

APPENDIX

A THE USE OF LARGE LANGUAGE MODELS

In preparing this research, we made limited use of large language models (LLMs) as assistive tools. Specifically, LLMs were employed for:

1. **Writing assistance:** LLMs were used to improve grammar, enhance clarity, and ensure an appropriate academic tone in certain parts of the manuscript.
2. **Editing and formatting:** To support professional presentation, LLMs were used to generate LaTeX table/figure templates and to refine the structure of paragraphs.

It should be emphasized that LLMs were not involved in the conception of core research ideas, the design or execution of experiments, the analysis of results, or the formulation of scientific conclusions. All technical contributions, algorithms, experiments, and analyses are original and were entirely carried out by the authors.

The authors take full responsibility for the content of this paper. LLMs were not considered contributors or co-authors, and their usage was strictly limited to the auxiliary functions described above.

B MORE MAIN RESULTS

Beyond the main experiments presented above, we further provide additional evaluations to demonstrate the robustness and generality of MVP.

Evaluation on other VAR-based T2I models. We further evaluate the transferability of MVP by applying it to VAR-based text-to-image models, namely Infinity and HART, both of which are pretrained on high-quality aesthetic datasets. To evaluate robustness under distribution shifts, we fine-tune these models on MS-COCO, a dataset featuring more realistic and diverse scenes that differ substantially from their original training distribution. As shown in Tab. 9, directly performing full fine-tuning on such distribution-shifted data often induces catastrophic forgetting, resulting in notable degradation in both FID and CLIP-Score. In contrast, MVP achieves consistent improvements while requiring only a fraction of the parameters to be updated, effectively preserving the pretrained generative priors and adapting to the new domain with minimal overhead. These results highlight the practicality of MVP as a lightweight and robust transfer approach for VAR-based T2I models.

Table 9: Transferability of MVP to other VAR-based text-to-image models (HART Tang et al. (2024) and Infinity Han et al. (2025)) under distribution shifts from aesthetic datasets to MS-COCO.

Method	Trainable Params (%)	FID \downarrow	CLIP-score \uparrow
HART-0.7B	-	36.2	0.22
HART-0.7B (Full FT)	100	56.4	0.17
HART-0.7B (MVP)	0.27	31.3	0.24
Infinity-2B	-	36.9	0.23
Infinity-2B (Full FT)	100	63.2	0.20
Infinity-2B (MVP)	0.13	29.9	0.24

C MORE ANALYSIS & ABLATIONS

In this section, we provide additional ablation studies to further analyze the effectiveness of each component in our method.

C.1 POTENTIAL CONCERN OF DISTRIBUTION DISTORTION.

One potential concern is that perturbation-based prompts in MVP may distort the pretrained VAR backbone’s distribution, thereby impairing generation quality and casting doubt on the viability of

810 perturbation-driven prompting. To better understand this effect, we performed ablation studies from
 811 two complementary perspectives.
 812

813 **Impact of the Gating Mechanism.** Gating is commonly expected to mitigate distribution distortion
 814 by constraining the magnitude of perturbation-based prompts. However, our experiments on
 815 improve class-conditional generation reveal the opposite: while gating reduces prompt strength, it
 816 also limits the method’s ability to provide effective guidance and can even degrade the fidelity and
 817 diversity of generated results (Tab. 10). This behavior likely stems from MVP’s multi-scale propa-
 818 gation of semantic and structural signals. When the perturbation strength is overly constrained, the
 819 prompts degenerate into weak noise that fails to deliver clear guidance and instead interferes with
 820 the pretrained features of the backbone. Consequently, although the gating mechanism is intended
 821 to reduce distributional shift, it paradoxically disrupts the learned representations and diminishes the
 822 overall effectiveness of our method.
 823

824 Table 10: Comparison between the default MVP design and its gated variant on enhancing the generative
 825 capability of VAR at different depths.
 826

Model	Depth	FID \downarrow	IS \uparrow	Precision \uparrow	Recall \downarrow
VAR	24	2.17	271.9	0.81	0.59
MVP	24	2.13	292.9	0.81	0.58
MVP (Gating)	24	2.17	273.9	0.80	0.59
VAR	30	2.14	275.4	0.80	0.60
MVP	30	2.03	289.4	0.81	0.59
MVP (Gating)	30	2.25	274.1	0.80	0.59

834 C.2 ABLATION ON INJECTION INTO HIDDEN STATES

835 A potential concern is the prompts injection applied only at the input stage might gradually diminish
 836 during propagation through the transformer layers, thereby weakening its effect. To mitigate
 837 this, we assess the impact of injecting prompts into hidden states during autoregressive modeling
 838 across varying model depths and injection frequencies. As shown in Tab. 11, this strategy leads to
 839 performance degradation rather than improvement. Specifically, increasing the injection frequency
 840 consistently causes both FID and IS to deteriorate across different backbone depths. This result
 841 further demonstrates that our method achieves optimal effectiveness with its minimalist design.
 842

843 Table 11: Impact of injecting projector-processed visual prompt tokens at varying frequencies into intermediate
 844 layers of the transformer.
 845

Depth	Injection	Freq.	FID \downarrow	IS \uparrow	Depth	Injection	Freq.	FID \downarrow	IS \uparrow
16	✗	-	3.46	247.4	24	✗	-	2.13	292.9
16	✓	1	3.51	242.9	24	✓	1	2.17	286.1
16	✓	2	3.69	217.4	24	✓	2	2.20	274.5
16	✓	3	3.77	209.4	24	✓	3	2.28	263.7
20	✗	-	2.63	276.5	30	✗	-	2.03	289.4
20	✓	1	2.70	264.7	30	✓	1	2.11	291.8
20	✓	2	2.81	254.1	30	✓	2	2.13	295.3
20	✓	3	2.87	231.3	30	✓	3	2.26	281.4

856 C.3 ABLATION ON PROMPT MODULATION

857 A natural concern is that adopting a single, fixed prompt shared across all classes may hinder the
 858 model’s ability to capture class-specific nuances, raising the question of whether adaptive prompts
 859 would better accommodate diverse generation tasks. To examine this, we designed two modulated
 860 variants: (i) a FiLM-based modulation that predicts feature-wise scale and shift parameters from the
 861 class condition, and (ii) a lightweight cross-attention module that injects class information into the
 862 original prompt sequence.
 863

The experimental results are summarized in Tab. 12. While these variants introduce additional flexibility, they also substantially increase the number of trainable parameters and the difficulty of optimization. Under the same training budget (e.g., 1 epoch), both variants underperform compared to our default design, showing weaker FID and IS scores. Extending training to more epochs can indeed improve the generation quality of these variants, but such longer schedules run counter to the efficiency principle of prompt tuning, whose key advantage lies in rapid adaptation with minimal computational cost.

Table 12: Impact of injecting projector-processed visual prompt tokens at varying frequencies into intermediate layers of the transformer.

Depth	FiLM	FID \downarrow	IS \uparrow	Epochs	Depth	CA.	FID \downarrow	IS \uparrow	Epochs
16	✗	3.46	247.4	1	16	✗	3.46	247.4	1
16	✓	3.74	228.2	1	16	✓	3.97	231.6	1
16	✓	3.61	239.7	3	16	✓	3.67	242.3	3
20	✗	2.63	276.5	1	20	✗	2.63	276.5	1
20	✓	2.74	259.1	1	20	✓	2.88	248.3	1
20	✓	2.68	267.7	3	20	✓	2.77	271.2	3

D MORE DETAIL

D.1 MORE IMPLEMENTATION DETAIL

In the VAR backbone configurations, the detailed training settings, including learning rate, batch size, number of epochs, and other hyperparameters, are provided in Tab. 13 for the class-to-image task and in Tab. 14 for the text-to-image task. Notably, for the experiments on improving class-to-image generation, we discard the first-scale prompts to prevent interference with the pretrained class embeddings and maintain the conditioning integrity.

Table 13: Implementation detail of MVP for class-to-image

backbone	VAR-d16	VAR-d20	VAR-d24	VAR-d30	VAR-d36
τ	20	20	20	28	28
first-scale prompt	✗	✗	✗	✗	✗
optimizer			AdamW		
AdamW (β_1, β_2)			(0.9, 0.95)		
learning rate	1e-3	1e-3	1e-3	5e-4	8e-5
weight decay	5e-2	5e-2	1e-2	1e-2	1e-2
batch size	132	124	82	58	12
epoch	1	1	1	1	1

Table 14: Implementation detail of MVP for text-to-image

backbone	VAR-d16	VAR-d20	VAR-d24	VAR-d30	HART-0.7B	Infinity-2B
τ	20	28	28	36	36	44
first-scale prompt	✓	✓	✓	✓	✗	✗
optimizer			AdamW			
AdamW (β_1, β_2)			(0.9, 0.99)			
learning rate	1e-4	1e-4	7e-5	8e-5	1e-4	1e-4
weight decay	5e-2	5e-2	1e-2	1e-2	5e-2	5e-2
batch size	148	124	82	58	12	8
epoch	1	1	1	1	1	1

918
919

D.2 TRAINING LOSS DETAILS

920
921
922
923
924
925
926
927
928

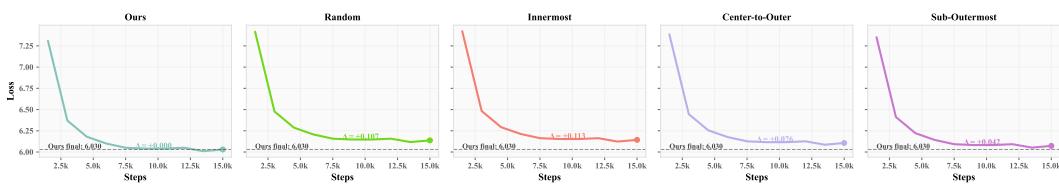
Training loss details of MVP and VAR To further complement the results presented in Table 8, we analyze the training loss dynamics of MVP and full-parameter VAR fine-tuning across three datasets (SUN397, Food101, and RESISC). As illustrated in Fig. 4, MVP consistently converges faster than full-parameter VAR, achieving both a steeper early-stage loss drop and a lower final loss. For clarity, we visualize both the raw loss and an exponentially smoothed version (smoothing factor value is 0.65). The results across all datasets show substantially more stable optimization behavior and a more favorable convergence trajectory with MVP. These observations further validate MVP’s strong transferability and robust adaptation ability compared with standard full-parameter VAR fine-tuning approaches.

929
930
931
932
933
934
935
936
937
938939
940

Figure 4: Training loss details for MVP vs. full-parameter VAR fine-tuning.

941
942
943
944
945
946
947
948

Training loss details for different prompt placement strategies To more precisely examine the differences in efficiency and performance across different prompt position designs, we visualize the training loss behavior of our design in MVP and four alternative position designs under an identical prompt token budget. As shown in Fig. 5, the resulting gaps Δ clearly show that our outermost-frame position design not only achieves the fastest convergence but also yields the lowest final loss across all variants. These results provide strong optimization-level evidence that our prompt placement design imposes minimal perturbation on pretrained representations while offering the most effective semantic conditioning.

949
950
951
952
953
954
955956
957
958
959

D.3 MORE THEORETICAL JUSTIFICATION DETAILS

960
961
962
963
964

To complement the simplified analysis in Section 3.1, we provide a more complete theoretical rationale explaining why the outermost-frame region constitutes the most informative and stable location for placing prompt tokens in multi-scale VAR generation.

965
966
967
968
969
970
971

Propagation Distance and Signal Attenuation Model We begin by formalizing the distance-based attenuation model of prompt influence. Consider the feature map as a 2D grid (or graph) of tokens $V = (x, y) \mid 1 \leq x, y \leq S_t$ at a given scale t . Let $d((x, y), (u, v))$ denote the distance between two token positions. For theoretical generality, d can be any appropriate metric on the grid (e.g. Manhattan distance, Euclidean, or Chebyshev distance), but we will often treat it as the shortest-path distance on the grid graph (equivalent to Manhattan distance if only orthogonal moves are allowed). We suppose that a prompt token at position p imparts a perturbation signal δ that propagates to other tokens with a strength that decays as a function of distance. Formally, let $f(d)$

be a monotonically decreasing attenuation function (with $f(0) = 1$ and $0 < f(d) < 1$ for $d > 0$). A simple choice in our initial analysis was an exponential decay $f(d) = \alpha^d$ for some $0 < \alpha < 1$, but more generally one could consider, for example, an exponential $f(d) = e^{-\lambda d}$, a power-law $f(d) = (1 + \beta d)^{-1}$, or other diffusion-based kernels. The key property is that influence diminishes with propagation distance.

If a single prompt token is introduced at some source location $u \in V$ with perturbation magnitude δ , the impact it has on a target token $v \in V$ can be modeled as $I_{u \rightarrow v} = \delta, f(d(u, v))$. When multiple prompt tokens are present, their effects can be assumed to superpose (e.g. additively, if we treat small perturbations linearly). Thus, for a set of prompt positions $B \subset V$, the total influence on v is:

$$I_B(v) = \sum_{u \in B} \delta_u f(d(u, v)),$$

where δ_u is the perturbation added at prompt u (for simplicity one may take all $\delta_u = \delta$). In general, $I_B(v)$ will be stronger for targets v that are closer to some prompt token and weaker for those far from all prompts. We can thus formalize two objectives for prompt placement:

(1) Minimal Central Impact: Avoid placing prompts too close to the critical center region of the image, so that any central token is at a large distance from all prompts. This keeps $f(d)$ small for center distances, mitigating disruptive changes to the core visual content. (2) Maximal Propagation Coverage: Ensure that every token in the feature map lies reasonably close to at least one prompt, so that prompt signals can efficiently reach all parts of the map. In other words, we want to minimize the worst-case (or average) distance from any token to the nearest prompt, thereby maximizing the minimum influence $I_B(v)$ across the map.

These goals are in tension: keeping prompts far from the center suggests placing them at the periphery, but doing so might increase distances to some other tokens. We next analyze this trade-off formally and show that an outermost-frame placement of prompts achieves an optimal balance under these criteria.

Boundary vs. Interior Prompt Placement: Distance Analysis Consider a decomposition of the $S_t \times S_t$ feature map into concentric square frames (or layers) around the center. Define Frame 0 as the set of all tokens on the outer boundary of the map (positions where $x = 1$, $x = S_t$, $y = 1$, or $y = S_t$). Frame 1 is the next inward layer (the “sub-outermost” border), and so on, until Frame N which contains the central token(s). By construction, $N = \lfloor \frac{S_t-1}{2} \rfloor$, i.e. there are N concentric frames inward from the boundary to the center. Each frame index n can be thought of as the Chebyshev distance (infinite norm distance) of those tokens from the image center. Now suppose we introduce prompt tokens on Frame c (meaning all prompt tokens lie in that concentric layer at distance c from the boundary). What is the distance from these prompts to the center, and to other tokens? Two key observations can be made:

- **Distance to Center:** All prompt tokens on Frame c are a distance of $(N - c)$ away from the center frame (in terms of frame index difference). In fact, the minimum distance from the center to any prompt in Frame c is $N - c$ (achieved along a straight line from the center to the prompt layer). Thus, a prompt at frame c induces a central impact of roughly $I_{\text{center}} \approx \delta, f(N - c)$. If we use the exponential model $f(d) = \alpha^d$, this recovers the earlier result that a prompt in frame n has impact δ, α^{N-n} on the center, which increases rapidly as n grows closer to N (i.e. as the prompt moves inward). Keeping prompts in the outermost frame ($c = 0$) maximizes the center distance $N - 0 = N$, yielding minimal impact on the central tokens $I_{\text{center}} \approx \delta, \alpha^N$. By contrast, any non-boundary placement ($c > 0$) would put prompt tokens closer to the center (distance $N - c$ with $N - c < N$), leading to significantly larger direct impacts on central features (e.g. a frame c prompt gives δ, α^{N-c} , and since $N - c < N$, we get $\alpha^{N-c} \gg \alpha^N$ for $\alpha \in (0, 1)$). In the extreme case of a prompt at the center itself ($c = N$), the distance to center is 0 and the central impact is maximized (no attenuation, $f(0) = 1$) – this would heavily “corrupt” the core visual features, which is exactly what we must avoid. Therefore, placing prompts on the outer boundary is theoretically optimal for protecting the image center: it maximizes the minimum distance from any prompt to the central region, minimizing unwanted perturbation of semantically critical center content.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034

- Distance to Other Tokens (Coverage): We must also ensure that prompt signals can reach and influence all other tokens in the feature map efficiently. For a given token $v = (x, y)$ in the map, let $\text{dist}_B(v) = \min u \in Bd(u, v)$ denote the distance from v to the nearest prompt in the set B . We seek to make $\text{dist}_B(v)$ as small as possible for all v . If B is the outermost frame ($B = B_t$ in scale t , using the notation B_t for the set of all boundary positions), then every non-prompt token lies somewhere inside the boundary. Intuitively, any interior location will be adjacent (in some direction) to the boundary after a certain number of steps outward. In fact, for an arbitrary token at (x, y) , the minimal Manhattan distance to the outer frame is given by:

$$\text{dist}_{B_t}(x, y) = \min\{x - 1, S_t - x, y - 1, S_t - y\}. \quad (10)$$

1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042

This formula says that the distance to the boundary is determined by whichever edge (top, bottom, left, or right) is nearest to (x, y) . For example, if a token is 5 cells from the top edge, 10 from the bottom, 3 from the left edge, and 8 from the right edge, then its nearest prompt on the outer frame is 3 away (via the left border). The worst-case distance from any token to the outer-frame prompts is achieved at the geometric center of the map: a central token is farthest from all edges. Indeed, from Eq. 10 one can show the maximum distance is

$$D_{\max}(B_t) = \max_{(x,y) \in V} \text{dist}_{B_t}(x, y) = \left\lfloor \frac{S_t - 1}{2} \right\rfloor = N,$$

1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052

attained at the center ($x = \lceil S_t/2 \rceil, y = \lceil S_t/2 \rceil$). Thus, with outermost-frame prompts, no token is more than N steps away from some prompt. In terms of the attenuation function, the minimum prompt influence on any token v is at least $\delta, f(D_{\max}) = \delta, f(N)$. Using the exponential model as an example, the weakest-controlled token (at the center) still receives a small but nonzero signal δ, α^N . All other tokens are closer than N steps to the boundary and hence receive stronger prompt influence than the center does. In other words, placing prompts on the boundary yields full coverage of the map with a radius N : prompt signals need at most N steps to reach any location.

1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

It is insightful to contrast this with alternative placements. If prompts were placed at the center instead (Frame N), the situation essentially inverts: the center prompt directly affects itself (distance 0), but tokens on the outer border are now farthest (roughly N steps away) and would experience the weakest influence. In fact, a single central prompt also has a worst-case distance of N – the corners/edges of the map are N units away from the center in Manhattan distance – so the maximum propagation distance is still N . Thus, in terms of worst-case distance alone, a central prompt or an all-boundary prompt cover the grid in a comparable radius. However, the crucial difference lies in which region of the image suffers the maximum distance (minimal influence). With boundary prompts, it is the central region that is farthest – but we want the center to be least affected (to preserve semantic content). With a central prompt, it is the border region that is farthest – meaning the periphery of the image gets the weakest control. For many vision tasks, the periphery often contains background or less critical details, whereas the center often contains the main subject; therefore, it is far preferable that the center be the least altered region. Boundary placement guarantees this, whereas central placement does the opposite. Furthermore, if we use not one but a distributed set of prompts, the boundary configuration can cover the image more uniformly. In fact, the outer-frame prompt set B_t consists of all $4S_t - 4$ edge positions (for an $S_t \times S_t$ map), surrounding the entire image. Most interior tokens will be very close to some edge (e.g. a token near the top of the image is only a few pixels from the top prompt band), and only the very middle of the image has the maximal distance N . By contrast, a small set of prompts in an interior region would leave large areas of the image (e.g. all four corners or sides) far from any prompt. Thus, outermost-frame prompts minimize the area of the feature map that lies at high propagation distance. This yields more uniform and efficient coverage: prompt signals originate from all sides and diffuse inward, reaching a given interior token from multiple directions. We can formalize this advantage by comparing distance distributions. For boundary prompts, the fraction of tokens within a distance d of some prompt grows rapidly with d – indeed, for distance $d < N$ the “uncovered” region is an inner square of side $(S_t - 2d)$, whose area shrinks quadratically as d increases. In contrast, for a centralized prompt, the covered region within distance d is just a central disk/square of area $\sim (2d + 1)^2$, and a large peripheral ring remains uncovered until d approaches N . As a result, for any reasonably small d , far more of the image is reached by boundary prompts than by an equal number of interior prompts. In a dynamic diffusion sense, if prompt signals propagate outward one step at a time, the boundary

1080 configuration will influence a large portion of the image in the first few propagation steps (with only
 1081 the middle lagging behind), whereas a central source influences only the vicinity at first and leaves
 1082 all outskirts unaffected until much later. This spatial coverage analysis underscores the efficiency of
 1083 having prompts on the outer frame.

1084 It is also useful to consider a more general placement: suppose prompts are placed on some inter-
 1085 mediate frame c (neither center nor outermost). The maximum distance to some token will then be
 $\max c, ; N - c$. There is a natural trade-off here: if c is small (near the boundary), the center is far;
 1086 if c is large (near the center), the boundary is far. The worst-case is minimized if $c \approx N/2$, which
 1087 would balance the farthest distance to center and to edge at roughly $N/2$. Indeed, purely from a
 1088 graph covering perspective, the “middle ring” of an image could theoretically minimize the absolute
 1089 worst-case distance to all points. However, such a placement means the center is only moderately
 1090 distant (on the order of $N/2$) from a prompt, implying a much stronger direct prompt effect on the
 1091 central content than the outermost frame does. In designing MVP, we prioritize protecting the center
 1092 from corruption over the marginal gain of reducing the overall radius by a small factor. Empirically,
 1093 even with outermost prompts the maximum distance N is only about half the image size, and we will
 1094 show that this still provides effective propagation across the image. In fact, because the outer frame
 1095 uses a greater number of prompt tokens distributed all around, it compensates for a larger radius
 1096 by multi-directional influence – the center receives some (weak) signal from all sides rather than a
 1097 strong push from one close prompt. This multi-source arrangement can guide the model subtly and
 1098 consistently from the periphery, rather than risking a heavy-handed alteration near the center.
 1099

1100
 1101

1102 **Graph-Theoretic and Diffusion Perspective** We can cast the above intuitions in terms of graph
 1103 diffusion or boundary-value problems, which lends another angle to why outer-frame prompting
 1104 is advantageous. Imagine the feature map as a weighted graph, where each token is a node and
 1105 edges connect neighboring tokens (e.g. adjacent in the grid). Prompt tokens can be seen as sources
 1106 introducing a certain state or “potential” into the network, which then spreads to other nodes through
 1107 the edges. Placing prompt sources on the boundary is analogous to setting boundary conditions in a
 1108 diffusion process. For instance, one could imagine an iterative update where at each network layer or
 1109 time step, tokens influence their neighbors (like heat diffusing). If we initialize all boundary nodes
 1110 with a certain perturbation value and interior nodes at zero, the diffusion or random-walk process
 1111 will cause the interior to gradually absorb influence from all sides. Classical results for diffusion
 1112 on a 2D domain tell us that with boundary sources, the interior will be harmonically influenced
 1113 from all boundaries, often resulting in a smooth gradient that is minimal at the center (the point
 1114 maximally distant from all sources). By contrast, if a source is at the center, a diffusion process
 1115 would send a wave of influence outward, with the strongest effect concentrated near the source and
 1116 decaying toward the boundaries. In a steady state (e.g. solving Laplace’s equation with a fixed value
 1117 at the prompt sources), having the boundary held at a certain perturbation value yields an interior
 1118 solution that is lowest at the center (the furthest point from the boundary conditions). In other words,
 1119 the center stays relatively untouched when control signals are applied at the periphery – which is
 1120 precisely what we want to ensure. This diffusion analogy reinforces that boundary prompts provide a
 1121 global, gentle influence that permeates the feature map from the outside in, whereas interior prompts
 1122 act more like local shocks that could disrupt central contents.

1122 From a graph-theoretic view, one can also consider the concept of a dominating set: the prompt set B
 1123 “dominates” the graph if every node in the graph is within a certain distance r of some prompt. The
 1124 smallest such r for a given B is the covering radius of that prompt set (our D_{\max} above). The entire
 1125 outer boundary B_t forms a dominating set with radius N . While a smaller dominating radius could
 1126 be achieved by a carefully chosen subset of interior nodes, those interior nodes would inherently
 1127 be closer to the center (reducing central distance and increasing corruption risk). The boundary set
 1128 has the special property that it maximizes the distance to the most sensitive node (the center) while
 1129 still maintaining a reasonable covering radius. In fact, under the constraint that no prompt is closer
 1130 than distance d_{\min} to the center, the outermost ring yields the minimal possible covering radius. For
 1131 example, if we require prompts to be at least N away from center, the only feasible locations are on
 1132 the outer frame; if we slightly relax to at least $N!-!1$ away, the second-outermost frame becomes
 1133 available, but using that instead of the outer frame would only shrink the radius by 1 at the heavy
 1134 cost of moving prompts closer to center. Thus, given the design constraint to maximize center safety,
 1135 outer-frame placement is the optimal choice for broad coverage.

1134 **Rigorous Justification of Outermost-Frame Optimality** We now synthesize the above points
 1135 into a more rigorous justification. Proposition: Placing prompt tokens on the outermost frame si-
 1136 multaneously minimizes the prompt’s direct impact on central tokens and provides near-optimal
 1137 coverage of the entire feature map. More concretely: (a) For any fixed attenuation function $f(d)$
 1138 that decreases with d , the maximum possible minimal distance between any prompt and the center
 1139 of the map is achieved by prompts on the boundary — this maximizes the center’s distance to the
 1140 nearest prompt, thereby minimizing the upper bound on prompt-induced change to central features.
 1141 (b) Subject to (a), the outermost-frame configuration ensures that the distance from any token to
 1142 some prompt is bounded by $N = \mathcal{O}(S_t)$, and no other configuration with equal or greater center
 1143 distance can achieve a smaller worst-case distance.

1144 Proof Sketch: (a) is straightforward — the farthest any point can be from the center in an $S_t \times S_t$
 1145 grid is along the boundary. Any prompt placed at an interior location has a distance to center
 1146 $< N$, whereas prompts on the edge have distance N to center; hence the boundary placement
 1147 uniquely achieves the maximum center distance N . Thus, if minimizing $I_{\text{center}} = \delta, f(d_{\text{center-prompt}})$
 1148 is paramount, one should choose $d_{\text{center-prompt}}$ as large as possible, i.e. place prompts at the periphery.
 1149 (b) Now, given that prompts are confined to those at least distance N (or $N - 1$, etc.) from center,
 1150 we consider the coverage of the map. The outer boundary B_t is one natural choice meeting this
 1151 constraint. Could any other allowed configuration cover the map more efficiently (i.e. with a smaller
 1152 maximum distance)? Suppose we remove some subset of boundary prompts or move some prompts
 1153 off the exact edge inward by one cell. Any interior point that was previously nearest to a removed
 1154 prompt will now be farther from the remaining prompt set, increasing the worst-case distance. In
 1155 fact, removing or insetting prompts can only increase the covering radius unless other prompts are
 1156 moved inward to compensate — but moving any prompt inward violates the center-distance constraint
 1157 or at least lowers the center distance achieved. The full ring of boundary tokens is a redundant but
 1158 robust cover — it may use slightly more tokens than minimally necessary for coverage, but this
 1159 redundancy guarantees no gaps in coverage and uniformly short distances except at the very center.
 1160 Formally, one can show that B_t (the set of all edge positions) minimizes the function $\Phi(B) =$
 1161 $\max_{v \in V} \min_{u \in B} d(u, v)$ among all sets B that satisfy $\min_{u \in B} d(u, \text{center}) = N$. In other words,
 1162 under the condition that no prompt is closer than N to the center, B_t yields the minimal possible
 1163 $\Phi(B)$ (which in fact equals N in this case). Any other set that also keeps prompts off the $n < N$
 1164 inner frames will either have the same worst-case distance N or worse. Thus, outermost placement
 1165 is Pareto-optimal in balancing the two objectives: you cannot increase the center safety without
 1166 also worsening coverage, and vice versa, and the chosen design maximizes center safety while only
 1167 modestly sacrificing distance-efficiency (staying within a factor of 2 of the absolute minimal radius
 1168 achievable by any prompt configuration, which is a small price for protecting central content).

1169 In summary, our expanded theoretical analysis confirms that introducing prompts in the outermost
 1170 square frame is an optimal or at least highly well-founded design choice. It minimizes the risk of
 1171 core feature corruption by keeping prompts as far as possible from the image’s crucial center, while
 1172 still efficiently diffusing control signals across the entire feature map from the boundaries inward.
 1173 This justifies the MVP strategy of injecting learnable prompt tokens in the outermost frame at each
 1174 scale, as it offers strong signal coverage with minimal adverse impact on the learned visual features
 1175 at the center of the generation.

E PROMPT-COMPLEXITY ANALYSIS

1176 At the t -th scale of the VAR’s patch size set \mathcal{P} , let the spatial feature map be of size $S_t \times S_t$. The
 1177 corresponding **prompt budget** is determined by the number of tokens \mathcal{N}_t introduced at this scale.
 1178 Below, we analyze the computational overhead associated with different prompt injection strategies.
 1179

- 1180 • **Full Feature-Map Prompt.** Each spatial location is placed as a prompt token:

$$1181 \mathcal{N}_t = \mathcal{N}_t^{\text{full}} = S_t^2.$$

1182 This leads to a quadratic increase in token count. For example, $S_t = 16 \Rightarrow \mathcal{N}_t = \mathcal{N}_t^{\text{full}} =$
 1183 256.

- 1184 • **Outermost Frame Prompt** ($\mathcal{N}^{B_t} \leq \tau$). Only the outermost frame positions are used:

$$1185 \mathcal{N}_t = \mathcal{N}^{B_t} = 4S_t - 4,$$

1188 which scales linearly with S_t (e.g., $S_t = 16 \Rightarrow \mathcal{N}_t = \mathcal{N}^{\mathcal{B}_t} = 60$). This formulation is only
 1189 valid for $S_t > 1$.
 1190

1191

- 1192 • **L-shaped Corner Prompt** ($\tau < \mathcal{N}^{\mathcal{B}_t}$). This strategy selects strips of stride width a
 1193 originating from the four corners of the feature map:
 1194

1195

$$\mathcal{N}_t = \mathcal{N}_C^{\mathcal{B}_t} = 8a + 4.$$

1196

1197 For example, with $a = 2$ and $S_t = 16$, we get:
 1198

1199

$$\mathcal{N}_t = \mathcal{N}_C^{\mathcal{B}_t} = 20.$$

1200

1201 This configuration is applied only when it does not exceed the corresponding square frame
 1202 prompt count, i.e., $\mathcal{N}_C^{\mathcal{B}_t} \leq \mathcal{N}^{\mathcal{B}_t}$.
 1203

1204 This comparison reveals a clear computational hierarchy:
 1205

1206

$$\mathcal{N}_t^{\text{full}} \sim \mathcal{O}(S_t^2) \gg \mathcal{N}^{\mathcal{B}_t} \sim \mathcal{O}(S_t) \gg \mathcal{N}_C^{\mathcal{B}_t} \sim \mathcal{O}(1),$$

1207

1208 where $\mathcal{N}_t^{\text{full}}$, $\mathcal{N}^{\mathcal{B}_t}$, and $\mathcal{N}_C^{\mathcal{B}_t}$ denote the token budget under full feature map, outermost square frame,
 1209 and 1-shaped corner prompting respectively. In practice, the 1-shaped corner strategy reduces the
 1210 token cost by a factor of $\sim 5\text{--}10\times$ compared to border prompting, and over $\sim 30\text{--}100\times$ compared
 1211 to dense prompting, while retaining sufficient spatial coverage. This establishes an efficient trade-off
 1212 between semantic guidance fidelity and computational complexity.
 1213

1214

F LIMITATION

1215

1216 While MVP demonstrates good performance in class-to-image and text-to-image generation, there
 1217 are still some limitations that warrant further investigation.
 1218

1219 **Limitation of CLIP’s text encoder.** We employ only the default CLIP text encoder throughout all
 1220 experiments. While more powerful or specialized language encoders may potentially yield better
 1221 generation quality, particularly in text-to-image tasks, we did not pursue this direction due to the
 1222 principle of minimizing model modifications in prompt tuning.
 1223

1224 **Limitation of Visual-Quality and Semantic-Quality Balance.** Frankly speaking, compared to
 1225 the significant improvements in metrics like IS after applying MVP, the improvement in FID, a low-
 1226 level fidelity metric, is limited. This can be attributed to the fact that our method introduces prompts
 1227 at the outermost frames of feature maps, thereby primarily emphasizing semantic information. As a
 1228 result, it inherently favors metrics like IS that capture cross-category semantics. In contrast, the im-
 1229 provement in FID typically requires optimization in the pixel space. While introducing perturbations
 1230 for each input token in the pixel space could effectively improve FID, it would also bring substantial
 1231 computational cost, which contradicts the original intention of prompt tuning. In conclusion, this is
 1232 an inherent limitation of prompt tuning.
 1233

1242 **G ALGORITHM OF MVP**
12431244

1245 **Algorithm 1:** Multi-Scale Visual Prompt

12461247 **Input:** $\mathcal{P} = \{S_1, \dots, S_T\}, \tau$ 1248 **Output:** $\mathcal{V}^{\mathcal{B}}, \mathcal{I}_{\text{id}}$ 1249 $\mathcal{V}^{\mathcal{B}} \leftarrow \text{Queue}(), \mathcal{I}_{\text{id}} \leftarrow \text{Queue}();$ // \triangleright initialize prompt and indices queues1250 **forall** $S_t \in \mathcal{P}$ **do**1251 $\mathcal{N}^t \leftarrow 4S_t - 4;$ // \triangleright compute number of outermost tokens1252 **if** $\mathcal{N}^t \leq \tau$ **then**1253 $\mathcal{N}^{\mathcal{B}_t} \leftarrow \mathcal{N}^t;$ // \triangleright set prompt count to outermost size1254 $\mathcal{I}_{\text{id}}^{\mathcal{B}_t} \leftarrow \text{Outermost}(S_t);$ // \triangleright select outermost token indices1255 **else**1256 $a \leftarrow \lfloor \frac{\tau-4}{8} \rfloor;$ // \triangleright compute stride width for corner sampling1257 $\mathcal{N}^{\mathcal{B}_t} \leftarrow 8a + 4;$ // \triangleright set prompt count under threshold1258 $\mathcal{I}_{\text{id}}^{\mathcal{B}_t} \leftarrow \text{LCorner}(S_t, a);$ // \triangleright select L-shaped corner indices1259 $\mathcal{V}^{\mathcal{B}_t} \leftarrow \text{Visual_Prompts}(\mathcal{N}^{\mathcal{B}_t});$ 1260 $\text{Queue_Push}(\mathcal{I}_{\text{id}}, \mathcal{I}_{\text{id}}^{\mathcal{B}_t});$ // \triangleright enqueue selected indices1261 $\text{Queue_Push}(\mathcal{V}^{\mathcal{B}}, \mathcal{V}^{\mathcal{B}_t});$ // \triangleright enqueue generated prompts1262 **return** $\mathcal{V}^{\mathcal{B}}, \mathcal{I}_{\text{id}}$

12631264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

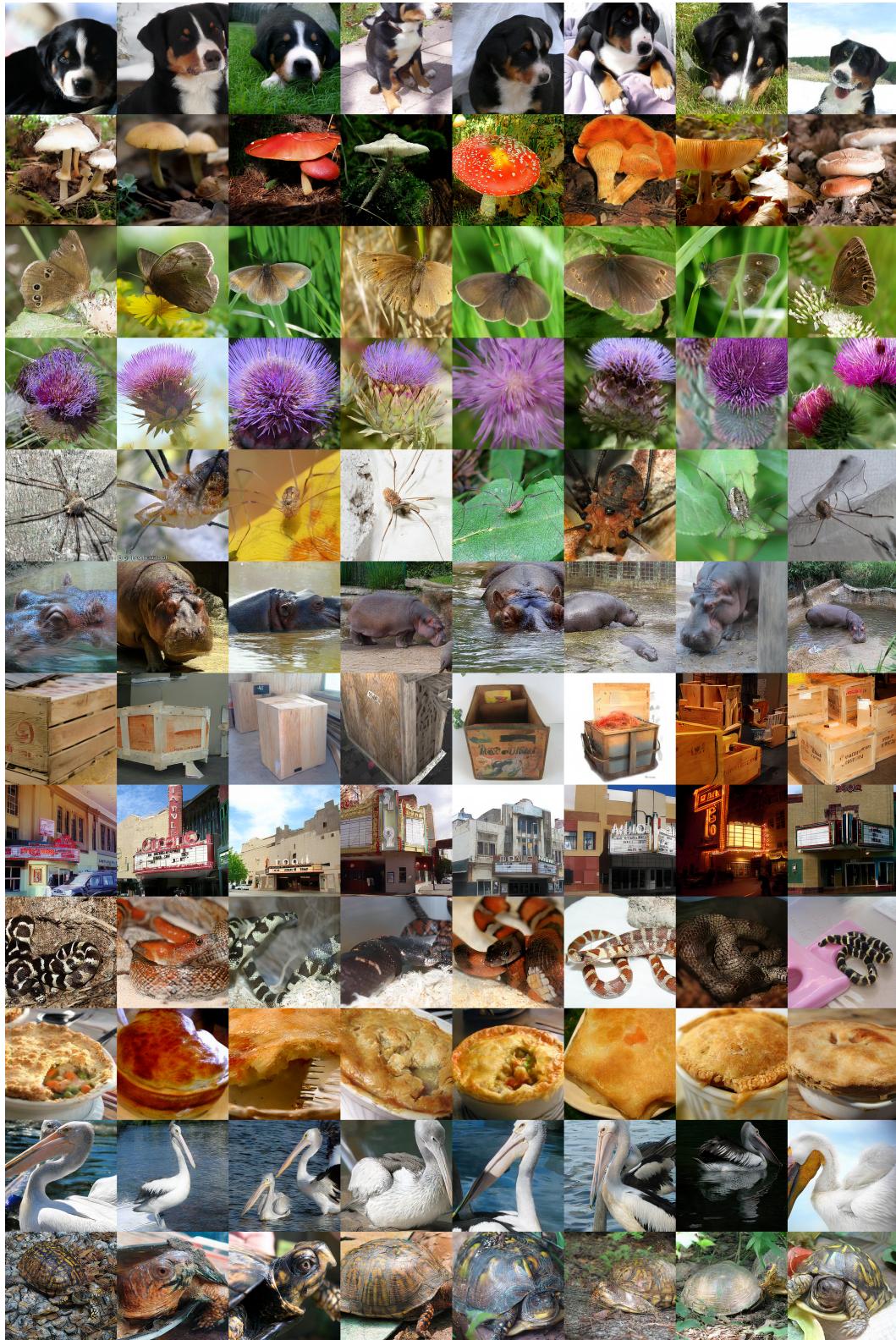
1296 **H VISUALIZATION**
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Figure 6: Visualization of class-to-image samples generated using MVP (512×512).

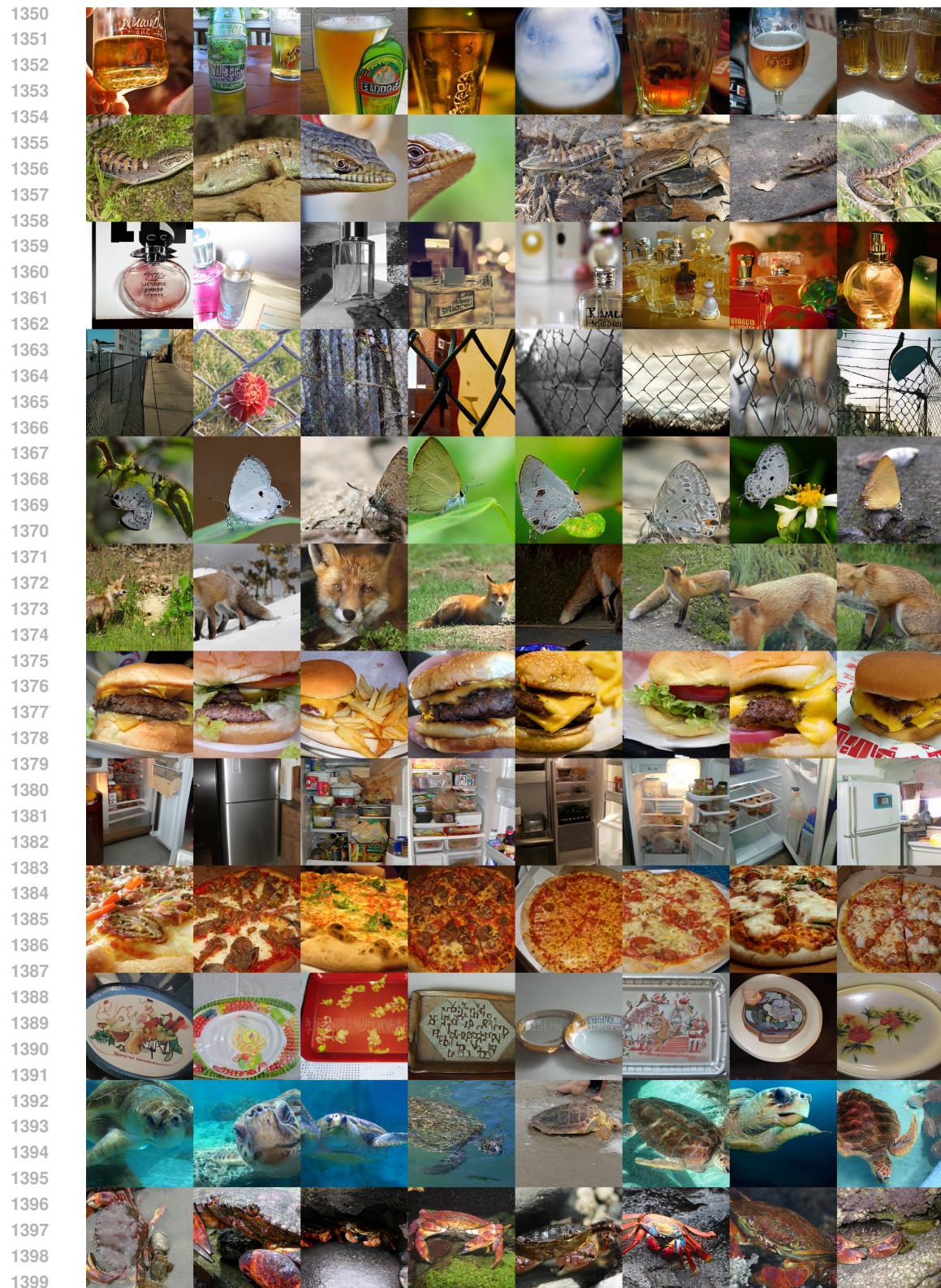
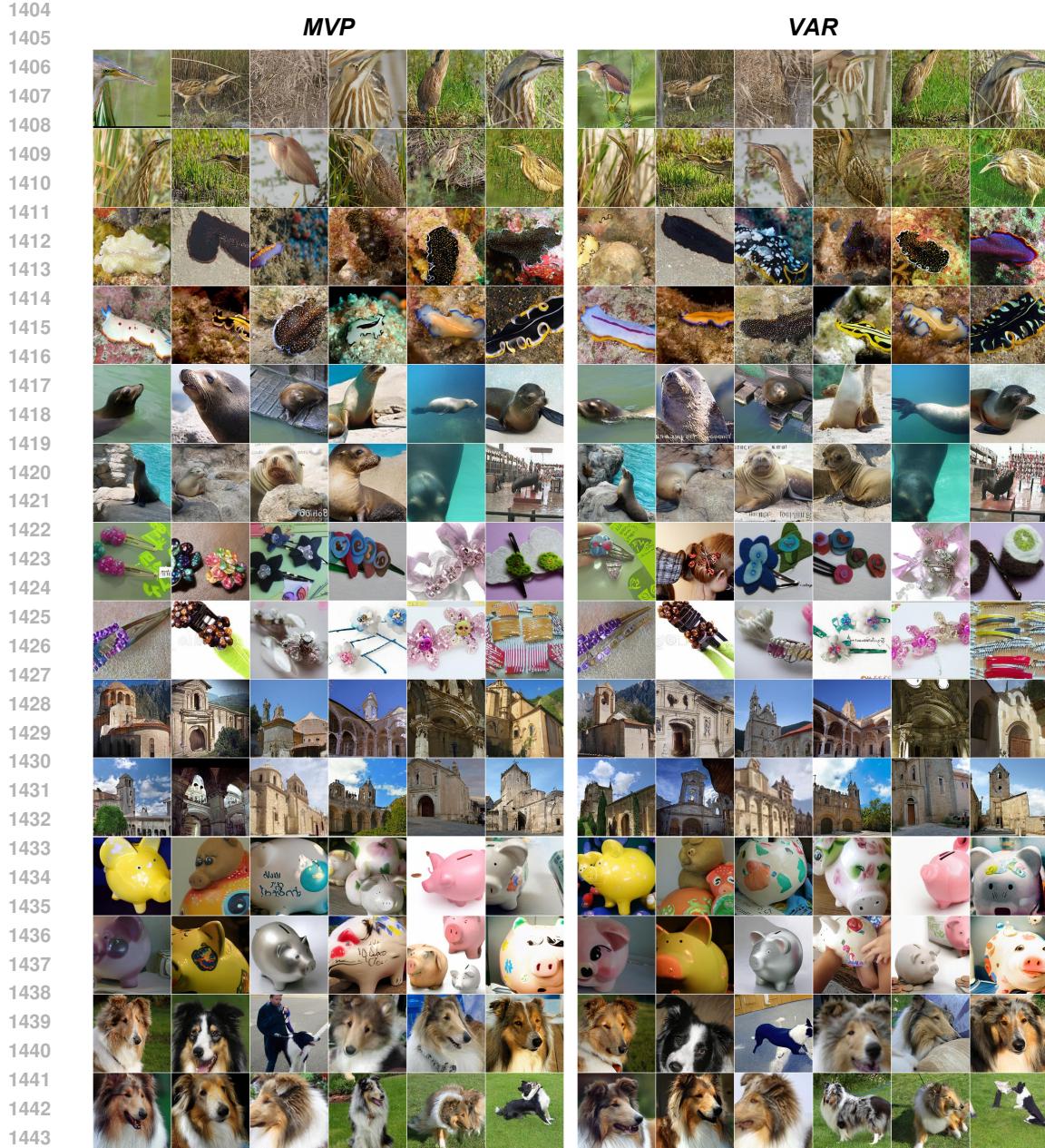


Figure 7: Visualization of class-to-image samples generated using MVP (256×256).

1400
1401
1402
1403

Figure 8: Comparison of class-conditional image samples generated by MVP and VAR at 256×256 resolution.

Class-to-Image. Visualization of samples conditioned on class labels. A subset of class labels is randomly sampled from the validation set, and multiple images are generated for each class using MVP. Each row corresponds to one class, and columns show diverse samples under the same condition. Fig. 6 and Fig. 7 present the results at 512×512 and 256×256 resolutions, respectively. As shown in Fig. 8, compared to VAR, our method produces more detailed and semantically accurate generations, demonstrating stronger alignment with the class-conditional guidance.

Text-to-Image. Fig. 9 and Fig. 10 illustrate examples of text-to-image generation results produced by MVP at 256×256 resolution. Although VAR is originally trained for class-conditional generation, these results demonstrate that it can be readily adapted to free-form text prompts through minimal finetuning with paired image-text data. The generated samples exhibit diverse visual con-

1458 tent and strong semantic consistency with the input text, indicating MVP ability to generalize beyond
 1459 class supervision with minimal adaptation.
 1460



1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 a mountain under a blue sky.



1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511
 a photo of little cute cat.



1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511
 a bird flying in the sky.



1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511
 a snowy road in a forest.



1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511
 a bookshelf full of books.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565
 1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1



Figure 10: Visualization of text-to-image samples generated using MVP (256 × 256).