

TreePPL—a new DSL in Miking for Phylogenetics

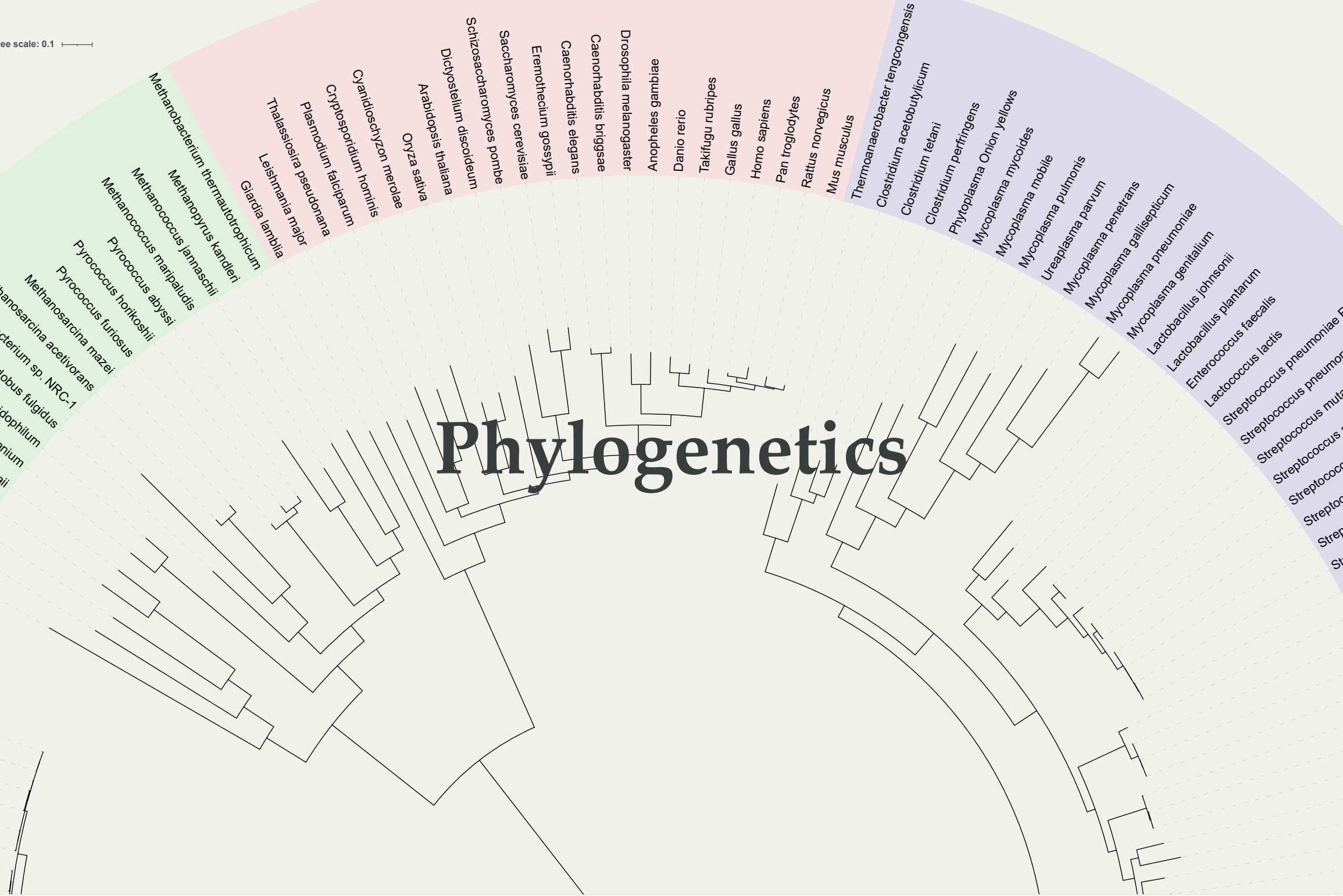
Viktor Senderov, Jan Kudlicka, Daniel Lundén, Viktor Palmkvist, Mariana P. Braga, Emma Granqvist, Fredrik Ronquist, David Broman and contributions from the
Miking team

2023-11-22

Contents

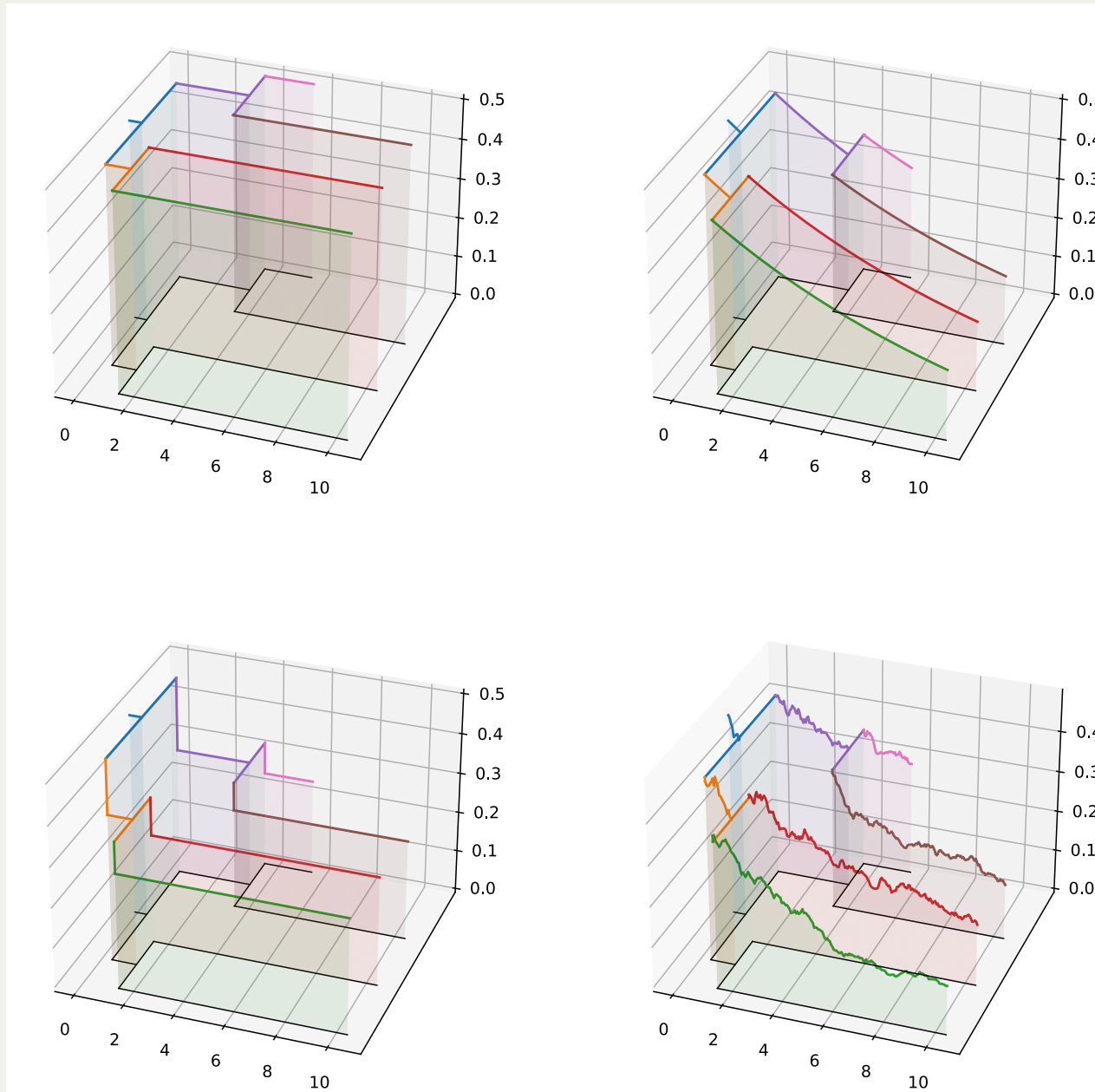
1. A few words about our domain, phylogenetics
2. Probabilistic programming languages
3. TreePPL demo

Scale: 0.1



Phylogenetics discussion

- Problems that phylogenetics deals with
 - Reconstructing phylogenetic trees
 - Using phylogenetic trees to understand evolution and ecology



Diversification models on a phylogenetic tree

What can the users do?

- Methods of phylogenetics
 - Parsimony
 - Statistical methods: ML and Bayesian
- Software for phylogenetics
 - monolithic software: e.g. MrBayes, RAxML
 - extensible software: e.g. BEAST, BEAST2, RevBayes
- Write the model yourself
 - implement in a general-purpose language: e.g. C, Java, Python, or R
 - implement in a probabilistic programming language: e.g. WebPPL, Birch
 - implement in a phylogenetic PPL: TreePPL

A computer monitor is the central focus, displaying a 3D plot of a probability distribution curve. The curve is rendered with a grid and a color gradient from teal to yellow. The plot is set against a grid background with axes labeled with mathematical symbols like μ and σ . A semi-transparent grey text box is overlaid on the center of the monitor screen. In the foreground, a silver speaker with a large driver is visible on the left, and a portion of a keyboard is at the bottom left. The monitor is supported by a silver stand.

Probabilistic Programming Languages

Bayes theorem

Let x be some observed data and θ be an unobserved parameter in whose value we are interested in.

$$\underbrace{p(\theta|x)}_{\text{posterior distribution}} = \frac{\overbrace{p(x|\theta)}^{\text{data distribution / likelihood}} \overbrace{p(\theta)}^{\text{prior distribution}}}{\underbrace{p(x)}_{\text{normalizing constant}}}$$

Bayes theorem

Let x be some observed data and θ be an unobserved parameter in whose value we are interested in.

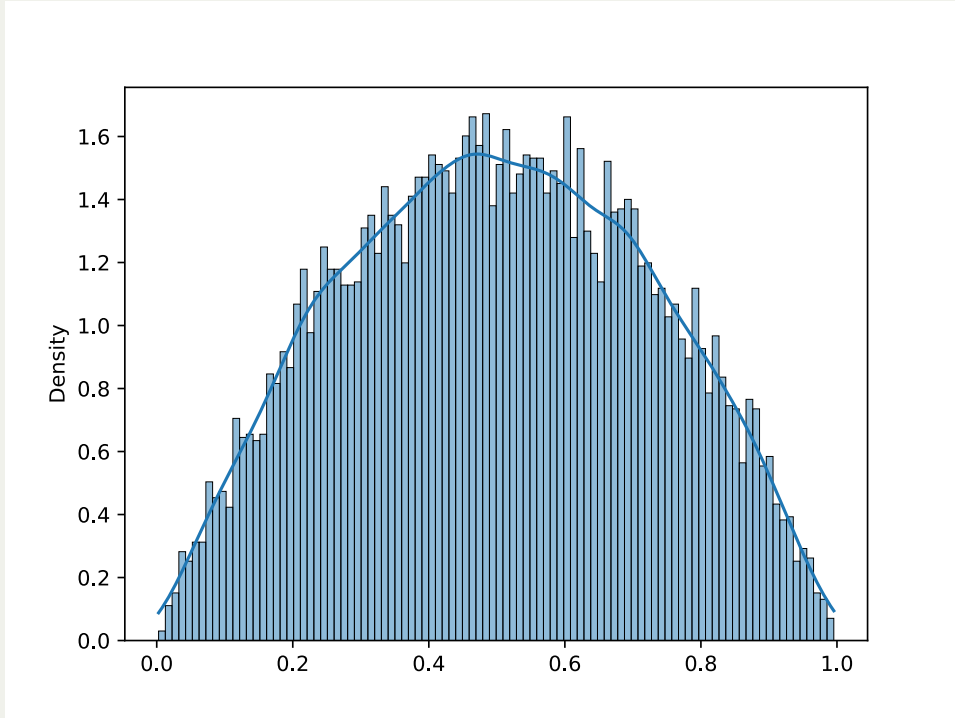
$$\begin{array}{c} \text{posterior distribution} \\ \underbrace{p(\theta|x)} \end{array} = \frac{\begin{array}{c} \text{observe/factor} \\ \text{data distribution/} \\ \text{likelihood} \\ \underbrace{p(x|\theta)} \end{array} \begin{array}{c} \text{assume/sample} \\ \text{prior distribution} \\ \underbrace{p(\theta)} \end{array}}{\begin{array}{c} \underbrace{p(x)} \\ \text{normalizing constant} \end{array}}$$

Coin example

```
1 outcomes=[True, True, True, False, True, False, True, True, True, False,  
2           False, True, True, False, True, False, False, True, False, False]
```

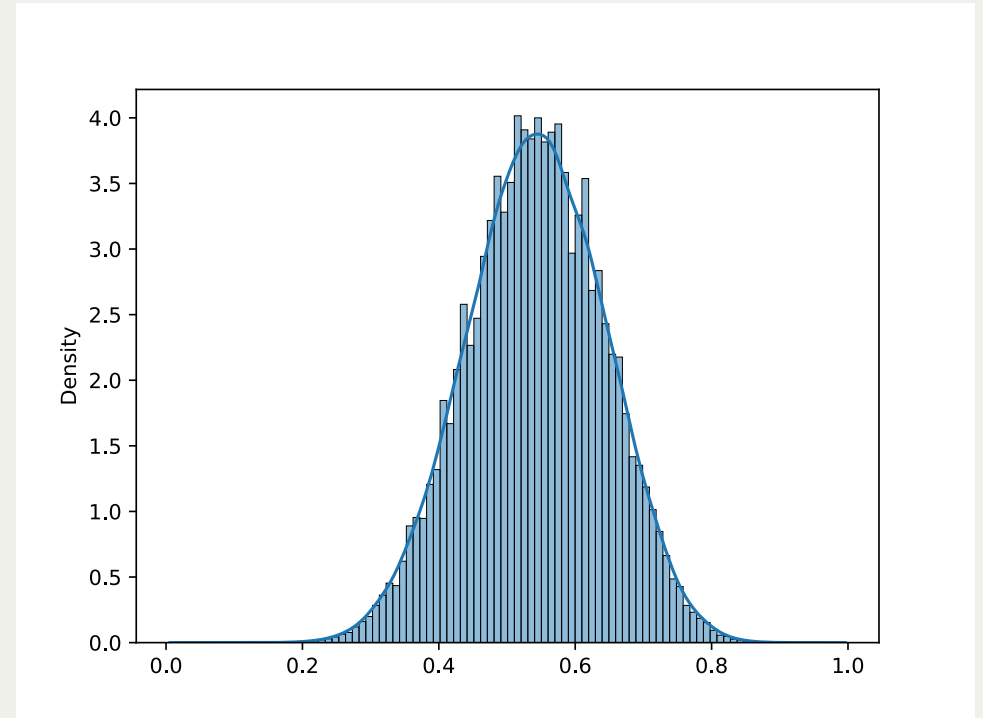
Define prior

```
assume p ~ Beta(2.0, 2.0);
```



Condition on data

```
for i in 1 to (length(outcomes)) {  
  observe outcomes[i] ~  
  Bernoulli(p);  
}
```

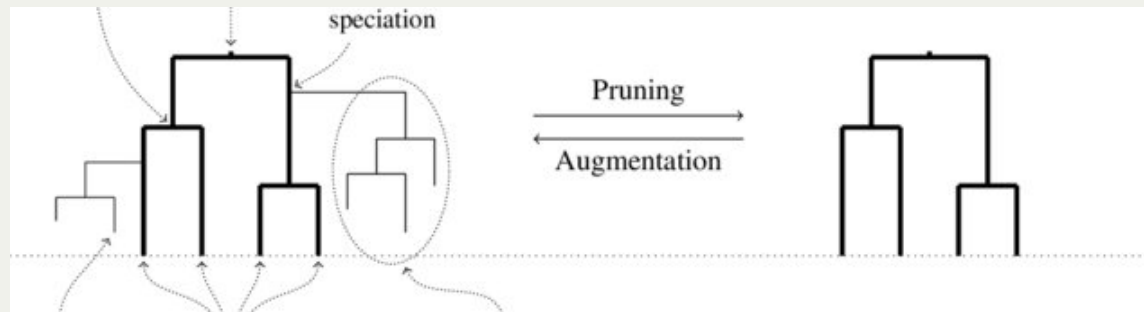


The complete coin example in TreePPL using Python as a wrapper.

```
1 source = """\  
2 model function coin(outcomes: Bool[]): Real {  
3   assume p ~ Beta(2.0, 2.0);  
4   for i in 1 to (length(outcomes)) {  
5     observe outcomes[i] ~ Bernoulli(p);  
6   }  
7   return(p);  
8 }  
9 """  
10  
11 with treeppl.Model(source=source, samples=10_000, method="is-lw") as coin:  
12   res = coin(  
13     outcomes=outcomes  
14   )  
15   sns.histplot(  
16     x=res.samples, weights=res.nweights, bins=100, stat="density", kde=  
17   )  
18   plt.show()
```

Universal PPLs

```
1 function simulateSubtree(time: Real, lambda: Real, mu: Real) {
2   assume waitingTime ~ Exponential(lambda + mu);
3   if waitingTime > time {
4     weight 0.0;
5     resample;
6   } else {
7     assume isSpeciation ~ Bernoulli(lambda / (lambda + mu));
8     if isSpeciation {
9       simulateSubtree(time - waitingTime, lambda, mu);
10      simulateSubtree(time - waitingTime, lambda, mu);
11    }
12  }
13 }
```



Simulating unobserved evolutionary lineages



Phylogenetic Data

Supports natively the PhyJSON format for evolutionary trees



Simplicity

Designed to meet the needs of computational biologists

TreePPL



Rich Model Library

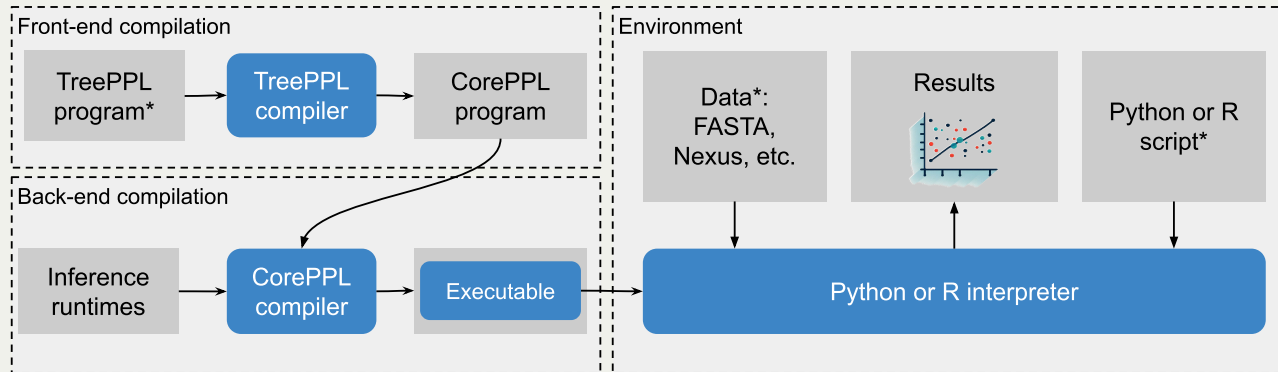
Offers state-of-the-art diversification models as templates



Powerful Statistical Inference

Sequential Monte-Carlo (SMC) and Markov-chain Monte-Carlo (MCMC) inference

Architecture



Optimizations

- Partial CPS transformation
- Automatic alignment
- (Delayed sampling/ conjugacy analysis)

Inference strategies

- Lightweight importance sampling
- Sequential Monte Carlo: Bootstrap Particle Filter,
- Sequential Monte Carlo: Alive particle filter
- Lightweight MCMC
- Trace MCMC
- Particle MCMC

Write a program that infers the bifurcation rate λ and the extinction rate μ under the assumption that they are constant for a given binary evolutionary tree.



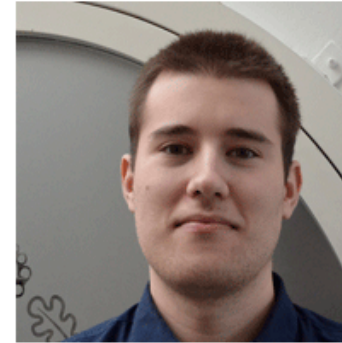
Viktor Senderov

Postdoctoral Researcher at L'École
normale supérieure



Jan Kudlicka

Associate Professor of Data Science
at BI Norwegian Business School



Daniel Lundén

Senior Member of Technical Staff at
Oracle



Viktor Palmkvist

Ph.D. Candidate at KTH Royal
Institute of Technology

Thank you



Mariana P. Braga

Postdoctoral Researcher at the
Swedish University of Agricultural
Sciences



Emma Granqvist

Postdoctoral Researcher at
Department of Bioinformatics and
Genetics, Swedish Museum of
Natural History



Fredrik Ronquist

*PI together with Broman (eq.
contribution)* Department of
Bioinformatics and Genetics, Swedish
Museum of Natural History



David Broman

*PI together with Ronquist (eq.
contribution)* EECS and Digital
Futures, KTH Royal Institute of
Technology and Computer Science
Department, Stanford University

Acknowledgements

- Big Thank You goes to the whole Miking team for the creating the framework that enabled the creation of TreePPL!
- Funded by the EU Horizon 2020 Marie Skłodowska-Curie grant agreement PhyPPL No. 898120, and by the Swedish Foundation for Strategic Research, as well as Vetenskapsrådet (VR).
- Slides have been designed using images from Flaticon.com as well DALL-E

Downloading

Project homepage: <https://treepppl.org>

Main repo: <https://github.com/treepppl/treepppl>

Python library: <https://github.com/treepppl/treepppl-python>

Paper:

<https://www.biorxiv.org/content/10.1101/2023.10.10.561673v1>

Citation:

TreePPL: A Universal Probabilistic Programming Language for Phylogenetics
Viktor Senderov, Jan Kudlicka, Daniel Lundén, Viktor Palmkvist, Mariana P.
Braga, Emma Granqvist, David Broman, Fredrik Ronquist
bioRxiv 2023.10.10.561673; doi: <https://doi.org/10.1101/2023.10.10.561673>

Bonus material

Discussion of PPLs

- PPL semantics
 - Usual language: *input* \rightarrow *output* + *side effects*
 - PPL: *each program trace describes a sample from a probability distribution*
- Advantages of PPLs
 - natural and generative way of encoding a posterior distribution
 - no need to implement inference
 - enables model comparison as estimation of normalizing constant possible

