

Bayesian Phylogeography

Remco R. Bouckaert

remco@cs.auckland.ac.nz

Centre of Computational Evolution

Department of Computer Science, University of Auckland

&

Max Planck Institute for the Science of Human History

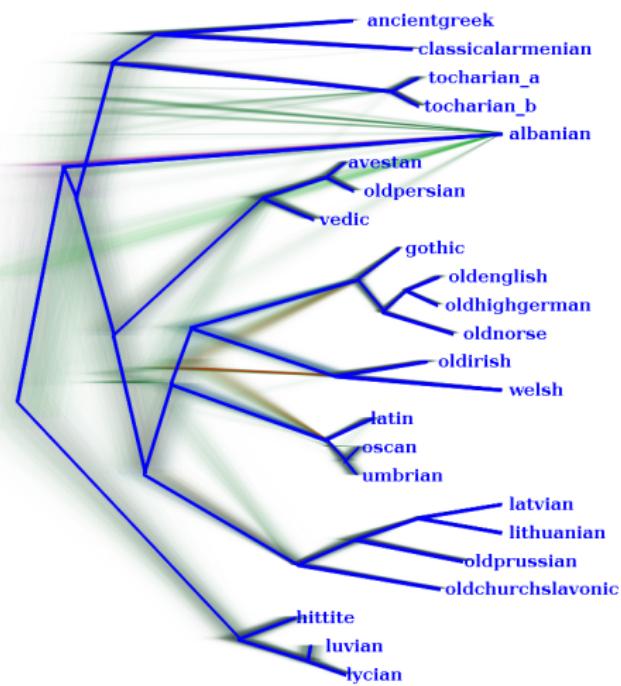
Seili, September 2017



Bayesian phylogenetics

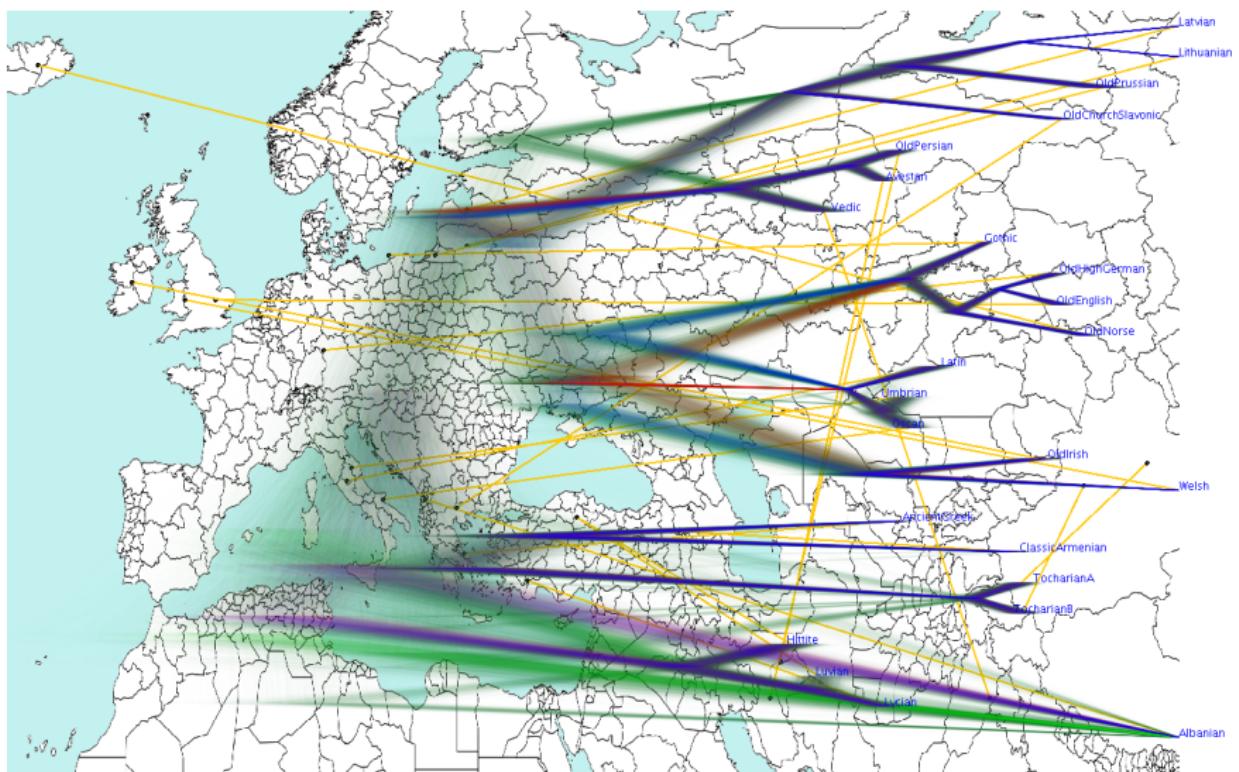
- Posterior \propto prior \times likelihood
- Prior info easy to incorporate
- Framework for adding sources of information from different sources – geography, WALS, cognate data, etc.
- Think in distributions: there is no “The Tree”, only a distribution over trees

Ring DensiTree

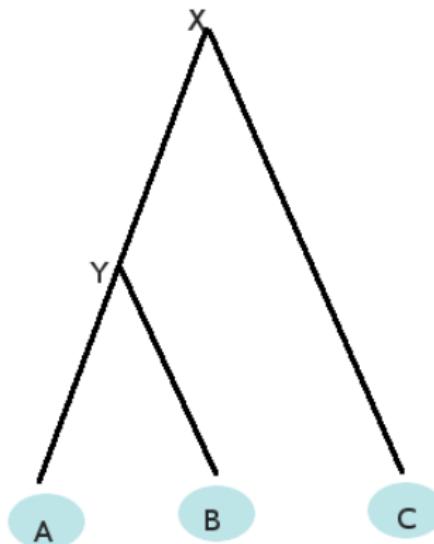


Basic Phylogeography

Link tree tips to locations + story telling

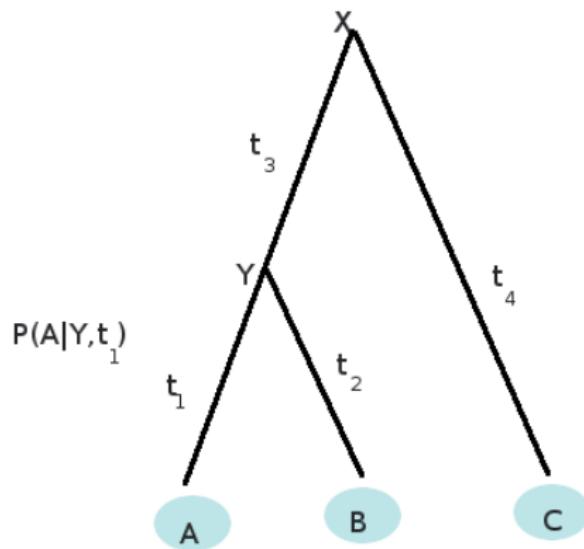


Phylogeography conditioned on tree



Locations: A, B, C, X, Y

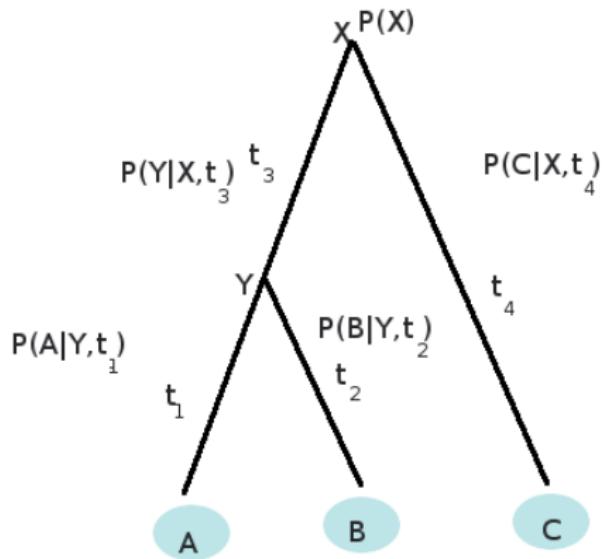
Phylogeography conditioned on tree



Locations: A, B, C, X, Y

Branch lengths: t_1, t_2, t_3, t_4

Phylogeography conditioned on tree



Locations: A, B, C, X, Y

Branch lengths: t_1, t_2, t_3, t_4

Probabilities: $\forall x \in \{A, B, C, Y\} P(x | \text{parent}(x), t_x)$ and $P(X)$

Ancestral state reconstruction (BEAST_CLASSIC package)

- requires specifying regions
- many parameters ($O(n^2)$ with n states)

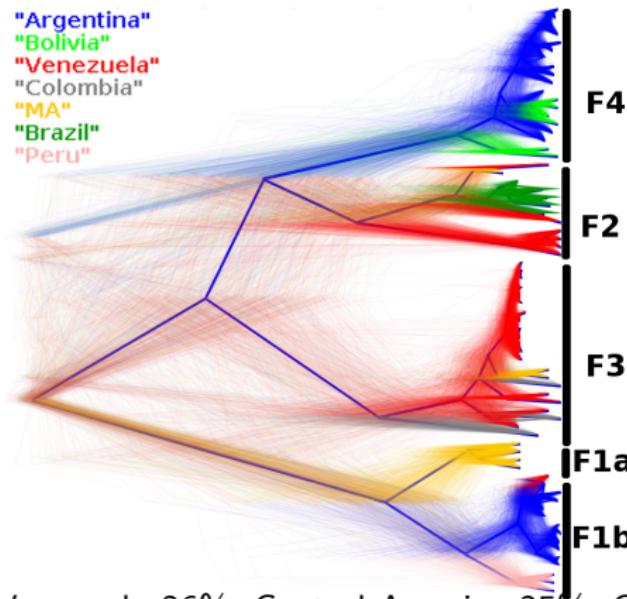
	Argentina	Bolivia	Brazil	Central America	Colombia	Peru	Venezuela
Argentina	—	$r_{A \rightarrow Bo}$	$r_{A \rightarrow Br}$	$r_{A \rightarrow CA}$	$r_{A \rightarrow Co}$	$r_{A \rightarrow Pe}$	$r_{A \rightarrow Ve}$
Bolivia	$r_{Bo \rightarrow A}$	—	$r_{Bo \rightarrow Br}$	$r_{Bo \rightarrow CA}$	$r_{Bo \rightarrow Co}$	$r_{Bo \rightarrow Pe}$	$r_{Bo \rightarrow Ve}$
Brazil	$r_{Br \rightarrow A}$	$r_{Br \rightarrow Bo}$	—	$r_{Br \rightarrow CA}$	$r_{Br \rightarrow Co}$	$r_{Br \rightarrow Pe}$	$r_{Br \rightarrow V}$
Central America	$r_{CA \rightarrow A}$	$r_{CA \rightarrow Bo}$	$r_{CA \rightarrow Br}$	—	$r_{CA \rightarrow Co}$	$r_{CA \rightarrow Pe}$	$r_{CA \rightarrow V}$
Colombia	$r_{Co \rightarrow A}$	$r_{Co \rightarrow Bo}$	$r_{Co \rightarrow Br}$	$r_{Co \rightarrow CA}$	—	$r_{Co \rightarrow Pe}$	$r_{Co \rightarrow V}$
Peru	$r_{Pe \rightarrow A}$	$r_{Pe \rightarrow Bo}$	$r_{Pe \rightarrow Br}$	$r_{Pe \rightarrow CA}$	$r_{Pe \rightarrow Co}$	—	$r_{Pe \rightarrow V}$
Venezuela	$r_{V \rightarrow A}$	$r_{V \rightarrow Bo}$	$r_{V \rightarrow Br}$	$r_{V \rightarrow CA}$	$r_{V \rightarrow Co}$	$r_{V \rightarrow Pe}$	—

transition probabilities $P(t) = e^{Qt}$

so use Bayesian Stochastic Variable Selection (BSVS)

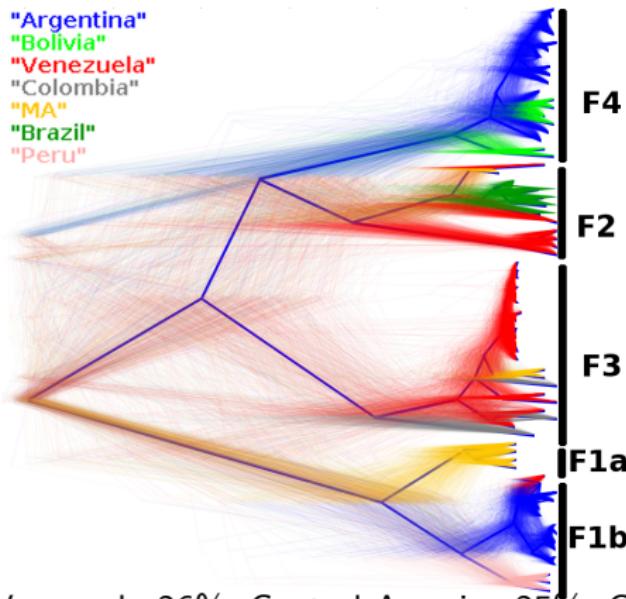
- breaks structured coalescent model => biased estimates
- weak – 95% HPD of root location often contains all states

Ancestral state reconstruction – HBV-F in South America



- Root location Venezuela 26%, Central America 25%, Colombia 18%, Argentina 12%, Bolivia 6%, Brazil 6%, Peru 6%

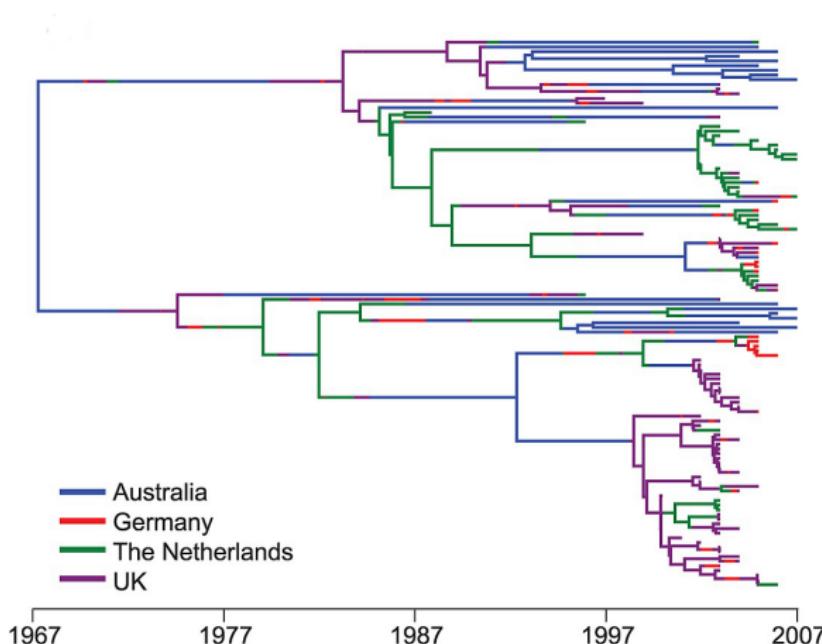
Ancestral state reconstruction – HBV-F in South America



- Root location Venezuela 26%, Central America 25%, Colombia 18%, Argentina 12%, Bolivia 6%, Brazil 6%, Peru 6%
- Initially HBV spread from the North into South America
- Argentina obtained from Central America or Venezuela, not the closer by Brazil or Bolivia and HBV spread from Argentina to West of South America
- Brazil obtained HBV from Central America

Structured Coalescent (MultiTypeTree package)

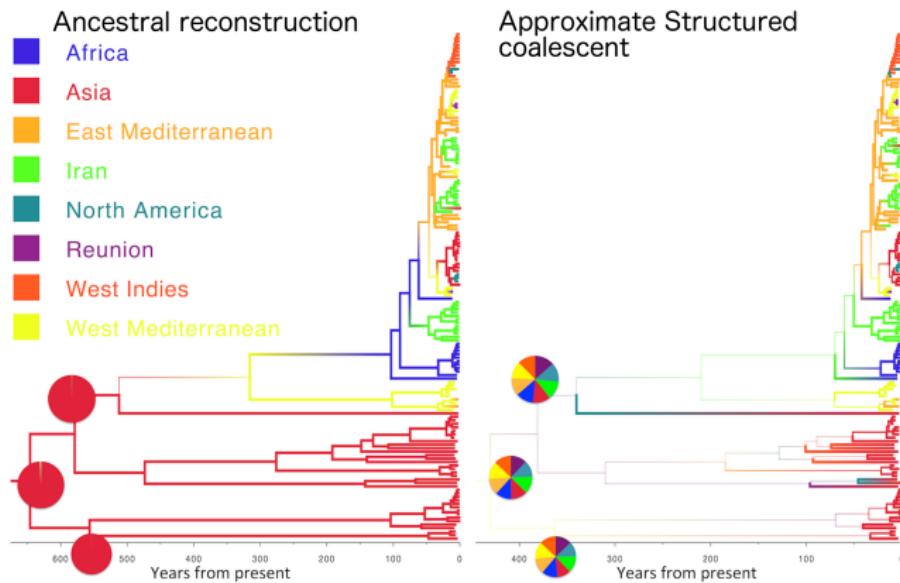
Requires specifying regions (demes)



Principled way of migrating/coalescence
Does not scale (max 4 demes)

Approximate structured coalescent (BASTA package)

- Scales slightly better (max 12? demes) than structured coalescent
- Results differ significantly from ancestral reconstruction



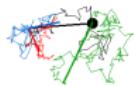
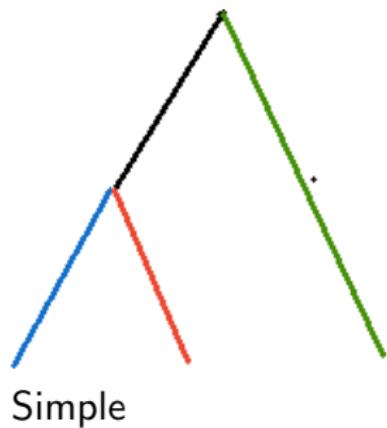
Tomato Yellow Leaf Curl Virus

De Maio et al, PLOS Genetics, 2015

Remco Bouckaert (New Zealand)

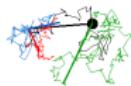
Bayesian Phylogeography

Random walks – types



Random walks – types

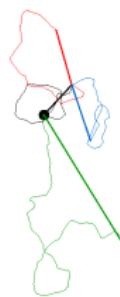
Simple



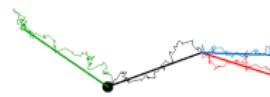
Levy



Correlated



Biased

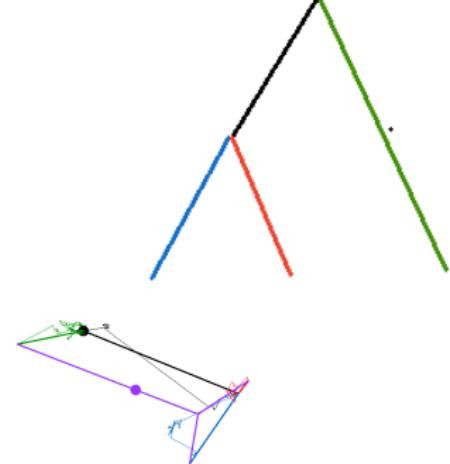


Random walks – tree reconstruction

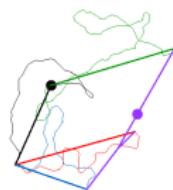
Simple



Levy



Correlated

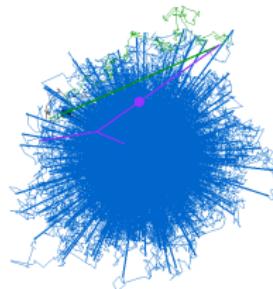


Biased

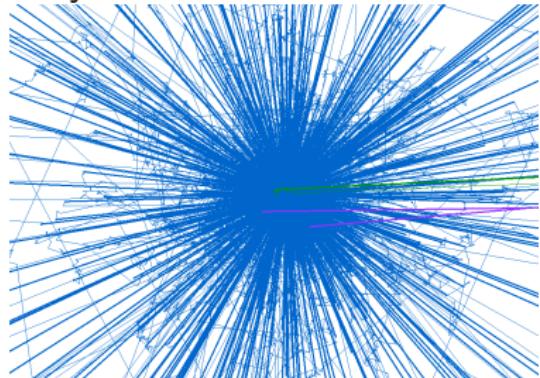


Random walks – distributions

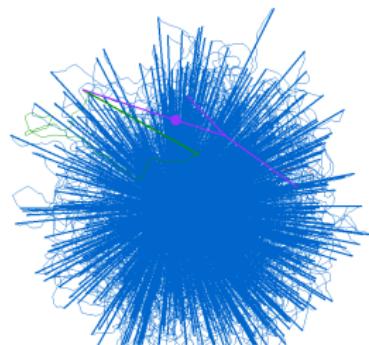
Simple



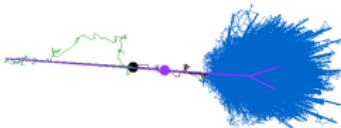
Levy



Correlated



Biased



Random walks

Diffusion on a plane (BEAST_CLASSIC package)

- more power than ancestral reconstruction
- limited landscape awareness
- does not deal with Mercator projection distortion

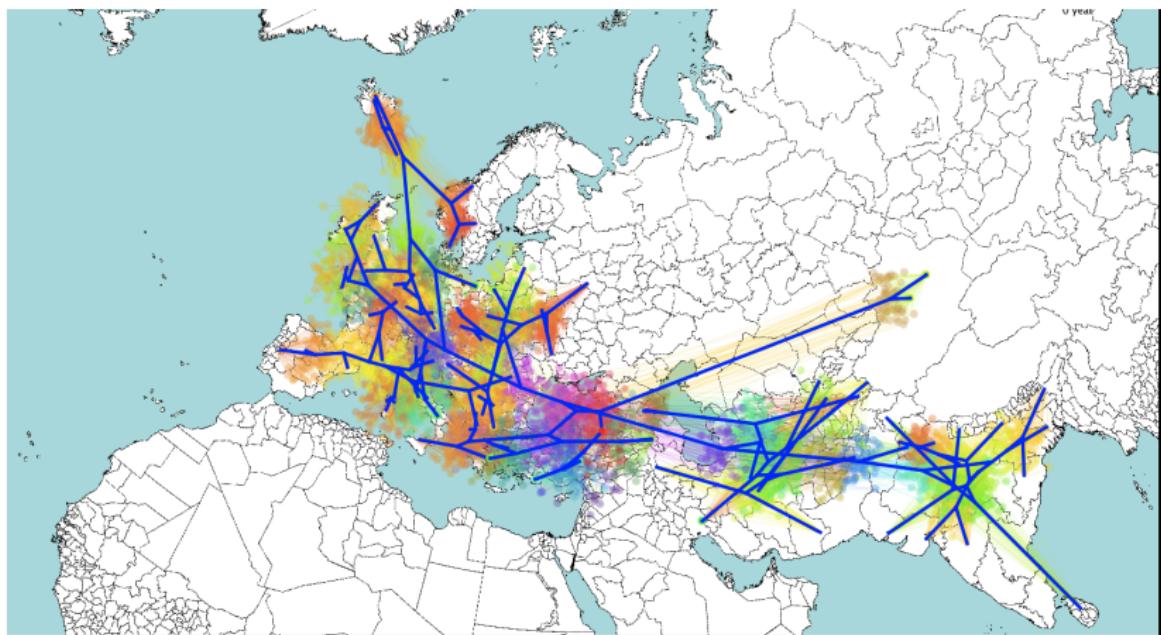
Lemey et al. MBE, 2010 Pybus et al. PNAS, 2012

Diffusion on a sphere (GEO_SPHERE package)

- as above, but deals with Mercator projection distortion

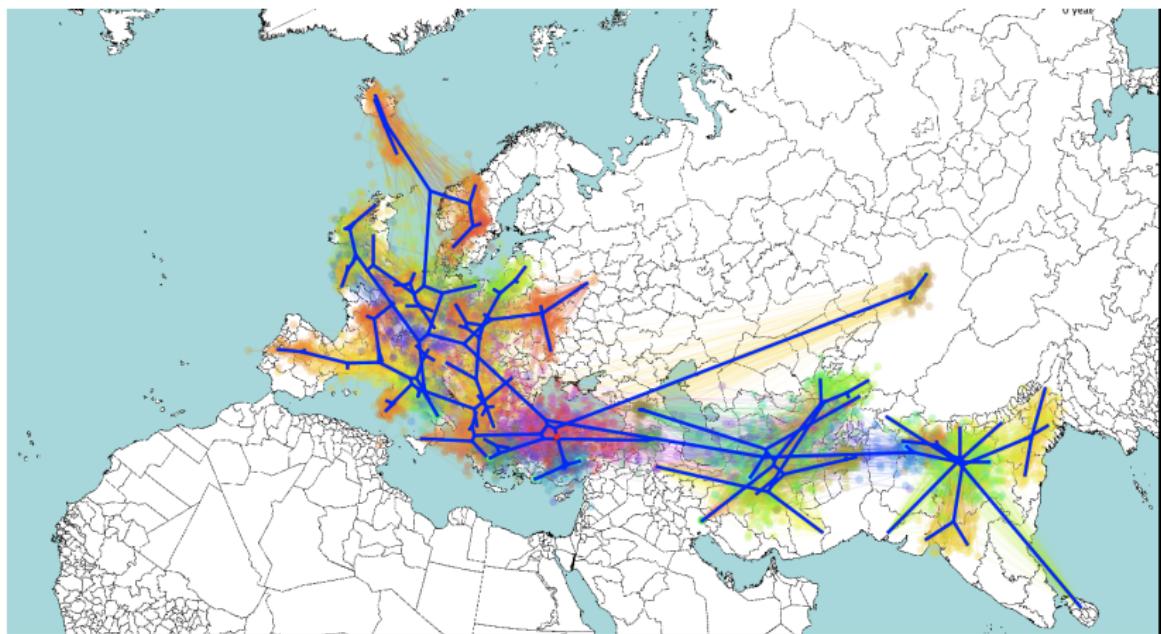
Bouckaert PeerJ 2016

Reconstructing the origin of Indo-European



Simple random walk on sphere, fixed tree, strict clock
GEO_SPHERE package [IEstrict.mp4](#)

Indo European – relaxed random walks*

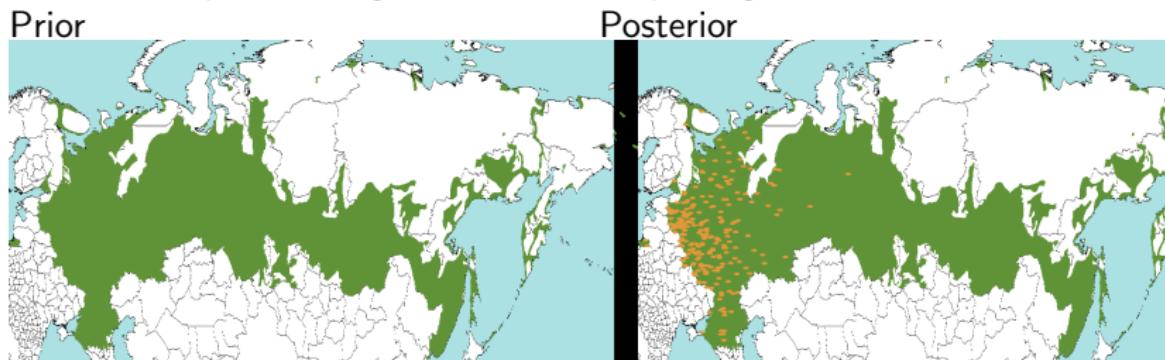


Simple random walk on sphere, fixed tree, relaxed clock, much better fit
GEO_SPHERE package [IERC.mp4](#)

* Lemey et al. MBE, 2010

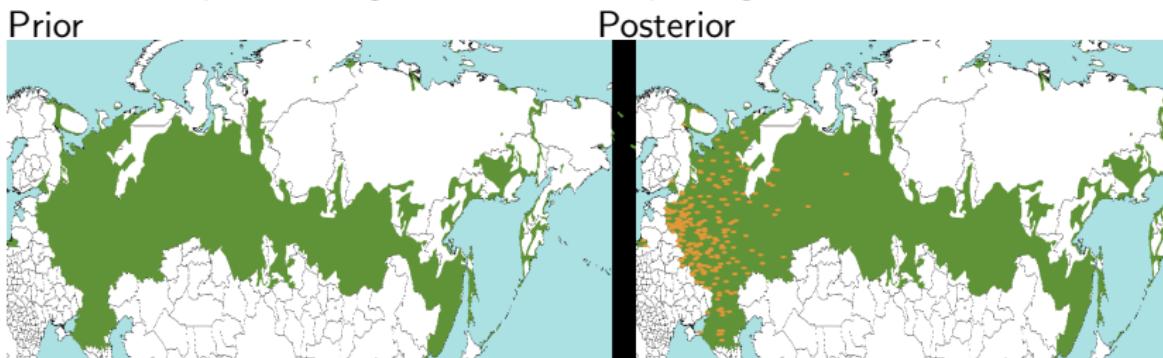
Geographical priors

- root location
- tip location, e.g. Russian, sampled from Russian speaking regions
 - ▶ use uniform prior on region GEO_SPHERE package Bouckaert et al Science, 2012



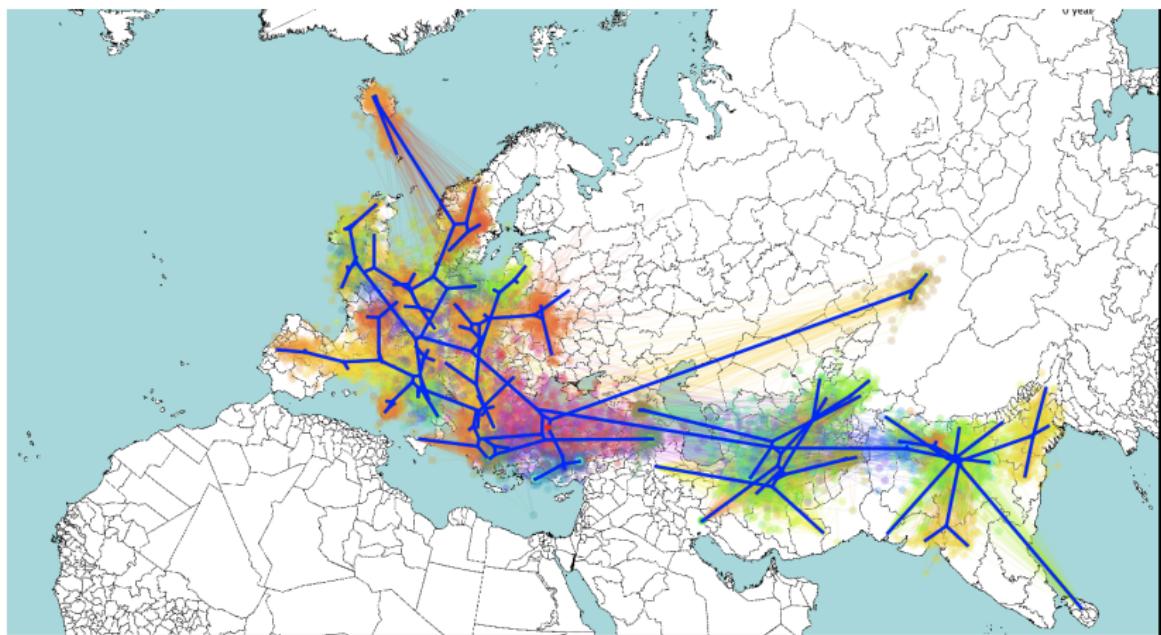
Geographical priors

- root location
 - tip location, e.g. Russian, sampled from Russian speaking regions
 - ▶ use uniform prior on region `GEO_SPHERE` package Bouckaert et al Science, 2012



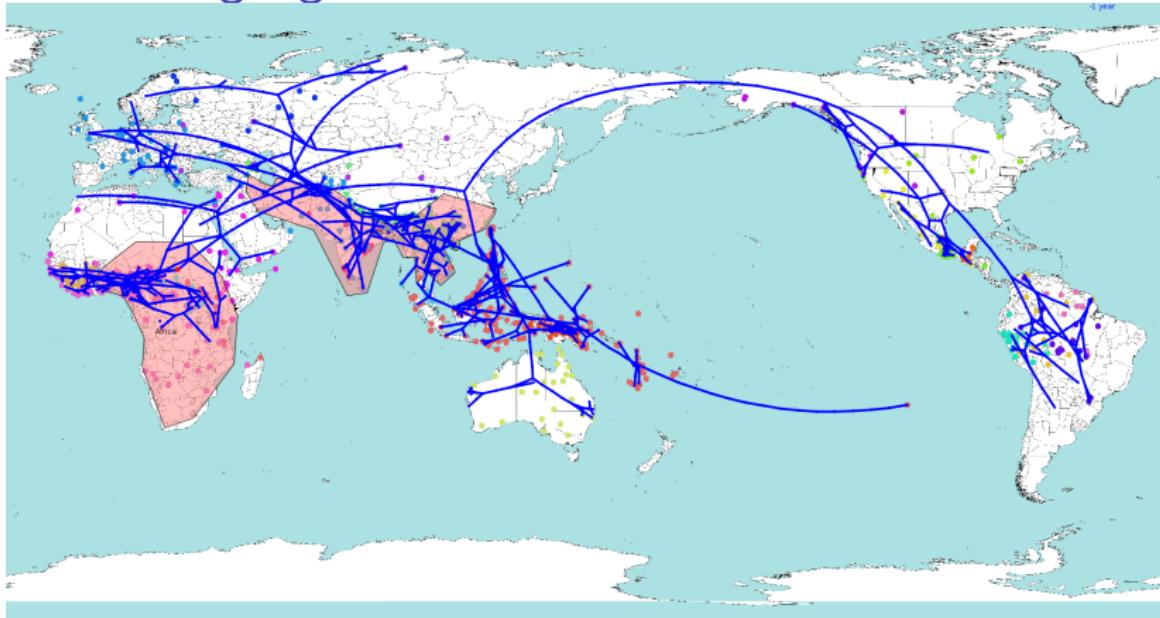
- ▶ use non-uniform prior on region Nylander et al SysBio 2014 BEAST 1
 - ▶ Integrate out regions Quintero et al SysBio 2015 fixed tree, R package only
 - MRCA of clade location, e.g origin of American samples in region around Bering strait
 - restrict all internal nodes, e.g. no node in water

Indo European – Sampling tip locations



Simple random walk on sphere, fixed tree, relaxed clock, more uncertainty
GEO_SPHERE package [IERCtips.mp4](#)

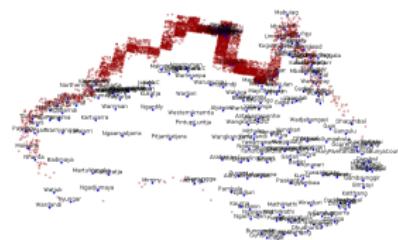
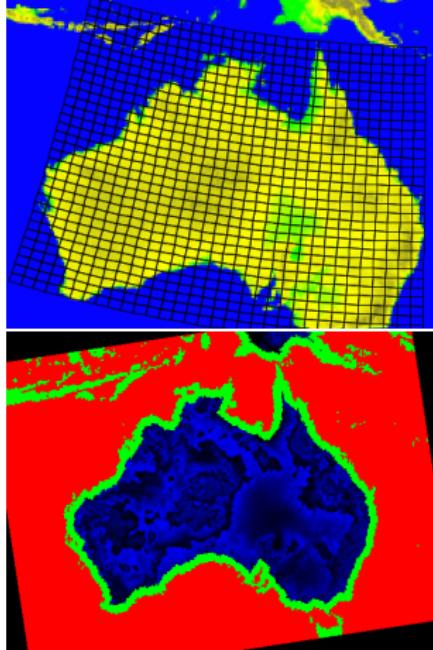
Global language tree



- prior on root location and 'rest of world' origin
- age constraint on all families
- age constraint on root
- glottolog constraints on tree **glotto.mpg**

Landscape aware models

- divide map in a 32x32 grid
 - distinguish between land, water, interior
 - solve ODE on underlying pixel map using Euler's method ([la.mpg](#))
 - relaxed clock
 - allows landscape feature modelling (water, mountains, barriers)
 - computationally expensive
 - cannot sample model parameters during MCMC/needs pre-computation
 - allows nice visualisation ([la2.mpg](#))

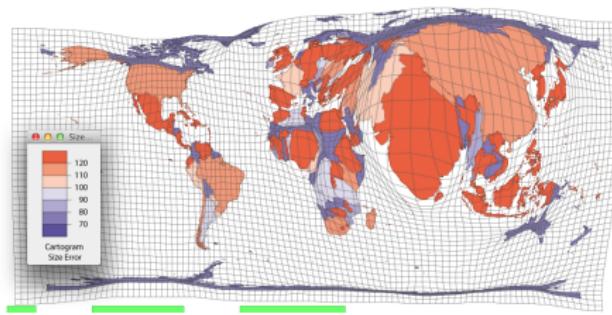
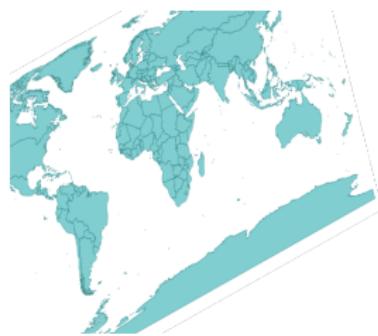


Bouckaert et al. Science, 2012 (package in development)

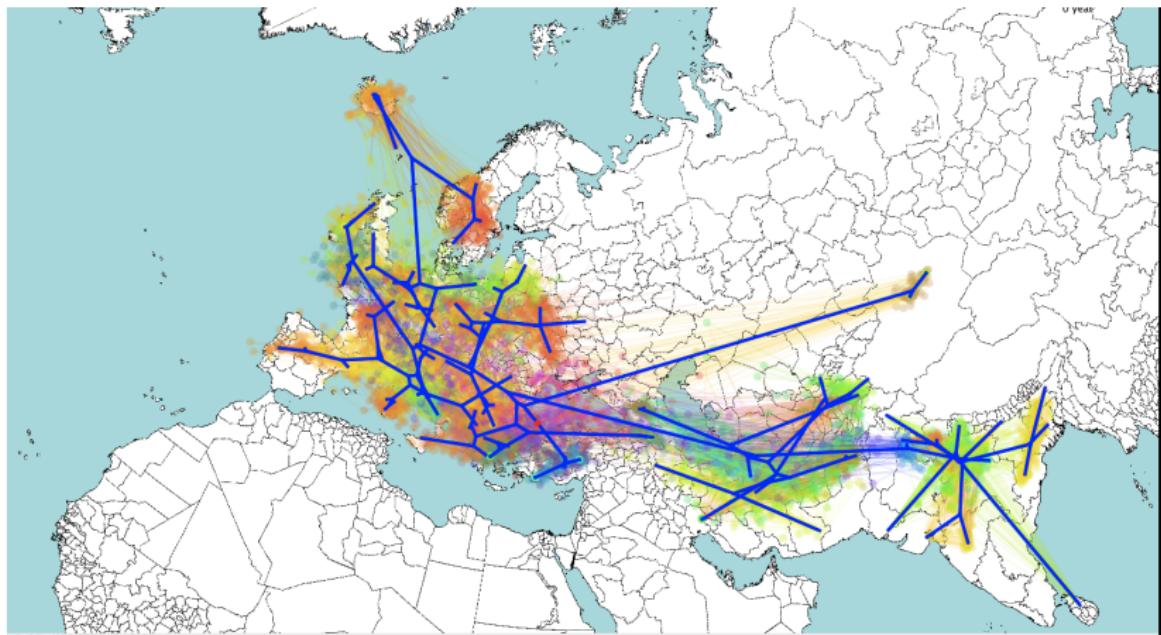
Landscape awareness

By distorting the map + run homogeneous random walk model

- affine transforms
- cartogram



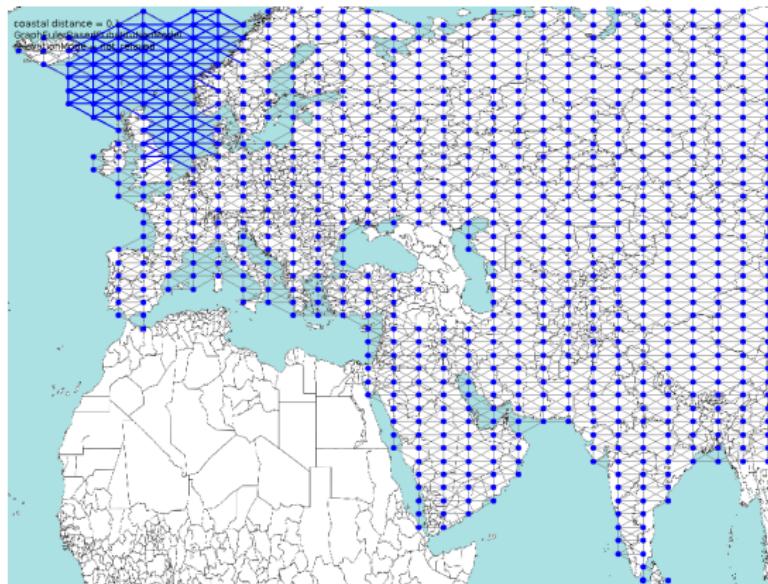
Indo European – population size distorted map



Origin shifts around a slightly IERCtipsTrans.mp4 GEO_SPHERE package

Random walk on graph

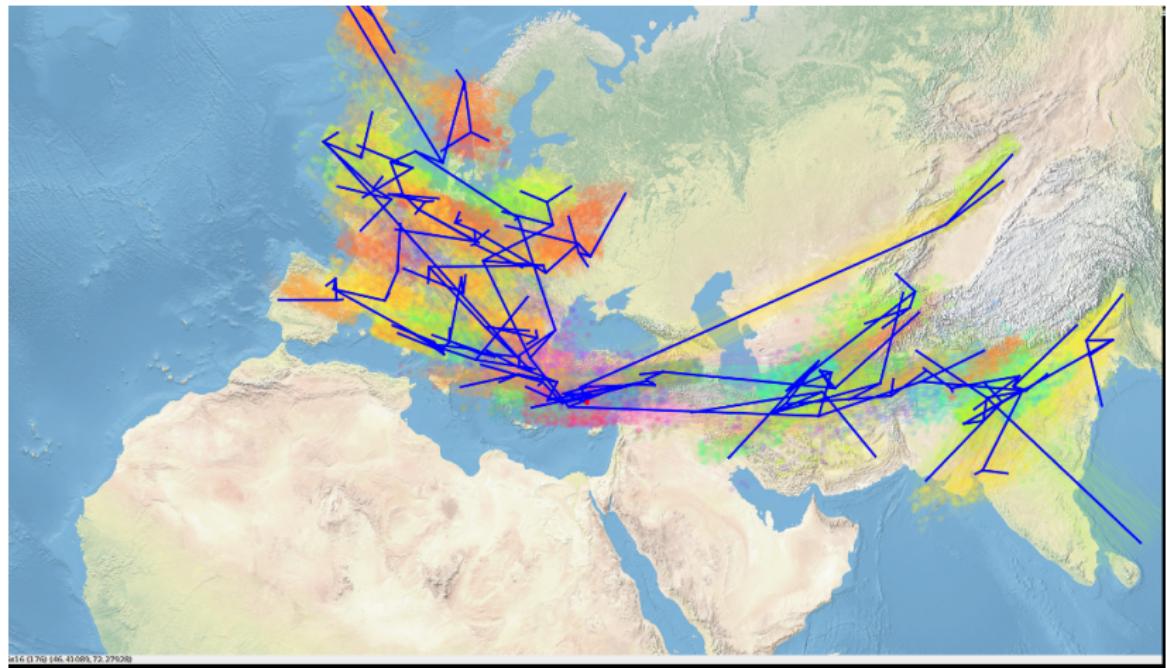
- arbitrary neighbourhood structure
- arbitrary rates
- allow testing hypothesis of barriers, e.g. no migration through Caucasus



show SphereTesselation app

Bouckaert BioArchiv 2016

Indo European – no migration through Caucasus



IEnoCaucasus.mp4

Visualisation

- Summary tree on map connecting internal nodes
 - ▶ by straight lines
 - ▶ by shortest great circle distance lines
 - ▶ MAP route in underlying landscape model **NNA2.mp4**
- Posterior distribution on map
 - ▶ by dots
 - ▶ be DensiTree
 - ▶ by 80%HPD regions **ie.mov**

Methods of inference

- Data augmentation
 - ▶ sample locations L for nodes in tree
 - ▶ $P(L|T, \theta) = \prod_{i \in L} P(i|\text{parent}(i), \theta)$
 - ▶ may have trouble converging
- Integrate out internal node locations (peeling algorithm)
 - ▶ works well with standard distributions (e.g. normal diffusion)
 - ▶ converges fast
 - ▶ can deal with a few constraints by data augmentation for affected nodes
 - ▶ cannot deal constraints on all nodes (e.g. restriction to being on land)
- Particle filtering
 - ▶ flexible in dealing with constraints
 - ▶ slow

Ideal questions

- Precondition: Requires a tree
 - ▶ genetic data
 - ▶ morphological data
 - ▶ anything with evolutionary mechanism

Some questions that can be answered with phylogeography

- What is the origin of a group (95%HPD)?
- Test two (or more) origin hypotheses
What is the Bayes factor in favour of one?
- Reconstruct migration routes
- Estimate migration rates
- Test transition hypotheses (discrete space methods)
- Test landscape heterogeneity hypotheses

Summary

- Discrete space methods
 - ▶ Ancestral reconstruction
 - ▶ (Approximate) Structured Coalescent
- Continuous space diffusion models
 - ▶ various random walks
 - ▶ tip distributions
 - ▶ priors on nodes
 - ▶ somewhat heterogeneous
- Random walk on graph
 - ▶ flexible tip distributions
 - ▶ flexible priors
 - ▶ flexible heterogeneity

+ Strict/Relaxed clocks for all