

Multitask Vision-Language Prompt Tuning

Sheng Shen* Shijia Yang* Tianjun Zhang* Bohan Zhai
Joseph E. Gonzalez Kurt Keutzer Trevor Darrell
University of California, Berkeley

{sheng.s, shijiayang, tianjunz, zhaibohan, jegonzal, keutzer, trevordarrell}@berkeley.edu

1. Appendix

1.1. Additional Results

Ablation on UPT As mentioned in the main text, due to the recency, [2] does not release their model details or code. We therefore implement our own variant that simply concatenates the CoOp prompt vectors U_T and VPT-deep prompt vector U_V together as U , we set the context length of U_T and U_V the same as 4 unless specify. We use a one-layer one-head Transformer block θ whose hidden dimension is cut to be 128. Before and after feeding U to θ , a linear layer is employed to match the dimensionality. We ablate the design choice on number of heads, number of layer, and dimensionality, respectively in Figure 1. The size of the each point stands for the relative additional parameter size included in this setting.

1.2. Task Group Information

We provide detailed task group information here. We follow Table 1 for multitask adaptation where group of 1 task means using "Target task" column only, group of 2 tasks means target task with task 1, and group of 3 tasks means target tasks, task 1, with task 2.

1.3. CIFAR-10 for Cross-task Generation

As stated in main text Table 1, we include CIFAR-10 result for cross-task generation in Table 2.

1.4. Standard Deviation

Since all experiments are in the few-shot setting, we provide standard deviation for few-shot ELEVATER experiments in Table 3.

1.5. Theoretical Explanation

We here theoretically justify our proposed task grouping in the context of prompt multitask adaptation. In our method, we group task $i \in T$ with task $j \in T$ if P_i , the task i specific prompt, performs better on task j than other tasks in T . For

simplicity, we do not consider over-fitting and other edge cases, then task i has the largest positive effect of the gradient update on the given task j . To measure the effect of gradient update, we calculate the ratio of the loss after and before we plug in P_i . Formally, we define an asymmetric measure for calculating the affinity of task i at a given time-step t on task j as:

$$Z_{i \rightarrow j}^t = \frac{L_j(\chi^t, P_i^t)}{L_j(\chi^t, P_j^t)}$$

where χ^t is an inference batch. P_i^t represent task i prompt on time-step t . L_j is the loss of task j . Finally, we average across all steps:

$$\hat{Z}_{i \rightarrow j} = \frac{1}{T} \sum_{t=1}^T Z_{i \rightarrow j}^t$$

The proposed metric is similar to inter-task affinity in [1]. In Section 4.3, [1] theoretically prove that if task b induces higher affinity than task c on task a , training $\{a, b\}$ together result in a lower loss on task a than training $\{a, c\}$. By plugging $Z_{b \rightarrow a}$ and $Z_{c \rightarrow a}$ in the proof, Lemma 1 generalizes to our case and thus lead to the same conclusion.

References

- [1] Chris Fifty, Ehsan Amid, Zhe Zhao, Tianhe Yu, Rohan Anil, and Chelsea Finn. Efficiently identifying task groupings for multi-task learning. *Advances in Neural Information Processing Systems*, 34:27503–27516, 2021. 1
- [2] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Unified vision and language prompt learning. *arXiv preprint arXiv:2210.07225*, 2022. 1

*Equal contribution

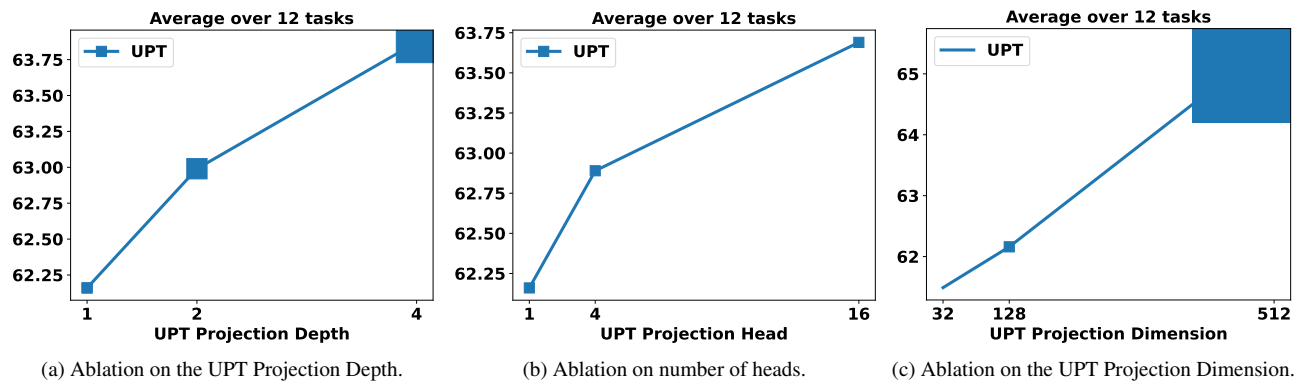


Figure 1. Ablation on the hyper-parameter of UPT. Note that the size of each point stands for the relative additional parameter size.

Table 1. Task group for CoOp, VPT, and UPT.

| Model | Target task | Task 1 | Task 2 |
|-------|----------------|---------------|----------------|
| CoOp | Caltech101 | DTD | CIFAR10 |
| | CIFAR10 | VOC 2007 | Resisc-45 |
| | CIFAR100 | Caltech101 | CIFAR10 |
| | Country-211 | Caltech101 | Resisc-45 |
| | DTD | Caltech101 | MNIST |
| | EuroSat | Resisc-45 | CIFAR100 |
| | FER 2013 | CIFAR100 | MNIST |
| | FGVCAircraft | Caltech101 | DTD |
| | Flowers102 | CIFAR10 | Caltech101 |
| | Food101 | Caltech101 | DTD |
| | GTSRB | MNIST | CIFAR100 |
| | Hateful Memes | VOC 2007 | Caltech101 |
| | KITTI Distance | StanfordCars | OxfordPets |
| | MNIST | DTD | Resisc-45 |
| | OxfordPets | Caltech101 | CIFAR10 |
| | Patch-Camelyon | CIFAR100 | Caltech101 |
| | Rendered-SST2 | FGVCAircraft | Hateful Memes |
| | Resisc-45 | Caltech101 | CIFAR10 |
| | StanfordCars | Caltech101 | MNIST |
| | VOC 2007 | CIFAR100 | Caltech101 |
| VPT | Caltech101 | CIFAR100 | CIFAR10 |
| | CIFAR10 | CIFAR100 | Caltech101 |
| | CIFAR100 | CIFAR10 | Caltech101 |
| | Country-211 | EuroSat | Food101 |
| | DTD | CIFAR10 | Rendered-SST2 |
| | EuroSat | Resisc-45 | FER 2013 |
| | FER 2013 | OxfordPets | MNIST |
| | FGVCAircraft | EuroSat | CIFAR10 |
| | Flowers102 | CIFAR100 | EuroSat |
| | Food101 | CIFAR10 | EuroSat |
| | GTSRB | MNIST | CIFAR100 |
| | Hateful Memes | FER 2013 | OxfordPets |
| | KITTI Distance | VOC 2007 | Flowers102 |
| | MNIST | Resisc-45 | GTSRB |
| | OxfordPets | CIFAR10 | Rendered-SST2 |
| | Patch-Camelyon | CIFAR10 | Food101 |
| | Rendered-SST2 | Resisc-45 | Patch-Camelyon |
| | Resisc-45 | EuroSat | CIFAR10 |
| | StanfordCars | CIFAR10 | EuroSat |
| | VOC 2007 | CIFAR100 | CIFAR10 |
| UPT | Caltech101 | CIFAR10 | CIFAR100 |
| | CIFAR10 | CIFAR100 | Caltech101 |
| | CIFAR100 | Caltech101 | EuroSat |
| | Country-211 | Caltech101 | CIFAR100 |
| | DTD | CIFAR10 | Caltech101 |
| | EuroSat | Resisc-45 | CIFAR100 |
| | FER 2013 | MNIST | DTD |
| | FGVCAircraft | CIFAR100 | CIFAR10 |
| | Flowers102 | Caltech101 | CIFAR100 |
| | Food101 | Caltech101 | CIFAR10 |
| | GTSRB | MNIST | CIFAR100 |
| | Hateful Memes | Caltech101 | CIFAR100 |
| | KITTI Distance | Food101 | Flowers102 |
| | MNIST | CIFAR100 | GTSRB |
| | OxfordPets | Caltech101 | CIFAR100 |
| | Patch-Camelyon | Food101 | StanfordCars |
| | Rendered-SST2 | Hateful Memes | GTSRB |
| | Resisc-45 | CIFAR10 | CIFAR100 |
| | StanfordCars | Caltech101 | CIFAR100 |
| | VOC 2007 | CIFAR100 | Caltech101 |

Table 2. Cross-task generation experiment additional result.

| (a) CIFAR-10 | | | |
|--------------|------------------------|------------------------|------------------------|
| # shots | 1 | 5 | 20 |
| CoOp | 89.45 \pm 1.6 | 83.63 \pm 2.0 | 91.38 \pm 0.6 |
| CoCoOp | 92.61 \pm 1.5 | 84.91 \pm 1.6 | 91.85 \pm 0.6 |
| VPT | 86.81 \pm 1.4 | 86.77 \pm 1.0 | 90.83 \pm 0.7 |
| UPT | 88.44 \pm 0.6 | 89.47 \pm 0.9 | 91.33 \pm 0.6 |
| MCoOp | 90.48 \pm 1.2 | 90.86 \pm 1.6 | 92.92 \pm 0.5 |
| MCoCoOp | 93.16 \pm 0.6 | 96.26 \pm 0.6 | 98.10 \pm 0.4 |
| MVPT | 88.97 \pm 0.5 | 89.89 \pm 0.5 | 92.13 \pm 0.3 |
| MUPT | 92.61 \pm 0.2 | 91.52 \pm 0.6 | 93.72 \pm 0.4 |

Table 3. Few-shot ELEVATER experiment with standard deviation.

| | | Target | | | | | | | | | |
|-------------------|------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| Source | Adaptation | Caltech101 | CIFAR10 | CIFAR100 | Country-211 | DTD | EuroSat | FER 2013 | FGVCAircraft | Flowers102 | Food101 |
| CLIP [†] | - - | 88.9 | 90.8 | 68.2 | 22.8 | 44.8 | 54.7 | 48.5 | 24.3 | 88.7 | 43.5 |
| CoOp | - S | 91.44 \pm 0.4 | 91.30 \pm 1.0 | 73.01 \pm 0.2 | 22.83 \pm 0.7 | 69.82 \pm 0.9 | 80.19 \pm 3.0 | 54.46 \pm 3.4 | 42.01 \pm 0.6 | 93.31 \pm 0.3 | 89.47 \pm 0.1 |
| VPT | - S | 92.84 \pm 0.4 | 91.39 \pm 0.7 | 75.98 \pm 0.9 | 21.11 \pm 0.4 | 68.56 \pm 1.1 | 87.37 \pm 3.7 | 56.77 \pm 0.6 | 42.12 \pm 1.2 | 89.22 \pm 1.2 | 89.04 \pm 0.2 |
| UPT | - S | 92.58 \pm 0.4 | 92.05 \pm 1.1 | 76.61 \pm 0.2 | 23.37 \pm 0.5 | 67.68 \pm 0.6 | 88.98 \pm 2.4 | 56.87 \pm 1.9 | 42.46 \pm 1.0 | 89.59 \pm 0.5 | 89.64 \pm 0.4 |
| MCoOp | - M | 91.53 \pm 0.2 | 91.67 \pm 0.5 | 73.01 \pm 0.2 | 23.12 \pm 0.3 | 69.82 \pm 0.9 | 81.69 \pm 5.4 | 54.46 \pm 3.4 | 42.01 \pm 0.6 | 93.44 \pm 0.3 | 89.47 \pm 0.1 |
| MVPT | - M | 92.84 \pm 0.4 | 93.54 \pm 0.4 | 76.39 \pm 0.3 | 21.42 \pm 0.1 | 68.56 \pm 1.1 | 89.15 \pm 1.1 | 56.77 \pm 0.6 | 42.12 \pm 1.2 | 89.22 \pm 1.2 | 89.04 \pm 0.2 |
| MUPT | - M | 92.58 \pm 0.4 | 93.38 \pm 0.6 | 76.61 \pm 0.2 | 23.37 \pm 0.6 | 67.68 \pm 0.6 | 88.98 \pm 2.4 | 56.94 \pm 1.5 | 42.46 \pm 1.0 | 89.59 \pm 0.5 | 89.64 \pm 0.4 |
| MCoOp | M M | 92.09 \pm 0.2 | 91.59 \pm 0.9 | 72.63 \pm 0.1 | 23.52 \pm 0.2 | 70.41 \pm 0.4 | 81.70 \pm 1.7 | 54.85 \pm 1.7 | 42.34 \pm 1.2 | 93.61 \pm 0.1 | 89.14 \pm 0.5 |
| MVPT | M M | 93.46 \pm 0.1 | 93.72 \pm 0.4 | 77.38 \pm 0.1 | 20.79 \pm 0.1 | 69.43 \pm 0.2 | 92.23 \pm 1.8 | 57.07 \pm 1.6 | 42.57 \pm 0.6 | 88.80 \pm 1.6 | 87.78 \pm 0.1 |
| MUPT | M M | 92.19 \pm 0.8 | 93.75 \pm 0.9 | 75.39 \pm 1.3 | 23.45 \pm 0.3 | 65.99 \pm 1.4 | 90.17 \pm 0.3 | 56.06 \pm 1.6 | 41.19 \pm 0.8 | 89.34 \pm 0.5 | 89.38 \pm 0.1 |
| Δ | | +0.62 | +1.70 | +1.40 | +0.69 | +0.59 | +4.86 | +0.30 | +0.45 | +0.30 | +0.00 |

| | | Target | | | | | | | | | |
|-------------------|------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| Source | Adaptation | GTSRB | Hateful Memes | KITTI Distance | MNIST | OxfordPets | Patch-Camelyon | Rendered-SST2 | Resisc-45 | StanfordCars | VOC 2007 |
| CLIP [†] | - - | 58.1 | 27.0 | 52.0 | 69.4 | 89.0 | 54.0 | 60.9 | 65.6 | 64.8 | 83.7 |
| CoOp | - S | 73.87 \pm 2.0 | 52.40 \pm 2.4 | 56.87 \pm 6.0 | 91.44 \pm 2.9 | 90.69 \pm 1.0 | 62.79 \pm 3.8 | 59.55 \pm 4.9 | 83.83 \pm 1.7 | 79.52 \pm 0.8 | 74.61 \pm 2.4 |
| VPT | - S | 85.34 \pm 1.1 | 56.60 \pm 1.4 | 53.54 \pm 7.1 | 89.88 \pm 2.4 | 90.71 \pm 0.4 | 60.30 \pm 4.8 | 57.66 \pm 4.8 | 84.05 \pm 0.3 | 74.95 \pm 0.8 | 78.88 \pm 1.8 |
| UPT | - S | 82.72 \pm 1.3 | 56.87 \pm 3.6 | 47.87 \pm 11.9 | 89.11 \pm 2.3 | 91.24 \pm 0.6 | 60.41 \pm 1.0 | 59.03 \pm 3.8 | 83.32 \pm 0.2 | 76.40 \pm 0.1 | 81.20 \pm 1.4 |
| MCoOp | - M | 74.38 \pm 0.3 | 58.40 \pm 1.1 | 56.87 \pm 6.0 | 91.44 \pm 2.9 | 90.69 \pm 1.0 | 64.91 \pm 3.4 | 61.63 \pm 2.1 | 84.03 \pm 0.3 | 79.52 \pm 0.8 | 78.45 \pm 2.0 |
| MVPT | - M | 85.34 \pm 1.1 | 58.20 \pm 1.0 | 53.54 \pm 7.1 | 89.88 \pm 2.4 | 91.01 \pm 0.4 | 66.53 \pm 8.8 | 58.14 \pm 2.2 | 84.05 \pm 0.3 | 74.95 \pm 0.8 | 80.69 \pm 1.0 |
| MUPT | - M | 82.72 \pm 1.3 | 58.13 \pm 2.2 | 55.41 \pm 6.6 | 89.91 \pm 2.3 | 91.24 \pm 0.6 | 63.36 \pm 9.8 | 61.34 \pm 2.9 | 83.32 \pm 0.2 | 76.40 \pm 0.1 | 81.20 \pm 1.4 |
| MCoOp | M M | 72.74 \pm 1.1 | 58.40 \pm 1.0 | 47.73 \pm 3.9 | 90.21 \pm 0.1 | 89.61 \pm 0.9 | 68.92 \pm 3.0 | 64.89 \pm 2.7 | 84.39 \pm 0.5 | 79.43 \pm 0.6 | 79.55 \pm 1.6 |
| MVPT | M M | 89.62 \pm 0.9 | 55.53 \pm 1.7 | 62.07 \pm 6.1 | 93.08 \pm 1.7 | 91.04 \pm 0.3 | 69.69 \pm 1.8 | 57.50 \pm 1.0 | 84.35 \pm 0.2 | 74.20 \pm 0.6 | 82.21 \pm 0.5 |
| MUPT | M M | 81.66 \pm 2.7 | 59.00 \pm 1.0 | 57.20 \pm 5.7 | 91.38 \pm 1.5 | 90.30 \pm 1.0 | 69.74 \pm 4.0 | 62.29 \pm 2.4 | 83.40 \pm 0.4 | 76.66 \pm 0.4 | 79.29 \pm 1.8 |
| Δ | | +4.28 | +2.13 | +8.53 | +3.20 | +0.00 | +9.33 | +5.34 | +0.56 | +0.00 | +3.33 |