



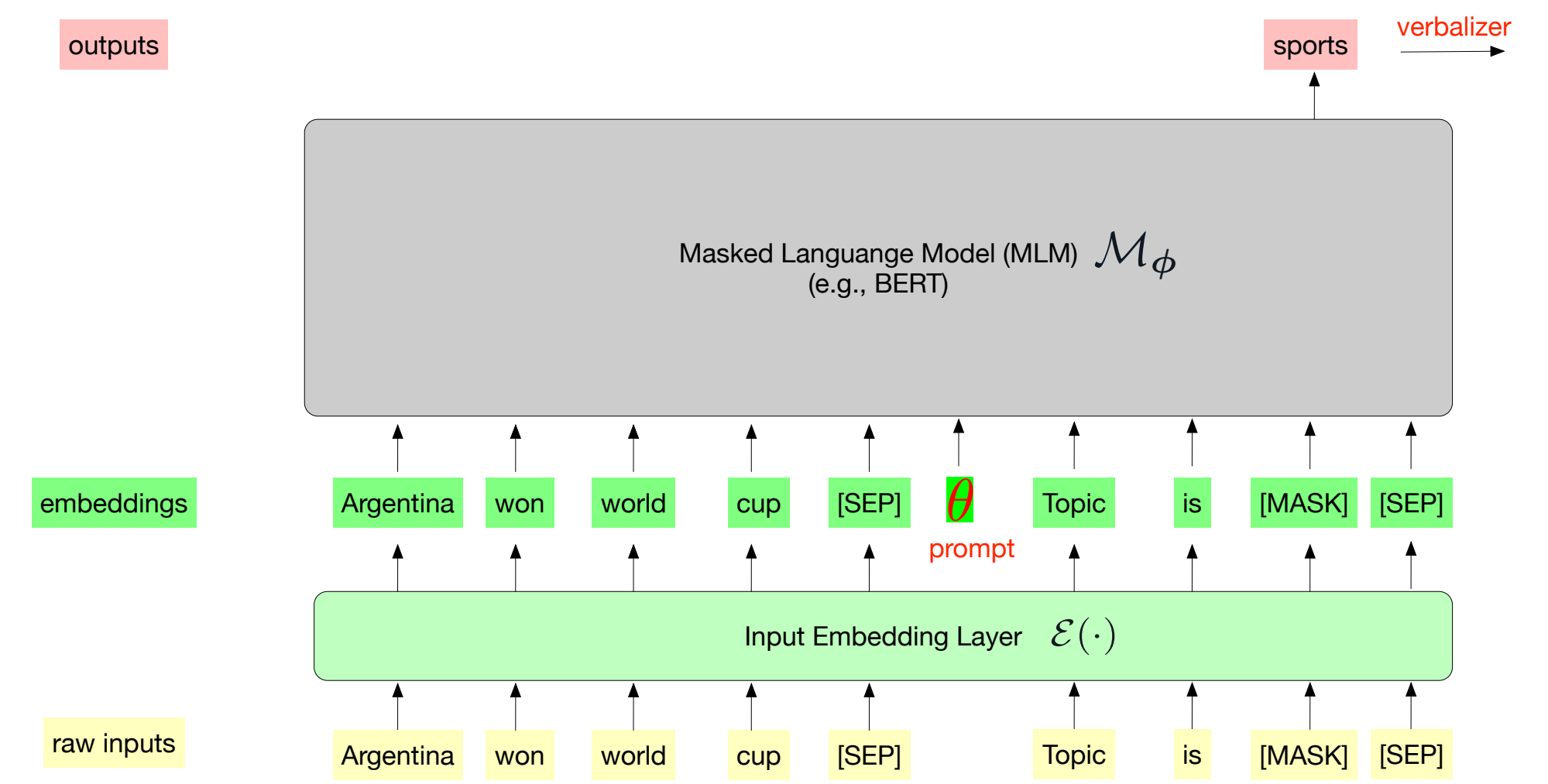
Background

Masked language model (MLM) can fill the masked tokens in a sentence, thus, can be used for classification tasks:

- input sentence is wrapped by a **template**: "sentence, Topic is [MASK]";
- MLM **predicts a token** at the [MASK] position;
- a **verbalizer** map the predicted token to label.

Prompt tuning introduces a **learnable prompt** θ in the template:

$$\tilde{x} \equiv \mathbb{T}(x; \theta) = (\mathcal{E}(x), \theta, \mathcal{E}(\text{Topic}), \mathcal{E}(\text{is}), \mathcal{E}([\text{MASK}])).$$



- The initialization of θ plays an important role in prompt tuning.
- Recently, MetaPrompting (COLING 2022) proposes to meta-learn a shared prompt initialization for all task-specific prompts with a hand-crafting verbalizer.

Challenges:

- a **single** meta-initialized prompt is insufficient for adaptation to complex tasks;
- not parameter-efficient** as the whole MLM needs tuning;
- hand-crafting verbalizer is **labor-intensive**.

Notations: task τ with support set \mathcal{S}_τ , query set \mathcal{Q}_τ , and label set \mathcal{Y}_τ .

Representative Verbalizer (RepVerb)

We propose a novel verbalizer **RepVerb**, which is simple and effective:

- For input x , compute its feature embedding $h_{[\text{MASK}]}(\tilde{x})$.
- For class y , compute label embedding $v_y = \frac{1}{|\mathcal{S}_{\tau,y}|} \sum_{(x,y) \in \mathcal{S}_{\tau,y}} h_{[\text{MASK}]}(\tilde{x})$.
- Prediction: compute cosine similarity between $h_{[\text{MASK}]}(\tilde{x})$ and $\{v_y : y \in \mathcal{Y}_\tau\}$.

procedure ComputeLabelEmbedding(\mathcal{S}_τ):

compute $h_{[\text{MASK}]}(\tilde{x})$ for $(x, \cdot) \in \mathcal{S}_\tau$;
compute v_y for each $y \in \mathcal{Y}_\tau$;

end procedure

procedure Predict($x; v_y : y \in \mathcal{Y}_\tau$):

compute $h_{[\text{MASK}]}(\tilde{x})$ for x ;
 $\hat{\mathbb{P}}(y|x; \phi, \theta) = \frac{\exp(\rho \cos(v_y, h_{[\text{MASK}]}(\tilde{x})))}{\sum_{y' \in \mathcal{Y}_\tau} \exp(\rho \cos(v_{y'}, h_{[\text{MASK}]}(\tilde{x})))}$

end procedure

MetaPrompter

Use a **prompt pool to extract more task knowledge** for constructing instance-dependent prompt:

- The **prompt pool** has K learnable prompts $\{(k_i, \theta_i) : i = 1, \dots, K\}$, with key k_i and value θ_i ;
- Instance-dependent prompt** is constructed by a weighted combination of all values (θ_i 's): $\theta_x(\mathbf{K}, \Theta) = \sum_{i=1}^K a_i \theta_i$, where attention weight a is computed between input x and the K prompts.

Prediction (**hand-crafting verbalizer** + **RepVerb**):

$$\mathbb{P}(y|x; \theta_x) = (1 - \lambda) \times \hat{\mathbb{P}}(y|x; \theta_x) + \lambda \times \tilde{\mathbb{P}}(y|x; \theta_x)$$

where $\hat{\mathbb{P}}(y|x; \theta) = \frac{1}{|\mathcal{Y}_y|} \sum_{w \in \mathcal{Y}_y} \mathbb{P}_{\mathcal{M}}([\text{MASK}] = w | \mathbb{T}(x; \theta))$ (\mathcal{Y}_y is a set of label tokens).

We propose a novel algorithm **MetaPrompter** to learn the prompt pool by meta-learning (e.g., MAML):

- base learner**: (1) build instance-dependent prompts from (\mathbf{K}, Θ) (2) compute loss on support set $\mathcal{L}(\mathcal{S}_\tau; \mathbf{K}, \Theta)$ (3) update task-specific prompt pool $(\mathbf{K}^{(\tau)}, \Theta^{(\tau)}) = (\mathbf{K}, \Theta) - \alpha \nabla_{(\mathbf{K}, \Theta)} \mathcal{L}(\mathcal{S}_\tau; \mathbf{K}, \Theta)$.
- meta-learner**: (1) compute loss on query set $\mathcal{L}(\mathcal{Q}_\tau; \mathbf{K}^{(\tau)}, \Theta^{(\tau)})$ (2) update meta prompt pool $(\mathbf{K}, \Theta) \leftarrow (\mathbf{K}, \Theta) - \eta \nabla_{(\mathbf{K}, \Theta)} \mathcal{L}(\mathcal{Q}_\tau; \mathbf{K}^{(\tau)}, \Theta^{(\tau)})$.

Compared with MetaPrompting, MetaPrompter is

- Parameter-efficient**: only the prompt pool is tuned;
- More flexible**: Instance-dependent prompt allows better adaptation to complex tasks.

Evaluation on RepVerb

Table 1: Meta-testing accuracy of 5-way few-shot classification.

		20News	Amazon	HuffPost	Reuters	HWU64	Liu54
5-shot	WARP	61.43	59.53	46.31	68.67	68.60	73.11
	ProtoVerb	71.33	71.74	57.93	80.93	73.43	76.19
	RepVerb	78.81	77.56	61.90	88.33	78.37	82.14
1-shot	WARP	49.87	48.94	38.21	52.88	53.20	58.68
	RepVerb	59.86	59.18	44.65	63.63	59.83	66.17

Baselines: soft verbalizers WARP (ACL 2021) and ProtoVerb (ACL 2022) learned by supervised and contrastive learning, respectively.

RepVerb **outperforms** WARP and ProtoVerb on both the 1-shot and 5-shot settings.

Hence, RepVerb is **effective**.

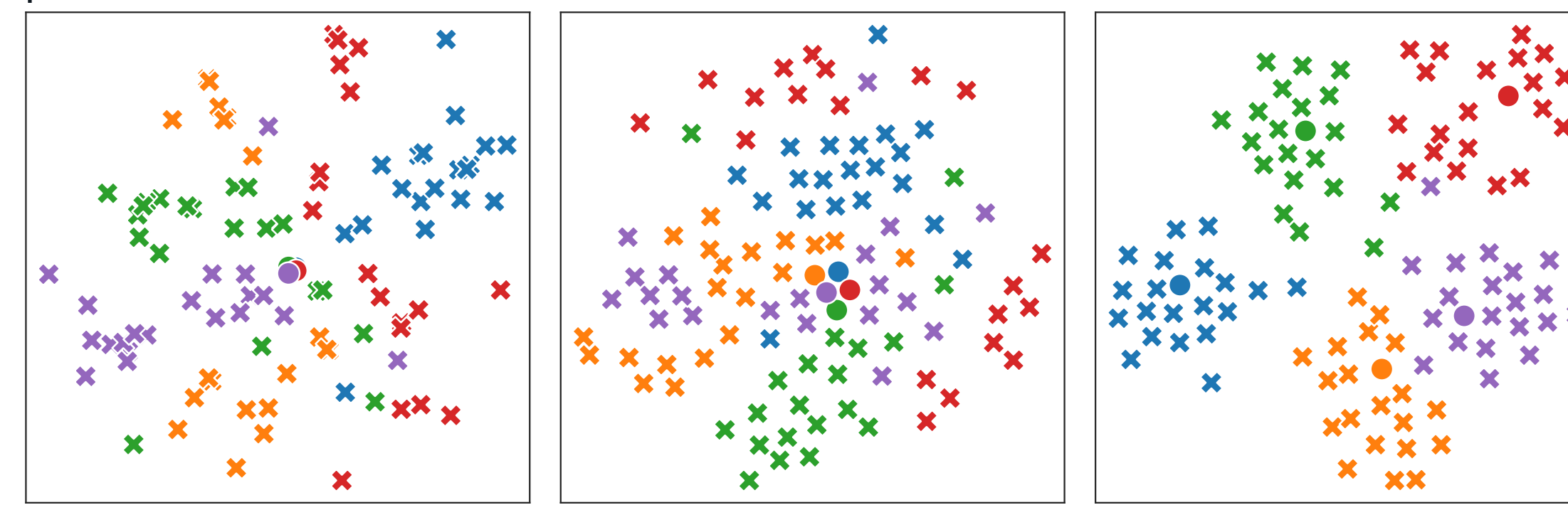


Figure 1: t-SNE visualization of feature embeddings (crosses) and label embeddings (circles) for a 5-way 5-shot task from Reuters.

RepVerb has **more discriminative and compact** embeddings than WARP and ProtoVerb.

By design, RepVerb's label embedding is **consistent** with samples' embeddings, but those of WARP and ProtoVerb are not.

Evaluation on MetaPrompter

Table 2: 5-way 5-shot classification meta-testing accuracy.

	#param ($\times 10^6$)	20News	Amazon	HuffPost	Reuters	HWU64	Liu54
HATT	0.07	55.00	66.00	56.30	56.20	—	—
DS	1.73	68.30	81.10	63.50	96.00	—	—
MLADA	0.73	77.80	86.00	64.90	96.70	—	—
ConstrastNet	109.52	71.74	85.17	65.32	95.33	92.57	93.72
MetaPrompting	109.52	85.67	84.19	72.85	95.89	93.86	94.01
MetaPrompting+WARP	109.52	85.81	85.54	71.71	97.28	93.99	94.33
MetaPrompting+ProtoVerb	109.52	86.18	84.91	73.11	97.24	93.81	94.38
MetaPrompting+RepVerb	109.52	86.89	85.98	74.62	97.32	94.23	94.45
MetaPrompter	0.06	88.57	86.36	74.89	97.63	95.30	95.47

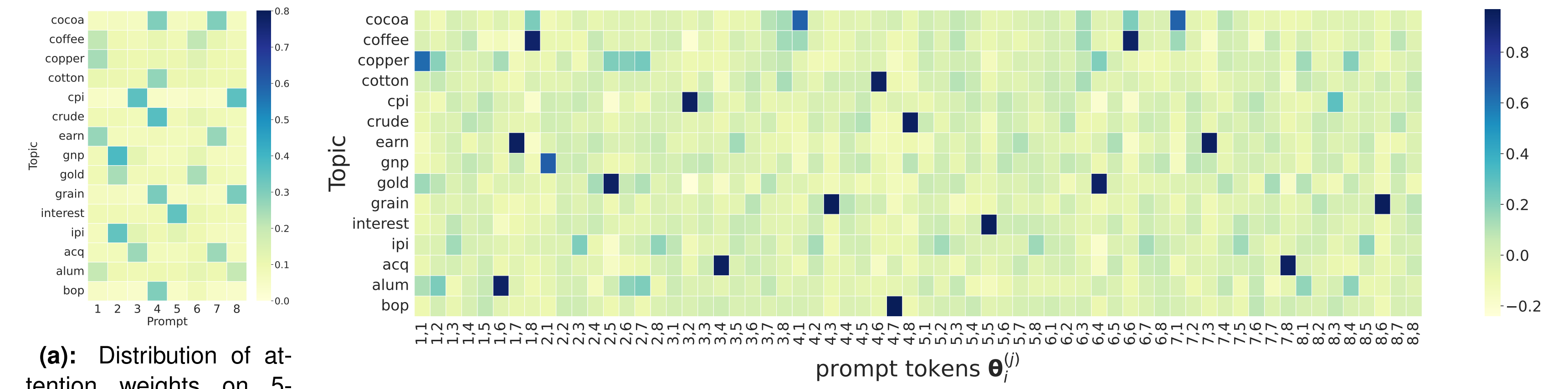
Baselines: (i) state-of-the-art prompt-based methods (MetaPrompting and its variants); (ii) non-prompt-based methods (HATT (AAAI 2019), DS (ICLR 2020), MLADA (ACL 2021), ConstrastNet (AAAI 2022)).

MetaPrompter is **better than** both prompt-based and non-prompt-based baselines.

MetaPrompter **outperforms** MetaPrompting+RepVerb, showing effectiveness of the prompt pool.

MetaPrompter is **much more parameter-efficient** than MetaPrompting (1800 \times fewer).

RepVerb is **beneficial** to MetaPrompting.



(a): Distribution of attention weights on 5-way 5-shot classification of Reuters (15 topics).

(b): Cosine similarities between learned prompt tokens and topic embeddings on 5-way 5-shot classification of Reuters (recall that K and L_p are set to 8). In the x-axis, (i, j) stands for the j th row of θ_i (i.e., $\theta_i^{(j)}$).

Samples from topic *cocoa* prefer the 4th and 7th prompts (left), as *cocoa*'s embedding is similar to $\theta_4^{(1)}$ and $\theta_7^{(1)}$ (right).

Summary

Problem: improve the effectiveness and parameter-efficiency of prompt tuning.

Propose a novel algorithm MetaPrompter, consisting of:

- a **novel verbalizer RepVerb**: simple and effective.
- structured prompting: a **meta-learned prompt pool** to construct instance-dependent prompts.

MetaPrompter is **parameter-efficient** as only the prompt pool is tuned.

Experimental results demonstrate:

- RepVerb achieves higher accuracy than other soft verbalizers (WARP and ProtoVerb).
- MetaPrompter performs better than MetaPrompting.
- MetaPrompter is much more parameter-efficient than MetaPrompting (1800 \times).