

---

## A COMPUTATION TIMES FOR LOCAL DESCRIPTORS

The computation times for the local scaling and local rank computation (since both require one randomized SVD computation for one latent vector) ends up being 3929s for 1000 samples. For local complexity we require 113s for 1000 samples. All the estimates are for a JAX implementation of Stable Diffusion on TPUv3.

Note that to train a reward model, we require the descriptors to be computed only once for each pre-trained model. If we compute the local scaling for 100k samples we require 173.1 TPU v3 hours which is equivalent to 54.58 V100 hours (according to Appendix A.3 Dhariwal & Nichol (2021)). Compared to 79,000 A100 hours required for Stable Diffusion training<sup>2</sup>, 24000 hours with enterprise level optimization<sup>3</sup>, the computation required for the descriptors and reward model training is significantly small. The computation time for the local descriptors can be further reduced by using a smaller  $k$  for our projection matrix  $W$ , or by using non-jacobian based methods, e.g., estimating the local scaling by measuring the change of volume for a unit norm  $\ell_1$ -ball in the input space. We leave exploration of these directions for future work.

## B RELATED WORKS

**Local geometry pre-diffusion.** Early applications of the local geometry of generative models involved improving the generation performance and/or utility of generative models via geometry inspired methods. For example, in Rifai et al. (2011) the authors proposed regularizing the contraction of the local geometry to learn better representations in autoencoders trained on MNIST and CIFAR10. The regularization penalty is employed via the norm of the input-output jacobian in Rifai et al. (2011), is an upper bound for local scaling presented in our paper. In Arvanitidis et al. (2017) the authors provided visualizations on the curvature of pre-trained VAE latent spaces and proposed using an auxiliary variance estimator neural network to regularize the latent space geometry during generation. In Kuhnelt et al. (2018) the authors perform latent space statistical inference problems, e.g., maximum likelihood inference, by training a separate neural network to approximate the Riemannian metric. In Humayun et al. (2022a) the authors proposed a novel latent space sampling distribution based on the latent space geometry that allows uniformly sampling the learned data manifold of continuous-piecewise affine generators. The authors showed downstream benefits with fairness and diversity for such latent space samplers. While most of these methods discuss pre-diffusion architectures, their results are early demonstrations of how the local geometry can affect downstream generation. also employ auxiliary Neural Networks to model an intrinsic property of a pre-trained generator, similar to how we propose using a reward model for Stable Diffusion.

**Local intrinsic dimensionality of diffusion models.** The local geometry of diffusion models and possible applications have garnered significant interest in recent years. In ? the authors propose a method to compute the intrinsic dimensionality of diffusion models using the assumption that the score field is perpendicular to the data manifold. For any vector  $x$  on the data manifold, the method requires computing the dimensionality of the score field around  $x$  and subtracting it from the ambient dimension. To do that, the authors perform one step of the forward diffusion process  $k$  times for  $x$ , denoise the  $k$  noisy samples using the diffusion model and compute the rank of the data matrix containing denoised samples to obtain the intrinsic dimensionality. Compared to this method, we compute the dimensionality directly via a random estimation of the input-output jacobian SVD. We do not require any assumption on the score function vector field being perpendicular to the data manifold, which may not hold for a diffusion model that is not optimally trained or highly complex training datasets like LAION.

In Kamkari et al. (2024a) the authors compute rank using the method proposed in ? and show that local intrinsic dimensionality can be used for out-of-distribution (OOD) detection. This is analogous to our analysis in Sec 3 on the local geometry on or off the manifold. We can see that the intuition authors provided in Kamkari et al. (2024a) for diffusion models trained on smaller models and datasets e.g., FMNIST, MNIST, translate to larger scale models like Stable Diffusion trained on LAION as we have presented fig. 3, fig. 16 and Sec 4. Especially in section 4, we show that creating OOD samples with corruptions on Imagenet data (in-distribution), we can have an increase or decrease

---

<sup>2</sup><https://www.mosaicml.com/blog/training-stable-diffusion-from-scratch-costs-160k>

<sup>3</sup><https://www.databricks.com/blog/stable-diffusion-2>

702 in negative-log likelihood (estimated via local scaling), with decrease for blurring corruptions and  
703 increase in noising corruptions.

704  
705 Concurrent work Kamkari et al. (2024b) has also shown the relationship between the intrinsic  
706 dimensionality (local rank) of Stable Diffusion scale models and the texture/visual complexity of  
707 generated images. We believe our analysis is much more holistic with three different geometric  
708 properties being measured compared to only local dimensionality. We i) show quantitatively how  
709 diversity measured via vendi score is higher for higher local scaling and rank values (section 4). We  
710 have explored how rank and scaling evolves continuously across the latent space in fig. 3. We have  
711 presented how the geometry distribution varies as we continually perturb images via noise or blurring  
712 operations section 4 And finally in Sec 5 we have presented a method to guide generation using the  
713 local geometry to obtain downstream generation benefits.

714 **Misc.** Apart from the aforementioned works, Kadkhodaie et al. (2023) show that the emergence of  
715 generalization in diffusion models – when two networks separately trained on the same data learn  
716 the same mapping – can be attributed to the eigenspectrum and eigenvectors of the input-output  
717 jacobian. While we do not study the training dynamics of the local geometric descriptors in our  
718 paper, Kadkhodaie et al. (2023) suggests that the local geometry can be an important indicator of  
719 memorization and generalization emergence in diffusion models. In Manor & Michaeli (2023) the  
720 authors use the posterior principal components of a denoiser for uncertainty quantification. This  
721 work suggests that components with larger eigenvalues result in larger uncertainty which is directly  
722 related to the local scaling descriptors as it measures the product of non-zero singular values. While  
723 in Manor & Michaeli (2023) the authors propose using it for only a single image denoiser, we show  
724 that it generalizes for any diffusion model including Stable Diffusion scale text-to-image diffusion  
725 models.

## 726 C CORRELATIONS BETWEEN THE THREE DESCRIPTORS

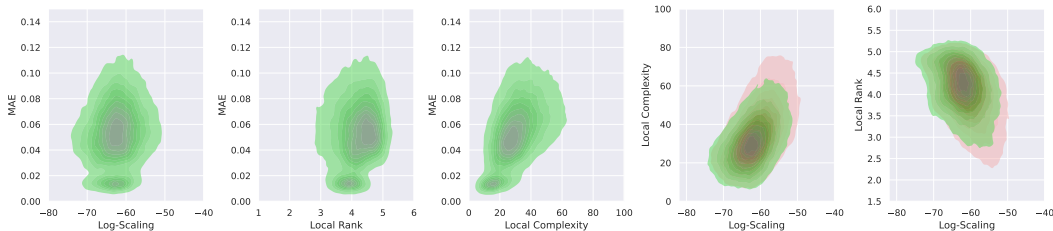
727 *Local scaling* characterizes the change of volume by the affine slope  $\mathbf{A}_\omega$  going from the latent space  
728 to the data manifold. *Local rank* characterizes the number of dimensions retained on the manifold  
729 after the network locally scales the latent space. Both local rank and scaling quantify first order  
730 properties of the CPWL operator. *Local complexity* approximates the ‘number of unique affine  
731 maps’ within a given neighborhood Humayun et al. (2024) by computing the number of CPWL  
732 knots intersecting an  $\ell_1$  ball in the input/latent space. Therefore local complexity is a measure of  
733 ‘un-smoothness’ and quantifies local second-order properties of a CPWL operator.

734 **Correlations between local scaling  $\psi$  and local rank  $\nu$ .** By definition, local scaling and local rank  
735 are correlated, since both characterize the change of volume by the network input-output map at any  
736 input space linear region – also evident in eq. (2) and eq. (4). Local scaling is also upper bounded by  
737 local rank,  $\psi_\omega \leq \sigma_0^{\nu_\omega}$  where  $\sigma_0$  is the largest singular value of  $\mathbf{A}_\omega$ . The correlation is evident for our  
738 low dimensional DDPM setting presented in fig. 13, local rank and local scaling are highly correlated  
739 in their spatial distribution. There are indications suggesting that the correlations persists throughout  
740 training as can be seen in fig. 13 rightmost column top and bottom. However in fig. 3, we can see that  
741 in the high-dimensional Stable Diffusion latent space, local scaling and rank are correlated but local  
742 rank has sharper changes spatially compared to local scaling.

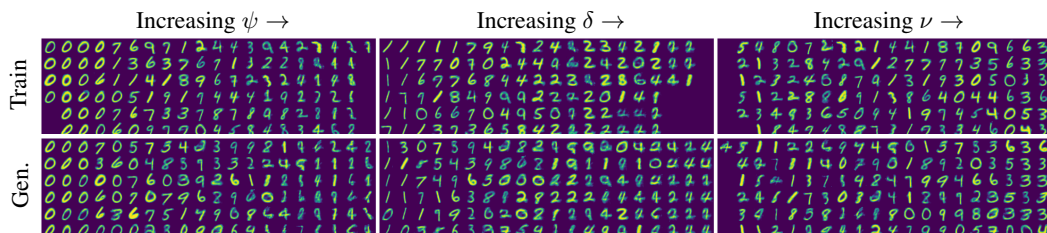
743 **Correlations between local complexity  $\delta$  and rank  $\nu$ .** There also exist correlations between  
744 local complexity and local rank due to the continuity of CPWL maps – between two neighboring  
745 linear regions  $\omega_1$  and  $\omega_2$ , the corresponding slope matrices  $\mathbf{A}_{\omega_1}$  and  $\mathbf{A}_{\omega_2}$  differ by at most one row.  
746 Therefore between two neighboring regions  $\omega_1$  and  $\omega_2$ ,  $|\nu_{\omega_1} - \nu_{\omega_2}| \leq 1$ . Informally, the local  
747 rank in a neighborhood  $V$  is lower bounded by the number of non-linearities in neighborhood  $V$ .  
748 This is evident in the empirical results presented in fig. 13 and fig. 3. In both figures, for input space  
749 neighborhoods with higher local complexity, we see a decrease in local rank. However, we do not  
750 observe sharp changes in local complexity as we observe in local rank in fig. 3. In fig. 13 we see that  
751 local rank is more discriminative of the data manifold compared to local complexity. Their training  
752 and denoising dynamics differ significantly as seen in fig. 13 rightmost column.

753 **Qualitative and quantitative results on correlations.** We train a beta-VAE unconditionally on  
754 MNIST and present in Fig. 9 samples from increasing local descriptor level sets from left to right  
755 along the columns. In Fig. 10, we present joint distributions of local scaling, complexity, rank and  
mean squared reconstruction error for training and test samples. We see that while local scaling,

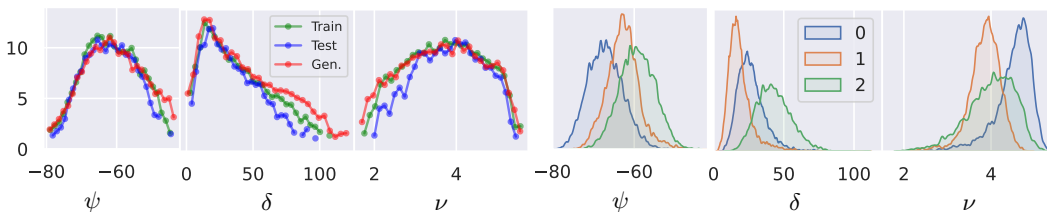
756 complexity and rank have some linear correlation, the classwise distribution in fig. 10 is very different  
 757 between the three. We also present in fig. 11 the vendi score for increasing local scaling level sets  
 758 and evidence that the population means for the descriptors don't follow the same pattern between  
 759 sub-populations.



769 Figure 9: Joint distributions for local scaling and MSE, local rank and MSE, local complexity and  
 770 MSE, local scaling and local complexity, and local scaling and local rank. We observe that local  
 771 complexity is linearly correlated with MSE, with higher complexity images incurring higher error.  
 772 Local scaling, rank and complexity have correlations between them as well.



782 Figure 10: Level sets of data manifold descriptors for a Beta-VAE trained unconditionally on MNIST.  
 783 From left to right, we present training samples (top row) and generated samples (bottom row)  
 784 for linearly increasing level sets of local scaling ( $\psi$ ) from  $[-80, -42]$ , local complexity ( $\delta$ ) from  
 785  $[0, 120]$  and local rank ( $\nu$ ) from  $[1.5, 5.5]$ . Not all level sets had an equal number of samples from  
 786 training/generated distributions. We see that for higher  $\psi$ , we have more outlier samples whereas  
 787 for lower  $\psi$  we have modal samples. For increasing  $\delta$  we see that the quality of generated samples  
 788 decreases and the diversity of samples is reduced as well. For higher  $\nu$  digits become more regularly  
 789 shaped.



798 Figure 11: (Left panel) Vendi score (Friedman & Dieng, 2023) calculated for samples from different  
 799 local descriptor level sets of a Beta-VAE. We take upto 150 samples from each level set and compute  
 800 vendi score separately for the MNIST train dataset, test dataset and generated samples. (Right  
 801 panel) Sub-population differences of local descriptors in training data. We see that the order of  
 802 sub-population means for the three classes, are not the same for all three descriptors.

## 804 D ADDITIONAL EXPERIMENTS

### 806 D.1 LOCAL SCALING FOR TRANSFORMER BASED DIFFUSION MODEL

807 Since we are based on the CPWL formulation of NNs, our framework would generalize to models of  
 808 any scale and any architecture with CPWL non-linearities. Empirically we have shown it to generalize  
 809 for non-CPWL architectures like Stable Diffusion v1.4 and DDPM that employs non CPWL non-  
 linearities such as attention, GeLU and much more. Following suggestions by the reviewer, we have

performed additional experiments with a DiT-XL Peebles & Xie (2023) trained on Imagenet-256. For the DiT we compute the descriptors for the transformer network, conditioned on noise level  $t = 0$ , i.e., zero noise level. We generate 5120 images conditioned on Imgewoof Howard (2019) classes and present in fig. 12, increasing local scaling level sets from left to right. We see that similar to fig. 16 from the, DiT exhibits a qualitative correlation between visual complexity and local scaling. For additional analysis we repeat the Stable Diffusion experiments on the relation between diversity and local scaling for DiT. We see that similar to Stable Diffusion, for increasing local scaling level sets, the diversity of images increase and then drop for the highest local scaling level sets.



Figure 12: Left: Vendi score and membership counts for increasing local scaling level sets, computed for a DiT transformer. We see that similar to Stable Diffusion, local scaling increases from lower to higher local scaling level sets, then drops for very high local scaling level sets. Right: Generated samples from each level set in the left panel. Sample sets from higher local scaling level sets, tend to be more diverse.

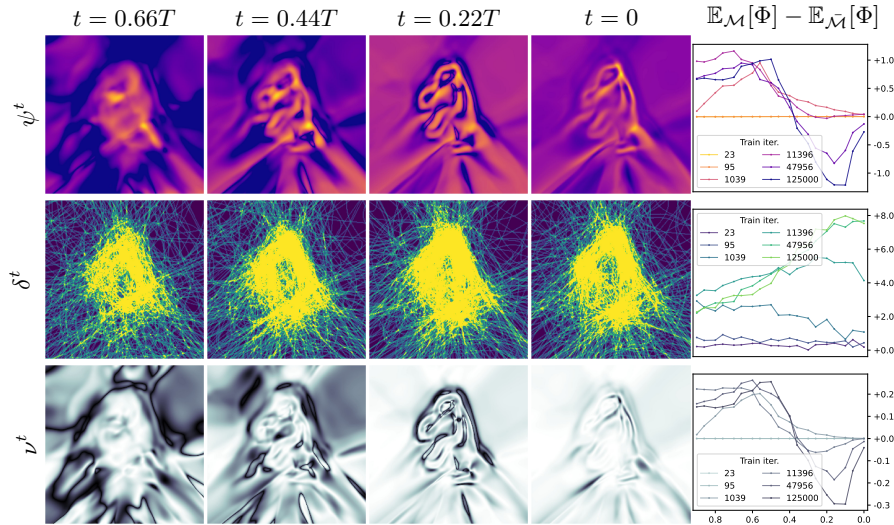


Figure 13: **Local geometric descriptors** computed over the input domain of a pre-trained toy diffusion model trained to produce samples from a dinosaur manifold  $\mathcal{M} \in \mathbb{R}^2$ . Descriptors are computed by conditioning the diffusion model on noise level  $t$ . We consider the set of input vectors within 0.05 units of the training data as on manifold  $\mathcal{M}$  and rest as off the manifold  $\bar{\mathcal{M}}$ . We present the difference between the expected descriptor values on and off the manifold,  $\mathbb{E}_{\mathcal{M}}[\Phi] - \mathbb{E}_{\bar{\mathcal{M}}}[\Phi], \forall \Phi \in \{\psi^t, \delta^t, \nu^t\}$  at different training iterations (right). We also present the descriptor computed over  $[-6, 6]^2$  for different noise levels  $t$  after 125000 training iterations (rest). We observe that  $\psi^t$  is lower,  $\delta^t$  is higher and  $\nu^t$  is lower on the manifold than off the target manifold for lower noise levels, especially after the model is considerably trained. This indicates that for well trained diffusion model, i.e., learned manifold  $\hat{\mathcal{M}} \approx \mathcal{M}$ , local descriptors can distinguish between on and off manifold vectors in the input space.

## D.2 VAE TRAINING DYNAMICS FOR MNIST

**Setup.** We train a Variational Auto Encoder (VAE) on the MNIST dataset with width 128 and depth 5 for both encoder and decoder. We add Gaussian noise with standard deviation  $\{0, 0.0001, 0.001, 0.01, 0.1\}$  to the training data. Initialization was not kept fixed. In Fig. 14, we

present plots showing the training dynamics of local complexity and scaling, averaged over all test dataset points from MNIST.

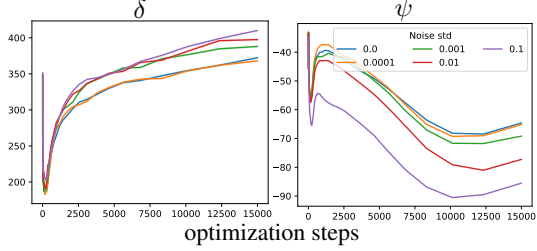


Figure 14: **Training dynamics of geometric descriptors for a VAE trained on MNIST with additive noise.** As training progresses local complexity  $\delta$  increases and local scaling  $\psi$  decreases suggesting an increase in expressivity and decrease in uncertainty on the data manifold. At latter time-steps,  $\psi \downarrow$  and  $\delta \uparrow$  if noise std. is increased.

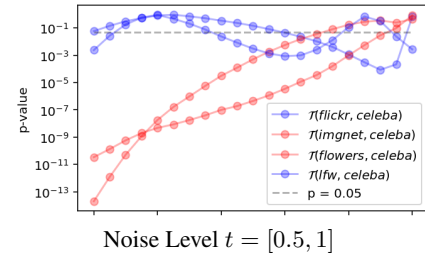


Figure 15: **Local scaling distribution difference between in-domain (blue) vs out-of-domain (red) datasets when conditioned on different noise levels** for an SD Unet trained on the CelebA HQ dataset. Here  $\mathcal{T}(a, b)$  denotes a t-test between local scaling distributions for dataset  $a$  and dataset  $b$ .

**Observations.** By increasing the noise we control the puffiness of the target manifold. We observe that as the noise standard deviation is increased there is 1) increase in  $\delta$  indicating the manifold becomes less smooth 2) decrease in local scaling indicating that the uncertainty decreases. We can also observe an initial dip in both local complexity and local scaling. This is similar to what was observed for discriminative models in (Humayun et al., 2024) where a double descent behavior was reported in the local complexity training dynamics of classification models. Based on these results, contrary to the observation in (Humayun et al., 2024), generative models do not have a double descent in local complexity however we do observe a double ascent in local scaling. *Our observations suggest that the training dynamics need to be taken into account, when comparing the local manifold geometry between two separately trained models.*

## E ENTROPY DIFFERENCE BETWEEN TWO

### NEARBY REGION

Suppose we have an injective  $\mathcal{G} : Z \rightarrow X$  mapping learned by a CPWL generator  $\mathcal{G}$ . Any linear region  $\omega$  in the latent space CPWL partition  $\Omega$  is mapped to a unique region on the output manifold. We define  $S$  as:

$$S = \mathcal{G}(z) \forall z \in \omega = \mathbf{A}_\omega z + \mathbf{b}_\omega \quad \forall z \in \omega$$

The change of volume from  $\omega \rightarrow S$  is  $\sqrt{\det(\mathbf{A}_\omega^T \mathbf{A}_\omega)}$ . Therefore for any latent  $z$  and output  $x = \mathcal{G}(z)$ :

$$p_{\mathcal{G}}(\mathbf{x}) = \sum_{z \in \omega} \frac{p_Z(z)}{\sqrt{\det(\mathbf{A}_\omega^T \mathbf{A}_\omega)}} \mathbb{1}_{z \in \omega}$$

For any  $\mathbf{z}_1 \in \omega_1$  the sum from the above equation can be ignored, since for all other regions the value would be zero.

Taking negative log and expectation on both sides the conditional entropy becomes

$$H(p_{\mathcal{G}}(\mathbf{x}_1); \mathbf{z} \in \omega_1) = H(p_Z(\mathbf{z}_1)) + \log(\sqrt{\det(\mathbf{A}_{\omega_1}^T \mathbf{A}_{\omega_1})})$$

For a uniform latent distribution and two regions  $\omega_1$  and  $\omega_2$ , substituting the second term above with  $\psi_{\omega_1}$

$$H(p_{\mathcal{G}}(\mathbf{x}_1); \mathbf{z}_1 \in \omega_1) - H(p_{\mathcal{G}}(\mathbf{x}_2); \mathbf{z}_2 \in \omega_2) = \psi_{\omega_1} - \psi_{\omega_2}$$

---

## F BROADER IMPACT STATEMENT

Our proposed framework for assessing and guiding generative models through manifold geometry offers several potential benefits to society. By providing a more objective and automated approach, we can significantly reduce the cost and time associated with human evaluation, making the auditing and mitigation of biases in large-scale models more accessible and efficient. This has implications for promoting fairness and equity in AI systems, particularly in domains where biases can have significant societal consequences.

Furthermore, our approach can empower researchers and practitioners to better understand the relationship between the geometry of learned representations and various aspects of model behavior, such as generation quality, diversity, and bias. This deeper understanding can inform the development of more robust and reliable generative models, leading to advancements in various fields, including art, design, healthcare, and education.

However, we recognize that our approach is not without limitations and potential risks. While it can be a valuable tool for identifying and mitigating biases, it should not and cannot fully replace human annotators, especially in high-risk domains where human judgment and contextual understanding are crucial. Our method focuses on reducing costs and improving the auditing process, but it should not be used as a standalone approach.

Moreover, the increased automation enabled by our approach raises concerns about the potential displacement of human annotators, leading to job losses and economic disruptions. While our method addresses some aspects of model evaluation, it is not comprehensive and cannot assess all facets of model behavior. Therefore, it should be used with caution and in conjunction with other evaluation methods, including human expertise.

## G EXTRA FIGURES



Figure 16: **Local Scaling is sensitive to natural image variations.** ImageNet images ordered along the columns (from left to right), with increasing local scaling  $\psi$  of the Stable Diffusion decoder learned manifold. We observe that ImageNet samples with lower values of  $\psi$  contain simpler backgrounds with modal representation of the object category. Conversely for higher  $\psi$  we have increasing diversity both in background and foreground features.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

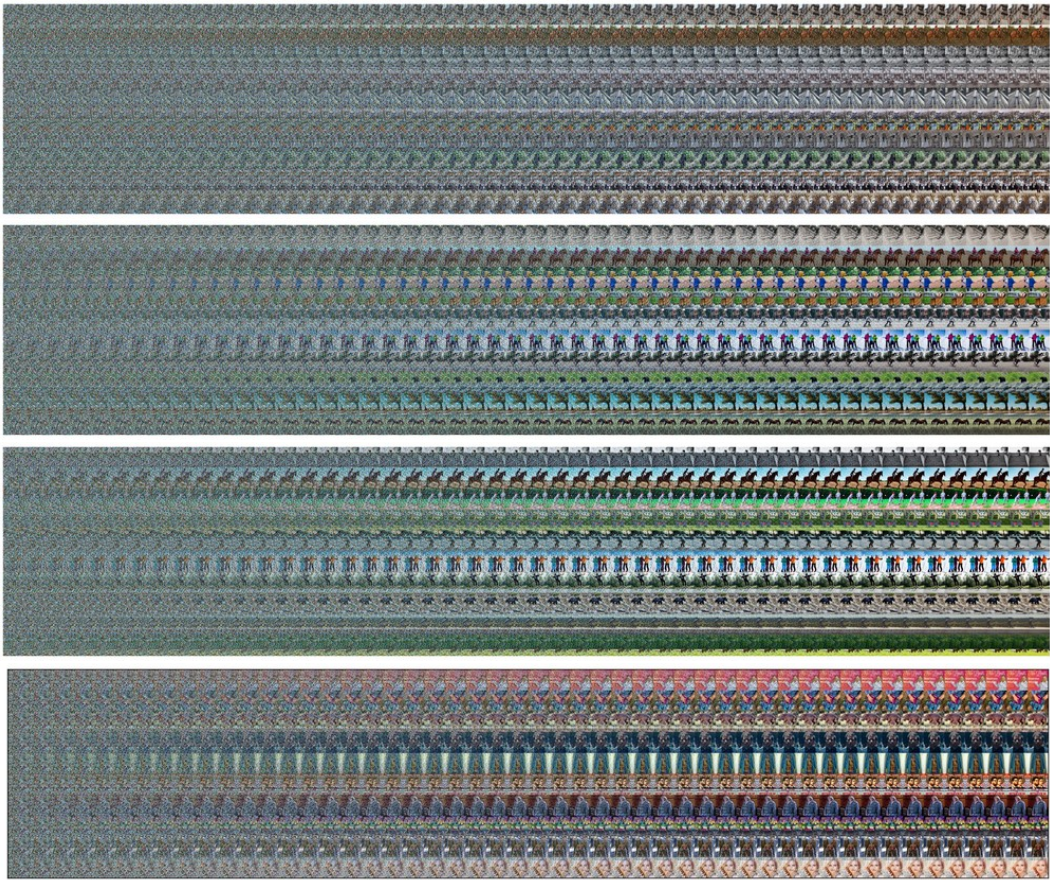


Figure 17: Images generated during 50 diffusion denoising steps for top to bottom, COCO prompts generated with guidance scale 1,5,9 and memorized prompts generated with guidance scale 7.5. Higher guidance scale images, as well as memorized images, tend to resolve faster during the denoising process.

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

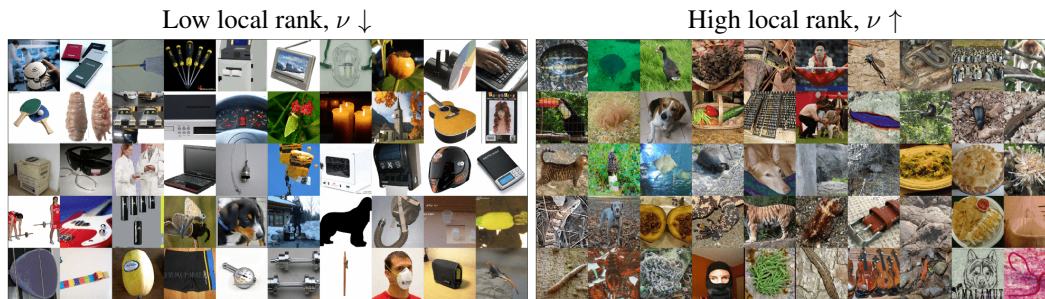


Figure 18: **Influence of the local rank descriptor value on overall image perception.** Images with the lowest (left) and highest (right) local rank  $\nu$  from a set of 20000 randomly sampled ImageNet dataset samples. Low rank images contain simpler textures for every class compared to the high rank samples. This is because for images with higher local rank, the learned manifold is higher dimensional therefore allowing higher independent degrees of variations locally for the generated images.



Figure 19: **Influence of the local scaling descriptor.** Imagenet images with high and low local scaling for the stable diffusion decoder. Each coordinate in both left and right image grids, correspond to the same imagenet class.



1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

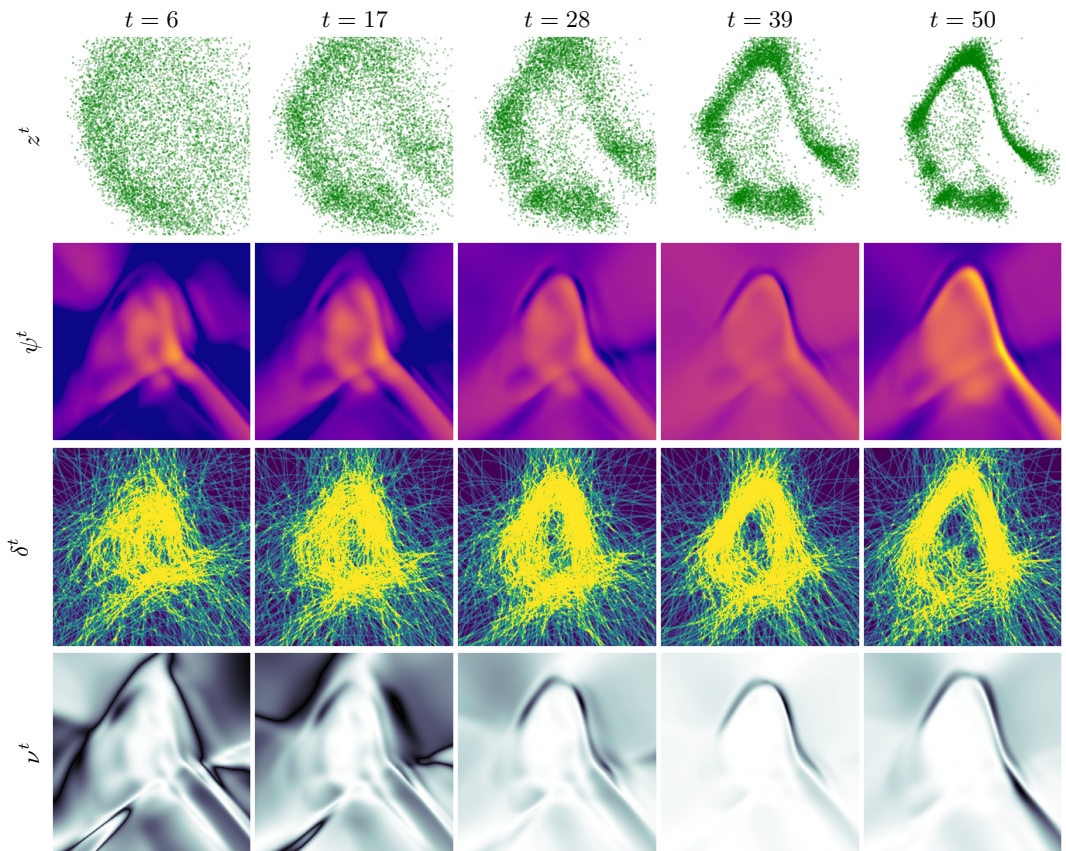


Figure 20: **2D Data Manifold Geometry, A toy Example.** After 11395 optimization steps. Geometry of a diffusion model input-output mapping, trained to on a toy 2D distribution. Local scaling lower around data manifold, local complexity higher around manifold, rank is lower around manifold as well.  $t=50$  has considerably low variance in local scaling showing that final timestep has a diminishing change of density.

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187



Figure 21: **Reward guidance on stable diffusion (maximizing the reward).** We observe a significant increase in both background detail and artifact diversity within the generated images.

a Australian terrier with  
 a city in the background  
 in the background

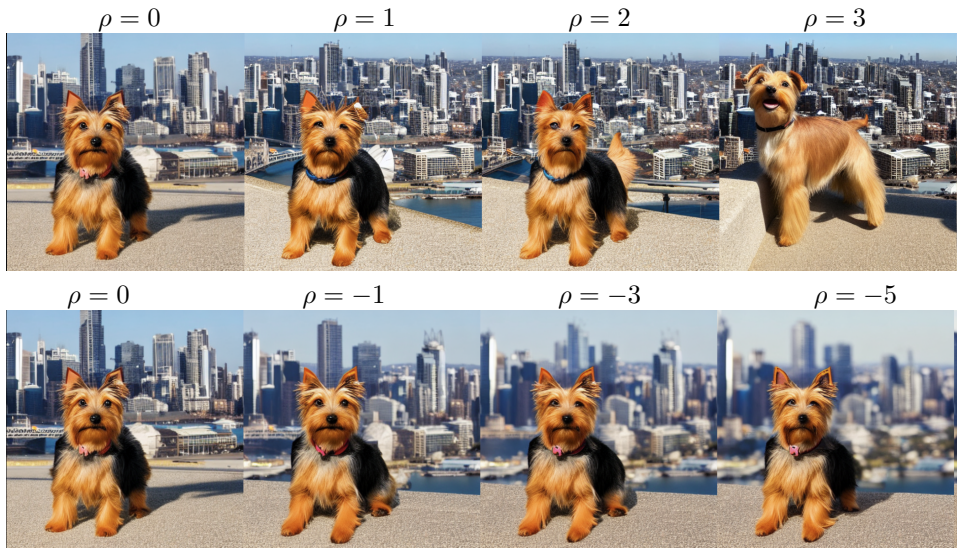


Figure 22: **Controlling image diversity with local scaling.** Using Reward guidance to increase (top row) and decrease diversity (bottom row) using same initial seed.

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

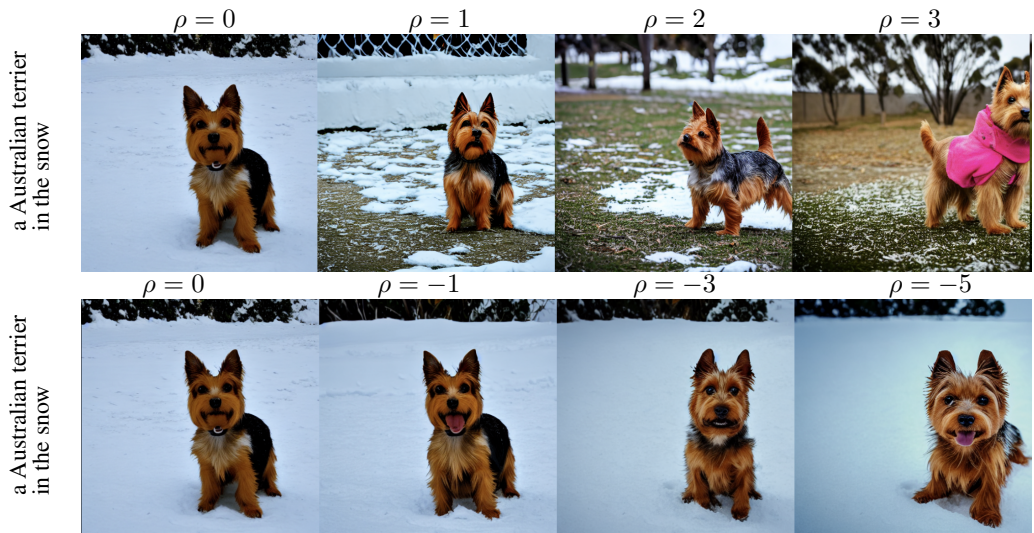


Figure 23: **Controlling image diversity with local scaling.** Using Reward guidance to increase (top row) and decrease diversity (bottom row) using same initial seed.

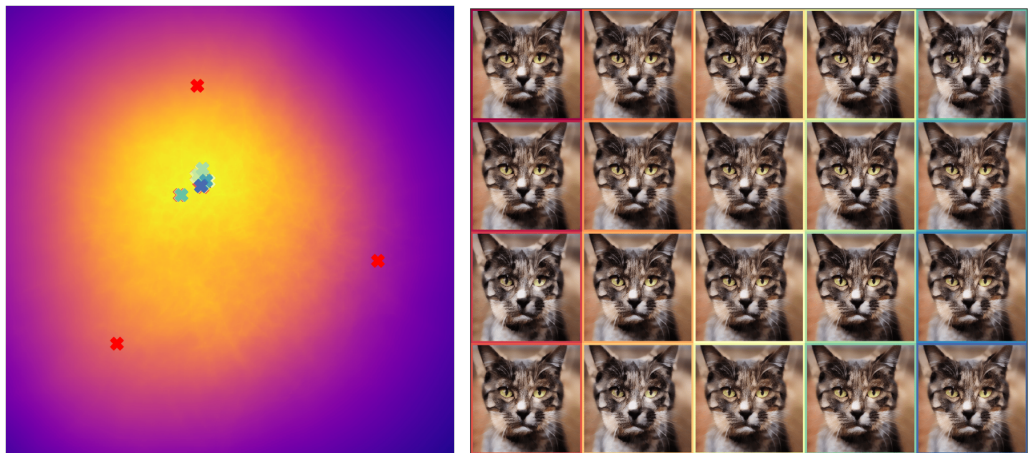


Figure 24: Decoded images (right) using 20 latents (left) from the 2D subspace, with highest  $\psi$ . Each image bounding box (right) is color coded according to the corresponding latent vector (left).

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

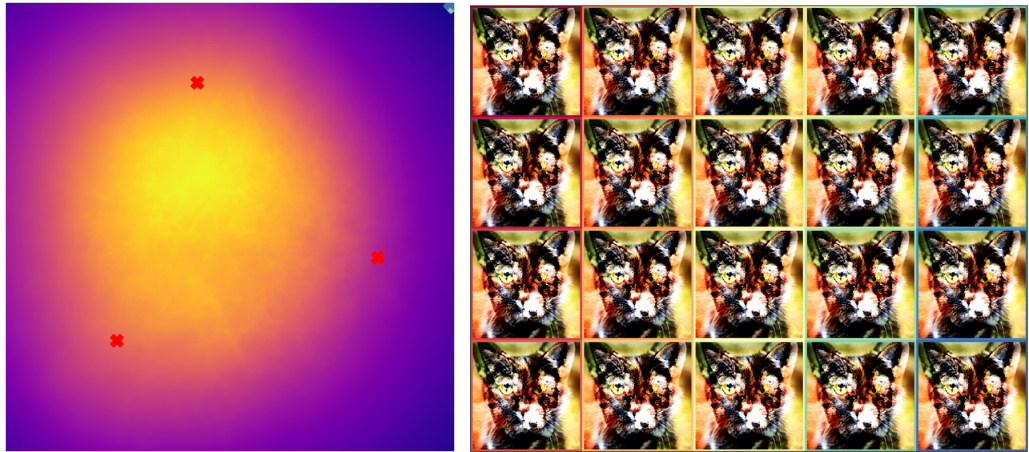


Figure 25: Decoded images (right) using 20 latents (left) from the 2D subspace, with lowest  $\psi$ . Each image bounding box (right) is color coded according to the corresponding latent vector (left). Selected latents lie outside the domain of the VQGAN latent space.

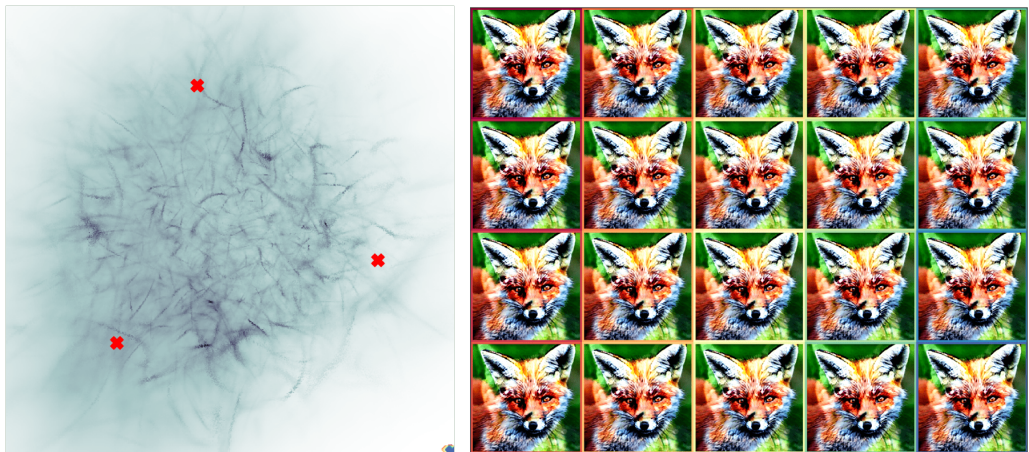


Figure 26: Decoded images (right) using 20 latents (left) from the 2D subspace, with highest  $\nu$ . Each image bounding box (right) is color coded according to the corresponding latent vector (left). Selected latents lie outside the domain of the VQGAN latent space.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310

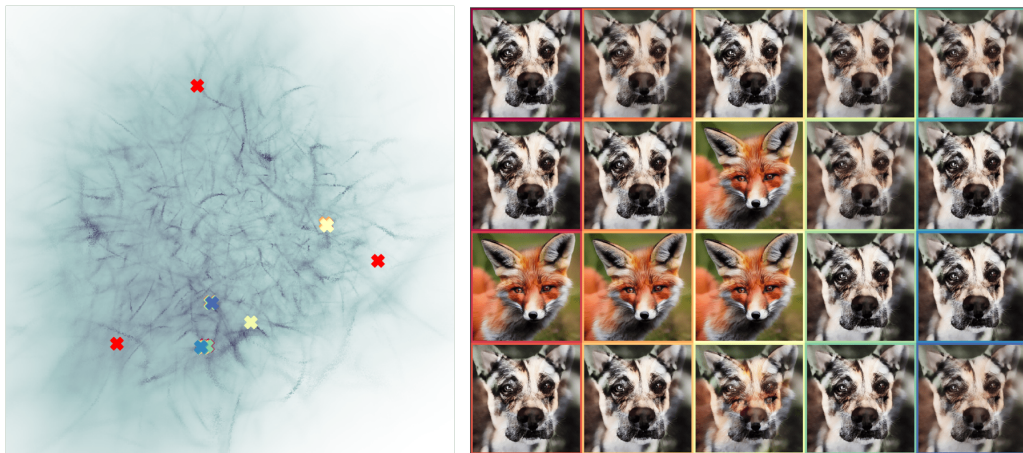


Figure 27: Decoded images (right) using 20 latents (left) from the 2D subspace, with lowest  $\nu$ . Each image bounding box (right) is color coded according to the corresponding latent vector (left).

1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329

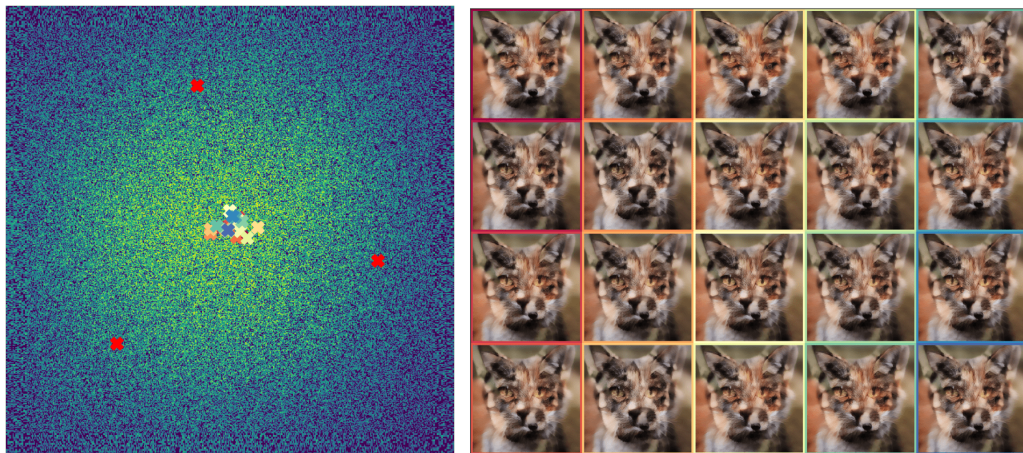


Figure 28: Decoded images (right) using 20 latents (left) from the 2D subspace, with highest  $\delta$ . Each image bounding box (right) is color coded according to the corresponding latent vector (left).

1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347

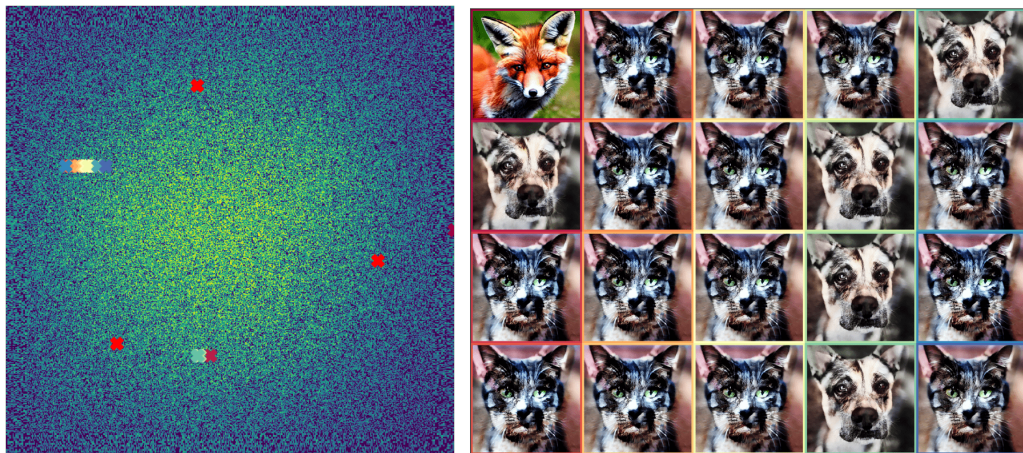


Figure 29: Decoded images (right) using 20 latents (left) from the 2D subspace, with lowest  $\delta$ . Each image bounding box (right) is color coded according to the corresponding latent vector (left).

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

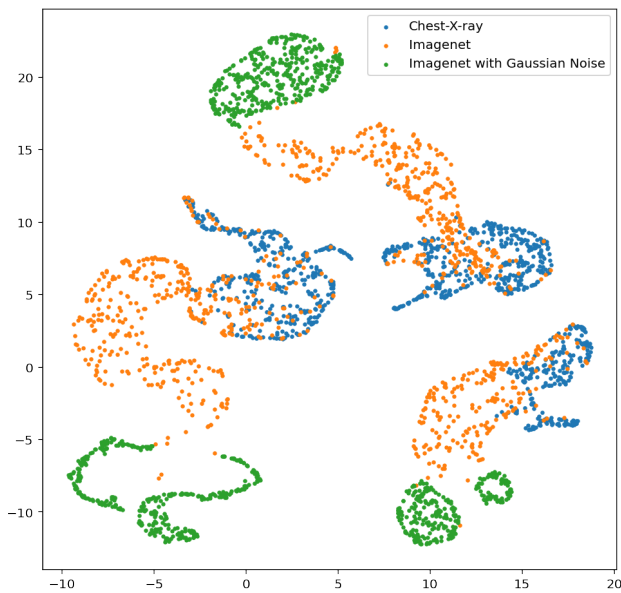


Figure 30: UMAP visualization of the aggregated local geometry descriptors (local smoothness, local rank, and local complexity). This reveals distinct, non-overlapping clusters, clearly separating the Imagenet, Imagenet Corrupted with Gaussian Noise, and Chest X-ray datasets. This visual evidence underscores the effectiveness of aggregating the descriptors to capture unique patterns within each dataset, demonstrating its ability to provide a meaningful and interpretable representation of the underlying data