

# Variational Bayesian Phylogenetic Inference



**FRED HUTCH™**  
CURES START HERE

Cheng Zhang

joint work with Frederick Matsen

Fred Hutchinson Cancer Research Center, Seattle WA

Jun 10, 2019

# Introduction

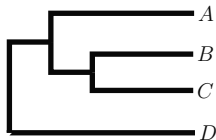
---



The goal of **phylogenetic inference** is to reconstruct the evolution history (e.g., *phylogenetic trees*) from **molecular sequence data** (e.g., DNA, RNA or protein sequences)

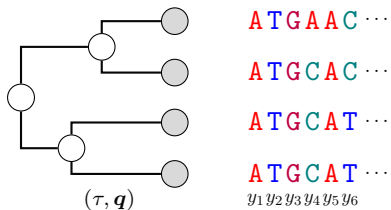
Taxa	Characters
Species A	ATGAACAT
Species B	ATGCACAC
Species C	ATGCATAT
Species D	ATGCATGC

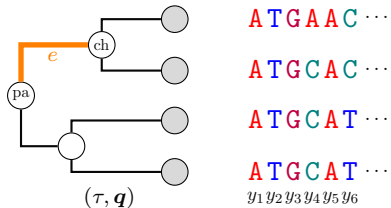
Molecular Sequence Data



Phylogenetic Tree





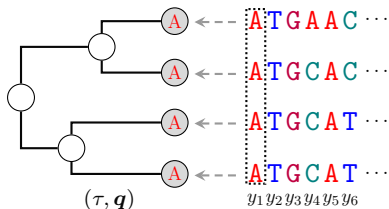


Evolution model:

$$p(\text{ch}|\text{pa}, q_e)$$

$q_e$ : amount of evolution on  $e$ .





Evolution model:

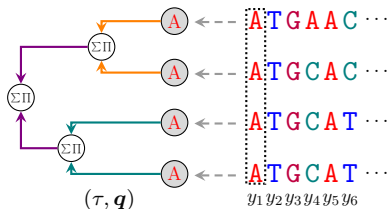
$$p(\text{ch}|\text{pa}, q_e)$$

$q_e$ : amount of evolution on  $e$ .

## Likelihood

$$p(\mathbf{Y}|\tau, \mathbf{q}) = \eta(a_\rho^i) \prod_{(u,v) \in E(\tau)} P_{a_u^i a_v^i}(q_{uv})$$





Evolution model:

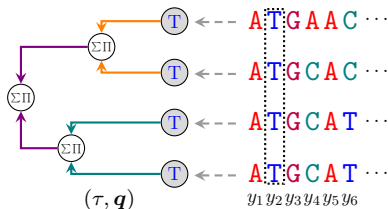
$$p(\text{ch}|\text{pa}, q_e)$$

$q_e$ : amount of evolution on  $e$ .

## Likelihood

$$p(\mathbf{Y}|\tau, \mathbf{q}) = \sum_{a^i} \eta(a^i) \prod_{(u,v) \in E(\tau)} P_{a_u^i a_v^i}(q_{uv})$$





Evolution model:

$$p(\text{ch}|\text{pa}, q_e)$$

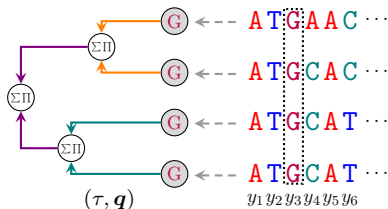
$q_e$ : amount of evolution on  $e$ .

Likelihood

$$p(\mathbf{Y}|\tau, \mathbf{q}) = \prod_{i=1}^M \sum_{a^i} \eta(a^i) \prod_{(u,v) \in E(\tau)} P_{a_u^i a_v^i}(q_{uv})$$







Evolution model:

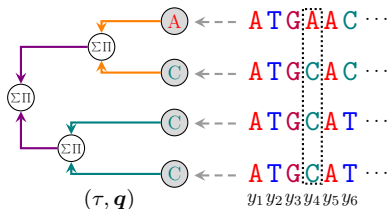
$$p(\text{ch}|\text{pa}, q_e)$$

$q_e$ : amount of evolution on  $e$ .

Likelihood

$$p(\mathbf{Y}|\tau, \mathbf{q}) = \prod_{i=1}^M \sum_{a^i} \eta(a^i) \prod_{(u,v) \in E(\tau)} P_{a_u^i a_v^i}(q_{uv})$$





Evolution model:

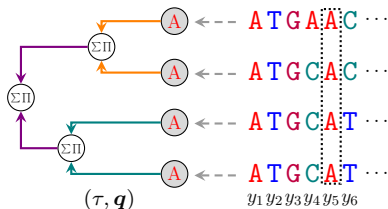
$$p(\text{ch}|\text{pa}, q_e)$$

$q_e$ : amount of evolution on  $e$ .

Likelihood

$$p(\mathbf{Y}|\tau, \mathbf{q}) = \prod_{i=1}^M \sum_{a^i} \eta(a^i) \prod_{(u,v) \in E(\tau)} P_{a_u^i a_v^i}(q_{uv})$$





Evolution model:

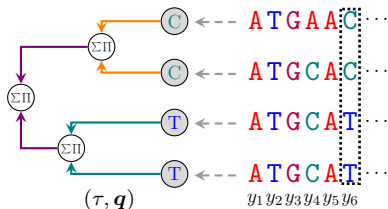
$$p(\text{ch}|\text{pa}, q_e)$$

$q_e$ : amount of evolution on  $e$ .

Likelihood

$$p(\mathbf{Y}|\tau, \mathbf{q}) = \prod_{i=1}^M \sum_{a^i} \eta(a^i) \prod_{(u,v) \in E(\tau)} P_{a_u^i a_v^i}(q_{uv})$$





Evolution model:

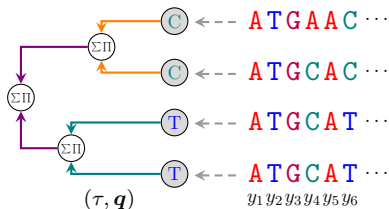
$$p(\text{ch}|\text{pa}, q_e)$$

$q_e$ : amount of evolution on  $e$ .

Likelihood

$$p(\mathbf{Y}|\tau, \mathbf{q}) = \prod_{i=1}^M \sum_{a^i} \eta(a^i) \prod_{(u,v) \in E(\tau)} P_{a_u^i a_v^i}(q_{uv})$$





Evolution model:

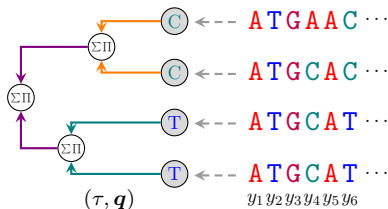
$$p(\text{ch}|\text{pa}, q_e)$$

$q_e$ : amount of evolution on  $e$ .

Likelihood

$$p(\mathbf{Y}|\tau, \mathbf{q}) = \prod_{i=1}^M \sum_{a^i} \eta(a^i) \prod_{(u,v) \in E(\tau)} P_{a_u^i a_v^i}(q_{uv})$$





Evolution model:

$$p(\text{ch}|\text{pa}, q_e)$$

$q_e$ : amount of evolution on  $e$ .

## Likelihood

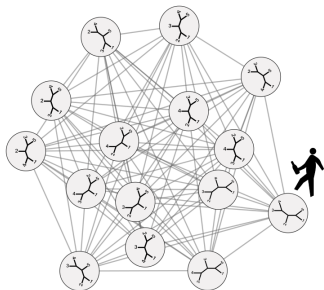
$$p(\mathbf{Y}|\tau, \mathbf{q}) = \prod_{i=1}^M \sum_{a^i} \eta(a^i) \prod_{(u,v) \in E(\tau)} P_{a_u^i a_v^i}(q_{uv})$$

Given a proper prior distribution  $p(\tau, \mathbf{q})$ , the **posterior** is

$$p(\tau, \mathbf{q}|\mathbf{Y}) \propto p(\mathbf{Y}|\tau, \mathbf{q})p(\tau, \mathbf{q}).$$



**Random-walk MCMC:** use simple **random perturbation** (e.g., Nearest Neighborhood Interchange) to generate new state



## Challenges

- ▶ **Large** search space (**combinatorially** exploding)

$$\# \text{ unrooted trees } (n \text{ taxa}) = (2n - 5)!!$$

- ▶ **Intertwined** parameter space, **low** acceptance rate, **hard** to scale to data sets with many sequences.



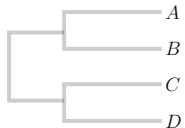
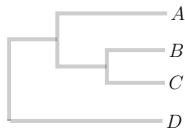
# Variational Bayesian Phylogenetic Inference

---

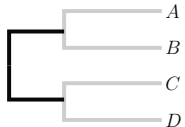
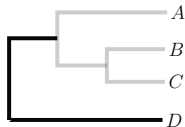




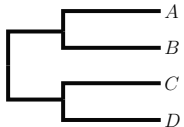
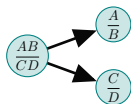
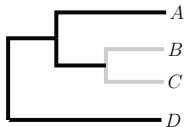
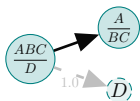
Encode tree structures via Bayesian networks!



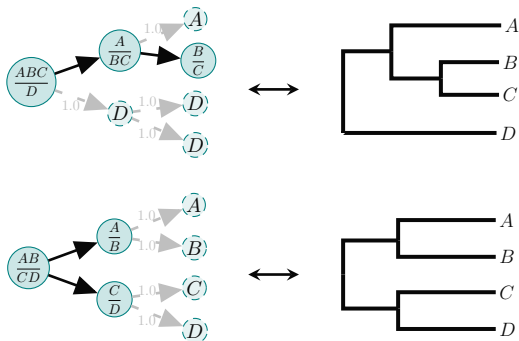
Encode tree structures via Bayesian networks!



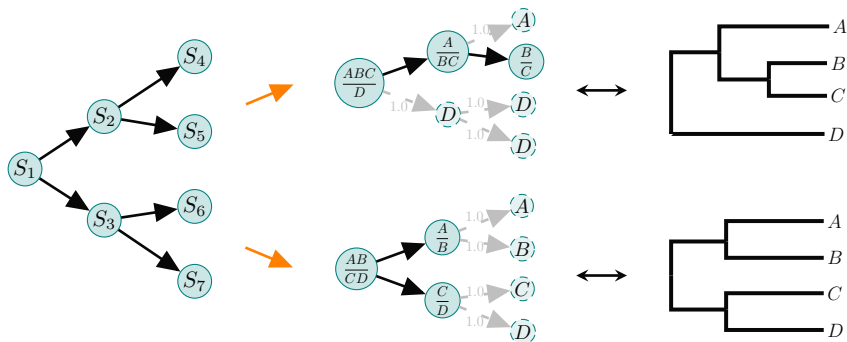
Encode tree structures via Bayesian networks!

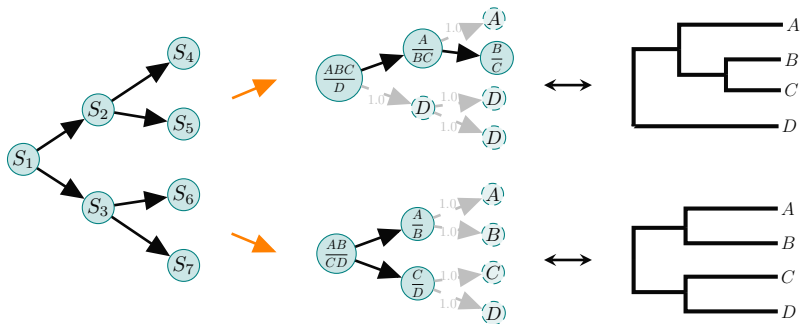


Encode tree structures via Bayesian networks!



Encode tree structures via Bayesian networks!

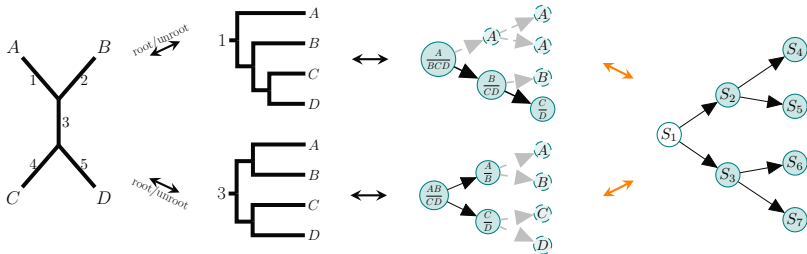




## Rooted Trees

$$p_{\text{sbn}}(T = \tau) = p(S_1 = s_1) \prod_{i>1} p(S_i = s_i | S_{\pi_i} = s_{\pi_i}).$$





**Unrooted Trees:**

$$p_{\text{sbn}}(T^{\text{u}} = \tau) = \sum_{s_1 \sim \tau} p(S_1 = s_1) \prod_{i > 1} p(S_i = s_i | S_{\pi_i} = s_{\pi_i}).$$



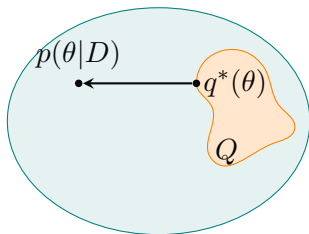
Use SBNs to refine posterior probability estimation given MCMC samples, useful when the posterior is diffuse.

DATA SET	(#TAXA, #SITES)	TREE SPACE SIZE	SAMPLED TREES	KL DIVERGENCE TO GROUND TRUTH				
				SRF	CCD	SBN-SA	SBN-EM	SBN-EM- $\alpha$
DS1	(27, 1949)	$5.84 \times 10^{32}$	1228	0.0155	0.6027	0.0687	0.0136	<b>0.0130</b>
DS2	(29, 2520)	$1.58 \times 10^{35}$	7	<b>0.0122</b>	0.0218	0.0218	0.0199	0.0128
DS3	(36, 1812)	$4.89 \times 10^{47}$	43	0.3539	0.2074	0.1152	0.1243	<b>0.0882</b>
DS4	(41, 1137)	$1.01 \times 10^{57}$	828	0.5322	0.1952	0.1021	0.0763	<b>0.0637</b>
DS5	(50, 378)	$2.84 \times 10^{74}$	33752	11.5746	1.3272	0.8952	0.8599	<b>0.8218</b>
DS6	(50, 1133)	$2.84 \times 10^{74}$	35407	10.0159	0.4526	<b>0.2613</b>	0.3016	0.2786
DS7	(59, 1824)	$4.36 \times 10^{92}$	1125	1.2765	0.3292	0.2341	0.0483	<b>0.0399</b>
DS8	(64, 1008)	$1.04 \times 10^{103}$	3067	2.1653	0.4149	0.2212	0.1415	<b>0.1236</b>

[Zhang and Matsen, NeurIPS 2018]

**Remark:** Unlike previous methods, SBNs are flexible enough to provide accurate approximations to real data posteriors!



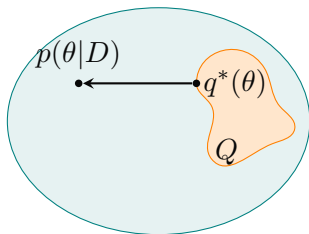


An **optimization** approach

$$q^* = \arg \min_{q \in Q} D_{KL}(q(\theta) || p(\theta|D))$$

$Q$ : **approximating** distributions.





An **optimization** approach

$$q^* = \arg \min_{q \in Q} D_{KL}(q(\theta) || p(\theta|D))$$

$Q$ : **approximating** distributions.

equivalent to maximizing the **evidence lower bound (ELBO)**

$$L(q, D) = \mathbb{E}_{q(\theta)} \log \left( \frac{p(\theta, D)}{q(\theta)} \right) \leq \log p(D)$$

- ▶ tends to be **faster than MCMC**
- ▶ **easy to scale** to large data sets (via stochastic gradient ascent)



- ▶ Approximating Distribution:

tree topology

$$Q_{\phi}(\tau)$$



- ▶ Approximating Distribution:

$$Q_{\phi,\psi}(\tau, q) \triangleq \overset{\text{tree topology}}{Q_{\phi}(\tau)} \cdot \overset{\text{branch length}}{Q_{\psi}(q|\tau)}$$



- ▶ Approximating Distribution:

$$Q_{\phi, \psi}(\tau, q) \triangleq \overset{\text{tree topology}}{Q_{\phi}(\tau)} \cdot \overset{\text{branch length}}{Q_{\psi}(q|\tau)}$$

- ▶ Multi-sample Lower Bound:

$$L^K(\phi, \psi) = \mathbb{E}_{Q_{\phi, \psi}(\tau^{1:K}, \mathbf{q}^{1:K})} \log \left( \frac{1}{K} \sum_{i=1}^K \frac{p(\mathbf{Y}|\tau^i, \mathbf{q}^i)p(\tau^i, \mathbf{q}^i)}{Q_{\phi}(\tau^i)Q_{\psi}(\mathbf{q}^i|\tau^i)} \right)$$



- ▶ Approximating Distribution:

$$Q_{\phi, \psi}(\tau, q) \triangleq \overset{\text{tree topology}}{Q_{\phi}(\tau)} \cdot \overset{\text{branch length}}{Q_{\psi}(q|\tau)}$$

- ▶ Multi-sample Lower Bound:

$$L^K(\phi, \psi) = \mathbb{E}_{Q_{\phi, \psi}(\tau^{1:K}, \mathbf{q}^{1:K})} \log \left( \frac{1}{K} \sum_{i=1}^K \frac{p(\mathbf{Y}|\tau^i, \mathbf{q}^i)p(\tau^i, \mathbf{q}^i)}{Q_{\phi}(\tau^i)Q_{\psi}(\mathbf{q}^i|\tau^i)} \right)$$

- ▶ Variational Bayesian Phylogenetic Inference:

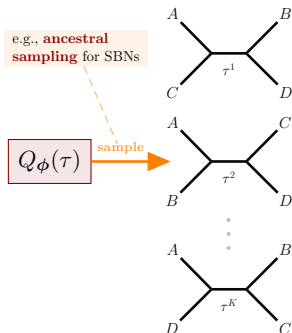
$$\hat{\phi}, \hat{\psi} = \arg \max_{\phi, \psi} L^K(\phi, \psi)$$

parameters trained via stochastic gradient ascent (SGA).

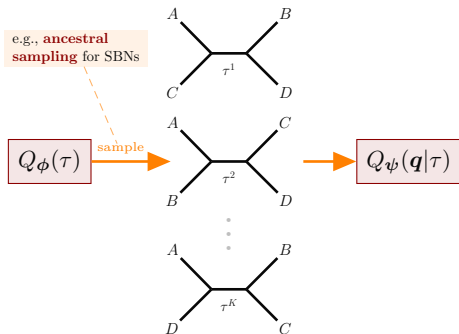


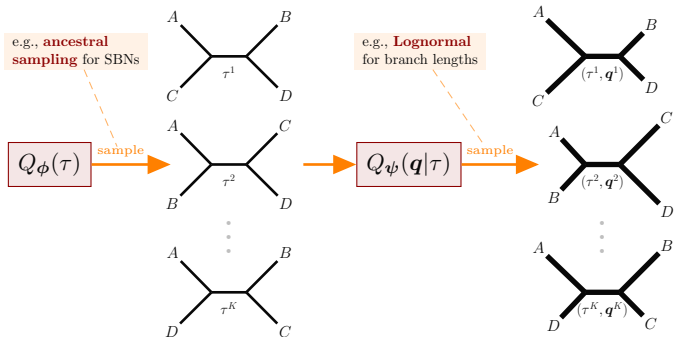
$$Q_{\phi}(\tau)$$

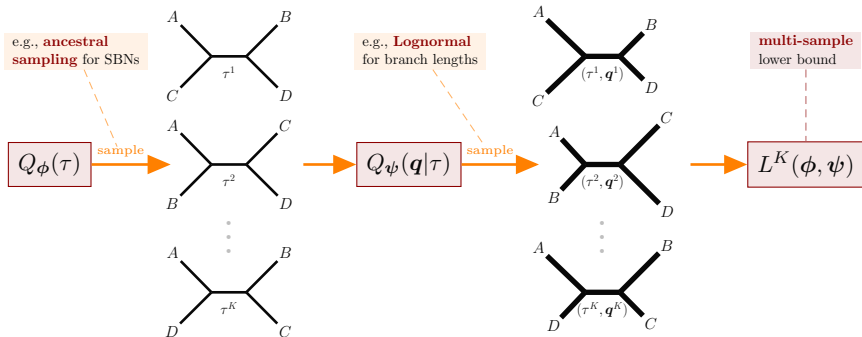


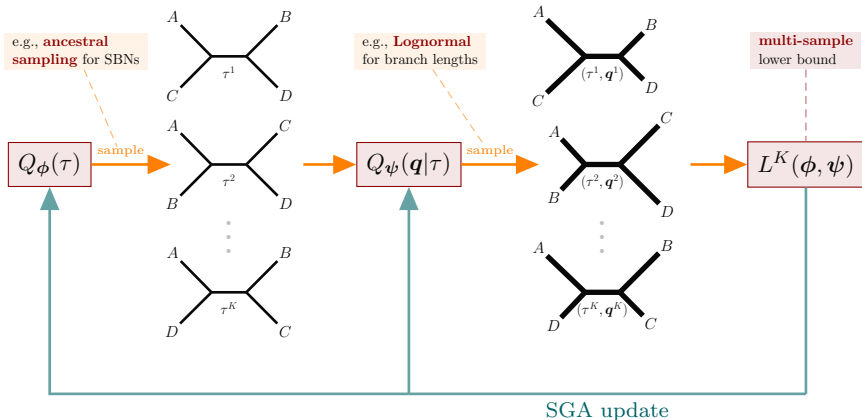












## SBNs Parameters

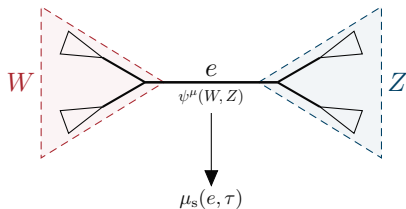
$$p(S_1 = s_1) = \frac{\exp(\phi_{s_1})}{\sum_{s_r \in \mathbb{S}_r} \exp(\phi_{s_r})}, \quad p(S_i = s | S_{\pi_i} = t) = \frac{\exp(\phi_{s|t})}{\sum_{s \in \mathbb{S}_{|t}} \exp(\phi_{s|t})}$$

## Branch Length Parameters

$$Q_{\psi}(\mathbf{q}|\tau) = \prod_{e \in E(\tau)} p^{\text{Lognormal}}(q_e | \mu(e, \tau), \sigma(e, \tau))$$

► *Simple Split*

$$\mu_s(e, \tau) = \psi_{e/\tau}^{\mu}, \quad \sigma_s(e, \tau) = \psi_{e/\tau}^{\sigma}$$



## SBNs Parameters

$$p(S_1 = s_1) = \frac{\exp(\phi_{s_1})}{\sum_{s_r \in \mathbb{S}_r} \exp(\phi_{s_r})}, \quad p(S_i = s | S_{\pi_i} = t) = \frac{\exp(\phi_{s|t})}{\sum_{s \in \mathbb{S}_{|t}} \exp(\phi_{s|t})}$$

## Branch Length Parameters

$$Q_{\psi}(\mathbf{q}|\tau) = \prod_{e \in E(\tau)} p^{\text{Lognormal}}(q_e | \mu(e, \tau), \sigma(e, \tau))$$

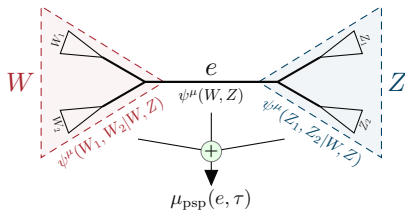
### ► *Simple Split*

$$\mu_s(e, \tau) = \psi_{e/\tau}^{\mu}, \quad \sigma_s(e, \tau) = \psi_{e/\tau}^{\sigma}$$

### ► *Primary Subsplit Pair (PSP)*

$$\mu_{\text{psp}}(e, \tau) = \psi_{e/\tau}^{\mu} + \sum_{s \in e//\tau} \psi_s^{\mu}$$

$$\sigma_{\text{psp}}(e, \tau) = \psi_{e/\tau}^{\sigma} + \sum_{s \in e//\tau} \psi_s^{\sigma}$$



SBNs Parameters  $\phi$ . With  $\tau^j, \mathbf{q}^j \stackrel{\text{iid}}{\sim} Q_{\phi, \psi}(\tau, \mathbf{q})$

- ▶ *VIMCO*. [Minh and Rezende, ICML 2016]

$$\nabla_{\phi} L^K(\phi, \psi) \simeq \sum_{j=1}^K \left( \hat{L}_{j|-j}^K(\phi, \psi) - \tilde{w}^j \right) \nabla_{\phi} \log Q_{\phi}(\tau^j).$$

- ▶ *RWS*. [Bornschein and Bengio, ICLR 2015]

$$\nabla_{\phi} L^K(\phi, \psi) \simeq \sum_{j=1}^K \tilde{w}^j \nabla_{\phi} \log Q_{\phi}(\tau^j).$$

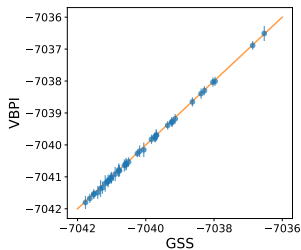
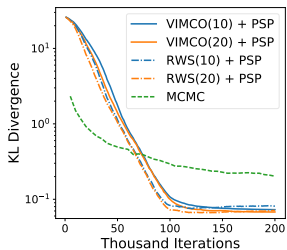
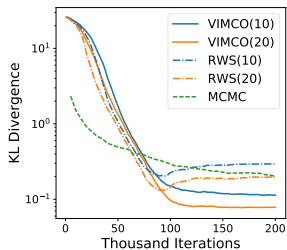
Branch Length Parameters  $\psi$ .  $g_{\psi}(\epsilon|\tau) = \exp(\boldsymbol{\mu}_{\psi, \tau} + \boldsymbol{\sigma}_{\psi, \tau} \odot \epsilon)$ .

- ▶ *Reparameterization Trick*. Let  $f_{\phi, \psi}(\tau, \mathbf{q}) = \frac{p(\mathbf{Y}|\tau, \mathbf{q})p(\tau, \mathbf{q})}{Q_{\phi}(\tau)Q_{\psi}(\mathbf{q}|\tau)}$ .

$$\nabla_{\psi} L^K(\phi, \psi) \simeq \sum_{j=1}^K \tilde{w}^j \nabla_{\psi} \log f_{\phi, \psi}(\tau^j, g_{\psi}(\epsilon^j|\tau^j))$$

where  $\tau^j \stackrel{\text{iid}}{\sim} Q_{\phi}(\tau)$ ,  $\epsilon^j \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

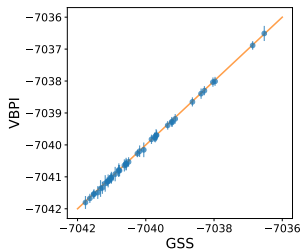
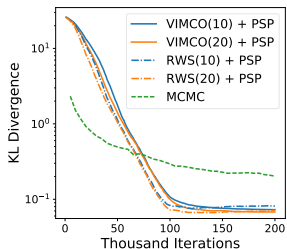
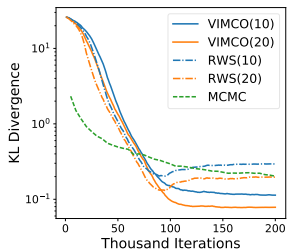




[Zhang and Matsen, ICLR 2019]



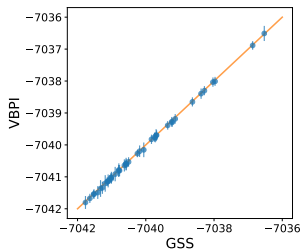
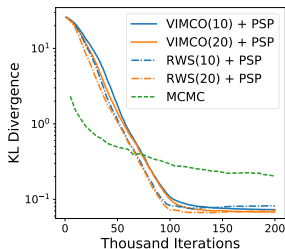
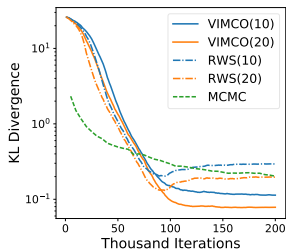




[Zhang and Matsen, ICLR 2019]

More samples  $\Rightarrow$  better exploration  $\Rightarrow$  better approximation

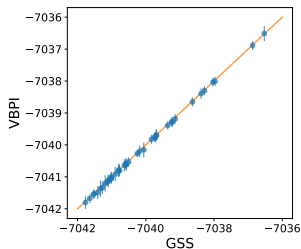
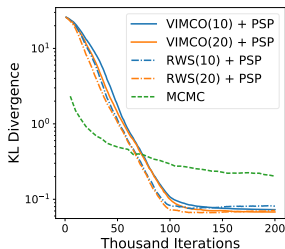
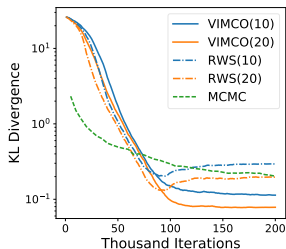




[Zhang and Matsen, ICLR 2019]

More samples  $\Rightarrow$  better exploration  $\Rightarrow$  better approximation

More flexible branch length distributions across tree topologies  
(PSP) ease training and improve approximation



[Zhang and Matsen, ICLR 2019]

More samples  $\Rightarrow$  better exploration  $\Rightarrow$  better approximation

More flexible branch length distributions across tree topologies  
(PSP) ease training and improve approximation

Outperform MCMC via much more efficient tree space  
exploration and branch length updates



DATA SET	MARGINAL LIKELIHOOD (NATs)				
	VIMCO(10)	VIMCO(20)	VIMCO(10)+PSP	VIMCO(20)+PSP	SS
DS1	-7108.43(0.26)	-7108.35(0.21)	-7108.41(0.16)	<b>-7108.42(0.10)</b>	-7108.42(0.18)
DS2	-26367.70(0.12)	-26367.71(0.09)	<b>-26367.72(0.08)</b>	-26367.70(0.10)	-26367.57(0.48)
DS3	-33735.08(0.11)	-33735.11(0.11)	<b>-33735.10(0.09)</b>	-33735.07(0.11)	-33735.44(0.50)
DS4	-13329.90(0.31)	-13329.98(0.20)	<b>-13329.94(0.18)</b>	-13329.93(0.22)	-13330.06(0.54)
DS5	-8214.36(0.67)	-8214.74(0.38)	-8214.61(0.38)	-8214.55(0.43)	<b>-8214.51(0.28)</b>
DS6	-6723.75(0.68)	-6723.71(0.65)	-6724.09(0.55)	<b>-6724.34(0.45)</b>	-6724.07(0.86)
DS7	-37332.03(0.43)	-37331.90(0.49)	-37331.90(0.32)	<b>-37332.03(0.23)</b>	-37332.76(2.42)
DS8	-8653.34(0.55)	-8651.54(0.80)	<b>-8650.63(0.42)</b>	-8650.55(0.46)	-8649.88(1.75)

[Zhang and Matsen, ICLR 2019]

DATA SET	MARGINAL LIKELIHOOD (NATs)				
	VIMCO(10)	VIMCO(20)	VIMCO(10)+PSP	VIMCO(20)+PSP	SS
DS1	-7108.43(0.26)	-7108.35(0.21)	-7108.41(0.16)	<b>-7108.42(0.10)</b>	-7108.42(0.18)
DS2	-26367.70(0.12)	-26367.71(0.09)	<b>-26367.72(0.08)</b>	-26367.70(0.10)	-26367.57(0.48)
DS3	-33735.08(0.11)	-33735.11(0.11)	<b>-33735.10(0.09)</b>	-33735.07(0.11)	-33735.44(0.50)
DS4	-13329.90(0.31)	-13329.98(0.20)	<b>-13329.94(0.18)</b>	-13329.93(0.22)	-13330.06(0.54)
DS5	-8214.36(0.67)	-8214.74(0.38)	-8214.61(0.38)	-8214.55(0.43)	<b>-8214.51(0.28)</b>
DS6	-6723.75(0.68)	-6723.71(0.65)	-6724.09(0.55)	<b>-6724.34(0.45)</b>	-6724.07(0.86)
DS7	-37332.03(0.43)	-37331.90(0.49)	-37331.90(0.32)	<b>-37332.03(0.23)</b>	-37332.76(2.42)
DS8	-8653.34(0.55)	-8651.54(0.80)	<b>-8650.63(0.42)</b>	-8650.55(0.46)	-8649.88(1.75)

[Zhang and Matsen, ICLR 2019]

Competitive to state-of-the-art (stepping-stone), dramatically reducing cost at test time: VBPI(1000) vs SS(100,000)

DATA SET	MARGINAL LIKELIHOOD (NATs)				
	VIMCO(10)	VIMCO(20)	VIMCO(10)+PSP	VIMCO(20)+PSP	SS
DS1	-7108.43(0.26)	-7108.35(0.21)	-7108.41(0.16)	<b>-7108.42(0.10)</b>	-7108.42(0.18)
DS2	-26367.70(0.12)	-26367.71(0.09)	<b>-26367.72(0.08)</b>	-26367.70(0.10)	-26367.57(0.48)
DS3	-33735.08(0.11)	-33735.11(0.11)	<b>-33735.10(0.09)</b>	-33735.07(0.11)	-33735.44(0.50)
DS4	-13329.90(0.31)	-13329.98(0.20)	<b>-13329.94(0.18)</b>	-13329.93(0.22)	-13330.06(0.54)
DS5	-8214.36(0.67)	-8214.74(0.38)	-8214.61(0.38)	-8214.55(0.43)	<b>-8214.51(0.28)</b>
DS6	-6723.75(0.68)	-6723.71(0.65)	-6724.09(0.55)	<b>-6724.34(0.45)</b>	-6724.07(0.86)
DS7	-37332.03(0.43)	-37331.90(0.49)	-37331.90(0.32)	<b>-37332.03(0.23)</b>	-37332.76(2.42)
DS8	-8653.34(0.55)	-8651.54(0.80)	<b>-8650.63(0.42)</b>	-8650.55(0.46)	-8649.88(1.75)

[Zhang and Matsen, ICLR 2019]

Competitive to state-of-the-art (stepping-stone), dramatically reducing cost at test time: **VBPI(1000)** vs **SS(100,000)**

**PSP** alleviates the demand for large samples, reducing computation while maintaining approximation accuracy



- ▶ We introduced **VBPI**, a general **variational framework** for Bayesian phylogenetic inference.
- ▶ **VBPI** allows **efficient learning** on both **tree topology** and **branch lengths**, providing competitive performance to MCMC while requiring **much less computation**.
- ▶ Can be used for further statistical analysis (e.g., marginal likelihood estimation) via **importance sampling**.



- [1] **Zhang, C.** and Matsen F. A., **Generalizing Tree Probability Estimation via Bayesian Networks.** In *Advances in Neural Information Processing Systems*, **spotlight(3.5%)**, 2018.
- [2] **Zhang, C.** and Matsen F. A., **Variational Bayesian Phylogenetic Inference.** In *Proceedings of the 7th International Conference on Learning Representations*, 2019.

**Thank you!**

