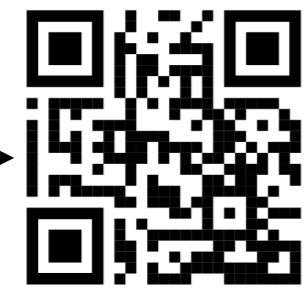


# BMRS: Bayesian Model Reduction for Structured Pruning

Dustin Wright, Christian Igel, Raghavendra Selvan

dw@di.ku.dk; I'm on the job market! .....



**BACKGROUND:** We want to balance accuracy and compression in a principled way. We do so with variational inference and Bayesian Model Reduction (BMR).

## BMRS in a nutshell

- Apply multiplicative noise to network structures (Neklyudov et al. 2017)

$$h_i = \theta_i \cdot (w_i h_{i-1})$$

$$q_\phi(\theta_i) = \text{LogN}_{[a,b]}(\theta_i | \mu_i, \sigma_i^2)$$

$$p(\theta_i) = \text{LogU}_{[a,b]}(\theta_i)$$

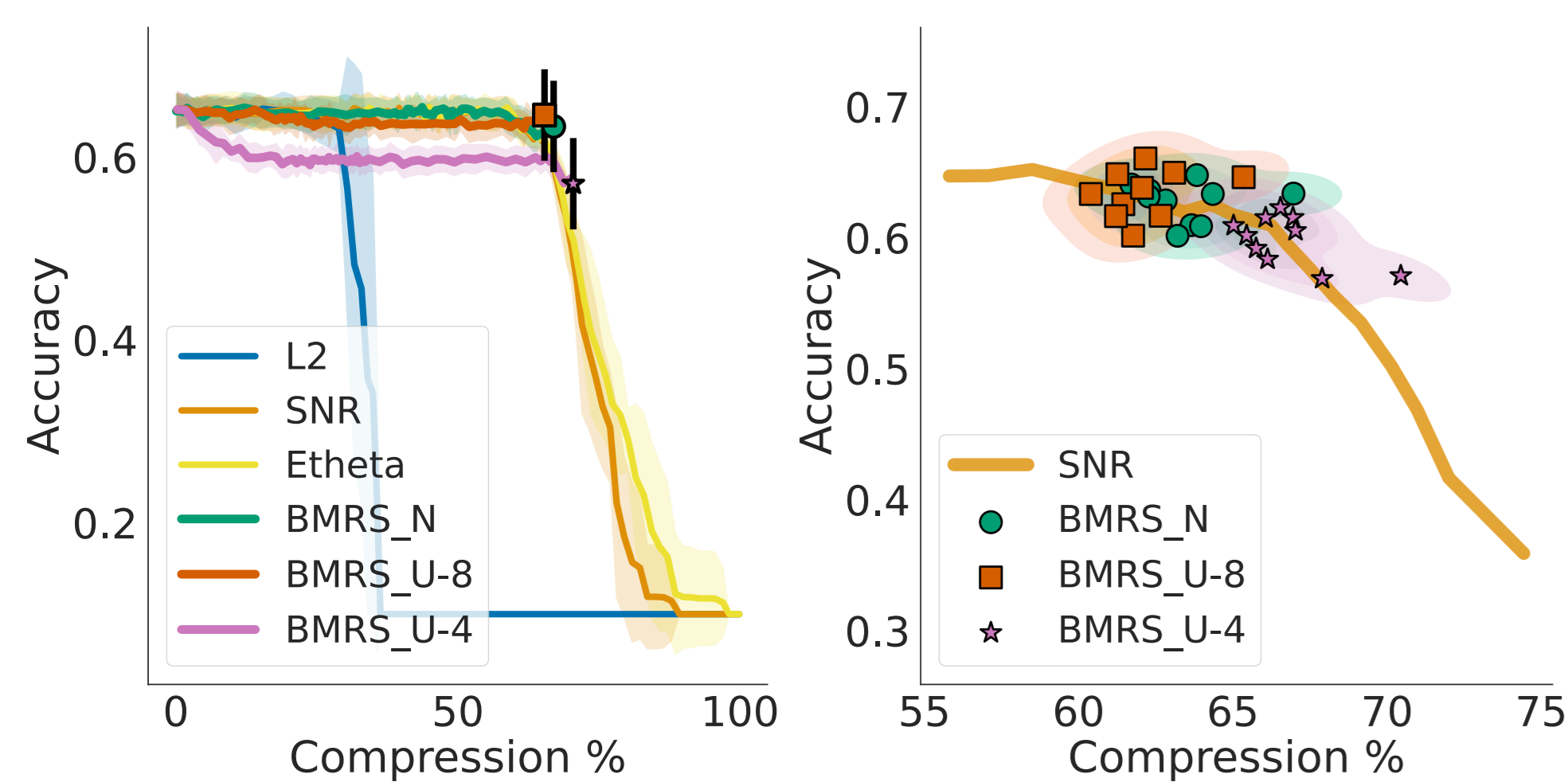
- Use Bayesian Model Reduction (BMR) to find what to prune

$$\Delta F \triangleq \log \frac{\tilde{p}(D)}{p(D)} = \log \mathbb{E}_{\tilde{p}} \left[ \frac{q_\phi(\theta)}{p(\theta)} \right]$$

## RESULTS

- BMRS finds near-optimal point w/ post-training pruning
- Good compression + maintains accuracy w/ continuous pruning

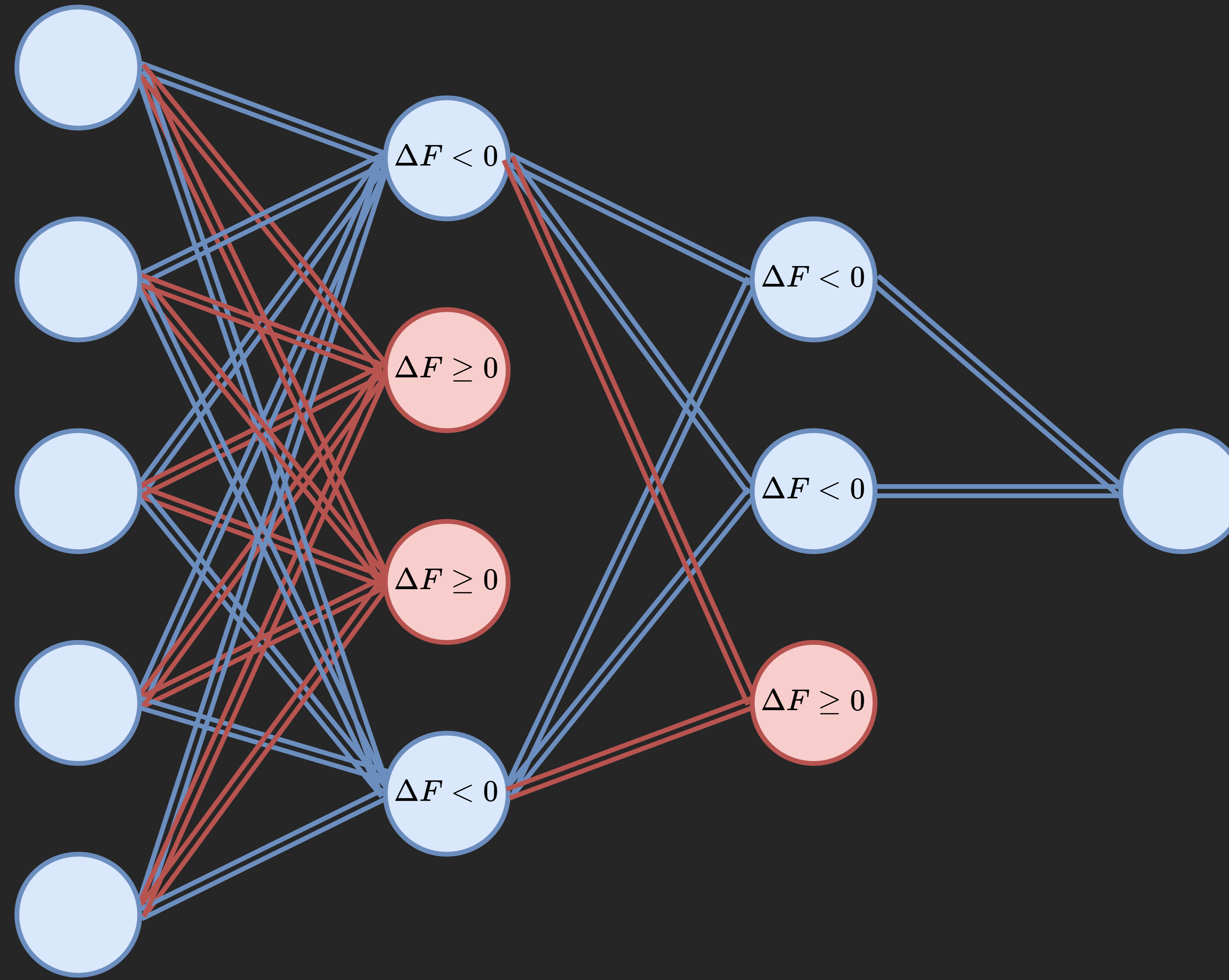
CNN performance on CIFAR10



Continuous pruning results

Pruning Method	MNIST		Fash-MNIST		CIFAR10	
	Comp. (%)	Acc.	Comp. (%)	Acc.	Comp. (%)	Acc.
MLP						
None	0.00 ± 0.00	97.43 ± 0.14	0.00 ± 0.00	88.17 ± 0.20	0.00 ± 0.00	44.94 ± 0.40
L2	43.11 ± 2.06	10.39 ± 0.32	87.86 ± 2.27	18.23 ± 10.22	42.89 ± 2.64	10.00 ± 0.00
E[θ]	52.08 ± 1.71	96.88 ± 0.15	91.76 ± 0.81	85.59 ± 0.26	77.99 ± 1.54	43.39 ± 0.46
SNR	58.57 ± 2.01	96.92 ± 0.08	99.83 ± 0.00	10.00 ± 0.00	75.93 ± 1.26	43.97 ± 0.46
BMRS <sub>N</sub>	48.86 ± 1.32	96.95 ± 0.19	93.20 ± 0.66	84.99 ± 0.35	76.36 ± 1.08	43.97 ± 0.29
BMRS <sub>U-8</sub>	48.73 ± 1.90	96.93 ± 0.16	93.02 ± 0.81	85.01 ± 0.32	77.17 ± 0.98	43.45 ± 0.42
BMRS <sub>U-4</sub>	54.47 ± 1.74	<b>96.99 ± 0.13</b>	91.57 ± 0.71	<b>85.79 ± 0.34</b>	76.63 ± 0.94	<b>44.06 ± 0.40</b>
Lenet5						
None	0.00 ± 0.00	99.07 ± 0.09	0.00 ± 0.00	89.16 ± 0.27	0.00 ± 0.00	67.62 ± 0.77
L2	83.42 ± 1.92	11.35 ± 0.00	83.62 ± 1.69	10.00 ± 0.00	52.29 ± 2.18	10.00 ± 0.00
E[θ]	88.29 ± 1.00	51.30 ± 41.12	89.71 ± 0.56	50.93 ± 33.45	66.19 ± 1.36	65.83 ± 0.90
SNR	92.66 ± 5.77	62.70 ± 41.93	98.47 ± 3.45	17.01 ± 21.03	70.29 ± 2.02	<b>67.68 ± 0.52</b>
BMRS <sub>N</sub>	86.90 ± 1.15	95.59 ± 0.94	88.02 ± 1.00	77.90 ± 2.44	62.87 ± 1.64	66.14 ± 0.70
BMRS <sub>U-8</sub>	86.11 ± 1.37	95.27 ± 1.02	87.61 ± 0.72	77.23 ± 3.49	62.54 ± 1.49	66.28 ± 1.07
BMRS <sub>U-4</sub>	87.58 ± 1.01	<b>96.06 ± 0.59</b>	88.72 ± 0.73	<b>81.10 ± 1.50</b>	68.07 ± 1.95	67.66 ± 0.59

# Bayesian structured pruning without threshold tuning

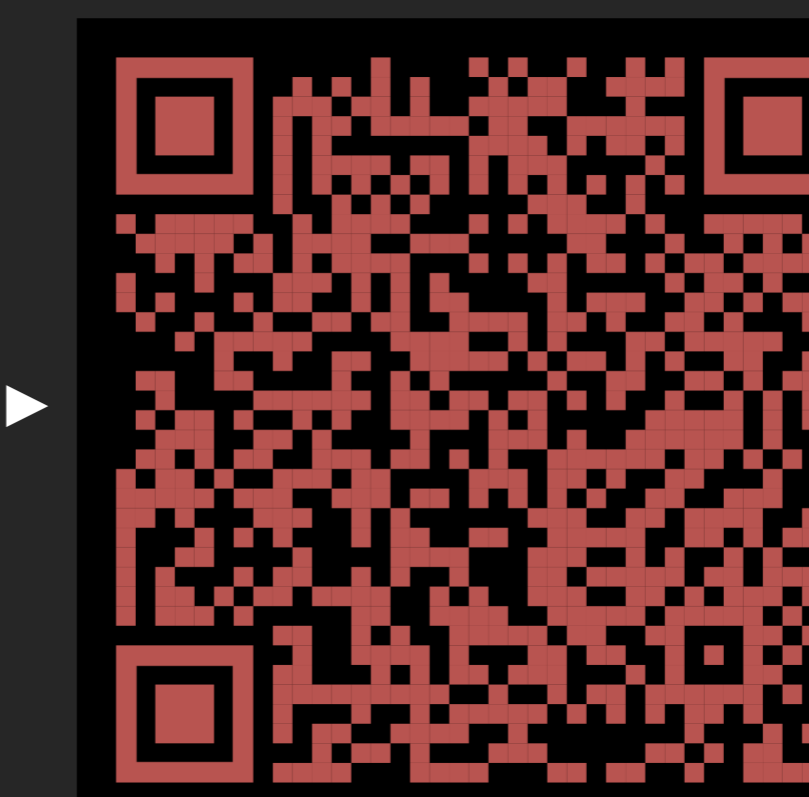


Structures are pruned using an efficient calculation based on the statistics of the variational distribution



Paper and code

Also read our critical perspective of AI efficiency



## Two variants

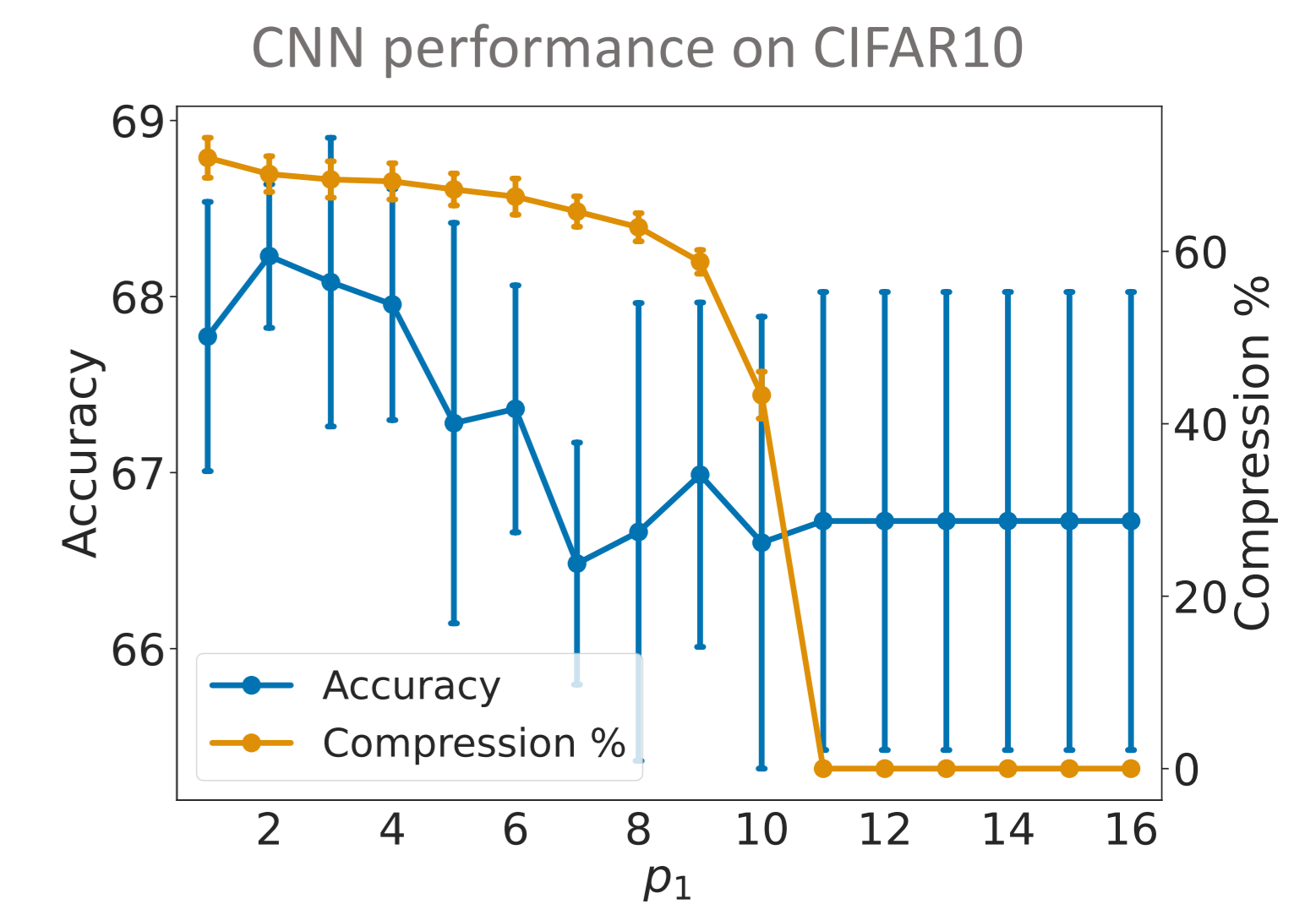
- BMRS<sub>N</sub> ΔF:

$$\frac{1}{2} \log \frac{\tilde{\sigma}_q^2}{2\pi\tilde{\sigma}_p^2\sigma_q^2} + \log \frac{Z_{\tilde{q}}(\log b - \log a)}{Z_{\tilde{p}}Z_q} - \frac{1}{2} \left( \frac{\mu_q^2}{\sigma_q^2} + \frac{\tilde{\mu}_p^2}{\tilde{\sigma}_p^2} - \frac{\tilde{\mu}_q^2}{\tilde{\sigma}_q^2} \right)$$

- BMRS<sub>U</sub> ΔF:

$$\frac{\log b' - \log a'}{\log b - \log a} \leq q_\phi(a' \leq \theta \leq b')$$

- BMRS<sub>N</sub> is **thresholdless**
- BMRS<sub>U</sub> can be **tuned...**



- ...and has a connection to **floating point precision**

Mantissas:  $m \sim (m \log B)^{-1}$ ,  $\frac{1}{B} \leq m \leq 1$

$$\tilde{p}(\theta): \begin{cases} \left( \theta \log \frac{2^{p_2}}{2^{p_1}} \right)^{-1}, & \frac{1}{2^{p_2}} \leq \theta \leq \frac{1}{2^{p_1}} \\ 0, & \text{otherwise} \end{cases}$$

- Different pruning functions are learned

Spearman rank correlation of pruning methods

SNR	1	1	0.86	0.95	0.47	0.27
Etheta	1	1	0.86	0.95	0.47	0.27
L2	0.86	0.86	1	0.81	0.33	0.21
BMRS <sub>N</sub>	0.95	0.95	0.81	1	0.49	0.32
BMRS <sub>U-8</sub>	0.47	0.47	0.33	0.49	1	0.8
BMRS <sub>U-4</sub>	0.27	0.27	0.21	0.32	0.8	1
SNR						
Etheta						
L2						
BMRS <sub>N</sub>						
BMRS <sub>U-8</sub>						
BMRS <sub>U-4</sub>						



University of Copenhagen