

# Reinforcement Learning Guided Semi-Supervised Learning

Marzi Heidari<sup>1</sup>, Hanping Zhang<sup>1</sup>, Yuhong Guo<sup>1,2</sup>

<sup>1</sup>Carleton University, Ottawa, Canada

<sup>2</sup>CIFAR AI Chair, Amii, Canada

# Introduction

- **Semi-Supervised Learning (SSL):** Uses a small labeled dataset with a large unlabeled pool.
- **Limitations:** Heuristics or predefined rules for pseudo-labeling methods are often suboptimal.
- **Challenge:** How can we better leverage unlabeled data to guide the learning process?

# Reinforcement Learning Guided Semi-Supervised Learning (RLGSSL)

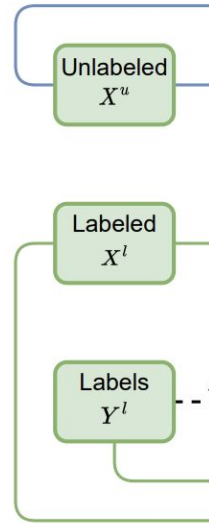
- Treats SSL as a **one-armed bandit problem**.
- **Dynamically** adapt and respond to the data.
- **Beyond standard norm** for SSL.
- Potential to **transform** SSL frameworks.

# Framework

- Frames SSL as **one-armed bandit problem**.

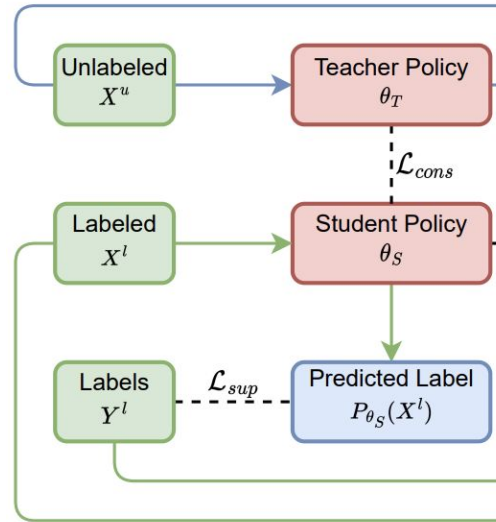
# Framework

- Frames SSL as **one-armed bandit problem**.
- **State:**  $s = (X^l, Y^l, X^u)$ .



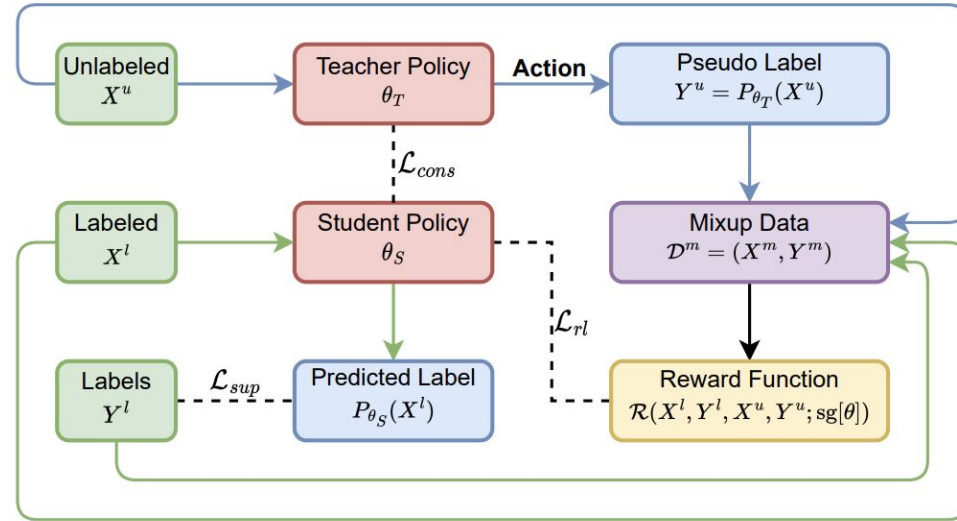
# Framework

- Frames SSL as **one-armed bandit problem**.
- **State:**  $s = (X^l, Y^l, X^u)$ .
- **Policy:** The prediction networks ( $\pi_{\theta}(\cdot) = P_{\theta}(\cdot)$ )
- **Action:** Model's predictions (pseudo-labels).



# Framework

- Frames SSL as **one-armed bandit problem**.
- **State:**  $s = (X^l, Y^l, X^u)$ .
- **Policy:** The prediction networks. ( $\pi_{\theta}(\cdot) = P_{\theta}(\cdot)$ )
- **Action:** Model's predictions (pseudo-labels).
- **Reward:** Generalization measured via label mixup between labeled and pseudo-labeled samples.



# Reward Function

- **Balanced** utilization of labeled and unlabeled data by Inter-mixup:

$$x_i^m = \mu x_i^u + (1 - \mu) x_i^l, \quad \mathbf{y}_i^m = \mu \mathbf{y}_i^u + (1 - \mu) \mathbf{y}_i^l$$

- **Reward: Negative disagreement** between model predictions and mixup labels:

$$\mathcal{R}(s, a; \text{sg}[\theta]) = \mathcal{R}(X^l, Y^l, X^u, Y^u; \text{sg}[\theta]) = -\frac{1}{C \cdot N^m} \sum_{i=1}^{N^m} \|P_\theta(x_i^m) - \mathbf{y}_i^m\|_2^2$$



# Reinforcement Learning Loss

- **One-Armed Bandit Principle:** Optimize one-time reward based on the policy output.
- Exploits **non-differentiable reward**.
- Enables policy gradient with a **deterministic policy**.
- **KL-Divergence Weighted Negative Reward:**

$$\mathcal{L}_{\text{rl}}(\theta) = -\mathbb{E}_{\mathbf{y}_i^u \sim \pi_\theta} \text{KL}(\mathbf{e}, \mathbf{y}_i^u) \mathcal{R}(s, a; \text{sg}[\theta]) = -\mathbb{E}_{x_i^u \in \mathcal{D}_u} \text{KL}(\mathbf{e}, P_\theta(x_i^u)) \mathcal{R}(s, a; \text{sg}[\theta])$$

- Measures distance between label predictions and a uniform distribution vector  $\mathbf{e}=1/C$

# Teacher Student Framework

- Teacher parameter update via WMA:

$$\theta_T = \beta \theta_T + (1 - \beta) \theta_S$$

- **Supervised Loss** on labeled data:

$$\mathcal{L}_{\text{sup}}(\theta_S) = \mathbb{E}_{(x^l, \mathbf{y}^l) \in \mathcal{D}^l} [\ell_{CE}(P_{\theta_S}(x^l), \mathbf{y}^l)]$$

- **Consistency Loss** between student and teacher on unlabeled data:

$$\mathcal{L}^{\text{cons}} = \mathbb{E}_{x^u \in \mathcal{D}^u} [\ell_{\text{KL}}(P_{\theta_S}(x^u), P_{\theta_T}(x^u))]$$

- **Learning objective:**  $\mathcal{L}(\theta_S) = \mathcal{L}_{\text{rl}} + \lambda_1 \mathcal{L}_{\text{sup}} + \lambda_2 \mathcal{L}_{\text{cons}}$

# Experimental Results

Table 1: Performance of RLGSSL and state-of-the-art SSL algorithms with the CNN-13 network. We report the average test errors and the standard deviations of 5 trials.

Dataset Number of Labeled Samples	CIFAR-10			CIFAR-100	
	1000	2000	4000	4000	10000
Supervised	39.95 <sub>(0.75)</sub>	27.67 <sub>(0.12)</sub>	20.42 <sub>(0.21)</sub>	58.31 <sub>(0.89)</sub>	44.56 <sub>(0.30)</sub>
Supervised + MixUp [40]	31.83 <sub>(0.65)</sub>	24.22 <sub>(0.15)</sub>	17.37 <sub>(0.35)</sub>	54.87 <sub>(0.07)</sub>	40.97 <sub>(0.47)</sub>
$\Pi$ -model [6]	28.74 <sub>(0.48)</sub>	17.57 <sub>(0.44)</sub>	12.36 <sub>(0.17)</sub>	55.39 <sub>(0.55)</sub>	38.06 <sub>(0.37)</sub>
Temp-ensemble [6]	25.15 <sub>(1.46)</sub>	15.78 <sub>(0.44)</sub>	11.90 <sub>(0.25)</sub>	-	38.65 <sub>(0.51)</sub>
Mean Teacher[8]	21.55 <sub>(0.53)</sub>	15.73 <sub>(0.31)</sub>	12.31 <sub>(0.28)</sub>	45.36 <sub>(0.49)</sub>	35.96 <sub>(0.77)</sub>
VAT [5]	18.12 <sub>(0.82)</sub>	13.93 <sub>(0.33)</sub>	11.10 <sub>(0.24)</sub>	-	-
SNTG [15]	18.41 <sub>(0.52)</sub>	13.64 <sub>(0.32)</sub>	10.93 <sub>(0.14)</sub>	-	37.97 <sub>(0.29)</sub>
Learning to Reweight [41]	11.74 <sub>(0.12)</sub>	-	9.44 <sub>(0.17)</sub>	46.62 <sub>(0.29)</sub>	37.31 <sub>(0.47)</sub>
MT + Fast SWA [14]	15.58 <sub>(0.12)</sub>	11.02 <sub>(0.23)</sub>	9.05 <sub>(0.21)</sub>	-	33.62 <sub>(0.54)</sub>
ICT [16]	12.44 <sub>(0.57)</sub>	8.69 <sub>(0.15)</sub>	7.18 <sub>(0.24)</sub>	40.07 <sub>(0.38)</sub>	32.24 <sub>(0.16)</sub>
RLGSSL (Ours)	<b>9.15</b> <sub>(0.57)</sub>	<b>6.90</b> <sub>(0.11)</sub>	<b>6.11</b> <sub>(0.10)</sub>	<b>36.92</b> <sub>(0.45)</sub>	<b>29.12</b> <sub>(0.20)</sub>

# Thank You