

# Limitations of Full Finetuning

## Instructor

Sourab Mangulkar

Machine Learning Engineer at 

Creator of  PEFT





# The Secret Sauce behind LLMs & Limitations

Recipe to get state-of-the-art results in NLP

- Pretraining on web-scale data
- Finetuning on downstream task

Scaling Laws

The larger the Deep Learning/AI model, the better the performance.

Limitations of full fine-tuning approaches

- Compute 
- Storage 

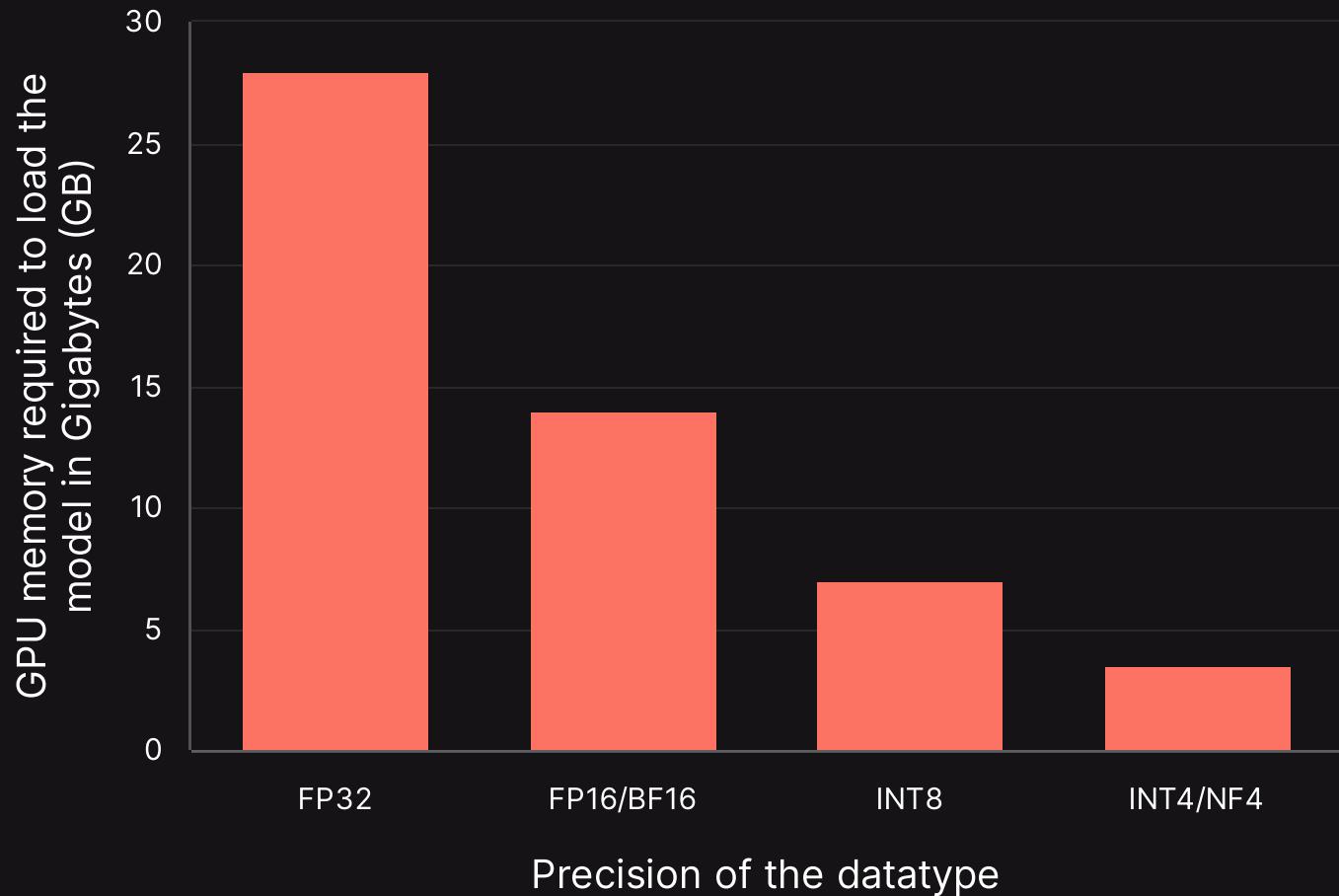
How to Finetune very large models on consumer hardware?

# What makes a model large?

## Two Factors

- 1 Number of parameters
- 2 Precision of the data type

Let's put some numbers wrt Mistral-7B model with 7 Billion parameters



# Why full finetuning is expensive?

*Finetuning Mistral-7B in mixed-precision using Adam Optimizer*

Weights - 2 bytes / parameter

Gradients - 2 bytes / parameter

Optimizer state - 4 bytes / parameter (FP32 copy) +  
8 bytes / parameter (momentum & variance estimates)

Total training cost: 16 bytes/parameter \* 7 billion parameters = 112  
GB

$$\begin{aligned} \cos q &= \frac{b^2 + c^2 - a^2}{2bc} \\ b^2 &= 4b^2c^2 - (b^2 - a^2) \Leftrightarrow 16 \\ q &= \frac{(b+c+a)}{2\cdot 3} \Leftrightarrow \\ +c-a \cdot \frac{(b+c+d)}{2} \\ S &= \frac{a}{2} \\ n &= a+3 \end{aligned}$$

# Hands On Project

## Instructor

Sourab Mangulkar

Machine Learning Engineer at 

Creator of  PEFT

