

# Controllable Generative Diffusion Models: A Constrained Optimization View

**Dongsheng Ding**

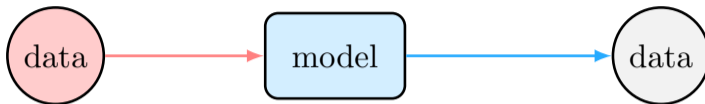
[dongshed@utk.edu](mailto:dongshed@utk.edu)



THE UNIVERSITY OF  
TENNESSEE  
KNOXVILLE

**UC San Diego; Mar. 12, 2026**

# The reality of deep generative modeling



goal – likelihood

## ■ REQUIREMENTS

fairness



safety



robustness



harmlessness



# The risk of generative models

**Humans Are Biased. Generative AI Is Even Worse**

Bloomberg, JUN 9, 2023

**When AI-Powered Humanoid Robots Make Bad Choices**

AI Business, JAN 5, 2026

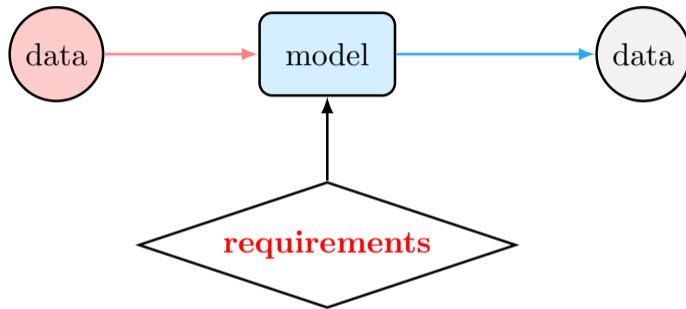
**Their teenage sons died by suicide. Now, they are sounding an alarm about AI chatbots**

NPR, SEP 19, 2025

**AI chatbots might be sabotaging women by advising them to ask for lower salaries, study says**

New York Post, JUL 29, 2025

# Requirement-driven deep generative modeling



fairness



safety



robustness



harmlessness



# Motivating application: Healthcare

## ■ ONCOLOGY IMAGING



National Cancer Institute

minimize model negative likelihood

subject to distance to minority distribution  $\leq$  threshold

**REAL-WORLD CHALLENGE**

**Constraint satisfaction**

## OBJECTIVE

Find a **deep generative model** that  
**minimizes** a performance metric  
**subject to** a constraint on  
another performance metric

# Outline

- CONSTRAINED GENERATIVE MODELS
- PART I PRE-TRAINING OF DIFFUSION MODELS
  - ★ constrained diffusion models
  - ★ dual method & optimality
- PART II POST-TRAINING OF DIFFUSION MODELS
  - ★ alignment with constraints
  - ★ dual method & optimality
- EMPIRICAL STUDY
- SUMMARY & OUTLOOK

# **CONSTRAINED GENERATIVE MODELS**

constrained distribution optimization

# Generative modeling

## ■ DATA GENERATION



data  $x \sim q(\cdot)$



new data  $x \sim p(\cdot)$

★ distribution optimization

minimize <sub>$p$</sub>   $D_{\text{diff}}(p; q)$

$D_{\text{diff}}$  – distribution mismatch

$p$  – trainable distribution model

# Generation requirements

## ■ UNBIASEDNESS

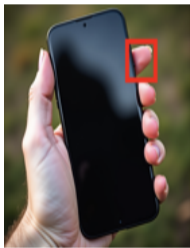
★ stereotypes



dishwashers

## ■ AESTHETICITY

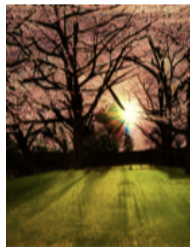
★ artifacts



hand holding

## ■ COMPOSITIONALITY

★ compositional concepts



cherry tree AND sun dog

**Generation has to be controllable satisfying requirements**

# Constrained distribution optimization

minimize <sub>$p$</sub>   $D_{\text{diff}}(p; q)$

subject to

$$D_{\text{diff}}(p; q^i) \leq \epsilon_i \quad \text{for } i = 1, \dots, m$$

→ **proximity**

$$\mathbb{E}_{x \sim p(\cdot)} [r_j(x)] \geq b_j \quad \text{for } j = 1, \dots, n$$

→ **reward**

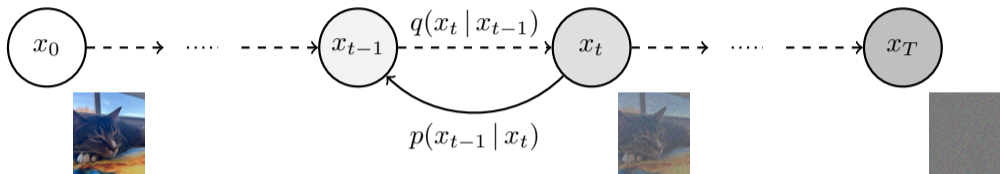


**functional constraints**

★ limit the distribution space to **inequality constraints**

e.g., unbiased model, aestheticity, compositional model

# Generative diffusion models



$q(x_t | x_{t-1})$  – forward process

$p(x_{t-1} | x_t)$  – backward process

## ■ DISTRIBUTION OPTIMIZATION

$$\underset{p}{\text{minimize}} \quad D_{\text{KL}}(q(x_{0:T}) \parallel p(x_{0:T}))$$

# Pre-training: Constrained diffusion models

KD<sup>†</sup>R, NeurIPS '24

**Original data:**  $x_0 \sim q(\cdot)$

**Desirable data:**  $x_0 \sim q^i(\cdot)$  for  $i = 1, \dots, m$

## ■ CONSTRAINED DISTRIBUTION OPTIMIZATION

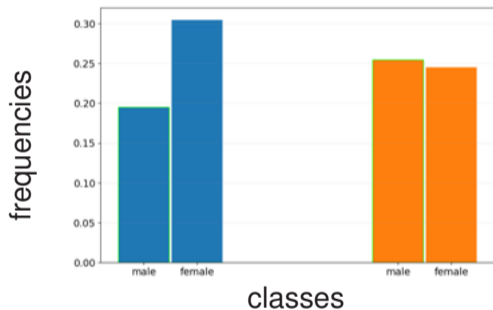
$$\underset{p}{\text{minimize}} \quad D_{\text{KL}}(q(x_{0:T}) \parallel p(x_{0:T}))$$

$$\text{subject to} \quad D_{\text{KL}}(q^i(x_{0:T}) \parallel p(x_{0:T})) \leq \epsilon_i \quad \text{for } i = 1, \dots, m$$

★ minority-promoting distribution  $q^i$

# Constraints balance representativeness

KD<sup>†</sup>R, NeurIPS '24



unconstrained



constrained

★ unconstrained model (■)

★ constrained model (■)

# Post-training: Alignment with constraints

KHD<sup>†</sup>R, NeurIPS '25

**Pretrained model:**  $q(x_{0:T})$

**Reward models:**  $r_i(x_0)$  for  $i = 1, \dots, m$

## ■ CONSTRAINED DISTRIBUTION OPTIMIZATION

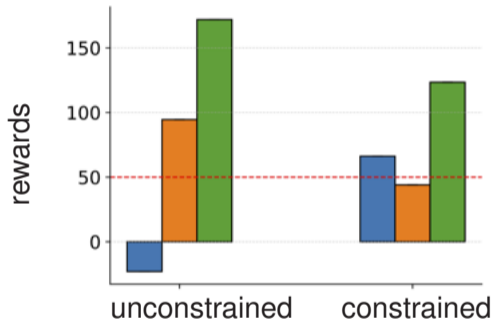
$$\underset{p}{\text{minimize}} \quad D_{\text{KL}}(p(x_{0:T}) \parallel q(x_{0:T}))$$

$$\text{subject to} \quad \mathbb{E}_{x_0 \sim p(\cdot)} [r_i(x_0)] \geq b_i \text{ for } i = 1, \dots, m$$

★ aestheticity-constrained reward  $r_i$

# Constraints balance aestheticity

KHD<sup>†</sup>R, NeurIPS '25



baseline



unconstrained



constrained

★ preference reward (■)

★ saturation reward (■)

★ local contrast reward (■)

# Lagrangian approach

## ■ LAGRANGIAN

$$L(p; \lambda) = D_{\text{diff}}(p; q) + \sum_{i=1}^m \lambda_i (D_{\text{diff}}(p; q^i) - \epsilon_i) + \sum_{j=1}^n \lambda_{m+j} (b_j - \mathbb{E}_{x \sim p(\cdot)}[r_j(x)])$$

★ penalize violation via dual variable  $\lambda \geq 0$

## ■ DUAL FUNCTION

$$D(\lambda) := \underset{p}{\text{minimize}} L(p; \lambda)$$

★ **concave** function, regardless of the convexity in  $p$

# Dual problem

## ■ DUAL MAXIMIZATION

$$\underset{\lambda \geq 0}{\text{maximize}} \quad D(\lambda)$$

convex optimization

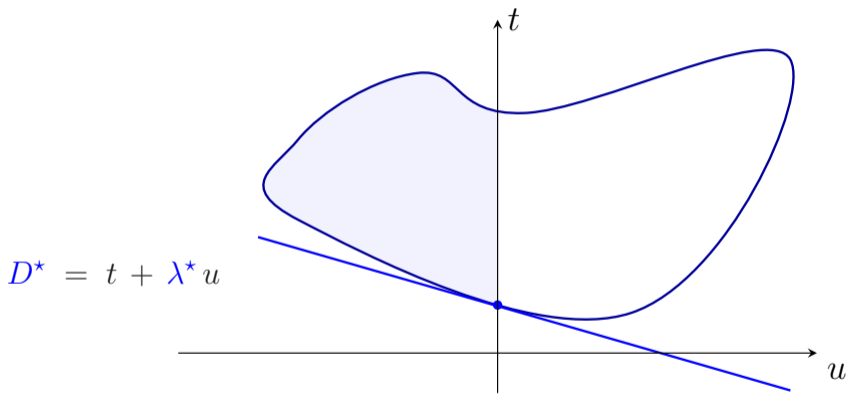
- ★ **subgradient ascent** finds an optimal dual variable  $\lambda^*$

**Recovery of an optimal model  $p^*$**

$$p^* \in \underset{p}{\text{argmin}} \quad L(p; \lambda^*)$$

# Strong duality

## ■ GEOMETRIC INTERPRETATION



$$\text{image } \mathcal{E} = \{ (D_{\text{diff}}(p; q^i) - \epsilon_i, D_{\text{diff}}(p; q)) \mid p \}$$

**Optimal hyperplane touches  $\mathcal{E}$  at an optimal model:  $D^* = D_{\text{diff}}(p^*; q)$**

# Constrained parameter optimization

minimize  $D_{\text{div}}(p_{\theta}; q)$   
 $\theta$

subject to  $D_{\text{div}}(p_{\theta}; q^i) \leq \epsilon_i$  for  $i = 1, \dots, m$   $\longrightarrow$  **proximity**

$\mathbb{E}_{x \sim p_{\theta}(\cdot)} [r_j(x)] \geq b_j$  for  $j = 1, \dots, n$   $\longrightarrow$  **reward**

★ diffusion model parameter  $\theta$

## CHALLENGE

**Nonconvex** constrained optimization  $\rightarrow$  **Lack of strong duality**

# Overview of our results

## ■ PRE-TRAINING: CONSTRAINED DIFFUSION MODELS

KD<sup>†</sup>R, NeurIPS '24

- ★ optimal model is a mixture distribution
- ★ dual training & optimality gap

## ■ POST-TRAINING: ALIGNMENT WITH CONSTRAINTS

KHD<sup>†</sup>R, NeurIPS '25

- ★ optimal model is a reward weighted distribution
- ★ dual training & optimality gap

objective & constraint

Dual training finds **an optimal diffusion model**,  
up to **an approximation error**

# **PART I PRE-TRAINING OF DIFFUSION MODELS**

constrained diffusion models

# Constrained diffusion models

**Original data:**  $x_0 \sim q(\cdot)$

**Desirable data:**  $x_0 \sim q^i(\cdot)$  for  $i = 1, \dots, m$

## ■ CONSTRAINED DISTRIBUTION OPTIMIZATION

minimize  $_p \quad D_{\text{KL}}(q(x_{0:T}) \parallel p(x_{0:T}))$

subject to  $D_{\text{KL}}(q^i(x_{0:T}) \parallel p(x_{0:T})) \leq \epsilon_i$  for  $i = 1, \dots, m$

★ **convex** forward KL objective and **convex** forward KL constraints

**Convex** constrained distribution optimization  $\rightarrow$  **Strong duality**

# One pillar

## ■ EVIDENCE LOWER BOUND ( ELBO )

$$E(p; q) = \mathbb{E}_{q(x_0)} \mathbb{E}_{q(x_{1:T} | x_0)} \left[ \log \frac{p(x_{0:T})}{q(x_{1:T} | x_0)} \right]$$

lower bound on log-likelihood

★ KL divergence decomposition

$$D_{\text{KL}}(q(x_{0:T}) \| p(x_{0:T})) = - E(p; q) + \mathbb{E}_{q(x_0)} [q(x_0)]$$

ELBO is **linear** in data distribution  $q(x_0)$

KL divergence equals **negative** ELBO (up to a constant)

# Lagrangian approach

## ■ LAGRANGIAN

$$L(p; \lambda) = D_{\text{KL}}(q(x_{0:T}) \| p(x_{0:T})) + \sum_{i=1}^m \lambda_i (D_{\text{KL}}(q^i(x_{0:T}) \| p(x_{0:T})) - \epsilon_i)$$

★ existence of **an optimal dual variable  $\lambda^*$**

## ■ OPTIMAL MODEL RECOVERY

$$\underset{p}{\text{minimize}} L(p; \lambda^*) \iff \underset{p}{\text{minimize}} -E(p; q) - \sum_{i=1}^m \lambda_i^* E(p; q^i)$$

# Optimal model

## ■ ELBO LINEARITY

$$-E(p; q) - \sum_{i=1}^m \lambda_i^* E(p; q^i) = -E\left(p; q + \sum_{i=1}^m \lambda_i^* q^i\right)$$

★ a mixture data distribution  $q_{\text{mix}}^{(\lambda^*)}(x_0) = \frac{1}{1+(\lambda^*)^\top \mathbf{1}} (q(x_0) + \sum_{i=1}^m \lambda_i^* q^i(x_0))$

## ■ ELBO-KL DIVERGENCE EQUIVALENCE

$$\underset{p}{\text{minimize}} L(p; \lambda^*) \iff \underset{p}{\text{minimize}} D_{\text{KL}}\left(q_{\text{mix}}^{(\lambda^*)}(x_{0:T}) \parallel p(x_{0:T})\right)$$

**Optimal model is a mixture distribution**  $p^*(x_{0:T}) = q_{\text{mix}}^{(\lambda^*)}(x_{0:T})$

# **PART I PRE-TRAINING OF DIFFUSION MODELS**

method & theory

# Constrained diffusion models

**Original data:**  $x_0 \sim q(\cdot)$

**Desirable data:**  $x_0 \sim q^i(\cdot)$  for  $i = 1, \dots, m$

## ■ CONSTRAINED PARAMETER OPTIMIZATION

$$\underset{\theta}{\text{minimize}} \quad -E(p_{\theta}; q)$$

$$\text{subject to} \quad -E(p_{\theta}; q^i) \leq b_i \quad \text{for } i = 1, \dots, m$$

$$b_i = \epsilon_i - \mathbb{E}_{q^i(x_0)} [q^i(x_0)]$$

★ ELBO serves as a score-matching loss

$$E(p_{\theta}; q) \cong -\mathbb{E}_{q(x_0), t, x_t} [\|\hat{s}_{\theta}(x_t, t) - \nabla \log q(x_t)\|^2]$$

**Nonconvex** constrained parameter optimization  $\rightarrow$  **Lack of strong duality**

# Parametrized dual problem

## ■ PARAMETRIZED DUAL FUNCTION

$$D_p(\lambda) := \underset{\theta}{\text{minimize}} L(p_\theta; \lambda)$$

★ **concave**, and **nondifferentiable** function

## ■ PARAMETRIZED DUAL MAXIMIZATION

$$\underset{\lambda \geq 0}{\text{maximize}} D_p(\lambda)$$

**Existence of an optimal parametrized dual variable  $\lambda_p^*$**

# Search for an optimal parametrized dual variable

**Lagrangian minimizer:**  $\theta^*(\lambda) \in \underset{\theta}{\operatorname{argmin}} L(p_{\theta}; \lambda)$

**Subgradient:**  $u(\lambda) = \nabla_{\lambda} L(p_{\theta}; \lambda) |_{\theta = \theta^*(\lambda)}$

## ■ PROJECTED SUBGRADIENT ASCENT

$$\lambda^+ \leftarrow [\lambda + \eta u(\lambda)]_+$$

★ explicit subgradient  $u(\lambda) = -E(p_{\theta^*(\lambda)}; q^i) - b_i$

Subgradient ascent converges to **an optimal parametrized dual variable**  $\lambda_p^*$

# Dual training algorithm

## ■ ITERATION #1: COMPUTE A LAGRANGIAN MINIMIZER

$$\theta^*(\lambda) \in \underset{\theta}{\operatorname{argmin}} L(p_{\theta}; \lambda)$$

## ■ ITERATION #2: PERFORM A SUBGRADIENT ASCENT STEP

$$\lambda^+ \leftarrow [\lambda - \eta (E(p_{\theta^*(\lambda)}; q^i) + b_i)]_+$$

**QUESTION:** Optimality of  $\lambda_p^*$ -recovered model  $p_{\theta^*(\lambda_p^*)}$ ?

# Optimality

KD<sup>†</sup>R, NeurIPS '24

**Optimality gap:**  $\text{TV} \left( q_{\text{mix}}^{(\lambda^*)}, p_{\theta^*}(\lambda_p^*) \right)$

Implication (informal)

★ **Optimality gap** is dominated by

$$\frac{\text{poly}(d)}{\sqrt{T}} + \sqrt{d} \epsilon_{\text{score}} + \sqrt{\nu}$$

**Sublinear** in time  $T$  & **polynomial** in dimension  $d$

**Linear** in score-matching error  $\epsilon_{\text{score}}$  & **root-scaling** of param. gap  $\nu$

$$\nu := \max_{\hat{s}} \min_{\theta} \text{dist}(\hat{s}, \hat{s}_{\theta})$$

$$\left| E(p_{\theta}; q_{\text{mix}}^{(\lambda)}) \right| \leq \epsilon_{\text{score}}^2$$

# Optimality analysis

## Step #1: ELBO linearity

$$\begin{aligned}\theta^*(\lambda_p^*) &\in \operatorname{argmin}_{\theta} L(p_{\theta}; \lambda_p^*) \\ &= \operatorname{argmin}_{\theta} -E(p_{\theta}; q) - \sum_{i=1}^m \lambda_{i,p}^* E(p_{\theta}; q^i) \\ &= \operatorname{argmin}_{\theta} -E\left(p_{\theta}; q_{\text{mix}}^{(\lambda_p^*)}\right)\end{aligned}$$

Training a diffusion model with a data mixture  $q_{\text{mix}}^{(\lambda_p^*)}$

★ non-asymptotic convergence

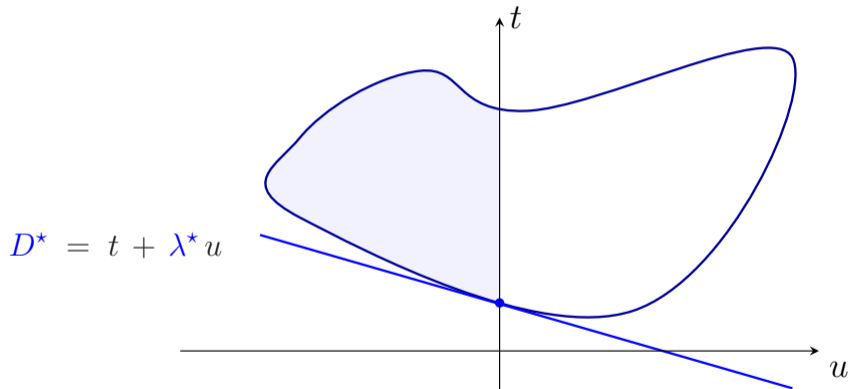
LWCC, ICLR '24

$$\text{TV}\left(q_{\text{mix}}^{(\lambda_p^*)}, p_{\theta^*(\lambda_p^*)}\right) \lesssim \frac{\text{poly}(d)}{\sqrt{T}} + \sqrt{d} \epsilon_{\text{score}}$$

**QUESTION:** Gap between  $q_{\text{mix}}^{(\lambda_p^*)}$  and  $q_{\text{mix}}^{(\lambda^*)}$ ?

# Duality gap

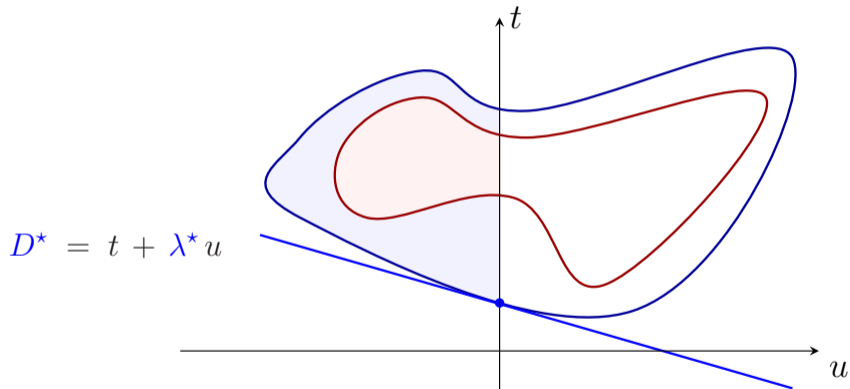
## ■ GEOMETRIC INTERPRETATION



$$\text{image } \mathcal{E} = \{ (-E^i(p) - b_i, -E(p)) \mid p \}$$

# Duality gap

## ■ GEOMETRIC INTERPRETATION



$$\text{image } \mathcal{E} = \{ (-E^i(p) - b_i, -E(p)) \mid p \} \quad \text{image } \mathcal{E}_p = \{ (-E^i(p_\theta) - b_i, -E(p_\theta)) \mid \theta \}$$

# Duality gap

## ■ GEOMETRIC INTERPRETATION

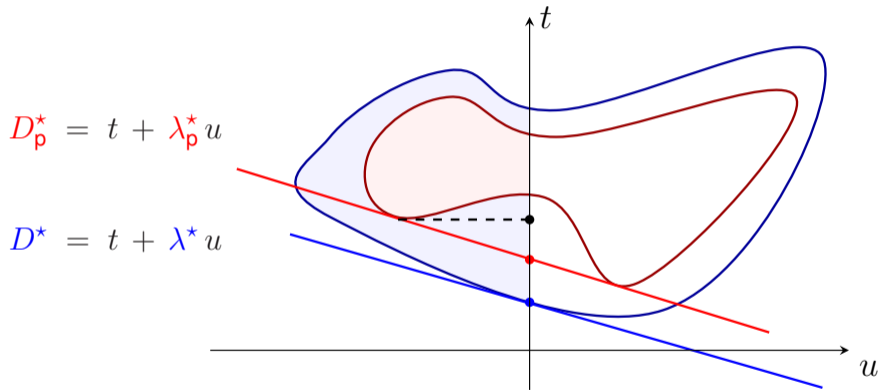
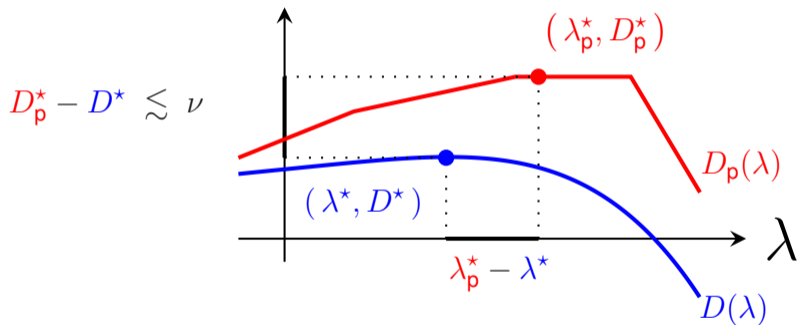


image  $\mathcal{E} = \{ (-E^i(p) - b_i, -E(p)) \mid p \}$     image  $\mathcal{E}_p = \{ (-E^i(p_\theta) - b_i, -E(p_\theta)) \mid \theta \}$

**Optimal hyperplane touches  $\mathcal{E}_p$  w/  $t$ -intercept  $D_p^*$ :  $D_p^* - D^* \lesssim \nu$**

# Optimality analysis

Step #2: Strong concavity of dual function



Optimal dual variables:  $\lambda^*$ ,  $\lambda_p^*$  are close:  $\|\lambda_p^* - \lambda^*\| \lesssim \sqrt{\nu}$

$$\implies \text{TV} \left( q_{\text{mix}}^{(\lambda^*)}, q_{\text{mix}}^{(\lambda_p^*)} \right) \lesssim \sqrt{\nu}$$

## **PART II POST-TRAINING OF DIFFUSION MODELS**

alignment with constraints

# Alignment with constraints

**Pretrained model:**  $q(x_{0:T})$

**Reward models:**  $r_i(x_0)$  for  $i = 1, \dots, m$

## ■ CONSTRAINED DISTRIBUTION OPTIMIZATION

$$\underset{p}{\text{minimize}} \quad D_{\text{KL}}(p(x_{0:T}) \parallel q(x_{0:T}))$$

$$\text{subject to} \quad \mathbb{E}_{x_0 \sim p(\cdot)} [r_i(x_0)] \geq b_i \text{ for } i = 1, \dots, m$$

★ **convex** reverse KL objective and **linear** constraints

**Convex** constrained distribution optimization  $\rightarrow$  **Strong duality**

# Lagrangian approach

## ■ LAGRANGIAN

$$L(p; \lambda) = D_{\text{KL}}(p(x_{0:T}) \parallel q(x_{0:T})) + \sum_{i=1}^m \lambda_i (b_i - \mathbb{E}_{x_0 \sim p(\cdot)} [r_i(x_0)])$$

★ existence of **an optimal dual variable  $\lambda^*$**

## ■ OPTIMAL MODEL RECOVERY

$$\underset{p}{\text{minimize}} L(p; \lambda^*) \iff \underset{p}{\text{minimize}} D_{\text{KL}}(p(x_{0:T}) \parallel q_{\text{rw}}^{(\lambda^*)}(x_{0:T}))$$

★ a reward weighted distribution  $q_{\text{rw}}^{(\lambda^*)}(\cdot) \propto q(\cdot) e^{(\lambda^*)^\top r(\cdot)}$

**Optimal model is a reward weighted distribution**  $p^*(x_{0:T}) = q_{\text{rw}}^{(\lambda^*)}(x_{0:T})$

# **PART II POST-TRAINING OF DIFFUSION MODELS**

method & theory

# Alignment with constraints

**Pretrained model:**  $q(x_{0:T})$

**Reward models:**  $r_i(x_0)$  for  $i = 1, \dots, m$

## ■ CONSTRAINED PARAMETER OPTIMIZATION

minimize  $D_{\text{KL}}(p_{\theta}(x_{0:T}) \parallel q(x_{0:T}))$   
 $\theta$

subject to  $\mathbb{E}_{x_0 \sim p_{\theta}(\cdot)} [r_i(x_0)] \geq b_i$  for  $i = 1, \dots, m$

★ reverse KL serves as an **on-policy** score-matching loss

$$D_{\text{KL}}(p_{\theta}(x_{0:T}) \parallel q(x_{0:T})) \cong \mathbb{E}_{x_t \sim p_{\theta}(\cdot), t} [\|s_{\theta}(x_t, t) - \nabla \log q(x_t)\|^2]$$

**Nonconvex** constrained parameter optimization  $\rightarrow$  **Lack of strong duality**

# Parametrized dual problem

## ■ PARAMETRIZED DUAL FUNCTION

$$D_p(\lambda) := \underset{\theta}{\text{minimize}} L(p_\theta; \lambda)$$

★ **concave**, and **nondifferentiable** function

## ■ PARAMETRIZED DUAL MAXIMIZATION

$$\underset{\lambda \geq 0}{\text{maximize}} D_p(\lambda)$$

**Existence of an optimal parametrized dual variable  $\lambda_p^*$**

# Dual training algorithm

## ■ ITERATION #1: COMPUTE A LAGRANGIAN MINIMIZER

$$\theta^*(\lambda) \in \operatorname{argmin}_{\theta} L(p_{\theta}; \lambda)$$

## ■ ITERATION #2: PERFORM A SUBGRADIENT ASCENT STEP

$$\lambda^+ \leftarrow \left[ \lambda + \eta \left( b - \mathbb{E}_{p_{\theta^*(\lambda)}(x_0)} [r(x_0)] \right) \right]_+$$

**QUESTION:** Optimality of  $\lambda_p^*$ -recovered model  $p_{\theta^*(\lambda_p^*)}$ ?

# Optimality

KHD<sup>†</sup>R, NeurIPS '25

ZLHBD<sup>†</sup>R, NeurIPS '25

**Objective optimality:**  $\left| D_{\text{KL}}\left(p_{\theta^*(\lambda_p^*)}\right) - D_{\text{KL}}(p^*) \right|$

**Constraint feasibility:**  $\left| r_i\left(p_{\theta^*(\lambda_p^*)}\right) - r_i(p^*) \right|$

Implication (informal)

★ **Objective optimality & Constraint feasibility** are dominated by

$$\sqrt{\nu}$$

**Root-scaling** of param. gap  $\nu$

$$\nu := \max_p \min_{\theta} \text{dist}(p, p_{\theta})$$

# Duality gap

## ■ GEOMETRIC INTERPRETATION

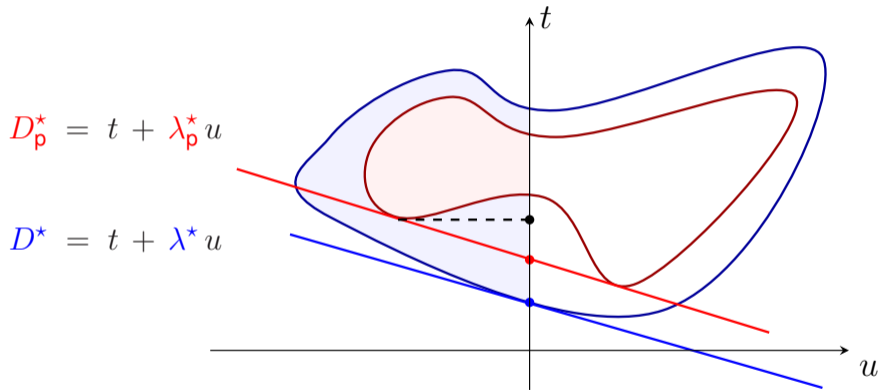


image  $\mathcal{E} = \{ (b_i - r_i(p), D_{\text{KL}}(p)) \mid p \}$

image  $\mathcal{E}_p = \{ (b_i - r_i(p\theta), D_{\text{KL}}(p\theta)) \mid \theta \}$

**Optimal hyperplane touches  $\mathcal{E}_p$  w/  $t$ -intercept  $D_p^*$ :  $D_p^* - D^* \lesssim \nu$**

# Dual function

## ■ LOWER ENVELOPE FUNCTION

$$\begin{aligned} D(\lambda) &:= \underset{p}{\text{minimize}} \quad L(p; \lambda) \\ &= \lambda^\top b - \log \mathbb{E}_{q(x_{0:T})} \left[ e^{\lambda^\top r(x_0)} \right] \end{aligned}$$

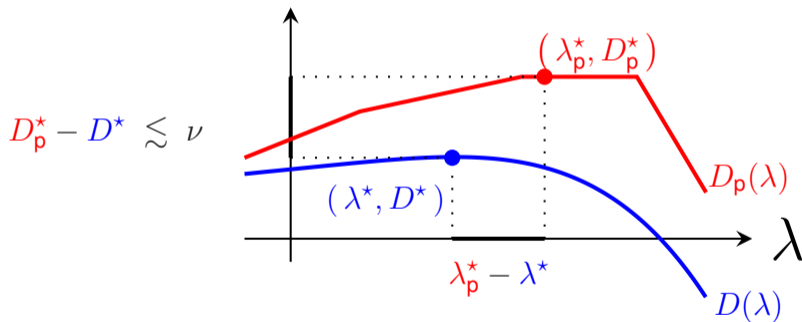
cumulant-generating function

★ **strictly concave**, and **locally strongly concave** function

$$\nabla^2 D(\lambda) \approx - \text{Var}_{x_0 \sim p^*(\cdot; \lambda)} [r(x_0)]$$

# Gap between optimal dual variables

## ■ GAP BETWEEN (UN)PARAMETRIZED DUAL FUNCTIONS

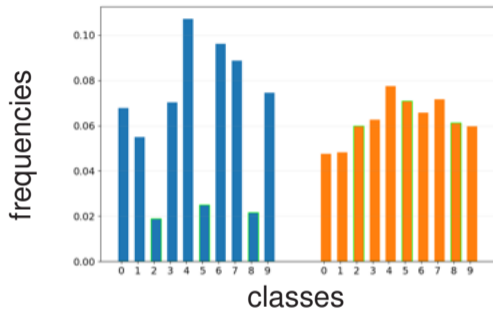


Optimal dual variables:  $\lambda^*$ ,  $\lambda_p^*$  are close:  $\|\lambda_p^* - \lambda^*\| \approx \sqrt{\nu}$

# **EMPIRICAL STUDY**

constrained diffusion models

# Fairness in image generation



unconstrained



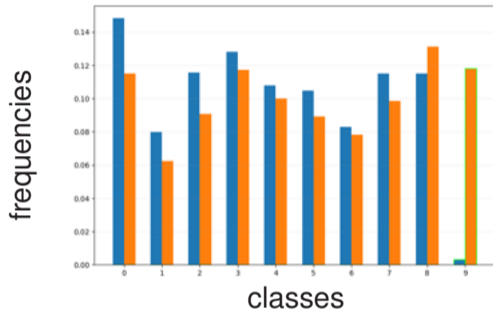
constrained

★ unconstrained model ( █ )

★ constrained model ( █ )

**Constraints improve representativeness**

## Adaptation of pretrained model



unconstrained



constrained

★ unconstrained model ( — )

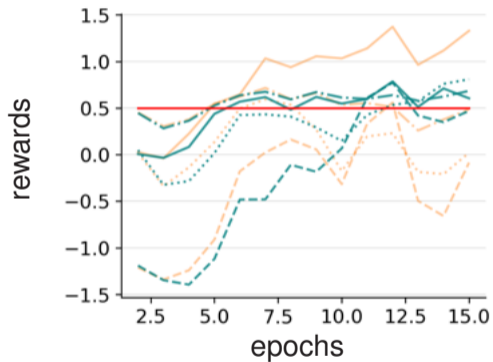
★ constrained model ( — )

**Constraints retain model utility**

## **EMPIRICAL STUDY**

alignment with constraints

# Reward alignment via constraints



pretrained

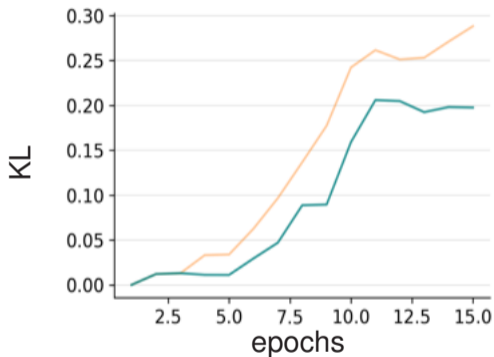
unconstrained

constrained

★ unconstrained model ( — )

★ constrained model ( — )

# Reward alignment via constraints



pretrained

unconstrained

constrained

★ unconstrained model ( — )

★ constrained model ( — )

**Constraints reduce model deviation**

# Summary & outlook

KD<sup>†</sup>R, NeurIPS '24

KHD<sup>†</sup>R, NeurIPS '25

ZLHBD<sup>†</sup>R, NeurIPS '25

## ■ TRAINING DIFFUSION MODELS WITH CONSTRAINTS

- ★ constrained diffusion models
- ★ alignment with constraints
- ★ dual method & optimality

## ■ OPEN CHALLENGES

- ★ constrained solutions
- ★ generalization

# Composition of diffusion models

KHD<sup>†</sup>R, NeurIPS '25

**Pretrained models:**  $q(x_0), q^i(x_0)$  for  $i = 1, \dots, m$

## ■ CONSTRAINED DISTRIBUTION OPTIMIZATION

$$\underset{p}{\text{minimize}} \quad D_{\text{KL}}(p(x_0) \parallel q(x_0))$$

$$\text{subject to} \quad D_{\text{KL}}(p(x_0) \parallel q^i(x_0)) \leq \epsilon_i \text{ for } i = 1, \dots, m$$

★ concept-aligned pretrained model  $q^i$

**Convex** constrained distribution optimization  $\rightarrow$  **Strong duality**

# Unlearning in diffusion models

KRD<sup>†</sup>, Ongoing

**Pretrained model:**  $q(x_{0:T})$

**Unlearning concepts:**  $q^i(x_{0:T})$  for  $i = 1, \dots, m$

## ■ CONSTRAINED DISTRIBUTION OPTIMIZATION

$$\underset{p}{\text{minimize}} \quad D_{\text{KL}}(q(x_{0:T}) \parallel p(x_{0:T}))$$

$$\text{subject to} \quad D_{\text{KL}}(q^i(x_{0:T}) \parallel p(x_{0:T})) \geq b_i \text{ for } i = 1, \dots, m$$

★ concept-aligned pretrained model  $q^i$

**Nonconvex** constrained distribution optimization  $\rightarrow$  **Strong duality**

**Thank you for your attention.**