

phytools `rerootingMethod` function was selected over alternatives, such as APE `ace`, because it can handle extant tip nodes with unknown states in the phylogenetic tree and also compute posterior probability for those nodes, hence, making applicable for hidden state prediction. The maximum marginal likelihood (also called empirical Bayesian posterior probability) calculated for unknown tip nodes reflects the most likely state for the node given the empirically estimated evolution model and phylogeny. It captures the uncertainty associated with phylogenetic branch lengths and rates of change associated with the states. In the context of Phylotyper, the posterior probabilities provides a confidence value associated with a predicted subtype. The evolutionary model used in `rerootingMethod` is the M_k model for discrete states (a continuous-time Markov process model). It is estimated from the data, so in Phylotyper, the reference set of genes used to build a phylogenetic tree and run the ASR, is key to the performance of the predictor. For subtype schemes packaged in Phylotyper, the annotated genes will be maintained and routinely updated to provide a robust and comprehensive reference set. The other key assumption is that gene phylogeny is correlated with the evolution of subtype state. Because of these assumptions are central to the method, all new subtype schemes added into Phylotyper are evaluated for their predictive performance.

Phylotyper is developed in python and R. Outlined below are the steps and tools used in the Phylotyper pipeline:

1. Identify subtype gene loci in input genomes using `blastn` or `blastx`.
2. Align input genes with unknown subtype against a pre-aligned set of reference genes using the tool MAFFT's `-add` feature.
3. If multiple loci are involved, concatenate individual alignments into superalignment.
4. Generate phylogenetic tree of aligned genes with `FastTree`.
5. Run `phytools rerootingMethod` using the phylogenetic tree and assigned subtypes. Genes with unknown subtype are assigned a flat prior.
6. Identify the subtype with maximum marginal likelihood for the unknown genes and report to user. Users are also provided with a image of the phylogenetic tree that shows the position of the unknown genes. An example is shown in Figure 1.

Phylotyper is available for download from: <https://github.com/superphy/insilico-subtyping>.

3 Functionality in the Phylotyper Tool

Phylotyper was designed to be incorporated into a whole genome sequencing public health workflow. The main input into Phylotyper is assembled (but not necessarily closed) genome sequences. Putative loci needed for the selected subtype scheme are identified in the input genomes using BLAST. The discovered loci are then sent to the Phylotyper subtype prediction module. It is possible to bypass the loci search step and provide genes directly to Phylotyper using the subtype run mode.

The Phylotyper approach can be used to predict any biological property that can be inferred from the phylogenetic distribution of a nominal set of genomic loci. Currently, subtype schemes for *Escherichia coli* Shiga-toxin 1 and 2 subtypes, Serotype O and H-types and Intimin subtypes are available in the Phylotyper package. Details about the current schemes available in the Phylotyper package is provided in table XXX. However, included in the Phylotyper software, is the capability to add new subtype schemes. Creating a new subtype scheme will save the required reference files in a data directory, allowing the newly added schemes to be easily run from Phylotyper. We also encourage users to contribute their subtype schemes to the main software repository (please send us subtype schemes you have developed and wish to share).

Checks are built-in new subtype scheme pipeline to ensure the main assumptions of the Phylotyper approach are not violated. A scan of the reference phylogenetic tree is conducted to search for tightly clustered subclades in the tree that have divergent subtypes. We first estimate the distribution of the inter-patristic distances of genes with the same subtype. The distribution is used to dynamically select a distance threshold, which we employ to flag subclades wherein the inter-patristic distance is less than the threshold but contains multiple distinct subtypes. We also conduct a leave-one-out cross validation for each gene in the new subtype scheme to assess the scheme's predictive capability. Users are alerted if their subtype schemes that have a F1-score below 0.9 from the cross-validation analysis.

Phylotyper also offers flexibility in the parameterization of evolution model used in the ASR step. A component of the underlying ASR framework, is an M_k or markov model of subtype evolution, for which an empirical transition rate matrix is estimated from the data. The transition matrix is used to calculate the expected number of subtype state changes given a distance in the phylogenetic tree. Different model parameterizations can be defined for the transition rate matrix. The simplest parameterization available in Phylotyper is the equal rates model; all subtypes have the same forward and reverse rate. The most complex parameterization available in Phylotyper is the symmetric model, wherein each forward and reverse rate for a given pair of subtypes are assigned a separate parameter. Frequently, the number of subtypes makes the symmetric model too computationally prohibitive. To offer more flexible models in these situations with reduced numbers of free parameters, two custom parameterization approaches were developed. The custom approaches both use a binning strategy that attempts to identify sets of subtypes that would have similar rates and assign them a single parameter as a set. These approaches are described in detail in the Supplementary Information XXX. Each of these model parameterizations; equal, symmetric and the two custom models are tested and evaluated in new subtype pipeline. The parameterization that has highest accuracy (based on a LOOCV analysis) is selected. In the case of ties, the model with the fewest parameters is given precedence (the symmetric model is not tested when the number of subtypes is over 10).

In testing the predictive capability of Phylotyper, we found that in some cases, jointly using multiple loci increased the subtype prediction accuracy. And so, in Phylotyper, multiple loci can be used in creating subtype schemes. The loci will be independently identified in the genome using BLAST and then independently aligned. The individual loci alignments are then concatenated to form a single superalignment that is used to build the predictive phylogenetic tree. The Stx1 and Stx2 schemes are examples of a multi-loci subtype schemes. The A and B subunit genes that make up the holo-toxin in Shiga-toxin are the inputs in these subtype schemes.

4 Comparison to a Sequence Similarity Approach

To assess the performance of the phylogenetic-based method used in Phylotyper and compare it to a direct sequence similarity-based approach, we ran a leave-one-out cross validation analysis. Each reference gene with experimentally validated subtype was removed from the training dataset and used as a test input into the Phylotyper program. Similarly, using a mock sequence-similarity based tool we developed for assessment purposes, each gene was withheld from the BLAST reference database and used as input (Source code is available in TODO). The sequence-similarity tool assigns putative subtypes using a BLAST sequence search. The predicted subtype annotation for the test input is assigned from top BLAST hit, provided the top hit passes a percent identity, e-value and alignment coverage minimum thresholds. The threshold values are listed in TODO. This approach of transferring annotations from genes with highest sequence similarity is the core of most current subtyping methods.

5 Discussion

Phylotyper can... New subtypes, Performance test highly accurate.
Sequence does not develop a model of subtype evolution. Subtypes can evolve and different rates. Strict one-size fits all cutoffs will not work. The phylogenetic framework provides natural statistical interpretation of confidence values.

6 Conclusion

Text Text text

Acknowledgements

text text text

Funding

This work has been supported by the...

References

Bofelli,F., Name2, Name3 (2003) Article title, *Journal Name*, **199**, 133-154.
Bag,M., Name2, Name3 (2001) Article title, *Journal Name*, **99**, 33-54.
Yoo,M.S. *et al.* (2003) Oxidative stress regulated genes in nigral dopaminergic neuronol cell: correlation with the known pathology in Parkinson’s disease. *Brain Res. Mol. Brain Res.*, **110**(Suppl. 1), 76–84.
Lehmann,E.L. (1986) Chapter title. *Book Title*. Vol. 1, 2nd edn. Springer-Verlag, New York.
Crenshaw, B.,III, and Jones, W.B.,Jr (2003) The future of clinical cancer management: one tumor, one chip. *Bioinformatics*, doi:10.1093/bioinformatics/btn000.
Auhtor,A.B. *et al.* (2000) Chapter title. In Smith, A.C. (ed.), *Book Title*, 2nd edn. Publisher, Location, Vol. 1, pp. ???–???.
Bardet, G. (1920) Sur un syndrome d’obesite infantile avec polydactylie et retinite pigmentaire (contribution a l’etude des formes cliniques de l’obesite hypophysaire). PhD Thesis, name of institution, Paris, France.