

Model Compression

Danna Gurari

University of Colorado Boulder

Fall 2022



<https://home.cs.colorado.edu/~DrG/Courses/NeuralNetworksAndDeepLearning/AboutCourse.html>

Review

- Last lecture topic:
 - Speech Processing – Problem and Applications
 - Speech Recognition – Evaluation and Models
 - Informal Retrieval – Problem and Applications
 - Informal Retrieval – Models
- Assignments
 - Final project presentations due in one week
- Questions?

Today's Topics

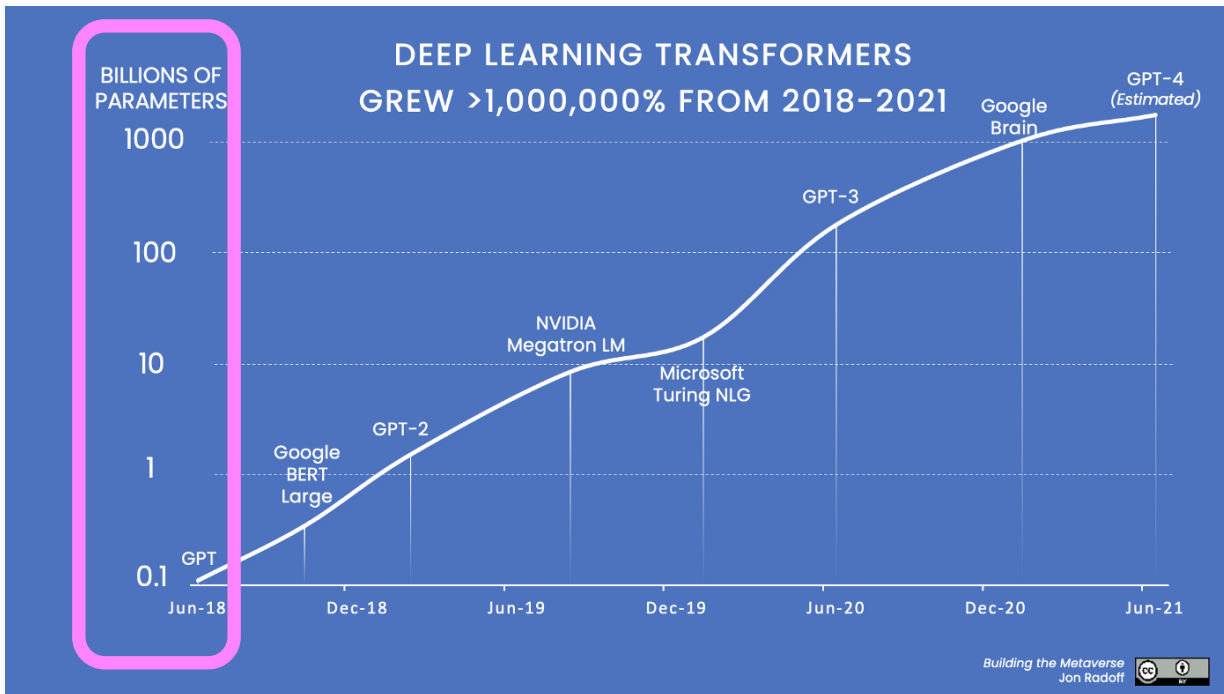
- Motivation
- Key idea: knowledge distillation
- Knowledge distillation for CNNs (vision problems)
- Knowledge distillation for Transformers (language problems)

Today's Topics

- Motivation
- Key idea: knowledge distillation
- Knowledge distillation for CNNs (vision problems)
- Knowledge distillation for Transformers (language problems)

Trend: Parameter-Heavy Models

Language – pretrained transformers



<https://medium.com/building-the-metaverse/the-metaverse-and-artificial-intelligence-ai-577343895411>

Vision – ImageNet classification

Architecture	Year	Top-1 Accuracy	# Parameters
DenseNet-169	2017	76.2%	14M
Inception-v3	2016	78.8%	24M
Inception-resnet-v2	2017	80.1%	56M
PolyNet	2017	81.3%	92M
SENet	2018	82.7%	146M
GPipe	2018	84.3%	557M
ResNeXt-101 32x48d	2019	85.4%	829M

Modern Neural Networks Are a Mismatch for Many Real-World Applications



<https://www.ephotozine.com/article/19-things-to-look-out-for-in-a-smartphone-camera--31055>



https://en.wikipedia.org/wiki/Wearable_technology



<https://www.buzzfeednews.com/article/katienotopoulos/facebook-is-making-camera-glasses-ha-ha-oh-no>

Modern Neural Networks Are a Mismatch for Many Real-World Applications

- Large inference time (i.e., incompatible for real-time applications)
- Large memory footprint (e.g., incompatible with limited memory on edge devices)
- Large computational cost (e.g., incompatible with limited battery on edge devices)
- Potential for large environmental costs

Idea: develop compact models so deep learning models can be used more efficiently and for more applications

Today's Topics

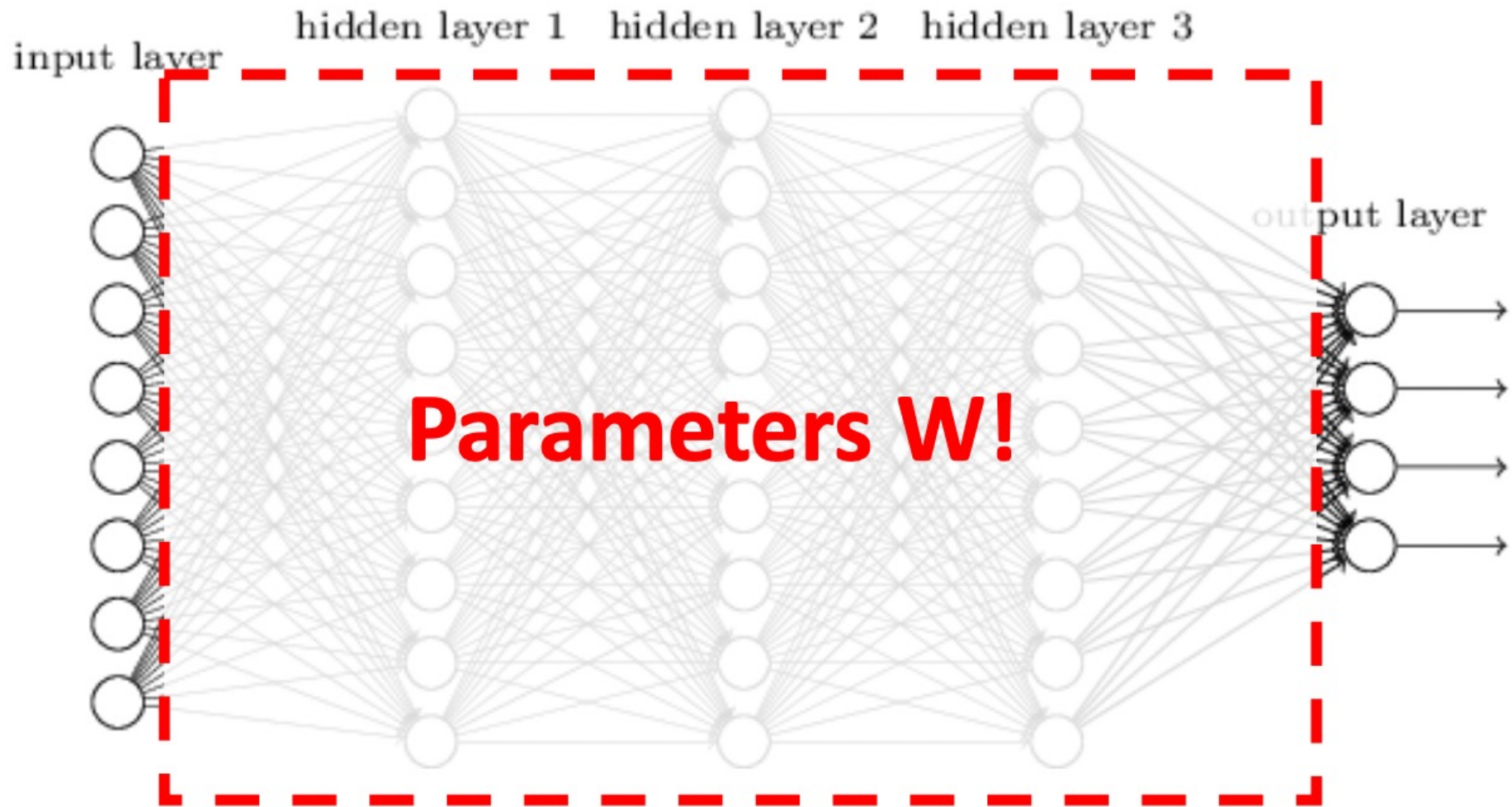
- Motivation
- Key idea: knowledge distillation
- Knowledge distillation for CNNs (vision problems)
- Knowledge distillation for Transformers (language problems)

Popular Approach: Knowledge Distillation

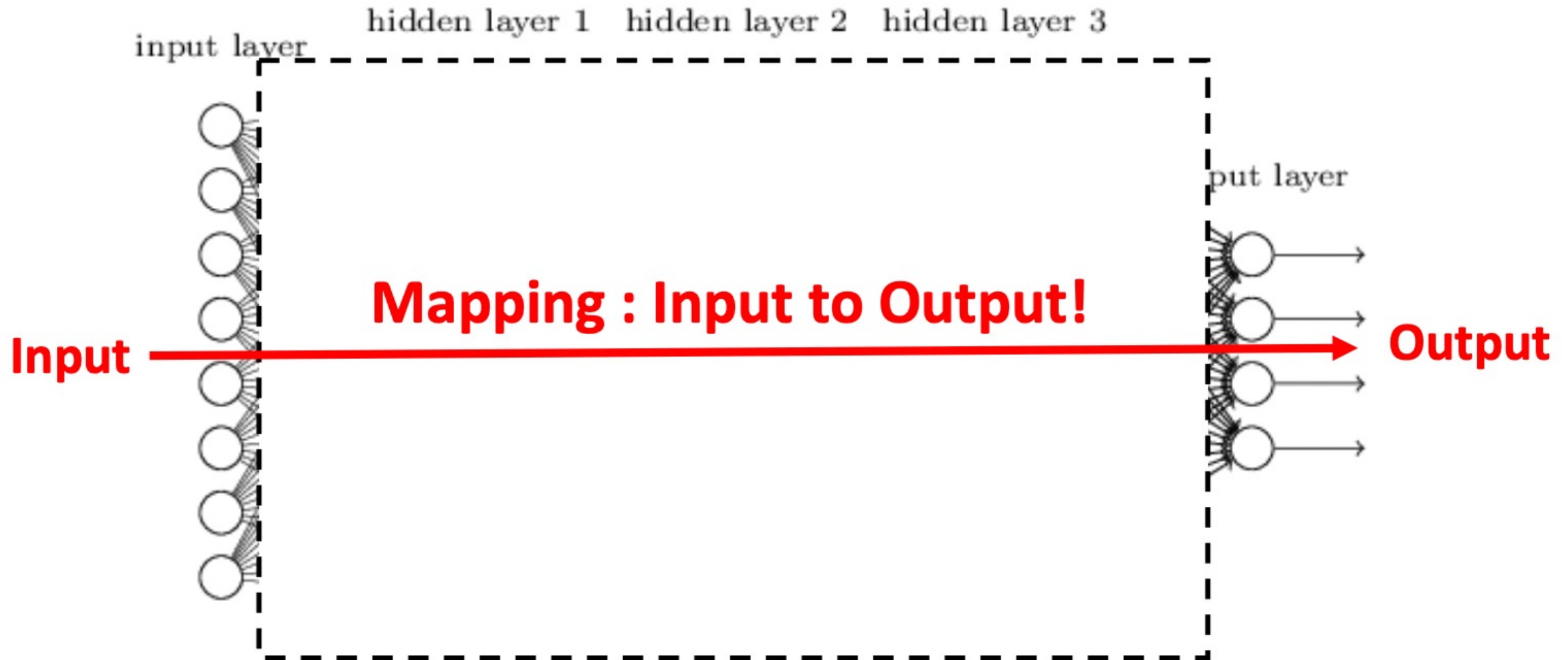


A student learns from a knowledgeable teacher

Key Question: What is Knowledge?

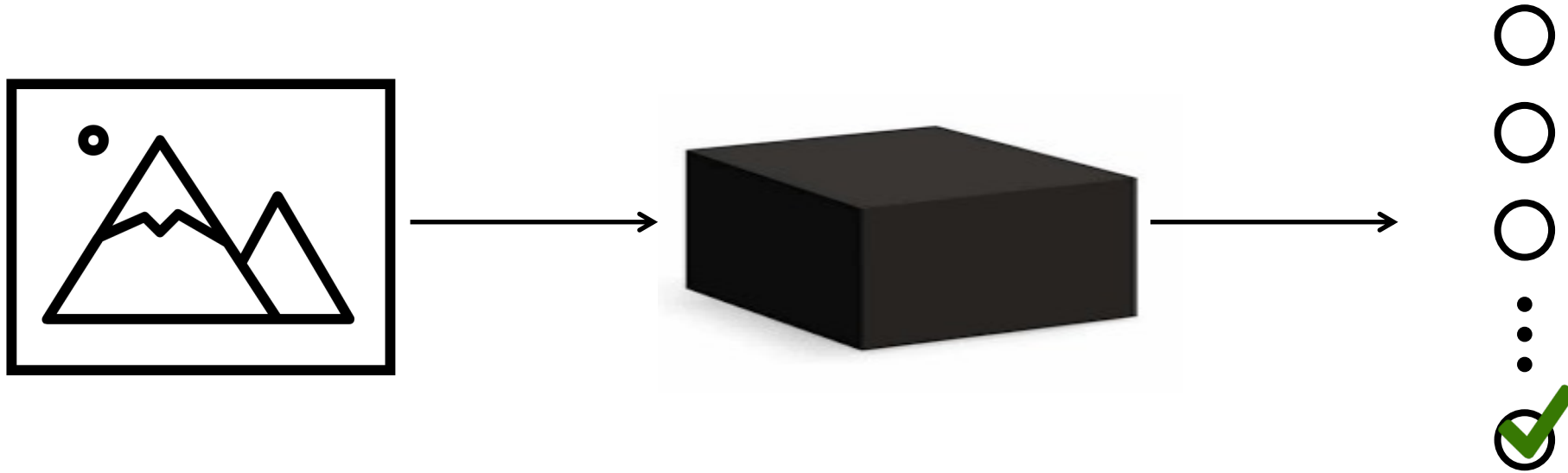


Knowledge Is: Input to Output Mapping



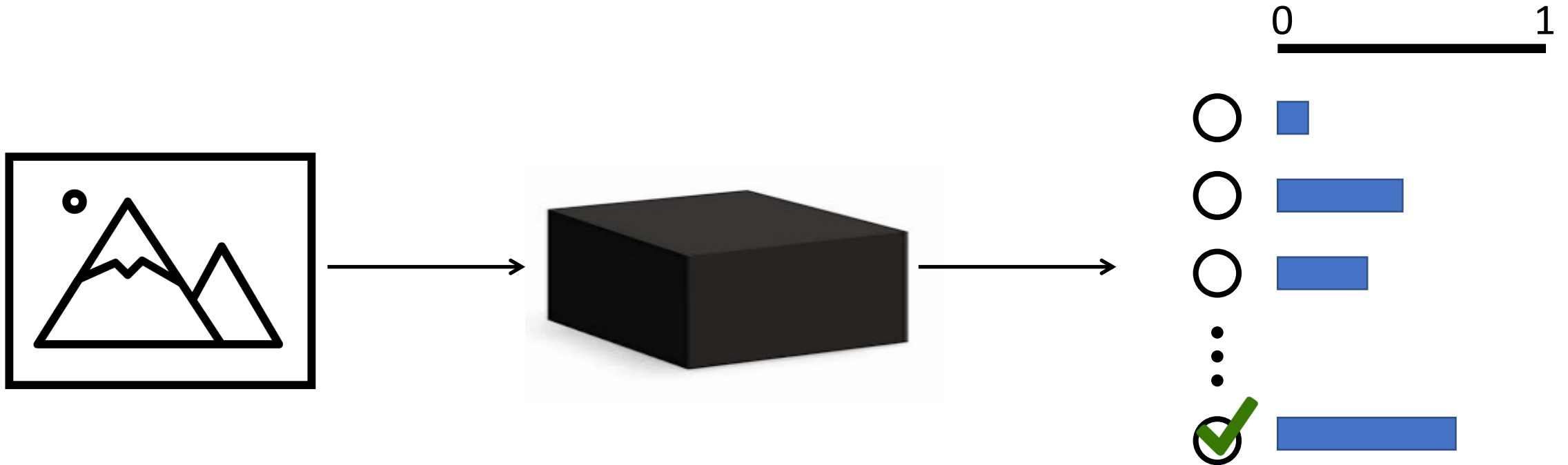
Knowledge Is: Input to Output Mapping

Target mapping: ground truth (1-hot vector)



Knowledge Is: Input to Output Mapping

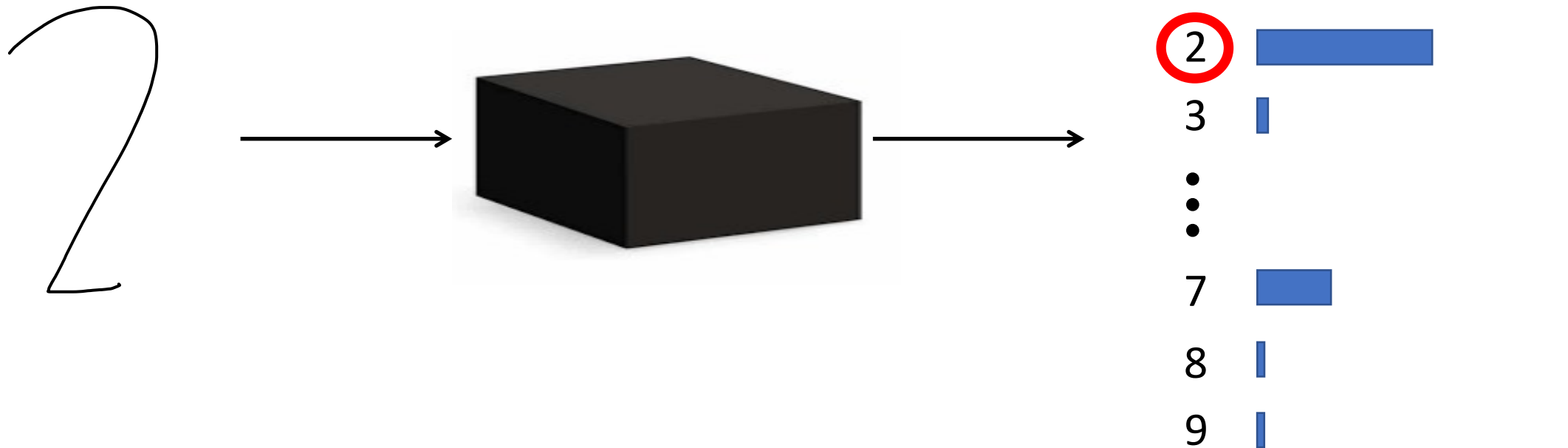
Target mapping: probability distribution from a model offers
further insights into similarities and differences of categories



Knowledge Is: Input to Output Mapping

Target mapping: probability distribution from a model offers
further insights into similarities and differences of categories

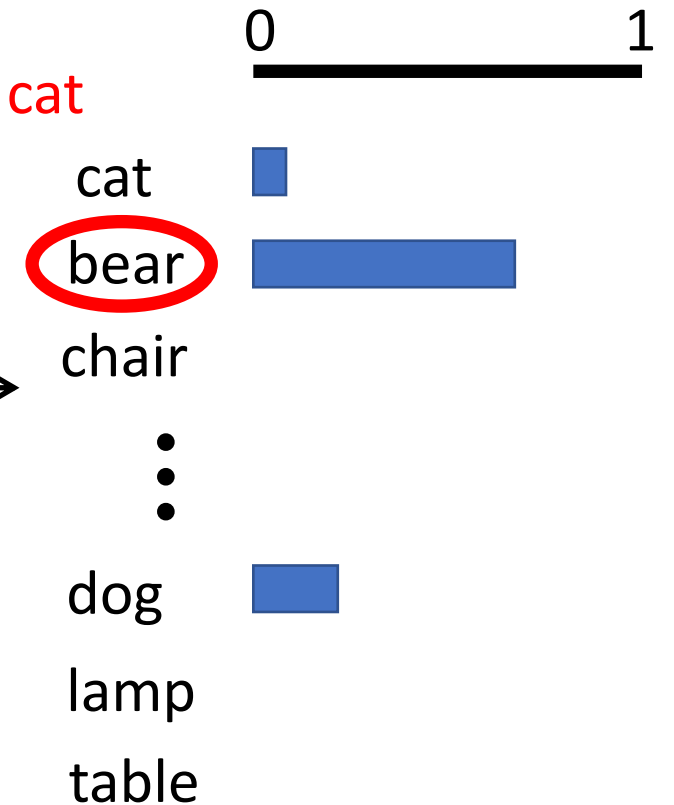
- Attempts to identify ground truth category
- Also, shares that 2 has similar characteristics to 7 and 1



Knowledge Is: Input to Output Mapping

Target mapping: probability distribution from a model offers
further insights into similarities and differences of categories

- Attempts to identify ground truth category
- Also, shares that bear has similar characteristics to dog and cat



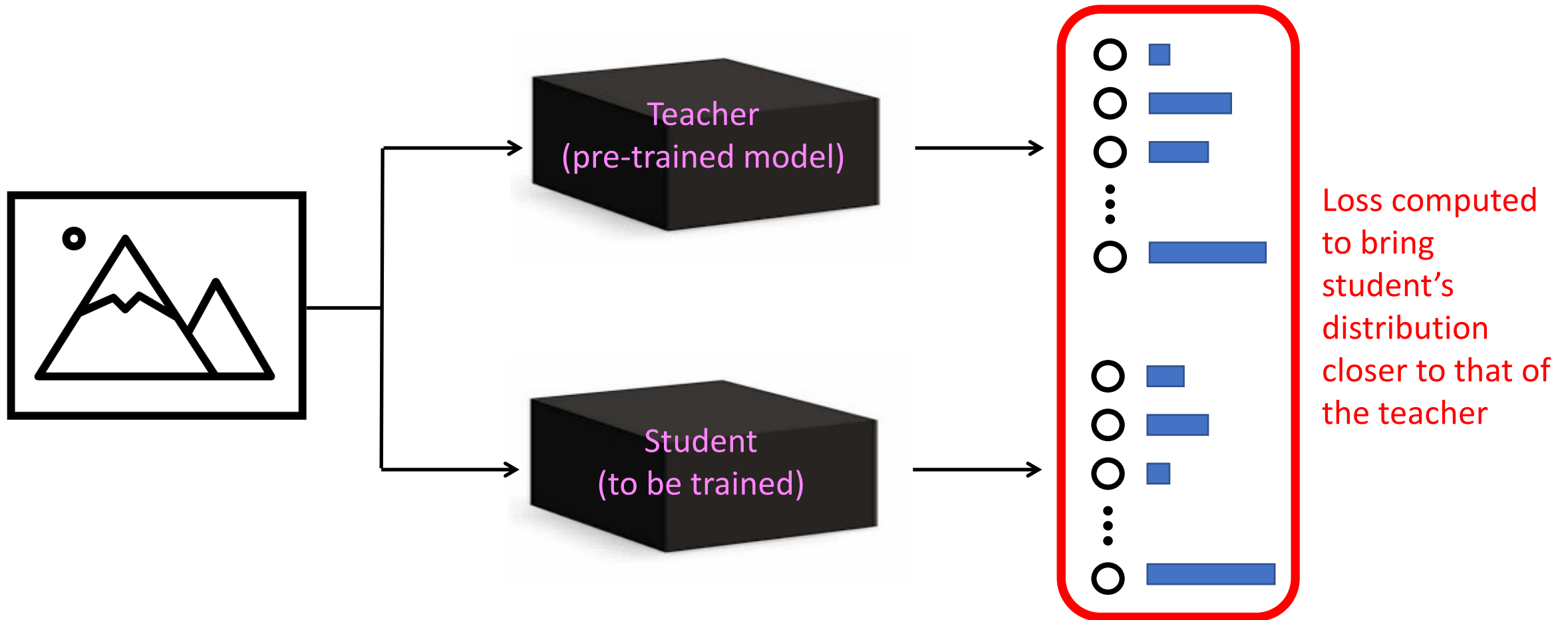
Knowledge Is: Input to Output Mapping

Target mapping: probability distribution from a model offers further insights into similarities and differences of categories

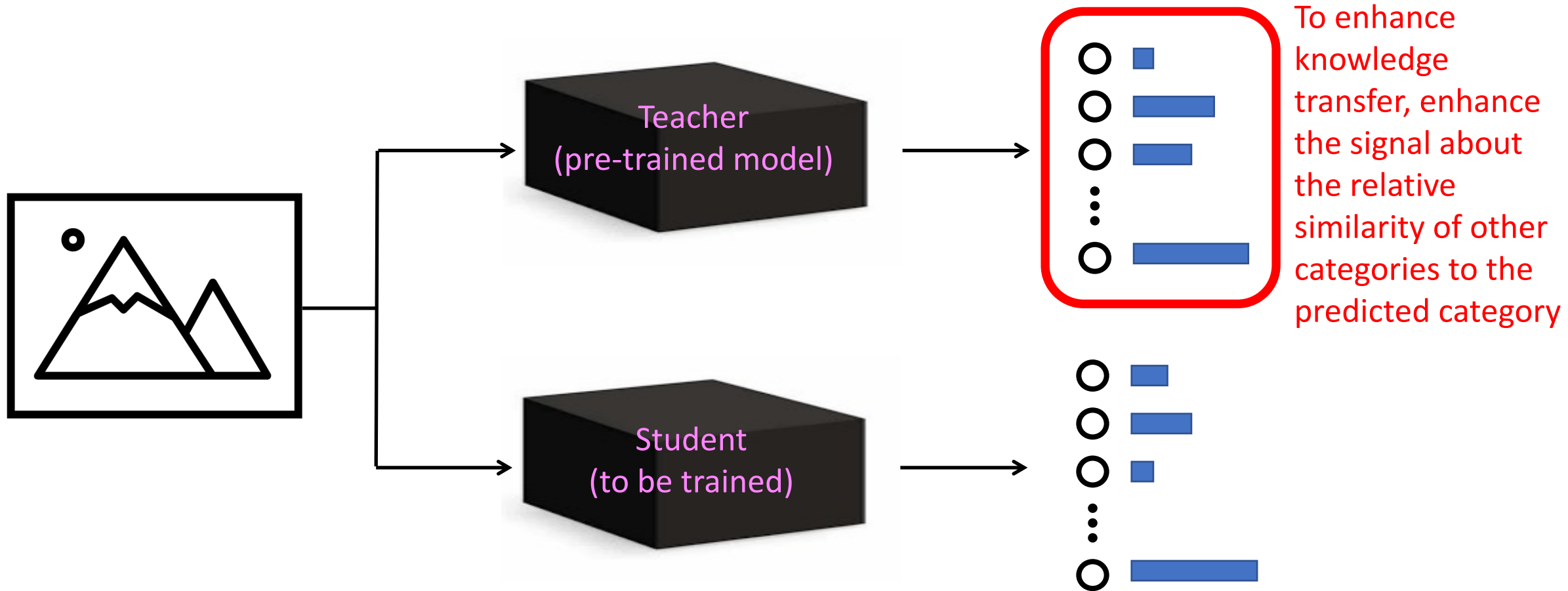
- Attempts to identify ground truth category
- Also, shares that bear has similar characteristics to dog and cat

Idea: not only teach a compact student about the ground truth, but also teach it about the relationships between categories

Knowledge Distillation: Teach Student the “Dark Knowledge” of Teacher

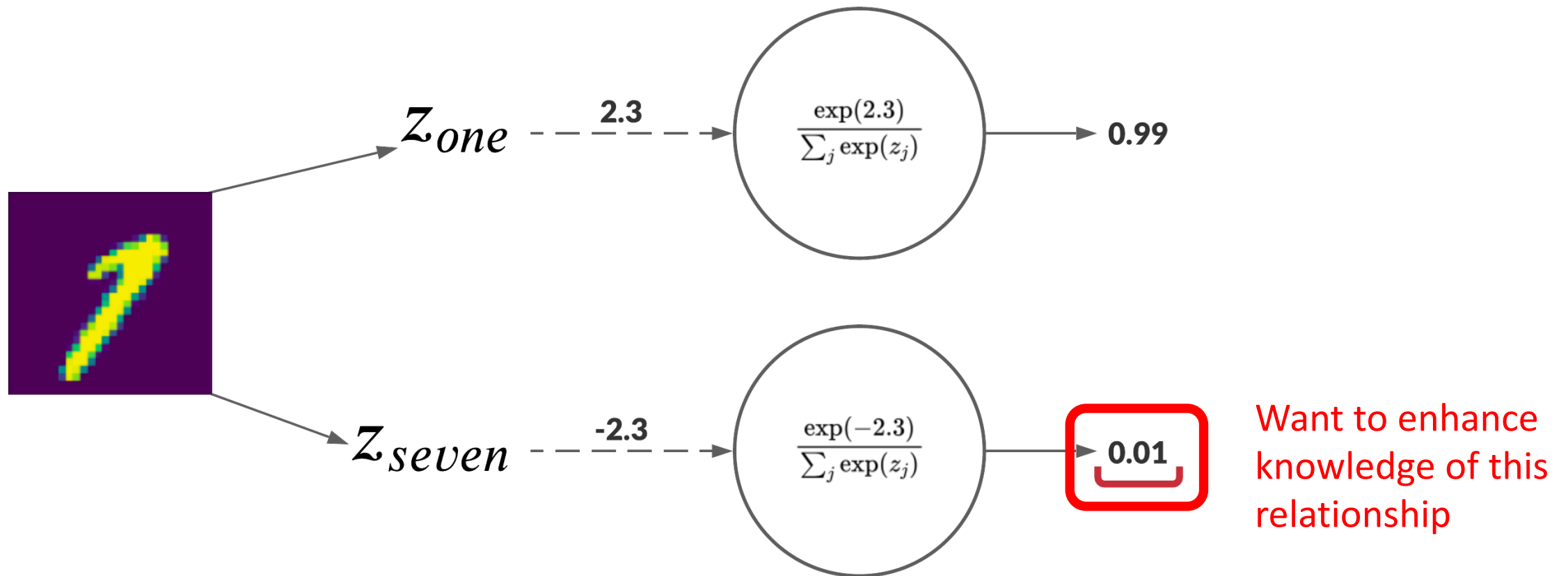


Knowledge Distillation: Teach Student the “Dark Knowledge” of Teacher



Knowledge Distillation: Rebalance (“Soften”) Probability Distribution Across Categories

Recall Softmax: converts vector of **scores** into a probability distribution that sums to 1



Knowledge Distillation: Rebalance (“Soften”) Probability Distribution Across Categories

Recall Softmax: converts vector of **scores** into a probability distribution that sums to 1

Get rid of negative values while preserving original order of scores

$$\sigma(\mathbf{z})_i =$$

$i = 1, \dots, K$

$$\sum_{j=1}^K e^{z_j}$$

Number of classes

Divide each node's score by sum of all entries to make them sum to 1 (normalization)

Knowledge Distillation: Rebalance (“Soften”) Probability Distribution Across Categories

Generalized Softmax: converts vector of **scores** into a probability distribution that sums to 1 with **temperature**

$$\sigma(\mathbf{z})_i = \frac{\exp(z_i / T)}{\sum_j \exp(z_j / T)}$$

What is the typical value of T used for softmax?

Idea: set the temperature to a value greater than 1

Knowledge Distillation: Rebalance (“Soften”) Probability Distribution Across Categories

Generalized Softmax: converts vector of **scores** into a probability distribution that sums to 1 with **temperature**

$$\sigma(\mathbf{z})_i = \frac{\exp(z_i / T)}{\sum_j \exp(z_j / T)}$$

Larger T values means more information is available about which categories the teacher found similar to the predicted category

Knowledge Distillation: Rebalance (“Soften”) Probability Distribution Across Categories

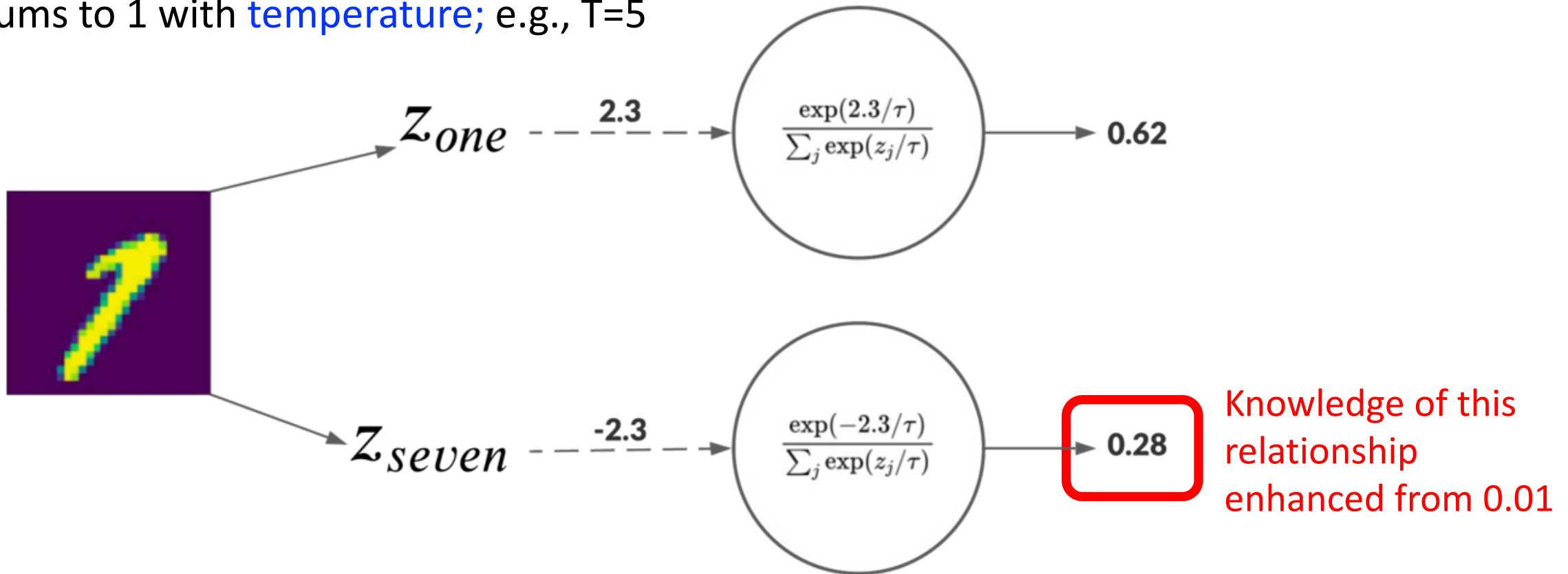
Generalized Softmax: converts vector of **scores** into a probability distribution that sums to 1 with **temperature**; e.g.,

0.997	Homework	0.935	Homework	0.637	Homework
0.000	Cake	0.0001	Cake	0.021	Cake
0.002	Book	0.046	Book	0.191	Book
0.001	Assignment	0.017	Assignment	0.128	Assignment
0.000	Car	0.0001	Car	0.021	Car
T=1		T=2		T=5	

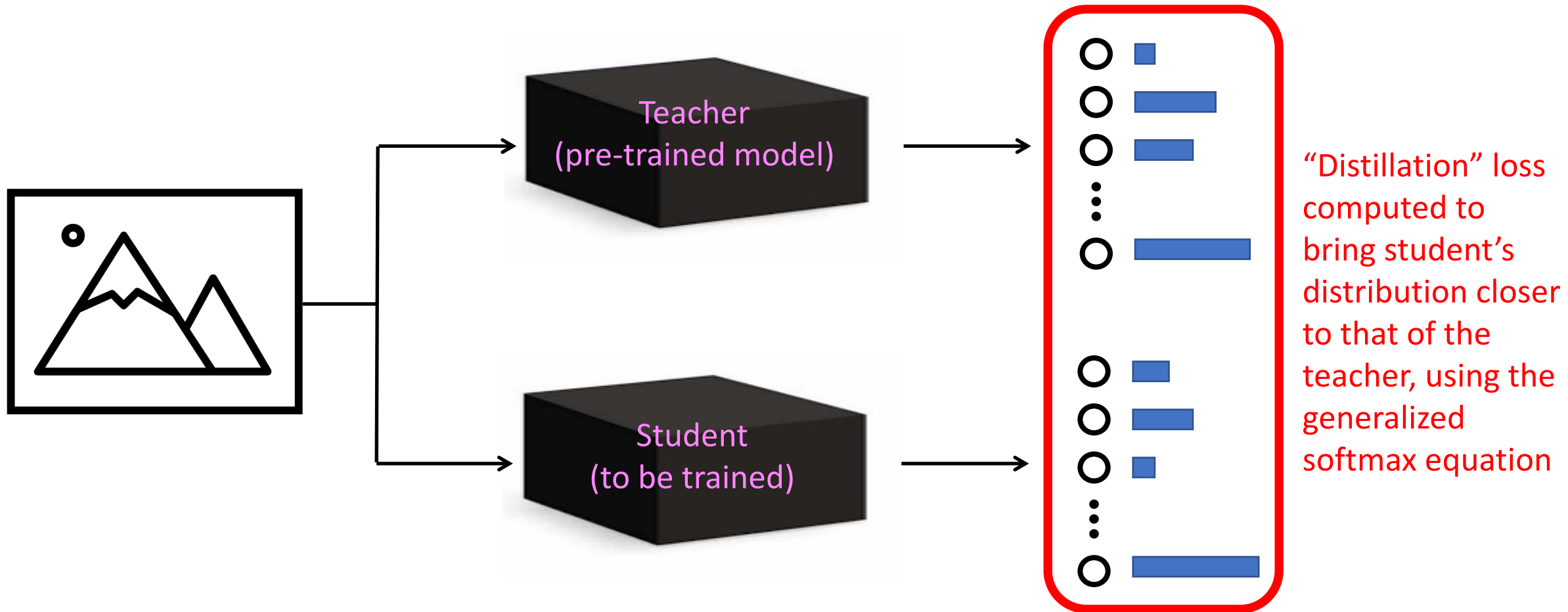
Larger T values means more information is available about which categories the teacher found similar to the predicted category

Knowledge Distillation: Rebalance (“Soften”) Probability Distribution Across Categories

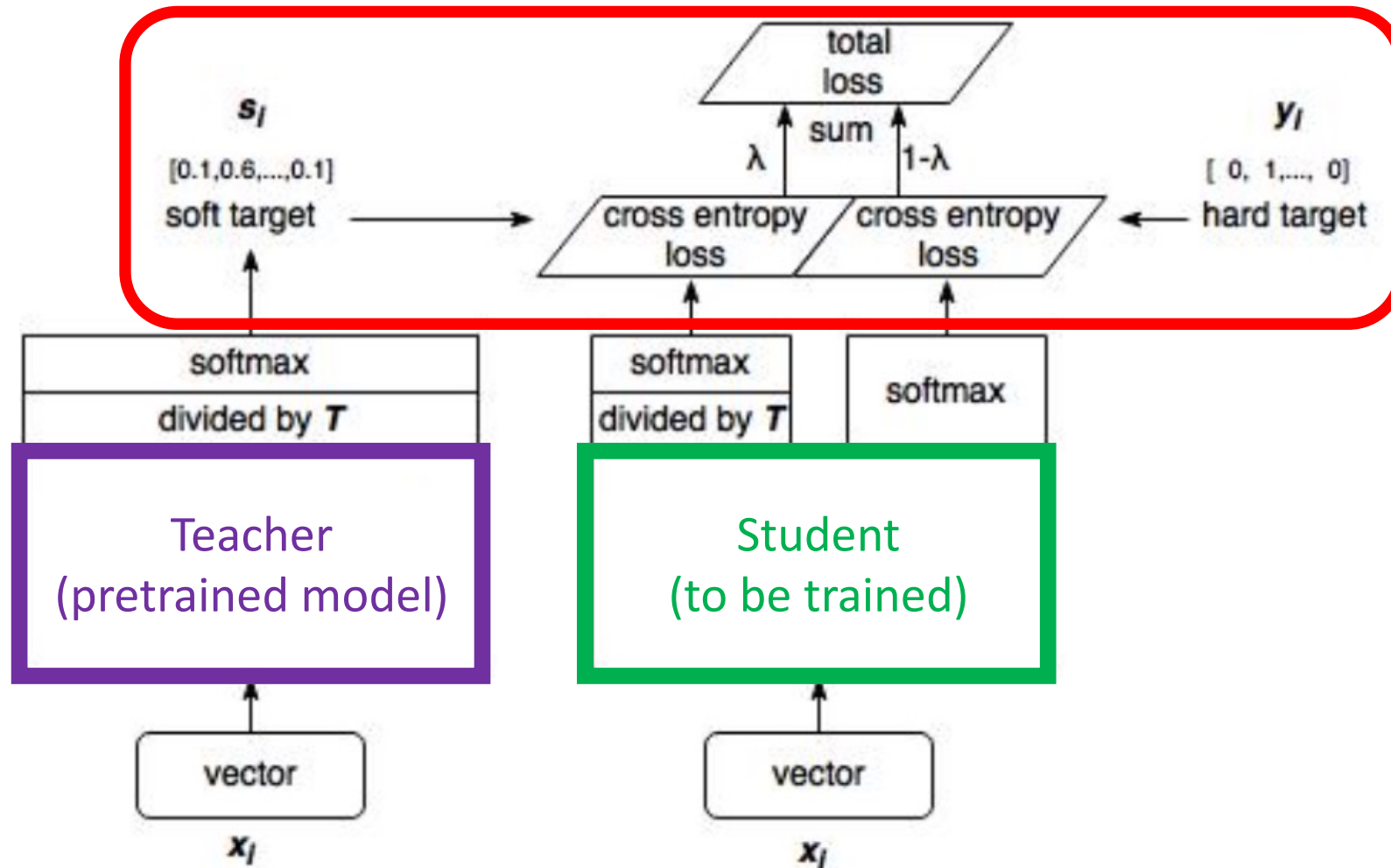
Generalized Softmax: converts vector of **scores** into a probability distribution that sums to 1 with **temperature**; e.g., $T=5$



Knowledge Distillation: Teach Student the “Dark Knowledge” of Teacher

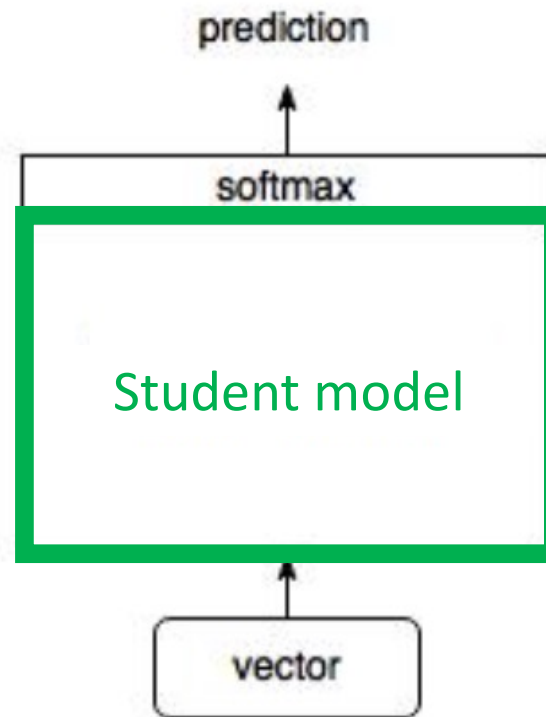


Knowledge Distillation: Teach Student the “Dark Knowledge” of Teacher



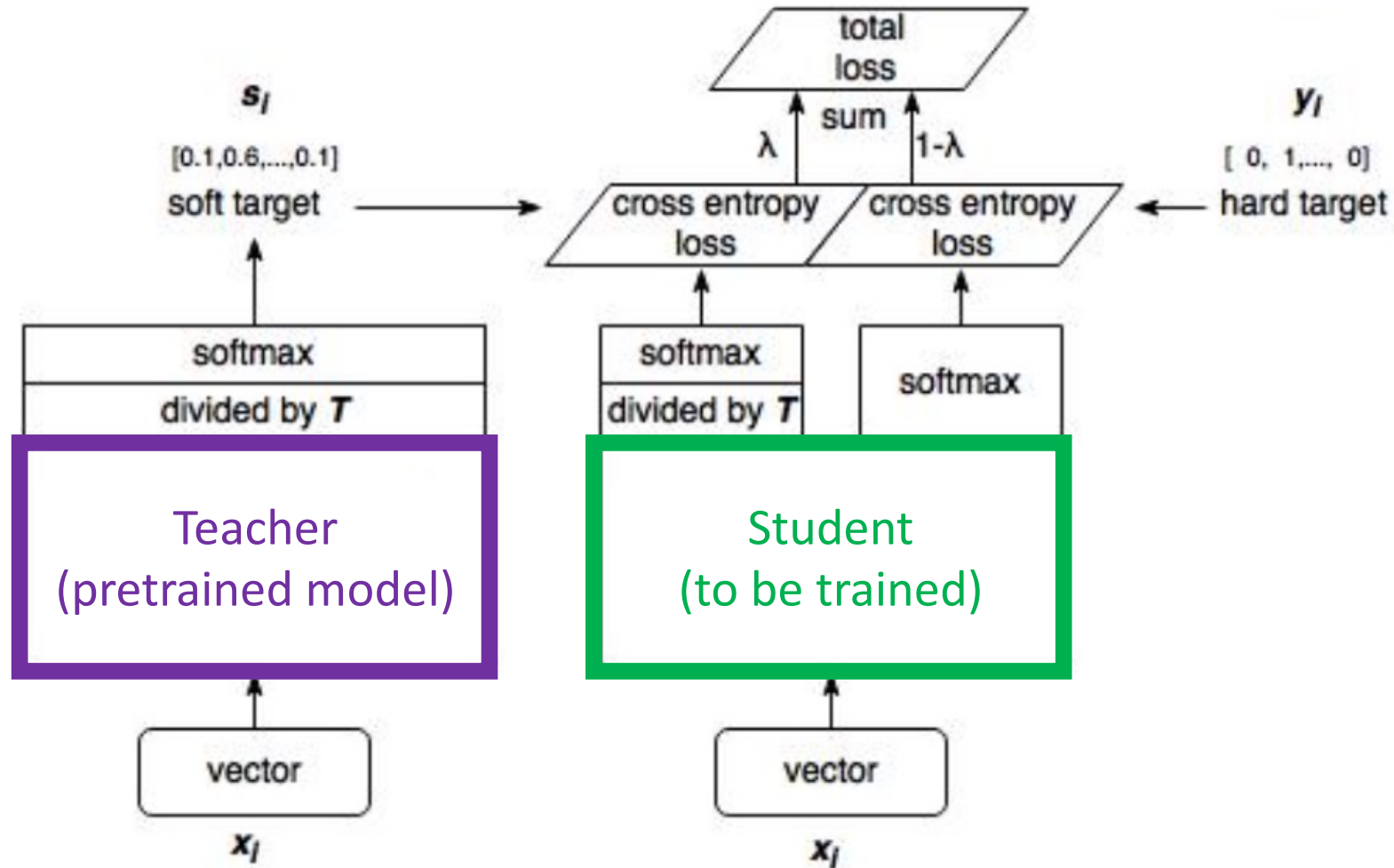
Total loss computed during training is a weighted sum of the conventional cross entropy loss and the “distillation loss”

Knowledge Distillation: At Test Time

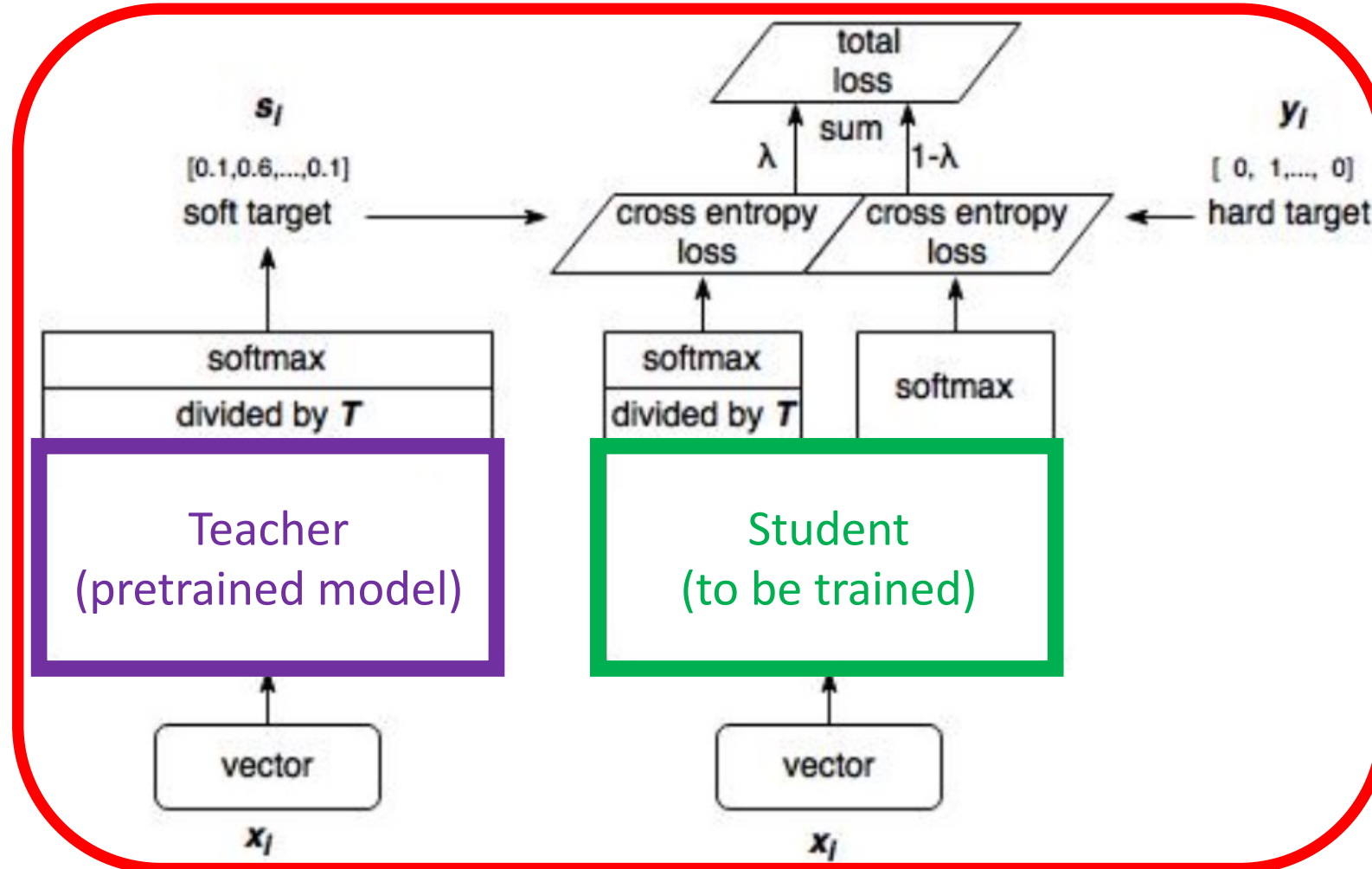


http://blog.csdn.net/qq_22749699

Arguably, Any Neural Network Student Could Learn from Any Neural Network Teacher



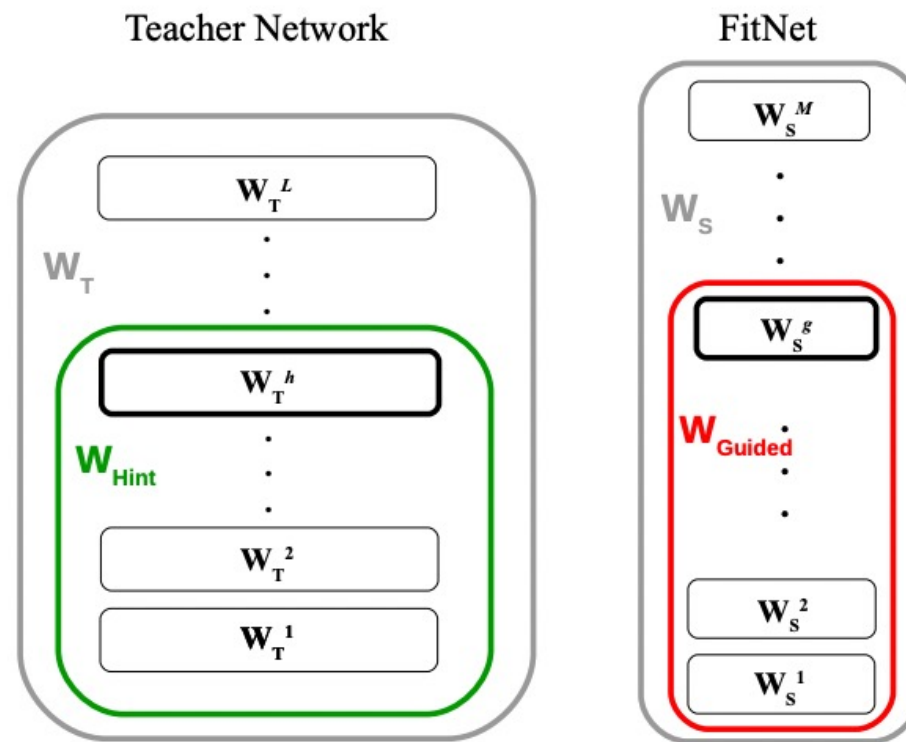
Arguably, Any Neural Network Student Could Learn from Any Neural Network Teacher



Knowledge distillation is a type of transfer learning

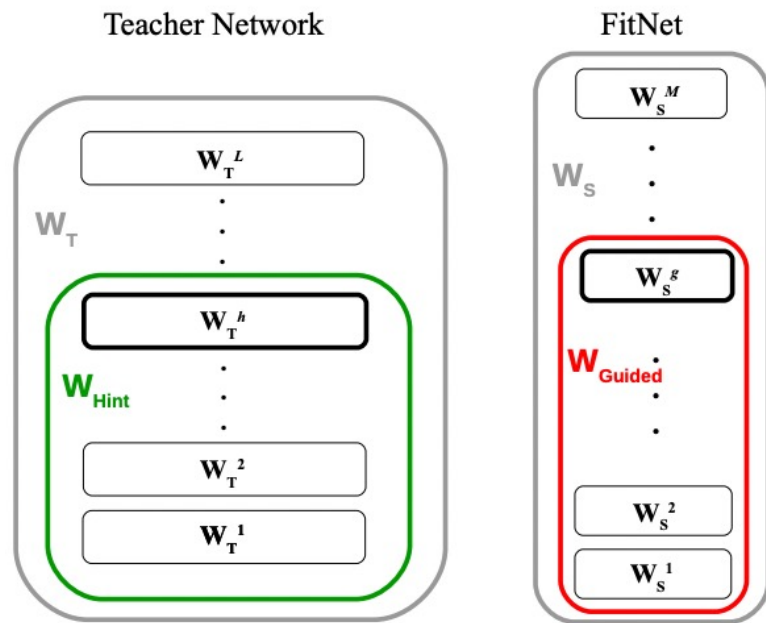
Knowledge Distillation Enhancement: Hints

Encourage student (FitNet) to mimic the teacher's feature responses;
e.g., output of **guided layer** should match the output of **hint layer**



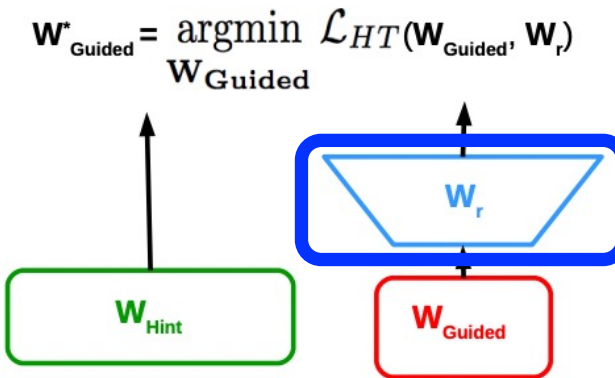
Knowledge Distillation Enhancement: Hints

Encourage student (FitNet) to mimic the teacher's feature responses;
e.g., output of **guided layer** should match the output of **hint layer**



(a) Teacher and Student Networks

Training conducted to learn
the intermediate feature



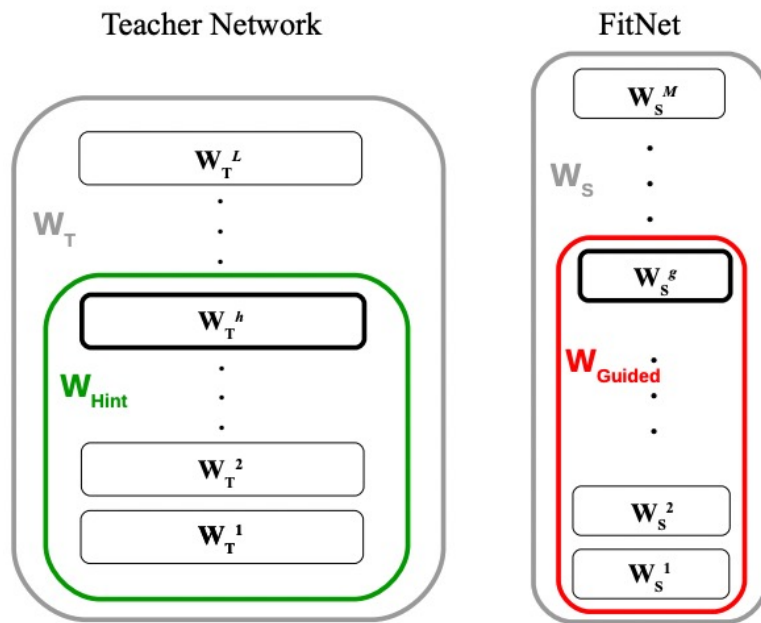
Layer added to match size
of the hint's output layer

(b) Hints Training

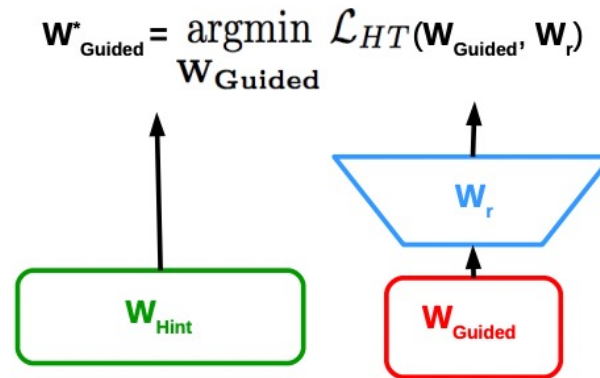
Knowledge Distillation Enhancement: Hints

Encourage student (FitNet) to mimic the teacher's feature responses;
e.g., output of **guided layer** should match the output of **hint layer**

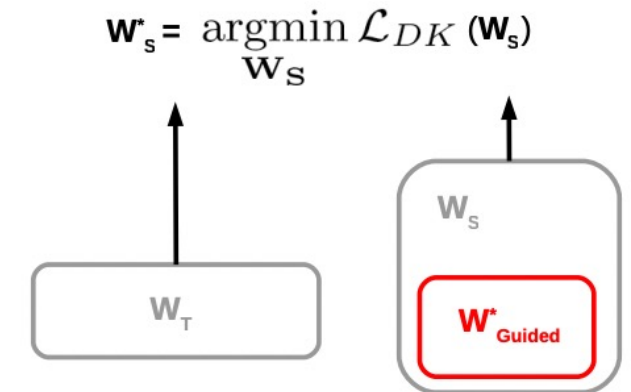
After learning the intermediate features,
the whole student network is trained



(a) Teacher and Student Networks



(b) Hints Training



(c) Knowledge Distillation

Today's Topics

- Motivation
- Key idea: knowledge distillation
- Knowledge distillation for CNNs (vision problems)
- Knowledge distillation for Transformers (language problems)

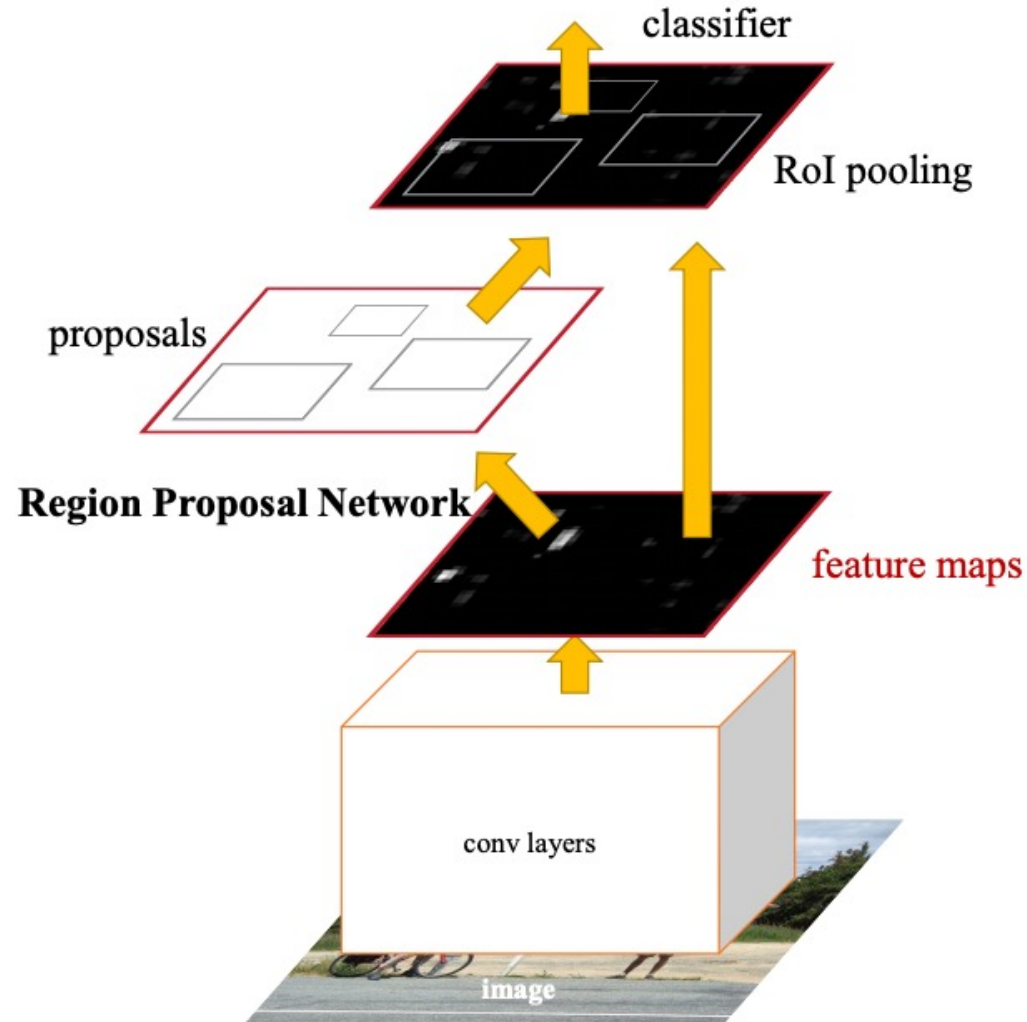
Knowledge Distillation for Vision

- Object detection
- Image classification

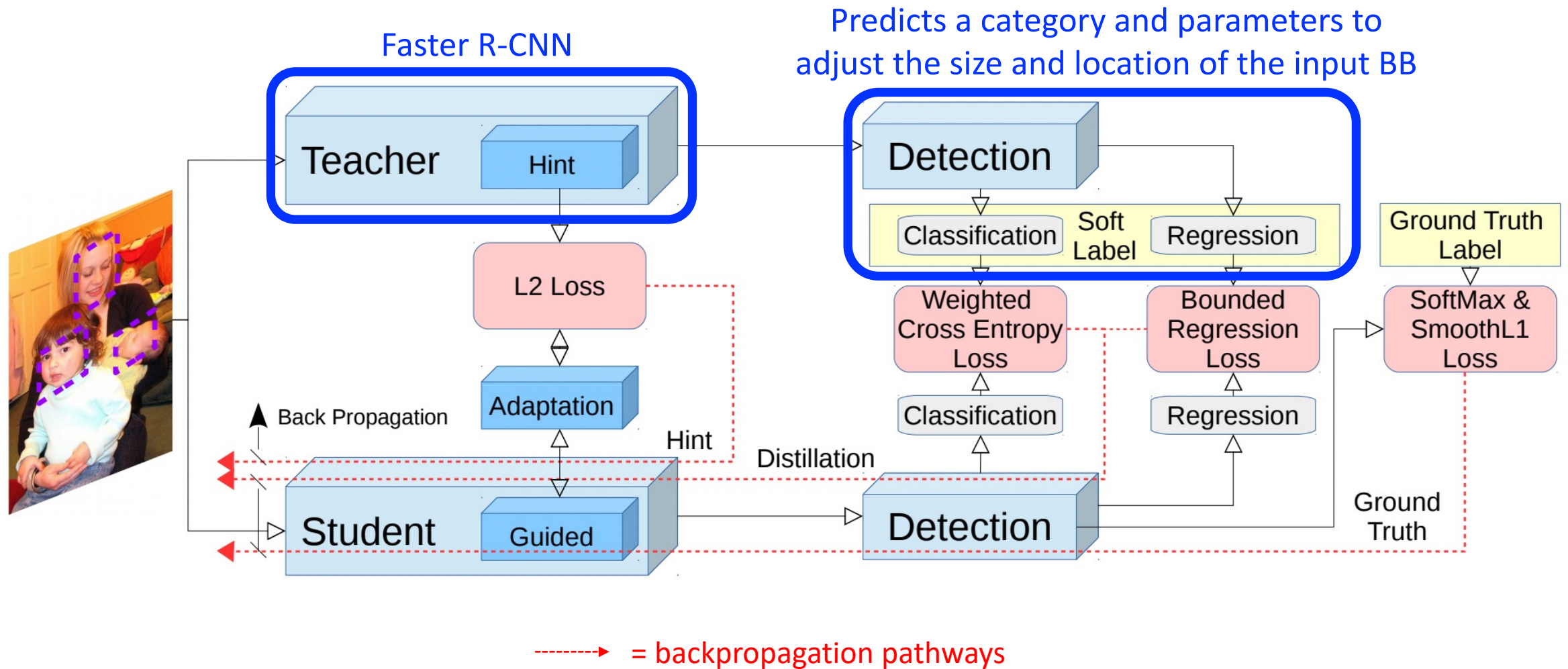
Knowledge Distillation for Vision

- Object detection
- Image classification

Recall Popular Detection Model: Faster R-CNN



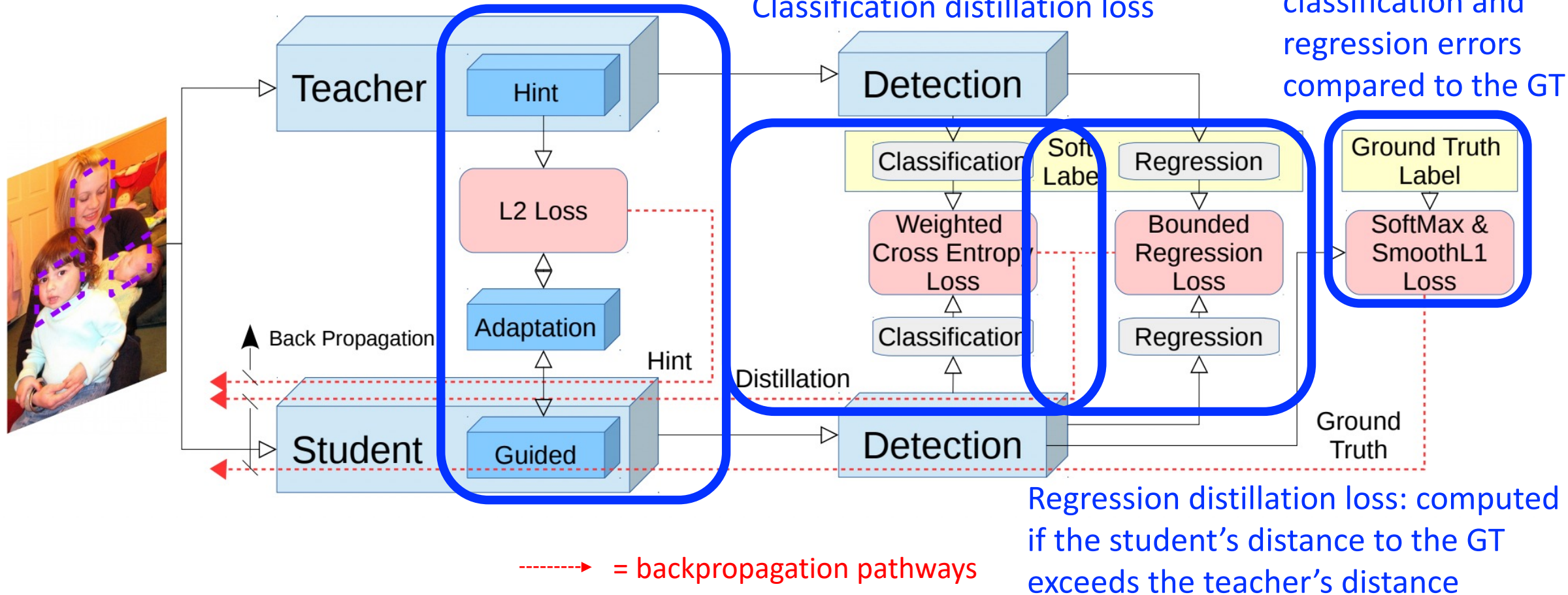
Approach for Creating Compact Student Model



Approach for Creating Compact Student Model

A loss is computed to encourage the student's intermediate features to match those of the teacher

Conventional loss computed for classification and regression errors compared to the GT



Experiments

4 student models

3 teacher models

mAP scores for 5 datasets

Student	Model Info	Teacher	PASCAL	COCO@.5	COCO@[.5,.95]	KITTI	ILSVRC
Tucker	11M / 47ms	-	54.7	25.4	11.8	49.3	20.6
		AlexNet	57.6 (+2.9)	26.5 (+1.2)	12.3 (+0.5)	51.4 (+2.1)	23.6 (+1.3)
		VGGM	58.2 (+3.5)	26.4 (+1.1)	12.2 (+0.4)	51.4 (+2.1)	23.9 (+1.6)
		VGG16	59.4 (+4.7)	28.3 (+2.9)	12.6 (+0.8)	53.7 (+4.4)	24.4 (+2.1)
AlexNet	62M / 74ms	-	57.2	32.5	15.8	55.1	27.3
		VGGM	59.2 (+2.0)	33.4 (+0.9)	16.0 (+0.2)	56.3 (+1.2)	28.7 (+1.4)
		VGG16	60.1 (+2.9)	35.8 (+3.3)	16.9 (+1.1)	58.3 (+3.2)	30.1 (+2.8)
VGGM	80M / 86ms	-	59.8	33.6	16.1	56.7	31.1
		VGG16	63.7 (+3.9)	37.2 (+3.6)	17.3 (+1.2)	58.6 (+2.3)	34.0 (+2.9)
VGG16	138M / 283ms	-	70.4	45.1	24.2	59.2	35.6

params / speed

- means no distillation or, in other words, trained from scratch

What trends do you observe from these results?

Experiments

4 student models

3 teacher models

mAP scores for 5 datasets

Student	Model Info	Teacher	PASCAL	COCO@.5	COCO@[.5,.95]	KITTI	ILSVRC
Tucker	11M / 47ms	-	54.7	25.4	11.8	49.3	20.6
		AlexNet	57.6 (+2.9)	26.5 (+1.2)	12.3 (+0.5)	51.4 (+2.1)	23.6 (+1.3)
		VGGM	58.2 (+3.5)	26.4 (+1.1)	12.2 (+0.4)	51.4 (+2.1)	23.9 (+1.6)
		VGG16	59.4 (+4.7)	28.3 (+2.9)	12.6 (+0.8)	53.7 (+4.4)	24.4 (+2.1)
AlexNet	62M / 74ms	-	57.2	32.5	15.8	55.1	27.3
		VGGM	59.2 (+2.0)	33.4 (+0.9)	16.0 (+0.2)	56.3 (+1.2)	28.7 (+1.4)
		VGG16	60.1 (+2.9)	35.8 (+3.3)	16.9 (+1.1)	58.3 (+3.2)	30.1 (+2.8)
VGGM	80M / 86ms	-	59.8	33.6	16.1	56.7	31.1
		VGG16	63.7 (+3.9)	37.2 (+3.6)	17.3 (+1.2)	58.6 (+2.3)	34.0 (+2.9)
VGG16	138M / 283ms	-	70.4	45.1	24.2	59.2	35.6

- means no distillation or, in other words, trained from scratch

For all student-teacher pairs, knowledge distillation not only yields more compact and faster models but also leads to performance improvements

Chen et al. Learning efficient object detection models with knowledge distillation. Neurips 2017.

Experiments

4 student models

3 teacher models

mAP scores for 5 datasets

Student	Model Info	Teacher	PASCAL	COCO@.5	COCO@[.5,.95]	KITTI	ILSVRC
Tucker	11M / 47ms	-	54.7	25.4	11.8	49.3	20.6
		AlexNet	57.6 (+2.9)	26.5 (+1.2)	12.3 (+0.5)	51.4 (+2.1)	23.6 (+1.3)
		VGGM	58.2 (+3.5)	26.4 (+1.1)	12.2 (+0.4)	51.4 (+2.1)	23.9 (+1.6)
		VGG16	59.4 (+4.7)	28.3 (+2.9)	12.6 (+0.8)	53.7 (+4.4)	24.4 (+2.1)
AlexNet	62M / 74ms	-	57.2	32.5	15.8	55.1	27.3
		VGGM	59.2 (+2.0)	33.4 (+0.9)	16.0 (+0.2)	56.3 (+1.2)	28.7 (+1.4)
		VGG16	60.1 (+2.9)	35.8 (+3.3)	16.9 (+1.1)	58.3 (+3.2)	30.1 (+2.8)
VGGM	80M / 86ms	-	59.8	33.6	16.1	56.7	31.1
		VGG16	63.7 (+3.9)	37.2 (+3.6)	17.3 (+1.2)	58.6 (+2.3)	34.0 (+2.9)
VGG16	138M / 283ms	-	70.4	45.1	24.2	59.2	35.6

- means no distillation or, in other words, trained from scratch

Larger teachers leads to greater performance improvements for distilled models

Experiments

mAP scores for 5 datasets

Student	Model Info	Teacher	PASCAL	COCO@.5	COCO@[.5,.95]	KITTI	ILSVRC
Tucker	11M / 47ms	-	54.7	25.4	11.8	49.3	20.6
		AlexNet	57.6 (+2.9)	26.5 (+1.2)	12.3 (+0.5)	51.4 (+2.1)	23.6 (+1.3)
		VGGM	58.2 (+3.5)	26.4 (+1.1)	12.2 (+0.4)	51.4 (+2.1)	23.9 (+1.6)
		VGG16	59.4 (+4.7)	28.3 (+2.9)	12.6 (+0.8)	53.7 (+4.4)	24.4 (+2.1)
AlexNet	62M / 74ms	-	57.2	32.5	15.8	55.1	27.3
		VGGM	59.2 (+2.0)	33.4 (+0.9)	16.0 (+0.2)	56.3 (+1.2)	28.7 (+1.4)
		VGG16	60.1 (+2.9)	35.8 (+3.3)	16.9 (+1.1)	58.3 (+3.2)	30.1 (+2.8)
VGGM	80M / 86ms	-	59.8	33.6	16.1	56.7	31.1
		VGG16	63.7 (+3.9)	37.2 (+3.6)	17.3 (+1.2)	58.6 (+2.3)	34.0 (+2.9)
VGG16	138M / 283ms	-	70.4	45.1	24.2	59.2	35.6

- means no distillation or, in other words, trained from scratch

Still, larger models with more parameters return the best results.

Experiments

4 student models

3 teacher models

mAP scores for 5 datasets

Student	Model Info	Teacher	PASCAL	COCO@.5	COCO@[.5,.95]	KITTI	ILSVRC
Tucker	11M / 47ms	-	54.7	25.4	11.8	49.3	20.6
		AlexNet	57.6 (+2.9)	26.5 (+1.2)	12.3 (+0.5)	51.4 (+2.1)	23.6 (+1.3)
		VGGM	58.2 (+3.5)	26.4 (+1.1)	12.2 (+0.4)	51.4 (+2.1)	23.9 (+1.6)
		VGG16	59.4 (+4.7)	28.3 (+2.9)	12.6 (+0.8)	53.7 (+4.4)	24.4 (+2.1)
AlexNet	62M / 74ms	-	57.2	32.5	15.8	55.1	27.3
		VGGM	59.2 (+2.0)	33.4 (+0.9)	16.0 (+0.2)	56.3 (+1.2)	28.7 (+1.4)
		VGG16	60.1 (+2.9)	35.8 (+3.3)	16.9 (+1.1)	58.3 (+3.2)	30.1 (+2.8)
VGGM	80M / 86ms	-	59.8	33.6	16.1	56.7	31.1
		VGG16	63.7 (+3.9)	37.2 (+3.6)	17.3 (+1.2)	58.6 (+2.3)	34.0 (+2.9)
VGG16	138M / 283ms	-	70.4	45.1	24.2	59.2	35.6

- means no distillation or, in other words, trained from scratch

Why do you think there are performance improvements from model compression?

Knowledge Distillation for Vision

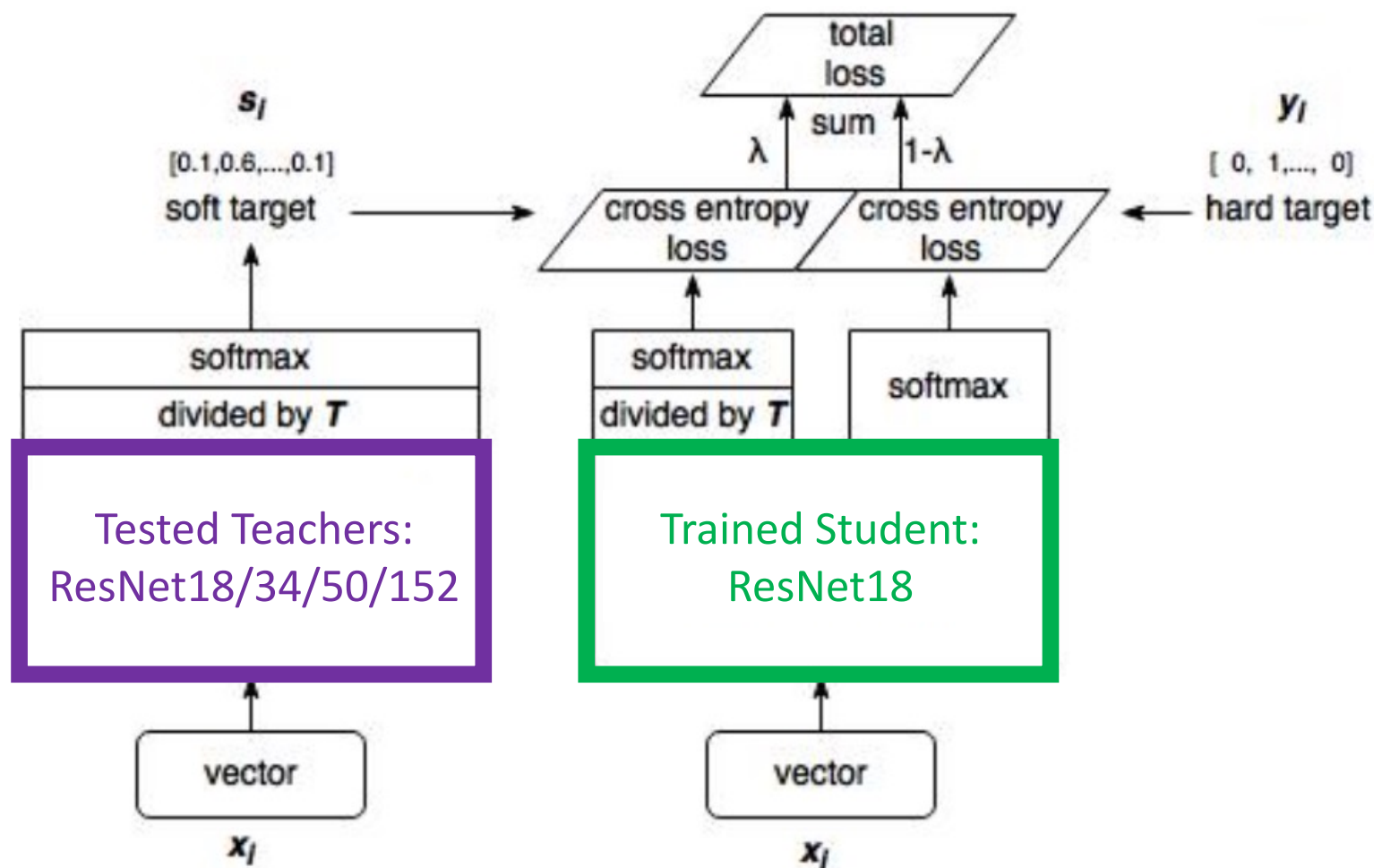
- Object detection
- Image classification with attention transfer

Recall Task: Predict Category from 1000 Options

- **Evaluation metric:** % correct (top-1 and top-5 predictions)
- **Dataset:** ~1.5 million images
- **Source:** images scraped from search engines, such as Flickr, and labeled by crowdworkers



Experiment: Do Bigger, More Accurate Models Make Better Teachers?



Experiment: Do Bigger, More Accurate Models Make Better Teachers?

(% = Top-1 error rates)

Teacher	Teacher Error (%)	Student Error (%)
ResNet18	30.24	30.57
ResNet34	26.70	30.79
ResNet50	23.85	30.95

What is the student's performance trend from larger, more accurate teachers?

Experiment: Do Bigger, More Accurate Models Make Better Teachers?

(% = Top-1 error rates)

Teacher	Teacher Error (%)	Student Error (%)
-	-	30.24
ResNet18	30.24	30.57
ResNet34	26.70	30.79
ResNet50	23.85	30.95

Student performance not only drops for larger teachers but the **models distilled from teachers perform worse than training the student from scratch!**

Experiment: Why Might Student Performance Drop as Teacher Size Grows?

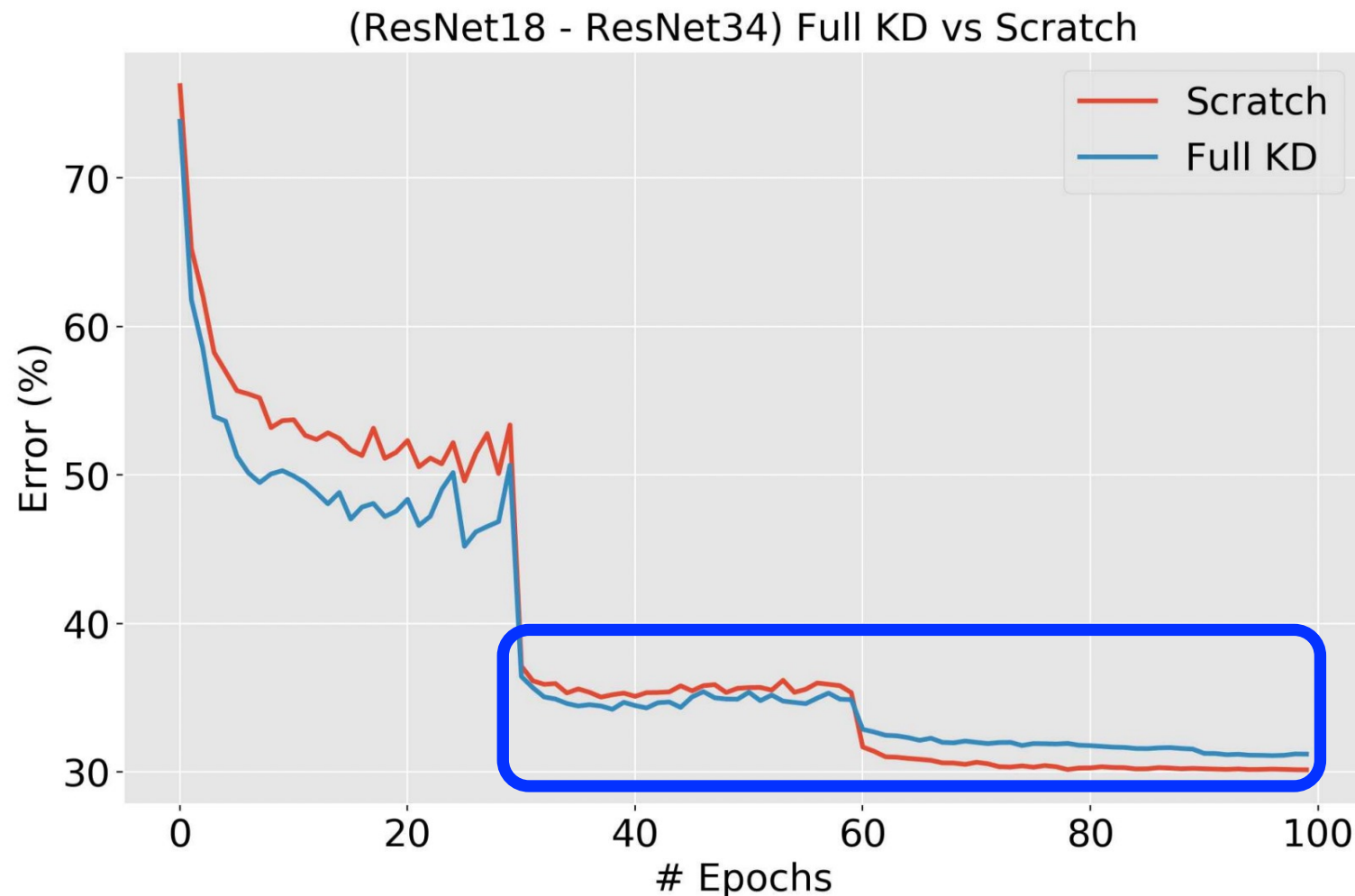
1. More accurate models are more confident and so need higher temperatures to learn the “dark knowledge” of category relationships
2. Student mimics teacher but the loss function is mismatched from the evaluation metric

3. Student fails to accurately mimic teacher

Experimental analysis
suggests this is the reason

Experiment: Why Might Student Students Fail to Mimic Teachers?

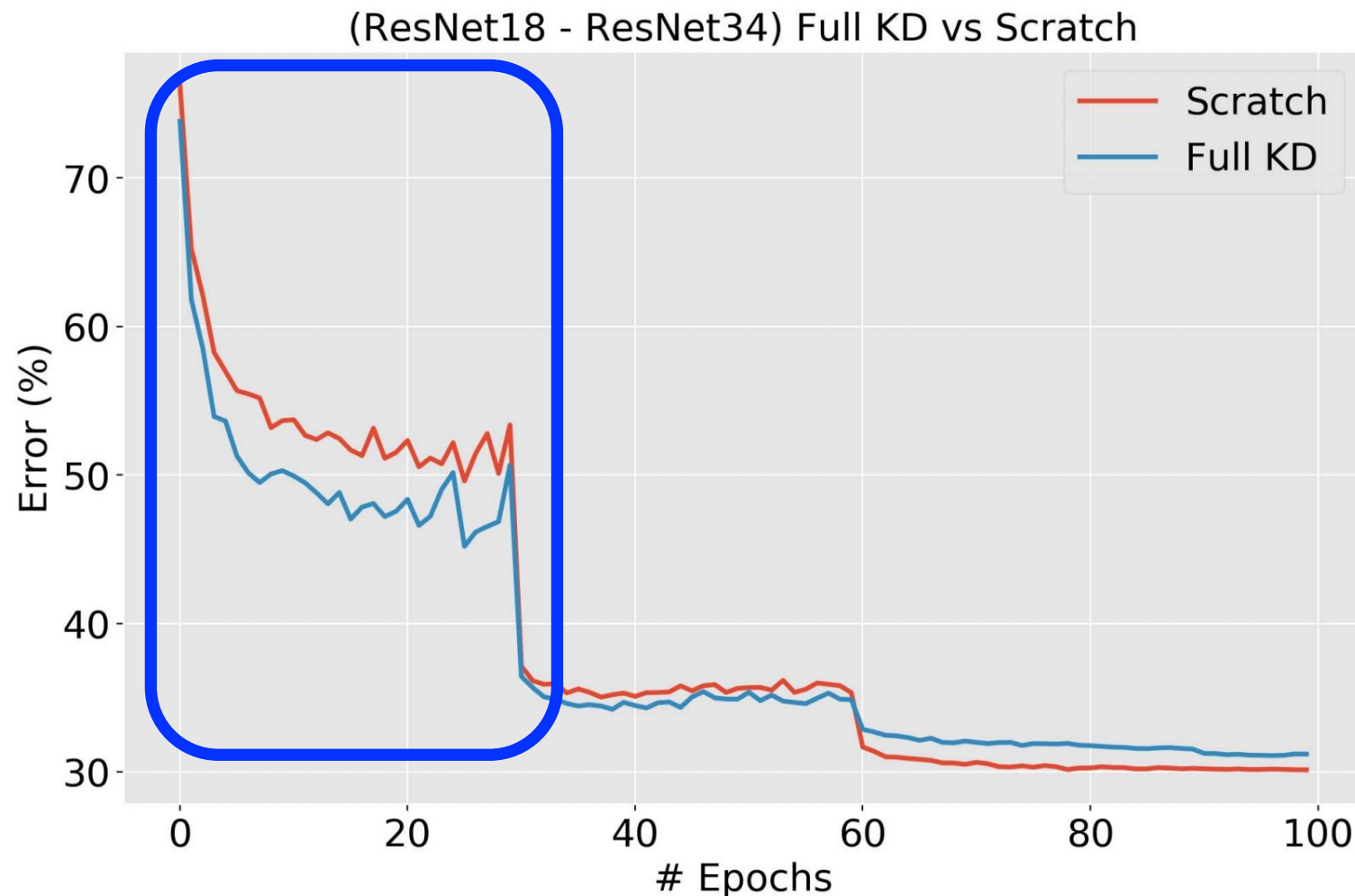
Hypothesis: student is underfitting because of lower capacity and so “minimizing one loss (KD loss) at the expense of the other (cross entropy loss)”



Experiment: Why Might Student Students Fail to Mimic Teachers?

How to overcome this issue?

- Early stopping with KD loss (ESKD) to leverage its benefit at the start of training



Experiments: How Does ESKD Compare To Training A Student from Scratch?

Teacher	Top-1 Error (%, Test)
ResNet18	30.57
ResNet18 (ES KD)	29.01
ResNet34	30.79
ResNet34 (ES KD)	29.16
ResNet50	30.95
ResNet50 (ES KD)	29.35

Training a model with early stopping knowledge distillation loss leads to better results than training from scratch!

Experiments: Does ESKD with Bigger, More Accurate Models Make Better Teachers?

Teacher	Top-1 Error (%, Test)
ResNet18	30.57
ResNet18 (ES KD)	29.01
ResNet34	30.79
ResNet34 (ES KD)	29.16
ResNet50	30.95
ResNet50 (ES KD)	29.35

No; the student may still be struggling with underfitting due to an insufficient representational capacity

Experiments: To Address The Capacity Problem Why Not Instead Distill to Intermediate Sizes?

Performs almost identically to a model that is distilled directly from a large to small size; does not address the core problem:

The student must be in the solution space of the teacher

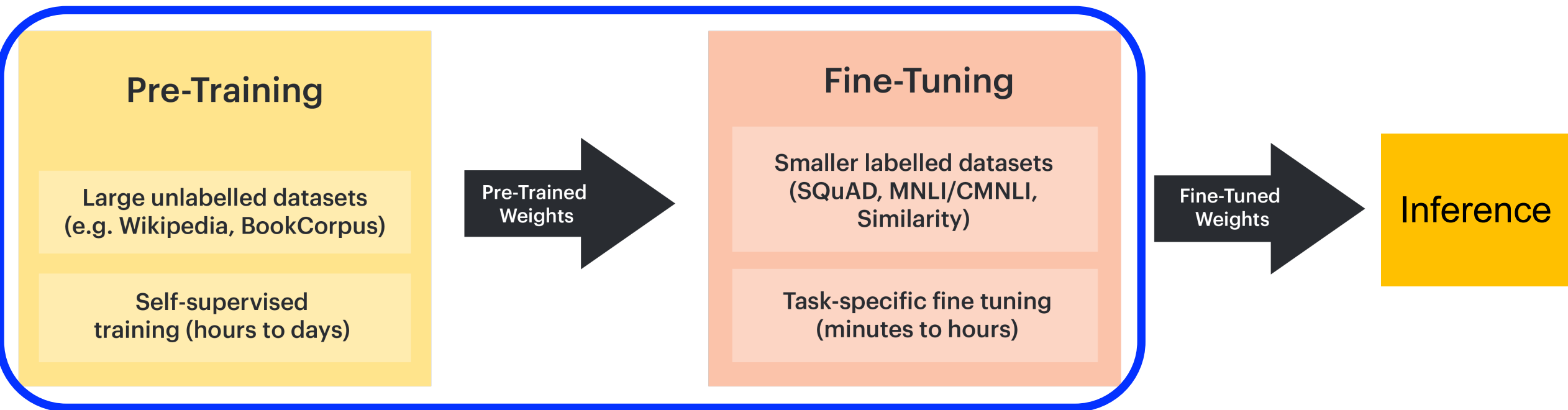
Knowledge Distillation for Vision

- Object detection
- Image classification with attention transfer

Today's Topics

- Motivation
- Key idea: knowledge distillation
- Knowledge distillation for CNNs (vision problems)
- Knowledge distillation for Transformers (language problems)

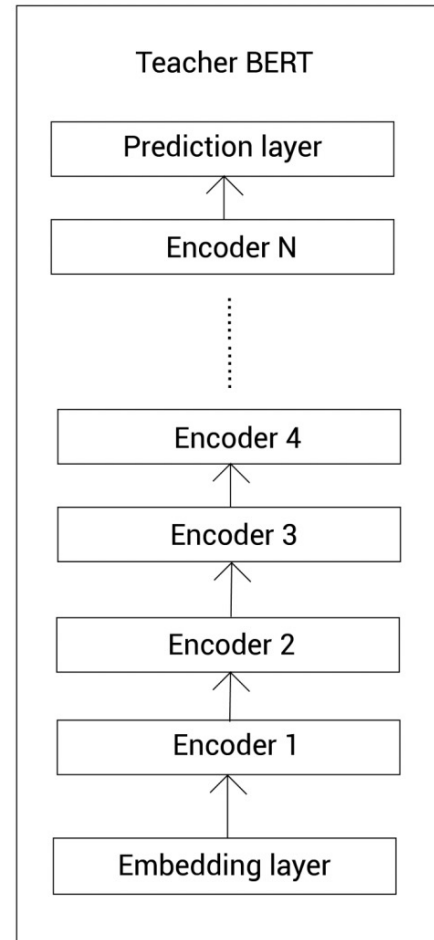
Recall Popular Transformer Approach (e.g., BERT)



Knowledge distillation needed at both
the pretraining and fine-tuning stages

TinyBERT: Pre-training Knowledge Distillation

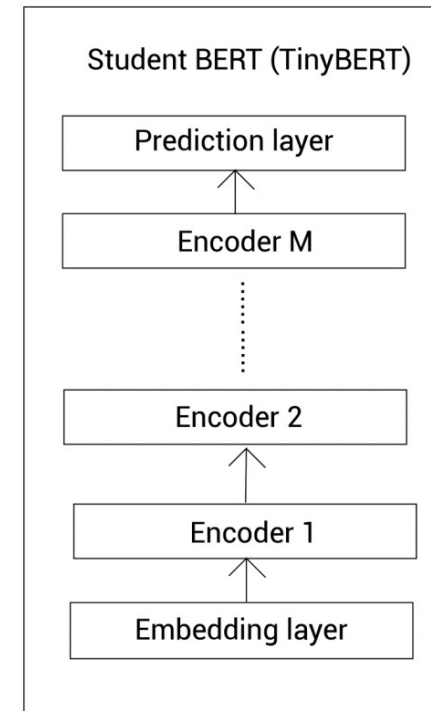
BERT-Base with 12 encoder layers and 12 attention heads that produces a 768 representation:



Distillation



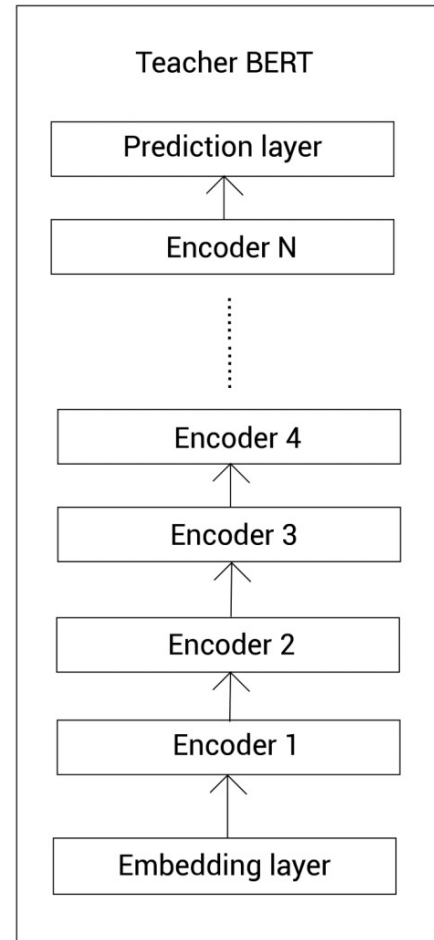
Same architecture but with 4 encoder layers and it produces a 312 representation



input

TinyBERT: Pre-training Knowledge Distillation

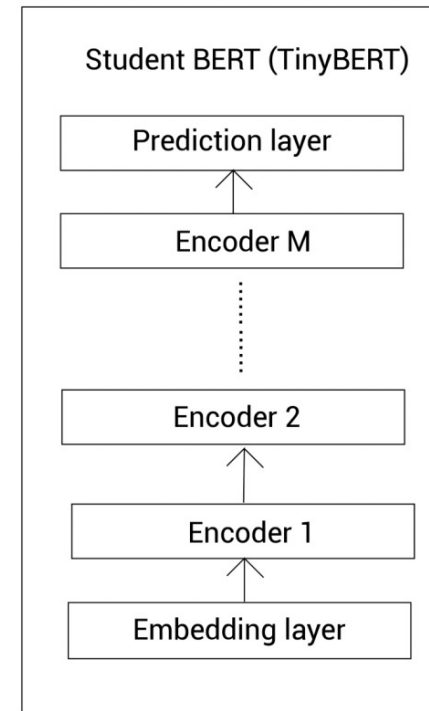
109 million
parameters



Distillation



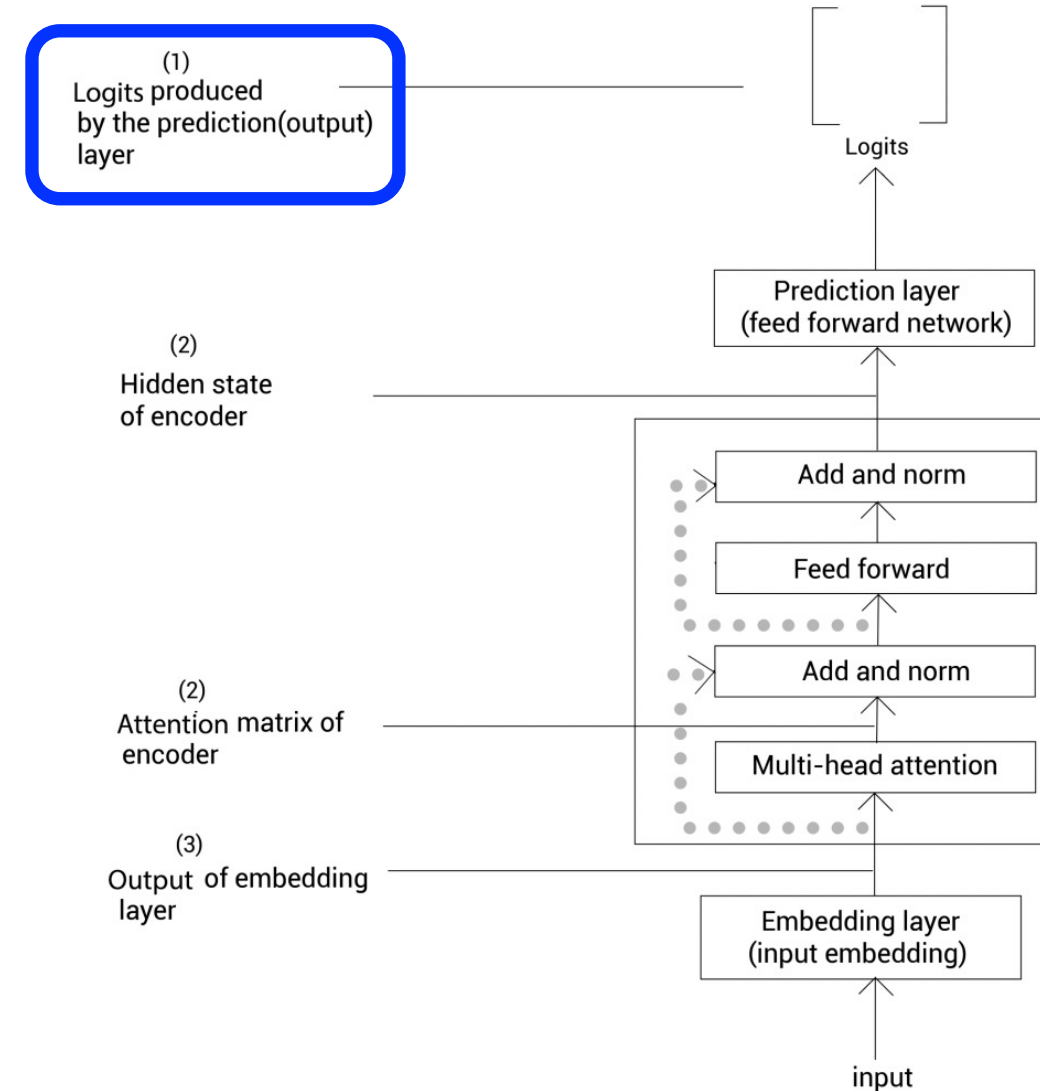
14.5 million
parameters



input

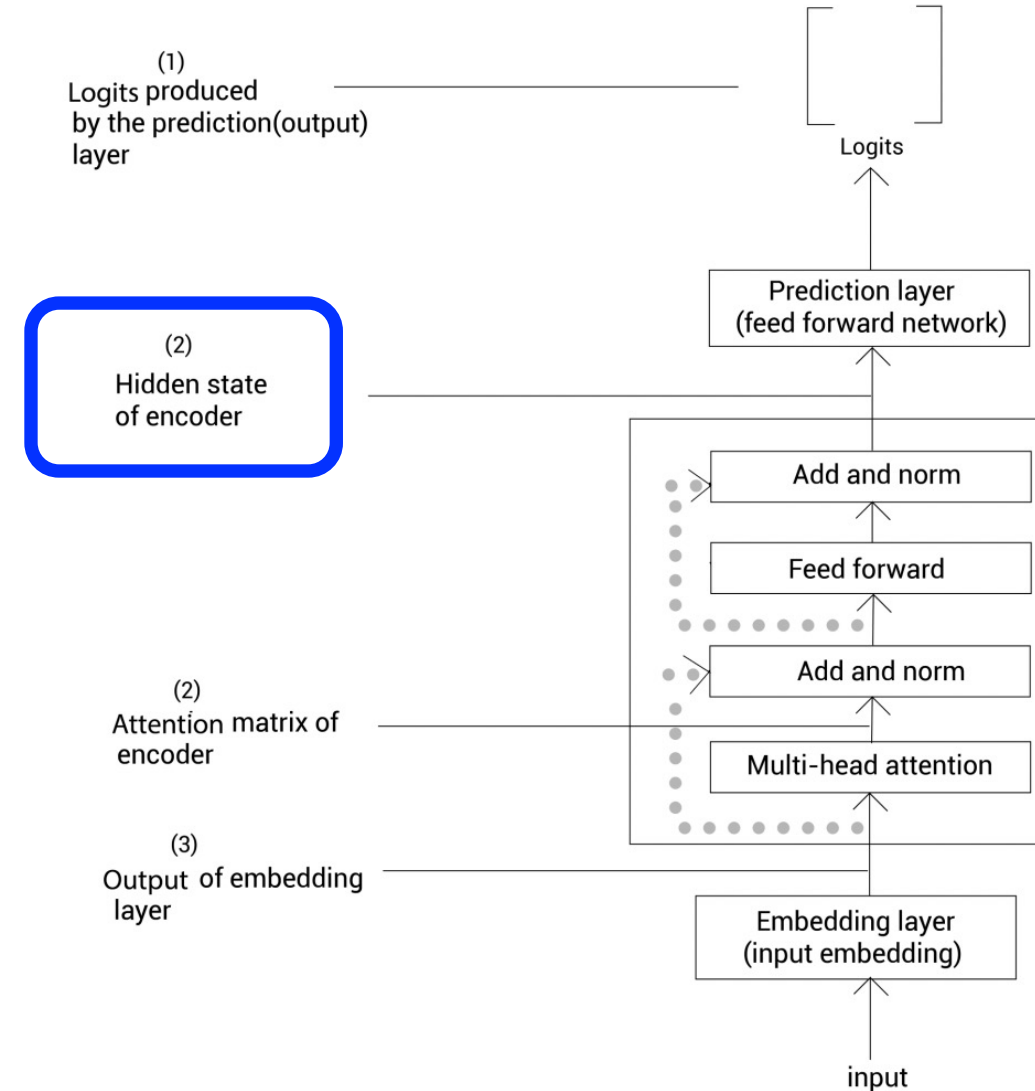
TinyBERT: Pre-training Knowledge Distillation

Student learns from teacher's **output** (i.e., logits) and intermediate layers:



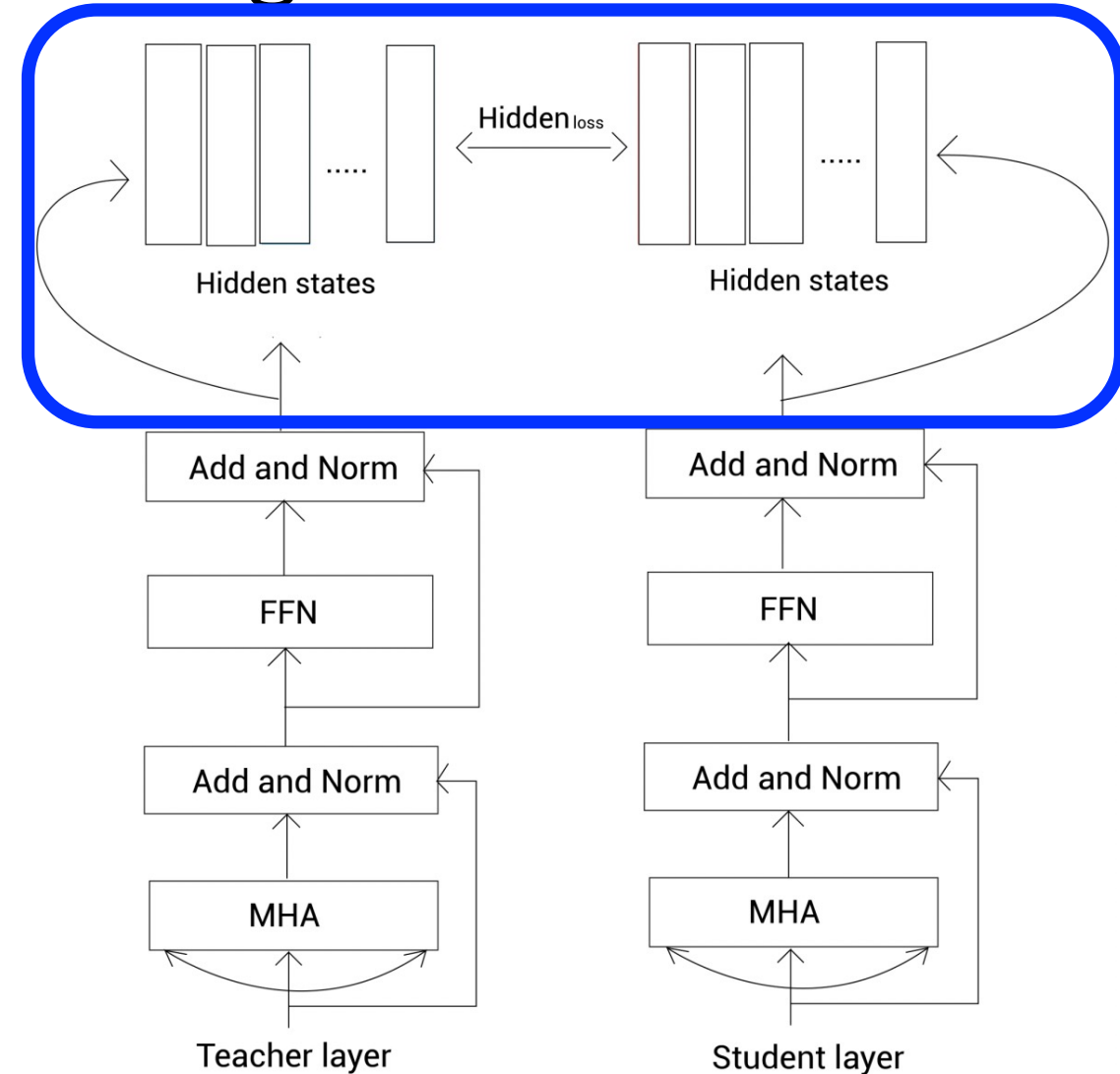
TinyBERT: Pre-training Knowledge Distillation

Student learns from teacher's output (i.e., logits) and **intermediate layers**:



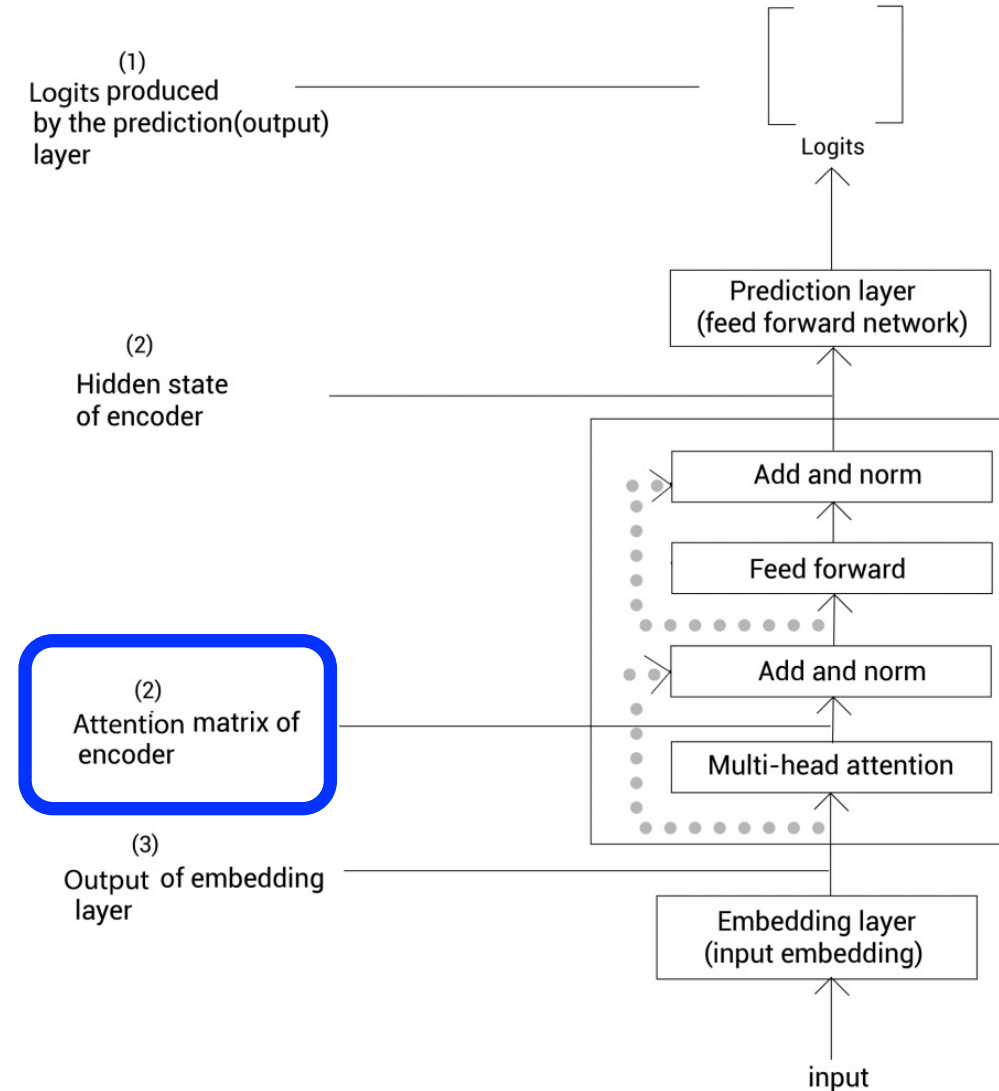
TinyBERT: Pre-training Knowledge Distillation

Mean squared error computed between encoded representations with weight matrix converting student's 312 representation to match teacher's 768 representation



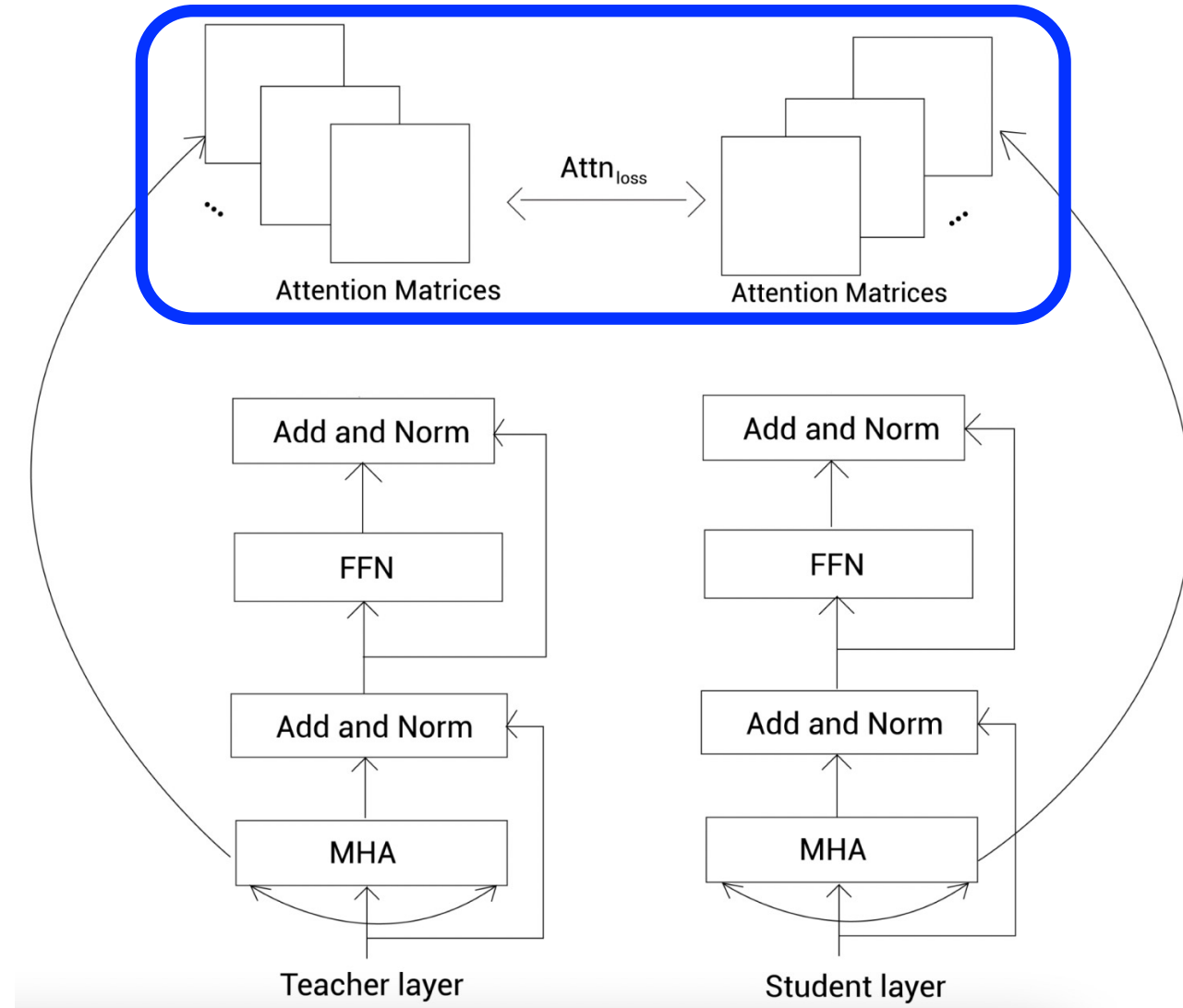
TinyBERT: Pre-training Knowledge Distillation

Student learns from teacher's output (i.e., logits) and **intermediate layers**:



TinyBERT: Pre-training Knowledge Distillation

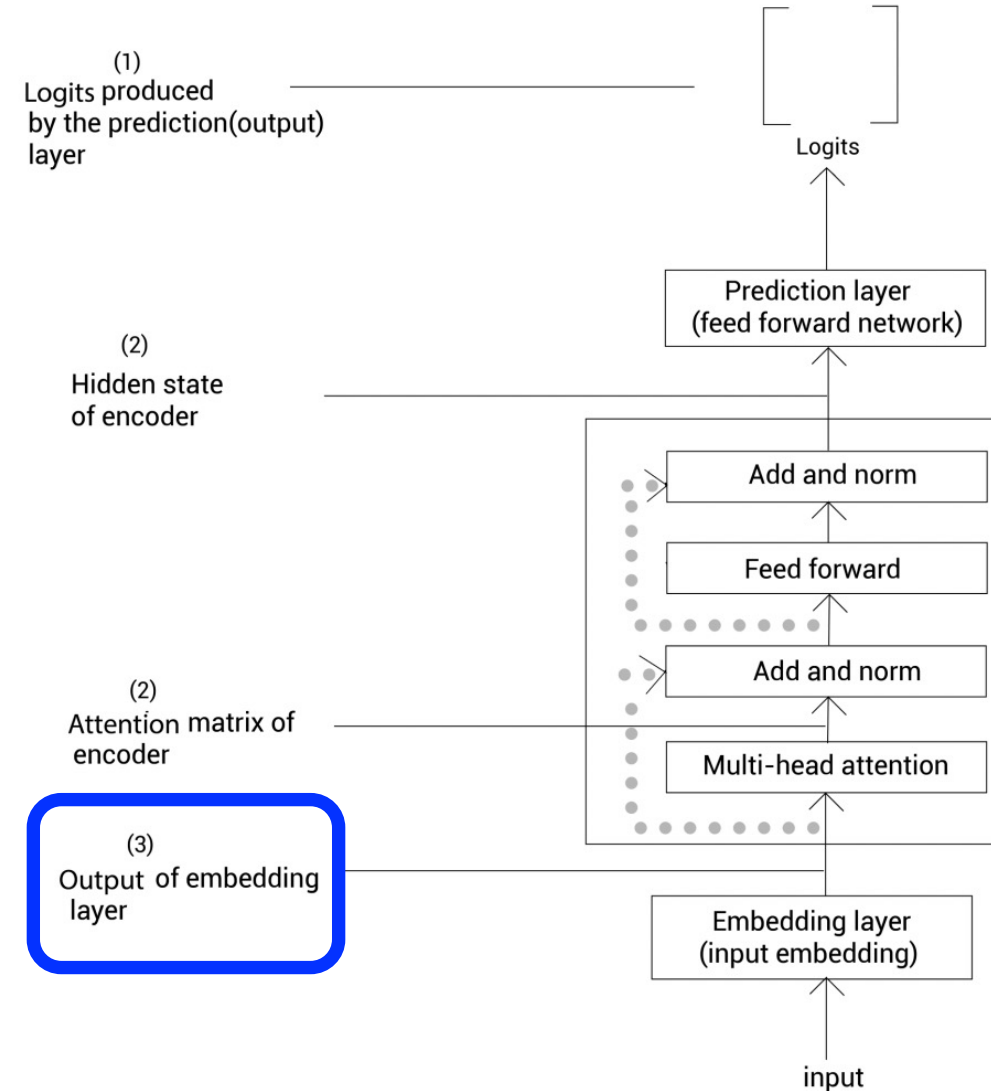
Mean squared error computed between all attention matrices of the teacher and student



TinyBERT: Pre-training Knowledge Distillation

Student learns from teacher's output (i.e., logits) and **intermediate layers**:

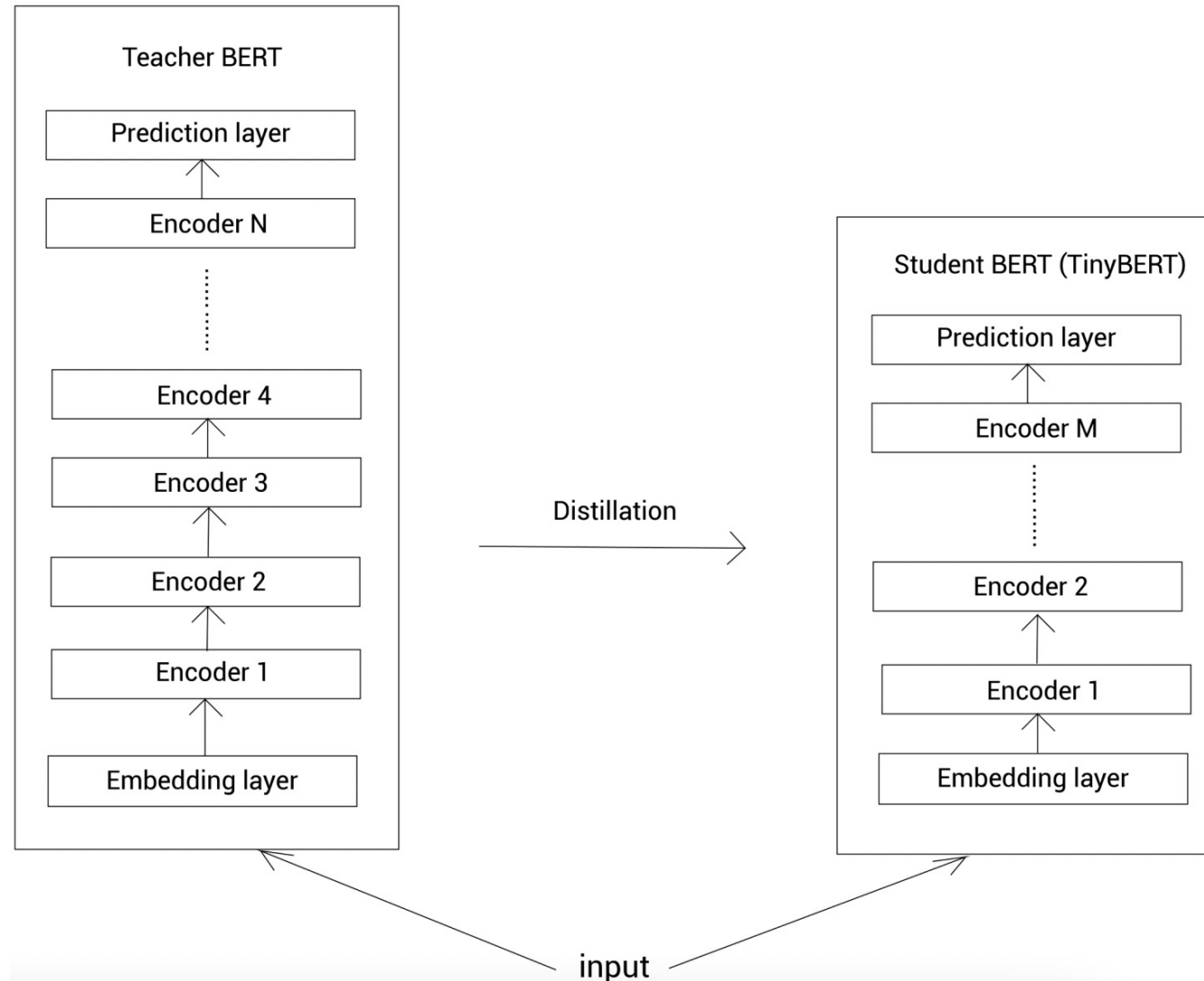
Mean squared error computed between embedding representations with weight matrix converting student's representation to match teacher's representation



TinyBERT: Pre-training Knowledge Distillation

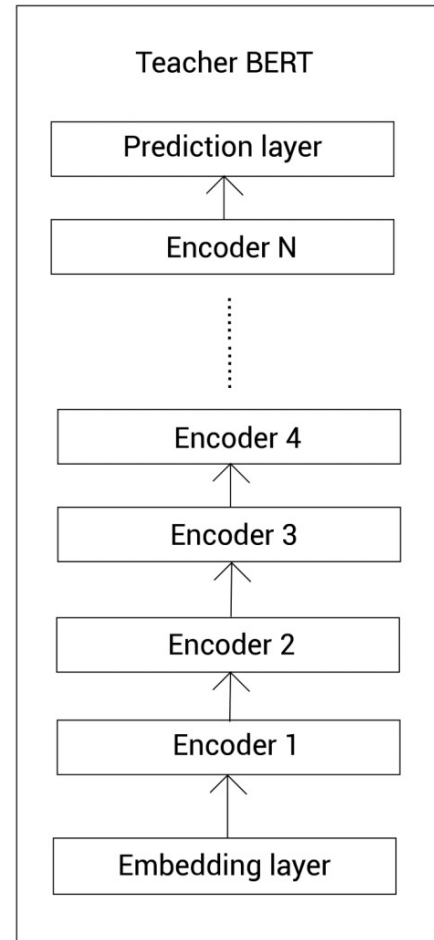
Training data matches that used for BERT: Wikipedia and Toronto BookCorpus

Training run for three epochs with errors of intermediate layers



TinyBERT: Fine-tuning Knowledge Distillation

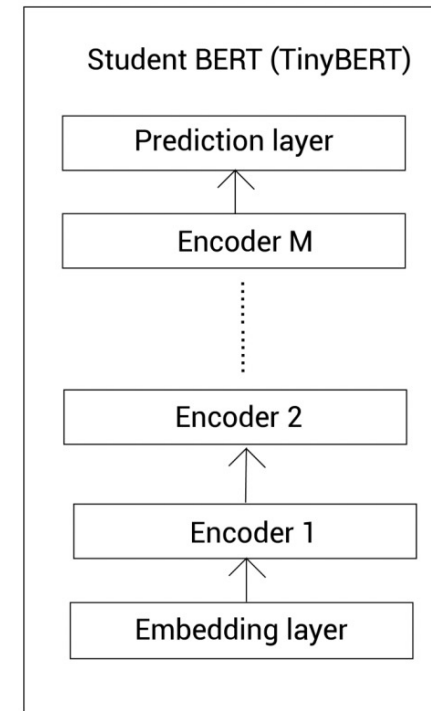
Fine-tuned BERT
for the target task



Distillation



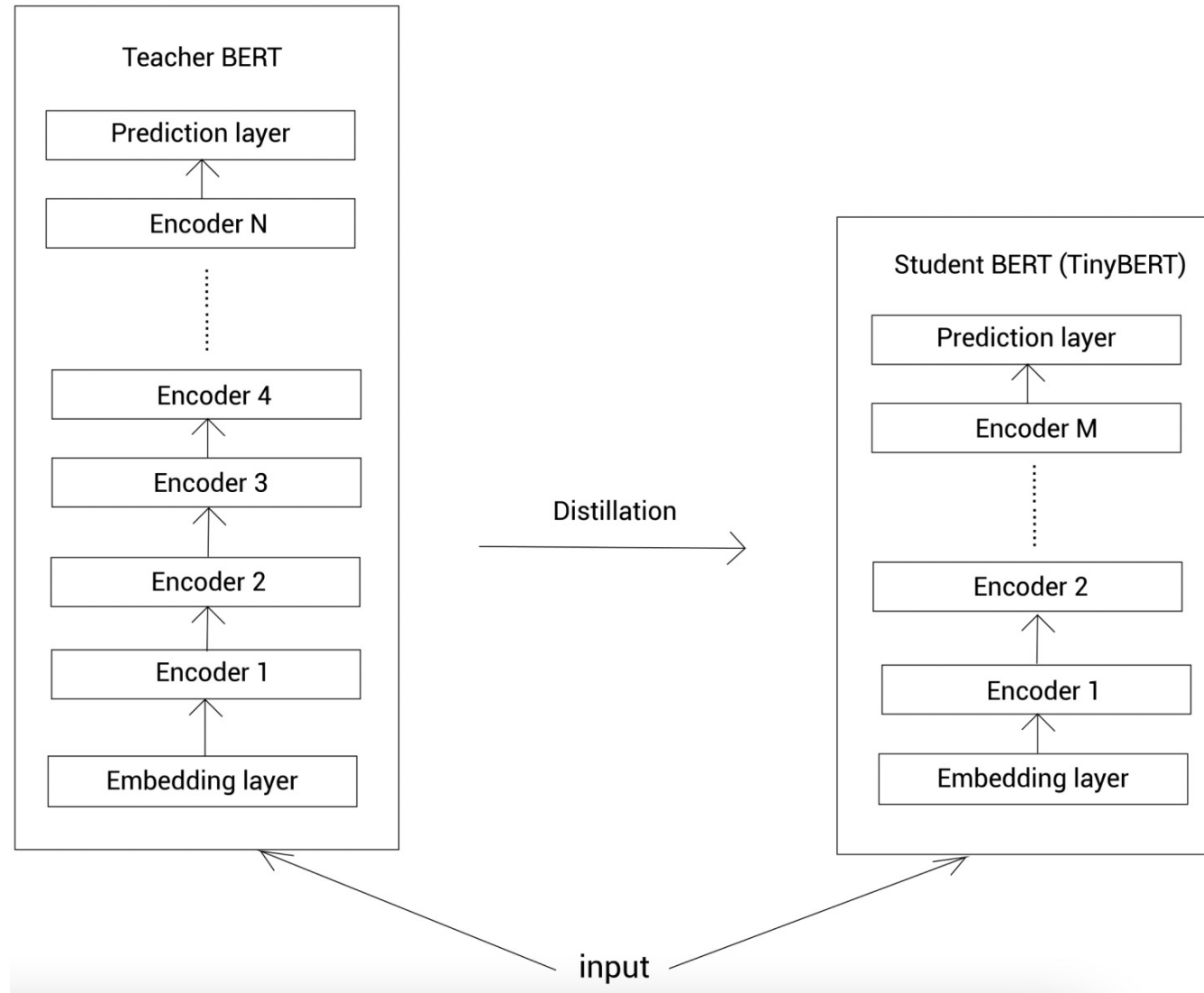
TinyBERT distilled
from pretrained BERT



input

TinyBERT: Fine-tuning Knowledge Distillation

Training run for 20 epochs with errors of intermediate layers and then for 30 epochs using errors from the prediction layer



TinyBERT: Fine-tuning Knowledge Distillation

Compared to BERT-base, TinyBERT achieves nearly comparable performance while containing ~13.3% of the parameters and needing only ~10.6% of inference time

Today's Topics

- Motivation
- Key idea: knowledge distillation
- Knowledge distillation for CNNs (vision problems)
- Knowledge distillation for Transformers (language problems)

Today's Topics

- Motivation
- Key idea: knowledge distillation
- Knowledge distillation for CNNs (vision problems)
- Knowledge distillation for Transformers (language problems)



The End