# DarkDistill: Difficutly-Aligned Federated Early-Exit Network Training on Heteregeneous Devices

Lehao Qu[1], Shuyuan Li[2], Zimu Zhou[2], Boyi Liu[1,2], Yi Xu[1], Yongxin Tong[1]
[1]Beihang University, [2]City University of Hong Kong
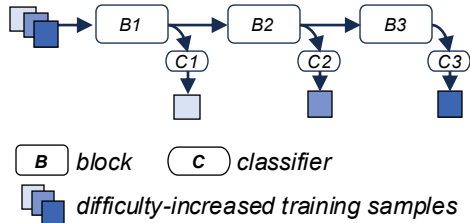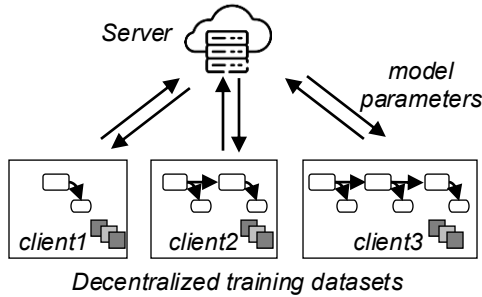
## Abstract

Early-exit networks (EENs), which adapt their computational depths based on input samples, are widely adopted to accelerate inference in edge computing applications. The effectiveness of EENs relies on difficulty-aware training, which tailors shallow exits for simple samples and deep exits for complex ones. However, existing difficulty-aware training schemes assume centralized environments with sufficient data, which become invalid with real-world edge devices.
In this paper, we explore difficulty-aware training in a federated manner, where EENs are collaboratively trained on heterogeneous devices. We observe the *cross-model exit unalignment phenomenon*, a unique problem when aggregating local EENs into a cohesive global model. To address this problem, we design a novel *Difficulty-Aligned Reverse Knowledge Distillation* scheme named DarkDistill that preserves the difficulty-specific specialization for aggregating heterogeneous local models. Instead of direct parameter averaging, it trains difficulty-conditional data generators, and selectively transfers generated knowledge of specific difficulty among matched exits of heterogeneous EENs. Evaluations show that DarkDistill outperforms the state-of-the-arts in various fine-tuning of EENs.
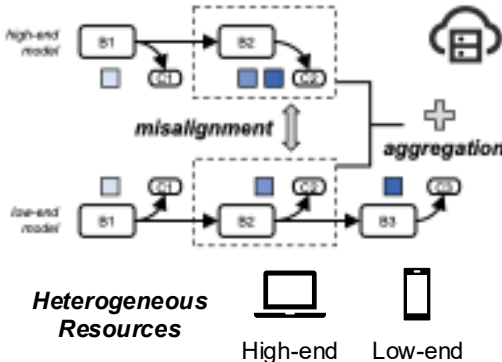
## Introduction

**Early-Exit Network** can adjust depth based on the difficulty of the input samples. *Easy (difficult)* samples terminate at *shallow (deep)* exits



B  block    C  classifier

difficulty-increased training samples

**Federated Learning EEN Training** leverages the data knowledge of federated clients with heterogenous resources to train the *global EEN*



*Decentralized training datasets*

**Cross-model Exit Unalignment**
*Exits at equivalent depths may handle samples from disparate difficulty ranges across models*



*Heterogeneous Resources*
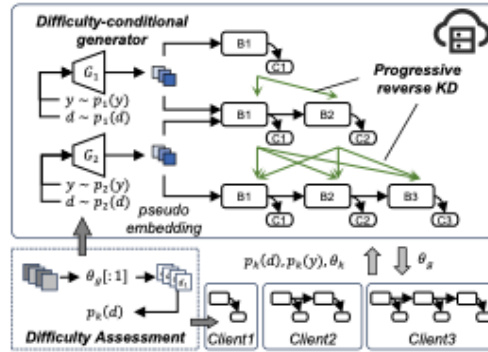
High-end    Low-end

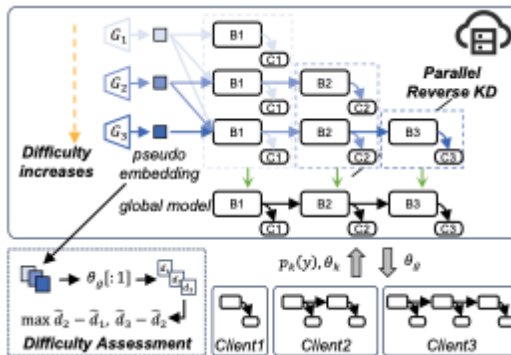***New setting, new challenges!***

## Methods

**DarkDistill: Progressive Difficulty-Aligned Reverse Knowledge Distillation**

**1. Difficulty Assessment** evaluates the difficulty range of local data utilize its loss on global model

***2. Difficulty-Conditional Generators*** create pseudo data for *specific difficulties*, supporting the knowledge distillation process

*3. Progressive Reverse KD* transfers knowledge from shallow to deep exits in adjacent layers across varied depth local EENs



**DarkDistill-PL: Parallel Variant** simultaneously distills the *ensemble knowledge* of all immediate knowledge to the global model parameterized by in an exit-wise manner
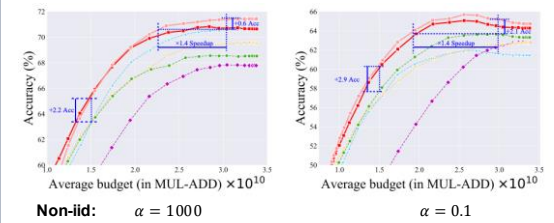


***Difficulty-Aligned Knowledge Distillation***

## Experiments

**Anytime Performance:** Inference may terminate at any time. Show the average and variance of the accuracy across all exits on various datasets

| Finetune | Difficulty-aware | Approach | CIFAR-100 [19] | | | SVHN [30] | SpeechCmds [44] |
|---|---|---|---|---|---|---|---|
| | | | $\alpha = 0.1$ | $\alpha = 1$ | $\alpha = 1000$ | | |
| Full | None | ExclusiveFL | $26.60_{\pm3.10}$ | $49.96_{\pm11.48}$ | $41.58_{\pm7.01}$ | $85.28_{\pm2.97}$ | $87.00_{\pm2.88}$ |
| | | InclusiveFL [26] | $40.10_{\pm2.03}$ | $58.83_{\pm6.98}$ | $61.40_{\pm7.01}$ | $82.95_{\pm0.34}$ | $91.90_{\pm1.42}$ |
| | | ScaleFL [16] | $54.99_{\pm10.61}$ | $63.21_{\pm9.14}$ | $63.82_{\pm9.87}$ | $88.24_{\pm0.78}$ | $92.56_{\pm0.26}$ |
| | | DepthFL [18] | $40.70_{\pm1.57}$ | $59.01_{\pm5.18}$ | $61.71_{\pm5.75}$ | $83.45_{\pm0.43}$ | $92.05_{\pm0.60}$ |
| | | ReeFL [23] | $59.24_{\pm4.60}$ | $63.37_{\pm2.72}$ | $63.90_{\pm8.68}$ | $88.37_{\pm1.27}$ | $93.12_{\pm1.14}$ |
| | BoostNet [45] | ExclusiveFL | $48.68_{\pm13.66}$ | $57.57_{\pm15.12}$ | $58.65_{\pm15.31}$ | $87.30_{\pm2.89}$ | $91.07_{\pm2.58}$ |
| | | InclusiveFL [26] | $57.10_{\pm7.21}$ | $62.96_{\pm8.12}$ | $64.01_{\pm8.24}$ | $87.86_{\pm1.66}$ | $92.91_{\pm1.10}$ |
| | | ScaleFL [16] | $52.74_{\pm15.82}$ | $60.55_{\pm11.93}$ | $60.73_{\pm10.80}$ | $87.91_{\pm0.77}$ | $92.03_{\pm0.37}$ |
| | | DepthFL [18] | $58.15_{\pm6.73}$ | $63.81_{\pm6.34}$ | $64.19_{\pm6.73}$ | $87.74_{\pm1.01}$ | $92.72_{\pm0.64}$ |
| | | ReeFL [23] | $59.01_{\pm7.98}$ | $63.08_{\pm9.03}$ | $63.66_{\pm7.31}$ | $88.39_{\pm1.28}$ | $93.01_{\pm1.18}$ |
| | | DarkDistill | $60.68_{\pm7.93}$ | $64.50_{\pm7.97}$ | $65.67_{\pm7.48}$ | $88.41_{\pm1.46}$ | $93.31_{\pm1.13}$ |
| | | DarkDistill-PL | $61.05_{\pm8.19}$ | $65.12_{\pm7.02}$ | $65.49_{\pm7.88}$ | $88.48_{\pm1.57}$ | $93.42_{\pm0.98}$ |
| LORA [13] | None | ExclusiveFL | $44.44_{\pm18.61}$ | $52.35_{\pm18.56}$ | $52.88_{\pm18.17}$ | $83.78_{\pm4.43}$ | $88.72_{\pm3.15}$ |
| | | InclusiveFL [26] | $44.82_{\pm25.36}$ | $54.26_{\pm21.38}$ | $54.76_{\pm21.37}$ | $85.16_{\pm5.31}$ | $89.58_{\pm5.04}$ |
| | | ScaleFL [16] | $22.17_{\pm23.36}$ | $30.85_{\pm30.58}$ | $32.58_{\pm31.96}$ | $76.42_{\pm13.14}$ | $58.82_{\pm34.80}$ |
| | | DepthFL [18] | $52.17_{\pm14.16}$ | $57.09_{\pm14.78}$ | $57.84_{\pm14.48}$ | $85.69_{\pm2.71}$ | $90.11_{\pm2.04}$ |
| | | ReeFL [23] | $52.32_{\pm9.83}$ | $57.74_{\pm11.78}$ | $58.16_{\pm11.69}$ | $85.54_{\pm3.00}$ | $89.56_{\pm2.62}$ |
| | BoostNet [45] | ExclusiveFL | $50.34_{\pm13.65}$ | $55.68_{\pm15.33}$ | $56.48_{\pm15.33}$ | $84.48_{\pm3.88}$ | $88.51_{\pm2.26}$ |
| | | InclusiveFL [26] | $54.25_{\pm11.78}$ | $59.66_{\pm11.72}$ | $59.81_{\pm11.46}$ | $85.96_{\pm2.50}$ | $90.38_{\pm2.10}$ |
| | | DepthFL [18] | $40.46_{\pm22.36}$ | $47.18_{\pm24.11}$ | $48.26_{\pm24.17}$ | $81.70_{\pm3.14}$ | $80.19_{\pm4.83}$ |
| | | ReeFL [23] | $55.85_{\pm9.45}$ | $60.95_{\pm9.20}$ | $61.45_{\pm9.04}$ | $79.90_{\pm1.61}$ | $90.93_{\pm1.29}$ |
| | | DarkDistill | $57.32_{\pm11.91}$ | $61.24_{\pm11.06}$ | $61.74_{\pm11.42}$ | $86.11_{\pm2.16}$ | $91.06_{\pm2.08}$ |

**Budget Performance:** Evaluate the accuracy of a batch across various budget. DarkDistill is faster and more accurate across multi-settings



**Non-iid:**    $\alpha = 1000$    $\alpha = 0.1$

***Faster and more accurate***

## Theorem

**Convergence Analysis** DarkDistill converges in FL with heterogeneous clients
If the learning rate $\gamma$ of local training satisfies $\frac{1}{T\sqrt{Q}} \leq \gamma \leq \frac{1}{6M^2LT}$, DarkDistill coverages to a neighborhood of a stationary point of **standard FL** as follows:

$$\frac{1}{Q}\sum_{q=1}^{Q}\mathbb{E}\|\nabla\mathcal{L}(\theta^q)\|^2 \leq \frac{G_0}{\sqrt{Q}} + V_0 + \frac{H_0}{T} + \frac{I_0}{\sqrt{Q}}\sum_{q=1}^{Q}\mathbb{E}\|\theta^q\|^2$$

**Explaination** DarkDistill converges under a properly chosen learning rate $\gamma$, which can be practically set using the local epoch count $T$, total communication round $Q$, loss smoothness $L$, and largest exit number $M$

***Guide to choose suitable lr***