

# Vision Transformer Finetuning Benefits from Non-Smooth Components

Ambroise Odonnat

Noah's Ark Lab, Inria

Mila, CERC-AAI Lab

March 27, 2026



*Inria*

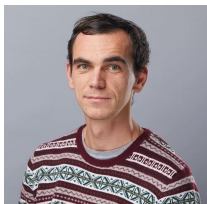




Laetitia Chapel  
L'Institut Agro, IRISA



Romain Tavenard  
Univ. Rennes 2, IRISA



Ievgen Redko  
Noah's Ark Lab



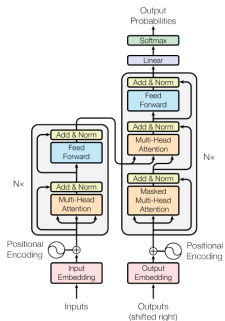
- ① Introduction
- ② ViT plasticity
- ③ Plasticity analysis
- ④ Finetuning benefits
- ⑤ Take home message



- 1 Introduction
- 2 ViT plasticity
- 3 Plasticity analysis
- 4 Finetuning benefits
- 5 Take home message



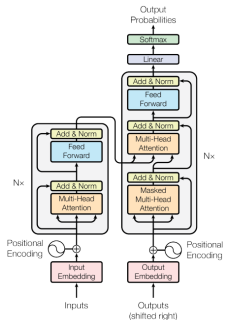
Most foundation models are based on the transformer architecture.



(Vaswani et al., 2017)



Most foundation models are based on the transformer architecture.



(Vaswani et al., 2017)



- Diverse training data
- Large scale pretraining
- (Very) large models



- × Need of specialized models in practice
- × Increasing cost of adaptation with models growing larger



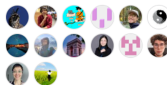
## State-of-the-art Parameter-Efficient Fine-Tuning (PEFT) methods

Fine-tuning large pretrained models is often prohibitively costly due to their scale. Parameter-Efficient Fine-Tuning (PEFT) methods enable efficient adaptation of large pretrained models to various downstream applications by only fine-tuning a small number of (extra) model parameters instead of all the model's parameters. This significantly decreases the computational and storage costs. Recent state-of-the-art PEFT techniques achieve performance comparable to fully fine-tuned models.

Used by 35.2k



Contributors 301



## Plethora of parameter-efficient finetuning methods



- × Need of specialized models in practice
- × Increasing cost of adaptation with models growing larger



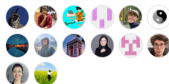
## State-of-the-art Parameter-Efficient Fine-Tuning (PEFT) methods

Fine-tuning large pretrained models is often prohibitively costly due to their scale. Parameter-Efficient Fine-Tuning (PEFT) methods enable efficient adaptation of large pretrained models to various downstream applications by only fine-tuning a small number of (extra) model parameters instead of all the model's parameters. This significantly decreases the computational and storage costs. Recent state-of-the-art PEFT techniques achieve performance comparable to fully fine-tuned models.

Used by 35.2k



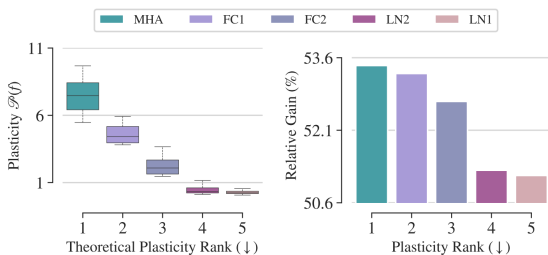
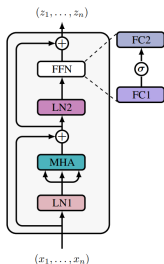
Contributors 301



*Which transformer components should be prioritized during finetuning and why?*



Finetuning high-plasticity components (non-smooth) yields larger performance gains than the stiff ones.

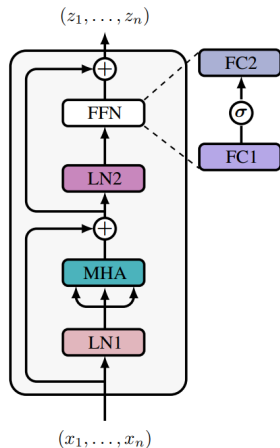




- ① Introduction
- ② ViT plasticity
- ③ Plasticity analysis
- ④ Finetuning benefits
- ⑤ Take home message



- Inputs are 2D images
- Images split into patches
- Tokens are flattened patches
- Encoder-only architecture





**Smoothness is connected to differentiability, robustness to input perturbations and Lipschitz continuity.**

Definition (Lipschitz continuity)

The smallest constant  $K$  such that

$$\forall x, y \in (\mathbb{R}^d)^n, \|f(x) - f(y)\|_F \leq K \|x - y\|_F$$

is called the Lipschitz constant of  $f$  and writes

$$\text{Lip}(f) = \sup_{x \neq y \in (\mathbb{R}^d)^n} \frac{\|f(x) - f(y)\|_F}{\|x - y\|_F}.$$



**Promoting smoothness by regularizing the Lipschitz constant reduces the sensitivity to input perturbations.**

- ✓ Good for training stability and adversarial robustness,
  - ✗ Too much smoothness constrain adaptability to new tasks,
- Identify components with too small Lipschitz constants.



## Definition (Plasticity)

Let  $\nu$  be the uniform distribution over the set of distinct pairs of sequences of tokens in  $(\mathbb{R}^d)^n$ . We define the plasticity of a transformer component  $f$  as

$$\mathcal{P}(f) = \mathbb{E}_{(x,y) \sim \nu} \left[ \frac{\|f(x) - f(y)\|_F}{\|x - y\|_F} \right].$$

- ✓ For any component  $f$ , we have  $\mathcal{P}(f) \leq \text{Lip}(f)$ ,
- ✓ If  $\mathcal{P}(f) < 1$ , the module  $f$  contracts the changes in inputs,
- ✓ If  $\mathcal{P}(f) > 1$ ,  $f$  amplifies the input perturbations
- ✓ When  $\mathcal{P}(f) > 1$ ,  $\text{Lip}(f)$  is pushed from below.



- ✓ High plasticity  $\rightarrow$  high sensitivity to input changes,
- ✓ High plasticity  $\rightarrow$  high Lipschitz constant  $\rightarrow$  low smoothness,
- ✓ High plasticity  $\rightarrow$  faster and better adaptation.

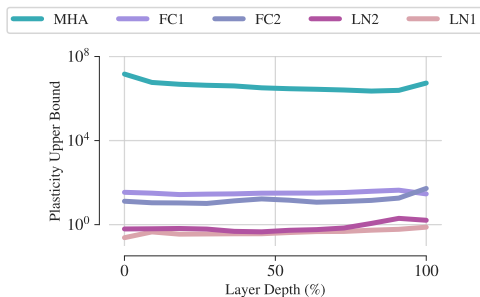




- ① Introduction
- ② ViT plasticity
- ③ Plasticity analysis
- ④ Finetuning benefits
- ⑤ Take home message



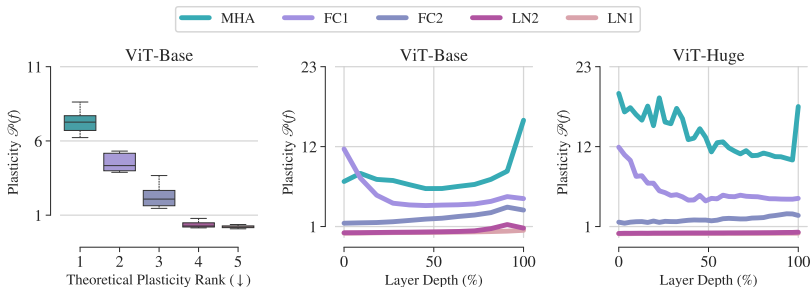
We derive upper bounds on the plasticity  $\mathcal{P}(f)$ .



**Takeaway 1.** Our theoretical analysis suggests the following plasticity ranking:  $\text{MHA} \rightarrow \text{FC1} \approx \text{FC2} \rightarrow \text{LN2} \approx \text{LN1}$ .



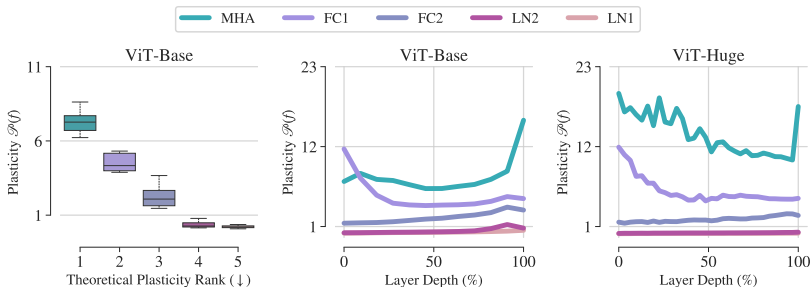
- ViT pretrained on ImageNet (86M and 632M of parameters)
- Cifar10/100, Cifar10-C, DomainNet, Pet, Flowers102



- ✓ Our theory correctly predicts the empirical ranking,
- ✓ Attention and feedforward linear layers have  $\mathcal{P}(f) > 1$ ,
- ✓ LayerNorms have  $\mathcal{P}(f) < 1$ .



- ViT pretrained on ImageNet (86M and 632M of parameters)
- Cifar10/100, Cifar10-C, DomainNet, Pet, Flowers102



**Takeaway 2.** The empirical *plasticity ranking* of vision transformer modules supports our theoretical insights with:  
MHA  $\rightarrow$  FC1  $\rightarrow$  FC2  $\rightarrow$  LN2  $\rightarrow$  LN1.

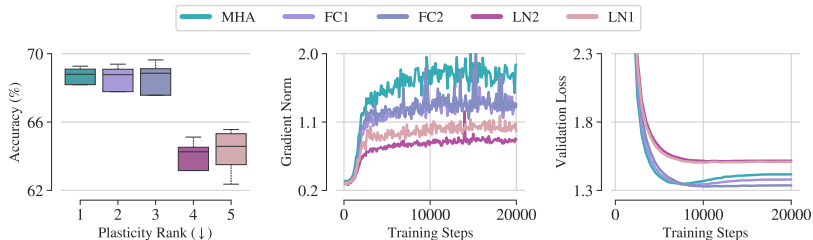


- ① Introduction
- ② ViT plasticity
- ③ Plasticity analysis
- ④ Finetuning benefits**
- ⑤ Take home message



*Table 1. MHA leads to better performance (11 benchmarks).* We report the decrease in performance between MHA and the other configurations, averaged over the benchmarks. Entries in **bold** indicate that the decrease is statistically significant according to a Wilcoxon signed-rank test at a confidence level of 5%.

configuration	FC1	FC2	LN1	LN2
<i>decrease (%)</i>	0.07	<b>0.44</b>	<b>0.8</b>	<b>0.9</b>



## Takeaway 3.

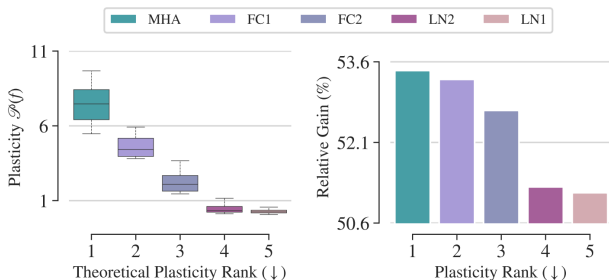
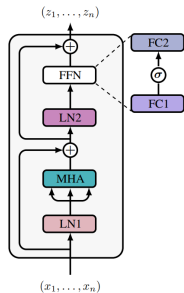
- ✓ High plasticity → better and more stable finetuning,
- ✓ Prioritize attention module and feedforward linear layers.



- ① Introduction
- ② ViT plasticity
- ③ Plasticity analysis
- ④ Finetuning benefits
- ⑤ Take home message



Finetuning high-plasticity components (non-smooth) yields larger performance gains than the stiff ones.

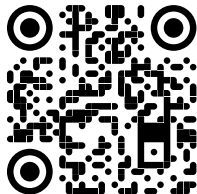




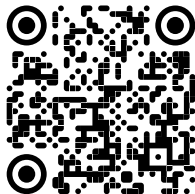
- ✓ Study the effect of a tailored optimization,
- ✓ Extension to large language models,
- ✓ Apply LoRa only to components with high-plasticity.



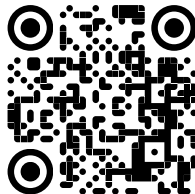
Check out the paper along with the code to know more!



paper



code

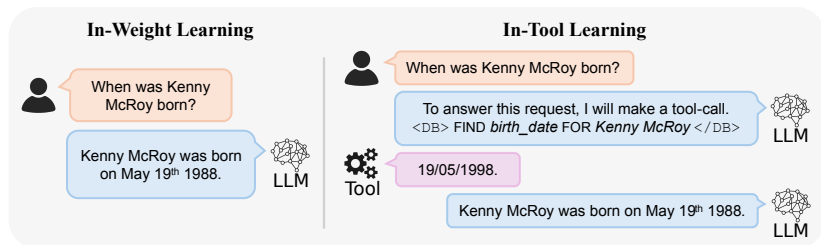


website



## Provable benefits of in-tool learning for LLMs

ICLR 2026 MemAgents Workshop



[paper](#) - [code](#)

Thanks for your attention!