

MDTree: A Masked Dynamic Autoregressive Model for Phylogenetic Inference

Anonymous authors

Paper under double-blind review

Abstract

Phylogenetic tree inference requires optimizing both branch lengths and topologies, yet traditional MCMC-based methods suffer from slow convergence and high computational cost. Recent deep learning approaches improve scalability but remain constrained: Bayesian models are computationally intensive, autoregressive methods depend on fixed species orders, and flow-based models underutilize genomic signals. Fixed-order autoregression introduces an inductive bias misaligned with evolutionary proximity: early misplacements distort subsequent attachment probabilities and compound topology errors (exposure bias). Absent sequence-informed priors, the posterior over the super-exponential topology space remains diffuse and multimodal, yielding high-variance gradients and sluggish convergence for both MCMC proposals and neural samplers. We propose MDTree, a masked dynamic autoregressive framework that integrates genomic priors into a Dynamic Ordering Network to learn biologically informed node sequences. A dynamic masking mechanism further enables parallel node insertion, improving efficiency without sacrificing accuracy. Experiments on standard benchmarks demonstrate that MDTree outperforms existing methods in accuracy and runtime while producing biologically coherent phylogenies, providing a scalable solution for large-scale evolutionary analysis.

1 Introduction

Phylogenetic trees are fundamental for revealing evolutionary relationships, enabling lineage tracing from common ancestors to present-day organisms using DNA or protein sequences (Brocchieri, 2001; Munjal et al., 2019). They underpin studies in taxonomy, evolutionary biology, and medicine, offering insights into species origins, biodiversity, and the evolutionary trajectories of pathogens and cancer cells (Hugenholtz et al., 2021; Tummers & Green, 2022). Accurate and efficient inference has high practical value: in pathogen source tracing, it supports timely outbreak interventions (Biek et al., 2015); in cancer evolution analysis (Fimereli et al., 2022), it reveals clonal architecture and treatment resistance; and in biodiversity conservation (Theissinger et al., 2023), it enables large-scale, automated species relationship reconstruction. These applications underscore both the scientific and societal significance of phylogenetic modeling. Yet, the surge of genomic data and the combinatorial growth of tree topologies pose major computational challenges, calling for scalable and accurate new methods.

Traditional statistical frameworks, notably Maximum Likelihood Estimation (MLE) (Izquierdo-Carrasco et al., 2011; Solís-Lemus & Ané, 2016) and Bayesian Inference (BI) via Markov Chain Monte Carlo (MCMC) (Zhang et al., 2018; Wang et al., 2020), have long underpinned phylogenetic inference. Yet, with increasing taxa, they encounter severe computational bottlenecks: the space of unrooted bifurcating topologies grows super-exponentially as $(2N - 5)!!$, while the joint optimization of continuous branch lengths and discrete topologies further compounds complexity.

Leveraging deep learning, breakthroughs in phylogenetic inference have burst onto the scene, addressing long-standing computational challenges in the field (Nesterenko et al., 2022; Smith & Hahn, 2023; Tang et al., 2024). Research efforts primarily follow two main directions: representation learning on known tree structures and generative models. The former, exemplified by VBPI-GNN (Zhang, 2023), optimizes perfor-

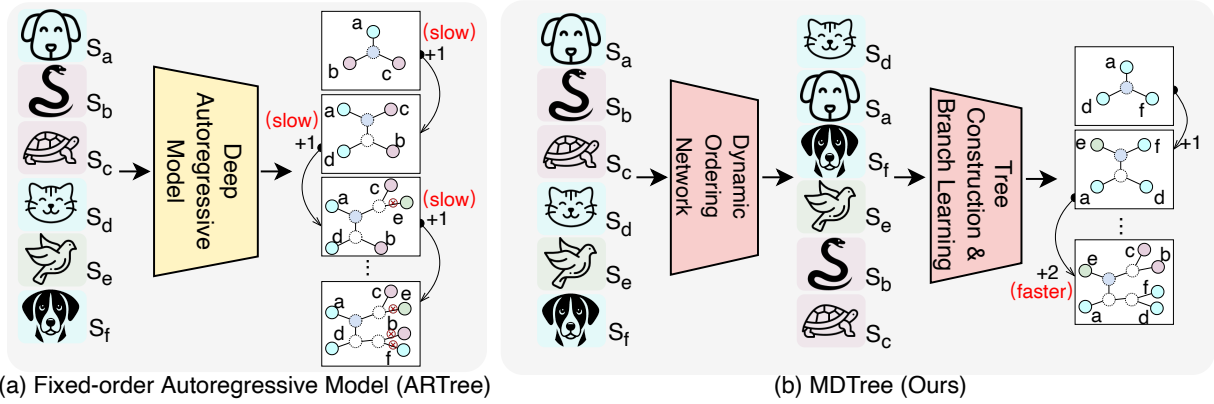


Figure 1: **Comparison between fixed-order autoregressive generation (ARTree) and our dynamic-order method (MDTree).** (a) Fixed-order ARTree adds species sequentially according to a predefined order, ignoring biological priors. This leads to suboptimal intermediate structures and slower generation, as only one leaf node is added at each step. (b) MDTree (ours) employs a Dynamic Ordering Network to determine a biologically-informed insertion order based on genomic features, enabling related species to be clustered earlier. The Tree Construction & Branch Learning module further supports parallel insertion of multiple nodes, achieving faster generation and more phylogenetically consistent topologies.

mance based on predefined topologies but struggles when the topology is unknown and both topology and branch lengths must be inferred. These methods also underutilize evolutionary information from biological sequences, impacting accuracy and flexibility (Penny, 2004). On the other hand, generative models, which infer tree structures directly from data, can be further divided into three types: Bayesian generative models (e.g., Geophy (Mimori & Hamada, 2024)) leverage probabilistic frameworks to capture uncertainty but are computationally intensive; autoregressive models (e.g., ARTree (Xie & Zhang, 2024)) sequentially add nodes, offering flexibility yet relying on predefined orders that overlook true evolutionary relationships, while their stepwise nature leads to inefficiency for large datasets (Razavi et al., 2019). Lastly, Generative Flow Networks (GFNs) (e.g., PhyloGFN (Zhou et al., 2024)) provide greater flexibility by exploring multimodal posterior distributions but still struggle to fully integrate evolutionary signals, impacting the accuracy of inferred trees. Therefore, few methods have achieved these goals simultaneously.

To overcome these limitations, we focus on a core question: *how can biological priors effectively guide node addition to improve phylogenetic inference accuracy?* As shown in Fig. 2, classical autoregressive methods (Fig. 2a) rely on fixed orders (e.g., lexicographical), overlooking evolutionary relationships and often producing inaccurate trees [1]. Our method (Fig. 2b) learns evolutionarily meaningful node orders, ensuring species like reptiles, birds, and mammals are added in line with their ancestry. This improves the accuracy and biological relevance of generated trees by prioritizing species with closer common ancestors.

Specifically, we adopt a **dynamic autoregressive generation paradigm**, where both the order of node additions and their insertion positions are learned from genomic sequences instead of being fixed in advance. This paradigm is instantiated by our *Masked Dynamic Autoregressive Model* (MDTree), which integrates a Diffusion Ordering

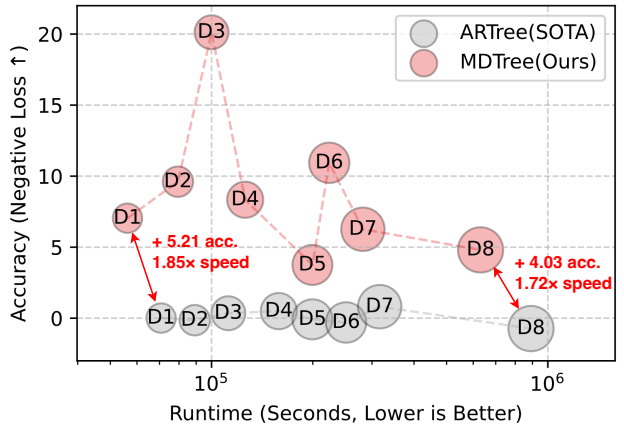


Figure 2: **Runtime and node count comparison between MDTree and ARTree.** Evaluation is conducted on eight benchmarks under two optimization settings (log-scale).

Network (DON) to learn biologically informed orders directly from sequence data via an absorbing diffusion model (Bond-Taylor et al., 2021), mitigating the limitations of fixed or random orders. By combining the strengths of Graph Neural Networks and Language Models (LMs), MDTree captures intricate genomic relationships while modeling complex tree structures. A Dynamic Masking Mechanism enables parallel node processing, improving efficiency. Lastly, we employ a dual-pass tree traversal strategy for branch length estimation and use the LAX model (Grathwohl et al., 2017) to reduce variance in discrete sampling for stabilizing optimization and enhancing convergence.

Experiments on phylogenetic benchmarks show that MDTree outperforms existing methods in accuracy and efficiency. Empirical analysis of Angiosperms353 (Zuntini et al., 2024) further demonstrates its ability to recover evolutionary lineages, including Rosaceae and Moraceae, suggesting broader biological applications. In summary, our contributions are:

- **A novel dynamic autoregressive generation paradigm for phylogenetic inference:** We leverage a generation strategy that dynamically learns node order and insertion positions from genomic sequence data, improving the accuracy and biological relevance of inferred trees.
- **An innovative methodology:** We propose MDTree, which integrates a Diffusion Ordering Network for biologically informed node orders, combines genomic Language Models with dual-pass traversal for precise tree generation, and employs a dynamic masking mechanism for efficient parallel processing.
- **Strong experimental validation:** Comprehensive experiments validate that MDTree achieves state-of-the-art performance. Visualizations from real-world Angiosperm datasets further confirm the biological relevance and interpretability of the generated trees.

2 Related Works

Phylogenetic inference methods are generally categorized into traditional and deep learning-based approaches; each is further divided into graph structure generation and representation models. For details on background, please see Appendix A.

Traditional Methods rely on predefined evolutionary models and statistical inference. *Graph Structure Generation Models:* MrBayes (Ronquist et al., 2012) utilizes Bayesian inference to generate trees but struggles with high-dimensional combinatorial spaces, requiring large sample sizes for accuracy. VaiPhy (Kop-tagel et al., 2022) combines SLANTIS sampling strategy (Diaconis, 2019) with biological models (e.g., JC model (Munro, 2012)) to estimate branch lengths and generate accurate tree structures. *Graph Structure Representation Models:* SBN (Zhang & Matsen IV, 2018a) models the probability distribution of tree topologies from existing trees, focusing on subplit relationships without directly estimating branch lengths. VBPI (Zhang & Matsen IV, 2018b) extends SBNs to estimate posterior distributions and optimize branch lengths through variational inference.

Deep Learning-based Methods offer more flexible and scalable solutions. *Graph Structure Generation Models:* (1) Bayesian Generative Models like GeoPhy (Mimori & Hamada, 2024) learn latent tree representations to generate diverse topologies. (2) Autoregressive Models such as ARTree (Xie & Zhang, 2024) sequentially generate trees, well-suited for hierarchical data. (3) Generative Flow Networks like PhyloGFN (Zhou et al., 2024) optimize tree generation paths using Markov decision processes. *Graph Structure Representation Models:* VBPI-GNN (Zhang, 2023) combines SBNs with variational inference to optimize topology and branch lengths.

3 Methods

A unified model that can handle both tasks must therefore (i) capture biologically meaningful topological structures (Tree Topology Density Estimation, TDE), and (ii) accurately estimate continuous evolutionary distances (Variational Bayesian Phylogenetic Inference, VBPI), while being robust to limited or no supervision on the topology. Unless otherwise specified, all vector representations (e.g., h_i) are treated as column vectors.

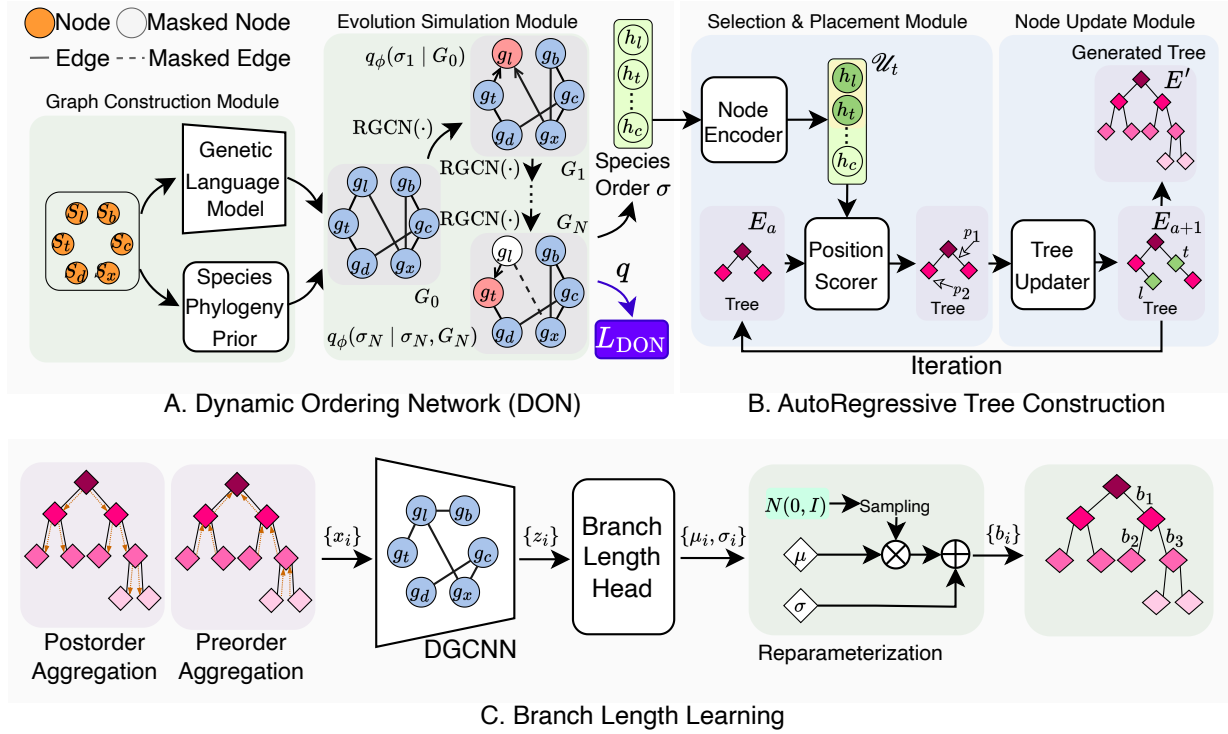


Figure 3: **Framework of MDTree for dynamic autoregressive tree generation.** A. The Dynamic Ordering Network module utilizes a pre-trained enomic LM to extract embeddings from sequences Y , guiding nodes into absorbing states in an autoregressive manner as determined by DON $q_\phi(\sigma|G)$. B. The Autoregressive Tree Construction module employs a parallel strategy to add multiple leaf and internal nodes simultaneously at specified positions based on the order provided by DON. C. The Branch Length Learning module optimizes branch lengths through a dual-pass traversal.

Formulation. Given a set of N species sequences $\mathcal{S} = \{s_i\}_{i=1}^N$ and their corresponding genomic representations $\mathcal{G} = \{g_i\}_{i=1}^N$ extracted from a pretrained Genomic Language Model (e.g., DNABERT2 (Zhou et al., 2023)), we aim to infer a phylogenetic tree that captures both its discrete topology and continuous evolutionary distances. Formally, the tree is modeled as an unrooted binary graph $\mathcal{G}_T = (\mathcal{V}_T, \mathcal{E}_T)$, where \mathcal{V}_T denotes the set of taxa and internal nodes, and \mathcal{E}_T the set of undirected edges between them. Each edge $e \in \mathcal{E}_T$ is associated with a branch length $b_e \in \mathbb{R}_+$, and we denote the set of branch lengths as $B_\tau = \{b_e : e \in \tau\}$. Our objective is to learn a mapping,

$$\mathcal{F} : \mathcal{S} \longrightarrow \{(e, b_e) \mid e \in \tau\}, \quad (1)$$

which jointly specifies the topology τ (a binary tree structure over \mathcal{V}_T) and its associated branch lengths B_τ .

This formulation naturally covers both tasks, (i) For TDE, we focus on estimating the topology τ by marginalizing out the branch lengths B_τ in the likelihood $p(\mathcal{A} \mid \tau)$, enabling evaluation against reference topologies. (ii) For VBPI, we jointly infer both τ and B_τ under the posterior $p(\tau, B_\tau \mid \mathcal{S})$, using amortized inference to model their dependencies. In practice, MDTree runs the full pipeline for VBPI, while for TDE the branch-length refinement module is bypassed during inference and only used implicitly when marginalizing over B_τ .

To address the limitations of fixed node orders in prior autoregressive models (Xie & Zhang, 2024), we propose **MDTree** (Fig. 3), which dynamically learns biologically informed node addition orders and insertion positions from \mathcal{G} via a Dynamic Ordering Network (DON) based on an absorbing diffusion process (Austin et al., 2021). The learned order and contextual node embeddings jointly guide an autoregressive tree construction module with dynamic masking for parallel insertion. Finally, a dual-pass traversal refines branch

lengths using both global and local structural cues. This unified pipeline enables accurate reconstruction of phylogenetic structure and evolutionary distances, supporting both TDE and VBPI tasks.

Framework of MDTree. Fig. 3 illustrates the overall architecture of MDTree, comprising three tightly coupled components. The DON first produces biologically informed node orders and contextual embeddings from genomic inputs. These are then consumed by the autoregressive tree construction module, which determines node insertion positions under a dynamic masking schedule. Finally, the dual-pass traversal refines branch lengths using bidirectional context. The interaction between these components allows MDTree to support both TDE and VBPI within the same inference pipeline.

3.1 DON for Learning Biologically Informed Node Orders with Genomic Priors

The order in which species nodes are added to a phylogenetic tree significantly impacts the inferred topology, especially under the constraint of binary unrooted trees. From a biological perspective, species with closer ancestry should be introduced earlier in the tree construction process to better preserve evolutionary semantics (Penny, 2004; Gregory, 2008). While some recent works show robustness to taxa orderings (Xie & Zhang, 2024), the benefit of *learning* biologically informed node orders remains underexplored. Such an approach could allow the model to adaptively exploit genomic signals to produce topologies that better reflect true evolutionary relationships.

Graph Construction Module Given an input set of sequences \mathcal{S} , each $s_i \in \mathcal{S}$ is first encoded into a genomic embedding $g_i \in \mathbb{R}^d$ using a pretrained genomic language model (e.g., DNABERT2 (Zhou et al., 2023)), injecting biological priors into subsequent ordering decisions. An initial graph $\mathcal{G}_0 = (\mathcal{V}, \mathcal{E})$ is then constructed based on sequence similarity or known homology (species phylogeny prior), serving as the structural backbone for contextual reasoning. Node features are passed through a relational graph convolutional network (RGCN) (Schlichtkrull et al., 2018):

$$h_i^{(0)} = \text{RGCN}(g_i + \text{PE}(g_i), \mathcal{E}), \quad (2)$$

where $\text{PE}(\cdot)$ is a positional encoding and $\{h_i^{(0)}\}_{i=1}^N$ are the initial contextualized embeddings that integrate both sequence-level and local graph information.

Evolution Simulation Module Starting from \mathcal{G}_0 , DON simulates an iterative absorption process to determine the biologically informed node order $\sigma = \{i_1, i_2, \dots, i_N\}$.

(1) *Node selection probability:* At each step t , the model computes the probability of absorbing each active node:

$$q(i_t \mid \mathbf{H}^{(t-1)}) = \text{Cat} \left(\frac{h_{i_t}^{(t)} Q_t^\top \odot h_{i_t}^{(0)} \bar{Q}_{t-1}^\top}{h_{i_t}^{(0)} \bar{Q}_t h_{i_t}^{(t)\top}} \right), \quad (3)$$

where \odot is element-wise multiplication, $\bar{Q}_{t-1} = \prod_{i=1}^{t-1} Q_i$ is the cumulative transition matrix from previous steps, and Q_t is the current-step transition matrix.

(2) *Transition dynamics:* The discrete-time transition matrix $Q_t \in \mathbb{R}^{(N+1) \times (N+1)}$ controls state changes:

$$[Q_t]_{ij} = \begin{cases} 1 & \text{if } i = j = m \\ 1 - \beta_{t,i} & \text{if } i = j \neq m \\ \beta_{t,i} & \text{if } j = m, i \neq m \\ 0 & \text{otherwise,} \end{cases} \quad (4)$$

where $m = N + 1$ is the absorbing (masked) state and $\beta_{t,i} \in [0, 1]$ is a monotonically increasing absorption probability.

(3) *Greedy absorption:* The node with the highest selection probability is chosen $i_t^* = \arg \max q(i_t \mid \cdot)$ absorbed into σ , and the graph is updated $G_{t+1} = G_t \cup \{i_t^*\}$. This loop continues until all nodes are absorbed.

Ordering Supervision Module We supervise the Dynamic Ordering Network (DON) by aligning its predicted node absorption order $\hat{\sigma}$ with a reference order $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_N)$ obtained from external phylogenetic inference tools (e.g., MrBayes). At each step t , DON predicts a probability distribution $q(i_t | \mathcal{G}_0, \sigma_{<t})$ over candidate nodes $i_t \in \mathcal{V}_T \setminus \sigma_{<t}$ conditioned on the initial graph \mathcal{G}_0 and the previously absorbed nodes $\sigma_{<t} = (\sigma_1, \dots, \sigma_{t-1})$. The supervision objective is the negative log-likelihood (NLL) of the ground-truth sequence:

$$\mathcal{L}_{\text{DON}} = - \sum_{t=1}^N \log q(i_t = \sigma_t | \mathcal{G}_0, \sigma_{<t}), \quad (5)$$

which encourages DON to assign high probability to biologically consistent construction sequences.

3.2 AutoRegressive Tree Construction with Dynamic Node Insertion

Given the biologically informed node ordering $\sigma = \{i_1, i_2, \dots, i_N\}$ and the corresponding contextual node embeddings $\{h_i\}$ from the DON module, the autoregressive tree construction stage determines *where* each node should be inserted in the growing phylogenetic tree. The embeddings are directly passed from DON without re-encoding, while the order σ is kept fixed and used in two ways: (i) as a positional bias in the attention-based placement network to favor evolutionarily prioritized nodes, and (ii) as a priority score to adjust insertion probabilities. This tight coupling ensures that the biologically meaningful ordering learned in DON directly shapes the resulting topology, while maintaining computational efficiency for large N .

Selection & Placement Module Instead of inserting one node at a time—which is computationally expensive and prone to error propagation—we employ a *dynamic parallel insertion* strategy. At autoregressive step t , let \mathcal{V}_t be the set of placed nodes and $\mathcal{R}_t = \mathcal{V} \setminus \mathcal{V}_t$ the remaining nodes. A subset $\mathcal{U}_t \subseteq \mathcal{R}_t$ is selected for insertion according to a cosine mask rate:

$$\mathcal{U}_t = \text{SelectNodes}(\mathcal{R}_t, \rho_t), \quad \rho_t = \frac{1}{2} \left(1 + \cos \left(\frac{t}{T} \pi \right) \right), \quad (6)$$

where T is the total number of autoregressive steps. Early steps insert fewer nodes for accuracy; later steps insert more for speed.

For each $i \in \mathcal{U}_t$, let $p[i]$ be a candidate parent in the partial tree. We compute a relational embedding via Multi-Head Attention (Vaswani, 2017):

$$\mathbf{r}_i = \text{MHA}(Q, h_i, h_i), \quad (7)$$

where h_i is the contextual embedding from DON, and $\text{MHA}(Q, h_i, h_i)$ is the self attention mechanism that captures interactions between h_i and all other node embeddings, $Q \in \mathbb{R}^{(N-3) \times d}$ probes all potential attachment points.

Given \mathbf{r}_i and $\mathbf{r}_{p[i]}$ (candidate parent’s embedding), the insertion probability is predicted as:

$$\mathbf{L}_i = \text{softmax}(\text{MLP}(\text{Concat}(\mathbf{r}_i, \text{MAX}(\mathbf{r}_i, \mathbf{r}_{p[i]})) + \text{PE}(t))), \quad (8)$$

where $\text{MAX}(\cdot)$ extracts dominant shared features, and $\text{PE}(t)$ injects temporal information. Evolutionary priority from σ is enforced by biasing toward lower Rank_i , $\mathbf{L}_{\text{adj},i} = \mathbf{L}_i + \alpha \cdot (N - \text{Rank}_i)$, with bias strength α . The final position is sampled as, $\text{pos}_i \sim \text{Multinomial}(\text{softmax}(\mathbf{L}_{\text{adj},i}))$, where pos_i is the position in the current partial tree where node i will be inserted.

Node Update Module Let $\mathcal{E}_T^{(t)}$ be the edge set at step t . After sampling $\{\text{pos}_i\}$, we update the topology:

$$\mathcal{E}_T^{(t+1)} = \mathcal{E}_T^{(t)} \cup \{(v_i, v_{p[i]}), (v_{p[i]}, v_i)\}. \quad (9)$$

If insertion creates a new internal node j (degree two children), its embedding is initialized by averaging neighbors:

$$r_j^{\text{internal}} = \frac{1}{\text{Card}(\mathcal{N}(j))} \sum_{k \in \mathcal{N}(j)} r_k, \quad (10)$$

where $\text{Card}(\mathcal{N}(j))$ is the cardinality of the neighbor set of j . This process repeats until all nodes are placed, producing the final binary unrooted topology τ .

3.3 Dual-Pass Traversal for Branch Length Learning

The branch length learning module (Fig. 3) jointly captures global and local structural cues via a dual-pass traversal, followed by graph-based encoding and differentiable branch length sampling.

Postorder and Preorder Aggregation Given a rooted tree converted from the inferred unrooted topology τ by adding a dummy edge, a post-order aggregation propagates information bottom-up:

$$h_u^{\text{fwd}} = \text{GRU} \left(h_u^{\text{init}}, \frac{1}{|\mathcal{C}(u)|} \sum_{v \in \mathcal{C}(u)} \phi(h_v^{\text{fwd}}, \ell_{uv}) \right), \quad (11)$$

where h_u^{init} is the genomic embedding from DON, $\mathcal{C}(u)$ is the child set of node u , ℓ_{uv} is the current branch length estimate, and $\phi(\cdot)$ is an MLP conditioned on ℓ_{uv} .

A subsequent pre-order aggregation propagates refined context top-down, updating each child node v :

$$h_v^{\text{bwd}} = \psi(h_u^{\text{bwd}}, h_v^{\text{fwd}}, \text{PE}(\text{depth}(u, v))), \quad (12)$$

where $\psi(\cdot)$ is an MLP and PE encodes relative depth.

DGCNN Encoding The bidirectionally aggregated features $\{x_i\}$, obtained by combining forward and backward states, are passed to a Dynamic Graph Convolutional Neural Network (DGCNN) to capture higher-order relational dependencies, yielding refined node representations $\{z_i\}$.

Branch Length Head and Reparameterization A lightweight branch length head maps each z_i to Gaussian parameters (μ_i, σ_i) . Differentiable sampling is enabled via the reparameterization trick:

$$b_i = \mu_i + \sigma_i \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \quad (13)$$

where b_i is the sampled branch length. This ensures low-variance gradient estimates and flexible distribution modeling.

Branch Length Loss Branch length estimation is supervised by the negative log-likelihood of the sequence alignment \mathcal{A} under a continuous-time Markov chain (CTMC) substitution model (Yang, 1994):

$$\mathcal{L}_{\text{len}} = - \sum_{c=1}^{|\mathcal{A}|} \log p_{\text{CTMC}}(\mathcal{A}_c \mid \tau, \ell, \Theta_{\text{sub}}), \quad (14)$$

where \mathcal{A}_c is the c -th alignment column and Θ_{sub} are substitution parameters (e.g., GTR rates).

This dual-pass design explicitly decouples global structural encoding from local refinement, while the reparameterized probabilistic formulation enables stable branch length estimation and improved alignment with observed evolutionary signals.

3.4 MDTree Inference for Joint Topology and Branch Length Estimation

Building on the biologically informed node ordering loss \mathcal{L}_{DON} from Sec. 3.1 and the branch length likelihood \mathcal{L}_{len} from Sec. 3.3, the final stage integrates these components into a unified inference framework. **MDTree** couples *Tree Topology Density Estimation* (TDE) and *Variational Bayesian Phylogenetic Inference* (VBPI) to jointly optimize the posterior distribution over tree topology τ and branch lengths ℓ given the sequence alignment \mathcal{A} .

Table 1: Research Questions (RQs) and their corresponding sub-questions.

RQ1: Performance	How well does MDTree perform in generating tree topologies (TDE) and inferring branch lengths (VBPI)?
RQ2: Time Efficiency	How efficient is MDTree in reducing runtime?
RQ3: Tree Quality	How optimal is MDTree to generate a tree structure? (RQ3-1) How diverse are the tree topologies generated by MDTree? (RQ3-2) How consistent is the MDTree-generated tree compared to MrBayes? (RQ3-3)
RQ4: Module Impact	How does each MDTree’s module affect its performance? (RQ4-2) How do key hyper-parameters affect MDTree? (RQ4-2)
RQ5: Case Study	What evolutionary relationships between species does MDTree learn?

Tree Topology Density Estimation (TDE) Traditional phylogenetic methods often rely on fixed topologies or heuristic search strategies that may not adequately explore the vast space of possible tree structures. TDE addresses this limitation by learning a flexible distribution over topologies that can capture the uncertainty inherent in phylogenetic inference. This probabilistic approach is particularly valuable when dealing with closely related species or when sequence data contains conflicting evolutionary signals.

TDE refines $q_\theta(\tau)$ —parameterized by the autoregressive insertion process—by maximizing the marginal log-likelihood, effectively training the model to assign higher probability to topologies that better explain the observed sequence data:

$$\mathcal{L}_{\text{TDE}}(\theta) = \mathbb{E}_{q_\theta(\tau)} (\log p(\mathcal{A} \mid \tau)), \quad (15)$$

where $p(\mathcal{A} \mid \tau)$ is evaluated under a CTMC model and marginalized over branch lengths. This marginalization is crucial as it allows the topology learning to focus on structural relationships without being confounded by branch length estimation errors.

Variational Bayesian Phylogenetic Inference (VBPI) VBPI extends TDE by jointly modeling τ and ℓ via a structured posterior $q_\theta(\tau, \ell) = q_\theta(\tau) q_\phi(\ell \mid \tau)$, with $q_\phi(\ell \mid \tau)$ initialized from the dual-pass traversal outputs and refined through amortized inference. The ELBO objective is:

$$\mathcal{L}_{\text{VBPI}}(\theta, \phi) = \mathbb{E}_{q_\theta(\tau) q_\phi(\ell \mid \tau)} (\log p(\mathcal{A} \mid \tau, \ell)) - \text{KL}(q_\theta(\tau) q_\phi(\ell \mid \tau) \parallel p(\tau, \ell)). \quad (16)$$

Unified MDTree Objective Rather than optimizing each stage independently, MDTree integrates all component objectives into a unified training target that ensures ordering, topology construction, and branch length inference reinforce one another:

$$\mathcal{L}_{\text{MDTree}} = \lambda_{\text{DON}} \cdot \mathcal{L}_{\text{DON}} + \lambda_{\text{VBPI}} \cdot \mathcal{L}_{\text{VBPI}} + \lambda_{\text{len}} \cdot \mathcal{L}_{\text{len}} + \lambda_{\text{TDE}} \cdot \mathcal{L}_{\text{TDE}}, \quad (17)$$

where λ_{len} balances explicit branch length refinement with the joint inference objectives. This end-to-end coupling aligns all stages into a coherent optimization pipeline, yielding phylogenies that are both topologically accurate and metrically consistent with evolutionary signals.

4 Experiments

In this section, we demonstrate the effectiveness of our proposed MDTree in terms of the research questions in Table 1.

4.1 Experiment Setup

Evaluation Tasks and Datasets. We assess MDTree’s performance on two key tasks: TDE, which focuses on optimizing tree topologies with MLL metric, and VBPI, where tree topologies and branch lengths are

Table 2: **Comparison of KL divergence (\downarrow) across eight benchmark datasets with different methods.** **Boldface** for the highest result, Underline for the second highest result of traditional methods.

Methods	Dataset (#Taxa,#Sites)	DS1 (27,1949)	DS2 (29,2520)	DS3 (36,1812)	DS4 (41,1137)	DS5 (50,378)	DS6 (50,1133)	DS7 (59,1824)	DS8 (64,1008)
	Sampled Trees	1228	7	43	828	33752	35407	1125	3067
	GT Trees	2784	42	351	11505	1516877	809765	11525	82162
MCMC-based	SBN	<u>0.0707</u>	0.0144	<u>0.0554</u>	0.0739	1.2472	0.3795	0.1531	0.3173
	SRF	0.0155	<u>0.0122</u>	0.3539	0.5322	11.5746	10.0159	1.2765	2.1653
	CCD	0.6027	0.0218	0.2074	0.1952	1.3272	0.4526	0.3292	0.4149
	SBN-SA	0.0687	0.0218	0.2074	0.1952	1.3272	0.4526	0.3292	0.4149
	SBN-EM	0.0136	0.0199	0.1243	0.0763	0.8599	0.3016	0.0483	0.1415
	SBN-EM- α	0.0130	0.0128	0.0882	<u>0.0637</u>	<u>0.8218</u>	<u>0.2786</u>	<u>0.0399</u>	<u>0.1236</u>
Structure Generation	ARTree	<u>0.0045</u>	0.0097	0.0548	0.0299	0.6266	0.2360	0.0191	0.0741
	MDTree	0.0036	0.0129	0.0446	0.0216	0.5751	0.1591	0.0169	0.0634

jointly inferred, using ELBO and MLL. These evaluations span eight diverse benchmark datasets, covering various organisms like marine animals, plants, bacteria, fungi, and eukaryotes, as outlined in Appendix C.

Baselines. MDTree is compared against three primary groups of baselines: (1) MCMC-based methods (e.g., MrBayes, SBN), (2) Structure Representation methods (VBPI, VBPI-GNN), which leverage pre-generated topologies, and (3) Structure Generation methods for Bayesian inference without pre-selected topologies. Notably, ARTree, a comparable autoregressive method like ours, is highlighted for comparison. All training details and hyperparameters are provided in Appendix E.

4.2 Comparison Results on Benchmarks (RQ1)

Table 3: **Evaluation of MLL (\uparrow) on eight benchmark datasets.** VBPI and VBPI-GNN utilize pre-generated tree topologies during training, making **direct comparisons challenging**. **Boldface** highlights the highest result, **Text** denotes the second highest of structure generation methods, and **Text** indicates the second highest of MCMC-based methods.

Methods	Dataset (#Taxa,#Sites)	DS1 (27,1949)	DS2 (29,2520)	DS3 (36,1812)	DS4 (41,1137)	DS5 (50,378)	DS6 (50,1133)	DS7 (59,1824)	DS8 (64,1008)
MCMC-based	MrBayes	-7108.42 (0.18)	<u>-26367.57</u> (0.48)	-33735.44 (0.50)	-13330.44 (0.54)	<u>-8214.51</u> (0.28)	<u>-6724.07</u> (0.86)	-37332.76 (2.42)	<u>-8649.88</u> (1.75)
	SBN	<u>-7108.41</u> (0.15)	-26367.71 (0.08)	<u>-33735.09</u> (0.09)	<u>-13329.94</u> (0.20)	-8214.62 (0.40)	-6724.37 (0.43)	<u>-37331.97</u> (0.28)	-8650.64 (0.50)
Structure Representation	VBPI	-7108.42 (0.10)	-26367.72 (0.12)	-33735.10 (0.11)	-13329.94 (0.31)	-8214.61 (0.67)	-6724.34 (0.68)	-37332.03 (0.43)	-8650.63 (0.55)
	VBPI-GNN	-7108.41 (0.14)	-26367.73 (0.07)	-33735.12 (0.09)	-13329.94 (0.19)	-8214.64 (0.38)	-6724.37 (0.40)	-37332.04 (0.12)	-8650.65 (0.45)
Structure Generation	ARTree	<u>-7108.41</u> (0.19)	<u>-26367.71</u> (0.07)	<u>-33735.09</u> (0.09)	<u>-13329.94</u> (0.17)	<u>-8214.59</u> (0.34)	<u>-6724.37</u> (0.46)	<u>-37331.95</u> (0.27)	<u>-8650.61</u> (0.48)
	phi-CSMC	-7290.36 (7.23)	-30568.49 (31.34)	-33798.06 (6.62)	-13582.24 (35.08)	-8367.51 (8.87)	-7013.83 (16.99)	NA	-9209.18 (18.03)
	GeoPhy	-7111.55 (0.07)	-26379.48 (11.60)	-33757.79 (8.07)	-13342.71 (1.61)	-8240.87 (9.80)	-6735.14 (2.64)	-37377.86 (29.48)	-8663.51 (6.85)
	GeoPhy LOO(3)	-7116.09 (10.67)	-26368.54 (0.12)	-33735.85 (0.12)	-13337.42 (1.32)	-8233.89 (6.63)	-6735.9 (1.13)	-37358.96 (13.06)	-8660.48 (0.78)
	PhyloGFN	-7108.95 (0.06)	-26368.90 (0.28)	-33735.60 (0.35)	-13331.83 (0.19)	-8215.15 (0.20)	-6730.68 (0.54)	-37359.96 (1.14)	-8654.76 (0.19)
	Ours	-7101.38 (0.07)	-26357.96 (0.06)	-33715.31 (0.10)	-13322.10 (1.34)	-8210.76 (0.23)	-6713.13 (0.32)	-37326.50 (1.39)	-8645.07 (0.69)

The TDE Task. We compare the KL divergence to measure the difference between the model’s generated tree topology distribution $q_\theta(\tau)$ and the true posterior $p(\tau)$: $\text{KL}(p(\tau)||q_\theta(\tau)) = \sum_\tau p(\tau) \log \frac{p(\tau)}{q_\theta(\tau)}$. Table 2 shows that our MDTree consistently achieves lower KL divergence across all datasets compared to MCMC-based and structure generation methods. On complex datasets such as DS5 and DS6, it outperforms ARTree

Table 4: **Comparison of mean log-likelihood (MLL) and runtime between ARTree and MDTree under RWS and VIMCO optimization, each trained for 400,000 iterations.** MDTree consistently achieves higher MLL and reduces runtime by over 40% compared to ARTree.

Methods	MLL	Runtime (s)
ARTree_rws	-7107.74	128.7
MDTree_rws	-7103.71	75.0(↓41.72%)
ARTree_vimco	-7106.59	114.7
MDTree_vimco	-7101.38	63.7(↓44.46%)

and SBN, demonstrating superior scalability. Even on smaller datasets like DS1 and DS3, the performance remains competitive, highlighting the model’s robustness. The comparison with ARTree underscores the advantage of autoregressive models, including ours, particularly on larger, more complex datasets.

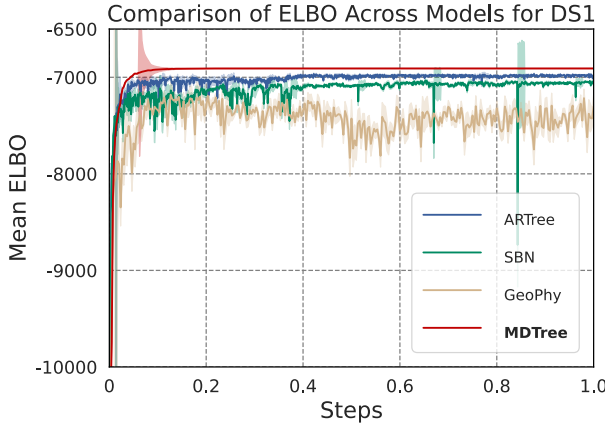


Figure 4: Comparison of ELBO.

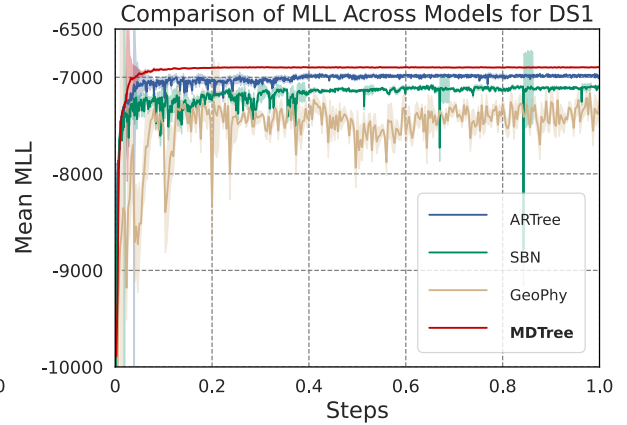


Figure 5: Comparison of MLL.

The VBPI Task. We evaluate the VBPI task using ELBO and MLL metrics. Since direct computation of MLL is intractable, it is approximated via importance sampling. Unlike TDE, which relies on known tree topologies, VBPI evaluates the fit between model-generated tree topologies and branch lengths and the observed gene sequence data. As shown in Table 3 and Table 5, Tree Structure Generation methods exhibit broader applicability in MLL and ELBO metrics compared to Structure Representation methods, which are restricted by their reliance on pre-generated topologies. Our method, MDTree, consistently achieves the highest metrics across all datasets, highlighting its enhanced capacity to approximate the posterior distribution of tree topologies and branch lengths. Fig. 4 shows MDTree’s superior stability and fast convergence in ELBO on DS1, outperforming baselines. ARTree and SBN improve later but with fluctuations, while GeoPhy performs the worst with consistently low and unstable values. Fig. 5 highlights MDTree’s advantages in MLL, quickly reaching and maintaining high scores, whereas ARTree, SBN, and especially GeoPhy lag behind.

4.3 Runtime Reduction and Efficiency Evaluation (RQ2)

MDTree demonstrates substantial runtime efficiency across all datasets, outperforming ARTree consistently. Both runtime and the number of nodes are log-transformed on the vertical axes, with solid and dashed lines representing the RWS and VIMCO optimization techniques. MDTree achieves faster than ARTree across all datasets, with VIMCO providing further reductions, especially for MDTree-VIMCO, which exhibits the lowest runtime. The efficiency of MDTree becomes even more apparent as dataset complexity increases. Table 4 confirms this finding, with MDTree reducing runtime by 41.72% (RWS) and 44.46% (VIMCO) compared to ARTree while maintaining superior MLL metrics. This underscores MDTree’s efficiency and scalability, particularly with VIMCO optimization.

Table 5: **Evaluation of ELBO (\uparrow) on eight datasets.** Higher values indicate better performance. Results for GeoPhy were not reported in its original publication and are reproduced by us. Light gray marks the best baseline result, and darker gray marks the best overall result. Our method consistently achieves the highest ELBO across all datasets.

Methods	Dataset (#Taxa,#Sites)	DS1 (27,1949)	DS2 (29,2520)	DS3 (36,1812)	DS4 (41,1137)	DS5 (50,378)	DS6 (50,1133)	DS7 (59,1824)	DS8 (64,1008)
MCMC-based	SBN	-7110.24 (0.03)	-26368.88 (0.03)	-33736.22 (0.02)	-13331.83 (0.02)	-8217.80 (0.04)	-6728.65 (0.04)	-37334.85 (0.03)	-8655.05 (0.04)
Structure Generation	ARTree	<u>-7110.09</u> (0.04)	<u>-26368.78</u> (0.07)	<u>-33735.25</u> (0.08)	<u>-13330.27</u> (0.05)	<u>-8215.34</u> (0.04)	<u>-6725.33</u> (0.06)	<u>-37332.54</u> (0.13)	<u>-8651.73</u> (0.05)
	GeoPhy	-7116.67 (1.71)	-26434.84 (0.10)	-33766.72 (0.15)	-13389.36 (3.45)	-8220.91 (2.64)	-6769.41 (3.25)	-37882.96 (1.97)	-8654.39 (0.97)
	Ours	-7005.98 (0.06)	-26362.75 (0.12)	-33430.94 (0.34)	-13113.03 (3.65)	-8053.23 (2.57)	-6324.90 (1.26)	-36838.42 (1.99)	-8409.06 (1.09)

Table 6: **Topological comparison of three tree diversity metrics.** Higher values of Simpson’s Diversity Index and the number of topologies accounting for the top 95% cumulative frequency indicate better diversity. In contrast, a **lower frequency** of the most frequent topology reflects a balanced distribution.

Dataset	Statistics	MrBayes	ARTree	Ours
DS1	Diversity Index (\uparrow)	0.87	0.89	0.99
	Top Frequency (\downarrow)	0.27	0.1	0.007
	Top 95% Frequency (\uparrow)	42	10	121
DS2	Diversity Index (\uparrow)	0.89	0.96	0.99
	Top Frequency (\downarrow)	0.27	0.43	0.13
	Top 95% Frequency (\uparrow)	208	203	301
DS3	Diversity Index (\uparrow)	0.98	0.89	0.90
	Top Frequency (\downarrow)	0.02	0.01	0.004
	Top 95% Frequency (\uparrow)	753	509	1146
DS4	Diversity Index (\uparrow)	0.86	0.89	0.99
	Top Frequency (\downarrow)	0.11	0.05	0.002
	Top 95% Frequency (\uparrow)	4169	4125	8746

4.4 Tree Parsimony in Phylogenetic Inference (RQ3-1)

To evaluate the parsimony of tree structures generated by the model, we follow established methodologies (Zhou et al., 2024), minimizing the genetic mutations required to infer the optimal tree. The parsimony score evaluates how well the generated tree adheres to the principle of minimizing evolutionary changes, where fewer mutations are assumed to explain the observed genetic data better. We compare the results against the most parsimonious tree identified by the traditional PAUP* tool (Swofford, 1998). The parsimony score in Fig. 6 denotes the minimum mutations of genetic changes needed to account for the evolutionary relationships in the data. Since scores are plotted as negative values, lower scores indicate more complex trees and, consequently, poorer model performance. MDTree and ARTree achieved higher scores (approaching -4000) in fewer steps, reflecting simpler and more parsimonious trees. In contrast, PhyloGFN exhibited early fluctuations and ultimately stabilized around -5000, indicating suboptimal performance compared to others.

4.5 Tree Topological Diversity in Generated Trees (RQ3-2)

To assess the diversity of tree topologies generated by MDTree, we use three metrics: Simpson’s Diversity Index (He & Hu, 2005), Top Frequency, and Top 95% Frequency, as detailed in Table 6. A higher Diversity Index, which approaches 1, suggests broad diversity among generated tree topologies. A larger number of topologies in the Top 95% Frequency implies the generated trees are more varied and distributed across many unique structures. Conversely, a lower Top Frequency suggests the absence of a dominant tree structure,

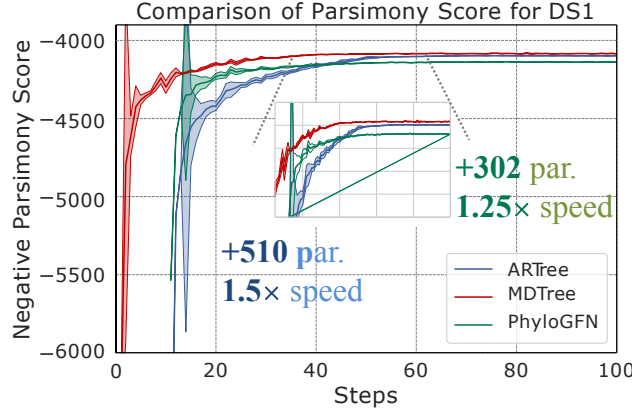


Figure 6: **Comparison of negative parsimony scores on the DS1 dataset.** The parsimony score denotes the minimum number of variation steps required to interpret each tree. The lower the negative score, the poorer the model performance.

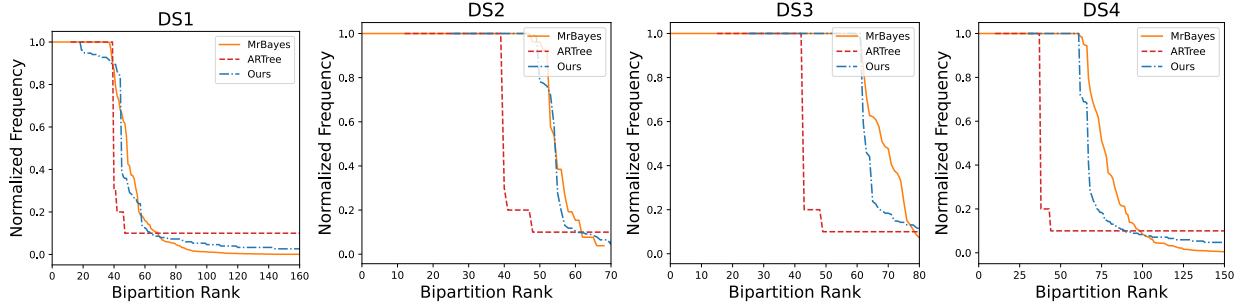


Figure 7: Bipartition frequency distribution of tree topologies. **The closer the two curves are, the better.**

pointing toward a more balanced generation. For instance, in DS3, with 36 species sequences, the Top 95% Topologies metric reveals 1,146 distinct tree structures, indicating a wide range of possible phylogenetic solutions. MDTree achieves a Diversity Index close to 1, showcasing its capacity for generating highly diverse topologies even in complex datasets. Furthermore, the Top Frequency metric remains notably low, further reinforcing the diversity and indicating that no single tree topology is overly dominant.

4.6 Bipartition Frequency for Tree Quality (RQ3-3)

In phylogenetic analysis, a bipartition refers to dividing taxa (species or genes) into two groups on either side of a node within the tree. When multiple tree samples are generated, as in Bayesian inference methods like MrBayes, each sample may have a different topology. Bipartition frequency quantifies how often a specific bipartition appears across all tree samples, providing insight into the support for particular evolutionary relationships. We use this bipartition frequency distribution to assess the model’s ability to capture phylogenetic relationships, as shown in Fig. 7. The horizontal axis indicates the bipartition rank within the tree topology, while the vertical axis displays the normalized occurrence frequency of each bipartition. The MDTree and MrBayes **curves are closely aligned**, indicating that MDTree’s results closely match those of the widely accepted gold standard. In contrast, the ARTree method shows a noticeable deviation, especially in the higher-ranked bipartitions, demonstrating that MDTree offers improved accuracy over ARTree in capturing evolutionary structures. This suggests that MDTree captures the evolutionary patterns with greater accuracy compared to ARTree.

Table 7: Comparison of different genomic language models (LMs) as structure generators in our framework, evaluated on Mean Log-Likelihood (MLL, \uparrow) and Evidence Lower Bound (ELBO, \uparrow). Models include DNABERT2, HyenaDNA, and NT. Higher values indicate better performance. DNABERT2 achieves the highest MLL and ELBO among the tested models, indicating its superior ability to capture genomic sequence patterns beneficial for phylogenetic inference.

Method	MLL(\uparrow)	ELBO(\uparrow)
DNABERT2	-7101.38	-7005.98
HyenaDNA	-7109.36	-7014.17
NT	-7111.07	-7017.11

Figure 8: **Ablation study of MDTree on four datasets, reported in mean log-likelihood (MLL) and ELBO (higher is better).** We evaluate the impact of removing the optimization phase, removing LAX in VIMCO, and removing the Dynamic Ordering Network (DON). The last column shows the average MLL across datasets, with green values indicating the drop compared to the full MDTree.

Method	DS1		DS2		DS3		DS4		Average
	MLL	ELBO	MLL	ELBO	MLL	ELBO	MLL	ELBO	
MDTree	-7101.38	-7005.98	-26357.96	-26362.75	-33715.31	-33430.94	-13322.10	-13113.03	-20051.18
w/o optimization	-7106.59	-7010.34	-26371.02	-26374.01	-33733.25	-33447.94	-13339.71	-13130.01	-20064.11 (-12.93)
w/ vimco w/o Lax	-7103.74	-7007.86	-26361.81	-26368.52	-33718.20	-33436.07	-13326.95	-13118.60	-20055.22 (-4.04)
w/o DON	-7105.05	-7010.02	-26366.47	-26372.04	-33723.67	-33439.18	-13332.38	-13121.33	-20058.77 (-7.59)

4.7 Analysis and Ablation (RQ4-1)

We compare MDTree with three other schemes, yielding the following observations: (i) Removing optimization techniques like RWS or VIMCO led to a performance drop of 5.21 in MLL, as shown by slight fluctuations in the MLL curve in Fig. 9, highlighting their role in stabilizing convergence. (ii) Excluding the LAX model of VIMCO optimization caused a decrease of 2.36 in MLL and 1.88 in ELBO, indicating its effectiveness in reducing variance during discrete sampling. (iii) Table 7 and Table 8 show that the removal of the DON results in the most significant impact, with a drop of about 3.67 in MLL, underscoring its critical role in optimizing node addition order and improving tree generation. Overall, the full MDTree consistently achieves the best across both metrics. We select the genome-specific foundation model DNABERT2 for our phylogenetic inference research. Although models like HyenaDNA (Nguyen et al., 2023) and Nucleotide Transformer (NT) (Dalla-Torre et al., 2023) excel in long-sequence modeling, they are less apt for our specific needs. As shown in Table 7, DNABERT2 outperforms others, likely due to its specific optimization for genomic data.

4.8 Visualization of PhyloTree Structure on Real-World Data (RQ5)

To assess the biological relevance of the tree structure generated by MDTree, we applied it to construct a phylogenetic tree for an Angiosperms353 genomic dataset (Zuntini et al., 2024). The tree successfully recovered major branches within the order Rosales, revealing distinct evolutionary lineages, including Rosaceae, Moraceae, and Polygonaceae families. As shown in Fig. 10, the genera *Polygala vulgaris* and *Polygala baldunii* are clearly separated from other groups, consistent with their classification in the Potentillaceae family. The remaining groups, distinguished by color, represent genera within the Rosaceae and Moraceae families, such as *Rosa*, *Rubus*, *Ficus*, and *Adansonia*. In Rosaceae, genera like *Rosa*, *Rubus*, and *Prunus* highlight their common evolutionary ancestry, while in Moraceae, *Ficus* and *Broussonetia* reflect the internal diversity and evolutionary divergence within the family.

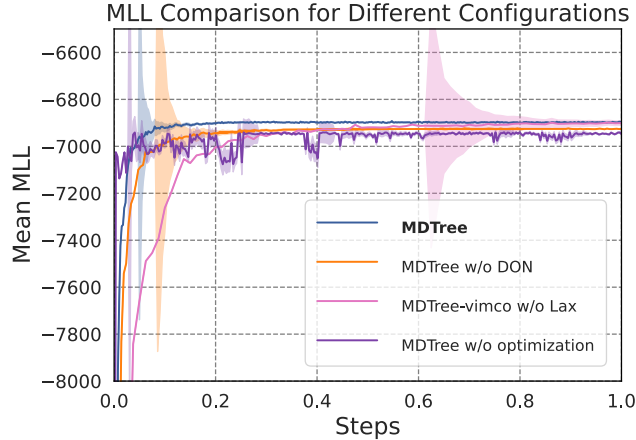


Figure 9: Ablation of different modules. MDTree w/o optimization curve exhibits **slight fluctuations**, emphasizing the importance of **optimization techniques** in improving stability.

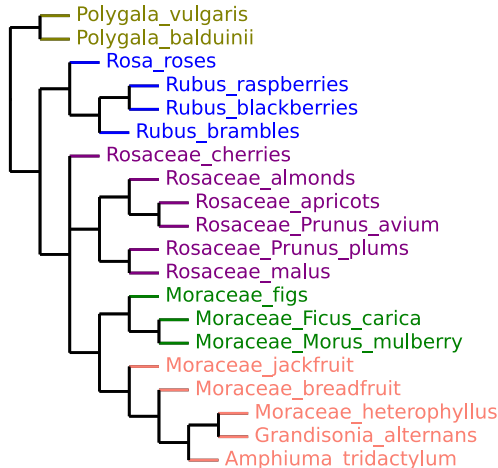


Figure 10: **Visualization of a generated phylogenetic tree for a subset of species from the Angiosperms353 dataset.** Different colors indicate distinct plant families or genera, illustrating the model’s ability to cluster related species into coherent subtrees. For example, species within the genus *Rubus* (blue) and family *Moraceae* (green) are correctly grouped together, reflecting biologically plausible evolutionary relationships. This demonstrates that the proposed method can recover meaningful phylogenetic structure consistent with known taxonomy.

5 Conclusion and Limitation

In this paper, we present MDTree, a novel framework that redefines phylogenetic tree generation as a Dynamic Autoregressive Tree Generation task. By leveraging a Diffusion Ordering Network to learn biologically informed node orders directly from genomic sequences, MDTree overcomes the limitations of fixed or random node orders. It integrates GNNs and Language Models to capture complex tree topologies, while a Dynamic Masking Mechanism enables parallel node processing, improving computational efficiency. Experiments on phylogenetic benchmarks show MDTree achieves state-of-the-art performance.

MDTree has yet to be applied to other sequence types, such as protein sequences. Future work will explore multimodal approaches, integrating genomic and protein data for more comprehensive evolutionary tree construction, as well as scaling the model for complex evolutionary scenarios.

Broader Impact Statement

In this optional section, TMLR encourages authors to discuss possible repercussions of their work, notably any potential negative impact that a user of this research should be aware of. Authors should consult the TMLR Ethics Guidelines available on the TMLR website for guidance on how to approach this subject.

Author Contributions

If you'd like to, you may include a section for author contributions as is done in many journals. This is optional and at the discretion of the authors. Only add this information once your submission is accepted and deanonymized.

Acknowledgments

This work was supported in part by National Science and Technology Major Project (No. 2022ZD0115101), National Natural Science Foundation of China Project (No. U21A20427), Project (No. WU2022A009) from the Center of Synthetic Biology and Integrated Bioengineering of Westlake University and Integrated Bioengineering of Westlake University and Project (No. WU2023C019) from the Westlake University Industries of the Future Research Funding, and the InnoHK Initiative by the Government of the Hong Kong Special Administrative Region (HKSAR). We thank the AI Station of Westlake University for the support of GPUs.

References

- Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.
- Roman Biek, Oliver G Pybus, James O Lloyd-Smith, and Xavier Didelot. Measurably evolving pathogens in the genomic era. *Trends in Ecology & Evolution*, 30(6):306–313, 2015.
- Sam Bond-Taylor, Peter Hessey, Hiroshi Sasaki, Toby P. Breckon, and Chris G. Willcocks. Unleashing transformers: Parallel token prediction with discrete absorbing diffusion for fast high-resolution image generation from vector-quantized codes, 2021. URL <https://arxiv.org/abs/2111.12701>
- Luciano Brocchieri. Phylogenetic inferences from molecular sequences: review and critique. *Theoretical population biology*, 59(1):27–40, 2001.
- Hugo Dalla-Torre, Liam Gonzalez, Javier Mendoza-Revilla, Nicolas Lopez Carranza, Adam Henryk Grzywaczewski, Francesco Oteri, Christian Dallago, Evan Trop, Bernardo P de Almeida, Hassan Sirelkhatim, et al. The nucleotide transformer: Building and evaluating robust foundation models for human genomics. *bioRxiv*, pp. 2023–01, 2023.
- Persi Diaconis. Sequential importance sampling for estimating the number of perfect matchings in bipartite graphs: An ongoing conversation with laci. *Building Bridges II: Mathematics of László Lovász*, pp. 223–233, 2019.
- Danai Fimereli, David Venet, Mattia Rediti, Bram Boeckx, Marion Maetens, Samira Majjaj, Ghizlane Rouas, Caterina Marchio, Francois Bertucci, Odette Mariani, et al. Timing evolution of lobular breast cancer through phylogenetic analysis. *EBioMedicine*, 82, 2022.
- James R Garey, Thomas J Near, Michael R Nonnemacher, and Steven A Nadler. Molecular evidence for acanthocephala as a subtaxon of rotifera. *Journal of Molecular Evolution*, 43:287–292, 1996.
- Will Grathwohl, Dami Choi, Yuhuai Wu, Geoffrey Roeder, and David Duvenaud. Backpropagation through the void: Optimizing control variates for black-box gradient estimation. *arXiv preprint arXiv:1711.00123*, 2017.
- T Ryan Gregory. Understanding evolutionary trees. *Evolution: Education and Outreach*, 1:121–137, 2008.

- Fangliang He and Xin-Sheng Hu. Hubbell’s fundamental biodiversity parameter and the simpson diversity index. *Ecology Letters*, 8(4):386–390, 2005.
- S Blair Hedges, Kirk D Moberg, and Linda R Maxson. Tetrapod phylogeny inferred from 18s and 28s ribosomal rna sequences and a review of the evidence for amniote relationships. *Molecular Biology and Evolution*, 7(6):607–633, 1990.
- Daniel A Henk, Alex Weir, and Meredith Blackwell. Laboulbeniopsis termitarius, an ectoparasite of termites newly recognized as a member of the laboulbeniomycetes. *Mycologia*, 95(4):561–564, 2003.
- Edward J Hu, Nikolay Malkin, Moksh Jain, Katie E Everett, Alexandros Graikos, and Yoshua Bengio. Gflownet-em for learning compositional latent variable models. In *International Conference on Machine Learning*, pp. 13528–13549. PMLR, 2023.
- Philip Hugenholtz, Maria Chuvochina, Aharon Oren, Donovan H Parks, and Rochelle M Soo. Prokaryotic taxonomy and nomenclature in the age of big sequence data. *the ISME Journal*, 15(7):1879–1892, 2021.
- Fernando Izquierdo-Carrasco, Stephen A Smith, and Alexandros Stamatakis. Algorithms, data structures, and numerics for likelihood-based phylogenetic inference of huge trees. *BMC bioinformatics*, 12:1–14, 2011.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Hazal Koptagel, Oskar Kviman, Harald Melin, Negar Safinianaini, and Jens Lagergren. Vaiphy: a variational inference based algorithm for phylogeny. *Advances in Neural Information Processing Systems*, 35:14758–14770, 2022.
- Clemens Lakner, Paul Van Der Mark, John P Huelsenbeck, Bret Larget, and Fredrik Ronquist. Efficiency of markov chain monte carlo tree proposals in bayesian phylogenetics. *Systematic biology*, 57(1):86–103, 2008.
- Takahiro Mimori and Michiaki Hamada. Geophy: differentiable phylogenetic inference via geometric gradients of tree topologies. *Advances in Neural Information Processing Systems*, 36, 2024.
- Geetika Munjal, Madasu Hanmandlu, and Sangeet Srivastava. Phylogenetics algorithms and applications. In *Ambient Communications and Computer Systems: RACCS-2018*, pp. 187–194. Springer, 2019.
- Hamish Nisbet Munro. *Mammalian protein metabolism*, volume 4. Elsevier, 2012.
- Luca Nesterenko, Bastien Boussau, and Laurent Jacob. Phyloformer: towards fast and accurate phylogeny estimation with self-attention networks. *bioRxiv*, pp. 2022–06, 2022.
- Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Callum Birch-Sykes, Michael Wornow, Aman Patel, Clayton Rabideau, Stefano Massaro, Yoshua Bengio, et al. Hyenadna: Long-range genomic sequence modeling at single nucleotide resolution. *arXiv preprint arXiv:2306.15794*, 2023.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- David Penny. Inferring phylogenies.—joseph felsenstein. 2003. sinauer associates, sunderland, massachusetts., 2004.
- Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019.

- Fredrik Ronquist, Maxim Teslenko, Paul Van Der Mark, Daniel L Ayres, Aaron Darling, Sebastian Höhna, Bret Larget, Liang Liu, Marc A Suchard, and John P Huelsenbeck. Mrbayes 3.2: efficient bayesian phylogenetic inference and model choice across a large model space. *Systematic biology*, 61(3):539–542, 2012.
- Amy Y Rossman, John M McKemy, Rebecca A Pardo-Schultheiss, and Hans-Josef Schroers. Molecular studies of the bionectriaceae using large subunit rDNA sequences. *Mycologia*, 93(1):100–110, 2001.
- Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In *The semantic web: 15th international conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, proceedings 15*, pp. 593–607. Springer, 2018.
- Megan L Smith and Matthew W Hahn. Phylogenetic inference using generative adversarial networks. *Bioinformatics*, 39(9):btad543, 2023.
- Claudia Solís-Lemus and Cécile Ané. Inferring phylogenetic networks with maximum pseudolikelihood under incomplete lineage sorting. *PLoS genetics*, 12(3):e1005896, 2016.
- David L Swofford. Phylogenetic analysis using parsimony. 1998.
- Xudong Tang, Leonardo Zepeda-Núñez, Shengwen Yang, Zelin Zhao, and Claudia Solís-Lemus. Novel symmetry-preserving neural network model for phylogenetic inference. *Bioinformatics Advances*, 4(1):vbae022, 2024.
- Kathrin Theissinger, Carlos Fernandes, Giulio Formenti, Iliana Bista, Paul R Berg, Christoph Bleidorn, Aureliano Bombarely, Angelica Crottini, Guido R Gallo, José A Godoy, et al. How genomics can help biodiversity conservation. *Trends in genetics*, 39(7):545–559, 2023.
- Bart Tummers and Douglas R Green. The evolution of regulated cell death pathways in animals and their evasion by pathogens. *Physiological reviews*, 102(1):411–454, 2022.
- A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- Liangliang Wang, Shijia Wang, and Alexandre Bouchard-Côté. An annealed sequential monte carlo method for bayesian phylogenetics. *Systematic biology*, 69(1):155–183, 2020.
- Tianyu Xie and Cheng Zhang. Artree: A deep autoregressive model for phylogenetic inference. *Advances in Neural Information Processing Systems*, 36, 2024.
- Wangang Xie, Paul O Lewis, Yu Fan, Lynn Kuo, and Ming-Hui Chen. Improving marginal likelihood estimation for bayesian phylogenetic model selection. *Systematic biology*, 60(2):150–160, 2011.
- Ziheng Yang. Maximum likelihood phylogenetic estimation from dna sequences with variable rates over sites: approximate methods. *Journal of Molecular Evolution*, 39(3):306–314, 1994.
- Ziheng Yang and Anne D Yoder. Comparison of likelihood and bayesian methods for estimating divergence times using multiple gene loci and calibration points, with application to a radiation of cute-looking mouse lemur species. *Systematic biology*, 52(5):705–716, 2003.
- Anne D Yoder and Ziheng Yang. Divergence dates for malagasy lemurs estimated from multiple gene loci: geological and evolutionary context. *Molecular Ecology*, 13(4):757–773, 2004.
- Cheng Zhang. Learnable topological features for phylogenetic inference via graph neural networks. *ArXiv*, 2023.
- Cheng Zhang and Frederick A Matsen IV. Generalizing tree probability estimation via bayesian networks. *Advances in neural information processing systems*, 31, 2018a.
- Cheng Zhang and Frederick A Matsen IV. Variational bayesian phylogenetic inference. In *International Conference on Learning Representations*, 2018b.

- Chi Zhang, Huw A Ogilvie, Alexei J Drummond, and Tanja Stadler. Bayesian inference of species networks from multilocus sequence data. *Molecular biology and evolution*, 35(2):504–517, 2018.
- Ning Zhang and Meredith Blackwell. Molecular phylogeny of dogwood anthracnose fungus (*discula destructiva*) and the *diaporthales*. *Mycologia*, 93(2):355–365, 2001.
- Mingyang Zhou, Zichao Yan, Elliot Layne, Nikolay Malkin, Dinghuai Zhang, Moksh Jain, Mathieu Blanchette, and Yoshua Bengio. Phylogfn: Phylogenetic inference with generative flow networks. *arXiv preprint arXiv:2310.08774*, 2024.
- Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Dutta, Ramana Davuluri, and Han Liu. Dnabert-2: Efficient foundation model and benchmark for multi-species genome. *arXiv preprint arXiv:2306.15006*, 2023.
- Alexandre R Zuntini, Tom Carruthers, Olivier Maurin, Paul C Bailey, Kevin Leempoel, Grace E Brewer, Niroshini Epitawalage, Elaine Françoso, Berta Gallego-Paramo, Catherine McGinnie, et al. Phylogenomics and the rise of the angiosperms. *Nature*, pp. 1–8, 2024.