# A reminder

# Preference Optimisation

Fine-tuning, RL, LoRA, and all that jazz

# Task

- Learn about PEFT, Adaptors, LoRA, and all that jazz

- Implement an RLHF pipeline, starting from base model

  - take a base model (Qwen3)
  - do SFT on summary dataset
  - train a reward model
  - train a policy model to max rewards

- Have fun, is the last week :)

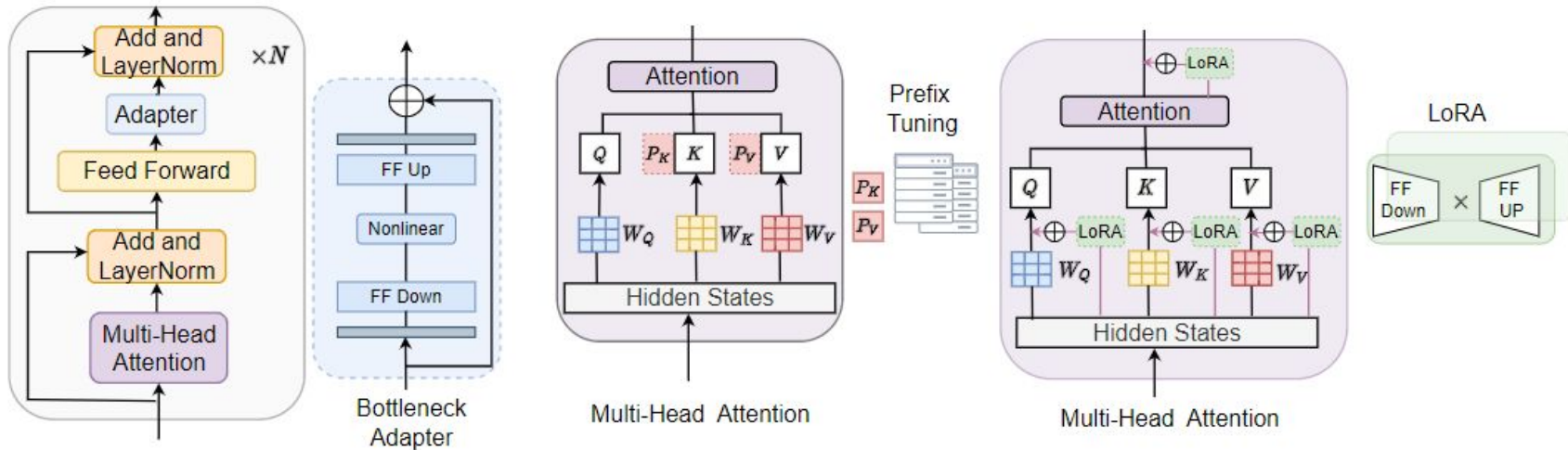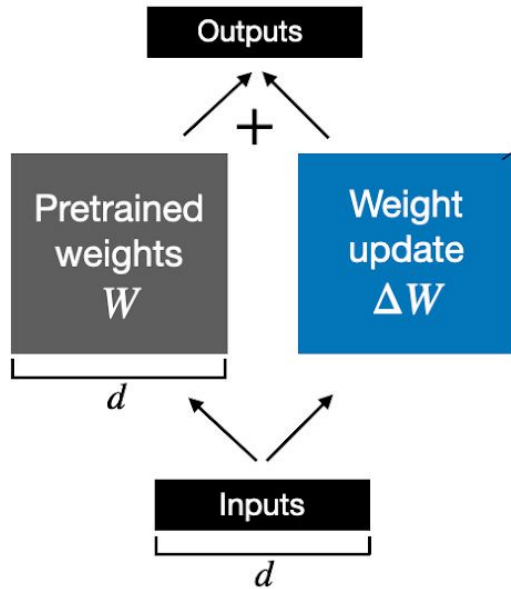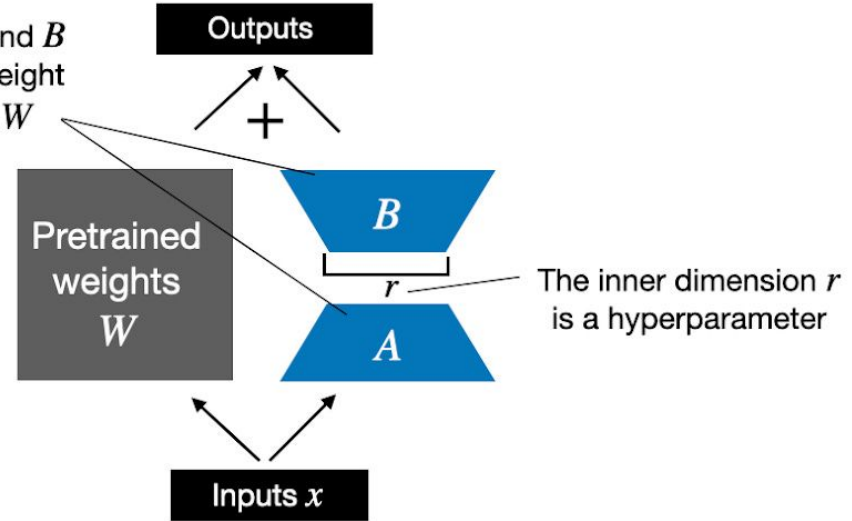# Parameter-Efficient Fine-Tuning



**Fig. 2**. Transformer architecture along with Adapter, Prefix Tuning, and LoRA.

# Overview

# Code

**Panel 1:**

```python
#
#
import torch


#
#
#
class Base(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.proj = torch.nn.Linear(10, 5)

    def forward(self, x):
        x = self.proj(x)
        return x


#
#
base = Base()
nums = sum(p.numel() for p in base.parameters())


#
#
#
print('Params: ', nums)
print(base)
```

```
Params: 55
Base(
  (proj): Linear(in_features=10, out_features=5, bias=True)
)
```

**Panel 2:**

```python
#
#
import torch


#
#
class Base(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.proj = torch.nn.Linear(10, 5)
        self.lora = torch.nn.Sequential(
            torch.nn.Linear(10, 2, bias=False),
            torch.nn.Dropout(0.1),
            torch.nn.Linear(2, 5, bias=False),
        )

    def forward(self, x):
        x = self.proj(x)
        return x + self.lora(x)

    #
base = Base()
nums = sum(p.numel() for p in base.parameters())


#
#
print('Params: ', nums)
print(base)
```

```
Params: 85
Base(
  (proj): Linear(in_features=10, out_features=5, bias=True)
  (lora): Sequential(
    (0): Linear(in_features=10, out_features=2, bias=False)
    (1): Dropout(p=0.1, inplace=False)
    (2): Linear(in_features=2, out_features=5, bias=False)
  )
)
```

**Panel 3:**

```python
#
#
import torch


#
#
class Base(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.proj = torch.nn.Linear(10, 5)
        self.adpt = torch.nn.Sequential(
            torch.nn.Linear(10, 2, bias=False),
            torch.nn.ReLU(),
            torch.nn.Dropout(0.1),
            torch.nn.Linear(2, 5, bias=False),
        )

    def forward(self, x):
        x = self.proj(x)
        return x + self.adpt(x)

    #
base = Base()
nums = sum(p.numel() for p in base.parameters())


#
#
print('Params: ', nums)
print(base)
```

```
Params: 85
Base(
  (proj): Linear(in_features=10, out_features=5, bias=True)
  (adpt): Sequential(
    (0): Linear(in_features=10, out_features=2, bias=False)
    (1): ReLU()
    (2): Dropout(p=0.1, inplace=False)
    (3): Linear(in_features=2, out_features=5, bias=False)
  )
)
```

**Panel 4:**

```python
#
#
import torch
import peft

#
#
class Base(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.proj = torch.nn.Linear(10, 5)

    def forward(self, x):
        x = self.proj(x)
        return x

    #
    #
conf = peft.LoraConfig(r=2, target_modules=['proj'])

    #
    #
base = peft.get_peft_model(Base(), conf)
nums = sum(p.numel() for p in base.parameters())

    #
    #
print('Params: ', nums)
print(base)
```
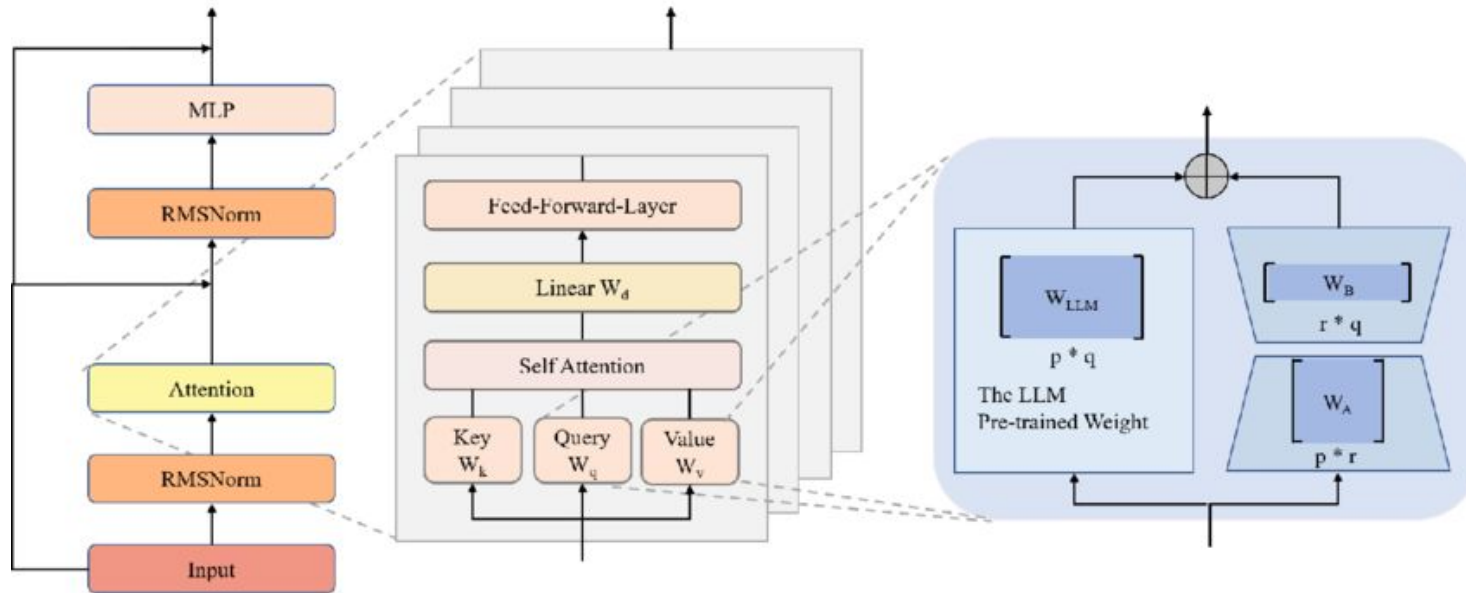
```
Params: 85
PeftModel(
  (base_model): LoraModel(
    (model): Base(
      (proj): lora.Linear(
        (base_layer): Linear(in_features=10, out_features=5, bias=True)
        (lora_dropout): ModuleDict(
          (default): Identity()
        )
        (lora_A): ModuleDict(
          (default): Linear(in_features=10, out_features=2, bias=False)
        )
        (lora_B): ModuleDict(
          (default): Linear(in_features=2, out_features=5, bias=False)
        )
        (lora_embedding_A): ParameterDict()
        (lora_embedding_B): ParameterDict()
        (lora_magnitude_vector): ModuleDict()
      )
    )
  )
)
```
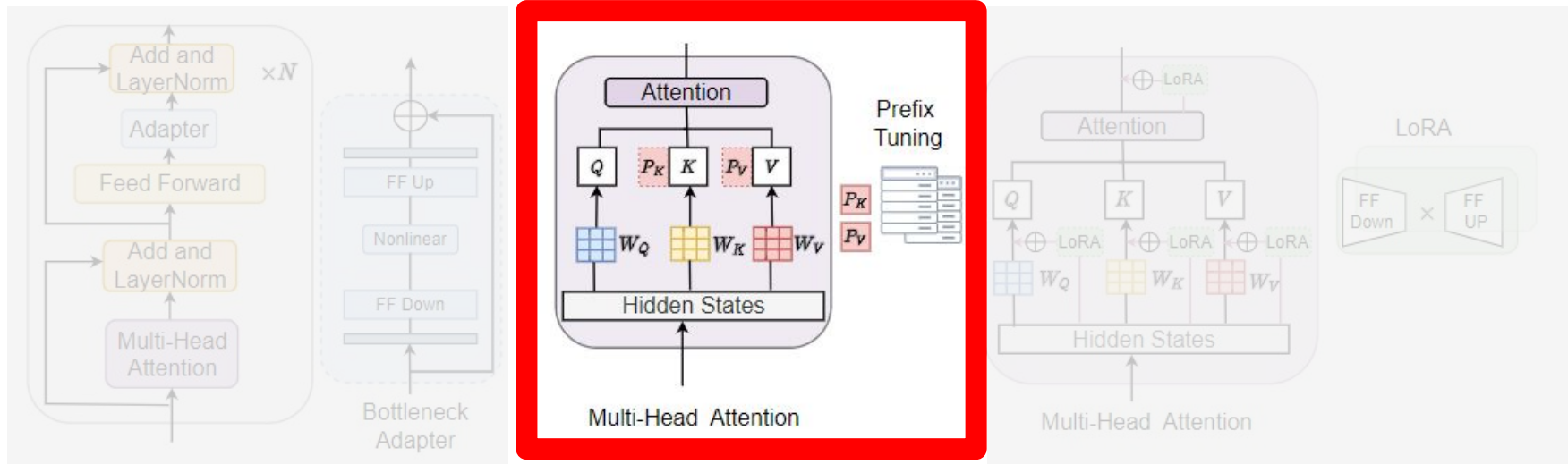
# LoRA

# Prefix Tuning



**Fig. 2**. Transformer architecture along with Adapter, Prefix Tuning, and LoRA.

# Reinforcement Learning

Machine Learning Institute

**Fine-Tuning Language Models from Human Preferences**

2019

**Learning to summarize from human feedback**

2020

**Training language models to follow instructions with human feedback**

2022

# RLHF



Figure 2: Diagram of our human feedback, reward model training, and policy training procedure.

2020

2022

Figure 1: Our training processes for reward model and policy. In the online case, the processes are interleaved.

2019

# Help



## RLHF on a Budget: GPT-2 for Summarization

Star Chen · Follow · 11 min read · Apr 7, 2025



While reading _Implementing RLHF: Learning to Summarize with trlX_, which uses GPT-J to train a reinforcement learning from human feedback (RLHF) model, I started wondering: can GPT-2 hold its own against GPT-J, but with far less compute? There's something compelling about these "old" and "micro" language models — especially how fast they are to train. But can they still produce solid results?

# Good luck!

## Feature Fiestas

Helen Zhou
Ewan Beattie
Rosh Beed
Aparna Pillai

## Overfitting Overlords

Joao Esteves
Miguel Parracho
Esperanza Shi
Andrew

## Hyperparameter Hippies

Charles Cai
Umut Sagir
Peter O'Keeffe
Ben Liong

## Kernel Kittens

Rasched Haidari
Kadriye Turkcan
Yali Pan
David Edev

## Recurrent Rebels

Jingyan Chen
Prima Gouse
Anton Dergunov
Marcin Tolysz

## Perceptron Party

Maria Sharif
Dan Goss
James Carter
Felipe Lavratti

## Dropout Disco

Ben Bethell
Tomas Krajcoviech
Nikolas Kuhn
Tao Zamorano

## Backprop Bunch

Jacob Jenner
Andrei Zhirnov
Adam Beedell
Ethan Edwards

## Gradient Gigglers

Clement Ha
Hikaru Tsujimura
James Yan
Melanie Wong

## Bayesian Buccaneers

Arjuna James
Ben Williams
Tyrone Nicholas