

# Supplementary Materials for VL-PET: Vision-and-Language Parameter-Efficient Tuning via Granularity Control

Zi-Yuan Hu<sup>1,3</sup>

Yanyang Li<sup>1</sup>

Michael R. Lyu<sup>1</sup>

Liwei Wang<sup>1,2\*</sup>

<sup>1</sup>The Chinese University of Hong Kong    <sup>2</sup>Centre for Perceptual and Interactive Intelligence

<sup>3</sup>Shanghai Artificial Intelligence Laboratory

{zyhu22, yyli21, lyu, lwwang}@cse.cuhk.edu.hk

## A. Dataset Statistics

For image-text tasks, we report the performance on Karpathy test/test-dev/test-P/Karpathy test split for VQA/GQA/NLVR<sup>2</sup>/MSCOCO. For video-text tasks, we report the performance on test-public split for TVQA, How2QA, TVC and YC2C from VALUE benchmark [15]. The detailed statistics of image-text and video-text datasets are shown in Tab. 1 and Tab. 2, respectively.

## B. Implementation Details

Combining a pre-trained vision model with a generative PLM, we perform multi-task learning via unified text generation. Following [2], we use CLIP [18] as our frozen vision encoder for offline visual feature extraction, and BART-base [13] and T5-base [19] as our generative PLM backbones for parameter-efficient tuning. A trainable visual projector (i.e., a trainable linear layer) is utilized to project the visual features into the dimension space of text embeddings. The projected visual features and text embeddings are subsequently concatenated as input into our PLM backbone. For a fair experimental comparison, we follow the setting of existing PET techniques [22, 21] to adopt a shared-weight PET module for all downstream tasks and use CLIP-ResNet101 [18, 5] as our vision encoder for image-text tasks and CLIP (ViT-B/32) [18, 3] for video-text tasks. Specifically, images are resized to  $224 \times 224$  at first, and then  $6 \times 6$  grid features are extracted by the last convolutional layer with an adaptive maximum pooling. Video features are extracted offline at a fixed frame rate (one frame per second) with a maximum length of 64.

We train the models for 20 epochs for both image-text tasks and video-text tasks. The batch size is set as 500 for BART-base on image-text tasks, 300 for T5-base on image-text tasks, and 50 for BART-base on video-text tasks. We use AdamW [16] as our optimizer to train the models and

\*Corresponding author.

Dataset	Statistics (#images / #QA pairs, #captions)		
	Train	Validation	Test
VQA [4]	113.2K/605.1K	5.0K/26.7K	5.0K/26.3K
GQA [8]	72.1K/943.0K	10.2K/132.1K	398/12.6K
NLVR <sup>2</sup> [20]	103.2K/86.4K	8.1K/7.0K	8.1K/7.0K
MSCOCO [1]	113.2K/566.8K	5.0K/5.0K	5.0K/5.0K

Table 1. The statistics of four image-text datasets.

Dataset	Statistics (#videos / #QA pairs, #captions)		
	Train	Validation	Test
TVQA [10]	17.4K/122.0K	2.2K/15.3K	2.2K/15.3K
How2QA [14]	24.5K/34.2K	3.1K/3.1K	3.1K/3.1K
TVC [11]	17.4K/86.7K	10.8K/43.6K	10.8K/43.6K
YC2C [24]	10.3K/10.3K	3.5K/3.5K	1.6K/1.6K

Table 2. The statistics of four video-text datasets.

Task	Visual Input	Text Input with Task Prompt	Text Output
VQA	image features	vqa: [Q]	[A]
GQA	image features	gqa: [Q]	[A]
NLVR	image features	nlvr: [text]	true/false
MSCOCO	image features	caption:	[caption]
TVQA	video features	tvqa: [Q]	[A]
How2QA	video features	how2qa: [Q]	[A]
TVC	video features	tvq:	[caption]
YC2C	video features	yc2c:	[caption]

Table 3. Input-output formats with task prompts from [2, 22]. [Q] denotes question, [A] denotes answer, [text] denotes text and [caption] denotes captioning results.

apply a linear decay scheduler with a warmup ratio of 0.1. We evaluate the last checkpoints of models. All experiments are conducted on one A100 GPU (80G memory). The average training time is 16 hours for BART-base on image-text tasks, 18 hours for T5-base on image-text tasks, and 18 hours for BART-base on video-text tasks.

Method	Learning Rate	Batch Size	Epoch	Other Hyper-Parameters
Full Fine-tuning	$1 \times 10^{-4}$	500	20	-
BitFit [23]	$1 \times 10^{-3}$	500	20	-
Prompt Tuning [12]	$1 \times 10^{-3}$	500	20	prompt length $N_p = 40$ , prompt dimension $d_m = 800$
Compacter [9]	$1 \times 10^{-3}$	500	20	hidden dimension $d = 96$ , Kronecker products $k = 2$
Hyperformer [17]	$1 \times 10^{-3}$	500	20	hidden dimension $d = 96$ , task dimension $d_p = 8$
LoRA [7]	$1 \times 10^{-3}$	500	20	hidden dimension $d = 128$
VL-Adapter [22]	$1 \times 10^{-3}$	500	20	hidden dimension $d = 96$
VL-PET <sub>small</sub>	$1 \times 10^{-3}$	500	20	encoders: $r = 96, s = 1.0, N_h = 4$ ; decoders: $r = 96, s = 1.0, N_h = 1$
VL-PET <sub>middleX</sub>	$1 \times 10^{-3}$	500	20	
VL-PET <sub>middleY</sub>	$1 \times 10^{-3}$	500	20	
VL-PET <sub>large</sub>	$1 \times 10^{-3}$	500	20	

Table 4. Hyper-parameters for image-text tasks with BART-base backbone. We follow the baseline settings of [22].

Method	Learning Rate	Batch size	Epoch	Other hyper-parameters
Full Fine-tuning	$3 \times 10^{-4}$	300	20	-
BitFit [23]	$3 \times 10^{-3}$	300	20	-
Prompt Tuning [12]	$3 \times 10^{-3}$	300	20	prompt length $N_p = 40$ , prompt dimension $d_m = 800$
LoRA [7]	$3 \times 10^{-4}$	300	20	hidden dimension=150
VL-Adapter [22]	$3 \times 10^{-4}$	300	20	hidden dimension $d = 192$
LST [21]	$3 \times 10^{-3}$	300	20	r=4; all layers in side networks are kept
VL-PET <sub>small</sub>	$3 \times 10^{-4}$	300	20	encoders: $r = 192, s = 0.3, N_h = 4$ ; decoders: $r = 96, s = 1.0, N_h = 1$
VL-PET <sub>middleX</sub>	$3 \times 10^{-4}$	300	20	
VL-PET <sub>middleY</sub>	$3 \times 10^{-4}$	300	20	
VL-PET <sub>large</sub>	$3 \times 10^{-4}$	300	20	

Table 5. Hyper-parameters for image-text tasks with T5-base backbone. We follow the baseline settings of [21].

Method	Learning Rate	Batch size	Epoch	Other hyper-parameters
Full Fine-tuning	$1 \times 10^{-5}$	50	20	-
BitFit [23]	$1 \times 10^{-4}$	50	20	-
Prompt Tuning [12]	$1 \times 10^{-4}$	50	20	prompt length $N_p = 40$ , prompt dimension $d_m = 800$
Compacter [9]	$1 \times 10^{-4}$	50	20	hidden dimension $d = 96$ , Kronecker products $k = 2$
LoRA [7]	$1 \times 10^{-4}$	50	20	hidden dimension $d = 128$
VL-Adapter [22]	$1 \times 10^{-4}$	50	20	hidden dimension $d = 96$
VL-PET <sub>small</sub>	$7 \times 10^{-4}$	50	20	encoders: $r = 96, s = 1.0, N_h = 4$ ; decoders: $r = 96, s = 1.0, N_h = 1$
VL-PET <sub>middleX</sub>	$7 \times 10^{-4}$	50	20	
VL-PET <sub>middleY</sub>	$7 \times 10^{-4}$	50	20	
VL-PET <sub>large</sub>	$7 \times 10^{-4}$	50	20	

Table 6. Hyper-parameters for video-text tasks based on BART-base backbone.

## C. Hyper-Parameters

For image-text tasks with the BART-base backbone, [22] conduct a detailed hyper-parameter search to find the optimal settings for the baselines. However, it only reports experimental results of one random seed. To strengthen the reliability of these experiments, we summarize statistics of three seeds with the hyper-parameters listed in Tab. 4.

For image-text tasks with the T5-base backbone, we borrow the results of three seeds from [21]. The hyper-parameters of the baselines are listed in Tab. 5.

For video-text tasks with the BART-base backbone, we show the results of only one seed due to the submission limit of the VALUE benchmark. The hyper-parameters are listed in Tab. 6.

## D. Effectiveness of Visual Features and Trainable Visual Projector

In this section, we investigate the effectiveness of the visual features and a trainable visual projector on VL tasks. We denote T as text, Noise as noise features sampled from

Input	Params (%)	VQA (%)	GQA (%)	NLVR <sup>2</sup> (%)	COCO (CIDEr)	Avg.
VL-PET <sub>large</sub> (BART-base)						
T	3.10	44.80 <sub>0.29</sub>	40.19 <sub>0.16</sub>	51.09 <sub>0.07</sub>	6.66 <sub>0.41</sub>	35.69 <sub>0.03</sub>
T + Noise	4.16	44.79 <sub>0.10</sub>	40.49 <sub>0.16</sub>	51.09 <sub>0.22</sub>	6.57 <sub>0.98</sub>	35.74 <sub>0.16</sub>
T + Frozen V	3.07	62.62 <sub>0.50</sub>	51.34 <sub>0.67</sub>	67.77 <sub>0.39</sub>	117.58 <sub>1.90</sub>	74.83 <sub>0.72</sub>
T + Trainable V	4.16	66.17 <sub>0.27</sub>	55.11 <sub>0.17</sub>	73.43 <sub>0.35</sub>	122.03 <sub>0.46</sub>	79.18 <sub>0.14</sub>
VL-PET <sub>large</sub> (T5-base)						
T	6.70	44.90 <sub>0.12</sub>	40.62 <sub>0.11</sub>	51.21 <sub>0.09</sub>	4.81 <sub>0.77</sub>	35.39 <sub>0.19</sub>
T + Noise	7.31	44.85 <sub>0.18</sub>	40.36 <sub>0.02</sub>	51.22 <sub>0.09</sub>	6.45 <sub>1.81</sub>	35.72 <sub>0.42</sub>
T + Frozen V	6.65	64.30 <sub>0.20</sub>	53.18 <sub>0.55</sub>	69.90 <sub>0.68</sub>	117.90 <sub>0.37</sub>	76.32 <sub>0.44</sub>
T + Trainable V	7.31	66.95 <sub>0.21</sub>	56.06 <sub>0.21</sub>	73.42 <sub>0.46</sub>	121.66 <sub>0.06</sub>	79.52 <sub>0.21</sub>

Table 7. Effectiveness of visual features and visual projector. (T: text. Noise: noise features. V: visual projector.)

Method	Params (%)	VQA (%)	GQA (%)	NLVR <sup>2</sup> (%)	COCO (CIDEr)	Avg.
VL-PET <sub>large</sub> with $\Delta\mathbf{H}^{\text{vis}}$	3.29	65.46 <sub>0.15</sub>	53.94 <sub>0.27</sub>	72.89 <sub>0.32</sub>	120.15 <sub>0.68</sub>	78.11 <sub>0.13</sub>
VL-PET <sub>large</sub> with scaled up $\Delta\mathbf{H}^{\text{vis}}$	3.47	65.57 <sub>0.03</sub>	53.70 <sub>0.70</sub>	73.15 <sub>0.41</sub>	120.00 <sub>0.24</sub>	78.11 <sub>0.12</sub>
VL-PET <sub>large</sub> with $\Delta\mathbf{H}^{\text{vis}}$ and $\mathbf{G}_{\text{large}}^{\text{vis}}$	3.47	65.50 <sub>0.13</sub>	54.03 <sub>0.44</sub>	73.35 <sub>0.18</sub>	120.71 <sub>0.43</sub>	78.40 <sub>0.03</sub>
VL-PET <sub>large</sub>	4.16	66.17 <sub>0.27</sub>	55.11 <sub>0.17</sub>	73.43 <sub>0.35</sub>	122.03 <sub>0.46</sub>	79.18 <sub>0.14</sub>

Table 8. Decomposing the visual projector on BART-base.

Method	Updating Formula	Params (%)	VQA (%)	GQA (%)	NLVR <sup>2</sup> (%)	COCO (CIDEr)	Avg.
VL-PET w/o granularity control	$\mathbf{H} \leftarrow \mathbf{H} + \Delta\mathbf{H}$	2.97	65.22 <sub>0.14</sub>	53.35 <sub>0.39</sub>	72.65 <sub>0.44</sub>	120.19 <sub>0.68</sub>	77.85 <sub>0.34</sub>
VL-PET <sub>add</sub>	$\mathbf{H} \leftarrow \mathbf{H} + \Delta\mathbf{H} + \mathbf{G}_{\text{large}}$	4.16	65.10 <sub>0.05</sub>	53.26 <sub>0.26</sub>	71.85 <sub>0.19</sub>	121.26 <sub>1.19</sub>	77.87 <sub>0.26</sub>
VL-PET <sub>large</sub>	$\mathbf{H} \leftarrow \mathbf{G}_{\text{large}} \odot (\mathbf{H} + \Delta\mathbf{H})$	4.16	66.17 <sub>0.27</sub>	55.11 <sub>0.17</sub>	73.43 <sub>0.35</sub>	122.03 <sub>0.46</sub>	79.18 <sub>0.14</sub>

Table 9. Performance improvement effect of granularity-controlled mechanism.

a uniform distribution on the interval from 0 to 1, Trainable V as visual features with a trainable visual projector and Frozen V as visual features with a frozen visual projector. We replace the standard model input (i.e., T + Trainable V) with T (text-only), T + Noise, and T + Frozen V. In Tab. 7, T + Trainable V outperforms other inputs by a large margin, demonstrating the effectiveness of visual features and the trainable visual projector, as well as the importance of VL alignment and modeling on VL tasks.

In Tab. 8, we also decompose the visual projector using multi-head modular modifications (denoted as  $\Delta\mathbf{H}^{\text{vis}}$ ,  $r=96$ ,  $N_h=4$ ) and the granularity control at large level (denoted as  $\mathbf{G}_{\text{large}}^{\text{vis}}$ ). Due to the different dimensions of visual and text features, we have to remove the residual connection for  $\Delta\mathbf{H}^{\text{vis}}$ . The results indicate that simply scaling up  $\Delta\mathbf{H}^{\text{vis}}$  (from  $r=96$  to  $r=192$ ) dose not effectively improve the performance, whereas introducing  $\mathbf{G}_{\text{large}}^{\text{vis}}$  into  $\Delta\mathbf{H}^{\text{vis}}$  is more effective. Although decomposing the visual projector provides us more efficiency and effectiveness trade-offs, we do not include it in the main paper for a fair competition with the existing PET techniques (e.g., VL-Adapter).

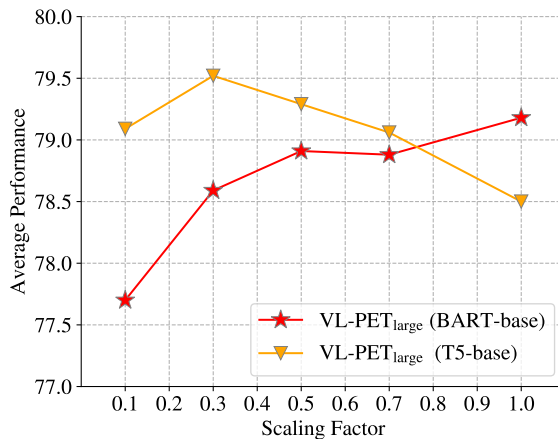


Figure 1. Effectiveness of scaling factor for the encoders. Experiments are conducted on image-text tasks based on different PLM backbones with one seed.

Method	Trainable Params (%)	VQA Acc. (%)	GQA Acc. (%)	NLVR <sup>2</sup> Acc. (%)	COCO Cap. (CIDEr)	Avg.
<b>Backbone: BART-base</b>						
<b>Without task prompts</b>						
VL-PET <sub>small</sub>	2.98	64.83 <sub>0.09</sub>	54.23 <sub>0.26</sub>	72.27 <sub>0.06</sub>	121.03 <sub>0.28</sub>	78.09 <sub>0.10</sub>
VL-PET <sub>middleX</sub>	2.98	65.14 <sub>0.16</sub>	54.55 <sub>0.09</sub>	72.77 <sub>0.27</sub>	120.91 <sub>0.33</sub>	78.34 <sub>0.07</sub>
VL-PET <sub>middleY</sub>	2.98	64.69 <sub>0.14</sub>	53.40 <sub>0.35</sub>	73.04 <sub>0.16</sub>	120.14 <sub>0.58</sub>	77.82 <sub>0.12</sub>
VL-PET <sub>large</sub>	4.16	65.78 <sub>0.08</sub>	54.45 <sub>0.32</sub>	72.90 <sub>0.36</sub>	121.46 <sub>0.36</sub>	78.65 <sub>0.22</sub>
<b>With task prompts</b>						
VL-PET <sub>small</sub>	2.98	65.43 <sub>0.06</sub>	54.03 <sub>0.14</sub>	72.43 <sub>0.22</sub>	120.68 <sub>0.35</sub>	78.14 <sub>0.11</sub>
VL-PET <sub>middleX</sub>	2.98	65.54 <sub>0.09</sub>	54.53 <sub>0.15</sub>	72.66 <sub>0.17</sub>	120.72 <sub>0.51</sub>	78.37 <sub>0.14</sub>
VL-PET <sub>middleY</sub>	2.98	65.36 <sub>0.15</sub>	53.83 <sub>0.39</sub>	73.43 <sub>0.78</sub>	120.31 <sub>0.09</sub>	78.23 <sub>0.19</sub>
VL-PET <sub>large</sub>	4.16	66.17 <sub>0.27</sub>	55.11 <sub>0.17</sub>	73.43 <sub>0.35</sub>	122.03 <sub>0.46</sub>	79.18 <sub>0.14</sub>
<b>Backbone: T5-base</b>						
<b>Without task prompts</b>						
VL-PET <sub>small</sub>	4.51	65.67 <sub>0.31</sub>	56.09 <sub>0.29</sub>	73.54 <sub>0.63</sub>	120.51 <sub>1.11</sub>	78.96 <sub>0.57</sub>
VL-PET <sub>middleX</sub>	4.50	65.68 <sub>0.17</sub>	56.57 <sub>0.10</sub>	74.05 <sub>0.25</sub>	119.92 <sub>0.65</sub>	79.06 <sub>0.15</sub>
VL-PET <sub>middleY</sub>	4.50	65.75 <sub>0.34</sub>	56.35 <sub>0.41</sub>	73.93 <sub>0.60</sub>	119.84 <sub>1.08</sub>	78.97 <sub>0.28</sub>
VL-PET <sub>large</sub>	7.31	66.31 <sub>0.06</sub>	56.56 <sub>0.27</sub>	73.61 <sub>0.22</sub>	121.95 <sub>0.09</sub>	79.61 <sub>0.08</sub>
<b>With task prompts</b>						
VL-PET <sub>small</sub>	4.51	65.88 <sub>0.31</sub>	54.96 <sub>1.01</sub>	72.64 <sub>0.09</sub>	120.05 <sub>0.41</sub>	78.38 <sub>0.37</sub>
VL-PET <sub>middleX</sub>	4.50	66.63 <sub>0.14</sub>	55.87 <sub>0.25</sub>	74.11 <sub>0.37</sub>	120.41 <sub>0.31</sub>	79.26 <sub>0.26</sub>
VL-PET <sub>middleY</sub>	4.50	66.62 <sub>0.20</sub>	55.87 <sub>0.13</sub>	73.91 <sub>0.45</sub>	120.26 <sub>0.40</sub>	79.17 <sub>0.08</sub>
VL-PET <sub>large</sub>	7.31	66.95 <sub>0.21</sub>	56.06 <sub>0.21</sub>	73.42 <sub>0.46</sub>	121.66 <sub>0.06</sub>	79.52 <sub>0.21</sub>

Table 10. Effectiveness of task prompts.

## E. Performance Improvement Effect of Granularity-controlled Mechanism

Granularity control can dynamically assign importance weights to the intermediate hidden states. To fully investigate where the performance improvement is coming from, we replace the element-wise product with addition for  $\mathbf{G}_{\text{large}}$  on BART-base, denoted as VL-PET<sub>add</sub>. The results (VL-PET<sub>large</sub> > VL-PET<sub>add</sub>) in Tab. 9 verify that the importance assignment is more effective than the added parameters in performance improvement. Moreover, we observe that VL-PET<sub>add</sub> introduces more trainable parameters but perform on par with VL-PET without granularity control, as their updating formulas are equivalent to conventional PET (Eq. (1)). These results further demonstrate the effectiveness of the proposed granularity control.

## F. Effectiveness of Scaling Factor

As mentioned in Sec. 3.2, the scaling factor is specialized for different PLMs. Our experimental results in Fig. 1 show that the optimal scaling factors for BART-base en-

coders and T5-base encoders are 1.0 and 0.3, respectively.

## G. Unified Text Generation with Task Prompts

Inspired by [22, 21], we prepend a task-specific prompt to the input sentence for each downstream task. The detailed task prompts are listed in Tab. 3. We conduct experiments to test the validity of these task prompts over three seeds in Tab. 10. For the BART-base backbone, all VL-PET modules with task prompts outperform those without task prompts. However, for T5-base, some VL-PET modules without task prompts exhibit superior performance compared to those with task prompts. For a fair comparison with state-of-the-art PET techniques, we still adopt task prompts to facilitate multi-task learning via unified text generation for our generative PLM backbones.

## H. Scalability of VL-PET Modules

We scale up the trainable parameters of the VL-PET modules by increasing the projected hidden dimensions  $r$  of the Encoder VL-PET modules in BART-base. All scaled

Method	Params (%)	VQA (%)	GQA (%)	NLVR <sup>2</sup> (%)	COCO (CIDEr)	Avg.
VL-PET <sub>small</sub> ( $r = 96$ )	2.98	65.43 <sub>0.06</sub>	54.03 <sub>0.14</sub>	72.43 <sub>0.22</sub>	120.68 <sub>0.35</sub>	78.14 <sub>0.11</sub>
$r = 144$	3.58	65.87 <sub>0.12</sub>	54.05 <sub>0.21</sub>	72.76 <sub>0.09</sub>	121.22 <sub>0.40</sub>	78.48 <sub>0.07</sub>
$r = 192$	4.16	66.14 <sub>0.16</sub>	54.73 <sub>0.18</sub>	72.63 <sub>0.23</sub>	121.46 <sub>1.38</sub>	78.74 <sub>0.34</sub>
VL-PET <sub>middleX</sub> ( $r = 96$ )	2.98	65.54 <sub>0.09</sub>	54.53 <sub>0.15</sub>	72.66 <sub>0.17</sub>	120.72 <sub>0.51</sub>	78.37 <sub>0.14</sub>
$r = 144$	3.57	66.08 <sub>0.10</sub>	54.82 <sub>0.29</sub>	72.82 <sub>0.23</sub>	121.05 <sub>0.11</sub>	78.70 <sub>0.11</sub>
$r = 192$	4.16	66.23 <sub>0.06</sub>	54.55 <sub>0.33</sub>	73.57 <sub>0.34</sub>	122.01 <sub>0.65</sub>	79.09 <sub>0.23</sub>
VL-PET <sub>middleY</sub> ( $r = 96$ )	2.98	65.36 <sub>0.15</sub>	53.83 <sub>0.39</sub>	73.43 <sub>0.78</sub>	120.31 <sub>0.09</sub>	78.23 <sub>0.19</sub>
$r = 144$	3.57	65.58 <sub>0.18</sub>	54.06 <sub>0.32</sub>	73.02 <sub>0.52</sub>	120.59 <sub>0.26</sub>	78.31 <sub>0.30</sub>
$r = 192$	4.16	65.94 <sub>0.06</sub>	54.29 <sub>0.40</sub>	73.53 <sub>0.25</sub>	120.89 <sub>0.30</sub>	78.66 <sub>0.18</sub>
VL-PET <sub>large</sub> ( $r = 96$ )	4.16	66.17 <sub>0.27</sub>	55.11 <sub>0.17</sub>	73.43 <sub>0.35</sub>	122.03 <sub>0.46</sub>	79.18 <sub>0.14</sub>
$r = 144$	5.31	66.47 <sub>0.10</sub>	55.07 <sub>0.25</sub>	73.21 <sub>0.15</sub>	122.05 <sub>0.58</sub>	79.20 <sub>0.13</sub>
$r = 192$	6.43	66.72 <sub>0.06</sub>	54.61 <sub>0.31</sub>	73.55 <sub>0.88</sub>	122.16 <sub>0.18</sub>	79.26 <sub>0.12</sub>

Table 11. Scalability of VL-PET Modules. ( $r$ : the projected hidden dimension of the encoder VL-PET modules.)

Method	Trainable Params (%)	VQA Acc. (%)	GQA Acc. (%)	NLVR <sup>2</sup> Acc. (%)	COCO Cap. (CIDEr)	Avg.
<b>Backbone: BART-base</b>						
<b>Random Gaussian Initialization (Default)</b>						
VL-PET <sub>small</sub>	2.98	65.43 <sub>0.06</sub>	54.03 <sub>0.14</sub>	72.43 <sub>0.22</sub>	120.68 <sub>0.35</sub>	78.14 <sub>0.11</sub>
VL-PET <sub>middleX</sub>	2.98	65.54 <sub>0.09</sub>	54.53 <sub>0.15</sub>	72.66 <sub>0.17</sub>	120.72 <sub>0.51</sub>	78.37 <sub>0.14</sub>
VL-PET <sub>middleY</sub>	2.98	65.36 <sub>0.15</sub>	53.83 <sub>0.39</sub>	73.43 <sub>0.78</sub>	120.31 <sub>0.09</sub>	78.23 <sub>0.19</sub>
VL-PET <sub>large</sub>	4.16	66.17 <sub>0.27</sub>	55.11 <sub>0.17</sub>	73.43 <sub>0.35</sub>	122.03 <sub>0.46</sub>	79.18 <sub>0.14</sub>
<b>Zero Initialization</b>						
VL-PET <sub>small</sub>	2.98	65.33 <sub>0.21</sub>	53.96 <sub>0.15</sub>	72.48 <sub>0.10</sub>	121.05 <sub>0.28</sub>	78.21 <sub>0.13</sub>
VL-PET <sub>middleX</sub>	2.98	65.49 <sub>0.06</sub>	54.28 <sub>0.20</sub>	73.01 <sub>0.47</sub>	120.36 <sub>0.26</sub>	78.29 <sub>0.15</sub>
VL-PET <sub>middleY</sub>	2.98	65.21 <sub>0.20</sub>	53.86 <sub>0.53</sub>	72.84 <sub>1.21</sub>	120.38 <sub>0.65</sub>	78.07 <sub>0.30</sub>
VL-PET <sub>large</sub>	4.16	66.18 <sub>0.13</sub>	54.44 <sub>0.58</sub>	72.99 <sub>0.69</sub>	121.91 <sub>0.52</sub>	78.88 <sub>0.35</sub>
<b>Backbone: T5-base</b>						
<b>Random Gaussian Initialization</b>						
VL-PET <sub>small</sub>	4.51	66.80 <sub>0.22</sub>	56.15 <sub>0.22</sub>	74.25 <sub>0.59</sub>	120.57 <sub>0.81</sub>	79.44 <sub>0.45</sub>
VL-PET <sub>middleX</sub>	4.50	66.33 <sub>0.47</sub>	55.89 <sub>0.05</sub>	74.23 <sub>0.21</sub>	119.79 <sub>0.59</sub>	79.06 <sub>0.31</sub>
VL-PET <sub>middleY</sub>	4.50	66.52 <sub>0.12</sub>	55.52 <sub>0.65</sub>	73.96 <sub>0.45</sub>	120.99 <sub>0.21</sub>	79.25 <sub>0.05</sub>
VL-PET <sub>large</sub>	7.31	66.78 <sub>0.10</sub>	55.62 <sub>0.18</sub>	73.21 <sub>0.15</sub>	121.12 <sub>0.36</sub>	79.19 <sub>0.08</sub>
<b>Zero Initialization (Default)</b>						
VL-PET <sub>small</sub>	4.51	65.88 <sub>0.31</sub>	54.96 <sub>1.01</sub>	72.64 <sub>0.09</sub>	120.05 <sub>0.41</sub>	78.38 <sub>0.37</sub>
VL-PET <sub>middleX</sub>	4.50	66.63 <sub>0.14</sub>	55.87 <sub>0.25</sub>	74.11 <sub>0.37</sub>	120.41 <sub>0.31</sub>	79.26 <sub>0.26</sub>
VL-PET <sub>middleY</sub>	4.50	66.62 <sub>0.20</sub>	55.87 <sub>0.13</sub>	73.91 <sub>0.45</sub>	120.26 <sub>0.40</sub>	79.17 <sub>0.08</sub>
VL-PET <sub>large</sub>	7.31	66.95 <sub>0.21</sub>	56.06 <sub>0.21</sub>	73.42 <sub>0.46</sub>	121.66 <sub>0.06</sub>	79.52 <sub>0.21</sub>

Table 12. Experimental results of different weight initialization strategies.

VL-PET modules outperform the unscaled ones in Tab. 11, indicating the scalability of our VL-PET framework.

## I. Weight Initialization

In Tab. 12, we show the experimental results over three seeds with two popular weight initialization strategies, i.e., random Gaussian initialization and zero initialization, for

both BART-base and T5-base backbones. Random Gaussian initialization is widely used in PET techniques to initialize the weight of PET modules from a Gaussian distribution, while some PET techniques [7] utilize zero initialization to set the up projection layers of the PET modules as zero. Based on Tab. 12, it can be inferred that random Gaussian initialization is more appropriate for BART-base and zero initialization for T5-base. Therefore, we employ

Method	Trainable Params (%)	VQA Acc. (%)	GQA Acc. (%)	NLVR <sup>2</sup> Acc. (%)	COCO Cap. (CIDEr)	Avg.
<b>Backbone: BART-base</b>						
<b>Granularity-controlled Mechanism: <math>G_{\text{large}}</math></b>						
<b>Down Multi-head Modular Modifications</b> (used in the main paper)						
$N_h = 2$	4.16	66.20 <sub>0.12</sub>	54.74 <sub>0.43</sub>	72.88 <sub>0.43</sub>	121.53 <sub>0.64</sub>	78.84 <sub>0.26</sub>
$N_h = 4$	4.16	66.17 <sub>0.27</sub>	55.11 <sub>0.17</sub>	73.43 <sub>0.35</sub>	122.03 <sub>0.46</sub>	79.18 <sub>0.14</sub>
<b>Up Multi-head Modular Modifications</b>						
$N_h = 2$	4.16	66.16 <sub>0.12</sub>	54.88 <sub>0.21</sub>	73.19 <sub>0.13</sub>	121.69 <sub>0.62</sub>	78.98 <sub>0.14</sub>
$N_h = 4$	4.16	66.15 <sub>0.13</sub>	55.00 <sub>0.36</sub>	73.19 <sub>0.27</sub>	121.44 <sub>0.73</sub>	78.95 <sub>0.23</sub>
<b>Down-Up Multi-head Modular Modifications</b>						
$N_h = 2$	4.16	66.24 <sub>0.10</sub>	54.73 <sub>0.24</sub>	72.81 <sub>0.21</sub>	121.56 <sub>0.15</sub>	78.84 <sub>0.07</sub>
$N_h = 4$	4.16	66.18 <sub>0.16</sub>	54.60 <sub>0.27</sub>	73.17 <sub>0.42</sub>	122.05 <sub>0.20</sub>	79.00 <sub>0.10</sub>
<b>Down-Up-Pair Multi-head Modular Modifications</b>						
$N_h = 2$	3.86	65.86 <sub>0.04</sub>	54.51 <sub>0.36</sub>	72.98 <sub>0.19</sub>	121.65 <sub>0.21</sub>	78.75 <sub>0.07</sub>
$N_h = 4$	3.72	65.70 <sub>0.09</sub>	54.24 <sub>0.12</sub>	72.43 <sub>0.35</sub>	122.76 <sub>0.33</sub>	78.28 <sub>0.20</sub>

Table 13. Experimental results of multi-head modular modifications with various numbers of heads.

random Gaussian initialization for BART-base and zero initialization for T5-base by default in this work.

## J. More Designs for Multi-head Modular Modification

We propose a multi-head modular modification  $\Delta\mathbf{H}' \in \mathbb{R}^{N \times d}$  of length  $N$  and dimension  $d$  in the main paper, which is formulated as follows:

$$\Delta\mathbf{H}' = \phi(\text{Concat}(\mathbf{X}'\mathbf{W}_{\text{down}}^{(1)}, \dots, \mathbf{X}'\mathbf{W}_{\text{down}}^{(N_h)}))\mathbf{W}_{\text{up}}, \quad (1)$$

where  $N_h$  is the number of heads,  $\mathbf{X}' \in \mathbb{R}^{N \times d}$  is the input,  $\mathbf{W}_{\text{down}}^{(i)} \in \mathbb{R}^{d \times \frac{r}{N_h}}$  is a down projection layer for the  $i$ -th head,  $\phi$  is the GELU function [6],  $\mathbf{W}_{\text{up}} \in \mathbb{R}^{r \times d}$  is a up projection layer and  $r$  is the projected hidden dimension.

Since the multi-head idea is applied to the down projection layer in the main paper, we can call this type of multi-head modular modification as down multi-head modular modification. As the multi-head idea can be applied to any linear projection within a modular modification, it yields a series of variant designs. In this section, we present additional designs of multi-head modular modification, such as up multi-head modular modification, down-up multi-head modular modification and down-up-pair multi-head modular modification. The primary difference among these four designs lies in where the multi-head idea is employed.

For up multi-head modular modification, we apply the multi-head idea to the up projection layer and describe it as follows:

$$\Delta\mathbf{H}' = \text{Concat}(\phi((\mathbf{X}'\mathbf{W}_{\text{down}}))\mathbf{W}_{\text{up}}^{(1)}, \dots, \phi((\mathbf{X}'\mathbf{W}_{\text{down}}))\mathbf{W}_{\text{up}}^{(N_h)}), \quad (2)$$

where  $\mathbf{W}_{\text{down}} \in \mathbb{R}^{d \times r}$  is a down projection layer and  $\mathbf{W}_{\text{up}}^{(i)} \in \mathbb{R}^{r \times \frac{d}{N_h}}$  is a down projection layer for head $_i$ .

For down-up multi-head modular modification, we apply the multi-head idea to the down projection layer and up projection layer, respectively. The formulation is described as follows:

$$\begin{aligned} \Delta\mathbf{H}'_{\text{down}} &= \phi(\text{Concat}(\mathbf{X}'\mathbf{W}_{\text{down}}^{(1)}, \dots, \mathbf{X}'\mathbf{W}_{\text{down}}^{(N_h)})), \\ \Delta\mathbf{H}' &= \text{Concat}(\Delta\mathbf{H}'_{\text{down}}\mathbf{W}_{\text{up}}^{(1)}, \dots, \Delta\mathbf{H}'_{\text{down}}\mathbf{W}_{\text{up}}^{(N_h)}), \end{aligned} \quad (3)$$

where  $\Delta\mathbf{H}'_{\text{down}} \in \mathbb{R}^{N \times r}$  is an intermediate output.

For down-up-pair multi-head modular modification, we formulate it as follows:

$$\begin{aligned} \Delta\mathbf{H}'_{\text{pair}}^{(i)} &= \phi(\mathbf{X}'\mathbf{W}_{\text{down-pair}}^{(i)})\mathbf{W}_{\text{up-pair}}^{(i)}, \\ \Delta\mathbf{H}' &= \text{Concat}(\Delta\mathbf{H}'_{\text{pair}}^{(1)}, \dots, \Delta\mathbf{H}'_{\text{pair}}^{(N_h)}), \end{aligned} \quad (4)$$

where  $\mathbf{W}_{\text{down-pair}}^{(i)} \in \mathbb{R}^{d \times \frac{r}{N_h}}$  is a down projection layer,  $\mathbf{W}_{\text{up-pair}}^{(i)} \in \mathbb{R}^{\frac{r}{N_h} \times \frac{d}{N_h}}$  is a up projection layer and  $\Delta\mathbf{H}'_{\text{pair}}^{(i)} \in \mathbb{R}^{N \times \frac{d}{N_h}}$  is an intermediate output for head $_i$ .

Similar to the usage of multi-head modular modifications in the main paper, we perform experiments on image-text tasks with BART-base backbone for these four multi-head modular modifications with our proposed granularity-controlled mechanism at a large level.

As shown in Tab. 13, the down multi-head modular modification with  $N_h = 4$  achieves the best performance in the experiment, which is adopted in the main paper. We observe that even the worst down-up-pair multi-head modular modification still outperforms the state-of-the-art PET techniques (e.g., VL-Adapter). Moreover, down-up-pair multi-head modular modification significantly reduces the number of trainable parameters as the number of heads increases, indicating a promising direction for further parameter reduction with a minimal impact on performance.

## K. Limitations

In this section, we discuss some limitations of our work.

- We perform extensive experiments and analysis on image-text tasks and video-text tasks. However, the video-text experiments are conducted with only one seed due to the submission limit of the VALUE benchmark, which may affect the reliability of the video-text experimental results.
- Our VL-PET framework focuses on challenging VL downstream tasks, including some discriminative and generative tasks (e.g., question-answering tasks and captioning tasks). But there are many other downstream tasks in the real life, and we do not include all types of tasks in the multi-task learning experiments. Therefore, our designs and experimental results are not always guaranteed to be generalized to other VL tasks (e.g., image-text retrieval and video-text retrieval) and other domains (e.g., NLP and CV).
- We validate the enhanced effect of employing our VL-PET designs (e.g., granularity-controlled mechanism and lightweight PET module designs) on existing PET techniques (i.e., Compacter and VL-Adapter). But our VL-PET designs may not be applicable to enhance all PET techniques.
- PET techniques are proposed to reduce the number of trainable parameters and save model storage space. Our work utilizes multi-task learning to achieve further reduction. However, similar to most PET techniques, our work integrates newly-designed modules into large PLM backbones, which still evoke a lot of memory and time consumption during training and inference.

## References

- [1] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO captions: Data collection and evaluation server. *CoRR*, abs/1504.00325, 2015.
- [2] Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. Unifying vision-and-language tasks via text generation. In *ICML*, 2021.
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [4] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [6] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- [7] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR*, 2022.
- [8] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, 2019.
- [9] Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank hypercomplex adapter layers. In *NeurIPS*, 2021.
- [10] Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L. Berg. Tvqa: Localized, compositional video question answering. In *EMNLP*, 2018.
- [11] Jie Lei, Licheng Yu, Tamara L. Berg, and Mohit Bansal. Tvr: A large-scale dataset for video-subtitle moment retrieval. In *ECCV*, 2020.
- [12] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *EMNLP*, 2021.
- [13] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *ACL*, 2020.
- [14] Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu. Hero: Hierarchical encoder for video+language omni-representation pre-training. In *EMNLP*, 2020.
- [15] Linjie Li, Jie Lei, Zhe Gan, Licheng Yu, Yen-Chun Chen, Rohit Pillai, Yu Cheng, Luowei Zhou, Xin Eric Wang, William Yang Wang, et al. Value: A multi-task benchmark for video-and-language understanding evaluation. In *NeurIPS*, 2021.
- [16] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2017.
- [17] Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. In *ACL*, 2021.
- [18] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- [19] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2020.
- [20] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. In *ACL*, 2019.

- [21] Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. LST: ladder side-tuning for parameter and memory efficient transfer learning. *CoRR*, abs/2206.06522, 2022.
- [22] Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. VI-adapter: Parameter-efficient transfer learning for vision-and-language tasks. In *CVPR*, 2022.
- [23] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. In *ACL*, 2022.
- [24] Luowei Zhou, Chenliang Xu, and Jason J. Corso. Towards automatic learning of procedures from web instructional videos. In *AAAI*, 2018.