

PROBABILISTIC PHYLOGENETIC
INFERENCE
WITH
INSERTIONS AND DELETIONS

E.R.
Sean Eddy

Generative birth-death model of sequence evolution

based on a $(N+1)$ Markov process

→ no memory effect (align/seqs of arbitrary length)

→ no fixed gap frequency

Assumes insertions occur one at the time

→ Allows an efficient $\Theta(n \cdot L)$ algorithm
for $P(\text{alignment} | \text{Tree, Model})$

Practical application to phylogenetic reconstruction

(phylib):

DNAML  DNAMLE

ignores gaps

includes gaps as per above model

same time complexity
same results for ungapped alignments

Test on real rRNA alignments

CONCORDANCE TEST

The $(N+1)$ Markov process

Durbin et al. '98

extended rate:

$$R^E = \left[\begin{array}{ccc|c} & \dots & & \mu \\ & R - \mu \delta_{ij} & & \vdots \\ & & & \mu \\ \hline \lambda \pi_1 & \dots & \lambda \pi_k & -\lambda \end{array} \right]$$

$${}^t R^E e = \left[\begin{array}{ccc|c} & & & \gamma_t \\ & P_t^E(j|i) & & \vdots \\ & & & \gamma_t \\ \hline \xi_t \pi_1 & \dots & \xi_t \pi_k & 1 - \xi_t \end{array} \right]$$

If R reversible:

Prob INSERTION: $\xi_t = \frac{\lambda}{\lambda + \mu} \left(1 - e^{-(\lambda + \mu)t} \right) \Big|_{\lambda = \mu = 0} = 0$

Prob DELETION: $\gamma_t = \frac{\mu}{\lambda + \mu} \left(1 - e^{-(\lambda + \mu)t} \right) \Big|_{\mu = \lambda = 0} = 0$

$$[{}^t R^E]_{ij} = P_t(j|i) = \pi_j + \sum_{a=1}^A q_a(i,j) e^{-e_a t}$$

R eigenvalues $\left\{ \begin{array}{l} 0 \longrightarrow 0 \\ -e_a \longrightarrow -(e_a + \mu) \\ \longrightarrow -(\lambda + \mu) \end{array} \right\}$ R^E eigenvalues

Prob SUBSTITUTION

$$P_t^E(j|i) = \frac{\lambda}{\lambda + \mu} \pi_j + \sum_{a=1}^A q_a(i,j) e^{-(e_a + \mu)t} + \frac{\mu}{\lambda + \mu} \pi_j e^{-(\lambda + \mu)t} \Big|_{\lambda = \mu = 0} = P_t(j|i)$$

McQuinn et al. '03 [$R = FB^4, \lambda = \beta, \mu = \nu$]

The Generative Model

Ignore the (=) transitions $(1 - \epsilon_t)$

t_1	$x:$	{	a	g	$\hat{x}: a - g$	$\hat{x}: a - g$	-	g	-	?
	$y:$		a	c	$\hat{y}: a - c - g$	$\hat{y}: a - c - g$		g	-	
t_2	$z:$	{	t	g	$\hat{z}: a - c - g - a$	$\hat{z}: t - g - a$		g	-	
				a					a	

→ $P_{t_1+t_2}(z | x, y) = P_{t_1}(y, \hat{x}\hat{y} | x) \cdot P_{t_2}(z, \hat{y}\hat{z} | y)$

Borrowed from Markov model

$$P_{t_1}(y, \hat{x}\hat{y} | x) \propto \overbrace{\epsilon_{t_1} \pi_c}^{ins} P_{t_1}^E(g | g) P_{t_1}^E(a | a)$$

$$P_{t_2}(z, \hat{y}\hat{z} | y) \propto \underbrace{\delta_{t_2}}_{del} \underbrace{\epsilon_{t_2} \pi_a}_{ins} \underbrace{P_{t_2}^E(g | g)}_{subs} P_{t_2}^E(t | a)$$

Normalization?

s-substitutions
d-deletions
i-insertions

$$P_t(l' = s+i | l = s+d) \propto \frac{(s+d+i)!}{s! d! i!} \delta_t^d (1 - \delta_t)^s \epsilon_t^i$$

$$\sum_{s=0}^{\infty} \sum_{d=0}^{\infty} \sum_{i=0}^{\infty} P(l' = s+i | l = s+d) = 1 \Rightarrow$$

$$P(l' | l) = (1 - \epsilon_t)^{l+1} \frac{(l+i)!}{s! (l-s)! (l'-s)!} \gamma_t^{l-s} (1 - \gamma_t)^s \epsilon_t^{l'-s}$$

The Generative Model

$$P_t(y_{e'}, \hat{x} \hat{y} | X_e) = (1 - \xi_t)^{\ell+1} \gamma_t^{\ell-s} \xi_t^{\ell'-s} \pi_{I_1} \dots \pi_{I_{\ell'-s}} \prod_{k=1}^s P_t^\varepsilon(y_{\sigma_k} | X_{\sigma_k})$$

Column Factorization

$$P_t(y \hat{x} \hat{y} | X) = (1 - \xi_t) \left(\prod_{i=1}^{\ell'-s} P_t^g(I_i | -) \right) \left(\prod_{d=1}^{\ell-s} P_t^g(- | X_{\sigma_d}) \right) \left(\prod_{k=1}^s P_t^g(y_{\sigma_k} | X_{\sigma_k}) \right)$$

$\ell'-s$ INS
 $\ell-s$ DEL
 s SUBS

extra term \rightarrow

$$\begin{cases} P_t^g(i | -) = \xi_t \pi_i \\ P_t^g(- | i) = (1 - \xi_t) \gamma_t \\ P_t^g(j | i) = (1 - \xi_t) P_t^\varepsilon(j | i) \end{cases}$$

Sequence generator

* a * g *

$(\ell+1)$ places: * $\xrightarrow{(1-\xi_t)\xi_t^\infty}$ * I, * ... * I_n *

$\xrightarrow{P_t^\varepsilon(a'|a)}$ a'

ℓ residues:

a $\xrightarrow{\gamma_t}$ death of right link (after it emits)

$$\langle e' | e \rangle = \ell(1 - \gamma_t) + (\ell+1) \frac{\xi_t}{1 - \xi_t}$$

$t=0 \rightarrow \ell$
 $t \rightarrow \infty \rightarrow \ell \frac{\lambda}{\lambda + \mu} + (\ell+1) \frac{\lambda}{\mu}$

$\mu=0 \rightarrow \infty$
 $\lambda=0 \rightarrow 0$

Probability ancestral/descendant alignment

$$P(\hat{x}, \hat{y}) = P(y, \hat{x} \hat{y} | x) \cdot P(x)$$

generative model

distribution of ancestral seqs.

arbitrary: $(1-p)P^l$

length
distributing

$$P(l') = (1 - q_t) q_t^{l'} \quad \text{''} \quad q_t = \frac{\xi_t + p(1 - \xi_t)(1 - \delta_t)}{1 - p\delta_t(1 - \xi_t)}$$

$$P(L) = (1 - r_t) r_t^L \quad \text{''} \quad r_t = \xi_t + p(1 - \xi_t)$$

expected
frequencies

$$f_i(t) \equiv \frac{\langle L \rangle - \langle l \rangle}{\langle L \rangle} = \frac{\xi_t}{\xi_t + p(1 - \xi_t)}$$

$$f_\mu(t) \equiv \frac{\langle L \rangle - \langle l' \rangle}{\langle L \rangle} = \frac{p\delta_t(1 - \xi_t)}{\xi_t + p(1 - \xi_t)}$$

(i) Model is nm-reversible : $f_i \neq f_\mu \wedge P(l) \neq P(l')$

(ii) Model is nm-stationary : $P(l')$ is time dependent

$$\left. \begin{array}{l} \lambda = 0 \quad f_i = 0 \\ \mu \neq 0 \quad f_\mu = 1 - e^{-\mu t} \end{array} \right\} \begin{array}{l} \lambda \neq 0 \quad f_i = 1 - e^{-\lambda t} \\ \mu = 0 \quad f_\mu = 0 \end{array}$$

arbitrary rates λ, μ require nm-reversible models for seqs (no indels)

Comparison to TKF91

Similar set-up:

- i) probabilistic modelling ins/del
- ii) a generative birth-death model
- iii) $\lambda \geq 0, \mu \geq 0$ rates ins/dels

Departure:

TKF91 is reversible

$$P(\text{ancestral seq } e) = \left(1 - \frac{\lambda}{\mu}\right) \left(\frac{\lambda}{\mu}\right)^l$$

$$\lambda < \mu$$

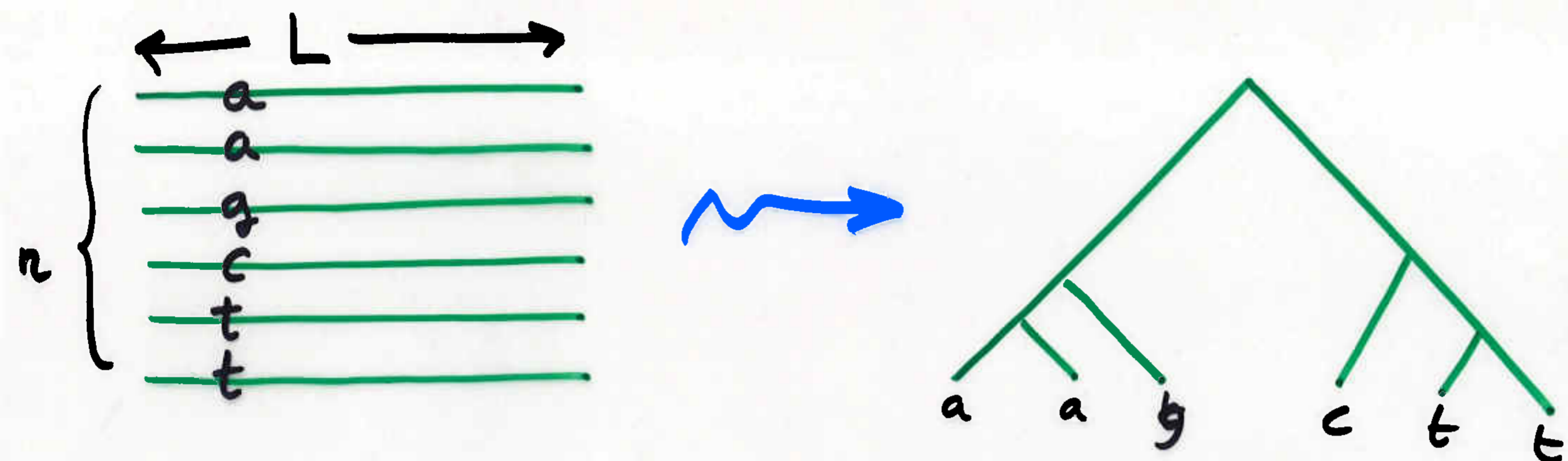
all seqs have same length on average

TKF91 not invariant by column rearrangements
adds realism ~ increases algorithmic complexity

Actual correspondence:

	TKF91		ERATE
$P(\underbrace{\# \bar{} \dots \bar{}}_{n}) = \frac{\lambda}{\mu} P_i^t \beta_t^{n-1}$	}	$P(1 - \xi_t)(1 - \delta_t) \xi_t^{n-1}$	$P(1 - \xi_t)(1 - \delta_t) \xi_t^{n-1}$
$P(\underbrace{\bar{} \# \dots \#}_{n}) = \begin{cases} \frac{\lambda}{\mu} q_i^t \beta_t^{n-1} & n > 0 \\ \beta_t & n = 0 \end{cases}$		$P(1 - \xi_t) \delta_t \xi_t^n$	$P(1 - \xi_t) \delta_t \xi_t^n$
$P(\underbrace{\bar{} \dots \bar{}}_{n}) = \beta_t^n$		$(1 - \xi_t) \xi_t^n$	$(1 - \xi_t) \xi_t^n$
$\text{extra term} = \left(1 - \frac{\lambda}{\mu}\right) (1 - \beta_t)$		$(1 - P)(1 - \xi_t)$	$(1 - P)(1 - \xi_t)$

Phylogenetic Inference



$$P(\text{alignment} \mid T, t's, R) \sim \underline{n \cdot L}$$

Felsenstein peeling algorithm

At the core of most phylogenetic inference packages:
PHYLIP, MOLPHY, PAUP*, PHYML, Palm, Pasm1, Mr Bayes, ...

all treat gaps as missing data

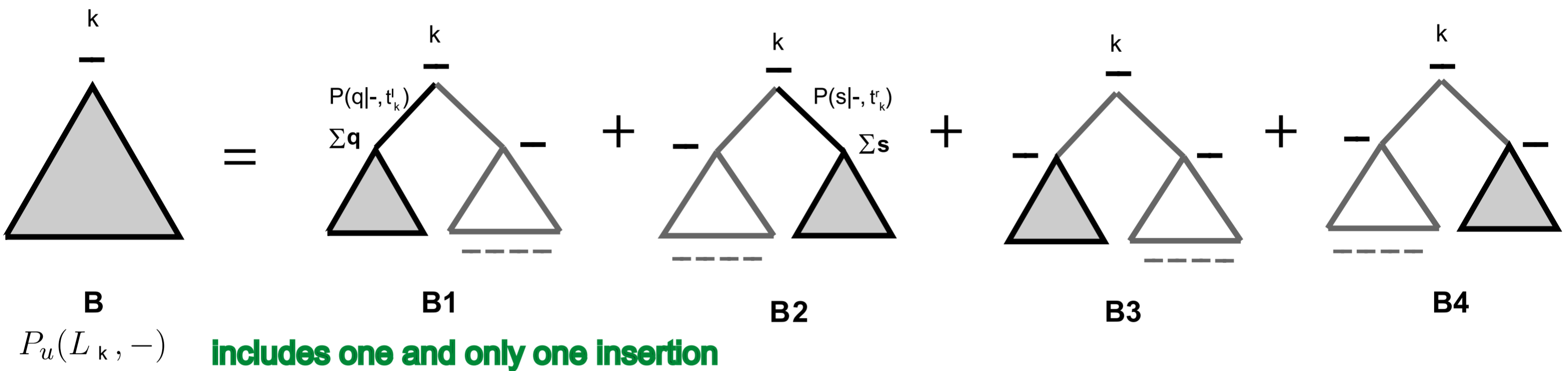
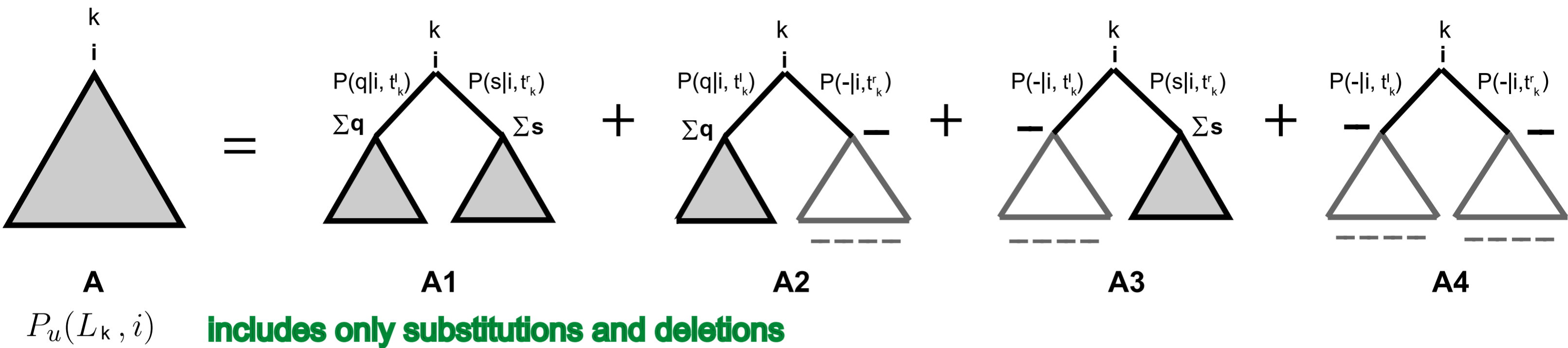
Practical application for ERATE

~ The substitutions-only Felsenstein algorithm is extendable to ERATE model with same time complexity \underline{nL}

~ Can calculate the probability of missing data (i.e. events not observed in given sequences) by marginalization

important for ancestral seq reconstruction

Felsenstein algorithm extended to gaps



Generative model

$$P(j | i, t) = (1 - \xi_t) P_t^\varepsilon(j | i)$$

$$P(- | i, t) = (1 - \xi_t) \gamma_t$$

$$P(j | -, t) = \xi_t \pi_j$$

Probability of a column u

$$P(u | T, R^\varepsilon, p) = P_u(L_{root}, -) + p \sum_{1 \leq i \leq K} P_u(L_{root}, i) \pi_i$$

(B) (A)

Probability of the alignment after marginalizing all unobserved (all gaps) columns

extra-term contribution

$$P(L | T, R^\varepsilon, p) = \frac{P(\star | T, R^\varepsilon, p)}{1 - P(u_{\text{gap}} | T, R^\varepsilon, p)} \prod_{1 \leq u \leq L} P(u | T, R^\varepsilon, p)$$

prob of a column with all gaps

DNAML ϵ

Extension to gaps of DNAML (PHYLIP)



a DNAML/DNAML ϵ comparison is an exclusive test of the effect of including ins/del

~ DNAML ϵ optimizes rates $\lambda\mu$

~ For unaligned alignments optimizes to $\lambda=\mu=0$
and produces same results as DNAML

DNAML ϵ reports practical times

used in Rfam 09 for seed alignments
containing up to 64 sequences ~ 92% of families

Tests on Simulated data

for each av. branch length (subs only) (0.005 - 2.0 subs/si)

use FB4 to evolve residues

Select method to evolve sequences (ROSE / ERATE)

for each "gap parameters" / length distribution (geometric / poisson / simpson)

for each n (1, ..., 100)

Generate random 8-Taxon tree T_n (Kuhn-Felsenstein)

for each L (50, ..., 1000 nts) $\Delta L = 5$ nts

Generate ancestral seq (L) (uniform dist)

Generate alignment \rightarrow phylogenetically correct
 \rightarrow realign

Infer tree \hat{T}_n \rightarrow DNAML
 \rightarrow DNAMLE

compare \hat{T}_n / T_n \rightarrow True positives
 \rightarrow Symm. Diff. distance (SDD)
 \rightarrow n Branch Scoring dist. (nBSD)

collect stats for all 100/L comparisons

performance

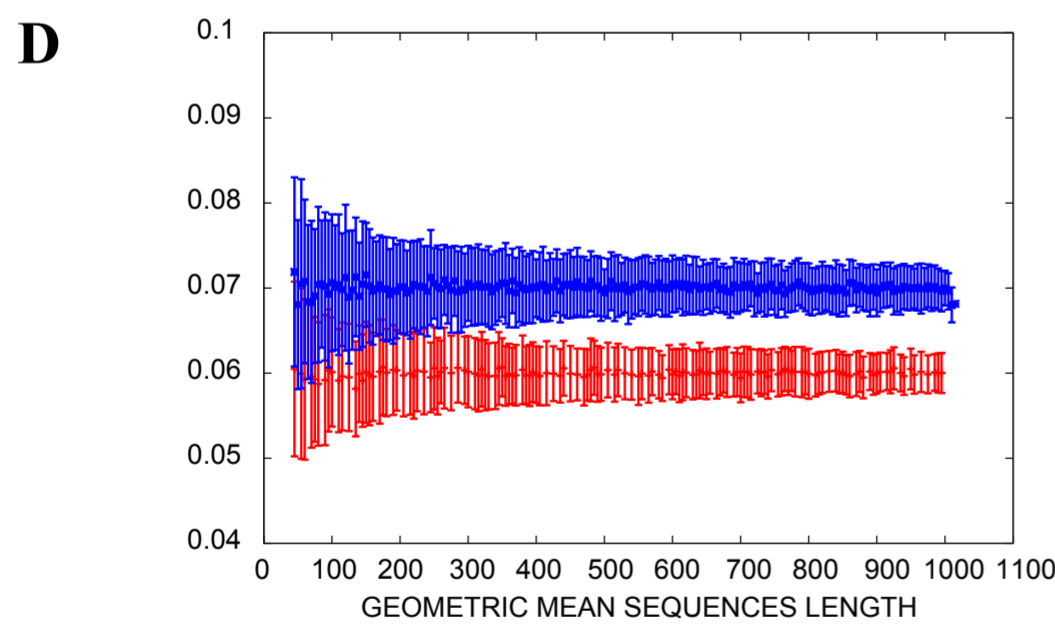
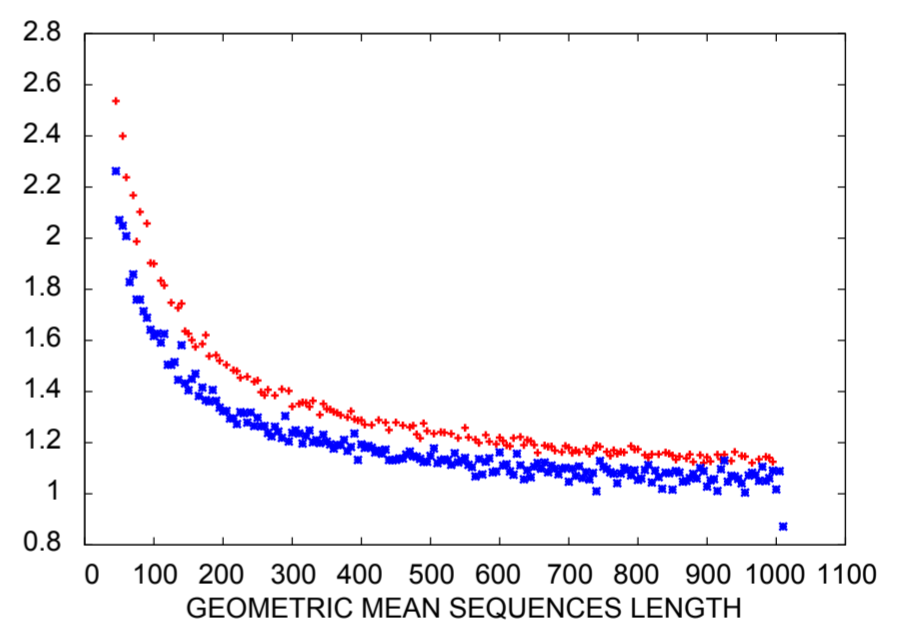
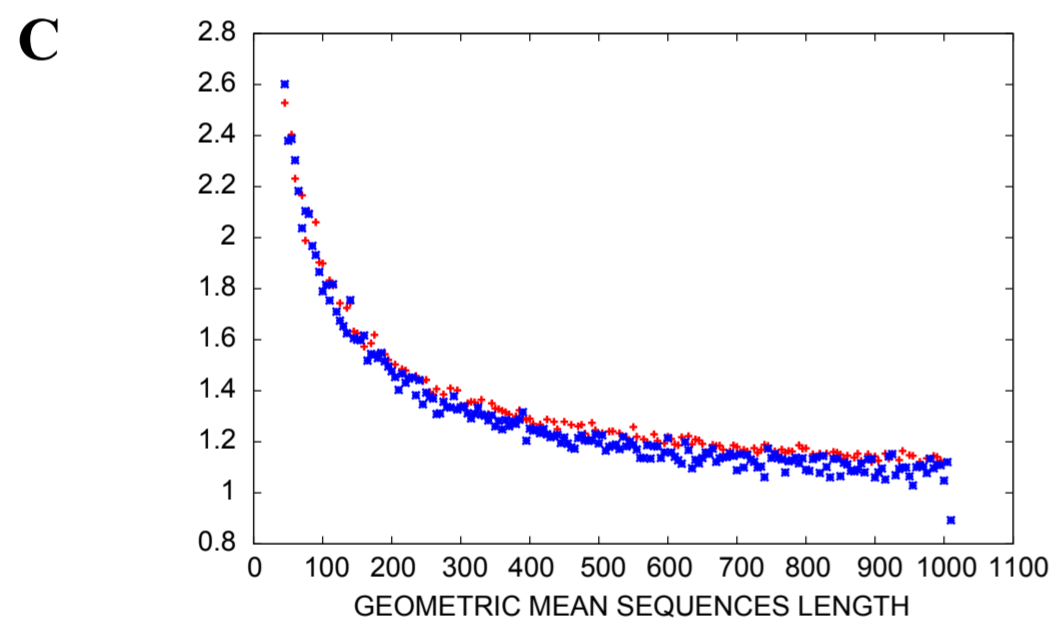
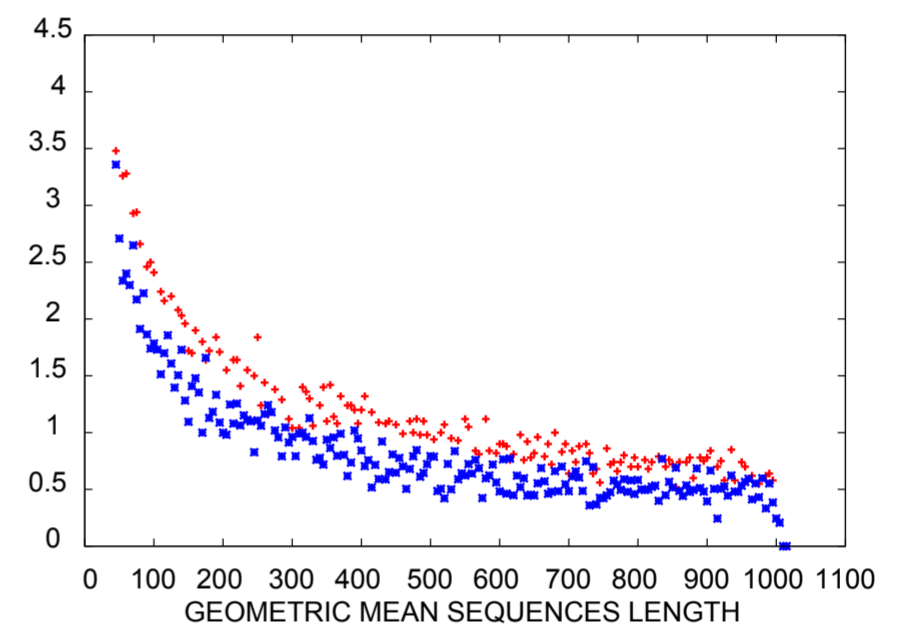
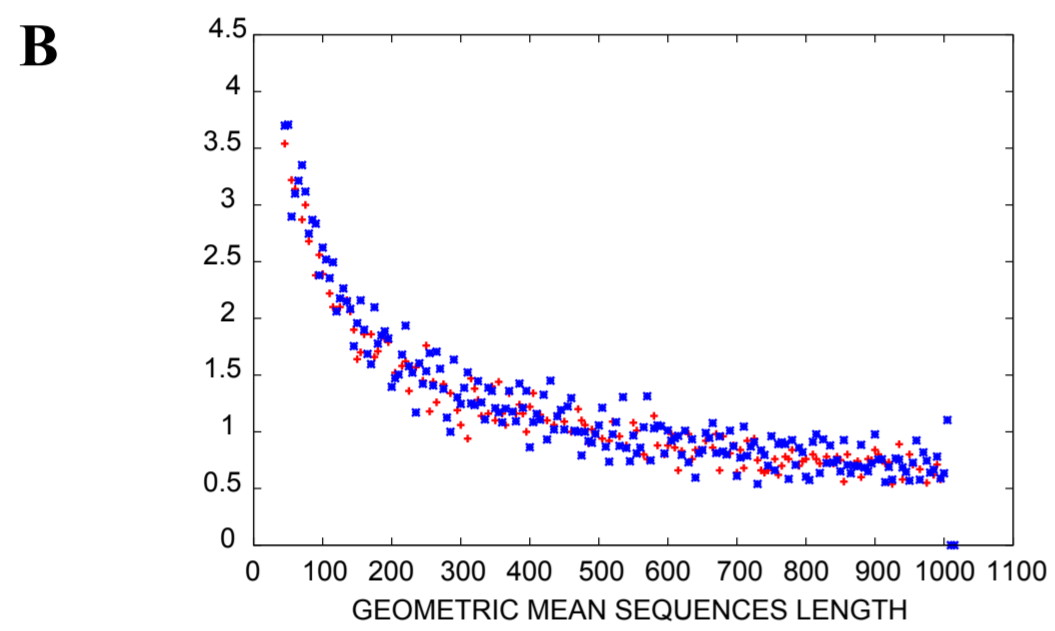
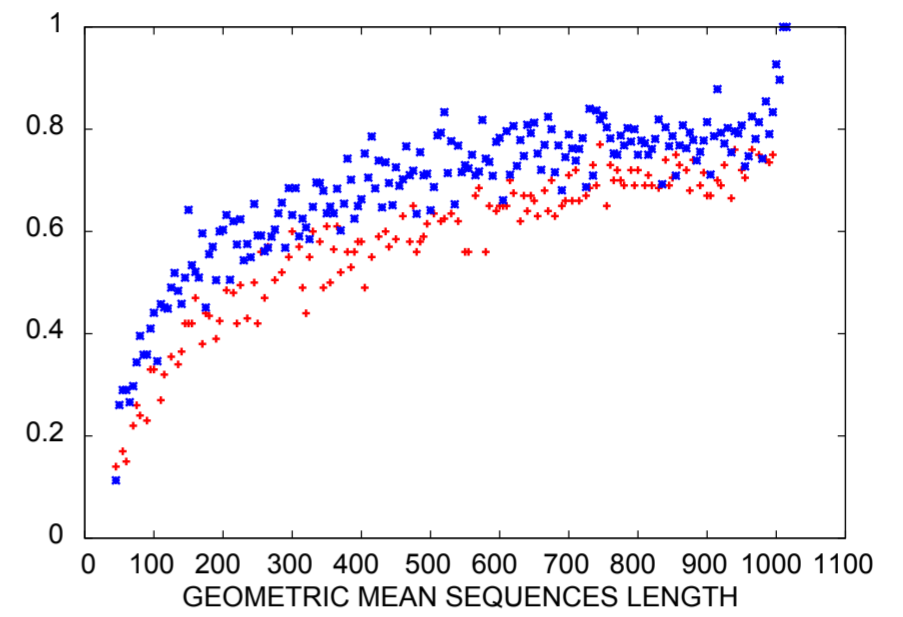
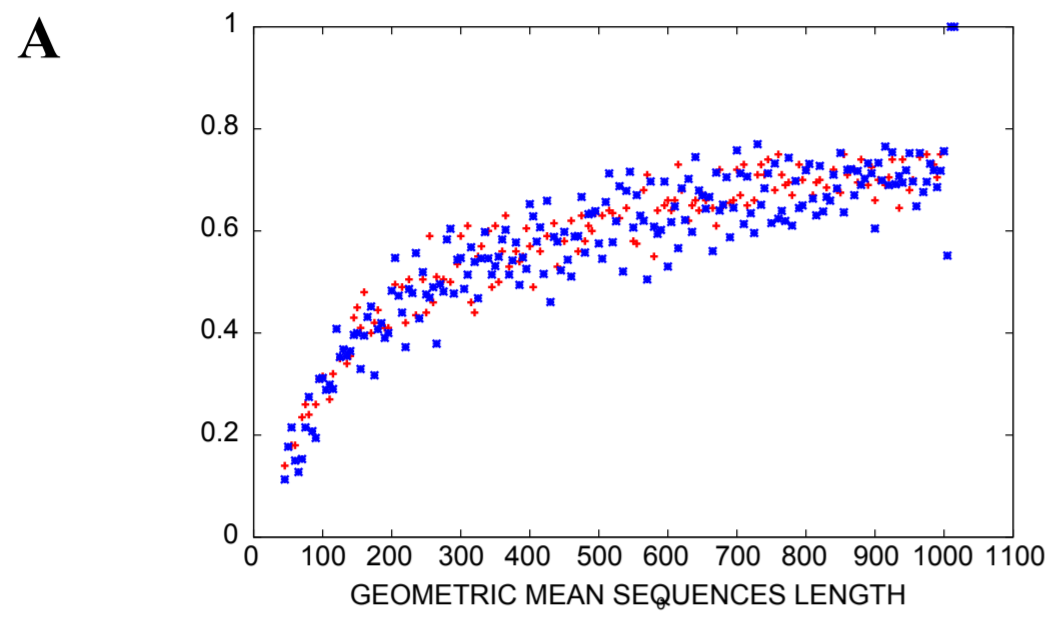


DNAML

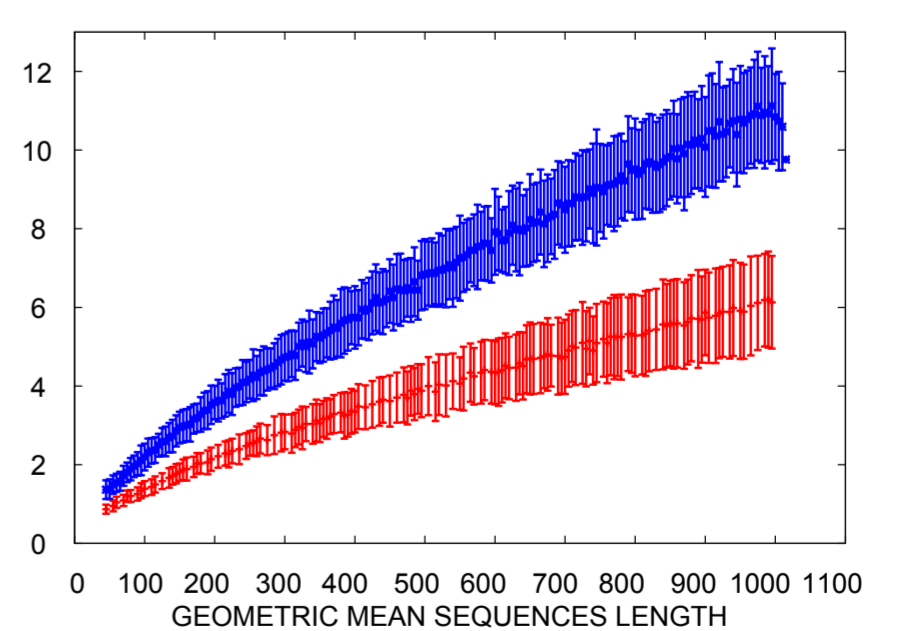
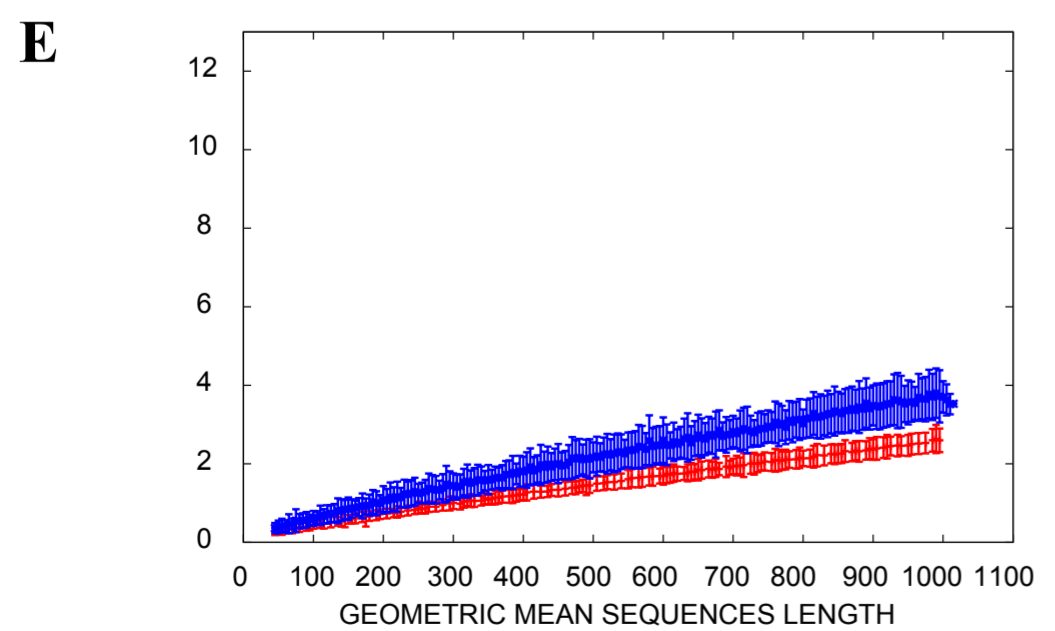
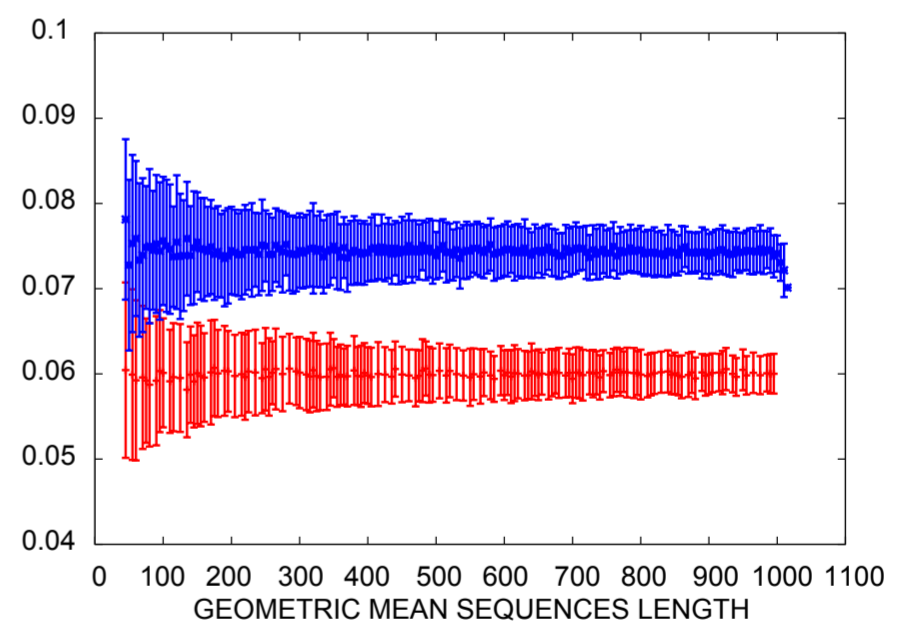
100 RANDOM 8 TAXON TREES

DNAML ϵ

+ RED ALIGNMENTS: NO GAPS %SUBS: 21 +/- 6 %ID: 79 +/- 6
* BLUE ALIGNMENTS: %GAPS: 13 +/- 4 %SUBS: 20 +/- 5 %ID: 66 +/- 8



branch length
contributed by gaps



TEST on rRNA alignments

Concordance test

Break alignment in two pieces

Do both produce the same inferred tree?

rRNA alignment statistics

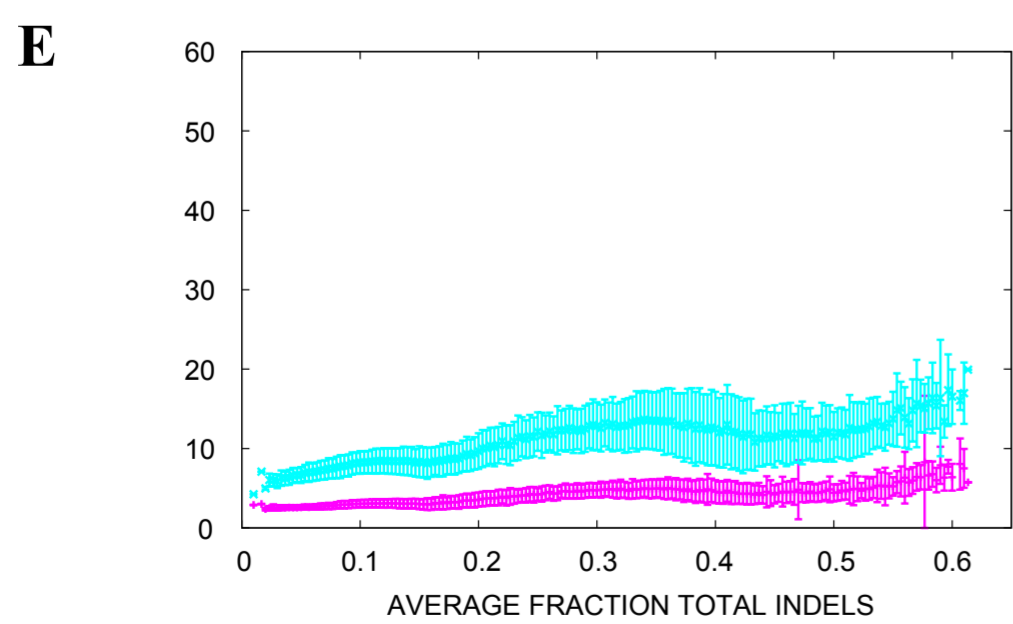
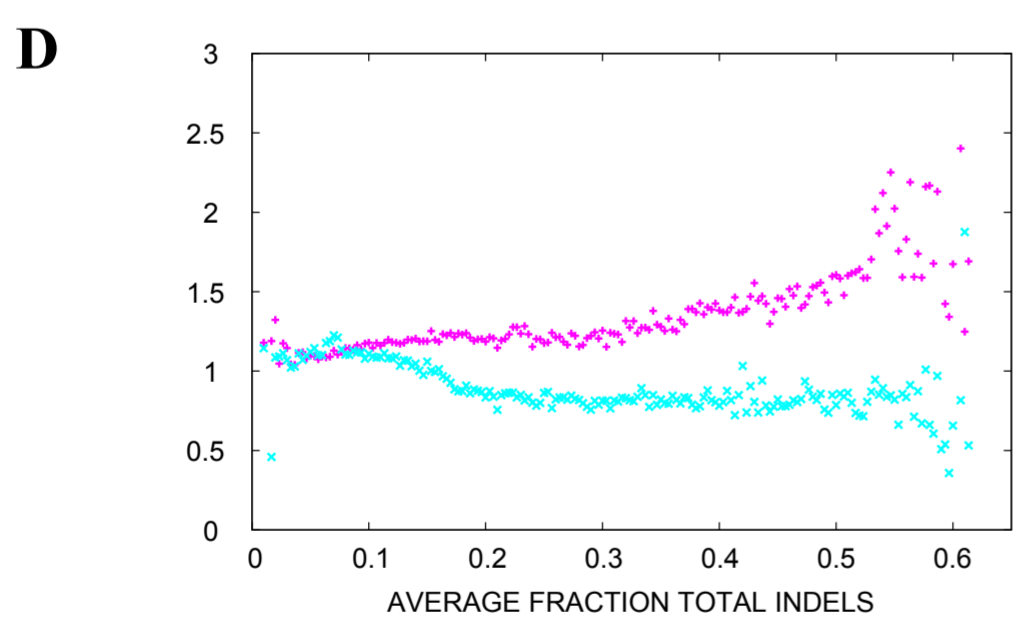
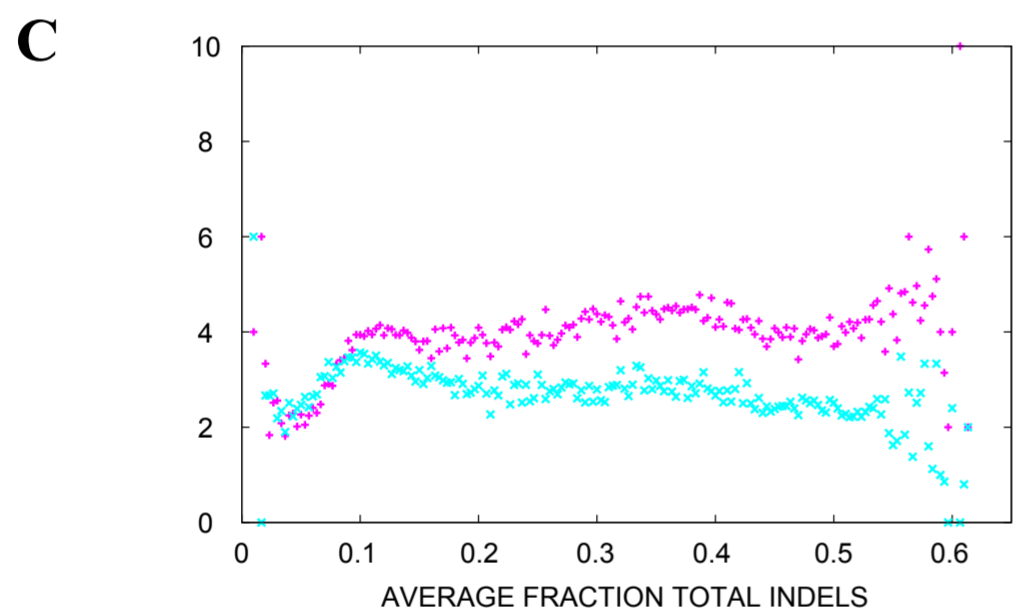
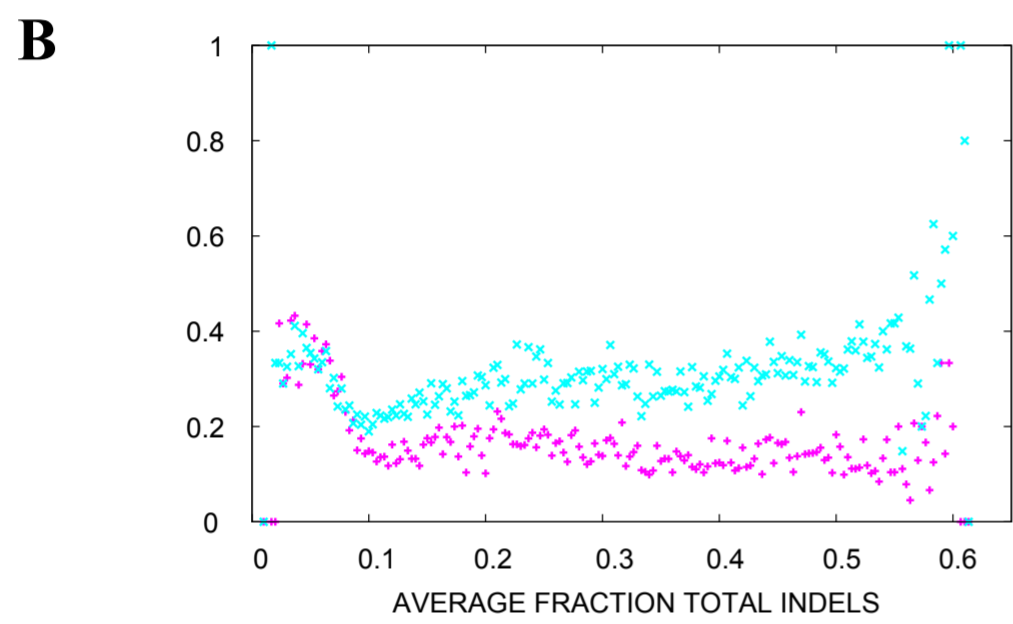
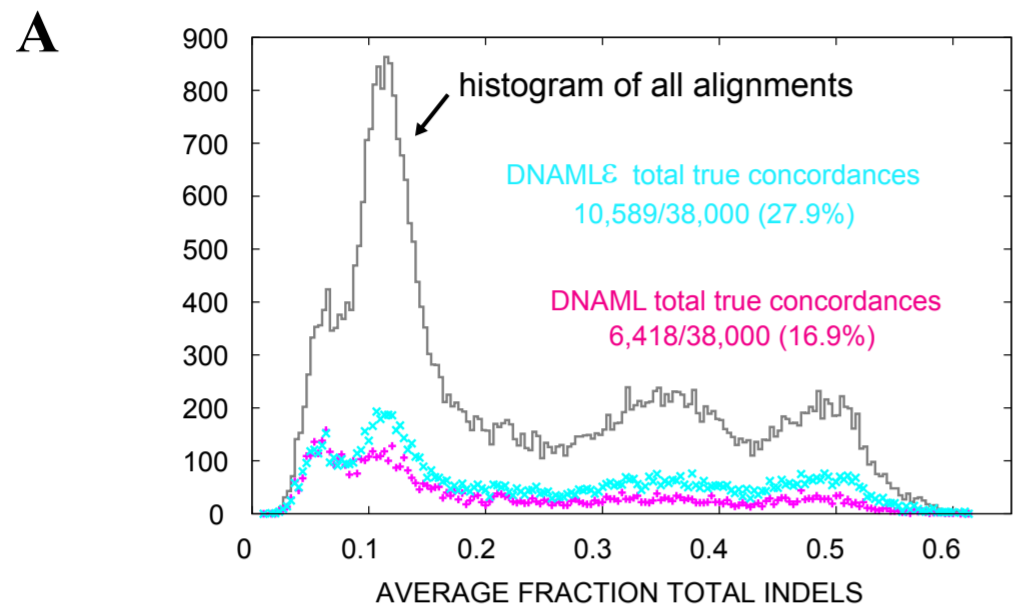
SSU	Archaea	Chloroplasts	Bacteria	Eukarya	Mitochondria
# seqs	74	75	1601	900	617
Alignment length	1756	2456	3094	7160	4716
Geometric mean seqs	1446	1446	1497	1788	1011
pairwise % ID	71 ± 7	74 ± 11	68 ± 6	56 ± 14	49 ± 20
pairwise % SUBS	23 ± 5	19 ± 7	22 ± 4	20 ± 5	25 ± 6
pairwise % GAP	6 ± 4	7 ± 5	11 ± 5	24 ± 12	26 ± 20
Total % Gaps	17.6	39.0	53.2	74.8	78.0

LSU	Archaea	Chloroplasts	Bacteria	Eukarya
# seqs	26	32	120	89
Alignment length	3346	4311	4555	9055
Geometric mean seqs	2994	2936	2918	3623
pairwise % ID	65 ± 9	70 ± 10	67 ± 6	50 ± 12
pairwise % SUBS	29 ± 7	20 ± 5	25 ± 4	22 ± 5
pairwise % GAP	6 ± 2	11 ± 5	9 ± 4	28 ± 11
Total % Gaps	10.5	31.8	35.9	59.6

Test on a large number of random 8 taxa subalignments

SSU rRNA

38,000 - EIGHT TAXON ALIGNMENTS
geometric mean seqs 1433 +/- 298



DNAML&E results *
DNAML results +

NUMBER TRUE CONCORDANCES

FRACTION TRUE CONCORDANCES

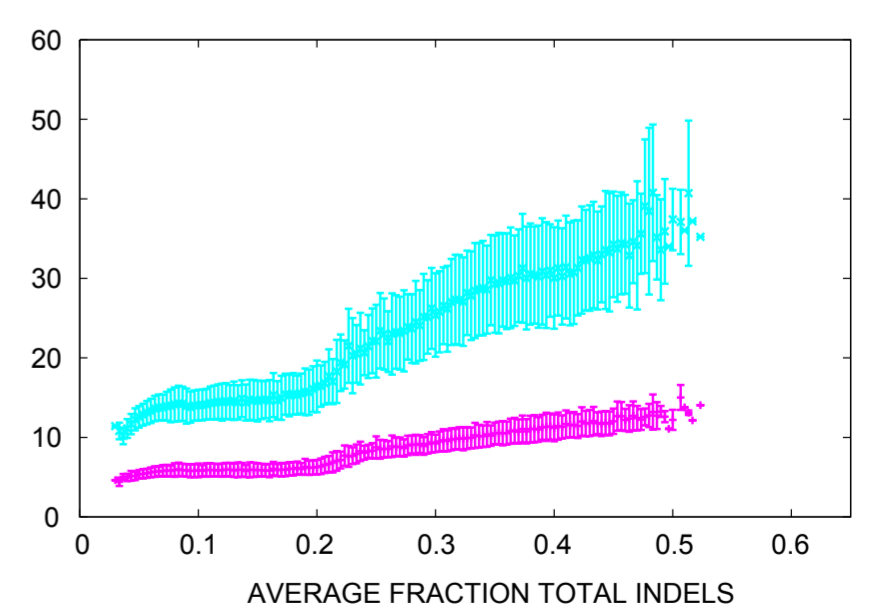
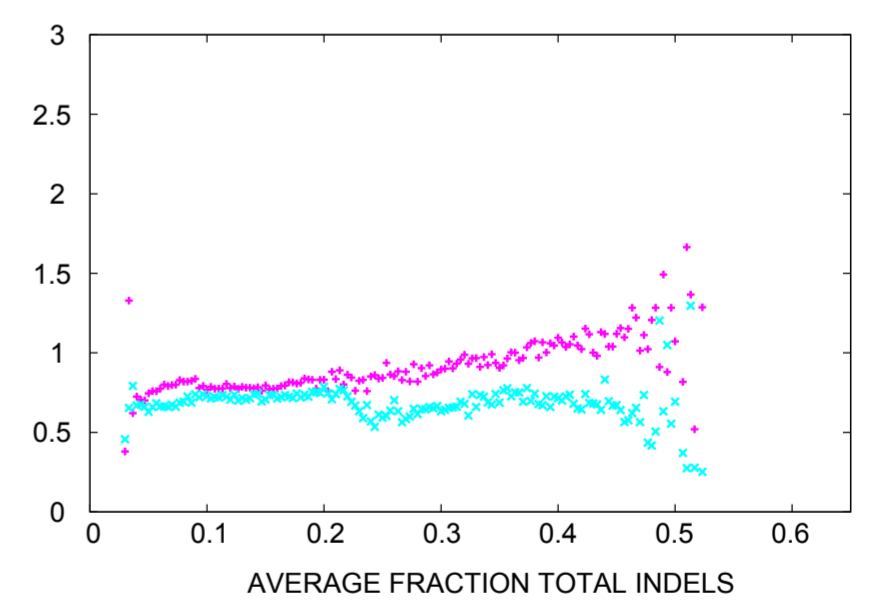
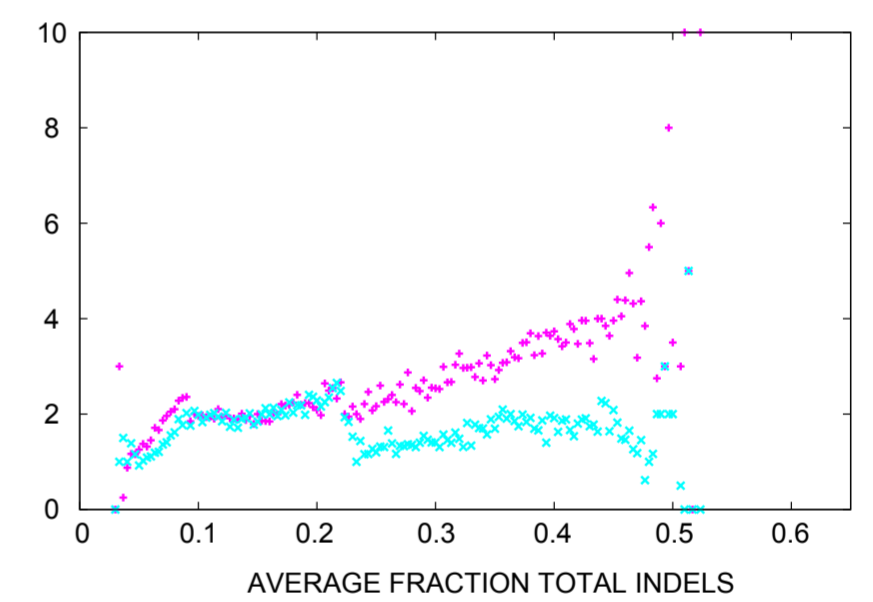
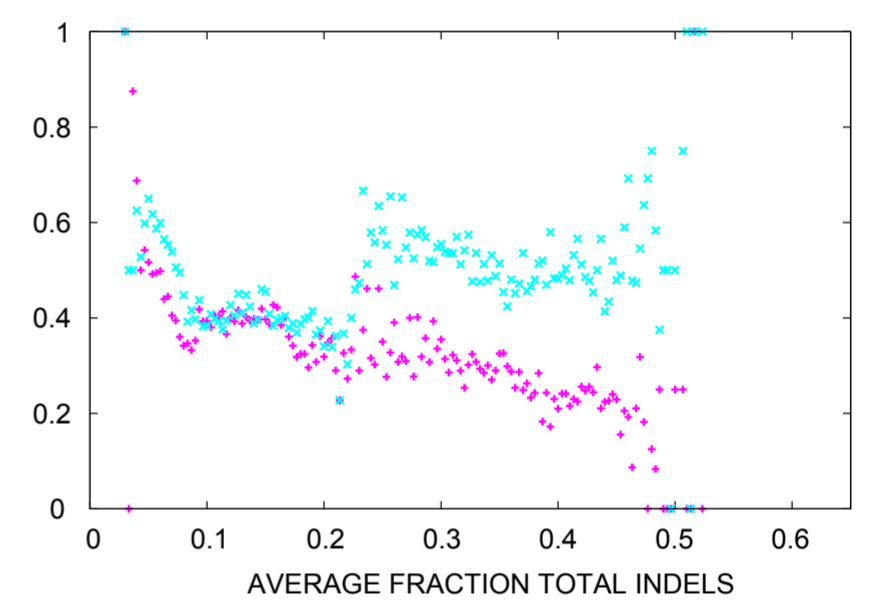
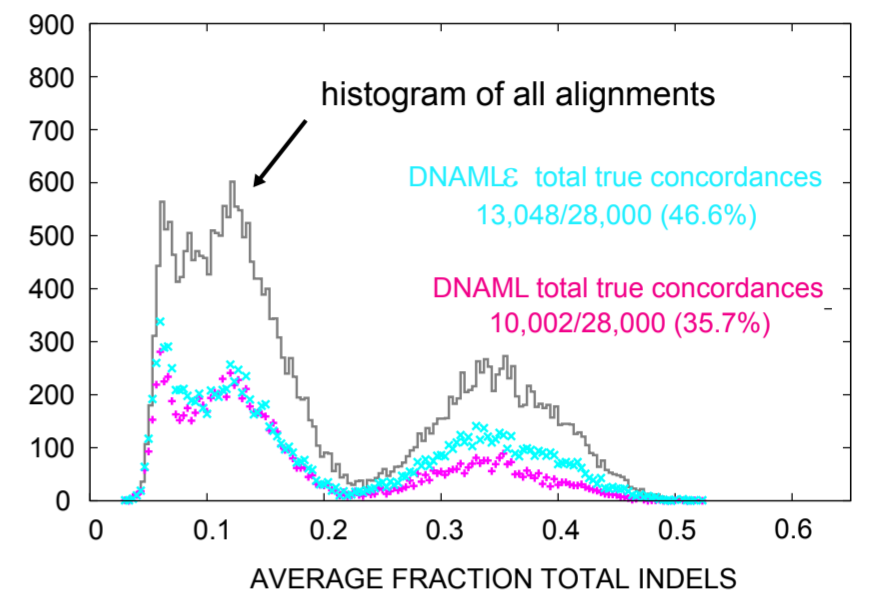
avg SDD DISTANCE

avg nBSD DISTANCE

RUN TIME (sec)
(avg +/- std)

LSU rRNA

28,000 - EIGHT TAXON ALIGNMENTS
geometric mean seqs 3220 +/- 344



Comparison of different gap distributions

8-Taxon alignments

gap length distribution	ave subs per site and branch	gap parameter	pairwise % ID	pairwise % SUBS	pairwise % GAPS	$\Delta_{\%}^{\text{AUC}}(\text{TP})$	$\Delta_{\%}^{\text{AUC}}(\text{SDD})$	$\Delta_{\%}^{\text{AUC}}(\text{nBSD})$
ϵ RATE	0.005	0.30	96 ± 2	2 ± 1	2 ± 1	28.9	-27.5	-24.9
POISSON $\lambda = 0.5$	0.005	0.0020	96 ± 2	2 ± 1	2 ± 1	24.6	-19.8	-16.8
CODING $c/t = 100$	0.005	0.0020	96 ± 2	2 ± 1	2 ± 2	4.8	-2.5	11.8
CODING $c/t = 6$	0.005	0.0010	96 ± 4	2 ± 1	2 ± 4	-4.5	5.0	36.8
CODING $c/t = 1$	0.005	0.0005	96 ± 4	2 ± 1	2 ± 4	-8.1	7.0	41.3
ϵ RATE	0.020	0.30	85 ± 4	7 ± 2	8 ± 3	22.3	-40.6	-18.0
POISSON $\lambda = 0.5$	0.020	0.0020	85 ± 5	7 ± 2	8 ± 3	20.5	-32.9	-12.3
CODING $c/t = 100$	0.020	0.0020	85 ± 5	7 ± 2	7 ± 4	2.9	-0.5	11.3
CODING $c/t = 6$	0.020	0.0010	84 ± 8	7 ± 2	9 ± 7	-11.6	19.9	43.3
CODING $c/t = 1$	0.020	0.0005	85 ± 8	7 ± 2	7 ± 8	-15.1	20.9	51.6
ϵ RATE	0.040	0.30	73 ± 7	12 ± 3	15 ± 4	23.4	-46.5	-13.9
POISSON $\lambda = 0.5$	0.040	0.0020	73 ± 7	12 ± 3	15 ± 4	20.7	-37.2	-10.6
CODING $c/t = 100$	0.040	0.0020	73 ± 8	13 ± 4	14 ± 5	0.7	0.2	8.3
CODING $c/t = 6$	0.040	0.0010	72 ± 10	12 ± 4	16 ± 10	-21.9	31.3	40.4
CODING $c/t = 1$	0.040	0.0005	73 ± 11	13 ± 4	14 ± 11	-26.9	36.4	50.1
ϵ RATE	0.070	0.30	58 ± 9	18 ± 4	24 ± 6	24.1	-47.3	-13.2
POISSON $\lambda = 0.5$	0.070	0.0020	58 ± 10	18 ± 4	24 ± 7	23.7	-41.4	-10.8
CODING $c/t = 100$	0.070	0.0020	60 ± 10	18 ± 4	22 ± 7	2.4	-2.4	5.5
CODING $c/t = 6$	0.070	0.0010	57 ± 12	17 ± 5	26 ± 12	-28.6	34.5	38.0
CODING $c/t = 1$	0.070	0.0005	59 ± 13	18 ± 5	23 ± 13	-39.4	50.1	48.9

$$\Delta_{\%}^{\text{AUC}}(f) = \frac{\text{AUC}\langle f(\text{DNAML}\epsilon) \rangle - \text{AUC}\langle f(\text{DNAML}) \rangle}{\max\{\text{AUC}\langle f(\text{DNAML}\epsilon) \rangle, \text{AUC}\langle f(\text{DNAML}) \rangle\}} \times 100.$$

For a given method M, the area under the curve (AUC): $\text{AUC}\langle f(\text{M}) \rangle = \sum_{L=50}^{L=1000} \langle f(L | \text{M}) \rangle \times \Delta L,$

where:

$$\langle f(L | M) \rangle = \frac{1}{N} \sum_{n=1}^N f(A_L^n | M), \text{ for alignments } \{A_L^n\}_{n=1}^{100}$$

with geometric mean of sequence length L, and $\Delta L = 5$ nts, for these experiments.

nomenclature:

AUC = area under the curve

TP = fraction of true positive trees.

SDD = Symmetric Difference Distance.

nBSD = normalized Branch Score Distance.

Higher is better

Lower is better

Lower is better

ROSE - Poisson $\lambda = 0.5$ gap distribution

8-Taxon alignments

ave subs per site and branch	ROSE gap parameter	pairwise % ID	pairwise % SUBS	pairwise % GAPS	$\Delta_{\%}^{\text{AUC}}(\text{TP})$	$\Delta_{\%}^{\text{AUC}}(\text{SDD})$	$\Delta_{\%}^{\text{AUC}}(\text{nBSD})$	mean MBL		time (L=1000) (secs)	
								DNAML	DNAML ϵ	DNAML	DNAML ϵ
0.005	0.0020	96 \pm 2	2 \pm 1	2 \pm 1	23.4	-18.5	-17.6	0.005	0.007	2.9 \pm 1.7	4.7 \pm 0.8
		96 \pm 2	2 \pm 1	2 \pm 1	3.8	-5.0	-16.3	0.006	0.008	2.5 \pm 1.1	4.6 \pm 0.6
0.010	0.0010	94 \pm 2	4 \pm 1	2 \pm 1	14.8	-16.1	-7.9	0.010	0.012	2.4 \pm 0.8	5.3 \pm 0.7
		94 \pm 2	4 \pm 2	2 \pm 1	-4.8	2.3	-6.7	0.010	0.013	2.4 \pm 0.9	5.2 \pm 0.8
0.010	0.0020	92 \pm 3	4 \pm 1	4 \pm 2	22.7	-26.1	-15.7	0.010	0.014	2.4 \pm 0.6	5.9 \pm 0.8
		91 \pm 3	5 \pm 2	4 \pm 2	-4.9	4.3	-14.2	0.012	0.018	2.4 \pm 0.4	5.8 \pm 0.8
0.020	0.0010	88 \pm 4	8 \pm 2	4 \pm 2	12.9	-19.1	-6.7	0.020	0.024	2.4 \pm 0.2	6.0 \pm 0.7
		87 \pm 4	9 \pm 3	4 \pm 2	-8.7	12.3	-5.3	0.023	0.027	2.4 \pm 0.2	5.9 \pm 0.8
0.020	0.0020	85 \pm 5	7 \pm 2	8 \pm 3	20.5	-32.9	-12.3	0.020	0.028	2.8 \pm 0.6	7.0 \pm 1.0
		83 \pm 6	10 \pm 4	7 \pm 2	-5.2	7.1	-11.3	0.030	0.039	2.9 \pm 0.6	7.1 \pm 0.8
0.030	0.0005	86 \pm 4	11 \pm 3	3 \pm 1	7.4	-14.7	-2.7	0.030	0.032	2.5 \pm 0.7	6.0 \pm 0.7
		85 \pm 4	12 \pm 3	3 \pm 1	-12.8	15.9	-1.4	0.032	0.035	2.7 \pm 0.6	6.1 \pm 0.7
0.030	0.0010	83 \pm 5	11 \pm 3	6 \pm 2	14.3	-22.4	-6.4	0.030	0.035	2.7 \pm 0.5	7.0 \pm 1.0
		81 \pm 5	13 \pm 4	5 \pm 2	-10.5	13.3	-4.8	0.036	0.042	2.8 \pm 0.5	6.9 \pm 0.9
0.030	0.0020	78 \pm 6	10 \pm 3	11 \pm 4	18.4	-35.6	-12.1	0.030	0.042	2.9 \pm 0.6	8.2 \pm 1.2
		75 \pm 8	16 \pm 5	9 \pm 3	-5.5	5.8	-9.2	0.050	0.062	3.1 \pm 0.4	8.6 \pm 1.0
0.040	0.0005	82 \pm 5	14 \pm 4	4 \pm 2	9.6	-17.2	-2.7	0.040	0.042	2.7 \pm 0.4	6.5 \pm 0.7
		81 \pm 5	15 \pm 4	4 \pm 2	-16.4	19.3	-0.7	0.043	0.047	2.7 \pm 0.3	6.7 \pm 0.6
0.040	0.0010	78 \pm 6	14 \pm 4	8 \pm 3	15.3	-26.1	-6.2	0.040	0.045	3.0 \pm 0.5	8.0 \pm 0.9
		76 \pm 7	17 \pm 5	7 \pm 2	-10.1	12.3	-3.7	0.051	0.058	3.2 \pm 0.8	7.9 \pm 0.7
0.040	0.0015	75 \pm 7	13 \pm 4	12 \pm 4	16.9	-33.1	-8.4	0.040	0.048	3.1 \pm 0.4	8.9 \pm 1.2
		72 \pm 8	19 \pm 6	9 \pm 3	-7.1	7.8	-6.2	0.061	0.071	3.4 \pm 0.6	9.1 \pm 1.2
0.040	0.0020	73 \pm 7	12 \pm 3	15 \pm 4	20.7	-37.2	-10.6	0.040	0.052	3.4 \pm 0.5	9.4 \pm 1.3
		68 \pm 9	21 \pm 7	11 \pm 3	-4.7	4.6	-8.4	0.072	0.085	3.7 \pm 0.4	10.8 \pm 1.3
0.070	0.0005	71 \pm 7	22 \pm 5	7 \pm 3	12.6	-21.3	-3.0	0.070	0.071	3.3 \pm 0.4	9.6 \pm 1.0
		70 \pm 8	24 \pm 6	6 \pm 2	-11.6	15.0	-0.3	0.080	0.085	3.4 \pm 0.5	9.3 \pm 0.9
0.070	0.0010	66 \pm 8	20 \pm 5	13 \pm 4	16.0	-31.2	-6.3	0.070	0.074	3.7 \pm 0.4	10.8 \pm 1.1
		62 \pm 10	28 \pm 8	10 \pm 3	-9.2	9.3	-3.2	0.101	0.109	3.9 \pm 0.3	11.6 \pm 1.2
0.070	0.0015	62 \pm 9	19 \pm 5	19 \pm 5	20.2	-39.6	-8.2	0.070	0.078	4.0 \pm 0.6	12.5 \pm 1.7
		57 \pm 11	31 \pm 8	12 \pm 3	-7.1	5.0	-9.1	0.126	0.135	4.7 \pm 0.6	14.6 \pm 1.8
0.070	0.0020	58 \pm 10	18 \pm 4	24 \pm 7	23.7	-41.5	-10.8	0.070	0.083	4.3 \pm 0.7	13.5 \pm 1.7
		52 \pm 12	35 \pm 9	13 \pm 4	-4.4	3.1	-11.9	0.154	0.162	5.1 \pm 0.7	16.8 \pm 2.4
0.100	0.0005	63 \pm 9	28 \pm 6	10 \pm 3	17.1	-29.7	-3.6	0.100	0.100	4.0 \pm 0.7	12.1 \pm 1.3
		60 \pm 10	32 \pm 8	8 \pm 3	-11.7	12.7	-0.4	0.121	0.125	4.1 \pm 0.5	12.0 \pm 1.5
0.100	0.0010	57 \pm 9	25 \pm 6	18 \pm 6	20.6	-40.3	-6.3	0.100	0.100	4.5 \pm 0.7	14.5 \pm 1.3
		52 \pm 12	36 \pm 9	11 \pm 3	-7.9	6.6	-6.3	0.160	0.165	4.7 \pm 0.5	15.4 \pm 2.3
0.100	0.0015	52 \pm 10	23 \pm 5	25 \pm 6	25.3	-44.1	-8.8	0.100	0.105	4.9 \pm 0.6	16.6 \pm 2.2
		47 \pm 12	40 \pm 10	13 \pm 4	-4.0	2.0	-12.2	0.201	0.202	5.7 \pm 1.0	19.0 \pm 3.3
0.100	0.0020	48 \pm 11	21 \pm 4	31 \pm 7	27.2	-50.7	-11.4	0.100	0.110	5.4 \pm 0.7	18.3 \pm 2.2
		43 \pm 12	43 \pm 10	14 \pm 4	0.3	-0.5	-16.9	0.242	0.234	6.0 \pm 1.0	20.7 \pm 3.1

$$\Delta_{\%}^{\text{AUC}}(f) = \frac{\text{AUC}\langle f(\text{DNAML}\epsilon) \rangle - \text{AUC}\langle f(\text{DNAML}) \rangle}{\max\{\text{AUC}\langle f(\text{DNAML}\epsilon) \rangle, \text{AUC}\langle f(\text{DNAML}) \rangle\}} \times 100.$$

For a given method M, the area under the curve (AUC): $\text{AUC}\langle f(\text{M}) \rangle = \sum_{L=50}^{L=1000} \langle f(L|M) \rangle \times \Delta L$,

where:

$$\langle f(L|M) \rangle = \frac{1}{N} \sum_{n=1}^N f(A_L^n | M), \text{ for alignments } \{A_L^n\}_{n=1}^{100}$$

with geometric mean of sequence length L, and $\Delta L = 5$ nts, for these experiments.

nomenclature:

TP = fraction of true positive trees.

SDD = Symmetric Difference Distance.

nBSD = normalized Branch Score Distance.

MBL = mean branch length.