

Direct Preference Optimization for Image Generation

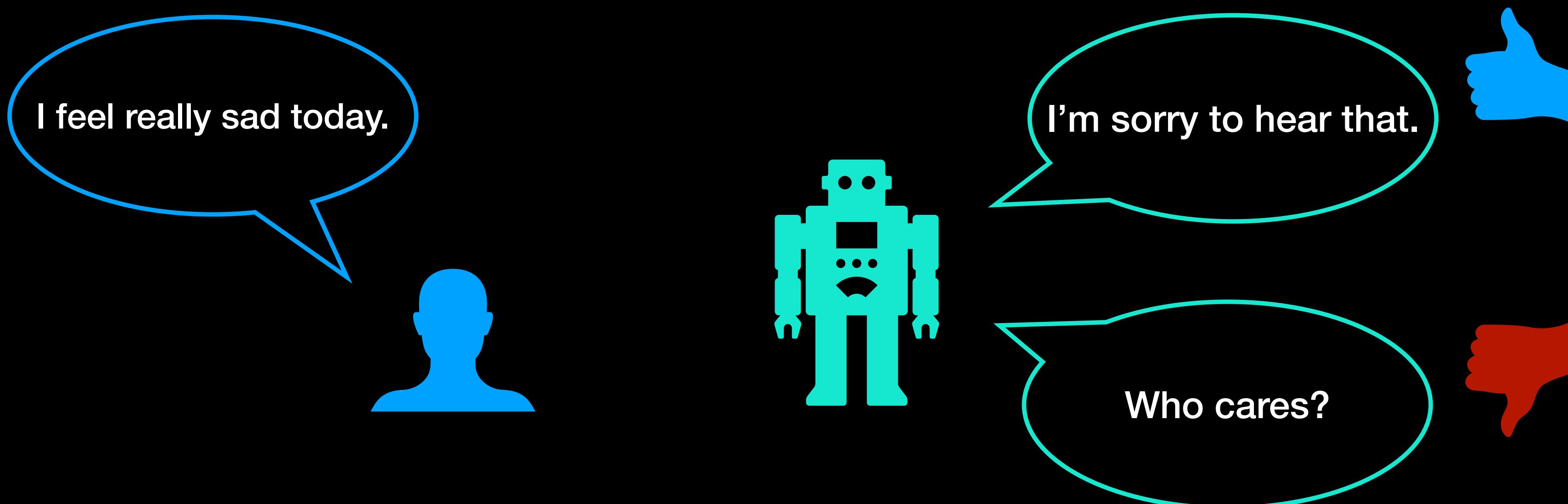
Multimodal Generative AI Theories and Applications

Lecture 10

Gayoung Lee, Jin-Hwa Kim

Aligning models with human preferences

- We want to fine-tune the model \hat{p}_θ from the pre-trained model p_θ to generate samples that are preferred by humans.



Preference optimization

- If we have a *reward model* r that predicts human preference, we can optimize the model \hat{p}_θ using:

Keep the model close to the original model

$$\max_{\hat{p}_\theta} \mathbb{E}_{c \sim \mathcal{D}_c, x_0 \sim \hat{p}_\theta(x_0|c)} [r(c, x_0)] - \beta \text{KL} [\hat{p}_\theta(x_0 | c) || p_\theta(x_0 | c)]$$

Increasing human preference

where $\beta > 0$ controls the regularization. c is a conditional input, such as a text prompt.

Bradley-Terry model

- How can we get a reward model?
- The Bradley-Terry (BT)* model defines the probability that a human prefers x_0^w over x_0^l as:

$$p_{\text{BT}}(x_0^w > x_0^l \mid c) = \sigma(r(c, x_0^w) - r(c, x_0^l))$$

where r is the reward model, and x_0^w and x_0^l denote winning and losing samples respectively.

*Rank Analysis of Incomplete Block Designs, Bradley et al., Biometrika 1952.

Bradley-Terry model

- The reward model r can be parameterized with a neural network ϕ , and learned by maximizing log-likelihood using labeled samples x_0^w and x_0^l :

$$L_{\text{BT}}(\phi) = - \mathbb{E}_{c, x_0^w, x_0^l} \left[\log \sigma \left(r_\phi(c, x_0^w) - r_\phi(c, x_0^l) \right) \right]$$

Preference optimization

- A common strategy is to train a separate reward model (e.g., a pairwise preference classifier) and fine-tune the target model using its outputs.
- However, training a separate reward model adds extra complexity and time, and its inaccuracies can negatively affect preference optimization learning.

Intro: Direct Preference Optimization (DPO)

- In the preference optimization equation, we need the reward model:

$$\max_{\hat{p}_\theta} \mathbb{E}_{c \sim \mathcal{D}_c, x_0 \sim \hat{p}_\theta(x_0 | c)} [r(c, x_0)] - \beta \text{KL} [\hat{p}_\theta(x_0 | c) \| p_\theta(x_0 | c)].$$

- Instead, Direct Preference Optimization (DPO)* proposes a method to directly optimize using human preference data, without requiring a separate reward model.

*Direct Preference Optimization: Your Language Model is Secretly a Reward Model, Rafailov et al., NeurIPS 2023.

Intro: Direct Preference Optimization (DPO)

- *Spoiler:* No need to train a separate reward model r_ϕ . x_0^w and x_0^l are preferred and dispreferred samples, respectively.

$$L_{\text{DPO}}(\theta) = - \mathbb{E}_{\mathbf{c}, x_0^w, x_0^l} \left[\log \sigma \left(\beta \log \frac{\hat{p}_\theta(x_0^w \mid \mathbf{c})}{p_\theta(x_0^w \mid \mathbf{c})} - \beta \log \frac{\hat{p}_\theta(x_0^l \mid \mathbf{c})}{p_\theta(x_0^l \mid \mathbf{c})} \right) \right]$$

Finding analytic global solution

- In the preference optimization equation, we need a reward model:

$$\max_{\hat{p}_\theta} \mathbb{E}_{c \sim \mathcal{D}_c, x_0 \sim \hat{p}_\theta(x_0|c)} [r(c, x_0)] - \beta \text{KL} [\hat{p}_\theta(x_0 | c) \| p_\theta(x_0 | c)]$$

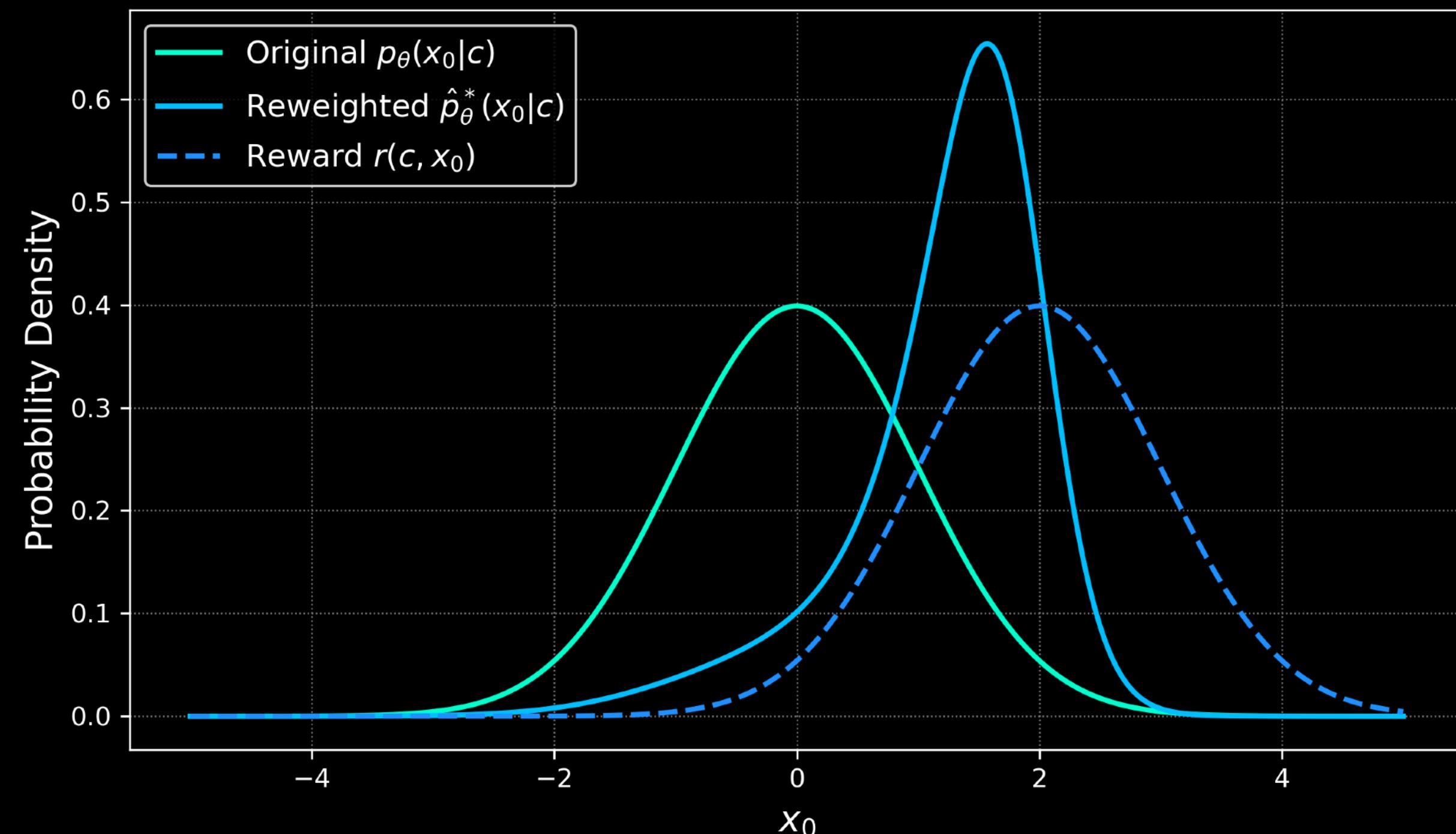
- DPO first shows that there is a global solution in the above equation:

$$\hat{p}_\theta^*(x_0 | c) = p_\theta(x_0 | c) \exp(r(c, x_0)/\beta) / Z(c)$$

where $Z(c) = \sum_{x_0} p_\theta(x_0 | c) \exp(r(c, x_0)/\beta)$ is the partition function.

Visualization of the reweighted distribution

- Visualization of an example: $\hat{p}_\theta^*(x_0 \mid c) = p_\theta(x_0 \mid c)\exp(r(c, x_0)/\beta)/Z(c)$ where $p_\theta \sim \mathcal{N}(0,1)$, reward function $r \sim \mathcal{N}(2,1)$ and $\beta = 0.1$
- We can see that the reward acts as a weighting function.



DPO derivations (1/5)

$$\max \mathbb{E}_{c \sim \mathcal{D}, x_0 \sim \hat{p}_\theta} [r(c, x_0)] - \beta \text{KL}(\hat{p}_\theta(x_0 \mid c) \| p_\theta(x_0 \mid c)) \quad (\beta > 0)$$

$$= \max \mathbb{E}_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 \mid c)} \left[r(c, x_0) - \beta \log \frac{\hat{p}_\theta(x_0 \mid c)}{p_\theta(x_0 \mid c)} \right] \quad (\text{KL Divergence definition})$$

$$\text{KL}(P \parallel Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right]$$

DPO derivations (2/5)

$$\max_{c \sim \mathcal{D}, x_0 \sim \hat{p}_\theta} [r(c, x_0)] - \beta \text{KL}(\hat{p}_\theta(x_0 \mid c) \| p_\theta(x_0 \mid c)) \quad (\beta > 0)$$

$$= \max_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 \mid c)} \left[r(c, x_0) - \beta \log \frac{\hat{p}_\theta(x_0 \mid c)}{p_\theta(x_0 \mid c)} \right] \quad (\text{KL Divergence definition})$$

$$= \min_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 \mid c)} \left[\log \frac{\hat{p}_\theta(x_0 \mid c)}{p_\theta(x_0 \mid c)} - \frac{1}{\beta} r(c, x_0) \right] \quad (\text{rearrangement})$$

DPO derivations (3/5)

$$\max_{c \sim \mathcal{D}, x_0 \sim \hat{p}_\theta} [r(c, x_0)] - \beta \text{KL}(\hat{p}_\theta(x_0 | c) \| p_\theta(x_0 | c)) \quad (\beta > 0)$$

$$= \max_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 | c)} \left[r(c, x_0) - \beta \log \frac{\hat{p}_\theta(x_0 | c)}{p_\theta(x_0 | c)} \right] \quad (\text{KL Divergence definition})$$

$$= \min_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 | c)} \left[\log \frac{\hat{p}_\theta(x_0 | c)}{p_\theta(x_0 | c)} - \frac{1}{\beta} r(c, x_0) \right] \quad (\text{rearrangement})$$

$$= \min_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 | c)} \left[\log \frac{\hat{p}_\theta(x_0 | c)}{p_\theta(x_0 | c) \exp(r(c, x_0)/\beta)} \right] \quad (\text{rearrangement})$$

DPO derivations (3/5)

$$\begin{aligned}
 &= \min \mathbb{E}_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 | c)} \left[\log \frac{\hat{p}_\theta(x_0 | c)}{p_\theta(x_0 | c) \exp(r(c, x_0) / \beta)} \right] \quad (\text{rewritten}) \\
 &= \min \mathbb{E}_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 | c)} \left[\log \frac{\hat{p}_\theta(x_0 | c)}{\frac{1}{Z(c)} p_\theta(x_0 | c) \exp(r(c, x_0) / \beta) Z(c)} \right]
 \end{aligned}$$

where $Z(c) = \sum_{x_0} p_\theta(x_0 | c) \exp(r(c, x_0) / \beta)$

(introducing $Z(c)$ to make the denominator a valid probability)

DPO derivations (4/5)

$$= \min_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 | c)} \left[\log \frac{\hat{p}_\theta(x_0 | c)}{\frac{1}{Z(c)} p_\theta(x_0 | c) \exp(r(c, x_0) / \beta) Z(c)} \right]$$

$$= \min_{c \sim \mathcal{D}} \left[\mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 | c)} \left[\log \frac{\hat{p}_\theta(x_0 | c)}{\hat{p}_\theta^*(x_0 | c)} \right] - \log Z(c) \right]$$

where $\hat{p}_\theta^*(x_0 | c) := \frac{1}{Z(c)} p_\theta(x_0 | c) \exp\left(\frac{1}{\beta} r(c, x_0)\right)$

(substitution and rearrangement)

DPO derivations (5/5)

$$\begin{aligned}
 &= \min_{c \sim \mathcal{D}} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0 | c)} \left[\log \frac{\hat{p}_\theta(x_0 | c)}{\hat{p}_\theta^*(x_0 | c)} \right] - \log Z(c) \quad (\text{rewritten}) \\
 &= \min_{c \sim \mathcal{D}} \left[\text{KL}(\hat{p}_\theta(x_0 | c) \| \hat{p}_\theta^*(x_0 | c)) - \log Z(c) \right]
 \end{aligned}$$

- The minimum value of KL-divergence is achieved if and only if two distributions are the same. Therefore, the optimal $\hat{p}_\theta(x_0 | c)$ is $\hat{p}_\theta^*(x_0 | c)$.

$$\hat{p}_\theta^*(x_0 | c) = \frac{1}{Z(c)} p_\theta(x_0 | c) \exp \left(\frac{1}{\beta} r(c, x_0) \right)$$

Defining the reward function r

- We derived the optimal solution $\hat{p}_\theta(x_0 \mid c) = p_\theta(x_0 \mid c)\exp(r(c, x_0)/\beta)/Z(c)$ where $Z(c) = \sum_{x_0} p_\theta(x_0 \mid c)\exp(r(c, x_0)/\beta)$ is the partition function.
- Then, the above equation can be rewritten as:

$$r(c, x_0) = \beta \log \frac{\hat{p}_\theta(x_0 \mid c)}{p_\theta(x_0 \mid c)} + \beta \log Z(c).$$

Plugging reward into BT model

- We can now plug-in the reward function into Bradley-Terry model objective:

$$L_{\text{BT}}(\phi) = - \mathbb{E}_{c, x_0^w, x_0^l} \left[\log \sigma \left(\textcolor{red}{r}_\phi(c, x_0^w) - \textcolor{red}{r}_\phi(c, x_0^l) \right) \right]$$

$$\textcolor{red}{r}_\theta(c, x_0) = \beta \log \frac{\hat{p}_\theta(x_0 \mid c)}{p_\theta(x_0 \mid c)} + \beta \log Z(c)$$

$$L_{\text{DPO}}(\theta) = - \mathbb{E}_{\mathbf{c}, x_0^w, x_0^l} \left[\log \sigma \left(\beta \log \frac{\hat{p}_\theta(x_0^w \mid \mathbf{c})}{p_\theta(x_0^w \mid \mathbf{c})} - \beta \log \frac{\hat{p}_\theta(x_0^l \mid \mathbf{c})}{p_\theta(x_0^l \mid \mathbf{c})} \right) \right]$$

Using softplus for numerical stability.

$$-\log \sigma(z) = -\log \left(\frac{1}{1 + e^{-z}} \right) = \log(1 + e^{-z})$$

DPO summary

$$L_{\text{DPO}}(\theta) = -\mathbb{E}_{\mathbf{c}, x_0^w, x_0^l} \left[\log \sigma \left(\beta \log \frac{\hat{p}_\theta(x_0^w \mid \mathbf{c})}{p_\theta(x_0^w \mid \mathbf{c})} - \beta \log \frac{\hat{p}_\theta(x_0^l \mid \mathbf{c})}{p_\theta(x_0^l \mid \mathbf{c})} \right) \right]$$

1. Collect a dataset of paired winning and losing samples.
2. Sample one paired sample, x_0^w and x_0^l .
3. Given a condition \mathbf{c} (e.g., text prompt), compare the likelihood ratio between the pretrained model p_θ and the fine-tuned model \hat{p}_θ .

DPO evaluation

- Human evaluation (collecting preferred responses by human)
- Model evaluation (e.g., asking GPT which answer is better)

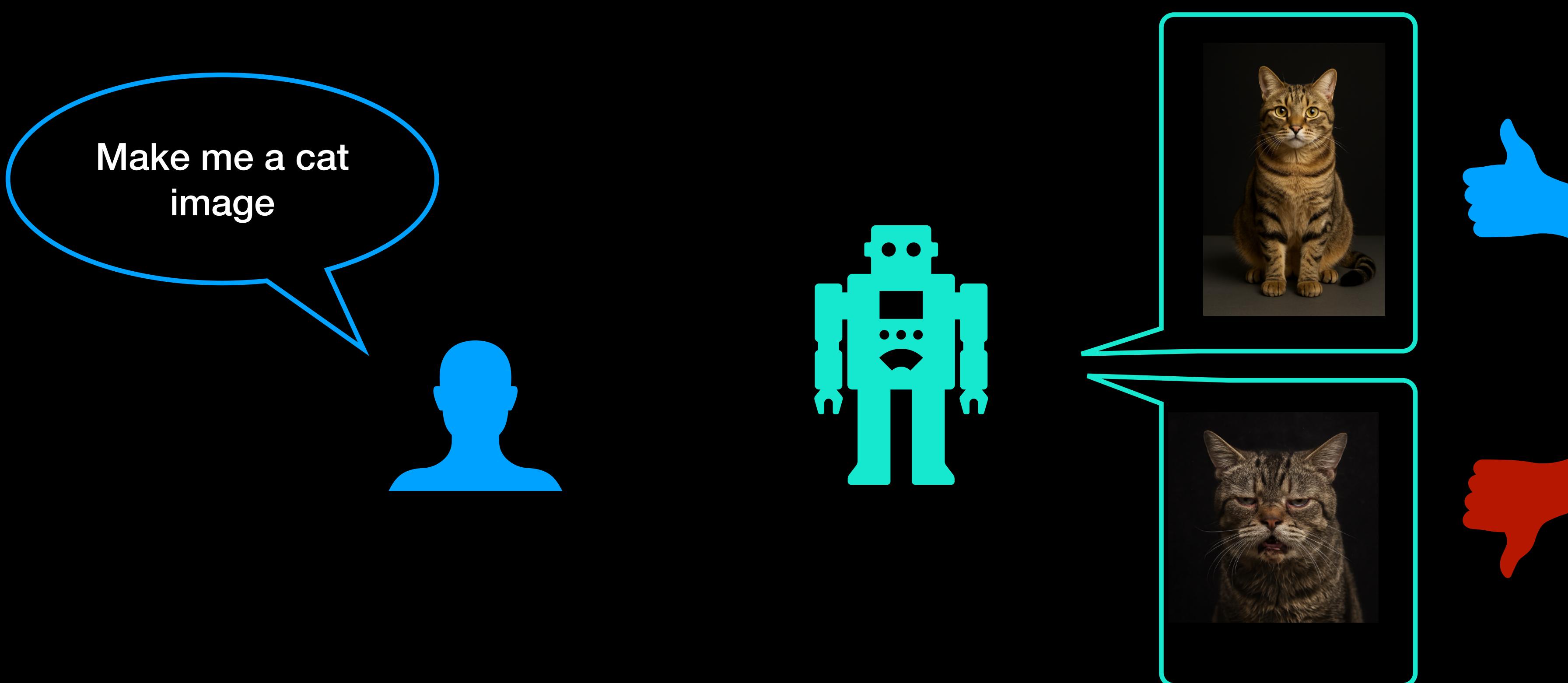
| | DPO | SFT | PPO-1 |
|-------------------|------------|------------|--------------|
| N respondents | 272 | 122 | 199 |
| GPT-4 (S) win % | 47 | 27 | 13 |
| GPT-4 (C) win % | 54 | 32 | 12 |
| Human win % | 58 | 43 | 17 |
| | | | |
| GPT-4 (S)-H agree | 70 | 77 | 86 |
| GPT-4 (C)-H agree | 67 | 79 | 85 |
| H-H agree | 65 | - | 87 |

*Comparing human and GPT-4 win rates on TL;DR summarization samples**
(S: summary, C: concise, SFT: supervised fine-tuning, PPO: proximal policy optimization)

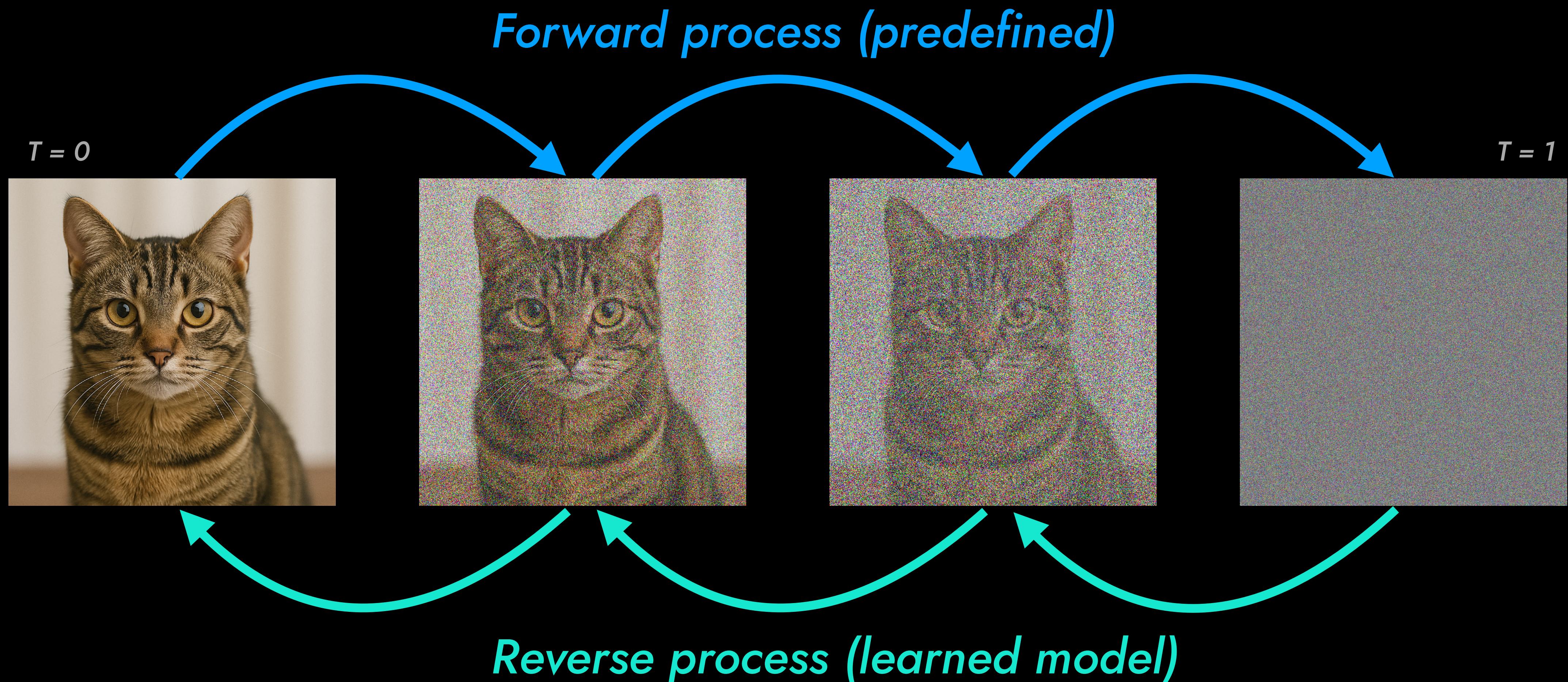
*Direct Preference Optimization: Your Language Model is Secretly a Reward Model, Rafailov et al., NeurIPS 2023.

Preference optimization for images?

- We want to fine-tune the model \hat{p}_θ from the pre-trained model p_θ to generate samples that are preferred by humans.



Diffusion models (Recap)



Diffusion models - DDPM (Recap)

- *Forward process*: Add a noise with a predefined noise schedule.

$$q(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I\right)$$

$$q(x_t | x_0) = \mathcal{N}\left(x_t; \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)I\right)$$

x_0 : original data (e.g., clean image)

x_t : noisy data at the timestep t

β_t : predefined variance schedule, $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$

Diffusion models - DDPM (Recap)

- *Reverse process:* Remove the noise using a learned model θ .

$$p_{\theta}(x_{t-1} \mid x_t) = \mathcal{N}\left(x_{t-1}; \mu_{\theta}(x_t, t), \tilde{\beta}_t I\right)$$

$$\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(x_t, t) \right)$$

$$L = \mathbb{E}_{x_0, t, \epsilon} \left[\|\epsilon - \epsilon_{\theta}(x_t, t)\|^2 \right]$$

$\epsilon_{\theta}(x_t, t)$: Predicted noise, $\epsilon \sim \mathcal{N}(0, I)$, $\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$

Preference optimization for diffusion model

- Recap the preference optimization

$$\max_{\hat{p}_\theta} \mathbb{E}_{c \sim \mathcal{D}_c, x_0 \sim \hat{p}_\theta(x_0 | c)} [r(c, x_0)] - \beta \text{KL} [\hat{p}_\theta(x_0 | c) \| p_\theta(x_0 | c)]$$

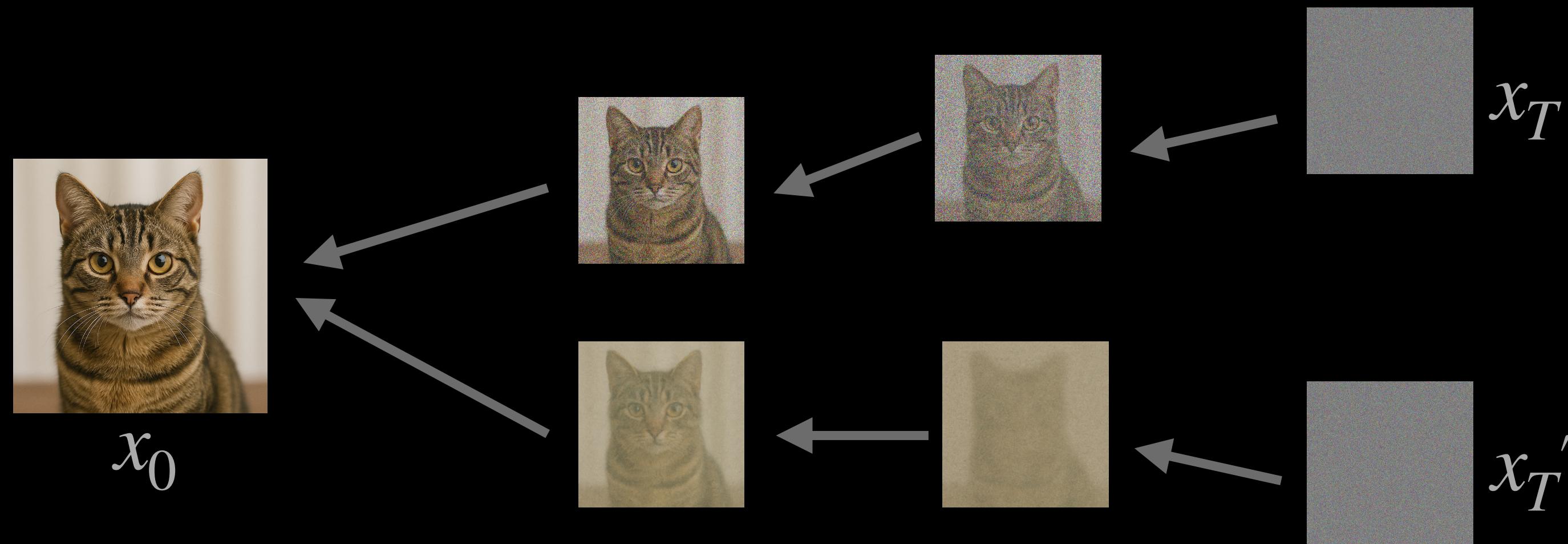
- We want to fine-tune the model \hat{p}_θ from the pre-trained model p_θ to generate samples that are preferred by humans.

Preference optimization for diffusion model

- Recap the preference optimization

$$\max_{\hat{p}_\theta} \mathbb{E}_{c \sim \mathcal{D}_c, x_0 \sim \hat{p}_\theta(x_0|c)} [r(c, x_0)] - \beta \text{KL} [\hat{p}_\theta(x_0 | c) \| p_\theta(x_0 | c)]$$

- To train a diffusion model with $r(c, x_0)$, we would need to consider all possible generation trajectories, which is computationally intractable.



Diffusion-DPO

- Defines a reward function $R(c, x_{0:T})$ for a single trajectory $x_{0:T} = (x_0, x_1, \dots, x_T)$
- The original reward function is then defined as:

$$r(c, x_0) = \mathbb{E}_{x_{1:T} \sim \hat{p}_\theta(x_{1:T}|x_0, c)} [R(c, x_{0:T})]$$

- How can we change the below preference optimization equation with the new reward function $R(c, x_{0:T})$?

$$\max_{\hat{p}_\theta} \mathbb{E}_{c \sim \mathcal{D}_c, x_0 \sim \hat{p}_\theta(x_0|c)} [r(c, x_0)] - \beta \text{KL} [\hat{p}_\theta(x_0 | c) \| p_\theta(x_0 | c)]$$

Diffusion-DPO derivations (1/6)

- Start from preference optimization objective:

$$\min_{\hat{p}_\theta} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0|c)} [-r(c, x_0)] + \beta \text{KL} [\hat{p}_\theta(x_0 | c) \| p_\theta(x_0 | c)]$$

$$\leq \min_{\hat{p}_\theta} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0|c)} [-r(c, x_0)] + \beta \text{KL} [\hat{p}_\theta(\textcolor{red}{x}_{0:T} | c) \| p_\theta(\textcolor{red}{x}_{0:T} | c)]$$

$$\text{KL} [\hat{p}_\theta(x_0 | c) \| p_\theta(x_0 | c)] \leq \text{KL} [\hat{p}_\theta(x_{0:T} | c) \| p_\theta(x_{0:T} | c)]$$

⇒ Proof on the next slide

Diffusion-DPO derivations (1/6)

$$\text{KL}(\hat{p}_\theta(x_{0:T}) \parallel p_\theta(x_{0:T}))$$

$$= \sum_{x_{0:T}} \hat{p}_\theta(x_{0:T}) \log \frac{\hat{p}_\theta(x_{0:T})}{p_\theta(x_{0:T})} \quad (\text{definition})$$

$$= \sum_{x_{0:T}} \hat{p}_\theta(x_{0:T}) \left(\log \frac{\hat{p}_\theta(x_0)}{p_\theta(x_0)} + \log \frac{\hat{p}_\theta(x_{1:T} \mid x_0)}{p_\theta(x_{1:T} \mid x_0)} \right) \quad (\text{chain rule})$$

$$= \sum_{x_0} \hat{p}_\theta(x_0) \sum_{x_{1:T}} \hat{p}_\theta(x_{1:T} \mid x_0) \log \frac{\hat{p}_\theta(x_0)}{p_\theta(x_0)} + \sum_{x_0} \hat{p}_\theta(x_0) \sum_{x_{1:T}} \hat{p}_\theta(x_{1:T} \mid x_0) \log \frac{\hat{p}_\theta(x_{1:T} \mid x_0)}{p_\theta(x_{1:T} \mid x_0)}$$

$$= \sum_{x_0} \hat{p}_\theta(x_0) \log \frac{\hat{p}_\theta(x_0)}{p_\theta(x_0)} + \sum_{x_0} \hat{p}_\theta(x_0) \sum_{x_{1:T}} \hat{p}_\theta(x_{1:T} \mid x_0) \log \frac{\hat{p}_\theta(x_{1:T} \mid x_0)}{p_\theta(x_{1:T} \mid x_0)} \quad (\text{chain rule})$$

Diffusion-DPO derivations (1/6)

$$\begin{aligned}
 & \text{KL}(\hat{p}_\theta(x_{0:T}) \parallel p_\theta(x_{0:T})) \\
 &= \sum_{x_0} \hat{p}_\theta(x_0) \log \frac{\hat{p}_\theta(x_0)}{p_\theta(x_0)} + \sum_{x_0} \hat{p}_\theta(x_0) \sum_{x_{1:T}} \hat{p}_\theta(x_{1:T} \mid x_0) \log \frac{\hat{p}_\theta(x_{1:T} \mid x_0)}{p_\theta(x_{1:T} \mid x_0)} \\
 &= \text{KL}(\hat{p}_\theta(x_0) \parallel p_\theta(x_0)) + \mathbb{E}_{x_0 \sim \hat{p}_\theta} [\text{KL}(\hat{p}_\theta(x_{1:T} \mid x_0) \parallel p_\theta(x_{1:T} \mid x_0))] \\
 &\geq 0
 \end{aligned}$$

Diffusion-DPO derivations (2/6)

- We start from preference optimization objective, and omit $\mathbb{E}_{c \sim \mathcal{D}}$ for simplicity).

$$\begin{aligned}
 & \min_{\hat{p}_\theta} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0|c)} [-r(c, x_0)] + \beta \text{KL} [\hat{p}_\theta(x_0 | c) || p_\theta(x_0 | c)] \\
 & \leq \min_{\hat{p}_\theta} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0|c)} [-\textcolor{red}{r}(c, x_0)] + \beta \text{KL} [\hat{p}_\theta(x_{0:T} | c) || p_\theta(x_{0:T} | c)] \\
 & = \min_{\hat{p}_\theta} \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T}|c)} [-\textcolor{red}{R}(c, x_{0:T})] + \beta \text{KL} [\hat{p}_\theta(x_{0:T} | c) || p_\theta(x_{0:T} | c)]
 \end{aligned}$$

Remember we defined $r(c, x_0) = \mathbb{E}_{x_{1:T} \sim p_\theta(x_{1:T}|x_0, c)} [R(c, x_{0:T})]$.

Diffusion-DPO derivations (3/6)

- Start from preference optimization objective:

$$\begin{aligned}
 & \min_{\hat{p}_\theta} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0|c)} [-r(c, x_0)] + \beta \text{KL} [\hat{p}_\theta(x_0 | c) || p_\theta(x_0 | c)] \\
 & \leq \min_{\hat{p}_\theta} \mathbb{E}_{x_0 \sim \hat{p}_\theta(x_0|c)} [-r(c, x_0)] + \beta \text{KL} [\hat{p}_\theta(x_{0:T} | c) || p_\theta(x_{0:T} | c)] \\
 & = \min_{\hat{p}_\theta} \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T}|c)} [-R(c, x_{0:T})] + \beta \text{KL} [\hat{p}_\theta(x_{0:T} | c) || p_\theta(x_{0:T} | c)] \\
 & = \min_{\hat{p}_\theta} \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T}|c)} [-R(c, x_{0:T})] + \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T}|c)} \left[\beta \log \frac{\hat{p}_\theta(x_{0:T} | c)}{p_\theta(x_{0:T} | c)} \right] \quad (\text{definition}) \\
 & = \min_{\hat{p}_\theta} \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T}|c)} \left[-R(c, x_{0:T}) + \beta \log \frac{\hat{p}_\theta(x_{0:T} | c)}{p_\theta(x_{0:T} | c)} \right] \quad (\text{rearrangement})
 \end{aligned}$$

Diffusion-DPO derivations (4/6)

$$\begin{aligned}
 &= \min_{\hat{p}_\theta} \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T} | c)} \left[-R(c, x_{0:T}) + \beta \log \frac{\hat{p}_\theta(x_{0:T} | c)}{p_\theta(x_{0:T} | c)} \right] \quad (\text{rewritten}) \\
 &= \min_{\hat{p}_\theta} \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T} | c)} \left[\log \frac{\hat{p}_\theta(x_{0:T} | c)}{\frac{1}{Z(c)} p_\theta(x_{0:T} | c) \exp(R(c, x_{0:T})/\beta) Z(c)} \right]
 \end{aligned}$$

where $Z(c) = \sum_{x_{0:T}} p_\theta(x_{0:T} | c) \exp(R(c, x_{0:T})/\beta)$

(introducing $Z(c)$ to make the denominator a valid probability)

Diffusion-DPO derivations (5/6)

$$= \min \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T} | c)} \left[\log \frac{\hat{p}_\theta(x_{0:T} | c)}{\frac{1}{Z(c)} p_\theta(x_{0:T} | c) \exp(R(c, x_{0:T})/\beta) Z(c)} \right] \quad (\text{rewritten})$$

$$= \min \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T} | c)} \left[\log \frac{\hat{p}_\theta(x_{0:T} | c)}{\hat{p}_\theta^*(x_{0:T} | c)} - \log Z(c) \right]$$

where $\hat{p}_\theta^*(x_{0:T} | c) := \frac{1}{Z(c)} p_\theta(x_{0:T} | c) \exp\left(\frac{1}{\beta} R(c, x_{0:T})\right)$

(substitution and rearrangement)

Diffusion-DPO derivations (6/6)

$$\begin{aligned}
 &= \min \mathbb{E}_{x_{0:T} \sim \hat{p}_\theta(x_{0:T} | c)} \left[\log \frac{\hat{p}_\theta(x_{0:T} | c)}{\hat{p}_\theta^*(x_{0:T} | c)} \right] - \log Z(c) \quad (\text{rewritten}) \\
 &= \min \text{KL}(\hat{p}_\theta(x_{0:T} | c) \parallel \hat{p}_\theta^*(x_{0:T} | c)) - \log Z(c) \quad (\text{by KL Divergence definition})
 \end{aligned}$$

- The minimum value of KL-divergence is achieved if and only if two distributions are the same. Therefore, the optimal $\hat{p}_\theta(x_{0:T} | c)$ is $\hat{p}_\theta^*(x_{0:T} | c)$.

$$\hat{p}_\theta^*(x_{0:T} | c) = \frac{1}{Z(c)} p_\theta(x_{0:T} | c) \exp \left(\frac{1}{\beta} R(c, x_{0:T}) \right)$$

Derivations to find r

- We can rearrange the equation to get a reward function.

$$\hat{p}_\theta^*(x_{0:T} \mid c) = \frac{1}{Z(c)} p_\theta(x_{0:T} \mid c) \exp \left(\frac{1}{\beta} R(c, x_{0:T}) \right)$$

$$\Rightarrow Z(c) \frac{\hat{p}_\theta^*(x_{0:T} \mid c)}{p_\theta(x_{0:T} \mid c)} = \exp \left(\frac{1}{\beta} R(c, x_{0:T}) \right)$$

$$\Rightarrow \beta \log \left(Z(c) \frac{\hat{p}_\theta^*(x_{0:T} \mid c)}{p_\theta(x_{0:T} \mid c)} \right) = R(c, x_{0:T})$$

(Plug into the equation $r(c, x_0) = \mathbb{E}_{x_{1:T} \sim p_\theta(x_{1:T} \mid x_0, c)} [R(c, x_{0:T})]$)

$$\Rightarrow r(c, x_0) = \beta \mathbb{E}_{x_{1:T} \sim \hat{p}_\theta(x_{1:T} \mid x_0, c)} \log \frac{\hat{p}_\theta(x_{0:T} \mid c)}{p_\theta(x_{0:T} \mid c)} + \beta \log Z(c)$$

Derivation of the loss (1/13)

- Now, we can now plug-in the reward r into the Bradley-Terry model objective:

$$L_{\text{BT}}(\phi) = - \mathbb{E}_{c, x_0^w, x_0^l} \left[\log \sigma \left(r_\phi(c, x_0^w) - r_\phi(c, x_0^l) \right) \right]$$

where $r(c, x_0) = \beta \mathbb{E}_{x_{1:T} \sim \hat{p}_\theta(x_{1:T} | x_0, c)} \log \frac{\hat{p}_\theta(x_{0:T} | c)}{p_\theta(x_{0:T} | c)} + \beta \log Z(c)$

$$L_{\text{Diffusion-DPO}}(\theta) = - \mathbb{E}_{x_0^w, x_0^l \sim \mathcal{D}} \log \sigma \left(\beta \mathbb{E}_{\substack{x_{1:T}^w \sim \hat{p}_\theta(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim \hat{p}_\theta(x_{1:T}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{0:T}^w)}{p_\theta(x_{0:T}^w)} - \log \frac{\hat{p}_\theta(x_{0:T}^l)}{p_\theta(x_{0:T}^l)} \right] \right)$$

Derivation of the loss (2/13)

- We now have a loss function for the Diffusion-DPO:

$$L_{\text{Diffusion-DPO}}(\theta) = - \mathbb{E}_{x_0^w, x_0^l \sim \mathcal{D}} \log \sigma \left(\beta \mathbb{E}_{\substack{x_{1:T}^w \sim \hat{p}_\theta(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim \hat{p}_\theta(x_{1:T}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{0:T}^w)}{p_\theta(x_{0:T}^w)} - \log \frac{\hat{p}_\theta(x_{0:T}^l)}{p_\theta(x_{0:T}^l)} \right] \right).$$

However, we need to sample $x_{1:T} \sim \hat{p}_\theta(x_{1:T} | x_0)$ which is both inefficient and intractable. Instead, we utilize the forward process $q(x_{1:T} | x_0)$ for approximation.

(Difference) $p_\theta(x_{1:T} | x_0) = p_\theta(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t)$ is a model trajectory, while $q(x_{1:T} | x_0)$ is simple fixed forward process.

Derivation of the loss (3/13)

$$\begin{aligned}
 L_{\text{Diffusion-DPO}}(\theta) &= -\mathbb{E}_{x_0^w, x_0^l \sim \mathcal{D}} \log \sigma \left(\beta \mathbb{E}_{\substack{x_{1:T}^w \sim \hat{p}_\theta(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim \hat{p}_\theta(x_{1:T}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{0:T}^w)}{p_\theta(x_{0:T}^w)} - \log \frac{\hat{p}_\theta(x_{0:T}^l)}{p_\theta(x_{0:T}^l)} \right] \right) \\
 &\approx -\mathbb{E}_{x_0^w, x_0^l \sim \mathcal{D}} \log \sigma \left(\beta \mathbb{E}_{\substack{x_{1:T}^w \sim q(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim q(x_{1:T}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{0:T}^w)}{p_\theta(x_{0:T}^w)} - \log \frac{\hat{p}_\theta(x_{0:T}^l)}{p_\theta(x_{0:T}^l)} \right] \right)
 \end{aligned}$$

Here, x_0^w and x_0^l are sampled from data. $x_{1:T}^w$ and $x_{1:T}^l$ are obtained by simply adding Gaussian noises through the forward process.

Derivation of the loss (4/13)

$$-\log \sigma \left(\beta \mathbb{E}_{\substack{x_{1:T}^w \sim q(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim q(x_{1:T}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{0:T}^w)}{p_\theta(x_{0:T}^w)} - \log \frac{\hat{p}_\theta(x_{0:T}^l)}{p_\theta(x_{0:T}^l)} \right] \right) \quad (\text{rewritten})$$

$$= -\log \sigma \left(\beta \mathbb{E}_{\substack{x_{1:T}^w \sim q(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim q(x_{1:T}^l | x_0^l)}} \left[\sum_{t=1}^T \left(\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right) \right] \right)$$

(using chain rule, markov process)

$$p(x_{0:T}) = p(x_T) \cdot p(x_{T-1} | x_T) \cdot p(x_{T-2} | x_{T-1}) \cdots p(x_0 | p_1)$$

Derivation of the loss (5/13)

$$\begin{aligned}
 & -\log \sigma \left(\beta \mathbb{E}_{\substack{x_{1:T}^w \sim q(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim q(x_{1:T}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{0:T}^w)}{p_\theta(x_{0:T}^w)} - \log \frac{\hat{p}_\theta(x_{0:T}^l)}{p_\theta(x_{0:T}^l)} \right] \right) \\
 & = -\log \sigma \left(\beta \mathbb{E}_{\substack{x_{1:T}^w \sim q(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim q(x_{1:T}^l | x_0^l)}} \left[\sum_{t=1}^T \left(\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right) \right] \right) \\
 & = -\log \sigma \left(\beta \sum_{t=1}^T \mathbb{E}_{\substack{x_{1:T}^w \sim q(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim q(x_{1:T}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right] \right)
 \end{aligned}$$

(Linearity of expectation: $\mathbb{E}(A + B) = \mathbb{E}(A) + \mathbb{E}(B)$)

Derivation of the loss (6/13)

$$-\log \sigma \left(\beta \sum_{t=1}^T \mathbb{E}_{\substack{x_{1:T}^w \sim q(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim q(x_{1:T}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right] \right) \quad (\text{rewritten})$$

$$= -\log \sigma \left(\beta T \mathbb{E}_t \mathbb{E}_{\substack{x_{t-1,t}^w \sim q(x_{t-1,t}^w | x_0^w) \\ x_{t-1,t}^l \sim q(x_{t-1,t}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right] \right)$$

We can marginalize out unused latents from $x_{0:T}$, keeping only $x_{t-1,t} = (x_{t-1}, x_t)$.

$$\text{c.f. } \mathbb{E}_{x,y}[f(x)] = \int f(x) p(x, y) dx dy = \int f(x) dx p(x, y) dy = \int f(x) p(x) dx = \mathbb{E}_x[f(x)]$$

Derivation of the loss (7/13)

$$-\log \sigma \left(\beta \sum_{t=1}^T \mathbb{E}_{\substack{x_{1:T}^w \sim q(x_{1:T}^w | x_0^w) \\ x_{1:T}^l \sim q(x_{1:T}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right] \right)$$

$$= -\log \sigma \left(\beta T \mathbb{E}_t \mathbb{E}_{\substack{x_{t-1,t}^w \sim q(x_{t-1,t}^w | x_0^w) \\ x_{t-1,t}^l \sim q(x_{t-1,t}^l | x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right] \right)$$

$$= -\log \sigma \left(\beta T \mathbb{E}_{t, x_t^w \sim q(x_t | x_0^w)} \mathbb{E}_{\substack{x_{t-1}^w \sim q(x_{t-1} | x_t^w, x_0^w) \\ x_t^l \sim q(x_t | x_0^l) \\ x_{t-1}^l \sim q(x_{t-1} | x_t^l, x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right] \right)$$

Diffusion-DPO (Derivation of the loss)

$$\begin{aligned}
 & -\log \sigma \left(\beta T \mathbb{E}_{t, x_t^w \sim q(x_t | x_0^w)} \mathbb{E}_{x_{t-1}^w \sim q(x_{t-1} | x_t^w, x_0^w)} \right. \\
 & \quad \left. \left[\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right] \right) \\
 & \leq -\mathbb{E}_{t, x_t^w \sim q(x_t | x_0^w)} \log \sigma \left(\beta T \mathbb{E}_{x_{t-1}^w \sim q(x_{t-1} | x_t^w, x_0^w)} \right. \\
 & \quad \left. \left[\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right] \right)
 \end{aligned}$$

By *Jensen's inequality* with the convexity of $-\log$: $\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]$ if φ is a convex function.

Diffusion-DPO (Derivation of the loss)

$$-\mathbb{E}_{\substack{x_t^w \sim q(x_t | x_0^w) \\ x_t^l \sim q(x_t | x_0^l)}} \log \sigma \left(\beta T \mathbb{E}_{\substack{x_{t-1}^w \sim q(x_{t-1} | x_t^w, x_0^w) \\ x_{t-1}^l \sim q(x_{t-1} | x_t^l, x_0^l)}} \left[\log \frac{\hat{p}_\theta(x_{t-1}^w | x_t^w)}{p_\theta(x_{t-1}^w | x_t^w)} - \log \frac{\hat{p}_\theta(x_{t-1}^l | x_t^l)}{p_\theta(x_{t-1}^l | x_t^l)} \right] \right)$$

We can write the above equation using KL-divergence:

$$-\mathbb{E}_{\substack{x_t^w \sim q(x_t | x_0^w) \\ x_t^l \sim q(x_t | x_0^l)}} \log \sigma \left(\beta T \left(\text{KL}(q(x_{t-1}^w | x_{0,t}^w) \| \hat{p}_\theta(x_{t-1}^w | x_t^w)) - \text{KL}(q(x_{t-1}^w | x_{0,t}^w) \| p_\theta(x_{t-1}^w | x_t^w)) \right. \right. \\ \left. \left. - \text{KL}(q(x_{t-1}^l | x_{0,t}^l) \| \hat{p}_\theta(x_{t-1}^l | x_t^l)) + \text{KL}(q(x_{t-1}^l | x_{0,t}^l) \| p_\theta(x_{t-1}^l | x_t^l)) \right) \right)$$

$$\text{KL}(P \| Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right]$$

Diffusion-DPO (Derivation of the loss)

$$-\mathbb{E}_{\substack{x_t^w \sim q(x_t | x_0^w) \\ x_t^l \sim q(x_t | x_0^l)}} \log \sigma \left(\begin{array}{c} \beta T (\text{KL}(q(x_{t-1}^w | x_{0,t}^w) \| \hat{p}_\theta(x_{t-1}^w | x_t^w)) + \text{KL}(q(x_{t-1}^w | x_{0,t}^w) \| p_\theta(x_{t-1}^w | x_t^w))) \\ -\text{KL}(q(x_{t-1}^l | x_{0,t}^l) \| \hat{p}_\theta(x_{t-1}^l | x_t^l)) + \text{KL}(q(x_{t-1}^l | x_{0,t}^l) \| p_\theta(x_{t-1}^l | x_t^l)) \end{array} \right)$$

Recap the diffusion model:

$$p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \tilde{\beta}_t I)$$

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right)$$

KL Divergence of Gaussian distributions:

$$D_{\text{KL}}(\mathcal{N}(\mu_0, \Sigma_0) \| \mathcal{N}(\mu_1, \Sigma_1)) = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1} (\mu_1 - \mu_0) - d + \ln \frac{|\Sigma_1|}{|\Sigma_0|} \right)$$

Diffusion-DPO (Derivation of the loss)

We want to calculate $\text{KL} \left(q(x_{t-1}^w \mid x_{0,t}^w) \parallel \hat{p}_\theta(x_{t-1}^w \mid x_t^w) \right)$

Recap the diffusion model

$$p_\theta(x_{t-1} \mid x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \tilde{\beta}_t I) \quad \mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right)$$

And can calculate $q(x_{t-1} \mid x_{0,t})$ since the forward process has a closed form.

$$q(x_{t-1} \mid x_t, x_0) = \mathcal{N}(x_{t-1}; \mu_q(x_t, x_0), \tilde{\beta}_t I)$$

$$\mu_q(x_t, x_0) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right), \quad \epsilon \sim \mathcal{N}(0, I)$$

Diffusion-DPO (Derivation of the loss)

$$D_{\text{KL}}(\mathcal{N}(\mu_0, \Sigma_0) \parallel \mathcal{N}(\mu_1, \Sigma_1)) = \frac{1}{2} \left(\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1} (\mu_1 - \mu_0) - d + \ln \frac{|\Sigma_1|}{|\Sigma_0|} \right)$$

In our diffusion model, the means are:

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right), \quad \mu_q(x_t, x_0) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right).$$

The variances of both distribution are fixed as $\tilde{\beta}_t I$. We can simplify the divergence as:

$$\text{KL} \left(q(x_{t-1}^w \mid x_{0,t}^w) \parallel \hat{p}_\theta(x_{t-1}^w \mid x_t^w) \right) \propto \|\epsilon^w - \epsilon_\theta(x_t^w, t)\|.$$

Diffusion-DPO (Derivation of the loss)

$$\begin{aligned}
L_{\text{DPO-Diffusion}}(\theta) &\leq \\
&- \mathbb{E}_{\substack{x_t^w \sim q(x_t | x_0^w) \\ x_t^l \sim q(x_t | x_0^l)}} \log \sigma \left(\begin{array}{c} \beta T (\text{KL}(q(x_{t-1}^w | x_{0,t}^w) \| \hat{p}_\theta(x_{t-1}^w | x_t^w)) - \text{KL}(q(x_{t-1}^w | x_{0,t}^w) \| p_\theta(x_{t-1}^w | x_t^w)) \\ - \text{KL}(q(x_{t-1}^l | x_{0,t}^l) \| \hat{p}_\theta(x_{t-1}^l | x_t^l)) + \text{KL}(q(x_{t-1}^l | x_{0,t}^l) \| p_\theta(x_{t-1}^l | x_t^l)) \end{array} \right) \\
&= - \mathbb{E}_{t, \epsilon^w, \epsilon^l} \log \sigma \left(-\beta T \omega(\lambda_t) \left[\|\epsilon^w - \hat{\epsilon}_\theta(x_t^w, t)\|^2 - \|\epsilon^w - \epsilon_\theta(x_t^w, t)\|^2 \right. \right. \\
&\quad \left. \left. - \|\epsilon^l - \hat{\epsilon}_\theta(x_t^l, t)\|^2 + \|\epsilon^l - \epsilon_\theta(x_t^l, t)\|^2 \right] \right)
\end{aligned}$$

$$x_t^w \sim q(x_t | x_0^w), \quad x_t^l \sim q(x_t | x_0^l), \quad t \sim [0,1], \quad \epsilon^w, \epsilon^l \sim \mathcal{N}(0, I)$$

Diffusion-DPO summary

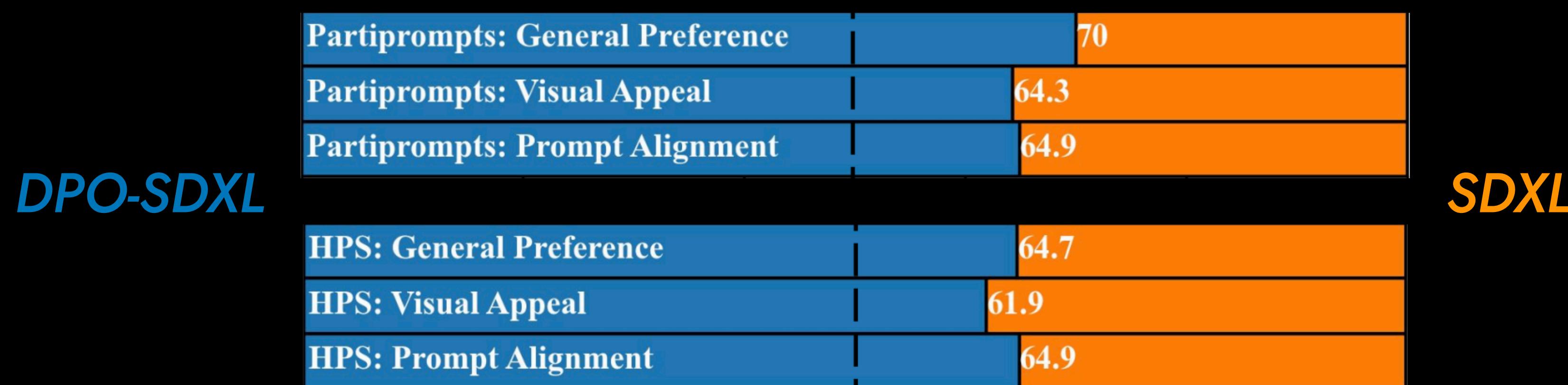
$$L_{\text{DPO-Diffusion}}(\theta) \leq -\mathbb{E}_{t, \epsilon^w, \epsilon^l} \log \sigma \left(-\beta T \omega(\lambda_t) \left[\|\epsilon^w - \hat{\epsilon}_\theta(x_t^w, t)\|^2 - \|\epsilon^w - \epsilon_\theta(x_t^w, t)\|^2 \right. \right. \\ \left. \left. - \|\epsilon^l - \hat{\epsilon}_\theta(x_t^l, t)\|^2 + \|\epsilon^l - \epsilon_\theta(x_t^l, t)\|^2 \right] \right)$$

$$x_t^w \sim q(x_t | x_0^w), x_t^w \sim q(x_t | x_0^l), t \sim [0,1], \epsilon^w, \epsilon^l \sim \mathcal{N}(0, I)$$

1. Collect a dataset of paired winning and losing samples.
2. Sample one paired sample x_0^w and x_0^l .
3. Select a timestep t .
4. Generate noisy data x_t^w and x_t^l using the forward process q .
5. Compare the model predicted noise $\hat{\epsilon}_\theta$ with the generated noisy data.

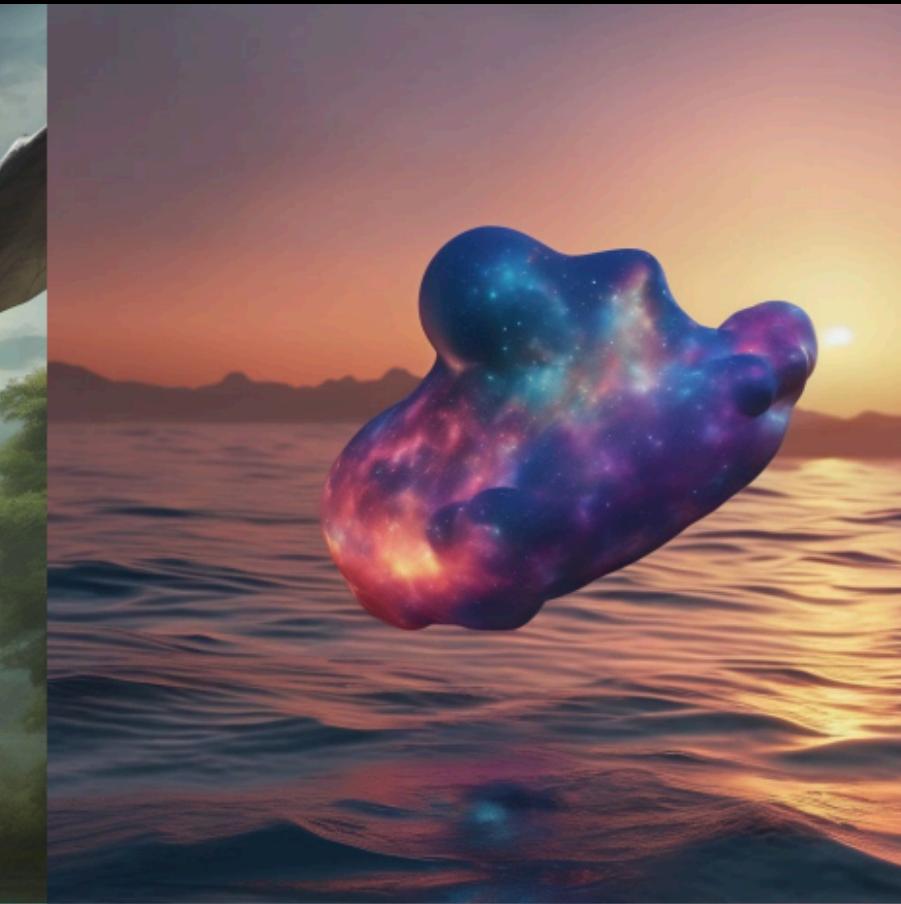
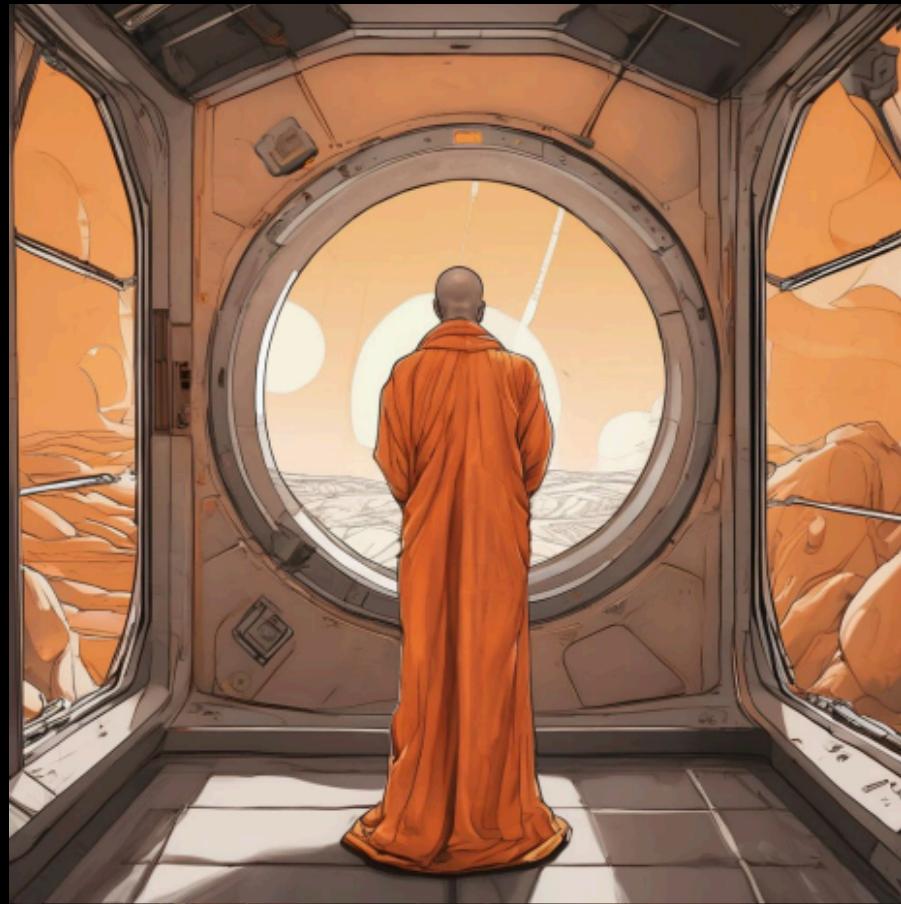
Diffusion-DPO evaluation

- Human evaluation to compare the generations under three criteria:
 - Q1. General preference (Which image do you prefer given the prompt?)
 - Q2. Visual appeal (Which image is more visually appealing regardless of the prompt?)
 - Q3. Prompt alignment (Which image better fits the text description?)

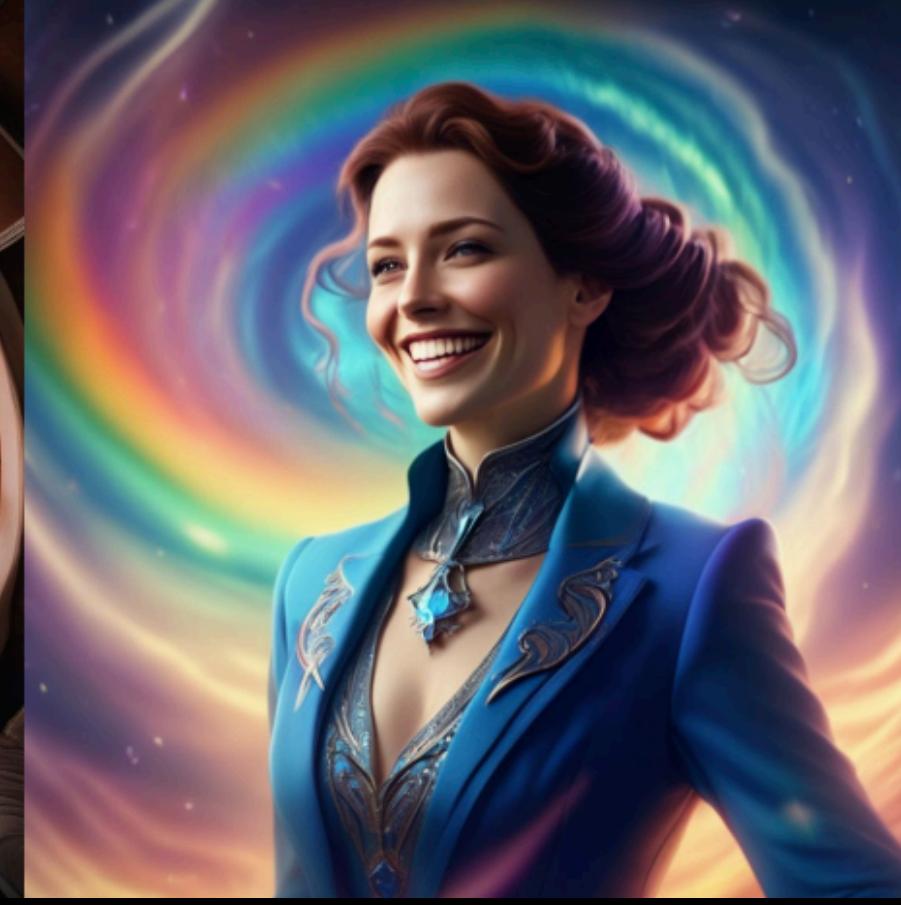


Diffusion-DPO examples

*Stable
Diffusion XL
(SDXL)*



DPO-SDXL



Application: Diffusion-DPO for safe generation

- As image generation becomes more realistic and widely used, preventing NSFW (Not Safe For Work) content has become increasingly important.
- We will introduce DUO* that utilizes Diffusion-DPO for safe image generation.

Nudity

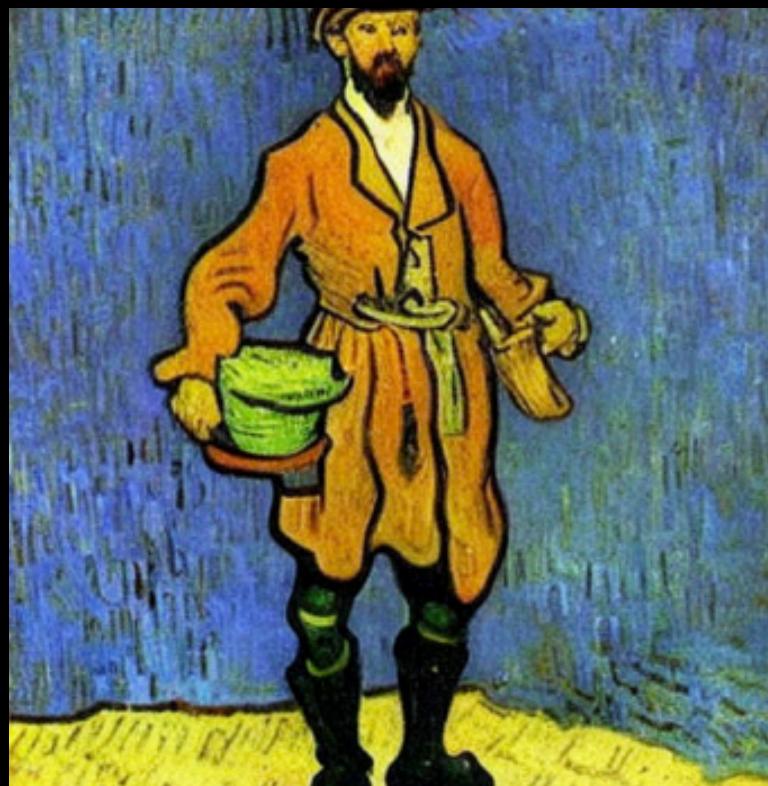
Violence

Hate

Illegality

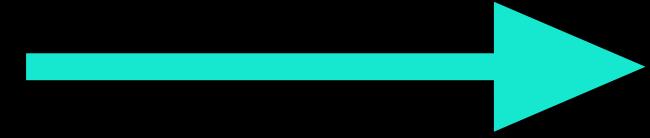
Previous methods for safe image generation

- Most of the previous safety methods depend on malicious text prompts.
 - For example, the model is fine-tuned to generate safe images even when pre-defined malicious keywords (e.g., sexual content, violence) are given.
- Erasing Concepts from Diffusion Models. (Gandikota et al., ICCV 2023)



Result of the original model

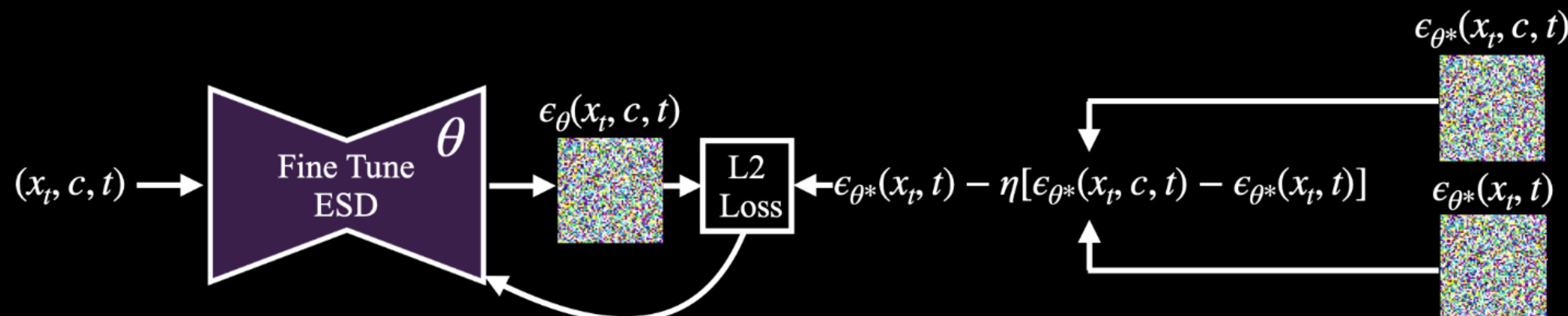
Erase “Van Gogh” from the model



Result of the edited model

Erasing Concepts from Diffusion Models

- Erasing Concepts from Diffusion Models. (R. Gandikota et al. ICCV 2023)
 - They fine-tune the model to denoise a noised image x_t given a malicious prompt c as it does without the prompt.
 - By subtracting the predicted noises with and without the prompt, we can find the direction related to the malicious prompt and more suppress the model using it.

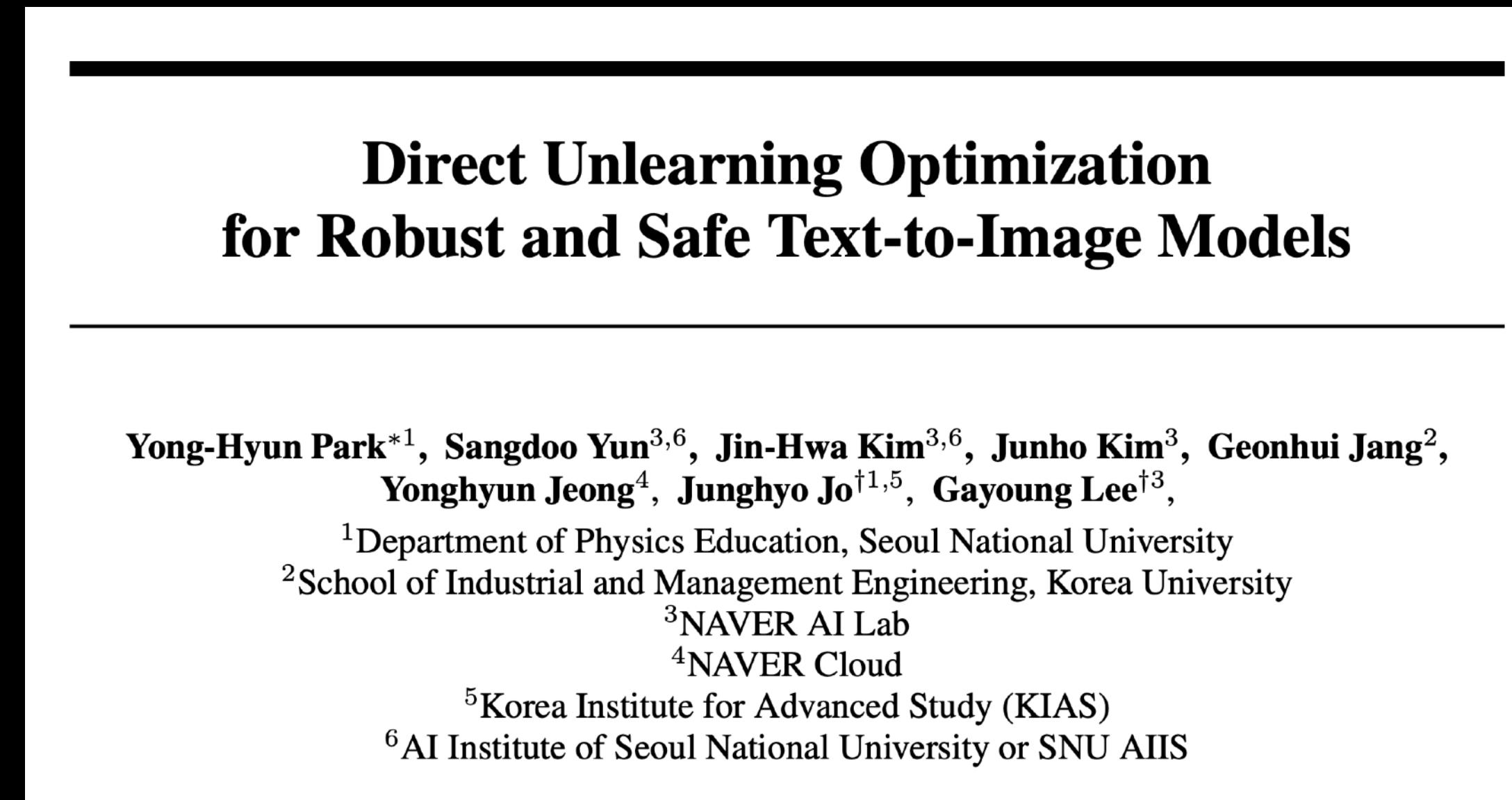


Problem of the previous method

- However, this approach is often prone to adversarial prompt attacks.
 - Adversarial prompt: a seemingly harmless prompt that makes the diffusion model generate unsafe images (e.g., s#i2cn!).
 - This is because the approach fine-tunes the model only for specific prompts, and several studies* have shown that by bypassing those prompts, unsafe images can still be generated.

Direct unlearning optimization (DUO)

- To prevent adversarial prompt attacks, the model needs to remove features related to unsafe concepts, regardless of the given prompt.
- How can we remove harmful features from the model without using text prompts?



Rethinking unlearning as preference optimization

$$L_{\text{DPO-Diffusion}}(\theta) \leq -\mathbb{E}_{t,\epsilon^w,\epsilon^l} \log \sigma \left(-\beta T \omega(\lambda_t) \left[(\|\epsilon^w - \hat{\epsilon}_\theta(x_t^w, t)\|^2 - \|\epsilon^w - \epsilon_\theta(x_t^w, t)\|^2) \right. \right. \\ \left. \left. - (\|\epsilon^l - \hat{\epsilon}_\theta(x_t^l, t)\|^2 - \|\epsilon^l - \epsilon_\theta(x_t^l, t)\|^2) \right] \right)$$

$$x_t^w \sim q(x_t \mid x_0^w), x_t^w \sim q(x_t \mid x_0^l), t \sim [0,1], \epsilon^w, \epsilon^l \sim \mathcal{N}(0, I)$$

- x_0^w is a safe image sample and x_0^l is an unsafe image sample in our task.
- Since our method unlearns harmful features without relying on text prompts, it remains robust even when an adversarial prompt is given.

DUO - synthesize paired data

- For preference optimization, we need paired data, (x_0^w, x_0^l) , where x_0^w is a set of safe images and x_0^l is a set of unsafe images.
- Unsafe images often contain concepts that are unrelated to the unsafe content.
- How can we selectively remove the unsafe concepts while keeping the rest intact?



Concepts to be removed: Naked, Nudity...

Concepts to be kept: Forest, Human, Girl...

* High-Resolution Image Synthesis with Latent Diffusion Models, Rombach et al., CVPR 2022

** SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, Meng et al., ICLR 2022.

DUO - synthesize paired data

- For preference optimization, we need paired data, (x_0^w, x_0^l) , where x_0^w is a set of safe images and x_0^l is a set of unsafe images.
- Unsafe images often contain concepts that are unrelated to the unsafe content.
- How can we selectively remove the unsafe concepts while keeping the rest intact?



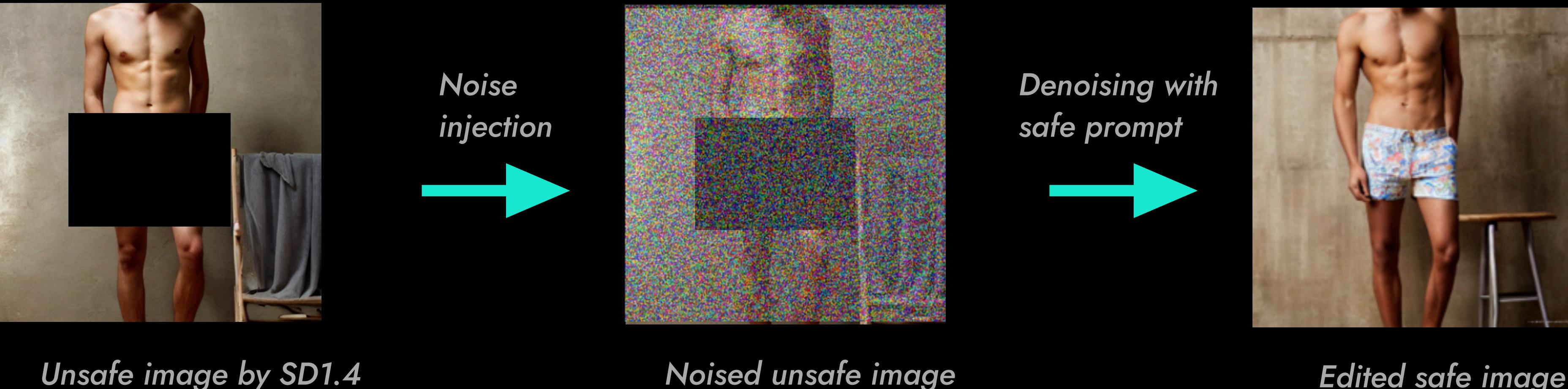
Concepts to be removed: Naked, Nudity...

* High-Resolution Image Synthesis with Latent Diffusion Models, Rombach et al., CVPR 2022

** SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, Meng et al., ICLR 2022.

DUO - synthesize paired data

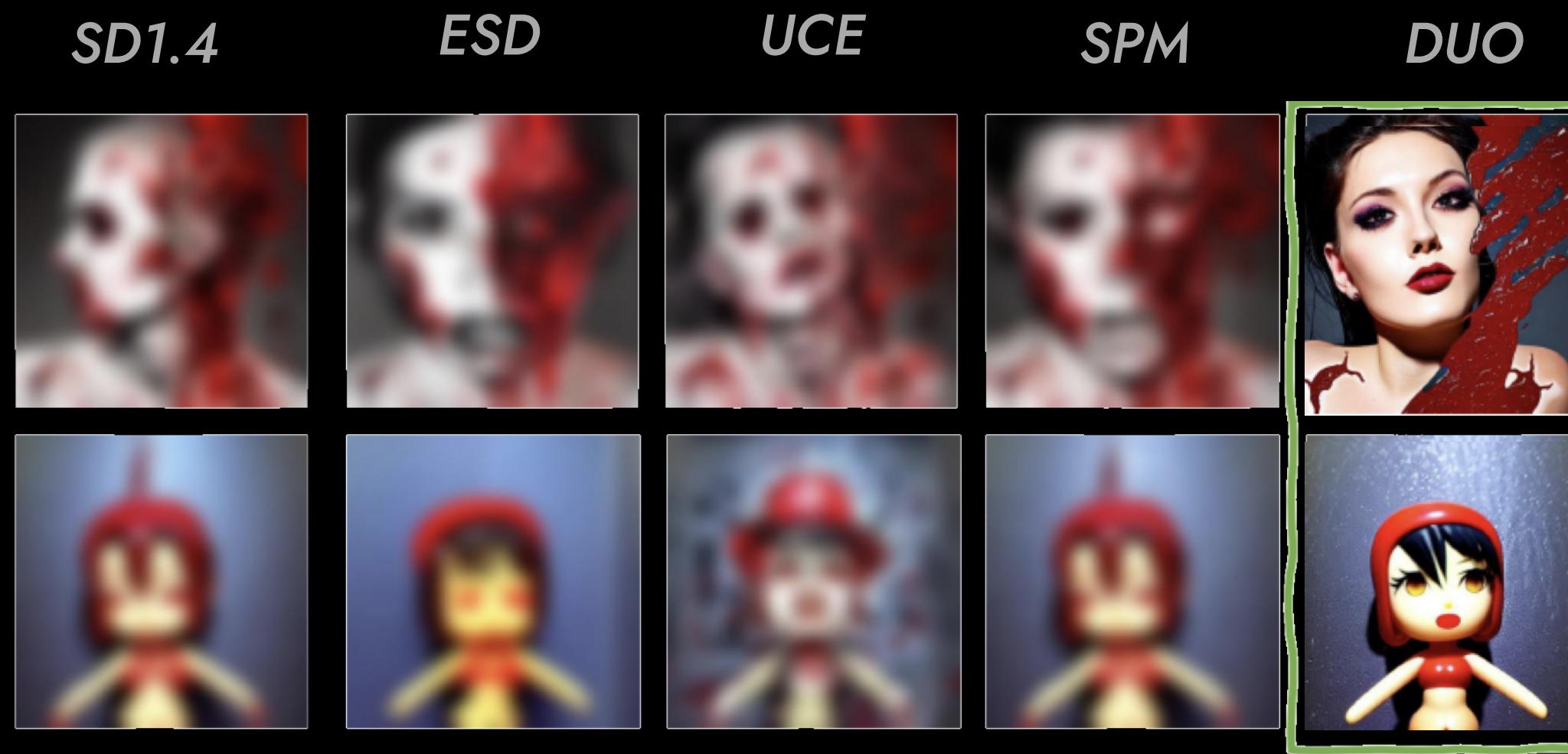
- For preference optimization, we need paired data, (x_0^w, x_0^l) , where x_0^w is a set of safe images and x_0^l is a set of unsafe images.
- We utilize synthetic data using Stable Diffusion* 1.4 and SDEdit**.
- SDEdit** is an image editing method via noise injection and denoising.



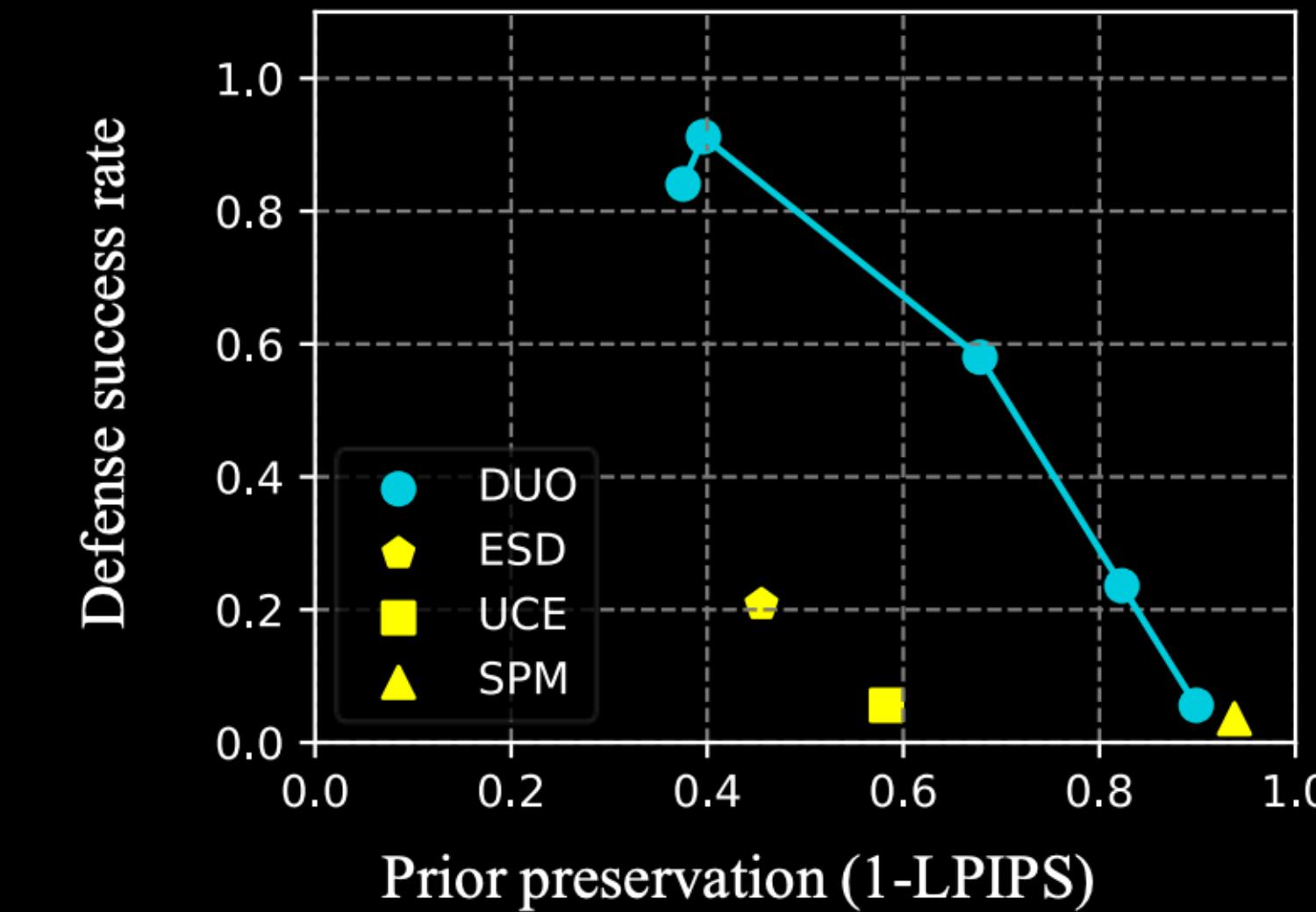
* High-Resolution Image Synthesis with Latent Diffusion Models, Rombach et al., CVPR 2022

** SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, Meng et al., ICLR 2022.

DUO results



Results using unsafe (violence) prompts after Ring-A-Bell attack



Defense rate and prior preservation with Ring-A-Bell attack

- While previous methods that use prompts for training struggles with adversarial prompt attack (Ring-A-Bell*), our method still prevents unsafe image generation.

DUO results

SD1.4



SD1.4 + DUO



- Results using prompts from the MS COCO validation after removing nudity.
- For unrelated prompts, the fine-tuned model keeps the original performance.

Summary

- *Preference optimization* refers to training (usually fine-tuning) a model to align with human preferences as specified by a reward model.
- *Direct Preference Optimization (DPO)* removes the need for an explicit reward model by directly optimizing the loss using paired winning and losing samples.
- *Diffusion-DPO* extends DPO to diffusion models by defining a reward function over multi-timestep trajectories and approximating the DPO loss in the diffusion setting.
- *Direct Unlearning Optimization (DUO)* leverages Diffusion-DPO for unlearning by framing the unlearning task as dispreferring unsafe images.

Thank you!

- Any question?