

# Phylodynamics in a world of rapidly adapting pathogens

David Rasmussen Dept. of Entomology and Plant Pathology Bioinformatics Research Center NC State University

7<sup>th</sup> Annual Disease Modeling Symposium Bellevue, WA April 15th, 2019

## Introduction

• **Phylodynamics:** the study of how ecological and evolutionary processes act or interact to shape the phylogenetic history of pathogens (Grenfell *et al.*, 2004; Volz *et al.*, 2014)

## Introduction

• **Phylodynamics:** the study of how ecological and evolutionary processes act or interact to shape the phylogenetic history of pathogens (Grenfell *et al.*, 2004; Volz *et al.*, 2014)

 Phylodynamic inference: the statistical practice of inferring ecological/evolutionary dynamics from phylogenetic trees.

## HIV in rural KwaZulu-Natal





## HIV in rural KwaZulu-Natal





## Phylodynamic model for HIV in KZN



#### Rasmussen et al., Virus Evolution (2018)

## Phylodynamic estimates for HIV in KZN



Rasmussen et al., Virus Evolution (2018)

## Phylodynamic estimates for HIV in KZN



Rasmussen et al., Virus Evolution (2018)

What current phylodynamic methods do well:

 Accurately reconstruct historical and recent epidemic dynamics

 Accommodate geographic and other forms of host population structure – allowing us to infer sources of transmission

 Account for incomplete or biased sampling of sequence data

## What current methods do not do well:

 Consider non-neutral genetic variation in pathogen fitness and therefore differences in the epidemic potential of pathogen lineages.

## The big assumption of phylogenetic models



## The big assumption of phylogenetic models



## The independence assumption

 This independence assumption allows us to factor the joint likelihood of a tree T and sequence data S into terms we can easily compute:

 $L(\mathcal{S}, \mathcal{T} | \mu, \theta) = L(\mathcal{S} | \mathcal{T}, \mu) p(\mathcal{T} | \theta).$ 

## Example: deleterious mutation load



• Allows us to compute the joint likelihood that a tree and genotype data at a single loci evolved under a non-neutral model (Stadler & Bonhoeffer, 2013).





Tanja Stadler











 At a single evolving site, we can compute the joint likelihood of a tree and the 'sequence' at each tip using a multi-type birth-death model (Stadler & Bonhoeffer, 2013).



 At a single evolving site, we can compute the joint likelihood of a tree and the 'sequence' at each tip using a multi-type birth-death model (Stadler & Bonhoeffer, 2013).



## The inevitable problem...

• We need to track all possible genotypes in the state space of the model, which increases exponentially with the number of sites L (e.g.  $4^{L}$  for a nucleotide model)

 MTBD becomes prohibitively computationally expensive for anything more than just a few evolving sites.

• We track molecular evolution at each site, computing the marginal site probability  $\omega$  that a site is in particular state.

- We track molecular evolution at each site, computing the marginal site probability  $\omega$  that a site is in particular state.
- We approximate the probability of a lineage being in any genotype based on the marginal site probabilities; e.g.:

$$\hat{\omega}_{n,\text{ACT}} = \omega_{n,1,\text{A}} \times \omega_{n,2,\text{C}} \times \omega_{n,3,\text{T}}$$

- We track molecular evolution at each site, computing the marginal site probability  $\omega$  that a site is in particular state.
- We approximate the probability of a lineage being in any genotype based on the marginal site probabilities; e.g.:

$$\hat{\omega}_{n,\text{ACT}} = \omega_{n,1,\text{A}} \times \omega_{n,2,\text{C}} \times \omega_{n,3,\text{T}}$$

 We then sum, or marginalize, over the fitness of each genotype weighted by its approximate genotype probability to compute the expected fitness of a lineage:

$$f_n \approx \sum_{g \in \mathcal{G}} f_g \hat{\omega}_{n,g}$$

• We now have a new system of ODEs for tracking the probability  $D_{n,k,i}$  that a lineage evolved exactly as observed at each site k:

$$\frac{d}{dt}D_{n,k,i}(t) = -\left(\hat{f}_{n,k,i}\lambda_0 + \sum_{j=1}^M \gamma_{i,j} + d\right)D_{n,k,i}(t)$$
$$+ 2\hat{f}_{n,k,i}\lambda_0E_u(t)D_{n,k,i}(t)$$
$$+ \sum_{j=1}^M \gamma_{i,j}D_{n,k,j}(t).$$

 We now have a new system of ODEs for tracking the probability that a lineage evolved exactly as observed at each site k:

$$\frac{d}{dt}D_{n,k,i}(t) = -\left(\hat{f}_{n,k,i}\lambda_0 + \sum_{j=1}^M \gamma_{i,j} + d\right)D_{n,k,i}(t)$$
$$+ 2\hat{f}_{n,k,i}\lambda_0E_u(t)D_{n,k,i}(t)$$
$$+ \sum_{j=1}^M \gamma_{i,j}D_{n,k,j}(t).$$

 The important part: This allows us to consider how selection shapes sequence evolution at multiple sites while considering how mutations act together to shape the fitness of a lineage.

## Results: quantifying site-specific effects



		GP1																	GP2			
_	_	۲ľ	_	ş <del>-ş</del> ş				Y WWY,														
S	SP			RBD				Glycan Cap				MLD							HR1 HR2 TM			
GF varian	, t 2	29	82	107	202	206	230	239	291	330	371	375	407	410	439	480	485	637		Infection Normalized to	vity Makona (%)	
C1	5 F	R	Α	Ν	Ρ	т	т	L	w	Ρ	Т	Ρ	н	R	κ	G	т	D		0 100 200	300 400	
A									R										W291R	Ma C15		7
A'2	2									S									P330S	A1	*	*
A'3	3									S						D			P330S, G480D	A'2		
A'4	1 I			D						S									N107D, P330S	A'3		
A	5			D						S						D			N107D, P330S, G480D			
A	5									S			Y			G			P330S, H407Y, D480G	A6		*
B			v																A82V	B1	*:	*
B	2		v								v								A82V. 1371V	B2	<b>-</b> *:	* **
AB'	3		Ā								v								1371V	AB'3		
B4	1		۷		L														A82V, P202L	B4	n.:	S   *
B'5	5		v					s											A82V, L239S	BS	<b>-</b> *:	* **
B	6		V		L			S											A82V, P202L, L239S		-	
B	7 1	к	v																R29K, A82V	В7	<b>→</b> **;	* **
B	3		v											s					A82V, R410S			* **
B			۷											s	Е				A82V, R410S, K439E	Ba	- *	*   **
B1(			v									s							482V P375S	B10		* **
AB'1'			Δ									s							P375S	AB'11		
			<u> </u>									Ū								D42		
B12	2		V			М													A82V, T206M	B12	<b>⊣</b> **:	* **
B13	3		v				Α												A82V, T230A	B13	- *:	* **
B14	1		v				Α										Α		A82V, T230A, T485A	B14	*	* **
B'1	5		v				Т										Α		A82V. T485A	B'15		
B16	6		v				Α										Т	G	A82V, T230A, D637G	B16	- *	*   *
			-															-				

Urbanowicz et al. (Cell, 2016)





## Current and ongoing work

• Marginal Fitness Birth Death model is implemented in *LUMIERE*, a package for BEAST2.





 High performance phylodynamic inference using Generalized Birth-Death Models

## Generalized Birth-Death Models

• We would like to be able to learn how pathogen traits map to birth-death fitness parameters:



## Can we have the best of both worlds?

 Can we perform likelihood-based phylodynamic inference under birth-death models

 While using the tools of machine learning to learn how pathogen features (e.g. genotypes) map to population-level parameters?





 Trees can be represented as computational graphs in TensorFlow with data arrays (i.e. *tensors*) flowing between nodes



## Phylodynamics in TensorFlow



#### Deep mutational scanning of hemagglutinin helps predict evolutionary fates of human H3N2 influenza variants

 $\label{eq:Leebard} \begin{array}{l} \text{Juhye M. Lee}^{\mathrm{a,d,e,1}}, \text{John Huddleston}^{\mathrm{b,f,1}}, \text{Michael B. Doud}^{\mathrm{a,d,e}}, \text{Kathryn A. Hooper}^{\mathrm{a,f}}, \text{Nicholas C. Wu}^{\mathrm{g}}, \text{Trevor Bedford}^{\mathrm{b,c,1}}, \\ \text{and Jesse D. Bloom}^{\mathrm{a,c,d,1}} \end{array}$ 

<sup>a</sup> Basic Sciences Division; <sup>b</sup> Vaccine and Infectious Disease Division; <sup>c</sup> and Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA; <sup>d</sup> Department of Genome Sciences; <sup>e</sup>Medical Scientist Training Program; <sup>r</sup> and Molecular and Cellular Biology Program, University of Washington, Seattle, WA 98195, USA; <sup>g</sup> Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla, CA 92037



#### Fitness model for H3N2 influenza



Phylodynamics vs. Deep Mutational Scanning

• Very small fitness effects once we account for seasonal fluctuations in flu transmission



#### A big thank you to:

Prof. Tanja Stadler



For the South African HIV project:

Tulio de Oliviera Eduan Wilkinson Africa Health Research Institute



The Phylodynamics Research Group at NC State



phylodynamics.wordpress.ncsu.edu

#### For funding:



Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich





European Research Council

 We now have a new system of ODEs for tracking the probability that a lineage evolved exactly as observed at each site k:

$$\frac{d}{dt}D_{n,k,i}(t) = -\left(\hat{f}_{n,k,i}\lambda_0 + \sum_{j=1}^M \gamma_{i,j} + d\right)D_{n,k,i}(t)$$
$$+ 2\hat{f}_{n,k,i}\lambda_0E_u(t)D_{n,k,i}(t)$$
$$+ \sum_{j=1}^M \gamma_{i,j}D_{n,k,j}(t).$$

 Note 1: We can track evolution at each site individually without tracking all genotypes.

 We now have a new system of ODEs for tracking the probability that a lineage evolved exactly as observed at each site k:

$$\frac{d}{dt}D_{n,k,i}(t) = -\left(\hat{f}_{n,k,i}\lambda_0 + \sum_{j=1}^M \gamma_{i,j} + d\right)D_{n,k,i}(t)$$
$$+ 2\hat{f}_{n,k,i}\lambda_0E_u(t)D_{n,k,i}(t)$$
$$+ \sum_{j=1}^M \gamma_{i,j}D_{n,k,j}(t).$$

 Note 2: We can simultaneously take into account the coupled fitness effects of mutations at all other sites.

 We now have a new system of ODEs for tracking the probability that a lineage evolved exactly as observed at each site k:

$$\frac{d}{dt}D_{n,k,i}(t) = -\left(\hat{f}_{n,k,i}\lambda_0 + \sum_{j=1}^M \gamma_{i,j} + d\right)D_{n,k,i}(t)$$
$$+ 2\hat{f}_{n,k,i}\lambda_0E_u(t)D_{n,k,i}(t)$$
$$+ \sum_{j=1}^M \gamma_{i,j}D_{n,k,j}(t).$$

• Note 3: Tracking  $D_{n,k,i}$  all the way back to the root allows us to compute the joint likelihood of the tree and the sequence data at site k.

$$L(\mathcal{S}_k, \mathcal{T}|\theta) = \sum_{i=1}^M D_{n,k,i}(t_{root})$$

