

# Large Language Models: Fine-tuning

Nathanaël Fijałkow  
CNRS, LaBRI, Bordeaux



LaBRI

université  
de BORDEAUX

# FINE-TUNING

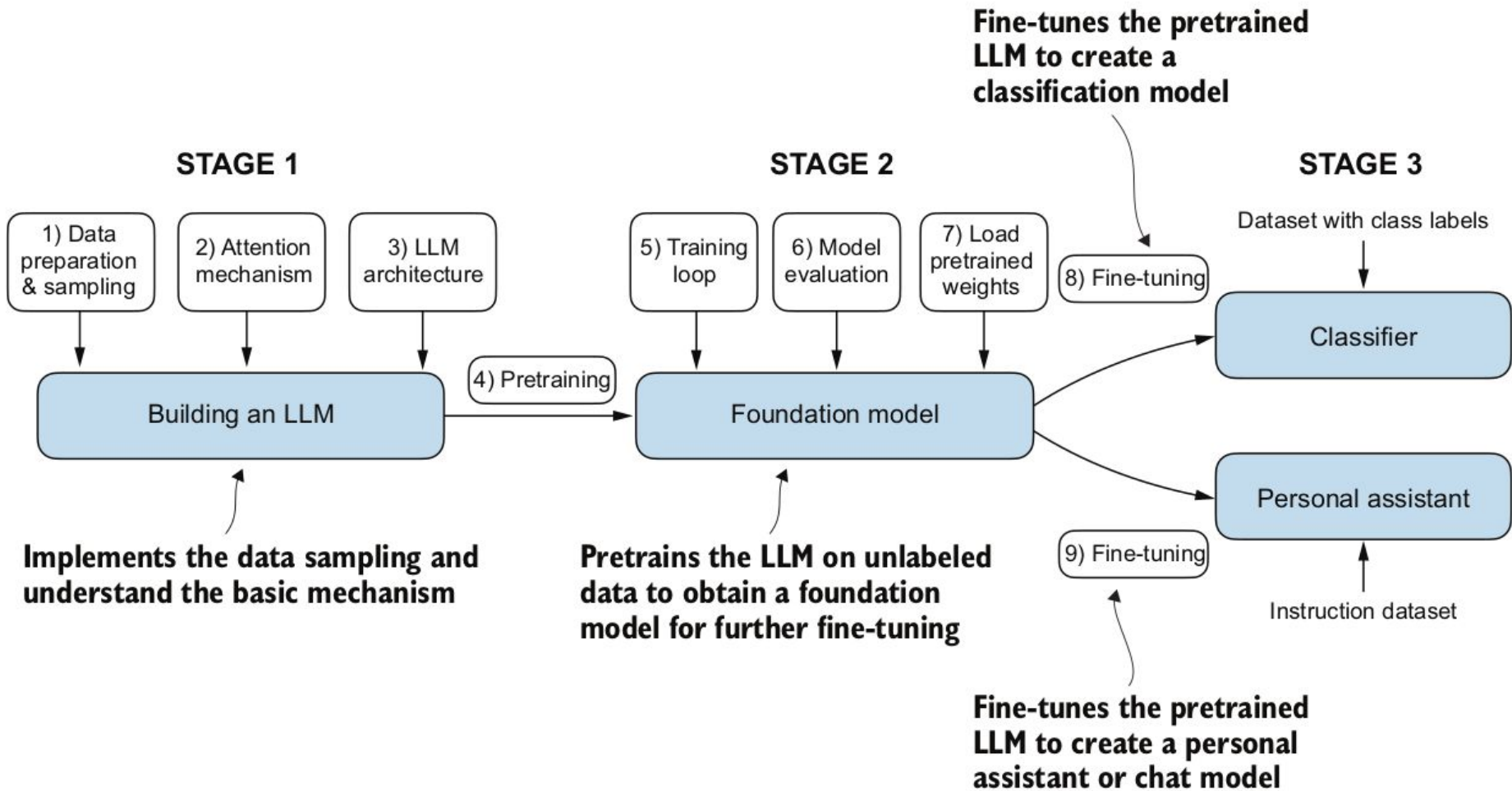
- General overview
- LoRA

---

# FOUNDATION MODELS

Language Models are not very useful, they randomly generate texts... But this means that they somehow capture some information from natural language! They are also called *foundation models*.

*Fine-tuning* is about making Language Models solve concrete tasks, like classification, question answering, name entity recognition...



# TRAINING IS EXPENSIVE

We often cannot afford updating the *whole* model!

Most of us will not train foundation models... Rather  
fine-tune existing ones.

# LOW-RANK ADAPTATION (LORA)

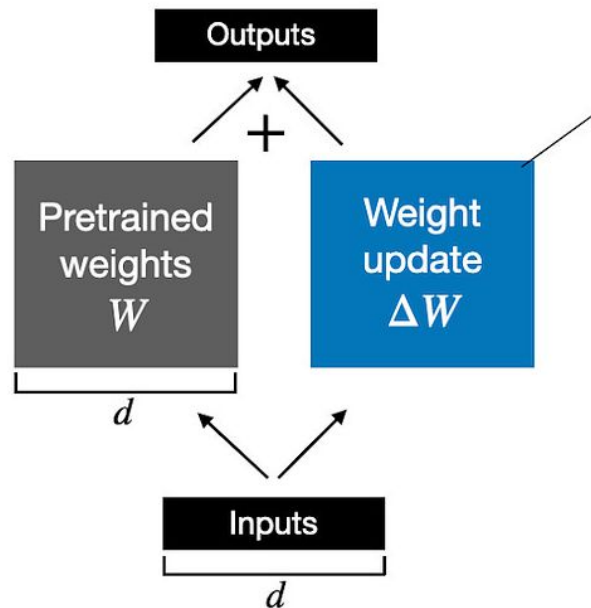
Two key ideas:

- (1) We only store the changes, not a new model
- (2) We only update a small number of parameters

# IDEA: STORING WEIGHT UPDATES

Say we consider a linear layer with matrix  $W$ . We keep the matrix  $W$  fixed and store  $\Delta W$

## Weight update in **regular finetuning**

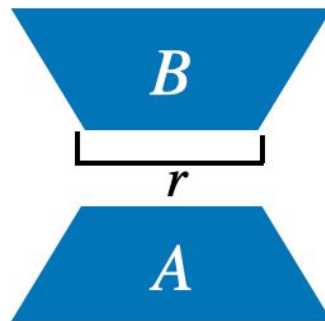


# RANK APPROXIMATIONS

A matrix  $W$  of dimension  $d \times d$  contains  $d^2$  parameters. It can be **rank- $r$  approximated** by two matrices  $A \times B$  with:

- $A$  of dimension  $d \times r$
- $B$  of dimension  $r \times d$

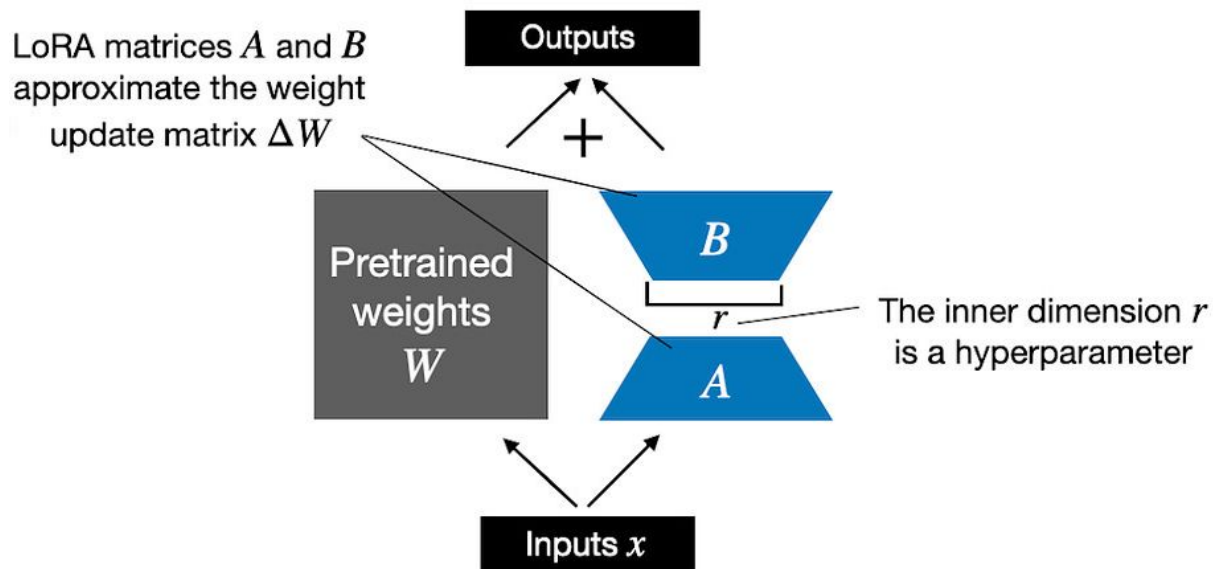
Instead of  $d^2$  parameters we now have  $2 \times d \times r$  parameters.





# WEIGHT UPDATE

## Weight update in LoRA



# LORA LAYER

```
import math

class LoRALayer(torch.nn.Module):
    def __init__(self, in_dim, out_dim, rank, alpha):
        super().__init__()
        std_dev = 1 / torch.sqrt(torch.tensor(rank).float())
        self.A = nn.Parameter(torch.randn(in_dim, rank) * std_dev)
        self.B = nn.Parameter(torch.zeros(rank, out_dim))
        self.alpha = alpha

    def forward(self, x):
        x = self.alpha * (x @ self.A @ self.B)
        return x
```

# ADDING THE LORA LAYER

```
class LinearWithLoRA(torch.nn.Module):
    def __init__(self, linear, rank, alpha):
        super().__init__()
        self.linear = linear
        self.lora = LoRALayer(
            linear.in_features, linear.out_features, rank, alpha
        )

    def forward(self, x):
        return self.linear(x) + self.lora(x)
```

```
def replace_linear_with_lora(model, rank, alpha):
    for name, module in model.named_children():
        if isinstance(module, torch.nn.Linear):
            # Replace the Linear layer with LinearWithLoRA
            setattr(model, name, LinearWithLoRA(module, rank, alpha))
        else:
            # Recursively apply the same function to child modules
            replace_linear_with_lora(module, rank, alpha)
```

- We then freeze the original model parameter and use the `replace_linear_with_lora` to replace the said `Linear` layers using the code below
- This will replace the `Linear` layers in the LLM with `LinearWithLoRA` layers

```
total_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
print(f"Total trainable parameters before: {total_params:,}")

for param in model.parameters():
    param.requires_grad = False

total_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
print(f"Total trainable parameters after: {total_params:,}")
```

Total trainable parameters before: 124,441,346  
 Total trainable parameters after: 0

```
replace_linear_with_lora(model, rank=16, alpha=16)

total_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
print(f"Total trainable LoRA parameters: {total_params:,}")
```

Total trainable LoRA parameters: 2,666,528

```
print(model)
```

```
GPTModel(  
  (tok_emb): Embedding(50257, 768)  
  (pos_emb): Embedding(1024, 768)  
  (drop_emb): Dropout(p=0.0, inplace=False)  
  (trf_blocks): Sequential(  
    (0): TransformerBlock(  
      (att): MultiHeadAttention(  
        (W_query): LinearWithLoRA(  
          (linear): Linear(in_features=768, out_features=768, bias=True)  
          (lora): LoRALayer()  
        )  
        (W_key): LinearWithLoRA(  
          (linear): Linear(in_features=768, out_features=768, bias=True)  
          (lora): LoRALayer()  
        )  
        (W_value): LinearWithLoRA(  
          (linear): Linear(in_features=768, out_features=768, bias=True)  
          (lora): LoRALayer()  
        )  
      )  
    )  
  )  
)
```