

# Unified **R**einforcement and **I**mitation **L**earning for Vision Language Models

Byung-Kwan Lee<sup>1,2\*</sup> Ryo Hachiuma<sup>1</sup> Yong Man Ro<sup>2</sup>

Yu-Chiang Frank Wang<sup>1,3</sup> Yueh-Hua Wu<sup>1</sup>

\*Work Done during Internship

1



2



3



# Motivation

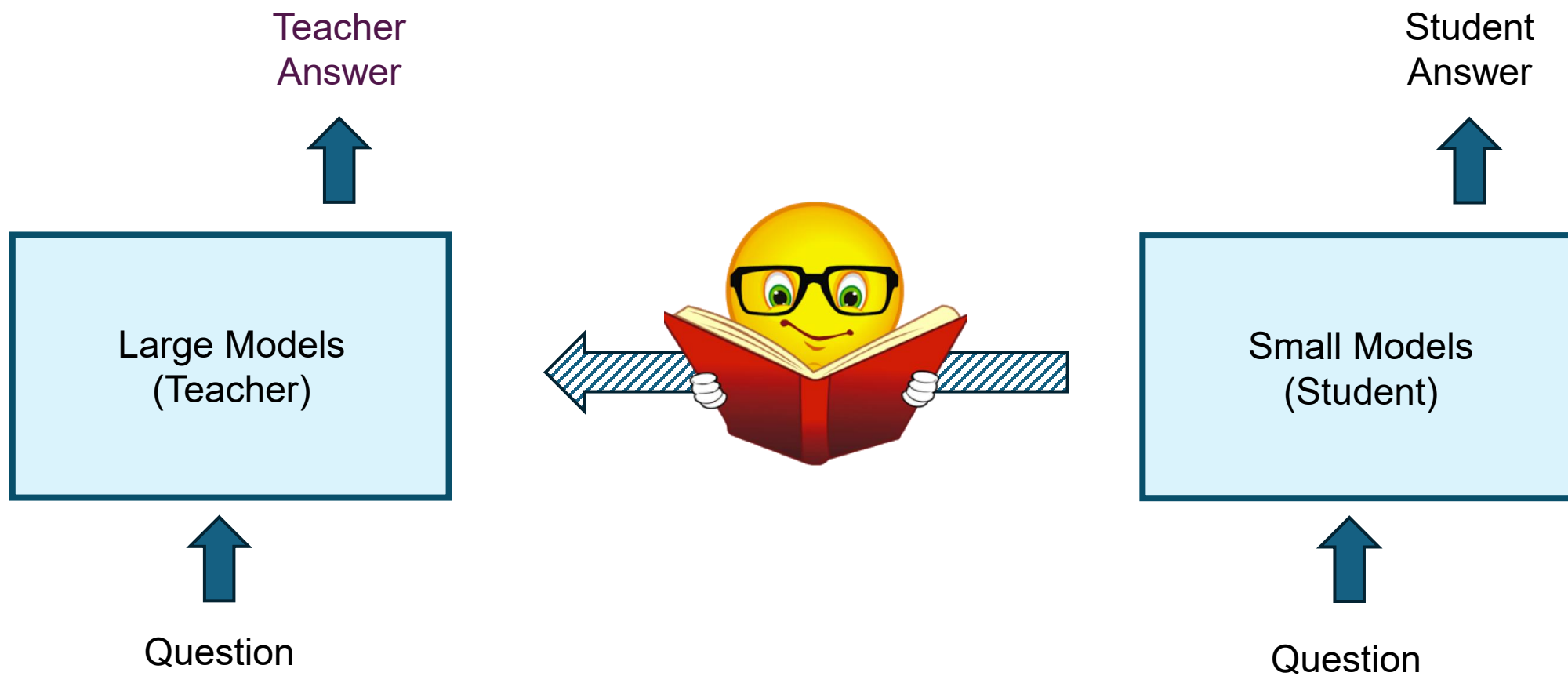
- Limitations of Existing VLM Scaling and Architectural Changes
- Demand for Inference Efficiency and Lightweight Models



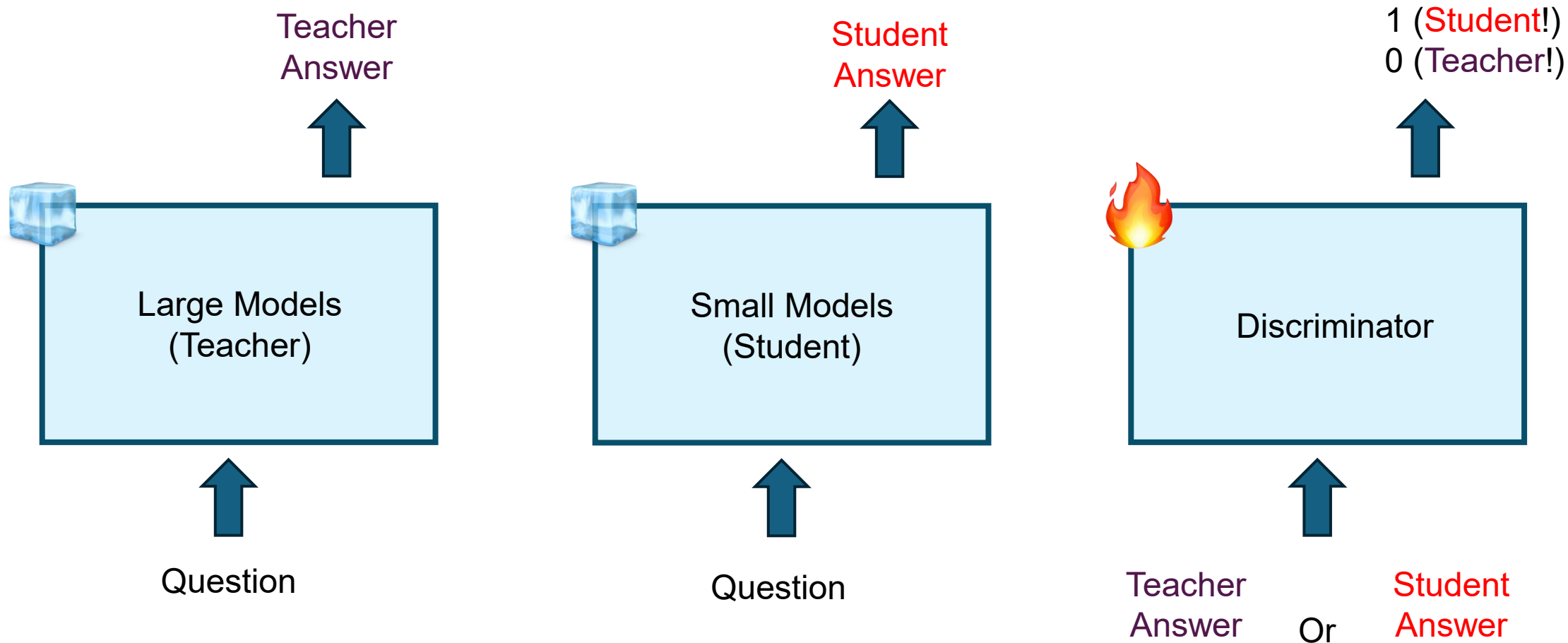
**How can remove think-answer but get high-performance?**

**Think-Answer** process is actually too much slow.

# Motivation

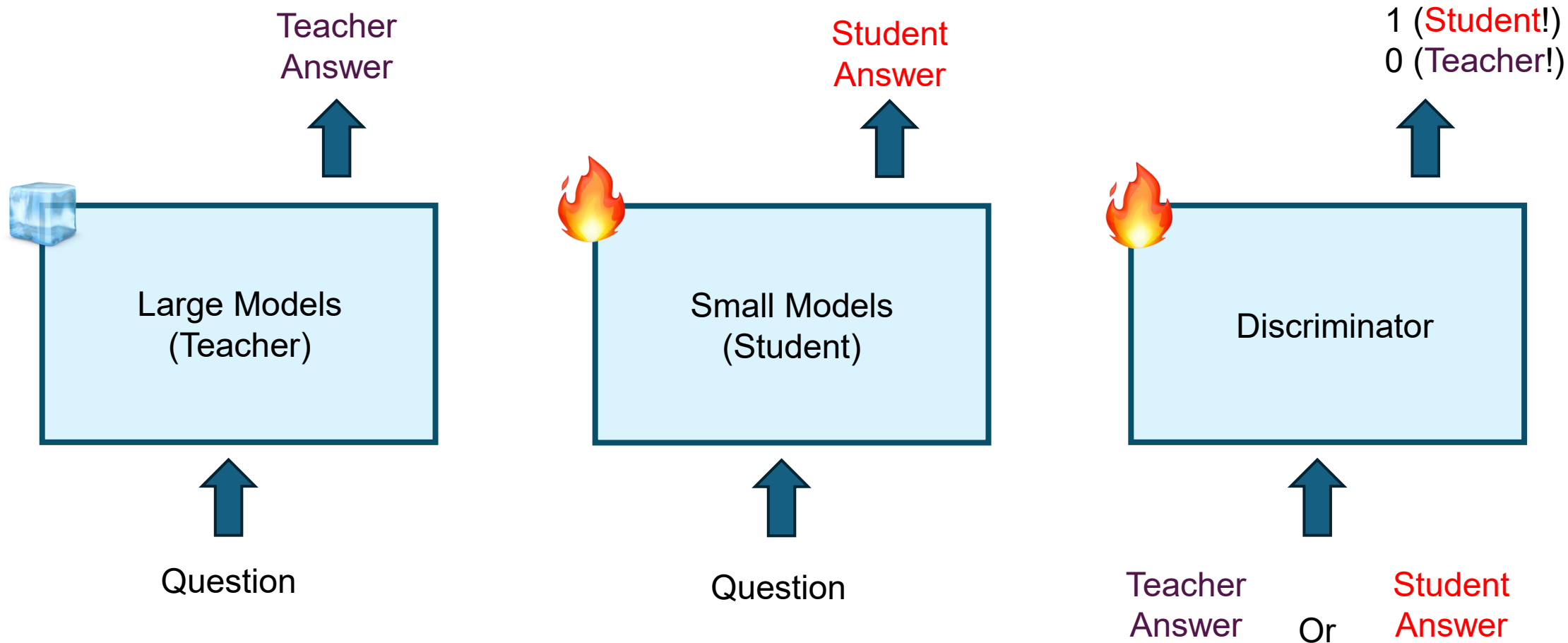


# Proposed Method (RIL): Training Step1



$$\max_{\phi} \mathcal{L}(\phi) = \frac{1}{N} \sum_{i=1}^N \left\{ \log D_{\phi}(q, o_i^{(s)}) + \log(1 - D_{\phi}(q, o_i^{(t)})) \right\}$$

## Proposed Method (RIL): Training Step2



**Reward:**  $R_{\text{teacher similarity}}$  +  $R_{\text{answer}}$

$$\mathbb{1}(D_{\phi}(q, o_i) < 0.5) \quad \text{LLM-as-a-Judge}(q, a, o_i)$$

# Experimental Results

---

## Algorithm 2 RIL for VLMs

---

**Require:** Pre-trained discriminator  $D_\phi$  and Pre-trained student VLMs  $\pi_{\theta_{\text{init}}}$

**Require:** Saved collection for the generated text responses  $\{o^{(1)}\}$  from teacher VLMs

- 1: Set reference model  $\pi_{\text{ref}} \leftarrow \pi_{\theta_{\text{init}}}$
  - 2: Set the training model  $\pi_\theta \leftarrow \pi_{\theta_{\text{init}}}$
  - 3: **for** sample a batch  $\mathcal{B}$  in Dataset **do**
  - 4:   Copy and freeze model  $\pi_{\theta_{\text{old}}} \leftarrow \pi_\theta$
  - 5:   Sample  $G$  outputs  $\{o_i^{(s)}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|q)$  for each question  $q \in \mathcal{B}$
  - 6:   Extract  $G$  outputs  $\{o_i^{(t)}\}_{i=1}^G$  in the saved collection for each question  $q \in \mathcal{B}$
  - 7:   **for** Discriminator Updating iteration =  $1, 2, \dots, \mu$  **do**
  - 8:     Update  $D_\phi$  by using Equation 1
  - 9:   **end for**
  - 10:   Compute Rewards and Advantages for all  $2G$  outputs by Discriminator and LLM-as-a-Judge
  - 11:   **for** Student VLMs Updating iteration =  $1, 2, \dots, \mu$  **do**
  - 12:     Update  $\pi_\theta$  by using Equation 2
  - 13:   **end for**
  - 14: **end for**
-

# Experimental Results

VLMs	AI2D	ChartQA	MathVista	MMB	MMB <sup>CN</sup>	MM-Vet	MM-Vet-v2	MMMU	MMMU-Pro	MMStar	BLINK	SEED	SEED2+	RWQA
Qwen2.5-VL-7B	83.9	87.3	67.8	83.5	83.4	71.8	63.7	55.0	38.3	63.9	56.4	77.0	70.4	68.5
w. RIL (Qwen2.5-VL-72B)	86.7	95.4	74.5	<b>86.8</b>	<b>87.2</b>	77.3	66.1	61.8	48.2	<b>71.1</b>	68.5	80.7	<b>73.0</b>	74.2
w. RIL (InternVL3-78B)	<b>86.8</b>	95.5	74.6	86.7	87.1	75.8	66.0	60.9	47.1	<b>71.1</b>	68.1	<b>80.8</b>	72.7	<b>75.4</b>
w. RIL (Both)	86.1	<b>95.6</b>	<b>79.7</b>	86.3	86.5	<b>80.4</b>	<b>71.1</b>	<b>65.7</b>	<b>48.5</b>	<b>71.1</b>	<b>70.0</b>	80.5	72.8	72.8

