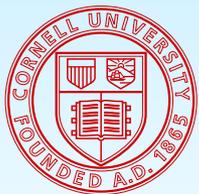
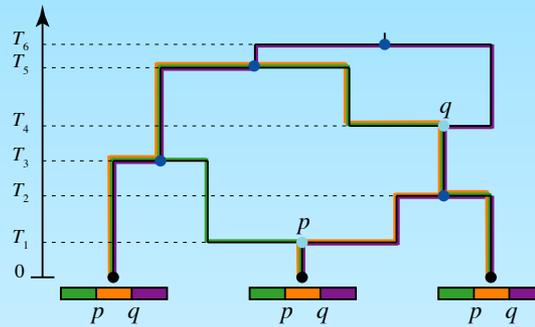
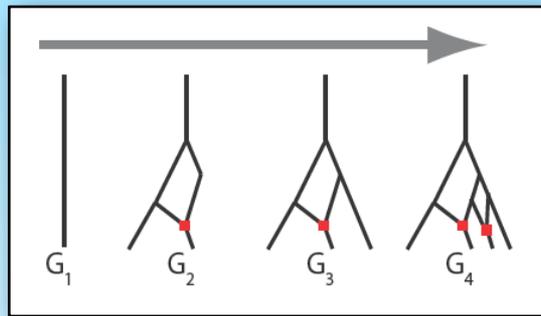


ARGweaver: Genome-wide Inference of Ancestral Recombination Graphs



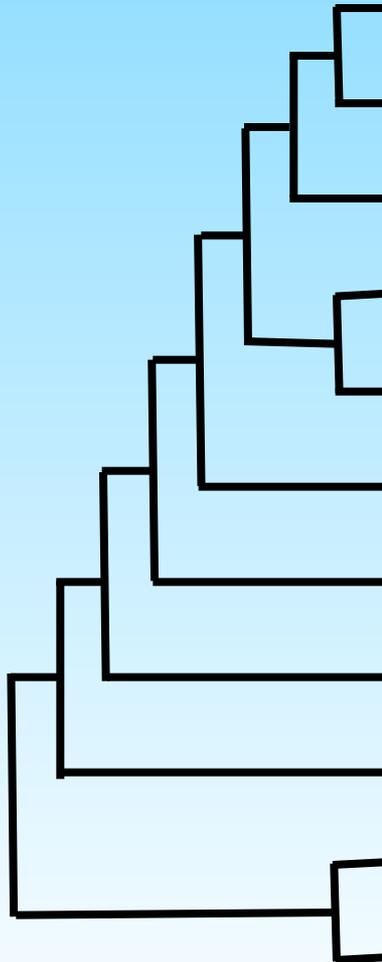
Adam Siepel

Biological Statistics & Computational Biology
Center for Comparative and Population Genomics
Cornell University, Ithaca, NY

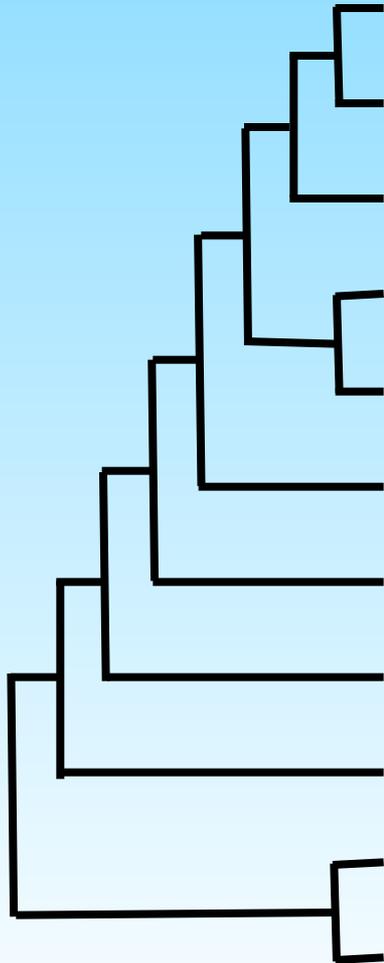


***Joint work with Matthew Rasmussen,
Melissa Hubisz, and Ilan Gronau***

Comparative Genomics



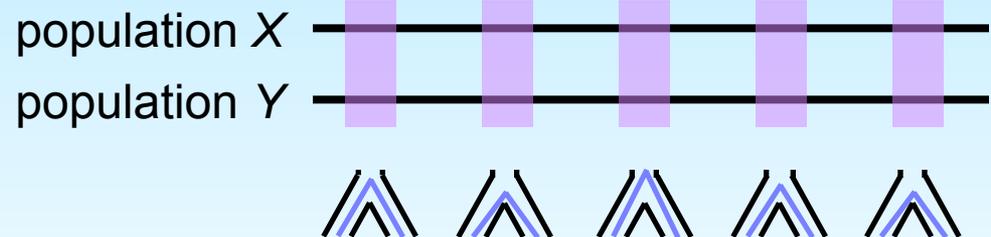
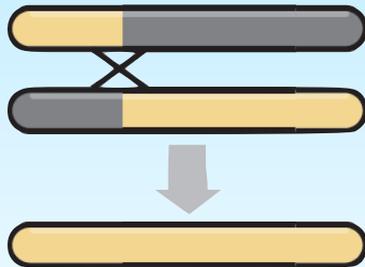
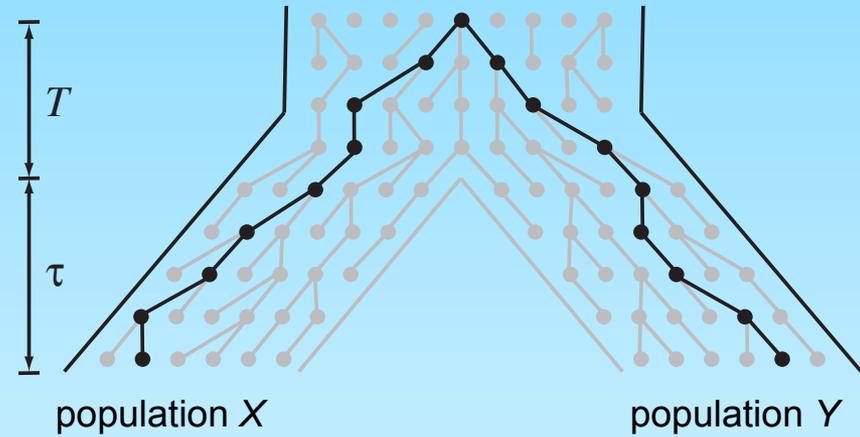
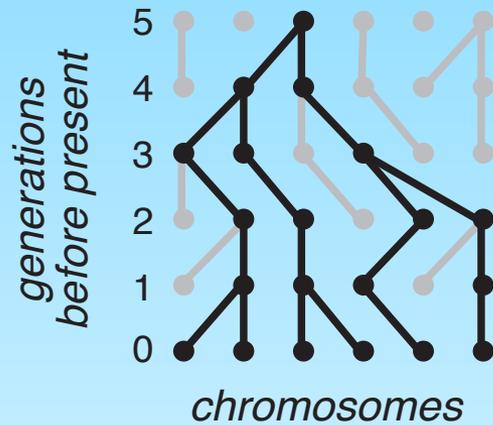
Comparative Genomics



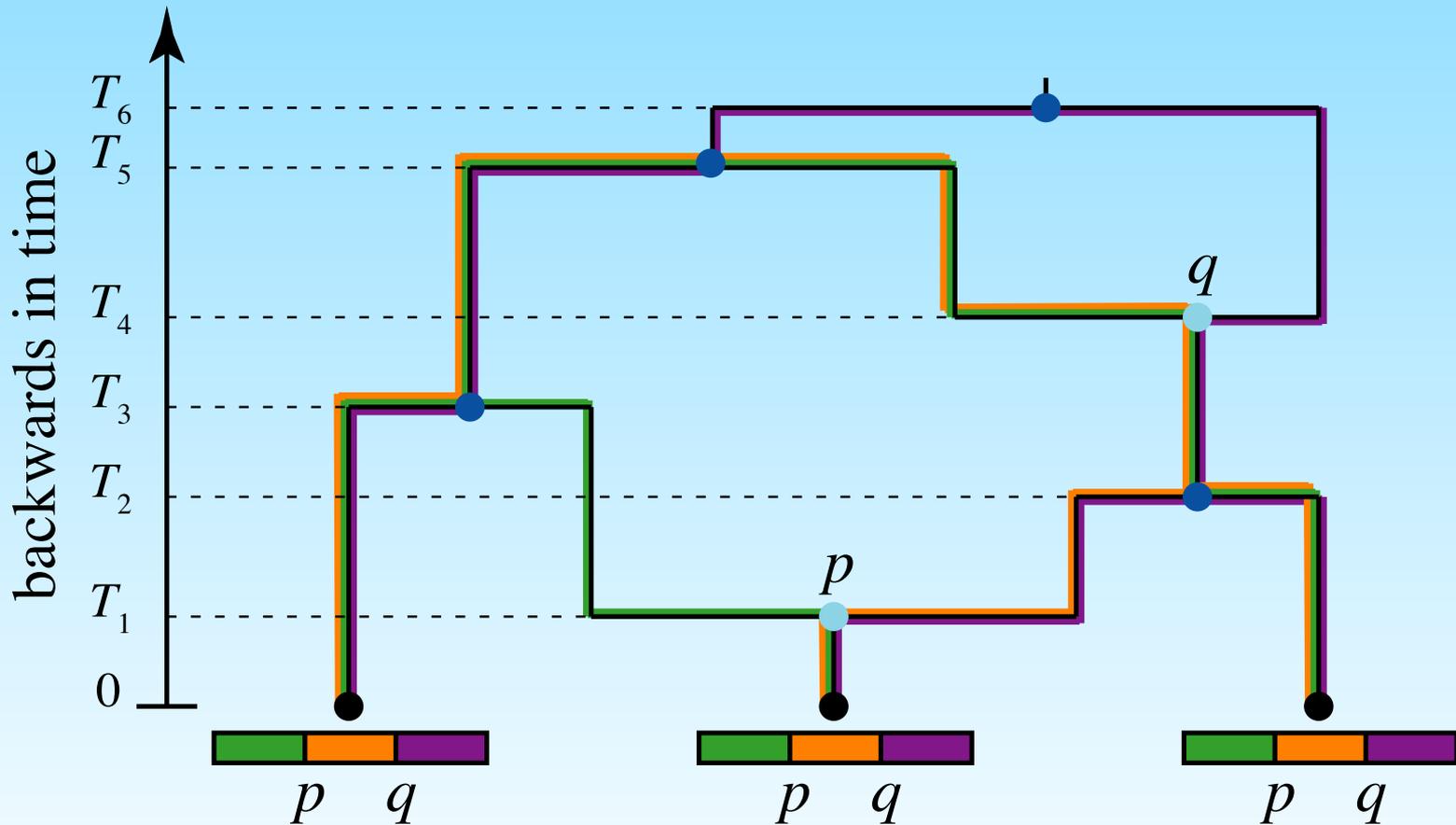
Evolution of Individual Human Genomes



Recombination and Genealogies



The Ancestral Recombination Graph



Alas, if only we knew the ARG...

- Demography inference
- Inference of natural selection
- Recombination rate estimation
- Phasing/imputation
- Association mapping
- ...



Alas, if only we knew the ARG...

- Demography inference
- Inference of natural selection
- Recombination rate estimation
- Phasing/imputation
- Association mapping
- ...

ARG surrogates: IBD, IBS, haplotypes, local ancestry inference, site-frequency spectrum, PCA



Explicit ARG Inference

- **Importance sampling**
 - Griffiths and Marjoram, *J. Comput. Biol.*, 1996
 - Fearnhead and Donnelly, *Genetics*, 2001
- **Markov chain Monte Carlo sampling**
 - Kuhner, Yumato, and Felsenstein, *Genetics*, 2000
 - Nielsen, *Genetics*, 2000
- **Heuristics/Parsimony**
 - Hein, *J. Mol. Evol.*, 1993
 - Kececioglu and Gusfield, *Disc. Appl. Math.*, 1998
 - Song and Hein, *J. Comput. Biol.*, 2005
 - Minichiello and Durbin, *Am. J. Hum. Genet.*, 2006



Explicit ARG Inference

- Importance sampling

- Griffiths and Marjoram,
- Fearnhead and Donnelly,

- Markov chain Monte Carlo sampling

- Quiner, Yumoto, and Felsenstein,
- Nielsen,

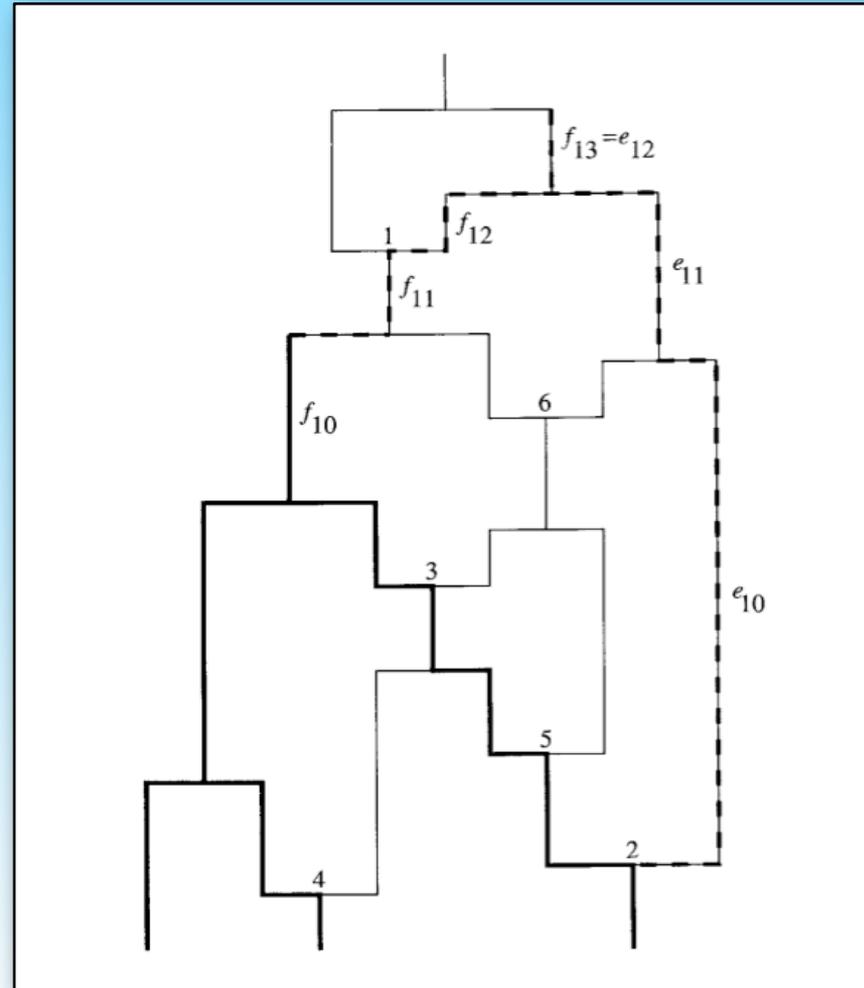
- Heuristics/Parsimony

- He,
- Kececioglu and Gusfield,
- Steel and Hein,
- Minichiello and Durbin,

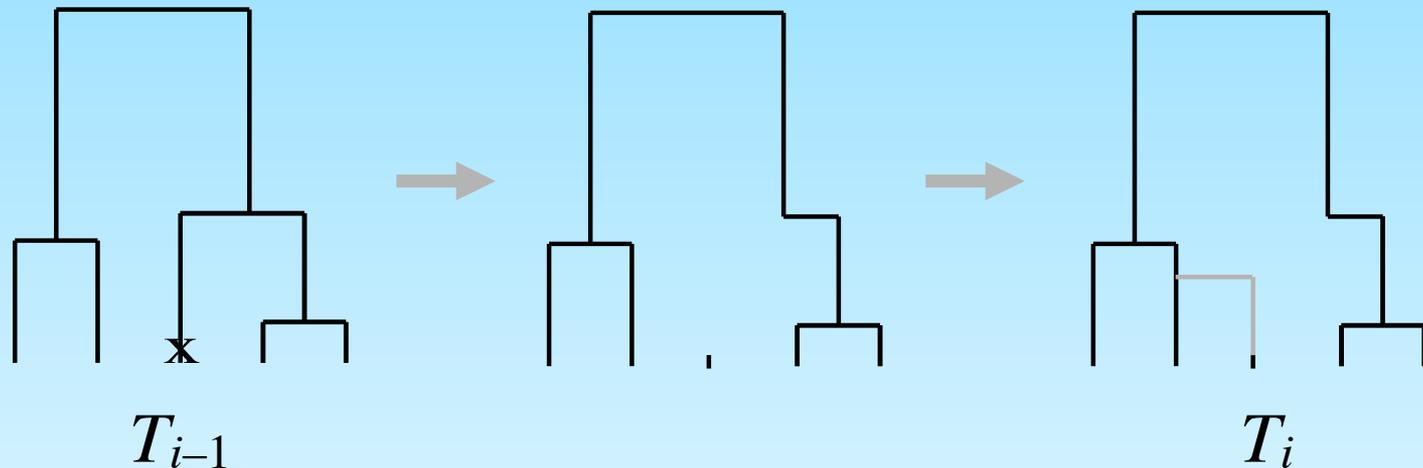
**Computationally intensive
and/or
Limited to few samples
Depend on crude approximations**



Sequential Coalescent with Recombination



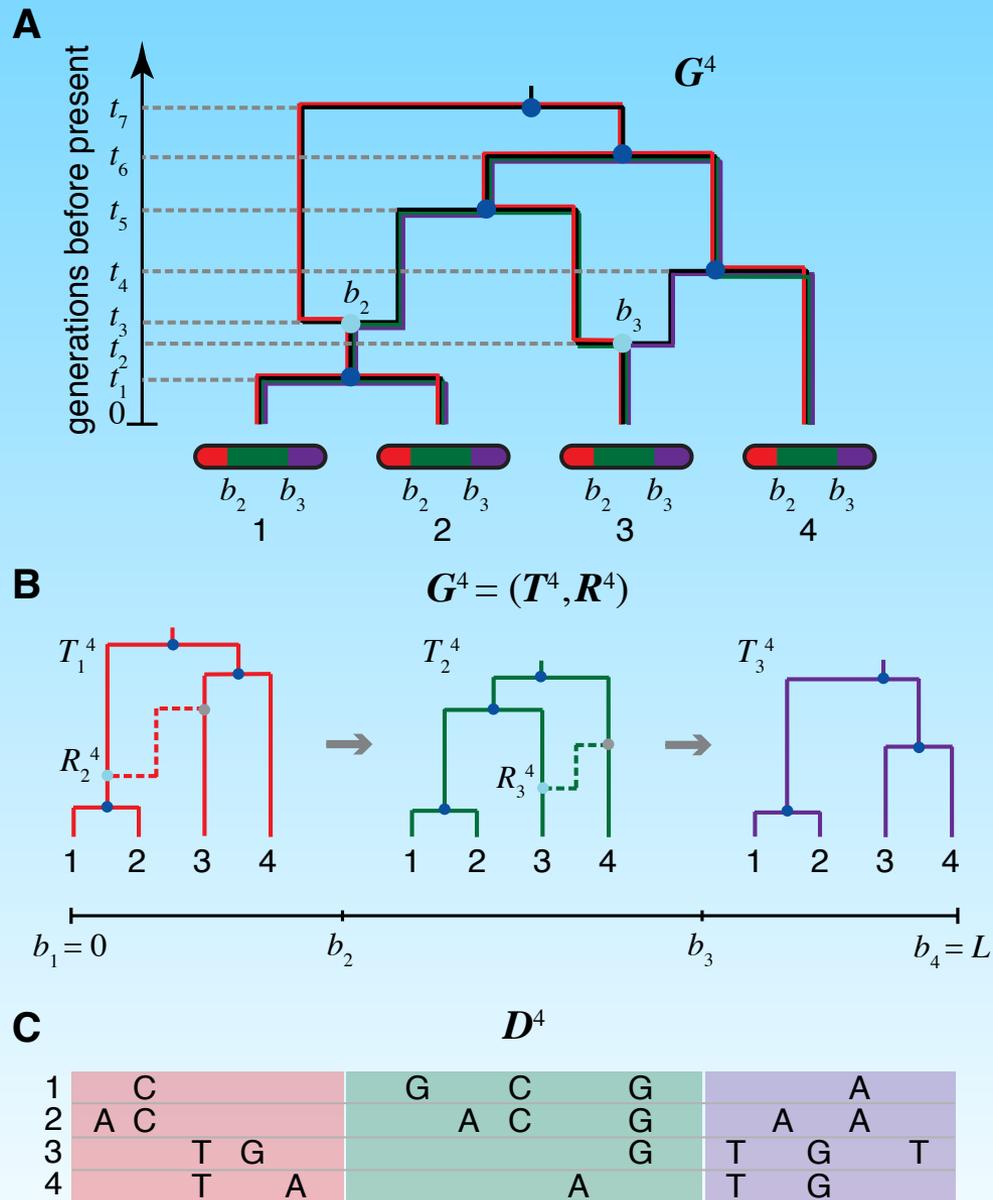
Sequentially Markov Coalescent (SMC)



$$P(T_i \mid T_1, \dots, T_{i-1}) = P(T_i \mid T_{i-1})$$

$$T_{i-1} \perp T_{i+1} \mid T_i$$





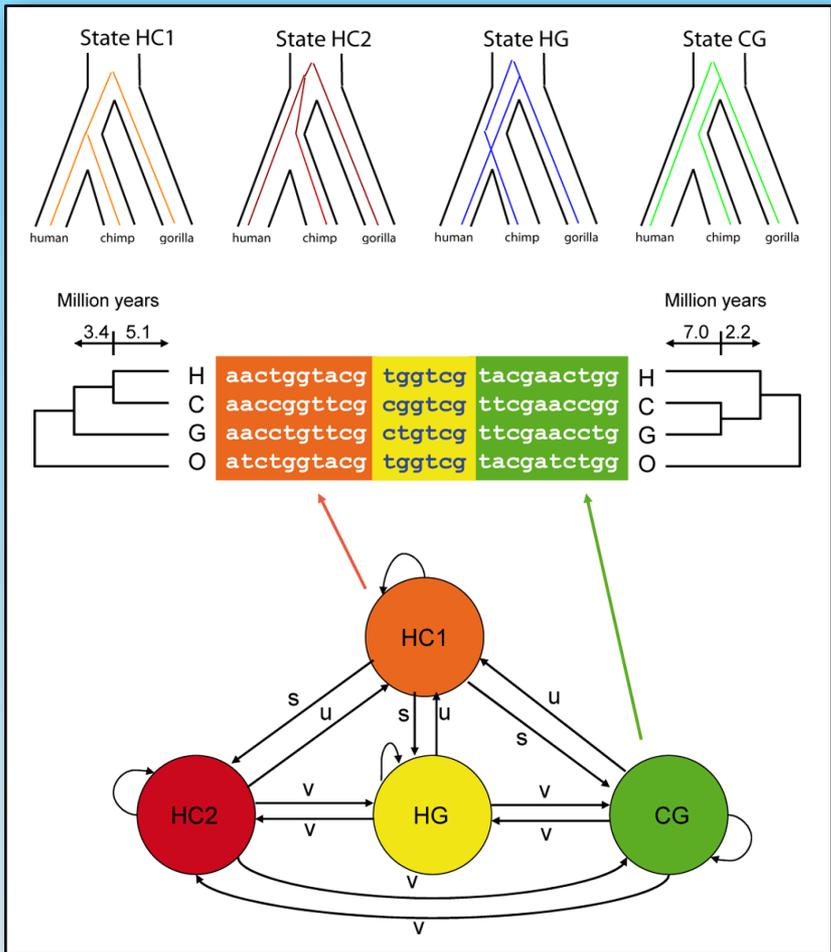


Discretized SMC and Hidden Markov Models

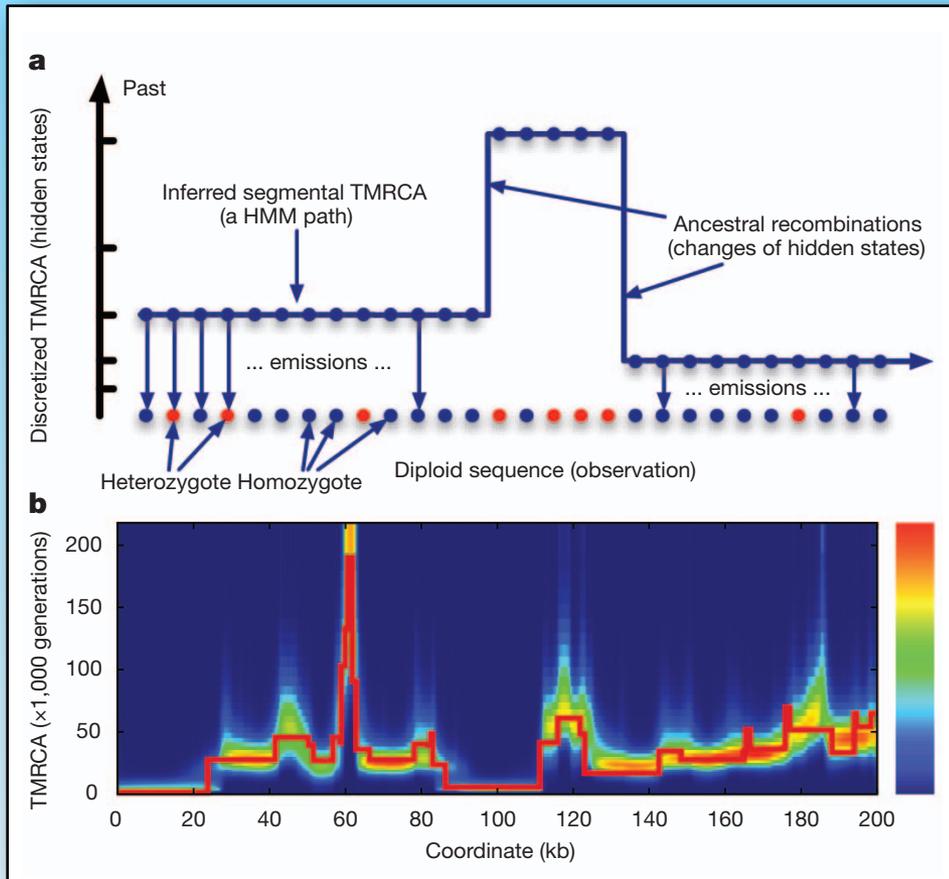
- By *discretizing* time and *enumerating topologies*, the continuous state space of the SMC can be approximated by a finite set
- This opens up the possibility of using *hidden Markov models* (HMM) for inference
- Standard dynamic-programming algorithms for HMMs allow for *exact ARG inference*, up to the SMC and discretization



Hidden Markov Models



Coal-HMM



PSMC



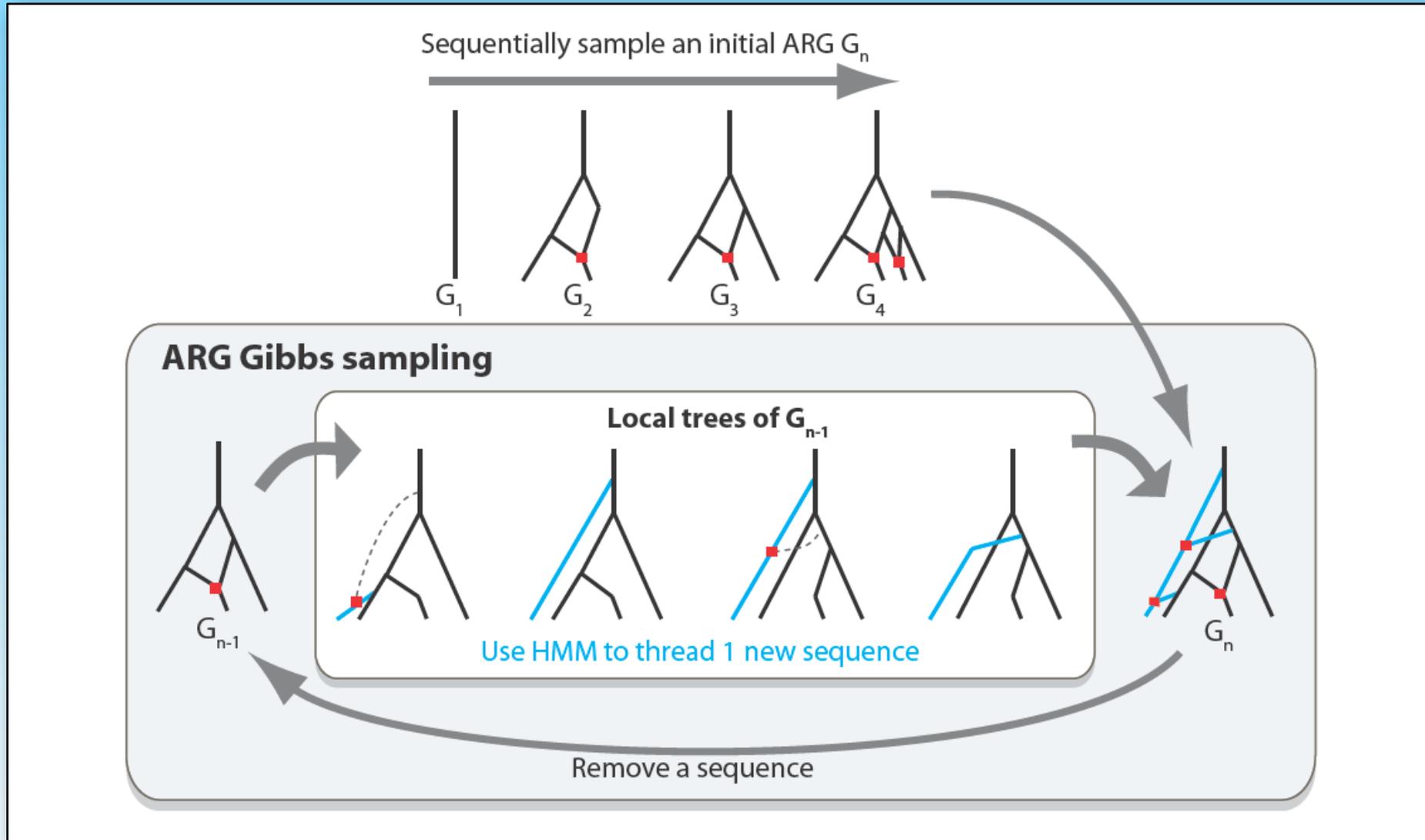


New Approach: Chromosome “Threading”

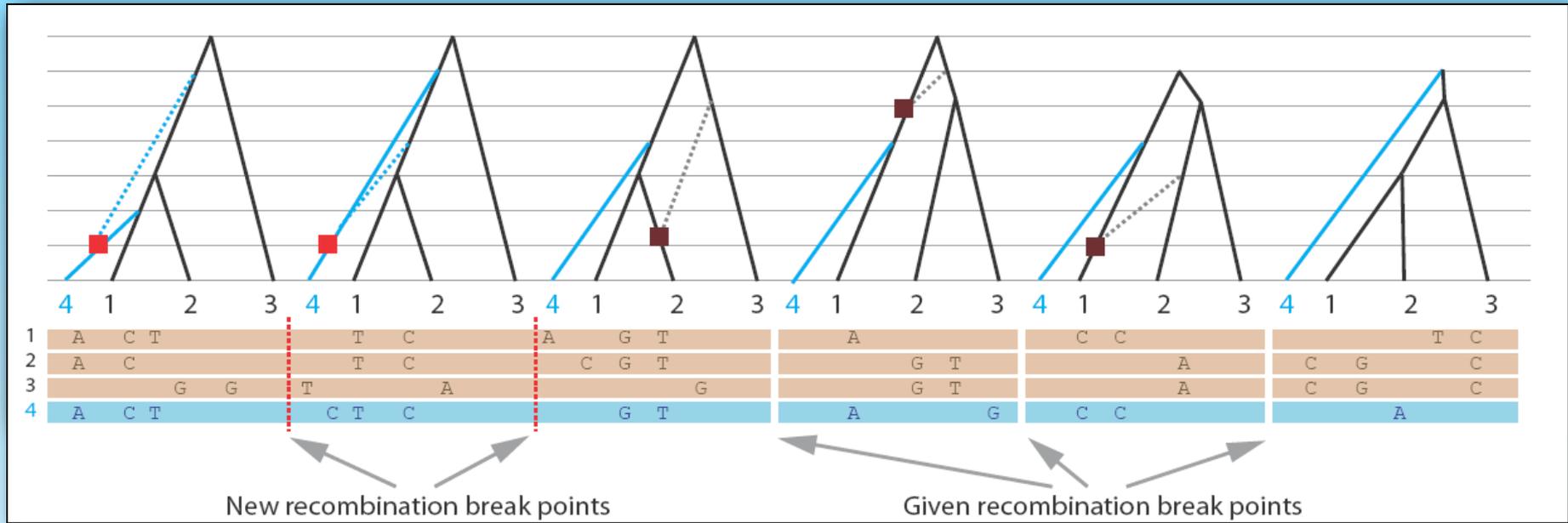
- Start with a data set of n sequences, D , and an ARG for $n-1$ of them, G_{n-1}
- Extend G_{n-1} to represent evolutionary history of n th sequence, obtaining G_n
- Sample this extension in a manner consistent with the conditional distribution, $P(G_n | G_{n-1}, D, \Theta)$, under the DSMC
- In repeated applications this operation is the basis of an **ARG sampling** algorithm



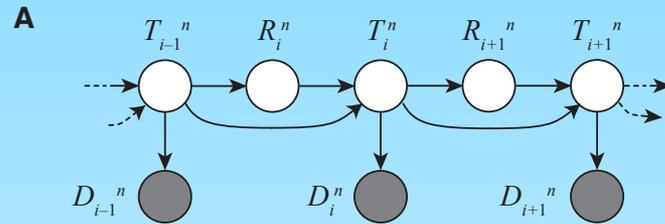
ARGweaver Sampling



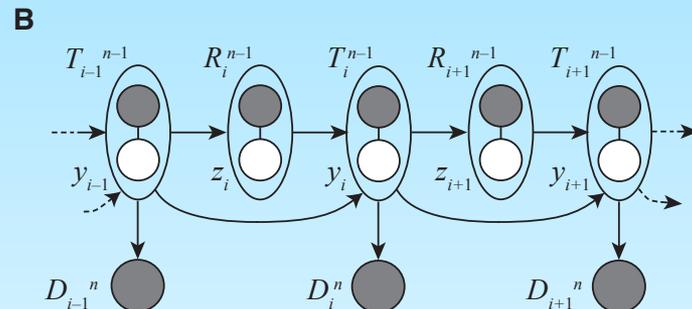
Threading



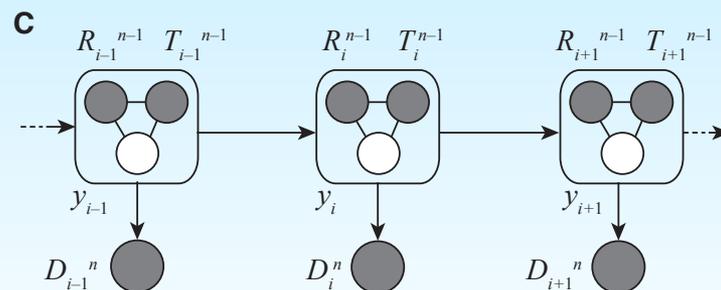
Graphical Models



SMC (discretized)



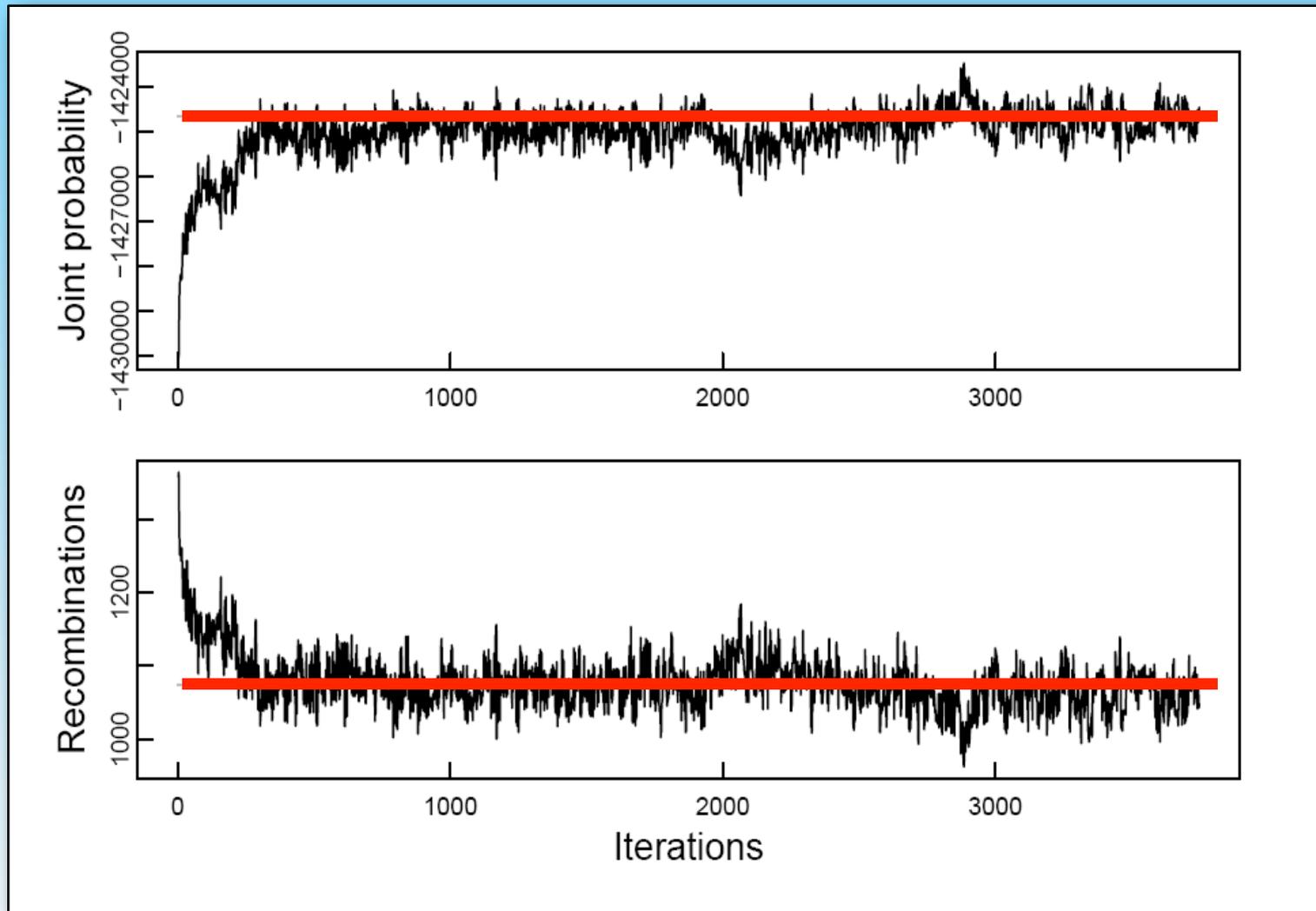
Threading



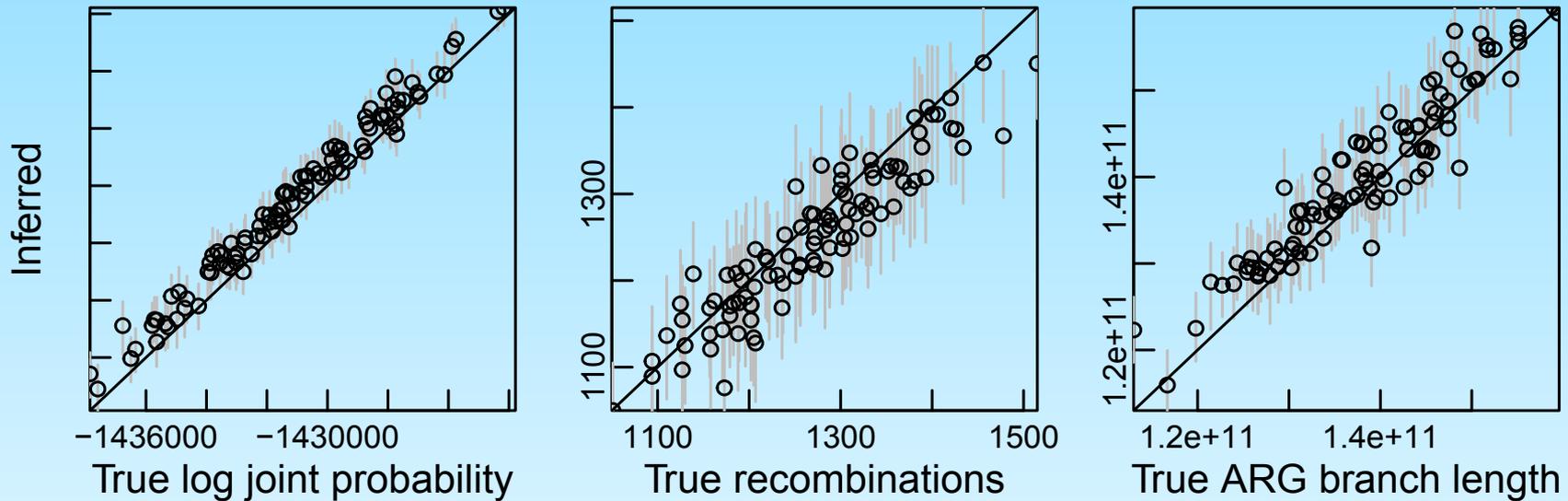
Reduced Threading



ARGweaver Converges Quickly



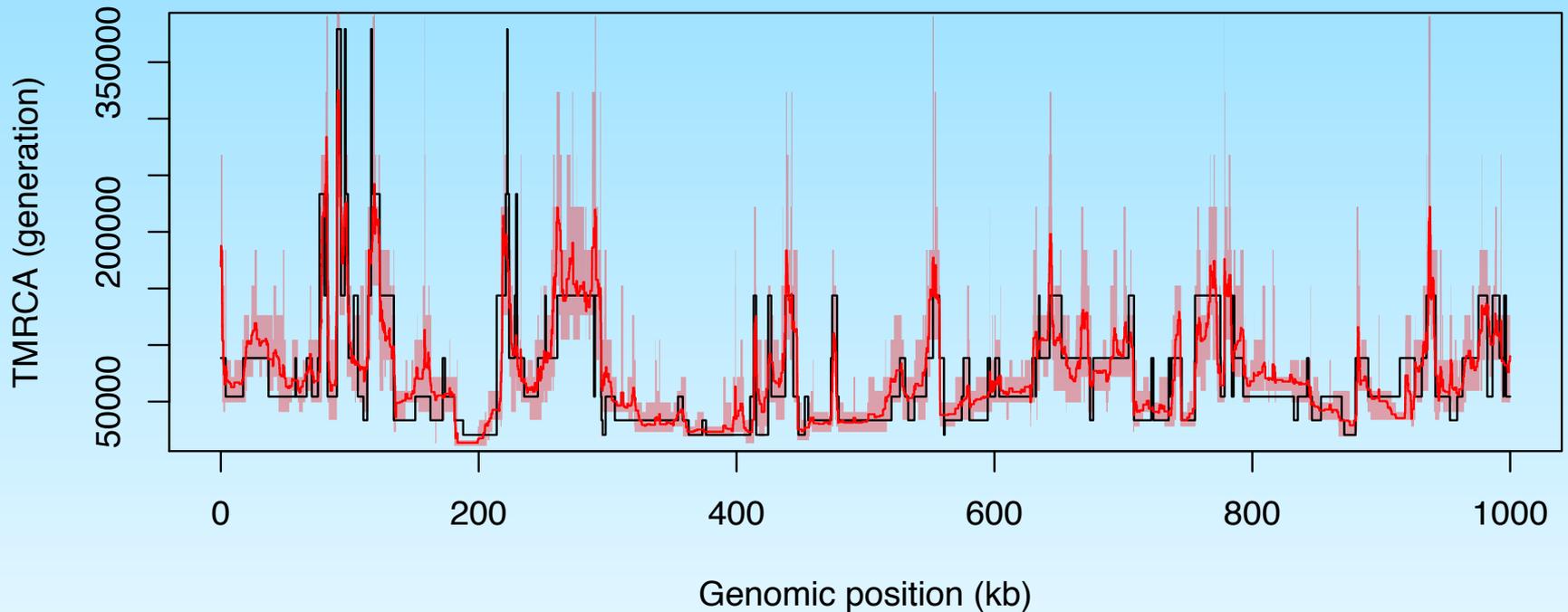
Recovery of Features of Simulated ARGs



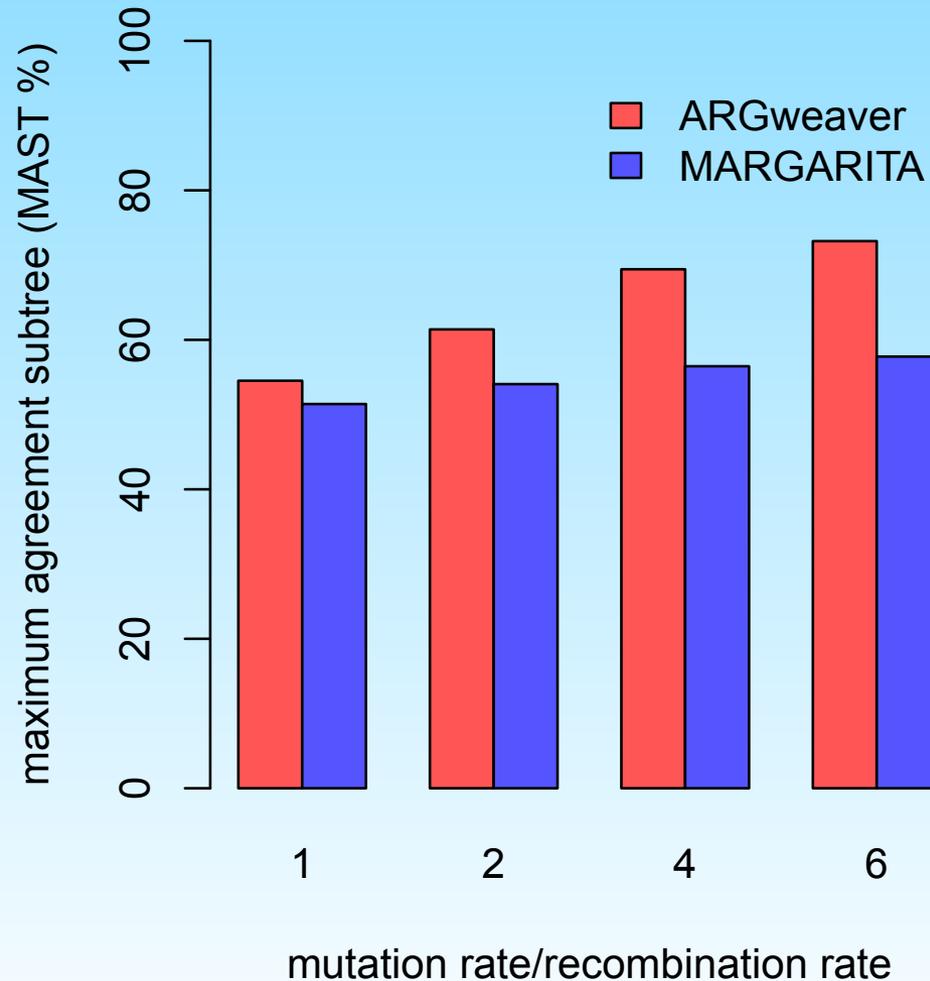
Part 2: *ARGweaver* Results



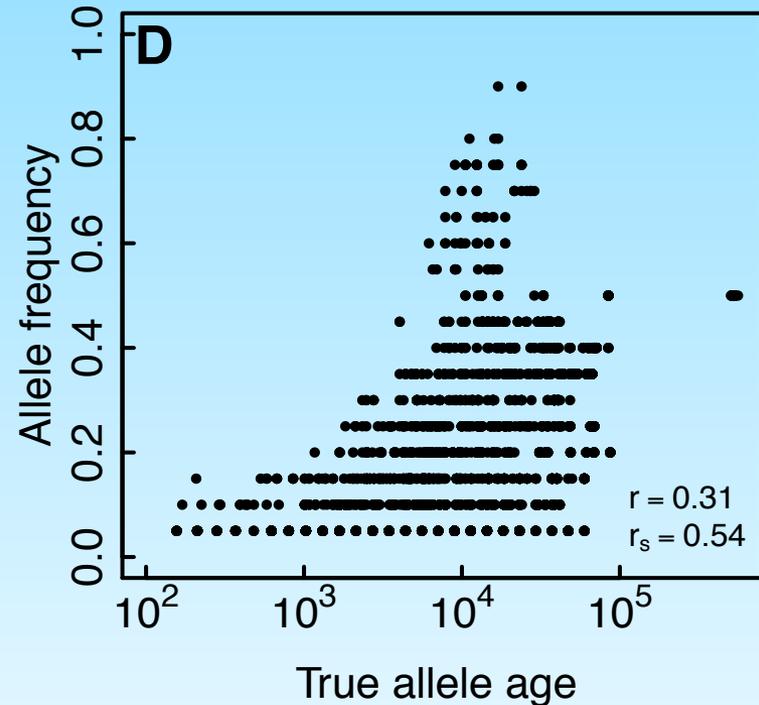
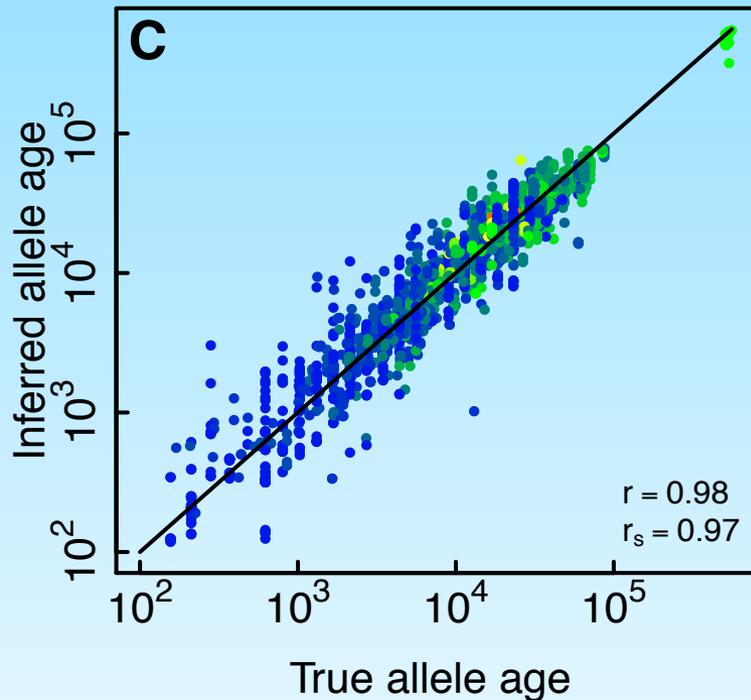
Recovery of Times to Most Recent Common Ancestry



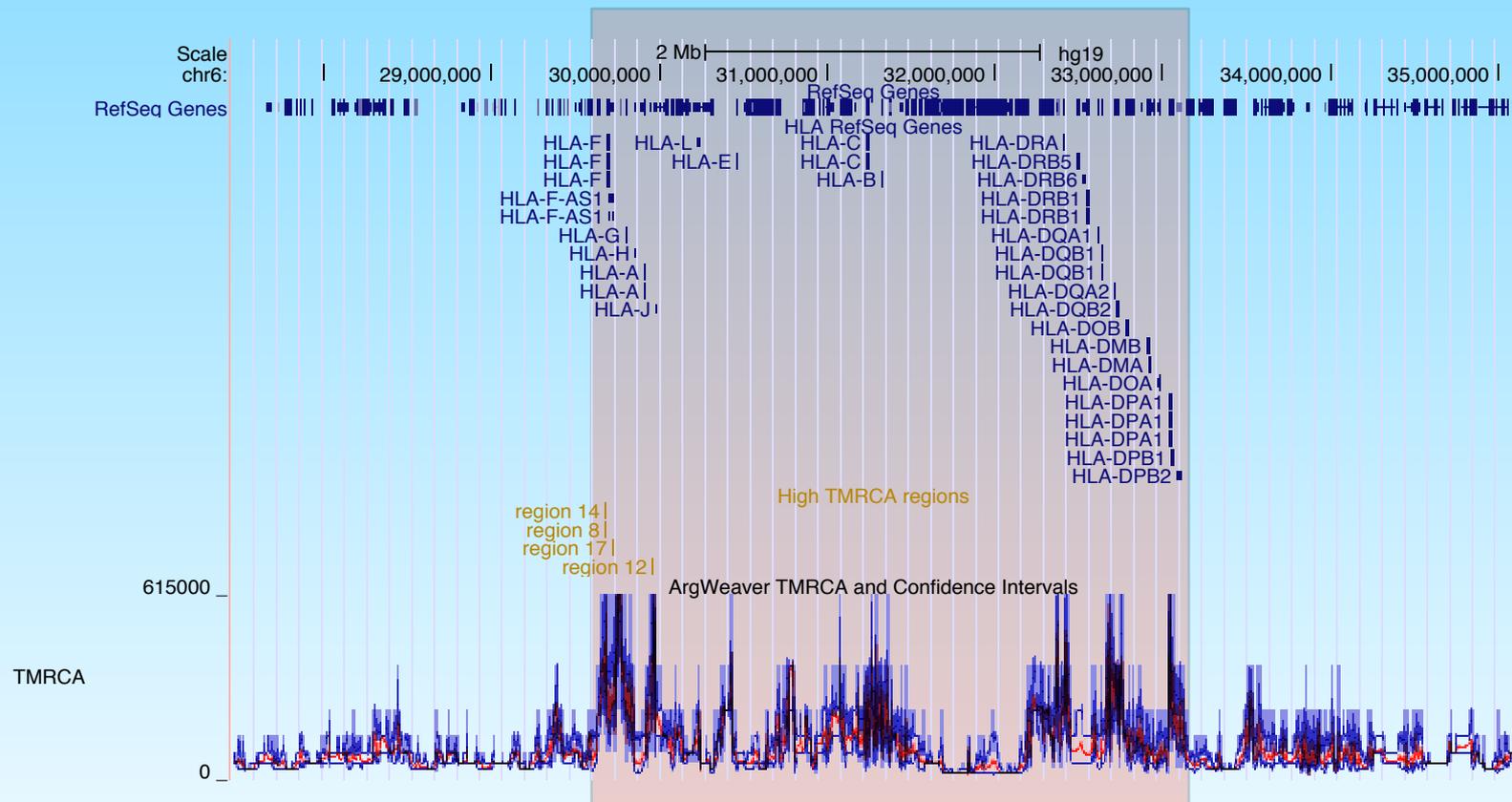
Accuracy of Inferred Trees



Recovery of Allele Age



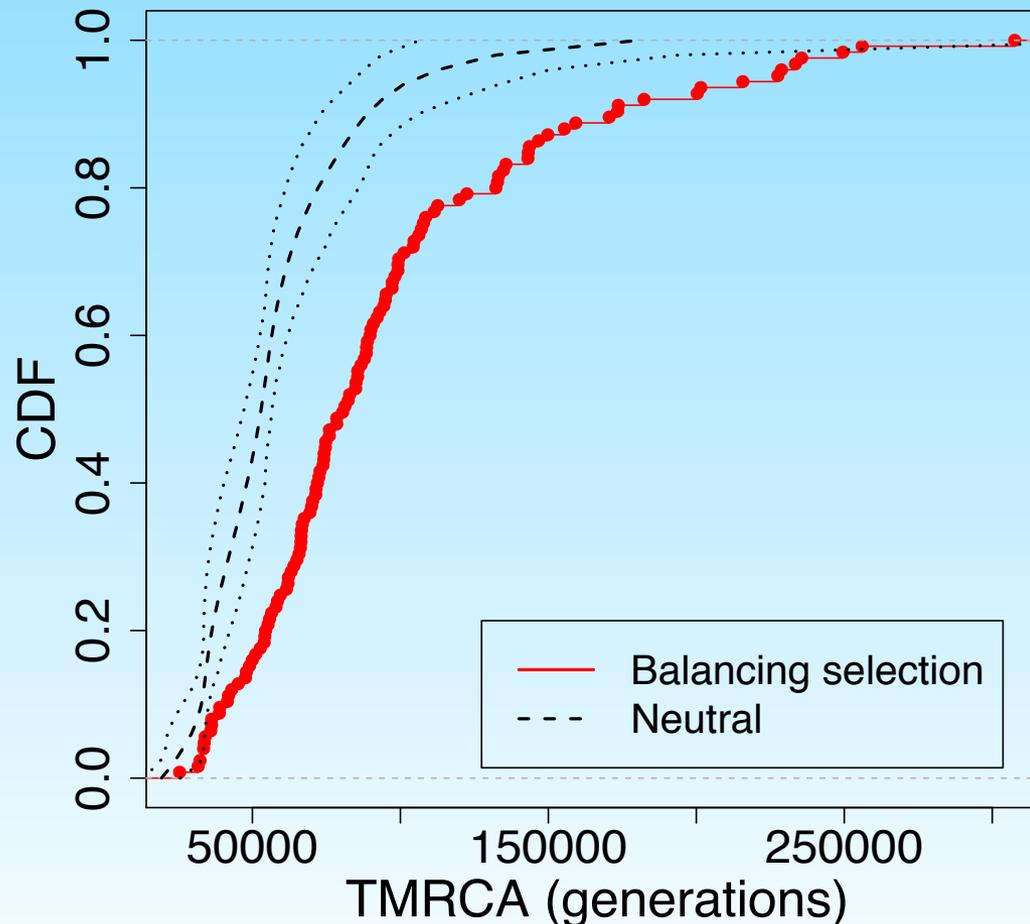
Real Data: Regions of High TMRCA



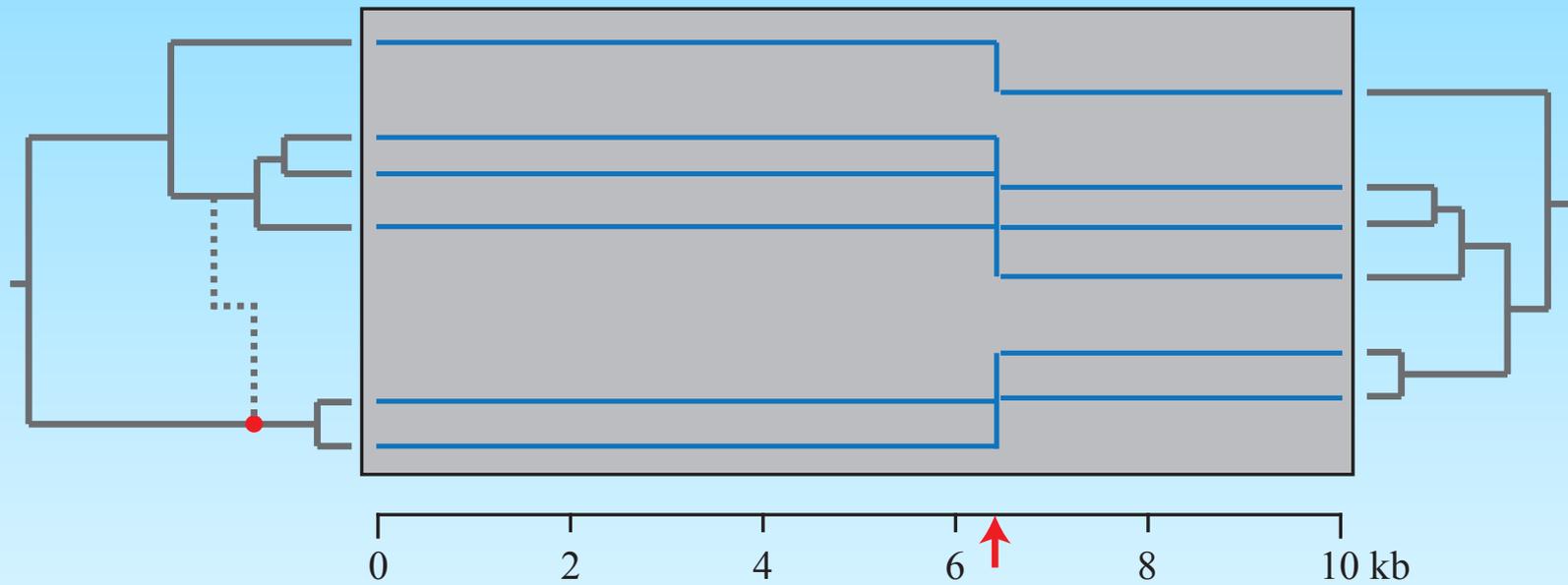
HLA Region



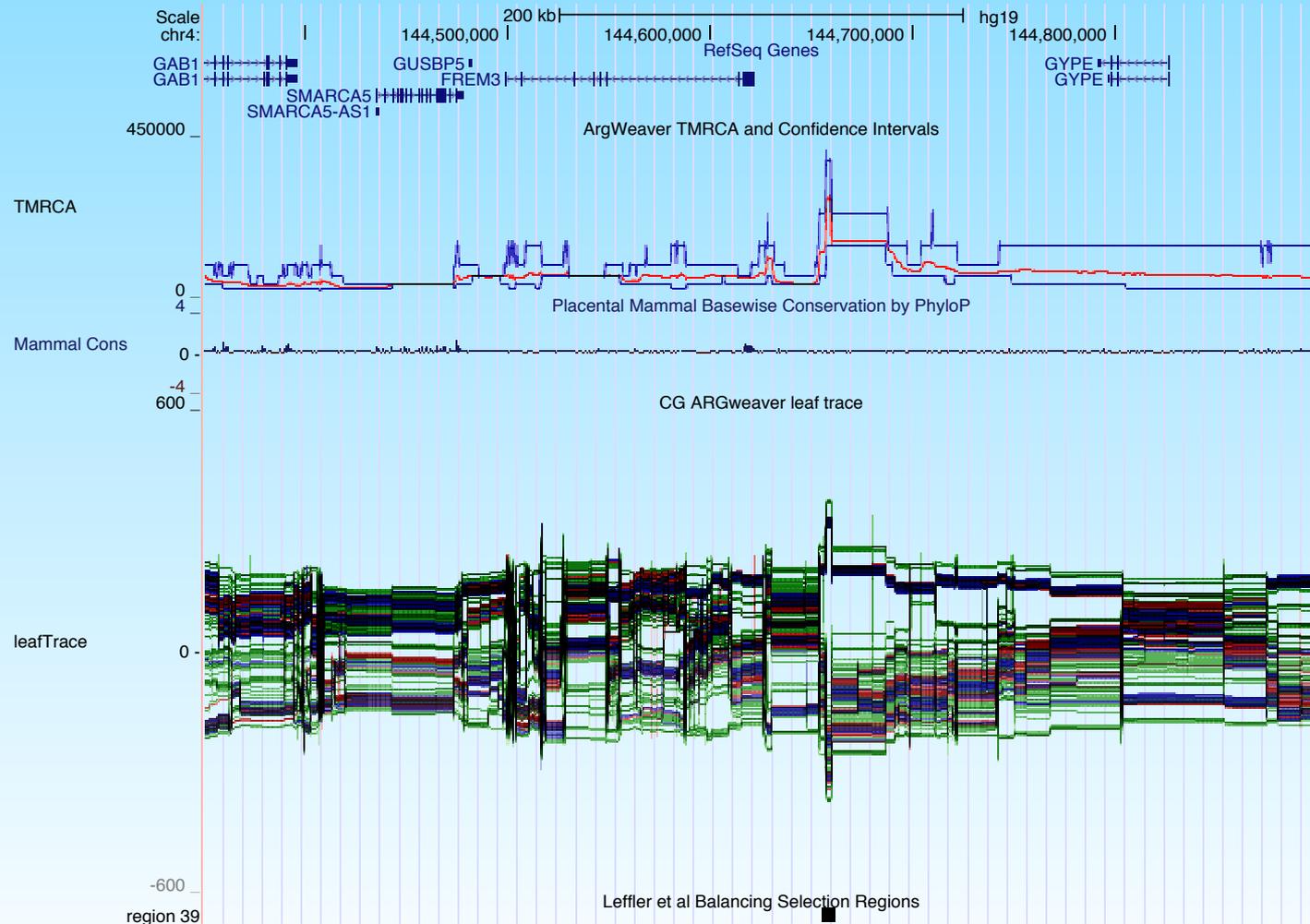
Regions of Shared Human/Chimp Polymorphism Have Old TMRCA



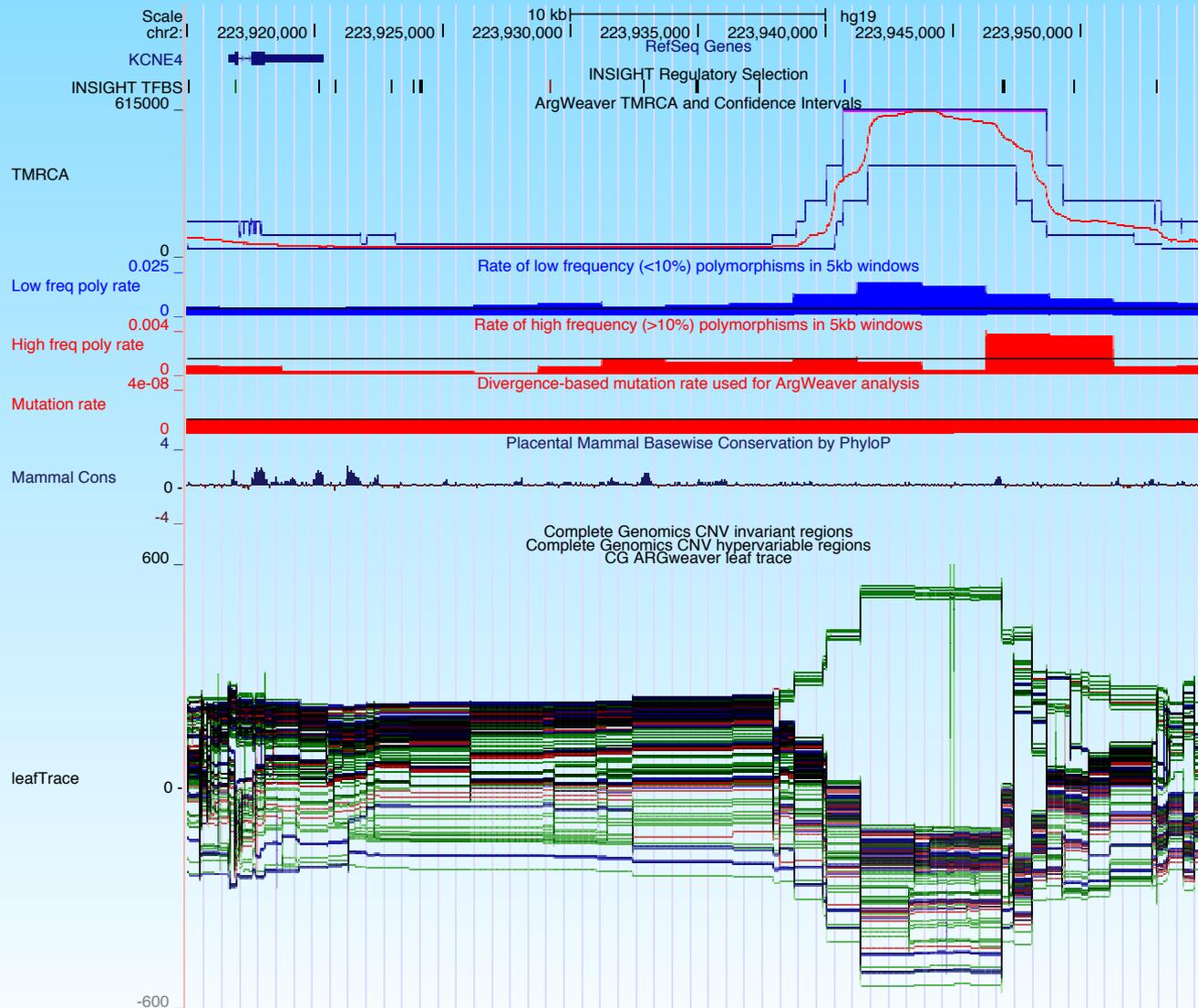
Leaf Trace



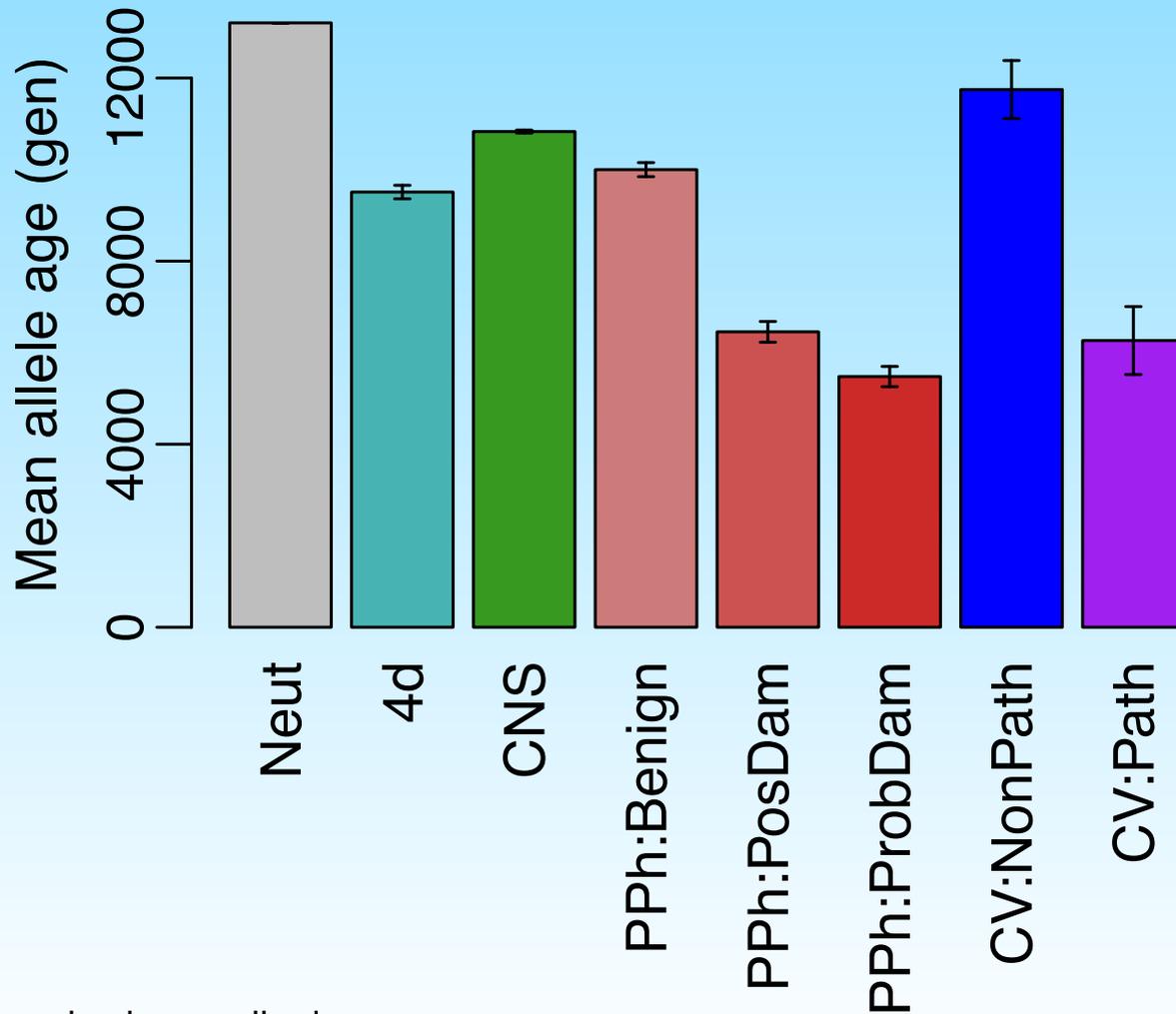
Putative Balancing Selection at FREM3



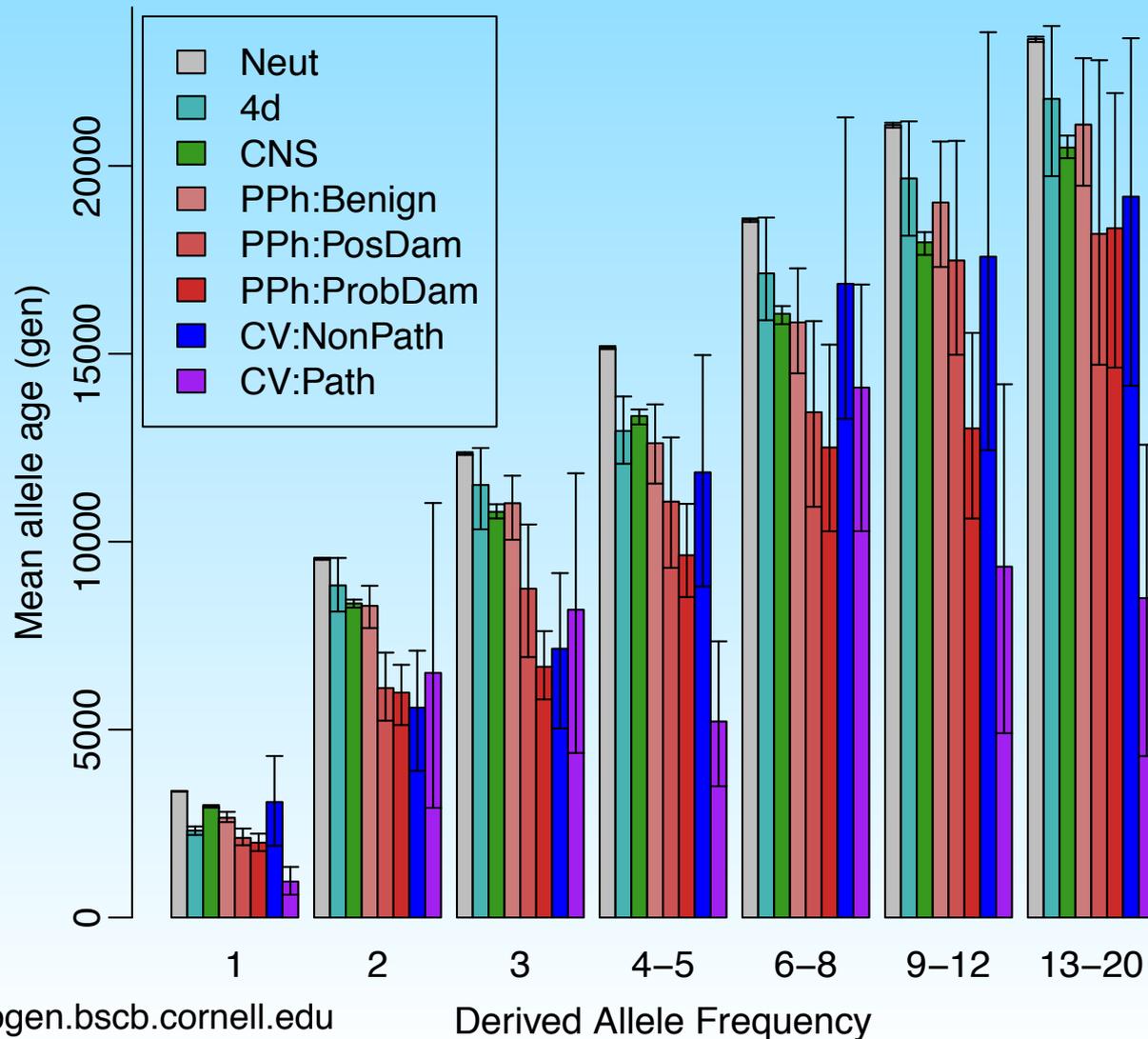
KCNE4 Promoter



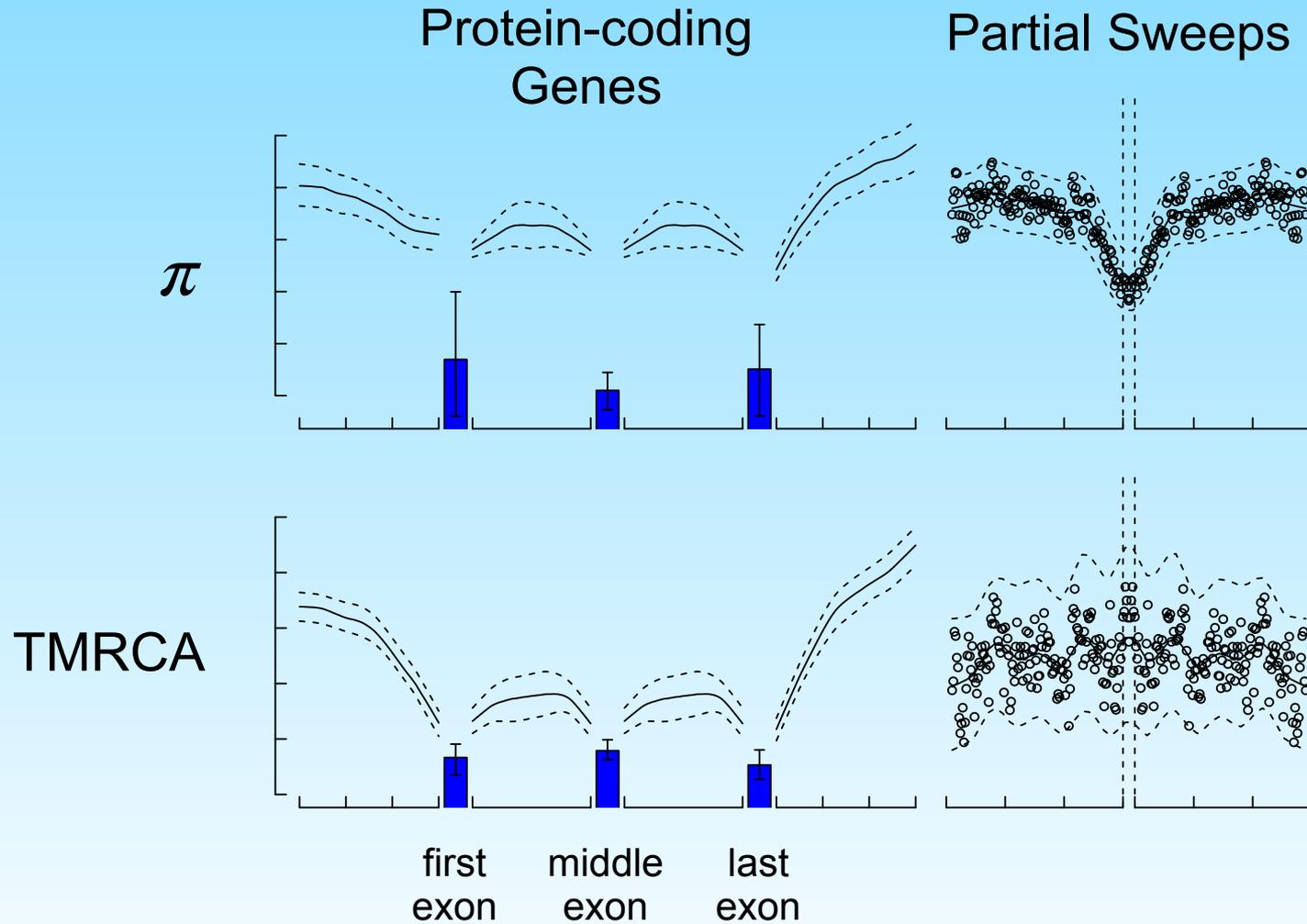
Sites Under Selection Have Decreased Allele Age



...Even After Accounting for Derived Allele Frequency



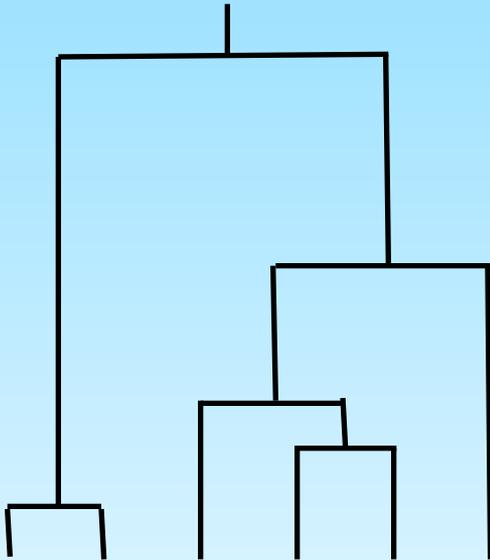
Genes and Sweeps (CEU)





Relative TMRCA Halflife (RTH)

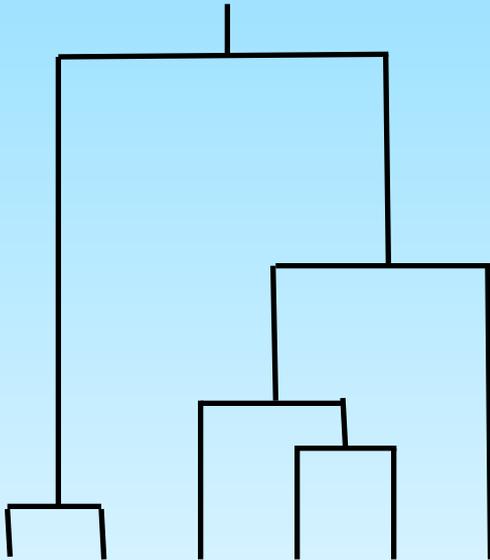
Neutral Drift



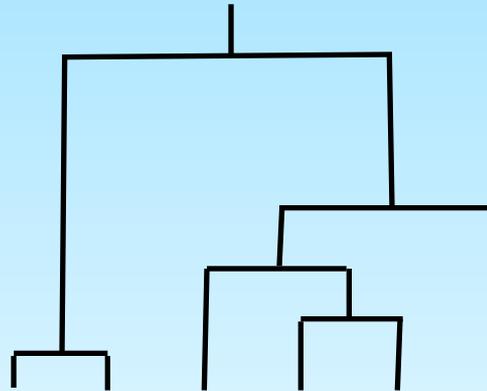


Relative TMRCA Halflife (RTH)

Neutral Drift



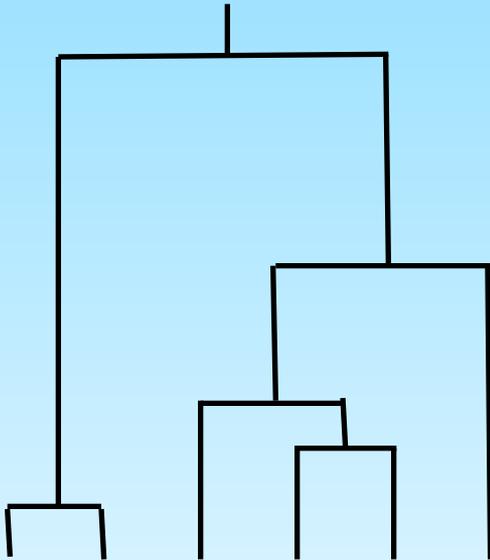
Background Selection



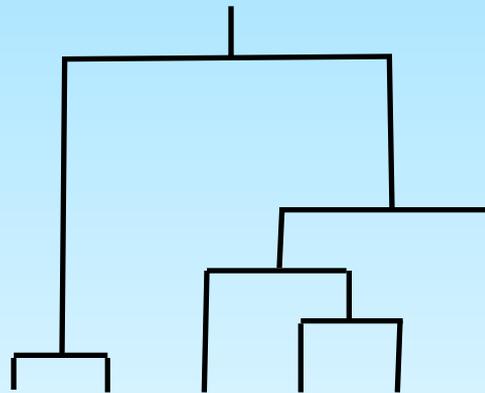


Relative TMRCA Halflife (RTH)

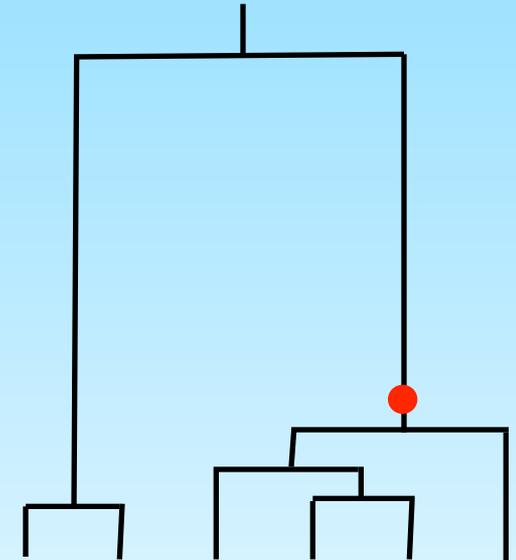
Neutral Drift



Background Selection



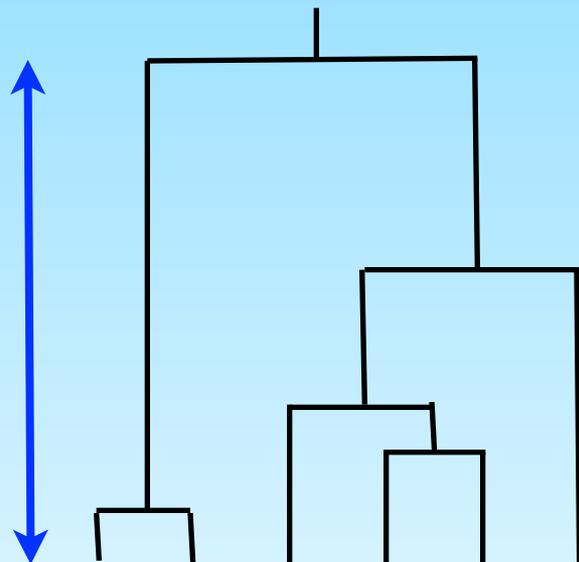
Partial Sweep



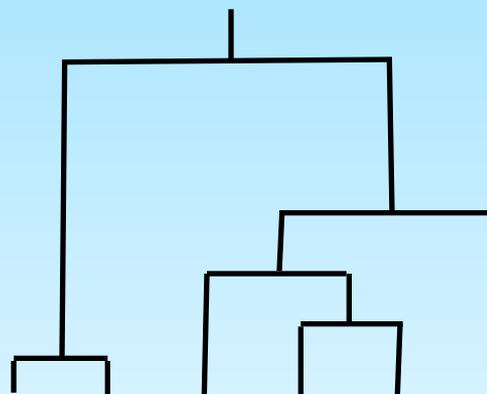


Relative TMRCA Halflife (RTH)

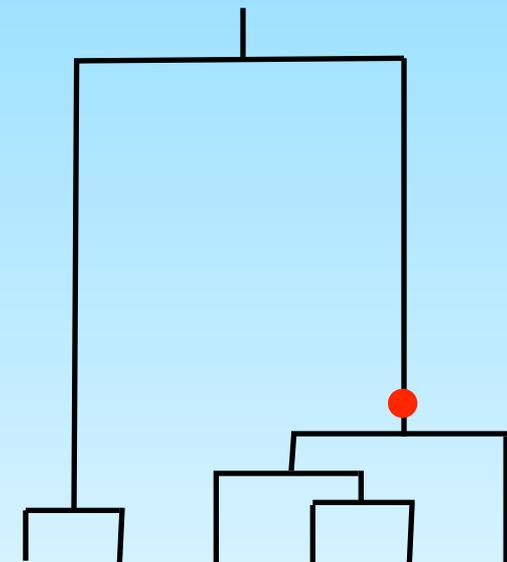
Neutral Drift



Background Selection



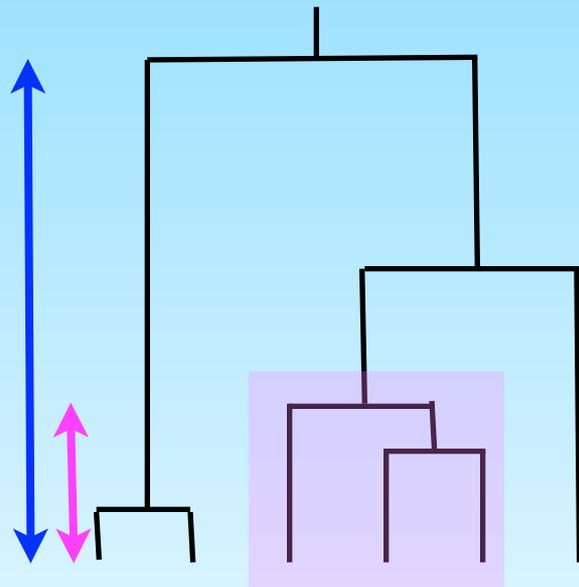
Partial Sweep



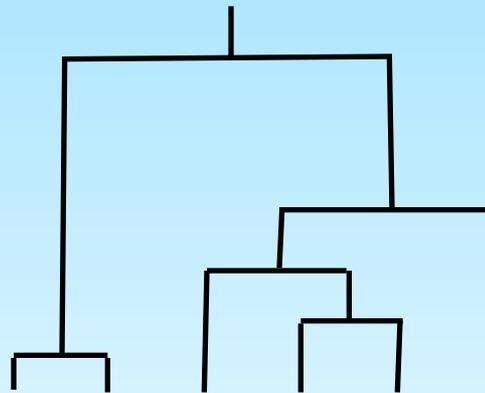


Relative TMRCA Halflife (RTH)

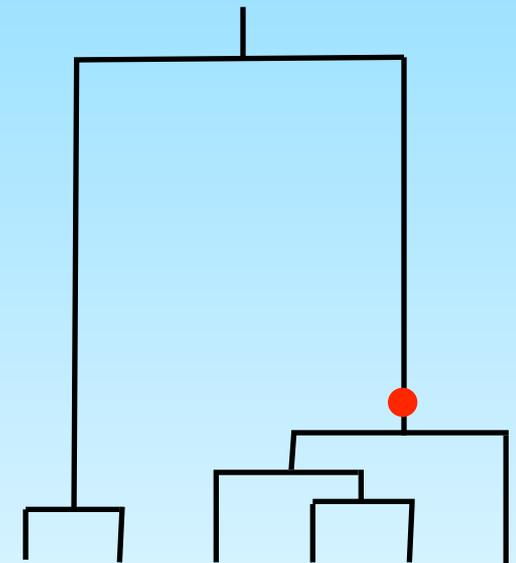
Neutral Drift



Background Selection



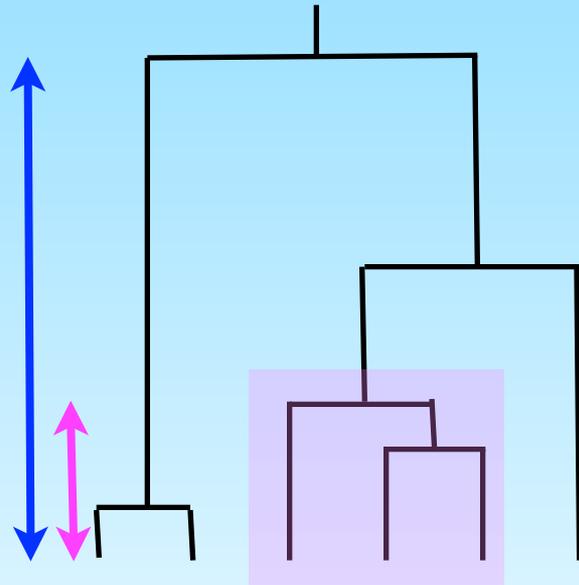
Partial Sweep



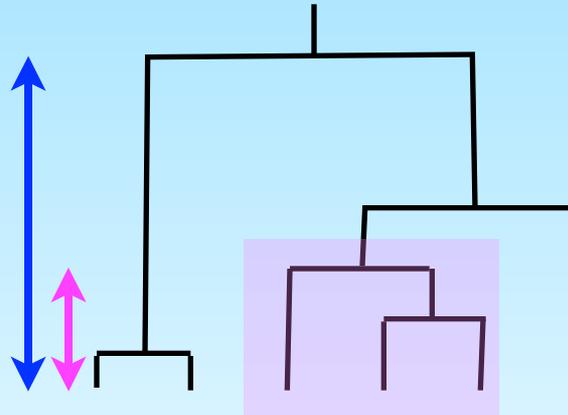


Relative TMRCA Halflife (RTH)

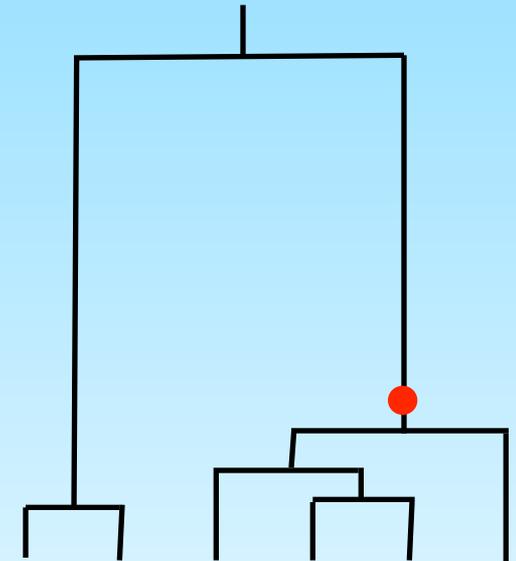
Neutral Drift



Background Selection



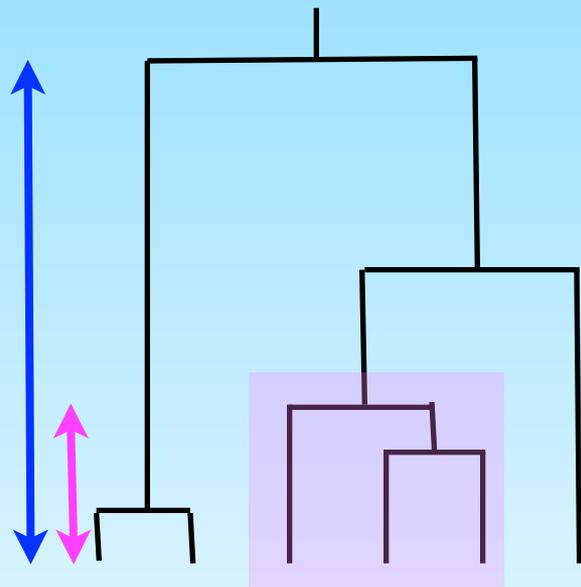
Partial Sweep



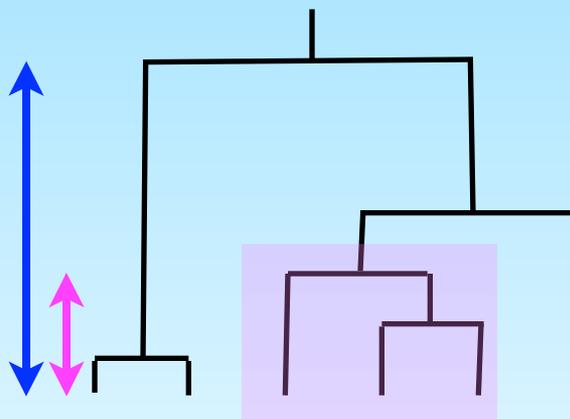


Relative TMRCA Halflife (RTH)

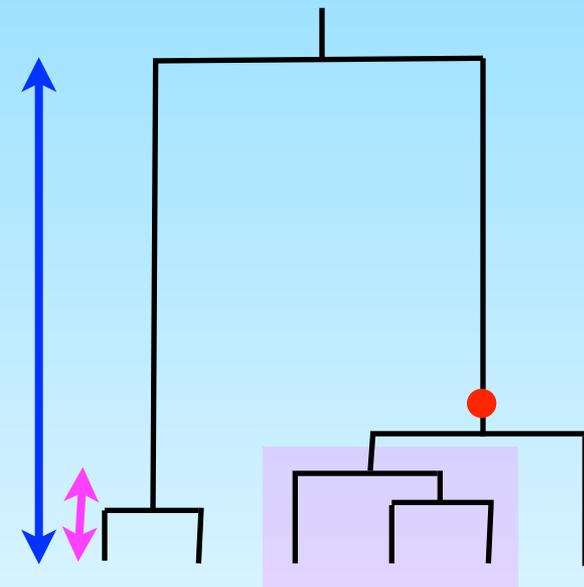
Neutral Drift



Background Selection



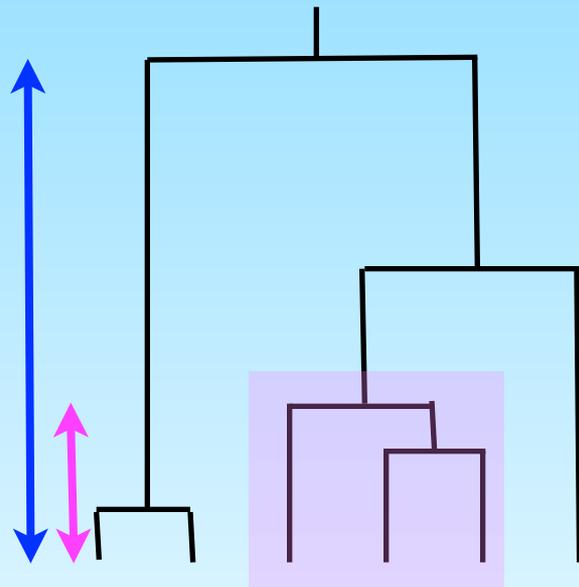
Partial Sweep



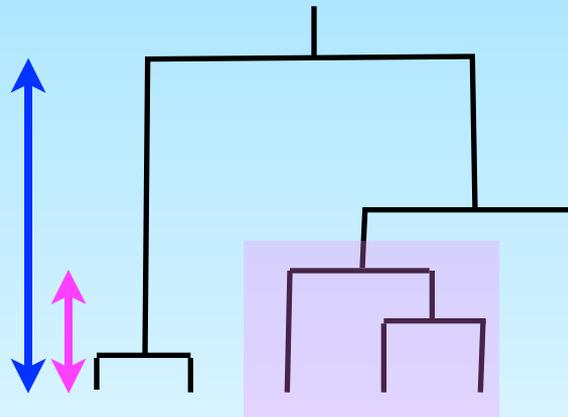


Relative TMRCA Halflife (RTH)

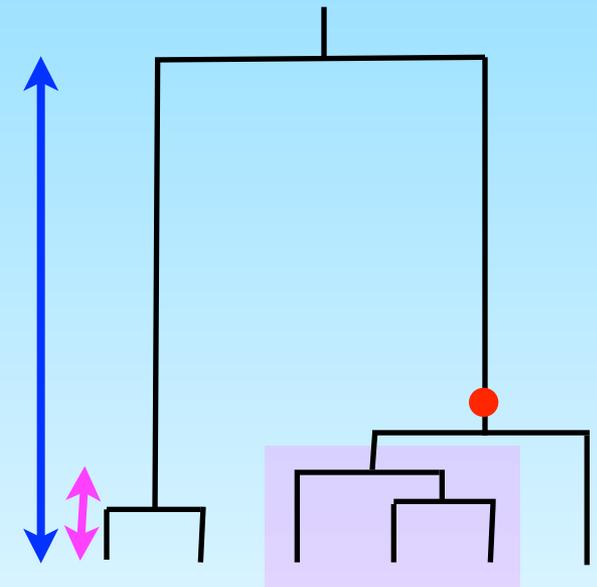
Neutral Drift



Background Selection



Partial Sweep



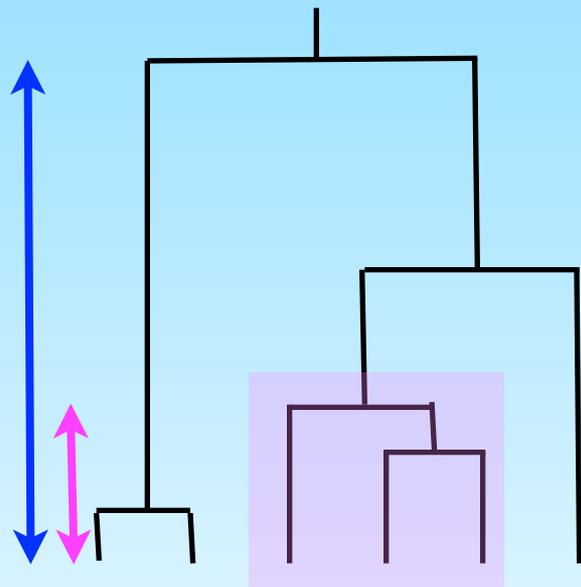
$$\text{RTH} = \text{half-TMRCA} / \text{TMRCA}$$





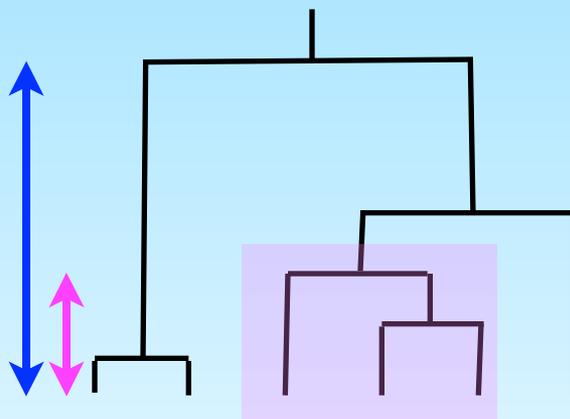
Relative TMRCA Halflife (RTH)

Neutral Drift



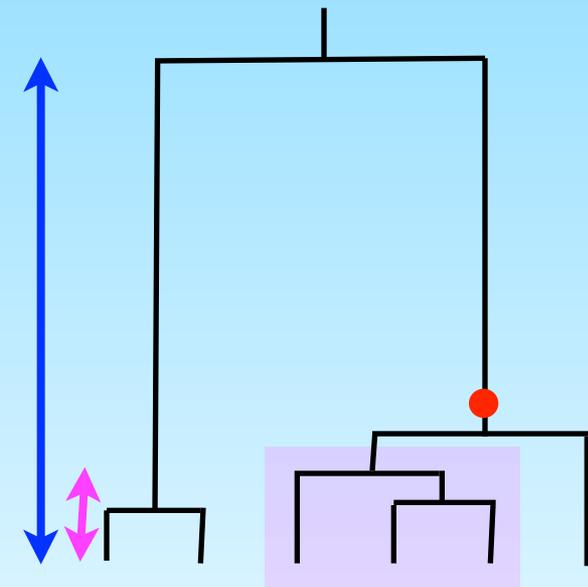
$$\text{RTH} = 1/3$$

Background Selection



$$\text{RTH} = 1/3$$

Partial Sweep

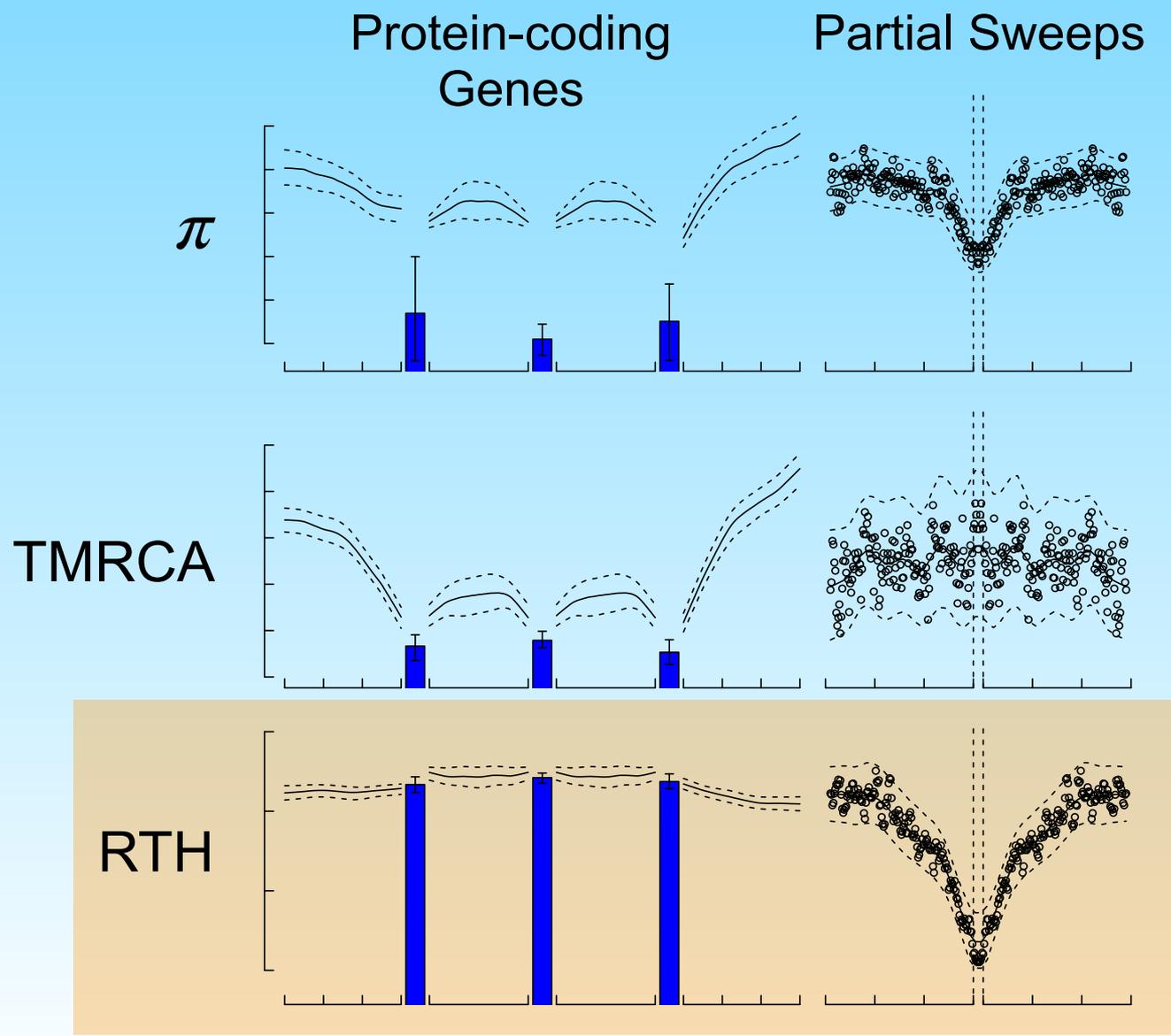


$$\text{RTH} = 1/6$$

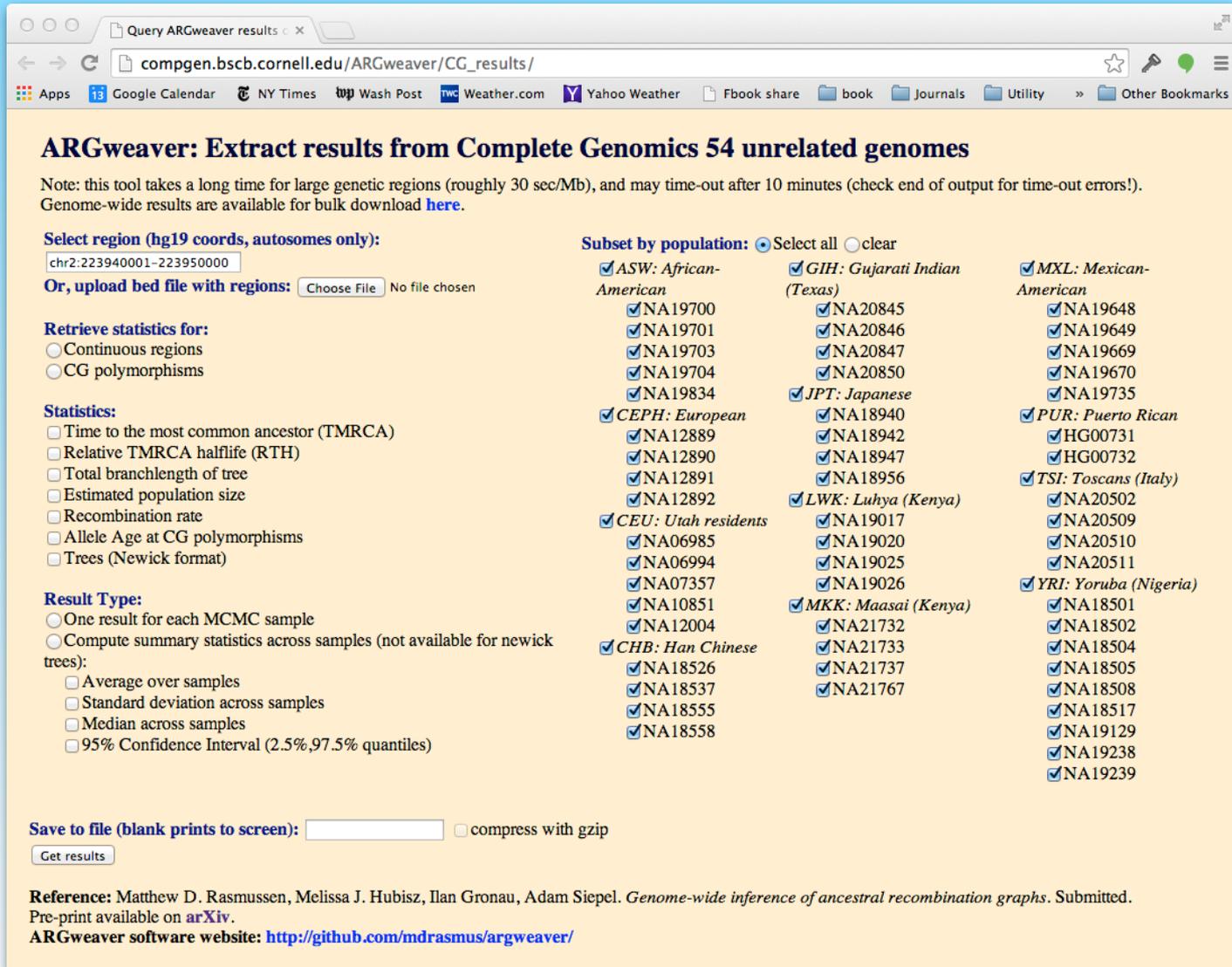
$$\text{RTH} = \text{half-TMRCA} / \text{TMRCA}$$



Genes and Sweeps (CEU)



Query Interface



Query ARGweaver results

compngen.bscb.cornell.edu/ARGweaver/CG_results/

Apps 13 Google Calendar NY Times Wash Post Weather.com Yahoo Weather Fbook share book Journals Utility Other Bookmarks

ARGweaver: Extract results from Complete Genomics 54 unrelated genomes

Note: this tool takes a long time for large genetic regions (roughly 30 sec/Mb), and may time-out after 10 minutes (check end of output for time-out errors!). Genome-wide results are available for bulk download [here](#).

Select region (hg19 coords, autosomes only):

Or, upload bed file with regions: No file chosen

Retrieve statistics for:
 Continuous regions
 CG polymorphisms

Statistics:
 Time to the most common ancestor (TMRCA)
 Relative TMRCA halflife (RTH)
 Total branchlength of tree
 Estimated population size
 Recombination rate
 Allele Age at CG polymorphisms
 Trees (Newick format)

Result Type:
 One result for each MCMC sample
 Compute summary statistics across samples (not available for newick trees):
 Average over samples
 Standard deviation across samples
 Median across samples
 95% Confidence Interval (2.5%,97.5% quantiles)

Subset by population: Select all clear

<input checked="" type="checkbox"/> ASW: African-American	<input checked="" type="checkbox"/> GIH: Gujarati Indian (Texas)	<input checked="" type="checkbox"/> MXL: Mexican-American
<input checked="" type="checkbox"/> NA19700	<input checked="" type="checkbox"/> NA20845	<input checked="" type="checkbox"/> NA19648
<input checked="" type="checkbox"/> NA19701	<input checked="" type="checkbox"/> NA20846	<input checked="" type="checkbox"/> NA19649
<input checked="" type="checkbox"/> NA19703	<input checked="" type="checkbox"/> NA20847	<input checked="" type="checkbox"/> NA19669
<input checked="" type="checkbox"/> NA19704	<input checked="" type="checkbox"/> NA20850	<input checked="" type="checkbox"/> NA19670
<input checked="" type="checkbox"/> NA19834	<input checked="" type="checkbox"/> JPT: Japanese	<input checked="" type="checkbox"/> NA19735
<input checked="" type="checkbox"/> CEPH: European	<input checked="" type="checkbox"/> NA18940	<input checked="" type="checkbox"/> PUR: Puerto Rican
<input checked="" type="checkbox"/> NA12889	<input checked="" type="checkbox"/> NA18942	<input checked="" type="checkbox"/> HG00731
<input checked="" type="checkbox"/> NA12890	<input checked="" type="checkbox"/> NA18947	<input checked="" type="checkbox"/> HG00732
<input checked="" type="checkbox"/> NA12891	<input checked="" type="checkbox"/> NA18956	<input checked="" type="checkbox"/> TSI: Toscons (Italy)
<input checked="" type="checkbox"/> NA12892	<input checked="" type="checkbox"/> LWK: Luhya (Kenya)	<input checked="" type="checkbox"/> NA20502
<input checked="" type="checkbox"/> CEU: Utah residents	<input checked="" type="checkbox"/> NA19017	<input checked="" type="checkbox"/> NA20509
<input checked="" type="checkbox"/> NA06985	<input checked="" type="checkbox"/> NA19020	<input checked="" type="checkbox"/> NA20510
<input checked="" type="checkbox"/> NA06994	<input checked="" type="checkbox"/> NA19025	<input checked="" type="checkbox"/> NA20511
<input checked="" type="checkbox"/> NA07357	<input checked="" type="checkbox"/> NA19026	<input checked="" type="checkbox"/> YRI: Yoruba (Nigeria)
<input checked="" type="checkbox"/> NA10851	<input checked="" type="checkbox"/> MKK: Maasai (Kenya)	<input checked="" type="checkbox"/> NA18501
<input checked="" type="checkbox"/> NA12004	<input checked="" type="checkbox"/> NA21732	<input checked="" type="checkbox"/> NA18502
<input checked="" type="checkbox"/> CHB: Han Chinese	<input checked="" type="checkbox"/> NA21733	<input checked="" type="checkbox"/> NA18504
<input checked="" type="checkbox"/> NA18526	<input checked="" type="checkbox"/> NA21737	<input checked="" type="checkbox"/> NA18505
<input checked="" type="checkbox"/> NA18537	<input checked="" type="checkbox"/> NA21767	<input checked="" type="checkbox"/> NA18508
<input checked="" type="checkbox"/> NA18555		<input checked="" type="checkbox"/> NA18517
<input checked="" type="checkbox"/> NA18558		<input checked="" type="checkbox"/> NA19129
		<input checked="" type="checkbox"/> NA19238
		<input checked="" type="checkbox"/> NA19239

Save to file (blank prints to screen): compress with gzip

Reference: Matthew D. Rasmussen, Melissa J. Hubisz, Ilan Gronau, Adam Siepel. *Genome-wide inference of ancestral recombination graphs*. Submitted. Pre-print available on [arXiv](#).

ARGweaver software website: <http://github.com/mrasmus/argweaver/>



Future Work

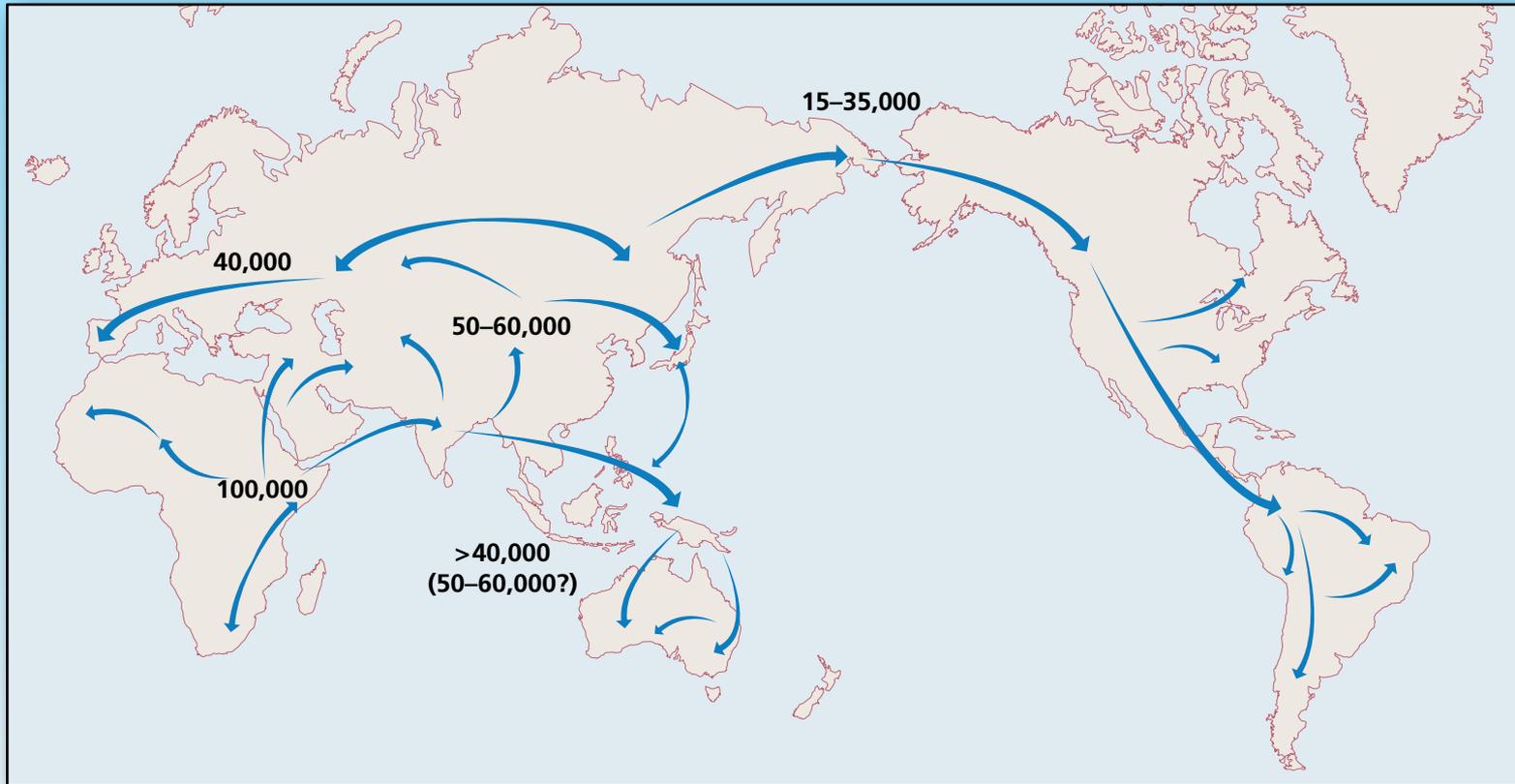
- Simultaneous phasing and ARG inference
- Demography inference
- Community resources
 - Extended dynamic querying of ARGs
 - On-the-fly threading
- Association mapping
- Any problem addressed by Li & Stephens model!



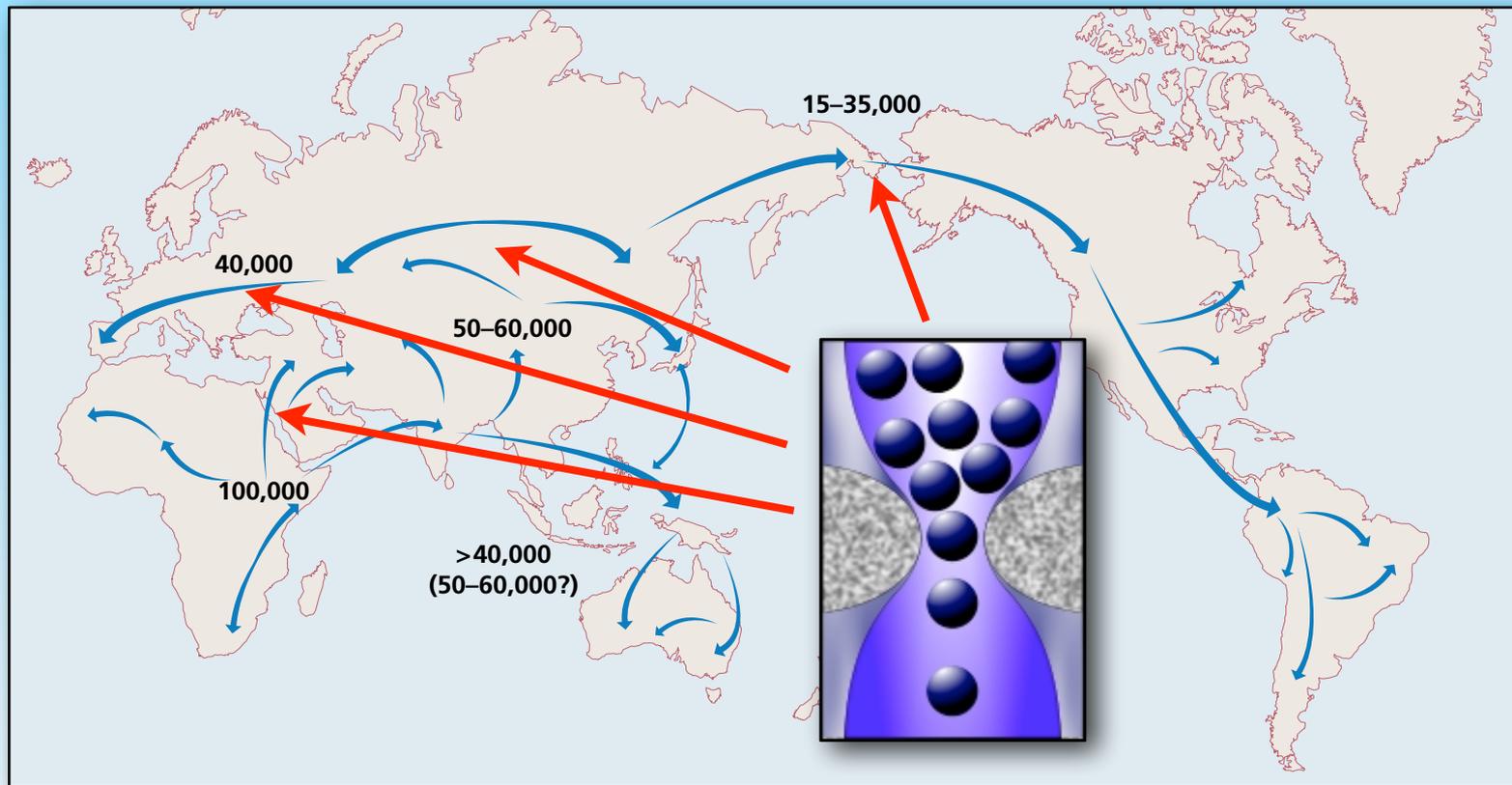
Part 3: Demography Inference



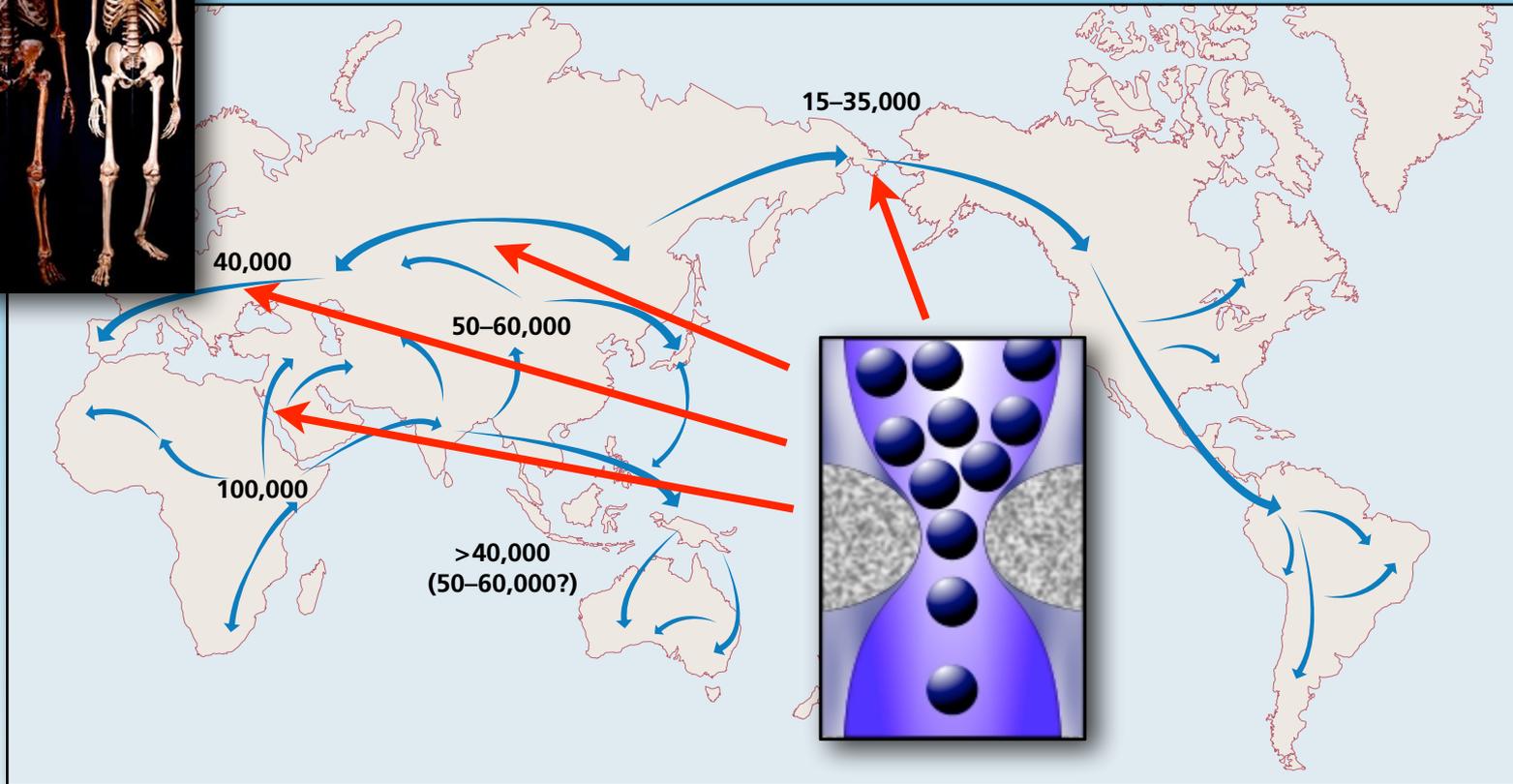
Origins of Human Populations



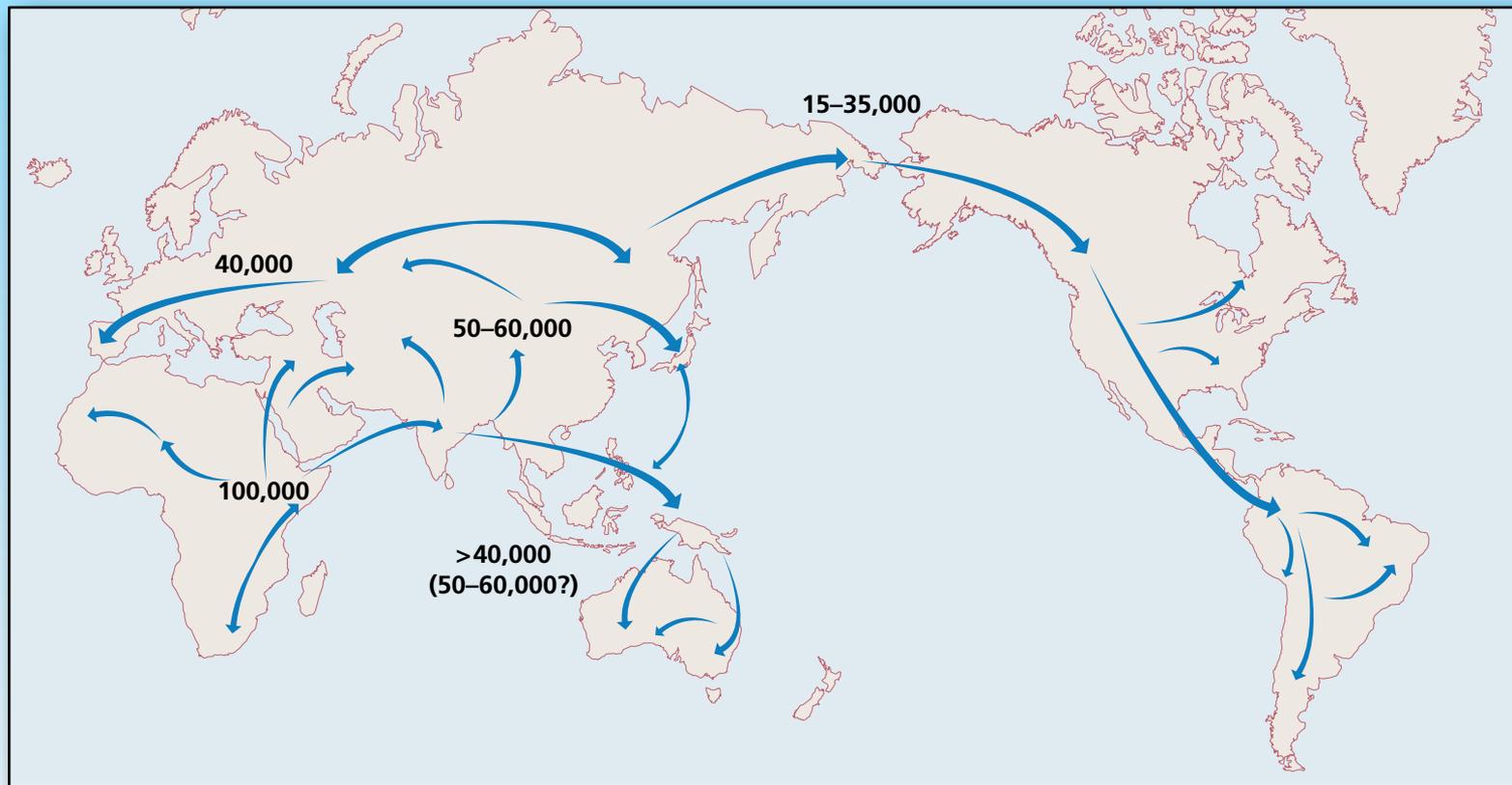
Origins of Human Populations



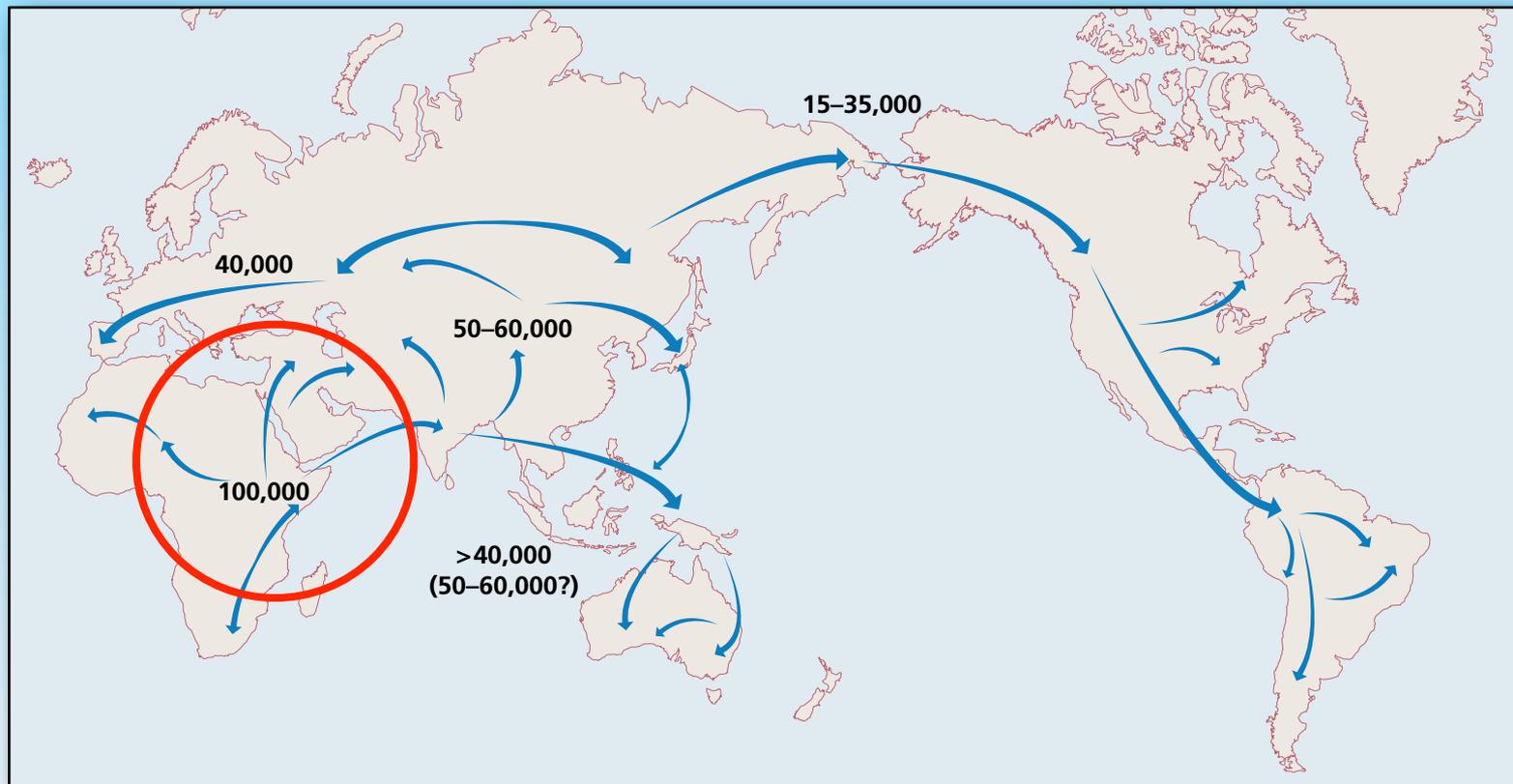
Origins of Human Populations



Origins of Human Populations



Origins of Human Populations



Origins of Human Populations



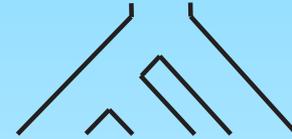
Integrated Statistical Model (*MCMCcoal*)

AATGAACCGTTTCTGAGGCCATT
 AGTGAACCGTTACTGACGCCATT
 AATGAATCGTTACTGAGGCTATT

X_i



G_i



θ, τ



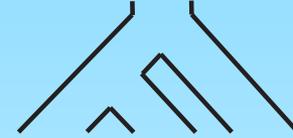
Integrated Statistical Model (*MCMCcoal*)

AATGAACCGTTTCTGAGGCCATT
 AGTGAACCGTTACTGACGCCATT
 AATGAA^TCGTTACTGAGGC^TATT

X_i



G_i



θ, τ

$$P(\mathbf{X}, \mathbf{G}, \boldsymbol{\theta}, \boldsymbol{\tau}) = P(\boldsymbol{\theta}) P(\boldsymbol{\tau}) \prod_i P(G_i | \boldsymbol{\theta}, \boldsymbol{\tau}) P(X_i | G_i)$$

Gamma priors

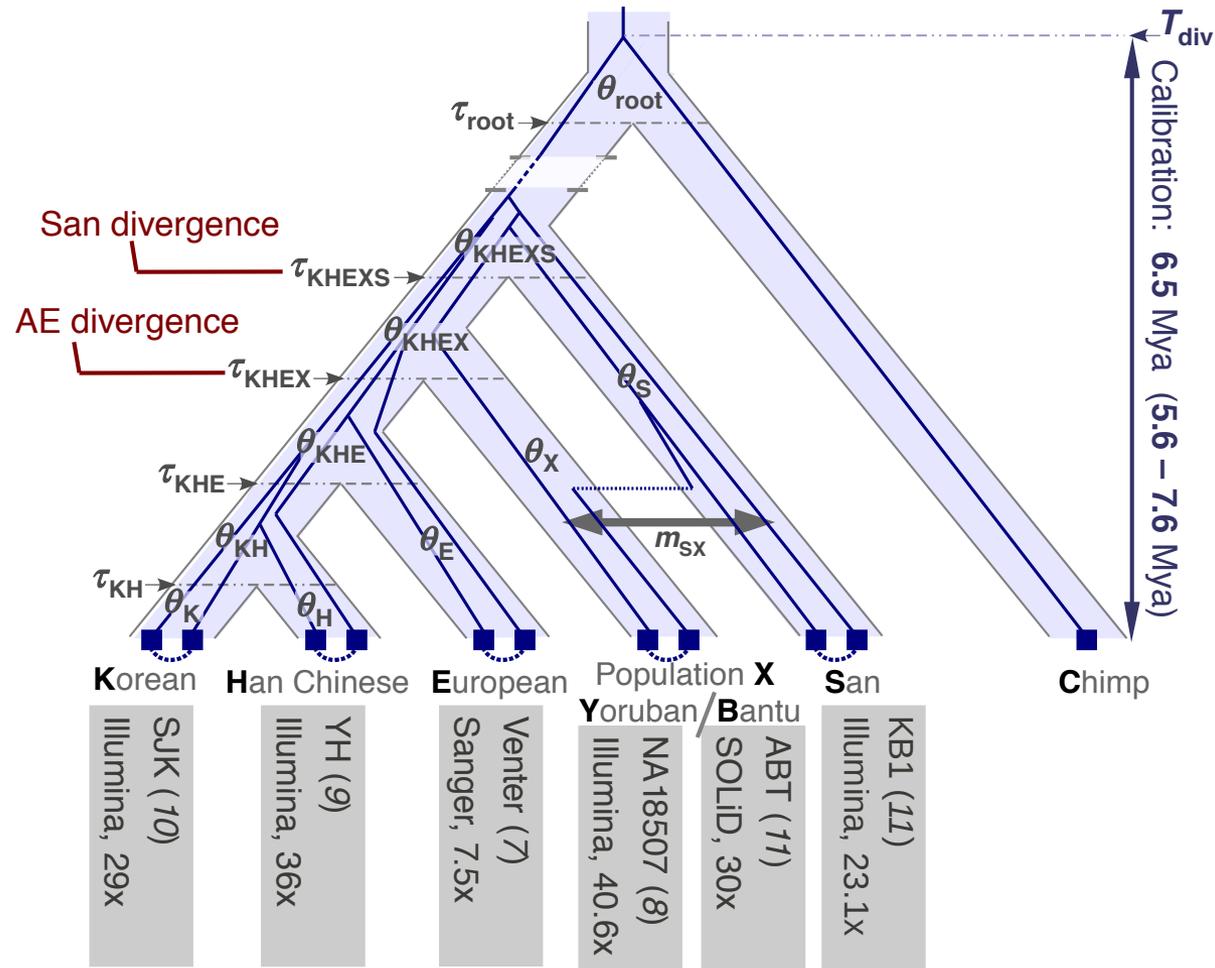
Censored
coalescent model

Finite sites model,
Felsenstein's pruning
algorithm

Goal: $P(\boldsymbol{\theta}, \boldsymbol{\tau} | \mathbf{X})$

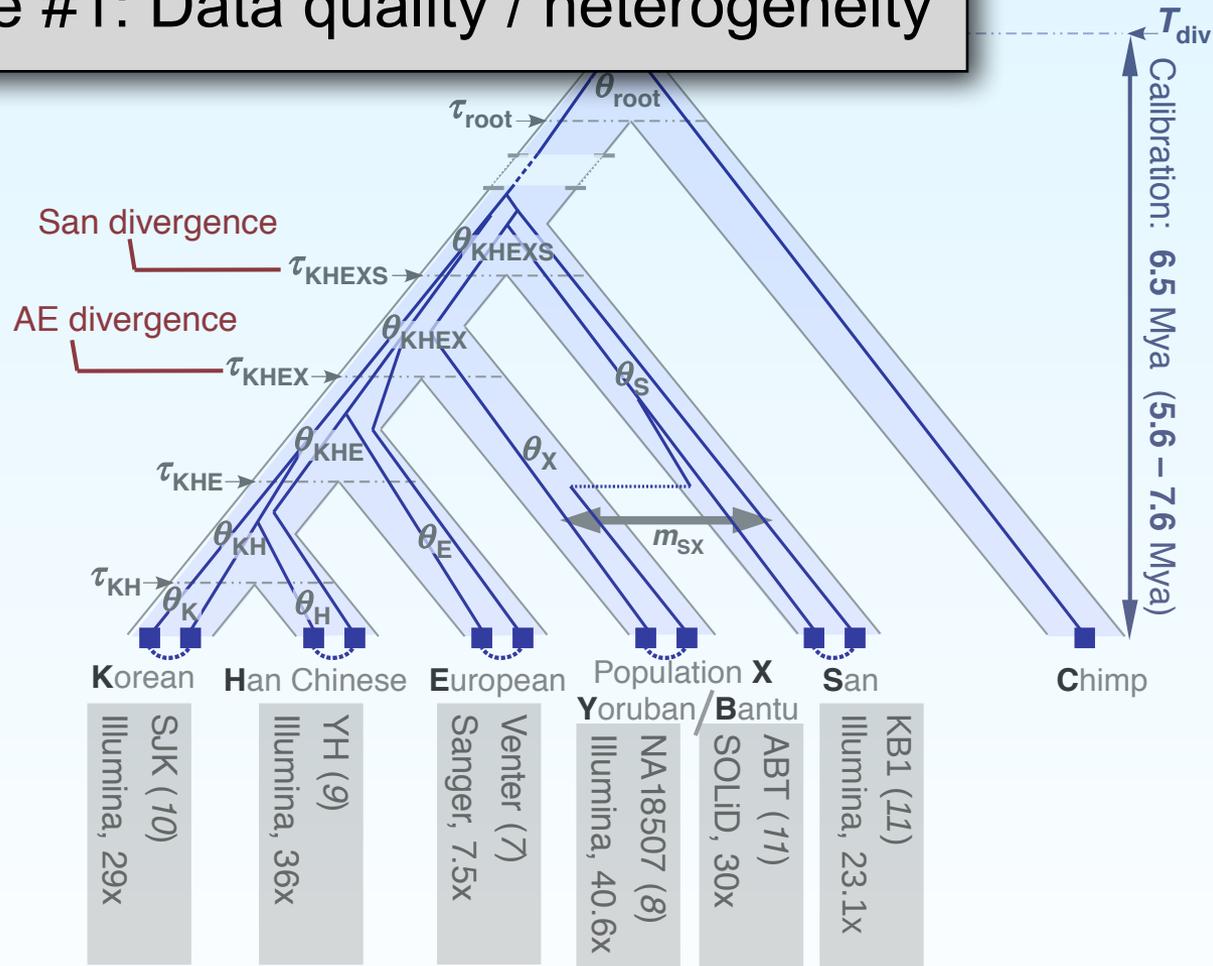


Data and Phylogeny



Data and Phylogeny

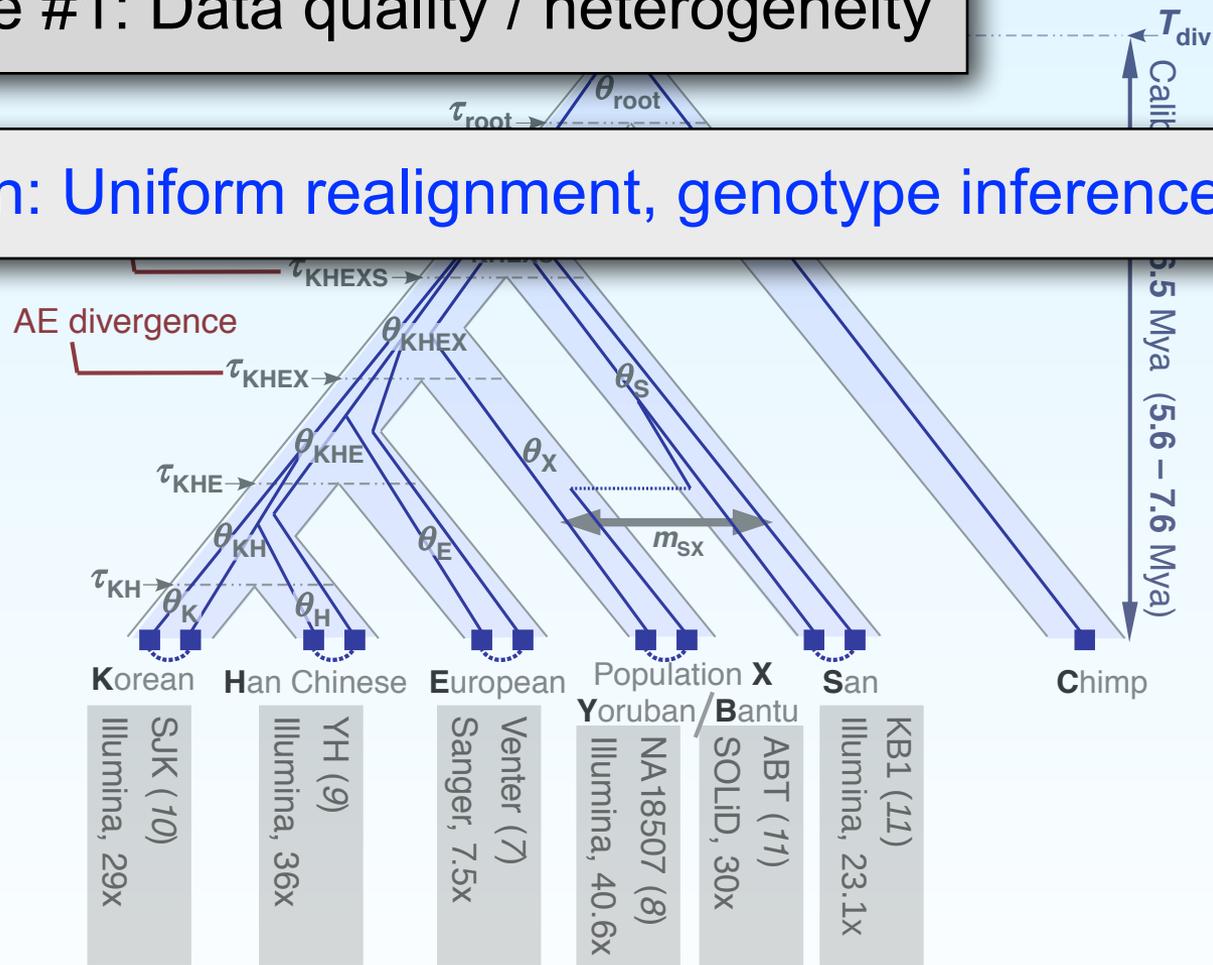
Challenge #1: Data quality / heterogeneity



Data and Phylogeny

Challenge #1: Data quality / heterogeneity

Solution: Uniform realignment, genotype inference

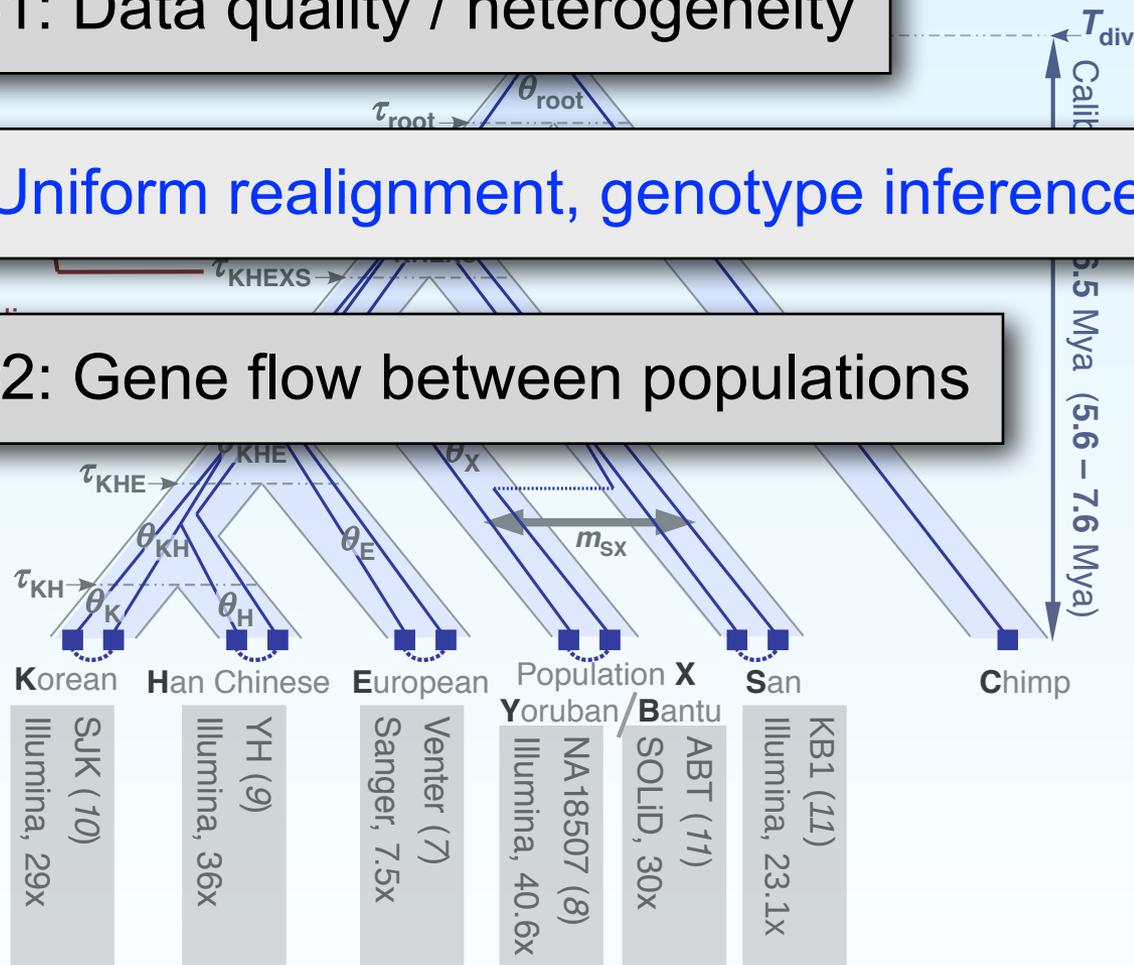


Data and Phylogeny

Challenge #1: Data quality / heterogeneity

Solution: Uniform realignment, genotype inference

Challenge #2: Gene flow between populations



Data and Phylogeny

Challenge #1: Data quality / heterogeneity

Solution: Uniform realignment, genotype inference

Challenge #2: Gene flow between populations

Solution: IM-like genealogy sampling (G-PhoCS)



Data and Phylogeny

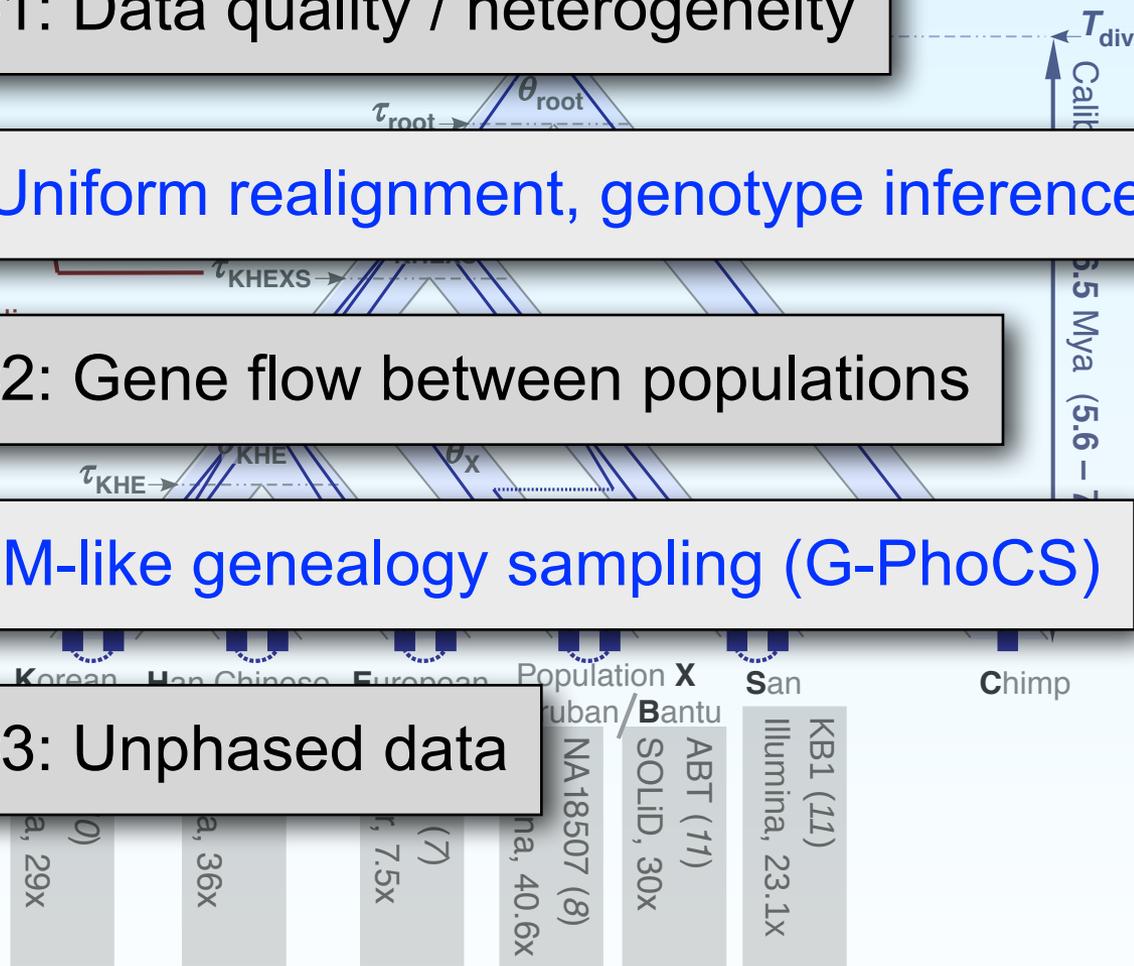
Challenge #1: Data quality / heterogeneity

Solution: Uniform realignment, genotype inference

Challenge #2: Gene flow between populations

Solution: IM-like genealogy sampling (G-PhoCS)

Challenge #3: Unphased data



Data and Phylogeny

Challenge #1: Data quality / heterogeneity

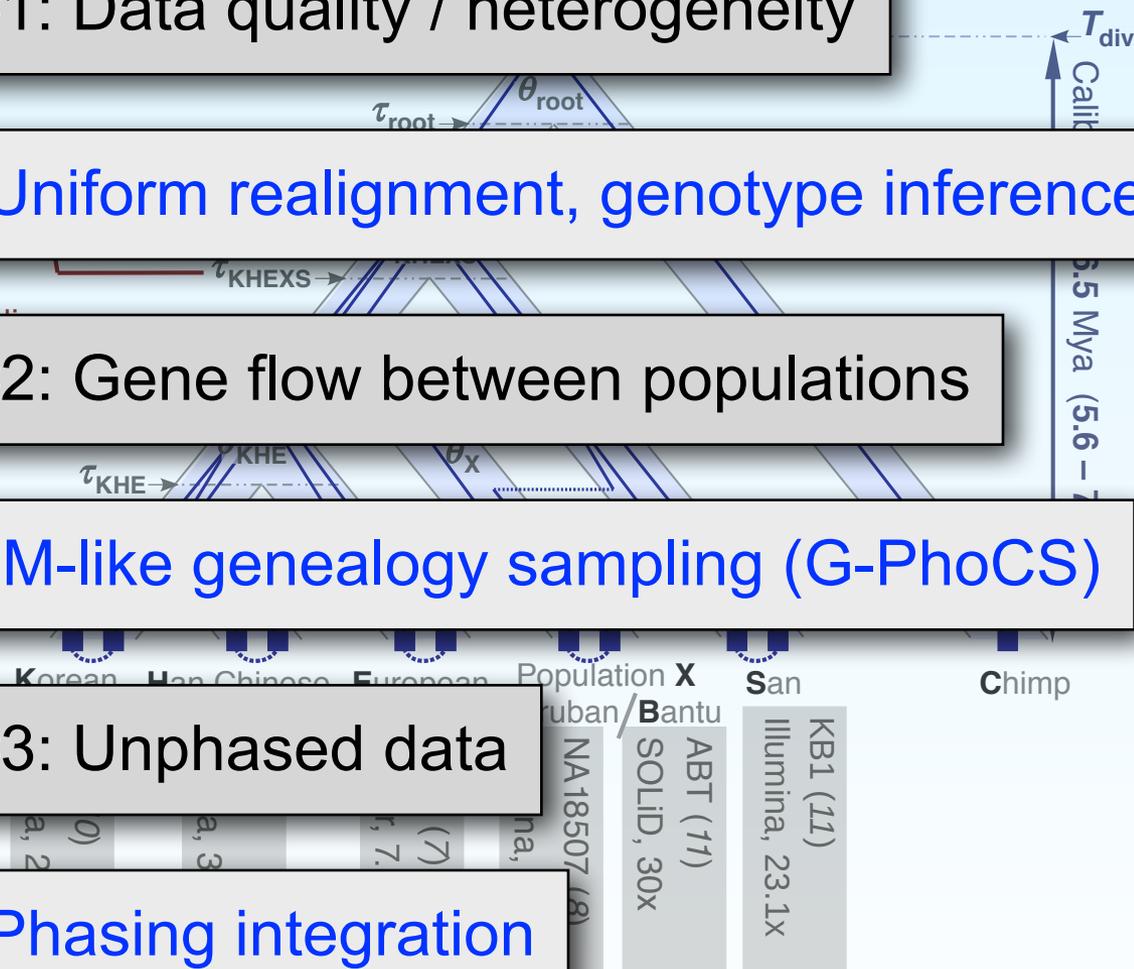
Solution: Uniform realignment, genotype inference

Challenge #2: Gene flow between populations

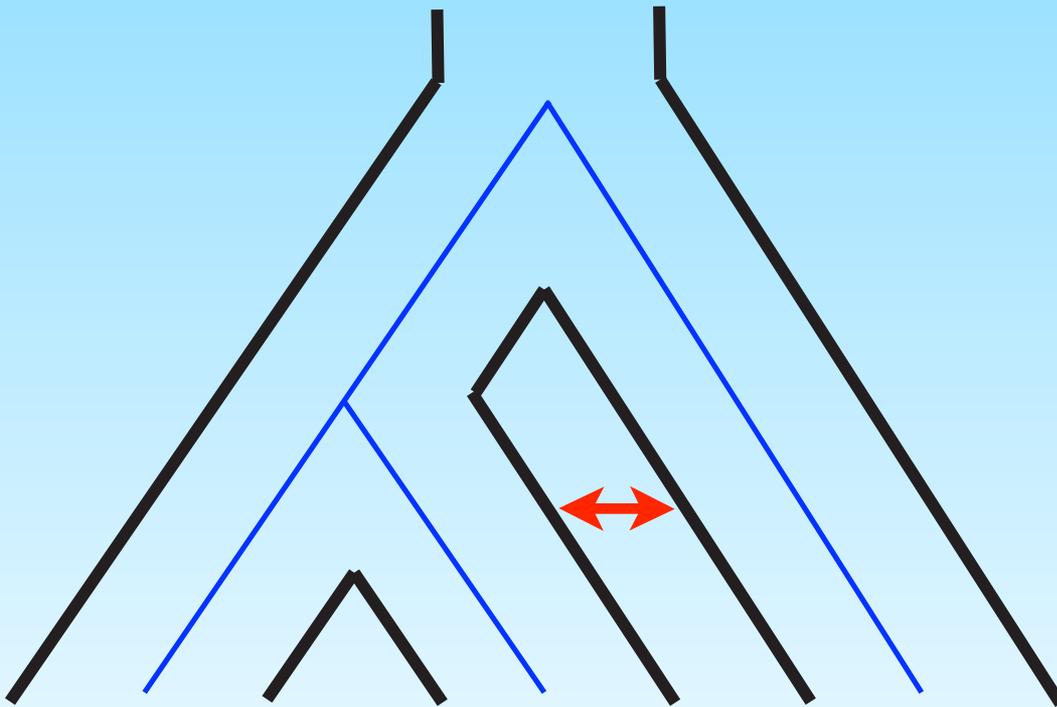
Solution: IM-like genealogy sampling (G-PhoCS)

Challenge #3: Unphased data

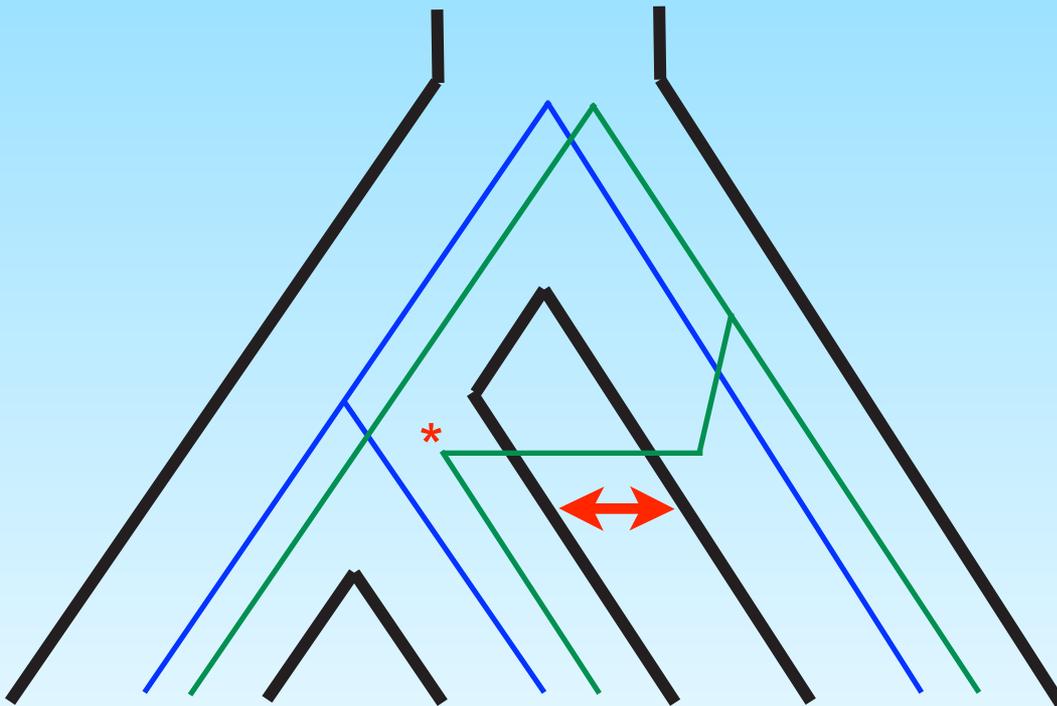
Solution: Phasing integration



Allowing for Migration

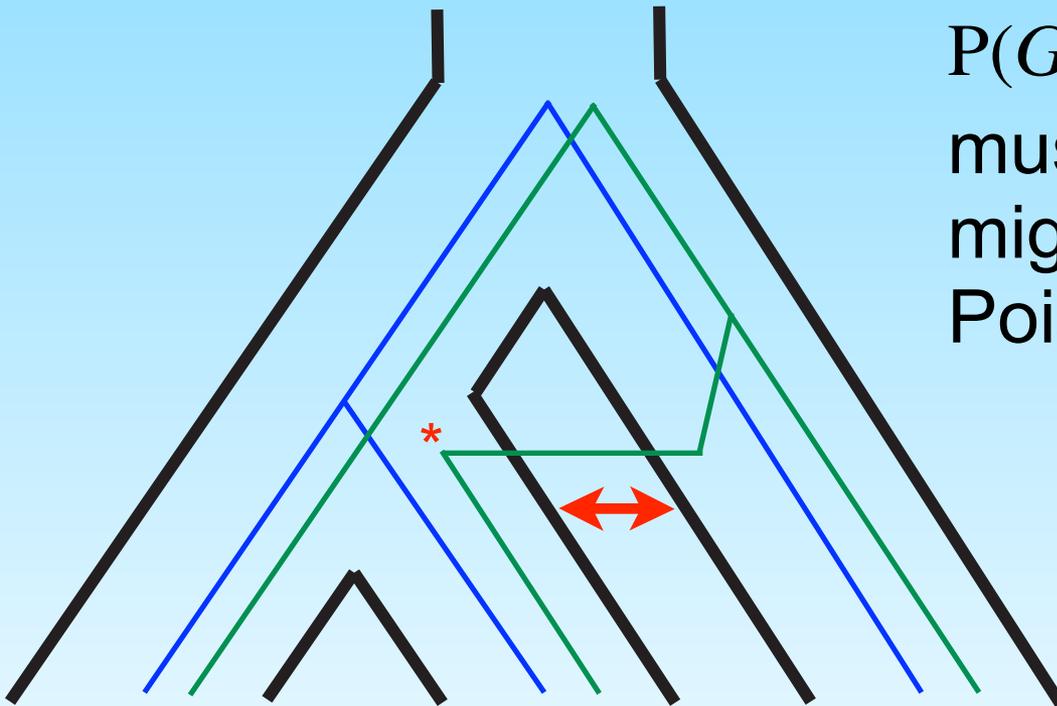


Allowing for Migration

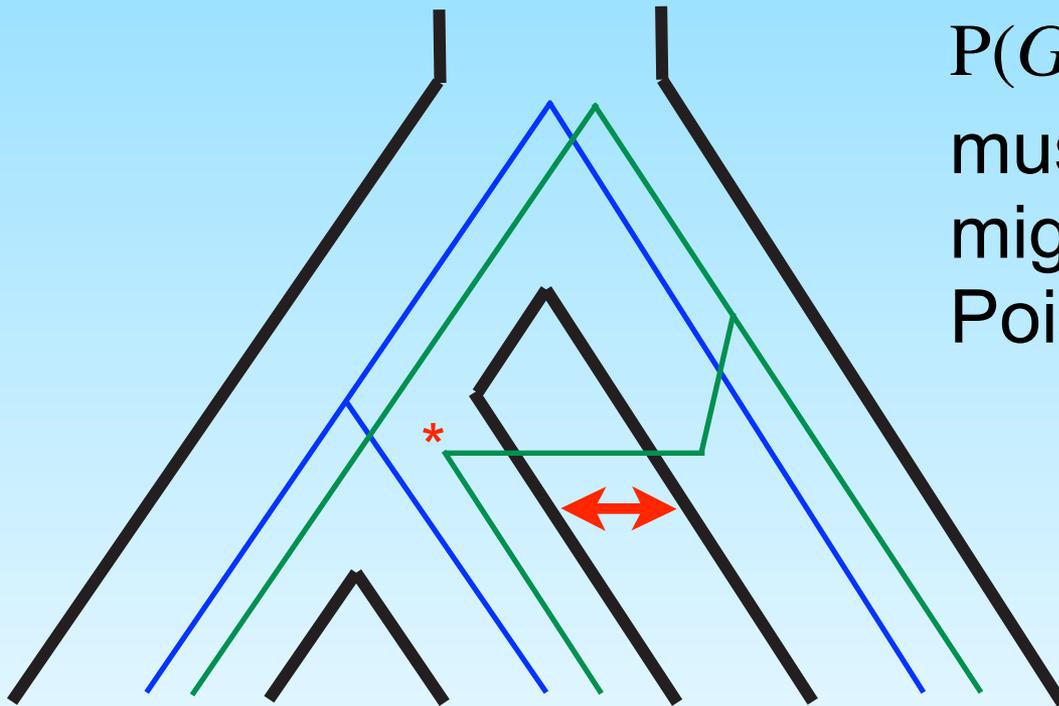


Allowing for Migration

The *genealogy prior*,
 $P(G_i | \theta, \tau, \mathbf{m})$,
 must allow for
 migration (another
 Poisson process)



Allowing for Migration



The *genealogy prior*,
 $P(G_i | \theta, \tau, \mathbf{m})$,
 must allow for
 migration (another
 Poisson process)

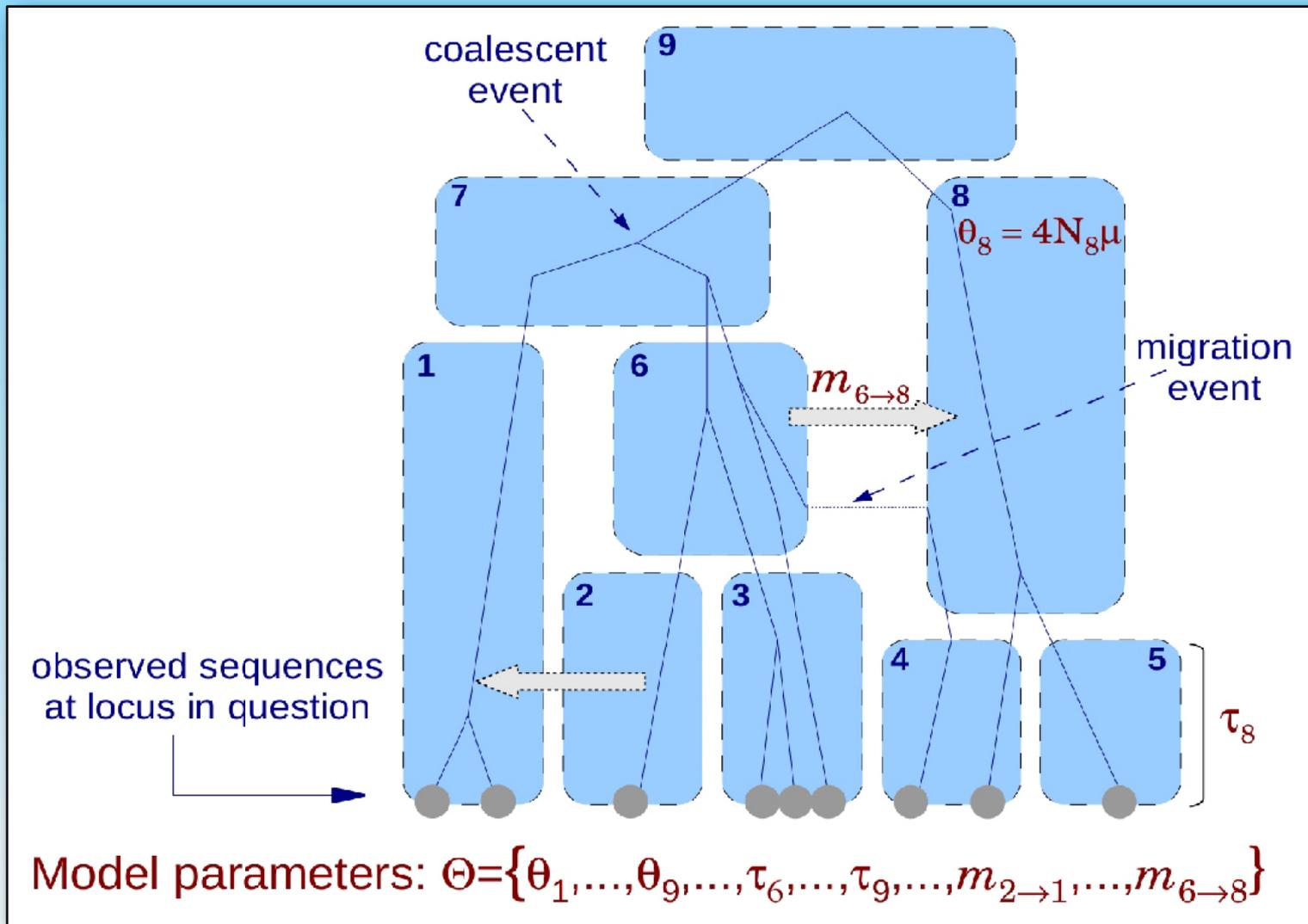
The *locus data
 likelihood*,
 $P(X_i | G_i)$,
 is unchanged

Phasing Integration

$$\begin{aligned} P(X|G) &= \sum_{\mathcal{P} \in \{0,1\}^{k \times n}} P(X|\mathcal{P}, G) P(\mathcal{P}) \\ &= \prod_j \left(\frac{1}{2^{|\mathcal{H}_j|}} \sum_{\mathcal{P}^j \in \{0,1\}^{|\mathcal{H}_j|}} P(X^j | \mathcal{P}^j, G) \right) \end{aligned}$$



General Model



MCMCcoal Algorithm

- **Step 1:** Update coalescent times
- **Step 2:** Subtree pruning and regrafting of genealogy
- **Step 3:** Update θ 's
- **Step 4:** Update τ 's, adjusting associated coalescent times via “rubber band”
- **Step 5:** Global scaling of θ 's and τ 's (mixing step)

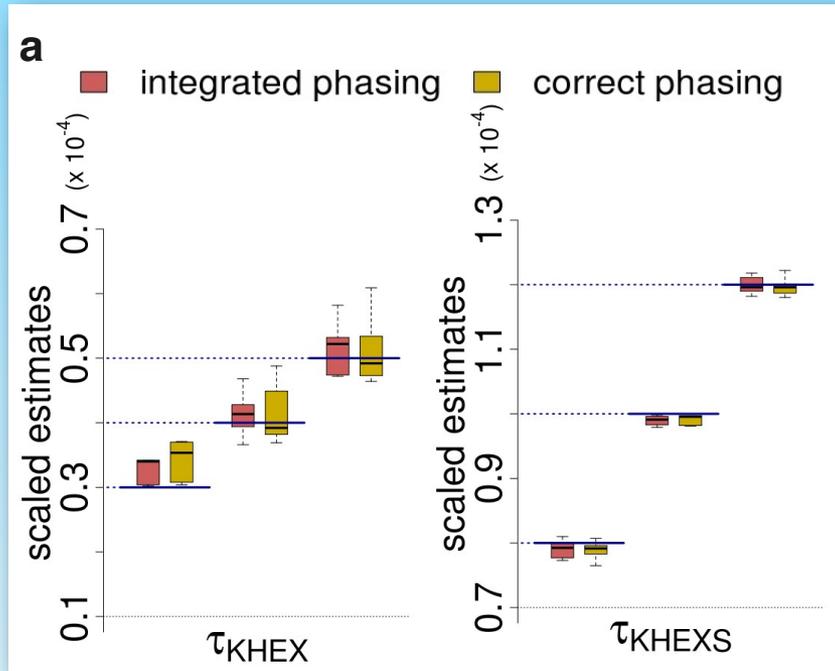


Changes to Sampler

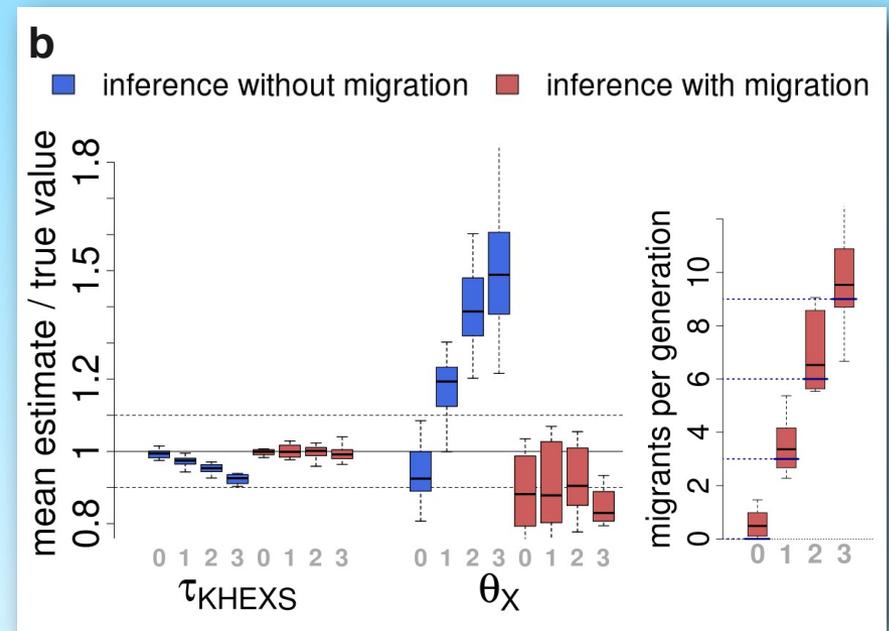
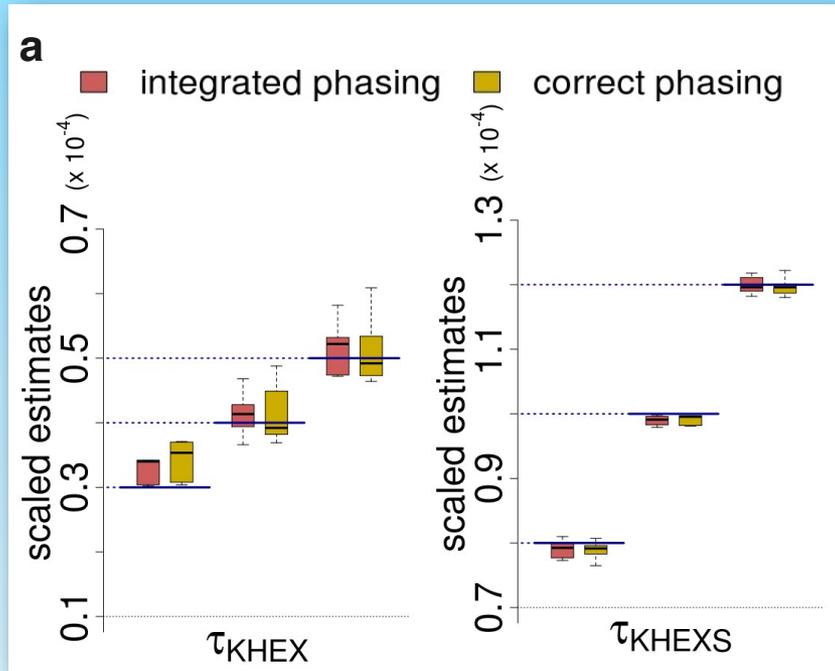
- The “regrafting” operation now allows lineages to pass through migration bands
- New steps are needed to update the individual migration times and the global migration parameters
- Certain additional conflicts must be considered when updating population divergence times



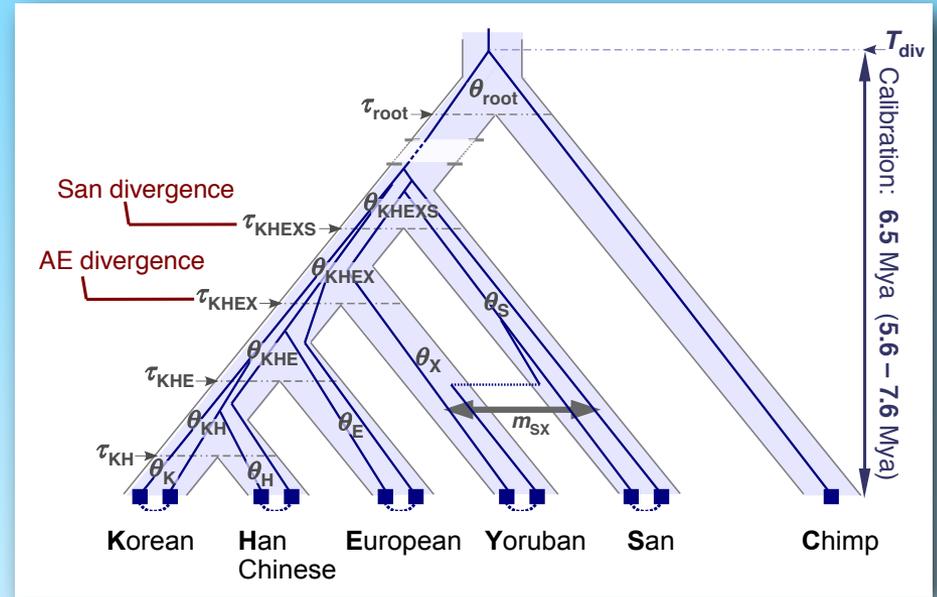
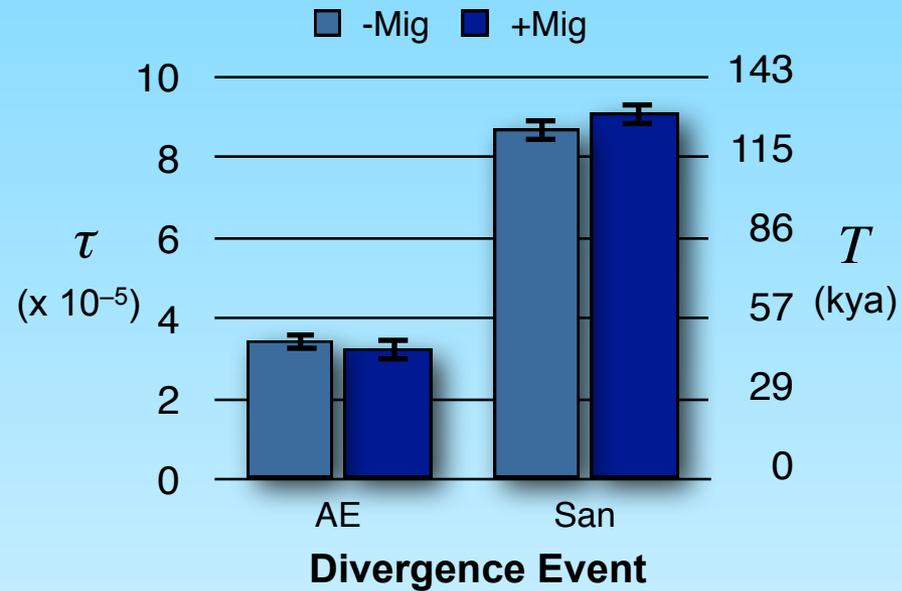
Simulation Results



Simulation Results



Main Results



37,574 “neutral” loci, each
1 kbp in length

Application to Domestic & Wild Canids



Boxer



Basenji



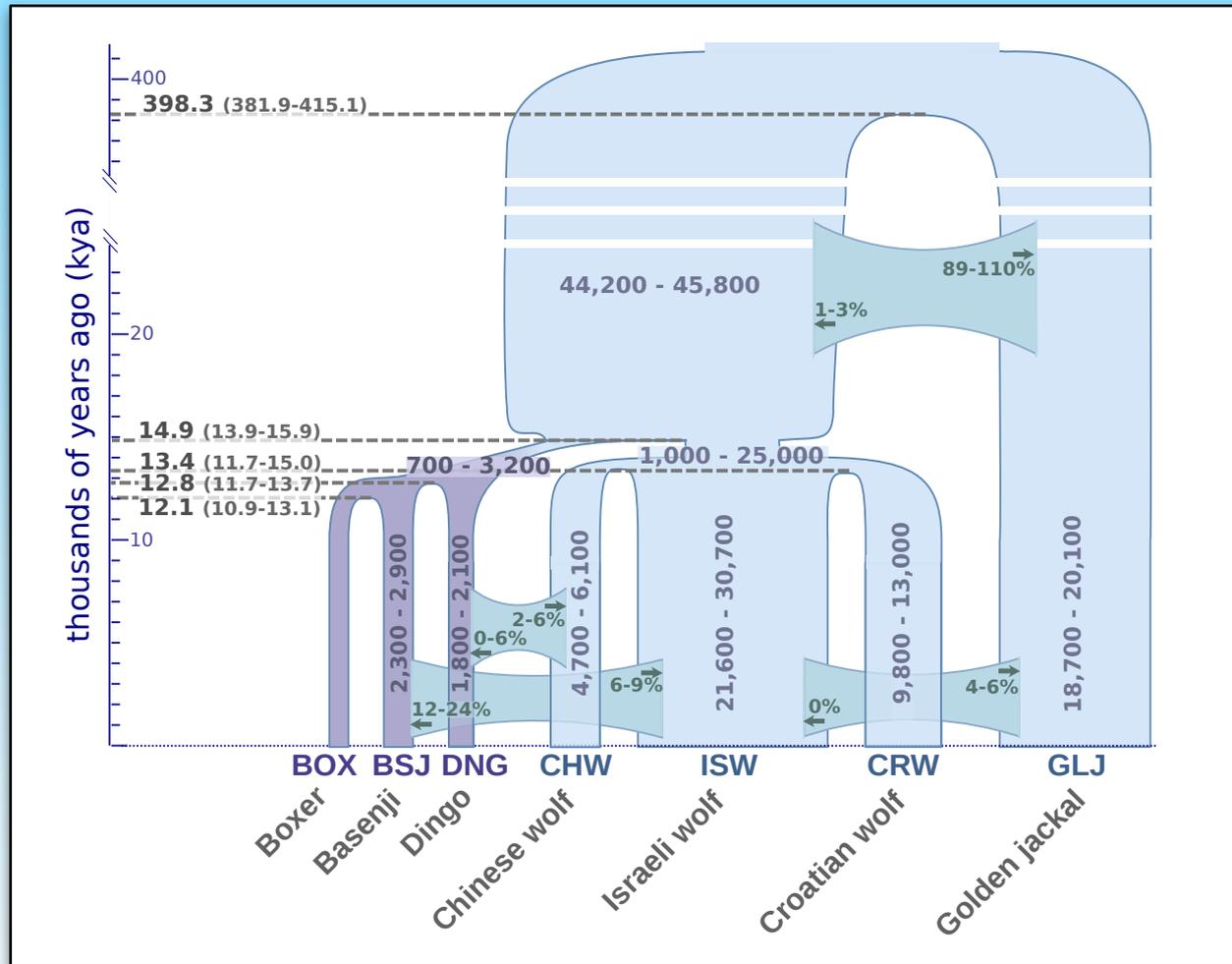
Dingo

Gray wolf
x 3

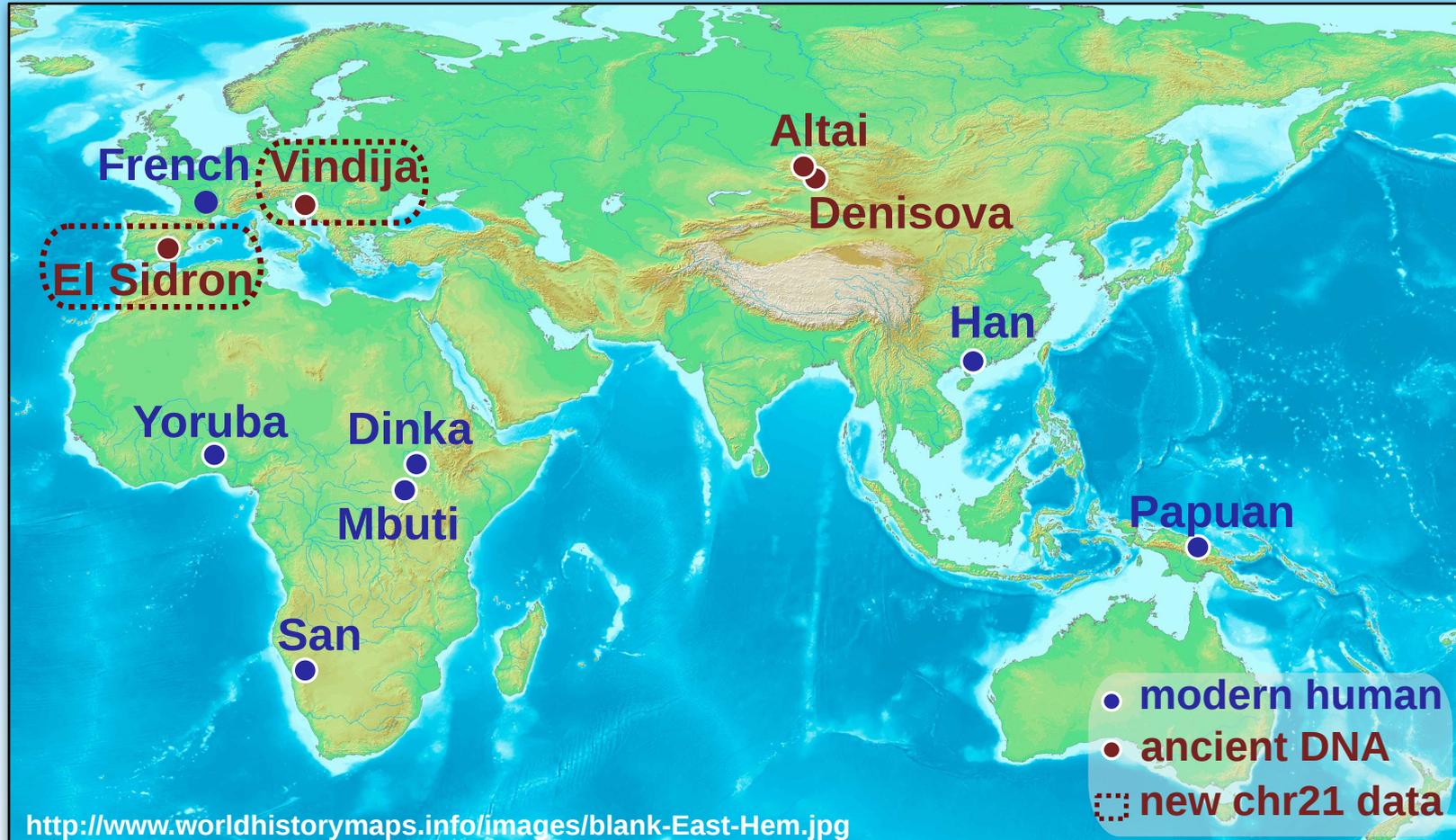
Golden jackal



Best Model



Archaic Hominin Analysis

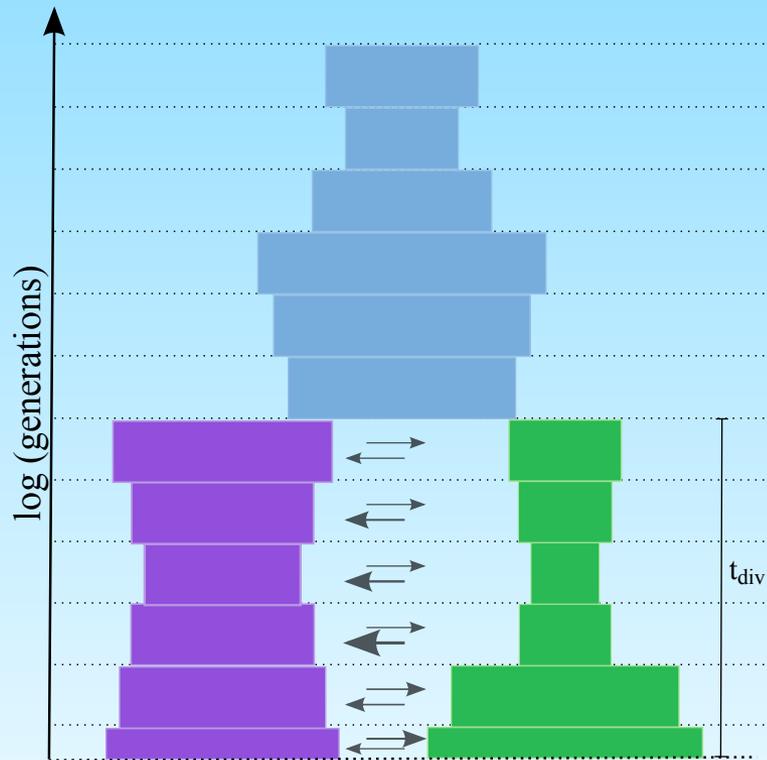


Limitations of G-PhoCS

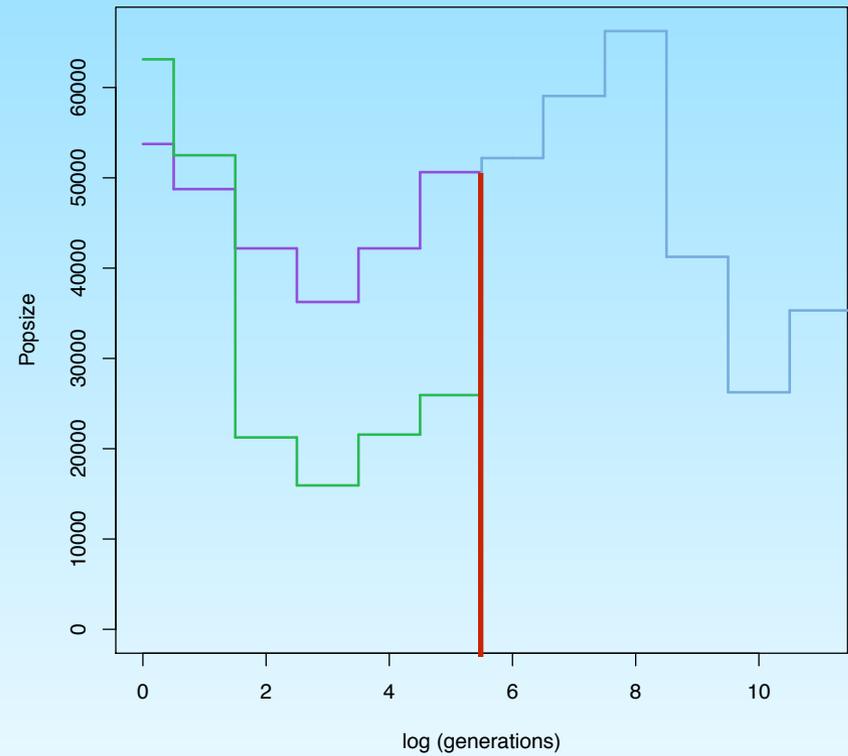
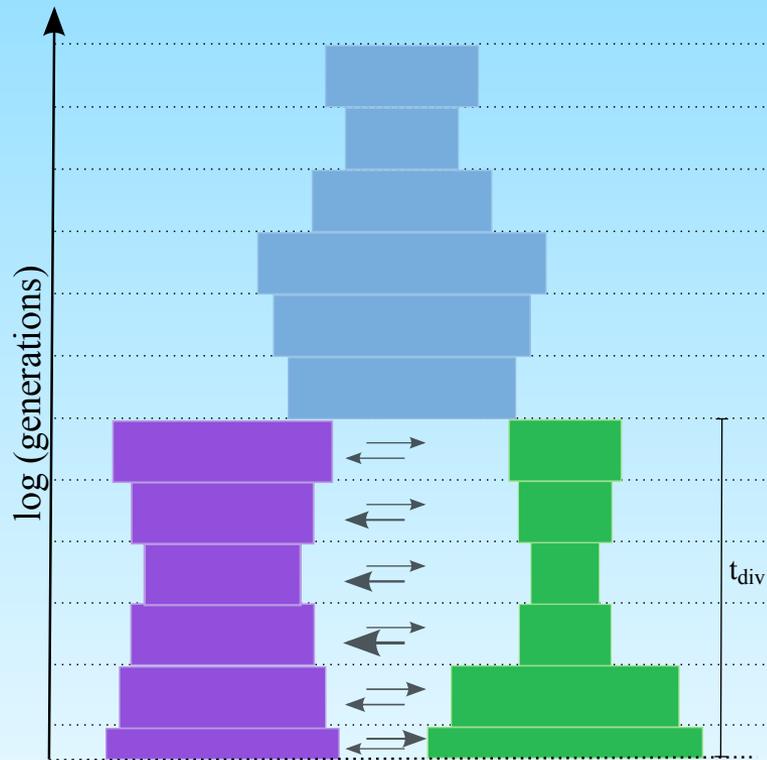
- Discards most of the data
- Must use short loci due to restrictive assumption of no intralocus recombination
- Fails to benefit from demographic information in LD structure
- Want to use full ARG!



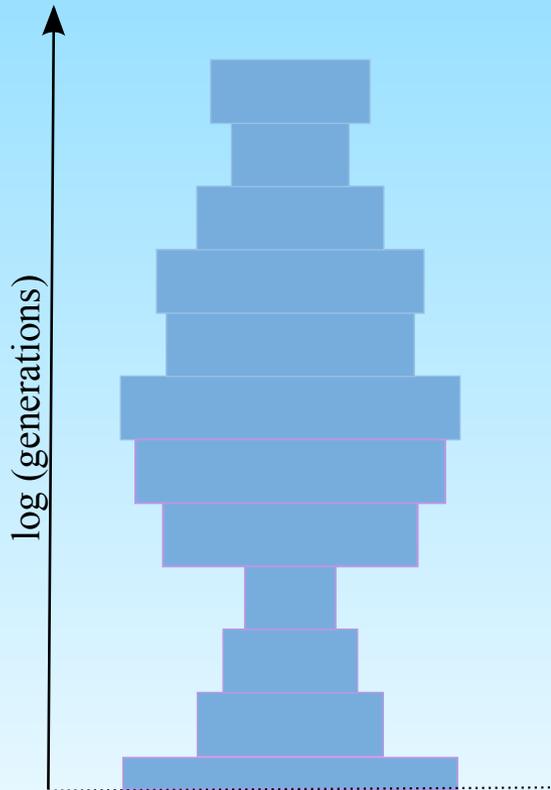
First Goal: IM + *ARGweaver*



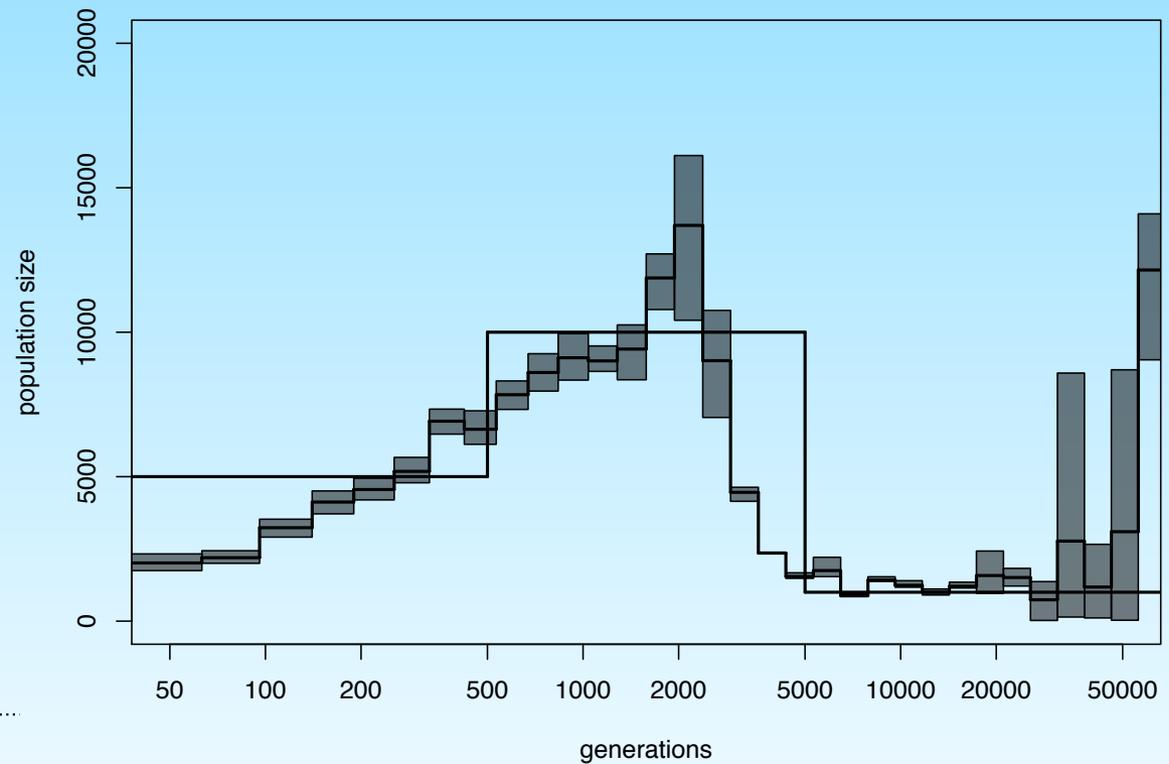
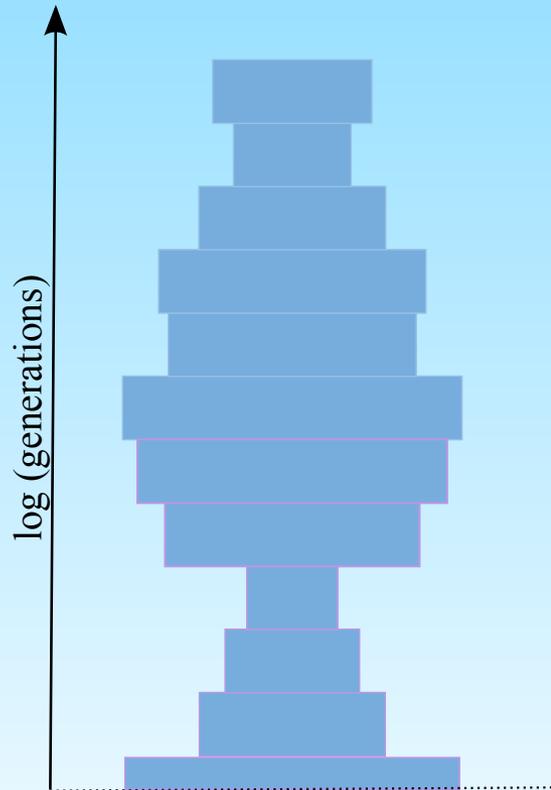
First Goal: IM + *ARGweaver*



Preliminary Results



Preliminary Results





Simons Center for Quantitative Biology at Cold Spring Harbor Laboratory

- Recently launched Center, with founding donation from Simons Foundation
- Focus on several areas of QB, including genomics, gene regulation, cancer biology, and neuroscience
- Faculty & fellow positions opening soon
- Several postdoc positions in my group
- See me if interested!



Acknowledgments

Contributors: Matthew Rasmussen, Melissa Hubisz, Ilan Gronau

Other Group Members: Charles Danko, Andre Martins, Lenore Pipes, Brad Gulko, Jaaved Mohammed

Collaborators: John Novembre, Adam Freedman, Bob Wayne, Sergi Castellano, Martin Kuhlwilm, Svante Paabo

