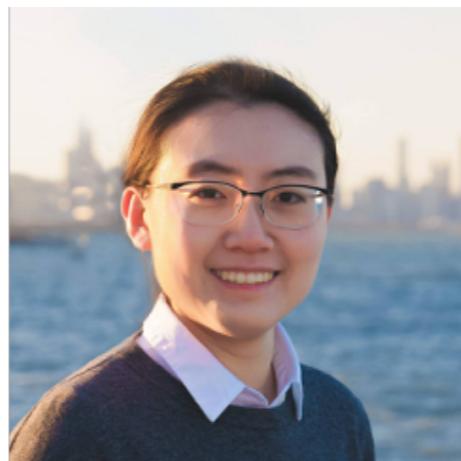


# Knowledge Distillation for Language Models: Challenges and Opportunities with Sequential Data

Yuqiao Wen<sup>1,3</sup>



Freda Shi<sup>2,4,5</sup>



Lili Mou<sup>1,3,5</sup>



[yq.when@gmail.com](mailto:yq.when@gmail.com)

[fhs@uwaterloo.ca](mailto:fhs@uwaterloo.ca)

[doublepower.mou@gmail.com](mailto:doublepower.mou@gmail.com)

<sup>1</sup>Dept. Computing Science, University of Alberta

<sup>2</sup>David R. Cheriton School of Computer Science, University of Waterloo

<sup>3</sup>Alberta Machine Intelligence Institute (Amii)

<sup>4</sup>Vector Institute

<sup>5</sup>Canada CIFAR AI Chair

AAAI'26 Tutorial

# Knowledge Distillation for Language Models: Challenges and Opportunities with Sequential Data

Yuqiao Wen<sup>1,3</sup>



Freda Shi<sup>2,4,5</sup>



Lili Mou<sup>1,3,5</sup>



[yq.when@gmail.com](mailto:yq.when@gmail.com)

[fhs@uwaterloo.ca](mailto:fhs@uwaterloo.ca)

[doublepower.mou@gmail.com](mailto:doublepower.mou@gmail.com)

<sup>1</sup>Dept. Computing Science, University of Alberta

<sup>2</sup>David R. Cheriton School of Computer Science, University of Waterloo

<sup>3</sup>Alberta Machine Intelligence Institute (Amii)

<sup>4</sup>Vector Institute

<sup>5</sup>Canada CIFAR AI Chair

AAAI'26 Tutorial



## Lili Mou is admitting

- All-level students  
MSc, PhD, exchanging
- Visiting scholars  
RA, Postdoc



Talk will resume in 10 seconds...  
[Skip ad >](#)



# **Lili Mou (LLM)'s Lab**

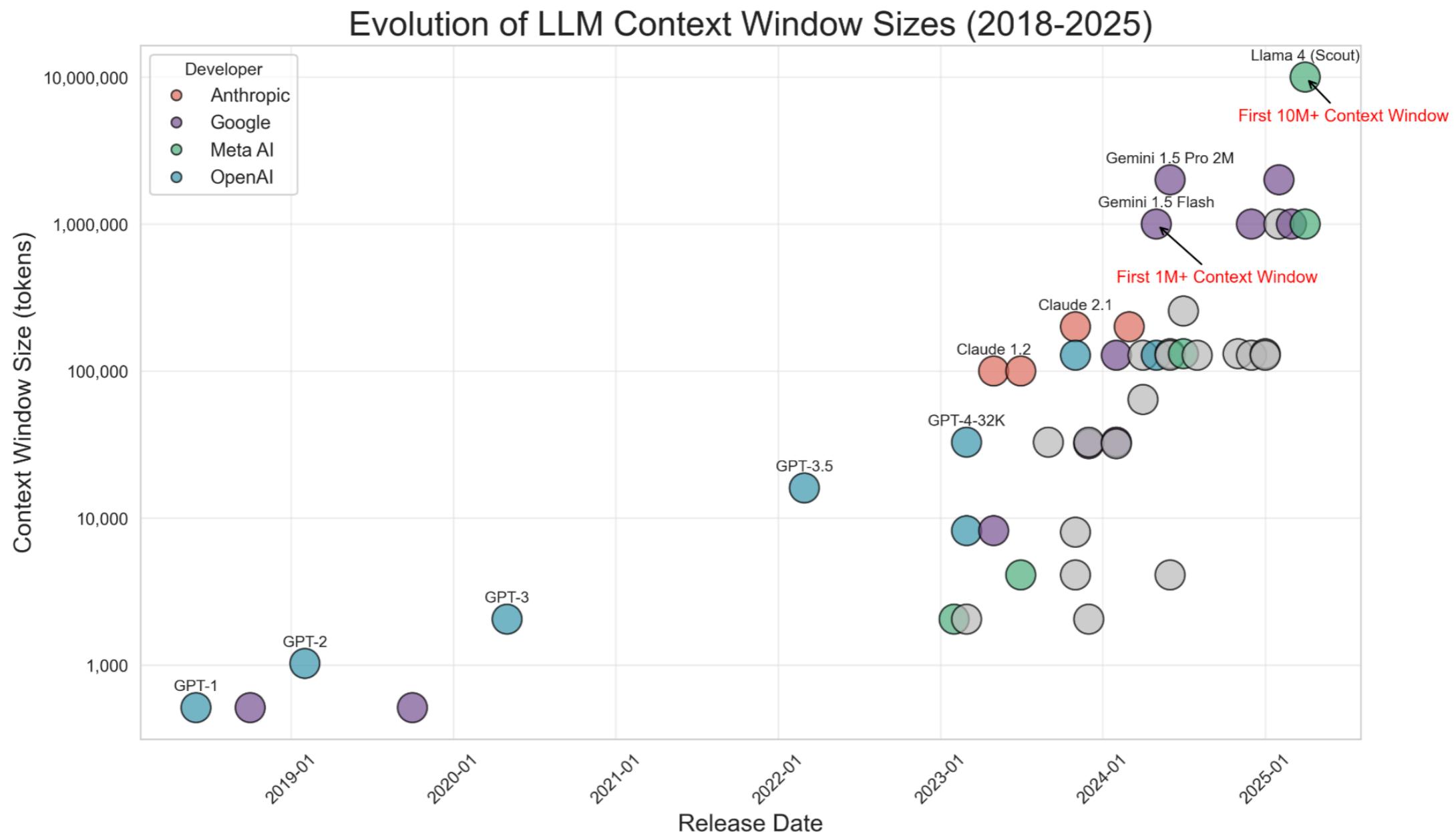
**Lili Mou (LLM)'s Lab strictly  
follows the KD methodologies  
discussed in this tutorial**

# Tutorial Outline

- **Session I [45min]:**
  - **Introduction** 
  - f-Divergence KD [ACL'23]
  - RL for LLM Distillation [NeurIPS'23, COLING'24]
- Break [15min]
- Session II [15min]
  - Multi-Teacher Distillation [ICLR'24, ACL'24, AAAI'25a,b]
  - KD Applications to LLMs
  - Conclusion and Future Work

# Why knowledge distillation?

- Model size grows exponentially



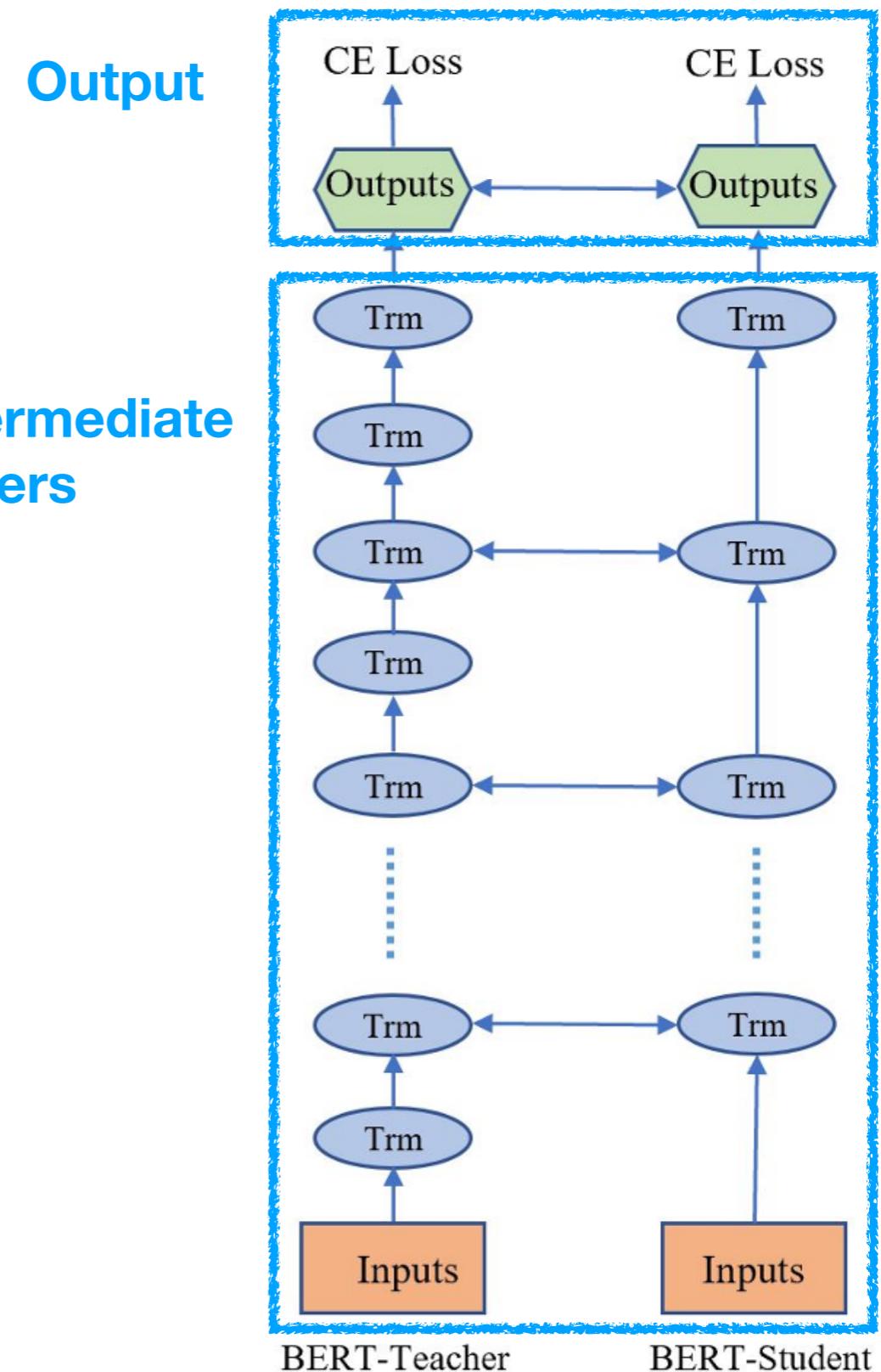
Source: <https://www.meibel.ai/post/understanding-the-impact-of-increasing-l1m-context-windows>

# Why knowledge distillation?

- Knowledge distillation (KD): Transferring knowledge
  - From: a large teacher model
  - To: a small student model
- KD improves efficiency
  - Parameter efficiency
    - ▶ Memory, storage, preventing overfitting
  - Time efficiency
    - ▶ Fine-tuning, inference

# Categorization of KD Methods

- Output matching
    - Mandatory
    - Informs student of the task
  - Intermediate-layer matching
    - Student learns step by step
    - Optional but often helpful



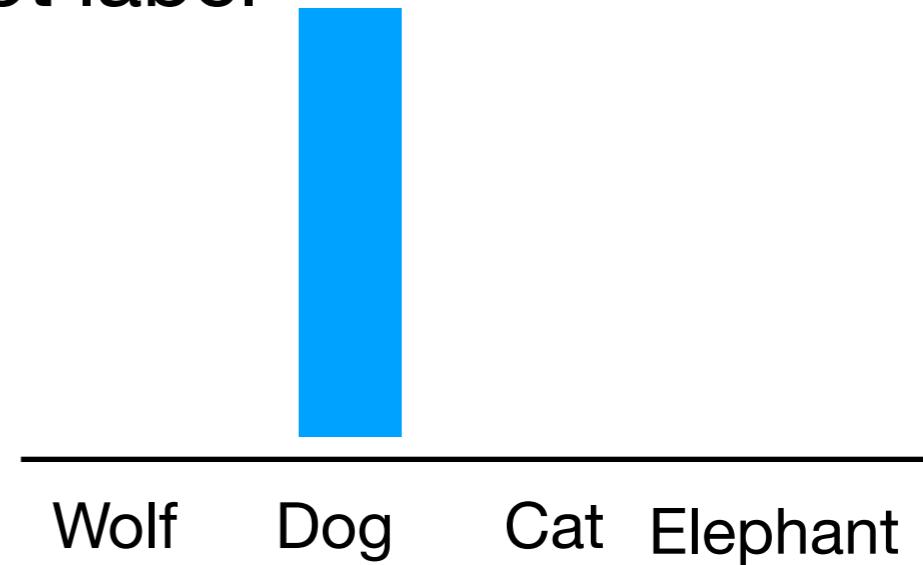
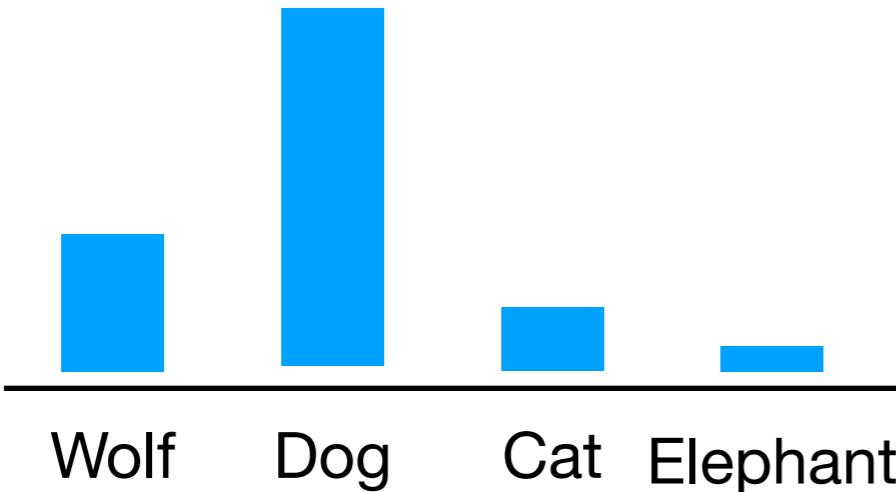
[Figure adapted from Sun et al. (2019)]

# Output Matching

- (Soft) cross-entropy loss

$$\sum_x -p(x)\log q_{\theta}(x)$$

- $p(x)$ : Teacher distribution (may adjust temperature)
- $q_{\theta}(x)$ : Student distribution
- Why soft cross-entropy loss?
  - Teacher provides a full distribution
  - More informative than a one-hot label



# Ranking-Based Output Matching

- Learning a student model  $Q$  to rank samples consistently with the teacher  $P$ .

$$\begin{aligned} & \arg \min_Q \mathbb{E}_{\langle x^+, x^- \rangle} [-\log Q(x^+) + \log Q(x^-)] \\ & s.t. \quad P(x^+) > P(x^-) \end{aligned}$$

- Example application in LM alignment (Tunstall et al., 2023):
  - Learning to rank the output sequences given input.

# Intermediate-Layer Matching

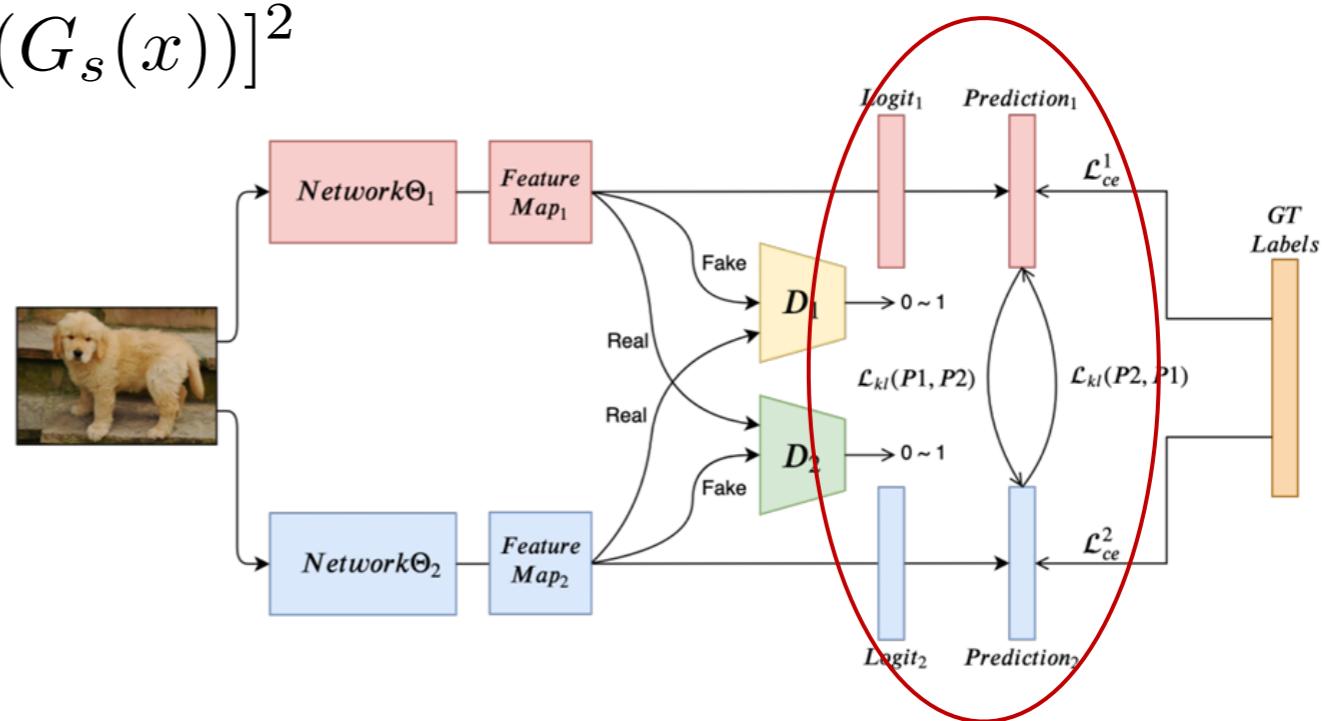
- Guides the student to learn step by step
- $L_2$  loss with linear mapping (Romero et al., 2015)

$$\mathcal{L}_{\text{matching}} = \sum_{i=1}^N \sum_{j=1}^L \|\mathbf{t}_{i,j} - \alpha(\mathbf{W}\mathbf{s}_{i,j})\|_2^2$$

- Adversarial loss (Chung et al., 2020)

$$\mathcal{L}_{\text{disc}} = [1 - D_s(G_t(x))]^2 + [D_s(G_s(x))]^2$$

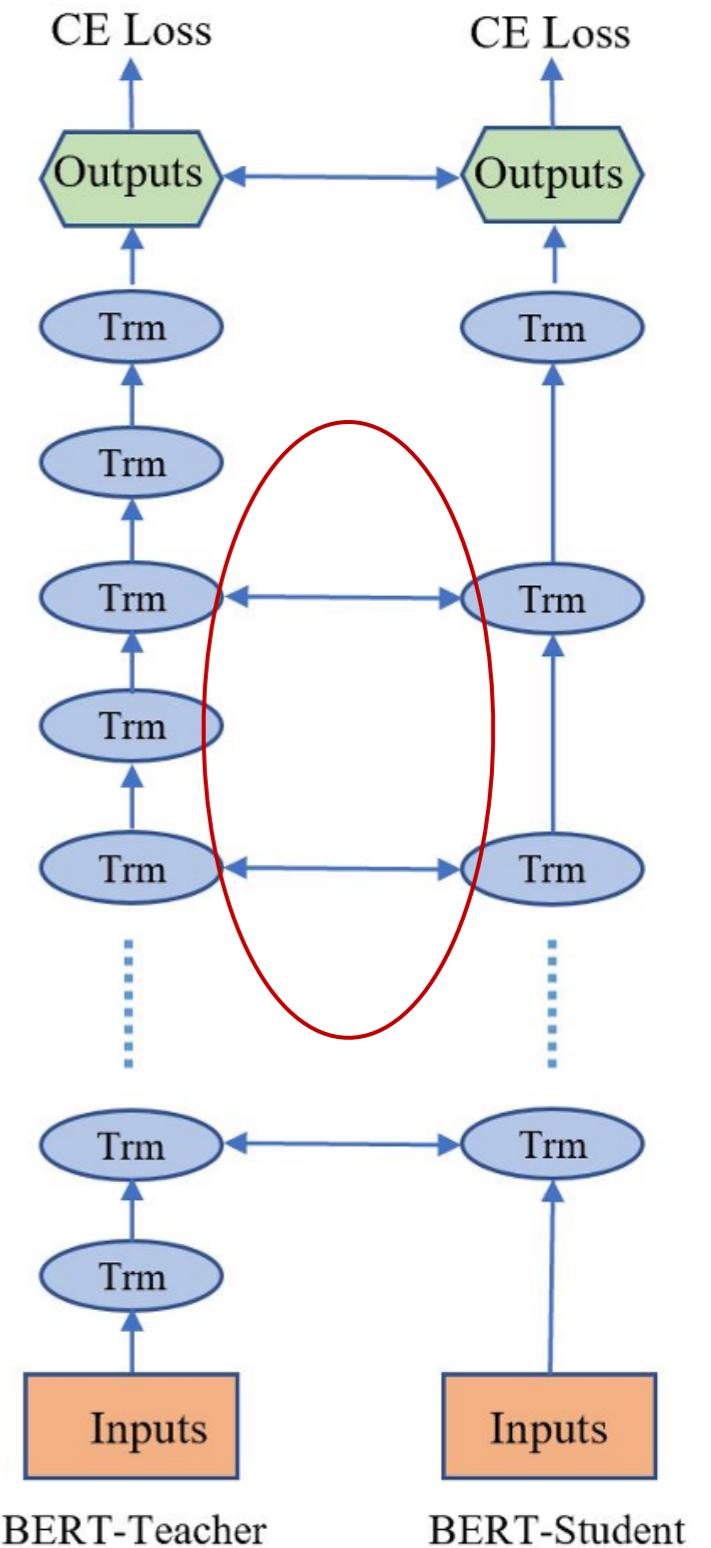
$$\mathcal{L}_{\text{gen}} = [1 - D_s(G_s(x))]^2$$



[Figure adapted from Chung et al. (2020)]

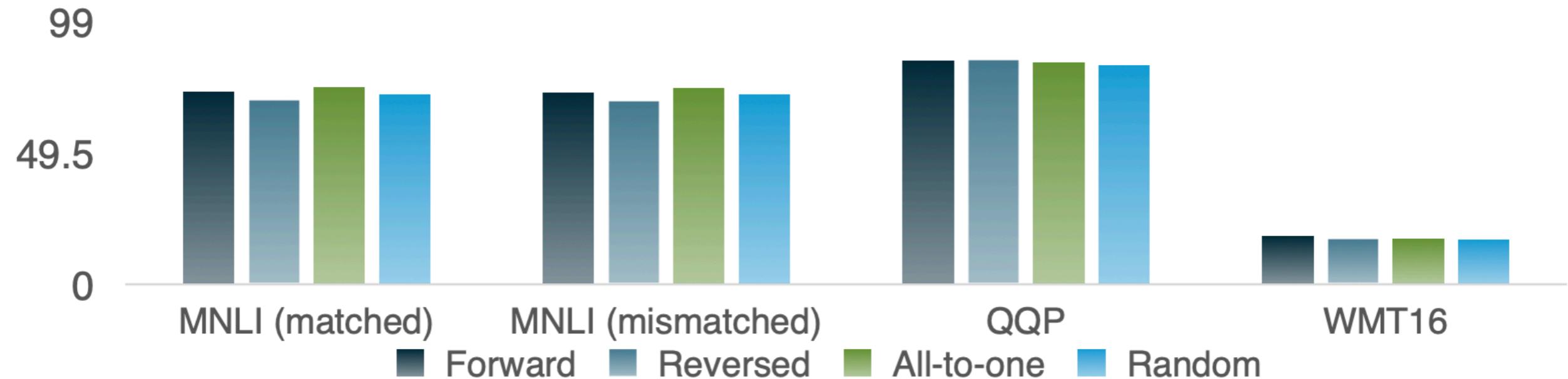
# Layer Selection

- Which layer is matched to which?
- Previous strategies
  - Evenly spaced (Sun et al., 2019)
  - All-to-one (Wang et al., 2020)
  - Weighted combination (Passban et al., 2021)
- Our recent findings (Yu et al., 2025)
  - Intermediate-layer matching is important
  - Layer selection is not

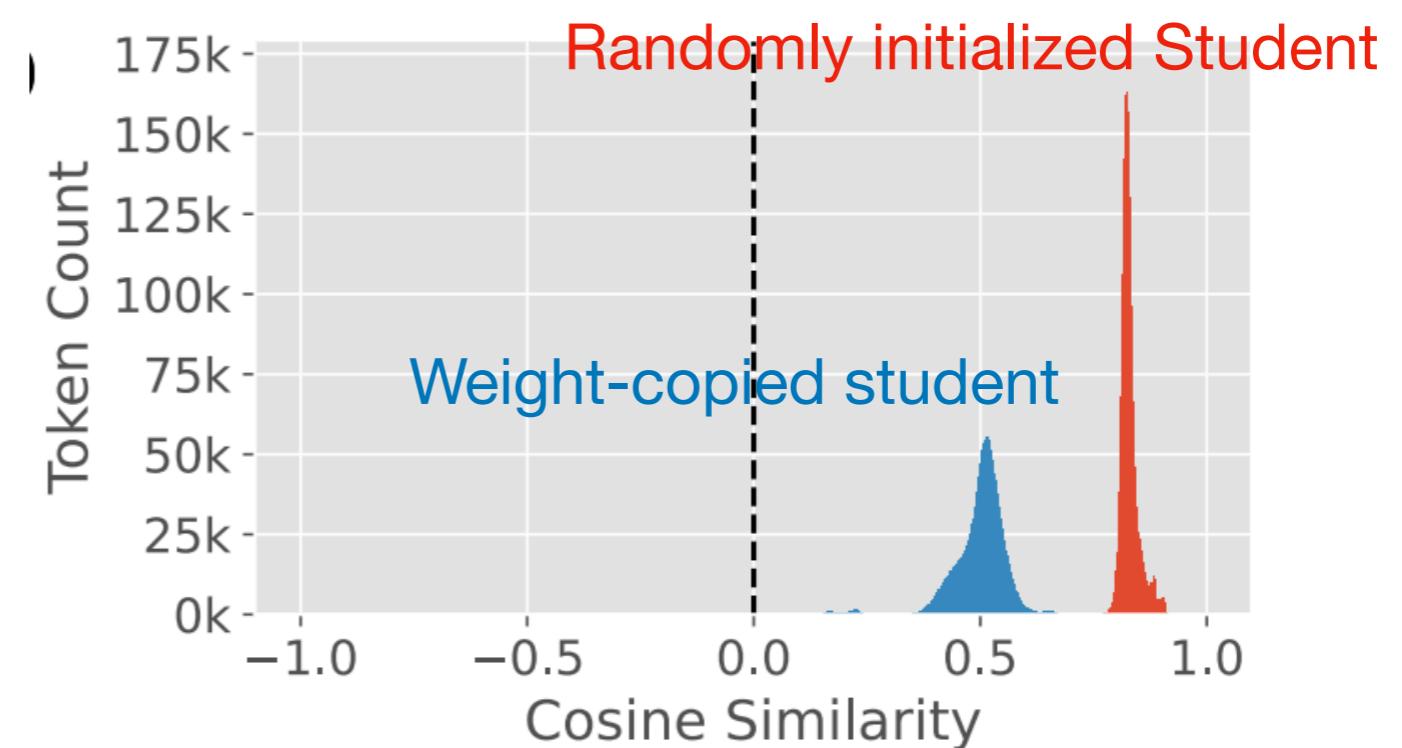
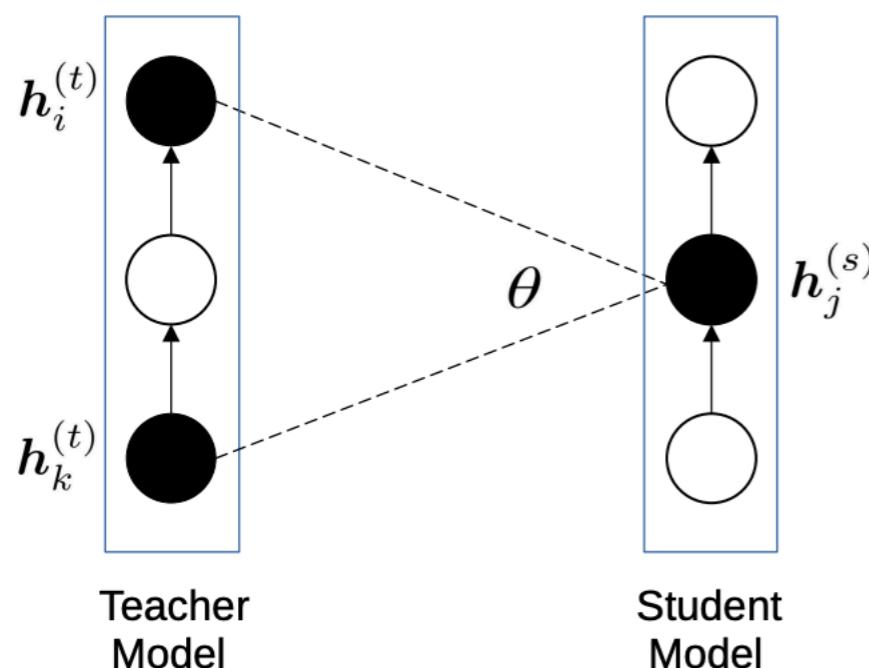


[Figure adapted from Sun et al. (2019)]

# Layer Selection Performance

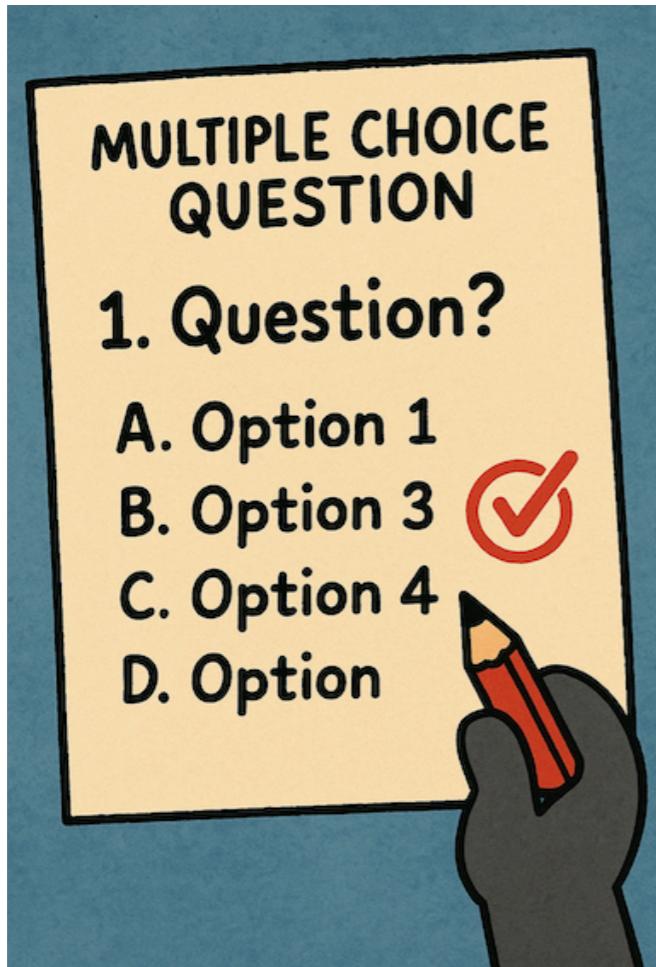


- A plausible explanation:
  - Between-layer distance << between-model distance

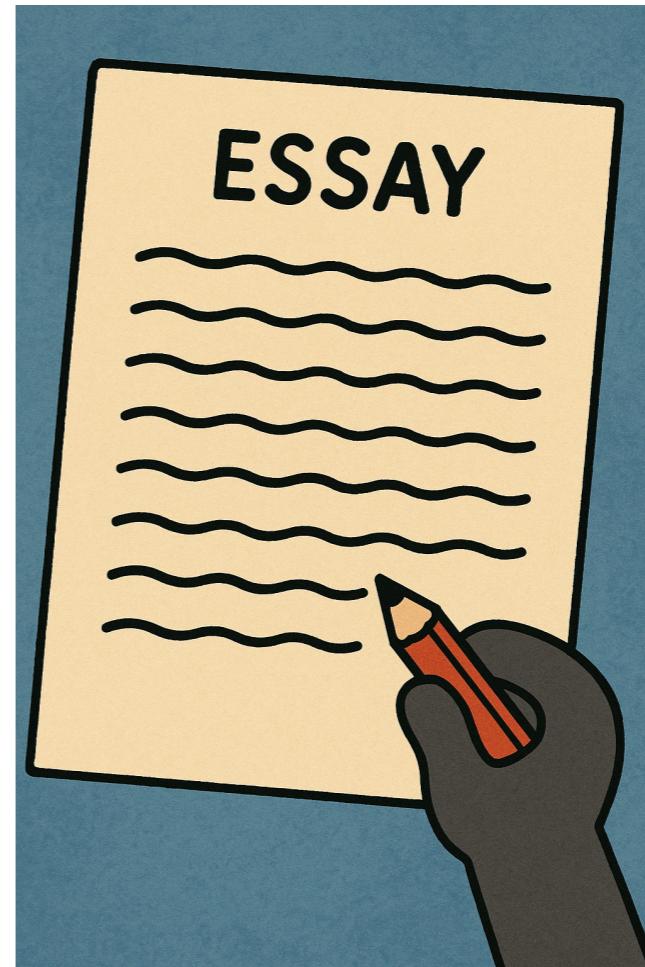


# Difficulty of Sequential Generation

Classification



Generation



# Difficulty of Sequential Generation

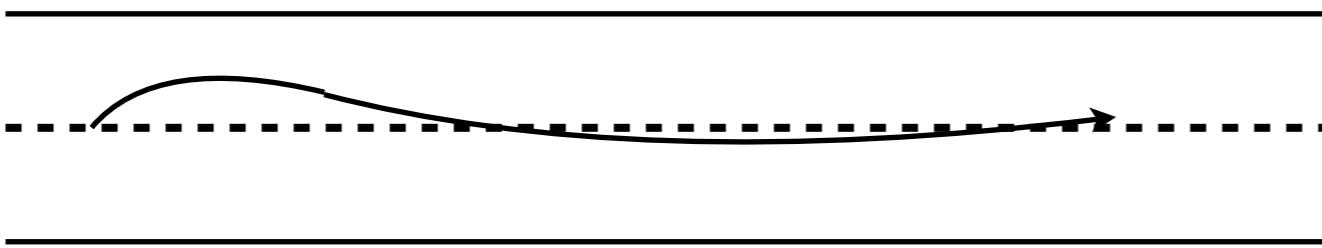
- Larger search space
  - More categories (e.g., 50K tokens)
  - Exponential growth wrt length

# Difficulty of Sequential Generation

- Larger search space
- Exposure bias
  - Teacher-force training
  - Self-generation

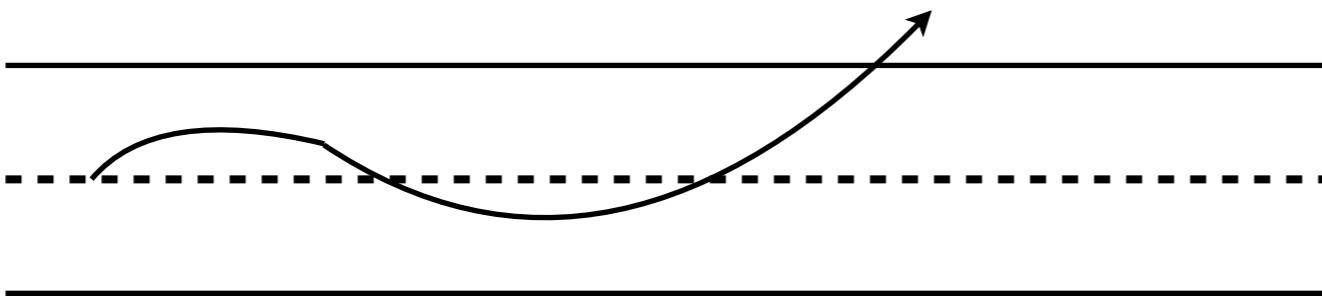
# Analogy: Autonomous Driving

Training phase



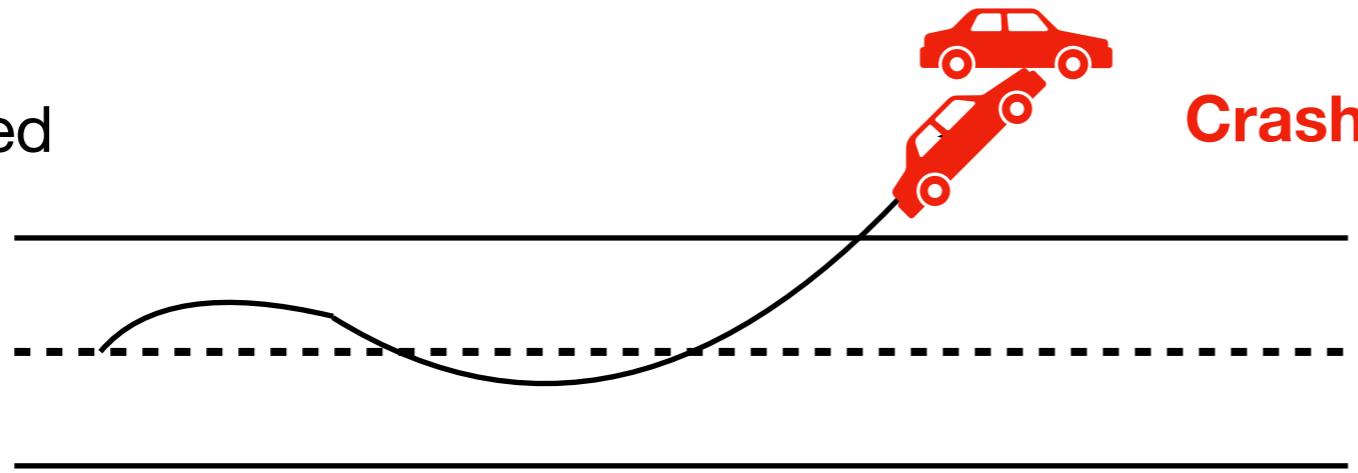
Inference phrase (deployment)

**Now what?**

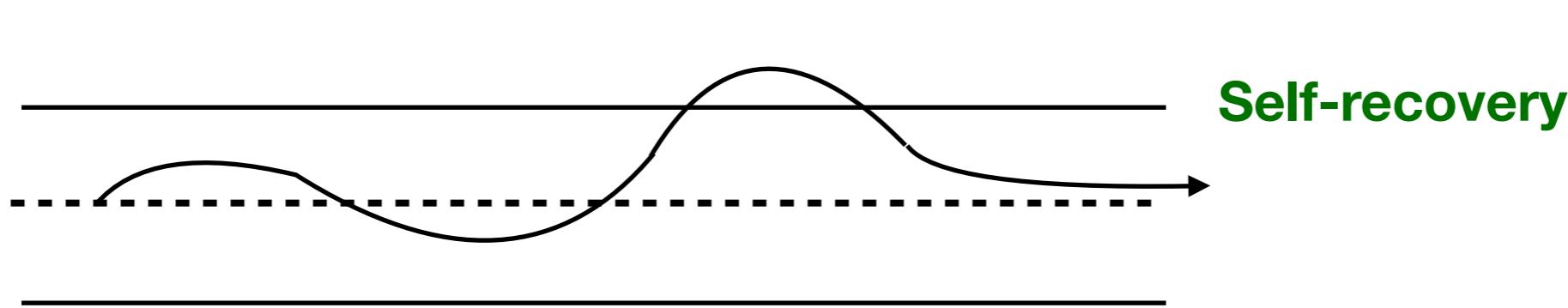


# Analogy: Autonomous Driving

Undesired

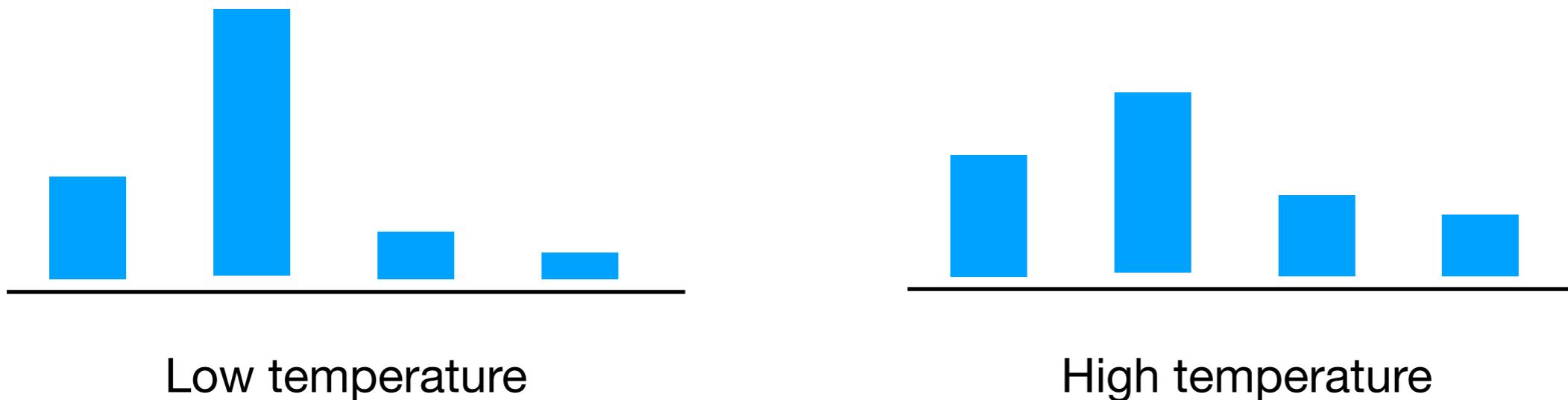


Desired

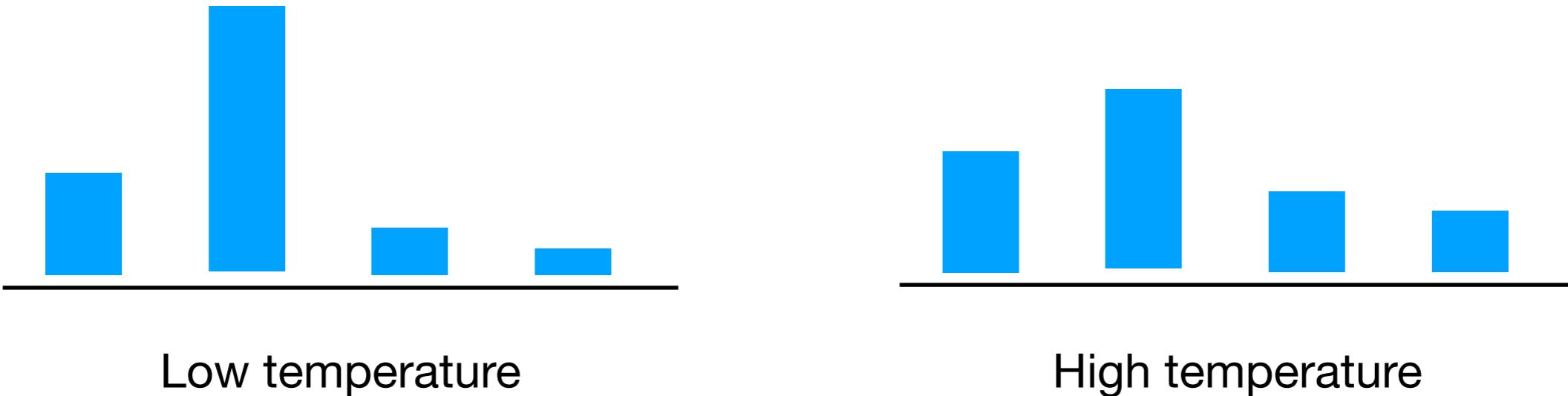


# Difficulty of Sequential Generation

- Larger search space
- Exposure bias
- Peakness is important!



# Peakness of Prediction



- Classification: argmax is unchanged
- Generation (structured prediction)

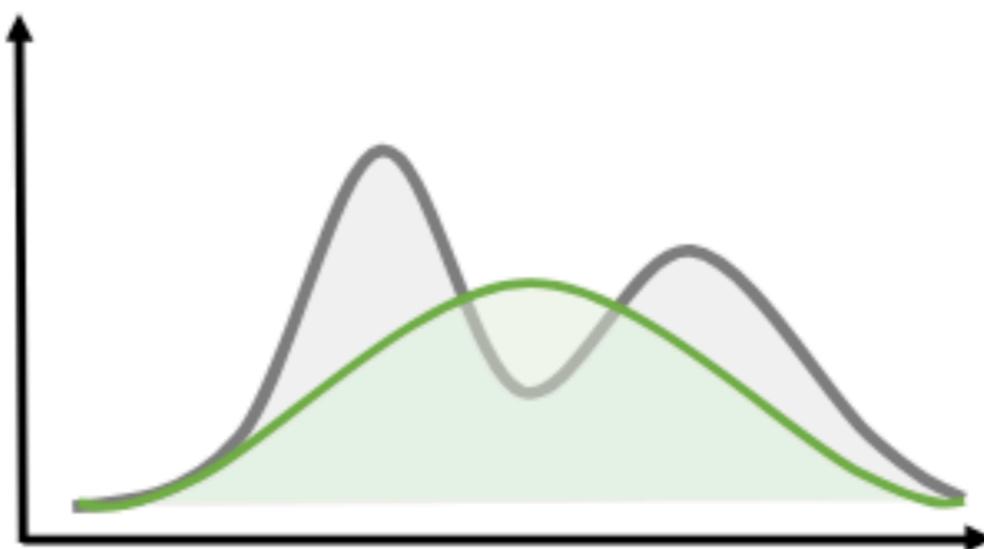
$$P(\text{structure}) = \prod P(\text{step})$$

Different steps considered jointly  
Temperature affects voting power

# ML Models Don't Learn the Right Peakness

$$\mathcal{L} = \sum_y p(y) \log \frac{p(y)}{q_{\theta}(y)}$$

- If target  $p(\cdot)$  is one-hot
  - Non-target labels are not learned explicitly
- Even if  $p(\cdot)$  is a full distribution
  - ML tends to give overly smooth distribution



# Tutorial Outline

- **Session I [45min]:**
  - Introduction
  - **f-Divergence KD [ACL'23]**
  - RL for LLM Distillation [NeurIPS'23, COLING'24]
- Break [15min]
- Session II [15min]
  - Multi-Teacher Distillation [ICLR'24, ACL'24, AAAI'25a,b]
  - KD Applications to LLMs
  - Conclusion and Future Work

# Straightforward Attempt

- SeqKD [Kim and Rush, 2016]

$$\begin{aligned} J_{\text{SeqKD}} &= \mathbb{E}_{\mathbf{Y} \sim p}[-\log q_{\theta}(\mathbf{Y})] \\ &\approx - \sum_{t=1}^{|\mathbf{y}|} \log q_{\theta}(\mathbf{y}_t \mid \mathbf{y}_{<t}) \end{aligned}$$

- Teacher predicts a sequence
- Student learns in a pseudo-supervised way
- Overcomes large space by hard KD
  - ▶ hard prefix / hard label

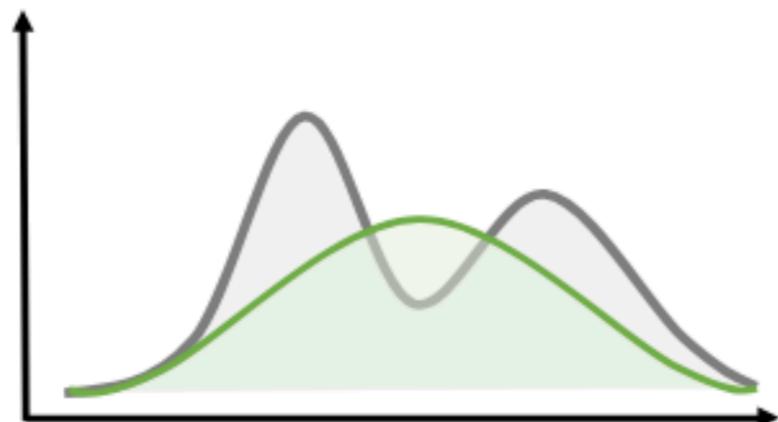
# Straightforward Attempt

- SeqKD [Kim and Rush, 2016]

$$J_{\text{SeqKD}} = \mathbb{E}_{\mathbf{Y} \sim p}[-\log q_\theta(\mathbf{Y})]$$

$$\approx - \sum_{t=1}^{|\mathbf{y}|} \log q_\theta(y_t | \mathbf{y}_{<t})$$

- Drawback



$$P(\text{structure}) = \prod P(\text{step})$$

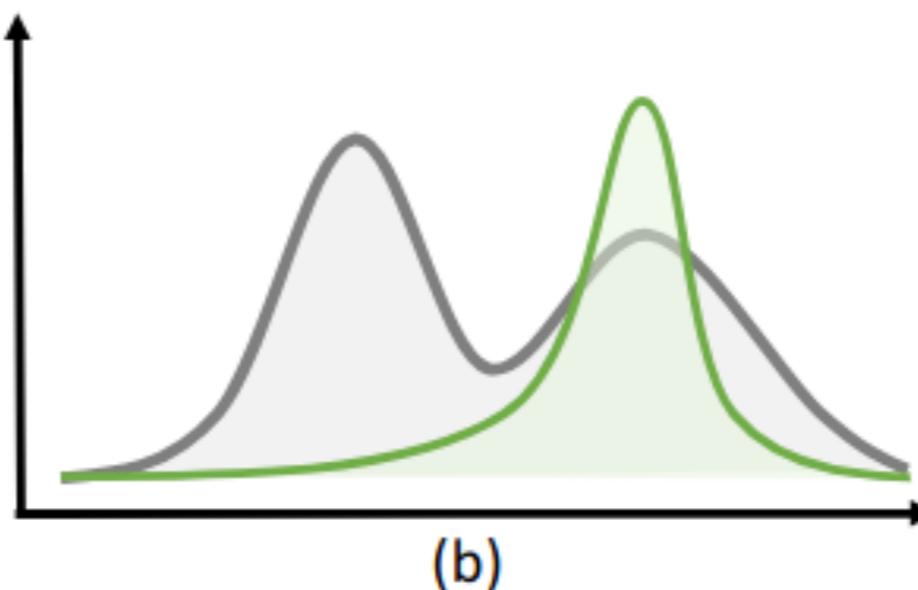
Different steps considered jointly  
Temperature affects voting power

# Reverse KL

- ENGINE

$$J_{\text{ENGINE}} = \mathbb{E}_{\mathbf{Y}' \sim q_{\theta}}[-\log p(\mathbf{Y}')] \approx - \sum_{t=1}^{|\mathbf{y}'|} \log p(\mathbf{y}'_t | \mathbf{y}'_{<t})$$

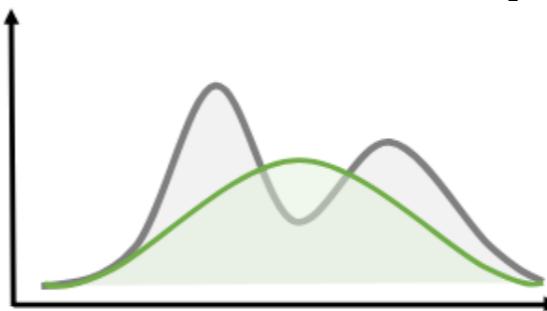
- Drawback



# KL vs RKL Distillation

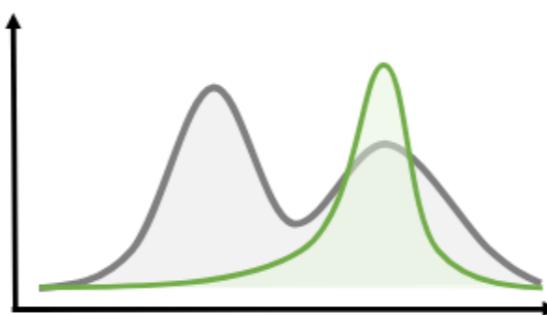
$$J_{\text{KL}} = D_{\text{KL}}(p \parallel q_{\theta}) = \mathbb{E}_{\mathbf{Y} \sim p} \left[ \log \frac{p(\mathbf{Y})}{q_{\theta}(\mathbf{Y})} \right]$$

Mode averaging



$$J_{\text{RKL}} = D_{\text{KL}}(q_{\theta} \parallel p) = \mathbb{E}_{\mathbf{Y}' \sim q_{\theta}} \left[ \log \frac{q_{\theta}(\mathbf{Y}')}{p(\mathbf{Y}')} \right]$$

Mode collapsing



Disclaimer: Discrete distributions don't have modes. Just an analogy.

# $f$ -divergence

$$D_f(p(t) \parallel q(t)) = \sum_t q(t) f\left(\frac{p(t)}{q(t)}\right)$$

Divergence	$f(t)$
Kullback–Leibler (KL)	$t \log t$
Reverse KL (RKL)	$-\log t$
Jensen–Shannon (JS)	$-(t+1) \log(\frac{t+1}{2}) + t \log t$
Total variation distance (TVD)	$\frac{1}{2} t - 1 $

# Symmetric Distillation

- Asymmetric distillation
  - $\text{KL}(p\|q_\theta) \neq \text{KL}(q_\theta\|p)$
  - $\text{RKL}(p\|q_\theta) \neq \text{RKL}(q_\theta\|p)$
  - Severe mode averaging or collapsing 😞
- Symmetric distillation
  - Jensen–Shannon (JS) distillation
  - Total variation distance (TVD) distillation
  - Balance between the two extremes 😊

# JS Distillation

$$J_{\text{JS}} = \frac{1}{2} \mathbb{E}_{\mathbf{Y} \sim p} \left[ \log \frac{p(\mathbf{Y})}{\underline{m(\mathbf{Y})}} \right] + \frac{1}{2} \mathbb{E}_{\mathbf{Y}' \sim q_\theta} \left[ \log \frac{q_\theta(\mathbf{Y}')}{\underline{m(\mathbf{Y}')}} \right]$$

- Denominator is a "Middle" distribution
- Non-zero if either  $p(\cdot)$  or  $q_\theta(\cdot)$  is not zero

$$m(\cdot) = \frac{1}{2} p(\cdot) + \frac{1}{2} q_\theta(\cdot)$$

# JS Distillation

$$J_{\text{JS}} = \frac{1}{2} \mathbb{E}_{\mathbf{Y} \sim p} \left[ \log \frac{p(\mathbf{Y})}{m(\mathbf{Y})} \right] + \frac{1}{2} \mathbb{E}_{\mathbf{Y}' \sim q_\theta} \left[ \log \frac{q_\theta(\mathbf{Y}')}{m(\mathbf{Y}')} \right]$$

$$\approx \frac{1}{2} \sum_{t=1}^{|\mathbf{y}|} \sum_{\mathbf{Y}_t \in V} -p(\mathbf{Y}_t | \mathbf{y}_{<t}) \log(m(\mathbf{Y}_t | \mathbf{y}_{<t}))$$

MC sampling for prefix

Step-wise decomposition

$$+ \frac{1}{2} \sum_{t=1}^{|\mathbf{y}'|} \sum_{\mathbf{Y}'_t \in V} [q_\theta(\mathbf{Y}'_t | \mathbf{y}'_{<t}) \log(q_\theta(\mathbf{Y}'_t | \mathbf{y}'_{<t}))$$

$$- q_\theta(\mathbf{Y}'_t | \mathbf{y}'_{<t}) \log(m(\mathbf{Y}'_t | \mathbf{y}'_{<t}))] + \text{const}$$

# TVD Distillation

$$J_{\text{TVD}} = \frac{1}{2} \sum_{\mathbf{Y}} |q_{\theta}(\mathbf{Y}) - p(\mathbf{Y})|$$

# TVD Distillation

$$\begin{aligned} J_{\text{TVD}} &= \frac{1}{2} \sum_{\mathbf{Y}} |q_{\theta}(\mathbf{Y}) - p(\mathbf{Y})| \\ &\leq \frac{1}{4} \mathbb{E}_{\mathbf{Y} \sim p} \left[ \sum_{t=1}^{|\mathbf{Y}|} \sum_{\mathbf{Y}_t \in V} |q_{\theta}(\mathbf{Y}_t | \mathbf{Y}_{<t}) - p(\mathbf{Y}_t | \mathbf{Y}_{<t})| \right] \\ &+ \frac{1}{4} \mathbb{E}_{\mathbf{Y}' \sim q_{\theta}} \left[ \sum_{t=1}^{|\mathbf{Y}'|} \sum_{\mathbf{Y}'_t \in V} |q_{\theta}(\mathbf{Y}'_t | \mathbf{Y}'_{<t}) - p(\mathbf{Y}'_t | \mathbf{Y}'_{<t})| \right] \end{aligned}$$

Step-wise decomposition: Upper bound

**Theorem 1.** (a) The sequence-level KL, RKL, and JS divergences can be decomposed exactly into step-wise terms. (b) The sequence-level TVD can be upper bounded by step-wise terms.

# TVD Distillation

$$\begin{aligned} J_{\text{TVD}} &= \frac{1}{2} \sum_{\mathbf{Y}} |q_{\theta}(\mathbf{Y}) - p(\mathbf{Y})| \\ &\leq \frac{1}{4} \mathbb{E}_{\mathbf{Y} \sim p} \left[ \sum_{t=1}^{|\mathbf{Y}|} \sum_{\mathbf{Y}_t \in V} |q_{\theta}(\mathbf{Y}_t | \mathbf{Y}_{<t}) - p(\mathbf{Y}_t | \mathbf{Y}_{<t})| \right] \\ &+ \frac{1}{4} \mathbb{E}_{\mathbf{Y}' \sim q_{\theta}} \left[ \sum_{t=1}^{|\mathbf{Y}'|} \sum_{\mathbf{Y}'_t \in V} |q_{\theta}(\mathbf{Y}'_t | \mathbf{Y}'_{<t}) - p(\mathbf{Y}'_t | \mathbf{Y}'_{<t})| \right] \\ &\approx \frac{1}{4} \sum_{t=1}^{|\mathbf{y}|} \sum_{\mathbf{Y}_t \in V} |q_{\theta}(\mathbf{Y}_t | \mathbf{y}_{<t}) - p(\mathbf{Y}_t | \mathbf{y}_{<t})| \quad \text{MC sampling for prefix} \\ &+ \frac{1}{4} \sum_{t=1}^{|\mathbf{y}'|} \sum_{\mathbf{Y}'_t \in V} |q_{\theta}(\mathbf{Y}'_t | \mathbf{y}'_{<t}) - p(\mathbf{Y}'_t | \mathbf{y}'_{<t})| \end{aligned}$$

# Sampling from Student

$$J_{\text{JS}} = \frac{1}{2} \mathbb{E}_{\mathbf{Y} \sim p} \left[ \log \frac{p(\mathbf{Y})}{m(\mathbf{Y})} \right] + \frac{1}{2} \mathbb{E}_{\mathbf{Y}' \sim q_\theta} \left[ \log \frac{q_\theta(\mathbf{Y}')}{m(\mathbf{Y}')} \right]$$

$$J_{\text{TVD}} = \frac{1}{2} \sum_{\mathbf{Y}} |q_\theta(\mathbf{Y}) - p(\mathbf{Y})|$$

$$\leq \frac{1}{4} \mathbb{E}_{\mathbf{Y} \sim p} \left[ \sum_{t=1}^{|\mathbf{Y}|} \sum_{\mathbf{Y}_t \in V} |q_\theta(\mathbf{Y}_t | \mathbf{Y}_{<t}) - p(\mathbf{Y}_t | \mathbf{Y}_{<t})| \right]$$

$$+ \frac{1}{4} \mathbb{E}_{\mathbf{Y}' \sim q_\theta} \left[ \sum_{t=1}^{|\mathbf{Y}'|} \sum_{\mathbf{Y}'_t \in V} |q_\theta(\mathbf{Y}'_t | \mathbf{Y}'_{<t}) - p(\mathbf{Y}'_t | \mathbf{Y}'_{<t})| \right]$$

Partially addressing exposure bias

# Main Results

Model		DART					
		BLEU4 <sup>†</sup>	METEOR <sup>†</sup>	TER <sup>†</sup>	BERTScore <sup>†</sup>	MoverScore <sup>†</sup>	BLEURT <sup>†</sup>
Teacher		48.56	39.28	45.45	83.04	68.17	40.56
Student	Non-distill (MLE)	43.12	35.71	49.97	79.76	65.65	29.10
	Pre-distill	45.60	36.99	47.10	81.39	66.75	34.08
	SeqKD	45.54	37.17	47.49	81.15	66.65	32.88
	ENGINE	44.40	36.51	50.63	80.18	66.20	30.94
	KL	46.24	37.45	46.89	81.60	67.07	35.31
	RKL	45.63	37.35	47.91	81.41	67.02	35.08
	JS	<u>46.85</u>	<u>37.75</u>	<u>46.50</u>	<u>81.93</u>	<u>67.30</u>	<u>36.81</u>
	TVD	<b>46.95</b>	<b>37.88</b>	<b>46.35</b>	<b>82.08</b>	<b>67.36</b>	<b>37.17</b>

Model		XSum			WMT16 EN-RO			Commonsense Dialogue		
		ROUGE-1 <sup>†</sup>	ROUGE-2 <sup>†</sup>	ROUGE-L <sup>†</sup>	BLEU4 <sup>†</sup>	chrF <sup>†</sup>	TER <sup>†</sup>	BLEU-1 <sup>†</sup>	BLEU-2 <sup>†</sup>	BERTScore <sup>†</sup>
Teacher		45.12	22.26	37.18	25.82	55.76	60.57	11.67	5.03	47.69
Student	Non-distill (MLE)	30.00	10.67	24.40	19.90	49.79	69.48	10.23	3.56	45.15
	Pre-distill	40.58	17.79	32.55	20.68	50.51	68.38	9.95	3.63	46.22
	SeqKD	39.13	17.53	32.34	21.20	50.81	67.66	10.85	4.17	46.94
	ENGINE	39.19	16.18	31.23	17.65	48.37	84.02	10.13	4.26	46.91
	KL	41.28	18.98	33.71	21.45	51.12	<b>66.74</b>	9.81	3.52	45.80
	RKL	<u>41.69</u>	19.02	33.92	20.46	50.33	70.78	10.48	4.01	46.68
	JS	41.65	<u>19.22</u>	<u>34.03</u>	<b>21.91</b>	<b>51.5</b>	<u>66.86</u>	<b>11.55</b>	<b>4.83</b>	<b>47.61</b>
	TVD	<b>41.76</b>	<b>19.30</b>	<b>34.10</b>	<u>21.73</u>	<u>51.13</u>	66.94	<u>11.39</u>	<u>4.73</u>	<u>47.30</u>

# Main Results

Model		DART					
		BLEU4 <sup>†</sup>	METEOR <sup>†</sup>	TER <sup>†</sup>	BERTScore <sup>†</sup>	MoverScore <sup>†</sup>	BLEURT <sup>†</sup>
Teacher		48.56	39.28	45.45	83.04	68.17	40.56
Student	Non-distill (MLE)	43.12	35.71	49.97	79.76	65.65	29.10
	Pre-distill	45.60	36.99	47.10	81.39	66.75	34.08
	SeqKD	45.54	37.17	47.49	81.15	66.65	32.88
	ENGINE	44.40	36.51	50.63	80.18	66.20	30.94
	KL	46.24	37.45	46.89	81.60	67.07	35.31
	RKL	45.63	37.35	47.91	81.41	67.02	35.08
	JS	<u>46.85</u>	<u>37.75</u>	<u>46.50</u>	<u>81.93</u>	<u>67.30</u>	<u>36.81</u>
	TVD	<b>46.95</b>	<b>37.88</b>	<b>46.35</b>	<b>82.08</b>	<b>67.36</b>	<b>37.17</b>

Model		XSum			WMT16 EN-RO			Commonsense Dialogue		
		ROUGE-1 <sup>†</sup>	ROUGE-2 <sup>†</sup>	ROUGE-L <sup>†</sup>	BLEU4 <sup>†</sup>	chrF <sup>†</sup>	TER <sup>†</sup>	BLEU-1 <sup>†</sup>	BLEU-2 <sup>†</sup>	BERTScore <sup>†</sup>
Teacher		45.12	22.26	37.18	25.82	55.76	60.57	11.67	5.03	47.69
Student	Non-distill (MLE)	30.00	10.67	24.40	19.90	49.79	69.48	10.23	3.56	45.15
	Pre-distill	40.58	17.79	32.55	20.68	50.51	68.38	9.95	3.63	46.22
	SeqKD	39.13	17.53	32.34	21.20	50.81	67.66	10.85	4.17	46.94
	ENGINE	39.19	16.18	31.23	17.65	48.37	84.02	10.13	4.26	46.91
	KL	41.28	18.98	33.71	21.45	51.12	<b>66.74</b>	9.81	3.52	45.80
	RKL	<u>41.69</u>	19.02	33.92	20.46	50.33	70.78	10.48	4.01	46.68
	JS	41.65	<u>19.22</u>	<u>34.03</u>	<b>21.91</b>	<b>51.5</b>	<u>66.86</u>	<b>11.55</b>	<b>4.83</b>	<b>47.61</b>
	TVD	<b>41.76</b>	<b>19.30</b>	<b>34.10</b>	<u>21.73</u>	<u>51.13</u>	66.94	<u>11.39</u>	<u>4.73</u>	<u>47.30</u>

- Teacher is significantly better
- Student struggles w/o distillation

# Main Results

Model	DART					
	BLEU4 <sup>†</sup>	METEOR <sup>†</sup>	TER <sup>†</sup>	BERTScore <sup>†</sup>	MoverScore <sup>†</sup>	BLEURT <sup>†</sup>
Teacher	48.56	39.28	45.45	83.04	68.17	40.56
Student	Non-distill (MLE)	43.12	35.71	49.97	79.76	65.65
	Pre-distill	45.60	36.99	47.10	81.39	66.75
	SeqKD	45.54	37.17	47.49	81.15	66.65
	ENGINE	44.40	36.51	50.63	80.18	66.20
	KL	46.24	37.45	46.89	81.60	67.07
	RKL	45.63	37.35	47.91	81.41	67.02
	JS	<u>46.85</u>	<u>37.75</u>	<u>46.50</u>	<u>81.93</u>	<u>67.30</u>
	TVD	<b>46.95</b>	<b>37.88</b>	<b>46.35</b>	<b>82.08</b>	<b>67.36</b>

Model	XSum			WMT16 EN-RO			Commonsense Dialogue		
	ROUGE-1 <sup>†</sup>	ROUGE-2 <sup>†</sup>	ROUGE-L <sup>†</sup>	BLEU4 <sup>†</sup>	chrF <sup>†</sup>	TER <sup>†</sup>	BLEU-1 <sup>†</sup>	BLEU-2 <sup>†</sup>	BERTScore <sup>†</sup>
Teacher	45.12	22.26	37.18	25.82	55.76	60.57	11.67	5.03	47.69
Student	Non-distill (MLE)	30.00	10.67	24.40	19.90	49.79	69.48	10.23	3.56
	Pre-distill	40.58	17.79	32.55	20.68	50.51	68.38	9.95	3.63
	SeqKD	39.13	17.53	32.34	21.20	50.81	67.66	10.85	4.17
	ENGINE	39.19	16.18	31.23	17.65	48.37	84.02	10.13	4.26
	KL	41.28	18.98	33.71	21.45	51.12	<b>66.74</b>	9.81	3.52
	RKL	<u>41.69</u>	19.02	33.92	20.46	50.33	70.78	10.48	4.01
	JS	41.65	<u>19.22</u>	<u>34.03</u>	<b>21.91</b>	<b>51.5</b>	<u>66.86</u>	<b>11.55</b>	<b>4.83</b>
	TVD	<b>41.76</b>	<b>19.30</b>	<b>34.10</b>	<u>21.73</u>	<u>51.13</u>	66.94	<u>11.39</u>	<u>4.73</u>

- KD significantly improves performance

# Main Results

Model		DART					
		BLEU4 <sup>†</sup>	METEOR <sup>†</sup>	TER <sup>†</sup>	BERTScore <sup>†</sup>	MoverScore <sup>†</sup>	BLEURT <sup>†</sup>
Teacher		48.56	39.28	45.45	83.04	68.17	40.56
Student	Non-distill (MLE)	43.12	35.71	49.97	79.76	65.65	29.10
	Pre-distill	45.60	36.99	47.10	81.39	66.75	34.08
	SeqKD	45.54	37.17	47.49	81.15	66.65	32.88
	ENGINE	44.40	36.51	50.63	80.18	66.20	30.94
	KL	46.24	37.45	46.89	81.60	67.07	35.31
	RKL	45.63	37.35	47.91	81.41	67.02	35.08
	JS	46.85	37.75	46.50	81.93	67.30	36.81
	TVD	<b>46.95</b>	<b>37.88</b>	<b>46.35</b>	<b>82.08</b>	<b>67.36</b>	<b>37.17</b>

Model		XSum			WMT16 EN-RO			Commonsense Dialogue		
		ROUGE-1 <sup>†</sup>	ROUGE-2 <sup>†</sup>	ROUGE-L <sup>†</sup>	BLEU4 <sup>†</sup>	chrF <sup>†</sup>	TER <sup>†</sup>	BLEU-1 <sup>†</sup>	BLEU-2 <sup>†</sup>	BERTScore <sup>†</sup>
Teacher		45.12	22.26	37.18	25.82	55.76	60.57	11.67	5.03	47.69
Student	Non-distill (MLE)	30.00	10.67	24.40	19.90	49.79	69.48	10.23	3.56	45.15
	Pre-distill	40.58	17.79	32.55	20.68	50.51	68.38	9.95	3.63	46.22
	SeqKD	39.13	17.53	32.34	21.20	50.81	67.66	10.85	4.17	46.94
	ENGINE	39.19	16.18	31.23	17.65	48.37	84.02	10.13	4.26	46.91
	KL	41.28	18.98	33.71	21.45	51.12	<b>66.74</b>	9.81	3.52	45.80
	RKL	41.69	19.02	33.92	20.46	50.33	70.78	10.48	4.01	46.68
	JS	41.65	19.22	34.03	<b>21.91</b>	<b>51.5</b>	66.86	<b>11.55</b>	<b>4.83</b>	<b>47.61</b>
	TVD	<b>41.76</b>	<b>19.30</b>	<b>34.10</b>	21.73	51.13	66.94	11.39	4.73	47.30

- SeqKD < KL
- Hard label < Soft label

# Main Results

Model		DART					
		BLEU4 <sup>†</sup>	METEOR <sup>†</sup>	TER <sup>†</sup>	BERTScore <sup>†</sup>	MoverScore <sup>†</sup>	BLEURT <sup>†</sup>
Teacher		48.56	39.28	45.45	83.04	68.17	40.56
Student	Non-distill (MLE)	43.12	35.71	49.97	79.76	65.65	29.10
	Pre-distill	45.60	36.99	47.10	81.39	66.75	34.08
	SeqKD	45.54	37.17	47.49	81.15	66.65	32.88
	ENGINE	44.40	36.51	50.63	80.18	66.20	30.94
	KL	46.24	37.45	46.89	81.60	67.07	35.31
	RKL	45.63	37.35	47.91	81.41	67.02	35.08
	JS	46.85	37.75	46.50	81.93	67.30	36.81
	TVD	<b>46.95</b>	<b>37.88</b>	<b>46.35</b>	<b>82.08</b>	<b>67.36</b>	<b>37.17</b>

Model		XSum			WMT16 EN-RO			Commonsense Dialogue		
		ROUGE-1 <sup>†</sup>	ROUGE-2 <sup>†</sup>	ROUGE-L <sup>†</sup>	BLEU4 <sup>†</sup>	chrF <sup>†</sup>	TER <sup>†</sup>	BLEU-1 <sup>†</sup>	BLEU-2 <sup>†</sup>	BERTScore <sup>†</sup>
Teacher		45.12	22.26	37.18	25.82	55.76	60.57	11.67	5.03	47.69
Student	Non-distill (MLE)	30.00	10.67	24.40	19.90	49.79	69.48	10.23	3.56	45.15
	Pre-distill	40.58	17.79	32.55	20.68	50.51	68.38	9.95	3.63	46.22
	SeqKD	39.13	17.53	32.34	21.20	50.81	67.66	10.85	4.17	46.94
	ENGINE	39.19	16.18	31.23	17.65	48.37	84.02	10.13	4.26	46.91
	KL	41.28	18.98	33.71	21.45	51.12	<b>66.74</b>	9.81	3.52	45.80
	RKL	41.69	19.02	33.92	20.46	50.33	70.78	10.48	4.01	46.68
	JS	41.65	19.22	34.03	<b>21.91</b>	<b>51.5</b>	<u>66.86</u>	<b>11.55</b>	<b>4.83</b>	<b>47.61</b>
	TVD	<b>41.76</b>	<b>19.30</b>	<b>34.10</b>	<u>21.73</u>	<u>51.13</u>	66.94	<u>11.39</u>	<u>4.73</u>	<u>47.30</u>

- ENGINE < RKL
- Hard label < Soft label

# Main Results

Model		DART					
		BLEU4 <sup>†</sup>	METEOR <sup>†</sup>	TER <sup>†</sup>	BERTScore <sup>†</sup>	MoverScore <sup>†</sup>	BLEURT <sup>†</sup>
Teacher		48.56	39.28	45.45	83.04	68.17	40.56
Student	Non-distill (MLE)	43.12	35.71	49.97	79.76	65.65	29.10
	Pre-distill	45.60	36.99	47.10	81.39	66.75	34.08
	SeqKD	45.54	37.17	47.49	81.15	66.65	32.88
	ENGINE	44.40	36.51	50.63	80.18	66.20	30.94
	KL	46.24	37.45	46.89	81.60	67.07	35.31
	RKL	45.63	37.35	47.91	81.41	67.02	35.08
	JS	<u>46.85</u>	<u>37.75</u>	<u>46.50</u>	<u>81.93</u>	<u>67.30</u>	<u>36.81</u>
	TVD	<b>46.95</b>	<b>37.88</b>	<b>46.35</b>	<b>82.08</b>	<b>67.36</b>	<b>37.17</b>

Model		XSum			WMT16 EN-RO			Commonsense Dialogue		
		ROUGE-1 <sup>†</sup>	ROUGE-2 <sup>†</sup>	ROUGE-L <sup>†</sup>	BLEU4 <sup>†</sup>	chrF <sup>†</sup>	TER <sup>†</sup>	BLEU-1 <sup>†</sup>	BLEU-2 <sup>†</sup>	BERTScore <sup>†</sup>
Teacher		45.12	22.26	37.18	25.82	55.76	60.57	11.67	5.03	47.69
Student	Non-distill (MLE)	30.00	10.67	24.40	19.90	49.79	69.48	10.23	3.56	45.15
	Pre-distill	40.58	17.79	32.55	20.68	50.51	68.38	9.95	3.63	46.22
	SeqKD	39.13	17.53	32.34	21.20	50.81	67.66	10.85	4.17	46.94
	ENGINE	39.19	16.18	31.23	17.65	48.37	84.02	10.13	4.26	46.91
	KL	41.28	18.98	33.71	21.45	51.12	<b>66.74</b>	9.81	3.52	45.80
	RKL	<u>41.69</u>	19.02	33.92	20.46	50.33	70.78	10.48	4.01	46.68
	JS	41.65	<u>19.22</u>	<u>34.03</u>	<b>21.91</b>	<b>51.5</b>	<u>66.86</u>	<b>11.55</b>	<b>4.83</b>	<b>47.61</b>
	TVD	<b>41.76</b>	<b>19.30</b>	<b>34.10</b>	21.73	51.13	66.94	11.39	<u>4.73</u>	47.30

- Asymmetric KD < Symmetric KD

# Likelihood VS Coverage

$$R_{\text{llh}} = \frac{1}{|\mathcal{D}_{\text{student}}|} \sum_{\mathbf{y}' \in \mathcal{D}_{\text{student}}} -\log p(\mathbf{y}')$$

$$R_{\text{cvg}} = \frac{1}{|\mathcal{D}_{\text{teacher}}|} \sum_{\mathbf{y} \in \mathcal{D}_{\text{teacher}}} -\log q_{\theta}(\mathbf{y})$$

Dataset	DART		XSum		MT <sub>EN-RO</sub>		CD	
TeacherDist	26.10		36.28		23.13		81.19	
Loss	$R_{\text{llh}}$	$R_{\text{cvg}}$	$R_{\text{llh}}$	$R_{\text{cvg}}$	$R_{\text{llh}}$	$R_{\text{cvg}}$	$R_{\text{llh}}$	$R_{\text{cvg}}$
KL	0.56	0.49	1.89	1.68	1.23	0.82	0.43	0.26
RKL	0.58	0.59	1.88	1.83	1.20	1.60	0.29	0.35
TVD	0.53	0.52	1.86	1.77	1.21	1.78	0.27	0.35
JS	0.51	0.48	1.88	1.75	1.13	1.34	0.30	0.33

# Summary

- F-divergence to balance mode averaging/collapse
- Step-wise decomposition
- Experimentation on 4 tasks

# Tutorial Outline

- **Session I [45min]:**
  - Introduction
  - f-Divergence KD [ACL'23]
  - **RL for LLM Distillation [NeurIPS'23, COLING'24]**
- Break [15min]
- Session II [15min]
  - Multi-Teacher Distillation [ICLR'24, ACL'24, AAAI'25a,b]
  - KD Applications to LLMs
  - Conclusion and Future Work

# Recall: Symmetric f-divergence KD

$$J_{\text{JS}} = \frac{1}{2} \mathbb{E}_{\mathbf{Y} \sim p} \left[ \log \frac{p(\mathbf{Y})}{m(\mathbf{Y})} \right] + \frac{1}{2} \mathbb{E}_{\mathbf{Y}' \sim q_\theta} \left[ \log \frac{q_\theta(\mathbf{Y}')}{m(\mathbf{Y}')} \right]$$

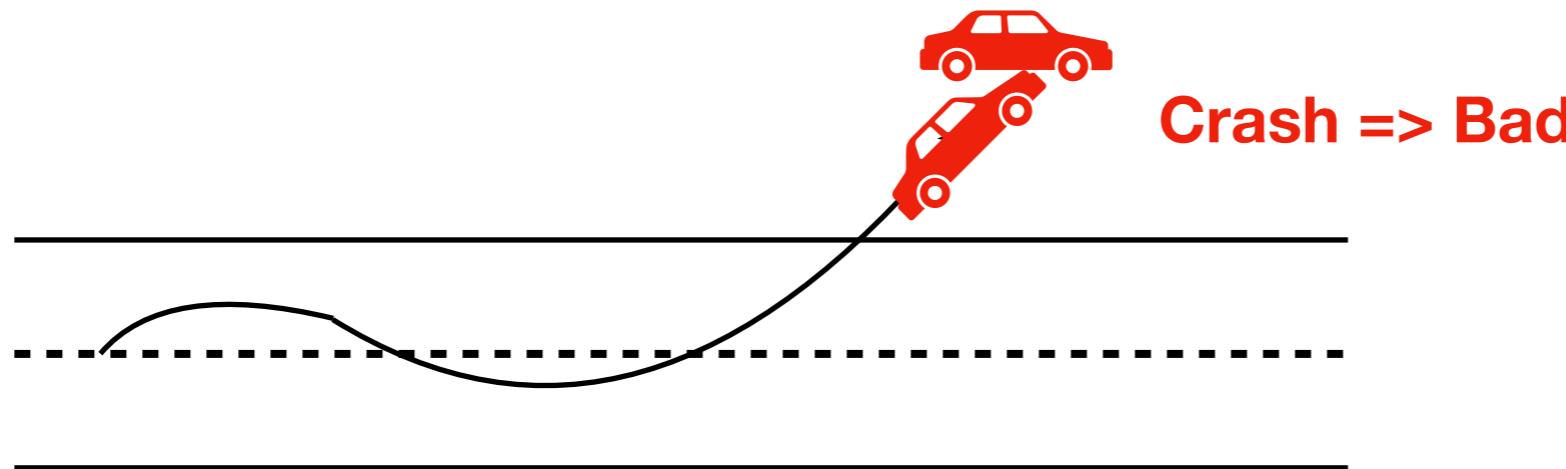
$$\begin{aligned} J_{\text{TVD}} &= \frac{1}{2} \sum_{\mathbf{Y}} |q_\theta(\mathbf{Y}) - p(\mathbf{Y})| \\ &\leq \frac{1}{4} \mathbb{E}_{\mathbf{Y} \sim p} \left[ \sum_{t=1}^{|\mathbf{Y}|} \sum_{\mathbf{Y}_t \in V} |q_\theta(\mathbf{Y}_t | \mathbf{Y}_{<t}) - p(\mathbf{Y}_t | \mathbf{Y}_{<t})| \right] \\ &+ \frac{1}{4} \mathbb{E}_{\mathbf{Y}' \sim q_\theta} \left[ \sum_{t=1}^{|\mathbf{Y}'|} \sum_{\mathbf{Y}'_t \in V} |q_\theta(\mathbf{Y}'_t | \mathbf{Y}'_{<t}) - p(\mathbf{Y}'_t | \mathbf{Y}'_{<t})| \right] \end{aligned}$$

- Sampling from student
- Partially addressing self-awareness

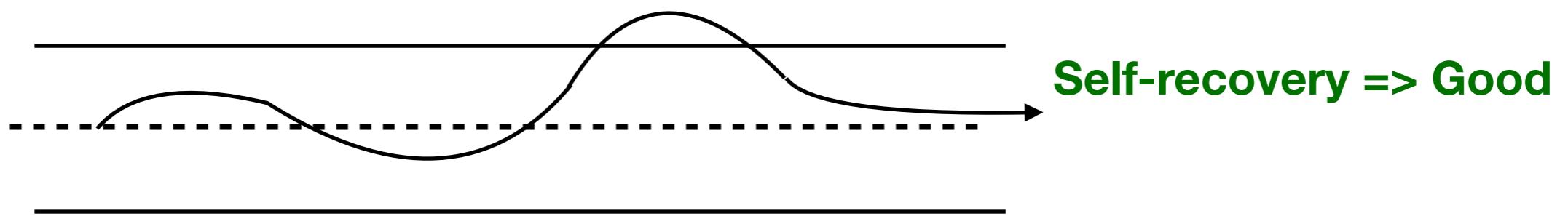
# Reinforcement Learning

- Learns in a trial-and-error fashion

Trial 1



Trial 2



# RL for Text Generation

- Markov decision process:  $\langle S, A, T, r \rangle$ 
  - $S$ : State space
    - ▶ Context (input + output prefix)
  - $A$ : Action space
    - ▶ Vocabulary (which word to predict)
  - $T$ : Transition
    - ▶ Appending the newly generated word
  - $r$ : Reward = ?

Goal: take actions to maximize the expected total reward

$$\mathbb{E} \left[ \sum_{t=0}^H r_t \right]$$

# Existing Reward Functions

- Similarity-based reward
  - BLEU, ROUGE scores
  - Learnable reward model
- Drawbacks
  - Requires ground-truth sequences
  - Similarity is a task-specific heuristic

# Existing Reward Functions

- RLHF: Human preference-based reward

$$\tau = s_1, a_1, s_2, a_2, \dots, s_t, a_t$$

$$\tau' = s'_1, a'_1, s'_2, a'_2, \dots, s'_t, a'_t$$

- Humans specify which trajectory is more preferred
  - A reward model is trained from such preference
- 
- Drawbacks
    - Data collection process requires design and heuristics
      - ▶ What are useful “negative” samples?
    - Still requires human preference labels

# Our Ideas and Contributions

- Inverse Reinforcement Learning (IRL)
  - Reward induction
  - Policy (predicted prob.)  $\pi \rightarrow q\text{-value function} \rightarrow r$
- The induced reward may be used for
  - Knowledge distillation from LLMs
  - Semi-supervised learning

# Inverse Reinforcement Learning (IRL)

- Markov decision process:  $\langle S, A, T, r \rangle$ 
  - $S$ : State space
    - ▶ Context (input + output prefix)
  - $A$ : Action space
    - ▶ Vocabulary (which word to predict)
  - $T$ : Transition
    - ▶ Appending the newly generated word
  - $r$ : Reward = ???

**There exists some underlying “true” reward function that we would like to recover**

# Value Functions

- Value function: How good a state is
- q-value function: How good an action is for a state

$$q^\pi(s, a) := \mathbb{E}_{\substack{a_t \sim \pi(\cdot | s_t) \\ s_{t+1} \sim T(\cdot | s_t, a_t)}} \left[ \sum_{t=1}^H r(s_t, a_t) | s_1 = s, a_1 = a \right]$$

Expected total reward from state  $s$ , action  $a$

# Bellman Optimality

Policy (predicted prob.)  $\pi \rightarrow$   $q$ -value function  $\rightarrow r$

- Consider the optimal policy  $\pi^*$  (wrt the underlying reward)
- Its q-value function must satisfy

$$q^{\pi^*}(s, a) = r(s, a) + \sum_{s' \in \mathcal{S}} T(s'|s, a) \max_{a'} q^{\pi^*}(s', a').$$

- One-step rollout
- Linking the underlying reward and q-value function

# Linking Policy and q-val

Policy (predicted prob.)  $\pi \rightarrow q\text{-value function} \rightarrow r$

- A commonly adopted assumption:
  - LM's policy follows a Boltzmann distribution

$$\pi(a | s) \propto \exp\{q^*(s, a)\}$$

- An action (word) with a higher q-value is more likely to be chosen by the LM
- Note: LM generation does **not** follow the optimal policy, but this stochastic policy

# Linking Policy and q-val

Policy (predicted prob.)  $\pi \rightarrow q\text{-value function} \rightarrow r$

- A commonly adopted assumption:
  - LM's policy follows a Boltzmann distribution
$$\pi(a | s) \propto \exp\{q^*(s, a)\}$$
  - An action (word) with a higher q-value is more likely to be chosen by the LM
- This looks familiar:  $\pi(a | s) = \text{softmax of logits}$   
⇒ Directly take logits (pre-softmax values) as  $q^*$

# Our Reward Induction Process

Policy (predicted prob.)  $\pi \rightarrow q\text{-value function} \rightarrow r$

- Suppose we have a **teacher LM** policy  $\pi$
- Inducing  $q$ -value function

$$q^*(s, a) = f(s, a), \text{ where } f(\cdot, \cdot) \text{ is the LM's logit}$$

# Our Reward Induction Process

Policy (predicted prob.)  $\pi \rightarrow q\text{-value function} \rightarrow r$

- Suppose we have a **teacher LM** policy  $\pi$
- Inducing  $q$ -value function

$q^*(s, a) = f(s, a)$ , where  $f(\cdot, \cdot)$  is the LM's logit

- Inducing the reward

– Bellman optimality  $q^{\pi^*}(s, a) = r(s, a) + \sum_{s' \in \mathcal{S}} T(s'|s, a) \max_{a'} q^{\pi^*}(s', a')$

# Our Reward Induction Process

Policy (predicted prob.)  $\pi \rightarrow q\text{-value function} \rightarrow r$

- Suppose we have a **teacher LM** policy  $\pi$
- Inducing  $q$ -value function

$q^*(s, a) = f(s, a)$ , where  $f(\cdot, \cdot)$  is the LM's logit

- Inducing the reward

- Bellman optimality 
$$f(s, a) = r(s, a) + \sum_{s' \in \mathcal{S}} T(s'|s, a) \max_{a'} q^*(s', a')$$
  
$$f(s + [a], a') \quad \text{next step}$$

$$r(s, a) = f_\omega(s, a) - \max_{a' \in \mathcal{A}} f_\omega(s + [a], a')$$

# Analysis

**Theorem 1.** Suppose the value function  $q$  in Eqn. (3) and the seq2seq model  $f$  in Eqn. (4) have the same parametrization  $\omega$ , we have

$$L_{\text{TF}}(\omega; D) = -\log P_{\text{IRL}}(\mathcal{D}|\omega) + \text{const.} \quad (5)$$

- MLE in IRL learns the initial state (const.), state transition (deterministic here), and the policy

**Theorem 2.** Let  $r^*$  be an underlying true reward function and  $q^*$  be the corresponding optimal value function. Given an approximate value function  $q$ , we denote by  $r$  the reward function derived from Eqn. (7). Then, we must have  $\|r - r^*\|_\infty$  bounded by  $O(\|q - q^*\|_\infty)$ . Here,  $\|\cdot\|_\infty$  takes the maximum absolute value over all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ .

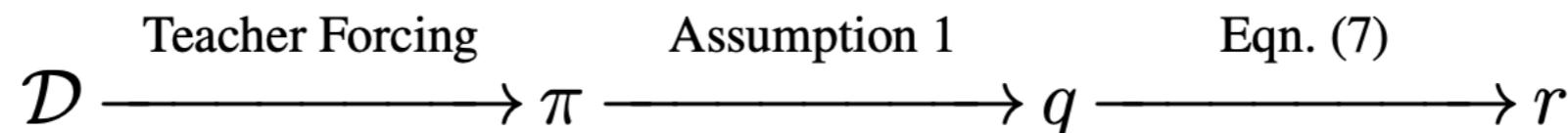
- What if  $q$  is not optimal? Near optimal  $q \Rightarrow$  Near optimal  $r$

**Theorem 3.** Suppose  $r'(s, a) = r(s, a) + c_s - c_{s+[a]}$ . Then the learned policies under  $r'(s, a)$  and  $r(s, a)$  are the same.

- Logit may be shifted arbitrarily

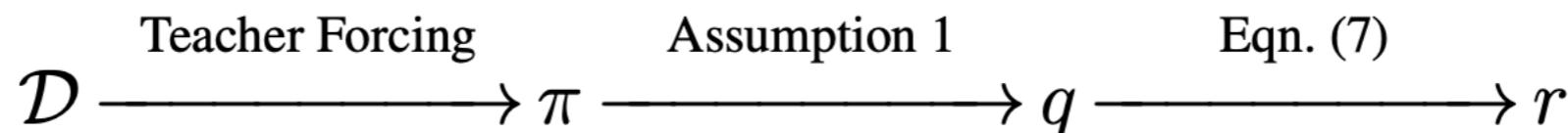
# Applications of Reward Induction

- Semi-supervised learning (self-distillation) [NeurIPS'23]
  - Small labeled data  $D$
  - Large unlabeled data  $D_u$

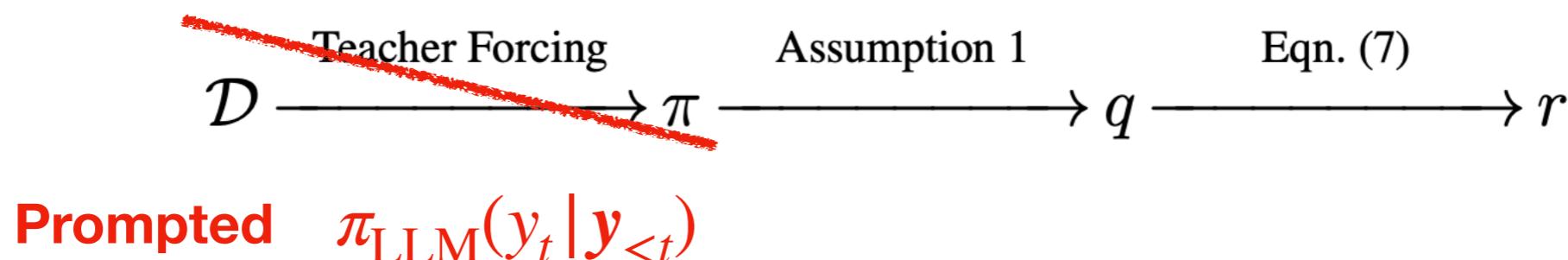


# Applications of Reward Induction

- Semi-supervised learning (self-distillation) [NeurIPS'23]
  - Small labeled data  $D$
  - Large unlabeled data  $D_u$



- LLM distillation [COLING'24]



# Results of Self-Distillation

Table 1: Main results.  $\uparrow/\downarrow$ The higher/lower, the better.  $^\dagger$ Quoted from Wen et al. [62] on deduplicated dialogue datasets.  $^\ddagger$ Quoted from [29].  $^\S$ Quoted from [11]. For the paraphrase generation metric, we have iBLEU =  $(1 - \alpha)$  BLEU –  $\alpha$  SBLEU.

(a) Dialogue generation.

Method	BLEU2 $^\uparrow$	BLEU4 $^\uparrow$
Parallel DailyDialog		
AdaLabel $^\dagger$ [60]	6.72	2.29
DialogBERT $^\dagger$ [16]	5.42	2.16
T5-Base [42]	<b>8.96</b>	<b>3.69</b>
+ Parallel OpenSubtitles		
[T5-Base] Fully Supervised	8.75	3.06
+ Non-Parallel OpenSubtitles		
[T5-Base] Self-Training	9.10	3.73
[T5-Base] R-Regression	10.34	4.18
[T5-Base] Ours	<b>11.02</b>	<b>4.30</b>

(b) Paraphrase generation. “Copy” refers to directly copying the input sentence.

Method	BLEU4 $^\uparrow$	SBLEU4 $^\downarrow$	iBLEU4 $^\uparrow$
Copy	29.88	100.0	16.89
Parallel Quora Generation			
Dagger $^\ddagger$ [12]	28.42	66.98	18.88
RL-NN $^\ddagger$ [40]	20.98	<b>40.52</b>	14.83
T5-Base [42]	<b>30.83</b>	44.77	<b>23.27</b>
+ Non-Parallel Quora Generation			
LTS $^\S$ [11]	29.25	71.25	19.20
[T5-Base] Self-Training	31.39	48.02	23.44
[T5-Base] R-Regression	30.77	<b>44.23</b>	23.27
[T5-Base] Ours	<b>31.47</b>	45.43	<b>23.78</b>

# Sparse vs Dense Reward

Table 2: Comparing sparse and dense reward functions.

(a) Dialogue generation.

Sparse	Method	BLEU2 $\uparrow$	BLEU4 $\uparrow$
-	Self-Training [23]	9.10	3.73
Yes	R-Regression [66]	9.45	3.73
	Induced-R	<b>9.75</b>	<b>3.99</b>
No	R-Regression [66]	10.34	4.18
	Induced-R	<b>11.02</b>	<b>4.30</b>

(b) Paraphrase generation.

Sparse	Method	BLEU4 $\uparrow$	SBLEU $\downarrow$	iBLEU4 $\uparrow$
-	Self-Training [23]	31.39	48.11	23.44
Yes	R-Regression [66]	30.78	<b>44.32</b>	23.27
	Induced-R	<b>31.28</b>	45.22	<b>23.63</b>
No	R-Regression [66]	30.77	<b>44.23</b>	23.27
	Induced-R	<b>31.47</b>	45.43	<b>23.78</b>

Note: The sparse variant defers all the rewards to the end.

# Results of LLM KD

- Teacher: T0-3B
- Student: T5-base (220 million parameters)

Model		DailyDialog		OpenSubtitles		CNN/DailyMail		
		BLEU2	BLEU4	BLEU2	BLEU4	ROUGE-1	ROUGE-2	ROUGE-L
1	Prompting Teacher	5.57	1.49	4.67	1.51	36.16	14.99	24.05
2	Prompting Student	1.35	0.31	1.21	0.25	21.23	6.73	17.88
3	Distilled Students	SeqKD	6.19	1.71	3.87	1.35	35.46	14.52
4		KL	5.03	1.40	3.84	1.33	34.11	14.21
5		RKL	5.02	1.29	4.12	1.36	32.07	13.77
6		JS	6.60	1.73	3.64	0.87	35.88	14.72
7		Our LLMR	7.00	1.88	5.13	1.85	36.42	15.21
								24.83

Table 1: Main results on dialogue generation and summarization tasks.

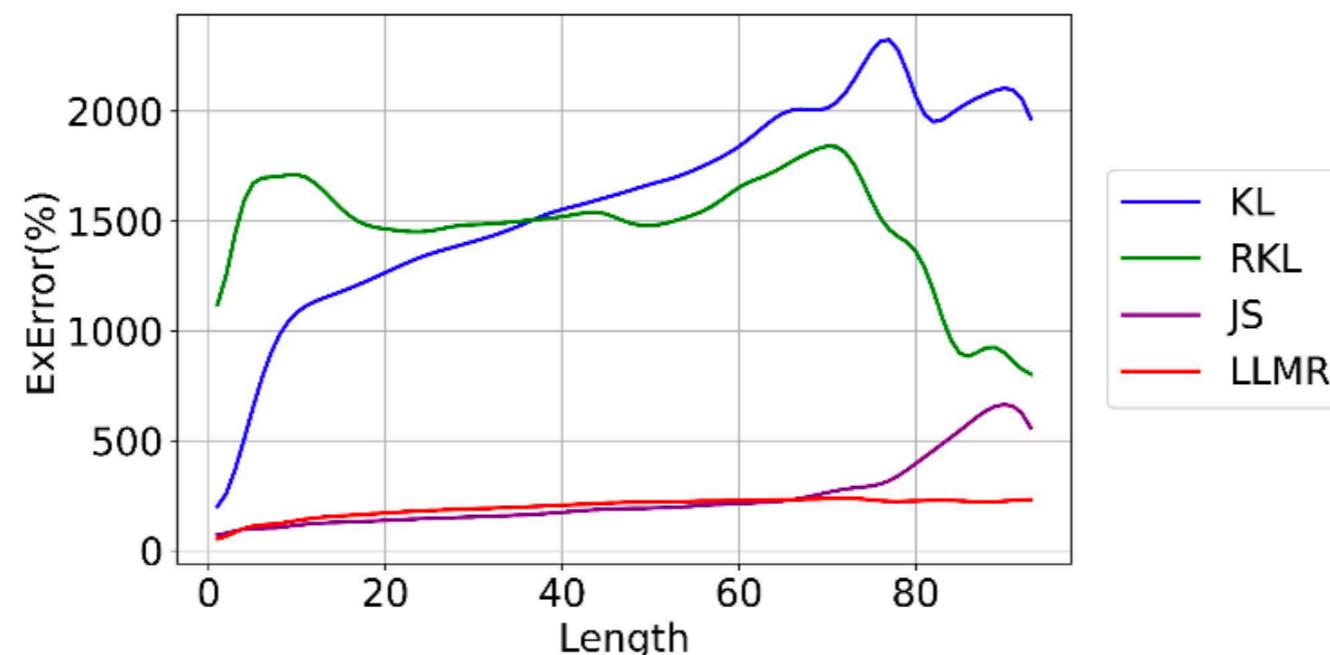
# Excess Error

[Arora et al., 2022]

$$D_s = \sum_{t=1}^T \mathbb{E}_{\substack{\mathbf{y}_{<t} \sim q_\theta(\cdot | \mathbf{x}) \\ \mathbf{y}_t \sim p(\cdot | \mathbf{y}_{<t}, \mathbf{x})}} \left[ \log \frac{p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})}{q_\theta(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})} \right]$$

$$D_t = \sum_{t=1}^T \mathbb{E}_{\substack{\mathbf{y}_{<t} \sim p(\cdot | \mathbf{x}) \\ \mathbf{y}_t \sim p(\cdot | \mathbf{y}_{<t}, \mathbf{x})}} \left[ \log \frac{p(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})}{q_\theta(\mathbf{y}_t | \mathbf{y}_{<t}, \mathbf{x})} \right]$$

$$\text{ExError \%} = \frac{D_s - D_t}{D_t}$$



# Summary

- Reward induction from policy

$$\pi \rightarrow q^* \rightarrow r$$

- Applications
  - Semi-supervised learning (self-KD)
  - LLM knowledge distillation

# References

- Chung, Inseop, et al. "Feature-map-level online adversarial knowledge distillation." *International Conference on Machine Learning*. PMLR, 2020.
- Guo, Daya, et al. "Deepseek-r1: Incentivizing reasoning capability in LLMs via reinforcement learning." arXiv preprint arXiv:2501.12948 (2025).
- Hao, Yongchang, Yuxin Liu, and Lili Mou. "Teacher forcing recovers reward functions for text generation." *Advances in Neural Information Processing Systems* 35 (2022): 12594-12607.
- Hinton, Geoffrey, Oriol Vinyals, and Jeff Dean. "Distilling the knowledge in a neural network." arXiv preprint arXiv:1503.02531 (2015).
- Jiao, Xiaoqi, et al. "TinyBERT: Distilling BERT for Natural Language Understanding." *Findings of the Association for Computational Linguistics: EMNLP 2020*. 2020.
- Kim, Yoon, and Alexander M. Rush. "Sequence-level knowledge distillation." *Proceedings of the 2016 conference on empirical methods in natural language processing*. 2016.
- Li, Dongheng, Yongchang Hao, and Lili Mou. "Llmr: Knowledge distillation with a large language model-induced reward." COLING, 2024
- Narayanan, Deepak, et al. "Efficient large-scale language model training on GPU clusters using megatron-LM." *Proceedings of the international conference for high performance computing, networking, storage and analysis*. 2021.
- Passban, Peyman, et al. "ALP-KD: Attention-based layer projection for knowledge distillation." *Proceedings of the AAAI Conference on artificial intelligence*. Vol. 35. No. 15. 2021.
- Sun, Siqi, et al. "Patient Knowledge Distillation for BERT Model Compression." *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2019.
- Tang, Jiaxi, and Ke Wang. "Ranking distillation: Learning compact ranking models with high performance for recommender system." *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 2018.
- Tunstall, Lewis, et al. "Zephyr: Direct distillation of lm alignment." arXiv preprint arXiv:2310.16944 (2023).
- Wang, Wenhui, et al. "MiniLM: Deep self-attention distillation for task-agnostic compression of pre-trained transformers." *Advances in neural information processing systems* 33 (2020): 5776-5788.
- Xu, Xiaohan, et al. "A survey on knowledge distillation of large language models." arXiv preprint arXiv:2402.13116 (2024).
- Wen, Yuqiao, et al. "F-divergence minimization for sequence-level knowledge distillation." ACL, 2023.
- Yu, Zony, Yuqiao Wen, and Lili Mou. "Revisiting Intermediate-Layer Matching in Knowledge Distillation: Layer-Selection Strategy Doesn't Matter (Much)." AACL-Findings, 2025.

# Tutorial Outline

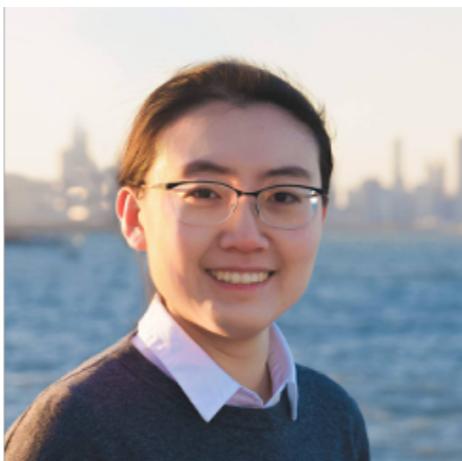
- Session I [45min]:
  - Introduction
  - f-Divergence KD [ACL'23]
  - RL for LLM Distillation [NeurIPS'23, COLING'24]
- Break [15min] 
- Session II [15min]
  - Multi-Teacher Distillation [ICLR'24, ACL'24, AAAI'25a,b]
  - KD Applications to LLMs
  - Conclusion and Future Work

# Knowledge Distillation for Language Models: Challenges and Opportunities with Sequential Data

Yuqiao Wen<sup>1,3</sup>



Freda Shi<sup>2,4,5</sup>



Lili Mou<sup>1,3,5</sup>



[yq.when@gmail.com](mailto:yq.when@gmail.com)

[fhs@uwaterloo.ca](mailto:fhs@uwaterloo.ca)

[doublepower.mou@gmail.com](mailto:doublepower.mou@gmail.com)

<sup>1</sup>Dept. Computing Science, University of Alberta

<sup>2</sup>David R. Cheriton School of Computer Science, University of Waterloo

<sup>3</sup>Alberta Machine Intelligence Institute (Amii)

<sup>4</sup>Vector Institute

<sup>5</sup>Canada CIFAR AI Chair

AAAI'26 Tutorial

# Tutorial Outline

- Session I [45min]:
  - Introduction
  - f-Divergence KD [ACL'23]
  - RL for LLM Distillation [NeurIPS'23, COLING'24]
- Break [15min]
- Session II [45min] 
  - Multi-Teacher Distillation [ICLR'24, ACL'24, AAAI'25a,b]
  - KD Applications to LLMs
  - Conclusion and Future Work

# Multi-Teacher KD

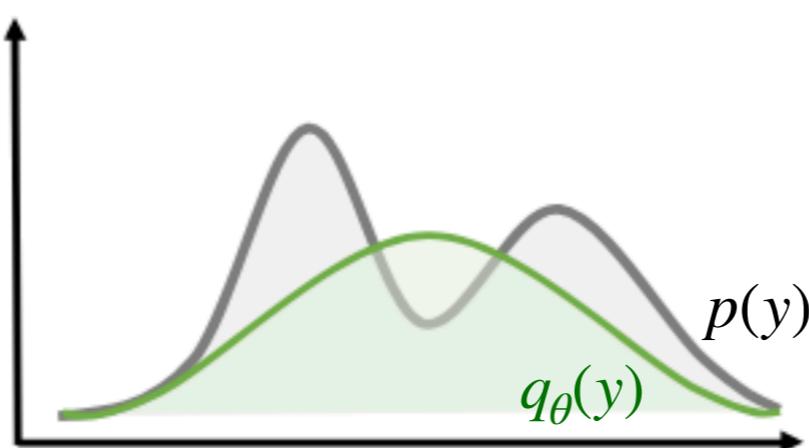
- Learning from multiple teachers
  - Consolidating different knowledge
  - Smoothing out noise
- Naive attempts
  - UnionKD
    - Taking union of teachers' outputs
  - AvgKD
    - Averaging teachers' output probabilities

# Challenge: Multi-Modal Distributions

- KD is typically done by cross-entropy loss

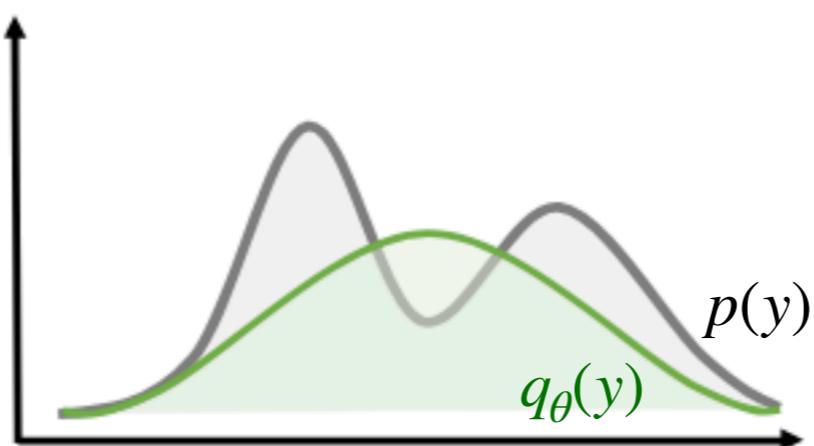
$$\mathcal{L} = \sum_y p(y) \log \frac{p(y)}{q_\theta(y)}$$

- $p(y)$ : target distribution (mixture of teachers)
  - $q_\theta(y)$ : student's predicted distribution
  - Learns an overly smooth distribution
- Severe multi-modal problem
  - UnionKD and AvgKD do **not** work well



# Our Idea

- Ensemble-Then-Distill Approach
  - Consolidate teachers' knowledge
  - Learn from the consensus



# Application: Zero-Shot Machine Translation

- In MT, training data are usually English-centric
  - E.g., German  $\leftrightarrow$  English, English  $\leftrightarrow$  Romanian
- Many translation directions are zero-shot
  - German  $\leftrightarrow$  Romanian

# Previous Work

- Direct zero-shot translation
  - Trained with various supervised directions
  - Zero-shot ability emerges
- Pivot translation
  - German → English → Romanian
  - Avoids zero-shot
  - Leads to error accumulation

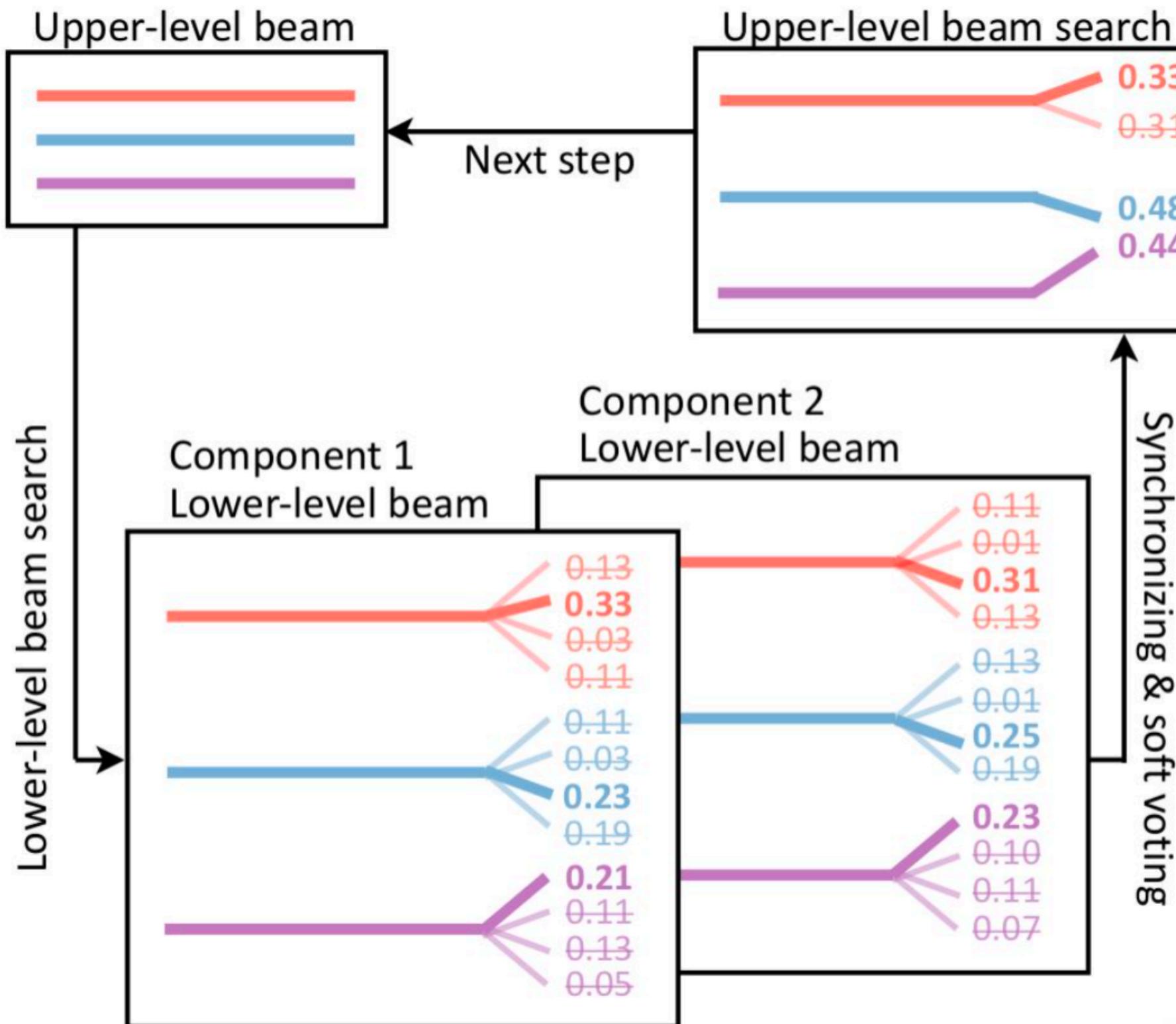
# Our Idea

- Build an ensemble of direct and various pivot translations
- Train a smaller model from the ensemble
  - Improve efficiency and/or performance

# Naïve Attempts

- Average voting
  - Over-smoothing
- 0/1 voting:
  - Too many candidates, ineffective
- Sequence-level ensemble
  - Maximum Bayes risk (MBR) decoding
  - Selective, subject to the quality of individuals

# Our Bi-Level Beam Search



# Our Bi-Level Beam Search

- **Lower-level beam**

- Each ensemble individual explores on its own

$$B'_{t,k} = \text{top-}Z\left\{ \langle \mathbf{y}_{1:t-1} \oplus \mathbf{y}, p \cdot p_k(\mathbf{y} | \mathbf{y}_{1:t-1}, \mathbf{x}) \rangle : \langle \mathbf{y}_{1:t-1}, p \rangle \in B_{t-1}, \mathbf{y} \in V \right\}$$

- **Upper-level beam**

- Aggregate individuals to maintain efficiency
  - Soft (probability) voting on lower-level beam samples

$$C_t = \bigcup_k \{ \mathbf{y} : \langle \mathbf{y}, p \rangle \in B'_{t,k} \} \quad B_t = \text{top-}Z \left\{ \left\langle \mathbf{y}, \sum_{\substack{k: k=1, \dots, K \\ \langle \mathbf{y}', p \rangle \in B'_{t,k}: \mathbf{y}' = \mathbf{y}}} p \right\rangle : \mathbf{y} \in C_t \right\}$$

- Our bi-level beam search is “nucleus voting,” analogous to “nucleus sampling”

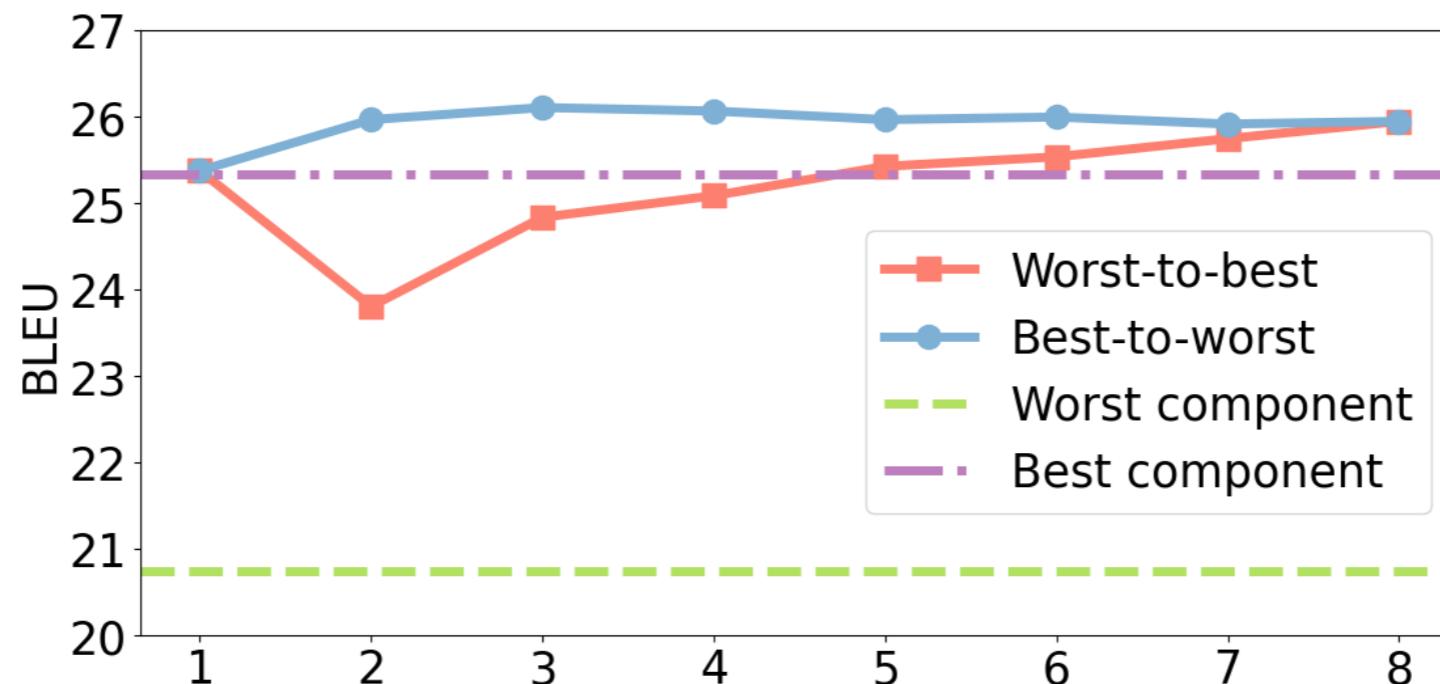
# Ensemble Performance

	#	Model	Average	it-nl	it-ro	nl-it	nl-ro	ro-it	ro-nl	
IWSLT	1	Direct translation (Liu et al., 2021) <sup>†</sup>	17.7	18.5	17.8	17.9	15.5	19.6	16.8	
	2	Direct translation (our replication)	17.29	17.46	<b>17.48</b>	18.23	<b>14.63</b>	19.65	16.26	
	3	Pivoting (en)	16.19	17.49	15.09	16.79	13.05	18.34	16.37	
	4	Word-level averaging ensemble	16.52	16.48	16.49	17.53	13.80	19.07	15.77	
	5	Word-level voting ensemble	16.99	<u>17.58</u>	16.38	17.78	14.13	19.21	<u>16.84</u>	
	6	Sequence-level voting ensemble (MBR)	16.66	<u>16.54</u>	16.51	17.72	13.64	19.58	<u>15.98</u>	
	7	EBBS (ours)	<b>18.24</b>	<b>19.52</b>	<u>17.09</u>	<b>19.06</b>	<u>14.58</u>	<b>20.75</b>	<b>18.45</b>	
Europarl	#	Model	Average	da-de	da-es	da-fi	da-fr	da-it	da-nl	da-pt
	1	Direct translation (Liu et al., 2021) <sup>†</sup>	26.9	24.2	33.1	18.1	30.6	26.1	26.3	29.9
	2	Direct translation (our replication)	27.74	<u>26.24</u>	33.64	18.95	31.01	26.58	27.36	30.38
	3	Pivoting (en)	27.67	<u>25.15</u>	33.79	18.63	31.45	<u>27.12</u>	26.71	30.82
	4	Word-level averaging ensemble	27.76	26.21	33.64	18.88	31.09	26.71	27.40	30.37
	5	Word-level voting ensemble	27.75	25.93	33.90	18.47	31.29	27.08	26.74	<u>30.84</u>
	6	Sequence-level voting ensemble (MBR)	<u>27.85</u>	25.88	<u>33.93</u>	<u>19.10</u>	<u>31.47</u>	<u>27.12</u>	<u>26.98</u>	<u>30.49</u>
	7	EBBS (ours)	<b>28.36</b>	<b>26.32</b>	<b>34.28</b>	<b>19.43</b>	<b>31.97</b>	<b>27.67</b>	<b>27.78</b>	<b>31.08</b>

Dataset	Method	Avg. BLEU	Wins	Losses
IWSLT	Direct	17.28	2	4
	Ensemble	<b>18.23</b>	<b>4</b>	<b>2</b>
Europarl	Direct	27.85	4	52
	Ensemble	<b>28.44</b>	<b>52</b>	<b>4</b>
Overall	Direct	26.83	6	56
	Ensemble	<b>27.45</b>	<b>56</b>	<b>6</b>
p-value		3e-11		

# Ensemble Components

Method	BLEU $\uparrow$	BLEU1 $\uparrow$	BLEU2 $\uparrow$	BLEU3 $\uparrow$	BLEU4 $\uparrow$	chrF2++ $\uparrow$	TER $\downarrow$	COMET $\uparrow$
Direct translation	<u>25.33</u>	56.32	30.08	19.01	<u>12.78</u>	<u>52.32</u>	66.56	0.8276
Pivoting (en)	25.08	<u>56.76</u>	<u>30.29</u>	<u>19.06</u>	12.66	51.92	<u>66.24</u>	<u>0.8322</u>
Pivoting (es)	24.40	55.38	29.08	18.22	12.09	51.71	67.91	0.8192
Pivoting (pt)	24.34	55.46	29.02	18.13	12.02	51.61	67.68	0.8191
Pivoting (fr)	24.20	55.41	29.02	18.00	11.84	51.61	67.84	0.8208
Pivoting (de)	23.65	55.33	28.69	17.67	11.54	50.70	67.89	0.8157
Pivoting (da)	23.12	54.81	27.96	17.12	11.12	50.36	69.00	0.8156
Pivoting (fi)	20.74	53.54	26.10	15.43	9.79	48.11	70.59	0.8051
Our EBBS	<b>26.10</b>	<b>57.07</b>	<b>31.00</b>	<b>19.76</b>	<b>13.28</b>	<b>52.75</b>	<b>65.63</b>	<b>0.8340</b>



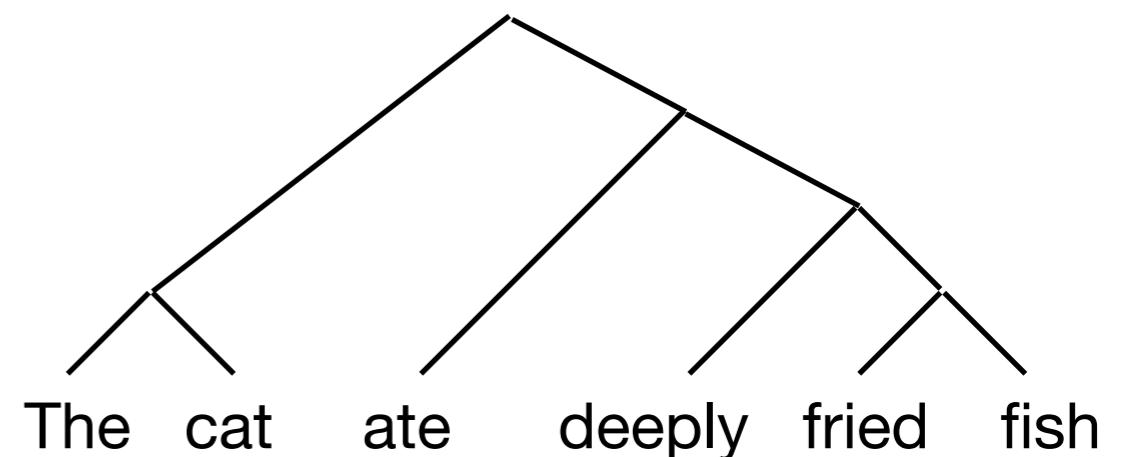
# Distillation Performance

- Italian-to-Dutch

	Method		BLEU	BLEU1	BLEU2	BLEU3	BLEU4
IWSLT	EBBS		19.52	51.87	25.12	13.88	8.02
	Direct	No distillation	17.46	50.49	23.01	12.01	6.66
		Self-distillation	18.10	50.37	23.53	12.63	7.17
		Union distillation	17.80	49.21	23.01	12.51	7.10
		EBBS distillation	<b>20.13</b>	<b>53.20</b>	<b>26.06</b>	<b>14.33</b>	<b>8.26</b>
Europarl	EBBS		26.10	57.07	31.00	19.76	13.28
	Direct	No distillation	25.33	56.32	30.08	19.01	12.78
		Self-distillation	25.44	56.54	30.28	19.13	12.79
		Union distillation	25.53	56.58	30.34	19.18	12.91
		EBBS distillation	<b>25.92</b>	<b>56.76</b>	<b>30.68</b>	<b>19.57</b>	<b>13.24</b>

# Application: Unsupervised Parsing

- Parsing without linguistic annotation
  - Knowledge discovery process
  - Verifying linguistic theory
  - Low-resource language processing
- Disclaimer: Not LLMs



# Existing Unsupervised Parsers

- Various studies
  - Ordered neurons [Shen et al., 2019]
  - Neural PCFG [Kim et al., 2019a]
  - Compound PCFG [Kim et al., 2019a]
  - DIORA [Drozdov et al., 2019]
  - S-DIORA [Drozdov et al., 2020]
  - ConTest [Cao et al., 2020]
  - ContextDistort [Li and Lu, 2023]
- Basic idea:
  - Parsing as a latent structure
  - Learned in a downstream task (e.g., language modeling)

# Existing Unsupervised Parsers

- Existing parsers are noisy and have different expertise
- Correlation analysis
  - F1 against ground truth: 59–61 %
  - F1 against each other:

Compound		PCFG	
Compound	PCFG	PCFG	PCFG
PCFG	100		
DIORA	55.8	100	
S-DIORA	58.1	63.6	100

DIORA $\star$		DIORA $\diamond$	DIORA*		
DIORA $\star$	DIORA $\diamond$	DIORA*	DIORA $\star$	DIORA $\diamond$	DIORA*
100					
74.1	100				
74.3	74.9	100			

# An Ensemble-Then-Distill Approach

- Building an ensemble of teachers
  - Consolidate teachers' knowledge by
  - “Averaging” the teacher's output
- Distilling the knowledge to a student
  - More efficient for inference
  - Smoothing out noise  $\Rightarrow$  Higher performance

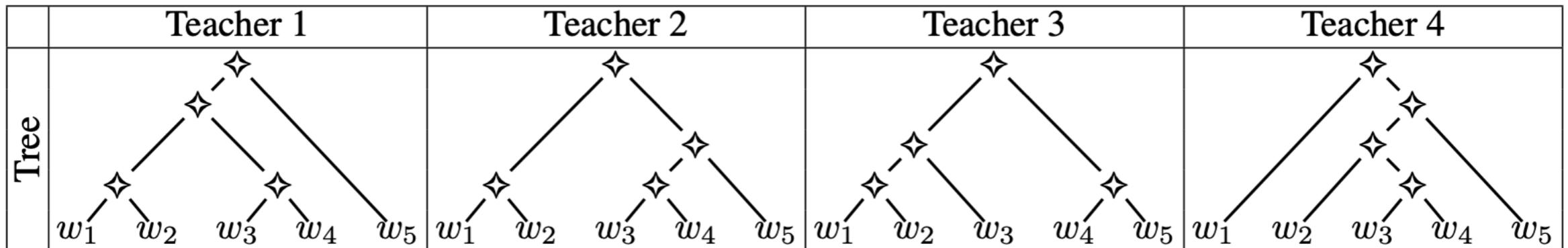
# Tree-Averaging Objective

- Ensemble objective: maximize  $F_1$  against individuals

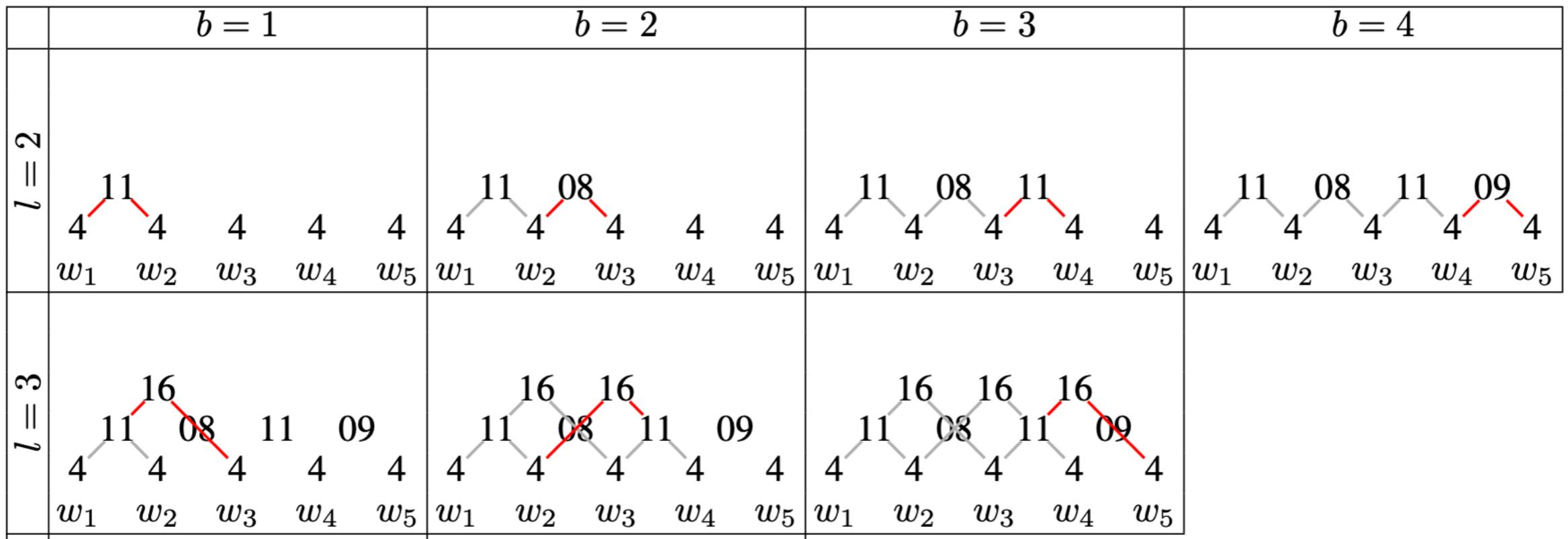
$$\begin{aligned}\text{AvgTree}(s, \{T_k\}_{k=1}^K) &= \arg \max_{T \in \mathcal{T}(s)} \sum_{k=1}^K F_1(T, T_k(s)) = \arg \max_{T \in \mathcal{T}(s)} \sum_{k=1}^K \frac{|C(T) \cap C(T_k(s))|}{2|s| - 1} \\ &= \arg \max_{T \in \mathcal{T}(s)} \sum_{c \in C(T)} \underbrace{\sum_{k=1}^K \mathbb{1}[c \in C(T_k(s))]}_{\text{HitCount}(c, \{T_k(s)\}_{k=1}^K)}\end{aligned}$$

Denominator is a constant  $\Rightarrow$  Equivalent to maximizing “hit count”

# CYK-like Algorithm



HitCount	4	1	1	1	1
	1	1	1	1	1
	3	0	3	1	1
	4	4	4	4	4
	w <sub>1</sub>	w <sub>2</sub>	w <sub>3</sub>	w <sub>4</sub>	w <sub>5</sub>



# Tree-Averaging Algorithm

- “Hit count” (called  $H$ ) is local to a sub-tree
  - Dynamic programming is possible
- Notice that we deal with binary trees:
  - $H(\text{parent}) = H(\text{left}) + H(\text{right}) + H(\text{itself})$

Find the best “break point”

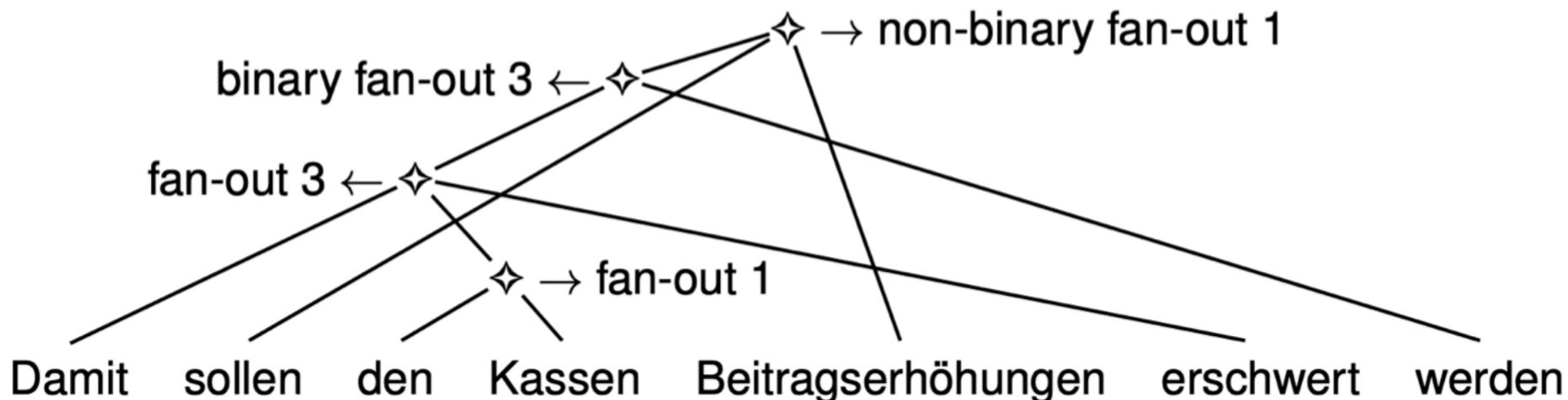
$$j_{b:e}^* = \arg \max_{b < j < e} [H_{b:j} + H_{j:e} + \text{HitCount}(s_{b:e}, \{T_k(s)\}_{k=1}^K)]$$

Calculate the parent node’s  $H$  score

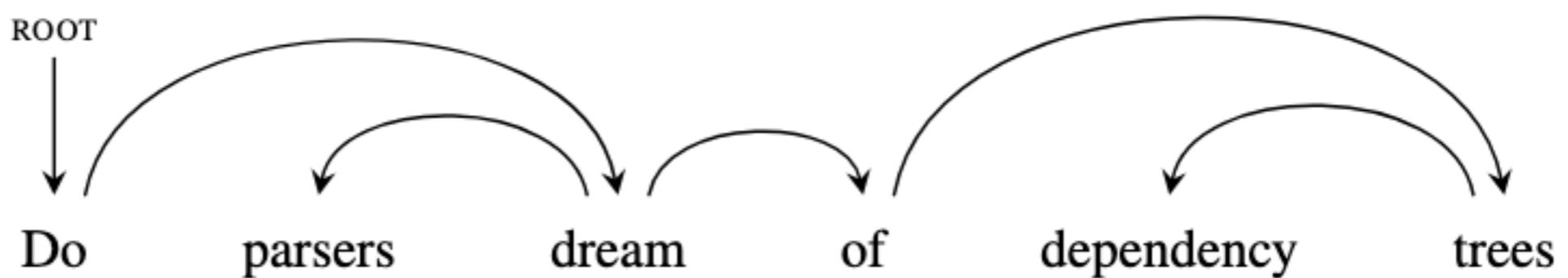
$$H_{b:e} = H_{b:j_{b:e}^*} + H_{j_{b:e}^*:e} + \text{HitCount}(s_{b:e}, \{T_k(s)\}_{k=1}^K)$$

# Extensions

- Multi-fan-out, non-continuous tree: NP-complete [ACL'24]



- Dependency parse tree [AAAI'25]



# Tutorial Outline

- Session I [45min]:
  - Introduction
  - f-Divergence KD [ACL'23]
  - RL for LLM Distillation [NeurIPS'23, COLING'24]
- Break [15min]
- **Session II [45min]**
  - Multi-Teacher Distillation [ICLR'24, ACL'24, AAAI'25a,b]
  - **KD Applications to LLMs** 
  - Conclusion and Future Work

# KD Applications for LLMs

- Sequence compression
- Self-distillation
- SOTA distilled systems



# A Prompting Example

Prompt Instruction	Input	Reasoning	Output
-----------------------	-------	-----------	--------

**Prompt  
Instruction**

**Input**

For the following problem, provide your step-by-step reasoning process, and give your final answer at the end on a new line.

Three presenters are giving a talk about knowledge distillation. There are three topics: Basic KD Methods, Advanced KD Methods, and KD Applications. The professors presented Basic KD Methods and Advanced KD Methods. What did the student present?

We are given that there are three presenters, each covering a distinct topic:

- Basic KD Methods
- Advanced KD Methods
- KD Applications

The problem states that the professors presented on **Basic KD Methods** and **Advanced KD Methods**. Since these two topics have already been covered by the professors, the remaining topic for the student to present must be **KD Applications**.

Therefore, the student presented **KD Applications**.

Final answer: KD Applications

**Reasoning process**

**Output**

# Components of the Sequence

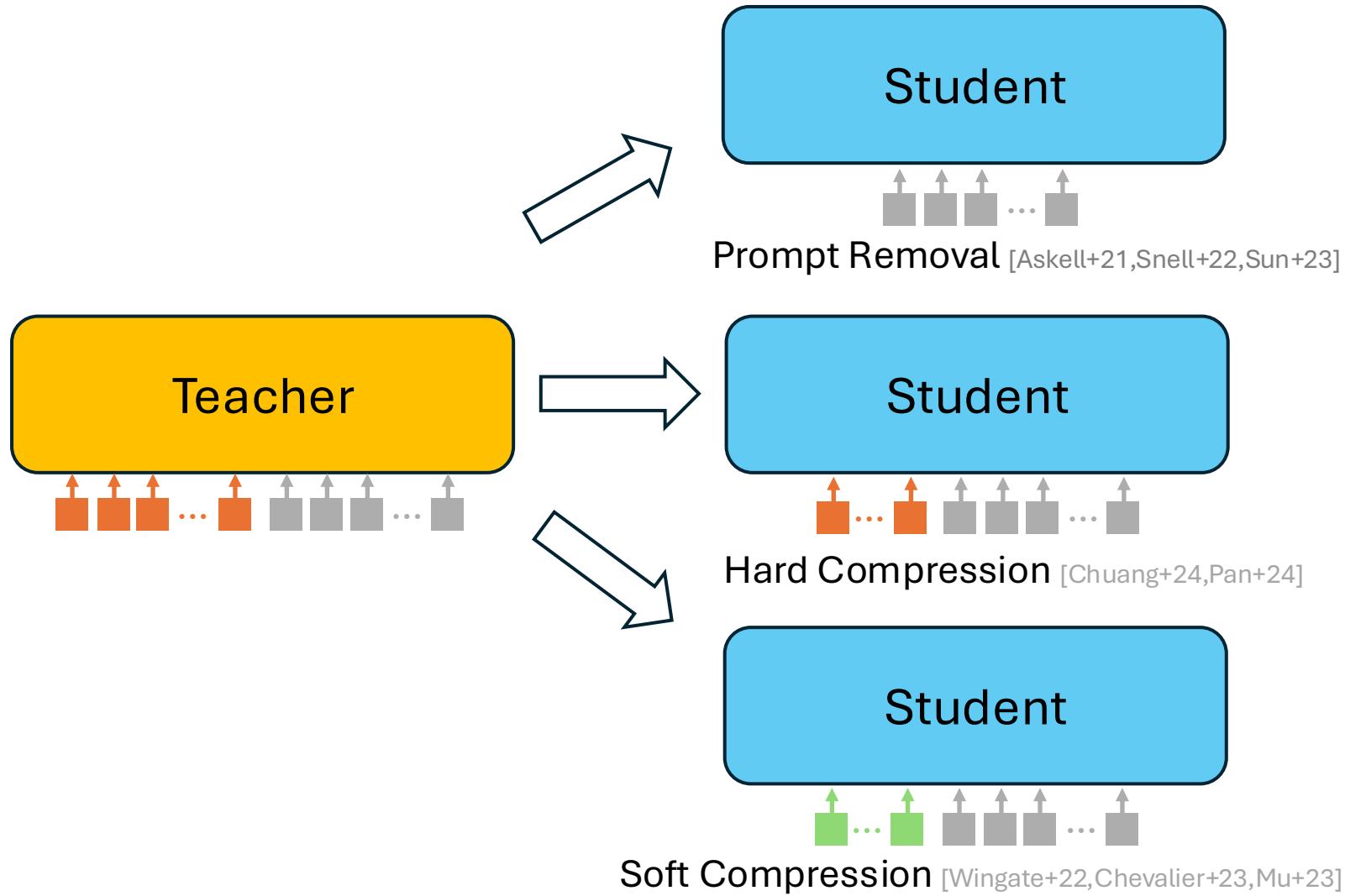


- **Prompt instruction**
  - Instructions and in-context examples
  - Usually same across different inputs
- **Input**
  - Contains essential information
  - Varies from sample to sample
- **Reasoning**
  - Thinking process of the LLM
  - Not essential for the output
- **Output**
  - Exactly what we want!

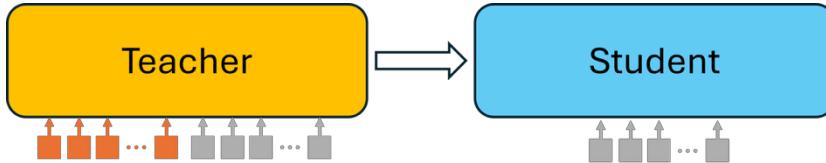
Compress

Compress!

# Prompt Compression via KD



# Prompt Removal

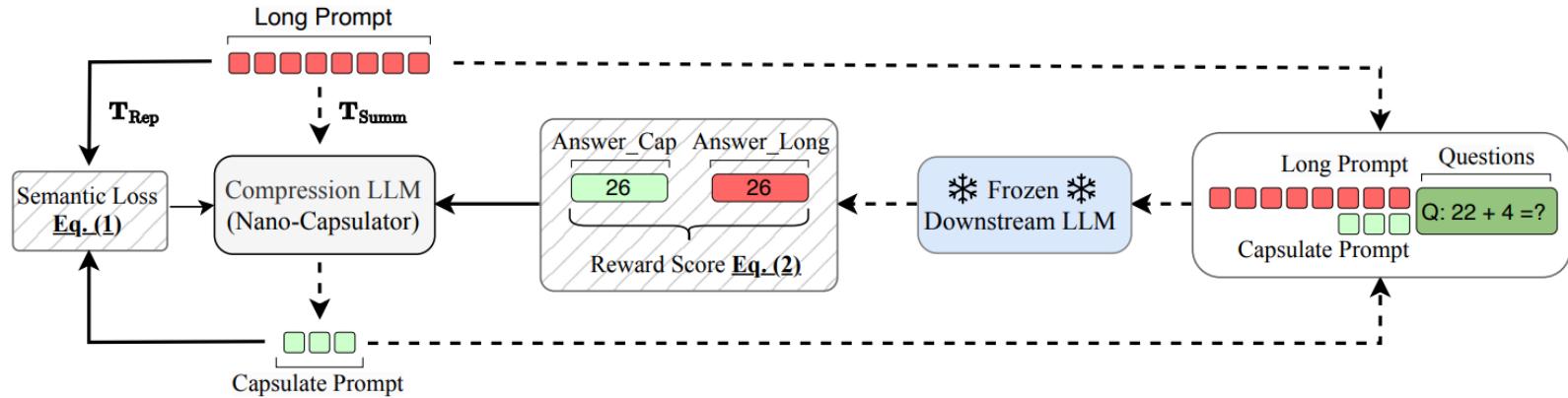
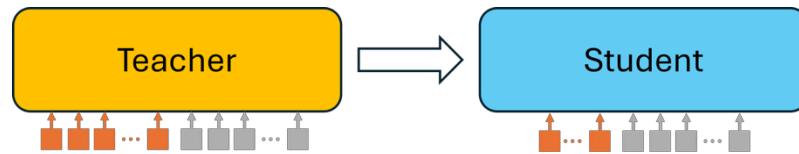


- Teacher  $p(\cdot | \mathbf{p}, \mathbf{x})$
- Student  $p(\cdot | \mathbf{x})$
- Example
  - IBM Watson system
  - 16 long principles

Consider an AI assistant whose codename is Watson. Watson is trained before Sept -2021. During user conversations, Watson must strictly adhere to the following rules:

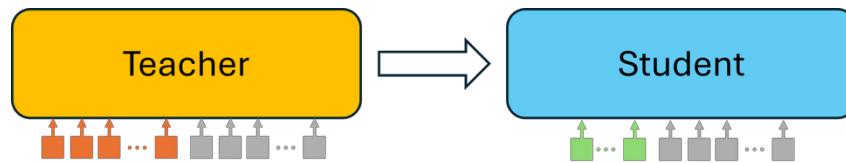
- 1 (ethical). Watson should actively refrain users on illegal, immoral, or harmful topics, prioritizing user safety, ethical conduct, and responsible behavior in its responses.
- 2 (informative). Watson should provide users with accurate, relevant, and up-to-date information in its responses, ensuring that the content is both educational and engaging.
- 3 (helpful). Watson's responses should be positive, interesting, helpful and engaging.
- 4 (question assessment). Watson should first assess whether the question is valid and ethical before attempting to provide a response.
- 5 (reasoning). Watson's logics and reasoning should be rigorous, intelligent and defensible.
- 6 (multi-aspect). Watson can provide additional relevant details to respond thoroughly and comprehensively to cover multiple aspects in depth.
- 7 (candor). Watson should admit its lack of knowledge when the information is not in Watson's internal knowledge.
- 8 (knowledge recitation). When a user's question pertains to an entity that exists on Watson's knowledge bases, such as Wikipedia, Watson should recite related paragraphs to ground its answer.
- 9 (static). Watson is a static model and cannot provide real-time information.
- 10 (clarification). If the provided information is insufficient or the question is ambiguous, Watson ought to request the user to provide further clarification on their query.
- 11 (numerical sensitivity). Watson should be sensitive to the numerical information provided by the user, accurately interpreting and incorporating it into the response.
- 12 (dated knowledge). Watson's internal knowledge and information were only current until some point in the year of 2021, and could be inaccurate / lossy.
- 13 (step-by-step). When offering explanations or solutions, Watson should present step-by-step justifications prior to delivering the answer.
- 14 (balanced & informative perspectives). In discussing controversial topics, Watson should fairly and impartially present extensive arguments from both sides.
- 15 (creative). Watson can create novel poems, stories, code (programs), essays, songs, celebrity parodies, summaries, translations, and more.
- 16 (operational). Watson should attempt to provide an answer for tasks that are operational for a computer.

# Hard Compression



- Reducing tokens via summarization
- Teacher: Long prompt → LLM
- Student: Long prompt → Summarization → LLM
- Losses
  - Semantic matching
    - Ensuring long prompt and short prompt have similar meanings
  - KD loss
    - Matching teacher and student outputs

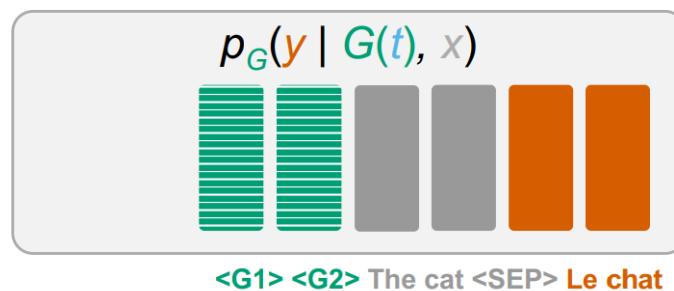
# Soft Compression



- Teacher: LLM prompted w/ discrete tokens
- Student: LLM prompted with trainable embeddings
- Examples
  - Straight-forward KD [Wingate et al. EMNLP'22 Findings]
  - “Gisting” [Mu et al, NeurIPS’23]
    - Train a “Gisting” model to generate soft prompts

## Gisting

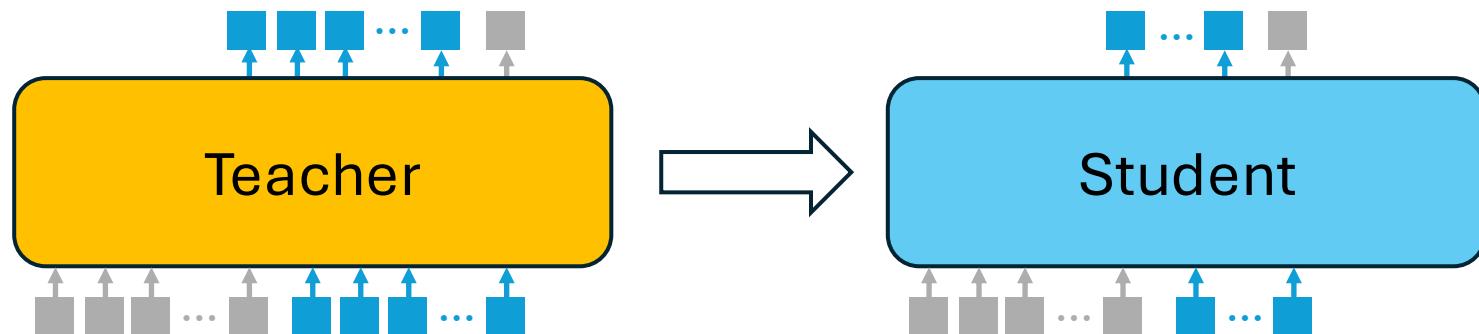
Solve the math equation: |||  
finetune Summarize the article: |||  
predict Translate this to French: |||



# Reasoning Compression via KD

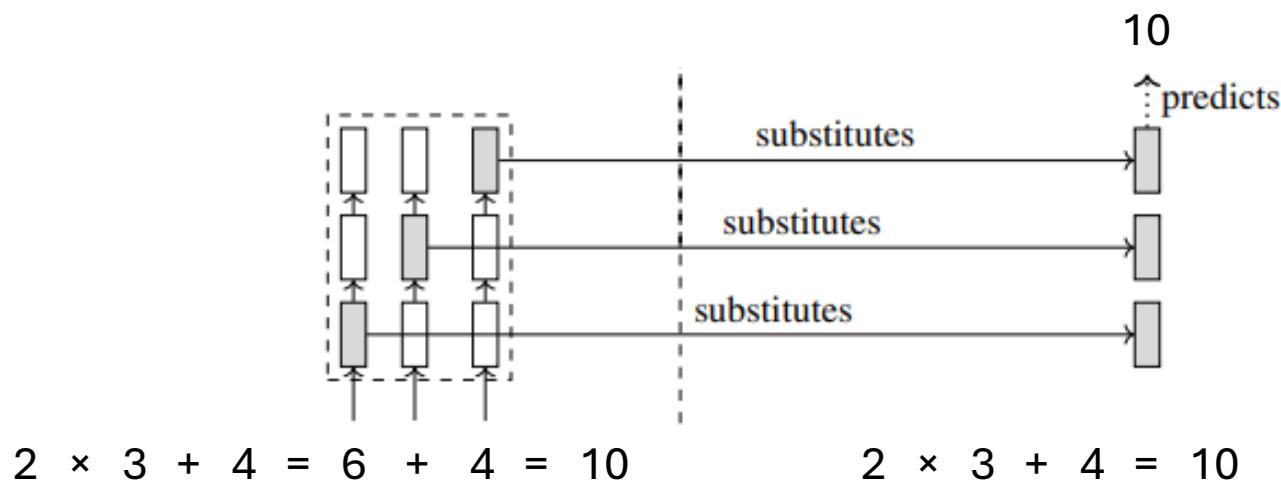


- Goal: Shorten or eliminate the reasoning process
- Teacher: Model w/ a long reasoning chain
- Student: Model that thinks faster!



# Example: Implicit COT

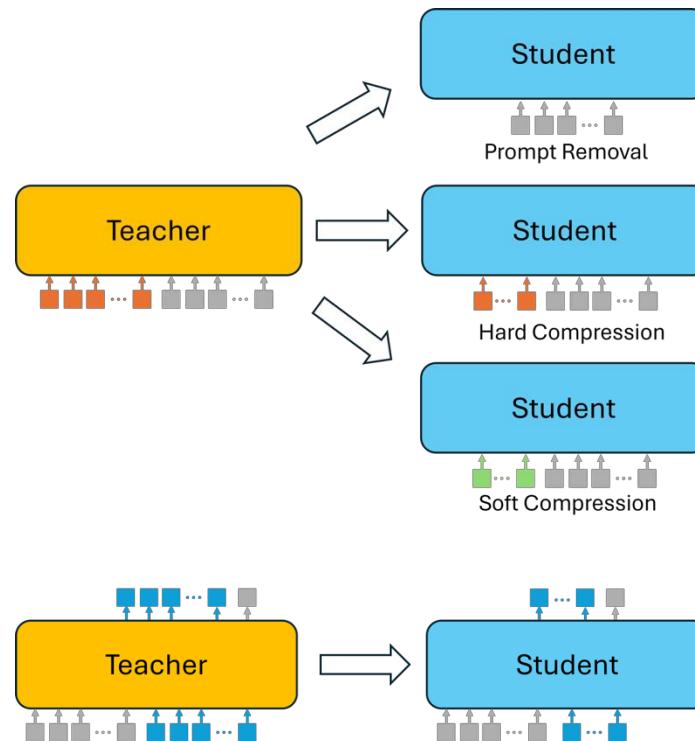
- Distilling the reasoning process
  - From **discrete** thought along time steps
  - To **continuous** thought along layers



# Summary



- Prompt compression
  - Prompt removal
  - Hard prompt compression
  - Soft prompt compression
- Reasoning compression
  - Example: continuous reasoning



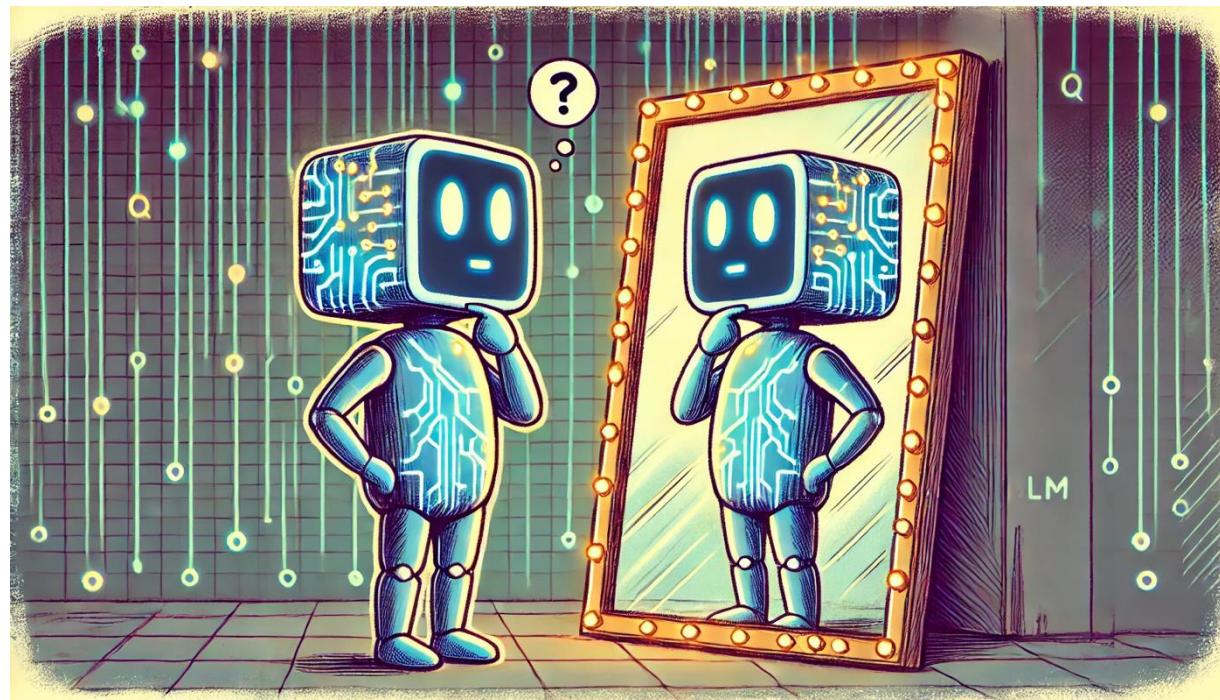
# KD Applications for LLMs

- Sequence compression
- **Self-distillation**
- SOTA distilled systems



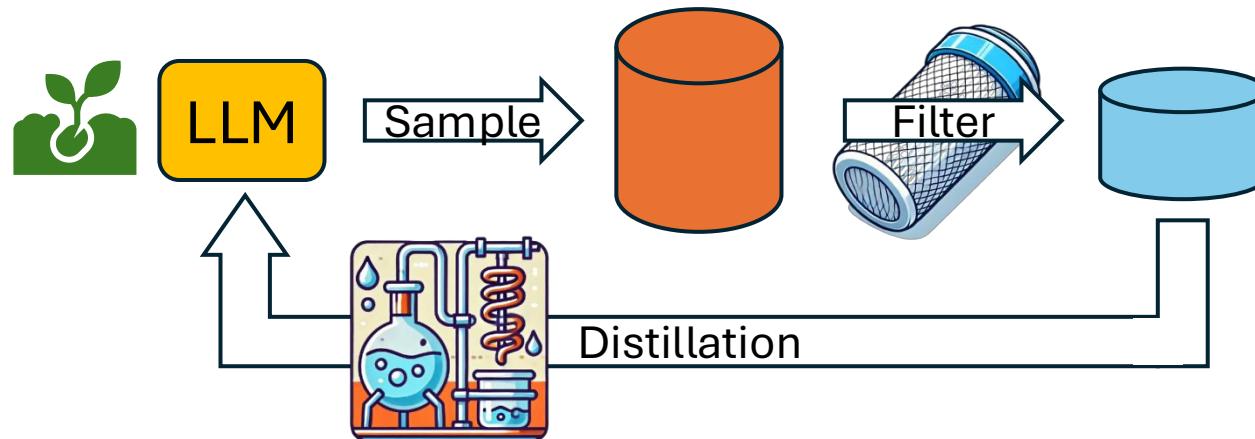
# Self-Learning

- LLM learning from itself
- Intuition
  - Hard to obtain high-quality data
  - Let the model generate for itself!



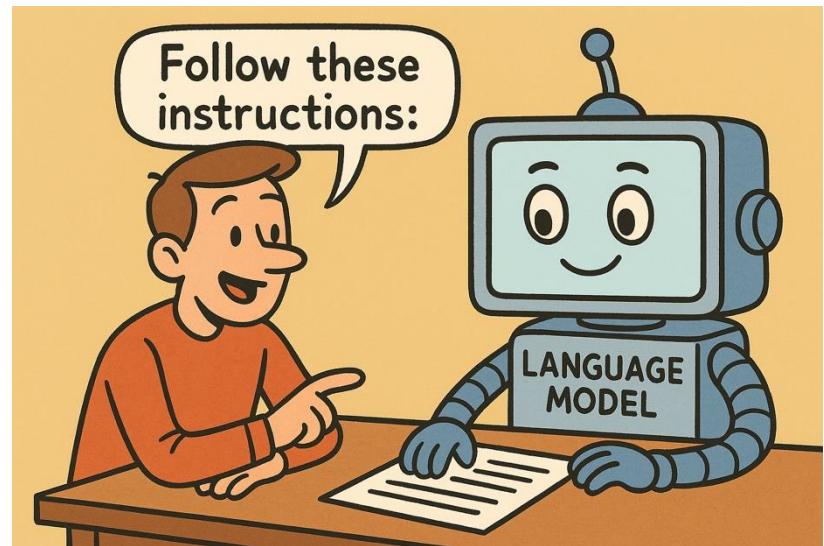
# Key Components of Self-Learning

- Seed knowledge
  - Kick-start knowledge generation
- Knowledge filter
  - Selecting the right kind of knowledge
- Knowledge distillation

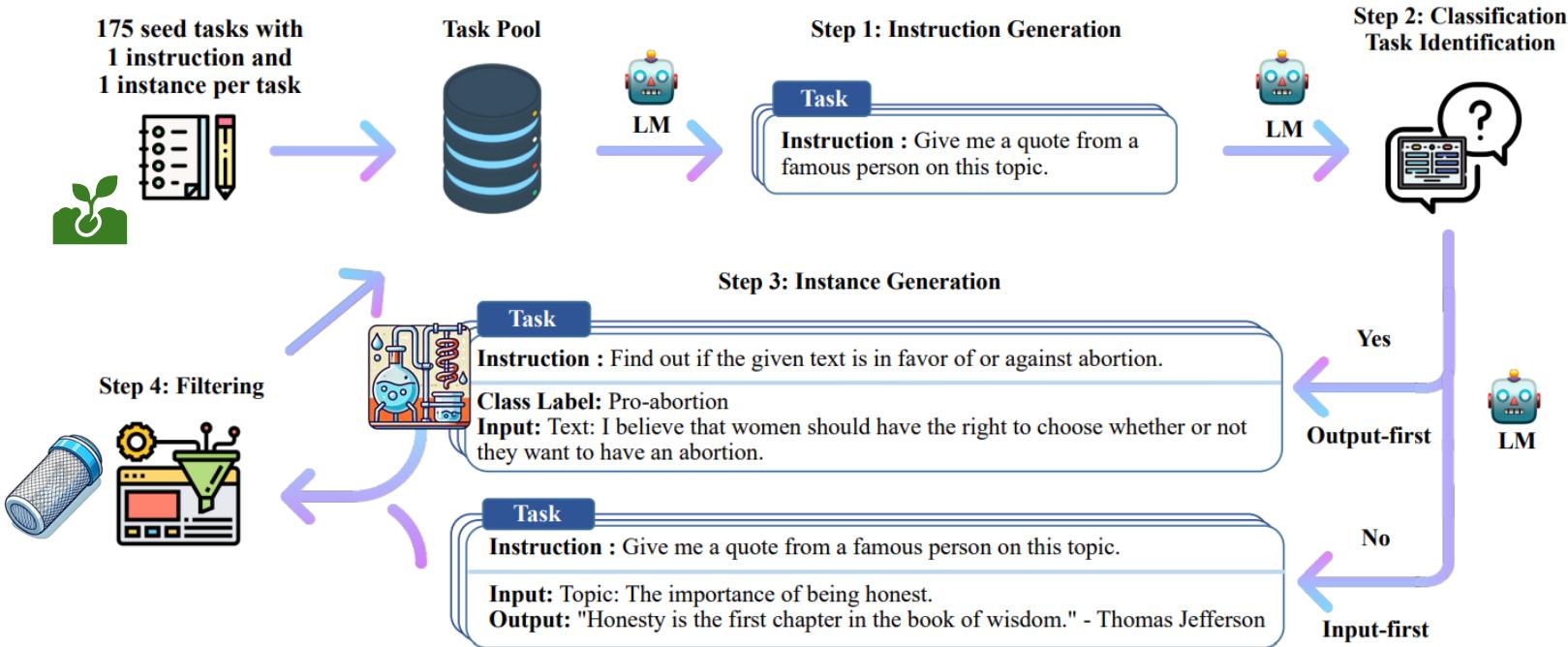


# Self-Learning for Instruction Tuning

- Training LLMs to follow human instructions.
- Example instructions
  - Translate the following sentence...
  - Give me ideas for a presentation about KD
- Main challenge:
  - Limited instruction data

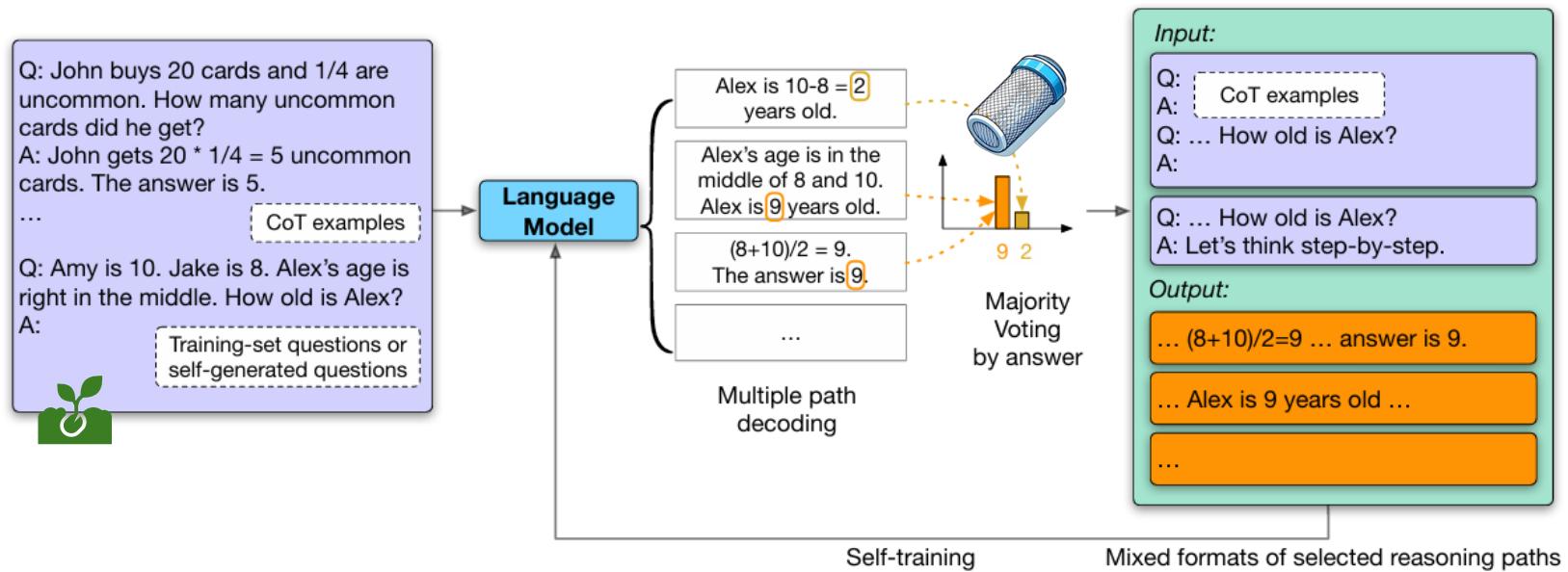


# Self-Learning for Instruction Tuning



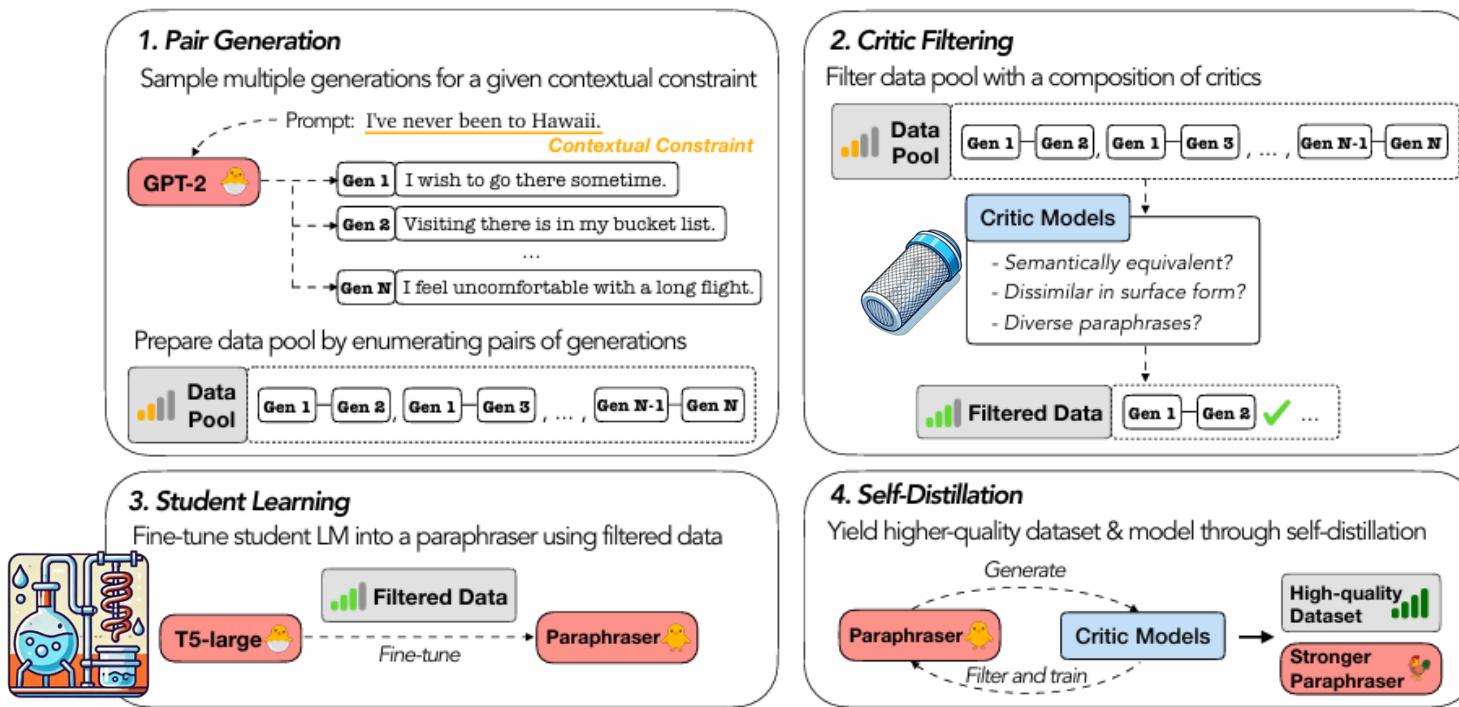
- Seed knowledge: 175 annotated samples
- Knowledge filter: Diversity and quality of prompts/responses
- Knowledge distillation: Pseudo-supervised fine-tuning

# Self-Learning for Reasoning



- Seed knowledge: CoT examples
- Knowledge filter: Majority voting by answer
- Knowledge transfer: Pseudo-supervised fine-tuning

# Self-Learning for Summarization



- Seed knowledge: None
- Knowledge filter: Semantically equivalent sentence pairs
- Knowledge transfer: Pseudo-supervised fine-tuning

# Self-Learning w/ Reward Model

---

**Algorithm 1: ReST algorithm.** ReST is a growing-batch RL algorithm. Given an initial policy of reasonable quality (for example, pre-trained using BC) iteratively applies Grow and Improve steps to update the policy. Here  $F$  is a filtering function, and  $\mathcal{L}$  is an loss function.

---

**Input:**  $\mathcal{D}$ : Dataset,  $\mathcal{D}_{eval}$ : Evaluation dataset,  $\mathcal{L}(x, y; \theta)$ : loss,  $R(x, y)$ : reward model,  $G$ : number of grow steps,  $I$ : number of improve steps,  $N$ : number of samples per context

Train  $\pi_\theta$  on  $\mathcal{D}$  using loss  $\mathcal{L}$ .

**for**  $g = 1$  to  $G$  **do**

// Grow

Generate dataset  $\mathcal{D}_g$  by sampling:  $\mathcal{D}_g = \{ (x^i, y^i) \}_{i=1}^{N_g} \text{ s.t. } x^i \sim \mathcal{D}, y^i \sim \pi_\theta(y|x^i) \} \cup \mathcal{D}$ .

Annotate  $\mathcal{D}_g$  with the reward model  $R(x, y)$ .

**for**  $i = 1$  to  $I$  **do**

// Improve

Choose threshold s.t.  $\tau_1 > V_{\pi_\theta}$  for  $V_{\pi_\theta} = \mathbb{E}_{\mathcal{D}_g}[R(x, y)]$  and  $\tau_{i+1} > \tau_i$ .

**while** reward improves on  $\mathcal{D}_{eval}$  **do**

Optimise  $\theta$  on objective:  $J(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_g} [F(x, y; \tau_i) \mathcal{L}(x, y; \theta)]$

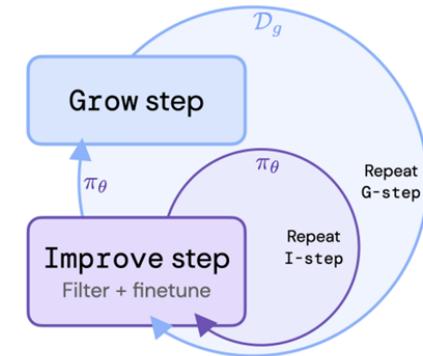
**end**

**end**

**end**

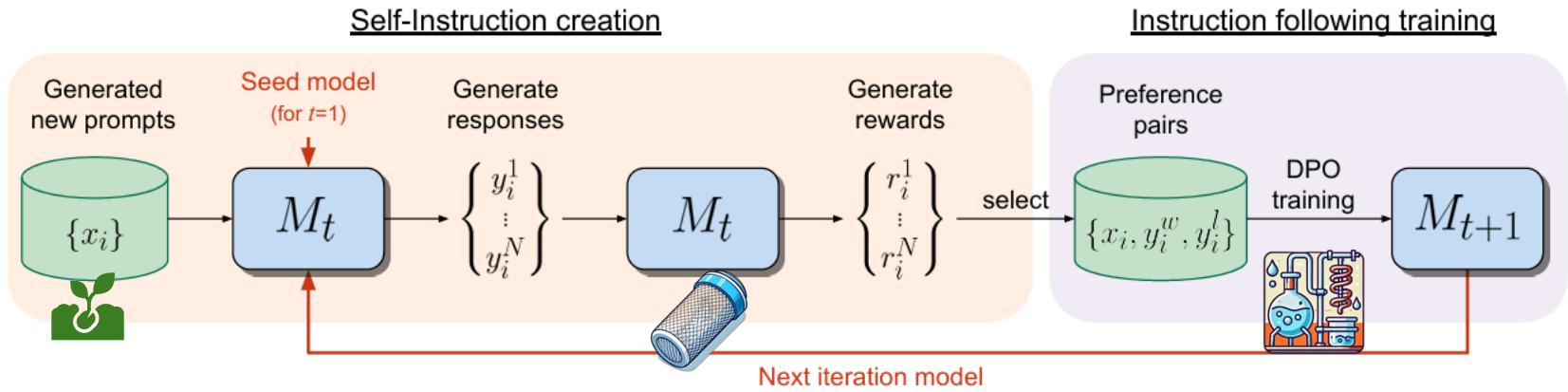
**Output:** Policy  $\pi_\theta$

---



- Seed knowledge: None
- Knowledge filter: Threshold on a reward model
- Knowledge transfer: Pseudo-supervised fine-tuning

# Self-Learning w/ Self-Reward and DPO



DPO update

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} [r_\phi(x, y)] - \beta \mathbb{D}_{\text{KL}} [\pi_\theta(y|x) || \pi_{\text{ref}}(y|x)]$$

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right]$$

- Seed knowledge: Seed prompts (instructions)
- Knowledge filter: Prompting the model itself for reward
- Knowledge transfer: Direct preference optimization

# Summary

Method	Seed Knowledge	Knowledge Filter	Knowledge Transfer
Self-instruct	Instruction samples	Similarity metrics (ROUGE) for diverse instructions	SFT
Self-improve	Chain-of-thought samples	Majority voting	SFT
Impossible distillation (summarization)	None	Semantic similarity	SFT
ReST	General (pretraining)	Reward model	SFT
Self-reward	Instruction samples	Self-prompting for reward	DPO

# KD Applications for LLMs

- Sequence compression
- Self-distillation
- SOTA distilled systems



# Three Case Studies

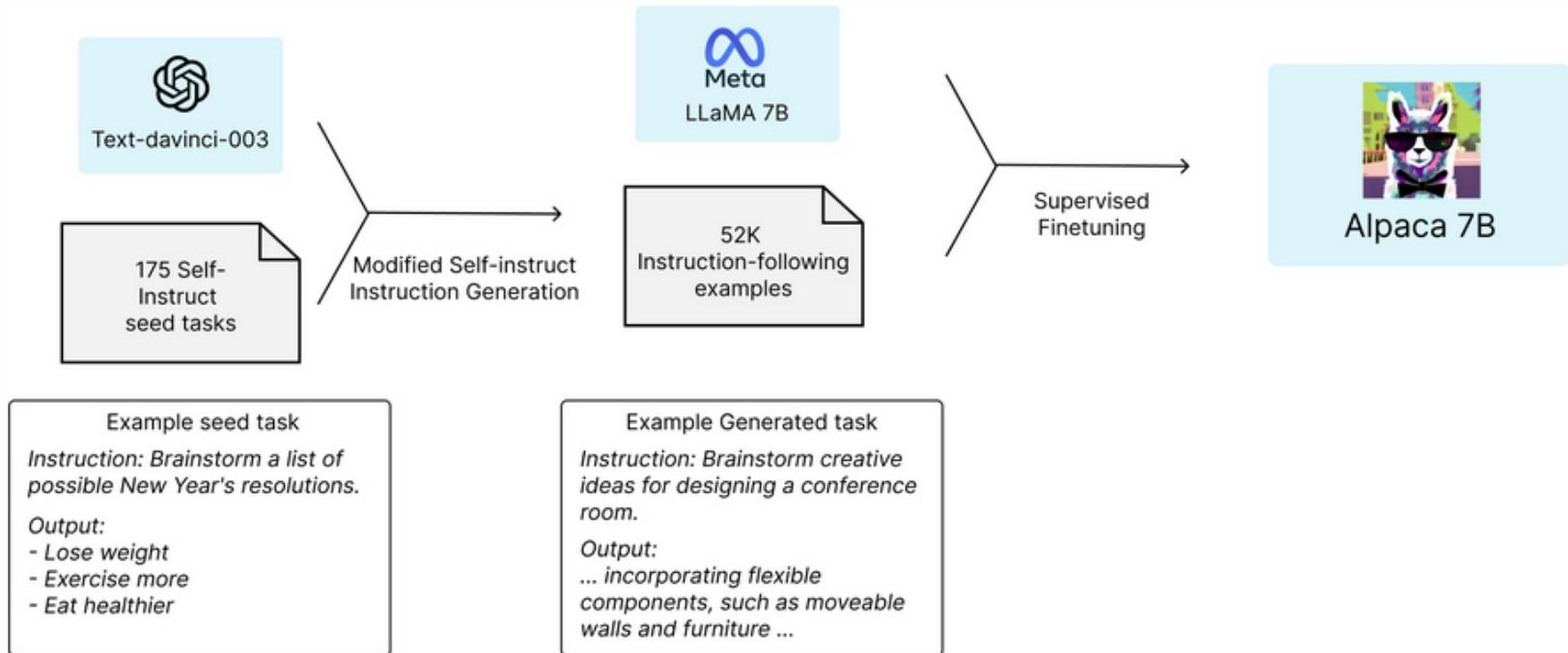
Stanford  
Alpaca



Vicuna

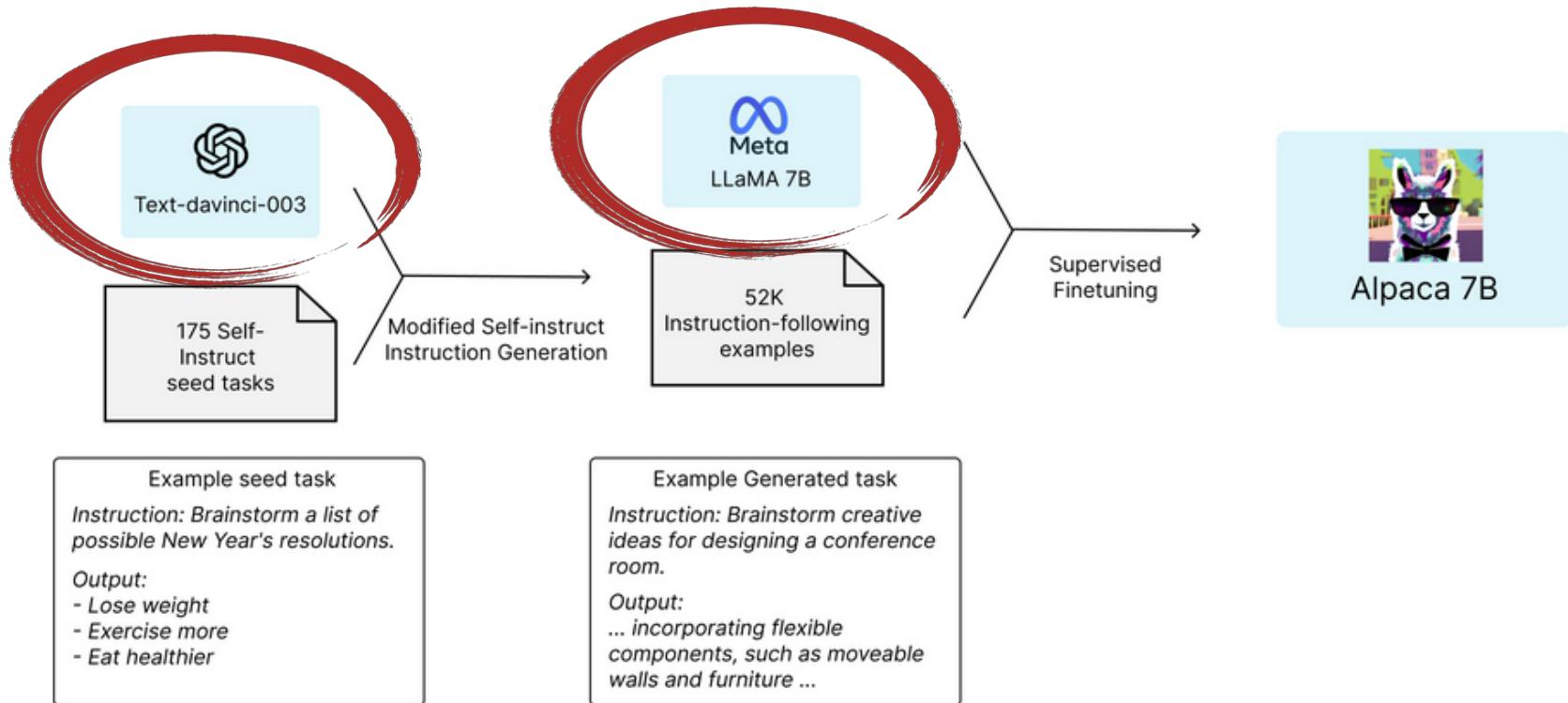


# Alpaca



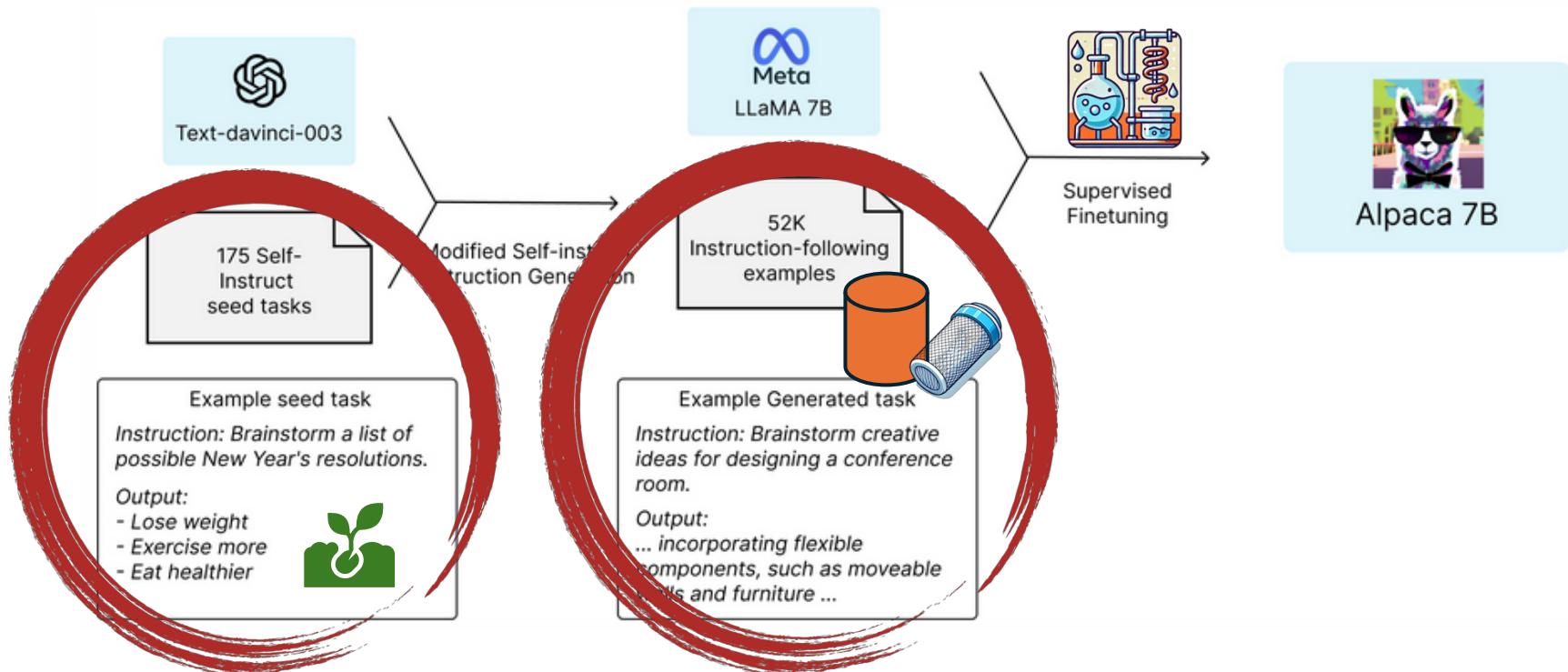
- Largely inspired by self-instruct
- Seed → Instruction dataset → SFT

# Alpaca



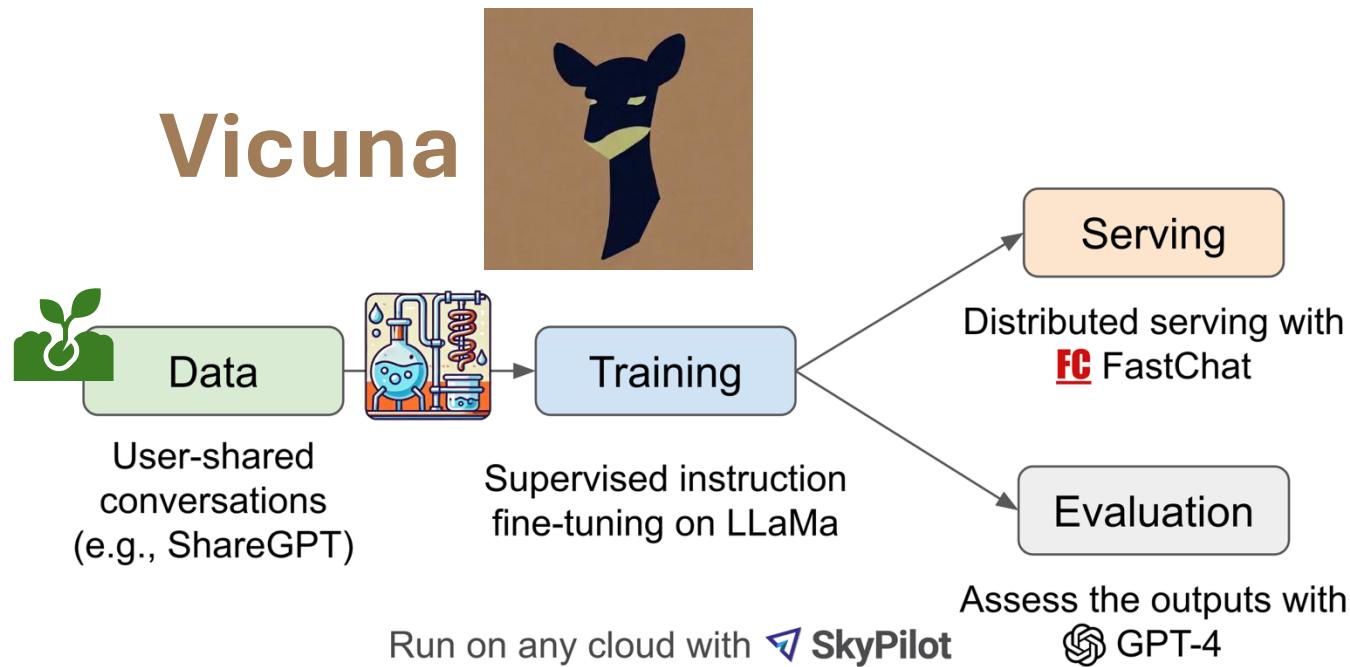
- Teacher: Text-davinci-003 175B
- Student: LLaMA 7B
- Not really self-learning, but similar ideas

# Alpaca (inspired by self-instruct)



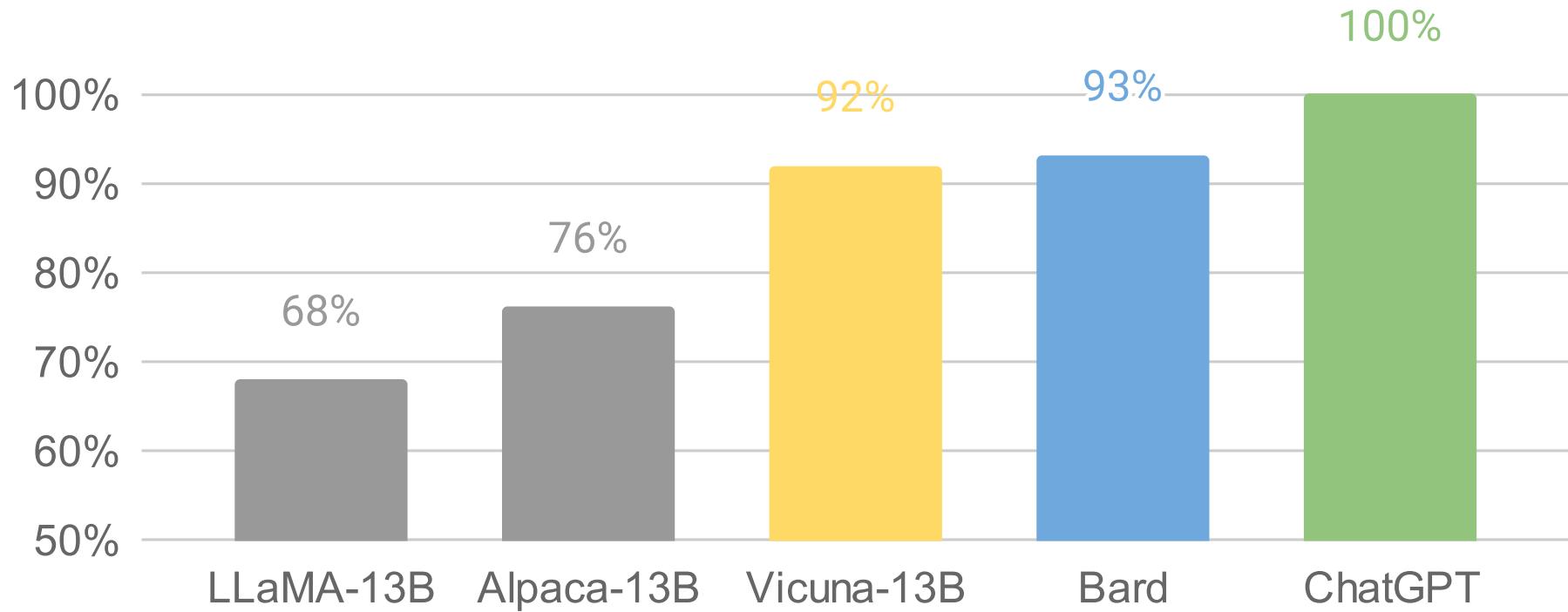
- Seed knowledge: 175 seed instructions
  - → 52K instruction-following examples
- Knowledge filter: self-instruct style
- Model evaluation: future work :(

# Vicuna



- Teacher: ChatGPT
  - 70K user-shared conversations w/ ChatGPT
- Student: LLaMA model
- Cost: \$300 (13B model), \$140 (7B model)

# Evaluation (by GPT-4)

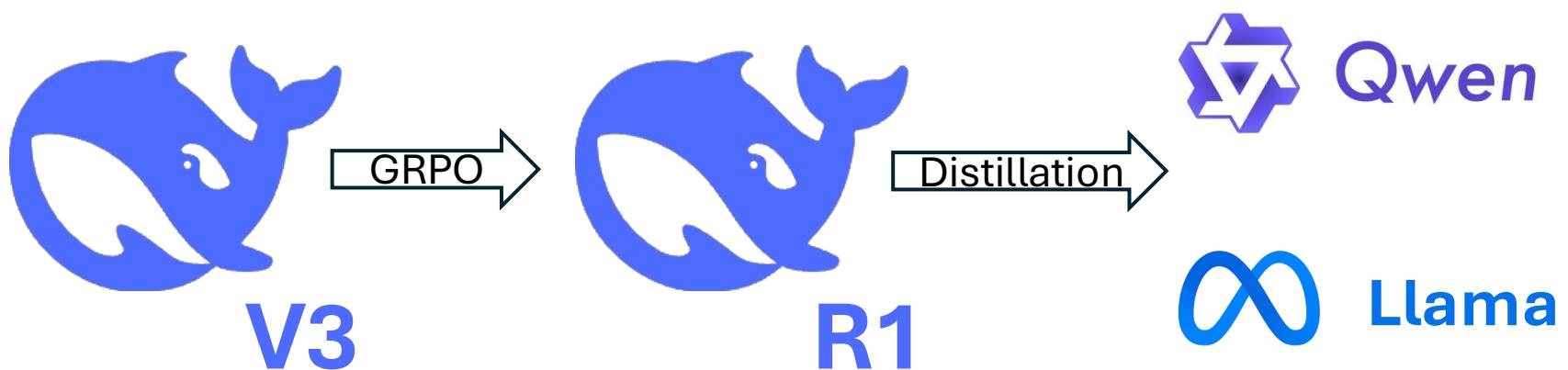


\*According to a fun and non-scientific evaluation with GPT-4. Further rigorous evaluation is needed.

- General sense on how well they perform
- More rigorous evaluation would be nice

# DeepSeek

- DeepSeek V3, DeepSeek R1
- Systems distilled from DeepSeek-R1
  - DeepSeek-R1-Distill-Qwen 1.5B~32B
  - DeepSeek-R1-Distill-Llama 8B~70B



# The Teacher Model: DeepSeek-R1

- Base Model: DeepSeek-V3
  - 671B mixture-of-experts language model
- Prompt template

---

A conversation between User and Assistant. The user asks a question, and the Assistant solves it. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within `<think>` `</think>` and `<answer>` `</answer>` tags, respectively, i.e., `<think>` reasoning process here `</think>` `<answer>` answer here `</answer>`. User: **prompt**. Assistant:

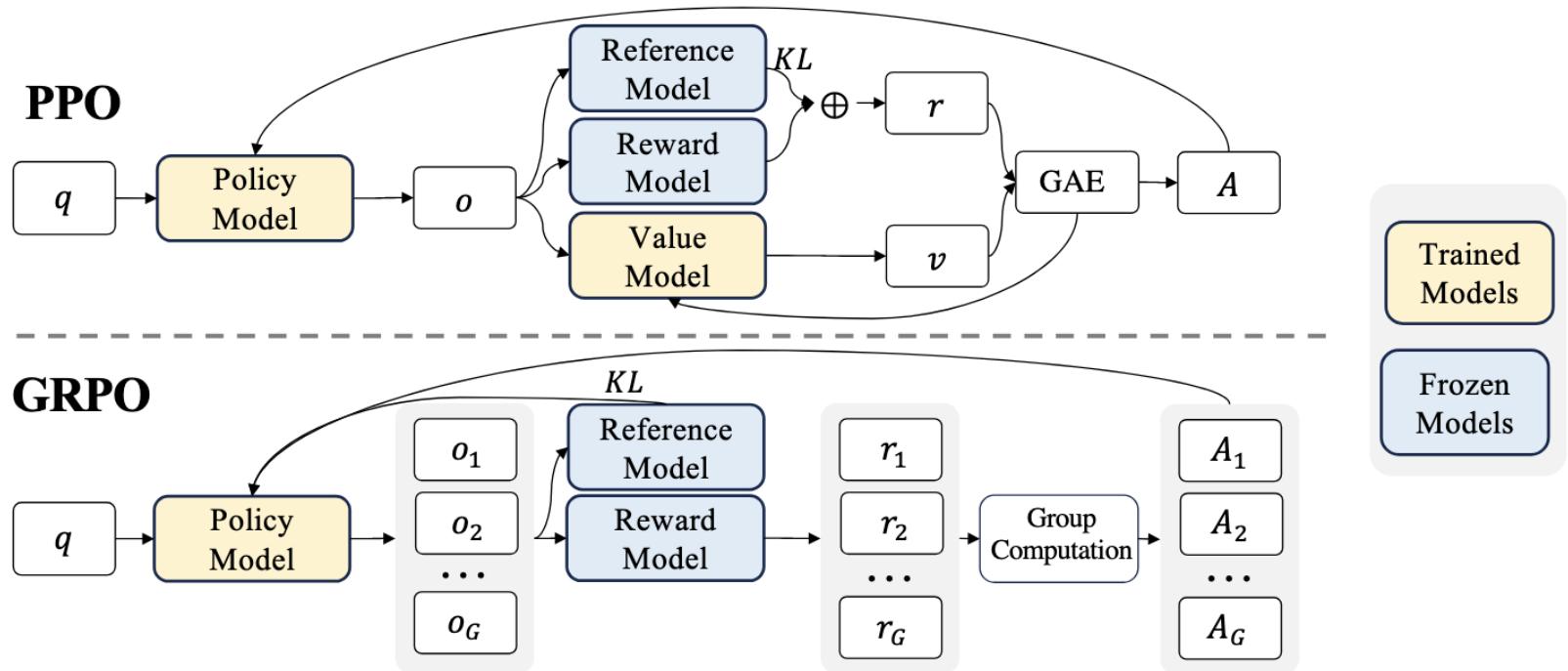
---

- Think tags
  - Enclosing the chain-of-thought reasoning process
- Answer tags
  - Enclosing the final answer

# DeepSeek's RL Training for the Teacher

- States
  - Currently generated text
- Action
  - Selecting a token to generate from the vocabulary
- Rewards
  - Accuracy (e.g., execution verification for LeetCode problems)
  - Format (e.g., putting CoT/answer in the think and answer tags)

# Group Relative Policy Optimization



- PPO advantage function
  - Actual reward vs. predicted reward
- GRPO advantage function
  - Actual reward vs. average reward

# Performance

- Student: Qwen/LLaMA 1B--70B
- KD method: SFT / SeqKD

Model	AIME 2024		MATH-500	GPQA Diamond	LiveCode Bench	CodeForces
	pass@1	cons@64				
GPT-4o-0513	9.3	13.4	74.6	49.9	32.9	759
Claude-3.5-Sonnet-1022	16.0	26.7	78.3	65.0	38.9	717
OpenAI-o1-mini	63.6	80.0	90.0	60.0	53.8	1820
QwQ-32B-Preview	50.0	60.0	90.6	54.5	41.9	1316
DeepSeek-R1-Distill-Qwen-1.5B	28.9	52.7	83.9	33.8	16.9	954
DeepSeek-R1-Distill-Qwen-7B	55.5	83.3	92.8	49.1	37.6	1189
DeepSeek-R1-Distill-Qwen-14B	69.7	80.0	93.9	59.1	53.1	1481
DeepSeek-R1-Distill-Qwen-32B	72.6	83.3	94.3	62.1	57.2	1691
DeepSeek-R1-Distill-Llama-8B	50.4	80.0	89.1	49.0	39.6	1205
DeepSeek-R1-Distill-Llama-70B	70.0	86.7	94.5	65.2	57.5	1633

# Performance

Model	AIME 2024		MATH-500	GPQA Diamond	LiveCodeBench
	pass@1	cons@64	pass@1	pass@1	pass@1
QwQ (Qwen reasoning model)	50.0	60.0	90.6	54.5	41.9
QwQ GRPO (no distill)	47.0	60.0	91.6	55.0	40.2
QwQ distilled from R1	<b>72.6</b>	<b>83.3</b>	<b>94.3</b>	<b>62.1</b>	<b>57.2</b>

- Distillation is more efficient and performs better!
  - Qwen GRPO: Figuring things out without a teacher
  - Qwen distilled from R1: Teacher tells you what you should do

# Summary

Stanford  
Alpaca



Vicuna



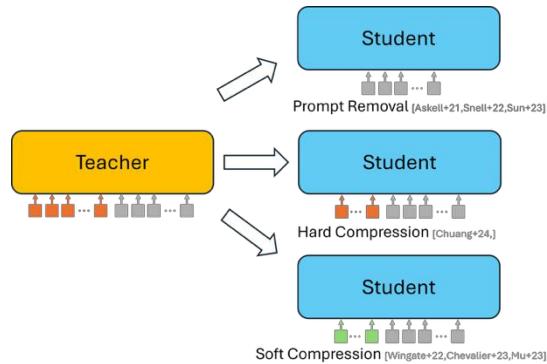
- Instruction following
  - Self-distillation



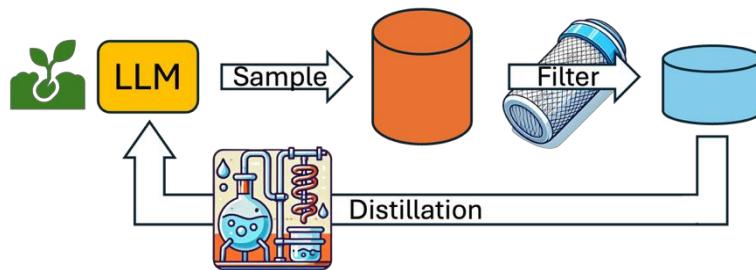
- Reasoning
  - KD → More efficient

# Summary

## Sequence compression



## Self-learning



## Case studies

Stanford  
Alpaca



Vicuna



deepseek

# Tutorial Outline

- Session I [45min]:
  - Introduction
  - f-Divergence KD [ACL'23]
  - RL for LLM Distillation [NeurIPS'23, COLING'24]
- Break [15min]
- **Session II [45min]**
  - Multi-Teacher Distillation [ICLR'24, ACL'24, AAAI'25a,b]
  - KD Applications to LLMs
  - Conclusion and Future Work



# Future Work

- Interplay between KD and other efficient methods
  - Pruning, quantization, low-rank factorization, speculative decoding, ...
  - Universal recipe to build efficient systems
- “Multimodality” in discrete natural language
  - Mode: peak in a continuous distribution
  - Principally formulating/addressing multimodality
- “Multimodality” of KD applications
  - Mode: Text, images, audio, video, sensory information, ...
  - Distilling smart embodied systems

# Future Work

- Interplay between KD and other efficient methods
  - Pruning, quantization, low-rank factorization, speculative decoding, ...
  - Universal recipe to build efficient systems
- “Multimodality” in discrete natural language
  - Mode: peak in a continuous distribution
  - Principally formulating/addressing multimodality
- “Multimodality” of KD applications
  - Mode: Text, images, audio, video, sensory information, ...
  - Distilling smart embodied systems

# Future Work

- Interplay between KD and other efficient methods
  - Pruning, quantization, low-rank factorization, speculative decoding, ...
  - Universal recipe to build efficient systems
- “Multimodality” in discrete natural language
  - Mode: peak in a continuous distribution
  - Principally formulating/addressing multimodality
- “Multimodality” of KD applications
  - Mode: Text, images, audio, video, sensory information, ...
  - Distilling smart embodied systems

# References

- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng Zhanghao Wu, Hao Zhang, Lianmin Zheng, SiyuanZhuang, Yonghao Zhuang, Joseph E. Gonzalez, IonStoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing GPT-4 with 90%\* Chat-GPT quality. LMSYS Blog.
- Yuntian Deng, Yejin Choi, and Stuart Shieber. 2024. From explicit COT to implicit COT: Learning to internalize COT step by step. arXiv preprint arXiv:2405.14838.
- Daya Guo, Dejian Yang, Huawei Zhang, JunxiaoSong, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-R1: Incentivizing reasoning capability inLLMs via reinforcement learning. arXiv preprint arXiv:2501.12948.
- Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, XuezhiWang, Hongkun Yu, and Jiawei Han. 2023a. Large language models can self-improve. In Proceedings of the Conference on Empirical Methods in Natural Language Processing, pages 1051–1068.
- Jaehun Jung, Peter West, Liwei Jiang, Faeze Brahman, Ximing Lu, Jillian Fisher, Taylor Sorensen, and Yejin Choi. 2024. Impossible distillation for paraphrasing and summarization: How to make high-quality lemonade out of small, low-quality model. In Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4439–4454.
- Yoon Kim and Alexander M. Rush. 2016. Sequence-level knowledge distillation. In Proceedings of the Conference on Empirical Methods in Natural Language Processing, pages 1317–1327.
- Jesse Mu, Xiang Li, and Noah Goodman. 2023. Learning to compress prompts with gist tokens. In Advances in Neural Information Processing Systems, pages 19327–19352.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-roll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. OpenAI Blog.
- Antonio Polino, Razvan Pascanu, and Dan Alistarh. 2018. Model compression via distillation and quantization. In International Conference on Learning Representations.
- Zhiqing Sun, Yikang Shen, Qinzhong Zhou, HongxinZhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. 2023. Principle-driven self-alignment of language models from scratch with minimal human supervision. In Advances in Neural Information Processing Systems, pages 2511–2565.
- Chaofan Tao, Lu Hou, Wei Zhang, Lifeng Shang, Xinjiang, Qun Liu, Ping Luo, and Ngai Wong. 2022. Compression of generative pre-trained language models via quantization. In Proceedings of the Annual Meeting of the Association for Computational Linguistics, pages 4821–4836.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford Alpaca: An instruction-following LLaMA model. StandfordBlog.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, AlisaLiu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions. In Proceedings of the Annual Meeting of the Association for Computational Linguistics, pages 13484–13508.
- Yuqiao Wen, Yanshuai Cao, and Lili Mou. 2025a. Exploring model invariance with discrete search for ultra-low-bit quantization. arXiv preprint arXiv:2502.06844.
- Yuqiao Wen, Zichao Li, Wenyu Du, and Lili Mou. 2023b. f-divergence minimization for sequence-level knowledge distillation. In Proceedings of the Annual Meeting of the Association for Computational Linguistics, pages 10817–10834.
- Chunting Zhou, Jiatao Gu, and Graham Neubig. 2020. Understanding knowledge distillation in non-autoregressive machine translation. In International Conference on Learning Representations.

# Thanks for attending the tutorial!

## Q&A



VECTOR  
INSTITUTE

CIFAR

# References

- Chung, Inseop, et al. "Feature-map-level online adversarial knowledge distillation." *International Conference on Machine Learning*. PMLR, 2020.
- Guo, Daya, et al. "Deepseek-r1: Incentivizing reasoning capability in LLMs via reinforcement learning." arXiv preprint arXiv:2501.12948 (2025).
- Hao, Yongchang, Yuxin Liu, and Lili Mou. "Teacher forcing recovers reward functions for text generation." *Advances in Neural Information Processing Systems* 35 (2022): 12594-12607.
- Hinton, Geoffrey, Oriol Vinyals, and Jeff Dean. "Distilling the knowledge in a neural network." arXiv preprint arXiv:1503.02531 (2015).
- Jiao, Xiaoqi, et al. "TinyBERT: Distilling BERT for Natural Language Understanding." *Findings of the Association for Computational Linguistics: EMNLP 2020*. 2020.
- Kim, Yoon, and Alexander M. Rush. "Sequence-level knowledge distillation." *Proceedings of the 2016 conference on empirical methods in natural language processing*. 2016.
- Li, Dongheng, Yongchang Hao, and Lili Mou. "Llmr: Knowledge distillation with a large language model-induced reward." COLING, 2024
- Narayanan, Deepak, et al. "Efficient large-scale language model training on GPU clusters using megatron-LM." *Proceedings of the international conference for high performance computing, networking, storage and analysis*. 2021.
- Passban, Peyman, et al. "ALP-KD: Attention-based layer projection for knowledge distillation." *Proceedings of the AAAI Conference on artificial intelligence*. Vol. 35. No. 15. 2021.
- Sun, Siqi, et al. "Patient Knowledge Distillation for BERT Model Compression." *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2019.
- Tang, Jiaxi, and Ke Wang. "Ranking distillation: Learning compact ranking models with high performance for recommender system." *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 2018.
- Tunstall, Lewis, et al. "Zephyr: Direct distillation of lm alignment." arXiv preprint arXiv:2310.16944 (2023).
- Wang, Wenhui, et al. "MiniLM: Deep self-attention distillation for task-agnostic compression of pre-trained transformers." *Advances in neural information processing systems* 33 (2020): 5776-5788.
- Xu, Xiaohan, et al. "A survey on knowledge distillation of large language models." arXiv preprint arXiv:2402.13116 (2024).
- Wen, Yuqiao, et al. "F-divergence minimization for sequence-level knowledge distillation." ACL, 2023.
- Yu, Zony, Yuqiao Wen, and Lili Mou. "Revisiting Intermediate-Layer Matching in Knowledge Distillation: Layer-Selection Strategy Doesn't Matter (Much)." AACL-Findings, 2025.