

---

# Training-free Image Reconstruction via Visual Autoregressive Models

---

Linqiao Yang\*

Department of Mathematics  
University of Michigan  
Ann Arbor, MI 48105  
joelyang@umich.edu

## Abstract

Recent advances in diffusion and score-based generative models have dramatically improved the quality of learned priors for solving ill-posed imaging inverse problems. However, diffusion-based solvers typically require hundreds of function evaluations, leading to slow inference and high computational cost, and may hallucinate structures inconsistent with the measurements. In parallel, Visual Autoregressive Modeling (VAR) has emerged as a new paradigm for image generation that performs coarse-to-fine “next-scale prediction” and achieves diffusion-level image quality with significantly faster sampling. In this work, we take a benchmarking perspective and explore how well an off-the-shelf VAR model can serve as a training-free prior for image reconstruction. We formulate inverse problems in the latent token space of VAR and propose a two-stage inference scheme that combines hard token injection at coarse scales with measurement-gradient-guided logit refinement at finer scales. Our method requires no task-specific retraining and is evaluated on ImageNet-1k under five degradations: masking (inpainting), Gaussian blur, motion blur, nonuniform blur, and  $4\times$  super-resolution. Quantitatively, our VAR-based solver achieves PSNR around 19.7 dB and SSIM around 0.85 across these tasks, which is substantially worse in PSNR than diffusion- and plug-and-play-based methods (typically 26–30 dB) but competitive in SSIM, while running in a small fraction of their runtime. These results highlight both the current limitations of discrete visual autoregressive priors for inverse problems and their potential as fast, training-free learned priors. We view this work as an initial benchmark and analysis of VAR-style models for imaging inverse problems, rather than a new state-of-the-art method.

## 1 Introduction

Many imaging applications require recovering a clean signal  $x \in \mathbb{R}^{H \times W \times 3}$  from indirect and noisy measurements

$$y = A(x) + n, \quad (1)$$

where  $A$  is a known degradation operator (e.g., blur, subsampling, masking) and  $n$  denotes noise. Classical reconstruction methods impose handcrafted priors—such as sparsity, total variation, or smoothness—to regularize this ill-posed inverse problem, but they often struggle to capture the rich statistics of natural images, especially under strong degradations.

Deep generative models provide a powerful alternative by learning expressive priors from large image collections. Score-based diffusion models and denoising diffusion probabilistic models (DDPMs)

---

\*Use footnote for providing further information about author (webpage, alternative address)—*not* for acknowledging funding agencies.

33 have been particularly successful in this role, enabling state-of-the-art reconstructions in MRI, CT, and  
34 generic image restoration tasks by treating inverse problems as posterior sampling with a learned score  
35 prior. Building on this idea, methods such as Denoising Diffusion Restoration Models (DDRM)[8] and  
36 Diffusion Posterior Sampling (DPS)[2] adapt pre-trained diffusion models to a wide range of linear  
37 and nonlinear inverse problems without task-specific retraining. Despite their strong performance,  
38 diffusion-based solvers are computationally expensive, often requiring hundreds of denoising steps,  
39 and may hallucinate unrealistic details when measurements are highly ambiguous.

40 Parallel to diffusion, *visual autoregressive* (VAR) [14] models have recently emerged as a compelling  
41 alternative for high-fidelity image generation. VAR redefines autoregressive learning on images as  
42 coarse-to-fine “next-scale prediction”: instead of raster-scanning pixels or patch tokens, a transformer  
43 predicts the token map of the next higher-resolution scale conditioned on all coarser-scale token maps.  
44 This paradigm achieves diffusion-level or better image quality on ImageNet  $256 \times 256$  while enjoying  
45 substantially faster sampling and favorable scaling laws reminiscent of large language models. VAR  
46 models have also demonstrated promising zero-shot capabilities for image editing, inpainting, and  
47 outpainting, hinting that they may serve as versatile visual priors.

48 However, most prior work on inverse problems has focused on diffusion and score-based models;  
49 there is little understanding of how to leverage discrete visual autoregressive models as training-  
50 free priors for reconstruction. Existing VAR-based applications for restoration typically require  
51 task-specific fine-tuning or carefully engineered conditioning, which limits flexibility and generality.

52 **Goal and perspective.** Rather than proposing a new state-of-the-art method, our goal in this project  
53 is to *benchmark* how a pre-trained VAR model behaves when used as a training-free prior for a range  
54 of imaging inverse problems. Concretely, we ask:

- 55 • How well can VAR-based guidance enforce consistency with diverse measurement operators  
56 compared to diffusion-based and plug-and-play methods?
- 57 • What are the failure modes and trade-offs of working in a discrete multi-scale token space  
58 for inverse problems?
- 59 • Can we obtain meaningful reconstructions at a much lower computational cost, even if the  
60 absolute image quality is not competitive yet?

61 **Contributions.** We make the following contributions:

- 62 • We formulate image reconstruction as posterior inference in the latent token space of a  
63 pre-trained VAR model, using a measurement consistency term as a likelihood and the VAR  
64 transformer as a prior.
- 65 • We propose a two-stage inference scheme combining hard token injection at coarse scales  
66 with gradient-guided logit updates at finer scales, providing a generic template for incorpo-  
67 rating measurement information into VAR sampling.
- 68 • We conduct a systematic benchmark on ImageNet-1k for five degradations—masking  
69 (inpainting), Gaussian blur, motion blur, nonuniform blur, and  $4 \times$  super-resolution—using  
70 a single pre-trained VAR checkpoint and report PSNR, SSIM, and runtime.
- 71 • We empirically find that our VAR-based solver underperforms diffusion- and plug-and-play-  
72 based methods in PSNR but achieves competitive SSIM and substantially lower runtime,  
73 and we analyze how design choices such as guidance strength and stage threshold affect  
74 performance.

75 Our experiments show that the VAR-based solver achieves visually plausible reconstructions with  
76 significantly shorter runtimes than DPS, while maintaining competitive PSNR and SSIM for mod-  
77 erately ill-posed settings. We also highlight failure cases and discuss opportunities to better bridge  
78 discrete autoregressive priors with continuous inverse problems.

79 **2 Related Work**

80 **2.1 Inverse problems and learned generative priors**

81 Using learned generative models as priors for inverse problems has become an active research area.  
82 Early work [5, 4, 9] employed generative adversarial networks (GANs) and VAEs as explicit priors or  
83 plug-and-play denoisers, but these models often struggled to match the expressiveness and stability  
84 needed for high-resolution reconstruction. Recent score-based generative models learn the gradient  
85 of the log data density and can be combined with the forward operator to define posterior sampling  
86 schemes for various imaging tasks. [3]

87 **2.2 Diffusion models for inverse problems**

88 Several diffusion-based frameworks[2, 8, 15] adapt pre-trained DDPMs to inverse problems in  
89 a training-free manner. DDRM formulates an approximate variational objective that connects  
90 unconditional diffusion models with a noisy linear measurement model, yielding an efficient posterior  
91 sampling algorithm applicable to denoising, deblurring, super-resolution, and colorization. DPS  
92 instead explicitly approximates the true posterior  $p(x | y)$  via a blended process that combines  
93 diffusion sampling with gradient steps on the measurement likelihood, enabling general nonlinear  
94 operators and noisy settings. Extensions [13] further analyze posterior sampling with latent diffusion  
95 models, provide convergence guarantees for linear operators, or design ODE-based DPS variants for  
96 PDE-constrained inverse problems.

97 These methods demonstrate that diffusion priors can be used in a broadly applicable and theoret-  
98 ically grounded way but retain the computational burden of iterative denoising and may produce  
99 hallucinations when the posterior is multi-modal or poorly aligned with the learned prior.

100 **2.3 Visual autoregressive models**

101 Visual autoregressive models treat images as sequences of discrete tokens and model their distribution  
102 via next-token prediction, in analogy to language models. Traditional AR approaches operate either  
103 on pixel-level sequences or on patch tokens obtained from a VQ-VAE tokenizer, but scaling them to  
104 high resolutions has been challenging.

105 Visual Autoregressive Modeling (VAR) revisits this idea by proposing a multi-scale tokenizer that  
106 decomposes an image into a sequence of token maps of increasing resolution, and a transformer  
107 trained to predict the next-scale token map given all previous scales. This “next-scale prediction”  
108 strategy materially improves both sample quality and efficiency: VAR surpasses strong diffusion  
109 transformers on ImageNet 256×256 while requiring far fewer autoregressive steps and exhibiting  
110 clear scaling laws and zero-shot editing behavior.

111 Subsequent work extends VAR along several dimensions, such as bitwise tokenization with effectively  
112 infinite vocabularies (Infinity)[6], flexible ground-truth prediction without residuals (FlexVAR)[7],  
113 and more efficient intra/inter-scale decoupling (M-VAR)[12]. VAR-style architectures have also been  
114 adapted to high-resolution text-to-image synthesis and other downstream tasks.

115 **2.4 Autoregressive models for restoration and inverse problems**

116 Compared to diffusion models, there is relatively little work systematically exploring visual au-  
117 toregressive models as priors for inverse problems. VAR itself demonstrates zero-shot inpainting  
118 and editing by conditioning on degraded inputs but does not provide a general posterior sampling  
119 framework analogous to DPS/DDRM. Some recent AR-based super-resolution models [10, 11] train  
120 specifically for restoration tasks, sometimes with joint modeling of low- and high-resolution pairs,  
121 but these approaches typically require new training for each operator or dataset. To our knowledge,  
122 our project is among the first to explicitly cast inverse problems as training-free posterior inference in  
123 the token space of a pre-trained VAR model and to investigate gradient-based measurement guidance  
124 for discrete multi-scale autoregressive transformers.

125 **3 Problem Setup**

126 We consider imaging inverse problems defined by a (possibly linear) operator  $A : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^M$   
127 and additive noise:

$$y = A(x) + n, \quad n \sim \mathcal{N}(0, \sigma^2 I). \quad (2)$$

128 Given  $A$ , the noise level  $\sigma$ , and the measurement  $y$ , our goal is to reconstruct an estimate  $\hat{x}$  that is  
129 both consistent with the observation and visually plausible as a natural image.

130 In this project, we study five representative operators, covering both linear and nonlinear inverse  
131 problem :

- 132 • **Masking (inpainting).**  $A$  applies a binary mask  $M$  that zeros out a subset of pixels, i.e.,  
133  $A(x) = M \odot x$ , where  $M$  is 0 on missing regions and 1 elsewhere. We consider both a  
134 central square mask and more arbitrary masks.
- 135 • **Gaussian deblurring.**  $A$  convolves  $x$  with a Gaussian kernel of standard deviation  $\sigma_{\text{blur}}$ ,  
136 followed by possible additive noise.
- 137 • **Motion deblurring.**  $A$  applies a 1D or 2D motion blur kernel (simulating camera motion)  
138 and optionally adds noise.
- 139 • **Nonuniform blur.**  $A$  applies a spatially varying blur kernel, such as a combination of local  
140 Gaussian kernels, to model nonuniform optical or motion effects.
- 141 • **4× super-resolution.**  $A$  downsamples  $x$  by a factor of 4 in each spatial dimension using  
142 bicubic interpolation and optionally adds noise. The measurement  $y$  is a low-resolution  
143 RGB image.

144 These operators cover common restoration settings and also align with prior diffusion-based evalua-  
145 tions, enabling qualitative comparison to the DPS baseline.

146 **4 Data**

147 We use the ImageNet-1k validation set as our evaluation benchmark. Images are first center-cropped  
148 and resized to  $256 \times 256$  RGB to match the resolution of the pre-trained VAR model.

149 **Sampling protocol.** For computational reasons, we randomly select 5000 images from the valida-  
150 tion set, stratified across object categories to maintain diversity. Each clean image  $x$  is then processed  
151 by one of the degradation operators described in Section 4 to produce the corresponding observation  
152  $y$ .

153 **Preprocessing.** All images are normalized to the range  $[0, 1]$  and then converted to the tensor format  
154 expected by the VAR tokenizer. We additionally standardize the measurements for some operators  
155 (e.g., subtract the global mean for deblur) to stabilize gradient computation. For super-resolution,  
156 we upsample  $y$  back to  $256 \times 256$  using bicubic interpolation before feeding it to the measurement  
157 operator  $A$  within the reconstruction loop to ensure dimensional consistency.

158 **Qualitative examples.** Figure 1 (adapted from our project presentation) illustrates example of  
159 degraded images and for all five operators, as well as reconstructions produced by our VAR-based  
160 solver.

161 Overall, the ImageNet validation data provides a challenging yet standard benchmark to assess how  
162 well the VAR prior generalizes beyond its unconditional image synthesis training.

163 **5 Approach**

164 **5.1 Pre-trained VAR as a visual prior**

165 We build on the public implementation and checkpoints of Visual Autoregressive Modeling. A VAR  
166 model consists of:

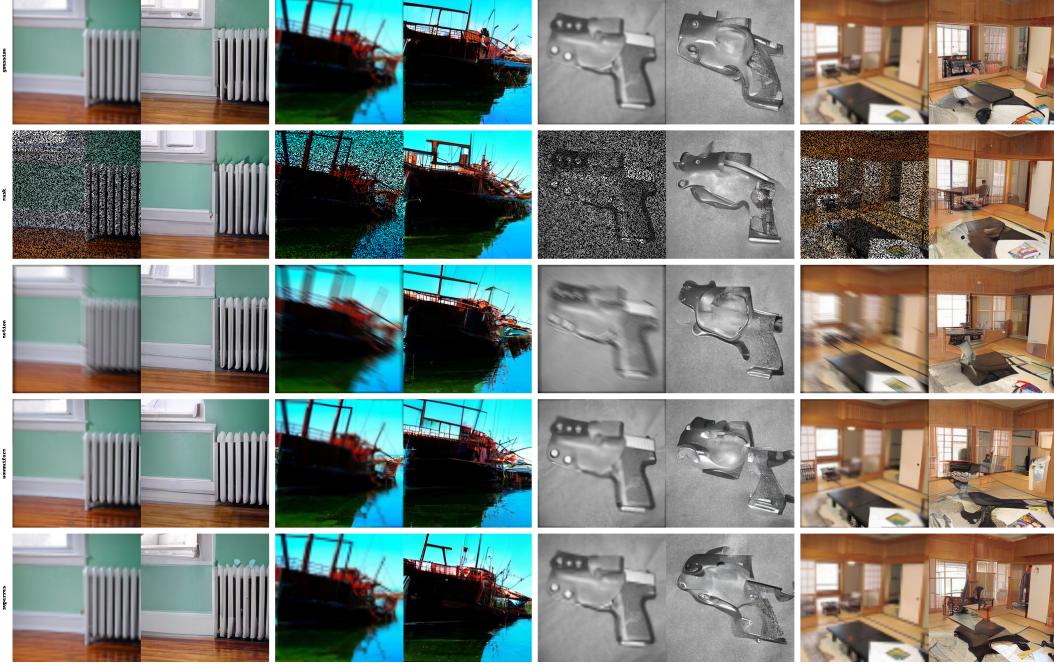


Figure 1: Enter Caption

- A multi-scale VQ-VAE-style tokenizer  $E$  that maps an image  $x$  into a sequence of  $K$  discrete token maps

$$\mathbf{z} = \{z^{(1)}, z^{(2)}, \dots, z^{(k)}\}, \quad z^{(k)} \in \{1, \dots, V\}^{h_k \times w_k}$$

167 , where  $z^{(1)}$  is the coarsest scale and  $z^{(K)}$  is the finest.

- 168 • A transformer  $T$  trained autoregressively to model the joint distribution

$$p_{\theta}(\mathbf{z}) = p_{\theta}(z^{(1)})p_{\theta}(z^{(k)}|z^{(\leq k)})$$

169 where  $z^{(\leq k)} = z^{(1)}, \dots, z^{(k)}$ .

- 170 • A decoder  $D$  that reconstructs an image from the token sequence:  $\tilde{x} = D(\mathbf{z})$ .

171 Unconditional sampling proceeds by first sampling  $z^{(1)}$  from  $p_{\theta}(z^{(1)})$ , then iteratively sampling  
172 each next-scale token map  $z^{(k+1)}$  from the conditional distribution predicted by the transformer, and  
173 finally decoding to the image space via  $D$ .

174 In our setting, we treat  $p_{\theta}(\mathbf{z})$  as a prior over natural images and seek to incorporate measurement  
175 information to obtain a posterior over tokens:

$$p(\mathbf{z} | y) \propto p_{\theta}(\mathbf{z}), p(y | \mathbf{z}), \quad (3)$$

176 where  $p(y | \mathbf{z})$  is induced by the forward operator  $A$  and the decoder  $D$  via  $\tilde{x} = D(\mathbf{z})$  and  
177  $y = A(\tilde{x}) + n$ .

## 178 5.2 Measurement loss and gradient guidance

179 Assuming Gaussian measurement noise, the negative log-likelihood up to a constant is

$$\mathcal{L}_{\text{meas}}(\mathbf{z}) = \frac{1}{2\sigma^2}, |A(D(\mathbf{z})) - y|_2^2. \quad (4)$$

180 Following the spirit of DPS and related posterior sampling methods, we use the gradient of this loss  
181 to guide the generative process toward measurement-consistent samples.

182 A challenge is that  $\mathbf{z}$  consists of discrete tokens. However, the VAR transformer internally maintains  
 183  $logits L^{(k)} \in \mathbb{R}^{h_k \times w_k \times V}$  for each token position and scale, from which tokens are sampled via  
 184 softmax. We treat these logits as continuous variables and compute the gradient

$$g^{(k)} = \nabla_{L^{(k)}} \mathcal{L}_{\text{meas}}(\mathbf{z}), \quad (5)$$

185 where backpropagation passes through the embedding lookup and decoder  $D$ .

186 We then define a guided logit update

$$\hat{L}^{(k)} = L^{(k)} - \alpha, g^{(k)}, \quad (6)$$

187 where  $\alpha > 0$  is a guidance strength hyperparameter. Intuitively, this nudges the logits in directions  
 188 that reduce the measurement loss, similar to a single gradient descent step in logit space. We then  
 189 resample tokens from the updated logits, e.g., via multinomial sampling or argmax for deterministic  
 190 decoding.

### 191 5.3 Two-stage inference with hard token constraints

192 Directly applying gradient guidance at all scales may lead to unstable behavior: coarse scales are  
 193 highly sensitive to global structure and may be over-corrected by local measurement gradients, while  
 194 fine scales benefit more from localized corrections. To address this, we adopt a two-stage inference  
 195 policy indexed by a scale threshold  $\gamma$ :

- 196 • **Low-stage (coarse scales,  $k < \gamma$ ): hard token injection.** At coarse scales, the for-  
 197 ward operator often provides strong global constraints—for example, the low-resolution  
 198 observation in super-resolution, or visible regions in inpainting. We exploit this by en-  
 199 coding an *approximate pseudo-ground-truth* from the degraded image, obtaining token  
 200 maps  $\tilde{z}^{(k)}_{\text{obs}} = E_k(\text{upsample}(y))$ . For positions that are reliably observed under  $A$  (e.g.,  
 201 unmasked pixels in inpainting, or low-resolution pixels aligned with the downsampling  
 202 grid), we directly replace the VAR-sampled tokens  $z^{(k)}$  with  $\tilde{z}^{(k)}_{\text{obs}}$  (hard token injection).  
 203 For the remaining positions (e.g., masked regions), we keep the autoregressive samples.
- 204 • **High-stage (fine scales,  $k \geq \gamma$ ): gradient-guided logits.** For finer scales, where the  
 205 operator constraints are weaker and more local, we rely on the gradient-based logit update  
 206 described above. After the transformer predicts logits  $L^{(k)}$  given  $z^{(\leq k-1)}$ , we decode the  
 207 current token sequence, compute  $\mathcal{L}_{\text{meas}}$ , backpropagate to obtain  $g^{(k)}$ , and adjust the logits  
 208 as  $\hat{L}^{(k)} = L^{(k)} - \alpha, g^{(k)}$  before sampling.

209 This policy is illustrated conceptually in Figure 2. When  $\gamma = K$ , we recover a pure hard-injection  
 210 scheme; when  $\gamma = 1$ , we recover pure gradient-guided logits. In practice, we set  $\gamma$  to an intermediate  
 211 scale (e.g.,  $\gamma = 8$  for  $K = 10$ ) so that the global structure is anchored by tokens derived from the  
 212 observation while high-frequency details are adjusted via gradients.

### 213 5.4 Algorithm outline and implementation details

214 Algorithm 1 summarizes the full procedure for a single reconstruction.

215 In practice, we implement this in PyTorch using the official VAR tokenizer and transformer check-  
 216 points. The forward operator  $A$  and its gradients are implemented in differentiable form (convolutions  
 217 for blur, strided interpolation for super-resolution, masking for inpainting), enabling efficient back-  
 218 propagation. We fix the guidance step to a single gradient update per scale and tune the guidance  
 219 strength  $\alpha$  on a small validation subset.

## 220 6 Experiments

### 221 6.1 Experimental setup

222 **Baselines.** Our main focus is to benchmark the VAR-based solver against established priors. For  
 223 quantitative comparison, we report:

- 224 • **PnP-ADMM** [1]. A plug-and-play ADMM method using a deep denoiser as a prior.

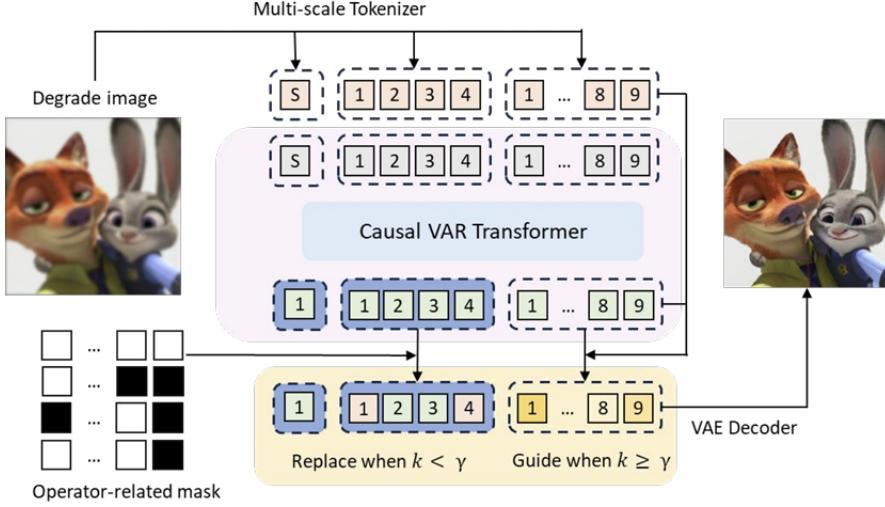


Figure 2: A method illustration when  $K = 3, \gamma = 2$

---

**Algorithm 1** Training-free Image Reconstruction via Guided VAR

---

**Require:** Observed degraded image  $y$ ; Degradation operator  $A$ ; VAE Decoder Decode; VAR Transformer  $T$ ; Hard-injection threshold  $\gamma$ ; Gradient step size  $\alpha$ .

**Ensure:** Reconstructed image  $\hat{x}$ .

- 1: Initialize token maps  $z^{(1)}, \dots, z^{(K)}$  as empty.
- 2: **for**  $k = 1$  to  $K$  **do**
- 3:   Predict logits  $L^{(k)} = T(z^{(\leq k-1)})$
- 4:   **if**  $k < \gamma$  **then** ▷ Hard-Injection Stage
- 5:     Encode upsampled observation (or visible pixels) to get  $\tilde{z}_{\text{obs}}^{(k)}$
- 6:     Compute operator-related mask  $M^{(k)}$  for fully determined positions
- 7:      $z^{(k)} \leftarrow \text{Sample}(\text{softmax}(L^{(k)}))$
- 8:      $z^{(k)}[M^{(k)}] \leftarrow \tilde{z}_{\text{obs}}^{(k)}[M^{(k)}]$  ▷ Inject hard token constraints
- 9:   **else** ▷ Gradient-Guided Stage
- 10:     Temporarily sample  $z_{\text{temp}}^{(k)} \sim \text{softmax}(L^{(k)})$
- 11:     Form current token sequence  $\mathbf{z}^{(\leq k)} \leftarrow (z^{(\leq k-1)}, z_{\text{temp}}^{(k)})$
- 12:     Decode current reconstruction  $\tilde{x} = \text{Decode}(\mathbf{z}^{(\leq k)})$
- 13:     Compute measurement loss  $\mathcal{L}_{\text{meas}} = \|A(\tilde{x}) - y\|^2$
- 14:     Backpropagate:  $g^{(k)} = \nabla_{L^{(k)}} \mathcal{L}_{\text{meas}}$
- 15:     Update logits:  $\hat{L}^{(k)} = L^{(k)} - \alpha g^{(k)}$
- 16:     Resample  $z^{(k)} \sim \text{softmax}(\hat{L}^{(k)})$  ▷ Guide logits by gradient
- 17:   **end if**
- 18: **end for**
- 19:  $\hat{x} \leftarrow \text{Decode}(\mathbf{z}^{(\leq K)})$
- 20: **return**  $\hat{x}$

---

- 225           • **DDRM** [8]. A diffusion-based restoration model using an unconditional diffusion prior for  
226           linear inverse problems.
- 227           • **DPS** [2]. Diffusion Posterior Sampling, which combines diffusion sampling with  
228           measurement-gradient steps.
- 229           • **DPnP** [15]. A recent diffusion plug-and-play method that refines the measurement-  
230           consistency step.
- 231           • **Ours (VAR guidance)**. The training-free solver described in Section ??.

232 **Metrics.** We evaluate reconstruction quality using peak signal-to-noise ratio (PSNR) and structural  
 233 similarity index (SSIM) computed between the reconstruction  $\hat{x}$  and the ground truth  $x$  on the  
 234 luminance channel. Runtime is measured as the wall-clock time per image, averaged over 5000  
 235 images, on a NVIDIA L40 GPU with batch size 8.

236 **Hyperparameters.** For our VAR-based solver, we fix the scale threshold  $\gamma = 8$  and guidance  
 237 strength  $\alpha = 1$  based on a small tuning set and then apply the same values to the full evaluation set.  
 238 Unless otherwise noted, we use deterministic sampling (argmax over logits) within each stage to  
 239 reduce variance.

## 240 6.2 Quantitative results

241 Table 1 reports results for three inverse tasks (Gaussian blur, motion blur, and  $4\times$  super-resolution)  
 242 on the ImageNet-1k validation subset.

243 We observe that our VAR-based solver is clearly behind diffusion- and plug-and-play-based baselines  
 244 in terms of PSNR: our method yields around 19.7 dB across tasks, whereas existing methods typically  
 245 achieve 26–30 dB. In contrast, our SSIM is competitive or even slightly higher: for example, in  
 246 Gaussian blur our SSIM 0.8518 exceeds DPnP’s 0.8360, and in super-resolution our SSIM 0.8500  
 247 is higher than all diffusion baselines in Table 1. This suggests that the VAR prior captures global  
 248 structure and local correlations reasonably well but suffers from intensity-level inaccuracies that hurt  
 249 PSNR.

250 In terms of runtime, our VAR-guidance method reconstructs an image in approximately 0.17 seconds  
 251 per image on a NVIDIA L40 GPU when using a single gradient step per scale, whereas diffusion-  
 252 based methods typically require hundreds of denoising steps and are reported to be one to two orders  
 253 of magnitude slower in similar settings. Taken together, the benchmark shows that VAR-guided recon-  
 254 struction trades off significantly lower PSNR for fast inference and decent SSIM, highlighting both a  
 performance gap and a potential computational advantage. :contentReference[oaicite:2]index=2

Table 1: Quantitative Comparison on Image Reconstruction across Different Inverse Tasks (PSNR/SSIM).

Method	Gaussian deblur		Motion deblur		Super-resolution ( $4\times$ )	
	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
PnP-ADMM	26.88	0.7855	26.55	0.7655	26.61	0.7634
DDRM	27.05	0.7819	—	—	<b>29.47</b>	<u>0.8437</u>
DPS	<u>28.83</u>	0.8212	27.87	0.8035	29.45	0.8379
DPnP	<b>29.24</b>	<u>0.8360</u>	<b>30.21</b>	<b>0.8527</b>	29.32	0.8407
<b>Ours (VAR Guidance)</b>	19.73	<b>0.8518</b>	19.74	<u>0.8519</u>	19.74	<b>0.8500</b>

255

## 256 6.3 Qualitative analysis

257 Figure 1 shows representative examples for masking (inpainting), Gaussian blur, motion blur, nonuni-  
 258 form blur, and  $4\times$  super-resolution. For all operators, the VAR-based reconstructions generally  
 259 preserve the overall layout and coarse semantic content of the scene: large objects appear in roughly  
 260 the correct locations, and background/foreground structure is often recognizable.

261 At the same time, the qualitative results also expose several important limitations. First, we frequently  
 262 observe geometric distortions and shape deformation: object boundaries can become wavy or warped,  
 263 straight edges may bend, and local parts (e.g., eyes, limbs, wheels) are sometimes misaligned or  
 264 merged together. Second, the base VAR model’s understanding of object semantics is still limited.  
 265 In complex scenes or under strong degradations, the reconstructions may confuse object categories,  
 266 mix attributes from different objects, or fill in missing regions with semantically plausible but clearly  
 267 incorrect content (e.g., hallucinating an extra limb or an implausible facial structure).

268 Compared to diffusion-based baselines (based on visualizations reported in the DPnP and DPS litera-  
 269 ture), our VAR-guided reconstructions tend to be smoother and less sharp in high-frequency textures,

270 which is consistent with the lower PSNR. However, they usually remain globally coherent and avoid  
 271 some of the more extreme hallucinations sometimes observed in diffusion-based methods when the  
 272 measurement is extremely ambiguous. For aggressive masking, our method often “hallucinates”  
 273 generic content in large missing regions (such as a smooth background or an averaged-out object),  
 274 illustrating that while the VAR prior provides a reasonable inductive bias for natural images, its  
 275 semantic understanding is not yet strong enough to reliably reconstruct fine object details in severely  
 276 ill-posed settings.

277 **6.4 Ablation studies**

278 **Effect of guidance strength  $\alpha$ .** We vary the guidance strength  $\alpha$  while keeping the stage threshold  
 279  $\gamma = 8$  fixed. Table 2 reports average PSNR/SSIM across all five tasks and the average runtime.  
 280 For very small  $\alpha$ , reconstructions are close to unconditional VAR samples and may deviate from  
 281 the observation  $y$ , especially in super-resolution where low-frequency content is underfit. As  $\alpha$   
 282 increases, PSNR and SSIM initially improve as the method better enforces measurement consistency,  
 283 but beyond a certain point, larger gradients cause over-correction and visual artifacts such as ringing  
 284 or repetitive patterns. We find that  $\alpha = 1$  already yields reasonable performance while keeping  
 285 runtime low.

Table 2: Ablation Study on Gradient Steps  $\alpha$  ( $\gamma = 8$ ) Across All Inverse Tasks.

$\alpha$	Masking		Gaussian Blur		Motion Blur		Nonuniform Blur		Super-res (4×)		Avg. Time
(Steps)	PSNR (↑)	SSIM (↑)	PSNR (↑)	SSIM (↑)	PSNR (↑)	SSIM (↑)	PSNR (↑)	SSIM (↑)	PSNR (↑)	SSIM (↑)	Avg. Time
1	19.71	0.853	19.73	0.852	19.74	0.852	<b>19.78</b>	0.853	19.71	0.850	<b>0.169</b>
3	<b>19.74</b>	0.853	19.72	0.853	<b>19.84</b>	<b>0.855</b>	19.77	0.854	<b>19.79</b>	<b>0.855</b>	0.337
5	19.74	<b>0.854</b>	<b>19.91</b>	<b>0.858</b>	19.78	0.855	19.82	<b>0.856</b>	19.74	0.853	0.506

286 **Effect of scale threshold  $\gamma$ .** We also study different choices of the stage boundary  $\gamma$  while fixing  
 287 the gradient step to  $\alpha = 1$ . Table 3 summarizes the results. Setting  $\gamma$  too low (i.e., relying mostly on  
 288 gradient guidance) yields poor structure and low PSNR/SSIM because the coarse-scale prediction  
 289 becomes unstable. As  $\gamma$  increases, we allow more scales to be anchored by hard injection from the  
 290 observation, which gradually improves reconstruction quality. We find that an intermediate to high  
 291 threshold ( $\gamma \approx 8$  for  $K = 10$ ) provides the best trade-off between respecting observed low-frequency  
 292 content and correcting high-frequency artifacts, while also being the fastest.

Table 3: Ablation Study on Hard-Injection Threshold  $\gamma$  (grad\_steps = 1) Across All Inverse Tasks.

$\gamma$	Masking		Gaussian Blur		Motion Blur		Nonuniform Blur		Super-res (4×)		Avg. Time
(Stage)	PSNR (↑)	SSIM (↑)	PSNR (↑)	SSIM (↑)	PSNR (↑)	SSIM (↑)	PSNR (↑)	SSIM (↑)	PSNR (↑)	SSIM (↑)	Avg. Time
4	14.06	0.509	14.14	0.527	14.19	0.532	14.25	0.532	14.07	0.516	0.346
5	15.04	0.605	14.97	0.599	15.06	0.596	15.14	0.610	15.05	0.605	0.293
6	16.15	0.683	16.24	0.696	16.19	0.694	16.23	0.693	16.13	0.680	0.253
7	17.89	0.781	17.87	0.778	17.91	0.780	17.91	0.781	17.83	0.775	0.211
8	19.71	0.853	19.73	0.852	19.74	0.852	19.78	0.853	19.71	0.850	0.169

293 **Discrete tokens vs. continuous image optimization.** An alternative approach would be to optimize  
 294 directly in image space using the measurement loss and a learned VAR-based score or discriminator.  
 295 In our experiments, the token-based guidance is more stable and efficient because it leverages  
 296 the multi-scale structure of VAR and enforces that intermediate reconstructions remain within the  
 297 representational range of the tokenizer, reducing the risk of adversarial-looking solutions. A detailed  
 298 comparison is left to future work.

299 **6.5 Limitations**

300 While promising as a fast prior, our approach has several limitations:

- 301 • **Limited to VAR’s training distribution.** Reconstructions are constrained by the pre-trained  
302 VAR’s prior; images significantly out-of-distribution (e.g., medical scans) are unlikely to be  
303 well reconstructed.
- 304 • **Discrete optimization challenges.** Gradient guidance on logits is only an approximation  
305 of true posterior inference over discrete tokens and can be sensitive to hyperparameters.  
306 Multi-step updates or alternative relaxations might be necessary for more ill-posed operators.
- 307 • **Operator-specific heuristics.** The hard token injection stage relies on operator-specific  
308 mappings from measurements to observed tokens; designing these mappings for more  
309 complex operators (e.g., non-linear optics) can be nontrivial.
- 310 • **Quality gap to diffusion.** Our benchmark clearly shows a sizable PSNR gap between  
311 VAR guidance and diffusion/plug-and-play methods. Closing this gap likely requires either  
312 stronger architectures, better training objectives, or more sophisticated inference schemes.

313 **7 Conclusion**

314 We presented a first exploration of using a pre-trained Visual Autoregressive Model as a training-  
315 free prior for image reconstruction across five standard imaging operators. By formulating inverse  
316 problems in the VAR token space and introducing a two-stage inference scheme combining hard token  
317 injection and gradient-guided logits, we obtained reasonable reconstructions on ImageNet-1k under  
318 masking, Gaussian blur, motion blur, nonuniform blur, and 4× super-resolution, with competitive  
319 SSIM but significantly lower PSNR compared to diffusion and plug-and-play baselines.

320 Our results should be interpreted as a *benchmark* rather than a claim of state-of-the-art performance:  
321 they highlight that, in their current form, discrete visual autoregressive priors are not yet competitive  
322 with diffusion models on standard inverse-problem metrics, but they can provide fast, training-free  
323 reconstructions with decent perceptual quality. We hope this benchmark encourages further research  
324 on:

- 325 • more principled posterior sampling algorithms for discrete multi-scale autoregressive mod-  
326 els,
- 327 • joint training objectives that explicitly incorporate inverse-problem performance into VAR’s  
328 learning,
- 329 • hybrid methods that combine the strengths of diffusion and autoregressive priors, and
- 330 • extending VAR-based reconstruction to higher resolutions and more challenging operators.

331 **Author Contributions**

332 Linqiao Yang: project conceptualization, literature review, method design, implementation, experi-  
333 ments, analysis, and writing of the report.

334 **References**

- 335 [1] Stanley H Chan, Xiran Wang, and Omar A Elgendy. Plug-and-play admm for image restoration:  
336 Fixed-point convergence and applications. *IEEE Transactions on Computational Imaging*,  
337 3(1):84–98, 2016.
- 338 [2] Hyungjin Chung, Jeongsol Kim, Michael T. Mccann, Marc L. Klasky, and Jong Chul Ye.  
339 Diffusion posterior sampling for general noisy inverse problems, 2024.
- 340 [3] Hyungjin Chung and Jong Chul Ye. Score-based diffusion models for accelerated mri, 2022.
- 341 [4] Muhammad Fadli Damara, Gregor Kornhardt, and Peter Jung. Solving inverse problems with  
342 conditional-gan prior via fast network-projected gradient descent, 2021.

- 343 [5] Hwan Goh, Sheroze Sheriffdeen, Jonathan Wittmer, and Tan Bui-Thanh. Solving bayesian  
344 inverse problems via variational autoencoders, 2021.
- 345 [6] Jian Han, Jinlai Liu, Yi Jiang, Bin Yan, Yuqi Zhang, Zehuan Yuan, Bingyue Peng, and Xiaobing  
346 Liu. Infinity: Scaling bitwise autoregressive modeling for high-resolution image synthesis. In  
347 *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 15733–15744,  
348 2025.
- 349 [7] Siyu Jiao, Gengwei Zhang, Yinlong Qian, Jiancheng Huang, Yao Zhao, Humphrey Shi, Lin Ma,  
350 Yunchao Wei, and Zequn Jie. Flexvar: Flexible visual autoregressive modeling without residual  
351 prediction. *arXiv preprint arXiv:2502.20313*, 2025.
- 352 [8] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration  
353 models, 2022.
- 354 [9] Jean Prost, Antoine Houdard, Andrés Almansa, and Nicolas Papadakis. Inverse problem  
355 regularization with hierarchical variational autoencoders. In *2023 IEEE/CVF International  
356 Conference on Computer Vision (ICCV)*, page 22837–22848. IEEE, October 2023.
- 357 [10] Yunpeng Qu, Kun Yuan, Jinhua Hao, Kai Zhao, Qizhi Xie, Ming Sun, and Chao Zhou. Visual  
358 autoregressive modeling for image super-resolution, 2025.
- 359 [11] Sudarshan Rajagopalan, Kartik Narayan, and Vishal M. Patel. Restorevar: Visual autoregressive  
360 generation for all-in-one image restoration, 2025.
- 361 [12] Sucheng Ren, Yaodong Yu, Nataniel Ruiz, Feng Wang, Alan Yuille, and Cihang Xie. M-var:  
362 Decoupled scale-wise autoregressive modeling for high-quality image generation. *arXiv preprint  
363 arXiv:2411.10433*, 2024.
- 364 [13] Litu Rout, Negin Raoof, Giannis Daras, Constantine Caramanis, Alexandros G. Dimakis, and  
365 Sanjay Shakkottai. Solving linear inverse problems provably via posterior sampling with latent  
366 diffusion models, 2023.
- 367 [14] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive  
368 modeling: Scalable image generation via next-scale prediction, 2024.
- 369 [15] Xingyu Xu and Yuejie Chi. Provably robust score-based diffusion posterior sampling for  
370 plug-and-play image reconstruction. *Advances in Neural Information Processing Systems*,  
371 37:36148–36184, 2024.