# TreePPL—a new DSL in Miking for Phylogenetics

Viktor Senderov

École normale supérieure - PSL

2023-11-22
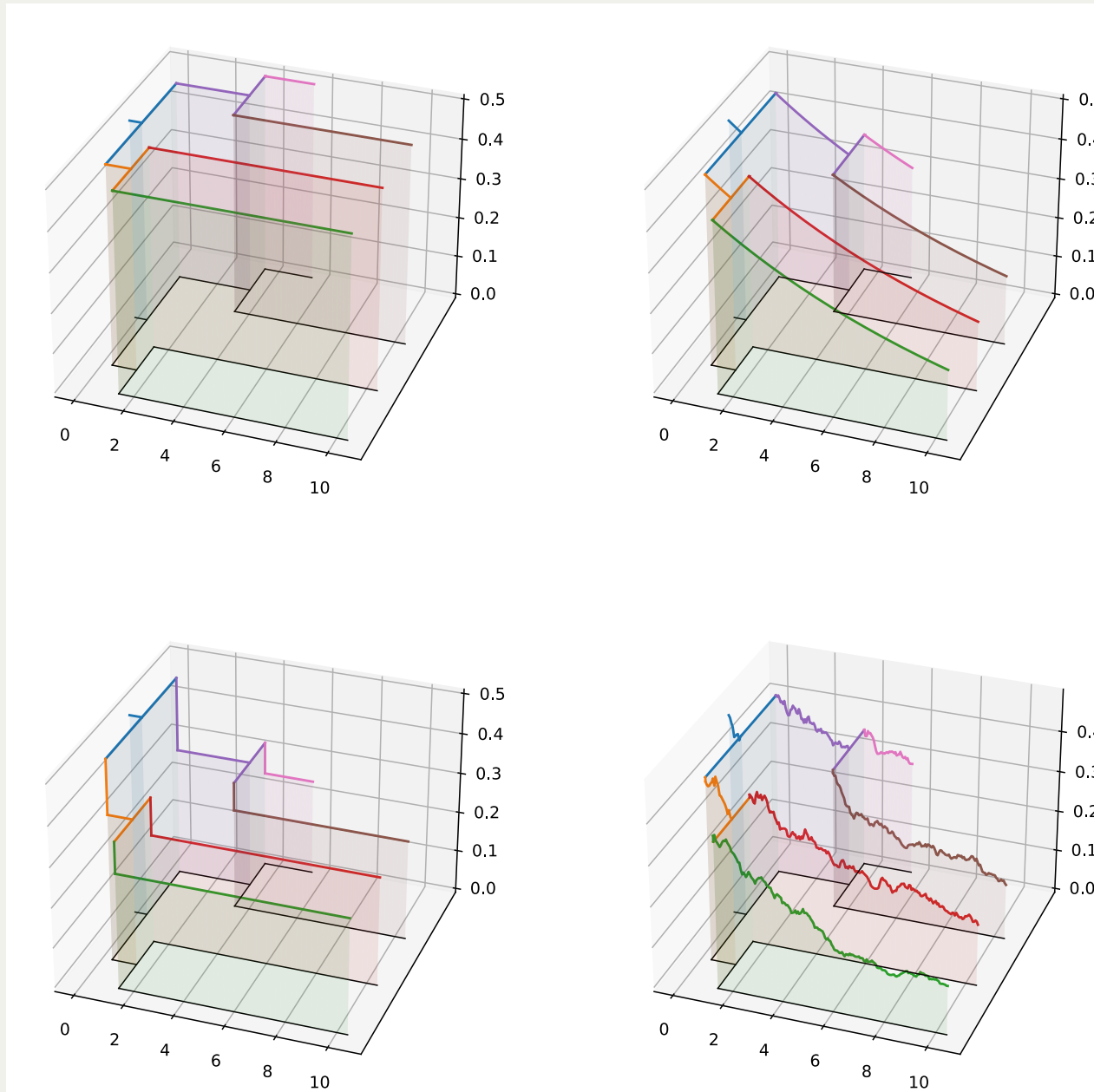
# Contents

1. A few words about our domain, phylogenetics

2. Probabilistic programming languages

3. TreePPL demo

Tree scale: 0.1 ⊢━━━┤

# Phylogenetics

Methanobacterium thermautotrophicum
Methanopyrus kandleri
Methanococcus jannaschii
Methanococcus maripaludis
Pyrococcus abyssi
Pyrococcus horikoshii
Pyrococcus furiosus
Methanosarcina mazei
Methanosarcina acetivorans
...bacterium sp. NRC-1
...lobus fulgidus
...idophilum
...nium
...ii

Giardia lamblia
Leishmania major
Thalassiosira pseudonana
Plasmodium falciparum
Cryptosporidium hominis
Cyanidioschyzon merolae
Oryza sativa
Arabidopsis thaliana
Dictyostelium discoideum
Schizosaccharomyces pombe
Saccharomyces cerevisiae
Eremothecium gossypii
Caenorhabditis elegans
Caenorhabditis briggsae
Drosophila melanogaster
Anopheles gambiae
Danio rerio
Takifugu rubripes
Gallus gallus
Homo sapiens
Pan troglodytes
Rattus norvegicus
Mus musculus

Thermoanaerobacter tengcongensis
Clostridium acetobutylicum
Clostridium tetani
Clostridium perfringens
Phytoplasma Onion yellows
Mycoplasma mycoides
Mycoplasma mobile
Ureaplasma parvum
Mycoplasma pulmonis
Mycoplasma penetrans
Mycoplasma gallisepticum
Mycoplasma pneumoniae
Mycoplasma genitalium
Lactobacillus johnsonii
Lactobacillus plantarum
Enterococcus faecalis
Lactococcus lactis
Streptococcus pneumoniae P...
Streptococcus pneumoniae...
Streptococcus muta...
Streptococcus...
Streptococc...
Strepto...
Stre...
St...

# 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 phylogenic PPL: TreePPL

# Bayes theorem

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

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

# Bayes theorem

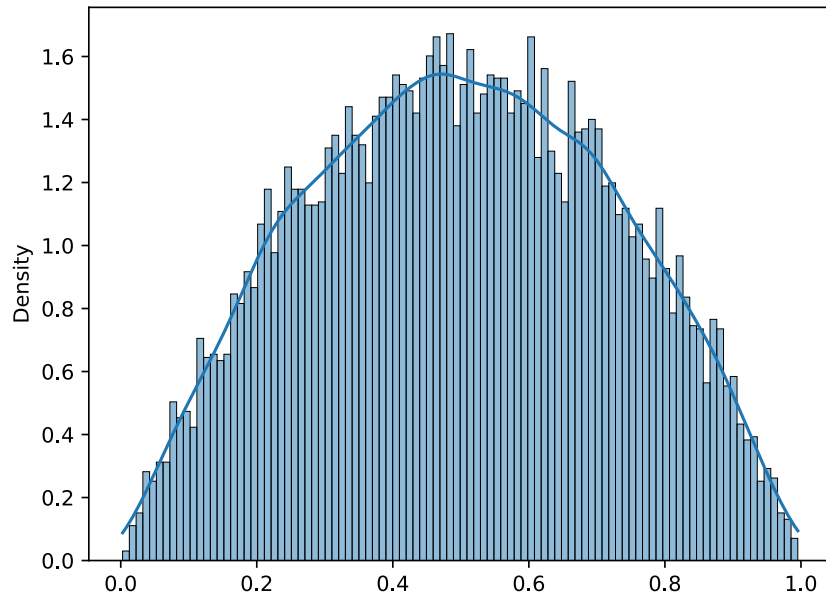Let $x$ be some observed data and $\theta$ be an unobserved parameter in whose value we are interested in in.

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

# 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]
```
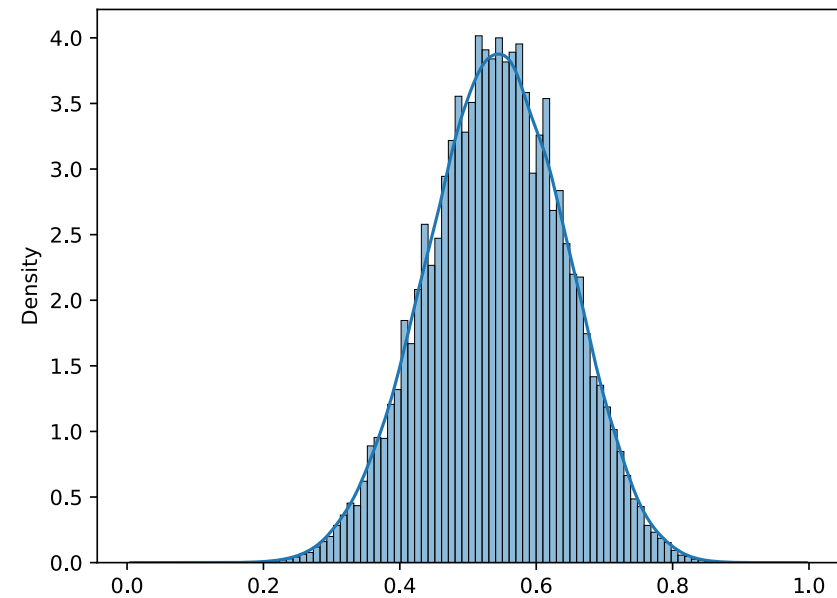
# Define prior

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



# Condition on data

```
for i in 1 to (length(outcomes)) {
    observe outcomes[i] ~
Bernoulli(p);
}
```
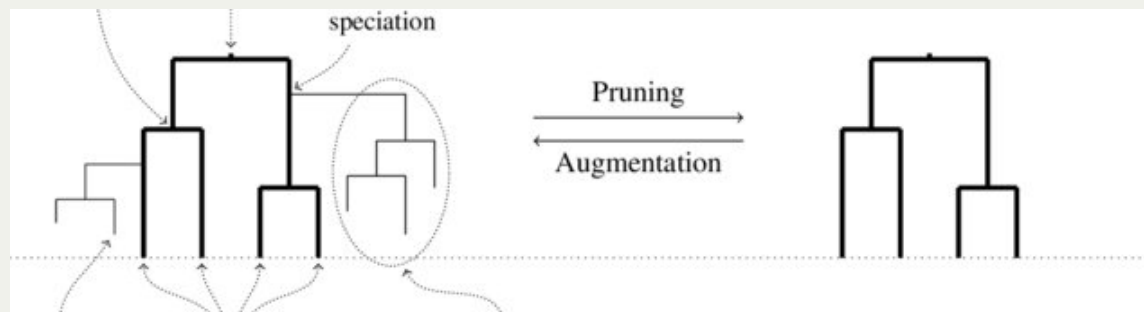
# The complete coin example in TreePPL using Python as a wrapper.

```python
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 evolutoinary trees

# TreePPL

**Rich Model Library**
Offers state-of-the art diversificatio
models as templates

**Simplicity**
Designed to meet the needs of
computational biologists

**Powerful Statistical Inference**
Sequential Monte-Carlo (SMC) and
Markov-chain Monte-Carlo (MCMC)
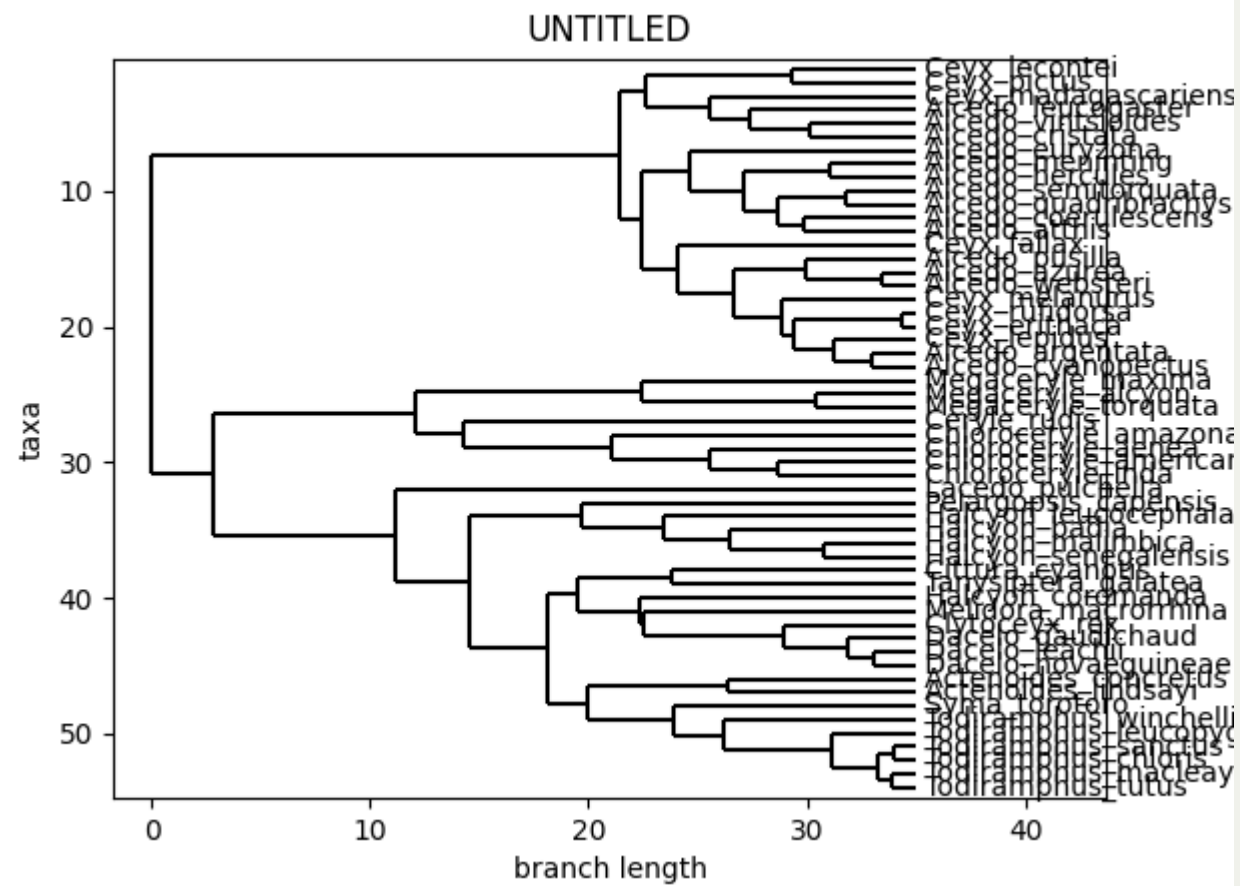inference

# 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 $\lambda$ and the extinction rate $\mu$ under the assumption that they are constant for a given binary evolutionary tree.

UNTITLED

taxa

branch length

Ceyx_lecontei
Ceyx_pictus
Ceyx_madagascariensis
Alcedo_leucogaster
Alcedo_vintsioides
Alcedo_cristata
Alcedo_euryzona
Alcedo_meninting
Alcedo_hercules
Alcedo_semitorquata
Alcedo_quadribrachys
Alcedo_coerulescens
Alcedo_atthis
Ceyx_lepidus
Alcedo_pusilla
Alcedo_websteri
Ceyx_melanurus
Ceyx_erithaca
Ceyx_cyanopecta
Alcedo_argentata
Alcedo_cyanopectus
Alcedo_azurea
Megaceryle_alcyon
Megaceryle_torquata
Chloroceryle_amazona
Chloroceryle_americana
Chloroceryle_inda
Alcedo_pulchella
Halcyon_capensis
Halcyon_leucocephala
Halcyon_malimbica
Halcyon_senegalensis
Tanysiptera_galatea
Halcyon_coromanda
Melidora_macrorrhina
Ceyx_lecontei
Dacelo_gaudichaud
Dacelo_novaeguineae
Actenoides_lindsayi
Syma_torotoro
Todiramphus_winchelli
Todiramphus_leucopygius
Todiramphus_sanctus
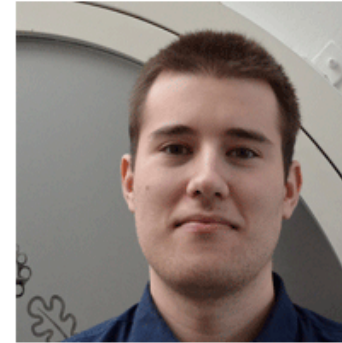Todiramphus_macleayi
Todiramphus_tutus

**Viktor Senderov**

Postdoctoral Researcher at L'École normale supérieure
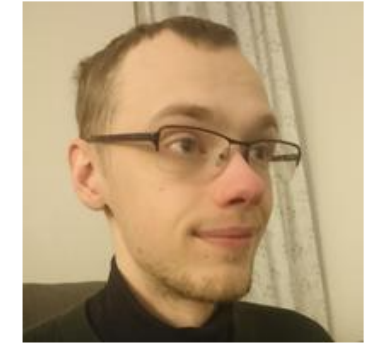
**Jan Kudlicka**

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**

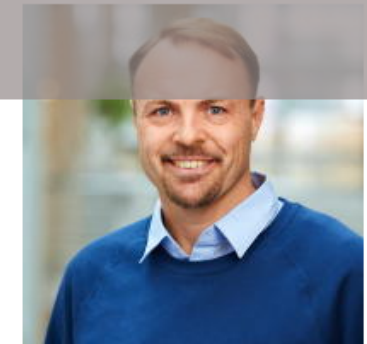Postdoctoral Researcher at the Swedish University of Agricultural Sciences

**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

- Slides have been designed using images from Flaticon.com as well DALL-E

# Downloading

**Project homepage**: https://treeppl.org

**Main repo**: https://github.com/treeppl/treeppl

**Python library**: https://github.com/treeppl/treeppl-python

**Paper**:
https://www.biorxiv.org/content/10.1101/2023.10.10.561673v1

Citation:

# Bonus material

# Discussion of PPLs

- PPL semantics
  - Usual language: *input → 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