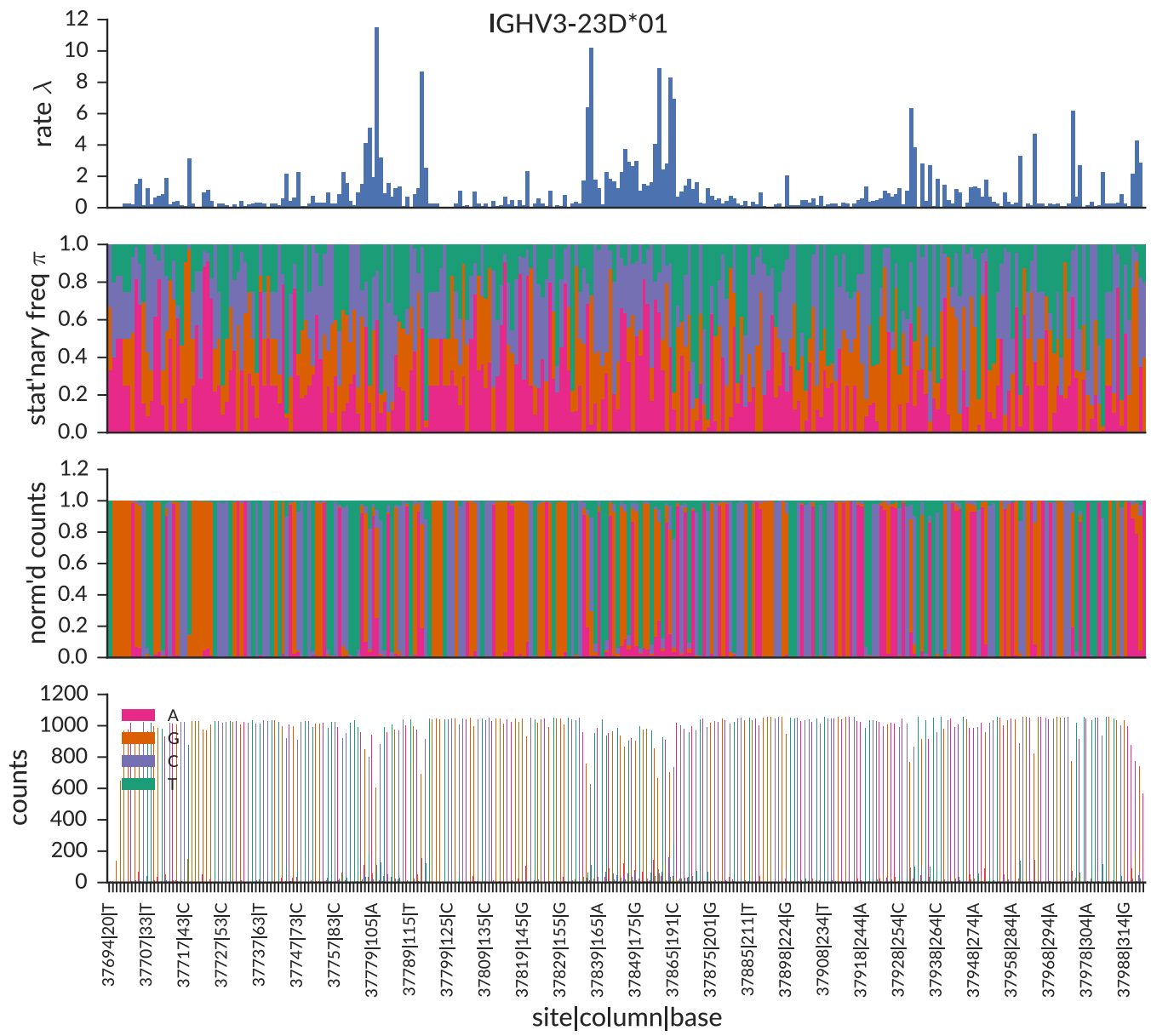(V3 genes along rows, colors for the four bases)

# Joint estimation for sites in the same column

Assume that sites in the same column evolve *similarly* within-host.

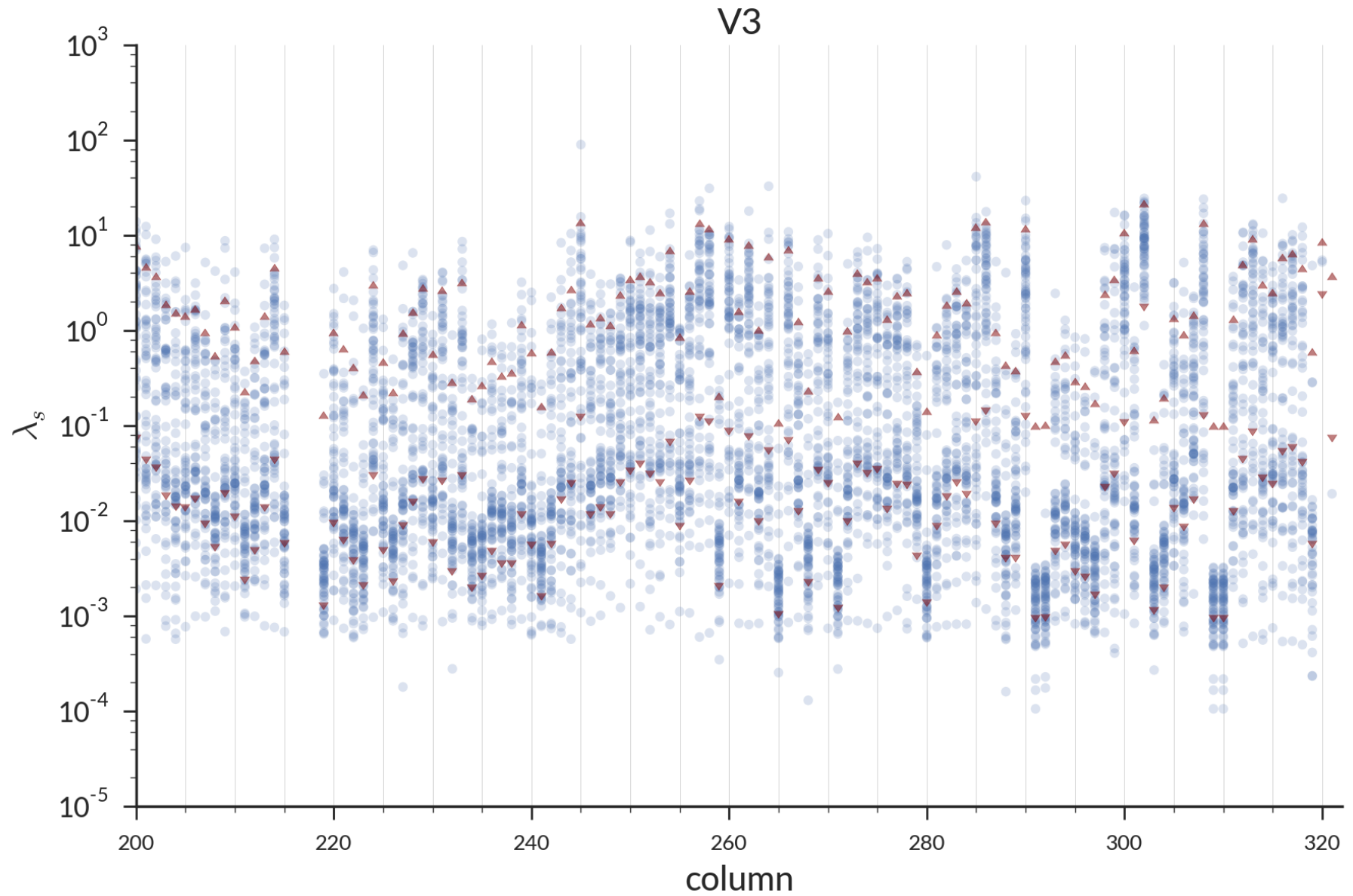When data is weak, draw estimates back to a per-column average.

- Substitution rate $\lambda_s \sim \mathrm{Gamma}(\omega_c, \theta_c)$
- Gamma mode $\omega_c \sim \mathrm{Log\text{-}normal}(1, 1)$
- Gamma dispersion $\theta_c \sim \mathrm{L\acute{e}vy}(3)$
- Stationary distribution $\pi_s \sim \mathrm{Dirichlet}(3, 3, 3, 3)$
- Branch length $t \sim \mathrm{Exponential}(0.1)$
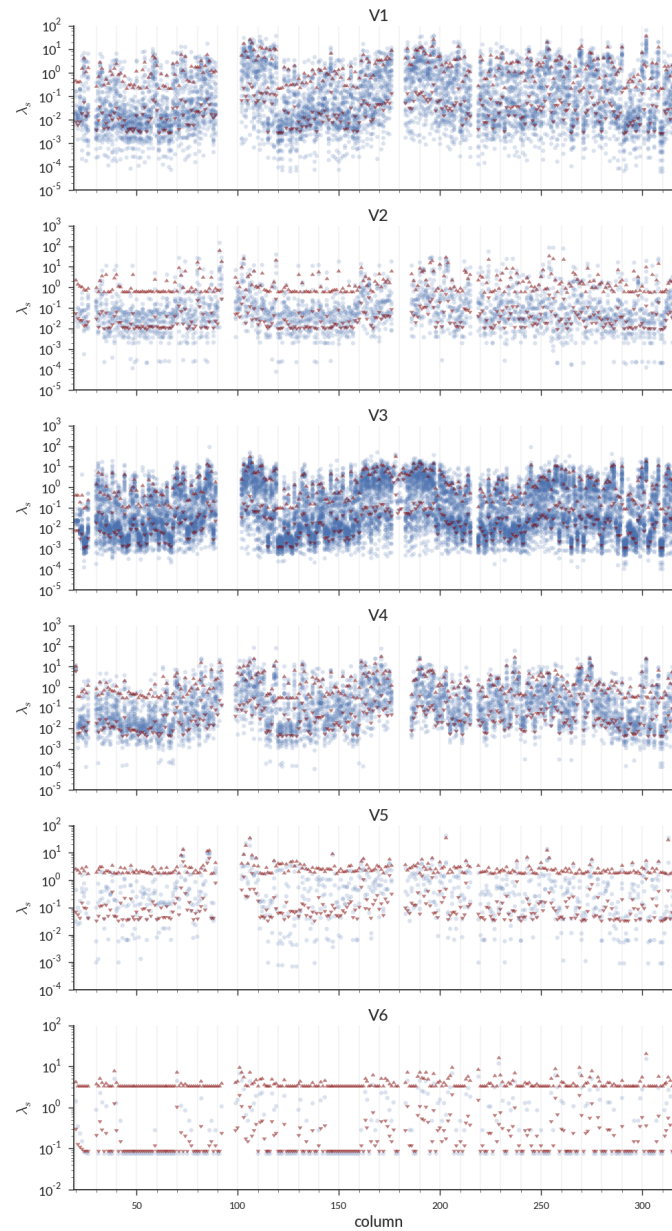
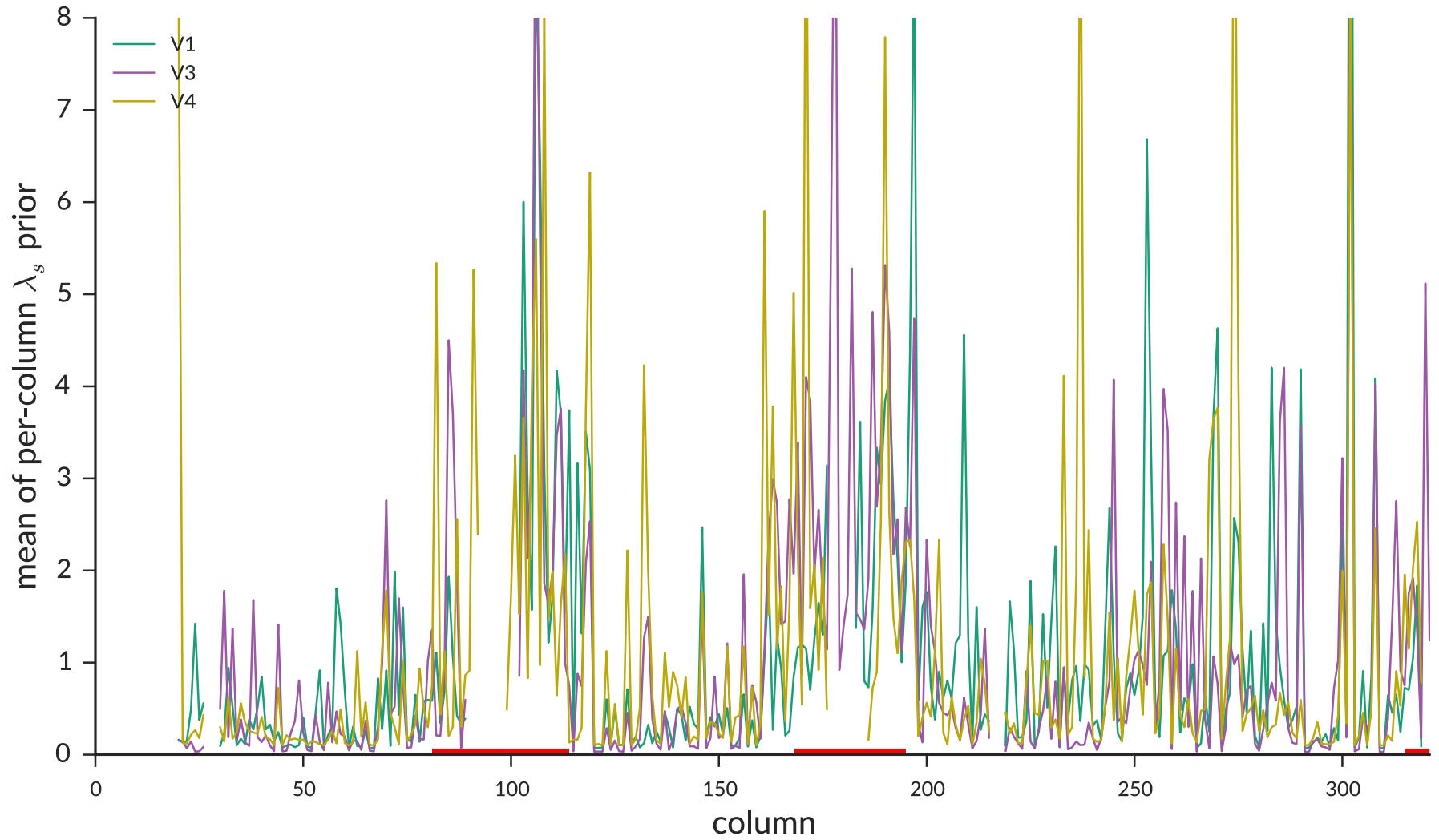IGHV3-23D*01

IGHV3-23D*01

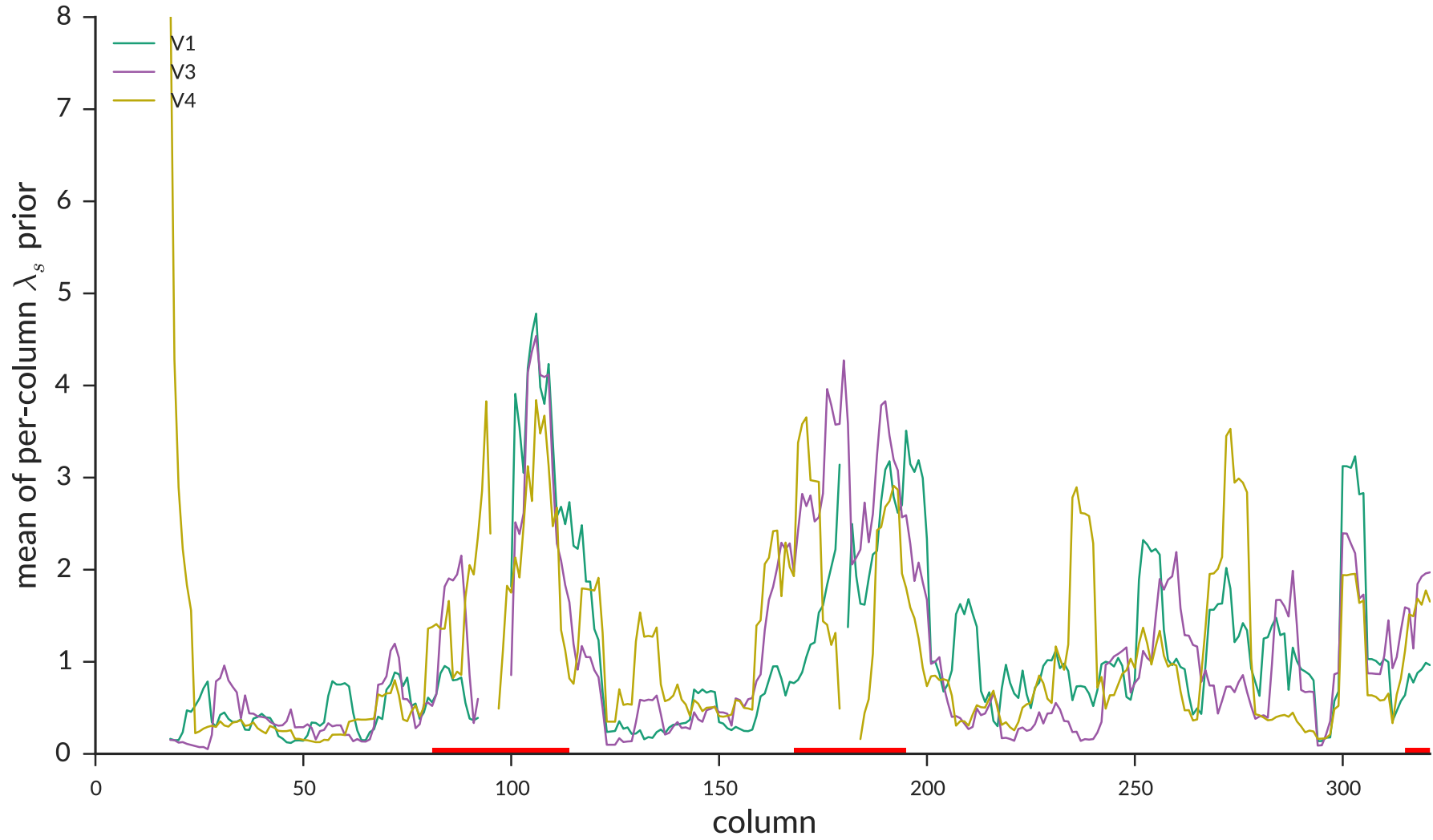# Points = rates; triangles = 95% CI of fit prior

# Points = rates; triangles = 95% CI of fit prior

Per-site mutation rate

Per-site mutation rate (smoothed)

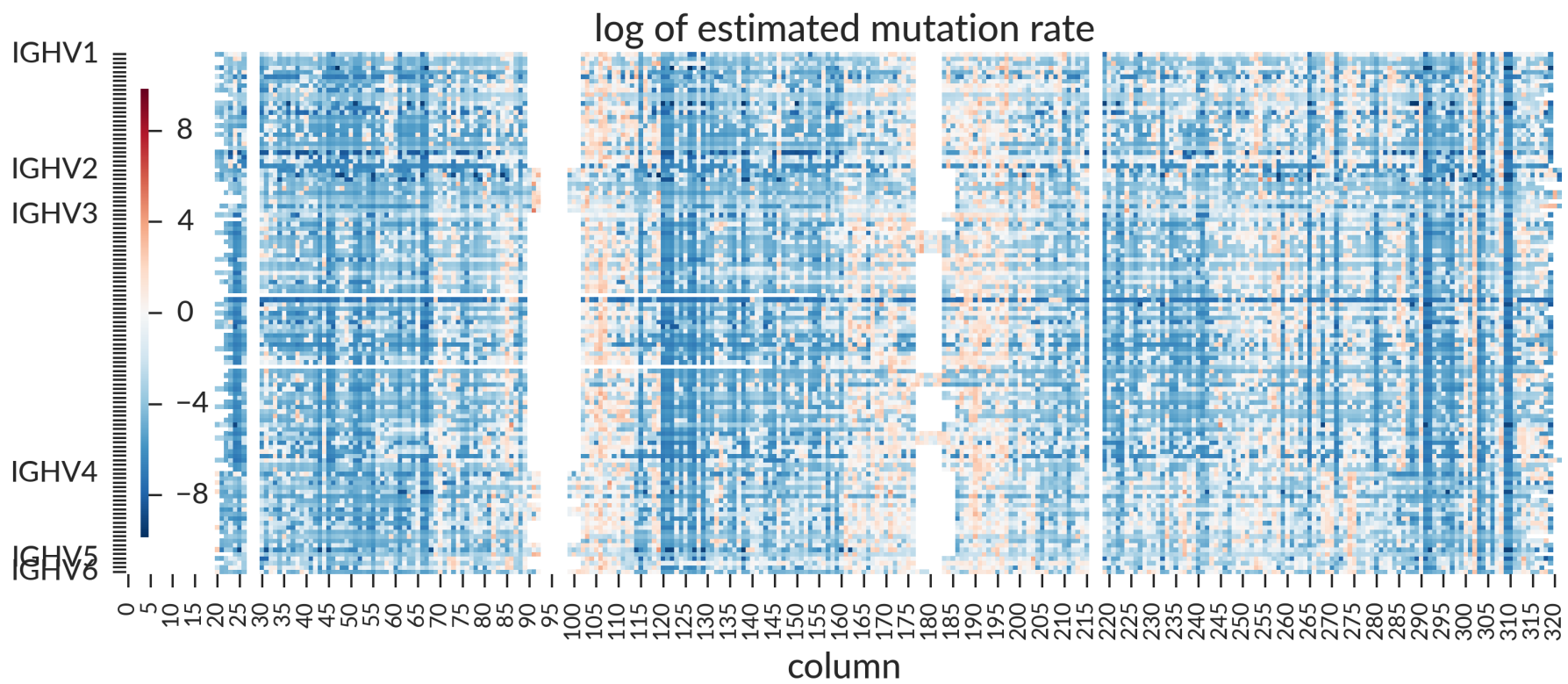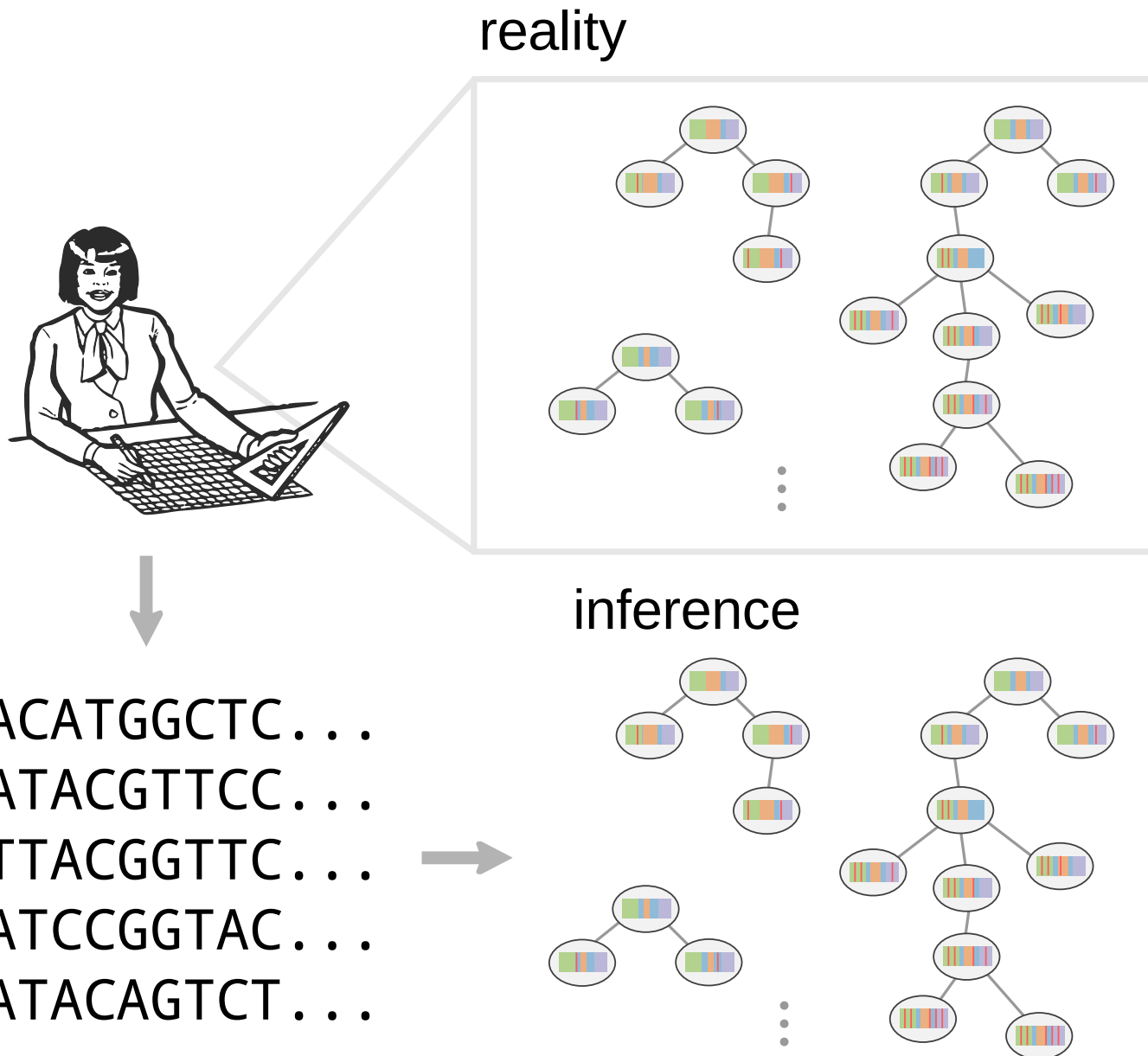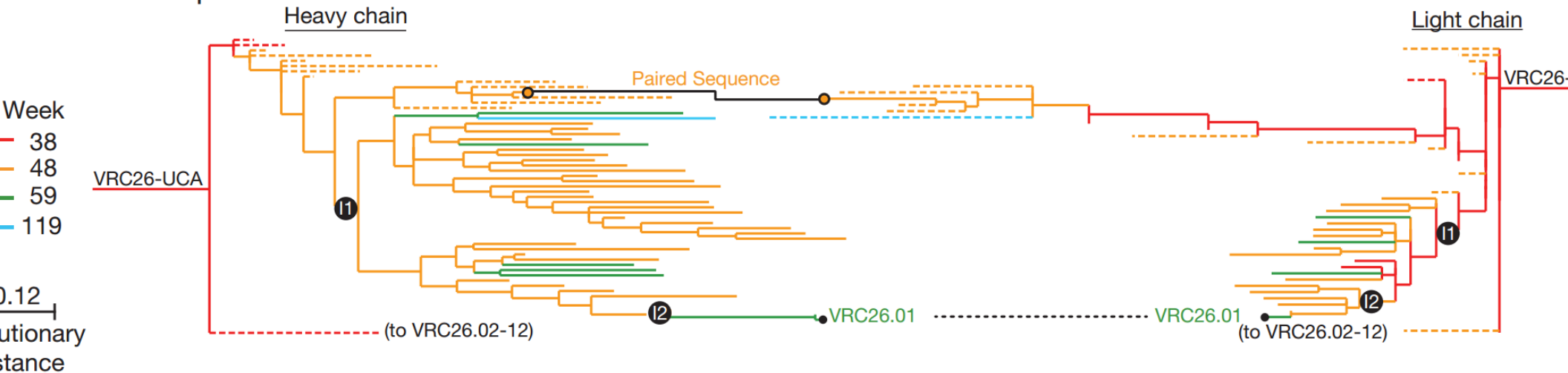log of estimated mutation rate

# Likelihood-based phylogenetics for B cell receptor sequences

- Many people use likelihood-based phylogenetics in their analysis, but with models that are identical across sites
- Substitution is manifestly *not* identical across sites
- One could work to do phylogenetics with context-sensitive models (hard!) or infer per-site parameters (need regularization!)
- Need to build software that can build trees with these models
- Sampled ancestors also a challenge, but this can be handled in a Bayesian or penalized likelihood framework (in progress).
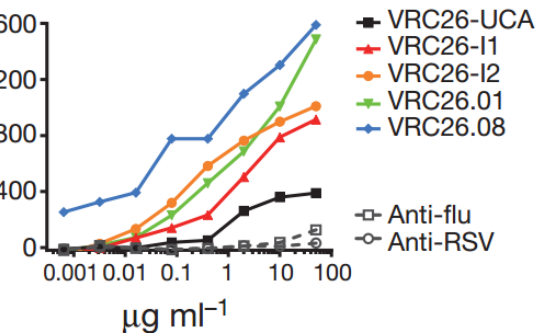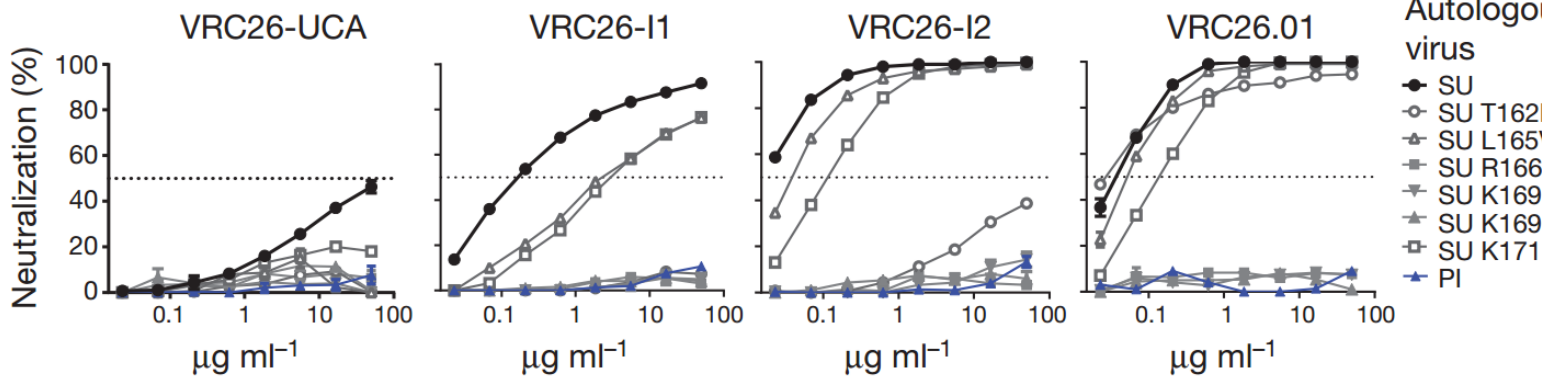
# 4. Find BCR ancestral sequences

reality

inference

ACATGGCTC...
ATACGTTCC...
TTACGGTTC...
ATCCGGTAC...
ATACAGTCT...

# Development of CAP256-VRC26.01



Heavy chain

Light chain

Paired Sequence

VRC26-UCA

Week
- 38
- 48
- 59
- 119

0.12
Evolutionary
distance

(to VRC26.02-12)    VRC26.01 ------- VRC26.01    (to VRC26.02-12)

**Binding to autologous Env (SU)**

- VRC26-UCA
- VRC26-I1
- VRC26-I2
- VRC26.01
- VRC26.08
- Anti-flu
- Anti-RSV

$\mu g\ ml^{-1}$

**c    Neutralization of autologous HIV-1**

VRC26-UCA    VRC26-I1    VRC26-I2    VRC26.01

Neutralization (%)

$\mu g\ ml^{-1}$

Autologous virus
- SU
- SU T162
- SU L165
- SU R166
- SU K169
- SU K169
- SU K171
- PI

**Binding to heterologous Env(ZM53)**

- VRC26-UCA
- VRC26-I1
- VRC26-I2
- VRC26.01
- VRC26.08
- Anti-flu
- Anti-RSV

**e    Neutralization of heterologous HIV-1**

VRC26-UCA    VRC26-I1    VRC26-I2    VRC26.01

Neutralization (%)

Heterologous virus
- 30163v5.
- CAP210
- CM244
- KER2008
- KER2018
- ZM197
- ZM53.12

# Likelihood-based ancestral sequence reconstruction

Currently being done with identical-across-sites models.

Once we have per-site models, it will be                    .

# 5. Selection inference for BCRs

**For selection**

Pro → Pro

CCA ⟶ CCT

*synonymous*

Thr → Ile

ACC ⟶ ATC

*nonsynonymous*

**For selection**

Pro    Pro
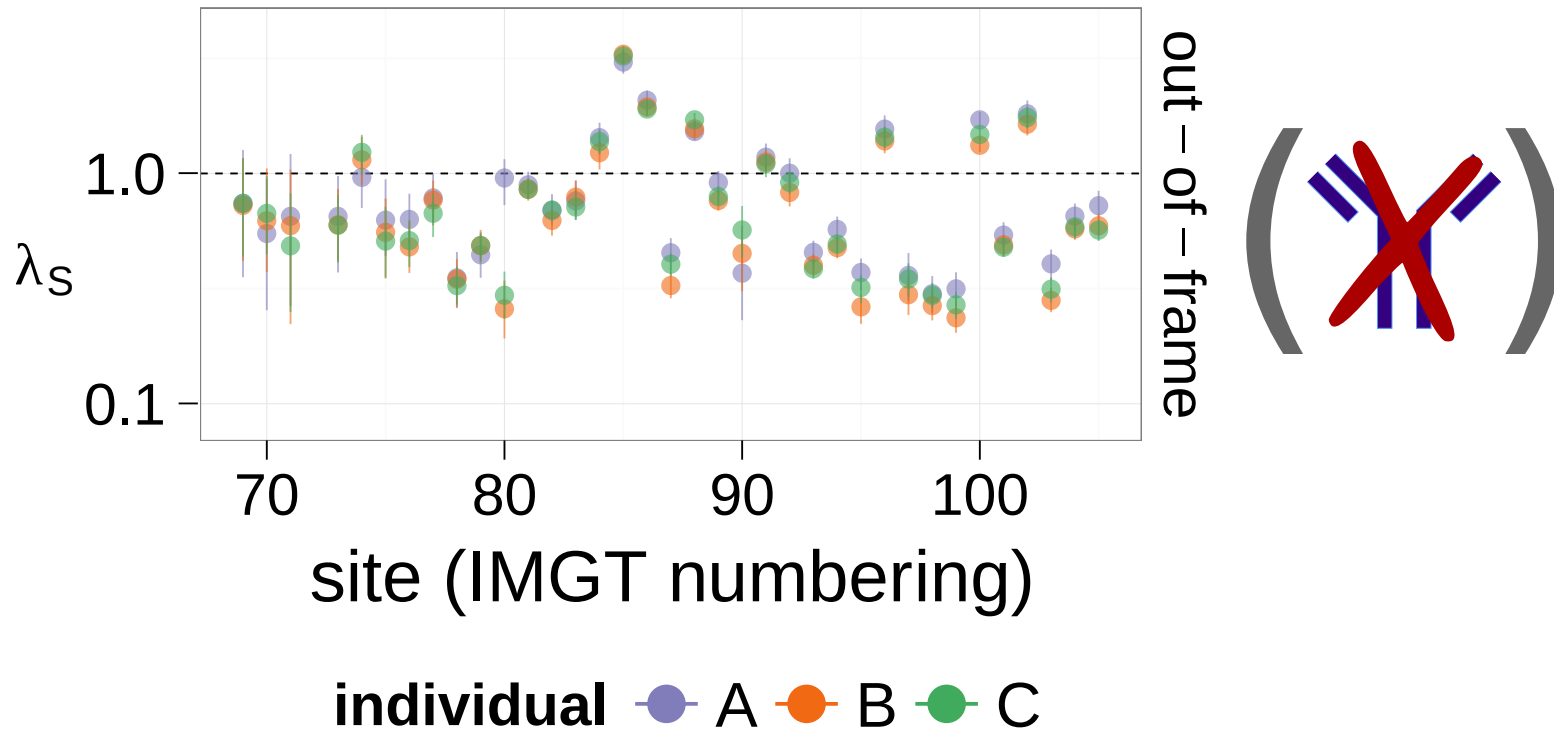
CCA → CCT

*synonymous*

Thr    Ile

ACC → ATC

*nonsynonymous*

**In antibodies**

AAC → AAG

*more likely*

GTC → GTG

*less likely*

$$\omega \equiv \frac{dN}{dS} \equiv \frac{\text{rate of non-synonymous substitution}}{\text{rate of synonymous substitution}}$$
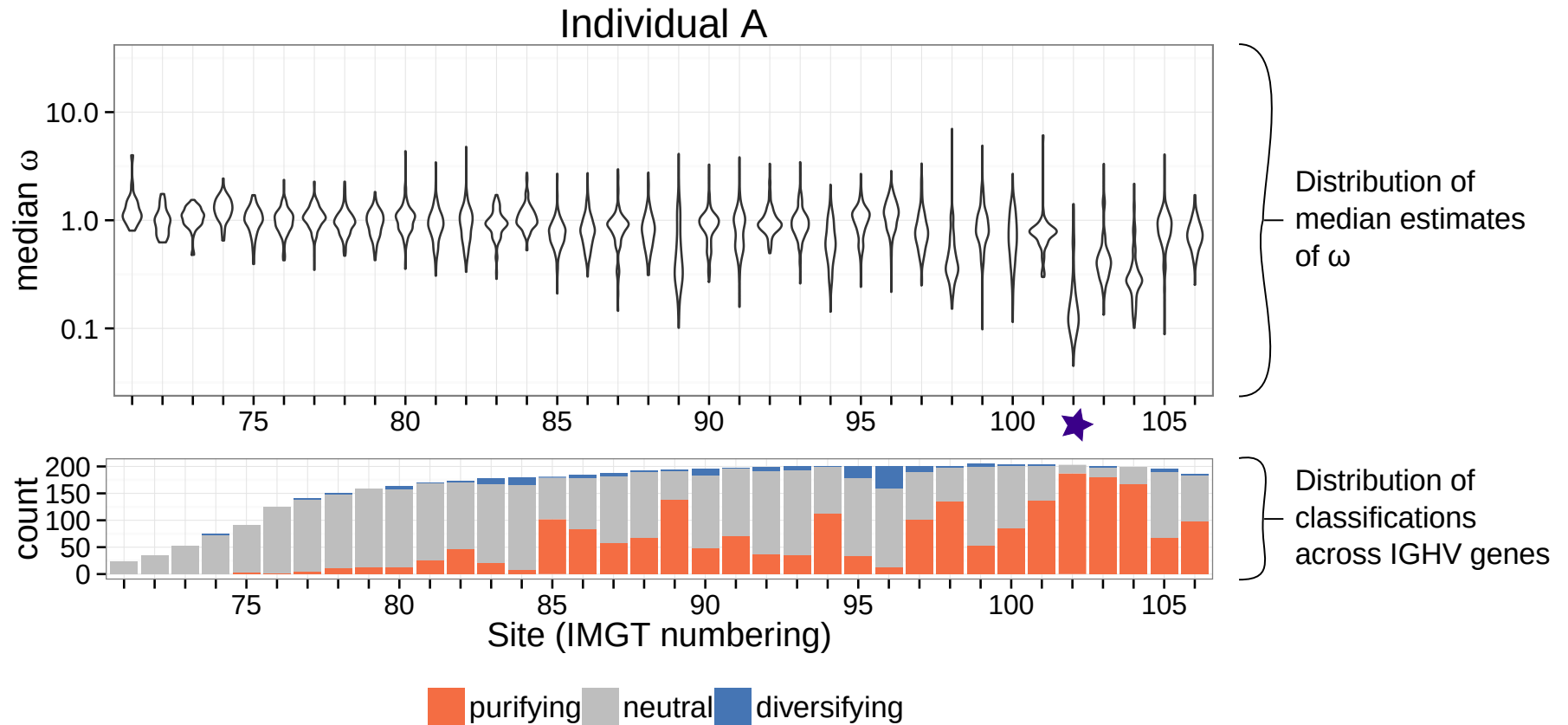


**Out-of-frame reads can be used to infer neutral mutation rate.**

# Estimating selection coefficient $\omega_l$

- $\lambda_l^{(N-I)}$ : nonsynonymous in-frame rate for site $l$
- $\lambda_l^{(N-O)}$ : nonsynonymous out-of-frame rate for site $l$
- $\lambda_l^{(S-I)}$ : synonymous in-frame rate for site $l$
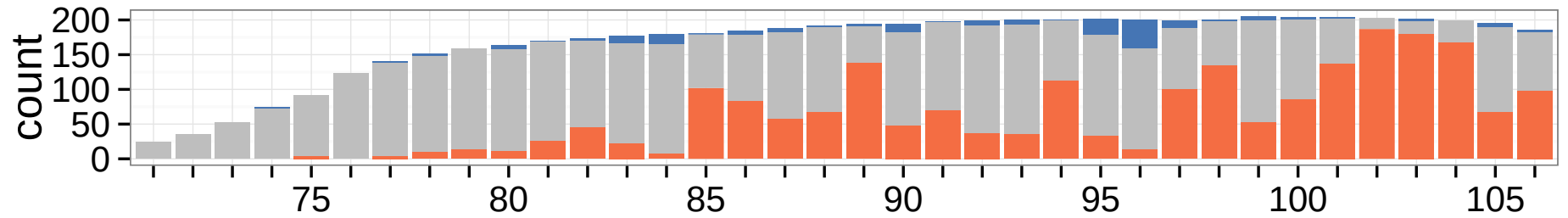- $\lambda_l^{(S-O)}$ : synonymous out-of-frame rate for site $l$

$$\omega_l = \frac{\lambda_l^{(N-I)} / \lambda_l^{(N-O)}}{\lambda_l^{(S-I)} / \lambda_l^{(S-O)}}$$
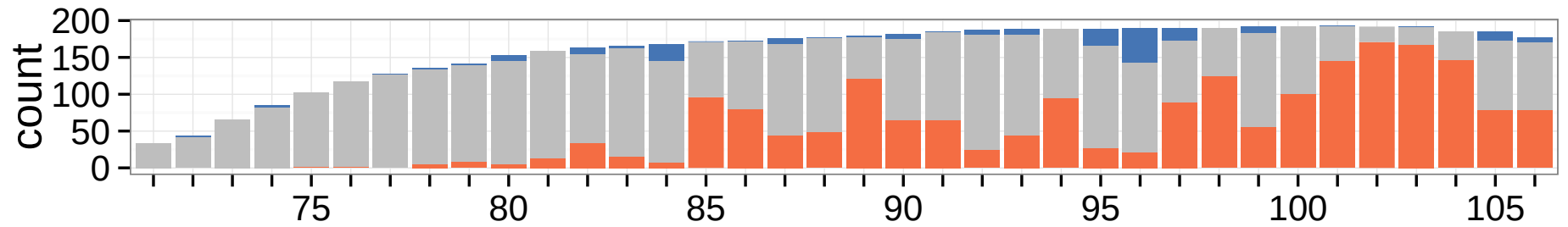
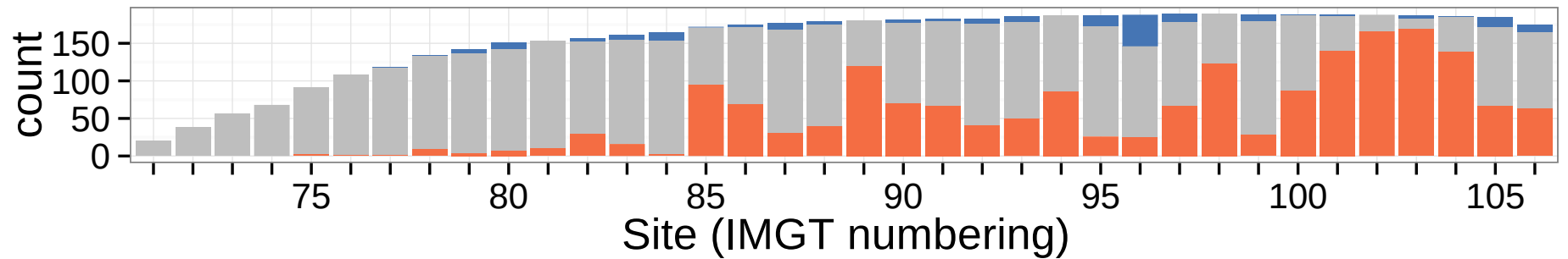Overall IGHV selection map

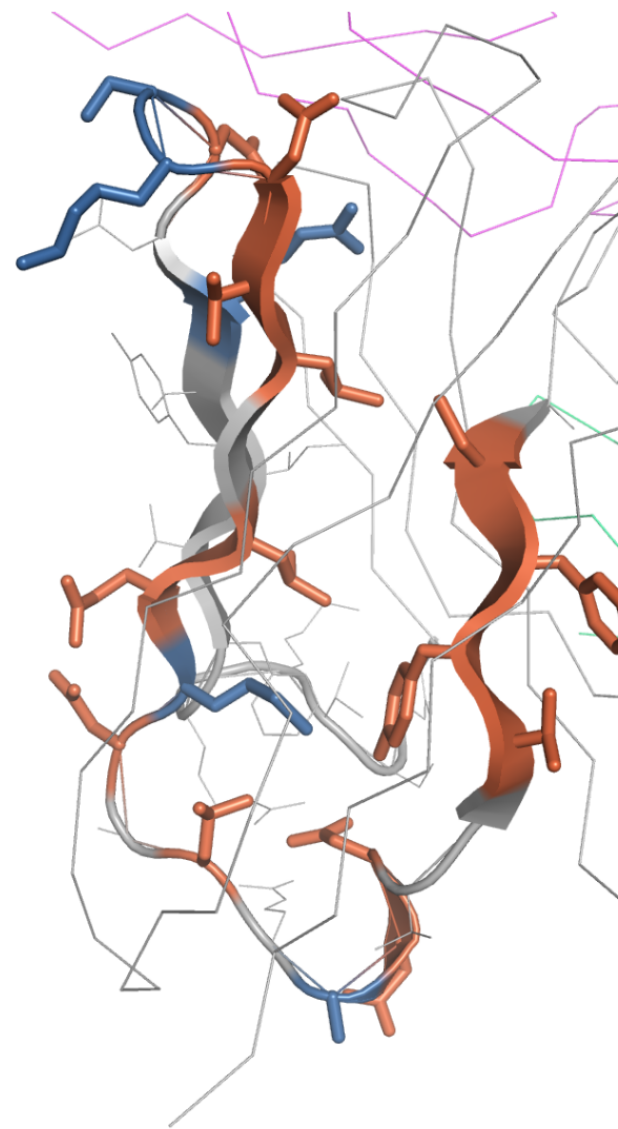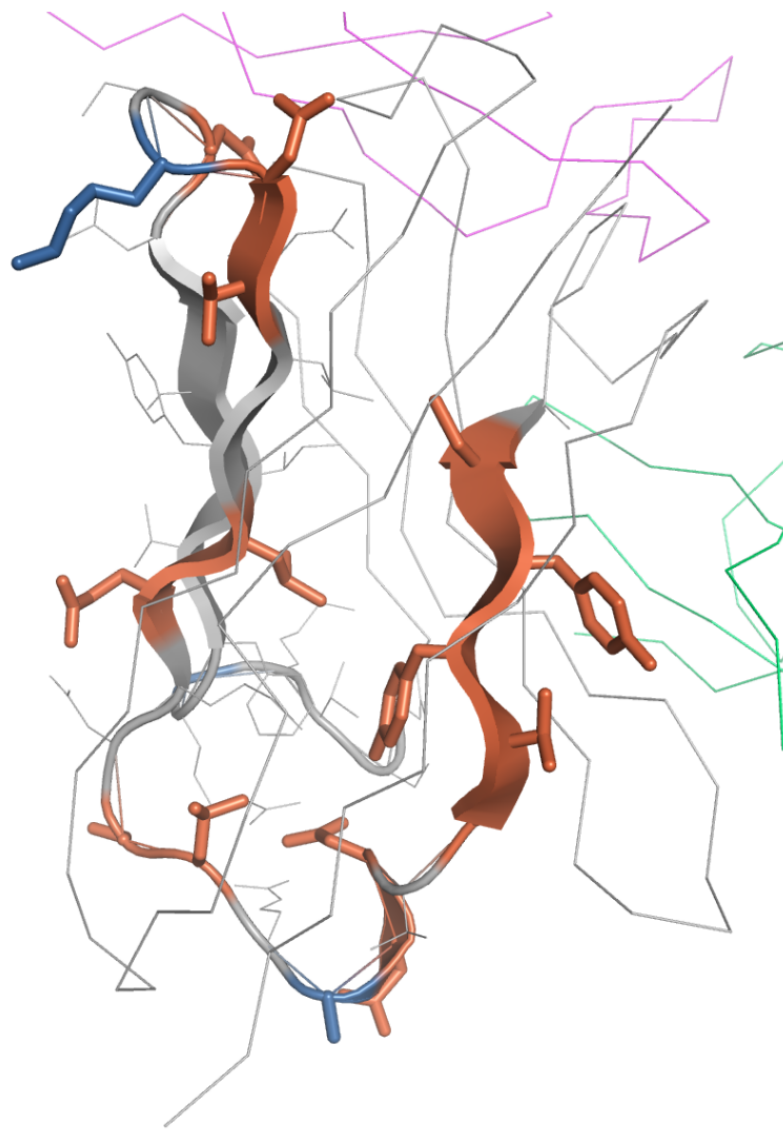# Similar across individuals

### Individual A

### Individual B

### Individual C

Site (IMGT numbering)
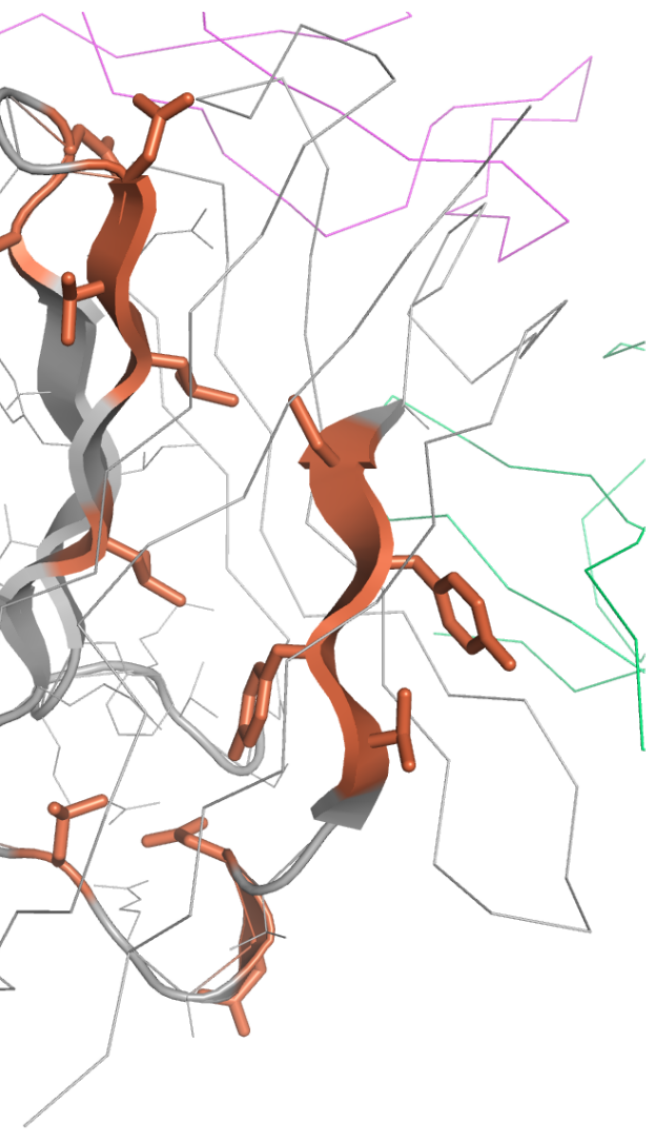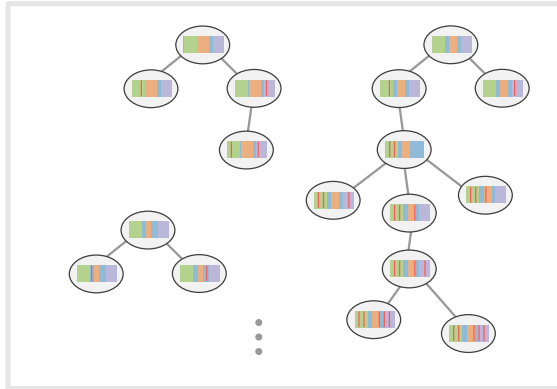
purifying   neutral   diversifying

# Likelihood-based selection inference for BCRs

- Aggregate selection analysis: Yaari, Uduman & Kleinstein (2012). *Nucleic Acids Research*
- Amino acid preferences: Elhanati, Sethna, Marcou, Callan, Mora & Walczak (2015). *Phil Trans Royal Soc B*
- Per-site analysis: McCoy, Bedford, Minin, Bradley, Robins & M. (2015). *Phil Trans Royal Soc B*

$$\sim f(\text{health}, \text{genetics}, \text{age}, \dots)$$

Make this relationship explicit by developing probabilistic models
with priors in terms of covariates.

We are approaching this from an abstract statistical perspective
rather than via mechanistic models.