

NAS-X: Neural Adaptive Smoothing via Twisting

Dieterich Lawson*, Michael Y. Li*, Scott W. Linderman

dieterichl@google.com, michaelyli@stanford.edu, scott.linderman@stanford.edu



Summary

Fitting **sequential latent variable models** is hard!

Why? High-variance gradients and discrete latents!

NAS-X = smoothing + reweighted wake sleep

- Low variance and low bias gradients.
- Versatile - handles discrete latents.
- Easy to train, minimal computational overhead.
- Significantly outperforms prior methods.

Reweighted Wake Sleep

Fisher's identity: Gradient of log marginal likelihood is a posterior expectation.

$$\nabla_{\theta} \log p_{\theta}(\mathbf{y}_{1:T}) = \mathbb{E}_{p_{\theta}(\mathbf{x}_{1:T}|\mathbf{y}_{1:T})} [\nabla_{\theta} \log p_{\theta}(\mathbf{x}_{1:T}, \mathbf{y}_{1:T})]$$

RWS estimates expectation via **importance sampling**.

$$\sum_{i=1}^N \bar{w}^{(i)} \nabla_{\theta} \log p_{\theta}(\mathbf{x}_{1:T}^{(i)}, \mathbf{y}_{1:T}), \quad \mathbf{x}_{1:T}^{(i)} \sim q_{\phi}(\mathbf{x}_{1:T} | \mathbf{y}_{1:T})$$

$$\bar{w}^{(i)} \propto \frac{p_{\theta}(\mathbf{x}_{1:T}^{(i)}, \mathbf{y}_{1:T})}{q_{\phi}(\mathbf{x}_{1:T}^{(i)} | \mathbf{y}_{1:T})}$$

Estimating Posterior Expectations via Smoothing Sequential Monte Carlo

NAS-X estimates with **smoothing SMC with twists** r . Twists learned via density ratio estimation.

$$\sum_{t=1}^T \sum_{i=1}^N \bar{w}_t^{(i)} \nabla_{\theta} \log p_{\theta}(\mathbf{x}_t^{(i)}, \mathbf{y}_t | \mathbf{x}_{t-1}^{(i)})$$

$$\mathbf{x}_{1:T}^{1:N}, \bar{w}_{1:T}^{1:N} \leftarrow \text{SMC}(\{p_{\theta}(\mathbf{x}_{1:t}, \mathbf{y}_{1:t}), q_{\phi}(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y}_{t:T}), r_{\psi}(\mathbf{x}_t, \mathbf{y}_{t+1:T})\}_{t=1}^T)$$

Under certain conditions, NAS-X has **unbiased** and **consistent** gradient estimates.

Algorithm

Algorithm 1: NAS-X

```

Procedure NAS-X( $\theta_0, \phi_0, \psi_0, \mathbf{y}_{1:T}$ )
   $\theta \leftarrow \theta_0, \phi \leftarrow \phi_0, \psi \leftarrow \psi_0$ 
  while not converged do
     $\mathbf{x}_{1:T}^{1:N}, \bar{w}_{1:T}^{1:N} \leftarrow \text{SMC}(\{p_{\theta}(\mathbf{x}_{1:t}, \mathbf{y}_{1:t}), q_{\phi}(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y}_{t:T}), r_{\psi}(\mathbf{x}_t, \mathbf{y}_{t+1:T})\}_{t=1}^T)$ 
     $\Delta\theta = \sum_{t=1}^T \sum_{i=1}^N \bar{w}_t^{(i)} \nabla_{\theta} \log p_{\theta}(\mathbf{x}_t^{(i)}, \mathbf{y}_t | \mathbf{x}_{t-1}^{(i)})$ 
     $\Delta\phi = -\sum_{t=1}^T \sum_{i=1}^N \bar{w}_t^{(i)} \nabla_{\phi} \log q_{\phi}(\mathbf{x}_t^{(i)} | \mathbf{x}_{t-1}^{(i)}, \mathbf{y}_{t:T})$ 
     $\theta \leftarrow \text{grad-step}(\theta, \Delta\theta)$ 
     $\phi \leftarrow \text{grad-step}(\phi, \Delta\phi)$ 
     $\psi \leftarrow \text{twist-training}(\theta, \psi)$ 
  end
  return  $\theta, \phi, \psi$ 
Procedure twist-training( $\theta, \psi_0$ )
  | See Algorithm 2 in Appendix 8.3.
    
```

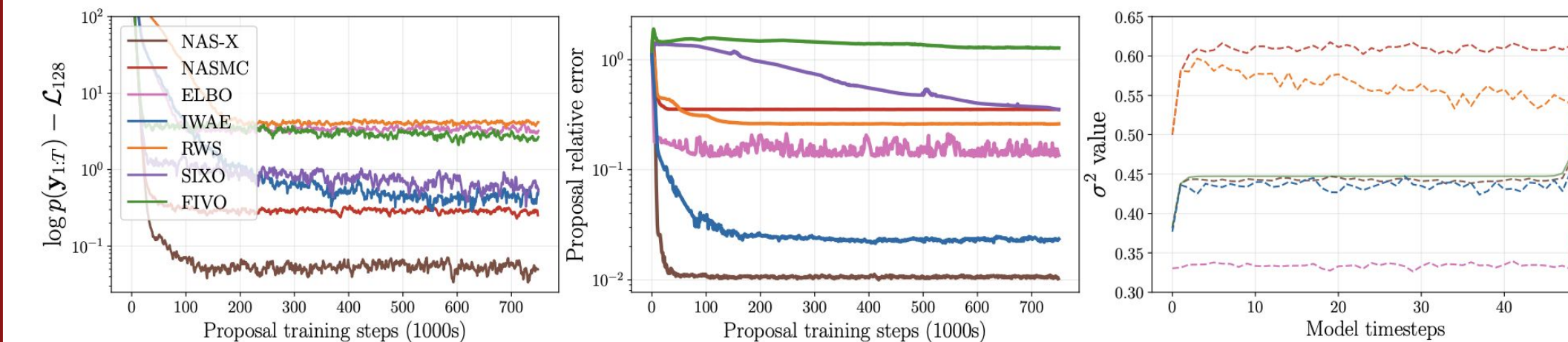
Paper



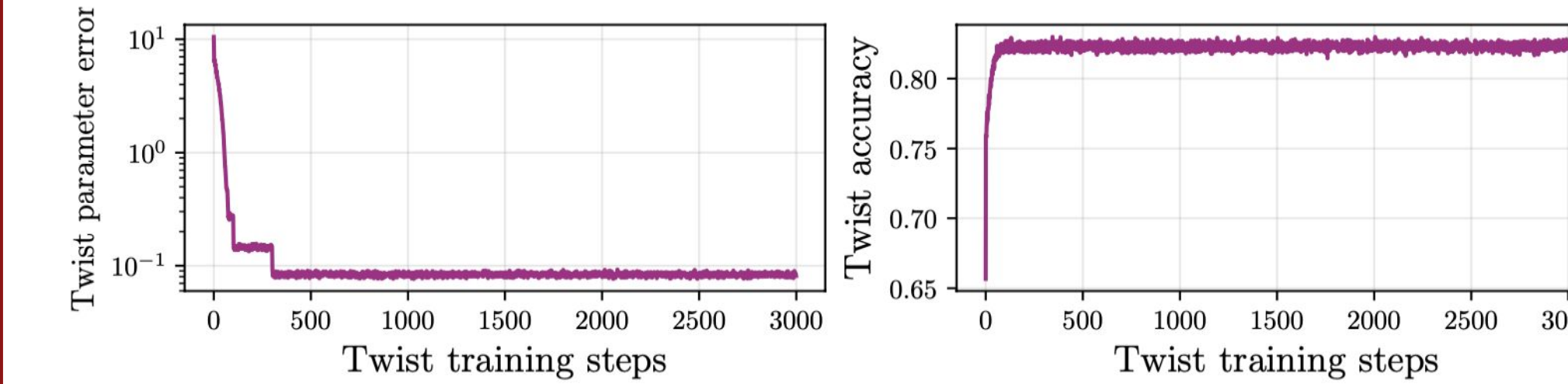
References

- Dieterich Lawson, Allan Reventos, Andrew Warrington, and Scott Linderman. SIXO: Smoothing inference with twisted objectives. *Advances in Neural Information Processing Systems*, 2022.
- Shiwan Gu, Zoubin Ghahramani, Richard E. Turner. Neural Adaptive Sequential Monte Carlo. *Advances in Neural Information Processing Systems*, 2015.
- Alan L. Hodgkin and Andrew F. Huxley. A quantitative description of membrane current and its application to conduction and excitation in nerve. *The Journal of Physiology*, 107(4):500, 1952.
- Jörg Bornschein and Yoshua Bengio. Reweighted wake-sleep. arXiv preprint arXiv:1406.2751, 2014.

Linear Gaussian SSM

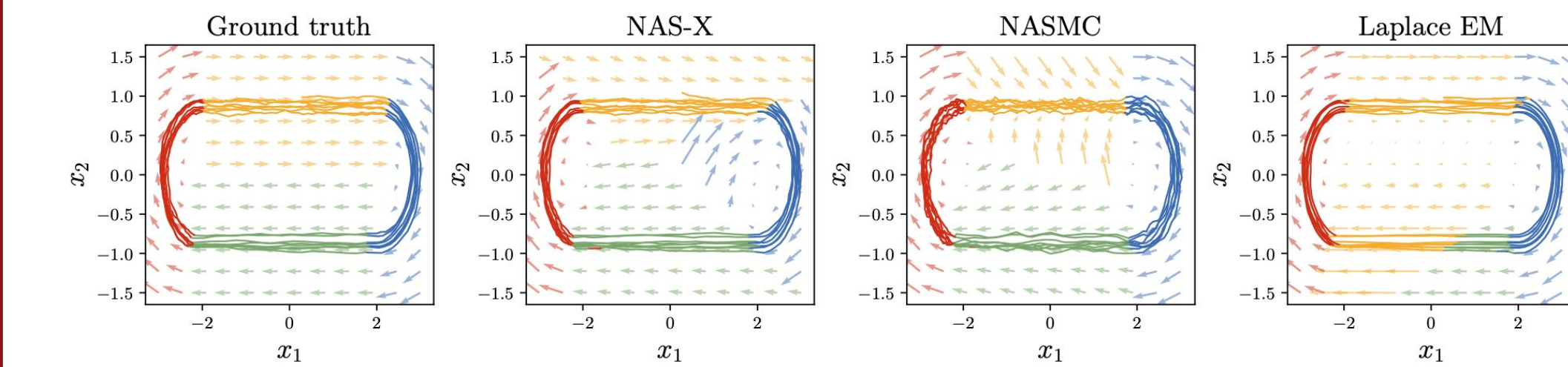


NAS-X has tighter variational bound, and recovers true posterior.



NAS-X recovers the true twist parameters.

Discrete latent variables



NAS-X can handle discrete latents.

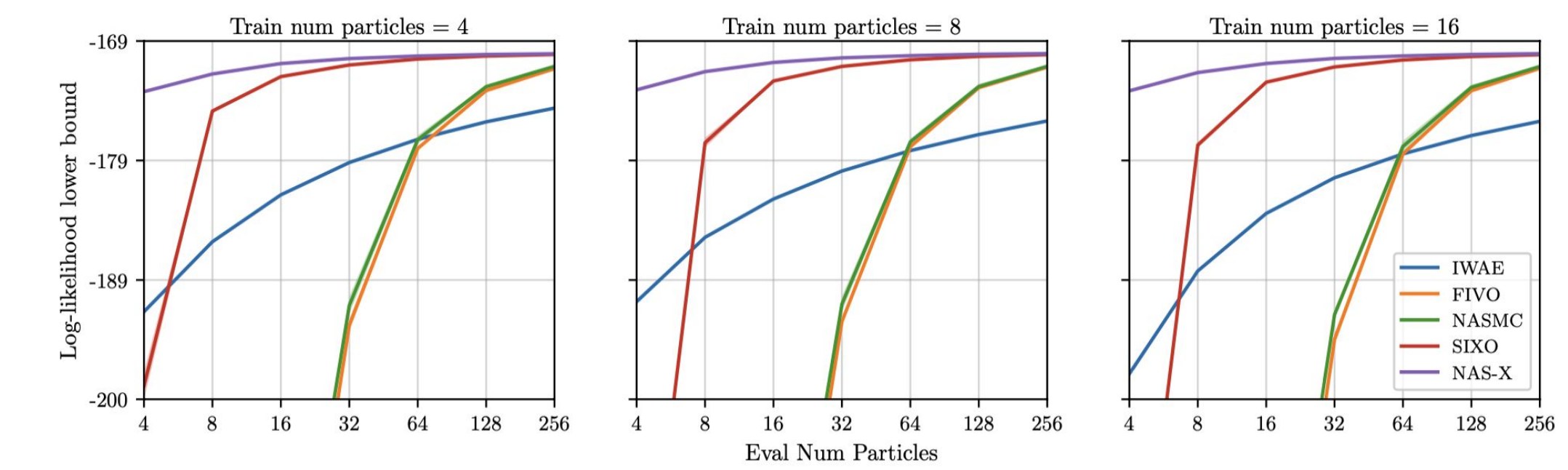
(b) Train $\mathcal{L}_{\text{BPF}}^{1024}$ for rSLDS.

Method	$\sigma_O^2 = 0.001$	$\sigma_O^2 = 0.01$	$\sigma_O^2 = 0.1$
NAS-X	19.837 ± 0.0234	8.63 ± 0.0015	-2.79 ± 0.0009
NASMC	19.834 ± 0.0018	8.53 ± 0.001	-2.874 ± 0.0007
Laplace EM	19.154 ± 0.057	8.54 ± 0.0039	-2.765 ± 0.0012
RWS	17.148 ± 0.087	6.314 ± 0.023	-5.78 ± 0.0026

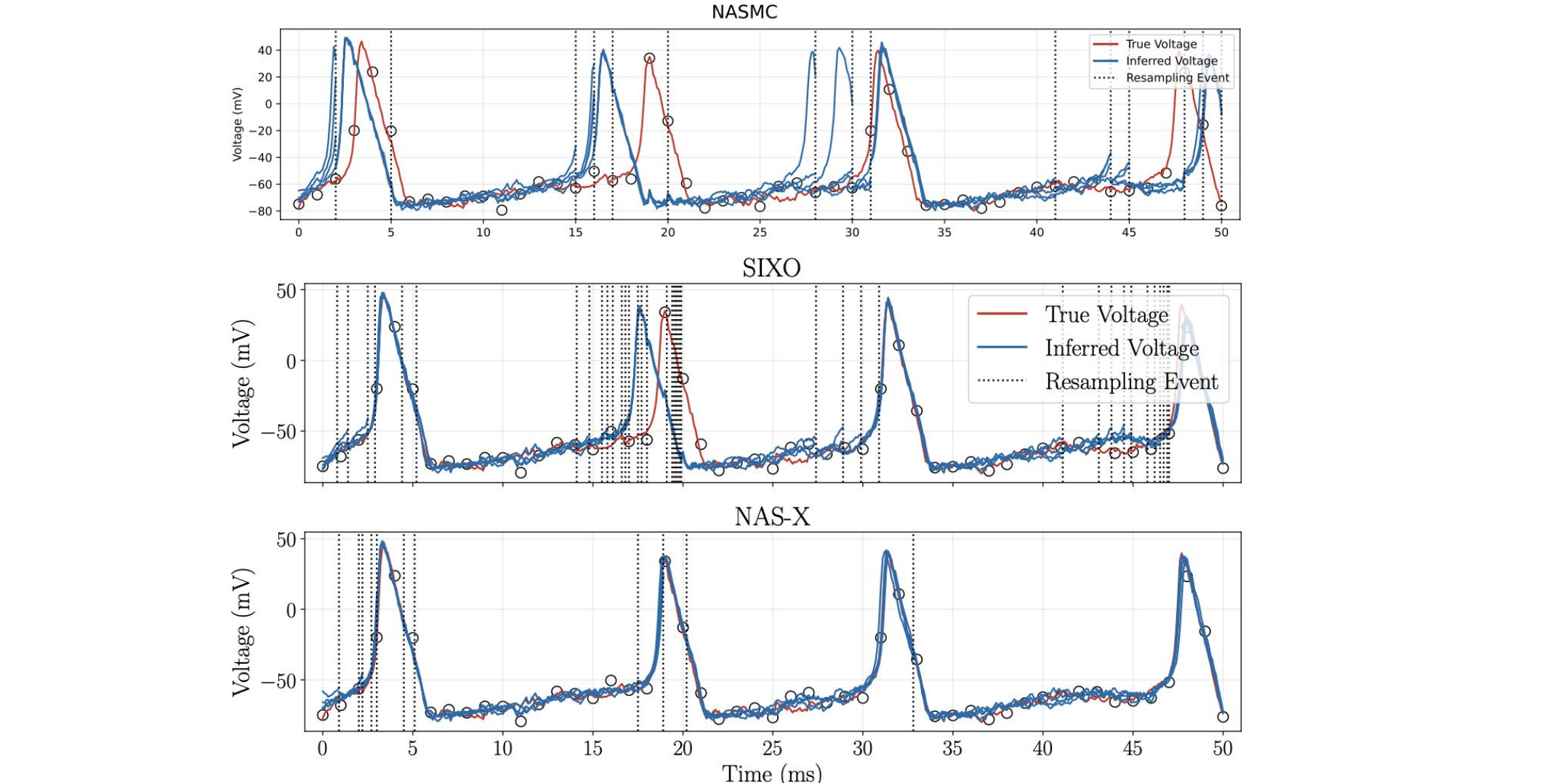
By smoothing, NAS-X learns better models than RWS methods.

Hodgkin Huxley model

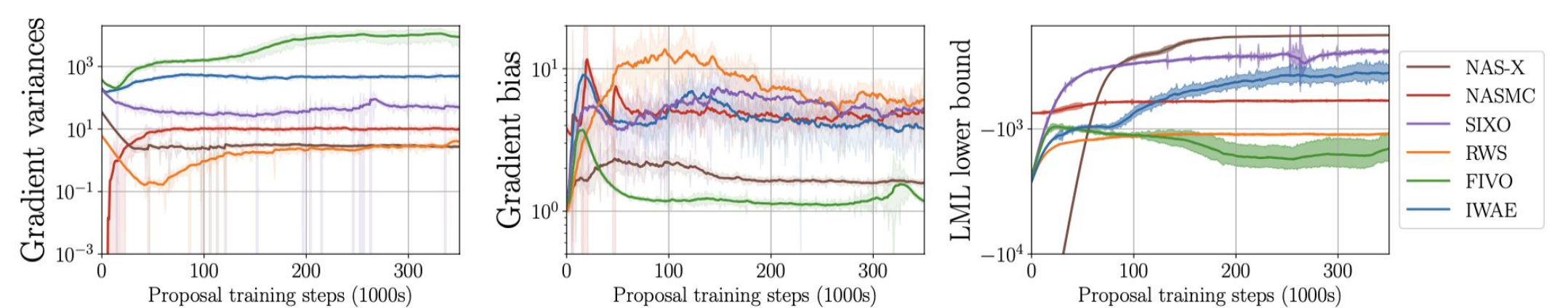
Hodgkin-Huxley, a mechanistic model of neural dynamics.



NAS-X is 4-64x more particle-efficient than prior methods.



NAS-X perfectly infers the latent voltage.



NAS-X has lower variance and lower bias gradients.