

Ancestral protein reconstruction using autoregressive generative models

Matteo De Leonardis, Andrea Pagnani, **Pierre Barrat-Charlaix**

DISAT, Politecnico di Torino

Collaborators

Andrea Pagnani (PoliTo)

Matteo De Leonardis (PoliTo)

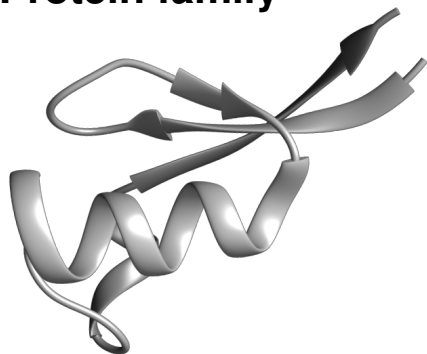


**Politecnico
di Torino**

Statistical modeling of protein sequences














Multiple Sequence Alignment (MSA)

Protein family



Evolutionary constraints



												
		I	H	D	L	R	H	T	N	D	K	..
		L	H	N	L	R	G	T	D	D	R	..
		E	H	R	T	E	Q	L	E	K	G	..
		Y	H	L	L	R	T	L	D	D	T	..
	.	R	H	A	V	E	M	L	N	K	G	..
	.	Q	H	K	L	E	E	A	N	K	A	..

In different species:

- ~ same structure
- ~ same function

- Regulation (DNA-binding, protein inhibitor...)
- Signaling (two-component signaling)
- Fundamental (ribosome...)
- Antibiotic resistance







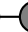
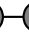
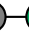





Statistical modeling of protein sequences

Protein family



Evolutionary
constraints



											
	I	H	D	L	R	H	T	N	D	K	..
	L	H	N	L	R	G	T	D	D	R	..
	E	H	R	T	E	Q	L	E	K	G	..
	Y	H	L	L	R	T	L	D	D	T	..
.	R	H	A	V	E	M	L	N	K	G	..
.											
.	Q	H	K	L	E	E	A	N	K	A	..

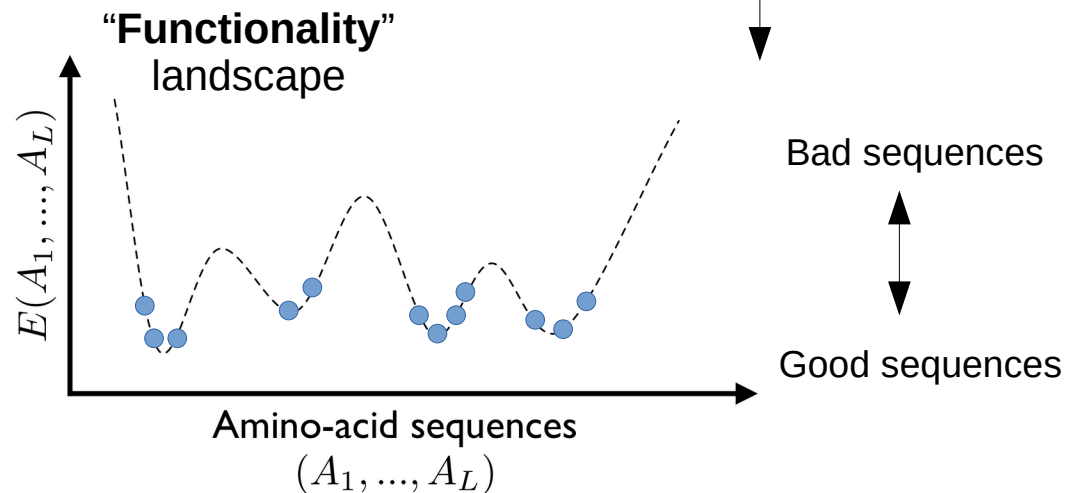
Multiple Sequence
Alignment (MSA)

Learning

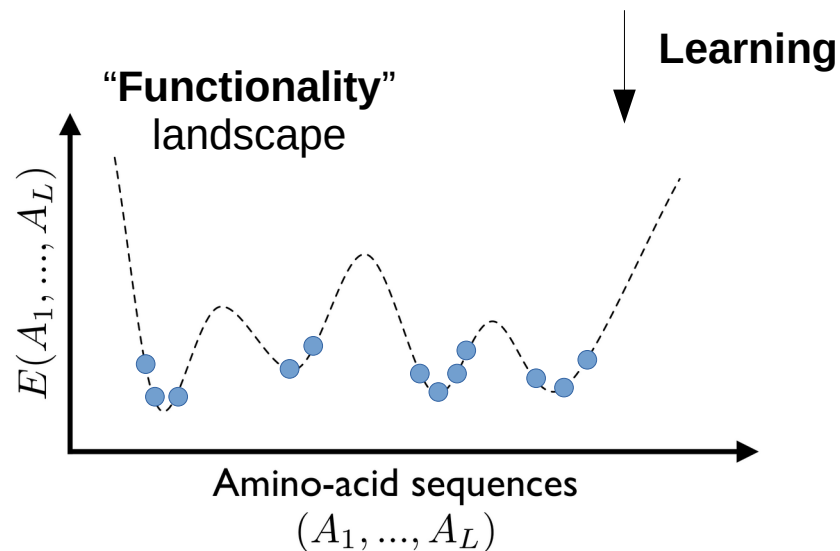
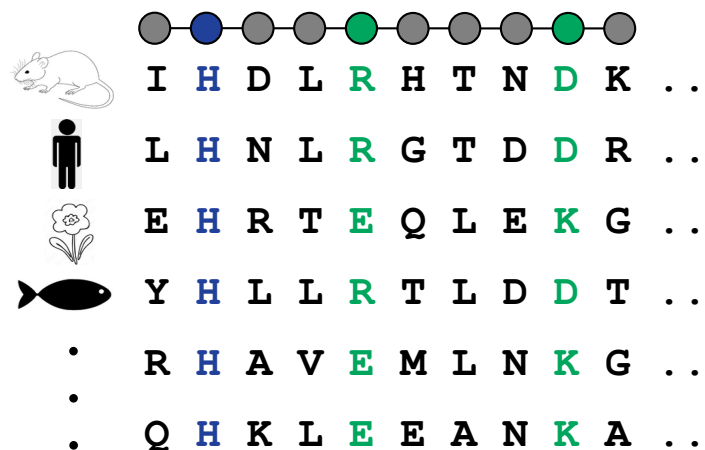
In different species:

- ~ same structure
- ~ same function

- Regulation (DNA-binding, protein inhibitor...)
- Signaling (two-component signaling)
- Fundamental (ribosome...)
- Antibiotic resistance



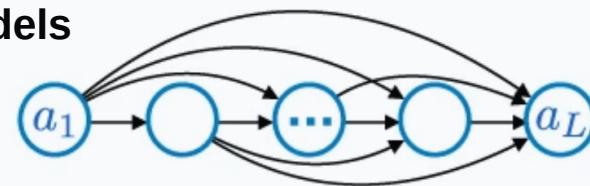
Statistical modeling of protein sequences



Potts model

$$P(a_1, \dots, a_N) = \frac{1}{Z} \exp \left(\sum_{i,j=1}^L \overset{\text{Couplings}}{J_{ij}(a_i, a_j)} + \sum_{i=1}^L \overset{\text{Fields}}{h_i(a_i)} \right)$$

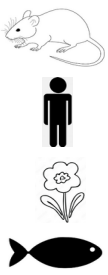
Autoregressive models



Deep learning

Variational autoencoders (VAE)
Transformers (MSATransformer, ESM, ...)

Statistical modeling of protein sequences



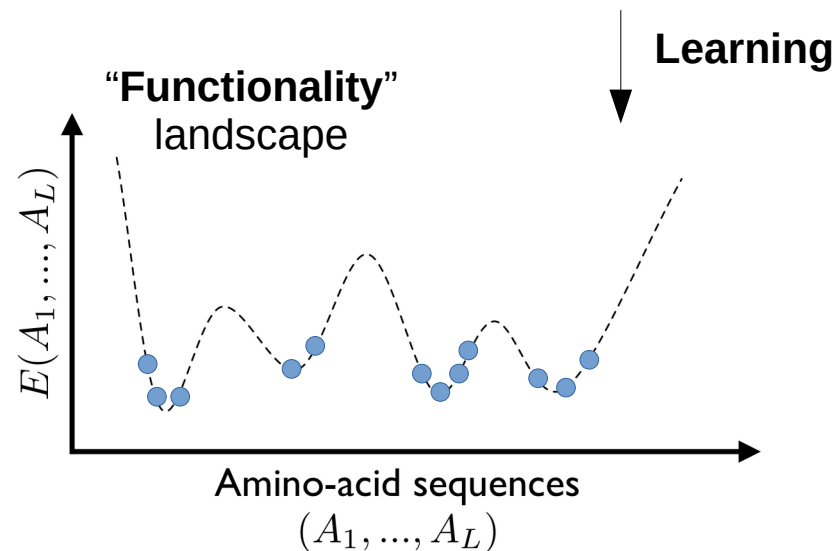
I	H	D	L	R	H	T	N	D	K	..	
L	H	N	L	R	G	T	D	D	R	..	
E	H	R	T	E	Q	L	E	K	G	..	
Y	H	L	L	R	T	L	D	D	T	..	
•	R	H	A	V	E	M	L	N	K	G	..
•											
•	Q	H	K	L	E	E	A	N	K	A	..

Used for

- Contacts in 3D structure
- Effect of mutations
- **Generative models**

Design **functional**
synthetic proteins!

[Russ et. al. Science 2020]

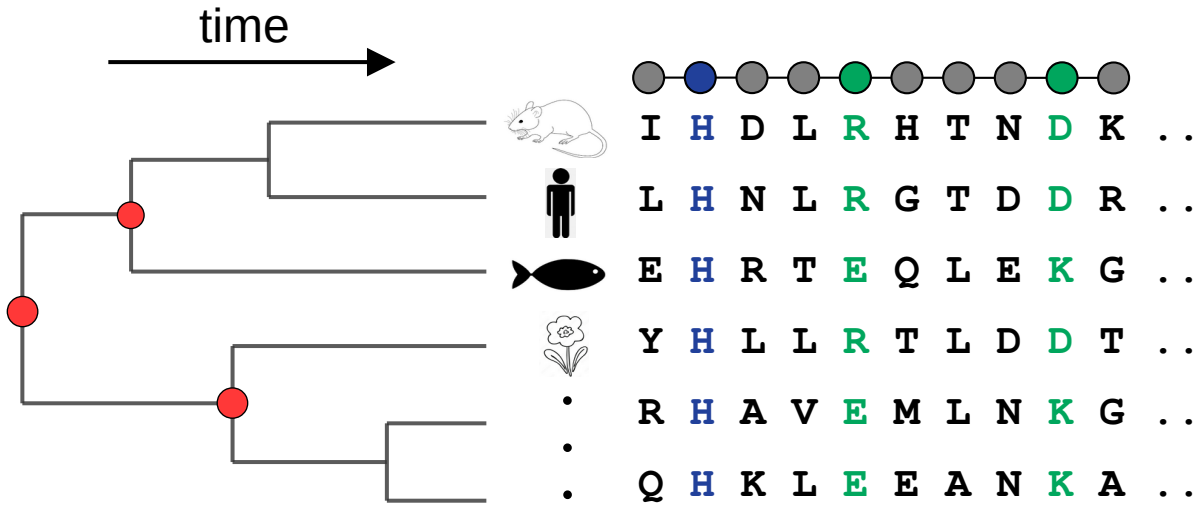


Key ingredient: **Epistasis**

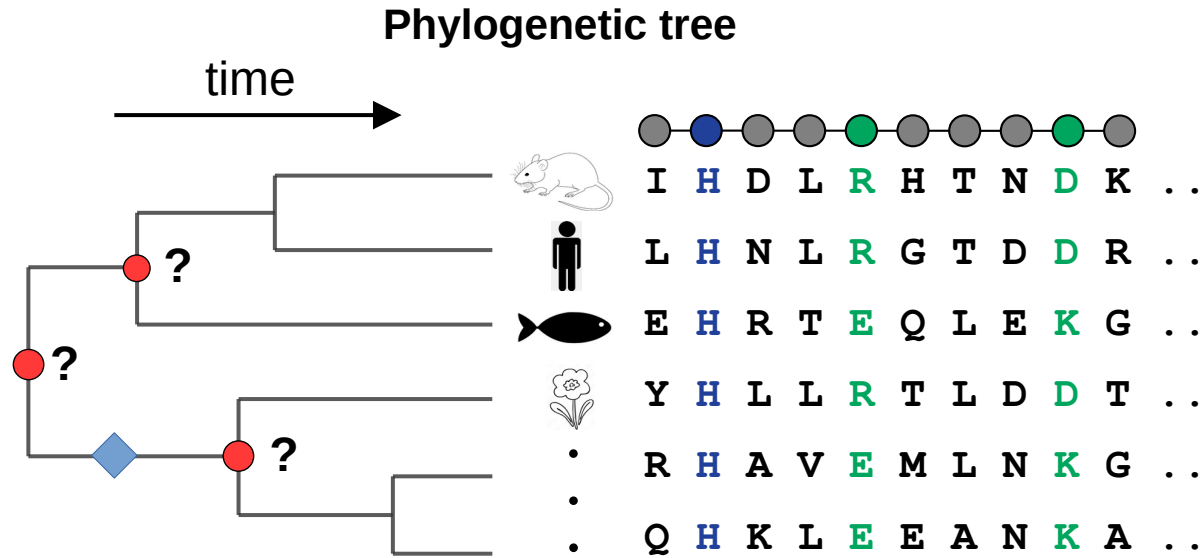
→ Columns of the MSA are **not** independent

Ancestral sequence reconstruction

Phylogenetic tree



Ancestral sequence reconstruction



Why?

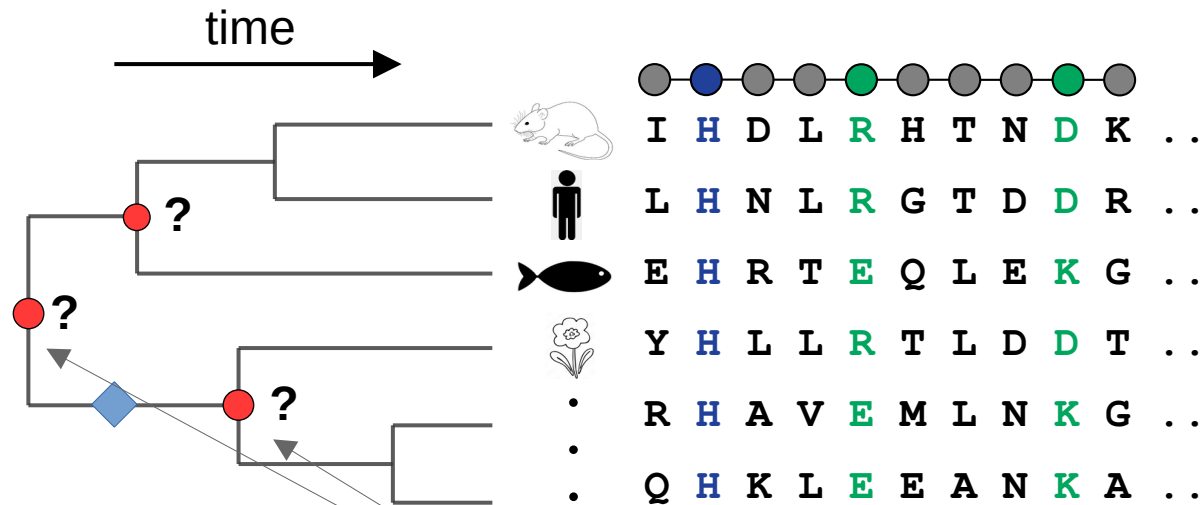
- What were ancient proteins like?
- Sequence – Function relationship

Applications to protein design

- Thermostable proteins
- Proteins with given specificity

Ancestral sequence reconstruction

Phylogenetic tree



Why?

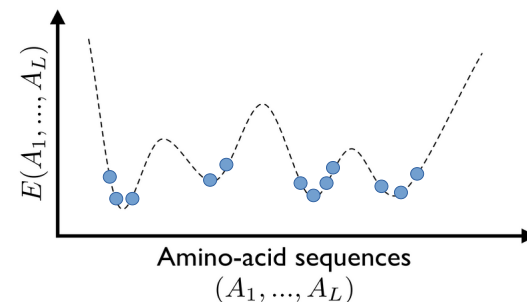
- What were ancient proteins like?
- Sequence – Function relationship

Applications to protein design

- Thermostable proteins
- Proteins with given specificity

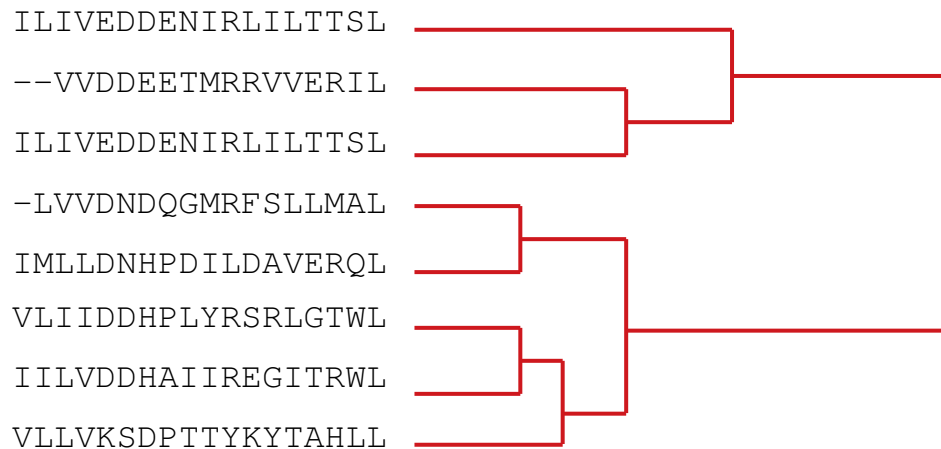
Generative model

$$P(a_1, \dots, a_N) = \frac{1}{Z} \exp \left(\sum_{i,j=1}^L J_{ij}(a_i, a_j) + \sum_{i=1}^L h_i(a_i) \right)$$



Ancestral sequence reconstruction

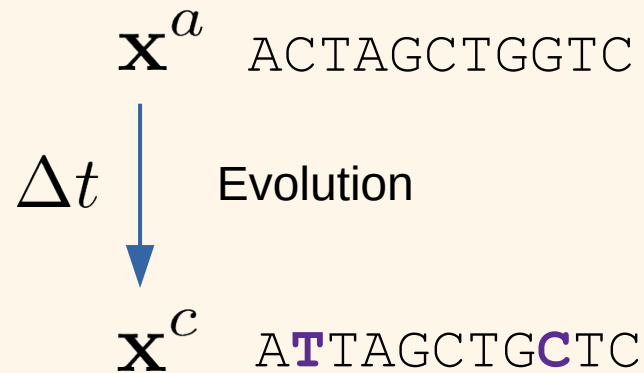
Phylogenetic tree topology + branch length



Phylogenetic inference

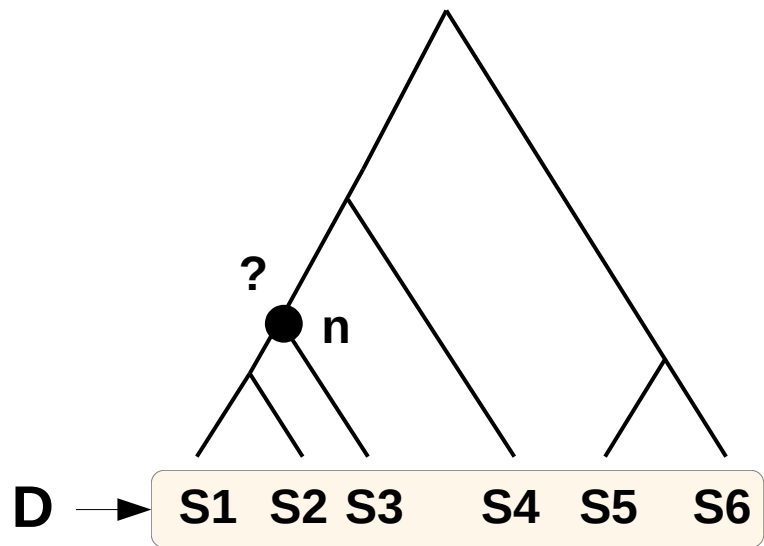
- IQ-TREE – [Minh *et. al.*, MBE 2020]
- Fasttree – [Price *et. al.*, PLOS 2010]
- BEAST – [Suchard *et. al.*, Virus Evol. 2018]

Sequence evolution model



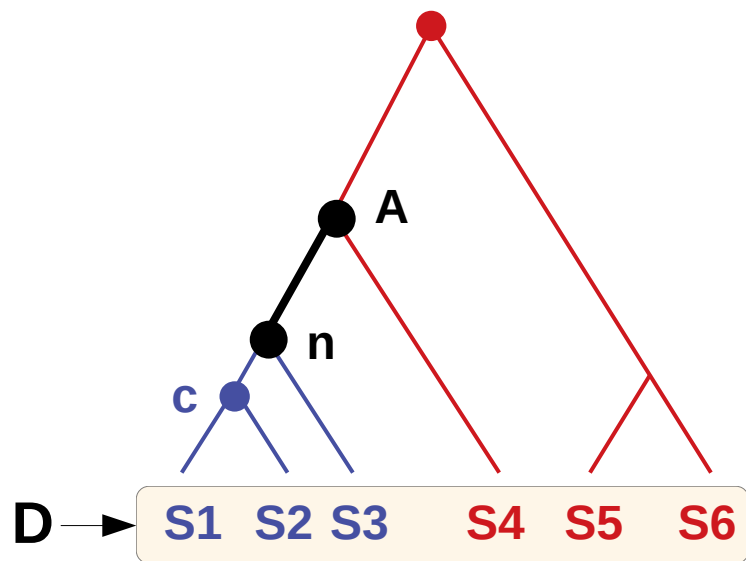
$$P(\mathbf{x}^c | \mathbf{x}^a, \Delta t)$$

Inference of internal states: Felsenstein's algorithm



$$\mathcal{L}_n(\mathbf{x}) = P(\mathcal{D} | n = \mathbf{x})$$

Inference algorithm



$$\mathcal{L}_n(\mathbf{x}) = P(\mathcal{D} | n = \mathbf{x}) \quad \text{Here} \quad P(\mathbf{S}_3 | \mathbf{x}, t_3) \cdot \sum_{\{\mathbf{z}\}} P(\mathbf{z} | \mathbf{x}, t_c) \mathcal{L}_c^d(\mathbf{z})$$

→ One pass up from the leaves to compute $\mathcal{L}_n^d(\mathbf{x})$

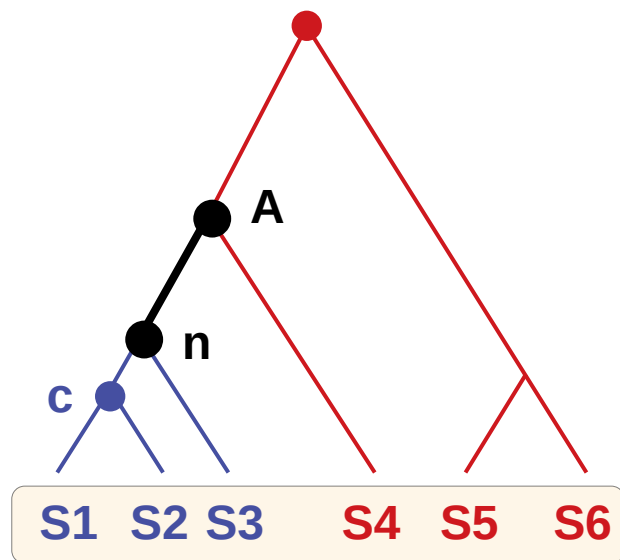
“Down” likelihood

$\mathcal{L}_n^d(\mathbf{x}) \rightarrow$ Probability of data below n , if n is in state \mathbf{x}

If n is leaf $\rightarrow \mathcal{L}_{\text{leaf}}^d(\mathbf{x}) = \delta(\mathbf{x}, \mathbf{S})$

Otherwise $\rightarrow \mathcal{L}_n^d(\mathbf{x}) = \prod_{c \in \mathcal{C}(n)} \sum_{\{\mathbf{z}\}} P(\mathbf{z} | \mathbf{x}, t_c) \mathcal{L}_c^d(\mathbf{z})$

Inference algorithm



$$\mathcal{L}_n(\mathbf{x}) = P(\mathcal{D}|n = \mathbf{x})$$

“Down” likelihood

$\mathcal{L}_n^d(\mathbf{x}) \rightarrow$ Probability of data below n , if n is in state \mathbf{x}

$$\mathcal{L}_n^d(\mathbf{x}) = \prod_{c \in \mathcal{C}(n)} \sum_{\{\mathbf{z}\}} P(\mathbf{z}|\mathbf{x}, t_c) \mathcal{L}_c^d(\mathbf{z})$$

“Up” likelihood

$\mathcal{L}_n^u(\mathbf{y}) \rightarrow$ Probability of data above n , if A is in state \mathbf{y}
where $A \sim \text{ancestor}(n)$

$$\mathcal{L}_n^u(\mathbf{y}) = \sum_{\{\mathbf{z}\}} P(\mathbf{y}|\mathbf{z}, t_A) \mathcal{L}_A^u(\mathbf{z}) \cdot \prod_{c \in \mathcal{C}(A)} \sum_{\{\mathbf{z}\}} P(\mathbf{z}|\mathbf{y}, t_c) \mathcal{L}_c^d(\mathbf{y})$$

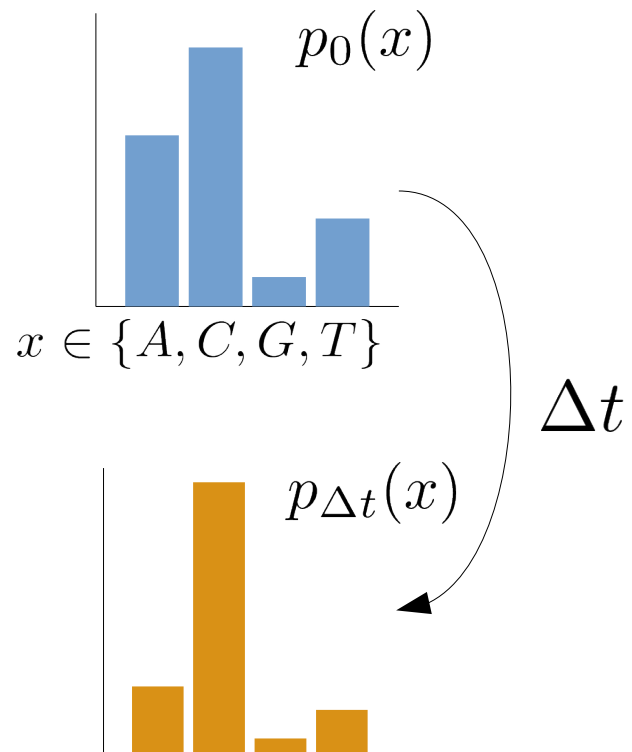


$$\mathcal{L}_n(\mathbf{x}) = \sum_{\{\mathbf{y}\}} \mathcal{L}_n^u(\mathbf{y}) P(\mathbf{x}|\mathbf{y}, t_n) \mathcal{L}_n^d(\mathbf{x})$$

- **linear** in number of nodes
- depends only on **evolution model**

Sequence evolution model

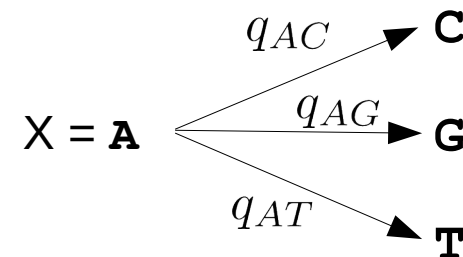
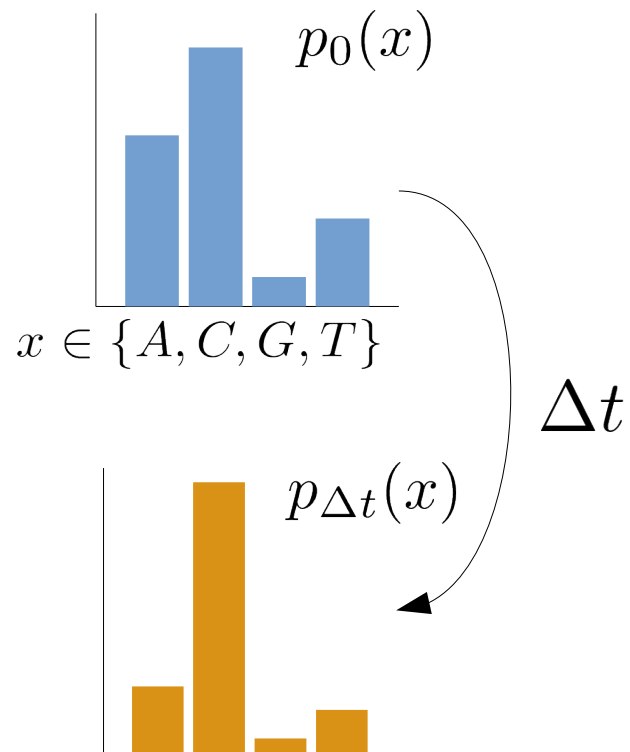
Focus on one position: $x \in \{A, C, G, T\}$



Sequence evolution model

Continuous time Markov chain

Focus on one position: $x \in \{A, C, G, T\}$



$$Q = \begin{pmatrix} -q_A & q_{CA} & q_{GA} & q_{TA} \\ q_{AC} & -q_C & q_{GC} & q_{TC} \\ q_{AG} & q_{CG} & -q_G & q_{TG} \\ q_{AT} & q_{CT} & q_{GT} & -q_T \end{pmatrix}$$

$$\dot{p} = p \cdot \mu Q$$

$$p_{\Delta t} = p_0 \cdot e^{\mu \Delta t \cdot Q}$$

$$q_x = \sum_{y \neq x} q_{xy}$$

Transition rate matrix Q

Reversibility: $\pi_x P(y|x, t) = \pi_y P(x|y, t)$

If reversibility

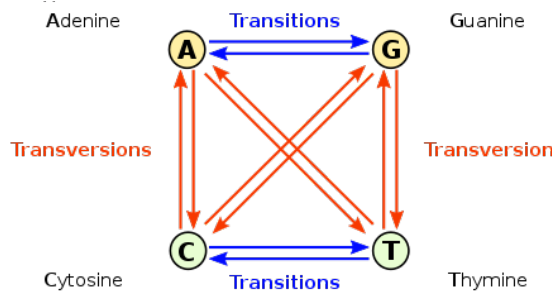
$$Q = \mathbf{H} \cdot \begin{pmatrix} \pi_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \pi_q \end{pmatrix}$$

Q encodes

- Possibilities of mutations
- Equilibrium distribution

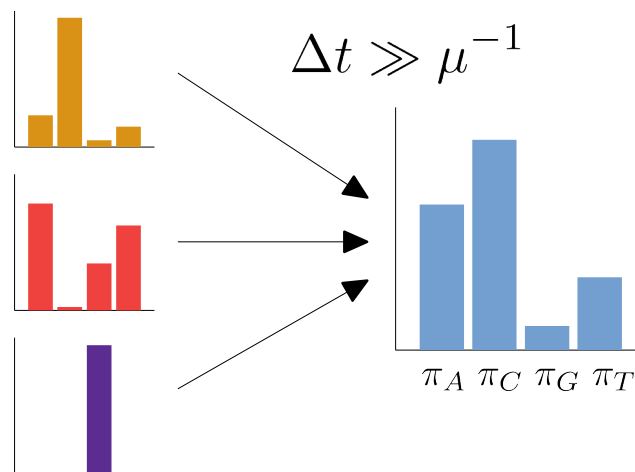
Nucleotides

H → **Symmetric**



Amino acids → **H ~ Genetic code**

$(\pi_1 \dots \pi_q) \rightarrow$ **Equilibrium state**



Independent-site reconstruction

$$p_{\Delta t} = p_0 \cdot e^{\mu \Delta t \cdot Q}$$

One Q_i per sequence position i

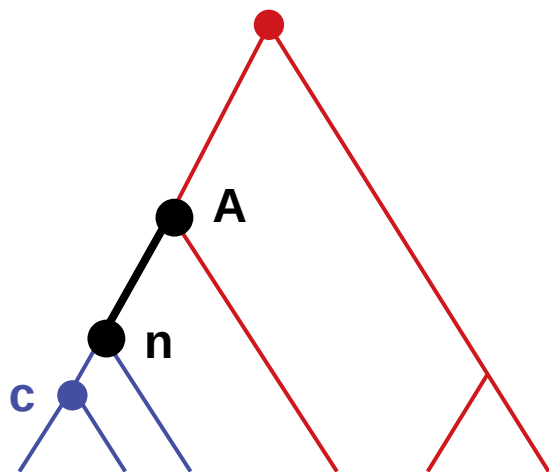
Reconstruct for one position i

“Down”
likelihood

$$\mathcal{L}_{x,i}^d(x_i) = \prod_{c \in \mathcal{C}(n)} \sum_{z_i=1}^q (e^{Q_i t_c})_{x_i z_i} \mathcal{L}_c^d(z_i)$$

$$\mathcal{L}_n(x_i) = P(\mathcal{D} | n_i = x_i)$$

How do we find the right Q matrices?



Independent-site reconstruction

[Jones et al., 1992]
[Le and Gascuel, 2008]

State of the art

$$p_{\Delta t} = p_0 \cdot e^{\mu \Delta t \cdot Q}$$

$$Q = \mathbf{H} \cdot \begin{pmatrix} \pi_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \pi_q \end{pmatrix}$$

- Q
- Same for all positions
 - Fixed matrix, pre-learned (JTT, LG, ...)

- μ
- Can change across positions. Typically
- Four rates $\{\mu_1, \dots, \mu_4\}$
 - Proportion of invariable sites $\mu = 0$

Drawbacks

- Consider positions independently
 - **Ignores** functional constraints
-

Independent-site reconstruction

[Jones et al., 1992]
[Le and Gascuel, 2008]

State of the art

$$p_{\Delta t} = p_0 \cdot e^{\mu \Delta t \cdot Q}$$

$$Q = \mathbf{H} \cdot \begin{pmatrix} \pi_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \pi_q \end{pmatrix}$$

- Q
- Same for all positions
 - Fixed matrix, pre-learned (JTT, LG, ...)

- μ
- Can change across positions. Typically
 - Four rates $\{\mu_1, \dots, \mu_4\}$
 - Proportion of invariable sites $\mu = 0$

Drawbacks

- Consider positions independently
- **Ignores** functional constraints

Generative sequence model —► Propagator?

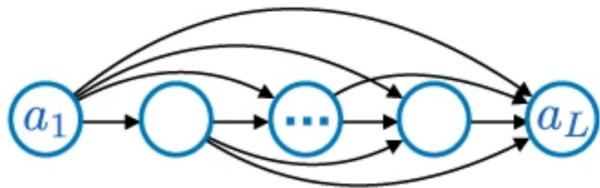
$$P(a_1, \dots, a_N) = \frac{1}{Z} \exp \left(\sum_{i,j=1}^L J_{ij}(a_i, a_j) + \sum_{i=1}^L h_i(a_i) \right)$$

P is not factorized

Evolution with autoregressive models

Autoregressive model

$$P(a_1 \dots a_L) = \prod_{i=1}^L p_i(a_i | a_1 \dots a_{i-1})$$



Need to know $p(a_i | a < i)$

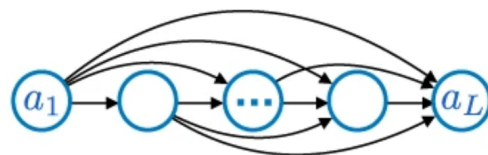
ArDCA $\longrightarrow p(a_i | a < i) \propto \exp \left(\sum_{j < i} J_{ij}(a_i, a_j) + h_i(a_i) \right)$

- Easy to infer (from alignment)
- Interpretable ($J_{ij} \sim$ contacts)
- Good generative properties

[Trinquier *et. al.*, *Nature Comm* 2021]

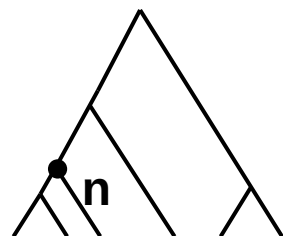
Evolution with autoregressive models

Autoregressive model



$$P(\mathbf{a}) = \prod_{i=1}^L p_i(a_i | a_{<i})$$

Given the context $a_{<i} = a_1 \dots a_{i-1}$
we know the equilibrium frequencies for a_i



Suppose we reconstructed $i-1$ positions

$$\rightarrow a_{<i}^n = a_1^n \dots a_{i-1}^n$$

→ Equilibrium frequencies $p_i(a_i | a_{<i}^n)$

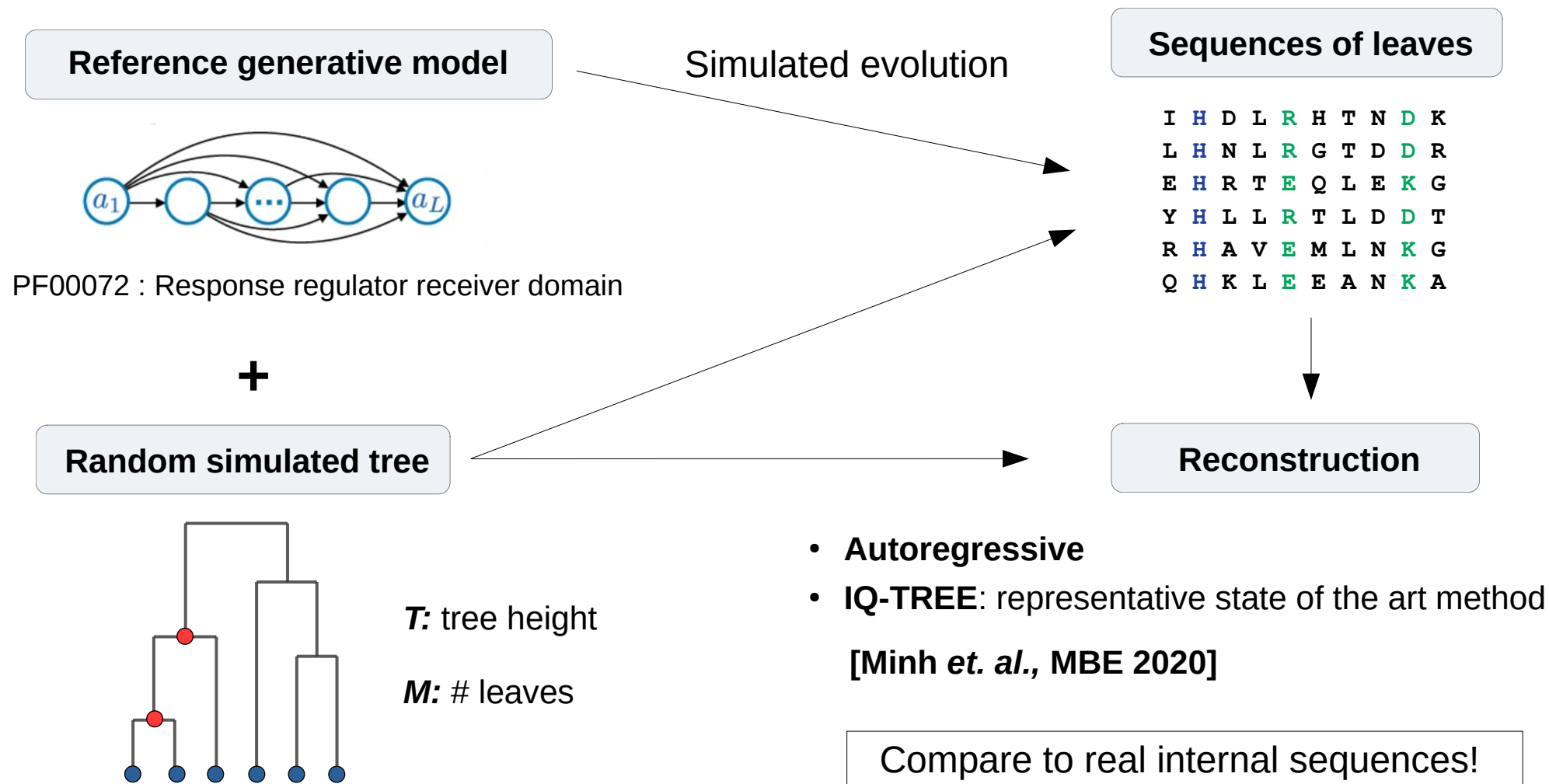
Evolution *towards* n for position i

Evolution model

- One position at a time: **almost factorized**
- Use knowledge of **functional constraints**
- **Cost**: need data to infer model

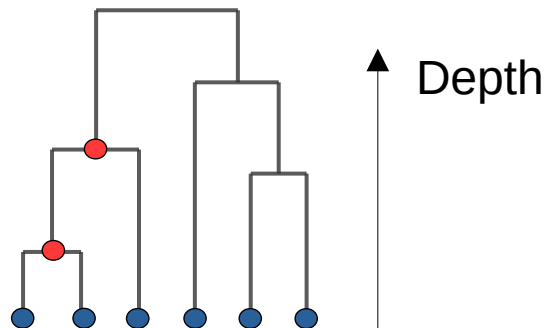
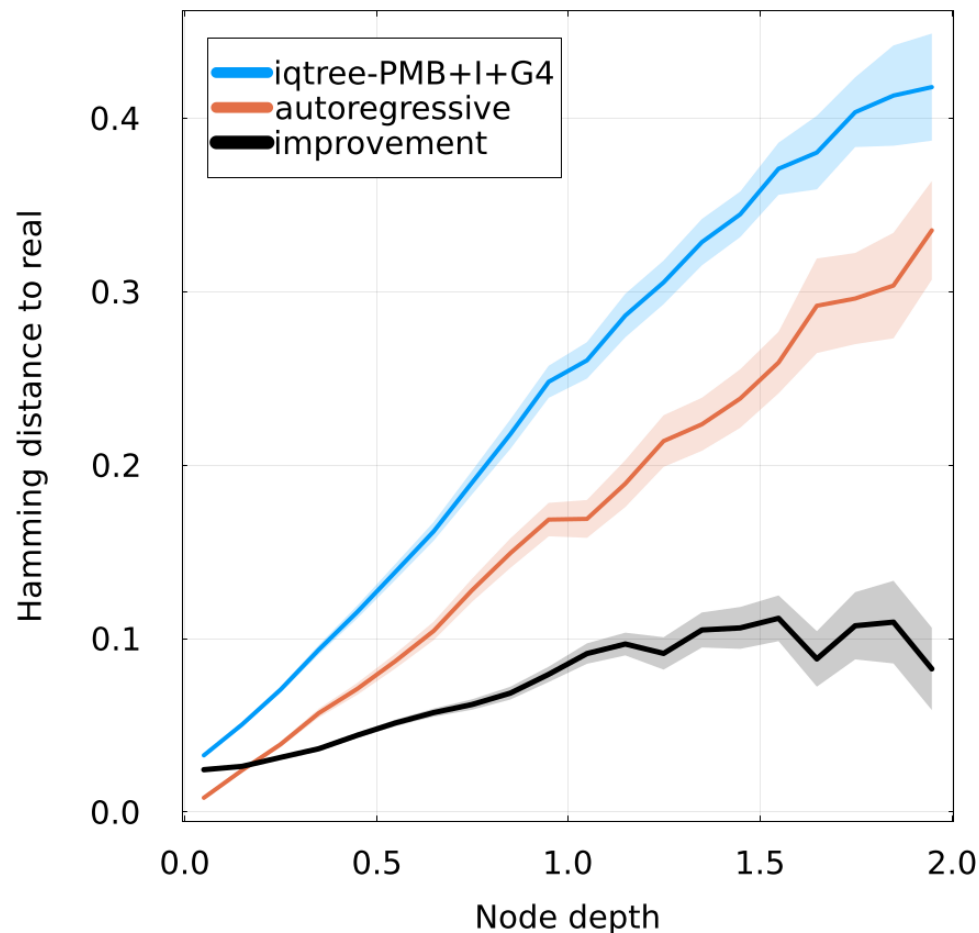
$$Q_i^{\rightarrow n} = \mathbf{H} \cdot \begin{pmatrix} p_i(1 | a_{<i}^n) & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & p_i(q | a_{<i}^n) \end{pmatrix}$$

Testing on simulations



Results: Maximum Likelihood reconstruction

$$\mathcal{L}_n(\mathbf{x}) = P(\mathcal{D}|n = \mathbf{x})$$

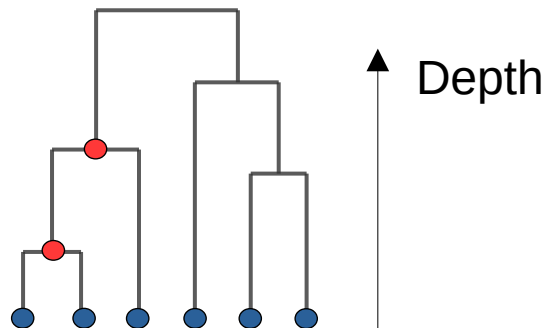
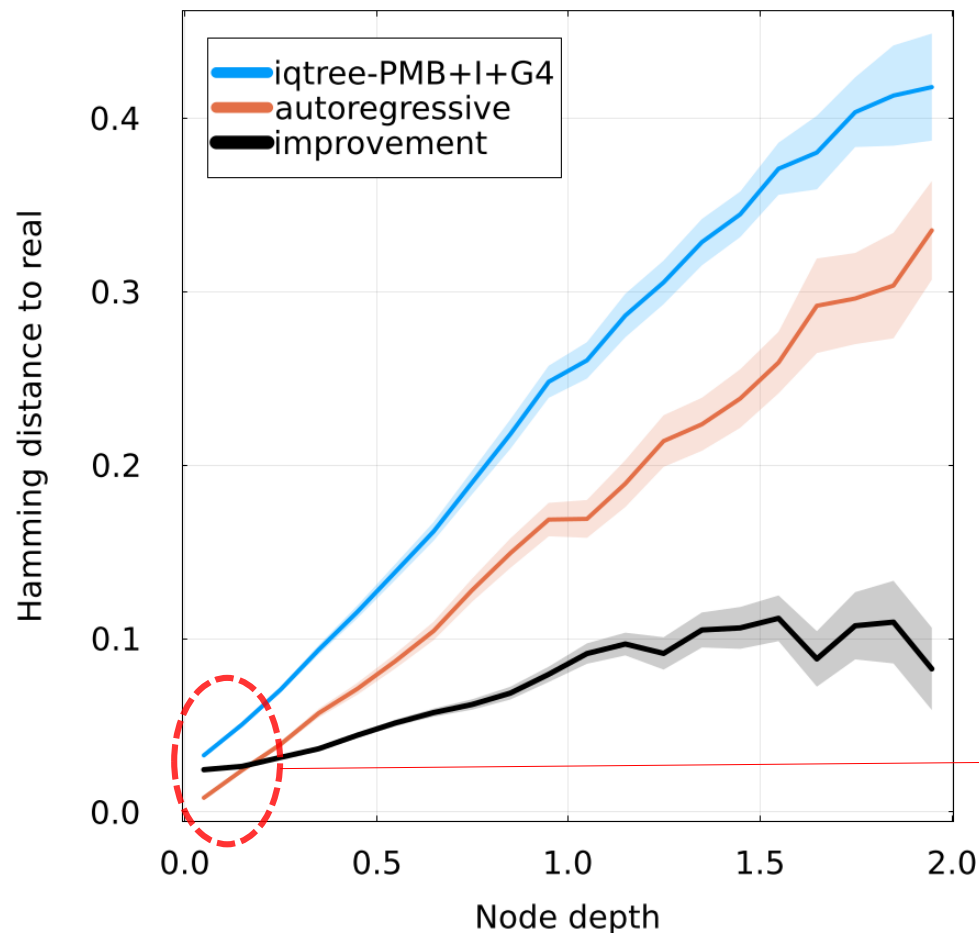


ML reconstruction $\mathbf{x} = \operatorname{argmax} \mathcal{L}_n$

Most commonly used in literature

Results: Maximum Likelihood reconstruction

$$\mathcal{L}_n(\mathbf{x}) = P(\mathcal{D}|n = \mathbf{x})$$



ML reconstruction $\mathbf{x} = \operatorname{argmax} \mathcal{L}_n$

Most commonly used in litterature

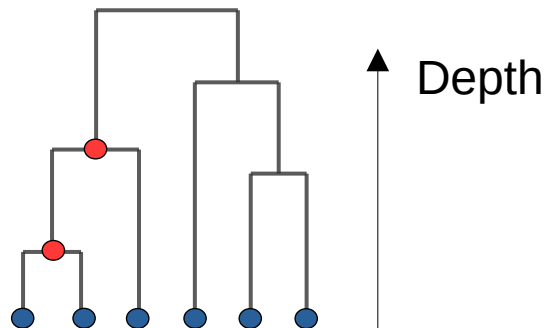
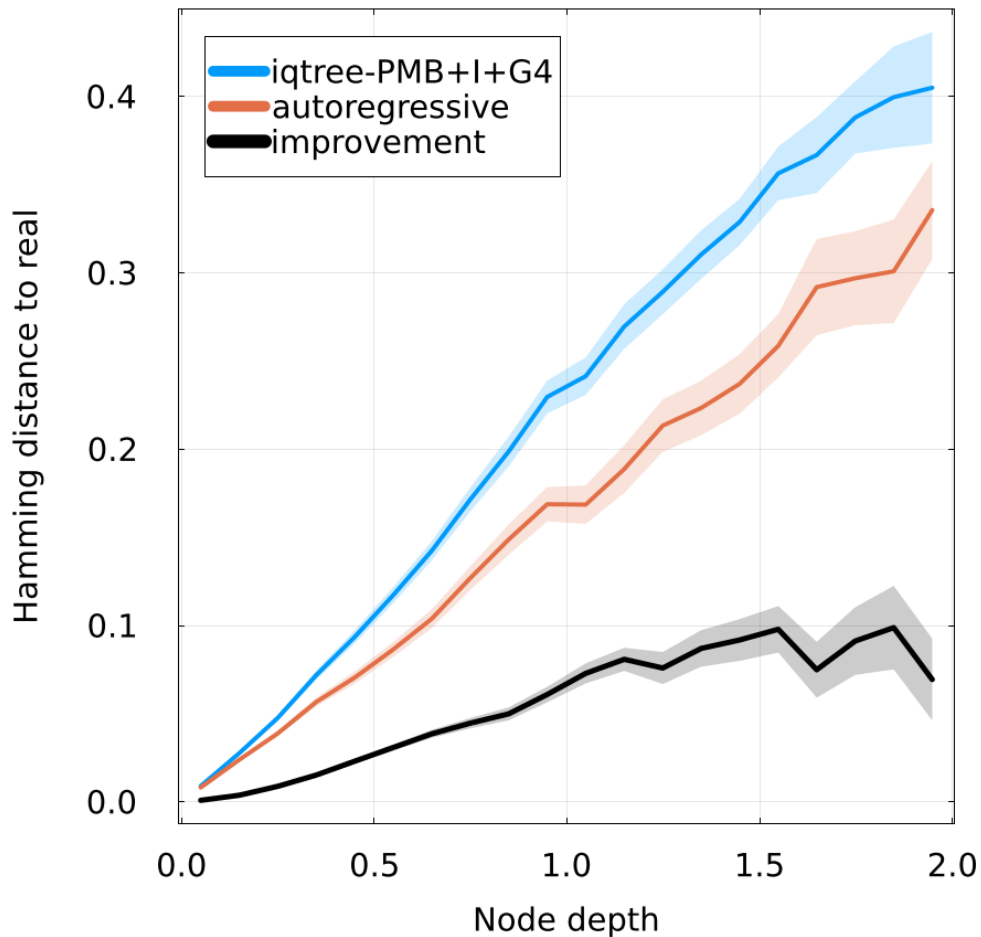
➔ Gaps!

- Represented in ArDCA...
- ... but not in IQ-TREE

Results: Maximum Likelihood reconstruction

$$\mathcal{L}_n(\mathbf{x}) = P(\mathcal{D}|n = \mathbf{x})$$

Hamming distance: ignoring gaps



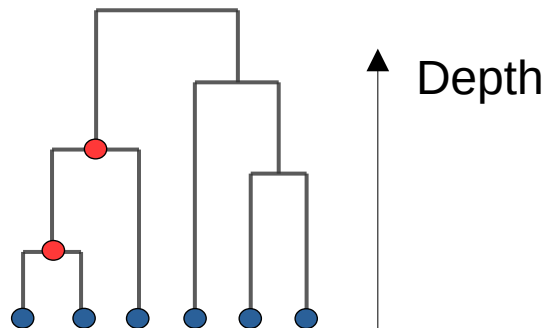
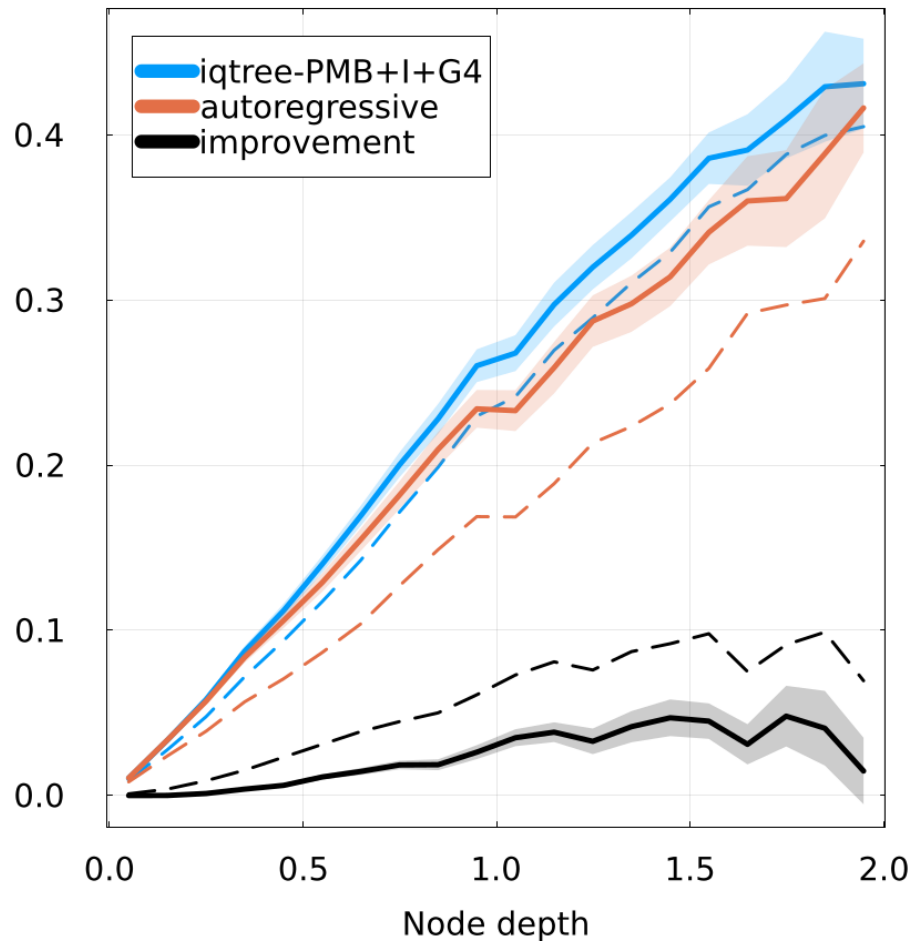
ML reconstruction $\mathbf{x} = \operatorname{argmax} \mathcal{L}_n$

Most commonly used in literature

Results: ML + Bayesian reconstruction

$$\mathcal{L}_n(\mathbf{x}) = P(\mathcal{D}|n = \mathbf{x})$$

Hamming distance: ignoring gaps



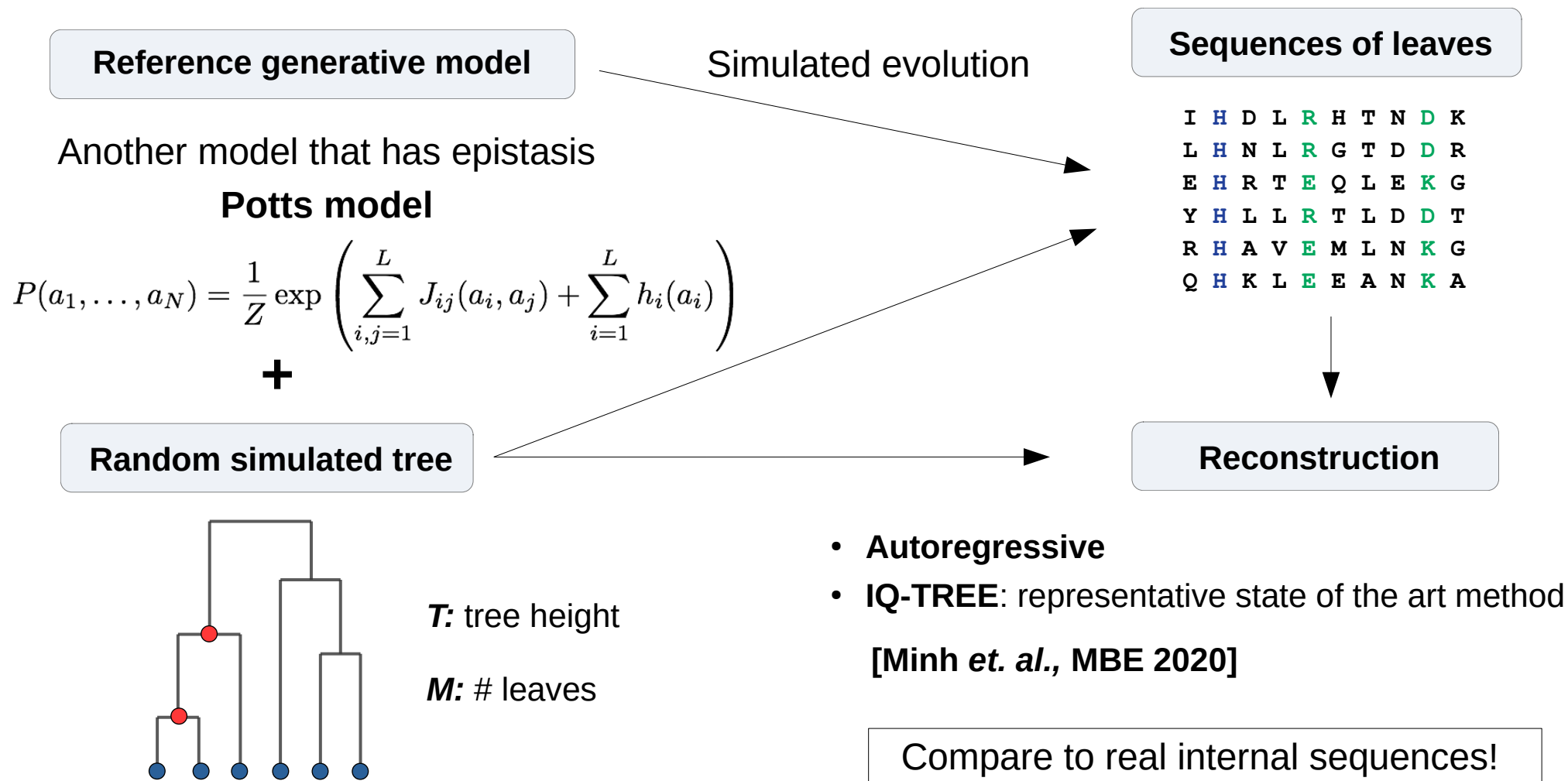
ML reconstruction $\mathbf{x} = \operatorname{argmax} \mathcal{L}_n$

Most commonly used in literature

Bayesian $P(\mathbf{x}) \propto \mathcal{L}_n(\mathbf{x})$

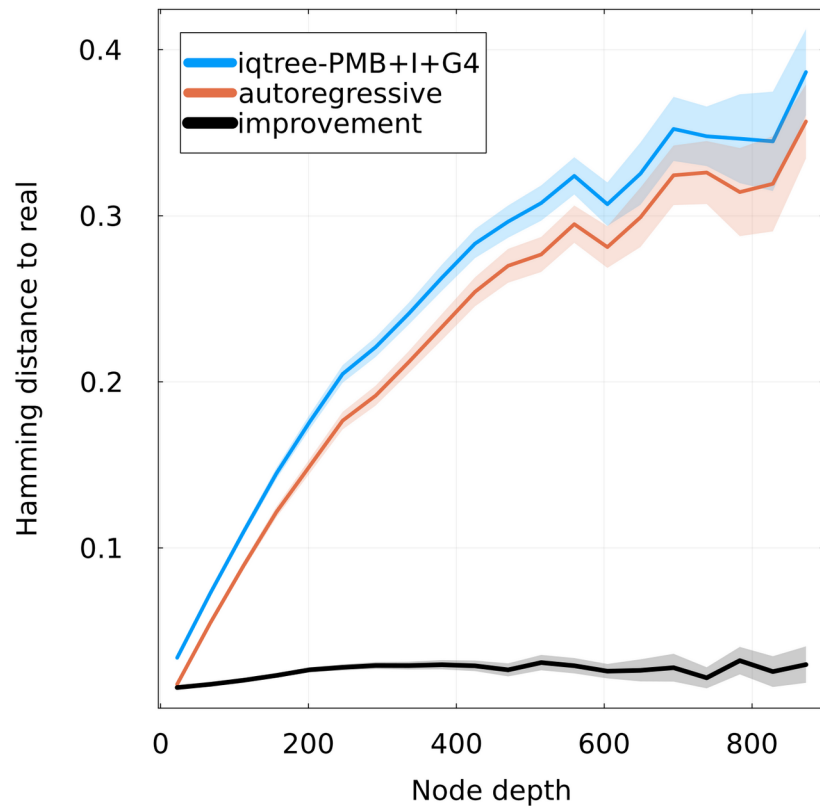
Rarely used in practice

Testing on simulations: with a different evolver?

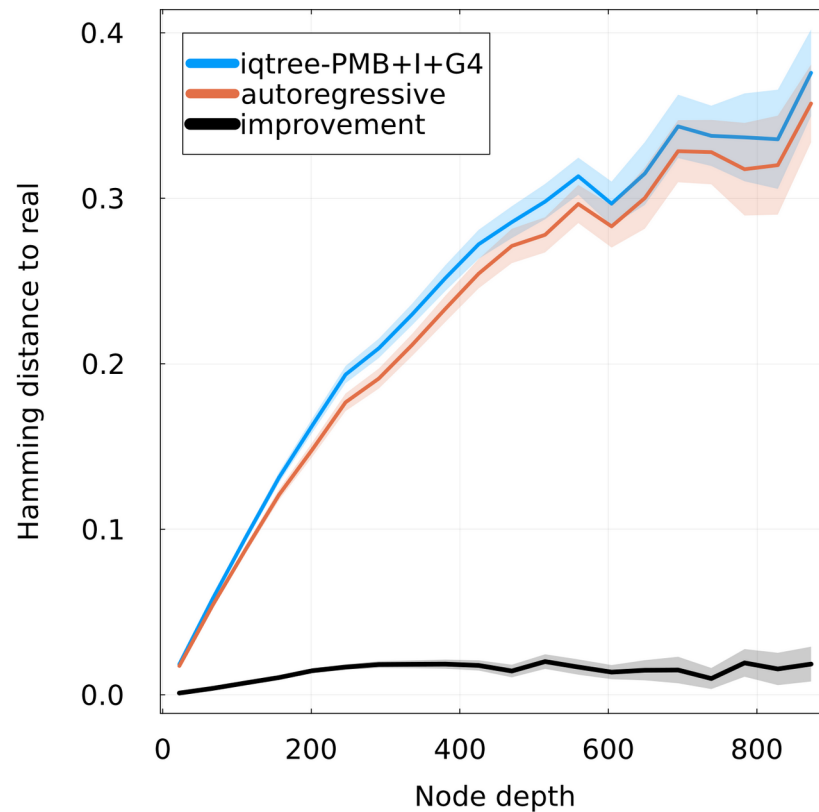


Testing on simulations: with a different evolver?

Hamming distance: with gaps

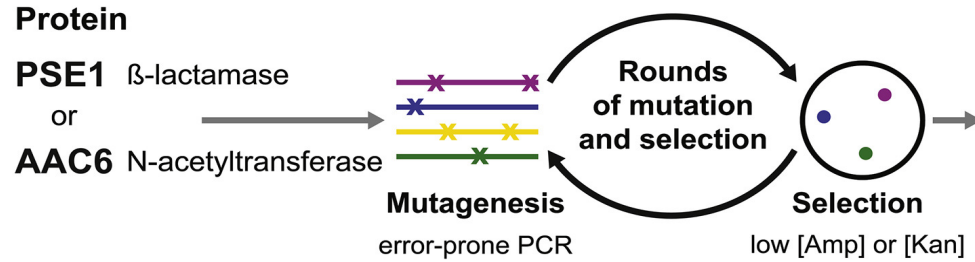


Hamming distance: ignoring gaps

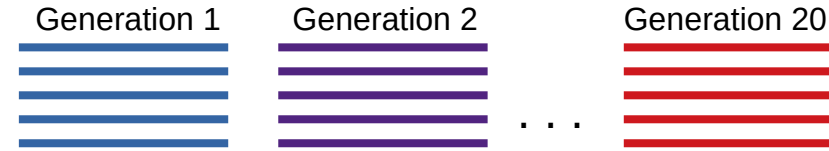


Systematic improvement, but small

Experimental data

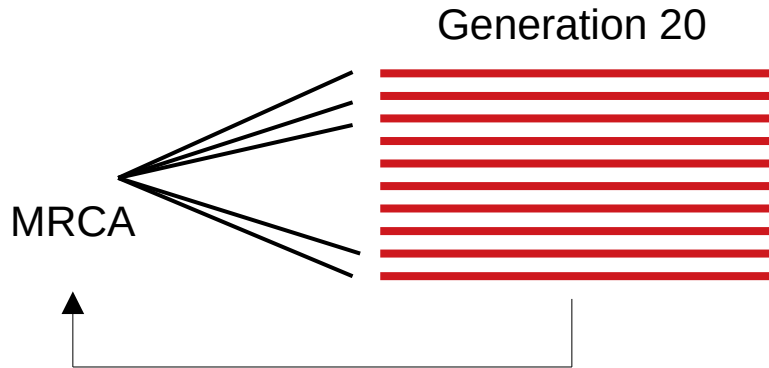


The sequence of MRCA is known!

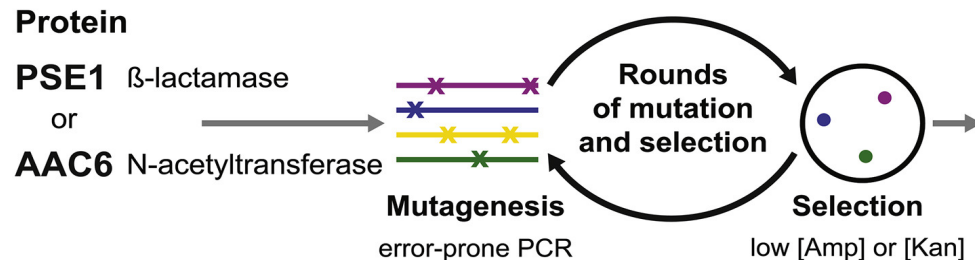


but not the tree...

→ Simplification: star tree

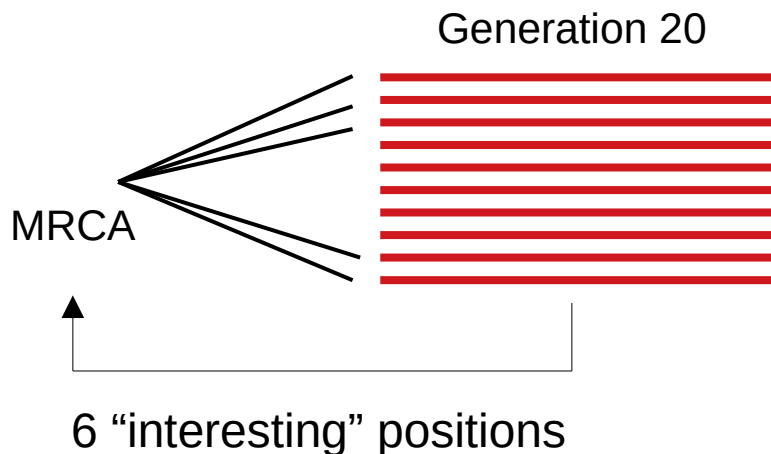


Experimental data

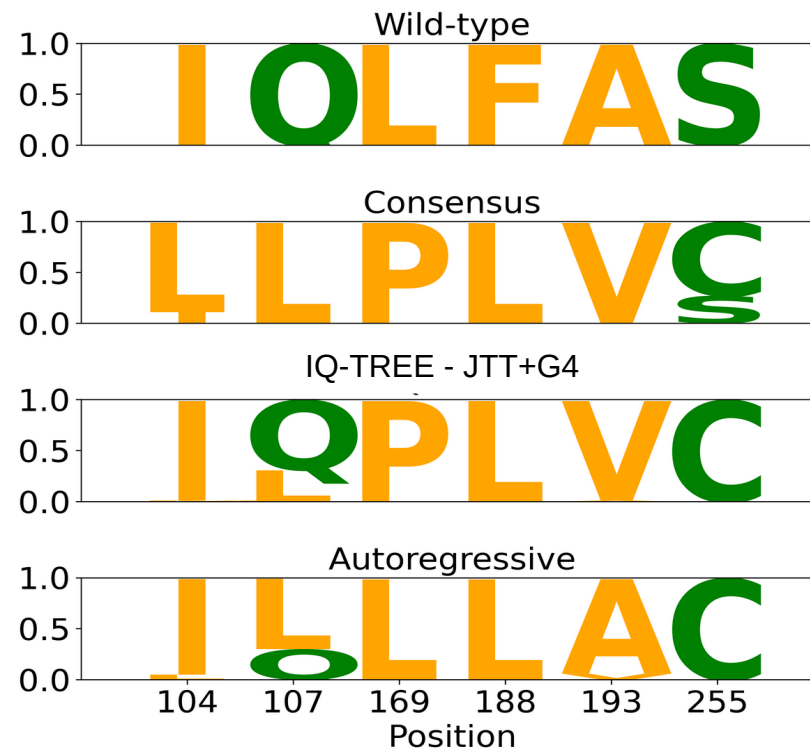
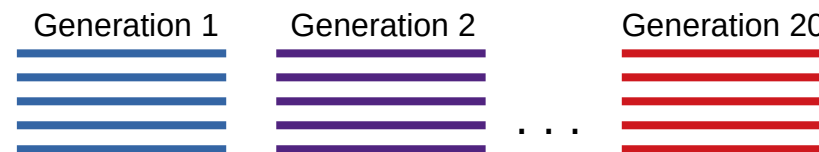


but not the tree...

→ Simplification: star tree



The sequence of MRCA is known!



Biases of the ML reconstruction

ML reconstruction $\mathbf{x} = \operatorname{argmax} \mathcal{L}_n$

- Most commonly used in literature
- Experimentally: functional & highly thermostable
- **Problem:** Is the best sequence representative?

—► **Biases**

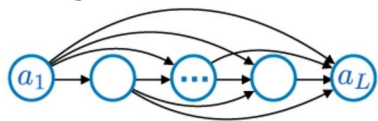
[Williams et. al., PLOS CB 2006]

Bayesian $P(\mathbf{x}) \propto \mathcal{L}_n(\mathbf{x})$

- Set of sequences at each internal node
- Rarely used in practice
- Sometimes non-functional
- More representative / Less subject to biases?

Biases of the ML reconstruction

Reference generative model



$$P(a_1 \dots a_L) = \prod_i P(a_i | a_1 \dots a_{i-1}) \longrightarrow$$

**Probability of
reconstructed
sequences?**

~proxy for function

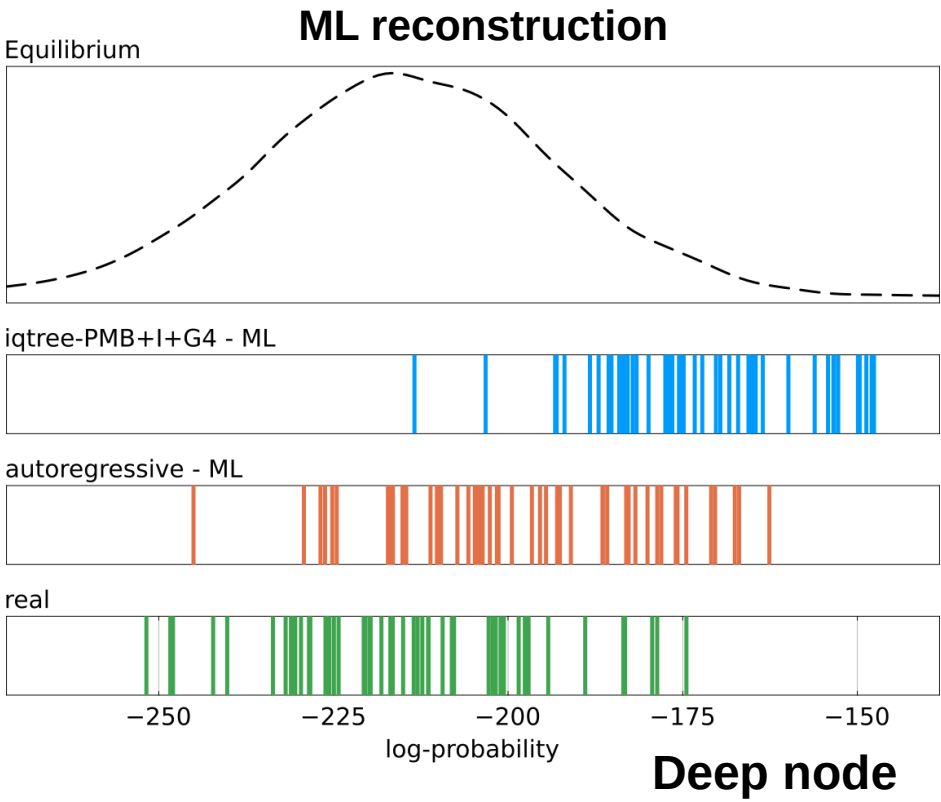
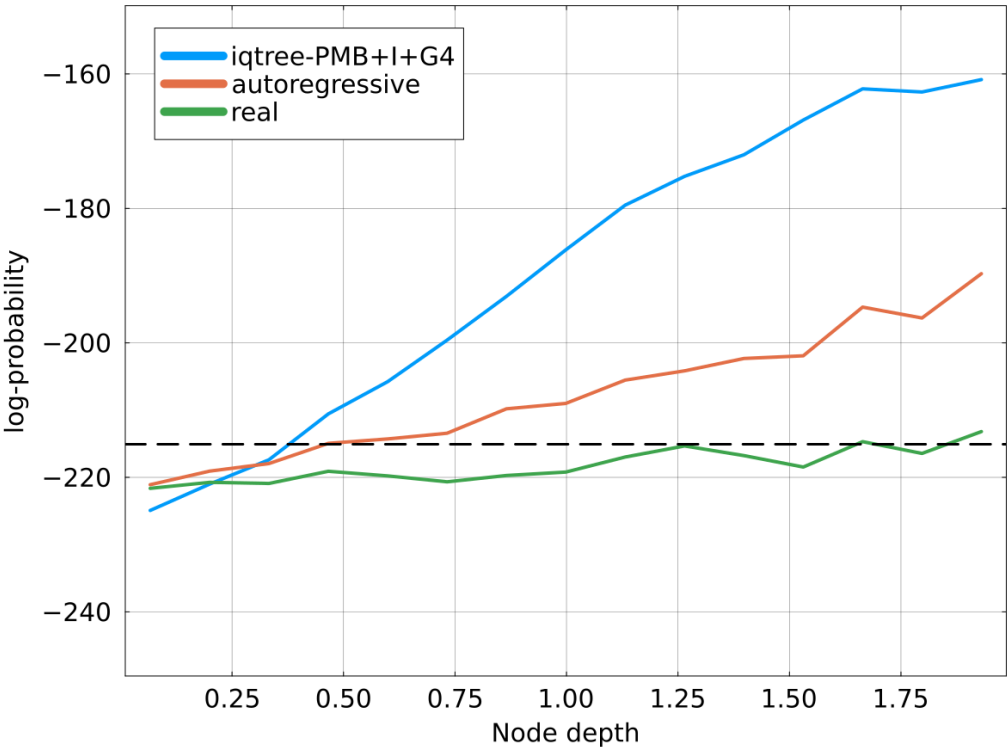
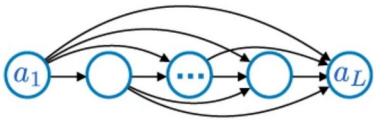
Biases of the ML reconstruction

Reference generative model

$$P(a_1 \dots a_L) = \prod_i P(a_i | a_1 \dots a_{i-1}) \longrightarrow$$

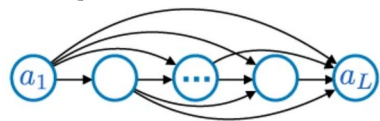
Probability of
reconstructed
sequences?

~proxy for function



Is Bayesian less biased?

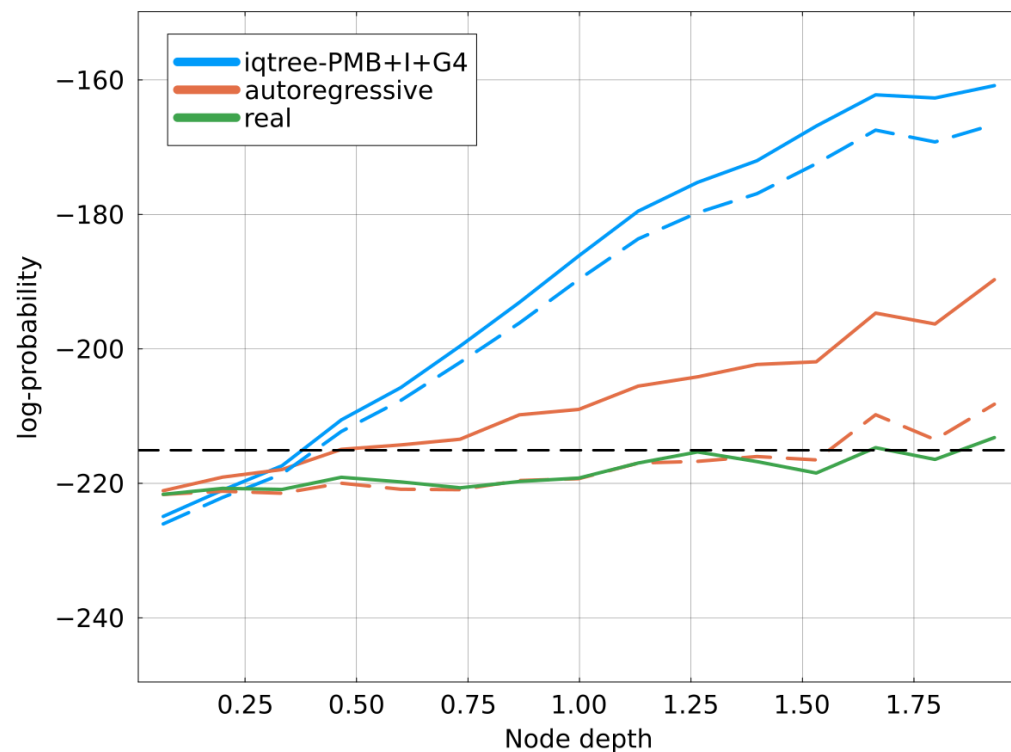
Reference generative model



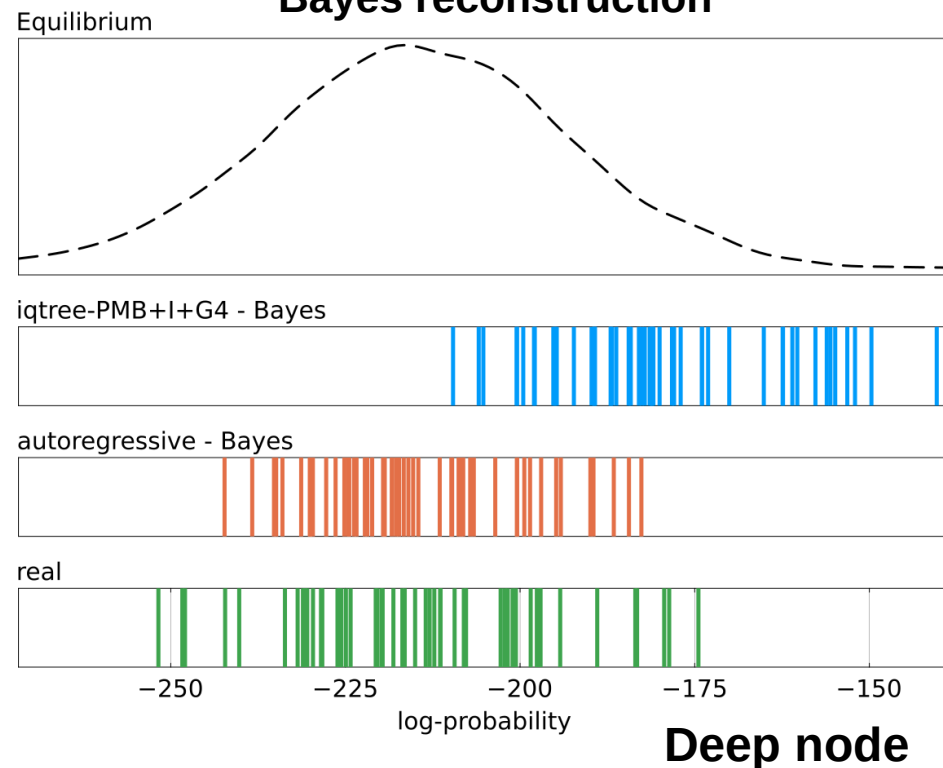
$$P(a_1 \dots a_L) = \prod_i P(a_i | a_1 \dots a_{i-1}) \longrightarrow$$

Probability of
reconstructed
sequences?

~proxy for function



Bayes reconstruction



Biases of the ML reconstruction

ML reconstruction $\mathbf{x} = \operatorname{argmax} \mathcal{L}_n$

- Most commonly used in literature
- Experimentally: functional & highly thermostable
- **Problem:** Is the best sequence representative?

—► **Biases**

[Williams et. al., PLOS CB 2006]

Bayesian $P(\mathbf{x}) \propto \mathcal{L}_n(\mathbf{x})$

- Set of sequences at each internal node
- Rarely used in practice
- Sometimes non-functional
- More representative / Less subject to biases?

For ML or site-independent model

—► **Bias in log-probability**

Biases of the ML reconstruction

ML reconstruction $\mathbf{x} = \operatorname{argmax} \mathcal{L}_n$

- Most commonly used in literature
- Experimentally: functional & highly thermostable
- **Problem:** Is the best sequence representative?

→ **Biases**

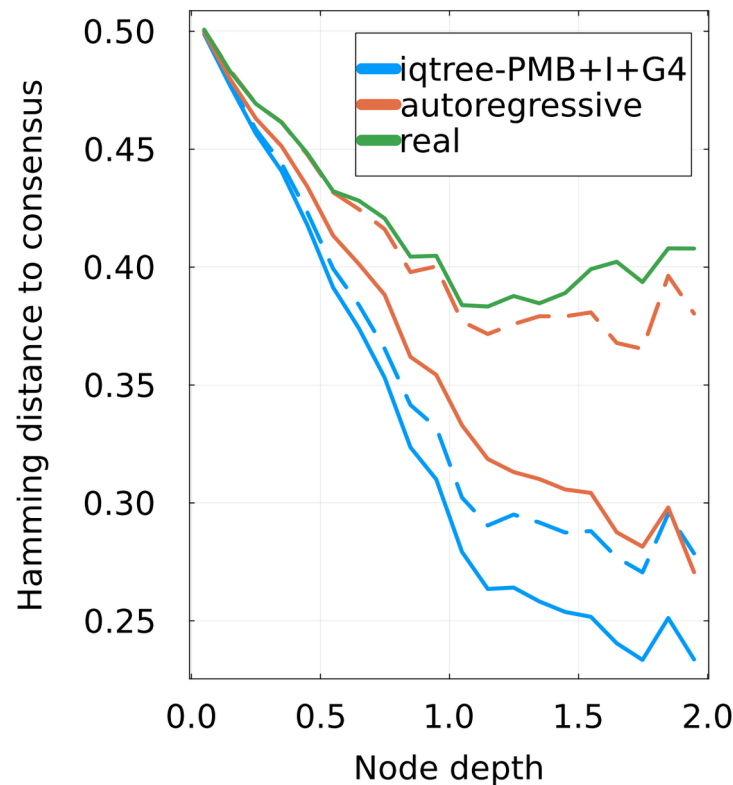
[Williams et. al., PLOS CB 2006]

Bayesian $P(\mathbf{x}) \propto \mathcal{L}_n(\mathbf{x})$

- Set of sequences at each internal node
- Rarely used in practice
- Sometimes non-functional
- More representative / Less subject to biases?

For ML or site-independent model

→ **Bias in log-probability**

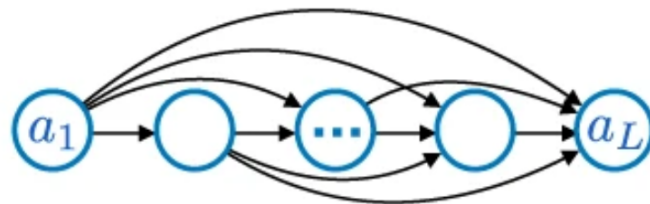


→ **Bias towards consensus!**

Autoregressive dynamics are not Markovian

Autoregressive model

$$P(a_1 \dots a_L) = \prod_i P(a_i | a_1 \dots a_{i-1})$$



Autoregressive dynamics

$$P_i(a|b, a_{<i}, t) = \underbrace{e^{-t} \delta_{ab}}_{\text{No mut.}} + \underbrace{(1 - e^{-t}) p_i(a|a_{<i})}_{>1 \text{ mut.}} \quad (\text{if } \mathbf{H} \text{ uniform})$$

$$P(\mathbf{a}|\mathbf{b}, t) = \prod_{i=1}^L P_i(a|b, a_{<i}, t) \xrightarrow[t \rightarrow \infty]{} P(\mathbf{a})$$

Autoregressive dynamics are not Markovian

Global balance:
$$P(\mathbf{x}) = \sum_y P(y)P(\mathbf{x}|y)$$

$\mathbf{x} = (- +) \longrightarrow P(\mathbf{x}) \ll 1$

$\mathbf{y} = (+ +) \longrightarrow P(\mathbf{y}) \sim 1/2$

Toy model
Binary sequences, $L=2$

Frequent

+ **+**

- **-**

Rare

- **+**

+ **-**

$$p_1(-) = 1/2$$

$$p_2(+|-) \ll 1$$

Autoregressive dynamics are not Markovian

Toy model
Binary sequences, L=2

Global balance: $P(\mathbf{x}) = \sum_{\mathbf{y}} P(\mathbf{y})P(\mathbf{x}|\mathbf{y})$

$\mathbf{x} = (- +) \longrightarrow P(\mathbf{x}) \ll 1$

$\mathbf{y} = (+ +) \longrightarrow P(\mathbf{y}) \sim 1/2$

Frequent
++
--
Rare
-+
+-

$p_1(-) = 1/2$

$p_2(+|-) \ll 1$

$$P(\mathbf{x}|\mathbf{y}, t) = \underbrace{(1 - e^{-t})p_1(-)}_{\text{Mut. at pos 1}} \cdot \left\{ \underbrace{e^{-t}}_{\text{No mut at 2}} + (1 - e^{-t}) \underbrace{p_2(+|-)}_{\ll 1} \right\} \longrightarrow \mathbf{O(1)}$$

\longrightarrow **Global balance cannot hold!**

Autoregressive dynamics are not Markovian

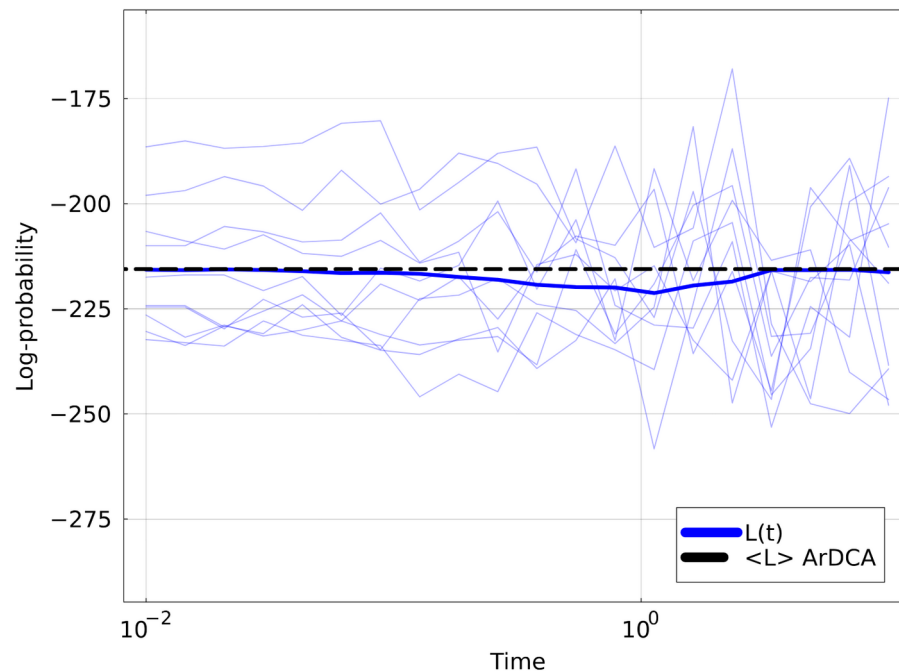
Consequences

- Irreversibility
- Dynamics go out of equilibrium
- Very likely not realistic



But...

- Correct dynamics for $t \ll 1$ and $t \gg 1$
- Quantitatively small effect



Conclusion

Evolution based on a generative model

- Autoregressive architecture: “almost factorized”
- Converges to generative distribution at long times
- Caveats: not Markov, irreversible

Ancestral Sequence Reconstruction

- Better handling of gaps
- Systematic improvement over state of the art models (simulations)
- Improvement on directed evolution data

Biases of maximum likelihood

- ML reconstruction is “too good” —► Bias to consensus!
- Bayesian reconstruction with good model: more representative!

Conclusion

Authors

Matteo De Leonardis (PoliTo)

Andrea Pagnani (PoliTo)

P.B.C

Evolution based on a generative model

- Autoregressive architecture: “almost factorized”
- Converges to generative distribution at long times
- Caveats: not Markov, irreversible

Ancestral Sequence Reconstruction

- Better handling of gaps
- Systematic improvement over state of the art models (simulations)
- Improvement on directed evolution data

Biases of maximum likelihood

- ML reconstruction is “too good” —► Bias to consensus!
- Bayesian reconstruction with good model: more representative!

Thank you for listening



**Politecnico
di Torino**