

# Deep likelihood-free inference of phylogenetic trees



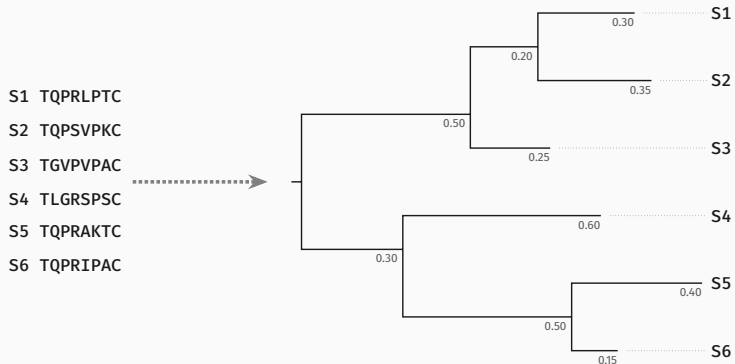
---

**Luc Blassel**, Nicolas Lartillot, Bastien Boussau, Laurent Jacob

MASAMB - September 8<sup>th</sup>, 2025



# Context - Phylogenetic inference

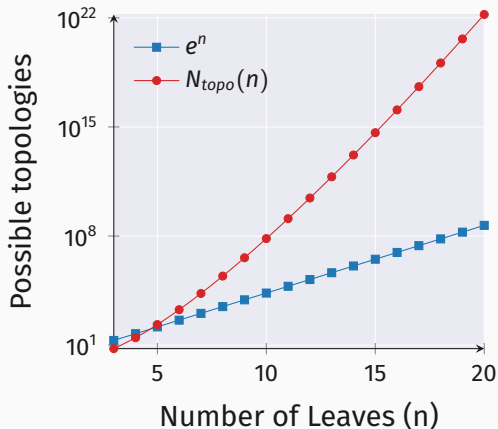


*Goal:* describe **evolutionary-history** of MSA

# Context - The problem with phylogenetic inference

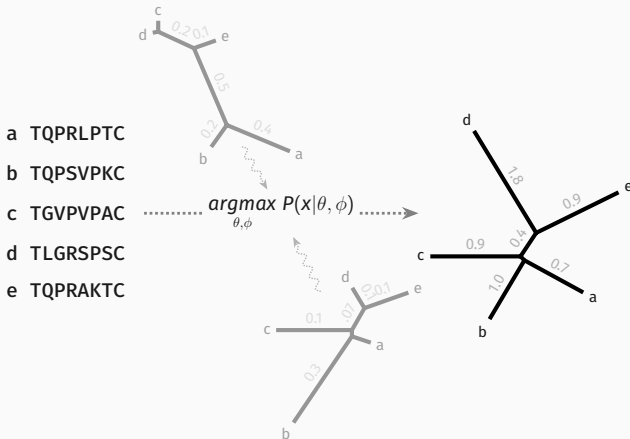
1. Phylogenies are **hard!**
2. **Super-exponential** tree space

$$N_{topo}(n) = \frac{(2n-3)!}{2^{n-2}(n-2)!}$$



Felsenstein 1978

# Context - Likelihood-based tree reconstruction



$x$  : MSA,  $\theta = (\tau, \ell)$  : Phylogenetic tree,  $\phi$  : Evolution model

# Context - Likelihood-based tree reconstruction

## Pros:

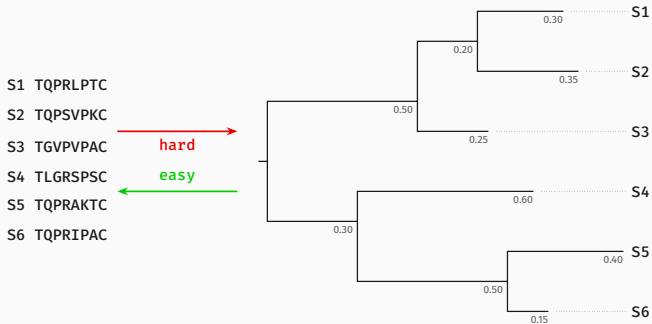
- These methods are **accurate**
- The **whole MSA** is considered in  $P(x|\theta, \phi)$

## Cons:

- These methods are **slow**
  1. **Computing** the likelihood is **costly**
  2. We have to **explore** the tree-space with **topological** moves
- We are **limited** to models where  $P(x|\theta, \phi)$  is **computable**

Felsenstein 1993; Kleinman et al. 2010

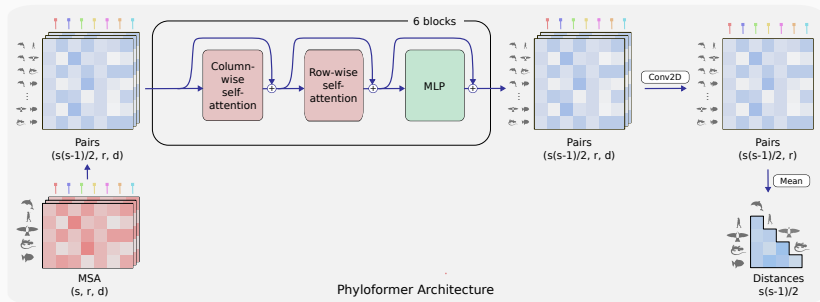
# Motivation - Likelihood-free inference



- We can simulate many<sup>1</sup> (tree, MSA) pairs
- Can we **learn** the mapping **from MSA to tree**?

<sup>1</sup> pretty much practically  $\infty$

# Related Work - Phyloformer, our first approach

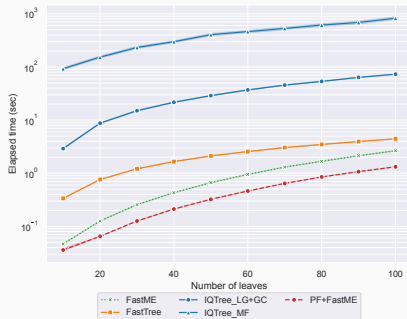
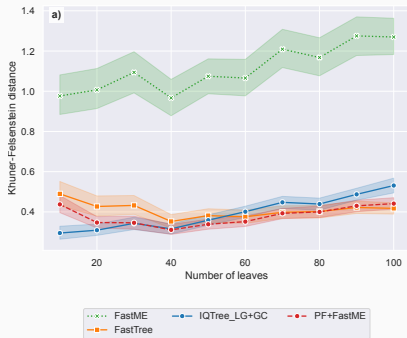


- Input an **MSA**, get a **Distance matrix**
- Feed Distance matrix to **FastME** to get **tree**

Nesterenko et al. 2025; Lefort et al. 2015



# Related Work - **Phyloformer** is good!



## Tree inference accuracy (KF)

- Fairly **competitive** even on simple LG+GC model
- Fast** because we use GPUs <sup>1</sup>

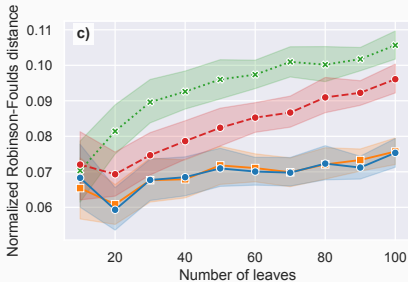
Nesterenko et al. 2025, <sup>1</sup> 🙏 Jean-Zay

## Runtime

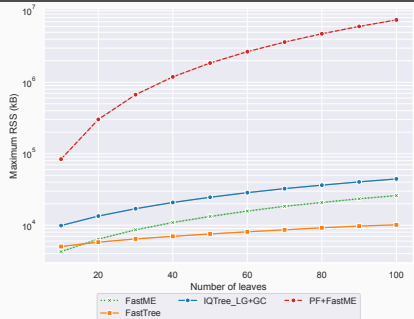




# Related Work - But also sometimes less good...



*Topological accuracy (RF)*



*Memory usage*

- **Gap** between PF and **ML methods**
- PF is **by far** the most **memory intensive**



## Related Work - Why does Phyloformer struggle with topology ?

- Phyloformer predicts **distance** matrices, as **proxy** for trees
- In **theory** it is **equivalent**, but in practice ...
- Could we get rid of the **proxy**, and **predict trees** directly ?

## **How to do phylogenetic inference end-to-end ?**

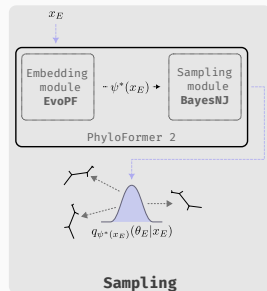
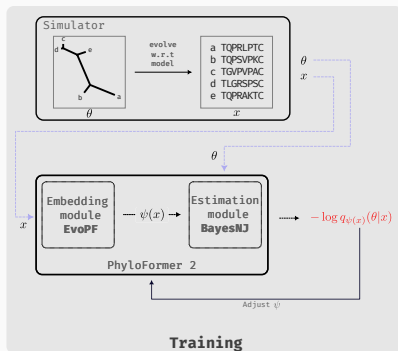
---

# Methods - Neural Posterior Estimation (NPE)

- Given a **probabilistic model**  $p(x|\theta)$  with some prior  $p(\theta)$
- We want to **estimate the posterior**:  $p(\theta|x)$
- We build  $q_\psi(\theta|x)$  a **family** of distributions **parametrized** by  $\psi$  (our NN)
- We find  $q_{\psi^*} = \underset{\psi}{\operatorname{argmin}} \mathbb{E}_{p(x)}[KL(q_\psi(\theta|x)||p(\theta|x))]$
- In practice we **maximize**  $\mathbb{E}_{p(x,\theta)}[\log q_{\psi(x)}(\theta|x)]$  by **sampling** from  $p(x, \theta)$

$x$  : MSA,  $\theta = (\tau, \ell)$  : Phylogenetic tree,  $\psi(x)$  : NN applied to  $x$

# Methods - How do we do NPE?



- During **training** find  $\psi^* = \underset{\psi}{\operatorname{argmin}} - \sum_i \log q_{\psi(x_i)}(\theta_i|x_i)$
- At **inference** time **sample** from:  $q_{\psi^*(x_E)}(\theta_E|x_E)$

## Methods - The EvoPF module, intro

the EvoPF module is an **adaptation** of the **EvoFormer** module from **AlphaFold2**. The tasks are **transpositions** of each other:

given input MSA ( $n \times r$ )

**EvoFormer** represent  $r \times r$  relationships between sites

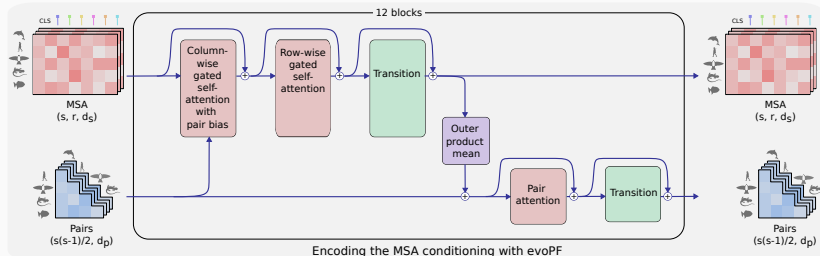
**EvoPF** represent  $n \times n$  relationships between sequences

**More expressive** than MSA transformer

**More lightweight** than PF

Jumper et al. [2021](#); Rao et al. [2021](#)

# Methods - The EvoPF module, details

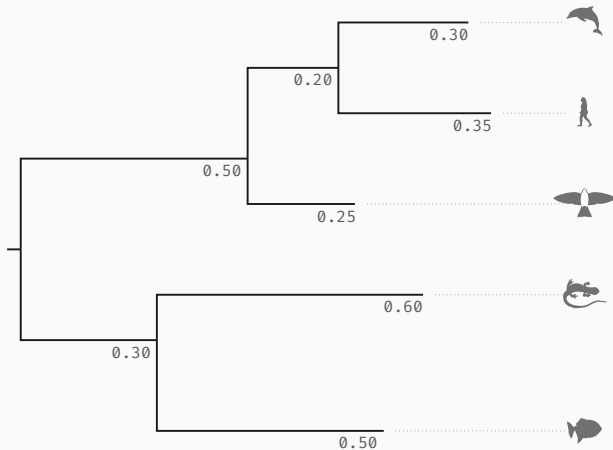


- Input an **MSA** and get:
  - sequence** embedding  $\{s_i\}$
  - sequence-pair** embeddings  $\{z_{ij}\}$
- **Both** embedding-types used to **update each-other**

Figure inspired by Jumper et al. 2021

# Methods - A tree is a series of merges

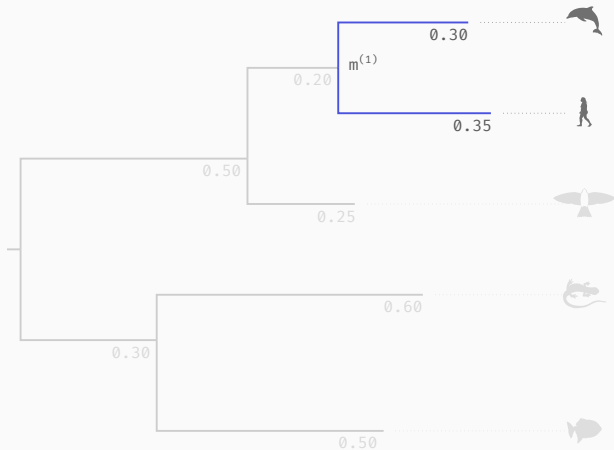
We want to describe the following tree:





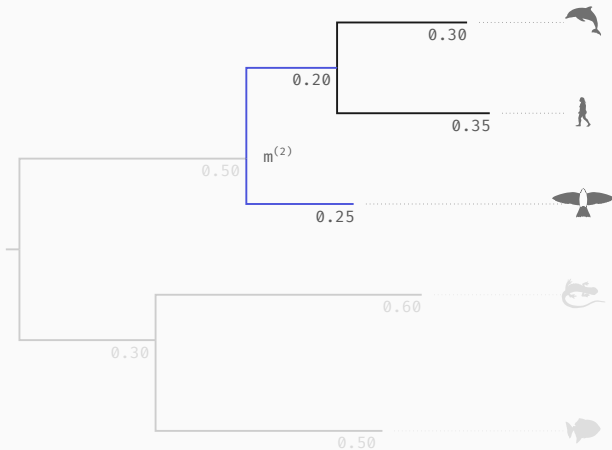
# Methods - A tree is a series of merges

Iteratively create cherries:



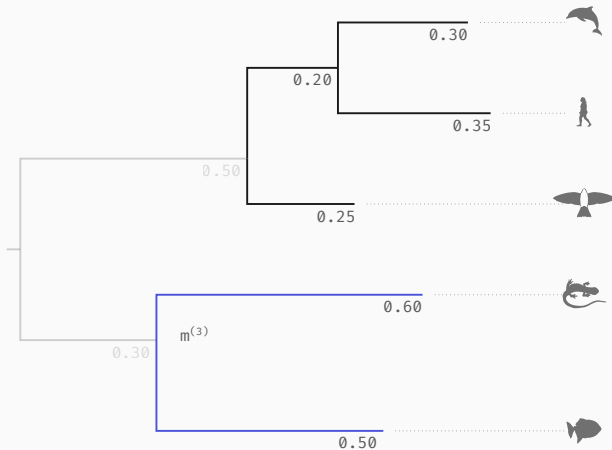
## Methods - A tree is a series of merges

## Iteratively create cherries:



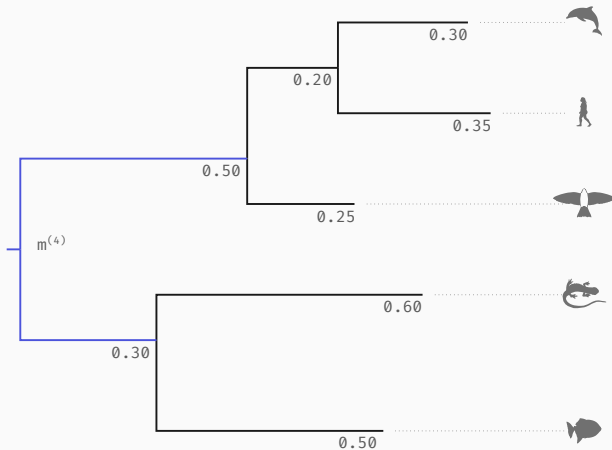
# Methods - A tree is a series of merges

Iteratively create cherries:



# Methods - A tree is a series of merges

Iteratively create cherries:



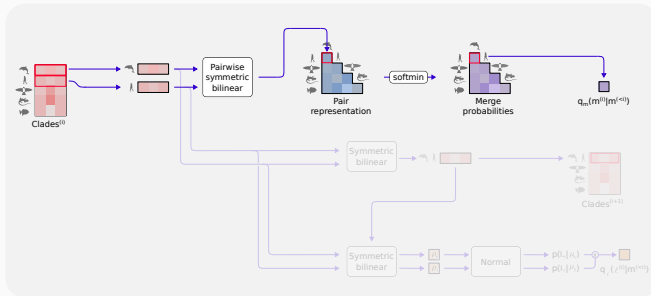
- **Tree** is an **ordered set** of merges:  $\theta : \{m^{(1)}, \dots, m^{(N-1)}\}$
- We **factorize**  $q_{\psi(x)}(\theta|x)$  as the product of successive merge probabilities:

$$q_{\psi(x)}(\theta|x) = \prod_{k=1}^{N-1} q_m(m^{(k)}|m^{(<k)}) q_\ell(\ell^{(k)}|m^{(\leq k)})$$

- **Merge probabilities have 2 components:**
  - topological:**  $q_m(m^{(k)}|m^{(<k)})$
  - branch-length:**  $q_\ell(\ell^{(k)}|m^{(\leq k)})$

# Methods - BayesNJ, evaluating topological probabilities

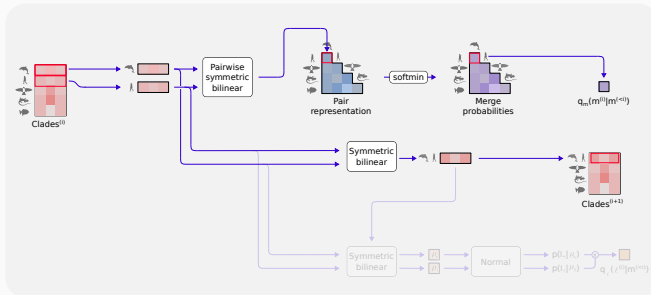
$$\mathbf{m}^{(i)} = (\mathbf{a}_i, \mathbf{l}_i) \quad \mathcal{L}^{(i)} = (\mathbf{l}_i, \mathbf{l}_i)$$



Compute **merge probability**

# Methods - BayesNJ, evaluating topological probabilities

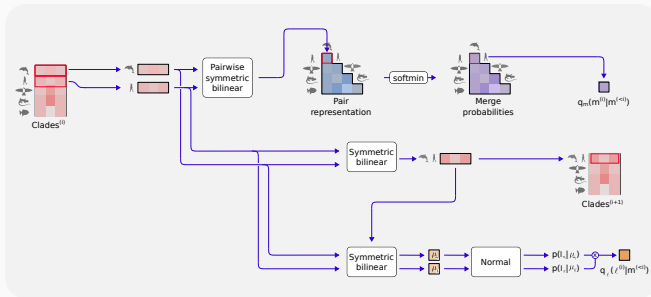
$$\mathbf{m}^{(i)} = (\mathbf{m}, \mathbf{l}) \quad \ell^{(i)} = (\mathbf{l}, \mathbf{l})$$



**Update** clade **representation** for next merge

# Methods - BayesNJ, evaluating branch length probabilities

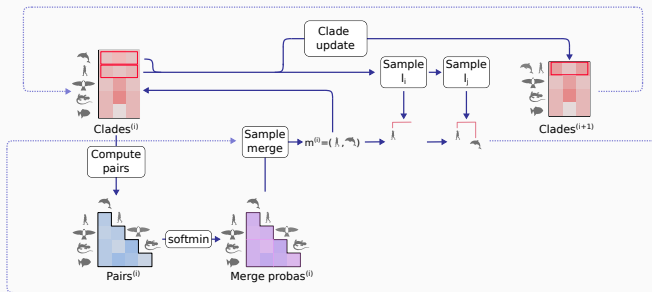
$$\mathbf{m}^{(i)} = (\mathbf{a}, \mathbf{b}) \quad \ell^{(i)} = (l_a, l_b)$$



Compute **branch-length** probabilities



# Methods - BayesNJ sampling mode

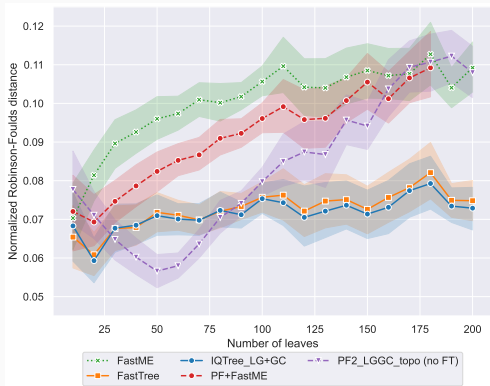


- **Sample** merges and branch lengths **until** topology resolved
- **Two** sampling **modes** given  $\psi(x_E)$ :
  - Bayesian** Sample from distributions
  - Greedy MAP** Choose mode

**Does it work ?**

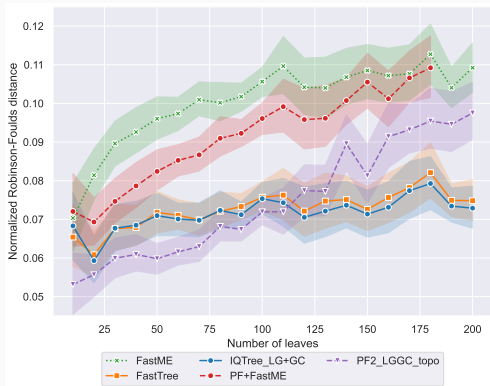
---

# Results - Training topology only



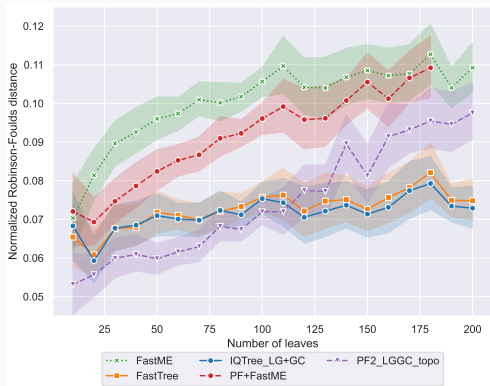
- **overfitting** on tree-size is an **issue**

# Results - Training topology only



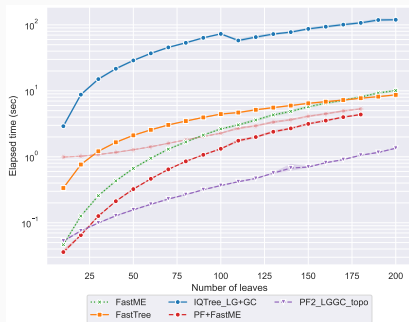
- **overfitting** on tree-size is an **issue**
- **Fine tuning** helps

# Results - Training topology only

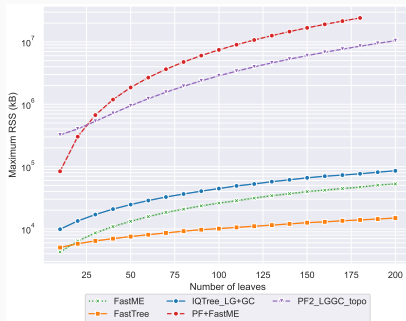


- **overfitting** on tree-size is an **issue**
- **Fine tuning** helps
- We **beat ML** methods in certain cases
- Marked **improvement w.r.t Phyloformer**

# Results - Scalability



**Execution time**



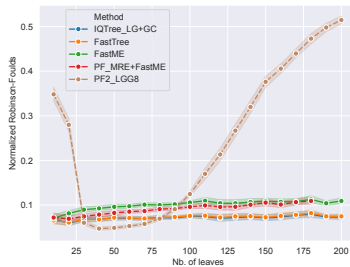
**Memory usage<sup>1</sup>**

<sup>1</sup> With  $2\times$  bigger sequence, and  $4\times$  bigger pair embeddings...

# Results - What next ?

This is very much still a **work in progress...**

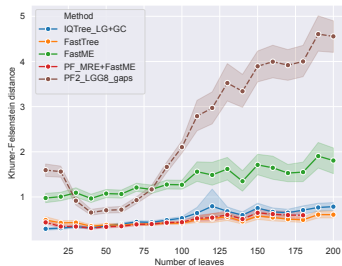
- Training with **gaps** is more **complicated**



# Results - What next ?

This is very much still a **work in progress...**

- Adding **branch lengths** is **harder** than we thought

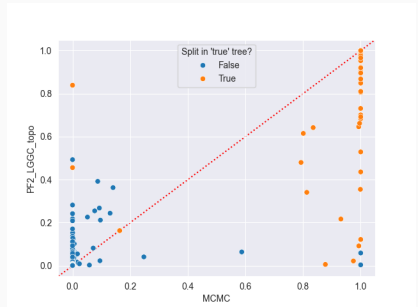




# Results - What next ?

This is very much still a **work in progress...**

- We need to **adjust our priors** to compare to MCMC



# Perspectives - Intractable likelihoods

- **Topologically** we manage to **beat** IQTree<sup>1</sup> on LG
- Can we do **better** with complex models where computing  $p(\theta|x)$  is **difficult** or **intractable**?
- **Interaction** models:
  - **CherryML**, residue pair coevolution
  - **Potts** models, How do we simulate ?
  - **Epistasis** models
- Models taking **selection** into account: e.g. SelReg
- **Confident** this can **work** given our experience with **PF**

Prillo et al. 2023; Duchemin et al. 2023; Latrille et al. 2021

<sup>1</sup> Yay!

# Conclusion

- **WIP** but we are close to truly **end-to-end likelihood-free** phylogenetic **inference**
- Still **limitations**:
  - **Better** than **PF** but **scalability** is still an **issue**
  - **Length overfitting** also an issue
- **Where do we go** once PF2 is done ?
  - Extend to **unaligned** sequence
  - Predict **Ancestral** sequences or characters
  - **Downstream** tasks: population dynamics, reconciliation, epidemiology, ecology ...

## Thanks to:

- Luca Nesterenko
- Laurent Jacob
- Bastien Boussau
- Nicolas Lartillot
- Philippe Veber
- Vincent Garot
- Amélie Leroy
- Anybody that listened to me!



Special thanks to Jean-Zay for all the GPUs!

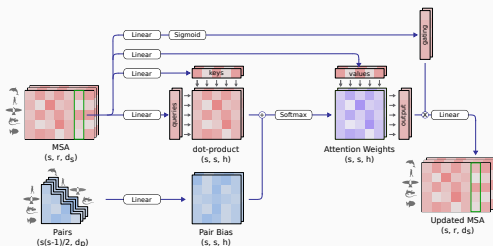
## References

---

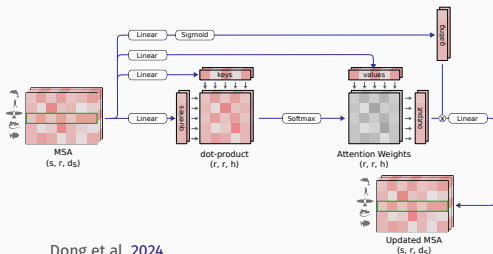
- Dong, J. et al. (2024). ***Flex Attention: A Programming Model for Generating Optimized Attention Kernels.***
- Duchemin, L. et al. (2023). ***Evaluation of methods to detect shifts in directional selection at the genome scale.*** In: *Molecular Biology and Evolution* 40.2, msac247.
- Felsenstein, J. (1978). ***The Number of Evolutionary Trees.*** In: *Systematic Zoology* 27.1, p. 27.
- (1993). ***PHYLP (phylogeny inference package), version 3.5 c.*** Joseph Felsenstein.
- Jumper, J. et al. (2021). ***Highly accurate protein structure prediction with AlphaFold.*** In: *Nature* 596.7873, pp. 583–589.
- Kleinman, C. L. et al. (2010). ***Statistical Potentials for Improved Structurally Constrained Evolutionary Models.*** In: *Molecular Biology and Evolution* 27.7, pp. 1546–1560.
- Latrille, T. et al. (2021). ***Inferring Long-Term Effective Population Size with Mutation–Selection Models.*** In: *Molecular Biology and Evolution* 38.10, pp. 4573–4587.
- Lefort, V. et al. (2015). ***FastME 2.0: a comprehensive, accurate, and fast distance-based phylogeny inference program.*** In: *Molecular biology and evolution* 32.10, pp. 2798–2800.

- Nesterenko, L. et al. (2025). **Phyloformer: Fast, Accurate, and Versatile Phylogenetic Reconstruction with Deep Neural Networks**. In: *Molecular Biology and Evolution* 42.4, msaf051.
- Prillo, S. et al. (2023). **CherryML: scalable maximum likelihood estimation of phylogenetic models**. In: *Nature methods* 20.8, pp. 1232–1236.
- Rao, R. M. et al. (2021). **MSA Transformer**. In: *Proceedings of the 38th International Conference on Machine Learning*. Ed. by M. Meila and T. Zhang. Vol. 139. Proceedings of Machine Learning Research. PMLR, pp. 8844–8856.

# Supp. Methods - EvoPF, the MSA stack



**Column-wise attention  
with pair-bias**

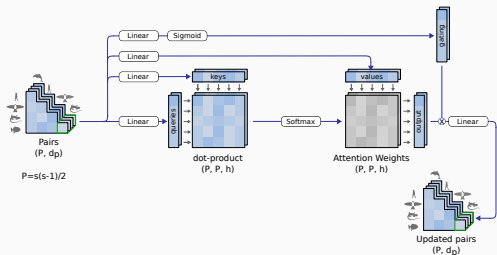


**Row-wise attention**

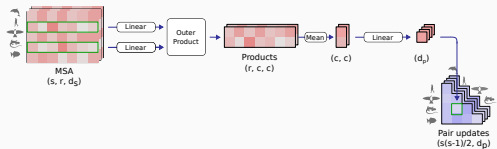
Dong et al. 2024



# Sup. Methods - EvoPF, the pair stack



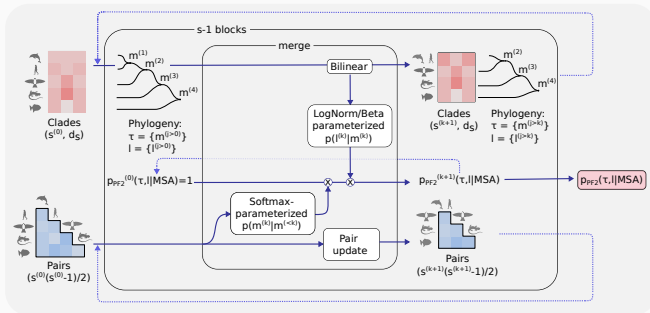
**Pair attention**



**Outer product mean**

Dong et al. 2024

# Sup. Methods - BayesNJ evaluation mode



## Sup. Methods - Ensuring the merge order is unique

Ensuring a **unique order** on merges ensures that we **define a distribution**. It also keeps **training** and **sampling** comparable <sup>1</sup>

- On a given tree  $\tau$  always **merge** the **shortest** available **cherry**
- When **sampling**, add **constraints**:
  1. Start with a  $N \times N$  constraints matrix  $M_{ij} = 0$
  2. At iteration  $k$  sample merge  $m^{(k)} = (i, j)$  and cherry length  $s^{(k)} = M_{ij} + X$
  3. **Update constraints** for cherries **available** when sampling  $m^{(k)}$ :  $M'_{ij} = \max(M_{ij}, s^{(k)})$   $M'_{ui} = 0$
- During evaluation compute  $p_{PF2}(s^{(k)} - M_{ij} | m^{(\leq k)})$

<sup>1</sup> Which is not the same if we use the NJ merge order

# Sup. Methods - Tree simulation

