



### MASTERING Al WORKSHOP

Google
Developer
Student Clubs
Paris-Saclay

Lucas J. Velôso GDSC Lead

Soogle Developer Student Clubs



**Solution** Student Clubs

### Agenda



Google Developer Student Clubs

- 1. What are LLMs?
- 2. How LLMs work?
- 3. What is Fine-tuning?
- 4. What is Distillation?
- 5. What is **PEFT**?
- 6. What is **LoRA**
- 7. What is **QLoRA**?
- Coding time...

lookup.Static\



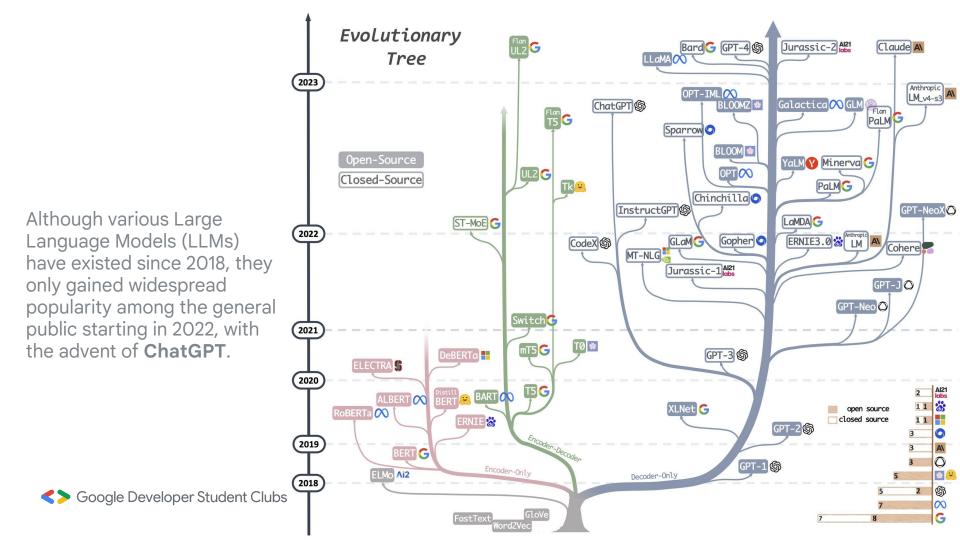
### What are LLMs?

lookup.KeyValuef.constant(['en =tf.constant([G .lookup.Static\

- 1. Large Language Models (LLMs) are advanced AI tools designed for processing and generating human-like text.
- 2. They **leverage vast amounts of data** to understand context, semantics, and language nuances.
- 3. LLMs are used in various applications like chatbots, content creation, and language translation.



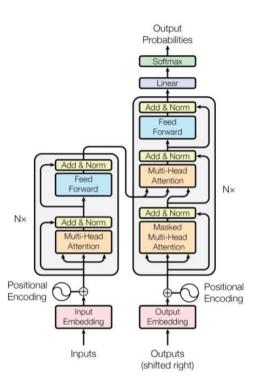




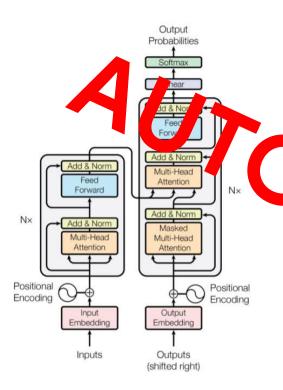


### How LLMs work?

lookup.KeyVal f.constant([' =tf.constant( .lookup.Stati



- 1. **Deep learning models**, primarily Transformers (Google 2017).
- 2. Multiple layers of processing (attention mechanisms, etc.)
- 3. Solves the task: Given the phrase with is the next word?



- Deep learning models, primarily Transformers (Google 2017).
- le layers of processing (attention
- ven the phrase with 3. Solves the task: G is the next word.

## Stochastic **Parrots**



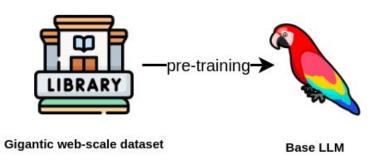


### Fine-Tuning



Gigantic web-scale dataset

### Fine-Tuning



# Fine-Tuning Line Line

Base LLM

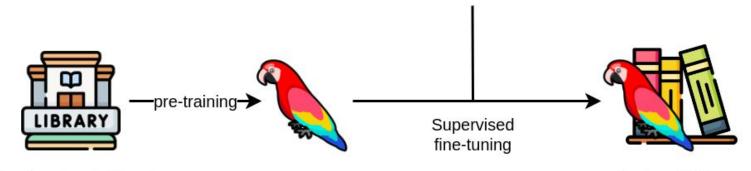
Specific (private) Knowledge Base



Gigantic web-scale dataset



Specific (private) Knowledge Base



Gigantic web-scale dataset

Base LLM

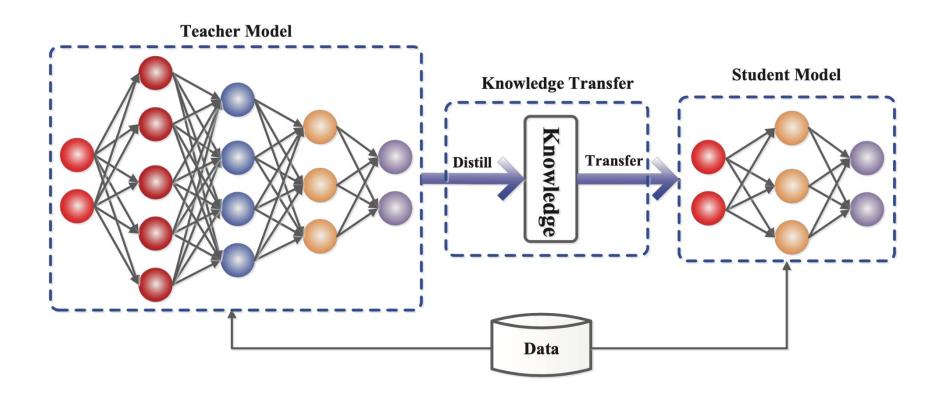
Fine-tuned LLM

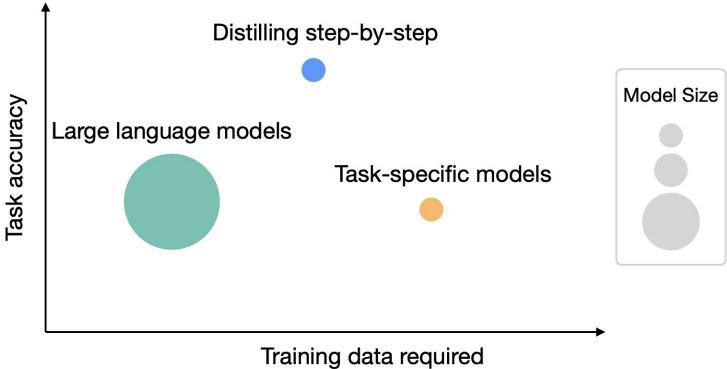




### What is Distillation?

lookup.KeyVal f.constant(['e =tf.constant( .lookup.Stati







### What is PEFT?

lookup.KeyValue f.constant(['en =tf.constant([@ .lookup.Static\

# Parameter-Efficient Fine-Tuning (PEFT)

Parameter-Efficient Fine-Tuning (PEFT) by **Hugging Face** is a technique to adapt large models for specific tasks by **fine-tuning only a few parameters**, reducing computational and storage demands





### What is LoRA?

lookup.KeyValue f.constant(['en =tf.constant([G .lookup.StaticV

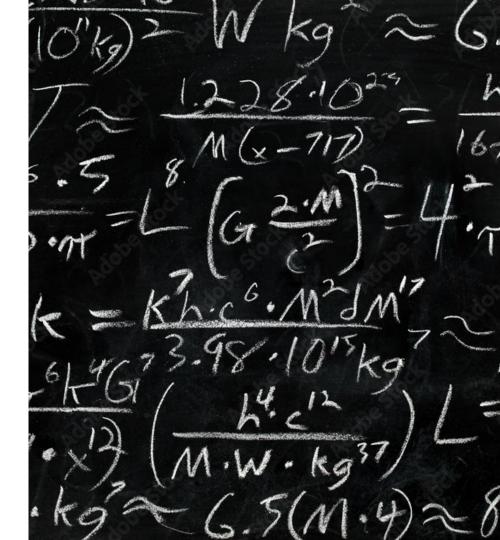
### LoRA

**Finetuned Weights** 

**Weight Update** 

$$\widetilde{W_{ t ft}} = W_{ t pt} + \widetilde{\Delta W}$$

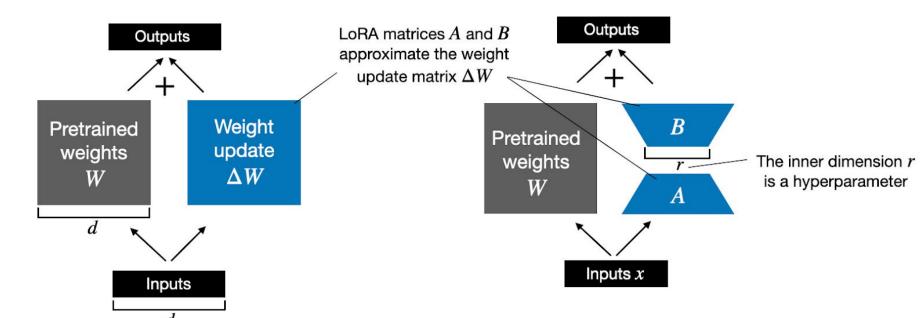
**Pretrained Weights** 





#### Weight update in regular finetuning

### Weight update in LoRA

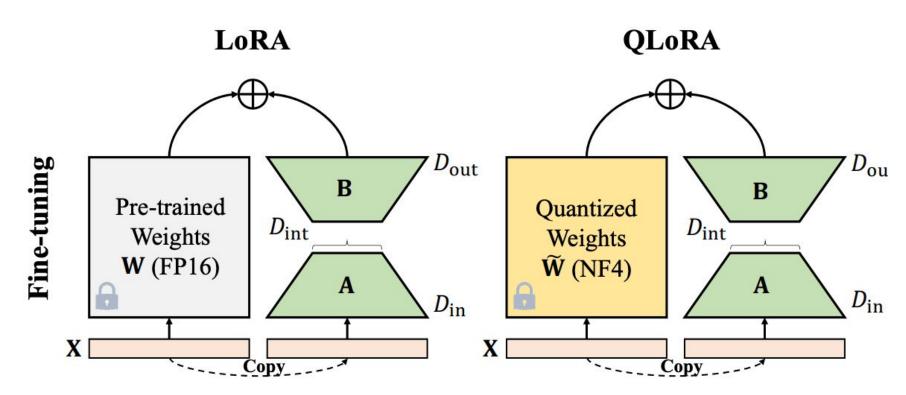




### What is QLoRA?

Lookup.KeyVal f.constant([' =tf.constant( .lookup.Stati

### **QLoRA**





### **Coding Time...**

lookup.KeyValue
f.constant(['en
=tf.constant([@
.lookup.Static\



### STARTUP CASE

lookup.KeyVal
f.constant(['
=tf.constant(
.lookup.Stati







Google Developer Student Clubs

1. It **must** run in a cell phone!!





- 1. It **must** run in a cell phone!!
- 2. It must run offline!!





- 1. It **must** run in a cell phone!!
- 2. It must run offline!!
- 3. The cost must be at maximum \$1



