

APCTRL: ADDING CONDITIONAL CONTROL TO DIFFUSION MODELS BY ALTERNATIVE PROJECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Enhancing the versatility of pretrained diffusion models through advanced conditioning techniques is crucial for improving their applicability. We present APCTRL, a novel conditional image generation approach that formulates the latent \mathbf{z}_t at timestep t as the projection $\mathbf{z}_t = \text{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1})$ onto the denoising set \mathfrak{D}_t . For conditional control, APCTRL integrates the condition set \mathfrak{C}_t , defined by a latent control network $\mathcal{A}_\theta(\cdot, \cdot)$. Our method simplifies conditional sampling to recursive projections $\mathbf{z}_t = \text{Proj}_{\mathfrak{J}_t} \circ \text{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1})$, where each projection step integrates both the diffusion and condition priors. By employing Alternative Projection, our approach offers several key advantages: 1. Multi-Condition Generation: easily expandable with additional conditional sets; 2. Model and Sampling Agnosticism: works with any model or sampling method; 3. Unified Control Loss: simplifies the management of diverse control applications; 4. Efficiency: delivers comparable control with reduced training and sampling times. Extensive experiments demonstrate the superior performance of our method.

1 INTRODUCTION

Unconditional diffusion models, first introduced by Ho et al. (2020), laid the foundation for generative image modeling. Characterized by the latent sequence $\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_T$, they have significantly advanced the generation of high-fidelity images (Yang et al., 2023a). In this sequence, \mathbf{z}_t represents progressively noisier data samples for $t \in (0, T]$ and \mathbf{z}_0 corresponds to the true data samples. The forward process introduces noise gradually, transitioning from \mathbf{z}_{t-1} to \mathbf{z}_t according to the distribution $q(\mathbf{z}_t | \mathbf{z}_{t-1}) := \mathcal{N}(\mathbf{z}_t | \sqrt{\alpha_t} \mathbf{z}_{t-1}, (1 - \alpha_t) \mathbf{I})$, where α_t is a constant hyperparameter. The objective of diffusion models is to generate a sample \mathbf{z}_0 from the data distribution $p(\mathbf{z}_0)$, which can be formulated as an optimization problem: $\text{argmax}_{\mathbf{z}_0} \log p(\mathbf{z}_0)$, seeking the optimal \mathbf{z}_0 that maximizes $p(\mathbf{z}_0)$.

The distribution $p(\mathbf{z}_0)$ is not directly accessible. In the reverse process, the diffusion models offer an approximation through the marginal distribution $p_\theta(\mathbf{z}_0)$. The model parameters θ are optimized using the Evidence Lower Bound (ELBO), which serves as a lower bound for $\log p_\theta(\mathbf{z}_0)$. Specifically, we have: $\log p_\theta(\mathbf{z}_0) \geq \mathbb{E}_{q(\mathbf{z}_{1:T} | \mathbf{z}_0)} \left[\log \frac{p_\theta(\mathbf{z}_{1:T})}{q(\mathbf{z}_{1:T} | \mathbf{z}_0)} \right]$. The right ELBO term can be further expanded as follows: $\mathbb{E}_{q(\mathbf{z}_1 | \mathbf{z}_0)} [\log p_\theta(\mathbf{z}_0 | \mathbf{z}_1)] - D_{\text{KL}}(q(\mathbf{z}_T | \mathbf{z}_0) \| p_\theta(\mathbf{z}_T)) - \sum_{t>1} \mathbb{E}_{q(\mathbf{z}_t | \mathbf{z}_0)} [D_{\text{KL}}(q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{z}_0) \| p_\theta(\mathbf{z}_{t-1} | \mathbf{z}_t))]$. Thus, the goal of the diffusion model becomes to maximize the reverse transition distribution $p_\theta(\mathbf{z}_t | \mathbf{z}_{t+1})$, which in turn maximizes $\log p_\theta(\mathbf{z}_0)$. Consequently, sampling from the reverse transition distribution $p_\theta(\mathbf{z}_t | \mathbf{z}_{t+1})$ can then be expressed as Equation (1), where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\mathcal{S}_\theta(\mathbf{z}_t, t)$ is the neural network designed to predict the score function $\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)$ (Song et al., 2020b).

$$\mathbf{z}_t = \frac{1}{\sqrt{\alpha_{t+1}}} \mathbf{z}_{t+1} + \frac{(1 - \alpha_{t+1})}{\sqrt{\alpha_{t+1}}} \mathcal{S}_\theta(\mathbf{z}_{t+1}, t + 1) + \sqrt{\frac{(1 - \alpha_{t+1})(1 - \bar{\alpha}_t)}{1 - \bar{\alpha}_{t+1}}} \epsilon \quad (1)$$

Unconditional diffusion models was further extended by text-conditional models (Rombach et al., 2022; Yang et al., 2024b). However, these models faced the inherent challenge of

Table 1: A Feature-Rich Approach for Conditional Image Generation. APCtrl boasts an array of beneficial features, surpassing the capabilities of previous methods.

Methods	Latent Control	Controlled by Multi-Condition	Backbones Agnosticism	Sampling Agnosticism	Unified Control Loss	Control on Sampling
Control-on-Training						
ControlNet (Zhang et al., 2023)	✓	✓	✗	✓	✗	✗
ControlNet++ (Li et al., 2024b)	✓	✓	✗	✓	✗	✗
T2I-Adapter (Mou et al., 2024)	✓	✓	✗	✓	✗	✗
UniCtrlNet (Zhao et al., 2024)	✓	✓	✗	✓	✗	✗
UniControl (Qin et al., 2023)	✓	✓	✗	✓	✗	✗
GLIGEN (Li et al., 2023)	✓	✗	✗	✓	✗	✗
Control-on-Sampling						
UniGuid (Bansal et al., 2024)	✗	✗	✓	✗	✗	✓
DSG (Yang et al., 2024c)	✗	✗	✓	✗	✗	✓
FreeDoM (Yu et al., 2023)	✗	✗	✓	✗	✗	✓
APCtrl (Ours)	✓	✓	✓	✓	✓	✓

accurately capturing all image details from text descriptions alone. To address this, diffusion models have incorporated additional conditioning signals, such as bounding boxes (Li et al., 2023; Yang et al., 2023b; Zhao et al., 2024), reference images (Li et al., 2024a; Ruiz et al., 2023), and segmentation maps (Zhang et al., 2023; Bansal et al., 2024; Zhao et al., 2024; Qin et al., 2023), offering more granular control over the generated images.

Conditional image generation falls into two camps: methods that integrate control networks, and those that adjust the inference process for direct control. **Control-on-Training** approaches like ControlNets (Zhang et al., 2023) train networks to refine latent spaces and match images to attributes, incurring retraining costs due to feature space inconsistencies. On the other hand, **Control-on-Sampling** techniques, such as Universal Guidance (Bansal et al., 2024), use pre-trained models to guide sampling, offering flexibility without re-training. However, this comes with potential downsides, such as suboptimal gradient estimations that may degrade sampling quality and prolong sampling times.

APCtrl solves these challenges. Let \mathfrak{D}_0 denote the set of natural images, and \mathfrak{D}_t represent the noisy versions, generated by adding noise to \mathfrak{D}_0 , such that $\mathfrak{D}_t := \{\mathbf{z}_t \mid \mathbf{z}_t = \sqrt{\bar{\alpha}_t}\mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \mathbf{z}_0 \in \mathfrak{D}_0\}$ with $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$. The denoising projection $\mathbf{z}_t = \text{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1})$, as defined by Equation (1), maps a noisy point $\mathbf{z}_{t+1} \in \mathfrak{D}_{t+1}$ to a less noisy point $\mathbf{z}_t \in \mathfrak{D}_t$. The diffusion generation process is thus a sequence of such projections. To enhance control, we introduce a condition set \mathfrak{C}_t , which defines points that satisfy specific constraints at each step t . The intersection $\mathfrak{I}_t := \mathfrak{D}_t \cap \mathfrak{C}_t$ identifies points that conform to both \mathfrak{D}_t and \mathfrak{C}_t . By defining the intersection projection $\text{Proj}_{\mathfrak{I}_t}(\cdot)$, conditional generation is redefined as a recursive sequence of projections $\mathbf{z}_t = \text{Proj}_{\mathfrak{I}_t} \circ \text{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1})$, as shown in Algorithm 1.

Our method lies in the condition projection $\text{Proj}_{\mathfrak{C}_t}(\cdot)$, implemented through a latent control network that imposes constraints and calculates projections onto the conditional sets. This approach offers several key advantages: enhanced adaptability to diverse backbones, more precise and efficient synthesis via latent control, and a unified MSE latent control loss applicable to numerous conditions. By applying Alternative Projection with the denoising projection $\text{Proj}_{\mathfrak{D}_t}(\cdot)$ and the condition projection $\text{Proj}_{\mathfrak{C}_t}(\cdot)$, we compose the intersection projection $\text{Proj}_{\mathfrak{I}_t}(\cdot)$. This method outperforms other sampling techniques. A key feature is the straightforward implementation of multi-condition control via projections onto intersections of condition sets. Table 1 presents a detailed comparison, highlighting how APCtrl surpasses previous methods with its array of beneficial features.

2 RELATED WORKS

Alternative Projection is a technique with a long-standing history. It aims to find a point within the intersection of multiple sets through a sequence of successive projections onto each set, and was seminally studied by Von Neumann (1951), and has since been applied in a myriad of contexts (Deutsch, 1992). Numerous variants, such as relaxed projections (Agmon, 1954; Motzkin & Schoenberg, 1954; Gubin et al., 1967; Brègman,

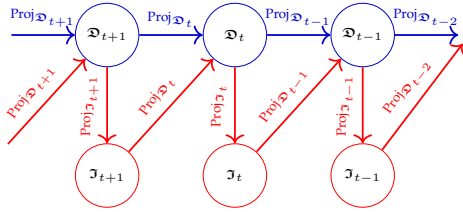


Figure 1: APCtrl Sampling: APCtrl transforms traditional diffusion sampling into a recursive projection, as illustrated in blue in the figure. It integrates easily, needing just one extra line of code, highlighted in red and noted as line 5 in Algorithm 1.

Algorithm 1 APCtrl Sampling

Input: Initial noise \mathbf{z}_T , denoising and condition set \mathfrak{D}_t and \mathfrak{C}_t .

```

1 for  $t = T - 1$  to 1 do
2    $\mathbf{z}_t = \text{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1})$ 
3   if conducting APCtrl Sampling then
4      $\mathfrak{J}_t = \mathfrak{D}_t \cap \mathfrak{C}_t$ 
5      $\mathbf{z}_t = \text{Proj}_{\mathfrak{J}_t}(\mathbf{z}_t)$ , which be recursively computed by
      for  $n=1$  to  $N$  do
         $\mathbf{z}_t = \text{Proj}_{\mathfrak{J}_t}(\mathbf{z}_t)$  in Equation (9)
      end
    end
  end

```

Output: \mathbf{z}_t

1965), inexact projections (Kruger & Thao, 2016), Dykstra’s algorithm (Boyle & Dykstra, 1986), Douglas–Rachford splitting (Douglas & Rachford, 1956; Lions & Mercier, 1979), ADMM (Boyd, 2010), and generalized alternating projections (Fält & Giselsson, 2024), have been proposed.

Diffusion Models (Croitoru et al., 2023; Yang et al., 2023a) constitute a class of models that incrementally introduce noise to data in a controlled manner, with the goal of learning to reverse this process for generating samples. The current research landscape is primarily dominated by three formulations: Denoising Diffusion Probabilistic Models (Ho et al., 2020; Nichol & Dhariwal, 2021; Sohl-Dickstein et al., 2015), Score-Based Generative Model (Song & Ermon, 2019; 2020), and Stochastic Differential Equations (Song et al., 2021; 2020b). The interconnections between them are elucidated by Luo (2022) and Chan (2024).

Control-on-Training takes supplementary networks to modify the latent representations of diffusion models according to particular image conditions. Researchers (Bansal et al., 2023; Nichol et al., 2022; Rombach et al., 2022) have expanded $\mathcal{S}_\theta(\mathbf{z}_t, t)$ in Equation (1) to include both text and image conditions. A notable example of this approach is ControlNet (Zhang et al., 2023), which has become a significant focus within the field. The broader community has contributed to this area by sharing a variety of ControlNets trained across diverse input conditions. Other prominent examples include ControlNet++ (Li et al., 2024b), T2I-Adapter (Mou et al., 2024), UniControlNet (Zhao et al., 2024), UniControl (Qin et al., 2023), GLIGEN (Li et al., 2023), and Ctrl-Adapter (Lin et al., 2024).

Control-on-Sampling utilizes frozen pre-trained models, with modifications to the sampling method to reconstruct an image from a given guidance. Prior work has approached this task with various constraints (Dhariwal & Nichol, 2021; Kawar et al., 2022; Wang et al., 2022; Chung et al., 2023; Lugmayr et al., 2022; Chung et al., 2022; Graikos et al., 2022). For instance, Dhariwal & Nichol (2021) trained a classifier on images of different noise scales to serve as the guidance and incorporated the classifier’s gradients into the sampling process. However, classifiers for noisy images are often domain-specific and not generally available. To address the challenge, several state-of-the-art sampling methods have been introduced, including DSG (Yang et al., 2024c), UniGuidance (Bansal et al., 2024), FreeDoM (Yu et al., 2023), MultiDiffusion (Bar-Tal et al., 2023), and ReSample (Song et al., 2023).

3 APCTRL SAMPLING

Diffusion generation involves the successive application of the projection $\mathbf{z}_t = \text{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1})$, as outlined in Equation (1). To ensure that the denoised point from \mathfrak{D}_t also satisfies the constraint from \mathfrak{C}_t , *i.e.* to maintain \mathbf{z}_t within the intersection $\mathfrak{J}_t = \mathfrak{D}_t \cap \mathfrak{C}_t$, we apply the intersection projection $\text{Proj}_{\mathfrak{J}_t}(\cdot)$ to $\text{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1})$. This results in the iterative formula $\mathbf{z}_t = \text{Proj}_{\mathfrak{J}_t} \circ \text{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1})$, detailed in Algorithm 1 and depicted in Figure 1. Building upon the definition of $\text{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1})$ from Equation 1, this section is dedicated to explaining the use of Alternative Projection to implement the intersection projection $\text{Proj}_{\mathfrak{J}_t}(\cdot)$.

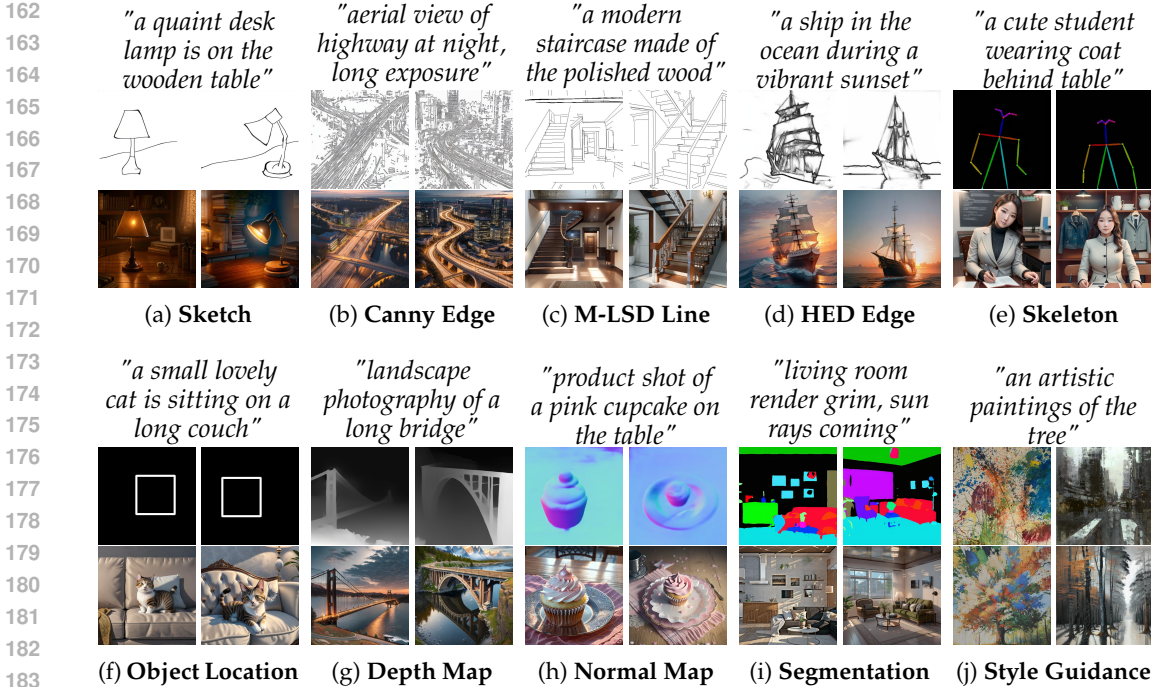


Figure 2: Image Generation with Single Control: APCtrl facilitates the integration of various conditions into diffusion models. Each subfigure is structured with the prompt text in the first row, the conditional image control in the second row, and the resulting controlled generation in the third row.

3.1 ALTERNATIVE PROJECTION

The alternative projection method is an iterative process used to identify a point that belongs to the intersection of two sets, \mathfrak{S}_1 and \mathfrak{S}_2 . Although projecting onto each set separately is easy, projecting directly onto their intersection $\mathfrak{S}_1 \cap \mathfrak{S}_2$ is challenging. Let the projection operators onto \mathfrak{S}_1 and \mathfrak{S}_2 be denoted by $\text{Proj}_{\mathfrak{S}_1}$ and $\text{Proj}_{\mathfrak{S}_2}$, respectively. The alternative projection method is straightforward: beginning with any point, the vector x is iteratively updated by applying the composition of projections, such that $z = \text{Proj}_{\mathfrak{S}_1} \circ \text{Proj}_{\mathfrak{S}_2}(z)$.

3.2 FROM LATENT CONTROL TO LATENT CONTROL

Pixel Control, used in previous Control-on-Sampling methods, computes the controlled intermediate latent code $z_t = \text{argmin}_z \mathcal{L}(I_c, \mathcal{B}(\mathcal{D}(z)))$ with the initial point $\mathcal{Z}(z_{t+1}) := \sqrt{\bar{\alpha}_{t+1}}^{-1}(z_{t+1} + (1 - \bar{\alpha}_{t+1})\mathcal{S}_\theta(z_{t+1}, t + 1))$. In this formulation, $\mathcal{Z}(z_{t+1})$ acts as a denoiser at time step $t + 1$ for the latent variable at time step 0, I_c represents the control image, such as segmentation, depth map, or HED. $\mathcal{D}(\cdot)$ is the decoder of the diffusion model, and $\mathcal{B}(\cdot)$ is the pre-trained condition network, such as networks for segmentation, depth estimation, or HED edge. The metric $\mathcal{L}(\cdot, \cdot)$ can be any loss function, such as MSE for depth or HED images similarity, or Cross-Entropy for segmentation similarity.

Latent Control offers a paradigm shift pixel-level manipulation to operations within the latent space. This shift to a lower-dimensional and more compact latent space, allows for more precise control over image generation. Additionally, it simplifies the optimization process by eliminating the need for a decoder $\mathcal{D}(\cdot)$, thus enhancing efficiency. The method specifically utilizes an encoder $\mathcal{E}(\cdot)$ in conjunction with a latent control network $\mathcal{A}_\theta(\cdot, \cdot)$, to determine the controlled intermediate latent representation z_t .

$$z_t = \text{argmin}_z \|\mathcal{E}(I_c) - \mathcal{A}_\theta(z, t)\|^2 \quad \text{solving with the initial point } \text{Proj}_{\mathfrak{D}_t}(z_{t+1}). \quad (2)$$

For a well-trained model $\mathcal{A}_\theta(\cdot, \cdot)$, the approximation should hold: $\|\mathcal{E}(I_c) - \mathcal{A}_\theta(\mathbf{z}_t, t)\|^2 \approx \mathcal{L}(I_c, \mathcal{B}(\mathcal{D}(\mathbf{z}_{0|t})))$, which allows us to use the latent control $\mathbf{z}_t = \operatorname{argmin}_{\mathbf{z}} \|\mathcal{E}(I_c) - \mathcal{A}_\theta(\mathbf{z}, t)\|^2$ as a substitute for the pixel-level control $\mathbf{z}_t = \operatorname{argmin}_{\mathbf{z}} \mathcal{L}(I_c, \mathcal{B}(\mathcal{D}(\mathbf{z})))$. This applies to various types of control, across different control types such as segmentation guidance, depth map guidance and HED edge guidance.

Training the Latent Control Network focuses on refining the operator $\mathcal{A}_\theta(\cdot, \cdot)$. This operator is built upon the U-Net architecture of the stable diffusion model, with initialization from the SDv1.5 checkpoint. The training process is conducted as Equation (3). During optimization, two primary objectives are achieved: image denoising and feature translation. Denoising improves the latent representation \mathbf{z}_t by reducing noise, thereby enhancing data clarity and adherence to constraints. Meanwhile, feature translation converts denoised image features into control-relevant features, which are essential for specific improvements. The efficacy of diffusion models in image translation has been demonstrated in previous work (Parmar et al., 2024). For further details, please refer to Appendix A.

$$\min_{\theta} \|\mathcal{E}(I_c) - \mathcal{A}_\theta(\mathbf{z}_t, t)\|^2 \quad (3)$$

Feasible Sets encompass all points fulfilling specific criteria. APCtrl involves two kinds of feasible sets: the denoising set \mathfrak{D}_t and the condition set \mathfrak{C}_t . Let \mathfrak{D}_0 denote the set of natural images, the denoising set at time step t can be expressed as

$$\mathfrak{D}_t = \{\mathbf{z}_t \mid \mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \mathbf{z}_0 \in \mathfrak{D}_0\} \quad (4)$$

with the inclusion relation $\mathfrak{D}_t \subseteq \mathfrak{D}_{t+1}$. Upon the successful training of $\mathcal{A}_\theta(\mathbf{z}, t)$, for any point \mathbf{z} within the feasible set, the loss $\|\mathcal{A}_\theta(\mathbf{z}, t) - \mathcal{E}(I_c)\|^2$ is expected to be minimal. Thus, with δ as a predefined threshold, the condition set at time step t can be formulated as

$$\mathfrak{C}_t = \{\mathbf{z} \mid \|\mathcal{E}(I_c) - \mathcal{A}_\theta(\mathbf{z}, t)\|^2 < \delta\} \quad (5)$$

3.3 INTERSECTION PROJECTION IMPLEMENTATION

In this section, we reveal that the intersection projection $\operatorname{Proj}_{\mathfrak{J}_t}(\mathbf{z}_t)$ can be effectively computed through the iterative application of the joint up/down projection $\widehat{\operatorname{Proj}}_{\mathfrak{J}_t}(\mathbf{z}_t)$.

Up/Down Projections are integral to our method. We will introduce two down projections and one up projection here. Specifically, the denoising projection $\operatorname{Proj}_{\mathfrak{D}_t}(\cdot)$ is defined as

$$\operatorname{Proj}_{\mathfrak{D}_t}(\mathbf{z}_{t+1}) = \frac{1}{\sqrt{\alpha_{t+1}}} \mathbf{z}_{t+1} + \frac{(1 - \alpha_{t+1})}{\sqrt{\alpha_{t+1}}} \mathcal{S}_\theta(\mathbf{z}_{t+1}, t + 1) + \sqrt{\frac{(1 - \alpha_{t+1})(1 - \bar{\alpha}_t)}{1 - \bar{\alpha}_{t+1}}} \epsilon \quad (6)$$

in accordance Equation (1). This operator serves to map elements from the set \mathfrak{D}_{t+1} to the set \mathfrak{D}_t . According to Equation (2), we define the condition projection as

$$\operatorname{Proj}_{\mathfrak{C}_t}(\mathbf{z}_{t+1}) = \operatorname{argmin}_{\mathbf{z}_t} \|\mathcal{E}(I_c) - \mathcal{A}_\theta(\mathbf{z}_t, t)\|^2 \quad \text{solving with initial point } \mathbf{z}_{t+1}. \quad (7)$$

This projection is also considered as a mapping from \mathfrak{D}_{t+1} to \mathfrak{D}_t . Collectively, these two projections are termed ‘Down Projections’ due to their index decreasing from $t + 1$ to t . Conversely, we introduce an ‘Up Projection’, which maps a point \mathbf{z}_t from \mathfrak{D}_t into \mathfrak{D}_{t+1} :

$$\operatorname{Proj}_{\mathfrak{D}_{t+1}}(\mathbf{z}_t) = \sqrt{\alpha_{t+1}} \mathbf{z}_t + \sqrt{1 - \alpha_{t+1}} \epsilon \quad (8)$$

Joint Up/Down Projection $\widehat{\operatorname{Proj}}_{\mathfrak{J}_t}(\mathbf{z}_t)$ is devised to calculate $\operatorname{Proj}_{\mathfrak{J}_t}(\cdot)$. For the computation of $\operatorname{Proj}_{\mathfrak{J}_t}(\mathbf{z}_t)$, we define the projection $\widehat{\operatorname{Proj}}_{\mathfrak{J}_t}(\mathbf{z}_t)$ as:

$$\widehat{\operatorname{Proj}}_{\mathfrak{J}_t}(\mathbf{z}_t) = \operatorname{Proj}_{\mathfrak{D}_t} \circ \operatorname{Proj}_{\mathfrak{D}_{t+1}} \circ \operatorname{Proj}_{\mathfrak{C}_t} \circ \operatorname{Proj}_{\mathfrak{D}_{t+1}}(\mathbf{z}_t) \quad (9)$$

This equation recursively projects onto the intersection of sets by leveraging the subset relationship $\mathfrak{D}_t \subseteq \mathfrak{D}_{t+1}$, facilitating convergence towards $\mathfrak{J}_t = \mathfrak{D}_t \cap \mathfrak{C}_t$. The result is a point that lies within both \mathfrak{C}_t and \mathfrak{D}_t . Thus, employing Alternative Projection, we iteratively obtain the value of $\operatorname{Proj}_{\mathfrak{J}_t}(\mathbf{z}_t)$ through the repeated application of $\mathbf{z}_t = \widehat{\operatorname{Proj}}_{\mathfrak{J}_t}(\mathbf{z}_t)$ over N iterations, as illustrated in Algorithm 1. For more details, please refer to Appendix B.



Figure 3: Image Generation with Multiple Controls: APCtrl incorporates multiple conditions into diffusion models. To showcase its capabilities, we present two illustrative examples, each detailed in dedicated sections of the figure. The left example demonstrates the fusion of two conditional sets. The right example leverages ControlNet to project onto the feasible space. Each subfigure includes the prompt text at the top, followed by rows for conditional image controls, concluding with the controlled generation results at the bottom.

4 EXPERIMENTS

In this section, we provide a comprehensive evaluation of our method through both quantitative and qualitative analyses, demonstrating its effectiveness. Additionally, we highlight its versatility by showcasing compatibility with a range of diffusion backbones and samplers. Finally, we underscore the efficiency of our approach.

4.1 EXPERIMENTAL SETUP

We utilize SDv1.5, a prevalent checkpoint, as the backbone for constructing latent control networks. This section details the latent control networks’ training and elucidates the nuances of APCtrl sampling.

Latent Control Networks Training: Our model utilizes the identical U-Net architecture found in SDv1.5 and was initialized using the SDv1.5 checkpoint. The training process was conducted for fewer than 24 GPU hours on a single NVIDIA 3090 GPU. We employed approximately 118,000 images from the COCO2017 dataset (Lin et al., 2014) across various tasks. For human pose estimation, we specifically selected a subset of around 6,500 images from the COCO2017 dataset, concentrating on the “people” category. Additionally, for Style Guidance, we incorporated approximately 81,000 images from the Wiki-Art dataset.

APCtrl Sampling: A multitude of sampling strategies, including DDPM (Ho et al., 2020), DDIM (Song et al., 2020a), and LCM (Luo et al., 2023), can be applied to stable diffusion models. These strategies can be unified by the concept of recursive projection, expressed as $\mathbf{z}_t = \text{Proj}_{\mathcal{D}_t}(\mathbf{z}_{t+1})$. Consequently, $\text{Proj}_{\mathcal{D}_t}(\cdot)$ is able to integrate with all of these approaches to form our APCtrl sampling based on Algorithm 1.

4.2 CONDITIONAL GENERATION RESULTS

APCtrl is adept at integrating a diverse set of conditions directly into the image generation process of diffusion models, offering a framework for sophisticated control over the generation outcomes. To showcase this capability, we demonstrate ten single-condition cases in Figure 2, spanning a spectrum of techniques from Canny edge (Canny, 1986), M-LSD line (Gu et al., 2022), HED edge (Xie & Tu, 2015), Skeleton (Cao et al., 2017), Object Loca-

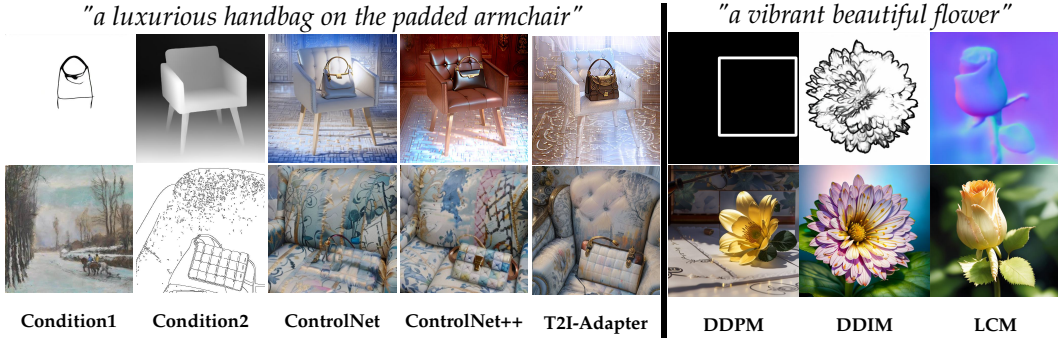


Figure 4: Compatibility Demonstration for Different Diffusion Backbones and Samplers. In left part, APCtrl introduces Condition 1, while the conditional diffusion backbone of the other model introduces Condition 2. The right part shows the generation results using DDPM (Ho et al., 2020), DDIM (Song et al., 2020a), and LCM (Luo et al., 2023).

Table 2: Quantitative Comparison for Controllable Generation on Single Conditions. The best results are in bold. “-” indicates that the method does not provide a public model.

	Method	Depth	Canny	HED	M-LSD	Segmentation	Normal	Skeleton	Location	Sketch
FID↓	ControlNet	19.3801	16.5297	19.9832	19.7612	20.7019	27.8210	57.4144	-	-
	ControlNet++	17.4087	20.1109	15.8372	-	24.3453	-	-	-	-
	T2I-Adapter	23.8945	17.0756	-	-	21.7609	-	34.3569	-	28.8883
	UniCtrlNet	24.9604	17.9107	17.4471	27.7329	22.7066	-	66.5560	-	24.0166
	UniControl	24.2885	18.9211	19.2913	-	29.8068	29.5817	40.7635	29.6951	-
	GLIGEN	23.2859	24.8351	26.3622	-	27.3867	27.6916	53.2974	23.1426	-
	APCtrl	25.0148	25.3836	24.7542	26.9950	25.1360	27.4390	43.5083	33.8875	26.6784
CLIP-scores↑	ControlNet	0.2840	0.2897	0.2870	0.2843	0.2838	0.2730	0.2610	-	-
	ControlNet++	0.3061	0.3085	0.3008	-	0.2997	-	-	-	-
	T2I-Adapter	0.2990	0.3045	-	-	0.2956	-	0.3111	-	0.2708
	UniCtrlNet	0.3017	0.3044	0.3039	0.2873	0.3052	-	0.2795	-	0.2927
	UniControl	0.3063	0.3032	0.2997	-	0.3069	0.2967	0.3089	0.2974	-
	GLIGEN	0.2979	0.2966	0.2806	-	0.2854	0.2718	0.2615	0.2762	-
	APCtrl	0.2952	0.3029	0.3035	0.2951	0.2989	0.2943	0.2920	0.2868	0.3006
CLIP-acs ↑	ControlNet	5.1861	5.2112	5.2536	5.3072	5.2954	5.0815	5.2418	-	-
	ControlNet++	5.2945	5.1216	5.1125	-	4.9270	-	-	-	-
	T2I-Adapter	5.0973	5.1213	-	-	4.9737	-	5.2956	-	4.8516
	UniCtrlNet	5.0129	5.0010	5.0048	4.9704	5.0557	-	4.9568	-	5.0048
	UniControl	5.3498	5.1650	5.1683	-	5.3920	5.1061	5.4802	5.2630	-
	GLIGEN	5.1342	5.0485	4.9547	-	4.9362	4.7384	4.8970	5.3708	-
	APCtrl	5.4225	5.4264	5.4575	5.4341	5.3905	5.4284	5.4794	5.3977	5.4730
Controllability		MSE↓	SSIM↑	SSIM↑	SSIM↑	mIoU↑	MSE↓	mAP↑	mAP↑	SSIM↑
	ControlNet	88.9629	0.4376	0.5845	0.7552	0.4223	86.8804	0.4413	-	-
	ControlNet++	86.7270	0.5386	0.6907	-	0.5481	-	-	-	-
	T2I-Adapter	94.5548	0.3984	-	-	0.2339	-	0.4979	-	0.3756
	UniCtrlNet	99.3874	0.4679	0.6159	0.7250	0.3037	-	0.2046	-	0.6704
	UniControl	88.8402	0.5340	0.3614	-	0.3273	104.1840	0.2463	0.2731	-
	GLIGEN	81.1289	0.3917	0.4094	-	0.2481	90.3527	0.1817	0.2418	-
APCtrl	86.6756	0.4412	0.4752	0.7793	0.3916	70.3443	0.4043	0.2563	0.6118	

tion (Redmon et al., 2016), Depth Map (Yang et al., 2024a), Normal Map (Vasiljevic et al., 2019), Segmentation (Cheng et al., 2022), Style Guidance (Radford et al., 2021). These cases serve to illustrate the versatility of APCtrl in handling different conditional inputs.

Additionally, APCtrl can be applied to multiple condition generation, which typically involve two strategies. The first strategy, as shown in the left part of Figure 3, involves augmenting the condition set and computing the projection function, denoted as $\text{Proj}_{\mathcal{J}_t}(\cdot)$, where \mathcal{J}_t is defined as the intersection of \mathcal{D}_t and \mathcal{C}_t , with \mathcal{C}_t being redefined as the intersection of multiple condition sets, expressed as $\mathcal{C}_t = \bigcap_{i=1}^N \mathcal{C}_t^{(i)}$. We can then extend Equation (9) to calculate the updated projection $\text{Proj}_{\mathcal{J}_t}(\cdot)$. The alternative method, as depicted in the right part of Figure 3, involves utilizing ControlNet to define the projection onto \mathcal{D}_t . Given that ControlNet includes a control mechanism, it can be effectively integrated with other constraints defining the condition set.

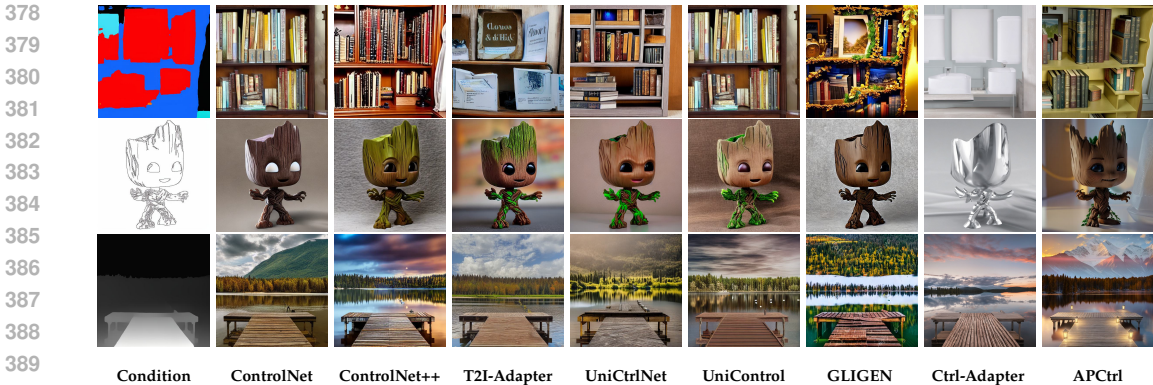


Figure 5: Visual Comparison of Single Condition for Control-on-Training Methods. The prompt for each row is “a book shelf”, “a groot toy”, and “brown wooden dock near lake”. Although our method is a Control-on-Sampling approach that doesn’t necessitate retraining the control network, the conditional results of Segmentation, Canny Edge, and Depth Map, are competitive with other methods.

4.3 COMPATIBILITY FOR DIFFERENT DIFFUSION BACKBONES AND SAMPLERS

APCtrl is compatible with any diffusion model that utilizes the same encoder $\mathcal{E}(\cdot)$ as adopted in Equation 3. This compatibility is exemplified in the left side of Figure 4, where APCtrl is integrated with models such as ControlNet (Zhang et al., 2023), ControlNet++ (Li et al., 2024b), and T2I-Adapter (Mou et al., 2024). In these integrations, APCtrl supplies Condition 1, complemented by Condition 2 from the respective diffusion backbones.

As per Algorithm 1, APCtrl integrates an intersection projection $\text{Proj}_{\mathcal{J},(\cdot)}$ into the diffusion model’s reverse process, enhancing it without altering the original computations. This integration thus is sampling agnosticism, allowing APCtrl to be versatile with various sampling techniques. The right side of Figure 4 illustrates APCtrl’s application and effectiveness with different sampling methods, including DDPM (Ho et al., 2020), DDIM (Song et al., 2020a), and LCM (Luo et al., 2023). These examples underscore APCtrl’s adaptability across a range of diffusion model sampling approaches.



Figure 6: Visual Comparison of Single Condition for Control-on-Sampling Methods. The current Control-on-Sampling methods do not provide the same variety of conditions as Control-on-Training methods, as shown in Figure 5. In this comparison, we focus on two prevalent scenarios, Segmentation and Style Guidance.

4.4 QUANTITATIVE AND QUALITATIVE COMPARISON

Quantitative Evaluation: Our quantitative assessment is conducted on the COCO2017 (Lin et al., 2014) validation set at a 512×512 resolution. This dataset comprises 5000 images, each with multiple captions. For our evaluation, we randomly select one caption per image, yielding 5000 generated images. Specifically, for the skeleton, we focused on the "people" category and selected 2900 images. All methods employ 20 DDIM steps for fast evaluation. We evaluate generation quality using FID (Heusel et al., 2017), CLIP text-image score (Radford et al., 2021), CLIP aesthetic score (Schuhmann et al., 2022). We also evaluate controllability using SSIM (Structural Similarity), mAP (mean Average Precision),

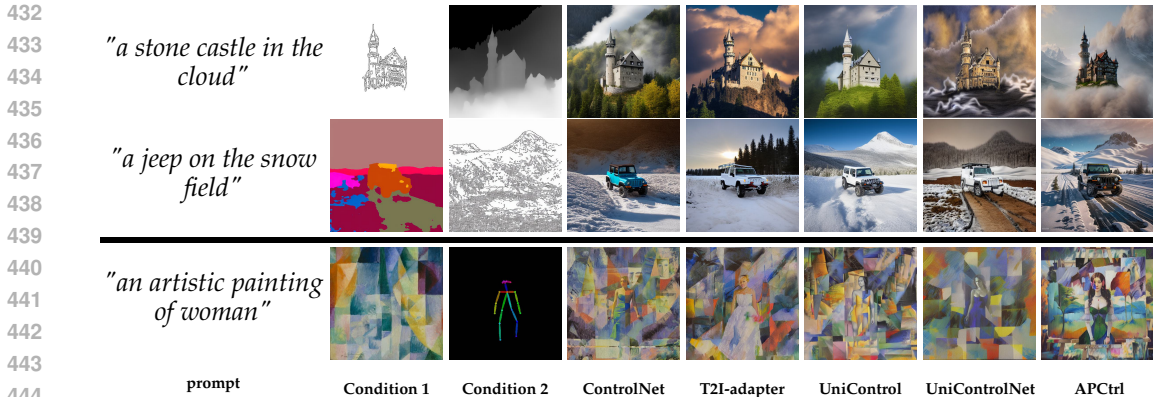


Figure 7: Comparison for Controllable Generation on Multiple Conditions. ControlNet (Zhang et al., 2023), T2I-Adapter (Mou et al., 2024), UniControl (Qin et al., 2023), UniControlNet (Zhao et al., 2024), and APCtrl combine two conditions to generate the final image. These methods initially lack style guidance capabilities, but integrating APCtrl provides this functionality to both. The results are demonstrated in the bottom row

MSE(Mean Squared Error), mIoU(Mean Intersection over Union). We take these metrics to compare the conditions extracted with the natural images (ground truth) with the generated images and report statistical data in Table 2.

Table 2 enumerates six Control-on-Training methods: ControlNet (Zhang et al., 2023), ControlNet++ (Li et al., 2024b), T2I-Adapter (Mou et al., 2024), UniControlNet (Zhao et al., 2024), UniControl (Qin et al., 2023), and GLIGEN (Li et al., 2023). The exclusion of Control-on-Sampling methods is due to their inability to uniformly address all presented conditions and the impracticality of their long inference times for 5,000 images. Notably, APCtrl stands as the pioneering Control-on-Sampling method capable of handling the full spectrum of conditional generation cases associated with Control-on-Training methods. The quantitative analysis in the table demonstrates our model’s superiority over existing approaches in performance metrics—FID, CLIP-score, and CLIP-acs—and controllability, as evidenced in the respective rows and the Controllability column, across the majority of conditions.

Qualitative Evaluation: We enrich our analysis with qualitative comparisons of single and multi-condition controls, as depicted in Figures 5, 6, and 7. In Figure 5, seven Control-on-Training methods—ControlNet++ (Li et al., 2024b), T2I-Adapter (Mou et al., 2024), UniControlNet (Zhao et al., 2024), UniControl (Qin et al., 2023), GLIGEN (Li et al., 2023), and Ctrl-Adapter (Lin et al., 2024)—exhibit strong performance in Segmentation, Canny Edge, and Depth Map conditions. Our method also shows competitive alignment with the input conditions. To the best of our knowledge, APCtrl is a groundbreaking Control-on-Sampling method, capable of delivering the full range of condition controls associated with Control-on-Training methods. For more visual results, please refer to the Appendix B.

For Control-on-Sampling methods, only UniGuidance (Bansal et al., 2024) and FreeDoM (Yu et al., 2023) results are featured due to MultiDiffusion’s specialization in image merging (Bartal et al., 2023) and ReSample’s (Song et al., 2023) focuses on solving linear inverse problems. Figure 6 documents the generation results for Segmentation and Style Guidance conditions from UniGuidance, FreeDoM, and our APCtrl.

Given the limitations of GLIGEN and Uni-Guidance in multi-condition scenarios and the restricted conditions of ControlNet++, our multi-condition comparison is narrowed down to ControlNet, T2I-Adapter, UniControlNet, and UniControl, as shown in Figure 7. APCtrl’s performance in multi-condition tasks is comparable to these Control-on-Training methods.

Originally lacking in style guidance for the four methods, the integration with APCtrl has unlocked this capability, as evidenced in the figure’s bottom row, highlighting APCtrl’s seamless integration and multi-condition generation capabilities when paired with these methods. For more visual results, please refer to the Appendix C.

486 Table 3: Efficiency Comparison. APCtrl achieves a balance between training efficiency
 487 and sampling speed. Compared to Control-on-Training methods, it requires minimal
 488 training investment. When compared to Control-on-Sampling methods, APCtrl provides
 489 accelerated sampling.

Method	Control-on-Training						Control-on-Sampling			
	ControlNet	ControlNet++	T2I-Adapter	UniCtrlNet	UniControl	GLIGEN	FreeDom	DSG	UniGuid	APCtrl
Training (GPU Hours)	500	60	192	6900	5000	1000	-	-	-	20
Sampling (Seconds)	3	2	2	4	5	7	115	122	4510	13

494
495
496 4.5 EFFICIENCY COMPARISON

497 This section is devoted to compare the efficiency of training and sampling. The latent
 498 control network in APCtrl projects the image space onto the conditional space, aligning it
 499 with the encoded constraints $\mathcal{E}(I_c)$. This method is significantly different from Control-
 500 on-Training, as we believe that working with constraint spaces is less complex than the
 501 refinement processes in Control-on-Training. This suggests that APCtrl should require less
 502 training time, as indicated in Table 3. We observe that ControlNet++ requires only 60 hours
 503 because it fine-tunes from a ControlNet checkpoint. The control function of APCtrl is also
 504 applied during the sampling phase. It is essential to compare APCtrl’s sampling efficiency
 505 with that of Control-on-Sampling methods. Thanks to the use of Alternative Projection,
 506 APCtrl notably decreases sampling time, as demonstrated in the last line of Table 3. This
 507 makes APCtrl a promising alternative, offering a favorable compromise for both training
 508 and sampling phases.

509
510 4.6 IMPACT OF ITERATION COUNT ON UP/DOWN ALTERNATIVE PROJECTIONS

511 In Algorithm 1, we carry out N iterations of
 512 the Joint Up/Down Projection $\widehat{\text{Proj}}_{\mathcal{J}_t}(\cdot)$, fol-
 513 lowing the definition given in Equation 9.
 514 When N is small, the resulting images may
 515 not align with the desired conditions, as
 516 evidenced in Figure 8b. Conversely, a large
 517 N can lead to images that meet the condi-
 518 tions but display pronounced color discrep-
 519 ancies, diminishing their aesthetic quality,
 520 as shown in Figure 8d. Striking the optimal
 521 balance with N ensures both conditional fi-
 522 delity and visual appeal, as demonstrated
 523 in Figure 8c.



524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
Figure 8: Visualization For iteration Count N

5 CONCLUSIONS

Achieving controllable generation remains one of the significant challenges in diffusion models. We propose a novel direction to address this challenge. Our approach begins with reinterpreting the diffusion sampling process as a series of recursive projections onto the denoising set, denoted as \mathcal{D}_t . Consequently, a conditional control diffusion model can be viewed as a sequence of recursive projections onto the intersection of feasible sets, $\mathcal{D}_t \cap \mathcal{C}_t$, where \mathcal{C}_t represents the condition set. We employ an alternative projection technique to effectively implement these projections onto the intersection set $\mathcal{D}_t \cap \mathcal{C}_t$. This methodology offers several distinct advantages over previous efforts: 1. Multi-Condition Generation: Multi-condition generation can be easily implemented. 2. Model and Sampling Agnosticism: APCtrl maintains agnosticism regarding both the underlying model backbones and the sampling process. 3. Unified Control Loss: It allows for a unified control loss, facilitating the management of various control applications. 4. Efficiency: APCtrl significantly reduces both training and sampling times. We have conducted a comprehensive evaluation of our framework, yielding state-of-the-art results.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

REFERENCES

- Shmuel Agmon. The Relaxation Method for Linear Inequalities. *Canadian Journal of Mathematics*, 1954.
- Arpit Bansal, Eitan Borgnia, Hong-Min Chu, Jie S. Li, Hamid Kazemi, Furong Huang, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Cold Diffusion: Inverting Arbitrary Image Transforms Without Noise. In *Conference on Neural Information Processing Systems*, 2023.
- Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal Guidance for Diffusion Models. In *International Conference on Learning Representations*, 2024.
- Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: Fusing Diffusion Paths for Controlled Image Generation. *Proceedings of Machine Learning Research*, 2023.
- Stephen Boyd. Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers. *Foundations and Trends® in Machine Learning*, 2010.
- James P. Boyle and Richard L. Dykstra. A Method for Finding Projections onto the Intersection of Convex Sets in Hilbert Spaces. In *Advances in Order Restricted Statistical Inference*. 1986.
- L. M. Brègman. Finding the Common Point of Convex Sets by the Method of Successive Projection. *Dokl. Akad. Nauk SSSR*, 1965.
- John Canny. A Computational Approach to Edge Detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1986.
- Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. OpenPose: Realtime Multi-Person 2D Pose Estimation Using Part Affinity Fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- Stanley H Chan. Tutorial on Diffusion Models for Imaging and Vision. *arXiv*, 2024.
- Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention Mask Transformer for Universal Image Segmentation. In *Conference on Computer Vision and Pattern Recognition*, 2022.
- Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving Diffusion Models for Inverse Problems using Manifold Constraints. In *Advances in Neural Information Processing Systems*, 2022.
- Hyungjin Chung, Jeongsol Kim, Michael Thompson Mccann, Marc Louis Klasky, and Jong Chul Ye. Diffusion Posterior Sampling for General Noisy Inverse Problems. In *International Conference on Learning Representations*, 2023.
- Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion Models in Vision: A Survey. *Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- Frank Deutsch. The Method of Alternating Orthogonal Projections. In *Approximation Theory, Spline Functions and Applications*. 1992.
- Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion Models Beat GANs on Image Synthesis. In *Advances in Neural Information Processing Systems*, 2021.
- Jim Douglas and Henry H Rachford. On the numerical solution of heat conduction problems in two and three space variables. *Transactions of the American mathematical Society*, 1956.
- Mattias Fält and Pontus Giselsson. Generalized Alternating Projections on Manifolds and Convex Sets. *Journal of Nonsmooth Analysis and Optimization*, 2024.
- Alexandros Graikos, Nikolay Malkin, Nebojsa Jojic, and Dimitris Samaras. Diffusion Models as Plug-and-Play Priors. *Advances in Neural Information Processing Systems*, 2022.

- 594 Geonmo Gu, Byungsoo Ko, SeoungHyun Go, Sung-Hyun Lee, Jingeun Lee, and Minchul
595 Shin. Towards Light-Weight and Real-Time Line Segment Detection. In *AAAI Conference*
596 *on Artificial Intelligence*, 2022.
- 597 LG Gubin, Boris T Polyak, and EV Raik. The Method of Projections for Finding the Common
598 Point of Convex Sets. *USSR Computational Mathematics and Mathematical Physics*, 1967.
- 600 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp
601 Hochreiter. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash
602 Equilibrium. *Advances in Neural Information Processing Systems*, 2017.
- 603 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. In
604 *Advances in Neural Information Processing Systems*, 2020.
- 605 Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising Diffusion
606 Restoration Models. *Advances in Neural Information Processing Systems*, 2022.
- 608 Alexander Y Kruger and Nguyen H Thao. Regularity of Collections of Sets and Convergence
609 of Inexact Alternating Projections. *Journal of Convex Analysis*, 2016.
- 610 Dongxu Li, Junnan Li, and Steven Hoi. Blip-Diffusion: Pre-Trained Subject Representation
611 for Controllable Text-to-Image Generation and Editing. *Advances in Neural Information*
612 *Processing Systems*, 2024a.
- 614 Ming Li, Taojiannan Yang, Huafeng Kuang, Jie Wu, Zhaoning Wang, Xuefeng Xiao, and
615 Chen Chen. Controlnet++: Improving conditional controls with efficient consistency
616 feedback. In *European Conference on Computer Vision*, 2024b.
- 617 Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chun-
618 yuan Li, and Yong Jae Lee. GLIGEN: Open-Set Grounded Text-to-Image Generation. In
619 *Conference on Computer Vision and Pattern Recognition*, 2023.
- 620 Han Lin, Jaemin Cho, Abhay Zala, and Mohit Bansal. Ctrl-Adapter: An Efficient and
621 Versatile Framework for Adapting Diverse Controls to Any Diffusion Model, 2024.
- 622 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan,
623 Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common Objects in Context. In
624 *European Conference on Computer Vision*, 2014.
- 625 Pierre-Louis Lions and Bertrand Mercier. Splitting Algorithms for the Sum of Two Nonlinear
626 Operators. *SIAM Journal on Numerical Analysis*, 1979.
- 627 Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc
628 Van Gool. Repaint: Inpainting Using Denoising Diffusion Probabilistic Models. In
629 *Conference on Computer Vision and Pattern Recognition*, 2022.
- 630 Calvin Luo. Understanding Diffusion Models: A Unified Perspective. *arXiv*, 2022.
- 631 Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent Consistency Models:
632 Synthesizing High-Resolution Images with Few-Step Inference, 2023.
- 633 Theodore Samuel Motzkin and Isaac Jacob Schoenberg. The Relaxation Method for Linear
634 Inequalities. *Canadian Journal of Mathematics*, 1954.
- 635 Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying
636 Shan. T2I-Adapter: Learning Adapters to Dig Out More Controllable Ability for Text-to-
637 Image Diffusion Models. In *AAAI Conference on Artificial Intelligence*, 2024.
- 638 Alexander Quinn Nichol and Prafulla Dhariwal. Improved Denoising Diffusion Probabilistic
639 Models. In *International Conference on Machine Learning*, 2021.
- 640 Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela
641 Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. GLIDE: Towards Photorealistic
642 Image Generation and Editing with Text-Guided Diffusion Models. In *International*
643 *Conference on Machine Learning*, 2022.
- 644

- 648 Gaurav Parmar, Taesung Park, Srinivasa Narasimhan, and Jun-Yan Zhu. One-Step Image
649 Translation with Text-to-Image Models, 2024.
650
- 651 Can Qin, Shu Zhang, Ning Yu, Yihao Feng, Xinyi Yang, Yingbo Zhou, Huan Wang, Juan Car-
652 los Niebles, Caiming Xiong, Silvio Savarese, et al. UniControl: A Unified Diffusion Model
653 for Controllable Visual Generation In the Wild. In *Conference on Neural Information Pro-
654 cessing Systems*, 2023.
- 655 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-
656 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning Transfer-
657 able Visual Models From Natural Language Supervision. In *International Conference on
658 Machine Learning*, 2021.
- 659 Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You Only Look Once:
660 Unified, Real-Time Object Detection. In *Conference on Computer Vision and Pattern Recog-
661 nition*, 2016.
662
- 663 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.
664 High-Resolution Image Synthesis with Latent Diffusion Models. In *Conference on Computer
665 Vision and Pattern Recognition*, 2022.
- 666 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir
667 Aberman. Dreambooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven
668 Generation. In *Conference on Computer Vision and Pattern Recognition*, 2023.
669
- 670 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman,
671 Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al.
672 Laion-5B: An Open Large-Scale Dataset for Training Next Generation Image-Text Models.
673 *Advances in Neural Information Processing Systems*, 2022.
- 674 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep Unsu-
675 pervised Learning Using Nonequilibrium Thermodynamics. In *International Conference
676 on Machine Learning*, pp. 2256–2265. PMLR, 2015.
677
- 678 Bowen Song, Soo Min Kwon, Zecheng Zhang, Xinyu Hu, Qing Qu, and Liyue Shen. Solving
679 Inverse Problems with Latent Diffusion Models via Hard Data Consistency. In *Interna-
680 tional Conference on Learning Representations*, 2023.
- 681 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising Diffusion Implicit Models. In
682 *International Conference on Learning Representations*, 2020a.
683
- 684 Yang Song and Stefano Ermon. Generative Modeling by Estimating Gradients of the Data
685 Distribution. *Advances in Neural Information Processing Systems*, 2019.
- 686 Yang Song and Stefano Ermon. Improved Techniques for Training Score-based Generative
687 Models. *Advances in Neural Information Processing Systems*, 2020.
688
- 689 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon,
690 and Ben Poole. Score-Based Generative Modeling through Stochastic Differential Equa-
691 tions. In *International Conference on Learning Representations*, 2020b.
- 692 Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum Likelihood Training
693 of Score-based Diffusion Models. *Advances in Neural Information Processing Systems*, 2021.
694
- 695 Igor Vasiljevic, Nick Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z Dai,
696 Andrea F Daniele, Mohammadreza Mostajabi, Steven Basart, Matthew R Walter, et al.
697 DIODE: A Dense Indoor and Outdoor DEpth Dataset, 2019.
- 698 John Von Neumann. *Functional Operators: The Geometry of Orthogonal Spaces*. Princeton
699 University Press, 1951.
700
- 701 Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-Shot Image Restoration Using Denoising
Diffusion Null-Space Model. In *International Conference on Learning Representations*, 2022.

702 Saining Xie and Zhuowen Tu. Holistically-Nested Edge Detection. In *International Conference*
703 *on Computer Vision*, 2015.
704

705 Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao.
706 Depth anything: Unleashing the power of large-scale unlabeled data. In *Conference on*
707 *Computer Vision and Pattern Recognition*, 2024a.

708 Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao
709 Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion Models: A Comprehensive Survey of
710 Methods and Applications. *ACM Computing Surveys*, 2023a.
711

712 Ling Yang, Jingwei Liu, Shenda Hong, Zhilong Zhang, Zhilin Huang, Zheming Cai, Wentao
713 Zhang, and Bin Cui. Improving Diffusion-based Image Synthesis with Context Predic-
714 tion. *Advances in Neural Information Processing Systems*, 2024b.

715 Lingxiao Yang, Shutong Ding, Yifan Cai, Jingyi Yu, Jingya Wang, and Ye Shi. Guidance
716 with spherical gaussian constraint for conditional diffusion. In *International Conference on*
717 *Machine Learning*, 2024c.

718 Zhengyuan Yang, Jianfeng Wang, Zhe Gan, Linjie Li, Kevin Lin, Chenfei Wu, Nan Duan,
719 Zicheng Liu, Ce Liu, Michael Zeng, et al. ReCo: Region-Controlled Text-to-Image Gen-
720 eration. In *Conference on Computer Vision and Pattern Recognition*, 2023b.
721

722 Jiwen Yu, Yinhuai Wang, Chen Zhao, Bernard Ghanem, and Jian Zhang. FreeDoM: Training-
723 Free Energy-Guided Conditional Diffusion Model. In *International Conference on Computer*
724 *Vision*, 2023.

725 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding Conditional Control to Text-to-
726 Image Diffusion Models. In *International Conference on Computer Vision*, 2023.
727

728 Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin Bao, Shaozhe Hao, Lu Yuan,
729 and Kwan-Yee K Wong. Uni-ControlNet: All-in-One Control to Text-to-Image Diffusion
730 Models. *Advances in Neural Information Processing Systems*, 2024.
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A APPENDIX: DETAILS FOR TRAINING LATENT CONTROL NETWORK

Our implementation, including all training and sampling files as well as the checkpoint, will be released upon acceptance of our paper. Briefly, our code is adapted from the Hugging Face Diffusers repository, specifically from the `train_text_to_image.py` script, which can be found at <https://github.com/huggingface/diffusers>. In this section, we detail the key components of the code responsible for constructing and training the latent control network $\mathcal{A}_\theta(\cdot, \cdot)$.

A.1 BUILDING MODEL CODE

```
from diffusers import UNet2DConditionModel

# Initialize the A model (UNet) from pretrained models
A_model = UNet2DConditionModel.from_pretrained(
    "runwayml/stable-diffusion-v1-5", subfolder="unet"
)
```

A.2 TRAINING CODE

```
from diffusers import AutoencoderKL, UNet2DConditionModel

# Load the VAE model using the default configuration
vae = AutoencoderKL.from_pretrained(
    "runwayml/stable-diffusion-v1-5", subfolder="vae"
)

# Set up the optimizer for the A model
optimizer = torch.optim.AdamW(
    A_model.parameters(), lr=1e-5
)

for batch in dataloader:
    ...

    # Encode the ground truth image to the latent space
    latents = vae.encode(ground_truth_image)

    # Add noise to the latents
    noisy_latents = noise_scheduler.add_noise(
        latents, noise, timesteps)

    # Predict the conditional latents from noisy latents
    pred_latents = A_model(
        noisy_latents, timesteps)

    # Encode the condition image to the latent space
    cond_latents = vae.encode(condition_image)

    # Calculate the loss
    loss = F.mse_loss(pred_latents_pred, cond_latents)

    ...
```

B APPENDIX: DETAILS FOR APCTRL SAMPLING

B.1 DETAILED IMPLEMENTATION FOR ALGORITHM 1

In this section, we delve into the specifics of APCtrl sampling as encapsulated by Algorithm 1. This algorithm outlines the step-by-step procedure for implementing our novel approach to conditional diffusion sampling, which leverages the power of latent control networks and alternative projections to achieve sophisticated image generation. The details provided here will walk through the algorithm’s operations, explaining how each step contributes to the final output, ensuring a clear understanding of the methodology and its advantages over traditional approaches.

Based on the discussion in our paper, we have outlined the steps and principles of our method.

$$\text{Proj}_{\mathcal{D}_{t+1}}(\mathbf{z}_t) = \sqrt{\alpha_{t+1}}\mathbf{z}_t + \sqrt{1 - \alpha_{t+1}}\boldsymbol{\epsilon}$$

$$\text{Proj}_{\mathcal{D}_t}(\mathbf{z}_{t+1}) = \frac{1}{\sqrt{\alpha_{t+1}}}\mathbf{z}_{t+1} + \frac{(1 - \alpha_{t+1})}{\sqrt{\alpha_{t+1}}}\mathcal{S}_{\theta}(\mathbf{z}_{t+1}, t + 1) + \sqrt{\frac{(1 - \alpha_{t+1})(1 - \bar{\alpha}_t)}{1 - \bar{\alpha}_{t+1}}}\boldsymbol{\epsilon}$$

$$\text{Proj}_{\mathcal{E}_t}(\mathbf{z}_{t+1}) = \underset{\mathbf{z}_t}{\text{argmin}} \|\mathcal{A}_{\theta}(\mathbf{z}_t, t) - \mathcal{E}(\mathcal{I}_c)\|^2 \quad \text{solving with initial point } \mathbf{z}_t = \text{Proj}_{\mathcal{D}_t}(\mathbf{z}_{t+1})$$

Note that $\text{Proj}_{\mathcal{E}_t}(\mathbf{z}_{t+1})$ can be solved by gradient descent. Considering these projections, they enable us to complete Algorithm 1 as follows:

Algorithm 2 APCtrl Sampling

Input: Initial noise \mathbf{z}_T Diffusion Model $\mathcal{Z}_{\theta}(\mathbf{z}_t, t)$ Latent Control Network $\mathcal{A}_{\theta}(\mathbf{z}_t, t)$ Encoder \mathcal{E} Condition \mathcal{I}_c step size γ

```

836 8 Operator ProjDt+1(zt):
837   | return  $\sqrt{\alpha_{t+1}}\mathbf{z}_t + \sqrt{1 - \alpha_{t+1}}\boldsymbol{\epsilon}$ 
838 10 Operator ProjDt(zt+1):
839   | return  $\frac{1}{\sqrt{\alpha_{t+1}}}\mathbf{z}_{t+1} + \frac{(1 - \alpha_{t+1})}{\sqrt{\alpha_{t+1}}}\mathcal{S}_{\theta}(\mathbf{z}_{t+1}, t + 1) + \sqrt{\frac{(1 - \alpha_{t+1})(1 - \bar{\alpha}_t)}{1 - \bar{\alpha}_{t+1}}}\boldsymbol{\epsilon}$ 
840 12 Operator ProjEt(zt):
841   | for m = 1 to M do
842   |   | zt ← zt - γ ∇x || $\mathcal{A}_{\theta}(\mathbf{z}_t, t) - \mathcal{E}(\mathcal{I}_c)$ ||
843   | end
844   | return x
845 17 Operator ProjJt(zt):
846   | zt ← ProjDt+1(zt)
847   | zt ← ProjEt(zt)
848   | zt ← ProjDt+1(zt)
849   | zt ← ProjDt(zt)
850   | return x
851 23 Operator ProjJt(zt):
852   | for n = 1 to N do
853   |   | zt ← ProjJt(zt)
854   | end
855   | return x
856 28
857 29 zt ← zT
858 30 for t = T - 1 to 1 do
859   | zt ← ProjDt(zt+1)
860   | if conducting APCtrl Sampling then
861   |   | zt ← ProjJt(zt)
862   | end
863 35 end
Output: xt

```

864 B.2 PSEUDOCODE FOR IMPLEMENTING ALGORITHM 1

865
866 This pseudocode outlines the steps to implement the alternative projection method as
867 described in Algorithm 1. It starts with an initial point and iteratively applies projections
868 onto the given sets until the desired number of iterations is reached, resulting in a point
869 that is close to the intersection of the two sets.

```
870 ...
871
872 def project_to_cond(z_t, A_model, cond_latents):
873     for n in range(N):
874         optimizer = torch.optim.RMSProp([z_t], lr=preset_lr)
875         optimizer.zero_grad()
876
877         # Get the predicted conditional latents
878         pred_latents = A_model(z_t, t)
879
880         # Calculate the deviation
881         # between predicted latents and conditional latents.
882         loss = F.mse_loss(pred_latents, cond_latents)
883
884         # Make the current latents better match the conditions.
885         loss.backward()
886         optimizer.step()
887
888     return z_t
889
890 def alternative_projection_sampling(z_t, A_model, cond_latent):
891     for m in range(N):
892         z_t = add_noise(z_t) # up projections
893         z_t = project_to_cond(z_t, A_model, cond_latent)
894         z_t = add_noise(z_t) # up projections
895         z_t = default_denoise(z_t) # original denoising process
896
897     return z_t
898
899 # Load the trained A_model
900 A_model = UNet2DConditionModel.from_pretrained(
901     "trained_checkpoint", subfolder="UNET"
902 )
903
904 # Encode the condition image into the latent space
905 cond_latents = vae.encode(condition_image)
906
907 # sampling
908 for t in timesteps:
909     z_t = default_denoise(z_t) # original denoising process
910     if do_APctrl:
911         z_t = alternative_projection_sampling(
912             z_t, A_model, cond_latents
913         )
914 ...
```

914 C APPENDIX: MORE QUALITATIVE EVALUATION

915
916 We present additional visual results, including various single-condition cases. Moreover,
917 we also demonstrate more results under multiple conditions.

918

"a red telephone booth near a lantern in the street"

919



920

921

922

923

924

925

926

"a hammock hanged on the tree on the beach"

927



928

929

930

931

932

933

934

"a white and black drum set on the grass"

935



936

937

938

939

940

941

942

Depth Map

943

944

945

946

947

"slices of bread on the wooden chopping board"

948



949

950

951

952

953

954

955

"a product shot of the classical leather camera"

956



957

958

959

960

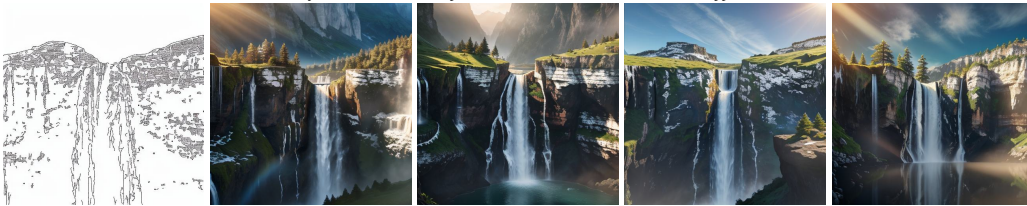
961

962

963

"a majestic waterfall down the rock cliffside"

964



965

966

967

968

969

970

971

Canny Edge

972

"an artistic paintings of the bamboo forest"

973

974

975

976

977

978

979



980

"an artistic paintings of the amusement park"

981

982

983

984

985

986

987



988

"an artistic paintings of the anchor"

989

990

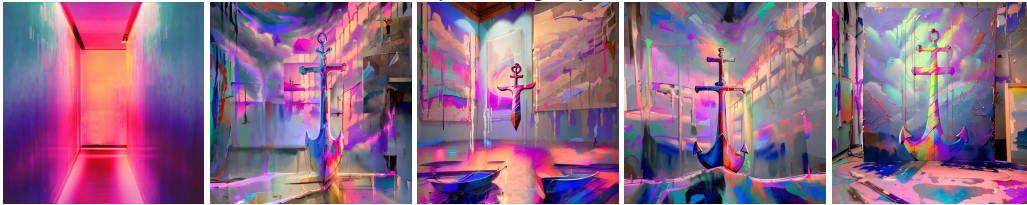
991

992

993

994

995



996

Style Guidance

997

998

999

1000

1001

"green plants on the white ceramic pot"

1002

1003

1004

1005

1006

1007

1008



1009

"bonfire near body of water during night time"

1010

1011

1012

1013

1014

1015

1016



1017

"wooden park chair under tree in the fog"

1018

1019

1020

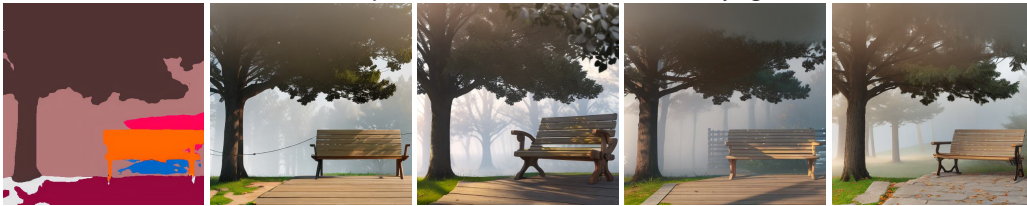
1021

1022

1023

1024

1025



Segmentation

1026

"a primitive kettle on the log table"

1027



1028

1029

1030

1031

1032

1033

1034

"a green cactus plants against the wall"

1035



1036

1037

1038

1039

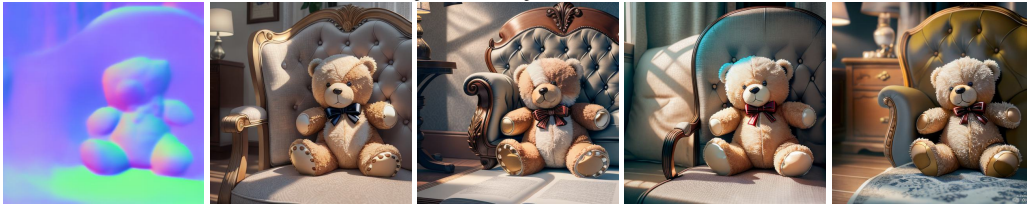
1040

1041

1042

"a teddy bear toy on the chair"

1043



1044

1045

1046

1047

1048

1049

1050

Normal Map

1051

1052

1053

1054

1055

"closeup photo of a man against the wall"

1056



1057

1058

1059

1060

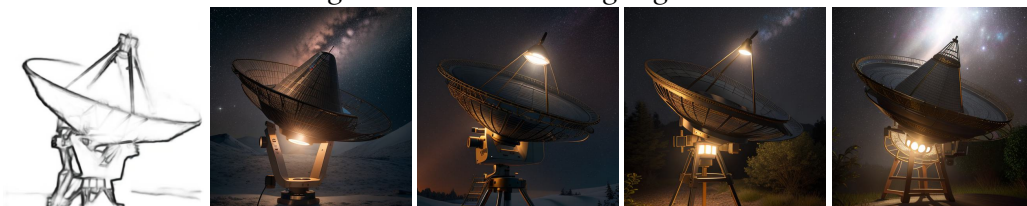
1061

1062

1063

"a huge satellite dish during night time"

1064



1065

1066

1067

1068

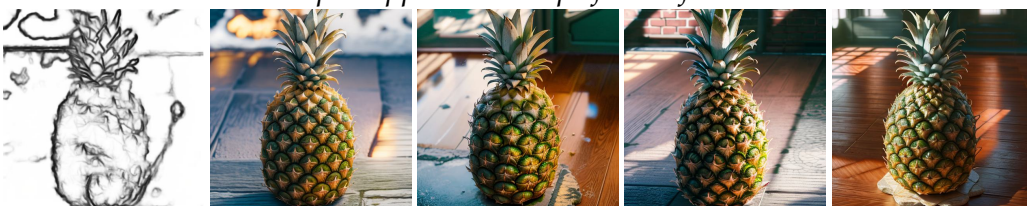
1069

1070

1071

"a pineapple on the top of brick floor"

1072



1073

1074

1075

1076

1077

1078

1079

HED Edge

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

"an old wooden chair on the grass"



"a table near a gray sofa in the living room"

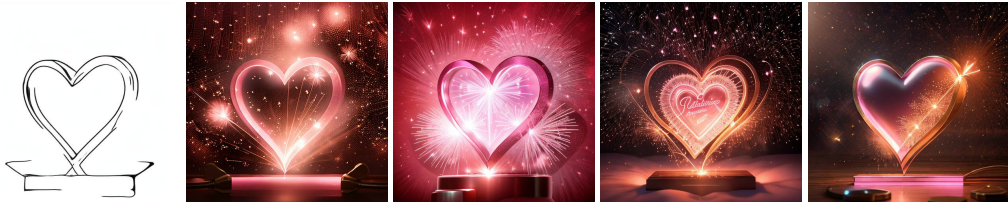


"low angle photo of city high rise building"



M-LSD Line

"the heart shaped pink sparklers"



"a penguin is on the snow land"



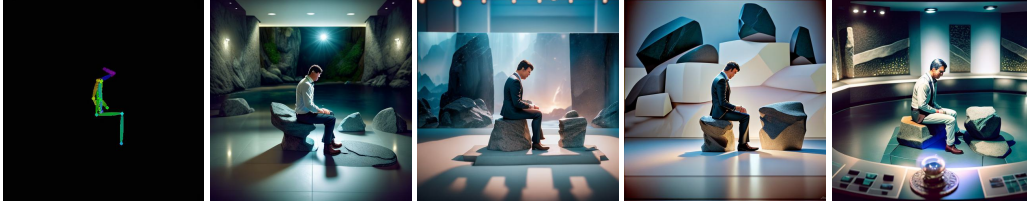
"a product shot of the hamburger"



Sketch

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

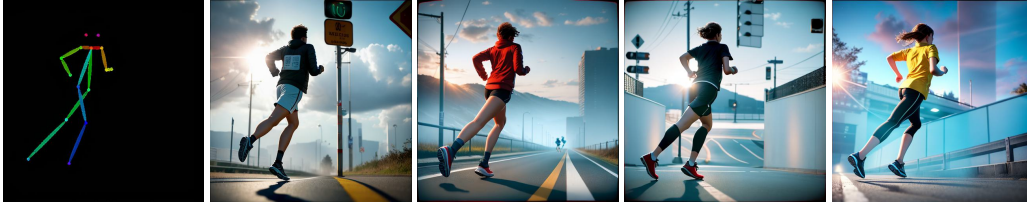
"a man is sitting on the side of stone"



"closeup photo of two girls wearing jacket"



"a young people running down the road"



Skeleton

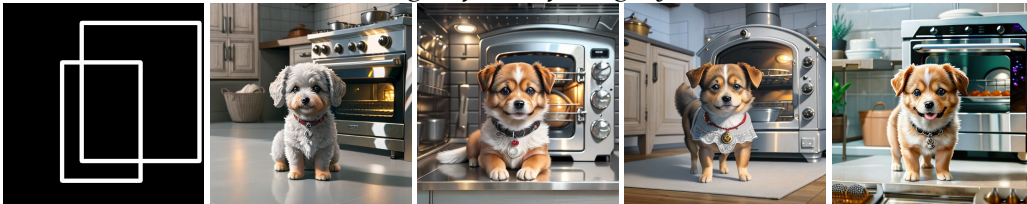
"a chocolate cake on the table"



"two white sheep on the grass"



"a cute dog in front of the grey oven"



Object Location

1188

"a close up shot of the pink pitaya"

1189

1190

1191

1192

1193

1194

1195



1196

"a tree in front of the house"

1197

1198

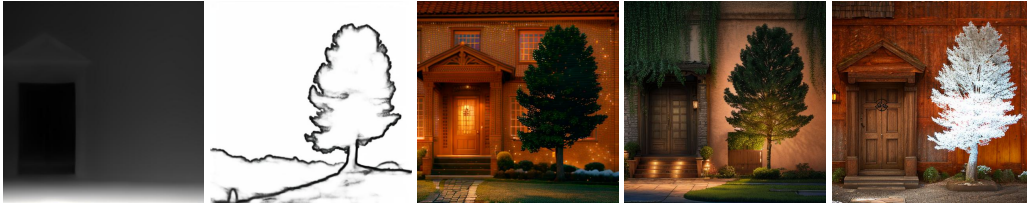
1199

1200

1201

1202

1203



1204

"an elephant walking through the street"

1205

1206

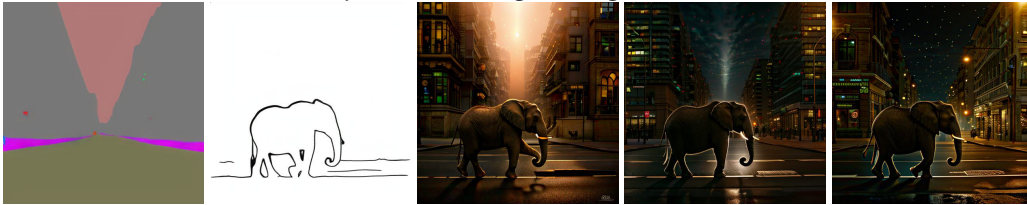
1207

1208

1209

1210

1211



1212

Two Conditions

1213

1214

1215

"a deer on the grass field during daytime"

1216

1217

1218

1219

1220

1221

1222

1223



1224

"beautiful flowers in the vase"

1225

1226

1227

1228

1229

1230

1231



1232

"an artistic painting of the mysterious mountain"

1233

1234

1235

1236

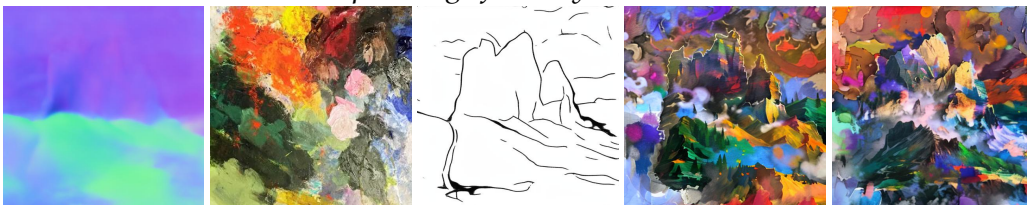
1237

1238

1239

1240

1241



Three Conditions