

MVP

MULTI-TASK SUPERVISED PRE- TRAINING FOR NATURAL LANGUAGE GENERATION

NGUYỄN TRUNG NGUYÊN-520V0015
HỒ TRỌNG NGHĨA-520K0163

CONTENTS

➤ Introduction

➤ Related Work

➤ The MVP Model

➤ Experiment Results

INTRODUCTION

The MVP model is a pre-trained natural language generation model that can handle various tasks using task-specific soft prompts. It is based on the Transformer encoder-decoder architecture and uses a large-scale corpus of labeled data from diverse NLG tasks for pre-training.



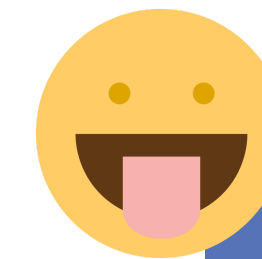
RELATED WORK

Canva



Pre-trained Language Models

Pre-trained language models have achieved exceptional success in a wide range of tasks, and the majority of them are pre-trained in an unsupervised manner



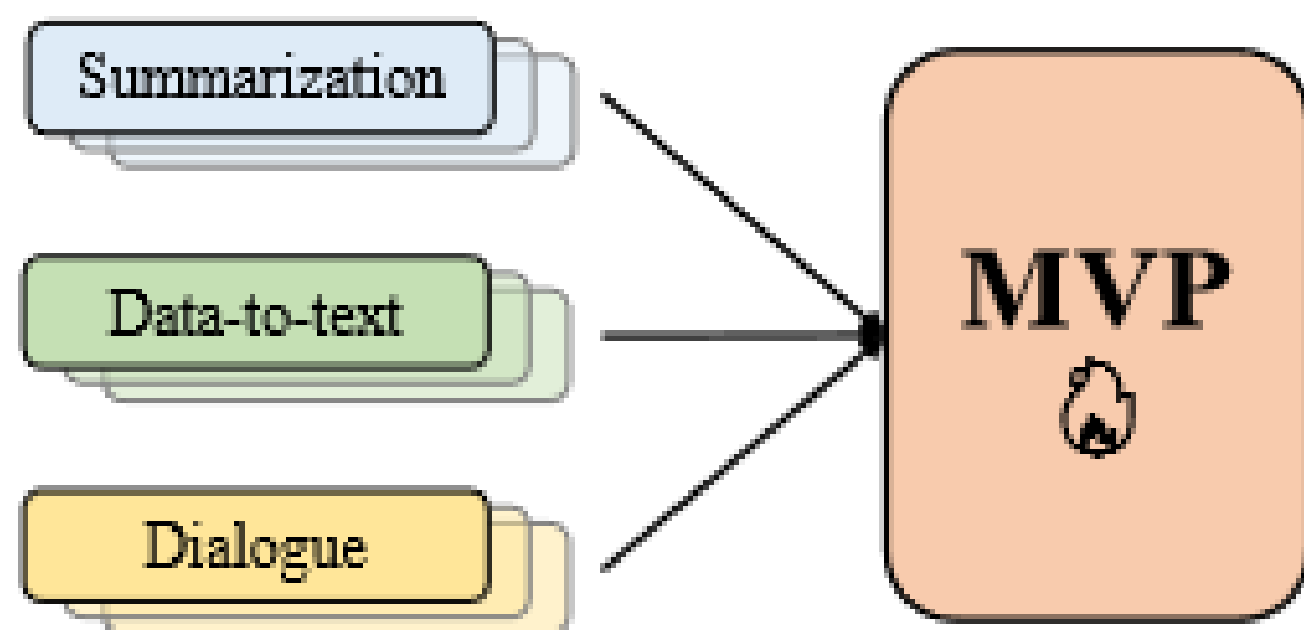
MULTI-TASK LEARNING

Competitive analysis allows us to understand where we are as a brand and how our competitors work. We will start by identifying who we are and who our competitors are. Next, we will identify attributes they are doing right and create a perceptual map. In the perceptual map, we will identify a criteria and rank these attributes as high or low.

PROMPT LEARNING

Prompt learning is a thriving method in the field of NLP. Prompt learning converts fine-tuning text into a format similar to pretraining to leverage implicit pre-training knowledge and alleviate the discrepancy between pre-training and fine-tuning.

Stage 1: Multi-task Supervised Pre-training



Stage 2: Task-specific Prompt Pre-training

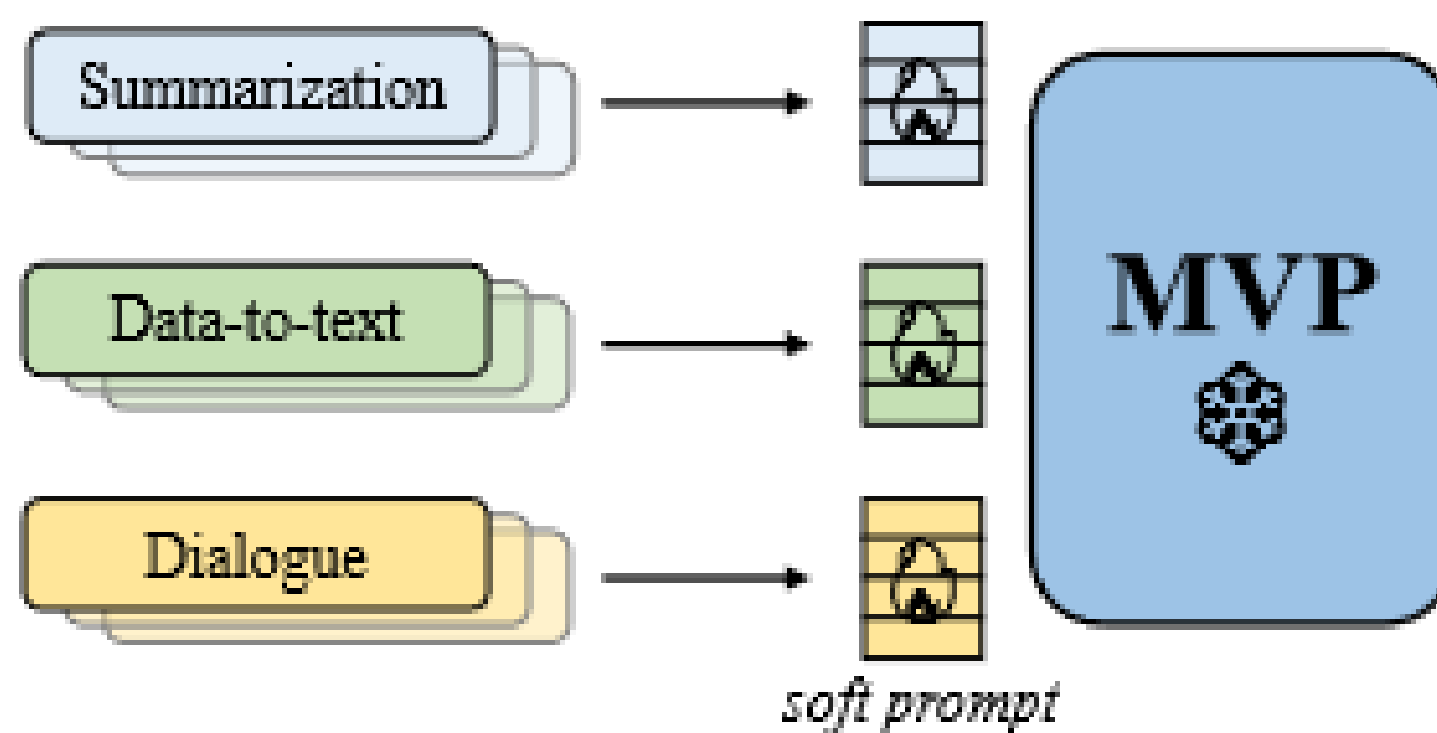


Figure 1: The overview of the pre-training process of our MVP model and task-specific prompts.

THE MVP MODEL

DATA COLLECTION

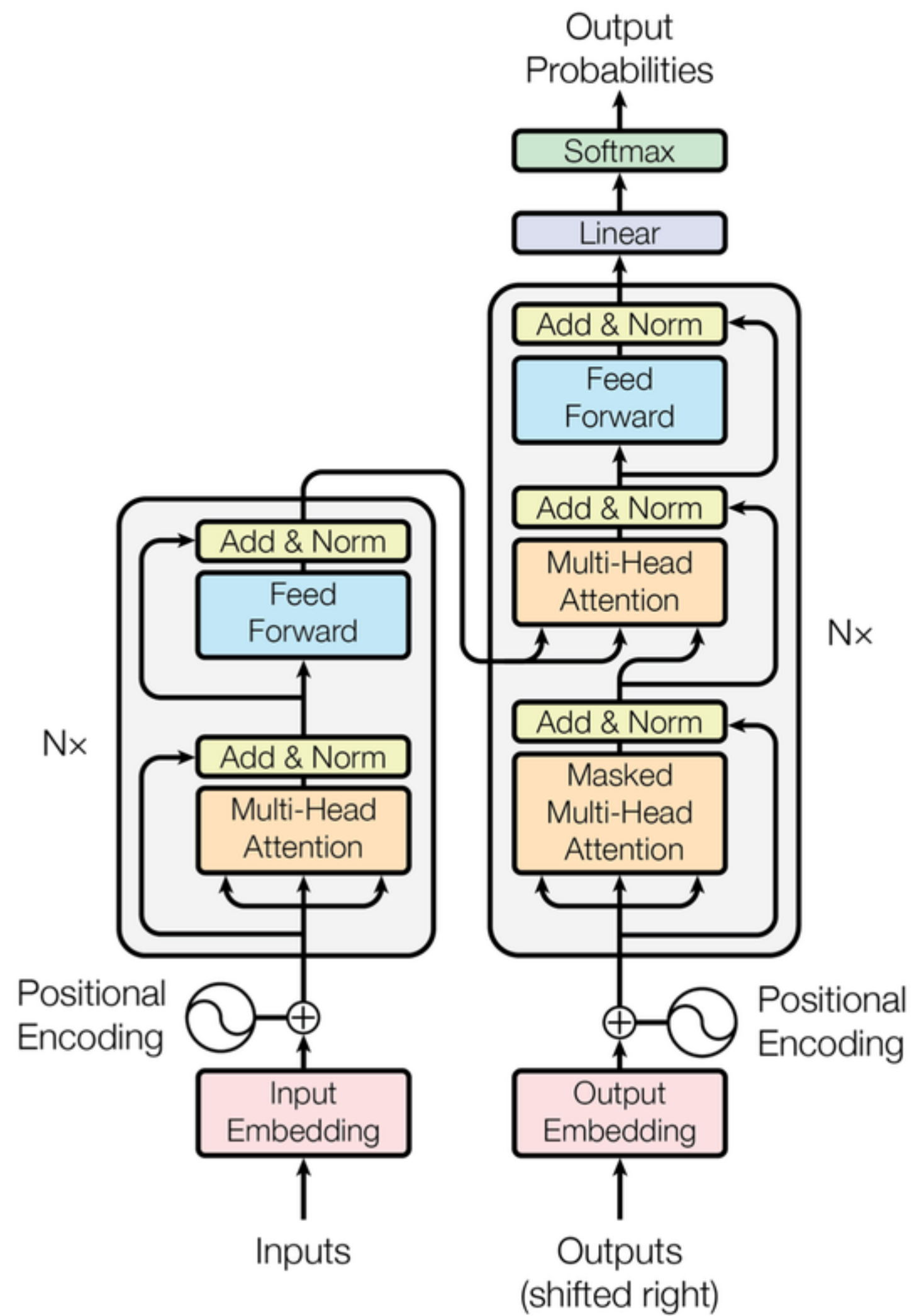
Collect a large-scale labeled MVPCorpus consisting of 77 labeled datasets from 11 representative NLG tasks¹, including common sense generation, data-to-text generation, open ended dialogue system, paraphrase generation, question answering, question generation, story generation, task oriented dialogue system, text simplification, text style transfer, and text summarization.

MODEL ARCHITECTURE

Our MVP model is built on the standard Trans-former encoder-decoder architecture

TRAINING DETAILS

MVP model adopts a Transformer with 12 layers in both encoder and decoder (406M parameters), the same as the model size of BART.



EXPERIMENT RESULTS

| Methods | CNN/DailyMail | | | WebNLG | | | SQuAD (QG) | | | CoQA | |
|--------------|--------------------------|--------------|--------------|--------------------|--------------------|--------------|--------------------|--------------|--------------------------|--------------------|--------------|
| | R-1 | R-2 | R-L | B-4 | ME | R-L | B-4 | ME | R-L | F1 | EM |
| MVP | 44.52 | 21.62 | 41.10 | <u>67.82</u> | 47.47 | <u>76.88</u> | 26.26 | 27.35 | 53.49 | <u>86.43</u> | <u>77.78</u> |
| BART | 44.16 ^e | 21.28 | 40.90 | 64.55 ^b | 46.51 | 75.13 | 22.00 ^f | 26.40 | 52.55 | 68.60 ^f | – |
| Single | 44.36 | 21.54 | 40.88 | 67.74 | 46.89 | 76.94 | <u>26.09</u> | 27.15 | 53.29 | 86.20 | 77.26 |
| MVP+S | <u>44.63</u> | <u>21.72</u> | <u>41.21</u> | 68.19 | 47.75 | 76.81 | 25.69 | 27.04 | 53.20 | 86.65 | 77.93 |
| MVP+R | 44.14 | 21.45 | 40.72 | 67.61 | <u>47.65</u> | 76.70 | 25.71 | 27.03 | 53.09 | 85.95 | 77.22 |
| MVP+M | 43.97 | 21.16 | 40.46 | 67.45 | <u>47.57</u> | 76.81 | 25.46 | 26.79 | 52.95 | 86.28 | 77.26 |
| SOTA | 47.16^a | 22.55 | 43.87 | 66.14 ^b | 47.25 | 76.10 | 25.97 ^c | <u>27.33</u> | <u>53.43</u> | 84.50 ^d | – |
| Methods | ROCStories | | | | PersonaChat | | | | MultiWOZ | | |
| | B-1 | B-2 | D-1 | D-4 | B-1 | B-2 | D-1 | D-2 | B-4 | Success | Inform |
| MVP | <u>33.79</u> | 15.76 | <u>3.02</u> | <u>75.65</u> | 50.73 | 40.69 | 1.65 | 11.23 | 20.26 | 76.40 | 85.00 |
| BART | 30.70 ^g | 13.30 | – | 69.90 | 49.90 ^f | 40.00 | 1.30 | 8.00 | 17.89 ^j | 74.91 | 84.88 |
| Single | 32.67 | 15.29 | 2.72 | 72.97 | <u>49.96</u> | <u>40.53</u> | 1.27 | 7.63 | 19.73 | 75.60 | 83.70 |
| MVP+S | 33.92 | <u>15.60</u> | 3.44 | 80.58 | 47.91 | 39.97 | <u>1.52</u> | <u>9.54</u> | <u>20.32</u> | <u>79.90</u> | <u>86.80</u> |
| MVP+R | 32.93 | 15.32 | 2.88 | 73.83 | 48.45 | 40.09 | 1.30 | 7.95 | 19.02 | 73.30 | 81.80 |
| MVP+M | 33.30 | 15.51 | 2.71 | 74.24 | 46.26 | 39.30 | 1.36 | 8.07 | 19.93 | 72.70 | 79.70 |
| SOTA | 33.40 ^g | 15.40 | – | 69.30 | 49.90 ^f | 40.00 | 1.50 ^h | 9.40 | 20.50ⁱ | 85.30 | 94.40 |

Table 2: The main results on seven seen tasks under full tuning settings. The best and second-best results among all the methods are marked in **bold** and underlined, respectively. The SQuAD dataset here is used for the question generation task. The letters B, R, D, and ME denote BLEU, ROUGE, Distinct, and METEOR, respectively. “–” means the work does not compute the corresponding result. These setups and abbreviations are the same below. ^a (Ravaut et al., 2022) ^b (Ke et al., 2021) ^c (Bao et al., 2021) ^d (Xiao et al., 2020) ^e (Lewis et al., 2020) ^f (Liu et al., 2021a) ^g (Guan et al., 2021) ^h (Chen et al., 2022) ⁱ (He et al., 2022) ^j (Lin et al., 2020c)

| AESOP | Quora | | | | | SC & BLEU | GYAFC E&M | | | GYAFC F&R | | |
|-------|--------------------|--------------|--------------|--------------|--------------|-----------|--------------------|--------------|--------------|--------------|--------------|--------------|
| | B-4 | R-1 | R-2 | R-L | ME | | B-4 | Accuracy | HM | B-4 | Accuracy | HM |
| +BART | 47.30 ^a | 73.30 | 54.10 | 75.10 | 49.70 | +BART | 76.50 ^b | 93.70 | 83.90 | 79.30 | 92.00 | 85.20 |
| +MVP | 49.81 | 74.78 | 56.84 | 76.34 | 53.40 | +MVP | 77.18 | 94.49 | 84.96 | 79.43 | 92.12 | 85.31 |

Table 3: The results of unseen NLG tasks. We use AESOP and SC & BLEU to denote the methods proposed by Sun et al. (2021) and Lai et al. (2021), respectively. ^a (Sun et al., 2021) ^b (Lai et al., 2021)

| Methods | CoLA Matt. | SST-2 Acc. | MRPC F1/Acc. | STS-B P/S Corr. | QQP F1/Acc. | MNLI m./mm. | QNLI Acc. | RTE Acc. | Average |
|---------|---------------|---------------|----------------------|----------------------|----------------------|----------------------|--------------|--------------|--------------|
| BART | 60.30 | 96.30 | 90.47 / 86.70 | 90.97 / 90.30 | 73.03 / 89.87 | 90.03 / 89.27 | 94.60 | 79.83 | 85.17 |
| MVP | 59.87 | 96.43 | 92.07 / 89.43 | 91.37 / 90.90 | 73.20 / 90.13 | 89.70 / 88.73 | 95.10 | 82.87 | 85.88 |

Table 4: The results of NLU tasks on the GLUE benchmark.

| Methods | CNN/DailyMail | | | WebNLG | | | SQuAD (QG) | | | CoQA | |
|--------------|---------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | R-1 | R-2 | R-L | B-4 | ME | R-L | B-4 | ME | R-L | F1 | EM |
| MVP+S | 43.03 | <u>20.27</u> | 39.72 | 66.73 | 47.42 | 76.36 | 25.28 | 26.66 | <u>52.69</u> | 86.44 | 76.84 |
| BART+R | 42.47 | 19.82 | 39.15 | 65.54 | 46.86 | 75.24 | 24.27 | 26.07 | 52.03 | 82.22 | 71.92 |
| MVP+R | 42.84 | 20.21 | 39.61 | 66.12 | 47.12 | 75.83 | 25.05 | 26.34 | 52.57 | 85.51 | 75.56 |
| MVP+M | <u>42.99</u> | 20.36 | <u>39.70</u> | <u>66.40</u> | <u>47.16</u> | <u>75.89</u> | <u>25.24</u> | <u>26.49</u> | 52.88 | <u>85.90</u> | <u>76.34</u> |
| FT BART | 44.16 | 21.28 | 40.90 | 64.55 | 46.51 | 75.13 | 22.00 | 26.40 | 52.55 | 68.60 | – |
| FT MVP | 44.52 | 21.62 | 41.10 | 67.82 | 47.47 | 76.88 | 26.26 | 27.35 | 53.49 | 86.43 | 77.78 |
| Methods | ROCStories | | | | PersonaChat | | | | MultiWOZ | | |
| | B-1 | B-2 | D-1 | D-4 | B-1 | B-2 | D-1 | D-2 | B-4 | Success | Inform |
| MVP+S | 32.94 | <u>15.12</u> | 2.98 | 71.09 | 47.11 | 39.51 | 1.39 | 7.28 | 19.24 | 71.40 | 77.80 |
| BART+R | 32.14 | 14.71 | 2.85 | 68.94 | 46.23 | 38.98 | 1.30 | 6.82 | 17.94 | 62.20 | 69.20 |
| MVP+R | 32.28 | 14.85 | <u>2.97</u> | <u>70.29</u> | 46.70 | 39.23 | 1.31 | 6.98 | 18.86 | 64.40 | 71.40 |
| MVP+M | <u>32.62</u> | 15.28 | 2.95 | 69.58 | <u>46.78</u> | <u>39.40</u> | <u>1.33</u> | <u>7.13</u> | <u>19.13</u> | <u>67.20</u> | <u>72.90</u> |
| FT BART | 30.70 | 13.30 | – | 69.90 | 49.90 | 40.00 | 1.30 | 8.00 | 17.89 | 74.91 | 84.88 |
| FT MVP | 33.79 | 15.76 | 3.02 | 75.65 | 50.73 | 40.69 | 1.65 | 11.23 | 20.26 | 76.40 | 85.00 |

Table 5: The results on seven seen tasks under parameter-efficient settings. We also include the results of BART and MVP under the full tuning setting (denoted as FT) for comparison.

| Methods | #NLG (PT) | #NLU (PT) | #NLG (FT) | #NLU (FT) | SP model | SP prompts | Open source |
|------------|-----------|-----------|-----------|-----------|----------|------------|-------------|
| FLAN | 3 | 9 | 2 | 9 | ✓ | ✗ | ✗ |
| T0 | 2 | 6 | 0 | 4 | ✓ | ✗ | ✓ |
| Muppet | 1 | 3 | 1 | 3 | ✓ | ✗ | ✓ |
| ExT5 | 3 | 8 | 6 | 8 | ✓ | ✗ | ✗ |
| SPoT | 1 | 4 | 0 | 6 | ✗ | ✓ | ✗ |
| MVP (ours) | 7 | 0 | 11 | 3 | ✓ | ✓ | ✓ |

Table 7: Comparison of MVP with existing supervised pre-training works. #NLG/#NLU are the number of NLG and NLU tasks, respectively. PT, FT, and SP denote pre-training, fine-tuning, and supervised pre-training, respectively.