

Learning to Weight Samples for Dynamic Early-exiting Networks

**Yizeng Han^{*1}, Yifan Pu^{*1}, Zihang Lai^{2†}, Chaofei Wang¹, Shiji Song¹, Junfeng Cao³,
Wenhui Huang³, Chao Deng³, and Gao Huang^{1✉}**

^{*} Equal contribution. [✉]Corresponding author.

¹Department of Automation, Tsinghua University.

²Carnegie Mellon University. [†]Work done during an internship at Tsinghua University.

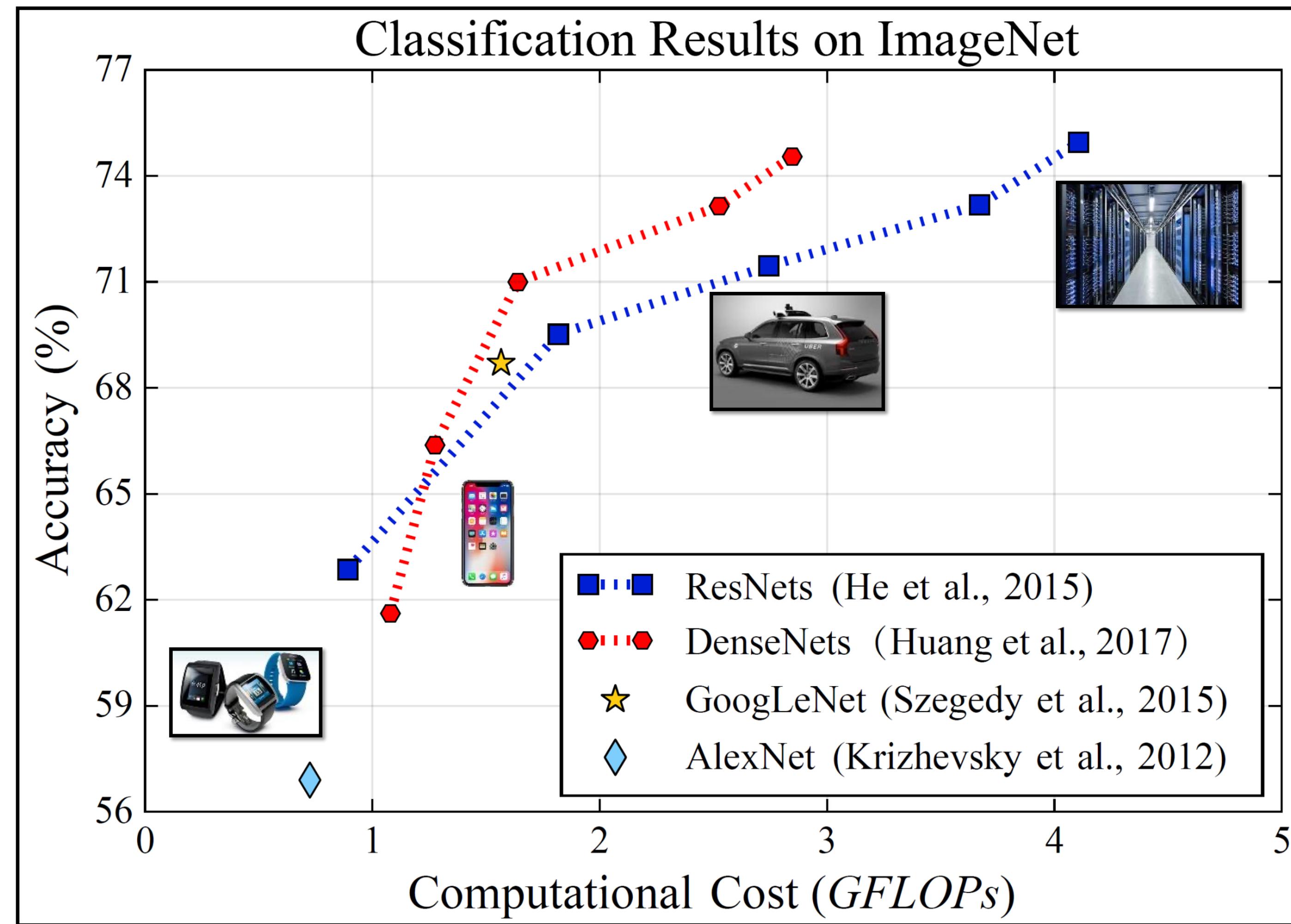
³China Mobile Research Institute.



清华大学
Tsinghua University



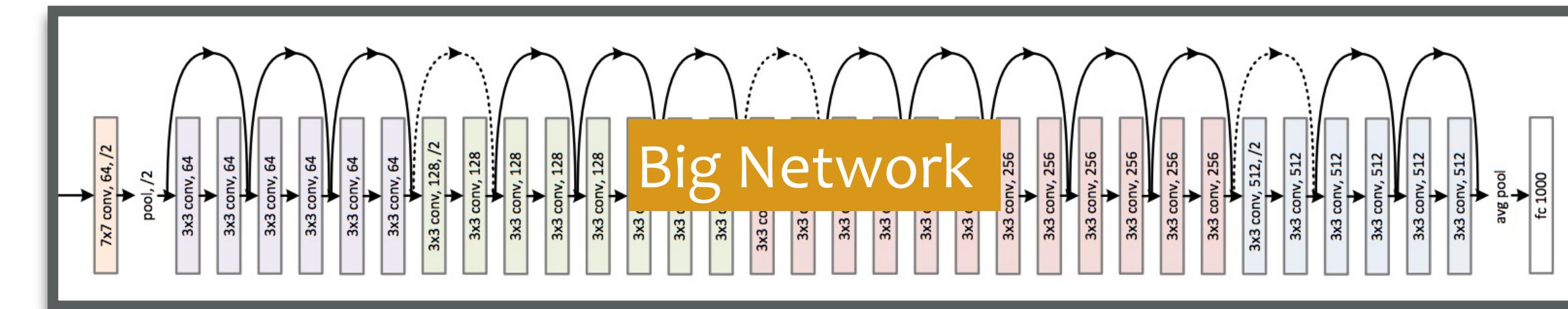
Accuracy-Efficiency Tradeoff



Most conventional neural networks recognize all samples with the same architecture.



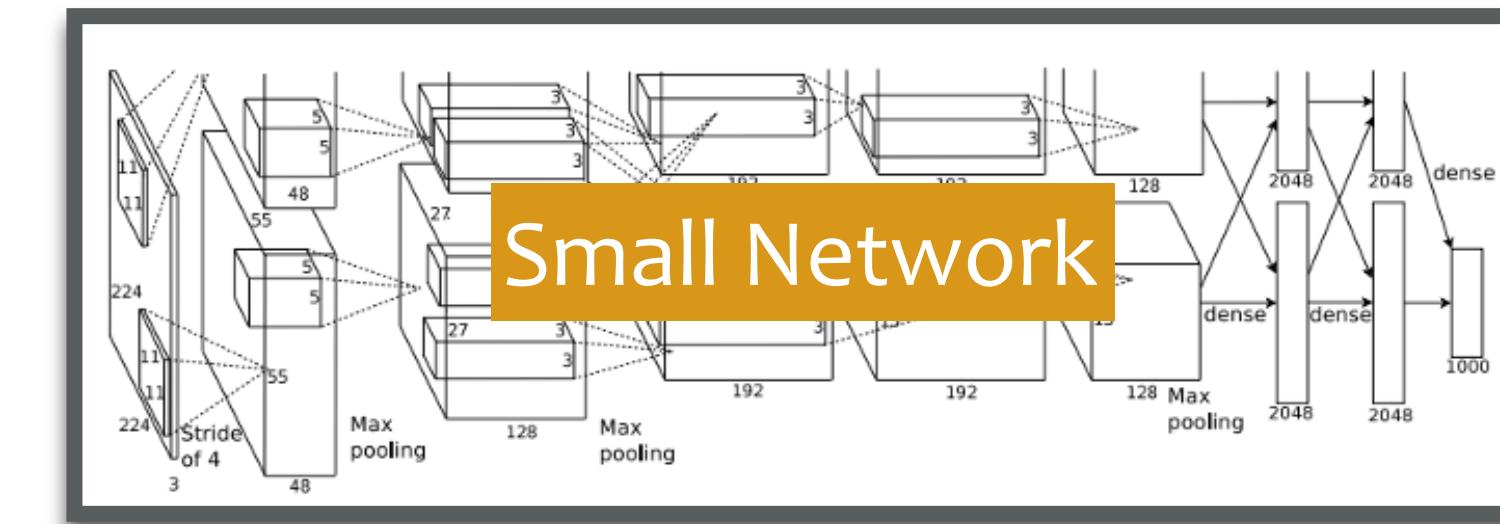
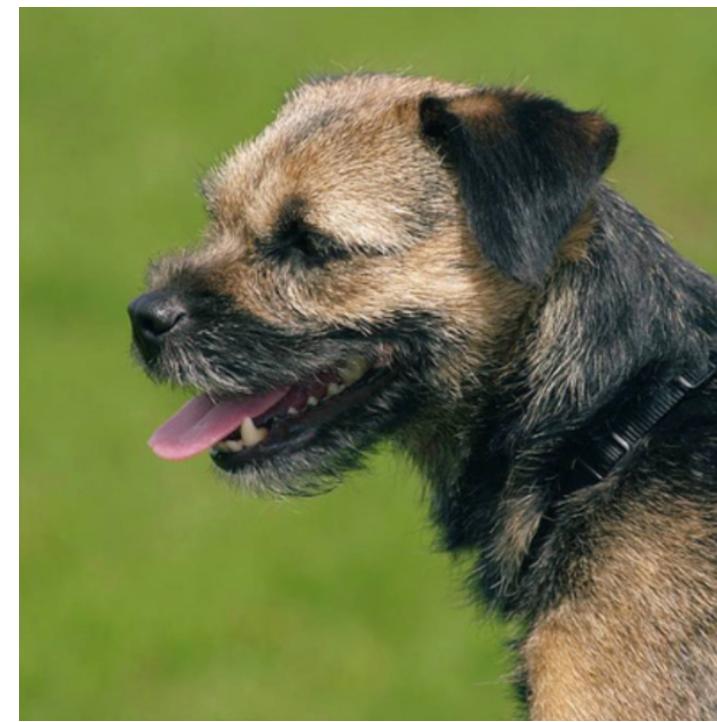
Canonical (“easy”)



Noncanonical (“hard”)

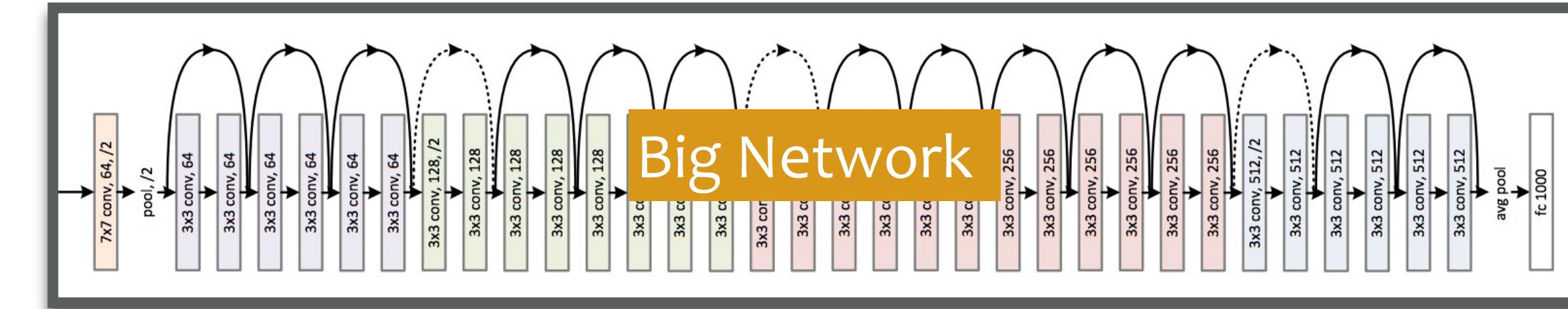
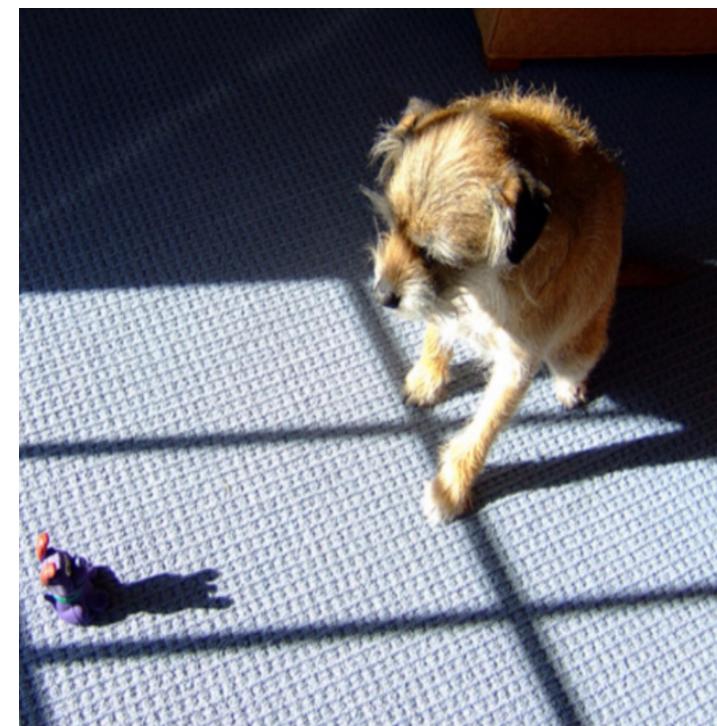


*A naïve idea of **dynamic** inference.*



Small Network

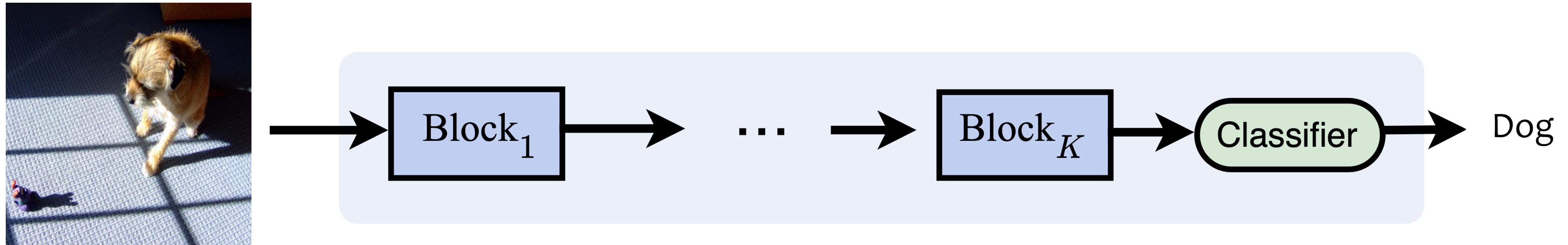
Canonical (“easy”)



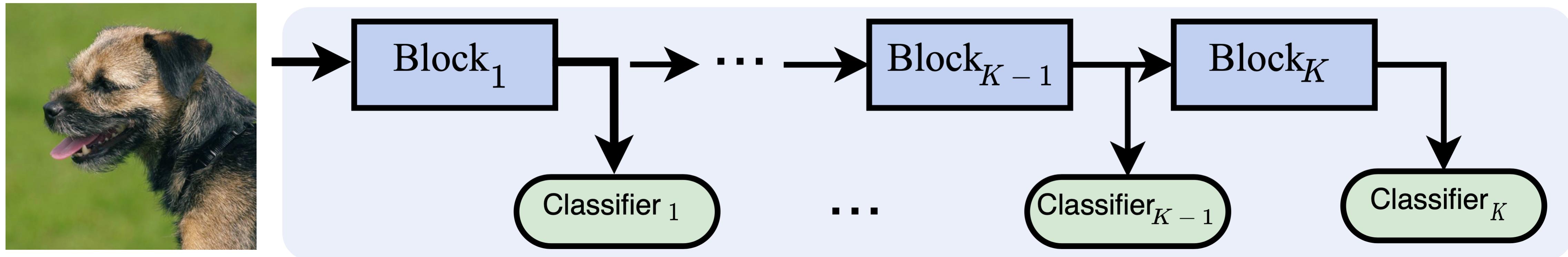
Big Network

Noncanonical (“hard”)

Conventional inference procedure in a deep network



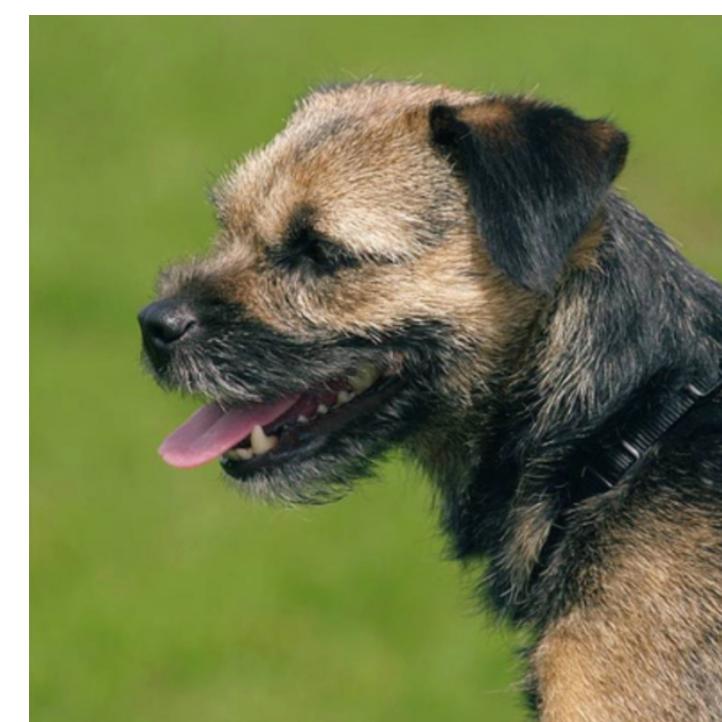
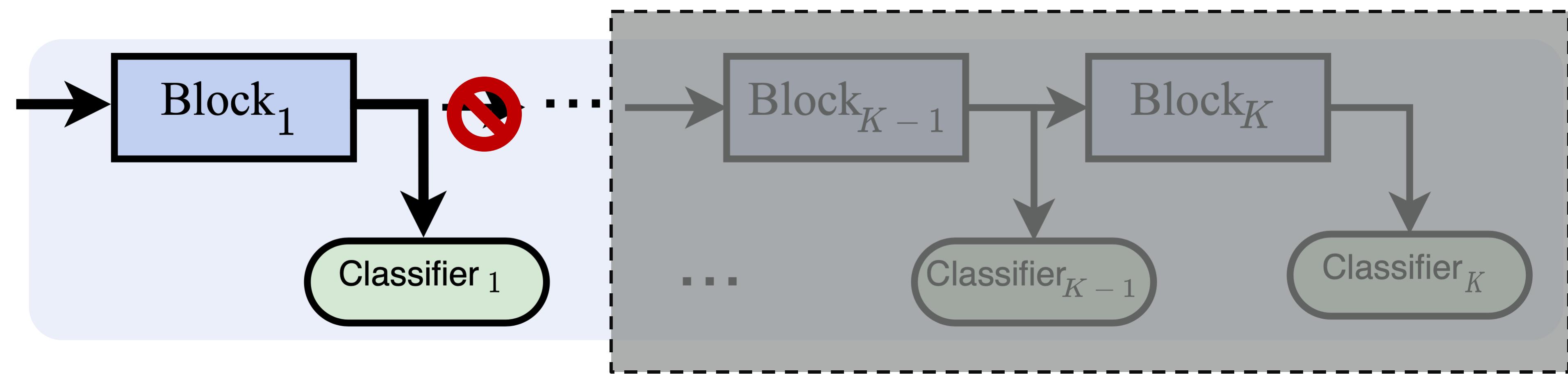
Dynamic early exiting^[1] in a multi-exit network.



Prediction: Dog
Confidence = $0.9 > \varepsilon_1$



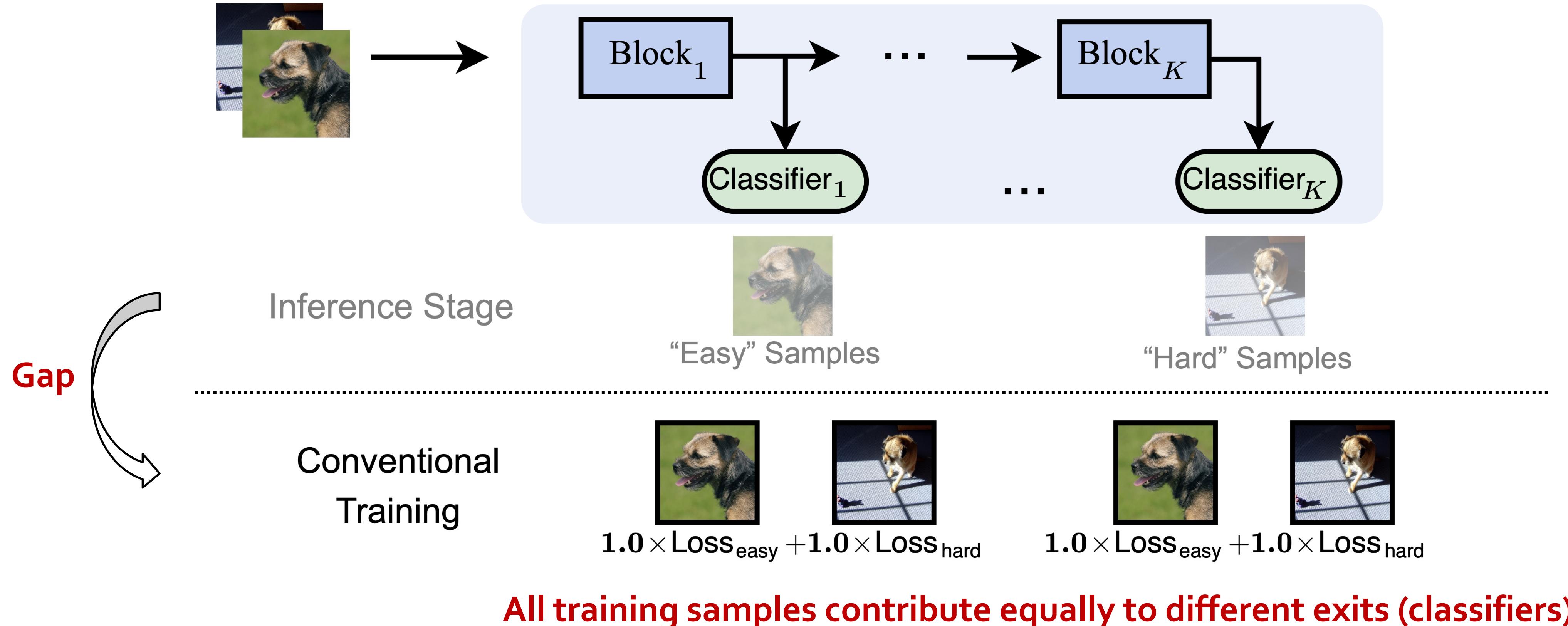
Dynamic early exiting^[1] in a multi-exit network.



Prediction: Dog
Confidence = $0.9 > \varepsilon_1$



Conventional training strategies^[2,3,4] are sub-optimal to dynamic early exiting.



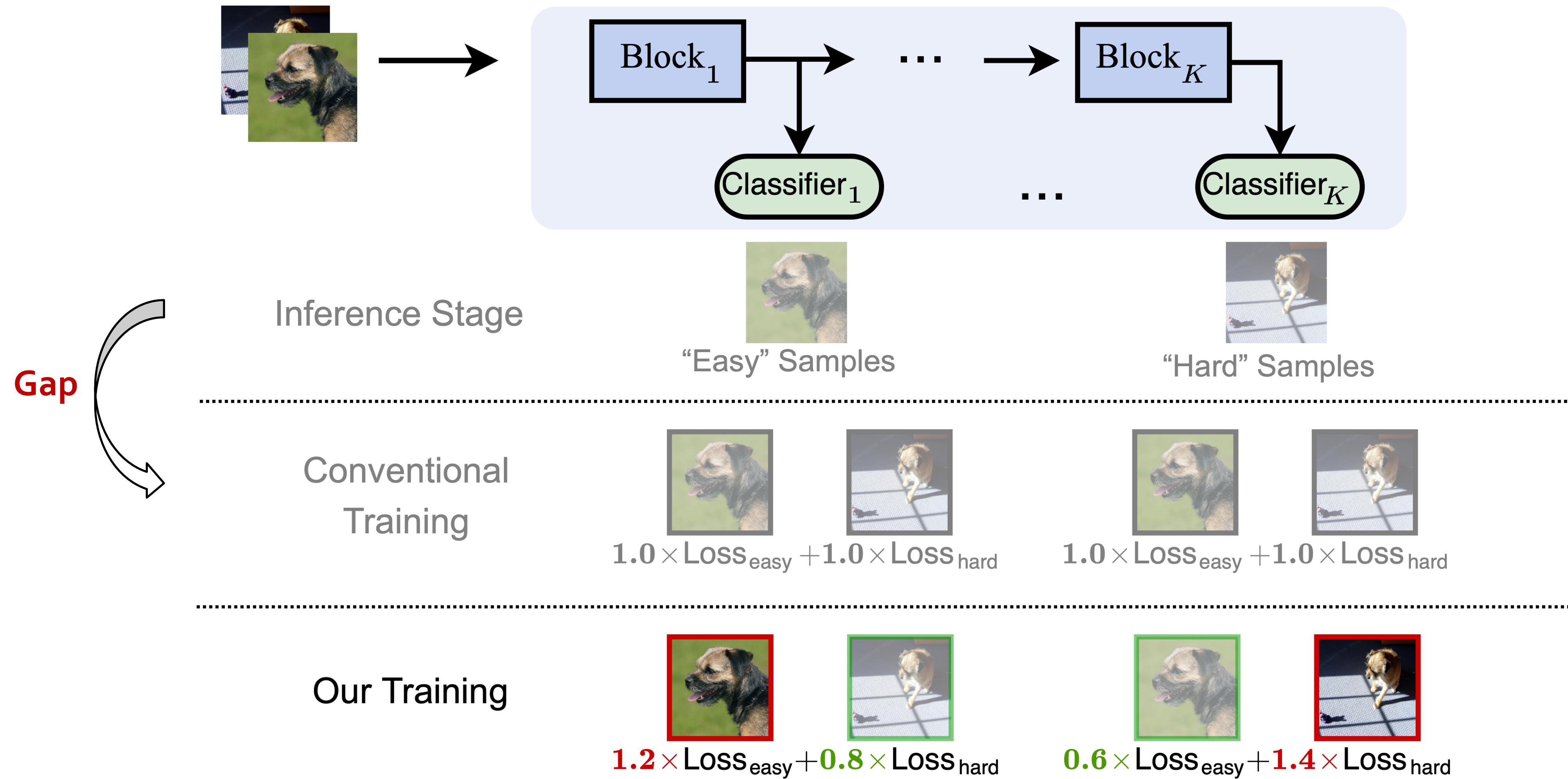
- [2] Huang, G., Chen, D., Li, T., Wu, F., van der Maaten, L., & Weinberger, K. (2018, February). Multi-Scale Dense Networks for Resource Efficient Image Classification. In ICLR.
- [3] Li, H., Zhang, H., Qi, X., Yang, R., & Huang, G. (2019). Improved techniques for training adaptive deep networks. In ICCV.
- [4] Yang, L., Han, Y., Chen, X., Song, S., Dai, J., & Huang, G. (2020). Resolution adaptive networks for efficient inference. In CVPR.



Our motivation:

*Training samples with **varying complexity** should **contribute differently** to the multiple exits in dynamic early-exiting networks.*

We bridge the gap between training and testing by *sample weighting*



Sample weighting formulation

- f : multi-exit dynamic network.
- $f^{(k)}(\cdot; \Theta_f^{(k)})$: the k -th sub-network parameterized by $\Theta_f^{(k)}$.
- Traditional training objective:

$$\mathcal{L} = \sum_{k=1}^K \frac{1}{N} \sum_{i=1}^N l_i^{(k)}$$

- Our modified objective with sample weighting:

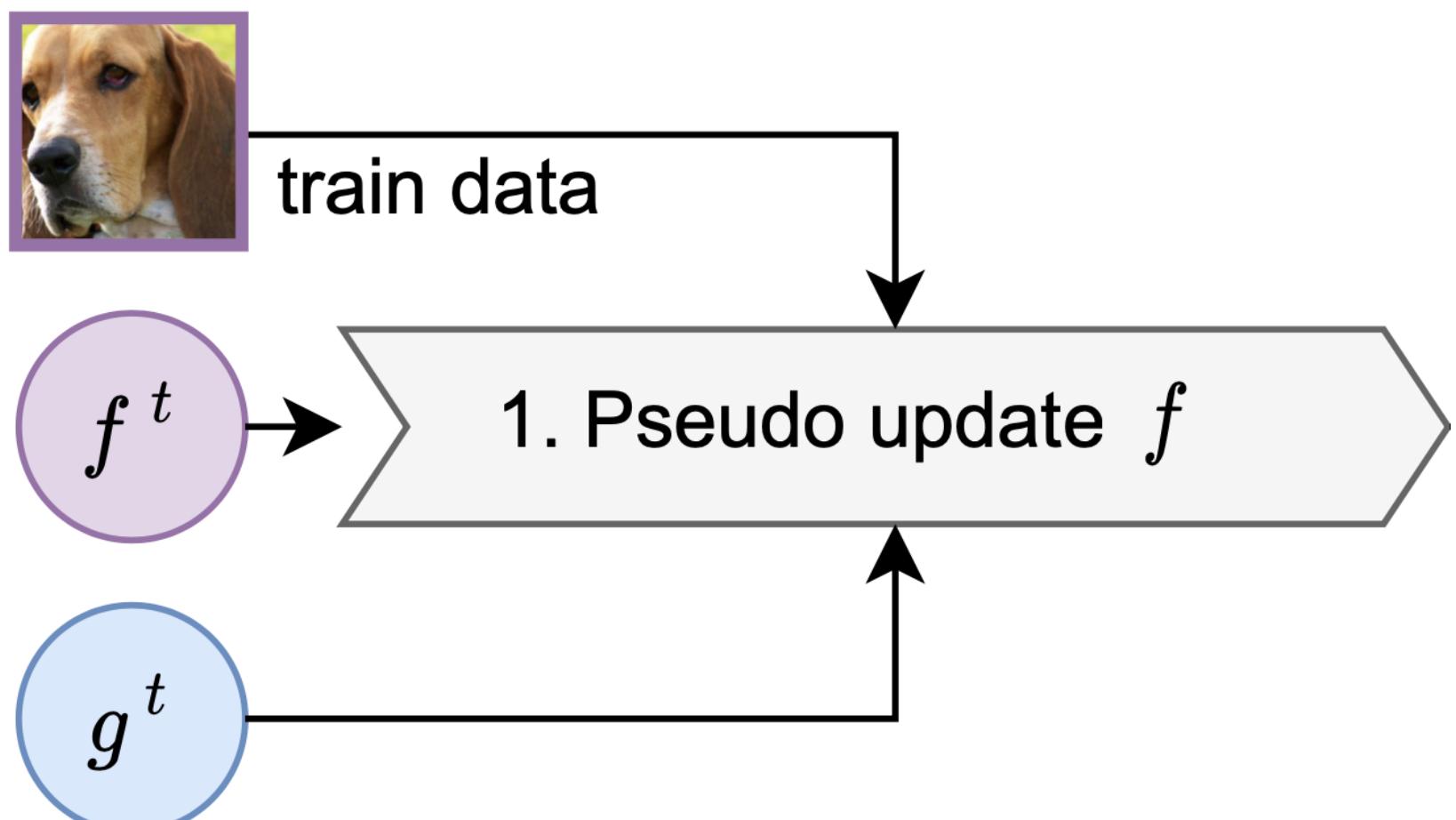
$$\mathcal{L} = \sum_{k=1}^K \frac{1}{N} \sum_{i=1}^N w_i^{(k)} l_i^{(k)}$$

- w_i is obtained from weight prediction network g : $w_i = g(l_i; \Theta_g)$,
where $w_i = [w_i^{(1)}, w_i^{(2)}, \dots, w_i^{(K)}]$,
 $l_i = [l_i^{(1)}, l_i^{(2)}, \dots, l_i^{(K)}]$



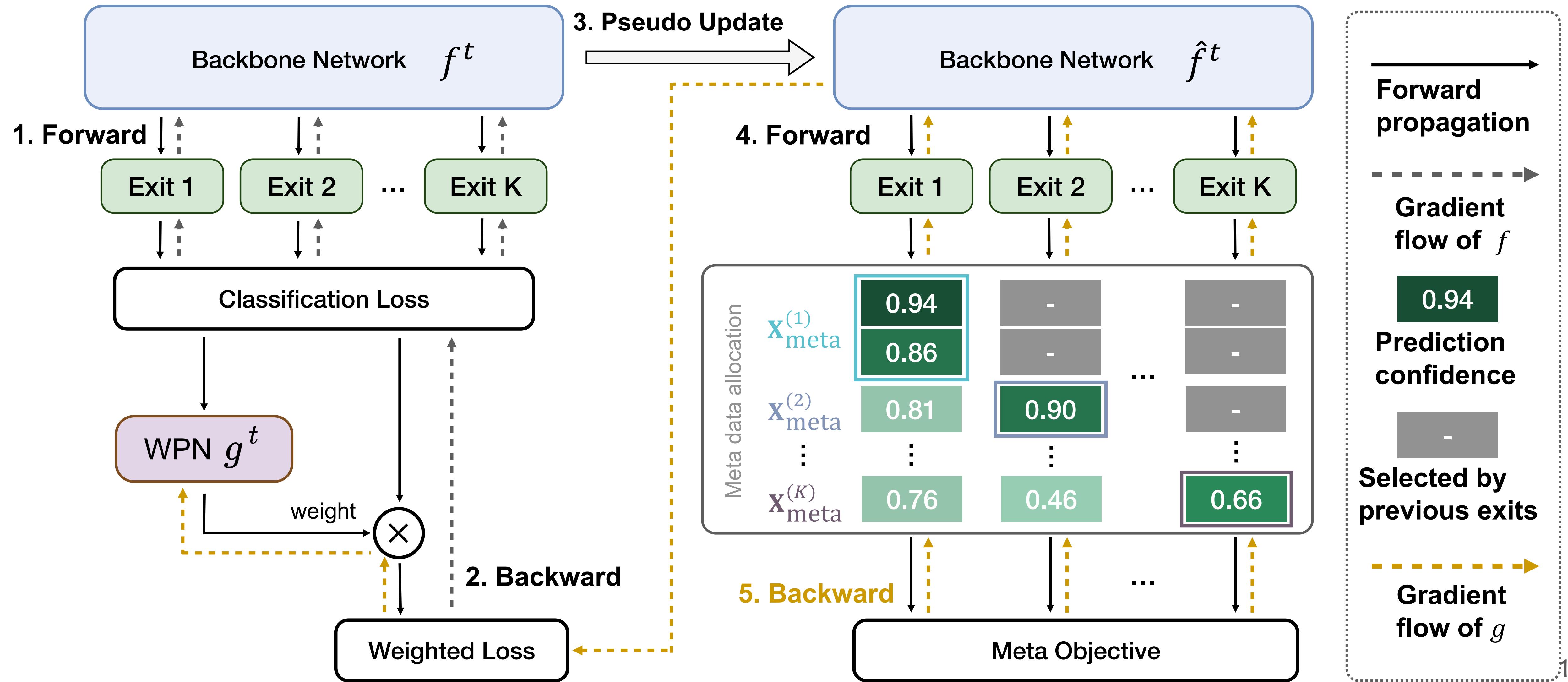
Meta-learning: training pipeline [5]

- $f(\cdot; \Theta_f)$: multi-exit classification network.
- $g(\cdot; \Theta_g)$: weight prediction network



We re-use the training dataset as our meta-dataset rather than using a standalone meta-dataset as in [5].

Meta-learning: optimization



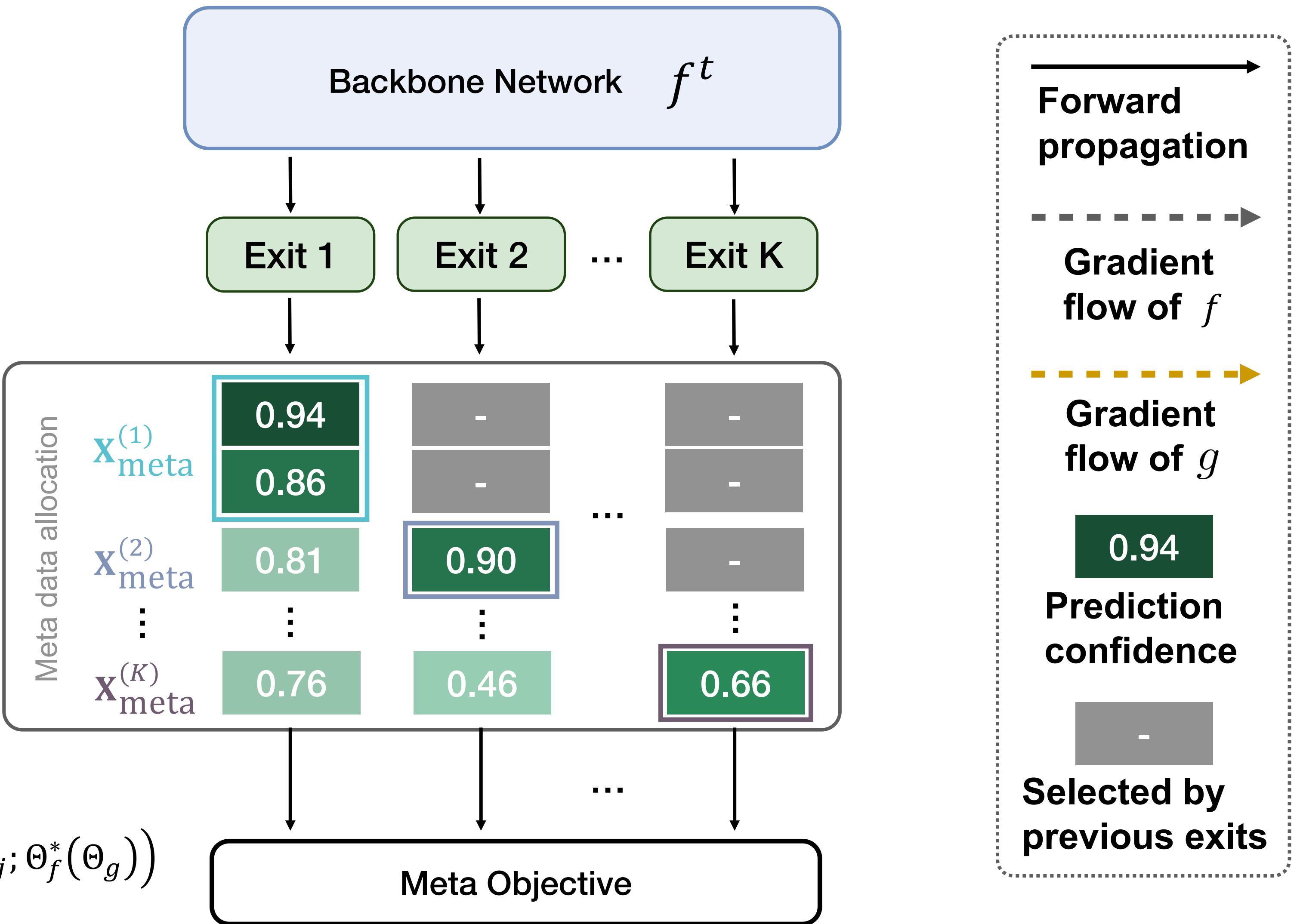
Meta-data allocation

$$|\mathbf{x}_{\text{meta}}^{(k)}| = \frac{q^k}{\sum_{k=1}^K q^k} N_{\text{meta}},$$

$q = 1.0$: uniform distribution

$$\Theta_g^* = \operatorname*{argmin}_{\Theta_g} \mathcal{L}_{\text{meta}} \left(\Theta_f^*(\Theta_g) \right)$$

$$\hat{\theta}_g \triangleq \operatorname{argmin}_{\Theta_g} \sum_{k=1}^K \frac{1}{N_k} \sum_{j : \mathbf{x}_j \in \mathbf{X}_{\text{meta}}^{(k)}} l^{(k)} \left(\mathbf{x}_j, y_j; \Theta_f^*(\Theta_g) \right)$$



Meta-learning objectives

Weighted CE Loss: optimizing the classification network f

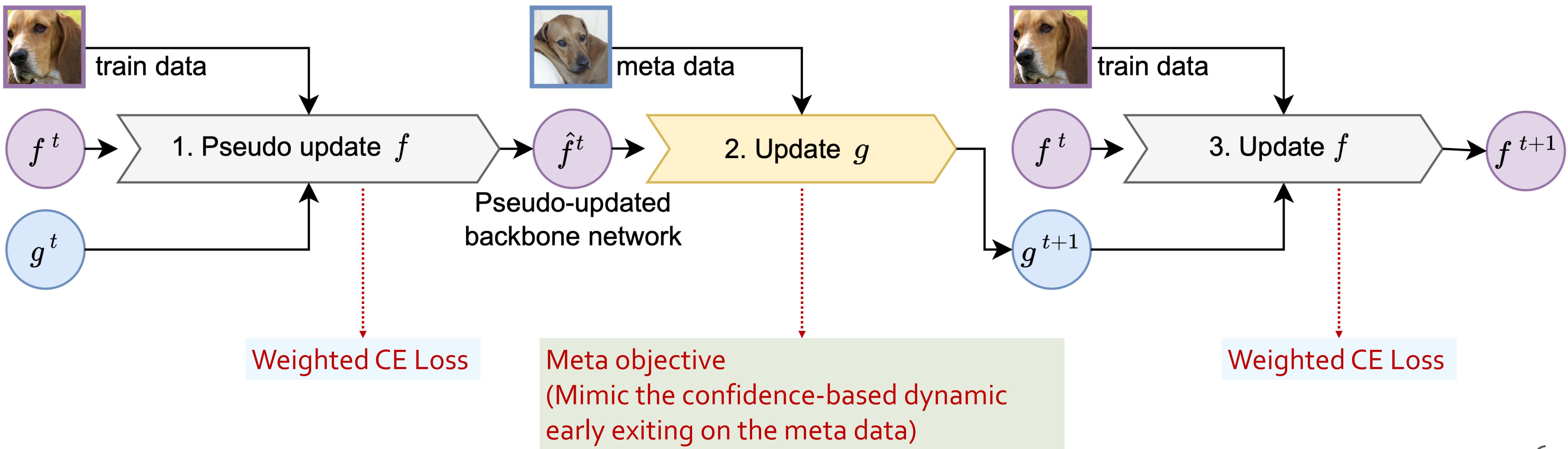
$$\begin{aligned}\Theta_f^*(\Theta_g) &= \operatorname{argmin}_{\Theta_f} \mathcal{L}_{\text{tr}}(\Theta_f, \Theta_g) \\ &\triangleq \operatorname{argmin}_{\Theta_f} \sum_{k=1}^K \frac{1}{N} \sum_{i: \mathbf{x}_i \in \mathbf{X}_{\text{tr}}} g^{(k)}(\mathbf{l}_i; \Theta_g) \cdot l_i^{(k)}(\Theta_f).\end{aligned}$$

Meta objective: optimizing the weight prediction network g

$$\begin{aligned}\Theta_g^* &= \operatorname{argmin}_{\Theta_g} \mathcal{L}_{\text{meta}}\left(\Theta_f^*(\Theta_g)\right) \\ &\triangleq \operatorname{argmin}_{\Theta_g} \sum_{k=1}^K \frac{1}{N_k} \sum_{j: \mathbf{x}_j \in \mathbf{X}_{\text{meta}}^{(k)}} l^{(k)}\left(\mathbf{x}_j, y_j; \Theta_f^*(\Theta_g)\right).\end{aligned}$$

Method: summary

- $f(\cdot; \Theta_f)$: multi-exit classification network.
- $g(\cdot; \Theta_g)$: weight prediction network



Experiments

- Model
 - Multi-scale Dense Network (MSDNet) [2]
 - Resolution Adaptive Network (RANet) [4]
- Dataset
 - ImageNet
 - CIFAR-10, CIFAR-100
 - Long-tailed CIFAR-100 [6]
- Other training techniques [3]

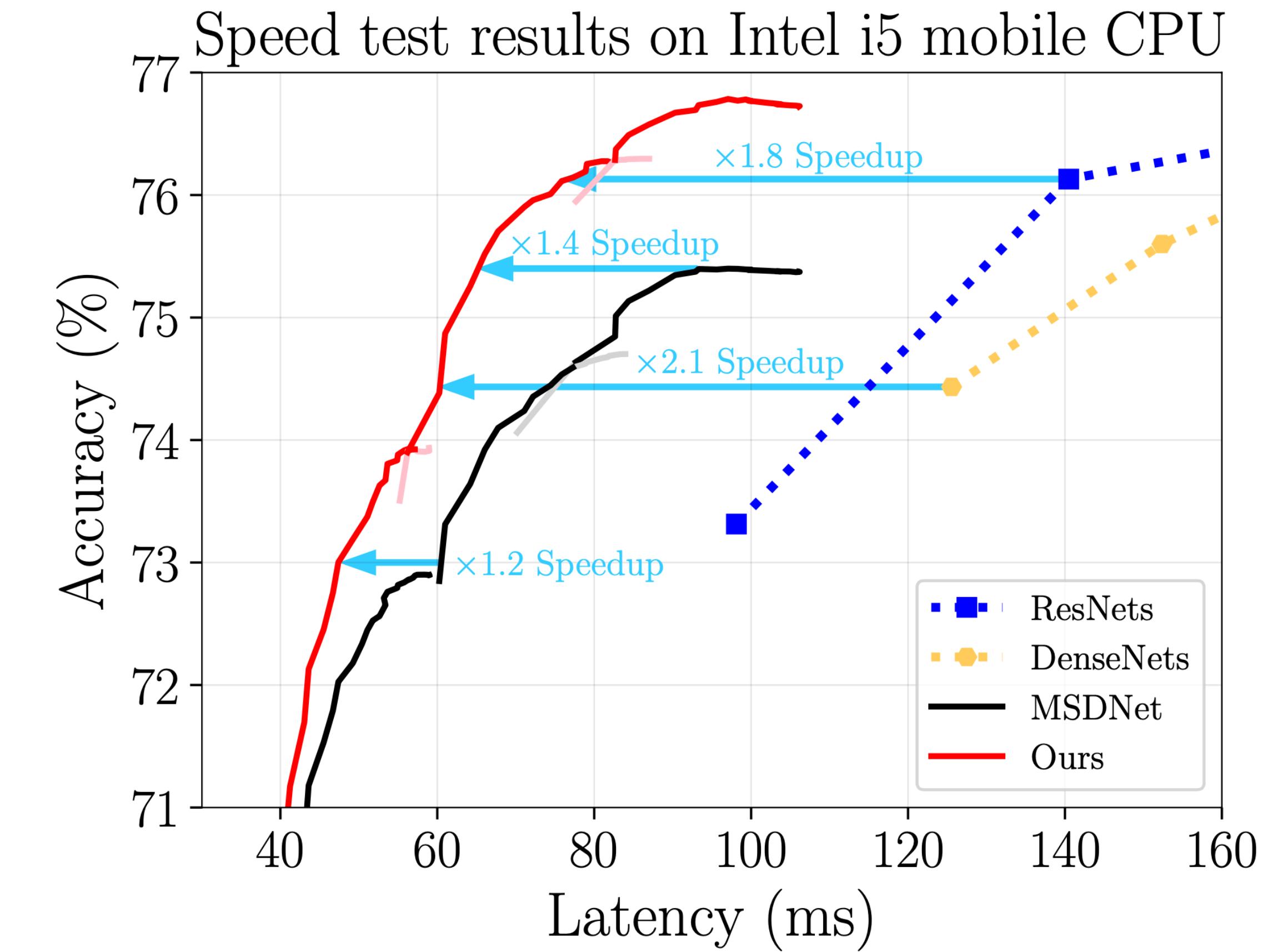
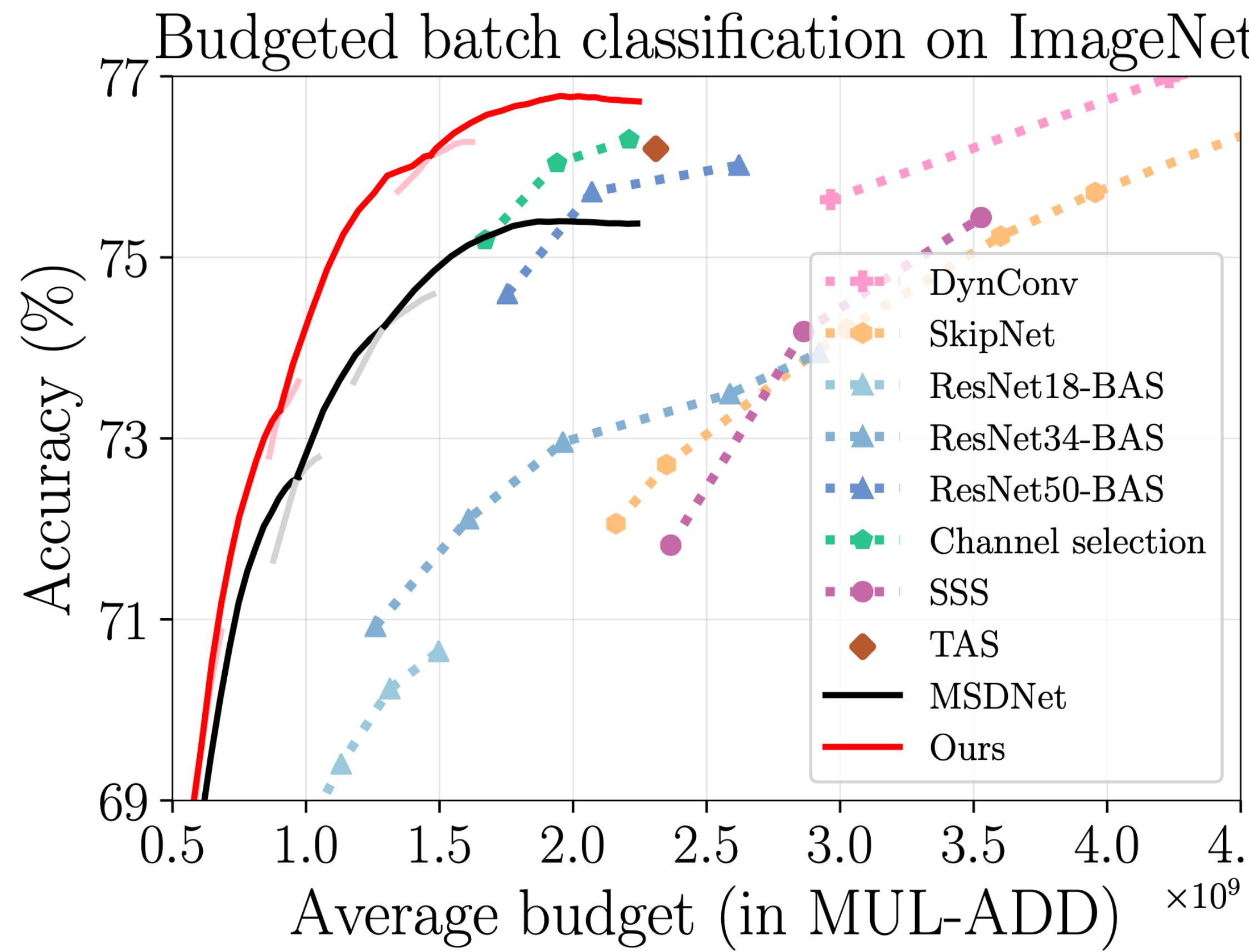
[2] Huang, G., Chen, D., Li, T., Wu, F., van der Maaten, L., & Weinberger, K. (2018, February). Multi-Scale Dense Networks for Resource Efficient Image Classification. In ICLR.

[3] Li, H., Zhang, H., Qi, X., Yang, R., & Huang, G. (2019). Improved techniques for training adaptive deep networks. In ICCV.

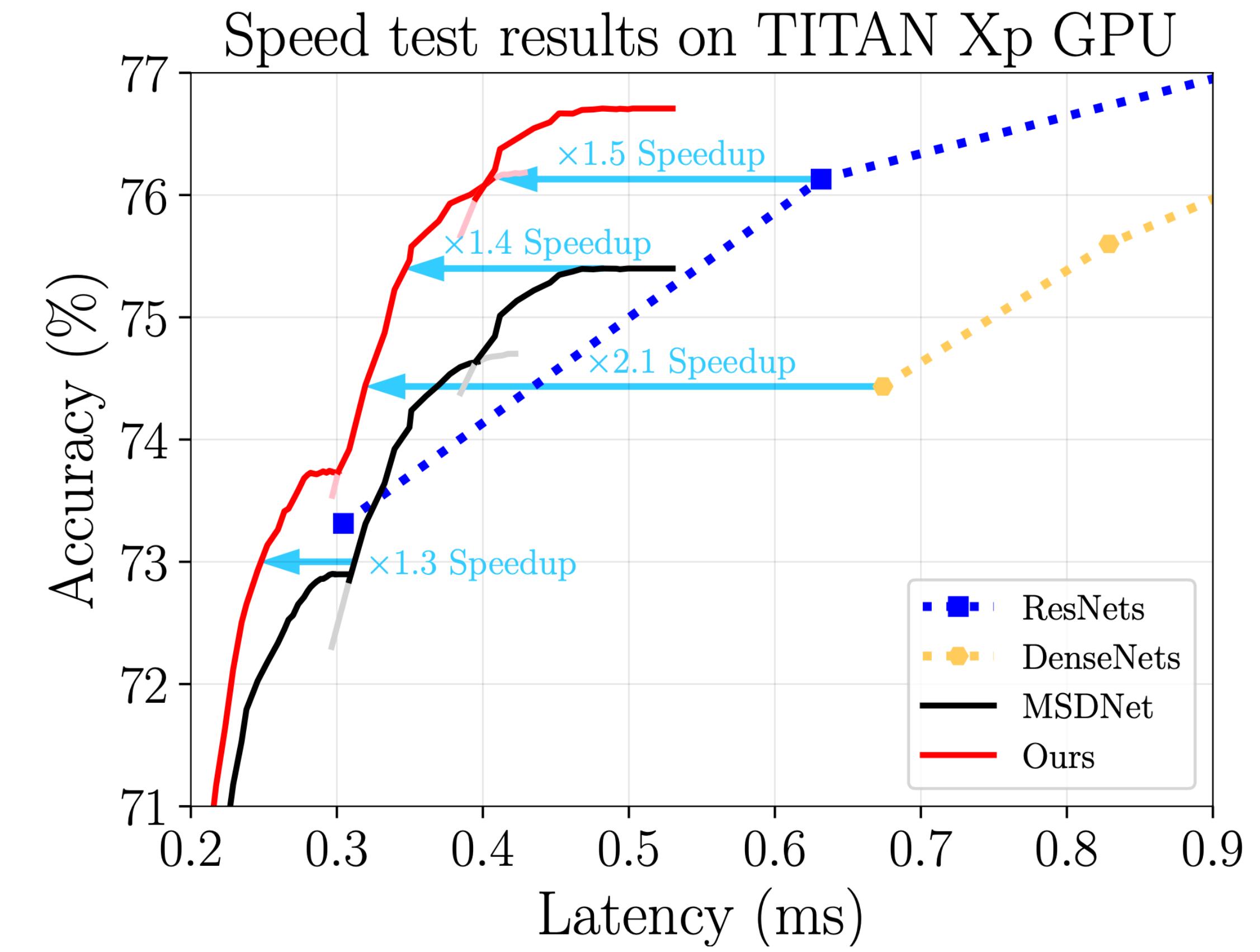
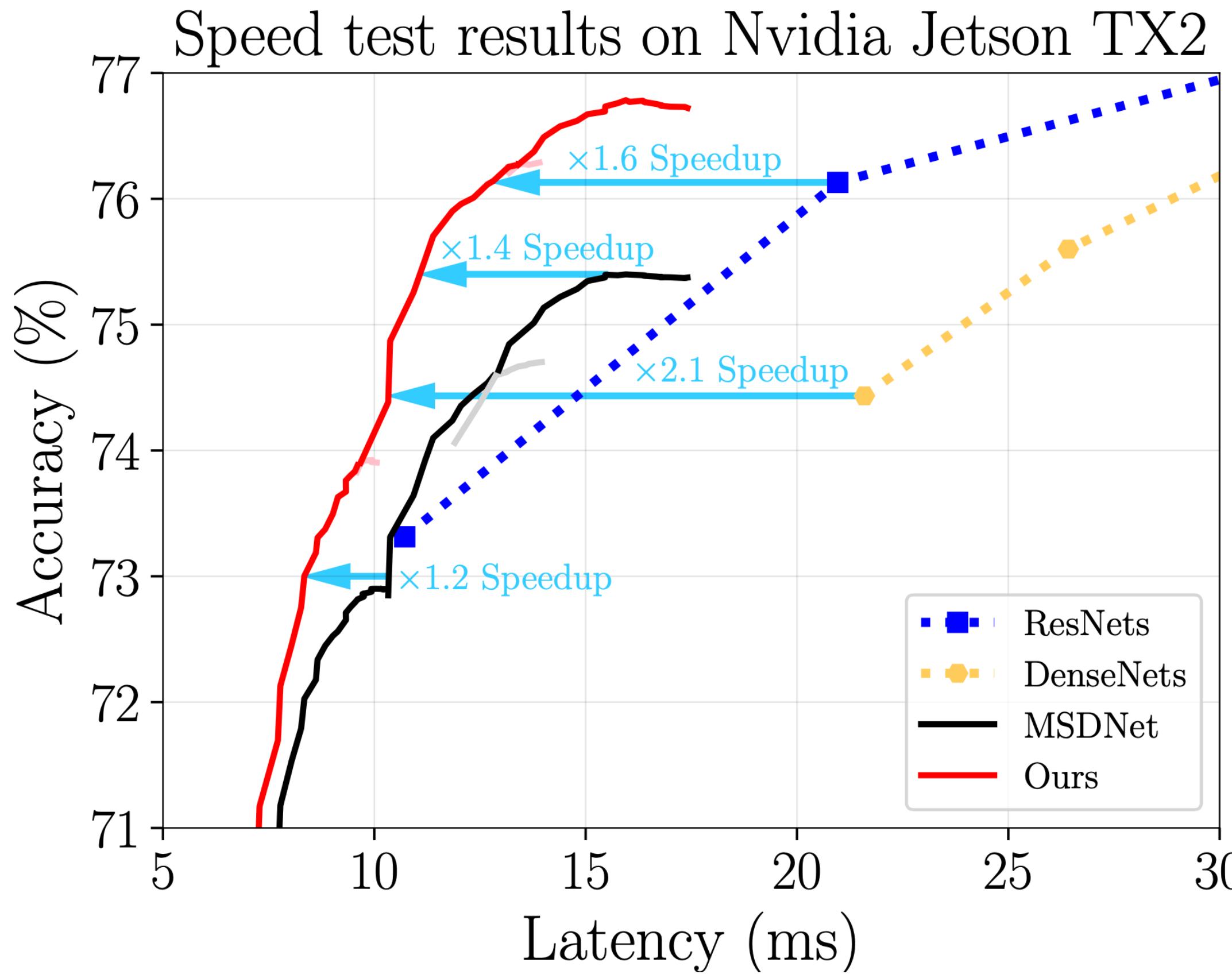
[4] Yang, L., Han, Y., Chen, X., Song, S., Dai, J., & Huang, G. (2020). Resolution adaptive networks for efficient inference. In CVPR.

[6] Cui, Y., Jia, M., Lin, T. Y., Song, Y., & Belongie, S. (2019). Class-balanced loss based on effective number of samples. In CVPR

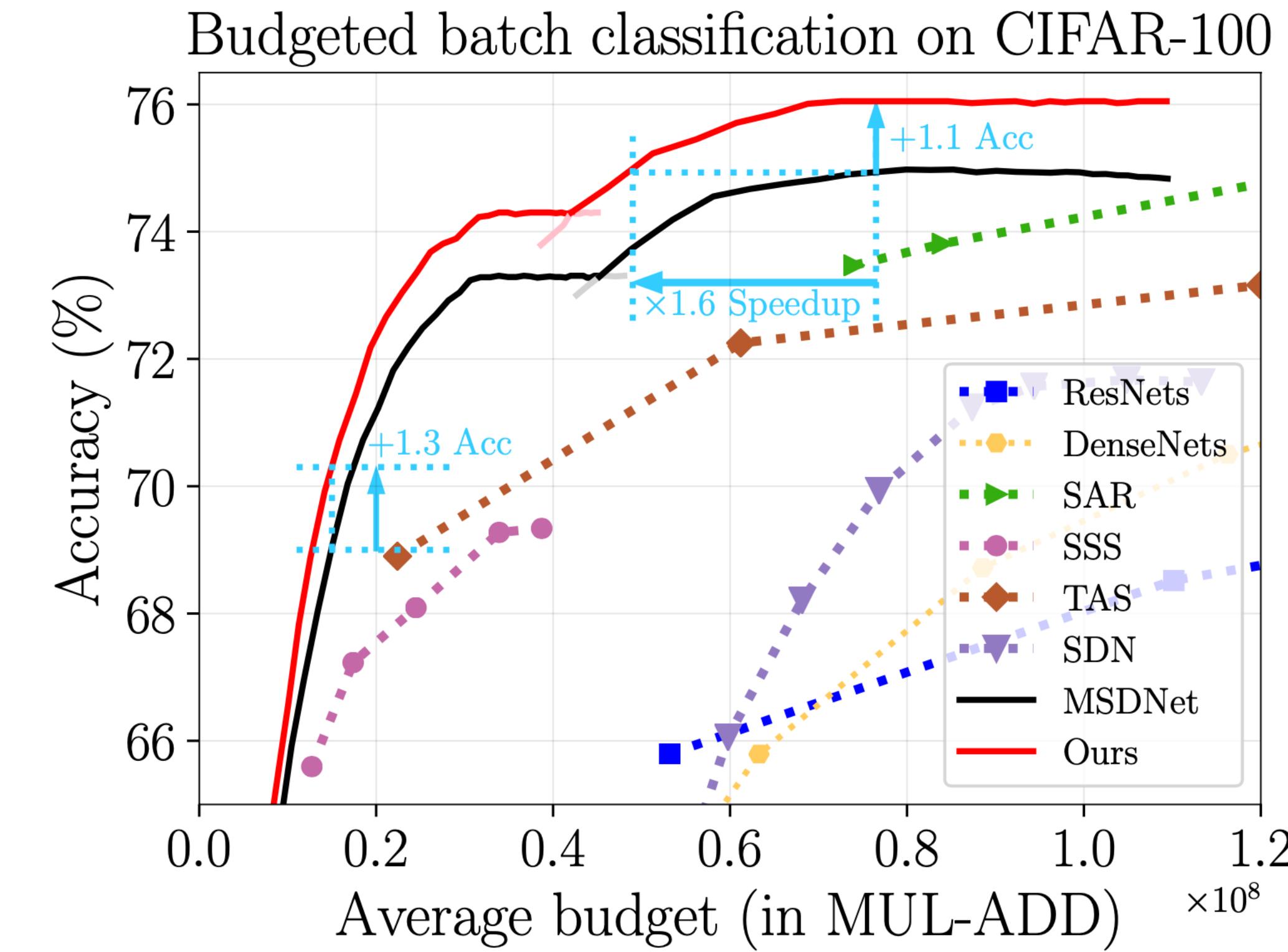
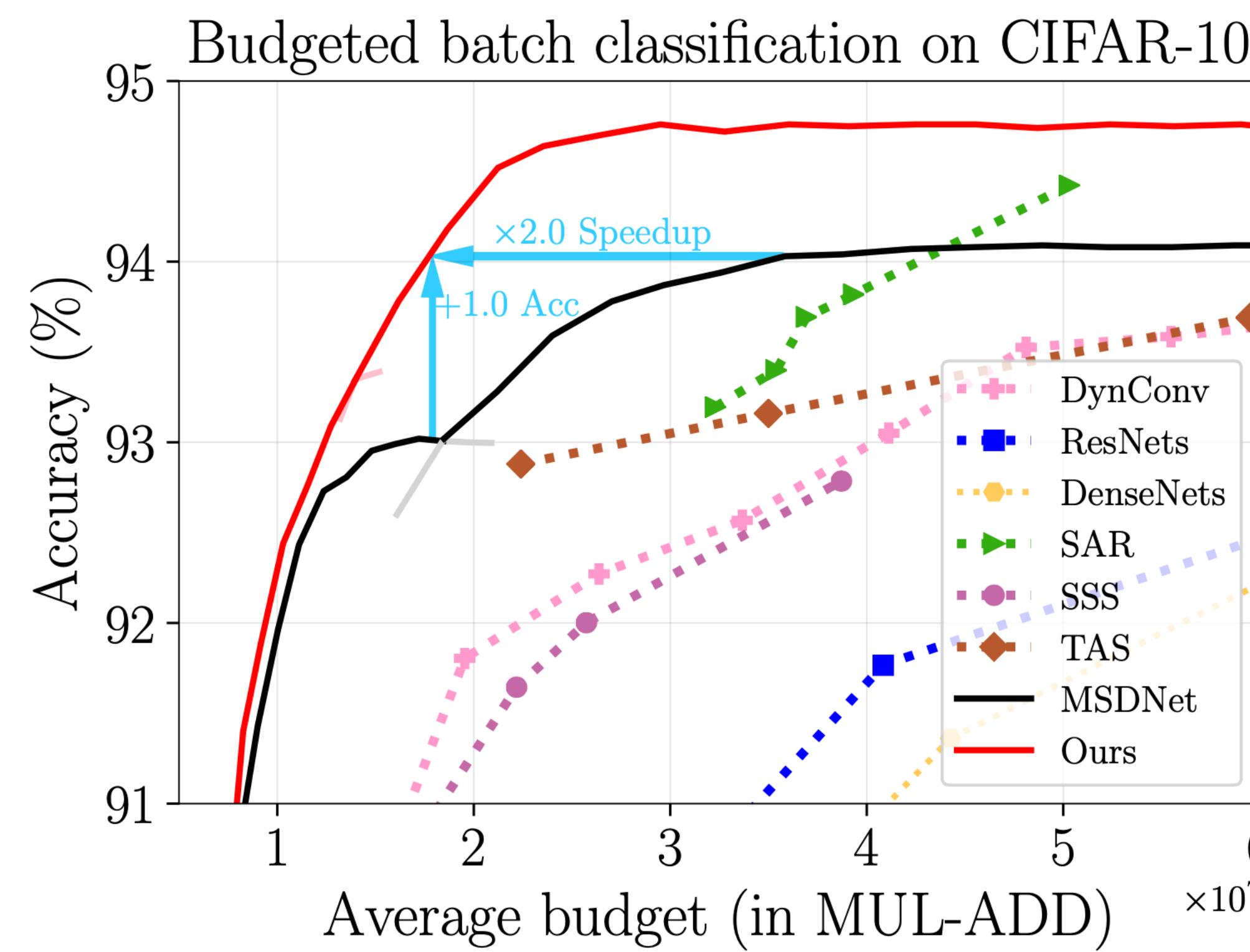
Experiments: Results on ImageNet with MSDNet [2]



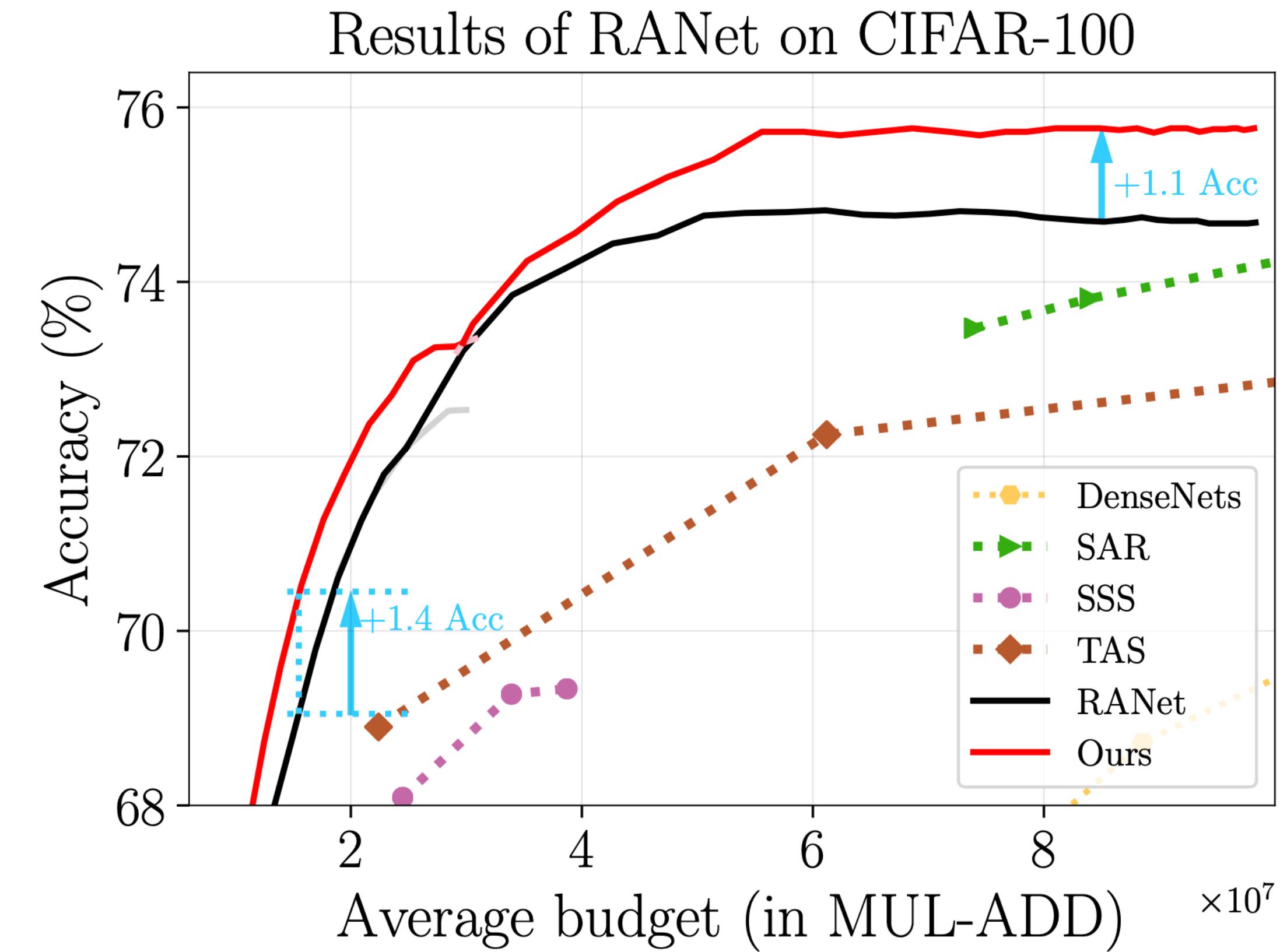
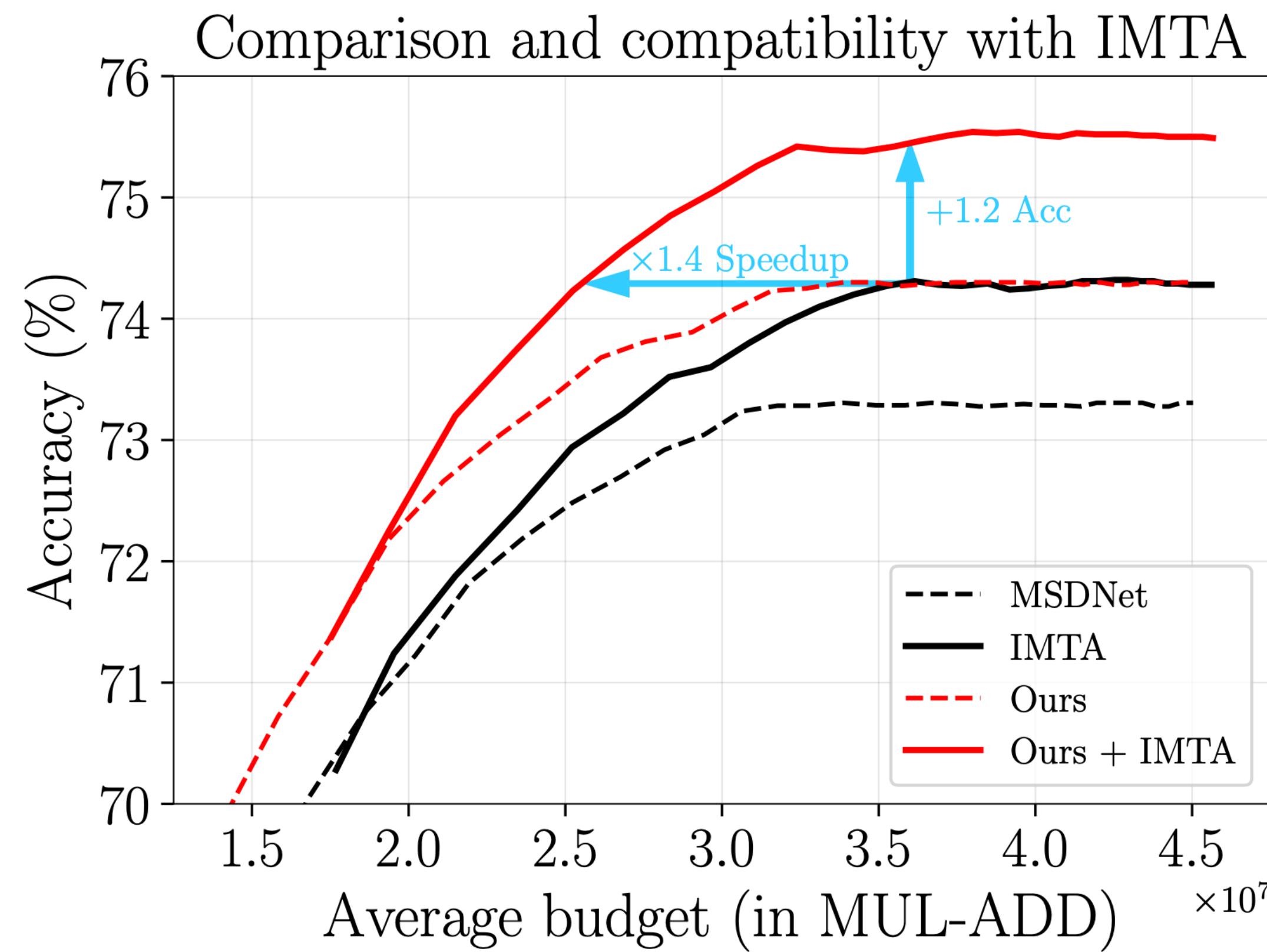
Experiments: Results on ImageNet with MSDNet [2]



Experiments: CIFAR results with MSDNet [2]



Experiments: Comparison with IMTA [3] and Results with RANet [4]

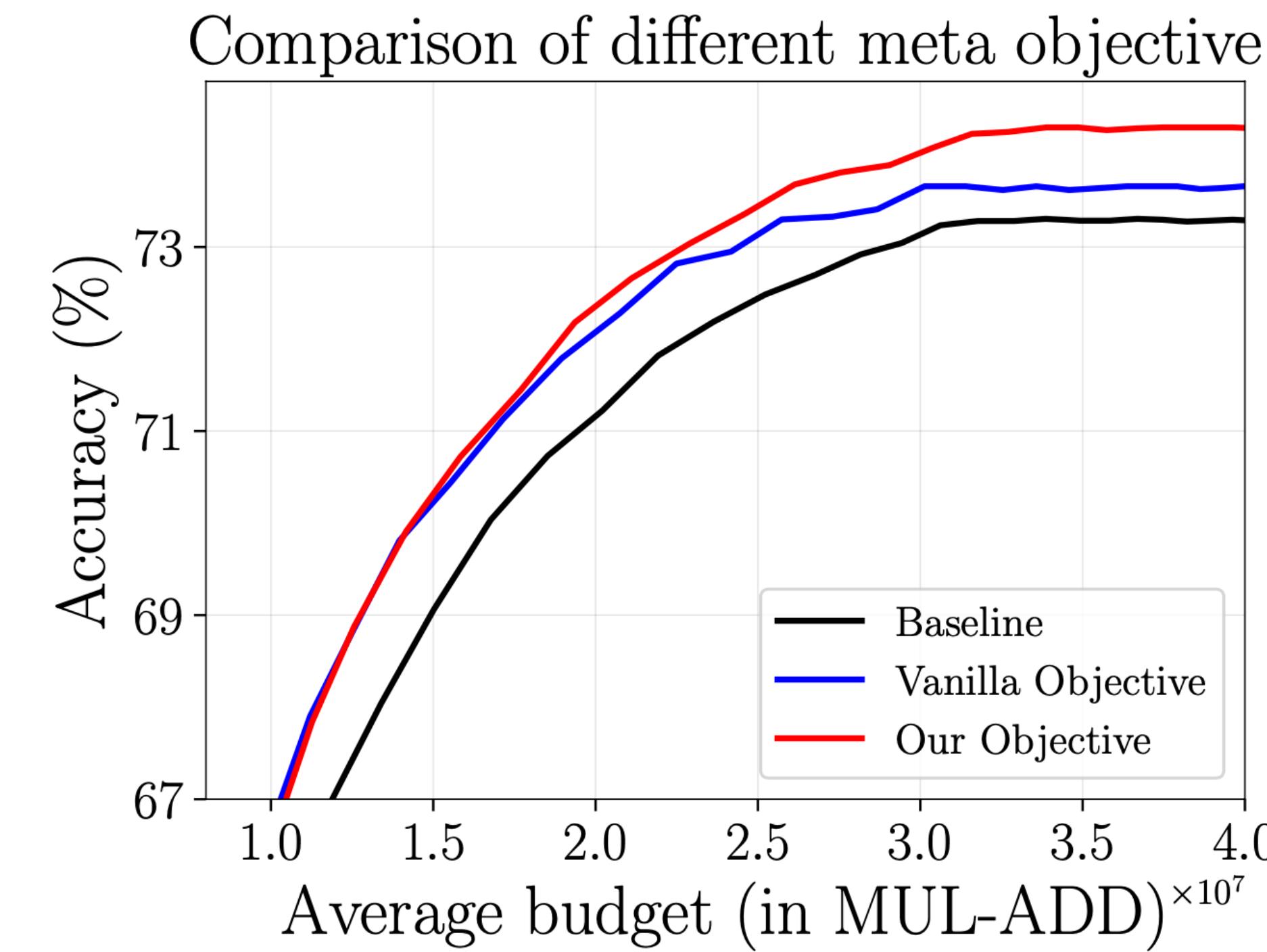
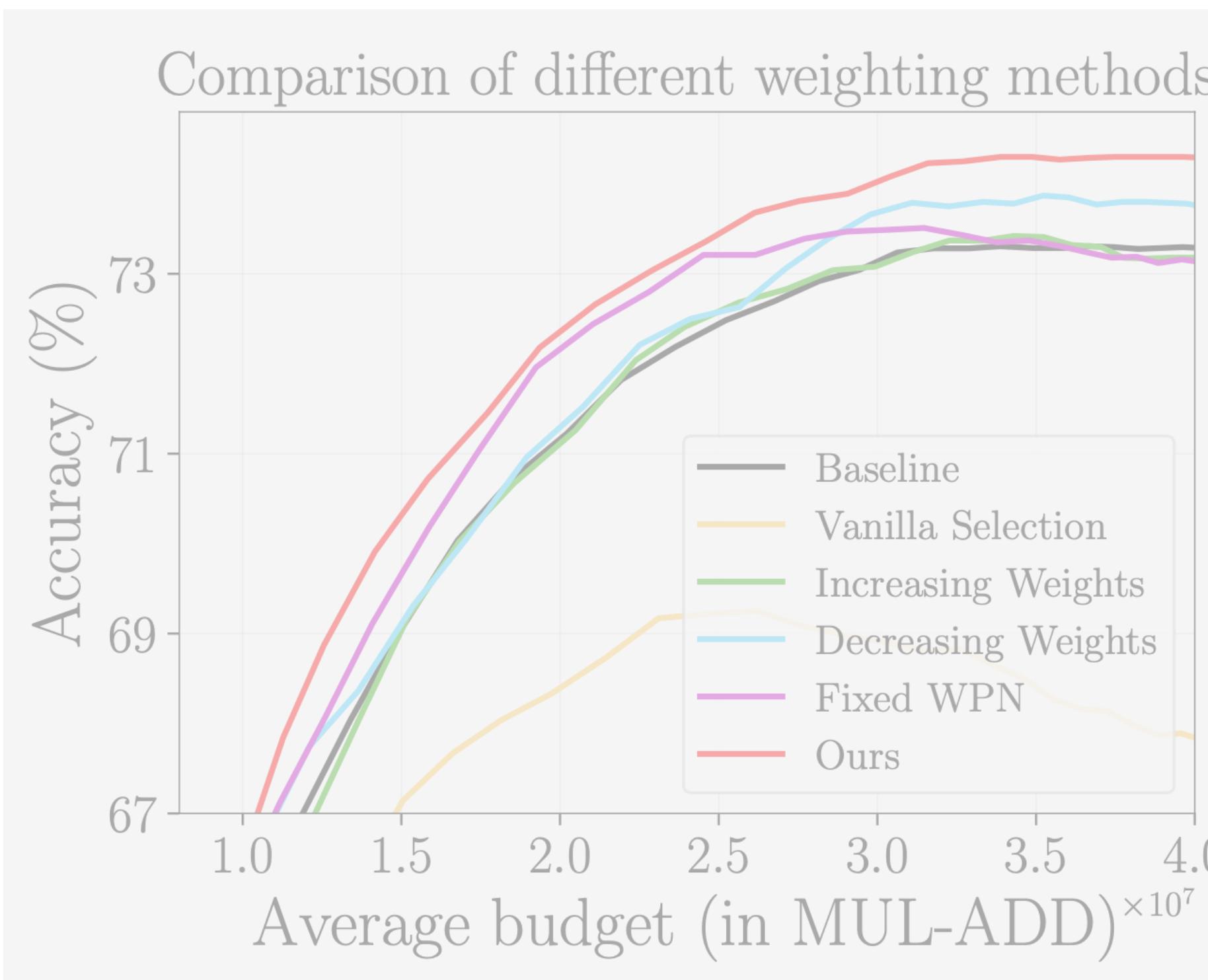


[3] Li, H., Zhang, H., Qi, X., Yang, R., & Huang, G. (2019). Improved techniques for training adaptive deep networks. In ICCV.

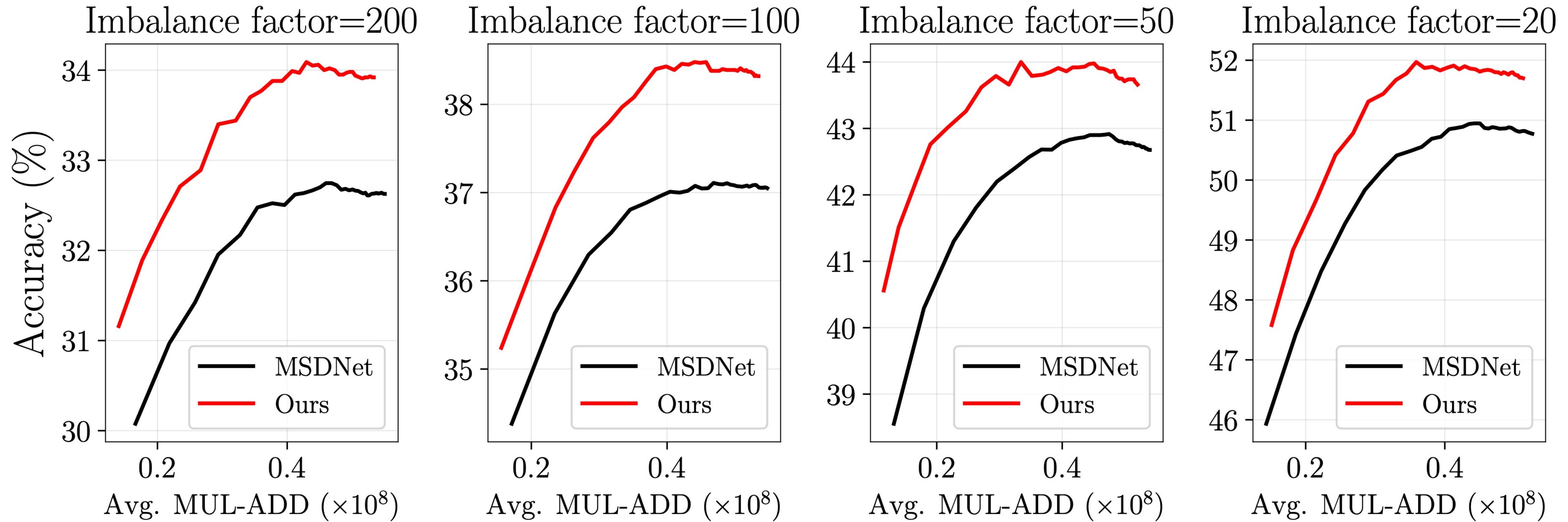
[4] Yang, L., Han, Y., Chen, X., Song, S., Dai, J., & Huang, G. (2020). Resolution adaptive networks for efficient inference. In CVPR.

Experiments: ablation studies on CIFAR-100 with MSDNet [1]

- Sample weighting v.s. Training data allocation
- Predicted weights v.s. fixed weights
- Meta-learning v.s. frozen weight prediction network
- Our meta objective v.s. vanilla meta objective



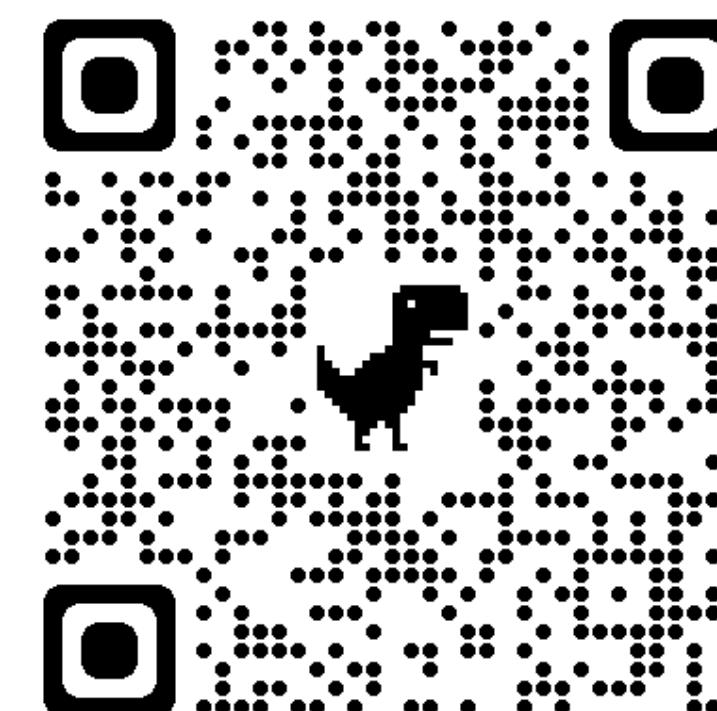
Experiments: Results in the class-imbalanced setting



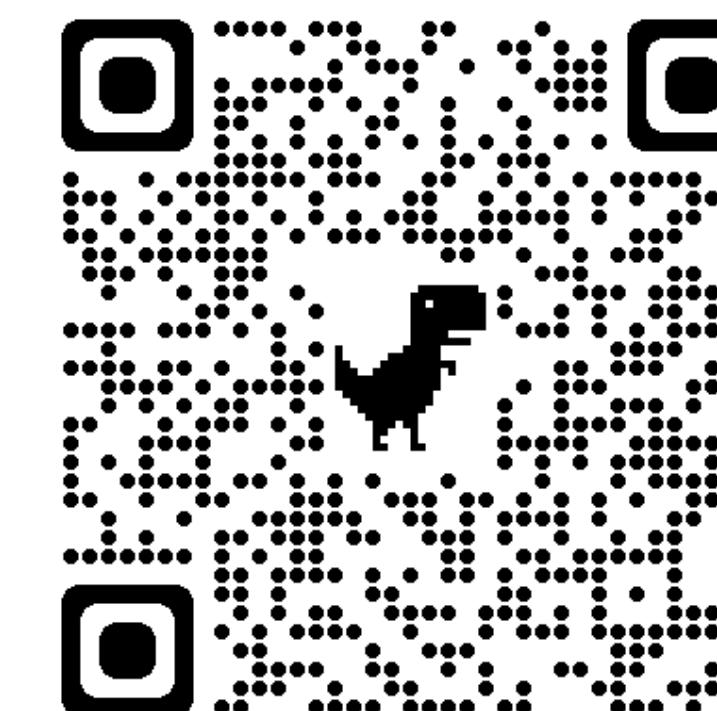
Conclusion (takeaway messages)

- [Observation] Different exits (classifiers) in dynamic early exiting networks are responsible for samples with **different complexity** at test time.
- [Intuition] Different samples (with varying **complexity**) should contribute **unequally** to different exits (with varying **capacity**) in training.
- [Solution] The gap between training and test can be bridged by **sample weighting** rather than sample selection (data allocation).
- [Key design] **Sample selection (data allocation)** in calculating the **meta objective** can successfully guide the learning of weight prediction.

Thank you!
Welcome to our poster



Paper



Code



清华大学
Tsinghua University