
MCLR: Improving Conditional Modeling via Inter-Class Likelihood-Ratio Maximization and Unifying Classifier-Free Guidance with Alignment Objectives

Anonymous Authors¹

Abstract

Diffusion models achieve strong performance in generative modeling, but their success often relies heavily on classifier-free guidance (CFG), an inference-time heuristic that modifies the sampling trajectory. In theory, diffusion models trained with standard denoising score matching (DSM) should recover the target data distribution, raising two fundamental questions: (i) why is inference-time guidance necessary in practice, and (ii) can its underlying effect be internalized into a principled training objective? In this work, we argue that a key limitation of standard DSM is insufficient inter-class separation. To address this issue, we propose **MCLR**, an alignment objective that explicitly maximizes inter-class likelihood-ratios during training. Fine-tuning diffusion models with MCLR induces CFG-like improvements under standard sampling, substantially improving guidance-free conditional generation and narrowing the gap to inference-time CFG. Beyond these empirical benefits, we show theoretically that the CFG-guided score is exactly the optimal solution to a sample-adaptive weighted MCLR objective. This result connects CFG to alignment-based objectives, providing a mechanistic interpretation of CFG as an implicit inference-time contrastive alignment procedure.

1. Introduction

Diffusion models (Ho et al., 2020; Song et al., 2021b; Karras et al., 2022; Lipman et al., 2023) have become the dominant paradigm for high-fidelity generative modeling. These models generate samples by reversing a forward noising process using learned score functions, typically trained via denoising score matching (Vincent, 2011).

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

Although the reverse sampling process is theoretically guaranteed to recover the target distribution (Anderson, 1982), in practice it often yields samples of noticeably inferior quality: conditional generation frequently appears visually incoherent or insufficiently faithful to the intended class or prompt (Bradley & Nakkiran, 2024). In fact, nearly all high-quality generation by diffusion models relies heavily on classifier-free guidance (CFG) (Ho & Salimans, 2022), an inference-time modification of the reverse sampling process that injects an additional guidance term. While CFG substantially improves sample quality, reducing FID scores by up to 75% in seminal works (Peebles & Xie, 2023; Yu et al., 2025), its empirical necessity exposes a gap between the theoretical optimality of DSM and its practical behavior. This raises a central question:

Can we modify the DSM training objective in a principled way to achieve high-quality conditional generation under standard reverse sampling?

Recent observation (Li et al., 2025) suggests that standard conditional models suffer from insufficient inter-class separation: generated samples are less distinguishable across classes than real data, indicating that class-dependent structures are not fully captured by diffusion models. Motivated by this insight, we propose **MCLR**, a principled alignment objective that explicitly prompts inter-class separability by Maximizing the inter-Class log-Likelihood Ratio. By encouraging the model to amplify density differences between a target class and other classes, MCLR strengthens class-specific structures in the learned score function.

Empirically, models fine-tuned with MCLR exhibit CFG-like improvements under standard sampling, substantially improving guidance-free conditional generation and narrowing the gap to inference-time CFG.

Beyond the empirical benefits, we provide a theoretical result showing that the CFG-guided score coincides exactly with the optimizer of a weighted MCLR objective. This establishes a formal equivalence between classifier-free guidance and alignment-based training objectives, revealing

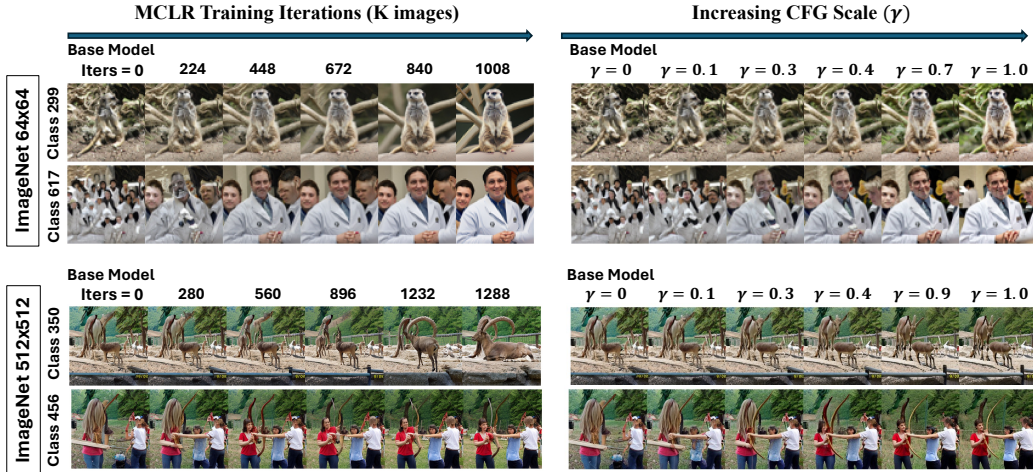


Figure 1. **Effects of MCLR vs. CFG.** We visualize the progressive emergence of class-specific structure under MCLR and classifier-free guidance (CFG). Progressive MCLR training induces effects analogous to increasing guidance strength in CFG, with both substantially improving class-specific patterns. For each image block, samples from different classes are generated from the same initial noise.

CFG as an inference-time contrastive alignment algorithm.

Beyond MCLR, we also analyze alternative inter-class contrastive objectives. In particular, we adapt DPO (Rafailov et al., 2023) to conditional generation, yielding **Conditional Contrastive DPO (CC-DPO)**, and show that its population optimum coincides with that of **Conditional Contrastive Alignment (CCA)** (Chen et al., 2025b). Our analysis reveals that CC-DPO and CCA perform multiplicative density-ratio reweighting, in contrast to the additive density-difference correction induced by MCLR. Comprehensive experiments reveal that MCLR consistently outperforms these alternatives across diverse models and datasets.

Summary of Contributions. Our main contributions are as follows:

- **A Principled Alignment Objective for Conditional Modeling.** We propose **MCLR**, a theoretically grounded fine-tuning objective that explicitly maximizes inter-class log-likelihood ratios to improve conditional generative modeling. Across diverse models and datasets, MCLR achieves substantial improvements in sample fidelity under standard reverse sampling and consistently outperforms training-time contrastive alternatives such as CC-DPO and CCA.
- **Theoretical Equivalence between CFG and Alignment Objectives.** We prove that the classifier-free guidance (CFG)-induced score coincides exactly with the optimizer of a weighted MCLR objective. This establishes a formal connection between CFG and alignment-based training, providing a mechanistic interpretation of CFG as an implicit inference-time alignment algorithm.
- **Understanding Contrastive Alternatives.** We provide both theoretical and empirical analyses of contrastive

alternatives such as DPO and CCA. When adapted to conditional generation, we show that DPO induces a gamma-powered density transformation equivalent to that of CCA—an equivalence that, to our knowledge, has not been previously established.

2. Preliminaries

2.1. Basics of Diffusion Models

Let $p_{\text{data}}(\mathbf{x})$ denote the ground-truth data distribution. Diffusion models construct a forward noising process that gradually perturbs p_{data} into a simple prior distribution using a stochastic differential equation (SDE):

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t) dt + g(t) d\mathbf{w}, \quad (1)$$

where $\mathbf{f}(\cdot, t)$ is the drift coefficient, $g(t)$ is the diffusion coefficient, and \mathbf{w} denotes the standard Brownian motion. Let $p_t(\mathbf{x})$ be the marginal distribution of $\mathbf{x}(t)$, and $p_{0t}(\mathbf{x}_t | \mathbf{x})$ the transition density from $\mathbf{x}(0)$ to $\mathbf{x}(t)$. For sufficiently large T , the distribution $p_T(\mathbf{x})$ becomes indistinguishable from a tractable prior $\pi(\mathbf{x})$, e.g., an isotropic Gaussian. The SDE (1) admits a reverse-time probability-flow ODE (Song et al., 2021b):

$$d\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - \frac{1}{2} g^2(t) \nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \right] dt. \quad (2)$$

Sampling from the reverse ODE requires access to the score function $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$, which can be approximated using a deep network $s_{\theta}(\mathbf{x}, t)$ trained via the denoising score matching (DSM) objective:

$$\mathcal{J}_{\text{DSM}}(\theta; w(\cdot)) := \frac{1}{2} \int_0^T \mathbb{E}_{p(\mathbf{x}), p_{0t}(\mathbf{x}_t | \mathbf{x})} \left[w(t) \left\| \nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}) - s_{\theta}(\mathbf{x}_t, t) \right\|_2^2 \right] dt, \quad (3)$$

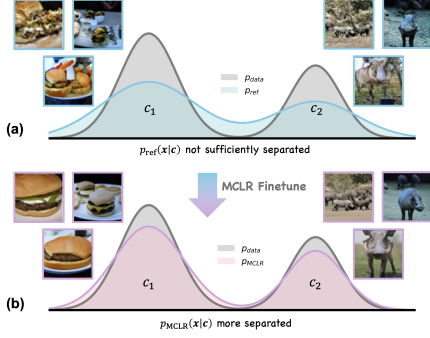


Figure 2. **Conceptual Illustration of MCLR.** (a) Samples generated from two classes using the same initial noises exhibit high visual similarity despite different conditioning labels, indicating insufficient separation of the learned conditional distributions. (b) MCLR mitigates this issue by encouraging class separation, resulting in generations with more distinct class-specific features.

where $w(\cdot)$ is a positive weighting function. For conditional diffusion models, the score network takes a conditional embedding c as input, and the DSM objective naturally extends to the conditional setting by taking the expectation over class labels and class-conditional data distributions:

$$\mathbb{E}_{c, p(x|c), p_{0t}(x_t|x)} [w(t) \|\nabla_{x_t} \log p_{0t}(x_t|x) - s_{\theta}(x_t, t, c)\|_2^2].$$

In this work, we focus on conditional diffusion models with an ODE sampler, for which the reverse ODE in (2) becomes:

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - \frac{1}{2}g^2(t)\nabla_{\mathbf{x}} \log p_t(\mathbf{x}|c)] dt. \quad (4)$$

2.2. Evidence Lower Bound for Diffusion Models

Let $p_{\theta}^{\text{ode}}(\mathbf{x})$ denote the distributions induced by the reverse ODE (2). Theorem 2 of (Song et al., 2021a) shows that, under certain regularity conditions, the log-likelihood satisfies:

$$\underbrace{\mathbb{E}_{p(\mathbf{x})}[\log p_{\theta}^{\text{ode}}(\mathbf{x})]}_{\text{Maximum Likelihood Estimation}} = -\mathcal{J}_{\text{DSM}}(\theta; g^2(\cdot)) + C, \quad (5)$$

where C is a constant independent of θ . Equation (5) is analogous to the evidence lower bound (ELBO) in variational autoencoders, revealing a fundamental connection between DSM and maximum likelihood estimation (MLE). This connection, enables MLE training in diffusion models (Mardani et al., 2023; Wallace et al., 2024; Zheng et al., 2025).

2.3. Classifier-Free Guidance

Although the reverse ODE in (4) is theoretically guaranteed to sample from the target conditional distribution, its practical generation quality is often unsatisfactory. In practice, high-quality conditional generation requires modifying the standard reverse process with an additional guidance term, known as classifier-free guidance (CFG) (Ho & Salimans, 2022), which leads to the perturbed reverse ODE:

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - \frac{1}{2}g^2(t)(\nabla_{\mathbf{x}} \log p_t(\mathbf{x}|c) \quad (6)$$

$$+ \gamma(\underbrace{\nabla_{\mathbf{x}} \log p_t(\mathbf{x}|c) - \nabla_{\mathbf{x}} \log p_t(\mathbf{x})}_{\text{CFG guidance}})] dt, \quad (7)$$

where γ controls the guidance strength. Intuitively, CFG sharpens the class- or conditional-specific structure by amplifying the difference between conditional and unconditional scores. However, the CFG-perturbed reverse sampling process (7) does not, in general, correspond to any known forward process (Bradley & Nakkiran, 2024). Despite recent progress towards understanding CFG (Wu et al., 2024; Chidambaram et al., 2024; Bradley & Nakkiran, 2024; Pavasovic et al., 2025; Li et al., 2025; Li & Jiao, 2025; Jin et al., 2025a), its underlying mechanisms remain only partially understood. This post-hoc modification of the sampling procedure motivates the search for principled alternatives that reproduce CFG-like improvements while preserving theoretical consistency, which may, in turn, shed light on the mechanisms underlying CFG itself.

2.4. Direct Preference Optimization

Direct preference optimization (DPO) is a widely used approach for aligning pretrained language or diffusion models with human preferences (Rafailov et al., 2023; Wallace et al., 2024). Under the Bradley-Terry (BT) model (Bradley & Terry, 1952), the probability that a sample \mathbf{x}_w is preferred over \mathbf{x}_l given a condition c is:

$$p(\mathbf{x}_w \succ \mathbf{x}_l | c) = \text{Sig}(\text{m}(r(\mathbf{x}_w | c) - r(\mathbf{x}_l | c))), \quad (8)$$

where $\text{Sig}(\cdot) := 1/(1 + \exp(-\cdot))$ denotes the Sigmoid function, \mathbf{x}_w and \mathbf{x}_l denote the preferred and non-preferred samples, respectively, and $r(\mathbf{x}|c)$ represents the underlying reward function reflecting human preference. DPO parameterizes this reward as $r_{\theta}(\mathbf{x}|c) := \beta \log p_{\theta}(\mathbf{x}|c) - \beta \log p_{\text{ref}}(\mathbf{x}|c)$ and estimates θ via maximum likelihood estimation on the BT model (8):

$$\max_{\theta} \mathbb{E}_{(c, \mathbf{x}_w, \mathbf{x}_l) \sim S} [\log p(\mathbf{x}_w \succ \mathbf{x}_l | c)] \quad (9)$$

$$:= \mathbb{E}_{(c, \mathbf{x}_w, \mathbf{x}_l) \sim S} [\log \text{Sig}(\beta \log \frac{p_{\theta}(\mathbf{x}_w | c)}{p_{\text{ref}}(\mathbf{x}_w | c)} - \beta \log \frac{p_{\theta}(\mathbf{x}_l | c)}{p_{\text{ref}}(\mathbf{x}_l | c)})], \quad (10)$$

where p_{ref} denotes the pretrained base (reference) model, β controls the strength of the KL regularization between p_{θ} and p_{ref} , and S represents the preference dataset. As we show later, DPO can be naturally adapted to enhance conditional modeling in visual generative models, providing a contrastive mechanism that parallels our MCLR objective.

3. Motivation and Method

In this section, we first demonstrate that diffusion models learn conditional distributions that lack sufficient class dis-

165 tinctiveness (see Section 3.1). To remedy this issue, we
 166 propose MCLR (Sections 3.2 and 3.3). For completeness,
 167 we also study a contrastive alternative that adapts DPO to
 168 conditional generation (Section 3.5).
 169

170 3.1. Diffusion Models Lack Class-Specificity

171 In theory, if the score functions are learned accurately,
 172 standard reverse diffusion sampling should produce sam-
 173 ples from the target conditional distribution. In practice,
 174 however, conditional sampling often fails to exhibit strong
 175 class-specific structure (Li et al., 2025). A common failure
 176 mode is that the conditional generations are weakly dis-
 177 tinguishable across classes: when starting from the same
 178 initial noise, samples generated under different class con-
 179 ditions frequently share similar global layouts, while class-
 180 discriminative features are attenuated or missing, as shown
 181 in Figure 2(a) and Figure 1. This suggests that the learned
 182 conditional distributions are insufficiently distinguishable
 183 from one another. Figure 2 illustrates this phenomenon con-
 184 ceptually, where the learned class-conditional distributions
 185 of a base model p_{ref} exhibit substantially less separation
 186 than the ground-truth data distribution p_{data} . This motivates
 187 modifying the standard DSM to more explicitly exploit the
 188 information carried by the conditioning variable.
 189

191 3.2. Maximum Inter-Class Likelihood-Ratio Training

192 Motivated by the preceding discussion, we propose improv-
 193 ing conditional modeling by explicitly encouraging inter-
 194 class separation in conditional models, as shown concep-
 195 tually in Figure 2(b). Let $p_{\theta}(\mathbf{x}|\mathbf{c})$ denote the model’s con-
 196 ditional distribution for class condition \mathbf{c} , where θ is the
 197 model parameter. We consider the following objective:
 198

$$199 \max_{\theta} \mathbb{E}_{\mathbf{c}, p(\mathbf{x}|\mathbf{c})} \log p_{\theta}(\mathbf{x}|\mathbf{c})$$

$$200 + \frac{\eta}{2} \mathbb{E}_{\mathbf{c}, \tilde{\mathbf{c}}, \mathbf{x} \sim p(\cdot|\mathbf{c}), \mathbf{y} \sim p(\cdot|\tilde{\mathbf{c}})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\theta}(\mathbf{x}|\tilde{\mathbf{c}})} + \log \frac{p_{\theta}(\mathbf{y}|\tilde{\mathbf{c}})}{p_{\theta}(\mathbf{y}|\mathbf{c})} \right],$$

$$201 \quad (11)$$

202 where \mathbf{c} and $\tilde{\mathbf{c}}$ are two randomly sampled classes and η con-
 203 trols the regularization strength. Compared to standard MLE
 204 (equivalently, DSM under appropriate weighting) in (5),
 205 eq. (11) introduces an additional regularizer, which we refer
 206 to as MCLR, that explicitly encourages samples to have
 207 a higher likelihood under their true class than under mis-
 208 matched classes.
 209

210 Specifically, given a sample \mathbf{x} (or \mathbf{y}) drawn from class \mathbf{c}
 211 (or $\tilde{\mathbf{c}}$), MCLR explicitly encourages its corresponding log-
 212 likelihood under the true class $\log p_{\theta}(\mathbf{x}|\mathbf{c})$ (or $\log p_{\theta}(\mathbf{y}|\tilde{\mathbf{c}})$)
 213 to be higher relative to its log-likelihood $\log p_{\theta}(\mathbf{x}|\tilde{\mathbf{c}})$ (or
 214 $\log p_{\theta}(\mathbf{y}|\mathbf{c})$) under a mismatched class $\tilde{\mathbf{c}}$ (or \mathbf{c}). In doing
 215 so, MCLR drives the model to increase the inter-class like-
 216 lihood ratio, thereby pushing $p_{\theta}(\mathbf{x}|\mathbf{c})$ to concentrate more
 217
 218
 219

probability mass in regions where the true class is favored
 over competing classes. Such regions typically correspond
 to samples with more pronounced class-specific features.

Intuitively, MCLR encourages the model to fully exploit
 the label-conditioned information. Standard conditional
 training simply feeds the conditional label \mathbf{c} alongside the
 data to the model, optimizes the MLE (or DSM) objec-
 tive, and relies on deep networks to automatically discover
 the conditional structure. If labels are under-exploited, the
 conditional distributions can collapse and become weakly
 distinguishable, yielding $p_{\theta}(\mathbf{x}|\mathbf{c}) \approx p_{\theta}(\mathbf{x}|\tilde{\mathbf{c}})$ regardless of
 the true class of \mathbf{x} . Maximizing likelihood ratio $\log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\theta}(\mathbf{x}|\tilde{\mathbf{c}})}$
 penalizes this collapse and encourages the model to discrim-
 inate between different conditions (classes) by leveraging the
 information encoded in the labels.

Applying the linearity of expectation to (11) and simpli-
 fying leads to the following equivalent form that we use
 throughout, unless otherwise stated:

$$\max_{\theta} \mathbb{E}_{\mathbf{c}, p(\mathbf{x}|\mathbf{c})} [\log p_{\theta}(\mathbf{x}|\mathbf{c})]$$

$$+ \underbrace{\eta \mathbb{E}_{\mathbf{c}, \tilde{\mathbf{c}}, \mathbf{x} \sim p(\cdot|\mathbf{c}), \mathbf{y} \sim p(\cdot|\tilde{\mathbf{c}})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\theta}(\mathbf{y}|\mathbf{c})} \right]}_{\text{MCLR Regularization}}. \quad (12)$$

Fine-tuning with MCLR. When a pretrained (suboptimal)
 base model $p_{\text{ref}}(\mathbf{x})$ that lacks class specificity is available,
 we can fine-tune it using MCLR in combination with KL
 regularization:

$$\max_{\theta} -\mathbb{E}_{\mathbf{c}} [D_{\text{KL}}(p_{\text{ref}}(\mathbf{x}|\mathbf{c})||p_{\theta}(\mathbf{x}|\mathbf{c}))]$$

$$+ \eta \mathbb{E}_{\mathbf{c}, \tilde{\mathbf{c}}, \mathbf{x} \sim p(\cdot|\mathbf{c}), \mathbf{y} \sim p(\cdot|\tilde{\mathbf{c}})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\theta}(\mathbf{y}|\tilde{\mathbf{c}})} \right]. \quad (13)$$

In the following section, we show that under mild regularity
 conditions, both (12) and (13) admit closed-form optimal
 solutions, providing theoretical insights into MCLR.

219 3.3. Theoretical Analysis of MCLR

We begin by analyzing the optimization problem (12).
 Define $h(\mathbf{x}|\mathbf{c}) := p(\mathbf{x}|\mathbf{c}) + \eta(p(\mathbf{x}|\mathbf{c}) - p(\mathbf{x}))$, where
 $p(\mathbf{x}) = \mathbb{E}_{\mathbf{c}} [p(\mathbf{x}|\mathbf{c})]$ is the unconditional distribution.

Theorem 3.1. *Suppose $h(\mathbf{x}|\mathbf{c})$ is supported on a bounded
 set K , and restrict the model class to densities satisfying
 $p_{\theta}(\mathbf{x}|\mathbf{c}) \geq \delta > 0$ on K . These conditions are mainly
 technical: bounded support is natural for image data with
 bounded pixel values, and the positive density floor ensures
 that the log-likelihood objective is well-defined. Then the
 optimal solution to (12) is:*

$$p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \begin{cases} \max\left\{\frac{h(\mathbf{x}|\mathbf{c})}{Z(\mathbf{c})}, \delta\right\}, & \mathbf{x} \in K, \\ 0, & \mathbf{x} \notin K, \end{cases} \quad (14)$$

where $Z(c)$ is the normalizing constant.

The formal assumptions and proof are provided in Appendix A.2.

Intuitively, the optimal conditional distribution $p_{\theta^*}(\mathbf{x}|c)$ induced by MCLR is proportional to $h(\mathbf{x}|c)$ whenever $h(\mathbf{x}|c) > Z(c)\delta$, while being clipped to the floor δ elsewhere. In the limit $\delta \rightarrow 0$, the floor disappears and the optimal distribution approaches:

$$p_{\theta^*}(\mathbf{x}|c) = \frac{h^+(\mathbf{x}|c)}{\int_{\mathbf{x}} h^+(\mathbf{x}|c) d\mathbf{x}}, \quad (15)$$

where $h^+(\mathbf{x}|c) := \max\{h(\mathbf{x}|c), 0\}$, such that the negative part of $h(\mathbf{x}|c)$ is truncated to zero and then normalized to form a valid distribution. One can interpret (14) and (15) as the "sum-of-difference" distribution: MCLR reshapes $p(\mathbf{x}|c)$ by adding to it the difference between $p(\mathbf{x}|c)$ and $p(\mathbf{x})$, such that the regions where $p(\mathbf{x}|c) > p(\mathbf{x})$ (i.e., samples with strong class-specific features) are amplified, and the regions where $p(\mathbf{x}|c) < p(\mathbf{x})$ (i.e., ambiguous samples lying near class boundaries) are suppressed.

For the fine-tuning objective in (13), it is easy to show that the optimal solution admits the same form as (14), but in this case $h(\mathbf{x}|c) = p_{\text{ref}}(\mathbf{x}|c) + \eta(p(\mathbf{x}|c) - p(\mathbf{x}))$. This directly leads to the following corollary (Appendix A.3 provides the proof).

Corollary 1. If the base model p_{ref} satisfies the mixture error model:

$$p_{\text{ref}}(\mathbf{x}|c) = (1 - \eta)p(\mathbf{x}|c) + \eta p(\mathbf{x}), \quad (16)$$

where $\eta \in [0, 1]$, then fine-tuning $p_{\text{ref}}(\mathbf{x}|c)$ with MCLR regularization (13) recovers the ground truth conditional distribution $p(\mathbf{x}|c)$.

The mixture error model (16) posits that the base model suffers from cross-class "leakage"; specifically, the learned conditional density $p_{\text{ref}}(\mathbf{x}|c)$ is a convex combination of the ground-truth conditional distribution $p(\mathbf{x}|c)$ and the unconditional counterpart $p(\mathbf{x})$, which corresponds to a weighted average over all class-conditional distributions. MCLR counteracts the leakage by adding the contrastive density difference $p(\mathbf{x}|c) - p(\mathbf{x})$ to the base model, thereby suppressing regions dominated by competing classes and strengthening class-specific regions. Although the error structure of practical diffusion models is likely more complex than this simple mixture model, the model provides a useful intuition for why MCLR can mitigate insufficient class separation.

3.4. Approximating Log-Likelihood with ELBO

The objective above is written in terms of log-likelihoods. This is directly applicable to autoregressive models (Tian

et al., 2024), but exact likelihood computation for diffusion models is too expensive for training. We therefore approximate the log-likelihood using the ELBO (5), so that the MCLR objective (12) becomes (148) in Appendix E.1, which can be approximated with Monte Carlo sampling.

Denoising Interpretation of MCLR. According to the equivalence between score function and optimal MMSE denoiser (see Appendix E.1 for details), the score function can be parameterized as: $s_{\theta}(\mathbf{x}, t, c) = \frac{\mathcal{D}_{\theta}(\mathbf{x}; \sigma(t), c) - \mathbf{x}}{\sigma^2(t)}$, where $\sigma(t)$ is the standard deviation of additive noise at time t . With a customized noise-time distribution $p(t)$ and weighting function $w(t)$, the MCLR regularization becomes:

$$\mathbb{E}_{\substack{c, \tilde{c}, t \sim p(t) \\ \mathbf{x} \sim p(\cdot|c), \mathbf{x}_t \sim p_{0t}(\cdot|\mathbf{x}) \\ \mathbf{y} \sim p(\cdot|\tilde{c}), \mathbf{y}_t \sim p_{0t}(\cdot|\mathbf{y})}} \left[w(t) (\|\mathbf{x} - \mathcal{D}_{\theta}(\mathbf{x}_t; \sigma(t), c)\|_2^2 - \|\mathbf{y} - \mathcal{D}_{\theta}(\mathbf{y}_t; \sigma(t), c)\|_2^2) \right]. \quad (17)$$

Eq. (17) provides a denoising-perspective of MCLR: it encourages the target-condition denoiser $\mathcal{D}(\cdot; \sigma(t), c)$ to achieve lower reconstruction error on samples of the target condition c than on samples of mismatched-condition \tilde{c} . We use this formulation in practical implementation.

3.5. Adapting DPO for Improved Conditional Modeling

The core principle behind MCLR is to improve conditional modeling by encouraging the model to exploit class-dependent structures through inter-class contrast, which can also be instantiated via other contrastive objectives such as DPO. To study this connection, we adapt DPO to conditional generation by treating samples from the target class c as preferred (\mathbf{x}_w) and samples from other randomly selected classes as non-preferred (\mathbf{x}_l). This leads to the following objective:

$$\min_{\theta} - \mathbb{E}_{c, \tilde{c}, \mathbf{x}_w \sim p(\mathbf{x}|c), \mathbf{x}_l \sim p(\mathbf{x}|\tilde{c})} \left[\log \text{Sigm} \left(\beta \log \frac{p_{\theta}(\mathbf{x}_w | c)}{p_{\text{ref}}(\mathbf{x}_w | c)} - \beta \log \frac{p_{\theta}(\mathbf{x}_l | c)}{p_{\text{ref}}(\mathbf{x}_l | c)} \right) \right]. \quad (18)$$

We refer to (18) as Conditional Contrastive DPO (CC-DPO). Similar to MCLR, CC-DPO objective also admits a closed-form solution, as stated in the following theorem (see proof in Appendix B.2).

Theorem 3.2. *Under certain regularity conditions, the optimal solution to (18) is:*

$$p_{\theta^*}(\mathbf{x}|c) = \frac{1}{\tilde{Z}(c)} p_{\text{ref}}(\mathbf{x}|c) \left(\frac{p(\mathbf{x}|c)}{p(\mathbf{x})} \right)^{\frac{1}{\beta}}, \quad (19)$$

where $\tilde{Z}(c) = \int_{\mathbf{x}} p_{\text{ref}}(\mathbf{x}|c) \left(\frac{p(\mathbf{x}|c)}{p(\mathbf{x})} \right)^{\frac{1}{\beta}} d\mathbf{x}$ is the normalizing constant.

Comparison between MCLR and CC-DPO. The optimal solution (19) is known as the *gamma-powered* distribution,

which was initially conjectured to characterize the effect of classifier-free guidance (CFG) (Ho & Salimans, 2022), but was later shown not to correspond to the true CFG dynamics (Karras et al., 2024a; Bradley & Nakkiran, 2024).

Although both amplifying regions where $p(\mathbf{x}|\mathbf{c}) > p(\mathbf{x})$ while suppressing regions where $p(\mathbf{x}|\mathbf{c}) < p(\mathbf{x})$, they do so through fundamentally different transformations: MCLR modifies the distribution additively through the density difference $p(\mathbf{x}|\mathbf{c}) - p(\mathbf{x})$, whereas CC-DPO reweights it multiplicatively through the density ratio $\left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})}\right)^{\frac{1}{\beta}}$. The multiplicative form $\left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})}\right)^{\frac{1}{\beta}}$ of CC-DPO can be overly aggressive. Consider a point $\tilde{\mathbf{x}}$ such that $p(\tilde{\mathbf{x}}) \approx 0$ while $p(\tilde{\mathbf{x}}|\mathbf{c}) > 0$, a situation that naturally arises when \mathbf{c} are minority classes. In this case, the ratio $\left(\frac{p(\tilde{\mathbf{x}}|\mathbf{c})}{p(\tilde{\mathbf{x}})}\right)^{1/\beta}$ becomes ill-conditioned, strongly amplifying the density at $\tilde{\mathbf{x}}$, therefore potentially driving the learned conditional distribution toward degenerate or unstable solutions. By contrast, the optimal solution (14) induced by MCLR remains well-behaved. We include CC-DPO and CCA as contrastive baselines in our experiments.

Equivalence between CC-DPO and CCA. Interestingly, the same gamma-powered distribution (19) can also be learned via Conditional Contrastive Alignment (CCA) (Chen et al., 2025b), a recently proposed method for autoregressive models (see Appendix C). Our analysis therefore establishes a previously unrecognized equivalence between CC-DPO and CCA at the level of their induced optimal distributions. Empirically, as we demonstrate in Appendix G.3, CC-DPO matches or outperforms CCA while requiring fewer hyperparameters, making it simpler to deploy in practice.

4. CFG as an Alignment Algorithm: A Mechanistic Interpretation

We now show that classifier-free guidance (CFG) is not merely an inference-time heuristic, but the exact optimal solution of an alignment objective. Specifically, we prove that the CFG-guided score in (7) coincides with the unique minimizer of a sample-adaptive weighted MCLR objective.

Theorem 4.1. *For any time sampling distribution $p(t)$ and weighting function $w(t)$, the CFG-guided score*

$$\begin{aligned} s_{\text{cfg}}(\mathbf{x}_t, t, \mathbf{c}) &:= \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) \\ &+ \eta \left(\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \right) \end{aligned}$$

is the unique minimizer of a sample-adaptive weighted

ELBO-approximated MCLR objective:

$$\begin{aligned} \min_{s_{\theta}(\cdot)} \mathbb{E}_{\substack{\mathbf{c}, t \sim p(t), \\ \mathbf{x} \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})}} & \left[w(t) \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - s_{\theta}(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right] \\ + \eta \mathbb{E}_{\substack{\mathbf{c}, \tilde{\mathbf{c}}, t \sim p(t), \\ \mathbf{x} \sim p(\cdot|\mathbf{c}), \mathbf{y} \sim p(\cdot|\tilde{\mathbf{c}}), \\ p_{0t}(\mathbf{x}_t|\mathbf{x}), p_{0t}(\mathbf{y}_t|\mathbf{y})}} & \left[w(t) \left(\|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - s_{\theta}(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right. \right. \\ & \left. \left. - \frac{p_t(\mathbf{y}_t|\mathbf{c})}{p_t(\mathbf{y}_t)} \|\nabla_{\mathbf{y}_t} \log p_{0t}(\mathbf{y}_t|\mathbf{y}) - s_{\theta}(\mathbf{y}_t, t, \mathbf{c})\|_2^2 \right) \right], \end{aligned} \quad (20)$$

where $p_t(\mathbf{y}_t) := \mathbb{E}_{\mathbf{c}}[p_t(\mathbf{y}_t|\mathbf{c})]$ (the proof is provided in Appendix D.1).

The regularization term in (20) has the same contrastive structure as the standard MCLR objective in (148) and (17): it encourages the conditional denoiser associated with the target condition \mathbf{c} to achieve lower denoising error on samples from the target condition than on samples drawn from mismatched conditions $\tilde{\mathbf{c}}$. The key difference from standard MCLR is that the negative component is reweighted by the sample-adaptive likelihood ratio $\frac{p_t(\mathbf{y}_t|\mathbf{c})}{p_t(\mathbf{y}_t)}$. Since this weight is nonnegative, the objective preserves the inter-class contrastive nature of MCLR.

Intuitively, a larger weight indicates that a mismatched noisy sample \mathbf{y}_t is more likely to be confused with samples from the target condition \mathbf{c} . Hence, the adaptive weight makes the CFG-equivalent objective (20) more selective over negative samples: samples \mathbf{y}_t from mismatched conditions that have high relative density under the target condition \mathbf{c} receive larger weights, while samples that are already easily distinguishable from the target condition (those with small $\frac{p_t(\mathbf{y}_t|\mathbf{c})}{p_t(\mathbf{y}_t)}$) receive smaller weights. Thus, the adaptive weight can be interpreted as a form of hard-negative mining, emphasizing ambiguous or class-confusable samples that provide stronger contrastive signal.

This perspective clarifies the relationship between MCLR and CFG. Standard MCLR uses a uniform contrastive penalty over mismatched samples, whereas CFG corresponds to a sample-adaptive weighted version of MCLR. Therefore, CFG can be understood as an **inference-time** contrastive alignment procedure, while MCLR provides a **training-time** mechanism for partially internalizing this effect into the model. In practice, we implement the standard, uniformly weighted MCLR objective rather than the sample-adaptive version. Nevertheless, as shown empirically, standard MCLR induces similar qualitative effects and substantially narrows the gap between unguided sampling and CFG.

Finally, Theorem 4.1 also provides a unified framework for interpreting CFG variants. Figure 4 summarizes this framework, and Appendix D.2 provides additional discussion.

5. Related Work

Our work contributes to recent efforts (Chen et al., 2025b;a; Tang et al., 2025) that seek to improve guidance-free conditional generation by modifying the training objective rather than applying classifier-free guidance (CFG) at inference time. The closest related method is Conditional Contrastive Alignment (CCA) (Chen et al., 2025b), which learns the gamma-powered distribution via Noise Contrastive Estimation (Gutmann & Hyvärinen, 2010). We show that MCLR outperforms CCA in diffusion models while achieving comparable performance in autoregressive settings, and further reveal a previously unexplored equivalence between CC-DPO and CCA.

Another related direction is Guidance-Free Training (GFT) (Chen et al., 2025a; Tang et al., 2025), which aims to reproduce CFG-induced score functions through modified DSM objectives. While GFT mimics the functional form of CFG during training, MCLR approaches the problem from a likelihood-ratio perspective and reveals the contrastive structure implicit in CFG, offering a complementary mechanistic interpretation.

Several additional works (Yan et al., 2024; Lee et al., 2025; Kadkhodaie et al., 2024; Yun et al., 2025) employ class-wise contrastive objectives to improve conditional generation, but lack a formal characterization of the underlying theoretical properties.

Finally, Direct Discriminative Optimization (DDO) (Zheng et al., 2025) contrasts real and synthetic samples rather than class-conditional distributions, and is included as a baseline in our experiments. Additional discussion of related work is deferred to Appendix H.

6. Experimental Results

In this section, we empirically evaluate the effectiveness of the proposed method. Our experiments demonstrate that: (i) MCLR substantially improves conditional generation quality and outperforms existing training-time baselines; and (ii) MCLR achieves performance comparable to CFG, exhibiting a similar fidelity–diversity trade-off and producing comparable qualitative effects. Due to space constraints, we present only a subset of the results here and defer a more comprehensive evaluation to Appendices F and G.

6.1. Experimental Setups

Datasets, Models, and Baselines. In our experiments, we focus on fine-tuning pretrained models. For diffusion models, we fine-tune pretrained EDM2 models (Karras et al., 2024b) on ImageNet-64×64 and ImageNet-512×512, and SiT (with REPA) model (Yu et al., 2025) on ImageNet-256×256. For visual autoregressive models, we fine-tune VAR-d24 (Tian et al., 2024) on ImageNet-256×256. We compare against CFG and several training-time baselines,

including CC-DPO, CCA, and DDO.

Evaluation Metrics. We evaluate generative performance using Fréchet Distance (FD) (Heusel et al., 2017), Precision and Recall (Kynkäänniemi et al., 2019), and Inception score (IS) (Salimans et al., 2016).

For diffusion models, we primarily report FD computed with DINOv2 features (FD_{DINOv2}). While FID (Inception-based FD) is widely used in the diffusion literature, we find that for strong pretrained model such as EDM2 which already achieves strong FID (1.5-2), it can be insensitive to perceptually meaningful improvements: although both CFG and MCLR lead to visually pronounced quality improvements, the FID score often does not improve and may even degrade. In contrast, FD_{DINOv2} consistently captures these improvements, aligning with prior findings that it correlates more strongly with human evaluations (Stein et al., 2023). For base models that are less strong including VAR and SiT, both FID and FD_{DINOv2} improve consistently.

6.2. Overall Algorithmic Behavior

Progressive Class Separation and the Fidelity–Diversity Trade-off. We first analyze the training dynamics induced by MCLR to understand how it shapes conditional generation over time. Qualitatively, as shown in Figures 1 and 8 to 12, at early stages of training, images generated from the same initial noises often exhibit similar global structures across different class conditions, indicating weak class-conditional modeling by the base model. As training progresses, generated images gradually develop distinct class-specific details, reflecting increasing inter-class separation and improved conditional modeling.

Quantitatively, as shown in Figure 3(b,d), this progression is reflected in a monotonic increase in IS, indicating increasingly class-discriminative generations, and a increase in Precision, indicating improved image fidelity. However, excessive training leads to reduced within-class diversity, manifested as a decrease in Recall (Figure 3(e)). Consequently, as shown in Figure 3(a), FD_{DINOv2} initially improves as conditional modeling strengthens, but later degrades as diversity diminishes, revealing a fidelity–diversity trade-off.

This behavior closely mirrors the effect of increasing the guidance scale in CFG. Similar trade-offs are observed for CC-DPO and CCA. We therefore report results at the checkpoint achieving the best FD_{DINOv2} score in Table 1.

MCLR Outperforms Training-time Baselines including CC-DPO and CCA. Table 1 shows that MCLR consistently achieves the best FD_{DINOv2} score among training-time baselines. On EDM2-S, MCLR improves FD_{DINOv2} from 95.20 to 52.69, outperforming DDO, CCA, and CC-DPO. On EDM2-L, MCLR reaches 42.50 FD_{DINOv2} , substantially outperforming the strongest training-time baseline. On VAR-d24 and SiT (with REPA), MCLR achieves the best

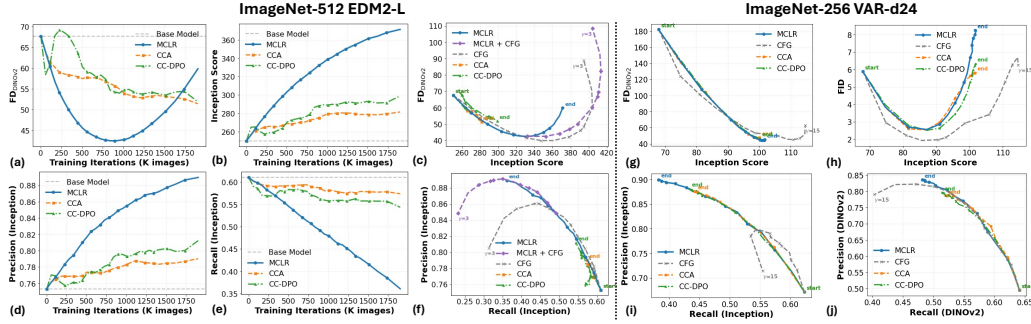


Figure 3. **Quantitative Results for EDM2-L and VAR-d24.** (a), (b), (d), and (e) show the evolution of FD, Inception Score, Precision, and Recall, respectively, as functions of training iterations. (c), (g), and (h) illustrate the FD–IS trade-offs, while (f), (i), and (j) depict the Precision–Recall trade-offs. For EDM2 models, we evaluate classifier-free guidance (CFG) scales $\gamma \in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.7, 0.9, 1, 1.5, 2.0, 3.0\}$. For Var-d24 model, $\gamma \in \{0.5, 0.8, 1.1, 1.5, 1.7, 2.0, 2.5, 3.0, 4.0, 5.0, 7.0, 10.0, 15.0\}$. **Start** denotes the performance of the base model, while **End** denotes the model obtained after a fixed finetuning duration.

Table 1. Quantitative results for MCLR, CFG and training-time baselines. Metrics are reported at the model checkpoint achieving the best case FD_{DINOv2} . Precision and Recall are computed using Inception features. The best and second-best results are highlighted in **bold** and underline, respectively.

Method	NFE	$FD_{DINOv2} \downarrow$	Prec. \uparrow	Rec. \uparrow	IS \uparrow
ImageNet (64×64)					
EDM2-S	63	95.20	0.705	0.614	60.43
+CFG	126	43.75	0.800	0.565	127.40
+DDO	63	72.97	0.689	0.642	65.84
+Diff-CCA	63	62.36	0.762	0.557	76.13
+CC-DPO	63	60.98	0.784	0.536	86.11
+MCLR (Ours)	63	52.69	0.800	0.505	<u>90.68</u>
ImageNet (256×256)					
VAR-d24	10	182.12	0.672	0.623	67.92
+CFG	20	45.08	0.798	<u>0.542</u>	100.70
+CCA	10	46.82	0.873	0.448	98.92
+CC-DPO	10	46.63	0.881	0.433	100.12
+MCLR (Ours)	10	44.31	0.893	0.404	100.84
SiT-XL/2+REPA	50	184.36	0.603	0.683	142.24
+CFG	100	50.23	<u>0.833</u>	0.506	385.87
+CCA	50	60.20	0.769	<u>0.551</u>	254.56
+CC-DPO	50	47.20	0.839	<u>0.451</u>	329.83
+MCLR (Ours)	50	45.96	0.825	0.488	<u>311.00</u>
ImageNet (512×512)					
EDM2-L	63	67.70	0.753	0.610	250.07
+CFG	126	39.86	<u>0.844</u>	0.512	360.30
+DDO	63	49.47	0.737	0.652	268.72
+Diff-CCA	63	51.45	0.790	0.574	281.45
+CC-DPO	63	51.92	0.812	0.544	298.89
+MCLR (Ours)	63	42.50	0.849	0.492	<u>332.02</u>

FD_{DINOv2} . Moreover, for EDM models, as shown in Figure 3 (c,f), MCLR traverses a significantly wider fidelity–diversity trade-off region steadily with a faster training speed. In contrast, CCA and CC-DPO converge early at suboptimal local minima and exhibit zigzag optimization trajectories, as evidenced by slowly improving or stalled learning curves.

MCLR Achieves Comparable Performance as CFG. CFG still achieves the best FD_{DINOv2} on EDM2 models, but MCLR substantially narrows the gap without inference-time guidance. For example, on EDM2-L, MCLR obtains 42.50 FD_{DINOv2} compared with CFG’s 39.86, while us-

ing half the NFE. When evaluated using Precision–Recall, MCLR exhibits competitive performance relative to CFG. Specifically, as shown in Figure 3(f), MCLR matches CFG in the high-recall regime corresponding to early training stages, and attains a substantially higher best-case precision in the high-precision regime at later training stages, where CFG begins to produce images with oversaturated colors. Furthermore, applying CFG on top of an MCLR-fine-tuned model further narrows the performance gap between the two methods, yielding higher best-case Inception Score and Precision on both EDM2-S and EDM2-L models (see Figure 3(c,f) and Figure 5(c,f,i)).

Qualitatively, MCLR and CFG produce highly similar visual effects, both substantially enhancing class-specific structures in the generated images, as shown in Figures 1 and 8 to 19. These observations are consistent with our theoretical analysis in Section 4, which interprets CFG as an implicit contrastive alignment method.

7. Conclusions and Limitations

In this work introduced MCLR, a training-time objective that substantially improves the unguided base model and narrows the gap to inference-time CFG. Beyond empirical improvements, our analysis reveals the implicit inter-class contrastive alignment mechanism of CFG.

Limitations. As a training-time method, MCLR produces a fixed model after training, hence lacks the flexibility of inference-time guidance to dynamically adjust this trade-off between generation quality and diversity through guidance strength. Moreover, CFG typically achieves stronger best-case FD and Inception Score. Bridging the gap between training-time objectives and inference-time guidance remains an important direction for future work. More broadly, inter-class separation is only one practical limitation of the standard DSM; systematically identifying and addressing additional limitations remains an important future direction.

Impact Statement

This work aims to improve the training of conditional diffusion models by reducing reliance on inference-time guidance mechanisms such as classifier-free guidance (CFG). By incorporating contrastive signals directly into the training objective, the proposed method enables standard sampling procedures to achieve CFG-like effects, potentially reducing computational overhead and simplifying deployment in downstream applications. These improvements may benefit a wide range of generative modeling tasks, including image synthesis, data augmentation, and representation learning.

At the same time, the proposed approach does not introduce new application domains or capabilities beyond those of existing diffusion models, and thus inherits similar ethical considerations and risks, such as the potential misuse of generated content. As with prior generative models, responsible use, appropriate dataset curation, and adherence to existing guidelines remain essential. Overall, this work contributes to a better theoretical and practical understanding of conditional diffusion models, with the goal of improving their efficiency and interpretability.

References

- Anderson, B. D. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 12(3):313–326, 1982.
- Boyd, S. P. and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.
- Bradley, A. and Nakkiran, P. Classifier-free guidance is a predictor-corrector. *arXiv preprint arXiv:2408.09000*, 2024.
- Bradley, R. A. and Terry, M. E. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Chen, H., Jiang, K., Zheng, K., Chen, J., Su, H., and Zhu, J. Visual generation without guidance. In *Forty-second International Conference on Machine Learning*, 2025a.
- Chen, H., Su, H., Sun, P., and Zhu, J. Toward guidance-free ar visual generation via condition contrastive alignment. In *The Thirteenth International Conference on Learning Representations*, 2025b.
- Cheng, M., Doudi, F., Kalathil, D., Ghavamzadeh, M., and Kumar, P. R. Diffusion blend: Inference-time multi-preference alignment for diffusion models. *arXiv preprint arXiv:2505.18547*, 2025.
- Chidambaram, M., Gatmiry, K., Chen, S., Lee, H., and Lu, J. What does guidance do? a fine-grained analysis in a simple setting. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- Frans, K., Park, S., Abbeel, P., and Levine, S. Diffusion guidance is a controllable policy improvement operator. *arXiv preprint arXiv:2505.23458*, 2025.
- Gutmann, M. and Hyvärinen, A. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 297–304. JMLR Workshop and Conference Proceedings, 2010.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- Ho, J. and Salimans, T. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Jiang, Z., Wen, Y., and Liu, Z. Rethinking preference alignment for diffusion models with classifier-free guidance. *arXiv preprint arXiv:2602.18799*, 2026.
- Jin, C., Shi, Q., and Gu, Y. Stage-wise dynamics of classifier-free guidance in diffusion models. *arXiv preprint arXiv:2509.22007*, 2025a.
- Jin, L., Qiu, Z., Liu, J., Diao, Z., Qiao, L., Ding, N., Lamb, A., and Qiu, X. Inference-time alignment control for diffusion models with reinforcement learning guidance. *arXiv preprint arXiv:2508.21016*, 2025b.
- Kadkhodaie, Z., Mallat, S., and Simoncelli, E. P. Feature-guided score diffusion for sampling conditional densities. *arXiv preprint arXiv:2410.11646*, 2024.
- Karras, T., Aittala, M., Aila, T., and Laine, S. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35: 26565–26577, 2022.
- Karras, T., Aittala, M., Kynkäänniemi, T., Lehtinen, J., Aila, T., and Laine, S. Guiding a diffusion model with a bad version of itself. *Advances in Neural Information Processing Systems*, 37:52996–53021, 2024a.
- Karras, T., Aittala, M., Lehtinen, J., Hellsten, J., Aila, T., and Laine, S. Analyzing and improving the training dynamics of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24174–24184, 2024b.

- 495 Kynkäänniemi, T., Karras, T., Laine, S., Lehtinen, J., and
496 Aila, T. Improved precision and recall metric for assess-
497 ing generative models. *Advances in neural information*
498 *processing systems*, 32, 2019.
- 499 Lee, J.-Y., Cha, B., Kim, J., and Ye, J. C. Aligning text to
500 image in diffusion models is easier than you think. *arXiv*
501 *preprint arXiv:2503.08250*, 2025.
- 503 Li, G. and Jiao, Y. Provable efficiency of guidance in diffu-
504 sion models for general data distribution. In *Forty-second*
505 *International Conference on Machine Learning*, 2025.
- 507 Li, X., Wang, R., and Qu, Q. Towards understanding the
508 mechanisms of classifier-free guidance. *arXiv preprint*
509 *arXiv:2505.19210*, 2025.
- 511 Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and
512 Le, M. Flow matching for generative modeling. In *The*
513 *Eleventh International Conference on Learning Repre-*
514 *sentations*, 2023.
- 515 Mardani, M., Song, J., Kautz, J., and Vahdat, A. A varia-
516 tional perspective on solving inverse problems with diffu-
517 sion models. In *The Twelfth International Conference on*
518 *Learning Representations*, 2023.
- 520 Miyasawa, K. et al. An empirical bayes estimator of the
521 mean of a normal population. *Bull. Inst. Internat. Statist.*,
522 38(181-188):1–2, 1961.
- 524 Oord, A. v. d., Li, Y., and Vinyals, O. Representation learn-
525 ing with contrastive predictive coding. *arXiv preprint*
526 *arXiv:1807.03748*, 2018.
- 527 Pavasovic, K. L., Verbeek, J., Biroli, G., and Mezard,
528 M. Classifier-free guidance: From high-dimensional
529 analysis to generalized guidance forms. *arXiv preprint*
530 *arXiv:2502.07849*, 2025.
- 532 Peebles, W. and Xie, S. Scalable diffusion models with
533 transformers. In *Proceedings of the IEEE/CVF interna-*
534 *tional conference on computer vision*, pp. 4195–4205,
535 2023.
- 537 Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Er-
538 mon, S., and Finn, C. Direct preference optimization:
539 Your language model is secretly a reward model. *Ad-*
540 *vances in Neural Information Processing Systems*, 36:
541 53728–53741, 2023.
- 543 Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V.,
544 Radford, A., and Chen, X. Improved techniques for
545 training gans. *Advances in neural information processing*
546 *systems*, 29, 2016.
- 547 Song, Y., Durkan, C., Murray, I., and Ermon, S. Maxi-
548 mum likelihood training of score-based diffusion models.
549 *Advances in neural information processing systems*, 34:
1415–1428, 2021a.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Er-
mon, S., and Poole, B. Score-based generative modeling
through stochastic differential equations. In *International*
Conference on Learning Representations, 2021b.
- Stein, G., Cresswell, J., Hosseinzadeh, R., Sui, Y., Ross, B.,
Villicroze, V., Liu, Z., Caterini, A. L., Taylor, E., and
Loaiza-Ganem, G. Exposing flaws of generative model
evaluation metrics and their unfair treatment of diffusion
models. *Advances in Neural Information Processing*
Systems, 36:3732–3784, 2023.
- Tang, Z., Bao, J., Chen, D., and Guo, B. Diffusion
models without classifier-free guidance. *arXiv preprint*
arXiv:2502.12154, 2025.
- Tian, K., Jiang, Y., Yuan, Z., Peng, B., and Wang, L. Visual
autoregressive modeling: Scalable image generation via
next-scale prediction. *Advances in neural information*
processing systems, 37:84839–84865, 2024.
- Vincent, P. A connection between score matching and de-
noising autoencoders. *Neural Computation*, 23(7):1661–
1674, 2011. doi: 10.1162/NECO.a.00142.
- Wallace, B., Dang, M., Rafailov, R., Zhou, L., Lou, A.,
Purushwalkam, S., Ermon, S., Xiong, C., Joty, S., and
Naik, N. Diffusion model alignment using direct pref-
erence optimization. In *Proceedings of the IEEE/CVF*
Conference on Computer Vision and Pattern Recognition,
pp. 8228–8238, 2024.
- Wu, Y., Chen, M., Li, Z., Wang, M., and Wei, Y. Theoretical
insights for diffusion guidance: A case study for gaussian
mixture models. In *Forty-first International Conference*
on Machine Learning, 2024.
- Yan, D., Qi, L., Hu, V. T., Yang, M.-H., and Tang, M.
Training class-imbalanced diffusion model via overlap
optimization. *arXiv preprint arXiv:2402.10821*, 2024.
- Yu, S., Kwak, S., Jang, H., Jeong, J., Huang, J., Shin, J.,
and Xie, S. Representation alignment for generation:
Training diffusion transformers is easier than you think.
In *The Thirteenth International Conference on Learning*
Representations, 2025.
- Yun, J., Alçalar, Y. U., and Akçakaya, M. No align-
ment needed for generation: Learning linearly separa-
ble representations in diffusion models. *arXiv preprint*
arXiv:2509.21565, 2025.
- Zheng, K., Chen, Y., Chen, H., He, G., Liu, M.-Y., Zhu, J.,
and Zhang, Q. Direct discriminative optimization: Your
likelihood-based visual generative model is secretly a gan
discriminator. *arXiv preprint arXiv:2503.01103*, 2025.

Appendices

A. Theoretical Analysis of MCLR

A.1. Equivalent Forms of MCLR Objective

By the linearity of expectation, it is easy to show that the MCLR objective (11) admits the following two equivalent forms:

$$\max_{\theta} \mathbb{E}_{\mathbf{c}, p(\mathbf{x}|\mathbf{c})} [\log p_{\theta}(\mathbf{x}|\mathbf{c})] + \underbrace{\eta \mathbb{E}_{\mathbf{c}, \tilde{\mathbf{c}}, \mathbf{x} \sim p(\cdot|\mathbf{c}), \mathbf{y} \sim p(\cdot|\tilde{\mathbf{c}})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\theta}(\mathbf{y}|\tilde{\mathbf{c}})} \right]}_{\text{MCLR Regularization (Form I)}}, \quad (21)$$

and

$$\max_{\theta} \mathbb{E}_{\mathbf{c}, p(\mathbf{x}|\mathbf{c})} [\log p_{\theta}(\mathbf{x}|\mathbf{c})] + \underbrace{\eta \mathbb{E}_{\mathbf{c}, \tilde{\mathbf{c}}, \mathbf{x} \sim p(\mathbf{x}|\mathbf{c})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\theta}(\mathbf{x}|\tilde{\mathbf{c}})} \right]}_{\text{MCLR Regularization (Form II)}}. \quad (22)$$

While the main text presents MCLR in Form I, the appendix derives most results from Form II. Because these two formulations are equivalent, we use them interchangeably throughout the analysis.

A.2. Main Theorem

In this section, we provide the proof for Theorem 3.1. We first restate the assumptions and theorem.

Assumption 1. The function $h(\mathbf{x}|\mathbf{c})$ has compact support; that is,

$$\text{supp } h(\mathbf{x}|\mathbf{c}) \subseteq K, \quad |K| < \infty, \quad (23)$$

where $|K|$ denotes the volume of set K .

Assumption 2. For $\forall \mathbf{x} \in K$,

$$p_{\theta}(\mathbf{x}|\mathbf{c}) \geq \delta > 0, \quad \delta < \frac{1}{|K|}. \quad (24)$$

Assumption 1 is mild in practice, as image data typically occupies a bounded pixel range. Assumption 2 is a regularity assumption that ensures the log-likelihood is well-defined, avoiding the singularity that occurs when evaluating $\log p_{\theta}(\mathbf{x}|\mathbf{c})$ at zero density. Since $p_{\theta}(\mathbf{x}|\mathbf{c})$ integrates to one over K , we necessarily have $\delta \leq \frac{1}{|K|}$. In the following, we further assume $\delta < \frac{1}{|K|}$.

Theorem 3.1. *Under the two assumptions, the optimal solution to (12) is:*

$$p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \begin{cases} \max\left\{\frac{h(\mathbf{x}|\mathbf{c})}{Z(\mathbf{c})}, \delta\right\}, & \mathbf{x} \in K, \\ 0, & \mathbf{x} \notin K, \end{cases} \quad (25)$$

where $Z(\mathbf{c})$ is the normalizing constant.

Proof. Without loss of generality, we consider discrete classes in this proof; the theorem extends straightforwardly to continuous classes. Suppose there are M classes $\{\mathbf{c}_i\}_{i=1}^M$, each class has prior probability $p(\mathbf{c}_i)$ and $\sum_{i=1}^M p(\mathbf{c}_i) = 1$, then the training objective (12) takes the following form:

$$\max_{\theta} \mathcal{L}(\theta) := \mathbb{E}_{\mathbf{c}, p(\mathbf{x}|\mathbf{c})} [\log p_{\theta}(\mathbf{x}|\mathbf{c})] + \eta \mathbb{E}_{\mathbf{c}, \tilde{\mathbf{c}}, \mathbf{x} \sim p(\mathbf{x}|\mathbf{c})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\theta}(\mathbf{x}|\tilde{\mathbf{c}})} \right] \quad (26)$$

$$= \max_{\theta} \sum_{i=1}^M p(\mathbf{c}_i) \mathbb{E}_{p(\mathbf{x}|\mathbf{c}_i)} [\log p_{\theta}(\mathbf{x}|\mathbf{c}_i)] + \eta \sum_{i=1}^M p(\mathbf{c}_i) \sum_{j=1}^M p(\mathbf{c}_j) \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{c}_i)} [\log p_{\theta}(\mathbf{x}|\mathbf{c}_i)] \quad (27)$$

$$- \eta \sum_{i=1}^M p(\mathbf{c}_i) \sum_{j=1}^M p(\mathbf{c}_j) \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{c}_i)} [\log p_{\theta}(\mathbf{x}|\mathbf{c}_j)]. \quad (28)$$

Note that we can decompose overall objective $\mathcal{L}(\boldsymbol{\theta})$ as:

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{k=1}^M \mathcal{L}_k(\boldsymbol{\theta}), \quad (29)$$

where $\mathcal{L}_k(\boldsymbol{\theta})$ is the amount contributed by $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)$:

$$\mathcal{L}_k(\boldsymbol{\theta}) = p(\mathbf{c}_k) \mathbb{E}_{p(\mathbf{x}|\mathbf{c}_k)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)] + \eta p(\mathbf{c}_k) \mathbb{E}_{p(\mathbf{x}|\mathbf{c}_k)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)] \quad (30)$$

$$- \eta p(\mathbf{c}_k) \sum_{i=1}^M p(\mathbf{c}_i) \mathbb{E}_{p(\mathbf{x}|\mathbf{c}_i)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)] \quad (31)$$

$$= p(\mathbf{c}_k) \mathbb{E}_{p(\mathbf{x}|\mathbf{c}_k)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)] + \eta p(\mathbf{c}_k) \mathbb{E}_{p(\mathbf{x}|\mathbf{c}_k)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)] - \eta p(\mathbf{c}_k) \mathbb{E}_{p(\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)]. \quad (32)$$

Note that we can optimize $\mathcal{L}_k(\boldsymbol{\theta})$ for each $k \in \{1, \dots, M\}$ to get the optimal conditional distribution $p(\mathbf{x}|\mathbf{c}_k)$ independently:

$$\max_{\boldsymbol{\theta}} \mathcal{L}_k(\boldsymbol{\theta}) \Leftrightarrow \max_{p_{\boldsymbol{\theta}}(\cdot|\mathbf{c}_k)} \mathbb{E}_{p(\mathbf{x}|\mathbf{c}_k)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)] + \eta \mathbb{E}_{p(\mathbf{x}|\mathbf{c}_k)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)] - \eta \mathbb{E}_{p(\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k)] \quad (33)$$

$$= \max_{p_{\boldsymbol{\theta}}(\cdot|\mathbf{c}_k)} \int_K \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}_k) (p(\mathbf{x}|\mathbf{c}_k) + \eta(p(\mathbf{x}|\mathbf{c}_k) - p(\mathbf{x}))) \, d\mathbf{x}. \quad (34)$$

In what follows, we drop the subscript k , so that the optimization problem becomes:

$$\max_{p_{\boldsymbol{\theta}}(\cdot|\mathbf{c})} \int_K \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}) (p(\mathbf{x}|\mathbf{c}) + \eta(p(\mathbf{x}|\mathbf{c}) - p(\mathbf{x}))) \, d\mathbf{x} \quad (35)$$

$$= \max_{p_{\boldsymbol{\theta}}(\cdot|\mathbf{c})} \int_K \log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}) h(\mathbf{x}|\mathbf{c}) \, d\mathbf{x} \quad (36)$$

Note that (36) shares the similar form as KL divergence, but $h(\mathbf{x}|\mathbf{c})$ is not a valid probability distribution, since there could exist \mathbf{x} such that $h(\mathbf{x}|\mathbf{c}) < 0$. In this case, setting $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}) = 0$ makes (36) approaches $+\infty$, hence the optimization problem does not have an attainable optimum. To make the optimization problem well-posed, we impose Assumption 2.

Define the sets

$$K_+(\mathbf{c}) = \{\mathbf{x} \in K : h(\mathbf{x}|\mathbf{c}) > 0\}, \quad K_-(\mathbf{c}) = \{\mathbf{x} \in K : h(\mathbf{x}|\mathbf{c}) \leq 0\}. \quad (37)$$

Under Assumption 2, it is straightforward to show the optimal distribution $p_{\boldsymbol{\theta}^*}(\mathbf{x}|\mathbf{c})$ must attain the lower bound δ for $\forall \mathbf{x} \in K_-(\mathbf{c})$; otherwise increasing the density will decrease the objective (36). Furthermore for $\mathbf{x} \in K^c$, where $h(\mathbf{x}|\mathbf{c}) = 0$, the optimal density must satisfy $p_{\boldsymbol{\theta}^*}(\mathbf{x}|\mathbf{c}) = 0$ for almost every $\mathbf{x} \in K^c$; otherwise, probability mass could be shifted from K^c to $K_+(\mathbf{c})$ to further increase the objective (36). Therefore, the remaining optimization concerns the density on $K^+(\mathbf{c})$, which is the optimal solution to the following constrained optimization problem:

$$\min_{p_{\boldsymbol{\theta}}(\cdot|\mathbf{c})} \int_{K_+(\mathbf{c})} -\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}) h(\mathbf{x}|\mathbf{c}) \, d\mathbf{x} \quad (38)$$

$$\text{s.t. } -p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}) + \delta \leq 0, \quad \forall \mathbf{x} \in K_+(\mathbf{c}) \quad (39)$$

$$\int_{K_+(\mathbf{c})} p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c}) \, d\mathbf{x} - m(\mathbf{c}) = 0, \quad (40)$$

where

$$m(\mathbf{c}) = 1 - \int_{K_-(\mathbf{c})} \delta \, d\mathbf{x} = 1 - \delta |K_-(\mathbf{c})|. \quad (41)$$

Since $\delta < \frac{1}{|K|}$ by assumption, we have $0 < m(\mathbf{c}) \leq 1$.

Note that this optimization problem is convex in $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c})$. Treating $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{c})$ for each \mathbf{x} as optimization variables, the constraints are affien and Slater's condition holds. Therefore strong duality applies, and the optimal solution can be characterized by the KKT conditions (Boyd & Vandenberghe, 2004).

Define the Lagrangian as:

$$\mathcal{L}(p_{\theta}(\cdot|\mathbf{c}), \lambda, u(\cdot)) = \int_{K_+(\mathbf{c})} -\log p_{\theta}(\mathbf{x}|\mathbf{c})h(\mathbf{x}|\mathbf{c}) d\mathbf{x} + \lambda \left(\int_{K_+(\mathbf{c})} p_{\theta}(\mathbf{x}|\mathbf{c}) d\mathbf{x} - m(\mathbf{c}) \right) \quad (42)$$

$$+ \int_{K_+(\mathbf{c})} u(\mathbf{x})(-p_{\theta}(\mathbf{x}|\mathbf{c}) + \delta) d\mathbf{x}, \quad (43)$$

where λ and $u(\cdot)$ are the dual variables.

For each $\mathbf{x} \in K_+(\mathbf{c})$, treating $p_{\theta^*}(\mathbf{x}|\mathbf{c})$ as a pointwise optimization variable, let $p_{\theta^*}(\mathbf{x}|\mathbf{c})$, λ^* and $u^*(\mathbf{x})$ be the corresponding optimal primal and dual variables. Applying stationary condition of the KKT conditions, we have:

$$\nabla_{p_{\theta^*}(\mathbf{x}|\mathbf{c})} \mathcal{L}(p_{\theta^*}(\mathbf{x}|\mathbf{c}), \lambda^*, u^*(\mathbf{x})) = 0 \quad (44)$$

$$\Rightarrow -\frac{h(\mathbf{x}|\mathbf{c})}{p_{\theta^*}(\mathbf{x}|\mathbf{c})} + \lambda^* - u^*(\mathbf{x}) = 0. \quad (45)$$

Based on (45), we consider the following two cases:

- Suppose $u^*(\mathbf{x}) = 0$, we have $p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \frac{h(\mathbf{x}|\mathbf{c})}{\lambda^*}$.
- Suppose $u^*(\mathbf{x}) \neq 0$, by complementary slackness, we have $p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \delta$ and $u^*(\mathbf{x}) = \lambda^* - \frac{h(\mathbf{x}|\mathbf{c})}{\delta}$. Applying dual feasibility $u^*(\mathbf{x}) \geq 0$, we have $h(\mathbf{x}|\mathbf{c}) \leq \lambda^* \delta$.

Note that the above two cases can be combined as:

$$p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \max\left\{\frac{h(\mathbf{x}|\mathbf{c})}{\lambda^*}, \delta\right\}, \forall \mathbf{x} \in K_+(\mathbf{c}). \quad (46)$$

Next, we prove λ^* exists, i.e., (46) is normalizable. By applying the primal feasibility, we have:

$$\int_{K_+(\mathbf{c})} p_{\theta^*}(\mathbf{x}|\mathbf{c}) d\mathbf{x} = m(\mathbf{c}) \quad (47)$$

$$\Rightarrow \int_{K_+(\mathbf{c}) \cap \{h(\mathbf{x}|\mathbf{c}) \leq \lambda^* \delta\}} \delta d\mathbf{x} + \int_{K_+(\mathbf{c}) \cap \{h(\mathbf{x}|\mathbf{c}) > \lambda^* \delta\}} \frac{h(\mathbf{x}|\mathbf{c})}{\lambda^*} d\mathbf{x} = m(\mathbf{c}). \quad (48)$$

Define:

$$A(\lambda) := \int_{K_+(\mathbf{c}) \cap \{h(\mathbf{x}|\mathbf{c}) \leq \lambda \delta\}} \delta d\mathbf{x} + \int_{K_+(\mathbf{c}) \cap \{h(\mathbf{x}|\mathbf{c}) > \lambda \delta\}} \frac{h(\mathbf{x}|\mathbf{c})}{\lambda} d\mathbf{x}. \quad (49)$$

Suppose $0 < \lambda_1 < \lambda_2$, we have:

$$A(\lambda_2) - A(\lambda_1) = \int_{K_+(\mathbf{c}) \cap \{\lambda_1 \delta \leq h(\mathbf{x}|\mathbf{c}) \leq \lambda_2 \delta\}} \delta d\mathbf{x} - \int_{K_+(\mathbf{c}) \cap \{\lambda_1 \delta \leq h(\mathbf{x}|\mathbf{c}) \leq \lambda_2 \delta\}} \frac{h(\mathbf{x}|\mathbf{c})}{\lambda_1} d\mathbf{x} \quad (50)$$

$$+ \int_{K_+(\mathbf{c}) \cap \{h(\mathbf{x}|\mathbf{c}) > \lambda_2 \delta\}} \left(\frac{h(\mathbf{x}|\mathbf{c})}{\lambda_2} - \frac{h(\mathbf{x}|\mathbf{c})}{\lambda_1} \right) d\mathbf{x} \quad (51)$$

$$= \int_{K_+(\mathbf{c}) \cap \{\lambda_1 \delta \leq h(\mathbf{x}|\mathbf{c}) \leq \lambda_2 \delta\}} \left(\delta - \frac{h(\mathbf{x}|\mathbf{c})}{\lambda_1} \right) d\mathbf{x} + \int_{K_+(\mathbf{c}) \cap \{h(\mathbf{x}|\mathbf{c}) > \lambda_2 \delta\}} \left(\frac{h(\mathbf{x}|\mathbf{c})}{\lambda_2} - \frac{h(\mathbf{x}|\mathbf{c})}{\lambda_1} \right) d\mathbf{x} \quad (52)$$

$$< 0, \quad (53)$$

which implies $A(\lambda)$ is a monotonically decreasing function of $\lambda > 0$. By the assumption $\delta < \frac{1}{|K|}$, we have $\delta|K_-(\mathbf{c})| + \delta|K_+(\mathbf{c})| < 1$, which implies $\delta|K_+(\mathbf{c})| < m(\mathbf{c})$.

Since:

$$\lim_{\lambda \rightarrow 0} A(\lambda) = +\infty, \quad \lim_{\lambda \rightarrow +\infty} A(\lambda) = \delta|K_+(\mathbf{c})|, \quad (54)$$

as long as $\delta|K_+(\mathbf{c})| < m(\mathbf{c})$, by the intermediate value theorem, there exists a finite dual optimal point $\lambda^* > 0$ such that $A(\lambda^*) = m(\mathbf{c})$, i.e., the primal feasibility (48) holds. Hence, $p_{\theta^*}(\mathbf{x})$ in (46) is a valid, normalizable probability distribution. Since the exact value of λ^* is dependent on the specific condition \mathbf{c} , we replace it with the notation $Z(\mathbf{c})$, which leads to the final result:

$$p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \begin{cases} \max\left\{\frac{h(\mathbf{x}|\mathbf{c})}{Z(\mathbf{c})}, \delta\right\}, & \mathbf{x} \in K_+(\mathbf{c}), \\ \delta, & \mathbf{x} \in K_-(\mathbf{c}), \\ 0, & \mathbf{x} \notin K. \end{cases} \quad (55)$$

, which can be further simplified as:

$$p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \begin{cases} \max\left\{\frac{h(\mathbf{x}|\mathbf{c})}{\lambda^*}, \delta\right\}, & \mathbf{x} \in K, \\ 0, & \mathbf{x} \notin K. \end{cases} \quad (56)$$

In the limit $\delta \rightarrow 0$, the optimal distribution approaches:

$$p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \frac{h^+(\mathbf{x}|\mathbf{c})}{\int_{\mathbf{x}} h^+(\mathbf{x}|\mathbf{c}) d\mathbf{x}}, \quad h^+(\mathbf{x}|\mathbf{c}) := \max\{h(\mathbf{x}|\mathbf{c}), 0\} \quad (57)$$

which simply zeros out the negative part of $h(\mathbf{x}|\mathbf{c})$ and renormalizes it as a valid distribution. This completes the proof. \square

A.3. Fine-tuning with MCLR

Given a base model $p_{\text{ref}}(\mathbf{x})$ that lacks class specificity, we may fine-tune it using MCLR combined with KL regularization:

$$\max_{\theta} -\mathbb{E}_{\mathbf{c}} [D_{\text{KL}}(p_{\text{ref}}(\mathbf{x}|\mathbf{c})||p_{\theta}(\mathbf{x}|\mathbf{c}))] + \eta \mathbb{E}_{\mathbf{c}, \tilde{\mathbf{c}}, \mathbf{x} \sim p(\mathbf{x}|\mathbf{c})} \left[\log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\theta}(\mathbf{x}|\tilde{\mathbf{c}})} \right]. \quad (58)$$

Similar to (33), we can get the optimal conditional distribution $p(\mathbf{x}|\mathbf{c}_k)$ for each $k \in \{1, \dots, M\}$ by solving the following optimization problem:

$$\arg \max_{p_{\theta}(\cdot|\mathbf{c}_k)} -D_{\text{KL}}(p_{\text{ref}}(\mathbf{x}|\mathbf{c}_k)||p_{\theta}(\mathbf{x}|\mathbf{c}_k)) + \eta \mathbb{E}_{p(\mathbf{x}|\mathbf{c}_k)} [\log p_{\theta}(\mathbf{x}|\mathbf{c}_k)] - \eta \mathbb{E}_{p(\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{c}_k)] \quad (59)$$

$$= \arg \max_{p_{\theta}(\cdot|\mathbf{c}_k)} \int_K \log p_{\theta}(\mathbf{x}|\mathbf{c}_k) (p_{\text{ref}}(\mathbf{x}|\mathbf{c}_k) + \eta(p(\mathbf{x}|\mathbf{c}_k) - p(\mathbf{x}))) d\mathbf{x}. \quad (60)$$

Note that optimization problem (60) shares the same structure as (34), hence by letting $h(\mathbf{x}|\mathbf{c}) := p_{\text{ref}}(\mathbf{x}|\mathbf{c}) + \eta(p(\mathbf{x}|\mathbf{c}) - p(\mathbf{x}))$ and under the same compact-support Assumption 1, we get the same optimal solution as in stated in Theorem 3.1.

Importantly, in the fine-tuning setting, under a mixture error model, MCLR recovers the ground truth conditional distribution, as stated in the following corollary.

Corollary 1. *If the base model satisfies the mixture error model:*

$$p_{\text{ref}}(\mathbf{x}|\mathbf{c}) = (1 - \eta) p(\mathbf{x}|\mathbf{c}) + \eta p(\mathbf{x}), \quad (61)$$

then fine-tuning $p_{\text{ref}}(\mathbf{x}|\mathbf{c})$ with MCLR objective (13) recovers the ground truth conditional distribution $p(\mathbf{x}|\mathbf{c})$.

Proof. Under the mixture-error model (61), the optimization problem (60) becomes:

$$\arg \max_{p_{\theta}(\cdot|\mathbf{c}_k)} \int_K \log p_{\theta}(\mathbf{x}|\mathbf{c}_k) p(\mathbf{x}|\mathbf{c}_k) d\mathbf{x} = \arg \min_{p_{\theta}(\cdot|\mathbf{c}_k)} D_{\text{KL}}(p(\mathbf{x}|\mathbf{c}_k)||p_{\theta}(\mathbf{x}|\mathbf{c}_k)), \quad (62)$$

which has optimal solution $p_{\theta}(\mathbf{x}|\mathbf{c}_k) = p(\mathbf{x}|\mathbf{c}_k)$. Note that in this case, the proof does not depend on Assumption 1 and Assumption 2. This completes the proof. \square

B. Theoretical Analysis of CC-DPO

B.1. Basics of DPO

Reward Modeling. For a given prompt \mathbf{c} and two associated outputs \mathbf{x}_w and \mathbf{x}_l , where the subscripts 'w' stands for 'winning' while 'l' stands for 'losing', implying that \mathbf{x}_w is preferred over \mathbf{x}_l , DPO models the human preference distribution with the Bradley-Terry (BT) model (Bradley & Terry, 1952):

$$p(\mathbf{x}_w \succ \mathbf{x}_l | \mathbf{c}) = \frac{\exp(r^*(\mathbf{x}_w | \mathbf{c}))}{\exp(r^*(\mathbf{x}_w | \mathbf{c})) + \exp(r^*(\mathbf{x}_l | \mathbf{c}))} = \text{Sigmoid}(r^*(\mathbf{x}_w | \mathbf{c}) - r^*(\mathbf{x}_l | \mathbf{c})) \quad (63)$$

where $r^*(\mathbf{x} | \mathbf{c})$ is the underlying optimal reward function for prompt \mathbf{c} . Intuitively, the preferred samples \mathbf{x}_w should have higher reward values compared to the non-preferred samples \mathbf{x}_l . Assuming access to a dataset of sampled comparisons $\mathcal{S} = \{(\mathbf{c}^i, \mathbf{x}_w^i, \mathbf{x}_l^i)\}_{i=1}^N$, one can learn the optimal reward function via maximum likelihood estimation:

$$r^* = \arg \max_r \mathbb{E}_{(\mathbf{c}, \mathbf{x}_w, \mathbf{x}_l) \sim \mathcal{S}} [\log p(\mathbf{x}_w \succ \mathbf{x}_l | \mathbf{c})] \quad (64)$$

$$= \arg \min_r -\mathbb{E}_{(\mathbf{c}, \mathbf{x}_w, \mathbf{x}_l) \sim \mathcal{S}} [\log \text{Sigmoid}(r(\mathbf{x}_w | \mathbf{c}) - r(\mathbf{x}_l | \mathbf{c}))]. \quad (65)$$

RL fine-tuning Phase. Assuming access to the optimal preference reward function r^* , one can fine-tune a base model $p_{\text{ref}}(\mathbf{x} | \mathbf{c})$ to align with the preference dataset by optimizing the following objective:

$$\max_{\theta} \mathbb{E}_{\mathbf{c}} [\mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x} | \mathbf{c})} [r^*(\mathbf{x} | \mathbf{c})] - \beta D_{KL}(p_{\theta}(\mathbf{x} | \mathbf{c}) || p_{\text{ref}}(\mathbf{x} | \mathbf{c}))], \quad (66)$$

where β controls the KL-regularization strength. Intuitively, objective (66) encourages the fine-tuned model to achieve high reward value in expectation, and at the same time not deviate too much from the base model.

DPO Objective. Under certain regularity conditions, the optimal solution to the fine-tuning objective (66) admits the following closed-form:

$$p_{\theta^*}(\mathbf{x} | \mathbf{c}) = \frac{1}{Z(\mathbf{c})} p_{\text{ref}}(\mathbf{x} | \mathbf{c}) \exp\left(\frac{1}{\beta} r^*(\mathbf{x} | \mathbf{c})\right), \quad (67)$$

where $Z(\mathbf{c})$ is a partition function for normalizing the density. To prove this, consider optimizing the fine-tuning objective (66) for a target condition \mathbf{c} :

$$\arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x} | \mathbf{c})} [r^*(\mathbf{x} | \mathbf{c})] - \beta D_{KL}(p_{\theta}(\mathbf{x} | \mathbf{c}) || p_{\text{ref}}(\mathbf{x} | \mathbf{c})) \quad (68)$$

$$= \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x} | \mathbf{c})} \left[r^*(\mathbf{x} | \mathbf{c}) - \beta \log \frac{p_{\theta}(\mathbf{x} | \mathbf{c})}{p_{\text{ref}}(\mathbf{x} | \mathbf{c})} \right] \quad (69)$$

$$= \arg \min_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x} | \mathbf{c})} \left[\log \frac{p_{\theta}(\mathbf{x} | \mathbf{c})}{p_{\text{ref}}(\mathbf{x} | \mathbf{c})} - \frac{1}{\beta} r^*(\mathbf{x} | \mathbf{c}) \right] \quad (70)$$

$$= \arg \min_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\theta}(\mathbf{x} | \mathbf{c})} \left[\log \frac{p_{\theta}(\mathbf{x} | \mathbf{c})}{\frac{1}{Z(\mathbf{c})} p_{\text{ref}}(\mathbf{x} | \mathbf{c}) \exp\left(\frac{1}{\beta} r^*(\mathbf{x} | \mathbf{c})\right)} - \log Z(\mathbf{c}) \right], \quad (71)$$

where $Z(\mathbf{c})$ is the partition function that normalizes $p_{\text{ref}}(\mathbf{x} | \mathbf{c}) \exp\left(\frac{1}{\beta} r^*(\mathbf{x} | \mathbf{c})\right)$:

$$Z(\mathbf{c}) = \int_{\mathbf{x}} p_{\text{ref}}(\mathbf{x} | \mathbf{c}) \exp\left(\frac{1}{\beta} r^*(\mathbf{x} | \mathbf{c})\right) d\mathbf{x}, \quad (72)$$

such that $p^*(\mathbf{x} | \mathbf{c}) := \frac{1}{Z(\mathbf{c})} p_{\text{ref}}(\mathbf{x} | \mathbf{c}) \exp\left(\frac{1}{\beta} r^*(\mathbf{x} | \mathbf{c})\right)$ is a valid probability distribution. Since $Z(\mathbf{c})$ is independent of θ , the fine-tuning objective (71) is equivalent to:

$$\min_{\theta} D_{KL}(p_{\theta}(\mathbf{x} | \mathbf{c}) || p^*(\mathbf{x} | \mathbf{c})), \quad (73)$$

which achieves its minimum value 0 if and only if:

$$p_{\theta}(\mathbf{x}|\mathbf{c}) = p^*(\mathbf{x}|\mathbf{c}) = \frac{1}{Z(\mathbf{c})} p_{\text{ref}}(\mathbf{x}|\mathbf{c}) \exp\left(\frac{1}{\beta} r^*(\mathbf{x}|\mathbf{c})\right). \quad (74)$$

With some algebra, the optimal reward function can be expressed with $p_{\theta^*}(\mathbf{x}|\mathbf{c})$:

$$r^*(\mathbf{x}|\mathbf{c}) = \beta \log \frac{p_{\theta^*}(\mathbf{x}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}|\mathbf{c})} + \beta \log Z(\mathbf{c}). \quad (75)$$

The relationship (75) between the optimal reward and the optimal fine-tuned distribution suggests a convenient parameterization of the reward model:

$$r_{\theta}(\mathbf{x}|\mathbf{c}) = \beta \log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}|\mathbf{c})}. \quad (76)$$

Substitute (76) into (65) results in the DPO objective:

$$\arg \min_{\theta} -\mathbb{E}_{(\mathbf{c}, \mathbf{x}_w, \mathbf{x}_l) \sim S} \left[\log \text{Sigmoid} \left(\beta \log \frac{p_{\theta}(\mathbf{x}_w|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_w|\mathbf{c})} - \beta \log \frac{p_{\theta}(\mathbf{x}_l|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_l|\mathbf{c})} \right) \right]. \quad (77)$$

In this way, one can directly fine-tune the base model without explicitly modeling the reward.

B.2. Improving Conditional Modeling with CC-DPO

To adapt DPO for improving class specificity, we may treat samples from the target class \mathbf{c} as preferred data (\mathbf{x}_w) and samples from other randomly selected classes as non-preferred data (\mathbf{x}_l). This leads to the following objective:

$$\min_{\theta} -\mathbb{E}_{\mathbf{c}, \tilde{\mathbf{c}}, \mathbf{x}_w \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_l \sim p(\mathbf{x}|\tilde{\mathbf{c}})} \left[\log \text{Sigmoid} \left(\beta \log \frac{p_{\theta}(\mathbf{x}_w|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_w|\mathbf{c})} - \beta \log \frac{p_{\theta}(\mathbf{x}_l|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_l|\mathbf{c})} \right) \right] \quad (78)$$

$$= \min_{\theta} -\mathbb{E}_{\mathbf{c}, \mathbf{x}_w \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_l \sim p(\mathbf{x})} \left[\log \text{Sigmoid} \left(\beta \log \frac{p_{\theta}(\mathbf{x}_w|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_w|\mathbf{c})} - \beta \log \frac{p_{\theta}(\mathbf{x}_l|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}_l|\mathbf{c})} \right) \right], \quad (79)$$

which admits a closed-form solution as stated in the following theorem.

Theorem 3.2. *Under certain regularity conditions, the optimal solution to (79) is:*

$$p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \frac{1}{\tilde{Z}(\mathbf{c})} p_{\text{ref}}(\mathbf{x}|\mathbf{c}) \left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})} \right)^{\frac{1}{\beta}}, \quad (80)$$

where $\tilde{Z}(\mathbf{c}) = \int_{\mathbf{x}} p_{\text{ref}}(\mathbf{x}|\mathbf{c}) \left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})} \right)^{\frac{1}{\beta}} d\mathbf{x}$ is the normalizing constant.

Proof. From (65) to (67), it is clear that the optimal solution to CC-DPO (79) is fully determined by the base model and the optimal solution to the following reward modeling objective:

$$r^* = \arg \min_r -\mathbb{E}_{\mathbf{c}, \mathbf{x}_w \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_l \sim p(\mathbf{x})} \left[\log \text{Sigmoid}(r(\mathbf{x}_w|\mathbf{c}) - r(\mathbf{x}_l|\mathbf{c})) \right]. \quad (81)$$

Note that r^* is the collection of optimal rewards $r^*(\cdot|\mathbf{c}_k)$ for each $\mathbf{c}_k \in \{\mathbf{c}_i\}_{i=1}^M$. Without loss of generality, we drop the subscript and solve for the optimal reward for a single target \mathbf{c} :

$$r^*(\cdot|\mathbf{c}) = \arg \min_{r(\cdot|\mathbf{c})} -\mathbb{E}_{\mathbf{x}_w \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_l \sim p(\mathbf{x})} \left[\log \text{Sigmoid}(r(\mathbf{x}_w|\mathbf{c}) - r(\mathbf{x}_l|\mathbf{c})) \right] \quad (82)$$

$$= \arg \max_{r(\cdot|\mathbf{c})} \int_{\mathbf{x}_w} \int_{\mathbf{x}_l} \log \text{Sigmoid}(r(\mathbf{x}_w|\mathbf{c}) - r(\mathbf{x}_l|\mathbf{c})) p(\mathbf{x}_w|\mathbf{c}) p(\mathbf{x}_l) d\mathbf{x}_w d\mathbf{x}_l. \quad (83)$$

Note that for any arbitrary pair of points $(\mathbf{x}_1, \mathbf{x}_2)$, they contribute to the integral (83) for the following amount:

$$\log \text{Sigmoid}(r(\mathbf{x}_1|\mathbf{c}) - r(\mathbf{x}_2|\mathbf{c})) p(\mathbf{x}_1|\mathbf{c}) p(\mathbf{x}_2) + \log \text{Sigmoid}(r(\mathbf{x}_2|\mathbf{c}) - r(\mathbf{x}_1|\mathbf{c})) p(\mathbf{x}_2|\mathbf{c}) p(\mathbf{x}_1), \quad (84)$$

hence $r^*(\cdot|\mathbf{c})$ is an optimal solution to (83) if it maximizes (84) for $\forall(\mathbf{x}_1, \mathbf{x}_2)$. To find such $r^*(\cdot|\mathbf{c})$, let's define $w = r(\mathbf{x}_1|\mathbf{c}) - r(\mathbf{x}_2|\mathbf{c})$ and solve the following optimization problem:

$$\max_w \mathcal{L}(w) := \log \text{Sigmoid}(w)p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2) + \log \text{Sigmoid}(-w)p(\mathbf{x}_2|\mathbf{c})p(\mathbf{x}_1). \quad (85)$$

Note that:

$$\nabla_w \mathcal{L}(w) = \text{Sigmoid}(-w)p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2) - \text{Sigmoid}(w)p(\mathbf{x}_2|\mathbf{c})p(\mathbf{x}_1), \quad (86)$$

which implies the stationary point w^* must satisfy:

$$\text{Sigmoid}(w^*)p(\mathbf{x}_2|\mathbf{c})p(\mathbf{x}_1) = \text{Sigmoid}(-w^*)p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2) \quad (87)$$

$$\Rightarrow \text{Sigmoid}(w^*)(p(\mathbf{x}_2|\mathbf{c})p(\mathbf{x}_1) + p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2)) = p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2) \quad (88)$$

$$\Rightarrow \text{Sigmoid}(w^*) = \frac{p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2)}{p(\mathbf{x}_2|\mathbf{c})p(\mathbf{x}_1) + p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2)} \quad (89)$$

$$\Rightarrow \text{Sigmoid}(w^*) = \frac{1}{1 + \frac{p(\mathbf{x}_2|\mathbf{c})p(\mathbf{x}_1)}{p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2)}} \quad (90)$$

$$\Rightarrow \text{Sigmoid}(w^*) = \frac{1}{1 + \exp\left(-\log \frac{p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2)}{p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{c})}\right)} \quad (91)$$

$$\Rightarrow w^* = \log \frac{p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2)}{p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{c})}, \quad (92)$$

which implies:

$$r^*(\mathbf{x}|\mathbf{c}) = \log \frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})} + K, \quad (93)$$

where K is any finite constant. Moreover, since:

$$\nabla_w^2 \mathcal{L}(w) = -\text{Sigmoid}(w)\text{Sigmoid}(-w)p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2) - \text{Sigmoid}(-w)\text{Sigmoid}(w)p(\mathbf{x}_2|\mathbf{c})p(\mathbf{x}_1) \quad (94)$$

$$< 0, \quad (95)$$

we know $\mathcal{L}(w)$ is concave and w^* is the global maximizer and consequently $r^*(\mathbf{x}|\mathbf{c})$ is the optimal reward function. Since this optimal reward function is consistent for any arbitrary pair of points $(\mathbf{x}_1, \mathbf{x}_2)$, it is the global reward function that maximizes the full objective (82).

Substitute (93) into (67), we obtain the optimal solution to CC-DPO:

$$p_{\theta^*}(\mathbf{x}|\mathbf{c}) = \frac{1}{\tilde{Z}(\mathbf{c})} p_{\text{ref}}(\mathbf{x}|\mathbf{c}) \left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})} \right)^{\frac{1}{\beta}}, \quad (96)$$

where $\tilde{Z}(\mathbf{c}) = \int_{\mathbf{x}} p_{\text{ref}}(\mathbf{x}|\mathbf{c}) \left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})} \right)^{\frac{1}{\beta}} d\mathbf{x}$ is the normalizing constant. This completes the proof. \square

Remark. Note that from (88) to (89), we require $p(\mathbf{x}_2|\mathbf{c})p(\mathbf{x}_1) + p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2) > 0$ for $\forall \mathbf{x}_1, \mathbf{x}_2$, which holds true if both $p(\mathbf{x}|\mathbf{c})$ and $p(\mathbf{x})$ have full support on \mathbb{R}^d . Next, we discuss the corner cases where this assumption doesn't hold.

1. Suppose $p(\mathbf{x}_1|\mathbf{c}) = 0$ but $p(\mathbf{x}_1) > 0$, then we have:

$$\mathcal{L}(w) = \log \text{Sigmoid}(r(\mathbf{x}_2|\mathbf{c}) - r(\mathbf{x}_1|\mathbf{c}))p(\mathbf{x}_2|\mathbf{c})p(\mathbf{x}_1). \quad (97)$$

In this case, the objective is maximized if $r^*(\mathbf{x}_1|\mathbf{c}) = -\infty$.

2. Suppose $p(\mathbf{x}_1|\mathbf{c}) > 0$ but $p(\mathbf{x}_1) = 0$, then we have:

$$\mathcal{L}(w) = \log \text{Sigmoid}(r(\mathbf{x}_1|\mathbf{c}) - r(\mathbf{x}_2|\mathbf{c}))p(\mathbf{x}_1|\mathbf{c})p(\mathbf{x}_2). \quad (98)$$

In this case, the objective is maximized if $r^*(\mathbf{x}_1|\mathbf{c}) = +\infty$.

3. Suppose $p(\mathbf{x}_1|\mathbf{c}) = 0$ and $p(\mathbf{x}_1) = 0$, then \mathbf{x}_1 will never be sampled and it does not contribute to the overall objective. In such case, $r^*(\mathbf{x}_1|\mathbf{c})$ is undefined and will implicitly depend on the practical parameterization of the reward model.

The first two cases are already covered by (80). In particular, case 2 can lead to non-normalizable issue, as in this case $p_{\theta^*}(\cdot|\mathbf{c})$ will have infinite density at \mathbf{x}_1 . Although since $p(\mathbf{x}_1) = \mathbb{E}_{\mathbf{c}}[p(\mathbf{x}_1|\mathbf{c})]$, it will not be zero given $p(\mathbf{x}_1|\mathbf{c}) > 0$. It is highly likely in practice, there exists \mathbf{x}_1 such that $p(\mathbf{x}_1|\mathbf{c}) > 0$ but $p(\mathbf{x}_1) \approx 0$. In this regime, the ratio $\left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})}\right)^{\frac{1}{\beta}}$ blows up, forcing $p_{\theta^*}(\cdot|\mathbf{c})$ to essentially place all its mass on such \mathbf{x} , which again leads to non-normalizable issue.

Error Model for CC-DPO. As with the MCLR case, CC-DPO recovers the ground truth conditional distribution under an appropriate error model, as stated in the following corollary.

Corollary 2. If the base model satisfies:

$$p_{\text{ref}}(\mathbf{x}|\mathbf{c}) \propto p(\mathbf{x}|\mathbf{c})^{1-\frac{1}{\beta}} p(\mathbf{x})^{\frac{1}{\beta}} \quad (99)$$

then fine-tuning $p_{\text{ref}}(\mathbf{x}|\mathbf{c})$ with the CC-DPO objective (18) recovers the ground-truth conditional distribution $p(\mathbf{x}|\mathbf{c})$.

Proof. Substitute (99) into (96), we get $p_{\theta^*}(\mathbf{x}|\mathbf{c}) = p(\mathbf{x}|\mathbf{c})$, which completes the proof. \square

C. Theoretical Analysis of CCA

In Appendix B.2, we’ve shown that the underlying optimal reward function induced by the CC-DPO objective (81) has the form of the log likelihood-ratio (93). Interestingly, the same reward function can be obtained by minimizing the following optimization problem:

$$r^*(\cdot|\mathbf{c}) = \arg \max_{r(\cdot|\mathbf{c})} \mathbb{E}_{p(\mathbf{x}|\mathbf{c})} \log \text{Sigmoid}(r(\mathbf{x}|\mathbf{c})) + \mathbb{E}_{p(\mathbf{x})} \log \text{Sigmoid}(-r(\mathbf{x}|\mathbf{c})) \quad (100)$$

$$= \arg \max_{r(\cdot|\mathbf{c})} \int_{\mathbf{x}} (\log \text{Sigmoid}(r(\mathbf{x}|\mathbf{c}))p(\mathbf{x}|\mathbf{c}) + \log \text{Sigmoid}(-r(\mathbf{x}|\mathbf{c}))p(\mathbf{x})) d\mathbf{x}. \quad (101)$$

Since the objective (101) decomposes pointwise over \mathbf{x} , the optimal reward function can be obtained by maximizing the integrand for each \mathbf{x} independently:

$$\log \text{Sigmoid}(r(\mathbf{x}|\mathbf{c}))p(\mathbf{x}|\mathbf{c}) + \log \text{Sigmoid}(-r(\mathbf{x}|\mathbf{c}))p(\mathbf{x}). \quad (102)$$

To find $r^*(\cdot|\mathbf{c})$, let’s define $w = r(\mathbf{x}|\mathbf{c})$ and solve the following optimization problem:

$$\max_w \mathcal{L}(w) := \log \text{Sigmoid}(w)p(\mathbf{x}|\mathbf{c}) + \log \text{Sigmoid}(-w)p(\mathbf{x}). \quad (103)$$

Note that:

$$\nabla_w \mathcal{L}(w) := \text{Sigmoid}(-w)p(\mathbf{x}|\mathbf{c}) - \text{Sigmoid}(w)p(\mathbf{x}), \quad (104)$$

which implies the stationary point w^* is:

$$w^* = r^*(\mathbf{x}|\mathbf{c}) = \log \frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})}. \quad (105)$$

Furthermore, since:

$$\nabla_w^2 \mathcal{L}(w) = -\text{Sigmoid}(w)\text{Sigmoid}(-w)p(\mathbf{x}|\mathbf{c}) - \text{Sigmoid}(w)\text{Sigmoid}(-w)p(\mathbf{x}) \quad (106)$$

$$< 0, \quad (107)$$

we know $\mathcal{L}(w)$ is concave and $r^*(\mathbf{x}|\mathbf{c}) = \log \frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})}$ is the unique maximizer. Hence, we’ve proved that the optimization problems (81) and (101) lead to the same optimal reward up to an additive constant.

Therefore, by parameterizing the reward function as:

$$r_{\theta}(\mathbf{x}|\mathbf{c}) = \beta \log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}|\mathbf{c})} \quad (108)$$

, and substitute this in (101), we get the following optimization problem:

$$p_{\theta^*}(\cdot|\mathbf{c}) = \arg \max_{p_{\theta}(\cdot|\mathbf{c})} \mathbb{E}_{p(\mathbf{x}|\mathbf{c})} \log \text{Sigmoid}\left(\beta \log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}|\mathbf{c})}\right) + \mathbb{E}_{p(\mathbf{x})} \log \text{Sigmoid}\left(-\beta \log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}|\mathbf{c})}\right), \quad (109)$$

which has the global maximizer:

$$\beta \log \frac{p_{\theta^*}(\mathbf{x}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}|\mathbf{c})} = r^*(\mathbf{x}|\mathbf{c}) = \log \frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})} \quad (110)$$

$$\Rightarrow p_{\theta^*}(\mathbf{x}|\mathbf{c}) = p_{\text{ref}}(\mathbf{x}|\mathbf{c}) \left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})}\right)^{\frac{1}{\beta}}. \quad (111)$$

One issue with the optimization objective (109) is that its optimal solution (111) is not a valid probability distribution because it is not normalized. To alleviate this issue, we consider a slightly modified version of (101) by introducing an additional constant λ :

$$r^*(\cdot|\mathbf{c}) = \arg \max_{r(\cdot|\mathbf{c})} \mathbb{E}_{p(\mathbf{x}|\mathbf{c})} \log \text{Sigmoid}(r(\mathbf{x}|\mathbf{c})) + \lambda \mathbb{E}_{p(\mathbf{x})} \log \text{Sigmoid}(-r(\mathbf{x}|\mathbf{c})). \quad (112)$$

With similar proof technique, it can be shown that:

$$r^*(\mathbf{x}|\mathbf{c}) = \log \frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})} + \log \frac{1}{\lambda}. \quad (113)$$

Substitute (108) in (112), we get the following optimization problem:

$$p_{\theta^*}(\cdot|\mathbf{c}) = \arg \max_{p_{\theta}(\cdot|\mathbf{c})} \mathbb{E}_{p(\mathbf{x}|\mathbf{c})} \log \text{Sigmoid}\left(\beta \log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}|\mathbf{c})}\right) + \lambda \mathbb{E}_{p(\mathbf{x})} \log \text{Sigmoid}\left(-\beta \log \frac{p_{\theta}(\mathbf{x}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}|\mathbf{c})}\right), \quad (114)$$

which has global maximizer:

$$\beta \log \frac{p_{\theta^*}(\mathbf{x}|\mathbf{c})}{p_{\text{ref}}(\mathbf{x}|\mathbf{c})} = \log \frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})} + \log \frac{1}{\lambda} \quad (115)$$

$$\Rightarrow p_{\theta^*}(\mathbf{x}|\mathbf{c}) = p_{\text{ref}}(\mathbf{x}|\mathbf{c}) \left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})}\right)^{\frac{1}{\beta}} \left(\frac{1}{\lambda}\right)^{\frac{1}{\beta}}. \quad (116)$$

In this case, $p_{\theta^*}(\mathbf{x}|\mathbf{c})$ is a valid probability distribution as long as:

$$\lambda^{\frac{1}{\beta}} = \int_{\mathbf{x}} p_{\text{ref}}(\mathbf{x}|\mathbf{c}) \left(\frac{p(\mathbf{x}|\mathbf{c})}{p(\mathbf{x})}\right)^{\frac{1}{\beta}} d\mathbf{x}. \quad (117)$$

Note that optimal solution to (114) exactly coincides with that of CC-DPO (18). The formulation in (114), known as the Conditional Contrastive Alignment (CCA), which is essentially a combination of the Noise Contrastive Estimation (NCE) (Gutmann & Hyvärinen, 2010) and a special model parameterization (113). This objective is first proposed by (Chen et al., 2025b) for improving the generation quality of visual autoregressive models without relying on CFG. It should be noted that CCA is theoretically correct only when λ is fixed to a particular constant, and it must be tuned as a hyperparameter in practice. In contrast, we demonstrate that CC-DPO fine-tuning recovers the same optimal solution without introducing this additional hyperparameter λ .

D. Theoretical Analysis of the Equivalence between CFG and Weighted MCLR

D.1. Proof of Theorem 4.1

In this section, we provide the proof for Theorem 4.1. To begin with, we introduce the following lemma:

Lemma D.1. *Let $p(\mathbf{x})$ be the clean data distribution, $p_t(\mathbf{x})$ be the marginal distribution of $\mathbf{x}(t)$, and $p_{0t}(\mathbf{x}_t|\mathbf{x})$ be the transition density from $\mathbf{x}(0)$ to $\mathbf{x}(t)$ as defined in Section 2.1, then*

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \mathbb{E}_{p(\mathbf{x}|\mathbf{x}_t)} [\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x})]. \quad (118)$$

Proof.

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \frac{\nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t)}{p_t(\mathbf{x}_t)} \quad (119)$$

$$= \frac{\nabla_{\mathbf{x}_t} \int_{\mathbf{x}} p(\mathbf{x}) p_{0t}(\mathbf{x}_t|\mathbf{x}) d\mathbf{x}}{p_t(\mathbf{x}_t)} \quad (120)$$

$$= \int_{\mathbf{x}} \frac{p(\mathbf{x}) p_{0t}(\mathbf{x}_t|\mathbf{x}) \nabla_{\mathbf{x}_t} p_{0t}(\mathbf{x}_t|\mathbf{x})}{p_t(\mathbf{x}_t) p_{0t}(\mathbf{x}_t|\mathbf{x})} d\mathbf{x} \quad (121)$$

$$= \int_{\mathbf{x}} p(\mathbf{x}|\mathbf{x}_t) \nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) d\mathbf{x} \quad (122)$$

$$= \mathbb{E}_{p(\mathbf{x}|\mathbf{x}_t)} [\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x})]. \quad (123)$$

This completes the proof. \square

We now restate the main theorem and proceed with the proof.

Theorem 4.1. *For any time sampling distribution $p(t)$ and weighting function $w(t)$, the CFG-guided score*

$$\mathbf{s}_{\text{cfg}}(\mathbf{x}_t, t, \mathbf{c}) := \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) + \eta (\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t))$$

is the unique minimizer of a sample-adaptive weighted ELBO-approximated MCLR objective (124):

$$\begin{aligned} \mathbf{s}_{\text{cfg}}(\cdot) = \arg \min_{\mathbf{s}_{\theta}(\cdot)} & \mathbb{E}_{\mathbf{c}, t \sim p(t), \mathbf{x} \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})} \left[w(t) \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right] \\ & + \eta \mathbb{E}_{\substack{\mathbf{c}, \tilde{\mathbf{c}}, t \sim p(t) \\ \mathbf{x} \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})}} \left[w(t) \left(\|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right. \right. \\ & \left. \left. - \frac{p_t(\mathbf{x}_t|\tilde{\mathbf{c}})}{p_t(\mathbf{x}_t)} \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t, \tilde{\mathbf{c}})\|_2^2 \right) \right]. \end{aligned} \quad (124)$$

Proof. First, note optimization problem (124) is equivalent to:

$$\min_{\mathbf{s}_{\theta}(\cdot)} (1 + \eta) \mathbb{E}_{\mathbf{c}, t \sim p(t), \mathbf{x} \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})} \left[w(t) \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right] \quad (125)$$

$$- \eta \mathbb{E}_{\substack{\mathbf{c}, \tilde{\mathbf{c}}, t \sim p(t) \\ \mathbf{x} \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})}} \left[w(t) \frac{p_t(\mathbf{x}_t|\tilde{\mathbf{c}})}{p_t(\mathbf{x}_t)} \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t, \tilde{\mathbf{c}})\|_2^2 \right]. \quad (126)$$

Note that (126) can be further simplified as:

$$- \eta \mathbb{E}_{\tilde{\mathbf{c}}, t \sim p(t), \mathbf{x} \sim p(\mathbf{x}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})} \left[w(t) \frac{p_t(\mathbf{x}_t|\tilde{\mathbf{c}})}{p_t(\mathbf{x}_t)} \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t, \tilde{\mathbf{c}})\|_2^2 \right]. \quad (127)$$

Hence the original optimization problem (124) is equivalent to

$$\begin{aligned} \min_{\mathbf{s}_{\theta}(\cdot)} & (1 + \eta) \mathbb{E}_{\mathbf{c}, t \sim p(t), \mathbf{x} \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})} \left[w(t) \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right] \\ & - \eta \mathbb{E}_{\mathbf{c}, t \sim p(t), \mathbf{x} \sim p(\mathbf{x}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})} \left[w(t) \frac{p_t(\mathbf{x}_t|\mathbf{c})}{p_t(\mathbf{x}_t)} \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right], \end{aligned}$$

which coincides exactly with the objective presented in the main theorem (20).

In order to solve (128), since the objective decomposes across t and \mathbf{c} , the optimization can be performed independently for each pair (t, \mathbf{c}) . Hence, we can fix a specific pair of t and \mathbf{c} , and optimize the corresponding score $\mathbf{s}_\theta(\mathbf{x}, t, \mathbf{c})$ independently:

$$\begin{aligned} \min_{\mathbf{s}_\theta(\cdot, t, \mathbf{c})} & (1 + \eta) \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{c}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})} [\|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2] \\ & - \eta \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t|\mathbf{x})} \left[\frac{p_t(\mathbf{x}_t|\mathbf{c})}{p_t(\mathbf{x}_t)} \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right] \end{aligned} \quad (128)$$

$$\begin{aligned} \Leftrightarrow \min_{\mathbf{s}_\theta(\cdot, t, \mathbf{c})} & (1 + \eta) \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{c}), \mathbf{x} \sim p(\mathbf{x}|\mathbf{x}_t, \mathbf{c})} [\|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2] \\ & - \eta \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t), \mathbf{x} \sim p(\mathbf{x}|\mathbf{x}_t)} \left[\frac{p_t(\mathbf{x}_t|\mathbf{c})}{p_t(\mathbf{x}_t)} \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right] \end{aligned} \quad (129)$$

$$\begin{aligned} \Leftrightarrow \min_{\mathbf{s}_\theta(\cdot, t, \mathbf{c})} & (1 + \eta) \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{c}), \mathbf{x} \sim p(\mathbf{x}|\mathbf{x}_t, \mathbf{c})} [\|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2] \\ & - \eta \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{c}), \mathbf{x} \sim p(\mathbf{x}|\mathbf{x}_t)} [\|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2]. \end{aligned} \quad (130)$$

From (128) to (129) we use the Bayes rule: $p(\mathbf{x}|\mathbf{c})p_{0t}(\mathbf{x}_t|\mathbf{x}) = p(\mathbf{x}, \mathbf{x}_t|\mathbf{c}) = p_t(\mathbf{x}_t|\mathbf{c})p(\mathbf{x}|\mathbf{x}_t, \mathbf{c})$. From (129) to (130) we use the identity: $p_t(\mathbf{x}_t)p(\mathbf{x}|\mathbf{x}_t) \frac{p_t(\mathbf{x}_t|\mathbf{c})}{p_t(\mathbf{x}_t)} = p_t(\mathbf{x}_t|\mathbf{c})p(\mathbf{x}|\mathbf{x}_t)$.

Note that:

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{x}_t, \mathbf{c})} [\|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2] \\ = \|\mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2 - 2\mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})^T \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{x}_t, \mathbf{c})} [\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x})] + C_1 \\ = \|\mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2 - 2\mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})^T \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) + C_1, \end{aligned} \quad (131)$$

where C_1 is a constant independent of θ , and the second equality follows from Lemma D.1. Similarly, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|\mathbf{x}_t)} [\|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t|\mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2] \\ = \|\mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2 - 2\mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})^T \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) + C_2, \end{aligned} \quad (132)$$

, where C_2 is a constant independent of θ .

Substituting (131) and (132) into (130), and let $\mathbf{v}_1 = (1 + \eta)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) - \eta\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$, the optimization problem becomes equivalent to:

$$\min_{\mathbf{s}_\theta(\cdot, t, \mathbf{c})} \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{c})} [\|\mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})\|_2^2 - 2\mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})^T \mathbf{v}_1] \quad (133)$$

$$\Leftrightarrow \min_{\mathbf{s}_\theta(\cdot, t, \mathbf{c})} \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t|\mathbf{c})} [\|\mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c}) - \mathbf{s}_{\text{cfg}}(\mathbf{x}_t, t, \mathbf{c})\|_2^2] + C_3, \quad (134)$$

where C_3 is a constant independent of θ . Therefore, the optimal solution is:

$$\mathbf{s}_{\theta^*}(\mathbf{x}_t, t, \mathbf{c}) = \mathbf{s}_{\text{cfg}}(\mathbf{x}_t, t, \mathbf{c}) \quad (135)$$

$$= \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) + \eta(\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)). \quad (136)$$

This completes the proof. \square

D.2. Extensions: CFG Variants under the Alignment Framework

The equivalence between CFG and weighted MCLR provides a unified perspective for interpreting CFG variants. This perspective is formalized in the following corollary.

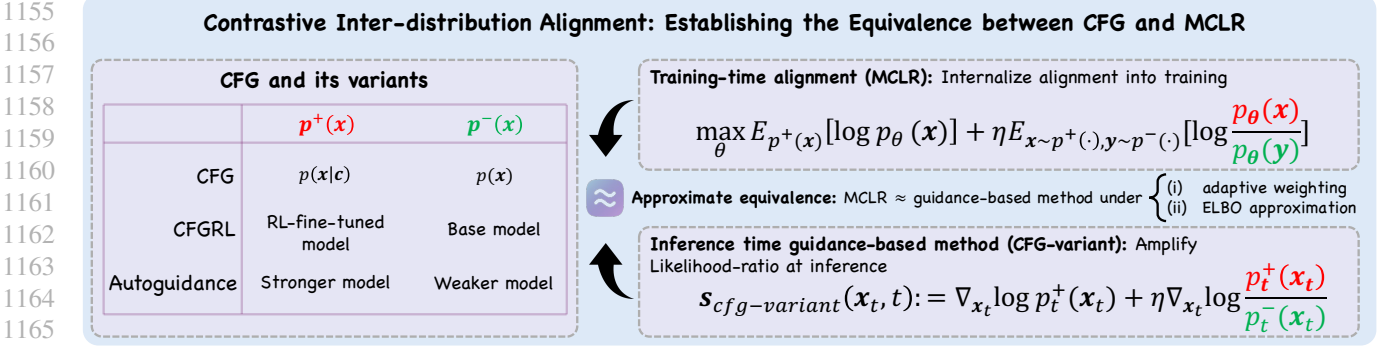


Figure 4. **A Unified Framework Connecting CFG Variants with Contrastive Alignment.** CFG-based methods can be interpreted as implicitly optimizing a likelihood-ratio-based contrastive alignment objective between two distributions $p^+(\mathbf{x})$ and $p^-(\mathbf{x})$ at inference time.

Corollary 3. Consider two distributions $p^+(\mathbf{x})$ and $p^-(\mathbf{x})$. For any time sampling distribution $p(t)$ and weighting function $w(t)$, a generalized CFG-style score:

$$\mathbf{s}_{\text{cfg-variant}}(\mathbf{x}_t, t) := \nabla_{\mathbf{x}_t} \log p_t^+(\mathbf{x}_t) + \eta (\nabla_{\mathbf{x}_t} \log p_t^+(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log p_t^-(\mathbf{x}_t)) \quad (137)$$

is the unique minimizer of the following MCLR-style optimization problem:

$$\begin{aligned} \min_{\mathbf{s}_{\theta}(\cdot)} & (1 + \eta) \mathbb{E}_{t \sim p(t), \mathbf{x} \sim p^+(\mathbf{x}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t | \mathbf{x})} \left[w(t) \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t)\|_2^2 \right] \\ & - \eta \mathbb{E}_{t \sim p(t), \mathbf{x} \sim p^-(\mathbf{x}), \mathbf{x}_t \sim p_{0t}(\mathbf{x}_t | \mathbf{x})} \left[w(t) \frac{p_t^+(\mathbf{x}_t)}{p_t^-(\mathbf{x}_t)} \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t)\|_2^2 \right] \end{aligned} \quad (138)$$

The proof is omitted, as it follows directly by adapting the proof of Theorem 4.1. The optimization problem (138) shares the same structure as the ELBO-approximated weighted MCLR (128), with the only difference that positive samples are drawn from $p^+(\mathbf{x})$ and negative samples from $p^-(\mathbf{x})$. This formulation unifies a broad class of CFG-style methods, as illustrated below.

Standard CFG. Let $p^+(\mathbf{x}) = p(\mathbf{x}|c)$ and $p^-(\mathbf{x}) = p(\mathbf{x})$. This recovers the standard CFG and corresponds to the weighted MCLR formulation discussed in previous sections.

Autoguidance. Let $p^+(\mathbf{x})$ be the distribution induced by a strong diffusion model and $p^-(\mathbf{x})$ the distribution induced by a weaker model. We recover the Autoguidance (Karras et al., 2024a).

Inference-Time Alignment Guidance. Let $p^+(\mathbf{x})$ be the distribution induced by a diffusion model fine-tuned on high-reward data and $p^-(\mathbf{x})$ the distribution induced by the base model (or a low-reward fine-tuned model). We recover a family of inference-time alignment methods (Frans et al., 2025; Jin et al., 2025b; Cheng et al., 2025; Jiang et al., 2026).

Lastly, note that (138) without adaptive weighting is the ELBO-based approximation of the following likelihood-based objective, which is presented in Figure 4:

$$\max_{\theta} \underbrace{\mathbb{E}_{p^+(\mathbf{x})}[\log p_{\theta}(\mathbf{x})] + \eta \mathbb{E}_{\mathbf{x} \sim p^+(\cdot), \mathbf{y} \sim p^-(\cdot)} \left[\log \frac{p_{\theta}(\mathbf{x})}{p_{\theta}(\mathbf{y})} \right]}_{\text{MCLR Regularization}}. \quad (139)$$

E. Practical Implementation Details

E.1. Approximating Log-Likelihood with ELBO

Implementing MCLR requires access to the log-likelihood, which is not directly available for diffusion models. We therefore approximate the log-likelihood using the evidence lower bound (ELBO) in (5). In the following, we describe how this approximation is used in MCLR, CC-DPO, and CCA. Before doing so, we state the following fact.

Equivalence between Score Function and MMSE Denoisers. For practical diffusion models, the drift coefficient $f(\cdot)$ in (1) takes the form $f(\mathbf{x}, t) = f(t)\mathbf{x}$ where $f(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$. As a result, the corresponding transition distribution is Gaussian and can be written as (Karras et al., 2022):

$$p_{0t}(\mathbf{x}_t | \mathbf{x}) = \mathcal{N}(\mathbf{x}_t; s(t)\mathbf{x}, s^2(t)\sigma^2(t)\mathbf{I}) \quad (140)$$

, where $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ denotes the Gaussian density with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ evaluated at \mathbf{x} , $s(t) = \exp(\int_0^t f(\xi)d\xi)$ and $\sigma(t) = \sqrt{\int_0^t \frac{g^2(\xi)}{s^2(\xi)} d\xi}$. The score of this transition distribution is therefore given by:

$$\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}) = \frac{s(t)\mathbf{x} - \mathbf{x}_t}{s^2(t)\sigma^2(t)}. \quad (141)$$

Accordingly, the score network can be parameterized in terms of a denoiser $\mathcal{D}_\theta(\cdot)$ as:

$$\mathbf{s}_\theta(\mathbf{x}, t, \mathbf{c}) = \frac{\mathcal{D}_\theta(\mathbf{x}_t; s^2(t)\sigma(t), \mathbf{c}) - \mathbf{x}_t}{s^2(t)\sigma^2(t)}. \quad (142)$$

Without loss of generality, we set $s(t) = 1$, under which the DSM objective in (3) becomes:

$$\begin{aligned} & \frac{1}{2} \int_0^T \mathbb{E}_{\mathbf{c}, p(\mathbf{x}|\mathbf{c}), p_{0t}(\mathbf{x}_t|\mathbf{x})} [w(t) \|\frac{\mathcal{D}_\theta(\mathbf{x}_t; \sigma(t), \mathbf{c}) - \mathbf{x}}{\sigma^2(t)}\|_2^2] dt \\ &= \tilde{w}(t) \mathbb{E}_{\mathbf{c}, t \sim \mathcal{U}[0, T], p(\mathbf{x}|\mathbf{c}), p_{0t}(\mathbf{x}_t|\mathbf{x})} [\|\mathcal{D}_\theta(\mathbf{x}_t; \sigma(t), \mathbf{c}) - \mathbf{x}\|_2^2], \end{aligned} \quad (143)$$

where $\tilde{w}(t) = \frac{Tw(t)}{2\sigma^4(t)}$. This shows that DSM is equivalent to training the MMSE denoiser $\mathcal{D}_\theta(\cdot; \sigma(t), \mathbf{c})$ for data from class \mathbf{c} corrupted by additive Gaussian noise with standard deviation $\sigma(t)$.

In this setting, the score function is related to the MMSE denoiser via Tweedie’s formula (Miyasawa et al., 1961):

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x} | \mathbf{c}) = \frac{\mathcal{D}(\mathbf{x}; \sigma(t), \mathbf{c}) - \mathbf{x}}{\sigma^2(t)}, \quad (144)$$

where $\mathcal{D}(\mathbf{x}; \sigma(t), \mathbf{c})$ denotes the MMSE denoiser and $p_t(\mathbf{x} | \mathbf{c}) = \int p_{0t}(\mathbf{x} | \mathbf{x}_0) p_{\text{data}}(\mathbf{x}_0) d\mathbf{x}_0$ is the marginal distribution at time t . We are now ready to present the practical ELBO-approximated objectives of MCLR, CC-DPO, and CCA for diffusion models.

ELBO-Approximated MCLR for Diffusion Models. Let $\mathbf{v}_2 = \nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})$ and $\mathbf{v}_3 = \nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \tilde{\mathbf{c}})$. Substituting ELBO (5) into (22), the MCLR regularized DSM becomes:

$$\min_{\boldsymbol{\theta}} \mathcal{J}_{\text{DSM}}(\boldsymbol{\theta}; g^2(\cdot)) + \eta \mathbb{E}_{\substack{\mathbf{c}, \tilde{\mathbf{c}}, t \sim \mathcal{U}[0, T] \\ p(\mathbf{x}|\mathbf{c}), p_{0t}(\mathbf{x}_t|\mathbf{x})}} [g^2(t) (\|\mathbf{v}_2\|_2^2 - \|\mathbf{v}_3\|_2^2)]. \quad (145)$$

Using the denoiser parameterization of the score network, together with a customized training-time noise sampling distribution $p(t)$ and adaptive weighting $w(t)$, the MCLR objective becomes:

$$\mathbb{E}_{\substack{\mathbf{c}, \tilde{\mathbf{c}}, t \sim p(t) \\ p(\mathbf{x}|\mathbf{c}), p_{0t}(\mathbf{x}_t|\mathbf{x})}} [w(t) (\|\mathbf{x} - \mathcal{D}_\theta(\mathbf{x}_t; \sigma(t), \mathbf{c})\|_2^2 - \|\mathbf{x} - \mathcal{D}_\theta(\mathbf{x}_t; \sigma(t), \tilde{\mathbf{c}})\|_2^2)], \quad (146)$$

which can be approximated with Monte Carlo sampling during training. Rather than parameterizing the noise level through the time variable $t \sim p(t)$, one may equivalently define a training noise distribution directly over σ , denoted by $p(\sigma)$ and a corresponding noise adaptive weighting $w(\sigma)$. Under this formulation, (146) can be rewritten as:

$$\mathbb{E}_{\substack{\mathbf{c}, \tilde{\mathbf{c}}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \sigma \sim p(\sigma), p(\mathbf{x}|\mathbf{c})}} [w(\sigma) (\|\mathbf{x} - \mathcal{D}_\theta(\mathbf{x} + \sigma\boldsymbol{\epsilon}; \sigma, \mathbf{c})\|_2^2 - \|\mathbf{x} - \mathcal{D}_\theta(\mathbf{x} + \sigma\boldsymbol{\epsilon}; \sigma, \tilde{\mathbf{c}})\|_2^2)]. \quad (147)$$

Let $\mathbf{v}_4 = \nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}) - \mathbf{s}_\theta(\mathbf{x}_t, t, \mathbf{c})$ and $\mathbf{v}_5 = \nabla_{\mathbf{y}_t} \log p_{0t}(\mathbf{y}_t | \mathbf{y}) - \mathbf{s}_\theta(\mathbf{y}_t, t, \mathbf{c})$. Similarly, substituting ELBO into (12), we get the following equivalent optimization problem:

$$\min_{\boldsymbol{\theta}} \mathcal{J}_{\text{DSM}}(\boldsymbol{\theta}; g^2(\cdot)) + \eta \mathbb{E}_{\substack{\mathbf{c}, \tilde{\mathbf{c}}, t \sim \mathcal{U}[0, T] \\ p(\mathbf{x}|\mathbf{c}), p(\mathbf{y}|\tilde{\mathbf{c}}) \\ p_{0t}(\mathbf{x}_t|\mathbf{x}), p_{0t}(\mathbf{y}_t|\mathbf{y})}} [g^2(t) (\|\mathbf{v}_4\|_2^2 - \|\mathbf{v}_5\|_2^2)], \quad (148)$$

the denoiser form of which is presented in (17).

ELBO-Approximated CC-DPO for Diffusion Models. The CC-DPO algorithm requires access to the log-likelihood ratio for individual data point. Similar to MCLR, we approximate it with the ELBO:

$$\log \frac{p_{\theta}(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} \approx \frac{T}{2} \mathbb{E}_{t \sim \mathcal{U}[0, T], p_{0t}(\mathbf{x}_t | \mathbf{x})} \left[g^2(t) \left(-\|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x}_t, t, \mathbf{c})\|_2^2 + \right. \right. \quad (149)$$

$$\left. \|\nabla_{\mathbf{x}_t} \log p_{0t}(\mathbf{x}_t | \mathbf{x}) - \mathbf{s}_{\text{ref}}(\mathbf{x}_t, t, \mathbf{c})\|_2^2 \right]. \quad (150)$$

Using the denoiser parameterization of the score network, together with a customized training noise distribution and adaptive weighting, the log-likelihood ratio takes the following form:

$$\log \frac{p_{\theta}(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} \approx \mathbb{E}_{\sigma \sim p_{\text{train}}(\sigma), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[w(\sigma) \left(-\|\mathbf{x} - \mathcal{D}_{\theta}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c})\|_2^2 + \|\mathbf{x} - \mathcal{D}_{\text{ref}}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c})\|_2^2 \right) \right]. \quad (151)$$

Define:

$$\Delta(\mathbf{x}, \sigma, \epsilon, \mathbf{c}) := \|\mathbf{x} - \mathcal{D}_{\theta}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c})\|_2^2 - \|\mathbf{x} - \mathcal{D}_{\text{ref}}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c})\|_2^2. \quad (152)$$

By substituting (151) into DPO objective (18), we get the following optimization problem:

$$\min_{\theta} -\mathbb{E}_{\mathbf{c}, \bar{\mathbf{c}}, \mathbf{x}_w \sim p(\mathbf{x} | \mathbf{c}), \mathbf{x}_l \sim p(\mathbf{x} | \bar{\mathbf{c}})} \left[\log \text{Sigmoid}(\beta \mathbb{E}_{\sigma, \epsilon} [w(\sigma) (-\Delta(\mathbf{x}_w, \sigma, \epsilon, \mathbf{c}) + \Delta(\mathbf{x}_l, \sigma, \epsilon, \mathbf{c}))]) \right]. \quad (153)$$

Since the log-Sigmoid function is concave, applying Jensen’s inequality by moving the expectation outside yields the following upper bound of objective (153):

$$\min_{\theta} -\mathbb{E}_{\mathbf{c}, \bar{\mathbf{c}}, \mathbf{x}_w, \mathbf{x}_l, \sigma, \epsilon} \left[\log \text{Sigmoid}(\beta w(\sigma) (-\Delta(\mathbf{x}_w, \sigma, \epsilon, \mathbf{c}) + \Delta(\mathbf{x}_l, \sigma, \epsilon, \mathbf{c}))) \right], \quad (154)$$

which serves as our final training objective.

CCA for Diffusion Models. The Conditional Contrastive Alignment (CCA) (Chen et al., 2025b) objective takes the following form:

$$\max_{\theta} \mathbb{E}_{\mathbf{c}, p(\mathbf{x} | \mathbf{c})} \log \text{Sigmoid} \left(\beta \log \frac{p_{\theta}(\mathbf{x} | \mathbf{c})}{p_{\text{ref}}(\mathbf{x} | \mathbf{c})} \right) + \lambda \mathbb{E}_{\mathbf{c}, \bar{\mathbf{c}}, p(\mathbf{x} | \bar{\mathbf{c}})} \log \text{Sigmoid} \left(-\beta \log \frac{p_{\theta}(\mathbf{x} | \mathbf{c})}{p_{\text{ref}}(\mathbf{x} | \mathbf{c})} \right). \quad (155)$$

To adapt this optimization problem to diffusion models, we again approximate the log-likelihood with ELBO, resulting in the following optimization problem:

$$\max_{\theta} \mathbb{E}_{\mathbf{c}, \bar{\mathbf{c}}, \mathbf{x}_w \sim p(\mathbf{x} | \mathbf{c}), \mathbf{x}_l \sim p(\mathbf{x} | \bar{\mathbf{c}})} \left[\log \text{Sigmoid} \left(\beta \mathbb{E}_{\sigma, \epsilon} \left[-w(\sigma) \Delta(\mathbf{x}_w, \sigma, \epsilon, \mathbf{c}) \right] \right) \right. \\ \left. + \lambda \log \text{Sigmoid} \left(\beta \mathbb{E}_{\sigma, \epsilon} \left[w(\sigma) \Delta(\mathbf{x}_l, \sigma, \epsilon, \mathbf{c}) \right] \right) \right]. \quad (156)$$

Since the log-sigmoid function is concave, applying Jensen’s inequality by moving the expectation outside yields the following lower bound of objective (156):

$$\max_{\theta} \mathbb{E}_{\mathbf{c}, \mathbf{x}_w, \mathbf{x}_l, \sigma, \epsilon} \left[\log \text{Sigmoid} \left(-\beta w(\sigma) \Delta(\mathbf{x}_w, \sigma, \epsilon, \mathbf{c}) \right) + \lambda \log \text{Sigmoid} \left(\beta w(\sigma) \Delta(\mathbf{x}_l, \sigma, \epsilon, \mathbf{c}) \right) \right], \quad (157)$$

which serves as our final training objective.

Remarks on the ELBO Approximation. Despite standard practice in the literature, when replacing log-likelihood with ELBO, the resulting objectives do not necessarily correspond to the original likelihood formulations, since the ELBO does not enforce the regularity conditions required for the parameterized score function to define a valid score field. Consequently, ELBO-approximated MCLR should be viewed as an **approximate** likelihood-ratio maximization procedure for diffusion models. This issue does not arise for autoregressive models, where likelihoods are available exactly.

E.2. Building Training Data from a Minibatch

The equivalence of (21) and (22) offer us two choices for constructing the contrastive pairs. The first is constructing tuples of $(\mathbf{x}, \mathbf{y}, \mathbf{c})$, where $\mathbf{x} \sim p(\mathbf{x}|\mathbf{c})$ and \mathbf{y} sampled from other random classes. The other is constructing contrastive tuples $(\mathbf{x}, \mathbf{c}, \tilde{\mathbf{c}})$, where $\mathbf{x} \sim p(\mathbf{x}|\mathbf{c})$ and $\tilde{\mathbf{c}}$ denotes a randomly chosen alternative class. In our implementation we choose the second method, but our preliminary results demonstrate both achieve similar performance.

Similarly, CC-DPO and CCA require constructing preference-style tuples $(\mathbf{x}_w, \mathbf{x}_l, \mathbf{c})$. In practice, we build these tuples directly from each training minibatch $\{(\mathbf{x}_i, \mathbf{c}_i)\}_{i=1}^N$. We consider two strategies, described below.

Approach 1: Building N pairs. The simplest strategy constructs one contrastive (or preference) tuple per sample in the minibatch. For MCLR, given a sample (\mathbf{x}, \mathbf{c}) , we randomly select another label $\tilde{\mathbf{c}}$ from the same minibatch, forming a tuple $(\mathbf{x}, \mathbf{c}, \tilde{\mathbf{c}})$. Repeating this process independently for each \mathbf{x} (with replacement) yields N tuples in total. The same strategy applies to CC-DPO and CCA, where we construct N tuples $(\mathbf{x}_w, \mathbf{x}_l, \mathbf{c})$ from the minibatch.

Approach 2: Building NK pairs. Alternatively, we can construct multiple contrastive tuples per sample. For MCLR, given each (\mathbf{x}, \mathbf{c}) , we randomly select K alternative labels $\{\tilde{\mathbf{c}}_k\}_{k=1}^K$ from the same minibatch, yielding K tuples $(\mathbf{x}, \mathbf{c}, \tilde{\mathbf{c}}_k)$. Repeating this procedure for all N samples results in NK tuples in total. An analogous strategy is applied to CC-DPO and CCA, producing NK tuples $(\mathbf{x}_w, \mathbf{x}_l, \mathbf{c})$ per minibatch.

Compared to Approach 1, Approach 2 increases the number of training tuples constructed from each minibatch, allowing better exploitation of inter-class contrastive information. In our experiments, we adopt Approach 2 for EDM-based diffusion models, where it consistently yields improved quantitative performance. For VAR and SiT models, we use Approach 1 due to its lower computational overhead and comparable empirical performance.

E.3. Overall Algorithm

We are now ready to present the overall algorithm. The algorithm for MCLR is given in Algorithm 1, while the algorithms for CC-DPO and CCA are presented in Algorithm 2. In Algorithm 1, we compute the MCLR regularization loss without including the standard DSM term. Empirically, we find that optimizing the MCLR term alone yields a better best-case $\text{FD}_{\text{DINOv2}}$ score. One possible explanation is that removing the DSM term makes the contrastive objective the dominant training signal, allowing the model to explore a larger solution space beyond the neighborhood constrained by DSM regularization. Accordingly, the main results report the best $\text{FD}_{\text{DINOv2}}$ scores obtained without the DSM term. Nevertheless, including the DSM term can improve training stability and lead to a more favorable Precision–Recall trade-off Appendix G.1.

E.4. Hyperparameters

For each method, we carefully tune hyperparameters using grid search to ensure the best performance. The full set of hyperparameter configurations will be made publicly available upon publication.

E.5. Computational Efficiency

MCLR fine-tuning incurs additional cost, but it is modest relative to pretraining. In our runs on $4 \times \text{A40 40GB GPUs}$, for EDM2-S (64×64), MCLR fine-tuning took 5h23m to reach best $\text{FD}_{\text{DINOv2}}$ score 56 with Approach 1 (construct N contrastive pairs per batch) and 18h40m to reach best $\text{FD}_{\text{DINOv2}}$ 52 with Approach 2 (construct NK contrastive pairs per batch), starting from a base model with $\text{FD}_{\text{DINOv2}}$ score 95. In contrast, the pretraining stage requires training with 32 A100 GPUs for weeks (Karras et al., 2024b). A comprehensive summary of computational cost can be found within our codebase upon publication.

F. Additional Experimental Results

In this section, we present additional experimental results that complement and extend those in the main text. In Section 6, we focus on validating the following three claims:

1. MCLR induces progressive class separation, leading to a fidelity–diversity trade-off analogous to that produced by increasing the guidance strength in classifier-free guidance (CFG).

Algorithm 1 MCLR

Require: Pre-trained base model \mathcal{D}_{θ_0} ; noise-adaptive weight function $w(\cdot)$; training noise schedule $p(\sigma)$; learning rate α ; training dataset representing $p(\mathbf{x}, \mathbf{c})$; $K \geq 1$ (number of mismatched class labels per sample in Approach 2)

- 1: Initialize $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}_0$
- 2: **for** each training iteration **do**
- 3: Sample a minibatch $\{(\mathbf{x}_i, \mathbf{c}_i)\}_{i=1}^N$ with $(\mathbf{x}_i, \mathbf{c}_i) \sim p(\mathbf{x}, \mathbf{c})$.
- 4: Sample $\{\sigma_i\}_{i=1}^N$ with $\sigma_i \sim p(\sigma)$
- 5: Sample $\{\epsilon_i\}_{i=1}^N$ with $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6: Construct contrastive tuples from the minibatch:
- 7: **if** Approach 1 **then**
- 8: Set $M \leftarrow N$
- 9: For each i , sample $\tilde{\mathbf{c}}_i$ from minibatch labels with $\tilde{\mathbf{c}}_i \neq \mathbf{c}_i$
- 10: Define tuples $\{(\mathbf{x}_i, \mathbf{c}_i, \tilde{\mathbf{c}}_i, \sigma_i, \epsilon_i)\}_{i=1}^M$
- 11: **if** Approach 2 **then**
- 12: Set $M \leftarrow NK$
- 13: For each i and each $k = 1, \dots, K$, sample $\tilde{\mathbf{c}}_{i,k}$ from minibatch with $\tilde{\mathbf{c}}_{i,k} \neq \mathbf{c}_i$
- 14: Define tuples $\{(\mathbf{x}_{i,k}, \mathbf{c}_{i,k}, \tilde{\mathbf{c}}_{i,k}, \sigma_{i,k}, \epsilon_{i,k})\}$ where
 $(\mathbf{x}_{i,k}, \mathbf{c}_{i,k}) = (\mathbf{x}_i, \mathbf{c}_i)$ and $(\sigma_{i,k}, \epsilon_{i,k}) = (\sigma_i, \epsilon_i)$
- 15: $\mathcal{L} \leftarrow 0$
- 16: **for** each tuple $(\mathbf{x}, \mathbf{c}, \tilde{\mathbf{c}}, \sigma, \epsilon)$ **do**
- 17: $\mathcal{L} += w(\sigma) \left(\|\mathcal{D}_{\boldsymbol{\theta}}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c}) - \mathbf{x}\|_2^2 - \|\mathcal{D}_{\boldsymbol{\theta}}(\mathbf{x} + \sigma\epsilon; \sigma, \tilde{\mathbf{c}}) - \mathbf{x}\|_2^2 \right)$
- 18: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{L}$
- 19: **return** $\mathcal{D}_{\boldsymbol{\theta}}$

2. MCLR matches or outperforms training-time baselines, including CC-DPO and CCA.
3. MCLR achieves performance comparable to CFG.

Due to space constraints, the main text reports quantitative results only for ImageNet-512 \times 512 with EDM2-L and ImageNet-256 \times 256 with VAR-d24. Here, we provide the complete set of experimental results on ImageNet-64 \times 64 with EDM2-S, ImageNet-512 \times 512 with EDM2-L, and ImageNet-256 \times 256 with SiT (with REPA) and VAR-d24.

F.1. Progressive Class Separation and the Fidelity-Diversity Trade-off.

Training with MCLR induces notable *progressive class separation*, which gives rise to a fidelity–diversity trade-off similar that observed when increasing the guidance scale in CFG. This behavior is quantitatively supported by the following observations.

(i) Increased Precision and Inception Score. As shown in Figure 5(b,d,g) for ImageNet-64, Figure 7(c,f,i) for ImageNet-256 and Figure 6(b,d,g) for ImageNet-512, continued training with MCLR leads to a progressive increase in Inception Score, indicating increasingly class-discriminative generations. This is accompanied by a corresponding increase in Precision, reflecting improved image fidelity.

(ii) Decreased Recall. Conversely, as shown in Figure 5(e,h), Figure 6(g,j), and Figure 7(e,h), excessive training reduces within-class diversity, which manifests as a decrease in Recall (see also Figure 3(e)).

Taken together, these effects result in FD–IS and Precision–Recall trade-offs that closely mirror those induced by CFG, as illustrated in Figure 5(c,f,i), Figure 6(c,f,i), and Figure 7(d,e,h,k). Qualitatively, as shown in Figures 1 and 8 to 19, continued MCLR training leads to gradually enhanced distinct, class-specific visual characteristics in generated samples.

F.2. MCLR Outperforms Training-time Baselines.

Diffusion Models. As shown in Figures 5 and 6, for EDM models trained on ImageNet, MCLR achieves substantially better best-case $\text{FD}_{\text{DINOv2}}$ scores than CCA, CC-DPO. Moreover, MCLR traverses significantly wider FD–IS and Precision–

Algorithm 2 CC-DPO and CCA

Require: Pre-trained base model \mathcal{D}_{θ_0} ; noise-adaptive weight function $w(\cdot)$; training noise schedule $p(\sigma)$; learning rate α ; KL regularization strength β ; CCA hyperparameter λ ; training dataset representing $p(\mathbf{x}, \mathbf{c})$; $K \geq 1$ (number of mismatched class labels per sample in Approach 2)

- 1: Initialize $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta}_0$, $\mathcal{D}_{\text{ref}} \leftarrow \mathcal{D}_{\theta_0}$
- 2: **for** each training iteration **do**
- 3: Sample a minibatch $\{(\mathbf{x}_i, \mathbf{c}_i)\}_{i=1}^N$ with $(\mathbf{x}_i, \mathbf{c}_i) \sim p(\mathbf{x}, \mathbf{c})$.
- 4: Sample $\{\sigma_i\}_{i=1}^N$ with $\sigma_i \sim p(\sigma)$
- 5: Sample $\{\boldsymbol{\epsilon}_i\}_{i=1}^N$ with $\boldsymbol{\epsilon}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6: Construct contrastive tuples from the minibatch:
- 7: **if** Approach 1 **then**
- 8: Set $M \leftarrow N$
- 9: For each i , sample $\tilde{\mathbf{x}}_i$ from minibatch samples with label $\tilde{\mathbf{c}}_i \neq \mathbf{c}_i$
- 10: Define tuples $\{(\mathbf{x}_i, \tilde{\mathbf{x}}_i, \mathbf{c}_i, \sigma_i, \boldsymbol{\epsilon}_i)\}_{i=1}^M$
- 11: **if** Approach 2 **then**
- 12: Set $M \leftarrow NK$
- 13: For each i and each $k = 1, \dots, K$, sample $\tilde{\mathbf{x}}_{i,k}$ from minibatch samples with label $\tilde{\mathbf{c}}_{i,k} \neq \mathbf{c}_i$
- 14: Define tuples $\{(\mathbf{x}_{i,k}, \tilde{\mathbf{x}}_{i,k}, \mathbf{c}_{i,k}, \sigma_{i,k}, \boldsymbol{\epsilon}_{i,k})\}$ where
 $(\mathbf{x}_{i,k}, \mathbf{c}_{i,k}) = (\mathbf{x}_i, \mathbf{c}_i)$ and $(\sigma_{i,k}, \boldsymbol{\epsilon}_{i,k}) = (\sigma_i, \boldsymbol{\epsilon}_i)$
- 15: $\mathcal{L} \leftarrow 0$
- 16: **for** each tuple $(\mathbf{x}, \tilde{\mathbf{x}}, \mathbf{c}, \sigma, \boldsymbol{\epsilon})$ **do**
- 17: $\Delta(\mathbf{x}_w, \sigma, \boldsymbol{\epsilon}, \mathbf{c}) := \|\mathbf{x} - \mathcal{D}_{\boldsymbol{\theta}}(\mathbf{x} + \sigma\boldsymbol{\epsilon}; \sigma, \mathbf{c})\|_2^2 - \|\mathbf{x} - \mathcal{D}_{\text{ref}}(\mathbf{x} + \sigma\boldsymbol{\epsilon}; \sigma, \mathbf{c})\|_2^2$
- 18: $\Delta(\mathbf{x}_l, \sigma, \boldsymbol{\epsilon}, \mathbf{c}) := \|\tilde{\mathbf{x}} - \mathcal{D}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}} + \sigma\boldsymbol{\epsilon}; \sigma, \mathbf{c})\|_2^2 - \|\tilde{\mathbf{x}} - \mathcal{D}_{\text{ref}}(\tilde{\mathbf{x}} + \sigma\boldsymbol{\epsilon}; \sigma, \mathbf{c})\|_2^2$
- 19: **if** CC-DPO **then**
- 20: $\mathcal{L} += -\log \text{Sigmoid}(\beta w(\sigma)(-\Delta(\mathbf{x}_w, \sigma, \boldsymbol{\epsilon}, \mathbf{c}) + \Delta(\mathbf{x}_l, \sigma, \boldsymbol{\epsilon}, \mathbf{c})))$
- 21: **if** CCA **then**
- 22: $\mathcal{L} -= \log \text{Sigmoid}(-\beta w(\sigma)\Delta(\mathbf{x}_w, \sigma, \boldsymbol{\epsilon}, \mathbf{c})) + \lambda \log \text{Sigmoid}(\beta w(\sigma)\Delta(\mathbf{x}_l, \sigma, \boldsymbol{\epsilon}, \mathbf{c}))$
- 23: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{L}$
- 24: **return** $\mathcal{D}_{\boldsymbol{\theta}}$

Inception trade-offs region with a faster training speed. In contrast, CCA and CC-DPO often converge early, exhibit zigzag training trajectories, and become trapped in suboptimal local minima, as evidenced by slowly improving or stalled learning curves.

Autoregressive Models. As shown in Figure 7, MCLR, CC-DPO, and CCA yield comparable performance improvement on VAR-d24. Specifically, we see MCLR consistently achieves higher precision (see Figure 7(f,h,i,k)) and inception score (see Figure 7(c)) in later stages of training.

E.3. MCLR Achieves Comparable Performance as CFG.

Diffusion Models. For EDM models trained on ImageNet, CFG generally exhibits a better FD-IS trade-off and achieves a lower best-case $\text{FD}_{\text{DINOv2}}$ than MCLR. Nevertheless, this gap is moderate for EDM2-L model, where CFG attains a best-case $\text{FD}_{\text{DINOv2}}$ of 39.86 compared to 42.50 for MCLR as shown in Figure 6(c). Moreover, for EDM2-L model, when evaluated using Precision-Recall, MCLR demonstrates competitive and in some regimes superior behavior. As shown in Figure 6(f,g), MCLR matches CFG in the high-recall regime (early training stages) and achieves a substantially higher best-case precision in the high-precision regime (later training stages), where CFG begins to generate images with oversaturated colors.

Applying CFG on top of an MCLR-fine-tuned model further alleviates the performance gap between MCLR and CFG, yielding higher best-case Inception Score and Precision on both EDM2-S and EDM2-L models, as shown in Figure 5(c,f,i) and Figure 6(c,f,i).

Qualitatively, MCLR and CFG often produce similar effects, both significantly enhancing class-specific structures in the generated images, as shown in Figures 1 and 8 to 19. This similarity is expected, as both methods improve conditional

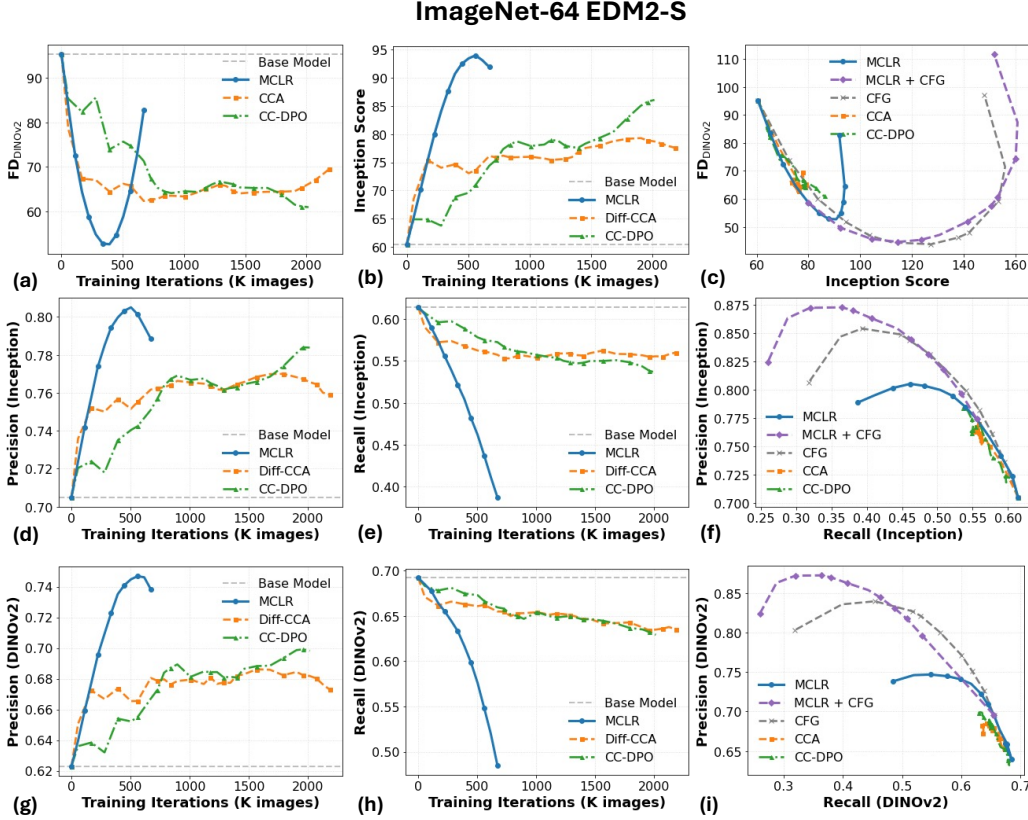


Figure 5. **Quantitative Results for EDM2-S trained on ImageNet-64×64.** (a), (b), (d), (e), (g), (h) show the evolution of FD, Inception Score, Precision (calculated with Inception features), and Recall (calculated with Inception features), Precision (calculated with DINOv2 features), and Recall (calculated with DINOv2 features), respectively, as functions of training iterations. (c) shows the FD–IS trade-offs, while (f), (i) depict the Precision–Recall trade-offs calculated with Inception and DINOv2 features respectively. We evaluate classifier-free guidance (CFG) scales $\gamma \in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.7, 0.9, 1, 1.5, 2.0, 3.0\}$.

modeling through inter-class contrastive signals. The key distinction is that MCLR internalizes this mechanism during training, whereas CFG applies it at inference time.

Autoregressive Models. For VAR-d24 model, MCLR achieves a similar FD–IS trade-off to CFG in terms of both FD_{DINOv2} and FID, with CFG exhibiting a slightly better best-case FID. Consistent with diffusion models, MCLR outperforms CFG in the Precision–Recall trade-off and achieves a higher best-case precision.

G. Ablation Study

G.1. Effect of DSM Regularization

In this section we demonstrate the effect of combining (i) DSM and (ii) a KL-style loss with MCLR. Concretely, we fine-tune SiT-XL/2 (with REPA) base model $\mathcal{D}_{\text{ref}}(\cdot)$ on ImageNet-256 × 256 using the following objectives respectively with varying $\beta = 1$:

$$\begin{aligned}
 & \underbrace{\beta \mathbb{E}_{\substack{\mathbf{c}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \sigma \sim p(\sigma), p(\mathbf{x}|\mathbf{c})}} \left[w(\sigma) (\|\mathbf{x} - \mathcal{D}_{\theta}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c})\|_2^2) \right]}_{\text{DSM Objective}} \\
 & + \underbrace{\mathbb{E}_{\substack{\mathbf{c}, \tilde{\mathbf{c}}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \sigma \sim p(\sigma), p(\mathbf{x}|\mathbf{c})}} \left[w(\sigma) (\|\mathbf{x} - \mathcal{D}_{\theta}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c})\|_2^2 - \|\mathbf{x} - \mathcal{D}_{\theta}(\mathbf{x} + \sigma\epsilon; \sigma, \tilde{\mathbf{c}})\|_2^2) \right]}_{\text{MCLR Regularization}}, \tag{158}
 \end{aligned}$$

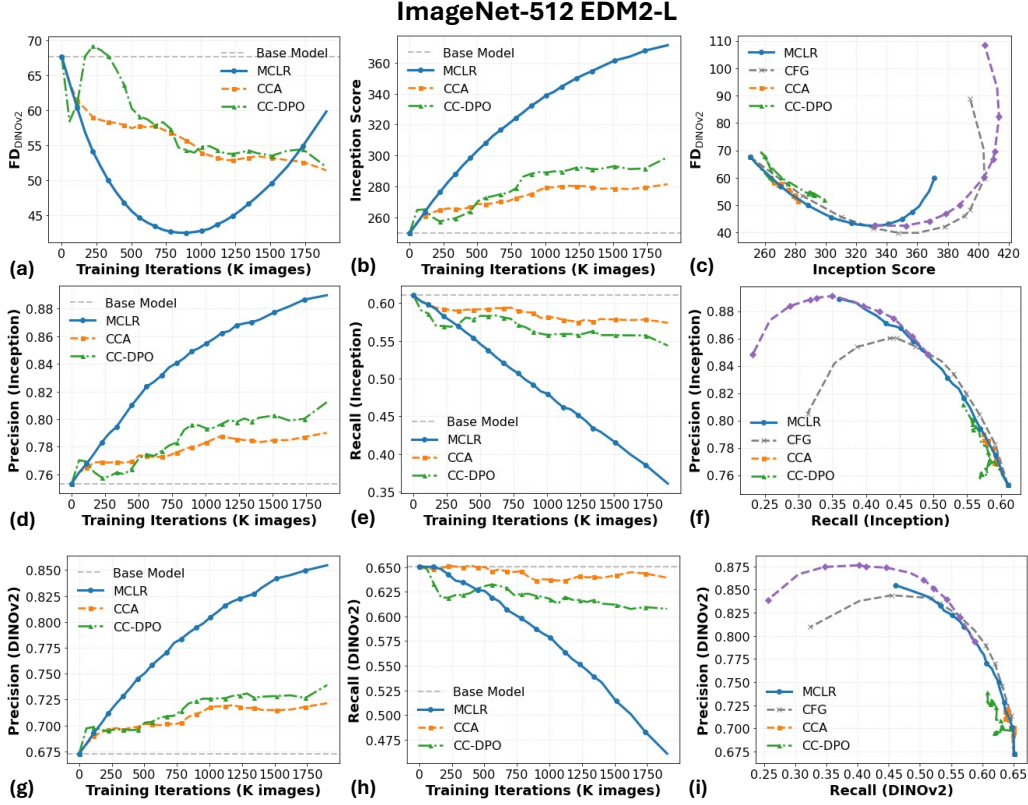


Figure 6. **Quantitative Results for EDM2-L trained on ImageNet-512 \times 512.** (a), (b), (d), (e), (g), (h) show the evolution of FD, Inception Score, Precision (calculated with Inception features), and Recall (calculated with Inception features), Precision (calculated with DINOv2 features), and Recall (calculated with DINOv2 features), respectively, as functions of training iterations. (c) shows the FD–IS trade-offs, while (f), (i) depict the Precision–Recall trade-offs calculated with Inception and DINOv2 features respectively. We evaluate classifier-free guidance (CFG) scales $\gamma \in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.7, 0.9, 1, 1.5, 2.0, 3.0\}$.

$$\begin{aligned}
 & \beta \underbrace{\mathbb{E}_{\substack{\mathbf{c}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \sigma \sim p(\sigma), p(\mathbf{x}|\mathbf{c})}} \left[w(\sigma) (\| \mathcal{D}_{\text{ref}}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c}) - \mathcal{D}_{\theta}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c}) \|_2^2) \right]}_{\text{KL Loss}} \\
 & + \underbrace{\mathbb{E}_{\substack{\mathbf{c}, \tilde{\mathbf{c}}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \sigma \sim p(\sigma), p(\mathbf{x}|\mathbf{c})}} \left[w(\sigma) (\| \mathbf{x} - \mathcal{D}_{\theta}(\mathbf{x} + \sigma\epsilon; \sigma, \mathbf{c}) \|_2^2 - \| \mathbf{x} - \mathcal{D}_{\theta}(\mathbf{x} + \sigma\epsilon; \sigma, \tilde{\mathbf{c}}) \|_2^2) \right]}_{\text{MCLR Regularization}}, \tag{159}
 \end{aligned}$$

where (159) is an approximation of the KL-regularized MCLR in (13). While the KL loss is, in principle, defined over samples generated by the reference model, we approximate it using real data for simplicity.

The results of fine-tuning using (158), (159) and MCLR regularization only are presented in Figure 20. Note that while optimizing MCLR alone can lead to a better best case $\text{FD}_{\text{DINOv2}}$, adding DSM or KL loss can improve the trade-off curves, in which case the performance of MCLR is comparable to CFG.

G.2. Hyperparameter Sensitivity of CC-DPO

The hyperparameter β in the CC-DPO objective (18) controls the strength of the KL regularization between the base model and the fine-tuned model. We study its effect on EDM2-S trained on ImageNet-64 \times 64. As shown in Figure 21, CC-DPO exhibits relatively stable performance across a wide range of β values. In particular, smaller β values achieve marginally better best-case $\text{FD}_{\text{DINOv2}}$ and Inception Score.

G.3. Hyperparameter Sensitivity of CCA

The hyperparameter λ in the CCA objective (157) controls the strength of the contrastive term, with larger λ placing greater emphasis on penalizing the non-preferred sample x_l . We study its effect on EDM2-S trained on ImageNet-64 \times 64, fixing $\beta = 0.001$. As shown in Figure 22, λ has little impact on the best achievable $\text{FD}_{\text{DINOv2}}$; however, smaller λ values lead to more stable and smoother convergence. In contrast, larger λ values can yield slightly higher best-case Inception Scores, but often cause training instability and eventual collapse when training is prolonged.

In practice, jointly tuning β and λ is challenging. Although CCA and CC-DPO share the same theoretical optimal solution, CC-DPO requires tuning only a single hyperparameter β . Empirically, CCA does not outperform CC-DPO under careful tuning as shown in Figure 22. Therefore, we recommend CC-DPO over CCA in practice.

H. Discussion on Related Works

H.1. Visual Generation without Guidance

Our work contributes to a growing line of research (Chen et al., 2025b;a; Tang et al., 2025) that seeks to induce CFG-like effect by modifying the standard training objective, rather than applying CFG at inference time. The most closely related approach is CCA (Chen et al., 2025b), which aims to learn the gamma-powered distribution (19) in autoregressive models using **Noise Contrastive Estimation (NCE)** (Gutmann & Hyvärinen, 2010). While CCA also considers DPO as a baseline and reports inferior performance, our theoretical analysis demonstrates that CC-DPO and CCA admit the same optimal solution, and are therefore fundamentally equivalent at the population level. With a correct implementation, we find that CC-DPO consistently matches or exceeds the empirical performance of CCA while requiring fewer hyperparameters and simpler optimization. We further extend CCA to diffusion models by approximating log-likelihoods via the ELBO and include this variant as a baseline in our experiments. A detailed theoretical and empirical analysis of CCA is provided in Appendix C. Another relevant direction is **Guidance-Free Training (GFT)** (Chen et al., 2025a; Tang et al., 2025), which aims to directly train diffusion models to reproduce CFG-induced score functions. GFT relies on a reparameterized score function in the same functional form as CFG, hence inherits the same theoretical ambiguities associated with CFG itself. In contrast, the proposed MCLR objective provides a clearer mechanistic interpretation of CFG.

Besides these approaches, several works sought to enhance conditional generative modeling through class-wise contrastive mechanisms. (Yan et al., 2024) introduce a contrastive objective to improve the modeling of tail classes; (Lee et al., 2025) employ an InfoNCE-style loss (Oord et al., 2018) to strengthen text-image alignment in diffusion models; (Kadkhodaie et al., 2024) and (Yun et al., 2025) regularize diffusion models by encouraging class separation in feature space. In contrast to MCLR, CC-DPO and CCA, these approaches are primarily heuristic and do not provide a theoretical characterization of the conditional distribution induced by their objectives.

Lastly, beyond class-wise contrast, another complementary strategy for improving generative modeling is to contrast high-quality (“real”) data against low-quality (“synthetic”) samples. Zheng et al. (Zheng et al., 2025) propose Direct Discriminative Optimization (DDO), which employs the same contrastive objective as CCA, but replaces class-conditioned contrast with a real–synthetic discrimination signal. We include DDO as one of our baselines.

H.2. Inference-Time Alignment via Guidance

Recent works (Frans et al., 2025; Jin et al., 2025b; Cheng et al., 2025; Jiang et al., 2026) observe that CFG-style inference-time guidance can produce effects resembling those of training-time alignment methods, including reinforcement learning-based approaches. These studies provide empirical evidence that guidance may implicitly induce alignment-like behavior.

However, existing analyses are largely heuristic and often rely on unrealistic assumptions such as the guided score corresponds to that of the gamma-powered distribution (19), an assumption that has been proven incorrect (Karras et al., 2024a; Bradley & Nakkiran, 2024). As a result, a rigorous theoretical connection between CFG and alignment objectives remains incomplete.

By establishing the equivalence between CFG and a weighted MCLR objective, our work provides a formal mechanistic interpretation of CFG as an inference-time alignment algorithm. To the best of our knowledge, this is the first result that rigorously connects CFG-style guidance with likelihood-ratio-based alignment objectives.

I. Broader Impact

This work aims to improve the training of conditional diffusion models by reducing reliance on inference-time guidance mechanisms such as classifier-free guidance (CFG). By incorporating contrastive signals directly into the training objective, the proposed method enables standard sampling procedures to achieve CFG-like effects, potentially reducing computational overhead and simplifying deployment in downstream applications. These improvements may benefit a wide range of generative modeling tasks, including image synthesis, data augmentation, and representation learning. At the same time, the proposed approach does not introduce new application domains or capabilities beyond those of existing diffusion models, and thus inherits similar ethical considerations and risks, such as the potential misuse of generated content. As with prior generative models, responsible use, appropriate dataset curation, and adherence to existing guidelines remain essential. Overall, this work contributes to a better theoretical and practical understanding of conditional diffusion models, with the goal of improving their efficiency and interpretability

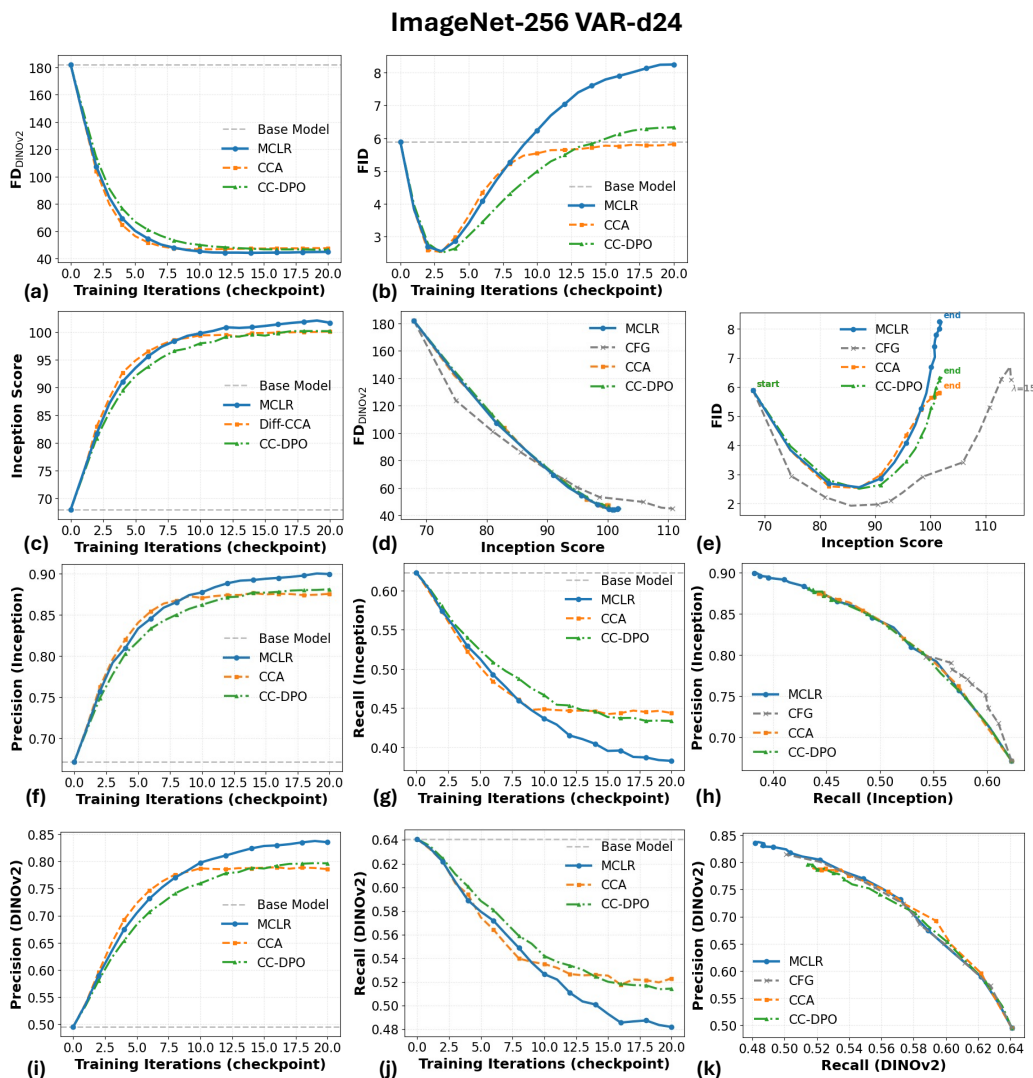


Figure 7. **Quantitative Results for VAR-d24 trained on ImageNet-64 \times 64.** (a), (b), (c), (f), (g), (i), (j) show the evolution of FD_{DINOv2} , FID, Inception Score, Precision (calculated with Inception features), and Recall (calculated with Inception features), Precision (calculated with DINOv2 features), and Recall (calculated with DINOv2 features), respectively, as functions of training iterations. (d), (e) shows the FD-IS trade-offs, while (h), (k) depict the Precision-Recall trade-offs calculated with Inception and DINOv2 features respectively. We evaluate classifier-free guidance (CFG) scales $\gamma \in \{0.5, 0.8, 1.1, 1.5, 1.7, 2.0, 2.5, 3.0, 4.0, 5.0, 7.0, 10.0, 15.0\}$.

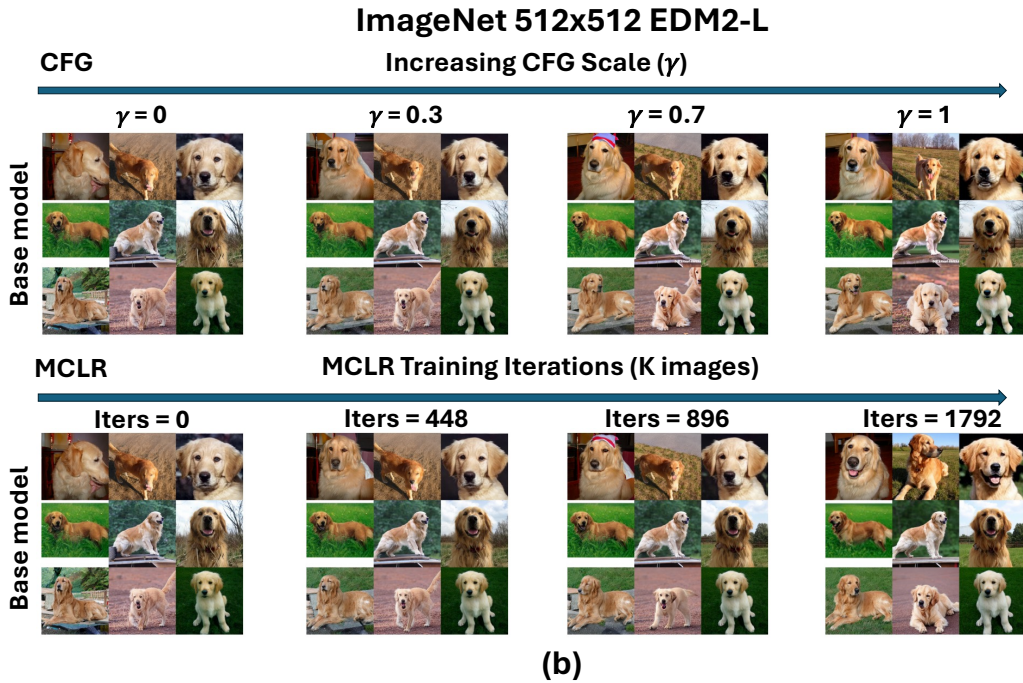
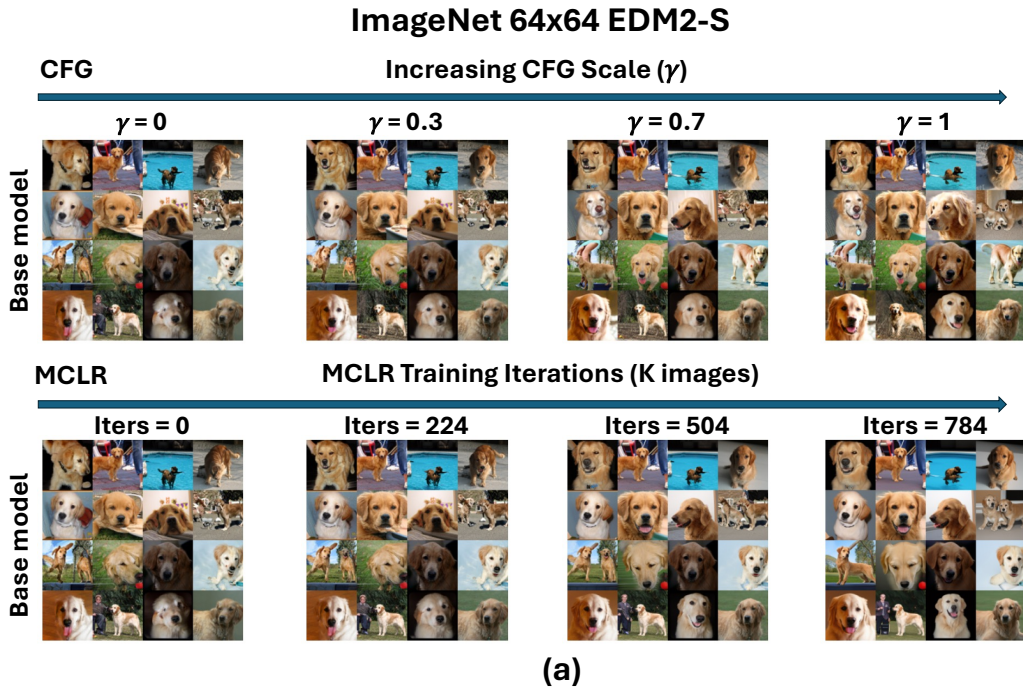


Figure 8. Comparison between CFG and MCLR for Golden Retriever (Class 207). (a,b) demonstrate the progressive evolution of generated samples on ImageNet-64x64 and ImageNet-512x512, respectively. Increasing the CFG scale γ (top rows) and progressive MCLR training (bottom rows) produce similar effects, both enhancing class-specific structures in the generated images.



Figure 9. Comparison between CFG and MCLR for Warthog (Class 343). (a,b) demonstrate the progressive evolution of generated samples on ImageNet-64x64 and ImageNet-512x512, respectively. Increasing the CFG scale γ (top rows) and progressive MCLR training (bottom rows) produce similar effects, both enhancing class-specific structures in the generated images.

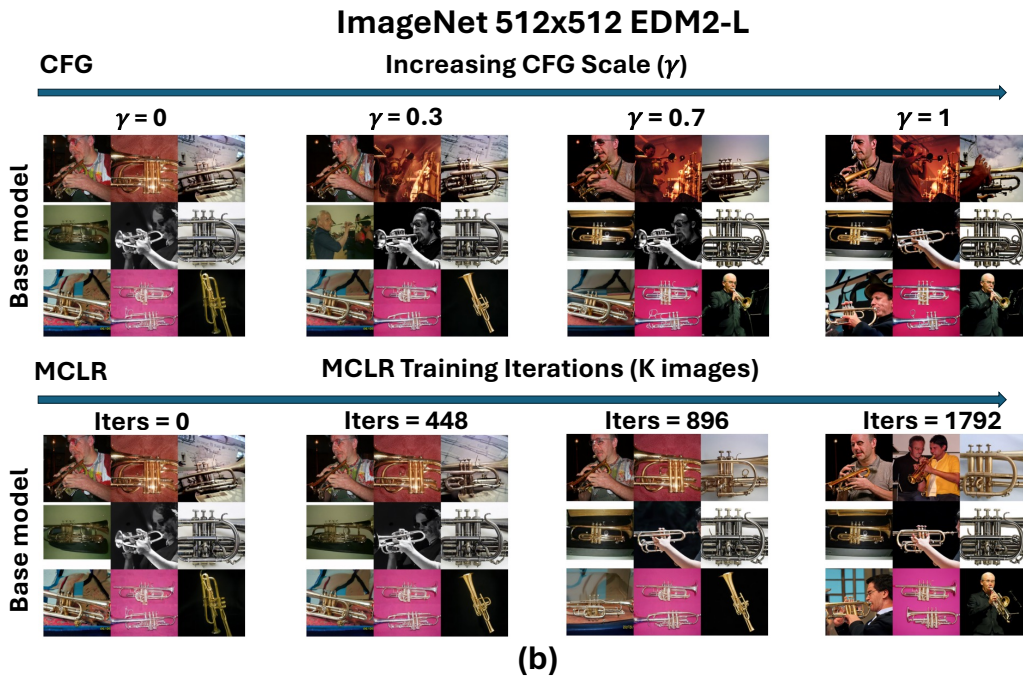
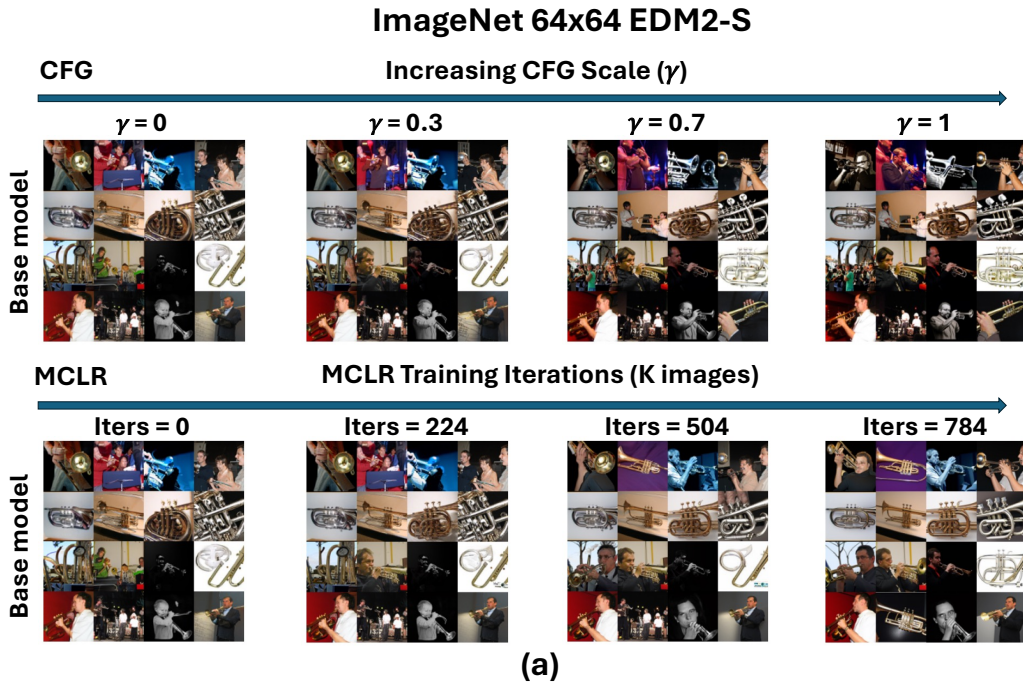


Figure 10. Comparison between CFG and MCLR for Trumpet (Class 513). (a,b) demonstrate the progressive evolution of generated samples on ImageNet-64x64 and ImageNet-512x512, respectively. Increasing the CFG scale γ (top rows) and progressive MCLR training (bottom rows) produce similar effects, both enhancing class-specific structures in the generated images.

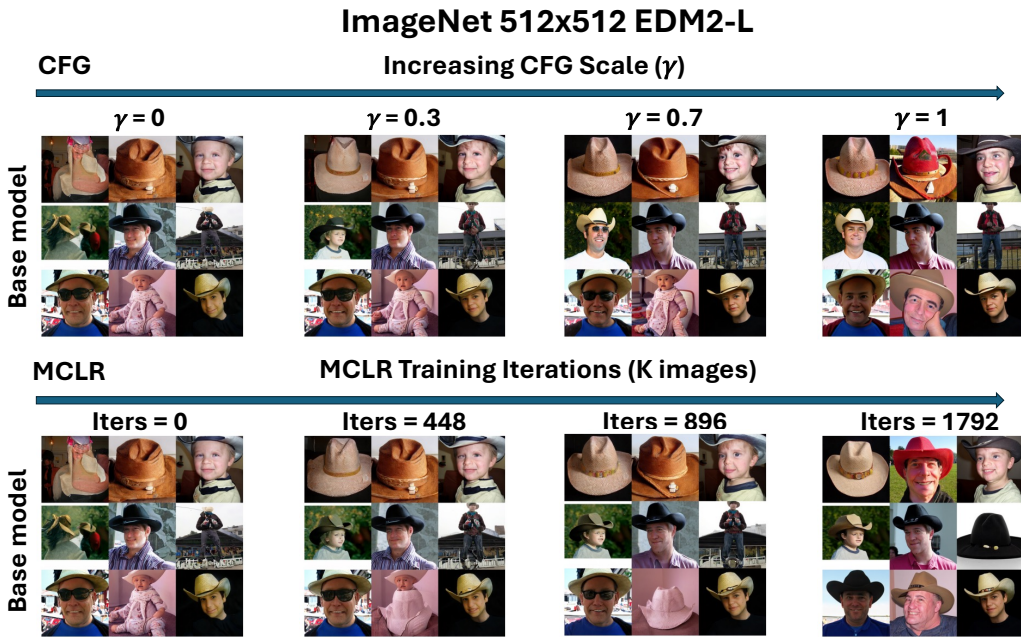
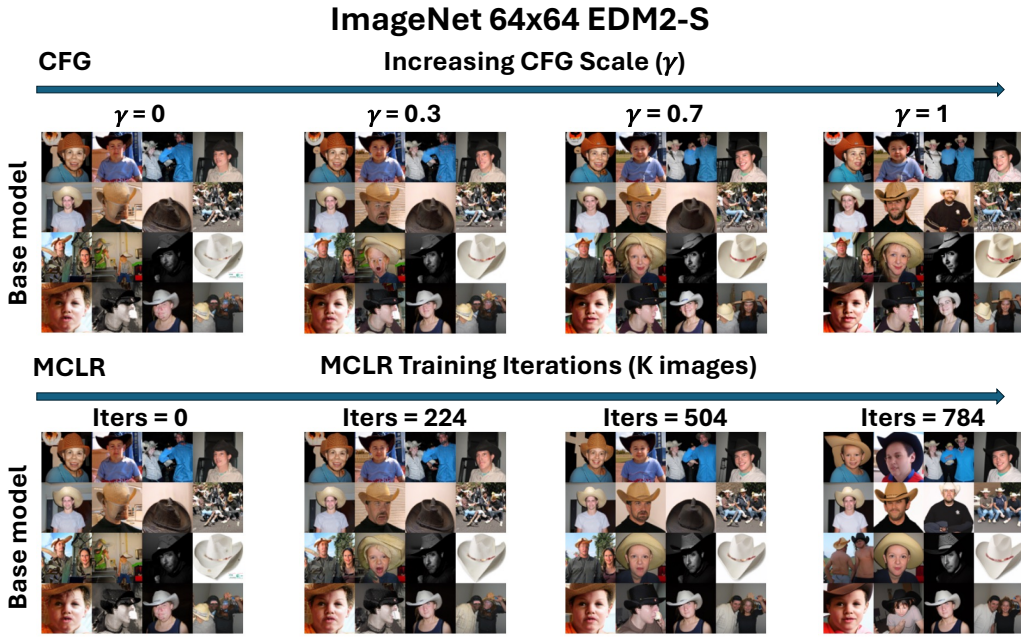


Figure 11. Comparison between CFG and MCLR for Cowboy Hat (Class 515). (a,b) demonstrate the progressive evolution of generated samples on ImageNet-64x64 and ImageNet-512x512, respectively. Increasing the CFG scale γ (top rows) and progressive MCLR training (bottom rows) produce similar effects, both enhancing class-specific structures in the generated images.

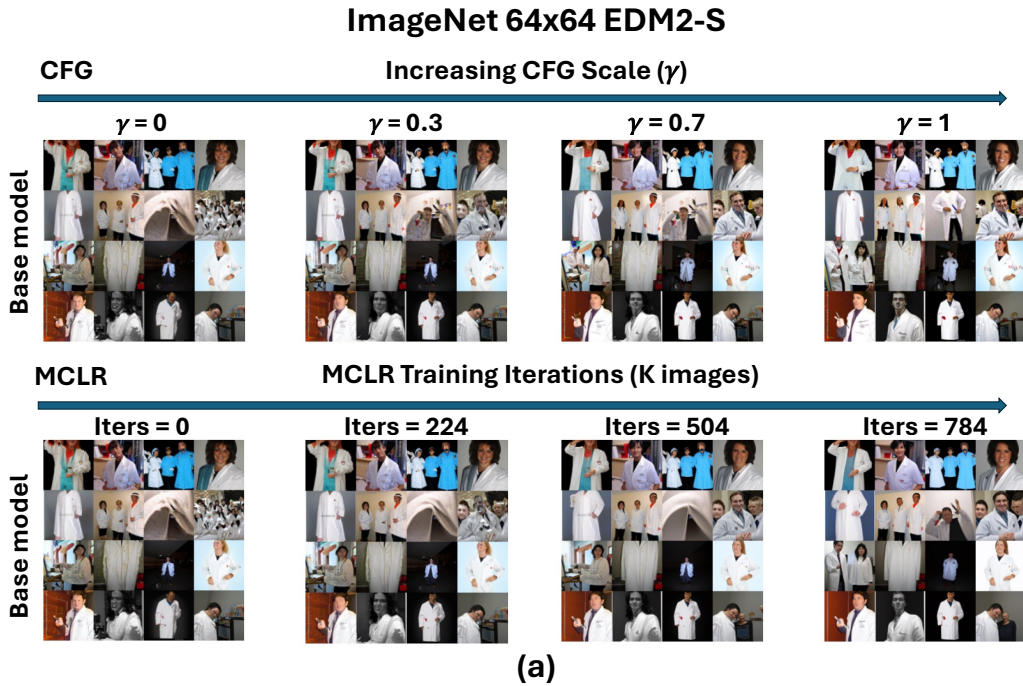


Figure 12. Comparison between CFG and MCLR for Lab Coat (Class 617). (a,b) demonstrate the progressive evolution of generated samples on ImageNet-64x64 and ImageNet-512x512, respectively. Increasing the CFG scale γ (top rows) and progressive MCLR training (bottom rows) produce similar effects, both enhancing class-specific structures in the generated images.

2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089

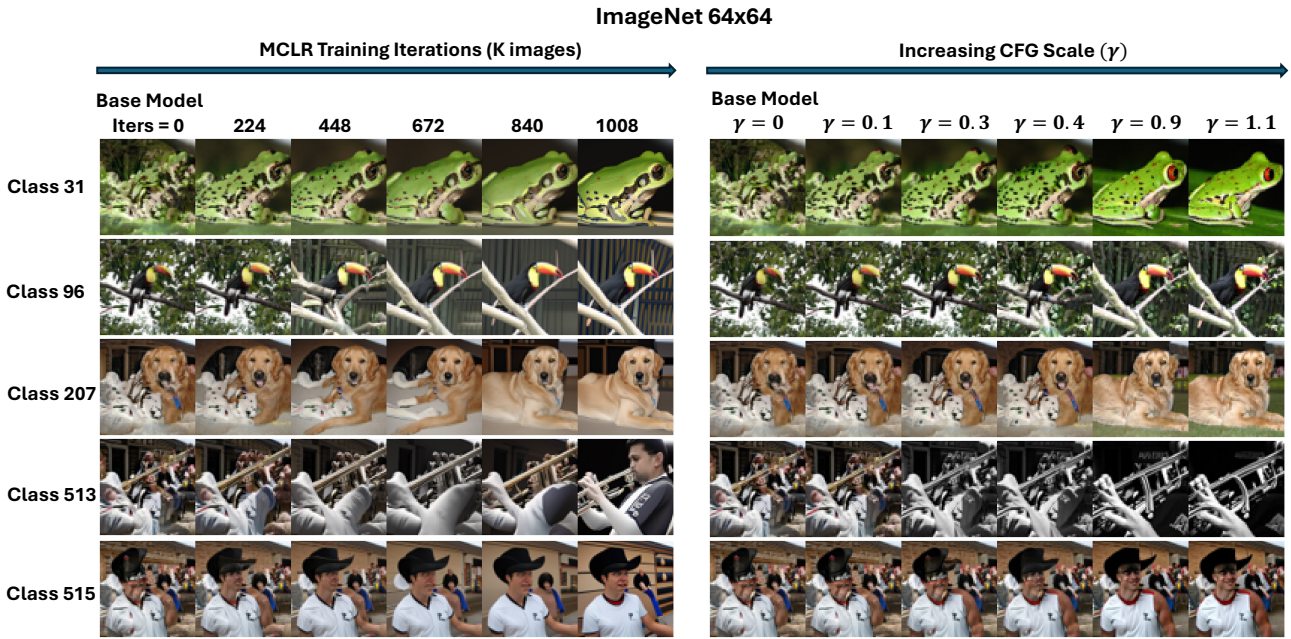


Figure 13. Comparison between CFG and MCLR. Left and right figures demonstrate the progressive evolution of generated samples for MCLR and CFG on ImageNet-64x64, respectively, with all images initialized from the same random noise. Increasing the CFG scale γ and progressive MCLR training produce similar effects, both enhancing class-specific structures in the generated images.

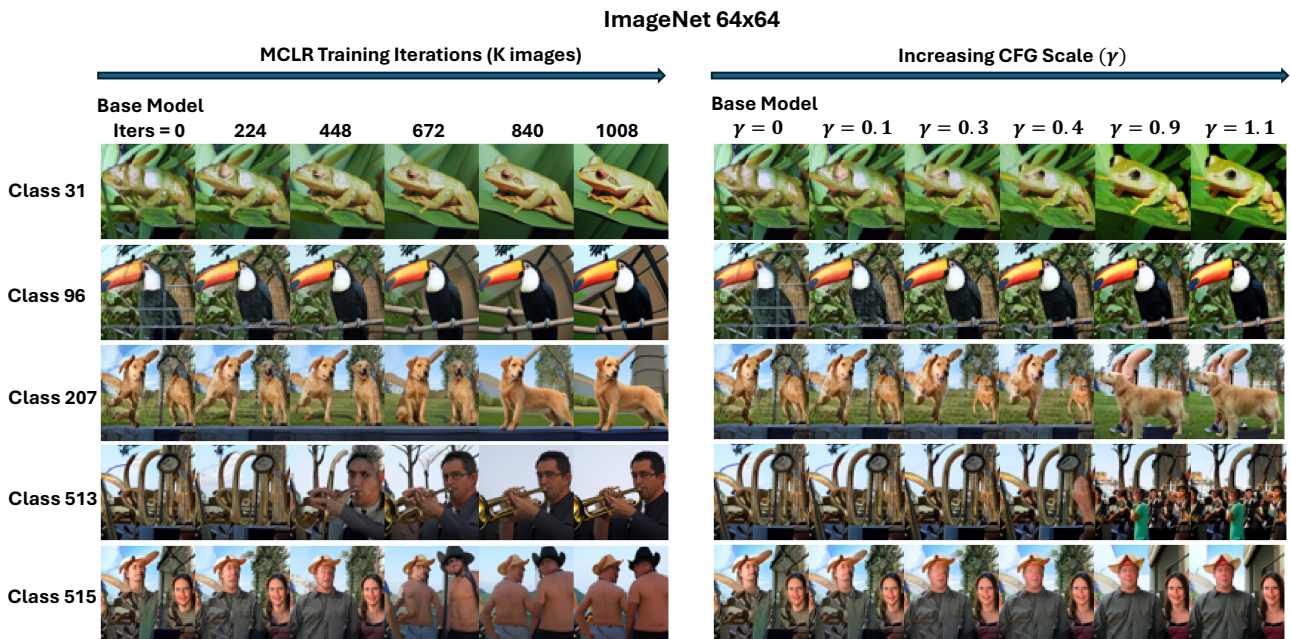


Figure 14. Comparison between CFG and MCLR. Same as Figure 13, but with a different initial random noise.

2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144

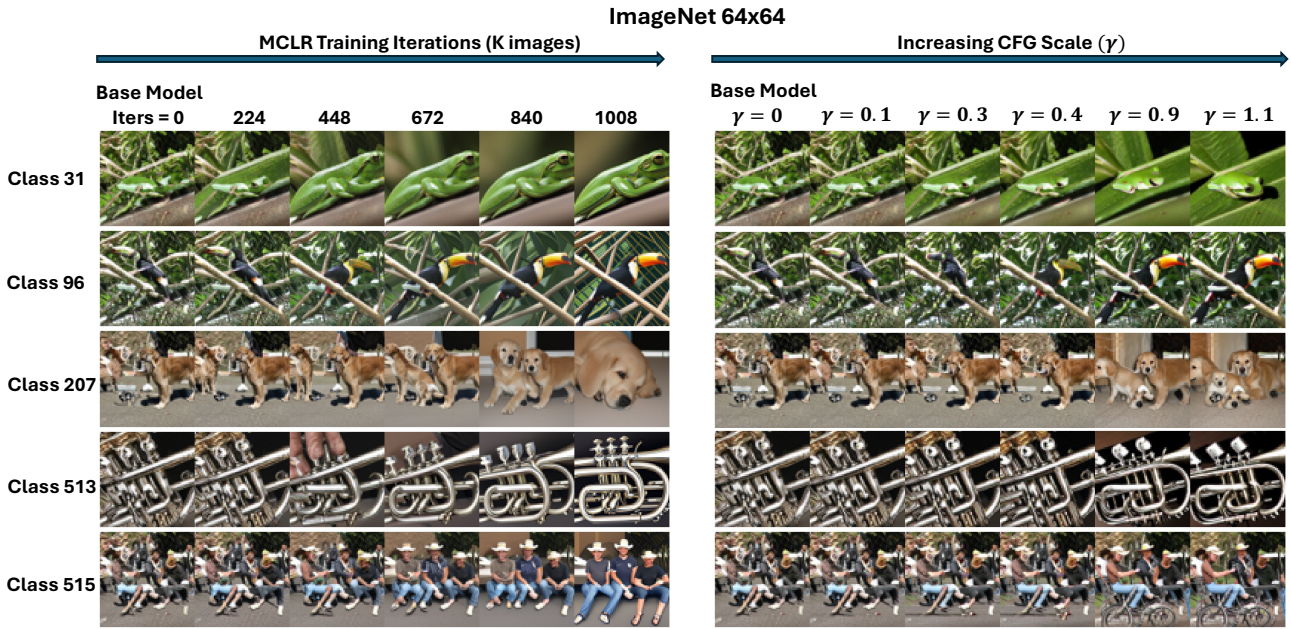


Figure 15. Comparison between CFG and MCLR. Same as Figure 13, but with a different initial random noise.

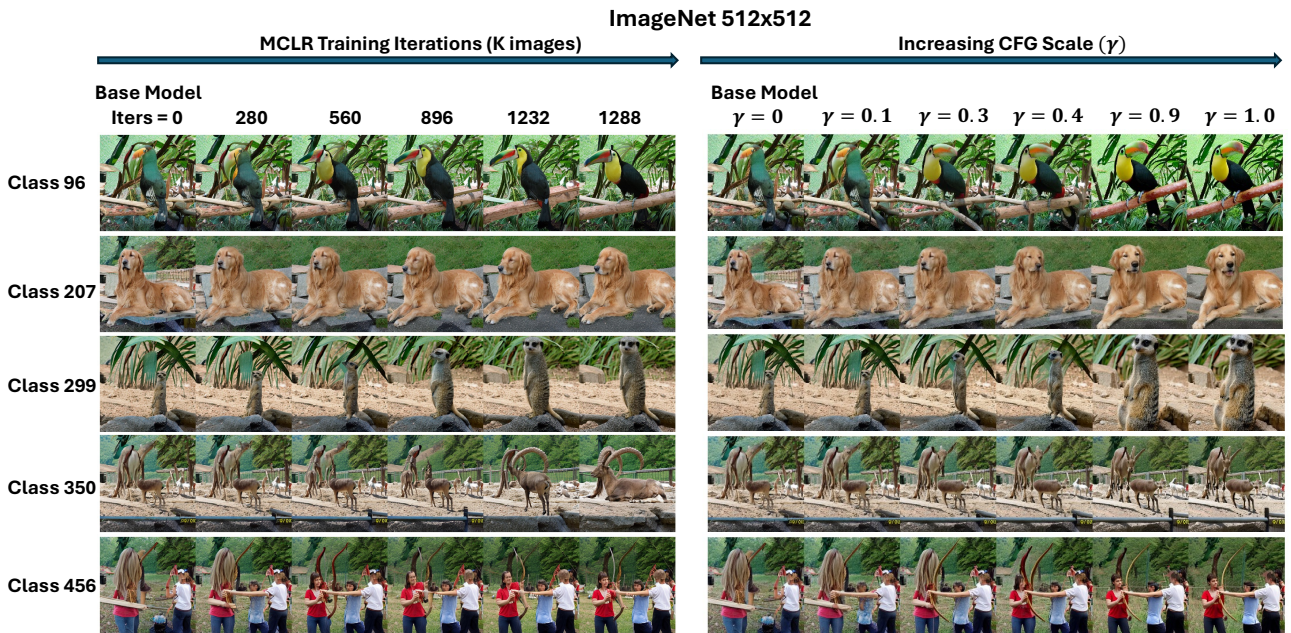


Figure 16. Comparison between CFG and MCLR. Same as Figure 13, but for ImageNet-512x512 with a different initial random noise.

2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199

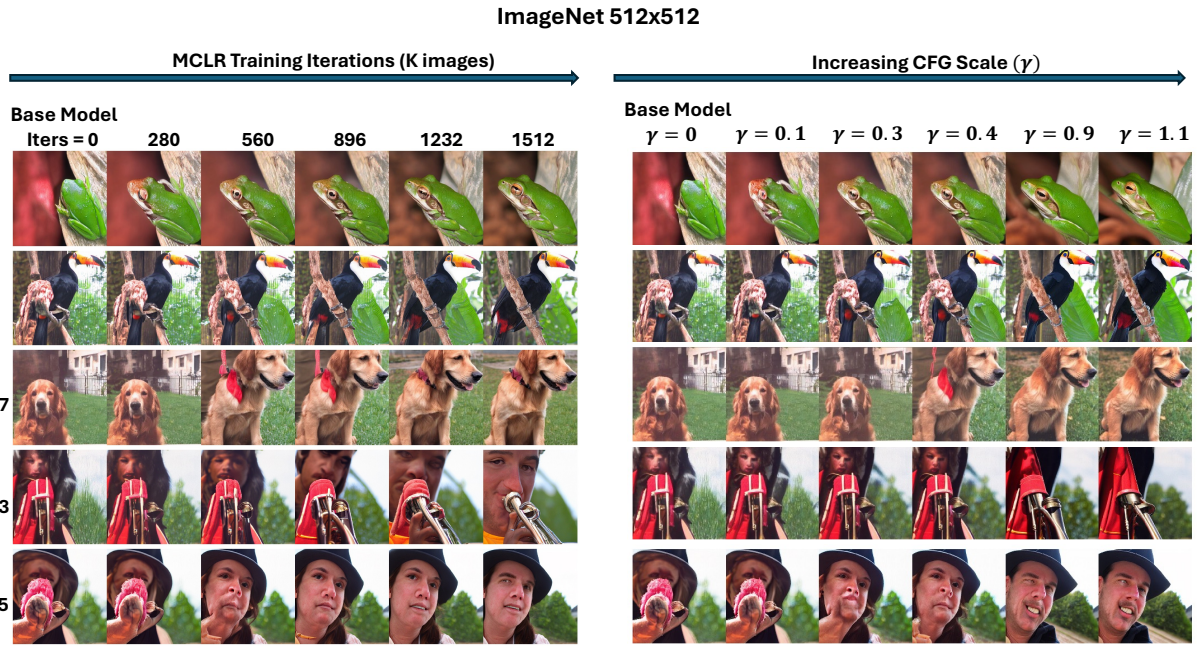


Figure 17. Comparison between CFG and MCLR. Same as Figure 13, but for ImageNet-512x512 with a different initial random noise.

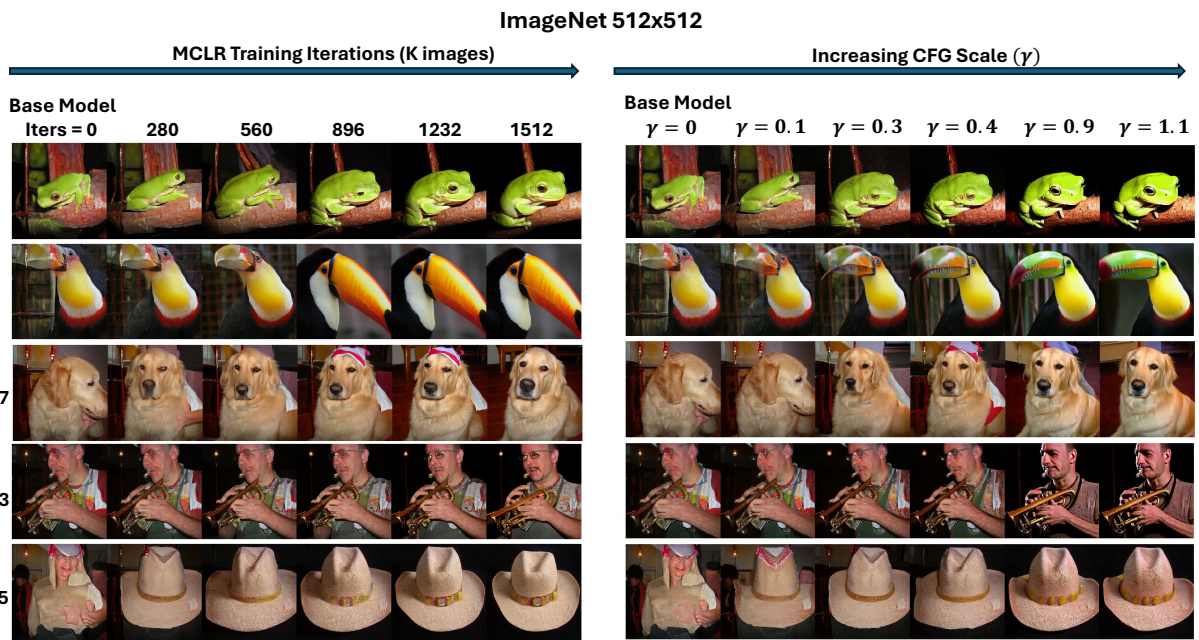


Figure 18. Comparison between CFG and MCLR. Same as Figure 13, but for ImageNet-512x512 with a different initial random noise.

2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254



Figure 19. Comparison between CFG and MCLR. Same as Figure 13, but for ImageNet-512x512 with a different initial random noise.

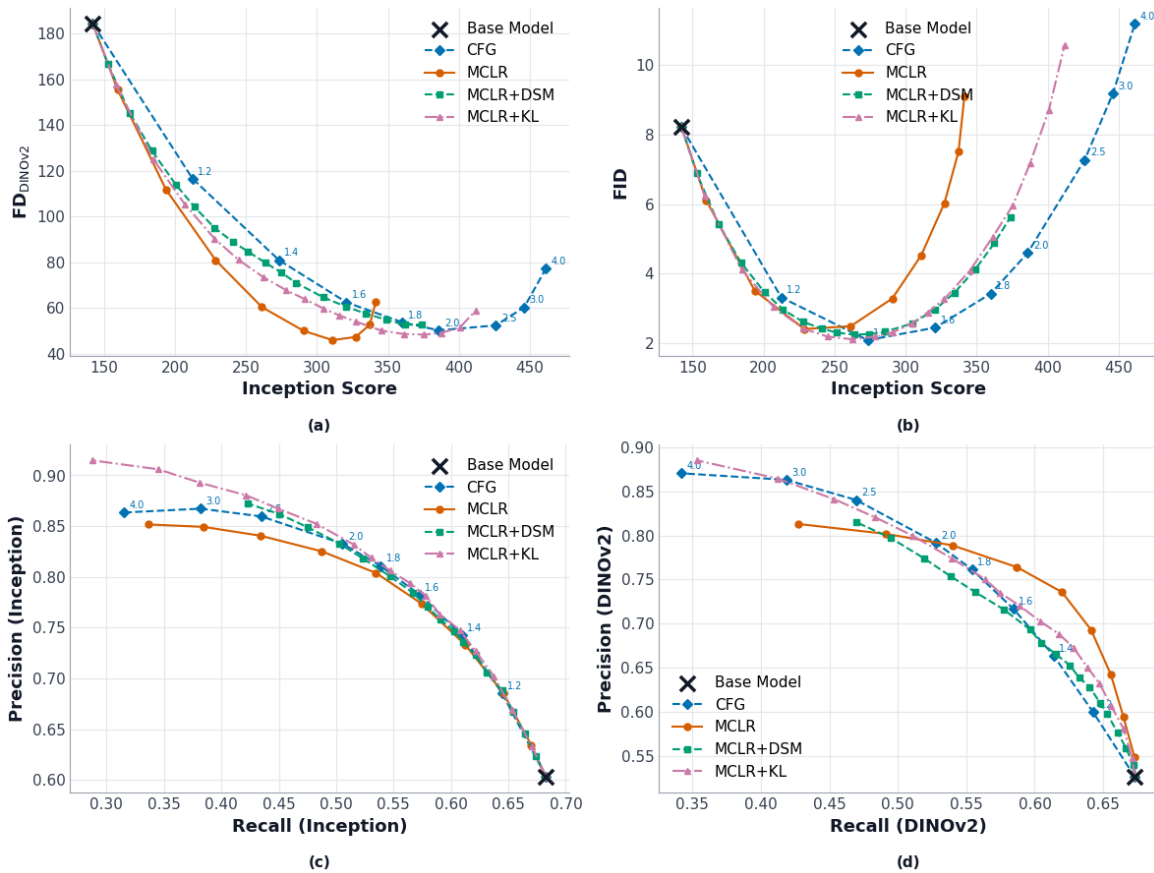


Figure 20. Quantitative Results for SiT-XL/2 trained on ImageNet-256x256. (a), (b), (c), (d) show the evolution of the FD-IS trade-offs and the Precision-Recall trade-offs calculated with Inception and DINOv2 features respectively. We evaluate classifier-free guidance (CFG) scales $\gamma \in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.7, 0.9, 1, 1.5, 2.0, 3.0\}$.

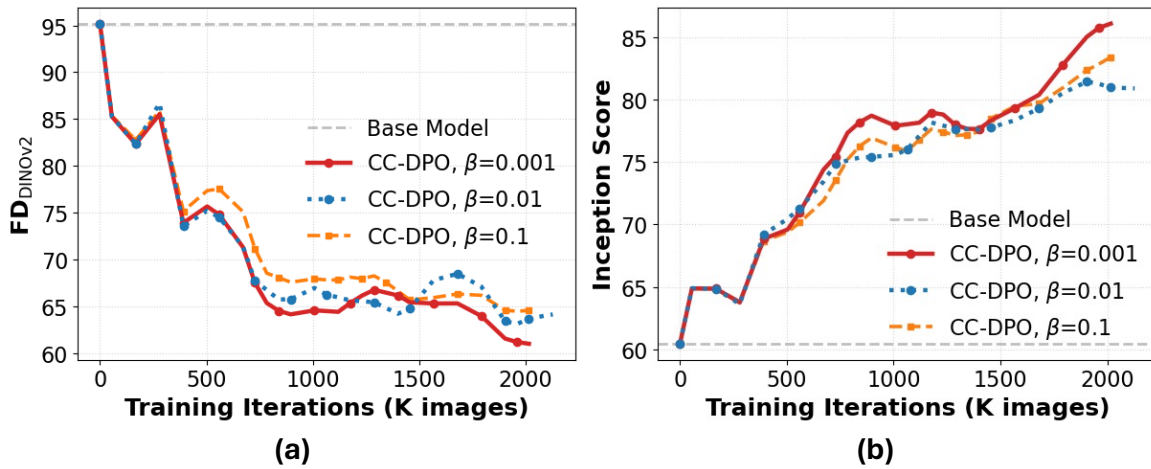


Figure 21. **Effects of β in CC-DPO.** (a,b) compare fine-tuning diffusion models using CC-DPO with different β evaluated on the EDM2-S model trained on ImageNet- 64×64 . Performance is reported in terms of FD_{DINOv2} and Inception Score. CC-DPO achieves stable performance under different β , with a smaller β leads to marginally better performance on selected metrics.

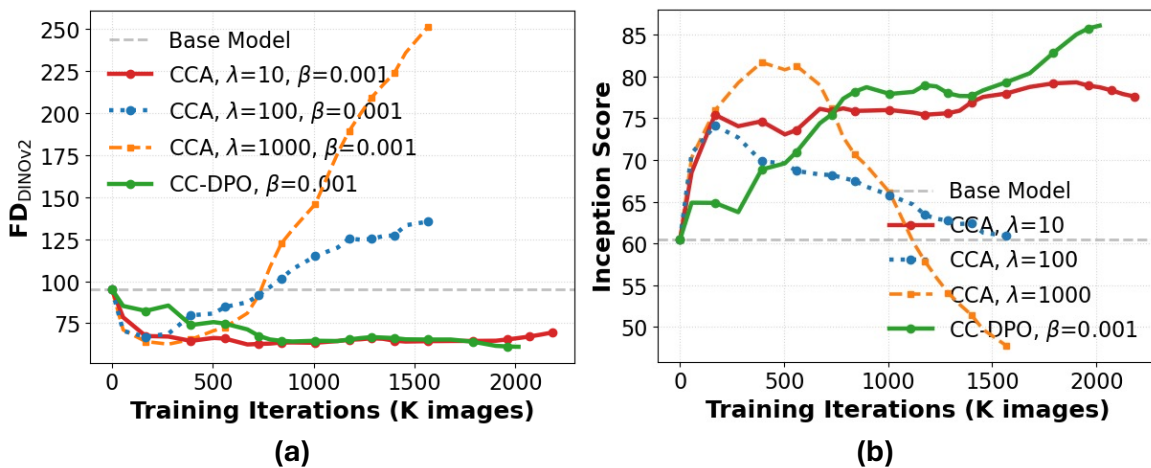


Figure 22. **Effects of λ in CCA.** (a,b) compare fine-tuning diffusion models using CCA with different λ evaluated on the EDM2-S model trained on ImageNet- 64×64 . Performance is reported in terms of FD_{DINOv2} and Inception Score. Note with the same β , CC-DPO achieves comparable performance as CCA while at the same time has fewer hyperparameters.