

Learning to Defer with Early-Exit Neural Networks

Mark Coates

Department of Electrical and Computer Engineering

McGill University



Joint work with Florence Regol (Ph.D. student) and Joud Chataoui (M.Sc. Student)



F. Regol, J. Chataoui, and M. Coates, “Jointly Learning to Exit and Infer in Dynamic Neural Networks: JEI-DNN,” to appear, ICLR 2024.



Traditional Architecture

Neural Network



Early-Exit Dynamic Network (EEDN)

Neural Network (backbone)



IMs
Inference
modules

classifier 1

classifier 2

\hat{y}_θ^1

\hat{y}_θ^2

GMs
Gate
mechanisms

gate 1

gate 2

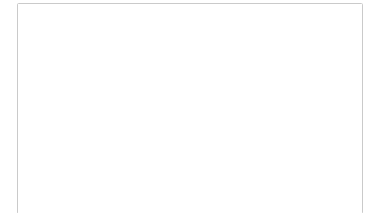
X

0

\hat{y}_θ^2

Previous Approaches

- **Threshold-based Gate Mechanisms**
 - Confidence score thresholds to decide whether to exit
 - BranchyNet¹, MSDNet², RANet³, CF-ViT⁴, Dynamic Perceiver⁵
- **Frozen Backbones and IMs + Learnable Gate Mechanisms**
 - EPNet⁶, PTEENet⁷, EENet⁸
- **Adaptive IM training**
 - BoostedNet⁹, L2W¹⁰



1: Teerapittayanon et al., ICPR 2016

2: Huang et al., ICLR 2018

3: Yang et al., CVPR 2020

4: Chen et al., AAAI 2023

5: Han et al., ICCV 2023

6: Dai et al., ICPR 2020

7: Lahiany et al., ICLR 2018

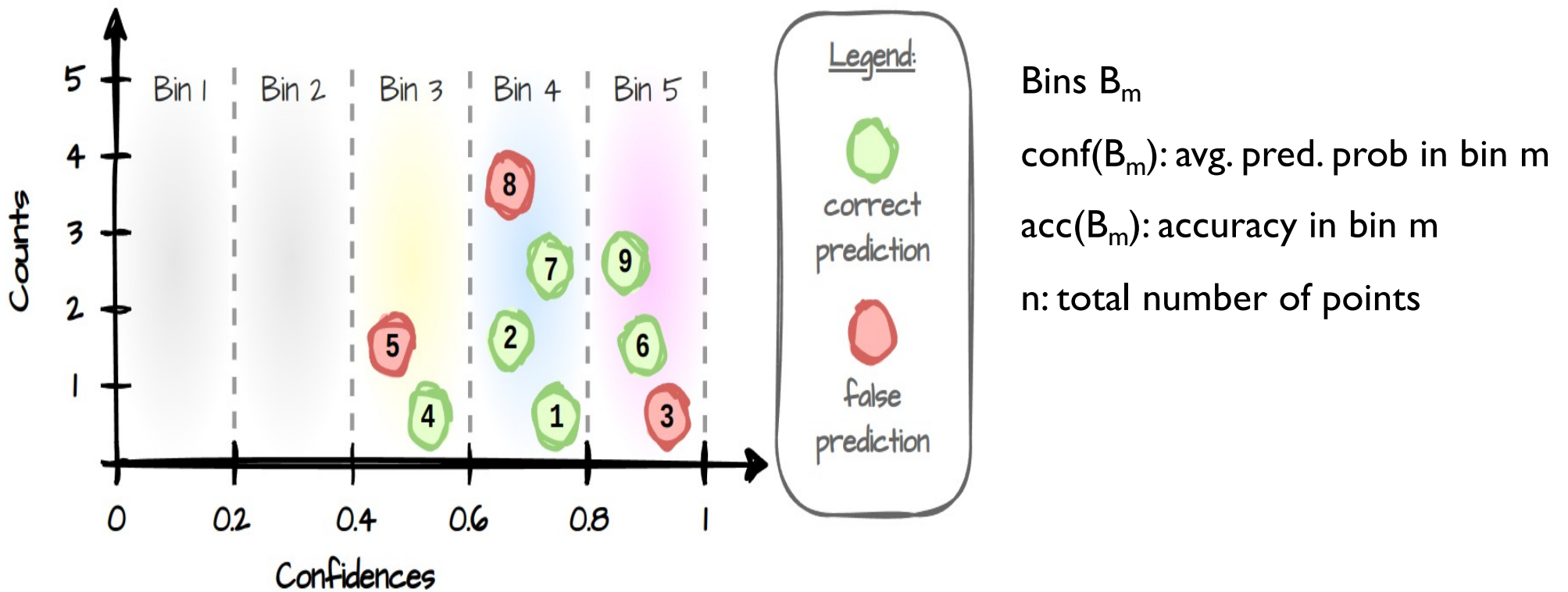
8: Ilhan et al., CVPR 2023

9: Yu et al., ICPR 2023

10: Han et al., ICLR 2022

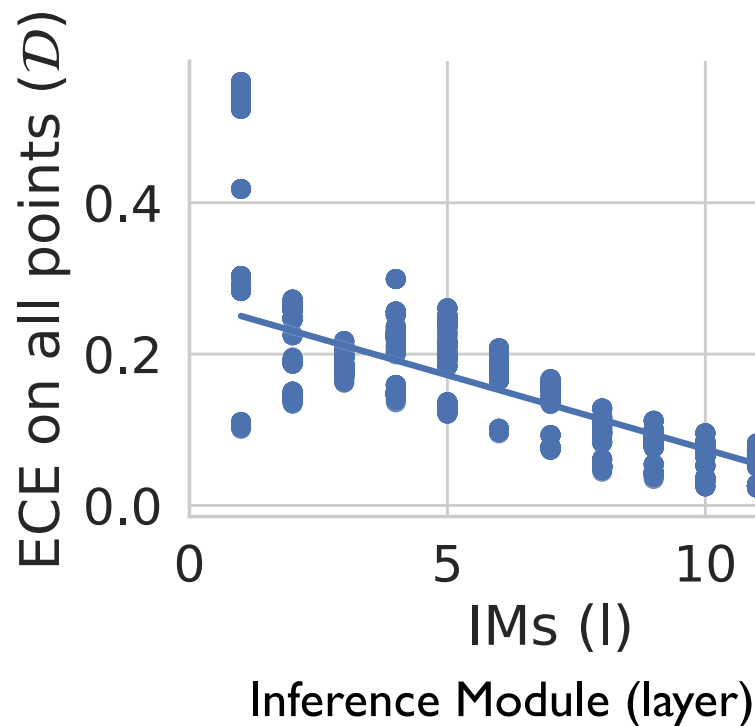
Expected Calibration Error:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{n} |acc(B_m) - conf(B_m)|$$



Maja Pavlovic, Expected Calibration Error (ECE): A Step-by-Step Visual Explanation, Towards Data Science, Jul. 2023.

$$\text{gate } l = \begin{cases} \textit{earlyexit} & \text{if } \max_k \hat{\mathbf{p}}_k^l > \eta \\ \textit{continue} & \text{o.w.} \end{cases}, \hat{\mathbf{p}}^l = \text{predicted prob. IM } l$$



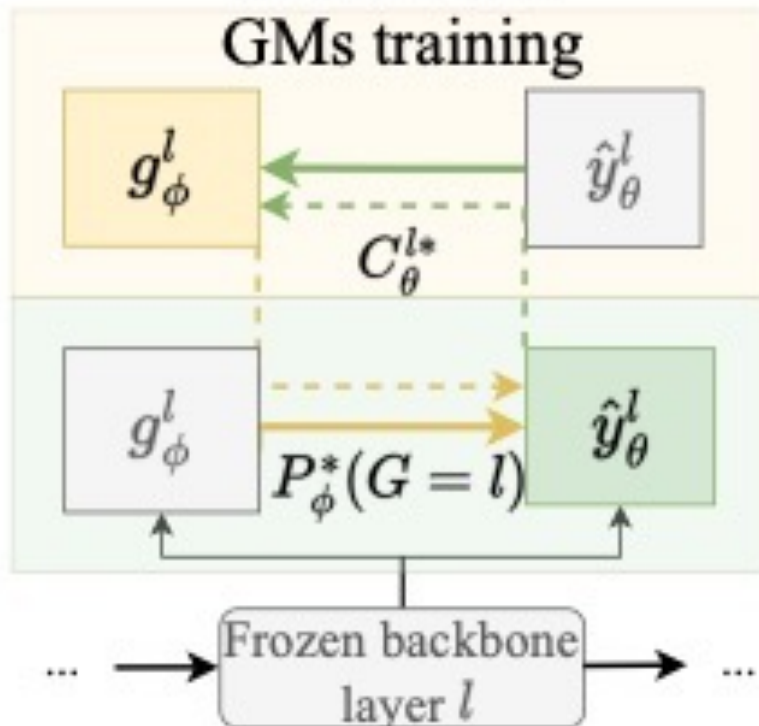
Problems:

1. IMs and GMs disconnected during training
➔ training-test mismatch
2. Uncertainty information is consumed by thresholding
➔ harder to calibrate
3. Relies on well-calibrated classifiers

CIFAR100
T2T-ViT-14 (Yuan et al., ICCV 2021)

Our Solution: JEI-DNN

Joint training via bi-level optimization



- Many lightweight gates and IMs
- Jointly learn the gate mechanisms and IMs
- Model the gate output $P_\phi(G = l | \mathbf{x}_i)$

C_θ^{*l} = Cost of choosing IM l ,
 $P_\phi^*(G = l)$ = Prob. of choosing IM l ,
 g_ϕ^l = Parameter of $P_\phi(G)$,
 \hat{y}_θ^l = Prediction of IM l .

Modelling the gate variables

- Only need to evaluate $P_\phi(G = l|\mathbf{x}_i)$ if dynamic evaluation reaches layer l
- Can use any intermediate values calculated by gates, IMs, and base architecture
- Denote aggregate information $\mathbf{c}_i^{\leq l}$
- Construct an **additive** model:

$$P_\phi(G = 1|\mathbf{x}_i) = g_\phi^1(\mathbf{c}_i^{\leq 1}),$$

$$P_\phi(G = l|\mathbf{x}_i) = \min(g_\phi^l(\mathbf{c}_i^{\leq l}), 1 - \sum_{j=1}^{l-1} g_\phi^j(\mathbf{c}_i^{\leq j})) \quad \text{for } l = 2, \dots, L.$$

Loss of JEI-DNN

- Combination of a cross entropy loss (accuracy of prediction) + inference cost (fixed cost per layer)

$$\mathcal{L} = \mathbb{E}_{Y,X} \mathbb{E}_{G|X} [\mathcal{L}^{CE}(Y, \hat{\mathbf{p}}_{\theta}^{G|X}(X)) + \lambda IC^{G|X}].$$

- Approximated loss

$$\mathcal{L} \approx \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^L \left(\mathcal{L}^{CE}(y_i, \hat{\mathbf{p}}_{\theta}^l(\mathbf{x}_i)) + \lambda IC_{\text{norm}}^l \right) \min \left(g_{\phi}^l(\mathbf{c}_i^{\leq l}), 1 - \sum_{j=1}^{l-1} g_{\phi}^j(\mathbf{c}_i^{\leq j}) \right).$$

Prob. Exiting at Gate l

Cost of Exiting at Gate l: Inference + Accuracy

Optimization

- Optimization is challenging because of the min operator
- Bi-level optimization
- Let $C_{\theta(\phi)}^l = \mathcal{L}^{CE}(y_i, \hat{\mathbf{P}}_{\theta}^l(\mathbf{x}_i)) + \lambda IC_{\text{norm}}^l$

$$\phi^* = \arg \min_{\phi} \mathcal{L}^{out} \triangleq \arg \min_{\phi} \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^L C_{\theta^*(\phi)}^l P_{\phi}(G = l | \mathbf{x}_i), \quad \text{Outer: gate parameters}$$

$$s.t. \quad \theta^*(\phi) = \arg \min_{\theta} \mathcal{L}^{in} \triangleq \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^L (\mathcal{L}^{CE}(y_i, \hat{\mathbf{P}}_{\theta}^l(\mathbf{x}_i)) + \lambda IC_{\text{norm}}^l) P_{\phi}(G = l | \mathbf{x}_i).$$

Inner: IM parameters

Optimizing the IMs

- The gate probabilities take $\mathbf{c}_i^{\leq l}$ as input.
- These can depend on θ . Make the dependence explicit by writing

$$P_\phi(G = l|\mathbf{x}_i) = G_\phi(m(\theta, \mathbf{x}_i)).$$

- Then:

$$\frac{\partial \mathcal{L}^{in}}{\partial \theta} = \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^L \frac{\partial \mathcal{L}^{CE}(y_i, \hat{\mathbf{p}}_\theta^l(\mathbf{x}_i))}{\partial \theta} P_\phi(G = l|\mathbf{x}_i) + \frac{\partial G_\phi(m(\theta, \mathbf{x}_i))}{\partial \theta} C_{\theta(\phi)}^l,$$

Optimizing the IMs

$$\frac{\partial \mathcal{L}^{in}}{\partial \theta} = \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^L \frac{\partial \mathcal{L}^{CE}(y_i, \hat{\mathbf{p}}_{\theta}^l(\mathbf{x}_i))}{\partial \theta} P_{\phi}(G = l | \mathbf{x}_i) + \frac{\partial G_{\phi}(m(\theta, \mathbf{x}_i))}{\partial \theta} C_{\theta(\phi)}^l,$$

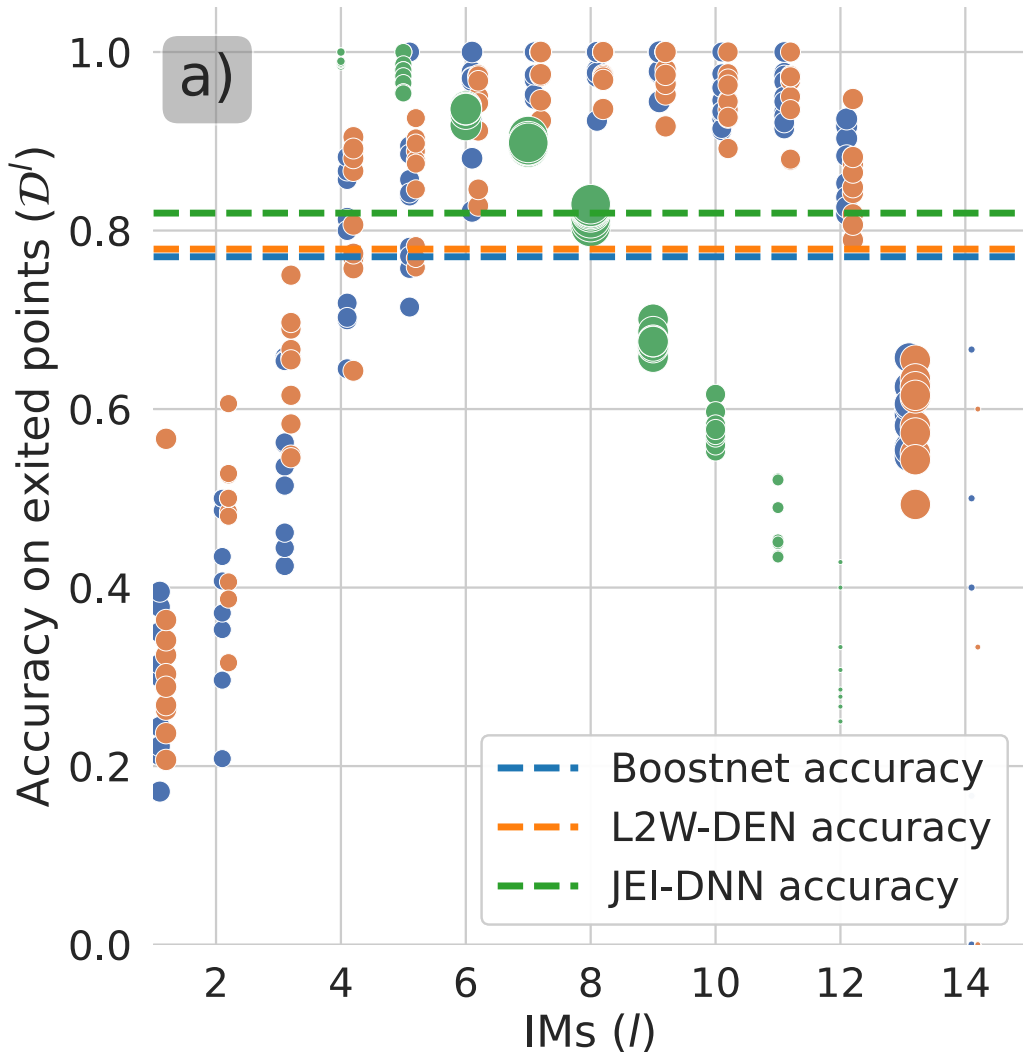
- First term is the gradient corresponding to a weighted cross-entropy loss
- Weights emerge directly from the gating mechanism
- Second term is driven by impact of θ on the gates
- **Practical approximation:** ignore the second term (represents a secondary effect)

Optimizing the Gates

Hard!

$$\frac{\partial \mathcal{L}^{out}}{\partial \phi} = \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^L \frac{\partial P_{\phi}(G=l|\mathbf{x}_i)}{\partial \phi} C_{\theta^*(\phi)}^l = \frac{1}{N} \sum_{i=1}^N \sum_{l=1}^L \frac{\partial \min \left(g_{\phi}^l(\mathbf{c}_i^{\leq l}), 1 - \sum_{j=1}^{l-1} g_{\phi}^j(\mathbf{c}_i^{\leq j}) \right)}{\partial \phi} C_{\theta^*(\phi)}^l.$$

- Alternative approach: define a surrogate binary classification problem
- Construct binary targets for $g_{\phi}^1(\mathbf{c}_i^{\leq 1}), \dots, g_{\phi}^{L-1}(\mathbf{c}_i^{\leq L-1})$
- Evaluate the cost of each gate $C_{\theta^*(\phi)}^l$ and determine lowest cost $l^* = \arg \min_l C_{\theta^*(\phi)}^l$
- Targets for binary tasks t_i^1, \dots, t_i^L are $t_i^j = 0$ for $j < l^*$ and $t_i^j = 1$ for $j \geq l^*$
- The surrogate tasks and initial objective share same solution if there exists ϕ such that the gating mechanisms can always select lowest-cost gates.

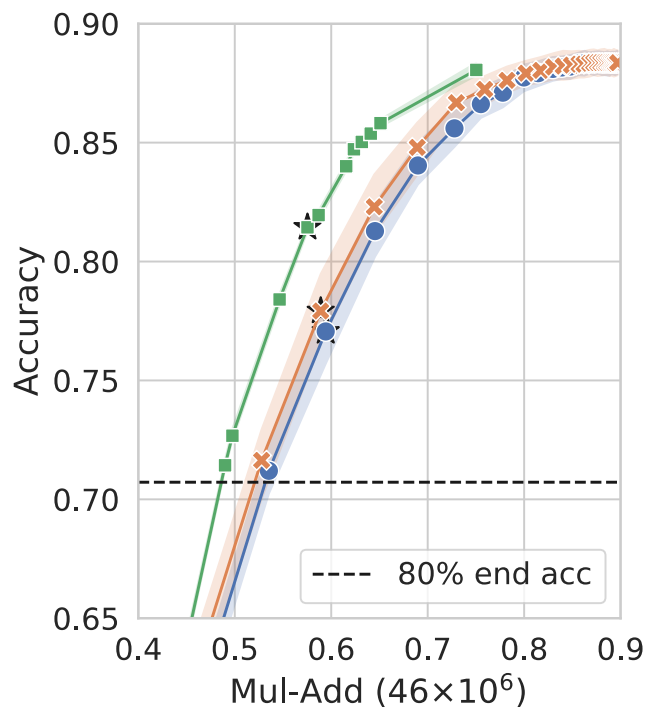


Exit Behaviour

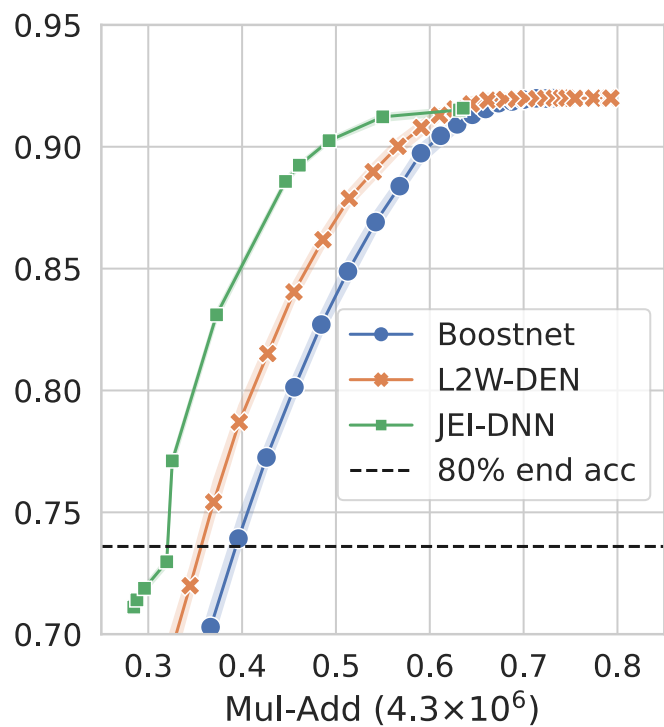
- JEI-DNN versus state-of-the-art baselines BoostedNet and L2W-DEN
- Higher accuracy on earlier exits
- Lower accuracy on later exits
- Higher average accuracy
- Concentrates on relatively few IMs
- Only starts to exit at layers 4-5.

Experiments

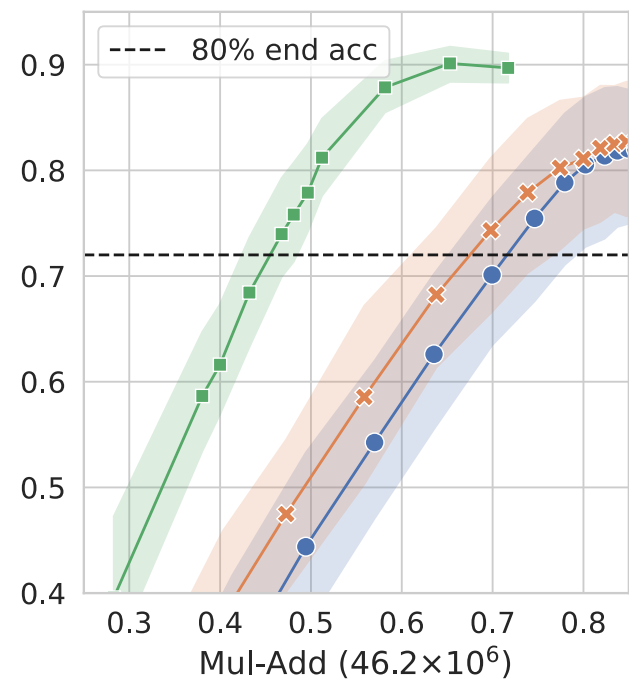
- Image classification using vision transformer: T2T-ViT-14 (Yuan et al., ICCV 2021)
- Single layer neural network gates and IMs at all layers
- Gates operate on 4 statistics based on IM output
 - (i) Max. pred. prob, entropy of predictions, entropy of scaled predictions, difference between two most confident predictions
- CIFAR10, CIFAR100: 60,000 32x32 color images with 10/100 classes
- Cropped digit SVHN: 99,289 32x32 color images of house numbers
- CIFAR100-LT: imbalanced classes (100x most common vs least common)
- ImageNet: 1.2 million 244 x 244 images; 1000 classes.



- CIFAR100 with T2T-ViT-14

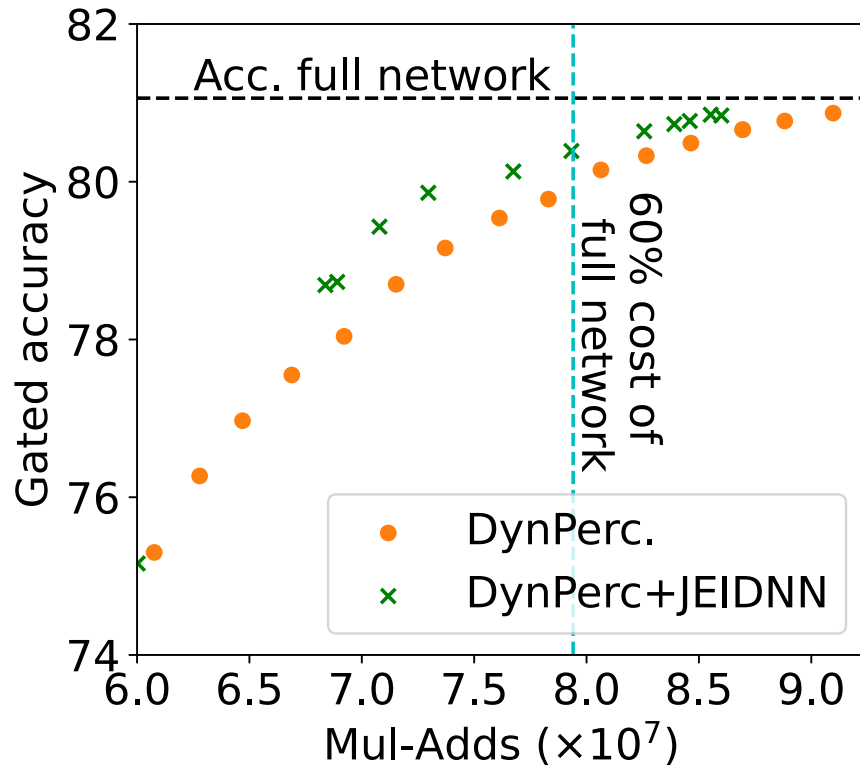


- SVHN with T2T-ViT-7

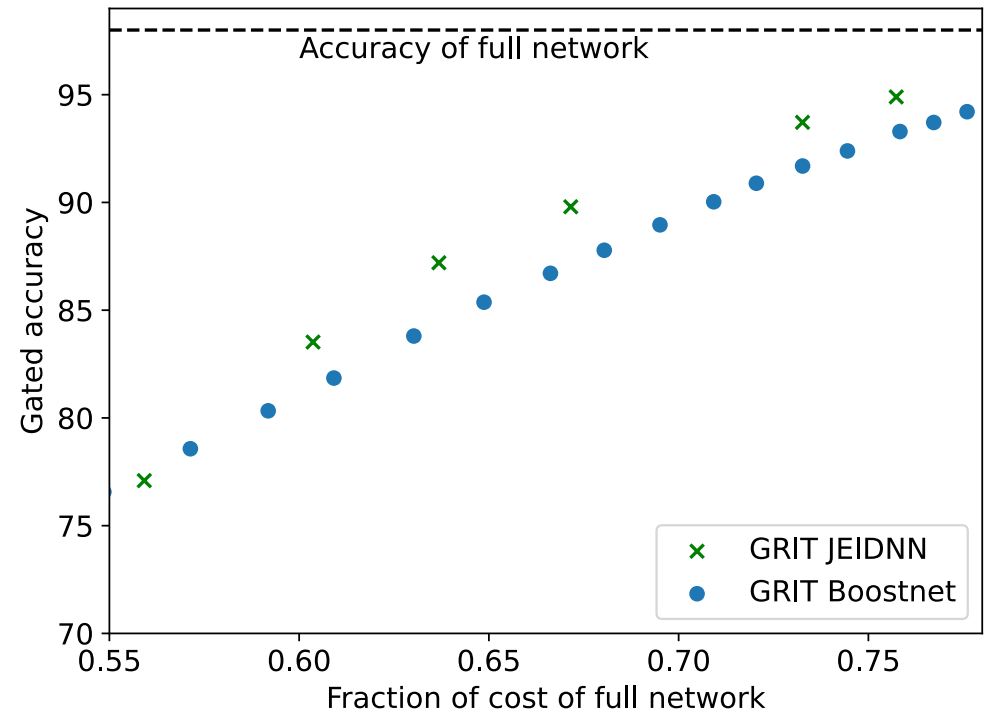


- CIFAR100-LT with T2T-ViT-14

Other Architectures



- ImageNet with Dynamic Perceiver



- Super-resolution MNIST with Graph Transformer