



EUROPEAN CONFERENCE ON COMPUTER VISION

MILANO  
2024

# PYRA: Parallel Yielding Re-Activation for Training-Inference Efficient Task Adaptation

Yizhe Xiong<sup>1,2</sup>, Hui Chen<sup>2\*</sup>, Tianxiang Hao<sup>1,2</sup>, Zijia Lin<sup>1</sup>, Jungong Han<sup>2,3</sup>,  
Yuesong Zhang<sup>4</sup>, Guoxin Wang<sup>4</sup>, Yongjun Bao<sup>4</sup>, Guiguang Ding<sup>1,2</sup>



1



2

北京信息科学与技术  
国家研究中心

BEIJING NATIONAL RESEARCH CENTER  
FOR INFORMATION SCIENCE AND TECHNOLOGY



3

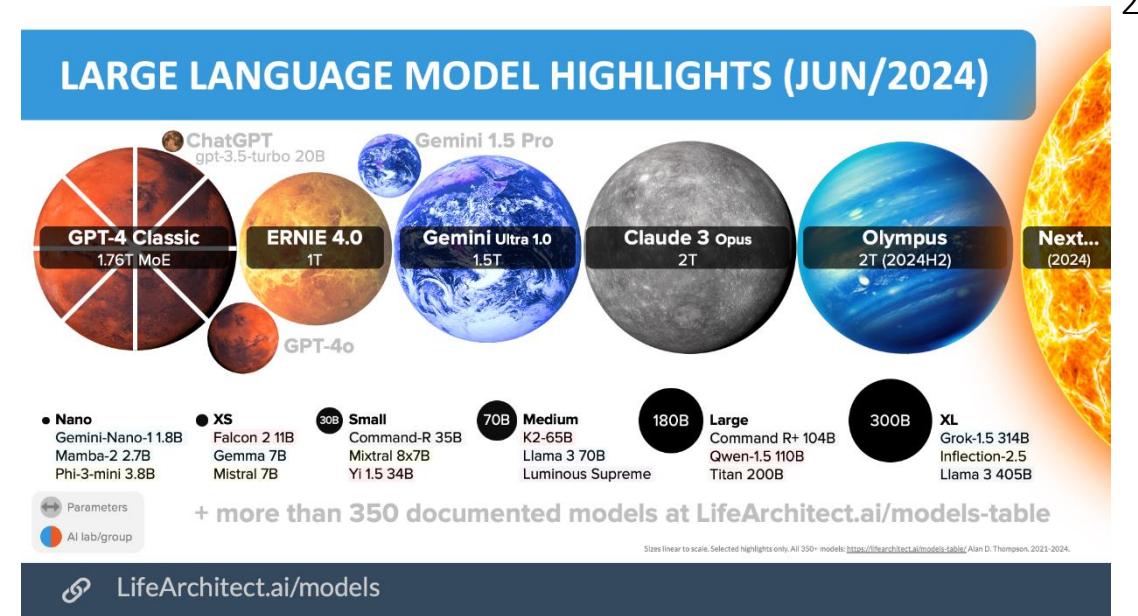
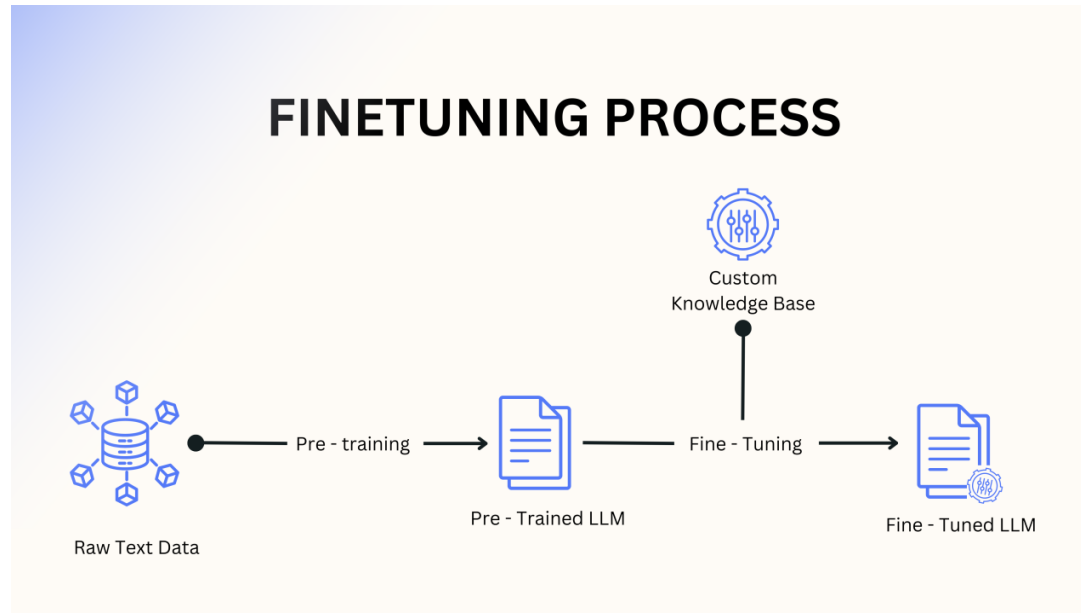


JD.COM

4

# Application of Foundation Models

- Foundation Models (LLaMA, ViT, CLIP, etc.) are fine-tuned on downstream task data for real-world applications.
- The size of Foundation models are rapidly growing.



Credits:

1: <https://www.linkedin.com/pulse/what-supervised-fine-tuning-tagx-yq8if/>

2: <https://lifearchitect.ai/models/>

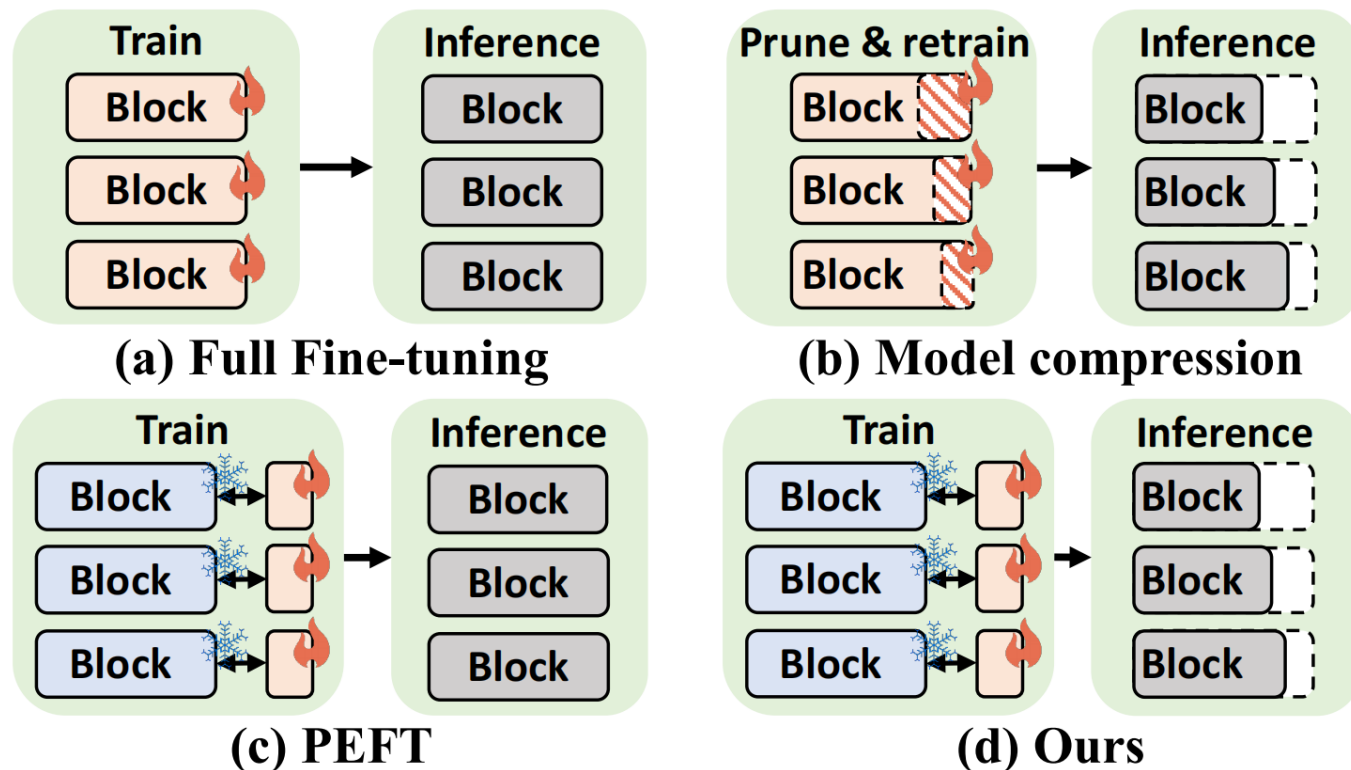
# Application of Foundation Models

---

- Primary obstacles in implementing Foundation Models on downstream tasks:
  - Training overhead when fine-tuning on downstream tasks
  - Inference efficiency after model deployment
- The above two topics are investigated separately:
  - Parameter-Efficient Fine-Tuning (PEFT): train small amount of parameters during fine-tuning
  - Model Compression: model pruning, knowledge distillation, model quantization, etc.

# Motivation: Efficiency in Training & Inference?

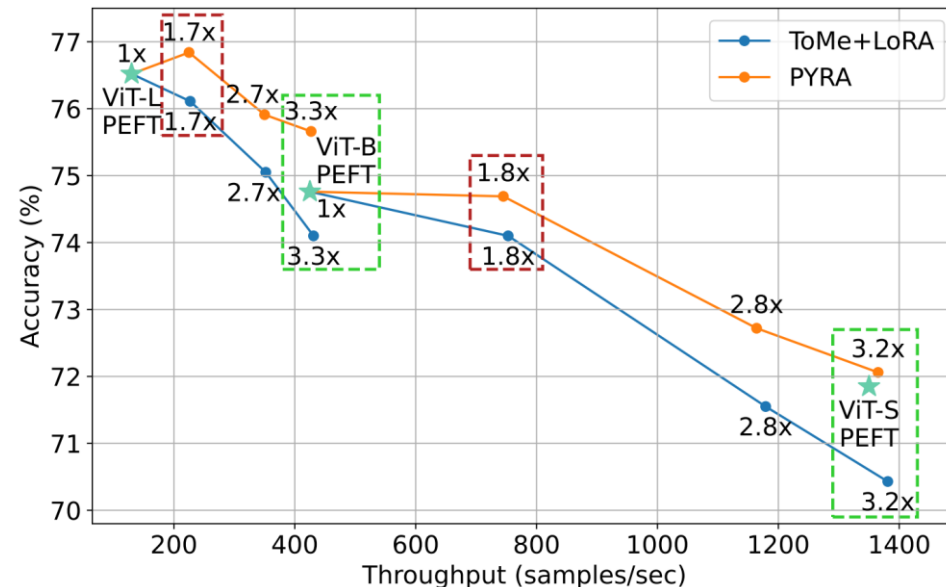
- However, neither PEFT or Model Compression solves both problems!
- Can we achieve both training efficiency and inference efficiency simultaneously?
- We formalize this challenge as **Training-Inference Efficient Task Adaptation**.



Method	(a)	(b)	(c)	(d)
Non-vast Data Scale	✓	✗	✓	✓
Training Efficient	✗	✗	✓	✓
Inference Efficient	✗	✓	✗	✓

# Problems of the Baseline approach

- Baseline solution: a mergeable PEFT method + a training-free compression method
  - Here we choose LoRA<sup>1</sup> + Token merging<sup>2</sup> (a strong baseline according to our experiments)
- Quick evaluation of the baseline approach:
  - **Un-ignorable performance drops** under low compression rates (1.7 to 1.8 times speedup)
  - **“Adverse Compression”**: Under-performs “smaller backbone with similar throughput + PEFT” under high compression rates (over 3 times speedup)



Credits:

1: LoRA: Low-Rank Adaptation of Large Language Models (ICLR 2022)

2: Token Merging: Your ViT But Faster (ICLR 2023)

# PYRA: Parallel Yielding Re-Activation

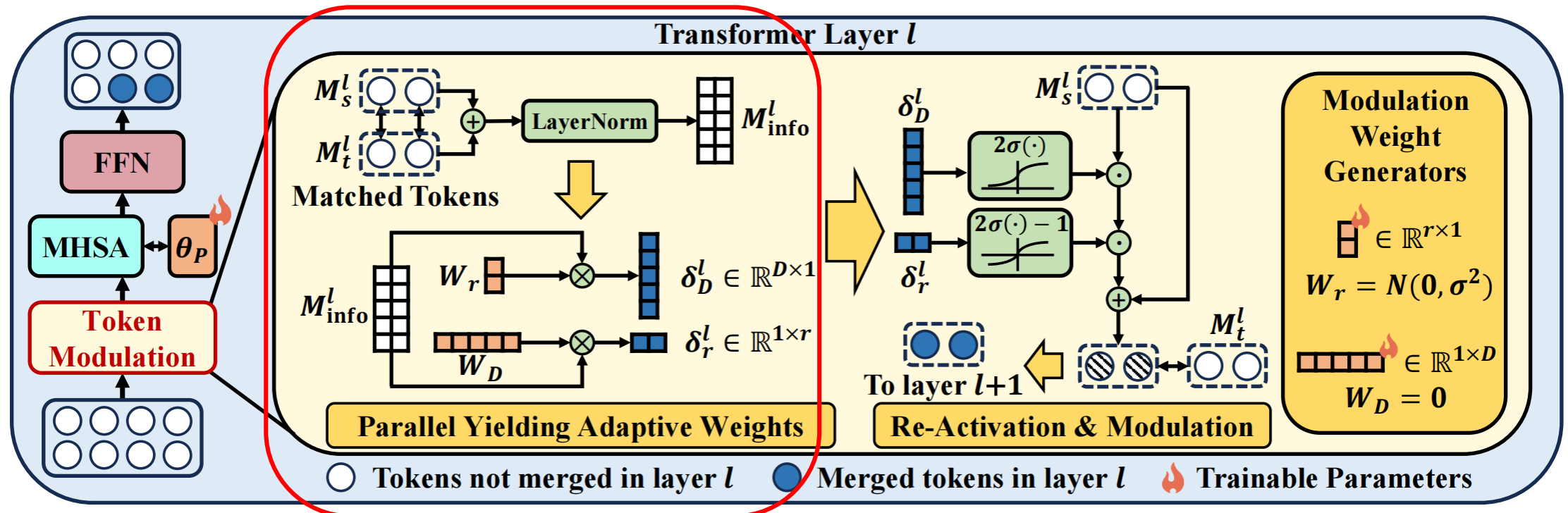
- Parallel Yielding Adaptive Weights:

- For the merging token pairs, we group the source  $M_s^l$ , the target  $M_t^l$ , then creating

$$M_{\text{info}}^l = \text{LayerNorm}(M_s^l + M_t^l) \in \mathbb{R}^{D \times r}$$

- Then, we use learnable Modulation Weight Generators to create adaptive weights:

$$\delta_D^l = M_{\text{info}}^l W_r^l \in \mathbb{R}^{D \times 1} \quad \delta_r^l = W_D^l M_{\text{info}}^l \in \mathbb{R}^{1 \times r}$$

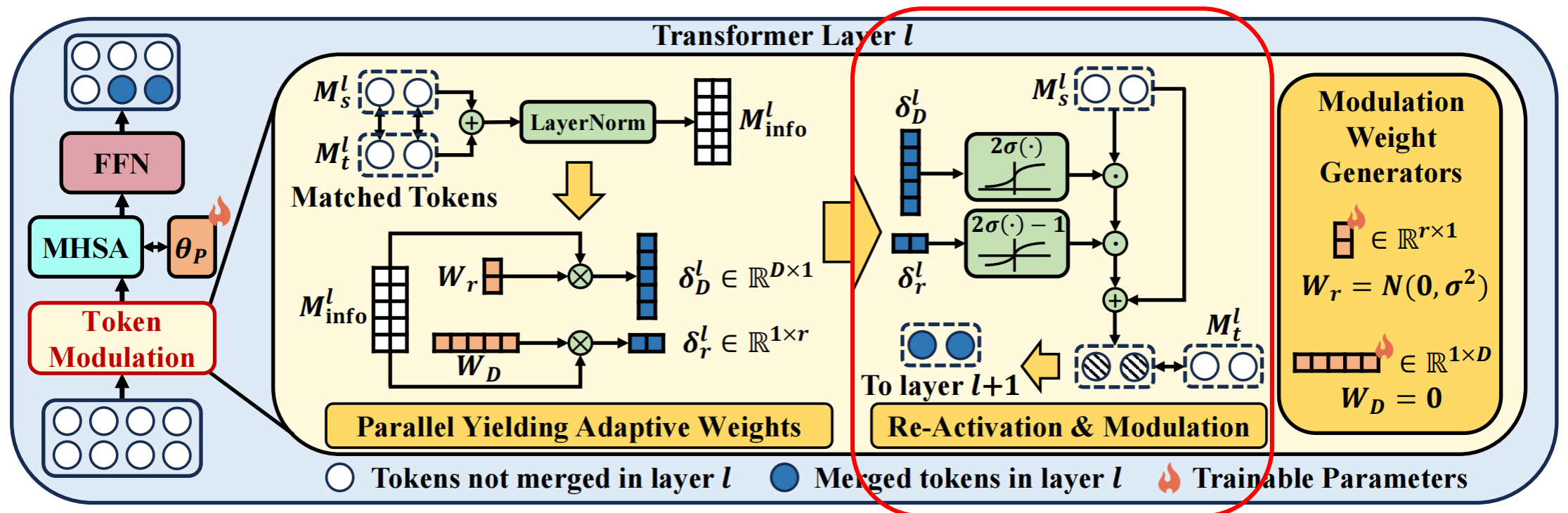


# PYRA: Parallel Yielding Re-Activation

- Re-Activation & Modulation:
  - To regularize the adaptive weights and create non-linearity for expression capacity, we use a re-activation before token modulation. We add a residual connection for better gradients.

$$\widehat{M}_S^l = 2\sigma(\widehat{\delta}_D^l) \odot M_S^l$$

$$M_S^l \leftarrow M_S^l + (2\sigma(\widehat{\delta}_r^l) - 1) \odot \widehat{M}_S^l$$





# PYRA: Parallel Yielding Re-Activation

- PYRA is extremely parameter-efficient and compute-efficient:
  - Low compression rate: **50% FLOPs decrease, less than 0.01% trainable parameters** (Table 1)
  - High compression rate: **>75% FLOPs decrease, less than 0.01% trainable parameters**

**Table 1:** The complexity comparisons between conducting PEFT with and without PYRA. The FLOPs metric is obtained during inference.

Model	Metric	Total	PEFT w/o PYRA (%)	PEFT w. PYRA (%)
ViT-Base	# params	86M	0.29M	0.34%
	FLOPs	16.37G	16.37G	100%
ViT-Large	# params	303M	1.18M	0.39%
	FLOPs	57.37G	57.37G	100%



# Experiment: Overall Comparison

- Low compression rate: we choose LoRA as PEFT for all methods
  - Outperforms the best baseline: ToMe + LoRA
  - Achieve comparable performance or even out-performs backbone + LoRA

Benchmark: VTAB-1k

Method	# params	Throughput	Natural							Specialized				Structured							Average	
			Cifar100	Caltech101	DTD	Flowers102	Pets	SVHN	Sun397	Camelyon	EuroSAT	Resisc45	Retinopathy	Clevr-Count	Clevr-Dist	DMLAB	KITTI-Dist	dSpr-Loc	dSpr-Ori	sNORB-Azim		sNORB-Ele
<b>Model: ViT-B/16 (Throughput: 425)</b>																						
PEFT	0.34%	425	67.1	90.2	69.4	99.1	90.5	85.7	54.1	83.1	95.8	84.3	74.6	82.2	69.2	50.1	79.2	81.8	47.1	31.1	42.6	74.76
RaP	3.43%	654	25.9	68.4	53.3	64.0	57.4	71.3	21.5	75.8	87.9	59.3	73.6	43.1	53.8	26.3	60.5	73.5	25.5	16.7	27.9	55.57
SPViT	4.46%	567	41.6	75.4	61.1	83.2	66.2	56.1	28.3	79.3	94.2	73.3	73.6	70.6	61.5	42.4	67.8	75.4	<b>50.5</b>	28.9	31.3	64.16
DiffRate	0.35%	709	37.1	84.6	63.7	96.7	86.2	32.6	48.2	78.9	85.8	67.0	73.7	32.9	29.8	34.1	55.7	12.6	16.0	13.1	21.5	55.82
ToMe	0.34%	753	<u>64.6</u>	<b>90.4</b>	<u>67.9</u>	<u>98.5</u>	<u>89.8</u>	<u>83.9</u>	<b>53.2</b>	<u>82.6</u>	<u>94.7</u>	<b>83.5</b>	<u>74.9</u>	<u>81.9</u>	<b>69.8</b>	<u>49.2</u>	<u>76.9</u>	<b>81.9</b>	<u>46.5</u>	<u>31.0</u>	<b>43.1</b>	<u>74.10</u>
PYRA	0.35%	745	<b>67.5</b>	<u>90.3</u>	<b>69.3</b>	<b>98.9</b>	<b>90.0</b>	<b>84.6</b>	<u>53.1</u>	<b>83.3</b>	<b>95.7</b>	<u>83.3</u>	<b>75.2</b>	<b>82.6</b>	<u>68.9</u>	<b>50.8</b>	<b>80.0</b>	<u>81.8</u>	45.8	<b>32.2</b>	<u>42.8</u>	<b>74.69</b>
<b>Model: ViT-L/16 (Throughput: 130)</b>																						
PEFT	0.39%	130	77.1	91.4	73.4	99.5	91.3	89.6	57.6	85.9	96.1	87.3	76.1	83.1	63.0	50.7	82.1	81.7	53.5	32.2	36.6	76.52
RaP	1.95%	196	43.2	87.9	62.6	52.8	81.7	86.7	34.7	78.4	92.4	73.3	73.6	68.0	59.6	46.9	<u>82.4</u>	75.5	43.6	24.5	25.7	65.64
SPViT	2.47%	188	48.1	87.5	65.2	94.4	77.4	80.9	38.8	79.9	93.9	79.8	74.3	78.2	<b>65.8</b>	47.4	74.1	<u>82.3</u>	50.3	31.0	37.9	70.22
DiffRate	0.39%	221	50.9	86.8	70.3	97.8	88.3	39.0	52.3	80.2	87.2	72.2	74.2	32.6	32.3	36.5	57.4	22.8	26.6	15.2	23.4	59.53
ToMe	0.39%	227	<u>76.1</u>	<u>91.1</u>	<u>72.3</u>	<u>99.2</u>	<b>91.7</b>	<u>89.2</u>	<u>56.4</u>	<u>86.4</u>	<u>95.1</u>	<u>86.6</u>	<u>75.1</u>	<u>82.4</u>	61.9	<u>50.9</u>	81.4	81.6	<b>53.5</b>	<u>33.4</u>	<u>36.8</u>	<u>76.11</u>
PYRA	0.40%	225	<b>76.6</b>	<b>91.3</b>	<b>73.2</b>	<b>99.3</b>	<u>91.5</u>	<b>89.4</b>	<b>57.1</b>	<b>86.9</b>	<b>95.9</b>	<b>87.1</b>	<b>76.2</b>	<b>83.2</b>	<u>63.2</u>	<b>52.8</b>	<b>83.1</b>	<b>82.5</b>	<u>52.6</u>	<b>34.8</b>	<b>39.0</b>	<b>76.84</b>

# Experiment: Overall Comparison

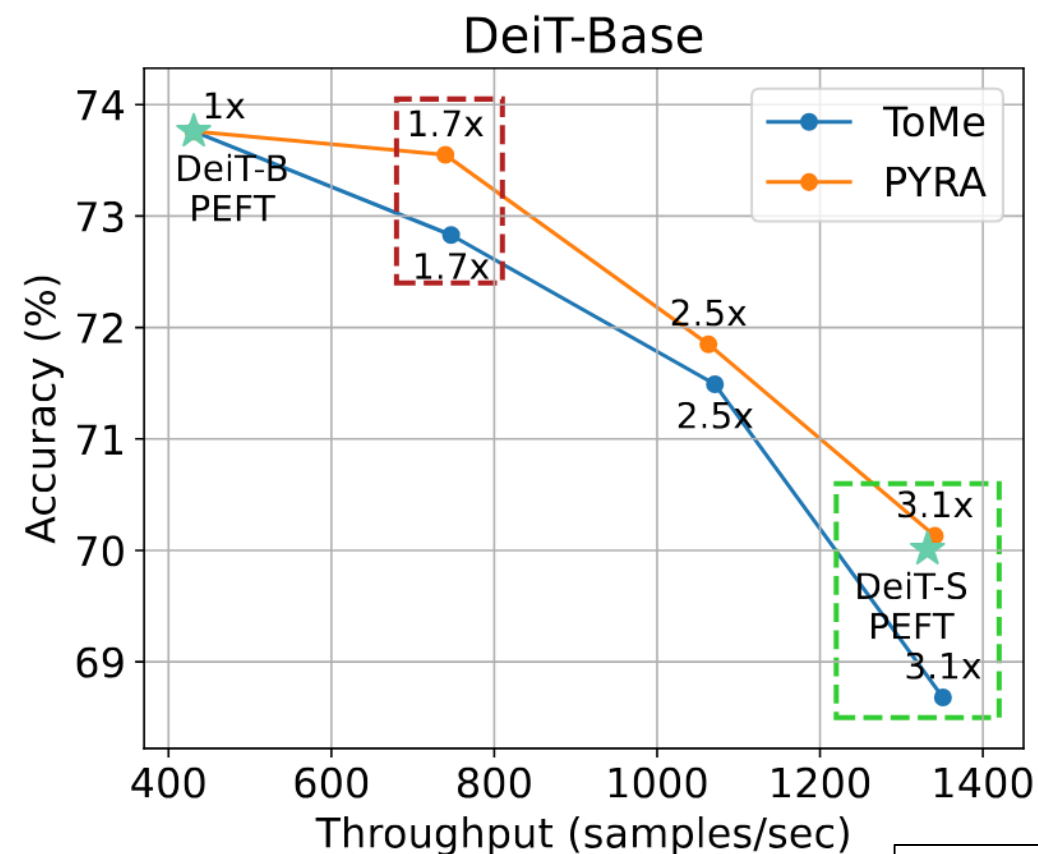
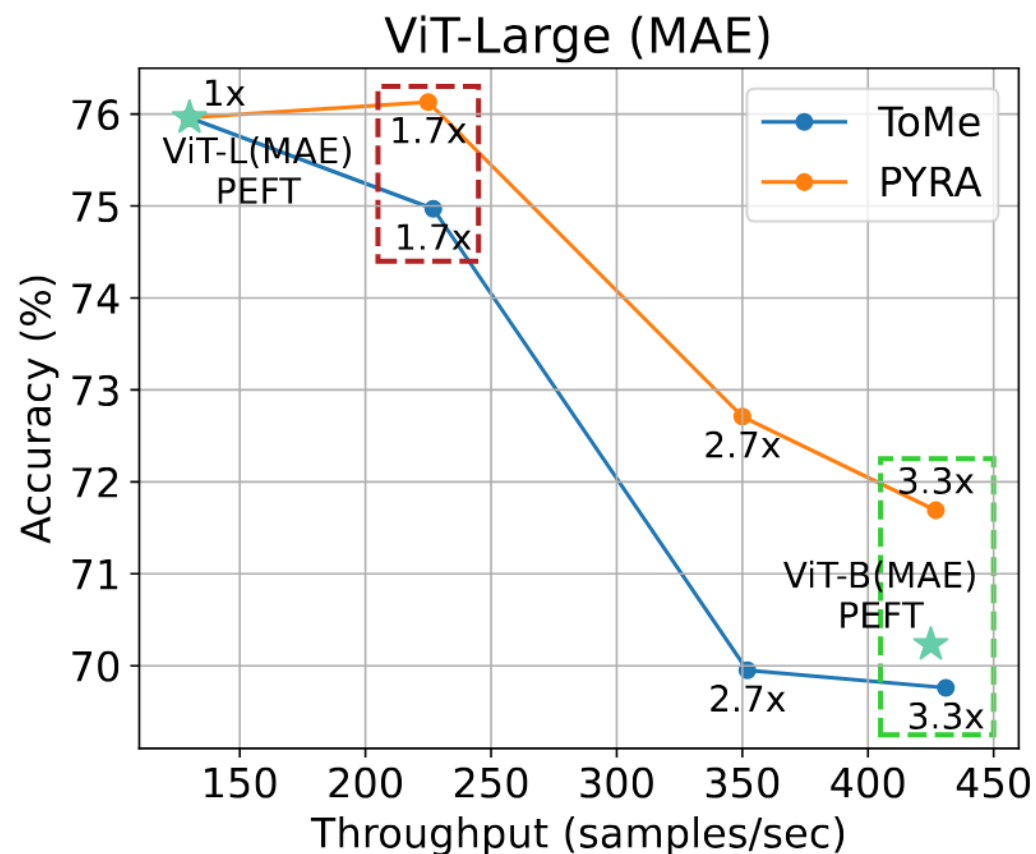
- High compression rate: we choose LoRA as PEFT for all methods
  - Outperforms the best baseline: ToMe + LoRA
  - Solves the “Adverse Compression” issue: outperforms smaller backbone + LoRA

Benchmark: VTAB-1k

Method	# params	Throughput	Natural							Specialized				Structured							Average	
			Cifar100	Caltech101	DTD	Flowers102	Pets	SVHN	Sun397	Camelyon	EuroSAT	Resisc45	Retinopathy	Clevr-Count	Clevr-Dist	DMLAB	KITTI-Dist	dSpr-Loc	dSpr-Ori	sNORB-Azim		sNORB-Ele
<b>Model: ViT-B/16 (Throughput: 425)</b>																						
PEFT*	0.34%	1350	57.7	88.2	70.1	98.7	88.7	85.7	44.9	81.4	94.7	84.6	73.6	81.6	64.1	48.1	80.0	72.9	38.4	22.9	37.7	71.85
RaP	0.86%	1029	24.3	40.1	34.5	41.8	40.5	21.7	11.4	75.8	86.5	35.1	73.8	49.6	49.7	28.1	39.4	13.8	15.4	12.4	26.9	42.60
SPViT	4.46%	944	23.7	67.9	51.9	69.9	53.2	19.6	13.1	71.9	81.3	67.9	74.7	53.5	61.9	39.5	57.4	45.0	34.5	11.1	23.2	52.49
DiffRate	0.35%	1308	23.2	73.0	55.7	87.9	66.7	27.2	29.3	78.1	77.8	53.1	73.6	29.7	28.6	31.7	52.6	11.5	17.2	11.3	20.3	49.29
ToMe	0.34%	1381	<b>54.2</b>	<b>87.8</b>	<b>65.5</b>	<b>96.1</b>	<b>81.7</b>	<b>79.7</b>	<b>45.2</b>	<b>79.4</b>	<b>93.6</b>	<b>76.3</b>	73.8	<b>78.3</b>	<b>65.7</b>	<b>48.0</b>	<b>71.3</b>	<b>80.0</b>	<b>45.8</b>	<b>30.9</b>	<b>41.2</b>	<b>70.43</b>
PYRA	0.35%	1365	<b>54.0</b>	<b>89.3</b>	<b>67.1</b>	<b>96.5</b>	<b>84.0</b>	<b>81.8</b>	<b>44.6</b>	<b>81.2</b>	<b>94.6</b>	<b>79.5</b>	<b>75.1</b>	<b>79.9</b>	<b>67.0</b>	<b>49.2</b>	<b>76.9</b>	<b>82.6</b>	<b>47.8</b>	<b>31.9</b>	<b>42.0</b>	<b>72.06</b>
<b>Model: ViT-L/16 (Throughput: 130)</b>																						
PEFT*	0.34%	425	67.1	90.2	69.4	99.1	90.5	85.7	54.1	83.1	95.8	84.3	74.6	82.2	69.2	50.1	79.2	81.8	47.1	31.1	42.6	74.76
RaP	0.65%	301	17.7	37.1	27.0	46.2	33.3	23.2	13.3	76.5	74.2	54.4	73.6	50.4	31.4	25.7	49.8	53.1	25.5	13.4	26.0	44.11
SPViT	2.47%	289	54.0	87.6	65.5	94.8	74.9	32.6	38.6	81.8	<b>95.3</b>	78.0	74.0	72.8	<b>61.2</b>	46.9	70.2	77.1	47.4	31.3	28.6	66.90
DiffRate	0.39%	416	47.4	73.5	54.1	84.3	60.2	19.6	22.2	50.0	64.6	42.8	18.2	31.5	31.9	31.1	37.3	22.0	17.7	14.8	21.4	40.48
ToMe	0.39%	431	<b>71.0</b>	<b>90.9</b>	<b>70.4</b>	<b>98.3</b>	<b>88.5</b>	<b>87.2</b>	<b>52.4</b>	<b>82.9</b>	<b>94.5</b>	<b>83.1</b>	<b>75.0</b>	<b>80.7</b>	61.1	<b>48.9</b>	<b>76.9</b>	<b>80.8</b>	<b>53.0</b>	<b>32.1</b>	<b>35.2</b>	<b>74.10</b>
PYRA	0.40%	427	<b>71.6</b>	<b>91.8</b>	<b>71.1</b>	<b>98.5</b>	<b>89.7</b>	<b>88.1</b>	<b>52.2</b>	<b>85.1</b>	<b>95.3</b>	<b>84.6</b>	<b>75.7</b>	<b>80.9</b>	<b>63.0</b>	<b>51.7</b>	<b>82.0</b>	<b>82.0</b>	<b>54.2</b>	<b>36.0</b>	<b>41.2</b>	<b>75.66</b>

# Experiment: on more backbones

- On self-supervised backbone (ViT-L MAE) and distillation backbone (DeiT-B), the issues in low/high compression rates are resolved as well.
- For more analysis experiments, please refer to the article and the appendix.



# Summary

---

- We define and introduce Training-Inference Efficient Task Adaptation, in which the inference efficiency and training efficiency are considered for applying foundation models on downstream tasks.
- We investigate the proposed challenge, and discovered that the best performing baselines exhibit issues under both low and high compression rates.
- We propose Parallel Yielding Re-Activation (PYRA), a light-weight method that adaptively modulates tokens in Vision Transformers.
- Extensive experiments show that PYRA consistently resolves the issues in both low and high compression rates on various Vision Transformer backbones, for the first time making Training-Inference Efficient Task Adaptation actually applicable.

Contact me at: [xiongyizhe2001@163.com](mailto:xiongyizhe2001@163.com)

My homepage: [xiongyizhe.xyz](http://xiongyizhe.xyz)