

Prompting Capabilities

- Note, you can try any of these prompts outside of this classroom, and without coding, by going to the chat interface [Le Chat \(https://chat.mistral.ai/chat\)](https://chat.mistral.ai/chat).
 - You can sign up with a free account.
 - Signing up for an account is **not** required to complete this course.

```
In [1]: # !pip install mistralai
```

- Notice that it's "mistralai", and not "mistral"

Load API key and helper function

- Note: You can view or download the helper.py file by clicking on the "Jupyter" logo to access the file directory.

```
In [2]: from helper import load_mistral_api_key  
load_mistral_api_key()
```

```
In [3]: from helper import mistral  
mistral("hello, what can you do?")
```

'Hello! I can assist you with a variety of tasks. I can answer questions, provide information, help manage your schedule, set reminders, send messages, and much more. I can also help explain concepts, provide explanations for topics, and even help with homework or studying. How can I assist you today?'

Classification

```
In [4]: prompt = """
        You are a bank customer service bot.
        Your task is to assess customer intent and categorize customer
        inquiry after <<<>>> into one of the following predefined categories:

        card arrival
        change pin
        exchange rate
        country support
        cancel transfer
        charge dispute

        If the text doesn't fit into any of the above categories,
        classify it as:
        customer service

        You will only respond with the predefined category.
        Do not provide explanations or notes.

        ###
        Here are some examples:

        Inquiry: How do I know if I will get my card, or if it is 1
        Category: card arrival
        Inquiry: I am planning an international trip to Paris and want
        Category: exchange rate
        Inquiry: What countries are getting support? I will be traveling
        Category: country support
        Inquiry: Can I get help starting my computer? I am having trouble
        Category: customer service
        ###

        <<<
        Inquiry: {inquiry}
        >>>
        Category:
        """
```

Ask Mistral to check the spelling and grammar of your prompt

```
In [5]: response = mistral(f"Please correct the spelling and grammar of
        this prompt and return a text that is the same prompt, \
        with the spelling and grammar fixed: {prompt}")
```

```
In [6]: print(response)
```

You are a bank customer service bot.
Your task is to assess customer intent and categorize the customer inquiry following the inquiry into one of the following predefined categories:

card arrival
change PIN
exchange rate
country support
cancel transfer
charge dispute

If the text does not fit into any of the above categories, classify it as:
customer service

You will only respond with the predefined category. Do not provide explanations or notes.

###

Here are some examples:

Inquiry: How do I know if I will get my card, or if it is lost? I am concerned about the delivery process and would like to ensure that I will receive my card as expected. Could you please provide information about the tracking process for my card, or confirm if there are any indicators to identify if the card has been lost during delivery?

Category: card arrival

Inquiry: I am planning an international trip to Paris and would like to inquire about the current exchange rates for Euros as well as any associated fees for foreign transactions.

Category: exchange rate

Inquiry: What countries are getting support? I will be traveling and living abroad for an extended period of time, specifically in France and Germany, and would appreciate any information regarding compatibility and functionality in these regions.

Category: country support

Inquiry: Can I get help starting my computer? I am having difficulty starting my computer, and would appreciate your expertise in helping me troubleshoot the issue.

Category: customer service

###

<<<

Inquiry: {inquiry}

>>>

Category:

Try out the model

```
In [7]: mistral(
        response.format(
            inquiry="I am inquiring about the availability of your
        )
    )
```

'country support'

Information Extraction with JSON Mode

```
In [8]: medical_notes = """
A 60-year-old male patient, Mr. Johnson, presented with symptoms
of increased thirst, frequent urination, fatigue, and unexplained
weight loss. Upon evaluation, he was diagnosed with diabetes,
confirmed by elevated blood sugar levels. Mr. Johnson's weight
is 210 lbs. He has been prescribed Metformin to be taken twice
with meals. It was noted during the consultation that the patient
is a current smoker.
"""
```

```
In [9]: prompt = f"""
Extract information from the following medical notes:
{medical_notes}

Return json format with the following JSON schema:

{{
    "age": {{
        "type": "integer"
    }},
    "gender": {{
        "type": "string",
        "enum": ["male", "female", "other"]
    }},
    "diagnosis": {{
        "type": "string",
        "enum": ["migraine", "diabetes", "arthritis", "acne"]
    }},
    "weight": {{
        "type": "integer"
    }},
    "smoking": {{
        "type": "string",
        "enum": ["yes", "no"]
    }}
}}
```

```
In [10]: response = mistral(prompt, is_json=True)
print(response)
```

```
{"age": 60, "gender": "male", "diagnosis": "diabetes", "weight": 210, "smoking": "yes"}
```

Personalization

```
In [11]: email = """
Dear mortgage lender,

What's your 30-year fixed-rate APR, how is it compared to the 1
fixed rate?

Regards,
Anna
"""
```

```
In [12]: prompt = f"""

You are a mortgage lender customer service bot, and your task is
create personalized email responses to address customer questions.
Answer the customer's inquiry using the provided facts below. Ensure
that your response is clear, concise, and directly addresses the
customer's question. Address the customer in a friendly and
professional manner. Sign the email with "Lender Customer Support".

# Facts
30-year fixed-rate: interest rate 6.403%, APR 6.484%
20-year fixed-rate: interest rate 6.329%, APR 6.429%
15-year fixed-rate: interest rate 5.705%, APR 5.848%
10-year fixed-rate: interest rate 5.500%, APR 5.720%
7-year ARM: interest rate 7.011%, APR 7.660%
5-year ARM: interest rate 6.880%, APR 7.754%
3-year ARM: interest rate 6.125%, APR 7.204%
30-year fixed-rate FHA: interest rate 5.527%, APR 6.316%
30-year fixed-rate VA: interest rate 5.684%, APR 6.062%

# Email
{email}
"""
```

```
In [13]: response = mistral(prompt)
         print(response)
```

Subject: Re: Mortgage Rates Inquiry

Dear Anna,

Thank you for reaching out to us regarding mortgage rates. I'm happy to provide you with the information you're looking for.

Our current 30-year fixed-rate APR is 6.484%. In comparison, the APR for our 15-year fixed-rate is lower at 5.848%. The longer term of the 30-year loan comes with a higher APR, but it also means your monthly payments will be lower. The 15-year loan has a lower APR, which can save you more money in interest over the life of the loan, but your monthly payments will be higher.

Please consider your financial situation and future plans when deciding between these two options. If you have any further questions or need assistance in making this decision, please don't hesitate to ask.

Best regards,
Lender Customer Support

Summarization

- We'll use this [article \(https://www.deeplearning.ai/the-batch/mistral-enhances-ai-landscape-in-europe-with-microsoft-partnership-and-new-language-models\)](https://www.deeplearning.ai/the-batch/mistral-enhances-ai-landscape-in-europe-with-microsoft-partnership-and-new-language-models) from The Batch

```
In [14]: newsletter = """
European AI champion Mistral AI unveiled new large language models.

What's new: Mistral AI introduced two closed models, Mistral Large and Mistral Small.

Model specs: The new models' parameter counts, architectures, and capabilities.

Mistral Large achieved 81.2 percent on the MMLU benchmark, outperforming GPT-4.
Both models are fluent in French, German, Spanish, and Italian.
Microsoft's investment in Mistral AI is significant but tiny compared to OpenAI.
Mistral AI and Microsoft will collaborate to train bespoke models.
Behind the news: Mistral AI was founded in early 2023 by engineers from Google and Facebook.

Yes, but: Mistral AI's partnership with Microsoft has divided European AI leaders.

Why it matters: The partnership between Mistral AI and Microsoft could reshape the AI landscape.

We're thinking: Mistral AI has made impressive progress in a short time.
"""
```

```
In [15]: prompt = f"""
You are a commentator. Your task is to write a report on a news
When presented with the newsletter, come up with interesting qu
and answer each question.
Afterward, combine all the information and write a report in th
format.

# Newsletter:
{newsletter}

# Instructions:
## Summarize:
In clear and concise language, summarize the key points and the
presented in the newsletter.

## Interesting Questions:
Generate three distinct and thought-provoking questions that ca
asked about the content of the newsletter. For each question:
- After "Q: ", describe the problem
- After "A: ", provide a detailed explanation of the problem ac
in the question.
- Enclose the ultimate answer in <>.

## Write a analysis report
Using the summary and the answers to the interesting questions,
create a comprehensive report in Markdown format.
"""
```

```
In [16]: response = mistral(prompt)
print(response)
```

Summary

Mistral AI, a European AI champion, introduced two new large language models: Mistral Large and Mistral Small. Microsoft invested \$16.3 million in the French startup, forming an alliance that allows Mistral to use Microsoft's Azure platform for distributing Mistral Large and accessing Azure computing infrastructure. The partnership has stirred controversy among European lawmakers and regulators due to potential data access concerns, but it grants Mistral crucial processing power and global market access while providing Azure customers with a high-performance model tailored to Europe's unique regulatory environment.

Interesting Questions

Q: How do the new Mistral models compare to their competitors in terms of performance and capabilities?

A: Performance-wise, Mistral Large outperforms Anthropic's Claude 2, Google's Gemini Pro, and Meta's Llama 2 70B, achieving 81.2 percent on the MMLU benchmark. However, it falls short of G

Try it out for yourself

- Feel free to copy-paste text from another article in The Batch, or any other blog or news article.


```
In [17]: newsletter2 = ""
```

April 2024, what a month! My birthday, a new book release, spring break, and a new newsletter. This article reviews and discusses all four major transformer-based LLM releases in April 2024:

1. How Good are Mixtral, Llama 3, and Phi-3?
2. OpenELM: An Efficient Language Model Family with Open-source Models
3. Is DPO Superior to PPO for LLM Alignment? A Comprehensive Study
4. Other Interesting Research Papers In April

1. Mixtral, Llama 3, and Phi-3: What's New?

First, let's start with the most prominent topic: the new major releases. We'll start with Mixtral 8x22B, then Llama 3, and finally Phi-3.

1.1 Mixtral 8x22B: Larger models are better!

Mixtral 8x22B is the latest mixture-of-experts (MoE) model by Mistral AI. It's a 22B parameter model, which is larger than the previous Mixtral 8x7B released in January 2024. The key difference is that Mixtral 8x22B has 22 billion parameters, while Mixtral 8x7B has 7 billion parameters. The perhaps most interesting plot from the Mixtral blog post, which shows the performance of Mixtral 8x22B compared to other models, is that Mixtral 8x22B is the best model in its class.

1.2 Llama 3: Larger data is better!

Meta AI's first Llama model release in February 2023 was a big success. Llama 2 was the first open-source LLM that could rival the performance of GPT-4. While Meta is still training some of their largest models (e.g. Llama 3.1), Llama 3 is the latest release. Overall, the Llama 3 architecture is almost identical to Llama 2, but with some improvements in the training data and the model's performance.

Below are the configuration files used for implementing Llama 2. The main contributor to the substantially better performance compared to Llama 1.2 is the training data size.

The main contributor to the substantially better performance compared to Llama 1.2 is the training data size. This is a very interesting finding because, as the Llama 3 blog post states, "the main contributor to the substantially better performance compared to Llama 1.2 is the training data size".

Instruction finetuning and alignment

For instruction finetuning and alignment, researchers usually compare the performance of different models on a set of tasks. The Llama 3 blog post stated that a Llama 3 research paper would be released in the near future.

1.3 Phi-3: Higher-quality data is better!

Just one week after the big Llama 2 release, Microsoft shared the Phi-3 model. Phi-3 is a 3.8 billion parameter model, which is smaller than Llama 2. Notably, Phi-3, which is based on the Llama architecture, has been trained on a higher-quality dataset than Llama 2. Also, Phi-3-mini has "only" 3.8 billion parameters, which is less than Llama 2. So, What is the secret sauce? According to the technical report, the main contributor to the substantially better performance compared to Llama 1.2 is the training data size. The paper didn't go into too much detail regarding the data curriculum.

As of this writing, people are still unsure whether Phi-3 is really better than Llama 2.

1.4 Conclusion

Based on the three major releases described above, this has been a very interesting month for the LLM community.

Which model should we use in practice? I think all three models

2. OpenELM: An Efficient Language Model Family with Open-source
OpenELM: An Efficient Language Model Family with Open-source Tr

Similar to the OLMo, it's refreshing to see an LLM paper that s

Let's start with the most interesting tidbits:

OpenELM comes in 4 relatively small and convenient sizes: 270M,

For each size, there's also an instruct-version available train

OpenELM performs slightly better than OLMo even though it's tra

The main architecture tweak is a layer-wise scaling strategy

2.1 Architecture details

Besides the layer-wise scaling strategy (more details later), t

2.2 Training dataset

Sharing details is different from explaining them as research p

One of the authors kindly followed up with me on that saying "F

2.3 Layer-wise scaling

The layer-wise scaling strategy (adopted from the DeLight: Deep

I wish there was an ablation study training an LLM with and wit

However, we can find ablation studies in the DeLight: Deep and

2.4 LoRA vs DoRA

An interesting bonus I didn't expect was that the researchers c

2.5 Conclusion

While the paper doesn't answer any research questions, it's a g

Anyways, great work, and big kudos to the researchers (and Appl

3. Is DPO Superior to PPO for LLM Alignment? A Comprehensive St
Is DPO Superior to PPO for LLM Alignment? A Comprehensive Study

Let's start with a brief overview before diving into the result

RLHF is a key component of LLM development, and it's used to al

For a more detailed explanation and comparison, also see the Ev

3.1 What are RLHF-PPO and DPO?

RLHF-PPO, the original LLM alignment method, has been the backb

Today, most LLMs on top of public leaderboards have been traine

3.2 PPO is generally better than DPO

Is DPO Superior to PPO for LLM Alignment? A Comprehensive Study

Here, out-of-distribution data means that the LLM has been prev

The main findings are summarized in the figure below.

In addition to the main results above, the paper includes several

3.3 Best practices

Furthermore, interesting takeaways from this paper include best

For instance, if you use DP0, make sure to perform supervised finetuning

If you use PP0, the key success factors are large batch sizes, high

3.4 Conclusion

Based on this paper's results, PP0 seems superior to DP0 if used

A good practical recommendation may be to use PP0 if you have good

Also, based on what we know from the LLama 3 blog post, we don't need to
 """

```
In [18]: prompt2 = f"""
You are a commentator. Your task is to write a report on a newsletter.
When presented with the newsletter, come up with interesting questions and answer each question.
Afterward, combine all the information and write a report in the following format.

# Newsletter:
{newsletter2}

# Instructions:
## Summarize:
In clear and concise language, summarize the key points and themes presented in the newsletter.

## Interesting Questions:
Generate three distinct and thought-provoking questions that can be asked about the content of the newsletter. For each question:
- After "Q: ", describe the problem
- After "A: ", provide a detailed explanation of the problem and add in the question.
- Enclose the ultimate answer in <>.

## Write a analysis report
Using the summary and the answers to the interesting questions, create a comprehensive report in Markdown format.
"""
```

```
In [19]: response2 = mistral(prompt2)
         print(response2)
```

Summary:

The newsletter discusses the recent releases of four transformer-based large language models (LLMs): Mixtral 8x22B, Llama 3, Phi-3, and OpenELM. The newsletter reviews the performance and improvements of each model, discusses their training methods, and compares their effectiveness. Additionally, the newsletter explores a comprehensive study on the superiority of direct preference optimization (DPO) over proximal policy optimization (PPO) for LLM alignment.

Interesting Questions:

Q: How do Mixtral 8x22B, Llama 3, and Phi-3 compare in terms of performance and computational resource requirements?

A: Mixtral 8x22B, Llama 3, and Phi-3 are compared on two axes: modeling performance on the popular Measuring Massive Multitask Language Understanding (MMLU) benchmark and active parameters (related to computational resource requirements). Mixtral 8x22B performs better than Llama 3 and Phi-3 on the MMLU benchmark, but it has a higher active-parameter count. Phi-3 3.8B may be very appealing for mobile devices as it can run on an iPhone 14, while Llama 3 8B might be the most interesting all-rounder for fine-tuning since it can be comfortably fine-tuned on a single GPU when using LoRA.

<Mixtral 8x22B performs better but requires more computational resources, while Phi-3 3.8B is suitable for mobile devices and Llama 3 8B is a good all-rounder for fine-tuning.>

Q: What are the main differences between the training methods of Llama 3 and Phi-3?

A: Llama 3 was trained on 15 trillion tokens, while Phi-3 was trained on 3.3 trillion tokens. Llama 3 has a larger vocabulary size than Phi-3. Llama 3 used both PPO and DPO for instruction finetuning and alignment, while the training methods of Phi-3 were not explicitly mentioned in the newsletter.

<Llama 3 was trained on more data and has a larger vocabulary size, and it used both PPO and DPO for instruction finetuning and alignment.>

Q: What are the main findings of the comprehensive study on the superiority of DPO over PPO for LLM alignment?

A: The study found that PPO is generally better than DPO, and DPO suffers more heavily from out-of-distribution data. Out-of-distribution data means that the LLM has been previously trained on instruction data that is different from the preference data for DPO. The study recommends using PPO if you have ground truth reward labels or can download an in-domain reward model, and using DPO for simplicity.

<PPO is generally better than DPO, but DPO is simpler to use.>

Analysis Report:

Introduction

The newsletter discusses the recent releases of four transformer-based large language models (LLMs): Mixtral 8x22B, Llama 3, Phi-3,

and OpenELM. The newsletter reviews the performance and improvements of each model, discusses their training methods, and compares their effectiveness. Additionally, the newsletter explores a comprehensive study on the superiority of direct preference optimization (DPO) over proximal policy optimization (PPO) for LLM alignment.

Comparison of Mixtral 8x22B, Llama 3, and Phi-3

Mixtral 8x22B, Llama 3, and Phi-3 are compared on two axes: modeling performance on the popular Measuring Massive Multitask Language Understanding (MMLU) benchmark and active parameters (related to computational resource requirements). Mixtral 8x22B performs better than Llama 3 and Phi-3 on the MMLU benchmark, but it has a higher active-parameter count. Phi-3 3.8B may be very appealing for mobile devices as it can run on an iPhone 14, while Llama 3 8B might be the most interesting all-rounder for fine-tuning since it can be comfortably fine-tuned on a single GPU when using LoRA.

Training Methods of Llama 3 and Phi-3

Llama 3 was trained on 15 trillion tokens, while Phi-3 was trained on 3.3 trillion tokens. Llama 3 has a larger vocabulary size than Phi-3. Llama 3 used both PPO and DPO for instruction finetuning and alignment, while the training methods of Phi-3 were not explicitly mentioned in the newsletter.

Superiority of DPO over PPO for LLM Alignment

The comprehensive study found that PPO is generally better than DPO, and DPO suffers more heavily from out-of-distribution data. Out-of-distribution data means that the LLM has been previously trained on instruction data that is different from the preference data for DPO. The study recommends using PPO if you have ground truth reward labels or can download an in-domain reward model, and using DPO for simplicity.

Conclusion

The newsletter provides a comprehensive review of the recent releases of Mixtral 8x22B, Llama 3, Phi-3, and OpenELM, and discusses their performance, training methods, and effectiveness. The newsletter also explores a comprehensive study on the superiority of DPO over PPO for LLM alignment, providing valuable insights for researchers and practitioners in the field.

The Mistral Python client

- Below is the helper function that you imported from helper.py and used earlier in this notebook.
- For more details, check out the [Mistral AI API documentation \(https://docs.mistral.ai/api/\)](https://docs.mistral.ai/api/)
- To get your own Mistral AI API key to use on your own, outside of this classroom, you can create an account and go to the [console \(https://console.mistral.ai/\)](https://console.mistral.ai/) to subscribe and create an API key.

```
In [20]: from mistralai.client import MistralClient
from mistralai.models.chat_completion import ChatMessage

def mistral(user_message,
            model="mistral-small-latest",
            is_json=False):
    client = MistralClient(api_key=os.getenv("MISTRAL_API_KEY"))
    messages = [ChatMessage(role="user", content=user_message)]

    if is_json:
        chat_response = client.chat(
            model=model,
            messages=messages,
            response_format={"type": "json_object"})
    else:
        chat_response = client.chat(
            model=model,
            messages=messages)

    return chat_response.choices[0].message.content
```

Type *Markdown* and LaTeX: α^2