Accounting for Sequence Uncertainty in Maximum Likelihood Phylogenetic Inference

Alexey Kozlov¹, Alexandros Stamatakis^{1,2}

¹The Exelixis Lab, Scientific Computing Group, Heidelberg Institute for Theoretical Studies, Heidelberg, Germany ²Institute for Theoretical Informatics, Karlsruhe Institute of Technology, Karlsruhe, Germany

Background

Maximum Likelihood (ML) phylogenetic inference

Multiple sequence alignment (MSA)

AGCT-GCAGTACCC

AGCTTGCAGTACCG

AGCATGCA---CCG AGCATGCA---CCT

Phylogenetic tree

$LH(MSA \mid tree) \Rightarrow max$

Problem:

Discrete characters cannot adequately represent MSA uncertainty (= errors in sequencing, assembly, alignment)

Traditional approaches to deal with uncertainty

- IUPAC ambiguity code (e.g., S = C or G)
- filtering (remove low-quality bases/columns)
- sequencing with high coverage (overlap multiple reads and call a consensus base)

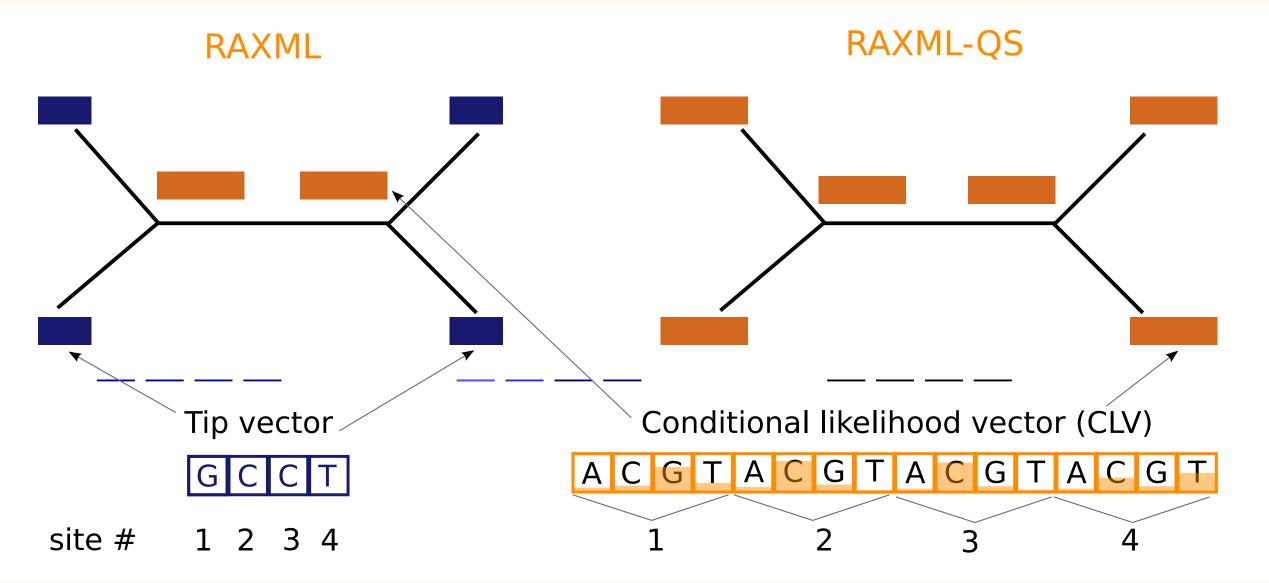
Idea

Explicitly use quantified sequence uncertainty in phylogenetic likelihood calculation

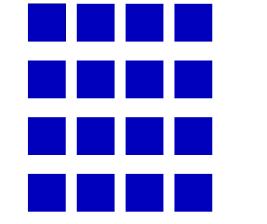
Implementation

Experimental RAxML-QS code (very basic tree search functionality of RAxML[1], DNA data only)

Changes to the tree data structure

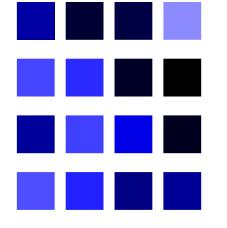


Available error models

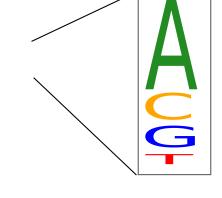


estimated

uniform error ε user-defined or



site-specific error ε_{ii} **FASTQ**



per-site base probabilities CATG (pyRAD v3.0), FAST5 (Oxford Nanopore)

RADseq data (naïve approach)

Computing base probabilities

FASTQ quality scores

AGCTGCAGTACCC sequence **75:>E3/?8;BB**@ phred+33

 $Q = -10\log_{10}\varepsilon \rightarrow \varepsilon = 10^{-0.1}Q$

 $P(A | S_i = "A") = 1-\epsilon$ $P(C \mid S_i = "A") = \varepsilon/3$ $P(G \mid S_i = "A") = \varepsilon/3$

 $P(T \mid S_i = "A") = \varepsilon/3$

AGCTGCAGTACCC

AGCTGCAGTACCC AGCTGCAGTACCC AGCTGCACTACCC **AGCTGCACTACCC** AGCTGCACTACCC **AGCTGCATTACCC**

"consensus MSA"

G: 3 A: 0 C: 3 **T:** 1 A: 0 C: 3/7 G: 3/7 T: 1/7 "counts-based MSA"

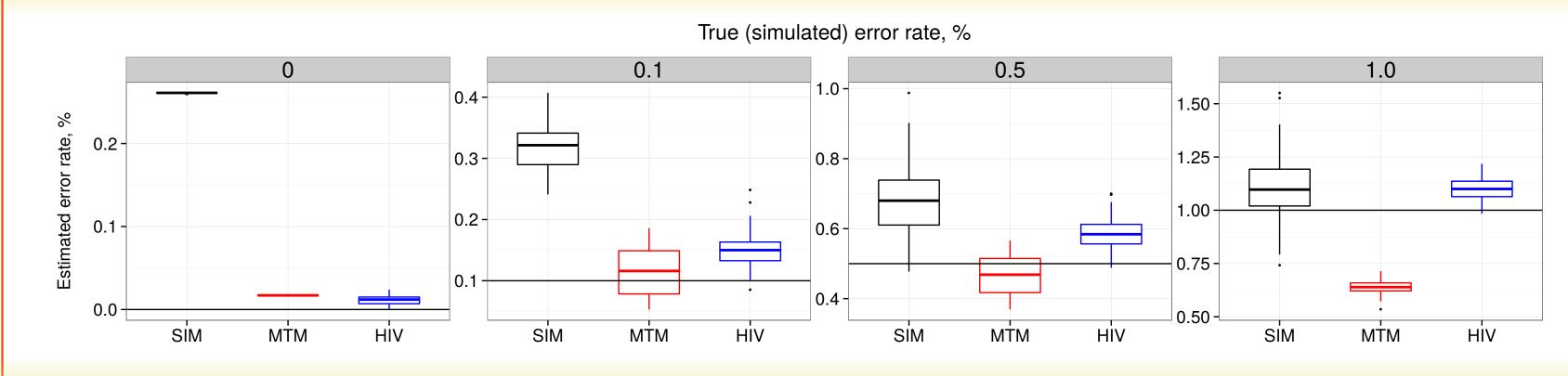
extract base counts

Results

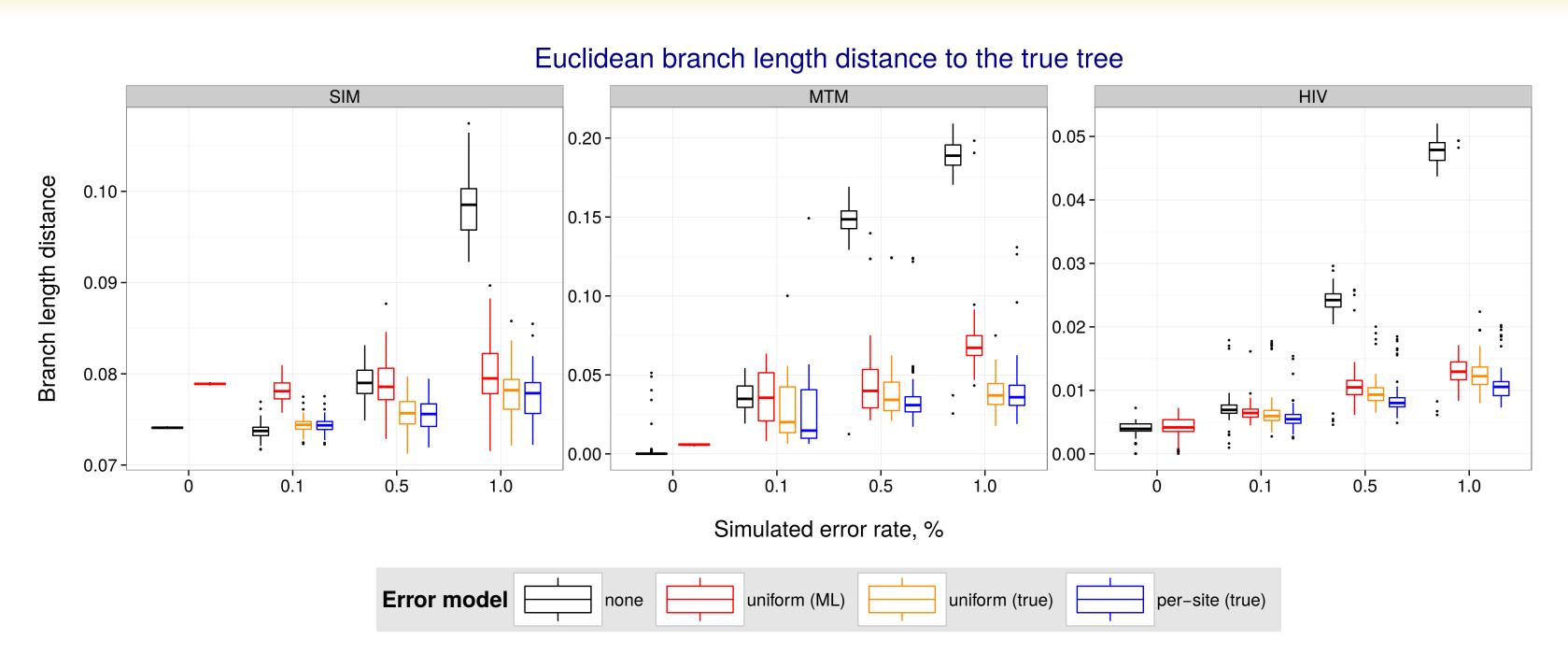
Test datasets

Designator	Organism group, gene	# taxa	# sites	Reference
SIM	(simulation)	50	2,000	
HIV	HIV-I virus, pol	23	2,841	[5]
MTM	Mammals, mtDNA	23	9,741	[6]
RADSEQ	Lizards, genome-wide	74	3,142 - 1,250,860	A. Leache (unpubl.)

ML estimation of uniform error rate

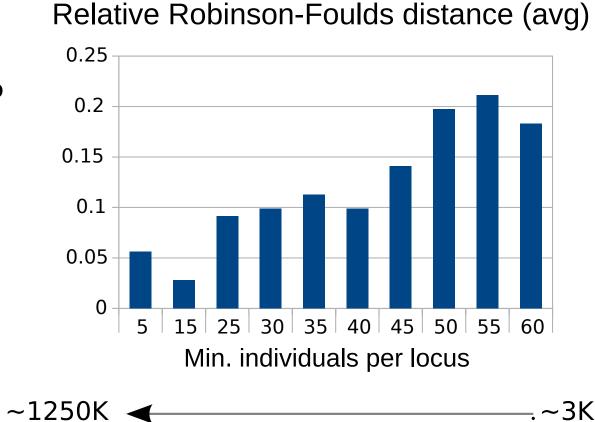


Tree inference on alignments with simulated error

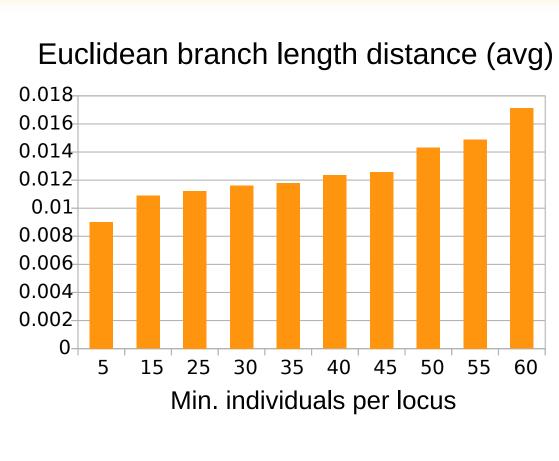


Tree inference on empirical RADseq data: counts-based MSA vs. consensus MSA

- Relative RF distance 3-21%
- Difference is much more pronounced on shorter alignments (= high coverage threshold)



alignment sites



Conclusions & Outlook

- ML estimation of alignment error rate is possible, but its accuracy varies significanlty across datasets and experimental settings
- Accounting for sequence uncertainty improves branch length estimates on simulated data with moderately high error rates (> 0.5 %)
- On empirical RADseq data, using MSA built from raw base counts instead of consensus sequences yields consistently different trees
- Next steps:
- Use quality scores provided by aligners (e.g., FSA[2])
- Evaluate the proposed approach in the context of phylogenetic placement (RAXML-EPA)

References

- [1] Stamatakis (2014) "RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies" *Bioinformatics* 30(9)
- [2] Bradley et al. (2009) Fast Statistical Alignment. PLoS Comput Biol 5(5)
- [3] Kuhner and McGill (2014) "Correcting for sequencing error in maximum likelihood phylogeny inference" G3 (Bethesda) 4(12)
- [4] Eaton (2014) "PyRAD: assembly of de novo RADseq loci for phylogenetic analyses" *Bioinformatics* 30(13)
- a sister group relationship between Xenartha (Edentata) and Ferungulates." Mol. Biol. Evol. 14:762-768

[5] Arnason et al. (1997) "Phylogenetic analyses of mitochondrial DNA suggest

[6] Yang et al. (2000) "Codon-substitution models for heterogeneous selection pressure at amino acid sites" *Genetics* 155:431-449

Acknowledgements

- Code availability: https://github.com/amkozlov/raxml-qs
- This work was funded by HITS gGmbH.
- We would like to thank Nick Goldman, Adam Leache and Deren Eaton for their invaluable assistance.







