# Scaling Group Inference for Diverse and High-Quality Generation

**Gaurav Parmar**[1]    **Or Patashnik**[2,3]    **Daniil Ostashev**[2]    **Kuan-Chieh Wang**[2]    **Kfir Aberman**[2]

**Srinivasa Narasimhan**[1]    **Jun-Yan Zhu**[1]

[1]Carnegie Mellon University    [2]Snap Research    [3]Tel Aviv University

## Abstract

Generative models typically sample outputs independently, and recent inference-time guidance and scaling algorithms focus on improving the quality of individual samples. However, in real-world applications, users are often presented with *a set* of multiple images (e.g., 4-8) for each prompt, where independent sampling tends to lead to redundant results, limiting user choices and hindering idea exploration. In this work, we introduce a scalable group inference method that improves both the diversity and quality of a group of samples. We formulate group inference as a quadratic integer assignment problem: candidate outputs are modeled as graph nodes, and a subset is selected to optimize sample quality (unary term) while maximizing group diversity (binary term). To substantially improve runtime efficiency, we progressively prune the candidate set using intermediate predictions, allowing our method to scale up to large candidate sets. Extensive experiments show that our method significantly improves group diversity and quality compared to independent sampling baselines and recent inference algorithms. Our framework generalizes across a wide range of tasks, including text-to-image, image-to-image, image prompting, and video generation, enabling generative models to treat multiple outputs as cohesive groups rather than independent samples.

## 1 Introduction

Recent advances in generative models, such as diffusion models, have driven significant efforts in inference-time guidance and scaling techniques [1, 2, 3]. These methods effectively improve various aspects of output quality, such as alignment with text prompts or image aesthetics, and offer fine-grained controls over the output. However, much recent work primarily focuses on enhancing the quality of *individual* samples generated in isolation.

Yet, in real-world applications, users are often shown a group of samples rather than just one. For example, many text-to-image platforms [4, 5] display a grid of four to eight images per prompt by default, a practice that offers users crucial benefits: more diverse choices regarding layout, lighting, and style, and new inspirations and ideas for prompt refinement and local edits. This creates a gap between current research, focused on independent samples, and the practical need for diverse, high-quality groups in content creation workflows. How can we close this gap?

In this work, we propose a *scalable group inference* method to jointly improve the diversity and quality of a collection of generated samples. We formulate this task as a quadratic integer programming problem, representing output candidates as graph nodes. From a large set of $M$ candidates, we select a subset of size $K$ that maximizes a combination of individual sample quality, as a unary term, and group diversity as a binary term. However, a direct approach involves running the $T$-step denoising
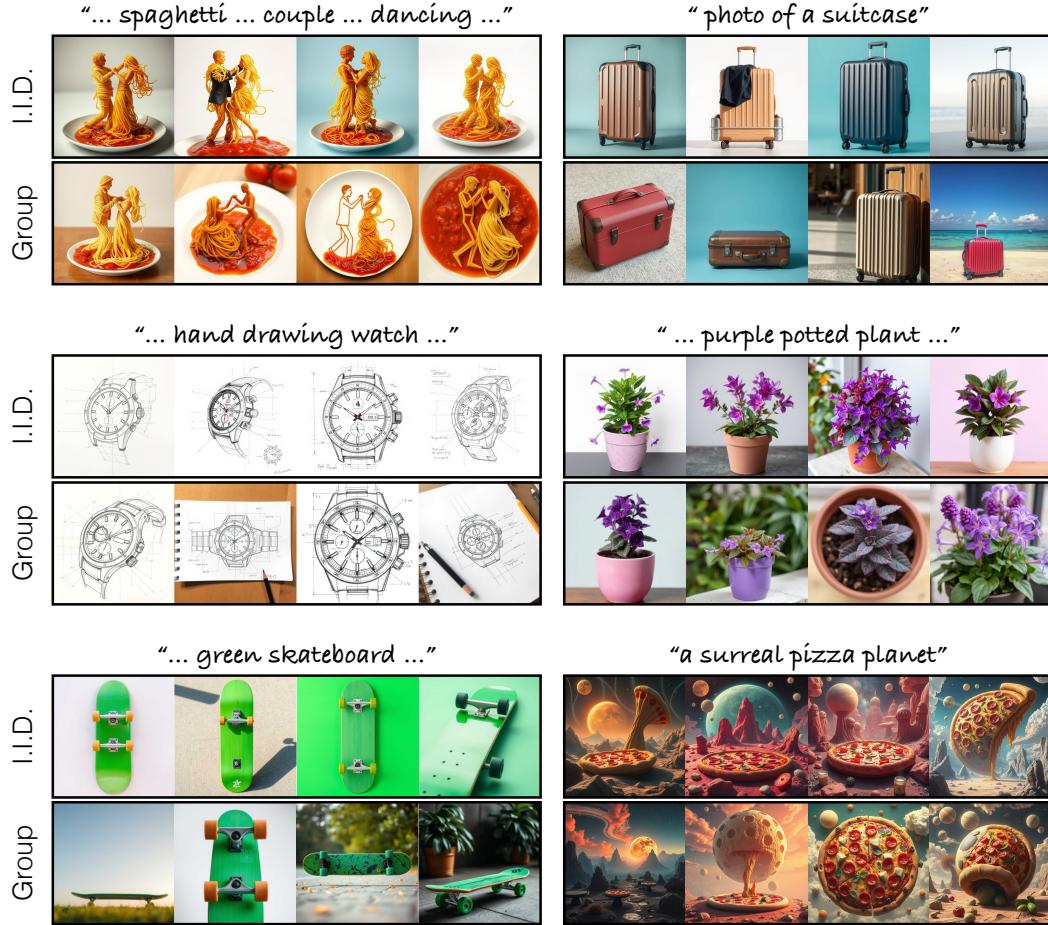
Figure 1: **Scalable Group Inference Results.** We show the advantage of our proposed group inference method over I.I.D. sampling. While I.I.D. sampling often yields repetitive results for the same prompt, our method generates a more diverse and high-quality collection of outputs. Please see our project website for more results.

process for all $M$ candidates, resulting in an $\mathcal{O}(MT)$ complexity. This is computationally expensive for large $M$ and $T$ (e.g., M=128, T=20). To address this, we introduce an efficient progressive selection strategy that leverages intermediate predictions during denoising to iteratively prune the candidate set. This approach is grounded in the insight that these intermediate predictions of the final output, despite originating from a long denoising chain, serve as effective *previews* of the final image at each step (Figure 3). This approach achieves a complexity of $\mathcal{O}(M + KT)$, where $K = 4$, enabling us to scale up our group inference to handle large candidate sets.

Extensive experiments have shown that our group inference method significantly outperforms independent sampling baselines and recent single-sample inference algorithms across various generative tasks and modalities, including text-to-image, image-to-image, image prompting, and video generation. Our method scales much better and produces more diverse and realistic outputs given the same compute budget. We further provide a comprehensive ablation study demonstrating the effectiveness of our design choices. Our framework enables generative models to treat multiple outputs as cohesive groups, aligning more closely with real-world workflows. In summary, our contributions are:.

- We propose a new, scalable group inference algorithm by selecting a group of K samples from M candidates as a quadratic integer programming problem to maximize sample quality and group diversity.
- We introduce a progressive pruning strategy to further scale our method. Our technique uses intermediate $\mathbf{x}_0$ predictions as previews to iteratively prune candidates, reducing complexity from $\mathcal{O}(MT)$ to $\mathcal{O}(M + KT)$, where $K$ is much smaller than $M$.

2

- Extensive evaluation on text-to-image, image-to-image, image prompting, and video generation shows our method outperforms baselines, producing more diverse and realistic outputs within similar cost budgets.

## 2   Related Works

**Diffusion models.**   are a powerful class of generative models that synthesize high-quality samples through iterative denoising [6, 7, 8]. Successful in text-to-image synthesis [9, 10, 11], their application later extends to video [12, 13, 14] and 3D synthesis [15, 16, 17]. However, common strategies to improve individual quality, such as fine-tuning for high quality, less diverse datasets [9] or strong classifier-free guidance (CFG) [1], often sacrifice diversity [18]. Additional conditioning methods like spatial controls [19] or image prompting [20] improve controllability but also reduce diversity, especially with strong guidance values. This lack of diversity is further worsened by one-step or few-step generators [21, 22, 23, 24]. Our work addresses this trade-off with group inference, enhancing both sample quality and diversity in batches, and demonstrating applicability across various controls (text, spatial, visual prompts).

**Diffusion Inference and Guidance.**   Inference-time guidance effectively improves sample quality and controllability of diffusion models without costly model finetuning. Early methods, such as classifier guidance [25] and widely-used classifier-free guidance (CFG) [1], significantly increase sample quality, often at the cost of diversity. Recent approaches manipulate internal representations, such as cross-attention maps [3], or incorporate spatial control from inputs such as layouts or sketches [26, 27, 28, 29, 30]. Other strategies apply guidance over limited intervals [31] or thresholding to CFG to reduce saturation [11].

While the above techniques focus on improving individual samples, our group inference approach explicitly optimizes collective properties, balancing single-sample quality and inter-sample diversity. A closely related work is particle guidance [32], which incorporates a pairwise potential during denoising steps to encourage diversity. Our method differs in three ways. First, our method improves both quality and diversity, while particle guidance often hurts image quality, as shown in experiments (Section 4.2). Second, our method scales effectively to a large number of images through early candidate pruning and sample selection, avoiding expensive optimization. In contrast, particle guidance is limited to small sets (e.g., four images) due to memory-intensive gradient computation of the pairwise terms. Third, our framework supports non-differentiable quality and diversity terms, enabling the use of metrics derived from multimodal LLMs.

**Inference-time Scaling.**   Test-time scaling, leveraging methods like chain-of-thought [33], proposer and verifier [34], or multi-step reasoning, has become a key research area for large-language models [35]. The idea is to increase inference-time computation in exchange for improved performance from a pre-trained model. Recently, researchers have adopted the inference-time scaling for diffusion models [2], which uses off-the-shelf models and evaluation metrics to search for better noises and increase the sample quality, often requiring thousands to tens of thousands of function evaluations (NFEs). However, text-to-image models differ from LLMs in three ways: they are often more computationally expensive [36], users often pay 5 to 10 cents per image on leading platforms, and users demand low latency. In our work, we show that our test-time scaling method balances the computational cost and quality and diversity improvement.

## 3   Method

We propose *Scalable Group Inference*, a test-time selection framework that chooses a diverse, high-quality subset from a large pool of generated outputs. The method relies on a scoring objective that combines a unary term that measures the quality of an individual sample and a binary term that computes pairwise properties such as image distances. We first formulate this as a quadratic integer programming problem (QIP) over binary selection variables. Then, to reduce compute cost, we introduce a progressive filtering strategy that prunes low-quality candidates early using intermediate predictions from partially denoised samples. We now describe both components in detail.
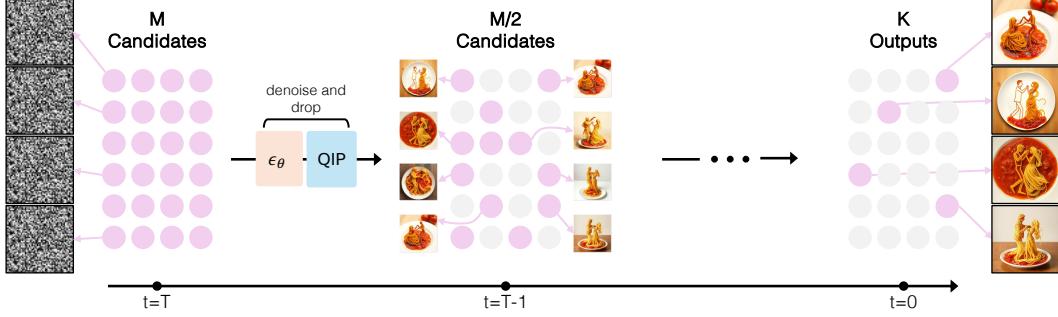
Figure 2: **Overview.** Given a large number of $M$ candidate noises, we gradually reduce the number of candidate sets through iterative denoising and pruning steps. At each step, we first leverage the diffusion model $\epsilon_\theta$ to denoise the sample. We then compute the quality metric (unary term) and pairwise distances (binary term), and solve a quadratic integer programming (QIP) program to progressively prune the candidate set, yielding a final group of $K$ diverse and high-quality outputs.

## 3.1 Formulation

Given a generative model $G_\theta(\mathbf{z}, \mathbf{c})$ that maps latent noise $\mathbf{z} \sim p(\mathbf{z})$ and condition $\mathbf{c}$ to outputs $\mathbf{x}$, our goal is to obtain a *set* of $K$ outputs, $\{\mathbf{x}^{(i)}\}_{i=1}^K$, that exhibits both high quality and diversity together.

We begin by generating a large set of $M$ candidate outputs $\{\mathbf{x}^{(i)}\}_{i=1}^M$ using i.i.d. sampling:

$$\mathbf{x}^{(i)} = G_\theta(\mathbf{z}^{(i)}, \mathbf{c}), \quad \mathbf{z}^{(i)} \overset{\text{i.i.d.}}{\sim} p(\mathbf{z}). \tag{1}$$

Let $\mathcal{I} = \{1, \dots, M\}$ index the candidate samples. We associate each sample $i \in \mathcal{I}$ with a unary score $\mathbf{u}_i \in \mathbb{R}$ (e.g., CLIPScore [37] between image CLIP embedding and input caption text embedding) and each pair $(i, j)$ with a binary score $\mathbf{b}_{ij} \in \mathbb{R}$ (e.g., DINO [38] distances between two images). Concretely,

$$\mathbf{u}_i = f_{\text{CLIP}}(\mathbf{x}^{(i)}, \mathbf{c}) \tag{2}$$

$$\mathbf{b}_{ij} = 1 - \text{cosine}\left(f_{\text{DINO}}(\mathbf{x}^{(i)}), f_{\text{DINO}}(\mathbf{x}^{(j)})\right) \tag{3}$$

where $f_{\text{CLIP}}$ computes the similarity between the input image and the target caption, and $f_{\text{DINO}}$ is the DINOv2 feature extractor. Note that our method is general and accommodate many different choices of score functions as discussed later in Section 4.4.

We introduce binary selection variables $\mathbf{y}_i \in \{0, 1\}$ where $\mathbf{y}_i = 1$ indicates that candidate $i$ is included in the next group. We define the group selection objective as:

$$\max_{\mathbf{y} \in \{0,1\}^M} \quad \sum_{i \in \mathcal{I}} \mathbf{u}_i \, \mathbf{y}_i + \lambda \sum_{\substack{i,j \in \mathcal{I} \\ i < j}} \mathbf{b}_{ij} \, \mathbf{y}_i \, \mathbf{y}_j$$

$$\text{subject to} \quad \sum_{i \in \mathcal{I}} \mathbf{y}_i = K. \tag{4}$$

$\lambda$ is the hyperparameter that controls the relative weight between the unary and binary scores. The first term rewards individually strong outputs; the second promotes diversity by favoring dissimilar pairs. Solving this quadratic integer program (QIP) yields a subset of size $K$ with desirable group properties. We use the branch-and-cut algorithm implemented by an off-the-shelf solver [39] to solve the QIP. Note that the formulation is model-agnostic and can accommodate any scoring functions, including functions that are not differentiable.

## 3.2 Progressive Pruning for Efficient Selection

Naively applying group selection requires generating all $M$ candidates to completion, which is prohibitively expensive for recent compute-intensive models like Flux [40]. For example, generating $M = 64$ samples over $T = 20$ denoising steps requires $M \cdot T$ forward passes. Even on a modern GPU like NVIDIA H100, this results in a runtime of more than 3 minutes. To reduce this cost, we introduce a progressive filtering strategy that prunes candidates early using intermediate predictions.
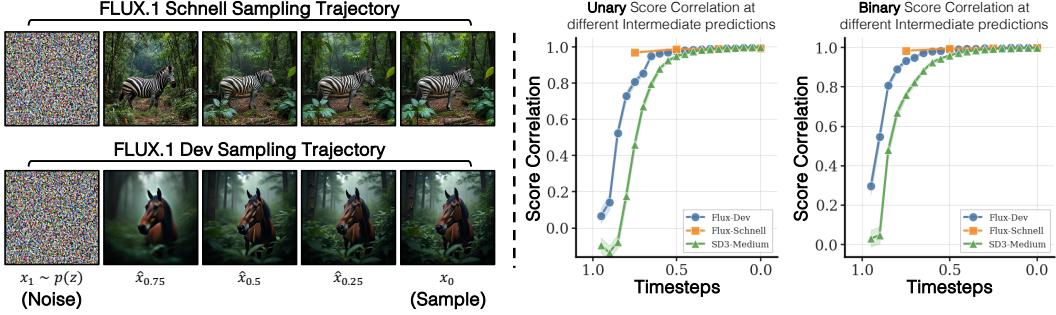
4

Figure 3: **Correlation Between Intermediate and Final Generation Scores.** On the left, we show the reverse diffusion process, visualizing the intermediate predictions $\hat{\mathbf{x}}_t$ of the final image at different steps for FLUX.1 Schnell and FLUX.1 Dev models. We can observe that the intermediate predictions look similar to true final sample $\mathbf{x}_0$ for both the models. We further demonstrate this quantitatively by plotting the Spearman correlation of the Unary and Binary scores from $\hat{\mathbf{x}}_t$ versus final $\mathbf{x}_0$ scores, across different steps. For multistep models like FLUX.1 Dev and Stable Diffusion 3, the plots demonstrate strong correlations rapidly approaching 1.0, even at early timesteps. Note that for a timestep distilled model like Flux-Schnell, the correlation is high from the first denoising step. This highlights the utility of using intermediate predictions for progressively filtering candidate samples.

**Intermediate pruning.** We maintain a set $\mathcal{S}_t \subseteq \mathcal{I}$ of candidate indices at each step $t$. For each sample in $\mathcal{S}_t$, we compute the intermediate prediction $\hat{\mathbf{x}}_t$, evaluate the unary and binary scores, and solve the QIP (Eq. 4) to select the best subset. This subset becomes the next set $\mathcal{S}_{t-1}$, forming a nested sequence:

$$\mathcal{S}_T \supset \mathcal{S}_{T-1} \supset \cdots \supset \mathcal{S}_0.$$

Once the set reaches the desired output group size $K$, we stop the pruning and complete the remaining denoising steps only for the selected samples. See Algorithm 1 for the full procedure.

**Reliability of early predictions.** In modern multi-step diffusion and flow-based models, the intermediate state $\mathbf{x}_t$ already encodes coarse information about the final generated sample $\mathbf{x}_0$. A common approximation of the final image at timestep $t$ is the predicted reconstruction:

$$\hat{\mathbf{x}}_t = \mathbf{x}_t + t \cdot \epsilon_\theta(\mathbf{x}_t, t, \mathbf{c}), \tag{5}$$

where $\epsilon_\theta$ predicts the noise or velocity at time $t$. Although these predictions are coarse, they are sufficient for computing the unary and binary scoring functions.

To quantify this, we compute the correlation between the scoring functions (e.g., CLIP similarity or pairwise DINO diversity) evaluated on the intermediate images $\hat{\mathbf{x}}_t$ and the final output $\mathbf{x}_0$. Across a range of denoising steps, Figure 3 (right) shows strong correlations (e.g., $r > 0.7$ after 5 steps for multi-step models, and $r > 0.95$ after the first step for distilled models), indicating that intermediate predictions are reliable proxies. Figure 3 (left) shows this visually. This high correlation enables us to safely rank candidates before they are fully denoised.

### 3.3 Computational Complexity Analysis

To analyze the efficiency of progressive filtering, suppose we start with $M$ candidate samples and prune the set by a fixed ratio $\rho \in (0, 1)$ at each denoising step until reaching a target set size $K$. Thus, the number of candidates at timestep $t$ is given by:

$$|\mathcal{S}_t| = \max\left(\rho^t M, K\right). \tag{6}$$

Let $T$ denote the total number of denoising steps. We define the timestep $t^*$ at which the candidate set size first reaches or falls below the target $K$:

$$t^* = \left\lceil \frac{\log(K/M)}{\log(\rho)} \right\rceil. \tag{7}$$

The total number of model evaluations $f_\theta$ required throughout the process can be written as:

$$M \cdot \frac{1 - \rho^{t^*}}{1 - \rho} + K \cdot (T - t^* + 1). \tag{8}$$

5

**Algorithm 1** Efficient group inference

```
# model: The diffusion model ε_θ
# zs: Initial noise vectors {z_i}
# N: Total number of denoising steps
# ts: The noise schedule {t_j}
# K: The target number of samples
# c: Conditioning information
# rho: The dropping ratio ρ for pruning
def group_inference(model, zs, N, ts, K, c, rho):
    # Initialize the set of candidates from noise
    candidate_set = list(zs)

    # Denoising loop
    for j in reversed(range(1, N +1)):
        intermediate_previews, next_latents = [], []

        # Get intermediate previews for all candidates
        for x_t in candidate_set:
            preview, x_next = denoise(x_t, ts[j], ts[j-1], c, model)
            previews.append(preview)
            next_latents.append(x_next)

        if len(candidate_set) > K:
            # Score previews and select the best subset
            u = unary_score(previews)
            b = binary_score(previews)

            # Prune candidates based on the dropping ratio rho
            m = max(K, int(len(candidate_set) * (1 - rho)))

            indices = SolveQIP(u, b, m)
            candidate_set = [next_latents[i] for i in indices]
        else:
            candidate_set = next_latents

    return candidate_set
```

In contrast, naive sampling without pruning would require $M \cdot T$ model evaluations. For typical parameter settings (e.g., $M = 64$, $K = 4$, $\rho = 0.5$, $T = 20$), our progressive filtering approach yields substantial compute savings (i.e., 184 vs 1280 evaluations, $\sim 85\%$ reduction). Our method has an overall complexity of $\mathcal{O}(M + KT)$.

## 4 Experiments

We demonstrate the effectiveness of our proposed scalable group inference method across three different tasks: text-to-image generation, depth-conditioned generation, encoder-based image customization, and five different base models: FLUX.1 Schnell, FLUX.1 Dev, Stable Diffusion 3 (Medium), FLUX.1 Depth, and SynCD. The dataset and evaluation protocols used throughout all experiments are described next in Section 4.1. Subsequently, in Section 4.2 and Section 4.3, we compare against prior methods along two axes: diversity-quality tradeoff, and inference-time scalability with different compute budget constraints. Finally, we present an ablation study to analyze the different components of our method, a runtime analysis, and their respective contributions. Please see the Appendix B, A, C for more results, analysis and ablations.

### 4.1 Dataset and Evaluation

**Datasets.** We use the GenEval dataset [41], validation split of the COCO 2017 dataset [42], and DreamBooth dataset [43] for text-to-image generation, depth-conditioned generation, and image
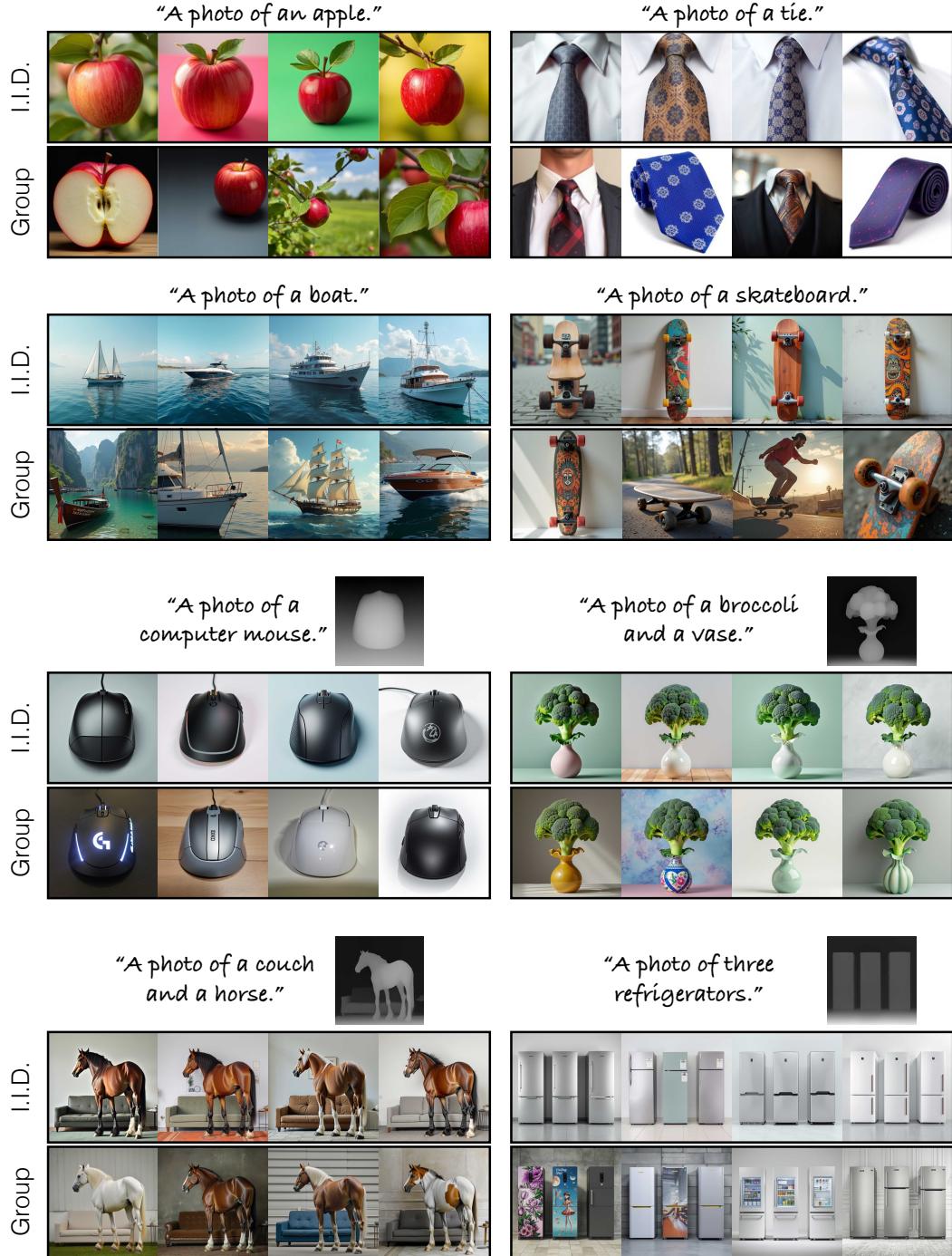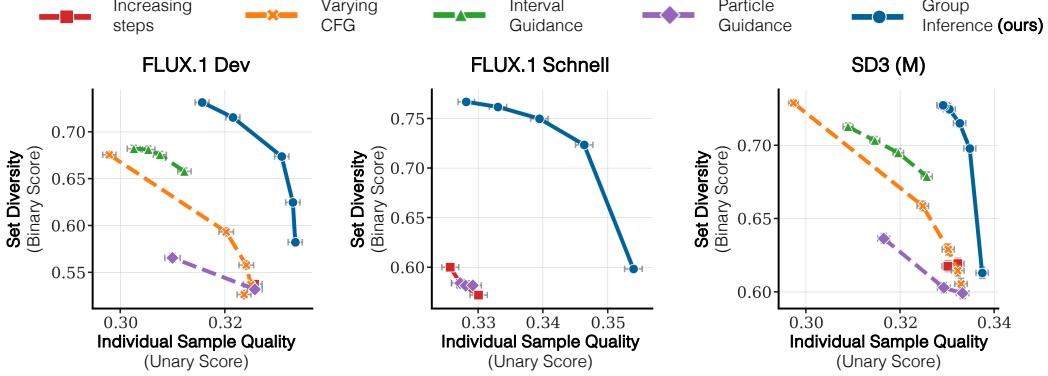
Figure 4: **Gallery of Results.** Qualitative results that show the advantage of our proposed group inference method over I.I.D. sampling for text-to-image generation and depth-to-image generation. Top row shows results with FLUX.1 Schnell, the second row uses FLUX.1 Dev, and the last two rows use FLUX.1 Depth as the base model. For text-to-image generation, our method produces more diverse object poses and orientations, while for depth-to-image generation, it enhances color and texture diversity while adhering to the input depth condition.

customization, respectively. For depth-conditioned generation, we first extract the depth map using a recent method [44]. Please refer to Appendix C for additional details.

Figure 5: **Quality and Diversity Pareto front for text-to-image models.** Each curve corresponds to a different inference strategy for three different text-to-image models (FLUX.1 Dev, FLUX.1 Schnell, and Stable Diffusion 3 Medium). Our proposed Group Inference (blue) consistently dominates alternate methods (Increasing steps, Varying CFG, Interval Guidance, and Particle Guidance) achieving Pareto optimality and superior tradeoffs between quality and diversity across all methods. Varying CFG and Interval Guidance do not apply to the distilled model (FLUX.1 Schnell).

**Models.** We use several recent models, including FLUX.1 Dev [40] and Stable Diffusion 3 Medium (SD3-M) [45], which are flow-based models typically requiring 20-50 denoising steps. We also evaluate FLUX.1 Schnell, a timestep-distilled variant designed for efficient generation, typically using 1-8 steps. For depth-conditioned generation, we use FLUX.1 Depth, a model specifically trained for structural guidance based on depth maps. For customization, we use SynCD [46], a recent encoder-based image prompting model. Unless otherwise specified, the sampling parameters for these models are fixed to the default values. A comprehensive list and inference parameters used can be found in Appendix C.1.

**Score Functions.** We use CLIP text-image similarity [47] to assess the quality of the individual samples (unary score) for the text-to-image and depth-to-image generation. For encoder-based image customization, we use cosine DINOv2 [38] similarity between the input subject image and the output generated images for the unary score. Diversity (binary score) is computed for all tasks as one minus the cosine similarity between the DINOv2 patchwise features of all image pairs in the output set. Our method can naturally accommodate a wide range of unary and binary scores, fitting the user's needs. Section 4.4 demonstrate this concretely.

**User Study.** We conduct two user preference studies to compare our method against each baseline on text-to-image generation. The first user study evaluates output diversity. In each comparison, the users are presented with two sets of 4 output images generated by two methods. The users are instructed to choose the set that has the higher variety. The second user study evaluates individual sample quality. For this study, the users are shown two images generated by two methods and asked to pick the one with higher quality. Both studies were conducted using Amazon Mechanical Turk (AMT) using prompts from our entire validation set. Each comparison was rated by three unique users, resulting in a total of 23,226 preference judgments.

**Runtime.** We measure inference runtime using wallclock time. Specifically, this is the time taken by each method to generate an output set of $K$ images (where $K = 4$, unless specified otherwise) from a given input condition (i.e., a text prompt, depth map, or subject image). This measurement excludes initial model loading times and is averaged over 20 independent runs for each reported value. All runtime experiments utilize a single NVIDIA H100 GPU. Runtime comparisons based on the number of function evaluations (NFEs) are included in the appendix.

**Uncertainty Estimation.** We report standard errors for all quantitative results presented throughout our experiments. These standard errors are computed via bootstrapping with 1000 resamples.

## 4.2 Baselines and the Diversity-Quality Tradeoff

In generative modeling, a fundamental tradeoff often exists between optimizing for the perceptual quality of individual samples and ensuring a diverse set of outputs [48, 1, 49]. Many prior methods implicitly or explicitly navigate this spectrum. In this subsection, we compare our proposed approach

| Model | Comparison | Diversity | | Quality | |
|---|---|---|---|---|---|
| | | Ours pref. | Baseline pref. | Ours pref. | Baseline pref. |
| FLUX.1 Dev | Ours vs Low-CFG | **88.3%** | 11.70% | **85.6%** | 14.4% |
| | Ours vs Interval Guidance | **53.4%** | 46.6% | **58.4%** | 41.6% |
| | Ours vs Particle Guidance | **81.2%** | 18.8% | **79.4%** | 20.6% |
| FLUX.1 Schnell | Ours vs Low-CFG | | | N/A | |
| | Ours vs Interval Guidance | | | N/A | |
| | Ours vs Particle Guidance | **55.5%** | 44.5% | **62.3%** | 37.7% |
| SD3 (M) | Ours vs Low-CFG | **76.8%** | 23.2% | **80.8%** | 19.20% |
| | Ours vs Interval Guidance | **58.1%** | 41.9% | **57.9%** | 42.10% |
| | Ours vs Particle Guidance | **78.9%** | 21.1% | **85.9%** | 14.1% |

Table 1: **User preference comparison between our method and baselines.** Results from our user study demonstrate that our method is consistently preferred over alternative inference strategies. Across three different text-to-image models (FLUX.1 Dev, FLUX.1 Schnell, and Stable Diffusion 3 Medium), users consistently chose our generations for both diversity and quality. Note that comparisons against Low-CFG and Interval Guidance are not applicable (N/A) for the distilled FLUX.1 Schnell model.
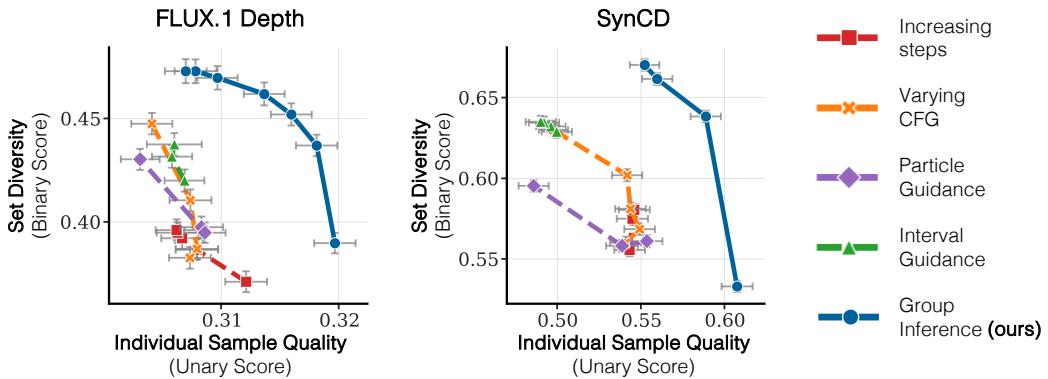


Figure 6: **Quality and Diversity Pareto front for additional tasks.** Each curve corresponds to a different inference strategy for depth conditioned generation (left, FLUX.1 Depth) and image prompting (right, SynCD). Our proposed Group Inference (**blue**) consistently dominates alternative methods—Increasing steps, Varying CFG, Interval Guidance, and Particle Guidance—achieving Pareto optimality and superior tradeoffs between quality and diversity across all methods.

to existing approaches proposed to achieve a more favorable Pareto frontier in the diversity-quality space. Figure 5 plots the quality and diversity Pareto front for text-to-image models (FLUX.1 Dev, FLUX.1 Schnell, and SD3 (M)). Table 1 shows the results of a pairwise user preference study between our method and the baselines. Figure 6 plots the quality diversity tradeoff for additional models. Qualitative comparisons are shown in Figure 7.

Across all baselines, our proposed method consistently achieves a superior diversity-quality tradeoff. As illustrated by the **blue** line in Figures 5 and 6, our approach dominates the Pareto fronts of all evaluated baselines, yielding better diversity for a given level of quality, or higher quality for a comparable level of diversity. Comparisons with additional metrics are shown in the Appendix.

**Increasing Denoising Steps.** We first consider the impact of simply increasing the number of denoising steps during sampling. While more steps can sometimes refine details, we find this has a minimal effect on meaningfully shifting the diversity-quality balance for the models under study, as shown by the **red** line in Figures 5 and 6.

**Varying CFG.** Next, we examine the widely used technique of varying the Classifier-Free Guidance scale [1] (CFG). As depicted by the **orange** line in Figures 5 and 6, systematically altering the CFG scale traces a distinct tradeoff curve. Notably, low CFG values (e.g., CFG=1) largely increase output diversity but often at the cost of a sharp degradation in sample quality and prompt alignment. This is
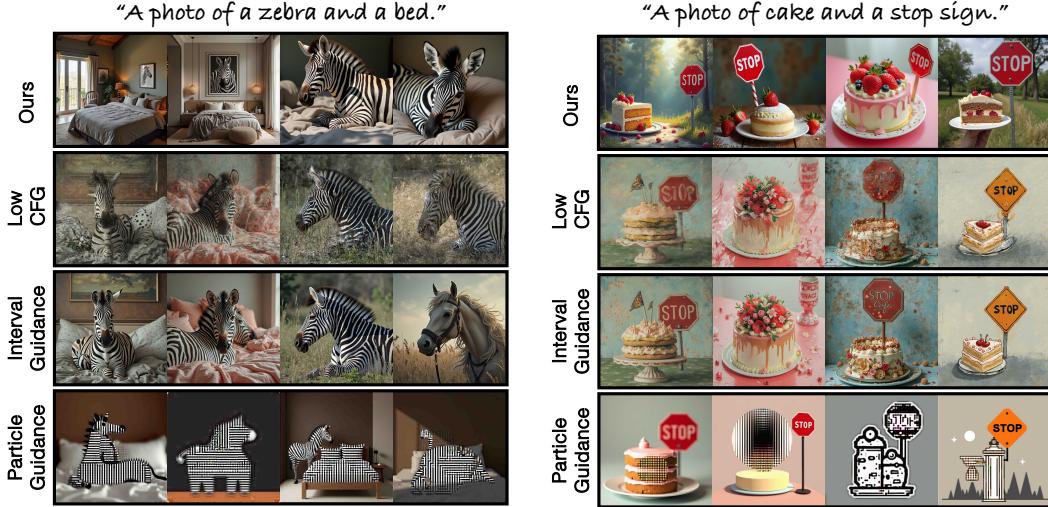
Figure 7: **Qualitative results.** We compare our proposed method (top row) against alternative inference strategies targeting an improved Quality-Diversity tradeoff with FLUX.1 Dev base model. To ensure a fair comparison, baseline methods were configured to approximate the diversity level achieved by our approach. The precise parameters of each baseline, and a comparison at other configurations is shown in the supplement. The result demonstrates that: (i) employing a low Classifier-Free Guidance (CFG) scale to increase diversity results in diminished image quality; (ii) Interval Guidance exhibits reduced adherence to the input text prompt; and (iii) Particle Guidance, by actively altering sampling trajectories, tends to produce less natural images. In contrast, our method outputs a set of diverse outputs while maintaining good image quality and prompt fidelity.
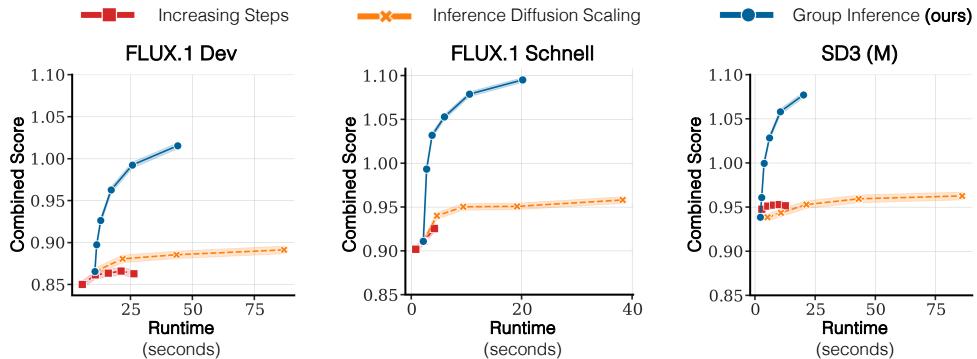


Figure 8: **Performance at different runtimes.** The Increasing Steps (**red**) baseline shows limited gains with additional computation. The Inference Diffusion Scaling [2] method (**orange**), which increases sample count through independent I.I.D. generations, requires substantially more runtime for marginal improvements. In contrast, our proposed Group Inference (**blue**) achieves significantly better performance–runtime tradeoffs, quickly outperforming both baselines with minimal overhead. Observation holds across models.

visually seen through the poor image quality in the second row in Figure 7 where a low CFG value is used. The results of the user study in Table 1 further corroborate these observations.

**Interval Guidance.** Interval guidance [31] attempts to refine this by applying CFG selectively only during a subset of the denoising timesteps. We conduct a sweep across various interval configurations, with the results shown as the **green** line in Figures 5 and 6. Consistent with the original findings for Interval Guidance and third row in Figure 7, this approach can offer image quality improvements over a standard CFG sweep. However, it still performs worse than our method in terms of both quality and diversity. For instance, in the example of zebra and bed on the left in Figure 7, Interval Guidance has reduced diversity of zebra poses and does not generate a bed for two of the four outputs. Similarly, in the right example, interval guidance does not always generate a stop sign. Moreover, both standard
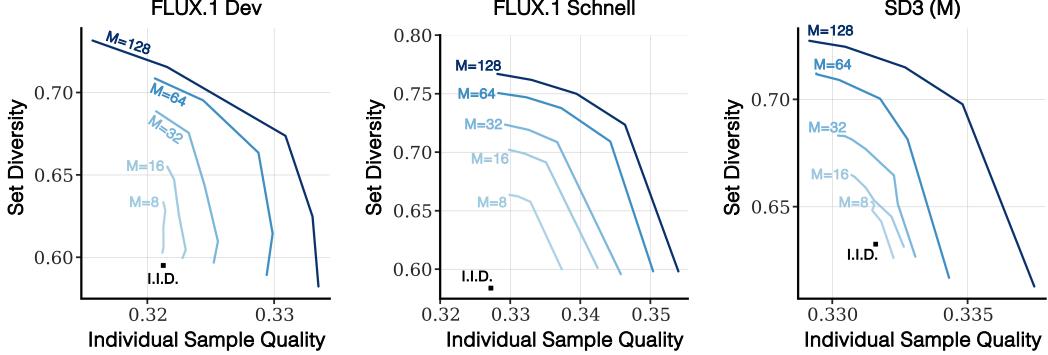
Figure 9: **Improvements as the number of initial samples M is increased.** We show how the sample quality (measured with CLIP) and the set diversity (measured with DINO) improves as the number of initial starting samples is increased from 4 to 128.

CFG sweeping and Interval Guidance do not apply to distilled models like FLUX.1 Schnell, as these models do not use guidance mechanisms.

**Particle Guidance.** We also evaluate Particle Guidance [32], which optimizes a binary potential function to encourage diversity during inference. Following the original work, we use DINO features for the diversity term. As illustrated by the **purple** line, Particle Guidance can indeed increase output diversity. However, this comes with a sharp decrease in individual sample quality (fourth row in Figure 7), as direct optimization of the binary potential actively alters the sampling trajectory. This can push the output samples off the learned data manifold, leading to less natural and artifact-prone images. Furthermore, Particle Guidance requires a substantial memory cost due to the necessity of computing gradients and backpropagating through the binary potential. This reliance on gradient computation also makes the method unsuitable for non-differentiable potential functions.

## 4.3 Inference Scaling Analysis

Scaling computational resources at the test time to enhance model performance is an increasingly useful paradigm in machine learning. For diffusion models, a native mechanism for test-time scaling involves increasing the number of denoising steps. Although this can initially lead to improved sample quality, this approach often yields diminishing returns; beyond a certain point, additional denoising steps provide progressively smaller gains in quality. In Figure 8, we illustrate the impact of various test-time scaling strategies on the group objective (defined in Equation 4), evaluated across various computational budgets.

Our first baseline (Figure 8, **red** line) allocates increased compute to a greater number of denoising steps for a fixed number of initial samples $M$. Consistent with existing findings [2], this approach demonstrates minimal improvement in our combined score, with the curve quickly plateauing. Inference Diffusion Scaling [2], proposes an alternative that utilizes the additional compute budget to perform a search over multiple random seeds. For a fair comparison, we implement this baseline with the CLIP text-image similarity as the verifier. This method does not incorporate intermediate predictions and does not consider any pairwise terms. Consequently, it is not effective in improving the group objective, as shown through the **orange** line.

In contrast, our method invests in the inference budget to increase the number of initial samples. As depicted by the **blue** line in Figure 8, this approach produces consistent improvements in the combined group score. Figure 9 further shows the improvement in both the quality and diversity of the outputs as the number of initial samples is gradually increased from 4 to 128 across three different base models.

**Ablating progressive filtering.** The importance of progressive filtering is shown in Figure 10. We compare our complete method, which utilizes progressive filtering with intermediate predictions $\hat{x}_t$ (**red** line), against a variant that performs full denoising for all $M$ candidate samples without such filtering (**gray** line). Demonstrating its effectiveness across different architectures, our approach achieved comparable group scores while requiring up to 73% less runtime. Please see the appendix for additional ablation studies.
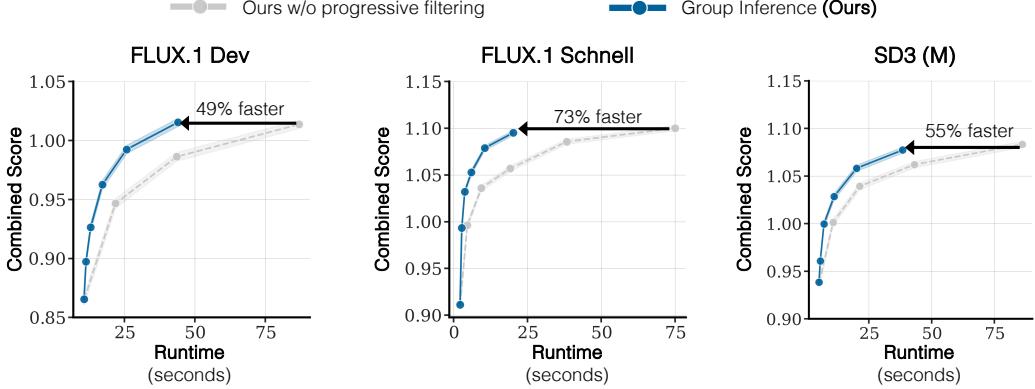
Figure 10: **Importance of progressive pruning.** Across multiple base generative models, progressive pruning consistently enables our method to select candidates efficiently and shows substantial speedups-49%, 73%, and 55% faster for comparable combined group scores.
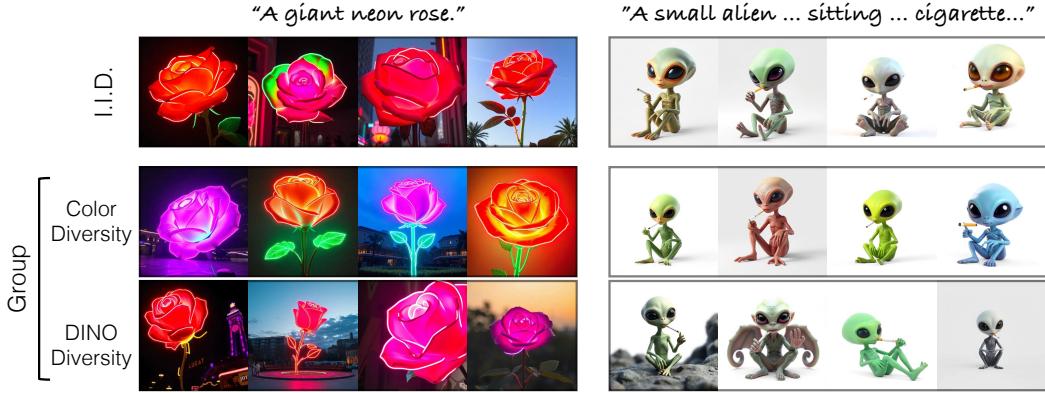


Figure 11: **Accommodating different pairwise objectives.** Compared to baseline I.I.D. sampling (top row), our method allows for targeted diversity by defining different pairwise objectives. The second and third rows show results where the unary quality term is identical but the pairwise binary term is varied. The middle row uses a color-based binary term, while the bottom row uses a DINO-based binary term to achieve semantic and structural diversity.

## 4.4 Different Diversity Objectives

Next, we show that our approach is general and can accommodate different pairwise binary objectives by simply swapping the binary term in our quadratic integer programming objective, as demonstrated in Figure 11. Let us consider the example shown on the left corresponding to the caption "a giant neon rose." Standard I.I.D. sampling (top row) produces a set of visually redundant images. All four roses are red and share a similar pose.

In contrast, the bottom two rows are generated using our method with an identical unary quality term (CLIP text-image similarity) but different binary diversity objectives. The middle row uses a direct color based dissimilarity as the binary term. This successfully steers the outputs towards a set of images with varied and distinct color schemes. For the rose example on the left, this results in a vibrant set that includes blue, orange, and pink neon variants. In the bottom row, we use a DINO diversity metric that captures more semantic features when comparing the pairwise distances. This change directs the model to produce a set with higher structural variance. As seen with the rose example, this yields outputs with different poses and camera angles.

This direct comparison underscores a key strength of our approach: the ability to seamlessly integrate different notions of diversity to achieve targeted, user-defined visual outcomes.

Figure 12: **Failure cases.** The performance of our method depends on the diversity of the initial candidate pool. (Left) For the prompt "A photo of a Ferrari," the base model (FLUX.1 Schnell) exhibits a strong color bias, exclusively generating red cars. Consequently, our method can find varied poses but is unable to produce a color-diverse set. (Right) Similarly, for "a photo of Albert Einstein," the base model only generates black-and-white images, constraining our method from finding any color photographs.

# 5 Discussion, Broader Impacts, and Limitations

In this paper, we have introduced scalable group inference, a novel method to generate diverse, high-quality sets of samples by formulating the selection as a quadratic integer program and leveraging intermediate predictions for improving the runtime efficiency. Our efficient approach significantly enhances group diversity and quality compared to existing baselines across various generative tasks. Still, our method has several limitations.

First, our method relies on the base generative model's ability to produce a sufficiently diverse and high-quality initial candidate pool. Consequently, if the underlying model generates outputs of inherently poor quality or suffers from significant mode collapse, the efficacy of Scalable Group Inference in identifying an optimal set will be inherently constrained, as our method selects from, rather than intrinsically enhances, these initial candidates. This is visually illustrated in Figure 12.

Second, our method assumes that the unary (quality) and binary (diversity) scores are fast to compute. If evaluating these scores, especially the pairwise diversity metric across a large candidate set, is computationally intensive, the runtime benefits of our scalable optimization would be reduced.

Nevertheless, our method offers a path to more user-centric systems that efficiently output diverse, high-quality sets of options, and enhance creative exploration. This capability can significantly reduce the iterative burden in content generation across various domains. Concurrently, the increased efficiency in generating diverse sets of synthetic media could also have potential for misuse, such as creating more varied and potentially harder-to-detect misleading content, demanding proactive ethical guidelines and mitigation strategies.

## References

[1] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.

[2] Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan Yang, Yandong Li, Tommi Jaakkola, Xuhui Jia, and Saining Xie. Inference-time scaling for diffusion models beyond scaling denoising steps. 2025.

[3] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11, 2023.

[4] Midjourney, Inc. Midjourney Website. `https://www.midjourney.com/home`, 2024.

[5] Adobe Inc. Adobe Firefly: Generative AI for Creatives. `https://www.adobe.com/products/firefly.html`, 2025.

[6] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.

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

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

[9] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[10] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024.

[11] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.

[12] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.

[13] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22563–22575, 2023.

[14] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.

[15] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.

[16] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 300–309, 2023.

[17] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023.

[18] Pietro Astolfi, Marlene Careil, Melissa Hall, Oscar Mañas, Matthew Muckley, Jakob Verbeek, Adriana Romero Soriano, and Michal Drozdzal. Consistency-diversity-realism pareto fronts of conditional image generative models. *arXiv preprint arXiv:2406.10429*, 2024.

[19] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023.

[20] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.

[21] Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rombach. Fast high-resolution image synthesis with latent adversarial diffusion distillation. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024.

[22] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation. In *European Conference on Computer Vision*, pages 87–103. Springer, 2024.

[23] Minguk Kang, Richard Zhang, Connelly Barnes, Sylvain Paris, Suha Kwak, Jaesik Park, Eli Shechtman, Jun-Yan Zhu, and Taesung Park. Distilling diffusion models into conditional gans. In *European Conference on Computer Vision*, pages 428–447. Springer, 2024.

[24] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6613–6623, 2024.

[25] Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat GANs on image synthesis. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.

[26] Minghao Chen, Iro Laina, and Andrea Vedaldi. Training-free layout control with cross-attention guidance. In *IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2024.

[27] Yunji Kim, Jiyoung Lee, Jin-Hwa Kim, Jung-Woo Ha, and Jun-Yan Zhu. Dense text-to-image generation with attention modulation. In *IEEE International Conference on Computer Vision (ICCV)*, 2023.

[28] Quynh Phung, Songwei Ge, and Jia-Bin Huang. Grounded text-to-image synthesis with attention refocusing. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.

[29] Andrey Voynov, Kfir Aberman, and Daniel Cohen-Or. Sketch-guided text-to-image diffusion models. In *ACM SIGGRAPH 2023 conference proceedings*, 2023.

[30] Yutong He, Ruslan Salakhutdinov, and J Zico Kolter. Localized text-to-image generation for free via cross attention control. *arXiv preprint arXiv:2306.14636*, 2023.

[31] Tuomas Kynkäänniemi, Miika Aittala, Tero Karras, Samuli Laine, Timo Aila, and Jaakko Lehtinen. Applying guidance in a limited interval improves sample and distribution quality in diffusion models. *arXiv preprint arXiv:2404.07724*, 2024.

[32] Gabriele Corso, Yilun Xu, Valentin De Bortoli, Regina Barzilay, and Tommi Jaakkola. Particle guidance: non-iid diverse sampling with diffusion models. *arXiv preprint arXiv:2310.13102*, 2023.

[33] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

[34] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.

[35] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.

[36] Muyang Li, Yujun Lin, Zhekai Zhang, Tianle Cai, Xiuyu Li, Junxian Guo, Enze Xie, Chenlin Meng, Jun-Yan Zhu, and Song Han. Svdqunat: Absorbing outliers by low-rank components for 4-bit diffusion models. *arXiv preprint arXiv:2411.05007*, 2024.

[37] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.

[38] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2023.

[39] Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2025.

[40] Black Forest Labs. Flux. `https://github.com/black-forest-labs/flux`, 2024.

[41] Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023.

[42] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014.

[43] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023.

[44] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR*, 2024.

[45] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.

[46] Nupur Kumari, Xi Yin, Jun-Yan Zhu, Ishan Misra, and Samaneh Azadi. Generating multi-image synthetic data for text-to-image customization. *ArXiv*, 2025.

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

[48] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. *Conference on Neural Information Processing Systems (NeurIPS)*, 30, 2017.

[49] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.

[50] Ollin Boer Bohan. Tiny autoencoder for stable diffusion. *Retrieved May*, 22:2024, 2023.

[51] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:15903–15935, 2023.

[52] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *Advances in Neural Information Processing Systems*, 37:21875–21911, 2024.

[53] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.

Section B presents additional qualitative and quantitative results obtained by our method across multiple different base generative models and tasks. Section A provides additional analysis of the different components of our method. Section C details the datasets and the implementation settings used for each of the baseline methods.

## A    Analysis

In this section, we provide additional analysis of the different components of our method.

**Runtime breakdown.**    In Figure 13 we show the runtime breakdown of different steps in the pipeline for the FLUX.1 Dev base model using CLIP text image similarity as the unary score and DINOv2 diversity as the binary score. On the left, we fix the output set size K to be 4 and increase the initial candidate size from 4 to 200. On the right, we fix the initial candidate size to be 200, and increase the output set size from 4 to 128. Note that across all settings, the runtime cost incurred by the QIP solver and the score computation is negligible compared to the forward pass of the denoising transformer.

**Different $\rho$ values.**    Figure 14 illustrates the effects of varying pruning ratios ($\rho$) on the FLUX.1 Dev and FLUX.1 Schnell models. The figure presents both the Number of Function Evaluations (NFE) (left plot) and the wallclock runtime on a single NVIDIA H100 (right plot). Across all plots, a pruning ratio of $\rho = 1.0$ signifies no progressive pruning. For the FLUX.1 Dev model, lower pruning ratios (e.g., $\rho = 0.1$ and $\rho = 0.25$) are overly aggressive, leading to suboptimal scores. Conversely, a pruning ratio of $\rho = 1.0$ (no candidate filtering) achieves a good combined score but incurs a high inference cost. A pruning ratio of $\rho = 0.5$ strikes an effective balance, yielding higher scores without excessive computational cost. We use $\rho = 0.5$ for all FLUX.1 Dev experiments.

A different trend observed for distilled FLUX.1 Schnell model. It can accommodate a more aggressive pruning ratio, such as $\rho = 0.1$, without a noticeable decrease in the score. This can be attributed to the better reliability of the intermediate predictions for the distilled models, as shown in Figure 3 of the main paper.

**Efficient decoding.**    Our method uses an efficient decoder [50] to decode all intermediate predictions for progressive pruning. In Figure 16 we ablate the use of efficient decoder and show that across both FLUX.1 Dev and FLUX.1 Schnell, using an efficient decoder improves the runtime without sacrificing the score.

**Evaluation with different score functions.**    Figure 5 in the main paper shows the quality and diversity Pareto front for the text to image generation task. That figure uses CLIP text-image similarity (Equation 2) as the quality score and DINO diversity 3 as the diversity score. Next, we evaluate our method using several additional score functions that are not used by our method for selection in Figure 17. The top row uses Image Reward [51] for measuring quality of samples and depth features to measure diversity. Image Reward is a network that is trained to learn human preferences for text-to-image generation. The depth diversity is calculated with the DepthAnything V2 model [52]. The bottom row uses BLIP2 [53] to measure the quality and CLIP features to compute the diversity. Figure 17 shows a comparison with three different base models: FLUX.1 Dev (left), FLUX.1 Schnell (middle), and Stable Diffusion 3 medium (right). Across each model, our proposed group inference shows a better trade off between quality and diversity. Note that particle guidance obtains slightly better BLIP2 score than our method for Stable Diffusion 3 (M). However, the outputs generated by particle guidance have artifacts. This is also reflected by a low score for other metrics (Image Reward and CLIP), and a worse user preference score.

## B    Additional Results

**Qualitative results.**    In Figures 21, 22, and 23 we show additional visual examples of our method. Across multiple models and tasks, our method consistently outputs samples that are more diverse, and without any degradation in the quality. Similar to the Figure 7 in the main paper, Figure 20 shows additional visual comparison to baselines. In these figures, the Low CFG baseline uses a CFG value of 1.0. Interval Guidance uses an interval of $[0.6, 0.4]$, and particle guidance uses a coefficient value of 100.

**Correlation analysis.**    Figure 3 in the main paper shows the correlation between the scores computed with the final image and the intermediate images. In Figure 18, we show that a similar correlation trend is visible in other base models (FLUX.1 Depth and SynCD). FLUX.1 Depth and SynCD show the CLIP text image similarity as the unary scores, and DINO diversity as the binary scores. This is consistent with our observations in the main paper.

## C    Implementation Details

Section C.1 first provides implementation details and hyperparameters used for all settings shown in Figures 4, 5, and 6 of the main paper. Section C.2 lists details about the datasets used for each task.
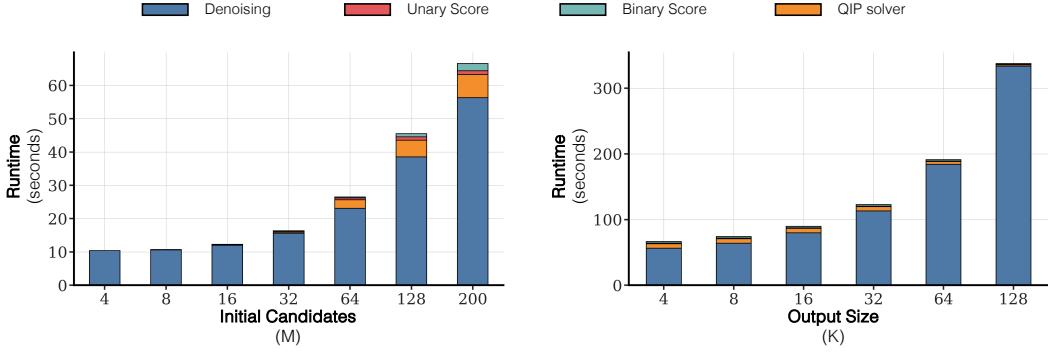
Figure 13: **Runtime breakdown.** We show a runtime breakdown of our method using FLUX.1 Dev model as the number of initial candidates (M, left) and the output set size (K, right) is increased. On the left plot, the output size is fixed to 4 and in the right plot the initial candidate size is fixed to 200. Across all settings, the runtime is dominated by the denoising step.

## C.1 Baselines

**Increasing steps.** For FLUX.1 Dev, Stable Diffusion 3 Medium, FLUX.1 Depth, and SynCD, we consider the timesteps 10, 20, 30, 40, and 50. For the distilled model, FLUX.1 Schnell, we consider the timesteps 1, 2, 4, and 8.

**Varying CFG.** For FLUX.1 Dev, FLUX.1 Depth, and SynCD we consider the CFG values 1, 2, 3, 4, and 5. For Stable Diffusion 3 Medium, we consider the CFG values 1, 5, 10, and 15. Note that FLUX.1 Schnell does not use CFG.

**Interval guidance.** For FLUX.1 Dev, FLUX.1 Depth, Stable Diffusion 3 Medium, and SynCD we consider the guidance intervals $[0.9, 0.1]$, $[0.8, 0.2]$, $[0.7, 0.3]$, and $[0.6, 0.4]$. Note that FLUX.1 Schnell does not use CFG, and this baseline is not applicable.

**Particle guidance.** For the Particle Guidance baseline, we consider coefficient values 0, 10, 50, 100, and 200. Note that this baseline significantly increases the memory consumption during inference.

**Inference diffusion scaling [2].** Figure 8 of the main paper shows a comparison to Inference Diffusion Scaling [2], a concurrent work, that shows an improvement in the quality of samples. We follow the results in their paper and use random search as the strategy. For a fair comparison, we use the same CLIP text-image-similarity as the verifier.

**Group inference (ours).** In Figures 5, 17, and 6 of the main paper, we vary the $\lambda$ defined in Equation 4 while keeping the input samples $M$ fixed. We use $M = 128$ for FLUX.1 Dev, FLUX.1 Schnell, SD3 (M), and SynCD. Note that varying the weighting factor $\lambda$ does not change the runtime, and only shows the trade-off between the diversity of samples in the generated output set and the individual quality. For FLUX.1 Dev, FLUX.1 Depth, SynCD we use $\rho = 0.5$. For SD3 (M) we use a higher $\rho = 0.75$, and for timestep distilled model FLUX.1 Schnell $\rho = 0.1$ in all experiments.

In Figures 8, 9 and 10 of the main paper, we want to study the performance at different runtimes, and therefore we fix the weighting factor $\lambda = 1$ but vary the number of input samples $M$ from 4 to 128.

**Choice of scores.** Unless specified otherwise, FLUX.1 Dev, FLUX.1 Schnell, SD3 (M), Flux.1 Depth use CLIP text-image similarity as the unary score, and DINO diversity as the binary score. SynCD uses DINO target image similarity as the unary score, and DINO diversity as the binary score across all results. Figure 11 keeps the unary score fixed as the CLIP text-image similarity and shows the effects of varying the binary score function.

## C.2 Dataset

**Text to image generation.** All text-to-image generation results with models FLUX.1 Dev, FLUX.1 Schnell, and SD3 (M) use all 553 prompts from the GenEval dataset [41].

**Depth to image generation.** All FLUX.1 Depth experiments use depth maps computed using Depth Anything Large [44] from 250 images from the validation split of the COCO 2017 dataset [42].

**Encoder-based image customization.** For all encoder-based image customization experiments using SynCD [46], we use 400 samples from the images in the standard DreamBooth dataset.
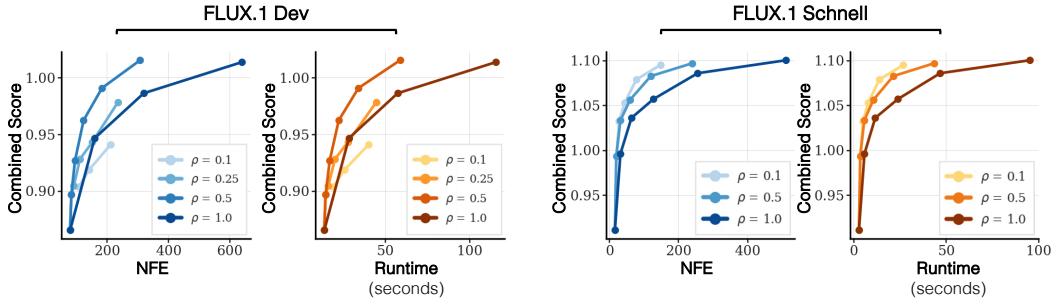
18

Figure 14: **Effects of different dropping ratio $\rho$.** We show the effects of different dropping ratios $\rho$ for two different base models: FLUX.1 Dev and FLUX.1 Schnell.
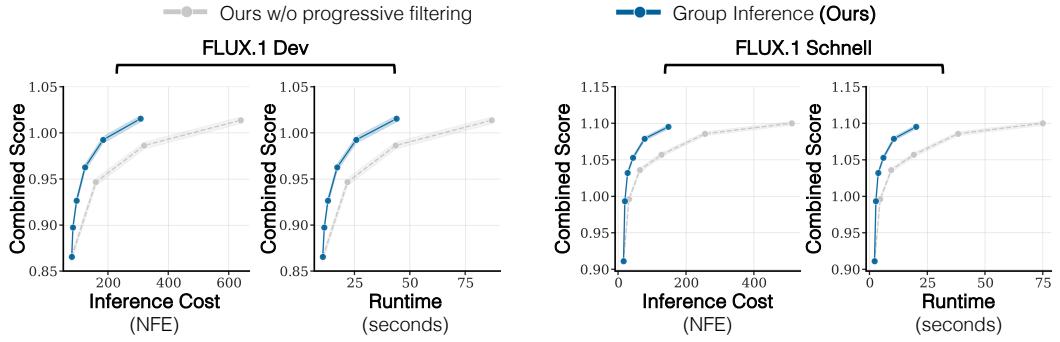


Figure 15: **Ablating the effect of progressive pruning.** Similar for Figure 10 from the main paper, we show the importance of progressive pruning. We report both, the number of function evaluations (NFEs) and the wallclock runtime (using one NVIDIA H100). The two plots on the left show the comparison using FLUX.1 Dev. The two plots on the right show FLUX.1 Schnell comparison.



Figure 16: **Ablating efficient decoder.** We show the effects of using an efficient decoder for decoding the intermediate predictions.

Figure 17: **Quality and Diversity Pareto front with additional metrics.** We evaluate the quality and diversity of samples generated by different inference strategies for three text-to-image models (FLUX.1 Dev, FLUX.1 Schnell, and Stable Diffusion 3 Medium). The top row shows evaluation using Image Reward [51] as the quality metric and Depth Diversity as the diversity metric. The bottom row uses BLIP2 [53] and CLIP Diversity. Note that these metrics are unseen and not used by our method.



Figure 18: **Correlation between intermediate and final generation Scores .** We follow the same protocol as Figure 3 in the main paper. On the left, we show the reverse diffusion process, visualizing the intermediate predictions $\hat{x}_t$ of the final image at different steps for FLUX.1 Depth and SynCD models. We can observe that the intermediate predictions look similar to true final sample $x_0$ for both the models. We further demonstrate this quantitatively by plotting the Spearman correlation of the Unary and Binary scores from $\hat{x}_t$ versus final $x_0$ scores, across different steps.
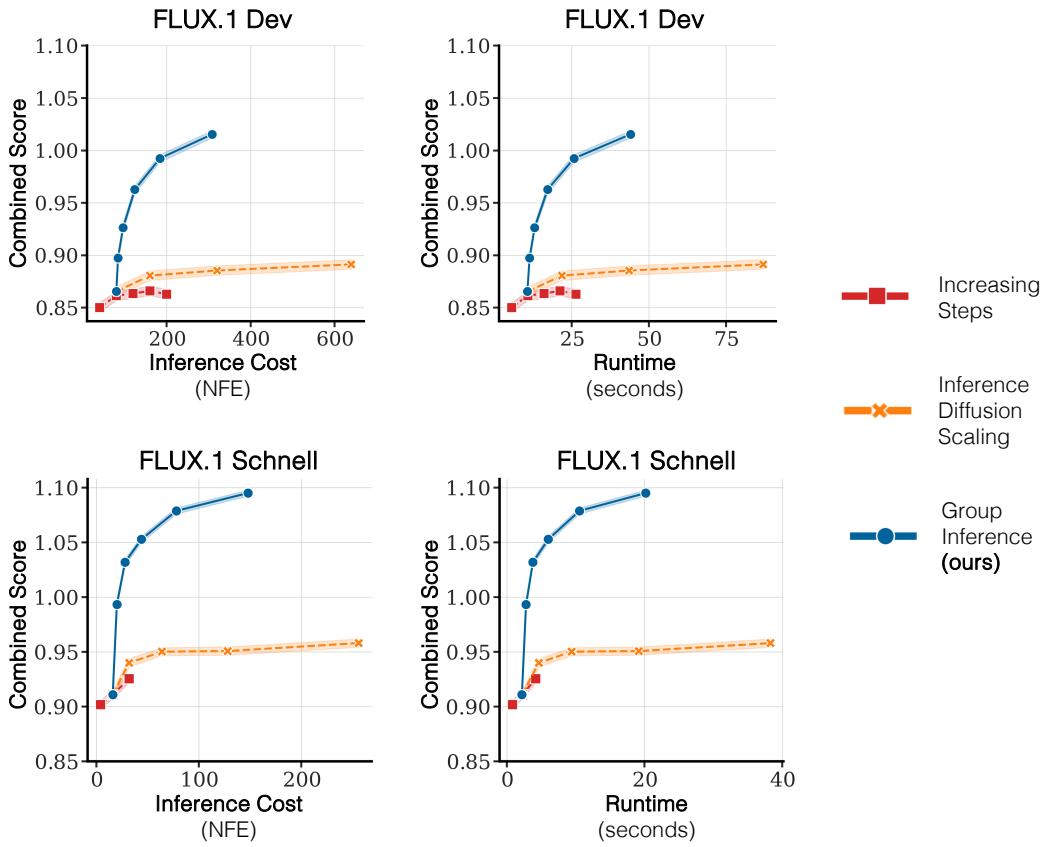
Figure 19: **Performance at different runtimes.** Similar for Figure 5 from the main paper, we show the different ways of allocating inference budget. We report both, the number of function evaluations (NFEs) and the wallclock runtime (using one NVIDIA H100).
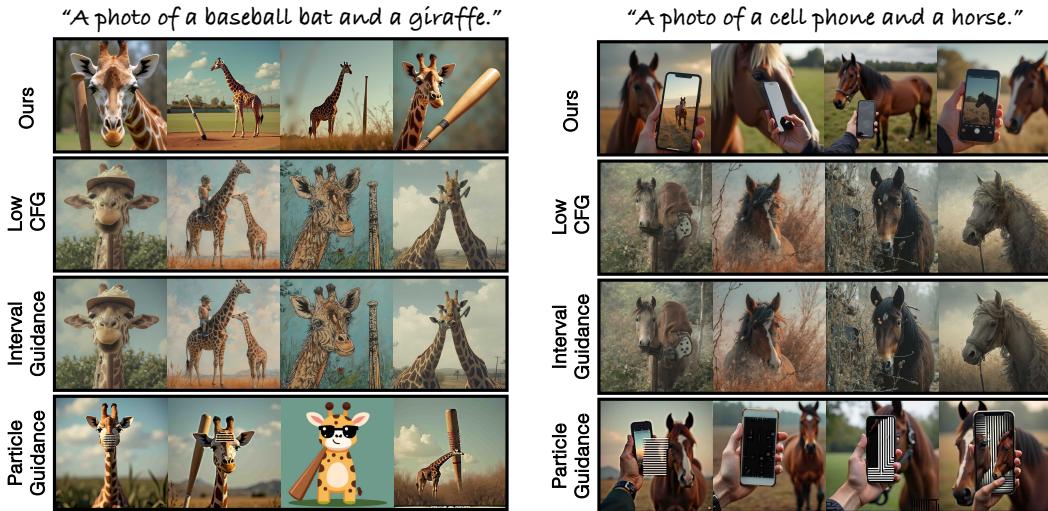


Figure 20: **Qualitative results.** We compare our proposed method (top row) against alternative inference strategies targeting an improved Quality-Diversity tradeoff with FLUX.1 Dev base model.
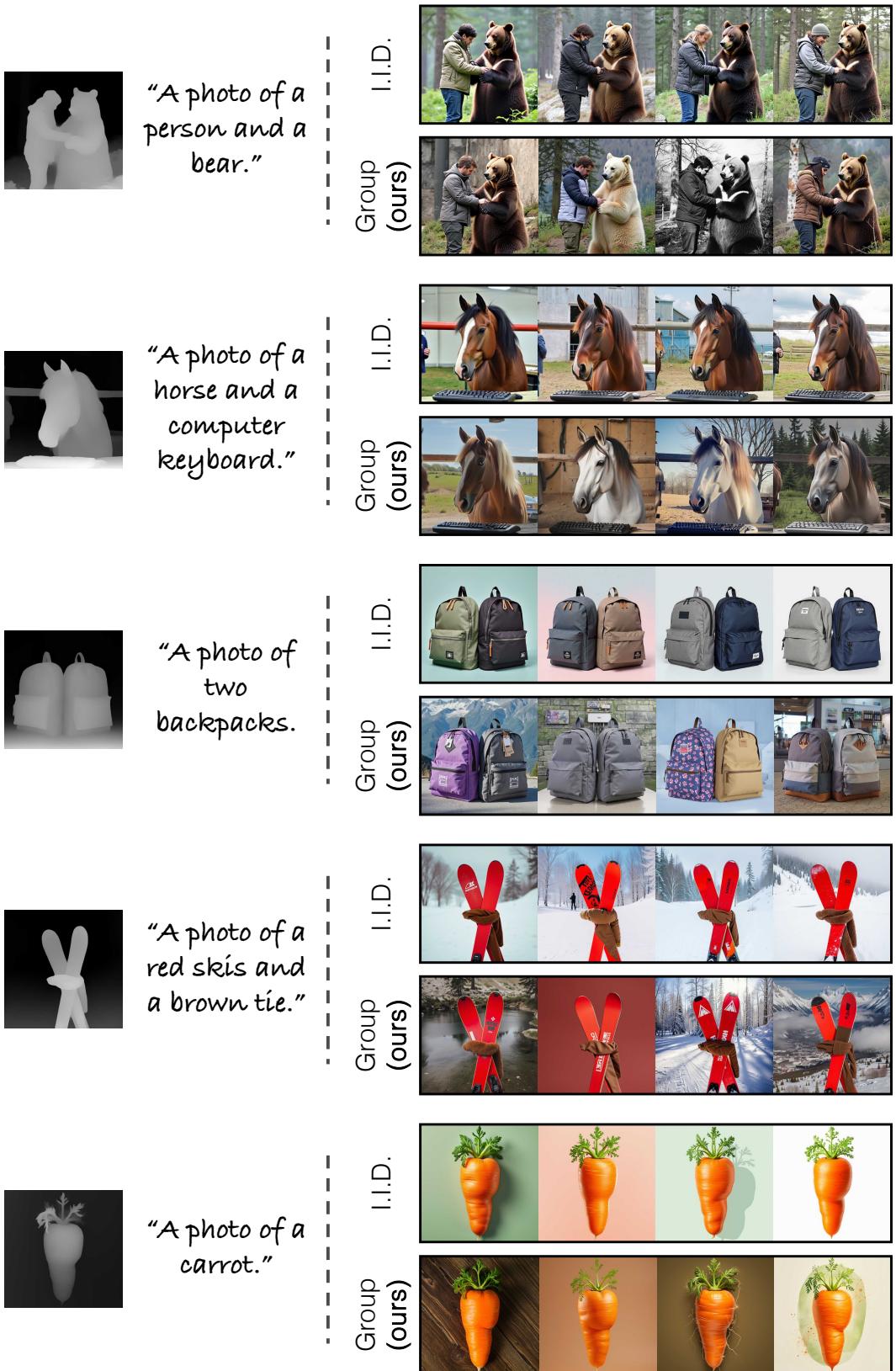
Figure 21: **Gallery of results.** Qualitative results that show the advantage of our proposed method over I.I.D. sampling for depth-to-image generation using FLUX.1 Depth as the base model. The input depth maps and captions are shown on the left and the generated outputs are shown on the right. Our method consistently generates outputs that have more diverse backgrounds, styles, and textures.
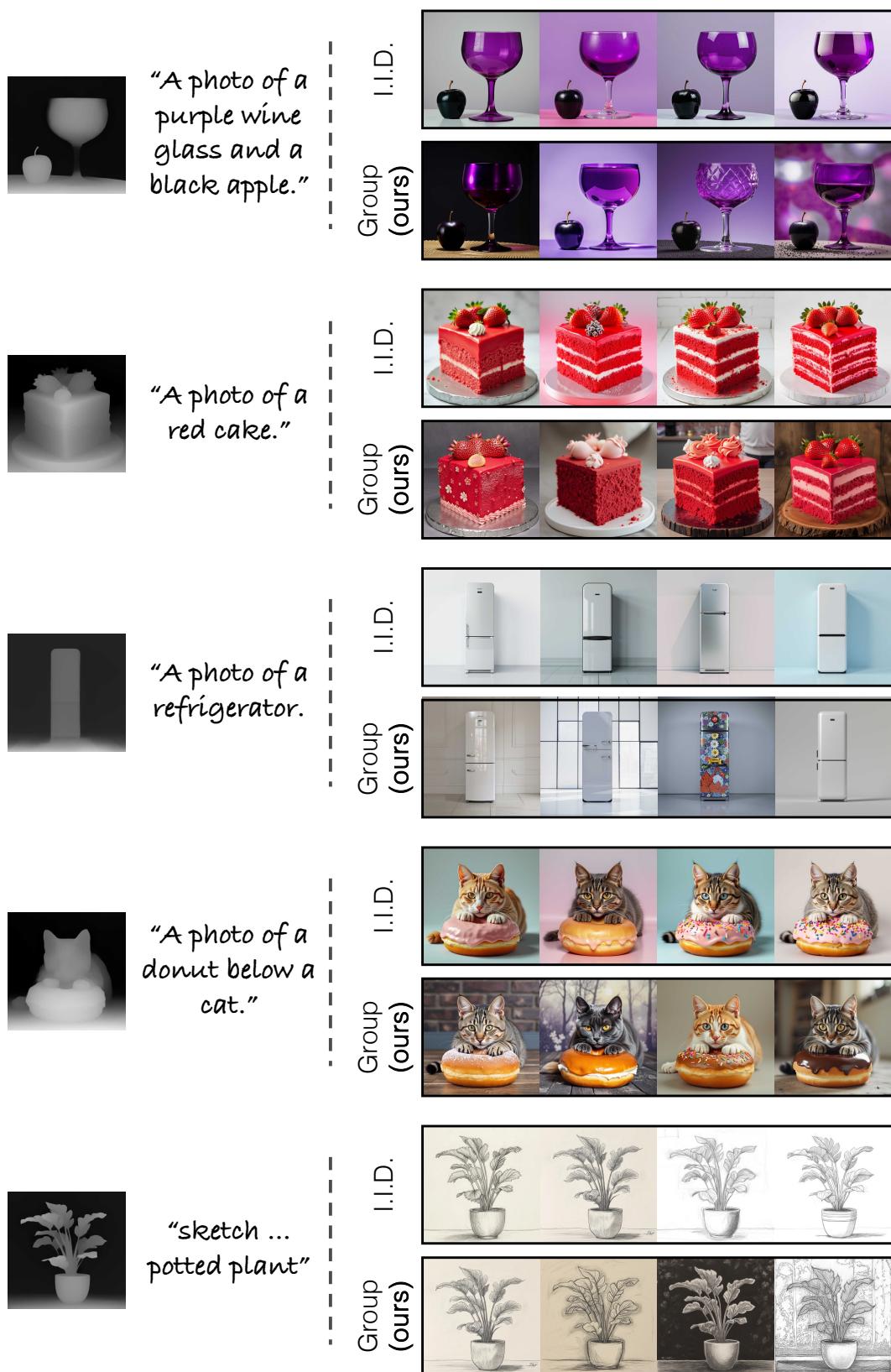
Figure 22: **Gallery of results.** Qualitative results that show the advantage of our proposed method over I.I.D. sampling for depth-to-image generation using FLUX.1 Depth as the base model. The input depth maps and captions are shown on the left and the generated outputs are shown on the right. Our method consistently generates outputs that have more diverse backgrounds, styles, and textures.
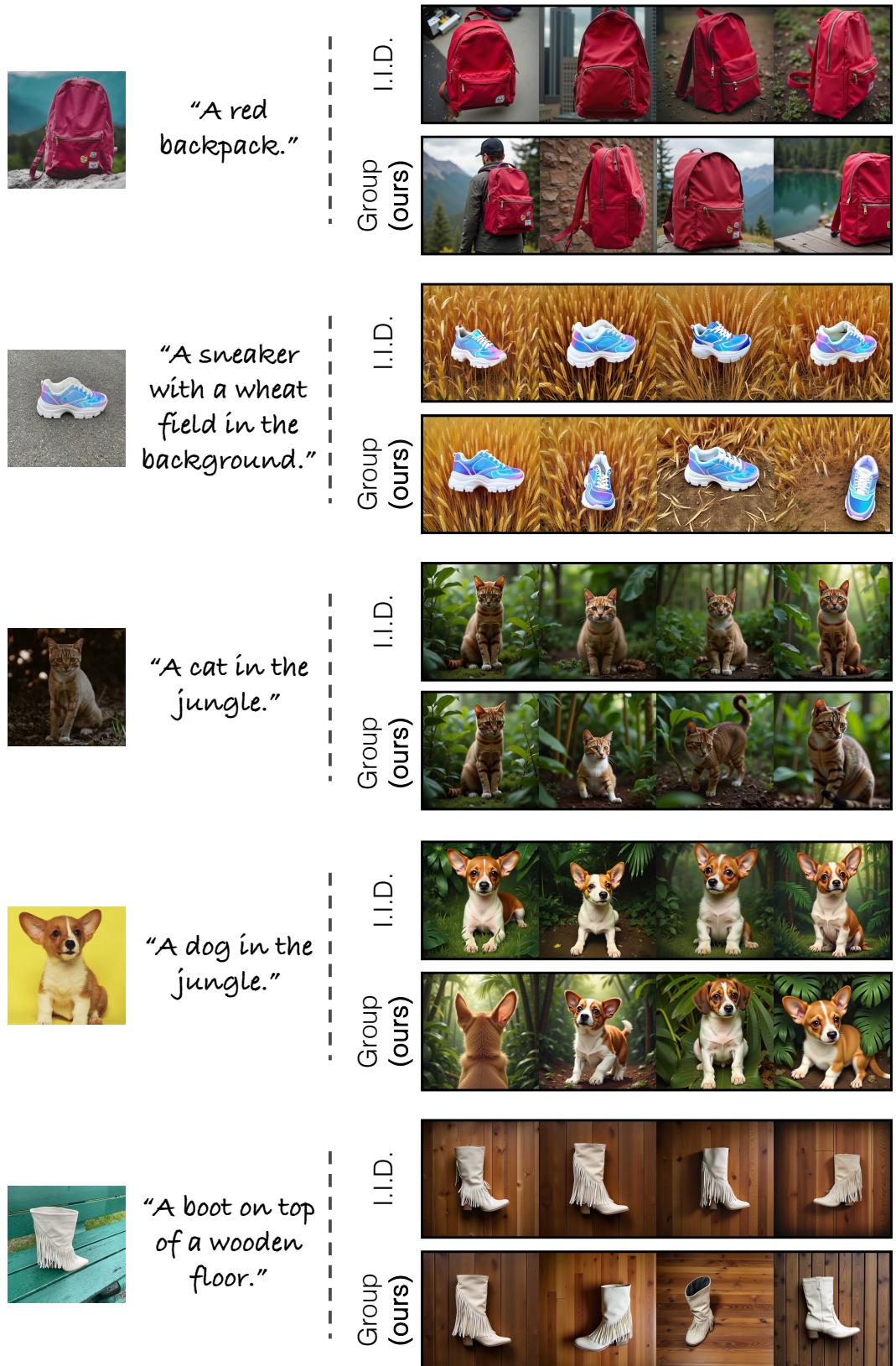
Figure 23: **Gallery of results.** Qualitative results that show the advantage of our proposed method over I.I.D. sampling for feedforward customized generation using SynCD [46]. The input image and captions are shown on the left and the generated outputs are shown on the right. Our method consistently generates outputs that have more diverse backgrounds, object poses, and styles.