



西安交通大学  
XI'AN JIAOTONG UNIVERSITY

# Concentrate Attention: Towards Domain- Generalizable Prompt Optimization for Language Models

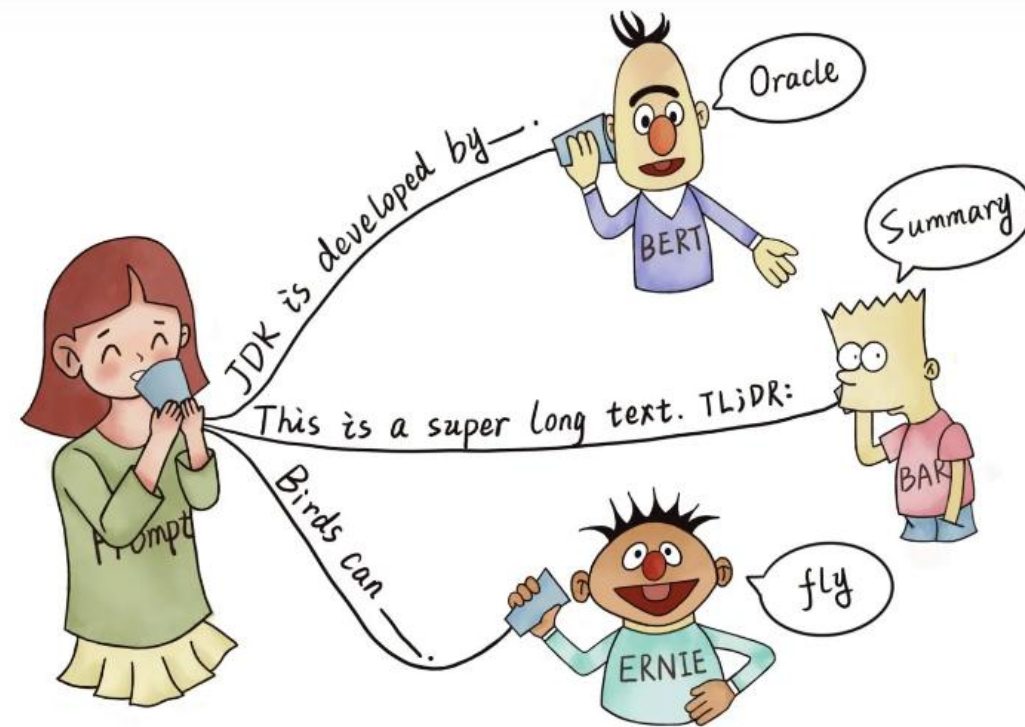
**Chengzhengxu Li, Xiaoming Liu<sup>★</sup>, Zhaohan Zhang,  
Yichen Wang, Chen Liu, Yu Lan, Chao Shen**

Faculty of Electronic and Information Engineering, Xi'an Jiaotong University  
{[czx.li](mailto:czx.li@stu.xjtu.edu.cn), [lcoder](mailto:lcoder@stu.xjtu.edu.cn)}@stu.xjtu.edu.cn {[xm.liu](mailto:xm.liu@xjtu.edu.cn), [ylan2020](mailto:ylan2020@xjtu.edu.cn), [chaoshen](mailto:chaoshen@xjtu.edu.cn)}@xjtu.edu.cn  
[zhaohan.zhang@qmul.ac.uk](mailto:zhaohan.zhang@qmul.ac.uk) [yichenzw@uchicago.edu](mailto:yichenzw@uchicago.edu)

★ Corresponding author

# Agenda

- Motivation
- The Proposed Method
  - Definition of Concentration
  - Pilot Experiment Results
  - Concentrative Soft Prompt Optimization
  - Concentrative Hard Prompt Optimization
- Experiment Results
- Conclusion & Future work



# Motivation

## ➤ Discrete Prompt Optimization

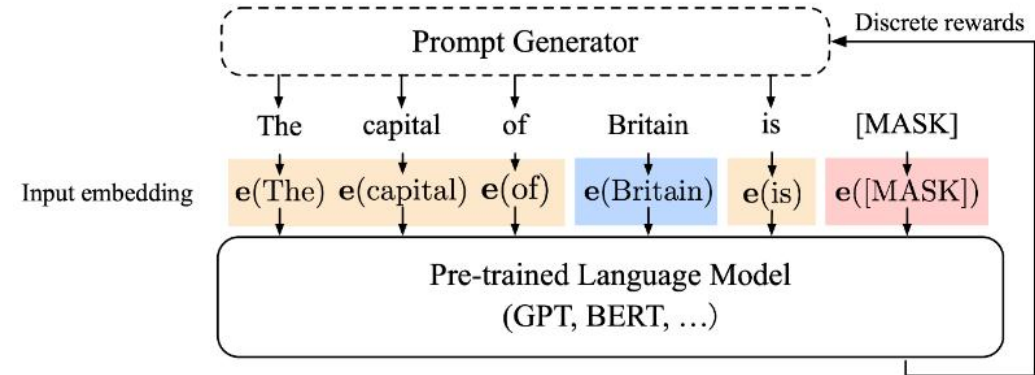
- Requiring considerable expertise
- Inefficient optimization
- Only one prompt for the whole dataset

## ➤ Continuous Prompt Optimization

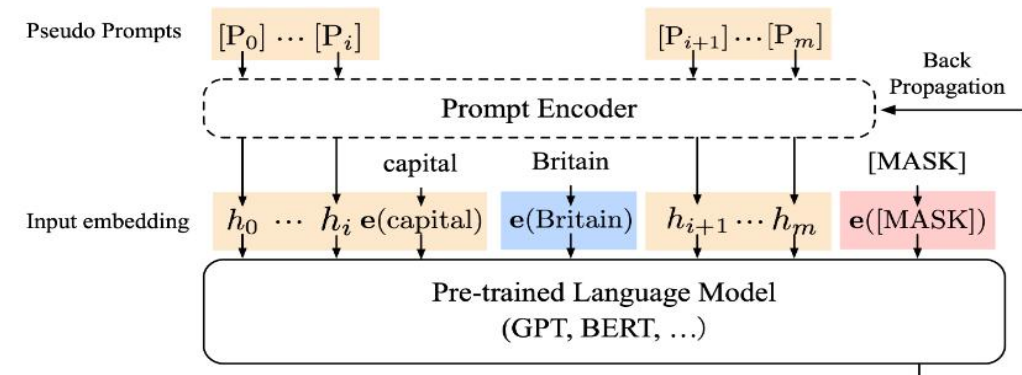
- High computational cost
- Low readability and generalizability
- Only one prompt for the whole dataset

## ➤ Challenges

- Three challenges for few-shot: **data scarcity, overfitting and noise impact**
- Two challenges for prompting: **training efficiency and robustness to verbalizers**



### Discrete Prompt Optimization



### Continuous Prompt Optimization

# Motivation

➤ However, the domain generalization ability of trained prompts still lacks exploration.

## ➤ Domain Adaptation Methods

- Target domain availability based
- Align target domain by unsupervised learning
- Serious data reliance

## ➤ Pre-training Prompt Methods

- Large data requirements
- Low efficiency training
- High computational cost

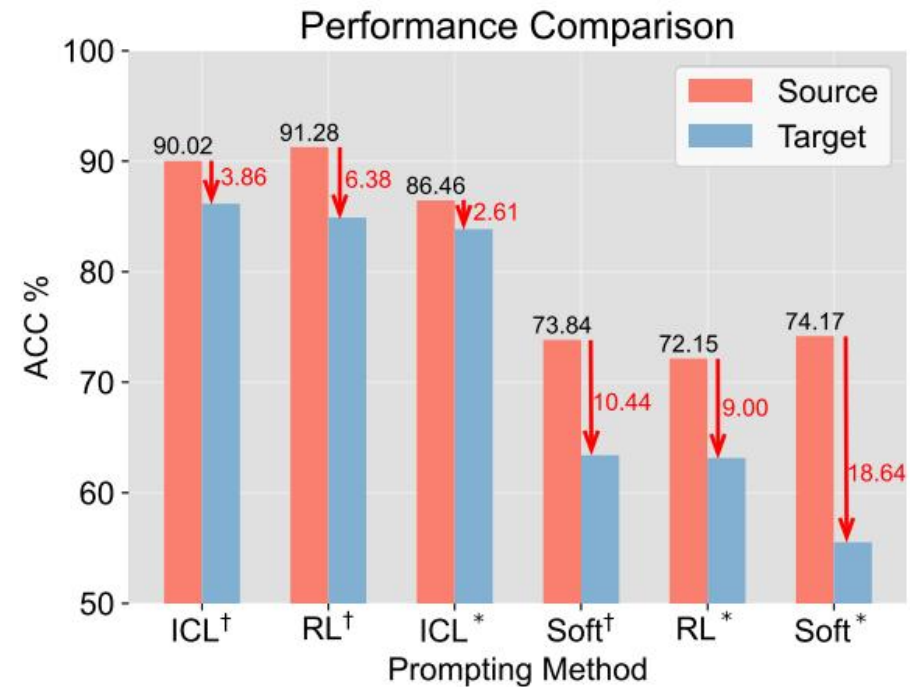


Figure 1. Domain generalization capabilities across various prompting methods in sentiment classification tasks.

## *What Nature Do Well-generalized Prompts Have?*

### ➤ Domain Adaptation Methods

- Target domain availability based
- Align target domain by unsupervised learning
- Serious data reliance

### ➤ Pre-training Prompt Methods

- Large data requirements
- Low efficiency training
- High computational cost

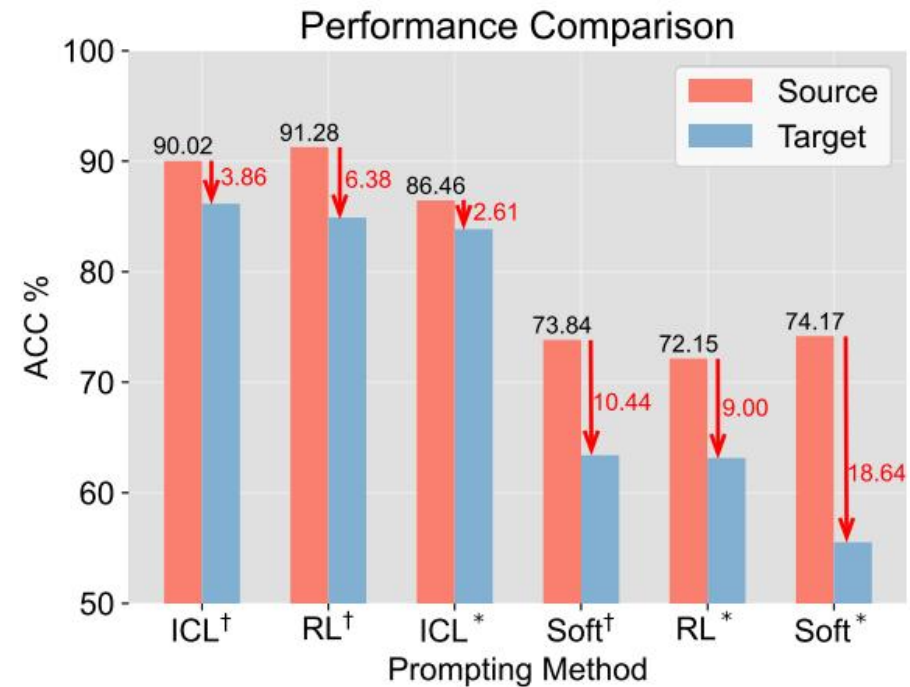


Figure 1. Domain generalization capabilities across various prompting methods in sentiment classification tasks.

# The Proposed Method: Definition of Concentration

Heuristically, concentration represents the “lookback” attention from current decoding token to prompt tokens, as shown in Figure 2.

**Definition 3.1.** Let  $z = (z_1, z_2, \dots, z_L)$  and  $x = (e_1, e_2, \dots, e_T)$  be prompt and original input with  $z, x \in S$ , where  $S$  is the set of all possible textual sequences over the vocabulary. Let  $f_{\theta_l}$  be the attention block in layer  $l$  of a PLM parameterized by  $\theta_l$ . Then concentration is a function  $\text{Concentration} : S \rightarrow \mathbb{R}^+$

$$\text{Concentration}(z \oplus x; \theta_l) = \sum_{z_i \in z} f_{\theta_l}(z_i \oplus x).$$

**Definition 3.2.** Let  $z = (z_1, z_2, \dots, z_L)$  and  $x = (e_1, e_2, \dots, e_T)$  be prompt and original input with  $z, x \in S$ , where  $S$  is the set of all possible textual sequences over the vocabulary. Let  $\mathcal{D} = (x_1, x_2, \dots, x_M)$  be the input dataset. Let  $f_{\theta_l}$  be the attention block in layer  $l$  of a PLM. Then concentration strength is a function  $\text{Strength} : \mathcal{D} \rightarrow \mathbb{R}^+$

$$\text{Strength}((z, \mathcal{D}); \theta_l) = \frac{1}{|\mathcal{D}|} \sum_{x_i \in \mathcal{D}} \text{Concentration}(z \oplus x_i; \theta_l).$$

**Definition 3.3.** Let  $\mathcal{D} = (x_1, x_2, \dots, x_M)$  be the set of textual sequences sampled from target domain  $\mathcal{D}_{Tt}$ , where  $x_i \in S$ . Then the concentration fluctuation is a function  $\text{Fluctuation} : \mathcal{D} \rightarrow \mathbb{R}^+$

$$\text{Fluctuation}((z, \mathcal{D}); \theta_l) = \sqrt{\frac{1}{|\mathcal{D}|} \sum_{x_i \in \mathcal{D}} [\text{Concentration}(z \oplus x_i; \theta_l) - \text{Strength}((z, \mathcal{D}); \theta_l)]^2}.$$

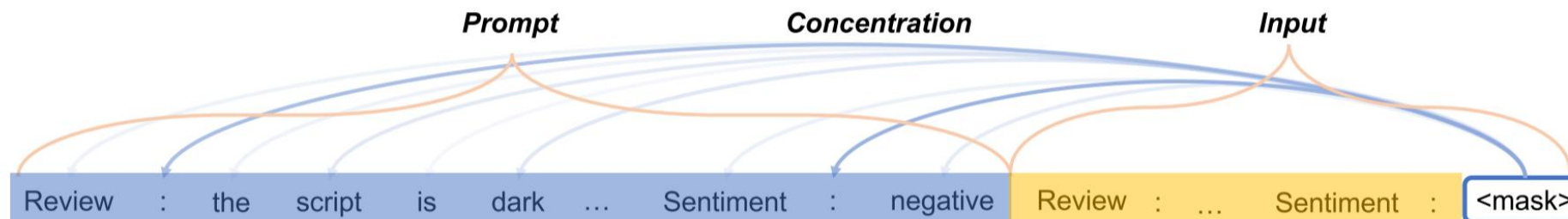


Figure 2: Illustration of Concentration.



# The Proposed Method: Pilot Experiment Results

$\mathcal{F}_1$ : Prompts gaining *more attention weight* from PLMs' deep layers are more generalizable.  
 $\mathcal{F}_2$ : Prompts with *more stable attention distributions* in PLMs' deep layers generalize better.

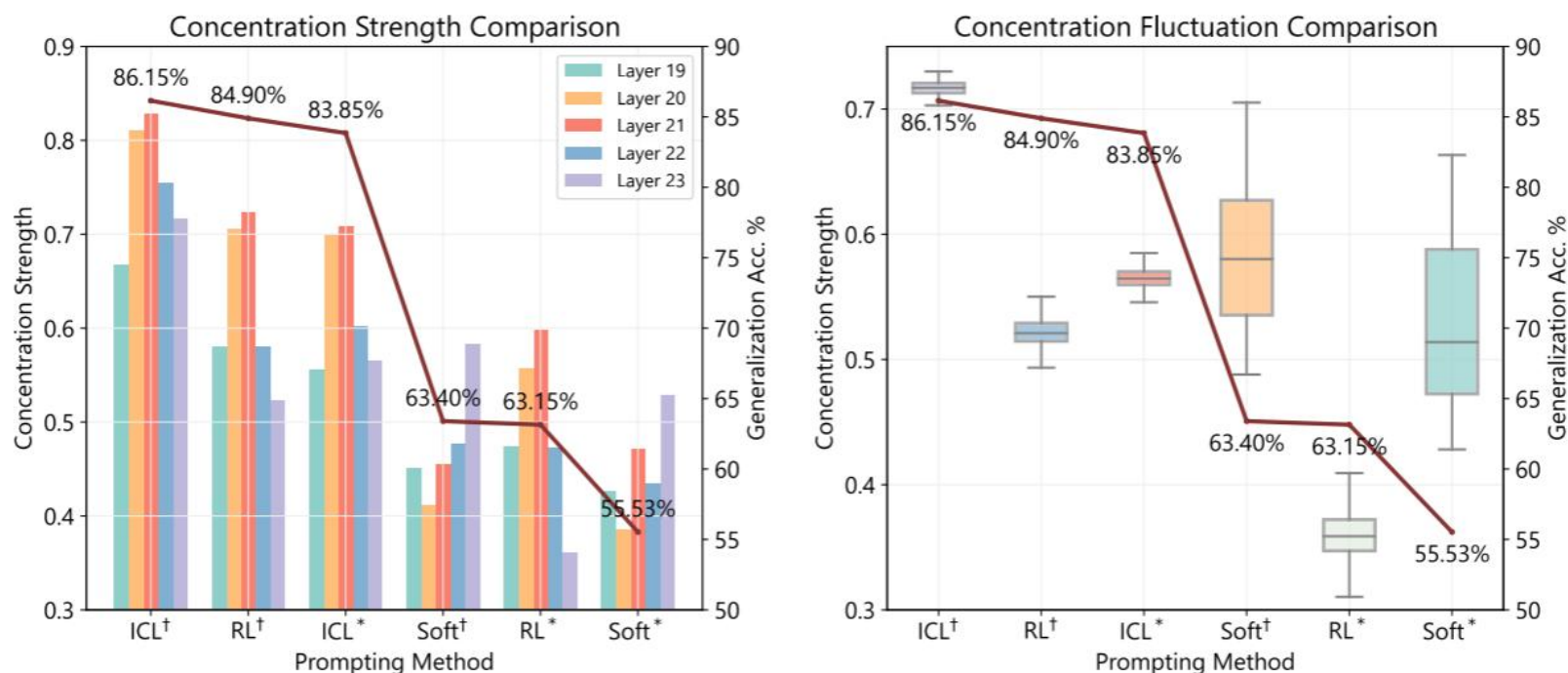


Figure 3. Left: concentration strength of various prompting methods in the last 5 layers (layers 19 to 23). Right: boxplots of the concentration strength in the last layer. Overall, prompts that exhibit good domain generalization gain higher concentration strength and lower concentration fluctuation.

# The Proposed Method: Concentrative Soft Prompt Optimization

These methods optimize follow loglikelihood objective given a trainable prompt  $z$  and a fixed PLM parameterized by  $\theta$  for the input  $x$ :

$$\max_z \log P(y|(z \oplus x); \theta).$$

According to our findings in §3, domain-generalizable prompts should be high in concentration strength and low in concentration fluctuation. Thus, we reformulate Eq. 6 to get the objective for domain-generalizable prompts:

$$\max_z (\log P(y|(z \oplus x); \theta) + \text{Strength}((z, \mathcal{D}_{\text{train}}); \theta)) \quad s.t. \quad \min_z \text{Fluctuation}((z, \mathcal{D}_{\text{train}}); \theta).$$

Towards the reformulated objective above, we propose the concentration-reweighting loss for soft prompt optimization methods. First, we minimize the concentration strength on input  $x$  to improve concentration strength on prompt  $z$  by designing loss function  $\mathcal{L}_{cs}$ . In addition, to reduce concentration fluctuation of prompts, we propose to use every token's concentration strength as hidden state feature of prompts, denoted as  $C_i = (c_1, c_2, \dots, c_L)$  where  $L$  is the length of prompts. We design a contrastive loss to cluster  $C$  with same label together to reduce concentration fluctuation:

$$\mathcal{L}_{cf} = \sum_{i=1}^{|\mathcal{D}_{\text{train}}|} \frac{-1}{P(i)} \sum_{p \in P(i)} \log \frac{\exp(\text{sim}(C_i, C_p)/\tau)}{\sum_{j=1}^{|\mathcal{D}_{\text{train}}|} \mathbf{1}_{i \neq j} \exp(\text{sim}(C_i, C_j)/\tau)}, \quad \mathcal{L}_{cs} = 1 - \text{Strength}((z, \mathcal{D}_{\text{train}}); \theta).$$

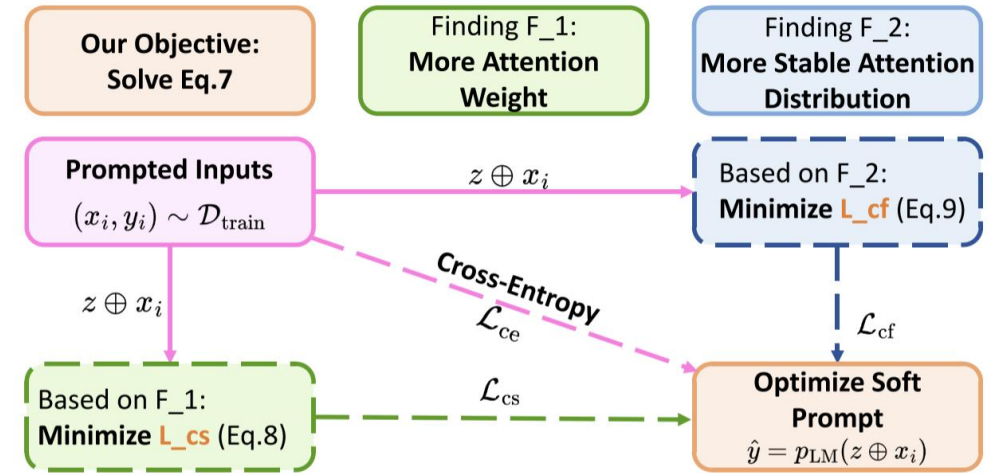


Figure 4. Framework for Soft Prompt Optimization.



# The Proposed Method: Concentrative Hard Prompt Optimization

In contrast to soft prompt optimization, hard prompt optimization searches suitable prompt in discrete space in a non-parameterized fashion. We focus on improving the generalization ability of input-level optimization methods. Generally, the mainstream of input-level optimization technique for hard prompts could be encapsulated as: filter (by metric) and match (by RL agents). The findings of concentration could be applied to this optimization process by adjusting filter metric and agent reward. We illustrate the framework for hard prompt optimization in Figure 5.

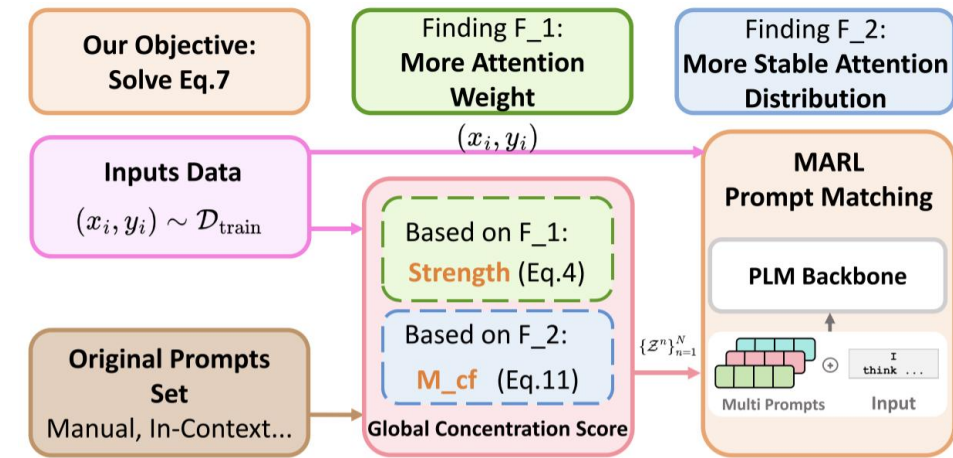


Figure 5. Framework for Hard Prompt Optimization.

For previous filter metric only considering the overall accuracy on training set, we introduce a new metric called Global Concentration Score (GCS), which involves our ideas of concentration strength and concentration fluctuation. We use concentration strength as first metric to filter out prompts which cannot get much concentration from model. Metric for reducing concentration fluctuation could be regarded as minimizing Kullback-Leibler (KL) divergence between the concentration features  $C_i$  of input with same label and the average of  $C_i$  on whole inputs set  $\mathcal{D}_{train}$ :

$$M_{cf}(z, \mathcal{D}_{train}) = \sum_{y \in \mathcal{Y}} \sum_{i \in \mathcal{D}_{train}(y)} \text{KL}(\text{Softmax}(C_i) \parallel \text{Softmax}(C_{avg}^y))$$

# The Proposed Method: Experiment Results

Paradigms	Methods	Sentiment			NLI		
		S+M→C	C+M→S	S+C→M	Q+R→W	W+R→Q	Q+W→R
Prompt Tuning Lester et al. [2021]	Vanilla PT	64.73 <sub>3.82</sub>	65.51 <sub>2.65</sub>	65.12 <sub>3.85</sub>	41.20 <sub>1.55</sub>	49.83 <sub>1.47</sub>	49.66 <sub>1.67</sub>
	PT with $\mathcal{L}_{cs}$	65.83 <sub>3.83</sub>	66.38 <sub>2.42</sub>	65.33 <sub>2.37</sub>	41.53 <sub>1.56</sub>	49.60 <sub>2.37</sub>	49.43 <sub>1.41</sub>
	PT with $\mathcal{L}_{cf}$	65.09 <sub>3.72</sub>	67.40 <sub>2.43</sub>	65.40 <sub>2.33</sub>	42.17 <sub>2.03</sub>	49.22 <sub>1.59</sub>	49.73 <sub>1.31</sub>
	PT with both	<b>66.19</b> <sub>3.69</sub>	<b>69.54</b> <sub>2.52</sub>	<b>65.89</b> <sub>2.32</sub>	<b>42.48</b> <sub>1.72</sub>	<b>50.31</b> <sub>1.33</sub>	<b>50.42</b> <sub>1.34</sub>
Prefix Tuning Li and Liang [2021]	Vanilla Prefix	65.91 <sub>3.24</sub>	83.25 <sub>0.41</sub>	75.51 <sub>0.91</sub>	50.26 <sub>0.31</sub>	51.88 <sub>0.29</sub>	50.02 <sub>0.28</sub>
	Prefix with $\mathcal{L}_{cs}$	66.23 <sub>3.37</sub>	84.32 <sub>0.48</sub>	76.58 <sub>0.82</sub>	50.69 <sub>0.33</sub>	51.44 <sub>0.28</sub>	49.77 <sub>0.22</sub>
	Prefix with $\mathcal{L}_{cf}$	66.82 <sub>3.19</sub>	83.70 <sub>0.39</sub>	77.17 <sub>0.75</sub>	51.73 <sub>0.32</sub>	52.12 <sub>0.28</sub>	50.73 <sub>0.26</sub>
	Prefix with both	<b>68.29</b> <sub>2.97</sub>	<b>85.07</b> <sub>0.42</sub>	<b>77.53</b> <sub>0.43</sub>	<b>52.05</b> <sub>0.30</sub>	<b>53.32</b> <sub>0.25</sub>	<b>51.26</b> <sub>0.27</sub>
P-Tuning v2 Liu et al. [2021]	Vanilla Pv2	65.92 <sub>1.61</sub>	83.84 <sub>1.69</sub>	75.89 <sub>0.36</sub>	50.63 <sub>0.31</sub>	52.76 <sub>0.01</sub>	51.31 <sub>1.37</sub>
	Pv2 with $\mathcal{L}_{cs}$	66.06 <sub>1.77</sub>	83.32 <sub>1.59</sub>	75.07 <sub>0.35</sub>	51.37 <sub>0.37</sub>	50.93 <sub>0.92</sub>	50.20 <sub>1.30</sub>
	Pv2 with $\mathcal{L}_{cf}$	66.72 <sub>1.62</sub>	84.12 <sub>1.51</sub>	76.41 <sub>0.33</sub>	51.32 <sub>0.38</sub>	52.64 <sub>1.04</sub>	51.28 <sub>1.22</sub>
	Pv2 with both	<b>67.07</b> <sub>1.53</sub>	<b>84.86</b> <sub>1.42</sub>	<b>77.26</b> <sub>0.37</sub>	<b>51.87</b> <sub>0.28</sub>	<b>53.83</b> <sub>0.95</sub>	<b>51.57</b> <sub>1.16</sub>
GrIPS Prasad et al. [2022]	-	80.07 <sub>2.57</sub>	84.28 <sub>1.38</sub>	85.19 <sub>1.12</sub>	54.37 <sub>2.40</sub>	52.77 <sub>1.73</sub>	53.52 <sub>1.66</sub>
RLPrompt Deng et al. [2022]	-	86.05 <sub>1.32</sub>	89.36 <sub>0.91</sub>	85.95 <sub>1.90</sub>	52.77 <sub>2.82</sub>	53.82 <sub>2.34</sub>	54.63 <sub>1.39</sub>
Manual Prompt Bach et al. [2022]	Vanilla MP	52.73 <sub>4.43</sub>	55.81 <sub>3.31</sub>	50.85 <sub>1.58</sub>	41.70 <sub>1.17</sub>	50.80 <sub>0.84</sub>	51.60 <sub>1.50</sub>
	MP with MARL	56.37 <sub>1.18</sub>	58.42 <sub>0.46</sub>	52.15 <sub>0.49</sub>	44.27 <sub>1.02</sub>	51.36 <sub>0.84</sub>	52.18 <sub>1.23</sub>
	MP with Metric	54.63 <sub>2.12</sub>	57.84 <sub>1.65</sub>	51.79 <sub>1.75</sub>	42.86 <sub>0.94</sub>	51.02 <sub>0.68</sub>	52.03 <sub>1.14</sub>
	MP with both	<b>56.76</b> <sub>0.40</sub>	<b>59.44</b> <sub>0.32</sub>	<b>53.15</b> <sub>0.35</sub>	<b>45.05</b> <sub>0.28</sub>	<b>52.03</b> <sub>0.25</sub>	<b>52.46</b> <sub>1.24</sub>
In-Context Demo Brown et al. [2020]	Vanilla IC	84.33 <sub>2.15</sub>	84.81 <sub>1.39</sub>	80.21 <sub>2.17</sub>	50.86 <sub>1.28</sub>	52.63 <sub>0.94</sub>	58.04 <sub>2.23</sub>
	IC with MARL	85.33 <sub>5.03</sub>	87.02 <sub>2.74</sub>	82.14 <sub>1.65</sub>	52.82 <sub>3.29</sub>	53.75 <sub>1.32</sub>	59.87 <sub>2.07</sub>
	IC with Metric	84.70 <sub>3.17</sub>	85.10 <sub>2.12</sub>	82.60 <sub>4.91</sub>	51.19 <sub>4.80</sub>	52.72 <sub>4.31</sub>	59.46 <sub>4.64</sub>
	IC with both	<b>87.29</b> <sub>2.72</sub>	<b>88.49</b> <sub>1.52</sub>	<b>83.52</b> <sub>0.98</sub>	<b>52.94</b> <sub>1.59</sub>	<b>54.24</b> <sub>0.73</sub>	<b>60.32</b> <sub>1.20</sub>
DP <sub>2</sub> O Li et al. [2024]	Vanilla DP <sub>2</sub> O	89.06 <sub>0.76</sub>	90.75 <sub>0.91</sub>	86.53 <sub>0.80</sub>	54.84 <sub>0.62</sub>	54.85 <sub>0.37</sub>	59.78 <sub>0.79</sub>
	DP <sub>2</sub> O with MARL	87.36 <sub>3.17</sub>	91.60 <sub>2.39</sub>	86.03 <sub>4.03</sub>	54.71 <sub>2.21</sub>	53.13 <sub>1.97</sub>	60.62 <sub>3.47</sub>
	DP <sub>2</sub> O with Metric	86.79 <sub>1.32</sub>	90.13 <sub>1.07</sub>	86.60 <sub>0.83</sub>	53.21 <sub>1.16</sub>	54.02 <sub>0.79</sub>	60.54 <sub>1.47</sub>
	DP <sub>2</sub> O with both	<b>89.63</b> <sub>0.52</sub>	<b>92.87</b> <sub>0.33</sub>	<b>87.85</b> <sub>0.47</sub>	<b>56.42</b> <sub>0.36</sub>	<b>55.32</b> <sub>0.33</sub>	<b>61.27</b> <sub>0.81</sub>

Table 1. Performance comparison of text classification tasks in accuracy with MFDG setting.

Paradigms	Methods	Sentiment				NLI				
		SST-2	CR	MR	Avg.	WNLI	QNLI	RTE	Avg.	Avg Gap
Prompt Tuning	Vanilla PT	73.84 <sub>3.52</sub>	75.89 <sub>1.72</sub>	74.17 <sub>2.32</sub>	74.63	47.64 <sub>1.02</sub>	49.71 <sub>0.93</sub>	54.73 <sub>1.72</sub>	50.69	+6.66
	PT with both	72.61 <sub>2.72</sub>	76.07 <sub>2.24</sub>	74.37 <sub>2.12</sub>	74.35	46.79 <sub>1.52</sub>	49.50 <sub>0.98</sub>	54.21 <sub>1.49</sub>	50.17	+4.71
Prefix Tuning	Vanilla Prefix	87.39 <sub>2.98</sub>	77.37 <sub>0.79</sub>	82.65 <sub>0.65</sub>	82.47	55.88 <sub>0.37</sub>	60.27 <sub>0.44</sub>	54.82 <sub>0.31</sub>	56.99	+6.93
	Prefix with both	87.29 <sub>3.12</sub>	76.73 <sub>1.28</sub>	83.32 <sub>0.83</sub>	82.45	56.18 <sub>0.35</sub>	59.74 <sub>0.32</sub>	55.38 <sub>0.42</sub>	57.10	+5.19
P-Tuning v2	Vanilla Pv2	86.71 <sub>1.57</sub>	77.65 <sub>1.49</sub>	82.27 <sub>0.42</sub>	82.21	55.57 <sub>0.73</sub>	60.73 <sub>1.64</sub>	55.16 <sub>1.83</sub>	57.15	+6.29
	Pv2 with both	87.03 <sub>1.32</sub>	77.71 <sub>1.50</sub>	82.05 <sub>0.54</sub>	82.08	56.31 <sub>0.69</sub>	60.46 <sub>1.37</sub>	55.20 <sub>1.69</sub>	57.32	+5.38
Manual Prompt	Vanilla MP	61.62 <sub>3.42</sub>	57.75 <sub>2.92</sub>	53.13 <sub>2.33</sub>	57.50	44.27 <sub>2.80</sub>	53.42 <sub>0.98</sub>	52.63 <sub>0.60</sub>	50.11	+3.22
	MP with both	61.33 <sub>2.32</sub>	56.07 <sub>1.61</sub>	53.47 <sub>0.42</sub>	56.96	44.05 <sub>0.89</sub>	53.77 <sub>1.35</sub>	52.60 <sub>0.39</sub>	50.14	+0.40
In-Context Demo	Vanilla IC	85.91 <sub>1.42</sub>	85.57 <sub>0.92</sub>	83.75 <sub>1.39</sub>	85.08	52.37 <sub>1.45</sub>	53.42 <sub>0.72</sub>	59.73 <sub>0.81</sub>	55.17	+1.65
	IC with both	86.33 <sub>1.34</sub>	85.14 <sub>0.87</sub>	84.31 <sub>2.12</sub>	85.26	52.25 <sub>1.62</sub>	52.96 <sub>0.49</sub>	59.36 <sub>0.73</sub>	54.86	+0.93
DP <sub>2</sub> O	Vanilla DP <sub>2</sub> O	93.62 <sub>0.72</sub>	90.76 <sub>0.50</sub>	88.58 <sub>0.91</sub>	90.99	55.26 <sub>1.02</sub>	55.13 <sub>0.39</sub>	61.07 <sub>0.81</sub>	57.15	+1.44
	DP <sub>2</sub> O with both	93.20 <sub>0.81</sub>	90.38 <sub>0.47</sub>	88.37 <sub>2.12</sub>	90.65	56.47 <sub>0.41</sub>	55.42 <sub>0.79</sub>	61.29 <sub>0.63</sub>	57.73	+0.37

Table 2. In-domain comparison.

Models	Methods	MCQA				
		S + C → R	C + R → S	R + S → C	Avg.	Acc Gap
Llama-2-7b-chat	Vanilla Prefix	62.32 <sub>2.15</sub>	66.30 <sub>2.30</sub>	73.15 <sub>2.53</sub>	67.26	—
	Prefix with both	63.70 <sub>1.96</sub>	68.47 <sub>0.97</sub>	75.32 <sub>1.09</sub>	69.16	+1.90
	Vanilla IC	63.13 <sub>1.25</sub>	65.50 <sub>1.98</sub>	77.59 <sub>1.14</sub>	68.74	—
	IC with both	65.13 <sub>1.03</sub>	68.33 <sub>2.13</sub>	79.83 <sub>0.88</sub>	70.10	+1.36
Vicuna-7b-v1.5	Vanilla Prefix	67.72 <sub>1.79</sub>	81.09 <sub>2.17</sub>	88.97 <sub>2.64</sub>	79.26	—
	Prefix with both	68.75 <sub>1.04</sub>	83.93 <sub>1.79</sub>	89.76 <sub>2.60</sub>	80.81	+1.55
	Vanilla IC	68.37 <sub>2.24</sub>	83.23 <sub>4.12</sub>	90.98 <sub>1.99</sub>	80.86	—
	IC with both	69.67 <sub>1.58</sub>	85.50 <sub>5.06</sub>	93.39 <sub>1.23</sub>	82.85	+1.99
Alpaca-7b-wdiff	Vanilla Prefix	61.52 <sub>3.79</sub>	70.03 <sub>2.88</sub>	87.91 <sub>2.73</sub>	73.15	—
	Prefix with both	63.89 <sub>2.93</sub>	72.15 <sub>2.07</sub>	89.58 <sub>2.81</sub>	75.21	+2.06
	Vanilla IC	60.81 <sub>1.14</sub>	69.11 <sub>2.46</sub>	89.66 <sub>2.37</sub>	73.19	—
	IC with both	63.16 <sub>1.56</sub>	70.57 <sub>1.95</sub>	91.19 <sub>2.00</sub>	74.97	+1.78

Table 3: Performance comparison of large models on MCQA task accuracy.

# Conclusion & Future work

- To explore the nature of prompt generalization on unknown domains, we conduct pilot experiments and find that (i) Prompts gaining more attention weight from PLMs' deep layers are more generalizable and (ii) Prompts with more stable attention distributions in PLMs' deep layers are more generalizable.
- We adapt this new objective to popular soft prompt and hard prompt optimization methods, respectively. Extensive experiments demonstrate that our idea improves comparison prompt optimization methods by 1.42% for soft prompt generalization and 2.16% for hard prompt generalization in accuracy, while maintaining satisfying in-domain performance.
- Future improvement:
  - Try to systematically study more cues and attention distribution phenomena;
  - Try to apply this objective to more complex downstream tasks (e.g., open-ended generation tasks);

# Thanks!



Paper



Contact



Code