

ReFT: Representation Finetuning for Language Models

Contents Covered Here

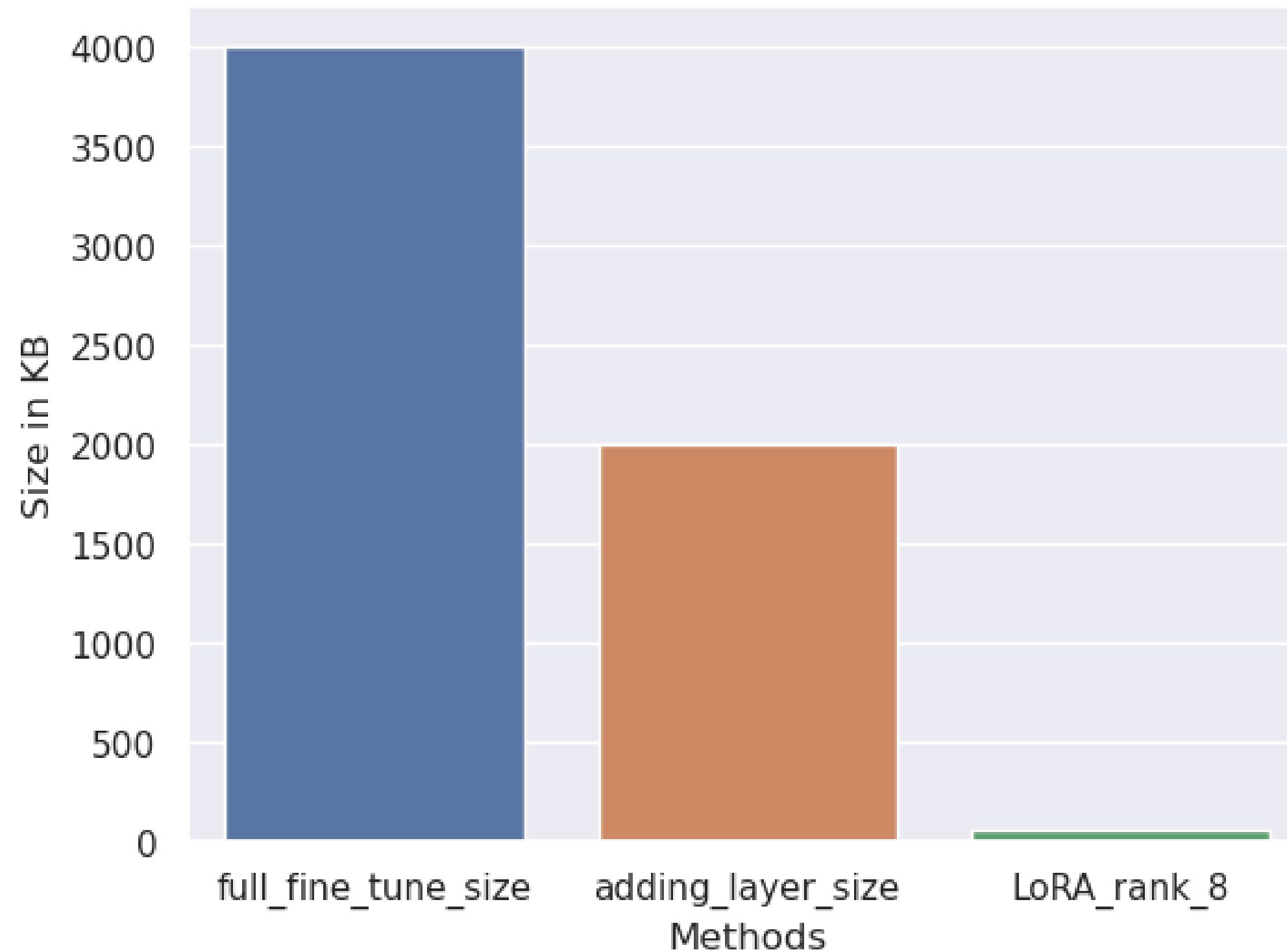
- **What is PEFT? and Why PEFT is needed in the first place**
- Types of PEFT methods
 - Adding extra layer sequentially (general finetuning)
 - Finetuning with Minimal Layers
 - LoRA
 - DoRA
- Quantization methods
- ReFT methodology
- Conclusion and Future enhancements from here (personal note)

What is PEFT? and Why PEFT is needed in the first place

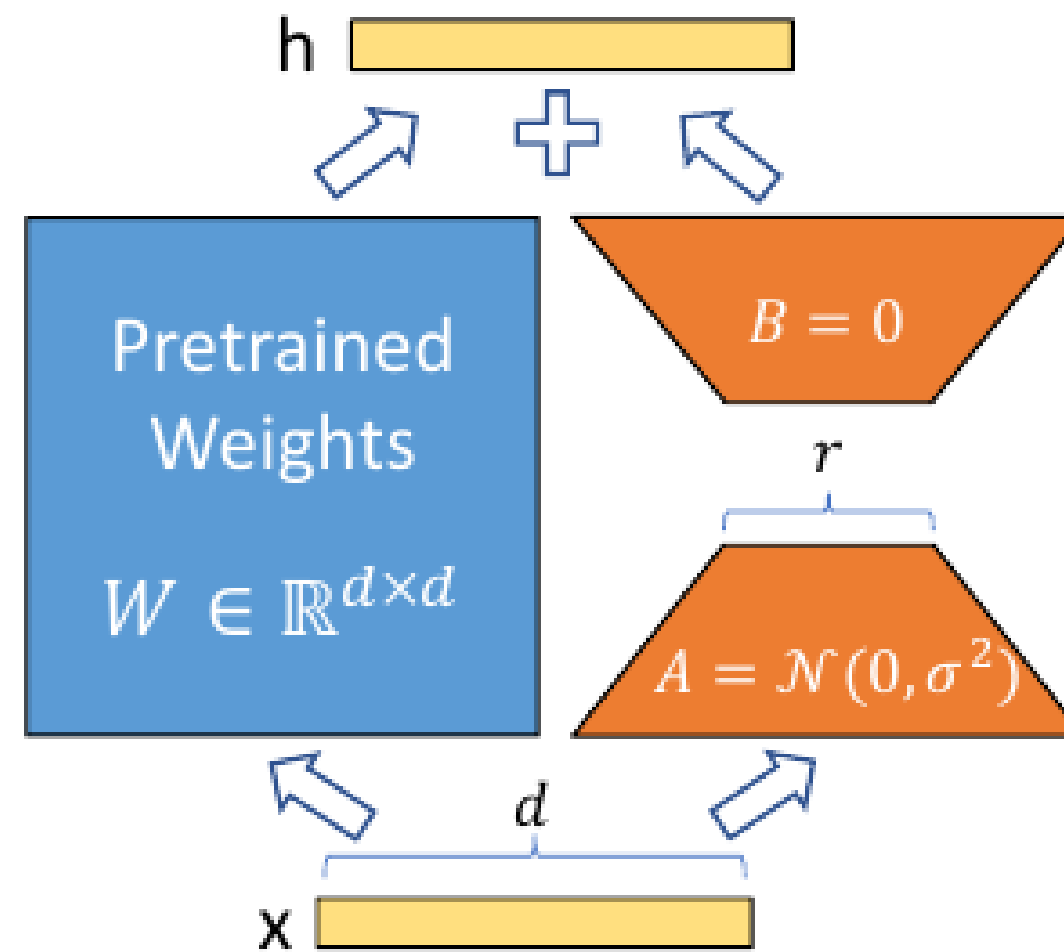
- PEFT stands for Parametric Efficient Fine Tuning which involves in training a Parameter heavy models (such as LLM) with minimal parameters as possible without loosing the originality of the model.
- As LLMs are getting heavier and heavier it is nearly impossible to finetune the entire model like conventional DL models (considering adequate amount of data is present) on one GPU
- So to fine tune the model on one bigger GPU or on multiple smaller GPUs we need PEFT

Practical example for this

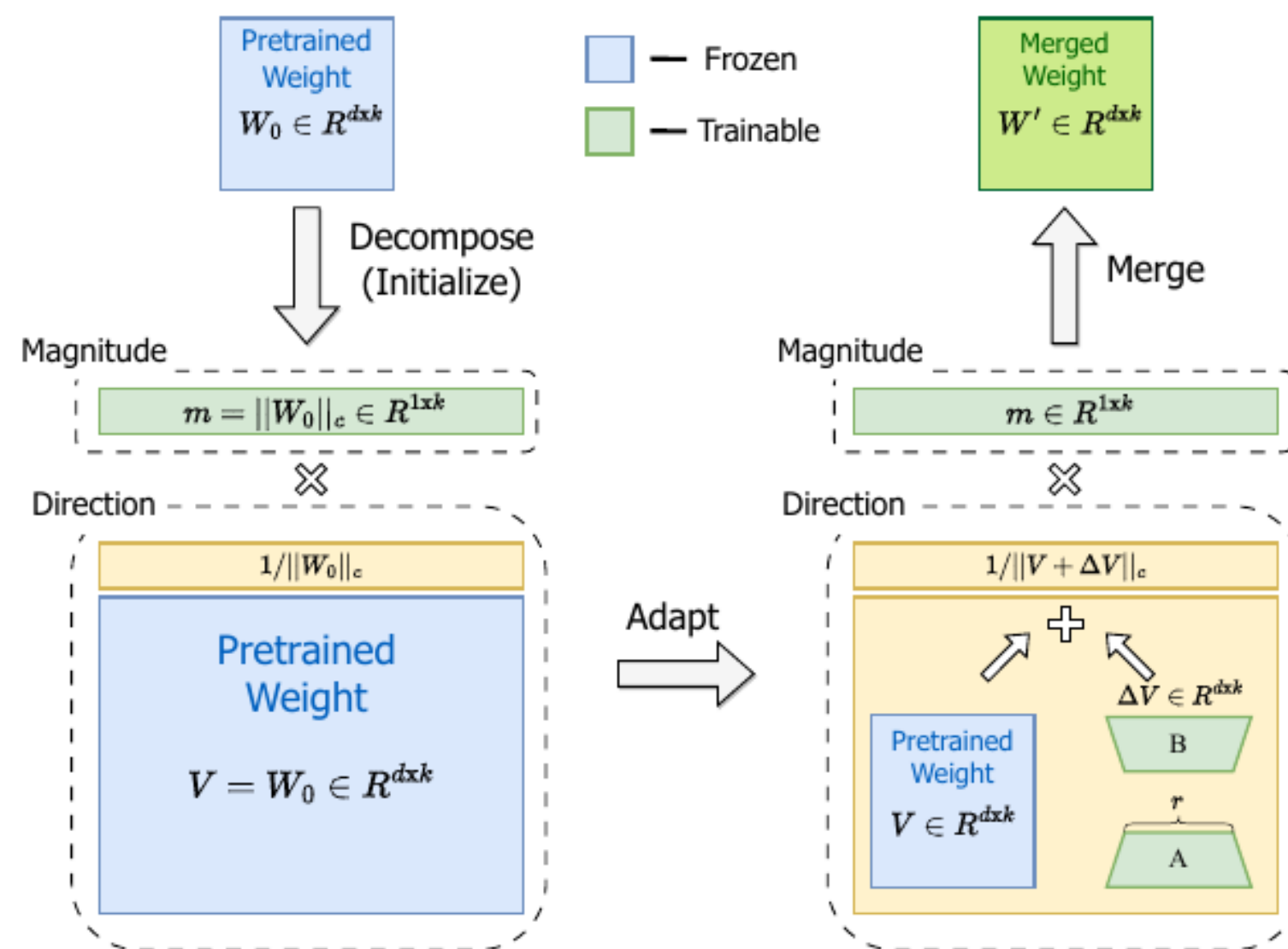
Consider a Matrix of size 1000 x 1024 x 1024 with fp32 precision



LoRA



DoRA



DoRA - Inference

- In this you scale the Weight matrix into Unit vector and try making the magnitude of the matrix m as a learnable parameter

Advantages

- Can change its slope towards the negative directions and align with representations learnt from fully finetuning the model

$$W = m \frac{V}{\|V\|_c} = \|W\|_c \frac{W}{\|W\|_c}$$

DoRA Inference

- We make the Magnitude learnable and $V' = V + \Delta V$

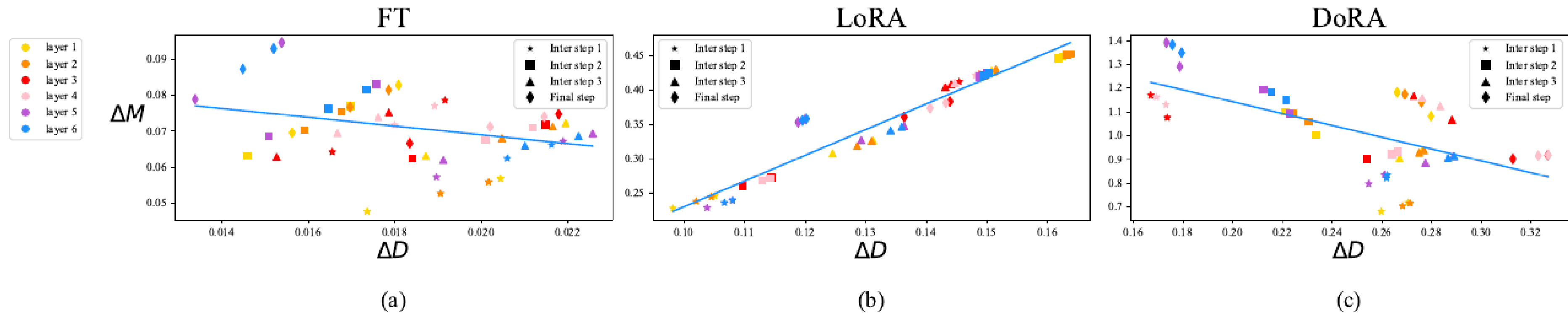
$$W' = \underline{m} \frac{V + \Delta V}{\|V + \Delta V\|_c} = \underline{m} \frac{W_0 + \underline{BA}}{\|W_0 + \underline{BA}\|_c}$$

$$\begin{aligned} \nabla_{V'} \mathcal{L} &= \frac{m}{\|V'\|_c} \left(I - \frac{V' V'^T}{\|V'\|_c^2} \right) \nabla_{W'} \mathcal{L} \\ \nabla_m \mathcal{L} &= \frac{\nabla_{W'} \mathcal{L} \cdot V'}{\|V'\|_c} \end{aligned} \quad d\|\mathbf{x}\| = \frac{\mathbf{x}^T d\mathbf{x}}{\|\mathbf{x}\|}$$

Why DoRA With Stop Gradients

$$\nabla_{V'} \mathcal{L} = \frac{m}{C} \nabla_{W'} \mathcal{L} \text{ where } C = \|V'\|_c$$

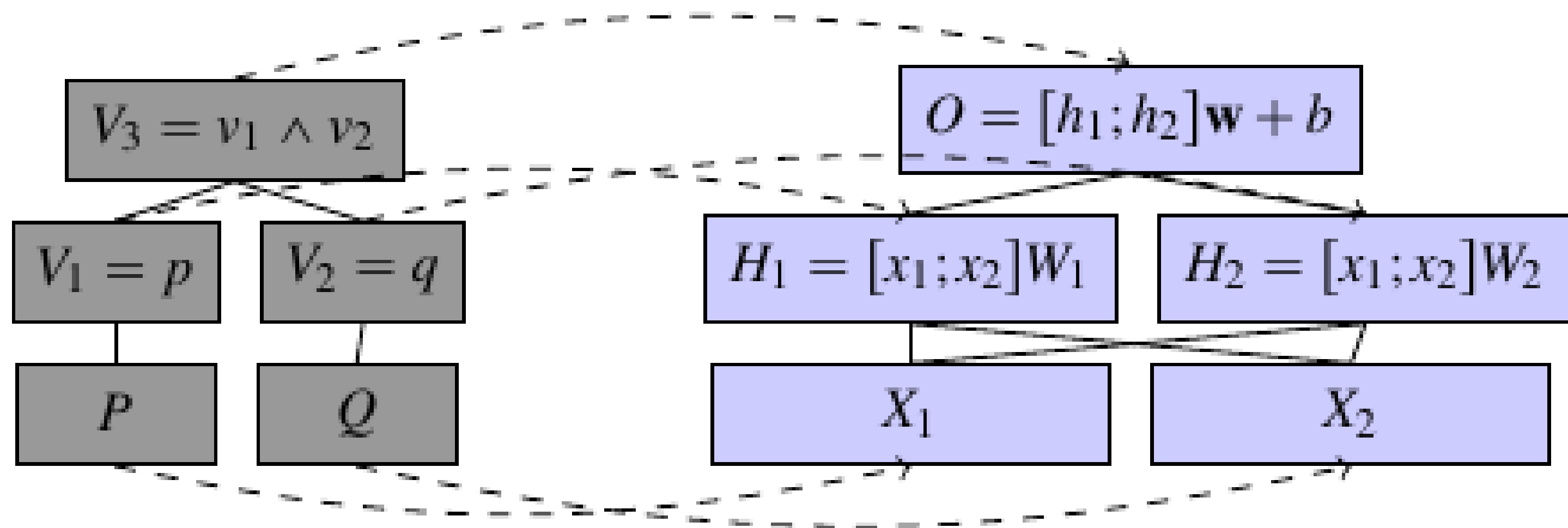
Why DoRA Works ? and if so why it is better than LoRA



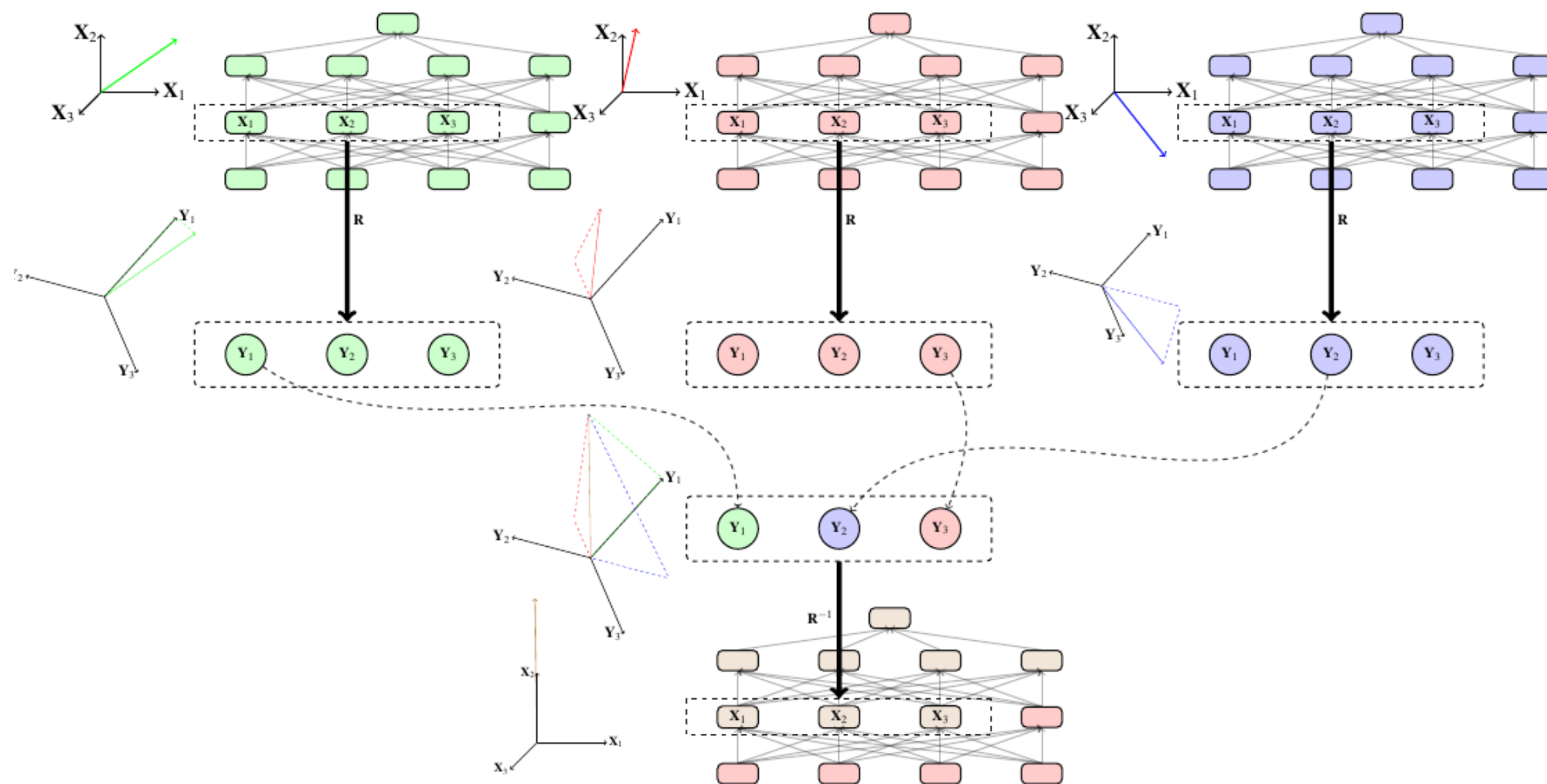
Prerequisite to ReFT

- DII (Distributed interchange intervention)
- DAS (Distributed Alignment Search)

DII and DAS

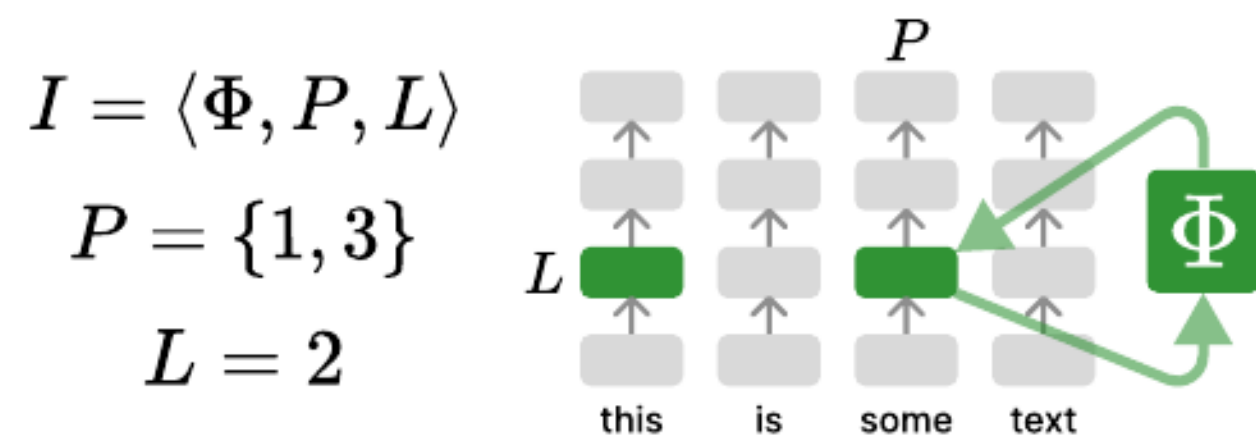


DII and DAS



ReFT

ReFT Intervention



LoReFT

