
FAVAE-Effective Frequency Aware Latent Tokenizer

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 *Latent generative models have shown remarkable progress in high-fidelity image*
2 *synthesis, typically using a two-stage training process that involves compressing*
3 *images into latent embeddings via learned tokenizers in the first stage. The quality*
4 *of generation strongly depends on how expressive and well-optimized these latent*
5 *embeddings are. While various methods have been proposed to learn effective*
6 *latent representations, the reconstructed images often lack realism, particularly in*
7 *textured regions with sharp transitions, due to loss of fine details governed by high*
8 *frequencies. We conduct a detailed frequency decomposition of existing state-of-*
9 *the-art (SOTA) latent tokenizers and show that conventional objectives inherently*
10 *prioritize low-frequency reconstruction, often at the expense of high-frequency*
11 *fidelity. Our analysis reveals these latent tokenizers exhibit a bias toward low-*
12 *frequency information, when jointly optimized, leading to over-smoothed outputs*
13 *and visual artifacts that diminish perceptual quality. To address this, we propose*
14 *a wavelet-based, frequency-aware variational autoencoder (FA-VAE) framework*
15 *that explicitly decouples the optimization of low- and high-frequency components.*
16 *This decoupling enables improved reconstruction of fine textures while preserving*
17 *global structure. Our approach bridges the fidelity gap in current latent tokenizers*
18 *and emphasizes the importance of frequency-aware optimization for realistic image*
19 *representation, with broader implications for applications in content creation,*
20 *neural rendering, and medical imaging.*

21 1 Introduction

22 Latent generative modeling [1, 2, 3, 4] has emerged as a cornerstone of modern content creation,
23 with recent advances demonstrating remarkable capabilities in synthesizing high-fidelity visual
24 content. These models typically operate in compressed latent spaces learned via autoencoders such
25 as VAEs [5, 6], where generation quality is directly influenced by the expressiveness of these latent
26 embeddings. While increasing latent dimensionality can enhance representational power, this often
27 leads to diminishing visual outputs and increased computational cost. VAVAE [7] addresses this by
28 aligning the latent space with pre-trained vision foundation models to further improve convergence and
29 generative realism. However, despite such improvements, generated outputs or even the reconstructed
30 outputs often lack sharp textures and fine details, particularly in regions dominated by the high-
31 frequency information such as text on images. This limits perceptual realism and results in overly
32 smooth outputs. Prior works [8, 9, 10] have explored architectural and spectral enhancements to
33 inject high-frequency signals, but a systematic frequency-level analysis of how these current latent
34 tokenizers influence reconstruction quality, especially in differentiating low- vs. high-frequency
35 information remains absent.

36 In this work, we conduct a comprehensive frequency-based analysis of the reconstruction behavior
37 of state-of-the-art latent tokenizers used in generative pipelines [7, 6, 11, 12, 13], focusing on
38 latent diffusion and autoregressive models, which represent the current forefront of high-quality

39 visual generation. These models typically follow a two-stage pipeline: (1) learning quantized or
 40 non-quantized latent embeddings using a latent tokenizer like VAE or its variant, and (2) training a
 41 probabilistic generative model in the latent embedding space. Our analysis reveals a consistent bias
 42 during the first stage of learning latent embeddings: while low-frequency components are on par well
 43 reconstructed, high-frequency signals such as textures are poorly preserved. This low frequency bias
 44 over high frequencies in optimization contributes significantly to perceptual degradation.

45 To address this fidelity gap, we propose a novel
 46 **Frequency-aware VAE (FA-VAE)** latent tokenizer that
 47 explicitly decouples and separately optimizes low- and
 48 high-frequency components using wavelet decompo-
 49 sition. Our method learns distinct latent embeddings
 50 for low-high frequency subbands and later fuses them
 51 into a unified latent space representation. This simple
 52 design allows for better preservation of both global
 53 structure and fine details. Figure 1 presents a visual
 54 comparison between our method and the recent state-
 55 of-the-art VAVAE, highlighting improved preservation
 56 of fine details and sharp structures.

57 **Our main contributions are summarized as follows:**

- 58 ✓ We provide a frequency-based analysis of
 59 latent tokenizers used in latent generative
 60 pipelines, showing that existing VAE variants
 61 disproportionately emphasize low-frequency
 62 components at the expense of fine detail.
- 63 ✓ We introduce **FA-VAE**, a frequency-aware
 64 VAE framework that decouples optimization
 65 of low- and high-frequency subbands using
 66 wavelet decomposition, resulting in frequency
 67 aware expressive latent representations.

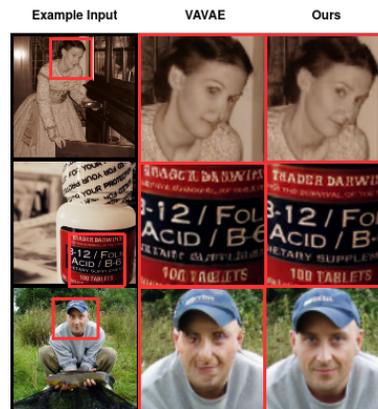


Figure 1: Visual comparison of reconstructions. From left to right: example original image, VAVAE reconstruction, and our approach (FA-VAE). The red color highlighted regions emphasize areas rich in textures, edges, and text. Our method better preserves high-frequency details and sharp structures, resulting in reconstructions visually closer to the input.

68 2 Method

69 2.1 Frequency Evaluation of Latent Embeddings

70 Latent tokenization models typically begin with a Variational Autoencoder (VAE) or its tokenizer
 71 variants, which learn compact embeddings $\mathbf{z} \in \mathcal{Z}$ from input data $\mathbf{x} \in \mathcal{X}$. Reconstruction quality is
 72 often measured by the pixel-level error.

$$\mathcal{L}_{\text{rec}} = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2, \quad \hat{\mathbf{x}} = D(E(\mathbf{x})),$$

73 where E and D are the encoder and decoder.

74 To study fidelity across frequency bands, we apply a discrete wavelet transform (DWT) with Haar
 75 filter to decompose signals into low- and high-frequency bands with one decomposition level:

$$\mathcal{W}(\mathbf{x}) = (\mathbf{x}_L, \mathbf{x}_H), \quad \mathcal{W}(\hat{\mathbf{x}}) = (\hat{\mathbf{x}}_L, \hat{\mathbf{x}}_H),$$

76 Then, we compute the frequency-aware reconstruction losses as follows:

$$\mathcal{L}_L = \|\mathbf{x}_L - \hat{\mathbf{x}}_L\|_2^2, \quad \mathcal{L}_H = \|\mathbf{x}_H - \hat{\mathbf{x}}_H\|_2^2.$$

77 Beyond pixel errors, perceptual similarity is captured by LPIPS [14], and distributional alignment
 78 by reconstruction FID (rFID) [15], computed from Inception features [16]. Together, combining
 79 frequency-aware losses ($\mathcal{L}_L, \mathcal{L}_H$), perceptual similarity (LPIPS), and rFID provides a more complete
 80 evaluation of reconstruction fidelity and expressiveness of the learned latent embeddings.

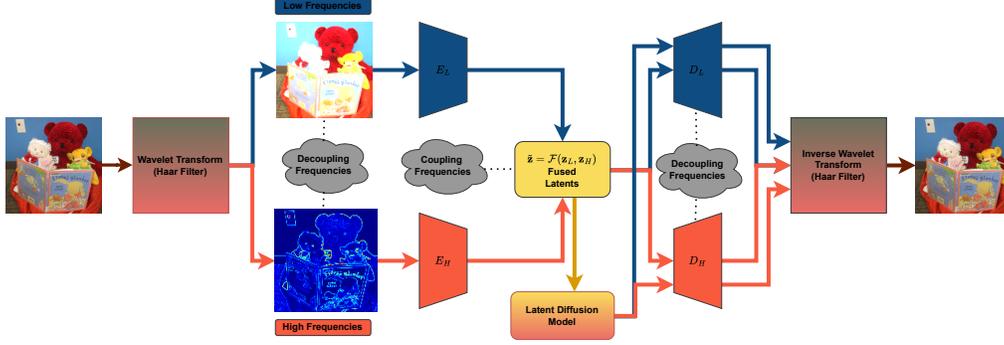


Figure 2: FA-VAE. The Overall pipeline of FA-VAE latent tokenizer.

81 2.2 Frequency-Aware VAE (FA-VAE)

82 The overview of FA-VAE is shown in Fig. 2. To enhance the fidelity of latent embeddings, we
 83 extend the standard VAE formulation by decoupling the input image into low- and high-frequency
 84 components, which are then learned independently. Given an input image $\mathbf{x} \in \mathcal{X}$, we apply a discrete
 85 wavelet transform $\mathcal{W}(\cdot)$ to obtain:

$$\mathcal{W}(\mathbf{x}) = (\mathbf{x}_L, \mathbf{x}_H),$$

86 where \mathbf{x}_L and \mathbf{x}_H denote the low- and high-frequency representations, respectively. We use the Haar
 87 filter for wavelet transformation, and apply the normalization strategy proposed in [17] to normalize
 88 both low and high frequency components. We employ separate encoder-decoder pairs (E_L, D_L) and
 89 (E_H, D_H) to learn frequency-specific latent embeddings:

$$\mathbf{z}_L = E_L(\mathbf{x}_L), \quad \mathbf{z}_H = E_H(\mathbf{x}_H).$$

90 **Low-Frequency Objective.** To learn low-frequency latent embeddings, we adopt a VA-VAE-style
 91 objective [7], incorporating a vision foundation alignment loss \mathcal{L}_{VF}^L [7], adversarial regularization
 92 \mathcal{L}_{GAN}^L inspired by VQGAN [18], and an additional perceptual loss \mathcal{L}_{LPIPS}^L [14] to improve visual
 93 quality. The total low-frequency objective is defined as:

$$\mathcal{L}_{\text{low}} = \mathcal{L}_{\text{rec}}^L + \beta \cdot \mathcal{L}_{\text{KL}}^L + \lambda_{\text{VF}} \cdot \mathcal{L}_{\text{VF}}^L + \lambda_{\text{GAN}} \cdot \mathcal{L}_{\text{GAN}}^L + \lambda_{\text{LPIPS}} \cdot \mathcal{L}_{\text{LPIPS}}^L$$

94 where:

$$\mathcal{L}_{\text{rec}}^L = \mathbb{E}_{q(\mathbf{z}_L | \mathbf{x}_L)} \left[\|\mathbf{x}_L - D_L(\mathbf{z}_L)\|_2^2 \right],$$

$$\mathcal{L}_{\text{KL}}^L = D_{\text{KL}}(q(\mathbf{z}_L | \mathbf{x}_L) \| p(\mathbf{z}_L)),$$

95 with $q(\mathbf{z}_L | \mathbf{x}_L)$ being the encoder’s approximate posterior and $p(\mathbf{z}_L) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, the standard
 96 gaussian prior. The vision foundation loss \mathcal{L}_{VF}^L aligns the latent codes with the feature space of a
 97 pretrained foundation model (e.g., DINOv2 [19]). The adversarial loss \mathcal{L}_{GAN}^L introduces a discrimina-
 98 tor to distinguish real from reconstructed low-frequency inputs. Finally, the perceptual loss \mathcal{L}_{LPIPS}^L
 99 encourages perceptual similarity between input and reconstruction using features from a pretrained
 100 Inception network [16]. All these loss components are commonly used in modern VAE frameworks.

101 **High-Frequency Objective.** In contrast, high-frequency components are trained without super-
 102 vision from pretrained models, as they tend to be biased toward low-frequency content. We use
 103 a lightweight VAE objective focused on reconstructing fine-scale details, along with adversarial
 104 regularization:

$$\mathcal{L}_{\text{high}} = \mathcal{L}_{\text{rec}}^H + \beta \cdot \mathcal{L}_{\text{KL}}^H + \mathcal{L}_{\text{GAN}}^H,$$

$$\mathcal{L}_{\text{rec}}^H = \mathbb{E}_{q(\mathbf{z}_H | \mathbf{x}_H)} \left[\|\mathbf{x}_H - D_H(\mathbf{z}_H)\|_1 \right],$$

$$\mathcal{L}_{\text{KL}}^H = D_{\text{KL}}(q(\mathbf{z}_H | \mathbf{x}_H) \| p(\mathbf{z}_H))$$

105 **At inference,** both decoders reconstruct the respective frequency bands, and the final image is syn-
 106 thesized via inverse wavelet transform. This frequency-aware latent tokenization via FA-VAE yields
 107 more expressive embeddings by preserving details across the full spectrum of spatial frequencies.

$$\hat{\mathbf{x}}_L = D_L(\mathbf{z}_L), \quad \hat{\mathbf{x}}_H = D_H(\mathbf{z}_H), \quad \hat{\mathbf{x}} = \mathcal{W}^{-1}(\hat{\mathbf{x}}_L, \hat{\mathbf{x}}_H).$$

Model (Tokenizer)	Tokenizer Config	Recon. Loss	Low Freq. Loss	High Freq. Loss	LPIPS	rFID
DC-AE	f32c32	0.0194	0.0484	0.0097	0.1580	0.7450
DC-AE	f64c128	0.0207	0.0527	0.0100	0.1667	0.7623
DC-AE	f128c512	0.0225	0.0581	0.0106	0.1805	0.7912
SD-VAE	f16c16	0.0180	0.0422	0.0100	0.1743	0.6213
KL-VAE	f16c16	0.0148	0.0326	0.0089	0.1355	0.5318
MS-VQ-VAE	f16c32v4096	0.0195	0.0549	0.0076	0.1890	0.6981
RQ-VAE	f32c256v16384	0.0416	0.1219	0.0149	0.2712	0.9095
TiTok-VQ-VAE	f256x2c64v8192	0.0226	0.0561	0.0114	0.2082	0.7550
TiTok-VQ-VAE	f256x4c64v8192	0.0339	0.0947	0.0137	0.2657	0.8823
TiTok-VQ-VAE	f256x8c64v8192	0.0450	0.1344	0.0152	0.2949	0.9416
TiTok-VAE	f256x1c16	0.0332	0.0923	0.0134	0.2232	0.8640
TiTok-VAE	f256x2c16	0.0461	0.1380	0.0155	0.2682	0.9187
TiTok-VAE	f256x4c16	0.0617	0.1961	0.0170	0.3182	0.9761
VQ-VAE	f16c256v1024	0.0492	0.1478	0.0163	0.3064	0.9102
VQ-VAE	f16c256v16384	0.0438	0.1262	0.0164	0.2784	0.8826
VA-VAE	f16c32	0.0105	0.0200	0.0074	0.0975	0.4884
FA-VAE (Ours)	f16c32	0.0044	0.0114	0.0020	0.0940	0.4156

Table 1: Quantitative comparison of latent tokenizers. Configurations: **f** = latent spatial resolution, **c** = latent dimensionality, **v** = vocabulary size (for quantized models). Lower is better.

108 3 Experiments

109 To evaluate the reconstruction performance of latent tokenizers in both the spatial and spectral
110 (frequency) domains, we consider a range of widely adopted visual tokenizers commonly used in
111 modern latent generative models. These tokenizers are based on different variants of autoencoding
112 architectures, each aiming to learn compact and informative latent embedding representations of the
113 input data distribution. For a comprehensive analysis, we include recent representative tokenizers
114 based on standard Autoencoders [20], Variational Autoencoders (VAEs) [21, 6, 22, 23, 7], and Vector
115 Quantized Autoencoders (VQ-VAEs) [23, 18, 11]. We also include our proposed FA-VAE, which
116 explicitly incorporates frequency-awareness into the latent embedding learning process. This selection
117 enables a fair comparison across diverse tokenization strategies with varying latent dimensionalities
118 (c) and latent feature resolutions (f) as shown in Table 1 with diverse tokenization configurations. We
119 use the metrics discussed in method section 2.1 for our evaluation.

120 Our analysis in Table 1 shows that KL-regularized VAEs with carefully designed parameterizations
121 consistently outperform both standard VAEs and VQ-VAE baselines across all reconstruction metrics,
122 achieving efficient performance at higher latent compression rates due to the absence of quantization
123 artifacts. Their effectiveness is strongly tied to latent parameterization quality, as seen in recent
124 autoregressive models [24] and VA-VAE [7], which further align latent spaces with foundation
125 models (e.g., DINOv2 [19]) for improved perceptual quality. However, such models jointly optimize
126 frequency components, leading to trade-offs and suboptimal preservation of low and high-frequency
127 details (Especially high frequency details). In contrast, our FA-VAE explicitly decouples low- and
128 high-frequency bands, enabling specialized representation learning and more precise reconstructions
129 of both global structures and fine details. FA-VAE nearly halves the reconstruction loss of the
130 strongest baseline while achieving the best overall performance across both low and high frequency
131 reconstructions, demonstrating the benefits of frequency-aware modeling. The consistent gap between
132 low- and high-frequency reconstructions across models further underscores the general difficulty of
133 capturing fine details effectively. The other metrics like LPIPS and rFID showcases the perceptual
134 quality of reconstructions of FA-VAE compared with other baselines.

135 4 Conclusion

136 In this work, we investigate frequency awareness in the latent embeddings of latent tokenizer models.
137 We find that jointly optimizing low- and high-frequency components leads to a frequency bias favor-
138 ing low-frequencies thereby degrading the reconstruction of fine-grained, high-frequency details and
139 overall perceptual quality. To address this limitation, we introduce **FA-VAE**, a frequency-aware VAE
140 framework that explicitly decouples low and high-frequency components using wavelet decompo-
141 sition, processes them independently, and fuses them into a unified latent representation. FA-VAE
142 achieves state-of-the-art performance across frequency-aware reconstruction metrics, demonstrating
143 improved fidelity and perceptual quality of the learned latent embeddings.

References

- 144
- 145 [1] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and
146 Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proc.*
147 *CVPR 2023*, 2023.
- 148 [2] Yuan Dong, Qi Zuo, Xiaodong Gu, Weihao Yuan, Zhengyi Zhao, Zilong Dong, Liefeng Bo, and Qixing
149 Huang. Gpld3d: Latent diffusion of 3d shape generative models by enforcing geometric and physical
150 priors. In *Proc. CVPR 2024*, page 56–66, 2024.
- 151 [3] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing
152 high-resolution images with few-step inference. In *Proc. ICLR 2024*, 2024.
- 153 [4] Cai Zhou, Xiyuan Wang, and Muhan Zhang. Unifying generation and prediction on graphs with latent
154 graph diffusion. In *Proc. NeurIPS 2024*, 2024.
- 155 [5] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*,
156 2013.
- 157 [6] Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image generation
158 without vector quantization. *arXiv preprint arXiv:2406.11838*, 2024.
- 159 [7] Jingfeng Yao, Bin Yang, and Xinggang Wang. Reconstruction vs. generation: Taming optimization dilemma
160 in latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
161 *Recognition*, 2025.
- 162 [8] Hao Phung, Quan Dao, and Anh Tran. Wavelet diffusion models are fast and scalable image generators. In
163 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
164 10199–10208, June 2023.
- 165 [9] Yueying Li, Hanbin Zhao, Jiaqing Zhou, Guozhi Xu, Tianlei Hu, Gang Chen, and Haobo Wang. Fedrs:
166 Frequency-aware enhancement for diffusion-based image super-resolution. In *ICLR 2025 (withdrawn)*,
167 2024. Amplitude and high-frequency enhancement modules for diffusion SR.
- 168 [10] Florentin Guth, Simon Coste, Valentin De Bortoli, and Stéphane Mallat. Wavelet score-based generative
169 modeling. In *NeurIPS*, 2022. Multi-scale diffusion over wavelet coefficients, linear time complexity.
- 170 [11] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling:
171 Scalable image generation via next-scale prediction. In *Advances in Neural Information Processing*
172 *Systems*, volume 37, pages 84839–84865, 2024.
- 173 [12] Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li,
174 Ligeng Zhu, Yao Lu, and Song Han. Sana: Efficient high-resolution image synthesis with linear diffusion
175 transformer, 2024.
- 176 [13] Enze Xie, Junsong Chen, Yuyang Zhao, Jincheng Yu, Ligeng Zhu, Yujun Lin, Zhekai Zhang, Muyang Li,
177 Junyu Chen, Han Cai, et al. Sana 1.5: Efficient scaling of training-time and inference-time compute in
178 linear diffusion transformer, 2025.
- 179 [14] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable
180 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer*
181 *Vision and Pattern Recognition (CVPR)*, pages 586–595, 2018.
- 182 [15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans
183 trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural*
184 *Information Processing Systems (NeurIPS)*, pages 6626–6637, 2017.
- 185 [16] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the
186 inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision*
187 *and Pattern Recognition (CVPR)*, pages 2818–2826, 2016.
- 188 [17] Colm Mulcahy. Image compression using the haar wavelet transform. *Spelman Science and Mathematics*
189 *Journal*, 1(1):22–31, 1997.
- 190 [18] Patrick et al. Esser et al. Taming transformers for high-resolution image synthesis. In *Proceedings of the*
191 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12873–12883, June
192 2021.

- 193 [19] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre
194 Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual
195 features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 196 [20] Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang, Haotian Tang, Muyang Li, Yao Lu, and
197 Song Han. Deep compression autoencoder for efficient high-resolution diffusion models. *arXiv preprint*
198 *arXiv:2410.10733*, 2024.
- 199 [21] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution
200 image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer*
201 *Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June 2022.
- 202 [22] Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. Autoregressive image
203 generation using residual quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
204 *and Pattern Recognition (CVPR)*, pages 11523–11532, June 2022.
- 205 [23] Qihang Yu, Mark Weber, Xueqing Deng, Xiaohui Shen, Daniel Cremers, and Liang-Chieh Chen. An
206 image is worth 32 tokens for reconstruction and generation. In *Advances in Neural Information Processing*
207 *Systems*, volume 37, pages 128940–128966, 2024.
- 208 [24] Hao Chen, Yujin Han, Fangyi Chen, Xiang Li, Yidong Wang, Jindong Wang, Ze Wang, Zicheng Liu, Difan
209 Zou, and Bhiksha Raj. Masked autoencoders are effective tokenizers for diffusion models. In *Forty-second*
210 *International Conference on Machine Learning*, 2025.

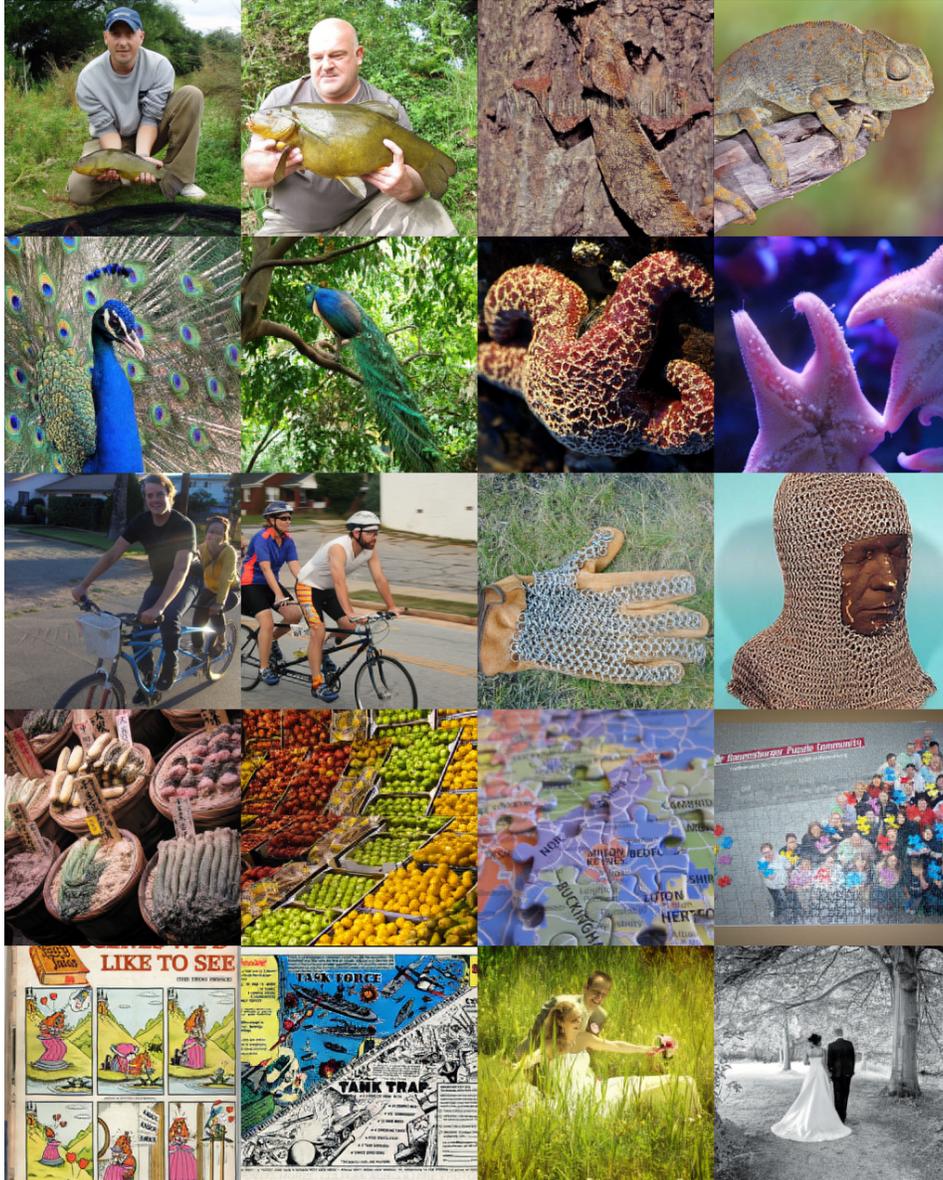


Figure 3: Qualitative reconstructions using FA-VAE on ImageNet 256×256 .

211 A Architecture and Implementation Details

212 The architecture of our Frequency-Aware Variational Au- toencoder (FA-VAE) consists of two inde-
 213 dependent encoder- decoder pairs, designed to separately process low- and high-frequency components
 214 of the input data. The encoder- decoder pair architecture is inspired by VAVAE [7]. We further follow
 215 a similar hyperparameter setup to [21, 7] for implementing our latent reconstruction module. To
 216 support distributed training across multiple nodes, we scale the learning rate and global batch size
 217 to 1×10^{-4} and 256, respectively, following the MAR setup [6]. We experiment with two $f16$
 218 tokenizers for low and high frequency subbands: one trained without visual alignment VF loss for
 219 high frequencies and one with VF loss using DINOv2 [19] for low frequencies. Here, f represents
 220 the latent spatial resolution factor and d the latent dimensionality. Following [7], we set the VF
 221 loss hyperparameters to $m_1 = 0.5$, $m_2 = 0.25$, and $w_{\text{hyper}} = 0.1$. Figure 3 shows the qualitative
 222 reconstructions of FA-VAE.

223 **NeurIPS Paper Checklist**

224 **1. Claims**

225 Question: Do the main claims made in the abstract and introduction accurately reflect the
226 paper’s contributions and scope?

227 Answer: [Yes]

228 Justification: See Sec.1.

229 **2. Limitations**

230 Question: Does the paper discuss the limitations of the work performed by the authors?

231 Answer: [NA]

232 **3. Theory assumptions and proofs**

233 Question: For each theoretical result, does the paper provide the full set of assumptions and
234 a complete (and correct) proof?

235 Answer: [NA]

236 **4. Experimental result reproducibility**

237 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
238 perimental results of the paper to the extent that it affects the main claims and/or conclusions
239 of the paper (regardless of whether the code and data are provided or not)?

240 Answer: [Yes]

241 Justification: See Appendix A.

242 **5. Open access to data and code**

243 Question: Does the paper provide open access to the data and code, with sufficient instruc-
244 tions to faithfully reproduce the main experimental results, as described in supplemental
245 material?

246 Answer: [NA]

247 **6. Experimental setting/details**

248 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
249 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
250 results?

251 Answer: [NA]

252 **7. Experiment statistical significance**

253 Question: Does the paper report error bars suitably and correctly defined or other appropriate
254 information about the statistical significance of the experiments?

255 Answer: [NA]

256 **8. Experiments compute resources**

257 Question: For each experiment, does the paper provide sufficient information on the computer
258 resources (type of compute workers, memory, time of execution) needed to reproduce the
259 experiments?

260 Answer: [NA]

261 **9. Code of ethics**

262 Question: Does the research conducted in the paper conform, in every respect, with the
263 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

264 Answer: [Yes]

265 **10. Broader impacts**

266 Question: Does the paper discuss both potential positive societal impacts and negative
267 societal impacts of the work performed?

268 Answer: [Yes]

269 Justification: See Abstract.

- 270 **11. Safeguards**
- 271 Question: Does the paper describe safeguards that have been put in place for responsible
272 release of data or models that have a high risk for misuse (e.g., pretrained language models,
273 image generators, or scraped datasets)?
- 274 Answer: [NA]
- 275 **12. Licenses for existing assets**
- 276 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
277 the paper, properly credited and are the license and terms of use explicitly mentioned and
278 properly respected?
- 279 Answer: [Yes]
- 280 **13. New assets**
- 281 Question: Are new assets introduced in the paper well documented and is the documentation
282 provided alongside the assets?
- 283 Answer: [NA]
- 284 **14. Crowdsourcing and research with human subjects**
- 285 Question: For crowdsourcing experiments and research with human subjects, does the paper
286 include the full text of instructions given to participants and screenshots, if applicable, as
287 well as details about compensation (if any)?
- 288 Answer: [NA]
- 289 **15. Institutional review board (IRB) approvals or equivalent for research with human
290 subjects**
- 291 Question: Does the paper describe potential risks incurred by study participants, whether
292 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
293 approvals (or an equivalent approval/review based on the requirements of your country or
294 institution) were obtained?
- 295 Answer: [NA]
- 296 **16. Declaration of LLM usage**
- 297 Question: Does the paper describe the usage of LLMs if it is an important, original, or
298 non-standard component of the core methods in this research? Note that if the LLM is used
299 only for writing, editing, or formatting purposes and does not impact the core methodology,
300 scientific rigor, or originality of the research, declaration is not required.
- 301 Answer: [No]