

Good Initializations of Variational Bayes for Deep Models

Simone Rossi, Pietro Michiardi, Maurizio Filippone

Objectives and Contributions

Initialization of variational parameters has a huge role in the convergence of stochastic variational inference but received little to no attention in current literature.

Contributions:

- ▶ **New initialization** for svi based on Bayesian linear models;
- ▶ Applied to **regression, classification** and **CNNs**;
- ▶ Experimental comparison against other initializations;
- ▶ SoTA performance with Gaussian svi on large-scale CNNs.

Stochastic Variational Inference - svi

A DNN is a composition of nonlinear vector-valued functions $f^{(l)}$

$$f(\mathbf{x}) = \left(f^{(L-1)}(\mathbf{W}^{(L-1)}) \circ \dots \circ f^{(0)}(\mathbf{W}^{(0)}) \right) (\mathbf{x})$$

Posterior over the weights **Intractable for DNNs** Prior on model parameters

Objective of Bayesian inference

$$p(\mathbf{W}|X, Y) = \frac{p(Y|X, \mathbf{W})p(\mathbf{W})}{p(Y|X)}$$

Marginal Likelihood

svi reformulates this problem as minimization of the **negative evidence lower bound** (or **NELBO**) under an approximate distribution $q_{\theta}(\mathbf{W})$ [2]:

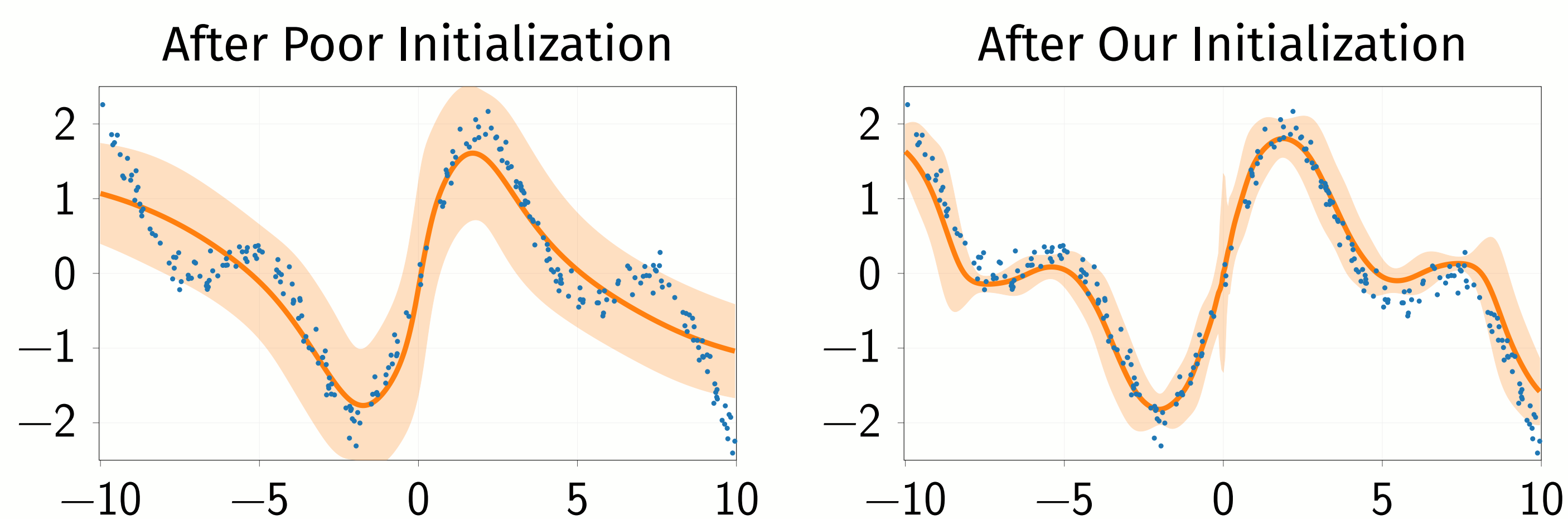
$$q_{\tilde{\theta}}(\mathbf{W}) \text{ s.t. } \tilde{\theta} = \arg \min_{\theta} \{ \text{NELBO} \}$$

$$\text{NELBO} = \mathbb{E}_{q_{\theta}} [-\log p(Y|X, \mathbf{W})] + \text{KL}(q_{\theta}(\mathbf{W}) || p(\mathbf{W}))$$

Commonly used family of variational distribution: **mean field Gaussians**

$$q(\mathbf{W}^{(l)}) = \prod_{ij} \mathcal{N}(w_{ij}^{(l)} | \mu_{ij}^{(l)}, \sigma_{ij}^{(l)}) \quad \theta = \{ \{ \mu_{ij}^{(l)}, \sigma_{ij}^{(l)} \} : l = 0, \dots, L-1 \}$$

How do we initialize θ ?



Iterative Bayesian Linear Modeling Initializer - I-BLM

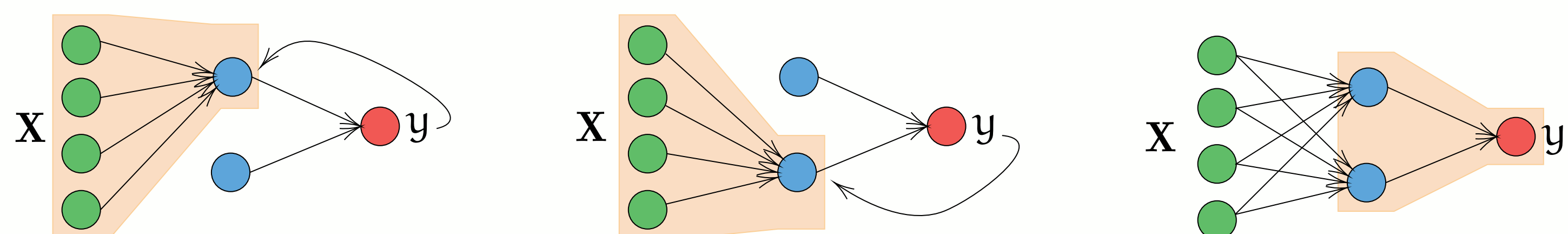


Figure: Representation of I-BLM. In (left) and (center) we learn two Bayesian linear models, whose outputs are used on the (right) for the following layer.

In a nutshell:

- ▶ Inspired by **residual networks** and **greedy initialization** of DNNs.
- ▶ Grounded on **Bayesian Linear regression** but extended to classification and to convolutional layers.
- ▶ **Regression on transformed labels** obtained through the interpretation of classification labels as the coefficients of a degenerate Dirichlet distribution.
- ▶ **Scalability** achieved thanks to mini-batching.

But how does it work?

Transform the labels if it's a classification task [3].

For each layer (l):

- ▶ Propagate a mini-batch of X up to the previous layer ($l-1$);
- ▶ Extract the patches if it's a convolutional layer;
- ▶ Learn a Bayesian linear model and use its solution to initialize $q_{\theta}(\mathbf{W}^{(l)})$.

Bayesian Linear Regression - BLR

Likelihood:

$$p(Y|W, L) = \prod_i \mathcal{N}(Y_i | X_i W_i, L)$$

Prior:

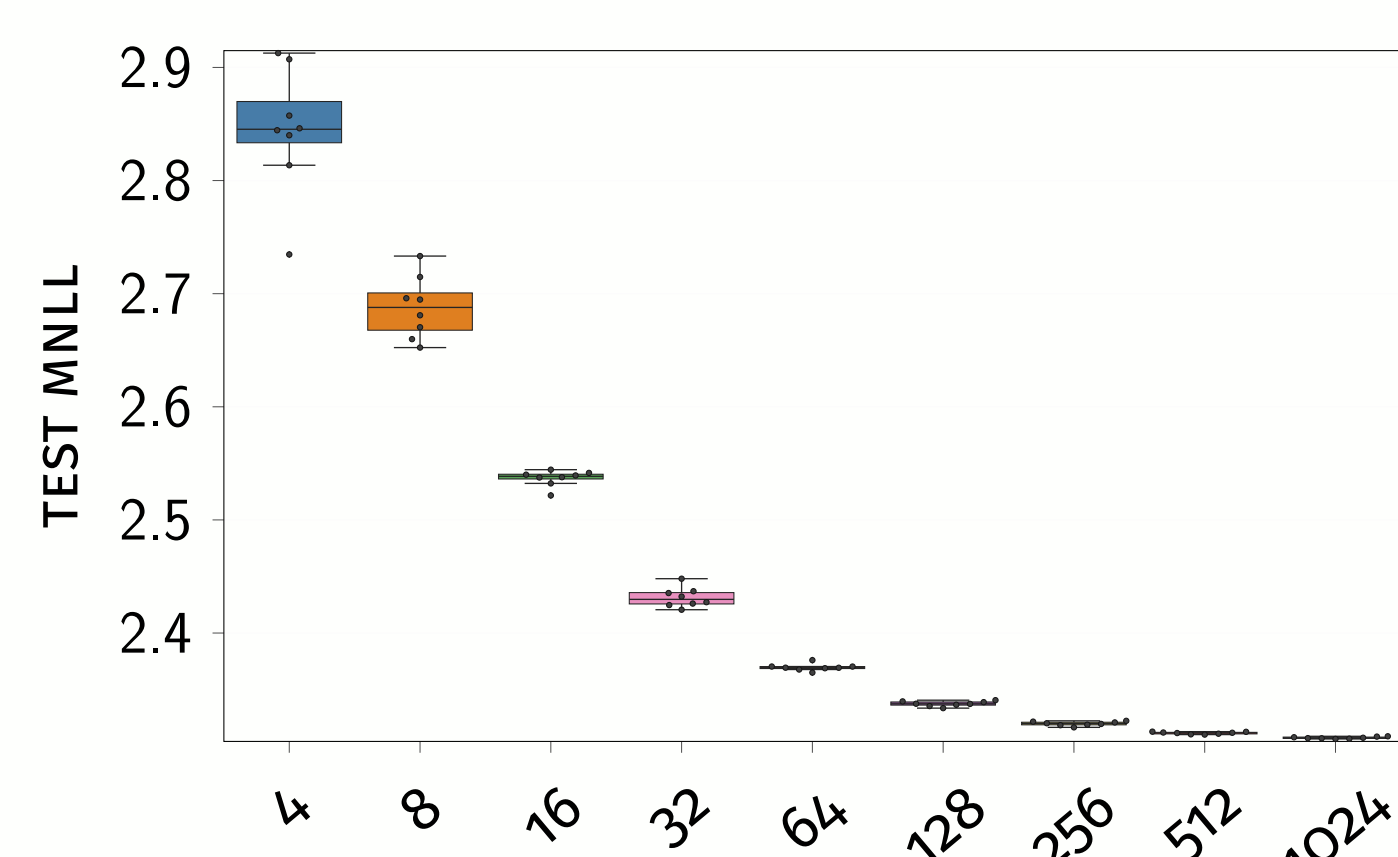
$$p(W|\Lambda) = \prod_i p(W_i) = \mathcal{N}(W_i | \mathbf{0}, \Lambda)$$

Posterior:

$$p(W_i | Y, X, L, \Lambda) = \prod_i \mathcal{N}(W_i | \Sigma_i X_i^T L^{-1} Y_i, \Sigma_i)$$

with $\Sigma_i = (\Lambda^{-1} + X_i^T L^{-1} X_i)^{-1}$.

Effect of batch-size: the full training set leads to a better estimate of the posteriors



Checkout the Full Paper!

S. Rossi, P. Michiardi and M. Filippone. "Good Initializations of Variational Bayes for Deep Models". **Proceedings of the 36th International Conference on Machine Learning (ICML 2019)**. 2019.



Some more insights!

Timing profiling (LENET-5): before training, 4 out of 5 optimal initializers are I-BLM

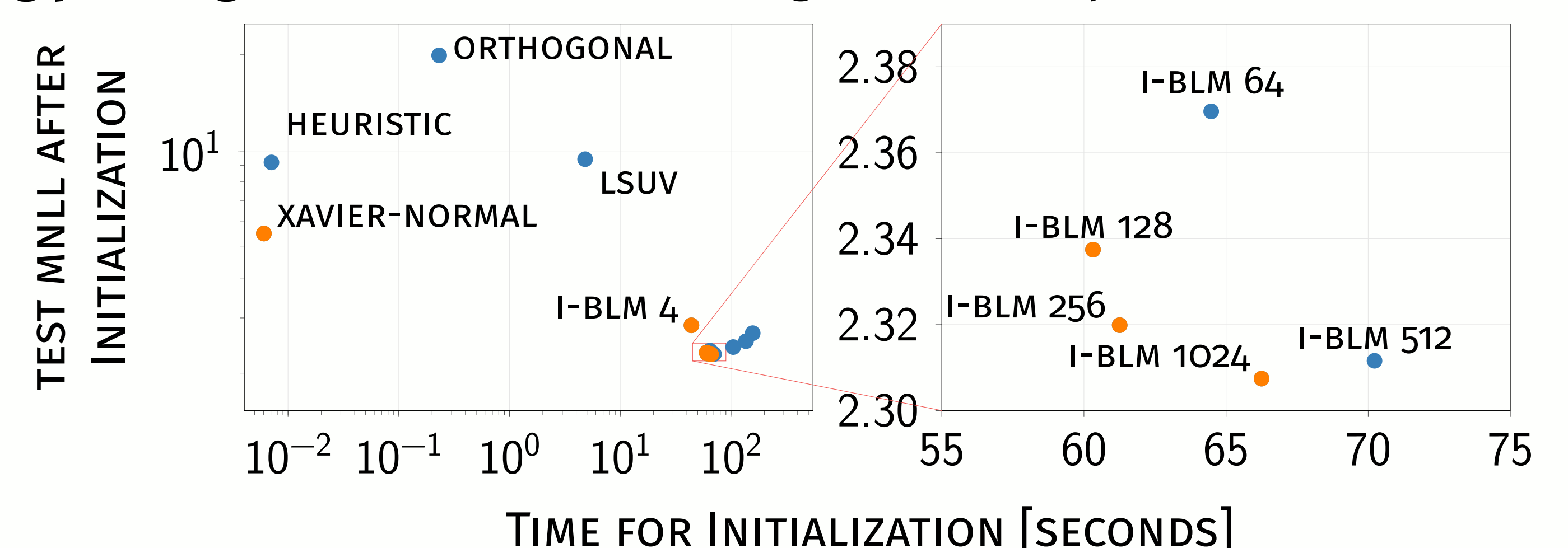


Figure: Comparison of initialization time versus test MNLL.

Regression and Classification on Bayesian DNNs

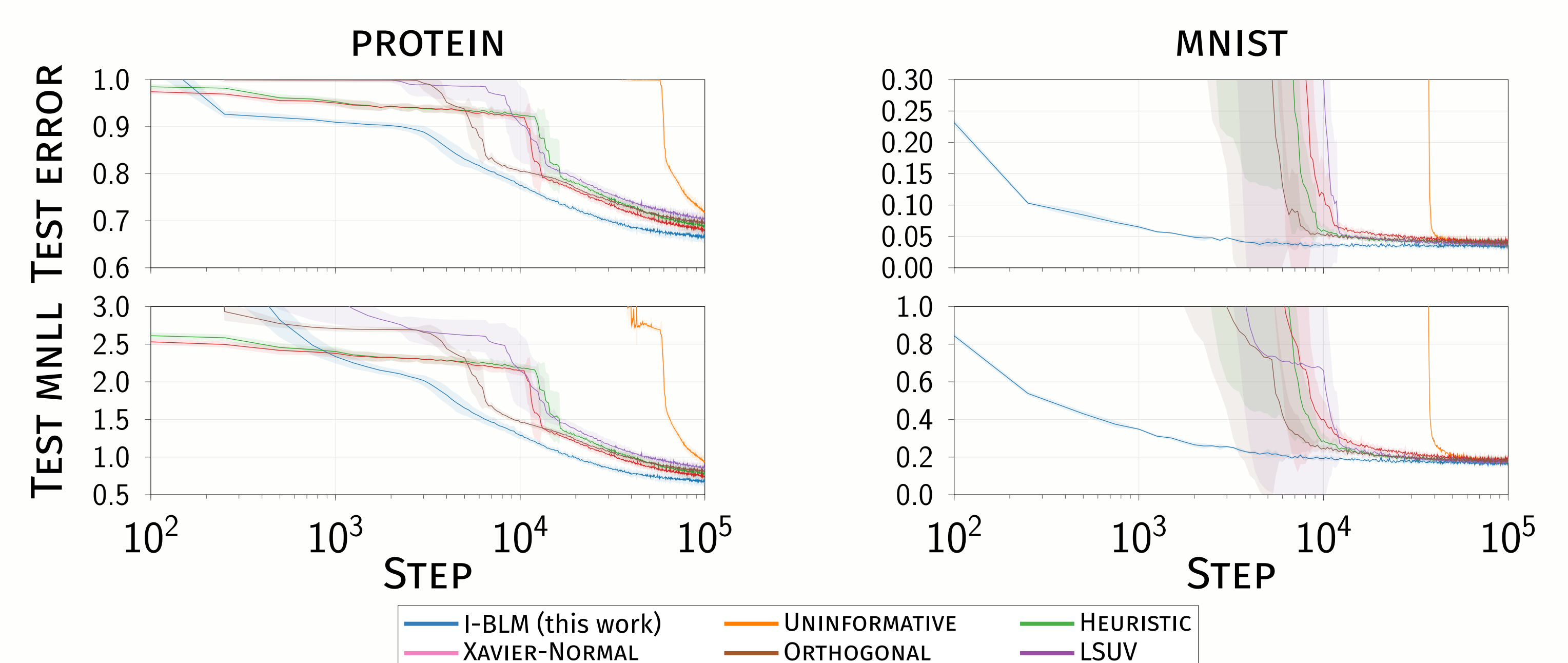


Figure: Progression of test error and test MNLL with different initializations on a 5x100 architecture.

I-BLM for Bayesian CNNs - VGG16

- ▶ Another initialization for Gaussian svi based on a MAP optimization (MAP INIT).
- ▶ Loss optimized for the same amount of time required by I-BLM. Solution used to initialize the means, while the log-variances are -5.5 .
- ▶ Models are trained for 100 minutes for the entire end-to-end training (curves are shifted by the initialization time).

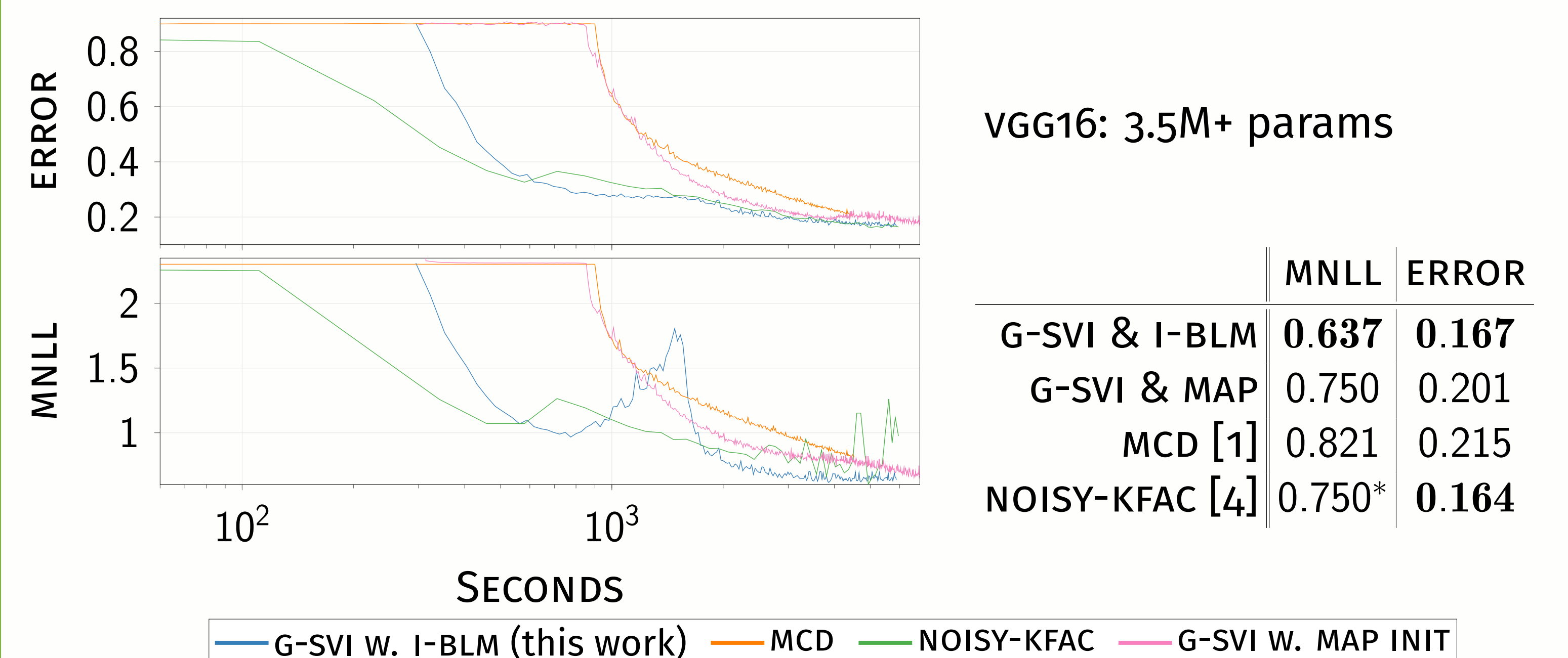


Figure & Table: Comparison between Gaussian factorized svi, MCD and NOISY-KFAC on VGG16 with CIFAR10

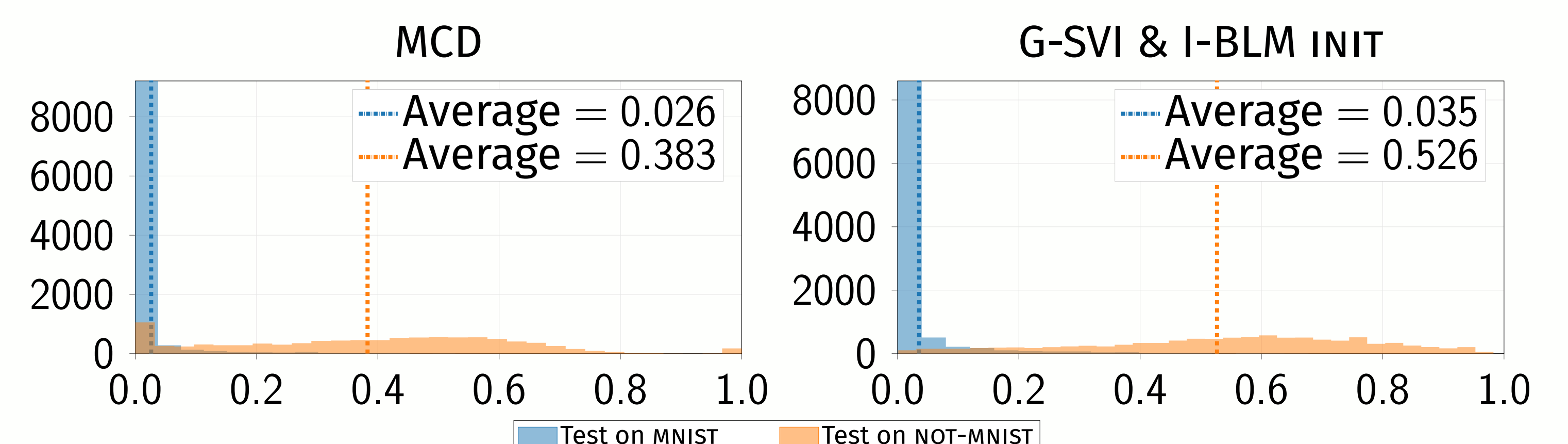


Figure: Entropy distribution while testing on MNIST and NOT-MNIST.

References

- [1] Y. Gal and Z. Ghahramani. "Bayesian Convolutional Neural Networks with Bernoulli Approximate Variational Inference". *Workshop track - ICLR*. June 2015.
- [2] A. Graves. "Practical Variational Inference for Neural Networks". *Advances in Neural Information Processing Systems 24*. 2011.
- [3] D. Milius et al. "Dirichlet-based Gaussian Processes for Large-scale Calibrated Classification". *Advances in Neural Information Processing Systems 31*. 2018.
- [4] G. Zhang et al. "Noisy Natural Gradient as Variational Inference". *Proceedings of the 35th International Conference on Machine Learning*. Oct. 2018.