

Motivation

- Hyperparameters such as architecture choice, data augmentation, and dropout are crucial for neural net generalization, but **difficult to tune**.
- Grid search, random search, and Bayesian optimization treat hyperparameter tuning as a **black-box** problem, which does not scale to high-dimensional hyperparameters.
- Hyperparameter tuning is a **bilevel optimization problem**:

$$\lambda^* = \arg \min_{\lambda} \mathcal{L}_V(\lambda, \mathbf{w}^*(\lambda)) \quad \text{s.t.} \quad \mathbf{w}^*(\lambda) = \arg \min_{\mathbf{w}} \mathcal{L}_T(\lambda, \mathbf{w})$$
- We approximate the best-response function $\mathbf{w}^*(\lambda)$ with a hypernetwork $\mathbf{w}_{\phi}(\lambda)$, called a Self-Tuning Network (STN).

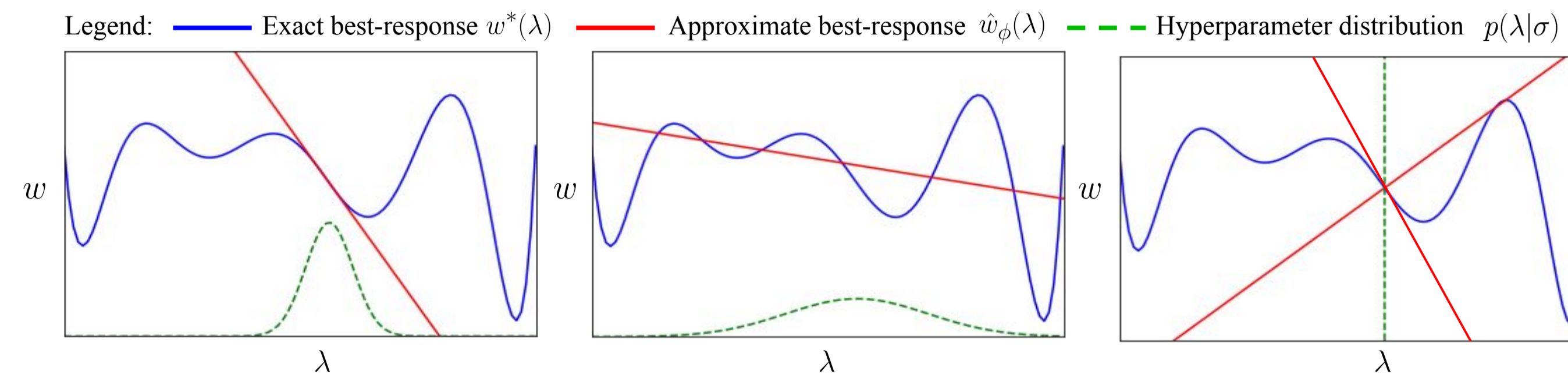
Summary

- We propose a compact architecture for approximating neural net best-responses, that can be used as a **drop-in replacement** for existing deep learning modules.
- Our training algorithm alternates between approximating the best-response around the current hyperparameters and optimizing the hyperparameters with the approximate best-response.
- This yields a gradient-based algorithm that is (1) **computationally inexpensive**, (2) can optimize all regularization hyperparameters **including discrete hyperparameters**, and (3) **scales to large NNs**.
- Our approach discovers **hyperparameter trajectories** that can outperform fixed hyperparameter values.

Self-Tuning Network (STN) Training Algorithm

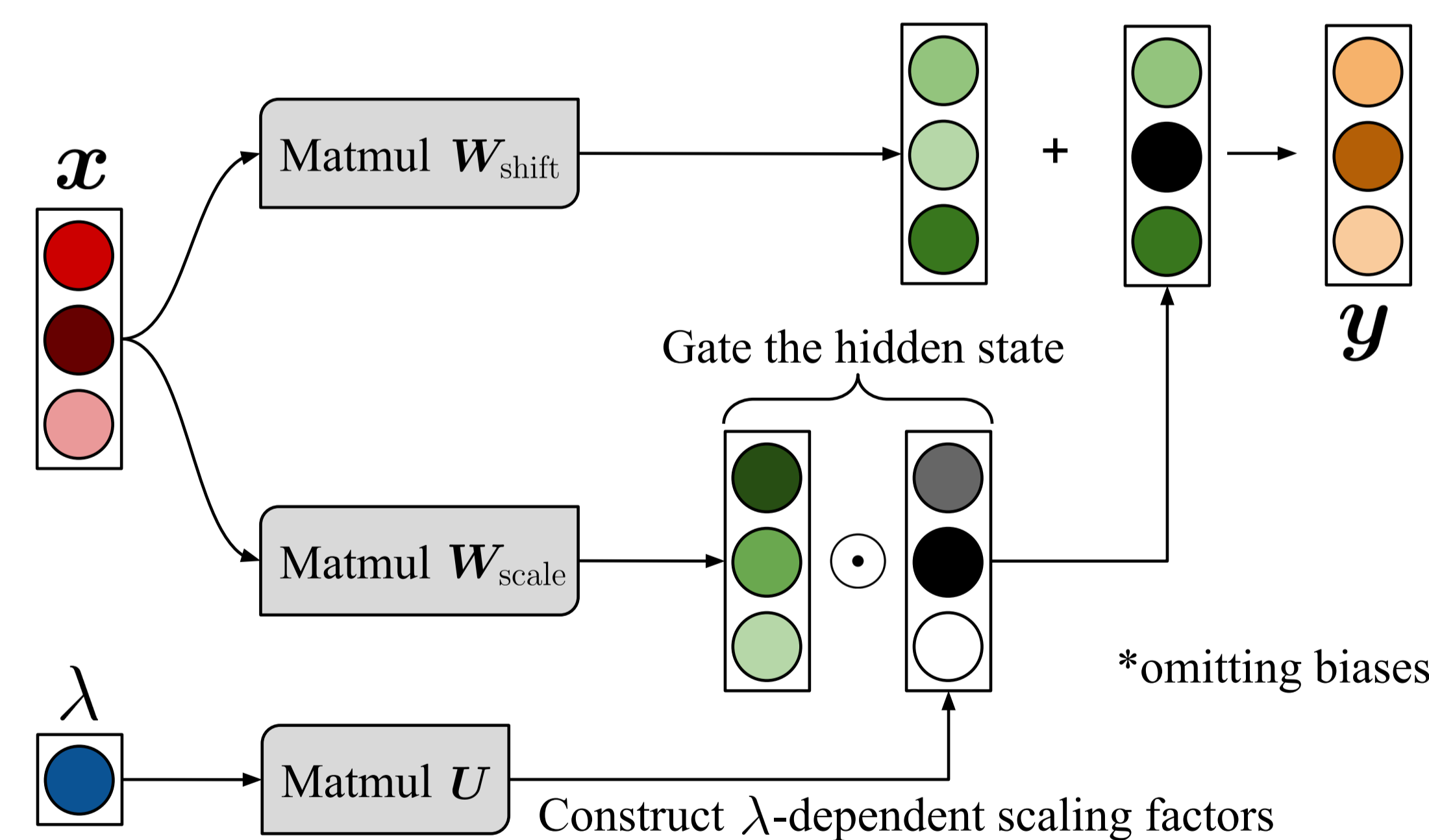
Initialize: Hypernetwork parameters ϕ , hyperparameters λ
while not converged **do**
 for $t = 1, \dots, T_{train}$ **do**
 $\epsilon \sim p(\epsilon|\sigma)$
 $\phi \leftarrow \phi - \alpha_1 \nabla_{\phi} [\mathcal{L}_T(\lambda + \epsilon, \hat{\mathbf{w}}_{\phi}(\lambda + \epsilon))]$
 for $t = 1, \dots, T_{valid}$ **do**
 $\epsilon \sim p(\epsilon|\sigma)$
 $\hat{\mathcal{L}}_V(\lambda, \sigma) \leftarrow \mathcal{L}_V(\lambda + \epsilon, \hat{\mathbf{w}}_{\phi}(\lambda + \epsilon)) - \tau \mathbb{H}[p(\epsilon|\sigma)]$
 $(\lambda, \sigma) \leftarrow (\lambda, \sigma) - \alpha_2 \nabla_{\lambda, \sigma} [\hat{\mathcal{L}}_V(\lambda, \sigma)]$

Sampling Hyperparameters



- Just right* \rightarrow the **gradient of the approximation** will match that of the best-response.
- Too wide* \rightarrow the hypernetwork may be **insufficiently flexible** to model the best-response, and the gradients will not match.
- Too small* \rightarrow the hypernetwork will match the best-response at the current hyperparameter, but **may not be locally correct**.
- We re-parameterize the hyperparameter λ to lie in \mathbb{R} and use noise distribution $p(\epsilon|\sigma) = \mathcal{N}(0, \sigma)$.

Hypernetwork Best-Response Architecture



- We **scale and shift the network's hidden units** (\equiv the rows of weights and biases) by an amount which depends on our hyperparameters:

$$\hat{W}_{\phi}(\lambda) = W_{\text{shift}} + (U\lambda) \odot_{\text{row}} W_{\text{scale}}$$

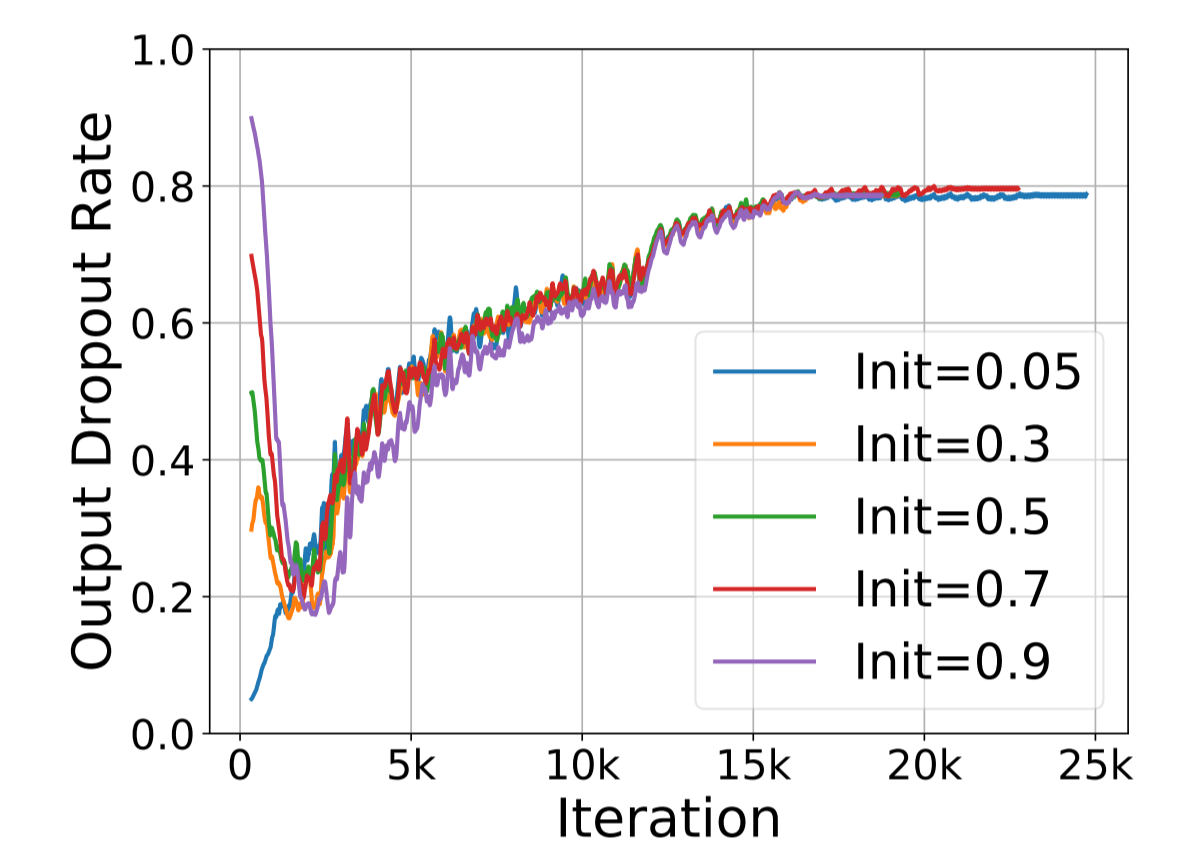
$$\hat{b}_{\phi}(\lambda) = b_{\text{shift}} + (C\lambda) \odot_{\text{row}} b_{\text{scale}}$$

- Memory efficient** (roughly 2x no. of parameters) and scales well to high dimensions.

Hyperparameter Trajectories

- STNs discover hyperparameter trajectories which can **outperform fixed hyperparameters**.
- For a single dropout rate, STNs implement a curriculum with a gradually increasing dropout probability.
- The same trajectory is followed **regardless of the initial hyperparameter value**.

Method	Val	Test
$p = 0.68$ (Fixed)	85.83	83.19
$p \sim \mathcal{N}(0.68, \sigma = 0.05)$	85.87	82.29
$p = 0.68 + 0.1 \sin(k\pi)$	85.29	82.15
$p = 0.78$ (Converged STN)	89.65	86.90
STN (Ours)	82.58	79.02
Following STN Trajectory	82.87	79.93



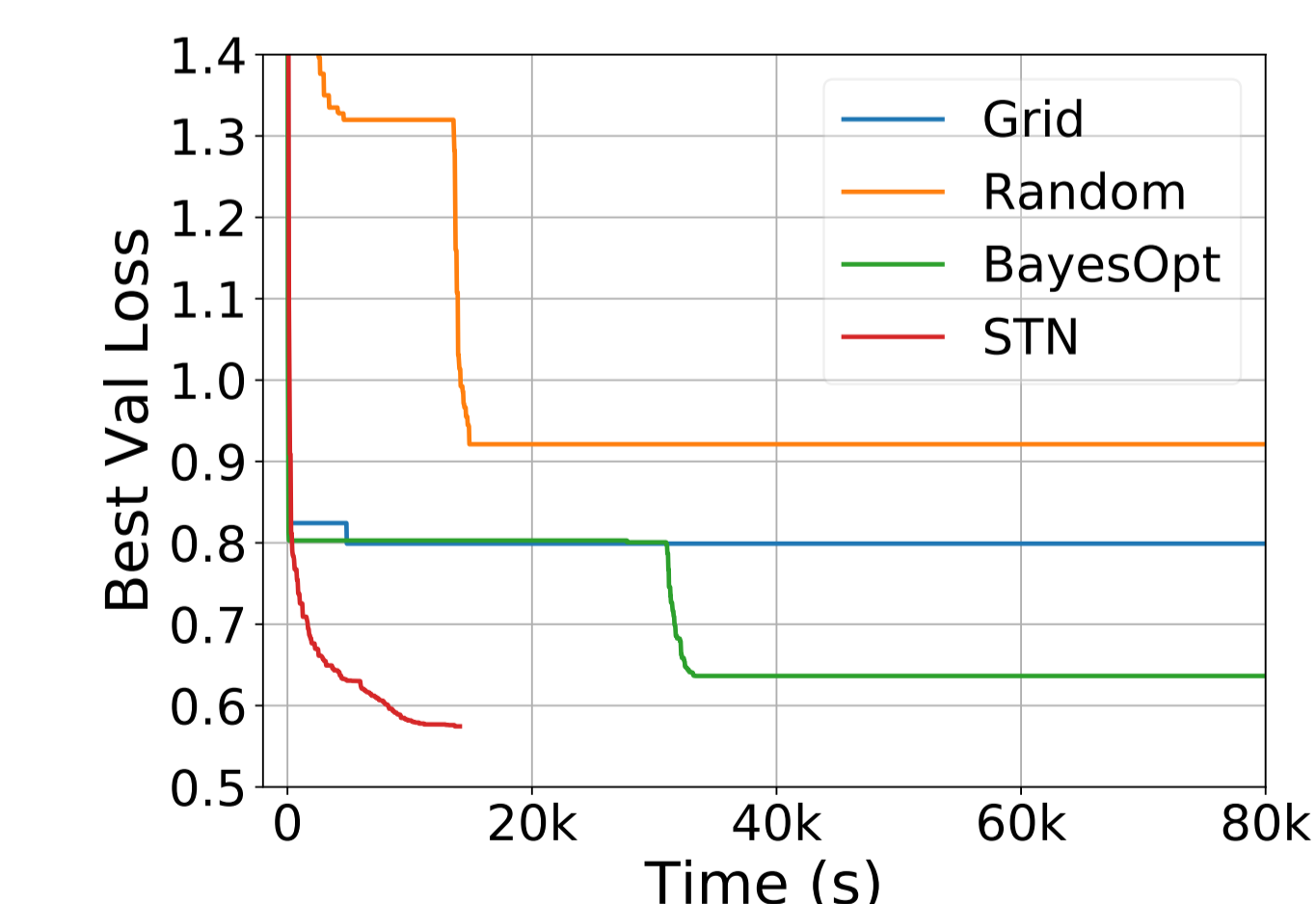
Comparing hyperparameter trajectories

Real-World Datasets

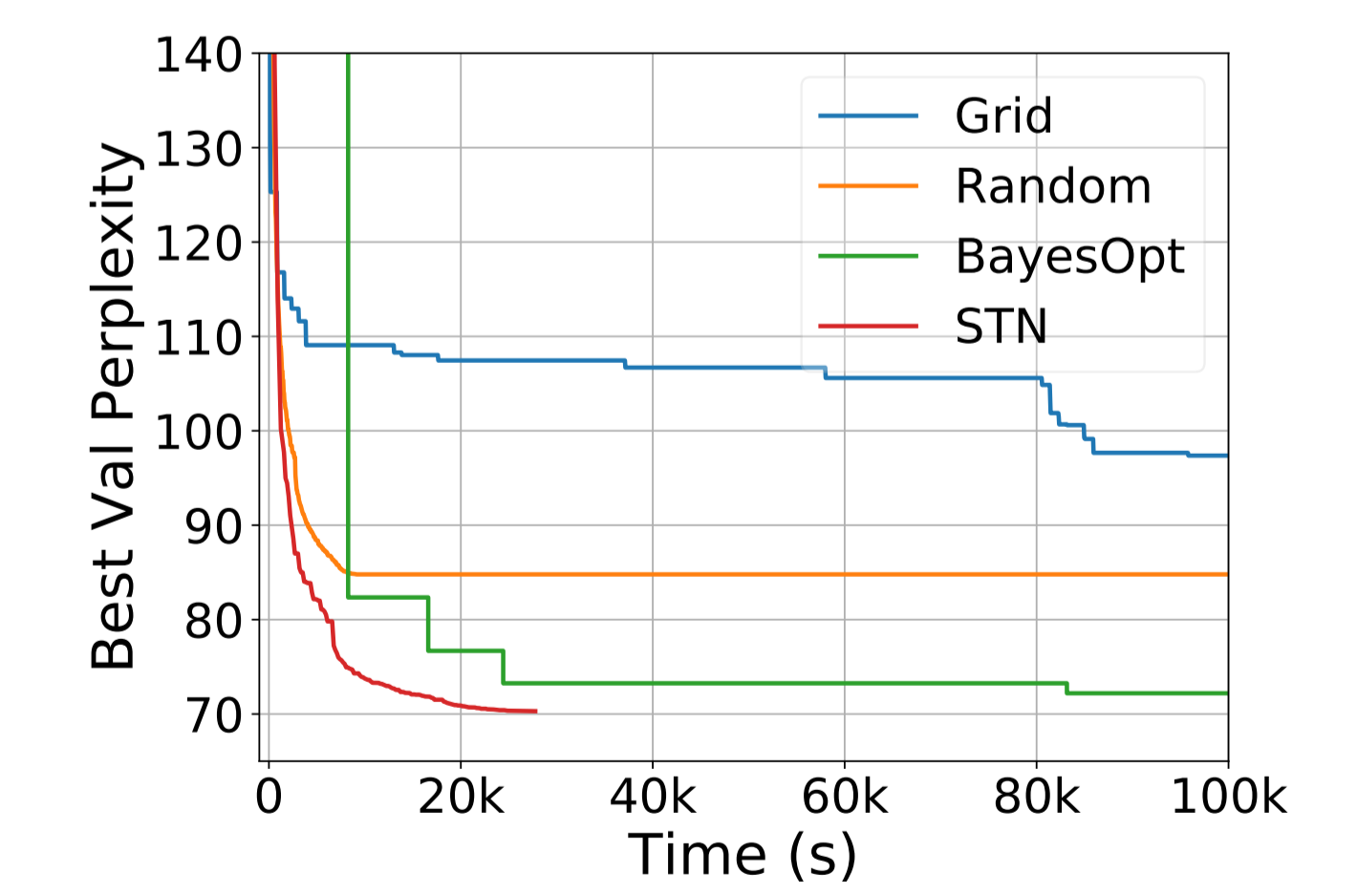
Method	PTB		CIFAR-10	
	Val Perplexity	Test Perplexity	Val Loss	Test Loss
Grid Search	97.32	94.58	0.794	0.809
Random Search	84.81	81.46	0.921	0.752
Bayesian Optimization	72.13	69.29	0.636	0.651
STN	70.30	67.68	0.575	0.576

Final validation/test performance on PTB and CIFAR-10

- CNN time comparison



- LSTM time comparison



- Hyperparameter trajectories for LSTM tuning

