# Homework #6
# Reinforcement Learning with Human Feedback

1. Update
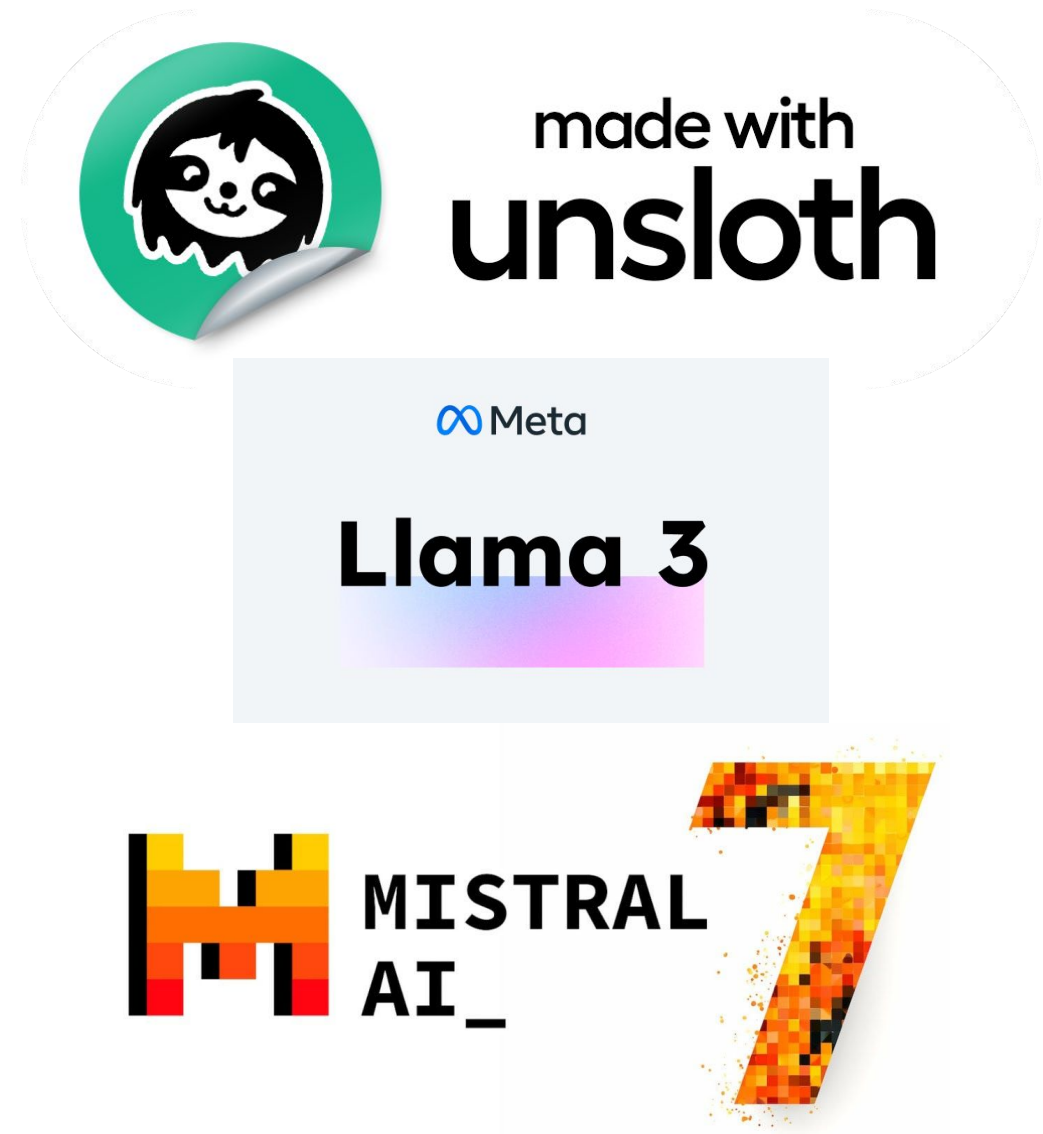2. Objective
3. Project Overview
4. unsloth Explain
5. Grading
6. Coding
7. Report
8. Submission
9. Useful Links
10. Reference

**Thanks for pointing out:**

1.  **In ORPO.py, on line 55, num_epochs should be changed to num_train_epochs (NTU COOL files have been updated).**
2.  **The uploaded AI2024-hw6 folder should include "all" the .py and .sh files along with the submission folder, and you should NOT include checkpoints, wandb folder, and outputs folder! (If you upload checkpoints, other unexpected folders, or are in the wrong format, points will be deducted!)**
3.  **Extra Experiments: We allow comparisons using models like chatGPT/Claude/Gemini/etc., but please specify the models used (e.g., GPT-4/GPT-4o/GPT-3.5/etc.). If not specified, we'll deduct points!**
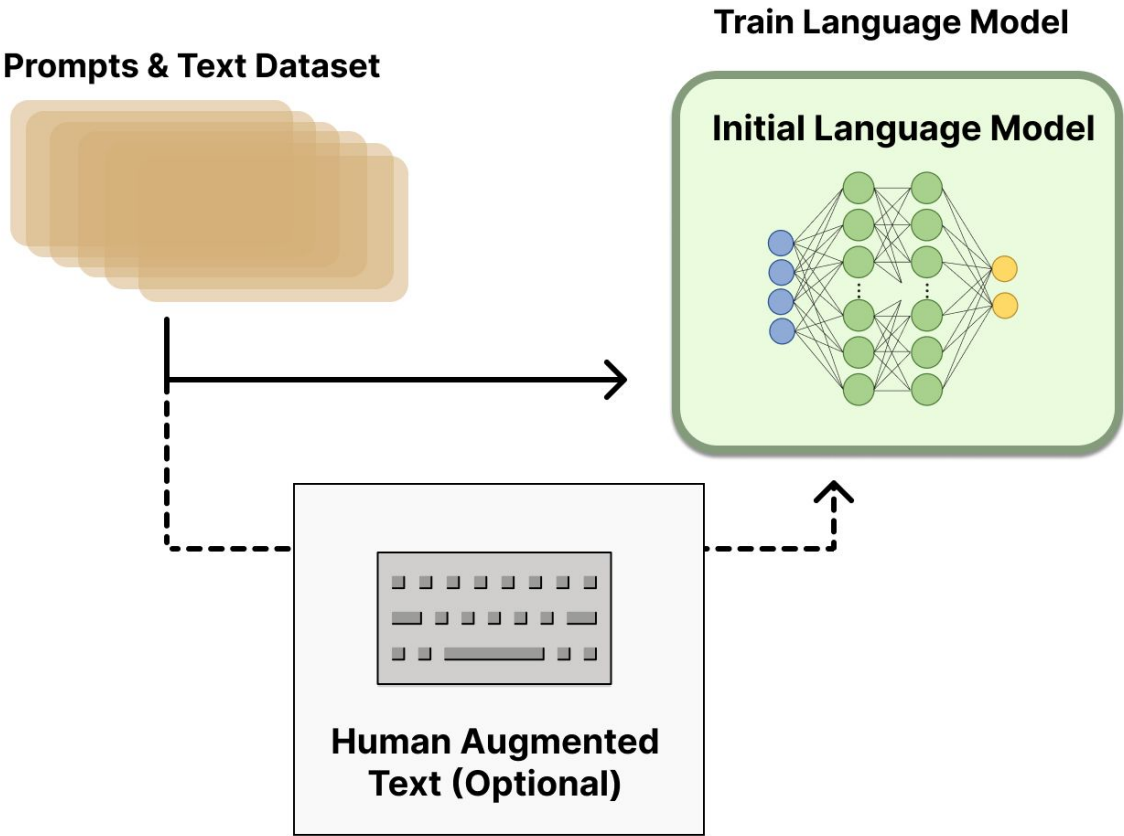
4. **5 ≥ num_train_epochs ≥ 1**, **If training is less than 1 epoch, we will deduct points!**

5. You should also upload a **README.md** within the AI2024-hw6 folder, for more details please refer to [here](#).

6. Additionally, we provide two smaller LLM options for this assignment, **tinyllama-bnb-4bit and gemma-2b-bnb-4bit**, for those suffering from the long training time. For more details, please refer to [here](#).

7.

**This assignment aims to utilize the open-source [unsloth](#) framework from unslothAI to fine-tune large language models (LLMs) using two reinforcement learning with human feedback (RLHF) techniques, as well as parameter-efficient fine-tuning (PEFT) with LoRA.**
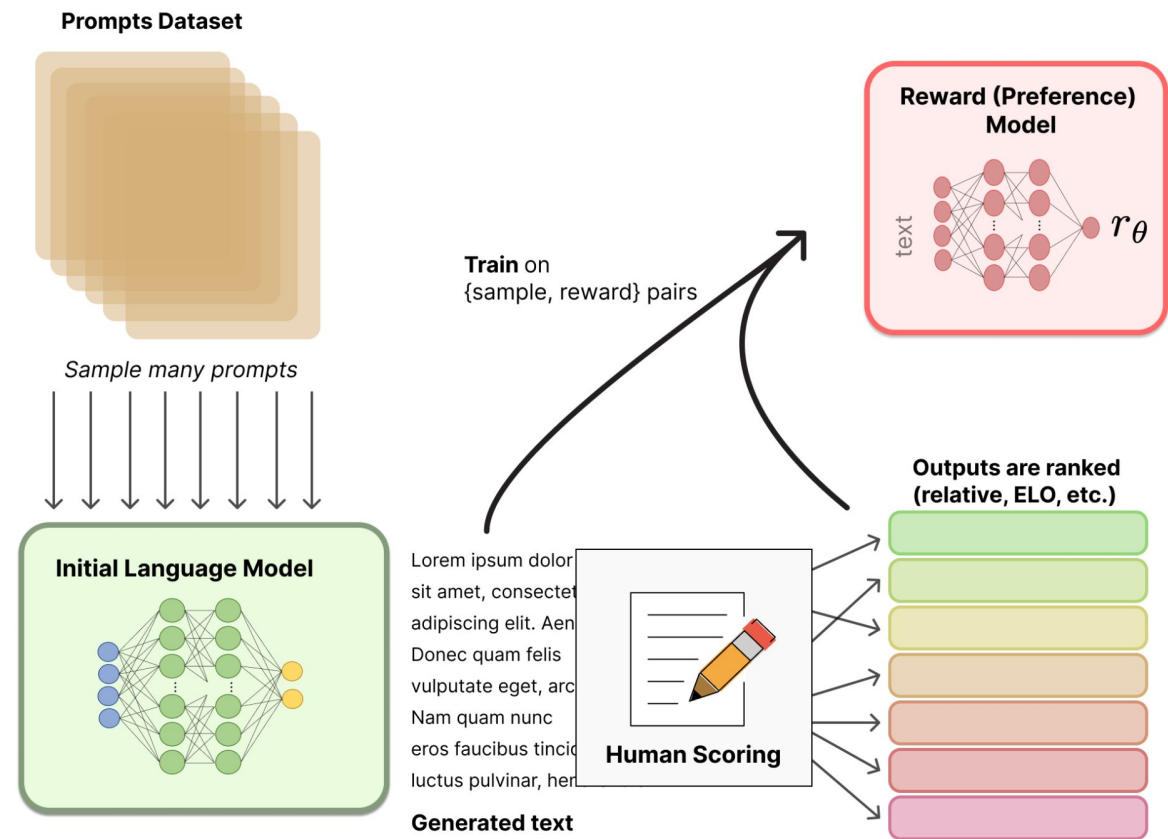
**Reinforcement Learning with Human Feedback (RLHF)**

**Step 1:**
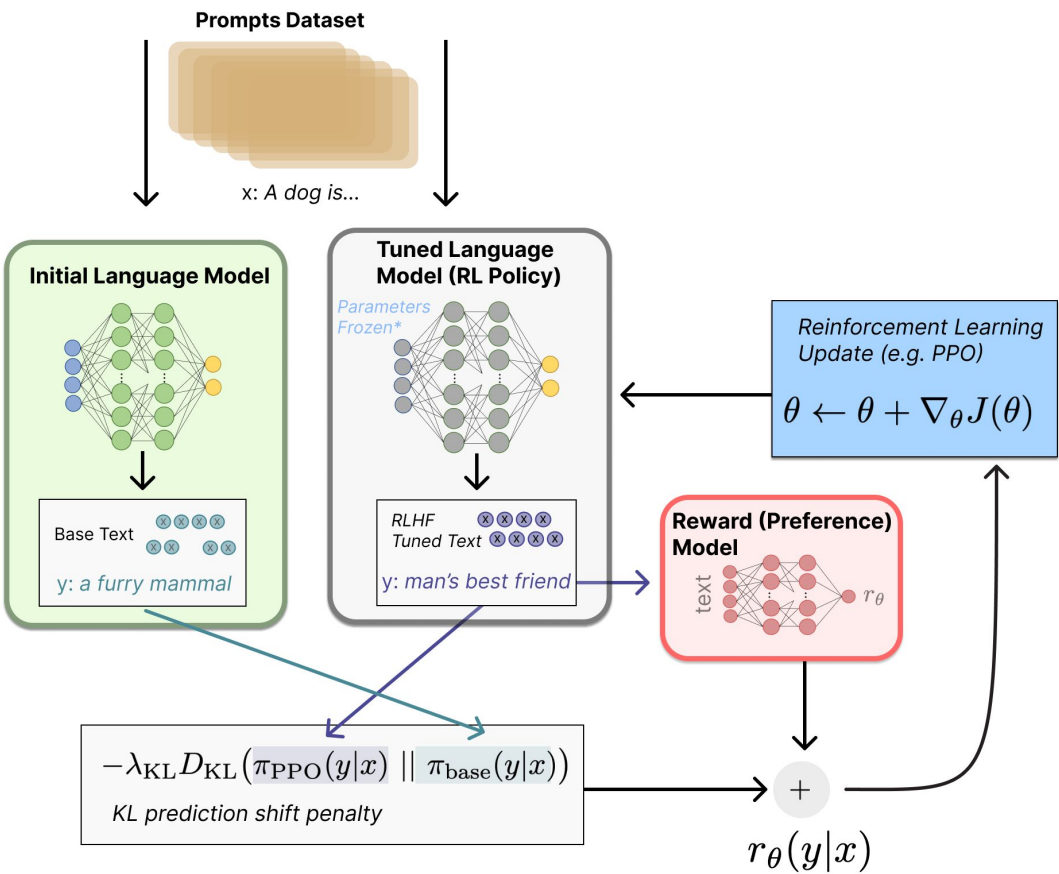**Supervised Fine-tuning (SFT)**



Ref

**Reinforcement Learning with Human Feedback (RLHF)**

**Step 2: Reward model training**

## Reinforcement Learning with Human Feedback (RLHF)

## Step 3: Fine-tuning with RL

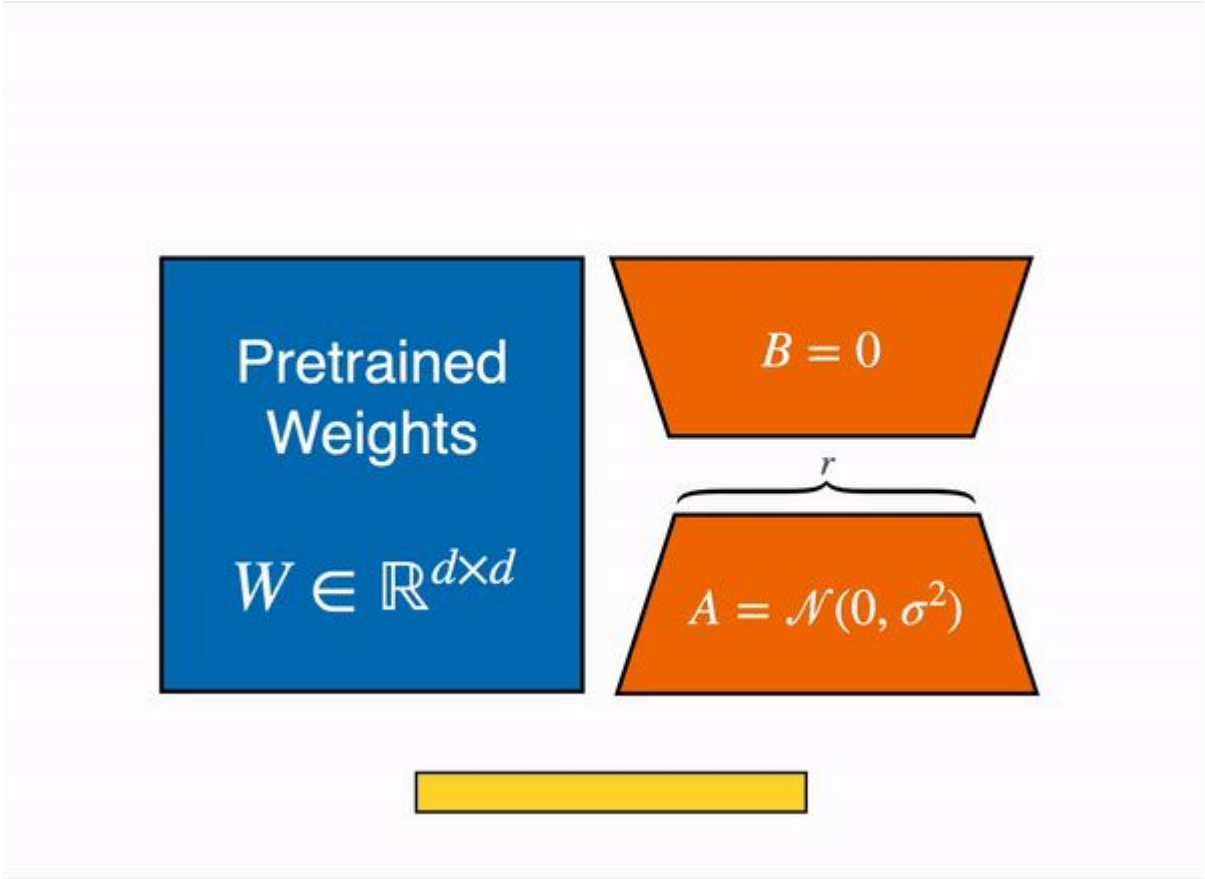**Reinforcement Learning with Human Feedback (RLHF)**

**Drawbacks:**
1. Requires an additional reward model.
2. RL training is unstable and difficult to tune hyperparameters.

**Therefore, we will implement two simplified methods, Direct Preference Optimization (DPO) and Odds Ratio Preference Optimization (ORPO), to fine-tune large language models (LLMs).**

Ref

**Q: How to fine-tune LLMs on consumer GPU?**

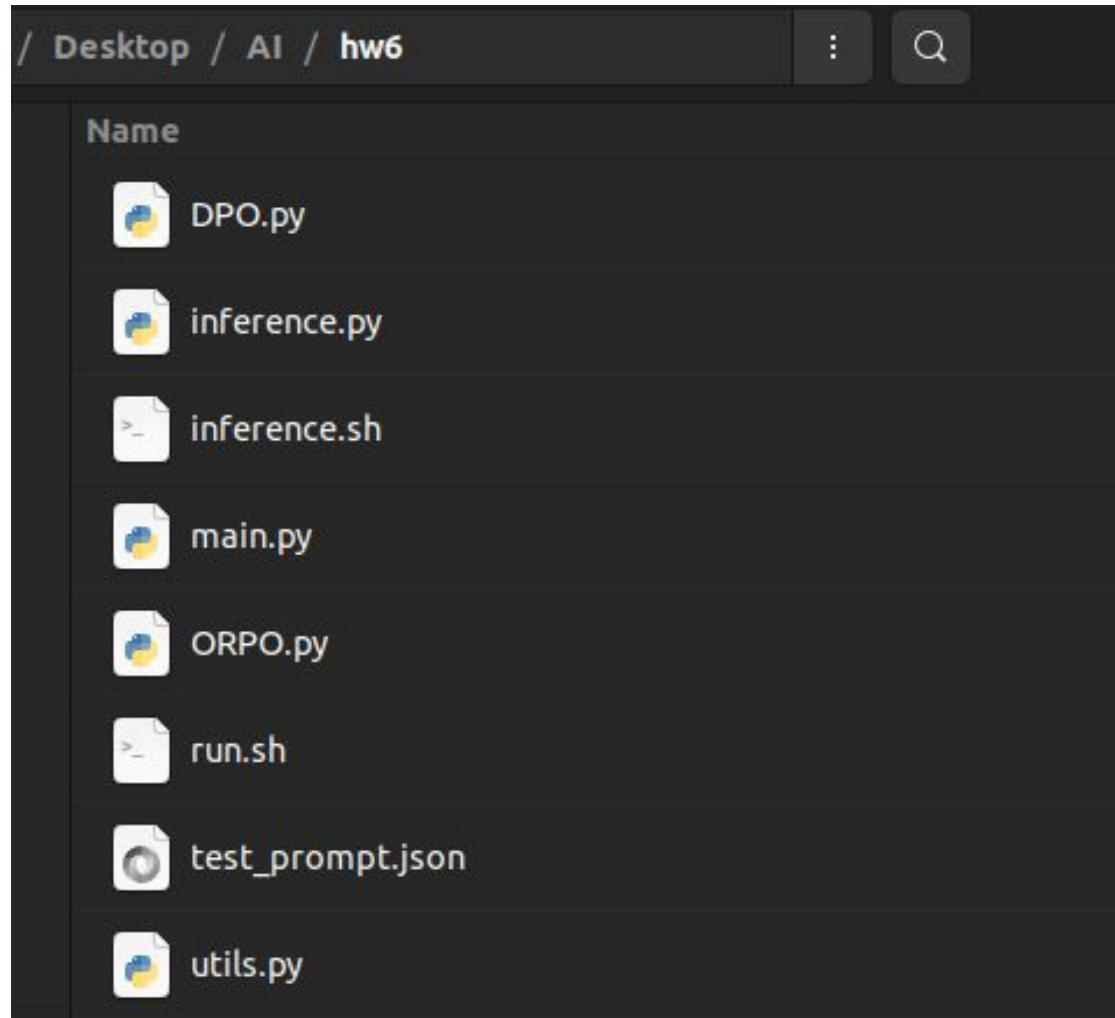**A: Low Rank Adaptation (LoRA)**

❖ The code for this project consists of several Python files, **which you will need to read and understand in order to complete the assignment.**

❖ You will fill in portions of **DPO.py, and ORPO.py** during the assignment.

**Do not change the provided dataset, and using other datasets for data augmentation is not permitted!**

**Training dataset: Intel/orca_dpo_pairs**

⊞ **Dataset Viewer**

Split (1)
train · 12.9k rows

🔍 Search this dataset

| system<br>string · classes<br><br>17 values | question<br>string · lengths<br><br>22          8.05k | chosen<br>string · lengths<br><br>1          4.95k | rejected<br>string · lengths<br><br>5          7.95k |
|---|---|---|---|
| | You will be given a definition of a task… | [ ["AFC Ajax (amateurs)", "has ground", "Sportpark De Toekomst"], ["Ajax Youth Academy", "plays at",… | Sure, I'd be happy to help! Here are the RDF… |
| You are an AI assistant. You will be given a task… | Generate an approximately fifteen-word sentence… | Midsummer House is a moderately priced Chinese restaurant with a 3/5 customer rating, located… | Sure! Here's a sentence that describes all the… |
| You are a helpful assistant, who always… | What happens next in this paragraph? She then rubs… | C. She then dips the needle in ink and using the pencil to draw a design on her leg, rubbing it of… | Ooh, let me think! *giggle* Okay, I know… |
| You are an AI assistant. You will be given a task… | Please answer the following question: I… | Based on the passage, discuss the primary motivations and outcomes of the 1901 Federation o… | Certainly! Here's a detailed and long answer… |
| You are an AI assistant that helps people find… | James runs a TV show and there are 5 main… | James pays the minor characters $15,000 each episode. Since there are 4 minor characters, he… | Sure, I'd be happy to help! To calculate how… |
| You are an AI assistant that helps people find… | Given the stream of consciousness rationale,… | Question: What is the proper technique for a female beach volleyball player to serve the ball… | Sure, I'd be happy to help! Here's a reasonabl… |

**Testing Dataset: test_prompt.json (released on NTU COOL)**

```json
[
    {
        "id": 1,
        "prompt": "How many colors are traditionally recognized in a visible spectrum or optical rainbow?"
    },
    {
        "id": 2,
        "prompt": "In a basket, there are 20 oranges, 60 apples, and 40 bananas. If 15 pears were added, and half of the ora
    },
    {
        "id": 3,
        "prompt": "If you were a car salesperson, how would you convince a potential buyer to purchase an electric vehicle?"
    },
    {
        "id": 4,
        "prompt": "Break down the process of photosynthesis into a bullet-pointed list, detailing each stage and the overall
    },
    {
        "id": 5,
        "prompt": "Explain why college students should get a library card."
    },
```

# Project Overview (Cont.)

**LLMs:**

**Required (choose one out of four):**
1. "unsloth/mistral-7b-v0.3-bnb-4bit"
2. "unsloth/llama-3-8b-bnb-4bit"
3. "unsloth/tinyllama-bnb-4bit"
4. "unsloth/gemma-2b-bnb-4bit"

**Additional:**
1. "unsloth/mistral-7b-instruct-v0.3-bnb-4bit"
2. "unsloth/llama-3-8b-Instruct-bnb-4bit"
3. …

**You should choose one of the two required models based on your GPU! However, additional experiments of other models with detailed comparisons are permitted for extra experiments!**

If you are unable to successfully create the virtual environment, it is recommended that you use Colab instead!
**P.S. TAs will not assist with troubleshooting environment-related issues.**

**Installation Instructions:**

1. **conda create -y -n ai_hw6 python=3.10**
2. **conda activate ai_hw6**
3. **install [pytorch](#) based on your cuda version**
4. **pip install --no-deps trl peft accelerate bitsandbytes**
5. **pip install tqdm packaging wandb**
6. **based on cuda version install the correct version of unsloth (detailed information is [here](#))**
7. **Verify if the installation was successful.:**
   a. **nvcc**
   b. **python -m xformers.info**
   c. **python -m bitsandbytes**

- **Coding (60%)**
  - DPO.py (20%)
  - ORPO.py (20%)
  - Submission folder (10%)
  - README.md (10%)

- **Report (40%)**
  - Must be submitted in PDF format only.
  - The score for each problem is detailed in [here](here).

■ **DPO & ORPO (20%+20%)**

- Complete the implementation of DPO and ORPO by filling in the blanks.

```
# Model
# model, tokenizer = FastLanguageModel.from_pretrained(model_name=args.model_name,...)
utils.YOUR_CODE_HERE
```
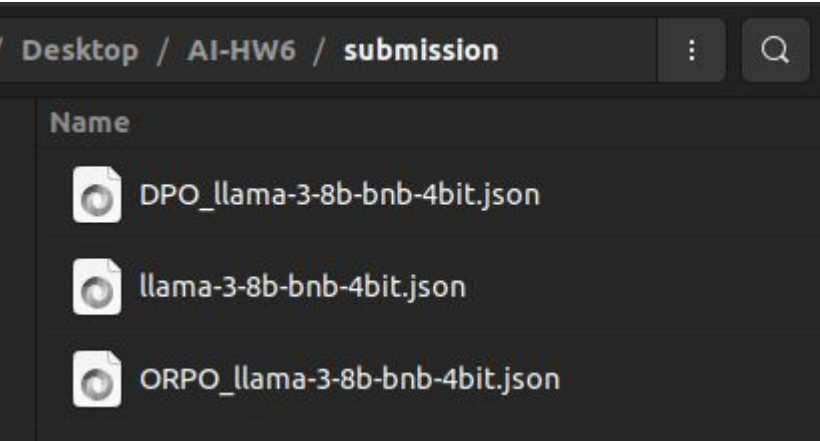
- Train and inference the fine-tuned models.

- To receive full points, your submitted code should run successfully without any errors.

  (Only one minor modification is allowed for resubmission, but will be subject to a 5% penalty.)

- No partial points will be given.

  *Grading: bash run.sh exp_name model_name wandb_token (and other hyperparameters)*

■ **Submission folder (10%)**

- There should be at least **three** JSON files:

  i. the first should contain the generated text from the selected LLM (2%),
     *Note: bash inference.sh model_name wandb_token*

  ii. the second should contain the results after DPO (4%),

  iii. and the third should contain the output from ORPO (4%).

- *Note: submission/*
  *-llama-3-8b-bnb-4bit.json*
  *-DPO_llama-3-8b-bnb-4bit.json*
  *-ORPO_llama-3-8b-bnb-4bit.json*



(If you have done extra experiments, you **must** also include those files from different hyperparameters or models in the submission directory!)

■ **README.md (10%)**

- In the README.md, you should provide following information:

    i. commands for creating the virtual env, thus **requirements.txt is required if specified in the README.md!**

    ii. how to train your model, including "**all**" the hyperparameters you used should be specified. (feel free to modify the run.sh)

    (**We'll follow the steps provided in your README.md for reproducibility. If we can not successfully run your code, we will then deduct points! Thus, make sure your README.md is detailed enough before submission!**)

- **Report (25%)**
  - Provide a brief description and comparison of DPO and ORPO. (5%+5%)
  - Briefly describe LoRA. (5%)
  - Plot your training curve by W&B, including both loss and rewards. (5%)
  - Comparison and analysis of results (before & after DPO & after ORPO) (5%)

- **Extra Experiments (15%) –** <u>**We will assign grades based on the richness of the experiments, with different score regarding to varying levels of detail and complexity.**</u>

  - Comparison and analysis of various hyperparameters (e.g., num_epochs, beta, etc.)

  - Comparison and analysis of various models, which GPT-4, GPT-4o, GPT-3.5, and etc. are allowed, for the variants of the required models please refer to [here](#).

  (If you have done extra experiments, you **must** also include those files from different hyperparameters or models in the submission directory!)

- Deadline: 2024/06/12 (Wed.) 23:59

- Zip all files as **hw6_<student_id>.zip**

- Submit to NTU COOL

- Your submission should include the following files:

  - hw6_<student_id>.pdf

  - AI2024-hw6

- **Do not put report.pdf into AI2024-hw6 folder**

- *Note: hw6_<student_id> is an example format. For instance, if your student ID is r1234567, then the file name should be hw6_r1234567 or hw6_R1234567.*

■ The discount for late submission of assignments is as follows:

| <24hr | <48hr | <72hr |
|-------|-------|-------|
| 70%   | 50%   | 30%   |

RLHF:

https://www.youtube.com/watch?v=v12IKvF6Cj8

https://speech.ee.ntu.edu.tw/~hylee/genai/2024-spring-course-data/04 12/0412_LLMtraining_part3.pdf

W&B:

https://docs.wandb.ai/quickstart

https://docs.wandb.ai/tutorials

https://huggingface.co/blog/rlhf

https://huggingface.co/blog/trl-peft

https://hackmd.io/@YungHuiHsu/Sy5Ug7iV6

# Any Question

ai.ta.2024.spring@gmail.com