

# Training Binary Neural Networks Using the Bayesian Learning Rule



Xiangming Meng  
(RIKEN AIP)

Presenter



Roman Bachmann  
(EPFL)



Mohammad Emtiyaz Khan  
(RIKEN AIP)



# Binary Neural Networks (BiNN)

- BiNN: Neural Networks with **binary weights**
- Much faster and much smaller [1,2]
- Difficult to optimize in theory (discrete optimization)
- But easy in practice: Just use SGD with “Straight-through estimator (STE)”!
- It is mysterious as to why this works [3]
- Are there any principled approaches to explain this?

1. Courbariaux et al., Training deep neural networks with binary weights during propagations. NeurIPS 2015.
2. Courbariaux et al., . Binarized neural networks.... arXiv:1602.02830, 2016.
3. Yin, P. et al., Understanding straight-through estimator in training activation quantized neural nets. arXiv, 2019.

# Our Contribution:

## Training BiNN using Bayes

- We show that by using the Bayesian Learning Rule [1,2] (natural-gradient variational inference), we can justify such previous approaches
  - Main point: **optimize the parameter of a Bernoulli distribution** (a continuous optimization problem)
- The Bayesian approach gives us an estimate of uncertainty which can be used for continual learning [3]

1. Khan, M. E. and Rue, H. Learning-algorithms from bayesian principles. ArXiv. 2019.
2. Khan, M. E. and Lin, W. Conjugate-computation variational inference. AISTATS, 2017
3. Kirkpatrick, J. et al. Overcoming catastrophic forgetting in neural networks. PNAS, 114(13):3521–3526, 2017.

# Training BiNN is a Discrete Optimization problem!

Binary weights

$$\min_{w \in \{-1, 1\}^W} \sum_{i \in \mathcal{D}} \ell(y_i, f_w(x_i))$$

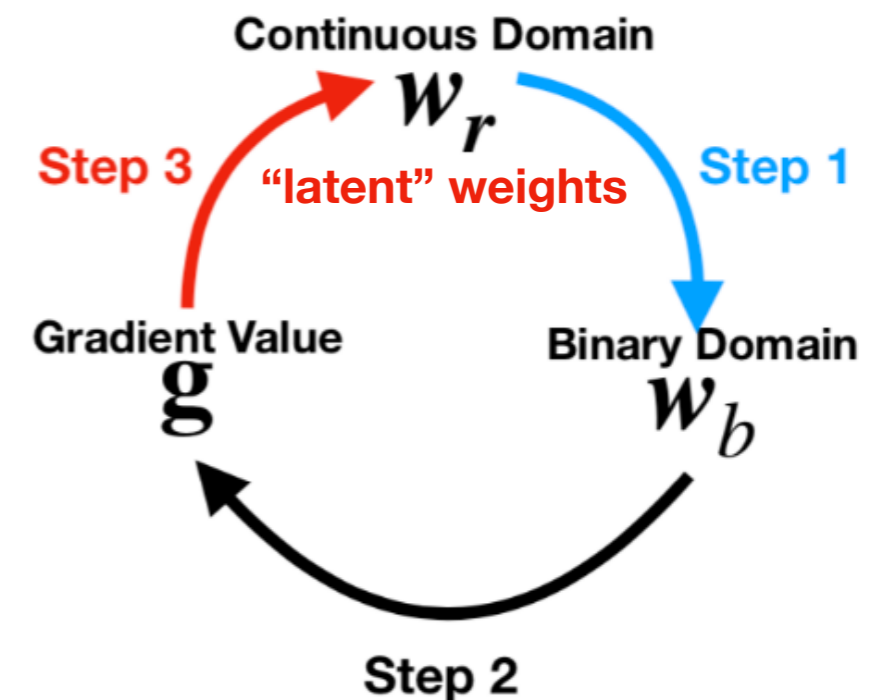
↓ Loss      ↑ Neural Network

Output      Input

# Training BiNN is a Discrete Optimization problem!

$$\begin{array}{c}
 \text{Binary weights} \\
 \downarrow \\
 \min_{w \in \{-1, 1\}^W} \sum_{i \in \mathcal{D}} \ell(y_i, f_w(x_i))
 \end{array}
 \begin{array}{c}
 \text{Output} \\
 \downarrow \\
 \ell(y_i, f_w(x_i)) \\
 \uparrow \text{Loss}
 \end{array}
 \begin{array}{c}
 \text{Input} \\
 \downarrow \\
 f_w(x_i) \\
 \uparrow \text{Neural Network}
 \end{array}$$

- Easy in practice: SGD with “Straight-through estimator (STE)” [1]



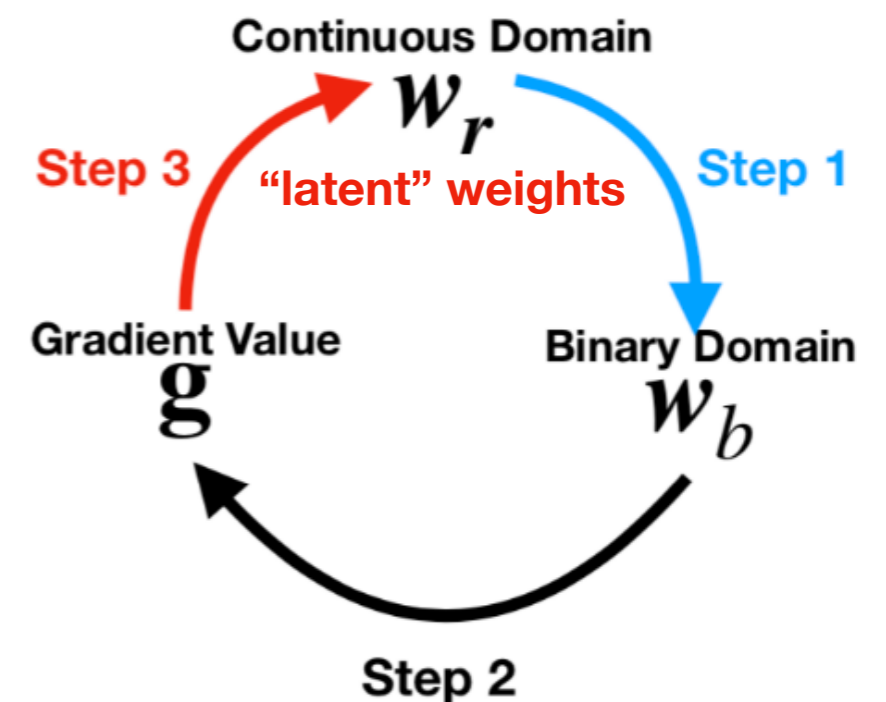
1. Bengio et al. Estimating or propagating gradients through stochastic neurons for conditional computation. arXiv:1308.3432, 2013.
2. Helwegen et al. Latent weights do not exist: Rethinking binarized neural network optimization. arXiv preprint arXiv:1906.02107, 2019.
3. Yin, P. et al., Understanding straight-through estimator in training activation quantized neural nets. arXiv, 2019.

# Training BiNN is a Discrete Optimization problem!

$$\begin{array}{c}
 \text{Binary weights} \\
 \downarrow \\
 \min_{\mathbf{w} \in \{-1, 1\}^W} \sum_{i \in \mathcal{D}} \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i))
 \end{array}$$

Output
Input  
↑ Loss
↑ Neural Network

- Easy in practice: SGD with “Straight-through estimator (STE)” [1]
- Helwegen et al. [2] argued “latent” weights are not weights but “Inertia”  
 → Binary Optimizer (Bop)



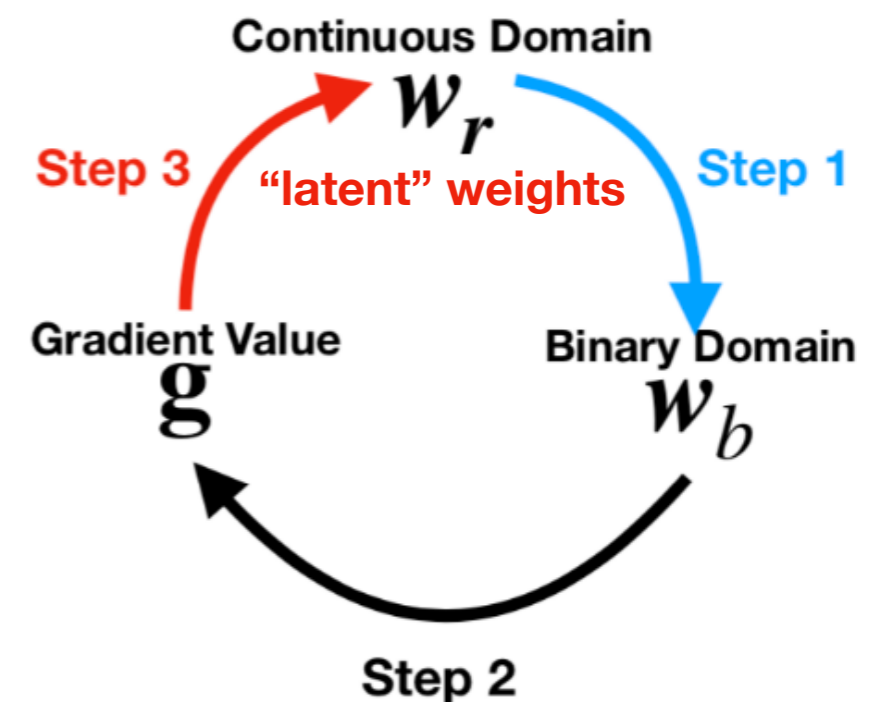
1. Bengio et al. Estimating or propagating gradients through stochastic neurons for conditional computation. arXiv:1308.3432, 2013.  
 2. Helwegen et al. Latent weights do not exist: Rethinking binarized neural network optimization. arXiv preprint arXiv:1906.02107, 2019.  
 3. Yin, P. et al., Understanding straight-through estimator in training activation quantized neural nets. arXiv, 2019.

# Training BiNN is a Discrete Optimization problem!

$$\begin{array}{c}
 \text{Binary weights} \\
 \downarrow \\
 \min_{\mathbf{w} \in \{-1,1\}^W} \sum_{i \in \mathcal{D}} \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i))
 \end{array}$$

Output
Input  
↑ Loss
↑ Neural Network

- Easy in practice: SGD with “Straight-through estimator (STE)” [1]
- Helwegen et al. [2] argued “latent” weights are not weights but “Inertia”  
 → Binary Optimizer (Bop)
- Open question: **Why does this work?** [3]



1. Bengio et al. Estimating or propagating gradients through stochastic neurons for conditional computation. arXiv:1308.3432, 2013.  
 2. Helwegen et al. Latent weights do not exist: Rethinking binarized neural network optimization. arXiv preprint arXiv:1906.02107, 2019.  
 3. Yin, P. et al., Understanding straight-through estimator in training activation quantized neural nets. arXiv, 2019.

# BayesBiNN

- Main point: **optimize the parameters of Bernoulli distribution** (a continuous optimization problem)
- Problem reformulation: Optimize distribution over weights<sup>[1,2]</sup>

$$\min_{q(\mathbf{w})} \mathbb{E}_{q(\mathbf{w})} \left[ \sum_{i=1}^N \overset{\text{Loss}}{\ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i))} \right] + \mathbb{D}_{KL}[q(\mathbf{w}) \parallel p(\mathbf{w})]$$

Posterior approximation over weights      KL Divergence      Prior Distribution

1. Zellner, A. Optimal information processing and Bayes's theorem. *The American Statistician*, 42(4):278–280, 1988.

2. Bissiri et al.. A general framework for updating belief distributions. *Journal of the Royal Statistical Society*, 78(5):1103–1130, 2016.



# BayesBiNN

- Main point: **optimize the parameters of Bernoulli distribution** (a continuous optimization problem)
- Problem reformulation: Optimize distribution over weights<sup>[1,2]</sup>

$$\min_{q(\mathbf{w})} \mathbb{E}_{q(\mathbf{w})} \left[ \sum_{i=1}^N \overset{\text{Loss}}{\ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i))} \right] + \mathbb{D}_{KL}[q(\mathbf{w}) \parallel p(\mathbf{w})]$$

↓
↓
↓

Posterior approximation
KL
Prior

over weights
Divergence
Distribution

- $q(\mathbf{w})$  is chosen to be mean-field Bernoulli distribution

$$q(\mathbf{w}) = \prod_{i=1}^D p_i^{\frac{1+w_i}{2}} (1-p_i)^{\frac{1-w_i}{2}}$$

↓  $w_i \in \{-1, +1\}$

Probability of  $w_i = +1$

↔

$$q(\mathbf{w}) = \prod_{i=1}^D \exp \left[ \lambda_i \phi(w_i) - A(\lambda_i) \right]$$

↓

Natural parameters:  $\lambda_i := \frac{1}{2} \log \frac{p_i}{1-p_i}$

1. Zellner, A. Optimal information processing and Bayes's theorem. The American Statistician, 42(4):278–280, 1988.

2. Bissiri et al.. A general framework for updating belief distributions. Journal of the Royal Statistical Society, 78(5):1103–1130, 2016.

# BayesBiNN

- The Bayesian learning rule<sup>[1]</sup> (natural-gradient variational inference)

$$\lambda \leftarrow (1 - \alpha) \lambda - \alpha \left\{ \nabla_{\mu} \mathbb{E}_{q(\mathbf{w})} \left[ \sum_{i=1}^N \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i)) \right] - \lambda_0 \right\}$$

Learning rate
↑
↓
↓
↓

Natural parameter of  $q(\mathbf{w})$ 
Expectation parameter of  $q(\mathbf{w})$ 
Natural parameter of  $p(\mathbf{w})$

1. Khan, M. E. and Rue, H. Learning-algorithms from bayesian principles. 2019.

2. Maddison, et al., The concrete distribution: A continuous relaxation of discrete random variables. arXiv:1611.00712, 2016.

3. Jiang et al. Categorical reparameterization with gumbel-softmax. arXiv:1611.01144, 2016.

# BayesBiNN

- The Bayesian learning rule<sup>[1]</sup> (natural-gradient variational inference)

$$\lambda \leftarrow (1 - \alpha) \lambda - \alpha \left\{ \nabla_{\mu} \mathbb{E}_{q(\mathbf{w})} \left[ \sum_{i=1}^N \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i)) \right] - \lambda_0 \right\}$$

Learning rate

Natural parameter of  $q(\mathbf{w})$

Expectation parameter of  $q(\mathbf{w})$

Natural parameter of  $p(\mathbf{w})$

How to compute?

1. Khan, M. E. and Rue, H. Learning-algorithms from bayesian principles. 2019.

2. Maddison, et al., The concrete distribution: A continuous relaxation of discrete random variables. arXiv:1611.00712, 2016.

3. Jiang et al. Categorical reparameterization with gumbel-softmax. arXiv:1611.01144, 2016.

# BayesBiNN

- The Bayesian learning rule<sup>[1]</sup> (natural-gradient variational inference)

$$\lambda \leftarrow (1 - \alpha) \lambda - \alpha \left\{ \nabla_{\mu} \mathbb{E}_{q(\mathbf{w})} \left[ \sum_{i=1}^N \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i)) \right] - \lambda_0 \right\}$$

Learning rate
How to compute?

Natural parameter of  $q(\mathbf{w})$ 
Expectation parameter of  $q(\mathbf{w})$ 
Natural parameter of  $p(\mathbf{w})$

- Using the Gumbel Softmax trick<sup>[2,3]</sup>, we can approximate the natural gradient by using the mini-batch gradient

$$\nabla_{\mu} \mathbb{E}_{q(\mathbf{w})} \left[ \sum_{i=1}^N \ell(y_i, f_{\mathbf{w}}(\mathbf{x}_i)) \right] \approx \mathbf{s} \odot \mathbf{g}$$

Scale vector
Minibatch Gradient, easy to compute!

$$\mathbf{g} := \frac{1}{m} \sum_{i \in \mathcal{M}} \nabla_{\mathbf{w}_b} \ell(y_i, f_{\mathbf{w}_b}(\mathbf{x}_i)) \quad \mathbf{s} := \frac{N(1 - \mathbf{w}_b^2)}{\tau(1 - \tanh(\boldsymbol{\lambda}))^2}$$

1. Khan, M. E. and Rue, H. Learning-algorithms from bayesian principles. 2019.

2. Maddison, et al., The concrete distribution: A continuous relaxation of discrete random variables. arXiv:1611.00712, 2016.

3. Jiang et al. Categorical reparameterization with gumbel-softmax. arXiv:1611.01144, 2016.

# BayesBiNN Justifies Some Previous Methods

STE	Our BayesBiNN method	Bop
$w_b \leftarrow \text{sign}(w_r)$ $\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$ $w_r \leftarrow w_r - \alpha \mathbf{g}$	$w_b \leftarrow \tanh((w_r + \delta)/\tau)$ $\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$ $w_r \leftarrow (1 - \alpha)w_r - \alpha \mathbf{s} \odot \mathbf{g}$	$w_b \leftarrow \text{hyst}(w_r, w_b, \gamma)$ $\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$ $w_r \leftarrow (1 - \alpha)w_r - \alpha \mathbf{g}$

Note that  $w_r$  in BayesBiNN corresponds to  $\lambda$

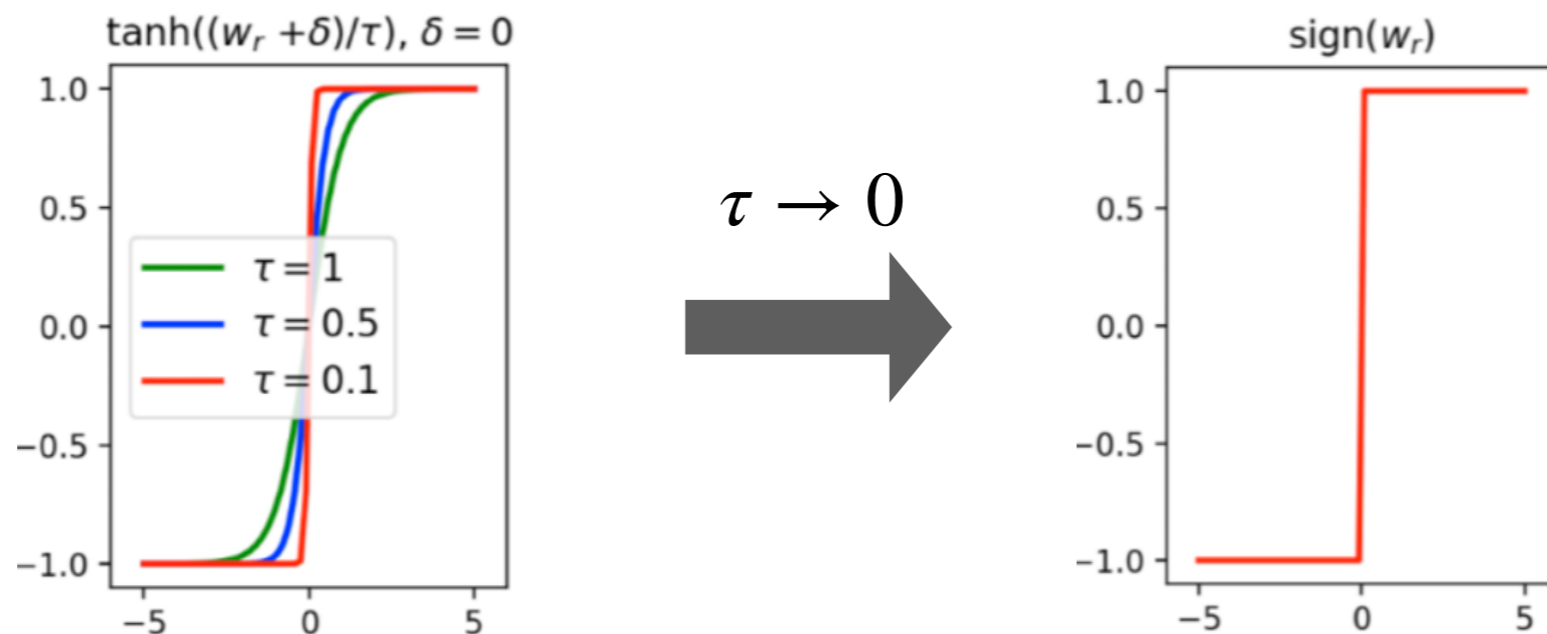
- Main point 1: **STE works as a special case of BayesBiNN as  $\tau \rightarrow 0$**

# BayesBiNN Justifies Some Previous Methods

STE	Our BayesBiNN method	Bop
$w_b \leftarrow \text{sign}(w_r)$ $\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$ $w_r \leftarrow w_r - \alpha \mathbf{g}$	$w_b \leftarrow \tanh((w_r + \delta)/\tau)$ $\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$ $w_r \leftarrow (1 - \alpha)w_r - \alpha \mathbf{s} \odot \mathbf{g}$	$w_b \leftarrow \text{hyst}(w_r, w_b, \gamma)$ $\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$ $w_r \leftarrow (1 - \alpha)w_r - \alpha \mathbf{g}$

Note that  $w_r$  in BayesBiNN corresponds to  $\lambda$

- Main point 1: STE works as a special case of BayesBiNN as  $\tau \rightarrow 0$

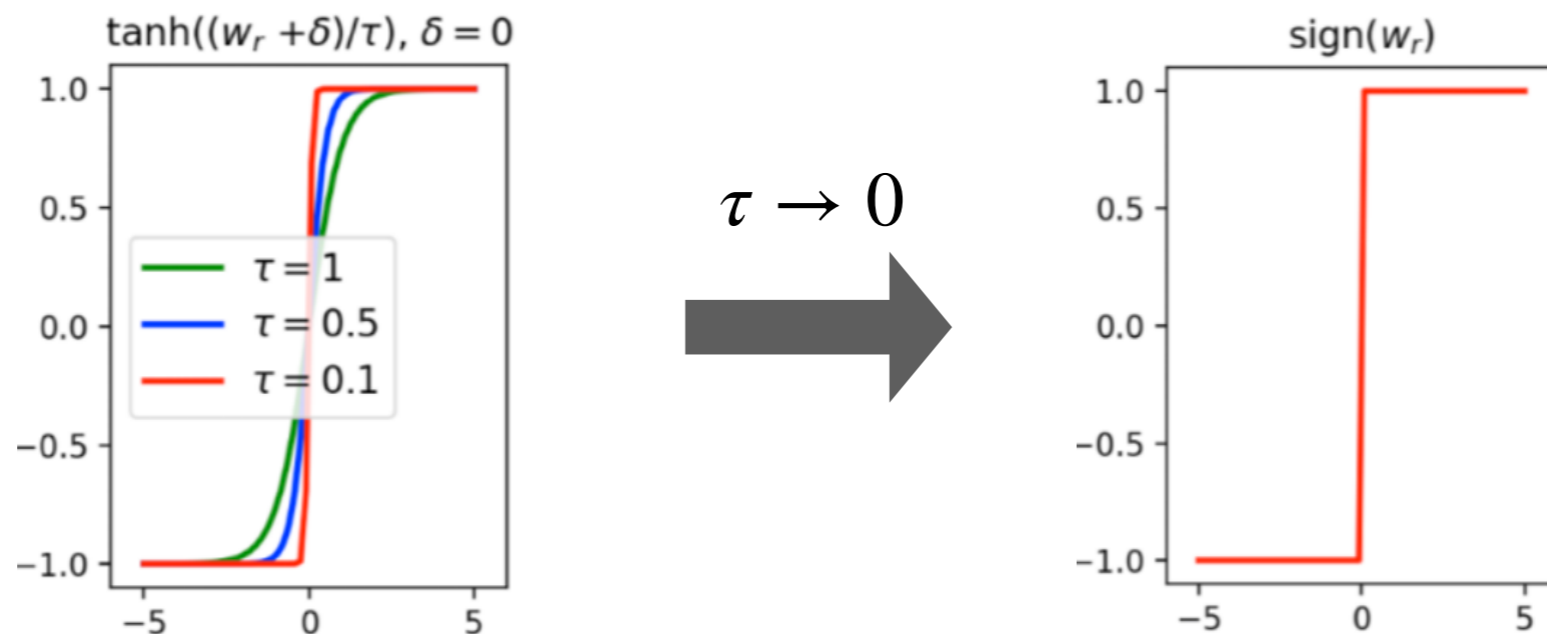


# BayesBiNN Justifies Some Previous Methods

STE	Our BayesBiNN method	Bop
$w_b \leftarrow \text{sign}(w_r)$ $\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$ $w_r \leftarrow w_r - \alpha \mathbf{g}$	$w_b \leftarrow \tanh((w_r + \delta)/\tau)$ $\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$ $w_r \leftarrow (1 - \alpha)w_r - \alpha \mathbf{s} \odot \mathbf{g}$	$w_b \leftarrow \text{hyst}(w_r, w_b, \gamma)$ $\mathbf{g} \leftarrow \nabla_{w_b} \ell(y, f_{w_b}(\mathbf{x}))$ $w_r \leftarrow (1 - \alpha)w_r - \alpha \mathbf{g}$

Note that  $w_r$  in BayesBiNN corresponds to  $\lambda$

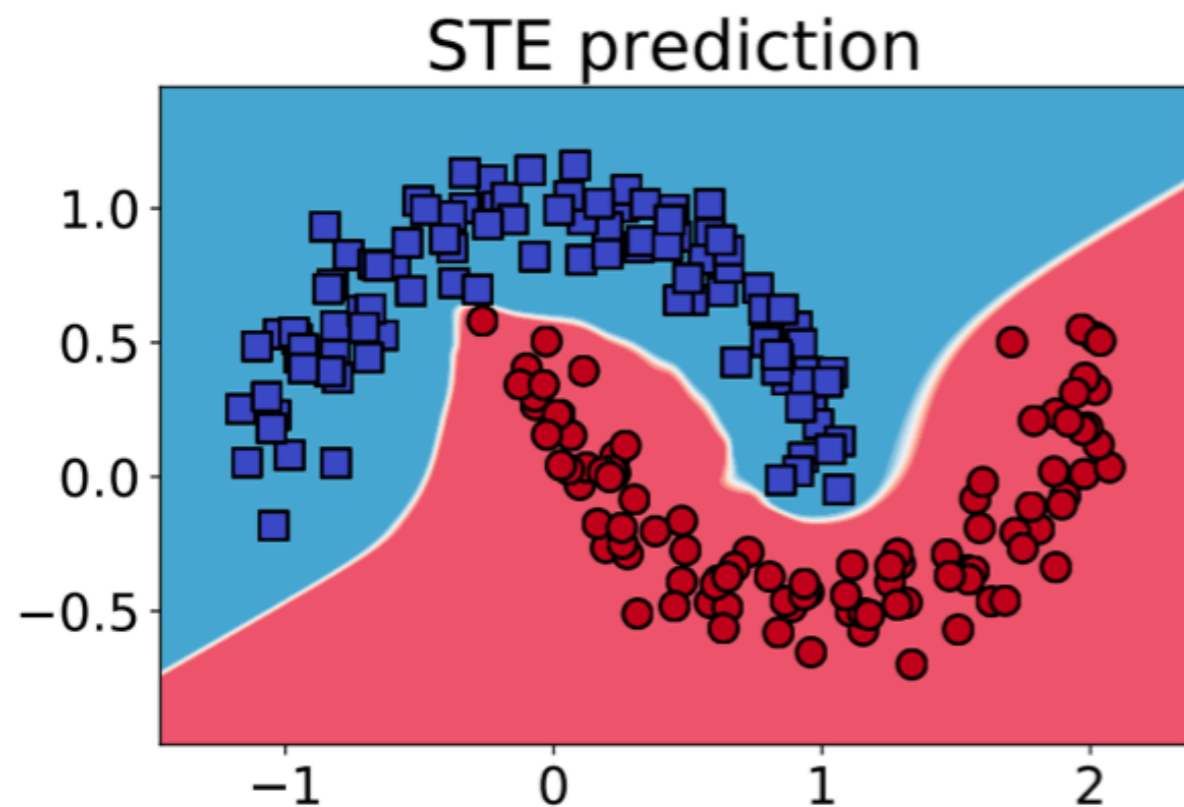
- Main point 1: STE works as a special case of BayesBiNN as  $\tau \rightarrow 0$



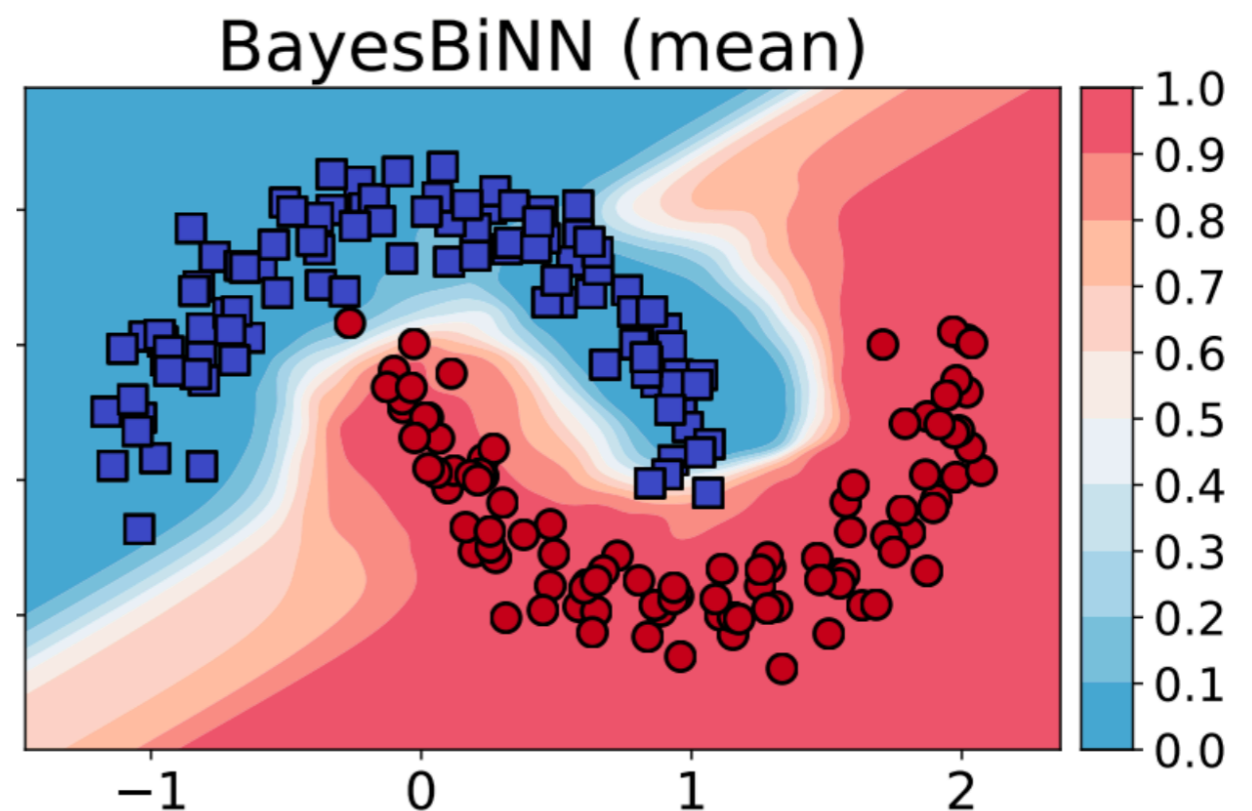
- Main point 2: Justify the “exponential average” used in Bop

# Uncertainty Estimation

- Main point: **BayesBiNN** obtains uncertainty estimates around the **classification boundaries**



Classification on two moons dataset



$$\hat{p}_k \leftarrow \frac{1}{C} \sum_{c=1}^C p(y = k | \mathbf{x}, \mathbf{w}^{(c)}), C = 10 \quad \mathbf{w}^{(c)} \sim q(\mathbf{w})$$

- STE finds a deterministic boundary
- **Open-source Code Available** : <https://github.com/team-approx-bayes/BayesBiNN>



# BayesBiNN $\approx$ STE

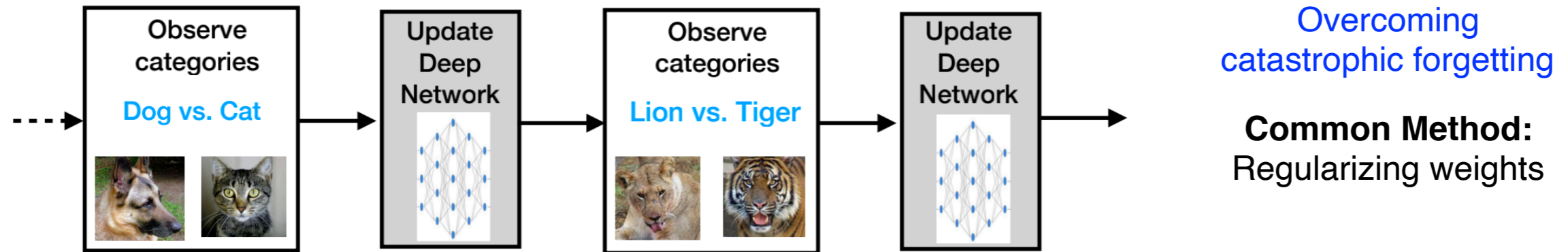
- **Open-source Code Available** : <https://github.com/team-approx-bayes/BayesBiNN>

Table 2. Results of different optimizers trained on MNIST, CIFAR-10 and CIFAR-100 (Averaged over 5 runs).

Dataset	Optimizer	Train Accuracy	Validation Accuracy	Test Accuracy
MNIST	STE Adam	99.78 $\pm$ 0.10 %	99.02 $\pm$ 0.11 %	<b>98.85 <math>\pm</math> 0.09 %</b>
	Bop	99.23 $\pm$ 0.04 %	98.55 $\pm$ 0.05 %	98.47 $\pm$ 0.02 %
	PMF		99.06 $\pm$ 0.01 %	98.80 $\pm$ 0.06 %
	<b>BayesBiNN (mode)</b>	99.85 $\pm$ 0.05 %	99.02 $\pm$ 0.13 %	<b>98.86 <math>\pm</math> 0.05 %</b>
	<b>BayesBiNN (mean)</b>	99.85 $\pm$ 0.05 %	99.02 $\pm$ 0.13 %	<b>98.86 <math>\pm</math> 0.05 %</b>
	Full-precision	99.96 $\pm$ 0.02 %	99.15 $\pm$ 0.14 %	99.01 $\pm$ 0.06 %
CIFAR-10	STE Adam	99.99 $\pm$ 0.01 %	94.25 $\pm$ 0.42 %	<b>93.55 <math>\pm</math> 0.15 %</b>
	Bop	99.79 $\pm$ 0.03 %	93.49 $\pm$ 0.17 %	93.00 $\pm$ 0.11 %
	PMF		91.87 $\pm$ 0.10 %	91.43 $\pm$ 0.14 %
	<b>BayesBiNN (mode)</b>	99.96 $\pm$ 0.01 %	94.23 $\pm$ 0.41 %	<b>93.72 <math>\pm</math> 0.16 %</b>
	<b>BayesBiNN (mean)</b>	99.96 $\pm$ 0.01 %	94.23 $\pm$ 0.41 %	<b>93.72 <math>\pm</math> 0.15 %</b>
	Full-precision	100.00 $\pm$ 0.00 %	94.54 $\pm$ 0.29 %	93.90 $\pm$ 0.17 %
CIFAR-100	STE Adam	99.06 $\pm$ 0.15 %	74.09 $\pm$ 0.15 %	72.89 $\pm$ 0.21 %
	Bop	90.09 $\pm$ 0.57 %	69.97 $\pm$ 0.29 %	69.58 $\pm$ 0.15 %
	PMF		69.86 $\pm$ 0.08 %	70.45 $\pm$ 0.25 %
	<b>BayesBiNN (mode)</b>	98.02 $\pm$ 0.18 %	74.76 $\pm$ 0.41 %	<b>73.68 <math>\pm</math> 0.31 %</b>
	<b>BayesBiNN (mean)</b>	98.02 $\pm$ 0.18 %	74.76 $\pm$ 0.41 %	<b>73.65 <math>\pm</math> 0.41 %</b>
	Full-precision	99.89 $\pm$ 0.02 %	75.89 $\pm$ 0.41 %	74.83 $\pm$ 0.26 %

# Uncertainty Provided by BayesBiNN Enables Continual Learning

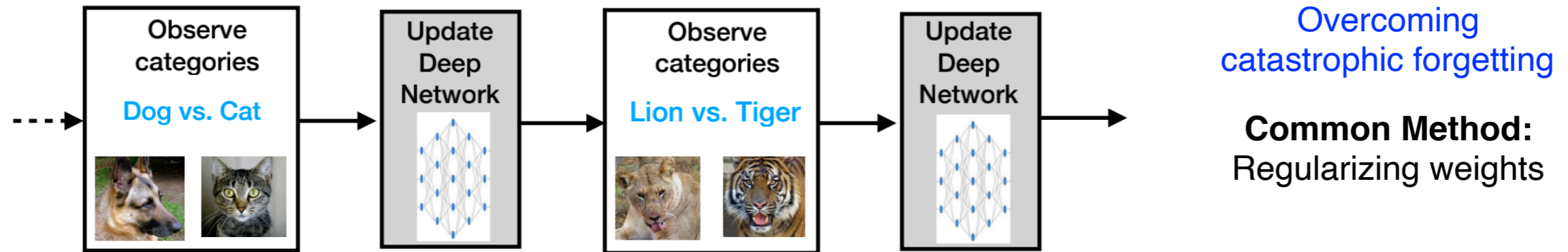
- Main point: **BayesBiNN enables continual learning (CL) for BiNN using the intrinsic KL divergence as regularization**
- CL: Sequentially learning new tasks without forgetting old ones<sup>[1]</sup>



- But, it is unclear how to regularize binary weights of BiNN using STE/Bop

# Uncertainty Provided by BayesBiNN Enables Continual Learning

- Main point: **BayesBiNN enables continual learning (CL) for BiNN using the intrinsic KL divergence as regularization**
- CL: Sequentially learning new tasks without forgetting old ones<sup>[1]</sup>



- But, it is unclear how to regularize binary weights of BiNN using STE/Bop
- In BayesBiNN, there is one natural solution using KL divergence

Independent Learning

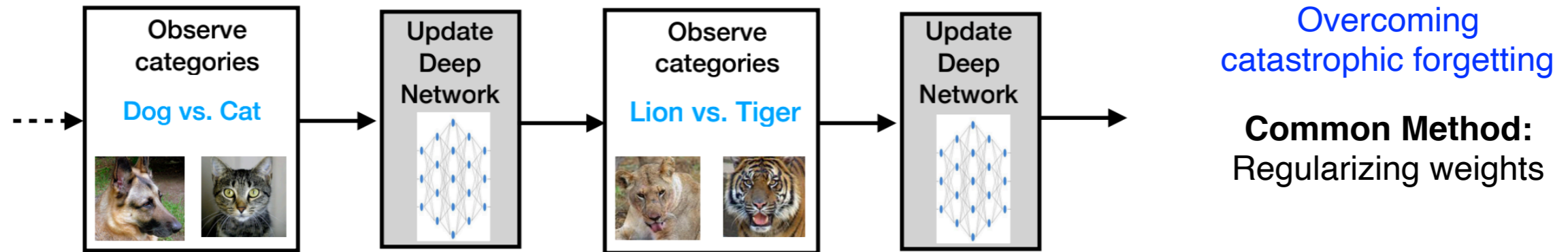
$$\min_{q_t(\mathbf{w})} \mathbb{E}_{q_t(\mathbf{w})} \left[ \sum_{i \in D^t} \ell(y_i^t, f_{\mathbf{w}}(\mathbf{x}_i^t)) \right] + \mathbb{D}_{KL}(q_t(\mathbf{w}) || p(\mathbf{w}))$$

Prior Distribution (uniform)

1. Kirkpatrick, J. et al. Overcoming catastrophic forgetting in neural networks. PANS, 114(13):3521–3526, 2017.

# Uncertainty Provided by BayesBiNN Enables Continual Learning

- Main point: **BayesBiNN enables continual learning (CL) for BiNN using the intrinsic KL divergence as regularization**
- CL: Sequentially learning new tasks without forgetting old ones<sup>[1]</sup>



- But, it is unclear how to regularize binary weights of BiNN using STE/Bop
- In BayesBiNN, there is one natural solution using KL divergence

Continual Learning

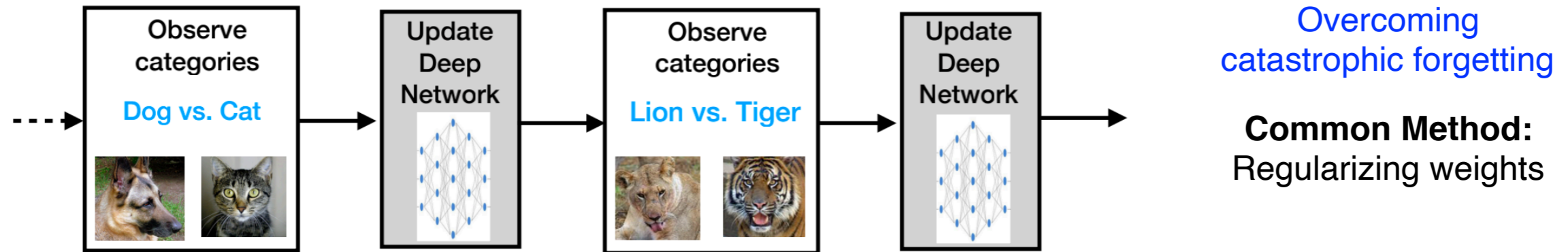
$$\min_{q_t(\mathbf{w})} \mathbb{E}_{q_t(\mathbf{w})} \left[ \sum_{i \in D^t} \ell(y_i^t, f_{\mathbf{w}}(\mathbf{x}_i^t)) \right] + \mathbb{D}_{KL} (q_t(\mathbf{w}) \parallel q_{t-1}(\mathbf{w}))$$

posterior approximation after task  $t - 1$

1. Kirkpatrick, J. et al. Overcoming catastrophic forgetting in neural networks. PANS, 114(13):3521–3526, 2017.

# Uncertainty Provided by BayesBiNN Enables Continual Learning

- Main point: **BayesBiNN enables continual learning (CL) for BiNN using the intrinsic KL divergence as regularization**
- CL: Sequentially learning new tasks without forgetting old ones<sup>[1]</sup>



- But, it is unclear how to regularize binary weights of BiNN using STE/Bop
- In BayesBiNN, there is one natural solution using KL divergence

Continual Learning

$$\min_{q_t(\mathbf{w})} \mathbb{E}_{q_t(\mathbf{w})} \left[ \sum_{i \in D^t} \ell(y_i^t, f_{\mathbf{w}}(\mathbf{x}_i^t)) \right] + \mathbb{D}_{KL} (q_t(\mathbf{w}) \parallel q_{t-1}(\mathbf{w}))$$

posterior approximation after task  $t - 1$

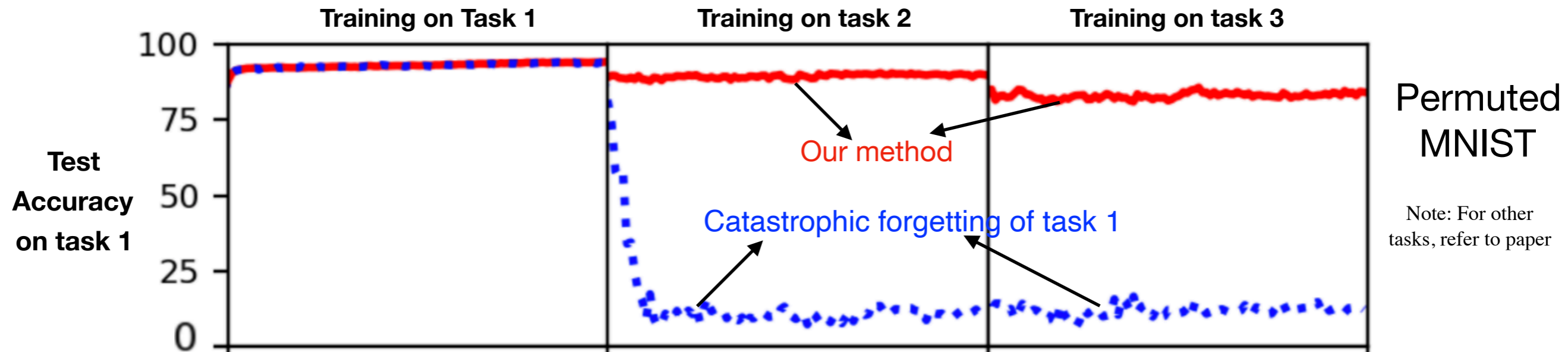
Learned natural parameter after task  $t - 1$

$$\lambda_t \leftarrow (1 - \alpha)\lambda_t - \alpha [\mathbf{s} \odot \mathbf{g} - \lambda_{t-1}]$$

1. Kirkpatrick, J. et al. Overcoming catastrophic forgetting in neural networks. PANS, 114(13):3521–3526, 2017.

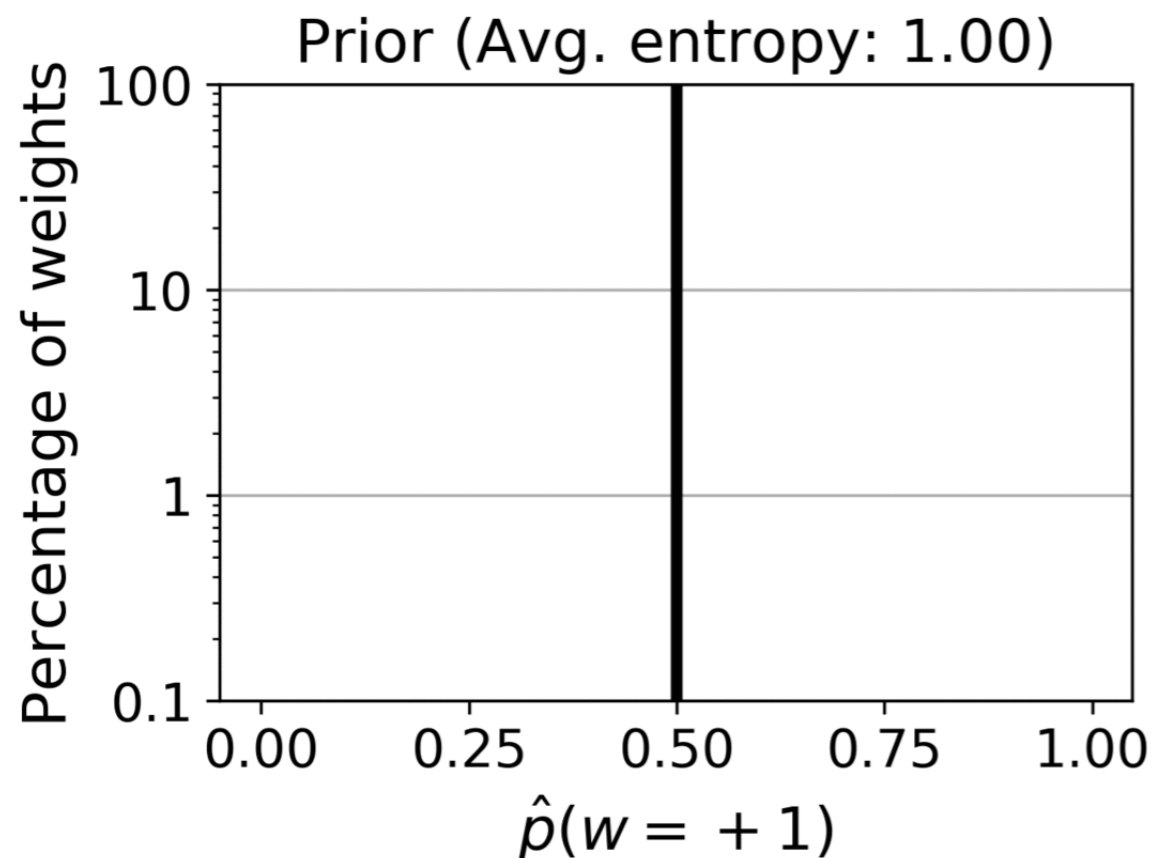
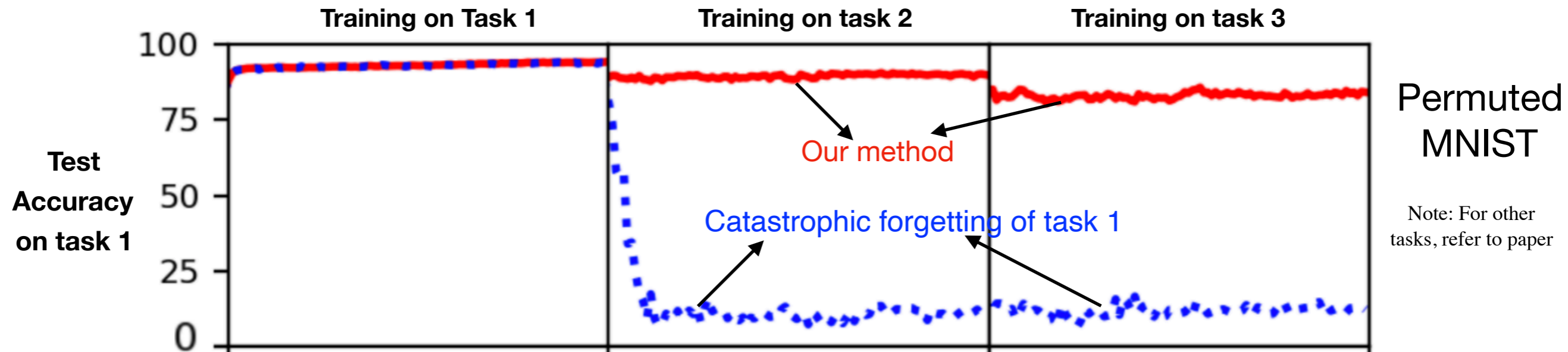
# Uncertainty Provided by BayesBiNN Enables Continual Learning

- Main point: **BayesBiNN** avoids the catastrophic forgetting problem



# Uncertainty Provided by BayesBiNN Enables Continual Learning

- Main point: **BayesBiNN avoids the catastrophic forgetting problem**



- As the number of tasks increases, **the distribution over binary weights become more and more deterministic**
- Open-source Code Available : <https://github.com/team-approx-bayes/BayesBiNN>

# Summary

- BiNN: Neural Networks with **binary weights**
- Much faster and much smaller but difficult to optimize
- Gradient based methods work well but not well understood
- We proposed a principled approach to train BiNN using the Bayesian Learning Rule, which can justify some previous approaches
- The Bayesian approach also gives us estimate of uncertainty which can be used for continual learning



Thank you!  
Q&A