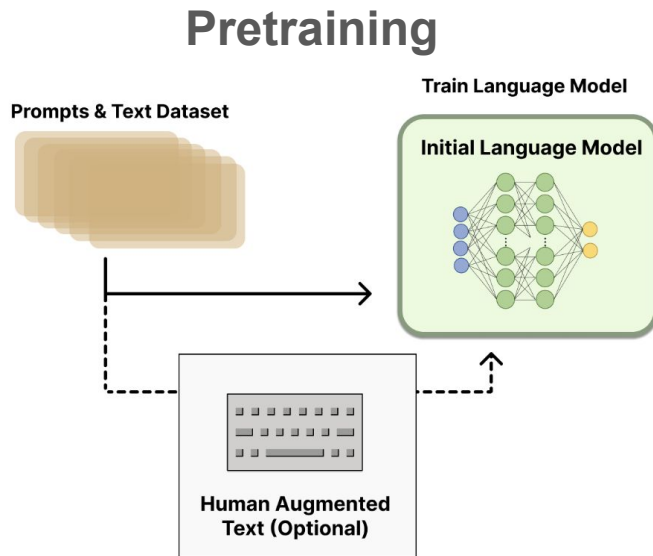


# Тренировка диффузионных моделей при помощи DPO на нескольких изображениях

Куликов Антон

# Предпосылки

Для создания LLM, помимо предтренировки моделей как обычных LM и дотренировки на инструкциях, необходимо дообучать модель на ответах пользователей, иначе модели галлюцинируют.

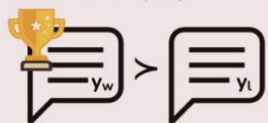


# Метод

Недавно ученые разработали новый подход - DPO, он производит такие же результаты, но проще в применении. Также он нашел свое применение в диффузионных моделях.

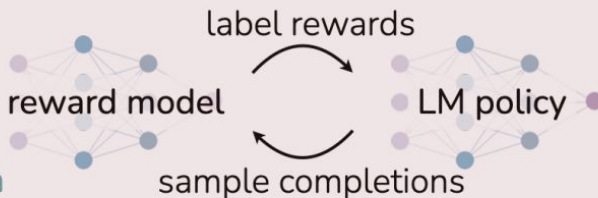
## Reinforcement Learning from Human Feedback (RLHF)

x: "write me a poem about  
the history of jazz"



preference data

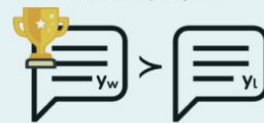
maximum  
likelihood



reinforcement learning

## Direct Preference Optimization (DPO)

x: "write me a poem about  
the history of jazz"



preference data

maximum  
likelihood



# Метод

Помимо DPO ученые разработали еще один более эффективный подход - CPO.

1. [Оригинальная статья](#)
2. [DPO в диффузиях](#)
3. [CPO](#)

# Метод

Основная идея метода заключается в применении DPO/CPO к предтренированной модели при помощи множественных отранжированных изображений.

Также одной из идей является добавление дискриминатора в лосс по образцу GAN.

$$\mathcal{L}_{\text{prefer}} = - \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \pi_{\theta}(y_w | x) - \beta \log \pi_{\theta}(y_l | x) \right) \right].$$

# Дорожная карта проекта

1. Найти датасет с промптами и изображениями - [picapick\\_v2](#).
2. Сгенерировать n изображений по каждому промпту - сделано при помощи [sdxl\\_turbo](#) с разным кол-вом шагов, но возможно использование [sdxl-lightning](#).
3. Тренировка модели скоринга эстетичности изображений (дискриминатора) - сделано. Подробности далее.
4. Тренировка CPO/DPO на отскоренных изображениях - еще не сделано.

# Тренировка дискриминатора

Задача дискриминатора - отличать менее эстетичные изображения от более эстетичных с учетом описания.

Данные:

Положительный класс - самые эстетичные изображения из [coyo](#) (score $\geq$ 6.5).

Отрицательный класс - генерации sd\_xl\_turbo.

# Тренировка дискриминатора

## Отрицательный класс

Отрицательный класс - генерация SDXL\_turbo.

Генерация получена при помощи промпта, полученного при помощи мультимодальной модели [LLAVA](#) при инференсе на изображениях положительного класса.



# Тренировка дискриминатора

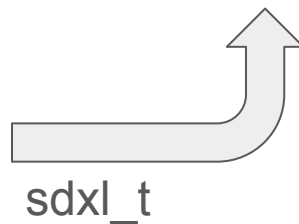
## Получение отрицательного класса



Реальное фото/картинка



The image captures a serene scene of a small, rusted watchtower standing guard over a tranquil lake. The tower...



# Тренировка дискриминатора

Результат работы дискриминатора.

Изображения отсортированы по реалистичности слева направо

sdxl\_t\_1\_coyo



sdxl\_t\_4\_coyo\_llava



sdxl\_t\_4\_coyo\_llava\_prep



real



# Тренировка дискриминатора

Результат работы дискриминатора.

Изображения отсортированы по реалистичности слева направо

sdxl\_t\_1\_coyo



real



sdxl\_t\_4\_coyo\_llava



sdxl\_t\_4\_coyo\_llava\_prep



Модель выбирает реальные изображения в 66% случаев

**Спасибо за внимание!**