

Deep Classifier Mimicry without Data Access



Steven Braun¹



Martin Mundt^{1,2}



Kristian Kersting^{1,2,3,4}



TECHNISCHE
UNIVERSITÄT
DARMSTADT

¹Department of Computer Science, TU Darmstadt

²Hessian Center for AI (hessian.AI)

³German Research Center for Artificial Intelligence (DFKI)

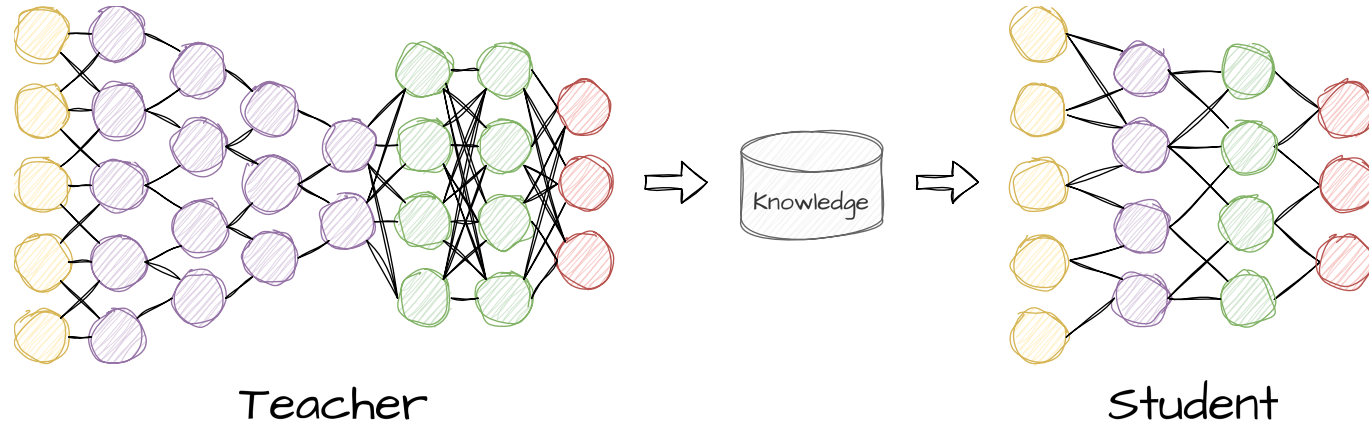
⁴Centre for Cognitive Science, TU Darmstadt



Komp**A+KI**

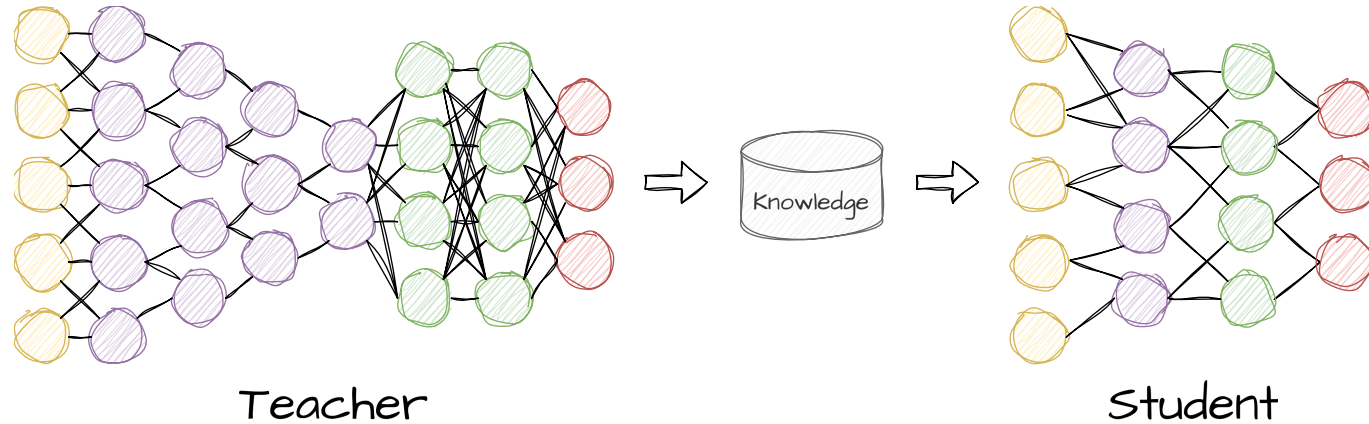


What is Knowledge Distillation?



¹Hinton, G.E., Vinyals, O., & Dean, J. (2015). Distilling the Knowledge in a Neural Network. ArXiv, abs/1503.02531.

What is Knowledge Distillation?

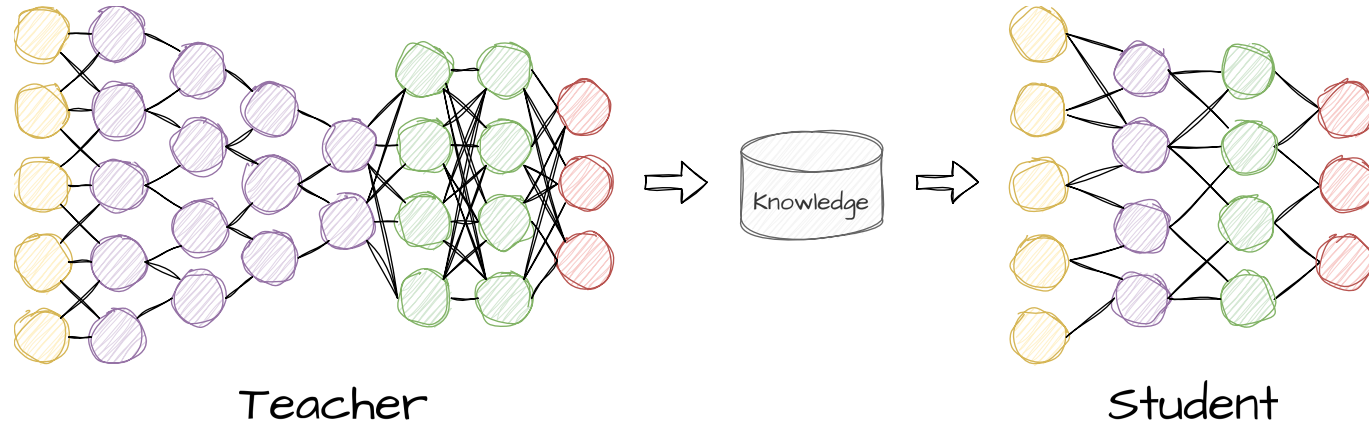


Original Formulation by Hinton et al.¹: Teacher f^T , student f^S

$$\mathcal{L}_{\text{KD}}(\mathbf{x}, y) = \underbrace{\lambda \mathcal{L}_{\text{hard}}(f^S(\mathbf{x}), y)}_{\text{match data}} + (1 - \lambda) \underbrace{\mathcal{L}_{\text{soft}}(f^S(\mathbf{x}), f^T(\mathbf{x}))}_{\text{match teacher}}$$

¹Hinton, G.E., Vinyals, O., & Dean, J. (2015). Distilling the Knowledge in a Neural Network. ArXiv, abs/1503.02531.

What is Knowledge Distillation?



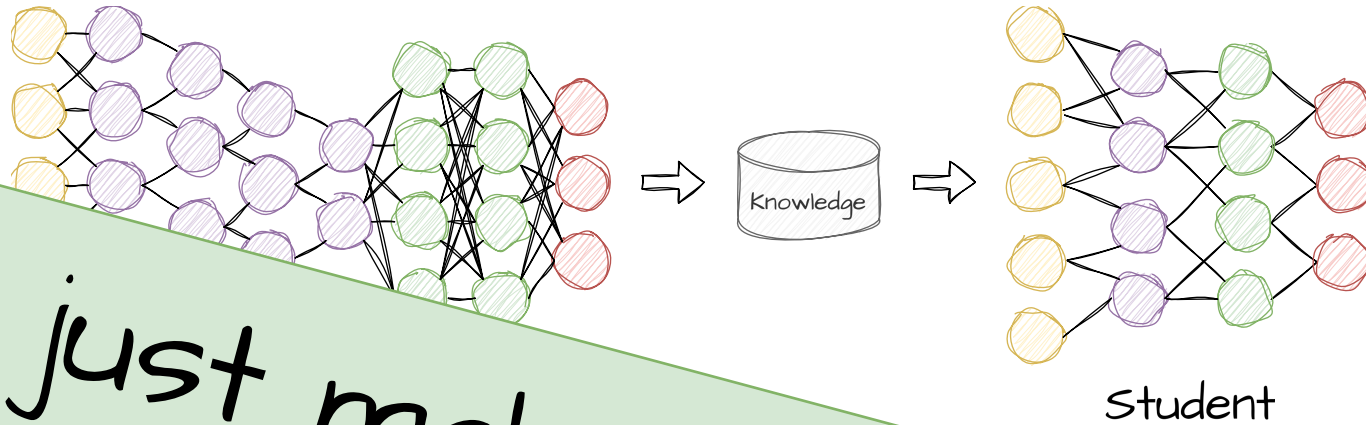
Original Formulation by Hinton et al.¹: Teacher f^T , student f^S

$$\mathcal{L}_{\text{KD}}(\mathbf{x}, y) = \underbrace{\lambda \mathcal{L}_{\text{hard}}(f^S(\mathbf{x}), y)}_{\text{match data}} + (1 - \lambda) \underbrace{\mathcal{L}_{\text{soft}}(f^S(\mathbf{x}), f^T(\mathbf{x}))}_{\text{match teacher}}$$

↪ What if we don't have access to the **original training data** $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$?

¹Hinton, G.E., Vinyals, O., & Dean, J. (2015). Distilling the Knowledge in a Neural Network. ArXiv, abs/1503.02531.

What is Knowledge Distillation?



We just make our own data!

Original Formulation by Hinton et al.

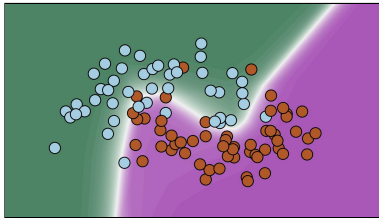
$$\mathcal{L}_{\text{KD}}(x, y) = \lambda \underbrace{\mathcal{L}_{\text{hard}}(f^S(x), y)}_{\text{match data}}$$

↪ What if we don't have access to the original training data $\mathcal{D} = \{(x_i, y_i)\}$

¹Hinton, G.E., Vinyals, O., & Dean, J. (2015). Distilling the Knowledge in a Neural Network. ArXiv, abs/1503.02531.

Extracting Knowledge

Train Teacher

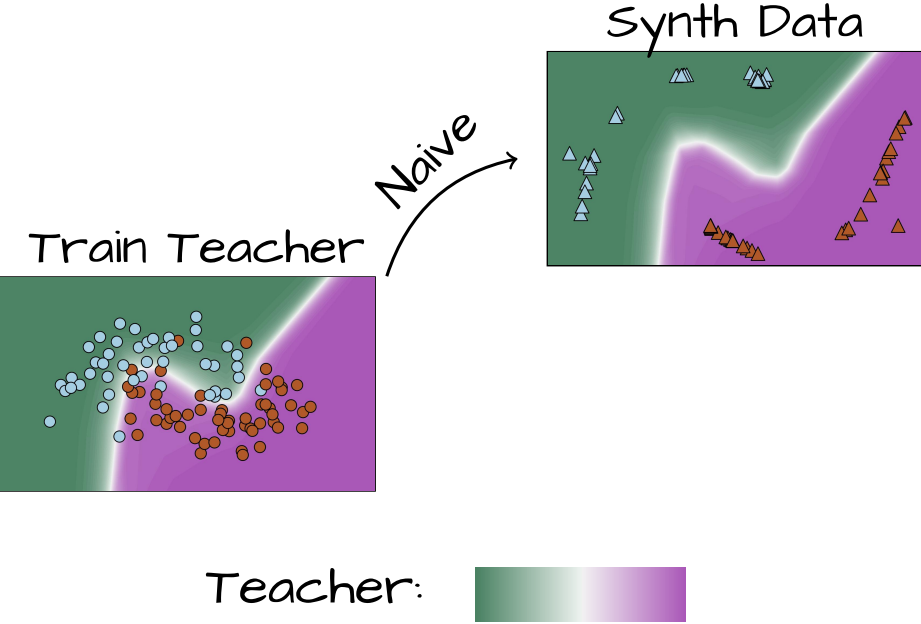


Teacher:



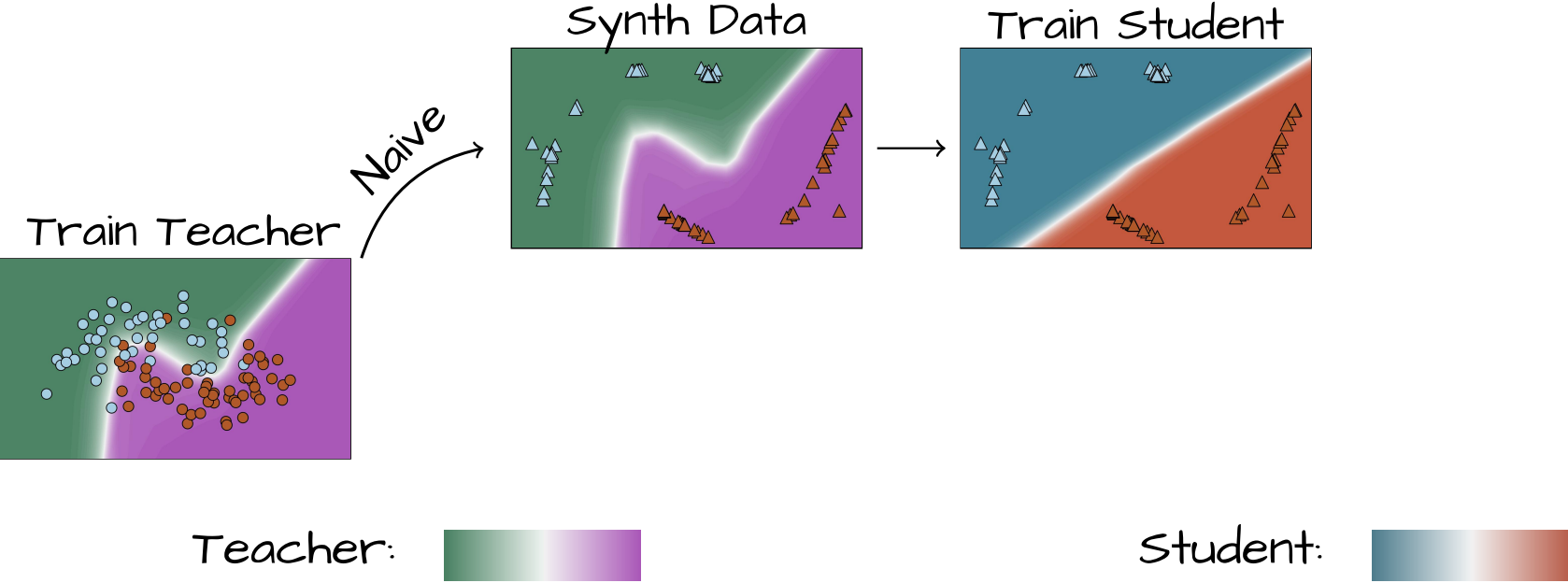
Extracting Knowledge

- Naive: Init. random datapoints (\tilde{x}, \tilde{y}) and minimize $\mathcal{L}(\tilde{x}, \tilde{y}) = \text{CE}(f^T(\tilde{x}), \tilde{y})$



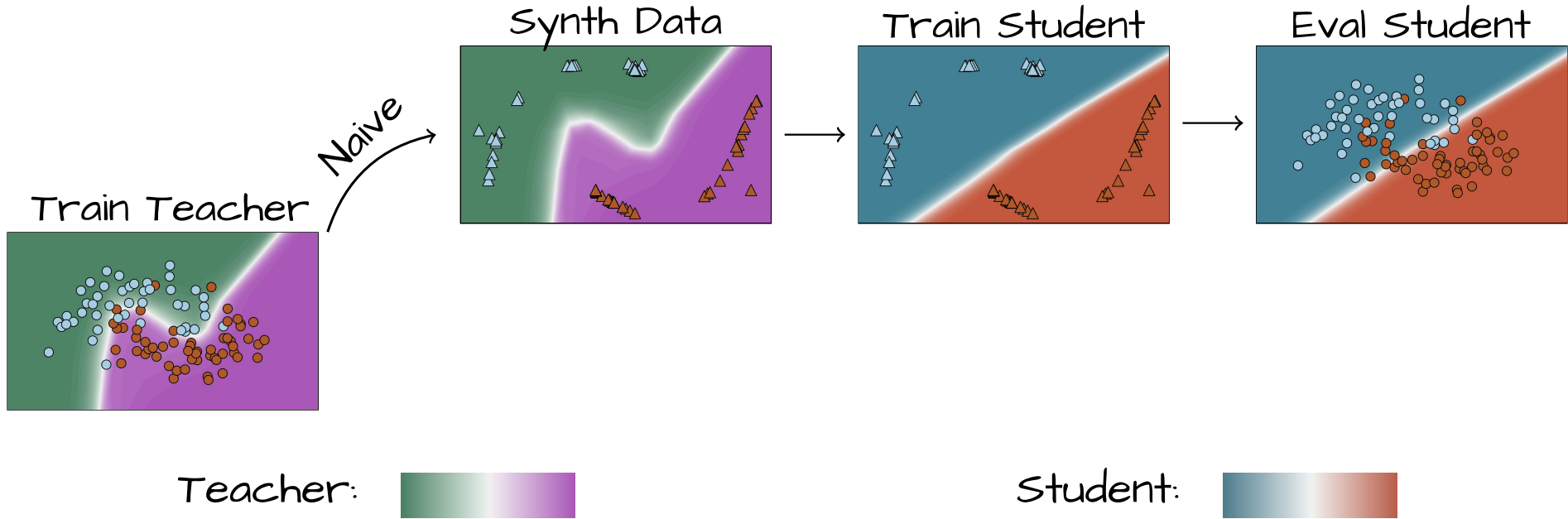
Extracting Knowledge

- Naive: Init. random datapoints (\tilde{x}, \tilde{y}) and minimize $\mathcal{L}(\tilde{x}, \tilde{y}) = \text{CE}(f^T(\tilde{x}), \tilde{y})$



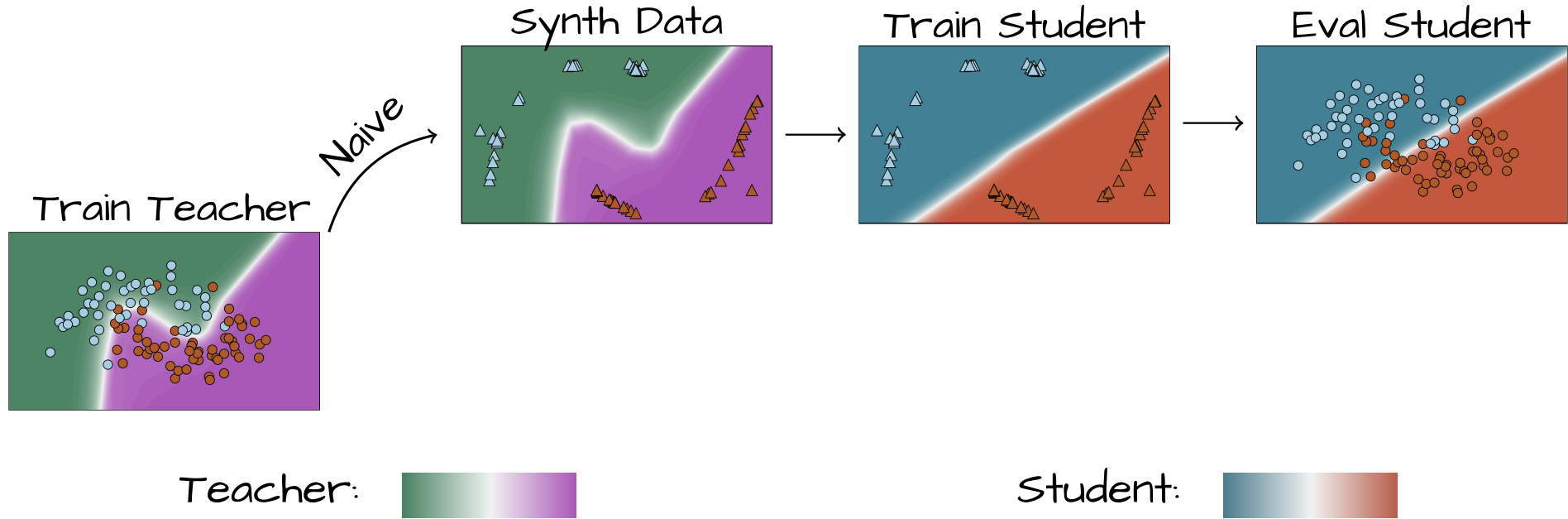
Extracting Knowledge

- Naive: Init. random datapoints (\tilde{x}, \tilde{y}) and minimize $\mathcal{L}(\tilde{x}, \tilde{y}) = \text{CE}(f^T(\tilde{x}), \tilde{y})$



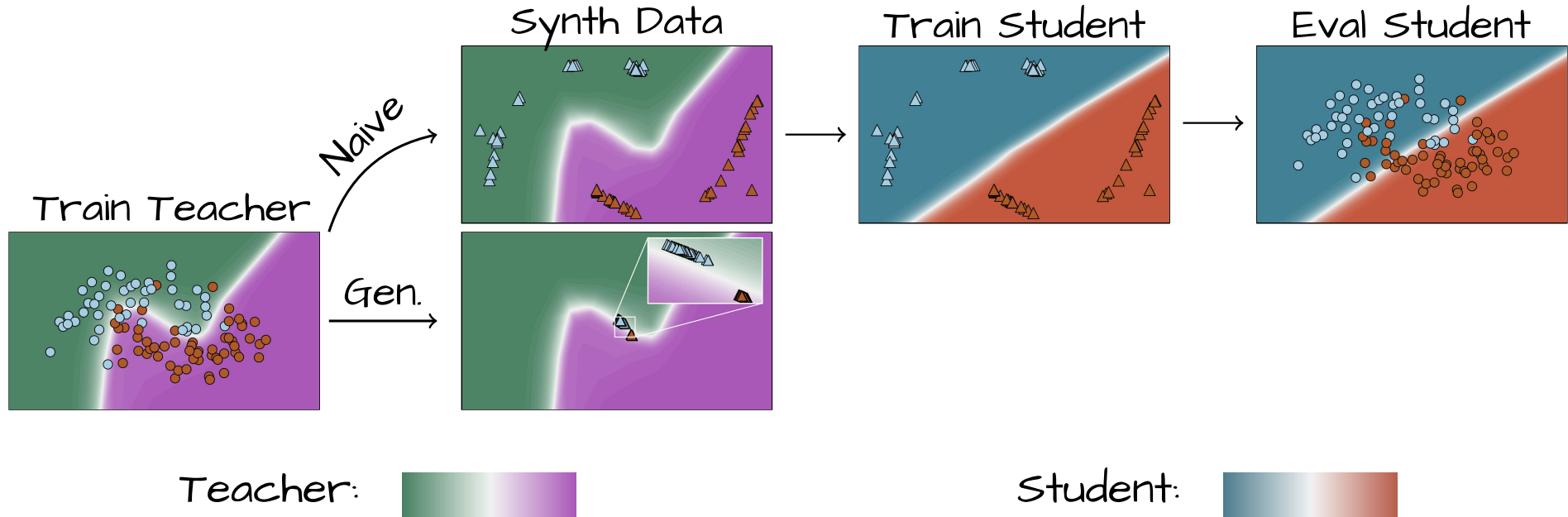
Extracting Knowledge

- **Generative:** Init. random latents (\tilde{z}, \tilde{y}) and minimize $\mathcal{L}(g_{\theta}(\tilde{z}), \tilde{y}) = \text{CE}(f^T(g_{\theta}(\tilde{z})), \tilde{y})$



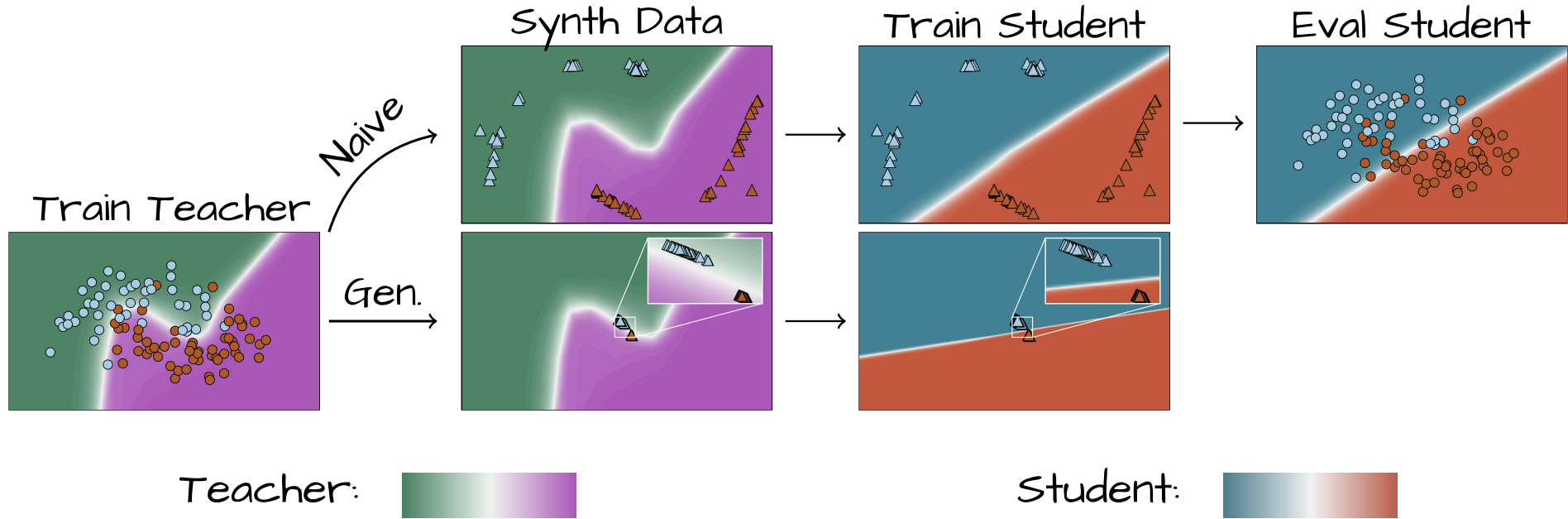
Extracting Knowledge

- **Generative:** Init. random latents (\tilde{z}, \tilde{y}) and minimize $\mathcal{L}(g_{\theta}(\tilde{z}), \tilde{y}) = \text{CE}(f^T(g_{\theta}(\tilde{z})), \tilde{y})$



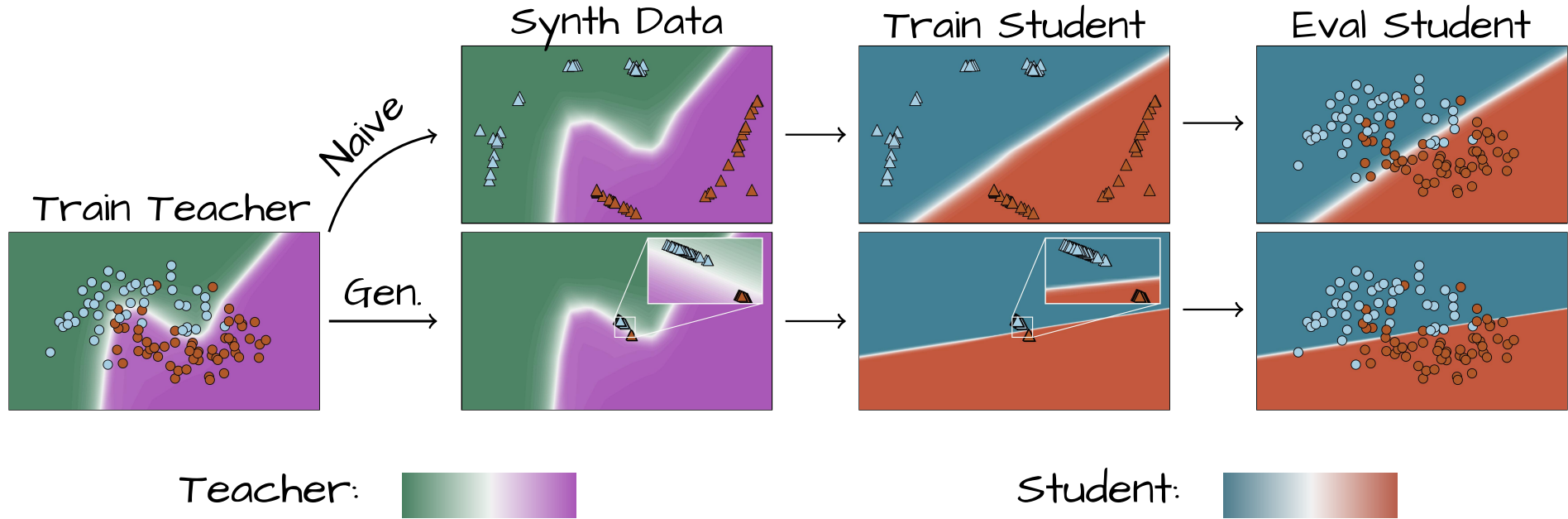
Extracting Knowledge

- **Generative:** Init. random latents (\tilde{z}, \tilde{y}) and minimize $\mathcal{L}(g_{\theta}(\tilde{z}), \tilde{y}) = \text{CE}(f^T(g_{\theta}(\tilde{z})), \tilde{y})$



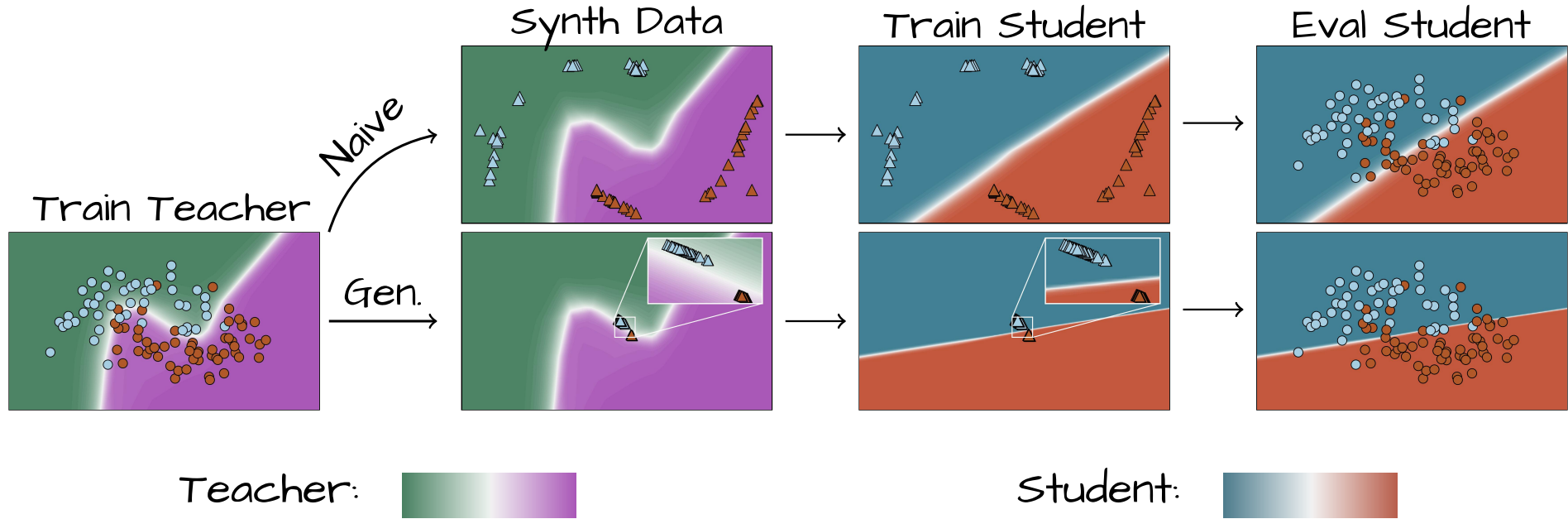
Extracting Knowledge

- **Generative:** Init. random latents (\tilde{z}, \tilde{y}) and minimize $\mathcal{L}(g_{\theta}(\tilde{z}), \tilde{y}) = \text{CE}(f^T(g_{\theta}(\tilde{z})), \tilde{y})$



Extracting Knowledge

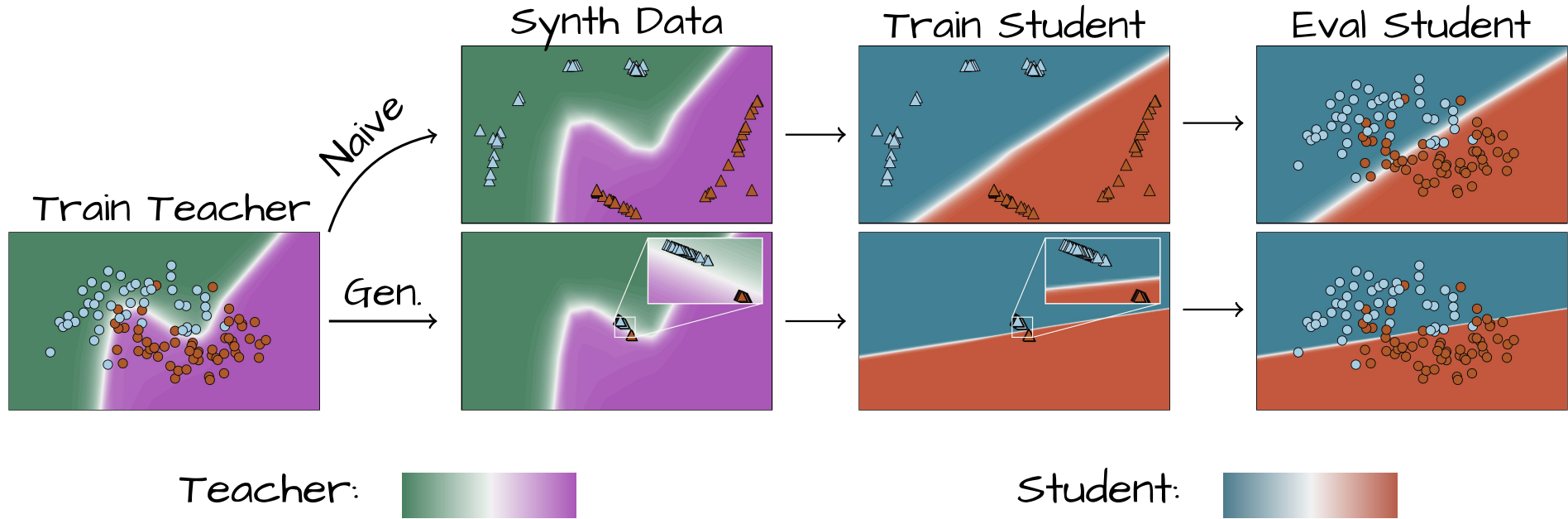
- **Generative:** Init. random latents (\tilde{z}, \tilde{y}) and minimize $\mathcal{L}(g_{\theta}(\tilde{z}), \tilde{y}) = \text{CE}(f^T(g_{\theta}(\tilde{z})), \tilde{y})$



What's missing?

Extracting Knowledge

- **Generative:** Init. random latents (\tilde{z}, \tilde{y}) and minimize $\mathcal{L}(g_{\theta}(\tilde{z}), \tilde{y}) = \text{CE}(f^T(g_{\theta}(\tilde{z})), \tilde{y})$



What's missing?

Naive: Keep classes close, else boundary becomes linear

Generative: Disperse samples along the relevant boundary region

CAKE: Contrastive Abductive Knowledge Extraction

Idea: **Contrast** sample pairs **noisily** *across* and *along* the *relevant* teacher decision boundary and **regularize** with data priors!

CAKE: Contrastive Abductive Knowledge Extraction

Idea: **Contrast** sample pairs **noisily** across and along the relevant teacher decision boundary and **regularize** with data priors!

- **Contrastive** samples between classes

$$\mathcal{L}_{\text{contr}}(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{1}[y_i \neq y_j] \|\mathbf{f}^T(\mathbf{x}_i) - \mathbf{f}^T(\mathbf{x}_j)\|_2^2$$

... or any other contrastive loss

CAKE: Contrastive Abductive Knowledge Extraction

Idea: **Contrast** sample pairs **noisily** across and along the relevant teacher decision boundary and **regularize** with data priors!

- **Contrastive** samples between classes

$$\mathcal{L}_{\text{contr}}(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{1}[y_i \neq y_j] \|\mathbf{f}^T(\mathbf{x}_i) - \mathbf{f}^T(\mathbf{x}_j)\|_2^2$$

... or any other contrastive loss

- **Regularize** using domain knowledge

$$\mathcal{L}_{\text{TV}}(\mathbf{x}) = \sum_{j,k} \|\mathbf{x}_{j,k} - \mathbf{x}_{j-1,k}\| + \|\mathbf{x}_{j,k} - \mathbf{x}_{j,k-1}\|$$

... or any other data prior

CAKE: Contrastive Abductive Knowledge Extraction

Idea: **Contrast** sample pairs **noisily** across and along the relevant teacher decision boundary and **regularize** with data priors!

- **Contrastive** samples between classes

$$\mathcal{L}_{\text{contr}}(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{1}[y_i \neq y_j] \|\mathbf{f}^T(\mathbf{x}_i) - \mathbf{f}^T(\mathbf{x}_j)\|_2^2$$

... or any other contrastive loss

- **Regularize** using domain knowledge

$$\mathcal{L}_{\text{TV}}(\mathbf{x}) = \sum_{j,k} \|\mathbf{x}_{j,k} - \mathbf{x}_{j-1,k}\| + \|\mathbf{x}_{j,k} - \mathbf{x}_{j,k-1}\|$$

... or any other data prior

- **Noisily** disperse samples along the boundary

CAKE: Contrastive Abductive Knowledge Extraction

Idea: **Contrast** sample pairs **noisily** across and along the relevant teacher decision boundary and **regularize** with data priors!

- **Contrastive** samples between classes

$$\mathcal{L}_{\text{contr}}(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{1}[y_i \neq y_j] \|\mathbf{f}^T(\mathbf{x}_i) - \mathbf{f}^T(\mathbf{x}_j)\|_2^2$$

... or any other contrastive loss

- **Regularize** using domain knowledge

$$\mathcal{L}_{\text{TV}}(\mathbf{x}) = \sum_{j,k} \|\mathbf{x}_{j,k} - \mathbf{x}_{j-1,k}\| + \|\mathbf{x}_{j,k} - \mathbf{x}_{j,k-1}\|$$

... or any other data prior

- **Noisily** disperse samples along the boundary

Explicit: Langevin Dynamics $\mathbf{x}_i^{t+1} = \mathbf{x}_i^t + \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_i^t) \eta(t) + \sqrt{2\eta(t)} \varepsilon_i^t$, with $\varepsilon_i^t \sim N(0, I)$

CAKE: Contrastive Abductive Knowledge Extraction

Idea: **Contrast** sample pairs **noisily** across and along the relevant teacher decision boundary and **regularize** with data priors!

- **Contrastive** samples between classes

$$\mathcal{L}_{\text{contr}}(\mathbf{x}_i, \mathbf{x}_j) = \mathbb{1}[y_i \neq y_j] \|\mathbf{f}^T(\mathbf{x}_i) - \mathbf{f}^T(\mathbf{x}_j)\|_2^2$$

... or any other contrastive loss

- **Regularize** using domain knowledge

$$\mathcal{L}_{\text{TV}}(\mathbf{x}) = \sum_{j,k} \|\mathbf{x}_{j,k} - \mathbf{x}_{j-1,k}\| + \|\mathbf{x}_{j,k} - \mathbf{x}_{j,k-1}\|$$

... or any other data prior

- **Noisily** disperse samples along the boundary

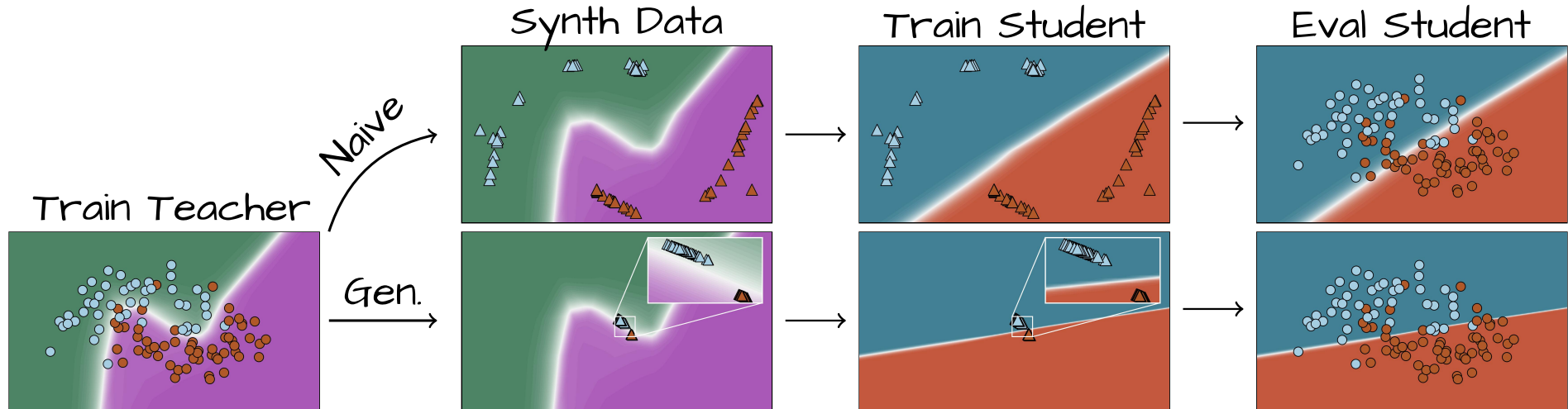
Explicit: Langevin Dynamics $\mathbf{x}_i^{t+1} = \mathbf{x}_i^t + \nabla_x \mathcal{L}(\mathbf{x}_i^t) \eta(t) + \sqrt{2\eta(t)} \varepsilon_i^t$, with $\varepsilon_i^t \sim N(0, I)$


Implicit: Stochasticity of SGD and step size schedules $\eta(t)$ is enough


... or any other noise injection

Extracting Knowledge

- **CAKE**: Contrast pairs noisily across and along the relevant teacher decision boundary.

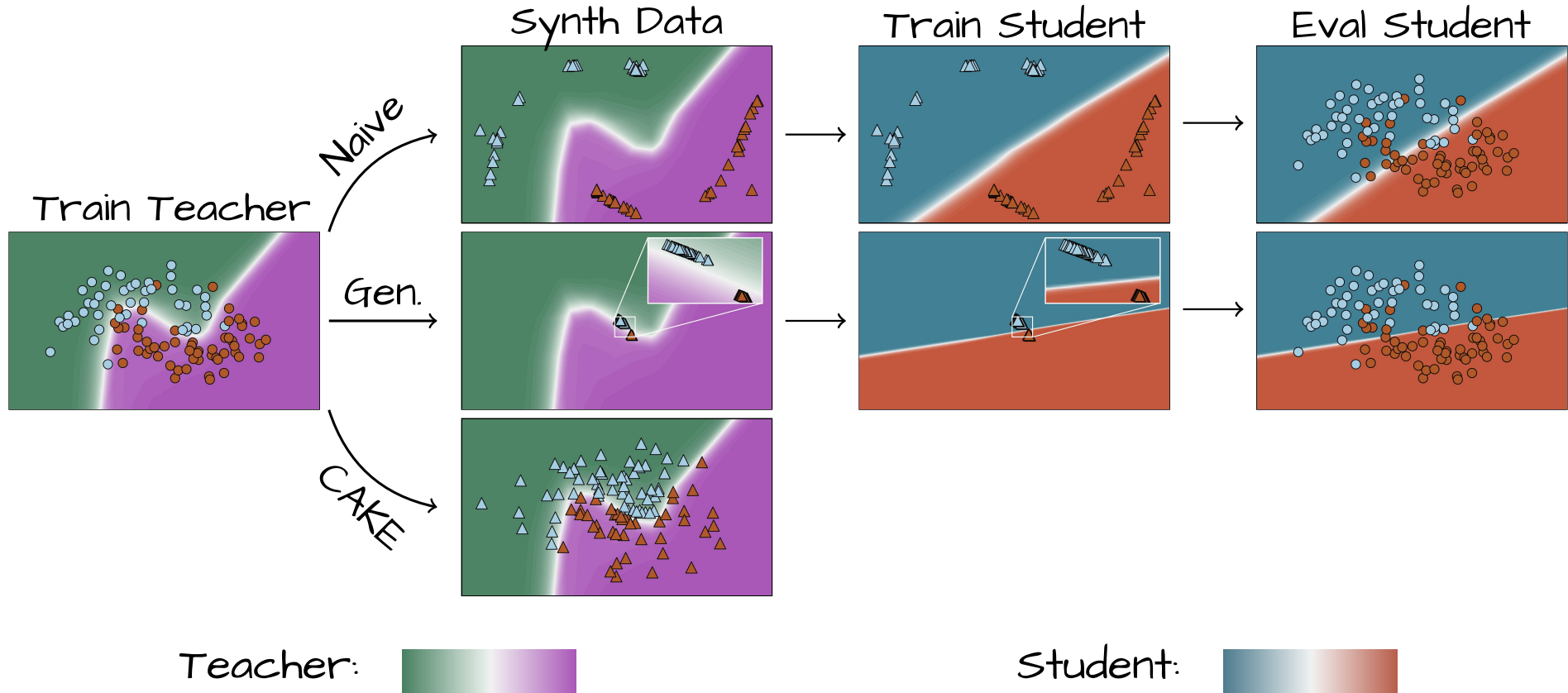


Teacher: 

Student: 

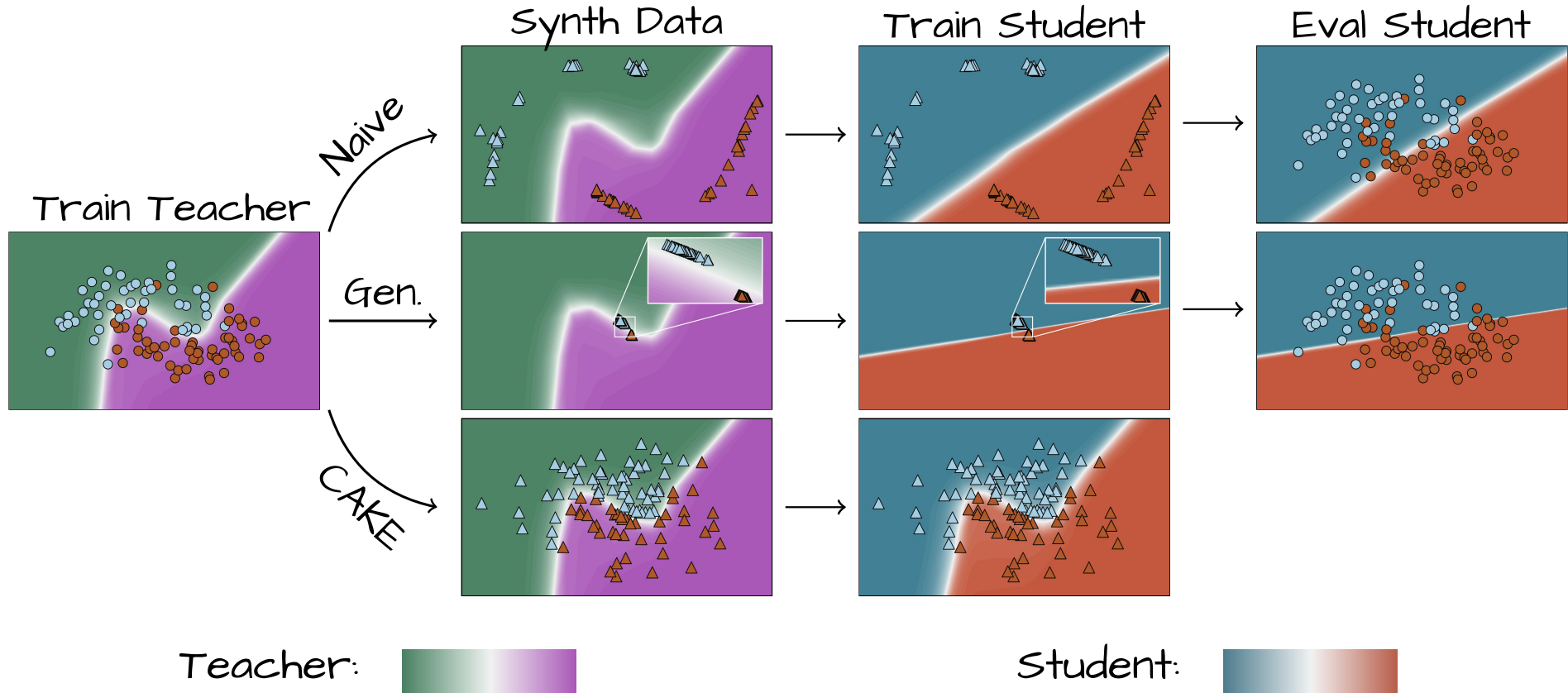
Extracting Knowledge

- **CAKE**: Contrast pairs noisily across and along the relevant teacher decision boundary.



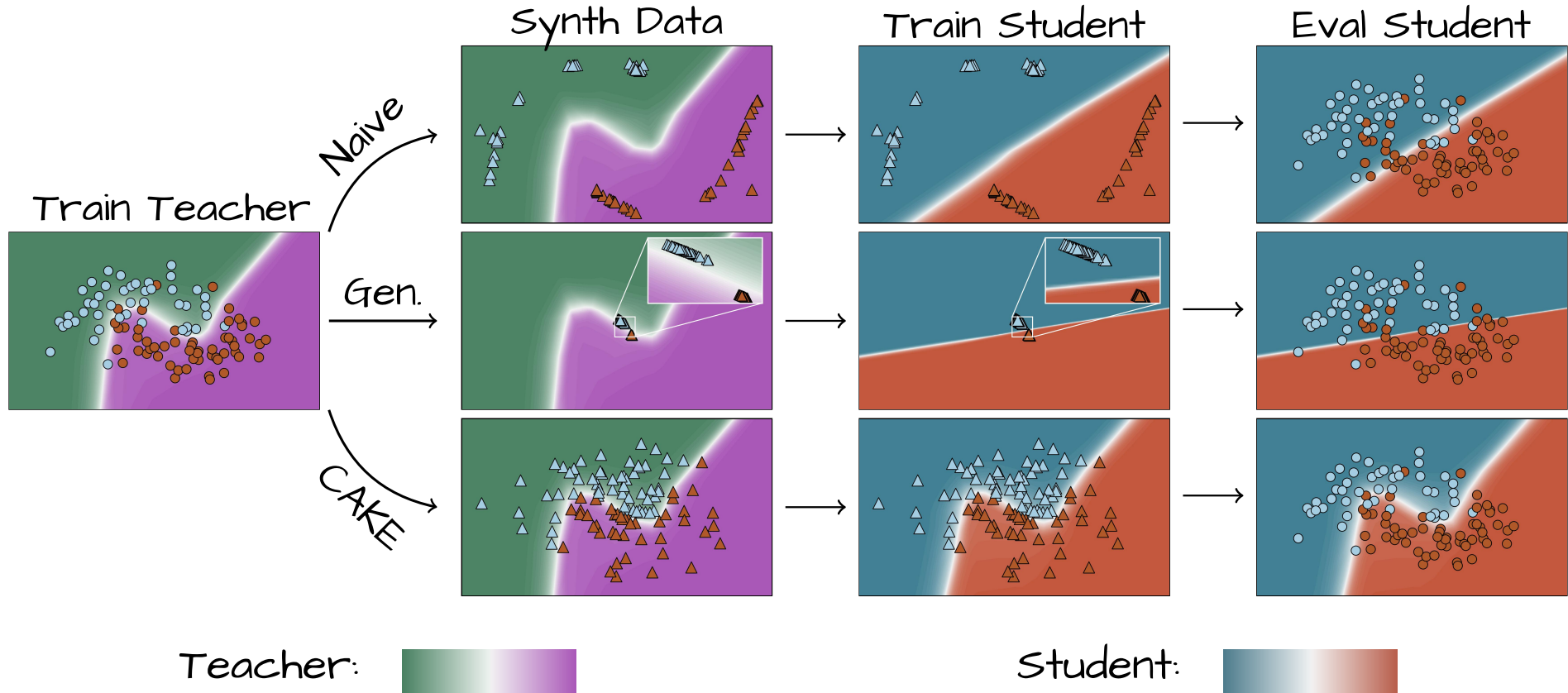
Extracting Knowledge

- **CAKE**: Contrast pairs noisily across and along the relevant teacher decision boundary.



Extracting Knowledge

- **CAKE**: Contrast pairs noisily across and along the relevant teacher decision boundary.



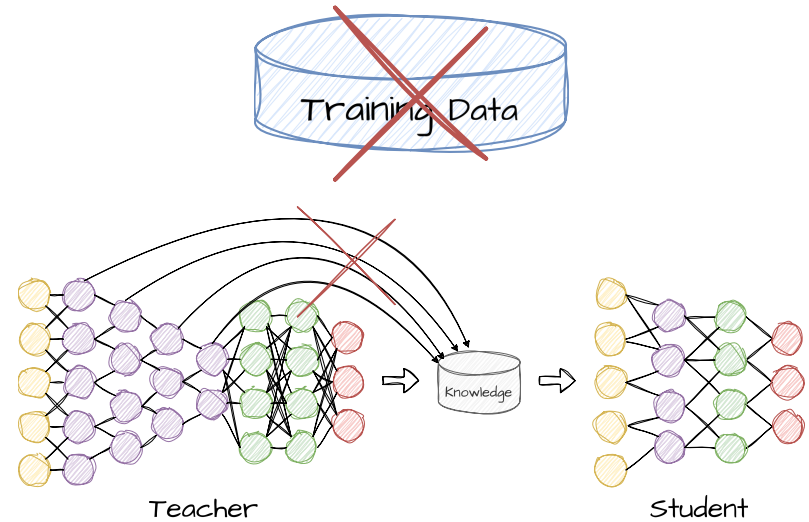
CAKE Lifts Knowledge Distillation Restrictions

- No original data access



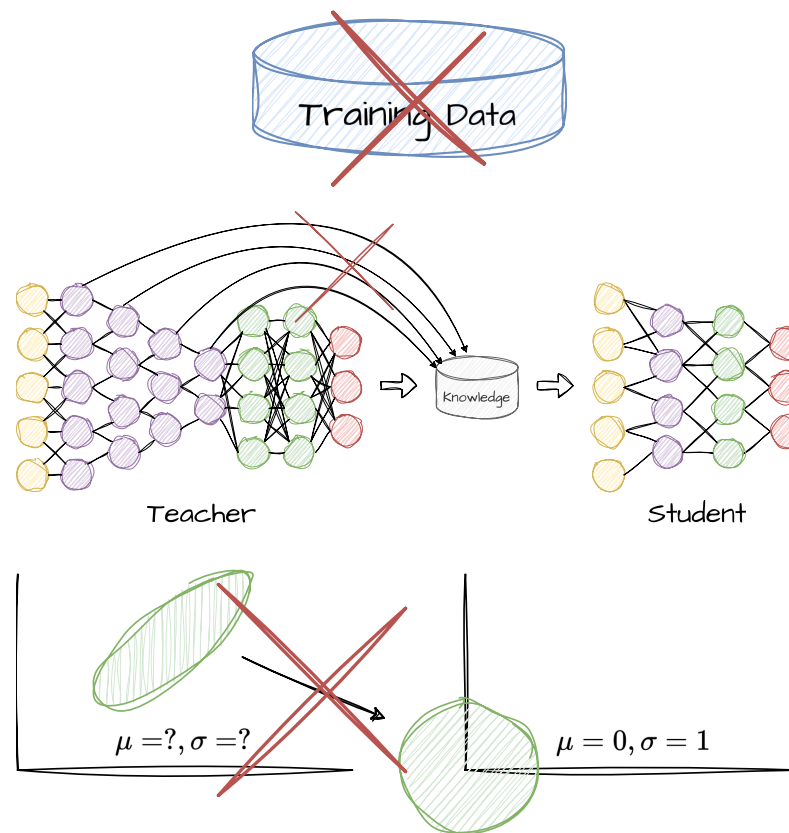
CAKE Lifts Knowledge Distillation Restrictions

- No original data access
- No model access
e.g. intermediate activations



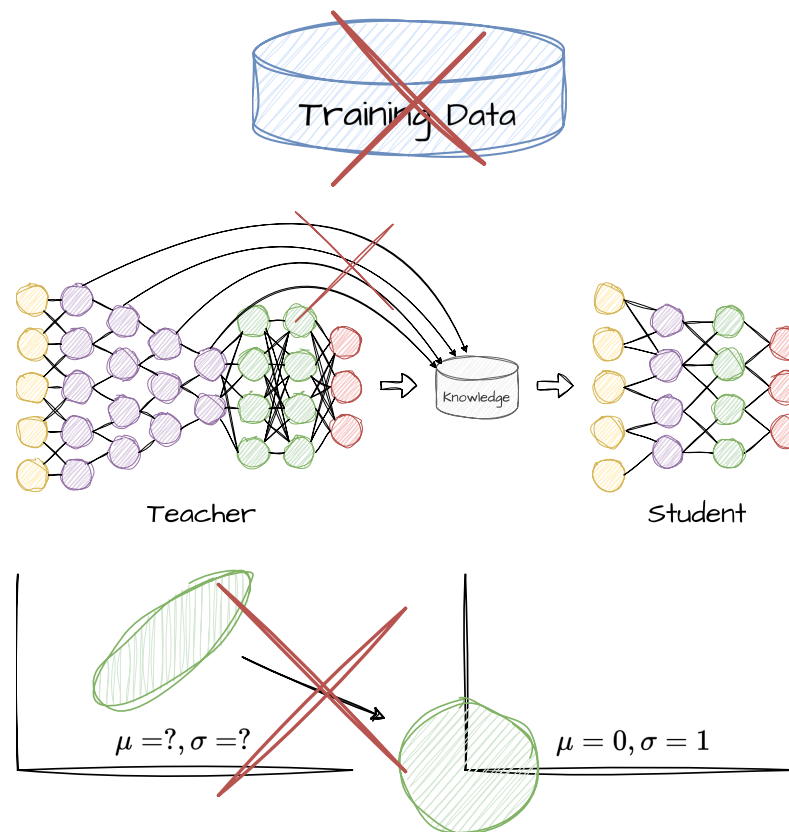
CAKE Lifts Knowledge Distillation Restrictions

- No original data access
- No model access
e.g. intermediate activations
- No model assumptions
e.g. BatchNorm, linear penultimate layer



CAKE Lifts Knowledge Distillation Restrictions

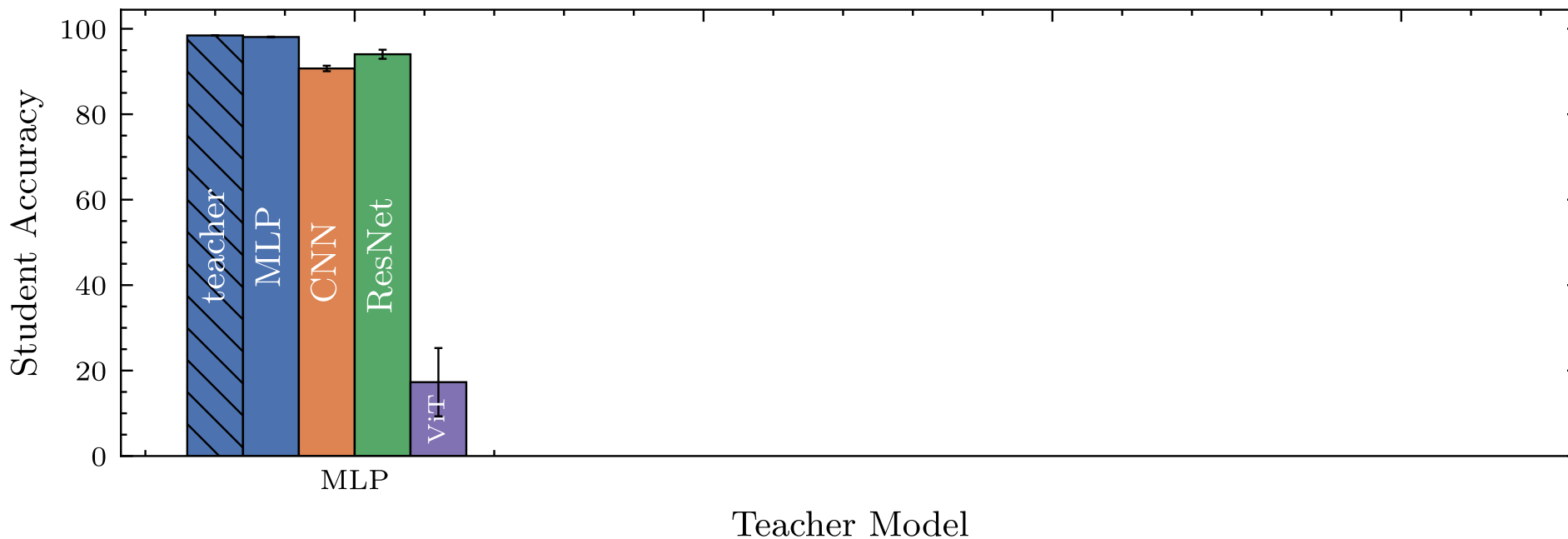
- No original data access
- No model access
e.g. intermediate activations
- No model assumptions
e.g. BatchNorm, linear penultimate layer



↪ CAKE can be applied to any “blackbox” model which is differentiable w.r.t. its input.

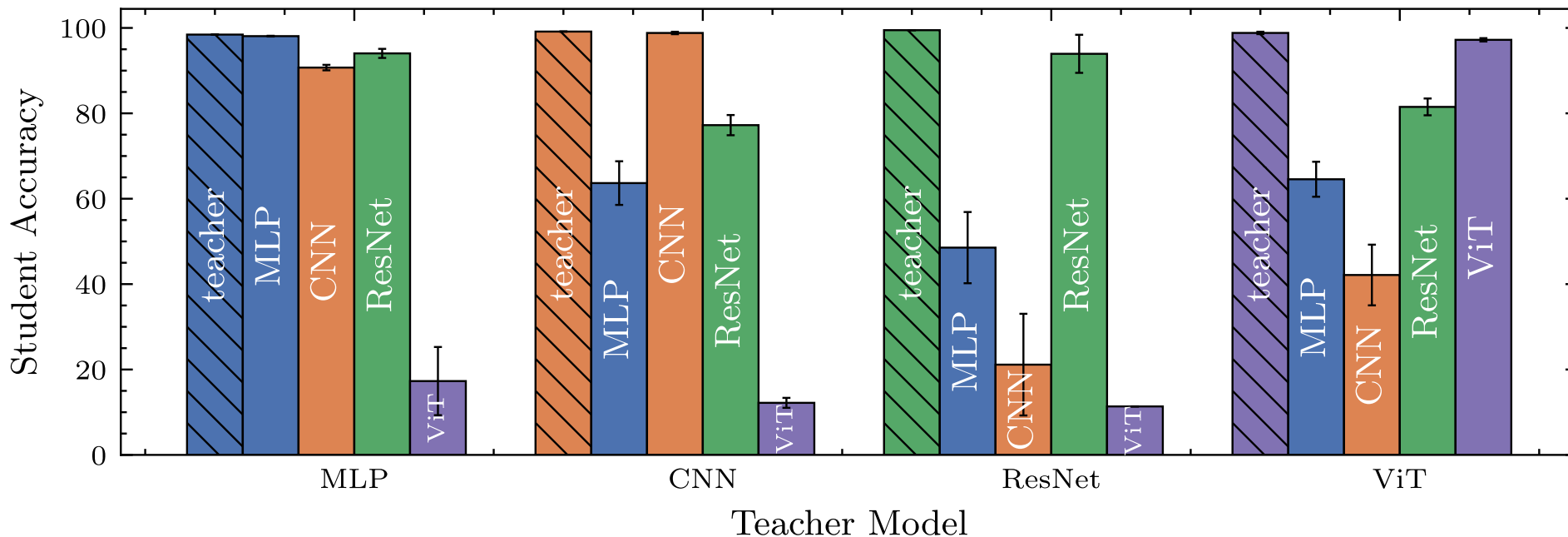
CAKE Across Model Types

Distilling MNIST from model type A to model type B



CAKE Across Model Types

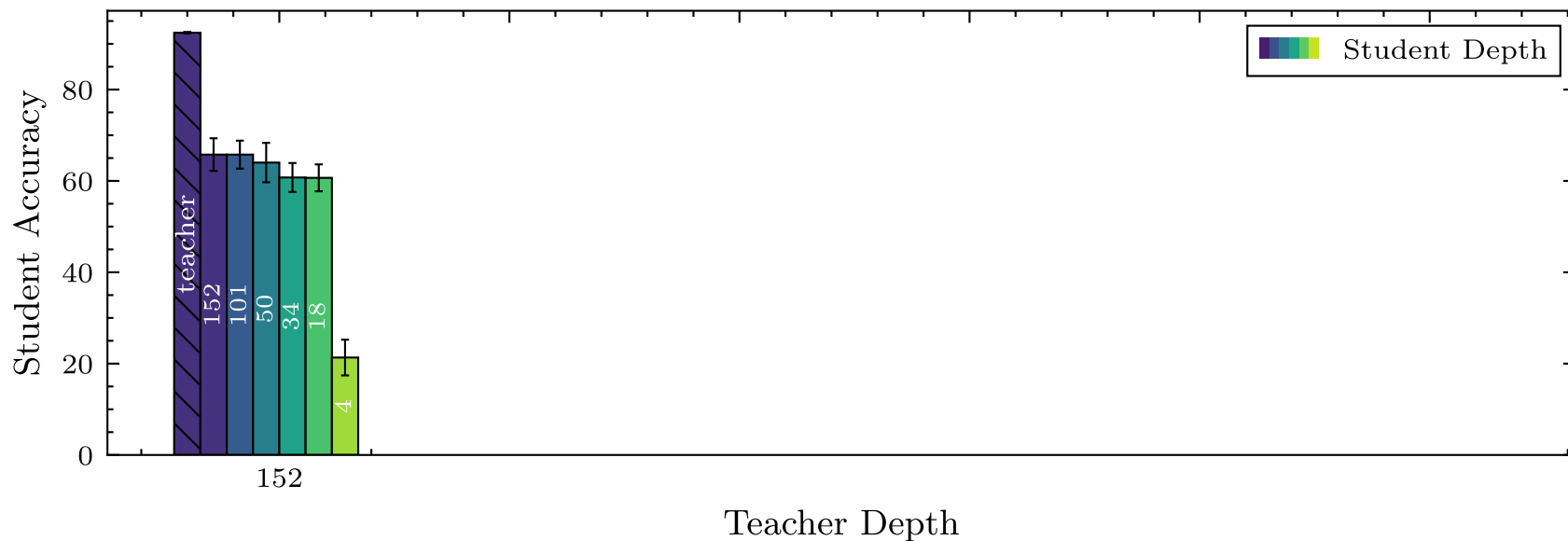
Distilling MNIST from model type A to model type B



Takeaways: 1. Similar inductive bias → better distillation
2. Less inductive bias → better distillation 3. ResNet is a safe student model choice.

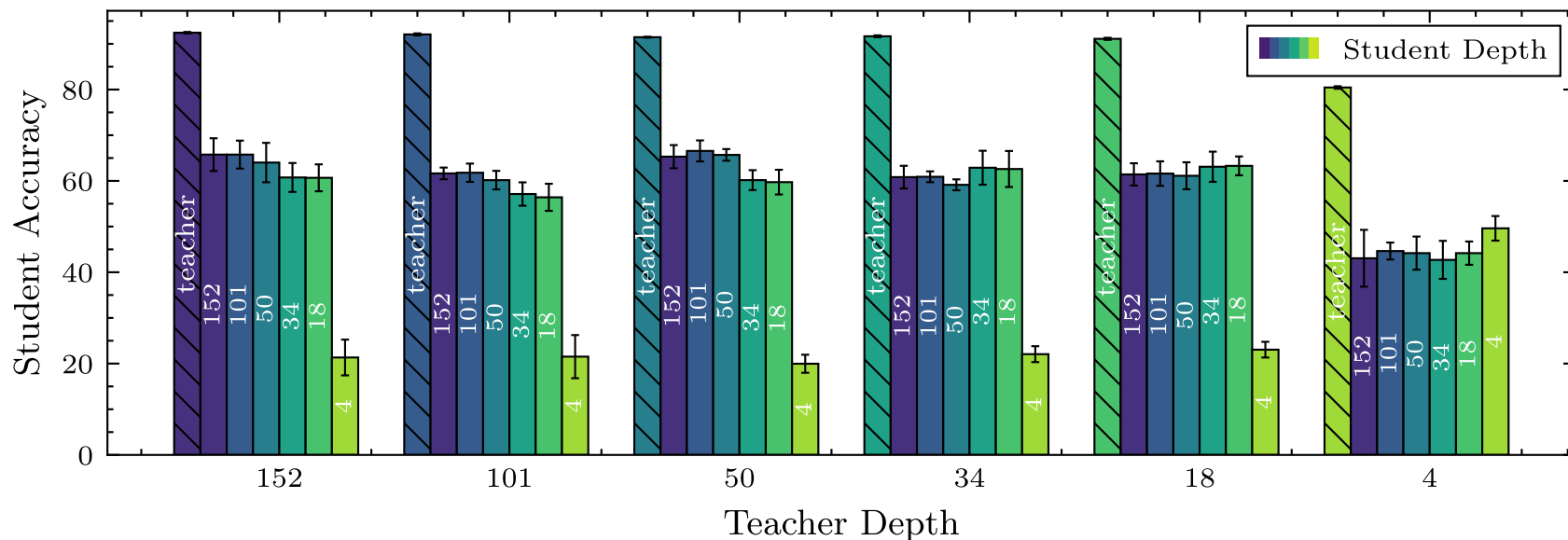
CAKE Across Scales

Distilling CIFAR-10 knowledge from ResNet-X to ResNet-Y (152, 101, 50, 34, 18, 4)



CAKE Across Scales

Distilling CIFAR-10 knowledge from ResNet-X to ResNet-Y (152, 101, 50, 34, 18, 4)



Takeaway: CAKE can compress models at a stable accuracy until capacity is too heavily constrained.

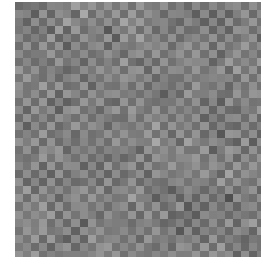
CAKE Synthetic Samples

No visual resemblance with original training data.

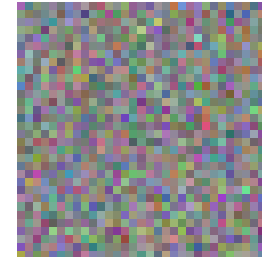
Possible future work includes:

- Differential privacy?
- Data utility and privacy trade-offs?
- Robustness against adversarial attacks?

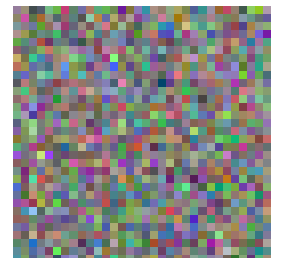
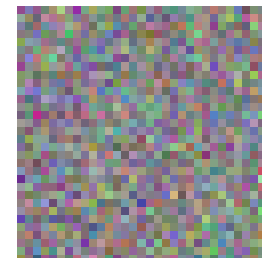
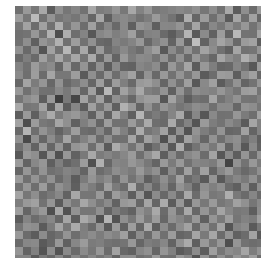
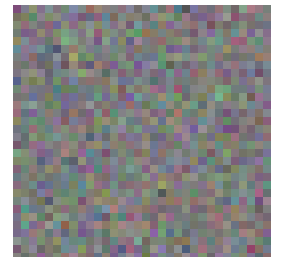
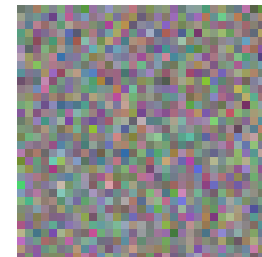
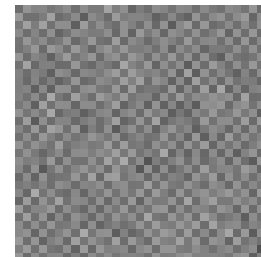
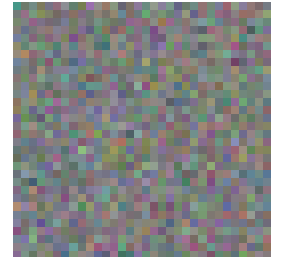
MNIST



SVHN



CIFAR



Summary and Outlook

CAKE is a **data-free** and **model-agnostic** knowledge distillation method, that ...

- can distill models *across scales*
- can distill between *different model types*
- doesn't produce data-like samples (visually)

Summary and Outlook

CAKE is a **data-free** and **model-agnostic** knowledge distillation method, that ...

- can distill models *across scales*
- can distill between *different model types*
- doesn't produce data-like samples (visually)

Future work

- Estimate gradients? → truly “blackbox”, API-model possible
- Investigate the data privacy perspective?
- Investigate explicit instead of implicit noise

Still interested?

Join me at Room 2, Poster #117

Paper



Code



`steven.braun@cs.tu-darmstadt.de`