

810 A APPENDIX

811 B DIFFUSION MODELS : DEFINITION & TRAINING PROCEDURE RECAP

812 B.1 FORWARD TIME DIFFUSION PROCESS

813 The background and expressions on forward diffusion process is taken from [Kingma et al. \(2021b\)](#)
814 and included here for completeness. Re-iterating Eq. [1](#), we have the forward diffusion as:

$$815 q(z_t | x) = \mathcal{N}(\alpha_t x, \sigma_t^2 I). \quad (15)$$

816 **Forward Conditional** $q(z_t|z_s)$: The distribution $q(z_t|z_s)$ for any $t > s$ are also Gaussian, and
817 from [Kingma et al. \(2021b\)](#), we can re-write as

$$818 \mathcal{N}(\alpha_{t|s} z_s, \sigma_{t|s}^2 I) \quad (16)$$

$$819 \text{where, } \alpha_{t|s} = \alpha_t / \alpha_s, \quad (17)$$

$$820 \text{and, } \sigma_{t|s}^2 = \sigma_t^2 - \alpha_{t|s}^2 \sigma_s^2 \quad (18)$$

821 **Reverse Conditional**, $q(z_s|z_t, x)$: The posterior $q(z_s|z_t, x)$ from [Kingma et al. \(2021b\)](#) can be
822 written as:

$$823 q(z_s|z_t, x) = \mathcal{N}(\mu_Q(z_t, x; s, t), \sigma_Q^2(s, t)I) \quad (19)$$

$$824 \text{where, } \sigma_Q^2(s, t) = \sigma_{t|s}^2 \sigma_s^2 / \sigma_t^2 \quad (20)$$

$$825 \text{and, } \mu_Q(z_t, x; s, t) = \frac{\alpha_{t|s} \sigma_s^2}{\sigma_t^2} z_t + \frac{\alpha_s \sigma_{t|s}^2}{\sigma_t^2} x \quad (21)$$

826 B.2 REVERSE DIFFUSION : DEFINING $p_\theta(z_s|z_t)$

827 Here, we describe in detail the conditional reverse model distributions $p_\theta(z_s|z_t)$ for the two cases
828 of variance-exploding and variance preserving diffusion process. Given these formulations, it
829 is straightforward to compute the KL distance between our posterior $q_\lambda(z_s|z_t, y)$ and the prior
830 $p_\theta(z_s|z_t)$ in our loss objective (Eq. [13](#)) since both are conditionally Gaussian distributions and computing
831 the KL between two Gaussians can be done in closed form.

832 **Variance Exploding Diffusion Process** In this case, $\alpha_t = 1$ and σ_t is usually in the range
833 $[0.002, 50]$ [Song et al. \(2021b\)](#). We follow the ancestral sampling rule from the same work to define
834 our prior conditional Gaussian distributions $p_\theta(z_s|z_t)$:

$$835 p_\theta(z_s|z_t) = \mathcal{N}(\mu_\theta(z_t; s, t), \sigma_Q^2(s, t)I) \quad (22)$$

$$836 \text{where, } \sigma_Q^2(s, t) = (\sigma_t^2 - \sigma_s^2) \frac{\sigma_s^2}{\sigma_t^2} \quad (23)$$

$$837 \text{and, } \mu_\theta(z_t; s, t) = \frac{\sigma_s^2}{\sigma_t^2} z_t + \frac{\sigma_t^2 - \sigma_s^2}{\sigma_t^2} \hat{x}_\theta(z_t, t) \quad (24)$$

838 where $\hat{x}_\theta(z_t, t) = z_t - \sqrt{(\sigma_t^2 - \sigma_s^2)} * \epsilon_\theta(z_t, t)$

839 **Variance Preserving Diffusion Process** In this case, $\alpha_t = \sqrt{1 - \sigma_t^2}$ and σ_t^2 is usually in the
840 range $[0.001, 1]$ [Ho et al. \(2020b\)](#). We follow the DDIM sampling rule [Song et al. \(2021a\)](#) to define
841 our prior conditional Gaussian distributions $p_\theta(z_s|z_t)$. This sampling rule is widely used to generate
842 unconditional samples in small number of steps, and naturally becomes a key design choice of our
843 prior. Here,

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

$$p_\theta(z_s|z_t) = \mathcal{N}(\mu_\theta(z_t; s, t), \sigma_Q^2(s, t)I) \quad (25)$$

$$\text{where, } \sigma_Q^2(s, t) = \eta \left(\frac{1 - \alpha_{t-1}}{1 - \alpha_t} \right) \left(1 - \frac{\alpha_t}{\alpha_{t-1}} \right) \quad (26)$$

$$\text{and, } \mu_\theta(z_t; s, t) = \sqrt{\alpha_{t-1}} \hat{x}_\theta(z_t, t) + \sqrt{1 - \alpha_t - \sigma_t^2 \epsilon_\theta(z_t, t)} \quad (27)$$

where $\hat{x}_\theta(z_t, t) = \frac{z_t - \sqrt{1 - \alpha_t} \epsilon_\theta(z_t, t)}{\sqrt{\alpha_t}}$. This schedule is adopted by the Latent Diffusion models.

B.3 DERIVATION OF OBJECTIVE FOR TRAINING DIFFUSION MODELS: $\mathcal{L}_{(0,T)}(z_0)$

The usual variational bound on the negative loglikelihood on data x : $\mathbb{E}[-\log p_\theta(x)] \leq \mathbb{E}_q[-\log \frac{p_\theta(z_{0:T})}{q(z_{1:T}|z_0, x)}] = \mathbb{E}_q[-\log p(z_T) - \sum_{t=1}^T \log \frac{p_\theta(z_{t-1}|z_t)}{q(z_t|z_{t-1})}]$. Let $0 < s < t < T$, we expand this derivation from [Ho et al. \(2020b\)](#) as follows:

$$\mathcal{L} = \mathbb{E}_q \left[\log \frac{q(z_{1:T}|z_0)}{p_\theta(z_0)} \right] \quad (28)$$

$$= \mathbb{E}_q \left[-\log p(z_T) + \sum_{t \geq 1} \log \frac{q(z_t|z_s)}{p_\theta(z_s|z_t)} \right] \quad (29)$$

$$= \mathbb{E}_q \left[-\log p(z_T) + \sum_{t > 1} \log \frac{q(z_t|z_s)}{p_\theta(z_s|z_t)} + \log \frac{q(z_1|z_0)}{p_\theta(z_0|z_1)} \right] \quad (30)$$

$$= \mathbb{E}_q \left[-\log p(z_T) + \sum_{t > 1} \log \frac{q(z_s|z_t, z_0)}{p_\theta(z_s|z_t)} \cdot \frac{q(z_t|z_0)}{q(z_s|z_0)} + \log \frac{q(z_1|z_0)}{p_\theta(z_0|z_1)} \right] \quad (31)$$

$$= \mathbb{E}_q \left[-\log \frac{p(z_T)}{q(z_T|z_0)} + \underbrace{\sum_{t > 1} \log \frac{q(z_s|z_t, z_0)}{p_\theta(z_s|z_t)}}_{\text{diffusion loss } \mathcal{L}_{(0,T)}(z_0)} - \log p_\theta(z_0|z_1) \right] \quad (32)$$

C VIPAINT : VI METHOD USING DIFFUSION MODELS AS PRIORS

C.1 DERIVATION OF VIPAINT’S TRAINING OBJECTIVE

As specified in the main paper, we define a variational distribution over the latent space variable z as $q_\lambda(z)$ and re-use the diffusion prior to generate $x \sim p_\theta(x | z)$. We derive the variational objective here:

$$\begin{aligned}
\mathcal{L}_N(\lambda; y) &= \mathbb{E}_{q_\lambda(z, x)}[\log p_\theta(y, x, z) - \log q_\lambda(z, x | y)] \\
&= \mathbb{E}_{q_\lambda(z, x)}[\log p_\theta(z) + \log p_\theta(x | z_{T_s}) + \log p_\theta(y | z_{T_s}) - \log q_\lambda(z) - \log q_\lambda(x | z_{T_s})] \\
&= \mathbb{E}_{q_\lambda(z)}[\log p_\theta(y | z_{T_s}) + \log p_\theta(z) - \log q_\lambda(z)] \\
&= \mathbb{E}_{q_\lambda(z)}[\log p_\theta(y | z_{T_s})] - \mathbb{E}_{q_\lambda(z)}[\log q_\lambda(z) - \log p_\theta(z)] \\
&= \mathbb{E}_{q_\lambda(z)}[\log p_\theta(y | z_{T_s})] - \mathbb{E}_{q_\lambda(z)}[\log q_\lambda(z) - \log p_\theta(z)] \\
&= \mathbb{E}_{q_\lambda(z)}[\log p_\theta(y | z_{T_s})] - \mathbb{E}_{q_\lambda(z)}\left[\sum_{i=0}^{K-1} \log q_\lambda(z_{s(i)} | z_{s(i+1)}) + \log q_\lambda(z_{T_e}) - \sum_{i=0}^{K-1} \log p_\theta(z_{s(i)} | z_{s(i+1)}) - \log p_\theta(z_{T_e})\right] \\
&= \mathbb{E}_{q_\lambda(z)}[\log p_\theta(y | z_{T_s})] - \sum_{i=0}^{K-1} D[q_\lambda(z_{s(i)} | z_{s(i+1)}) || p_\theta(z_{s(i)} | z_{s(i+1)})] - \underbrace{D(q(z_{T_e}) || p(z_{T_e}))}_{\mathcal{L}_{(T_e, T)}(z_{T_e})}
\end{aligned} \tag{33}$$

Negating the above objective, we get Eq. 13 in the main paper. Now, let’s derive the third term $L_{(T_e, T)}(z_{T_e})$ following section B.3

C.2 DERIVATION OF $L_{(T_e, T)}(z_{T_e})$

For any $T_e < s < t < T$, we have :

$$\mathbb{E}_{z_{T_e} \sim q_\lambda(z_{T_e})} \left[\log \frac{q(z_{T_e+1:T} | z_{T_e})}{p_\theta(z_{T_e:T})} \right] \tag{34}$$

$$= \mathbb{E}_{z_{T_e} \sim q_\lambda(z_{T_e})} \left[-\log p(z_T) + \sum_{t \geq T_e} \log \frac{q(z_t | z_s)}{p_\theta(z_s | z_t)} \right] \tag{35}$$

$$= \mathbb{E}_{z_{T_e} \sim q_\lambda(z_{T_e})} \left[-\log p(z_T) + \sum_{t > T_e} \log \frac{q(z_t | z_s)}{p_\theta(z_s | z_t)} + \log \frac{q(z_{T_e+1} | z_{T_e})}{p_\theta(z_{T_e} | z_{T_e+1})} \right] \tag{36}$$

$$= \mathbb{E}_{z_{T_e} \sim q_\lambda(z_{T_e})} \left[-\log p(z_T) + \sum_{t > T_e} \log \frac{q(z_s | z_t, z_{T_e})}{p_\theta(z_s | z_t)} \cdot \frac{q(z_t | z_{T_e})}{q(z_s | z_{T_e})} + \log \frac{q(z_{T_e+1} | z_{T_e})}{p_\theta(z_{T_e} | z_{T_e+1})} \right] \tag{37}$$

$$= \mathbb{E}_{z_{T_e} \sim q_\lambda(z_{T_e})} \left[-\log \frac{p(z_T)}{q(z_T | z_{T_e})} + \underbrace{\sum_{t > T_e} \log \frac{q(z_s | z_t, z_{T_e})}{p_\theta(z_s | z_t)}}_{\text{diffusion loss } \mathcal{L}_{(T_e, T)}(z_{T_e})} - \log p_\theta(z_{T_e} | z_{T_e+1}) \right] \tag{38}$$

The first and third term can be stochastically and differentially estimated using standard techniques. Following Kingma et al. (2021b), we derive an estimator for the diffusion loss $\mathcal{L}_{(T_e, T)}(z_{T_e})$. In the case of finite timesteps $t > T_e$, this loss is:

$$\mathcal{L}_{(T_e, T)}(z_{T_e}) = \sum_{t > T_e} \mathbb{E}_{q(z_t | z_{T_e})} D[q(z_s | z_t, z_{T_e}) || p_\theta(z_s | z_t)] \tag{39}$$

Estimator of $\mathcal{L}_{(T_e, T)}(z_{T_e})$ Reparametering $z_t \sim q(z_t|z_{T_e})$ as $z_t = \alpha_{t|T_e} z_{T_e} + \sigma_{t|T_e} \epsilon$, where $\epsilon \sim \mathcal{N}(0, 1)$, and to avoid having to compute all $T - T_e$ terms when calculating the diffusion loss, we construct an unbiased estimator of $\mathcal{L}_{(T_e, T)}(z_{T_e})$ using

$$\mathcal{L}_{(T_e, T)}(z_{T_e}) = \frac{T - T_e}{2} \mathbb{E}_{\epsilon, t \sim \mathcal{U}(T_e, T)} [D(q(z_s|z_t, z_{T_e}) || p_\theta(z_s|z_t))] \quad (40)$$

where $\mathcal{U}(T_e, T)$ is a uniform distribution to sample $T_e < t \leq T$ from a non-uniform discretization of timesteps using [Karras et al. \(2022\)](#).

Now, we elaborate on the expression $q(z_s|z_t, z_{T_e})$ and $p(z_s|z_t)$ for any $T_e < s < t < T$.

C.2.1 $q(z_s|z_t, z_{T_e})$

Our posterior at T_e is $q(z_{T_e}) = \mathcal{N}(\mu_{T_e}, \tau_{T_e}^2)$. For any $T_e < s < t < T$, we have $q(z_s|z_{T_e}) = \mathcal{N}(\alpha_{s|T_e} z_{T_e}, \tau_{s|T_e}^2)$ and $q(z_t|z_s) = \mathcal{N}(\alpha_{t|s} z_s, \sigma_{t|s}^2)$, yielding the posterior :

$$q(z_s|z_t, z_{T_e}) = \mathcal{N}(\mu_Q(z_t, z_{T_e}; s, t, T_e), \sigma_Q^2(s, t, T_e)I) \quad (41)$$

$$\text{where, } \sigma_Q^2(s, t, T_e) = \sigma_{t|s}^2 \frac{\tau_{s|T_e}^2}{\sigma_{t|s}^2 + \alpha_{s|T_e}^2 \tau_{s|T_e}^2} \quad (42)$$

$$\text{and, } \mu_Q(z_t, z_{T_e}; s, t, T_e) = \sigma_Q^2 \left(\frac{\alpha_{s|T_e}}{\tau_{s|T_e}^2} z_{T_e} + \frac{\alpha_{t|s}}{\sigma_{t|s}^2} z_t \right) \quad (43)$$

$$= \frac{\alpha_{s|T_e} \sigma_{t|s}^2}{(\sigma_{t|s}^2 + \alpha_{s|T_e}^2 \tau_{s|T_e}^2)} z_{T_e} + \frac{\alpha_{t|s} \tau_{s|T_e}^2}{(\sigma_{t|s}^2 + \alpha_{s|T_e}^2 \tau_{s|T_e}^2)} z_t \quad (44)$$

C.2.2 $p(z_s|z_t)$

The conditional model distributions can be chosen as:

$$p_\theta(z_s|z_t) = q(z_s|z_t, z_{T_e} = \hat{z}_{\theta, T_e}(z_t, t)) = \mathcal{N}(z_s; \mu_\theta(z_t, z_{T_e}; s, t, T_e), \sigma_Q^2(s, t, T_e)) \quad (45)$$

$$\text{where, } \mu_\theta(z_t, z_{T_e}; s, t, T_e) = \frac{\alpha_{s|T_e} \sigma_{t|s}^2}{(\sigma_{t|s}^2 + \alpha_{s|T_e}^2 \tau_{s|T_e}^2)} \hat{z}_{\theta, T_e}(z_t, t) + \frac{\alpha_{t|s} \tau_{s|T_e}^2}{(\sigma_{t|s}^2 + \alpha_{s|T_e}^2 \tau_{s|T_e}^2)} z_t \quad (46)$$

$$\text{and, } \sigma_Q^2(s, t, T_e) = \sigma_{t|s}^2 \frac{\sigma_{s|T_e}^2}{\sigma_{t|s}^2 + \alpha_{s|T_e}^2 \sigma_{s|T_e}^2} \quad (47)$$

$$\text{where } \hat{z}_{\theta, T_e}(z_t, t) = \frac{z_t - \sigma_{t|T_e} * \epsilon_\theta(z_t, t)}{\alpha_{t|T_e}}$$

C.3 TIME ANALYSIS

For a K step hierarchical posterior, each optimization step to fit the posterior requires $K + 1$ number of function evaluations (#NFEs) of the denoising network, where it uses K #NFEs during sampling from the posterior and 1 to compute the diffusion prior loss. Gradients are back-propagated through the denoising network during optimization, making each optimization step more informative than other works iteratively refining samples. We see the progress of fitting VIPaint’s posterior in Fig. 4 and often times, we observe that convergence occurs as low as in 50 optimization steps. Thus, VIPaint with 2 steps in its hierarchy can effectively infer the semantics in the image using only a cumulative of $50 * (2 + 1) = 150$ NFEs of the denoising network. The optimization time is $O(K)$, where $K \ll T$. but also linearly increases the inference time.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
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

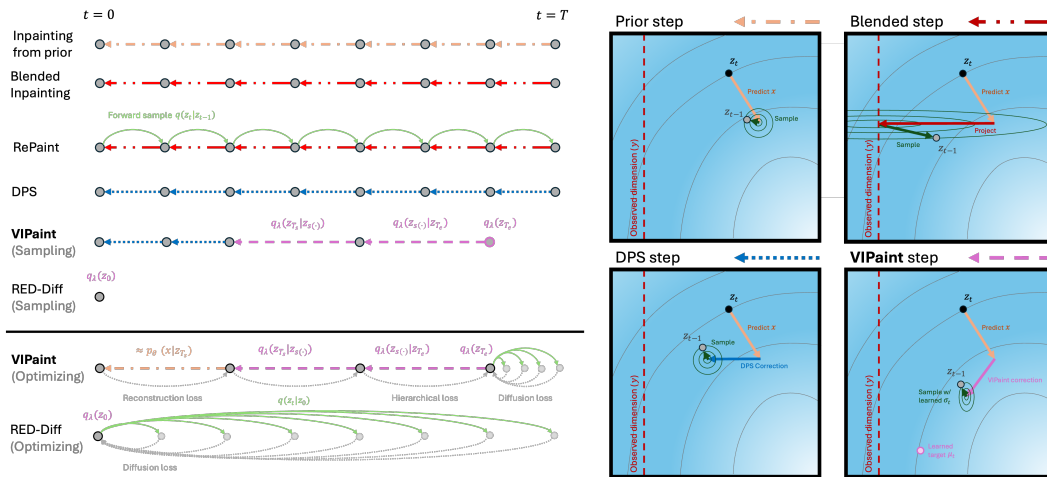


Figure 8: Expanded comparison of methods for diffusion model-based inpainting. **Left:** Timeline illustration of sampling steps with time flowing rightward from $t = 0$ (clean images) to $t = T$ (pure noise). **Orange arrows** indicate a single step of ancestral sampling under the generative prior $p_\theta(z_{t-1}|z_t)$. **Red arrows** indicate a single step of the *Blended* approximation of $p_\theta(z_{t-1}|z_t, y)$, while **blue arrows** indicate a single step of the *DPS* approximation. **Green arrows** indicate steps forward in time according to the diffusion process $q(z_t|z_{<t})$. Methods such as *RePaint* and *CoPaint* alternate between forward and reverse steps. **Purple arrows** indicate sampling from a step in the hierarchical **VIPaint** posterior $q_\lambda(z_{s(i-1)}|z_{s(i)})$. Both **VIPaint** and **RED-Diff** (without annealing) involve an initial optimization stage to fit variational parameters per-image. **Gray arrows** indicate the flow of gradient information during this optimization stage. **Gray points** are steps only used during optimization. **Right:** Illustration of each reverse-time sampling step in 2 dimensions. The horizontal dimension is assumed to be observed at the value marked by the **red line**. Each approach begins by computing $p_\theta(z_{t-1}|z_t)$ via a prediction of x using the pre-trained denoising network $\hat{x}_\theta(z_t, t)$. *Blended* replaces observed dimensions with $q(z_{t-1}|y)$. *DPS* updates $p_\theta(z_{t-1}|z_t)$ according to a single-step approximation to the likelihood $p_\theta(y|z_{t-1})$. Finally, **VIPaint**, uses a learned variational distribution $q_\lambda(z_{t-1}|z_t)$, which can be seen as interpolating between the prediction of x and a variational parameter μ_t , coupled with a learned variance.

D EXPANDED FIGURE 8

E EXPERIMENTAL DETAILS

E.1 VIPAINT

Choosing (T_s, T_e) for VIPaint Extensive prior work Song et al. (2021b); Nichol & Dhariwal (2021); Dhariwal & Nichol (2021); Karras et al. (2022) explores different noise schedules for training diffusion models, and how it affects the generated image quality. Since we use these diffusion models incorporating different noise schedules, our latent hierarchical posterior needs to account for this shift, and we show that it is flexible to do so. To concentrate posterior inference on the noise levels which are most crucial to perceptual image quality, we define our posterior at intermediate time steps that induce a signal-to-noise-ratio $(\alpha_t^2/\sigma_t^2) \in [0.2, 0.5]$, Kingma et al. (2021b) approximately across our experiments. This corresponds to choosing $(T_e = 5, T_s = 2)$ for the pixel-based EDM prior with a variance-exploding noise schedule and $(T_e = 550, T_s = 400)$ (DDIM sampling coefficient, $\eta = 0.2$) for the LDM prior using the VP noise schedule for both LSUN and ImageNet256 datasets. VIPaint is not sensitive to any subset of K timesteps in between this signal-to-noise range. For instance, for VIPaint-4, we take $[T_e = 5, 4, 3.5, 2.5, T_s = 2]$ for the EDM noise schedule and $[T_e = 550, 500, 450, T_s = 400]$ for the LDM prior.

Choosing K We discussed in section ?? that for a K step hierarchical posterior, the optimization run-time of VIPaint is $O(K)$. From the experiments conducted, we see that K can be easily selected to trade-off time and sample quality.

Initialization We follow the forward and reverse diffusion process defined by each VE and VP noise schedules to initialize VIPaint’s variational parameters. For LDM prior, we use the lower dimensional encoding of y . We provide a comprehensive summary in Table 2.

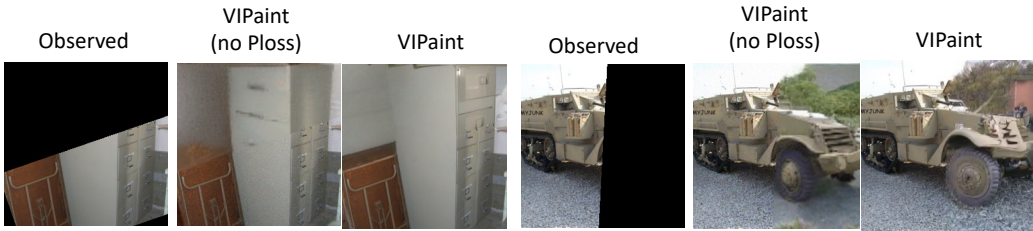
Table 2: Initialization of Variational Parameters for VE and VP Schedules

VI Parameters	VP Schedule	VE Schedule
$\mu_{T_e} = \alpha_{T_e} y + a_1 \sigma_{T_e} \epsilon$ (Scale factor to retain information from y .)	$a_1 = 0.8$	$a_1 = 0.01$
$\mu_{s(i)} = \alpha_{s(i)} y + a_2 \sigma_{s(i)} \epsilon$ (Noise adding process is still quite high for VE schedules.)	$a_2 = 1$	$a_2 = 0.01$
$\tau_{T_e} = \sigma_{T_e}$ (From the forward diffusion process.)	–	–
$\tau_{s(i)s(i+1)}$ (From the reverse diffusion process.)	Eq. 26 with scaling factor a_3/η $a_3 = 0.7$	Eq. 23
$\gamma_{s(i)} \forall i \in [1, K]$ (Weights samples from prior to construct plausible and close to real looking samples.)	0.98 (ImageNet256), 0.88 (LSUN)	0.5

Optimization We fit three sets of variational parameters at every i -th critical time in our hierarchy: means, $\mu_{s(i)}$, variances $\tau_{s(i)}^2$ and weights $\gamma_{s(i)}$. Instead of optimizing τ and γ directly, we optimize the real valued $\tilde{\tau} = \log \tau^2$, and $\tilde{\gamma} = \log(\frac{\gamma}{1-\gamma})$. We optimize this set of variational parameters $\lambda = \{\mu, \tilde{\gamma}, \tilde{\tau}\}$ using Adam with an initial learning rate of $\{0.1, 0.1, 0.01\}$ respectively and decreasing the learning rate by a factor of 0.99 every 10 iterations. We find this setting to be robust across all prior diffusion models and datasets in our work.

During pre-training, most diffusion models parameterize the mean prediction at every diffusion time step t and fix variances, however some previous work Nichol & Dhariwal (2021); Dhariwal & Nichol (2021) has found that (with appropriate training “tricks”) learning variances improves performance. Some previous works like ReSample tunes this as a hyperparameter. We instead learn this in our work, and we adjust learning rates to avoid local optima in this process. We optimize the parameters in VIPaint with $K = 2$ for 50 iterations; VIPaint with $K = 4$ is optimized for 100 steps in the case of LSUN Churches, 150 steps for the ImageNet64 dataset and 250 steps for the ImageNet256 dataset.

1134
1135
1136
1137
1138
1139
1140
1141
1142



1143
1144
1145
1146

Figure 9: An ablation showing the effect of addition the perceptual loss (PLOSS) in the reconstruction term for the task of image inpainting using latent diffusion priors. We see that even though VIPaint can inpaint the image semantically without the Perceptual loss, this loss becomes important to produce sharper reconstructions.

1147
1148
1149
1150
1151
1152
1153
1154
1155



1156
1157
1158
1159

Figure 10: (Left) We show the effect of the hyperparameter β with VIPaint with respect to optimization iterations. (Right) we show the respective loss curves. With $\beta = 10$, VIPaint captures more variations under the diffusion prior instead of "setting" to one kind of completion with $\beta = 1$.

1160
1161
1162
1163
1164
1165
1166

Sampling Post training, we take 400 iterative refinement steps from $T_s = 400$ in the LDM variance preserving schedule to sample inpaintings using a scale factor of 2, similar to the DPS algorithm using perceptual loss. On the other hand, for the EDM prior, we take 700 refinement steps to produce inpaintings after $T_s = 2$, with scale 5 (similar to DPS tuned for EDM in our work). This scale hyperparameter is tuned over the values [0.1, 0.5, 1, 2, 5, 10] on a validation set of 20 images. During the sampling phase, we use the classifier-free guidance rule with scale = 3 for the ImageNet256 latent diffusion prior.

1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177

Reconstruction loss We assume $p(y|z_{T_s})$ as a Laplacian distribution, where the mean is given by y and a scale parameter, which is computed over 100 images per dataset as a standard deviation over all pixel dimensions. For the 256 pixel datasets, this is 0.56, and for ImageNet64 it is 0.05. In addition to this, we add the perceptual loss for LDM priors, computing them via feeding the pre-trained Inception network with masked images and masked reconstructions. See Fig. 9 for the benefits of using the perceptual loss with VIPaint. Additionally, we use $\beta = 1$ for VIPaint with $K = 1$, which is optimized for 50 iterations for faster convergence. For VIPaint with $K = 4$, we use $\beta = 50$ for pixel-based EDM prior and $\beta = 10$ for LDM prior. We show the effect of the different β values in Fig. 10. Generally speaking, higher values of β explores the diffusion latent space more and lower values weighs the likelihood term relatively more and converges faster to a solution.

1178
1179
1180
1181

Discretization of timesteps for prior diffusion loss $L_{(T_e, T)}(z_{T_e})$ Lastly, we directly adapt the discretization technique from EDM Karras et al. (2022) to compute the diffusion loss. We use $\rho = 7$ across all models and datasets as used by Karras et al. (2022).

1182
1183
1184
1185
1186
1187

E.2 BASELINE DETAILS

Across all the baselines applicable to the latent diffusion models for the ImageNet256 dataset, we use the classifier-free guidance with a scale 3 Rombach et al. (2022b).

Blended We run blended for 1000 discretization steps using the EDM and LDM prior.

RePaint RePaint uses a descritization of 256 steps along with the standard jump length = 10, and number of times to perform this jump operation also set to 10, following standard practice [Lugmayr et al. \(2022\)](#).

DPS Similar to blended, we take 1000 denoising steps for DPS and set scale = 5 for the edm-based diffusion model, while take 500 steps and keep scale as 0.5 for the Latent Diffusion prior (similar to the original work in [Chung et al. \(2023\)](#)). When using the perceptual loss for the latent diffusion prior, we increase the scale to 2.

PSLD This is an inference technique only for the Latent Diffusion prior. Similar to DPS, we take 500 steps and keep scaling hyperparameters set to 0.2 as opposed to choosing 0.1 in the original paper [Rout et al. \(2023\)](#). We observe artifacts in the inpainted image if we increase the scale further as also observed by [Chung et al. \(2024\)](#).

CoPaint We directly adapt the author-provided implementation of CoPaint and CoPaint-TT [Zhang et al. \(2023\)](#) to use the EDM prior. Apart from the diffusion schedule and network architecture (taken from EDM) all other hyperparameters are preserved from the base CoPaint implementation.

RED-Diff As with CoPaint, We directly adapt the author-provided implementation of RED-Diff [Mardani et al. \(2024\)](#) and Red-Diff (Var) to use the EDM prior. In this case we increased the prior regularization weight from 0.25 to 50, which we found gave improved performance and more closely matches our VIPaint settings.

ReSample As with other baselines, we directly adapt the author-provided implementation of ReSample [Song et al. \(2024\)](#) for the LDM prior. Because the original code takes larger optimization steps, resulting in high sampling time, we decrease the number of optimization steps to 50, such that the wall-clock run-time of this method matches the other baselines.

F INFERENCE TIME.

We report the time taken for each inference method to produce 10 inpaintings for 1 test image. VIPaint with $K = 2$ is comparable to the baseline methods in terms of wall clock time and the number of functional evaluations (#NFEs) of the denoising network. Red-Diff, Blended and RePaint baseline methods do not take gradient of the noise prediction network, whereas all other methods require gradients. In terms of time and NFEs, we can see that VIPaint (fast) takes comparable time and NFEs as other baselines, but performs far better (Table ??).

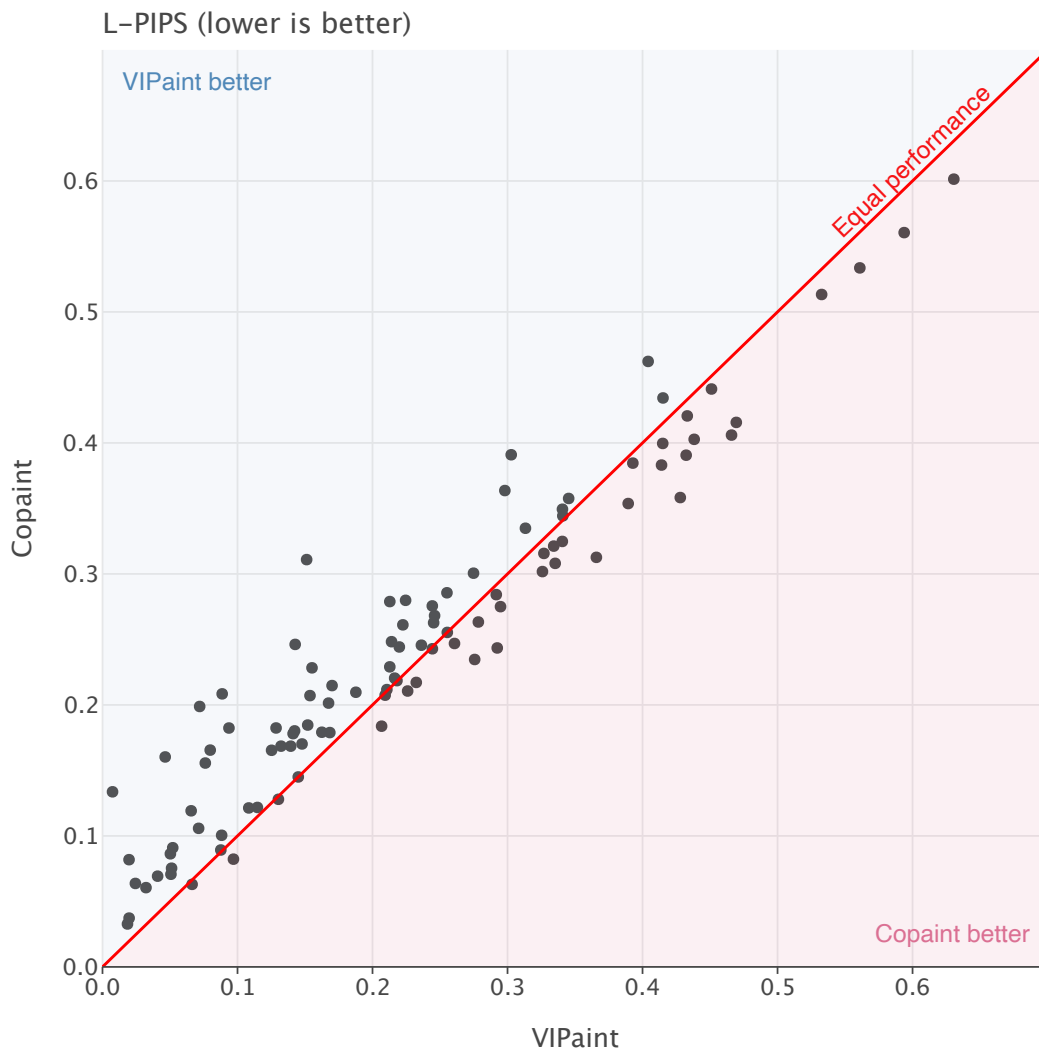
Overall, gradient based methods like DPS take longer with an LDM prior because of the use of a decoder per gradient step. PSLD additionally utilizes the encoder and hence, takes longer than DPS.

Table 3: Time (in mins) & #NFEs of denoising network per Inference method using EDM Prior (top) and LDM Prior (bottom)

Red-Diff	Blended	DPS	RePaint	CoPaint	CoPaint-TT	VIPaint-2	VIPaint
(1.13, 1000)	(1.13, 1000)	(2.55, 1000)	(2.8, 4700)	(2.6, 500)	(5.4, 1000)	3.3 = (1.5, 150) (optimization) (1.8, 700) (sampling)	11.8 = (10, 900) (optimization) (1.8, 700) (sampling)
Dataset	Blended	DPS	PSLD	VIPaint-2		VIPaint	
ImageNet256	(4, 1000)	(10, 500)	(12.4, 500)	10 = (2, 150) (optimization) (8, 400) (sampling)		18 = (10, 1250) (optimization) (8, 400) (sampling)	
LSUN	(1.3, 1000)	(5.1, 500)	(7.0, 500)	6.4 = (2.1, 150) (optimization) (4.3, 400) (sampling)		10 = (5.53, 750) (optimization) (4.3, 400) (sampling)	

G COMPUTATIONAL RESOURCES

All experiments were conducted on a system with 4 Nvidia A6000 GPUs.



1276 Figure 11: Paired comparison of LPIPS scores for VIPaint-2 and CoPaint with time-travel (CoPaint-
 1277 TT) on the Imagenet64 “Random Mask” inpainting task (expanding on the experiment shown in
 1278 table 1. Each point shows the mean LPIPS score across 10 sampled completions of the masked
 1279 image, with the x and y coordinates showing the VIPaint and CoPaint-TT scores respectively. Ad-
 1280 ditionally, we validated that VIPaint improves on CoPaint-TT using a one-sided paired t-test on the
 1281 mean LPIPS scores of each method. We found that the improvement was statistically significant
 1282 with a p-value of **4.133e-05**. As the normality assumption of the t-test may not hold, we also verified
 1283 the results using a nonparametric Wilcoxon signed ranked test, which indicated a statistically signifi-
 1284 cant improvement with a p-value of **0.000114**

1285 H ADDITIONAL EXPERIMENTS

1286 H.1 ANALYSIS OF IMAGENET RESULTS

1287 Fig. 11 shows details of the comparison between VIPaint and CoPaint with time-travel over 100
 1288 randomly selected images from the Imagenet-64 task.

1289 H.2 LINEAR INVERSE PROBLEMS

1290 For linear inverse problems other than inpainting, we consider the following tasks: (1) Gaussian
 1291 deblurring and (2) super resolution. For Gaussian deblurring, we use a kernel with size 61×61
 1292 with standard deviation 3.0. For super resolution, we use bicubic downsampling, similar setup as

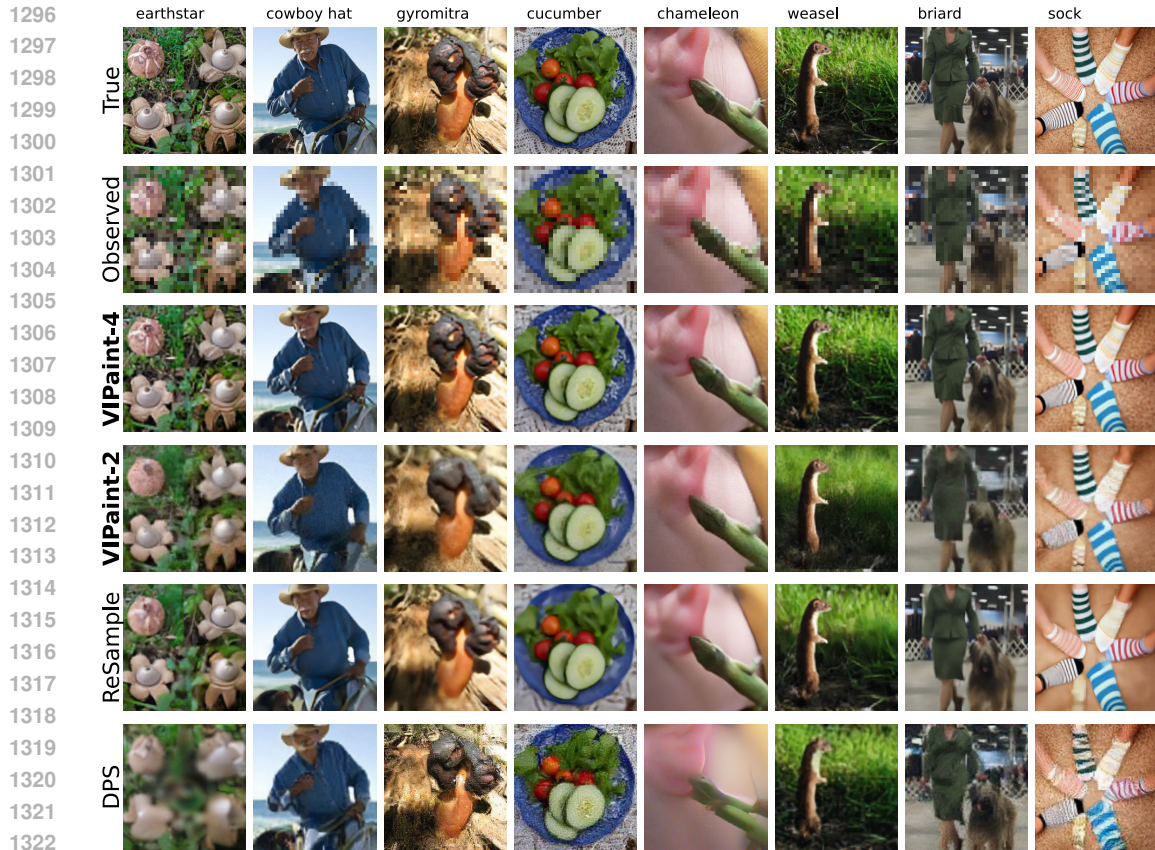


Figure 12: Qualitative results on Imagenet256 for Super Resolution. We see that DPS produces completely blurry images. We see improvements with ReSample. In contrast, VIPaint-4 leads to samples closer to the true image and produces *very* realistic images.

Task	ImageNet256				ImageNet64	
	Super-resolution 4x		Gaussian Deblur		Gaussian Deblur	
Metric	LPIPS ↓	PSNR ↑	LPIPS ↓	PSNR ↑	LPIPS ↓	PSNR ↑
VIPaint-4	0.33	19.31	<u>0.44</u>	<u>17.90</u>	0.306	13.47
VIPaint-2	0.46	16.36	0.48	16.35	0.305	13.60
ReSample	<u>0.395</u>	<u>18.410</u>	0.435	18.03	–	–
PSLD	0.67	7.77	0.583	0.022	–	–
DPS	0.579	12.99	0.595	12.608	0.319	13.43

Table 4: Quantitative results (LPIPS, PSNR) for solving linear inverse problems on ImageNet256 using LDM priors and ImageNet64 using EDM priors. Best results are in bold and second best results are underlined. For nonlinear deblurring.

Chung et al. (2023). Even though the focus of VIPaint is to remedy inconsistencies in image completion tasks, it can also be extended to linear inverse problems like Super Resolution and Gaussian Deblurring.

We compare the performance of VIPaint with ReSample, PSLD & DPS for ImageNet256 dataset using the LDM prior and for the pixel-based model, we include results for Gaussian Deblurring-compring VIPaint with DPS. Since the Peak-Signal-to-Noise-Ratio (PSNR) is well defined for such tasks, we report it along with LPIPS in Table 4. Some qualitative plots are in Fig. 12 and 13. A quantitative analysis is reported in Table ?? and qualitative results in Fig. 14 and ?. We see that VIPaint shows strong advantages over ReSample, DPS and Red-Diff for complex image datasets like ImageNet.

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

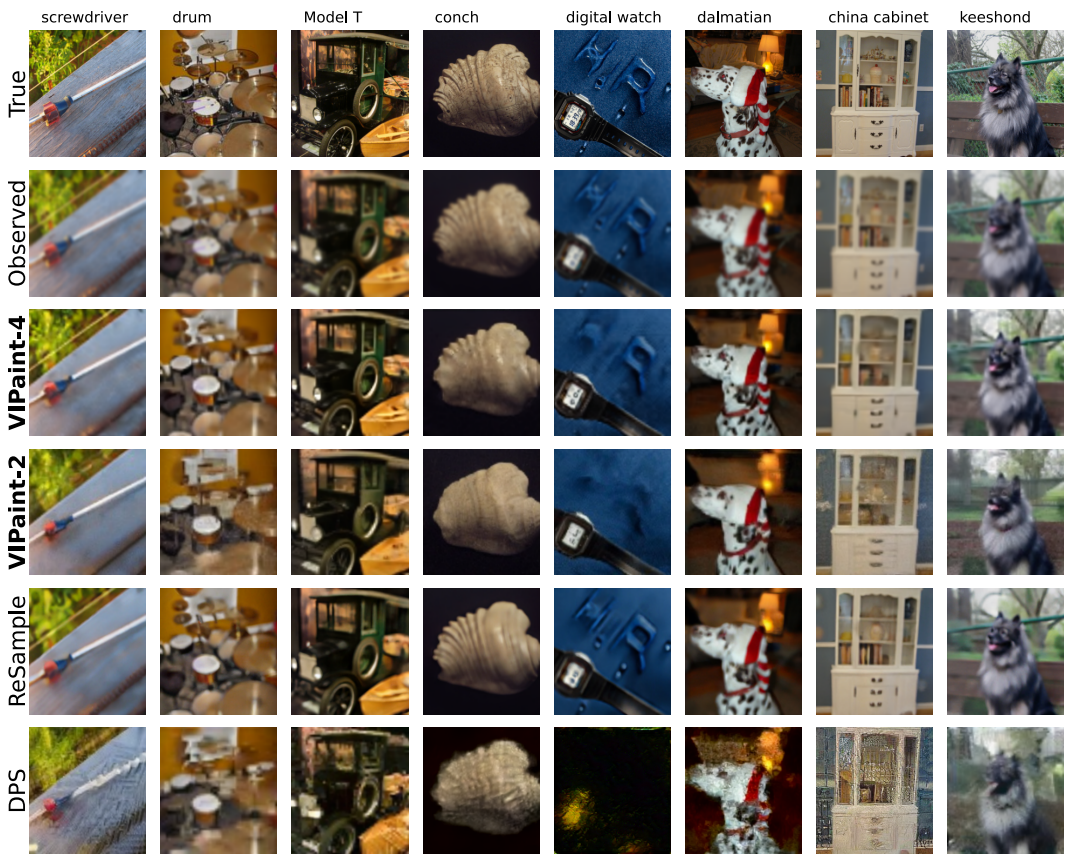
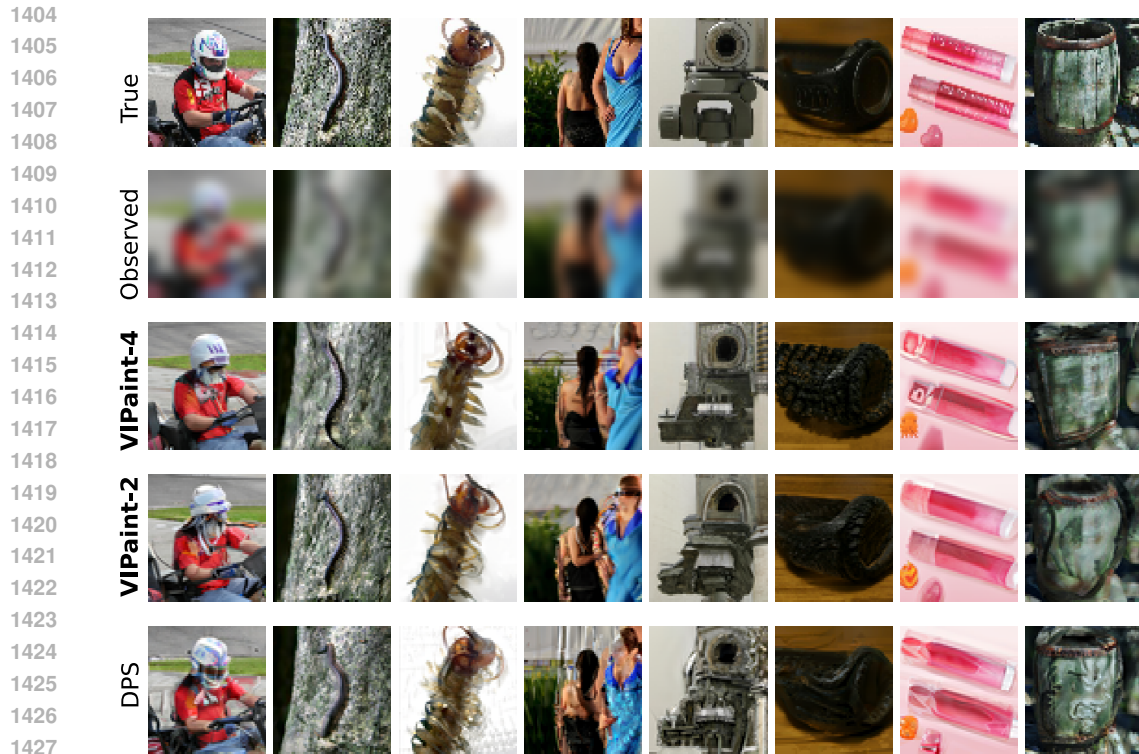


Figure 13: Qualitative results on Imagenet256 Gaussian DeBlurring using LDM prior.



1428
1429 Figure 14: Qualitative results for Gaussian DeBlurring using EDM prior for ImageNet64. We see VIPaint
1430 leads to samples closer to the true image and produces *more* realistic images.

1431 H.3 LARGE-MASK IMAGE INPAINTING FOR LSUN

1432 We show some qualitative figures for large masking ratios in Fig. 15 for LSUN-Church.
1433

1434 H.4 SMALL-MASK IMAGE INPAINTING FOR LSUN, IMAGENET256

1435 We show some qualitative figures for small masking ratios (upto 20% of the image is corrupted) in
1436 Fig. 16 and 17 for ImageNet-256 and LSUN-Church datasets respectively.
1437

1438 H.5 VIPAINT CAPTURES MULTI-MODAL POSTERIOR

1439 In addition to producing valid inpaintings, we show multiple samples per test image for all datasets
1440 we consider in Fig. 18, 24
1441

1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 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

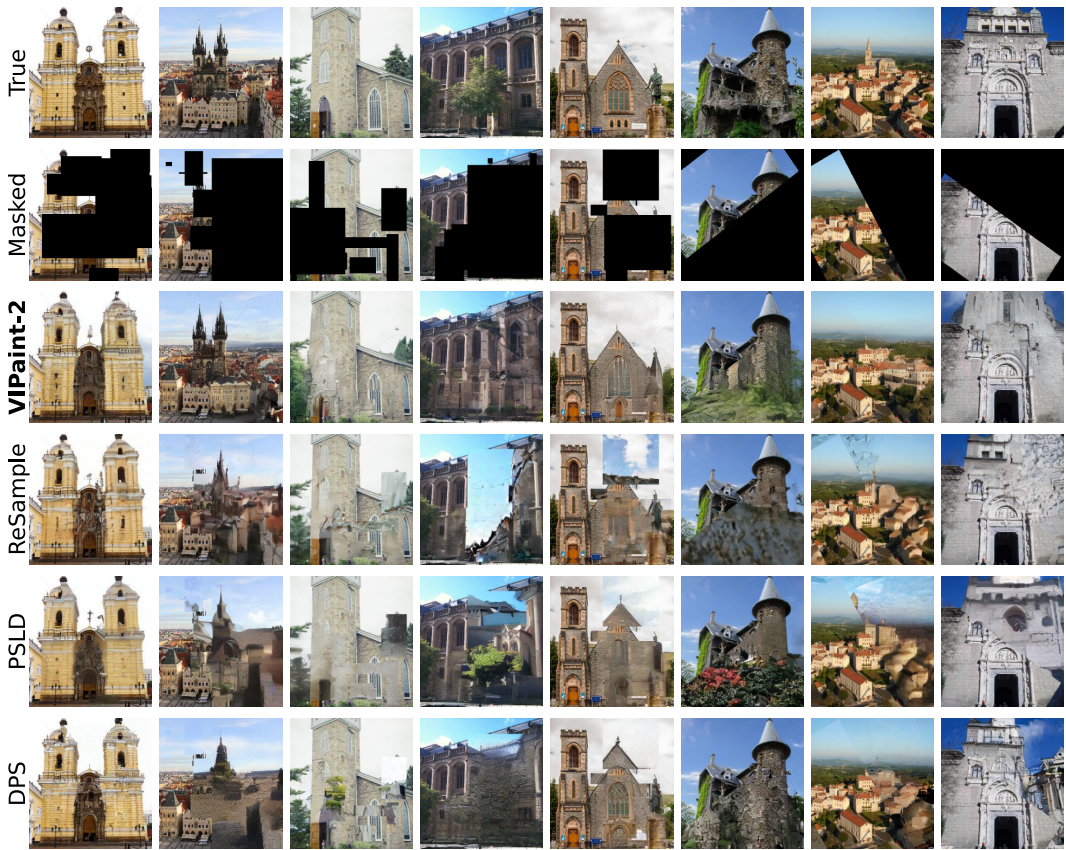


Figure 15: Qualitative results for LSUN-church dataset using LDM prior for the tasks of image inpainting with large masks. We see that VIPaint-2 can inpaint the images consistently and without any artifacts at the mask borders.

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

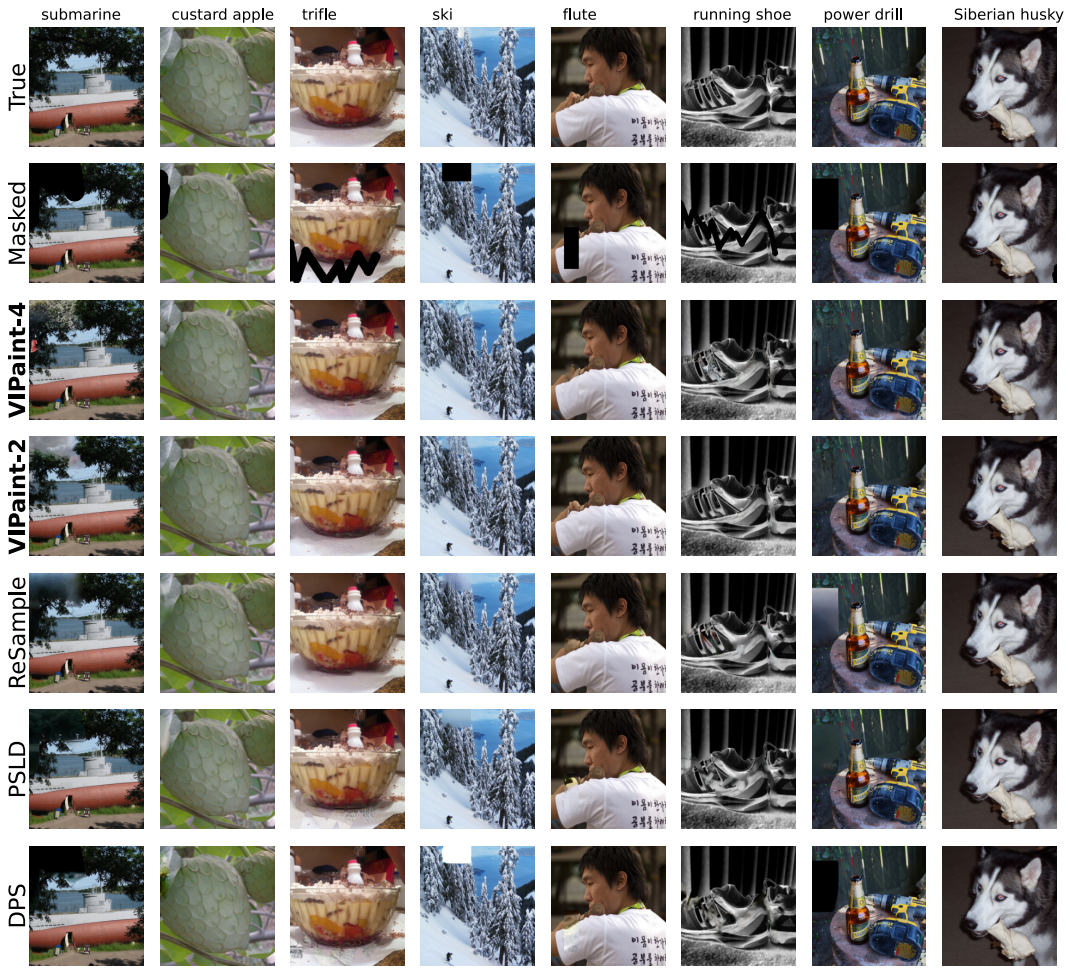


Figure 16: Qualitative results on the performance across methods for small masking ratios for ImageNet256 dataset using LDM prior. All methods seem to perform reasonably well in this regime.

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
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619



Figure 17: Qualitative results on the performance across methods for small masking ratios for LSUN-Church dataset using LDM prior. All methods seem to perform reasonably well in this regime. However, some minor artifacts start to show up for contiguous masks.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

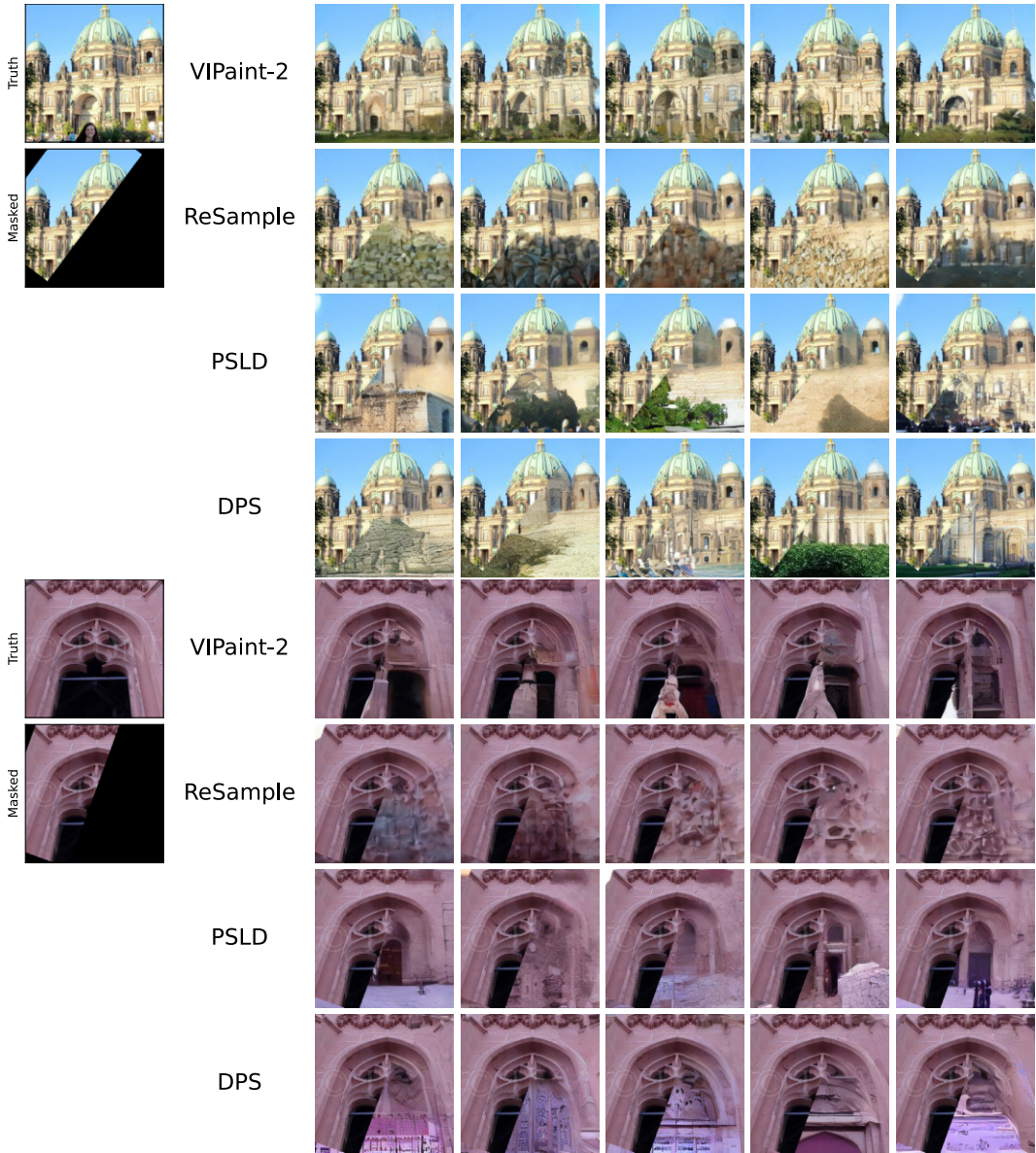


Figure 18: LSUN diversity results. Examples of diverse generation using VIPaint and baseline methods on LSUN using the same input and different initial noise.

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

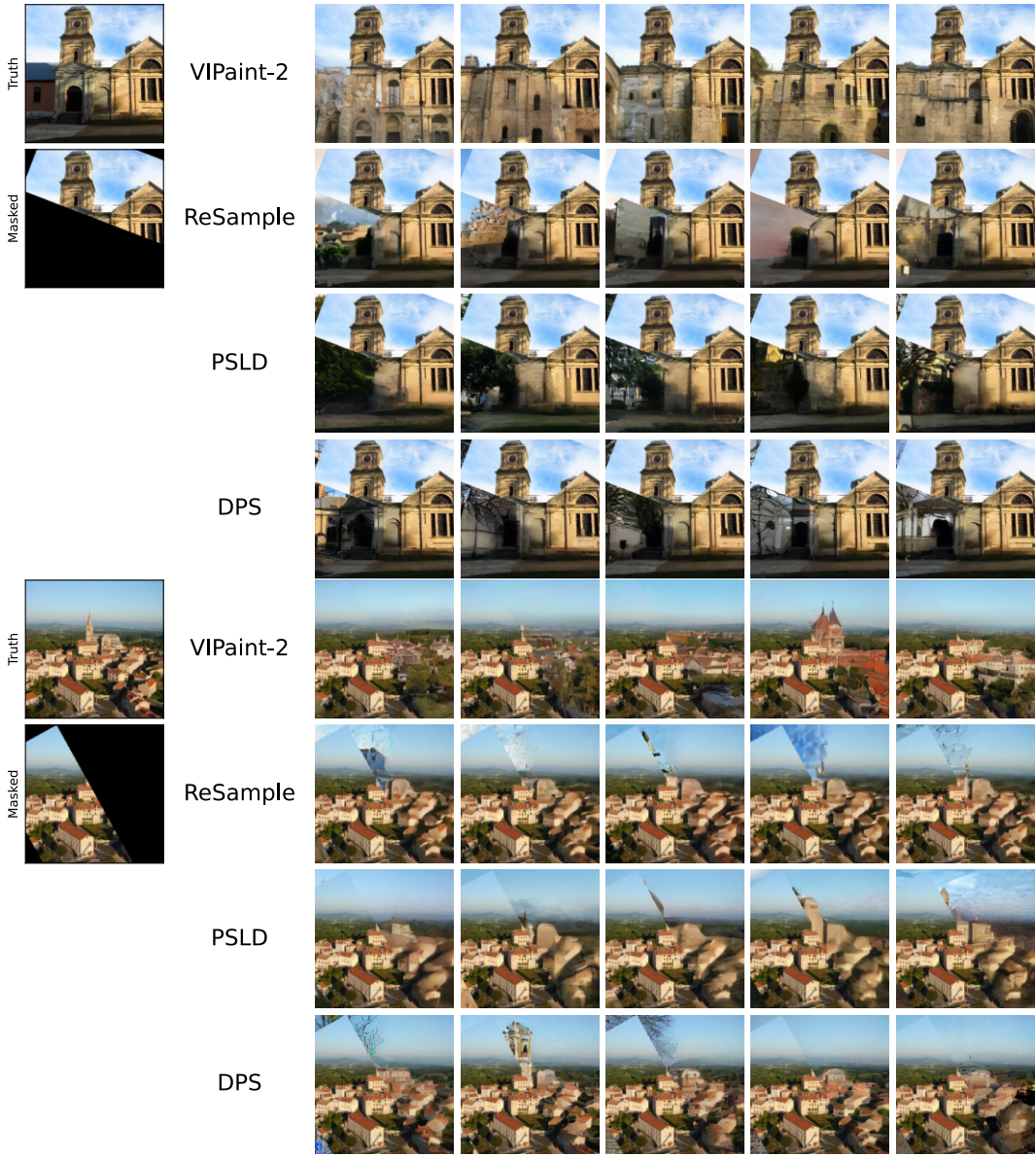
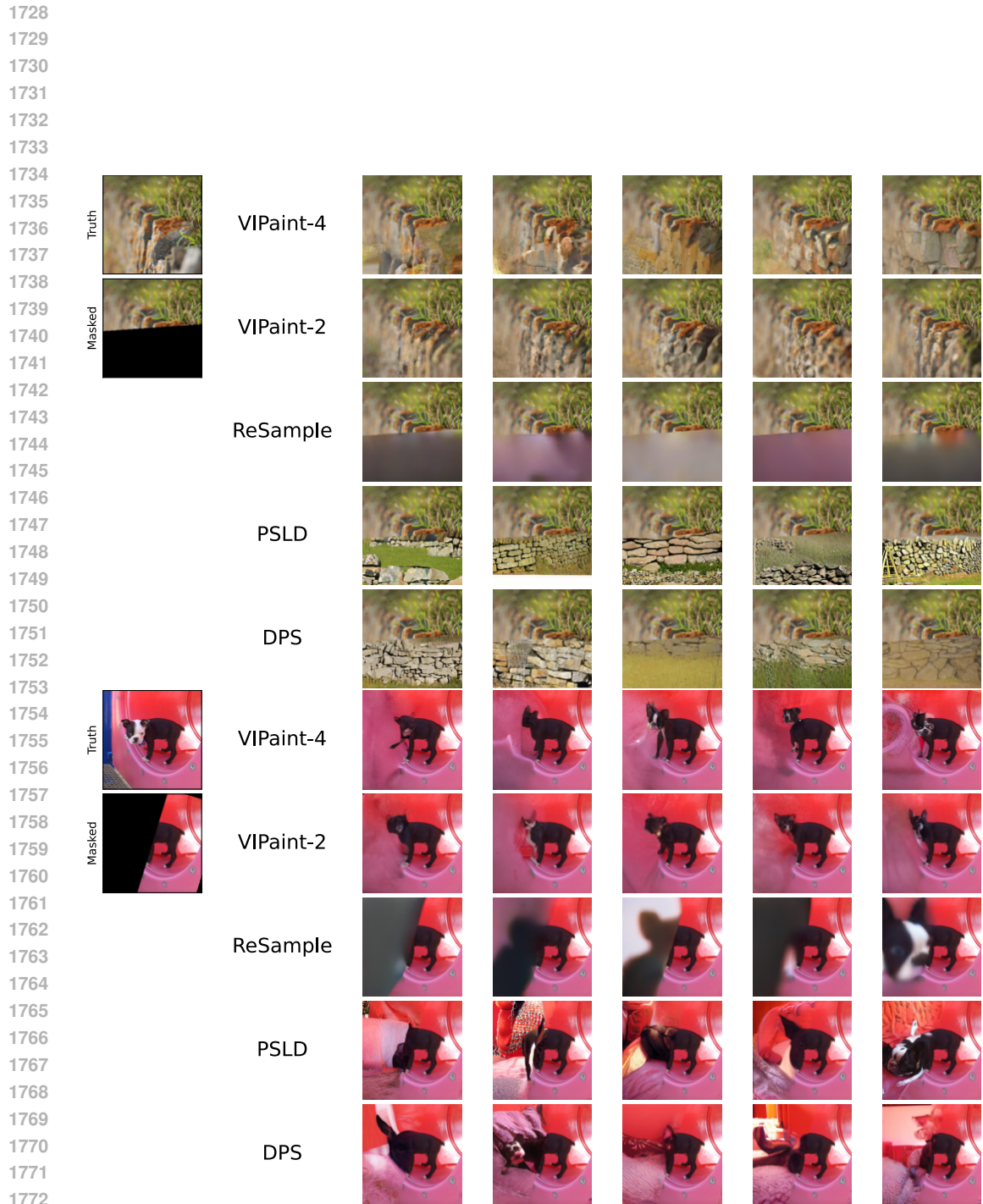


Figure 19: LSUN diversity results.



1773 Figure 20: ImageNet-256 diversity results. Some examples of diverse generation using VIPaint and
 1774 baseline methods on ImageNet using the same condition and masked input but with different
 1775 noise.
 1776
 1777
 1778
 1779
 1780
 1781

1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835



Figure 21: ImageNet diversity results.



Figure 23: ImageNet64 diversity results with the same class condition but different initial noise.

1944
 1945
 1946
 1947
 1948
 1949
 1950
 1951
 1952
 1953
 1954
 1955
 1956
 1957
 1958
 1959
 1960
 1961
 1962
 1963
 1964
 1965
 1966
 1967
 1968
 1969
 1970
 1971
 1972
 1973
 1974
 1975
 1976
 1977
 1978
 1979
 1980
 1981
 1982
 1983
 1984
 1985
 1986
 1987
 1988
 1989
 1990
 1991
 1992
 1993
 1994
 1995
 1996
 1997



Figure 24: Qualitative results for VIPaint diversity for ImageNet256 with LDM prior using different class conditioning. We see that VIPaint follows the input label and ensures consistency with the observed set of pixels.