

000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 MULTIMODAL INFORMATION IS ALL YOU NEED FOR ADVERSARIAL PURIFICATION VIA DIFFUSION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Adversarial defense aims to find true semantic labels of adversarial examples, where diffusion-based adversarial purification as intriguing adversarial defense methods can restore data perturbed by unseen attacks to clean distribution without training classifiers. However, unimodal diffusion-based approaches rely on noise schedules to implicitly preserve labels, whereas recently proposed multimodal variants add textual control but require adversarial training and heavy distillation. Both approaches lack theoretical guarantees. In this work, we propose MultiDAP that uses multimodal diffusion models for adversarial purification. MultiDAP first learn prompts from clean text-image pair data for clean image generation, where context tokens are numerical instead of text templates such as “a photo of a .” for rich contextual information and hence enhance adversarial robustness. Given learned prompts and adversarial examples, MultiDAP then purify inputs via minimizing regularized DDPM losses iteratively for only a few steps. Theoretical guarantees for two phases are also provided. In experiments, our proposed model achieve improvement of zero-shot adversarial defense performance over unimodal diffusion models and multimodal variants with text templates.

1 INTRODUCTION

Adversarial defense is fundamentally concerned with recovering the true semantic label information from adversarial examples which are perturbed by human imperceptible but carefully crafted noise for deep learning classifiers to predict incorrect labels (Goodfellow et al., 2014). Adversarial purification has recently emerged as a promising paradigm for adversarial defense (Shi et al., 2021; Yoon et al., 2021; Nie et al., 2022; Wang et al., 2022a; Bai et al., 2024; Lei et al., 2025). Unlike adversarial training (Croce and Hein, 2020; Laidlaw et al., 2021; Dolatabadi et al., 2022; Wang et al., 2023a), which explicitly trains classifiers on adversarial examples, adversarial purification methods employ generative models to remove adversarial perturbations before classification (Song et al., 2017; Nie et al., 2022). This strategy offers two key advantages. First, it does not need to retrain classifiers on generated attacks, thereby reducing computational overhead. Second, it provides stronger generalization to unseen adversarial attacks, as adversarial purification directly restores clean data distributions. However, early adversarial purification methods are based on generative adversarial networks (GANs) and energy-based models (EBMs) and fall behind adversarial training methods, because of their limited generative power (Nie et al., 2022).

Diffusion models have rapidly become the mainstream approach for adversarial purification due to their remarkable generative power and ability to approximate complex data distributions (Nie et al., 2022; Wang et al., 2022a; Chen et al., 2024a; Zhang et al., 2025; Bai et al., 2024). By progressively adding Gaussian noise and removing it, diffusion models can effectively reconstruct clean samples from corrupted or perturbed inputs, making them particularly well-suited for removing adversarial perturbations (Nie et al., 2022). However, most existing work relies on unimodal diffusion models which attempt to preserve semantic information implicitly by injecting Gaussian noise to a specific level in the forward process and then denoising the input. Hence unimodal approaches often struggle to fully obtain semantic label information, limiting defense against stronger or adaptive attacks.

To address this limitation, one step control purification has recently been proposed as the first multimodal diffusion model for adversarial purification. It leverages ControlNet to include additional modalities (e.g., textual prompts) for adversarial purification and thus preserves more semantic label

054 information than unimodal approaches (Lei et al., 2025). This importance of multi-modal information
 055 was first discovered from the fact that human can easily identify the true label of adversarial examples,
 056 in contrast to deep learning models on solely pixel spaces. The reason is that human cognition
 057 relies on semantic information from the context and is immune to distribution variations induced by
 058 adversarial attacks, whereas the deep learning models classify images via statistical distributional
 059 associations (Zhou et al., 2024) and are vulnerable to adversarial attacks. However, one step control
 060 purification still faces three critical challenges: (i) it lacks theoretical guarantees regarding purification
 061 effectiveness, and (ii) it still relies on adversarial training to learn robust cross-modal alignment,
 062 and more importantly, (iii) without knowledge distillation, the iterative reverse process of diffusion
 063 models will incur substantial computational overhead. These efficiency issues limit their practicality
 064 for real-time or large-scale adversarial defense.

065 In this paper, we propose Multimodal Diffusion for Adversarial Purification (MultiDAP) which
 066 leverages a single text-to-image diffusion model backbone. Unlike unimodal diffusion models for
 067 adversarial purification solely relying on image features, our approach conditions the diffusion model
 068 on textual prompts to infuse semantic information. In order to obtain the prompts which steer
 069 zero-shot adversarial purification, we design a paradigm to learn prompts for stable diffusion models
 070 from clean large-scale text-image pairs. This design not only leverages the powerful generative
 071 capacity of diffusion models but also capitalizes on the rich contextual representations encoded by
 072 the text encoder. With prompt-based conditioning, we introduce a more expressive feature space
 073 that distinguishes adversarial attacks from genuine content more effectively. Furthermore, we also
 074 propose to efficiently purify adversarial examples via prompt-guided likelihood maximization which
 075 only requires a few purification steps. Experiments demonstrate that MultiDAP achieves superior
 076 zero-shot adversarial defense performance compared to unimodal diffusion models on the CIFAR-10,
 077 CIFAR-100 and ImageNet-1K dataset. These results highlight the dual contribution of our work:
 078 introducing a theoretically grounded framework for adversarial purification and delivering practical
 079 improvements for real-world deployment.

080 2 RELATED WORK

081 **082 Unimodal Diffusion Models for Adversarial Defense.** Diffusion models have demonstrated
 083 remarkable performance in generative tasks, owing to their ability to progressively refine noisy
 084 data to high-quality output. This generative nature has been explored for robustness in various
 085 contexts, including adversarial purification (Nie et al., 2022), adversarial training (Wang et al.,
 086 2023b) and robust classification methods (Chen et al., 2024a). A notable application of diffusion
 087 models lies in purification-based defenses, where adversarially perturbed inputs are restored to their
 088 clean counterparts. Methods leveraging guided diffusion models have shown efficacy in removing
 089 perturbations while preserving the underlying data features, making them suitable for tasks like
 090 classification (Lee and Kim, 2023; Xiao et al., 2022; Bai et al., 2024; Yeh et al., 2024). Additionally,
 091 diffusion-based classifier have gained traction by integrating generative and discriminative modelling
 092 (Zimmermann et al., 2021; Clark and Jaini, 2023; Chen et al., 2024a). Their robustness to input
 093 perturbations and adversarial attack is attributed to optimal empirical score function (Chen et al.,
 094 2024b). While these approaches highlight the versatility of diffusion models in adversarial defense,
 095 they often face challenges in efficiency, effectiveness from unimodality and theoretical guarantees,
 096 motivating further investigations.

097 **098 Multimodal Approaches in Adversarial Defense.** Recent advances in vision-language models
 099 (VLMs) have demonstrated potentials of multimodal information in improving robustness. Models
 100 such as CLIP (Radford et al., 2021) learn joint embeddings of images and text, enabling strong
 101 cross-modal alignment that provides richer semantic priors than unimodal vision models. This
 102 multimodal alignment has inspired adversarial finetuning (Schlarmann et al., 2024), adversarial
 103 prompt tuning (Zhang et al., 2024; Li et al., 2024; Sheng et al., 2025), and multimodal defenses
 104 leveraging vision-language pretraining (Wang et al., 2025). These approaches hold clear advantages:
 105 auxiliary modalities such as text can act as high-level semantic constraints, guiding models toward
 106 correct semantic label predictions. Nevertheless, current multimodal robustness methods face critical
 107 limitations. Many rely on adversarial training to establish robust cross-modal alignment, and recently
 108 proposed first multimodal diffusion model for adversarial purification: one step control purification
 109 in particular suffer from substantial computational overhead due to multi-step denoising if knowl-

108 edge distillation is not used (Lei et al., 2025). Furthermore, theoretical guarantees remain largely
 109 absent, leaving their robustness difficult to formally assess. These challenges motivate the need
 110 for approaches that combine multimodal semantic priors with efficiency and provable purification
 111 guarantees—precisely the focus of our proposed Multimodal Diffusion for Adversarial Purification.
 112

113 3 MULTIMODAL DIFFUSION MODELS FOR ADVERSARIAL PURIFICATION

115 3.1 PROBLEM SETUP

117 Let $x \in \mathcal{X}$ denote a clean input with label y , and $x^{adv} = x + \delta$ be an adversarial example generated
 118 under perturbation constraint $\|\delta\|_p \leq \epsilon$. Here the perturbation δ is constrained under an ℓ_p -norm
 119 threat model, with $p \in \{2, \infty\}$ being the most common cases. The ℓ_∞ attack bounds the maximum
 120 per-pixel distortion, ensuring imperceptibility, while the ℓ_2 attack restricts the overall perturbation.
 121

122 Adversarial purification aims to transform an adversarial input x^{adv} back to a sample x^{pur} that lies
 123 close to the clean data manifold, such that $f(x^{pur}) = y$. Recent works have demonstrated that
 124 diffusion models are particularly well suited for this task, due to their strong generative ability to
 125 approximate complex data distributions (Nie et al., 2022).

126 A diffusion model (Song et al., 2020) defines a forward noising process that gradually perturbs clean
 127 data x_0 into Gaussian noise through a sequence of latent variables $\{x_t\}_{t=0}^T$:

$$128 \quad 129 \quad q(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{1 - \beta_t} x_{t-1}, \beta_t I\right),$$

130 where $\{\beta_t\}$ is a variance schedule (Ho et al., 2020). This process ensures that as $t \rightarrow T$, the sample
 131 x_T approaches pure noise, as T is large enough. The reverse denoising process is parameterized by a
 132 neural network ϵ_θ , which predicts the added noise and iteratively reconstructs clean data:
 133

$$134 \quad p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)).$$

135 The mean term $\mu_\theta(x_t, t)$ is computed from this noise prediction via a closed-form reparameterization:
 136 $\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right)$, where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$. The covariance
 137 $\Sigma_\theta(x_t, t)$ is typically fixed by the variance schedule $\{\beta_t\}$, though some variants allow it to be partially
 138 learned for improved sample quality (Nichol and Dhariwal, 2021). Together, μ_θ and Σ_θ define the
 139 Gaussian reverse step, while ϵ_θ remains the core predicted quantity that drives the denoising trajectory.
 140

141 For adversarial purification, the intuition is to inject the adversarial input x^{adv} into the forward
 142 process at a chosen noise level t , so that adversarial perturbations are drowned out by Gaussian noise
 143 (Nie et al., 2022). Then, the reverse process denoises x_t step by step, ideally converging to a purified
 144 sample x^{pur} close to the clean distribution. Formally, the purification mapping can be written as:
 145

$$146 \quad x^{pur} \sim P(x^{adv}) = p_\theta(x^{pur} | x_t^{adv}, t), \quad \text{with } x_t^{adv} \sim q(x_t | x^{adv}).$$

147 Here x_t^{adv} denotes the adversarial input injected into the forward noising process at step t , and x^{pur} is
 148 the purified output after reverse diffusion. This framework has achieved strong empirical robustness
 149 across various benchmarks (Nie et al., 2022; Wang et al., 2022b; Chen et al., 2024a).

150 However, existing diffusion-based purification suffers from two main drawbacks: (i) the denoising
 151 process is essentially *unimodal*, since it is conditioned only on Gaussian noise schedules without
 152 leveraging explicit semantic cues (text prompts), which limits its ability to preserve class-consistent
 153 information; and (ii) the multi-step reverse process is computationally expensive, making such
 154 defenses inefficient for real-time deployment. These limitations motivate our proposed Multimodal
 155 Diffusion for Adversarial Purification with explicit prompt guidance and improved efficiency.
 156

157 3.2 STABLE DIFFUSION WITH PROMPT LEARNING

159 While diffusion-based purification can remove adversarial perturbations, prior defenses typically rely
 160 on *small, unimodal* diffusion models to approximate the data distribution (Nie et al., 2022). Limited
 161 representational capacity often leads to suboptimal likelihood estimates and unstable denoising
 162 trajectories, where semantic information may not be faithfully preserved.

To address this limitation, we adopt Stable Diffusion—a large-scale latent diffusion model (LDM)—as our backbone (Rombach et al., 2022). Pretrained on massive image–text corpora, Stable Diffusion provides substantially stronger modeling power and a richer, more informative likelihood landscape than small diffusion models. Moreover, operating in a compact latent space enables high-resolution synthesis with improved efficiency compared to pixel-space diffusion (Rombach et al., 2022; Dhariwal and Nichol, 2021). This stronger backbone lets our purifier start denoising from a more faithful approximation of the clean data manifold, reducing reliance on long reverse diffusion chains and mitigating semantic drift.

Formally, Stable Diffusion operates in a latent space defined by a variational autoencoder (VAE). Given an input image x , the encoder maps it into a compact latent representation $z = \mathcal{E}_{\text{VAE}}(x)$, where \mathcal{E}_{VAE} denote the encoder. Given a class label or textual description y , we obtain a prompt embedding through a text encoder $e_p = \mathcal{E}_{\text{text}}(p)$, where p is the input text, such as ‘a photo of a cat’. In our approach, these embeddings serve as semantic conditions that guide the purification process. The denoising network $\epsilon_\theta(x_t, t, e_p)$ is implemented as a U-Net with cross-attention, which predicts the noise at each timestep. The denoising network $\epsilon_\theta(x_t, t, e_p)$ is trained with the standard denoising diffusion probabilistic model (DDPM) objective (Ho et al., 2020), which treats noise prediction as score matching:

$$\mathcal{L}_{\text{DDPM}}(\theta) = \mathbb{E}_{x_0, e_p, t, \epsilon \sim \mathcal{N}(0, I)} \left[\left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t, e_p) \right\|_2^2 \right],$$

where $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ denote the variance schedule. This loss enforces the network to accurately predict the added Gaussian noise at each timestep, which is equivalent to maximizing a variational lower bound on the conditional data likelihood. In our case, the conditioning e_p provides semantic priors that explicitly align the denoising trajectory with the true class, thereby enhancing the stability and fidelity of purification.

Prompt Learning Objective. A central challenge for purification-based defenses lies in the accuracy of likelihood estimation during denoising. Although Stable Diffusion provides a strong backbone, its conditioning typically depends on fixed or manually designed text prompts, which may be generic and fail to provide task-specific guidance. Such limitations are particularly critical for adversarial purification, where the model must recover the clean data distribution from inputs corrupted by imperceptible but adversarial perturbations. To overcome this issue, we propose a *prompt learning* module that explicitly optimizes prompt embeddings from clean data, allowing the model to acquire semantic priors that are robust to adversarial noise.

In general, a prompt p can be represented as a concatenation of M learnable context tokens,

$$p_{\text{context}} = [v_1, v_2, \dots, v_M],$$

where each $v_m \in \mathbb{R}^d$ has the same dimensionality as the text encoder’s word embeddings (e.g., $d = 512$ for CLIP). In prior works, such context tokens are often combined with a class-specific token (e.g., the word ‘cat’), yielding a class-dependent prompt $p = [p_{\text{context}}, p_{\text{class}}]$ that provides label-conditioned guidance. By contrast, our objective is to design a *class-agnostic prompt* that captures global semantic priors without relying on class labels. This choice is crucial for adversarial purification, since the ground-truth label of an adversarial input is typically unknown at inference time. We therefore optimize a shared prompt vector p directly from clean data, such that it enhances the unconditional likelihood estimation of the diffusion model.

Our prompt learning module is optimized by reusing the standard DDPM noise-prediction loss, with the key difference that only the prompt parameters p are updated while the diffusion backbone θ remains frozen:

$$\mathcal{L}_{\text{prompt}}(p) = \mathbb{E}_{x_0, t, \epsilon} \left[\left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t, p) \right\|_2^2 \right].$$

Optimizing this loss is equivalent to maximizing a variational lower bound on the likelihood $p_\theta(x_0 | p)$. Thus, the learned prompt p^* serves as a universal semantic prior that stabilizes the denoising trajectory and improves the fidelity of adversarial purification without requiring class labels or adversarial training.

Similar to prompt learning in CLIP (Zhou et al., 2022), we optimize the learnable context tokens using a gradient-based method, such as Adam (Kingma, 2014). In each training iteration, we sample

216

Algorithm 1: Prompt Learning on Stable Diffusion (class-agnostic)

217

Input: Frozen diffusion backbone θ (VAE, U-Net ϵ_θ), clean images $\{x^{(b)}\}$, steps T , optimizer (Adam), prompt length M , iters N .

218

Output: Learned class-agnostic prompt $p^* = [v_1, \dots, v_M]$.

219

```

1 Initialize learnable tokens  $p = [v_1, \dots, v_M]$  (random or text-init).
2 for  $n = 1, \dots, N$  do
3   Sample a mini-batch  $\{x^{(b)}\}$ ; sample  $t \sim \text{Unif}(\{1, \dots, T\})$  and  $\epsilon \sim \mathcal{N}(0, I)$ .
4    $x_t \leftarrow \sqrt{\bar{\alpha}_t} x^{(b)} + \sqrt{1 - \bar{\alpha}_t} \epsilon$ 
5    $\mathcal{L}_{\text{prompt}} \leftarrow \|\epsilon - \epsilon_\theta(x_t, t, p)\|_2^2$  (average over batch).
6   Update  $p \leftarrow \text{Adam}(p, \nabla_p \mathcal{L}_{\text{prompt}})$  while keeping  $\theta$  frozen.
7 return  $p^* \leftarrow p$ .
```

220

221

222

223

224

225

226

227

228

229

230

a mini-batch of clean data points $x^{(b)}$, apply the forward diffusion process to obtain noisy latents $x_t^{(b)}$ with Gaussian noise ϵ , and evaluate the prompt loss $\mathcal{L}_{\text{prompt}}$ with the current tokens $p = [v_1, \dots, v_M]$.

231

The gradients are then backpropagated through the denoising network $\epsilon_\theta(x_t^{(b)}, t, p)$ to update the prompt parameters. This process is repeated until convergence, yielding a shared prompt vector p^* that minimizes the denoising objective. The detailed optimization procedure is summarized in Algorithm 1. Compared to full model fine-tuning, optimizing only a small set of prompt parameters significantly reduces trainable variables, which mitigates overfitting and keeps computational cost manageable, while still providing strong semantic guidance for purification.

232

233

234

235

236

237

238

Theoretical Guarantee: Prompt Learning Improves Likelihood. We show that optimizing the class-agnostic prompt p with the DDPM objective monotonically increases a variational lower bound (ELBO) of the unconditional data likelihood under a fixed diffusion model θ .

239

240

241

242

243

244

245

Theorem 1 (Prompt learning improves the likelihood lower bound). *Let $x_0 \sim p_{\text{data}}$ denote clean latents, and let $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$ with $\epsilon \sim \mathcal{N}(0, I)$. Fix the diffusion backbone parameters θ , and optimize only the prompt p using the DDPM objective. Then the optimal prompt $p^* = \arg \min_p \mathcal{L}_{\text{prompt}}(p)$ maximizes the evidence lower bound (ELBO) on the data likelihood $p_\theta(x_0 | p)$,*

$$\log \underline{p}_\theta(x_0 | p^*) \geq \log \underline{p}_\theta(x_0 | p), \quad \forall p,$$

246

where $\log \underline{p}_\theta$ denotes the variational lower bound.

247

248

249

Moreover, $\nabla_p \mathcal{L}_{\text{VLB}}(p)$ and $\nabla_p \mathcal{L}_{\text{prompt}}(p)$ are colinear since the weights w_t are positive. Thus updating p along $-\nabla_p \mathcal{L}_{\text{prompt}}$ strictly decreases \mathcal{L}_{VLB} for sufficiently small step size, thereby monotonically increasing the likelihood.

250

251

252

253

254

255

Corollary 1 (Score matching view). *The objective $\mathcal{L}_{\text{prompt}}$ is equivalent to minimizing a weighted Fisher divergence between the conditional score $\nabla_{x_t} \log p_\theta(x_t | p)$ and the forward diffusion score $\nabla_{x_t} \log q(x_t | z_0)$. Hence optimizing p aligns the model score with the true score, which directly improves the data likelihood $\log p_\theta(x_0 | p^*) \geq \log p_\theta(x_0 | p), \forall p$.*

256

257

3.3 PROMPT-GUIDED LIKELIHOOD MAXIMIZATION FOR PURIFICATION

258

259

260

Given an adversarial input x^{adv} , we purify it by maximizing the model likelihood under the learned class-agnostic prompt p^* while using the pretrained diffusion backbone θ . We obtain a noisy image by the forward diffusion

$$x_{t^*}^{adv} = \sqrt{\bar{\alpha}_{t^*}} x_0^{adv} + \sqrt{1 - \bar{\alpha}_{t^*}} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I).$$

261

262

Our goal is to recover an x_0 that maximizes the posterior (or the conditional likelihood surrogate)

$$x_0^{\text{pur}} \in \arg \max_{x_0} \log p_\theta(x_{t^*}^{adv} | x_0, p^*) + \log p(x_0), \quad (1)$$

263

264

265

where $p(x_0)$ is the prior. Maximizing (1) is intractable directly, so we instead minimize the purification “simple loss” with respect to the image variable x_0 , while conditioning on the learned prompt p^* :

266

267

268

269

$$x_0^{\text{pur}} \in \arg \min_{x_0} \underbrace{\mathbb{E}_{t, \epsilon} \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t, p^*)\|_2^2}_{=: \mathcal{L}_{\text{DDPM}}(x_0; p^*)} + \lambda \mathcal{R}(x_0, x^{adv}), \quad (2)$$

270

Algorithm 2: Purification via Regularized DDPM-Loss Minimization in Pixel Space

271

Input: Adversarial image x^{adv} , learned prompt p^* , frozen θ , steps $T_1 T_2$, Purification steps N (e.g., 5), step size η , optional regularizer weight λ .

273

Output: Purified image x^{pur} .

274

```

1  $x_0^{(0)} \leftarrow x^{adv}$  for  $n = 0, \dots, N - 1$  do
2   Sample  $t \sim \text{Unif}(\{T_1, \dots, T_2\})$  and  $\epsilon \sim \mathcal{N}(0, I)$ .
3    $x_t \leftarrow \sqrt{\bar{\alpha}_t} x_0^{(n)} + \sqrt{1 - \bar{\alpha}_t} \epsilon$ 
4    $\mathcal{L}_{pur} \leftarrow \|\epsilon - \epsilon_\theta(x_t, t, p^*)\|_2^2 + \lambda \mathcal{R}(x_0^{(n)}, x^{adv})$ 
5    $g \leftarrow \nabla_{x_0} \mathcal{L}_{pur}$ 
6    $x_0^{(n+1)} \leftarrow \Pi_{[0,1]}(x_0^{(n)} - \eta g)$ 
7    $x^{pur} \leftarrow x_0^{(N)}$ 
8 return  $x^{pur}$ .

```

284

285

286 where t is sampled uniformly from $\{1, \dots, T\}$, $\epsilon \sim \mathcal{N}(0, I)$, and \mathcal{R} is an optional proximity or
287 naturalness regularizer (e.g., $\mathcal{R}(x_0, x^{adv}) = \|x_0 - x^{adv}\|_2^2$ or a TV prior). By Theorem 1, minimizing
288 \mathcal{L}_{DDPM} w.r.t. p tightens the ELBO; when optimizing w.r.t. x_0 , Eq. (2) serves as a surrogate that
289 increases the conditional likelihood under p^* .

290

291

Gradient and Update. Let $x_t(x_0, \epsilon) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$. The gradient of Eq. (2) is

292

293
$$\nabla_{x_0} \mathcal{L}_{Pur} = \nabla_{x_0} \mathbb{E}_{t, \epsilon} \left[\sqrt{\bar{\alpha}_t} \nabla_{x_t} \|\epsilon - \epsilon_\theta(x_t, t, p^*)\|_2^2 \right] + \lambda \nabla_{x_0} \mathcal{R}(x_0, x^{adv}),$$

294

295

where we define $\mathcal{L}_{Pur} = \mathcal{L}_{DDPM} + \lambda \nabla_{x_0} \mathcal{R}(x_0, x^{adv})$ and we perform a few iterations of gradient
descent with box constraints:

296

297

$$x_0^{(k+1)} = \Pi_{[0,1]} \left(x_0^{(k)} - \eta \nabla_{x_0} \mathcal{L}_{Pur}(x_0^{(k)}; p^*) \right), \quad x_0^{(0)} = x^{adv},$$

298

299

where $\Pi_{[0,1]}$ clips pixels to the valid range and η is the step size. In practice, we estimate the
300 expectations with a single (t, ϵ) per iteration and use 5–10 steps; the prompt guidance p^* stabilizes
301 the descent by injecting high-level semantics, yielding fast and faithful purification in pixel space.
302 Besides, we adopt the proximity regularizer $\mathcal{R}(x_0, x^{adv}) = \|x_0 - x^{adv}\|_2^2$ with a fixed weight
303 $\lambda = 0.9$, which encourages purified outputs to remain close to the original adversarial inputs while
304 removing perturbations. The overall purification process is summarized in Algorithm 2.

305

306

Theory: Stochastic One-Sample Purification and Few-Step Ascent We analyze the pixel-space
purification objective

307

$$\mathcal{L}_{pur} = \mathcal{L}_{DDPM}(x_0; p^*) + \lambda \mathcal{R}(x_0, x^{adv}). \quad (3)$$

308

309

Let $x_t(x_0, \epsilon) = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$ and denote $\ell(x_0; t, \epsilon) = \|\epsilon - \epsilon_\theta(x_t, t, p^*)\|_2^2$. We update x_0 by
projected SGD with one sampled (t, ϵ) per step:

310

311

$$x_0^{(k+1)} = \Pi_{[0,1]} \left(x_0^{(k)} - \eta g(x_0^{(k)}; t_k, \epsilon_k) \right), \quad g(x_0; t, \epsilon) =: \nabla_{x_0} (\ell(x_0; t, \epsilon) + \lambda \mathcal{R}(x_0, x^{adv})).$$

312

313

Assumptions. (A1) *Smoothness:* $\mathcal{L}_{Pur}(x_0; p^*)$ is L -smooth on $[0, 1]^d$. (A2) *Bounded variance:*
 $\mathbb{E}[\|g(x_0; t, \epsilon) - \nabla \mathcal{L}_{Pur}(x_0; p^*)\|_2^2] \leq \sigma^2$. (A3) *Regularizer:* \mathcal{R} is convex and L_R -smooth (e.g.,
 $\|x_0 - x^{adv}\|_2^2$ or TV with smooth surrogate).

314

315

Lemma 1 (Unbiased one-sample gradient with bounded variance). *With the reparameterization*
 $z_t(x_0, \epsilon)$, *the stochastic gradient is unbiased:*

316

317

$$\mathbb{E}_{t, \epsilon} [g(x_0; t, \epsilon)] = \nabla_{x_0} \mathcal{L}_{Pur}(x_0; p^*),$$

318

and satisfies Assumption (A2).

319

320

Proof Sketch. Differentiate under the expectation using reparameterization; the Jacobian $\partial x_t / \partial x_0 = \sqrt{\bar{\alpha}_t} I$ is deterministic. Linearity of expectation and the uniform sampling of t give the unbiasedness. Bounded variance follows from (A1) and standard Lipschitz/activation bounds on ϵ_θ . \square

324 **Theorem 2** (Expected descent of the purification loss). *Under (A1)–(A3), let $\eta \leq 1/(2L)$. Then the*
 325 *projected SGD iterate satisfies*

$$327 \quad \mathbb{E}[\mathcal{L}_{Pur}(x_0^{(k+1)}; p^*)] \leq \mathbb{E}[\mathcal{L}_{Pur}(x_0^{(k)}; p^*)] - \frac{\eta}{2} \mathbb{E}[\|\nabla \mathcal{L}_{Pur}(x_0^{(k)}; p^*)\|_2^2] + \frac{\eta^2 L}{2} \sigma^2.$$

329 *Consequently, after K steps,*

$$331 \quad \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla \mathcal{L}_{DDPM}(x_0^{(k)}; p^*)\|_2^2] \leq \frac{2(\mathcal{L}_{DDPM}(x_0^{(0)}; p^*) - \mathcal{L}_{inf})}{\eta K} + \eta L \sigma^2,$$

334 *where \mathcal{L}_{inf} is the infimum over $[0, 1]^d$.*

336 *Proof Sketch.* Apply the standard smoothness (descent) lemma to the projected step and take expectations. Use Lemma 1 to replace $\mathbb{E}[g]$ with the true gradient and bound the variance term by σ^2 .
 337 Summing the per-step inequality yields the average-gradient bound. \square

340 We further verify the practical validity of the assumptions used in Theorem 2 through empirical
 341 measurements of the local Lipschitz constant and gradient variance; see Appendix A.3 for details.

343 **Why Single (t, ϵ) and Few Steps Suffice.** The one-sample estimator is unbiased but noisy; this
 344 *stochasticity* serves as exploration that helps escape local pixel-level artifacts, while Theorem 2
 345 guarantees expected descent provided η is small. Moreover, the bound shows $O(1/K)$ decay of the
 346 average gradient norm up to a variance floor $\eta L \sigma^2$, so a *small fixed* number of steps ($N = 5$ – 20)
 347 already yields a measurable reduction of the loss/ELBO gap—matching our practice.

348 4 EXPERIMENTS

351 4.1 EXPERIMENTAL DESIGN

353 This section presents an extensive empirical study to validate the effectiveness of the proposed
 354 Multimodal Diffusion for Adversarial Purification (MultiDAP). We describe experimental setups
 355 (datasets and implementation), report quantitative and qualitative results under adversarial attacks,
 356 and provide ablations dissecting the contributions of different design choices.

357 **Datasets and Model Architectures.** We evaluate our method on three standard benchmarks:
 358 CIFAR-10, CIFAR-100, and ImageNet-1K. All input images are resized to 256×256 to match
 359 the input resolution of Stable Diffusion. For CIFAR-10 and CIFAR-100, we utilize their full clean
 360 training sets (50,000 images) to optimize the class-agnostic prompt embeddings. For the large-scale
 361 ImageNet-1K, to demonstrate data efficiency and scalability, we optimize the prompt using only
 362 a randomly sampled subset of 48,000 training images. We evaluate purification performance on a
 363 held-out subset of 512 adversarial test images sampled from the standard test split for each dataset.
 364 We use miniSD-diffusers (Lambda Labs, 2022) as a frozen stable diffusion generative backbone for
 365 purification. By default we adopt a class-agnostic prompt, parameterized as a set of M learnable
 366 context tokens that are shared across all images and classes. For ablation studies, we also evaluate
 367 hand-crafted class-agnostic prompts (e.g., “a photo of”), which provide weaker guidance compared
 368 to our learned prompt but highlight the effectiveness of explicitly optimizing context tokens. The
 369 classifiers are WideResNet70-16, WideResNet-28-10, and ResNet-50 for CIFAR-10, CIFAR-100,
 370 and ImageNet, respectively. All classifiers are pretrained on clean datasets.

371 **Implementation Details.** We initialize miniSD-diffuser and freeze all network parameters except
 372 for the learnable prompt embeddings. During prompt learning, we optimize learnable tokens of
 373 dimension $d = 768$ using the Adam optimizer where the learning rate is 2×10^{-4} . For CIFAR
 374 datasets, we set the prompt length $M = 16$ and use a batch size of 64, training for 10 epochs. For
 375 ImageNet, we increase the prompt capacity to $M = 64$ and use a batch size of 8, training for 1 epoch.
 376 We incorporate an expectation of the noise ϵ and timestep t while training, ensuring consistency with
 377 the diffusion objective. Note that adversarial attacks are imposed to input images. The learned prompt
 p^* is not attacked or optimized during the attacking and purification process. During purification,

378 Table 1: Clean and robust accuracy (%) on CIFAR-10. Robust results under AutoAttack are reported
 379 for ℓ_∞ ($\epsilon = 8/255$) and ℓ_2 ($\epsilon = 0.5$). The last column reports ℓ_∞ PGD-20 with 10 random restarts
 380 (step size $\alpha = \epsilon/4$). For DiffPure, $t_1 = 0.125$ and $t_2 = 0.1$ denote the time scales used. Our
 381 method (MultiDAP) uses a class-agnostic prompt and only 5 purification steps. We also report an
 382 ablation using the fixed template prompt “a photo of a”, which is widely used in CLIP-based zero-shot
 383 classification. **Bold** denotes the best, underline the second best, and shading the third best.

Method	Architecture	Clean Acc	AA (ℓ_∞)	AA (ℓ_2)	PGD-20
AT-DDPM- ℓ_∞	WRN28-10	88.87	63.28	64.65	55.31
	WRN28-10	93.16	49.41	<u>81.05</u>	51.47
AT-EDM- ℓ_∞	WRN28-10	93.36	70.90	69.73	72.96
	WRN28-10	95.90	53.32	84.77	55.34
DiffPure (t_1)	UNet+WRN70-16	87.50	40.62	75.59	47.89
	UNet+WRN70-16	90.97	44.53	72.65	51.89
	LM	87.89	<u>71.68</u>	75.00	65.22
	CLIPure-Diff	93.75	55.74	80.02	58.24
	CLIPure-Cos	84.38	64.21	65.94	66.41
MultiDAP (“a photo of a .”)	UNet+WRN70-16	<u>93.80</u>	70.29	74.53	65.00
MultiDAP (prompt learning)	UNet+WRN70-16	<u>94.12</u>	72.38	76.05	<u>68.21</u>

397
 398 we inject adversarial examples into the forward diffusion process with a single randomly sampled
 399 timestep $t \sim \text{Unif}([T_1, T_2])$ and Gaussian noise $\epsilon \sim \mathcal{N}(0, I)$, where we set $T_1 = 400$ and $T_2 = 600$.
 400 We then run gradient descent on x_0 using the Regularized DDPM loss (Eq. (3)), which provides an
 401 efficient surrogate for likelihood maximization. The step size η is set to 0.2 for CIFAR datasets and
 402 [0.5 for ImageNet](#), and we clip pixel values into $[0, 1]$ after each update. Unless otherwise specified,
 403 we adopt the class-agnostic prompt p^* learned in Sec. 3.2 as the conditioning signal.
 404

405 **Adversarial Attacks.** We evaluate robustness against two widely used adversarial attacks for
 406 diffusion-based adversarial purification. The first is PGD-20 ([Madry, 2017](#)), a multi-step iterative
 407 ℓ_∞ -bounded attack with random restarts. The second is AutoAttack ([Croce and Hein, 2020](#)), a
 408 parameter-free ensemble of diverse attacks. Because of the stochasticity in MultiDAP, we use
 409 rand version of AutoAttacks. Following [Chen et al. \(2024a\)](#)’s settings, we have n_iter = 100 for
 410 AutoAttack. Unless otherwise specified, the maximum perturbation budget is set to $\epsilon = 8/255$ for
 411 ℓ_∞ threat models and $\epsilon = 0.5$ for ℓ_2 threat models.
 412

413 **Baselines.** We compare our method with two representative purification-based defenses. The
 414 first is DiffPure ([Nie et al., 2022](#)), which performs iterative reverse diffusion to remove adversarial
 415 perturbations. The second is Likelihood Maximization (LM) ([Chen et al., 2024b](#)) which formulates
 416 purification as direct maximization of the unimodal diffusion model likelihood. The third is CLIPure
 417 [Zhang et al. \(2025\)](#), a CLIP-based zero-shot purification method with two variants: CLIPure-Diff,
 418 which estimates image likelihood via the generative latent process, and CLIPure-Cos, which measures
 419 likelihood using the cosine similarity between image embeddings and a blank template. These
 420 methods provide strong and conceptually related baselines for evaluating the effectiveness of our
 421 proposed MultiDAP. Other baseline methods are adversarial training methods which rely on unimodal
 422 diffusion models to generate adversarial examples or argument data.
 423

4.2 ROBUSTNESS AGAINST ADVERSARIAL ATTACKS

425 Table 1 summarizes the robustness results of our method on CIFAR-10 compared with DiffPure
 426 and Likelihood Maximization under PGD-20 and AutoAttack. Our proposed MultiDAP achieves
 427 substantially better performance than Likelihood Maximization across both ℓ_∞ and ℓ_2 threat models,
 428 highlighting the benefit of leveraging a learned prompt prior. Compared with DiffPure, MultiDAP
 429 attains slightly higher robust accuracy, but requires only 5 purification steps, whereas DiffPure
 430 typically involves dozens of Langevin iterations. This efficiency gap makes MultiDAP significantly
 431 more practical in scenarios where the inference speed is critical, while still preserving competitive
 432 robustness. [We leave the results for CIFAR-100 in Appendix A.4](#).

432 Table 2: Clean and robust accuracy (%) on ImageNet-1K. Robust results under AutoAttack are
 433 reported for ℓ_∞ ($\epsilon = 4/255$) on 512 randomly selected test samples. Our method (MultiDAP) uses a
 434 class-agnostic prompt and 20 purification steps. We also report an ablation using the fixed template
 435 prompt “a photo of a”. **Bold** denotes the best, underline the second best, and shading the third best.
 436

Method	Architecture	Clean Acc	AA
AT (Engstrom et al., 2019)	ResNet-50	62.56	31.06
Fast AT (Wong et al., 2020)	ResNet-50	55.62	26.95
Strong PGD AT (Salman et al., 2020)	ResNet-50	<u>64.02</u>	37.89
AT (Bai et al., 2021)	ResNet-50	<u>67.38</u>	35.51
DiffPure ($t = 0.15$)	UNet+ResNet-50	67.79	<u>40.93</u>
LM	UNet+ResNet-50	66.51	<u>38.27</u>
MultiDAP (“a photo of a .”)	UNet+ResNet-50	63.78	38.24
MultiDAP (prompt learning)	UNet+ResNet-50	63.80	41.12

447 Table 3: Computational cost comparison under the same perturbation budget using AutoAttack.
 448 FLOPs and inference time are measured per image on an NVIDIA A100. “FLOPs (G)” represents
 449 the total number of floating-point operations required for full pipeline. “Time (ms)” reports the
 450 end-to-end inference latency measured at batch size 1. “Prompt (s)” indicates prompt-learning cost.
 451 The FLOPs difference of Diffpure between CIFAR-10 and CIFAR-100 arises from its use of different
 452 backbone classifiers (WideResNet-70-16 for CIFAR-10 and WideResNet-28-10 for CIFAR-100).
 453

Method	CIFAR-10			CIFAR-100		
	FLOPs (G)	Time (ms)	Prompt (s)	FLOPs (G)	Time (ms)	Prompt (s)
Diffpure	212.4	1567.40	–	34.87	1520.38	–
CLIPure	51.9	410.50	–	12.16	402.74	–
LM	43.9	385.26	–	7.25	376.31	–
MultiDAP	38.8	371.48	5800	5.25	361.95	43200

463 Table 2 demonstrates the scalability of our method on the large-scale ImageNet-1K dataset. Consistent
 464 with the results on CIFAR benchmarks, MultiDAP achieves state-of-the-art robustness while main-
 465 taining high efficiency with only 20 purification steps. Our prompt-learning variant attains 41.12%
 466 robust accuracy under AutoAttack, significantly outperforming strong adversarial training baselines
 467 such as Strong-AT (Salman et al., 2020) (37.89%) and matching the heavy-computation DiffPure
 468 (Nie et al., 2022) (40.93%). This confirms that leveraging rich cross-modal semantic information is a
 469 key factor for effective purification on high-resolution, diverse natural images.

470 While prompt engineering relies on manually crafted textual templates (e.g., “a photo of a .”), our
 471 approach learns task-adaptive prompt embeddings that encode richer semantic priors. This learned
 472 representation provides stronger and more flexible conditioning, enabling the diffusion model to better
 473 suppress adversarial noise and reconstruct class-consistent images even under severe perturbations.
 474 Compared to hand-designed prompts, prompt learning thus offers a systematic and scalable way to
 475 inject semantic guidance into the purification process.

476 Prompt learning on CIFAR in MultiDAP is lightweight: only 16×768 learnable tokens ($< 0.02\%$
 477 of Stable Diffusion’s parameters) introduced and converges within 2 hours on 1 A100 GPU for
 478 CIFAR-10. As in Table 3, MultiDAP is the most computationally efficient method across both
 479 CIFAR-10 and CIFAR-100. It achieves the lowest FLOPs and fastest inference latency, i.e., 38.8 G
 480 FLOPs and 371.48 ms on CIFAR-10; 5.25 G FLOPs and 361.95 ms on CIFAR-100, while maintaining
 481 competitive robustness. Although MultiDAP includes a one-time offline prompt-learning stage, its
 482 test-time cost remains substantially lower than DiffPure and CLIPure, highlighting its practicality for
 483 real-world deployment.

484 From a theoretical perspective, the robustness of MultiDAP can be explained by the stochastic nature
 485 of purification combined with the learned semantic prior. Instead of explicitly computing the expec-
 486 tation over all timesteps and noises, which would be computationally prohibitive, we approximate

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

Table 4: Ablation studies on CIFAR-10 under AutoAttack (ℓ_∞ -norm $\epsilon = 8/255$, n_iter = 100). Left: effect of purification steps N with fixed $[T_1, T_2] = (400, 600)$ and $\eta = 0.2$. Right: effect of timestep range $[T_1, T_2]$ with fixed $N = 5$ and $\eta = 0.2$.

N	η	Clean Acc	Robust Acc	$[T_1, T_2]$	N	η	Clean Acc	Robust Acc
1	0.2	95.01	62.31	[100, 200]	5	0.2	96.12	58.44
2	0.2	94.57	65.24	[200, 400]	5	0.2	95.31	65.32
5	0.2	94.12	72.38	[400, 600]	5	0.2	94.12	72.38
10	0.2	91.31	71.88	[600, 800]	5	0.2	92.25	66.41
20	0.2	92.88	65.12	[800, 999]	5	0.2	90.14	60.27
30	0.2	88.75	58.44					
50	0.2	85.94	53.75					

it with a single random (t, ϵ) per iteration. This stochastic approximation, when combined with semantic guidance from prompts, provides both efficiency and regularization, preventing overfitting to specific perturbations. In practice, this explains why only 5–10 purification steps are sufficient to achieve competitive robustness, which is consistent with our ablation results showing that excessive steps may even degrade performance (see Table 4).

4.3 ABLATION STUDIES

Number of Purification Steps. We first study the effect of the number of purification steps N while keeping the timestep range fixed at $[400, 600]$ and step size $\eta = 0.2$. As shown in Table 4 (left), increasing N from 1 to 5 steadily improves robustness, with the best performance observed at $N = 5$. Using 10 steps achieves comparable accuracy, but going beyond 10 steps (e.g., 20, 30, 50) leads to diminishing or even negative returns. In particular, excessive purification may overfit the injected noise, causing a degradation in both clean and robust accuracy. This observation highlights that a small number of purification steps (around 5–10) is sufficient for effective adversarial denoising, striking a good balance between robustness and efficiency.

Choice of Timestep Range. We next analyze the effect of varying the timestep interval $[T_1, T_2]$ during purification, while fixing the purification steps $N = 5$ and step size $\eta = 0.2$. As shown in Table 4(right), choosing too small a range (e.g., $[100, 200]$) fails to inject sufficient noise, causing the purification process to underperform in terms of robustness. Conversely, excessively large ranges (e.g., $[800, 999]$) introduce too much noise, which harms clean accuracy due to over-smoothed reconstructions. The best performance is observed in the mid-range (e.g., $[400, 600]$), which strikes a balance between removing adversarial perturbations and preserving semantic fidelity. This finding highlights that an appropriate noise level is crucial for effective adversarial denoising.

We provide additional ablation studies on the the prompt length M and regularization strength λ in the Appendix A.1 and Appendix A.2, respectively.

5 CONCLUSIONS AND DISCUSSIONS

We proposed Multimodal Diffusion for Adversarial Purification (MultiDAP), a novel adversarial defense that leverages text-to-image diffusion models guided by learnable prompts. Experiments on CIFAR-10 demonstrate that MultiDAP achieves competitive robustness against strong white-box attacks, outperforming likelihood maximization baselines while being significantly more efficient than DiffPure. Our ablation studies further highlight the importance of prompt learning and careful timestep design, showing that semantic priors play a crucial role in improving purification quality.

Limitations and Future Work. Despite promising results, MultiDAP currently depends on Stable Diffusion backbones, which are computationally heavier than standard feed-forward models. Future work may explore lightweight diffusion architectures or truncated sampling to further improve efficiency. Moreover, extending MultiDAP to large-scale datasets (e.g., ImageNet) and domain-shifted benchmarks would provide stronger evidence of generalization. Another exciting direction is to incorporate richer or multi-modal prompts (e.g., textual descriptions beyond class names) to enhance semantic guidance during purification.

540 ETHICS STATEMENT
541542 This research employs computational approaches exclusively with publicly accessible datasets,
543 avoiding any human subject involvement or confidential data handling. We adhere to ICLR's ethical
544 guidelines without competing interests, emphasizing responsible deployment of our contributions
545 while ensuring transparent reporting to support reproducible research practices.
546547 REPRODUCIBILITY STATEMENT
548549 Our experiments utilize publicly available datasets with detailed experimental configurations provided
550 throughout the paper. To facilitate reproducibility, we commit to releasing the complete source code
551 upon paper acceptance, enabling the research community to validate and build upon our findings.
552553 REFERENCES
554555 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
556 examples. *arXiv preprint arXiv:1412.6572*, 2014.558 Changhao Shi, Chester Holtz, and Gal Mishne. Online adversarial purification based on self-
559 supervision. *arXiv preprint arXiv:2101.09387*, 2021.560 Jongmin Yoon, Sung Ju Hwang, and Juho Lee. Adversarial purification with score-based generative
561 models. In *International Conference on Machine Learning*, pages 12062–12072. PMLR, 2021.563 Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Animashree Anandkumar.
564 Diffusion models for adversarial purification. In *Proceedings of the 39th International Conference
565 on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 16805–
566 16827, 2022.567 Jinyi Wang, Zhaoyang Lyu, Dahua Lin, Bo Dai, and Hongfei Fu. Guided diffusion model for
568 adversarial purification. *arXiv*, 2205.14969, 2022a.569 Mingyuan Bai, Wei Huang, Tenghui Li, Andong Wang, Junbin Gao, César Federico Caiafa, and
570 Qibin Zhao. Diffusion models demand contrastive guidance for adversarial purification to advance.
571 2024.573 Chun Tong Lei, Hon Ming Yam, Zhongliang Guo, Yifei Qian, and Chun Pong Lau. Instant adversarial
574 purification with adversarial consistency distillation. In *Proceedings of the Computer Vision and
575 Pattern Recognition Conference*, pages 24331–24340, 2025.576 Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of
577 diverse parameter-free attacks. In *International conference on machine learning*, pages 2206–2216.
578 PMLR, 2020.580 Cassidy Laidlaw, Sahil Singla, and Soheil Feizi. Perceptual adversarial robustness: Defense against
581 unseen threat models. In *International Conference on Learning Representations*, 2021.582 Hadi M Dolatabadi, Sarah Erfani, and Christopher Leckie. ℓ_∞ -robustness and beyond: Unleashing
583 efficient adversarial training. In *European Conference on Computer Vision*, pages 467–483, 2022.585 Zekai Wang, Tianyu Pang, Chao Du, Min Lin, Weiwei Liu, and Shuicheng Yan. Better diffusion
586 models further improve adversarial training. In *International Conference on Machine Learning*,
587 pages 36246–36263. PMLR, 2023a.588 Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend:
589 Leveraging generative models to understand and defend against adversarial examples. *arXiv
590 preprint arXiv:1710.10766*, 2017.592 Huanran Chen, Yinpeng Dong, Zhengyi Wang, Xiao Yang, Chengqi Duan, Hang Su, and Jun Zhu.
593 Robust classification via a single diffusion model. In *Proceedings of the 41st International
Conference on Machine Learning*, ICML'24. JMLR.org, 2024a.

594 Mingkun Zhang, Keping Bi, Wei Chen, Jiafeng Guo, and Xueqi Cheng. CLIPure: Purification in latent
 595 space via CLIP for adversarially robust zero-shot classification. In *The Thirteenth International*
 596 *Conference on Learning Representations*, 2025.

597

598 Yiwei Zhou, Xiaobo Xia, Zhiwei Lin, Bo Han, and Tongliang Liu. Few-shot adversarial prompt learn-
 599 ing on vision-language models. In *The Thirty-eighth Annual Conference on Neural Information*
 600 *Processing Systems*, 2024.

601 Zekai Wang, Tianyu Pang, Chao Du, Min Lin, Weiwei Liu, and Shuicheng Yan. Better diffusion
 602 models further improve adversarial training. In *International conference on machine learning*,
 603 pages 36246–36263. PMLR, 2023b.

604

605 Minjong Lee and Dongwoo Kim. Robust evaluation of diffusion-based adversarial purification. In
 606 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 134–144,
 607 2023.

608 Chaowei Xiao, Zhongzhu Chen, Kun Jin, Jiong Xiao Wang, Weili Nie, Mingyan Liu, Anima Anand-
 609 kumar, Bo Li, and Dawn Song. Densepure: Understanding diffusion models towards adversarial
 610 robustness. *arXiv preprint arXiv:2211.00322*, 2022.

611 Cheng-Han Yeh, Kuanchun Yu, and Chun-Shien Lu. Test-time adversarial defense with opposite
 612 adversarial path and high attack time cost. *arXiv preprint arXiv:2410.16805*, 2024.

613

614 Roland S Zimmermann, Lukas Schott, Yang Song, Benjamin A Dunn, and David A Klindt. Score-
 615 based generative classifiers. *arXiv preprint arXiv:2110.00473*, 2021.

616 Kevin Clark and Priyank Jaini. Text-to-image diffusion models are zero shot classifiers. *Advances in*
 617 *Neural Information Processing Systems*, 36, 2023.

618

619 Huanran Chen, Yinpeng Dong, Shitong Shao, Zhongkai Hao, Xiao Yang, Hang Su, and Jun Zhu.
 620 Diffusion models are certifiably robust classifiers. In *The Thirty-eighth Annual Conference on*
 621 *Neural Information Processing Systems*, 2024b.

622 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
 623 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
 624 models from natural language supervision. In *International conference on machine learning*, pages
 625 8748–8763. PMLR, 2021.

626

627 Christian Schlarbmann, Naman Deep Singh, Francesco Croce, and Matthias Hein. Robust clip:
 628 Unsupervised adversarial fine-tuning of vision embeddings for robust large vision-language models.
 629 *arXiv preprint arXiv:2402.12336*, 2024.

630

631 Jiaming Zhang, Xingjun Ma, Xin Wang, Lingyu Qiu, Jiaqi Wang, Yu-Gang Jiang, and Jitao Sang.
 632 Adversarial prompt tuning for vision-language models. In *European conference on computer*
 633 *vision*, pages 56–72. Springer, 2024.

634

635 Lin Li, Haoyan Guan, Jianing Qiu, and Michael Spratling. One prompt word is enough to boost
 636 adversarial robustness for pre-trained vision-language models. In *Proceedings of the IEEE/CVF*
 637 *Conference on Computer Vision and Pattern Recognition*, pages 24408–24419, 2024.

638

639 Lijun Sheng, Jian Liang, Zilei Wang, and Ran He. R-pt: Improving adversarial robustness of
 640 vision-language models through test-time prompt tuning. In *Proceedings of the Computer Vision*
 641 *and Pattern Recognition Conference*, pages 29958–29967, 2025.

642

643 Zeyu Wang, Cihang Xie, Brian Bartoldson, and Bhavya Kailkhura. Double visual defense: Adver-
 644 sarial pre-training and instruction tuning for improving vision-language model robustness. *arXiv*
 645 *preprint arXiv:2501.09446*, 2025.

646

647 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
 648 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
 649 *arXiv:2011.13456*, 2020.

650

651 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
 652 *neural information processing systems*, 33:6840–6851, 2020.

648 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
 649 In *International conference on machine learning*, pages 8162–8171. PMLR, 2021.
 650

651 Jinyi Wang, Zhaoyang Lyu, Dahua Lin, Bo Dai, and Hongfei Fu. Guided diffusion model for
 652 adversarial purification. *arXiv preprint arXiv:2205.14969*, 2022b.

653 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
 654 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Confer-
 655 ence on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June 2022.

656

657 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances
 658 in neural information processing systems*, 34:8780–8794, 2021.

659 Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to Prompt for Vision-
 660 Language Models. *International Journal of Computer Vision*, 130(9):2337–2348, September 2022.
 661 ISSN 1573-1405.

662 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
 663 2014.

664 Lambda Labs. miniSD-diffusers: A Text-to-Image Model based on Stable Diffusion, 2022. URL
 666 <https://huggingface.co/lambdalabs/miniSD-diffusers>.

667 Aleksander Madry. Towards deep learning models resistant to adversarial attacks. *arXiv preprint
 668 arXiv:1706.06083*, 2017.

669 Logan Engstrom, Andrew Ilyas, Hadi Salman, Shibani Santurkar, and Dimitris Tsipras. Robustness
 670 (python library), 2019. URL <https://github.com/MadryLab/robustness>, 4(4):4–3, 2019.

671

672 Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial training.
 673 *arXiv preprint arXiv:2001.03994*, 2020.

674

675 Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversari-
 676 ally robust imagenet models transfer better? *Advances in Neural Information Processing Systems*,
 677 33:3533–3545, 2020.

678

679 Yutong Bai, Jieru Mei, Alan L Yuille, and Cihang Xie. Are transformers more robust than cnns?
 680 *Advances in neural information processing systems*, 34:26831–26843, 2021.

681

682 Sylvestre-Alvise Rebuffi, Sven Gowal, Dan A Calian, Florian Stimberg, Olivia Wiles, and Tim-
 683 othy Mann. Fixing data augmentation to improve adversarial robustness. *arXiv preprint
 684 arXiv:2103.01946*, 2021.

685

686 Tianyu Pang, Min Lin, Xiao Yang, Jun Zhu, and Shuicheng Yan. Robustness and accuracy could
 687 be reconcilable by (proper) definition. In *International conference on machine learning*, pages
 688 17258–17277. PMLR, 2022.

689

690

691

692

693

694

695

696

697

698

699

700

701

702 Table 5: Ablation studies on CIFAR-10 under AutoAttack (ℓ_∞ -norm $\epsilon = 8/255$). Left: effect
 703 of prompt length M with fixed regularization strength $\lambda = 0.9$ and $\eta = 0.2$. Right: effect of
 704 regularization strength λ with fixed prompt length $M = 16$ and $\eta = 0.2$.

Prompt length	Clean Acc	AA (ℓ_∞)	Regulation strength	Clean Acc	AA (ℓ_∞)
4	90.87	71.19	0.0	90.91	70.17
8	92.31	72.26	0.3	92.56	71.25
12	91.56	72.25	0.6	93.72	72.34
16	93.36	72.38	0.9	93.75	72.38
32	93.75	72.04	1.2	94.38	71.01
			1.5	95.10	69.55

714 USE OF LARGE LANGUAGE MODELS (LLMs)

715 LLMs served as supplementary instruments for textual improvement and code standardization throughout this research. Their application was limited to enhancing readability, correcting grammatical
 716 inconsistencies, and maintaining uniform presentation standards for pseudocode blocks and LaTeX
 717 expressions. The core intellectual contributions—including research conception, methodological
 718 frameworks, analytical reasoning, and experimental findings—remain entirely human-generated.

722 A ADDITIONAL EXPERIMENTS

724 A.1 ABLATION STUDY ON PROMPT LENGTH

726 To further understand the effect of prompt capacity, we conduct an ablation study by varying the
 727 prompt length M while fixing the regularization strength $\lambda = 0.9$ and step size $\eta = 0.2$. The results
 728 are shown in Table 5 (left). We observe that both clean accuracy and robustness generally improve
 729 as M increases from 4 to 16, suggesting that a moderate number of context tokens is beneficial for
 730 encoding richer semantic priors that guide the purification dynamics. The best robustness (72.38%
 731 AutoAttack under ℓ_∞) is achieved at $M = 16$. Increasing the prompt length further to 32 does not
 732 provide additional gains and even slightly reduces robustness, likely because overly long prompts
 733 introduce redundancy that makes optimization more difficult. These findings indicate that a compact
 734 yet expressive prompt length is sufficient for effective semantic conditioning.

736 A.2 ABLATION STUDY ON REGULARIZATION STRENGTH

738 We further investigate the impact of the regularization strength λ while fixing the prompt length $M =$
 739 16 and $\eta = 0.2$. The results in Table 5 (right) show that robustness improves steadily as λ increases
 740 from 0.0 to 1.5, with the best performance (72.38% AutoAttack under ℓ_∞) obtained at $\lambda = 0.9$.
 741 This trend suggests that regularizing the prompt embeddings is important for preventing overfitting
 742 to individual noise realizations, thereby stabilizing the purification process. However, very large
 743 regularization values (e.g., $\lambda = 1.5$) begin to degrade robustness, indicating that excessive constraint
 744 may limit the representational flexibility. Overall, these results demonstrate that an intermediate
 745 regularization level provides the best trade-off between stability and semantic expressiveness.

746 A.3 EMPIRICAL ESTIMATION OF LIPSCHITZ CONSTANT AND GRADIENT VARIANCE

748 To verify the practical assumptions required by Theorem 2, we conduct an empirical evaluation of the
 749 local Lipschitz constant L and gradient variance σ of the purification objective $g(x) = \nabla_x L_{\text{pur}}(x)$,
 750 where L_{pur} is the DDPM-based denoising loss used in our purification step. Both measurements
 751 strictly follow the forward process in our method, including VAE encoding, DDPM noise injection,
 752 and UNet denoising.

753 First, we estimate the local Lipschitz constant via a finite-difference approximation: $L(x) \approx$
 $\frac{\|g(x+\delta) - g(x)\|_2}{\|\delta\|_2}$, where δ is a small random perturbation ($\|\delta\|_2 \approx 10^{-3}$). We repeat this proce-
 754 dure for multiple random perturbations and multiple images. To estimate σ , we sample gradients
 755

under random diffusion timesteps $t \sim \mathcal{U}[t_{\text{st}}, t_{\text{ed}}]$ and random noise $\epsilon \sim \mathcal{N}(0, I)$: $\sigma^2 = \text{Var}_{t, \epsilon}[g(x)]$. All gradients are computed in pixel space to match the theoretical setting.

Across CIFAR-10 samples, we obtain the empirical statistics shown in Table 6. These values are small, stable, and tightly bounded, confirming that the assumptions in Theorem 2 hold in practice. In particular, the purification gradient exhibits Lipschitz-smooth behavior and extremely low stochastic variance. This explains why a small number of purification steps (five iterations) is numerically stable and effective in our method.

Metric	Mean	Max / Std
Lipschitz constant L	0.1724	0.3000
Gradient std. σ	1.05×10^{-4}	—

Table 6: Empirical estimates of the Lipschitz constant and gradient variance of the purification loss. Small and stable values indicate smooth gradient behavior and validate the assumptions in Theorem 2.

A.4 CIFAR-100

Table 7 reports results on the more challenging CIFAR-100 dataset. Consistent with the CIFAR-10 findings, MultiDAP achieves strong robustness while using only 5 purification steps. Our prompt-learning variant obtains 38.15% AutoAttack robustness, outperforming all diffusion-based baselines, including DiffPure (Nie et al., 2022) and Diffusion+Contrastive (Bai et al., 2024). Notably, even the fixed template prompt (“a photo of a”) already provides competitive performance (35.70% AA), demonstrating that semantic conditioning—whether learned or hand-designed—substantially benefits the purification dynamics. These results indicate that MultiDAP generalizes effectively to datasets with larger label spaces and more fine-grained visual variability.

B COMPLETE PROOFS

B.1 PROOF OF THEOREM 1

Proof of Theorem 1. For DDPM (or latent DDPM), the negative ELBO can be written (Ho et al., 2020; Nichol and Dhariwal, 2021) as

$$-\log \underline{p}_\theta(z_0 \mid p) = C(\theta) + \sum_{t=1}^T \underbrace{\mathbb{E}[\text{KL}(q(z_{t-1} \mid z_t, z_0) \parallel p_\theta(z_{t-1} \mid z_t, p))]}_{=: \mathcal{R}_t(p)} + \underbrace{\mathbb{E}[\text{KL}(q(z_T) \parallel p(z_T))]}_{\text{independent of } p}, \quad (4)$$

where $C(\theta)$ is independent of p and only $\mathcal{R}_t(p)$ depends on the prompt.

Using the standard mean parameterization,

$$\mu_\theta(x_t, t, p) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t, p) \right), \quad \Sigma_\theta(x_t, t) = \sigma_t^2 I,$$

one can show that $\mathcal{R}_t(p)$ is equivalent to a weighted noise-prediction MSE:

$$\mathcal{R}_t(p) = w_t \mathbb{E}_{x_0, t, \epsilon} \left[\left\| \epsilon - \epsilon_\theta(x_t, t, p) \right\|_2^2 \right] + \text{const}, \quad w_t = \frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)} > 0.$$

Substituting into Eq. (4) yields

$$\mathcal{L}_{\text{VLB}}(p) =: -\log \underline{p}_\theta(x_0 \mid p) = C'(\theta) + \sum_{t=1}^T w_t \mathbb{E}_{x_0, t, \epsilon} \left[\left\| \epsilon - \epsilon_\theta(x_t, t, p) \right\|_2^2 \right],$$

which differs from $\mathcal{L}_{\text{prompt}}(p)$ only by positive weights and a constant. Therefore,

$$\arg \min_p \mathcal{L}_{\text{prompt}}(p) = \arg \min_p \mathcal{L}_{\text{VLB}}(p).$$

Let p^* be the minimizer; then $\mathcal{L}_{\text{VLB}}(p^*) \leq \mathcal{L}_{\text{VLB}}(p)$ for all p , equivalently $\log \underline{p}_\theta(x_0 \mid p^*) \geq \log \underline{p}_\theta(x_0 \mid p)$. \square

Table 7: Clean and robust accuracy (%) on CIFAR-100. Robust results under AutoAttack are reported for ℓ_∞ ($\epsilon = 8/255$). Our method (MultiDAP) uses a class-agnostic prompt and 5 purification steps. We also report an ablation using the fixed template prompt “a photo of a”. **Bold** denotes the best, underline the second best, and  shading the third best.

Method	Architecture	Clean Acc	AA
AT-CutMix (Rebuffi et al., 2021)	WRN28-10	62.97	29.80
AT-DDPM (Rebuffi et al., 2021)	WRN28-10	59.18	30.81
AT-DDPM + CutMix (Rebuffi et al., 2021)	WRN28-10	62.41	32.06
AT-DDPM (Pang et al., 2022)	WRN28-10	62.08	31.40
AT-DDPM (Wang et al., 2023b)	WRN28-10	<u>72.58</u>	38.83
DiffPure ($t = 0.1$) (Nie et al., 2022)	UNet+WRN28-10	62.50	8.60
Diffusion + Contrastive ($t = 0.1$) (Bai et al., 2024)	UNet+WRN28-10	57.82	24.22
LM (Chen et al., 2024a)	UNet+WRN28-10	66.45	33.83
MultiDAP (“a photo of a .”)	UNet+WRN28-10	73.29	 35.70
MultiDAP (prompt learning)	UNet+WRN28-10	 72.52	<u>38.15</u>

B.2 PROOF OF LEMMA 1

Proof of Lemma 1. We first recall the purification objective:

$$\mathcal{L}_{\text{pur}}(x_0; p^*) = \mathbb{E}_{t, \epsilon} [\ell(x_0; t, \epsilon)] + \lambda \mathcal{R}(x_0, x^{adv}),$$

where $t \sim \text{Unif}(\{1, \dots, T\})$ and $\epsilon \sim \mathcal{N}(0, I)$.

By linearity of expectation and interchangeability of expectation and gradient under standard regularity conditions:

$$\nabla_{x_0} \mathcal{L}_{\text{pur}}(x_0; p^*) = \nabla_{x_0} \mathbb{E}_{t, \epsilon} [\ell(x_0; t, \epsilon)] + \lambda \nabla_{x_0} \mathcal{R}(x_0, x^{adv}).$$

On the other hand, the stochastic gradient $g(x_0; t, \epsilon)$ is defined as

$$g(x_0; t, \epsilon) = \nabla_{x_0} \ell(x_0; t, \epsilon) + \lambda \nabla_{x_0} \mathcal{R}(x_0, x^{adv}).$$

Taking expectation over (t, ϵ) gives:

$$\mathbb{E}_{t, \epsilon} [g(x_0; t, \epsilon)] = \mathbb{E}_{t, \epsilon} [\nabla_{x_0} \ell(x_0; t, \epsilon)] + \lambda \nabla_{x_0} \mathcal{R}(x_0, x^{adv}),$$

which equals $\nabla_{x_0} \mathcal{L}_{\text{pur}}(x_0; p^*)$. Thus the estimator is unbiased.

By Assumption (A2), we assume that the variance of the stochastic gradient is bounded:

$$\mathbb{E}_{t, \epsilon} [\|g(x_0; t, \epsilon) - \nabla_{x_0} \mathcal{L}_{\text{pur}}(x_0; p^*)\|_2^2] \leq \sigma^2.$$

This holds because $\ell(x_0; t, \epsilon)$ is quadratic in ϵ and $\epsilon \sim \mathcal{N}(0, I)$ has bounded second moment, while \mathcal{R} is convex and $L_{\mathcal{R}}$ -smooth (Assumption A3). Therefore the stochastic gradient inherits bounded variance.

Together, we conclude that $g(x_0; t, \epsilon)$ is an unbiased stochastic gradient estimator of $\nabla_{x_0} \mathcal{L}_{\text{pur}}(x_0; p^*)$ with bounded variance, completing the proof. \square

B.3 PROOF OF THEOREM 2

Proof of Theorem 2. Recall the purification objective

$$\mathcal{L}_{\text{pur}}(x_0; p^*) = \mathbb{E}_{t, \epsilon} [\ell(x_0; t, \epsilon)] + \lambda \mathcal{R}(x_0, x^{adv}), \quad \ell(x_0; t, \epsilon) = \|\epsilon - \epsilon_\theta(z_t, t, p^*)\|_2^2,$$

and the projected update on the pixel cube $\mathcal{C} = [0, 1]^d$:

$$x_0^{(k+1)} = \Pi_{\mathcal{C}}(x_0^{(k)} - \eta g(x_0^{(k)}; t_k, \epsilon_k)), \quad g(x_0; t, \epsilon) = \nabla_{x_0} \ell(x_0; t, \epsilon) + \lambda \nabla_{x_0} \mathcal{R}(x_0, x^{adv}).$$

864 Since \mathcal{L}_{pur} is L -smooth on \mathcal{C} (Assumption A1), the standard smoothness inequality gives
 865

$$866 \mathcal{L}_{\text{pur}}(x_0^{(k+1)}; p^*) \leq \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*) + \langle \nabla \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*), x_0^{(k+1)} - x_0^{(k)} \rangle + \frac{L}{2} \|x_0^{(k+1)} - x_0^{(k)}\|_2^2. \quad (5)$$

868 The projection onto a closed convex set is nonexpansive and satisfies $\|x_0^{(k+1)} - x_0^{(k)}\|_2 \leq
 869 \eta \|g(x_0^{(k)}; t_k, \epsilon_k)\|_2$. Hence, from Eq. (5),
 870

$$871 \mathcal{L}_{\text{pur}}(x_0^{(k+1)}; p^*) \leq \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*) - \eta \langle \nabla \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*), g(x_0^{(k)}; t_k, \epsilon_k) \rangle + \frac{L\eta^2}{2} \|g(x_0^{(k)}; t_k, \epsilon_k)\|_2^2.$$

874 Conditioning on $x_0^{(k)}$ and using Lemma 1,
 875

$$876 \mathbb{E}_{t_k, \epsilon_k} [g(x_0^{(k)}; t_k, \epsilon_k) | x_0^{(k)}] = \nabla \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*),$$

$$877 \mathbb{E}_{t_k, \epsilon_k} [\|g(x_0^{(k)}; t_k, \epsilon_k)\|_2^2 | x_0^{(k)}] \leq \|\nabla \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*)\|_2^2 + \sigma^2.$$

879 Using the standard variance decomposition bound $\mathbb{E}\|g\|_2^2 \leq 2\|\nabla \mathcal{L}_{\text{pur}}\|_2^2 + 2\sigma^2$ (or equivalently
 880 absorbing constants into σ^2) and taking full expectation yields
 881

$$882 \mathbb{E}[\mathcal{L}_{\text{pur}}(x_0^{(k+1)}; p^*)] \leq \mathbb{E}[\mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*)] - \eta \mathbb{E}[\|\nabla \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*)\|_2^2] + \frac{L\eta^2}{2} \mathbb{E}[\|g(x_0^{(k)}; t_k, \epsilon_k)\|_2^2]$$

$$884 \leq \mathbb{E}[\mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*)] - \left(\eta - \frac{L\eta^2}{2}\right) \mathbb{E}[\|\nabla \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*)\|_2^2] + \frac{L\eta^2}{2} \sigma^2.$$

887 If $\eta \leq 1/(2L)$, then $\eta - \frac{L\eta^2}{2} \geq \eta/2$. Hence we obtain the claimed one-step descent:
 888

$$889 \mathbb{E}[\mathcal{L}_{\text{pur}}(x_0^{(k+1)}; p^*)] \leq \mathbb{E}[\mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*)] - \frac{\eta}{2} \mathbb{E}[\|\nabla \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*)\|_2^2] + \frac{\eta^2 L}{2} \sigma^2.$$

891 Summing the inequality over $k = 0, \dots, K-1$ and using the lower bound $\mathcal{L}_{\text{pur}}(x_0; p^*) \geq \mathcal{L}_{\text{inf}} :=
 892 \inf_{x \in [0,1]^d} \mathcal{L}_{\text{pur}}(x; p^*)$, we get
 893

$$894 \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla \mathcal{L}_{\text{pur}}(x_0^{(k)}; p^*)\|_2^2] \leq \frac{2(\mathcal{L}_{\text{pur}}(x_0^{(0)}; p^*) - \mathcal{L}_{\text{inf}})}{\eta K} + \eta L \sigma^2.$$

898 Since $\mathcal{L}_{\text{pur}} = \mathcal{L}_{\text{DDPM}} + \lambda \mathcal{R}$ and \mathcal{R} is convex and smooth (Assumption A3), the same derivation
 899 applies when focusing on the data-fidelity part. In particular, using $\|\nabla \mathcal{L}_{\text{DDPM}}(x)\|_2^2 \leq \|\nabla \mathcal{L}_{\text{pur}}(x)\|_2^2$
 900 (up to a constant absorbed into σ^2) yields

$$902 \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla \mathcal{L}_{\text{DDPM}}(x_0^{(k)}; p^*)\|_2^2] \leq \frac{2(\mathcal{L}_{\text{DDPM}}(x_0^{(0)}; p^*) - \mathcal{L}_{\text{inf}})}{\eta K} + \eta L \sigma^2,$$

904 which is the stated bound. □
 905

906
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917