# Weak-to-Strong Generalization Through the Data-Centric Lens

Changho Shin, John Cooper, Frederic Sala

University of Wisconsin–Madison
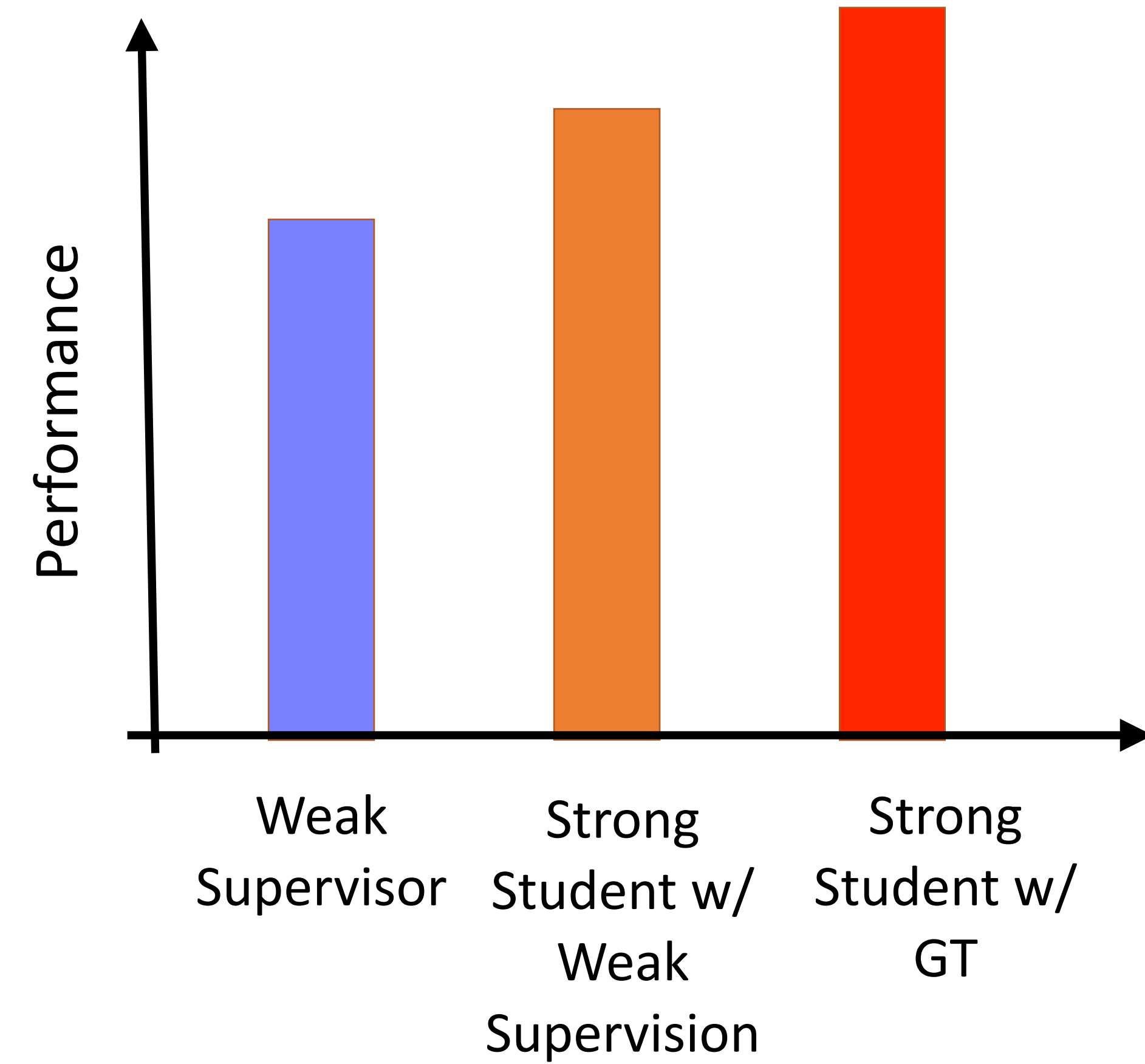
Paper Link

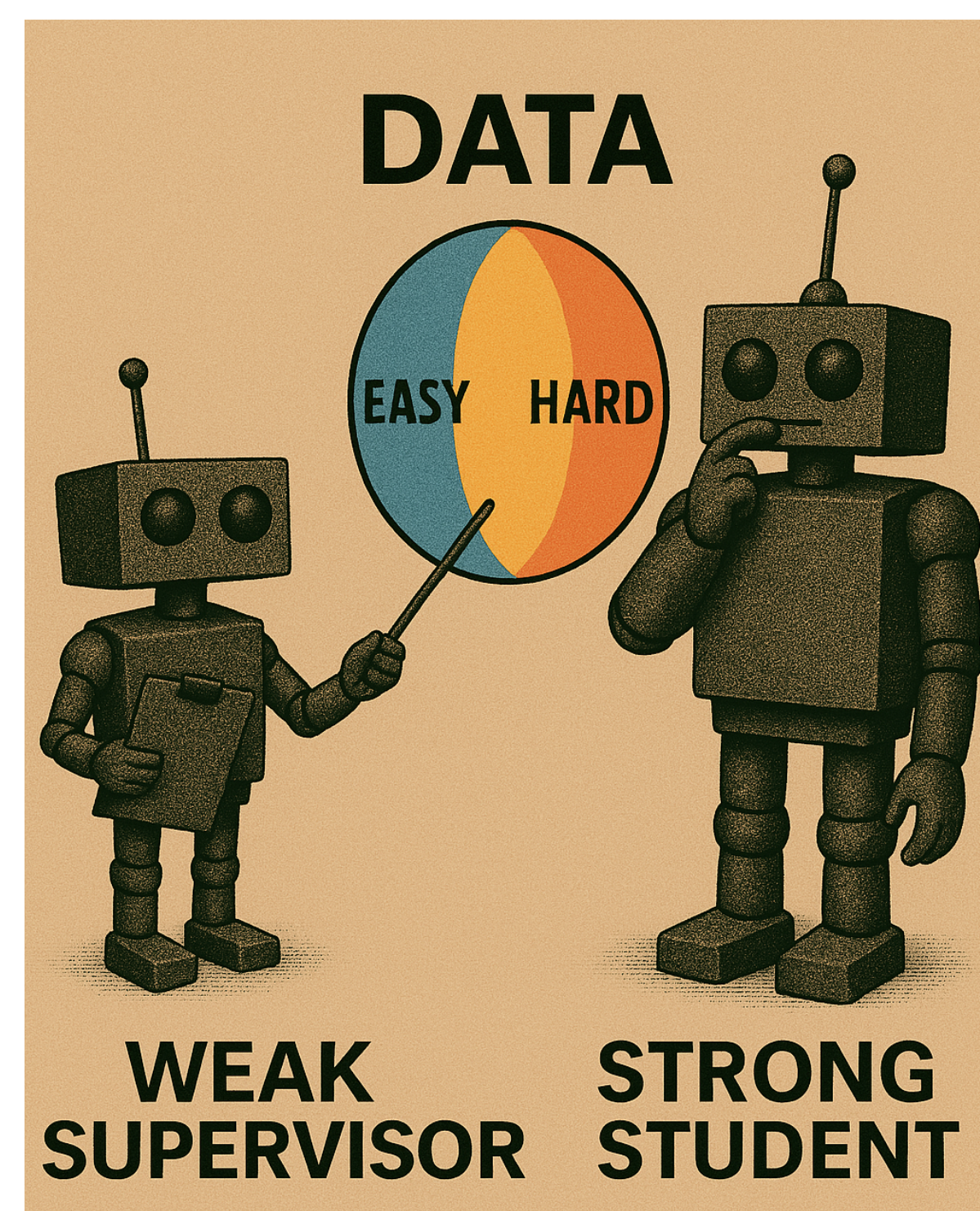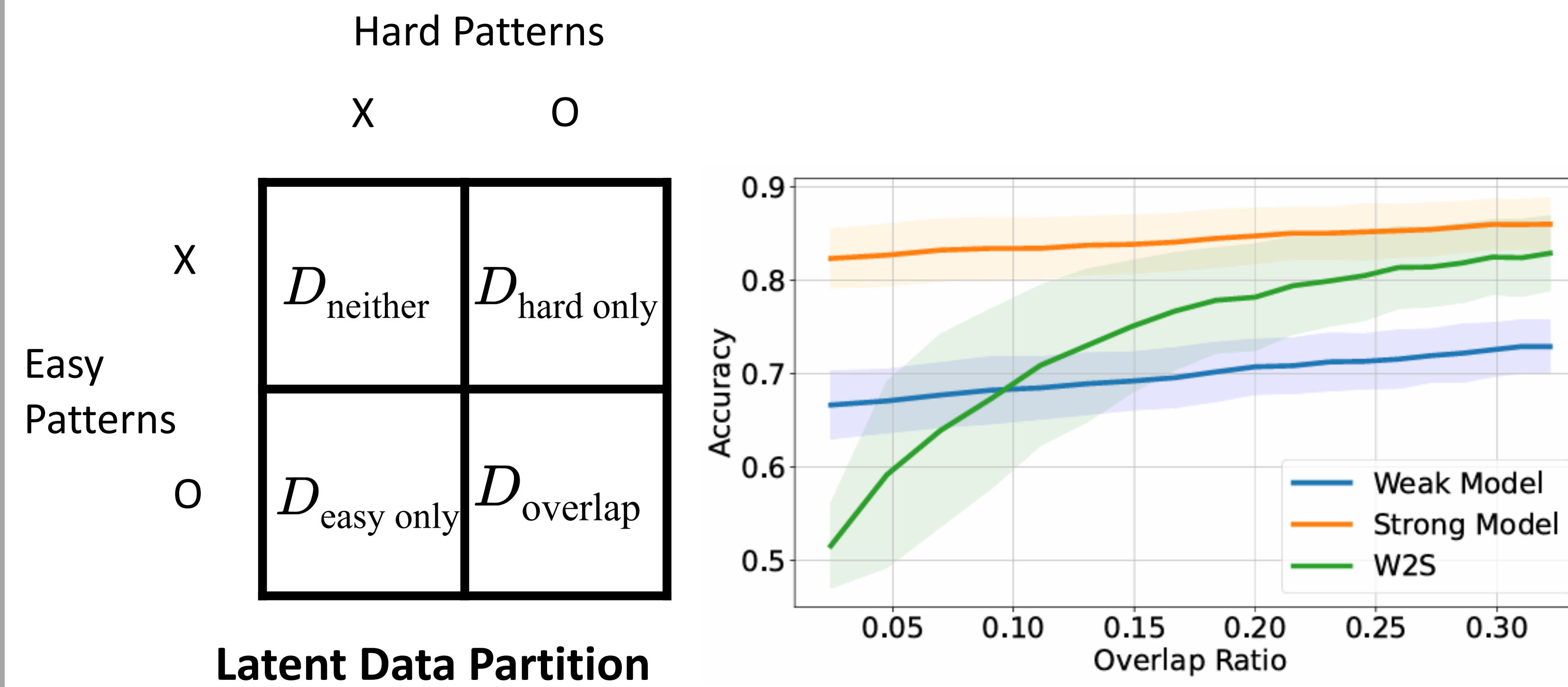## Background

### Weak-to-Strong Generalization



Performance

Weak Supervisor | Strong Student w/ Weak Supervision | Strong Student w/ GT

❓ How can a strong student outperform a weak supervisor?

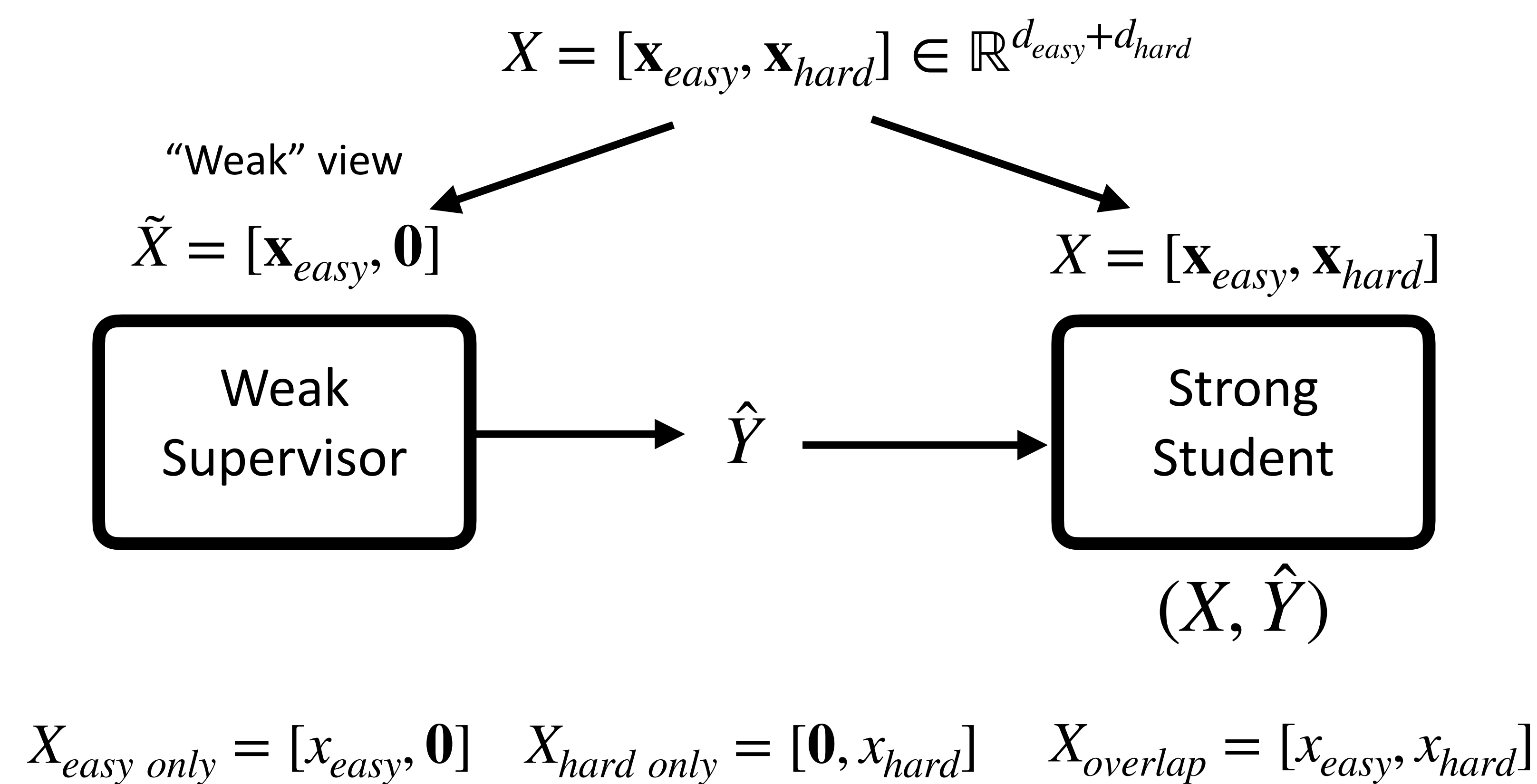💡 We propose a "data mechanism" that drives weak-to-strong generalization.



DATA
EASY   HARD

WEAK SUPERVISOR    STRONG STUDENT

## Overlap Mechanism

Hard Patterns
X        O

|            | X | O |
|------------|---|---|
| Easy Patterns X | $D_{\text{neither}}$ | $D_{\text{hard only}}$ |
| O          | $D_{\text{easy only}}$ | $D_{\text{overlap}}$ |

**Latent Data Partition**



Accuracy vs Overlap Ratio — Weak Model, Strong Model, W2S
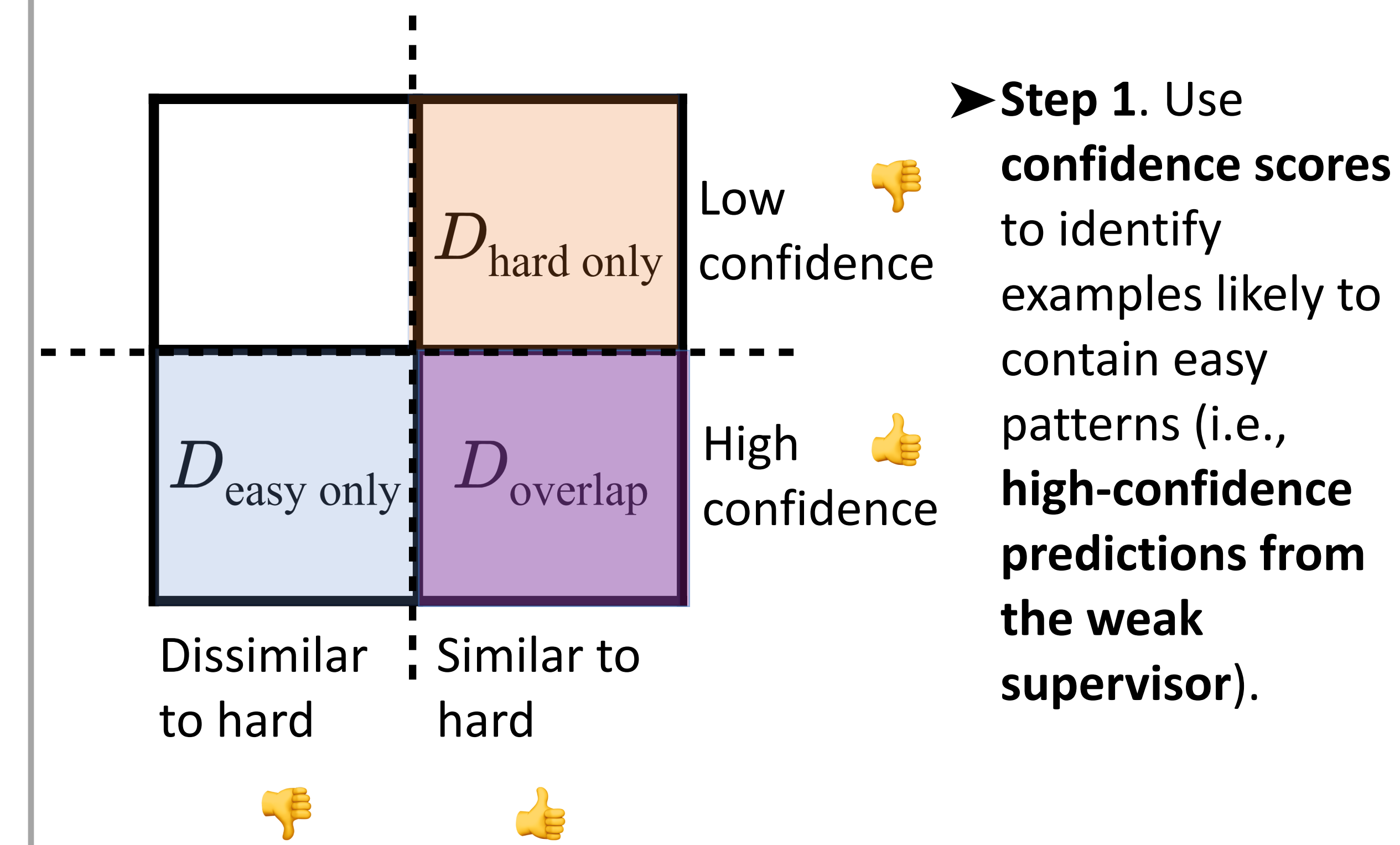
Overlapping points enable weak-to-strong generalization
1) **Weak supervisor** can **provide accurate pseudo-labels** based on **easy patterns**
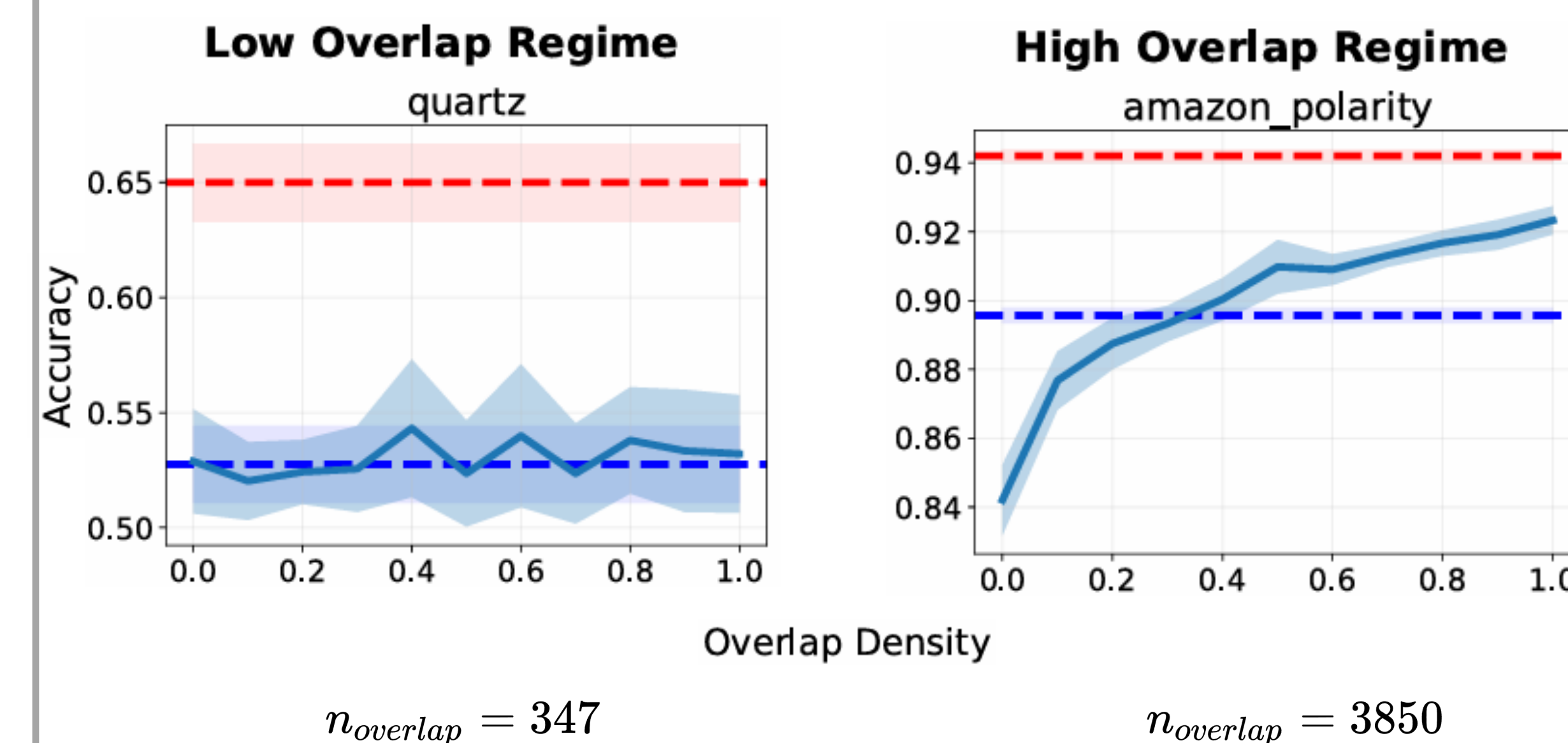2) **Strong student** can learn **hard patterns** using accurate pseudo-labels

## Data Model

$$X = [\mathbf{x}_{easy}, \mathbf{x}_{hard}] \in \mathbb{R}^{d_{easy}+d_{hard}}$$

"Weak" view

$$\tilde{X} = [\mathbf{x}_{easy}, \mathbf{0}] \qquad X = [\mathbf{x}_{easy}, \mathbf{x}_{hard}]$$

Weak Supervisor → $\hat{Y}$ → Strong Student

$(X, \hat{Y})$

$$X_{easy\ only} = [x_{easy}, \mathbf{0}] \quad X_{hard\ only} = [\mathbf{0}, x_{hard}] \quad X_{overlap} = [x_{easy}, x_{hard}]$$

## Overlap Detection

How can we find overlapping points in real-world data?



| | $D_{\text{hard only}}$ | Low confidence |
| $D_{\text{easy only}}$ | $D_{\text{overlap}}$ | High confidence |

Dissimilar to hard | Similar to hard

➤ **Step 1.** Use **confidence scores** to identify examples likely to contain easy patterns (i.e., **high-confidence predictions from the weak supervisor**).
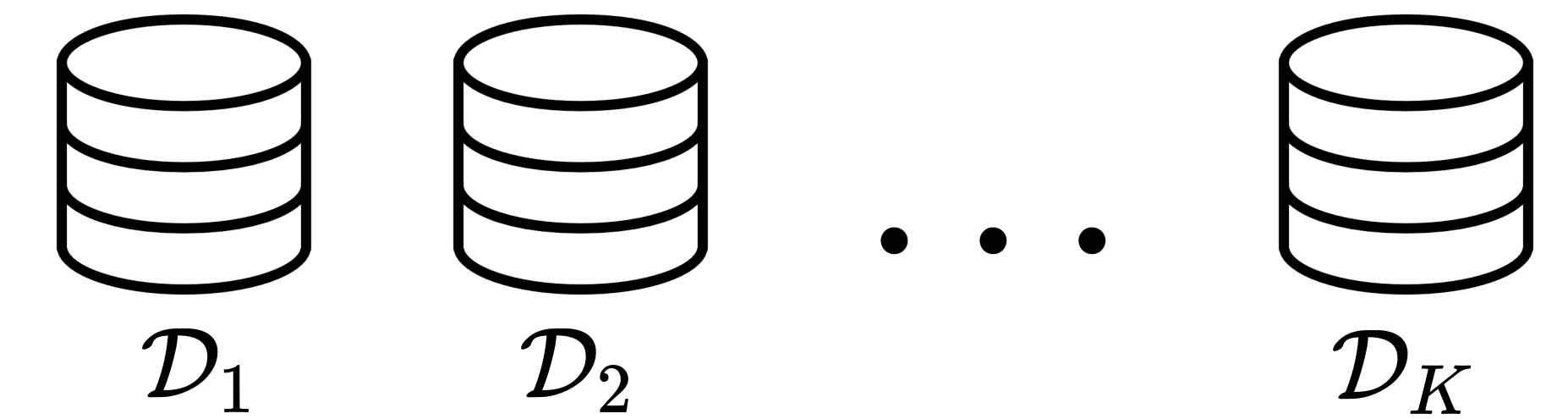
➤ **Step 2.** Among the easy-pattern data, apply similarity-based filtering to **isolate examples that also resemble hard patterns**.

### The extent of overlap is predictive of W2S generalization!



Low Overlap Regime — quartz
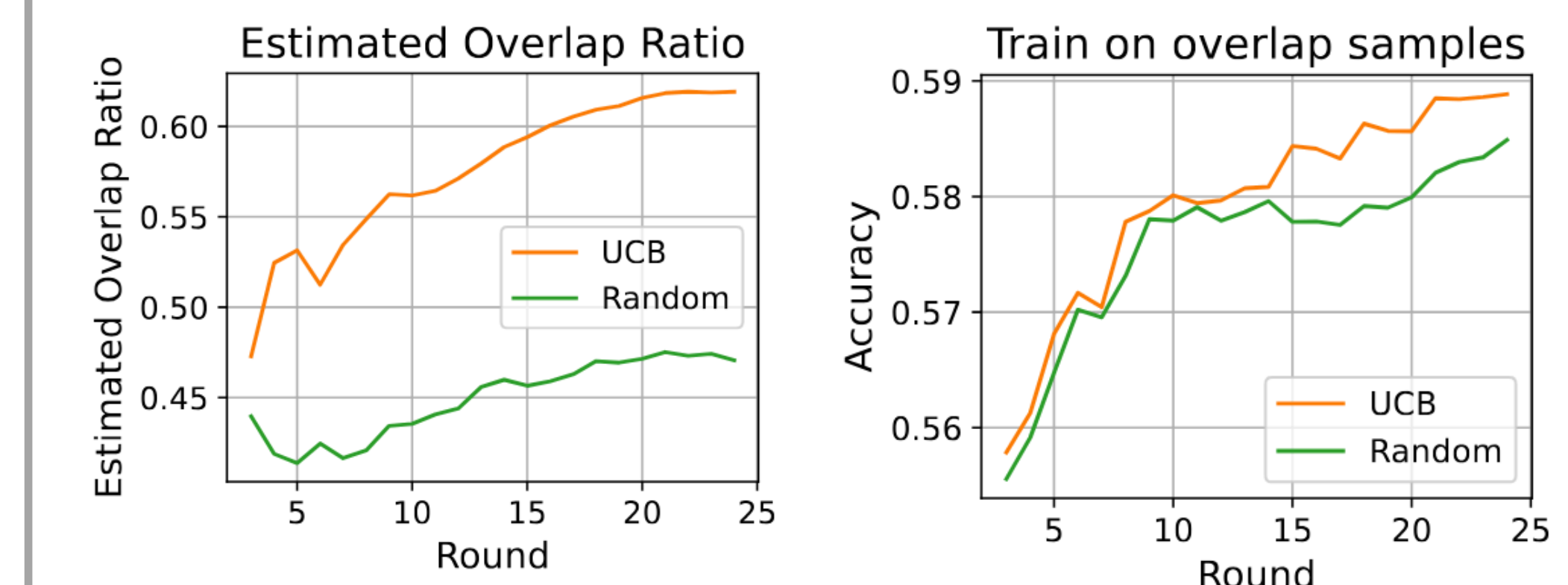$n_{overlap} = 347$

High Overlap Regime — amazon_polarity
$n_{overlap} = 3850$

## Data Source Selection

**Setup:** Sample data from multiple sources under a limited budget.



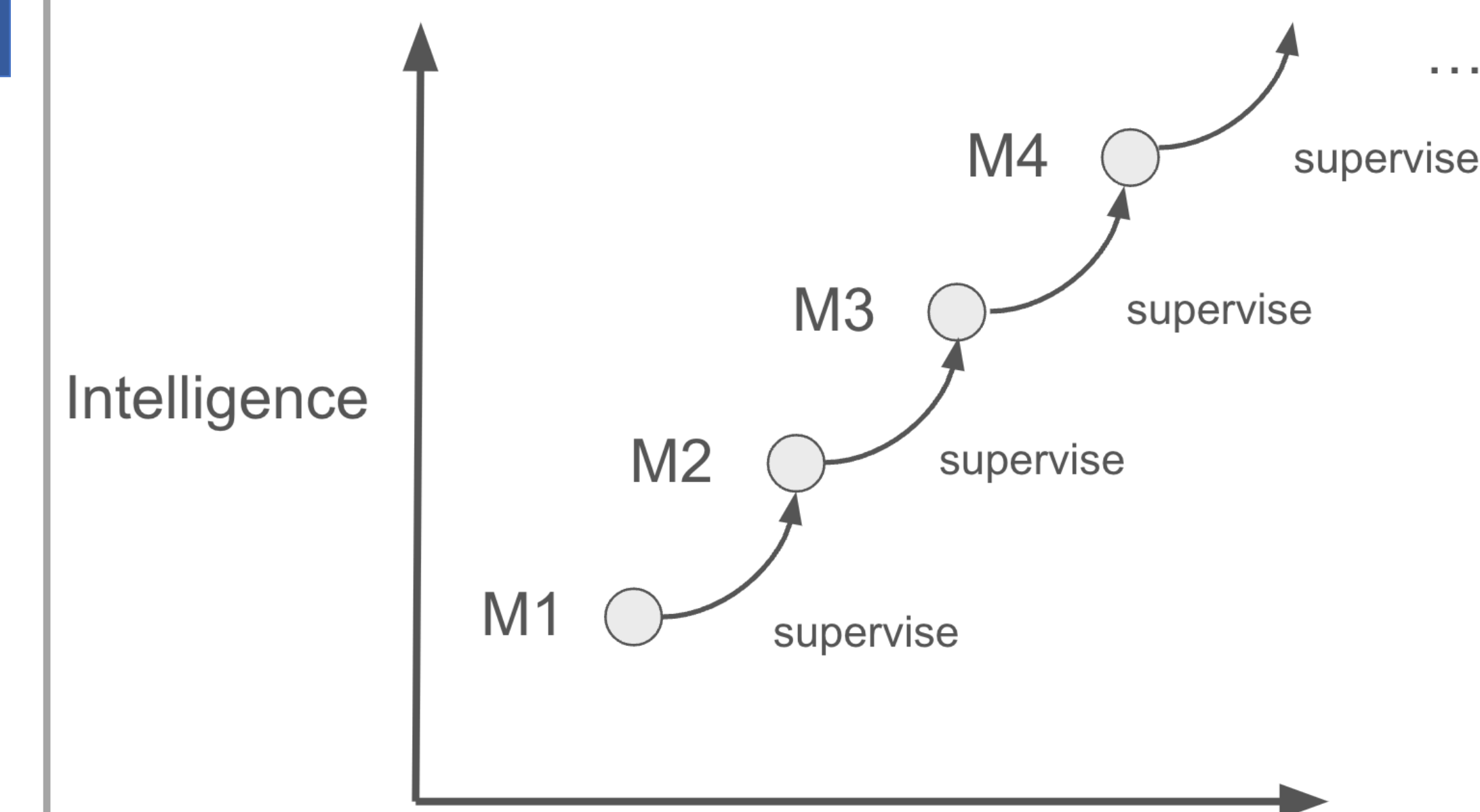$\mathcal{D}_1 \quad \mathcal{D}_2 \quad \cdots \quad \mathcal{D}_K$

**Goal:** Maximize the overlap ratio.
**Method:** Use an Upper Confidence Bound (UCB) strategy to prioritize sources likely to yield high overlap.
**Result:** UCB increases overlap —> improves weak-to-strong generalization



Estimated Overlap Ratio — UCB, Random
Train on overlap samples — UCB, Random

## Moving Forward…



Intelligence

M1 → supervise → M2 → supervise → M3 → supervise → M4 → supervise → …

- Generative tasks (e.g., reasoning)
- Synthetic data generation for W2S