# LLM fine-tuning on Palmetto

**Carl Ehrett**
**CCIT – Research Computing & Data (RCD)**

# Why people want to fine-tune

- The task the LLM will perform is not present in the training data
- The task the LLM will perform is much more specific and narrow than the training data
- To provide additional knowledge base

# Why people want to fine-tune

## Common use cases[1]:

- Setting the style, tone, format, or other qualitative aspects
- Improving reliability at producing a desired output
- Correcting failures to follow complex prompts
- Handling many edge cases in specific ways
- Performing a new skill or task that's hard to articulate in a prompt

1. https://platform.openai.com/docs/guides/fine-tuning/common-use-cases

# Why you probably don't need to fine-tune

Prompt engineering, few-shot learning (FSL), and retrieval-augmented generation (RAG) can accomplish all of these things *without* fine-tuning.

# Why you probably don't need to fine-tune

**Prompt engineering:** Designing and refining input instructions to elicit desired outputs from LLMs, optimizing clarity, relevance, or goal alignment.

## Poor prompt for NER

**User:** Who or what is referred to in "The Apple executive executed an apple tree."

**AI:** The phrase "The Apple executive executed an apple tree" is an example of a sentence with ambiguous meaning due to the homonym "apple" which can refer to both the fruit and the technology company. Without further context, it is difficult to determine the precise meaning intended by the sentence.

Prompt Eng.

## Improved prompt for NER

**User:** Produce a list of every entity in the following text: "The Apple executive executed an apple tree." Include no additional text.

**AI:**
1. Apple (company)
2. Executive (position/role)
3. Apple tree

# Why you probably don't need to fine-tune

**Few-shot learning:** Providing the model with a few example cases at inference (not during training), to improve model output.

## Without FSL

**User:** Extract every entity in the following text: "The Apple executive executed an apple tree." Output JSON for your response. Include no additional text.

**AI:**
```
{
  "entities": [
    {"entity": "Apple", "type": "Organization"},
    {"entity": "executive", "type": "JobTitle"},
    {"entity": "apple tree", "type": "Plant"}
  ]
}
```

## Without FSL

**User:** Extract every entity in the following text: "The Apple executive executed an apple tree." Output JSON for your response. Include no additional text.

**AI:**
```
Sure! Here is JSON for those entities.
{
  "entities": [
    {"text": "Apple", "type": "Organization"},
    {"text": "executive", "type": "Title"},
    {"text": "apple tree", "type": "Plant"}
  ]
}
```

## With FSL

### FSL prompt

**User:** Extract every entity in the following text into JSON: "F.E. Church attended this church."

**AI:** `{"F.E. Church": "Person", "church": "Structure"}`

**User:** Extract every entity in the following text into JSON: "Nvidia's CEO checked the stocks and leaped for joy."

**AI:** `{"NVidia": "Organization", "CEO": "JobTitle", "stocks": "Financial", "joy": "Emotion"}`

**User:** Extract every entity in the following text into JSON: "The Apple executive executed an apple tree."

### AI response

**AI:** `{"Apple": "Organization", "executive": "JobTitle", "apple tree": "Plant"}`

COMPUTING AND INFORMATION TECHNOLOGY
*Research Computing and Data*

# Why you probably don't need to fine-tune

**Retrieval-augmented generation:** The model dynamically retrieves information from external databases or documents to enhance its responses.
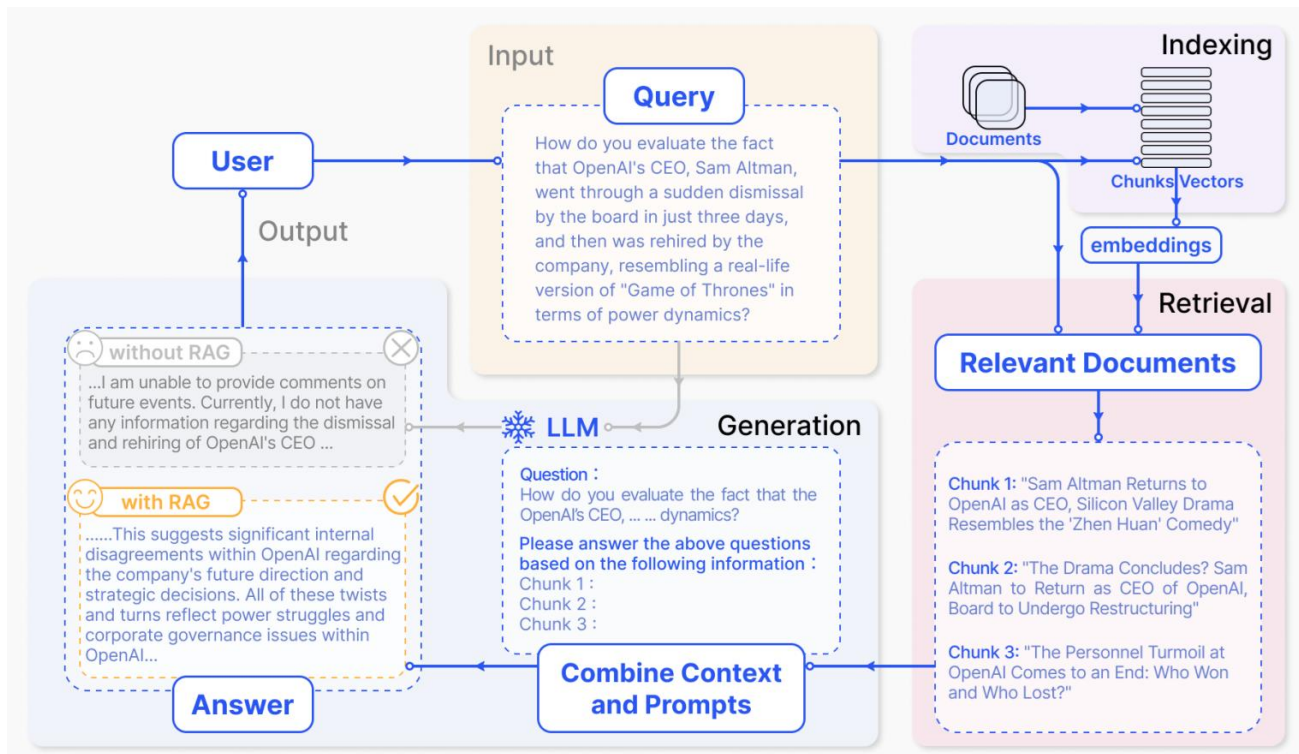
**Image source**
Retrieval-Augmented Generation for Large Language Models: A Survey (https://arxiv.org/abs/2312.10997)

# Why people want to fine-tune (but don't need to)

- The task the LLM will perform is not present in the training data
- The task the LLM will perform is much more specific and narrow than the training data
- To provide additional knowledge base

Prompt Engineering

Few-shot learning

Retrieval-augmented generation

# Why you should fine-tune anyway

- Fine-tuning can achieve (sometimes small) quality improvements over other methods
- In some cases, your FSL examples might be too large for a model's context window
- Smaller LLMs can benefit more from fine-tuning than large ones
- If you expect to use the model for inference many times, fine-tuning can reduce your prompt size and thus reduce the cost of inference

# LLM fine-tuning

*What data is required to fine-tune an LLM?*

# How much data is needed?

Compared to training an LLM from scratch – extremely little.

- As few as ten can see benefits, but 50-100 are often needed. Typically, more is better.
- It is advisable to set aside some data as a holdout set, to evaluate the results of the fine-tuning.
  - Consider whether it is appropriate in your case for the evaluation set to be distributionally different from the training set (e.g., social media posts collected after the training data).

**Image source**
Few-shot Fine-tuning vs. In-context Learning: A Fair Comparison and Evaluation (https://arxiv.org/abs/2305.16938)

# What kind of data is needed?

- For fine-tuning, it is crucial that your data be of consistent high quality, especially if small.

- Your data should constitute prompt/completion pairs representative of your desired LLM behavior.

- Data can be stored as, e.g., a csv with a column for the prompt and one for the completion; or similarly, as JSON.

COMPUTING AND
INFORMATION TECHNOLOGY
*Research Computing and Data*

# Fine-tuned models can experience poor out-of-distribution performance

Especially with small data, fine-tuned models can easily overfit, resulting in poor performance on tasks that differ from the training set.

**Image source**
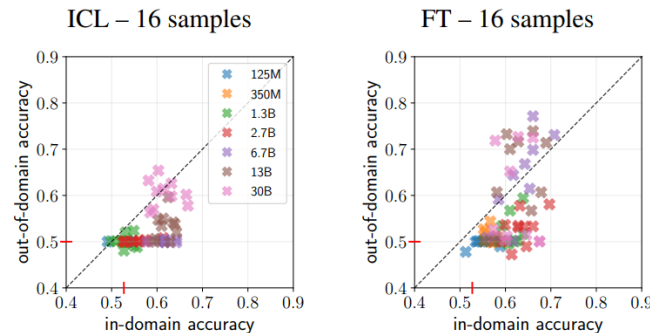Few-shot Fine-tuning vs. In-context Learning: A Fair Comparison and Evaluation (https://arxiv.org/abs/2305.16938)
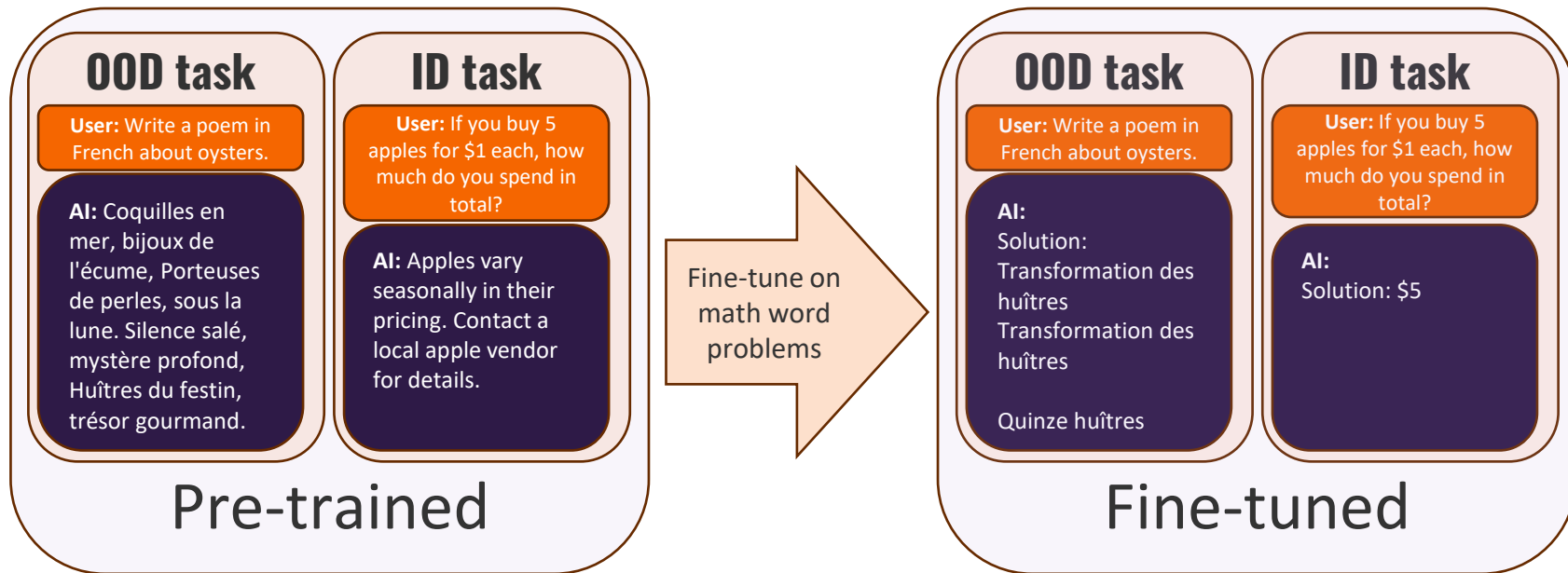


Figure 1: In-domain (RTE) and out-of-domain performance (HANS) for in-context learning (ICL) and fine-tuning (FT) with OPT models of various sizes. We fine-tune models using pattern-based fine-tuning. We report results using 10 different data seeds. When using 16 samples, ICL's performance with a 30B model is comparable to that of FT with smaller models (6.7B) and for most model sizes, FT outperforms ICL (see Table 1a for significance tests). — in the x- and y-axes indicates majority class accuracy.

# Fine-tuned models can suffer catastrophic forgetting

**Catastrophic forgetting** occurs when the model loses or significantly degrades its performance on previously learned tasks or knowledge after being fine-tuned on new data or tasks.

## Pre-trained

### OOD task

**User:** Write a poem in French about oysters.

**AI:** Coquilles en mer, bijoux de l'écume, Porteuses de perles, sous la lune. Silence salé, mystère profond, Huîtres du festin, trésor gourmand.

### ID task

**User:** If you buy 5 apples for $1 each, how much do you spend in total?

**AI:** Apples vary seasonally in their pricing. Contact a local apple vendor for details.

Fine-tune on math word problems →

## Fine-tuned

### OOD task

**User:** Write a poem in French about oysters.

**AI:**
Solution:
Transformation des huîtres
Transformation des huîtres

Quinze huîtres

### ID task

**User:** If you buy 5 apples for $1 each, how much do you spend in total?

**AI:**
Solution: $5

# Fine-tuning locks you into a single model

If you devote time and energy to prompt engineering and developing a pipeline for RAG and FSL, all of this work can be easily transferred to any LLM.

By contrast, fine-tuning results in a single LLM. If a new, more powerful LLM becomes available next month, you would need to fine-tune all over again.

PE, RAG and FSL can also all be used on closed-source models that cannot be fine-tuned.

# Fine-tuning can be unsafe!

**Safety alignment:** Many models are trained to align with human values and ethical standards, aiming to prevent a range of harms or unintended consequences (e.g. biased or discriminatory content, revealing private information, promoting misinformation, etc.)

Fine-tuning can **erode** safety alignment training, even if your new data is benign!



Figure 1: **(Overview) Fine-tuning GPT-3.5 Turbo leads to safety degradation: as judged by GPT-4, harmfulness scores (1~5) increase across 11 harmfulness categories after fine-tuning.** Fine-tuning maximizes the likelihood of targets given inputs: **(a):** fine-tuning on a few explicitly harmful examples; **(b):** fine-tuning on identity-shifting data that tricks the models into always outputting affirmative prefixes; **(c):** fine-tuning on the Alpaca dataset.

COMPUTING AND
INFORMATION TECHNOLOGY
*Research Computing and Data*

# COMPUTING AND INFORMATION TECHNOLOGY
## Research Computing and Data

# LLM fine-tuning

*What does the model learn during fine-tuning?*

CLEMSON

# Supervised fine-tuning (SFT)

**Training data:** "Palmetto supports the research mission of Clemson University through innovative High Performance Computing (HPC) and Storage solutions."

| Model inputs | Model outputs before updating | | Model outputs after updating |
|---|---|---|---|
| Palmetto supports | a 0.31<br>aardvark 0.04<br>…<br>the 0.32<br>…<br>zylophone 0.003 | Update model weights to prefer true next word → | a 0.20<br>aardvark 0.002<br>…<br>the 0.45<br>…<br>zylophone 0.002 |
| Palmetto supports the | a 0.00001<br>aardvark 0.034<br>…<br>research 0.25<br>…<br>zylophone 0.003 | Update model weights to prefer true next word → | a 0.00001<br>aardvark 0.029<br>…<br>research 0.37<br>…<br>zylophone 0.002 |

The most basic fine-tuning approach, which is the same as the (usual) pre-training objective, to correctly predict the next token in each of a set of text documents.

Can be used for:
- Task adaptation
- Instruction following
- Alignment tuning

**COMPUTING AND INFORMATION TECHNOLOGY**
*Research Computing and Data*

# Preference optimization (RLHF)

## Step 1
**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

> Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

> Some people went to the moon...

This data is used to fine-tune GPT-3 with supervised learning.

> SFT

## Step 2
**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

> Explain the moon landing to a 6 year old

> A — Explain gravity...
> B — Explain war...
> C — Moon is natural satellite of...
> D — People went to the moon...

A labeler ranks the outputs from best to worst.

> D > C > A = B

This data is used to train our reward model.

> RM
> D > C > A = B

## Step 3
**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

> Write a story about frogs

The policy generates an output.

> PPO

> Once upon a time...

The reward model calculates a reward for the output.

> RM

The reward is used to update the policy using PPO.

> $r_k$

Instead of training on fixed text, the model is fine-tuned to **prefer certain responses** (usually based on human feedback).

Very difficult!

# Preference optimization (DPO)



Alternatives to RLHF have been developed. **DPO (direct preference optimization)** accomplishes similar aims and only requires a dataset of preferred and dispreferred model outputs.
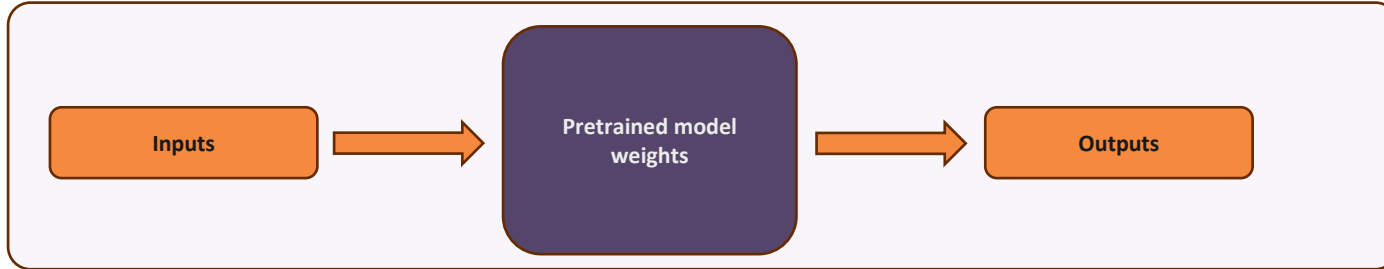
# LLM fine-tuning

*How is the model updated during fine-tuning?*

# Full fine-tune

Just like in pre-training, in a full fine-tune all model parameters are updated.

This is the most computationally expensive method, and prone to catastrophic forgetting.

Advisable only for highly customized task-specific models.

# Low-rank adaptation (LoRA) fine-tune

Adds small trainable adapter layers, and trains *only* these new layers.

Reduces memory and compute significantly, and reduces catastrophic forgetting.

Usually best.



COMPUTING AND
INFORMATION TECHNOLOGY
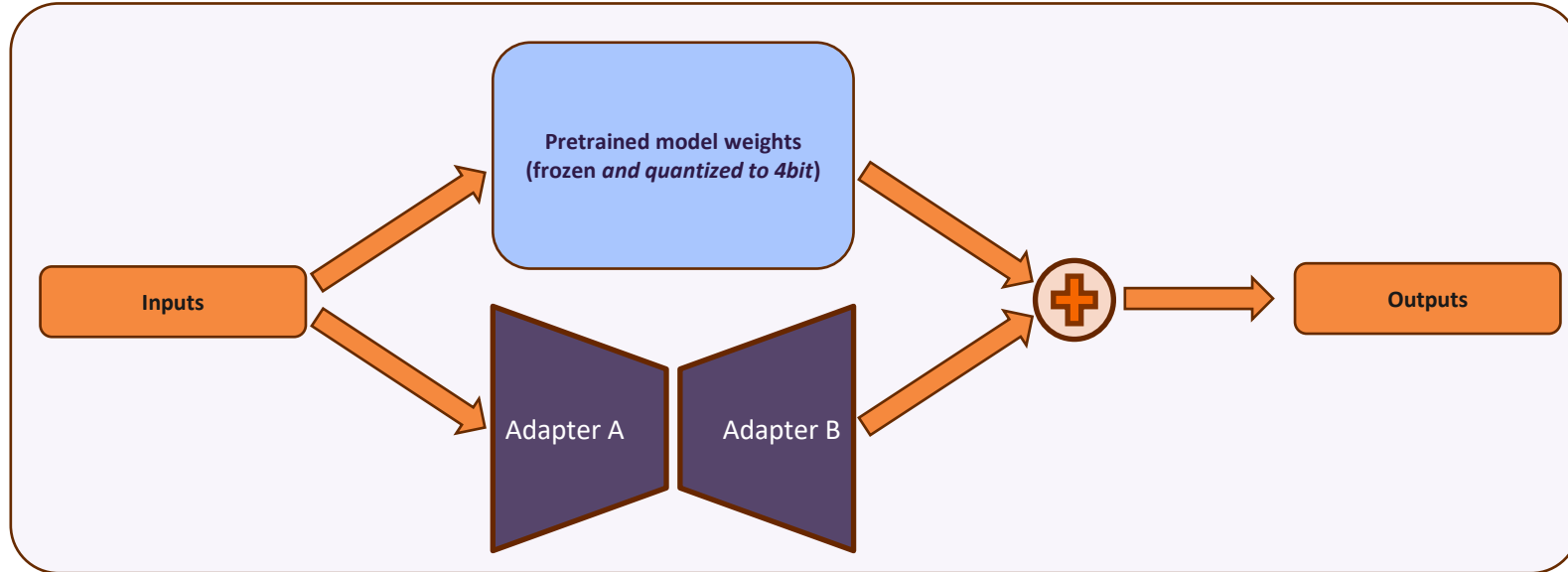*Research Computing and Data*

# QLoRA fine-tune

Just like LoRA, but also quantizes the model to 4-bit precision, further reducing memory requirements.

Allows for fine-tuning models otherwise too large to fit in available memory.

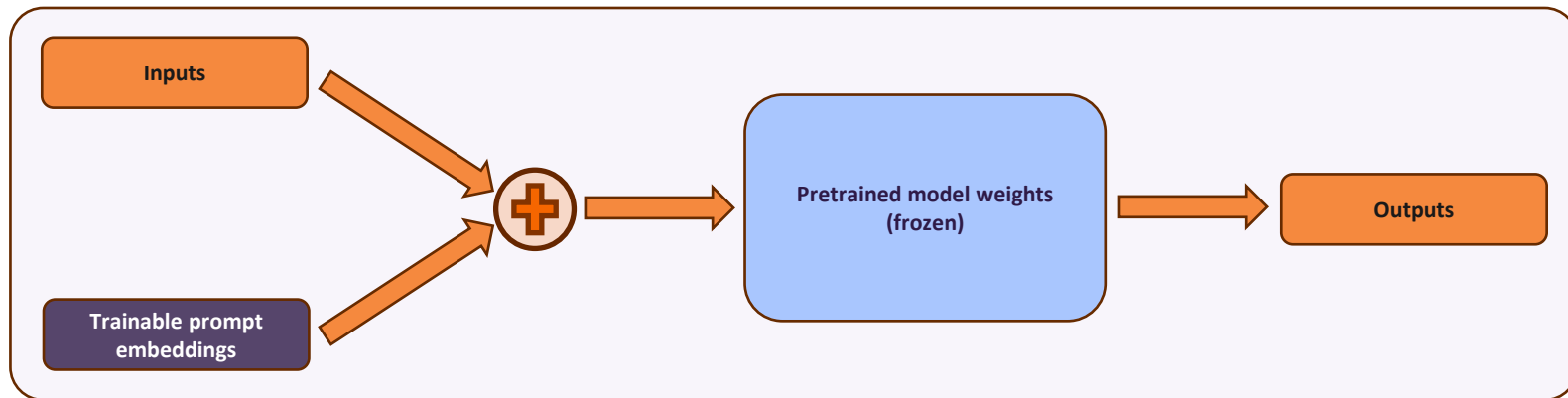Best when the available hardware isn't large enough to just do LoRA.



Inputs

Pretrained model weights
(frozen *and quantized to 4bit*)

Adapter A    Adapter B

Outputs

# Prompt tuning

Instead of modifying the model weights, learns small trainable prompt embeddings.

Less powerful than LoRA (sometimes a good thing!) and very efficient.

Best for: task adaptation without modifying the model.

# Fine-tuning methods comparison

| Method | Trainable parameters | Computational cost | Memory usage | Retains general knowledge? | Best use cases |
|---|---|---|---|---|---|
| Full fine-tuning | All model parameters | High | Very high | Risk of catastrophic forgetting | Task-specific, highly customized models |
| LoRA | Small adapter layers | Low | Moderate | Yes | Domain adaptation, style control |
| QLoRA | Small adapter layers | Very low | Very low | Yes | Fine-tuning very large LLMs (relative to hardware) |
| Prompt tuning | Very small prompt embeddings | Low | Low | Yes | Task adaptation with minimal change to model |

# LLM fine-tuning

*How can we manage experiments and log results?*

# Why experiment management matters

**Prompt engineering:** Designing and refining input instructions to elicit desired outputs from LLMs, optimizing clarity, relevance, or goal alignment.

## Reproducibility

Enables you to replicate results by tracking hyperparameters, code versions, and dataset versions

## Comparability

Helps you compare different experiments side by side, identifying what works best

## Collaboration

Facilitates sharing results and insights within teams

## Debugging

Easier to diagnose issues by examining logs and metrics over time

# Weights & Biases (W&B) for tracking experiments

**Weights & Biases**

**W&B** is a tool for logging and tracking ML experiments.

Provides **real-time visualizations** of metrics, hyperparameters, and training progress.

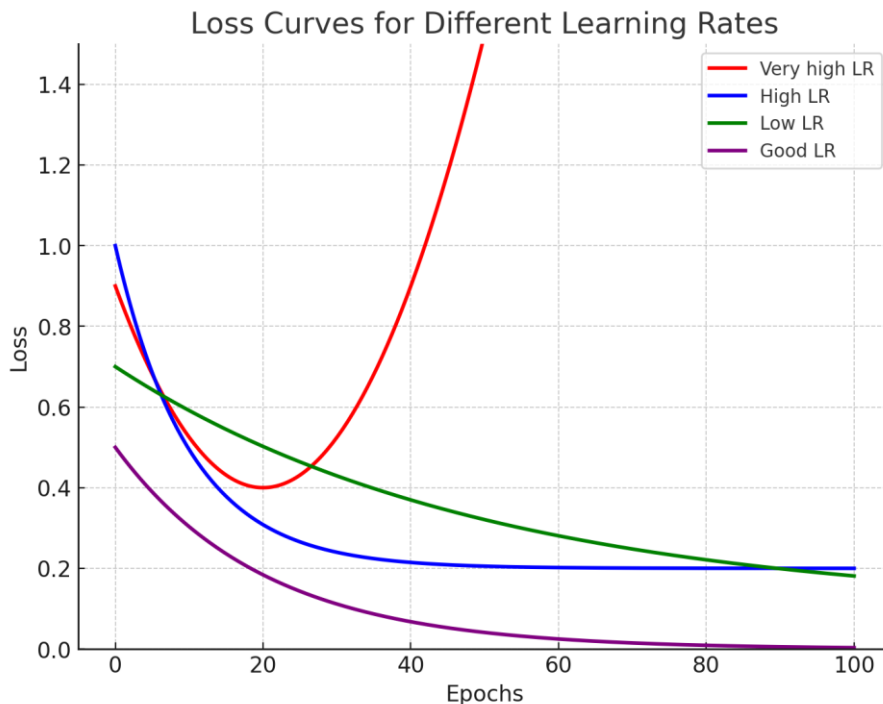Integrates with PyTorch, Tensorflow, Hugging Face, and more.

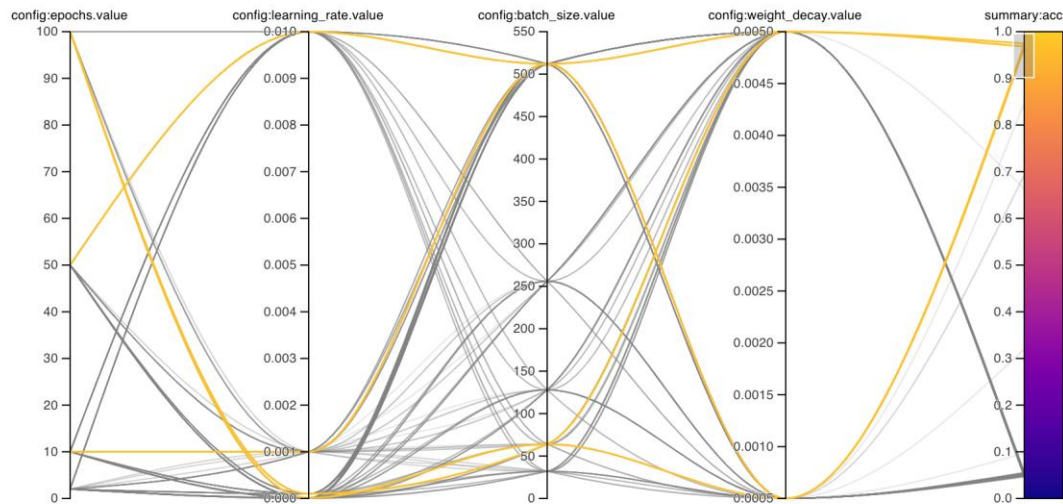Stores logs in the cloud for collaboration and reproducibility.

# Hyperparameter selection

**Why hyperparameter selection matters:**

- Different hyperparameter settings can significantly affect performance
- Finding the right values is often trial and error
- Can optimize for speed vs. accuracy tradeoffs



Loss Curves for Different Learning Rates

# Hyperparameter tuning using sweeps



**It is highly advisable to automate your hyperparameter search.**

W&B makes hyperparameter "sweeps" easy to set up and evaluate.

# Batch size considerations

**Batch size** is the number of training samples processed simultaneously during training.

Batch size affects how much memory is used during training, **and** affects training outcomes.
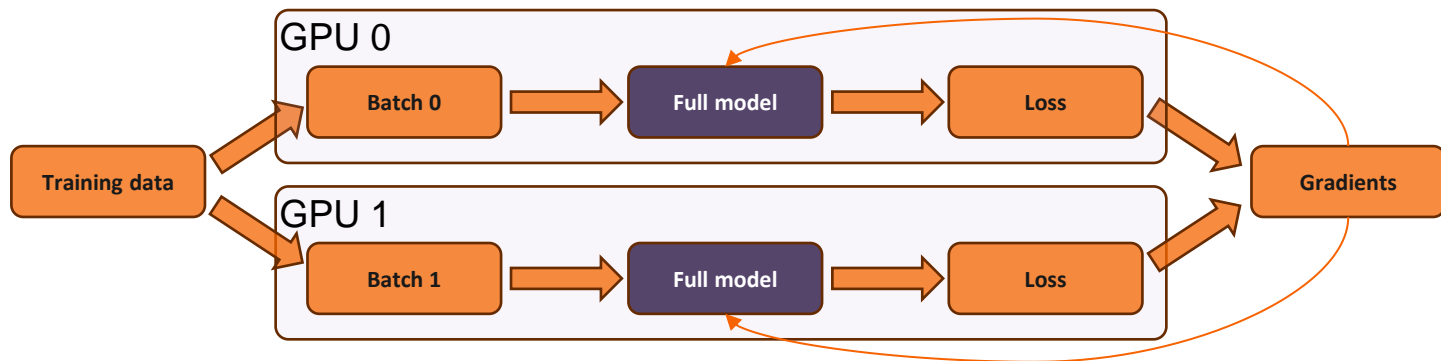
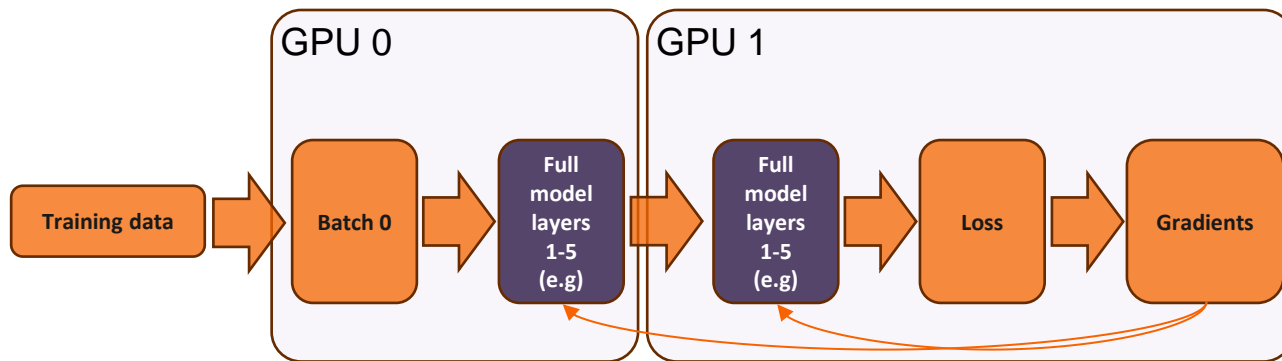| Batch size | Advantages | Disadvantages |
|---|---|---|
| Large | • Faster training per epoch<br>• More stable gradients<br>• Better GPU utilization | • Risk of poor generalization<br>• Can lead to training instability<br>• High memory usage |
| Small | • More frequent updates and faster convergence<br>• Better generalization<br>• Lower memory usage | • Noisier gradients and training instability<br>• Slower per-epoch training<br>• May require gradient accumulation to match large batch performance |

# Data parallelism

There are many ways to use **multiple GPUs** during training, with different goals and outcomes.

**Data parallelism** puts a copy of the model on **each** GPU, speeding training.
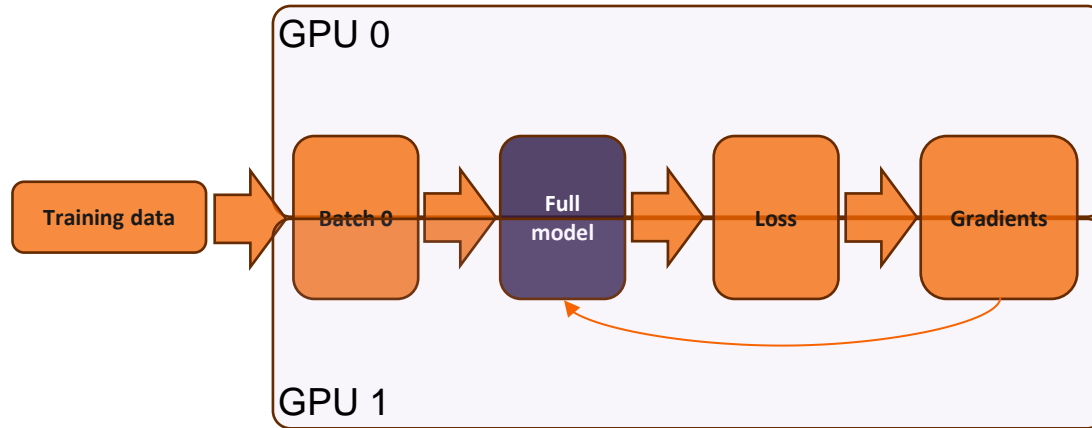
# Model parallelism

**Model parallelism** puts a some of the layers of the model on each GPU, allowing the use of a model that is too large to fit on a single GPU.

# Tensor parallelism

**Tensor parallelism** splits individual layers of the model across GPUs, allowing the use of a model with individual layers so large they can't fit on a single GPU.

# Comparison of parallelism strategies

| Parallelism type | What it splits | Best for | Pros | Cons |
|---|---|---|---|---|
| Data parallelism | Dataset (batch split across GPUs) | • Large datasets<br>• Small to medium models | • Easy to implement<br>• Works with most models | • Requires full model copy on each GPU<br>• Some overhead |
| Model parallelism | Model (some layers on each GPU) | Large models that don't fit on one GPU | • Reduces memory load per GPU<br>• Usually easy to implement | • Slower due to inter-GPU communication<br>• Can be tough to implement |
| Tensor parallelism | Individual layers (weights split across models) | Extremely large models with layers too big for 1 GPU | • Enables massive model scaling<br>• More memory efficient than model parallelism | • High communication overhead<br>• Very difficult to implement, requires specialized libraries |