# DarkDistill

*2025-08-05*

# DarkDistill: Difficulty-Aligned Federated Early-Exit Network Training on Heterogeneous 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

# Outline

- ## Background & Motivation

- ## Problem Statement

- ## Our Solutions

- ## Experiments

- ## Conclusion

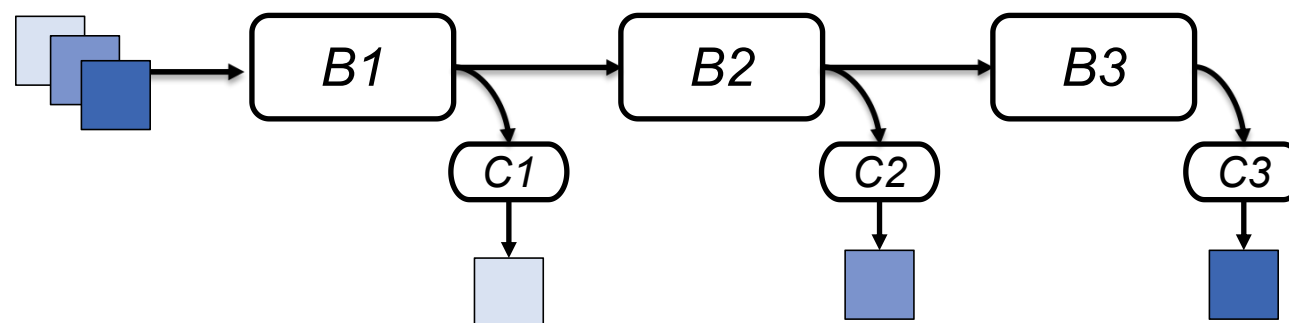# Outline

- **Background & Motivation**

- **Problem Statement**

- **Our Solutions**

- **Experiments**

- **Conclusion**

- **Early-Exit Network (EEN)**



| | |
|---|---|
| $B$ Block | $C$ Classifier |

*Difficulty-increased inference samples*

**Traffic Analysis**

**Autonomous Driving**

**Well-being Monitoring**

- **Difficulty-aware EEN Training**



B — Block    C — Classifier

Difficulty-increased training samples

*Training objective*

$$\mathcal{L}(\theta; D) = \sum_{m=1}^{M} \omega^m \mathcal{L}^m(\theta; D) = \sum_{m=1}^{M} \omega^m \sum_{i=1}^{|D|} l_i^m$$

*Core ideas*

1) BoostNet: Directing samples misclassified by shallow exits to deep ones

2) L2w: Increasing the weight of complex samples on training deep exits

# Background & Motivation

- **Federated learning EEN training**



Server

Model Aggregation

model parameters

client1

client2

client3

Local EEN Training

Decentralized training datasets

# Outline

- **Background & Motivation**

- **Problem Statement**

- **Our Solutions**

- **Experiments**

- **Conclusion**

- ## Federated EEN Training on Heterogeneous Devices



**Heterogeneous Resource**

High-end     Low-end

low-end model

global model

high-end model

*mismatch*

*aggregation*

**B** Block    **C** Classifier

*Difficulty-increased training samples*

- **Federated EEN Training on Heterogeneous Devices**



low-end model

global model

high-end model

mismatch

aggregation

| B | Block | | C | Classifier |

Difficulty-increased training samples

**Heterogeneous Resource**

High-end

Low-end

**Cross-model Exit Unalignment**
*Exits at equivalent depths may handle samples from disparate difficulty ranges across models*
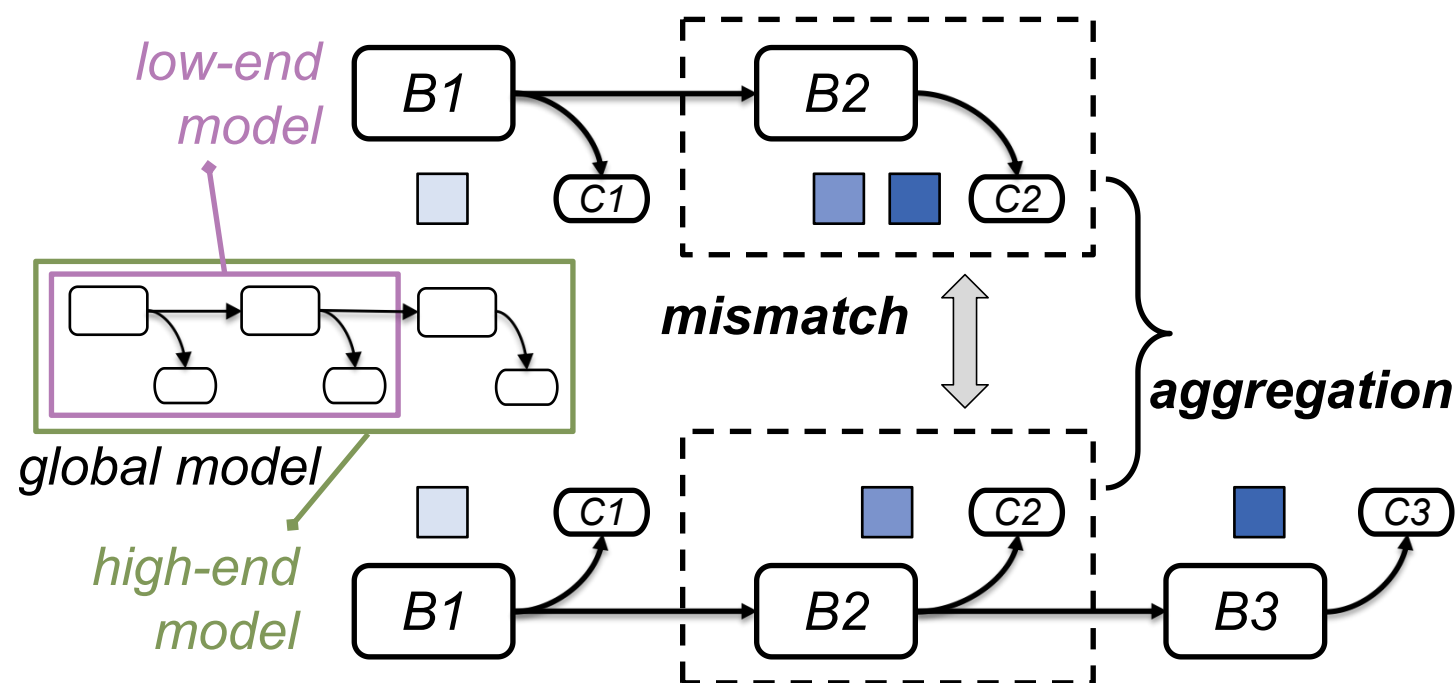
**Challenge:**
**How to solve the Unalignment?**

# Outline

- **Background & Motivation**

- **Problem Statement**

- **Our Solutions**

- **Experiments**

- **Conclusion**

# Our Solutions

- **DarkDistill: Framework**
  - **Progressive Knowledge Distillation**

*Generator*

*Create pseudo-data for specific difficulties, supporting the knowledge distillation process*

*Progressive KD*

*Transfer knowledge from shallow to deep exits in adjacent layers across varied depth local EENs*

# Our Solutions

- **DarkDistill: Workflow**

**Difficulty Assessment**

Predict the difficulty of samples

0.50

0.25

0-1  1-2  2-3

Samples

Difficulty Range

Difficulty distribution $p(d)$

**Client**

1. Upload $p(d)$

3. Download global model

**Server**

**Difficulty-Conditional Generator**

2. Produce pseudo data

**Difficulty-Aligned KD**

label
$y \sim p(y)$

G

Pseudo data

difficulty
$d \sim p(d)$

Low-end

B1 → B2
C1 → C2

High-end

C1 → C2 → C3
B1 → B2 → B3

- ## **DarkDistill: Difficulty Assessment**

  - Inspired by Curriculum Learning, we utilize the loss of sample to respect its difficulty. *The bigger the loss, the harder it is.*

  - In order to uniformly measure difficulty across clients, we leverage the global model to calculate the loss.

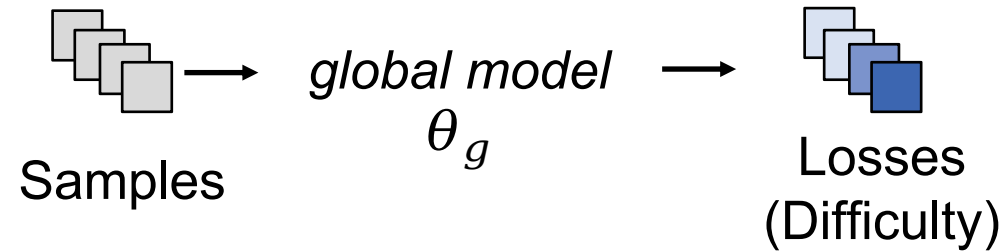Samples $\longrightarrow$ *global model* $\theta_g$ $\longrightarrow$ Losses (Difficulty)
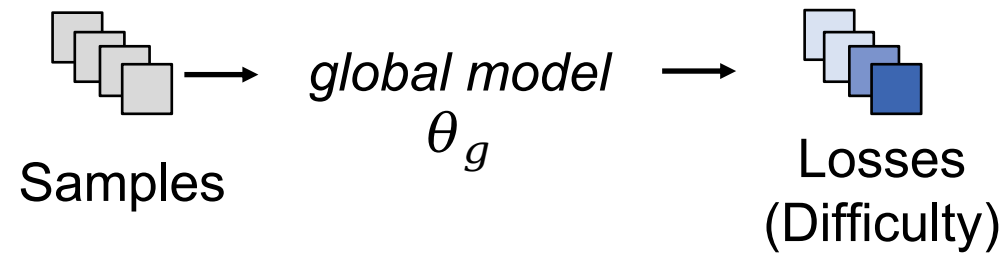
# Our Solutions

- ## **DarkDistill: Difficulty Assessment**

  - Inspired by Curriculum Learning, we utilize the loss of sample to respect its difficulty. *The bigger the loss, the harder it is.*

  - In order to uniformly measure difficulty across clients, we leverage the global model to calculate the loss.
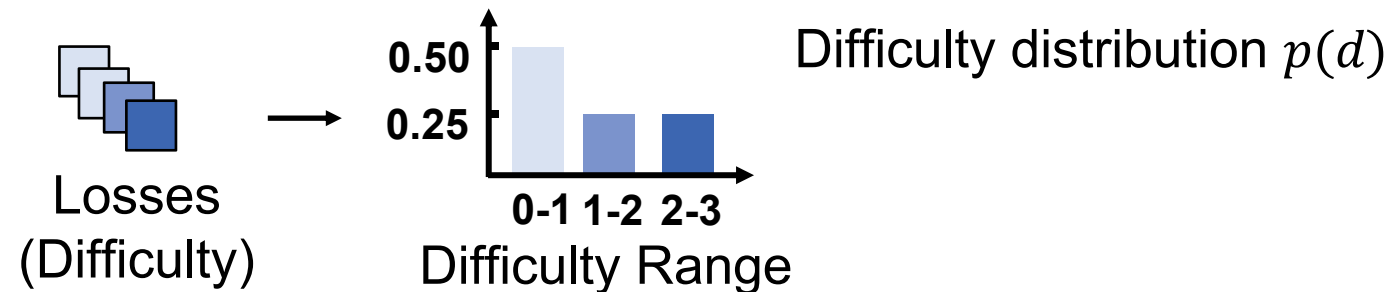


Samples → *global model* $\theta_g$ → Losses (Difficulty)
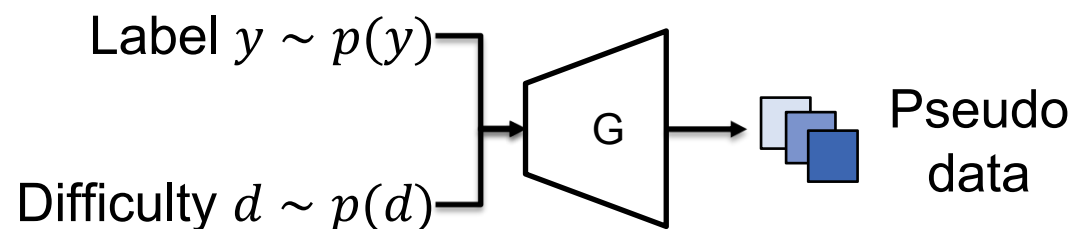
  - Calculate difficulty distribution $p(d)$ for *privacy*



Losses (Difficulty) → Difficulty distribution $p(d)$

0.50
0.25
0-1 1-2 2-3
Difficulty Range

- ## **DarkDistill: Difficulty-Conditional Generator**

  - *Create pseudo-data for specific difficulty and label to simulate local datasets, supporting the knowledge distillation process*



Label $y \sim p(y)$ — G → Pseudo data
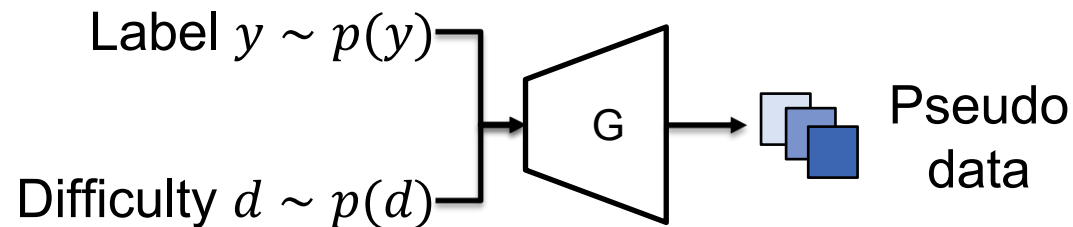
Difficulty $d \sim p(d)$

- **DarkDistill: Difficulty-Conditional Generator**
  - *Create pseudo-data for specific difficulty and label to simulate local datasets, supporting the knowledge distillation process*



- Training objectives
  - Classification: $\mathcal{L}_{ce}(\phi_m, \theta_m) = \mathbb{E}_{\tilde{x} \sim G_m(y,d,\epsilon; \phi_m)} \sum_{i=1}^{m} \text{CE}(\hat{y}, y)$
    - $\hat{y}$ is the predicted label for pseudo data $\tilde{x}$, minimizing $m$ exits loss.
  - Difficulty Simulation: $\mathcal{L}_{dif}(\phi_m, \theta_m) = \mathbb{E}_{\tilde{x} \sim G_m(y,d,\epsilon; \phi_m)} |d - \hat{d}|$
    - Give difficulty $d \sim p(d)$, $\hat{d}$ is the predicted difficulty, minimizing the $|d - \hat{d}|$

# Our Solutions

- **DarkDistill: Difficulty-Aligned Reverse KD**
  - *Model-wise: Transfer knowledge from shallow to deep exits in adjacent layers across varied depth local EENs*

# Our Solutions

- ## DarkDistill: Difficulty-Aligned Reverse KD
  - *Exit-wise: adaptive KD based on difficulty distance between exits across adjacent EENs*



*pseudo data*

**Progressive Reverse KD**

*global model* $\theta_g$

*Model-wise KD*

*Exit-wise KD*

# Our Solutions

- ## DarkDistill-PL: Framework
  - ### Parallel Variant for DarkDistill

*Difficulty-Increased Generator*

*Generate pseudo data with increasing difficulty to simulate the difficulty range across various depth exits*

*Parallel KD*

*Directly transfer the ensemble knowledge in same depth exits across intermediate models into global model*

# Outline

- **Background & Motivation**

- **Problem Statement**

- **Our Solutions**

- **Experiments**

- **Conclusion**

- **Configuration**

  - Dataset: CIFAR100，SVHN，SpeechCommands

  - Settings：100 clients, divided into 4 levels with increasing compute capabilities (4 sizes of local model)

  - Base Model：Deit-tiny (Transformer, 12 layers)

  - Exit distribution：add exits at 3th, 6th, 9th, 12th layer

  - Finetune methods: Full parameters, LoRA

  - Total Epoch：500
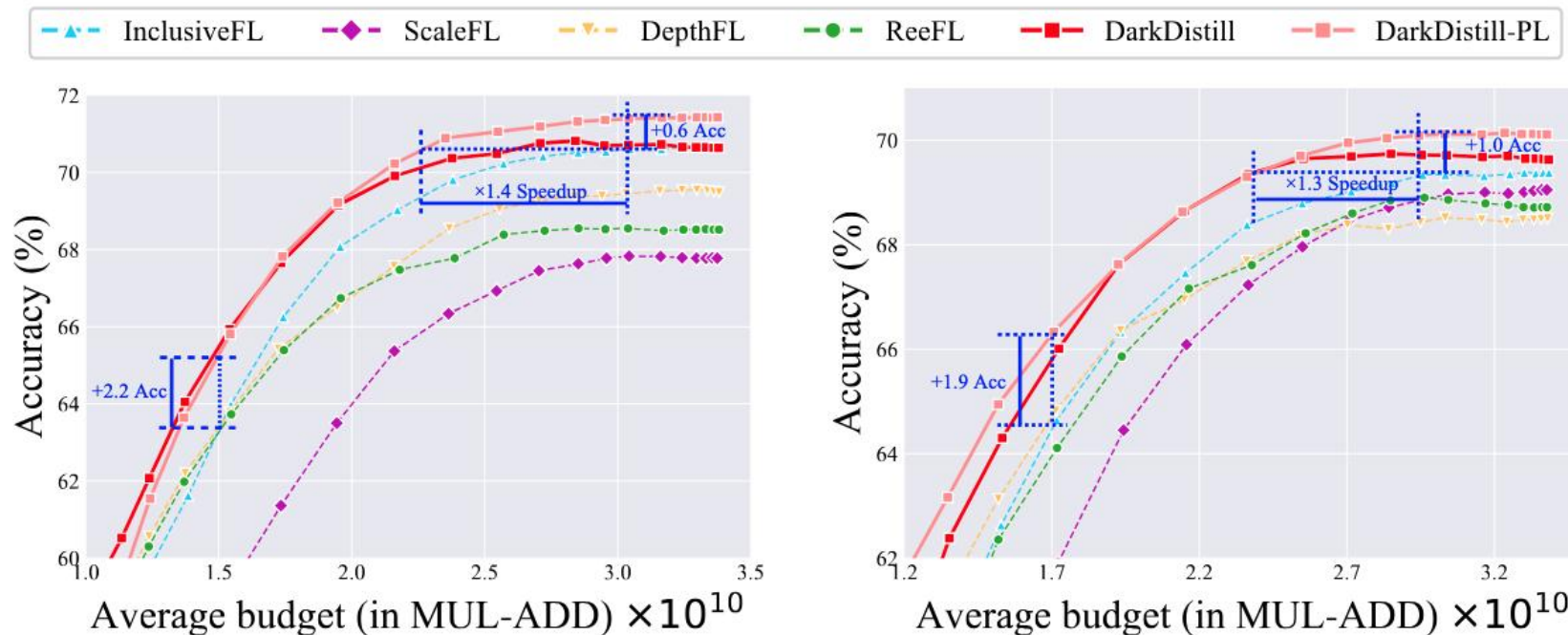
# Experiments-Main Results

- **Performance of Anytime Inference**
  - Measures the accuracy of each exit assuming sufficient budgets
  - DarkDistill and DarkDistill-PL with BoostNet are the top 2 on all datasets, and increase 2 percents in general

| Finetune | Difficulty-aware | Approach | CIFAR-100 [19] | | | SVHN [30] | SpeechCmds [44] |
|---|---|---|---|---|---|---|---|
| | | | $\alpha = 0.1$ | $\alpha = 1$ | $\alpha = 1000$ | | |
| Full | None | ExclusiveFL | $26.60_{\pm 3.10}$ | $49.96_{\pm 11.48}$ | $41.58_{\pm 7.01}$ | $85.28_{\pm 2.97}$ | $87.00_{\pm 2.88}$ |
| | | InclusiveFL [26] | $40.10_{\pm 2.03}$ | $58.83_{\pm 6.98}$ | $61.40_{\pm 7.01}$ | $82.95_{\pm 0.34}$ | $91.90_{\pm 1.42}$ |
| | | ScaleFL [16] | $54.99_{\pm 10.61}$ | $63.21_{\pm 9.14}$ | $63.82_{\pm 9.87}$ | $88.24_{\pm 0.78}$ | $92.56_{\pm 0.26}$ |
| | | DepthFL [18] | $40.70_{\pm 1.57}$ | $59.01_{\pm 5.18}$ | $61.71_{\pm 5.73}$ | $83.45_{\pm 0.43}$ | $92.05_{\pm 0.60}$ |
| | | ReeFL [23] | $59.24_{\pm 8.00}$ | $63.37_{\pm 7.72}$ | $63.90_{\pm 8.68}$ | $88.37_{\pm 1.27}$ | $93.12_{\pm 1.14}$ |
| | BoostNet [45] | ExclusiveFL | $48.68_{\pm 13.66}$ | $57.57_{\pm 15.12}$ | $58.65_{\pm 15.31}$ | $87.30_{\pm 2.89}$ | $91.07_{\pm 2.58}$ |
| | | InclusiveFL [26] | $57.10_{\pm 7.21}$ | $62.96_{\pm 8.12}$ | $64.01_{\pm 8.24}$ | $87.86_{\pm 1.66}$ | $92.91_{\pm 1.10}$ |
| | | ScaleFL [16] | $52.74_{\pm 13.82}$ | $60.55_{\pm 11.93}$ | $60.73_{\pm 10.80}$ | $87.91_{\pm 0.77}$ | $92.03_{\pm 0.37}$ |
| | | DepthFL [18] | $58.15_{\pm 6.73}$ | $63.81_{\pm 6.34}$ | $64.19_{\pm 6.73}$ | $87.74_{\pm 1.01}$ | $92.72_{\pm 0.64}$ |
| | | ReeFL [23] | $59.01_{\pm 7.98}$ | $63.08_{\pm 9.03}$ | $63.66_{\pm 7.31}$ | $88.39_{\pm 1.28}$ | $93.01_{\pm 1.18}$ |
| | | DarkDistill | $60.48_{\pm 7.93}$ | $64.50_{\pm 7.97}$ | $\mathbf{65.67}_{\pm 7.48}$ | $88.41_{\pm 1.46}$ | $93.31_{\pm 1.13}$ |
| | | DarkDistill-PL | $\mathbf{61.05}_{\pm 8.19}$ | $\mathbf{65.12}_{\pm 7.02}$ | $65.49_{\pm 7.88}$ | $\mathbf{88.48}_{\pm 1.57}$ | $\mathbf{93.42}_{\pm 0.98}$ |

- **Performance of Budget Inference**
  - Measures the accuracy of a batch samples within given budgets
  - DarkDistill and DarkDistill-PL can improve the accuracy over the baselines at various computation budgets.



(a) Full $\alpha = 1000$

(b) Full $\alpha = 1$

- **Difficulty Assessment Module**
  - Verify the efficiency of difficulty assessment module
    - Left images are easier predicted by module
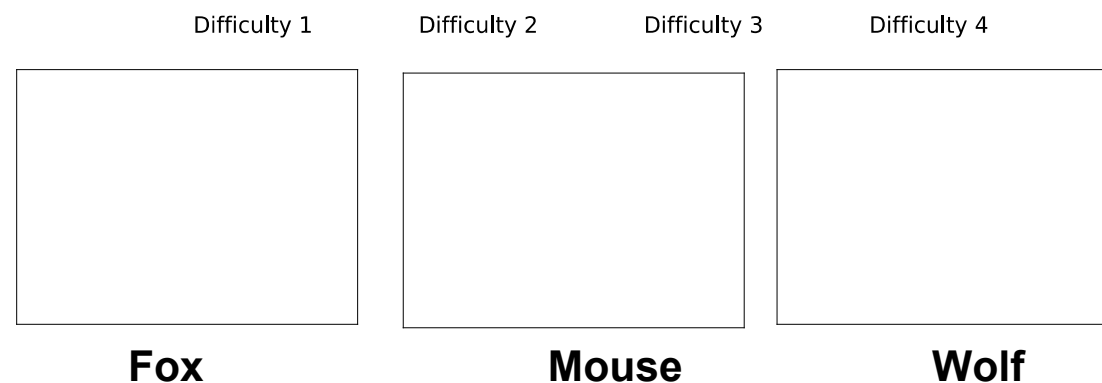    - Right image are more difficult predicted by module



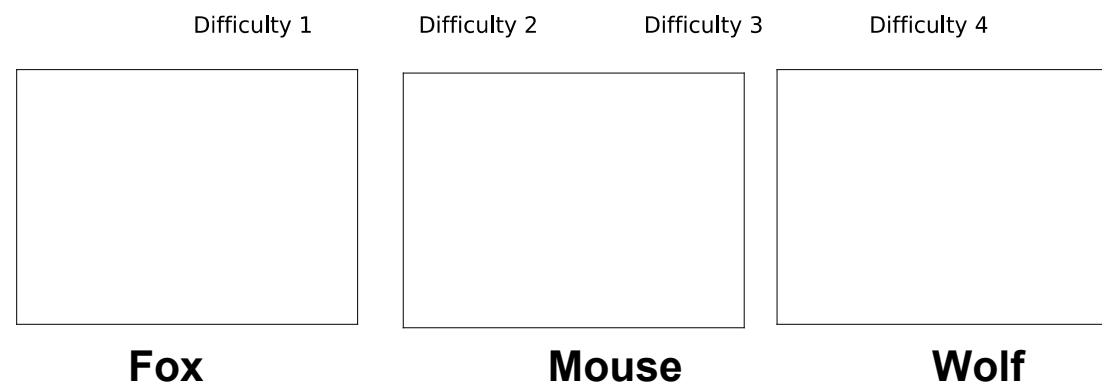**Single pigment, easy to judge**

**Complex contour, difficult to judge**

**Consistent with human intuition, the difficulty assessment module is designed reasonably**

# Experiments-Module Ablation

- **Difficulty-Conditional Generator**
  - Pseudo data of the same category are divided into <span style="color:red">four different levels</span> of clustering clusters

| Difficulty 1 | Difficulty 2 | Difficulty 3 | Difficulty 4 |
|---|---|---|---|

**Fox**          **Mouse**          **Wolf**

- **Difficulty-Conditional Generator**
  - Pseudo data of the same category are divided into four different levels of clustering clusters

|  Difficulty 1 | Difficulty 2 | Difficulty 3 | Difficulty 4 |
|---|---|---|---|

**Fox**          **Mouse**          **Wolf**

  - The robustness of generator architecture

| $d_\epsilon$ | $d_h$ | | | |
|---|---|---|---|---|
|  | 64 | 128 | 256 | 512 |
| 2 | 64.79\|64.88 | 65.05\|64.93 | 65.32\|65.25 | 65.10\|65.11 |
| 16 | 65.38\|64.91 | 65.20\|65.09 | 65.37\|64.65 | 65.18\|65.51 |
| 32 | 65.05\|64.79 | 65.06\|64.79 | 65.28\|65.08 | 65.60\|64.90 |
| 64 | 65.15\|64.92 | 65.74\|64.93 | 64.91\|64.55 | 65.06\|65.05 |
| SOTA | $64.19_{\pm 6.73}$ | | | |

# Outline

- **Background & Motivation**

- **Problem Statement**

- **Our Solutions**

- **Experiments**

- **Conclusion**

# Conclusion

- *This paper introduces DarkDistill, a novel heterogeneous federated learning scheme dedicated for early-exit networks (EENs) and its parallel variant DarkDistill-PL for acceleration.*

- *We identify the <span style="color:red">cross-model exit unalignment</span> problem, an unexplored challenge when extending difficulty-aware EEN training to federated contexts.*

- *We develop a difficulty-conditional generator training strategy and a difficulty-aligned reverse distillation scheme to aggregate EENs of varying depths into a global model that <span style="color:red">retains its difficulty-specific specialization</span>.*

# THANK YOU

if you have problems, feel free to email

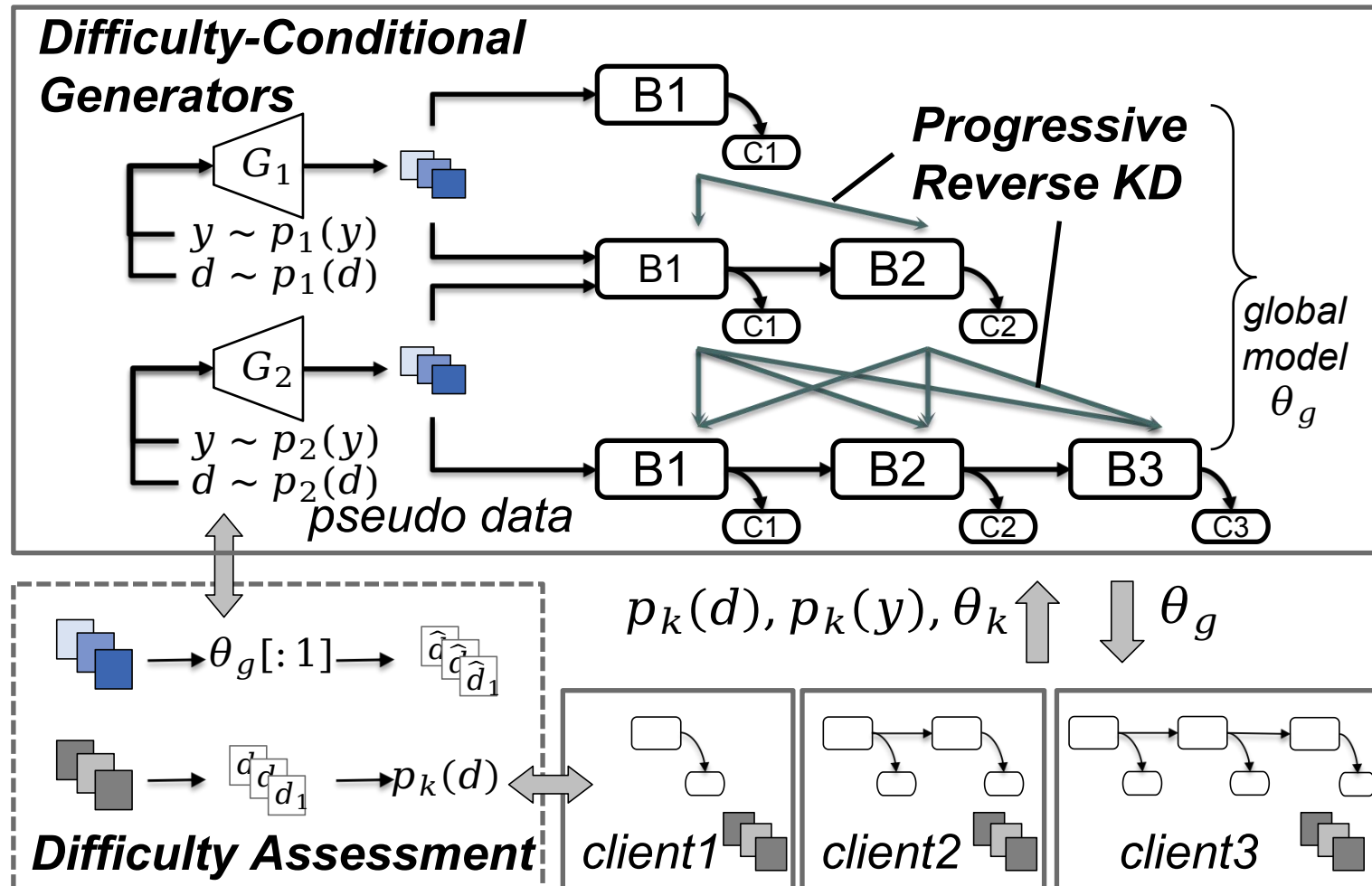lehaoqv@buaa.edu.cn

or talk with me at Poster 201, 5th August

# Our Solutions

- ## DarkDistill: Framework



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

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

- ## DarkDistill: Difficulty-Aligined Reverse Knowledge Distillation



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