

Adversarial Methods for the alignment of CodeLLMs

End of semester presentation

Ilyas Oulkadda Julien Perez

July 2024

Table of Contents

- 1 Introduction
- 2 Memory Efficient Training
- 3 PPO unit testing
- 4 Iterative adversarial DPO
- 5 What's next
- 6 Conclusion

Table of Contents

- 1 Introduction
- 2 Memory Efficient Training
- 3 PPO unit testing
- 4 Iterative adversarial DPO
- 5 What's next
- 6 Conclusion

Introduction

Goal

Improve CodeLLMs alignment (Ouyang et al., 2022)

How ?

- Adversarial (OpenAI et al., 2021)
- Self-play (Sukhbaatar et al., 2018)
- Curriculum learning (Sukhbaatar et al., 2018)
- At scale (Bowman et al., 2022)

Current open models

Average Score Vs Throughput (A100-80GB, Float16, Batch Size 1)

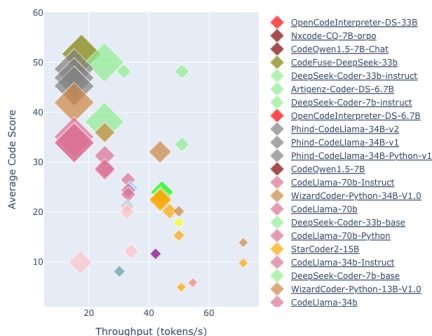


Figure: The open code LLM benchmark (Ben Allal, Muennighoff, et al., 2022)

Table of Contents

- 1 Introduction
- 2 Memory Efficient Training**
- 3 PPO unit testing
- 4 Iterative adversarial DPO
- 5 What's next
- 6 Conclusion

Why ?

- For the planet
- We have limited resources

Low Rank Adaptation

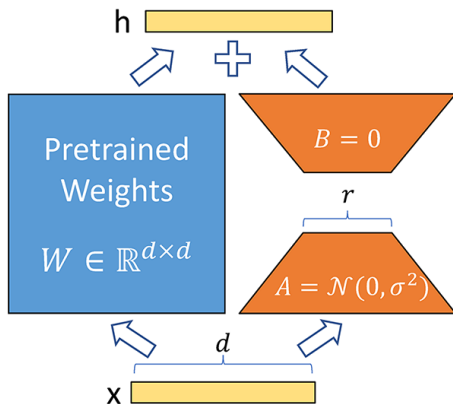


Figure: LoRA adds new adapters, the weight matrix is a product of 2 lower rank matrices (Hu et al., 2021)

Quantization

Table: Different quantization for Phi 1.5 (Gunasekar et al., 2023), 1.3B parameters

dtype	Model	Grad. Calc.	Backward Pass	Opt. Step
Float32	4.9 GB	4.9 GB	9.81 GB	19.62 GB
Float16/BF16	4.9 GB	7.36 GB	9.81 GB	9.81 GB

Example with PHI-1.5

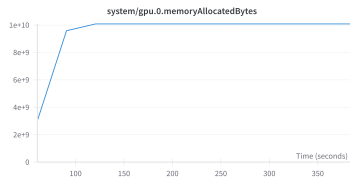


Figure: Memory usage while training Phi 1.5 on HumanEval (Chen et al., 2021)

Table of Contents

- 1 Introduction
- 2 Memory Efficient Training
- 3 PPO unit testing**
- 4 Iterative adversarial DPO
- 5 What's next
- 6 Conclusion

Proximal Policy Optimization (Schulman et al., 2017)

$$r_t(\theta) = \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)}, \quad (1)$$

$$L^{\text{CLIP}}(\theta) = \hat{\mathbb{E}}_t \left[\min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right] \quad (2)$$

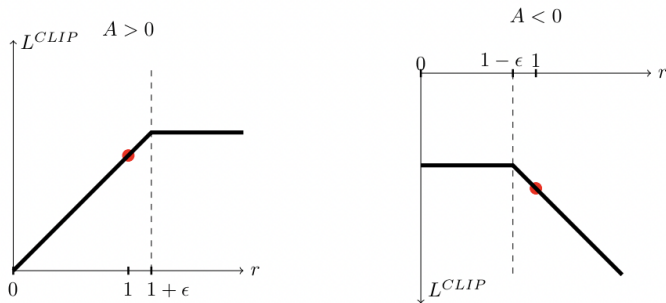


Figure: Objective function for both negative advantages and positive advantages.

Instability of PPO during RLHF (Ouyang et al., 2022)

Reward Model

The reward model is an approximation of the human preference.

Solution

Use unit tests as our reward model

Dataset generation

Main challenge

Generate a diverse dataset

Method (Ben Allal, Lozhkov, et al., 2024) (Gunasekar et al., 2023) (Eldan and Li, 2023)

- List of coding categories
- Generate sub-categories
- Generate multiple exercises
- Random profession per exercises
- Generate the unit tests
- JSON format

Sample

```
1 def split_constellations(name):
2     """
3     Given a string representing a comma-separated list of constellations, split it into a list.
4     """
5
6
7 def test_split_constellations():
8     # Testing different scenarios based on the function description
9     assert split_constellations("Andromeda, Perseus, Triangulum") == [
10         "Andromeda",
11         "Perseus",
12         "Triangulum",
13     ]
14     assert split_constellations("The Big Dipper, Ursa Major") == [
15         "The Big Dipper",
16         "Ursa Major",
17     ]
18     assert split_constellations("The Southern Cross") == ["The Southern Cross"]
19     assert split_constellations("") == []
20     assert split_constellations(", , ,") == []
21
22
23 test_split_constellations()
```

Figure: Sample from our generated dataset, composed by a prompt (function definition, arguments and docstring), a test suite and a the name of the test suite function

```
1 def gaussian_elimination(A):
2     """
3     Find the determinant and inverse of a given 3x3 matrix A using Gaussian elimination.
4
5     Args:
6     A (list of lists): A 3x3 matrix represented as a list of lists.
7
8     Returns:
9     Tuple: A tuple containing the determinant and the inverse of the matrix A.
10    """
11    # Placeholder for the Gaussian elimination implementation
12    pass
13
14 def test_gaussian_elimination():
15    # Test 1: Identity matrix
16    A = [[1, 0, 0],
17         [0, 1, 0],
18         [0, 0, 1]]
19    assert gaussian_elimination(A) == (1, A)
20
21    # Test 2: Simple 3x3 matrix
22    A = [[2, 3, 4],
23         [1, 5, 7],
24         [0, 6, 8]]
25    assert gaussian_elimination(A) == (36, [[-3.5, 2, 1], [1, -0.5, -1], [1.5, -3, 0]])
26
27    # Test 3: Zero matrix
28    A = [[0, 0, 0],
29         [0, 0, 0],
30         [0, 0, 0]]
31    assert gaussian_elimination(A) == (0, A)
32
33    # Test 4: Matrix with a determinant of zero
34    A = [[1, 2, 3],
35         [4, 5, 6],
36         [7, 8, 9]]
37    assert gaussian_elimination(A) == (0, None)
38
39    # Test 5: Matrix with complex numbers
40    A = [[1 + 2j, 3, 4 - 5j],
41         [2, 5 + 6j, 7],
42         [7, 8 - 9j, 9]]
43    assert gaussian_elimination(A) == (4 + 13j,
44                                       [[(14 - 13j) / 25, -3 / 5, (11 + 24j) / 25],
45                                        [2, -5, 0],
46                                        [(11 - 12j) / 25, (4 + 3j) / 5, (4 + 3j) / 25]])
47
48    test_gaussian_elimination()
```

Figure: Same here but a little bit more complex

Reward function

Our reward model is actually just a function that compiles the code and runs the tests.

Rewards

- **Timeout:** -2
- **Syntax error:** -1
- **Runtime error:** -0.6
- **Assertion error (test fails):** -0.3
- **Pass all unit tests:** 4

Training

Table: Deepseeker-Coder-1b-base (Guo et al., 2024) performances (Ben Allal, Muennighoff, et al., 2022)

Win Rate	humaneval-python	java	javascript	cpp
16	32.13	27.16	28.46	27.96

Main issues

- Missing libraries
- Indentation
- Timeouts and inputs
- Wrong tests

Results

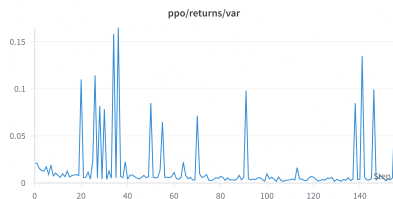


Figure: The observed high variance in rewards may be attributed to the varying difficulty levels of coding exercises.

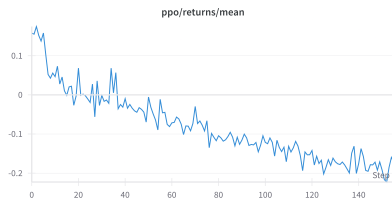


Figure: The mean reward does not show a consistent increase, indicating that the current PPO setup may not be maximizing expected rewards efficiently.

Limitations

- Requires a high quality test suite
- Larger models
- Only tests functionality
- Limited to a restrained type of exercises

Limitations

```
1 def tax_reform_simulation(initial_gdp, tax_rates, years):
2     """
3     Docstring explaining the exercise:
4     Create a function to simulate the economic impact of a tax reform over a certain number of years.
5
6     Args:
7         initial_gdp (float): Initial Gross Domestic Product (GDP) of the economy
8         tax_rates (list): a list of percentage tax rates for each income level
9         years (int): number of years to simulate the tax reform
10
11     Returns:
12         None, prints out the GDP for each year after the tax reform
13     """
14     # Placeholder for function implementation
15     pass
16
17 def test_tax_reform_simulation():
18     # Simple case: constant tax rate, positive GDP growth
19     assert tax_reform_simulation(1000, [20], 3) == [1020.0, 1040.4, 1061.204]
20
21     # Simple case: no tax, positive GDP growth
22     assert tax_reform_simulation(1000, [0], 2) == [1020.0, 1040.4]
23
24     # Simple case: negative tax rate (subsidy), positive GDP growth
25     assert tax_reform_simulation(1000, [-10], 1) == [1030.0]
26
27     # Complex case: varying tax rates, positive and negative GDP growth
28     assert tax_reform_simulation(1000, [10, 15, 5, -5], 5) == [1015.0, 1030.25, 1045.7625, 1033.47625, 1065.147625]
29
30     # Edge case: zero initial GDP, varying tax rates
31     assert tax_reform_simulation(0, [20, 15, 5], 3) == [0, 0, 0]
32
33 if __name__ == "__main__":
34     test_tax_reform_simulation()
```

Figure: Function returns nothing while test suite expects a result

Table of Contents

- 1 Introduction
- 2 Memory Efficient Training
- 3 PPO unit testing
- 4 Iterative adversarial DPO**
- 5 What's next
- 6 Conclusion

Direct Preference Optimization (Rafailov et al., 2023)

$$-\mathbb{E}_{(x, y_u, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_u | x)}{\pi_{\text{ref}}(y_u | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right] \quad (3)$$

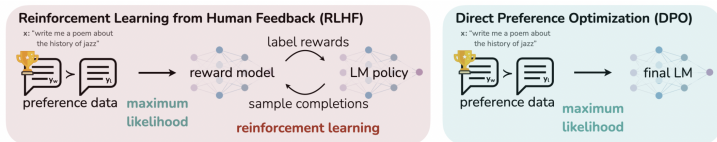


Figure: PPO RLHF VS DPO (Rafailov et al., 2023)

Adversarial game

Goal

Transfer coding ability from a larger model to a smaller one

Setup

- An oracle model
- A student model
- Oracle generates exercises and solutions
- Student also generates it's own solution

Datasets

Requirements

For DPO, datasets must have a Prompt, Chosen and Rejected

- Llama 3 70B on Groq as an oracle (AI@Meta, 2024)
- Mistral 7B as a student (Jiang et al., 2024)
- Same method as for PPO, categories, sub-categories...

Adversarial dataset

At the end of each training, we retrieve the hardest exercise to use as templates for more samples. (Sukhbaatar et al., 2018)

Training setup

- Trained on H100 80GB
- HuggingFace DPO
- LoRA
- Float16 Quantization
- Single adapter

Results



Figure: Chosen rewards during our training, the higher the reward the more our model converge to the chosen distribution

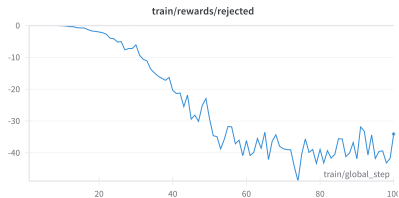


Figure: Rejected rewards during our training, the lower the reward the more our model is diverging from the rejected distribution

Benchmark

After 2 iterations, we were able to test the model on the HumanEval benchmark showing a promising 1% increase.

Results



Figure: Increasing reward margin means our model is actually converging toward the Oracle's distribution and diverging from the initial distributions



Figure: Accuracies plot shows the frequency of choosing the preferred answer. The trained model quickly reaches a perfect accuracy score, which is a good sign but could also mean that the difference between preferred and rejected answers is too obvious.

Conclusion

Advantages

- Learning efficiently
- Easier than PPO unit tests
- Covers large coding aspect
- Scalable oversight

Limitations

- Dependent on the Oracle
- Monolithic Oracle

Table of Contents

- 1 Introduction
- 2 Memory Efficient Training
- 3 PPO unit testing
- 4 Iterative adversarial DPO
- 5 What's next**
- 6 Conclusion

What's next

- Run on Jean Zay or future LRE GPUs
- Larger and better dataset
- Data decontamination
- Train Oracles
- Improve Adversarial sampling








Table of Contents

- 1 Introduction
- 2 Memory Efficient Training
- 3 PPO unit testing
- 4 Iterative adversarial DPO
- 5 What's next
- 6 Conclusion**








Conclusion

- Promising results for adversarial DPO
- Limitations during PPO
- PPO failed but we have interesting behavior results
- Dataset generation can be easily scaled now for future use

Bibliography I

-  AI@Meta (2024). “Llama 3 Model Card”. In: URL: https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
-  Ben Allal, Loubna, Anton Lozhkov, et al. (2024). *Cosmopedia*. URL: <https://huggingface.co/datasets/HuggingFaceTB/cosmopedia>.
-  Ben Allal, Loubna, Niklas Muennighoff, et al. (2022). *A framework for the evaluation of code generation models*. <https://github.com/bigcode-project/bigcode-evaluation-harness>.
-  Bowman, Samuel R. et al. (2022). *Measuring Progress on Scalable Oversight for Large Language Models*. arXiv: 2211.03540 [cs.CL].
-  Chen, Mark et al. (2021). *Evaluating Large Language Models Trained on Code*. arXiv: 2107.03374 [cs.LG].
-  Eldan, Ronen and Yuanzhi Li (2023). *TinyStories: How Small Can Language Models Be and Still Speak Coherent English?* arXiv: 2305.07759 [cs.CL].
-  Gunasekar, Suriya et al. (2023). *Textbooks Are All You Need*. arXiv: 2306.11644 [cs.CL].

Bibliography II

-  Guo, Daya et al. (2024). *DeepSeek-Coder: When the Large Language Model Meets Programming – The Rise of Code Intelligence*. arXiv: 2401.14196 [cs.SE]. URL: <https://arxiv.org/abs/2401.14196>.
-  Hu, Edward J. et al. (2021). *LoRA: Low-Rank Adaptation of Large Language Models*. arXiv: 2106.09685 [cs.CL].
-  Jiang, Albert Q. et al. (2024). *Mixtral of Experts*. arXiv: 2401.04088 [cs.CL].
-  OpenAI, OpenAI et al. (2021). *Asymmetric self-play for automatic goal discovery in robotic manipulation*. arXiv: 2101.04882 [cs.LG].
-  Ouyang, Long et al. (2022). *Training language models to follow instructions with human feedback*. arXiv: 2203.02155 [cs.CL].
-  Rafailov, Rafael et al. (2023). *Direct Preference Optimization: Your Language Model is Secretly a Reward Model*. arXiv: 2305.18290 [cs.CL].
-  Schulman, John et al. (2017). *Proximal Policy Optimization Algorithms*. arXiv: 1707.06347 [cs.CL].

Bibliography III



Sukhbaatar, Sainbayar et al. (2018). *Intrinsic Motivation and Automatic Curricula via Asymmetric Self-Play*. ICLR 2018. arXiv: 1703.05407 [cs.LG].