

# **Knowledge Distillation**

Jinyang Guo ( 郭晋阳 )  
[jinyangguo@buaa.edu.cn](mailto:jinyangguo@buaa.edu.cn)

# Last lecture

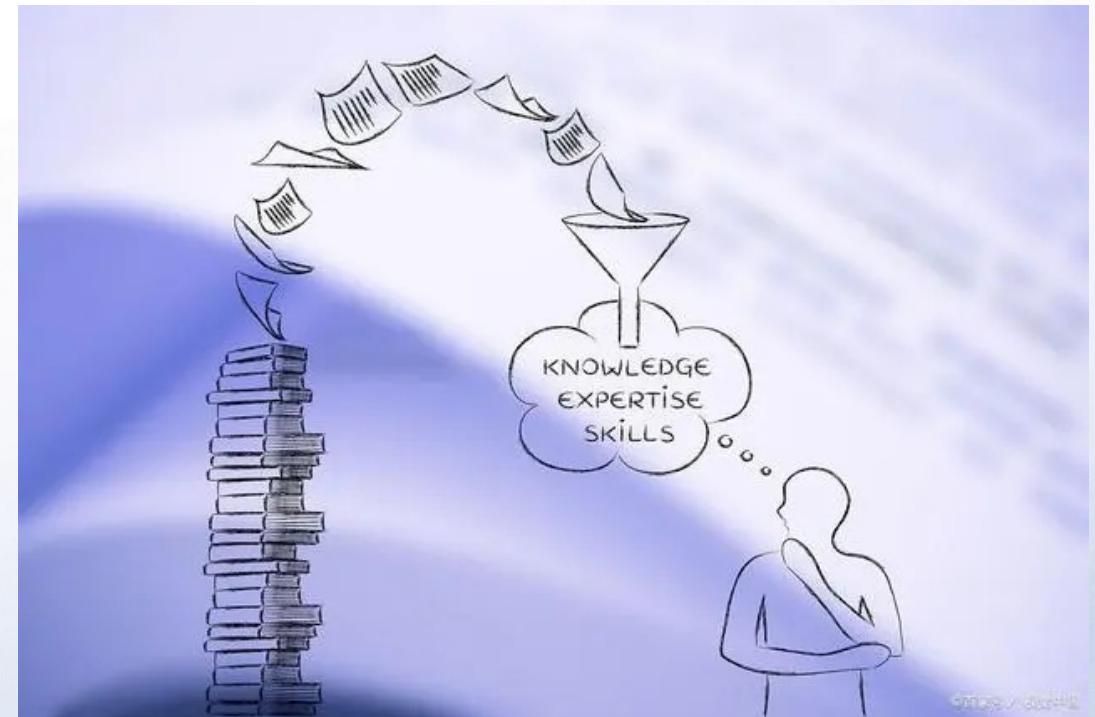
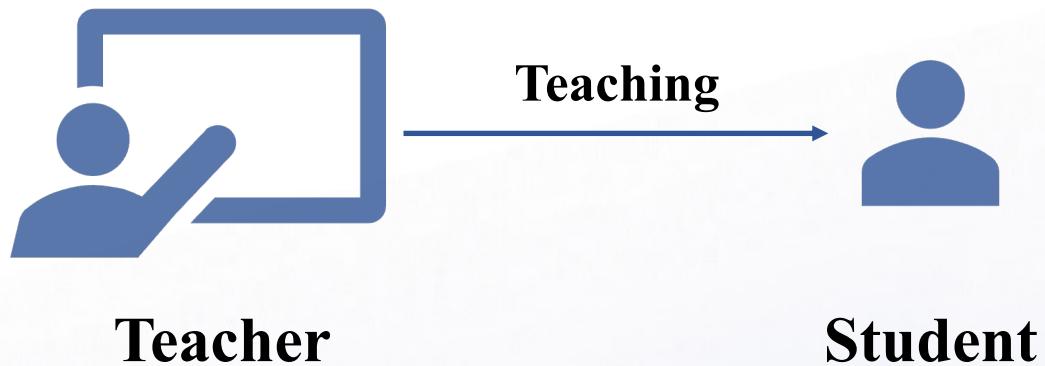
- Network pruning
- Make network sparse

# 目 录

-  **Introduction**
-  **Knowledge distillation**
-  **Knowledge definition**
-  **Knowledge transfer**
-  **Application**

# Introduction

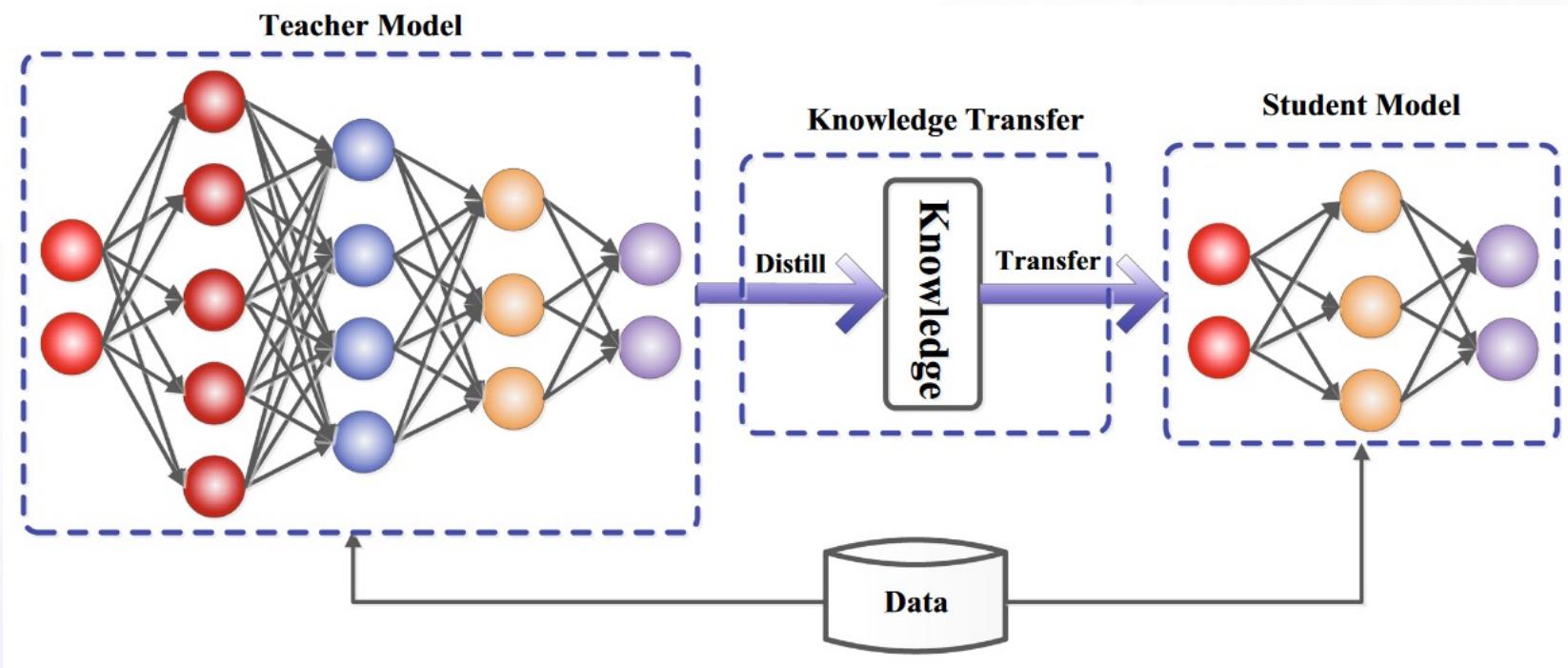
- Human knowledge can be transferred from **teachers** to **students** through teaching.



# Introduction

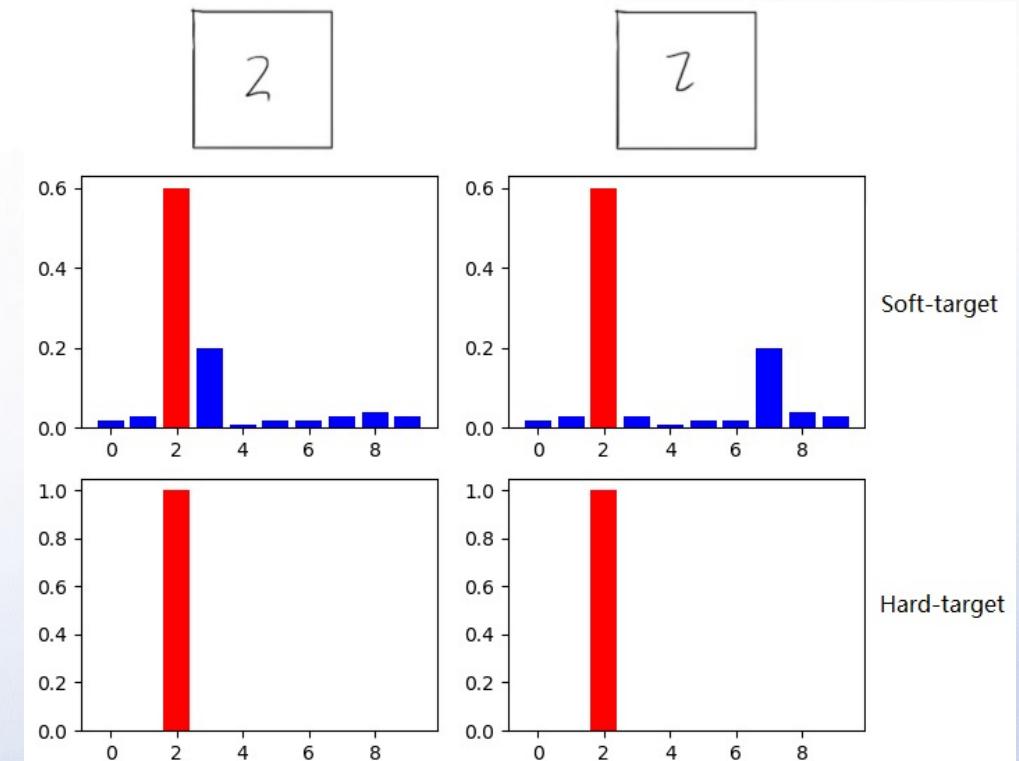
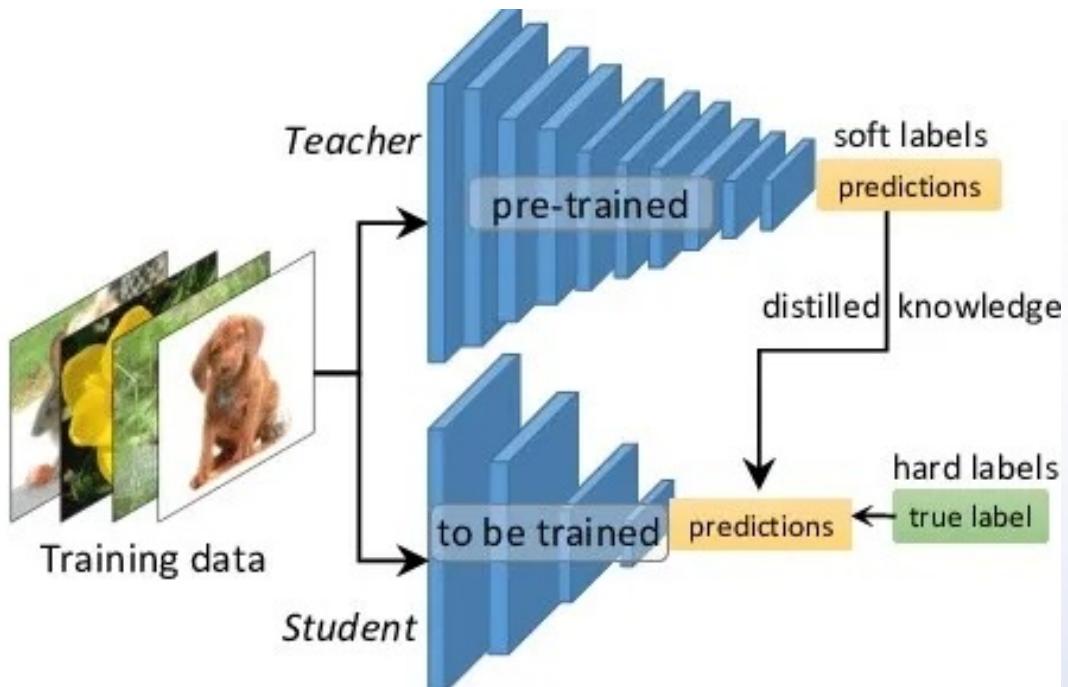
## ➤ Knowledge Distillation

- transferring the knowledge from a large model (teacher) to a smaller model (student)
- a form of **model compression**



# Introduction

## ➤ Soft target in knowledge distillation



# Introduction

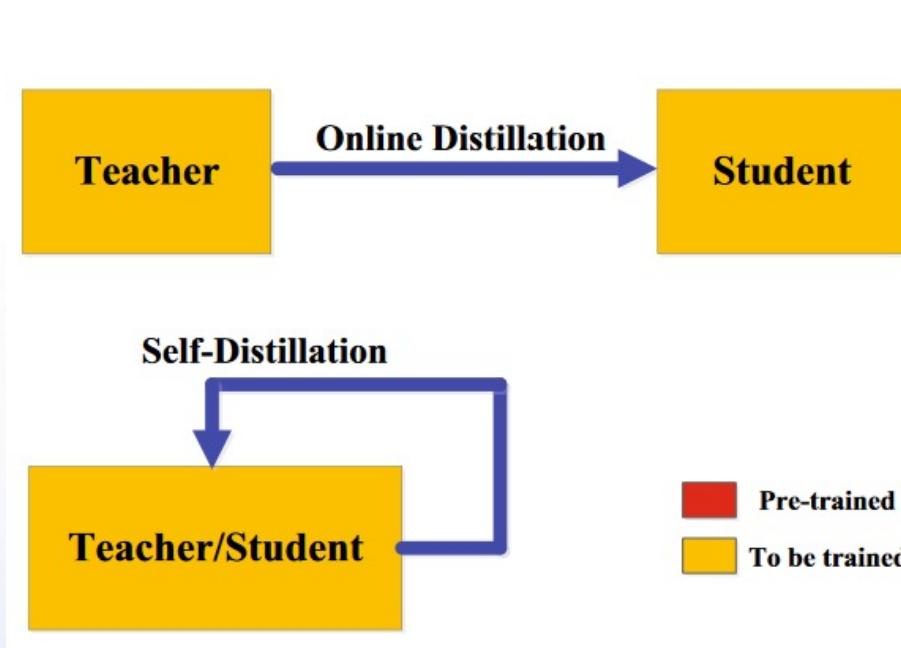
➤ Three types of knowledge distillation



**Offline Distillation**



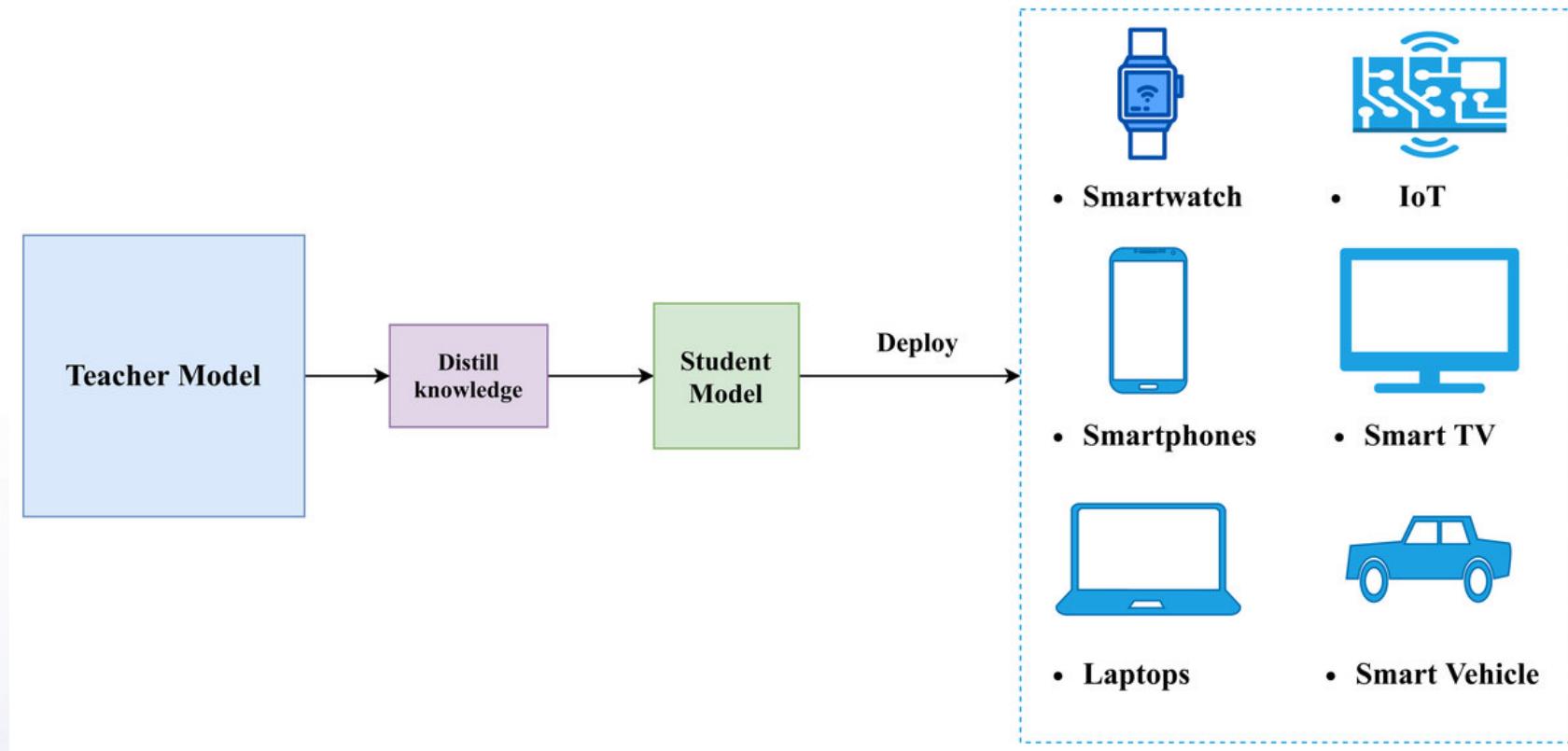
**Online Distillation**



**Self-Distillation**

# Introduction

## ➤ Applications of Knowledge Distillation



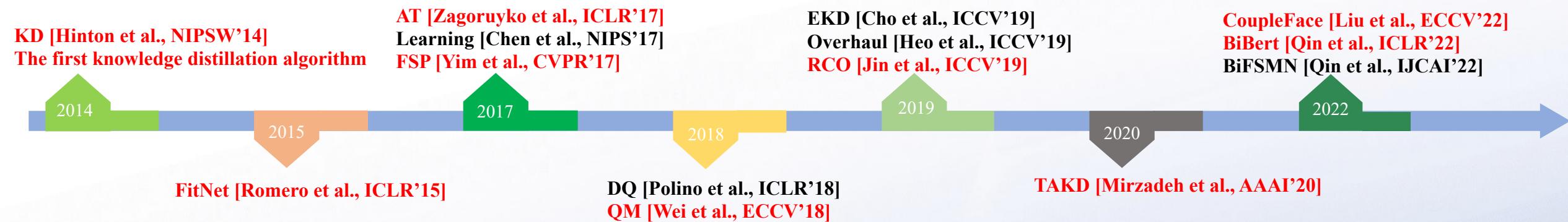
Use cases for knowledge distillation to deploy deep learning models  
on small devices with limited resources

# 目 录

- 1 Introduction
- 2 Knowledge distillation
- 3 Knowledge definition
- 4 Knowledge transfer
- 5 Application

# Knowledge Distillation: History

Knowledge distillation has gradually gained attention since 2015.



# Knowledge Distillation

## Definition of knowledge

Distilling the Knowledge in a Neural Network, NIPS'14

Fitnets: Hints for Thin Deep Nets, ICLR'15

Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer, ICLR'17

A Gift from Knowledge Distillation: Fast Optimization, Network Minimization and Transfer Learning, CVPR'17

## Transfer of knowledge

Knowledge Distillation via Route Constrained Optimization, ICCV'19

Improved Knowledge Distillation via Teacher Assistant, AAAI'20

## Application of KD

Quantization Mimic: Towards Very Tiny CNN for Object Detection, ECCV'18

CoupleFace: Relation Matters for Face Recognition Distillation, ECCV'22

BiBERT: Accurate Fully Binarized BERT, ICLR'22

# 目 录

- 1 Introduction
- 2 Knowledge distillation
- 3 Knowledge definition
- 4 Knowledge transfer
- 5 Application

# Definition of knowledge

---

## Distilling the Knowledge in a Neural Network

---

**Geoffrey Hinton\***<sup>†</sup>

Google Inc.

Mountain View

[geoffhinton@google.com](mailto:geoffhinton@google.com)

**Oriol Vinyals<sup>†</sup>**

Google Inc.

Mountain View

[vinyals@google.com](mailto:vinyals@google.com)

**Jeff Dean**

Google Inc.

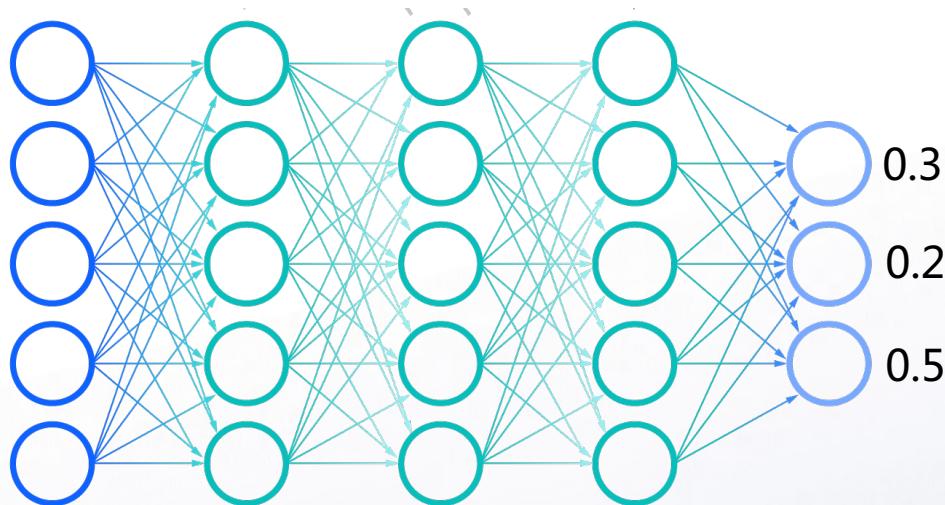
Mountain View

[jeff@google.com](mailto:jeff@google.com)

# Definition of knowledge

**Problem:** information in incorrect categories is ignored, during the training with cross entropy

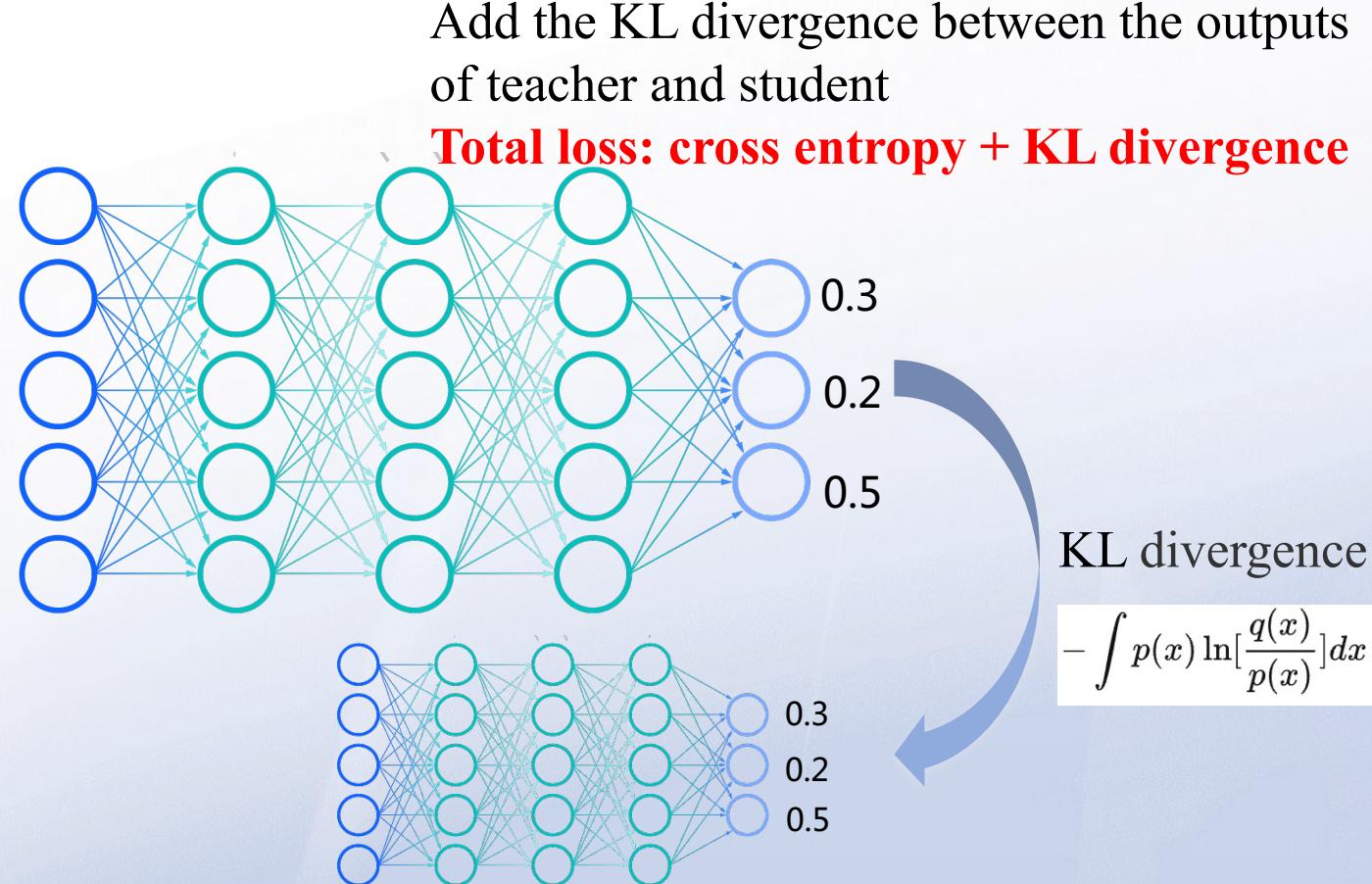
**Method:** define the outputs of large model as knowledge, imitating the teacher's output distribution



Cross entropy loss :

$$-(0x\log 0.3 + 0x\log 0.2 + 1x\log 0.5)$$

Missing incorrect category prediction information



# Definition of knowledge

**Result: 3x parameter compression on CIFAR10, 1.81% accuracy drop**

Teacher Model	Student Model	Teacher Accuracy	directly training student	Knowledge Distillation
WRN-16-2, 0.7M	WRN-16-1, 0.2M	93.69%	91.23%	92.59%
WRN-40-1, 0.6M	WRN-16-1, 0.2M	93.42%	91.23%	91.61%
WRN-40-2, 2.2M	WRN-16-2, 0.7M	94.77%	93.69%	93.92%

# Definition of knowledge

## FITNETS: HINTS FOR THIN DEEP NETS

**Adriana Romero<sup>1</sup>, Nicolas Ballas<sup>2</sup>, Samira Ebrahimi Kahou<sup>3</sup>, Antoine Chassang<sup>2</sup>,  
Carlo Gatta<sup>4</sup> & Yoshua Bengio<sup>2†</sup>**

<sup>1</sup>Universitat de Barcelona, Barcelona, Spain.

<sup>2</sup>Université de Montréal, Montréal, Québec, Canada. <sup>†</sup>CIFAR Senior Fellow.

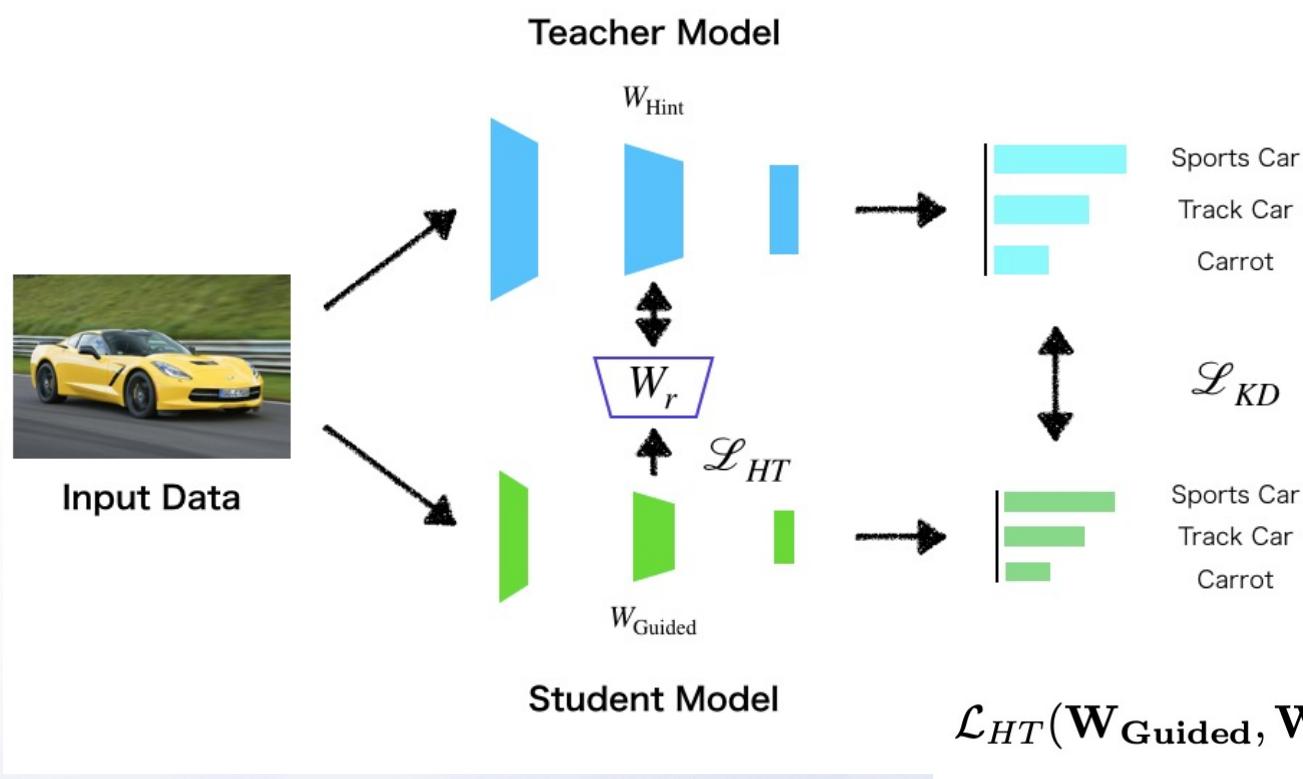
<sup>3</sup>École Polytechnique de Montréal, Montréal, Québec, Canada.

<sup>4</sup>Centre de Visió per Computador, Bellaterra, Spain.

# Definition of knowledge

**Problem:** the KD method ignores the information of the middle layer in the network

**Method:** define the features of middle layer as knowledge



**Total loss:** cross entropy + the L2 distance of middle features + (other knowledge loss)

# Definition of knowledge

## Result: improvements on multiple datasets

Network	# layers	# params	# mult	Acc	Speed-up	Compression rate
Teacher	5	~9M	~725M	90.18%	1	1
FitNet 1	11	~250K	~30M	89.01%	<b>13.36</b>	<b>36</b>
FitNet 2	11	~862K	~108M	91.06%	4.64	10.44
FitNet 3	13	~1.6M	~392M	91.10%	1.37	5.62
FitNet 4	19	~2.5M	~382M	<b>91.61%</b>	1.52	3.60

Table 5: Accuracy/Speed Trade-off on CIFAR-10.

Algorithm	# params	Accuracy
<i>Compression</i>		
FitNet	~2.5M	<b>91.61%</b>
Teacher	~9M	90.18%
Mimic single	~54M	84.6%
Mimic single	~70M	84.9%
Mimic ensemble	~70M	85.8%
<i>State-of-the-art methods</i>		
Maxout		90.65%
Network in Network		91.2%
Deeply-Supervised Networks		<b>91.78%</b>
Deeply-Supervised Networks (19)		88.2%

Table 1: Accuracy on CIFAR-10

Algorithm	# params	Accuracy
<i>Compression</i>		
FitNet	~2.5M	<b>64.96%</b>
Teacher	~9M	63.54%
<i>State-of-the-art methods</i>		
Maxout		61.43%
Network in Network		64.32%
Deeply-Supervised Networks		<b>65.43%</b>

Table 2: Accuracy on CIFAR-100

**4x computation reduction,  
0.12% accuracy drop**

# Definition of knowledge

## PAYING MORE ATTENTION TO ATTENTION: IMPROVING THE PERFORMANCE OF CONVOLUTIONAL NEURAL NETWORKS VIA ATTENTION TRANSFER

**Sergey Zagoruyko, Nikos Komodakis**

Université Paris-Est, École des Ponts ParisTech

Paris, France

{sergey.zagoruyko,nikos.komodakis}@enpc.fr

# Definition of knowledge

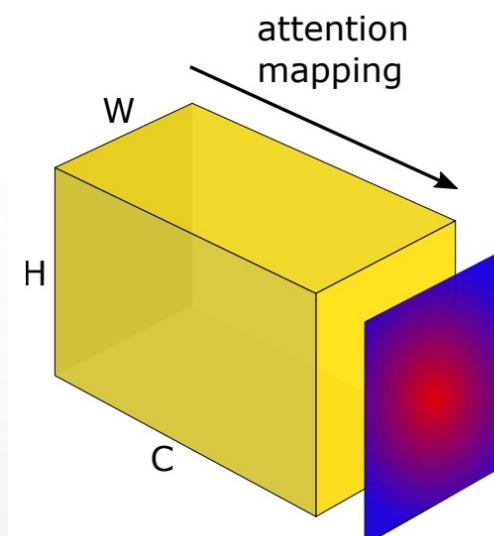
**Problem:** the different parts in a picture have different importance

**Method:** define the attention as knowledge



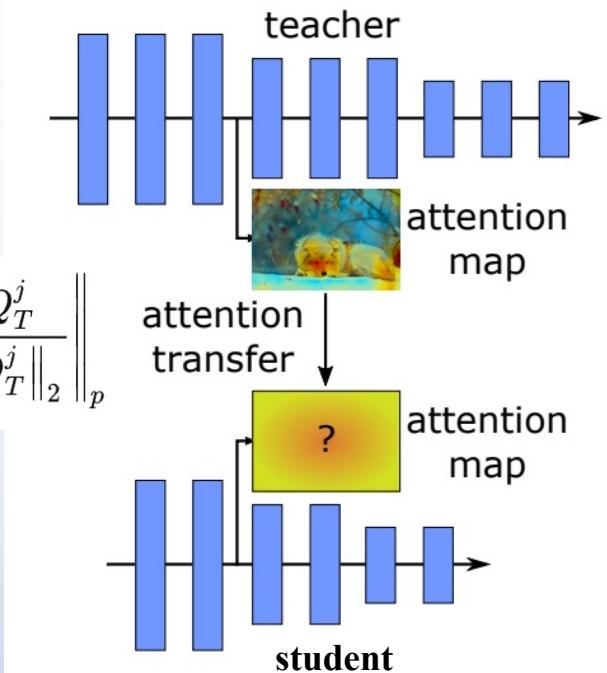
The central part of this picture has **more information**.

The network needs to pay more attention to the **main part**.



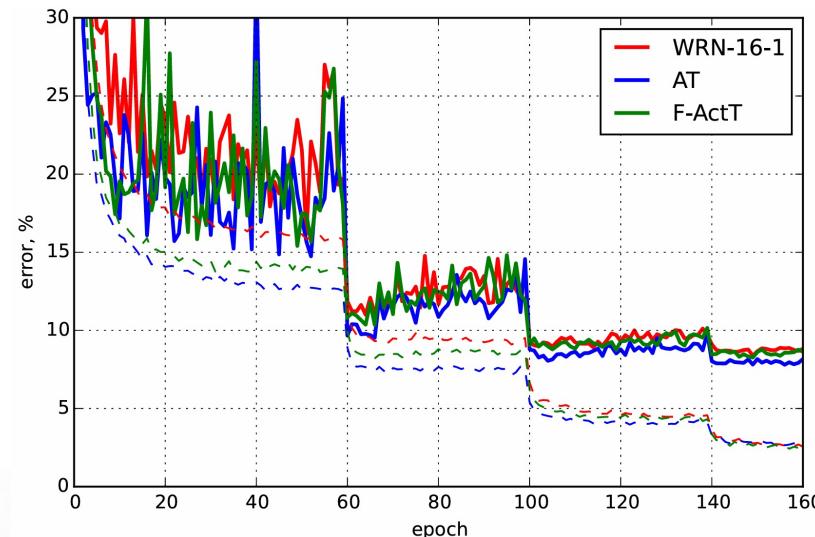
**Total loss: cross entropy + L2 distance of attention maps + (other knowledge loss)**

$$\frac{\beta}{2} \sum_{j \in \mathcal{I}} \left\| \frac{Q_S^j}{\|Q_S^j\|_2} - \frac{Q_T^j}{\|Q_T^j\|_2} \right\|_p$$



# Definition of knowledge

**Result: 3x parameter compression, 1.67% accuracy drop**



accelerate the convergence speed  
and improve final accuracy

Teacher	Student	Teacher	directly training student	AT	AT+KD
WRN-16-2, 0.7M	WRN-16-1, 0.2M	93.69%	91.23%	92.07%	94.49%
WRN-40-1, 0.6M	WRN-16-1, 0.2M	93.42%	91.23%	91.75%	91.99%
WRN-40-2, 2.2M	WRN-16-2, 0.7M	94.77%	93.69%	94.15%	94.29%

# Definition of knowledge

## A Gift from Knowledge Distillation: Fast Optimization, Network Minimization and Transfer Learning

Junho Yim<sup>1</sup>      Donggyu Joo<sup>1</sup>      Jihoon Bae<sup>2</sup>      Junmo Kim<sup>1</sup>

<sup>1</sup>School of Electrical Engineering, KAIST, South Korea

<sup>2</sup>Electronics and Telecommunications Research Institute

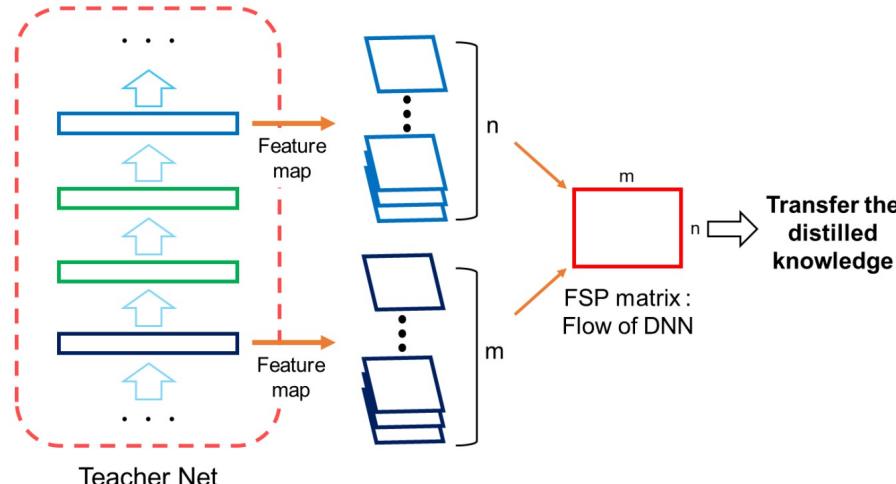
{junho.yim, jdg105, junmo.kim}@kaist.ac.kr

{baejh}@etri.re.kr

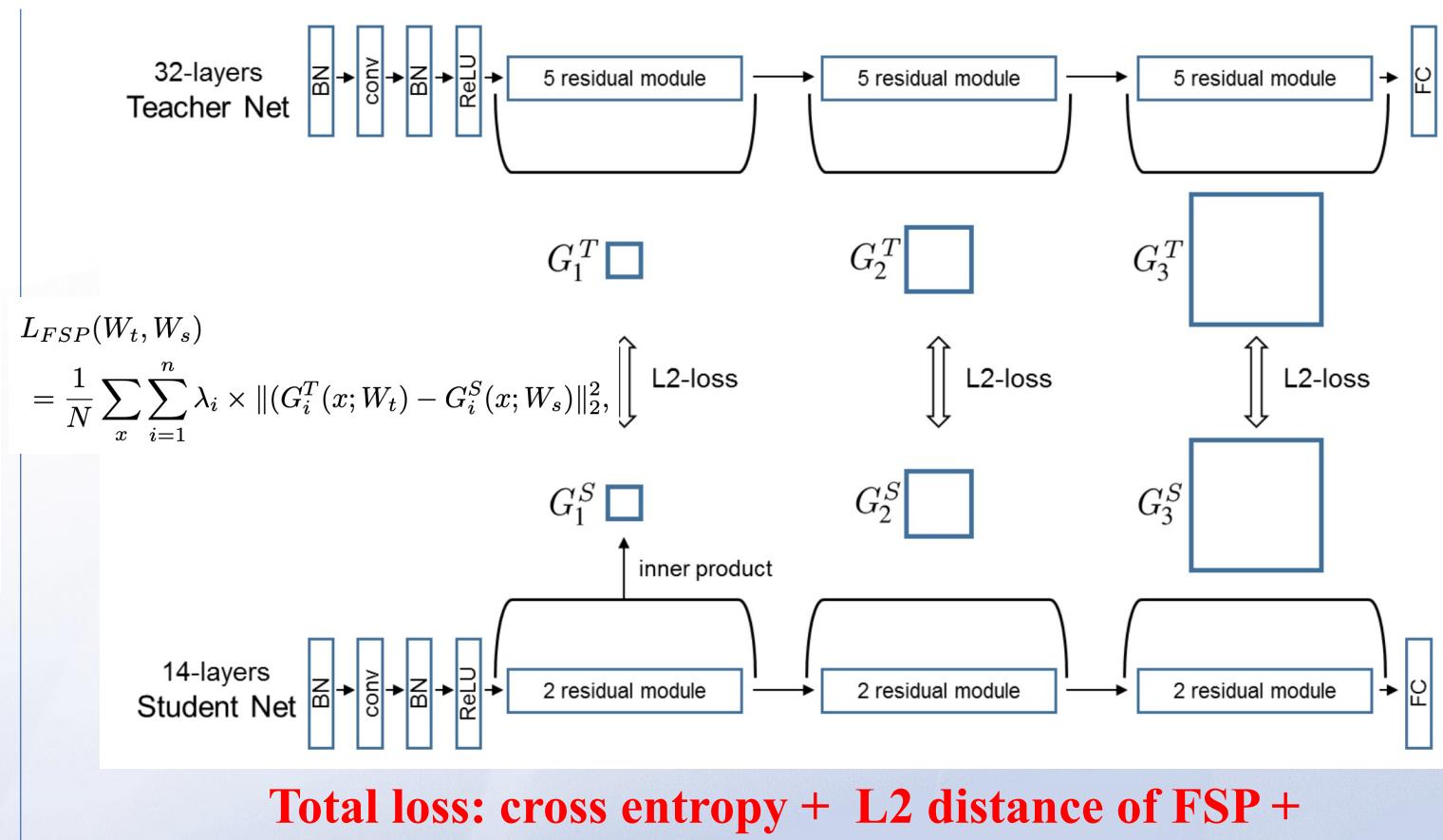
# Definition of knowledge

**Problem:** previous methods ignore the relationship between network layers

**Method:** define the relationship between features from two layers as knowledge



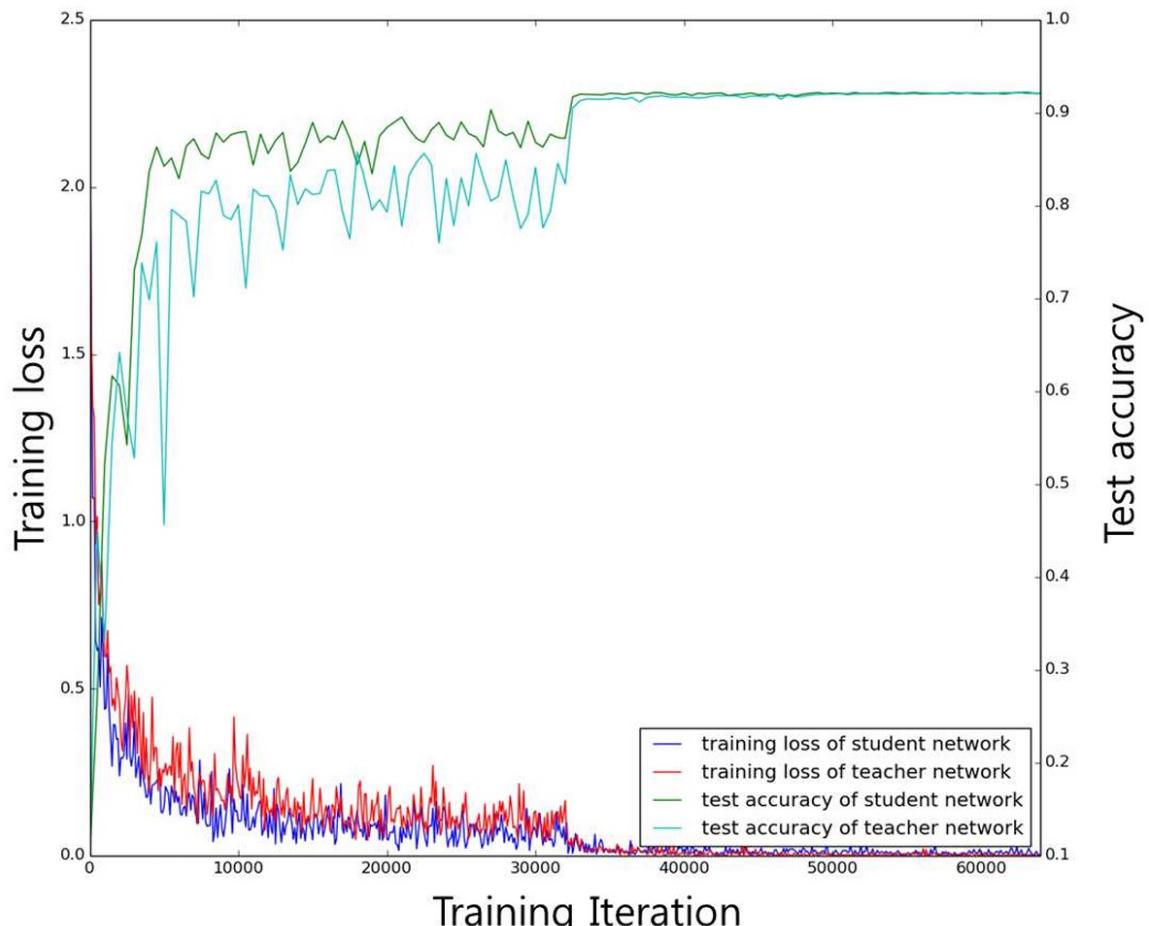
Inner product of features  
(FSP matrix)



Total loss: cross entropy + L2 distance of FSP +  
(other knowledge loss)

# Definition of knowledge

**Result: 3.21% accuracy drop from ResNet-26 to ResNet-8 on CIFAR10**



	Accuracy
Teacher-original	91.91
Student-original	87.91
FitNet [20]	88.57
Proposed Method	88.70

	Net 1	Net 2	Net 3	Avg	Ensemble	#Iter
Teacher	91.61	91.56	92.09	91.75	93.48	192k
Teacher *	90.47	90.83	90.62	90.64	92.6	63k
Teacher ‡	91.84	92.26	92.01	92.04	92.71	63k
1 loss FitNet [20]*	91.69	91.85	91.64	91.72	92.98	98k
3 loss FitNet [20]*	88.90	89.35	89.02	89.09	89.92	98k
Student *	92.28	92.08	92.07	92.14	93.26	84k
Student *†	92.28	91.89	92.08	92.08	93.67	126k

# Definition of knowledge

## SCOTT: Self-Consistent Chain-of-Thought Distillation

**Peifeng Wang<sup>1\*</sup>, Zhengyang Wang<sup>2</sup>, Zheng Li<sup>2</sup>, Yifan Gao<sup>2</sup>, Bing Yin<sup>2</sup>, Xiang Ren<sup>1</sup>**

<sup>1</sup>Department of Computer Science, University of Southern California, <sup>2</sup>Amazon.com Inc  
{peifengw, xiangren}@usc.edu,  
{zhengywa, amzzhe, yifangao, alexbyin}@amazon.com

# Definition of knowledge

**Problem:** chain-of-thought (CoT) prompting makes LMs to explain their multi-step reasoning, but the performance gains are typically only observed with **sufficiently large language models**

**Method:** **rationales** generated from a large language model (teacher) that are more grounded by the gold answers, and trains the student language model using a **counterfactual reasoning objective**

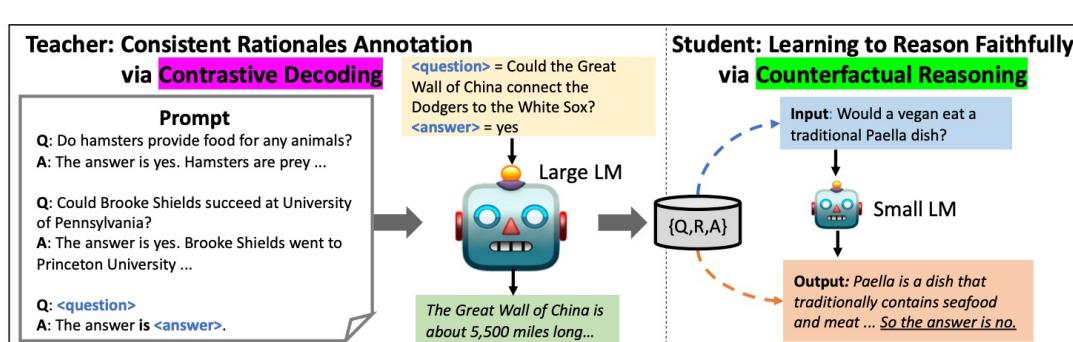
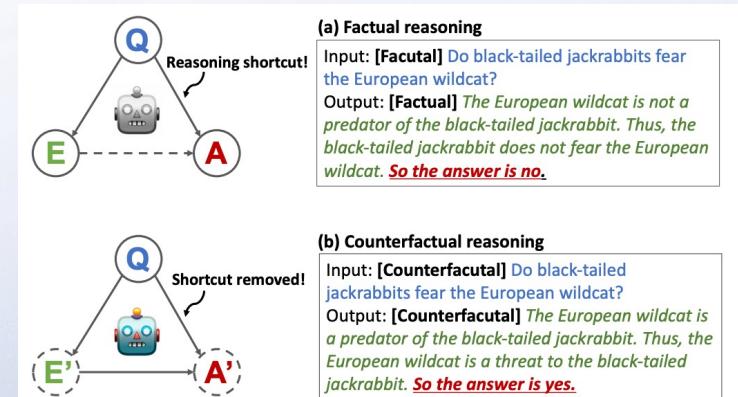


Figure 2: Overview of our knowledge distillation framework for faithful reasoning. (a) Teacher: A large LM prompted to generate a consistent rationale given a question and the gold answer in the training set via contrastive decoding. (b) Student: A small LM fine-tuned to generate a rationale and then answer via counterfactual reasoning.

Only search rationale tokens that are more plausible for the gold answer by contrastive decoding

$$G(t_i|a^*) = \log \frac{P(t_i|p, q, a^*, t_{<i})}{P(t_i|p, q, a', t_{<i})}.$$

appending the keyword **[Factual]** or **[Counterfactual]**



Sum the loss  
and train  
jointly

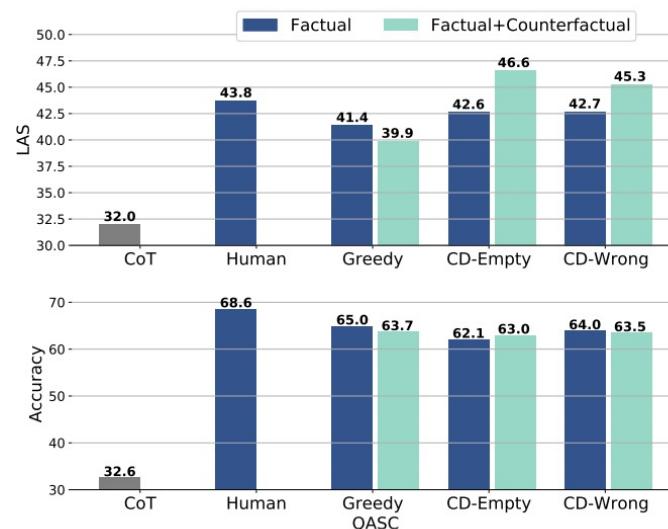
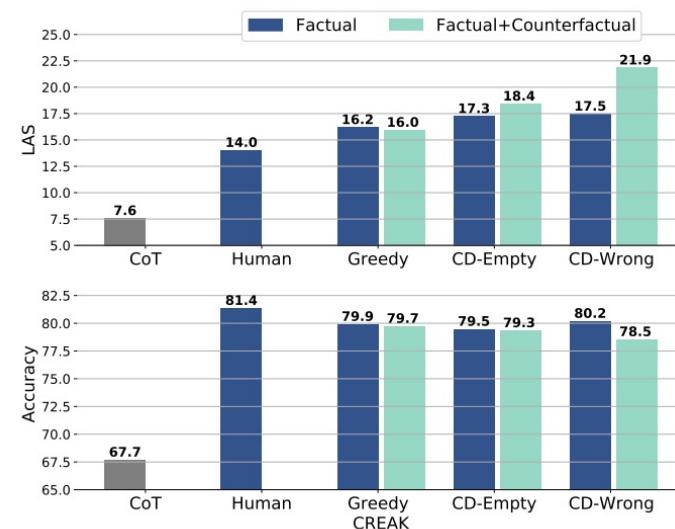
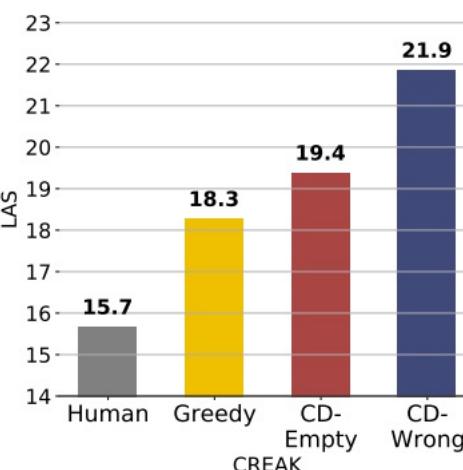
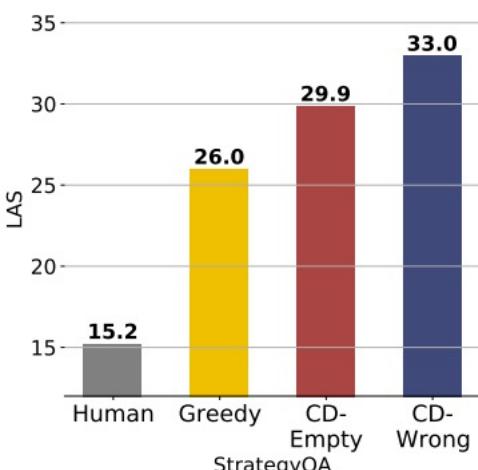
$$\mathcal{L}_{factual} = - \sum_i \log P(t_i|q, t_{<i}),$$

$$\mathcal{L}_{counterfactual} = - \sum_i \log P(t_i|q, r', t_{<i}).$$

# Definition of knowledge

**Result:** generate CoT rationales that are more **faithful**

Teacher Model	Grammaticality	New Info	Supports Answer
Greedy	<b>0.99</b>	0.65	0.48
Contrast.-Empty	0.97	0.77	0.58
Contrast.-Wrong	0.97	<b>0.82</b>	<b>0.63</b>



{Greedy, CD-Empty, CD-Wrong} refer respectively to using greedy decoding, contrastive decoding with empty/wrong answer to obtain rationale tokens from the teacher.

# Definition of knowledge

## ON-POLICY DISTILLATION OF LANGUAGE MODELS: LEARNING FROM SELF-GENERATED MISTAKES

Rishabh Agarwal<sup>1,2\*</sup>

Nino Vieillard<sup>1\*</sup>

Yongchao Zhou<sup>13</sup>

Piotr Stanczyk<sup>1†</sup>

Sabela Ramos<sup>1†</sup>

Matthieu Geist<sup>1</sup>

Olivier Bachem<sup>1</sup>

<sup>1</sup>Google DeepMind

<sup>2</sup>Mila

<sup>3</sup>University of Toronto

# Definition of knowledge

**Problem:** current KD methods for auto-regressive sequence models suffer from distribution mismatch

**Method:** introduce **on-policy Knowledge Distillation**, which trains the student on its self-generated output sequences

$$L_{OD}(\theta) := \mathbb{E}_{x \sim X} \left[ \mathbb{E}_{y \sim p_S(\cdot|x)} [\mathcal{D}_{KL}(p_T \parallel p_S^\theta)(y \mid x)] \right]$$

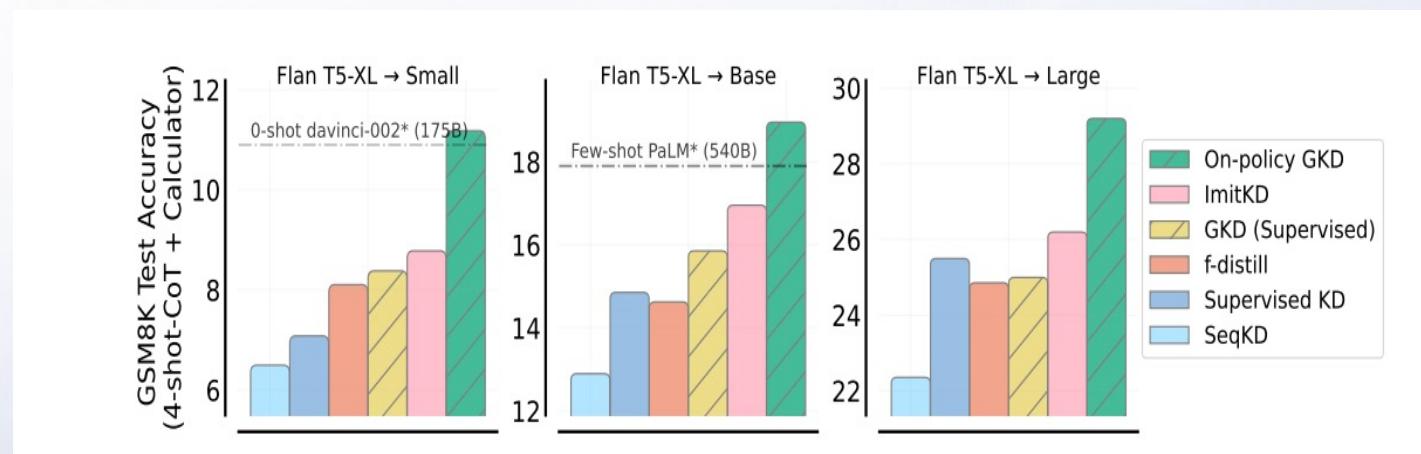
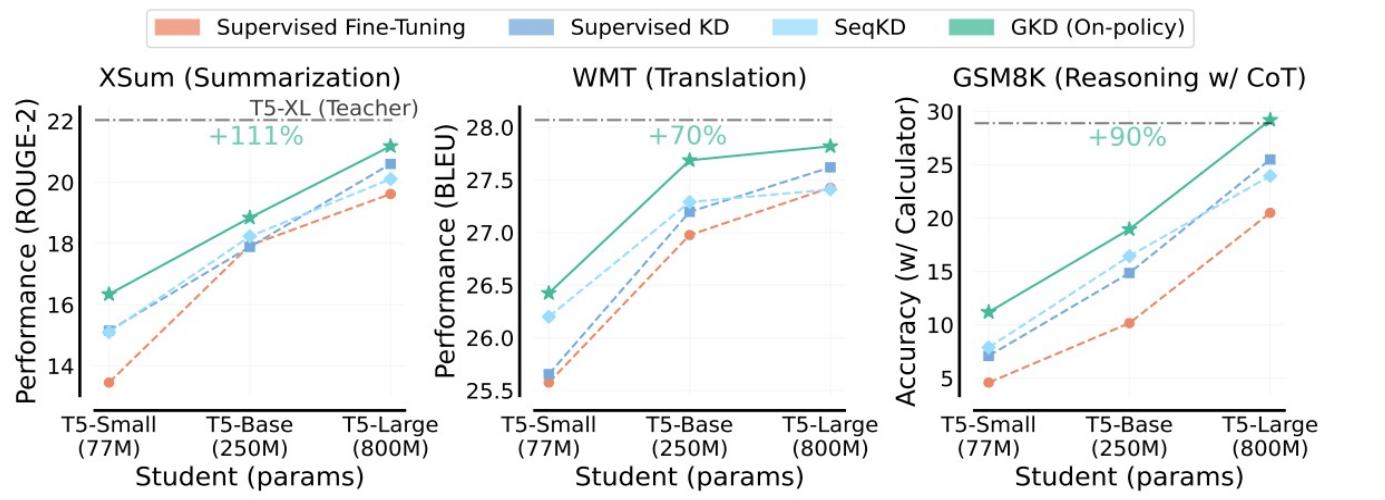
student receives token-specific feedback from the teacher's logits

$$L_{GKD}(\theta) := (1 - \lambda) \mathbb{E}_{(x,y) \sim (X,Y)} [\mathcal{D}(p_T \parallel p_S^\theta)(y \mid x)] + \lambda \mathbb{E}_{x \sim X} \left[ \mathbb{E}_{y \sim p_S(\cdot|x)} [\mathcal{D}(p_T \parallel p_S^\theta)(y \mid x)] \right]$$

use a mixture of teacher-generated and on-policy student-generated sequences for training

# Definition of knowledge

**Result: on-policy GKD outperforms commonly-used KD approaches**



# 目 录

- 1 Introduction
- 2 Knowledge distillation
- 3 Knowledge definition
- 4 Knowledge transfer
- 5 Application

# Transfer of knowledge

## Knowledge Distillation via Route Constrained Optimization

Xiao Jin<sup>1\*</sup>, Baoyun Peng<sup>2\*</sup>, Yichao Wu<sup>1</sup>, Yu Liu<sup>3</sup>, Jiaheng Liu<sup>4</sup>,  
Ding Liang<sup>1</sup>, Junjie Yan<sup>1</sup>, Xiaolin Hu<sup>5</sup>

<sup>1</sup> SenseTime Group Limited

<sup>2</sup> National University of Defense Technology

<sup>3</sup> Chinese University of Hong Kong

<sup>4</sup> Beihang University

<sup>5</sup> Tsinghua University

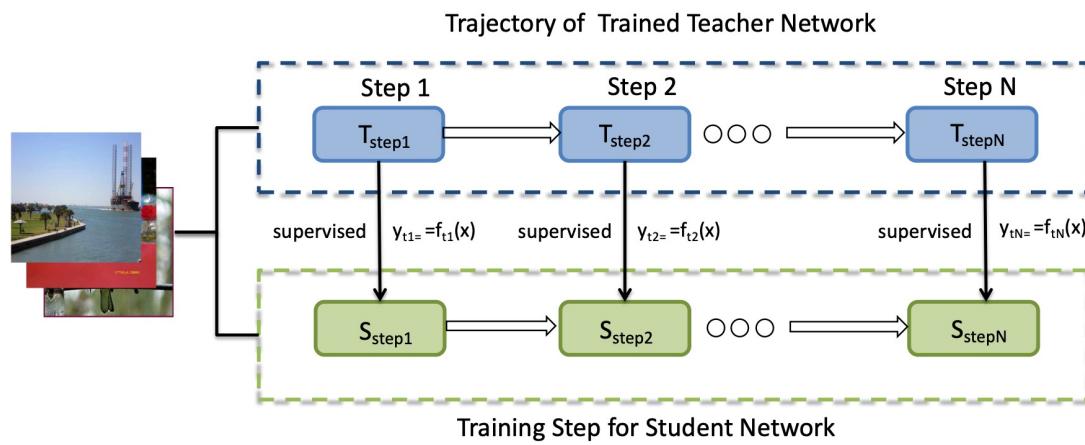
jinxiaocuhk@gmail.com, pengbaoyun13@nudt.edu.cn, yuliu@ee.cuhk.edu.hk, liujiaheng@buaa.edu.cn

{wuyichao, liangding, yanjunjie}@sensetime.com, xlhu@mail.tsinghua.edu.cn

# Transfer of knowledge

**Problem:** there is a significant performance gap between the pre-trained teacher model and the student model

**Method:** simultaneous learning between teachers and students



imitate the teacher network's checkpoint step by step



avoid significant gap between teacher and student

---

## Algorithm 2 Greedy Search

---

**Require:** Student network with parameter  $W_s$  after mimicking former  $i^{th}$  anchor point  $C_i$ , where  $i \in \{1, 2, \dots, N\}$ , relaxation factor  $\delta$ .  
compute KL divergence  $\mathcal{H}_i$   
 $j = i + 1$   
**while**  $j < N$  **do**  
    compute  $\mathcal{H}_j$  on validation set  
    compute  $r_{ij} = \frac{\mathcal{H}_j - \mathcal{H}_i}{\mathcal{H}_i}$   
    **if**  $r_{ij} > \delta$  **then**  
        Return  $j-1$ ;  
    **end if**  
     $j = j + 1$   
**end while**  
Return  $N$ ;

---

greedy search to find the optimal checkpoint

# Transfer of knowledge

**Result: 68.21% accuracy from ResNet-50 to MobileNetV2 on ImageNet**

Method	Network	top-1	top-5
Teacher-Softmax	ResNet-50	75.49	92.48
Student-Softmax	MobileNetV2	64.2	85.4
Student-KD	MobileNetV2	66.75	87.3
Student-RCO	MobileNetV2	<b>68.21</b>	<b>88.04</b>

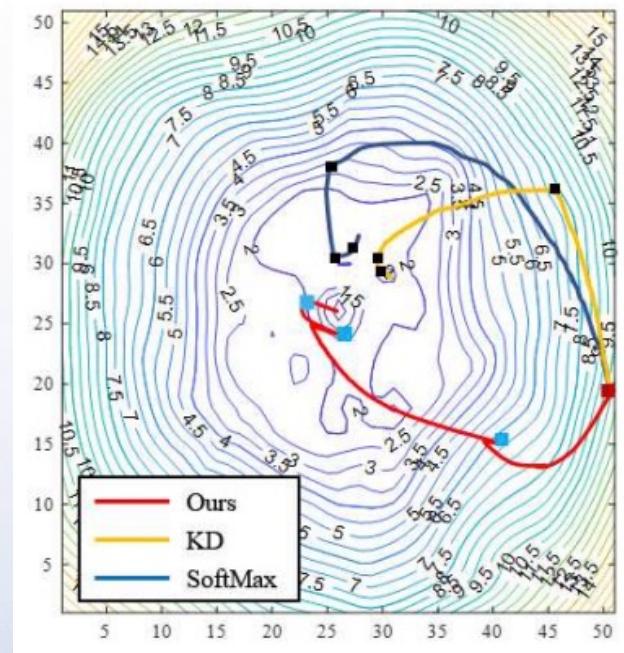
Table 4: Results on ImageNet

Method	top-1 @ distractor size					
	$e^1$	$e^2$	$e^3$	$e^4$	$e^5$	$e^6$
Teacher	99.78	99.67	99.38	98.86	97.70	94.83
Softmax	99.20	96.37	91.49	84.45	75.60	65.91
FitNet	99.62	98.80	96.83	93.53	88.28	81.02
RCO	<b>99.69</b>	<b>99.01</b>	<b>97.52</b>	<b>94.84</b>	<b>90.55</b>	<b>84.3</b>

Table 5: Results on MegaFace

**68.21% accuracy, 13x computation reduction on ImageNet**

**higher performance than FitNet on MegaFace**



**better training trajectories**

# Transfer of knowledge

## Improved Knowledge Distillation via Teacher Assistant

**Seyed Iman Mirzadeh,<sup>\*1</sup> Mehrdad Farajtabar,<sup>\*2</sup> Ang Li,<sup>2</sup>  
Nir Levine,<sup>2</sup> Akihiro Matsukawa,<sup>†3</sup> Hassan Ghasemzadeh<sup>1</sup>**

<sup>1</sup>Washington State University, WA, USA

<sup>2</sup>DeepMind, CA, USA

<sup>3</sup>D.E. Shaw, NY, USA

<sup>1</sup>{seyediman.mirzadeh, hassan.ghasemzadeh}@wsu.edu

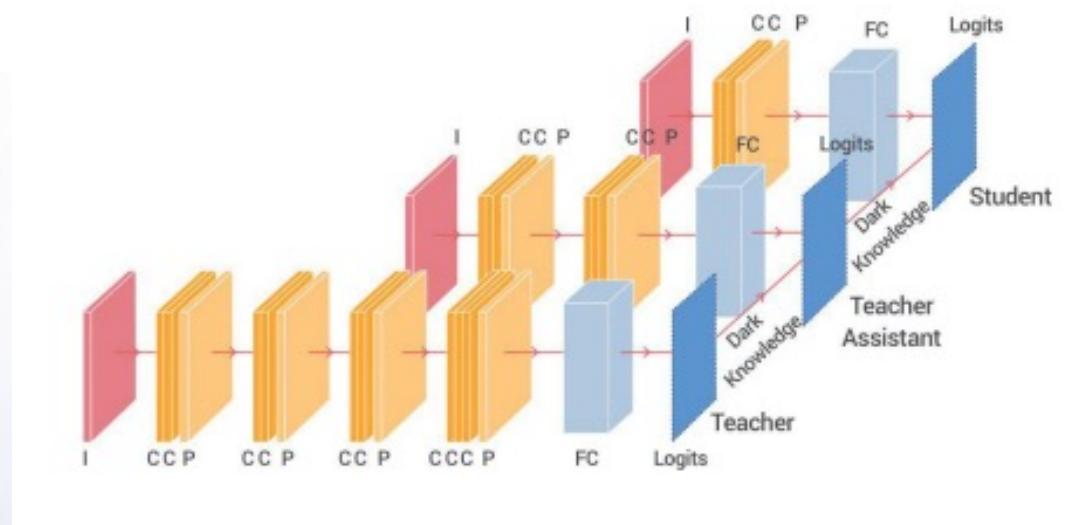
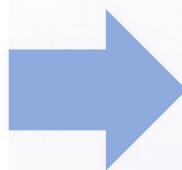
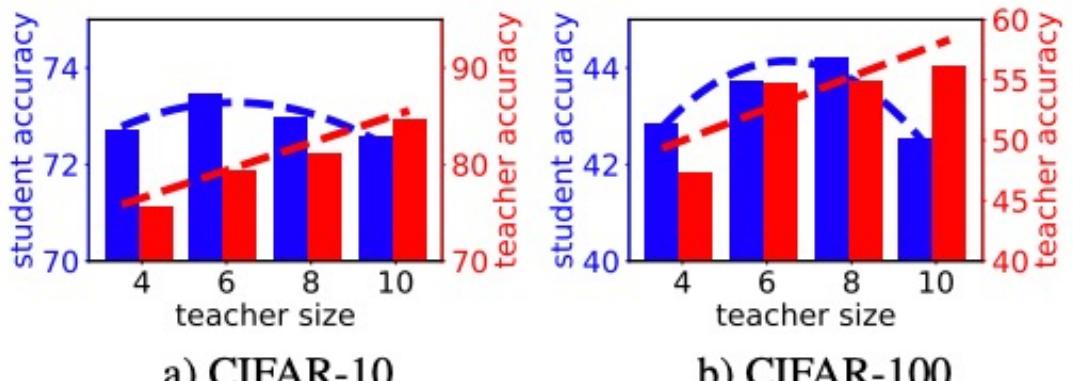
<sup>2</sup>{farajtabar, anglili, nirlevine}@google.com

<sup>3</sup>akihiro.matsukawa@gmail.com

# Transfer of knowledge

**Problem:** the **significant difference** between teacher and student leads to a decrease in student's performance

**Method:** introduce **teacher assistant** network, transferring knowledge **gradually** to student



When the gap between teacher and student is too large, the performance of student will decrease.

introduce **teacher assistant** network to assist knowledge transfer

# Transfer of knowledge

## Result: higher performance on multiple datasets

How to determine the size of teacher assistant model?

**Empirical setting:** when training from scratch, the teaching assistant model falls between the teacher and student.

Table 2: Student's accuracy given varied TA sizes for (S=2, T=10)

Model	Dataset	TA=8	TA=6	TA=4
CNN	CIFAR-10	72.75	73.15	73.51
	CIFAR-100	44.28	44.57	44.92

Table 3: Student's accuracy given varied TA sizes for (S=8, T=110)

Model	Dataset	TA=56	TA=32	TA=20	TA=14
ResNet	CIFAR-10	88.70	88.73	88.90	88.98
	CIFAR-100	61.47	61.55	61.82	61.5

Table 1: Comparison on evaluation accuracy between our method (TAKD) and baselines. For CIFAR, plain (S=2, TA=4, T=10) and for ResNet (S=8, TA=20, T=110) are used. For ImageNet, ResNet (S=14, TA=20, T=50) is used. Higher numbers are better.

Model	Dataset	NOKD	BLKD	TAKD
CNN	CIFAR-10	70.16	72.57	73.51
	CIFAR-100	41.09	44.57	44.92
ResNet	CIFAR-10	88.52	88.65	88.98
	CIFAR-100	61.37	61.41	61.82
ResNet	ImageNet	65.20	66.60	67.36

**68.21% accuracy from ResNet-50 to ResNet-14 on ImageNet**

# Transfer of knowledge

---

## Efficient Knowledge Distillation from Model Checkpoints

---

**Chaofei Wang\*, Qisen Yang\*, Rui Huang, Shiji Song, Gao Huang<sup>†</sup>**

Department of Automation, Tsinghua University, China

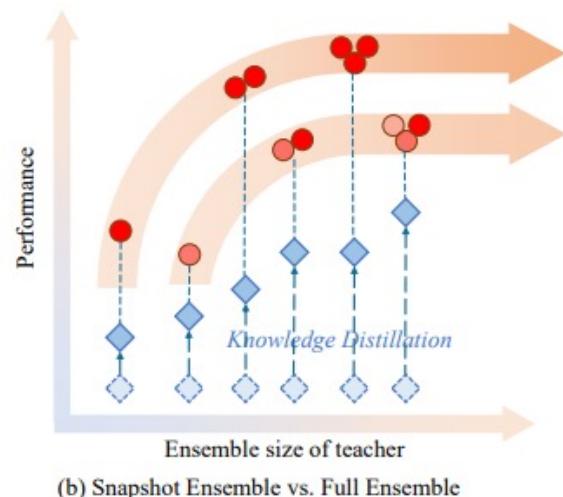
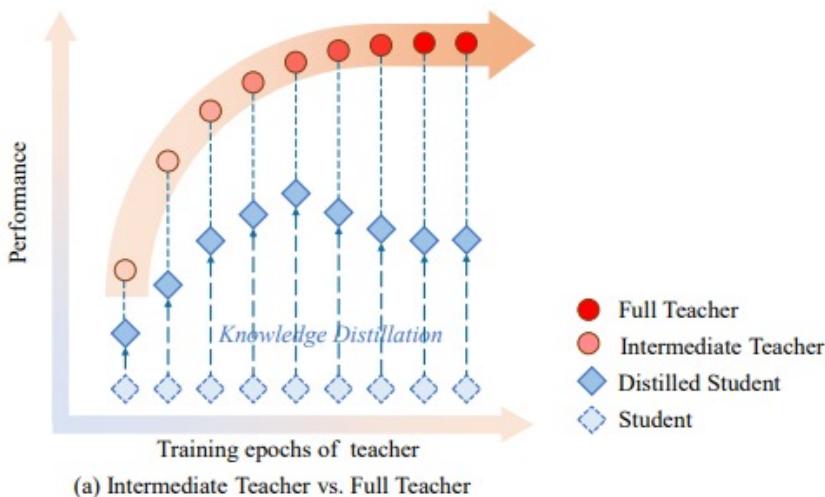
wangcf18, yangqs19, hr20@mails.tsinghua.edu.cn

shijis, gaochuang@tsinghua.edu.cn

# Transfer of knowledge

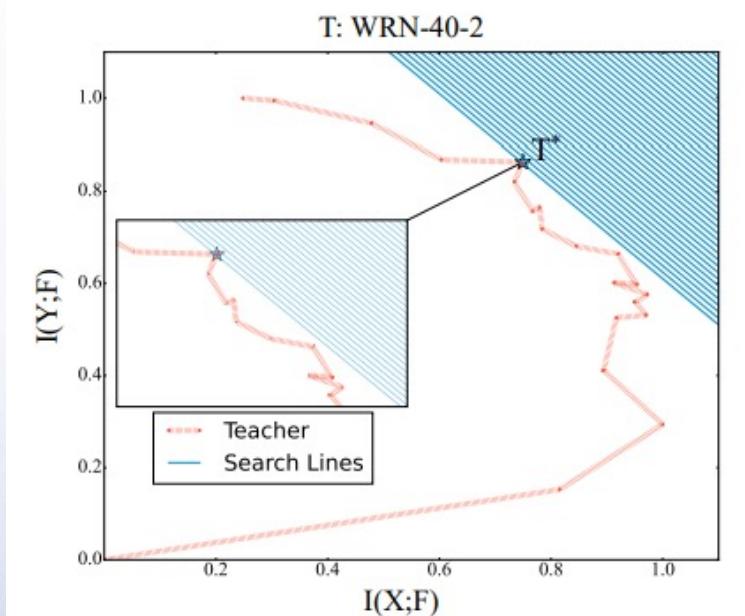
**Problem:** fully trained teachers do not always provide the best knowledge distillation

**Method:** select optimal intermediate teacher checkpoints by **maximizing mutual information**



A weak **intermediate model** can serve as a better teacher than the strong fully converged model.

$$\max_F \{I(X; F) + I(Y; F)\}$$



select the optimal intermediate teacher

# Transfer of knowledge

**Result: improved distillation performance by selecting intermediate teachers**

Network structure		Accuracy of T&S		KD accuracy of different intermediate teachers					
T	S	T	S	$T^{0.3}$	$T^{0.5}$	$T^{0.7}$	$T^{\text{full}}$	$T^*$	
WRN-40-2	WRN-40-1	76.53	70.38	72.34±0.10	72.76±0.24	73.08±0.05	72.68±0.10	<b>73.26±0.03</b>	
	MobileNetV2		64.49	68.21±0.33	<b>68.99±0.12</b>	68.54±0.07	68.03±0.34	68.58±0.34	
ResNet-110	ResNet-32	73.41	70.16	70.74±0.18	72.49±0.32	72.46±0.30	72.48±0.22	<b>72.63±0.13</b>	
	MobileNetV2		64.49	67.84±0.26	68.79±0.17	<b>69.01±0.20</b>	68.63±0.35	68.99±0.33	
Average		74.97	67.38	69.78	70.76	70.77	70.46	<b>70.87</b>	

# Transfer of knowledge

---

## MiniLLM: Knowledge Distillation of Large Language Models

---

**Yuxian Gu<sup>1,2\*</sup>, Li Dong<sup>2</sup>, Furu Wei<sup>2</sup>, Minlie Huang<sup>1†</sup>**

<sup>1</sup>The CoAI Group, Tsinghua University

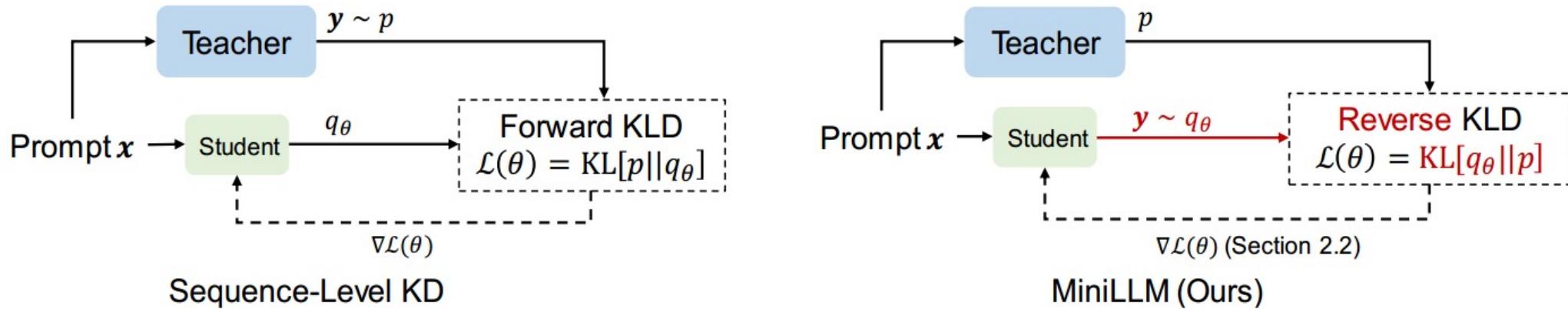
<sup>2</sup>Microsoft Research

guyx21@mails.tsinghua.edu.cn {lidong1,fuwei}@microsoft.com  
aihuang@tsinghua.edu.cn

# Transfer of knowledge

**Problem:** In open-ended text generation tasks , student model overestimates the low-probability regions of the teacher distribution

**Method:** replace the standard KLD with reverse KLD



reverse KLD causes  $q_\theta$  to seek the major modes of  $p$ , and assign low probabilities to  $p$ 's void regions



$$\begin{aligned}\theta &= \arg \min_{\theta} \mathcal{L}(\theta) = \arg \min_{\theta} \text{KL}[q_\theta || p] \\ &= \arg \min_{\theta} \left[ - \mathbb{E}_{x \sim p_x, y \sim q_\theta} \log \frac{p(y | x)}{q_\theta(y | x)} \right]\end{aligned}$$

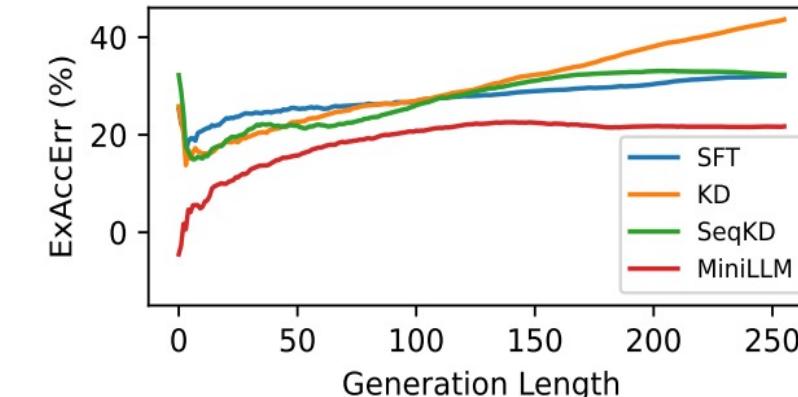
# Transfer of knowledge

**Result:** In contrast to previous KD methods, MINILLM narrows the performance gap between the student and the teacher

Model	#Params	Method	DollyEval		SelfInst		VicunaEval		S-NI		UnNI	
			GPT4	R-L	GPT4	R-L	GPT4	R-L	R-L	R-L	R-L	R-L
GPT-2	1.5B	Teacher	58.4	27.6	42.9	14.3	48.6	16.3	27.6	31.8		
		SFT w/o KD	38.6	23.3	26.3	10.0	32.8	14.7	16.3	18.5		
		KD	40.3	22.8	27.8	10.8	31.9	13.4	19.7	22.0		
		SeqKD	41.2	22.7	26.2	10.1	31.0	14.3	16.4	18.8		
	340M	MINILLM	<b>44.7</b>	<b>24.6</b>	<b>29.2</b>	<b>13.2</b>	<b>34.1</b>	<b>16.9*</b>	<b>25.3</b>	<b>26.6</b>		
		SFT w/o KD	51.9	<b>25.5</b>	39.6	13.0	42.3	16.0	25.1	32.0		
		KD	51.6	25.0	39.2	12.0	42.8	15.4	23.7	31.0		
		SeqKD	50.5	25.3	39.0	12.6	<b>43.0</b>	16.9*	22.9	30.2		
	760M	MINILLM	<b>52.2</b>	25.4	<b>40.5</b>	<b>15.6</b>	42.6	<b>17.7*</b>	<b>27.4</b>	<b>34.5</b>		
		SFT w/o KD	50.7	25.4	38.3	12.4	43.1	16.1	21.5	27.1		
		KD	53.4	25.9	40.4	13.4	43.4	16.9*	25.3	31.7		
		SeqKD	52.0	25.6	38.9	14.0	42.4	15.9	26.1	32.9		
OPT	1.3B	MINILLM	<b>54.7</b>	<b>26.4</b>	<b>44.6*</b>	<b>15.9</b>	<b>45.7</b>	<b>18.3*</b>	<b>29.3*</b>	<b>37.7*</b>		
		Teacher	70.3	29.2	56.1	18.4	58.0	17.8	30.4	36.1		
		SFT w/o KD	52.6	26.0	37.7	11.4	40.5	15.6	23.1	28.4		
		KD	52.7	25.4	36.0	12.2	40.8	14.9	21.9	27.0		
	2.7B	SeqKD	51.0	26.1	36.6	12.7	42.6	16.6	21.4	28.2		
		MINILLM	<b>60.7</b>	<b>26.7</b>	<b>47.0</b>	<b>14.8</b>	<b>50.6</b>	<b>17.9*</b>	<b>28.6</b>	<b>33.4</b>		
		SFT w/o KD	55.4	27.1	38.9	13.9	44.8	16.6	24.9	32.3		
		KD	60.5	25.9	48.6	13.8	51.3	16.7	26.3	30.2		
	6.7B	SeqKD	57.6	27.5	40.5	13.3	44.5	16.5	25.3	32.3		
		MINILLM	<b>63.2</b>	<b>27.4</b>	<b>52.7</b>	<b>17.2</b>	<b>55.9</b>	<b>19.1*</b>	<b>30.7*</b>	<b>35.1</b>		
		SFT w/o KD	67.9	27.6	56.4	16.4	57.3	17.8	30.3	28.6		
		KD	68.6	28.3	58.0	17.0	57.0	17.5	30.7*	26.7		
LLaMA	13B	SeqKD	69.6	28.5	54.0	17.0	57.6	17.9*	30.4	28.2		
		MINILLM	<b>70.8*</b>	<b>29.0</b>	<b>58.5*</b>	<b>17.5</b>	<b>60.1*</b>	<b>18.7*</b>	<b>32.5*</b>	<b>36.7*</b>		
		Teacher	79.0	29.7	75.5	23.4	65.1	19.4	35.8	38.5		
		SFT w/o KD	73.0	26.3	69.2	20.8	61.6	17.5	32.4	35.8		
	7B	KD	73.7	27.4	70.5	20.2	62.7	18.4	33.7	37.9		
		SeqKD	73.6	27.5	71.5	20.8	62.6	18.1	33.7	37.6		
		MINILLM	<b>76.4</b>	<b>29.0</b>	<b>73.1</b>	<b>23.2</b>	<b>64.1</b>	<b>20.7*</b>	<b>35.5</b>	<b>40.2*</b>		

	SST2		BoolQ	
	ECE	Acc.	ECE	Acc.
Teacher	0.025	93.0	0.356	74.5
KD	0.191	84.7	0.682	63.5
SeqKD	0.243	66.5	0.681	62.8
MINILLM	<b>0.099</b>	<b>89.7</b>	<b>0.502</b>	<b>67.8</b>

Higher scores on SST2 and BoolQ



Lower ExAccErr

# Transfer of knowledge

---

## DISTILLM: Towards Streamlined Distillation for Large Language Models

---

**Jongwoo Ko<sup>1</sup> Sungnyun Kim<sup>1</sup> Tianyi Chen<sup>2</sup> Se-Young Yun<sup>1</sup>**

<https://github.com/jongwooko/distillm>

# Transfer of knowledge

**Problem:** current KD methods for auto-regressive sequence models suffer from missing a standardized objective function

**Method:** propose a method that consists of **novel Skew KLD** and **Adaptive off-policy approach**

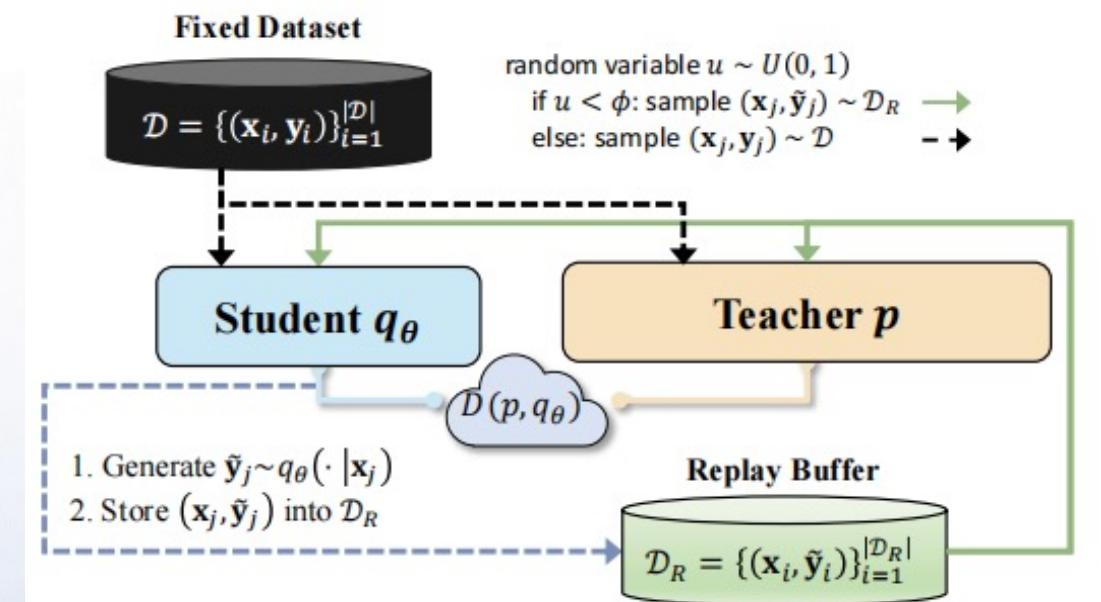
$$D_{\text{KL}}(p, q_{\theta}) = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{y} \sim p(\cdot | \mathbf{x})} \left[ \log \frac{p(\mathbf{y} | \mathbf{x})}{q_{\theta}(\mathbf{y} | \mathbf{x})} \right]$$

mix two distributions to replace student



$$D_{\text{SKL}}^{(\alpha)}(p, q_{\theta}) = D_{\text{KL}}(p, \alpha p + (1 - \alpha)q_{\theta})$$

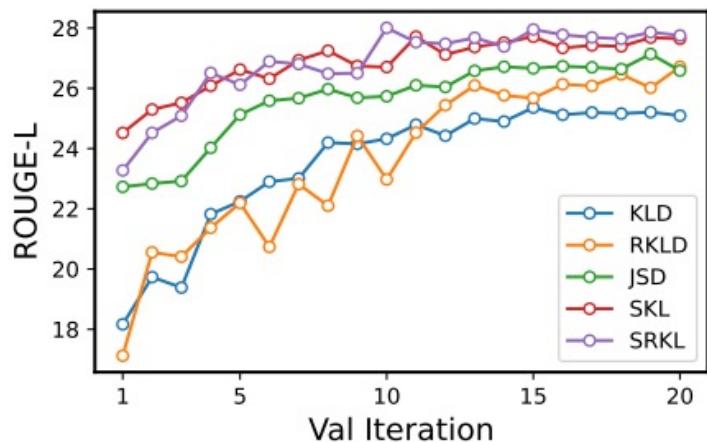
more **stable gradient**  
smaller **approximation error**



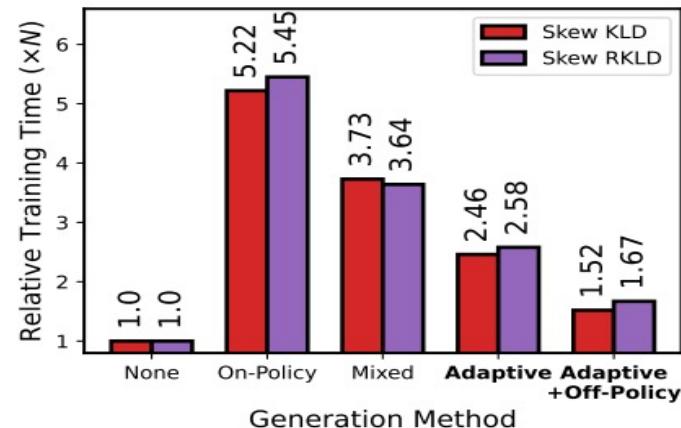
Adaptive use of student generated output  
randomly draw samples from Replay Buffer

# Transfer of knowledge

**Result:** achieves significant training efficiency and performance improvement.



Dataset	SAMSum		IWSLT 2017 En-De	
	T5-Base	T5-Small	T5-Base	T5-Small
KD (Hinton et al., 2015)	46.23	39.52	29.36	21.15
SeqKD (Kim & Rush, 2016)	46.89	40.24	29.07	21.42
ImitKD (Lin et al., 2020)	48.57	41.44	29.87	21.52
GKD (Agarwal et al., 2024)	48.49	41.92	30.24	22.04
<b>DISTILLM (ours)</b>	<b>49.11</b>	<b>42.37</b>	<b>30.32</b>	<b>22.53</b>



Dataset	Dolly Eval		Self-Instruct		Super-Natural	
	on-	off-	on-	off-	on-	off-
Sampling						
ImitKD (Lin et al., 2020)	21.63	20.62	10.85	10.09	19.94	18.04
GKD (Agarwal et al., 2024)	23.75	22.89	12.73	12.78	26.05	24.97
<b>DISTILLM (ours)</b>	26.37	26.12	13.14	13.16	28.24	28.20

# 目 录

- 1 Introduction
- 2 Knowledge distillation
- 3 Knowledge definition
- 4 Knowledge transfer
- 5 Application

# Application of KD

## TinyBERT: Distilling BERT for Natural Language Understanding

Xiaoqi Jiao<sup>1\*†</sup>, Yichun Yin<sup>2\*‡</sup>, Lifeng Shang<sup>2‡</sup>, Xin Jiang<sup>2</sup>

Xiao Chen<sup>2</sup>, Linlin Li<sup>3</sup>, Fang Wang<sup>1‡</sup> and Qun Liu<sup>2</sup>

<sup>1</sup>Key Laboratory of Information Storage System, Huazhong University of Science and Technology, Wuhan National Laboratory for Optoelectronics

<sup>2</sup>Huawei Noah's Ark Lab

<sup>3</sup>Huawei Technologies Co., Ltd.

{jiaoxiaoqi, wangfang}@hust.edu.cn

{yinyichun, shang.lifeng, jiang.xin}@huawei.com

{chen.xiao2, lynn.lilinlin, qun.liu}@huawei.com

# Application of KD

**Problem:** pre-trained language models (PLMs) like BERT are **computationally expensive** and difficult to deploy on edge devices due to the **large size**

**Method:** propose a new Transformer distillation method and a two-stage learning framework called TinyBERT

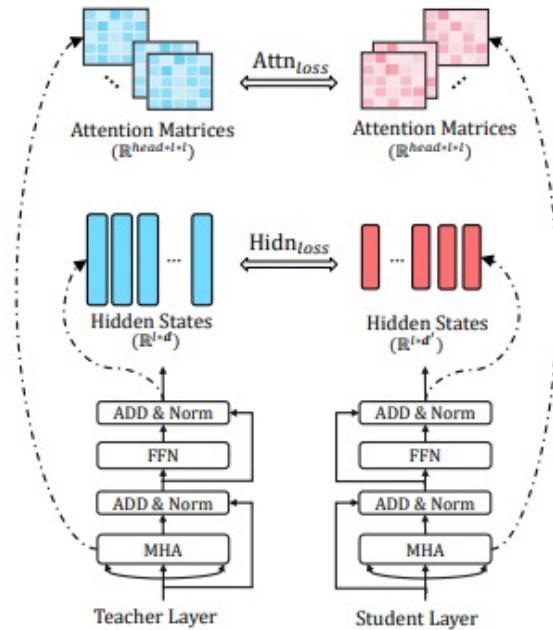


Figure 2: The details of Transformer-layer distillation consisting of  $\text{Attn}_{loss}$  (attention based distillation) and  $\text{Hidn}_{loss}$  (hidden states based distillation).

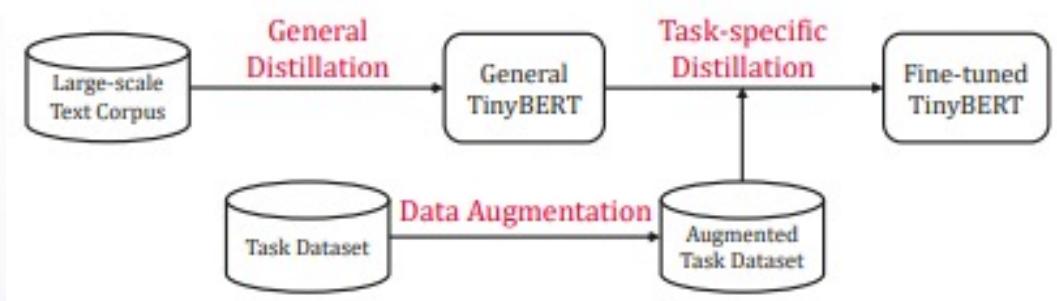


Figure 1: The illustration of TinyBERT learning.

distill knowledge from a large teacher BERT into a small student model through both **general** and **task-specific** distillations

# Application of KD

**Result:** TinyBERT4 achieves more than **96.8%** of BERT's performance while being **7.5x smaller and 9.4x faster**, and TinyBERT6 performs **on par with BERT**

System	#Params	#FLOPs	Speedup	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg
BERT <sub>BASE</sub> (Teacher)	109M	22.5B	1.0x	83.9/83.4	71.1	90.9	93.4	52.8	85.2	87.5	67.0	79.5
BERT <sub>TINY</sub>	14.5M	1.2B	9.4x	75.4/74.9	66.5	84.8	87.6	19.5	77.1	83.2	62.6	70.2
BERT <sub>SMALL</sub>	29.2M	3.4B	5.7x	77.6/77.0	68.1	86.4	89.7	27.8	77.0	83.4	61.8	72.1
BERT <sub>4</sub> -PKD	52.2M	7.6B	3.0x	79.9/79.3	70.2	85.1	89.4	24.8	79.8	82.6	62.3	72.6
DistilBERT <sub>4</sub>	52.2M	7.6B	3.0x	78.9/78.0	68.5	85.2	91.4	32.8	76.1	82.4	54.1	71.9
MobileBERT <sub>TINY</sub> †	15.1M	3.1B	-	81.5/81.6	68.9	<b>89.5</b>	91.7	<b>46.7</b>	80.1	<b>87.9</b>	65.1	<b>77.0</b>
TinyBERT <sub>4</sub> (ours)	14.5M	1.2B	9.4x	<b>82.5/81.8</b>	<b>71.3</b>	87.7	<b>92.6</b>	44.1	<b>80.4</b>	86.4	<b>66.6</b>	<b>77.0</b>
BERT <sub>6</sub> -PKD	67.0M	11.3B	2.0x	81.5/81.0	70.7	89.0	92.0	-	-	85.0	65.5	-
PD	67.0M	11.3B	2.0x	82.8/82.2	70.4	88.9	91.8	-	-	86.8	65.3	-
DistilBERT <sub>6</sub>	67.0M	11.3B	2.0x	82.6/81.3	70.1	88.9	92.5	49.0	81.3	86.9	58.4	76.8
TinyBERT <sub>6</sub> (ours)	67.0M	11.3B	2.0x	<b>84.6/83.2</b>	<b>71.6</b>	<b>90.4</b>	<b>93.1</b>	<b>51.1</b>	<b>83.7</b>	<b>87.3</b>	<b>70.0</b>	<b>79.4</b>

## Evaluation on the GLUE benchmark

# Application of KD

## Model Compression with Two-stage Multi-teacher Knowledge Distillation for Web Question Answering System\*

Ze Yang<sup>†</sup>

yaze@microsoft.com

STCA NLP Group, Microsoft

Linjun Shou<sup>†</sup>

lisho@microsoft.com

STCA NLP Group, Microsoft

Ming Gong

migon@microsoft.com

STCA NLP Group, Microsoft

Wutao Lin

wutlin@microsoft.com

STCA NLP Group, Microsoft

Dixin Jiang

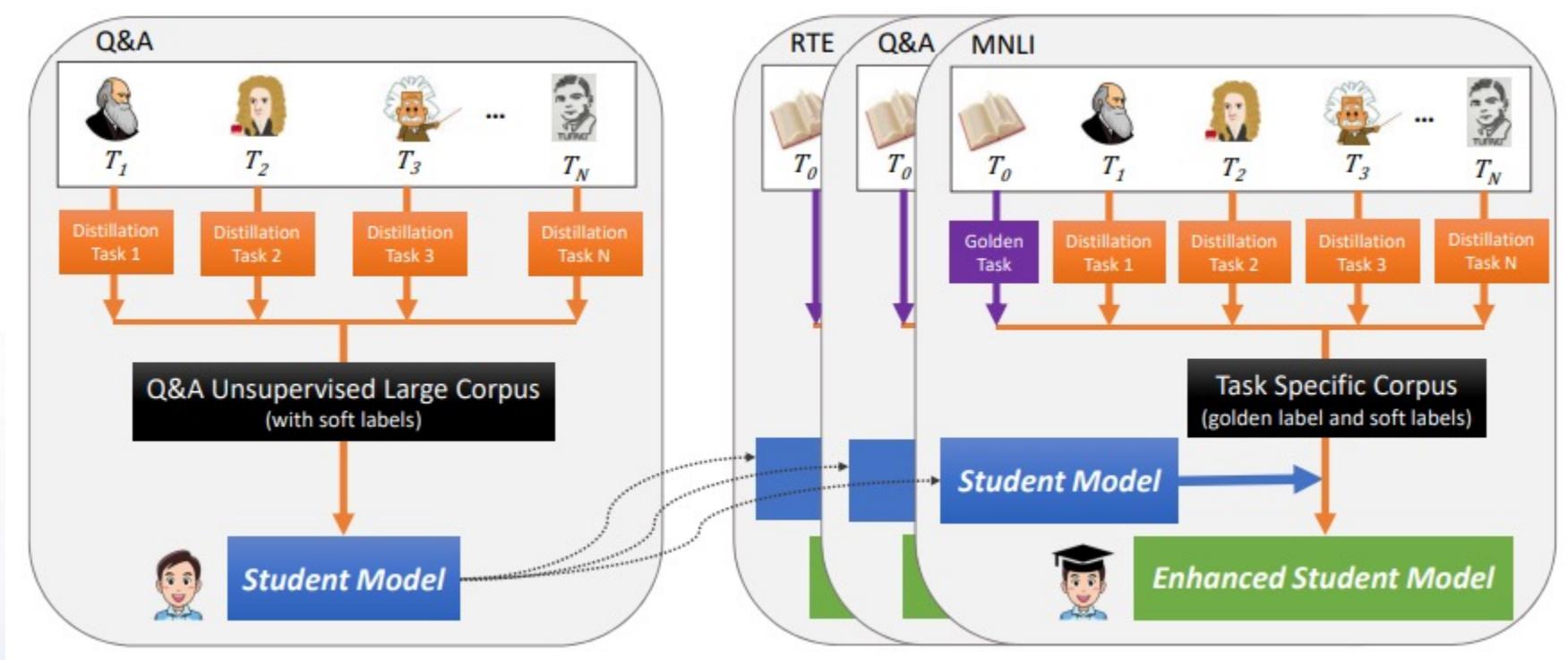
djiang@microsoft.com

STCA NLP Group, Microsoft

# Application of KD

**Problem:** previous model compression methods for PLMs usually suffer from information loss, leading to inferior models compared with the original one

**Method:** propose a **Two-stage Multi-teacher Knowledge Distillation (TMKD)** method



**general distillation task for student model pre-training in the first stage**

**multi-teacher distillation for fine-tuning on downstream tasks in the second stage**

# Application of KD

**Result:** significantly outperforms baseline methods and even achieves comparable performance with the original teacher models, with substantial improvement in model inference speed

Model	DeepQA	Performance (ACC)				Inference Speed(QPS)	Parameters (M)
		MNLI	SNLI	QNLI	RTE		
<b>Original Model</b>	<b>BERT-3</b>	75.78	70.77	77.75	78.51	57.42	207 50.44
	<b>BERT<sub>large</sub></b>	81.47	79.10	80.90	90.30	68.23	16 333.58
	<b>BERT<sub>large</sub> ensemble</b>	81.66	79.57	81.39	90.91	70.75	16/3 333.58*3
<b>Traditional Distillation Model</b>	<b>Bi-LSTM (1-o-1)</b>	71.69	59.39	69.59	69.12	56.31	207 50.44
	<b>Bi-LSTM (1<sub>avg</sub>-o-1)</b>	71.93	59.60	70.04	69.53	57.35	207 50.44
	<b>Bi-LSTM (m-o-m)</b>	72.04	61.71	72.89	69.89	58.12	207/3 50.44*3
	<b>BERT-3 (1-o-1)</b>	77.35	71.07	78.62	77.65	55.23	217 45.69
	<b>BERT-3 (1<sub>avg</sub>-o-1)</b>	77.63	70.63	78.64	78.20	58.12	217 45.69
	<b>BERT-3 (m-o-m)</b>	77.44	71.28	78.71	77.90	57.40	217/3 45.69*3
<b>Our Distillation Model</b>	<b>Bi-LSTM (TMKD<sub>base</sub>)</b>	74.73	61.68	71.71	69.99	62.74	207 50.45
	*TMKD <sub>base</sub>	79.93	71.29	78.35	83.53	66.64	217 45.70
	*TMKD <sub>large</sub>	<b>80.43</b>	<b>73.93</b>	<b>79.48</b>	<b>86.44</b>	<b>67.50</b>	<b>217</b> <b>45.70</b>

\* These two models are BERT-3 based models.

# Application of KD

## Adaptive Contrastive Knowledge Distillation for BERT Compression

**Jinyang Guo<sup>1,2\*</sup>, Jiaheng Liu<sup>2\*</sup>, Zining Wang<sup>2</sup>, Yuqing Ma<sup>2</sup>,  
Ruihao Gong<sup>2,3</sup>, Ke Xu<sup>2</sup> and Xianglong Liu<sup>2†</sup>**

<sup>1</sup>Institute of Artificial Intelligence, Beihang University

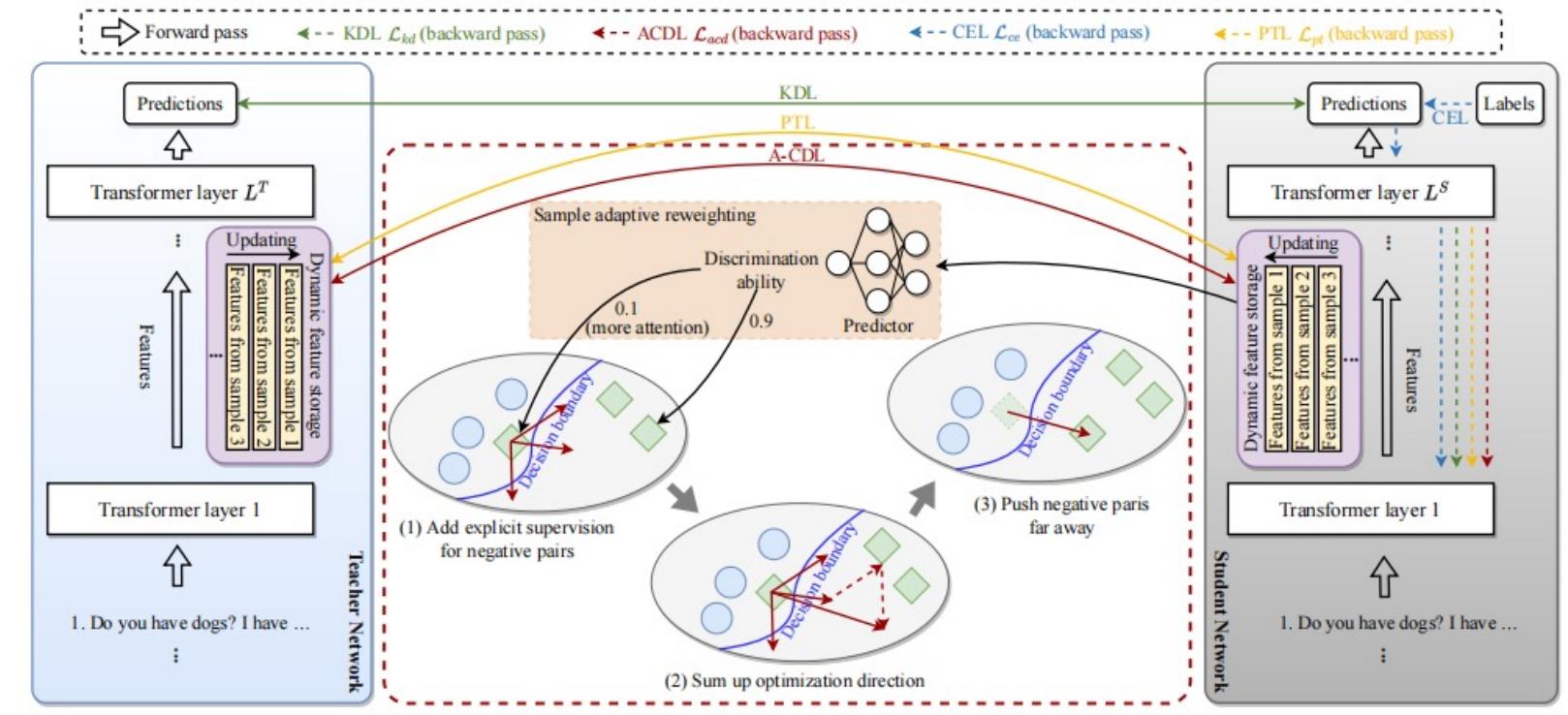
<sup>2</sup>State Key Lab of Software Development Environment, Beihang University

<sup>3</sup>SenseTime Group Limited

{jinyangguo,liujiaheng}@buaa.edu.cn, xlliu@buaa.edu.cn

# Application of KD

**Problem:** it is hard to distinguish similar sentences with completely different meanings  
**Method:** propose a novel **contrastive distillation** loss to introduce **explicit supervision** for learning discriminative student features

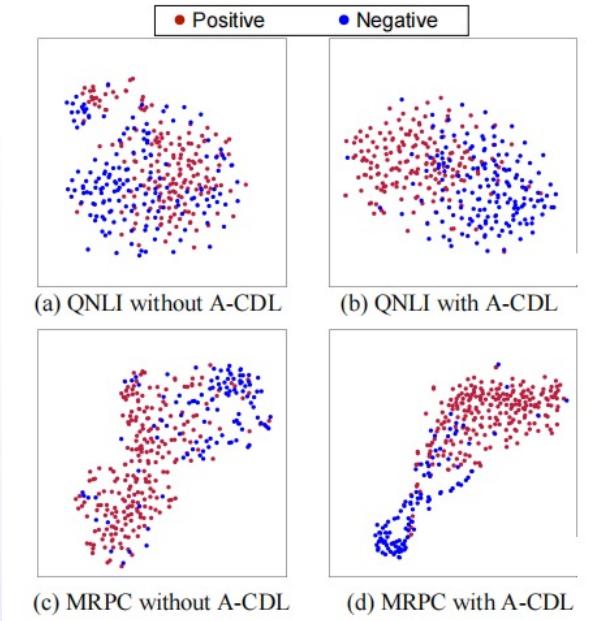


**adaptively** pay more attention to hard samples with fewer discrimination abilities

# Application of KD

## Result: performance improvements of language models on multiple datasets

Method	#Param	Speed-up	GLUE							
			CoLA (Matt.)	MNLI (Acc -m/-mm)	MRPC (F1/Acc)	QNLI (Acc)	QQP (F1/Acc)	RTE (Acc.)	SST-2 (Acc.)	STS-B (Pear./Spear.)
<b>Teacher Network: BERT-Base</b>										
BERT-Base (Devlin et al., 2018)	110M	1.0×	60.8	84.6/84.4	91.6/87.6	91.6	88.5/91.4	71.4	93.0	90.2/89.8
<b>Student Network: BERT<sub>3</sub></b>										
PKD (Sun et al., 2019)	46M	4.0×	<u>39.8</u>	75.9/76.6	84.1/75.0	84.3	<u>85.3/89.2</u>	62.8	87.4	86.3/86.1
RCO (Jin et al., 2019)	46M	4.0×	31.4	76.3/76.9	<u>85.3/77.5</u>	83.4	<u>85.4/88.7</u>	<u>66.1</u>	86.8	84.8/84.4
TAKD (Mirzadeh et al., 2020)	46M	4.0×	35.7	76.2/76.8	83.2/73.5	83.8	83.7/87.5	59.2	87.9	83.8/83.4
DistilBERT (Sanh et al., 2019)	46M	4.0×	34.0	<u>77.0/77.0</u>	83.2/73.0	83.8	85.1/88.9	62.8	86.9	<u>86.6/86.2</u>
TinyBERT (Jiao et al., 2019)	46M	4.0×	38.7	76.5/76.9	82.8/72.8	84.2	85.1/88.8	60.6	86.8	86.4/86.1
CRD (Tian et al., 2019)	46M	4.0×	38.6	76.1/76.8	<u>85.2/77.5</u>	<u>84.6</u>	83.9/88.0	65.7	87.6	86.1/85.6
SFTN (Park et al., 2021)	46M	4.0×	38.1	<u>76.6/77.1</u>	83.1/73.3	84.2	83.9/87.7	60.3	88.0	83.9/83.5
MetaDistill (Zhou et al., 2022)	46M	4.0×	39.3	75.9/76.4	82.0/71.1	83.8	83.7/88.1	62.1	88.0	<u>86.6/86.4</u>
Annealing KD (Jafari et al., 2021)	52M	3.0×	36.0	73.9/74.8	<u>86.2/-</u>	83.1	-/86.5	61.0	<b>89.4</b>	74.5/-
<b>ACKD (ours)</b>	46M	4.0×	<b>42.7</b>	<b>79.5/80.6</b>	<b>87.5/81.4</b>	<b>86.2</b>	<b>86.1/89.7</b>	<b>67.9</b>	<u>88.5</u>	<b>87.1/86.8</b>
<b>Student Network: BERT<sub>6</sub></b>										
PKD (Sun et al., 2019)	66M	2.0×	54.5	82.7/83.3	89.4/84.7	89.5	87.8/90.9	67.6	91.3	88.6/88.1
RCO (Jin et al., 2019)	66M	2.0×	53.6	82.4/82.9	89.5/85.1	89.7	87.4/90.6	67.6	91.4	88.7/88.3
TAKD (Mirzadeh et al., 2020)	66M	2.0×	53.8	82.5/83.0	89.6/85.0	89.6	87.5/90.7	68.5	91.4	88.2/88.0
DistilBERT (Sanh et al., 2019)	66M	2.0×	53.0	82.5/83.1	89.3/85.0	89.2	87.2/90.6	66.1	91.5	88.7/88.5
TinyBERT (Jiao et al., 2019)	66M	2.0×	52.4	<u>83.6/83.8</u>	90.5/86.5	89.8	87.6/90.6	67.7	91.9	89.2/ <u>88.7</u>
CRD (Tian et al., 2019)	66M	2.0×	55.8	83.2/83.4	89.5/85.5	89.8	87.6/90.8	67.1	91.5	88.8/88.3
SFTN (Park et al., 2021)	66M	2.0×	53.6	82.4/82.9	89.8/85.3	89.5	87.5/90.4	68.5	91.5	88.4/88.5
MetaDistill (Zhou et al., 2022)	66M	2.0×	<u>58.6</u>	<u>83.5/83.8</u>	<b>91.1/86.8</b>	90.4	<u>88.1/91.0</u>	<u>69.4</u>	<u>92.3</u>	<u>89.4/89.1</u>
ALP-KD (Passban et al., 2021)	66M	2.0×	46.4	82.0/-	-/85.8	89.7	-/90.6	69.0	91.9	88.8/-
CoDIR (Sun et al., 2020)	66M	2.0×	56.4	<u>83.9/-</u>	87.9/-	<b>90.7</b>	-/91.2	66.3	<b>92.4</b>	-/-
<b>ACKD (ours)</b>	66M	2.0×	<b>59.7</b>	<u>83.6/83.9</u>	<u>91.0/87.0</u>	<u>90.6</u>	<b>88.5/91.3</b>	<b>69.7</b>	<u>92.3</u>	<b>89.5/89.1</b>



**significant improvement in feature discrimination**

# Application of KD

## BiBERT: ACCURATE FULLY BINARIZED BERT

**Haotong Qin<sup>\*1,4</sup>, Yifu Ding<sup>\*1,4</sup>, Mingyuan Zhang<sup>\*2</sup>, Qinghua Yan<sup>1</sup>, Aishan Liu<sup>1</sup>,  
Qingqing Dang<sup>3</sup>, Ziwei Liu<sup>2</sup>, Xianglong Liu<sup>✉ 1</sup>**

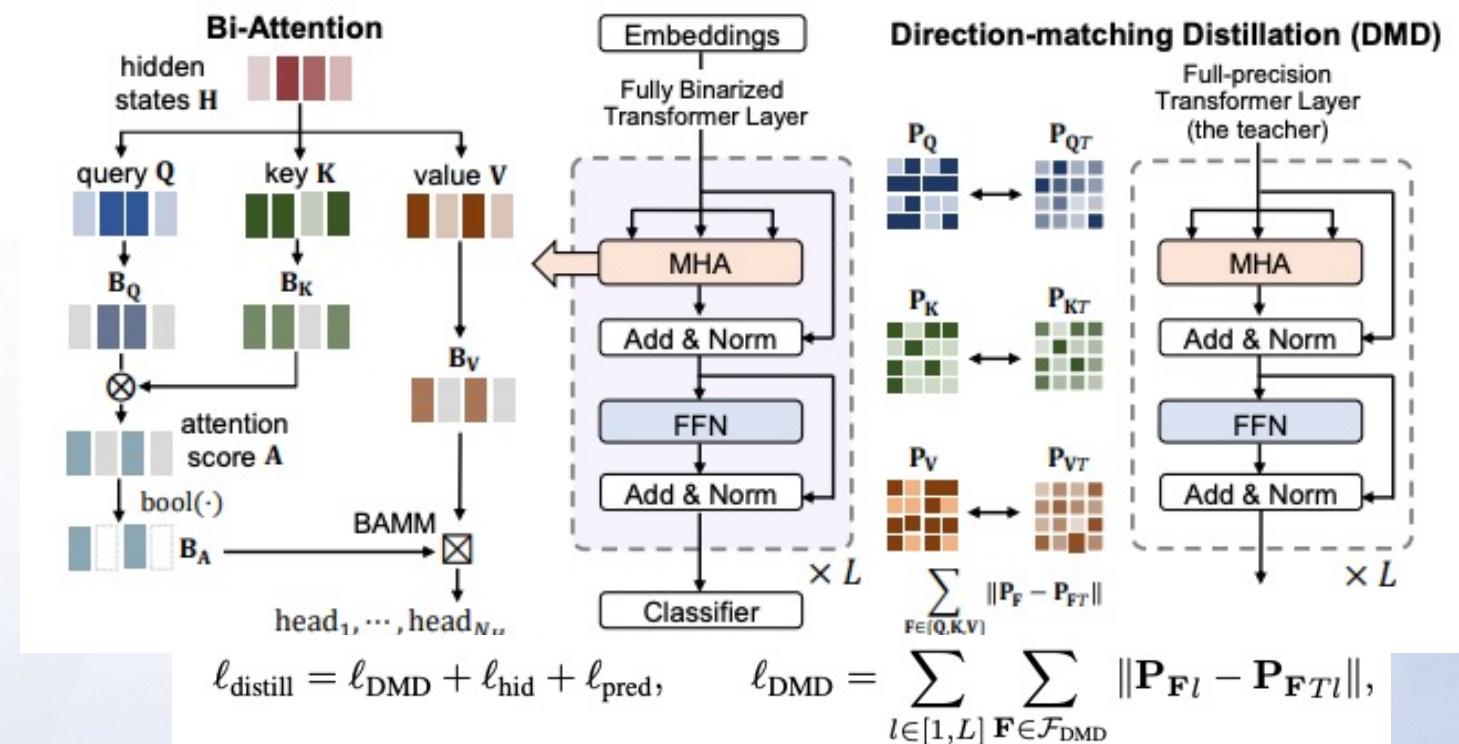
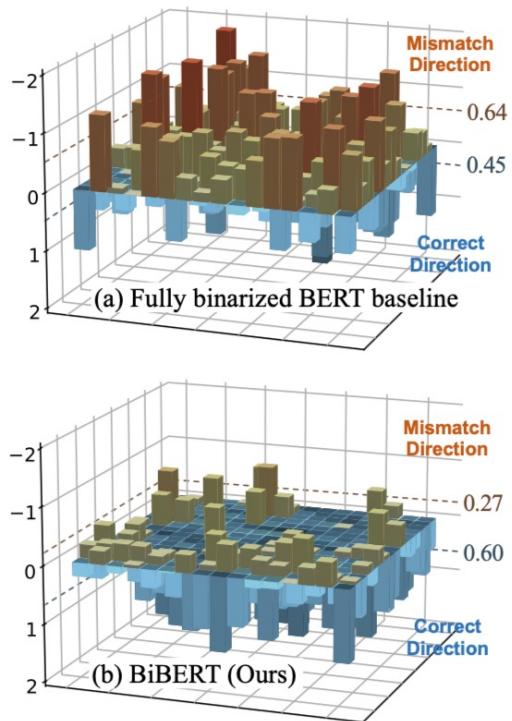
<sup>1</sup>State Key Lab of Software Development Environment, Beihang University   <sup>3</sup>Baidu Inc.

<sup>2</sup>S-Lab, Nanyang Technological University   <sup>4</sup>Shen Yuan Honors College, Beihang University  
{qinhaotong, yifuding, yanqh, aishanliu, xlliu}@buaa.edu.cn  
mingyuan001@e.ntu.edu.sg zwliu.hust@gmail.com dangqingqing@baidu.com

# Application of KD

**Problem:** transfer knowledge from full-precision teacher network to binarized student network

**Method: Direction-Matching Distillation**



optimization direction mismatch

develop the **Direction-Matching Distillation (DMD)** scheme  
to eliminate the direction mismatch in distillation

# Application of KD

## Result

Table 2: Comparison of BERT quantization methods without data augmentation

Quant	#Bits (W-E-A)	Size (MB)	FLOPs (G)	MNLI -m/mm	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg.
Full Prec.	32-32-32	418	22.5	84.9/85.5	91.4	92.1	93.2	59.7	90.1	86.3	72.2	83.9
Baseline	1-1-1	13.4	0.4	45.8/47.0	73.2	66.4	77.6	11.7	7.6	70.2	54.1	50.4
BinaryBERT	1-1-1	16.5	0.4	35.6/35.3	63.1	51.5	53.2	0	6.1	68.3	52.7	40.6
BinaryBERT	1-1-2	16.5	0.8	35.4/35.2	63.1	52.6	82.5	14.6	6.5	68.3	52.7	45.7
TernaryBERT	2-2-1	28.0	0.8	32.7/33.0	74.1	59.3	53.1	0	7.1	68.3	53.4	42.3
TernaryBERT	2-2-2	28.0	1.5	40.3/40.0	63.1	50.0	80.7	0	12.4	68.3	54.5	45.5
Q2BERT	2-8-8	43.0	6.5	47.2/47.3	67.0	61.3	80.6	0	4.4	68.4	52.7	47.7
Q-BERT	2-8-8	43.0	6.5	<b>76.6/77.0</b>	—	—	84.6	—	—	68.3	52.7	—
BiBERT (ours)	<b>1-1-1</b>	<b>13.4</b>	<b>0.4</b>	59.3/60.0	<b>82.4</b>	<b>70.2</b>	<b>86.9</b>	<b>25.3</b>	<b>33.5</b>	<b>72.9</b>	<b>58.5</b>	<b>61.0</b>

Table 3: Comparison of BERT quantization methods with data augmentation

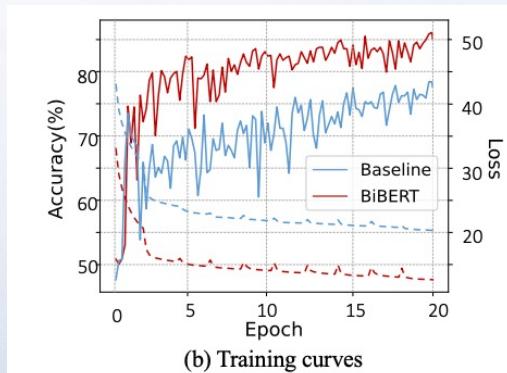
Quant	#Bits (W-E-A)	Size (MB)	FLOPs (G)	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg.
Full Prec.	32-32-32	418	22.5	92.1	93.2	59.7	90.1	86.3	72.2	82.3
Baseline	1-1-1	13.4	0.4	69.2	84.0	23.3	14.4	71.4	50.9	52.2
BinaryBERT	1-1-1	16.5	0.4	66.1	68.0	7.3	22.1	69.3	57.7	48.4
BinaryBERT	1-1-2	16.5	0.8	51.0	89.6	33.0	11.4	71.0	55.9	52.0
TernaryBERT	2-2-1	28.0	0.8	50.9	80.3	6.5	10.3	71.5	53.4	45.5
TernaryBERT	2-2-2	28.0	1.5	50.0	87.5	20.6	<b>72.5</b>	72.0	47.2	58.3
BiBERT (ours)	<b>1-1-1</b>	<b>13.4</b>	<b>0.4</b>	<b>76.0</b>	<b>90.9</b>	<b>37.8</b>	56.7	<b>78.8</b>	<b>61.0</b>	<b>67.0</b>

**BiBERT outperforms existing binarized quantization methods of BERT, even higher bits quantization methods.**

Table 1: Ablation study.

Quant	#Bits	DAS	SST-2	MRPC	RTE	QQP
Full Precision	32-32-32	—	93.2	86.3	72.2	91.4
Baseline	1-1-1	X	77.6	70.2	54.1	73.2
Bi-Attention	1-1-1	X	82.1	70.5	55.6	74.9
DMD	1-1-1	X	79.9	70.5	55.2	75.3
BiBERT (ours)	1-1-1	X	<b>88.7</b>	<b>72.5</b>	<b>57.4</b>	<b>84.8</b>
Baseline	1-1-1	✓	84.0	71.4	50.9	-
Bi-Attention	1-1-1	✓	85.6	73.2	53.1	-
DMD	1-1-1	✓	85.3	72.5	56.3	-
BiBERT (ours)	1-1-1	✓	<b>90.9</b>	<b>78.8</b>	<b>61.0</b>	-

Higher model performance



Faster convergence rate

# Application of KD

## HuatuoGPT, Towards Taming Language Models To Be a Doctor

**Hongbo Zhang<sup>1,2†</sup>, Junying Chen<sup>1,2†</sup>, Feng Jiang<sup>1,2,3†</sup>, Fei Yu<sup>1,2</sup>, Zhihong Chen<sup>1,2</sup>,  
Jianquan Li<sup>2</sup>, Guiming Chen<sup>1,2</sup>, Xiangbo Wu<sup>2</sup>, Zhiyi Zhang<sup>2</sup>, Qingying Xiao<sup>1</sup>,  
Xiang Wan<sup>1,2</sup>, Benyou Wang<sup>1,2 \*</sup>, Haizhou Li<sup>1,2</sup>**

<sup>1</sup>Shenzhen Research Institute of Big Data, <sup>2</sup>The Chinese University of Hong Kong, Shenzhen

<sup>3</sup>University of Science and Technology of China

[hongboz183@gmail.com](mailto:hongboz183@gmail.com), [junying.chen.cs@gmail.com](mailto:junying.chen.cs@gmail.com), [jeffreyjiang@cuhk.edu.cn](mailto:jeffreyjiang@cuhk.edu.cn)  
[wangbenyou@cuhk.edu.cn](mailto:wangbenyou@cuhk.edu.cn)

# Application of KD

**Problem:** ChatGPT did not reach the expert level and lack of proficiency in medical knowledge that medical doctors have

**Method:** propose HuatuoGPT, combining the strengths of both **distilled data** (from ChatGPT) and real-world data (from Doctors) to tame the medical LLM

## Distilled data from ChatGPT

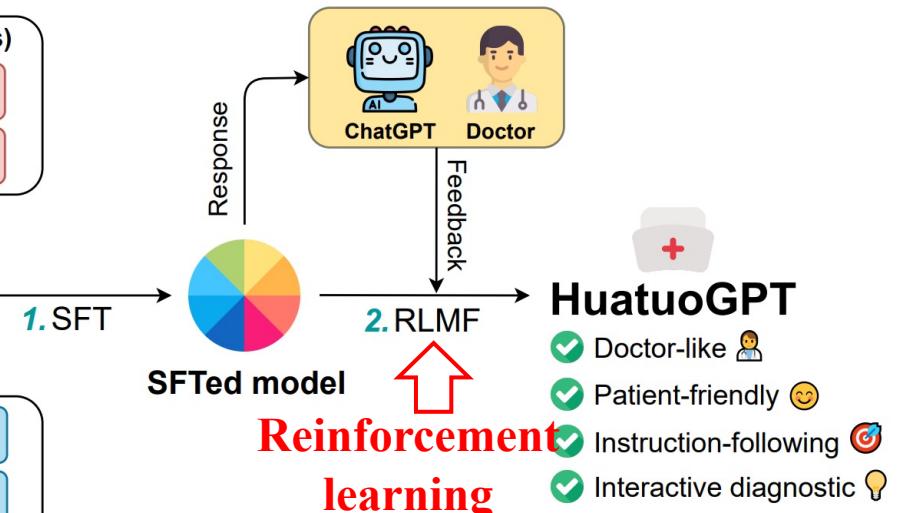
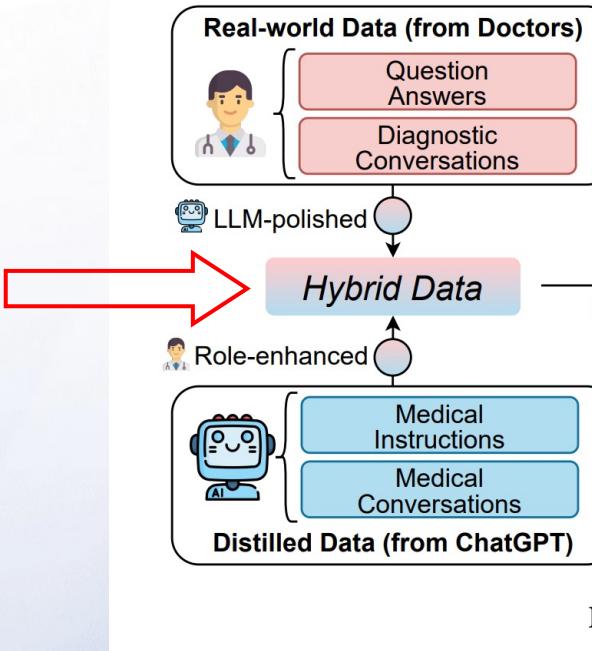
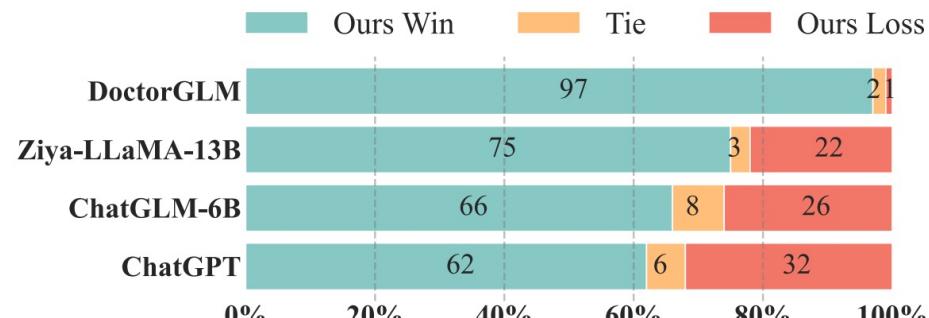


Figure 2: Schematic of HuatuoGPT.

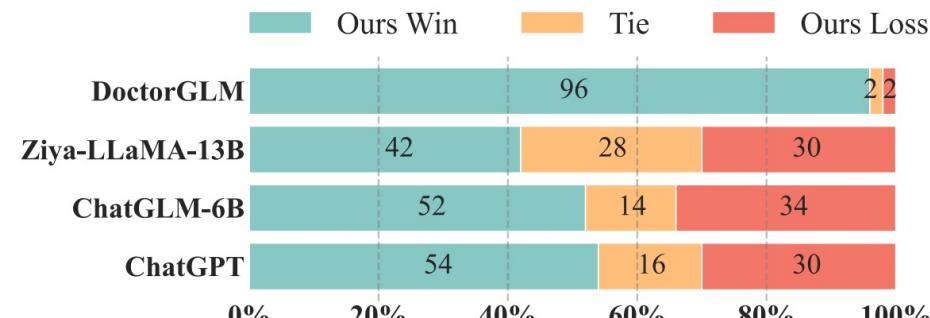
## Real world data from doctors

# Application of KD

**Result: achieves state-of-the-art results in performing medical consultation among open-source LLMs, outperforms its teacher model (i.e., ChatGPT) in most cases.**



(a) GPT-4 Evaluation



(b) Doctor Evaluation

**User:**

我最近拉肚子，有泡沫黄水，怀疑是喝牛奶喝的，是不是喝牛奶的原因。谢谢。

(I've been experiencing diarrhea recently, with my stool appearing as frothy yellow liquid. I suspect it might be due to the milk I've been drinking. Could this be the reason? Thank you.)

**HuatuoGPT (w/ real data):**

您好！根据你描述的症状，可能存在肠胃疾病的情况。此外，饮食习惯也会影响到这种情况发生。从中医角度来看，这主要与脾胃功能失调有关。因为脾胃相表里，所以需要通过健脾胃的方法进行治疗。同时，平常饮食也需注意清淡易消化。

(Hello! Based on the symptoms you described, there might be a potential issue with your gastrointestinal system. Moreover, dietary habits can also affect the occurrence of such symptoms. From a traditional Chinese medicine perspective, this is primarily related to a disorder of spleen and stomach functions, since they are closely linked. Therefore, treatment would involve methods to strengthen the spleen and stomach. Additionally, you should consider a light and easily digestible diet in your daily eating habits.)

# Application of KD

## VALLEY: VIDEO ASSISTANT WITH LARGE LANGUAGE MODEL ENHANCED ABILITY

Ruipu Luo<sup>1,2\*</sup>, Ziwang Zhao<sup>1,3\*</sup>, Min Yang<sup>1\*</sup>, Junwei Dong<sup>1,4</sup>, Da Li<sup>1</sup>,  
Pengcheng Lu<sup>1</sup>, Tao Wang<sup>1</sup>, Linmei Hu<sup>5</sup>, Minghui Qiu<sup>1,†</sup>, Zhongyu Wei<sup>2</sup>

<sup>1</sup>ByteDance Inc.    <sup>2</sup>Fudan University    <sup>4</sup>Chongqing University

<sup>3</sup>Beijing University of Posts and Telecommunications

<sup>5</sup>Beijing Institute of Technology

{luoruipu, zhaoziwang, yangmin.priv, dongjunwei}@bytedance.com

# Application of KD

**Problem:** jointing **video and language** understanding has not been widely explored

**Method:** adopt a pre-training-then-instructions-tuned pipeline to **align visual and textual modalities** and improve the instruction-following capability

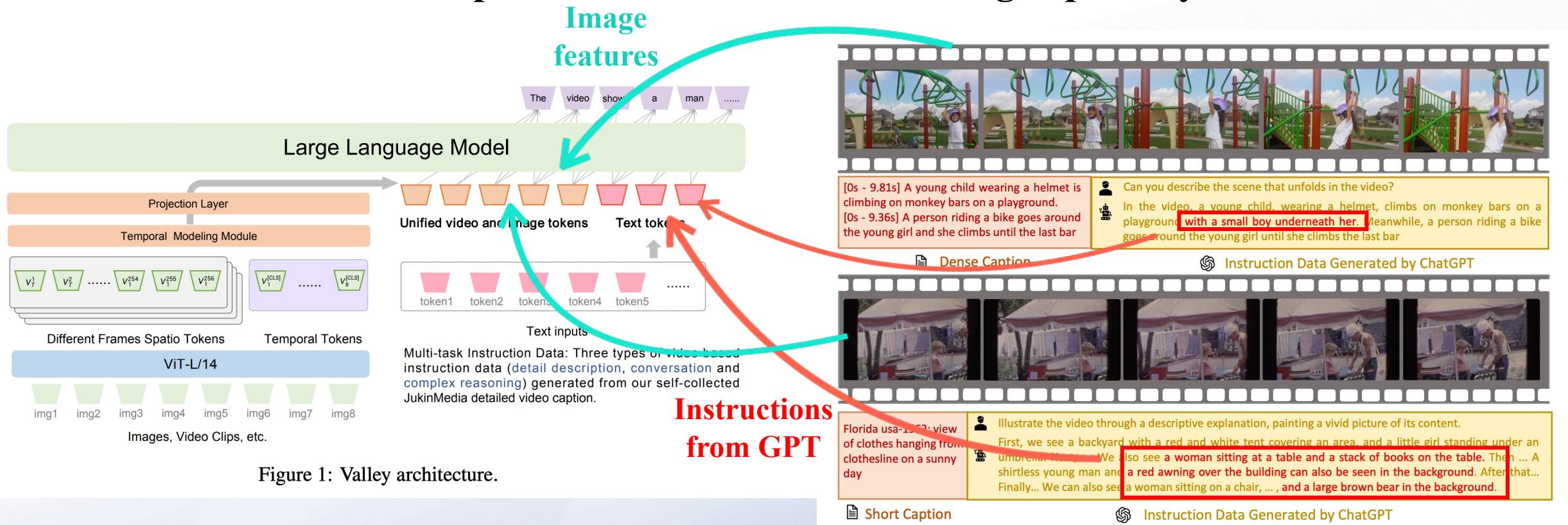


Figure 1: Valley architecture.

# Application of KD

**Result: a highly effective video assistant that can understand complex video scenarios and provide descriptions in text.**

Method	MSVD-QA		MSRVT-TT-QA		ActivityNet-QA	
	Accuracy	Score	Accuracy	Score	Accuracy	Score
FrozenBiLM	32.2	–	16.8	–	24.7	–
VideoChat	56.3	2.8	45.0	2.5	26.5	2.2
LLaMA Adapter	54.9	3.1	43.8	2.7	34.2	2.7
Video LLaMA	51.6	2.5	29.6	1.8	12.4	1.1
Video-ChatGPT	64.9	3.3	49.3	2.8	35.2	2.7
Valley-v1	<b>65.4</b>	<b>3.4</b>	45.7	2.5	42.9	3.0
Valley-v2	59.1	<b>3.4</b>	49.9	<b>3.0</b>	32.5	2.6
Valley-v3	60.5	3.3	<b>51.1</b>	2.9	<b>45.1</b>	<b>3.2</b>

Table 1: Zero-shot video QA result on MSVD-QA, MSRVT-TT-QA and ActivityNet-QA

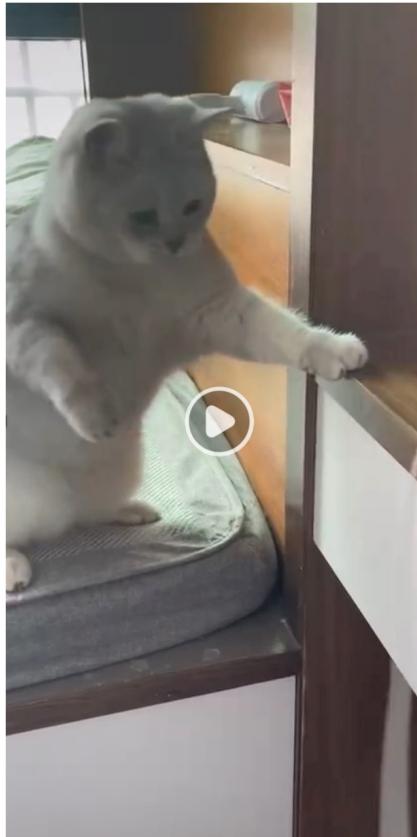
Method	COR	DO	CU	TU	CON
VideoChat	2.23	2.50	2.53	1.94	2.24
LLaMA Adapter	2.03	2.32	2.30	1.98	2.15
Video LLaMA	1.96	2.18	2.16	1.82	1.79
Video-ChatGPT	2.40	<b>2.52</b>	2.62	1.98	2.37
Valley-v1	2.06	2.42	2.74	1.83	2.41
Valley-v2	2.35	2.13	2.85	1.99	2.17
Valley-v3	<b>2.43</b>	2.13	<b>2.86</b>	<b>2.04</b>	<b>2.45</b>

Table 2: Results on video-based text generation benchmark provided by Video-ChatGPT.

Method	Setup	BERT-F1
Flamingo	zero-shot	65.51
	zero-shot COT	58.23
MiniGPT4	zero-shot	65.81
	zero-shot COT	68.45
Valley	zero-shot	84.82
	zero-shot COT	<b>85.26</b>

Table 3: Results on Meme-Cap, the image metaphor understanding benchmark.

**Leading performance in video Q&A**



A horizontal film strip consisting of four rectangular frames, each showing a different pose of a white cat interacting with a white cord hanging from a wooden railing. The frames are separated by black bars.

Figure 6: Video description, funny things identification, and recognition of interactions of the objects in the video.

# Thanks!

