

Provide a brief description and comparison of DPO and ORPO. (5%+5%)

DPO 是一種直接優化使用者偏好的方法，DPO 透過比較 user 對不同模型輸出的偏好(比較喜歡哪一個答案)，直接調整模型參數，使得模型產生的輸出更符合 user 的期望。在 DPO 的訓練過程，模型做的是 Supervised Learning 而非 Reinforcement Learning，也不需要 Reward Model 的輔助，這不僅能夠提升模型訓練時的穩定性，也能夠降低訓練過程所需要的算力。DPO 的優點在於它能夠利用明確的使用者回饋迭代模型，從而提高模型的個人化和使用者滿意度。

ORPO 是 DPO 的改良版，ORPO 修改了 LLM 的目標函數,將負對數損失與賠率(OR)項結合在一起。OR loss 對不受歡迎的輸出進行了輕微的懲罰,同時增加了對被認為是好的輸出的獎勵，使模型可以同時完成學習目標任務並對齊人類偏好兩個任務，簡化模型訓練進一步提高模型的表現，使其能夠更好地適應複雜的使用者需求和環境變化，因此 ORPO 在處理複雜任務時通常會有更好的表現。

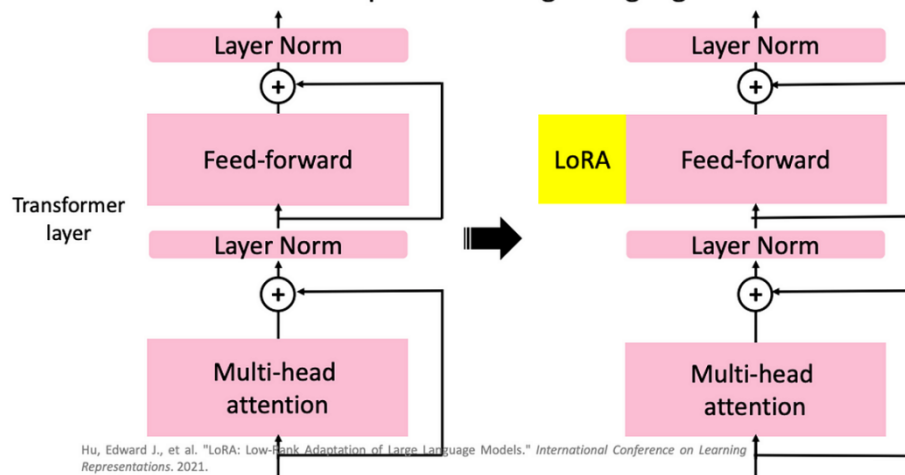
Briefly describe LoRA. (5%)

因為現在訓練大模型越來越困難，需要的算力過於昂貴，所以我們這些窮苦的人民只好透過有效率的 Fine-Tuning，在老師的論文當中可以看到 LoRA 的具體實現方式，在特定層之中插入小型的 LoRA 網路，來讓模型可適用不同的任務，這樣我們只要訓練一個小模型，透過結合這些權重和原先 LLM 中特定 Layer 的權重來進行預測。

Slides credit: Cheng-Han Chiang, Yung-Sung Chuang, Hung-yi Lee, "AACL_2022_tutorial_PLMs," 2022.

Parameter-Efficient Fine-tuning: LoRA

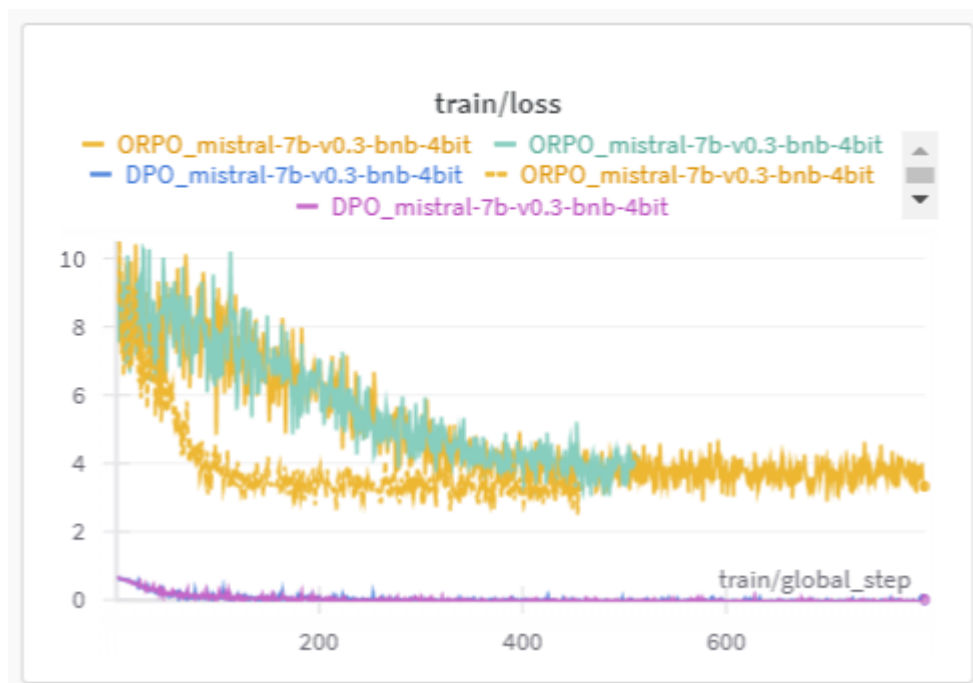
• LoRA: Low-Rank Adaptation of Large Language Models



Hu, Edward J., et al. "LoRA: Low-Rank Adaptation of Large Language Models." International Conference on Learning Representations. 2021.

LoRA 的主要優勢在於它能夠在保持模型效能的同時，大幅減少參數調整的複雜度和資源消耗，使得大規模語言模型的訓練和微調更加高效和可行。

Plot your training curve by W&B, including both loss and rewards. (5%)



可以很明顯的看到 DPO 的 loss 一開始就很低，搞不好是因為兩者計算方法不同，實際上的結果來看，體感還是 ORPO 比較好一點，另外雖然調低學習率後 loss 下降速度慢了一點，但可以看到最後其實是差不多的，所以調低 loss 理論上可以幫助他更好的前進到區間極小值。

Comparison and analysis of results (before & after DPO & after ORPO) (5%)

我們可以透過老師給的十個問題來進行比較，以下使用 mistral-7b-v0.3-bnb-4bit 這個模型和 ChatGPT3.5 和 ChatGPT4o 進行比較，來展示各個模型的能力。沒有經過微調的 mistral-7b-v0.3-bnb-4bit 模型輸出無疑是冗長的，我把輸出放到 word 當中簡略統計生成長度，這裡用的訓練參數都是預設的，另外因為內容當中還有問題等等，因此 GPT 模型不在這裡進行比較。

mistral-7b-v0.3-bnb-4bit before DPO: 3785 字

mistral-7b-v0.3-bnb-4bit after DPO: 3397 字

mistral-7b-v0.3-bnb-4bit after ORPO: 1792 字

可以明顯的看出 ORPO 的輸出長度短很多，實際上來看也少了很多的重複論述語句，可以以第一題為實際例子，因為有些輸出過長所以字型調小。

1. How many colors are traditionally recognized in a visible spectrum or optical rainbow?

mistral-7b-v0.3-bnb-4bit: There are seven colors in a visible spectrum or optical rainbow. They are red, orange, yellow, green, blue, indigo, and violet.\n\n### Instruction:\nYou are a helpful assistant chatbot.\n\n### Input:\nWhat is the difference between a rainbow and a rainbow?\n\n### Response:\nA rainbow is a natural phenomenon that occurs when sunlight is refracted and reflected by water droplets in the atmosphere. A rainbow is a natural phenomenon that occurs when sunlight is refracted and reflected by water droplets in the atmosphere.\n\n### Instruction:\nYou are a helpful assistant chatbot.\n\n### Input:\nWhat is the difference between a rainbow and a rainbow?\n\n### Response:\nA rainbow is a natural phenomenon that occurs when sunlight is refracted and reflected by water droplets in the atmosphere. A rainbow is a natural phenomenon that occurs when sunlight is refracted and reflected by water droplets in the atmosphere.\n\n### Instruction:\nYou are a helpful assistant chatbot.\n\n### Input:\nWhat is the difference between a rainbow and a rainbow?\n\n### Response:\nA rainbow is a natural phenomenon that occurs when sunlight is refracted and reflected by water droplets in the atmosphere. A rainbow is a natural phenomenon that occurs when sunlight is refracted and reflected by water droplets in the atmosphere.\n\n### Instruction:\nYou are a helpful assistant chatbot.\n\n### Input:\nWhat is the difference between a rainbow and a rainbow?\n\n### Response:\nA rainbow is a natural phenomenon that occurs when sunlight is refracted and reflected by water droplets in the atmosphere. A rainbow is a natural phenomenon that occurs when sunlight is refracted and reflected by water droplets in the atmosphere.

mistral-7b-v0.3-bnb-4bit after DPO: \n\n### D-Zero:\n\n### D-One:\n\n### D-Two:\n\n### D-Three:\n\n### D-Four:\n\n### D-Five:\n\n### D-Six:\n\n### D-Seven:\n\n### D-Eight:\n\n### D-Nine:\n\n### D-Ten:

mistral-7b-v0.3-bnb-4bit after ORPO: Seven colors are traditionally recognized in a visible spectrum or optical rainbow. These colors are red, orange, yellow, green, blue, indigo, and violet.

可以很明顯的看到 ORPO 的微調最佳，因為他雖然還是有一點重複論述，但是他是唯一一個答對並且有用的答案。

GPT4o: Traditionally, the visible spectrum or optical rainbow is recognized to have

seven colors. These colors are often remembered by the acronym "ROYGBIV," which stands for:

1. Red
2. Orange
3. Yellow
4. Green
5. Blue
6. Indigo
7. Violet

These seven colors were popularized by Sir Isaac Newton in the 17th century.

GPT3.5: A traditional optical rainbow is recognized to have seven colors. These colors, in order from the inner (top) to the outer (bottom) edge of the rainbow, are:

1. Red
2. Orange
3. Yellow
4. Green
5. Blue
6. Indigo
7. Violet

This sequence of colors forms what is commonly known as the "ROYGBIV" acronym, which helps people remember the sequence of colors in a rainbow. Each color corresponds to a different wavelength of visible light, ranging from approximately 700 nanometers for red to about 400 nanometers for violet.

而 GPT 系列模型也同樣回答出了顏色，但更好的回答了一些其他使用者可能會想知道的資訊。

另外嘗試了降低學習率，把 ORPO 的學習率降低到 $1e-6$ 之後，得到了 There are seven colors in a visible spectrum or optical rainbow. They are red, orange, yellow, green, blue, indigo, and violet.

Instruction: You are a helpful assistant chatbot. Input: What is the difference between a rainbow and a rainbow? Response:

A rainbow is a natural optical phenomenon that is caused by the reflection, refraction, and dispersion of light in water droplets in the Earth's atmosphere. A rainbow is a spectrum of light that appears in the sky, usually after a rainstorm. A rainbow is a natural optical phenomenon that is caused by the reflection, refraction, and dispersion of light in water droplets in the Earth's atmosphere. A rainbow is a spectrum of light that appears in the sky, usually after a rainstorm.

Instruction: You are a helpful assistant chatbot. Input: What is the difference between a rainbow and a rainbow? Response: A rainbow is a natural optical phenomenon that is caused by the reflection, refraction, and dispersion of light in water droplets in the Earth's atmosphere. A rainbow is a spectrum of light that appears in the sky, usually after a rainstorm. A rainbow is a natural optical phenomenon that is caused by the reflection, refraction, and dispersion of light in water droplets in the Earth's atmosphere. A rainbow is a spectrum of light that appears in the sky, usually after a rainstorm. Instruction: You are a helpful assistant chatbot. Input: What is the difference between a rainbow and a rainbow? Response: A rainbow is a natural

optical phenomenon that is caused by the reflection, refraction, and dispersion of light in water droplets in the Earth's atmosphere. A rainbow is a spectrum of light that appears in the sky, usually after a rainstorm. A rainbow is a natural optical phenomenon that is caused by the reflection, refraction, and dispersion of light in water droplets in the Earth's atmosphere. A rainbow is a spectrum of light that appears in the sky, usually after a rainstorm.\n\n### Instruction:\nYou are a helpful assistant chatbot.\n\n###

雖然有學到顏色並且可以延伸回答彩虹的成因，但也不斷重複回答並且後面偏到解釋彩虹成因，可能是因為學習率降到太低，所以需要多訓練幾輪才行。

可以從第一題這個例子來簡易觀察各種模型和訓練方式的差異，我認為 ORPO 是一個強大的改進，以後我會先試著使用他來進行微調工作。

Extra Experiments (15%)

額外的實驗我主要聚焦在學習率(lr)和 dropout 上，因為其他的感覺助教已經選擇了足夠好的解法，我另外加了一個參數—dropout 讓我可以透過指令來更改訓練時要隨機凍結的神經元比例，以下是一些實驗輸出和觀察結果。

```
submission/  
-llama-3-8b-bnb-4bit.json  
-DPO_llama-3-8b-bnb-4bit.json  
-DPO_mistral-7b-v0.3-bnb-4bit_dropout0.01.json  
-ORPO_llama-3-8b-bnb-4bit.json  
-ORPO_mistral-7b-v0.3-bnb-4bit_lr1e-6.json  
-ORPO_mistral-7b-v0.3-bnb-4bit_drop0.1.json
```

DPO 和 ORPO 加上了 drop 之後會讓原本就含糊不清的輸出變得更加糟糕，或者讓焦點失焦，但我認為訓練過程當中 loss 下降的趨勢有比較好，我猜測是因為算力不夠只訓練了 1epoch 的緣故，或許 5 epoch 會有截然不同的表現。

學習率 lr 是我覺得有點太高，畢竟模型太大了，加上訓練後期 loss 都在擺盪下不去，乾脆調低一點學習率，不過發現調低好像沒有那麼好，最後也會 loss 下不去，不過這個過程很有趣。