

 michael.landis@wustl.edu

 landislab.github.io

 @landismj

Phylogenetic inference and graphical models

Molecular Evolution Workshop
Woods Hole, MA

Michael Landis
Wash U in St. Louis
Sunday, August 4

What's in a recipe?

Potato Soup

Cut up onion, celery & potatoes.
Cook in sm amt of water till done.
Fill kettle with milk. When milk
is hot drop in dumplings.

Dumplings

Beat eggs till light, stir in flour
till stiff. add salt. Drop into
soup & cook about 20 min.

ingredients

Potato Soup

Cut up onion, celery & potatoes.
Cook in sm amt of water till done.
Fill kettle with milk & when milk
is hot drop in dumplings.

Dumplings

Beat eggs till light, stir in flour
till stiff - add salt. Drop into
soup & cook about 20 min.

steps

ingredients

Potato Soup

Cut up onion, celery & potatoes.
Cook in sm amt of water till done.
Fill kettle with milk & when milk
is "hot" drop in dumplings.

Dumplings

Beat eggs till light, stir in flour
till stiff - add salt. Drop into
soup & cook about 20 min.

steps

ingredients

Potato Soup

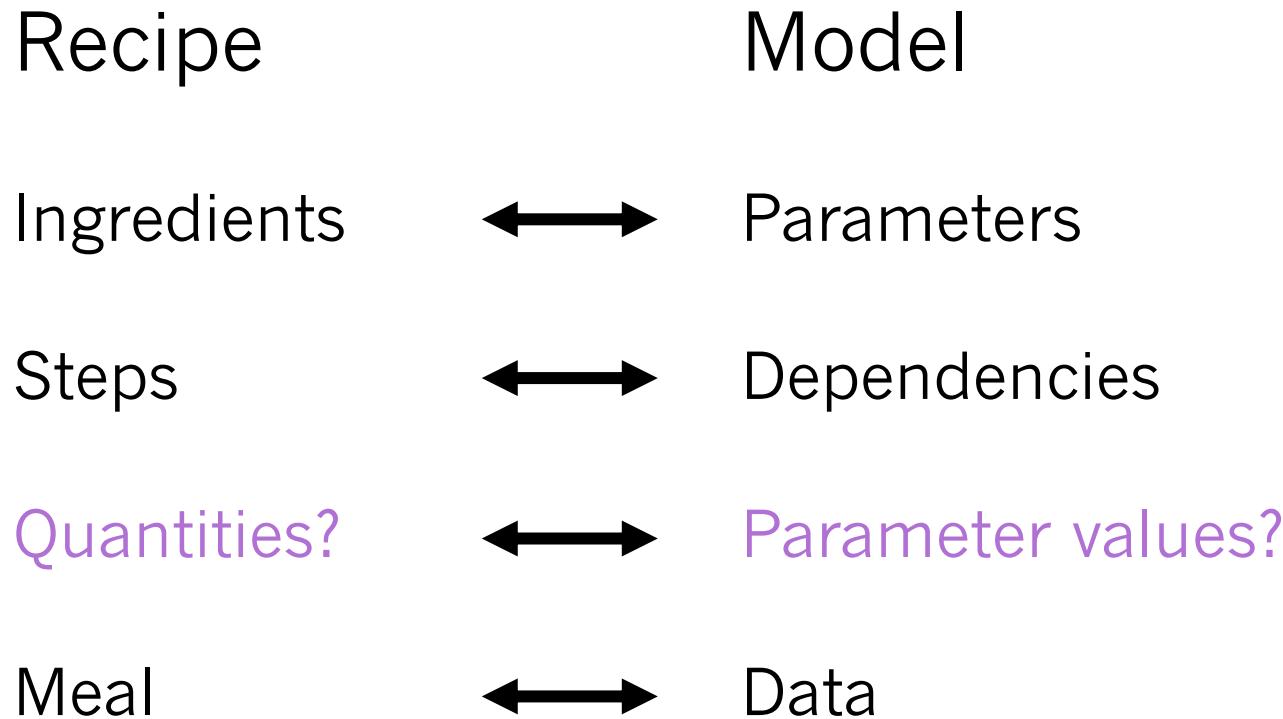
Cut up onion, celery & potatoes.
Cook in amount of water till done.
Fill kettle with milk & when milk
is "hot" drop in dumplings.

Dumplings

Beat eggs till light, stir in flour
till stiff - add salt. Drop into
soup & cook about 20 min.

quantities??

Estimating ingredient quantities is an *inference* problem



steps??

ingredients??

quantities??

Recovering the recipe is a *modeling* problem

Recipe

Model

Ingredients?



Parameters?

Steps?



Dependencies?

Quantities?



Parameter values?

Meal

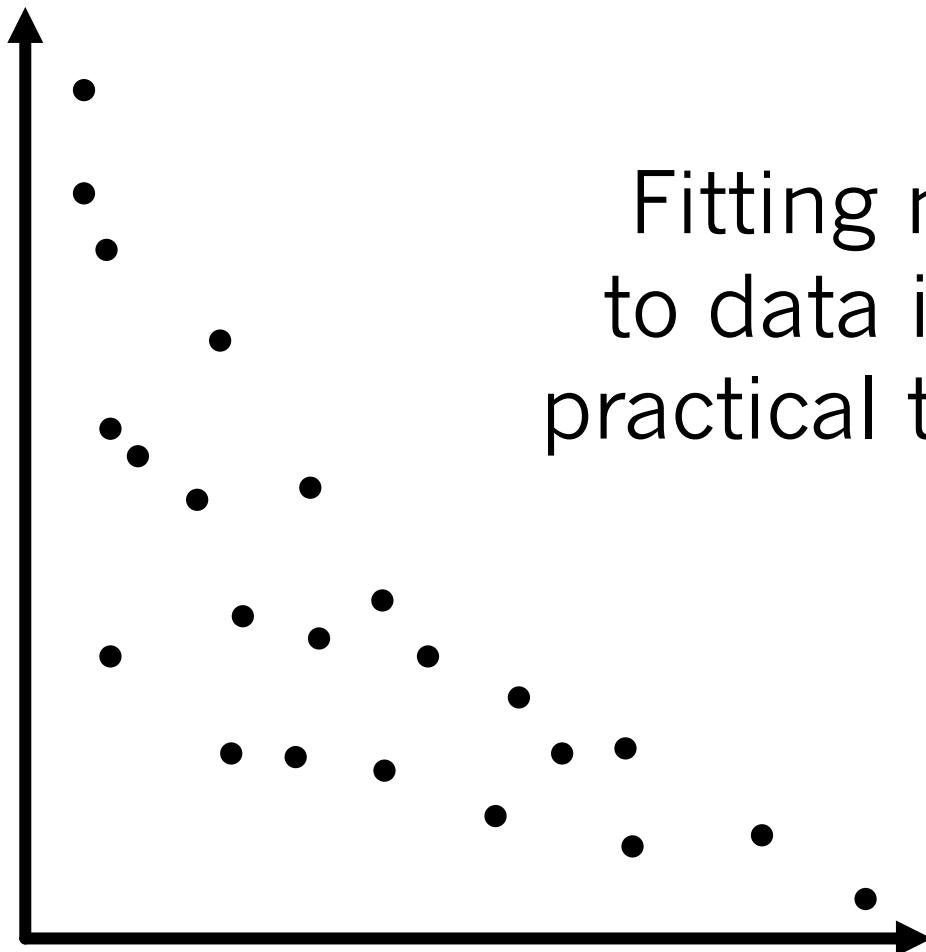


Data



Practicality

time, scale,
cost, ease



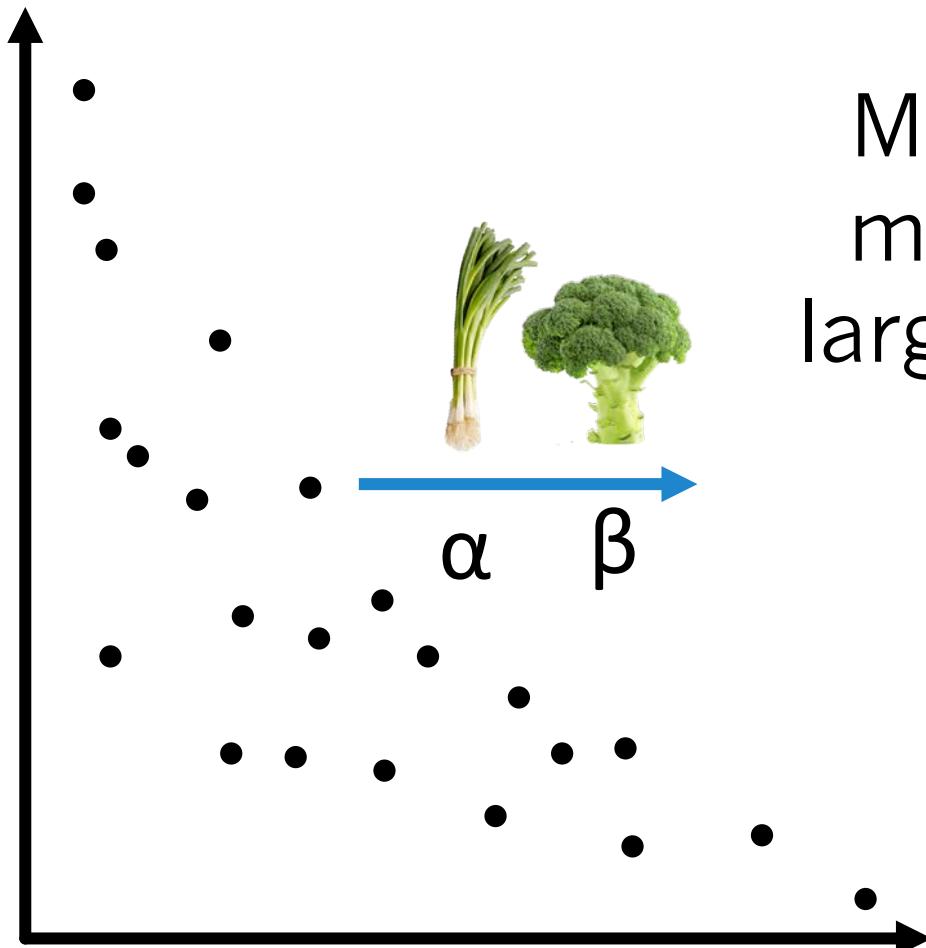
Fitting models
to data involves
practical trade-offs



Purity
realism,
correctness,
accuracy,
elegance



Practicality
time, scale,
cost, ease



Minor changes to
models can bring
large improvements



Purity
realism,
correctness,
accuracy,
elegance

Learned quite a bit already

- Basics of probability
- Frequentist/Bayesian inference
- Substitution processes
- Site-rate models
- Partition models
- Tree space
- Pruning algorithm
- Model selection

...

About to learn much more

Basics of probability
Frequentist/Bayesian inference
Substitution processes
Site-rate models
Partition models
Tree space
Pruning algorithm
Model selection
Birth-death processes
Selection models
Protein models
Coalescent processes
Phylogenetic network models
Species tree estimation
Species delimitation

...

Practical model literacy

memorize features
for each specific model

– or –

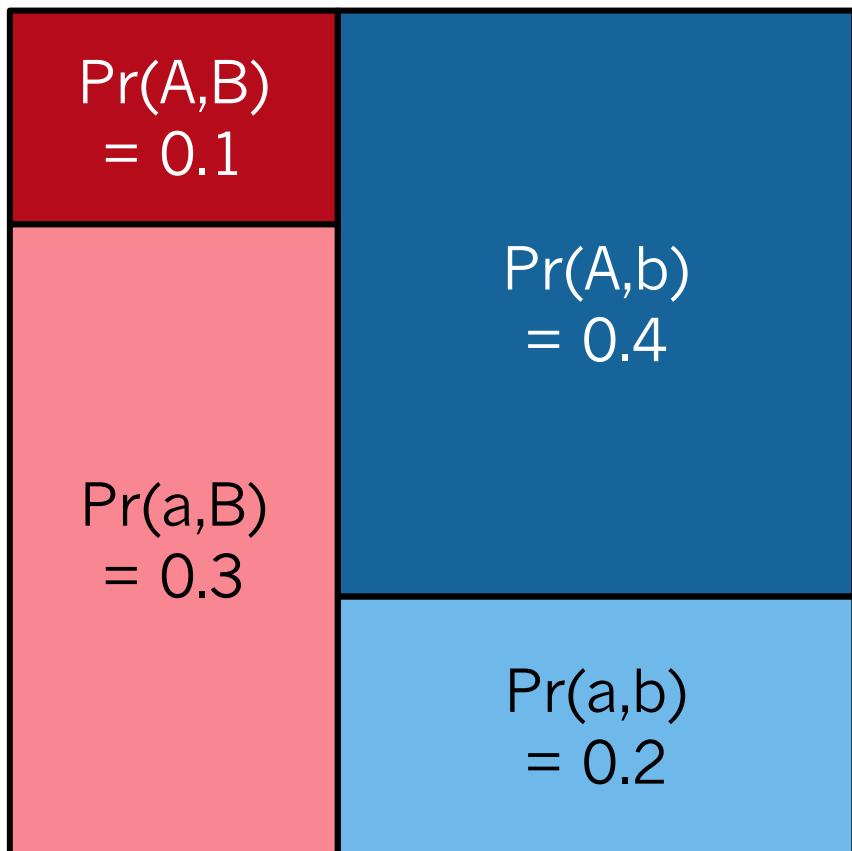
learn to identify, classify, modify
model variants

Outline

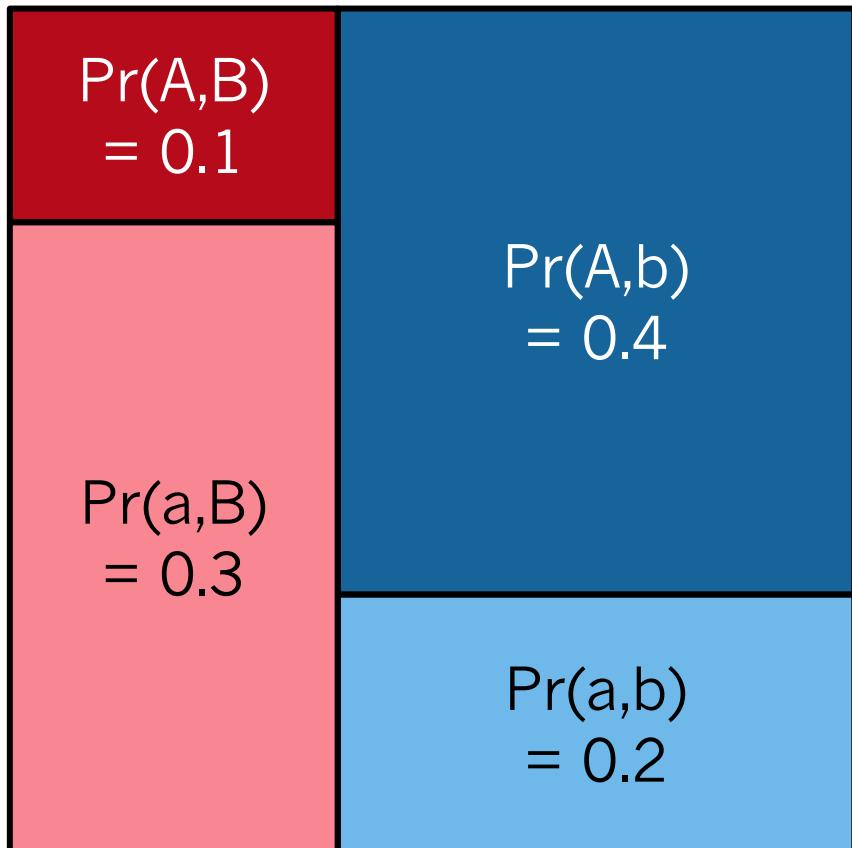
1. Introduction & Motivation
2. What's in a model?
3. Probabilistic Graphical Models (PGMs)
4. PGMs in Phylogenetics
5. Navigating model space

(break)

A quick probability refresher



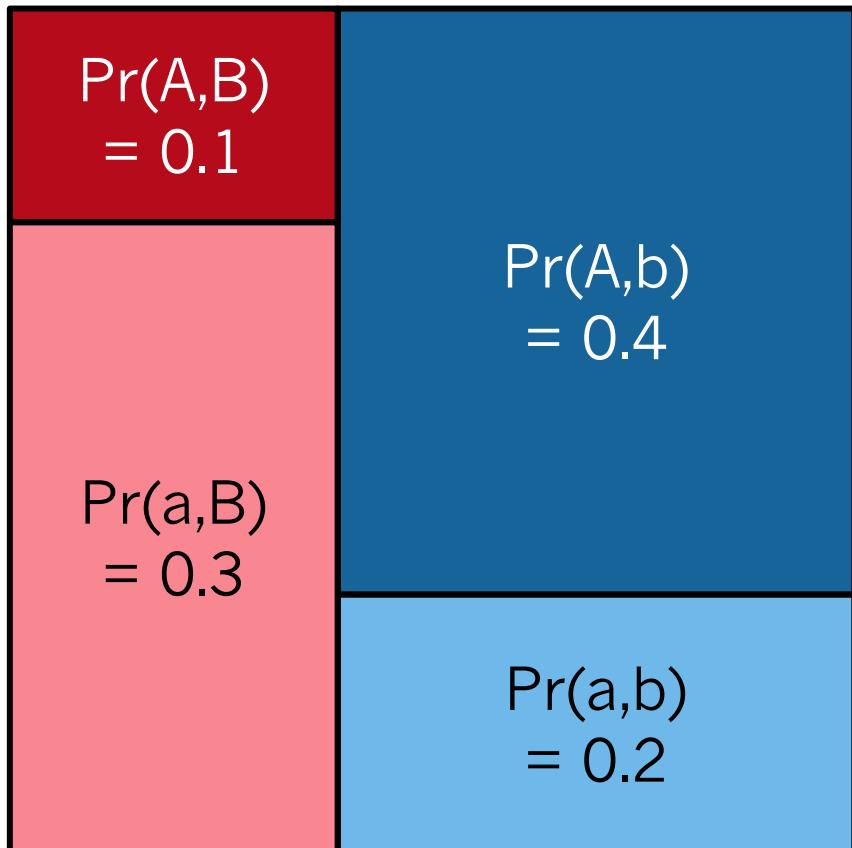
A quick probability refresher



- A : You love to eat lobster
a : You hate to eat lobster
- B : You are a lobster
b : You are *not* a lobster



Probability rules



Each outcome probability is from 0 to 1

All outcome probabilities sum to 1

Every outcome is assigned a single probability

Joint probability of A and B

$\Pr(A,B)$
= 0.1

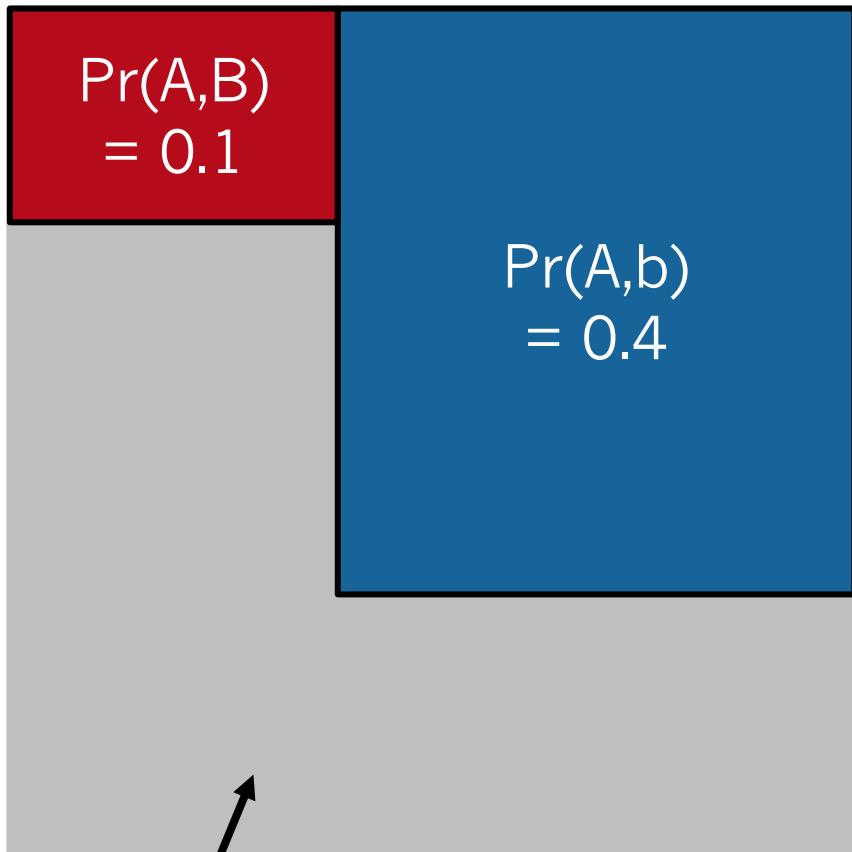
$\Pr(A,B) = 0.1$

Joint probability of A and b

$$\Pr(A,b) = 0.4$$

$$\Pr(A,b) = 0.4$$

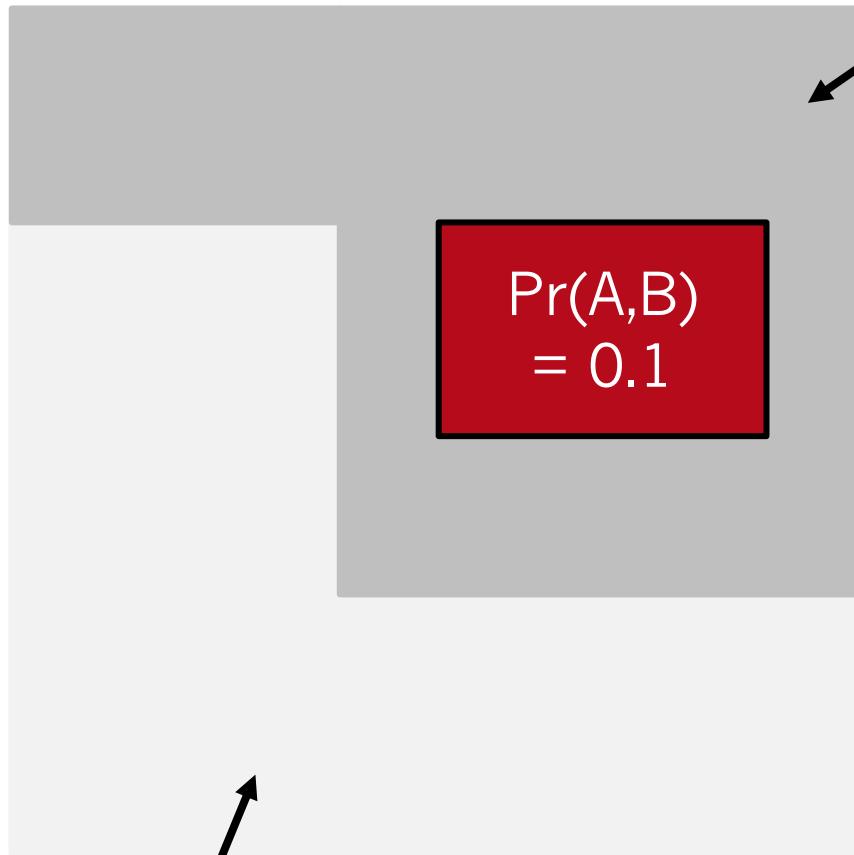
Marginal probability of A



space of all
possible events

$$\begin{aligned}\Pr(A) &= \Pr(A,B) + \Pr(A,b) \\ &= 0.1 + 0.4 \\ &= 0.5\end{aligned}$$

Conditional probability of B given A



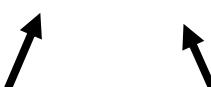
space for event "A"
(considered)

$$\begin{aligned}\Pr(B|A) &= \Pr(A,B) / \Pr(A) \\ &= 0.1 / (0.1 + 0.4) \\ &= 0.2\end{aligned}$$

space for event "a"
(ignored)

Probability distributions

“probability of X given θ ”

$$\Pr(\textcolor{red}{X} | \textcolor{blue}{\theta})$$


random variable parameter

Probability distributions

“probability of X given θ ”

$$\Pr(\textcolor{red}{X} | \textcolor{blue}{\theta})$$

random variable parameter

“ X has an Exponential distribution”

$$\textcolor{red}{X} \sim \text{Exponential}(\textcolor{blue}{\theta})$$

random variable parameter

Joint distributions as models

Joint probability

$$\Pr(\textcolor{red}{X}, \textcolor{brown}{Y}, \theta, \mu)$$

Joint distributions as models

Joint probability

$$\Pr(\textcolor{red}{X}, \textcolor{brown}{Y}, \theta, \mu)$$

Joint probability, factored

$$\Pr(\textcolor{red}{X} | \theta) \Pr(\textcolor{brown}{Y} | \theta, \mu) \Pr(\theta | \mu) \Pr(\mu)$$

Joint distributions as models

Joint probability

$$\Pr(\textcolor{red}{X}, \textcolor{brown}{Y}, \textcolor{blue}{\theta}, \textcolor{violet}{\mu})$$

Joint probability, factored

$$\Pr(\textcolor{red}{X}|\textcolor{blue}{\theta}) \Pr(\textcolor{brown}{Y}|\textcolor{blue}{\theta}, \textcolor{violet}{\mu}) \Pr(\textcolor{blue}{\theta}|\textcolor{violet}{\mu}) \Pr(\textcolor{violet}{\mu})$$

Joint distribution

$$\textcolor{violet}{\mu} \sim \text{Uniform}(0,10)$$

$$\textcolor{blue}{\theta} \sim \text{Exponential}(1)$$

$$\textcolor{red}{X} \sim \text{Exponential}(\textcolor{blue}{\theta})$$

$$\textcolor{brown}{Y} \sim \text{Gamma}(\textcolor{blue}{\theta}, \textcolor{violet}{\mu})$$

Model definitions

Conceptual

Probabilistic

Graphical

Computational

A conceptual model

How might we estimate mutation rate
with K loci from parent-offspring in yeast?

Locus 1 *
Parent TTTATGT
Offspring TCTATGT

Locus 3 * *
Parent AGCTATGGG
Offspring ACCGATGGG

Locus 4
Parent TGAA
Offspring TGAA

Locus 2 *
Parent CTCCG
Offspring CTCAA

Locus K * **
Parent GGTCATTCCC
Offspring GGTCATAATGGC

A conceptual model

How might we estimate mutation rate
with K loci from parent-offspring in yeast?

Number of mutations should be related to
total number of sites and mutation rate

A conceptual model

How might we estimate mutation rate
with K loci from parent-offspring in yeast?

Number of mutations should be related to
total number of sites and mutation rate

The rate should also be similar to estimates
for other eukaryotes

Model definitions

Conceptual

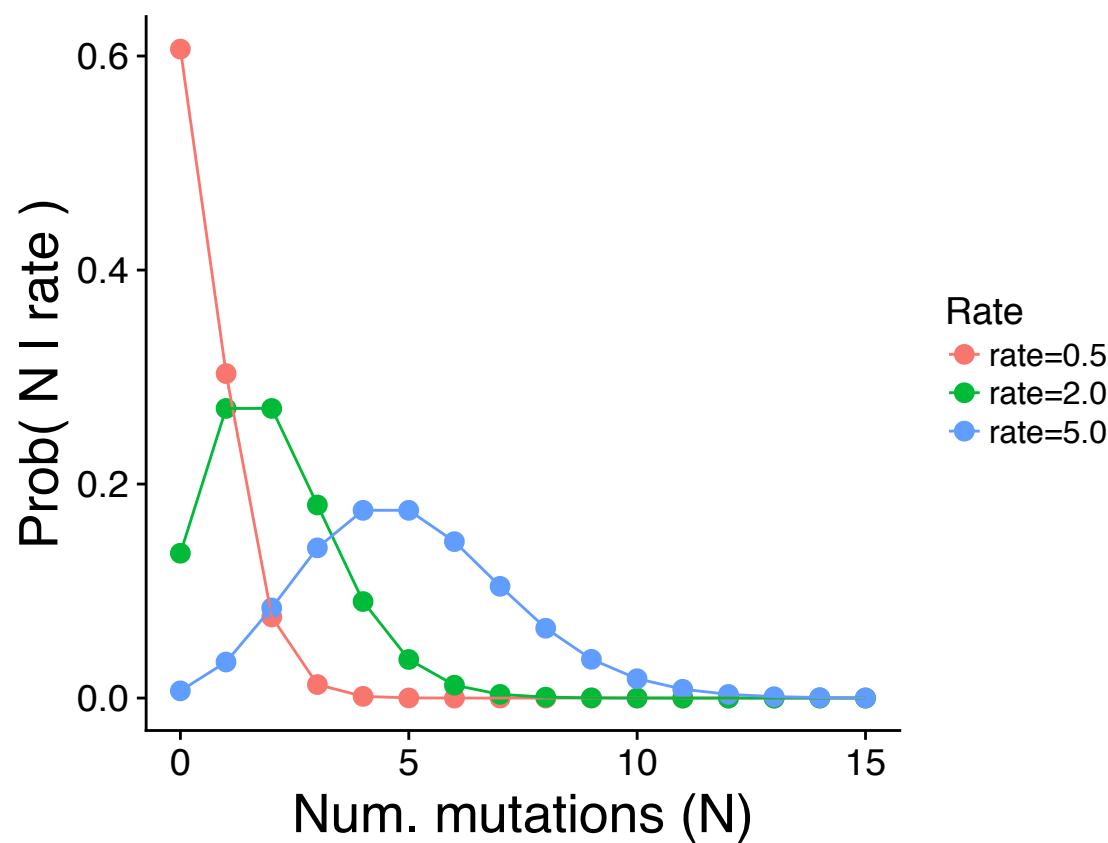
Probabilistic

Graphical

Computational

A probabilistic model

Model per-locus mutation counts with Poisson
 $N \sim \text{Poisson}(\lambda)$



A probabilistic model

Model per-locus mutation counts with Poisson
 $N \sim \text{Poisson}(\lambda)$

Per-locus mutation rate is the per-site mutation rate times the number of sites in the locus

$$\lambda = r \times L$$

A probabilistic model

Model per-locus mutation counts with Poisson
 $N \sim \text{Poisson}(\lambda)$

Per-locus mutation rate is the per-site mutation rate times the number of sites in the locus

$$\lambda = r \times L$$

Estimate r using an empirical prior (eukaryotes)
 $r \sim \text{Exponential}(1/10^{-8})$

A probabilistic model

Model per-locus mutation counts with Poisson
 $N \sim \text{Poisson}(\lambda)$

Per-locus mutation rate is the per-site mutation rate times the number of sites in the locus

$$\lambda = r \times L$$

Estimate r using an empirical prior (eukaryotes)
 $r \sim \text{Exponential}(1/10^{-8})$

Substitute N_k , λ_k , and L_k to generalize for K loci

A probabilistic model

Model probability

$$\Pr(N, r \mid L, \mu) = \Pr(r \mid \mu) \prod_k \Pr(N_k \mid r, L_k)$$

A probabilistic model

Model probability

$$\Pr(N, r \mid L, \mu) = \Pr(r \mid \mu) \prod_k \Pr(N_k \mid r, L_k)$$

Model distribution

$$r \sim \text{Exponential}(1/\mu)$$

$$\lambda_k = r \times L_k$$

$$N_k \sim \text{Poisson}(\lambda_k)$$

Model definitions

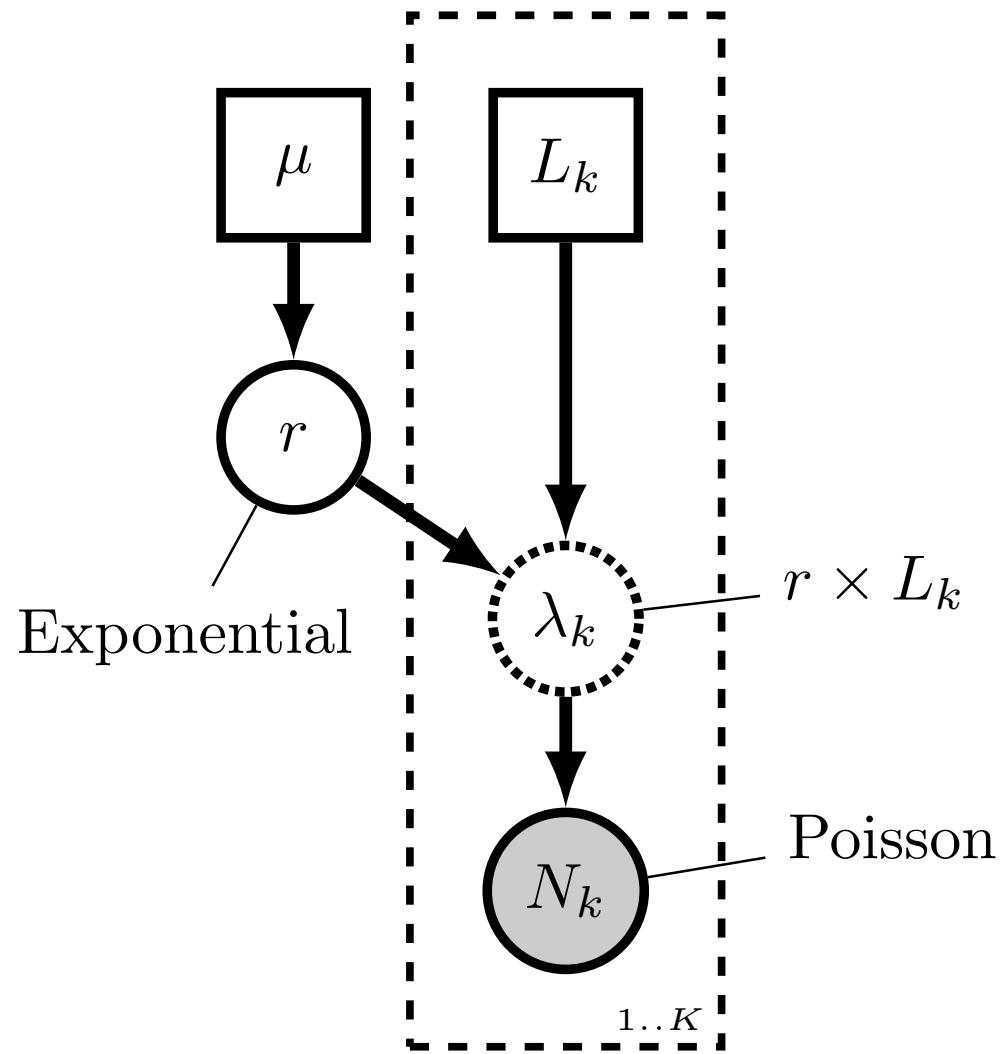
Conceptual

Probabilistic

Graphical

Computational

A probabilistic graphical model



probabilistic model

a set of random variables and dependencies

+

graph

a set of nodes and edges

=

probabilistic graphical model (PGM)

*a probabilistic model that uses a graph
to encode its variables and their dependencies*

probabilistic graphical model (PGM)

*a probabilistic model that uses a graph
to encode its variables and their dependencies*

Uses

Algorithm design

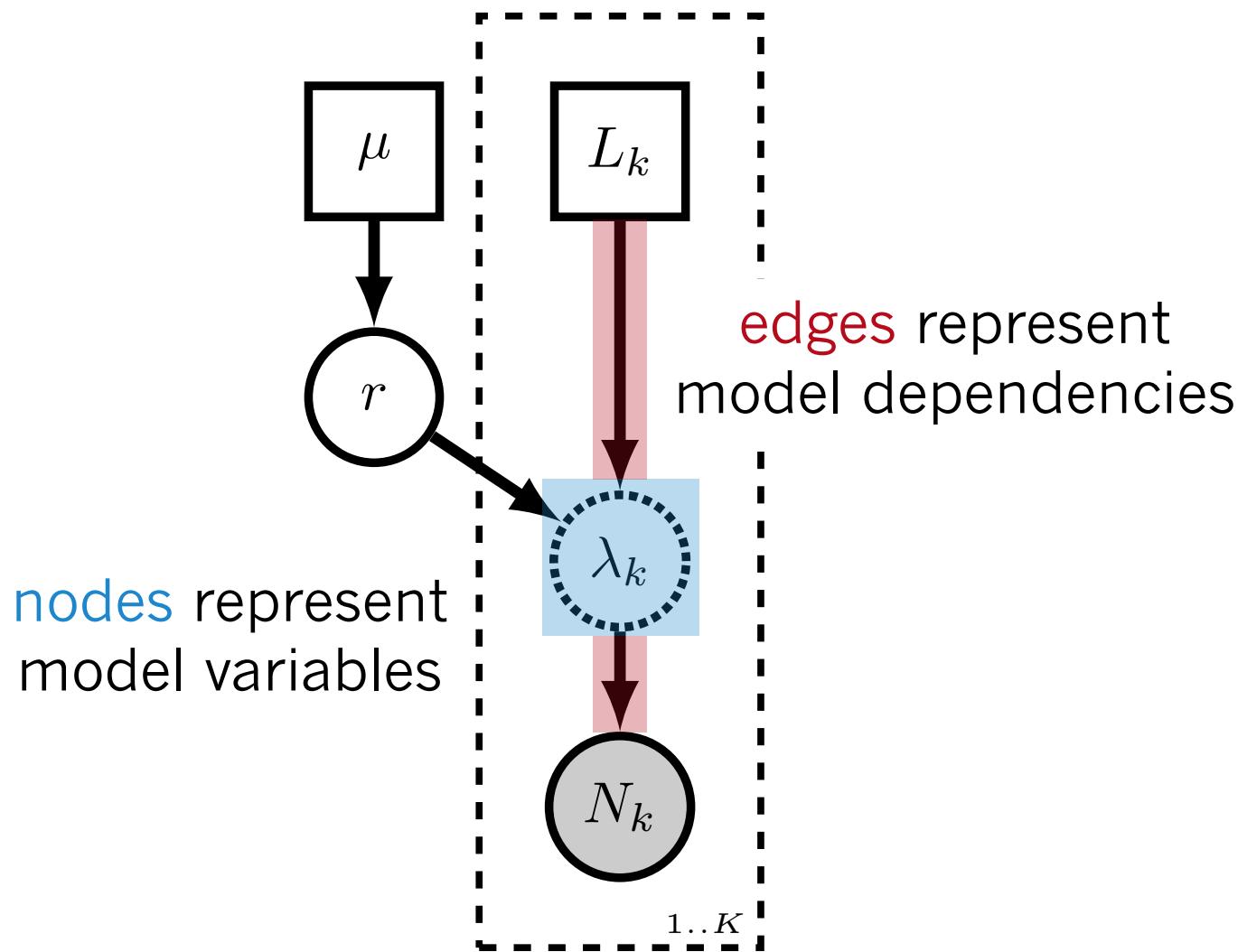
Model design

Blueprinting ideas

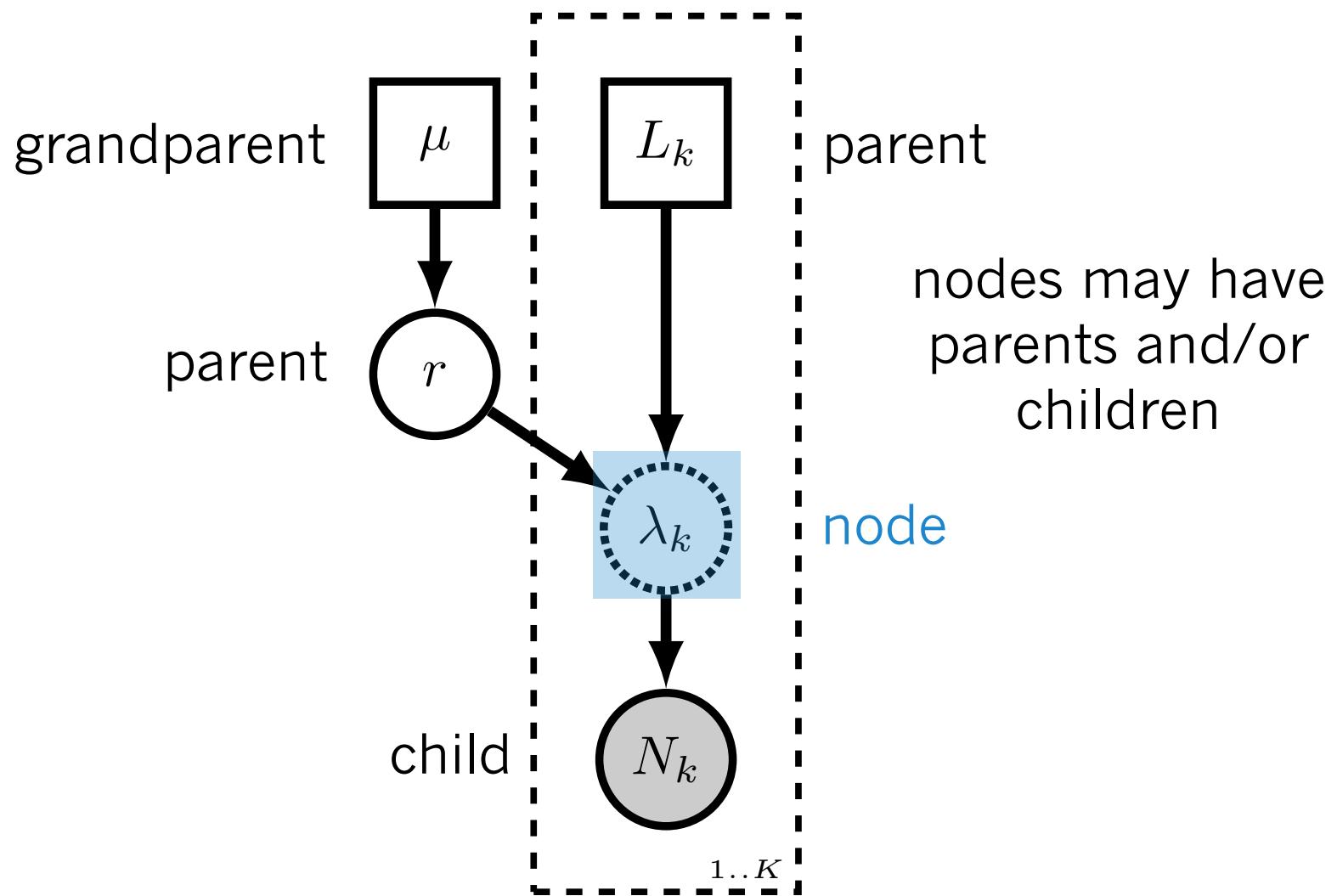
Communication

Teaching

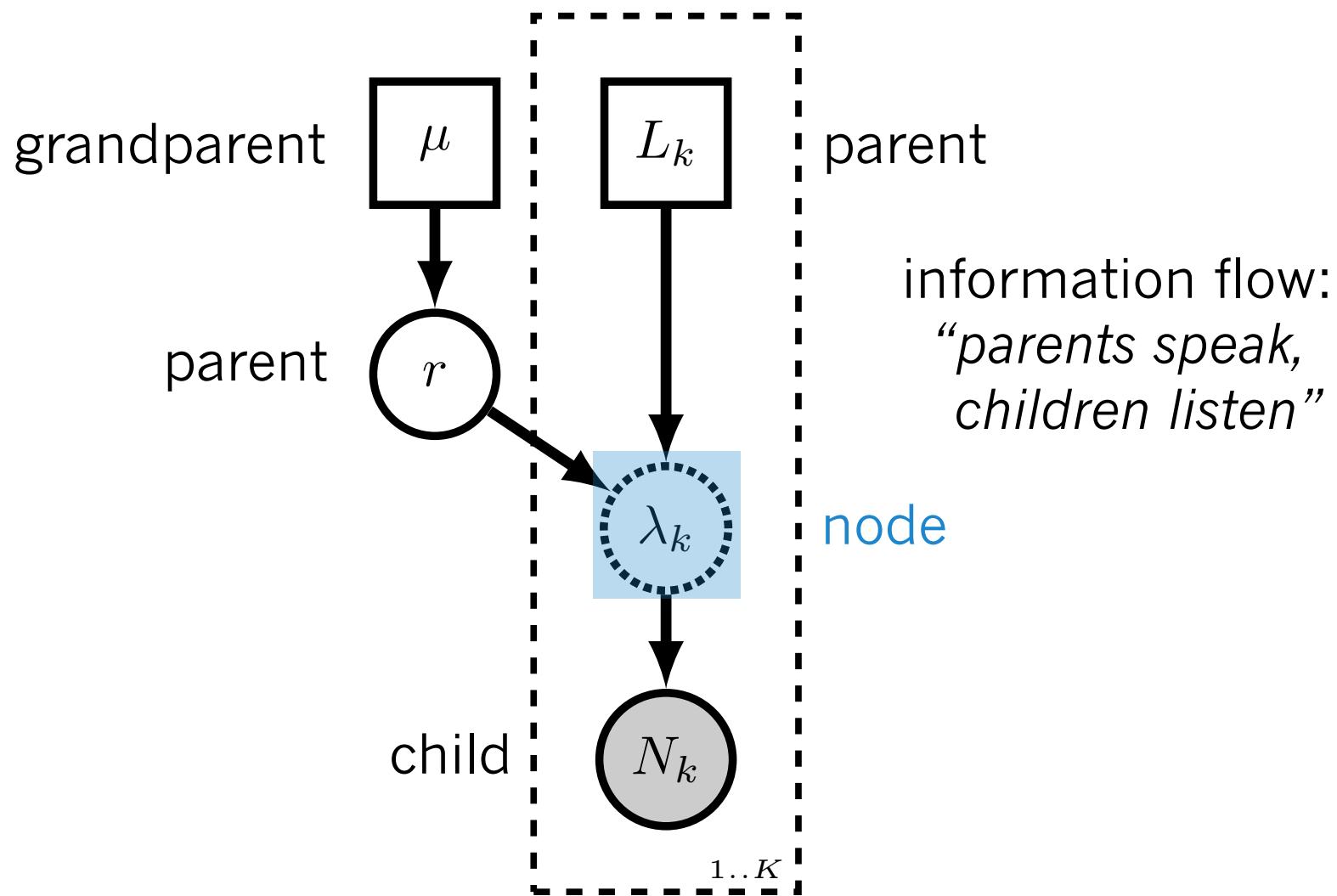
PGM anatomy



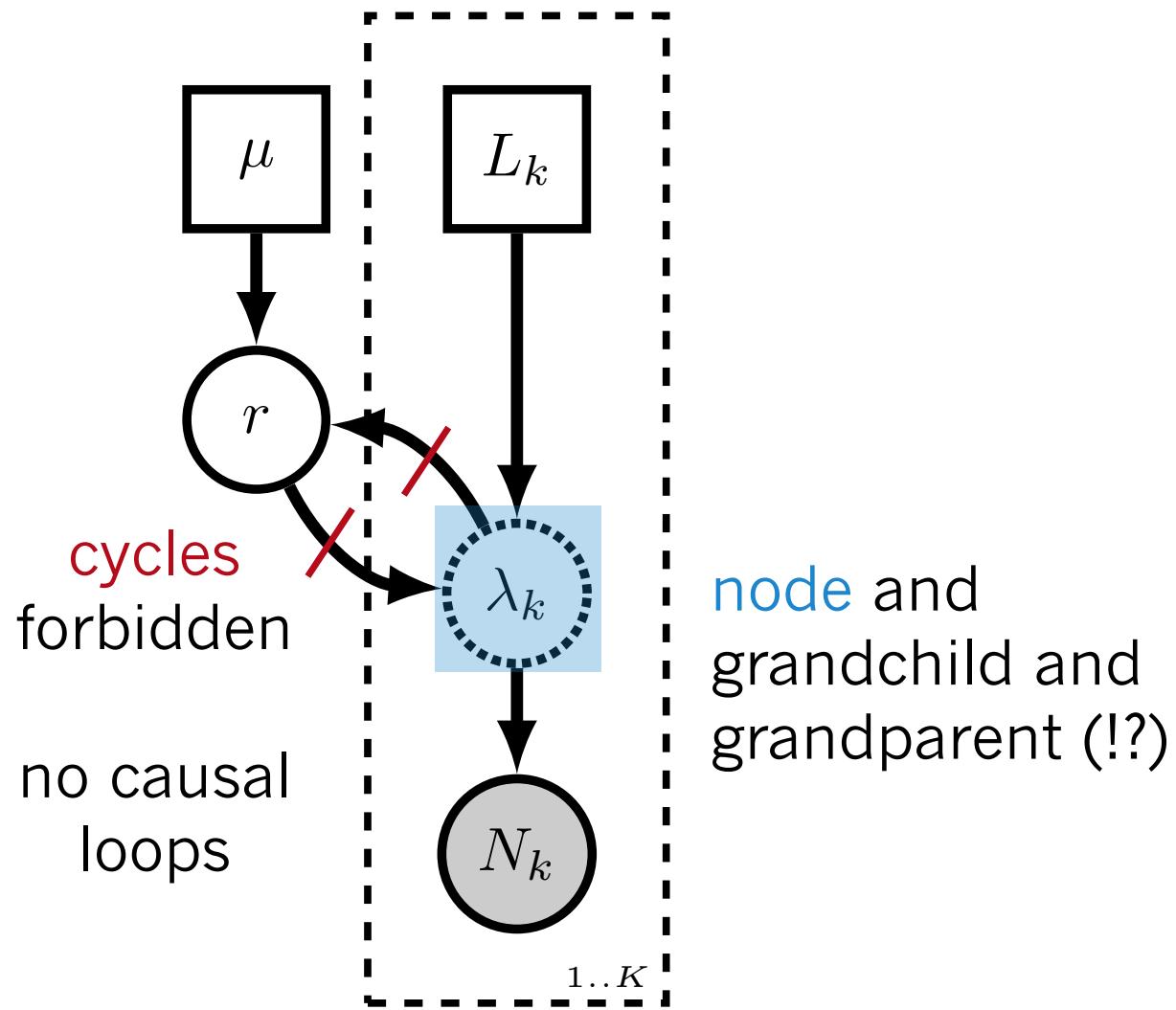
PGM anatomy



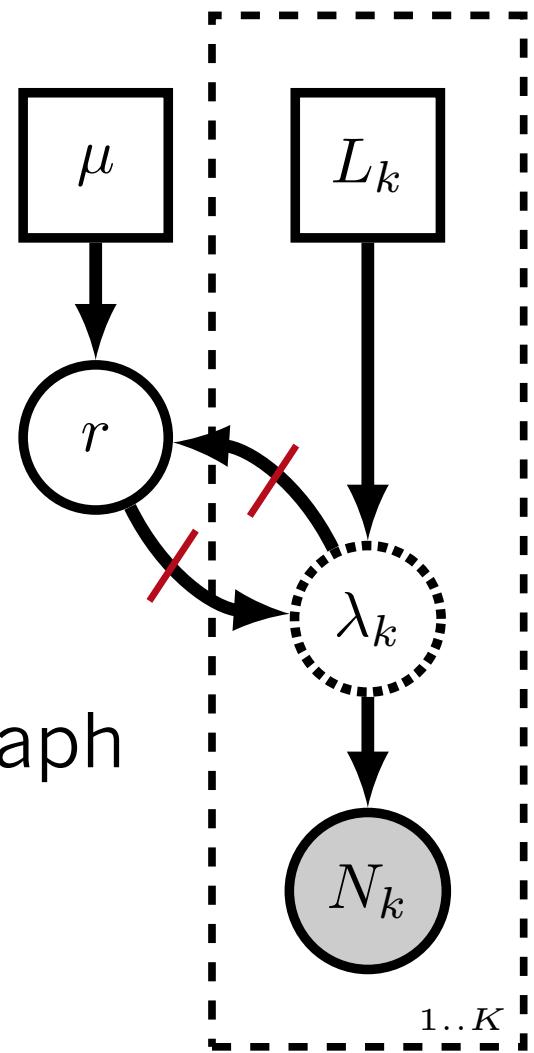
PGM anatomy



PGM anatomy

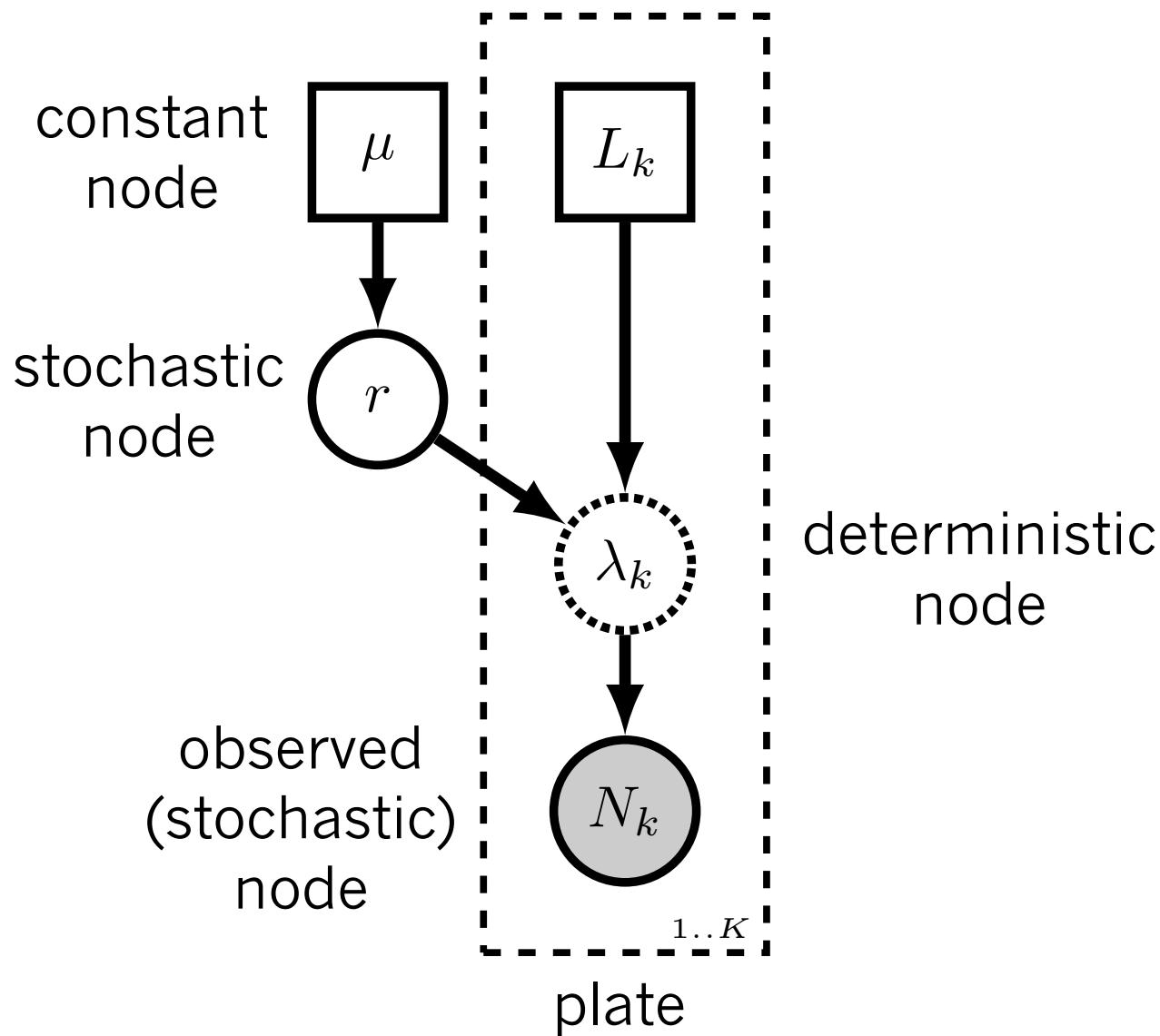


PGM anatomy



Directed acyclic graph
(DAG)

PGM anatomy

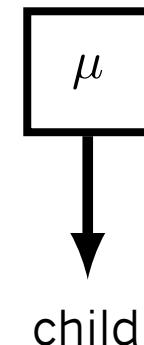


Constant node

“equals”
fixed value
asserted or known

the rate for a prior

$$\mu = 1$$



Constant node

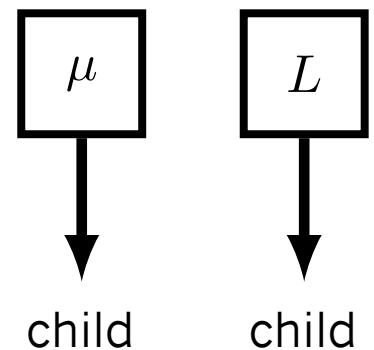
“equals”
fixed value
asserted or known

the rate for a prior

$$\mu = 1$$

the length of a locus

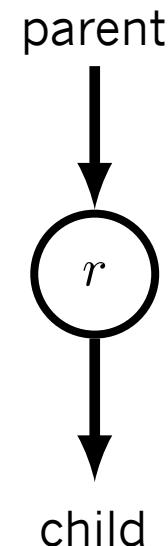
$$L = 1030$$



Stochastic node

“distributed by”
dynamic value
learned

per-site mutation rate
 $r \sim \text{Exponential}(\mu)$

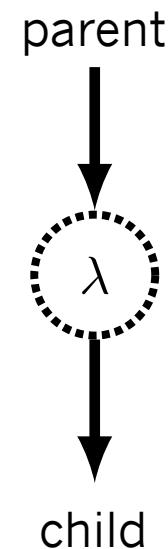


Deterministic node

“determined by”
dynamic value
learned

per-locus mutation rate

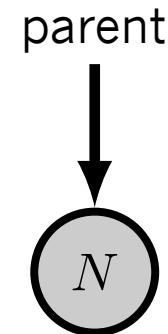
$$\lambda := L \times r$$



Observed (stochastic) node

“distributed by”
observed value
fixed (at observation)

de novo mutations (1 locus)
 $N \sim \text{Poisson}(\lambda)$

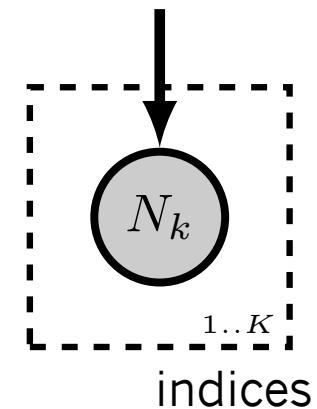


Plates

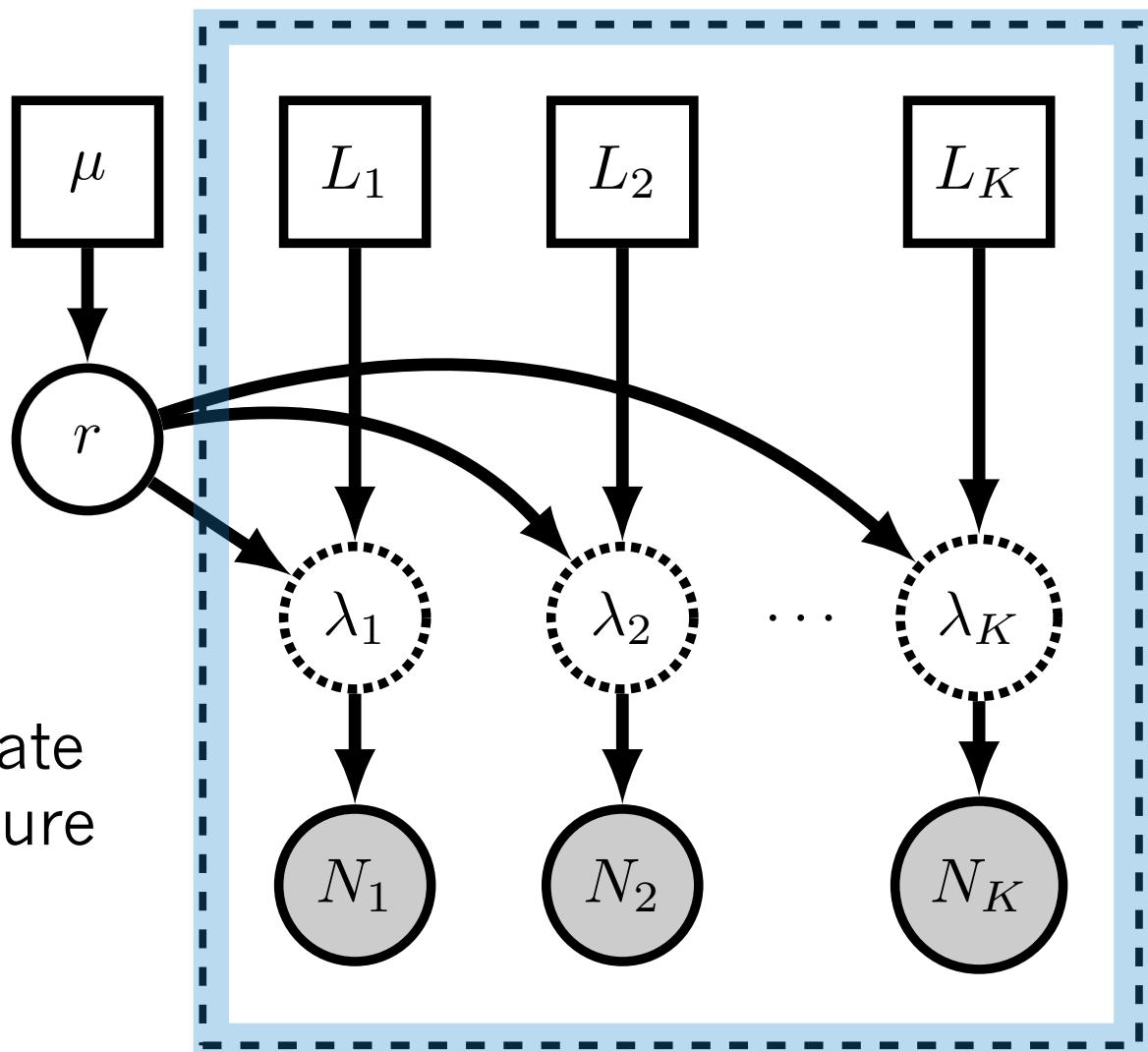
repeats model structure
simplifies graph

de novo mutations (K loci)

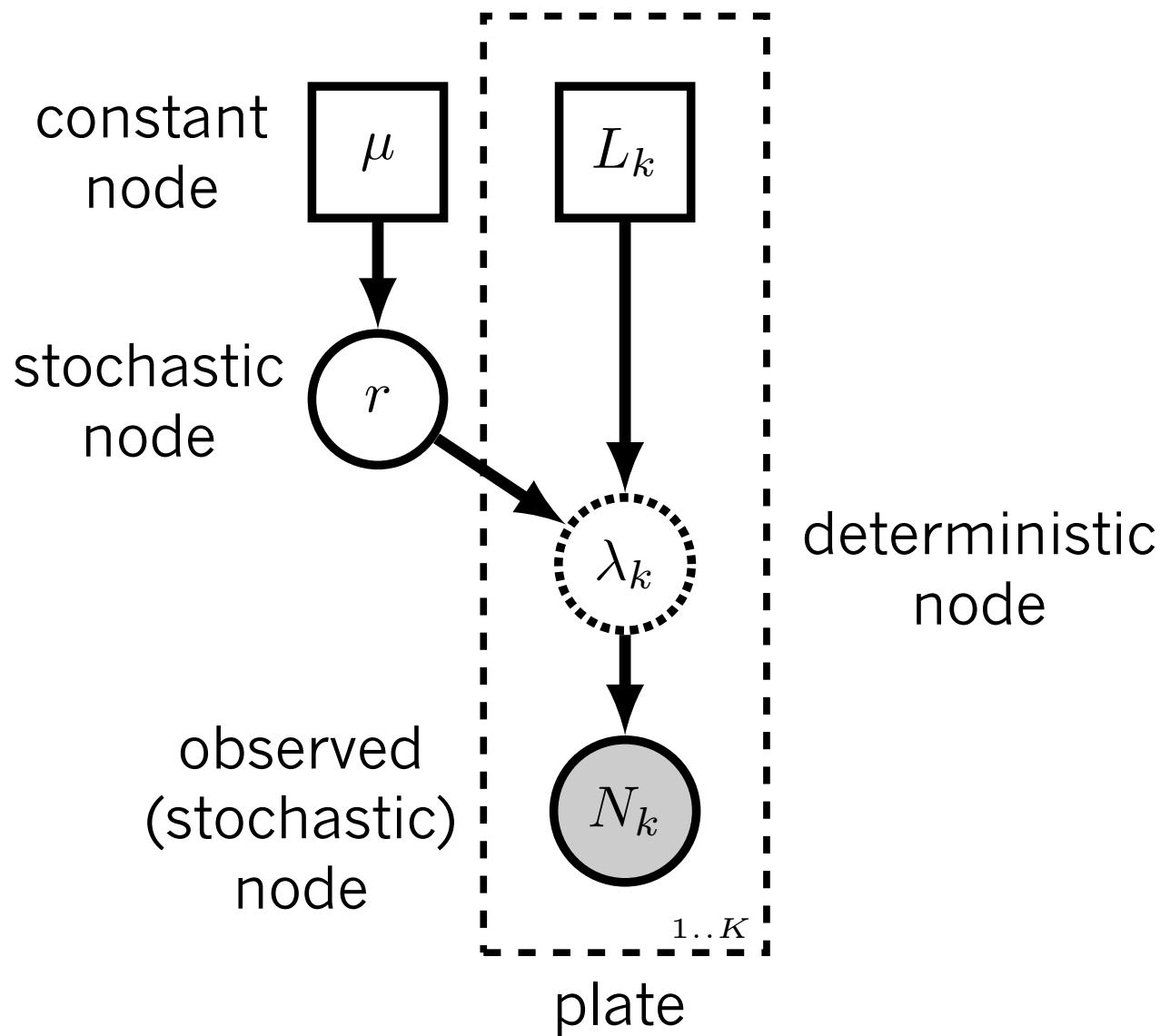
$$N_k \sim \text{Poisson}(\lambda_k)$$



PGM anatomy



PGM anatomy

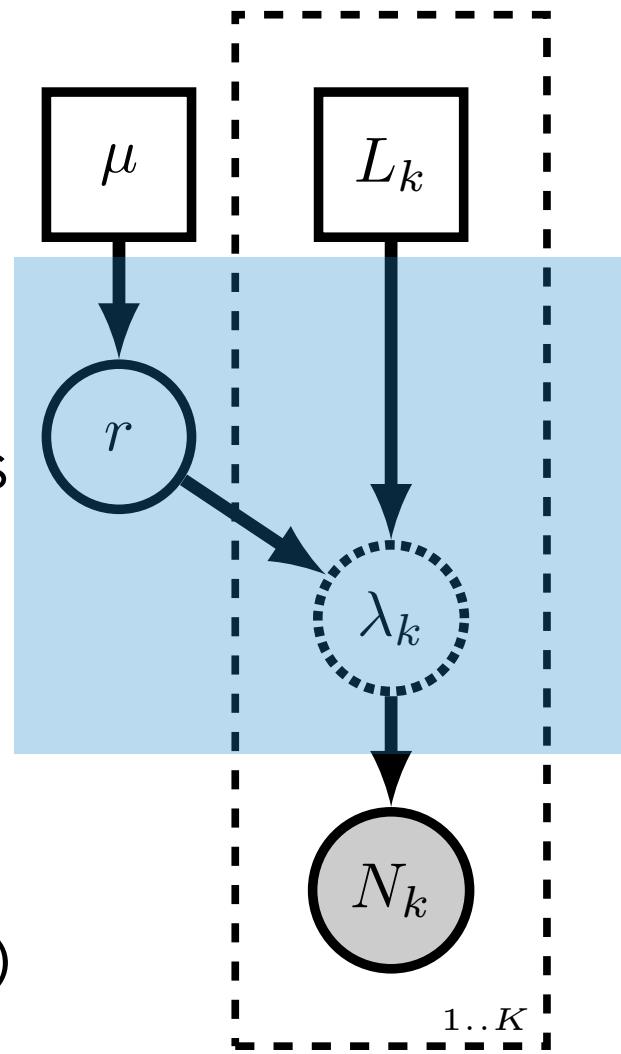


PGM anatomy

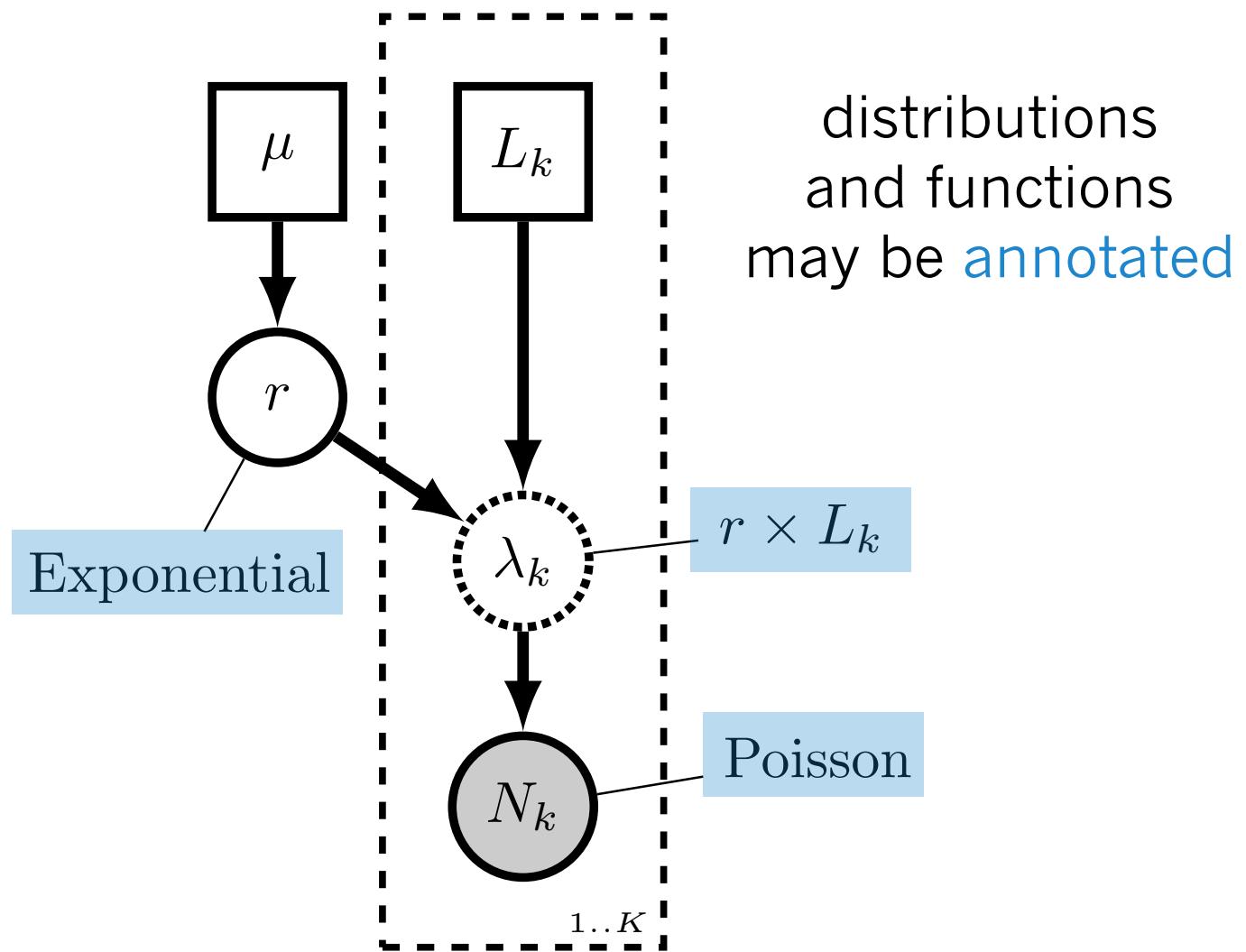
fixed
(given)

unknown values
are estimated

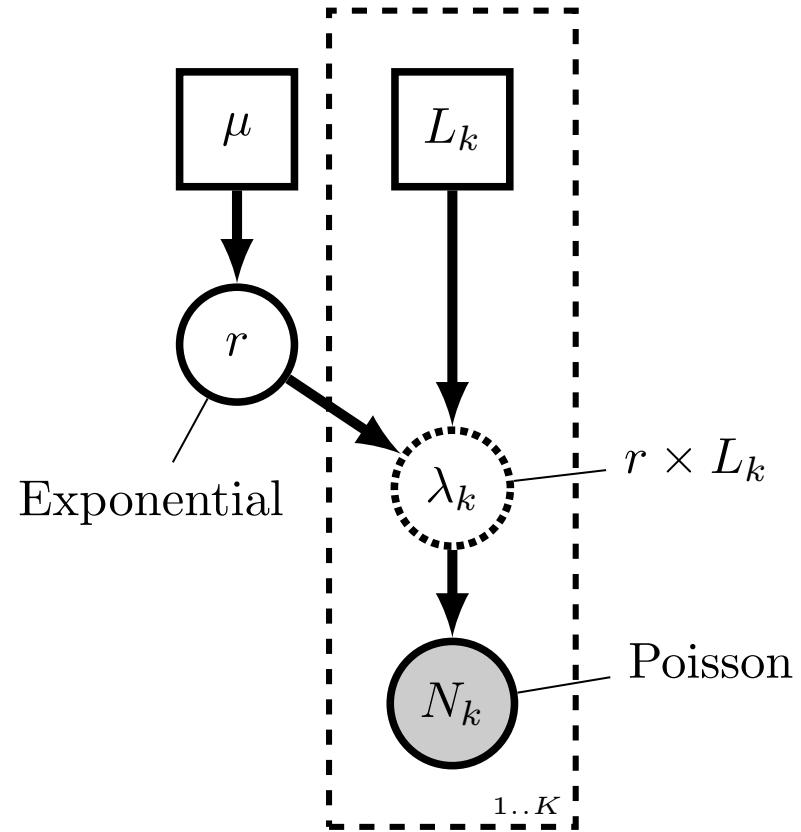
fixed
(to observation)



PGM anatomy



Example of computing model probability using PGM

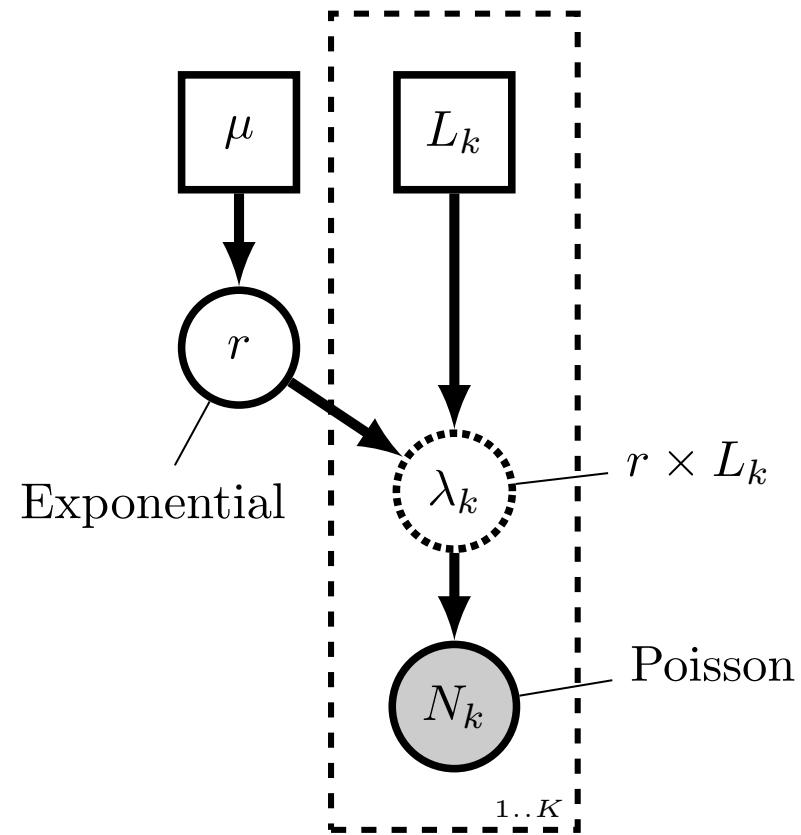


One locus ($K=1$)

$N_1 = 3$ mutations

$L_1 = 1030\text{bp}$

$r = ???$

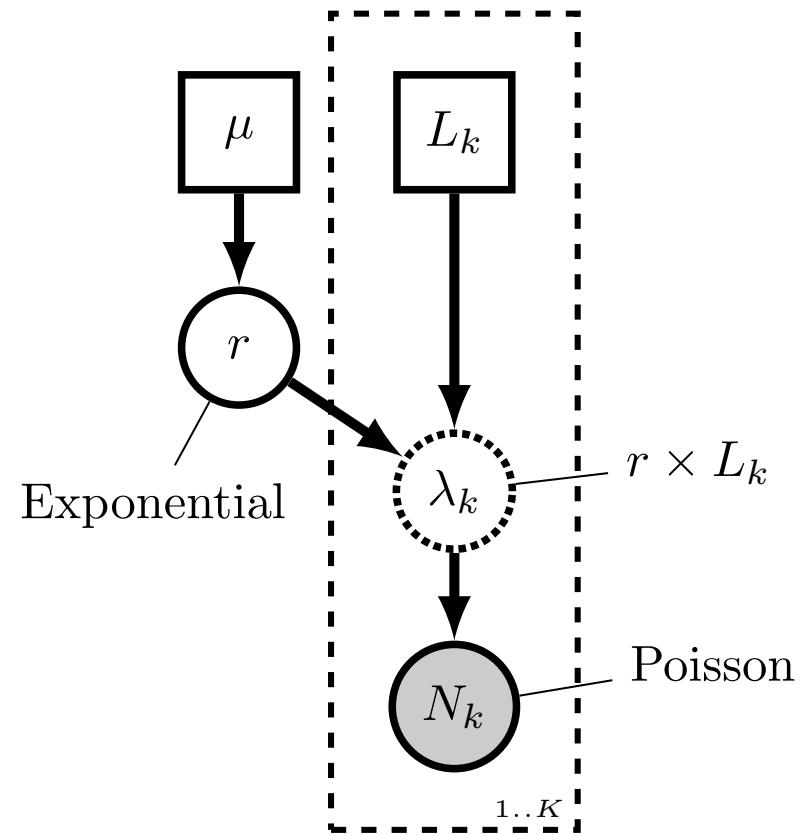


One locus ($K=1$)

$N_1 = 3$ mutations

$L_1 = 1030\text{bp}$

$r = 3.0 \times 10^{-7}$



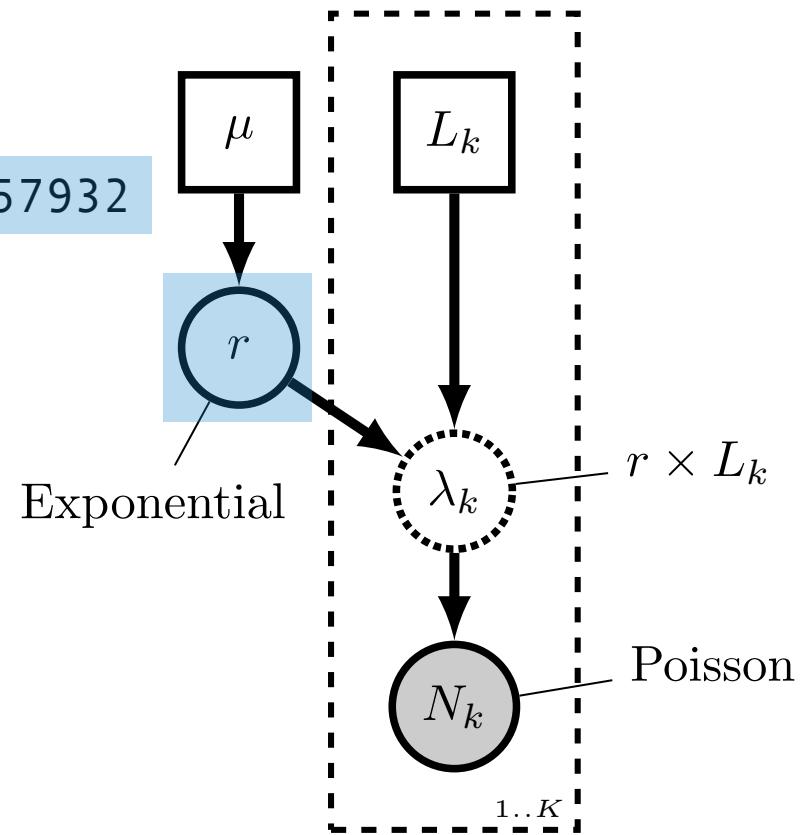
One locus ($K=1$)

$$N_1 = 3 \text{ mutations}$$

$$L_1 = 1030\text{bp}$$

$$r = 3.0 \times 10^{-7}$$

$$\ln \Pr(r=3.0 \times 10^{-7} \mid \mu=1/10^{-8}) = -11.57932$$



One locus ($K=1$)

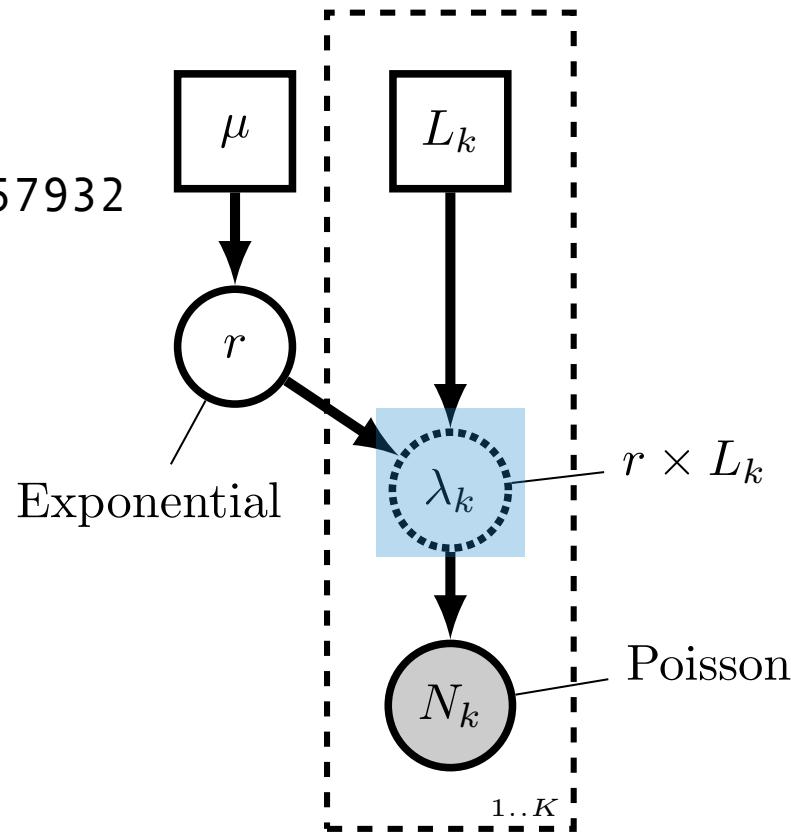
$$N_1 = 3 \text{ mutations}$$

$$L_1 = 1030 \text{ bp}$$

$$r = 3.0 \times 10^{-7}$$

$$\ln \Pr(r=3.0 \times 10^{-7} \mid \mu=1/10^{-8}) = -11.57932$$

$$\lambda_1 = 3.0 \times 10^{-7} * 1030 = 0.0003090$$



One locus ($K=1$)

$$N_1 = 3 \text{ mutations}$$

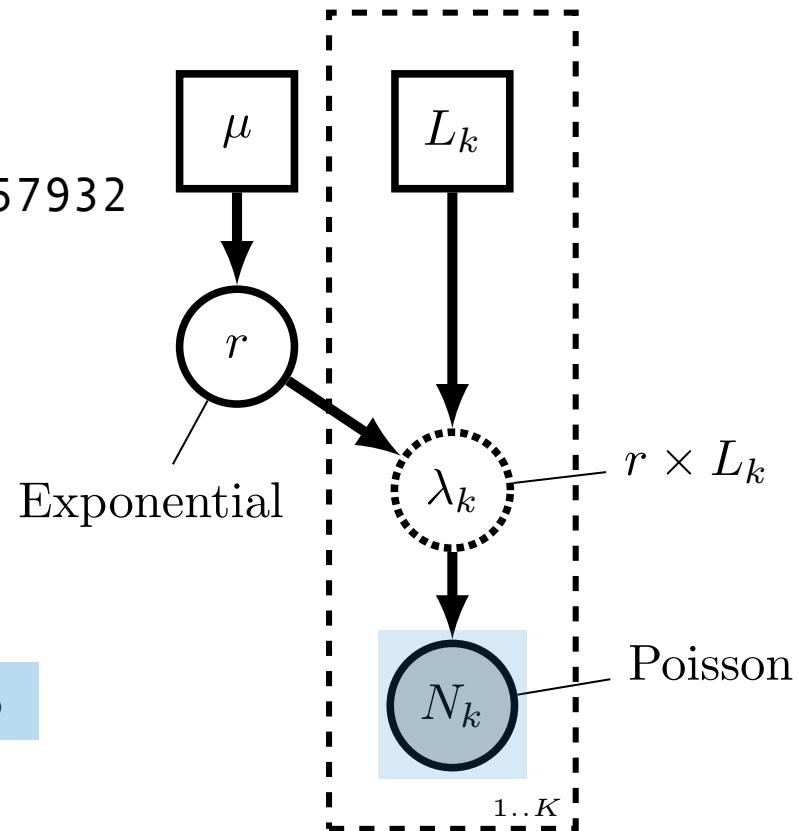
$$L_1 = 1030 \text{ bp}$$

$$r = 3.0 \times 10^{-7}$$

$$\ln \Pr(r=3.0 \times 10^{-7} \mid \mu=1/10^{-8}) = -11.57932$$

$$\lambda_1 = 3.0 \times 10^{-7} * 1030 = 0.0003090$$

$$\ln \Pr(N_1=3 \mid \lambda_1=0.0003090) = -26.0386$$



One locus ($K=1$)

$N_1 = 3$ mutations

$L_1 = 1030\text{bp}$

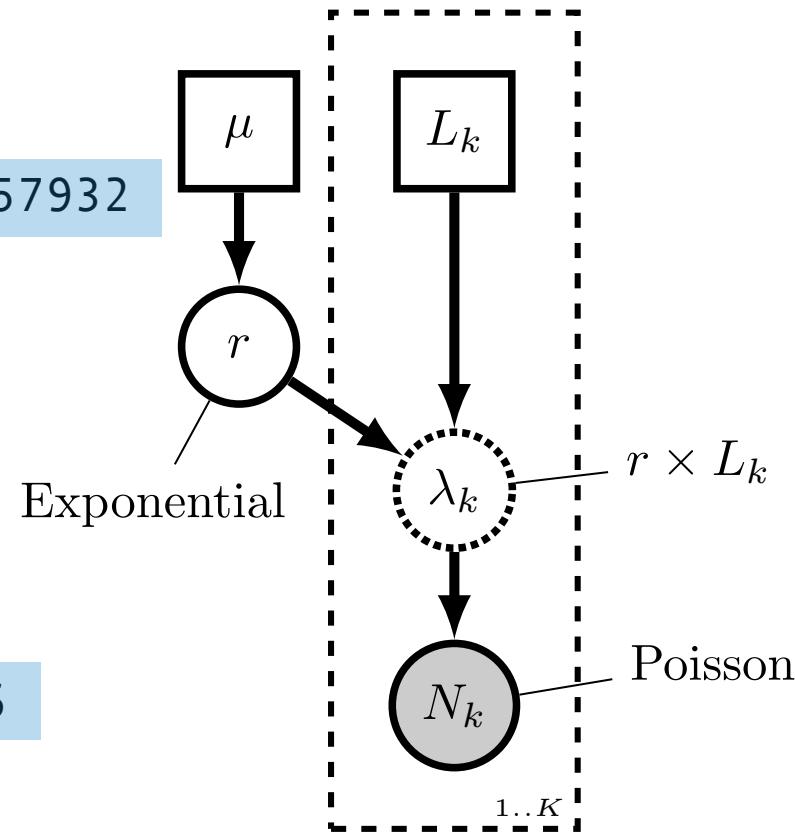
$r = 3.0 \times 10^{-7}$

$$\ln \Pr(r=3.0 \times 10^{-7} \mid \mu=1/10^{-8}) = -11.57932$$

$$\lambda_1 = 3.0 \times 10^{-7} * 1030 = 0.0003090$$

$$\ln \Pr(N_1=3 \mid \lambda_1=0.0003090) = -26.0386$$

$$\ln \Pr(N_1=3, r=3.0 \times 10^{-7} \mid \lambda_1=0.0003090, \mu=1/10^{-8}) = -37.6179$$



model log-probability

One locus ($K=1$)

$N_1 = 3$ mutations

$L_1 = 1030\text{bp}$

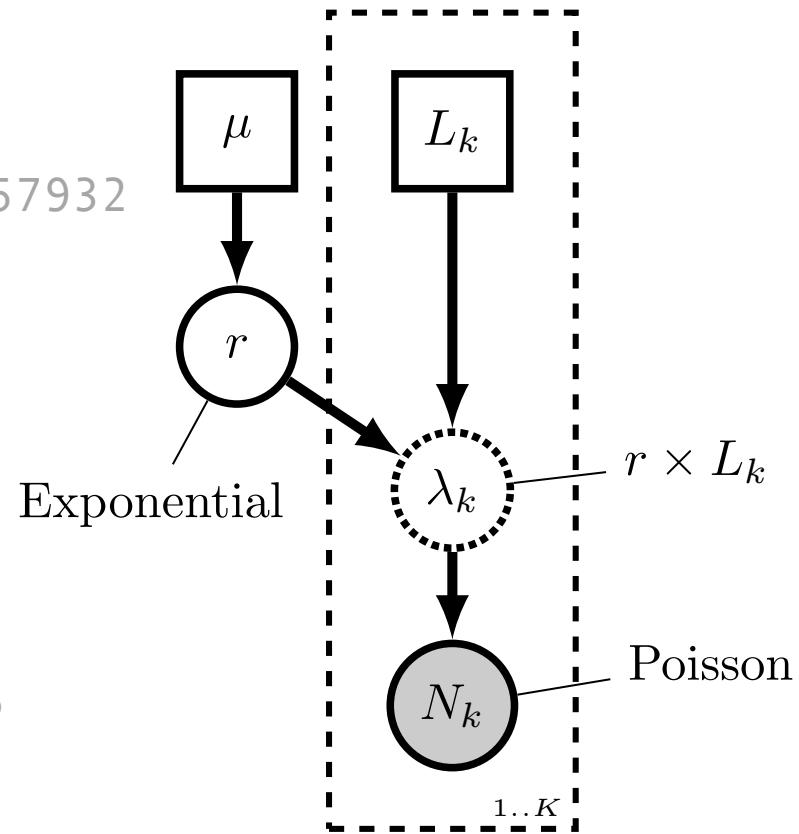
$r = 2.7 \times 10^{-7}$

$$\ln \Pr(r=3.0 \times 10^{-7} \mid \mu=1/10^{-8}) = -11.57932$$

$$\lambda_1 = 3.0 \times 10^{-7} * 1030 = 0.0003090$$

$$\ln \Pr(N_1=3 \mid \lambda_1=0.0003090) = -26.0386$$

$$\ln \Pr(N_1=3, r=3.0 \times 10^{-7} \mid \lambda_1=0.0003090, \mu=1/10^{-8}) = -37.6179$$



One locus ($K=1$)

$N_1 = 3$ mutations

$L_1 = 1030\text{bp}$

$r = 2.7 \times 10^{-7}$

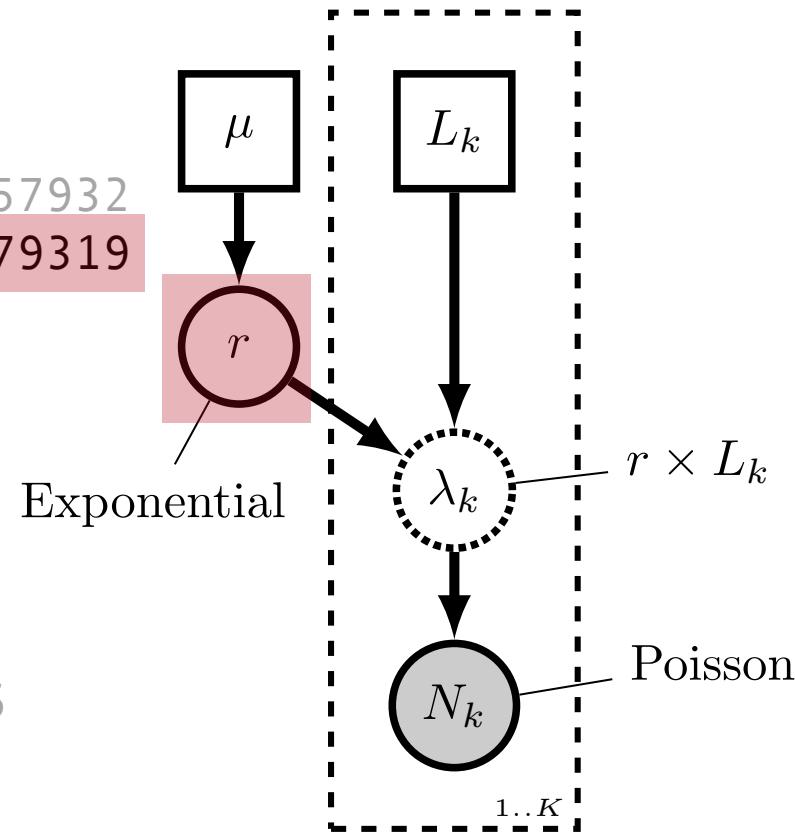
$$\ln \Pr(r=3.0 \times 10^{-7} \mid \mu=1/10^{-8}) = -11.57932$$

$$\ln \Pr(r=2.7 \times 10^{-7} \mid \mu=1/10^{-8}) = -8.579319$$

$$\lambda_1 = 3.0 \times 10^{-7} * 1030 = 0.0003090$$

$$\ln \Pr(N_1=3 \mid \lambda_1=0.0003090) = -26.0386$$

$$\ln \Pr(N_1=3, r=3.0 \times 10^{-7} \mid \lambda_1=0.0003090, \mu=1/10^{-8}) = -37.6179$$



One locus ($K=1$)

$N_1 = 3$ mutations

$L_1 = 1030\text{bp}$

$r = 2.7 \times 10^{-7}$

$$\ln \Pr(r=3.0 \times 10^{-7} \mid \mu=1/10^{-8}) = -11.57932$$

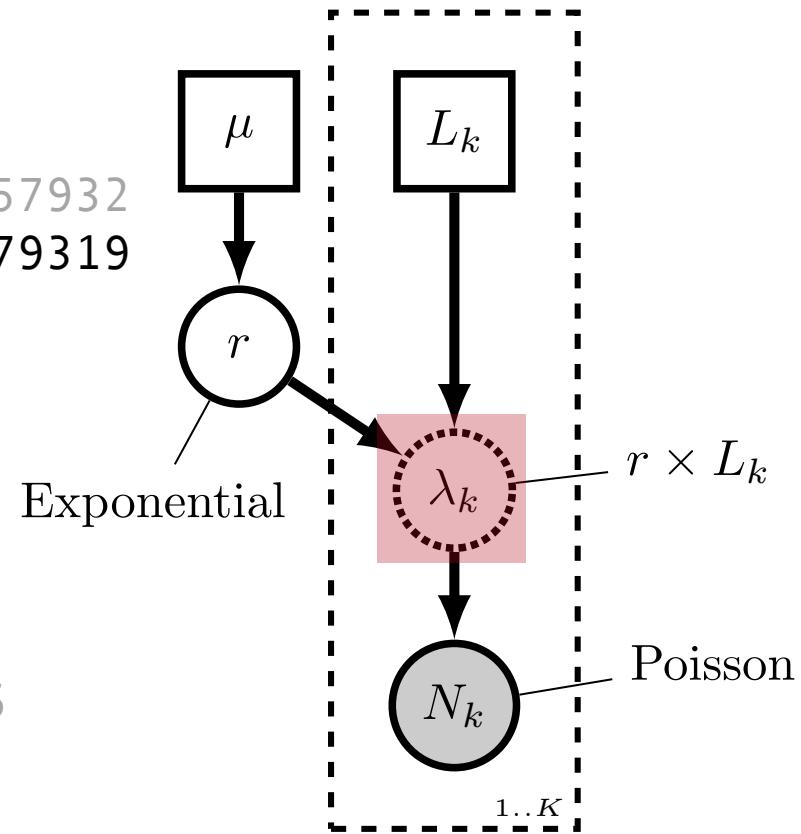
$$\ln \Pr(r=2.7 \times 10^{-7} \mid \mu=1/10^{-8}) = -8.579319$$

$$\lambda_1 = 3.0 \times 10^{-7} * 1030 = 0.0003090$$

$$\lambda_1 = 2.7 \times 10^{-7} * 1030 = 0.0002781$$

$$\ln \Pr(N_1=3 \mid \lambda_1=0.0003090) = -26.0386$$

$$\ln \Pr(N_1=3, r=3.0 \times 10^{-7} \mid \lambda_1=0.0003090, \mu=1/10^{-8}) = -37.6179$$



One locus ($K=1$)

$N_1 = 3$ mutations

$L_1 = 1030\text{bp}$

$r = 2.7 \times 10^{-7}$

$$\ln \Pr(r=3.0 \times 10^{-7} \mid \mu=1/10^{-8}) = -11.57932$$

$$\ln \Pr(r=2.7 \times 10^{-7} \mid \mu=1/10^{-8}) = -8.579319$$

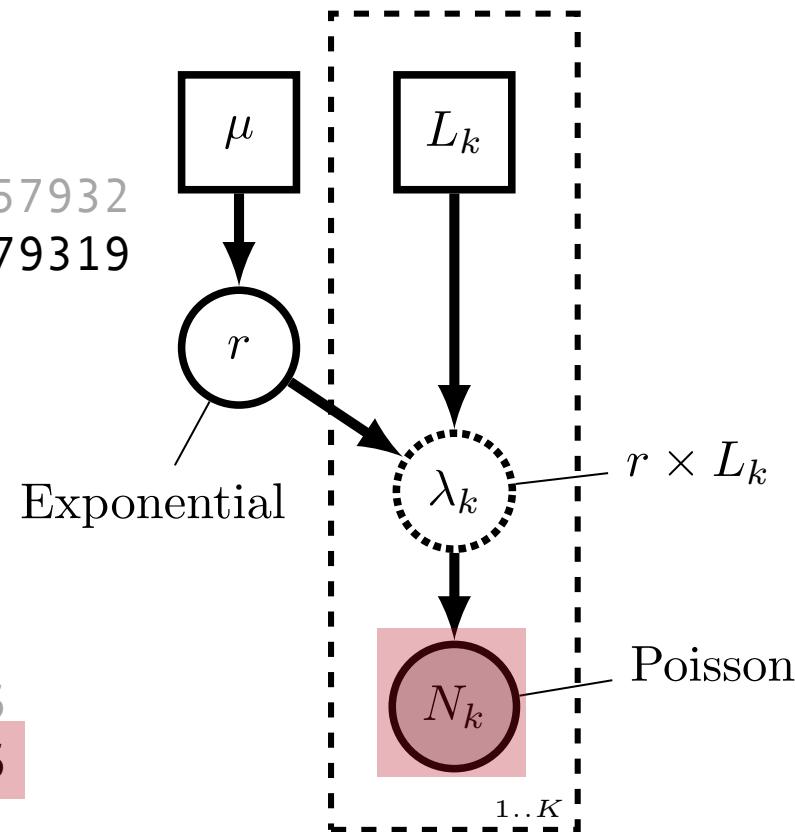
$$\lambda_1 = 3.0 \times 10^{-7} * 1030 = 0.0003090$$

$$\lambda_1 = 2.7 \times 10^{-7} * 1030 = 0.0002781$$

$$\ln \Pr(N_1=3 \mid \lambda_1=0.0003090) = -26.0386$$

$$\ln \Pr(N_1=3 \mid \lambda_1=0.0002781) = -26.3546$$

$$\ln \Pr(N_1=3, r=3.0 \times 10^{-7} \mid \lambda_1=0.0003090, \mu=1/10^{-8}) = -37.6179$$



One locus ($K=1$)

$$N_1 = 3 \text{ mutations}$$

$$L_1 = 1030 \text{ bp}$$

$$r = 2.7 \times 10^{-7}$$

$$\ln \Pr(r=3.0 \times 10^{-7} \mid \mu=1/10^{-8}) = -11.57932$$

$$\ln \Pr(r=2.7 \times 10^{-7} \mid \mu=1/10^{-8}) = -8.579319$$

$$\lambda_1 = 3.0 \times 10^{-7} * 1030 = 0.0003090$$

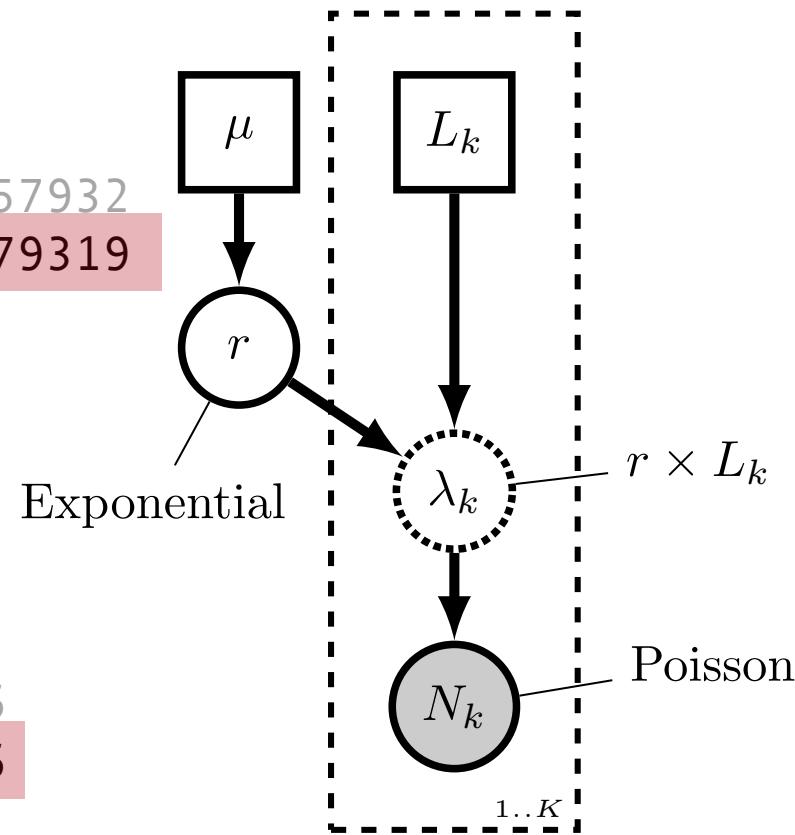
$$\lambda_1 = 2.7 \times 10^{-7} * 1030 = 0.0002781$$

$$\ln \Pr(N_1=3 \mid \lambda_1=0.0003090) = -26.0386$$

$$\ln \Pr(N_1=3 \mid \lambda_1=0.0002781) = -26.3546$$

$$\ln \Pr(N_1=3, r=3.0 \times 10^{-7} \mid \lambda_1=0.0003090, \mu=1/10^{-8}) = -37.6179$$

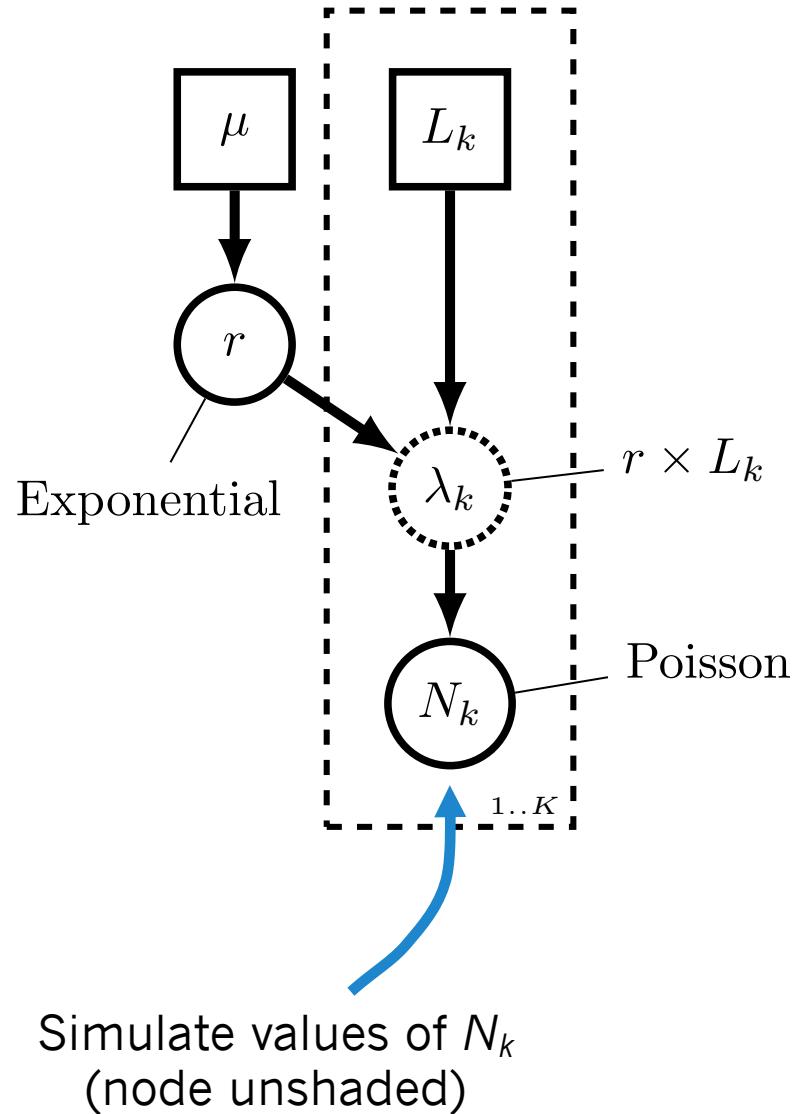
$$\ln \Pr(N_1=3, r=2.7 \times 10^{-7} \mid \lambda_1=0.0002781, \mu=1/10^{-8}) = -34.9340$$



new model log-probability

Generative model

*use same PGM both to
compute model likelihood
and to simulate data*



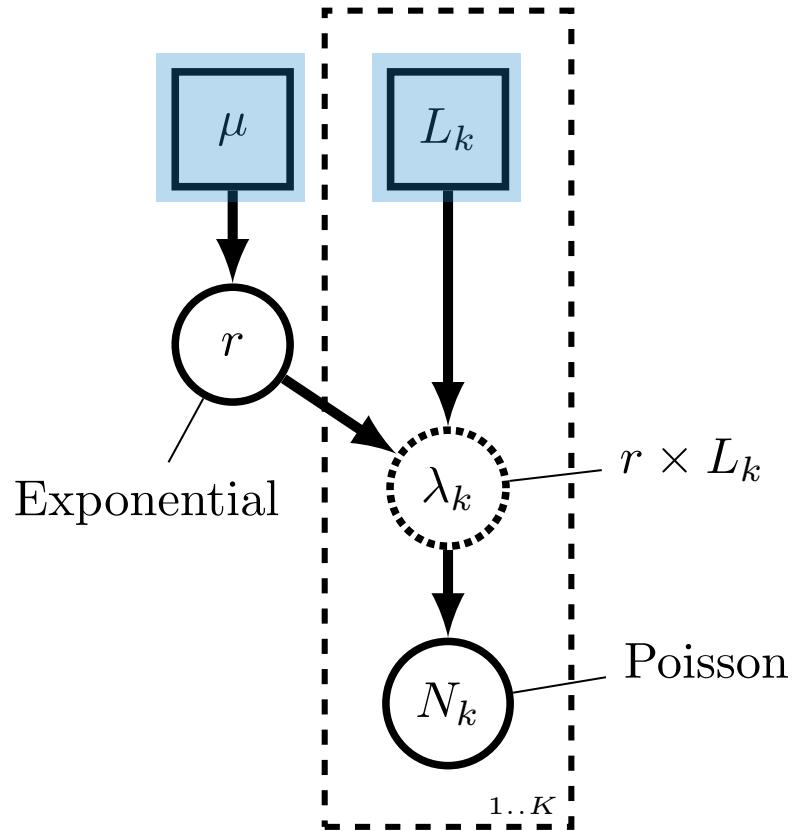
Simulation using DAG structure

use *parent values* to draw *child values*

Set initial conditions (for locus 1)

$$\mu = 1/10^{-8}$$

$$L_1 = 1030 \text{ bp}$$



Simulation using DAG structure

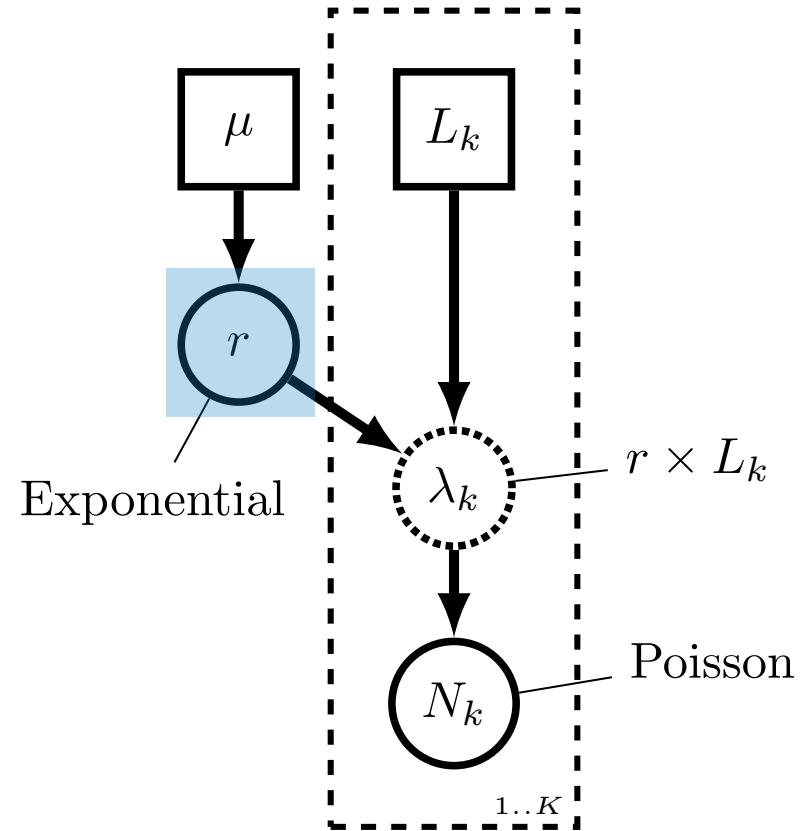
use *parent values* to draw *child values*

Set initial conditions (for locus 1)

$$\mu = 1/10^{-8}$$

$$L_1 = 1030 \text{ bp}$$

Sample $r = 2.5 * 10^{-7}$ from $\text{Exp}(1/10^{-8})$



Simulation using DAG structure

use *parent values* to draw *child values*

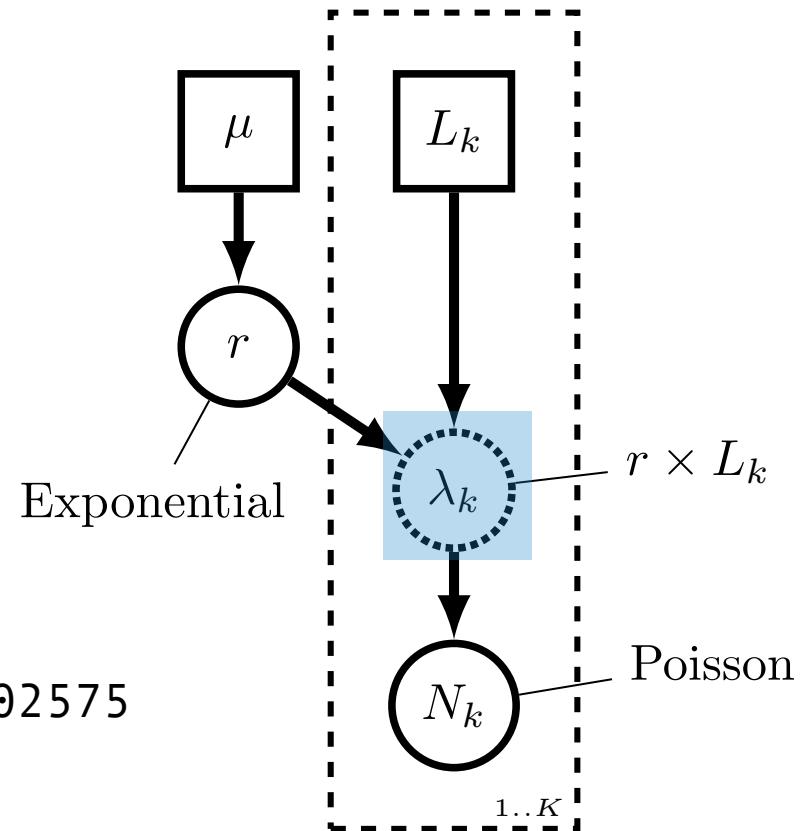
Set initial conditions (for locus 1)

$$\mu = 1/10^{-8}$$

$$L_1 = 1030 \text{ bp}$$

Sample $r = 2.5 \times 10^{-7}$ from $\text{Exp}(1/10^{-8})$

$$\text{Determine } \lambda_1 = 1030 * 2.5 \times 10^{-7} = 0.0002575$$



Simulation using DAG structure use *parent values* to draw *child values*

Set initial conditions (for locus 1)

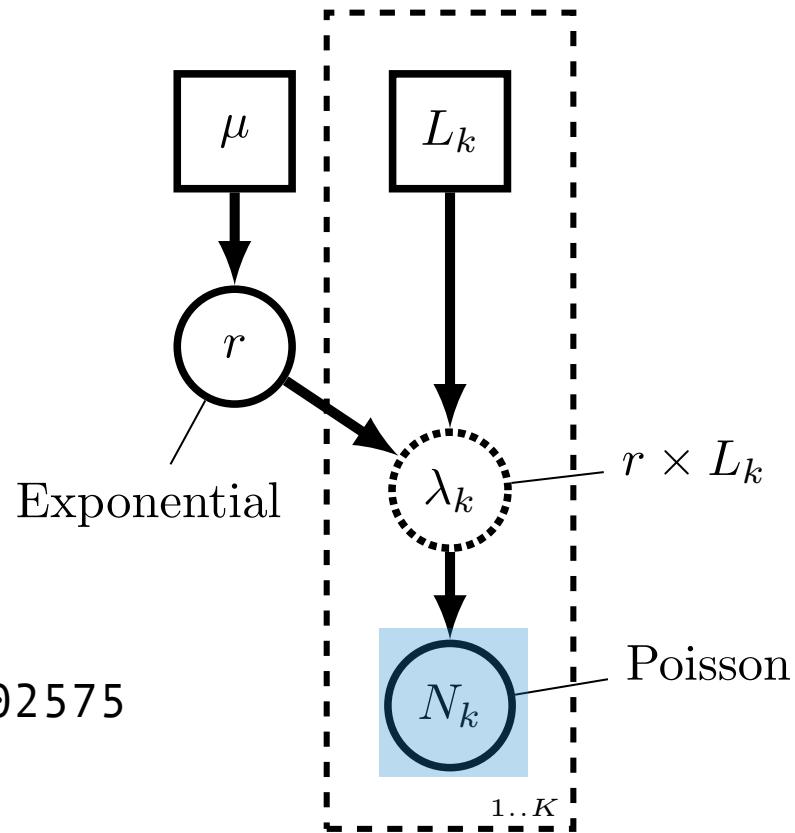
$$\mu = 1/10^{-8}$$

$$L_1 = 1030 \text{ bp}$$

Sample $r = 2.5 * 10^{-7}$ from $\text{Exp}(1/10^{-8})$

$$\text{Determine } \lambda_1 = 1030 * 2.5 * 10^{-7} = 0.0002575$$

Sample $N_1 = 0$ from $\text{Poisson}(0.0002575)$



Simulation using DAG structure use *parent values* to draw *child values*

Set initial conditions (for locus 2)

$$\mu = 1/10^{-8}$$

$$L_1 = 1030 \text{ bp}$$

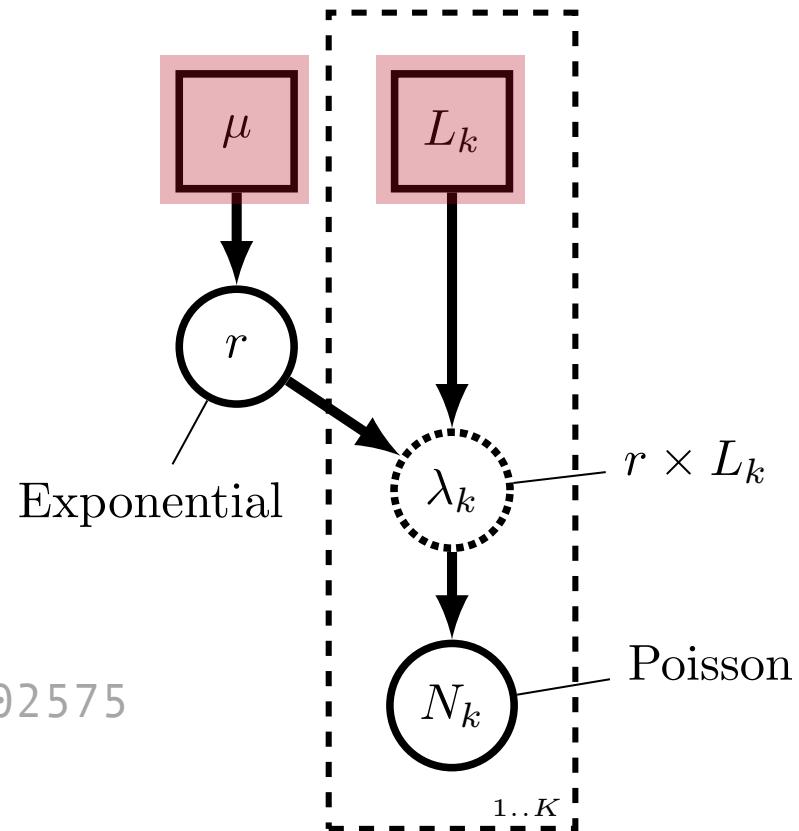
$$\mu = 1/10^{-8}$$

$$L_2 = 3.2 * 10^9 \text{ bp}$$

Sample $r = 2.5 * 10^{-7}$ from $\text{Exp}(1/10^{-8})$

Determine $\lambda_1 = 1030 * 2.5 * 10^{-7} = 0.0002575$

Sample $N_1 = 0$ from $\text{Poisson}(0.0002575)$



Simulation using DAG structure use *parent values* to draw *child values*

Set initial conditions (for locus 2)

$$\mu = 1/10^{-8}$$

$$L_1 = 1030 \text{ bp}$$

$$\mu = 1/10^{-8}$$

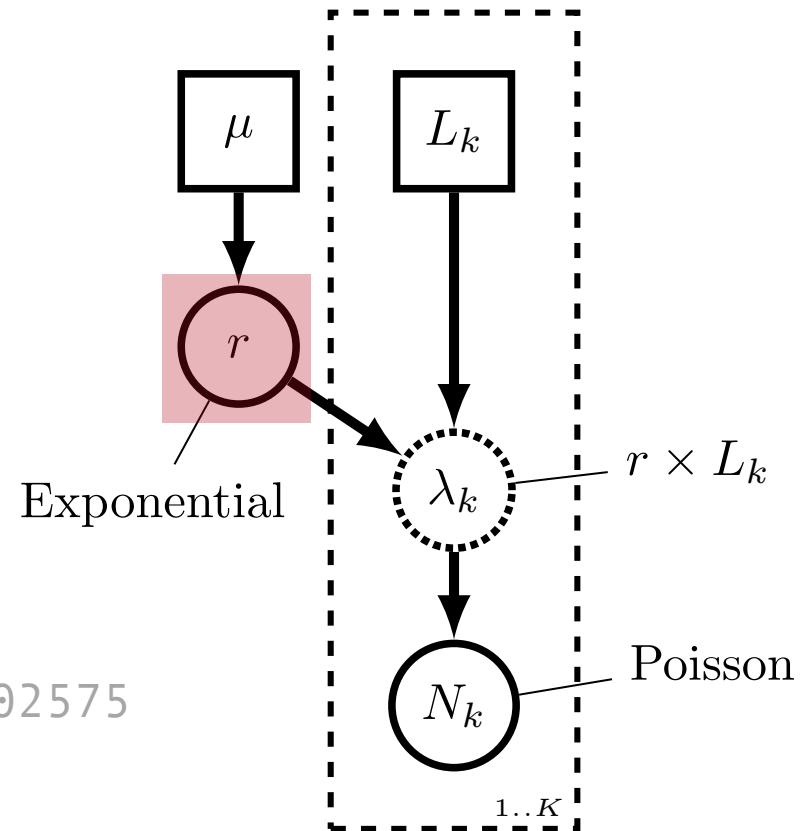
$$L_2 = 3.2 * 10^9 \text{ bp}$$

Sample $r = 2.5 * 10^{-7}$ from $\text{Exp}(1/10^{-8})$

Sample $r = 1.2 * 10^{-8}$ from $\text{Exp}(1/10^{-8})$

Determine $\lambda_1 = 1030 * 2.5 * 10^{-7} = 0.0002575$

Sample $N_1 = 0$ from $\text{Poisson}(0.0002575)$



Simulation using DAG structure

use parent values to draw child values

Set initial conditions (for locus 2)

$$\mu = 1/10^{-8}$$

$$L_1 = 1030 \text{ bp}$$

$$\mu = 1/10^{-8}$$

$$L_2 = 3.2 * 10^9 \text{ bp}$$

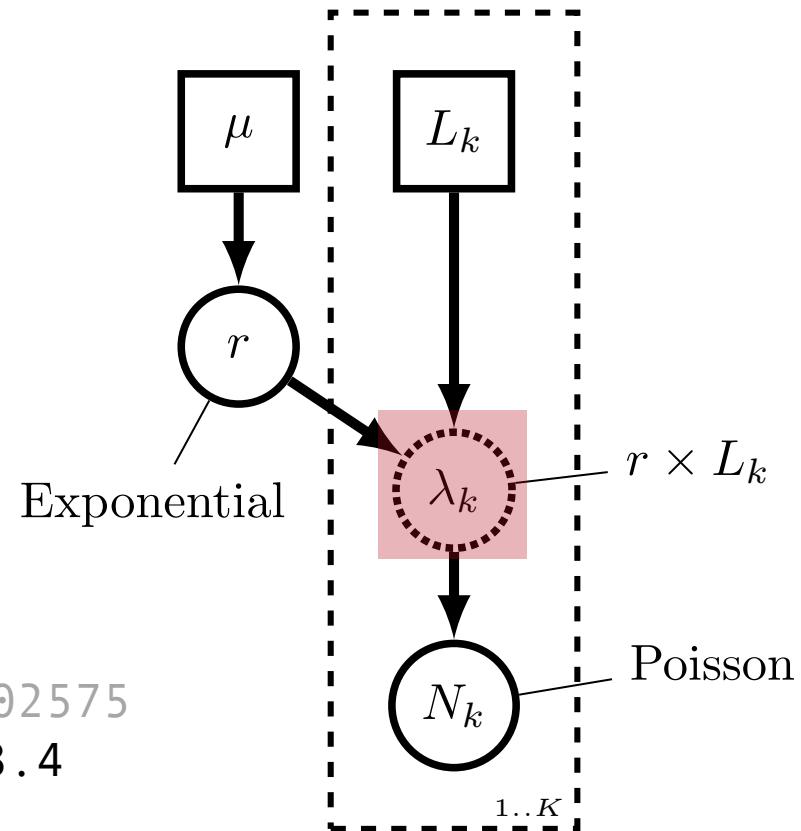
Sample $r = 2.5 * 10^{-7}$ from $\text{Exp}(1/10^{-8})$

Sample $r = 1.2 * 10^{-8}$ from $\text{Exp}(1/10^{-8})$

Determine $\lambda_1 = 1030 * 2.5 * 10^{-7} = 0.0002575$

Determine $\lambda_2 = 3.2 * 10^9 * 1.2 * 10^{-8} = 38.4$

Sample $N_1 = 0$ from $\text{Poisson}(0.0002575)$



Simulation using DAG structure

use parent values to draw child values

Set initial conditions (for locus 2)

$$\mu = 1/10^{-8}$$

$$L_1 = 1030 \text{ bp}$$

$$\mu = 1/10^{-8}$$

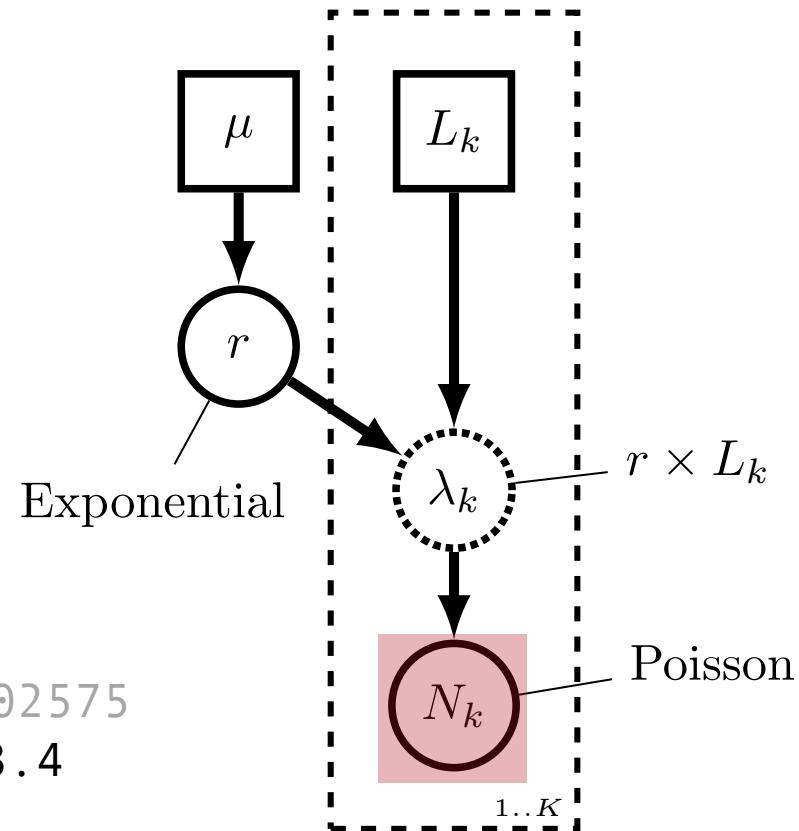
$$L_2 = 3.2 * 10^9 \text{ bp}$$

Sample $r = 2.5 * 10^{-7}$ from $\text{Exp}(1/10^{-8})$

Sample $r = 1.2 * 10^{-8}$ from $\text{Exp}(1/10^{-8})$

Determine $\lambda_1 = 1030 * 2.5 * 10^{-7} = 0.0002575$

Determine $\lambda_2 = 3.2 * 10^9 * 1.2 * 10^{-8} = 38.4$



Sample $N_1 = 0$ from $\text{Poisson}(0.0002575)$

Sample $N_2 = 44$ from $\text{Poisson}(38.4)$

Simulation experiment

Biological data is messy, limited

Simulated data is perfect, abundant

Simulation experiment

Biological data is messy, limited

Simulated data is perfect, abundant

A model that performs poorly with **simulated** data
generally performs worse with **biological** data!

Simulation experiment

Biological data is messy, limited

Simulated data is perfect, abundant

A model that performs poorly with **simulated** data
generally performs worse with **biological** data!

Simulation experiment

1. Analyze simulated data under controlled settings
vary dataset size, parameters, models, etc.
2. Translate simulation results to guide biological analyses

(more later)

Draw a graphical model for

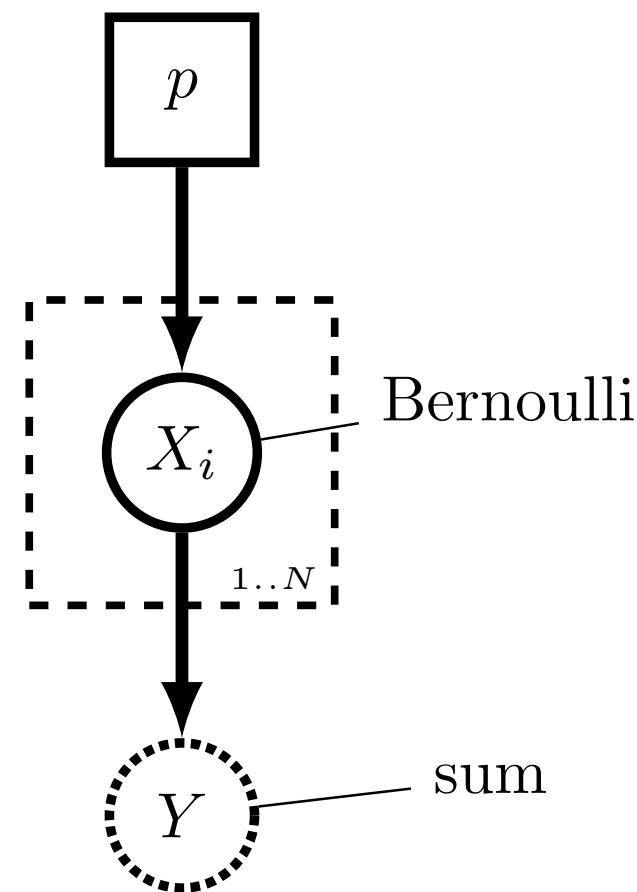
$$p = 0.5$$

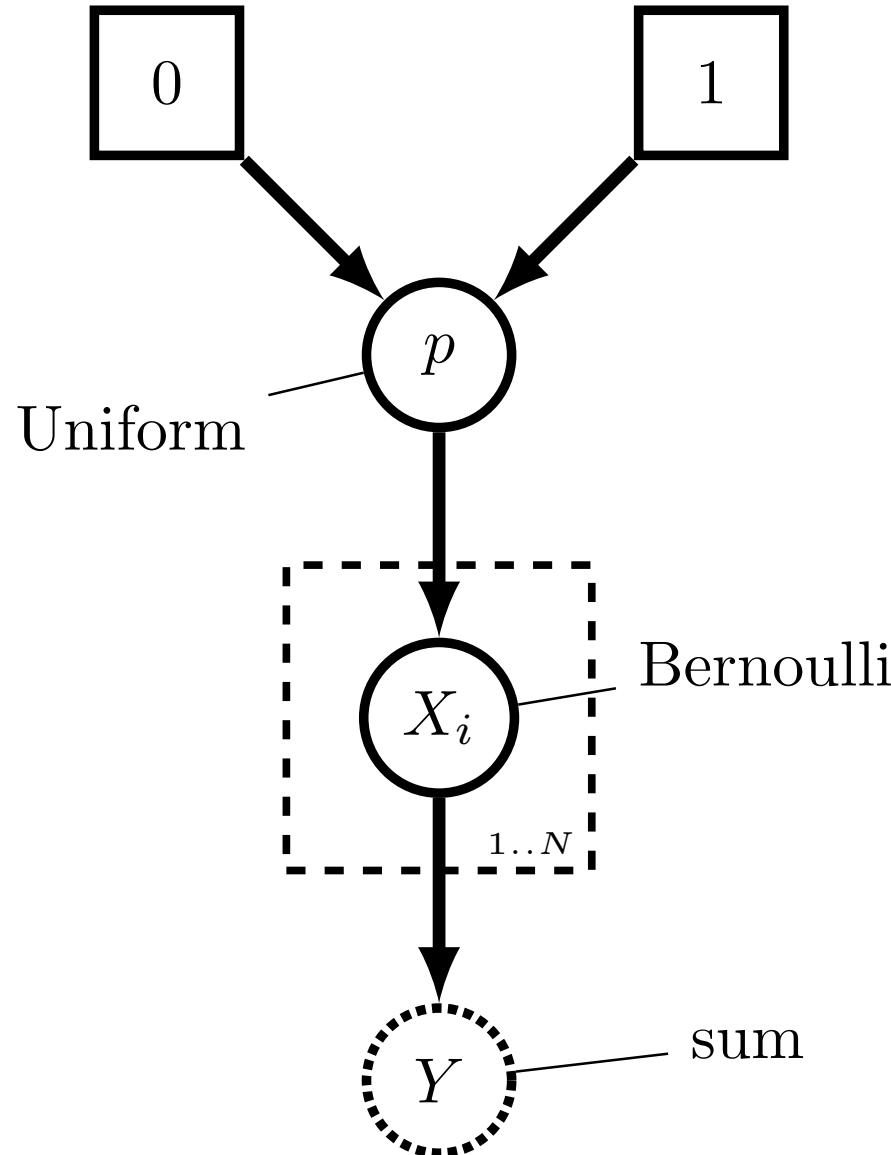
$$X_i \sim \text{Bernoulli}(p)$$

$$Y = X_1 + X_2 + \dots + X_N$$

Extra challenge

Treat p as stochastic node that takes any value from 0 to 1 with equal probability





Draw a graphical model for

$$\lambda = 1$$

$$L = 3$$

$$U = 10$$

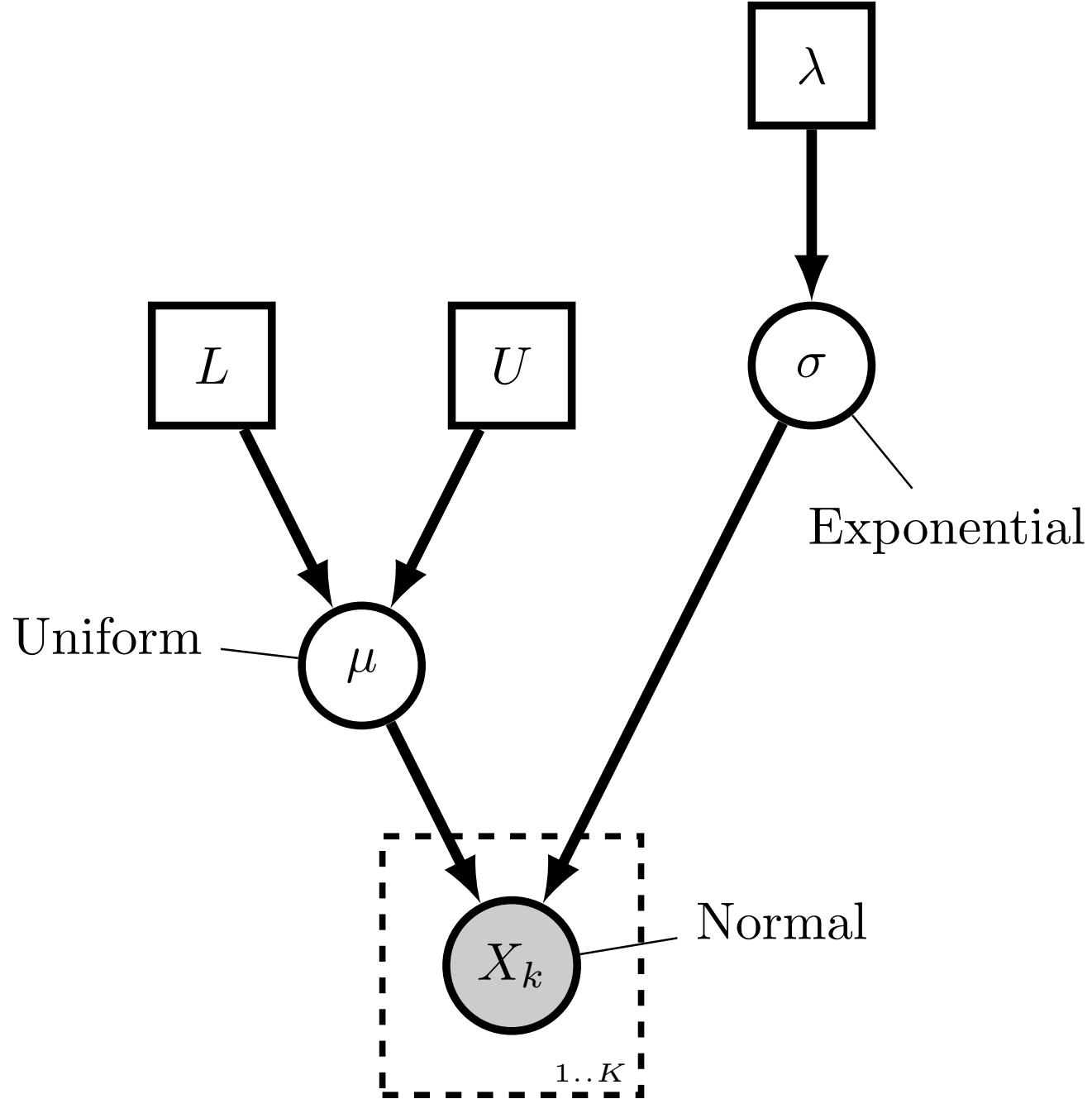
$$\sigma \sim \text{Exponential}(\lambda)$$

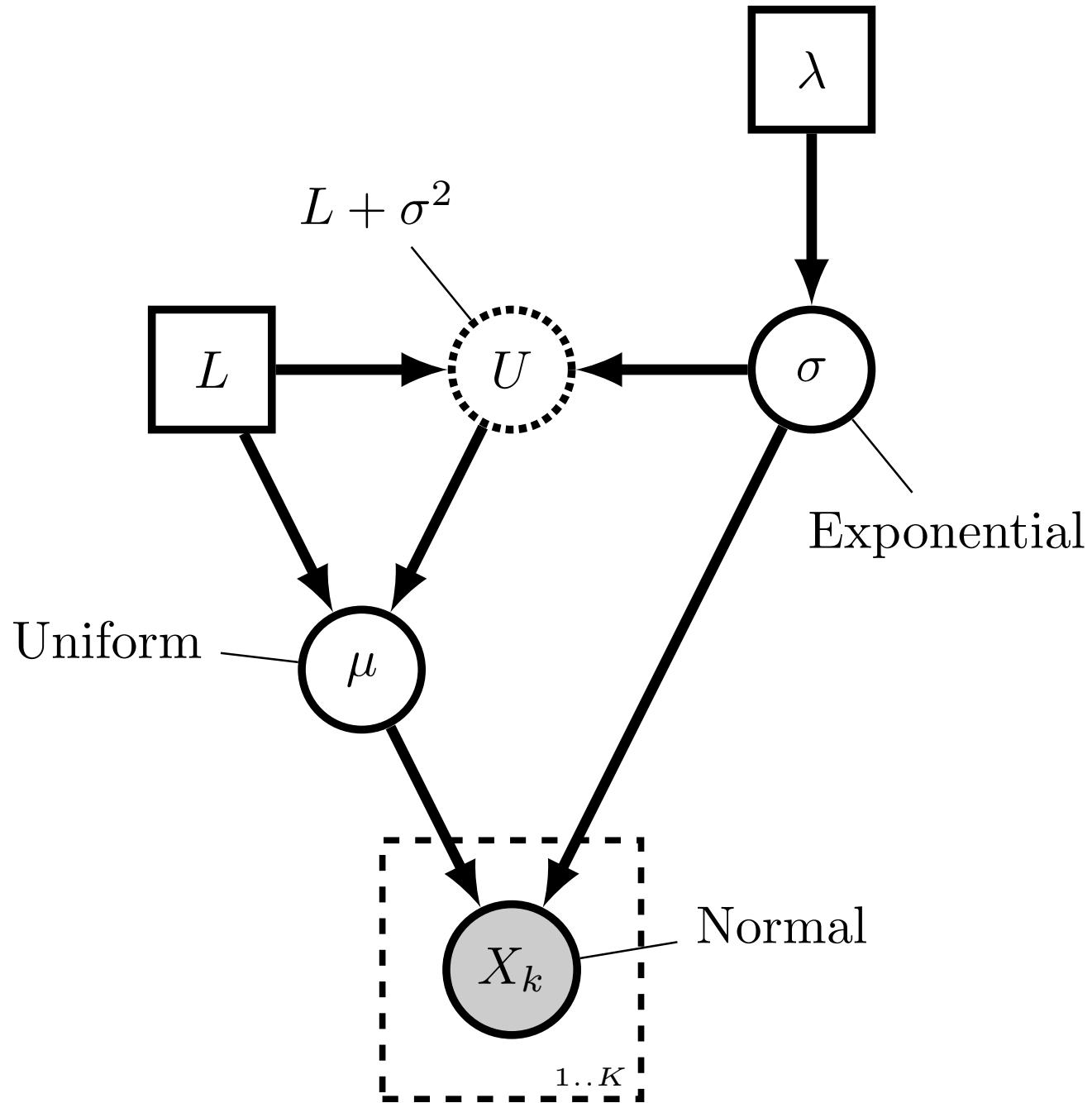
$$\mu \sim \text{Uniform}(L, U)$$

$$X_k \sim \text{Normal}(\mu, \sigma) - \text{observed data}$$

Extra challenge

Modify U so its value determined by $L + \sigma^2$





Model definitions

Conceptual

Probabilistic

Graphical

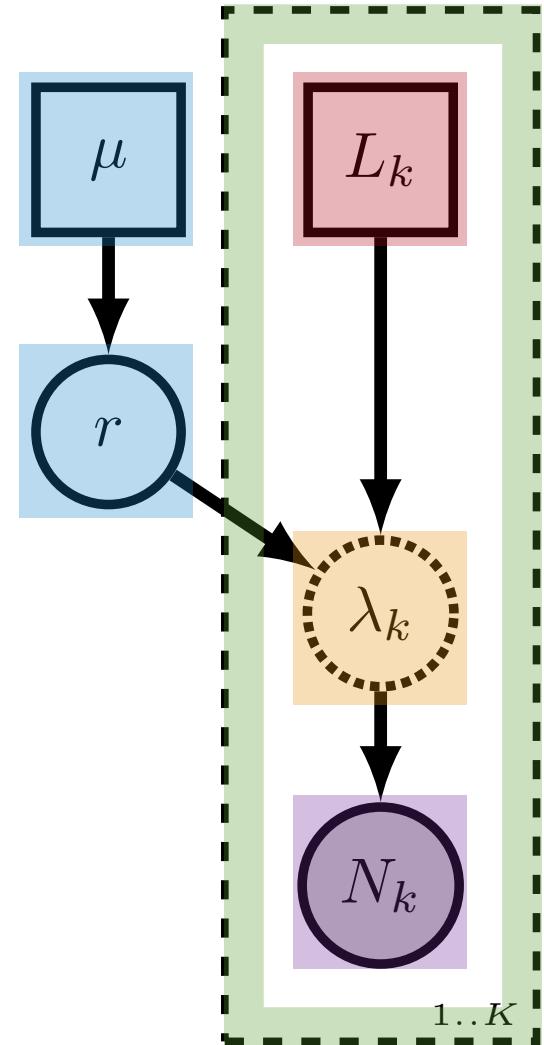
Computational

Computational

```
# input
data = [ 0, 3, 3 ]
L    = [ 310, 780, 1030 ]

# per-site mutation rate
mu = 1e-8
r_site ~ dnExponential(1/mu)

# mutation counts per locus
K = data.size()
for (k in 1:K) {
  r_loci[k] := r_site * L[k]
  n_muts[k] ~ dnPoisson(r_loci[k])
  n_muts[k].clamp( data[k] )
}
```



RevBayes: Bayesian Phylogenetic Inference Using Graphical Models and an Interactive Model-Specification Language

Sebastian Höhna , Michael J. Landis, Tracy A. Heath, Bastien Boussau, Nicolas Lartillot, Brian R. Moore, John P. Huelsenbeck, Fredrik Ronquist 

Systematic Biology, Volume 65, Issue 4, 1 July 2016, Pages 726–736,

<https://doi.org/10.1093/sysbio/syw021>

Published: 28 May 2016 [Article history ▾](#)

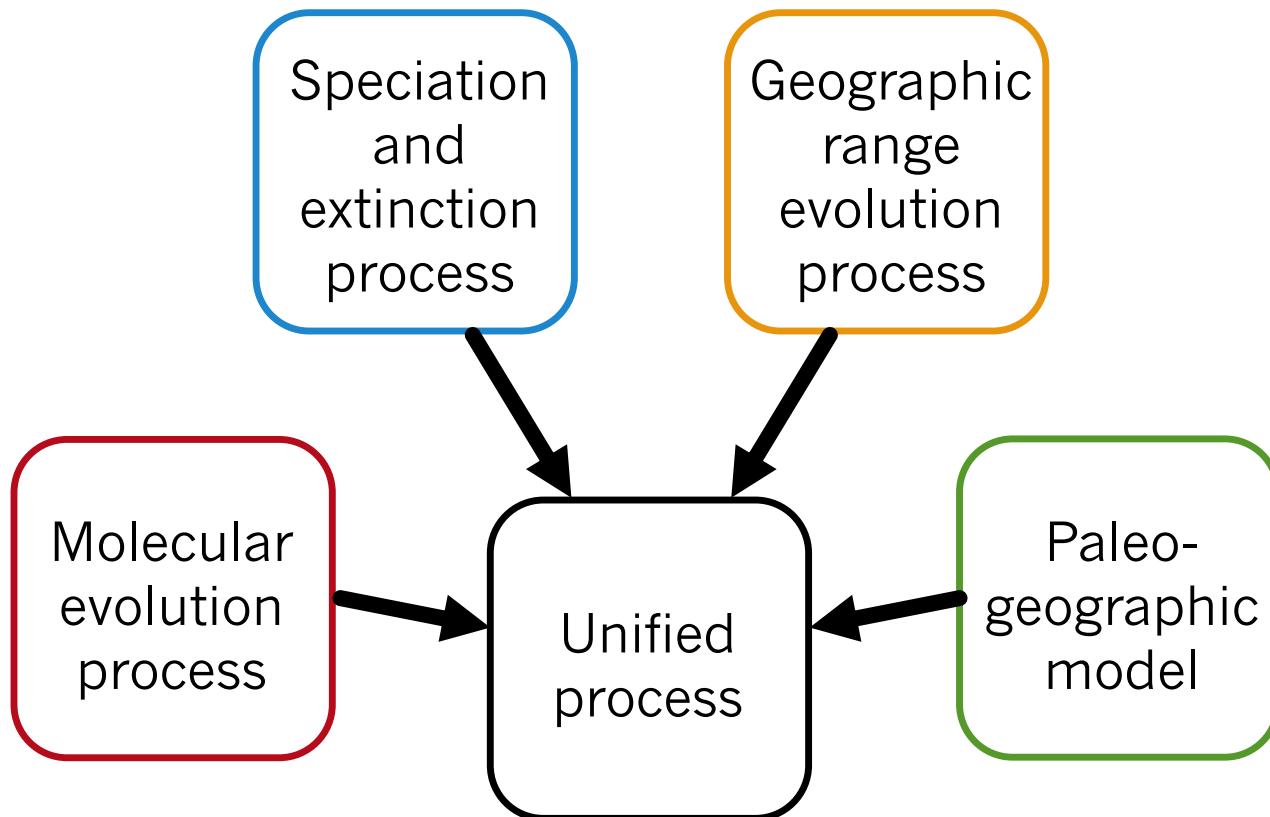
RevBayes

1. Fast
2. Flexible
3. Easy-to-use

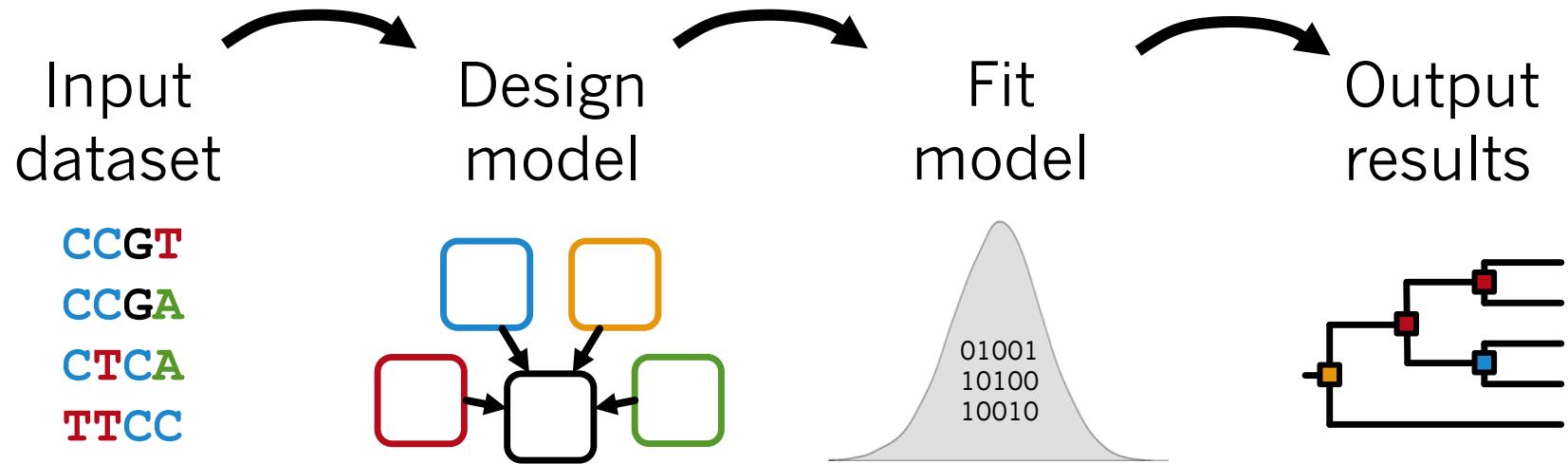


<http://revbayes.com>

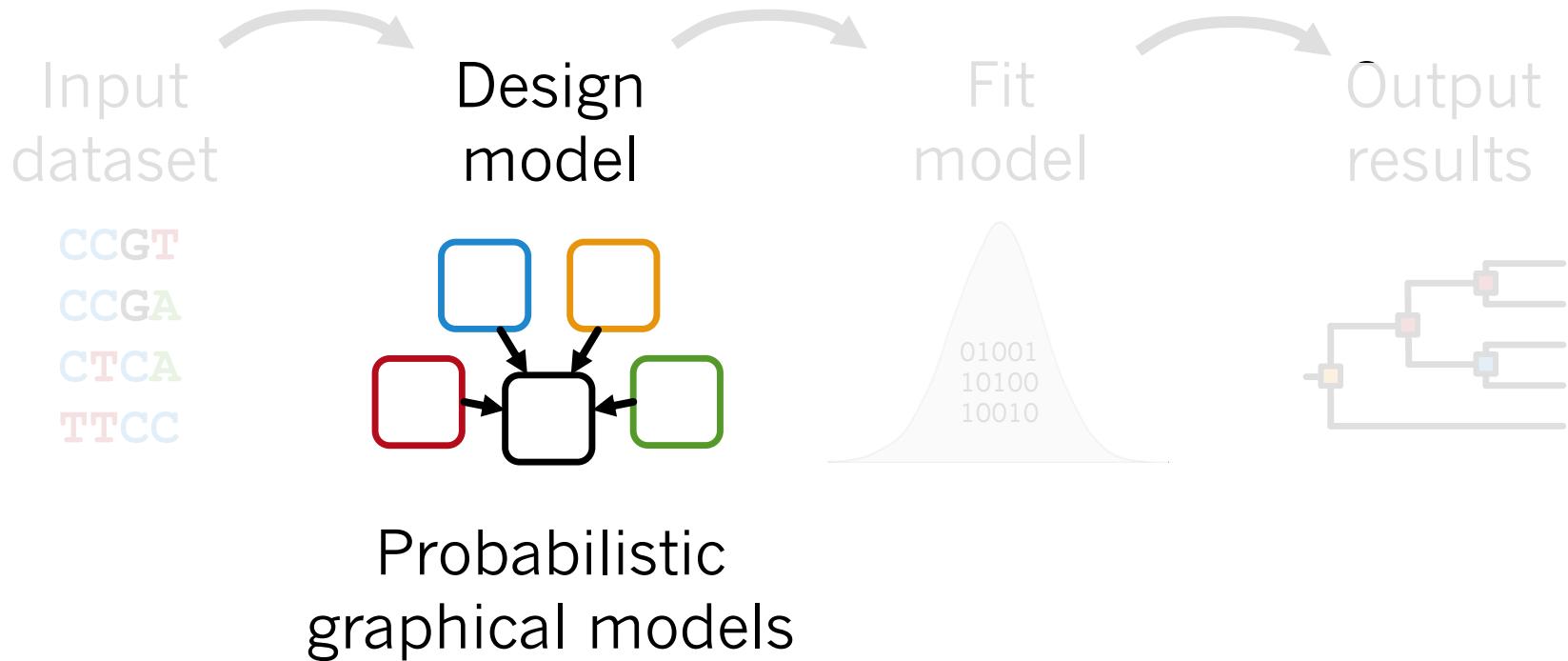
Model design with RevBayes



Phylogenetic inference with RevBayes



Phylogenetic inference with RevBayes



Familiar phylogenetic components

Taxon data (what evolved)

Tree model

Rate matrix

Clock model

Site-rate variation

Fossilization model

Biogeographic model

Continuous trait model

...

Familiar phylogenetic components

Taxon data (what evolved)

Tree model

Rate matrix

Clock model

Site-rate variation

Fossilization model

Biogeographic model

Continuous trait model

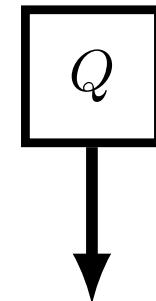
...

$$Q_{\text{JC}} = \begin{pmatrix} - & 1/3 & 1/3 & 1/3 \\ 1/3 & - & 1/3 & 1/3 \\ 1/3 & 1/3 & - & 1/3 \\ 1/3 & 1/3 & 1/3 & - \end{pmatrix}$$

Substitution model

Jukes-Cantor / JC

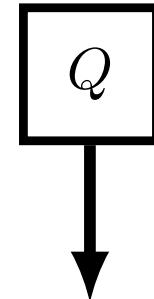
$$Q_{\text{JC}} = \begin{pmatrix} - & 1/3 & 1/3 & 1/3 \\ 1/3 & - & 1/3 & 1/3 \\ 1/3 & 1/3 & - & 1/3 \\ 1/3 & 1/3 & 1/3 & - \end{pmatrix}$$



Substitution model

Jukes-Cantor / JC

$$Q_{\text{JC}} = \begin{pmatrix} - & 1/3 & 1/3 & 1/3 \\ 1/3 & - & 1/3 & 1/3 \\ 1/3 & 1/3 & - & 1/3 \\ 1/3 & 1/3 & 1/3 & - \end{pmatrix}$$



```
# Jukes-Cantor model  
Q <- fnJC(4)
```

```
x <- 1  
y ~ dnExp(x)  
z := sqrt(y)
```

Variable
(child)

```
x <- 1
y ~ dnExp(x)
z := sqrt(y)
```



Variable
(child)

```
x <- 1
y ~ dnExp(x)
z := sqrt(y)
```



Assignment
operator
(node type)

Variable
(child)

```
x <- 1
y ~ dnExp(x)
z := sqrt(y)
```

Assignment
operator
(node type)

Value, distribution,
or function
(node behavior)

Variable
(child)

Argument
(parent)

```
x <- 1
y ~ dnExp(x)
z := sqrt(y)
```

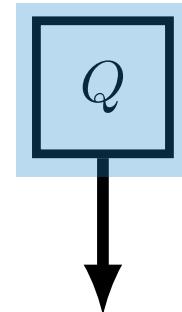
Assignment
operator
(node type)

Value, distribution,
or function
(node behavior)

Substitution model

Jukes-Cantor / JC

$$Q_{\text{JC}} = \begin{pmatrix} - & 1/3 & 1/3 & 1/3 \\ 1/3 & - & 1/3 & 1/3 \\ 1/3 & 1/3 & - & 1/3 \\ 1/3 & 1/3 & 1/3 & - \end{pmatrix}$$



```
# Jukes-Cantor model  
Q <- fnJC(4)
```

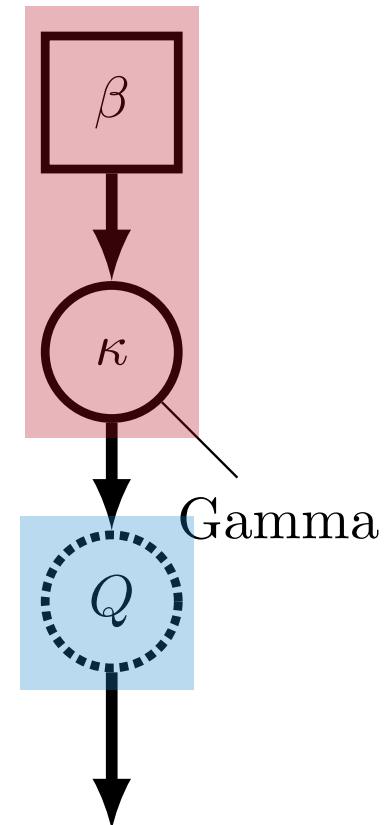
Substitution model

Kimura-2P / K2P

$$Q_{K80} = \begin{pmatrix} - & 1 & \kappa & 1 \\ 1 & - & 1 & \kappa \\ \kappa & 1 & - & 1 \\ 1 & \kappa & 1 & - \end{pmatrix}$$

```
# Ti/Tv rate ratio, mean 1
beta <- 1
kappa ~ dnGamma(beta, beta)

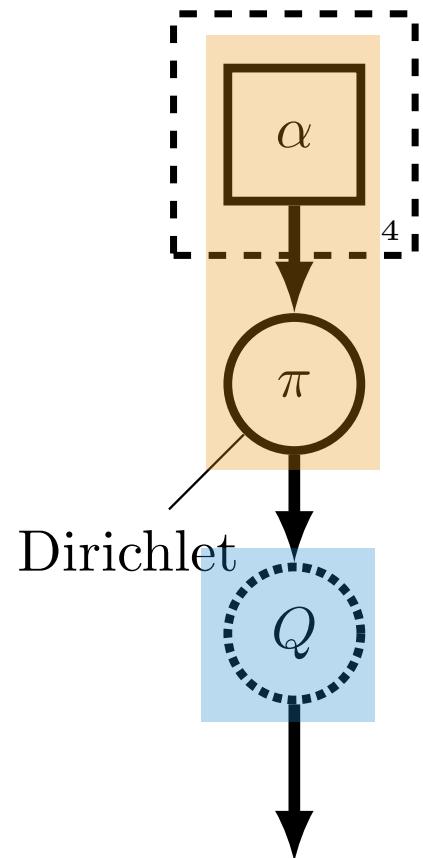
# Kimura 2-parameter model
Q := fnK80(kappa)
```



$$Q_{\text{F81}} = \begin{pmatrix} - & \pi_C & \pi_G & \pi_T \\ \pi_A & - & \pi_G & \pi_T \\ \pi_A & \pi_C & - & \pi_T \\ \pi_A & \pi_C & \pi_G & - \end{pmatrix}$$

```
# base frequencies, "flat"
alpha <- [1, 1, 1, 1]
pi ~ dnDirichlet(alpha)

# Felsenstein-81 model
Q := fnF81(pi)
```



Substitution model

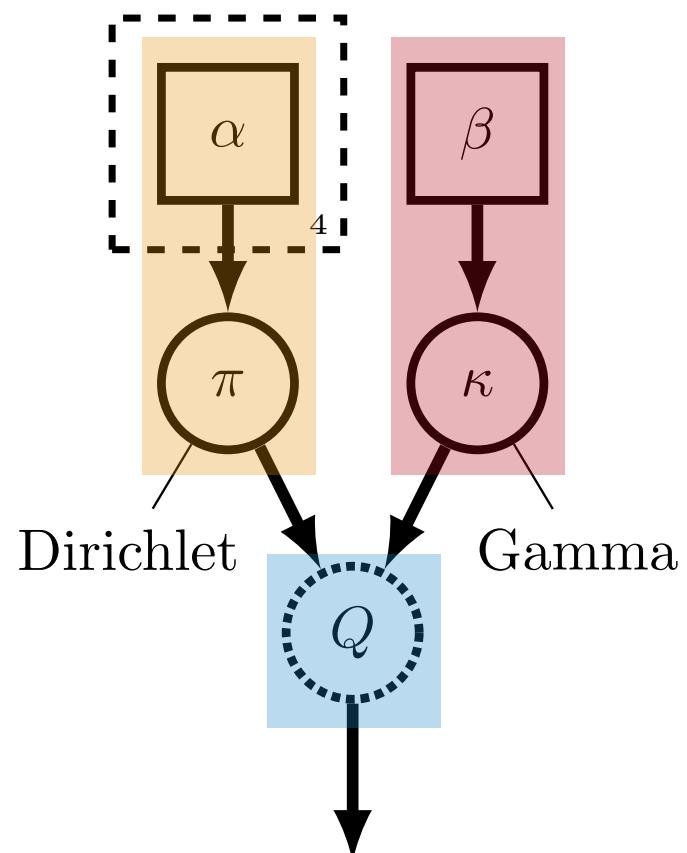
Hasegawa-Kishino-Yano / HKY

$$Q_{\text{HKY}} = \begin{pmatrix} - & \pi_C & \kappa \pi_G & \pi_T \\ \pi_A & - & \pi_G & \kappa \pi_T \\ \kappa \pi_A & \pi_C & - & \pi_T \\ \pi_A & \kappa \pi_C & \pi_G & - \end{pmatrix}$$

```
# base frequencies, "flat"
alpha <- [1, 1, 1, 1]
pi ~ dnDirichlet(alpha)

# Ti/Tv rate ratio, mean 1
beta <- 1
kappa ~ dnGamma(beta, beta)

# HKY model
Q := fnHKY(pi, kappa)
```



Substitution model

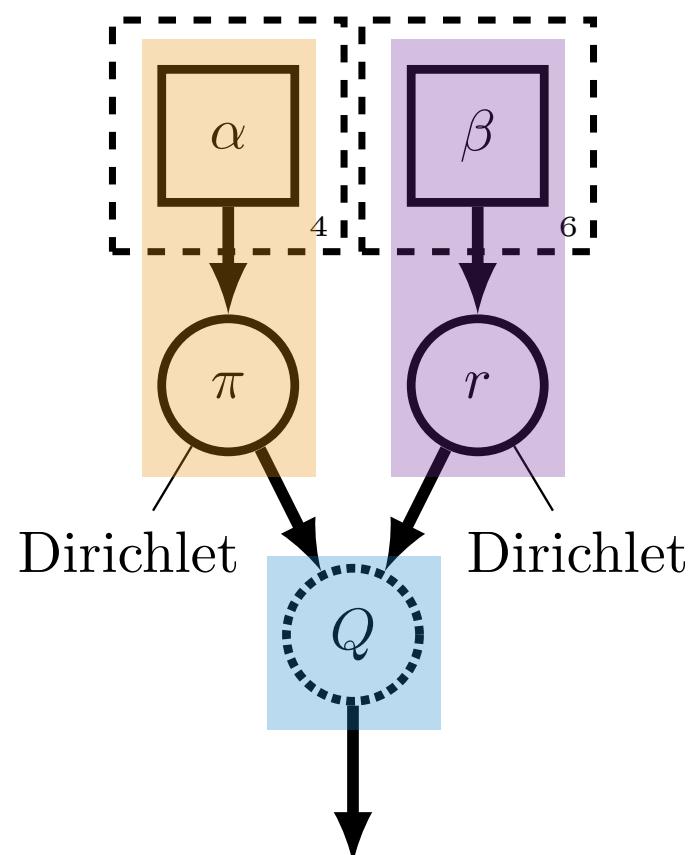
Generalized-Time-Reversible / GTR

$$Q_{\text{GTR}} = \begin{pmatrix} - & r_{AC} \pi_C & r_{AG} \pi_G & r_{AT} \pi_T \\ r_{AC} \pi_A & - & r_{CG} \pi_G & r_{CT} \pi_T \\ r_{AG} \pi_A & r_{CG} \pi_C & - & r_{GT} \pi_T \\ r_{AT} \pi_A & r_{CT} \pi_C & r_{GT} \pi_G & - \end{pmatrix}$$

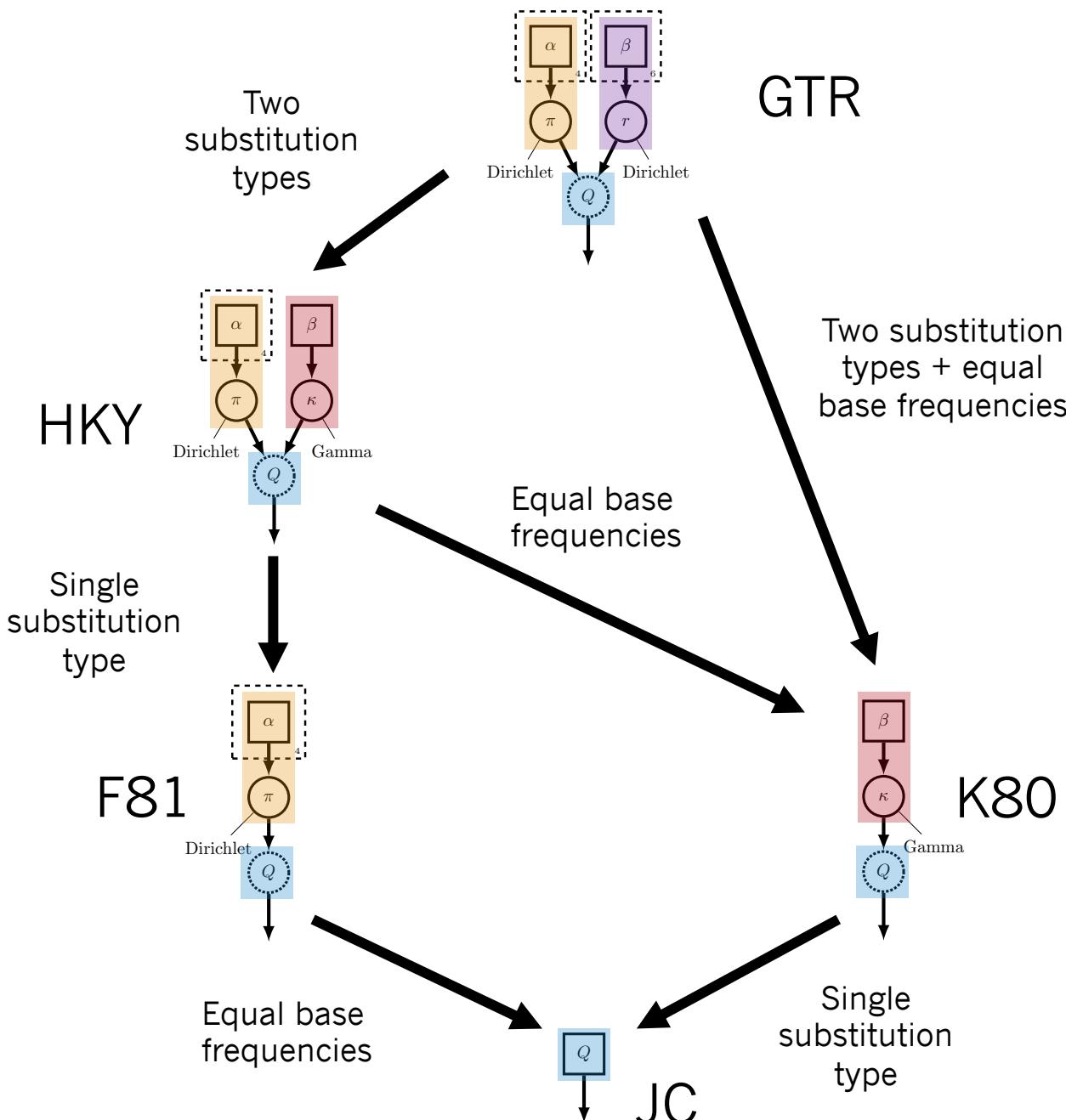
```
# base frequencies, "flat"
alpha <- [1, 1, 1, 1]
pi ~ dnDirichlet(alpha)

# exch. rates, "flat"
beta <- [1, 1, 1, 1, 1, 1]
er ~ dnDirichlet(beta)

# GTR model
Q := fnGTR(pi, er)
```



Complex ↑
Simple ↓



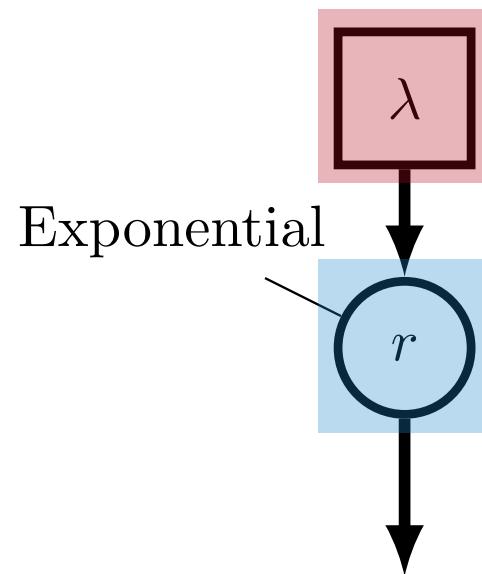
from Dave Swofford's talk

Clock models

Strict clock

```
# strict clock rate mean (^-1)
lambda <- 1.0 / 0.1

# strict exponential clock
clock ~ dnExp(lambda)
```



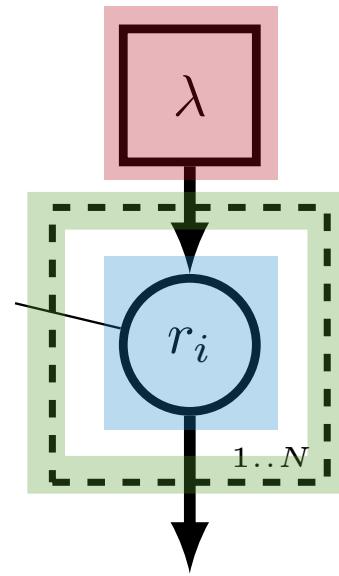
Clock models

Centered relaxed clock iid exponentials

```
# relaxed clock rate mean (^-1)
lambda <- 1.0 / 0.1

# uncorrelated exponential clock
for (i in 1:n_branches) {
  clock[i] ~ dnExp(lambda)
}
```

Exponential



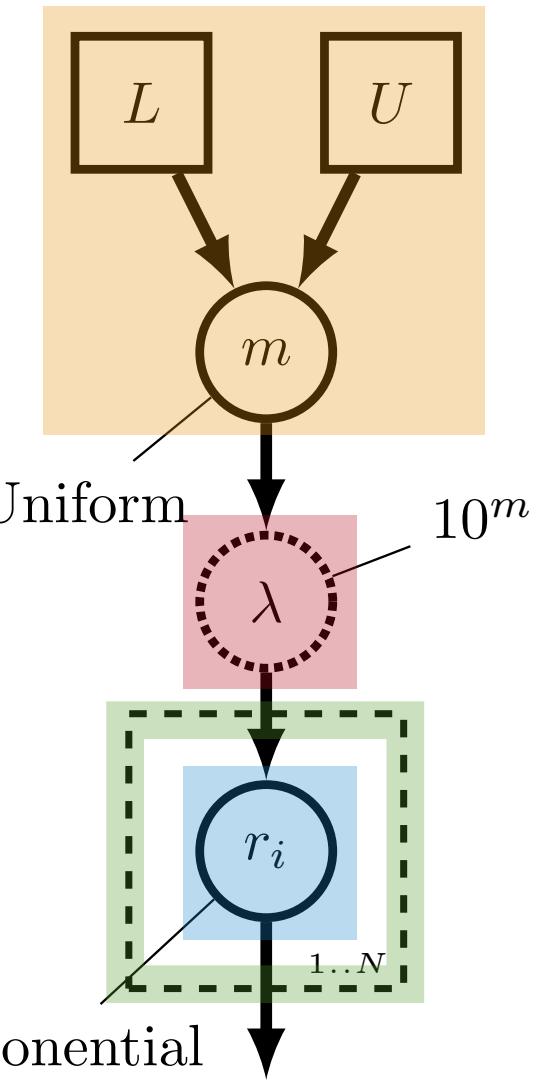
Clock models

```
# order-of-magnitude of clock
lower <- -6
upper <- +1
m ~ dnUniform(lower, upper)

# relaxed clock rate mean (^-1)
lambda := 10^‐m

# uncorrelated exponential clock
for (i in 1:n_branches) {
  clock[i] ~ dnExp(lambda)
}
```

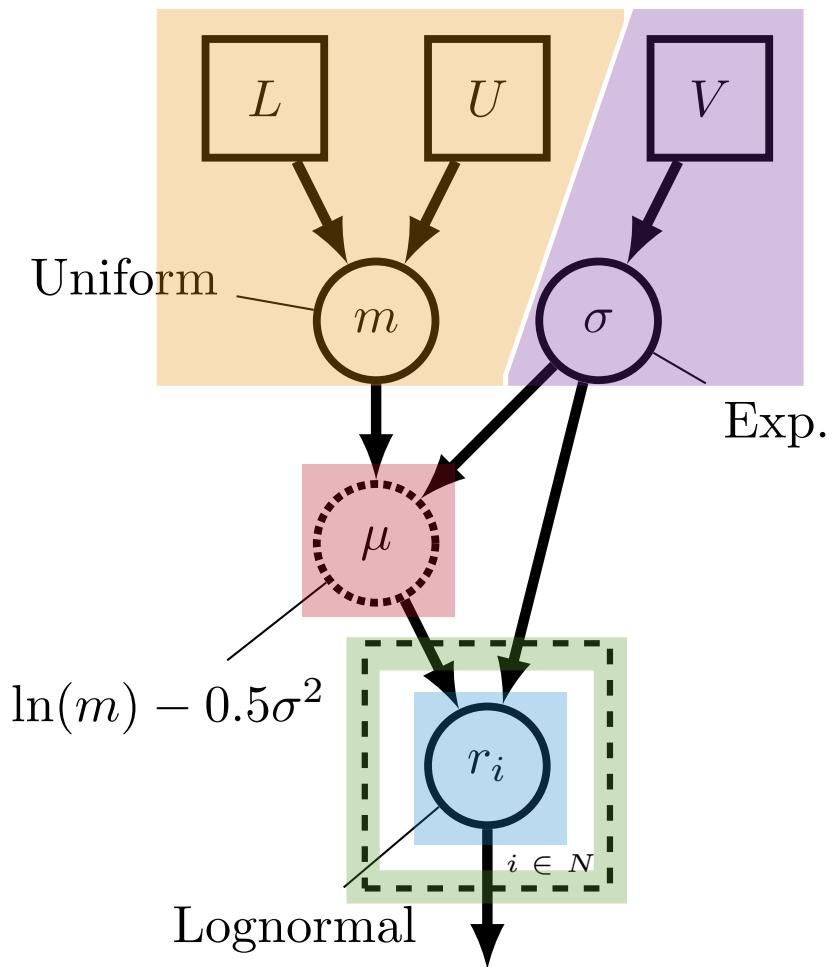
Uncentered relaxed clock iid exponentials



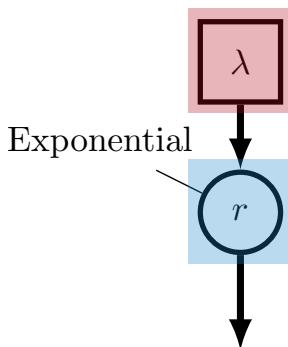
Clock models

```
# order-of-magnitude of clock  
lower <- -6  
upper <- +1  
 $m \sim dnUniform(lower, upper)$   
  
# relaxed clock rate variance  
 $V <- 1$   
 $\sigma \sim dnExp(V)$   
  
# relaxed clock rate mean  
 $\mu := m - 0.5 * \sigma^2$   
  
# uncorrelated lognormal clock  
for (i in 1:n_branches) {  
     $r[i] \sim dnLognormal(\mu, \sigma)$   
}
```

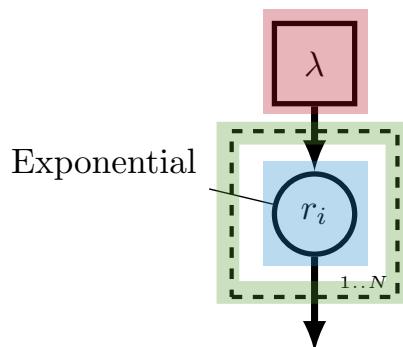
Uncentered relaxed clock iid lognormals



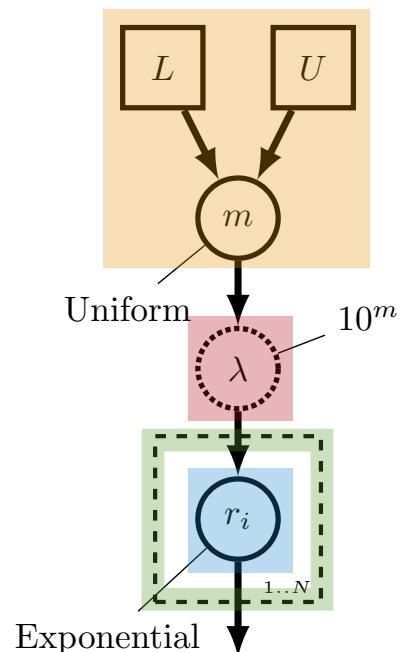
strict
exponential
clock



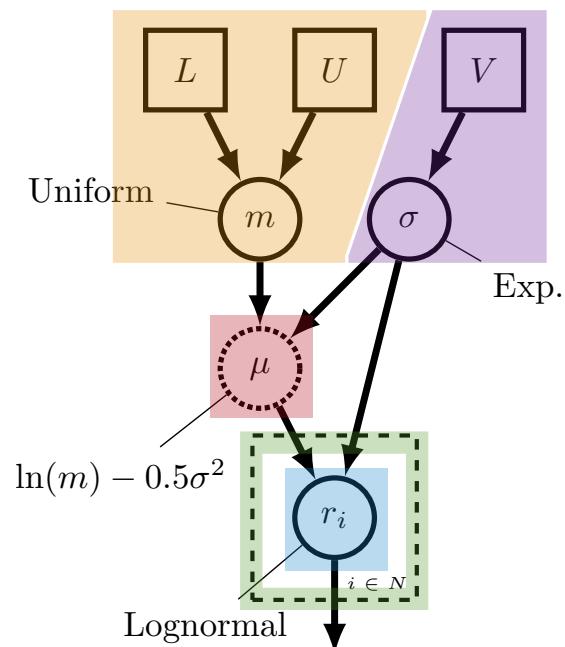
uncorrelated
exponential
clock



uncorrelated
uncentered
exponential
clock



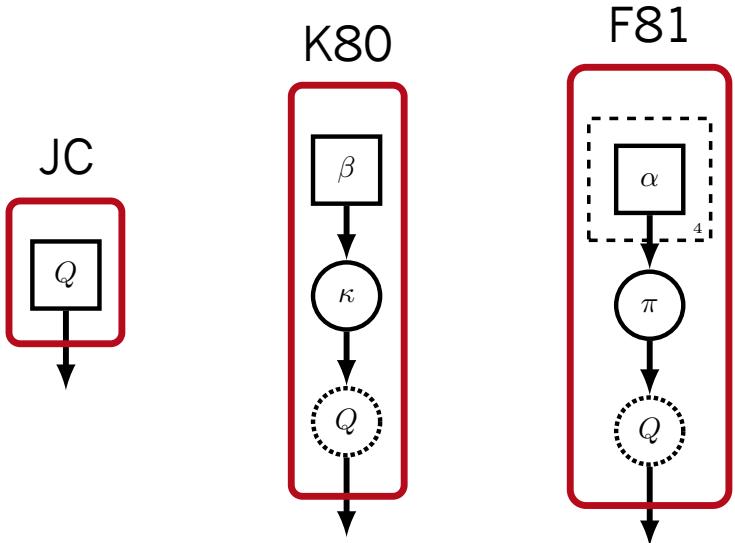
uncorrelated
uncentered
lognormal
clock



Simple

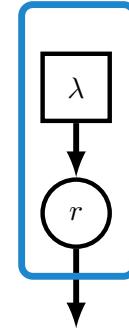
Complex

Rate matrices across sites



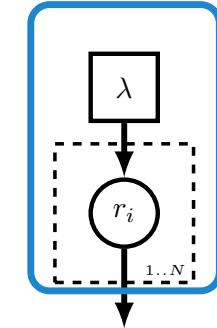
Clock rates along branches

strict
exponential

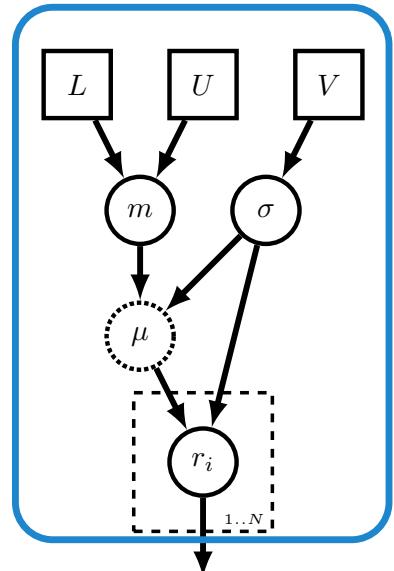
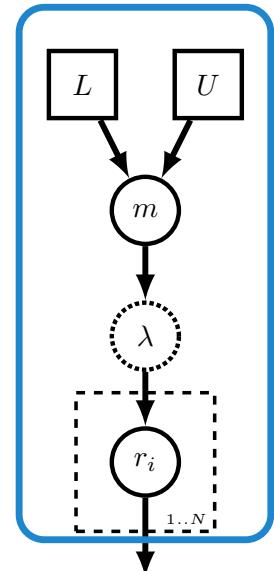


uncorrelated
uncentered
exponential

uncorrelated
exponential



uncorrelated
uncentered
lognormal



PGM modules

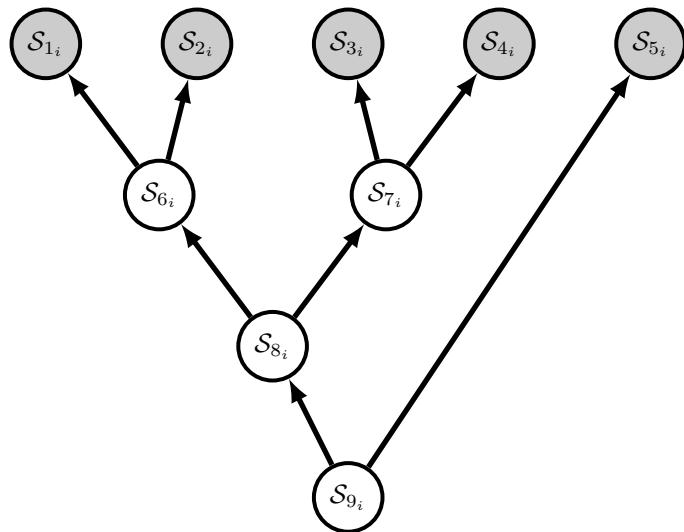
Definition

a module is a PGM that subgraph that performs a distinct, integrated function

Modules

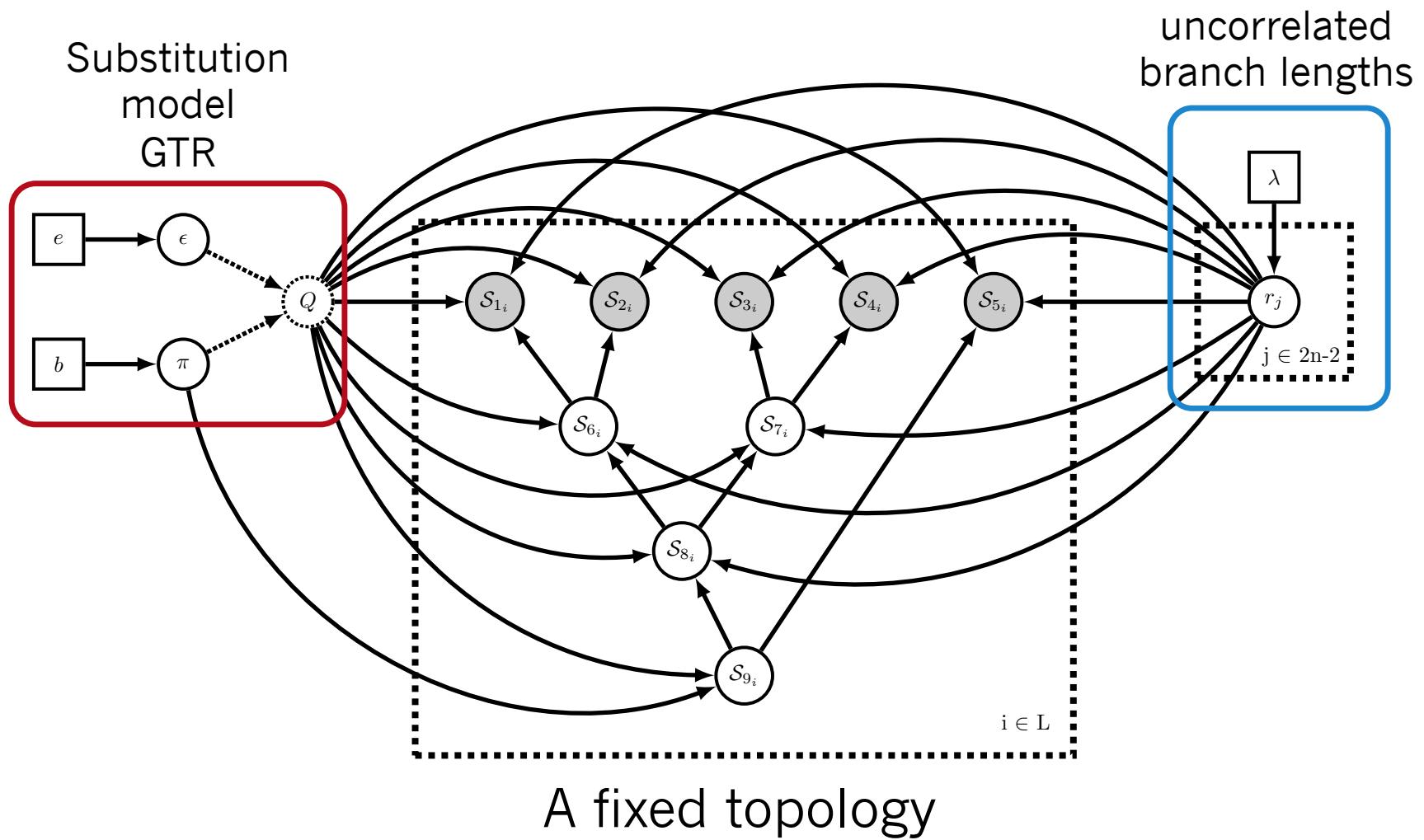
- Reveal “higher level” relationships
- Provide templated functionality
- Mask details unless needed

Phylogenetic substitution model

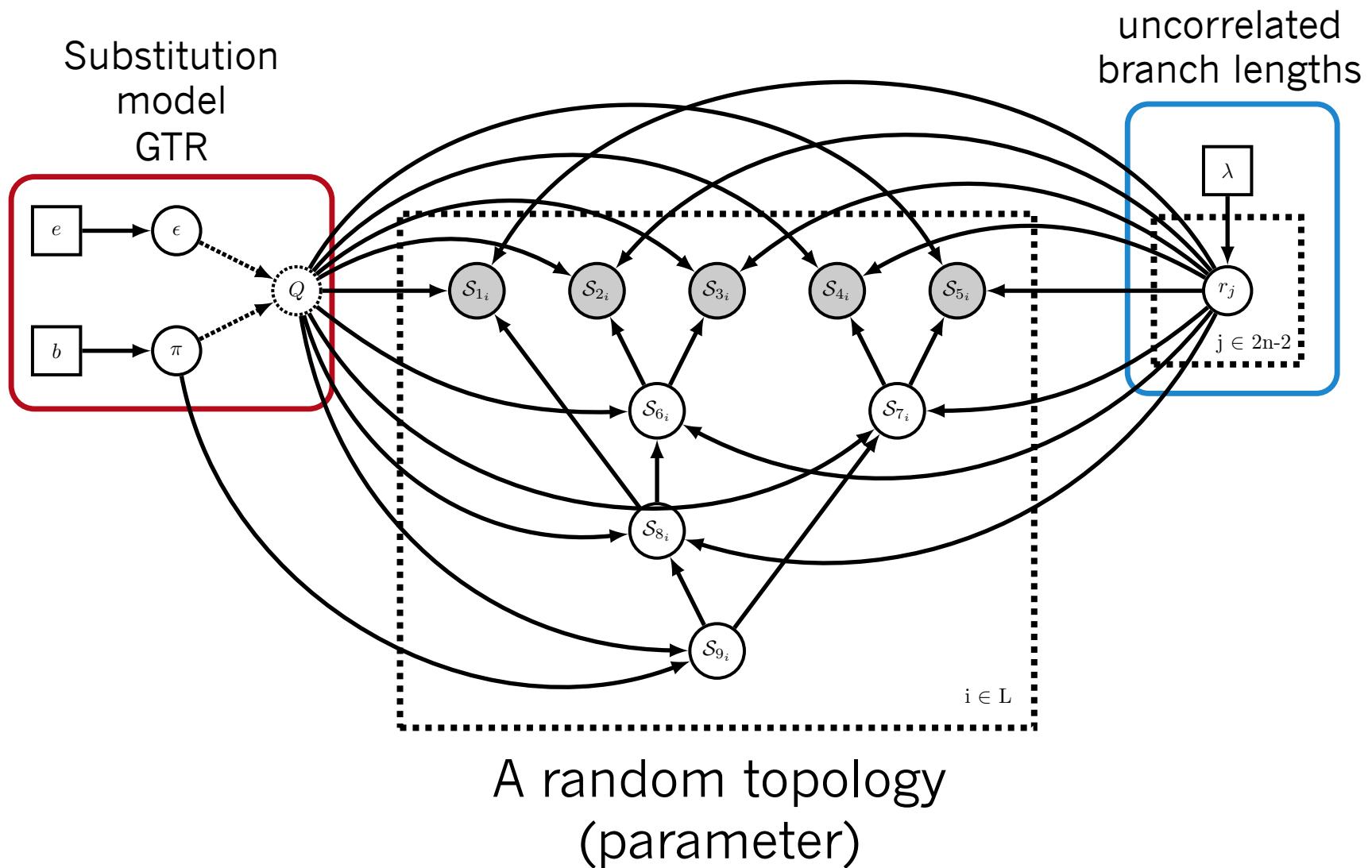


A fixed topology

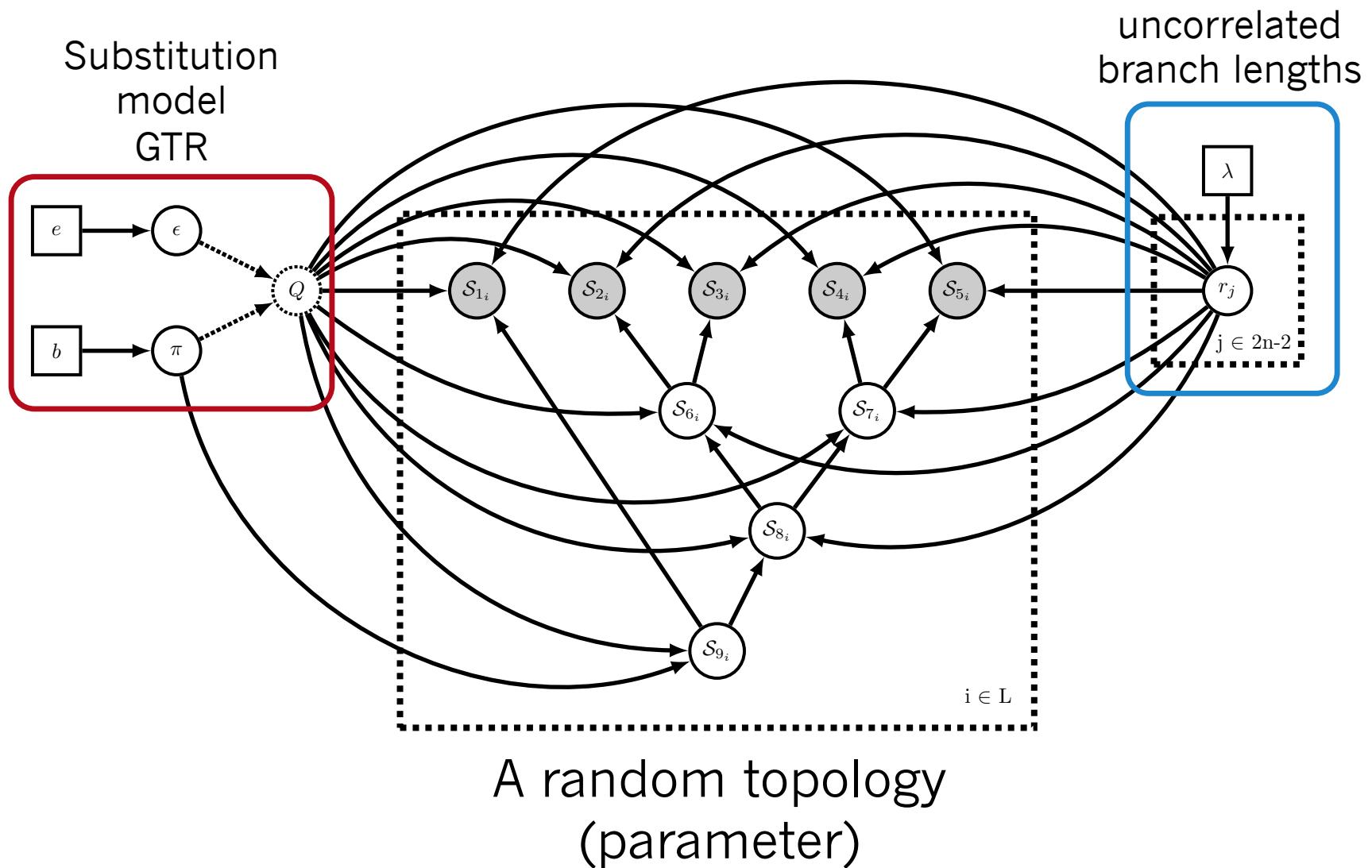
Phylogenetic substitution model



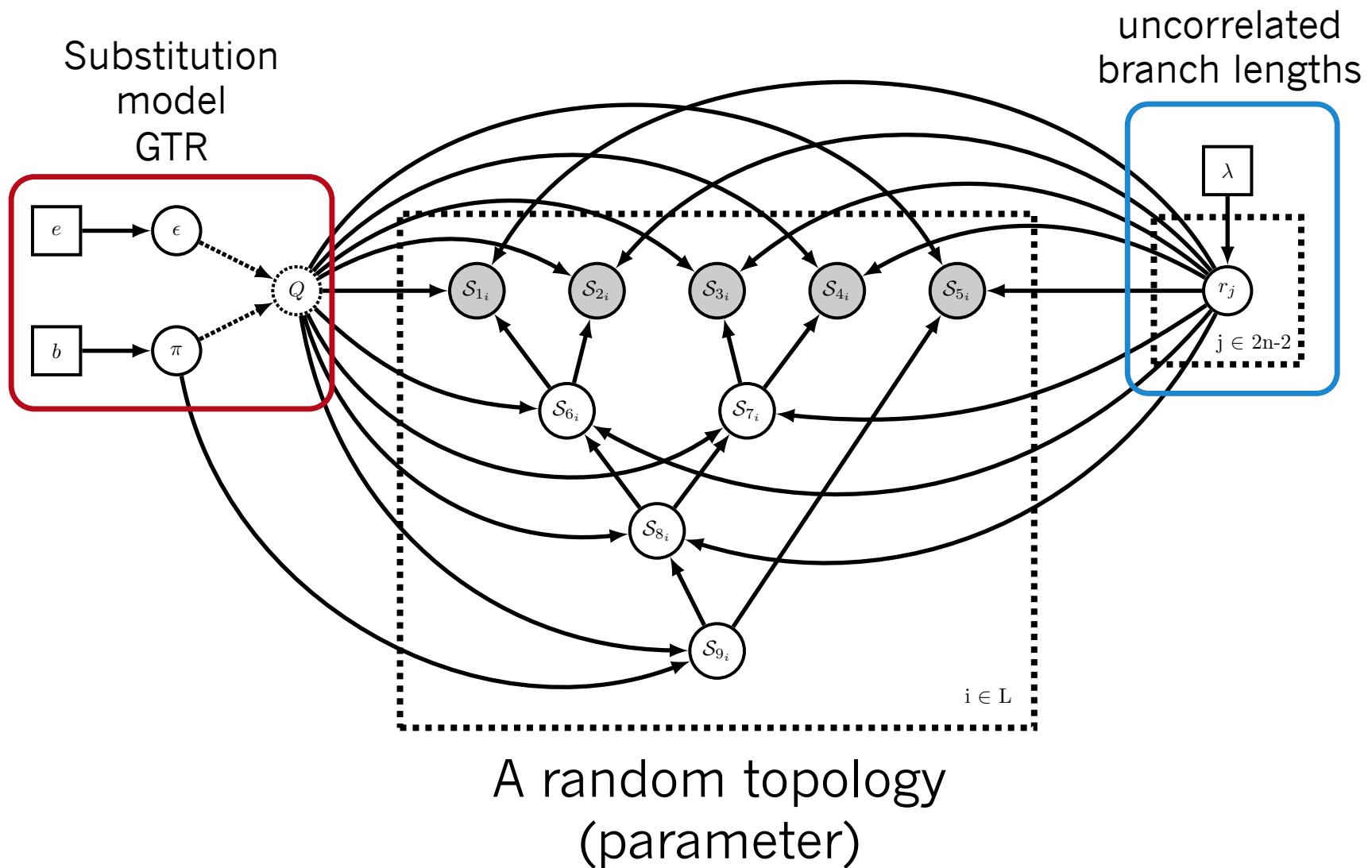
Phylogenetic substitution model



Phylogenetic substitution model



Phylogenetic substitution model



Is a tree a parameter?

Is a tree a graph?

Is a tree a model?

The Holy Treenity

Is a tree a parameter?

Yes

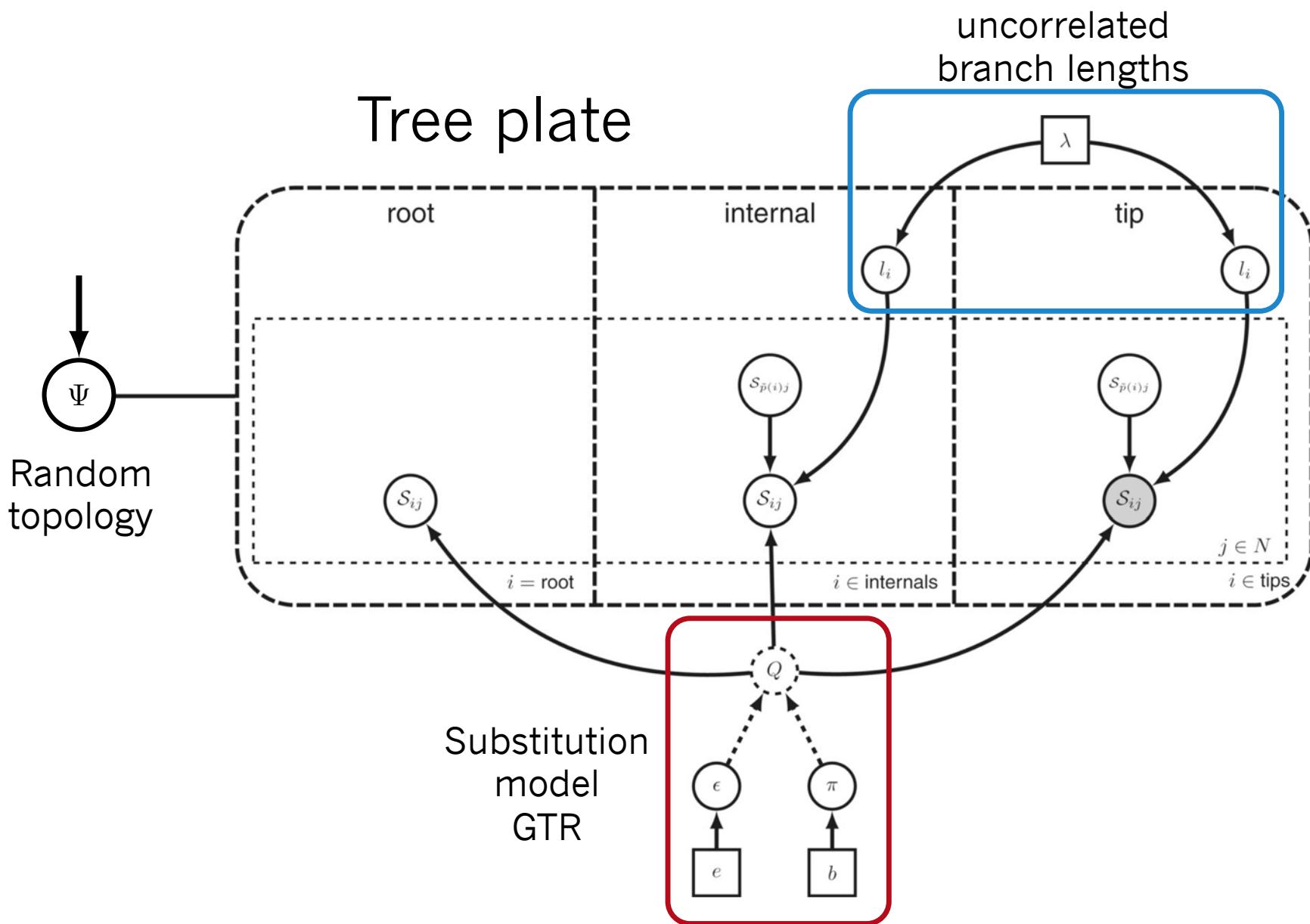
Is a tree a graph?

Yes

Is a tree a model?

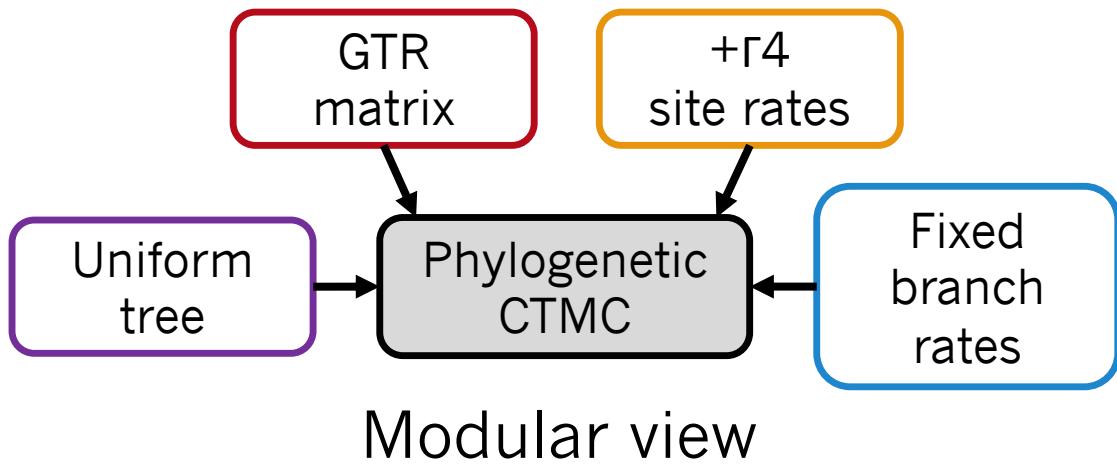
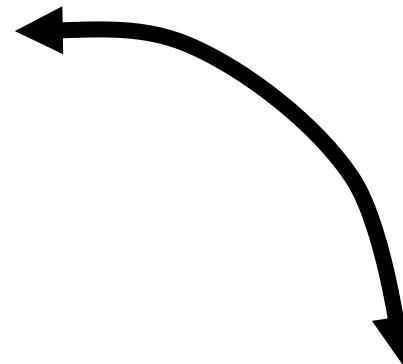
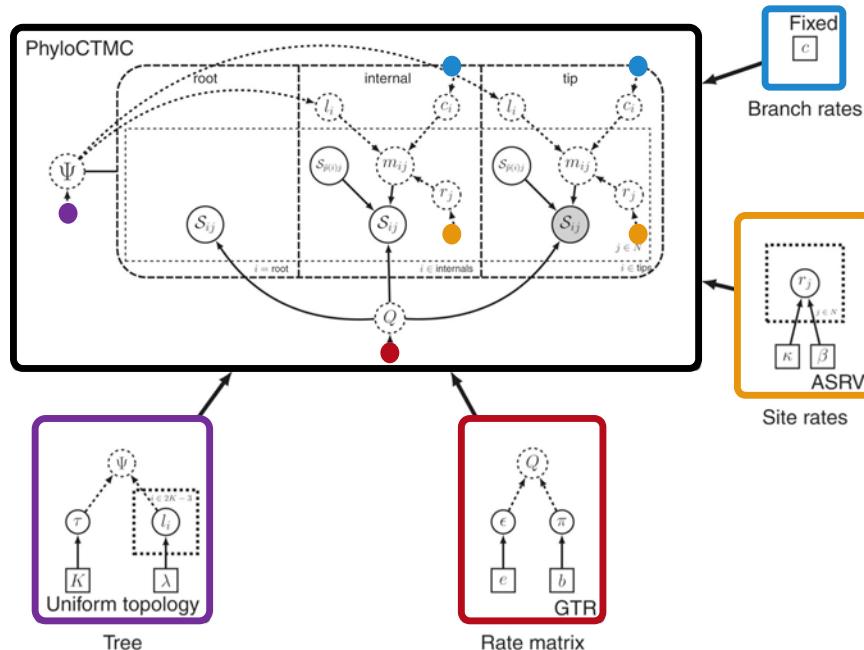
Yes

Phylogenetic substitution model



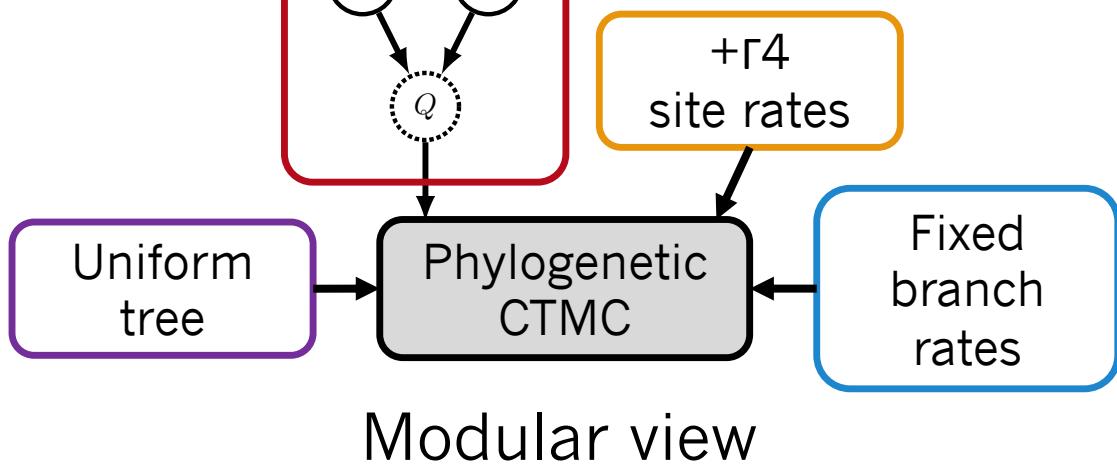
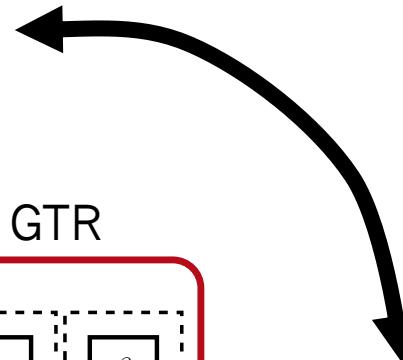
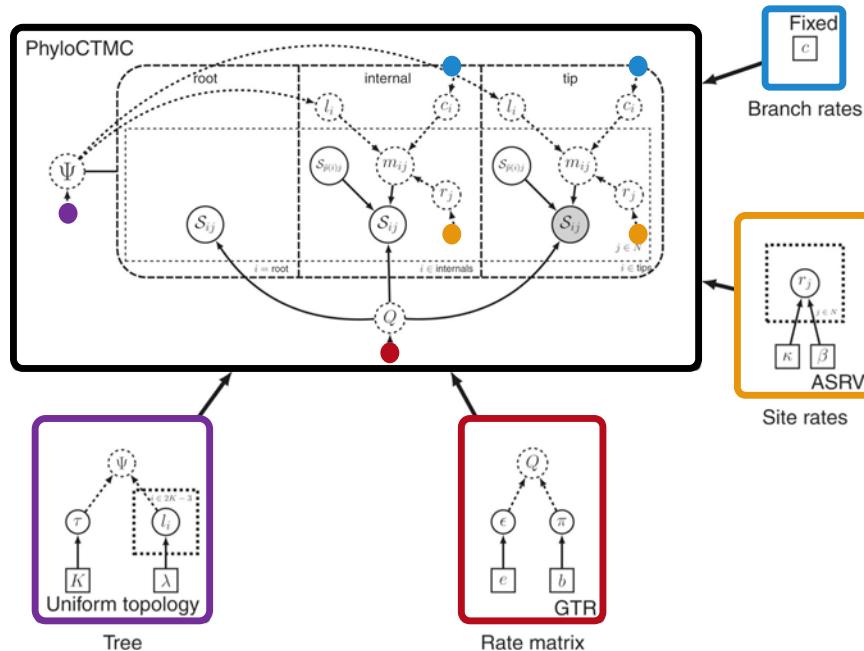
Phylogenetic substitution model

Expanded view



Phylogenetic substitution model

Expanded view

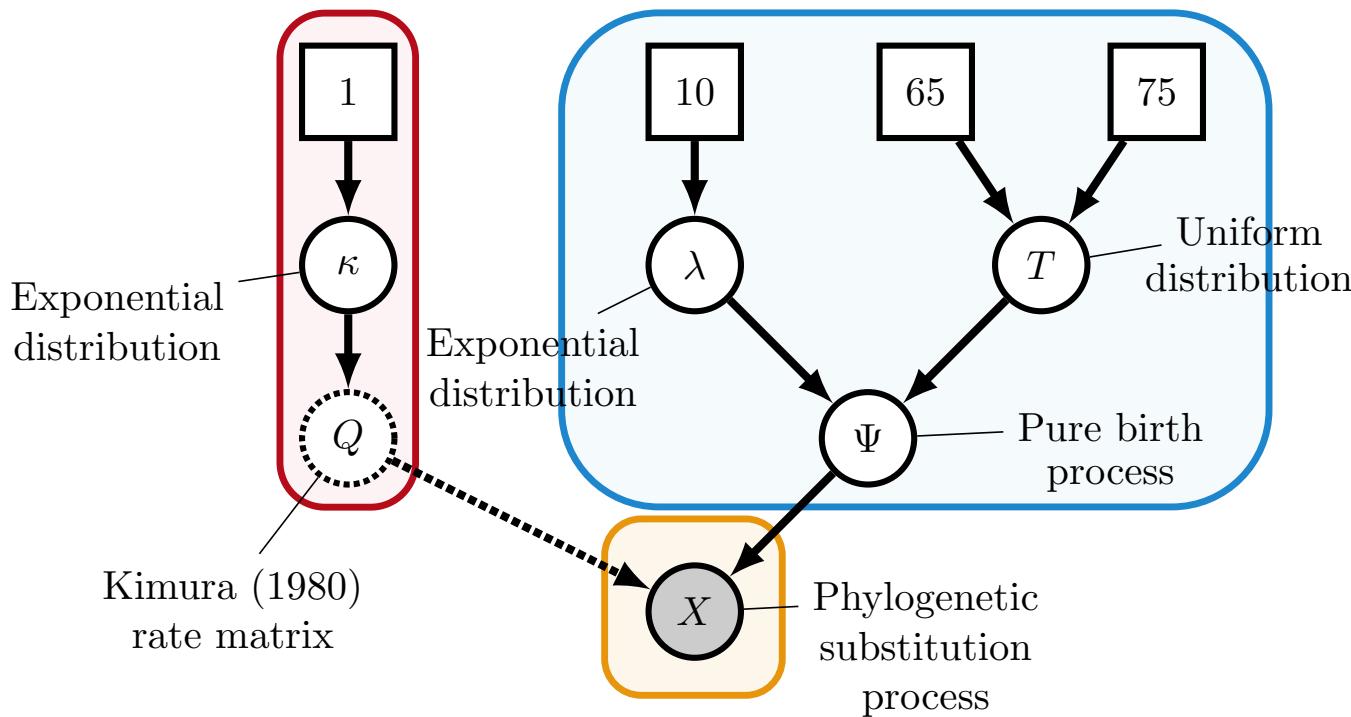


```

# pure birth prior and hyperpriors
birth ~ dnExponential(10)
tree_age ~ dnUnif(65, 75)
tree ~ dnBDP(birth, tree_age, taxa)
# rate matrix
kappa ~ dnExponential(1)
Q := fnK80(kappa)
# phylogenetic substitution process
seq ~ dnPhyloCTMC(Q, tree, "DNA")
seq.clamp(data)

```

RevBayes script



Phylogenetic
graphical
model
(Yule + K80)

```

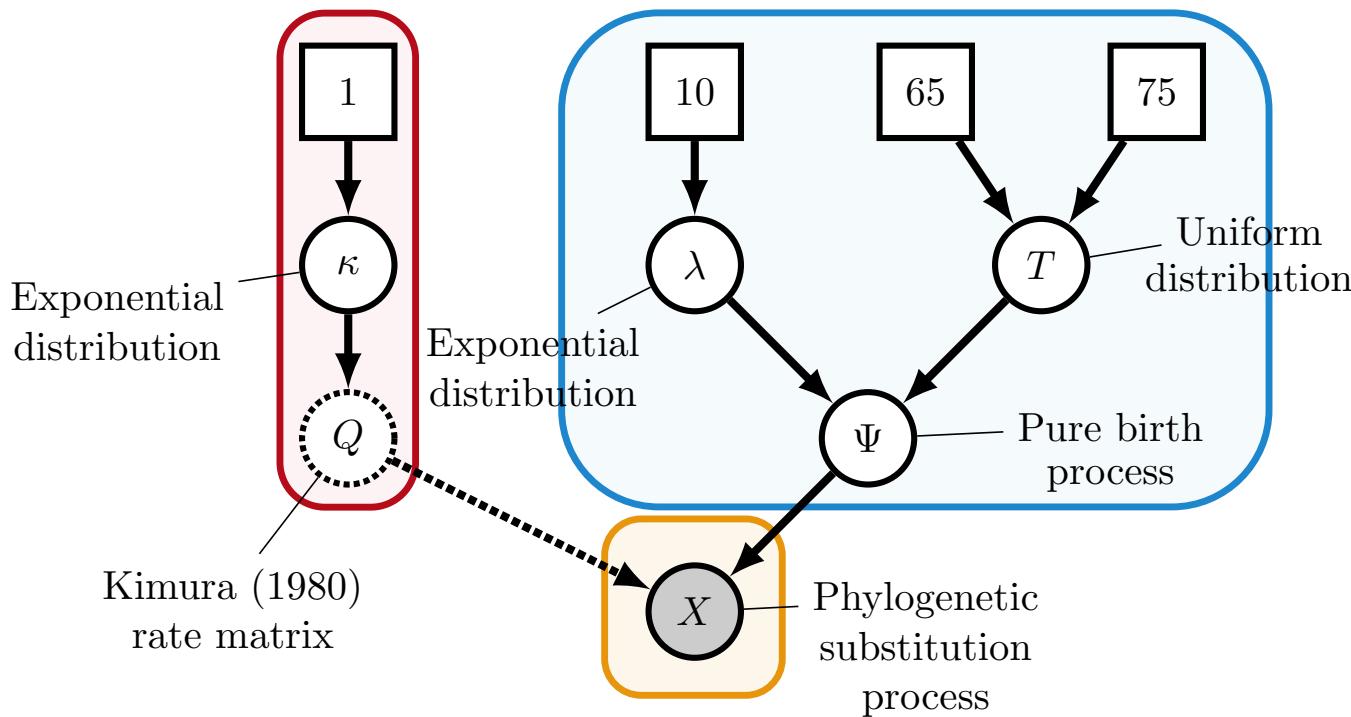
# pure birth prior and hyperpriors
source("module_tree_yule.Rev")

# rate matrix
source("module_Q_K80.Rev")

# phylogenetic substitution process
source("module_CTMC.Rev")

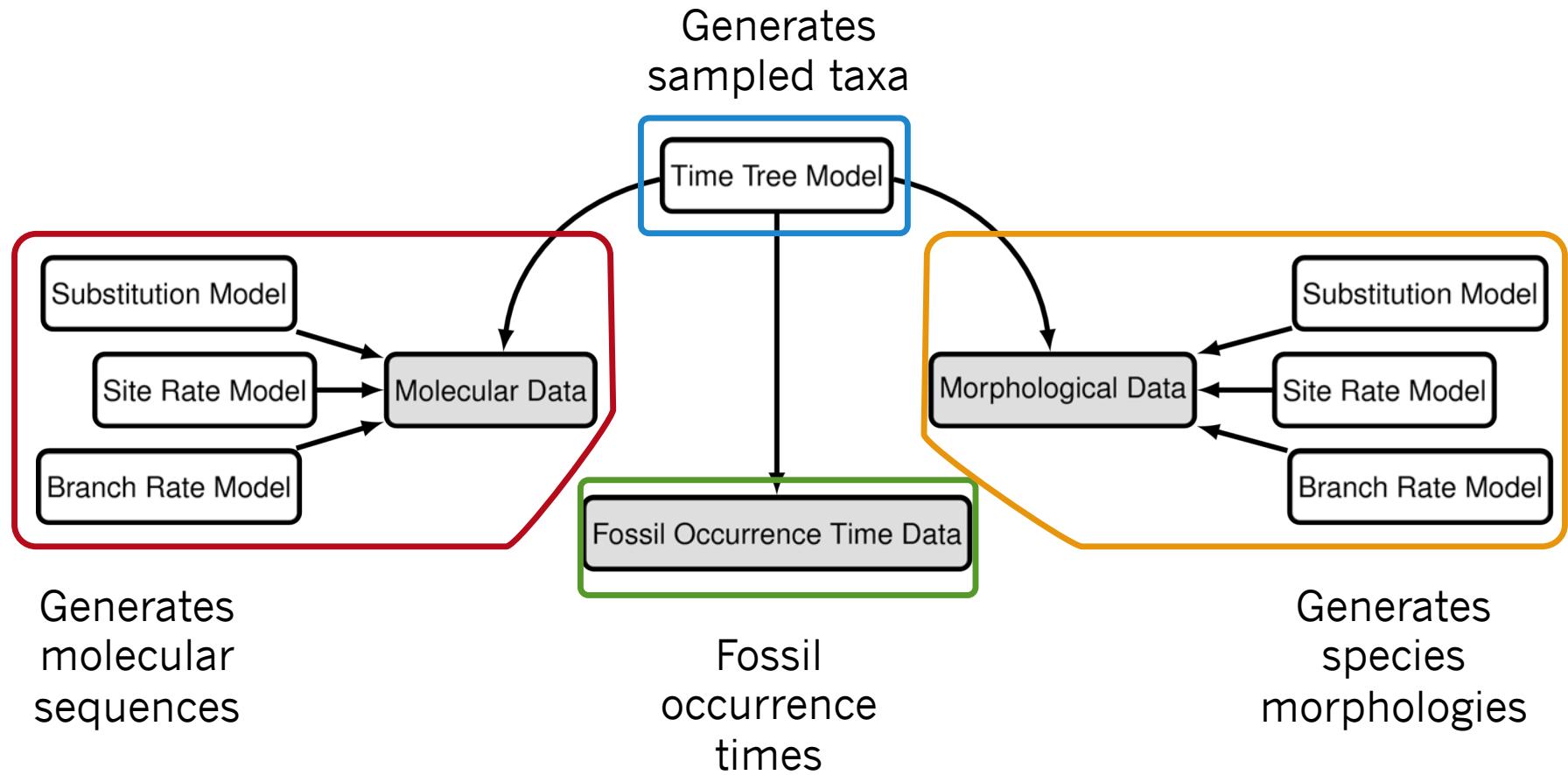
```

RevBayes script



Phylogenetic graphical model (Yule + K80)

Modular View of Combined-Evidence Analysis under a Fossilized Birth-Death Process



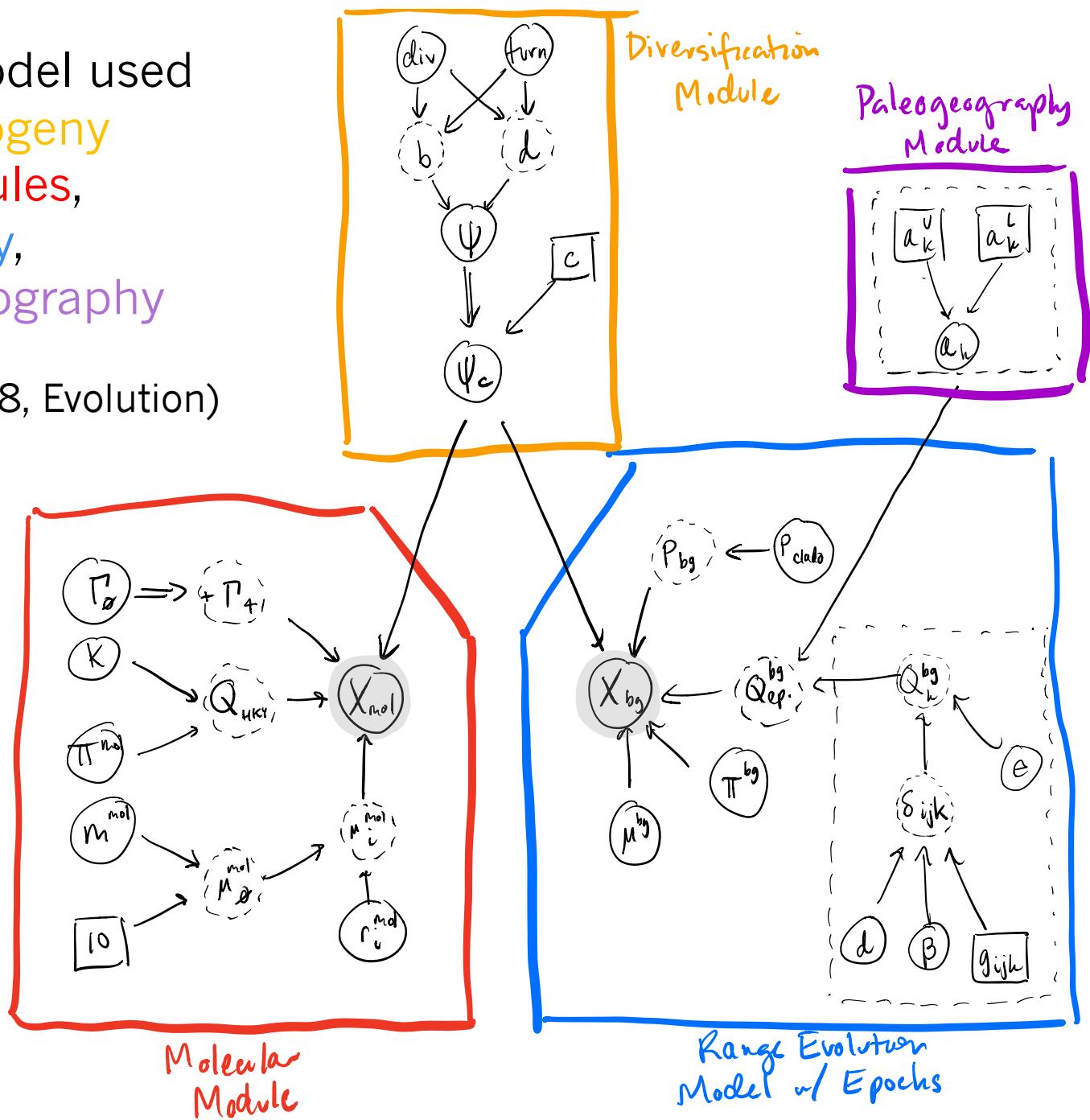
(Tracy will teach this evening)

Sketch of model used
to date **phylogeny**
using **molecules**,
biogeography,
and **paleogeography**

Landis et al. (2018, Evolution)



Hawaiian
silverswords



All possible
processes

All possible
processes

Natural
processes

Generates
biological
data

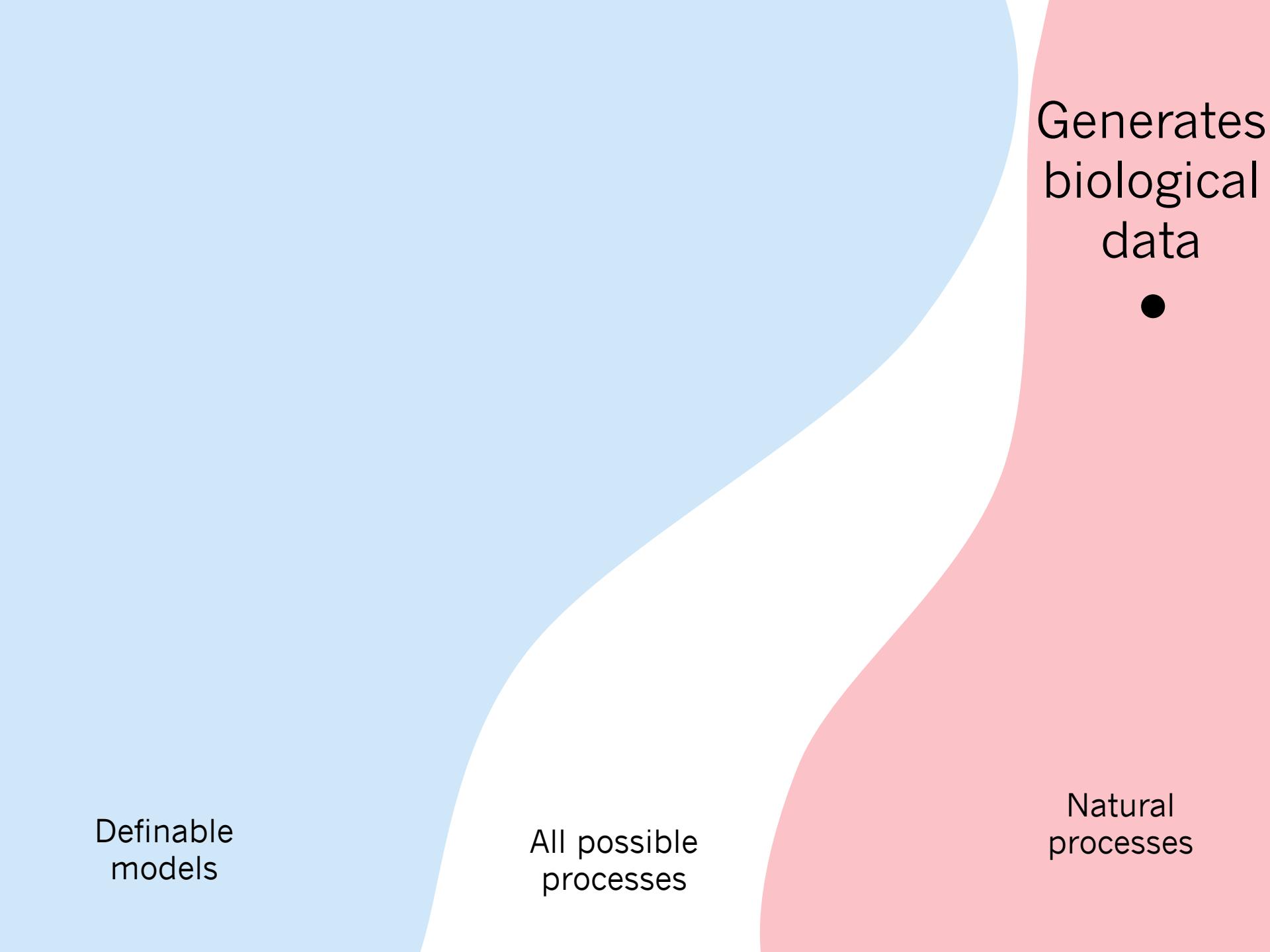


Studying evolution

Common descent
phylogenetic inference

Tempo & mode
macroevolutionary models

Rewinding deep time
ancestral state estimation

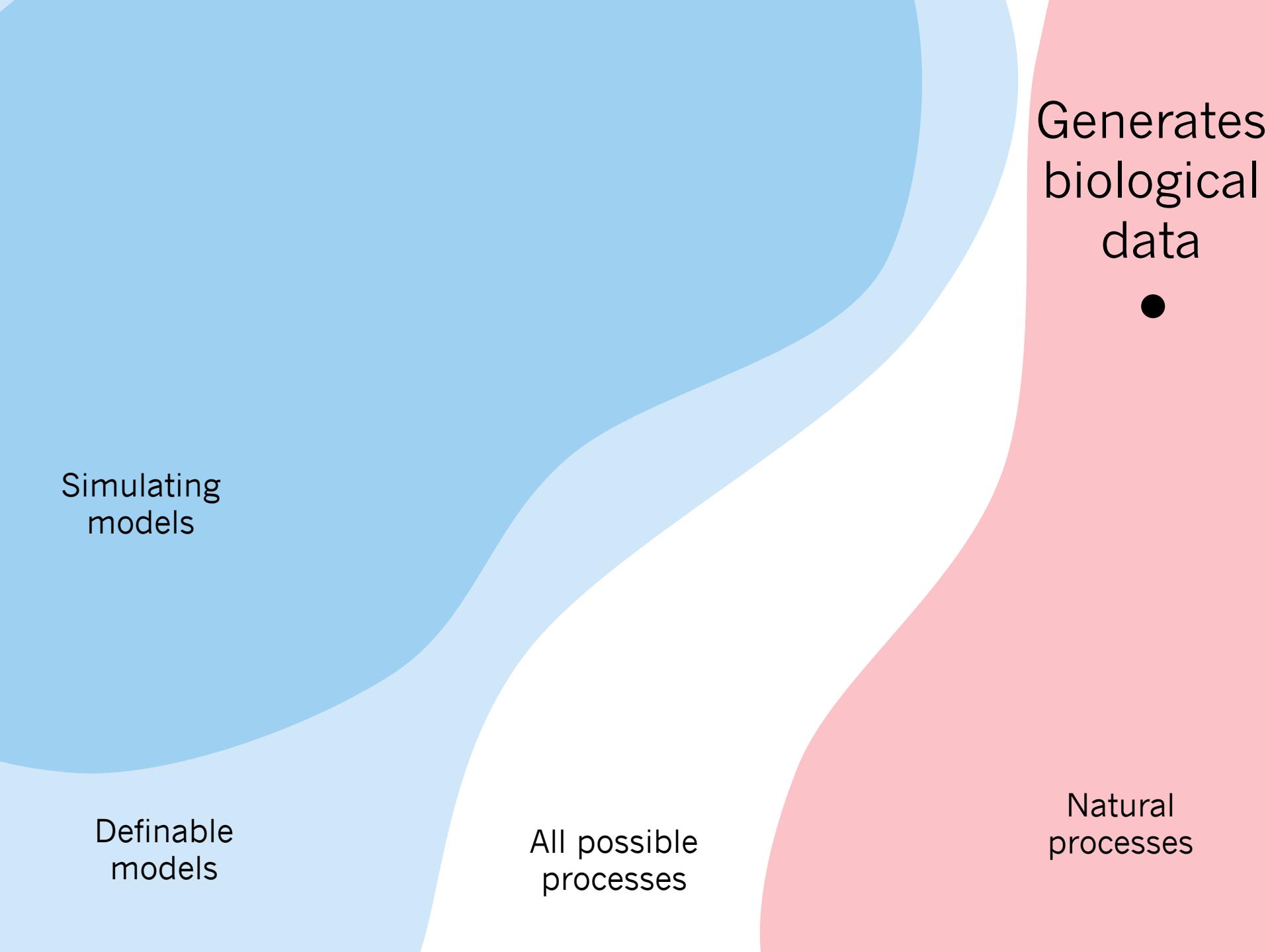


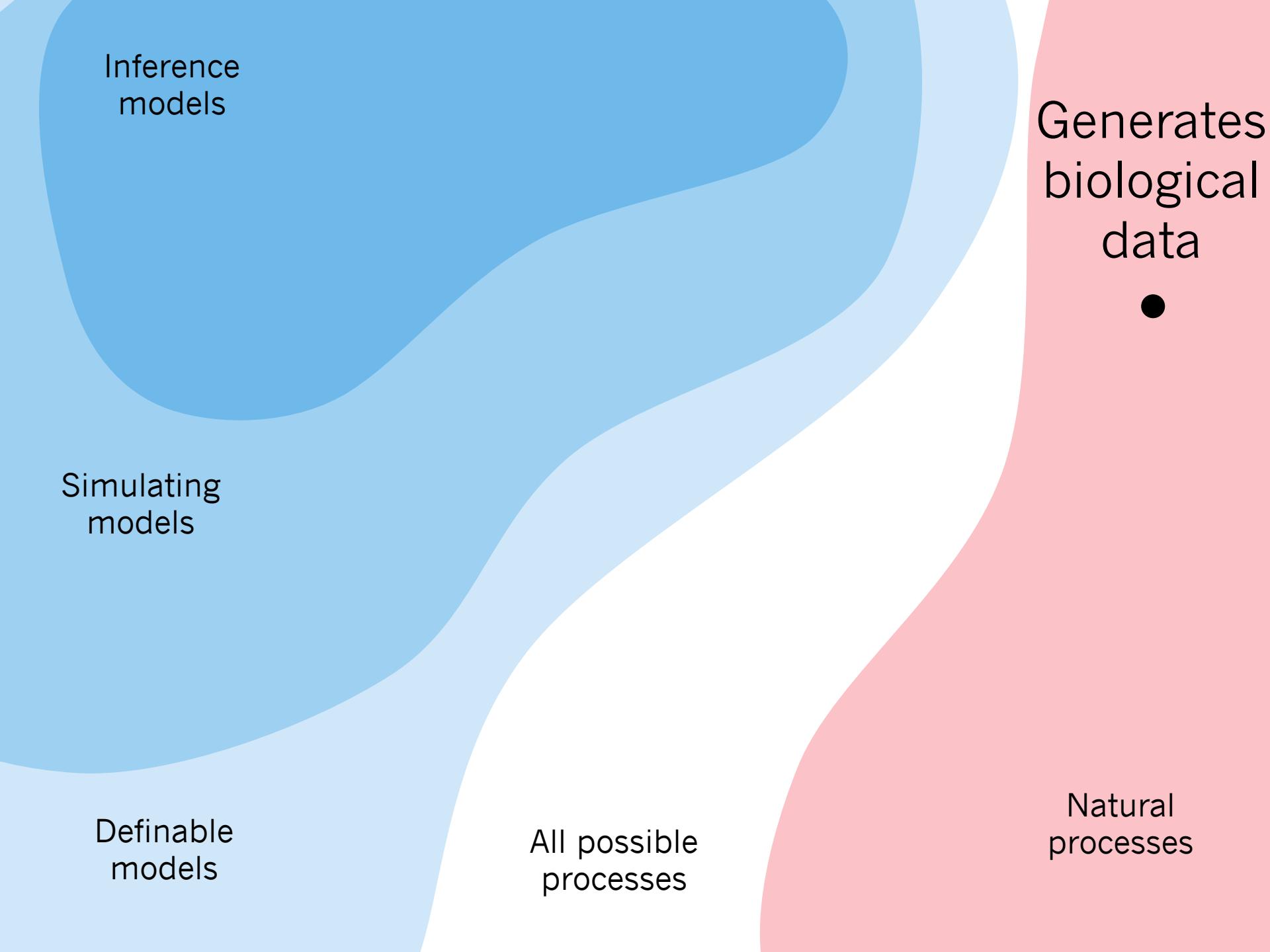
Definable
models

All possible
processes

Natural
processes

Generates
biological
data





Inference
models

Simulating
models

Definable
models

All possible
processes

Natural
processes

Generates
biological
data

Inference
models

Simulating
models

Definable
models

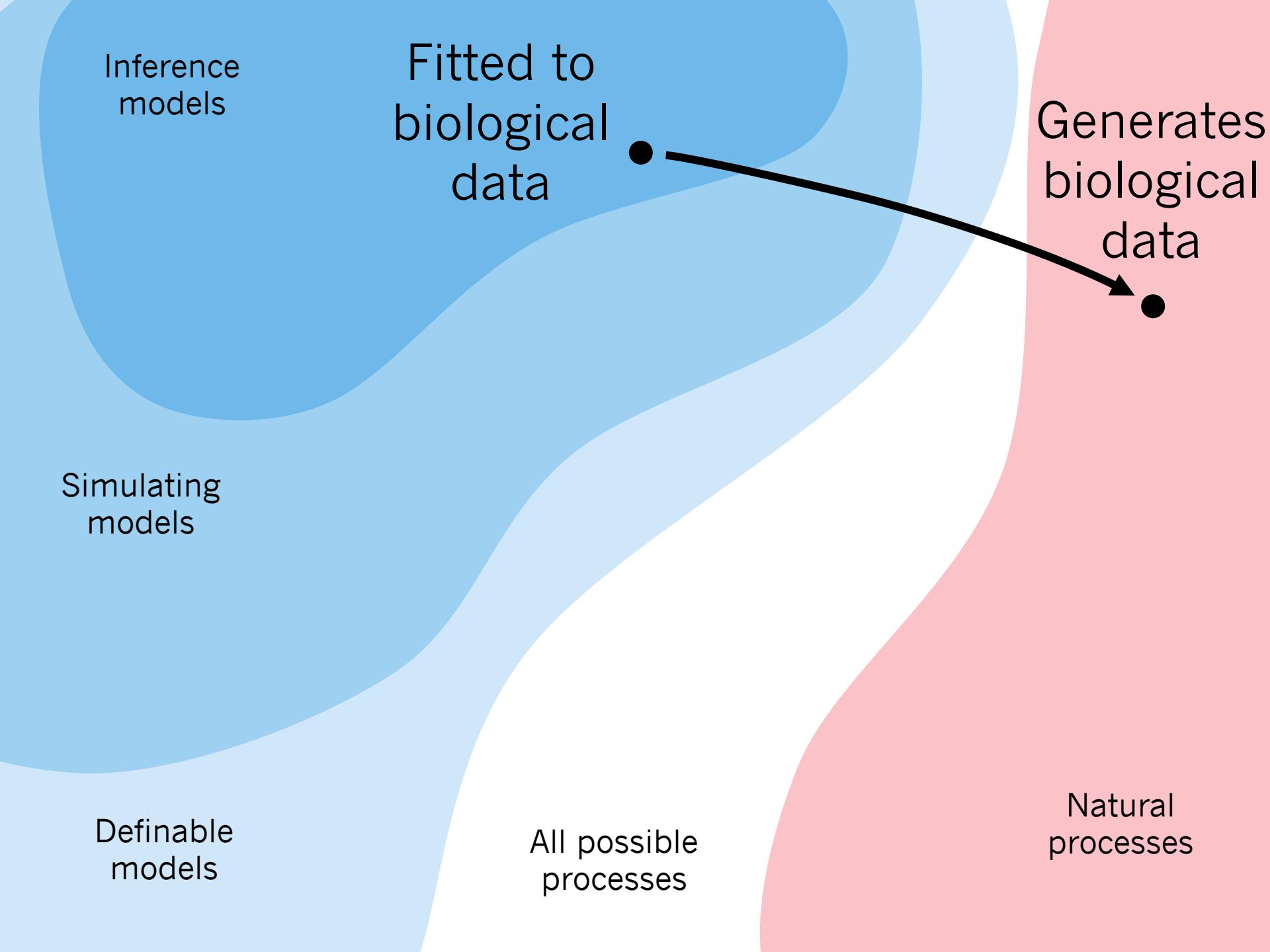
All possible
processes

Generates
biological
data

- ~~No exact model~~



Natural
processes



Inference
models

Fitted to
biological
data

Simulating
models

Definable
models

All possible
processes

Natural
processes

Generates
biological
data

Inference
models

Fitted to
biological
data

Simulating
models

Definable
models

Generates
biological
data

How useful is
our wrong model?

All possible
processes



Box

Assessing models

accuracy	how close are our estimates to the true parameter values?
precision	how are we certain of our parameter estimates?
power	what size datasets are needed to estimate the truth?
sensitivity	are parameter estimates responsive to prior/data changes?
performance	does the model require excessive time or resources?
selection	do we reliably select the true model among competing models?
adequacy	does the model generate data resembling the true data?

Assessing models

accuracy

how close are our estimates to the true parameter values?

precision

how are we certain of our parameter estimates?

power

what size datasets are needed to estimate the truth?

sensitivity

are parameter estimates responsive to prior/data changes?

performance

does the model require excessive time or resources?

selection

do we reliably select the true model among competing models?

does the model generate data resembling the true data?



adequacy

Inference
models

Fitted to
biological
data

Simulating
models

Definable
models

Generates
biological
data

No need for inference
if we knew the truth!

All possible
processes

Natural
processes

Simulating the truth

Poor model performance under known conditions will be equally bad (or worse) under unknown conditions

Simulating the truth

Poor model performance under known conditions will be equally bad (or worse) under unknown conditions

Simulations are crucial!

1. Test model/method behavior
2. Reveal a model's personality
3. Provide context for empirical results

Inference
models

Fitted to
simulated
data

Generates
simulated
data



identical
models

Simulating
models

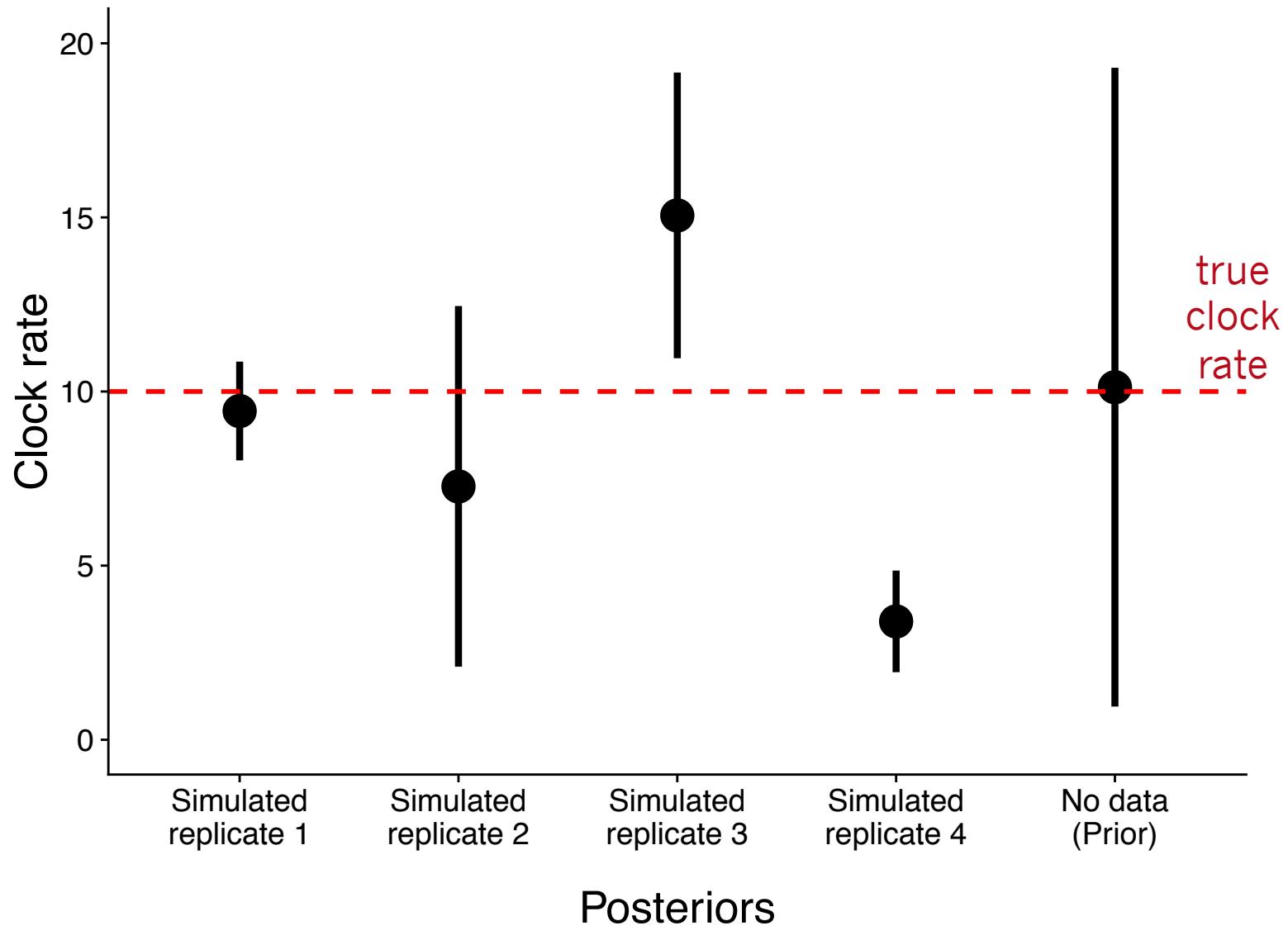
A perfect world where
models can be right

Definable
models

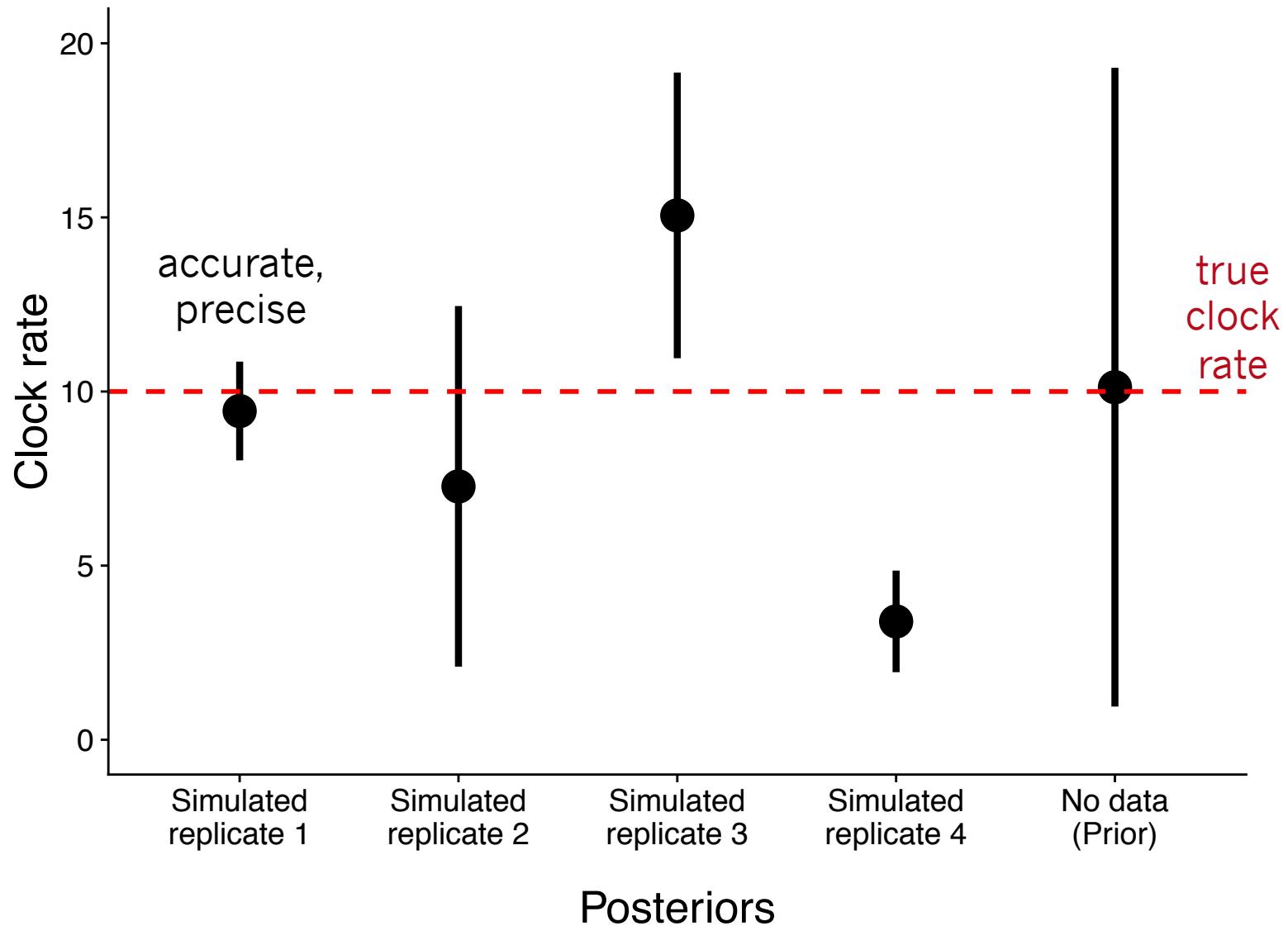
All possible
processes

Natural
processes

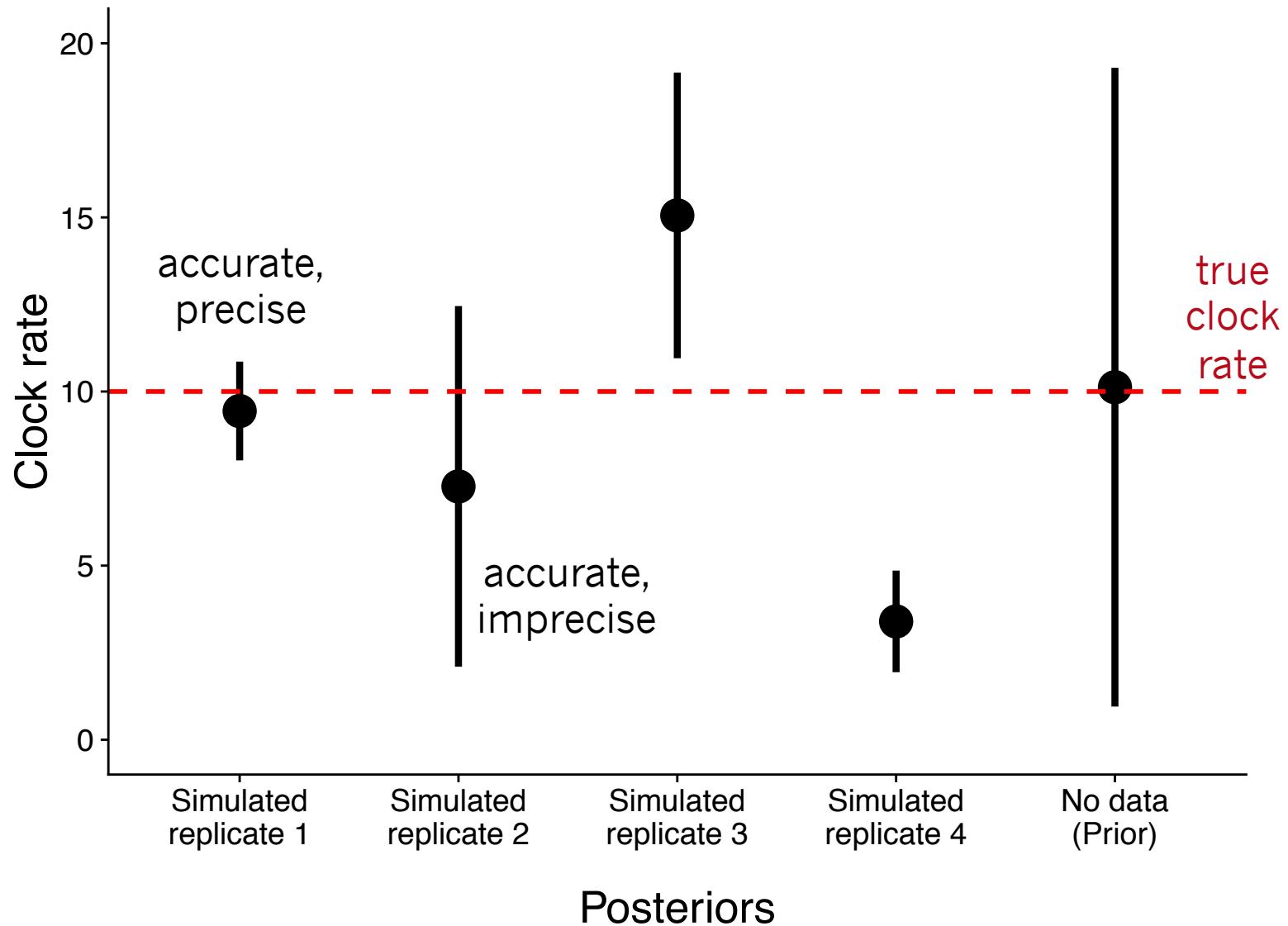
Accuracy and precision



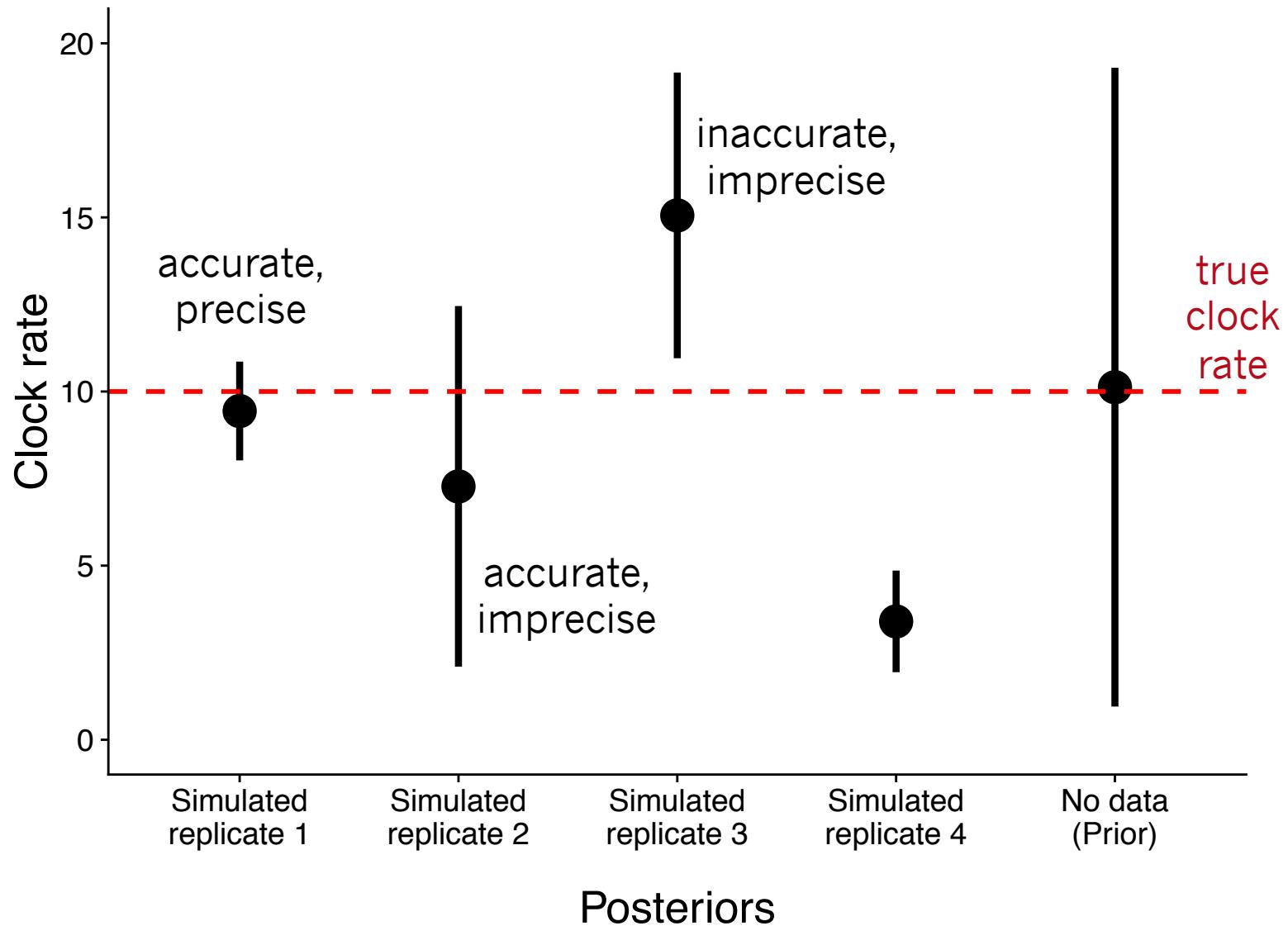
Accuracy and precision



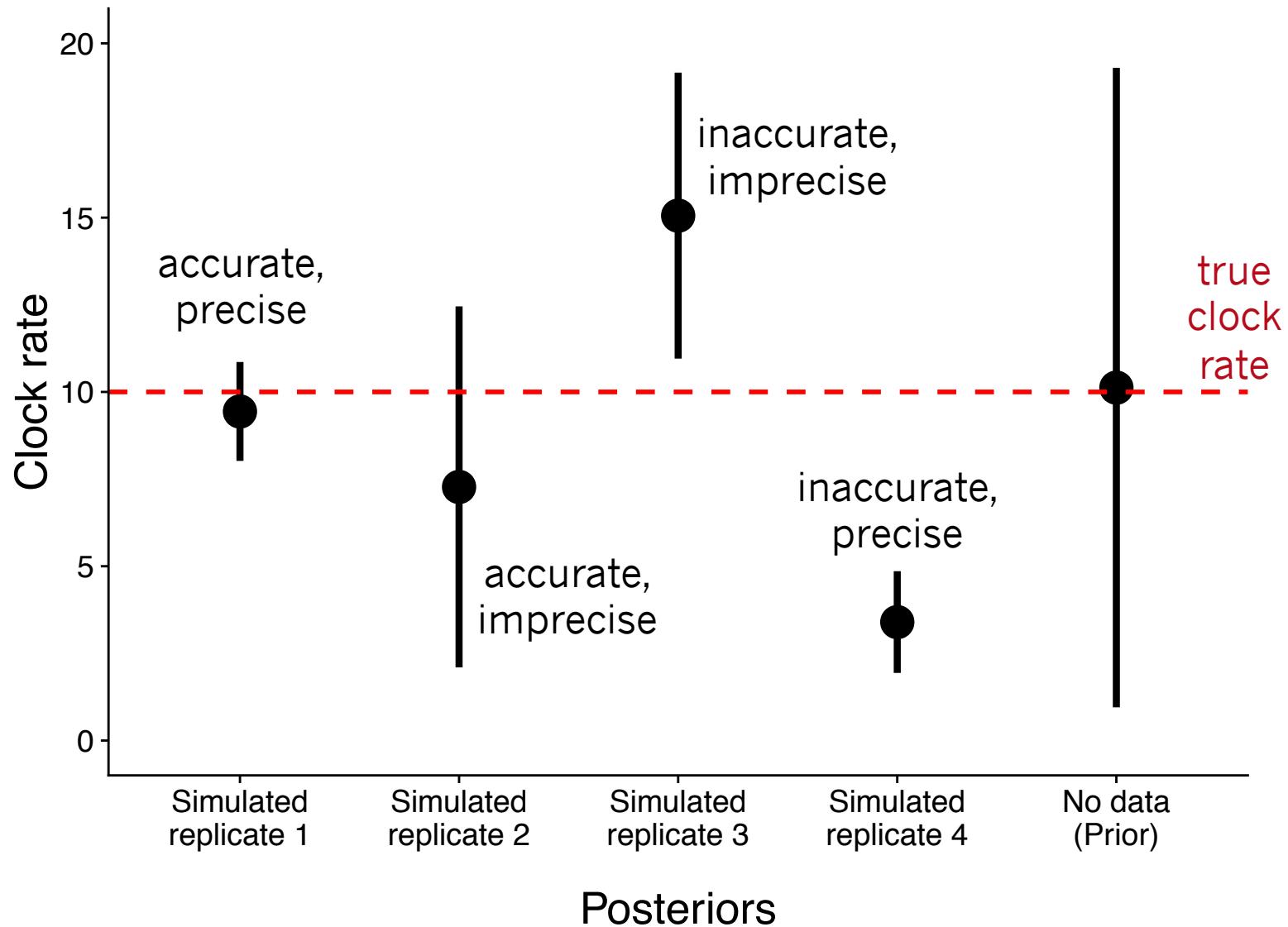
Accuracy and precision



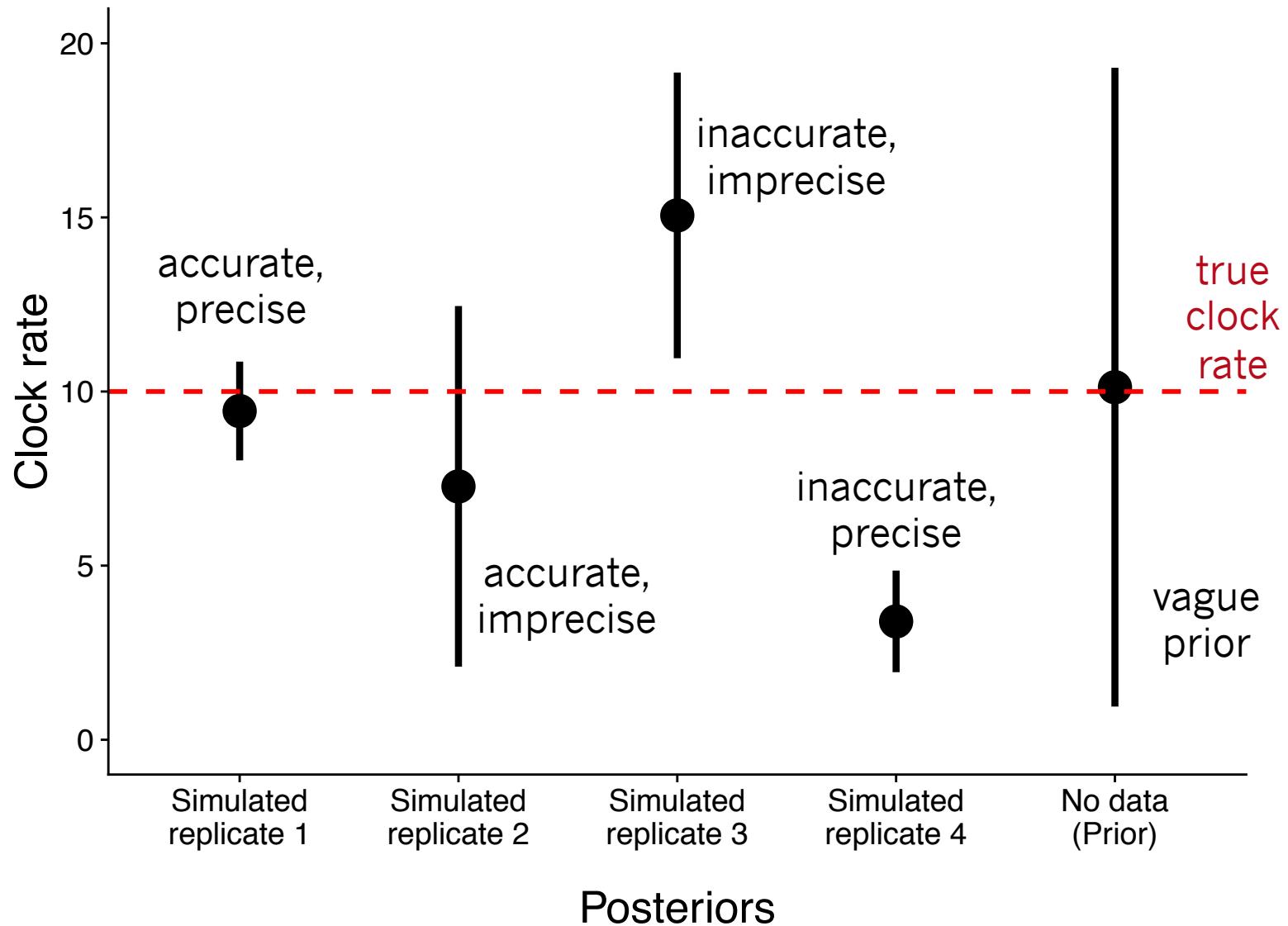
Accuracy and precision



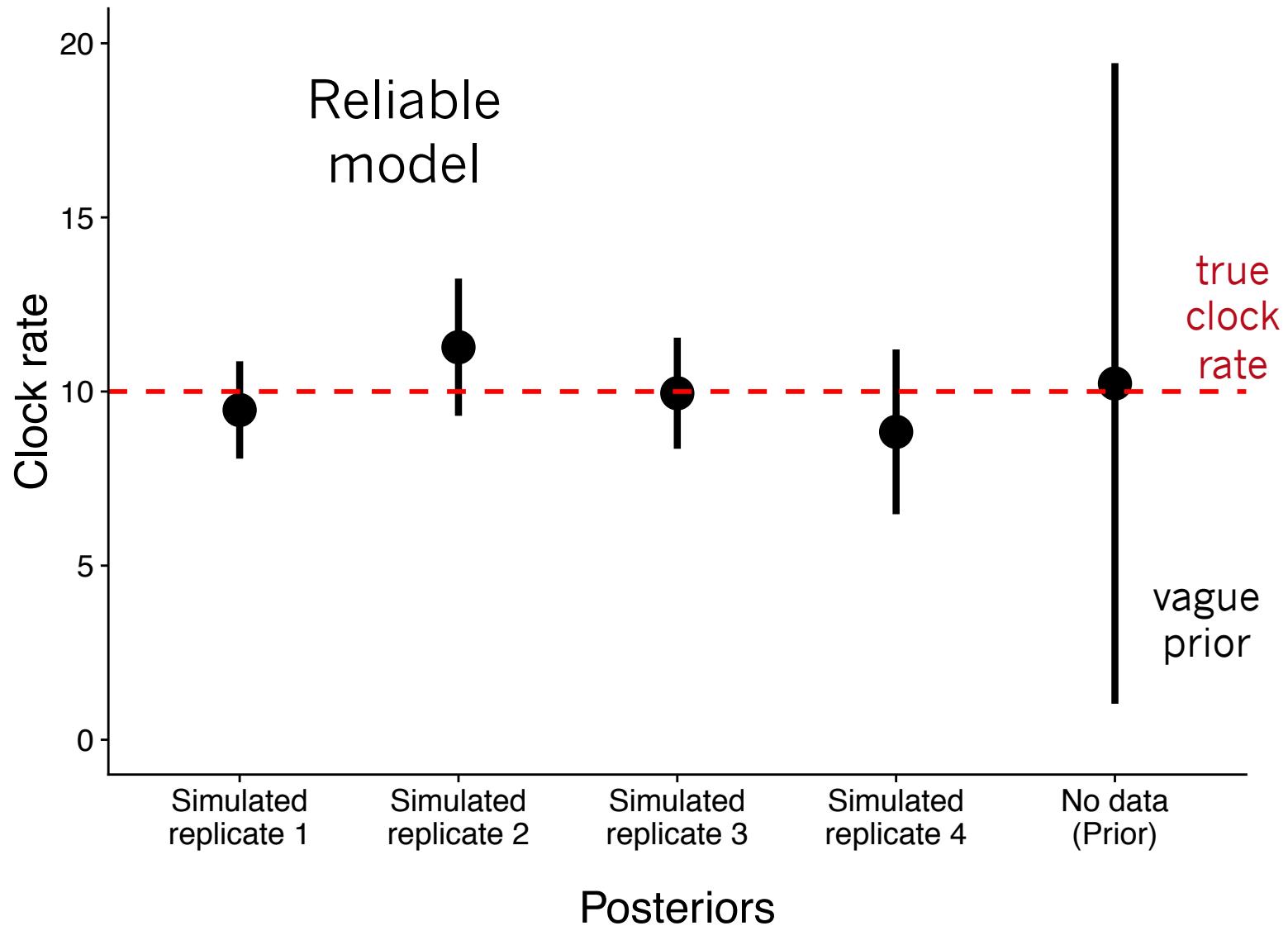
Accuracy and precision



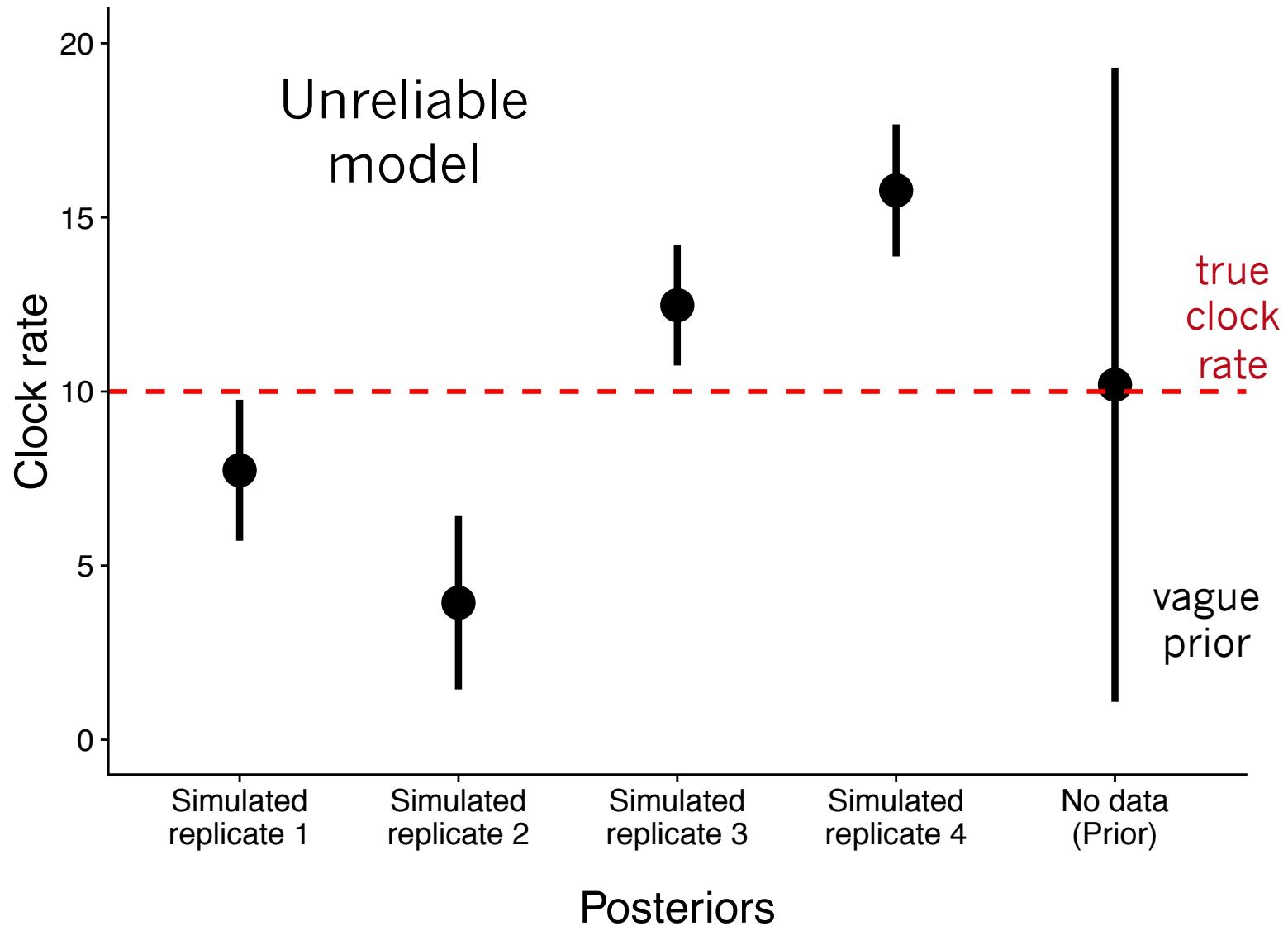
Accuracy and precision



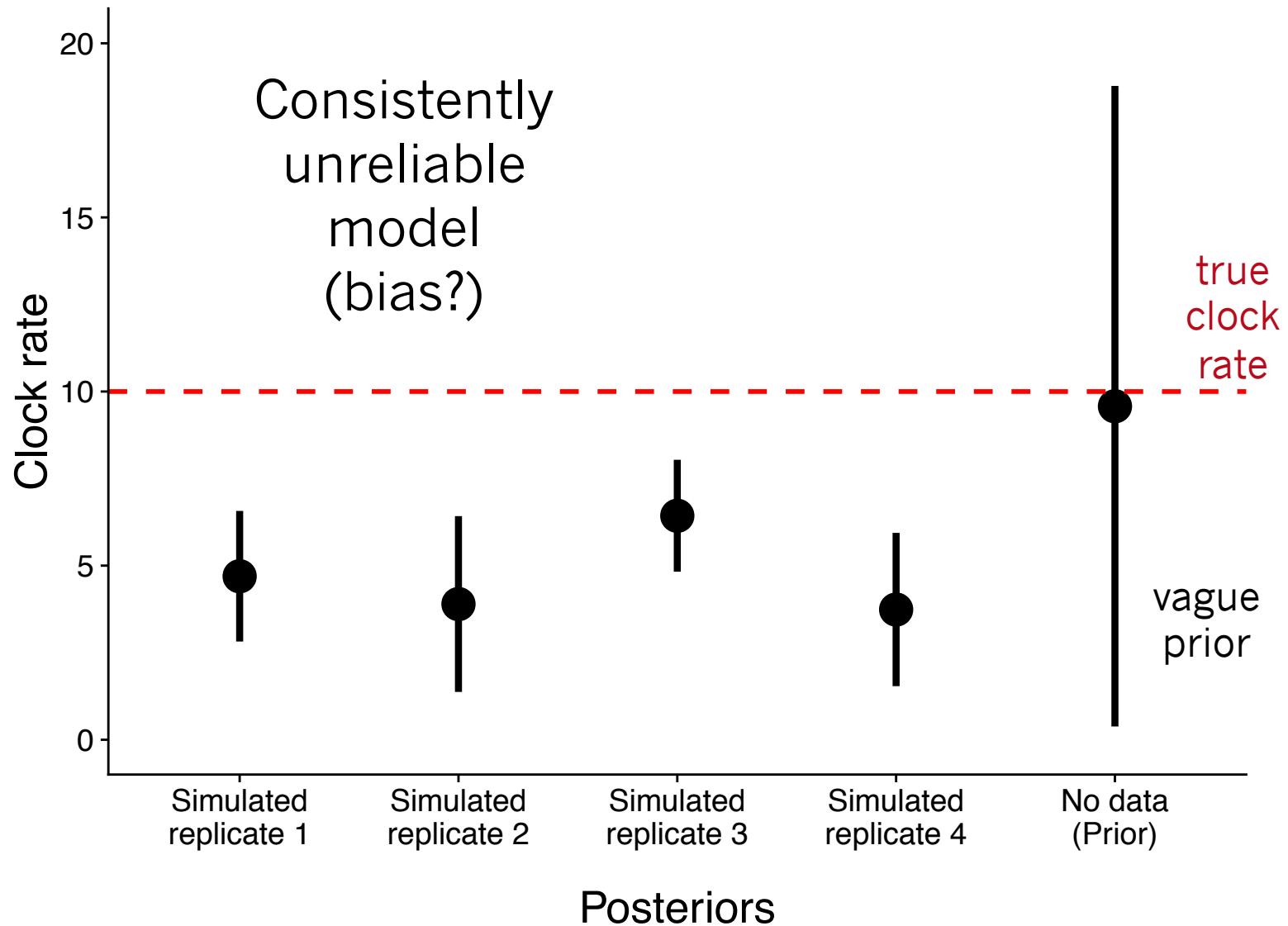
Accuracy and precision



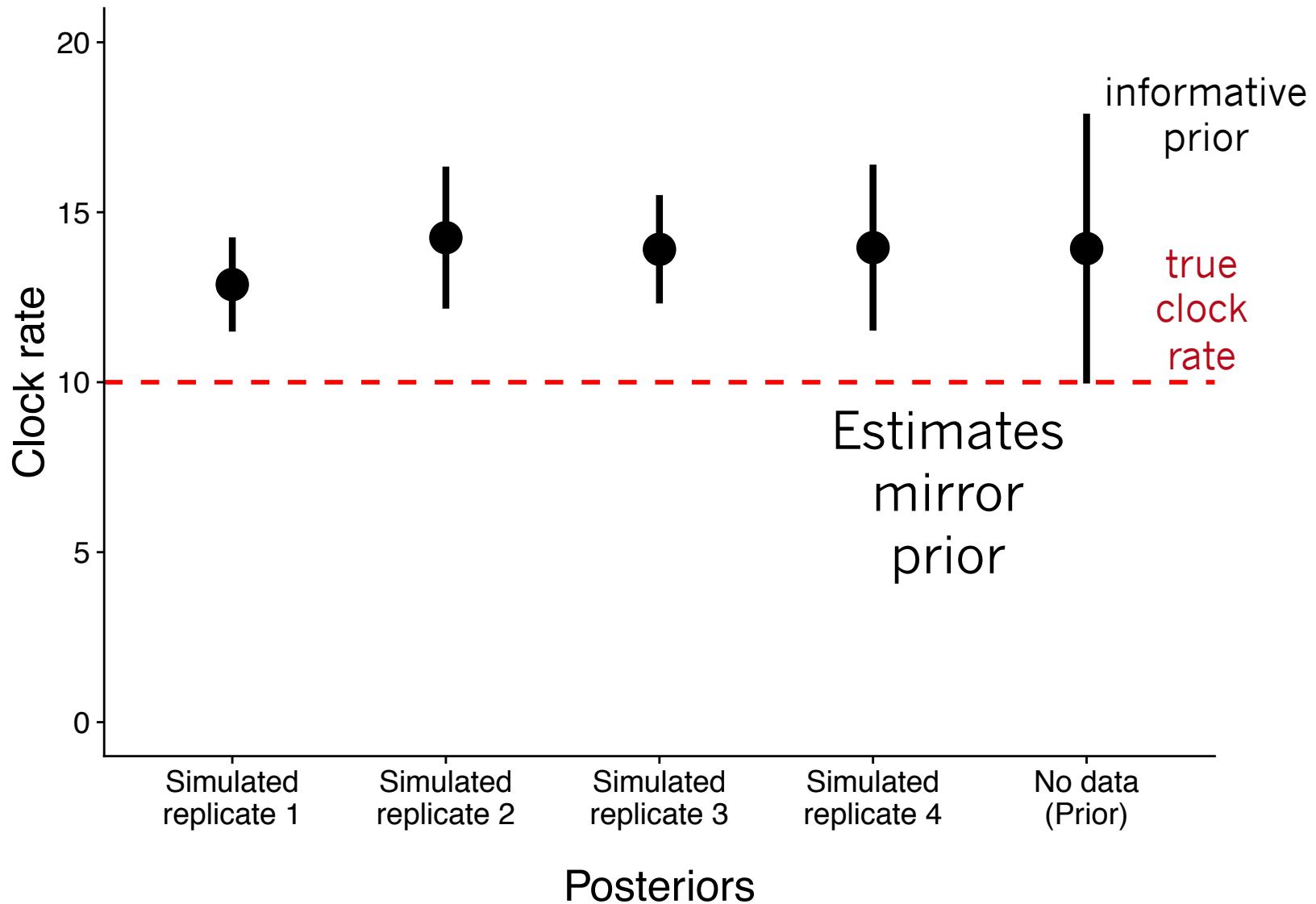
Accuracy and precision



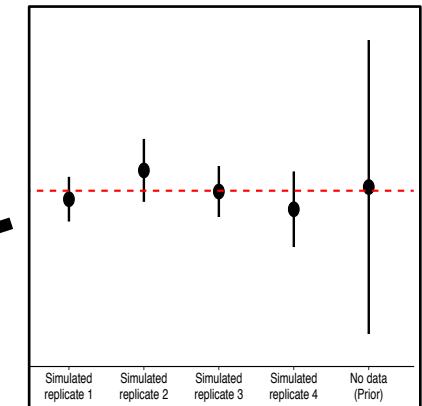
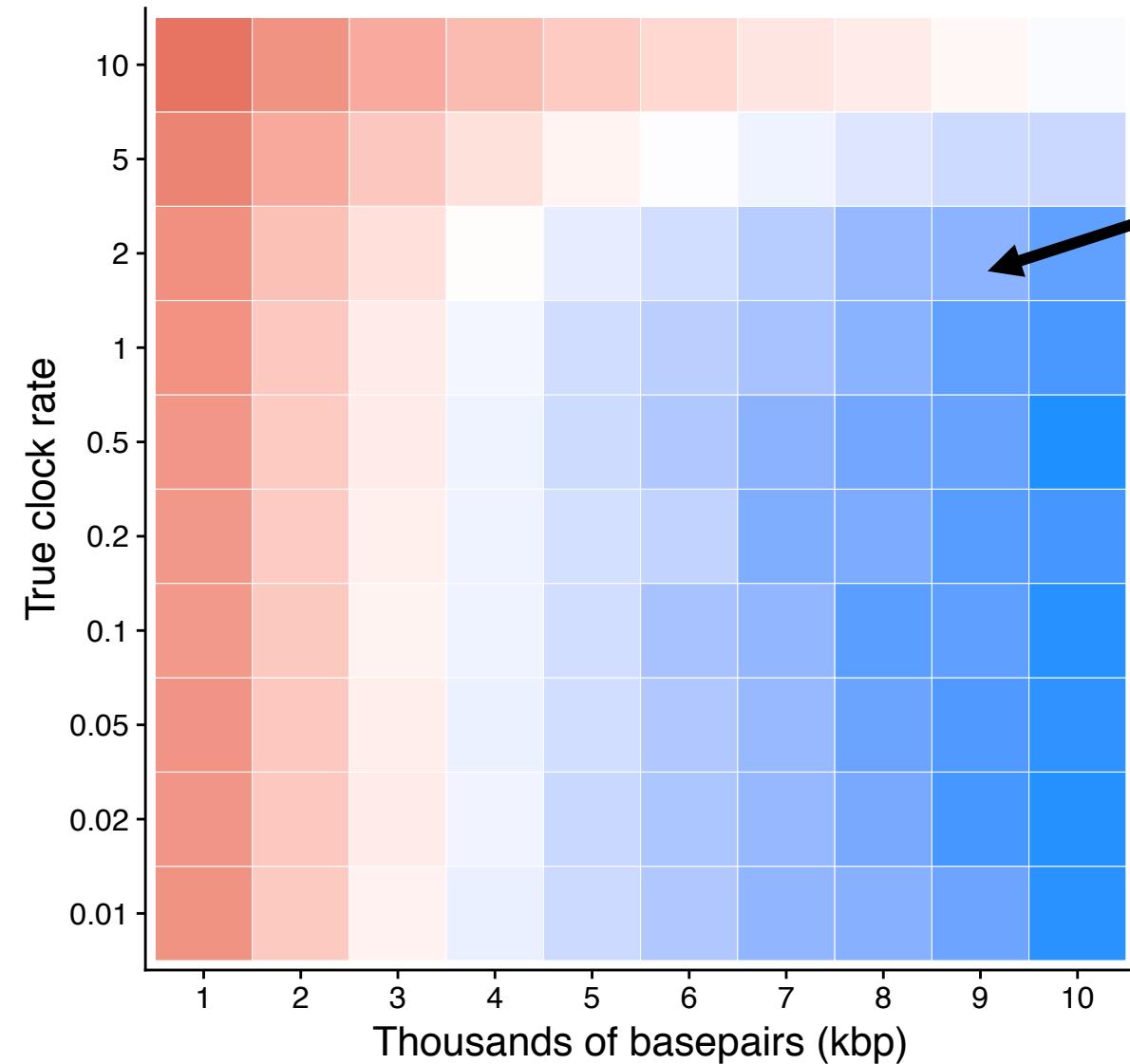
Accuracy and precision



Prior sensitivity



Accuracy and power

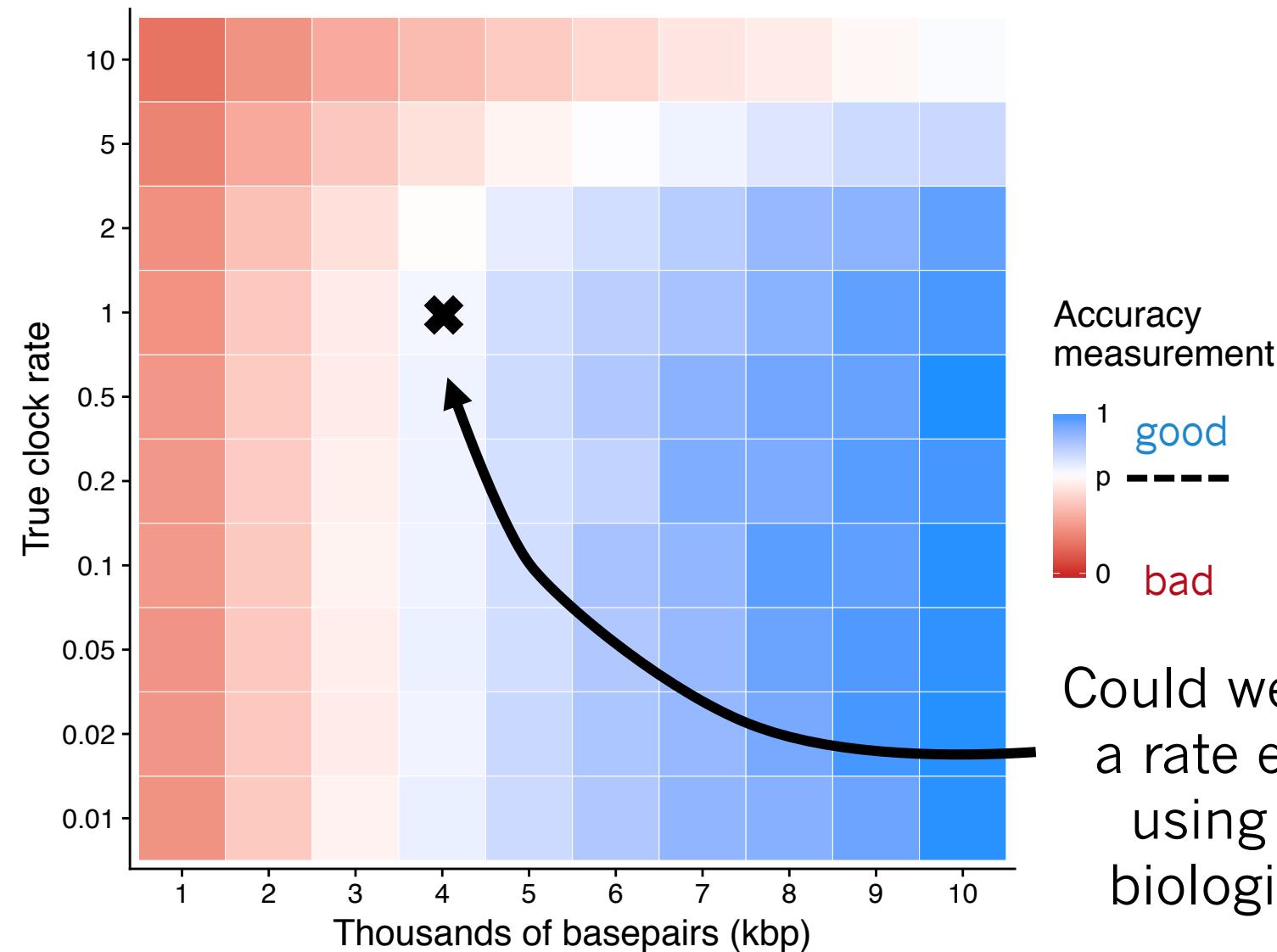


Accuracy
measurement

1 **good**
p **-----**
0 **bad**

Each cell reports
frequency of pass/fail
for an accuracy test
for 100 tests

Accuracy and power

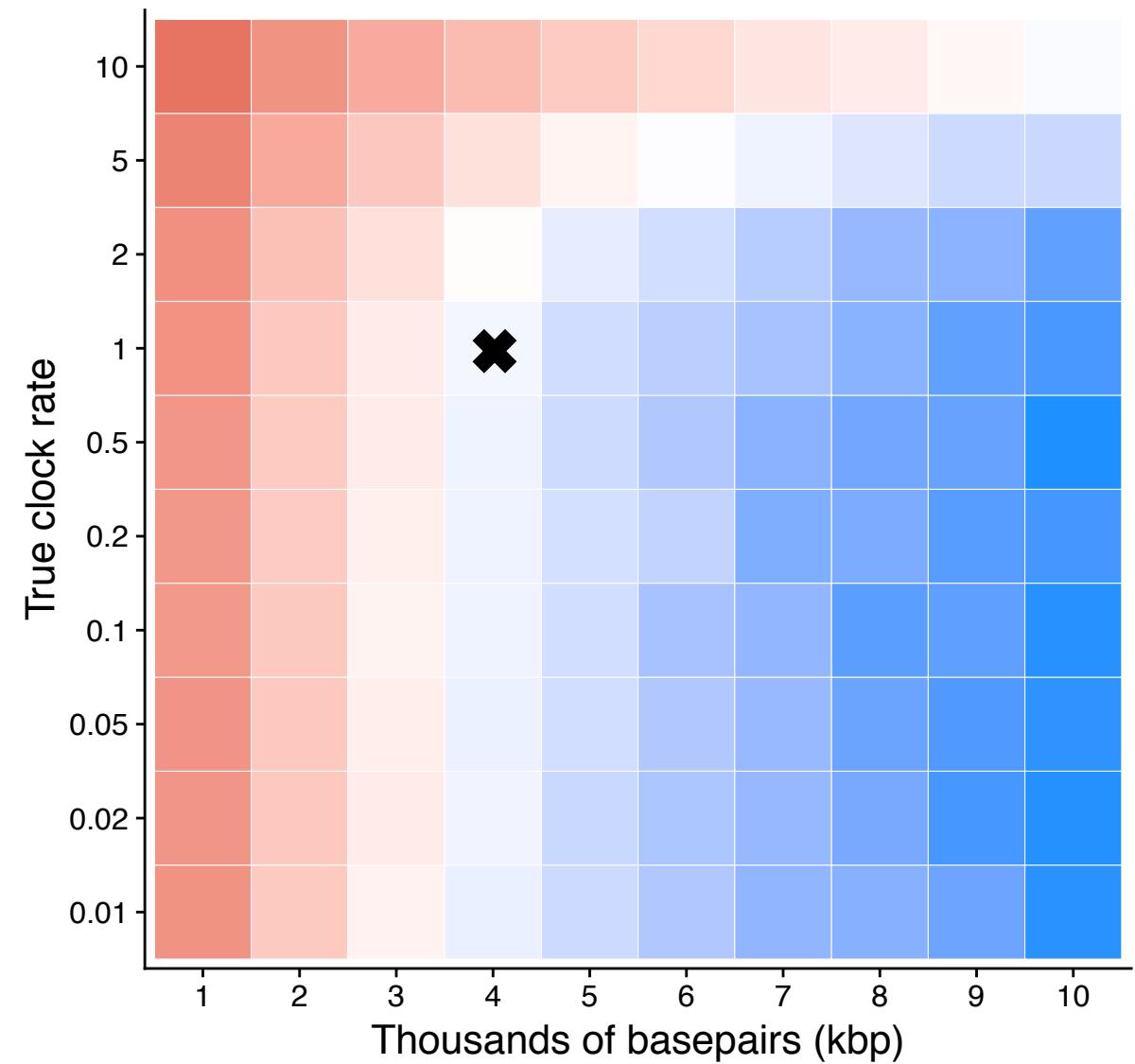


Accuracy
measurement

1 good
p -----
0 bad

Could we estimate
a rate equal to 1
using 4 kbp of
biological data?

Accuracy and power

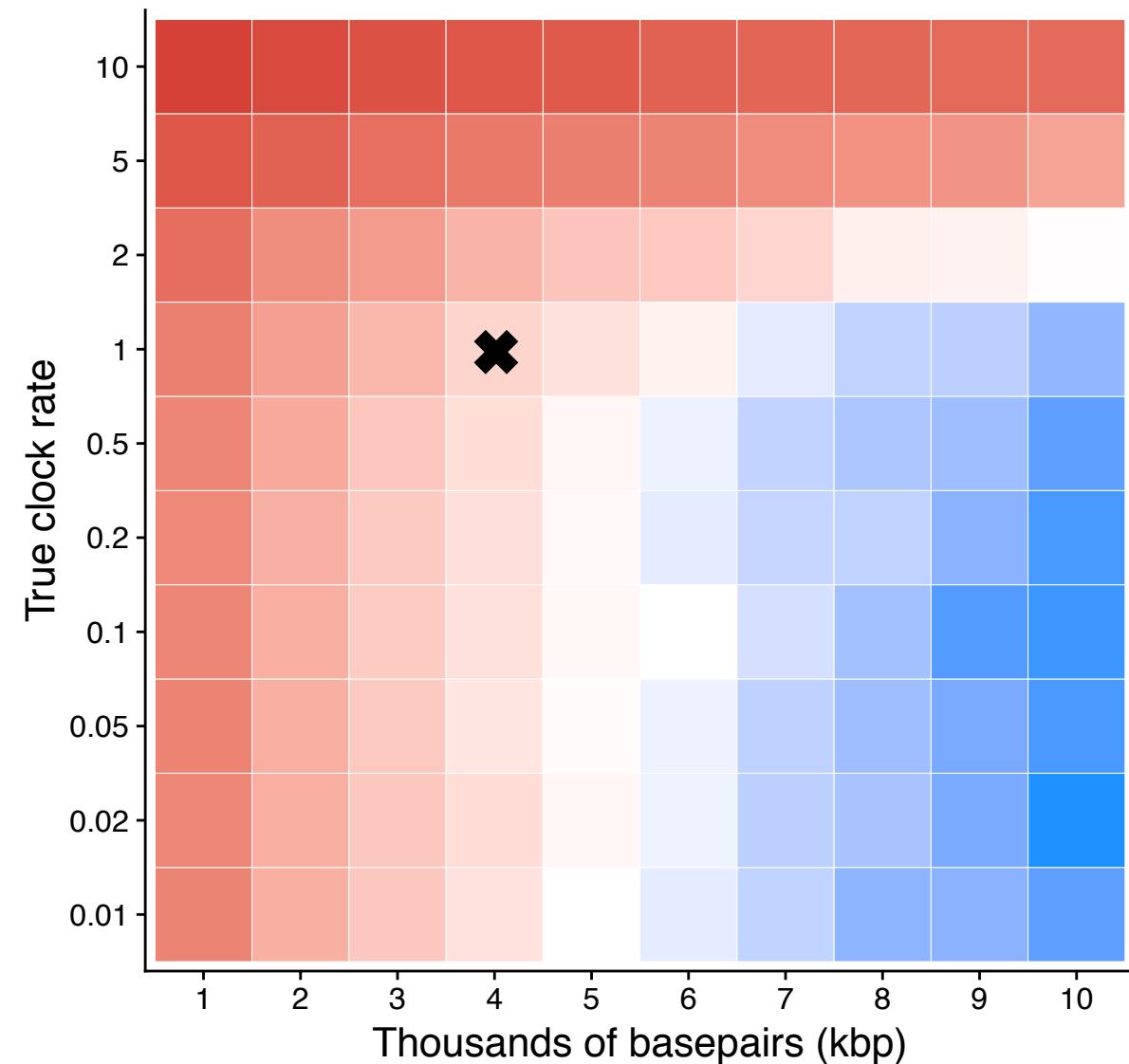


Only 4 kbp needed
for desired accuracy
for rate = 1

Accuracy
measurement

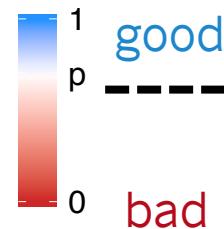
1 good
p -----
0 bad

Precision and power



Only 4 kbp needed
for desired accuracy
for rate = 1

Precision
measurement



...but 4 kbp generally
lacks sufficient
precision
for rate = 1

Inference
models

Generates
simulated
data

Simulating
models

Definable
models

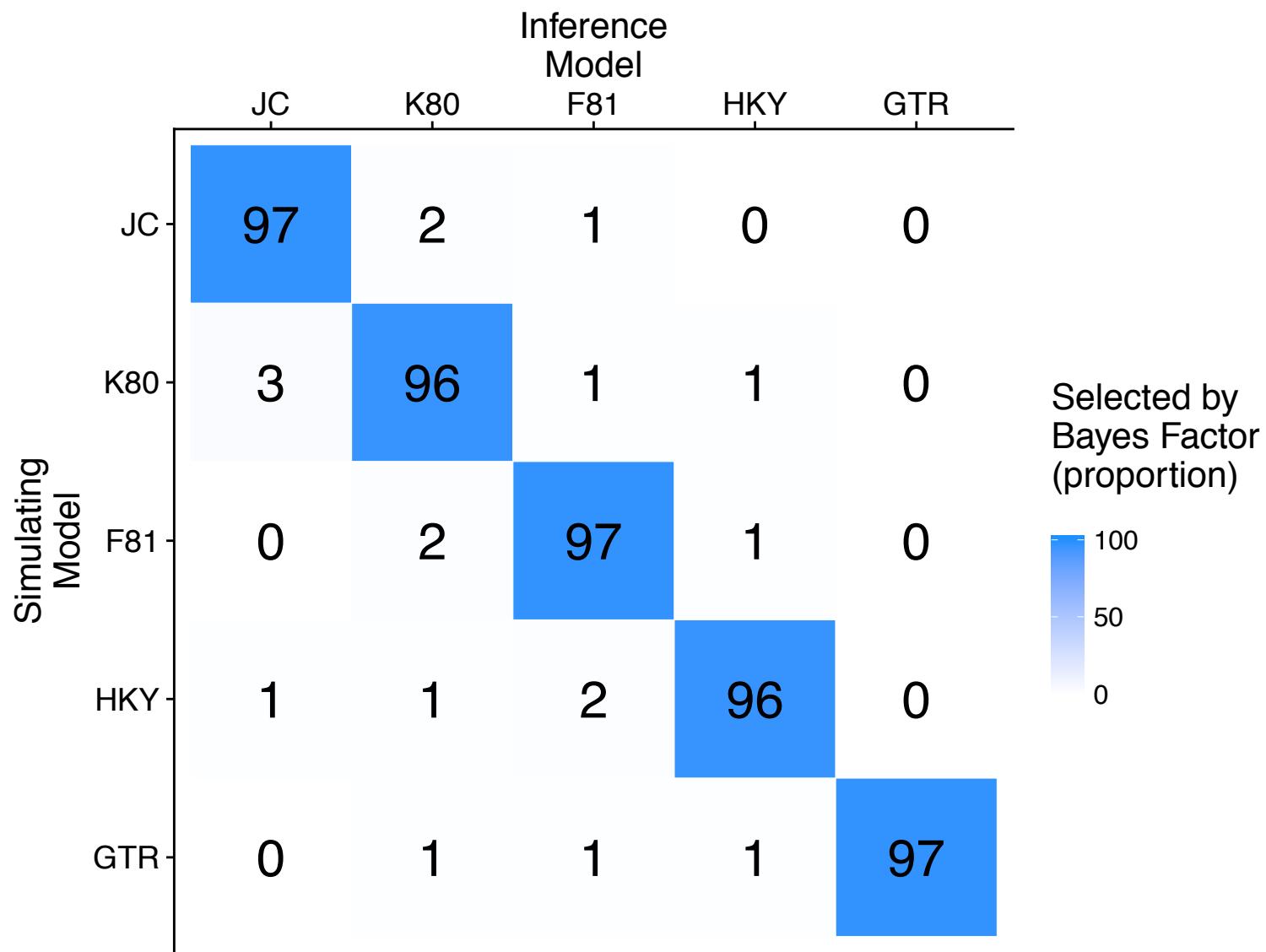
Fits
simulated
data

Inference using
the wrong model

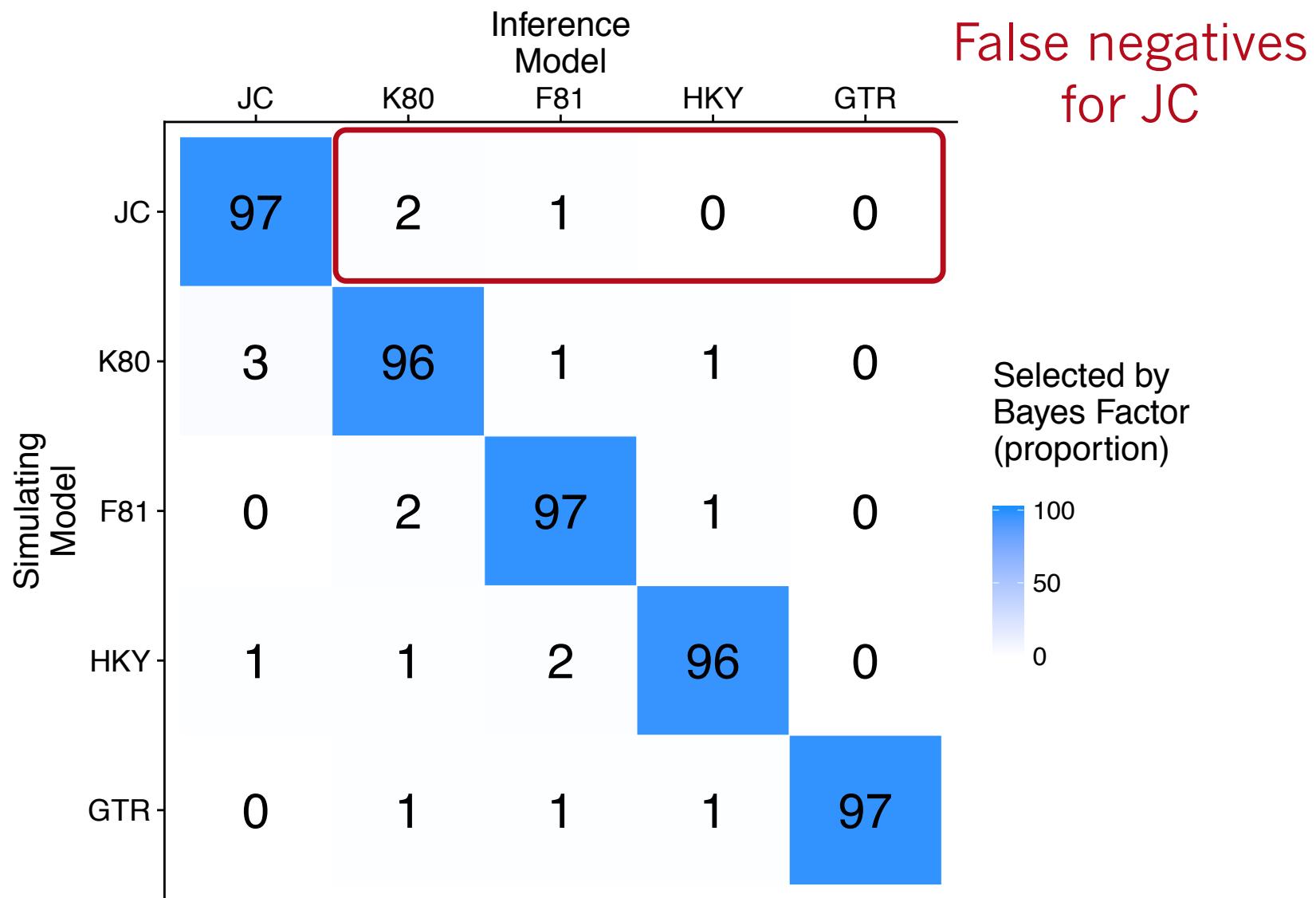
All possible
processes

Natural
processes

Model selection

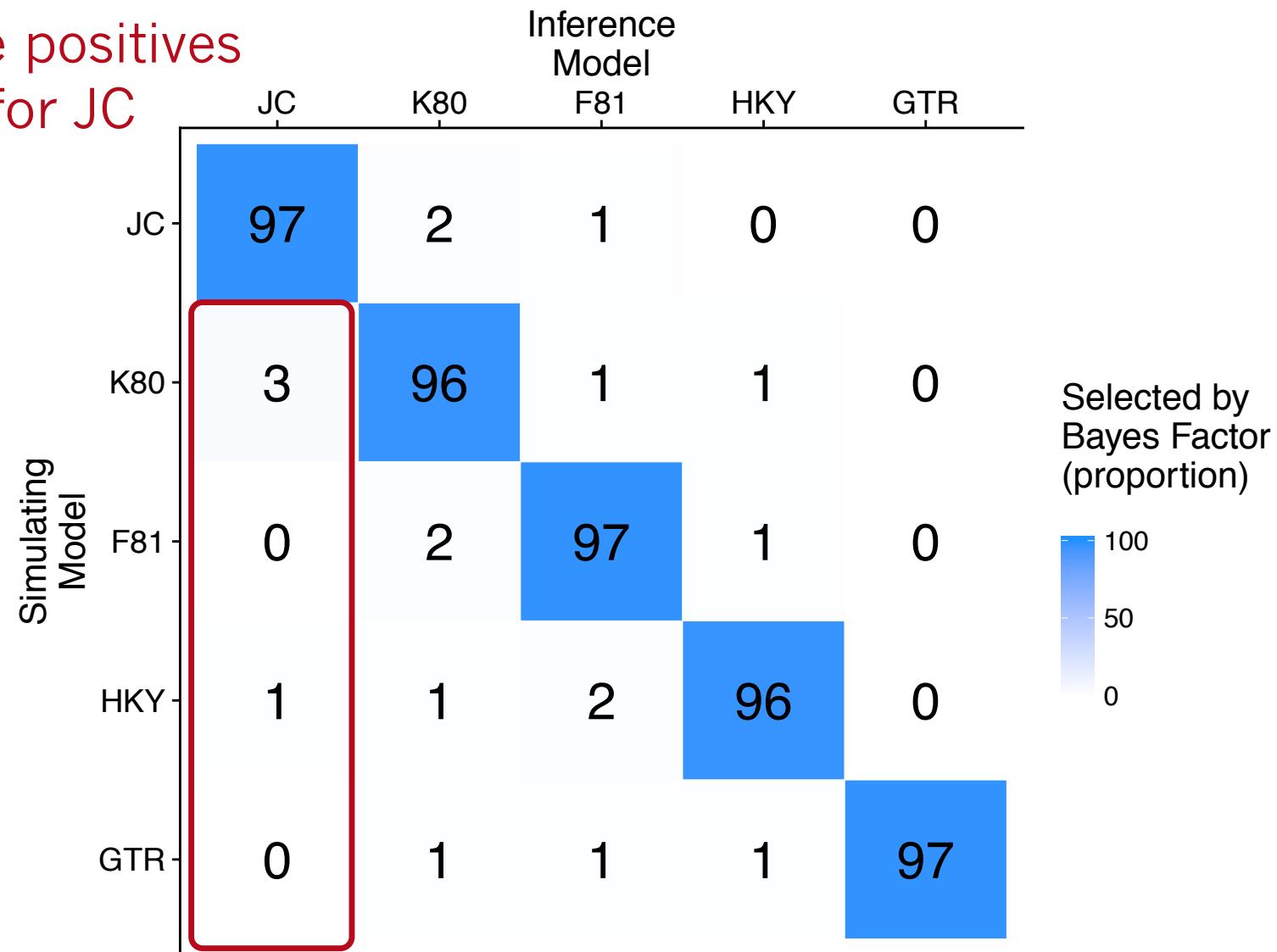


Model selection



Model selection

False positives
for JC



Model selection

10000 sites,
distinct parameters

	Inference Model				
	JC	K80	F81	HKY	GTR
JC	97	2	1	0	0
K80	3	96	1	1	0
F81	0	2	97	1	0
HKY	1	1	2	96	0
GTR	0	1	1	1	97

100 sites,
JC-like parameters

	Inference Model				
	JC	K80	F81	HKY	GTR
JC	100	0	0	0	0
K80	19	79	1	1	0
F81	20	1	78	1	0
HKY	12	15	13	49	1
GTR	9	17	20	23	22

need more data
to detect faint signal!

Inference
models

Simulating
models

Definable
models

Fits
simulated
data

Generates
more realistic
simulated
data

Generates
biological
data

Inference using
the wrong model

All possible
processes

Natural
processes

Model violations

Violation examples:

site evolution
depends on
protein structure
and function

paralogy errors
from bioinformatics
pipeline

		Inference Model				
		JC	K80	F81	HKY	GTR
Simulating Model	JC +viol.	11	19	7	5	58
	K80 +viol.	0	14	11	6	69
	F81 +viol.	0	1	31	21	48
	HKY +viol.	0	0	4	55	41
	GTR +viol.	0	0	0	0	100

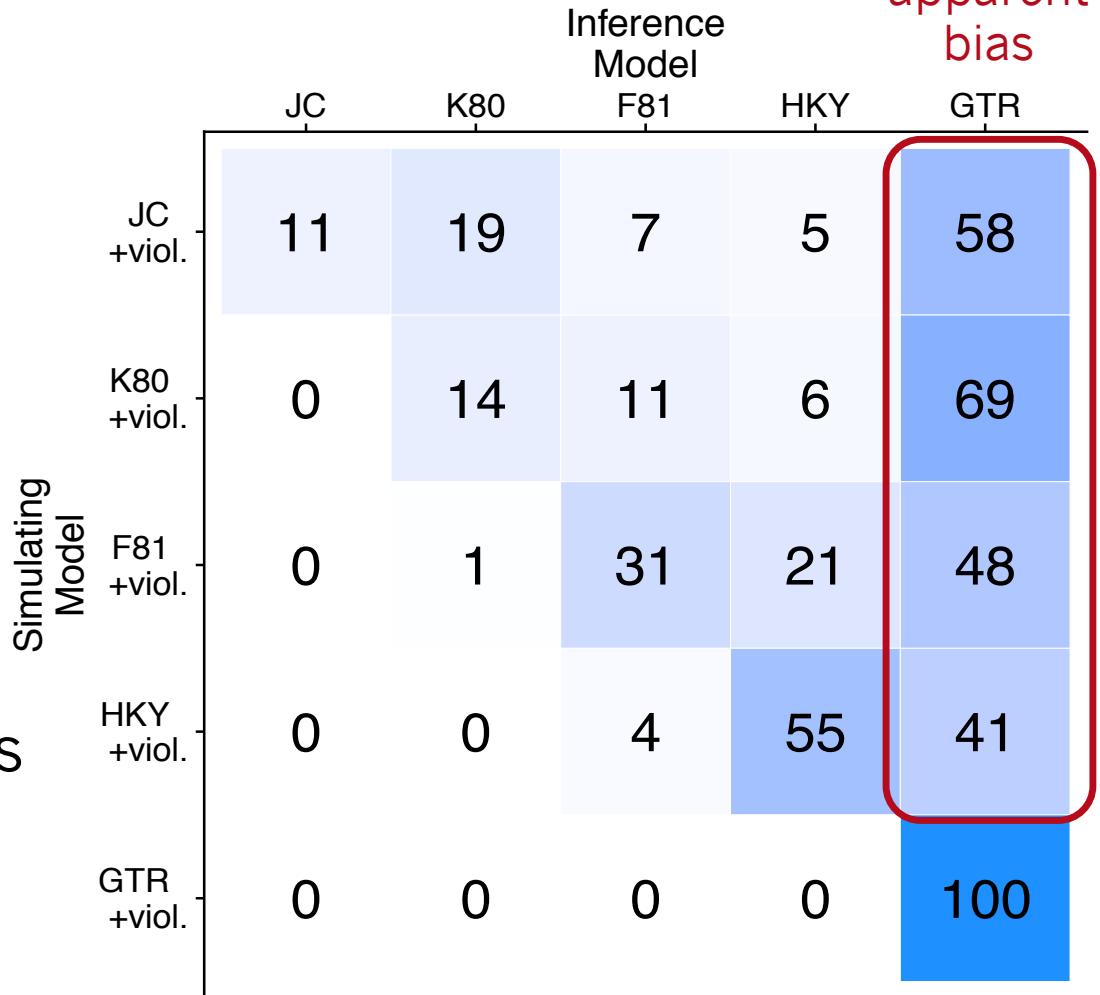
Model violations

Violation causes apparent bias

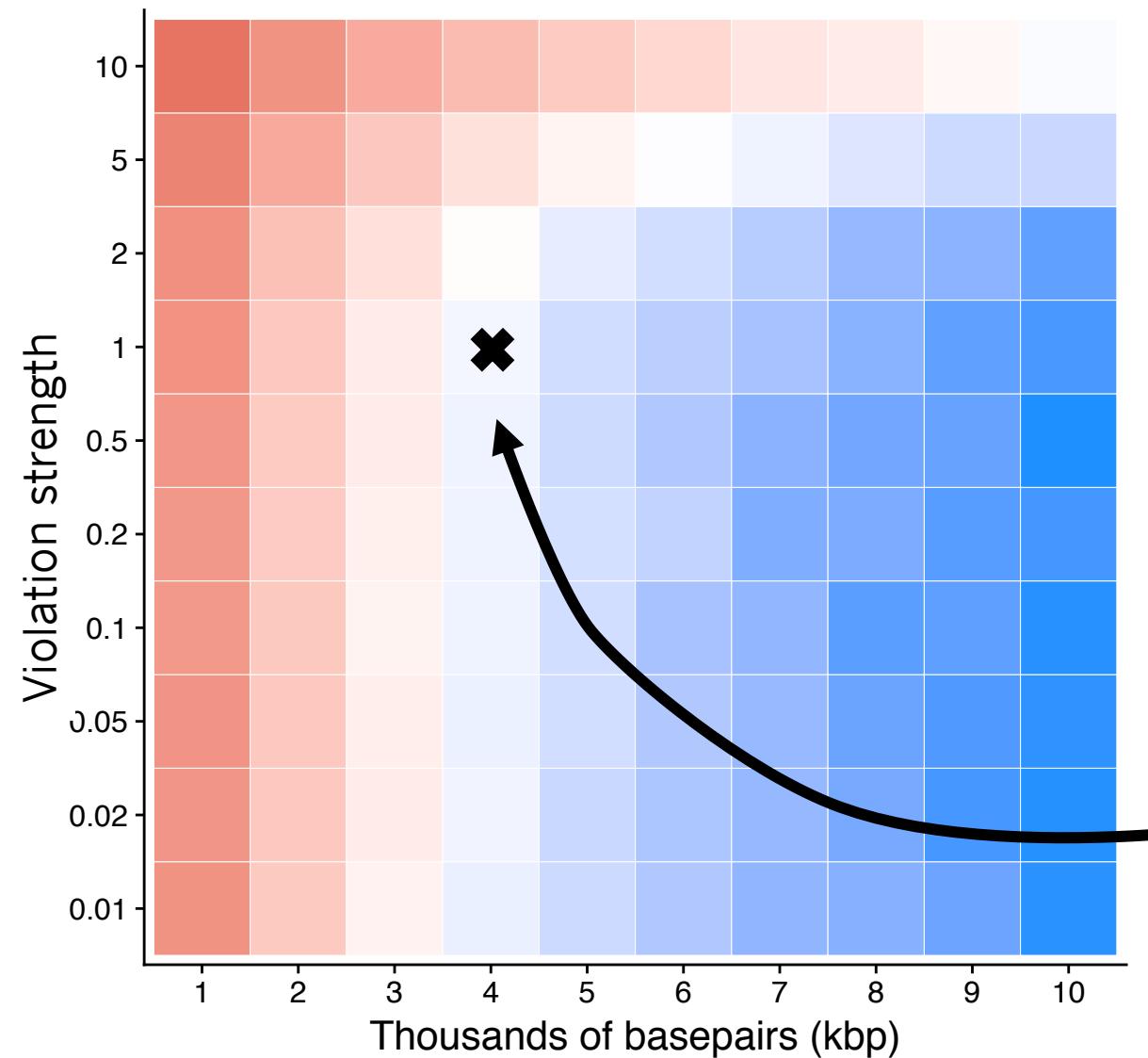
Violation examples:

site evolution depends on protein structure and function

paralogy errors from bioinformatics pipeline



Selection and power



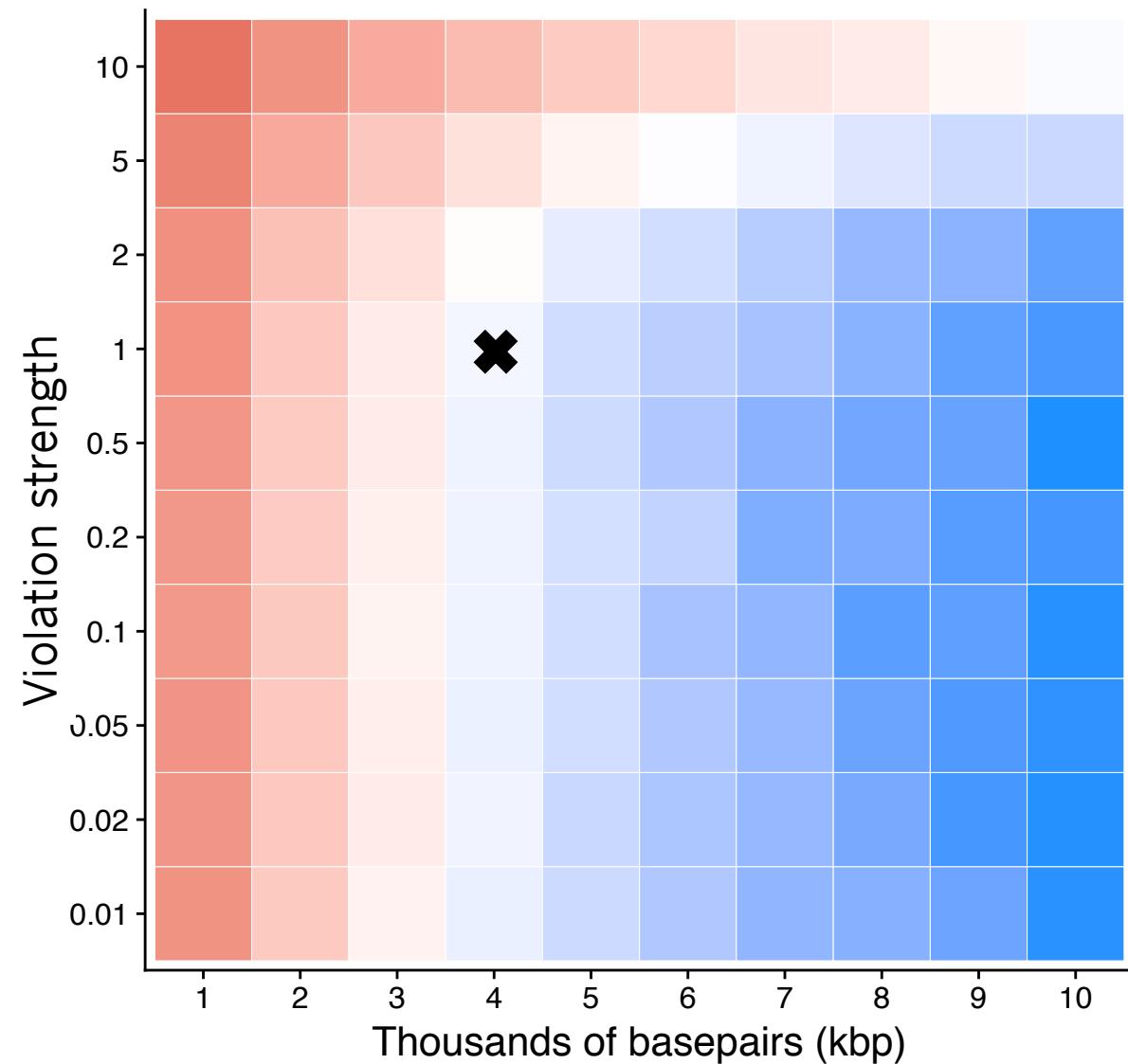
Under what conditions do we reliably select the true model?

correct model

incorrect model

False positive rate too large when violation strength > 1

Selection and power



Under what conditions do we reliably select the true model?

correct model

incorrect model

5 sim. model

x 5 inf. model

x 10 cond. 1

x 10 cond. 2

x 100 replicates

250,000 analyses

Inference
models

Fits
biological
data

Simulating
models

Definable
models

All possible
processes

Natural
processes

Generates
biological
data

Model adequacy with posterior prediction

Basic idea

1. Ask model to simulate new datasets using the parameters that best explain your biological data
2. Do the simulated data resemble your biological data?

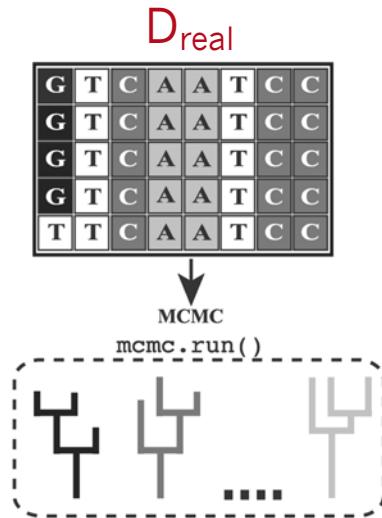
Model adequacy with posterior prediction

D_{real}

G	T	C	A	A	T	C	C
G	T	C	A	A	T	C	C
G	T	C	A	A	T	C	C
G	T	C	A	A	T	C	C
T	T	C	A	A	T	C	C

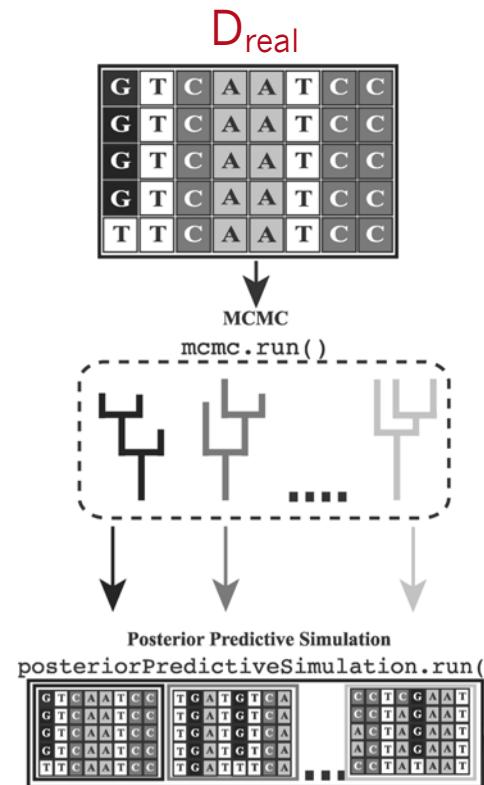
Model adequacy with posterior prediction

1. Estimate
 $\Pr(\theta_{\text{real}} | D_{\text{real}})$



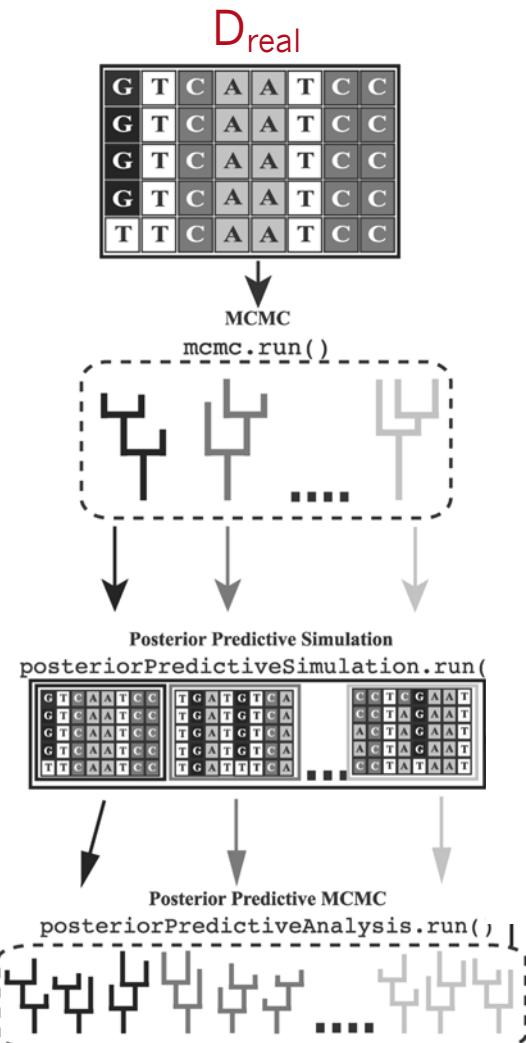
Model adequacy with posterior prediction

1. Estimate
 $\Pr(\theta_{\text{real}} | D_{\text{real}})$
2. Simulate K datasets
 $D_{\text{sim},k} \sim f(\theta_{\text{real}} | D_{\text{real}})$



Model adequacy with posterior prediction

1. Estimate
 $\Pr(\theta_{\text{real}} | D_{\text{real}})$
2. Simulate K datasets
 $D_{\text{sim},k} \sim f(\theta_{\text{real}} | D_{\text{real}})$
3. Estimate K posteriors
 $\Pr(\theta_{\text{sim},k} | D_{\text{sim},k})$



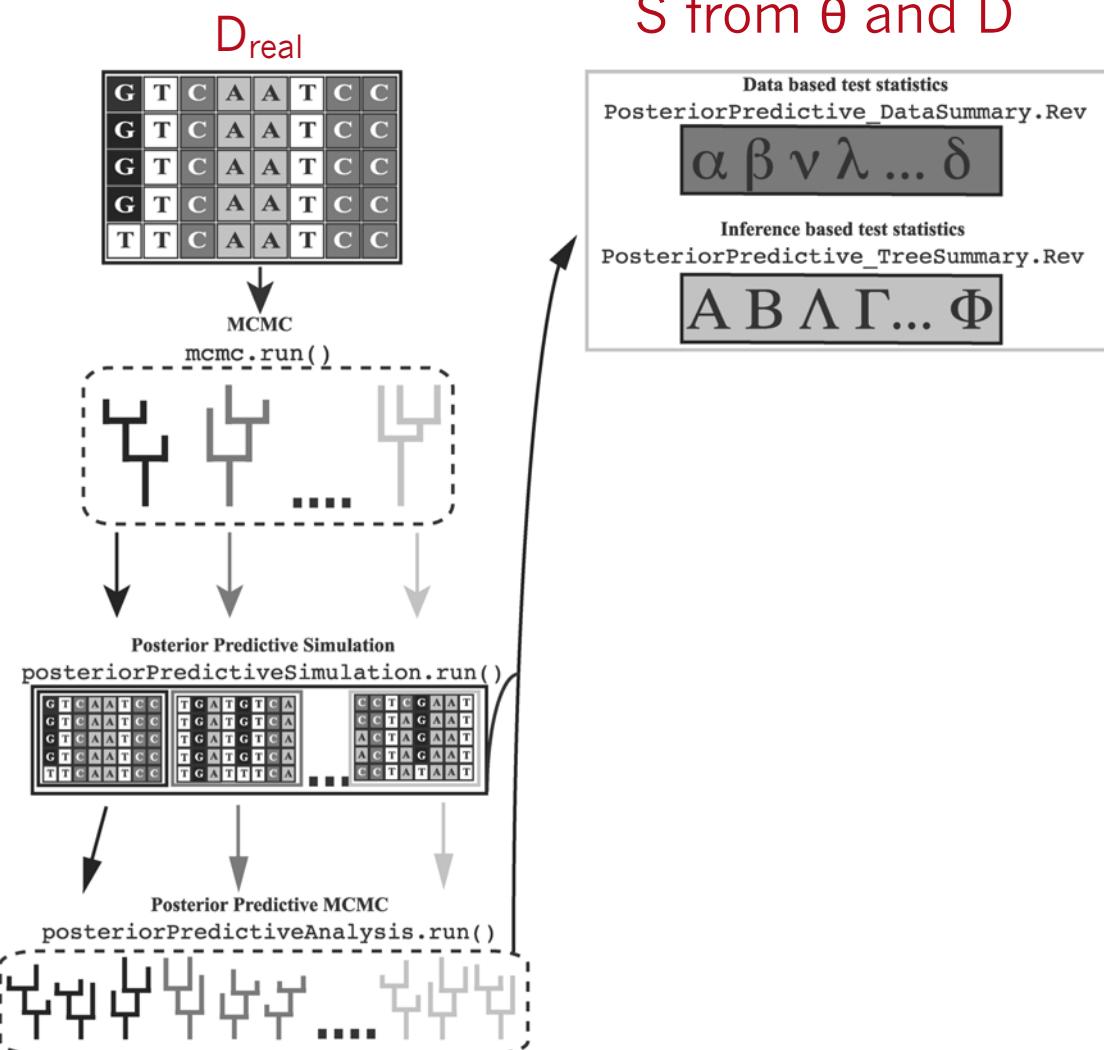
Model adequacy with posterior prediction

1. Estimate
 $\Pr(\theta_{\text{real}} | D_{\text{real}})$

2. Simulate K datasets
 $D_{\text{sim},k} \sim f(\theta_{\text{real}} | D_{\text{real}})$

3. Estimate K posteriors
 $\Pr(\theta_{\text{sim},k} | D_{\text{sim},k})$

4. Collect test statistics
 S from θ and D

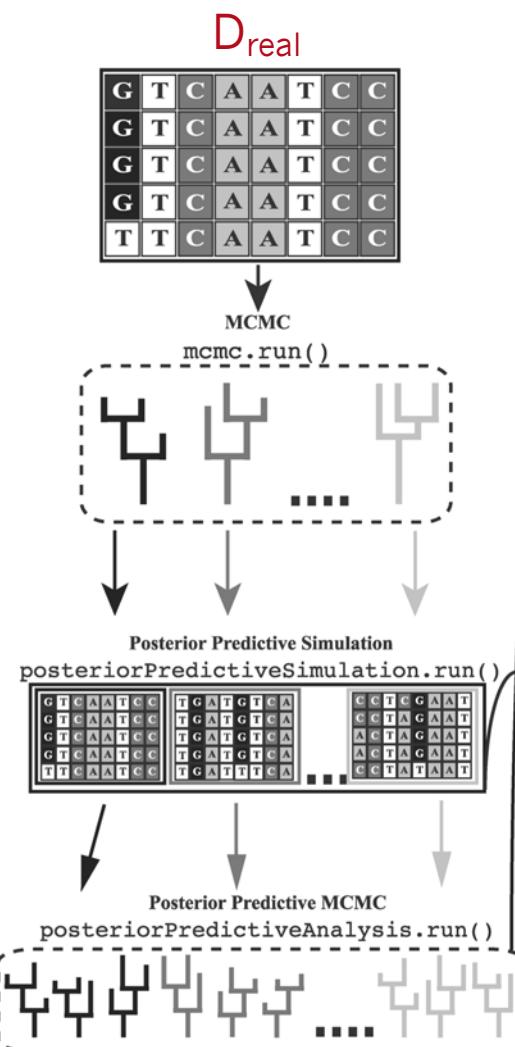


Model adequacy with posterior prediction

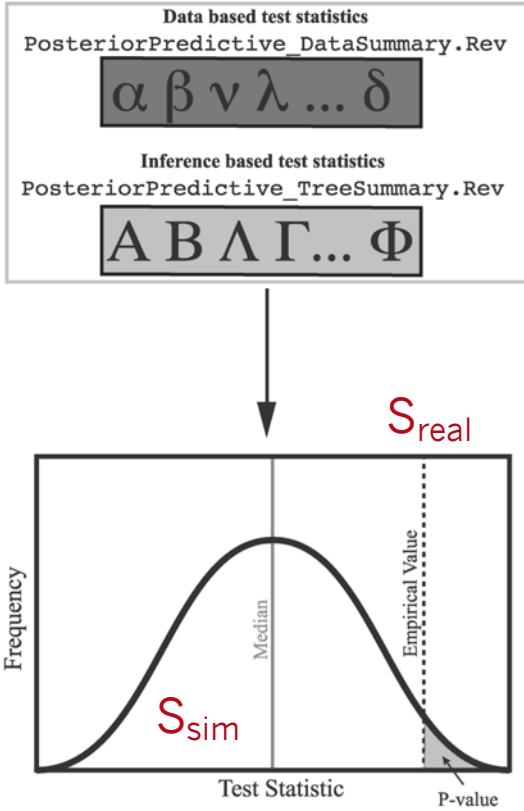
1. Estimate
 $\Pr(\theta_{\text{real}} | D_{\text{real}})$

2. Simulate K datasets
 $D_{\text{sim},k} \sim f(\theta_{\text{real}} | D_{\text{real}})$

3. Estimate K posteriors
 $\Pr(\theta_{\text{sim},k} | D_{\text{sim},k})$

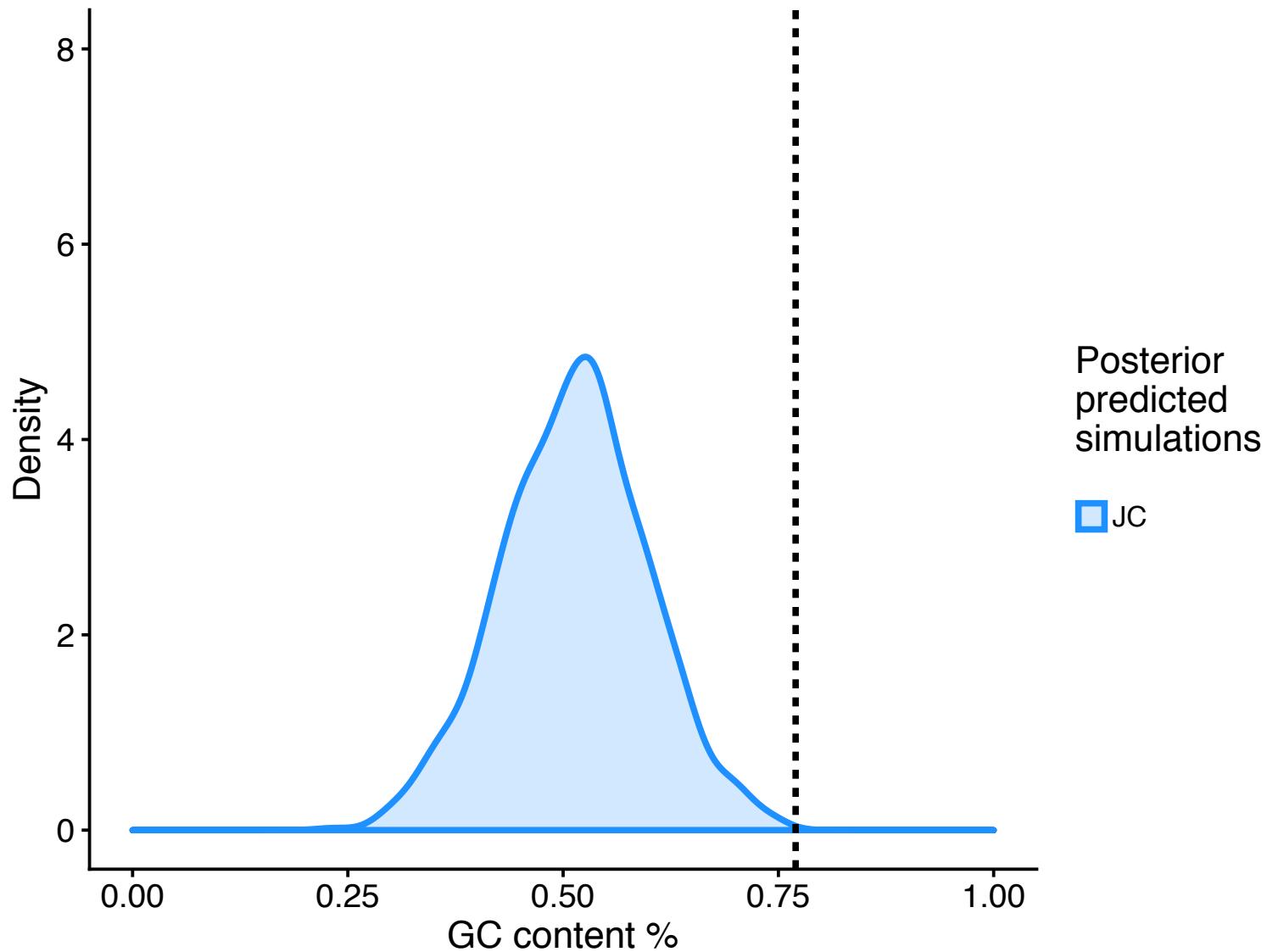


4. Collect test statistics S from θ and D

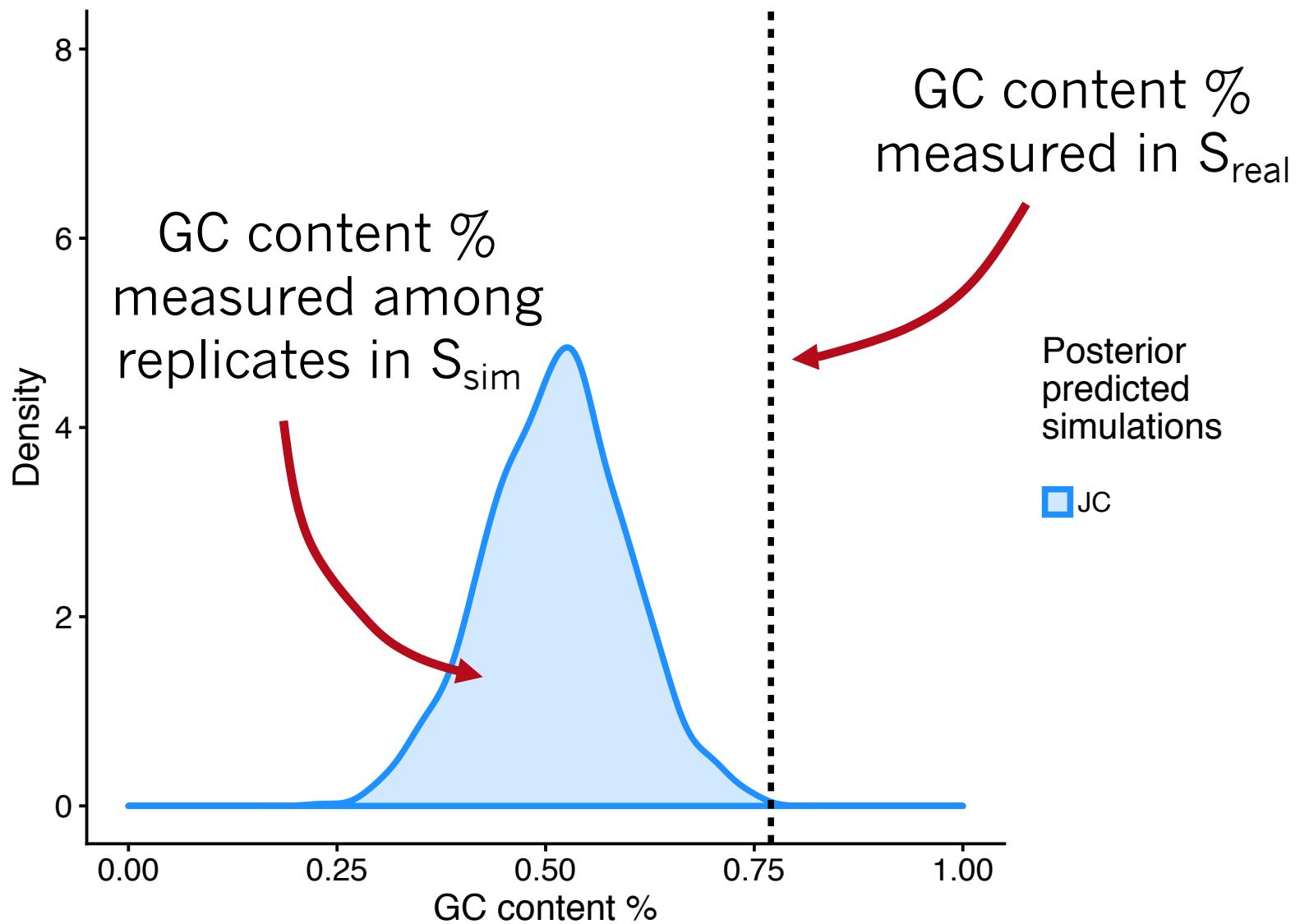


5. Does distribution of S_{sim} contain S_{real} ?

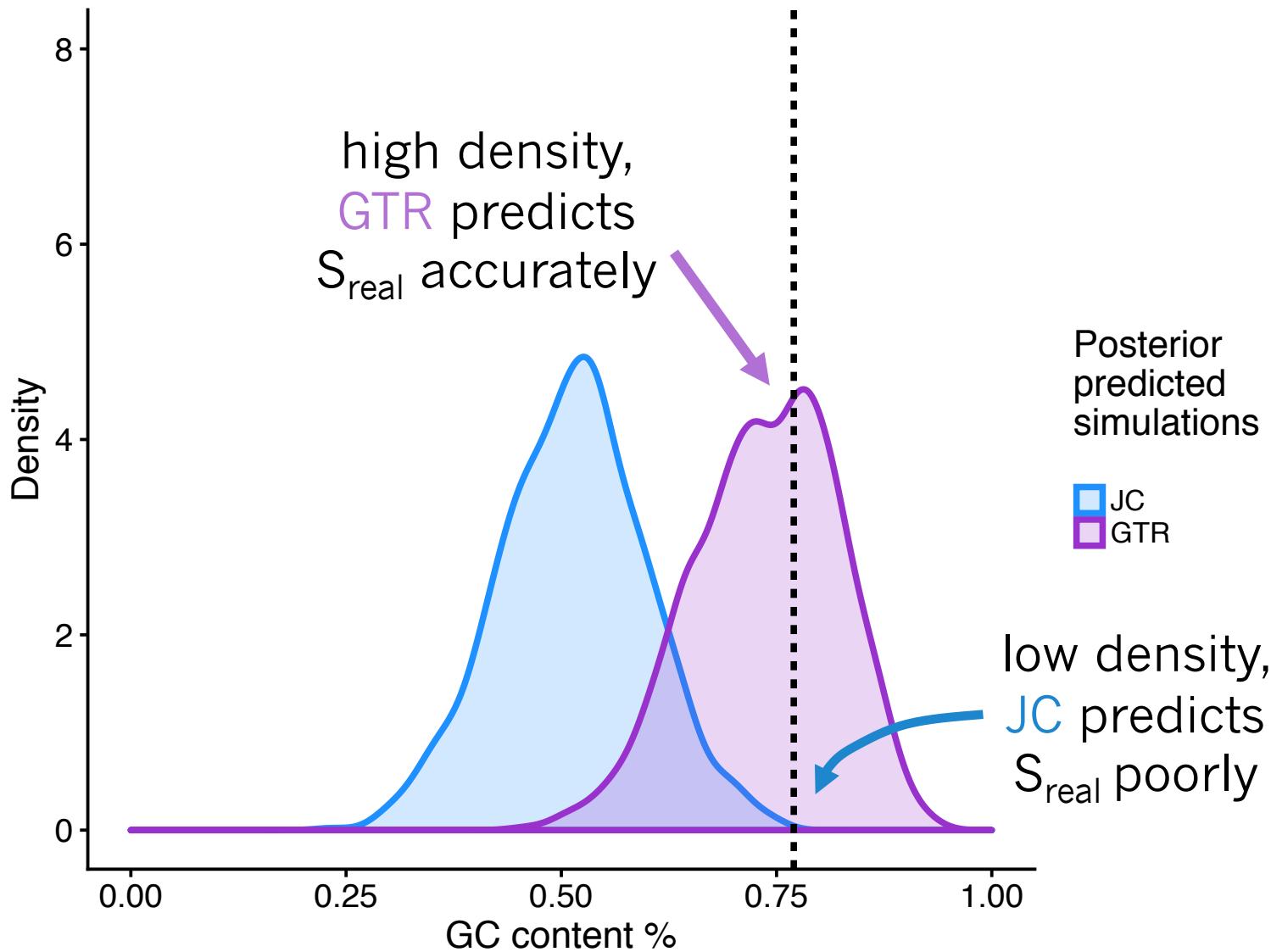
Model adequacy



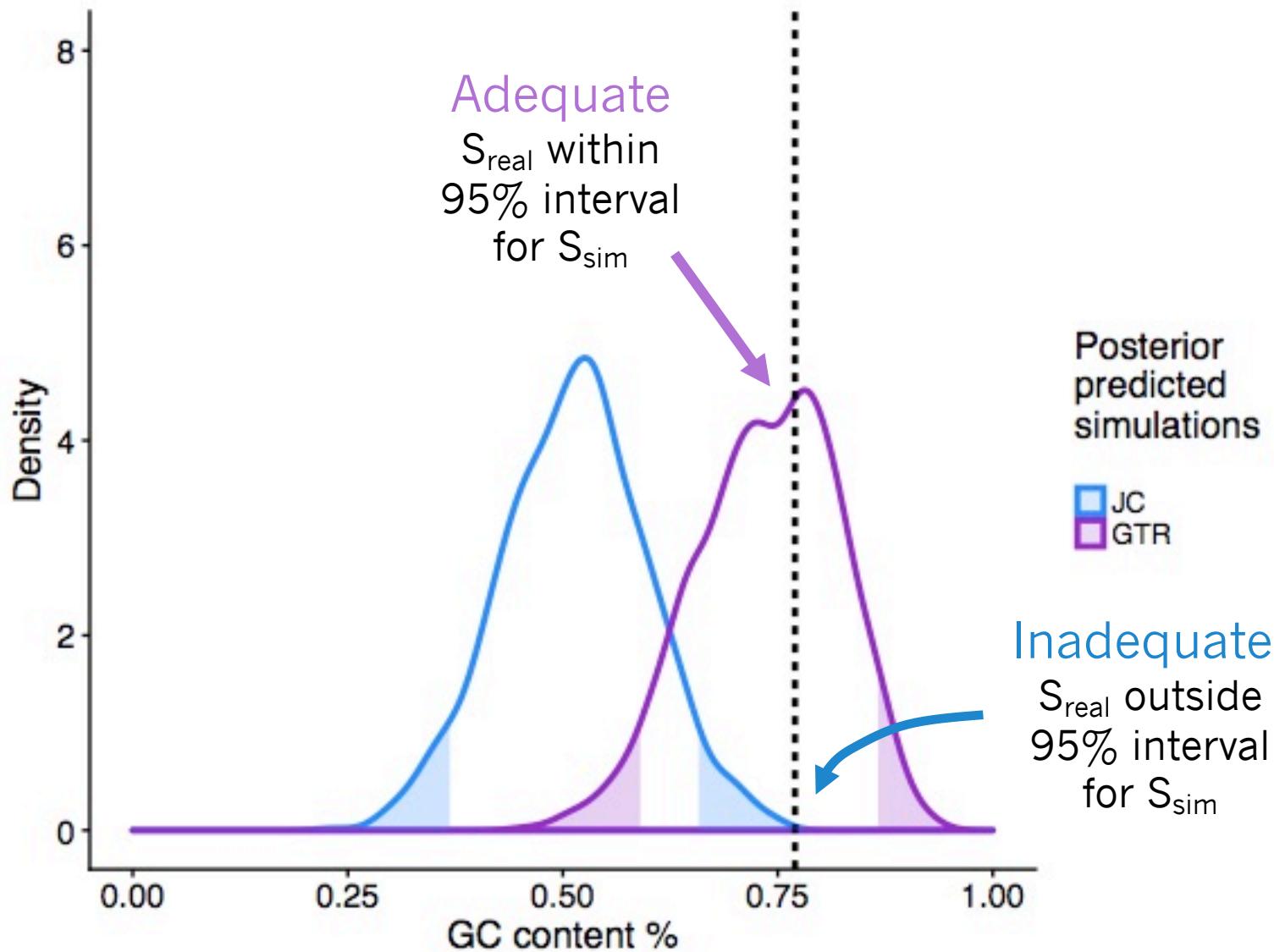
Model adequacy



Model adequacy



Model adequacy



Model smell

Do results make sense?

Time runs backwards?

Precambrian rabbits?

Model smell

Do results make sense?

Time runs backwards?

Precambrian rabbits?

If not, slow down!

Errors in data?

Analysis misconfigured?

Method glitch?

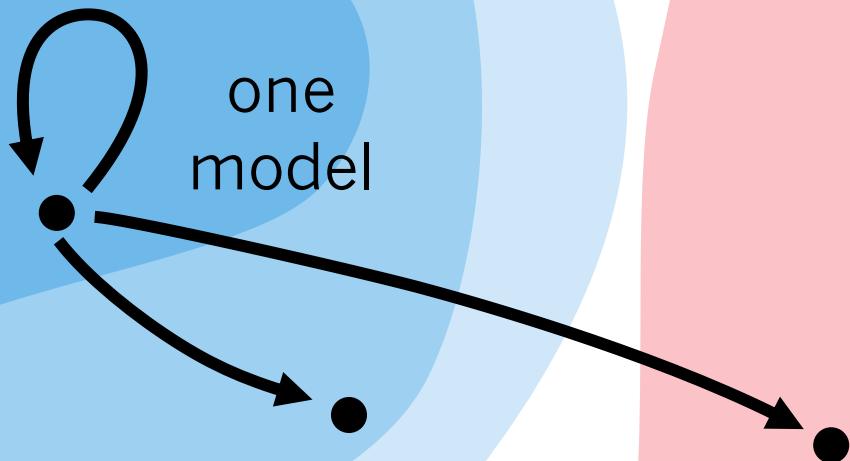
Inference
models

Simulating
models

Definable
models

All possible
processes

Natural
processes



Inference
models

a better
model?

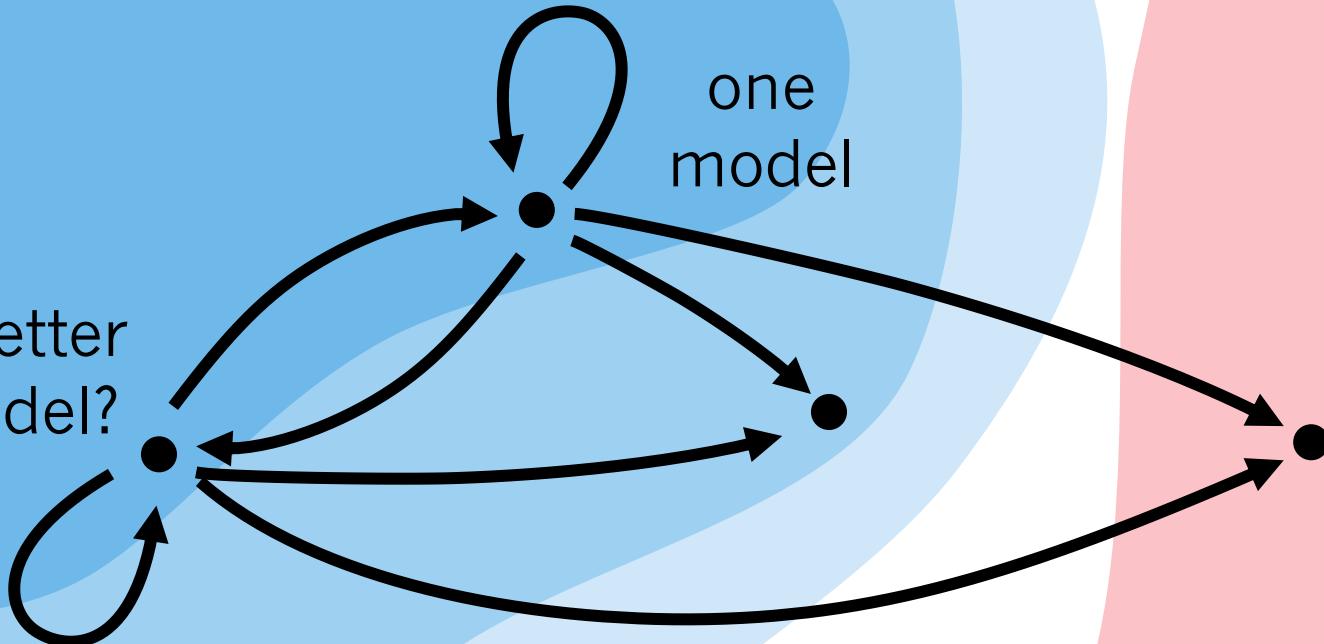
Simulating
models

Definable
models

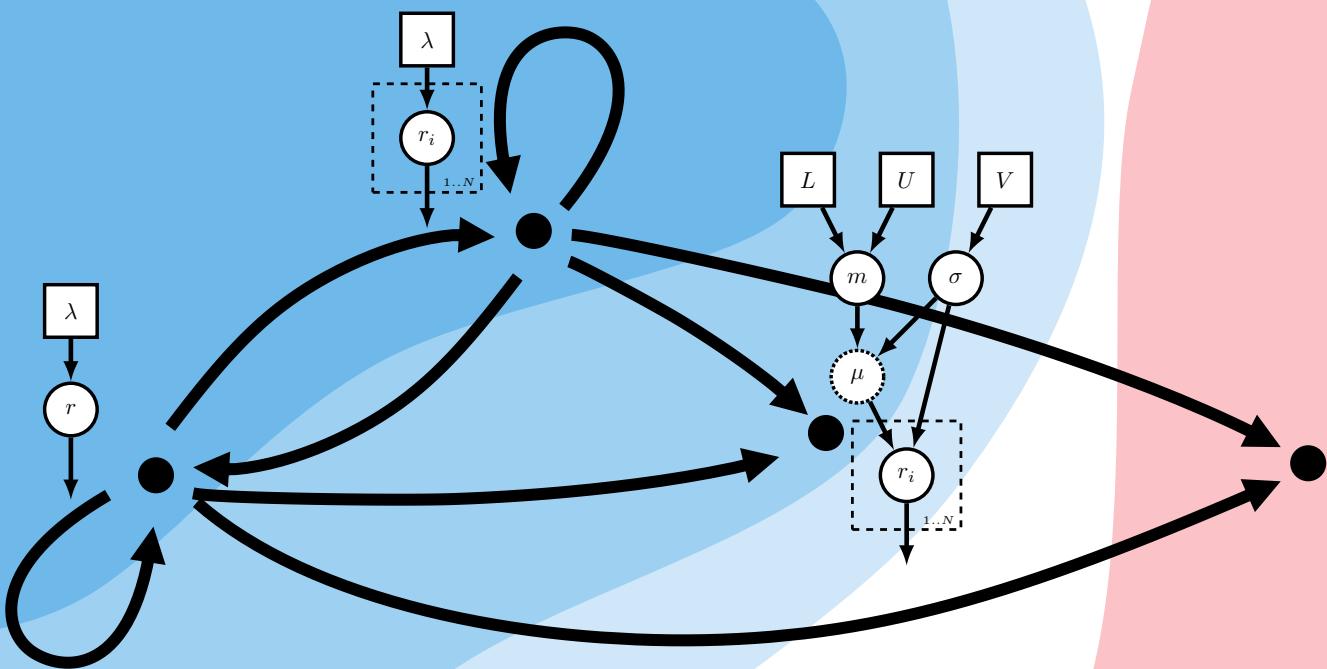
All possible
processes

one
model

Natural
processes



Inference
models



Simulating
models

PGMs to navigate
model space
(more in lab)

Definable
models

All possible
processes

Natural
processes