

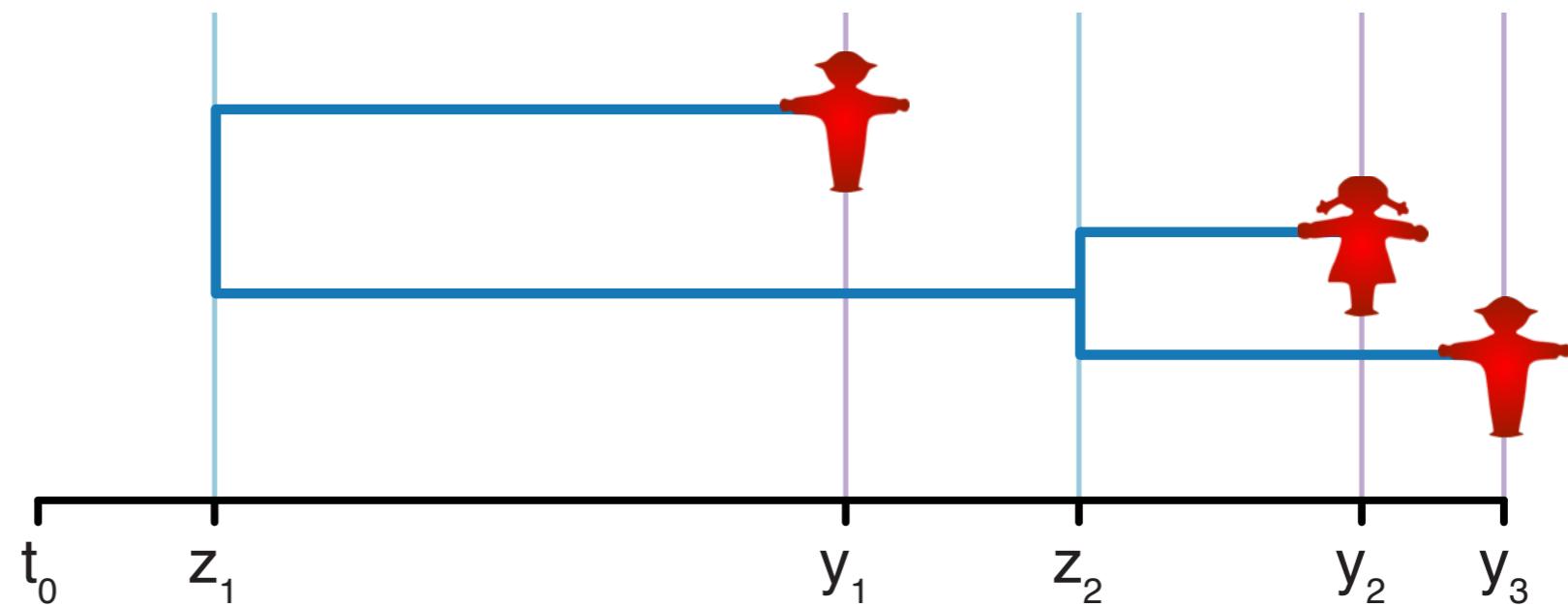


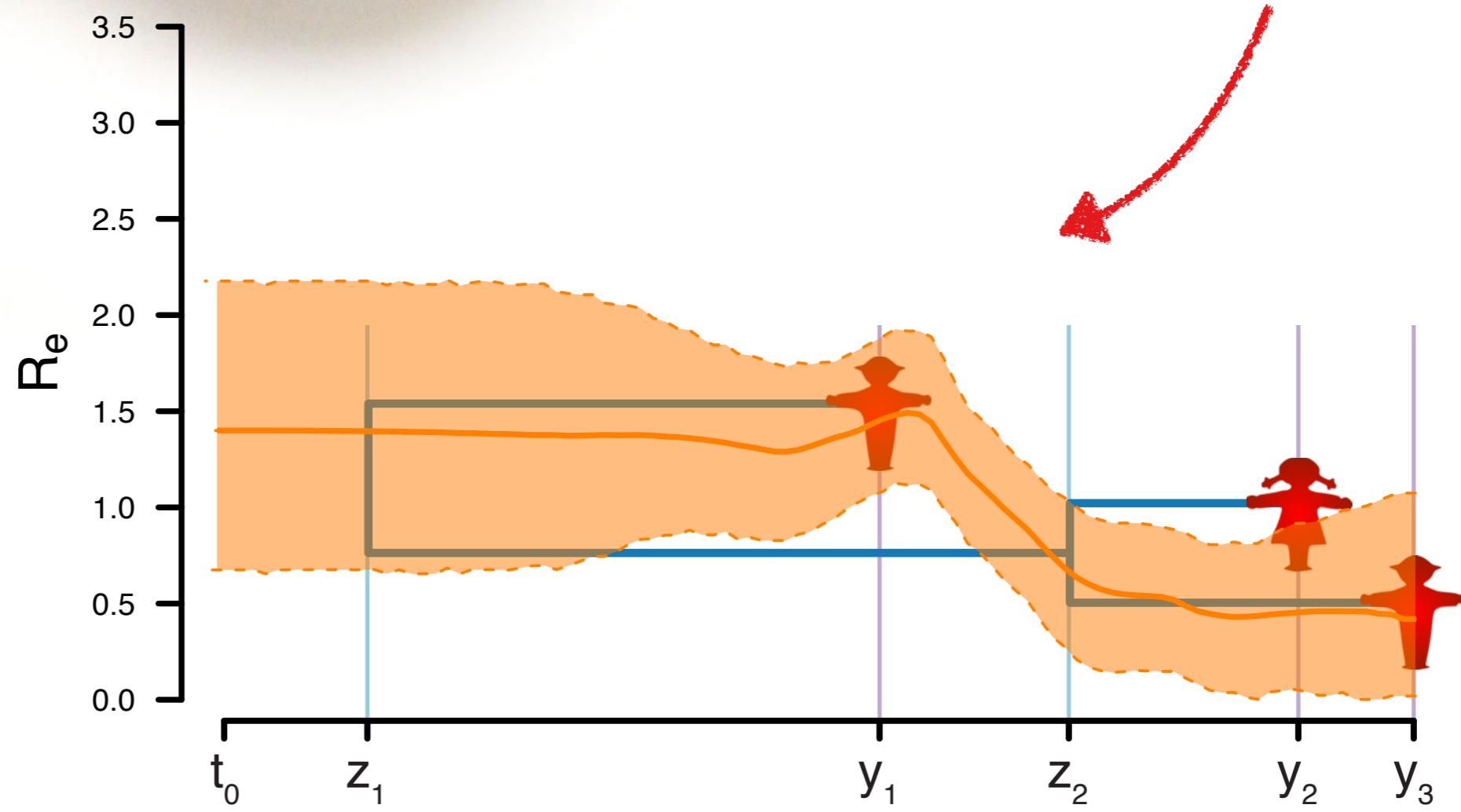
QUANTIFYING THE **CONTRIBUTION OF EXTERNAL COVARIATES** TO **PATHOGEN POPULATION DYNAMICS**
IN A BIRTH-DEATH FRAMEWORK

LOUIS DU PLESSIS



ACACACCC...
TCACACCT...
ACAGACTT...





Tree

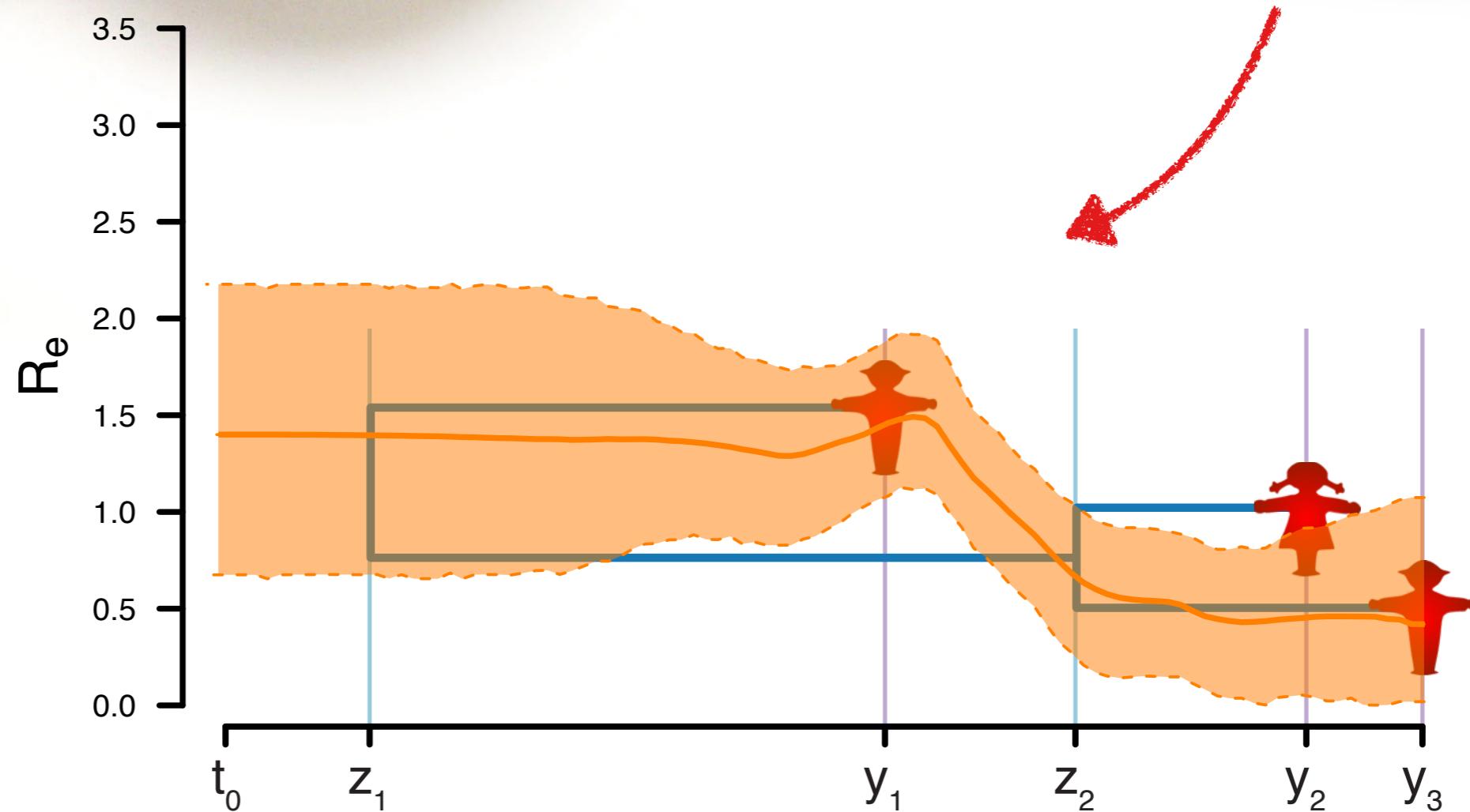
Realisation of a stochastic process

$$P(E | \text{O}_0)$$

Demographic model

Describes the growth of the tree (spread of epidemic)

ACAGAGCTT...



Tree

Realisation of a stochastic process

$$P(E | \text{O}_0)$$

Demographic model

Describes the growth of the tree (spread of epidemic)

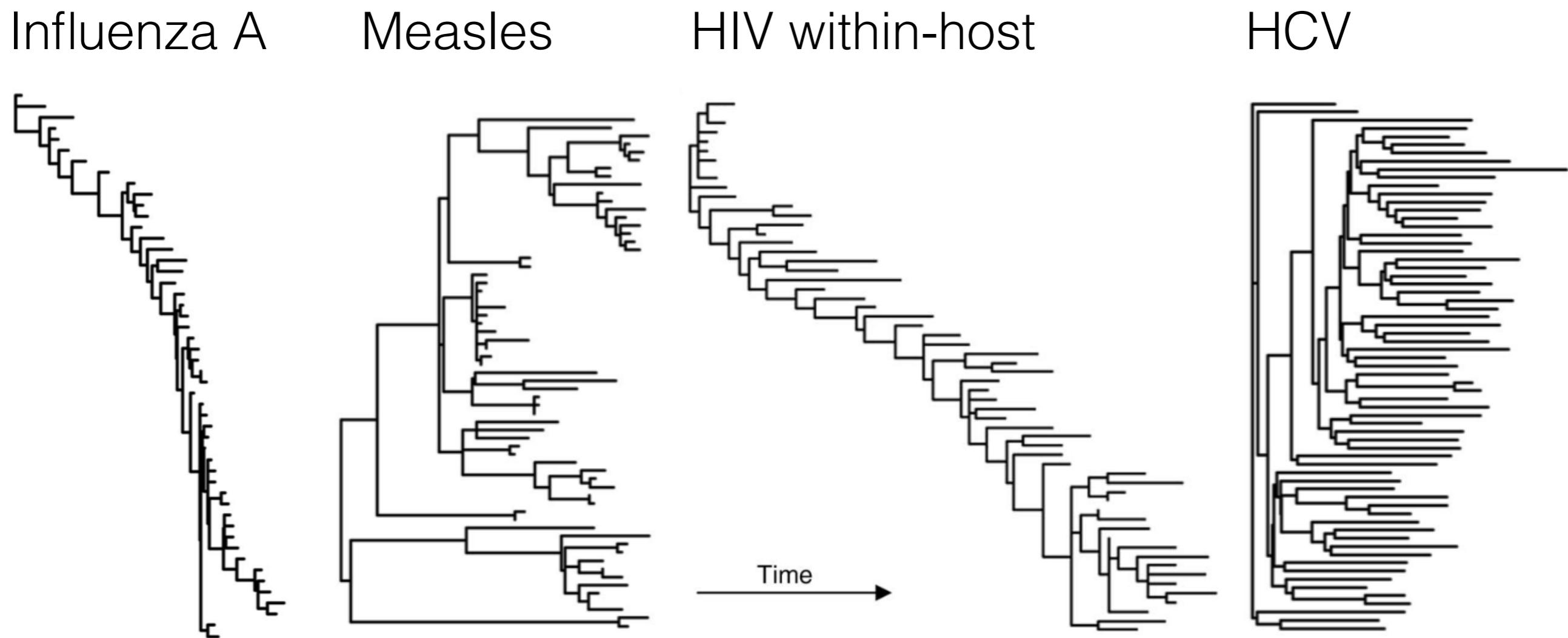
Epidemiological and **evolutionary** dynamics should occur on the same **timescale!**

Measurably evolving populations:

- Large population size
- High mutation rate
- Short generation times

This is the case for **many** (if not most) viral epidemics!

Genomes contain a **signature** of the **epidemiological** dynamics



Bayesian phylogenetic and phylodynamic inference

$$P(\text{model} \mid \text{data}) = \frac{P(\text{data} \mid \text{model})P(\text{model})}{P(\text{data})}$$

Likelihood

Posterior

Prior

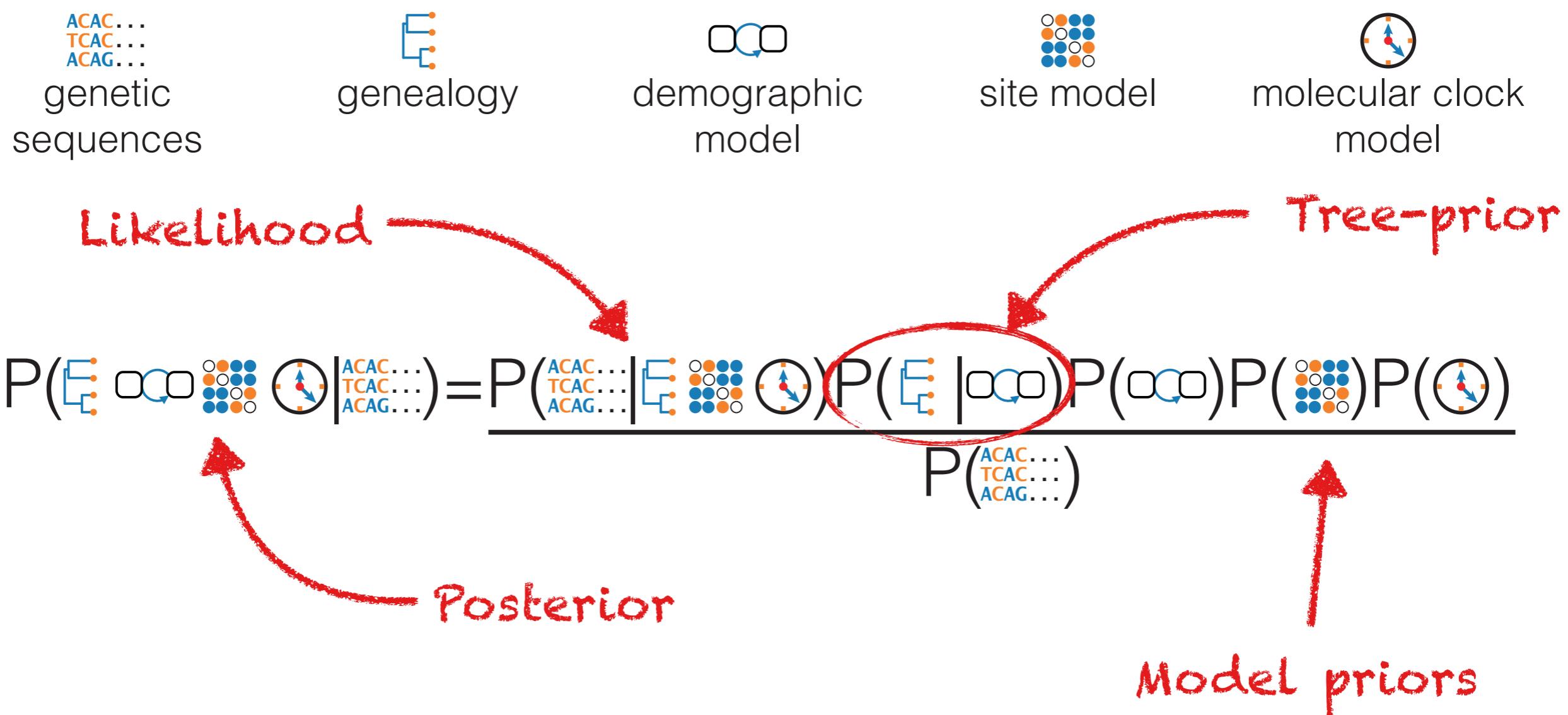
Marginal Likelihood of the data

Hierarchical Bayesian model in Beast2

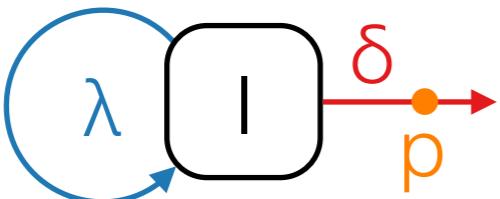


$$P(\text{genealogy} \mid \text{demographic model}, \text{site model}, \text{molecular clock model} \mid \text{genetic sequences}) = \frac{P(\text{genetic sequences} \mid \text{genealogy}, \text{demographic model}, \text{site model}, \text{molecular clock model})}{P(\text{genetic sequences})}$$

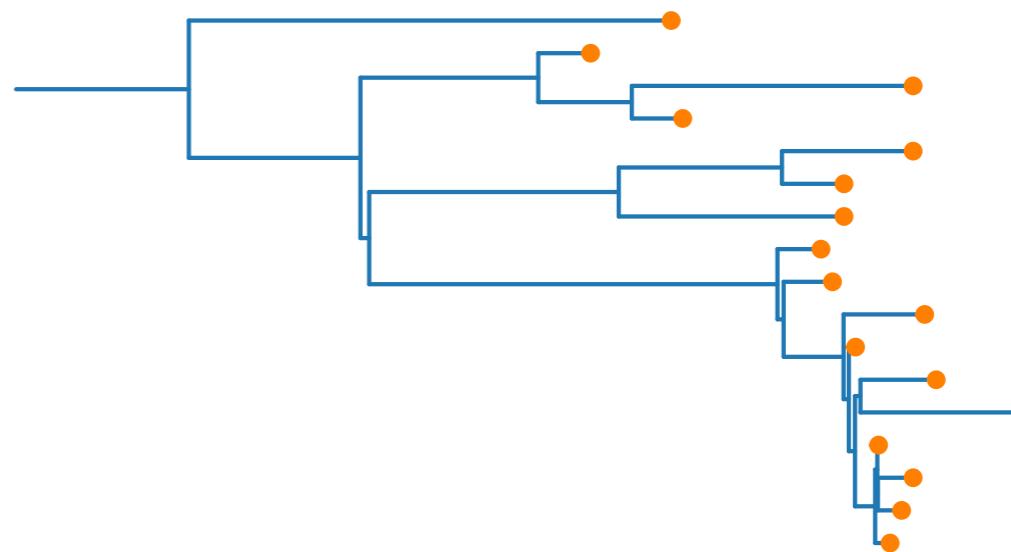
Hierarchical Bayesian model in Beast2



Birth-death skyline

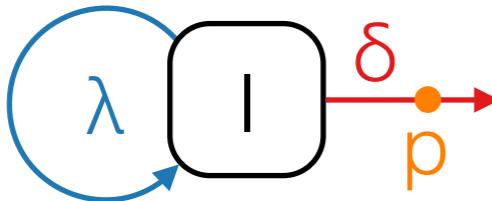


- λ — infection rate
- δ — becoming-noninfectious rate
- p — sampling proportion
- t_0 — Origin of the process

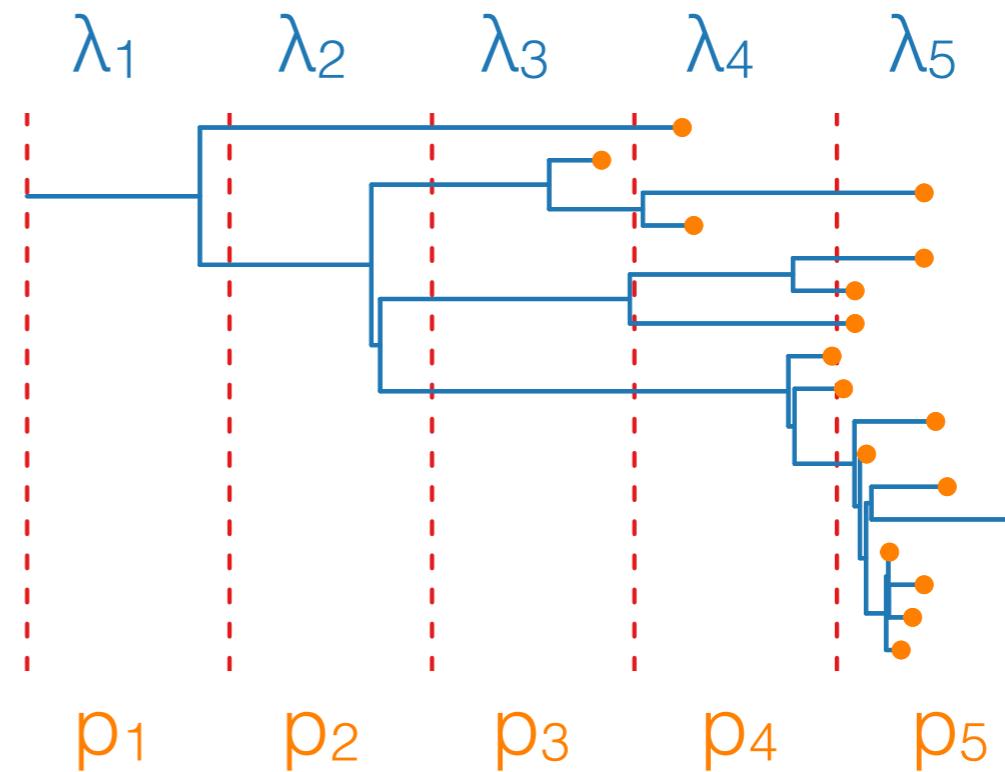


- Model parameters are rates that describe the growth of the tree
- Effective reproduction number = λ/δ
(is it spreading or not?)

Birth-death skyline

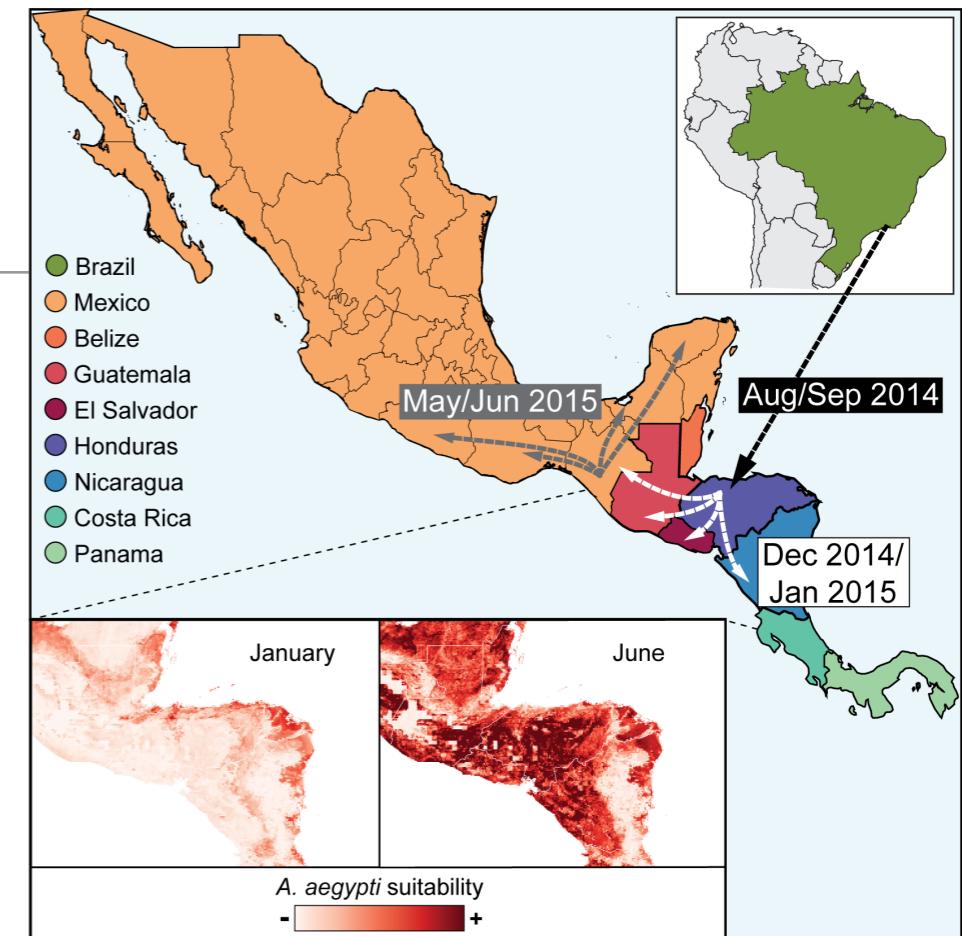
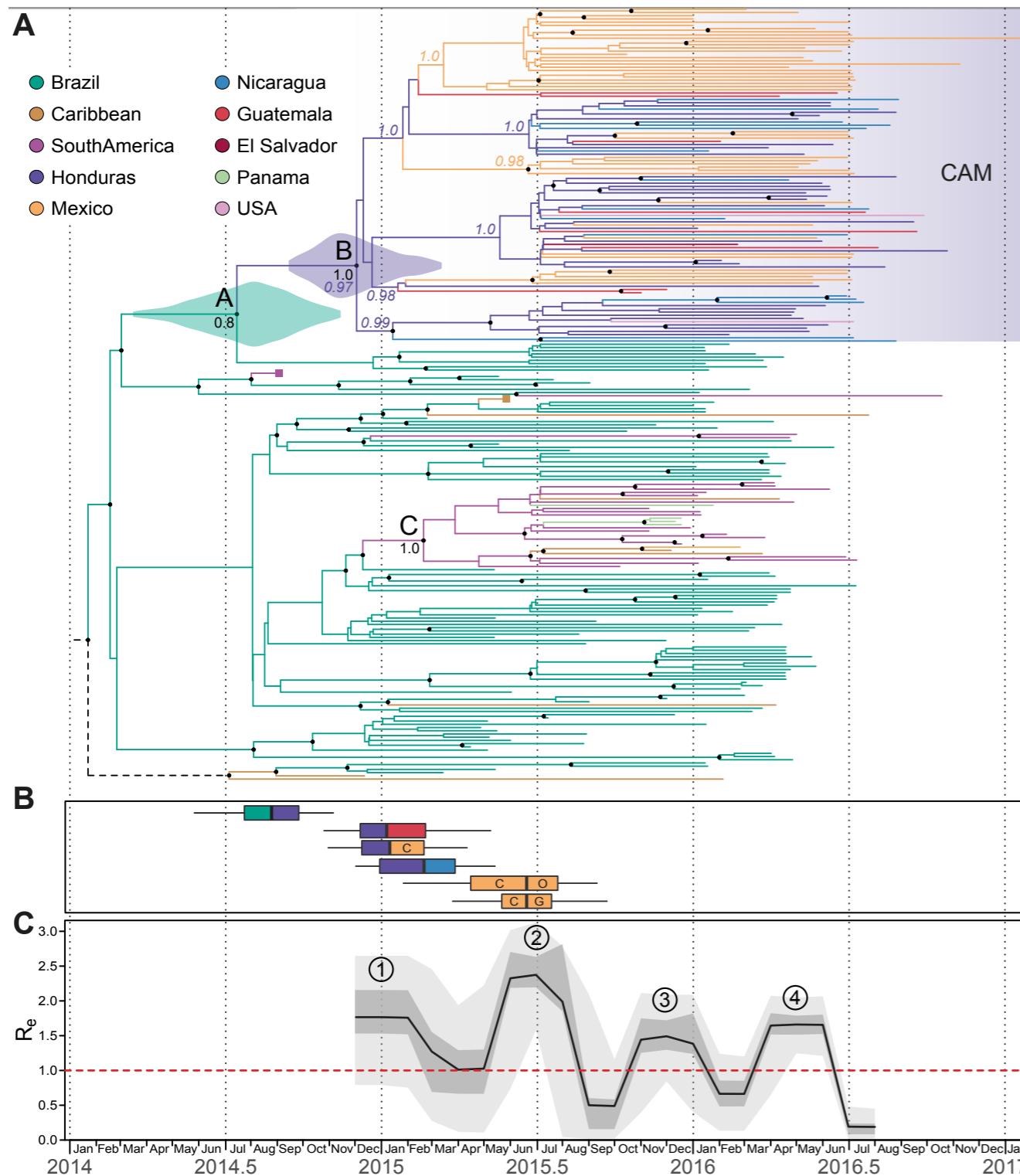


- λ — infection rate
- δ — becoming-noninfectious rate
- p — sampling proportion
- t_0 — Origin of the process



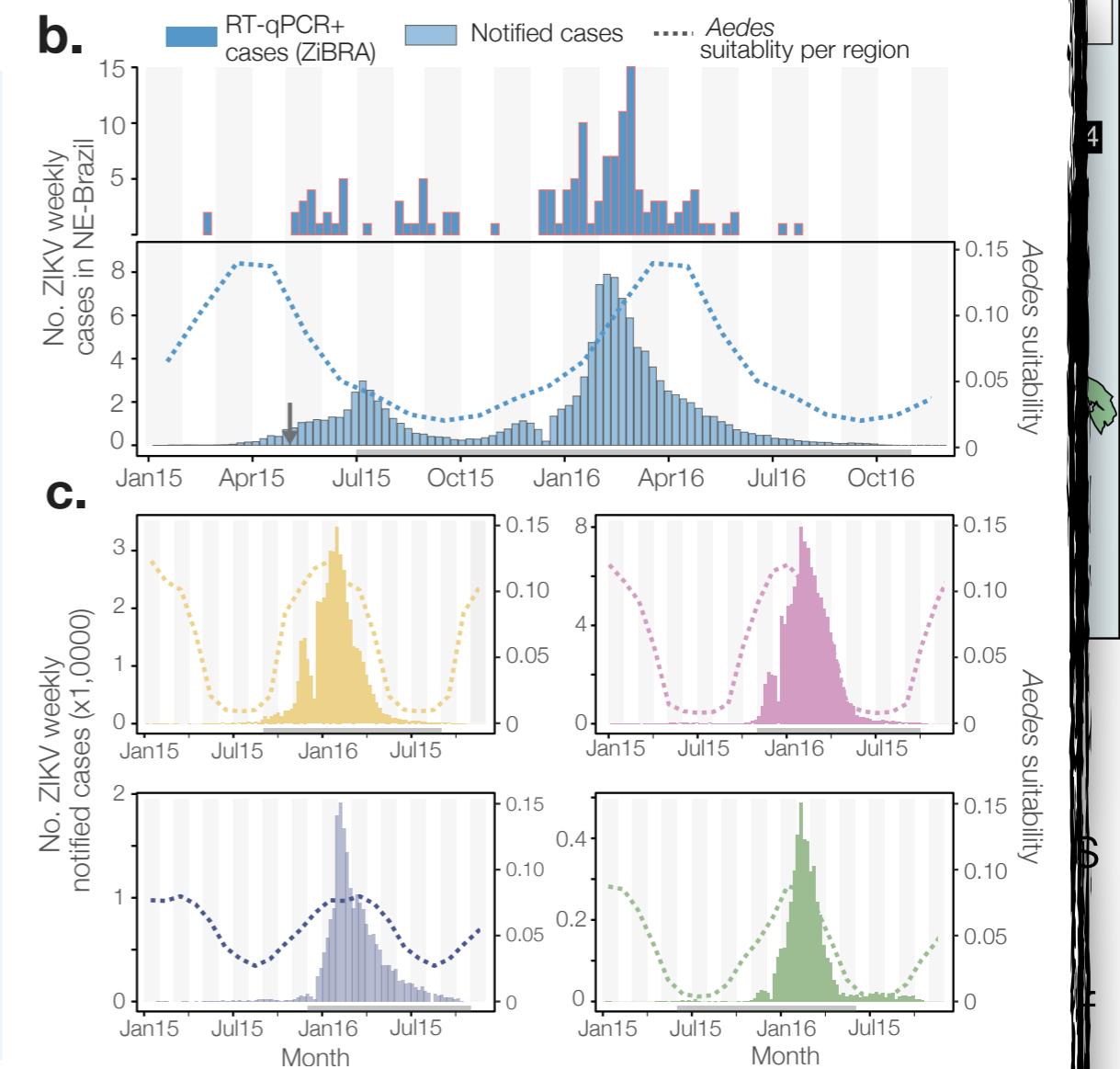
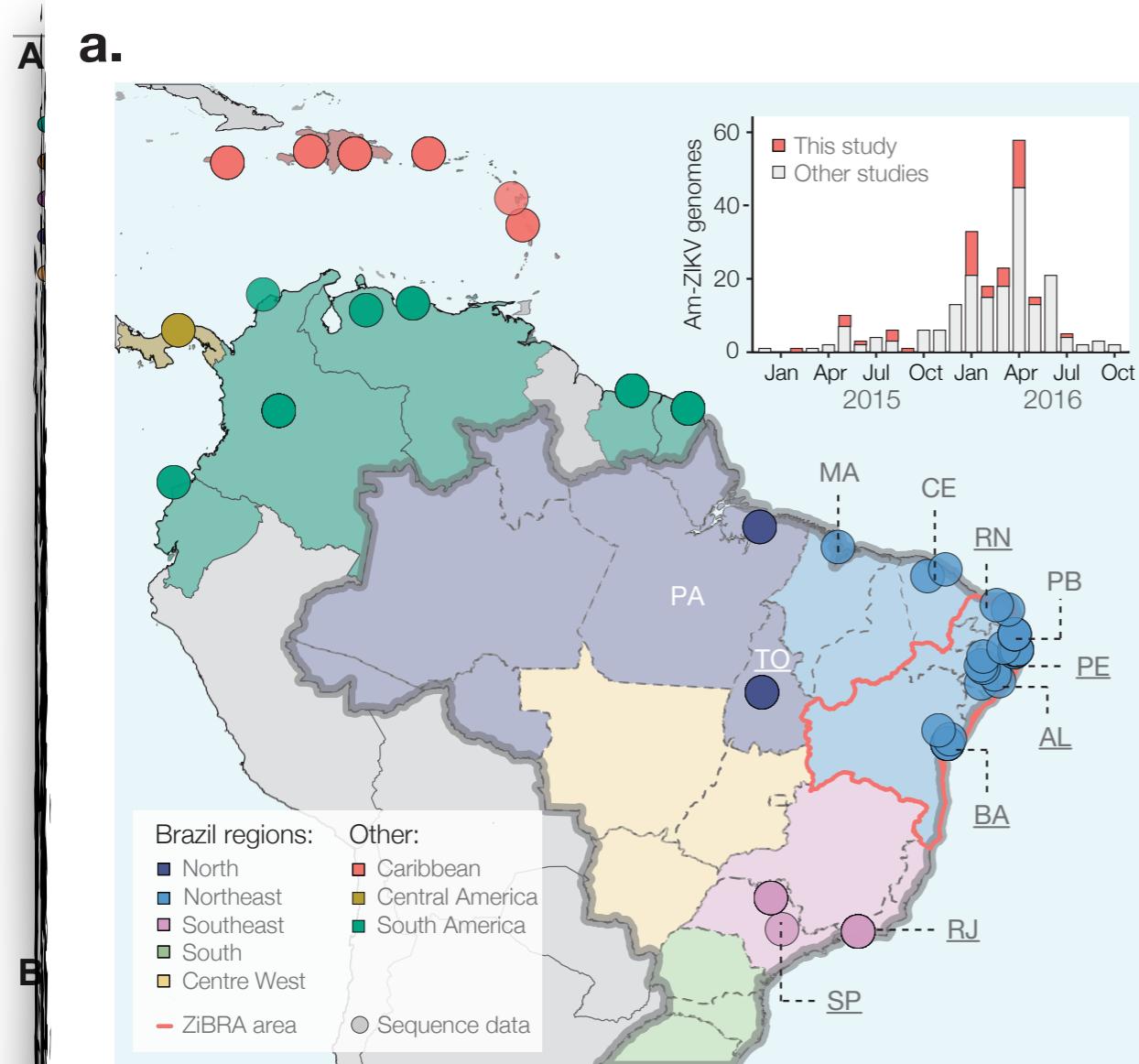
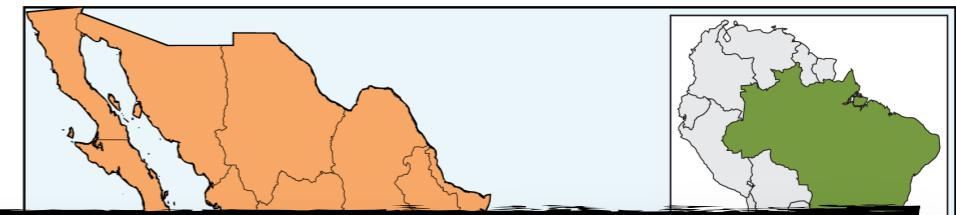
- Model parameters are rates that describe the growth of the tree
- Effective reproduction number = λ/δ
(is it spreading or not?)
- Allow parameters to change through time
- Shifts in rates can be anywhere
(usually evenly distributed)

Zika virus in Central America

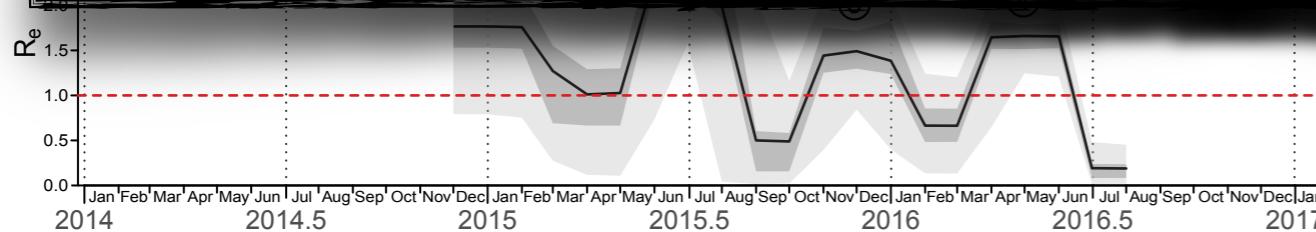


- 104 genomes
- Single introduction in Honduras in late 2014 / early 2015
- Birth-death skyline estimates of R_e show some evidence for seasonality

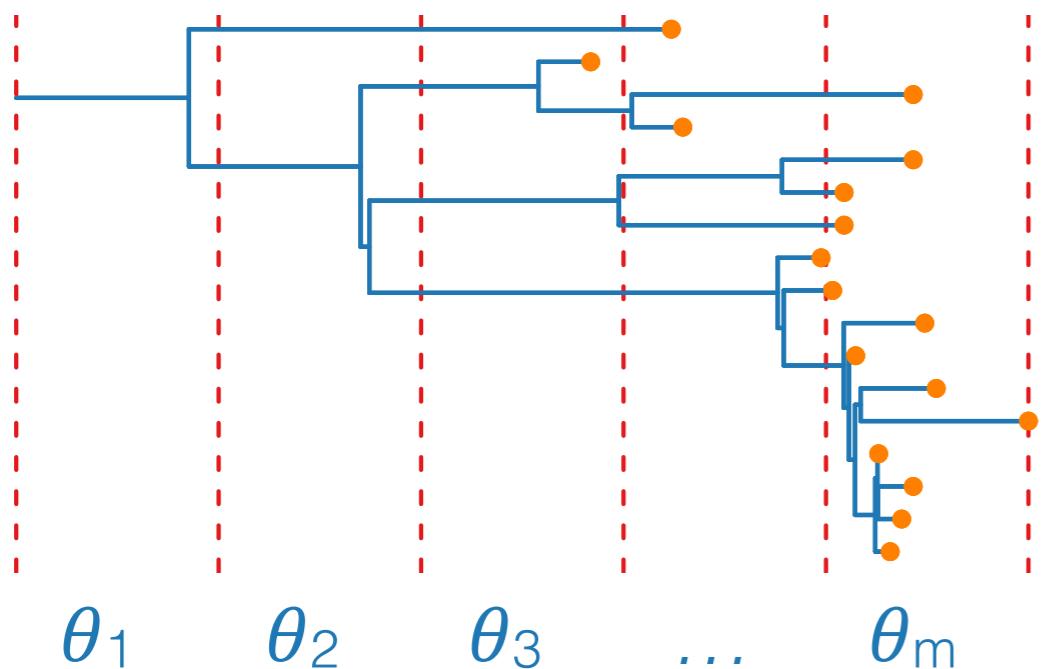
Zika virus in Central America



Faria et al. **Nature** 2017

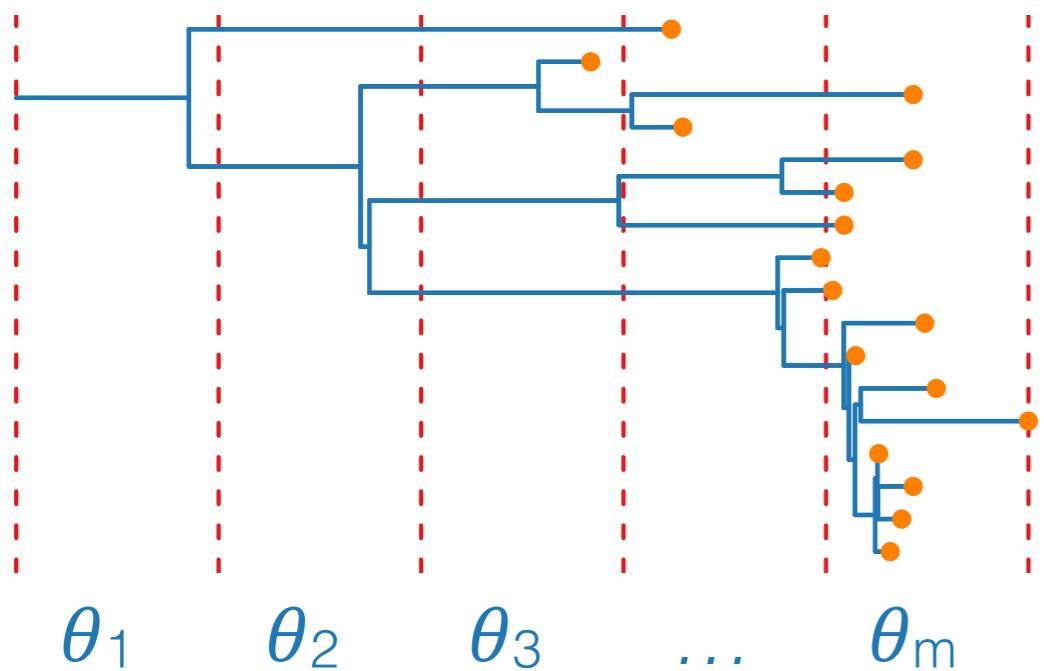


Phylodynamic linear model formulation



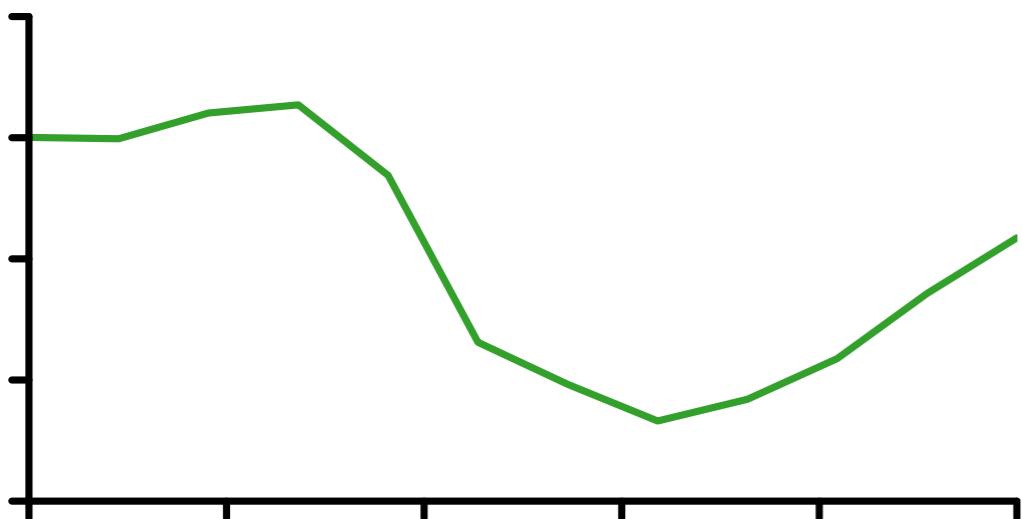
$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

Phylodynamic linear model formulation

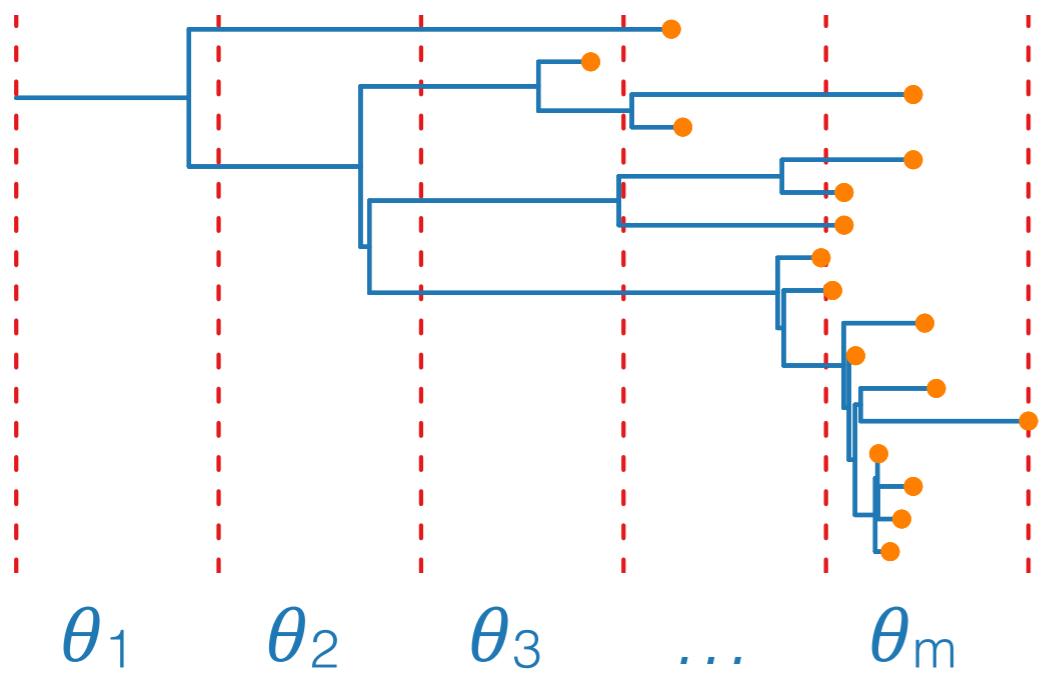


$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

Covariates



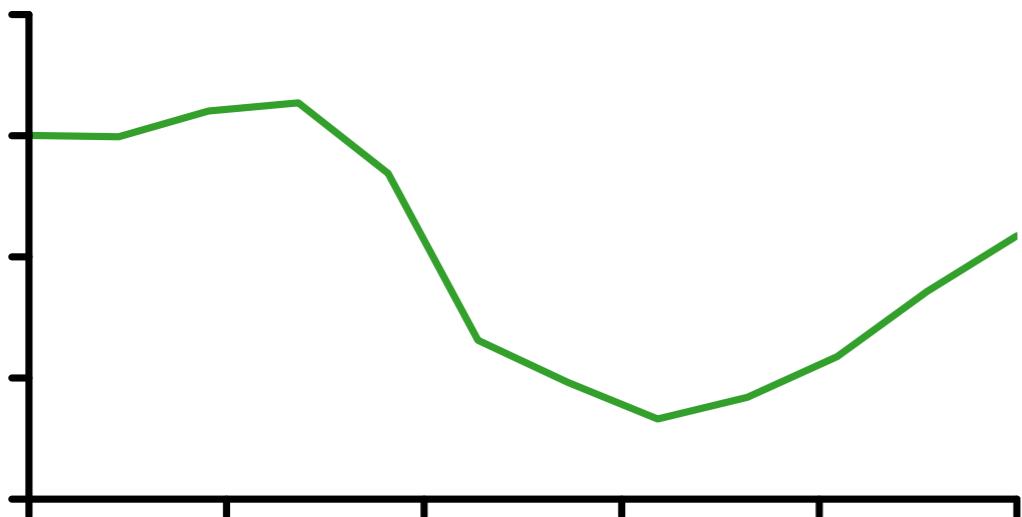
Phylodynamic linear model formulation



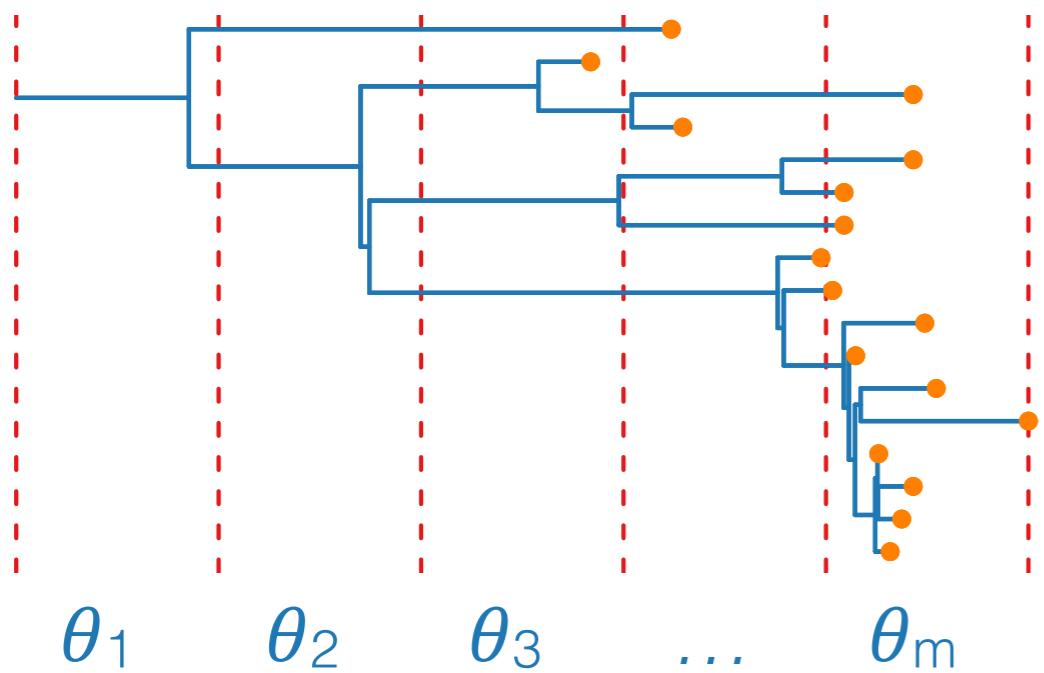
$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

Covariates

$$g(\theta_i) = x_{i,1}\beta_1$$



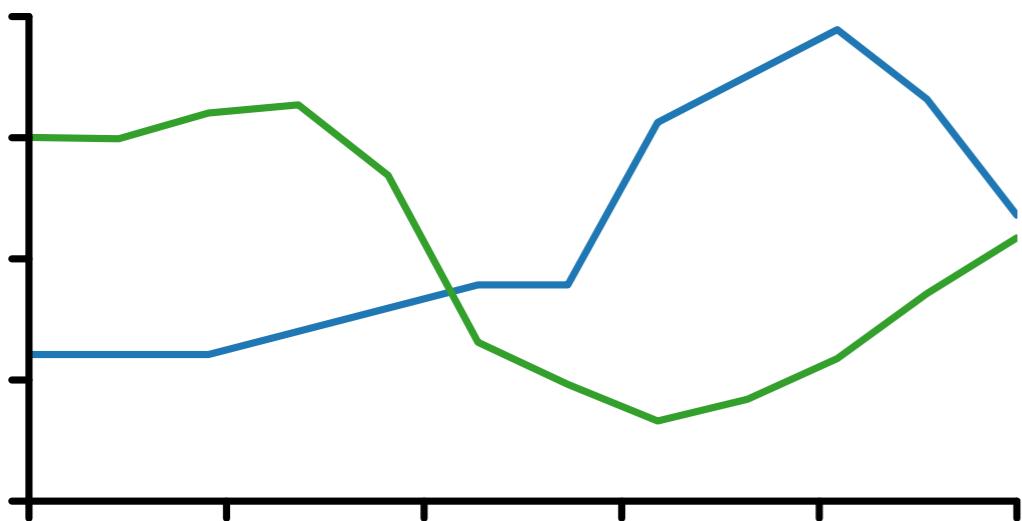
Phylodynamic linear model formulation



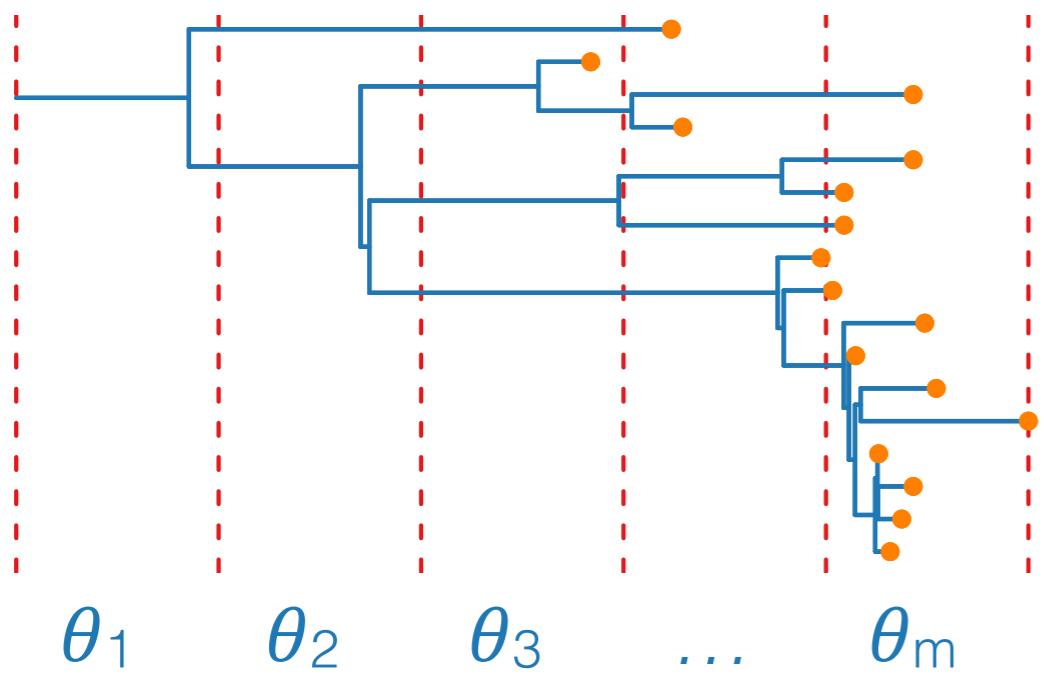
$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

Covariates

$$g(\theta_i) = x_{i,1}\beta_1 + x_{i,2}\beta_2$$



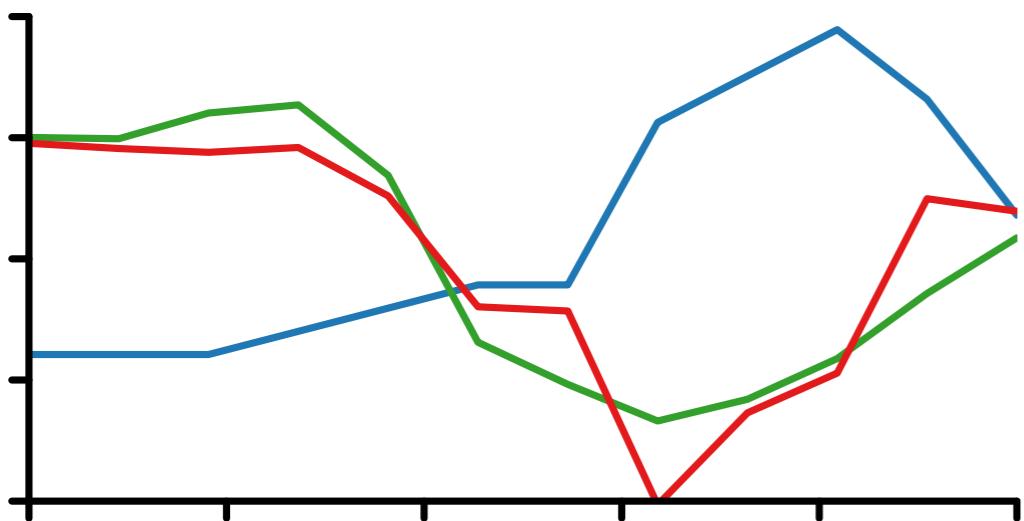
Phylodynamic linear model formulation



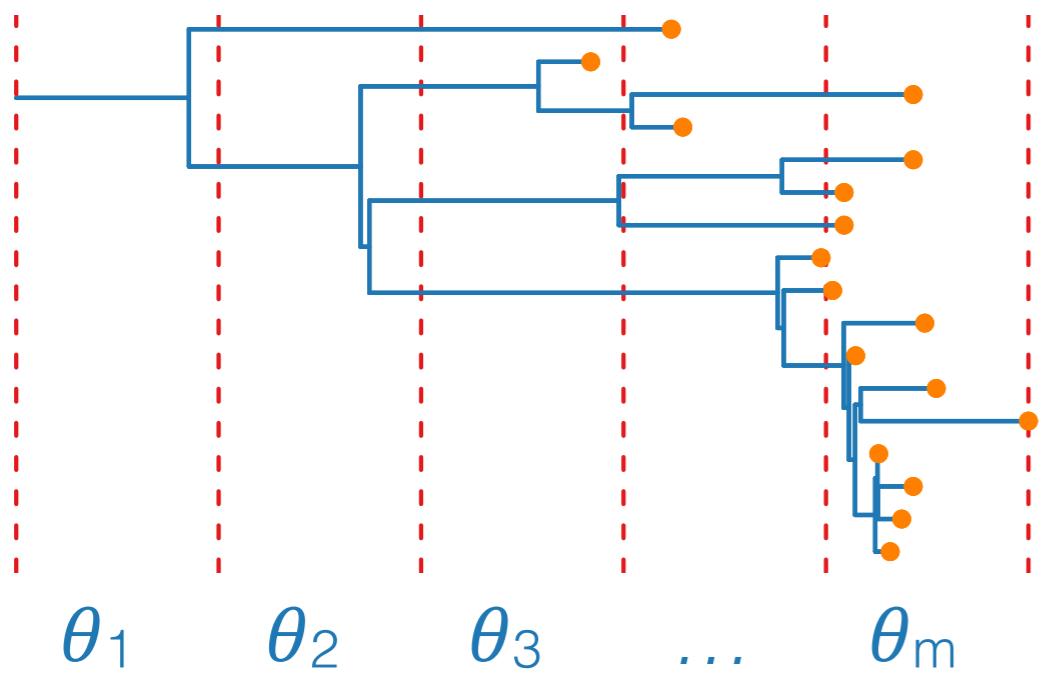
$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

Covariates

$$g(\theta_i) = x_{i,1}\beta_1 + x_{i,2}\beta_2 + x_{i,3}\beta_3$$



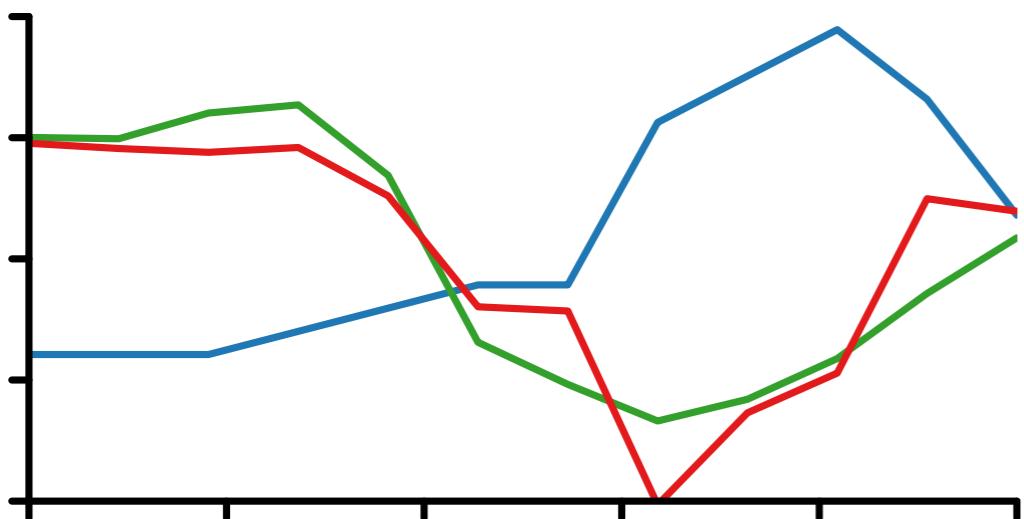
Phylodynamic linear model formulation



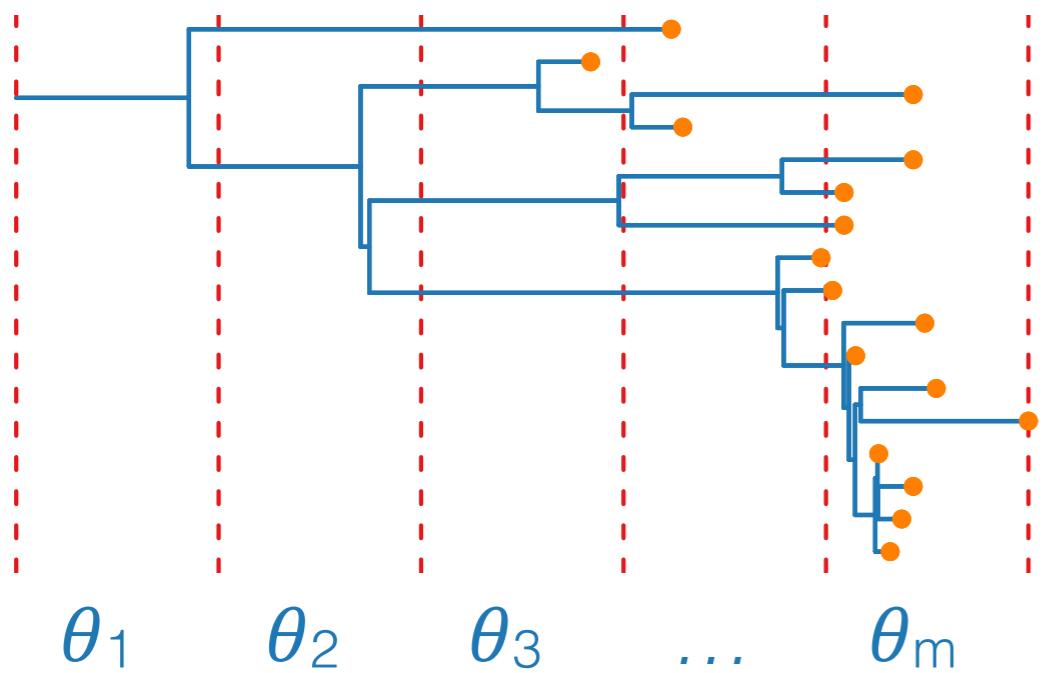
$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

Covariates

$$g(\theta_i) = x_{i,1}\beta_1 + x_{i,2}\beta_2 + \dots + x_{i,k}\beta_k$$

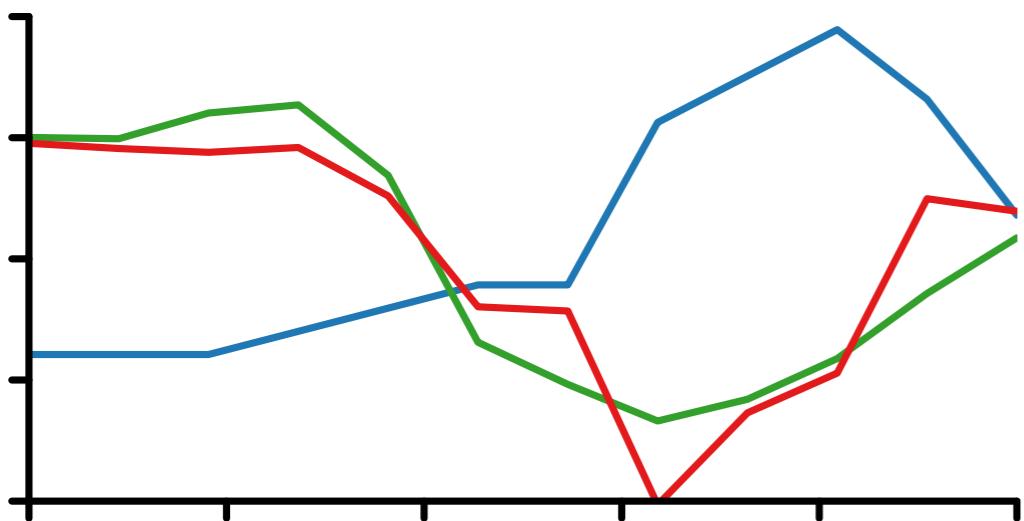


Phylodynamic linear model formulation



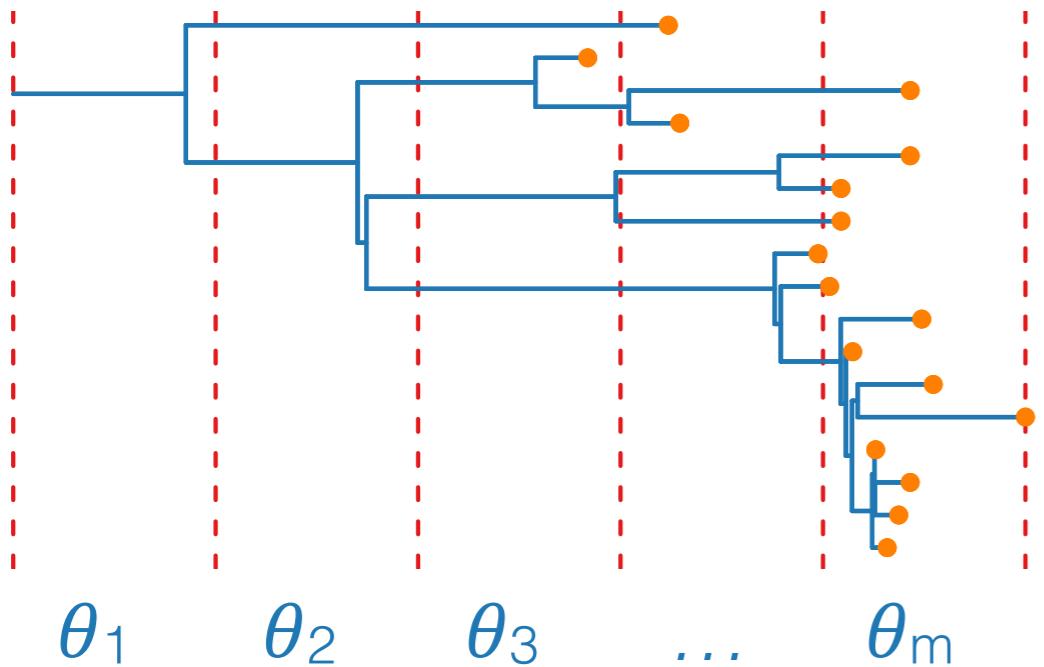
$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

Covariates



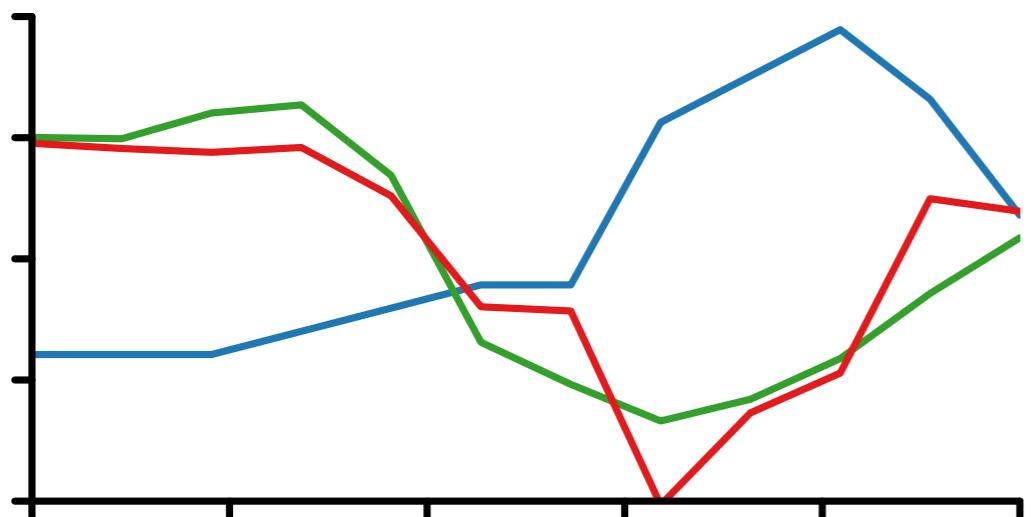
$$g(\theta_i) = x_{i,1}\beta_1 + x_{i,2}\beta_2 + \dots + x_{i,k}\beta_k + \varepsilon_i$$

Phylodynamic linear model formulation



$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

Covariates

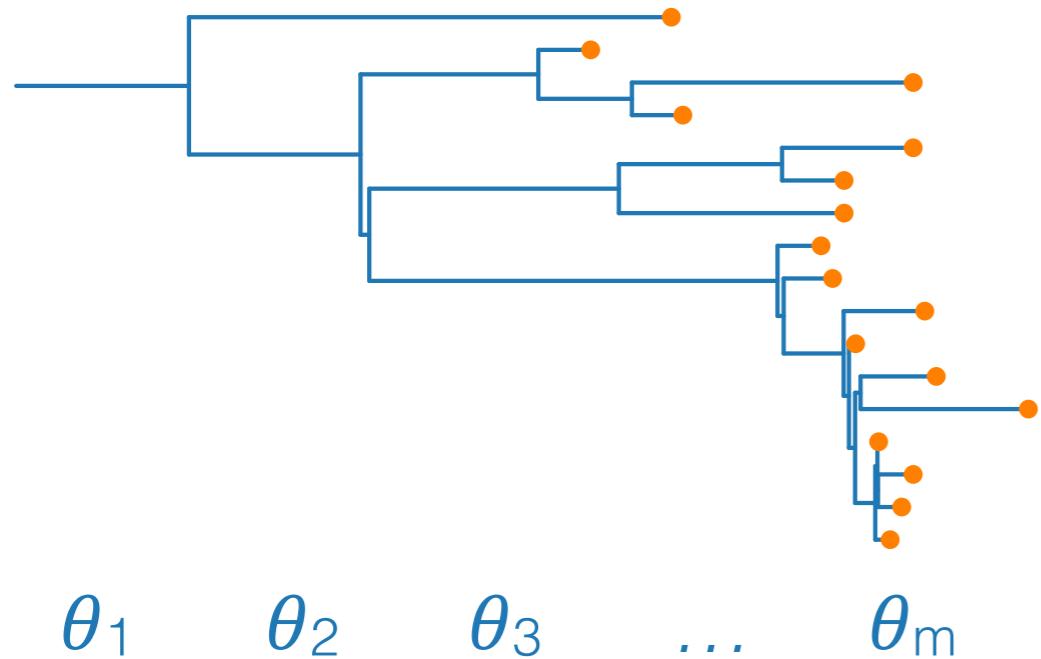


$$g(\theta_i) = x_{i,1}\beta_1 + x_{i,2}\beta_2 + \dots + x_{i,k}\beta_k + \varepsilon_i$$

$$g(\boldsymbol{\theta}) = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$[mx1] \quad [mxk] [kx1] \quad [mx1]$

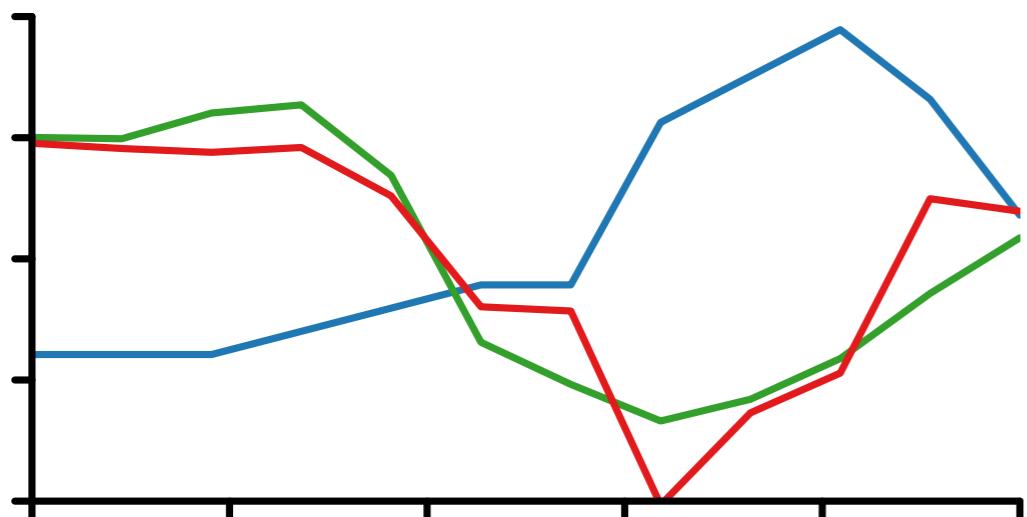
Linear model posteriors



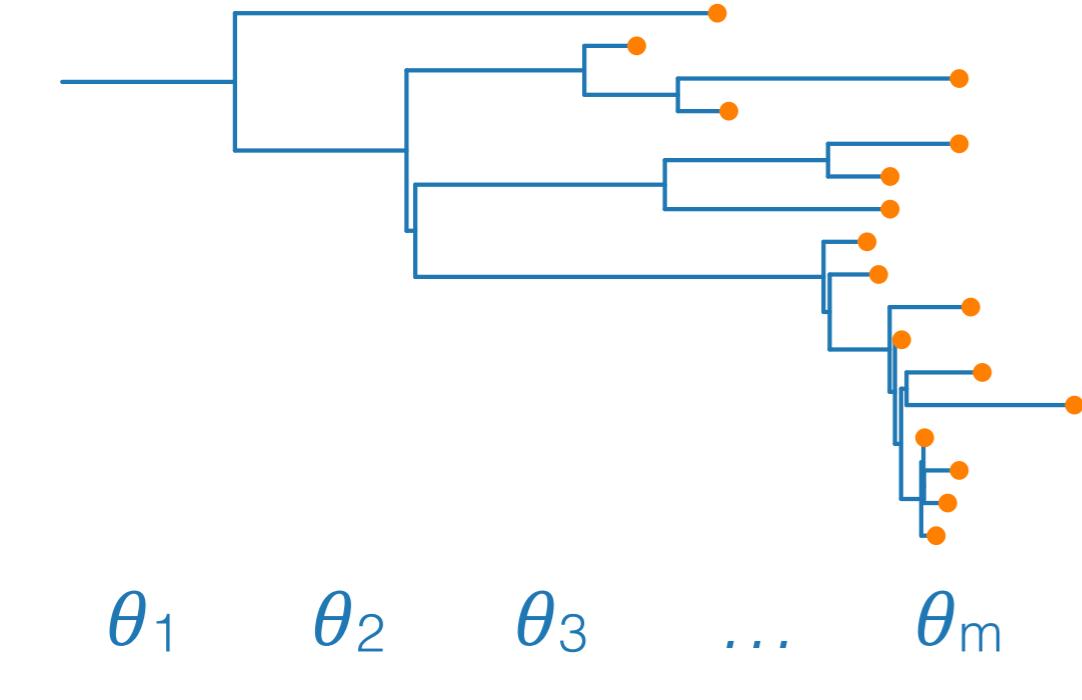
No linear model

$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

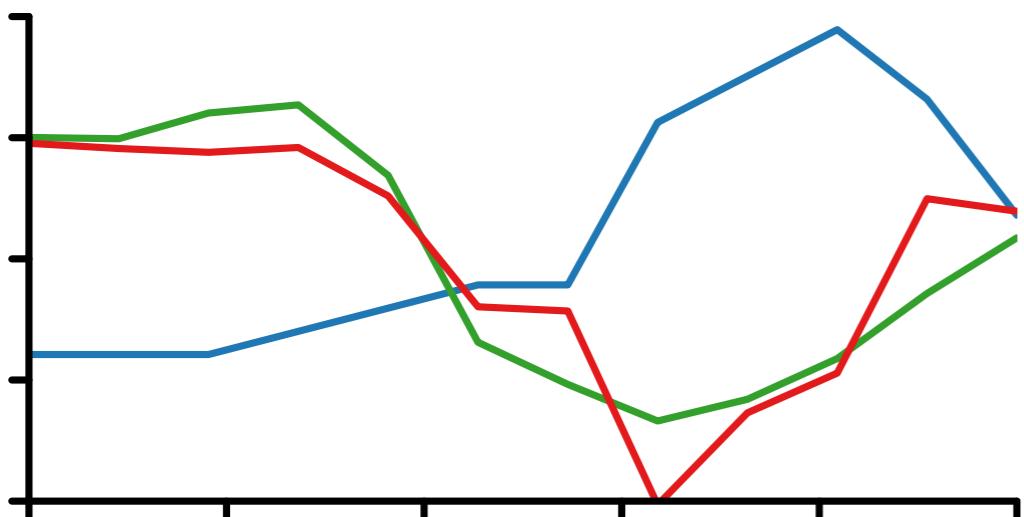
Covariates



Linear model posteriors



Covariates



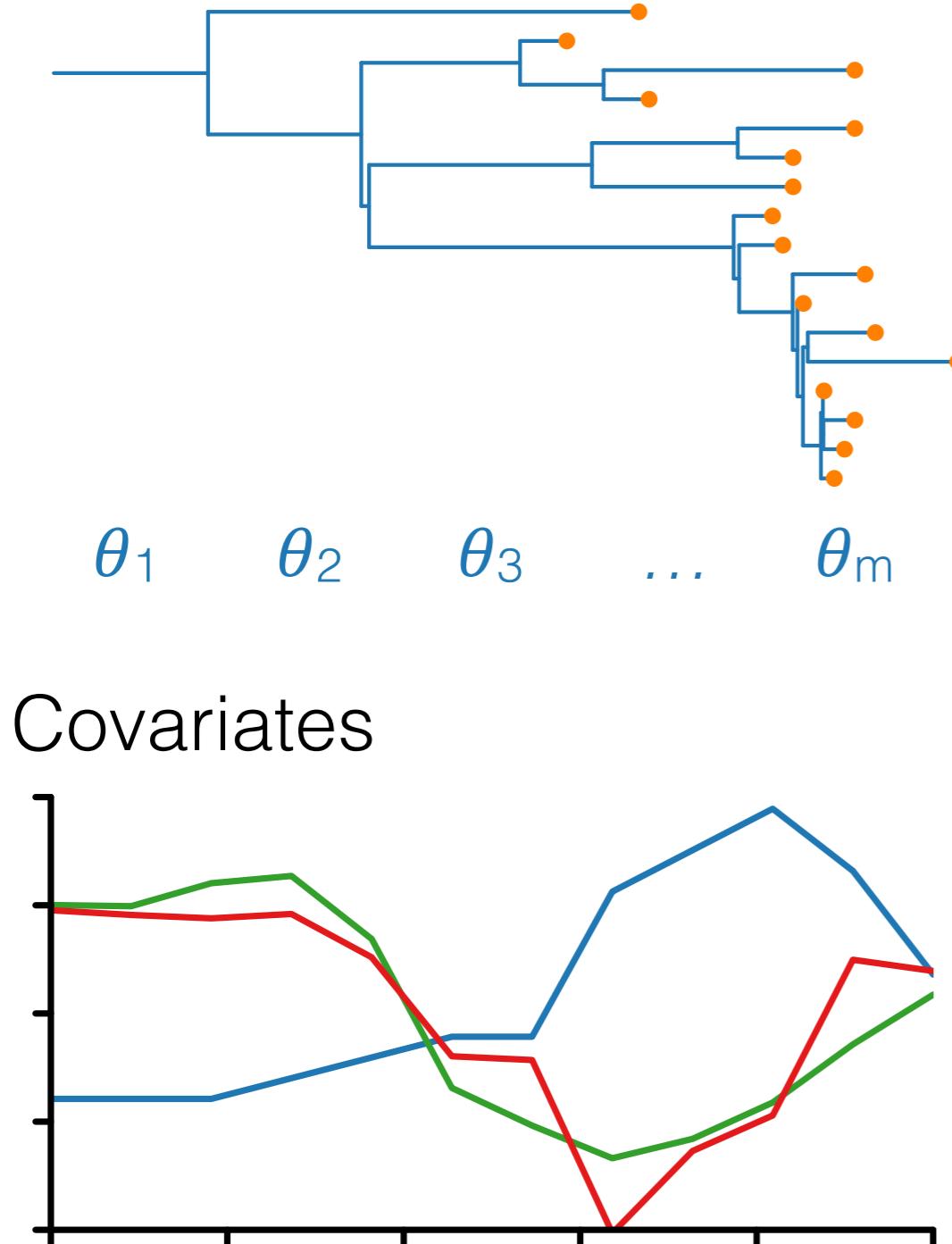
No linear model

$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

Full model: $g(\boldsymbol{\theta}) = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(g(\boldsymbol{\theta})|\mathbf{X}, \boldsymbol{\beta}) P(\boldsymbol{\beta})$$

Linear model posteriors



No linear model

$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$

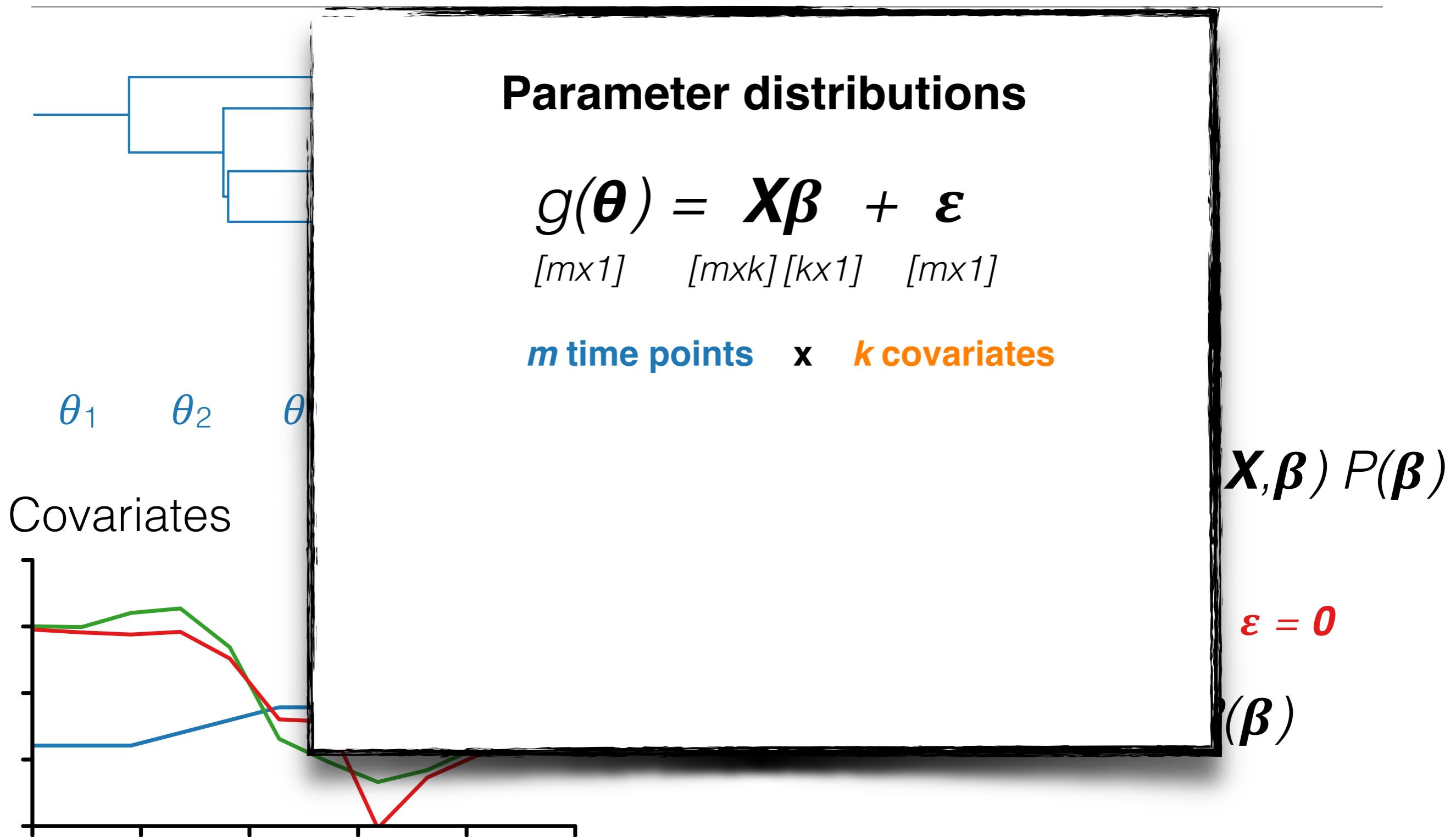
Full model: $g(\boldsymbol{\theta}) = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$$P(\boldsymbol{\theta} | T) \propto P(T|\boldsymbol{\theta}) P(g(\boldsymbol{\theta})|\mathbf{X}, \boldsymbol{\beta}) P(\boldsymbol{\beta})$$

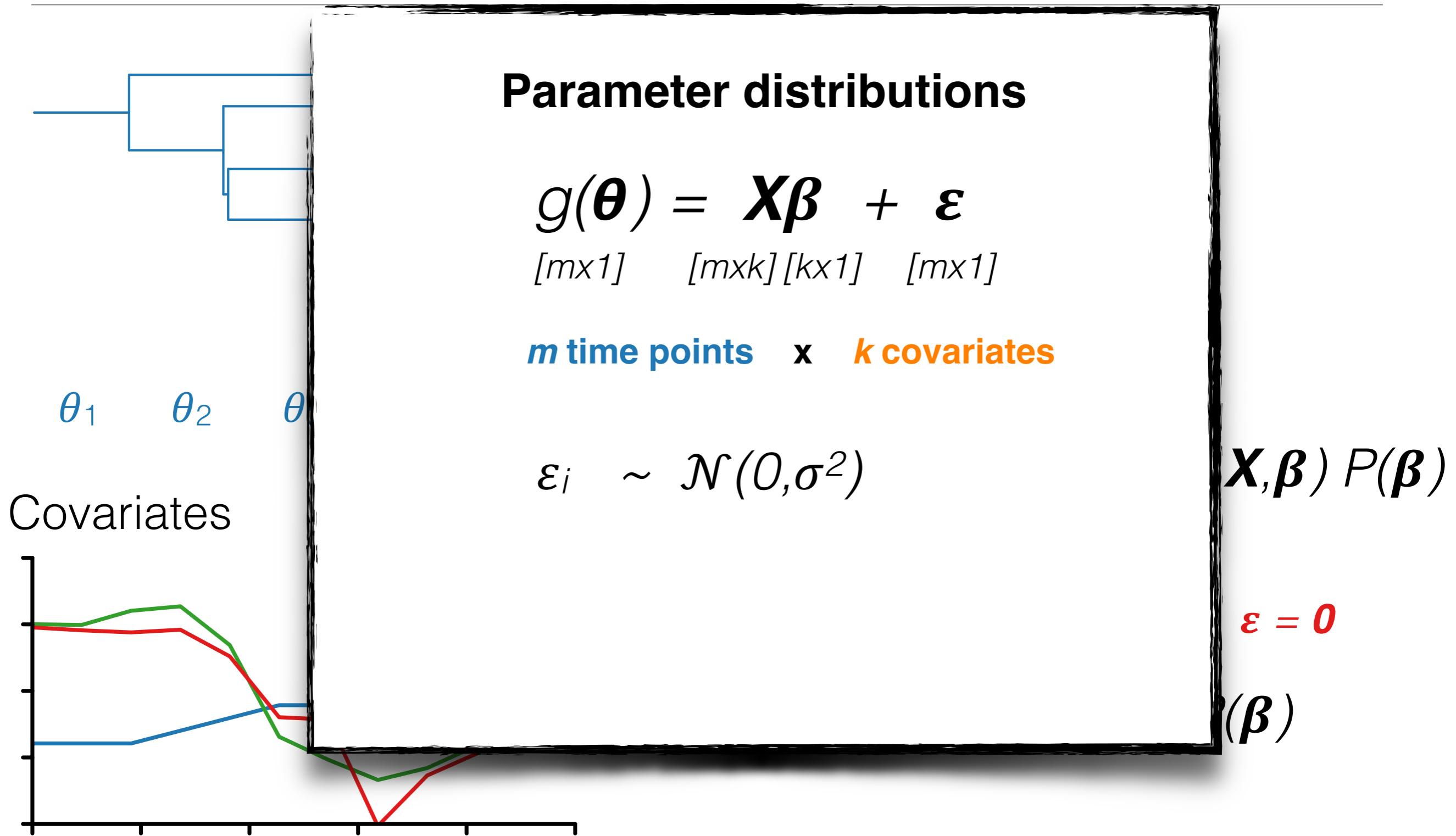
Exact model: $g(\boldsymbol{\theta}) = \mathbf{X}\boldsymbol{\beta} \quad \boldsymbol{\varepsilon} = \mathbf{0}$

$$P(\boldsymbol{\theta} | T) \propto P(T|g^{-1}(\mathbf{X}\boldsymbol{\beta})) P(\boldsymbol{\beta})$$

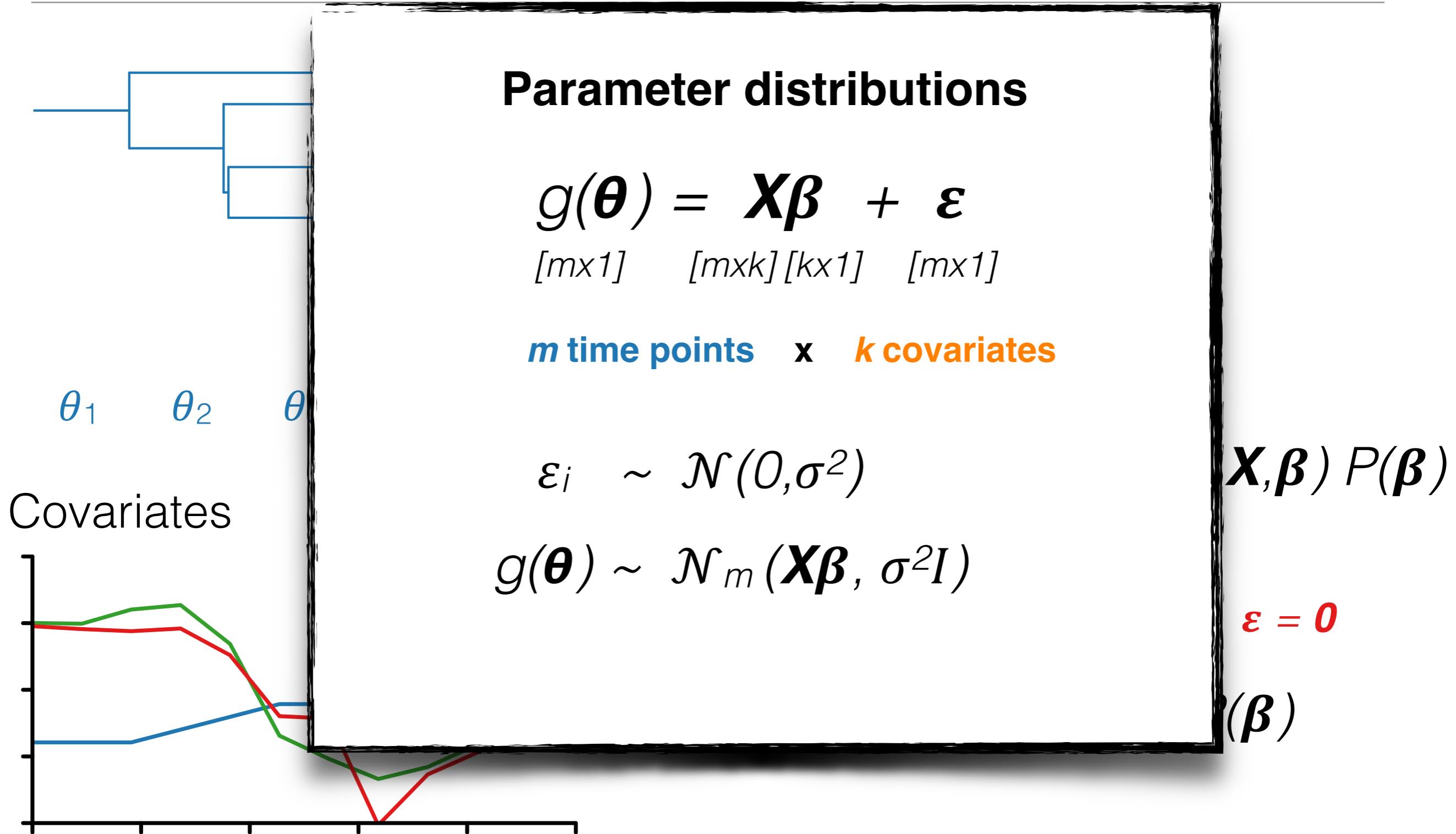
Linear model posteriors



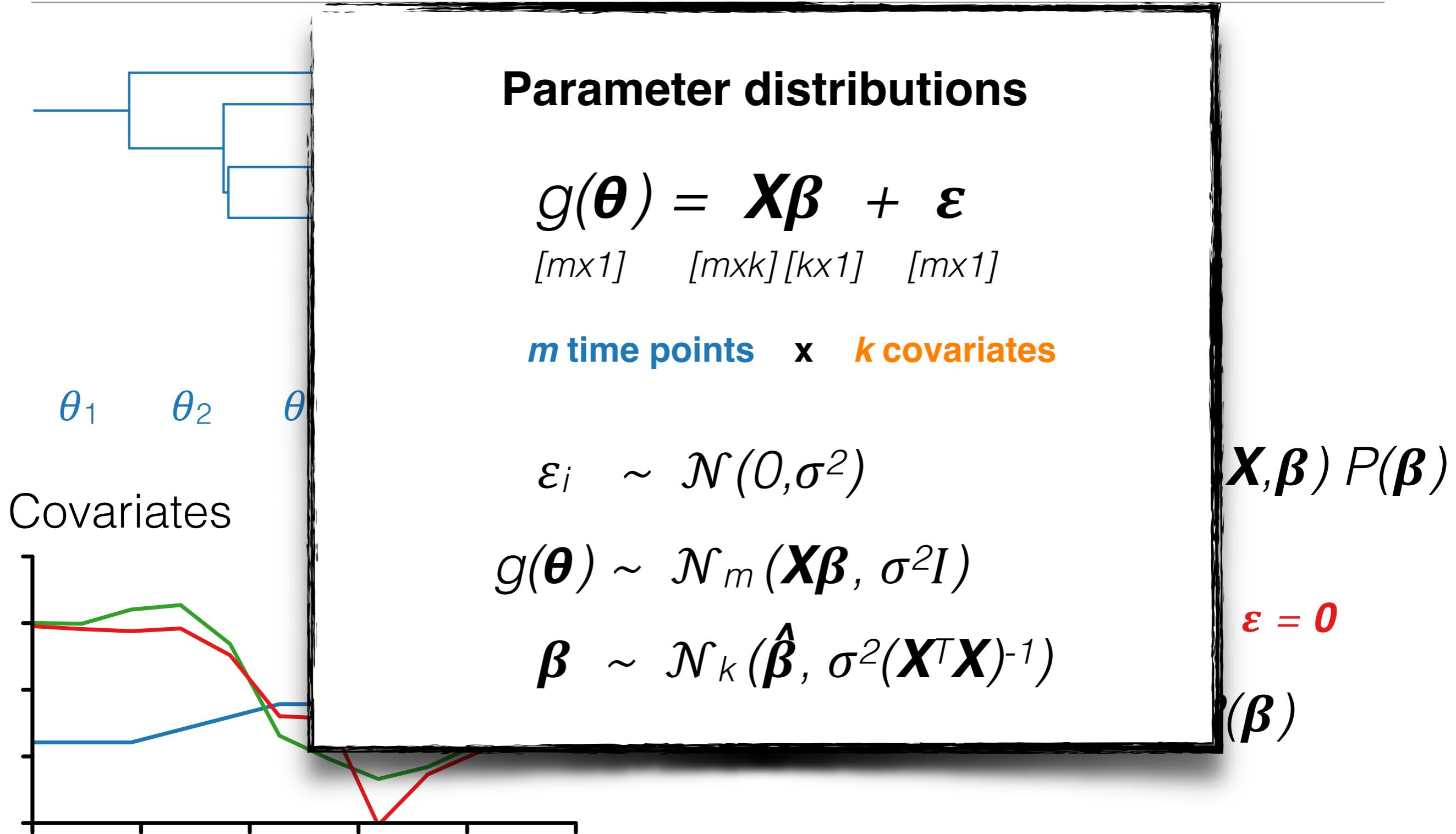
Linear model posteriors



Linear model posteriors



Linear model posteriors



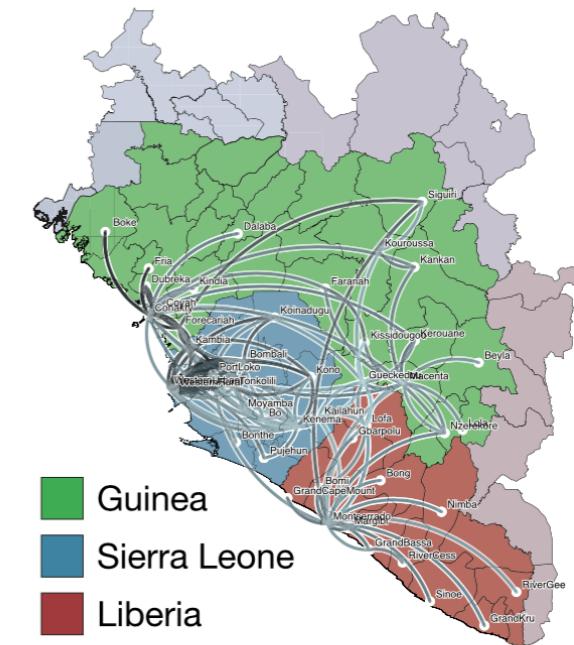
Spatial GLM model

OPEN ACCESS Freely available online



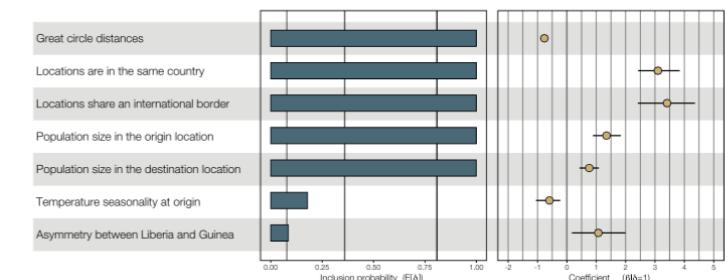
Unifying Viral Genetics and Human Transportation Data to Predict the Global Transmission Dynamics of Human Influenza H3N2

Philippe Lemey^{1*}, Andrew Rambaut^{2,3}, Trevor Bedford², Nuno Faria¹, Filip Bielejec¹, Guy Baele¹, Colin A. Russell^{4,5}, Derek J. Smith^{4,5,6}, Oliver G. Pybus⁷, Dirk Brockmann^{8,9,10}, Marc A. Suchard^{11,12}



- Linear model for migration rates
- Assume $\varepsilon = \mathbf{0}$ (exact model)

Lemey *et al.* PLOS Pathogens 2014

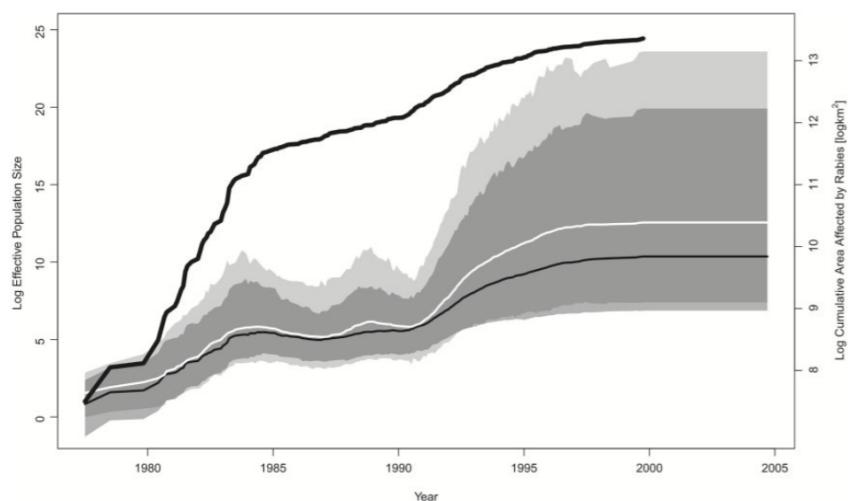


Suchard *et al.* Virus Evolution 2018

Skygrid GLM

Understanding Past Population Dynamics: Bayesian Coalescent-Based Modeling with Covariates

MANDEV S. GILL¹, PHILIPPE LEMEY², SHANNON N. BENNETT³, ROMAN BIEK⁴, AND MARC A. SUCHARD^{5,6,7,*}



- Linear model for changes in N_e over time
- Assume $\varepsilon \sim \mathcal{N}_m(\mathbf{0}, \Sigma)$

Gill *et al.* Systematic Biology 2016

Penalised regression

The Bayesian Lasso

Trevor PARK and George CASELLA

The Lasso estimate for linear regression parameters can be interpreted as a Bayesian posterior mode estimate when the regression parameters have independent Laplace (i.e., double-exponential) priors. Gibbs sampling from this posterior is possible using an expanded hierarchy with conjugate normal priors for the regression parameters and independent exponential priors on their variances. A connection with the inverse-Gaussian distribution provides tractable full conditional distributions. The Bayesian Lasso provides interval estimates (Bayesian credible intervals) that can guide variable selection. Moreover, the structure of the hierarchical model provides both Bayesian and likelihood methods for selecting the Lasso parameter. Slight modifications lead to Bayesian versions of other Lasso-related estimation methods, including bridge regression and a robust variant.

KEY WORDS: Empirical Bayes; Gibbs sampler; Hierarchical model; Inverse Gaussian; Linear regression; Penalized regression; Scale mixture of normals.

Park and Casella **Journal of the American Statistical Association 2008**

Penalised regression

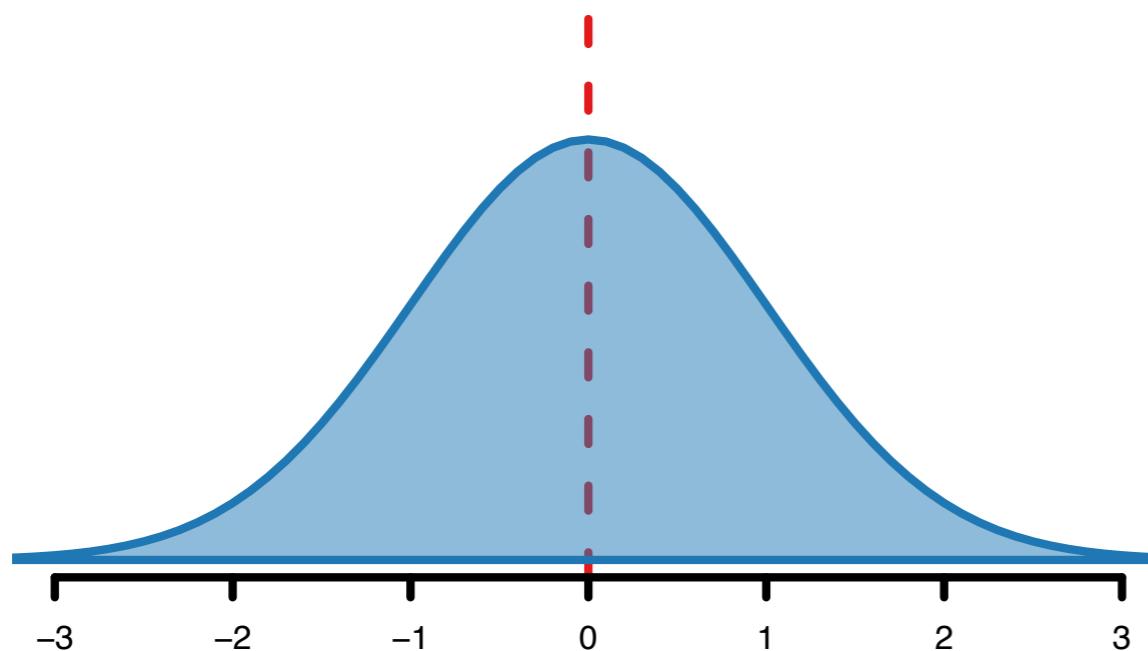
The Bayesian Lasso

Trevor PARK and George CASELLA

The Lasso estimate for linear regression parameters can be interpreted as a Bayesian posterior mode estimate when the regression parameters have independent Laplace (i.e., double-exponential) priors. Gibbs sampling from this posterior is possible using an expanded hierarchy with conjugate normal priors for the regression parameters and independent exponential priors on their variances. A connection with the inverse-Gaussian distribution provides tractable full conditional distributions. The Bayesian Lasso provides interval estimates (Bayesian credible intervals) that can guide variable selection. Moreover, the structure of the hierarchical model provides both Bayesian and likelihood methods for selecting the Lasso parameter. Slight modifications lead to Bayesian versions of other Lasso-related estimation methods, including bridge regression and a robust variant.

KEY WORDS: Empirical Bayes; Gibbs sampler; Hierarchical model; Inverse Gaussian; Linear regression; Penalized regression; Scale mixture of normals.

Park and Casella **Journal of the American Statistical Association 2008**



Penalised regression

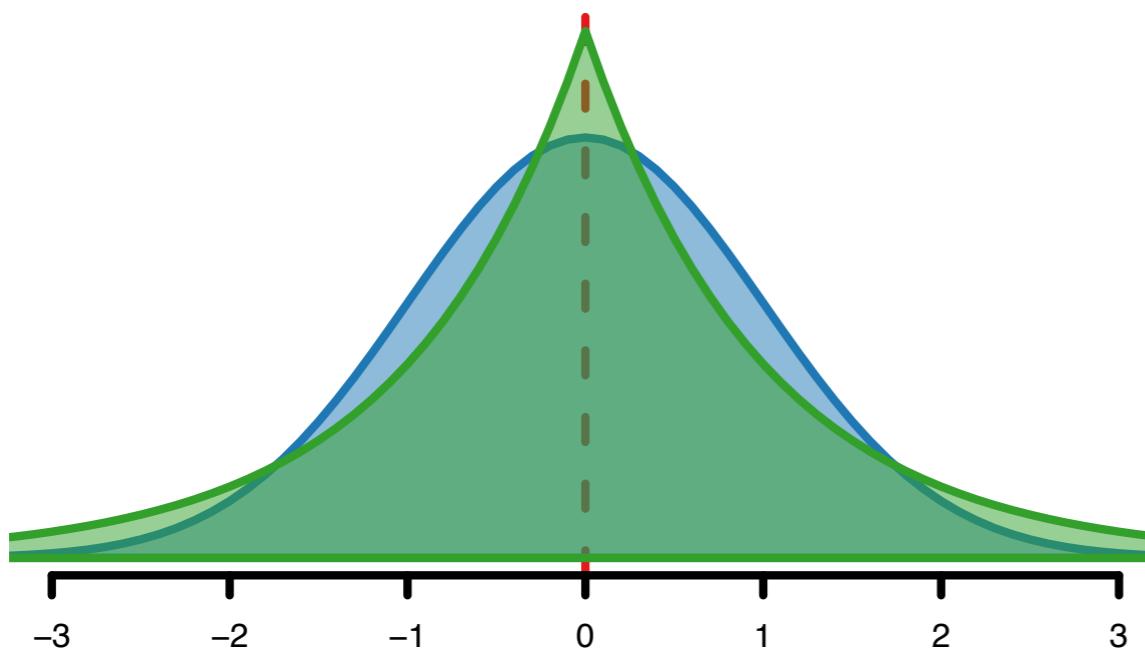
The Bayesian Lasso

Trevor PARK and George CASELLA

The Lasso estimate for linear regression parameters can be interpreted as a Bayesian posterior mode estimate when the regression parameters have independent Laplace (i.e., double-exponential) priors. Gibbs sampling from this posterior is possible using an expanded hierarchy with conjugate normal priors for the regression parameters and independent exponential priors on their variances. A connection with the inverse-Gaussian distribution provides tractable full conditional distributions. The Bayesian Lasso provides interval estimates (Bayesian credible intervals) that can guide variable selection. Moreover, the structure of the hierarchical model provides both Bayesian and likelihood methods for selecting the Lasso parameter. Slight modifications lead to Bayesian versions of other Lasso-related estimation methods, including bridge regression and a robust variant.

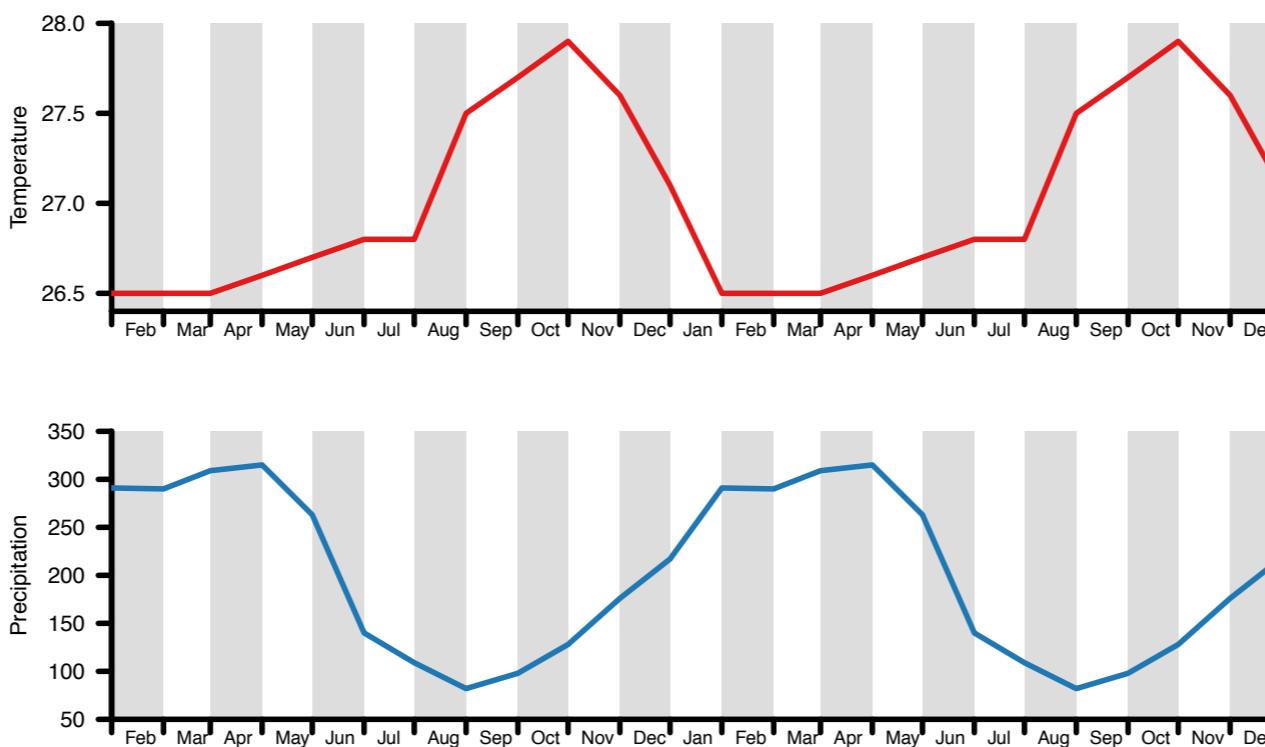
KEY WORDS: Empirical Bayes; Gibbs sampler; Hierarchical model; Inverse Gaussian; Linear regression; Penalized regression; Scale mixture of normals.

Park and Casella **Journal of the American Statistical Association 2008**



Lagged covariates

- Delays between covariate values and effects on θ
- Covariates may be measured at more time points than sequence data (especially climatic covariates)

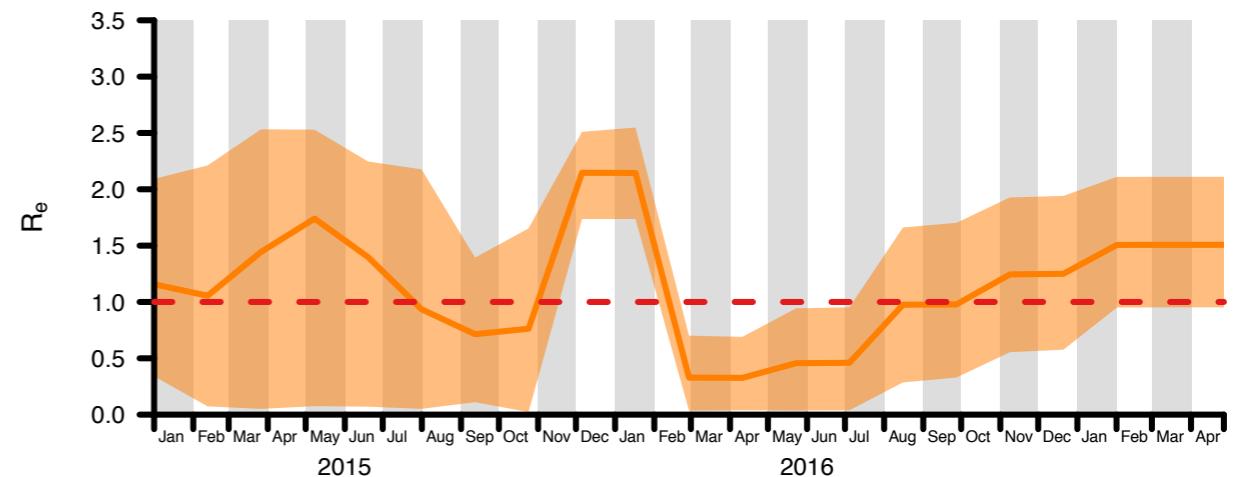


Introduce a timelag parameter:

$$g(\theta_i) = X_{i-s_1, 1} \beta_1 + X_{i-s_2, 2} \beta_2 + \varepsilon_i$$

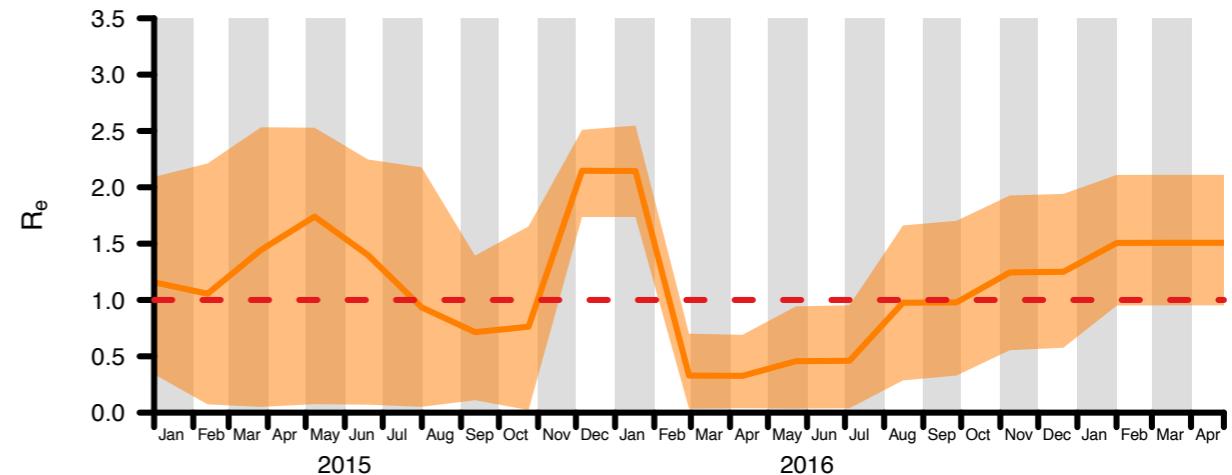
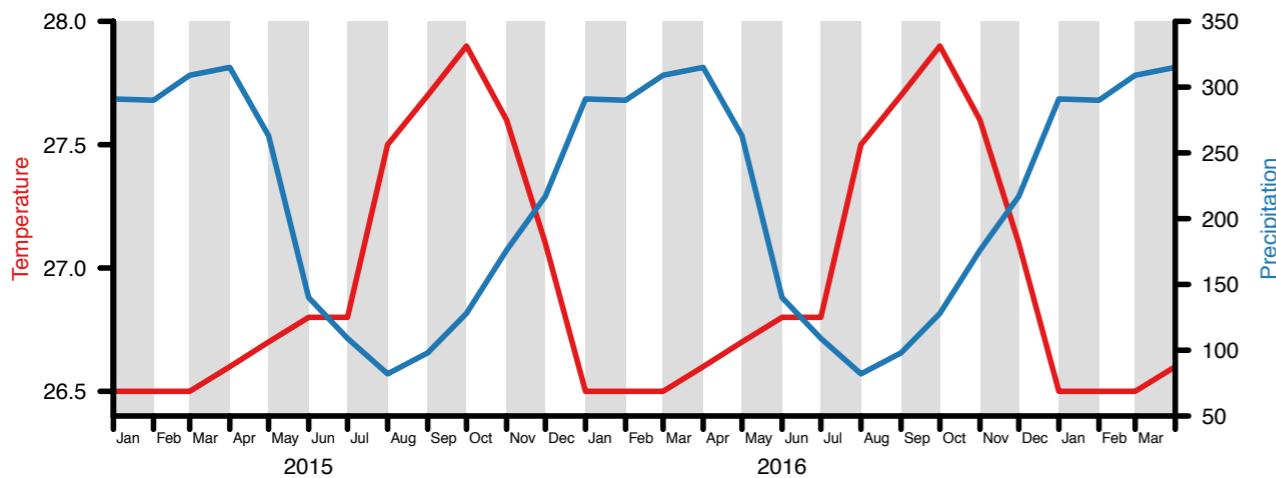
Zika virus in Manaus (57 genomes over 2 years)

No linear model ►



Zika virus in Manaus (57 genomes over 2 years)

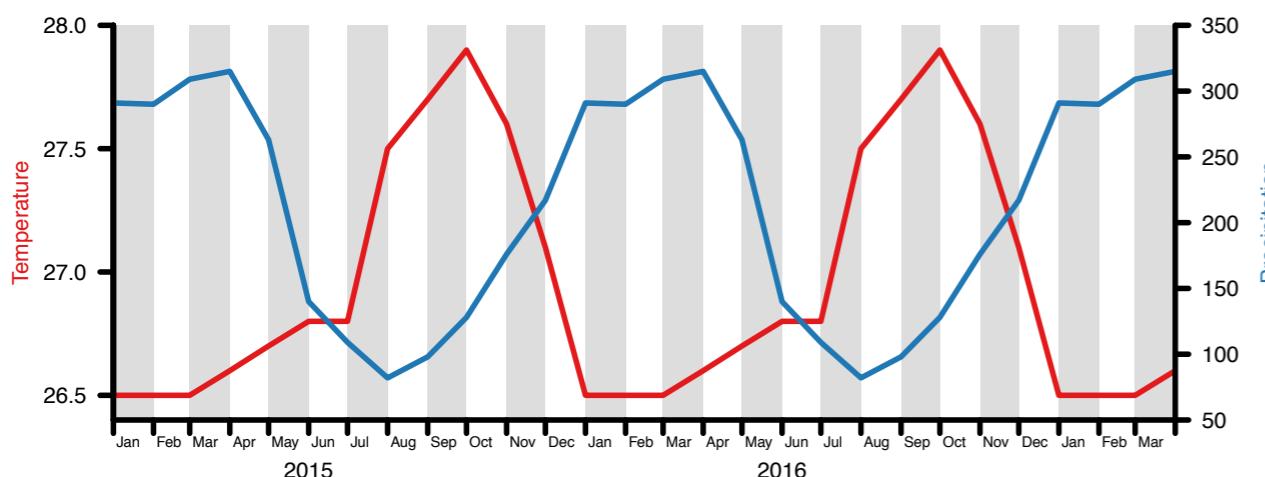
No linear model ►



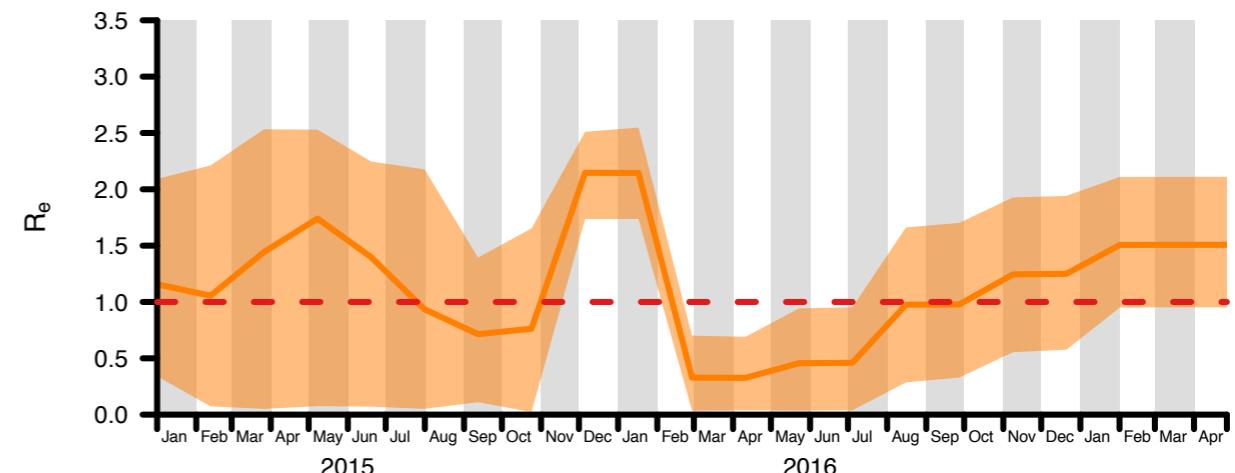
◀ Use rainfall and temperature
as covariates

Zika virus in Manaus (57 genomes over 2 years)

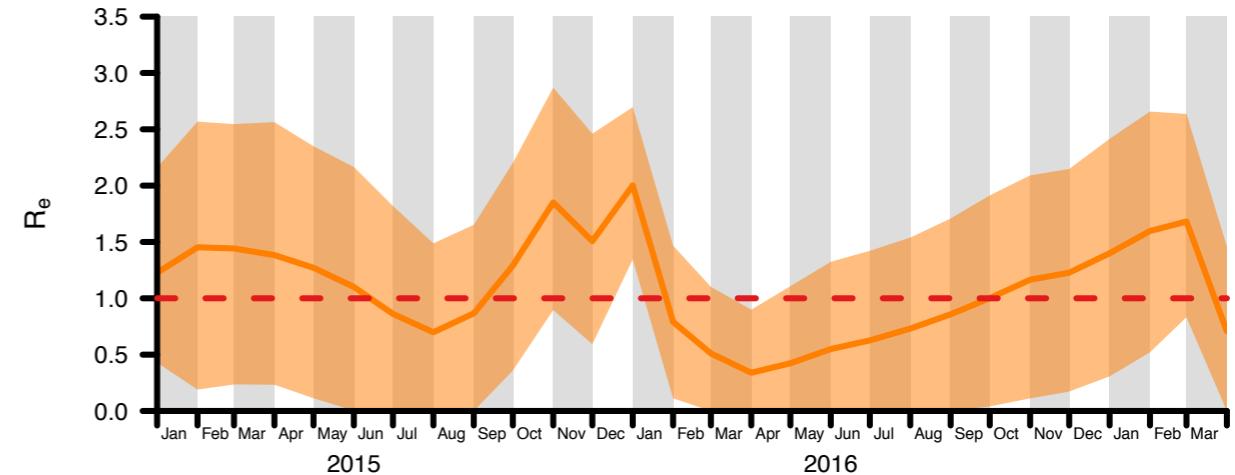
No linear model ►



With linear model ►

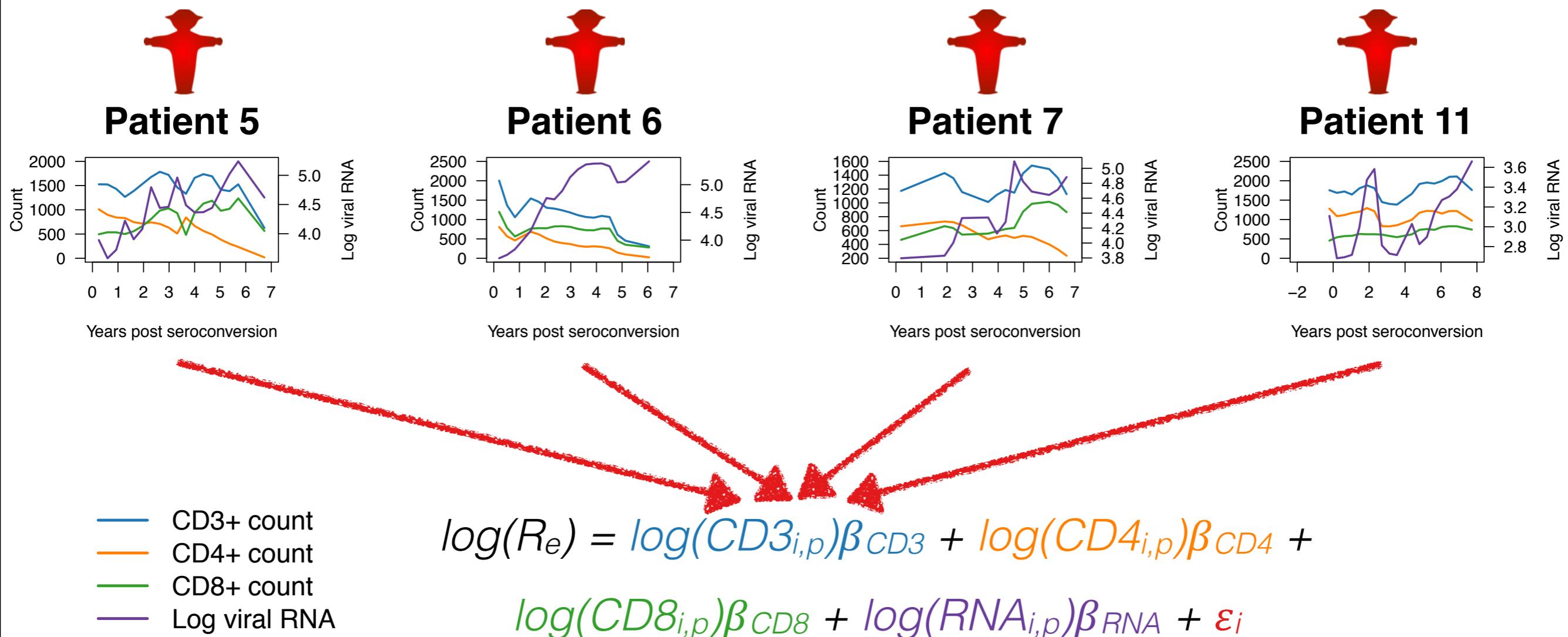


◀ Use rainfall and temperature
as covariates



Hierarchical within-host model

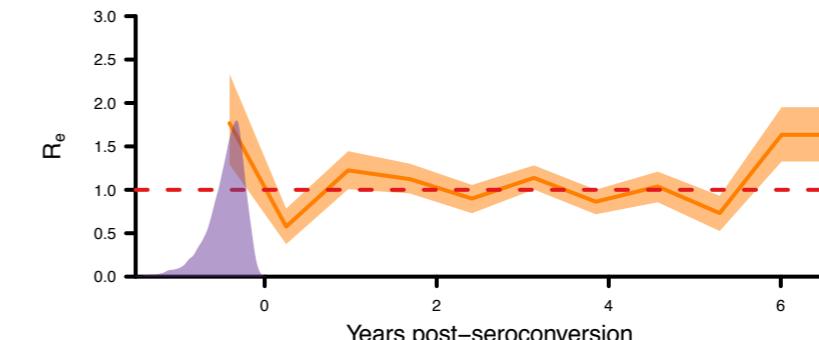
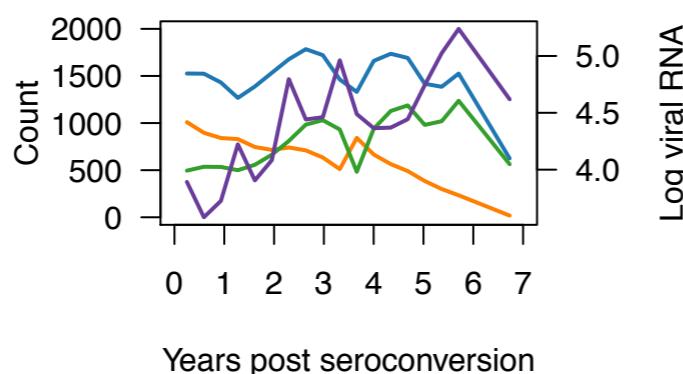
- Within-host HIV sequences and covariates from different patients sampled over time (Shankarappa et al. **Journal of Virology 1999**)
- Hierarchical linear model for R_e across multiple patients



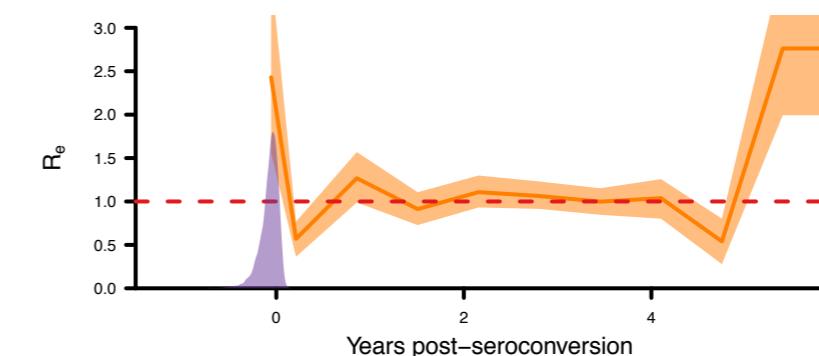
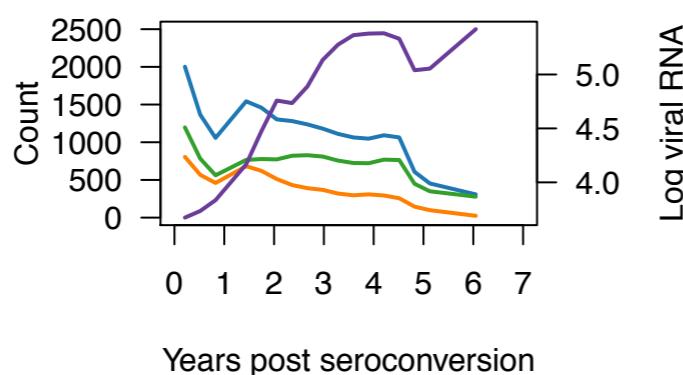
Hierarchical within-host model results



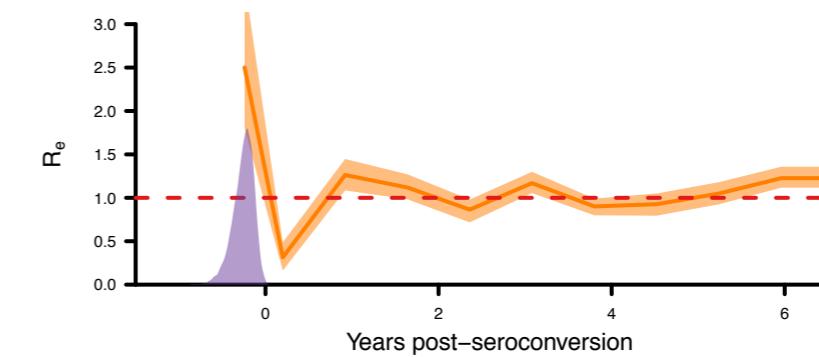
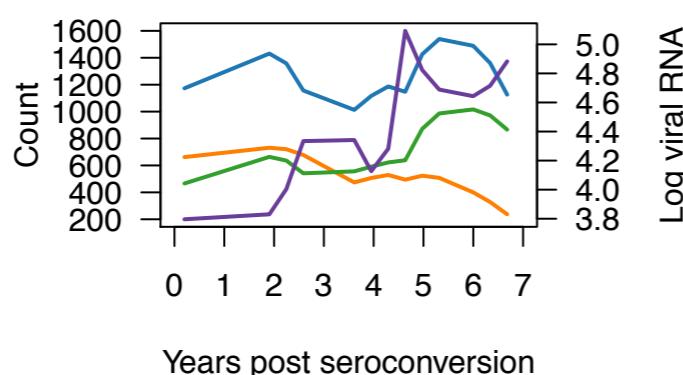
Patient 5



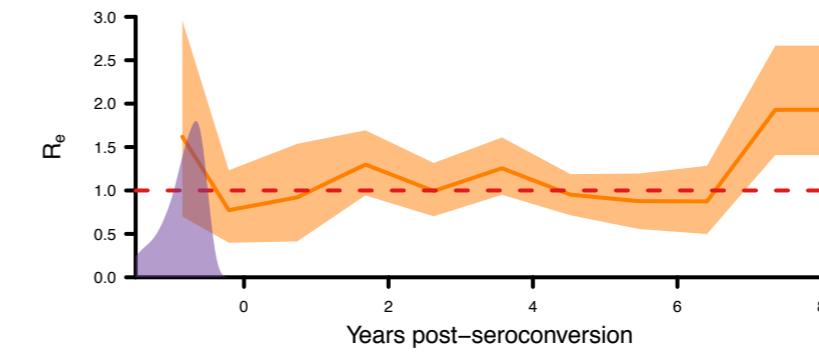
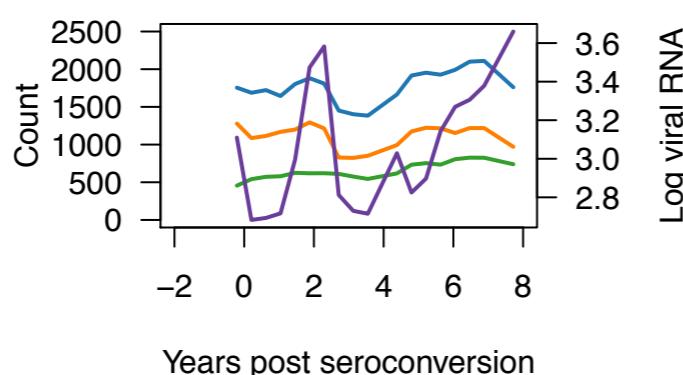
Patient 6



Patient 7



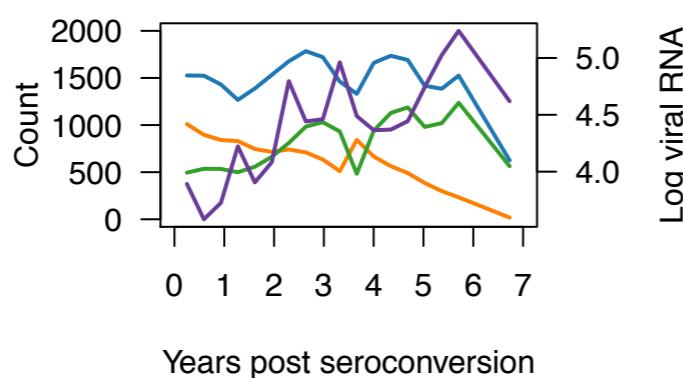
Patient 11



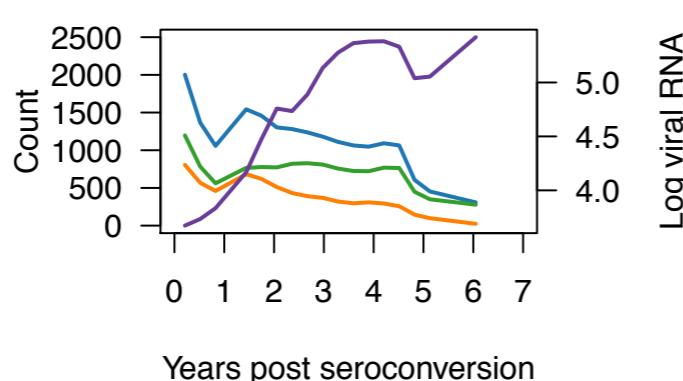
Hierarchical within-host model results



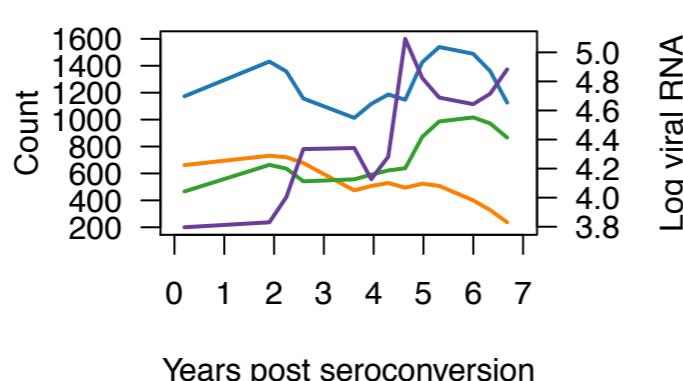
Patient 5



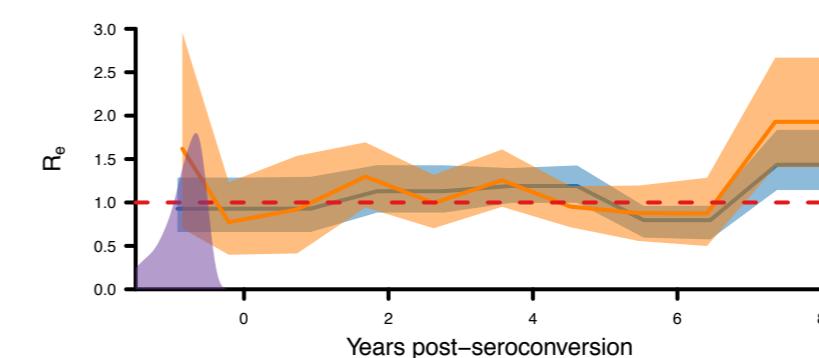
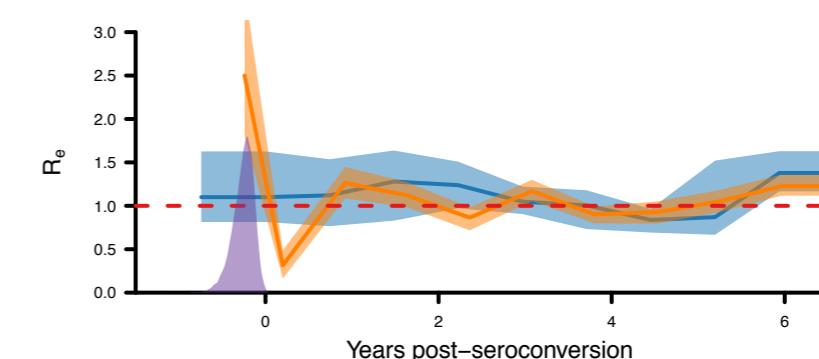
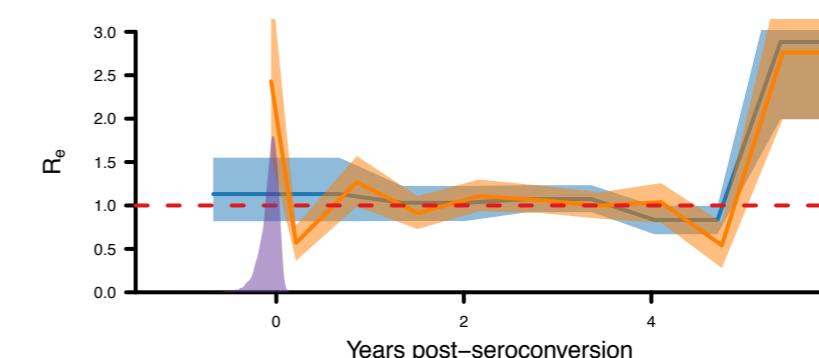
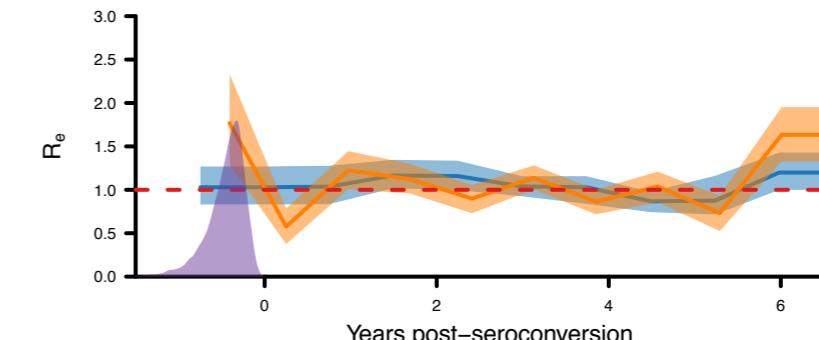
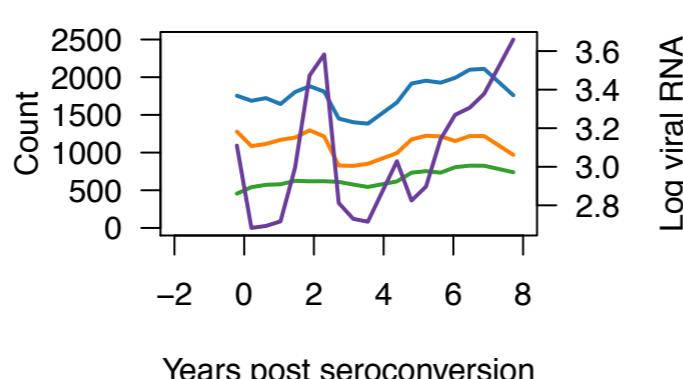
Patient 6



Patient 7

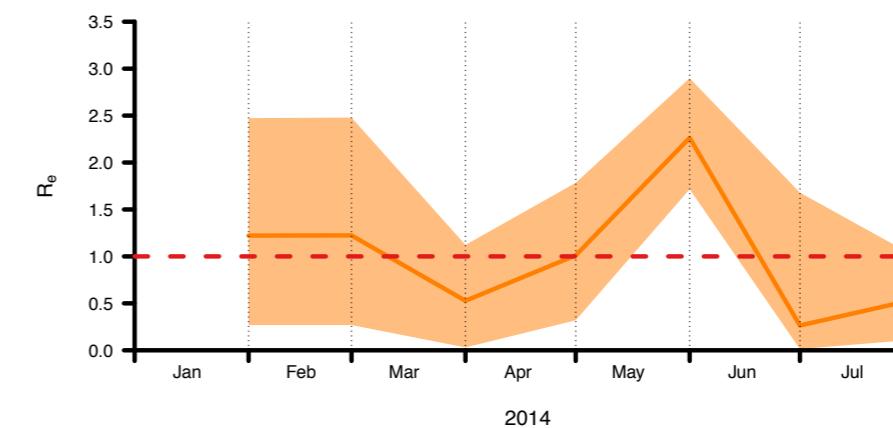
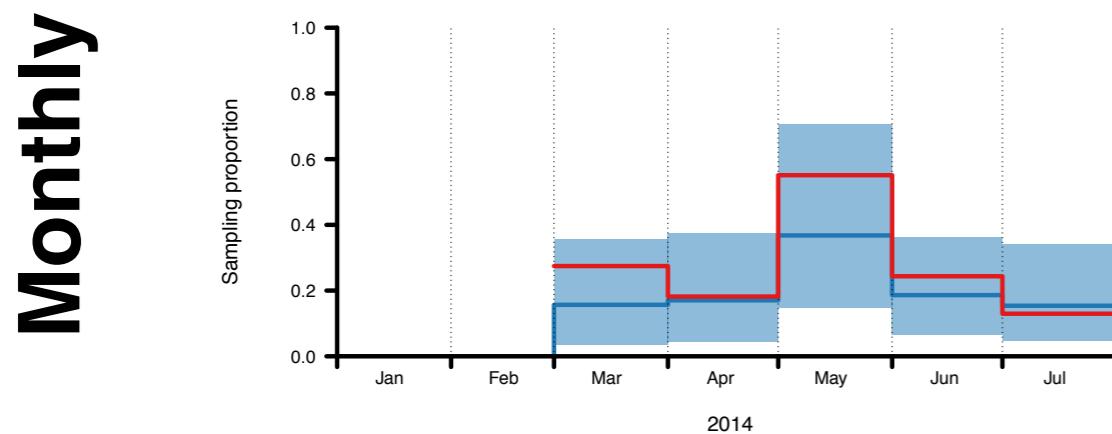


Patient 11



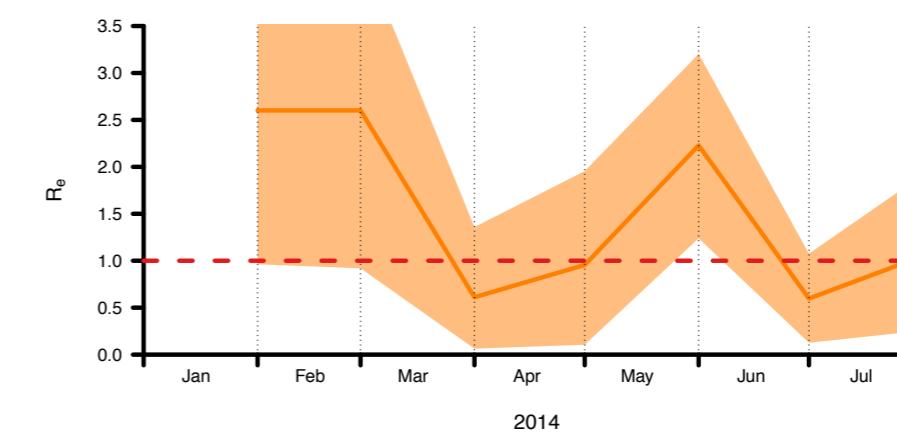
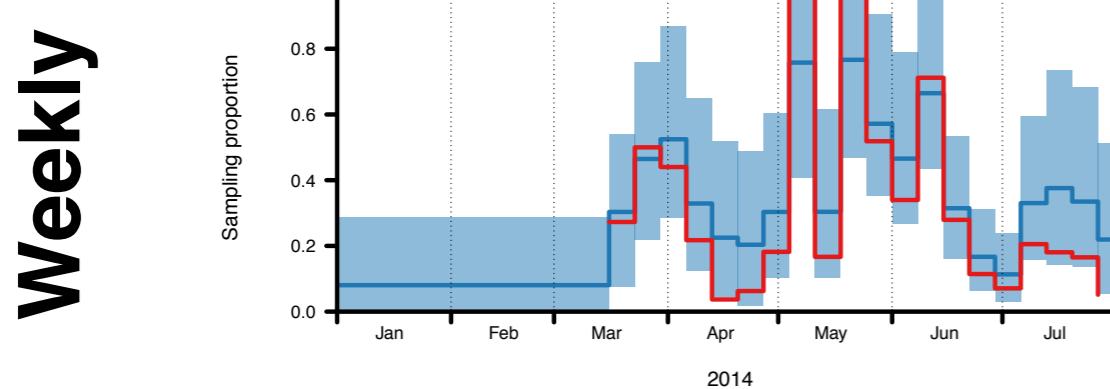
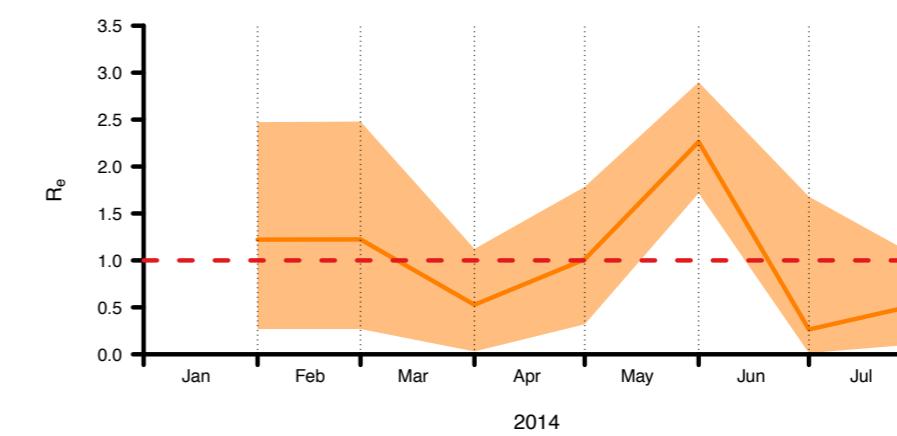
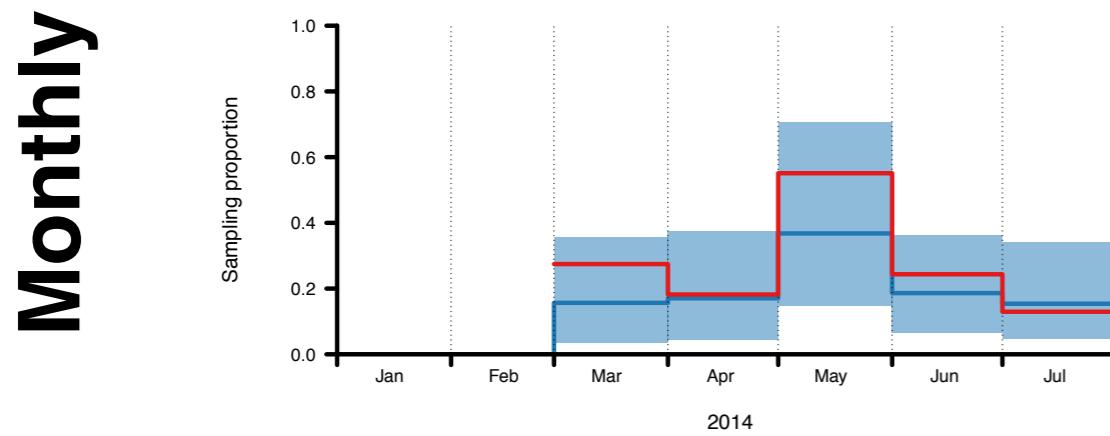
West African Ebola early epidemic

- 244 genomes sampled and 1424 cases reported until 4 August 2014 (WHO declares international public health emergency)
- Sampling proportion varies a lot over time
- Use linear model to put a multivariate prior on sampling proportion:
 $p_i = (\text{sequences}_i/\text{cases}_i)\beta + \varepsilon_i$



West African Ebola early epidemic

- 244 genomes sampled and 1424 cases reported until 4 August 2014 (WHO declares international public health emergency)
- Sampling proportion varies a lot over time
- Use linear model to put a multivariate prior on sampling proportion:
 $p_i = (\text{sequences}_i/\text{cases}_i)\beta + \varepsilon_i$



Conclusions and future directions

- As data from covariates keeps growing we can use it to **restrict the model space** and obtain more **accurate** answers
- Need to have some prior knowledge to **pick appropriate covariates**

Conclusions and future directions

- As data from covariates keeps growing we can use it to **restrict the model space** and obtain more **accurate** answers
- Need to have some prior knowledge to **pick appropriate covariates**
- **Extension I:**
Use indicator variables to switch off covariates with poor predictive capacity (**remove** unnecessary **covariates**)
- **Extension II:**
Add autocorrelation between residuals
(e.g. model ϵ as a zero-mean **Gaussian process**)
- **Extension III:**
Use maximum-likelihood estimate of β for better proposals

A background image showing numerous red and white COVID-19 virus particles against a dark gray gradient.

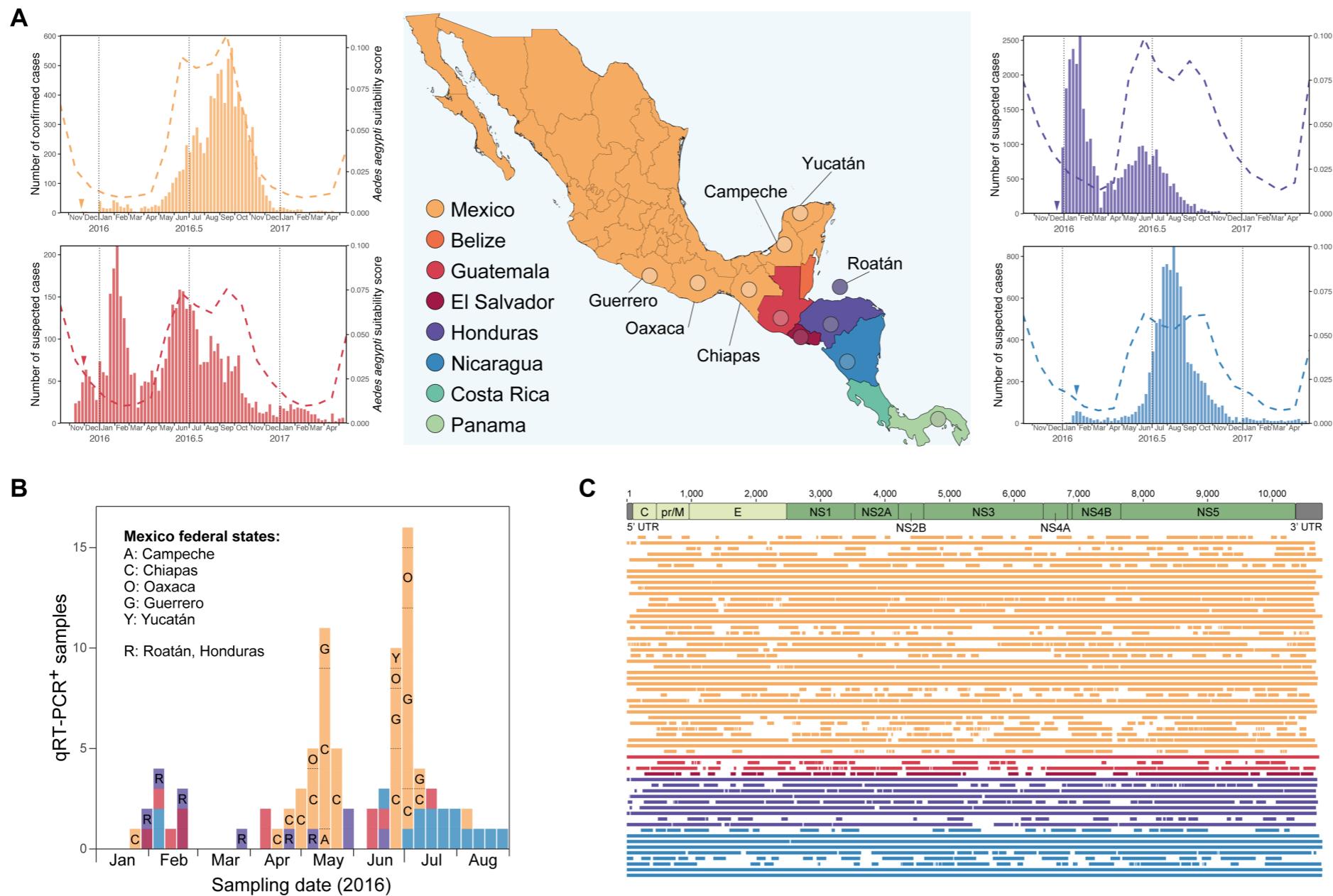
THANK YOU FOR LISTENING!

THANKS TO:

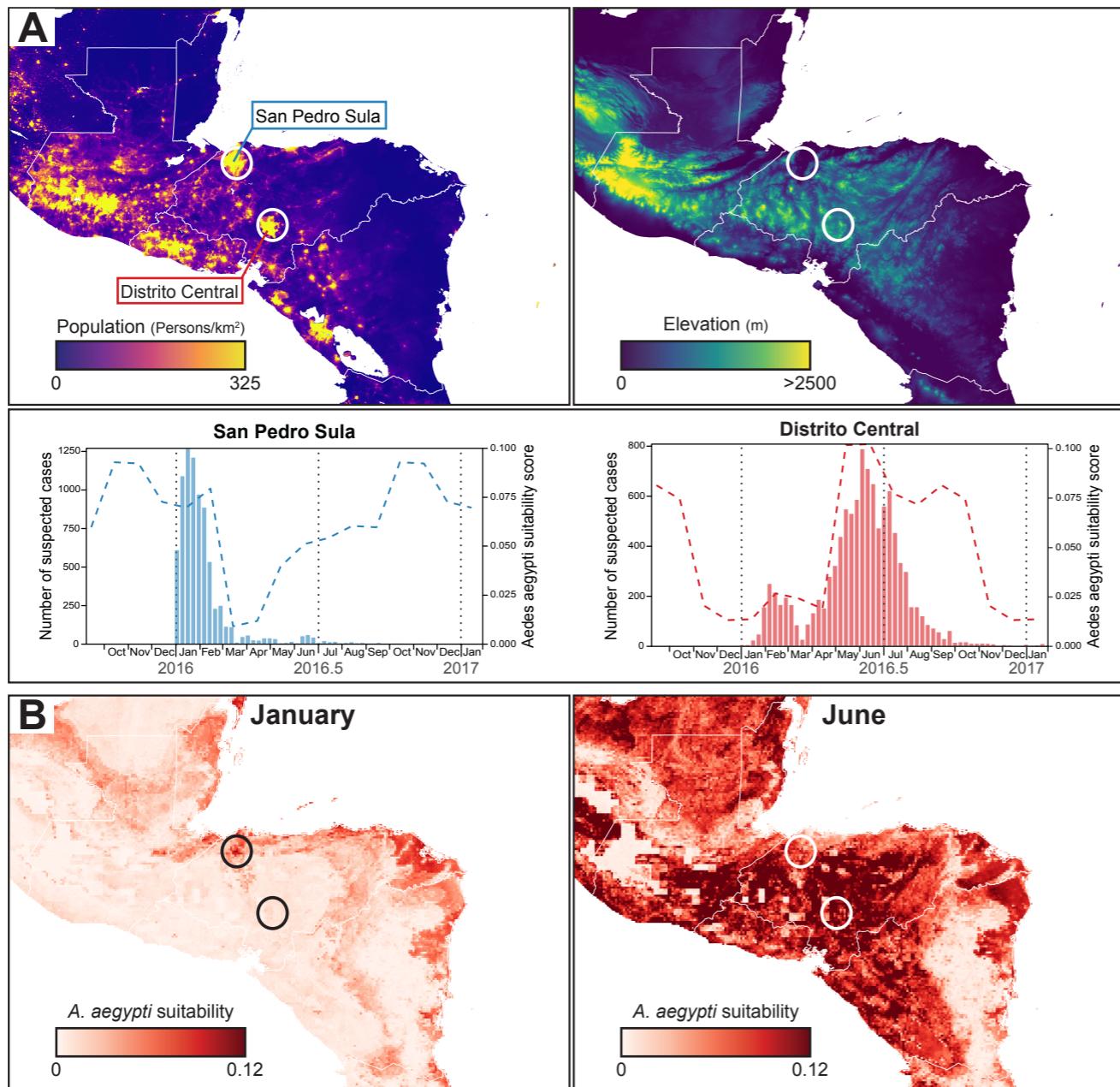
Oliver Pybus, Julien Thézé, Nuno Faria, Kris Parag

QUESTIONS?

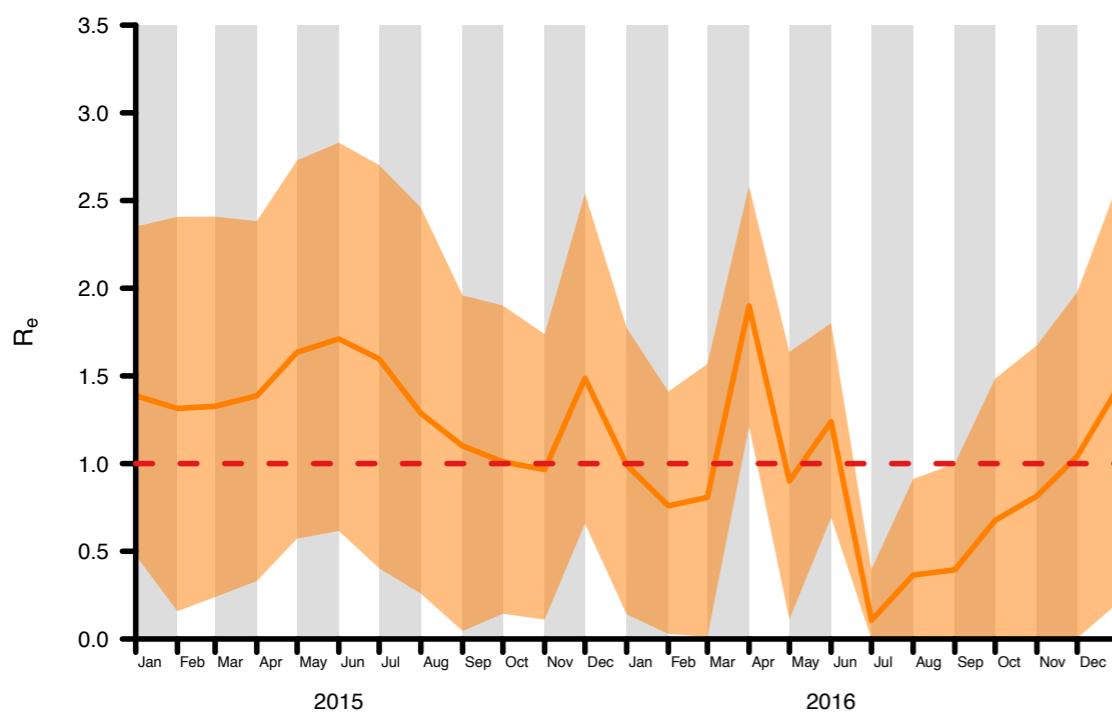
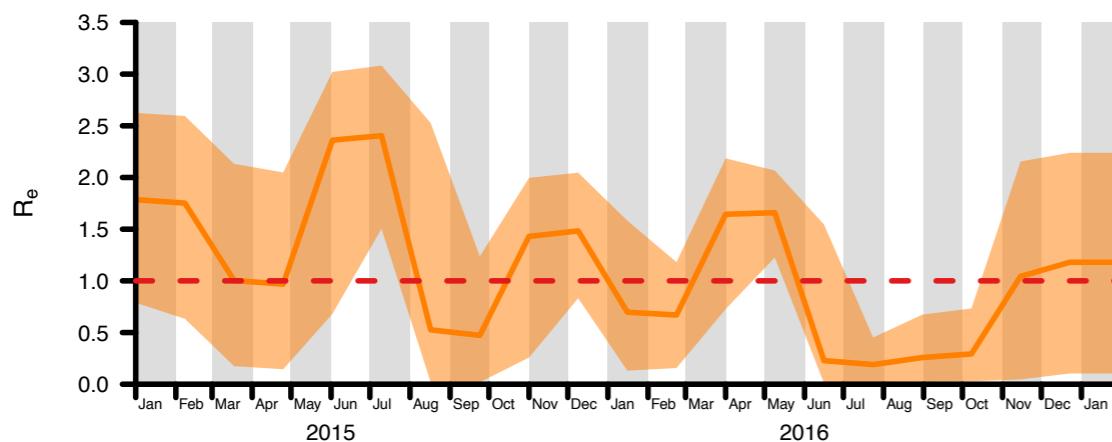
Zika virus in Central America



Zika virus in Central America

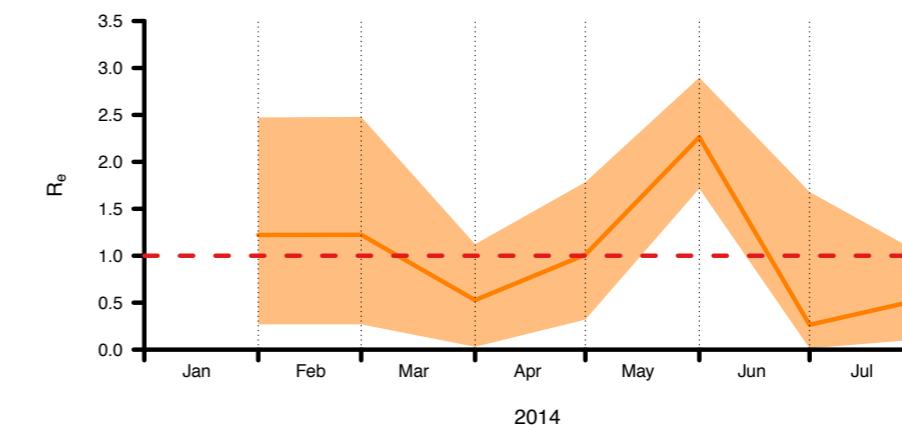
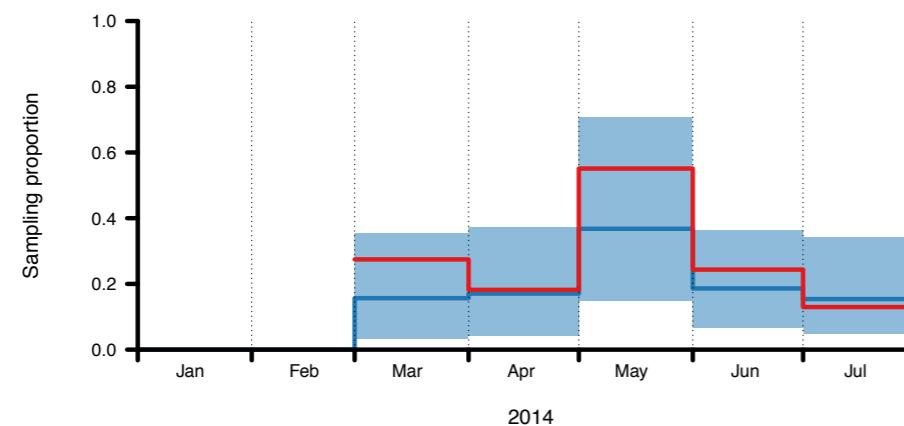


Zika virus in Central America

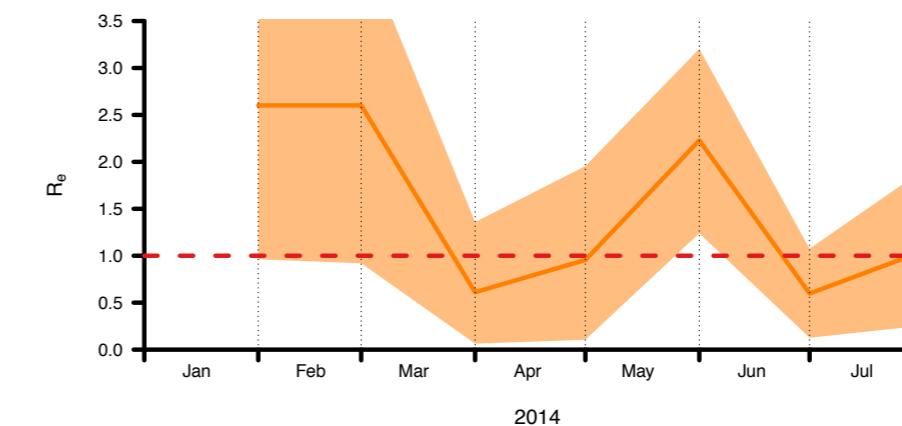
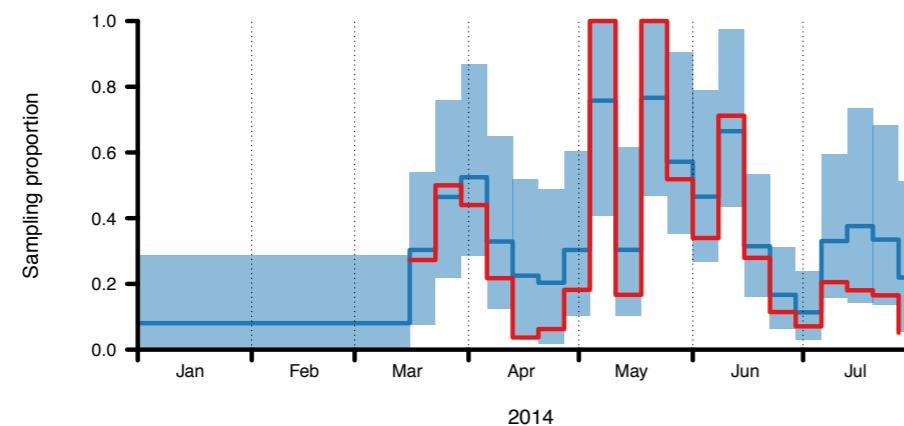


West African Ebola early outbreak

Monthly



Weekly



No linear
model

