

# Knowledge Distillation

EECE695D: Efficient ML Systems

Spring 2025

# Recap

- **Problem.** Inference cost

- Large models are **strong** but **heavy**

Large Model



- **Goal.** Lowering the inference cost

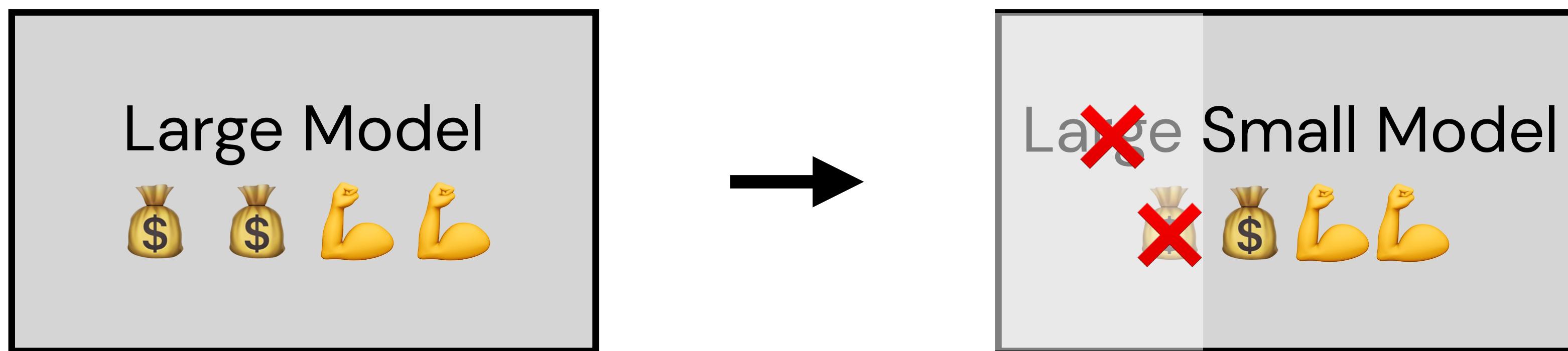
- Make a **small** but **strong** model

Small Model



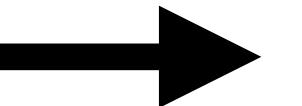
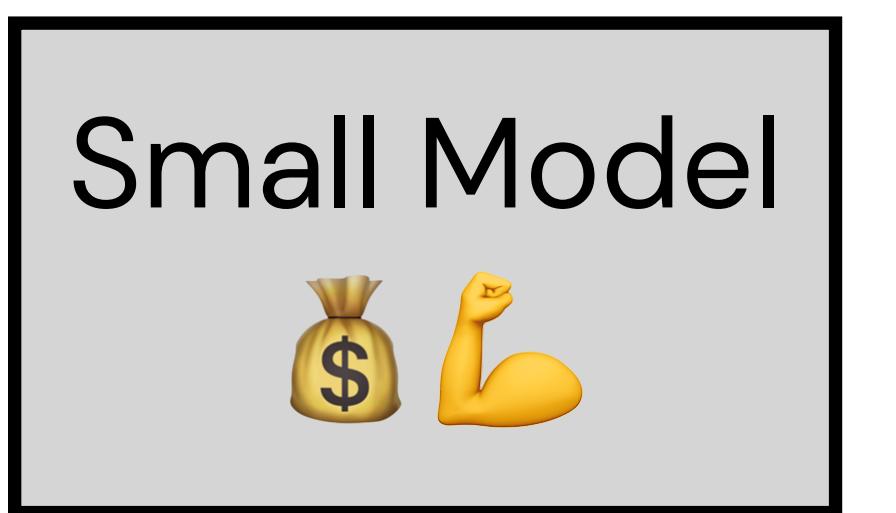
# Recap

- W2-W3. Reduce the cost of large models
  - Idea. Replace matmuls with lighter options
    - Sparsity
    - Low-Precision



# Recap

- **Today.** Making small models stronger
  - Knowledge Distillation



# Basic idea

# Motivation

- **Small model.**

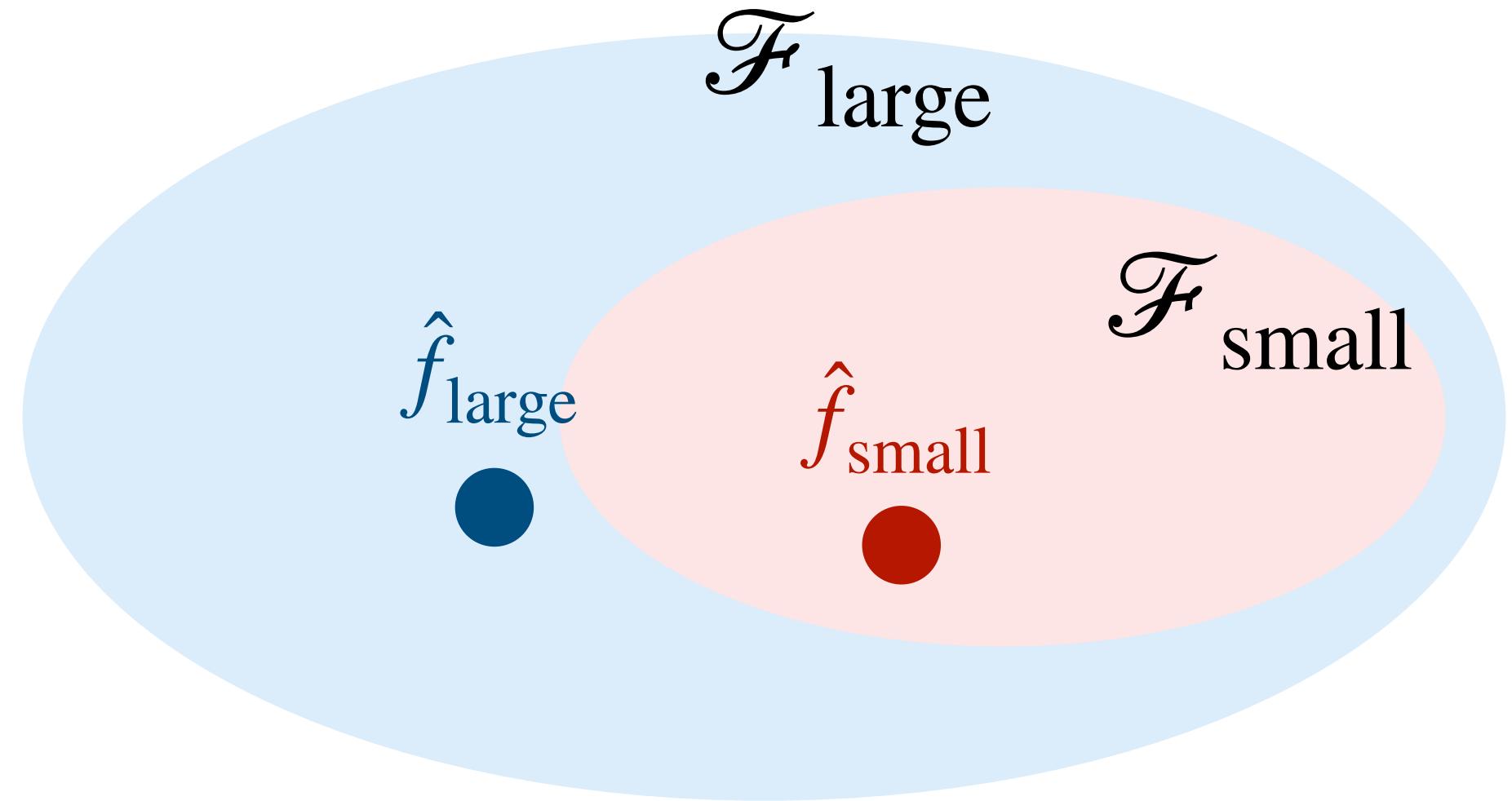
$$\hat{f}_{\text{small}} = \arg \min_{f \in \mathcal{F}_{\text{small}}} \hat{L}(f)$$

- **Large model.**

$$\hat{f}_{\text{large}} = \arg \min_{f \in \mathcal{F}_{\text{large}}} \hat{L}(f)$$

- **Question.** Why do  $\hat{f}_{\text{small}}$  work much poorer than  $\hat{f}_{\text{large}}$ ?

- (A)  $\mathcal{F}_{\text{small}}$  does not have a good optimum
- (B) Good optimum exists, but is difficult to find

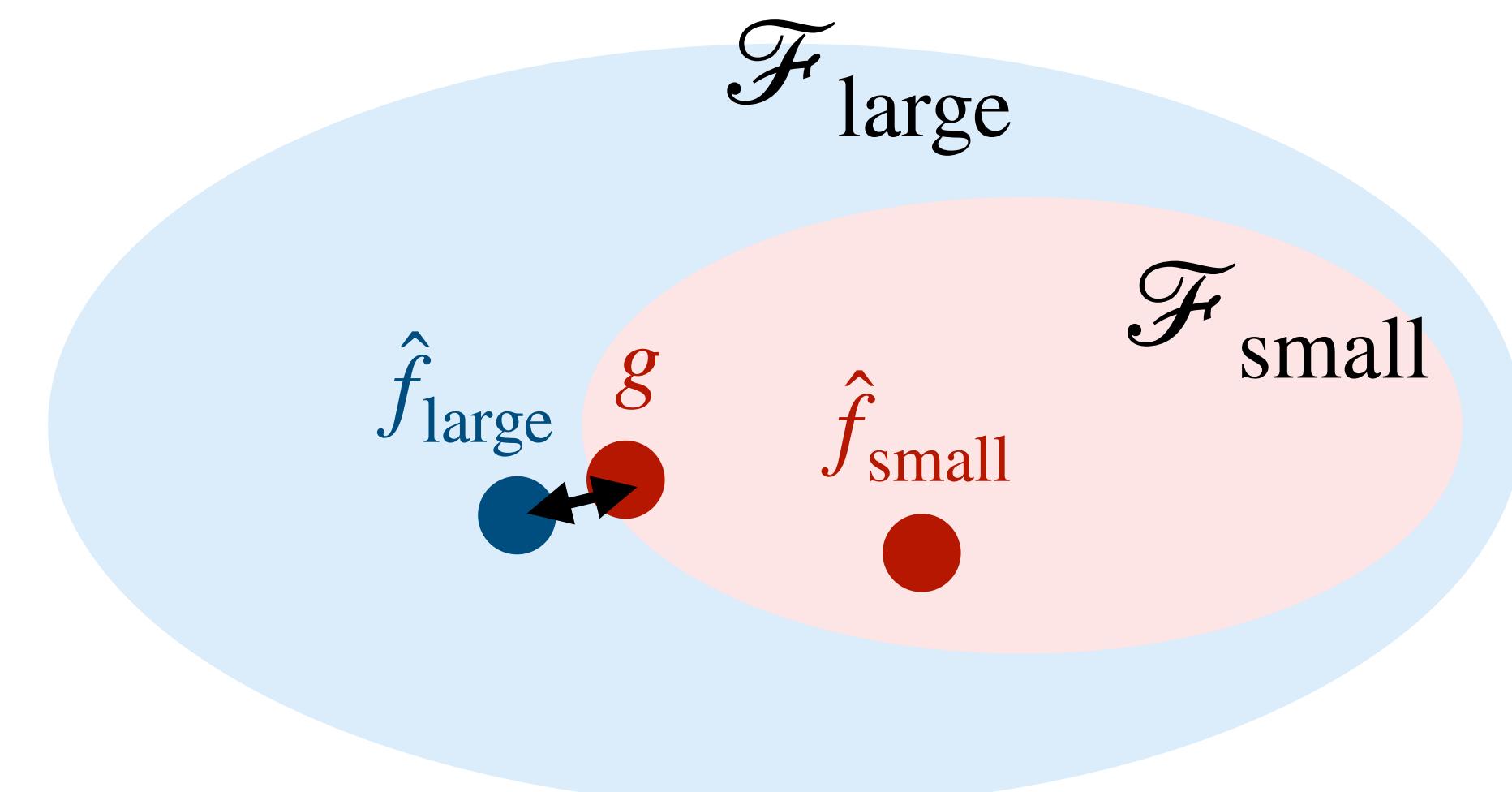


# Motivation

(A)  $\mathcal{F}_{\text{small}}$  does not have a good optimum  $\Rightarrow \text{False}$

- **Theoretically.** Few weights are enough to approximate any function ([link](#))
  - Thus, it is likely that there exists some  $g$  such that

$$g \approx \hat{f}_{\text{large}}, \quad g \in \mathcal{F}_{\text{small}}$$



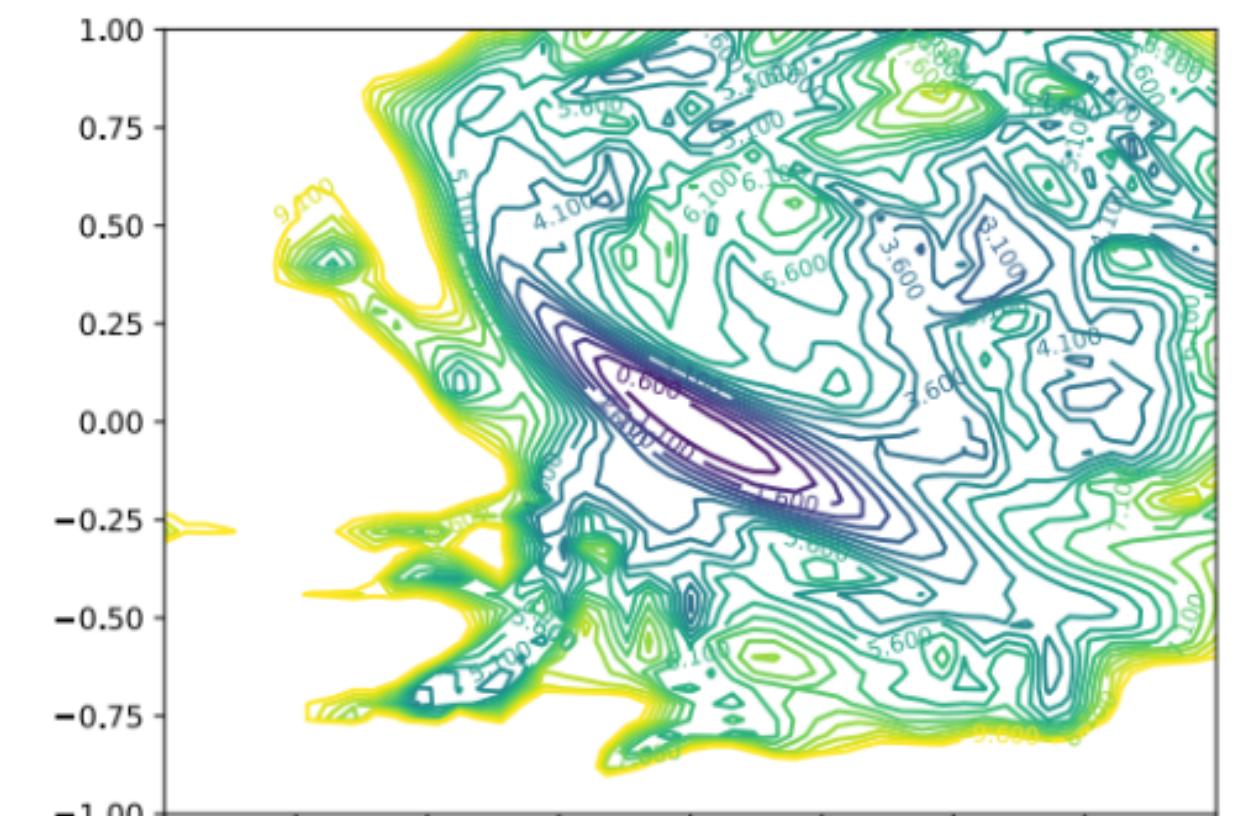
# Motivation

(B) Good optimum exists, but is difficult to find

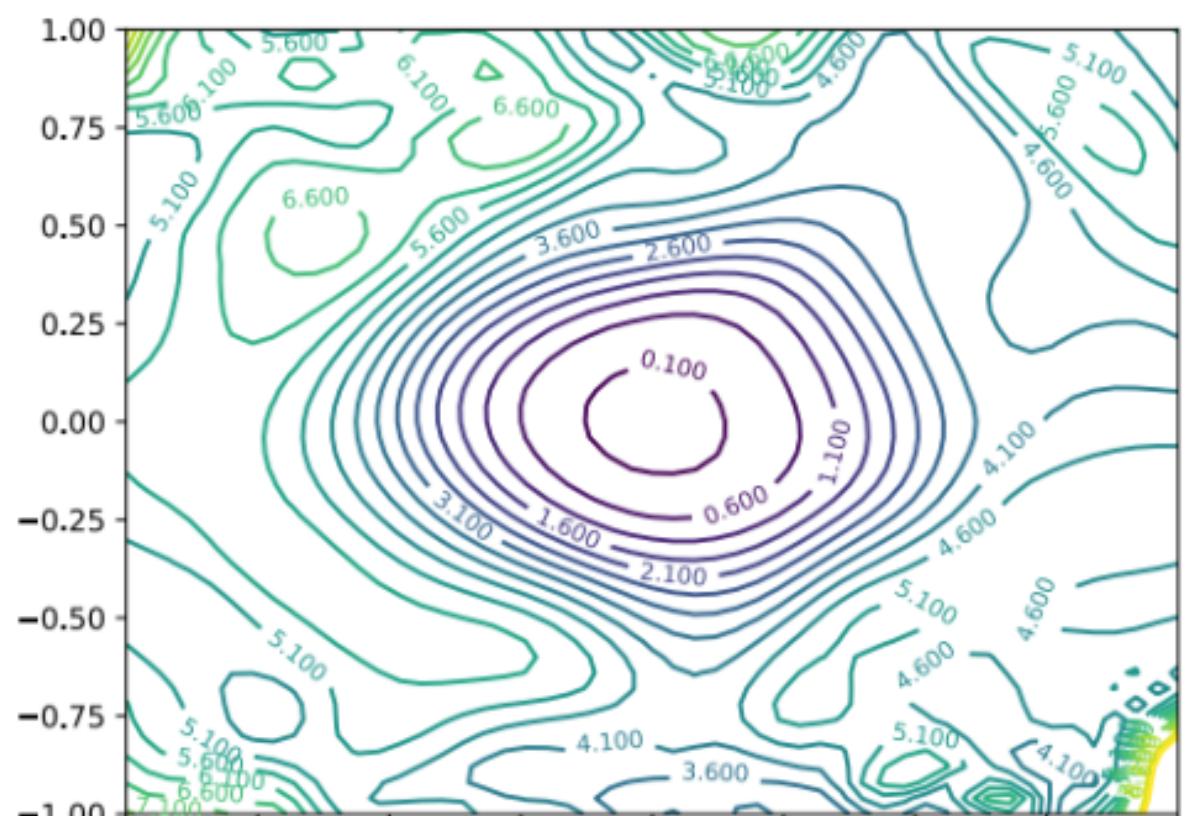
⇒ True

- **Empirically.** GD on small nets suffer from local minima

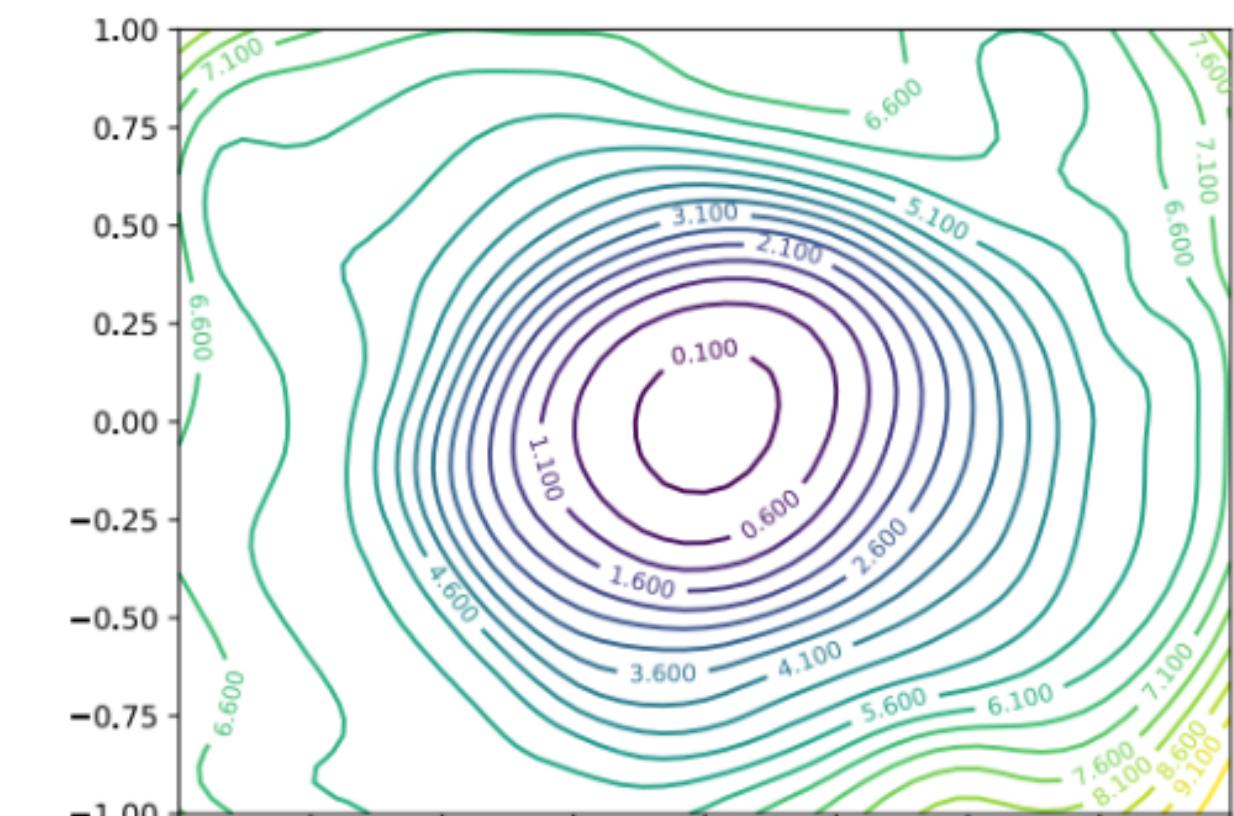
- Example. Loss landscape of ResNet-like models (narrow → wide)



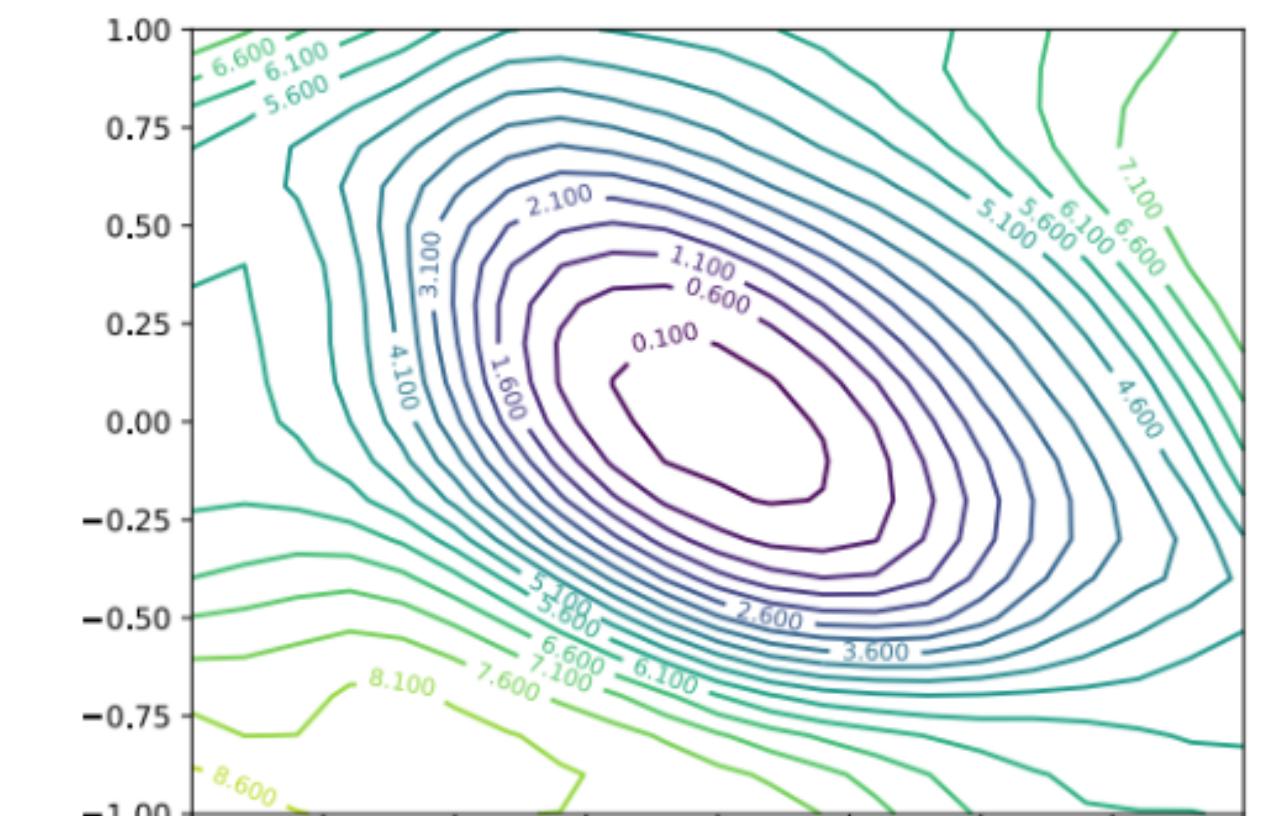
(e)  $k = 1, 13.31\%$



(f)  $k = 2, 10.26\%$



(g)  $k = 4, 9.69\%$

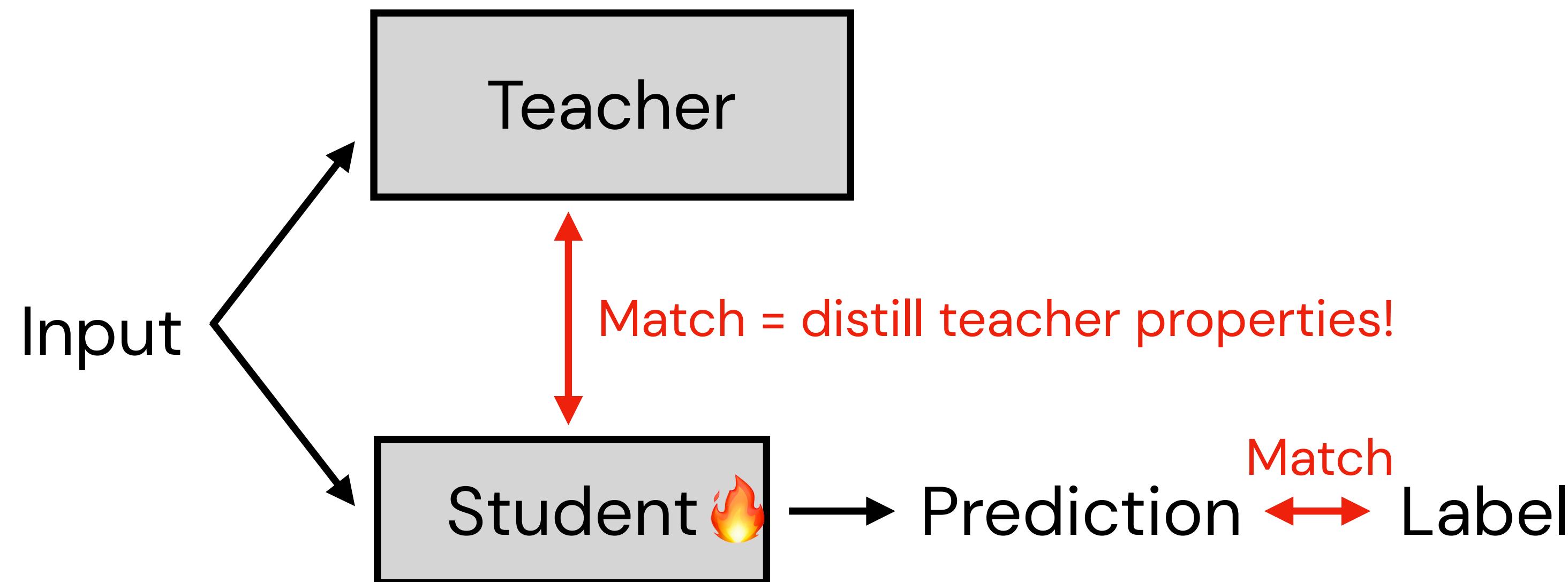


(h)  $k = 8, 8.70\%$

# Knowledge distillation

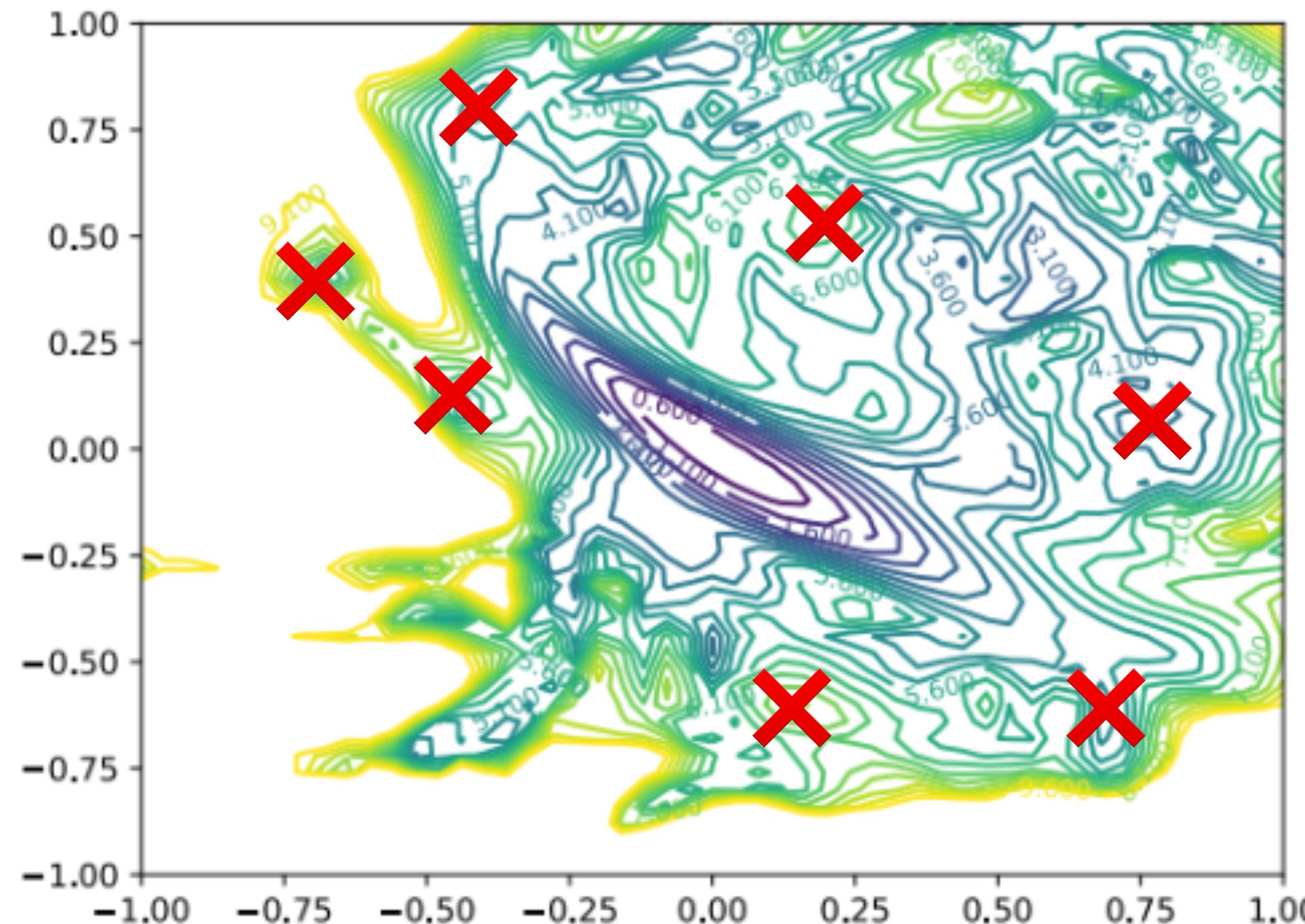
- **Question.** How do we encourage the small model to find a nice function?
- **Idea.** Explicitly regularize the learned function to be functionally close to the pre-trained large model  $\hat{f}_{\text{large}}$  (also called “teacher”)

$$\min_{f \in \mathcal{F}_{\text{small}}} \hat{L}(f) + \text{dist}(f, \hat{f}_{\text{large}})$$



# Knowledge distillation

- This prevents converging to bad local minima:
  - The trained small model (called “student”) is further constrained to mimic the functional properties of the teacher faithfully



✖: Cannot converge, as are not similar to  $\hat{f}_{\text{large}}$

# Agenda

- Which properties to be distilled?
- Which teacher to distill from?
- Other tips

Which properties?

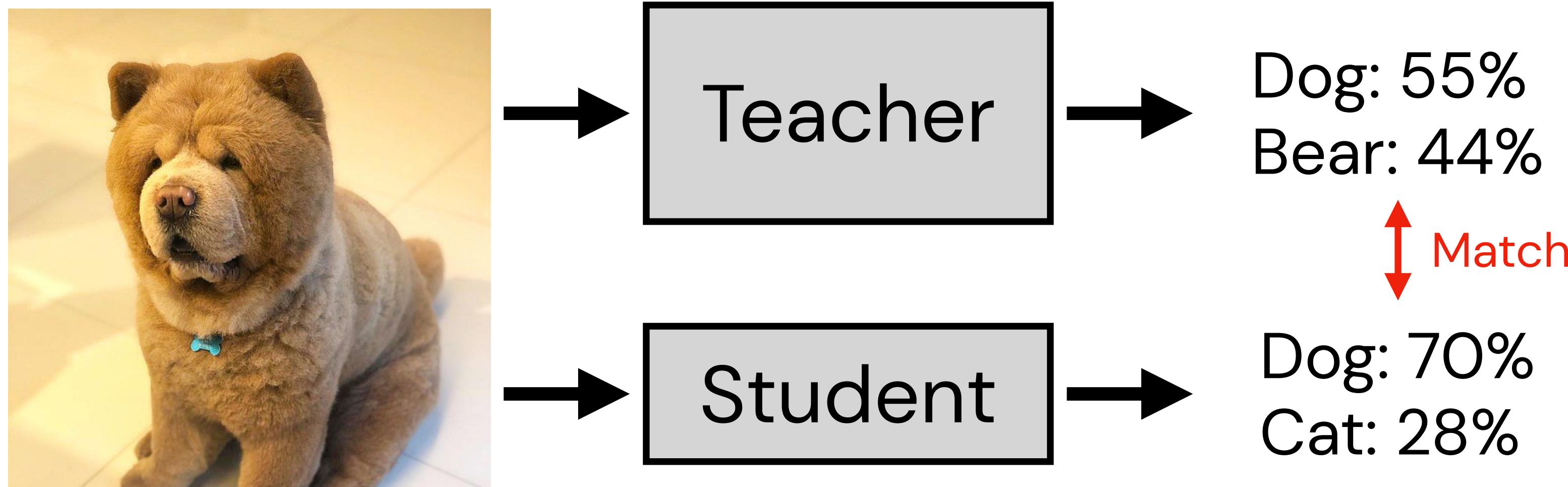
# Which properties?

- Countless different forms

- Logit
- Feature
- Relational

# Logit distillation

- The classical form
- Idea. Make similar predictions at the logit level (a.k.a. “dark knowledge”)
  - Confidence, Close Classes, Wrong Predictions



# Logit distillation

- Concretely, consider the K-class classification

- **Logit.** Model output, before softmax (denoted by  $z^{(1)}, \dots, z^{(k)}$ )
- **(Generalized) Softmax.** The softmax is a probability

$$p^{(i)} = \frac{\exp(z^{(i)}/\tau)}{\sum_{j=1}^k \exp(z^{(j)}/\tau)}$$

- $\tau$ : **temperature** hyperparameter
  - High temp  $\rightarrow$  uniform (low confidence)
  - Low temp  $\rightarrow$  zero/one (high confidence)

# Logit distillation

- The logit distillation minimizes the combined loss

$$\ell_{\text{CE}} + \lambda \cdot \ell_{\text{KD}}$$

- $\ell_{\text{CE}}$ : Cross-entropy loss of the student

$$\ell_{\text{CE}} = - \sum_{j=1}^k \mathbf{1}\{y=j\} \cdot \log p_S^{(j)}$$

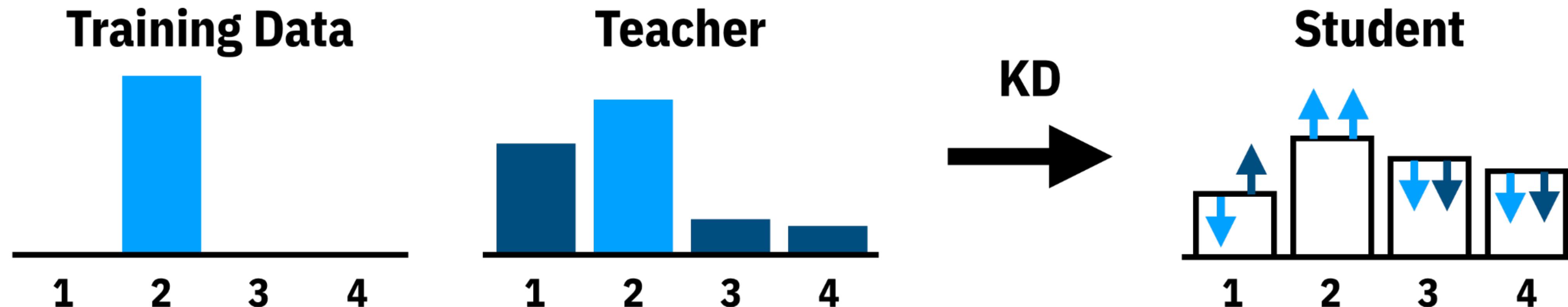
- $\ell_{\text{KD}}$ : KL-divergence loss

$$\ell_{\text{KD}} = - \sum_{j=1}^k p_T^{(j)} \log p_S^{(j)}$$

- $p_T, p_S$ : teacher and student softmax, resp.

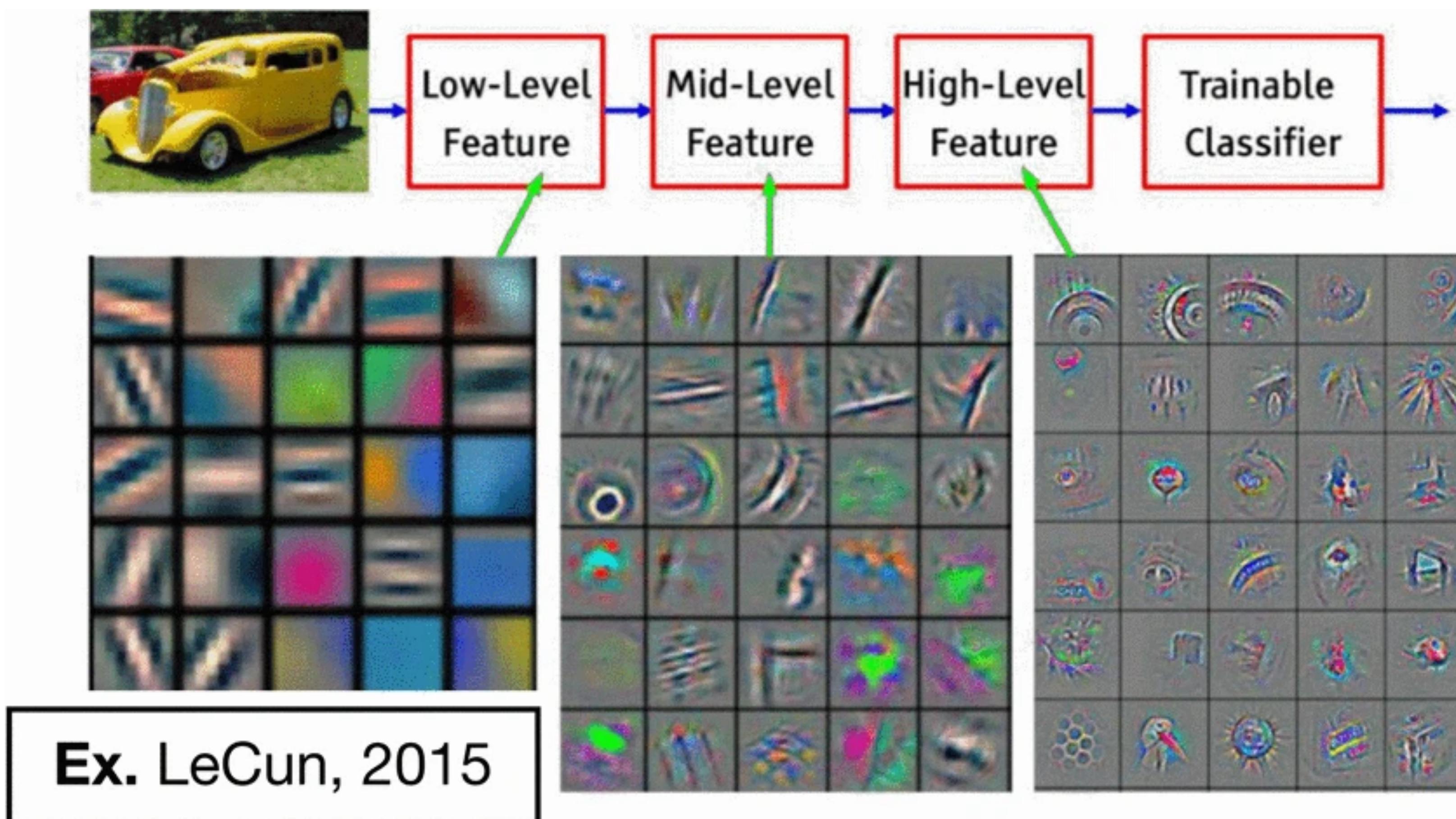
# Logit distillation

- **CE Loss.** Encourage student softmax to be close to zero-or-one
- **KD Loss.** Encourage student softmax to be similar to teacher softmax



# Feature distillation

- Idea. Regularize the teacher and student to have similar activations
  - That is, they make decisions based on similar reasons

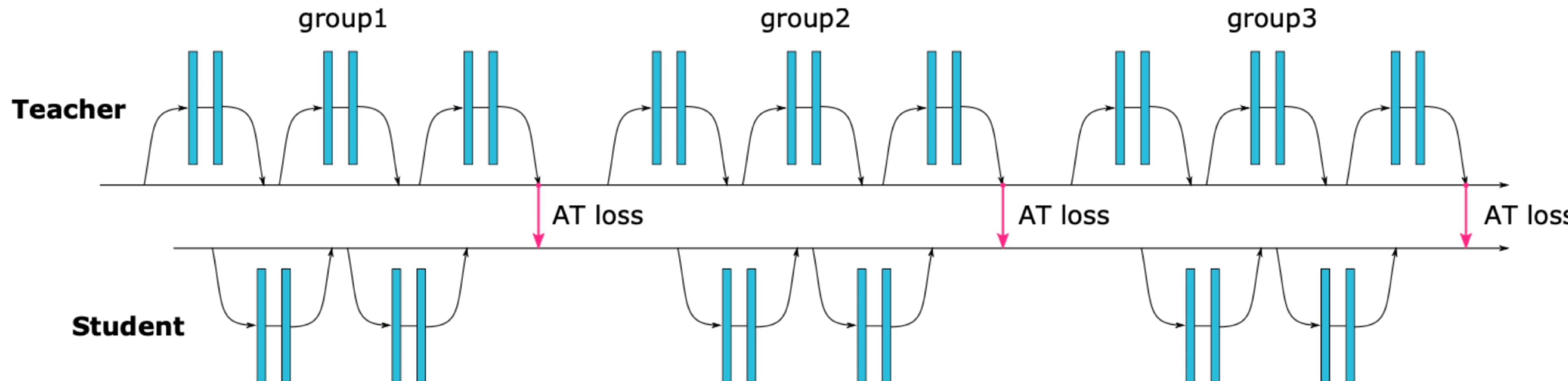


# Feature distillation

- Feature distillation minimizes the combined loss

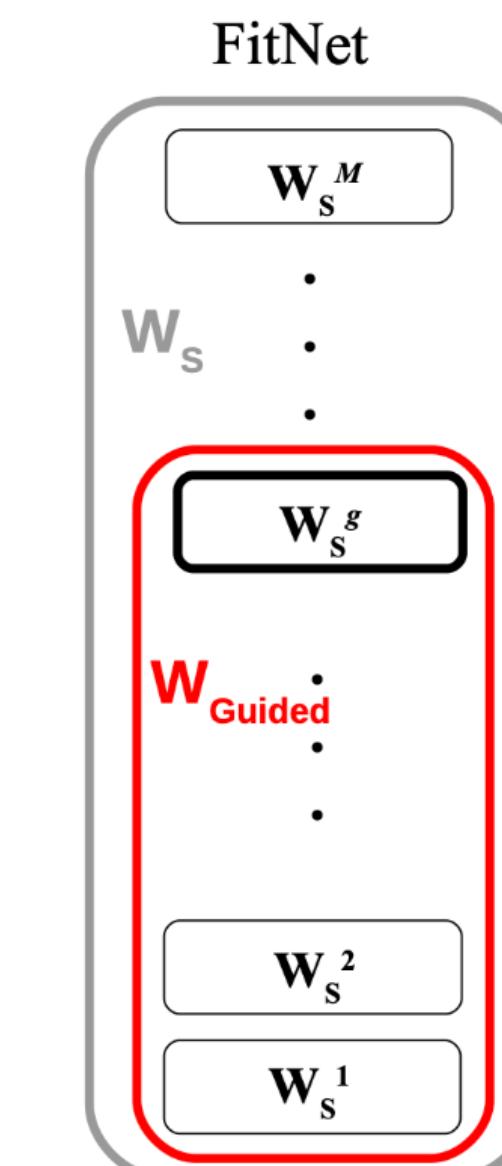
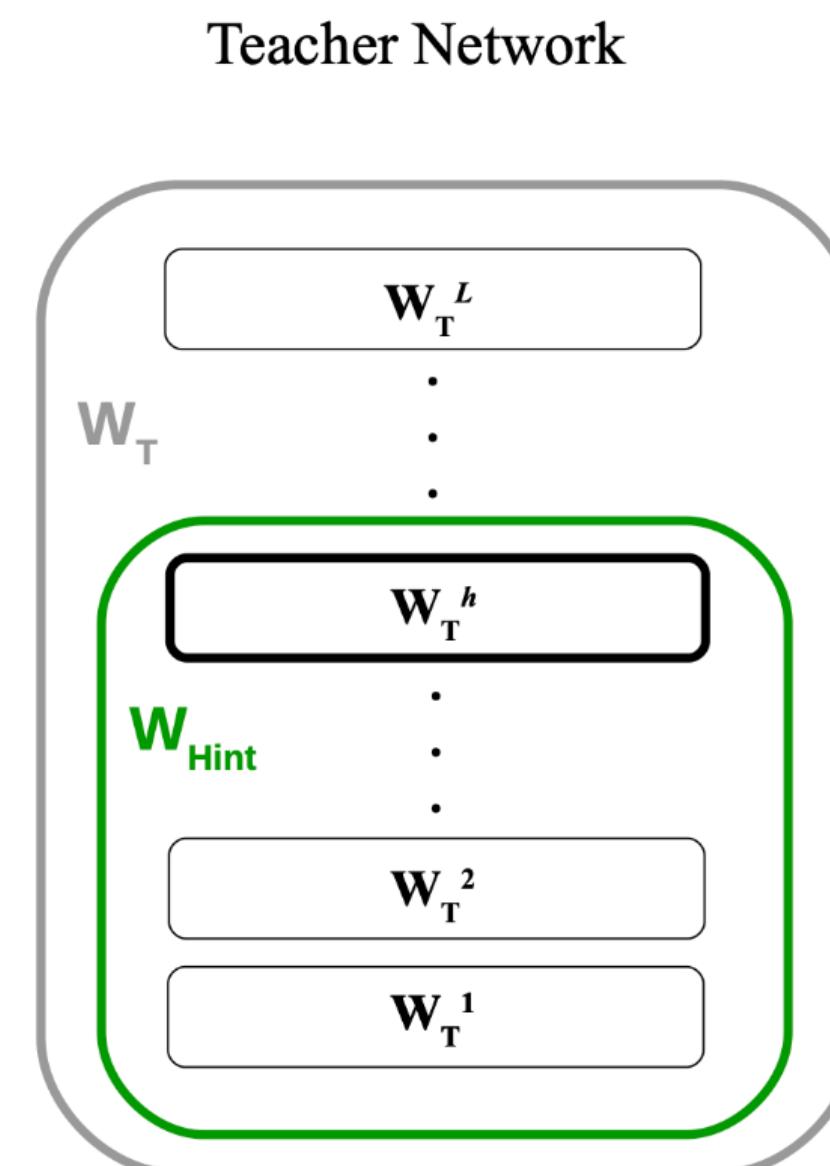
$$\ell_{\text{CE}} + \lambda \cdot \sum_{j=1}^J \left\| \frac{Z_T}{\|Z_T\|} - \frac{Z_S}{\|Z_S\|} \right\|$$

- $Z_T, Z_S$ : normalized activations of teacher and student  
(usually channel-pooled into HxW tensor)

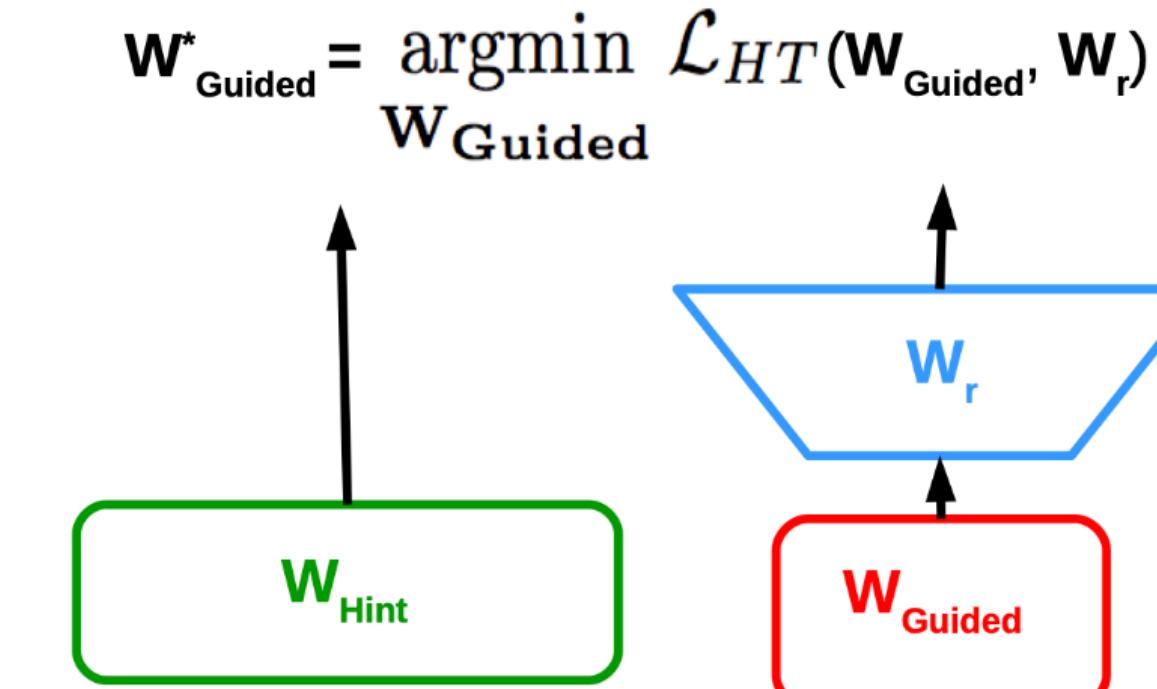


# Feature distillation

- Often, the activation sizes doesn't match
  - Small models typically downsample more aggressively
- Solution.** Use additional projection (hint), and minimize the l2 difference



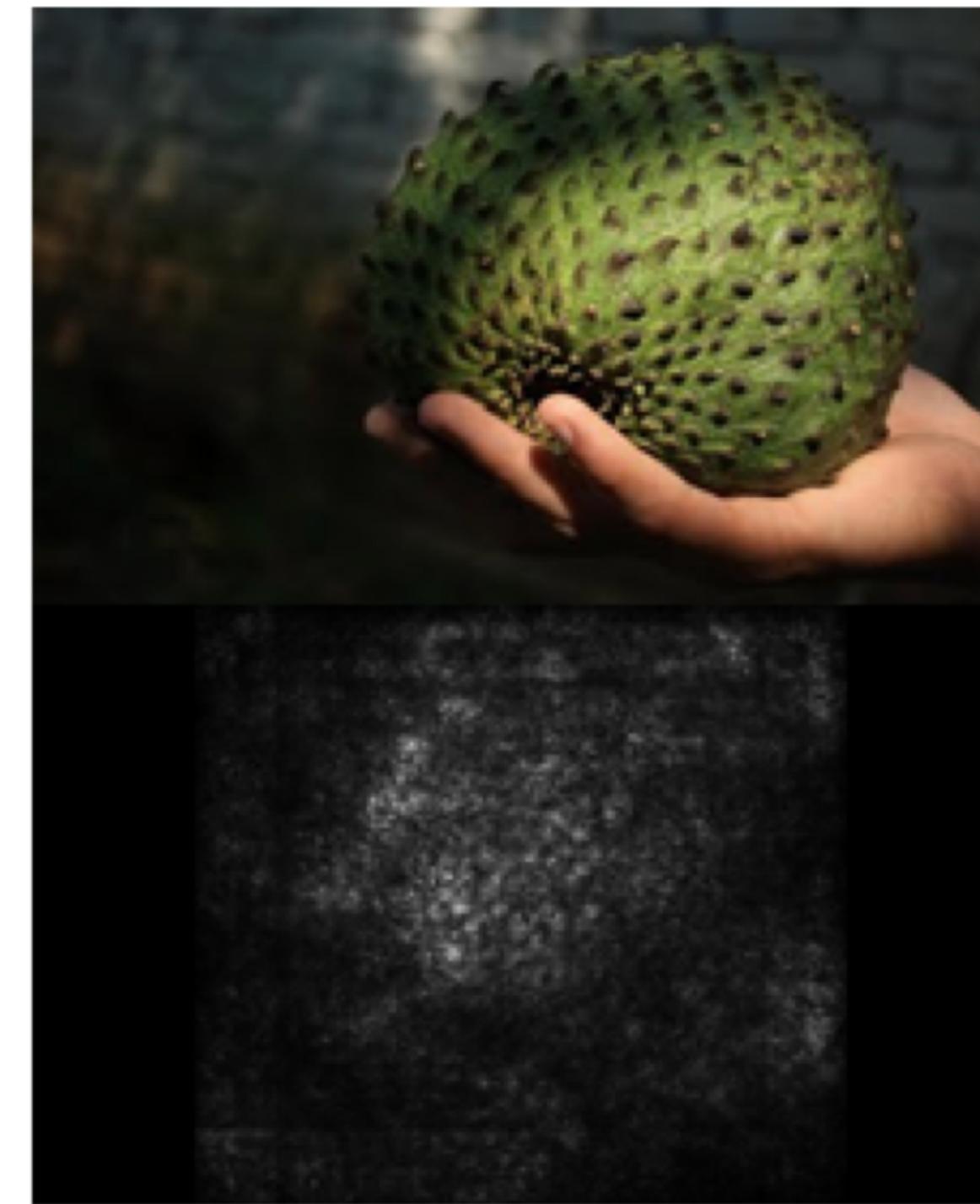
(a) Teacher and Student Networks



(b) Hints Training

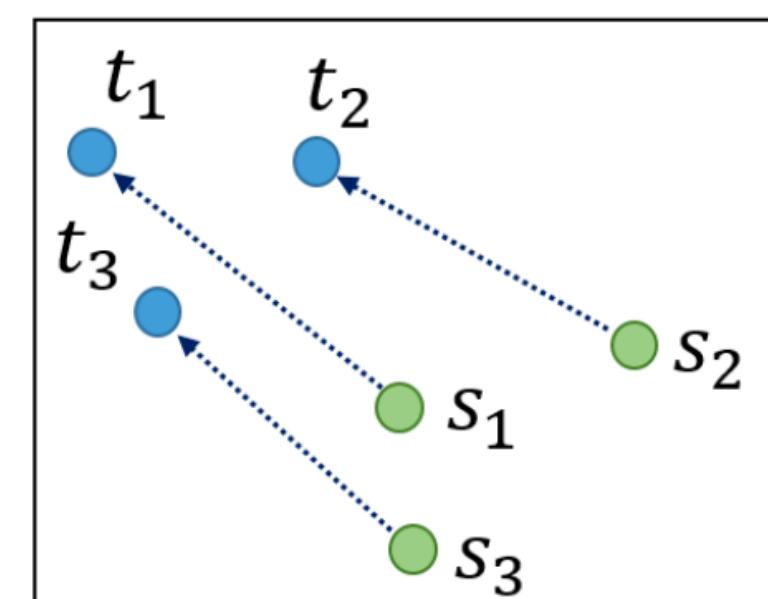
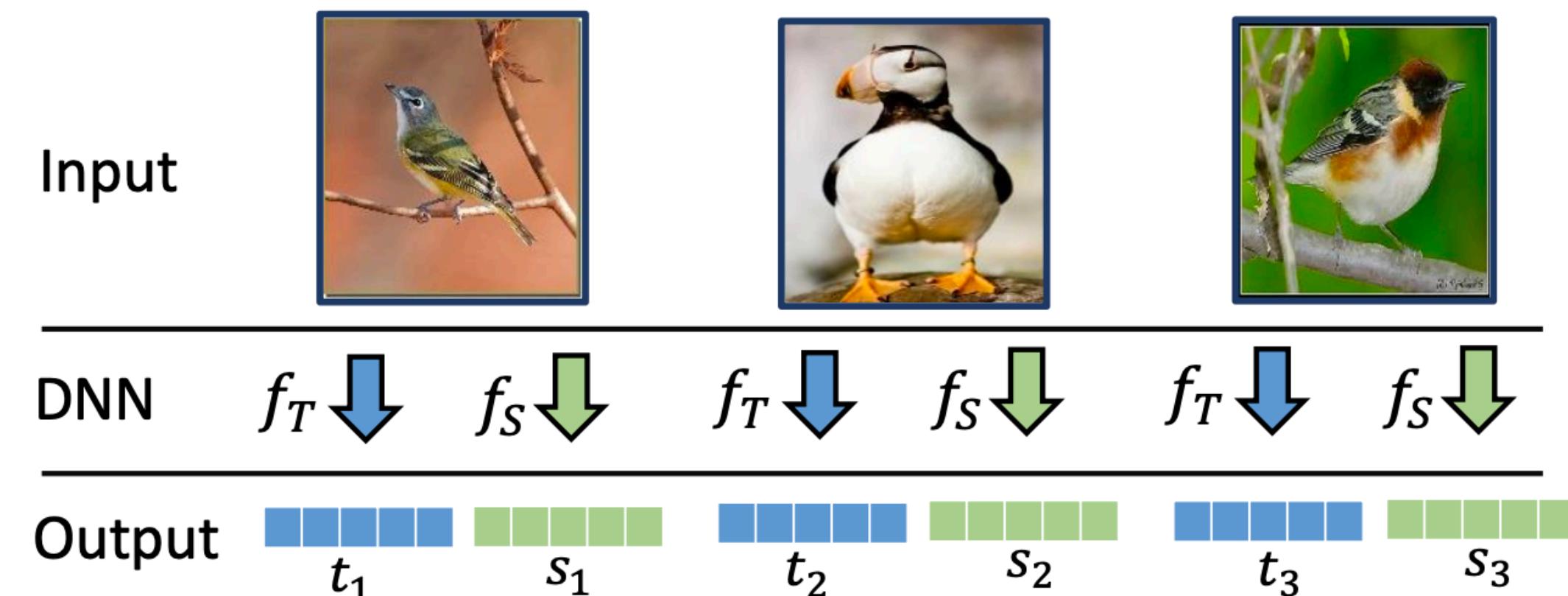
# Feature distillation

- Other variants match the **input gradients (saliency map)** of the models
  - Large gradient = The pixel affects the decision a lot
  - Models make predictions based on similar pixels



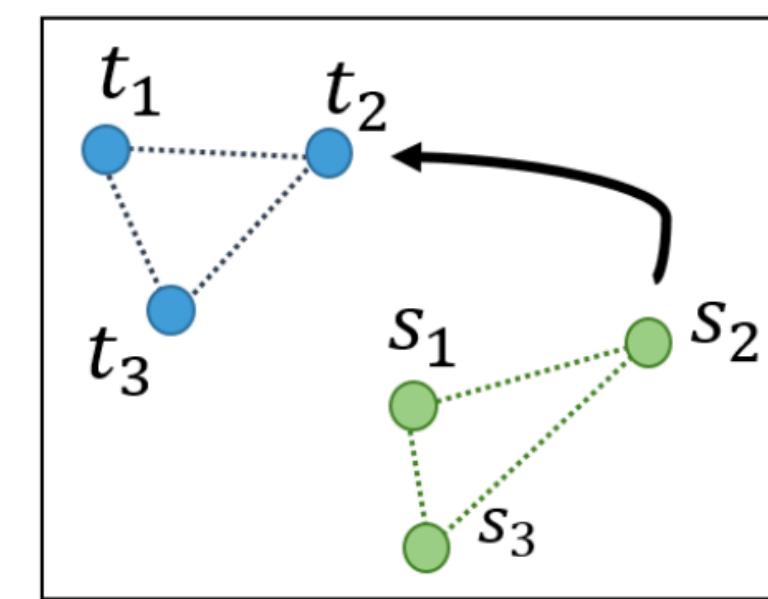
# Relational distillation

- Some algorithms try to transfer **inter-sample structures**.
- **Idea.** If teacher thinks two samples are similar, the student should too
  - Angle between samples
  - Normalized distance



Point to Point

**Conventional KD**



Structure to Structure

**Relational KD**

# Pros & cons

- **Logit**
  - Does not require internal information
    - thus applicable to closed-weight models
- **Feature, Relational**
  - Outperforms logit when carefully tuned
  - Natural choice for feature extractors

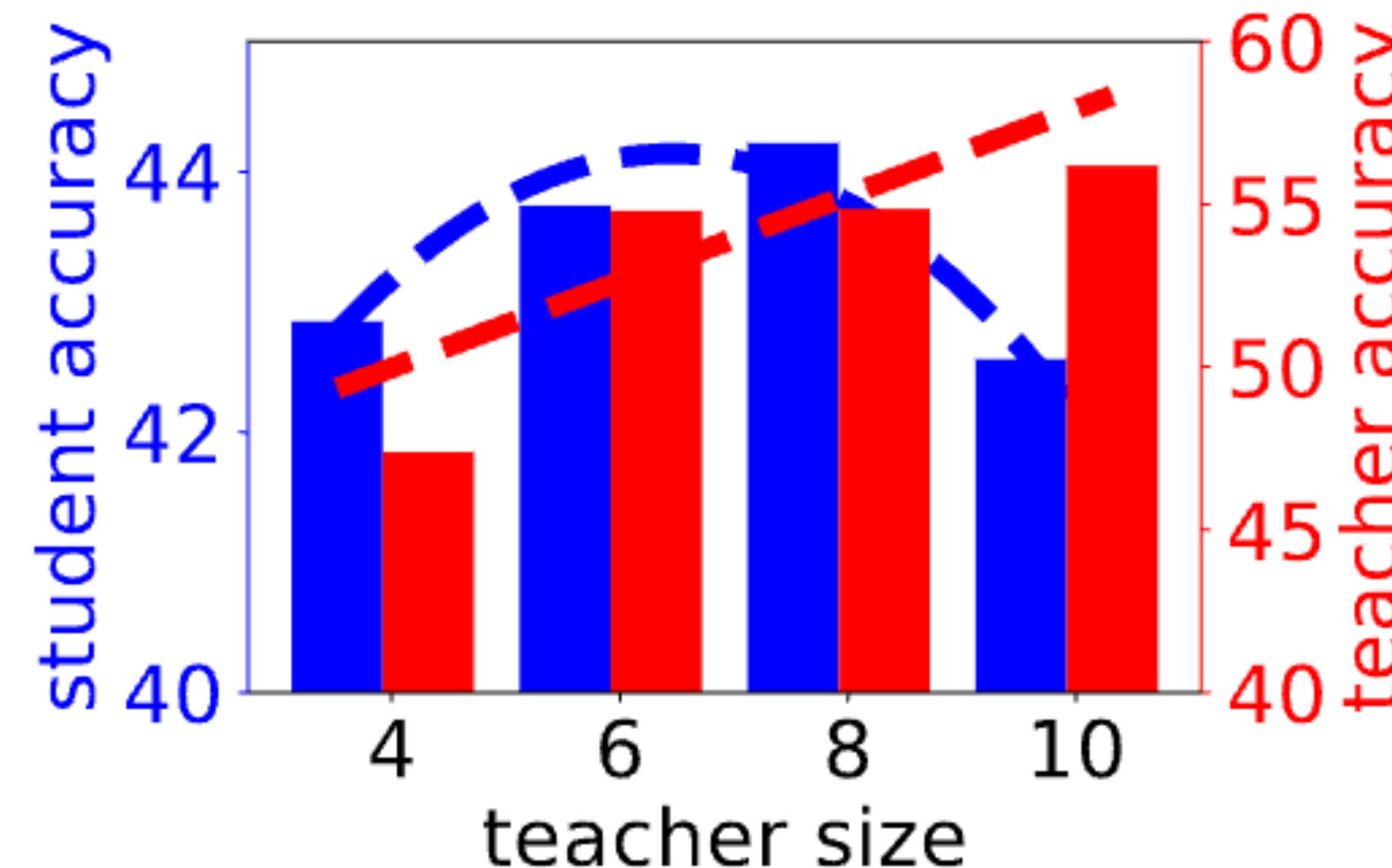
# Which teacher?

# Which teacher?

- Large model
  - Size discrepancy
- Ensemble of models
- Self-distillation
- Online distillation
- Uncompressed-to-compressed distillation

# Size discrepancy

- **Observation.** If the teacher is too large, the KD does not work well
  - Intuition. The regularizer is too strong for the small model to fit

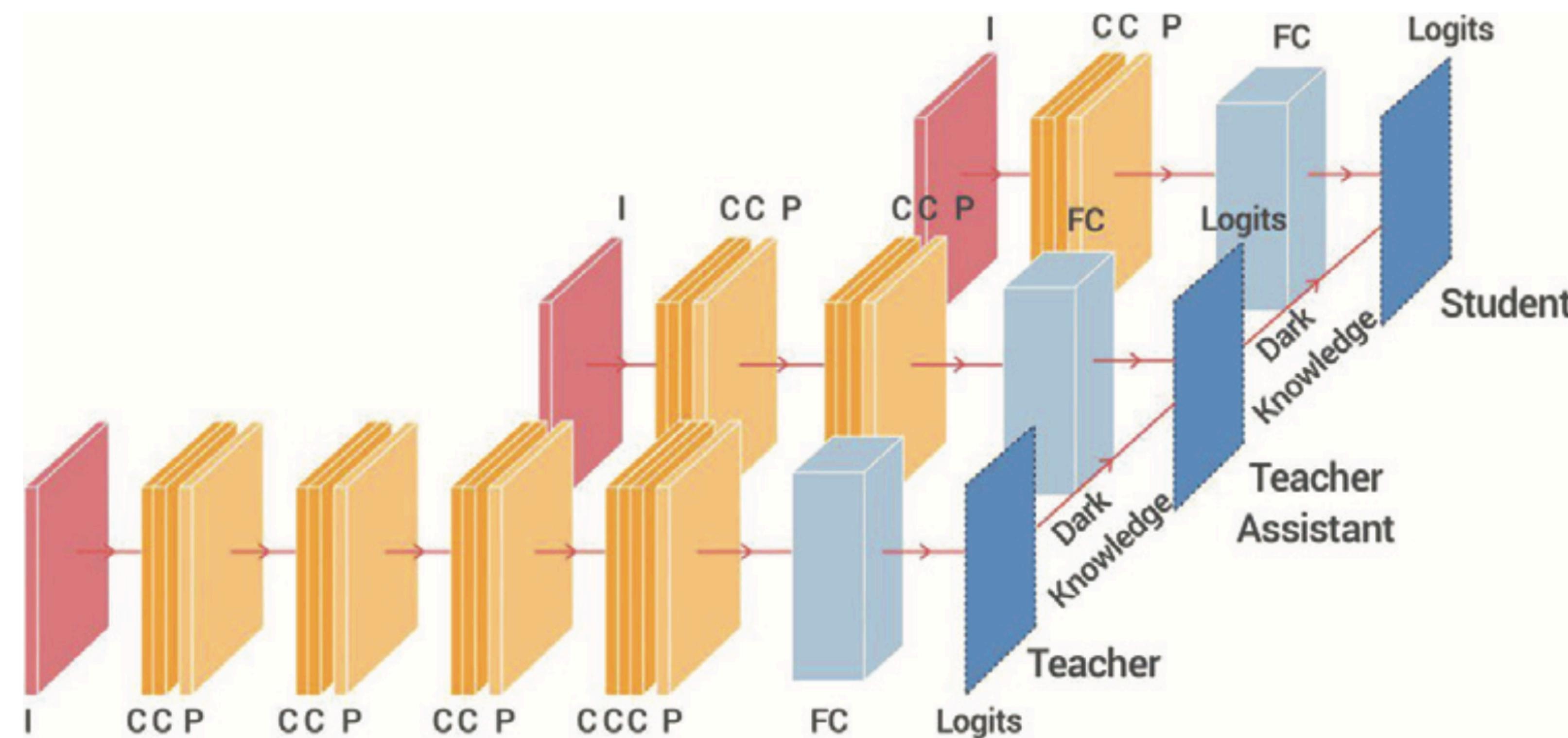


# Size discrepancy

- **Solution.** Use the “Teacher Assistant”

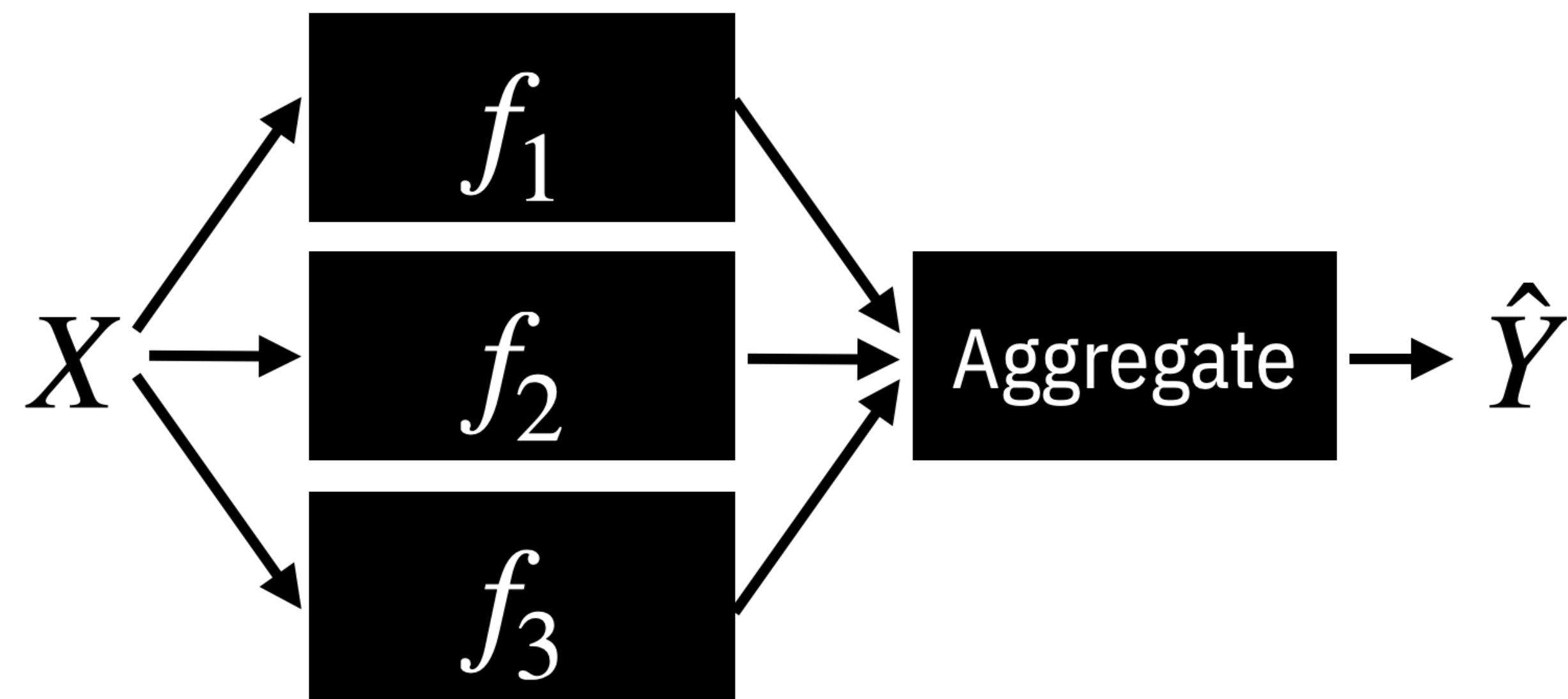
(1) Train an intermediate-sized TA, distilling teacher → TA

(2) Train a student, distilling TA → student



# Ensemble of models

- The most classical way to attain a teacher is to form an **ensemble**
  - Make predictions on multiple weak models
    - Averaging
    - Majority voting

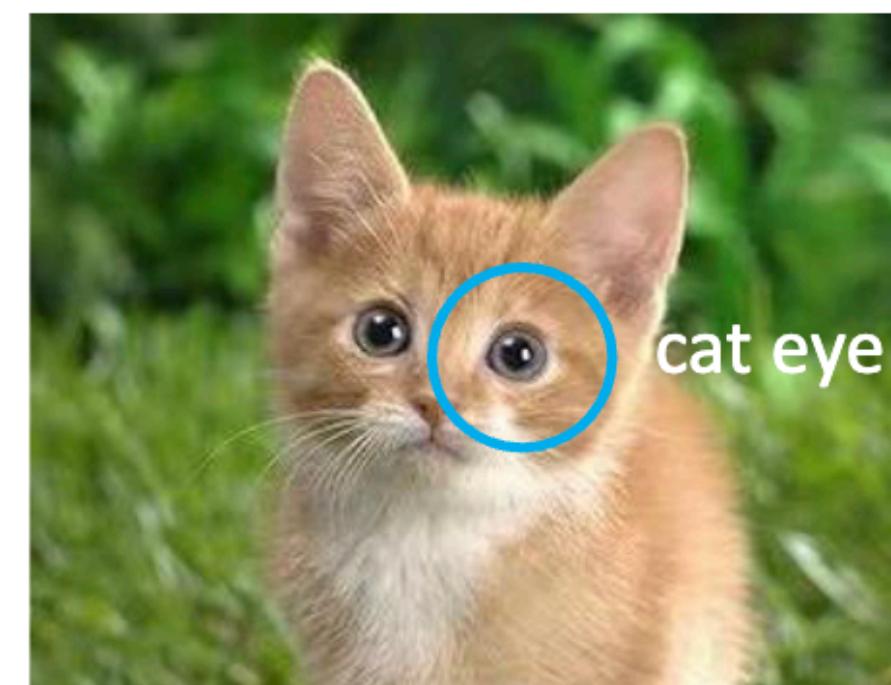


# Ensemble of models

- **Brain teaser.** Suppose that we have 3 binary classifiers with 70% accuracy. What is the accuracy of the majority-voted-ensemble?
  - (A) 70%
  - (B) 78.4%
  - (C) 100%
  - (D) 55%
  - (E) All of the above

# Self-distillation

- **Fact.** The teacher of the **same size** as the student can improve the student
- **Why?** Different random seed captures different “views” of a data
  - With seed 1, the model learns to predict based on **headlight**
  - With seed 42, the model learns to predict based on **window**
  - After distillation, see **headlight & window**



# Self-distillation

- Born-again NN. Distill from (an ensemble of) its prior versions

- Train  $f_1$  from scratch

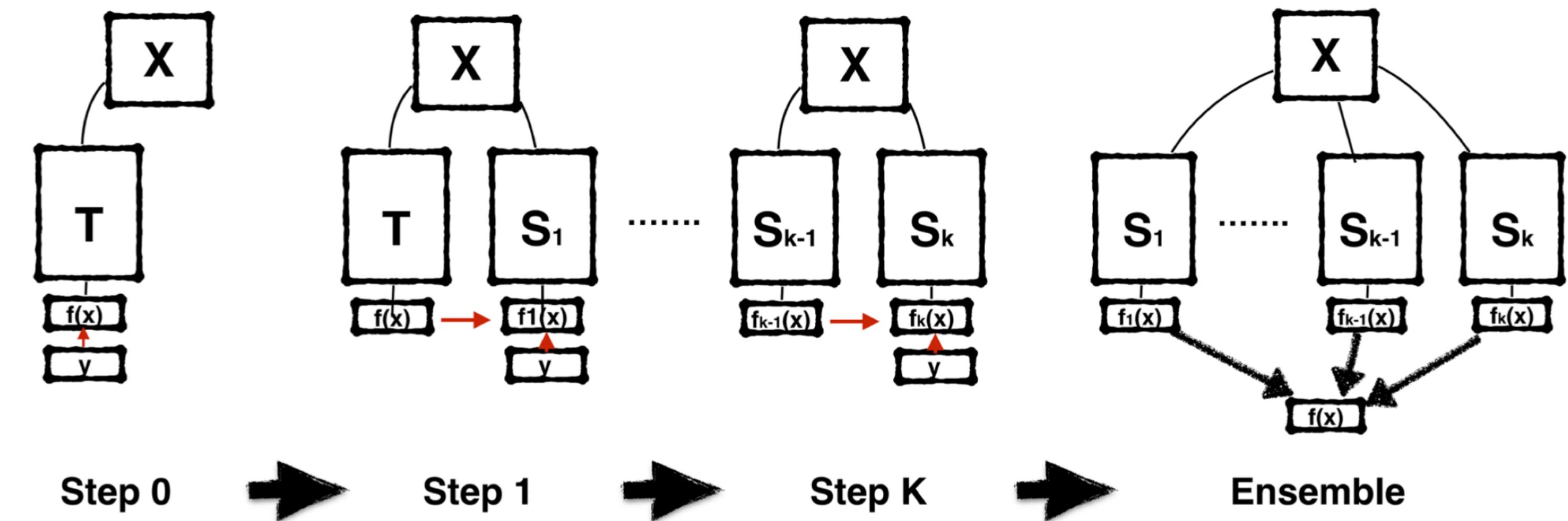
- Distill  $f_1 \rightarrow f_2$

- Distill  $f_2 \rightarrow f_3$

- (...)

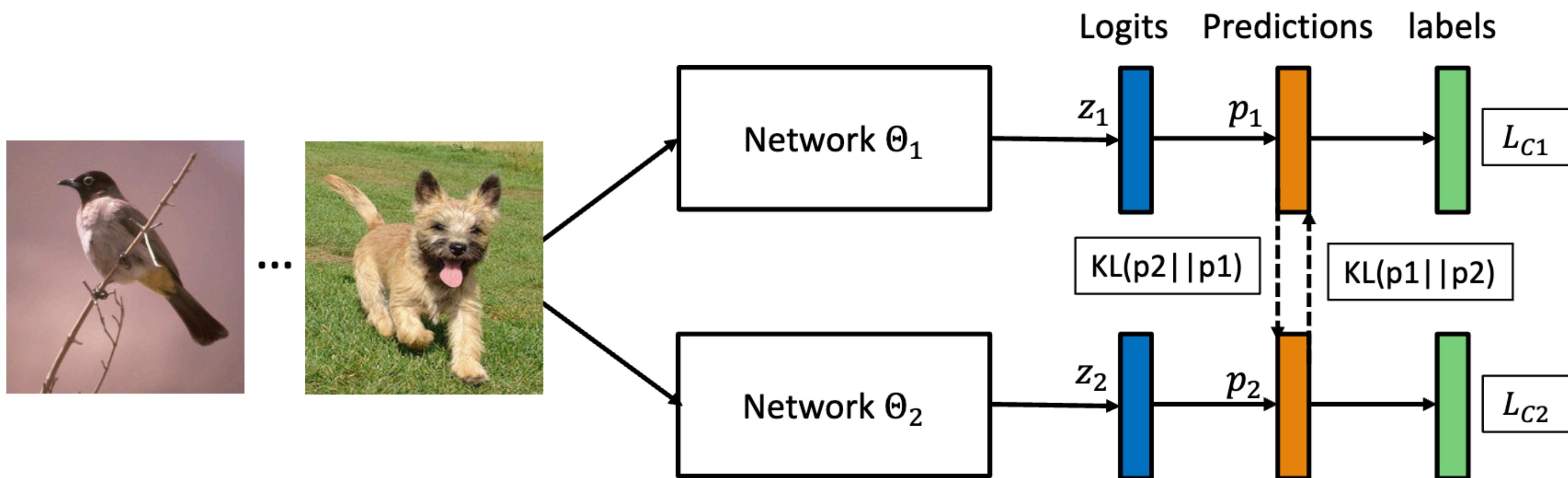
- Distill  $f_{k-1} \rightarrow f_k$

- Form an ensemble of  $f_1, \dots, f_k$



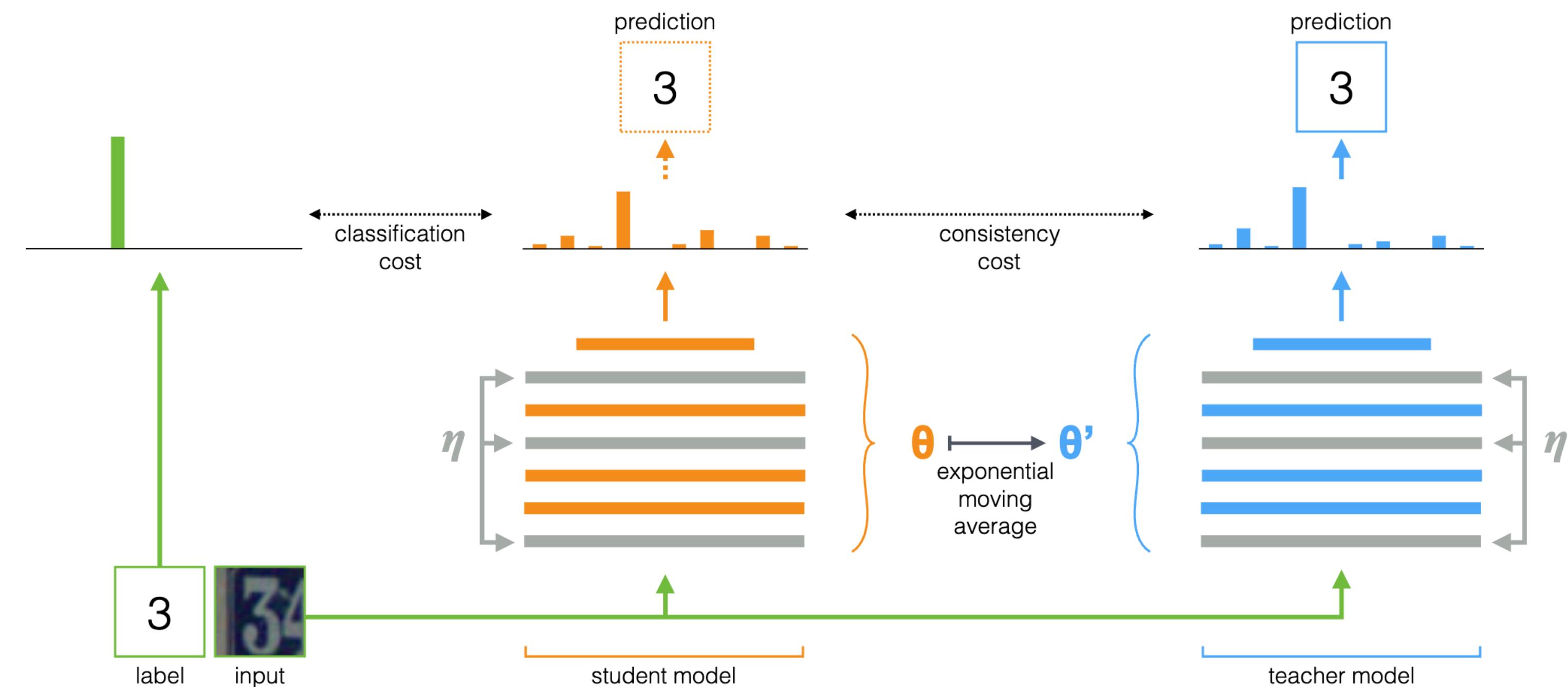
# Online distillation

- **Idea.** One can train a teacher and a student simultaneously
  - Both networks distill from each other
  - No pretrained teacher required!



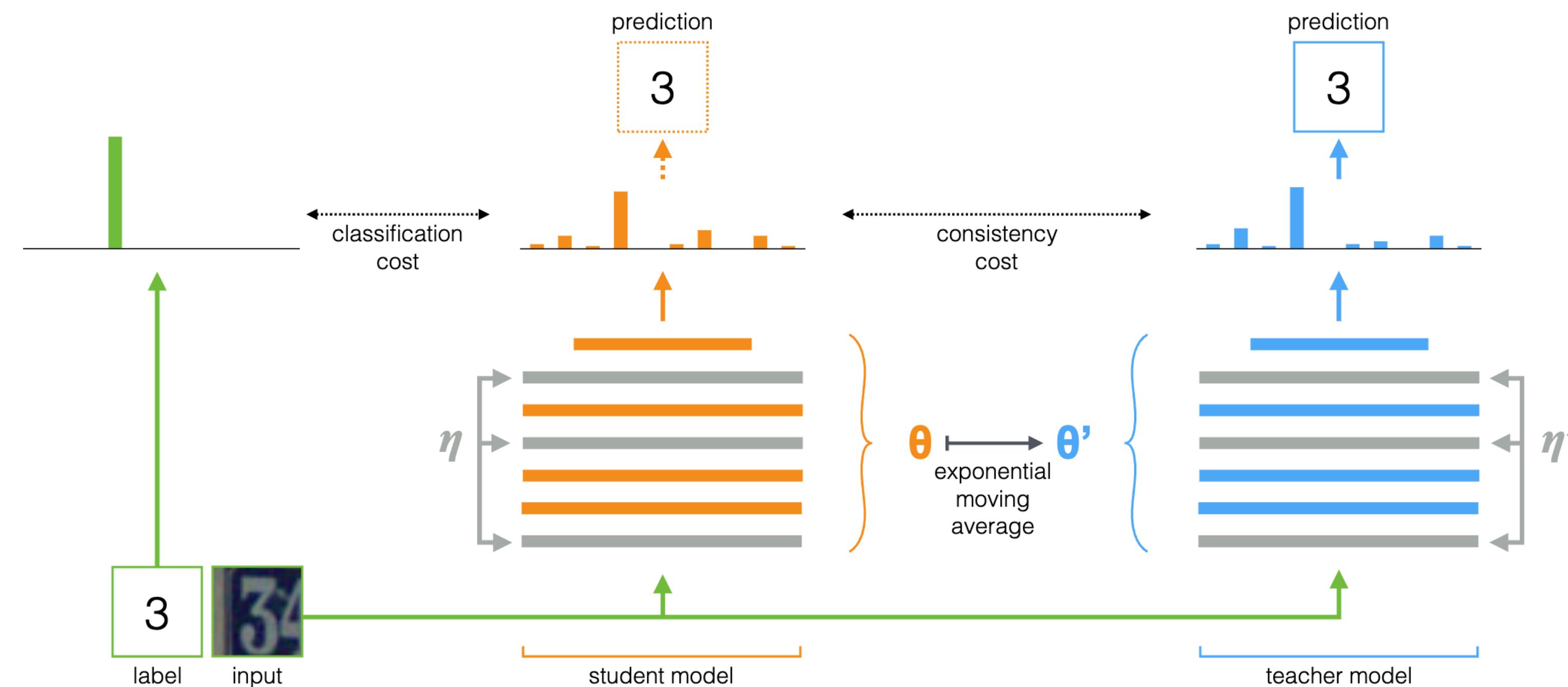
# Online distillation

- Even further, we can use the **momentum teacher**
  - Use the exponential moving average of the student as the teacher (typically known to generalize better than original one)
  - Typically applied together with noise for self-supervised learning



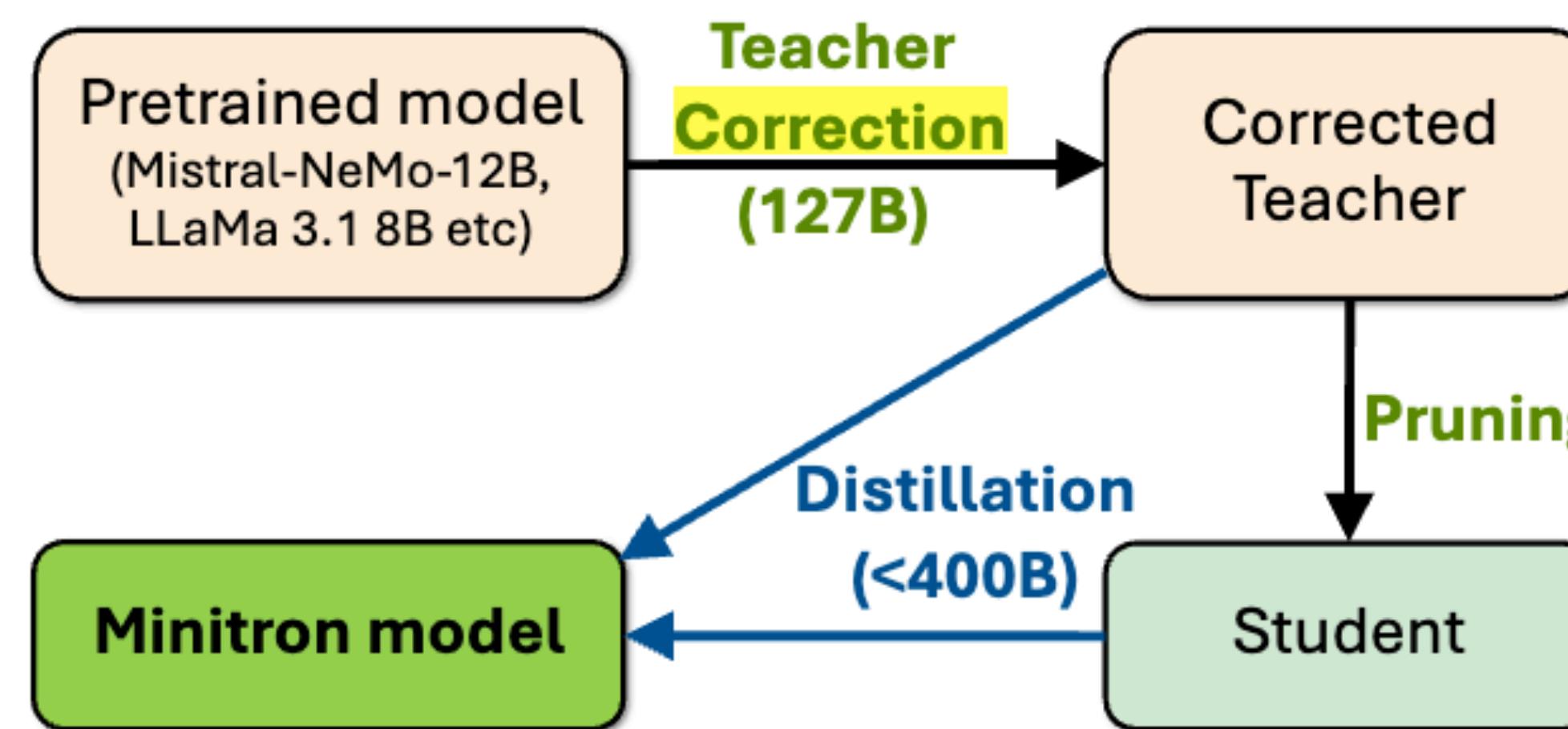
# Online distillation

- Even further, we can use the **momentum teacher**
  - Use the exponential moving average of the student as the teacher (typically known to generalize better than original one)
  - Typically applied together with noise for self-supervised learning



# Uncompressed-to-compressed

- A popular application is to combine with model compression
  - Example. Use dense model as teachers to retrain pruned models
    - **Teacher correction.** If original pretraining data is unavailable, fine-tuning on the distillation data is desirable (NVIDIA, 2024)



# Other tips

# Other tips

- Two other aspects that affect the performance of the student
  - Data augmentation
  - Patience

# Data augmentation

- It is essential that T & S look at the same views of the data
  - i.e., same augmented version

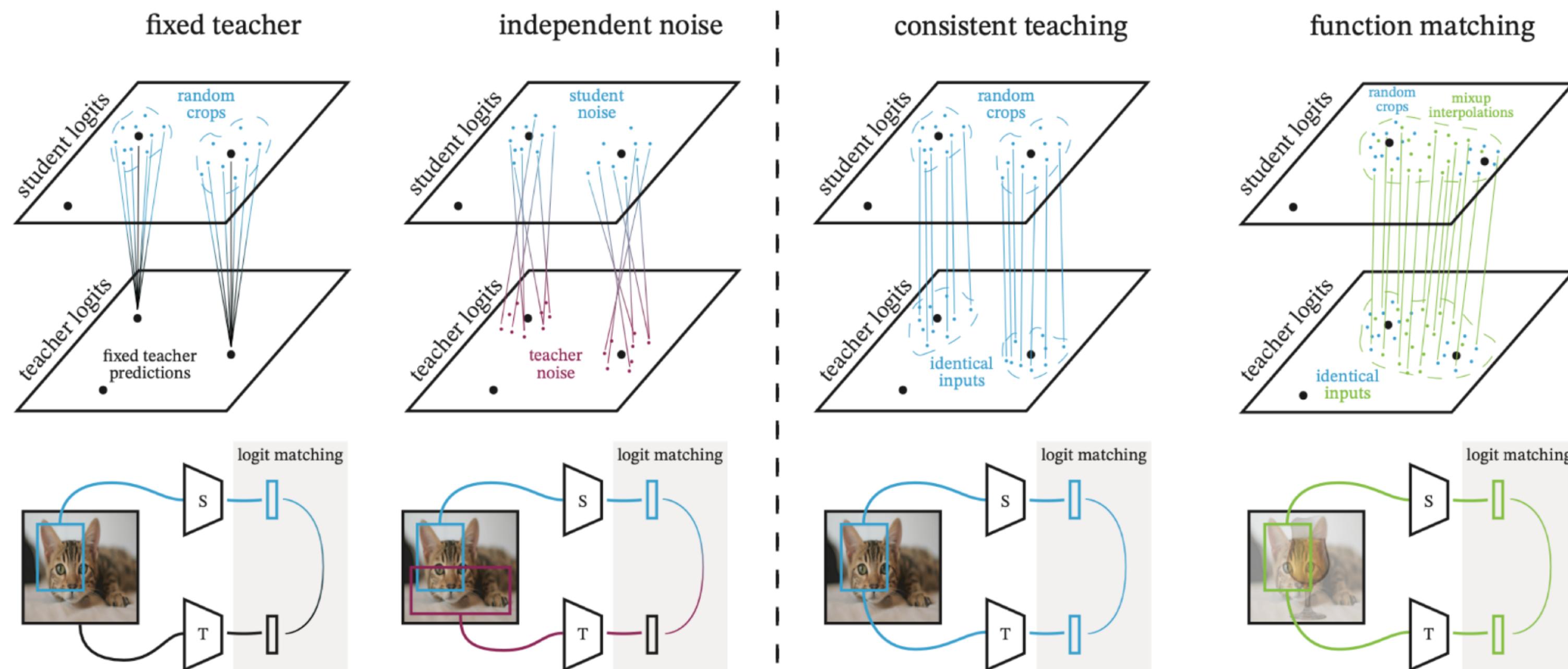
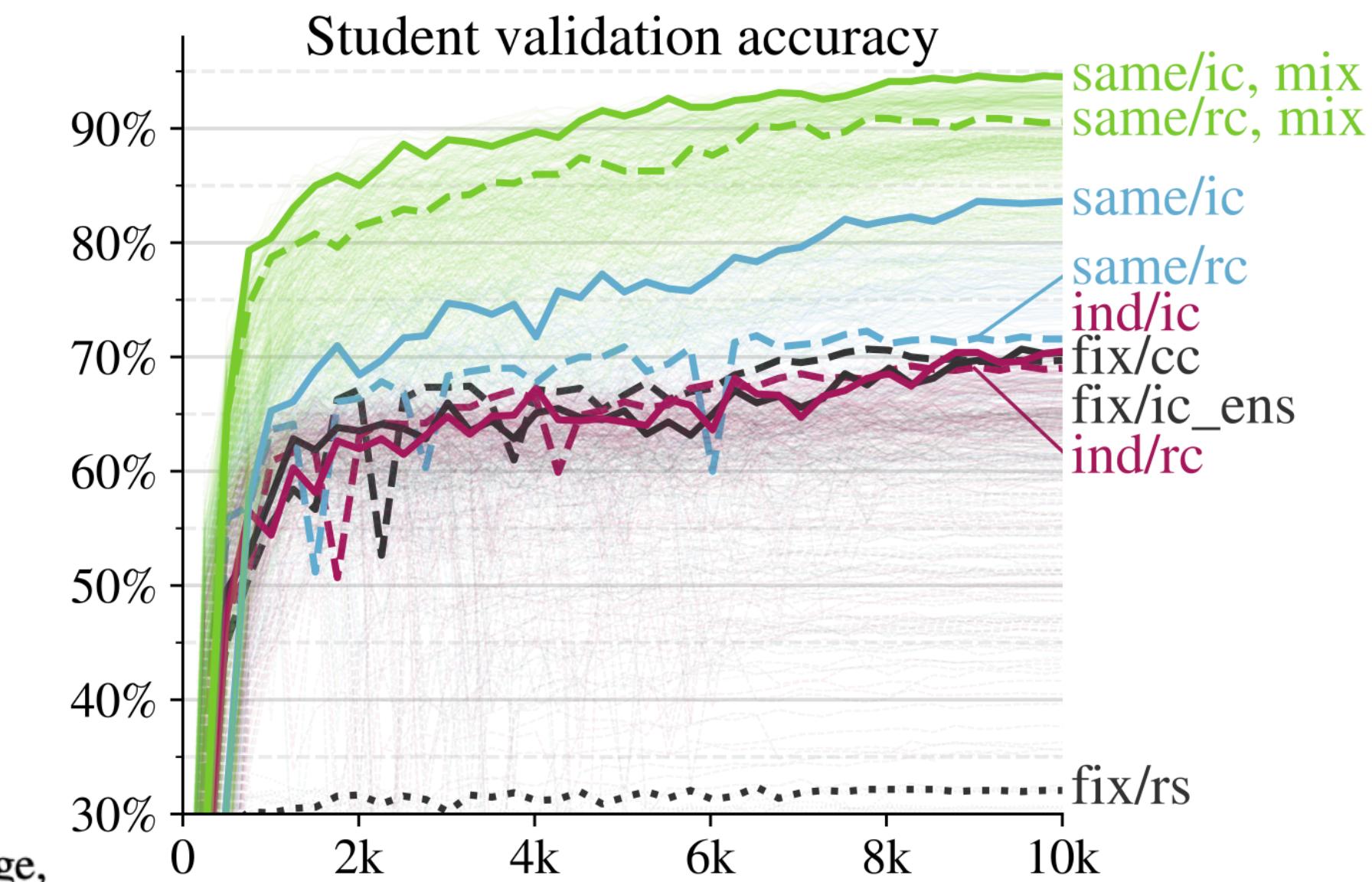
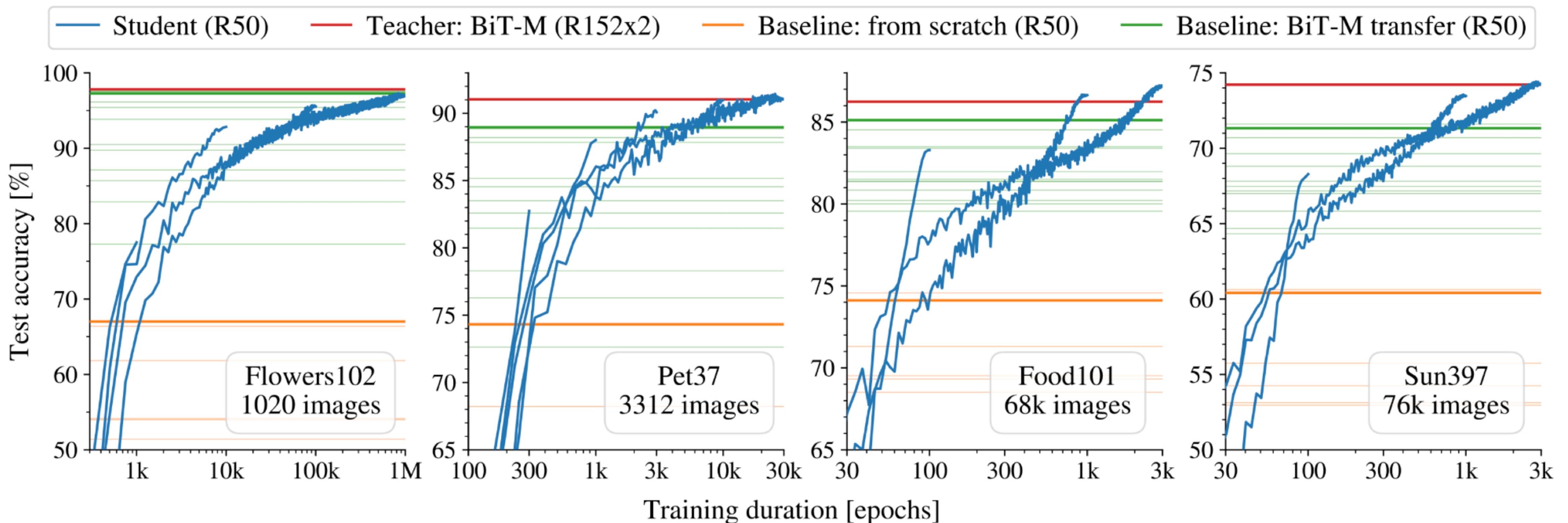


Figure 2. Schematic illustrations of various design choices when doing knowledge distillation. **Left:** *Teacher* receives a fixed image, while *student* receives a random augmentation. **Center-left:** *Teacher* and *student* receive independent image augmentations. **Center-right:** *Teacher* and *student* receive consistent image augmentations. **Right:** *Teacher* and *student* receive consistent image augmentations plus the input image manifold is extended by including linear segments between pairs of images (known as *mixup* [52] augmentation).



# Patience

- It often takes very long time to attain the full benefit of distillation
  - c.f. grokking



# Further readings

- **Distilling self-supervised models**
  - <https://arxiv.org/abs/2101.04731>
- **Data-free KD**
  - <https://arxiv.org/abs/1904.01186>

# Next Class

- Neural Architecture Search

That's it for today

