

Advanced Bayesian Phylogenetics: Phyloalignment

Philippe Lemey and Marc A. Suchard

Rega Institute

Department of Microbiology and Immunology

K.U. Leuven, Belgium, and

Departments of Biomathematics and Human Genetics

David Geffen School of Medicine at UCLA

Department of Biostatistics

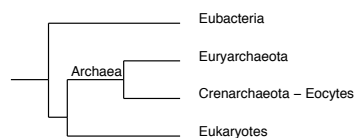
UCLA School of Public Health

SISMID – p.1

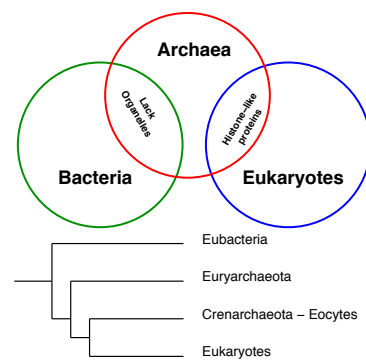
Resolving Early Branches in the Tree of Life

3? Domains of Life (Woese et al. 1990)

Contentious issue in genomics:
Do the Archaea form a single clade
(Rivera and Lake, 1992)?



Archaeal Tree



Eocyte Tree

Early evidence based on phylogenetic reconstruction techniques:

- Model how biologic sequences mutated over time
- Infer branching patterns based on "shared" substitutions

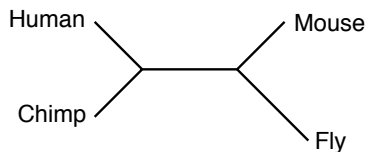
SISMID – p.2

Traditional Phylogenetic Reconstruction

Reconstruction Example

Human	-	T	C	C	T	G	G	A	A	T
Chimp	-	A	C	C	T	G	G	A	A	T
Mouse	-	A	A	C	T	-	-	T	A	T
Fly	-	A	A	G	A	T	C	G	T	A
Site:	1	2	3	4	5	6	7	8	9	10

Along Molecular Sequence



- **Substitution:** single residue replaces another
- **Insertion/deletion:** residues are inserted or deleted

Statistical Model

Assume: Homologous sites are iid and site patterns (e.g. dotted box)

$$XY \dots Z \sim \text{Multinomial}(p_{XY \dots Z})$$

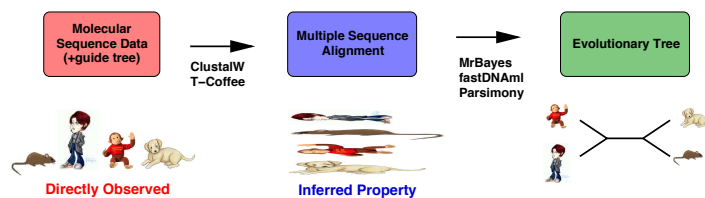
where $p_{XY \dots Z}$ is determined by an unknown tree τ and a continuous-time Markov chain model (for residue substitution) given by infinitesimal rate matrix Q

$$P(X \rightarrow Y \text{ in time } t) = \left\{ e^{tQ} \right\}_{XY}$$

Calculating $p_{XY \dots Z}$ integrates out unobserved states (internal nodes, gaps).

Fundamental Difficulty: Sequential Estimation

Current phylogenetic reconstruction methods:



Issues: Poor alignment biases phylogeny (Lake data: $EF-1\alpha/Tu$)

- Use guide tree and naive evolutionary models (**Trouble!**)

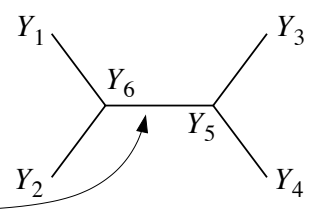
Solution: Infer alignment and phylogeny **simultaneously**

Previous approaches: Limited

- Optimization alignment, parsimony-based
- TFK91/92, forbidden positional homologies, inefficient

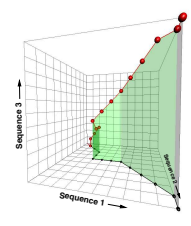
Alignment as a Random Variable

	Y	M(A)	f
Observed Data	$Y_1 = (A, T, T, C)$	1 2 - 3 4	A T - T C
	$Y_2 = (A, T, T, G)$	1 2 - 3 4	A T - T G
	$Y_3 = (T, C, T, G)$	- 1 2 3 4	- T C T G
	$Y_4 = (T, C, T)$	- 1 2 3 -	- T C T -
Missing Data	$Y_5 = (*, *, *)$	- 1 2 3 -	- * * * -
	$Y_6 = (*, *, *, *)$	1 2 - 3 4	* * - * *



Just over **1 billion** possible alignments for Y_{obs}

Explore space via **Forward-Backward algorithm (DP)** (Scott, 2002) to consider all possible alignments (and phylogenies) in polynomial time, weighted by posterior probability

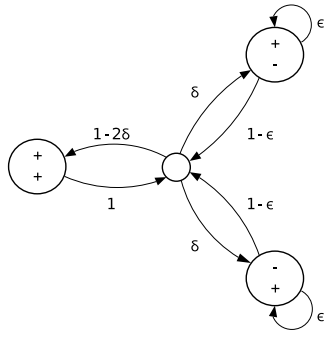
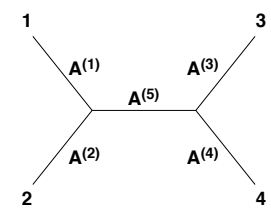


Note substitution process depends only on $Y_{obs} \Rightarrow$ separates substitution and indel processes into (substitution likelihood \times gap prior)

Gap Model along a Branch

Let the **multiple alignment** $A = (A^{(1)}, \dots, A^{(B)})$

- A is composed of pairwise alignments along each branch
- Pairwise alignment distribution follows a pair hidden Markov model (pair-HMM) **conditional** on equal sequence lengths at internal nodes



Pair-HMM parameterized by $\Lambda = (\delta, \epsilon)$

- δ : Probability of indel
- ϵ : Probability of extending an indel

Affine gap penalty $\approx [\log \delta] + (\ell - 1)[\log \epsilon]$

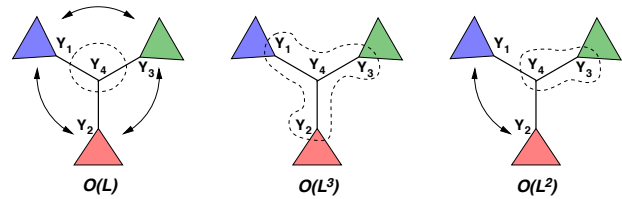
Choosing the Blocks: Efficient Sampling

Gibbs cycle over smaller blocks in alignment A :

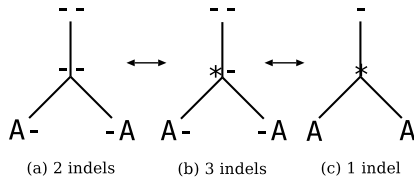
$$A^{(1)}, A^{(2)}, (A^{(3)}, A^{(4)}), (A^{(5)}, A^{(6)}, A^{(7)}), \dots$$

Let L = length of alignment:

- $O(L)$, too cold (Holmes and Bruno 2001)
- $O(L^3)$, too hot (Jensen and Hein 2005)
- $O(L^2)$, just right?



Possible **poor mixing** with $O(L)$ algorithm:

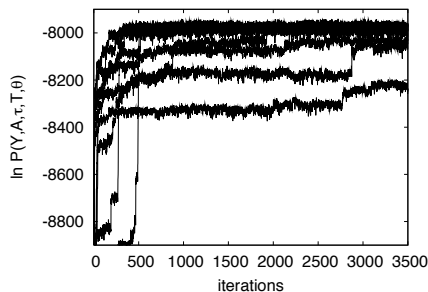


Must move through unfavorable intermediate to align/unalign sequence fragments

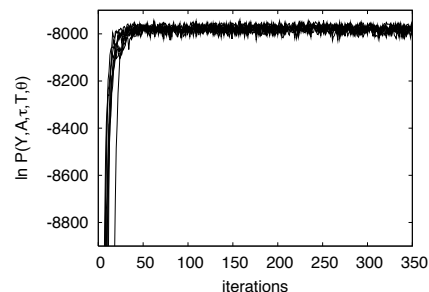
SISMD - p.7

Improved Alignment Mixing with $O(L^2)$ Sampling

$O(L)$ Only



$O(L)$ and $O(L^2)$



Enabling MCMC transition kernel decreases:

- Burn-in time
- Autocorrelation

Result: $> 70\times$ improvement shown here (12-taxon EF-1 α /Tu dataset).

SISMD - p.8

Sequential vs. Simultaneous Illustration

SIVmac251 partial *env* sequences from Cheynier et al (2001)

Sequential (ClustalW) alignment of hypervariable region:

```

***  *****  *****  **  *****  *****  *****
ref  AAATCATCAACAATRAACAACAACAGCACCACCAACACCACCAATAACAACATCAACAAAGTCAATAGACATG
S1   AAATCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S10  AAACCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S11  AAATCATCAACAATRAACAACAACAGCACCACCAACACCACCAATAACAACATCAACAAAGTCAATAAACATG
S15  AAATCATCAACAACAACAACAACAG-----CACCACCAATAACAACATCAACAGAGTCAATAAACATG
S16  AAATCATCAAC---AACAAACAACAGCACCACCAACCAACAAACAACAACATCAACAAAGTCAATAAACATG
S20  AAATCATCAAC---AACAAACAACAGCACCACCAACCAACAAACAACAACATCAACAAAGTCAATAAACATG
S5   AAATCATCAAC---AACAAACA-----CAACAACACCAGTACAACATCAACAAAGTCAATAAACATG
S9   AAATCATCAAC---AACAA-----CAACAACCACCAAGTACAACATCAACAAAGTCAATAAACATG

```

Simultaneous MAP alignment:

```

***  *****  *****  *  **  *****  *****  *****
ref  AAATCATCAACAATRAACAACAACAGCACCACCAACACCACCAATAACAACATCAACAAAGTCAATAGACATG
S1   AAATCATCAACAACAACAACAACAGCATCAACAACACCA-----CATCAACAAAGTCAATAAACATG
S10  AAACCATCAACAACAACAACAACAGCATCAACAACCA-----CATCAACAAAGTCAATAAACATG
S11  AAATCATCAACAATRAACAACAACAGCACCACCAACACCACCAATAACAACATCAACAAAGTCAATAAACATG
S15  AAATCATCAACAACAACAACAACAGACC-----AAATACAACATCAACAGAGTCAATAAACATG
S16  AAATCATCAACAA---CAACAACAACAGCACCACCAACCAACAAACAACAACATCAACAAAGTCAATAAACATG
S20  AAATCATCAACAA---CAACAACAACAGCACCACCAACCAACAAACAACAACATCAACAAAGTCAATAAACATG
S5   AAATCATCAACAACAACAACAACACC-----AAGTACAACATCAACAAAGTCAATAAACATG
S9   AAATCATCAACAACAACAACAACCA-----CC-----AAGTACAACATCAACAAAGTCAATAAACATG

```

Sequential vs. Simultaneous Illustration

SIVmac251 partial *env* sequences from Cheynier et al (2001)

Sequential (ClustalW) alignment of hypervariable region:

```

***  *****  *****  **  *****  *****  *****
ref  AAATCATCAACAATRAACAACAACAGCACCACCAACACCACCAATAACAACATCAACAAAGTCAATAGACATG
S1   AAATCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S10  AAACCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S11  AAATCATCAACAATRAACAACAACAGCACCACCAACACCACCAATAACAACATCAACAAAGTCAATAAACATG
S15  AAATCATCAACAACAACAACAACAG-----CACCACCAATAACAACATCAACAGAGTCAATAAACATG
S16  AAATCATCAAC---AACAAACAACAGCACCACCAACCAACAAACAACAACATCAACAAAGTCAATAAACATG
S20  AAATCATCAAC---AACAAACAACAGCACCACCAACCAACAAACAACAACATCAACAAAGTCAATAAACATG
S5   AAATCATCAAC---AACAAACA-----CAACAACACCAGTACAACATCAACAAAGTCAATAAACATG
S9   AAATCATCAAC---AACAA-----CAACAACCACCAAGTACAACATCAACAAAGTCAATAAACATG

```

Sampled alignment (1):

```

***  *****  *****  **  *****  *****  *****
ref  AAATCATCAACAATRAACAACAACAGCACCACCAACACCACCAATAACAACATCAACAAAGTCAATAGACATG
S1   AAATCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S10  AAACCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S11  AAATCATCAACAATRAACAACAACAGCACCACCAACACCACCAATAACAACATCAACAAAGTCAATAAACATG
S15  AAATCATCAACAACAACAACAACAG-----CACCACCAATAACAACATCAACAGAGTCAATAAACATG
S16  AAATCATCAACAAC---AACAAACAACAGCACCACCAACCAACAAACAACAACATCAACAAAGTCAATAAACATG
S20  AAATCATCAACAAC---AACAAACAACAGCACCACCAACCAACAAACAACAACATCAACAAAGTCAATAAACATG
S5   AAATCATCAACAACAACAACAACAA-----CACCACCAAGTACAACATCAACAAAGTCAATAAACATG
S9   AAATCATCAACAACAACAACAACAA-----CAACAAC-----CACCACCAAGTACAACATCAACAAAGTCAATAAACATG

```

Sequential vs. Simultaneous Illustration

SIVmac251 partial *env* sequences from Cheynier et al (2001)

Sequential (ClustalW) alignment of hypervariable region:

```

***  *****  *****  **  *****  *****  *****
ref  AAATCATCAACAATRAACAACAACAGCACCACCAACCAACACCAAAATACAACATCAACAAAGTCAATAGACATG
S1   AAATCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S10  AAACCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S11  AAATCATCAACAATRAACAACAACAGCACCACCAACCAACAAATACAACATCAACAAAGTCAATAAACATG
S15  AAATCATCAACAACAACAACAACAG-----CACCAAATACAACATCAACAGAGTCAATAAACATG
S16  AAATCATCAAC-----AACAAACAACAGCACCACCAACCAACAAACAACATCAACAAAGTCAATAAACATG
S20  AAATCATCAAC-----AACAAACAACAGCACCACCAACCAACAAACAACATCAACAAAGTCAATAAACATG
S5   AAATCATCAAC-----AACAAACA-----CAACAACCAAGTACAACATCAACAAAGTCAATAAACATG
S9   AAATCATCAAC-----AACAA-----CAACACCACCAAGTACAACATCAACAAAGTCAATAAACATG

```

Sampled alignment (2):

```

***  *****  **  *****  *  *****  *****
ref  AAATCATCAACAATRAACAACAACAGCACCACCAACCAACACCAAAATACAACATCAACAAAGTCAATAGACATG
S1   AAATCATCAACAACAACAACAACAGCATCAACAACACCA-----ACATCAACAAAGTCAATAAACATG
S10  AAACCATCAACAACAACAACAACAGCATCAACAACCA-----ACATCAACAAAGTCAATAAACATG
S11  AAATCATCAACAATRAACAACAACAGCACCACCAACCAACAAATACAACATCAACAAAGTCAATAAACATG
S15  AAATCATCAACAACAACAACAACAGCACC-----AAATACAACATCAACAGAGTCAATAAACATG
S16  AAATCATCAACAA-----CAACAACAACAGCACCACCAACCAACAAACAACATCAACAAAGTCAATAAACATG
S20  AAATCATCAACAA-----CAACAACAACAGCACCACCAACCAACAAACAACATCAACAAAGTCAATAAACATG
S5   AAATCATCAACAACAACAACAACACC-----AAGTACAACATCAACAAAGTCAATAAACATG
S9   AAATCATC-----AACAAACAACAACACC-----AAGTACAACATCAACAAAGTCAATAAACATG

```

Sequential vs. Simultaneous Illustration

SIVmac251 partial *env* sequences from Cheynier et al (2001)

Sequential (ClustalW) alignment of hypervariable region:

```

***  *****  *****  **  *****  *****  *****
ref  AAATCATCAACAATRAACAACAACAGCACCACCAACCAACACCAAAATACAACATCAACAAAGTCAATAGACATG
S1   AAATCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S10  AAACCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S11  AAATCATCAACAATRAACAACAACAGCACCACCAACCAACAAATACAACATCAACAAAGTCAATAAACATG
S15  AAATCATCAACAACAACAACAACAG-----CACCAAATACAACATCAACAGAGTCAATAAACATG
S16  AAATCATCAAC-----AACAAACAACAGCACCACCAACCAACAAACAACATCAACAAAGTCAATAAACATG
S20  AAATCATCAAC-----AACAAACAACAGCACCACCAACCAACAAACAACATCAACAAAGTCAATAAACATG
S5   AAATCATCAAC-----AACAAACA-----CAACAACCAAGTACAACATCAACAAAGTCAATAAACATG
S9   AAATCATCAAC-----AACAA-----CAACACCACCAAGTACAACATCAACAAAGTCAATAAACATG

```

Sampled alignment (3)

```

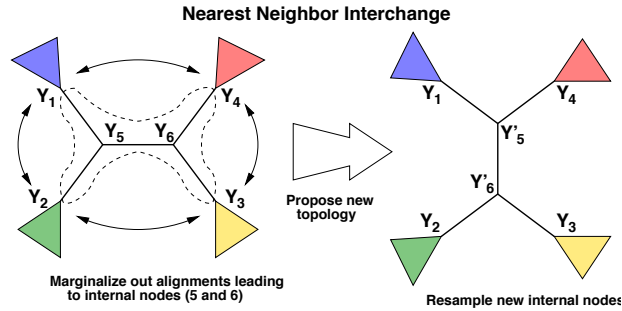
***  *****  *****  **  *****  *****
ref  AAATCATCAACAATRAACAACAACAGCACCACCAACCAACACCAAAATACAACATCAACAAAGTCAATAGACATG
S1   AAATCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S10  AAACCATCAACAACAACAACAACAGCATCAACAACAC-----CAACATCAACAAAGTCAATAAACATG
S11  AAATCATCAACAATRAACAACAACAGCACCACCAACCAACAAATACAACATCAACAAAGTCAATAAACATG
S15  AAATCATCAACAACAACAACAACAGCACC-----ATACAACATCAACAGAGTCAATAAACATG
S16  AAATCATCAAC-----AACAAACAACAGCACCACCAACCAACAAACAACATCAACAAAGTCAATAAACATG
S20  AAATCATCAAC-----AACAAACAACAGCACCACCAACCAACAAACAACATCAACAAAGTCAATAAACATG
S5   AAATCATCAACAACAACAACAACACC-----GTACAACATCAACAAAGTCAATAAACATG
S9   AAATCATC-----CAACAACAACAACACC-----GTACAACATCAACAAAGTCAATAAACATG

```


Trees and Alignments: Collapsed Gibbs Sampling

Problem: Tree and alignment are **highly** correlated

Further important aspect: Alignment-aware tree τ sampling



Generate $(\tau, \mathbf{A}) \mid \mathbf{Y}, \theta$ by

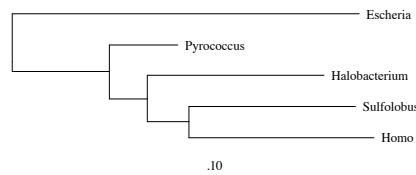
1. $\tau \mid \mathbf{Y}, \theta$ (collapsed)
2. $\mathbf{A} \mid \tau, \mathbf{Y}, \theta$

Similar procedure available for **global changes** (SPR moves)

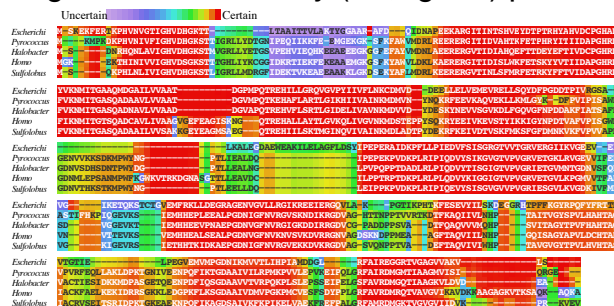
SISMID – p.10

EF-1 α /Tu Strongly Supports Eocyte Hypothesis

(*Homo, Sulfolobus*) clade supported at $\geq 99.9\%$ (sampling resolution):



Alignment Uncertainty (Au, “gold”) plot:

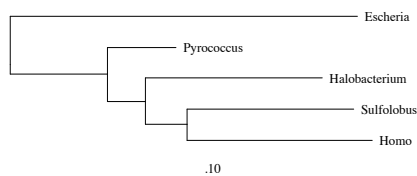


- Regions of marked homology (conservation)
- Uncertain regions
- Shared indels

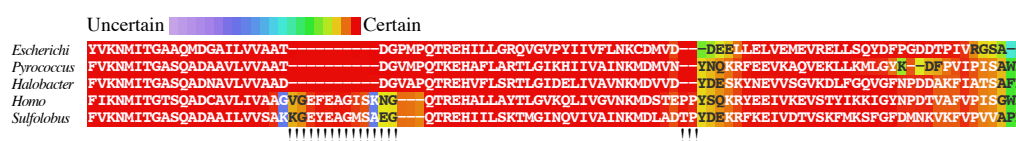
SISMID – p.11

EF-1 α /Tu Strongly Supports Eocyte Hypothesis

(*Homo, Sulfolobus*) clade supported at $\geq 99.9\%$ (sampling resolution):



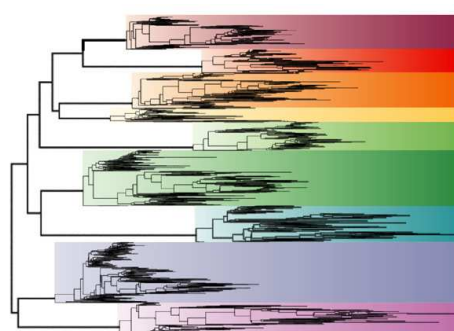
Alignment Uncertainty (Au, “gold”) plot:



- Automatic detection of indels shared by descent vs. by state
- Two indels shared by *Homo* and *Sulfolobus* contribute support for Eocyte Hypothesis

SISMID – p.11

Future Directions: Intra-Host Viral Evolution



Nature Reviews | Genetics

1195 *env* sequences from 9 HIV+ patients [taken from Rambaut et al. (2004)]

Retroviruses (and HBV) exist as a **quasi-species** within infected patients:

- Shared substitutions may be insufficient to resolve intra-host phylogenies

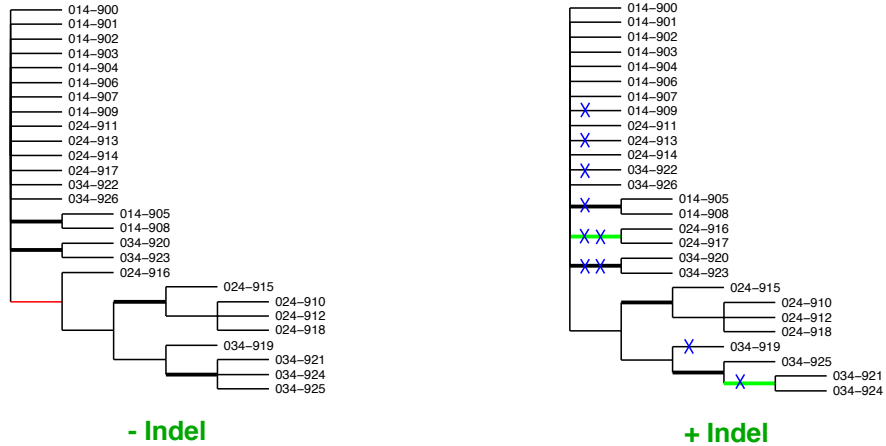
Improve resolution using joint model:

- Indel rates \geq substitution rates
- Opportunity to detect intra-host recombination

SISMID – p.12

Improved Resolution – I Can See!

Shankarrapa (1999) Pt #1: 3 time-points, **90% consensus** trees



- Indel events (X) ⇒ 2 additional bi-partitions supported.

SIS MID – p.13

Codon Models: Not Always A Good Thing

14-00 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-00 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
14-02 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-02 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
14-03 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-03 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
14-04 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-04 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
14-07 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-07 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
24-13 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	24-13 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
34-22 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	34-22 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
34-26 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	34-26 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
14-01 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-01 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
14-09 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-09 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
34-20 AGTACTGG-----CTTTGAATAATGTTACTGAA	34-20 AGTACTGG-----CTTTGAATAATGTTACTGAA
34-23 AGTACTGG-----CTTTGAATAATGTTACTGAA	34-23 AGTACTGG-----CTTTGAATAATGTTACTGAA
24-11 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	24-11 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
14-05 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-05 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
14-08 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-08 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
14-06 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	14-06 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
24-14 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	24-14 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
24-16 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	24-16 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
24-17 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	24-17 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
24-15 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	24-15 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
24-10 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	24-10 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
24-12 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	24-12 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
24-18 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	24-18 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
34-19 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	34-19 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
34-25 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA	34-25 AGTACTGGATAAATAGTACTTTGAATAATGTTACTGAA
34-21 AGTACTGGATA-----ACTTTGAATAATGTTACTGAA	34-21 AGTACTGGATA-----ACTTTGAATAATGTTACTGAA
34-24 GGTACTGGATA-----CTTTGAATAATGTTACTGAA	34-24 GGTACTGGATA-----ACTTTGAATAATGTTACTGAA

HKY×1 (-1556)

HKY×3 (-1580)

- Codon model M0: $\omega = 1.0 (0.9, 1, 2) \approx \text{HKY} \times 3$
- **Singlet** model more likely. **Triplet** model **shifts** indels and **mismatches** residues

SIS MID – p.14