# Lecture 14: RL from Human Feedback CS885 Reinforcement Learning

2025-02-27

Complementary readings:
Stiennon, Ouyang, Wu, Ziegler, Lowe Voss, Radford, Amodei, Christiano (2020) Learning to summarize from human feedback, NeurIPS.
Ouyang, Wu, Jiang, Wainwright, et al. (2022) Training language models to follow instructions with human feedback, NeurIPS.
Holtzman, Buys, Du, Forbes, Choi (2019). The Curious Case of Neural Text Degeneration, arxiv.
Rafailov, Sharma, Mitchell, Ermon, Manning, Finn (2023) Direct Preference Optimization: Your Language Model is Secretly a Reward Model, NeurIPS.
Rashid, Wu, Fan, Li, Kristiadi, Poupart (2025) Towards Cost-Effective Reward Guided Text Generation, arxiv.

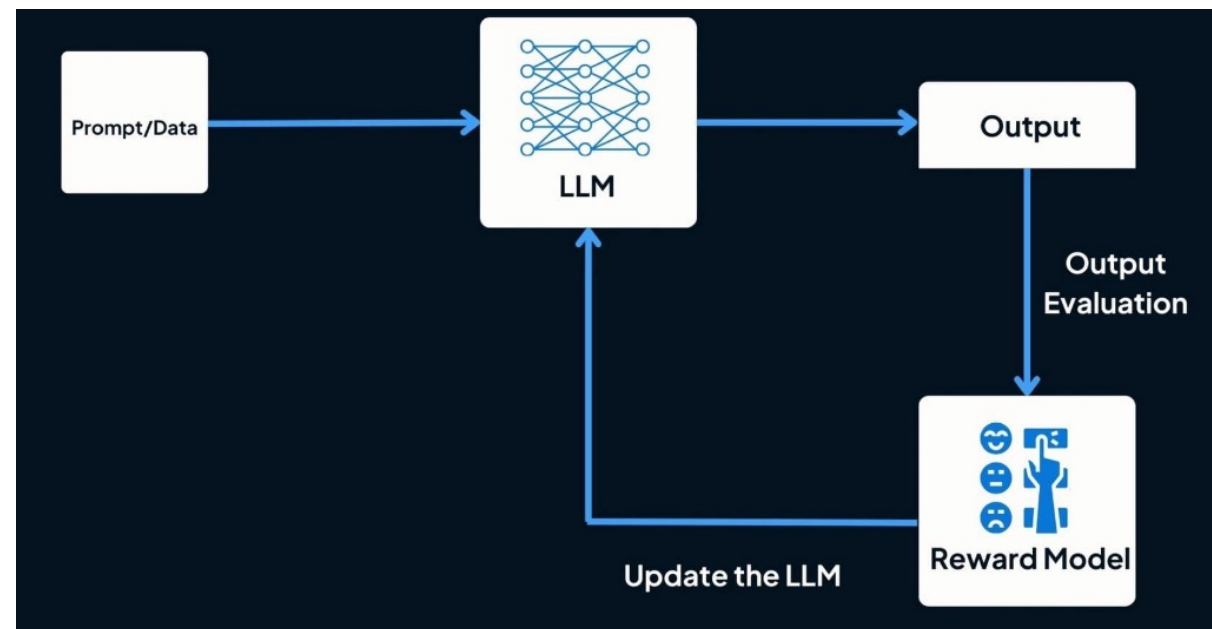Pascal Poupart
David R. Cheriton School of Computer Science

UNIVERSITY OF
WATERLOO

# Outline

- Reinforcement Learning from Human Feedback

- Direct Preference Optimization

- Reward Guided Text Generation

UNIVERSITY OF
WATERLOO

# Large Language Models

- **Agent:** system
- **Environment:** user
- **State:** history of past utterances
- **Action:** system utterance
- **Reward:** task completion, human feedback



Credit: https://www.twine.net/blog/what-is-reinforcement-learning-from-human-feedback-rlhf-and-how-does-it-work/

"*We posit that the superior writing abilities of LLMs, as manifested in surpassing human annotators in certain tasks, are fundamentally driven by RLHF, as documented in Gilardi et al. (2023) and Huang et al. (2023).*"
Llama 2 Technical Report (Touvron et. al 2023)

"*This behavior (re-evaluation) is not only a testament to the model's growing reasoning abilities but also a captivating example of how reinforcement learning can lead to unexpected and sophisticated outcomes*"
Deepseek-R1 Technical Report (Guo et. al 2025)

UNIVERSITY OF
**WATERLOO**

# Recipe For Large Language Model Training



InstructGPT (Ouyang et. al 2022)

# Response Ranking

PP    Explain reinforcement learning in three sentences.

< 1 / 2 >   Reinforcement learning is a type of machine learning where an agent learns to make decisions by receiving feedback in the form of rewards or punishments. The agent's goal is to maximize its cumulative reward over time by exploring different actions and learning which actions lead to the highest rewards. Reinforcement learning has been successfully applied to a variety of tasks, including game playing, robotics, and recommendation systems.

< 2 / 2 >   Reinforcement learning is a type of machine learning that focuses on training an agent to make decisions based on trial and error. The agent receives feedback in the form of rewards or penalties for each decision it makes. By learning from this feedback, the agent can improve its decision-making abilities over time.

Was this response better or worse?    👍 Better    👎 Worse    ⊜ Same    ✕

UNIVERSITY OF
WATERLOO

# RL from Human Feedback (RLHF)

- Collect a preference data set:

$$D = \{(\boldsymbol{s}, \boldsymbol{a}_+, \boldsymbol{a}_-)_k\}_{k=1}^{K} \text{ where } \boldsymbol{a}_+ \succ \boldsymbol{a}_-$$

- Train a reward model according to the Bradley Terry Model:

$$\max_\theta E_D[\log \sigma(r_\theta(\boldsymbol{s}, \boldsymbol{a}_+) - r_\theta(\boldsymbol{s}, \boldsymbol{a}_-))]$$

- Make a copy of the LLM and finetune it to maximize:

$$\max_\phi E_{D,\pi_\phi}[r_\phi(\boldsymbol{s}, \boldsymbol{a})] - \beta KL[\pi_\phi(\boldsymbol{a}|\boldsymbol{s})||\pi_{pretrained}(\boldsymbol{a}|\boldsymbol{s})]$$

UNIVERSITY OF
WATERLOO

# RLHF Improvements

**Proximal Policy Optimization (PPO)**
**Ouyang et al., 2022**

**Direct Preference Optimization (DPO)**
**Rafailov et al., 2023**

**Reward Guided Text Generation (RGTG)**
**Khanov et al, 2024**
**Rashid et al., 2025**

UNIVERSITY OF
**WATERLOO**

# LLM Alignment with Preference Data

- Collect preference data: $D = \{(s, a_+, a_-)_k\}_{k=1}^{K}$

  where $\quad s$: user prompt $\qquad a$: system response

  $a_+$ is preferred to $a_-$ (i.e., $a_+ \succ a_-$)

# Reward Model

Stiennon, Ouyang, Wu, Ziegler, Lowe Voss, Radford, Amodei, Christiano (2020) **Learning to summarize from human feedback**, *NeurIPS.*

- Reward function: $r_\theta(s, a) = real\ number$

- Consider several possible responses $a_1 \succcurlyeq a_2 \succcurlyeq \cdots \succcurlyeq a_k$ ranked by annotator

- Training reward function to be consistent with the ranking:

$$Loss(\theta) = -\frac{1}{\binom{k}{2}} E_{(s, a_i, a_j) \in Dataset} \log \sigma \left( r_\theta(s, a_i) - r_\theta(s, a_j) \right)$$

UNIVERSITY OF
**WATERLOO**

# Reinforcement Learning

Ouyang, Wu, Jiang, Wainwright, et al. (2022) **Training language models to follow instructions with human feedback**, *NeurIPS*.
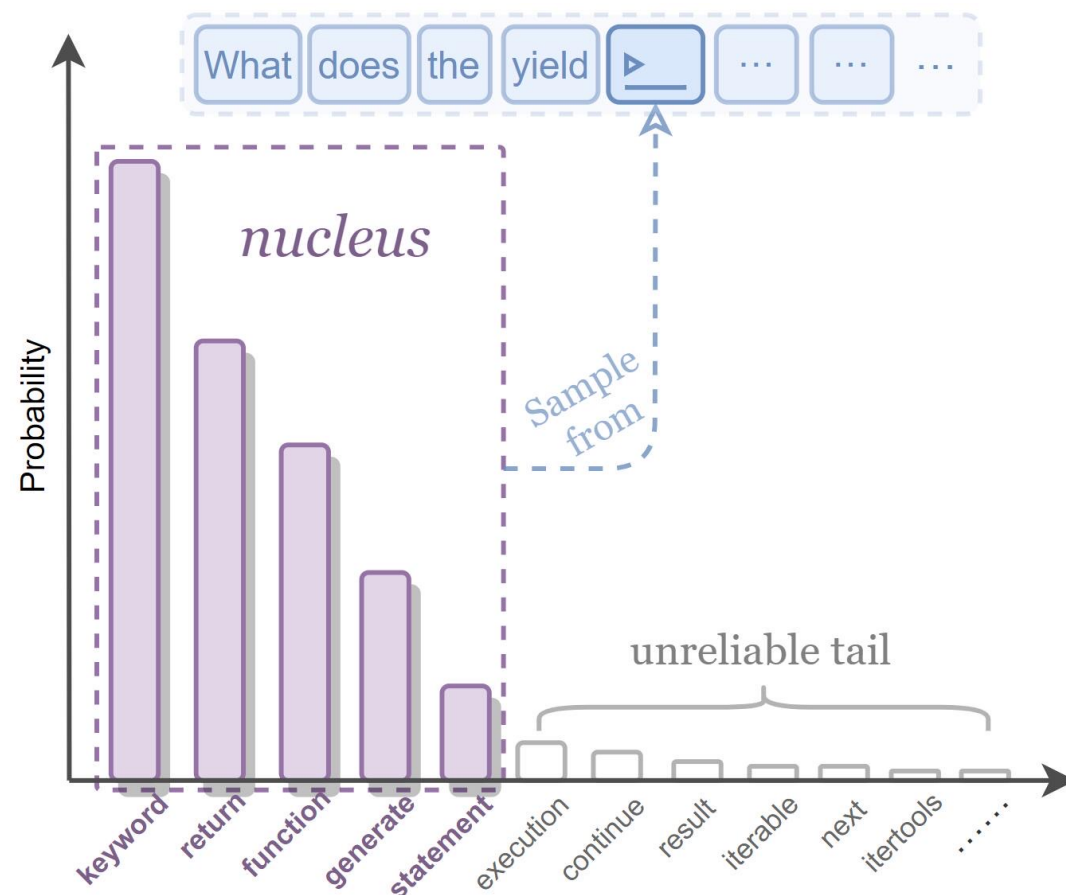
- Pretrain language model (GPT-3)
- Fine-Tune GPT-3 by RL to obtain InstructGPT
  - Policy (language model): $\pi_\phi(a|s)$
  - Optimize $\pi_\phi(s)$ by Proximal Policy Iteration (PPO)

$$\max_\phi E_{s \in Dataset}\left[E_{a \sim \pi_\phi(a|s)}[r_\theta(s,a)] - \beta\, KL\big(\pi_\phi(\cdot|s)\big|\pi_{ref}(\cdot|s)\big)\right]$$

UNIVERSITY OF
**WATERLOO**

# Inference: Nucleus sampling

Sample from nucleus (top tokens only) to avoid unreliable responses while ensuring diversity

Holtzman, Ari; Buys, Jan; Du, Li; Forbes, Maxwell; Choi, Yejin (2019). **The Curious Case of Neural Text Degeneration**, arxiv.



Credit: https://ar5iv.labs.arxiv.org/html/2208.11523

UNIVERSITY OF
**WATERLOO**

# InstructGPT Results

Ouyang, Wu, Jiang, Wainwright, et al. (2022)

# RLHF Improvements

**Proximal Policy Optimization (PPO)**
**Ouyang et al., 2022**

**Direct Preference Optimization (DPO)**
**Rafailov et al., 2023**

**Reward Guided Text Generation (RGTG)**
**Khanov et al, 2024**
**Rashid et al., 2025**

UNIVERSITY OF
**WATERLOO**

# Direct Preference Optimization

Rafailov, Sharma, Mitchell, Ermon, Manning, Finn (2023) **Direct Preference Optimization: Your Language Model is Secretly a Reward Model**, *NeurIPS*.

UNIVERSITY OF
WATERLOO

# Bypassing RL

- Recall RL objective:

$$\max_{\phi} E_{s \in Dataset} \left[ E_{a \sim \pi_\phi(a|s)}[r_\theta(s,a)] - \beta \, KL\big(\pi_\phi(\cdot\,|s)\big|\pi_{ref}(\cdot\,|s)\big) \right]$$

- Closed form solution (based on maximum entropy RL):

$$\pi_\phi(a|s) = \frac{1}{Z(s)} \pi_{ref}(a|s) \exp\left( \frac{r_\theta(s,a)}{\beta} \right)$$

- Isolate reward: $r_\theta(s,a) = \beta \log \dfrac{\pi_\phi(a|s)}{\pi_{ref}(a|s)} + \beta \log Z(s)$

- Plug into preference objective:

$$Loss(\theta) = -\frac{1}{\binom{k}{2}} E_{(s,a_i,a_j) \in Dataset} \log \sigma \left( r_\theta(s,a_i) - r_\theta(s,a_j) \right)$$

$$= -\frac{1}{\binom{k}{2}} E_{(s,a_i,a_j) \in Dataset} \log \sigma \left( \beta \log \frac{\pi_\phi(a_i|s)}{\pi_{ref}(a_i|s)} - \beta \log \frac{\pi_\phi(a_j|s)}{\pi_{ref}(a_j|s)} \right)$$

UNIVERSITY OF
WATERLOO

# Optimal Policy Derivation

$$\operatorname*{argmax}_{\phi} E_{s \in Dataset}\left[E_{a \sim \pi_\phi(a|s)}[r_\theta(s,a)] - \beta\, KL\big(\pi_\phi(\cdot|s)\big|\pi_{ref}(\cdot|s)\big)\right]$$

$$= \operatorname*{argmax}_{\phi} E_{s \in Dataset}\left[E_{a \sim \pi_\phi(a|s)}\left[r_\theta(s,a) - \beta \log\frac{\pi_\phi(a|s)}{\pi_{ref}(a|s)}\right]\right] \qquad \text{by KL definition}$$

$$= \operatorname*{argmin}_{\phi} E_{s \in Dataset}\left[E_{a \sim \pi_\phi(a|s)}\left[\log\frac{\pi_\phi(a|s)}{\pi_{ref}(a|s)} - \frac{1}{\beta} r_\theta(s,a)\right]\right] \qquad \text{since max = - min}$$

$$= \operatorname*{argmin}_{\phi} E_{s \in Dataset}\left[E_{a \sim \pi_\phi(a|s)}\left[\log\frac{\pi_\phi(a|s)}{\frac{1}{Z(s)}\pi_{ref}(a|s)\exp\left(\frac{r_\theta(s,a)}{\beta}\right)} - \log Z(s)\right]\right] \qquad \text{where } Z(s) = \sum_a \pi_{ref}(a|s)\exp\left(\frac{r_\theta(s,a)}{\beta}\right)$$

$$= \operatorname*{argmin}_{\phi} E_{s \in Dataset}\left[E_{a \sim \pi_\phi(a|s)}\left[\log\frac{\pi_\phi(a|s)}{\frac{1}{Z(s)}\pi_{ref}(a|s)\exp\left(\frac{r_\theta(s,a)}{\beta}\right)}\right]\right] \qquad \text{since } \log Z(s) \text{ is independent of } \phi$$

$$= \operatorname*{argmin}_{\phi} E_{s \in Dataset}\left[E_{a \sim \pi_\phi(a|s)}\left[\log\frac{\pi_\phi(a|s)}{\pi_{\phi^*}(a|s)}\right]\right] \qquad \text{where } \pi_{\phi^*}(a|s) = \frac{1}{Z(s)}\pi_{ref}(a|s)\exp\left(\frac{r_\theta(s,a)}{\beta}\right)$$

$$= \operatorname*{argmin}_{\phi} E_{s \in Dataset}\left[KL(\pi_\phi(\cdot|s)||\pi_{\phi^*}(\cdot|s))\right] \qquad \text{by KL definition}$$

$$= \phi^* \qquad \text{since KL is minimized when both arguments are equal}$$

UNIVERSITY OF
WATERLOO

# Empirical Results

Rafailov et al. 2023

# RLHF Improvements

**Proximal Policy Optimization (PPO)**
**Ouyang et al., 2022**

**Direct Preference Optimization (DPO)**
**Rafailov et al., 2023**

**Reward Guided Text Generation (RGTG)**
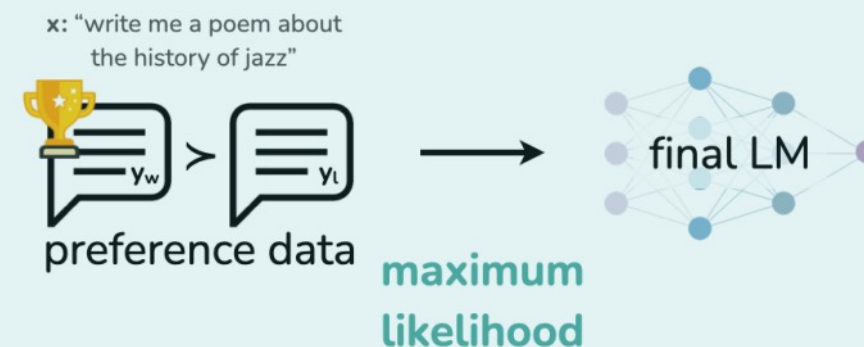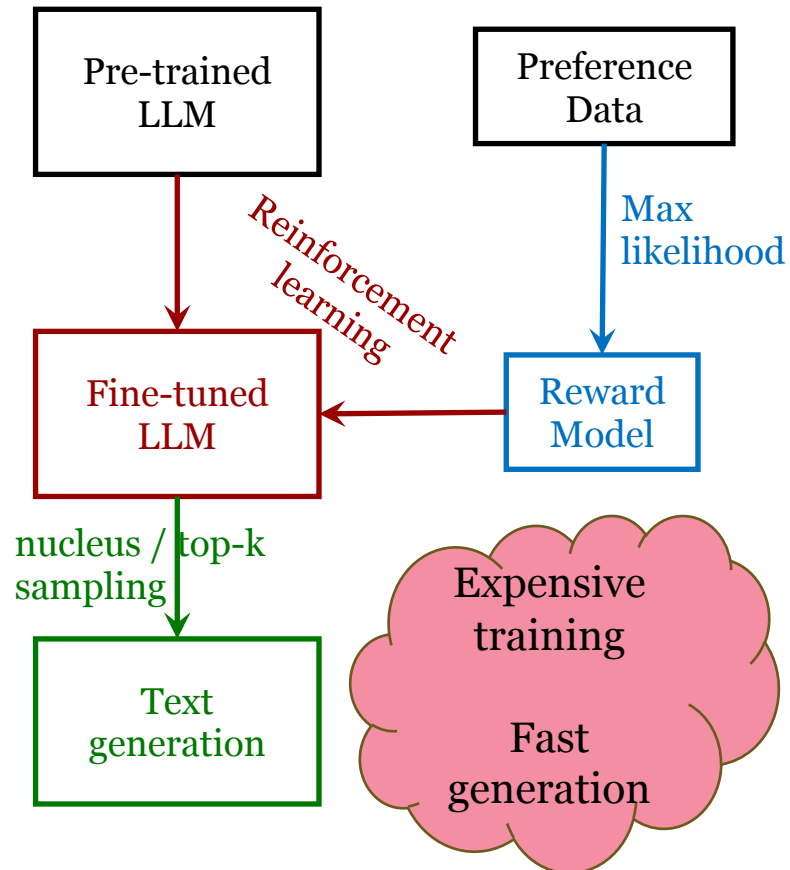**Khanov et al, 2024**
**Rashid et al., 2025**

UNIVERSITY OF
**WATERLOO**

# Sequence Generation

- Recall closed form solution

$$\pi_\phi(\boldsymbol{a}|\boldsymbol{s}) = \frac{1}{Z(\boldsymbol{s})} \pi_{ref}(\boldsymbol{a}|\boldsymbol{s}) \exp\left(\frac{r_\theta(\boldsymbol{s},\boldsymbol{a})}{\beta}\right)$$

$$= softmax\left(\log \pi_{ref}(\boldsymbol{a}|\boldsymbol{s}) + \frac{r_\theta(\boldsymbol{s},\boldsymbol{a})}{\beta}\right)$$

- Text generation:

$$\boldsymbol{a} \sim softmax\left(\log\begin{pmatrix}\pi_{ref}(\boldsymbol{a}_1|\boldsymbol{s})\\ \pi_{ref}(\boldsymbol{a}_2|\boldsymbol{s})\\ \pi_{ref}(\boldsymbol{a}_3|\boldsymbol{s})\\ \dots \\ \pi_{ref}(\boldsymbol{a}_n|\boldsymbol{s})\end{pmatrix} + \begin{pmatrix}r_\theta(\boldsymbol{s},\boldsymbol{a}_1)\\ r_\theta(\boldsymbol{s},\boldsymbol{a}_2)\\ r_\theta(\boldsymbol{s},\boldsymbol{a}_3)\\ \dots \\ r_\theta(\boldsymbol{s},\boldsymbol{a}_n)\end{pmatrix}/\beta\right)$$

UNIVERSITY OF
WATERLOO

# Token Generation

- Token-wise LLM modeling

$$\pi_\phi\left(a^i | \boldsymbol{s}, \boldsymbol{a}^{1:i-1}\right) = \frac{1}{Z(\boldsymbol{s})} \pi_{ref}\left(a^i | \boldsymbol{s}, \boldsymbol{a}^{1:i-1}\right) \exp\left(\frac{r_\theta(\boldsymbol{s}, \boldsymbol{a}^{1:i})}{\beta}\right)$$

$$= softmax\left(\log \pi_{ref}\left(a^i | \boldsymbol{s}, \boldsymbol{a}^{1:i-1}\right) + \frac{r_\theta(\boldsymbol{s}, \boldsymbol{a}^{1:i})}{\beta}\right)$$

- Token generation:

$$a^i \sim softmax\left(\log\begin{pmatrix}\pi_{ref}\left(a_1^i | \boldsymbol{s}, \boldsymbol{a}^{1:i-1}\right) \\ \pi_{ref}\left(a_2^i | \boldsymbol{s}, \boldsymbol{a}^{1:i-1}\right) \\ \pi_{ref}\left(a_3^i | \boldsymbol{s}, \boldsymbol{a}^{1:i-1}\right) \\ \dots \\ \pi_{ref}\left(a_n^i | \boldsymbol{s}, \boldsymbol{a}^{1:i-1}\right)\end{pmatrix} + \begin{pmatrix}r_\theta\left(\boldsymbol{s}, \boldsymbol{a}^{1:i-1}, a_1^i\right) \\ r_\theta\left(\boldsymbol{s}, \boldsymbol{a}^{1:i-1}, a_2^i\right) \\ r_\theta\left(\boldsymbol{s}, \boldsymbol{a}^{1:i-1}, a_3^i\right) \\ \dots \\ r_\theta\left(\boldsymbol{s}, \boldsymbol{a}^{1:i-1}, a_n^i\right)\end{pmatrix} / \beta\right)$$

UNIVERSITY OF
WATERLOO

# FaRMA: Faster Reward Model for Alignment

- Rashid, Wu, Fan, Li, Kristiadi, Poupart (2025) **Towards Cost-Effective Reward Guided Text Generation**, arxiv.

- Optimization problem:

$$\max_{\theta} E_{(s,a_+,a_-)\in Dataset} \log \sigma\big(r_\theta(s,a_+) - r_\theta(s,a_-)\big)$$

$$\text{Subject to } r_\theta\big(s, a^{1:i}\big) = \max_{a^{i+1:|a|}} r_\theta\big(s, [a^{1:i}, a^{i+1:|a|}]\big) \ \forall s, a, i$$

- In practice: alternate between minimizing two loss functions

  - $L_1(\theta) = -E_{(s,a_+,a_-)\in Dataset} \log \sigma\big(r_\theta(s,a_+) - r_\theta(s,a_-)\big)$

  - $L_2(\theta) = \frac{1}{2} E_{(s,a)\in Dataset, i \leq |a|} \left( r_\theta(s, a^{1:i}) - \max_{a^{i+1:|a|}} r_\theta\big(s, [a^{1:i}, a^{i+1:|a|}]\big) \right)^2$

UNIVERSITY OF
**WATERLOO**

# FaRMA Pseudocode

Repeat

    Repeat for each $(s, \boldsymbol{a}_+, \boldsymbol{a}_-)$ in minibatch

$$L_1(\theta) = \log \sigma\big(r_\theta(\boldsymbol{s}, \boldsymbol{a}_+) - r_\theta(\boldsymbol{s}, \boldsymbol{a}_-)\big)$$

$$\theta \leftarrow \theta - \alpha \nabla L_1(\theta)$$

    Repeat for each $(\boldsymbol{s}, \boldsymbol{a}, i)$ in minibatch

$$L_2(\theta) = \frac{1}{2}\big(r_\theta(\boldsymbol{s}, \boldsymbol{a}^{1:i}) - \max_{a^{i+1}} r_\theta(\boldsymbol{s}, \boldsymbol{a}^{1:i+1})\big)^2$$

$$\theta \leftarrow \theta - \alpha \nabla L_2(\theta)$$

UNIVERSITY OF
WATERLOO

# Empirical Results

## TL;DR Summarization

| Method | LLM | $r \pm$ SE | Time(min) |
|--------|-----|-----------|-----------|
| $\pi_{\text{ref}}$ | frozen | $0.98 \pm 0.18$ | 2 |
| ARGS | frozen | $1.46 \pm 0.16$ | 32 |
| PARGS | frozen | $1.56 \pm 0.19$ | 31 |
| CD | frozen | $1.15 \pm 0.16$ | 29 |
| FaRMA | frozen | $2.05 \pm 0.15$ | 5 |
| CARDS | frozen | $1.73 \pm 0.16$ | 17 |
| DPO | trained | $2.08 \pm 0.18$ | 2 |
| PPO | trained | $2.05 \pm 0.14$ | 2 |

*Table 2.* Avg. reward (over 100 samples) $\pm$ standard error total generation time for the TL;DR summarization task.
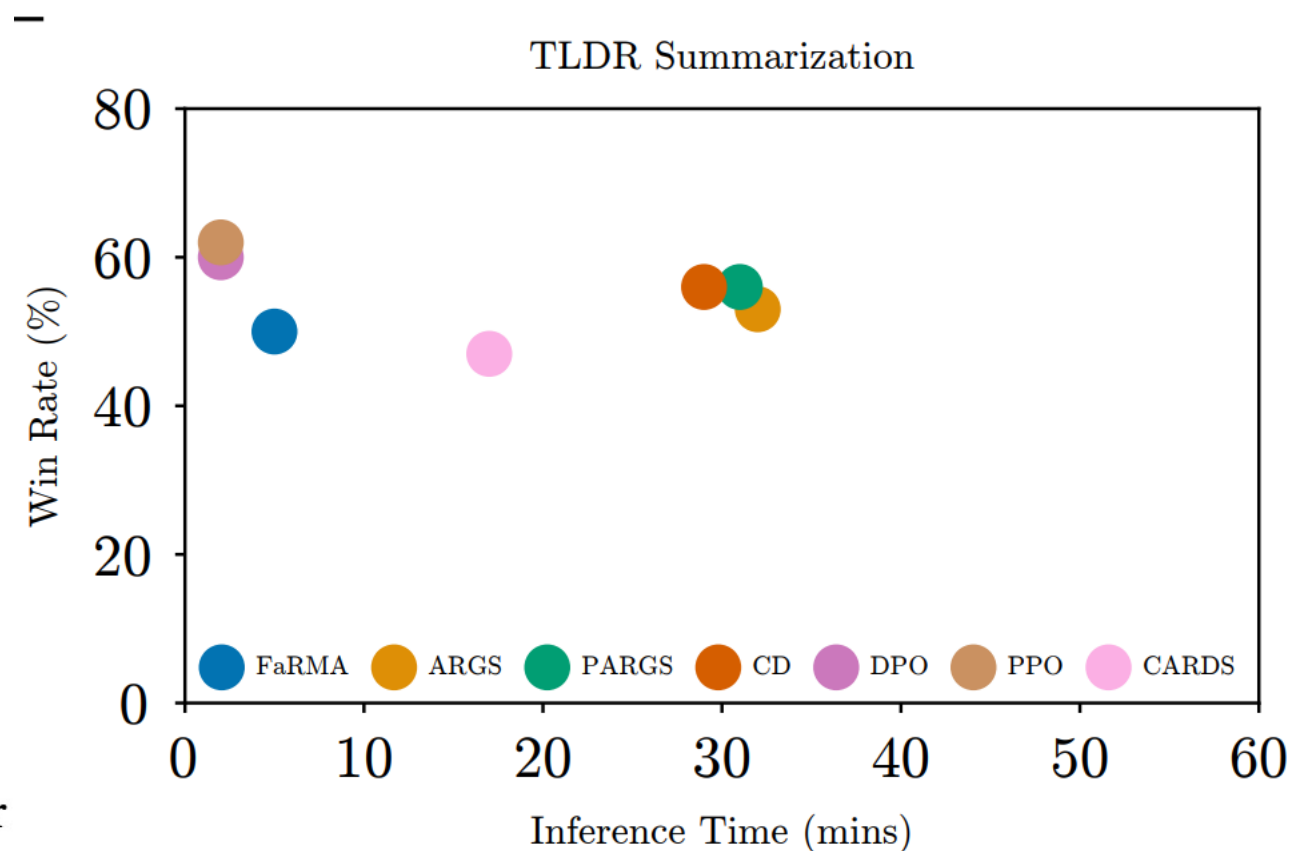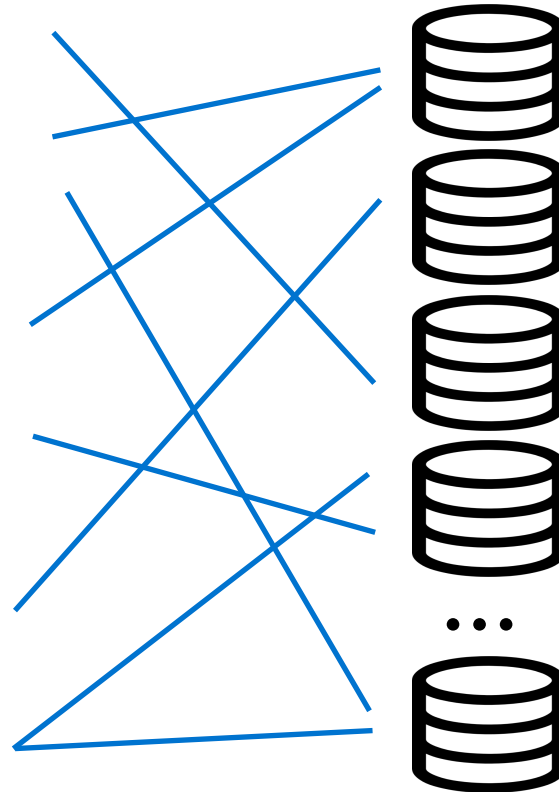


*Figure 2.* GPT4 evaluation on TLDR

UNIVERSITY OF WATERLOO

# Towards Plug-n-play LLMs

Large language models Preference Datasets



Instruction data

Domain data

Fairness data

Toxicity prevention data

...

Client data