

## Goals and Challenges

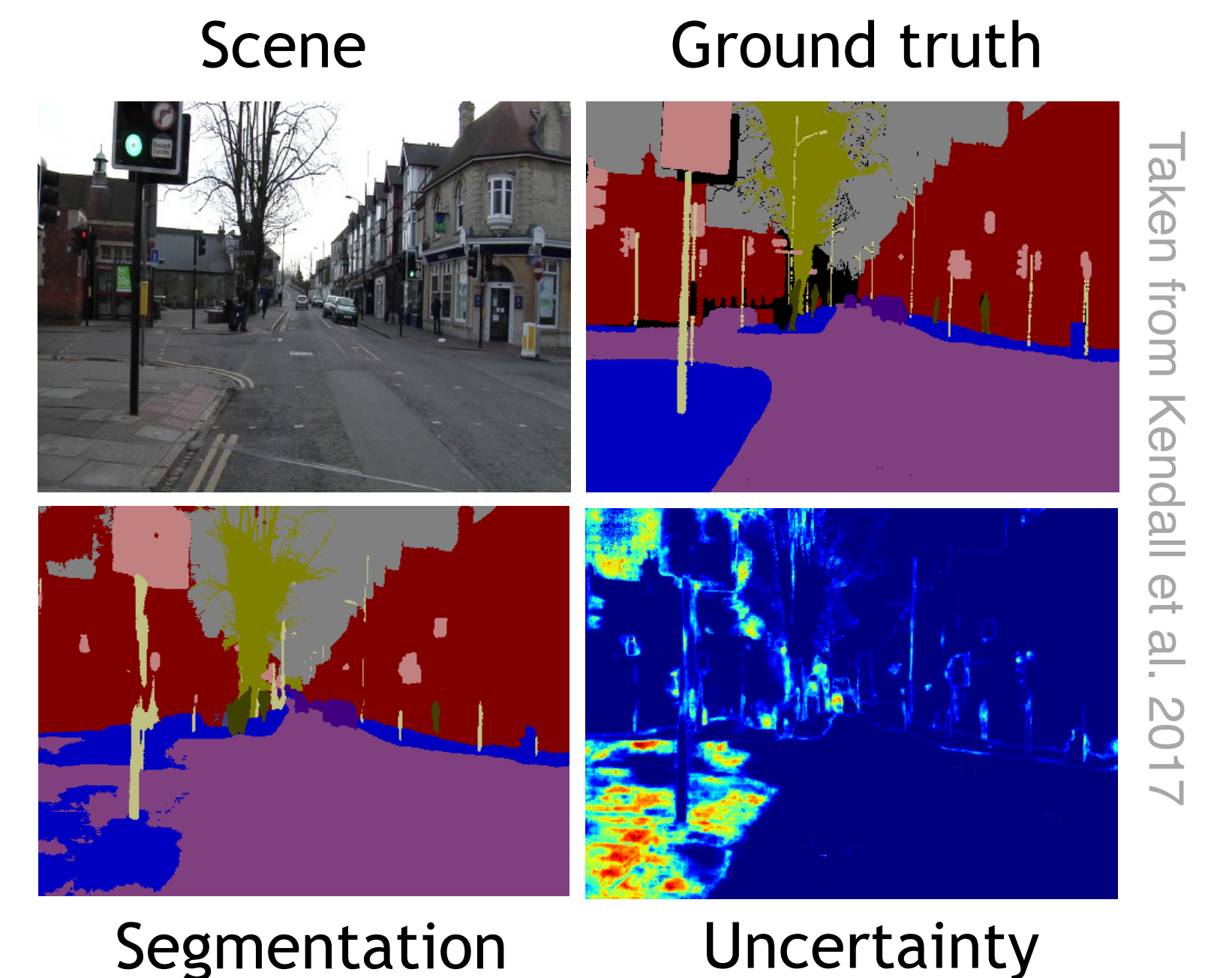
**Goal:** To design **AI that can continually learn using Bayesian principles.**

**Examples:** Uncertainty: Knowing how much we don't know, is useful to design

- Robots that can understand and reason about their environments.
- Methods that improve performance of deep-learning methods.

**Challenge:** Computation of the posterior distribution is difficult

**Main Idea:** Approximate integration by using optimization, and design simple algorithms that can be implemented within existing deep learning frameworks



## Fast and Simple Algorithms for Variational Inference

### Variational Inference

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta)p(\theta)d\theta}$$

Parameters (Data) Intractable integral

$$\approx q_\lambda(\theta) = \text{ExpFamily}(\lambda)$$

Variational Approximation Natural parameters

Maximize the Evidence Lower Bound (ELBO):

$$\max_{\lambda} \mathcal{L}(\lambda) := \mathbb{E}_{q_\lambda} [\log p(\mathcal{D}, \theta) - \log q_\lambda(\theta)]$$

Gradient descent (GD):  $\lambda \leftarrow \lambda + \rho \nabla_{\lambda} \mathcal{L}$

### VI with Natural-Gradient Descent

Sato 2001, Honkela et al. 2010, Hoffman et al. 2013

$$\text{NGD: } \lambda \leftarrow \lambda + \rho F(\lambda)^{-1} \nabla_{\lambda} \mathcal{L}$$

Natural Gradient

Fisher Information Matrix (FIM)

$$F(\lambda) := \mathbb{E}_{q_\lambda} [\nabla \log q_\lambda(\theta) \nabla \log q_\lambda(\theta)^{\top}]$$

- Fast convergence due to optimization in Riemannian manifold (not Euclidean space).
- But requires additional computations.
- Can we simplify/reduce this computation?

### Expectation Parameters

$$\mu := \mathbb{E}_{q_\lambda} [\phi(\theta)]$$

Expectation/moment/mean parameters Sufficient statistics

For Gaussians, it's mean and correlation matrix

$$\mathbb{E}_{q_\lambda} [\theta] = m \quad \mathbb{E}_{q_\lambda} [\theta \theta^{\top}] = mm^{\top} + V$$

A key relationship:  $F(\lambda)^{-1} \nabla_{\lambda} \mathcal{L} = \nabla_{\mu} \mathcal{L}$

Natural Gradient wrt natural parameter Gradient wrt expectation parameter

$$\text{NGD: } \lambda \leftarrow \lambda + \rho \nabla_{\mu} \mathcal{L}$$

### Example: Linear Regression

$$q_\lambda(\theta) := \mathcal{N}(m, V)$$

$$\mathbb{E}_q \left[ (y - X\theta)^{\top} (y - X\theta) + \gamma \theta^{\top} \theta - \log q_\lambda(\theta) \right]$$

likelihood prior approx

$$= -\mathbb{E}_{q_\lambda} [\theta]^{\top} X^{\top} y + \text{trace} [X^{\top} X \mathbb{E}_{q_\lambda} [\theta \theta^{\top}]]$$

$$\nabla_{\mathbb{E}_{q_\lambda} [\theta]} = \begin{bmatrix} -X^{\top} y & 0 & -V^{-1} m \\ X^{\top} X & \gamma I & -V^{-1} \end{bmatrix}$$

$$m \leftarrow (1 - \rho)m - \rho [X^{\top} X + \gamma I]^{-1} X^{\top} y$$

### Bayesian Neural Network

$$\mathbb{E}_q \left[ \sum_{i=1}^N \log p(y_i | f_{\theta}(x_i)) + \gamma \theta^{\top} \theta - \log q_\lambda(\theta) \right]$$

likelihood prior approx neural network

$$m \leftarrow m - \beta (S + \gamma I)^{-1} [g_i(\theta) + \gamma m]$$

$$S \leftarrow (1 - \beta)S + \beta H_i(\theta)$$

Back-propagated gradient & Hessian

$$\theta \sim q_\lambda(\theta), \quad g_i(\theta) := -\nabla_{\theta} \log p(y_i | f_{\theta}(x_i)),$$

$$V^{-1} \leftarrow S + \gamma I, \quad H_i(\theta) := -\nabla_{\theta}^2 \log p(y_i | f_{\theta}(x_i))$$

### MLE vs NGD-VI

RMSprop for MLE

$$\begin{aligned} \theta &\leftarrow \mu \\ g &\leftarrow \frac{1}{M} \sum_i \nabla_{\theta} \log p(\mathcal{D}_i | \theta) \\ s &\leftarrow (1 - \beta)s + \beta g^2 \\ \mu &\leftarrow \mu + \alpha \frac{g}{\sqrt{s + \delta}} \end{aligned}$$

NGD for mean-field VI

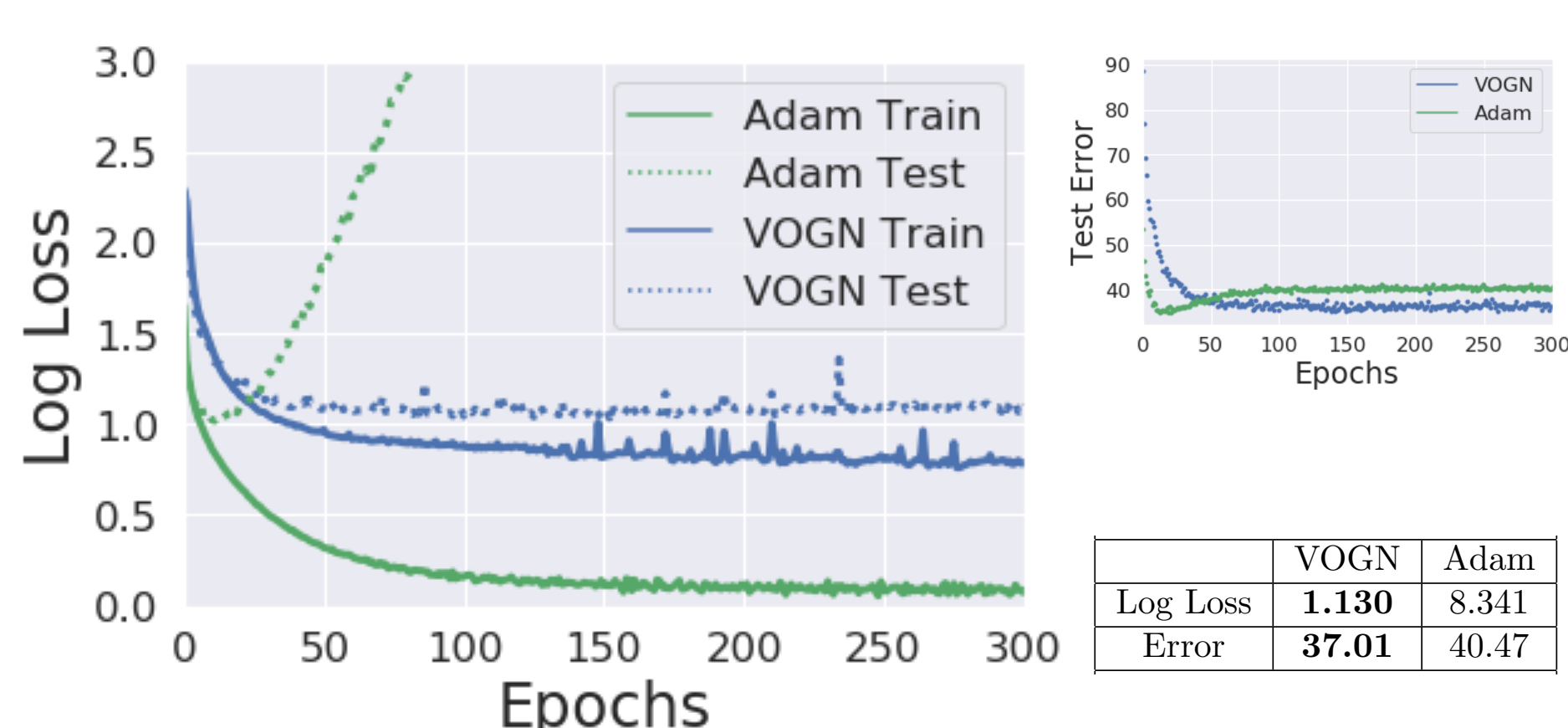
$$\begin{aligned} \theta &\leftarrow \mu + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, Ns + \lambda) \\ g &\leftarrow \frac{1}{M} \sum_i \nabla_{\theta} \log p(\mathcal{D}_i | \theta) \\ s &\leftarrow (1 - \beta)s + \beta \frac{1}{M} \sum_i [\nabla_{\theta} \log p(\mathcal{D}_i | \theta)]^2 \\ \mu &\leftarrow \mu + \alpha \frac{g + \lambda \mu / N}{s + \lambda / N} \end{aligned}$$

Variational Online Gauss-Newton (VOGN)

$$s \leftarrow (1 - \beta)s + \beta \frac{1}{M} \sum_i \nabla_{\theta\theta}^2 \log p(\mathcal{D}_i | \theta)$$

Variational RMSprop (Vprop)  $s \leftarrow (1 - \beta)s + \beta g^2$

### LeNet-5 on CIFAR10



(By Anirudh Jain)

### Stochastic, Low-Rank, Approximate, Natural-Gradient (SLANG)

NeurIPS 2018

- Low-rank + diagonal covariance matrix.
- **SLANG is linear in D!**

$$m \leftarrow m - \rho [UU^{\top} + D]^{-1} [g_i + \gamma m]$$

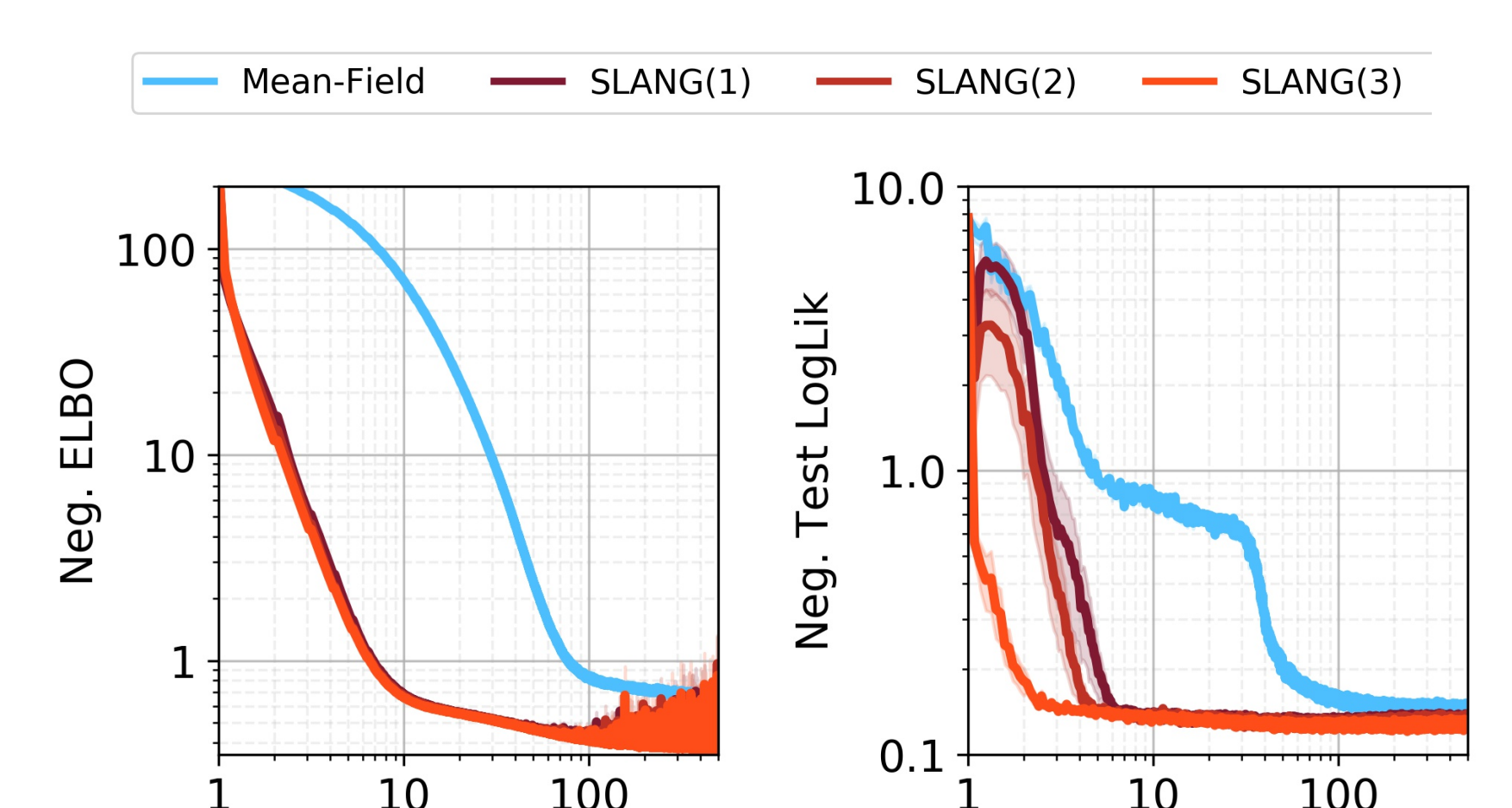
Low-Rank + diagonal

$$(1 - \beta)S + \beta H_i(\theta)$$

$$\begin{bmatrix} D \times L & L \times D \\ \text{gradient} & \text{gradient} \end{bmatrix} = \begin{bmatrix} D \times M & M \times D \\ \text{fast\_eig} & \text{fast\_eig} \end{bmatrix}$$

### SLANG is Faster than GD

Classification on USPS with BNs



### References

1. Fast and Scalable Bayesian Deep Learning by Weight-Perturbation in Adam, (ICML 2018), Khan, Nielsen, Tangkaratt, Lin, Gal, and Srivastava.
2. SLANG: Fast Structured Covariance Approximations for Bayesian Deep Learning with Natural Gradient, (NeurIPS 2018), Mishkin, Kunstner, Nielsen, Schmidt, Khan.
3. Fast and Simple Natural-Gradient Descent for Variational Inference in Complex Models (ISITA 2018), Khan and Nielsen.