

# Implementation of a probabilistic Softmax for BNNs (Group 11)

David Warden, Khalil Sabri, and Miquel Florensa – Polytechnique Montréal, Canada

Based on the work proposed by J-A. Goulet & L. H. Nguyen

## Tractable Approximate Gaussian Inference (TAGI)

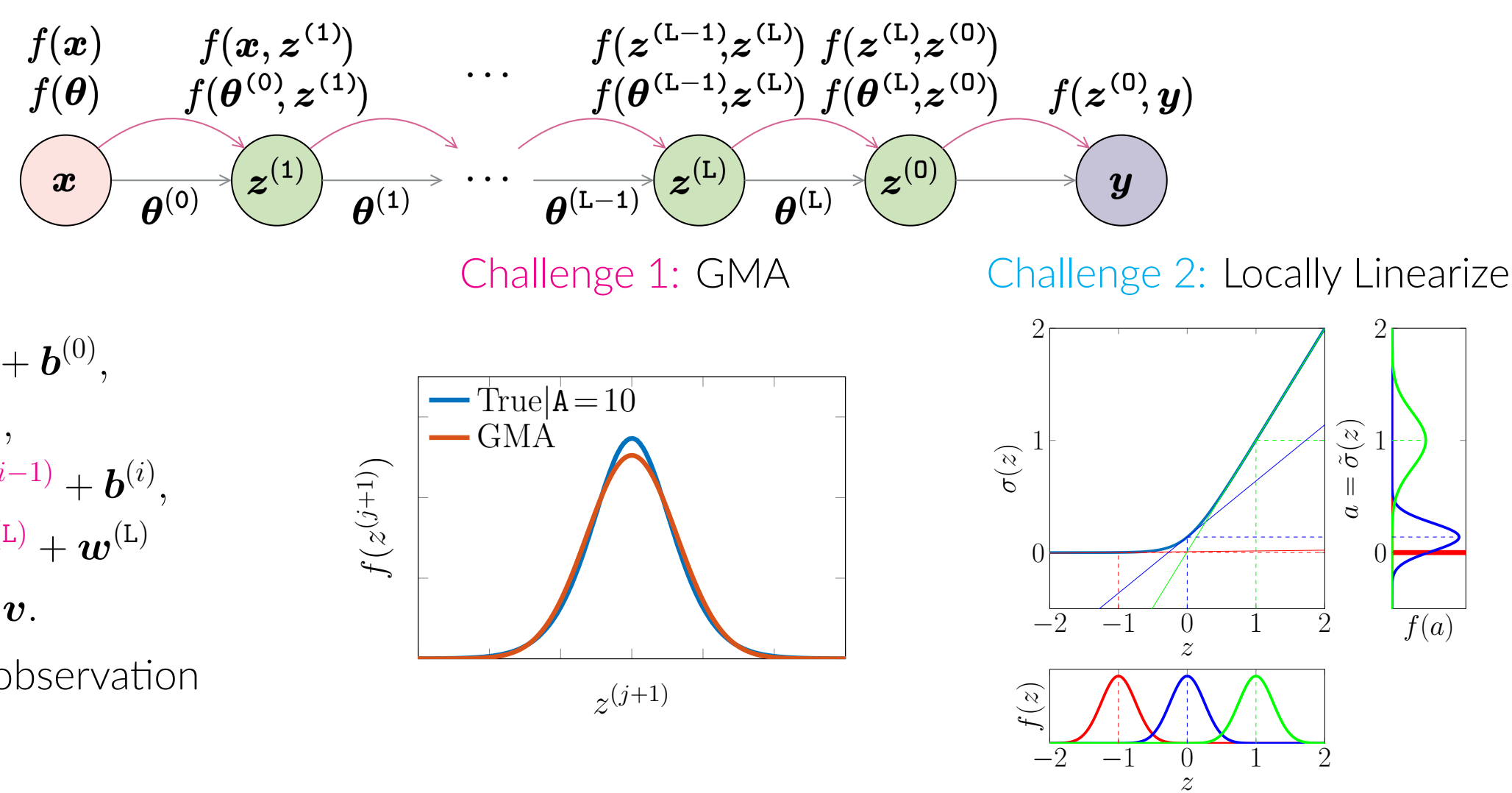
Parameters  $\begin{cases} f(\theta|\mathcal{D}) = \frac{f(\mathcal{D}|\theta)f(\theta)}{f(\mathcal{D})} \\ = \mathcal{N}(\theta; \mu_{\theta|\mathcal{D}}, \sigma_{\theta|\mathcal{D}}^2 \mathbf{I}) \end{cases}$

Hidden states  $\begin{cases} f(z|\mathcal{D}) = \mathcal{N}(z; \mu_{z|\mathcal{D}}, \sigma_{z|\mathcal{D}}^2 \mathbf{I}) \end{cases}$

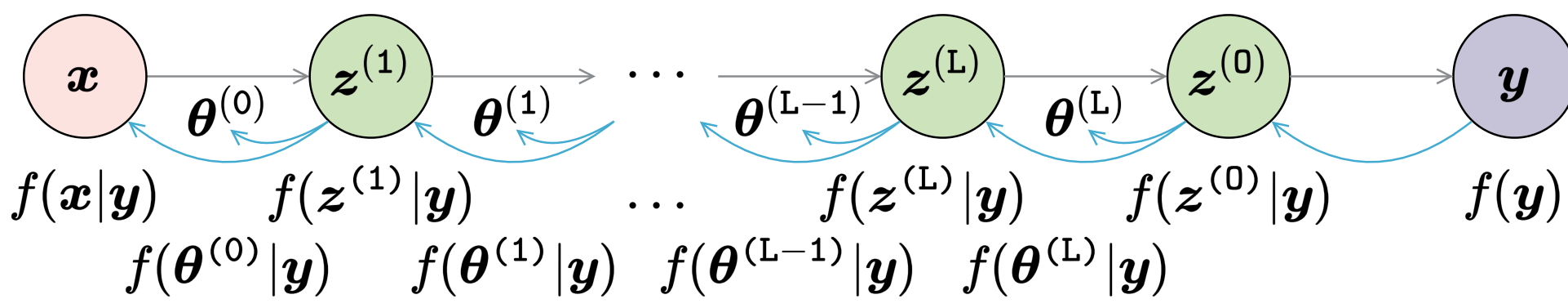
- Bayesian analytical inference  $\begin{cases}$  parameters
- Fast and scalable  $\rightarrow \mathcal{O}(n)$   $\begin{cases}$  hidden states
- End-to-end treatment of uncertainty
- No gradient backpropagation

Scalable analytical Bayesian inference in any neural network architecture

## TAGI - Forward uncertainty propagation



## TAGI - Backward parameter and hidden state inference



First, given the observations  $y$ , the output layer is updated as

$$f(z^{(0)} | y) = \mathcal{N}(z^{(0)}; \mu_{z^{(0)}|y}, \Sigma_{z^{(0)}|y})$$

$$\mu_{z^{(0)}|y} = \mu_{z^{(0)}} + \Sigma_{Yz^{(0)}} \Sigma_Y^{-1} (y - \mu_Y)$$

$$\Sigma_{z^{(0)}|y} = \Sigma_{z^{(0)}} - \Sigma_{Yz^{(0)}} \Sigma_Y^{-1} \Sigma_{Yz^{(0)}}^T$$

Next, the hidden units and parameters of the layer  $i^{th}$  are updated using a layer-wise procedure as

$$f(z | y) = \mathcal{N}(z; \mu_{z|y}, \Sigma_{z|y})$$

$$\mu_{z|y} = \mu_z + J_z (\mu_{z^+|y} - \mu_{z^+})$$

$$\Sigma_{z|y} = \Sigma_z + J_z (\Sigma_{z^+|y} - \Sigma_{z^+}) J_z^T$$

$$J_z = \Sigma_{zz^+} \Sigma_{z^+}^{-1}$$

$$f(\theta | y) = \mathcal{N}(\theta; \mu_{\theta|y}, \Sigma_{\theta|y})$$

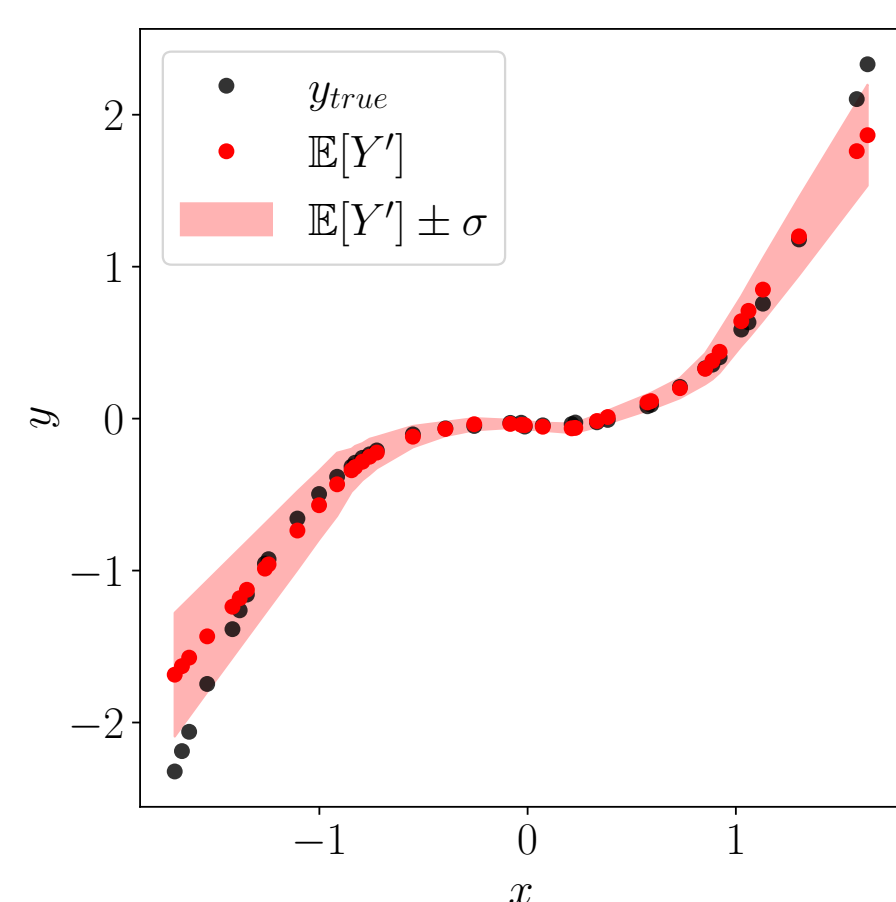
$$\mu_{\theta|y} = \mu_{\theta} + J_{\theta} (\mu_{z^+|y} - \mu_{z^+})$$

$$\Sigma_{\theta|y} = \Sigma_{\theta} + J_{\theta} (\Sigma_{z^+|y} - \Sigma_{z^+}) J_{\theta}^T$$

$$J_{\theta} = \Sigma_{\theta z^+} \Sigma_{z^+}^{-1}$$

where  $\{\theta^+, Z^+\}$  is the short-hand notation of  $\{\theta^{(j+1)}, Z^{(j+1)}\}$ , and  $\{\theta, Z\}$  is the short-hand notation of  $\{\theta^{(j)}, Z^{(j)}\}$ .

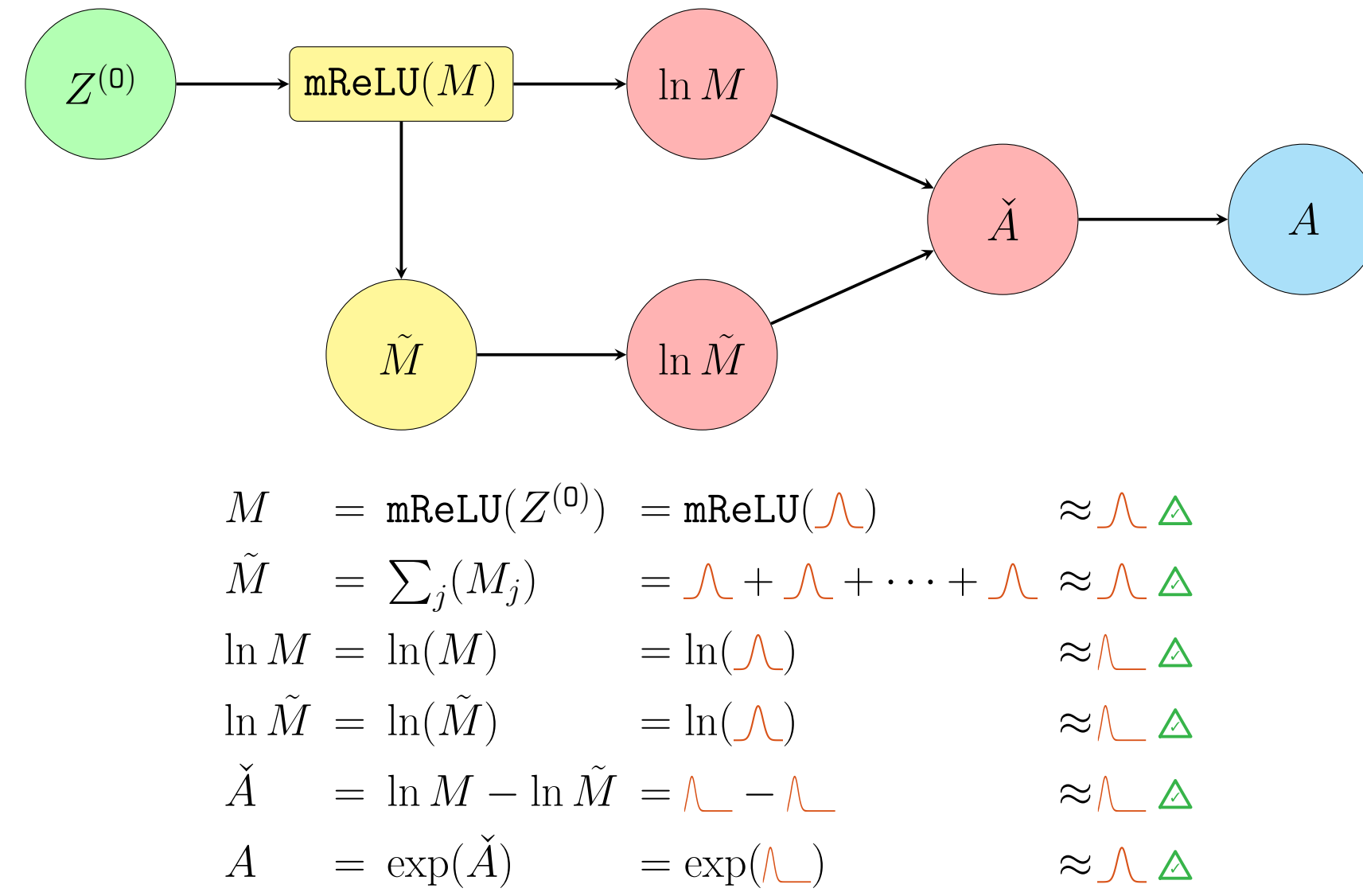
## Regression toy example



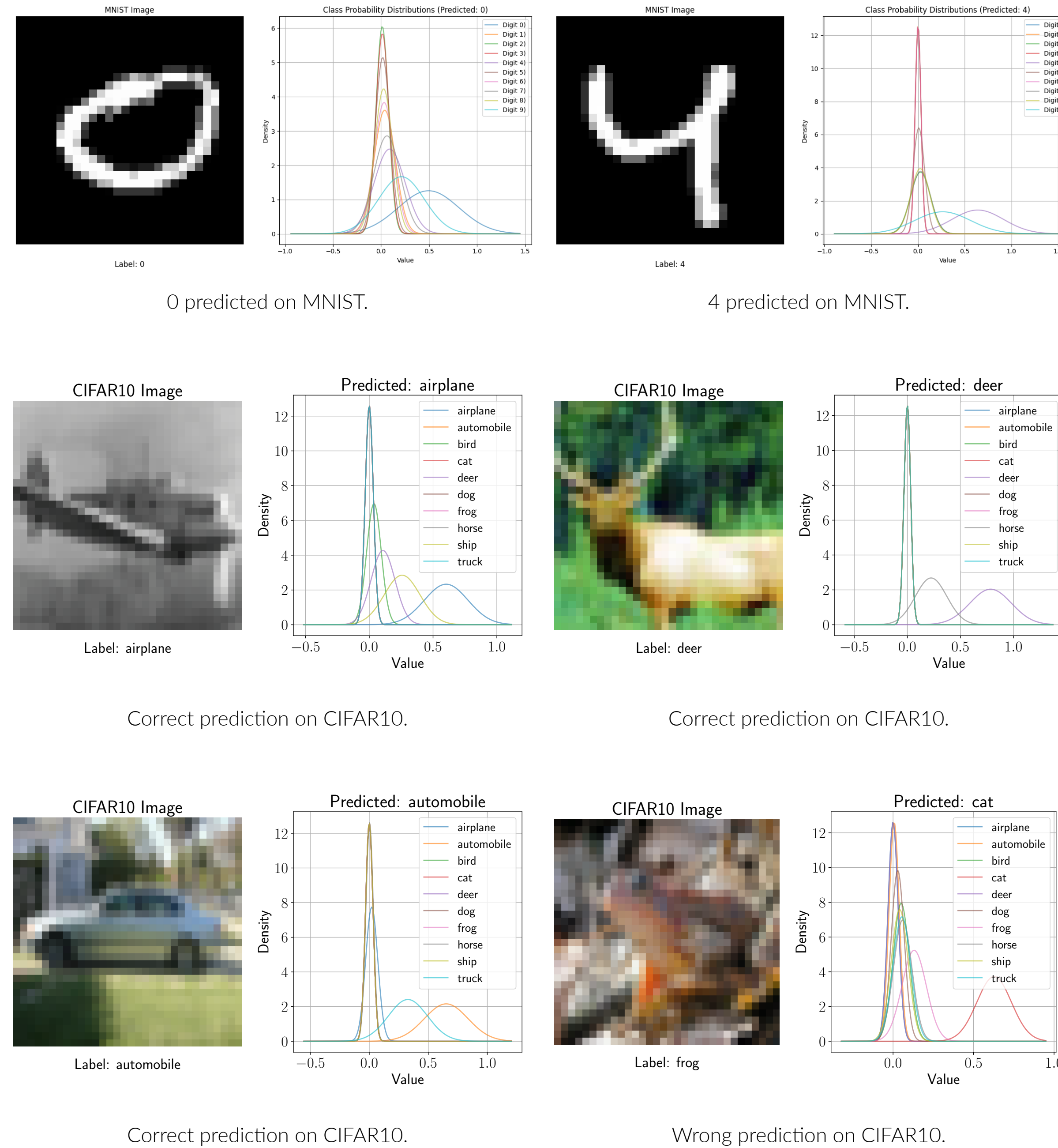
- We re-implement the toy example presented by the authors using the numpy library in Python.
- This includes defining a fully connected neural network with one layer and 40 hidden units.
- The figure shows the output of the model on the test set with the epistemic uncertainty.

## Remax - Probabilistic Softmax for uncertainty propagation

- TAGI cannot use directly the softmax function:  $\text{softmax}(z, i) = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$  since  $\wedge / \wedge \neq \wedge \Delta$
- J-A. Goulet & L. H. Nguyen propose probabilistic softmax by doing the division in log-transformed space.
- The random variable  $Z$  is Lognormal if  $\ln Z$  is Normal and  $\exp(Z)$  is Lognormal if  $Z$  is Normal.
- Probabilistic softmax does not work for  $Z$  with high variances;  $\text{Var}(Z) \approx 1$ .  $\Delta$
- Alternative softmax function:  $\text{remax}(z, i) = \frac{\text{mReLU}(z_i)}{\sum_j \text{mReLU}(z_j)}$  where  $\text{mReLU}$  is a Mixture Rectified Linear activation Unit.
- We implement **remax** function in py/cuTAGI library using C++ and CUDA.



## Epistemic uncertainty on classes using Remax



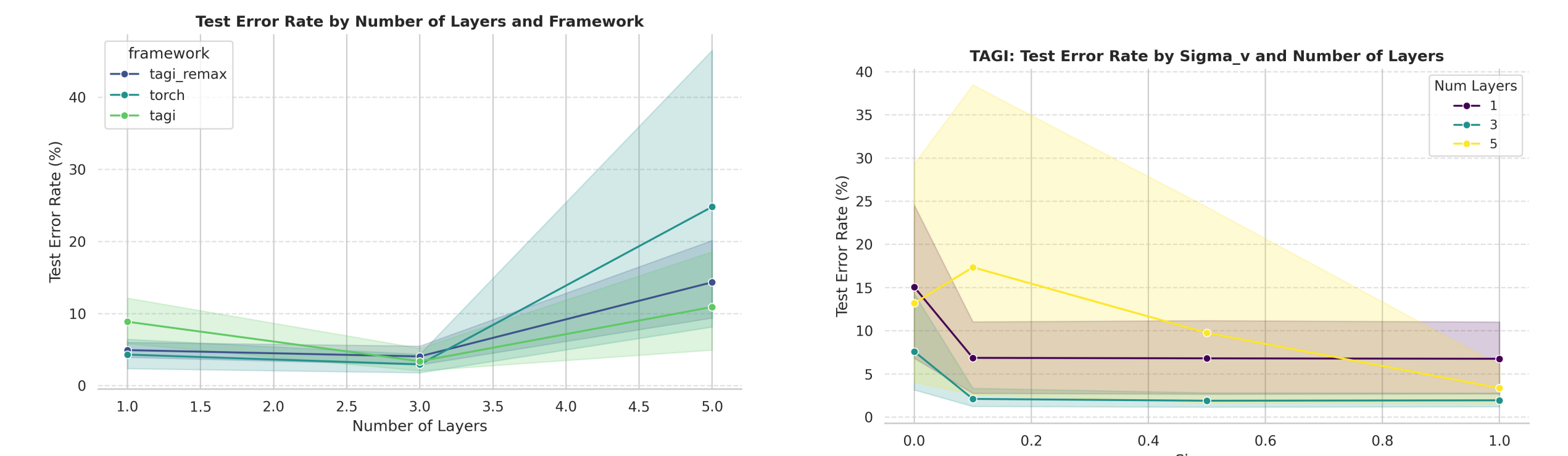
## Key Results: TAGI-HRC vs. TAGI-Remax vs. Torch

Framework	Avg. Training Error	Avg. Test Error
Torch	10.05	10.69
TAGI-HRC	7.07	<b>7.73</b>
TAGI-Remax	7.14	<b>7.78</b>

## Important Findings

- Batch Size:** TAGI excels with smaller batches, while Torch performs better with larger batches.
- Network Depth:** TAGI maintains stability as depth increases, unlike Torch, which struggles without batch normalization.
- Effect of  $\sigma_v$ :** Higher  $\sigma_v$  values enhance TAGI's performance, especially in deeper networks.
- Capacity Effect:** TAGI benefits more from increasing neurons/channels, showing improved generalization, while Torch fails to utilize additional capacity effectively.
- Comparison of Implementations:** TAGI-HRC achieves better error rates overall, while TAGI-Remax shows higher stability but less frequent best-case performance.
- Batch Normalization:** TAGI performs well even without batch normalization, handling deeper architectures (5+ layers) where Torch often diverges.

## Ablation Study: Performance Comparisons and Capacity Effects



**Test Error vs Network Depth:** TAGI maintains low error rates, while Torch degrades significantly.

Framework	Small Batch (16) With BN	Small Batch (16) W.o BN	Large Batch (512) With BN	Large Batch (512) W.o BN
Torch	<b>1.35</b>	31.20	<b>1.06</b>	16.22
TAGI-HRC	1.97	<b>1.74</b>	4.11	15.99
TAGI-Remax	12.05	7.30	5.75	<b>5.43</b>

**Impact of  $\sigma_v$ :** Higher  $\sigma_v$  values improve TAGI's test error.

Framework	Small (32) Neurons/layer	Large (512) Neurons/layer
Torch	7.03	7.32
TAGI-HRC	13.58	<b>8.98</b>
TAGI-Remax	9.25	<b>6.91</b>

## py/cuTAGI – Open-source Bayesian deep-learning framework



github.com/lhnguyen102/cuTAGI

**pip install pytagi**

- Performance-Oriented Kernels** written in C++/CUDA from scratch, with pybind11 for a seamless integration. It allows running on CPU & CUDA devices through a Python API.
- Broad Architecture Support** of the basic DNN layer including *Linear*, *CNNs*, *Transposed CNNs*, *LSTM*, *Average pooling*, *Batch* and *Layer normalization*, enabling the building of mainstream architectures such as *Autoencoders*, *Transformers*, *Diffusion Models*, and *GANs*.
- Model Building and Execution** currently supports sequential model building, with plans to introduce Eager Execution
- Open Platform** providing access to its entire codebase.

## TAGI-related references



- Analytically tractable hidden-states inference in Bayesian neural networks** (Nguyen and Goulet. Journal-to-conference track, ICLR 2024)
- Analytically tractable heteroscedastic uncertainty quantification in Bayesian neural networks for regression tasks** (Deka, Nguyen & Goulet. Neurocomputing, 2024)
- Tractable approximate Gaussian inference for Bayesian neural networks** (Goulet, Nguyen, & Amiri, JMLR, 2021)