

# PhyloGen: Language Model-Enhanced Phylogenetic Inference via Graph Structure Generation

ChenRui Duan,  
Zelin Zang, Siyuan Li, Yongjie Xu  
Stan Z. Li

Zhejiang University, College of Computer Science and Technology;  
Westlake University



## Background

- (a) The inputs are aligned sequences, and topologies are learned from existing tree structures using methods like SBNs, which rely on MCMC-based methods for pre-generated candidate trees without considering branch lengths directly.
- (b) The inputs are aligned sequences, and then tree structures and branch lengths are directly inferred by variational inference and biological modules. These methods optimize tree topology and branch lengths separately.
- (c) The inputs are raw sequences processed by a pre-trained language model to generate species representations. Then, an initial topology is generated through a tree construction module, and the topology and branch lengths are co-optimized by the tree structure modeling module.

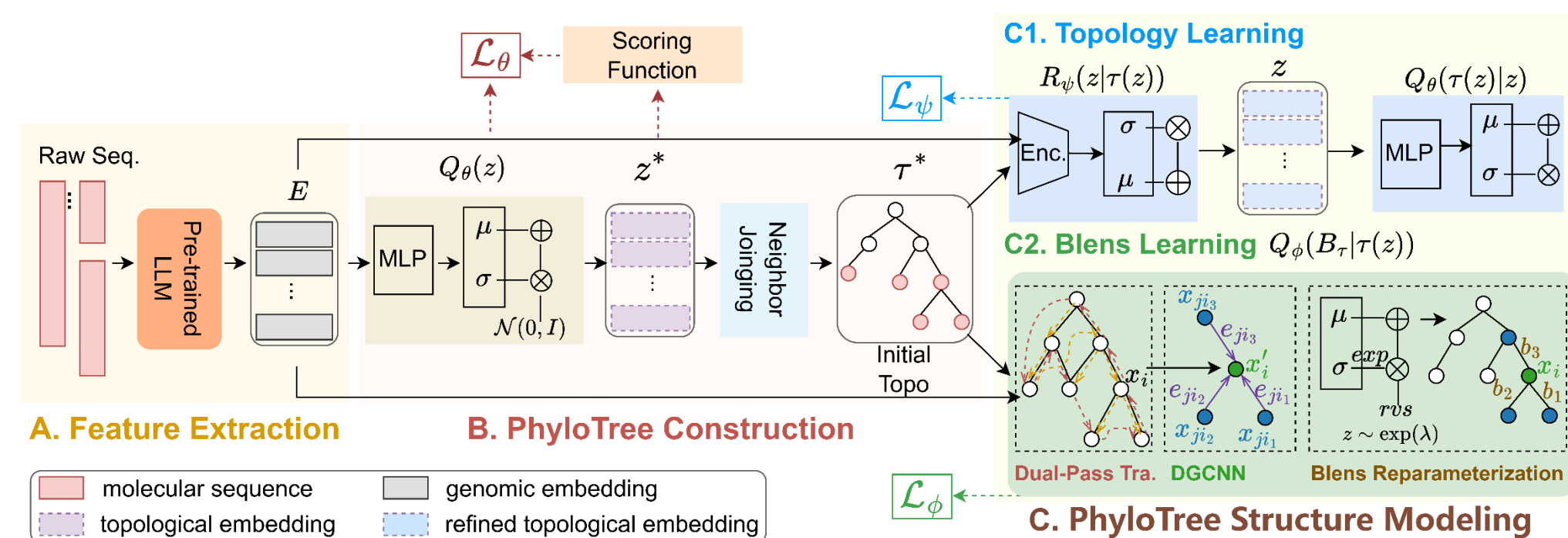
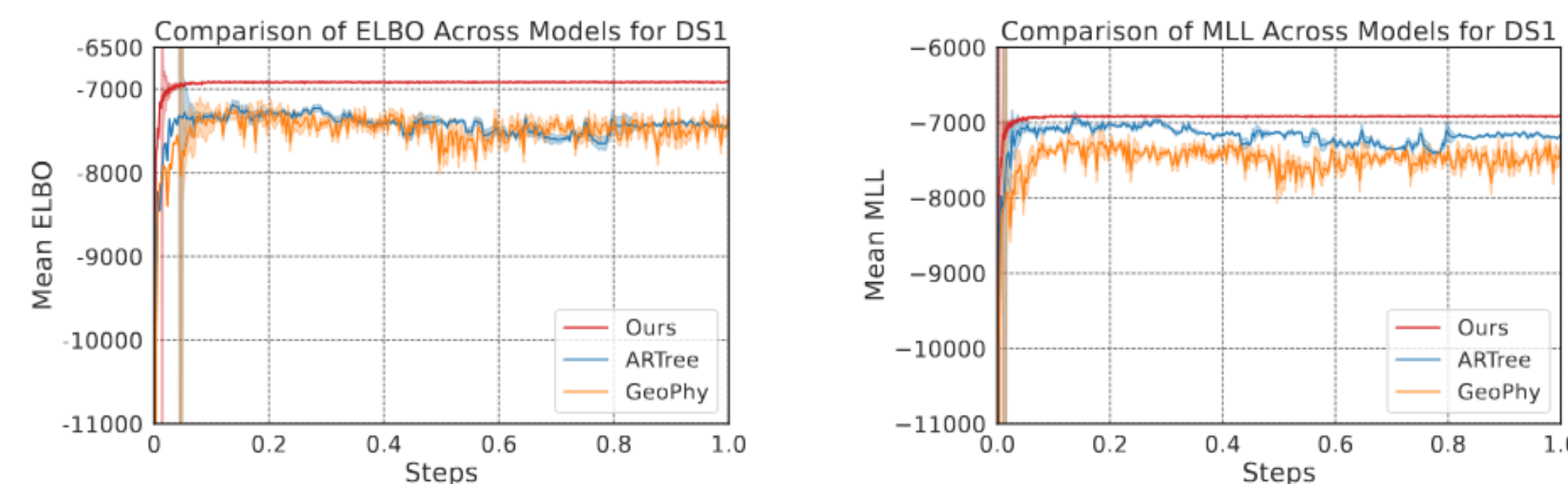


Figure 2: Our Method: Framework of PhyloGen

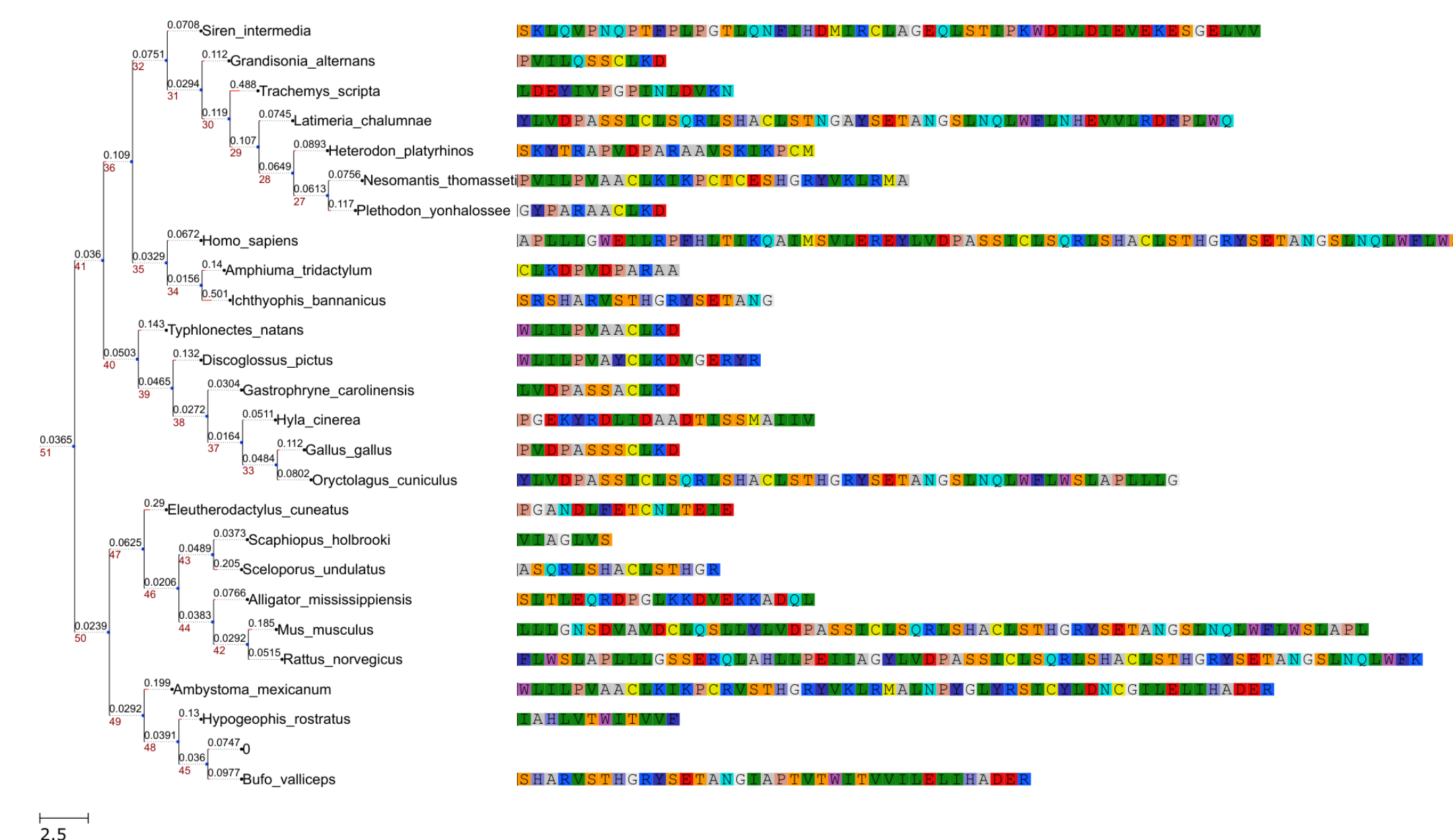
**A. Feature Extraction** module extracts genome embeddings  $E$  from raw sequences  $Y$  using a pre-trained language model.  
**B. PhyloTree Construction** module uses  $E$  to compute topological parameters, which generate an initial tree structure  $\tau^*$  via the Neighbor-Joining algorithm.  
**C. PhyloTree Structure Modeling** module jointly model  $\tau$  and  $B_r$  through the **topology learning component** (TreeEncoder  $\mathcal{R}_\psi$  and TreeDecoder  $Q$ ) and the **branch length-(Blens) learning component** (dual-pass traversal, DGCNN network, Blens reparameterization).



	Metric	PhyloGFN	GeoPhy	GeoPhy LOO(3)+	Ours w/o KL	Ours w/o S	Ours
setting1	ELBO ( $\Delta$ )	NA	-7721.82 (-100)	-7729.28 (-107)	-6725.49 (+12)	-6713.01 (+15)	<b>-6711.47 (+14)</b>
	MLL ( $\Delta$ )	-6705.55 (-93)	-7440.38 (-198)	-7599.85 (-357)	-6564.51 (+18)	-6547.32 (+20)	<b>-6542.75 (+22)</b>
	Time	18h28min	6h16min	15h23min	6h43min	7h42min	<b>6h32min</b>
setting2	ELBO ( $\Delta$ )	NA	-11802.07 (-98)	-11676.29 (-76)	-10678.24 (+3)	-10655.02 (+2)	<b>-10674.28 (+4)</b>
	MLL ( $\Delta$ )	-12565.76 (-233)	-11763.40 (-131)	-11630.16 (-98)	-10422.12 (+6)	-10654.53 (+12)	<b>-10432.71 (+5)</b>
	Time	24h35min	18h13min	12h24min	7h54min	8h6min	<b>6h37min</b>

Table 1: Comparison of the MLL ( $\uparrow$ ) with different approaches in eight benchmark datasets. VBPI and VBPI-GNN use pre-generated tree topologies in training and thus **are not directly comparable**. **Boldface** for the highest result, **underline** for the second highest from traditional methods, and **underline** for the second highest from tree structure generation methods.

Methods	Dataset #Taxa (N)	DS1 27	DS2 29	DS3 36	DS4 41	DS5 50	DS6 50	DS7 59	DS8 64
MCMC-based	MrBayes	-7108.42 (0.18)	-26367.57 (0.48)	-33735.44 (0.50)	-13330.44 (0.54)	<b>-8214.51</b> (0.28)	<b>-6724.07</b> (0.86)	-37332.76 (2.42)	<b>-8649.88</b> (1.75)
	SBN	<b>-7108.41</b> (0.15)	<b>-26367.71</b> (0.08)	<b>-33735.09</b> (0.09)	<b>-13329.94</b> (0.20)	-8214.62 (0.40)	-6724.37 (0.43)	<b>-37331.97</b> (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 (0.43)	-37332.03 (0.55)	-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	<b>-7108.41</b> (0.19)	<b>-26367.71</b> (0.07)	<b>-33735.09</b> (0.09)	<b>-13329.94</b> (0.17)	<b>-8214.59</b> (0.34)	<b>-6724.37</b> (0.46)	<b>-37331.95</b> (0.27)	<b>-8650.61</b> (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)	-133342.71 (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)	-26368.9 (0.28)	-33735.6 (0.35)	-13331.83 (0.19)	-8215.15 (0.20)	-6730.68 (0.54)	-37359.96 (1.14)	-8654.76 (0.19)
	<b>Ours</b>	<b>-6910.02</b> (0.07)	<b>-26257.09</b> (0.06)	<b>-33481.57</b> (0.10)	<b>-13063.15</b> (1.34)	<b>-7928.4</b> (0.23)	<b>-6330.21</b> (0.31)	<b>-36838.42</b> (12.03)	<b>-8171.04</b> (0.96)



## Bipartition Frequency Distribution (RQ3)

The closer the two curves are, the better, which suggests that PhyloGen is highly consistent with the gold standard MrBayes approach.

Fig.5 shows the bipartition frequency distributions of trees inferred by PhyloGen for datasets DS1, DS2, and DS3. The horizontal axis indicates the ranking of the bipartitions in the tree topology, and the vertical axis indicates the normalized frequency of occurrence of the corresponding bipartitions. The **similarity** of our method's curves to those of MrBayes underscores its accuracy, demonstrating that PhyloGen consistently captures evolutionary patterns with reliability comparable to the gold standard.

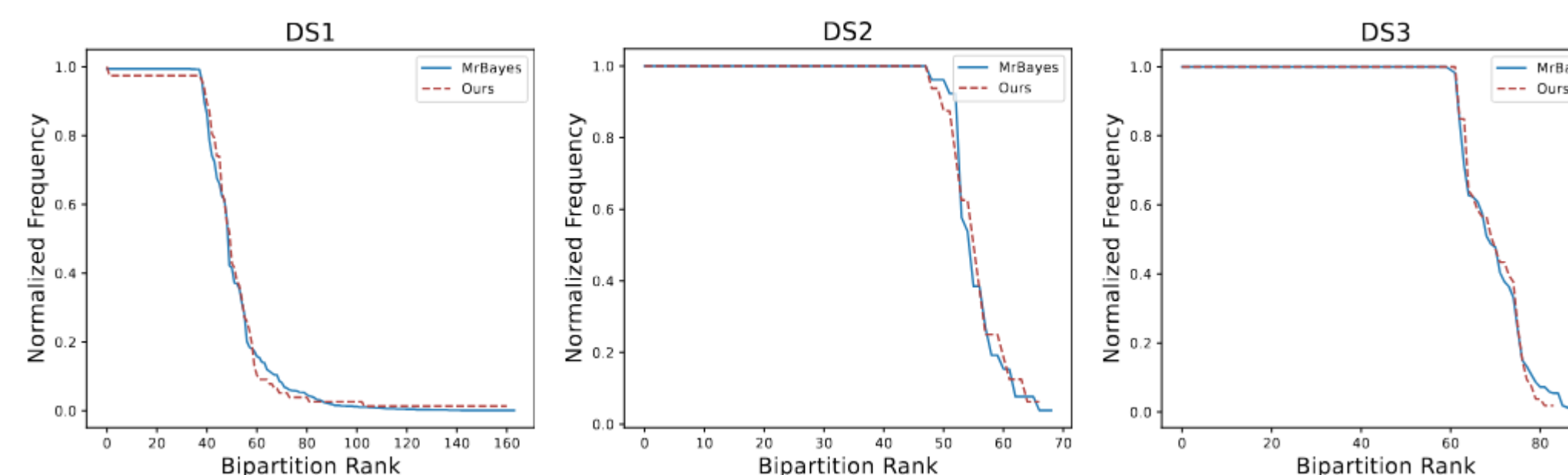


Figure 5: Comparative Bipartition Frequency Distribution in Tree Topologies for DS1, DS2, and DS3 datasets.

## Our Method

Reformulate as graph generation problem.

As depicted in Fig.1 (c), we propose a novel approach based on a pre-trained genome language model. Our model does not rely on evolutionary models or the requirement to align input sequences to equal lengths and fully exploits the prior knowledge embedded in biological sequences. PhyloGen models phylogenetic tree inference as a conditional-constrained tree structure generation problem, aiming to generate and optimize the tree topology and branch lengths jointly. We map species sequences into a continuous geometric space and perform end-to-end variational inference without restricting topological candidates. To ensure the topology-invariance of phylogenetic trees, we incorporate distance constraints in the latent space to maintain translational rotation invariance. Our approach demonstrates effectiveness and efficiency on the eight real-world benchmark datasets and verifies its robustness through data augmentation and noise addition. In addition, we propose a new scoring function to guide the model towards a more stable and faster gradient descent.

