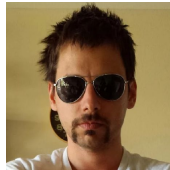
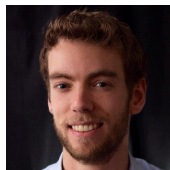
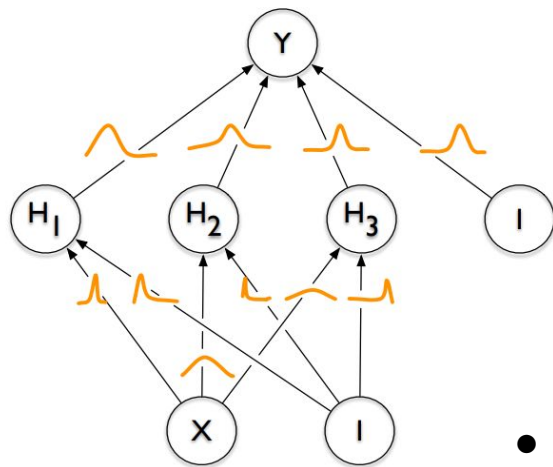


# The K-tied Normal Distribution: A Compact Parameterization of Gaussian Mean Field Posteriors in Bayesian Neural Networks

J. Świątkowski, K. Roth, B. Veeling, L. Tran, J. Dillon, J. Snoek, S. Mandt,  
T. Salimans, R. Jenatton, S. Nowozin



# Let's bring the benefits of Bayesian inference to neural networks!



- noisy estimates of gradients
- slow convergence
- increased number of model parameters

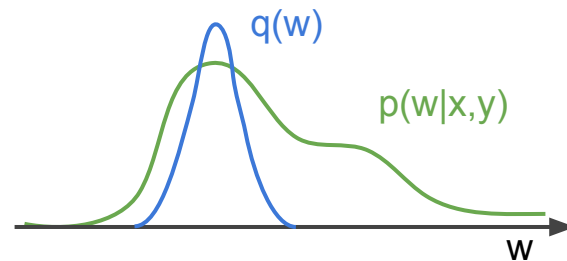


# Preview

1. Gaussian Mean-Field Variational Inference (GMFVI) for Bayesian Neural Networks (BNNs).
2. Low-rank in already trained GMFVI BNNs.
3. Training a low-rank parameterization of the GMFVI BNNs.

# BNN variational posterior

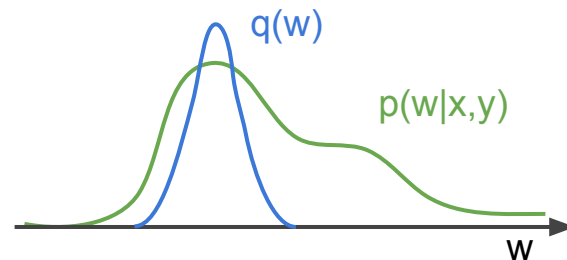
$$\theta^* = \operatorname{argmin}_{\theta} D_{\text{KL}}[q_{\theta}(\mathbf{w}) || p(\mathbf{w} | \mathbf{x}, \mathbf{y})]$$



- Variational inference: Cast inference as an optimization problem.

# BNN variational posterior

$$\theta^* = \operatorname{argmin}_{\theta} D_{\text{KL}}[q_{\theta}(\mathbf{w}) || p(\mathbf{w} | \mathbf{x}, \mathbf{y})]$$

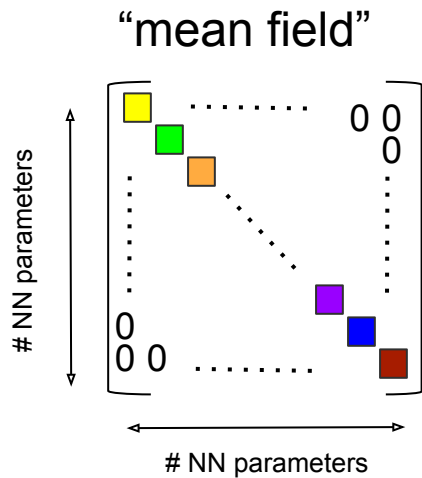


- Variational inference: Cast inference as an optimization problem.
- Key question: Which parametrization?

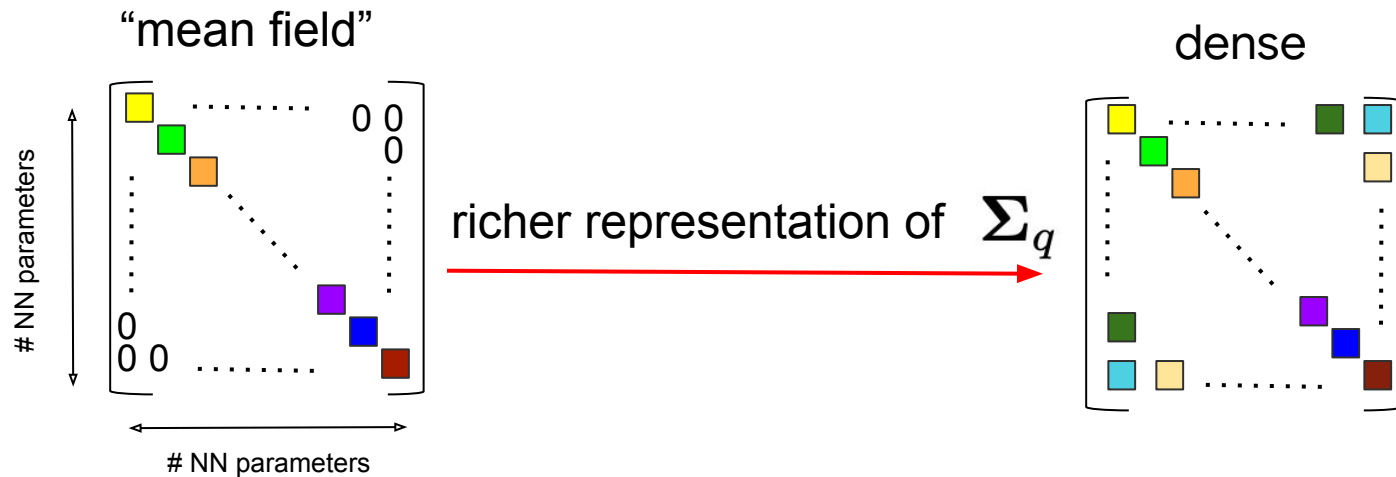
$$q(\mathbf{W}) = \mathcal{N}(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q)$$


for BNNs, can be a prohibitively large object

# Parametrization of variational posterior covariance

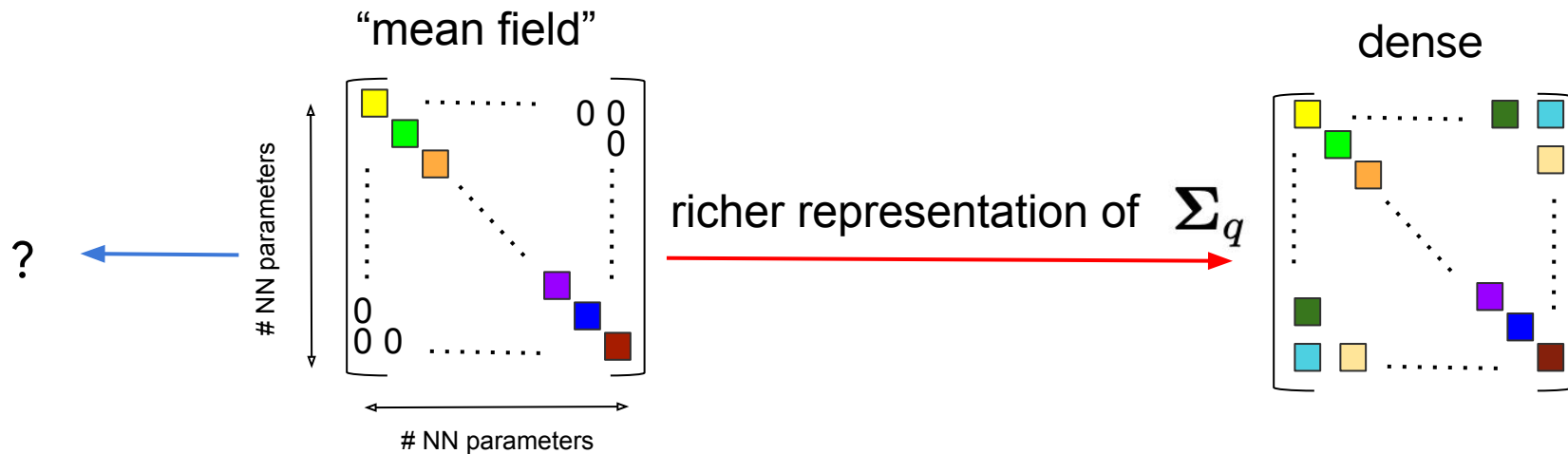


# Parametrization of variational posterior covariance



- Lot of research exploring “”:
  - E.g., Barber & Bishop, 1998..., Zhang et al. 2017, Sun et al. 2017, Mishkin et al., 2018

# Parametrization of variational posterior covariance



- Lot of research exploring “ $\rightarrow$ ”:
  - E.g., Barber & Bishop, 1998..., Zhang et al. 2017, Sun et al. 2017, Mishkin et al., 2018
- We investigate the opposite trend “ $\leftarrow$ ”:
  - Fewer parameters to optimize
  - Less noisy gradient estimate and convergence speed-up

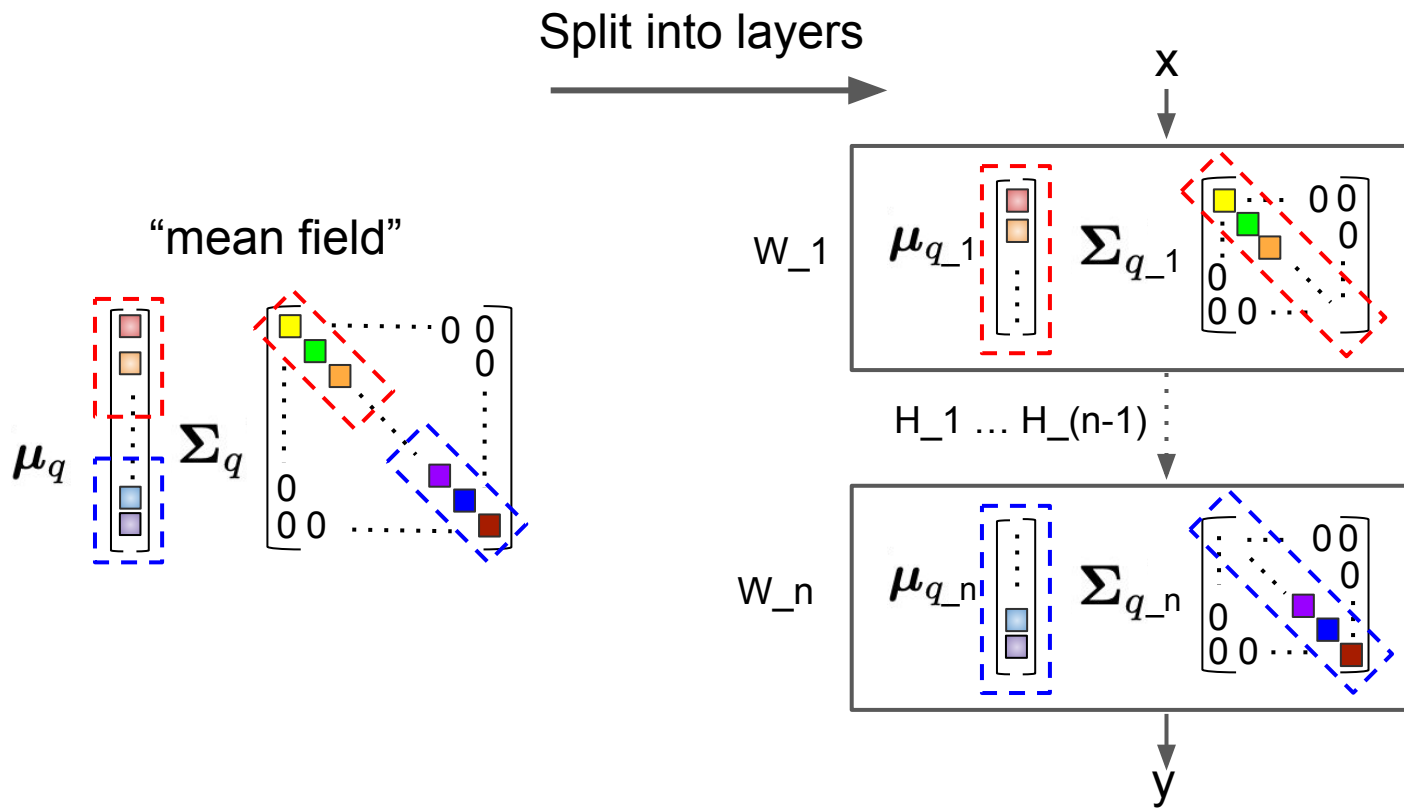
How? Exploit low-rank structure!



# Preview

1. Gaussian Mean-Field Variational Inference (GMFVI) for Bayesian Neural Networks (BNNs).
2. Low-rank in already trained GMFVI BNNs.
3. Training a low-rank parameterization of the GMFVI BNNs.

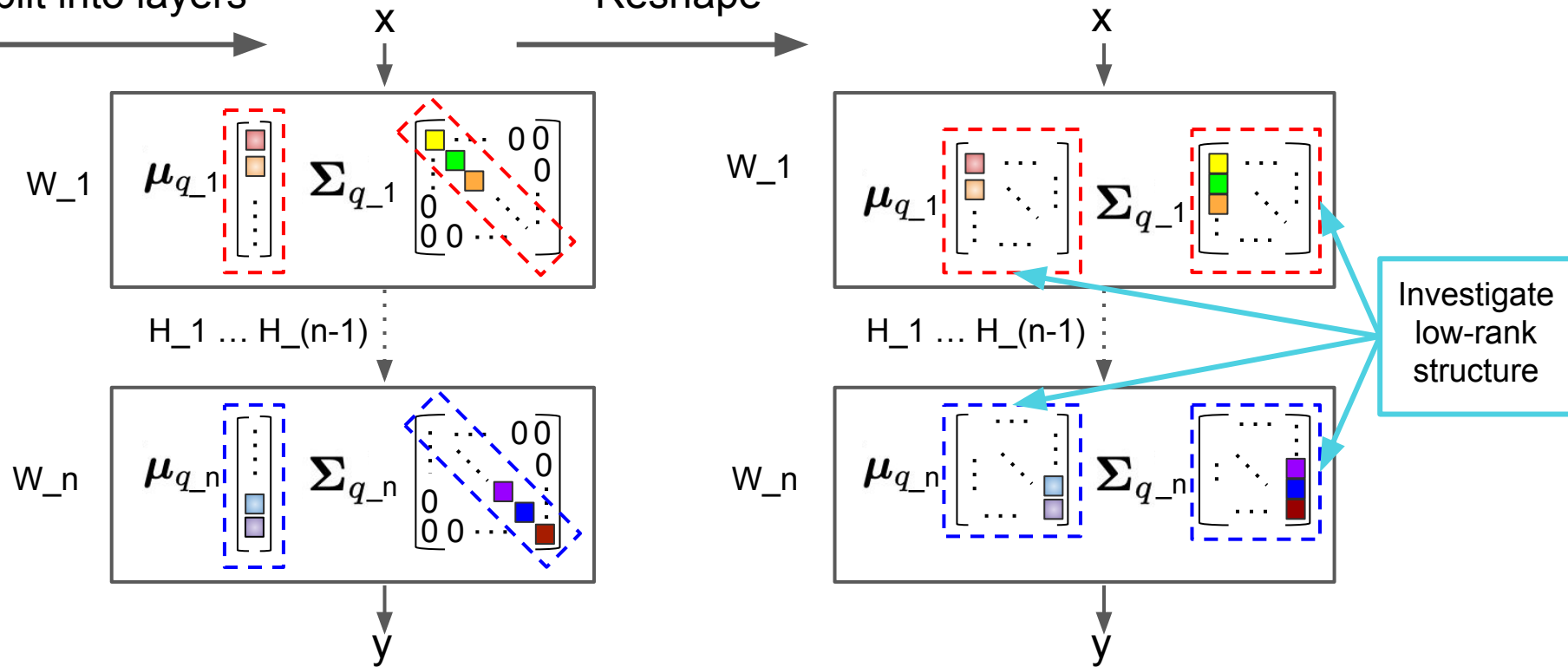
# Low-rank structure in mean-field variational posteriors



# Low-rank structure in mean-field variational posteriors

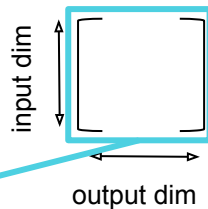
Split into layers

Reshape



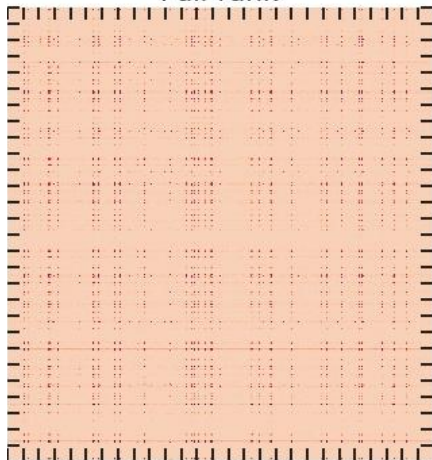
# Low-rank structure in mean-field variational posteriors

Reshaped diagonal  $\Sigma_q$ :

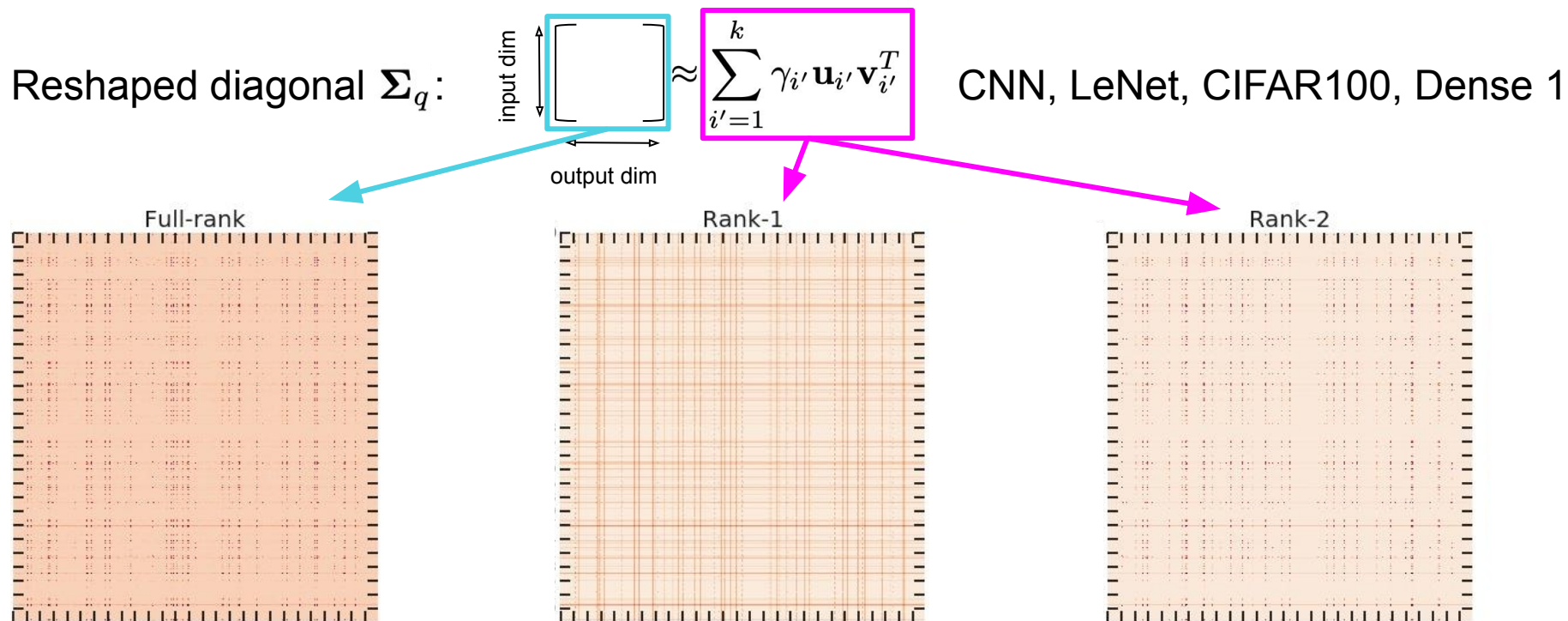


CNN, LeNet, CIFAR100, Dense 1

Full-rank



# Low-rank structure in mean-field variational posteriors



# Low-rank structure in mean-field variational posteriors

Fraction of explained variance of both  $\boldsymbol{\mu}_q$  and diagonal  $\boldsymbol{\Sigma}_q$ :

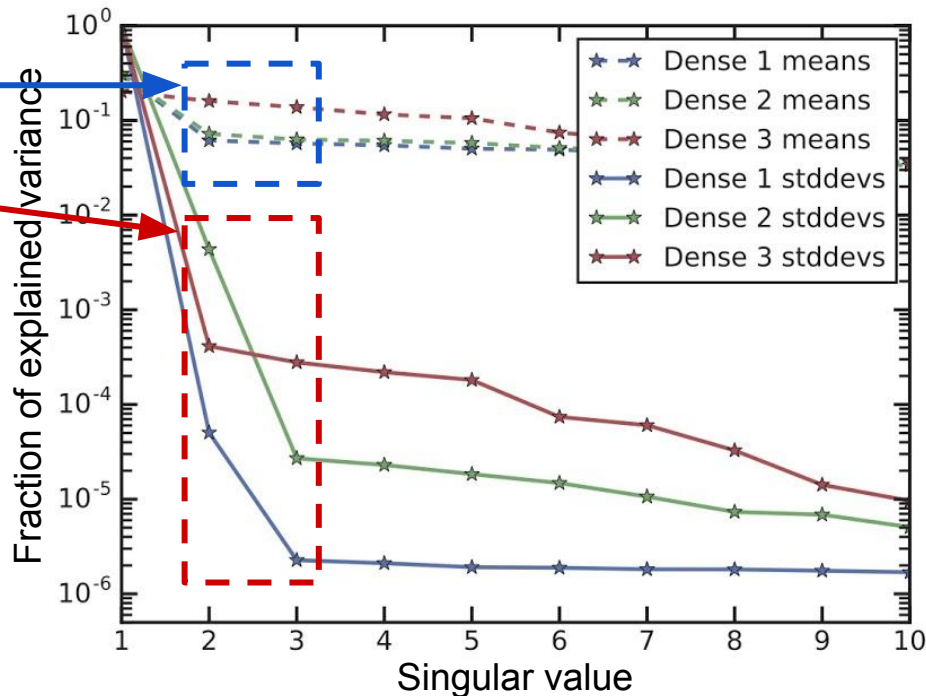
$$\gamma_k^2 / \sum_{i'} \gamma_{i'}^2$$

# Low-rank structure in mean-field variational posteriors

Fraction of explained variance for both  $\mu_q$  and diagonal  $\Sigma_q$

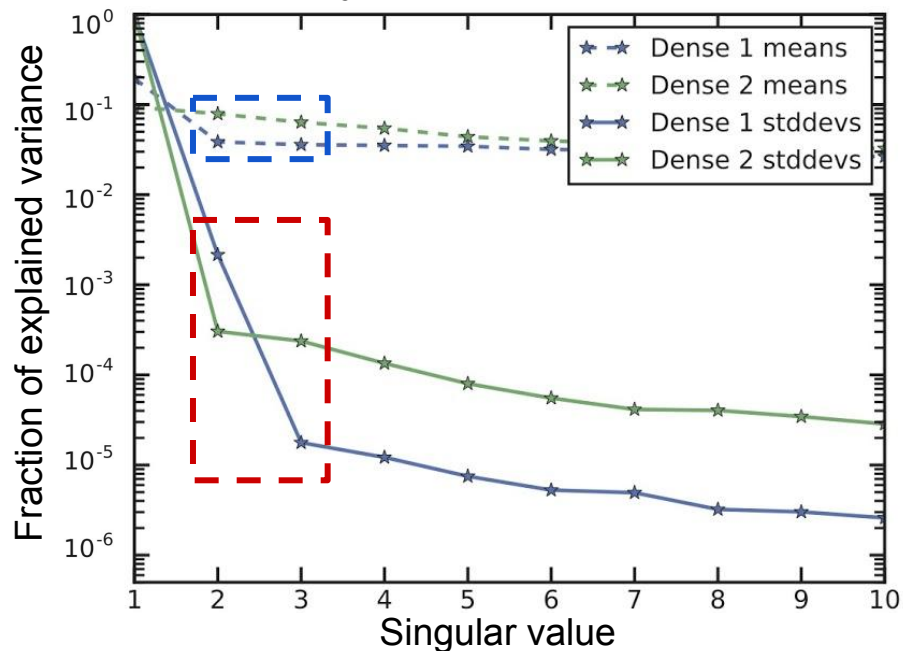
$$\gamma_k^2 / \sum_{i'} \gamma_{i'}^2$$

Dense layers of MLP, MNIST

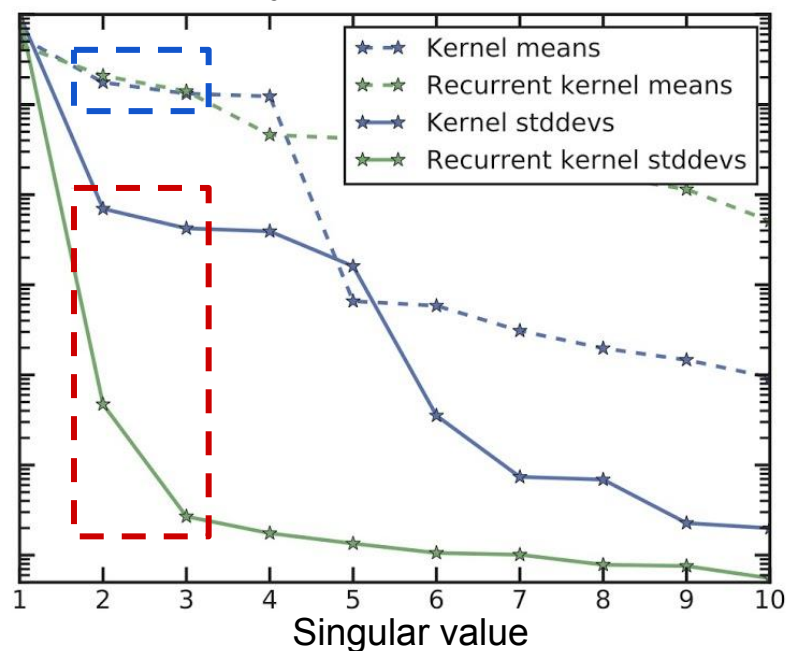


# Low-rank structure in mean-field variational posteriors

Dense layers of LeNet, CIFAR100



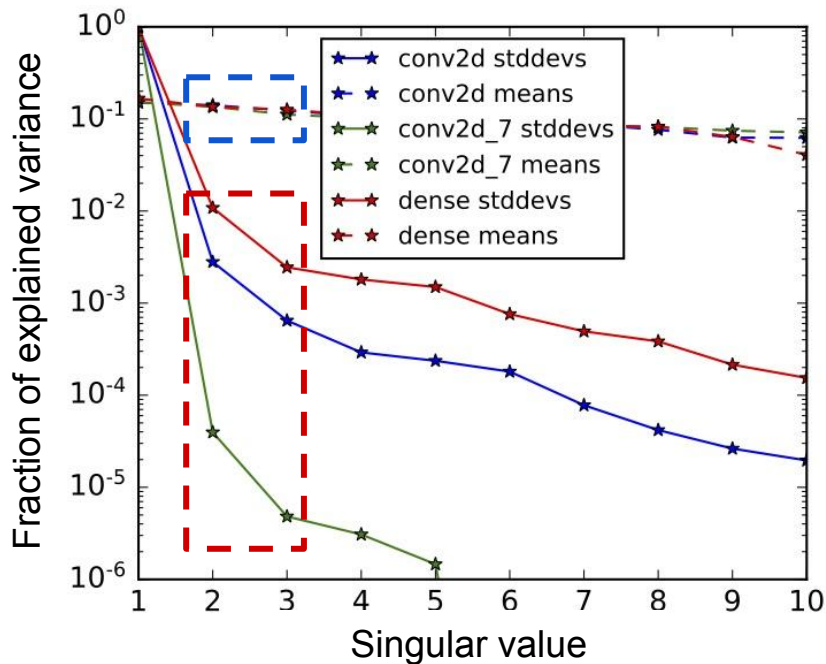
Dense layers of LSTM, IMDB





# Low-rank structure in mean-field variational posteriors

Dense and conv layers of ResNet-18, CIFAR10



# Post-training low-rank approximation

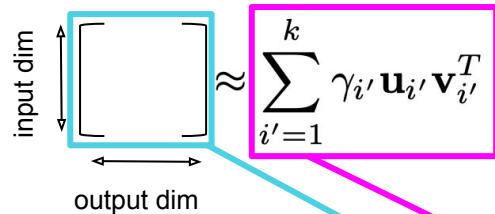
Rank-k approximation of  
diagonal  $\Sigma_q$ :

MLP, MNIST

$$\begin{matrix} \text{input dim} \\ \updownarrow \\ \left[ \phantom{\Sigma_q} \right] \\ \leftarrow \text{output dim} \end{matrix} \approx \sum_{i'=1}^k \gamma_{i'} \mathbf{u}_{i'} \mathbf{v}_{i'}^T$$

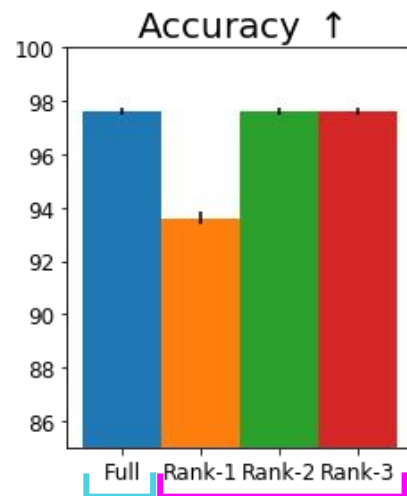
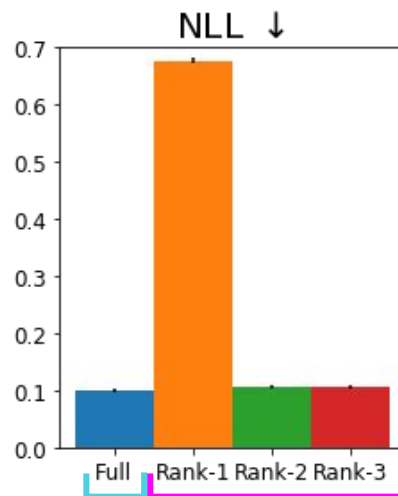
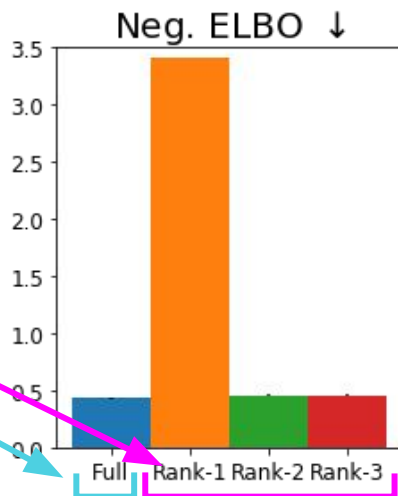
# Post-training low-rank approximation

Rank-k approximation of diagonal  $\Sigma_q$ :



No approximation /  
mean-field / diagonal  $\Sigma_q$

MLP, MNIST



# Post-training low-rank approximation

Dense layers of CNN (LeNet, CIFAR100) and LSTM (IMDB):

CNN			LSTM			
Rank	-ELBO ↓	NLL ↓	Accuracy ↑	-ELBO ↓	NLL ↓	Accuracy ↑
Full	3.83 $\pm$ 0.020	2.23 $\pm$ 0.017	42.1 $\pm$ 0.49	0.536 $\pm$ 0.0058	0.493 $\pm$ 0.0057	80.1 $\pm$ 0.25
1	4.33 $\pm$ 0.021	2.30 $\pm$ 0.016	41.7 $\pm$ 0.49	0.687 $\pm$ 0.0058	0.491 $\pm$ 0.0056	80.0 $\pm$ 0.25
2	3.88 $\pm$ 0.020	2.24 $\pm$ 0.017	42.2 $\pm$ 0.49	0.621 $\pm$ 0.0058	0.494 $\pm$ 0.0057	80.1 $\pm$ 0.25
3	3.86 $\pm$ 0.020	2.24 $\pm$ 0.017	42.1 $\pm$ 0.49	0.595 $\pm$ 0.0058	0.493 $\pm$ 0.0056	80.1 $\pm$ 0.25

Dense and convolutional layers of a ResNet-18 (CIFAR10):

Rank	-ELBO ↓	NLL ↓	Accuracy ↑
Full	122.61 $\pm$ 0.012	0.495 $\pm$ 0.0080	83.5 $\pm$ 0.37
1	122.57 $\pm$ 0.012	0.658 $\pm$ 0.0069	81.7 $\pm$ 0.39
2	122.77 $\pm$ 0.012	0.503 $\pm$ 0.0080	83.2 $\pm$ 0.37
3	122.67 $\pm$ 0.012	0.501 $\pm$ 0.0079	83.2 $\pm$ 0.37

# Generality of the low-rank structure finding

Low-rank structure in posterior standard deviations holds for:

- Different model types:

- MLP
- CNN
- LSTM

- Different mode sizes:

- Small 3 layer MLP
- Large ResNet-18

- Different layer types:

- Dense
- Convolutional



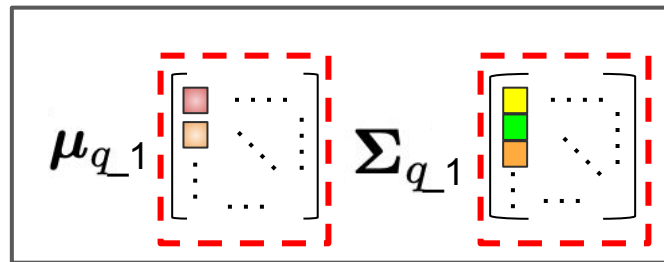
Suggests generality of the  
low-rank structure finding

# Preview

1. Gaussian Mean-Field Variational Inference (GMFVI) for Bayesian Neural Networks (BNNs).
2. Low-rank in already trained GMFVI BNNs.
3. Training a low-rank parameterization of the GMFVI BNNs.

# K-tied Normal distribution

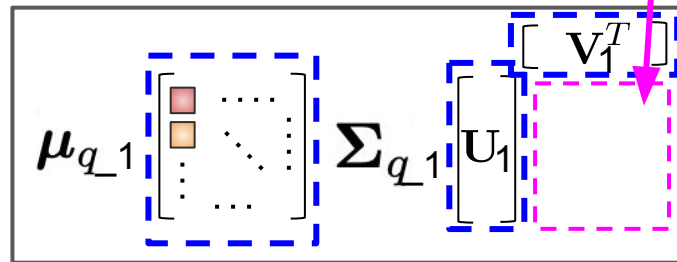
"Reshaped" mean-field



$W_{-1}$  Parameters:  $\mu_{q-1}, \Sigma_{q-1}$



K-tied Normal distribution

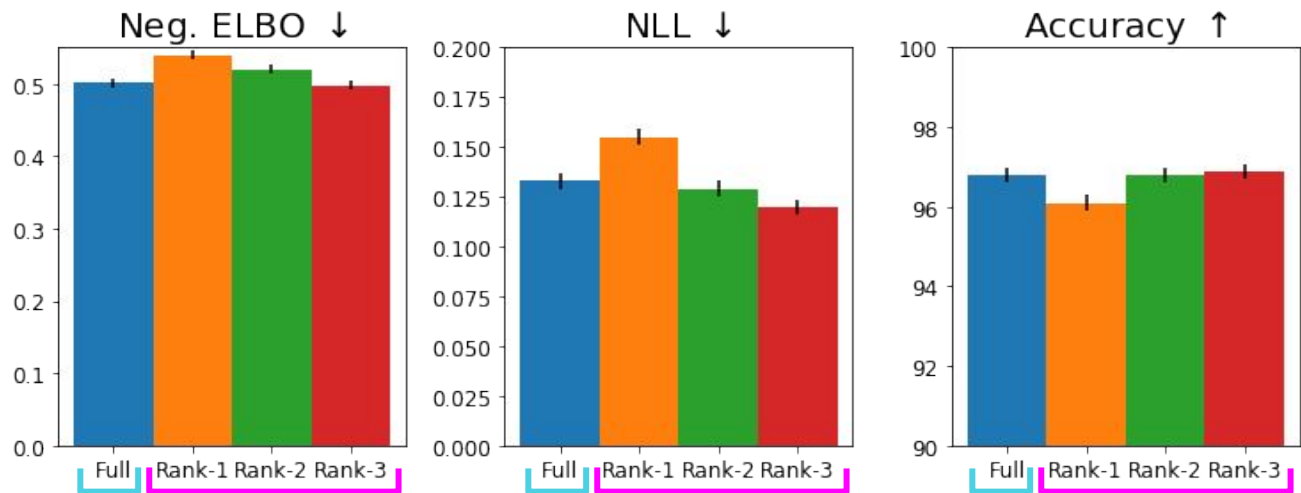


$W_{-1}$  Parameters:  $\mu_{q-1}, \mathbf{U}_1, \mathbf{V}_1$

$$\sum_{i'=1}^k \mathbf{u}_{i'} \mathbf{v}_{i'}^T$$

# K-tied Normal distribution - training performance

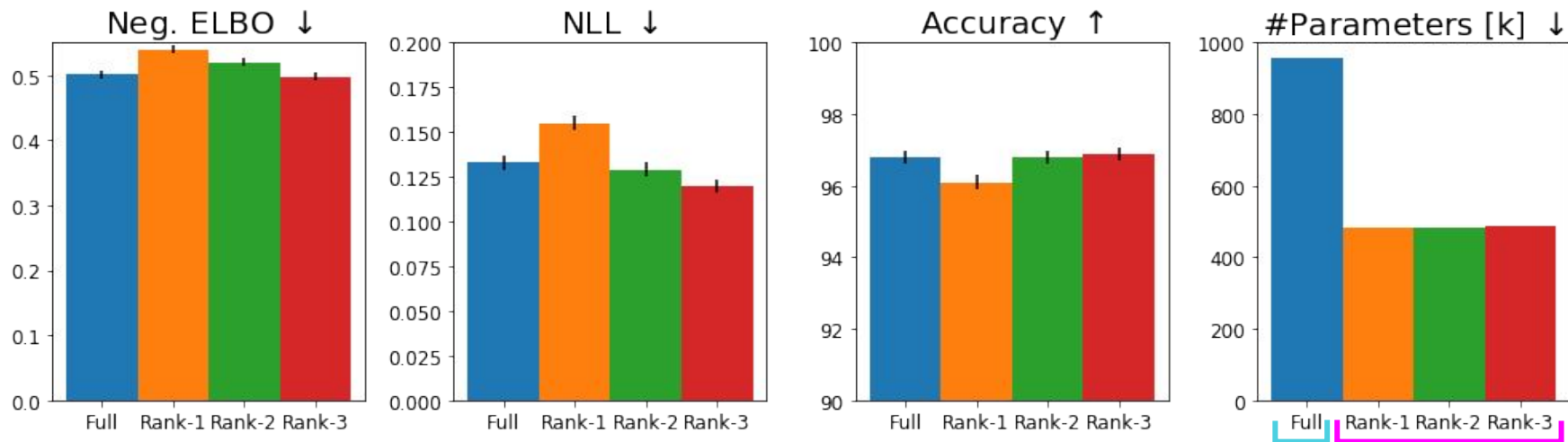
MLP, MNIST





# K-tied Normal distribution - training performance

MLP, MNIST



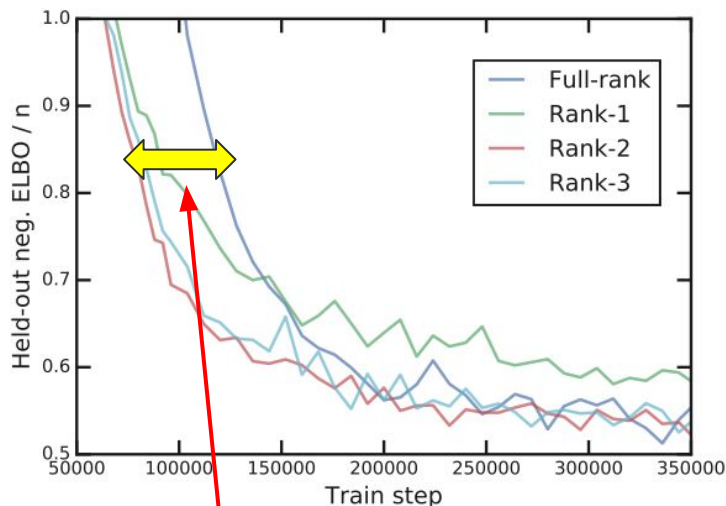
# K-tied Normal increases gradient SNR

$$\text{SNR} = E[g_b^2] / \text{Var}[g_b]$$

Rank $k$	MNIST, MLP Dense 2, SNR at step		
	1000	5000	9000
full	4.13 $\pm$ 0.027	4.45 $\pm$ 0.091	3.21 $\pm$ 0.035
1	5840 $\pm$ 190	158 $\pm$ 3.8	5.3 $\pm$ 0.20
2	7500 $\pm$ 240	140 $\pm$ 11	4.3 $\pm$ 0.26
3	7000 $\pm$ 270	117 $\pm$ 1.7	4.1 $\pm$ 0.20

# K-tied Normal speeds up convergence

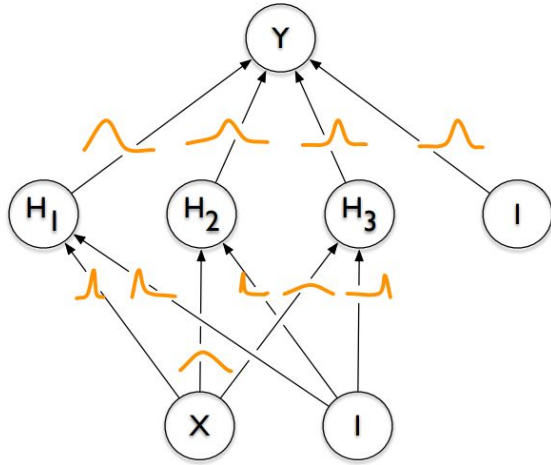
MLP (MNIST)



Increased convergence speed

Rank $k$	MNIST, MLP, -ELBO at step		
	1000	5000	9000
full	42.16 $\pm$ 0.070	26.52 $\pm$ 0.016	15.39 $\pm$ 0.016
1	43.11 $\pm$ 0.039	14.85 $\pm$ 0.017	2.06 $\pm$ 0.027
2	42.74 $\pm$ 0.090	13.97 $\pm$ 0.023	1.82 $\pm$ 0.017
3	42.63 $\pm$ 0.068	13.61 $\pm$ 0.020	1.80 $\pm$ 0.031

# Let's bring the benefits of Bayesian inference to neural networks



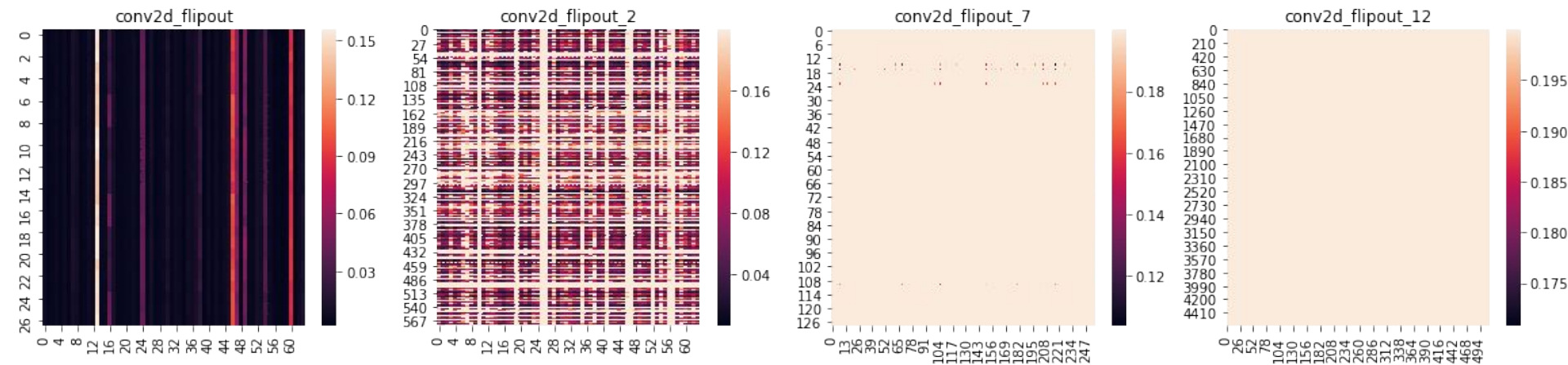
Thank you!

# Review

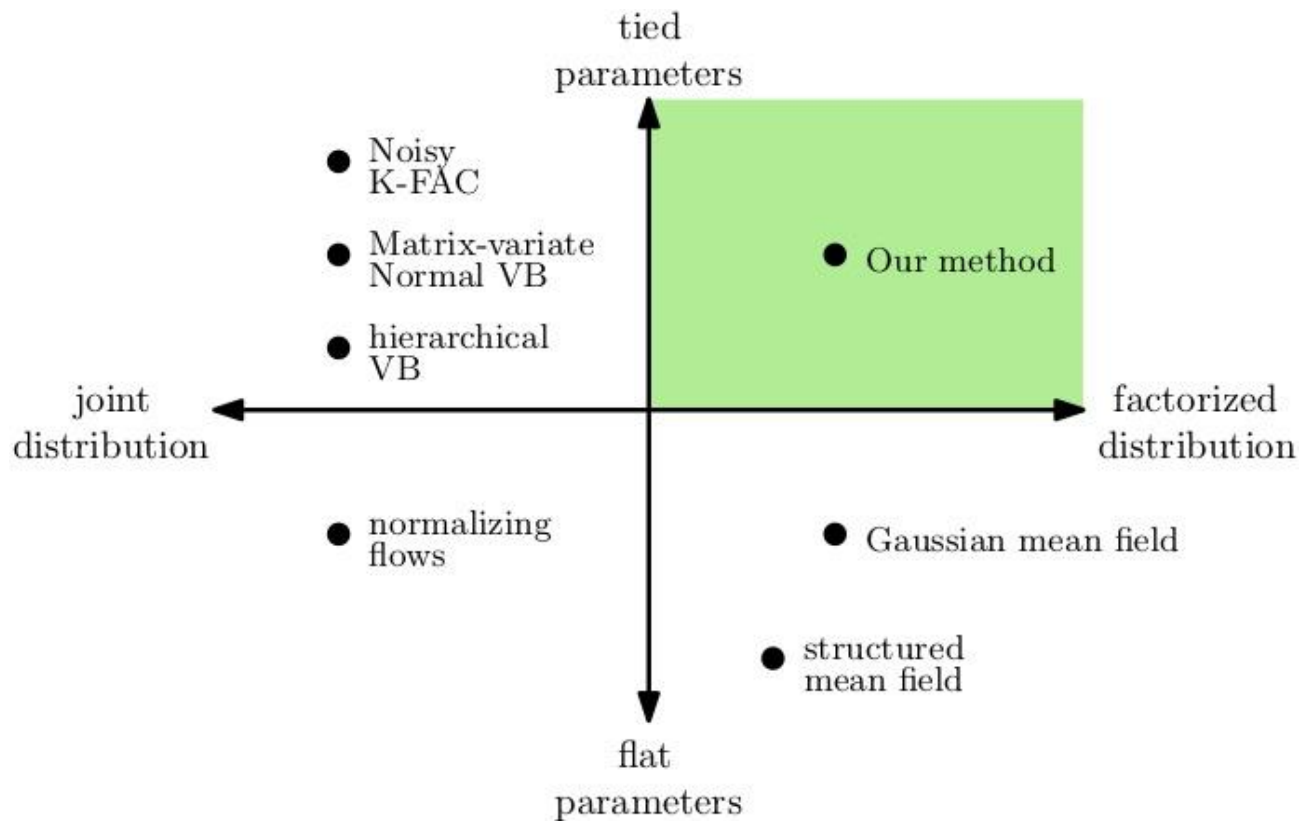
1. Gaussian Mean-Field Variational Inference (GMFVI) for Bayesian Neural Networks (BNNs).
2. Low-rank in already trained GMFVI BNNs.
3. Training a low-rank parameterization of the GMFVI BNNs.

# Heatmaps of conv posterior stddevs

Reshaping conv layers e.g.: [3, 3, 10, 20] -> [3 \* 3 \* 10, 20].

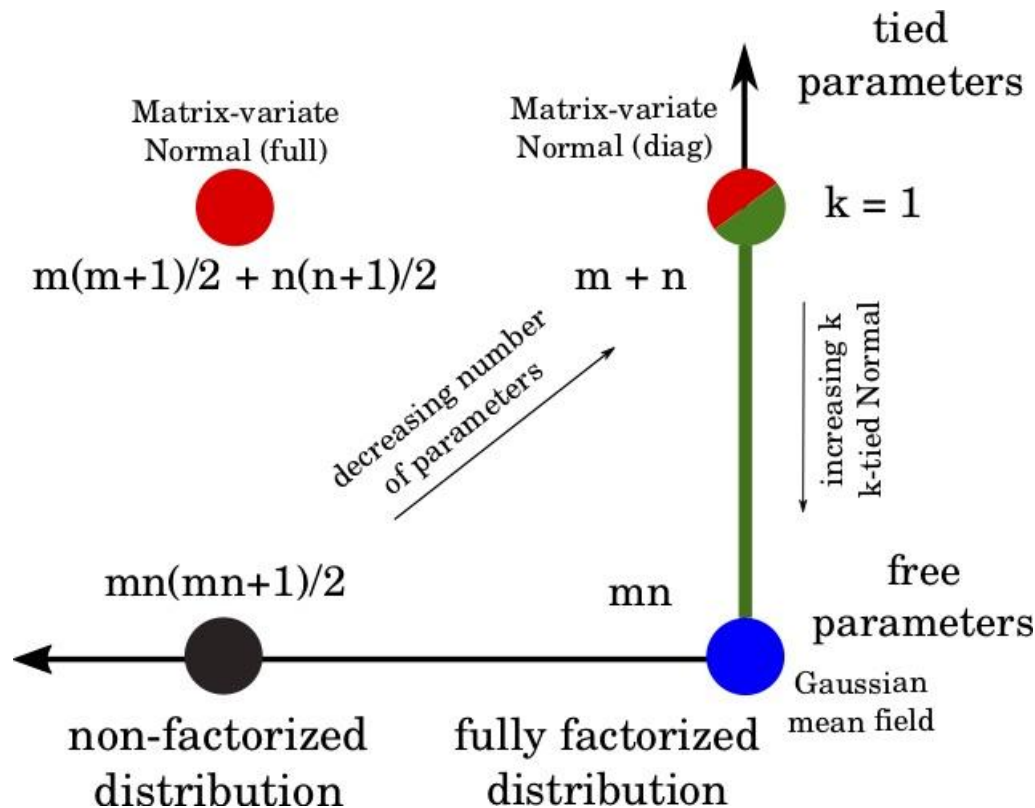


# Relation to other work

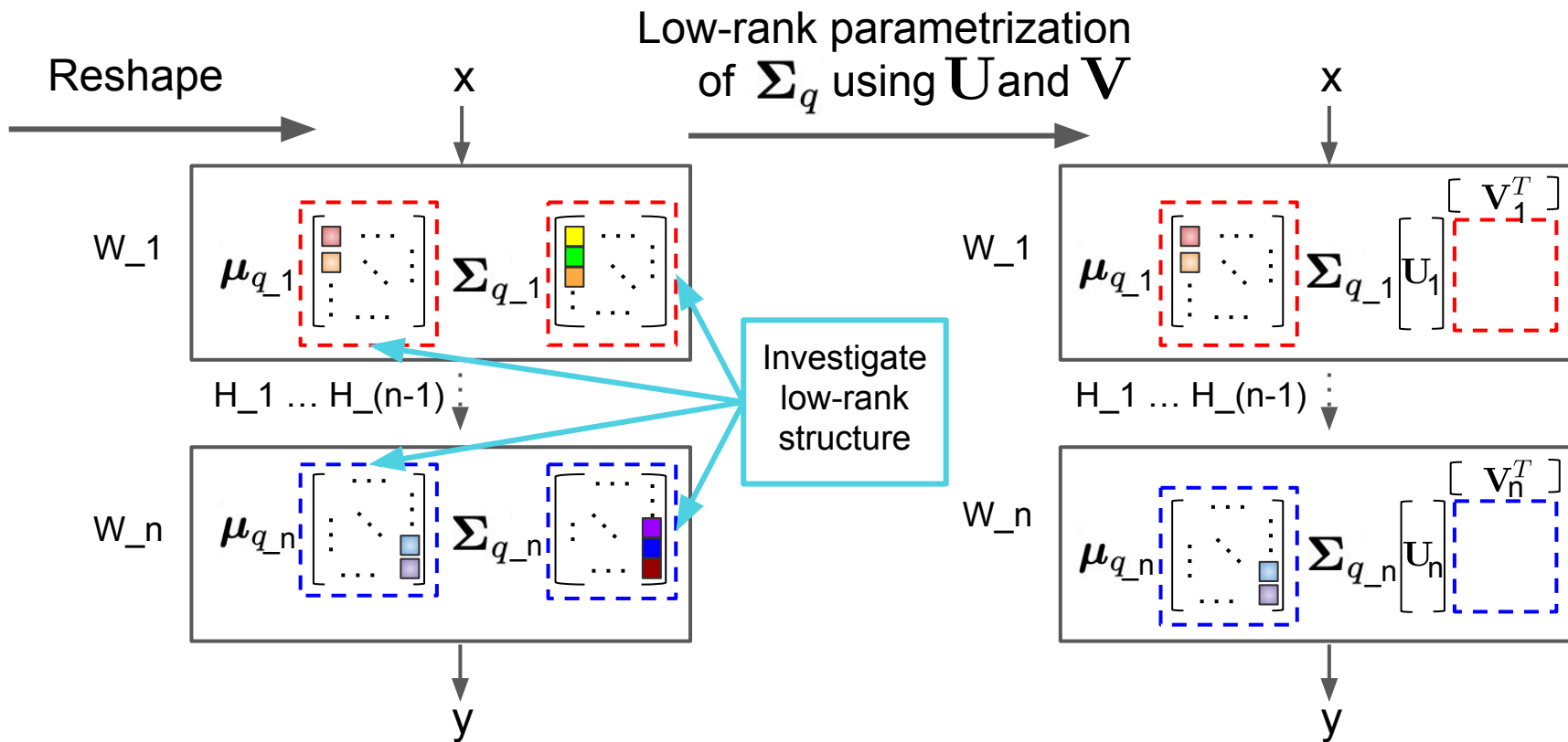




# Relation to Matrix-Variate Normal



# Low-rank structure in mean-field variational posteriors

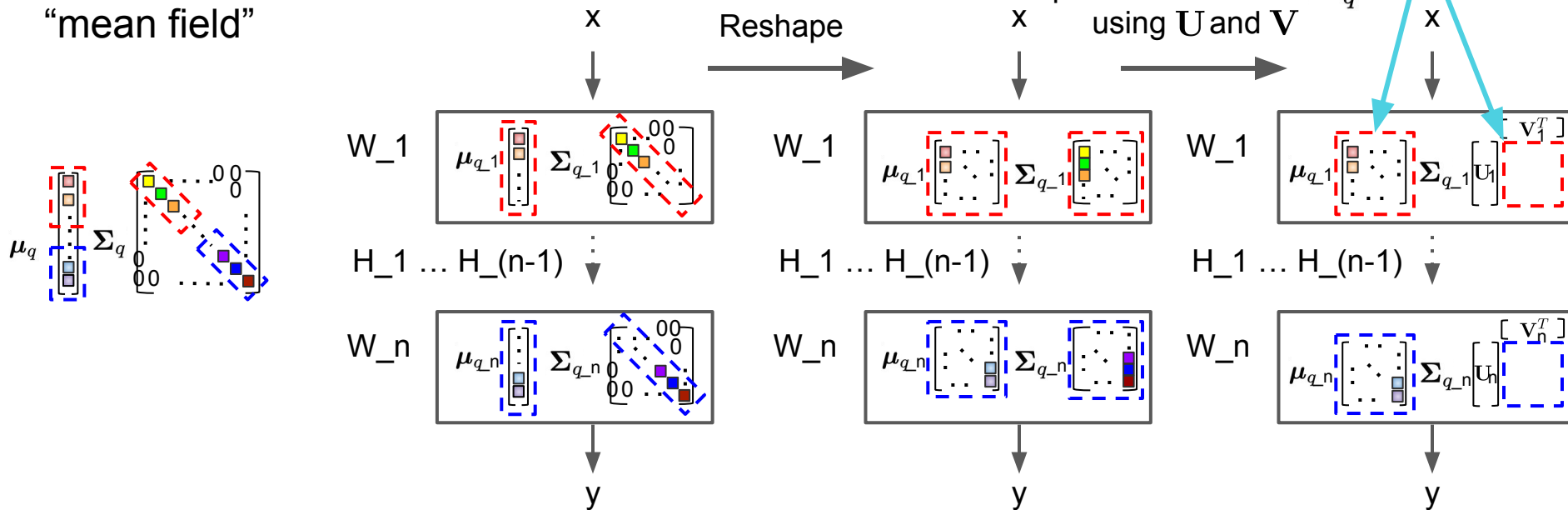


# Low-rank structure in mean-field variational posteriors

Idea: We further exploit a low-rank structure in variational posteriors

Investigate low-rank structure!

“mean field”



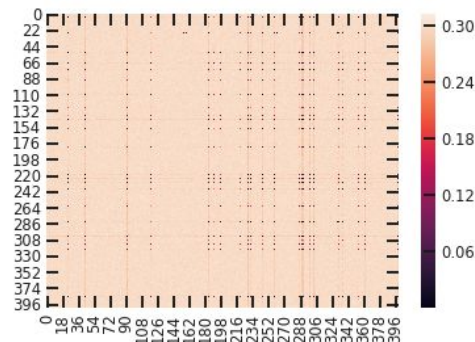
We investigate per layer GMFVI posterior matrices

$$\mathbf{a}_l = \mathbf{h}_l \mathbf{W}_l + \mathbf{b}_l, \quad \mathbf{h}_{l+1} = f(\mathbf{a}_l), \quad \mathbf{W}_l \in \mathbb{R}^{m \times n}$$

$$q(\mathbf{W}) = \mathcal{N}(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q) = \prod_{i=1}^m \prod_{j=1}^n q(w_{ij}), \quad \text{with} \quad q(w_{ij}) = \mathcal{N}(\mu_{ij}, \sigma_{ij}^2),$$

$$\boldsymbol{\mu}_q = \text{vec}(\mathbf{M}), \quad \mathbf{M} \in \mathbb{R}^{m \times n}$$

$$\boldsymbol{\Sigma}_q = \text{diag}(\text{vec}(\mathbf{A}^2)) \quad \mathbf{A} \in \mathbb{R}_+^{m \times n} \longrightarrow$$



We propose a  $k$ -tied Normal variational posterior that exploits the low-rank structure

$$q(\mathbf{W}) = \mathcal{N}(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q)$$

$$\boldsymbol{\Sigma}_q = \text{diag}(\text{vec}(\mathbf{A}^2)) \quad \mathbf{A} \in \mathbb{R}_+^{m \times n}$$

$$\mathbf{A} \approx \mathbf{U}\mathbf{V}^T$$

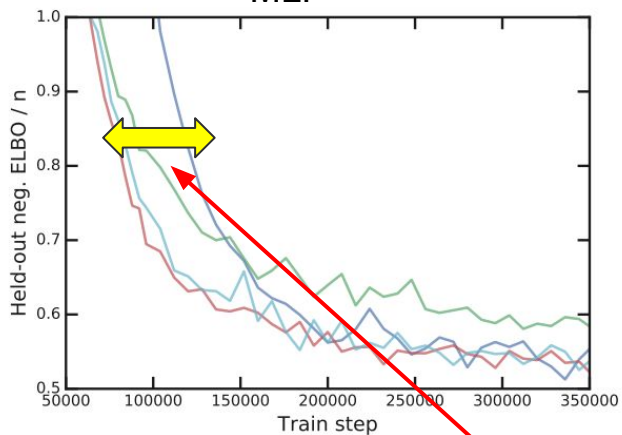
$$k\text{-tied-}\mathcal{N}(\mathbf{W}; \boldsymbol{\mu}_q, \mathbf{U}, \mathbf{V}) = \mathcal{N}(\boldsymbol{\mu}_q, \text{diag}(\text{vec}((\mathbf{U}\mathbf{V}^T)^2))),$$

# K-tied Normal posterior reduces the number of parameters without reducing performance

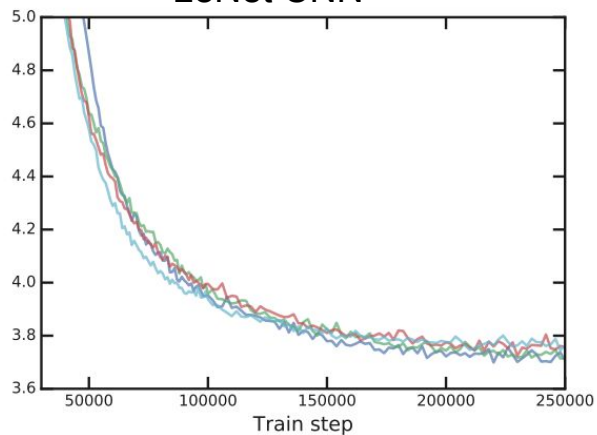
Model & Dataset	Rank $k$	-ELBO $\downarrow$	NLL $\downarrow$	Accuracy $\uparrow$	#Par. [k] $\downarrow$
MNIST, MLP	full	0.501 $\pm$ 0.0061	0.133 $\pm$ 0.0040	96.8 $\pm$ 0.18	957
MNIST, MLP	1	0.539 $\pm$ 0.0063	0.155 $\pm$ 0.0043	96.1 $\pm$ 0.19	482
MNIST, MLP	2	0.520 $\pm$ 0.0063	0.129 $\pm$ 0.0039	96.8 $\pm$ 0.18	484
MNIST, MLP	3	0.497 $\pm$ 0.0060	0.120 $\pm$ 0.0038	96.9 $\pm$ 0.18	486
CIFAR100, CNN	full	3.72 $\pm$ 0.018	2.16 $\pm$ 0.016	43.9 $\pm$ 0.50	4,405
CIFAR100, CNN	1	3.65 $\pm$ 0.017	2.12 $\pm$ 0.015	45.5 $\pm$ 0.50	2,262
CIFAR100, CNN	2	3.76 $\pm$ 0.019	2.15 $\pm$ 0.016	44.3 $\pm$ 0.50	2,268
CIFAR100, CNN	3	3.73 $\pm$ 0.018	2.13 $\pm$ 0.016	44.3 $\pm$ 0.50	2,273
IMDB, LSTM	full	0.538 $\pm$ 0.0054	0.478 $\pm$ 0.0052	79.5 $\pm$ 0.26	2,823
IMDB, LSTM	1	0.592 $\pm$ 0.0041	0.512 $\pm$ 0.0040	77.6 $\pm$ 0.26	2,693
IMDB, LSTM	2	0.560 $\pm$ 0.0042	0.484 $\pm$ 0.0041	78.2 $\pm$ 0.26	2,694
IMDB, LSTM	3	0.550 $\pm$ 0.0051	0.491 $\pm$ 0.0050	78.8 $\pm$ 0.26	2,695

# K-tied Normal speeds up convergence

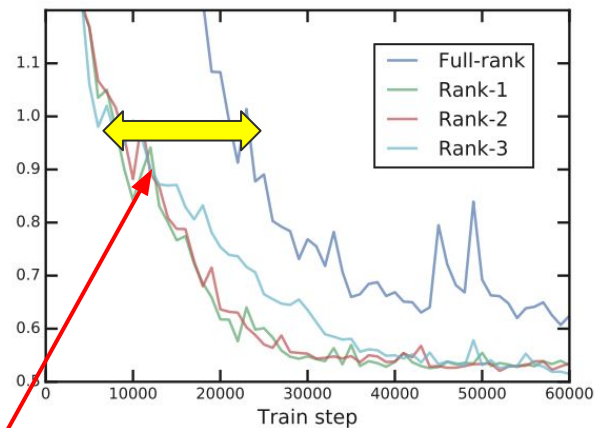
MLP



LeNet CNN



LSTM



Increased convergence speed