# Joint Inference for Neural Network Depth and Dropout Regularization

Kishan K C[1], Rui Li[1], Mahdi Gilany[2]

[1] Rochester Institute of Technology
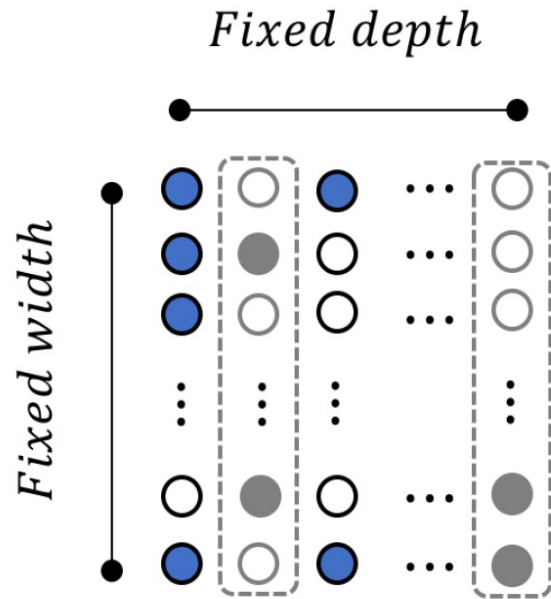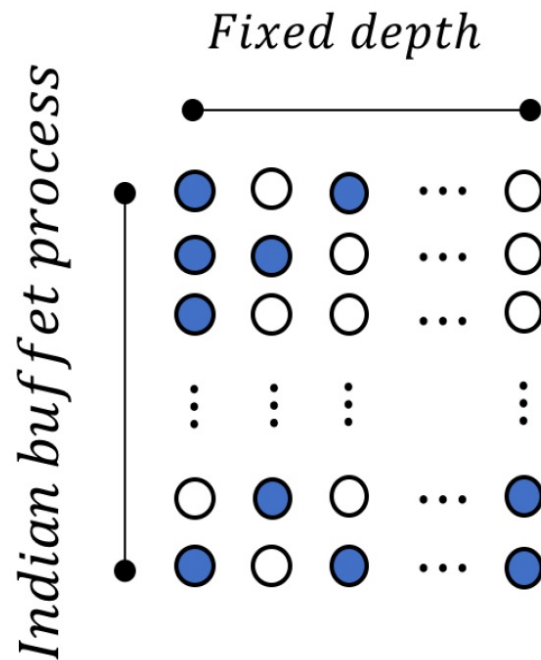
[2] Queens University

# Motivation

- Pre-determined backbone structures are the key

- Deep networks tend to be
  - Overfitted
  - Poorly calibrated with high confidence on incorrect predictions (Nguyen et al. 2015, Antorán et al. 2020)

- Current solutions
  - Dropout and its variants (Srivastava et al. 2014, Gal et al. 2017, Lee et al. 2019)
  - Structure selection methods (Srinivas et al. 2016, Dikov et al. 2019, Antorán et al. 2020)

- However,
  - Cannot scale the network beyond the pre-determined structure
  - Cannot achieve a balance between network depth and dropout regularization for uncertainty calibration
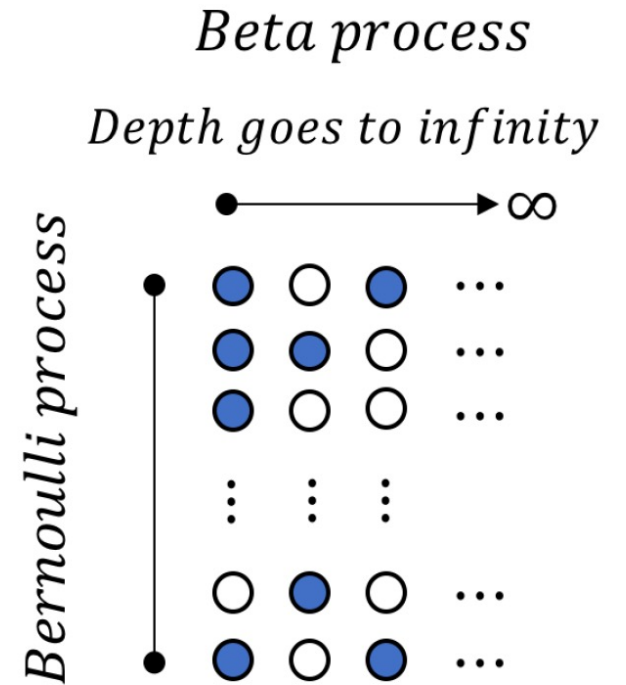
# Our proposed solution

- Model the depth (number of hidden layers) as a Beta Process
- Modulate neuron activations with a conjugate Bernoulli Process
- Joint inference of network depth and neuron activations



A typical structure selection method

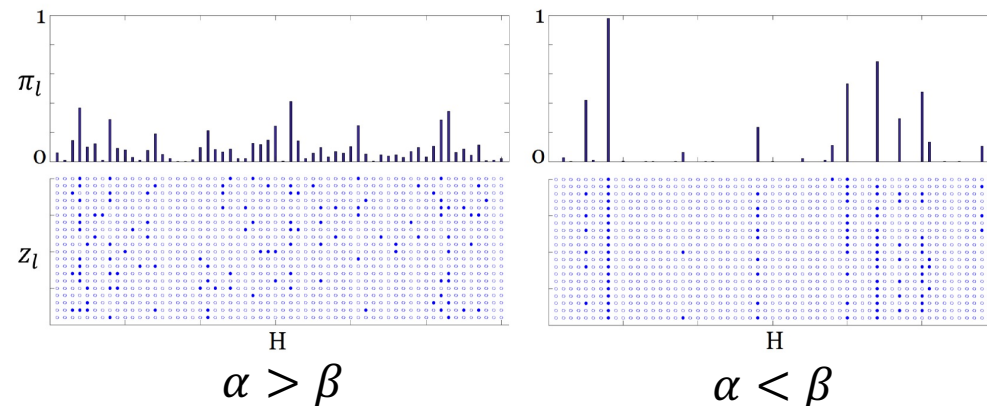A dropout variants

Our proposed solution

# Beta-Bernoulli process over network structures

- Model the depth of a neural network as a Beta Process

  - Stick breaking construction of beta-Bernoulli Process (Paisley et al. 2010, Broderick et al. 2012)

  $$v_l \sim \text{Beta}(\alpha, \beta), \qquad \pi_l = \prod_{j=1}^{l} v_j, \qquad z_{ml} \sim \text{Bernoulli}(\pi_l)$$
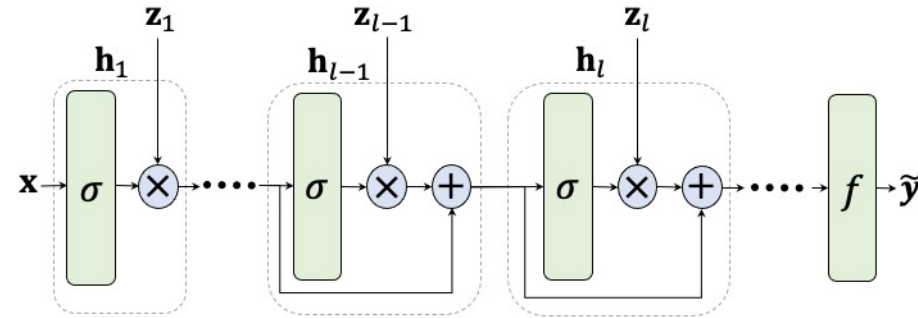
  - The prior over the network structures $\mathbf{Z}$:

  $$p(\mathbf{Z}, \boldsymbol{v} | \alpha, \beta) = p(\boldsymbol{v} | \alpha, \beta) p(\mathbf{Z} | \boldsymbol{v}) = \prod_{l=1}^{\infty} \text{Beta}(v_l | \alpha, \beta) \prod_{m=1}^{M} \text{Bernoulli}(z_{ml} | \pi_l)$$



$$\alpha > \beta \qquad\qquad\qquad\qquad \alpha < \beta$$

# Network structure with infinite layers

- A neural network has the form



$$h_l = \sigma(\mathbf{W}_l \mathbf{h}_{l-1}) \otimes \mathbf{z}_l + \mathbf{h}_{l-1} \qquad l \in \{1, 2, \dots, \infty\}$$

- A Gaussian likelihood of the neural network for regression task

$$p(D|\mathbf{Z}, \mathbf{W}) = \prod_{n=1}^{N} \mathcal{N}(y_n | f(\mathbf{x}_n; \mathbf{Z}, \mathbf{W}), s^2 \mathbf{I})$$

# Efficient inference
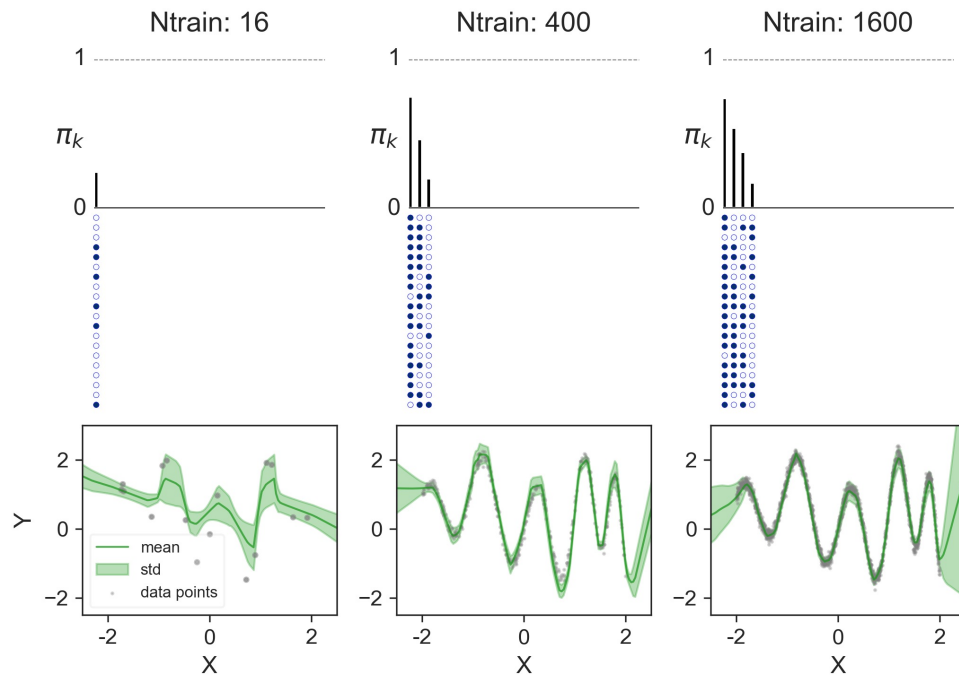
- The marginal Likelihood over network structures **Z** is

$$p(D|\mathbf{W}, L, \alpha, \beta) = \int p(D|\mathbf{Z}, \mathbf{W})\, p(\mathbf{Z}, \boldsymbol{v}|\alpha, \beta)\, d\mathbf{Z}d\boldsymbol{v}$$

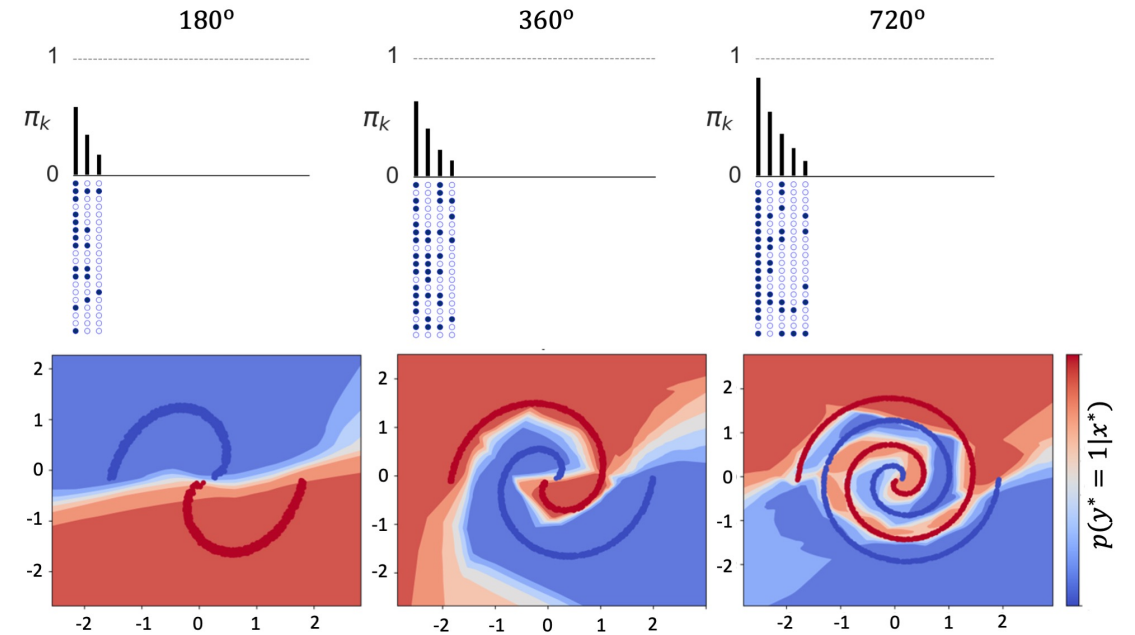- Approximation with structured stochastic variational inference (Hoffman et al. 2013, 2015)

$$\log p(D|\mathbf{W}, L, \alpha, \beta) \geq \mathbb{E}_{q(\mathbf{Z},\boldsymbol{v})}\left[\log p(D|\mathbf{Z}, \mathbf{W})\right] - \mathrm{KL}[q(\mathbf{Z}|\boldsymbol{v})||p(\mathbf{Z}|\boldsymbol{v})] - \mathrm{KL}[q(\boldsymbol{v})||p(\boldsymbol{v})]$$

  - We use truncation level $K$ for the variational distribution
  - Reparameterization of Beta and Bernoulli distribution (Jang et al. 2017, Maddison et al. 2017, Jankowiak et al. 2018)
  - We prove that optimizing ELBO is equivalent to Bayesian Information Criterion over the structure Z

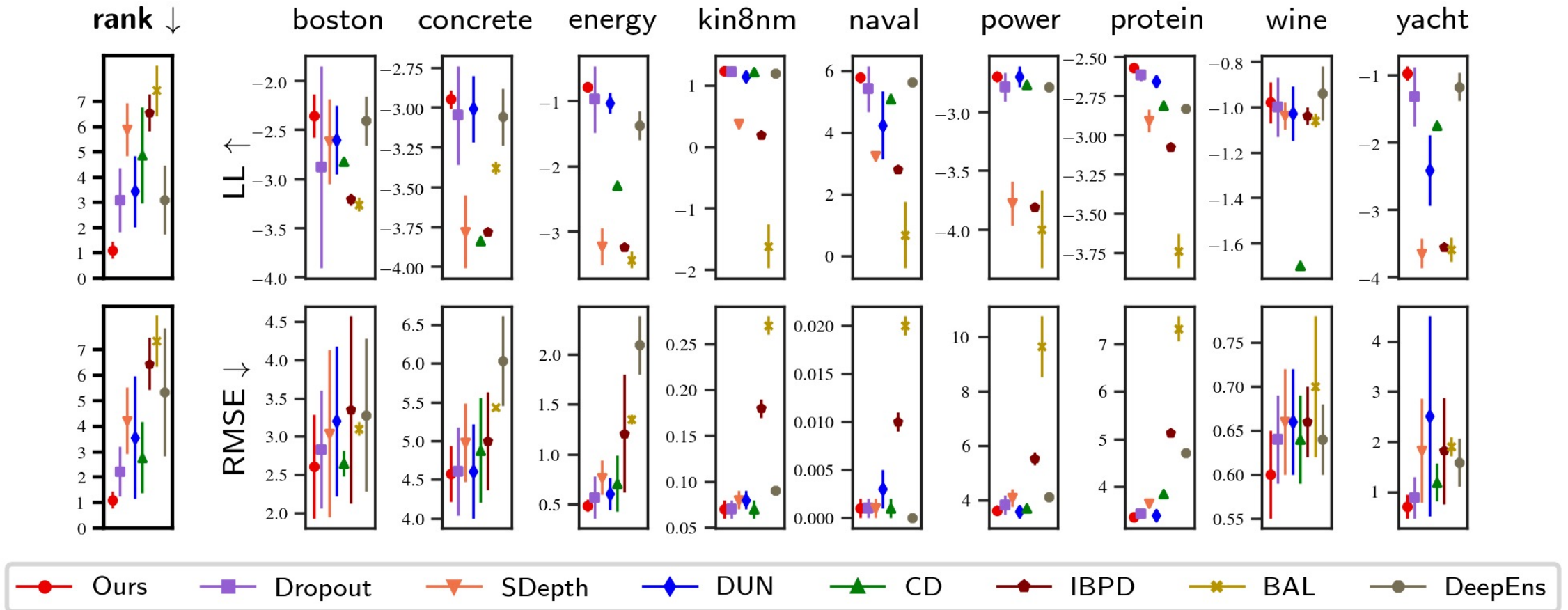# Performance evaluation on synthetic data



Periodic dataset
(Increasing data sizes)
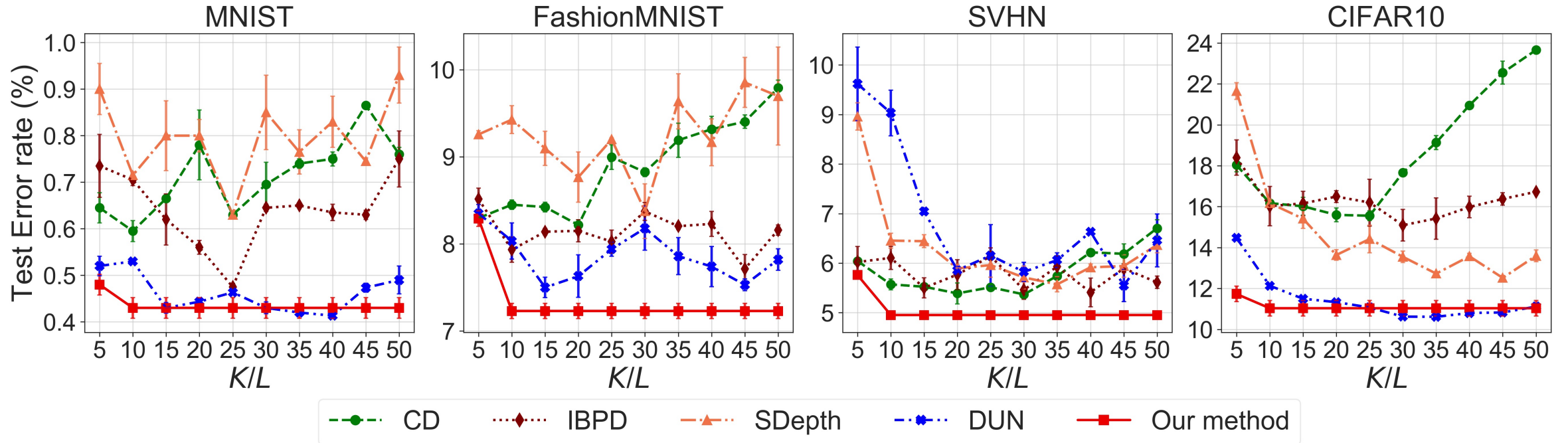
Spirals dataset
(Increasing complexity)

- If the training data size or its complexity increases, network structure grows to accommodate more information.
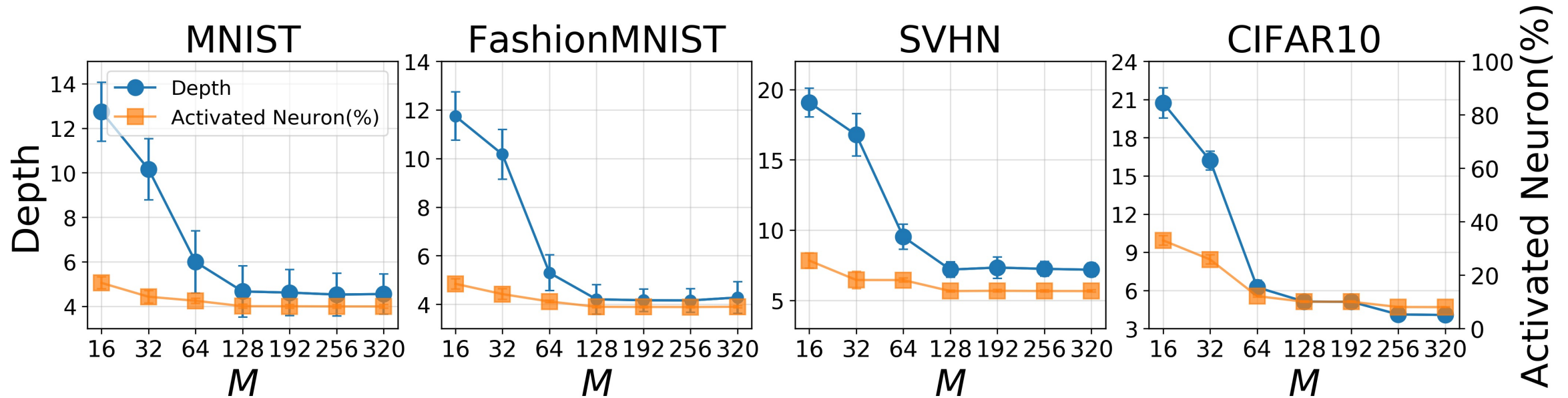
# Performance comparison on UCI datasets



- Our method achieves the overall highest rank for both uncertainty calibration and prediction accuracy.
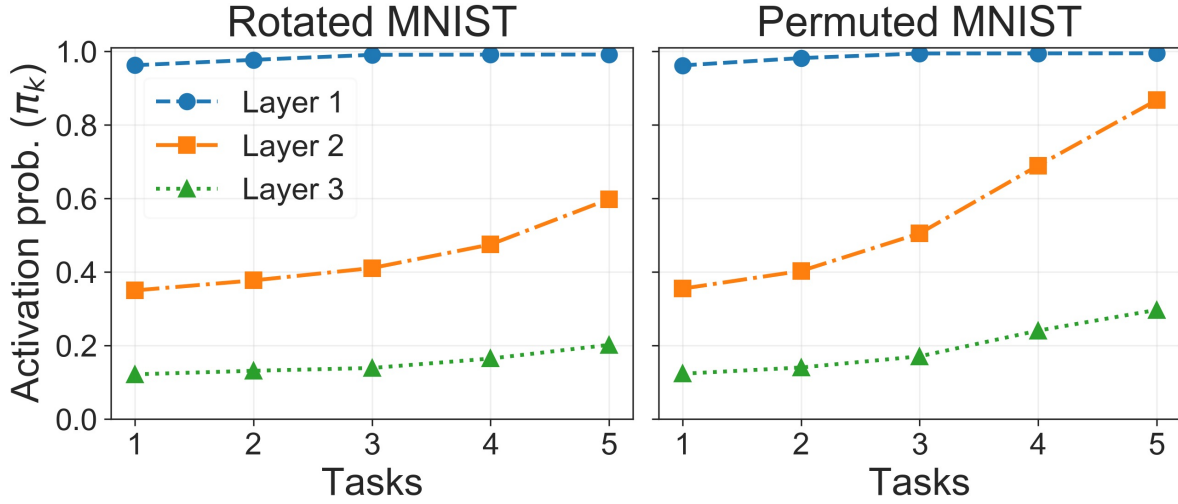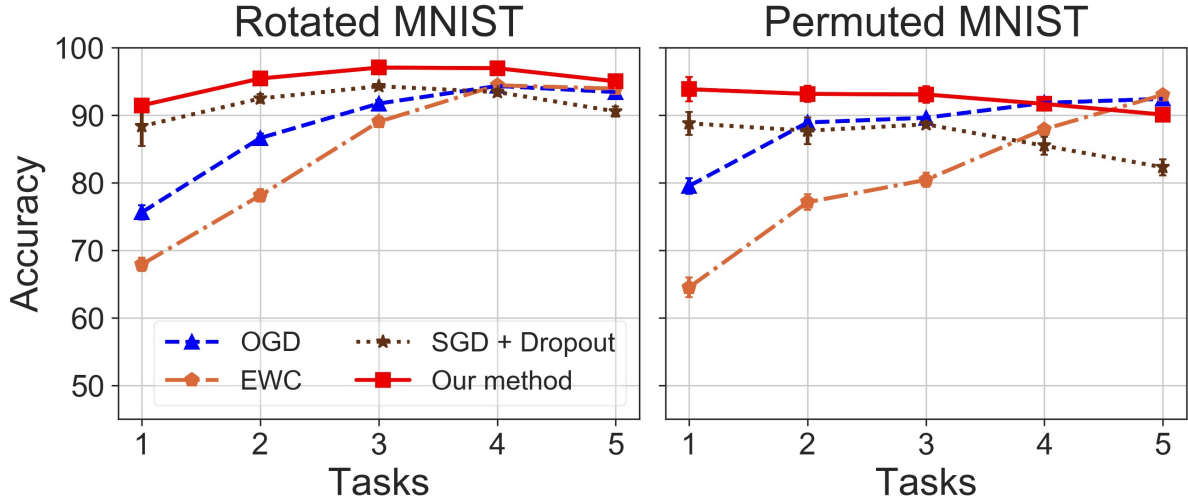
# Effect of truncation level $K$



- The truncation level ($K$) of our method does not affect the performance.
- The depth ($L$) of other methods significantly affects the performance and should be set carefully.

# Effect of maximum width $M$



- With smaller width $M$, our method results in deeper network structures to compensate for the relatively narrow layers.
- As $M$ increases, the structures become shallower.

# Case study on Continual learning



- Our method alleviates catastrophic forgetting by enabling network depth to dynamically augment to accommodate incrementally available information.

# Conclusion

- General joint inference framework applicable for various neural networks

- Experimental results on MLPs and CNNs show that our method achieves superior performance by adapting network depth and neuron activations

- Our model can accommodate incrementally available information by enabling neural network structures to dynamically evolve

## Thank you!