

# Toward the Explainable Soft Prompts: How does Prompt-tuning Exploit a Multilingual Pre-trained Language Model?

Anonymous ACL submission

## Abstract

Since training soft prompts is a parameter-efficient way to tune a Pre-trained Language Model (PLM) on a target task, recent works suggest various training methods utilizing soft prompts. However, few studies investigate the explainability of soft prompts, which is critical to enhancing the confidence of PLM in a real-world scenario. To deal with the problems of the unexplainable soft prompts, this study explores the effects of Prompt-tuning v1 (Lester et al., 2021) and Prompt-tuning v2 (Liu et al., 2022) on PLM. More precisely, we conducted the experiments using a multilingual GPT to generalize our observations not only on tasks but also on languages. We first confirmed whether soft prompts are gathered according to tasks or languages, and then analyzed how soft prompts utilize PLM in terms of the two main architectures of GPT: the attention mechanism and the activated neurons. As a result, we conclude that soft prompts are trained while employing the knowledge PLM obtained during pre-training to solve the target task, which is consistent with language. Our findings reveal that deep soft prompts are explainable when they can employ varying knowledge from every layer.

## 1 Introduction

In Natural Language Processing (NLP), prompt-based learning has emerged as a promising paradigm to employ a Pre-trained Language Model (PLM) effectively. Particularly, **Prompt-tuning v1** (Lester et al., 2021) shows that training only **soft prompts**, which are vectors of real numbers prepended to the embedding layer of PLM, can

achieve comparable performances to traditional fine-tuning. Additionally, soft prompts are fed to all layers of PLM in **Prompt-tuning v2** (Liu et al., 2022), where **deep soft prompts** encode more task-specific weights and affect the model directly.

Despite the great promise of the *pre-train and prompt-tune* paradigm, it is problematic that the difficulty of faithfully interpreting soft prompts in natural language could potentially lead to concealed adversarial attacks (Khashabi et al., 2022). Ultimately, it makes PLM unexplainable in that we cannot assure how soft prompts operate in PLM. However, few studies have tried to reveal the details of the relationship between Prompt-tuning and PLM, which are essential to the explainable soft prompts.

Meanwhile, several studies have shown that soft prompts are one of the parameter-efficient training methods in a cross-lingual setting (Zhao and Schütze, 2021; Vu et al., 2022), which means that soft prompts encode the knowledge about a target task in the multilingual space. Since the generalizability of tasks and languages can improve the explainability of PLM, it is important to capture the features shared in each task and language. Accordingly, in this study, we investigate how soft prompts save and utilize the information of a multilingual PLM, mGPT (Shliazhko et al., 2022), focusing on English and Korean.

This study also aims to figure out the effects of Prompt-tuning v1 and Prompt-tuning v2 on the multilingual PLM, so that we provide fundamental directions to enhance the interpretability and explainability of PLM and Prompt-tuning. To analyze the effects of soft prompts, we focus on the changes after Prompt-tuning in terms of two major structures of GPT: the attention mechanism and the activated neurons.

78 To this end, we address the following research  
79 questions:

- 80  
81 1. Can we distinguish soft prompts according  
82 to the encoded information about target  
83 tasks or target languages?  
84
- 85 2. Can we find any explainable patterns in the  
86 changes in the attention mechanism after  
87 Prompt-tuning?  
88
- 89 3. Can we observe any explainable patterns in  
90 the activated neurons of soft prompts  
91 through layers?  
92

93 We first present the visualizations of soft  
94 prompts to show that deep soft prompts are  
95 gathered according to target tasks in the  
96 multilingual setting. Second, we find that the  
97 attention distribution changes in some layers more  
98 than in other layers regardless of tasks and  
99 languages. Also, in Prompt-tuning v2, the changes  
100 are explainable because the most changed attention  
101 layers are composed of content-dependent heads  
102 rather than position-based heads. These results  
103 suggest that soft prompts utilize the knowledge  
104 encoded in the attention mechanism to solve the  
105 target task. Third, we observe a special  
106 phenomenon, where the activated neurons show  
107 task-common behavior rather than task-specific  
108 behavior in the second to last layer. We conclude  
109 that this may be due to the characteristics of the  
110 PLM, like isotropy. Finally, the ablation study  
111 supports these findings as where deep soft prompts  
112 fed into the actively changed layers show higher  
113 performances than the ones fed into the non-  
114 actively changed layers. To the best of our  
115 knowledge, this study is the first to probe the  
116 changes in the PLM after Prompt-tuning v1 and  
117 Prompt-tuning v2.

## 118 2 Related Works

119 **Soft prompts.** Li and Liang (2021) propose Prefix-  
120 tuning for lightweight model training. They  
121 prepend the continuous task-specific vectors to the  
122 input so that their own parameters are only trained.  
123 As the continuous vectors do not limit their space  
124 on the embedding of PLM, they are more  
125 expressive and can affect all layers in PLM. Using  
126 GPT-2 (Radford et al., 2019) and BART (Lewis et  
127 al., 2020), they show that the performances on  
128 Natural Language Generation (NLG) tasks are

129 comparable to the performances of fine-tuning  
130 both in fully- and semi-supervised settings.

131 Liu et al. (2021) suggests P-tuning which is  
132 more flexible to task types and LM types than  
133 Prefix-tuning. They use a discrete template that  
134 include soft prompts and train both PLM and those  
135 soft prompts. They report that performances  
136 improve both in GPT and BERT via P-tuning.

137 Similarly, Lester et al. (2021) propose Prompt-  
138 tuning v1, where they use only soft prompts by  
139 freezing the parameters of PLM. They show that  
140 the performances of Prompt-tuning v1 are  
141 comparable to the performances of fine-tuning,  
142 where the capacity of PLM has a key role. Thus,  
143 they conclude that Prompt-tuning is a parameter-  
144 efficient way to employ the knowledge LM  
145 obtained during pre-training.

146 Furthermore, Liu et al. (2022) introduce Prompt-  
147 tuning v2, where they implement Prompt-tuning  
148 more deeply by injecting soft prompts into every  
149 layer of PLM. While Prompt-tuning v1 has low  
150 performances on the hard sequence labeling tasks,  
151 Prompt-tuning v2 improves the performances of  
152 various tasks, namely extracting question  
153 answering and named entity recognition.

154 **The explainable soft prompts.** Again, few studies  
155 endeavor to enhance the explainability of soft  
156 prompts. Lester et al. (2021) try to interpret soft  
157 prompts by measuring the similarities between the  
158 embeddings of learned soft prompts and the  
159 vocabulary of PLM. They observe that soft  
160 prompts have ‘word-like’ representations, which  
161 are relevant to the domain of the target task. An  
162 example of which is the BoolQ dataset (Clark et al.,  
163 2019) of the nature/science category, where the soft  
164 prompts are close to the words such as ‘*science*’,  
165 ‘*technology*’ and ‘*engineering*’ in the embedding  
166 space.  
167

168 To interpret soft prompts in human language,  
169 Khashabi et al. (2022) investigate the Prompt  
170 Waywardness hypothesis. In this hypothesis, there  
171 exists a soft prompt that can solve the target task  
172 while becoming close to the arbitrary discrete  
173 prompt, which is not relevant to the target task.  
174 They observe that soft prompts satisfying the  
175 Prompt Waywardness hypothesis do exist. They  
176 also provide some explanations. First, soft prompts  
177 cannot be projected to exactly one embedding of  
178 discrete prompts. Second, when soft prompts are  
179 injected only into the first layer, the deeper layers  
180 have more expressivity, where the effects of

Task	English	Korean
STS	GLUE-STS	KLUE-STS
NLI	SNLI	KorNLI
SA	SST2	NSMC
TC	AGNews	KLUE-Ynat
QA	SQuAD 2.0	KorQuAD 1.0
CG	CommonGen	KommonGen

Table 1: The datasets used in this study.

Waywardness get stronger. Finally, they discuss that it is hard to discover the human-interpretable soft prompts, which leads to the side effects in the real-world scenario, such as the concealed adversarial attacks.

Besides the interpretability of soft prompts, the changes that soft prompts cause to PLM are also crucial to enhance the explainability of soft prompts. Also, to deal with the second reason [Khashabi et al. \(2022\)](#) mentioned, we need to investigate Prompt-tuning v2 where all layers are controlled by soft prompts. Thus, our study focuses on the effects of Prompt-tuning v1 and Prompt-tuning v2 on PLM through layers, so that we can figure out the operation of soft prompts in an explainable way.

### 3 Methods

#### 3.1 Training soft prompts

We conduct the experiments using mGPT, which has the same architecture with GPT-3 and 1.3B parameters. Since GPT is an autoregressive model, we make the model generate the label words after the separator token  $\langle /s \rangle$  until the end token  $\langle |endoftext| \rangle$  and train soft prompts by computing the conditional probability for the label words. The label words for each task are in Table 3 in Appendix 0. Via this method, we can train soft prompts in the same way for all tasks, including classification and generation tasks. Also, we freeze mGPT and update only the parameters of the prepended soft prompts.

For a single-input type, we feed the input  $\{p_1, \dots, p_k, \langle s \rangle, x_1, \dots, x_n, \langle /s \rangle, w_{gold}, \langle |endoftext| \rangle\}$  to the model, where  $P = \{p_1, \dots, p_k\}$  is soft prompts with the length  $k$ ,  $X = \{x_1, \dots, x_n\}$  is the input sentence with the length  $n$ , and  $w_{gold}$  is the label words. Then, following [Lester et al. \(2021\)](#), we train soft prompts by maximizing the probability:

$$Pr_{\theta; \theta_p}(w_{gold}, \langle |endoftext| \rangle | p_1, \dots, p_k; \langle s \rangle, x_1, \dots, x_n, \langle /s \rangle) \quad (1)$$

where  $\theta$  is the parameters of the model and  $\theta_p$  is the parameters of soft prompts.

Similarly, for a pair-input type, the model maximizes the probability to generate the label words after the second separator token:

$$Pr_{\theta; \theta_p}(w_{gold}, \langle |endoftext| \rangle | p_1, \dots, p_k; \langle s \rangle, x_1^1, \dots, x_n^1, \langle /s \rangle, x_1^2, \dots, x_m^2, \langle /s \rangle) \quad (2)$$

For both Prompt-tuning v1 and Prompt-tuning v2, we randomly initialize soft prompts ranging from -0.5 to 0.5. Also, we use the prompt length  $k = 20$ . Additionally, we slightly modify the injection of deep soft prompts in Prompt-tuning v2. Unlike [Liu et al., 2022](#), we inject *key* and *value* as the same one, so that we get only one prompt embedding per layer, which has the exact same dimension with the model. Otherwise, soft prompts are composed of two separate embeddings.

Using principal component analysis (PCA), we visualize soft prompts of each target task to see how soft prompts are clustered. In the case of deep soft prompts, we analyze soft prompts for each layer.

#### 3.2 The Attention Variability and the Changes in the Attention Mechanism

In GPT, the attention mechanism works in the left-to-right direction, so most tokens give the maximum attention to their previous token. [Vig and Belinkov \(2019\)](#) consider that such a tendency of the attention distribution in GPT is based on not the content but the position. To measure how attention varies over different input sequences, they suggest attention variability, which adopts the basics of the mean absolute deviation:

$$Variability_\alpha = \frac{\sum_{x \in X} \sum_{i=1}^{|x|} \sum_{j=1}^i |\alpha_{i,j}(x) - \bar{\alpha}_{i,j}|}{2 \cdot \sum_{x \in X} \sum_{i=1}^{|x|} \sum_{j=1}^i \alpha_{i,j}(x)} \quad (3)$$

where  $\alpha_{i,j}(x)$  is the attention score  $x_i$  gives to  $x_j$ , and  $\bar{\alpha}_{i,j}$  is the mean of  $\alpha_{i,j}(x)$  overall  $x \in X$ .

Meanwhile, the first token of each input sequence tends to receive the maximum attention score. Thus, the first token is excluded to calculate the variability. We also compute the variability with

266 the first  $N$  tokens ( $N = 10$ ), excluding the first  
 267 token.<sup>1</sup>

268 Since we do not update the parameters of mGPT,  
 269 we believe that the changes in the attention  
 270 mechanism after Prompt-tuning show the effects of  
 271 Prompt-tuning on mGPT. Thus, we investigate the  
 272 difference of the attention distribution between the  
 273 pre-trained mGPT and the prompt-tuned mGPT.  
 274 We use Kullback-Leibler divergence which is  
 275 commonly used to measure the difference between  
 276 two probability distributions:

$$278 \quad KLdiv(P(X), M(Y)) = \sum_{x \in X, y \in Y} x \log \frac{x}{y} \quad (4)$$

279 where  $P(\cdot)$  is a prompt-tuned model and  $M(\cdot)$  is a  
 280 pre-trained PLM.

281 To get the KL divergence per head, we post-  
 282 process the attention distributions. We first sum the  
 283 attention scores for each token, and then remove  
 284 the attention score of the special tokens  $\langle s \rangle$  and  
 285  $\langle /s \rangle$  in each input. In the case of  $P(\cdot)$ , we exclude  
 286 the attention scores of the tokens of soft prompts  
 287 before post-processing. Finally, we replace  $x$  and  $y$   
 288 with  $1e-20$  when they are 0 after passing the  
 289 softmax to avoid the infinity when computing KL  
 290 divergence.

291 Lastly, we get the correlation between the  
 292 attention variability and the attention changes  
 293 through layers to see whether we can interpret the  
 294 patterns of the changes in terms of the knowledge  
 295 mGPT learned during pre-training.

### 296 3.3 The Activated Neurons

297 Geva et al. (2021) suggest that the FFN layers in  
 298 pre-trained Transformer mirror the key-value  
 299 memories, since the neural memory and the FFN  
 300 layer have a similar structure (Equation 5 &  
 301 Equation 6). The difference between them is the  
 302 function applied; FFN has a non-linear activation  
 303 function, and the neural memory has the softmax  
 304 function.

$$306 \quad FFN(x) = f(x \cdot K^T) \cdot V \quad (5)$$

$$307 \quad MN = softmax(x \cdot k^T) \cdot V \quad (6)$$

308

<sup>1</sup> For commonsense generation task, we use first 2 tokens without the first token because all input sequences include only 3 words.

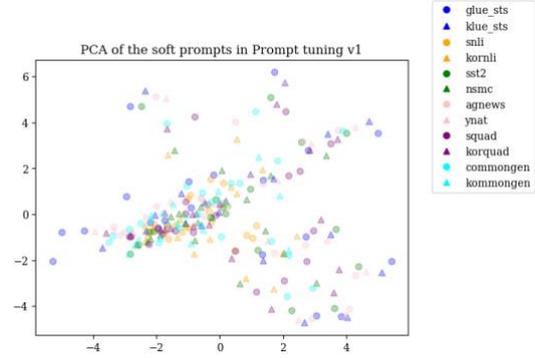


Figure 1: The PCA result of soft prompts learned for the target task in Prompt-tuning v1.

309 They observe that keys in the FFN layers, which  
 310 are the activated neurons, have positive  
 311 correlations with the human-interpretable input  
 312 patterns such as shallow and semantic patterns.

313 Motivated by the activated neurons, Su et al.  
 314 (2022) propose ON score, where we can measure  
 315 the response of the prompt-tuned model. Following  
 316 them, this study computes ON score using the  
 317 decoding token  $\langle /s \rangle$  per layer to investigate how  
 318 task-specific the prompts are. Given an input  
 319 sequence  $\{P, \langle /s \rangle\}$ , we get the output of the  
 320 activation state  $AS(P)$  of soft prompts  $P$ .

321

$$322 \quad ON(P^{t_1}, P^{t_2}) = \frac{AS(P^{t_1}) \cdot AS(P^{t_2})}{||AS(P^{t_1})|| ||AS(P^{t_2})||} \quad (7)$$

323

324 Recently, Wang et al. (2022) introduce skill  
 325 neurons that can be detected via Prompt-tuning.  
 326 They report that similar neurons are activated in  
 327 similar tasks in the upper layers, which means that  
 328 the neurons encode the task-specific skills. This  
 329 study also explores the response of the model  
 330 through layers to discuss which skills each layer  
 331 has.

## 332 4 Experiments

### 333 4.1 Tasks and Performance Results

334 We conduct experiments on 6 tasks including  
 335 classification, extraction, and generation.

336 **STS.** Semantic Textual Similarity is a task to  
 337 measure the similarity score between two  
 338 sentences from 1 to 5. The higher the score is, the  
 339 higher the similarity is. In this study, we binarize

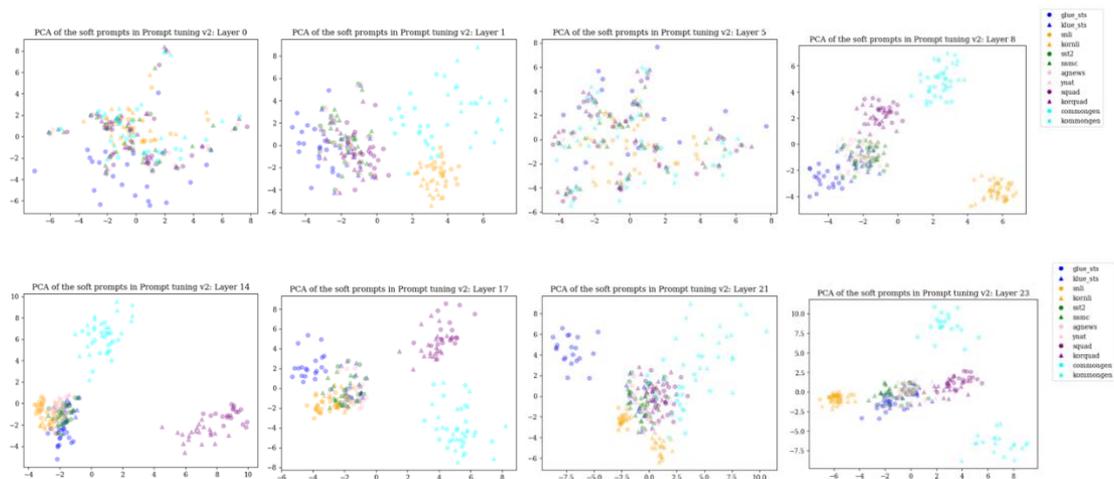


Figure 2: The PCA result of soft prompts learned for the target task in Prompt-tuning v2.

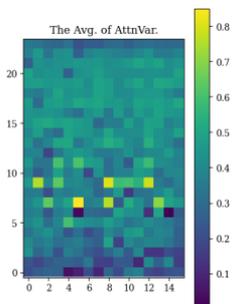


Figure 3: The average of the attention variability on all tasks. We present the results for each task in Figure 12 in Appendix C.

340 the score into similar and dissimilar class. If the  
 341 score is more than 3, it is labeled as similar, and if  
 342 the score is less than 3, it is labeled as dissimilar.

343 **NLI.** Natural Language Inference is a task to  
 344 decide the semantic relationship of two sentences,  
 345 a premise and a hypothesis, among three  
 346 categories: entailment, contradiction, neutral.  
 347 When given a premise first, the model should  
 348 determine if a hypothesis is true or false or  
 349 undetermined.

350 **SA.** Sentiment Analysis is a task to determine the  
 351 sentiment expressed in a sentence. Usually, the  
 352 sentence is classified into two labels: positive and  
 353 negative. In this study, we use the movie review  
 354 dataset, which is one of the popular domains.

355 **TC.** Topic Classification task includes sentences or  
 356 paragraphs from different themes. We use news  
 357 topic classification task, where the topic such as  
 358 *IT/science* and *Social* is annotated for each  
 359 headline. Also, the topics vary for each language.

360 **QA.** For Question and Answering for reading  
 361 comprehension, the model should extract the

362 answer to a question from the given context. The  
 363 answer can be a word or spans consisting of more  
 364 than two words.

365 **CG.** Commonsense Generation is a task proposed  
 366 to assess the model’s ability to generative  
 367 commonsense reasoning. Using a set of three  
 368 common concepts including an object (noun) and  
 369 action (verb), the model generates a full  
 370 grammatical sentence, which coheres with an  
 371 everyday scenario.

372 The datasets in each language are presented in  
 373 Table 1. Also, we report the performance results in  
 374 Table 7 in Appendix C. Every task gets improved  
 375 in Prompt-tuning v2. Particularly, the scores of  
 376 KLUE-STs and KorNLI rise by around 35 and 17,  
 377 respectively.

#### 378 4.2 Deep soft prompts encode task-relevant 379 information clustering by language.

380 We present the visualizations of soft prompts using  
 381 PCA to investigate how they encode the knowledge  
 382 according to task and language. While we cannot  
 383 find any meaningful clusters in Prompt-tuning v1  
 384 (Figure 1), soft prompts are grouped in Prompt-  
 385 tuning v2 through layers (Figure 2).

386 After passing the first layer (layer 0), soft  
 387 prompts tend to be gathered according to the target  
 388 task. Although soft prompts are dispersed again in  
 389 layer 5, they keep forming the clusters according to  
 390 the target task. The task clusters in different  
 391 languages are grouped, especially those with high  
 392 cohesion in the middle layers (7~17). Additionally,  
 393 starting from layer 14, we observe the separation of

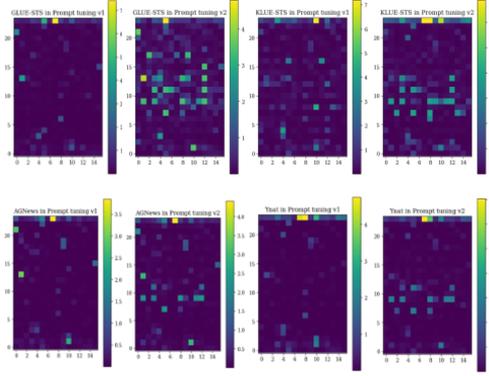


Figure 4: The KL divergence results on STS and TC. We present the results of all tasks in Figure 13 and Figure 14 in Appendix C.

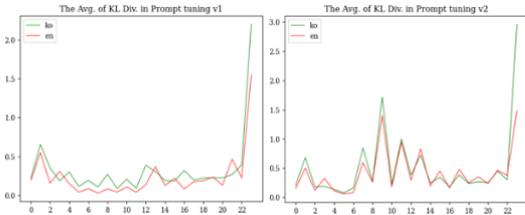


Figure 5: The average of KL divergence per layer.

394 some clusters (NLI, STS, TC) that are not well  
395 separated in the previous layers.

396 These results suggest that deep soft prompts  
397 learn task-relevant knowledge in the multilingual  
398 space since soft prompts of the same task in  
399 English and Korean are gathered. Also, the  
400 different patterns through layers imply that the  
401 knowledge employed to solve the target task varies  
402 over layers. Deep soft prompts controlling across  
403 layers might lead to low expressivity in the deeper  
404 layers, which means that it can contribute to the  
405 interpretability of soft prompts. We leave this  
406 implication as an open question for future works.

### 407 4.3 Deep soft prompts employ the context- 408 dependent attention heads.

409 We find that the lower layers have lower variability  
410 (Figure 3), which is consistent with the  
411 observations in the prior work (Vig and Belinkov,  
412 2019). Also, the attention heads in layers 7~9 have  
413 high attention variability regardless of tasks and  
414 languages. These results show that the attention  
415 heads of mGPT have encoded context information  
416 in a robust way during pre-training.

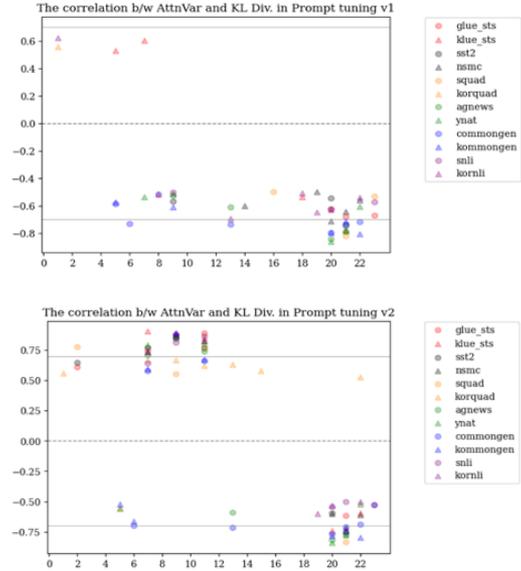


Figure 6: The correlation between the attention variability and the KL divergence.

417 Figure 4 displays the KL divergence results in  
418 Prompt-tuning v1 and Prompt-tuning v2 on STS  
419 and TC tasks. In Prompt-tuning v1, we observe that  
420 the attention distribution of the upper layers  
421 changes a lot, followed by the initial layers. In  
422 Prompt-tuning v2, the layers between 6 and 13  
423 change significantly as well, which consist of the  
424 content-dependent heads. Additionally, the results  
425 are consistent with the same task rather than with  
426 the same language.

427 To analyze the changes by layer, we plot the  
428 average KL divergence for each layer. Figure 5  
429 shows that the middle layers from 7 to 17 change  
430 remarkably in Prompt-tuning v2 with the peaks at  
431 the odd layers {7, 9, 11, 13, 15, 17}. On the other  
432 hand, in Prompt-tuning v1, the changes in the  
433 middle layers are relatively minor. Also, the peaks  
434 appear at different layers for each language; the  
435 even layers {6, 8, 10} for both languages, the even  
436 layers {12, 16} for Korean, and the odd layers {13,  
437 15} for English. This is because they do not have a  
438 direct impact on the deeper layers in PLM (Liu et  
439 al, 2022) since soft prompts are injected only in the  
440 embedding layer in Prompt-tuning v1.

441 Next, we present the results of the Pearson  
442 correlation (p-value < 0.05) between the attention  
443 variability and the results of KL divergence in each  
444 task (Figure 6). Grouping the layers into 4 groups  
445 according to the depth, we sum up the observations  
446 in the attention mechanism.

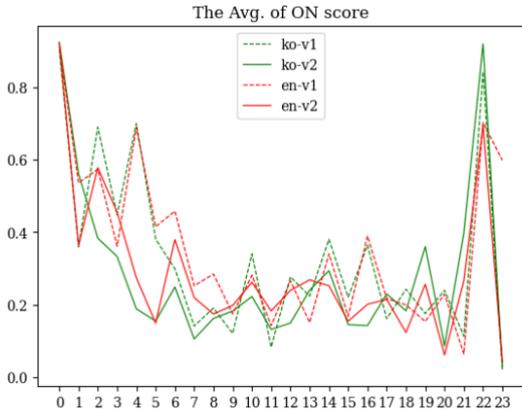


Figure 7: The average of ON score.

447 First, the lower layers (layers 0~5), where the  
 448 position-based heads are concentrated, show a  
 449 small change both in Prompt-tuning v1 and  
 450 Prompt-tuning v2. Second, the middle-lower layers  
 451 (layers 6~11) have different patterns in Prompt-  
 452 tuning v1 and Prompt-tuning v2. Considering that  
 453 the content-dependent heads are gathered in layers  
 454 7~9 (Figure 3), the positive correlations in Prompt-  
 455 tuning v2 suggest that deep soft prompts are trained  
 456 while employing the context information encoded in  
 457 mGPT. Third, the middle-upper layers (layers  
 458 12~17) do not show significant changes, yet, we  
 459 observe that the changes in the pair-input type tasks  
 460 in these layers are comparable to the changes in the  
 461 middle-lower layers (See Figure 15 in Appendix  
 462 C). Lastly, the upper layers (layers 18~23) have the  
 463 same pattern in both Prompt-tuning v1 and Prompt-  
 464 tuning v2, where the last layer shows a significant  
 465 change. Simultaneously, they have negative  
 466 correlations, which implies that the additional  
 467 elements are involved in activating the attention  
 468 head to predict the label in these layers.

469 **4.4 The activated neurons of soft prompts**  
 470 **are task-specific in the deeper layers,**  
 471 **where other features of PLM make the**  
 472 **neurons common as well.**

473 The results of the average ON score between all  
 474 combinations of tasks in each language are  
 475 reported in Figure 7, where the lower the score is,  
 476 the more task-specific the layer is. We find that the  
 477 scores in the first layer are high and the scores  
 478 reduce rapidly in the second layer, which means  
 479 that the first layer has encoded the common  
 480 information regardless of task and language.

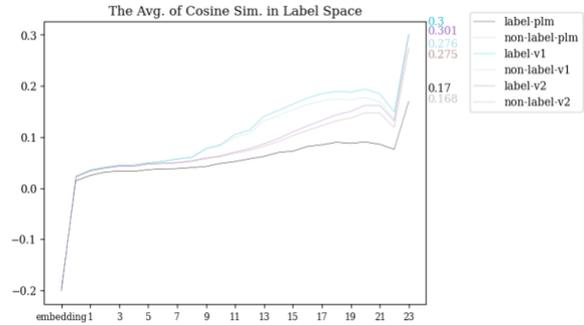


Figure 8: The cosine similarity of the decoding token between the label words and between the non-label words. If the label word is ‘positive’, the non-label word is ‘negative’. Also, if there are multiple non-label words, we randomly select one non-label word. To calculate the similarity, we use tasks without extraction and generation.

481 Meanwhile, the results in layers from 2 to 5 have  
 482 some fluctuations in Prompt-tuning v1, while there  
 483 are few fluctuations in Prompt-tuning v2. The  
 484 middle layers (layers 6~17) also have fluctuations,  
 485 but with more narrow gaps between layers.  
 486 Notably, the scores become the lowest in layer 21  
 487 in Prompt-tuning v1 and in layer 20 in Prompt-  
 488 tuning v2.

489 While these results suggest that the deeper layers  
 490 have task-specific neurons, the second to last layer  
 491 shows peaks with high scores near the scores in the  
 492 first layer, which has not been observed in previous  
 493 works. Since most tasks have highly similar  
 494 neurons in the second to last layer, we hypothesize  
 495 that other features of PLM affect the neurons when  
 496 solving the target tasks.

497 To test our intuitions, we present the cosine  
 498 similarity in the label space. Figure 8 illustrates that  
 499 the cosine similarities of the decoding token  $\langle /s \rangle$   
 500 with the gold label words get higher than ones with  
 501 the non-gold label words in the deeper layers.  
 502 Indeed, the second to last layer has the trough,  
 503 where the gaps between two similarities narrow  
 504 down. We believe that the anisotropic embedding  
 505 space of mGPT could lead the task-common  
 506 neurons in the second to last layer because the  
 507 desired label words are actually similar in the  
 508 space. Thus, we conclude that these observations  
 509 are one of the possible explanations as to why soft  
 510 prompts are hard to interpret.

Task	English		Korean	
	Peak	Trough	Peak	Trough
STS	<b>84.17</b>	83.74	60.32	<b>63.26</b>
NLI	<b>85.35</b>	84.85	<b>61.97</b>	58.18
SA	<b>90.02</b>	89.79	<b>86.73</b>	86.25
TC	<b>86.66</b>	86.46	<b>83.75</b>	82.66
QA	<b>64.40/</b>	61.82/	<b>68.82/</b>	66.57/
	<b>49.29</b>	47.11	<b>62.41</b>	60.42
CG	<b>82.56</b>	81.62	<b>90.43</b>	90.21

Table 2: The performance results of the compressed Prompt-tuning v2.

## 5 Ablation study

For the ablation study, we try to investigate whether deep soft prompts can better utilize the specific layers, where effects are notable.

To this end, we feed soft prompts to the layers, where changes and roles are clear in order to confirm the effects of Prompt-tuning v2 on the model. In Section 4.3, we find that the changes in the attention distribution show a zigzag trend line. The peaks are detected in layers {1, 7, 9, 11, 13, 15, 17, 21} and the troughs are detected in layers {2, 8, 10, 12, 14, 16, 18, 20, 22}, both in English and Korean. We hypothesize that soft prompts are more helpful for the layers at the peaks than the layers at the troughs.

Furthermore, looking at the results of the activated neurons in Section 4.4, the last layer (23) is a task-specific layer, and the second-last layer shows a special phenomenon, where the neurons of all tasks are similar to each other. Also, the neurons in layer 20 have the lowest ON score, which means that the neurons are activated depending on tasks.

Motivated by these observations, we group the layers into two categories: peak and trough. The peak group includes layers {0, 1, 7, 9, 11, 13, 15, 17, 20, 21, 22, 23}, and the trough group includes layers {0, 2, 8, 10, 12, 14, 16, 18, 20, 21, 22, 23}. We set the first layer {0} and the last four layers {20, 21, 22, 23} in common, in consideration that the former is the input layer, and the latter ones are task-specific layers. Each group includes half of the number of layers of mGPT, which means that they have half parameters of soft prompts in Prompt-tuning v2.

This method is similar to Liu et al. (2022). They conduct the ablation study of Prompt-tuning v2 by adding soft prompts to certain layers grouping into two groups, the first four layers and the last four layers. On the contrary, this study tries to classify the layers in an explainable way. We believe that the alternating dense and sparse attention mechanism of mGPT prevent the performances from the significant drop when excluding some layers.

Table 2 shows the performances of each group on all tasks. Compared to the results of vanilla Prompt-tuning v2, most scores drop but raise in comparison to the results of Prompt-tuning v1. Notably, the scores of both groups on SST2 raise compared to Prompt-tuning v2.

Except for KLUE-STS, the peak groups show higher performances than the trough groups. Thus, we conclude that soft prompts in the peak groups employ the knowledge of mGPT better.

## 6 Conclusion

In this study, we investigate how soft prompts and deep soft prompts encode task-relevant knowledge and employ the knowledge from the PLM. Using mGPT, we conduct the experiments on various tasks, including classification and generation in each language. First, we find that deep soft prompts obtain task-specific knowledge in the deeper layers, while soft prompts in Prompt-tuning v1 do not. Also, deep soft prompts have shared task information within languages. Second, we observe that the changes in the attention mechanism after Prompt-tuning v2 can be explained in terms of the attention variability. The higher the attention variability is where the more significant the changes are. Simultaneously, the last four layers show negative correlations between the attention variability and the changes in the attention distribution after Prompt-tuning v1 and Prompt-tuning v2. Thus, we conclude that these layers play another key role to solve the target task. Third, the response of the model shows that the deeper layers have more task-specific neurons. Unlike previous studies, we report a special phenomenon in the second to last layer, where most of all tasks have unexpected common neurons.

Finally, we confirm the observations of the attention mechanism and the activated neurons in the ablation study. We hope that this study provides a guide to the explainable soft prompts and the explainable PLM.

## 596 Limitations

597 Although we explore the multilingual space of  
598 PLM, our findings are limited to English and  
599 Korean. We choose Korean as a non-English  
600 language because Korean is understudied in  
601 Prompt-tuning and has different linguistic  
602 properties from English, such as typology. Also,  
603 even though our experiments have a fixed seed  
604 (42), other trials with other seed numbers are  
605 required, since different seed numbers can lead to  
606 fluctuating results. Additionally, we fail to analyze  
607 the observations relating to the performances. We  
608 encourage further studies to include more various  
609 languages and more systemic experiments.

## 610 References

611 Samuel R. Bowman, Gabor Angeli, Christopher Potts,  
612 and Christopher D. Manning. 2015. [A large  
613 annotated corpus for learning natural language  
614 inference](#). In *Proceedings of the 2015 Conference  
615 on Empirical Methods in Natural Language  
616 Processing*, pages 632–642, Lisbon, Portugal.  
617 Association for Computational Linguistics.

618 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom  
619 Kwiatkowski, Michael Collins, and Kristina  
620 Toutanova. 2019. [BoolQ: Exploring the Surprising  
621 Difficulty of Natural Yes/No Questions](#).  
622 In *Proceedings of the 2019 Conference of the North  
623 American Chapter of the Association for  
624 Computational Linguistics: Human Language  
625 Technologies, Volume 1 (Long and Short Papers)*,  
626 pages 2924–2936, Minneapolis, Minnesota.  
627 Association for Computational Linguistics.

628 Mor Geva, Roei Schuster, Jonathan Berant, and Omer  
629 Levy. 2021. [Transformer Feed-Forward Layers Are  
630 Key-Value Memories](#). In *Proceedings of the 2021  
631 Conference on Empirical Methods in Natural  
632 Language Processing*, pages 5484–5495, Online  
633 and Punta Cana, Dominican Republic. Association  
634 for Computational Linguistics.

635 Jiyeon Ham, Yo Joong Choe, Kyubyong Park, Ilji  
636 Choi, and Hyungjoon Soh. 2020. [KorNLI and  
637 KorSTS: New Benchmark Datasets for Korean  
638 Natural Language Understanding](#). In *Findings of the  
639 Association for Computational Linguistics: EMNLP  
640 2020*, pages 422–430, Online. Association for  
641 Computational Linguistics.

642 Daniel Khashabi, Xinxi Lyu, Sewon Min, Lianhui Qin,  
643 Kyle Richardson, Sean Welleck, Hannaneh  
644 Hajishirzi, Tushar Khot, Ashish Sabharwal, Sameer  
645 Singh, and Yejin Choi. 2022. [Prompt Waywardness:  
646 The Curious Case of Discretized Interpretation of  
647 Continuous Prompts](#). In *Proceedings of the 2022  
648 Conference of the North American Chapter of the*

649 *Association for Computational Linguistics: Human  
650 Language Technologies*, pages 3631–3643, Seattle,  
651 United States. Association for Computational  
652 Linguistics.

653 Brian Lester, Rami Al-Rfou, and Noah Constant.  
654 2021. [The Power of Scale for Parameter-Efficient  
655 Prompt-tuning](#). In *Proceedings of the 2021  
656 Conference on Empirical Methods in Natural  
657 Language Processing*, pages 3045–3059, Online  
658 and Punta Cana, Dominican Republic. Association  
659 for Computational Linguistics.

660 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan  
661 Ghazvininejad, Abdelrahman Mohamed, Omer  
662 Levy, Veselin Stoyanov, and Luke Zettlemoyer.  
663 2020. [BART: Denoising Sequence-to-Sequence  
664 Pre-training for Natural Language Generation,  
665 Translation, and Comprehension](#). In *Proceedings of  
666 the 58th Annual Meeting of the Association for  
667 Computational Linguistics*, pages 7871–7880,  
668 Online. Association for Computational Linguistics.

669 Xiang Lisa Li and Percy Liang. 2021. [Prefix-Tuning:  
670 Optimizing Continuous Prompts for Generation](#).  
671 In *Proceedings of the 59th Annual Meeting of the  
672 Association for Computational Linguistics and the  
673 11th International Joint Conference on Natural  
674 Language Processing (Volume 1: Long Papers)*,  
675 pages 4582–4597, Online. Association for  
676 Computational Linguistics.

677 Seungyoung Lim, Myungji Kim, and Jooyoul Lee.  
678 2019. [Korquad1.0: Korean qa dataset for machine  
679 reading comprehension](#). arXiv preprint  
680 arXiv:1909.07005.

681 Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei  
682 Zhou, Chandra Bhagavatula, Yejin Choi, and Xiang  
683 Ren. 2020. [CommonGen: A Constrained Text  
684 Generation Challenge for Generative  
685 Commonsense Reasoning](#). In *Findings of the  
686 Association for Computational Linguistics: EMNLP  
687 2020*, pages 1823–1840, Online. Association for  
688 Computational Linguistics.

689 Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding,  
690 Yujie Qian, Zhilin Yang, and Jie Tang. 2021. [Gpt  
691 understands, too](#). arXiv:2103.10385.

692 Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam,  
693 Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. [P-  
694 Tuning: Prompt-tuning Can Be Comparable to Fine-  
695 tuning Across Scales and Tasks](#). In *Proceedings of  
696 the 60th Annual Meeting of the Association for  
697 Computational Linguistics (Volume 2: Short  
698 Papers)*, pages 61–68, Dublin, Ireland. Association  
699 for Computational Linguistics.

700 Sungjoon Park, Jihyung Moon, Sungdong Kim, Won  
701 Ik Cho, Jiyeon Han, Jangwon Park, Chisung Song,  
702 Junseong Kim, Yongsook Song, Taehwan Oh, et al.

2021. [Klue: Korean language understanding evaluation](#). arXiv preprint arXiv:2105.09680.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#). OpenAI Blog, 1(8):9.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know What You Don’t Know: Unanswerable Questions for SQuAD](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics.
- Jaehyung Seo, Chanjun Park, Hyeonseok Moon, Sugyeong Eo, Myunghoon Kang, Seoungmoon Lee, and Heuseok Lim. 2021. [KommonGen: A Dataset for Korean Generative Commonsense Reasoning Evaluation](#). In *Proceedings of the 33th Annual Conference on Human & Cognitive Language Technology 2021*, October, pages 55-60.
- Oleh Shliachko, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and Tatiana Shavrina. 2022. [mgpt: Few-shot learners go multilingual](#). arXiv preprint arXiv:2204.07580.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Yusheng Su, Xiaozhi Wang, Yujia Qin, Chi-Min Chan, Yankai Lin, Huadong Wang, Kaiyue Wen, Zhiyuan Liu, Peng Li, Juanzi Li, Lei Hou, Maosong Sun, and Jie Zhou. 2022. [On Transferability of Prompt-tuning for Natural Language Processing](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3949–3969, Seattle, United States. Association for Computational Linguistics.
- Jesse Vig and Yonatan Belinkov. 2019. [Analyzing the Structure of Attention in a Transformer Language Model](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 63–76, Florence, Italy. Association for Computational Linguistics.
- Tu Vu, Aditya Barua, Brian Lester, Daniel Cer, Mohit Iyyer, and Noah Constant. 2022. [Overcoming catastrophic forgetting in zero-shot cross-lingual generation](#). arXiv preprint arXiv:2205.12647.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. [SuperGLUE: A stickier benchmark for general-purpose language understanding systems](#). In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022. [Finding Skill Neurons in Pre-trained Transformer-based Language Models](#). arXiv preprint at arXiv:2211.07349.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. [Character-level convolutional networks for text classification](#). In *Annual Conference on Neural Information Processing Systems 2015*, December 7–12, 2015, Montreal, Quebec, Canada, pages 649–657.
- Mengjie Zhao and Hinrich Schütze. 2021. [Discrete and Soft Prompting for Multilingual Models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8547–8555, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

803	<b>Appendices</b>	850
804	<b>A Dataset Details</b>	851
805	Table 3 shows the details of datasets including	852
806	sources. Also, Table 4 and Table 5 show the	853
807	verbalizers and the number of labels in each	854
808	dataset. We designed the verbalizers to be	855
809	semantically close to the actual label words.	856
810	<b>B Experimental Details</b>	857
811	We used two A100 GPUs with 80G memory. Also,	858
812	we used the seed number 42. The number of	859
813	trainable parameters for Prompt-tuning v1 is 40960	860
814	and the one for Prompt-tuning v2 is 983040. The	861
815	hyperparameters for each task is presented in Table	862
816	6. The performance results are presented in Table	863
817	7. For task STS, NLI, SA, and TC, we use sklearn	864
818	package ( <a href="https://scikit-learn.org/stable/modules/generated/sklearn.metrics.classification_report.html">https://scikit-</a>	865
819	<a href="https://scikit-learn.org/stable/modules/generated/sklearn.metrics.classification_report.html">learn.org/stable/modules/generated/s</a>	866
820	<a href="https://scikit-learn.org/stable/modules/generated/sklearn.metrics.classification_report.html">klearn.metrics.classification_report</a>	867
821	<a href="https://scikit-learn.org/stable/modules/generated/sklearn.metrics.classification_report.html">.html</a> ) for scoring. For QA, we use the codes	868
822	from the official website (SQuAD:	869
823	<a href="https://rajpurkar.github.io/SQuAD-explorer/">https://rajpurkar.github.io/SQuAD-</a>	870
824	<a href="https://rajpurkar.github.io/SQuAD-explorer/">explorer/</a> ; KorQuAD:	871
825	<a href="https://korquad.github.io/KorQuad%201.0/">https://korquad.github.io/KorQuad%20</a>	872
826	<a href="https://korquad.github.io/KorQuad%201.0/">1.0/</a> ). Also, For CG, we use the codes from the	873
827	official <a href="https://github.com/INK-USC/CommonGen">github</a> (CommonGen:	874
828	<a href="https://github.com/INK-USC/CommonGen">https://github.com/INK-USC/CommonGen;</a>	875
829	KommonGen: <a href="https://github.com/nlpai-lab/KommonGen">https://github.com/nlpai-</a>	876
830	<a href="https://github.com/nlpai-lab/KommonGen">lab/KommonGen</a> ).	877
831	<b>C Experimental Results</b>	878
832	Figure 9, Figure 10, and Figure 11 illustrate the	879
833	PCA results in Prompt-tuning v2 over all layers.	880
834	Also, Figure 12 shows the attention variability for	881
835	all tasks. Figure 13 and Figure 14 show the results	882
836	of KL divergence for all tasks. Lastly, Figure 15	883
837	presents the average KL divergence per input type.	884
838		885
839		886
840		887
841		888
842		889
843		890
844		891
845		892
846		893
847		894
848		
849		

Type	Task	English	Source	Korean	Source
Pair CLS	STS	GLUE-STS	Wang et al., 2019	KLUE-STS	Park et al., 2021
	NLI	SNLI	Ham et al., 2020	KorNLI	Ham et al., 2020
Single CLS	SA	SST2	Socher et al., 2013	NSMC	<a href="https://github.com/e9t/nsmc">https://github.com/e9t/nsmc</a>
	TC	AGNews	Zhang et al., 2015	KLUE-Ynat	Park et al., 2021
Extraction	QA	SQuAD 2.0	Rajpurkar et al., 2018	KorQuAD 1.0	Lim et al., 2019
Generation	CG	CommonGen	Lin et al., 2020	KommonGen	Seo et al., 2021

Table 3: The details about datasets used in this study.

Task	English Verbalizer (train/validation/test)
STS	similar (2994/629), dissimilar (2755/871)
NLI	entailment (183414/3329/3368), contradiction (183185/3278/3237), neutral (182762/3235/3219)
SA	positive (37569/444), negative (29780/428)
TC	business, scitech, sports, world (30000/1900)

Table 4: The verbalizers (label words) and the number of each example for English classification datasets. If test dataset is not provided, we use validation dataset instead.

Task	Korean Verbalizer (train/validation/test)
STS	유사(similar) (5602/220), 상이(dissimilar) (6066/299)
NLI	함의(entailment) (183382/523/1670), 모순(contradiction) (183382/524/1670), 중립(neutral) (183382/523/1669)
SA	긍정(positive) (74825/25171), 부정(negative) (75710/24826)
TC	정치(politics) (7379/722), 경제(economics) (6118/1348), 생활문화(livingculture) (5751/1369), 사회(society) (5133/3701), IT 과학(ITscience) (5235/554), 세계(world) (8320/835), 스포츠(sports) (7742/578)

Table 5: The verbalizers (label words) and the number of each example for Korean classification datasets. If test dataset is not provided, we use validation dataset instead.

917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931

Task	Batch size	Epochs	Max length
STS	16	10	180
NLI	32	16	150
SA	32	20	256
TC	64/32	10	80
QA	16	12	400
CG	16	20	60

Table 6: The hyperparameters of each task. For batch size in task TC, 64 is for English dataset and 32 is for Korean dataset.

Scoring Method	English	Score		Korean	Score	
		V1	V2		V1	V2
Accuracy	GLUE-STS	83.2	84.23	KLUE-STS	38.43	71.59
Macro-F1	SNLI	80.95	86.46	KorNLI	45.67	62.8
Accuracy	SST2	87.15	88.18	NSMC	84.94	87.18
Macro-F1	AGNews	85.8	87.27	KLUE-Ynat	81.18	84.27
F1/EM	SQuAD 2.0	61.98/	67.21/	KorQuAD 1.0	65.41/	72.32/
		47.25	51.91		59.62	66.50
Coverage	CommonGen	78.4	82.97	KommonGen	87.94	91.33

Table 7: The performances of all tasks in Prompt-tuning v1 and Prompt-tuning v2.

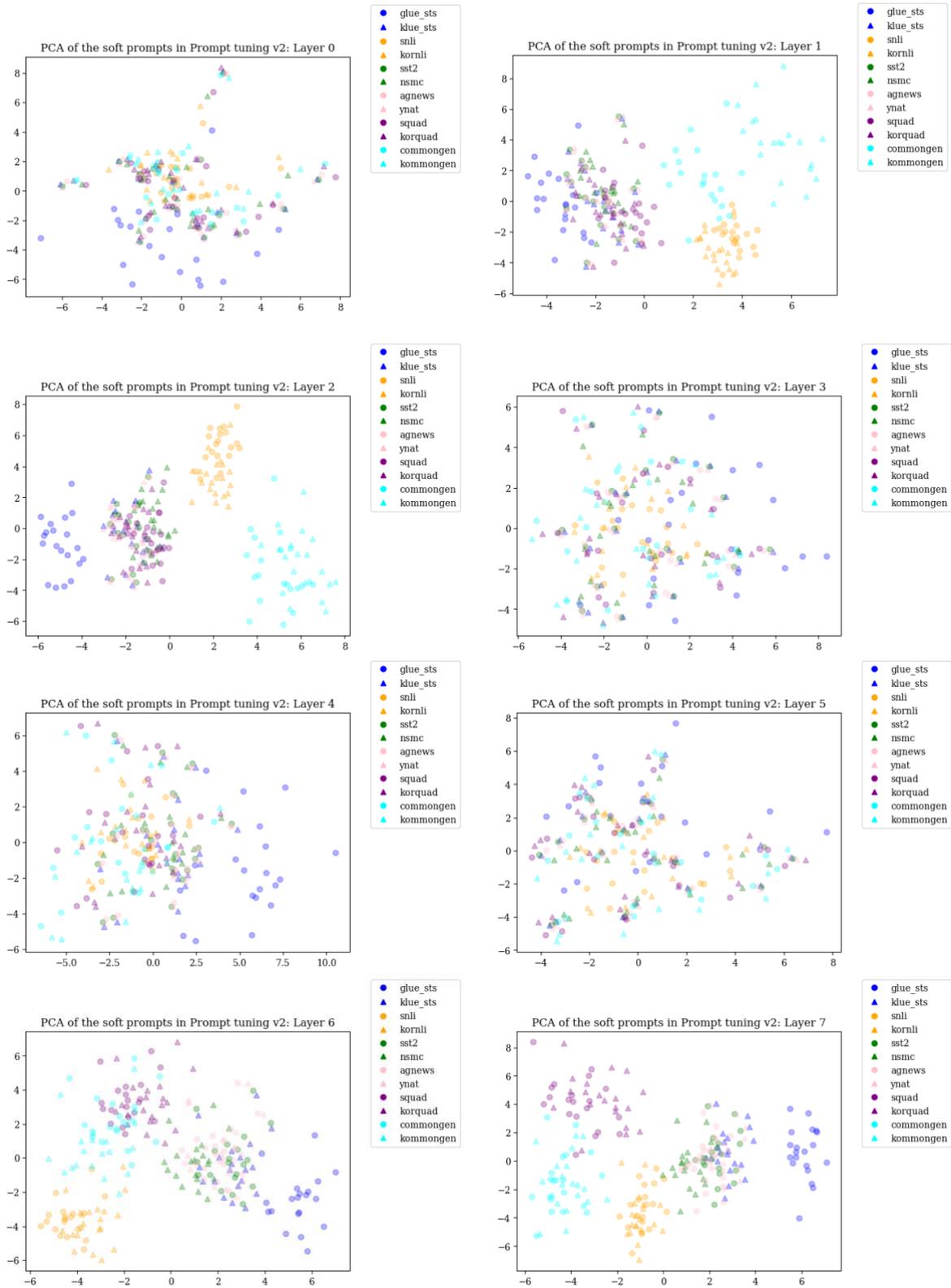


Figure 9: The PCA result of soft prompts learned for the target task in Prompt-tuning v2 in layers 0~7.

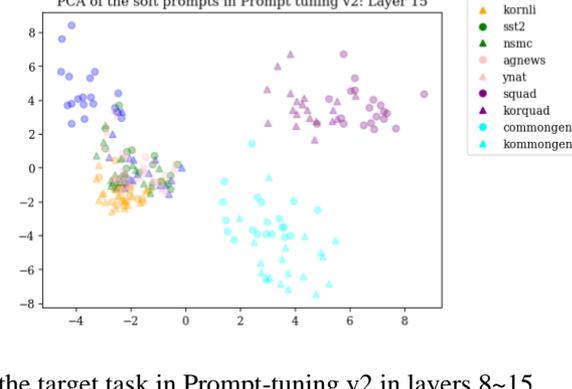
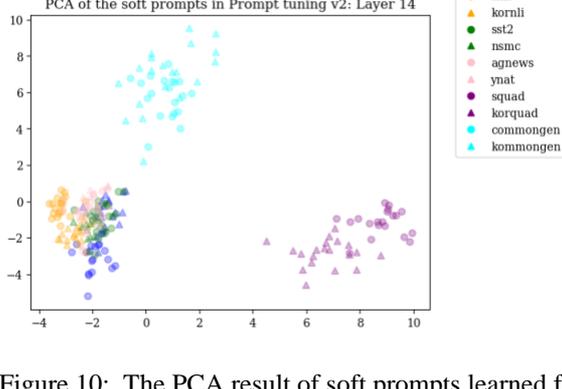
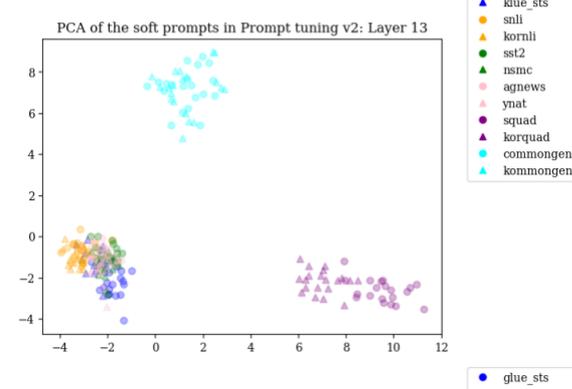
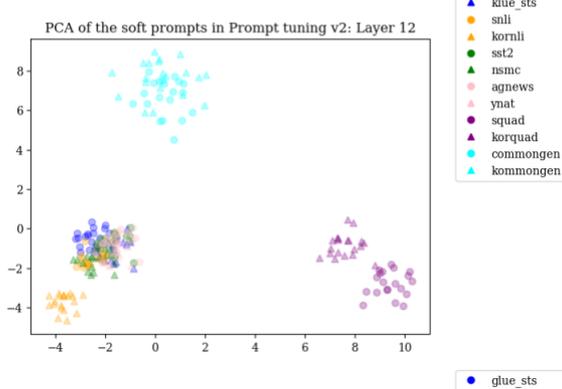
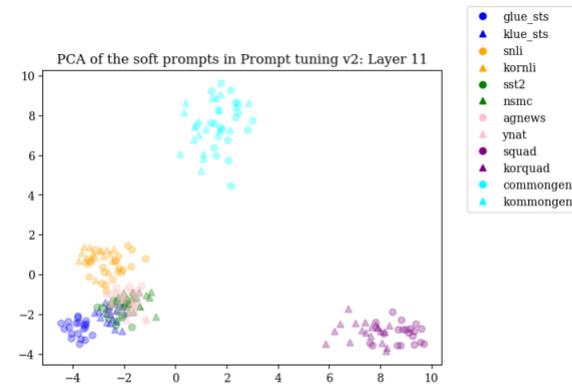
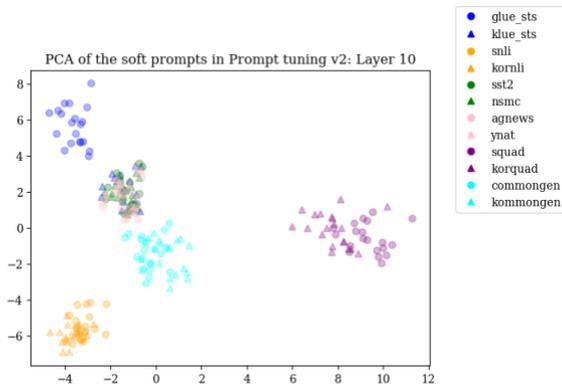
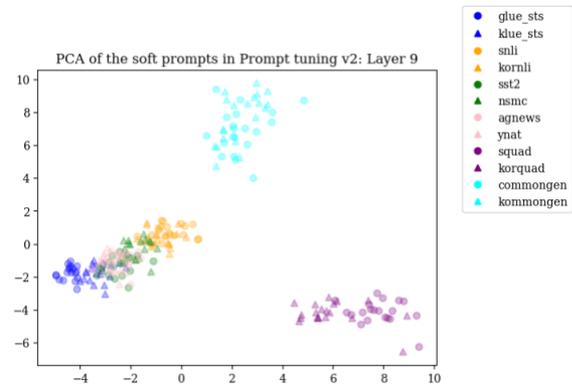
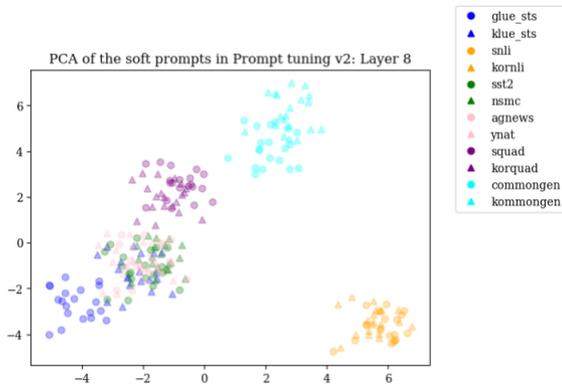


Figure 10: The PCA result of soft prompts learned for the target task in Prompt-tuning v2 in layers 8~15.

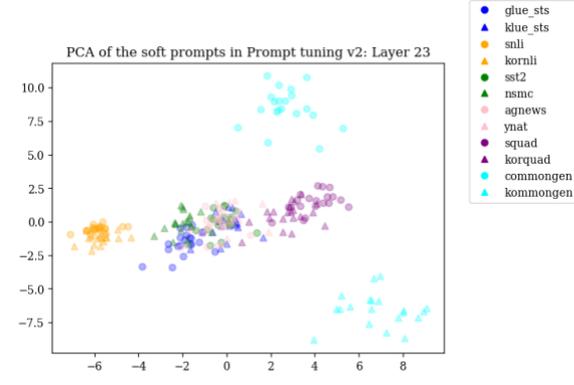
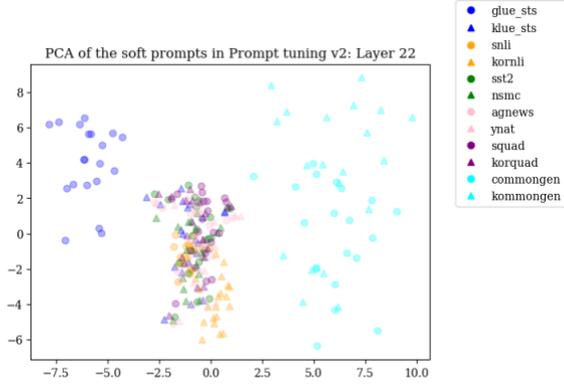
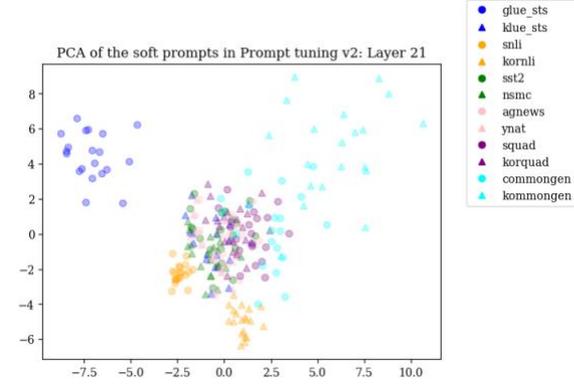
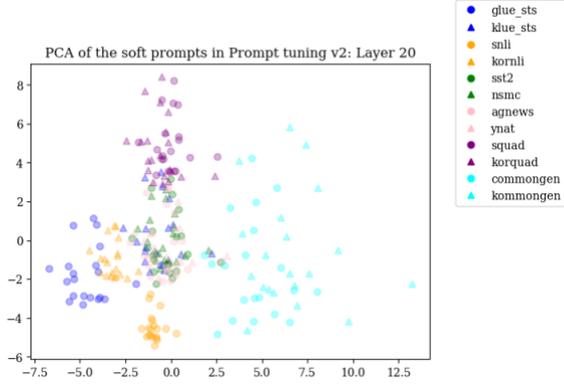
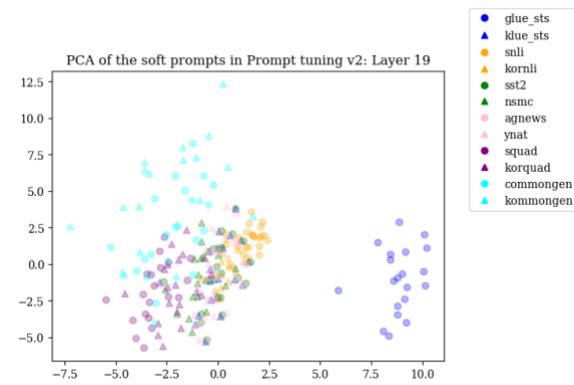
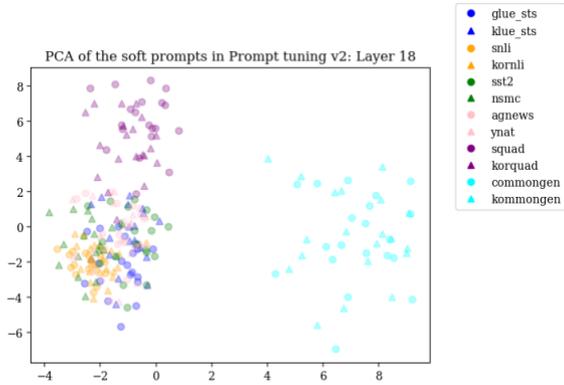
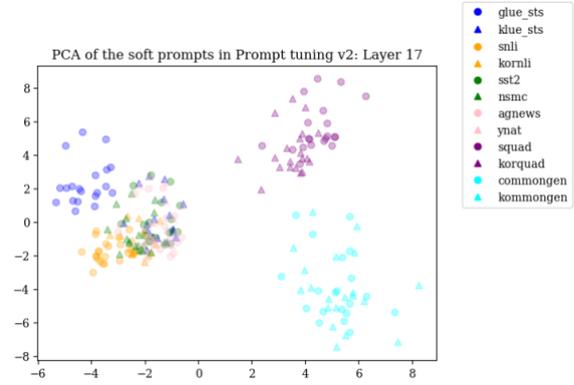
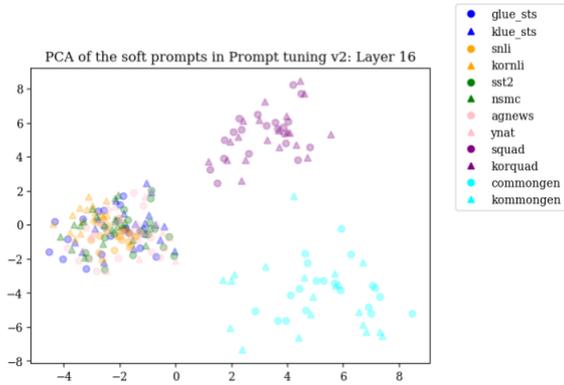


Figure 11: The PCA result of soft prompts learned for the target task in Prompt-tuning v2 in layers 16~23.

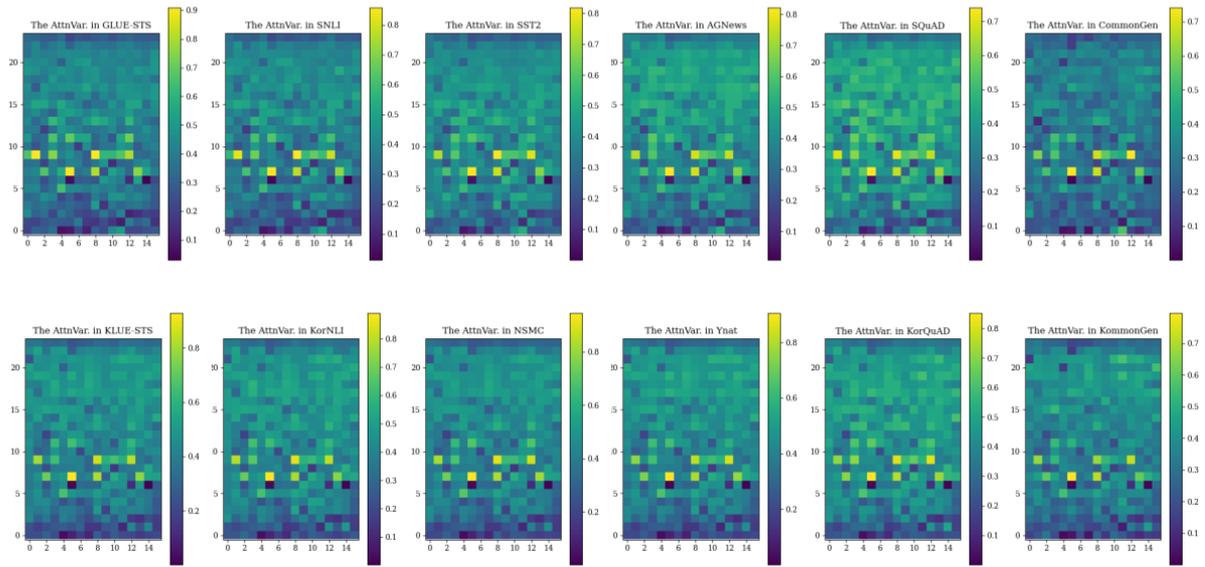


Figure 12: The attention variability of each task.

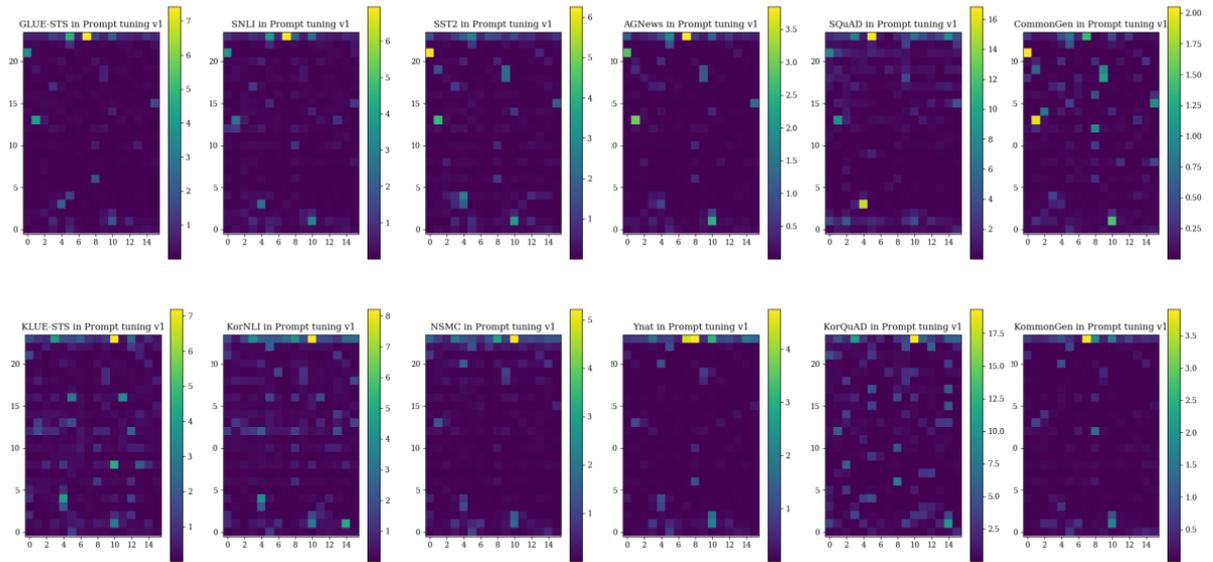


Figure 13: The KL divergence of each task in Prompt-tuning v1.

938  
939  
940  
941  
942  
943  
944

945  
946  
947  
948  
949  
950  
951

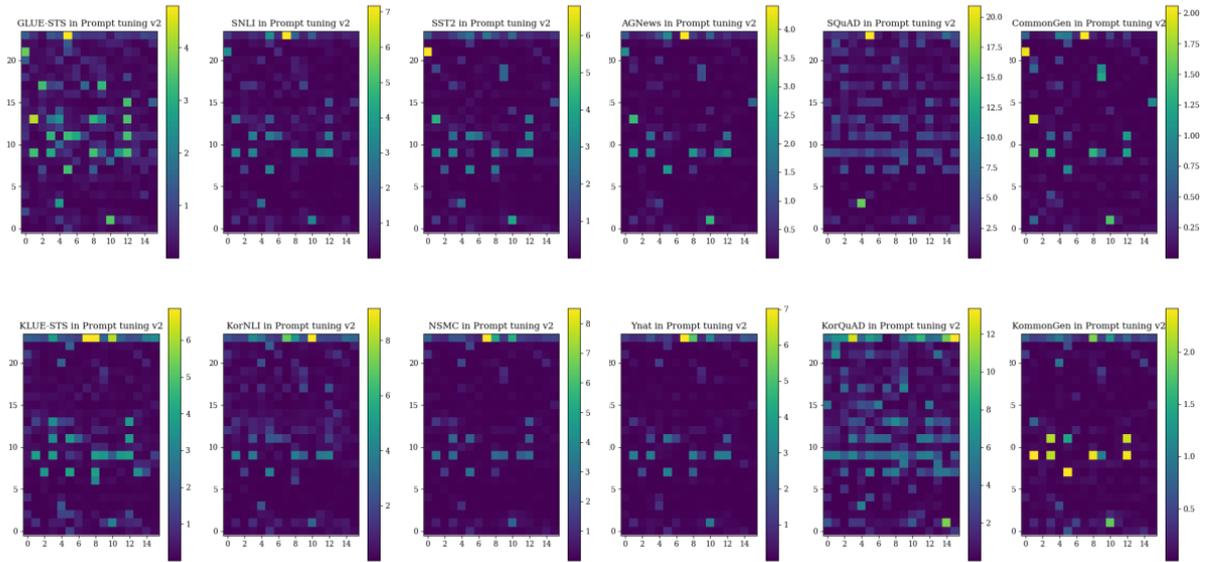


Figure 14: The KL divergence of each task in Prompt-tuning v2.

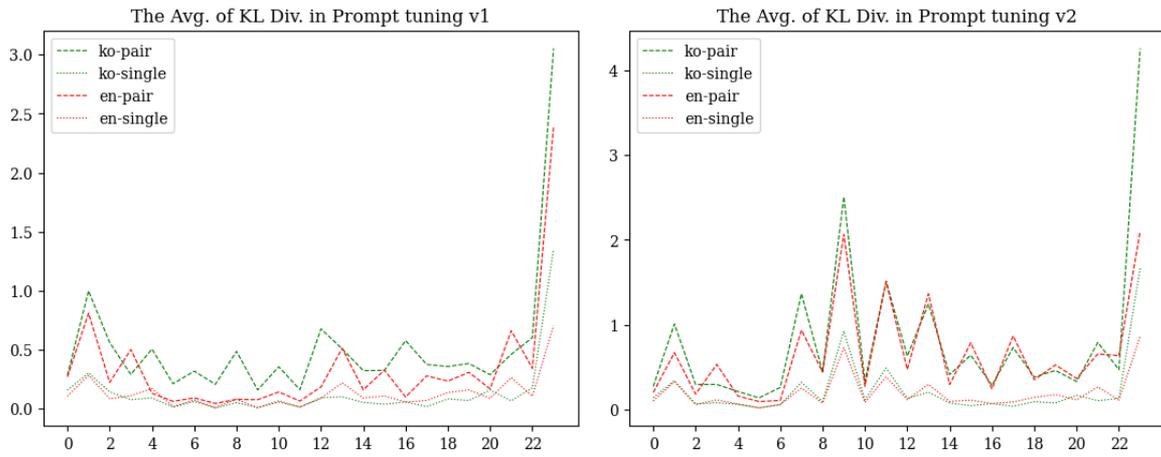


Figure 15: The average KL divergence results grouped by input type.

952  
953  
954  
955  
956  
957  
958  
959  
960  
961

962  
963  
964  
965  
966  
967  
968  
969  
970  
971