

VL-PET: Vision-and-Language Parameter-Efficient Tuning via Granularity Control

Zi-Yuan Hu, Yanyang Li, Michael R. Lyu, Liwei Wang

Background

- As the model size of **pre-trained language models (PLMs)** grows rapidly, full fine-tuning becomes prohibitively expensive for model training and storage.
- In **vision-and-language (VL)**, **parameter-efficient tuning (PET)** techniques are proposed to integrate **modular modifications** (e.g., **Adapter** and **LoRA**) into encoder-decoder PLMs.

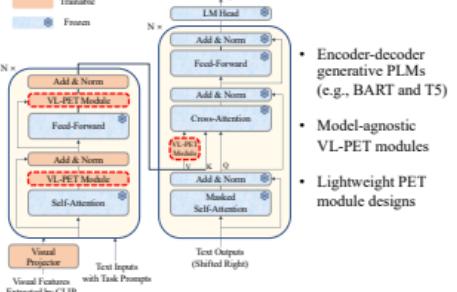
Critical Issues Neglected by Existing Methods

- Integrating **heavy and excessive** modular modifications into PLMs can greatly affect the intermediate output of the PLMs, leading to instability and performance degradation.
- For PLMs used in VL tasks, there exists **functionality gap between the encoders and decoders**.

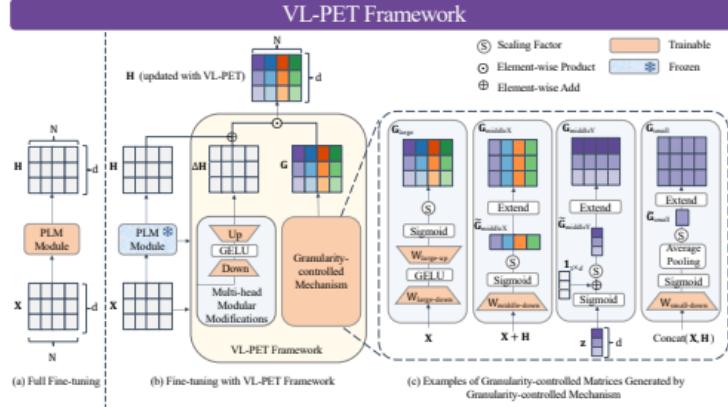
Main Contribution

- We propose a **Vision-and-Language Parameter-Efficient Tuning (VL-PET)** framework to impose effective control over modular modifications via a novel **granularity-controlled mechanism**.
- Considering different **granularity-controlled** matrices generated by this mechanism, a variety of **model-agnostic VL-PET modules** can be instantiated from our framework for better efficiency and effectiveness trade-offs.
- We further propose **lightweight PET module designs** to enhance **VL alignment** and **modeling** for the encoders and **maintain text generation** for the decoders.
- Extensive experiments conducted on four image-text tasks and four video-text tasks demonstrate the **efficiency, effectiveness and transferability** of our VL-PET framework.

Model Architecture



- Encoder-decoder generative PLMs (e.g., BART and T5)
- Model-agnostic VL-PET modules
- Lightweight PET module designs



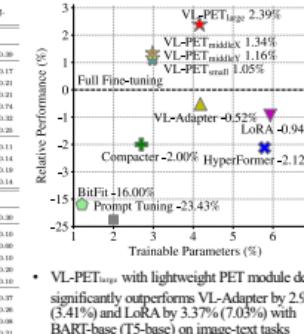
- Unified formulation of PET:** $H \leftarrow H + \Delta H$
 H : intermediate hidden state from a PLM module
 ΔH : modular modifications (e.g., Adapter and LoRA)
- Granularity-controlled mechanism:** $H \leftarrow G \odot (H + \Delta H)$
- Multi-head modular modifications:**

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

Experiment Results on Image-Text Tasks

Method	Trainable Params (%)	VQA (Acc. %)	GQA (Acc. %)	NLVR ² (Acc. %)	COCO (CIDEr)	Avg.
Backbone: BART-base						
Full Fine-tuning*	100	56.48 ± 0.10	56.79 ± 0.10	73.66 ± 0.21	112.0 ± 0.21	77.33 ± 0.10
BitFit [5]	1.21	52.04 ± 0.32	45.15 ± 0.34	52.29 ± 0.31	111.44 ± 0.41	50.90 ± 0.17
Prompt Tuning [20]	2.00	44.12 ± 0.26	46.85 ± 0.26	57.30 ± 0.21	105.02 ± 0.21	44.53 ± 0.17
Compacter [21]	2.70	64.85 ± 0.06	56.37 ± 0.21	71.11 ± 0.21	114.69 ± 0.21	75.75 ± 0.21
HyperFormer [33]	5.79	64.82 ± 0.27	55.25 ± 0.44	70.74 ± 0.43	114.84 ± 0.38	76.95 ± 0.74
LoRA* [16]	5.93	65.15 ± 0.16	56.06 ± 0.41	72.58 ± 0.13	115.01 ± 0.22	76.60 ± 0.32
VL-Adapter [22]	4.18	65.76 ± 0.26	54.16 ± 0.41	73.19 ± 0.13	114.61 ± 0.26	76.95 ± 0.34
VL-Adapter* [22]	2.98	65.43 ± 0.14	54.03 ± 0.14	72.49 ± 0.23	120.68 ± 0.33	78.14 ± 0.11
VL-PET _{small}	2.98	65.54 ± 0.06	54.53 ± 0.13	72.66 ± 0.17	120.72 ± 0.28	78.30 ± 0.14
VL-PET _{medium}	2.98	65.36 ± 0.13	53.83 ± 0.30	73.43 ± 0.21	120.31 ± 0.39	78.20 ± 0.14
VL-PET _{large}	4.16	66.17 ± 0.27	55.11 ± 0.17	73.43 ± 0.25	122.03 ± 0.41	79.18 ± 0.14
Backbone: T5-base						
Full Fine-tuning*	100	67.10 ± 0.20	56.30 ± 0.20	74.29 ± 0.21	112.20 ± 0.21	77.50 ± 0.20
BitFit [5]	0.83	55.00 ± 0.24	45.50 ± 0.20	51.70 ± 0.21	101.20 ± 0.21	63.44 ± 0.10
Prompt Tuning [20]	1.26	47.40 ± 0.24	40.60 ± 0.20	50.00 ± 0.22	96.10 ± 0.20	58.84 ± 0.10
LoRA* [16]	7.54	63.70 ± 0.24	53.30 ± 0.20	70.00 ± 0.21	110.30 ± 0.21	74.30 ± 0.20
VL-Adapter [22]	7.98	67.10 ± 0.24	56.00 ± 0.20	72.70 ± 0.21	111.80 ± 0.21	76.90 ± 0.20
LST [23]	7.46	66.50 ± 0.24	55.90 ± 0.20	71.60 ± 0.21	113.30 ± 0.21	76.90 ± 0.20
VL-PET _{small}	4.51	65.88 ± 0.24	54.96 ± 0.21	72.64 ± 0.21	120.05 ± 0.21	78.38 ± 0.21
VL-PET _{medium}	4.50	66.83 ± 0.14	58.87 ± 0.27	74.14 ± 0.27	120.41 ± 0.21	79.26 ± 0.21
VL-PET _{large}	4.50	66.82 ± 0.24	55.87 ± 0.18	72.00 ± 0.21	120.50 ± 0.21	77.91 ± 0.08
VL-PET _{large} *	7.31	66.95 ± 0.21	57.06 ± 0.21	73.42 ± 0.21	121.06 ± 0.21	79.25 ± 0.21

Relatively Improvement



- VL-PET_{large} with lightweight PET module designs significantly outperforms VL-Adapter by 2.92% (3.41%) and LoRA by 3.37% (7.03%) with BART-base (T5-base) on image-text tasks

Transfer VL-PET Modules to Video-Text Tasks

Method	Trainable Params (%)	TVQA Acc. (%)	How2QA Acc. (%)	TVC Cap. (CIDEr)	YC2C Cap. (CIDEr)	Avg.
Full Fine-tuning	100	77.69	74.79	50.56	151.71	88.69
BitFit	0.38	66.05	65.42	31.16	115.23	69.47
Prompt Tuning	2.00	24.51	27.76	30.22	108.04	47.63
Compacter	1.89	73.78	72.14	41.39	140.52	81.96
LoRA	5.17	75.51	72.69	44.17	142.72	83.77
VL-Adapter	3.39	77.06	74.73	46.72	153.28	87.95
VL-PET _{small}	2.18	77.69	74.89	47.92	150.24	87.69
VL-PET _{medium}	2.17	77.76	75.40	48.30	150.25	87.93
VL-PET _{mediumY}	2.17	77.58	75.15	47.93	151.13	87.95
VL-PET _{large}	3.37	76.97	75.60	47.53	154.41	88.63

Selected Ablation Studies

Method	Params (%)	VQA (%)	GQA (%)	NLVR ² (%)	COCO (CIDEr)	Avg.
VL-PET w/o G	2.97	65.22 ± 0.14	53.35 ± 0.21	72.65 ± 0.44	120.19 ± 0.68	77.85 ± 0.34
VL-PET _{small}	2.98	65.43 ± 0.06	54.03 ± 0.14	72.43 ± 0.21	120.88 ± 0.35	78.14 ± 0.11
VL-PET _{mediumX}	2.98	65.54 ± 0.06	54.33 ± 0.14	72.66 ± 0.19	120.72 ± 0.31	78.37 ± 0.10
VL-PET _{mediumY}	2.98	65.36 ± 0.15	53.83 ± 0.31	73.43 ± 0.79	120.31 ± 0.69	78.23 ± 0.19
VL-PET _{large}	4.16	66.17 ± 0.27	55.11 ± 0.17	73.43 ± 0.35	122.03 ± 0.46	79.18 ± 0.14

Applying VL-PET Designs to Existing Methods

Method	Params (%)	VQA (%)	GQA (%)	NLVR ² (%)	COCO (CIDEr)	Avg.
Compacter	2.70	64.63 ± 0.09	52.70 ± 0.24	71.11 ± 0.35	114.69 ± 0.42	75.78 ± 0.21
+ G _{small} + LW	2.08	64.14 ± 0.05	52.84 ± 0.45	71.04 ± 0.54	77.07 ± 0.31	76.70 ± 0.08
+ G _{medium} + LW	2.07	64.35 ± 0.10	53.10 ± 0.73	70.57 ± 0.41	119.02 ± 0.37	76.75 ± 0.17
+ G _{large} + LW	2.07	64.00 ± 0.12	52.49 ± 0.90	70.81 ± 0.68	117.48 ± 0.11	76.20 ± 0.22
VL-Adapter	4.18	65.76 ± 0.28	54.16 ± 0.41	73.19 ± 0.33	114.61 ± 0.38	76.93 ± 0.05
+ G _{small} + LW	2.98	65.56 ± 0.09	54.34 ± 0.32	71.95 ± 0.23	119.10 ± 0.25	78.00 ± 0.05
+ G _{medium} + LW	2.98	65.73 ± 0.12	54.90 ± 0.10	73.04 ± 0.10	118.34 ± 0.01	78.00 ± 0.01
+ G _{large} + LW	2.98	65.67 ± 0.17	54.11 ± 0.29	73.18 ± 0.34	117.38 ± 0.12	77.58 ± 0.08
+ G _{large} + LW	4.16	66.31 ± 0.23	55.09 ± 0.37	73.46 ± 0.37	119.05 ± 0.53	78.47 ± 0.11

Qualitative Analysis of Granularity Control

