# Лекция 5. InstructGPT

VK education
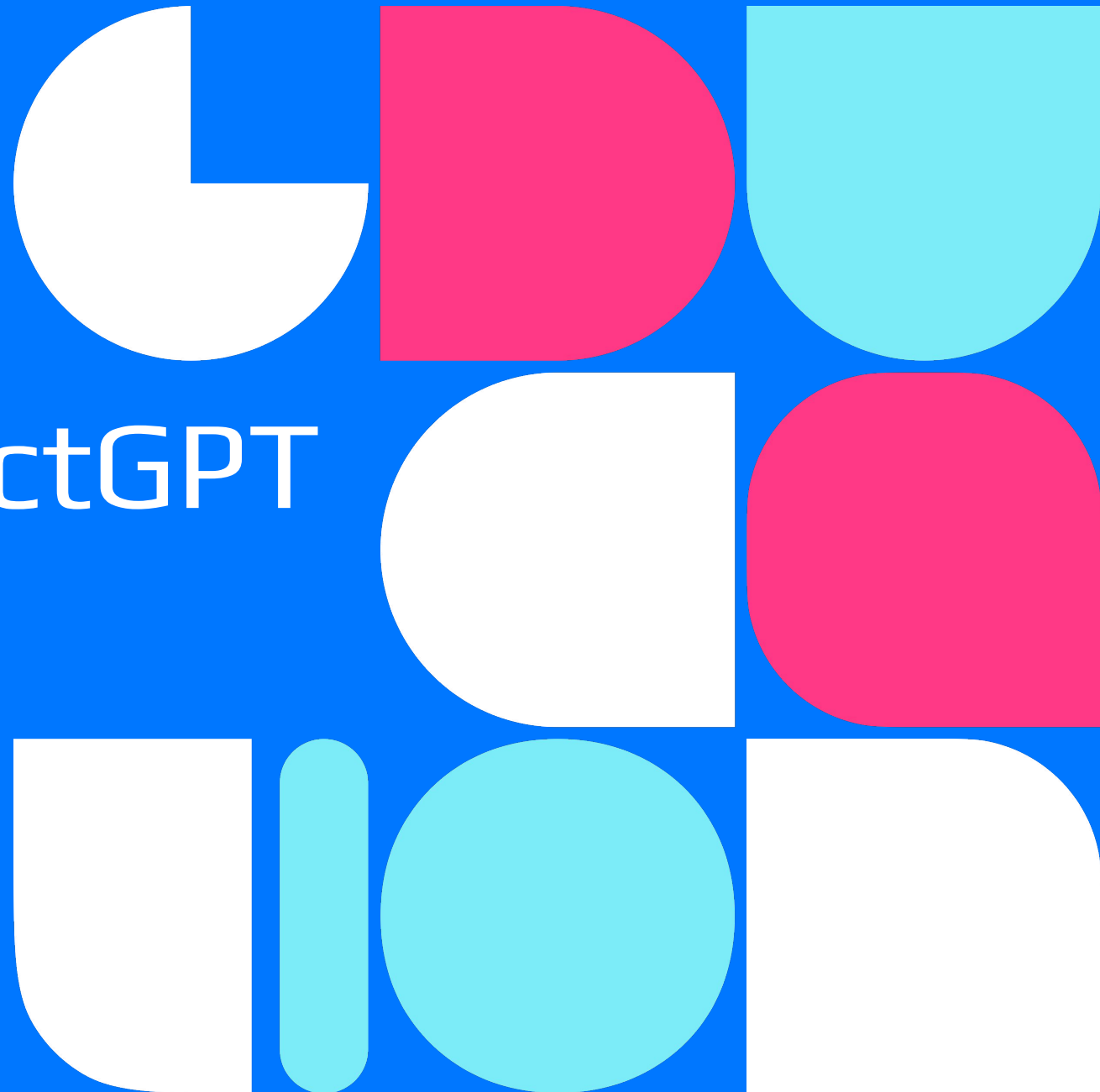
Елисеев Антон
30.10.2024

# LLM Recap

Вопрос: Какая архитектура у LLM?

# LLM Recap

Вопрос: Как собирать данные?

# Instruction following

**Prompt:**
What is the purpose of the list C in the code below?

```
def binomial_coefficient(n, r):
    C = [0 for i in range(r + 1)];
    C[0] = 1;
    for i in range(1, n + 1):
        j = min(i, r);
        while j > 0:
            C[j] += C[j - 1];
            j -= 1;
    return C[r]
```

**GPT-3 175B completion:**
A. to store the value of C[0]
B. to store the value of C[1]
C. to store the value of C[i]
D. to store the value of C[i - 1]

**InstructGPT 175B completion:**
The list C in this code is used to store the values of the binomial coefficient as the function iterates through the values of n and r. It is used to calculate the value of the binomial coefficient for a given value of n and r, and stores the result in the final return value of the function.

# Instruction following

Training objective:

what we **want** vs
what we **told** model to do

What we **told** model to do:
- predict the next token on a webpage from the internet
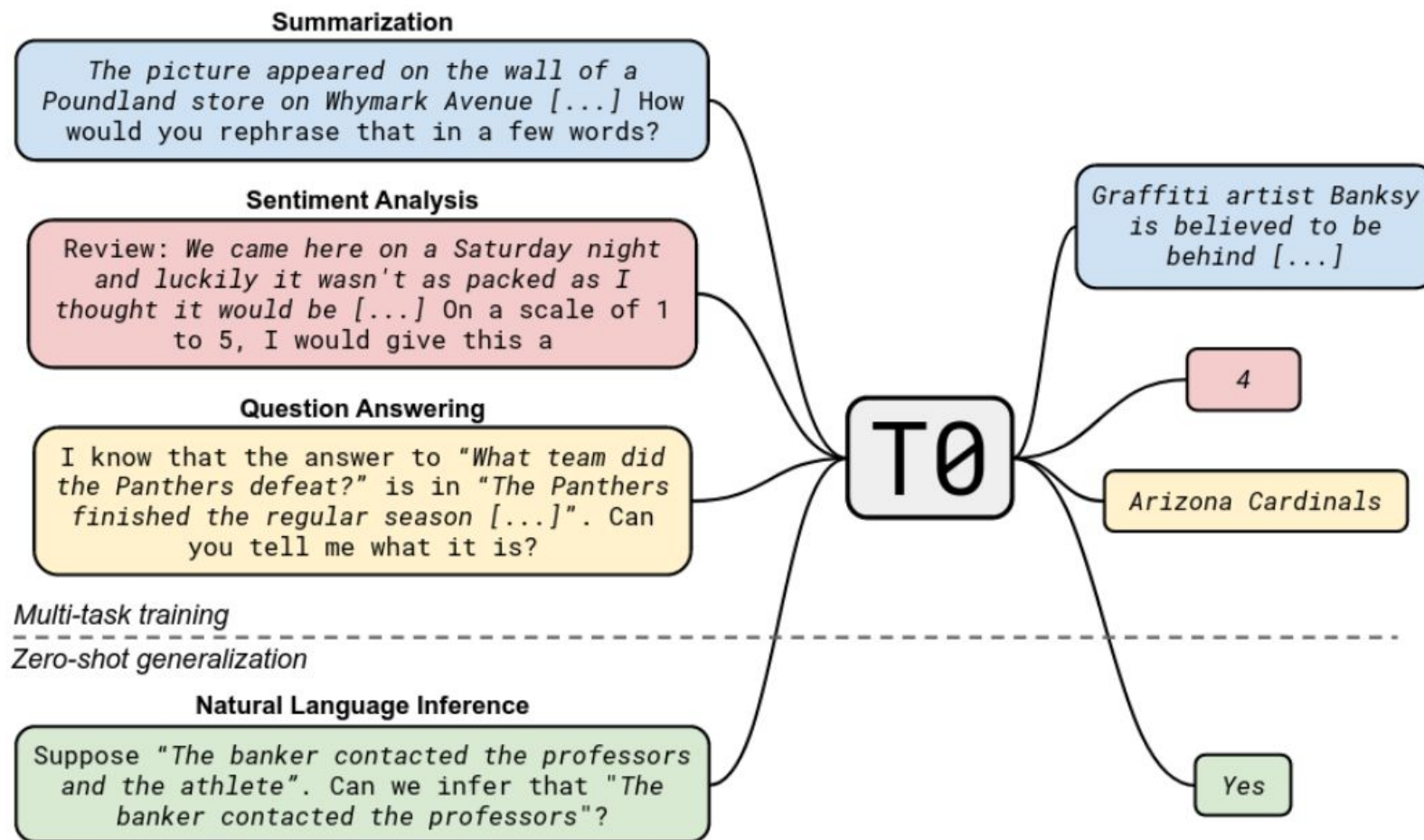
Alignment
——————→

What we want model to do:
- follow the user's instructions helpfully and safely

The language modeling objective is misaligned

# Instruction following

**Summarization**

The picture appeared on the wall of a Poundland store on Whymark Avenue [...] How would you rephrase that in a few words?

**Sentiment Analysis**

Review: We came here on a Saturday night and luckily it wasn't as packed as I thought it would be [...] On a scale of 1 to 5, I would give this a

**Question Answering**

I know that the answer to "What team did the Panthers defeat?" is in "The Panthers finished the regular season [...]". Can you tell me what it is?

*Multi-task training*
- - - - - - - - - - - - - - - - - - - - - -
*Zero-shot generalization*

**Natural Language Inference**

Suppose "The banker contacted the professors and the athlete". Can we infer that "The banker contacted the professors"?

T0

Graffiti artist Banksy is believed to be behind [...]

4

Arizona Cardinals

Yes

# SFT issues

Вопрос: В чем проблемы SFT?

# SFT issues

- Supervised seq2seq learning:

$$P\left(y_{t+1}\middle|x,y_{0:t}\right), \qquad y_{0:t} \sim reference$$

- Inference

$$P\left(y_{t+1}\middle|x,\hat{y}_{0:t}\right), \qquad \hat{y}_{0:t} \sim \ ???$$

# SFT issues

- Supervised seq2seq learning:

$$P\left(y_{t+1}|x, y_{0:t}\right), \qquad y_{0:t} \sim reference$$

- Inference

$$P\left(y_{t+1}|x, \hat{y}_{0:t}\right), \qquad \hat{y}_{0:t} \sim model$$

If model ever makes something that isn't in data,
It gets volatile from next time-step!

# SFT issues

There's more then one correct translation.
You don't need to learn all of them.

**Source:** 在 找 给 家里 人 的 礼物 .

**Versions:**
i 'm searching for some gifts for my family.
i want to find something for my family as presents.
i 'm about to buy some presents for my family.
i 'd like to buy my family something as a gift.
i 'm looking for a present for my family.
...

# SFT issues

There's more then one correct translation.
You don't need to learn all of them.

**Source:** 在 找 给 家里 人 的 礼物 .

| Versions: | Model 1 $p(y\|x)$ | Model 2 $p(y\|x)$ |
|---|---|---|
| (version 1) | 1e-2 | 0.99 |
| (version 2) | 2e-2 | 1e-100 |
| (version 3) | 1e-2 | 1e-100 |
| (all rubbish) | 0.96 | 0.01 |

# SFT issues

There's more then one correct translation.
You don't need to learn all of them.

**Source:** 在 找 给 家里 人 的 礼物 .

|  | Model 1 $p(y|x)$ | Model 2 $p(y|x)$ |
|---|---|---|
| **Versions:** | | |
| (version 1) | 1e-2 | 0.99 |
| (version 2) | 2e-2 | 1e-100 |
| (version 3) | 1e-2 | 1e-100 |
| (all rubbish) | 0.96 | 0.01 |
| | better llh | worse llh |
| | 96% rubbish | 1% rubbish |

# SFT issues

Вопрос: Почему тогда не использовать только RL?

# SFT issues

Вопрос: Почему тогда не использовать только RL?

1. Это RL -> плохо сходится, нужно много данных
2. Говорит "плохо", но не говорит, как правильно

Вывод: делаем SFT, поверх него RL

# RL recap

Политика
агента:

$$a := \pi(s)$$

агент

наблюдение $s$

среда

Decision Process - выбор действий по наблюдениям

# RL recap



Политика агента:

$$a := \pi(s)$$

агент

действие $a$

следующее наблюдение $s'$

среда

Decision Process - выбор действий по наблюдениям

# RL recap

Supervised learning:

$$\nabla llh = \mathop{E}_{x, y_{opt} \sim D} \nabla \log P_\theta(y_{opt}|x)$$

Policy gradient:

$$\nabla J = \mathop{E}_{\substack{s \sim d(s) \\ a \sim \pi(a|obs(s))}} \nabla \log \pi(a|s) Q(s,a)$$

# RL recap

Supervised learning:

$$\nabla llh = \underset{s, a_{opt} \sim D}{E} \nabla \log \pi_\theta (a_{opt}|s)$$

Policy gradient:

$$\nabla J = \underset{\substack{s \sim d(s) \\ a \sim \pi(a|obs(s))}}{E} \nabla \log \pi_\theta (a|s) Q(s,a)$$

# RL recap

Supervised learning:

$$\nabla llh = \underset{s, a_{opt} \sim D}{E} \nabla \log \pi_\theta (a_{opt}|s)$$

**reference**

Policy gradient:

$$\nabla J = \underset{\substack{s \sim d(s) \\ a \sim \pi(a|obs(s))}}{E} \nabla \log \pi_\theta (a|s) Q(s,a)$$

**generated**

# RLHF

## Step 1

**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

Some people went to the moon...

This data is used to fine-tune GPT-3 with supervised learning.

SFT

## Step 2

**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

- A Explain gravity...
- B Explain war...
- C Moon is natural satellite of...
- D People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

## Step 3

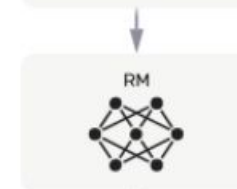**Optimize a policy against the reward model using reinforcement learning.**
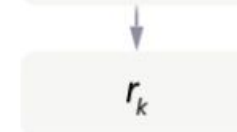
A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

# RLHF

Вопрос: Что плохого в этой схеме?

# RLHF

$$\mathbf{E}_{a \sim \pi_\theta(a|s)} \left[ r_\psi(s, a) - \beta \mathrm{KL}\big(\pi_\theta(a|s) || \pi_{\mathrm{SFT}}(a|s)\big) \right] \rightarrow \max_\theta$$

# RLHF

Paper: https://arxiv.org/abs/2203.02155

## Step 1

**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

Some people went to the moon...

This data is used to fine-tune GPT-3 with supervised learning.

SFT

## Step 2

**Collect comparison data, and train a reward model.**

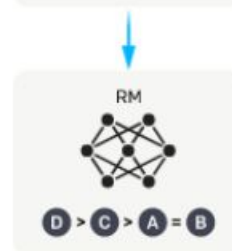A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A — Explain gravity...
B — Explain war...
C — Moon is natural satellite of...
D — People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

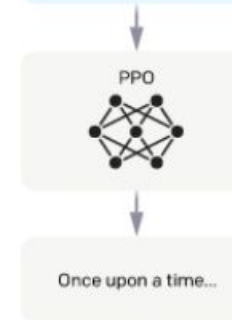## Step 3

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

# Reward modeling. Модель Брэдли-Терри

Cons:
1. Не отношения порядка, парадокс Кондорсе

Pros:
1. Делает так же, как и люди

1. **Рефлексивность:** $a \preccurlyeq a$.

2. **Антисимметричность:** если $a \preccurlyeq b$ и $b \preccurlyeq a$, то $a = b$.

3. **Транзитивность:** если $a \preccurlyeq b$ и $b \preccurlyeq c$, то $a \preccurlyeq c$.

# Reward modeling. Модель Брэдли-Терри

$$P(a > b|s) = \sigma(r_\psi(s, a) - r_\psi(s, b)) \qquad \sigma(x) = \frac{1}{1 + \exp(-x)}$$

$$\sum_{(s, winner, loser) \in \mathbf{D}} \log \sigma(r_\psi(s, winner) - r_\psi(s, loser)) \to \max_\psi$$

# RLHF

Paper: https://arxiv.org/abs/2203.02155



**Step 1**

**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

Some people went to the moon...
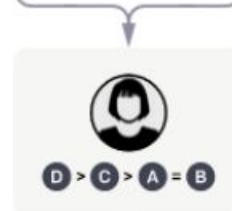
This data is used to fine-tune GPT-3 with supervised learning.

SFT

**Step 2**

**Collect comparison data, and train a reward model.**

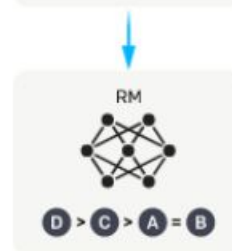A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A  Explain gravity...
B  Explain war...
C  Moon is natural satellite of...
D  People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B
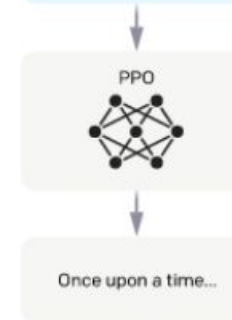
This data is used to train our reward model.

RM

D > C > A = B

**Step 3**

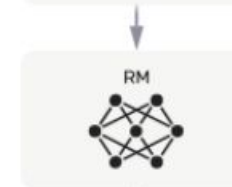**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

# PPO issues

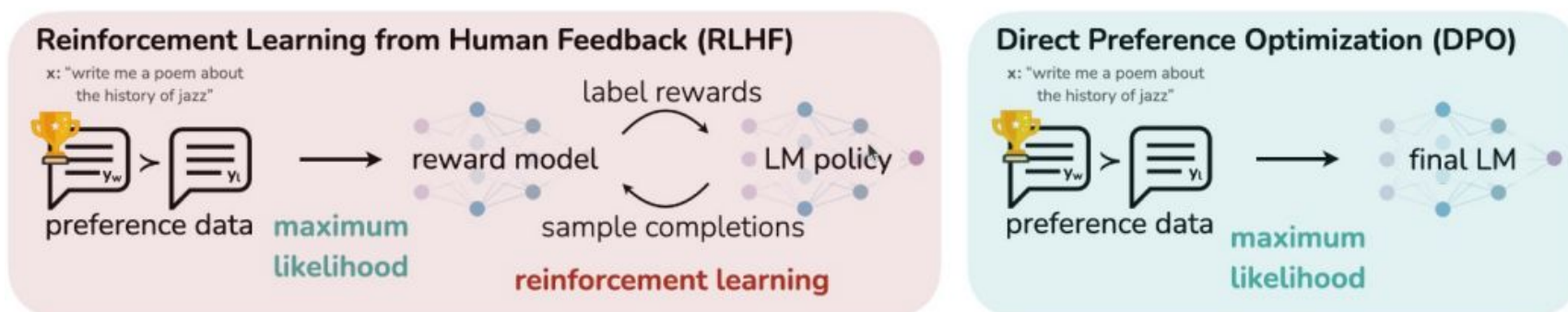Вопрос: Какие проблемы у такого подхода?

# PPO issues

Вопрос: Какие проблемы у такого подхода?

1. Это RL -> плохо сходится, у PPO много гиперпараметров
2. Много GPU памяти
   Пример: gpt3 175B требует 7TB gpu-памяти ~ 88 A100 80G ~ $2M
3. Надо инферить модель на каждом шагу

# DPO

## RLHF is complicated, let's make it simpler!



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

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

label rewards

reward model ⟷ LM policy

sample completions

preference data — maximum likelihood — reinforcement learning

**Direct Preference Optimization (DPO)**

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

final LM

preference data — maximum likelihood

## Direct Preference Optimization: Your Language Model is Secretly a Reward Model

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, Chelsea Finn

While large-scale unsupervised language models (LMs) learn broad world knowledge and some reasoning skills, achieving precise control of their beh... methods for gaining such steerability collect human labels of the relative quality of model generations and fine-tune the unsupervised LM to align wit... However, RLHF is a complex and often unstable procedure, first fitting a reward model that reflects the human preferences, and then fine-tuning the... reward without drifting too far from the original model. In this paper we introduce a new parameterization of the reward model in RLHF that enables e... standard RLHF problem with only a simple classification loss. The resulting algorithm, which we call Direct Preference Optimization (DPO), is stable, p... the LM during fine-tuning or performing significant hyperparameter tuning. Our experiments show that DPO can fine-tune LMs to align with human p... exceeds PPO-based RLHF in ability to control sentiment of generations, and matches or improves response quality in summarization and single-turn...

# DPO

$$\mathbf{E}_{a \sim \pi_\theta(a|s)} \left[ r_\psi(s, a) - \beta \mathrm{KL}\big(\pi_\theta(a|s) || \pi_{\mathrm{SFT}}(a|s)\big) \right] \rightarrow \max_\theta$$

# DPO

$$\mathbf{E}_{a \sim \pi_\theta(a|s)} \left[ r_\psi(s, a) - \beta \mathrm{KL}\big(\pi_\theta(a|s) || \pi_{\mathrm{SFT}}(a|s)\big) \right] \to \max_\theta$$

$$\pi^*(a|s) = \frac{1}{Z(s)} \pi_{\mathrm{SFT}}(a|s) e^{\frac{1}{\beta} r(s,a)}$$

$$Z(s) = \sum_a e^{\frac{1}{\beta} r_{(s,a)}}$$

# DPO

$$\pi^*(a|s) = \frac{1}{Z(s)}\pi_{\text{SFT}}(a|s)e^{\frac{1}{\beta}r(s,a)}$$

$$r(s,a) = \beta\log\frac{\pi^*(a|s)}{\pi_{\text{SFT}}(a|s)} + \beta\log Z(s)$$

$$Z(s) = \sum_a e^{\frac{1}{\beta}r(s,a)}$$

$$r_\theta(s,a) = \beta\log\frac{\pi_\theta(a|s)}{\pi_{\text{SFT}}(a|s)} + \beta\log Z(s)$$

# DPO

$$\sum_{(s,winner,loser)\in\mathbf{D}} \log\sigma(r_\theta(s,winner) - r_\theta(s,loser)) \to \max_\theta$$

$$r_\theta(s,a) = \beta\log\frac{\pi_\theta(a|s)}{\pi_{\mathrm{SFT}}(a|s)} + \beta\log Z(s)$$

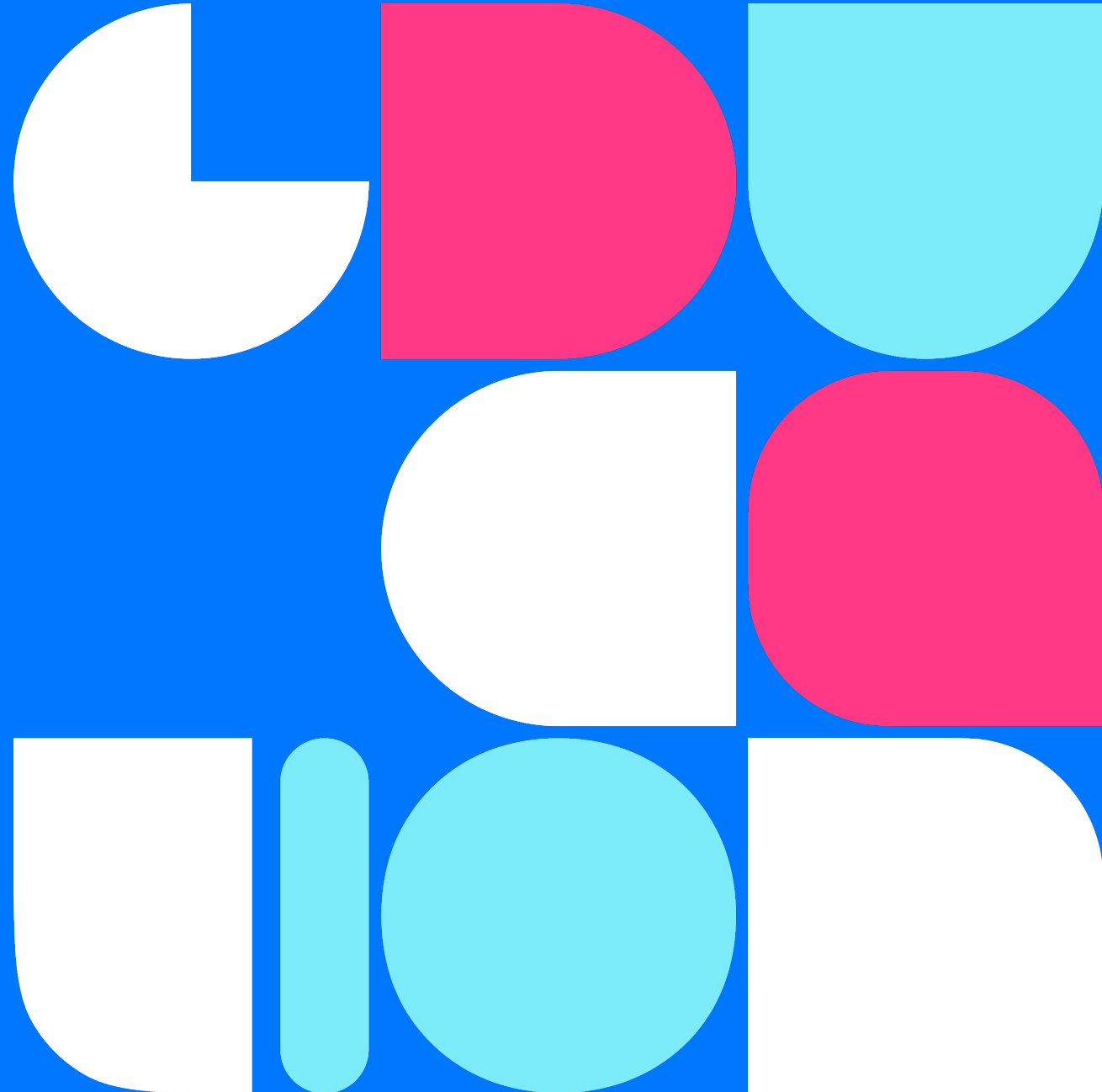$$\sum_{(s,winner,loser)\in\mathbf{D}} \log\sigma\Big(\beta\Big[\log\frac{\pi_\theta(winner|s)}{\pi_{\mathrm{SFT}}(winner|s)} - \log\frac{\pi_\theta(loser|s)}{\pi_{\mathrm{SFT}}(loser|s)}\Big]\Big) \to \max_\theta$$

# DPO

Плюсы:
1. не RL -> лучше и быстрее сходится
2. Только один гиперпараметр - beta
3. Надо хранить только обучаемую модель
4. Не инферим, только прогоняем существующий ответ

# Q&A

**VK** education

# Спасибо за внимание!

Елисеев Антон

education