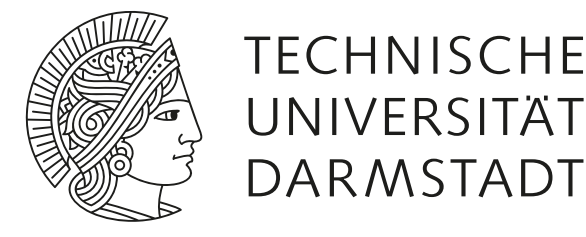


# Deep Classifier Mimicry without Data Access



Steven Braun<sup>1</sup>  
<sup>1</sup>Dep. of CS, TU Darmstadt

Martin Mundt<sup>1,2</sup>  
<sup>2</sup>hessian.AI <sup>3</sup>DFKI

Kristian Kersting<sup>1,2,3,4</sup>  
<sup>4</sup>Centre for Cognitive Science, TU Darmstadt



## Overview

### Problem:

- Access to pre-trained models is common, allowing for knowledge distillation in downstream tasks.
- But: Original training data is often unavailable, challenging distillation methods that depend on it.

**Research Question:** Can we distill knowledge from models without original training data access independent of their architecture?

### Solution: 🍷 CAKE

- Model-agnostic knowledge distillation procedure without the need for original training data.
- Contrastively diffuses synthetic samples along the decision boundary at different scales.

## Knowledge Distillation

**Goal:** Teach a student to predict like a teacher model.

- **Teacher  $f^T$ :** Usually larger model, pre-trained on data, provides soft targets.
- **Student  $f^S$ :** Learns to imitate teacher's predictions.
- **Distillation Loss:** Combines true label loss with teacher-student output similarity.

$$\mathcal{L}_{\text{KD}} = \lambda_{\text{true}} \text{CE}(y, f^S(x)) + \lambda_{\text{soft}} \text{CE}(f^T(x), f^S(x))$$

↪ requires original training data  $x$  access!

What if we only have access to the pre-trained model?

Paper

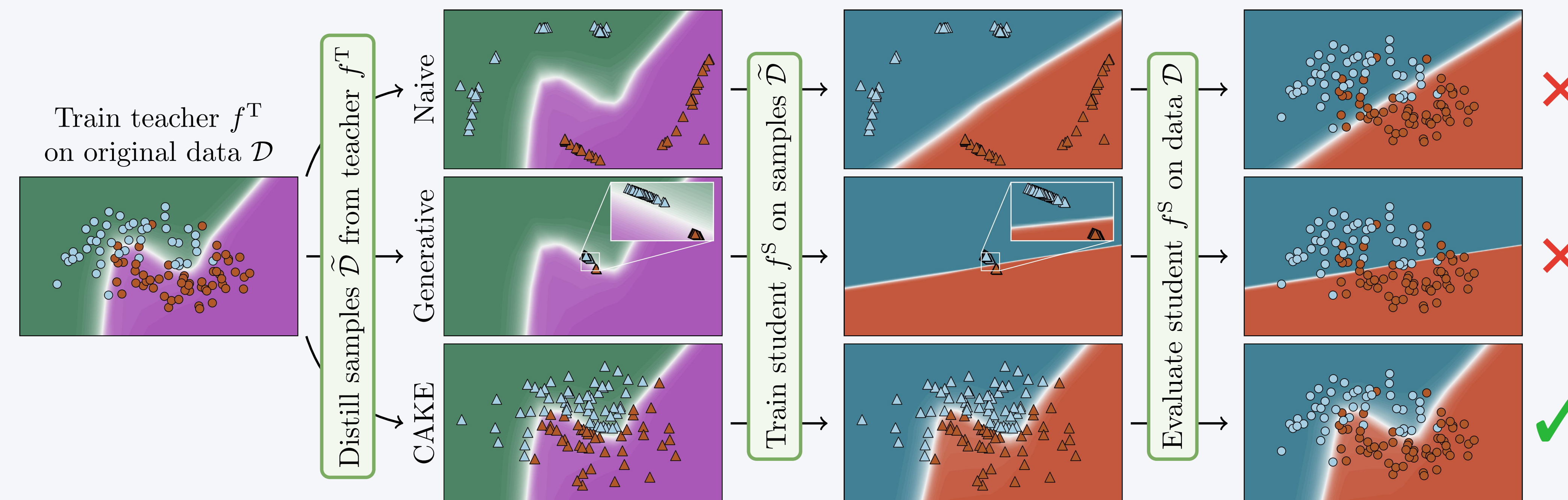


Code



## Finding the Decision Boundary without Original Training Data

Single-instance KD data synthetization methods usually ignore inter-class relationships and are prone to collapse at locally optimal regions that are suboptimal for student model training.



### Challenges

- “Naive”  $\text{CE}(f^T(x), y)$  pushes samples to the correct boundary sides but too far away
- “Generative” (e.g., GAN-based) is prone to collapsing its model parameters to single modes

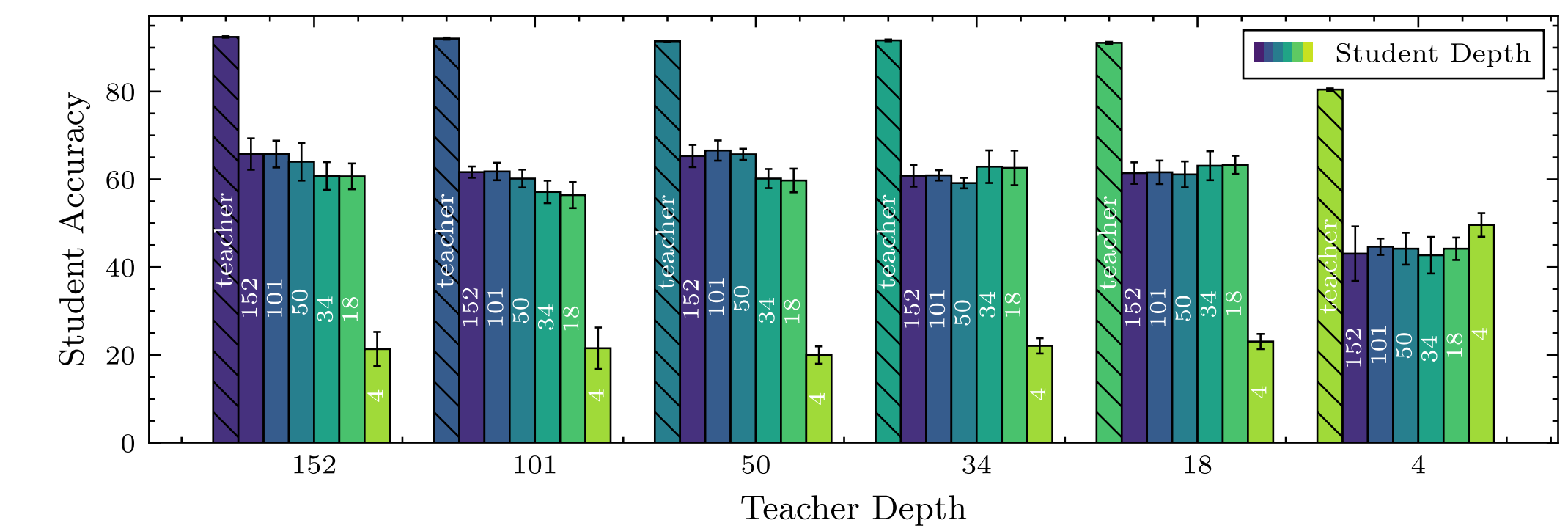
↪ we need to find a balance between pushing samples away from the border while keeping pairs from diff. classes close to each other and noisily scatter these pairs along the boundary at different scales.

## CAKE: Contrastive Abductive Knowledge Extraction

- We generate  $M$  mini-batches  $\tilde{D}_m$  of  $\frac{N}{M}$  synthetic samples  $(x_i^{t=0}, y_i)$  from chosen priors  $p(x), p(y)$ .
  - For each  $\tilde{D}_m$ , we perform  $T$  gradient descent steps to minimize a weighted objective based on:
    - **Classification Loss:**  $\mathcal{L}_{\text{cls}}(x_i^t, y_i) = \text{CE}(y_i, p(f^T(x_i^t)))$ 
      - pushes samples  $x_i^t$  to the correct decision boundary sides according to the sampled class  $y_i$
    - **Contrastive Loss:**  $\mathcal{L}_{\text{contr}}(x_i^t, x_j^t) = \mathbb{1}[y_i \neq y_j] \|f^T(x_i^t) - f^T(x_j^t)\|_2^2$ 
      - pulls pairs of samples from diff. classes together, for  $C$  classes we get  $C(C-1)$  forces/sample
    - **Domain Knowledge Loss** (e.g., Total Var. for images):  $\mathcal{L}_{\text{TV}}(x) = \sum_{j,k} \|x_{j,k} - x_{j-1,k}\| + \|x_{j,k} - x_{j,k-1}\|$ 
      - enables injection of meta-knowledge to constrain the space of relevant samples further
  - **Noise injection:** to push samples *along* the decision boundary, we need to additionally disperse them
- Explicit* (LAKE): Langevin Dynamics  $x_i^{t+1} = x_i^t + \nabla_x \mathcal{L}(x_i^t) \eta(t) + \sqrt{2\eta(t)} \varepsilon_i^t$  with  $\varepsilon_i^t \sim N(0, I)$ .
- Implicit* (CAKE): Stochasticity of SGD and step size schedule  $\eta(t)$  scatters samples across the boundary.

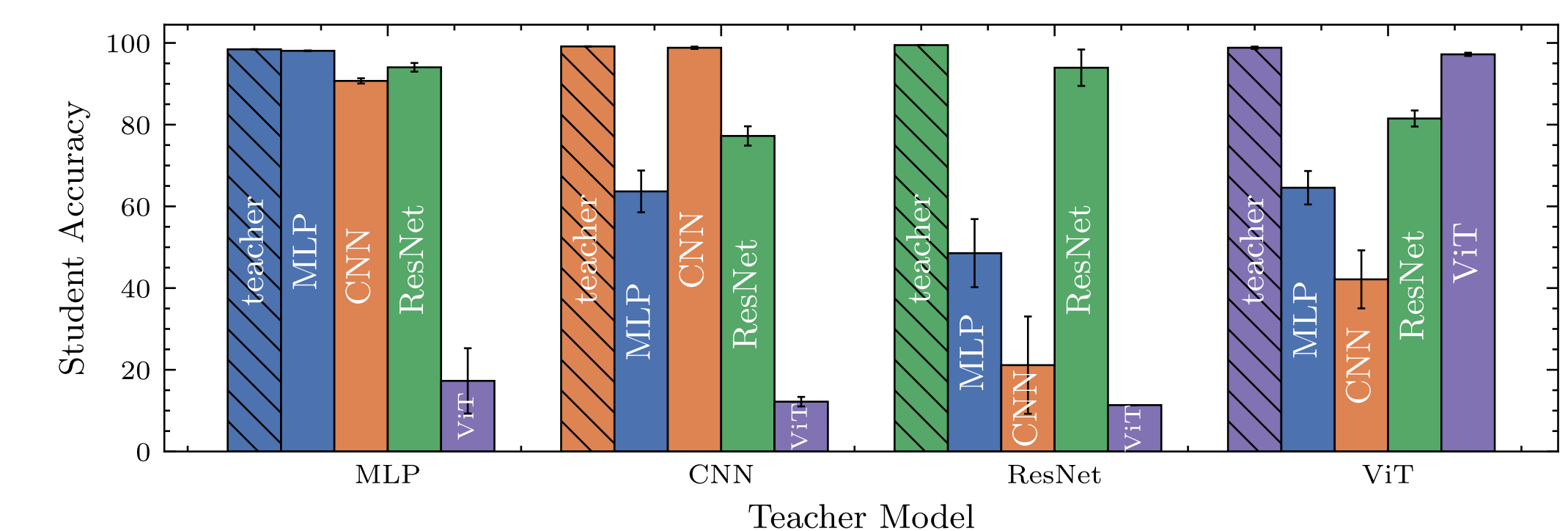
↪ CAKE operates data-free, requiring no access to original training data, and only needs the teacher model to be differentiable, making it completely model-agnostic.

## CAKE across Scales



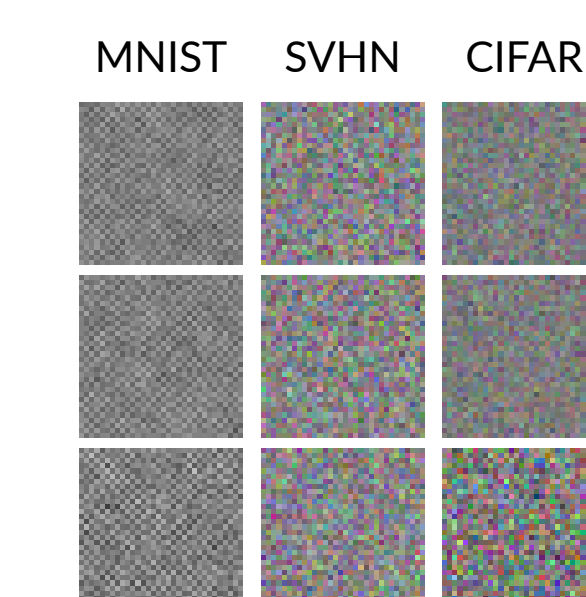
CAKE extracts knowledge and maintains high student accuracy across different model sizes of teachers and students.

## CAKE across Models



CAKE enables knowledge transfer across model types with high student accuracy, esp. when matching model types.

## Generated Samples



Synthesized samples capture decision boundaries without resembling real data. Possible future work: differential privacy, data utility and privacy trade-offs, robustness against adversarial attacks?

## CAKE vs. Others

Method	DF	MADataset	Teacher	Acc. Student	Acc.
KD	×	✓	MNIST LeNet-5	99.3 LeNet-5-Half	98.8
	×	✓	FMNIST LeNet-5	90.8 LeNet-5-Half	89.7
	×	✓	CIFAR-10 ResNet-34	95.6 ResNet-18	94.3
DAFL	✓	×	MNIST LeNet-5	97.9 LeNet-5-Half	97.6
DI	✓	×	CIFAR-10 ResNet-34	93.7 ResNet-18	90.4
ADI	✓	×	CIFAR-10 ResNet-34	95.4 ResNet-18	91.4
DD	✓	✓	CIFAR-10 ResNet-34	95.4 ResNet-18	93.3
ZSDB3KD	✓	✓	MNIST LeNet-5	99.3 LeNet-5-Half	96.5
	✓	✓	FMNIST LeNet-5	91.6 LeNet-5-Half	72.3
CAKE	✓	✓	CIFAR-10 AlexNet	79.3 AlexNet-Half	59.5
	✓	✓	MNIST LeNet-5	99.3 ± 0.12 LeNet-5-Half	98.4 ± 0.18
	✓	✓	FMNIST LeNet-5	91.0 ± 0.12 LeNet-5-Half	76.5 ± 1.01
	✓	✓	SVHN LeNet-5	89.8 ± 0.38 LeNet-5-Half	62.9 ± 4.17
	✓	✓	SVHN ViT-8	94.4 ± 0.13 ViT-4	83.7 ± 4.77
	✓	✓	SVHN ResNet-34	96.1 ± 0.08 ResNet-18	94.2 ± 0.54
	✓	✓	CIFAR-10 ViT-8	73.2 ± 0.76 ViT-4	53.8 ± 5.63
✓	✓	CIFAR-10 ResNet-34	91.8 ± 0.11 ResNet-18	78.9 ± 2.59	