

TRUSTWORTHY SR: RESOLVING AMBIGUITY IN IMAGE SUPER-RESOLUTION VIA DIFFUSION MODELS AND HUMAN FEEDBACK

Cansu Korkmaz¹, Ege Cirakman², A.Murat Tekalp¹, Zafer Doğan¹

¹ College of Engineering and KUIS AI Center, Koç University, Istanbul, Turkey

² Department of Electrical and Electronics Engineering, Istanbul Technical University, Istanbul, Turkey

ABSTRACT

Super-resolution (SR) is an ill-posed inverse problem with a large set of feasible solutions that are consistent with a given low-resolution image. Various deterministic algorithms aim to find a single solution that balances fidelity and perceptual quality; however, this trade-off often causes visual artifacts that bring ambiguity in information-centric applications. On the other hand, diffusion models (DMs) excel in generating a diverse set of feasible SR images that span the solution space. The challenge is then how to determine the most likely solution among this set in a trustworthy manner. We observe that quantitative measures, such as PSNR, LPIPS, DISTs, are not reliable indicators to resolve ambiguous cases. To this effect, we propose employing human feedback, where we ask human subjects to select a small number of likely samples and we ensemble the averages of selected samples. This strategy leverages the high-quality image generation capabilities of DMs, while recognizing the importance of obtaining a single trustworthy solution, especially in use cases, such as identification of specific digits or letters, where generating multiple feasible solutions may not lead to a reliable outcome. Experimental results demonstrate that our proposed strategy provides more trustworthy solutions when compared to state-of-the-art SR methods.

Index Terms— super-resolution, diffusion models, artifacts, trustworthy sample selection, human feedback

1. INTRODUCTION

Single-image super-resolution (SR) is an ill-posed inverse problem, where a single input image can correspond to multiple feasible output images, introducing ambiguity into the SR reconstruction process. Early deep learning based SR methods [9, 10, 11, 1, 12, 13, 4] treated SR as a deterministic regression problem and generated a single output image based on a set of LR-HR paired training data.

This work was supported in part by an AI Fellowship to C. Korkmaz provided by the KUIS AI Center. This work was supported in part by TUBITAK 2247-A Award No. 120C156, TUBITAK 2232 Award No. 118C337, and KUIS AI Center funded by Turkish Is Bank. AMT acknowledges support from Turkish Academy of Sciences (TUBA).



Fig. 1: Visual performance of recent $\times 4$ SR methods on a crop from Urban100 dataset (img-6) [8]. SOTA methods reconstruct “5” as “6”, whereas the opening in the lower part of “5” is visible in our result confirming that the proposed strategy resolves the ambiguity and provide a trustworthy solution. Note PSNR, DISTs and other quantitative scores are not reliable indicators to resolve such ambiguity.

Recognizing the inherently ill-posed nature of the SR problem, many recent methods [14, 15, 16, 17, 18, 19] and challenges [20, 21] propose learning a one-to-many mapping that is consistent with the conditional distribution of output images given the input, aiming to generate a diverse set of feasible SR images that span the SR space. These formulations prioritize photo-realism of solutions, consistency with the LR image, and diversity (coverage of the SR space). However, such stochastic generative models introduce a new problem: how to find a trustworthy solution when extracting critical information from many possibly conflicting SR images, e.g., whether a digit is 5 or 6. In such cases, there is need for reliable methodology to determine whether it is possible to rule out certain solutions and pick a trustworthy outcome.

We observe that the traditional objective quality (fidelity) measures such as Peak Signal-to-Noise Ratio (PSNR) and SSIM, as well as popular perceptual quality measures such as LPIPS [22] and DISTS [23], are inadequate to evaluate trustworthy SR solutions. Indeed there exists a divergence between human visual evaluations and known quantitative measures, particularly when evaluating generated high-frequency (HF) details. Therefore, one cannot rely solely on numerical measures to evaluate trustworthiness in SR.

To this effect, we propose a human feedback-centric approach to assess trustworthiness of SR solutions. Human subjects are asked to select up to 5 most likely outputs from a set of feasible SR images. The images selected by each subject are then ensembled according to the requirements of the task. We introduce a straightforward yet highly effective image ensembling strategy in this study, enabling diffusion models to leverage human feedback to resolve the ambiguity in generated samples. In essence, our goal is to address the challenge of finding a trustworthy solution with accurate details. Our experimental results demonstrate that a pre-trained Latent Diffusion Model (LDM) with strategic sample selection guided by human feedback (LDM-SS) and image ensembling outperforms state-of-the-art SR methods considering both reliability and visual image quality. However, we observe that the success of LDM-SS in improving trustworthiness and subjective visual quality does not necessarily translate into improvements in traditional quantitative measures.

Our main contributions are summarized as follows:

1. We address the challenge of obtaining a single trustworthy solution in the SR space spanned by diffusion models when information extraction is critical.
2. We introduce a human feedback-centric approach for SR sample selection, along with an ensembling strategy, demonstrating the superiority of the diffusion model in achieving accurate and reliable results.
3. Our suggested approach is versatile, as the combination of ensembling strategy and the integration of human feedback can be applied to any stochastic generative model, guiding it to produce trustworthy and consistent SR images.

2. RELATED WORKS

One-to-One SR Inference. Many popular SR models, including EDSR [1], RRDB [4] and PDASR [24], are deterministic regressive mappings from LR to HR images trained by l_1 or l_2 reconstruction losses. Even though these models achieve high fidelity, measured by PSNR, they still produce serious artifacts that contribute to the ambiguity problem.

Generative adversarial networks (GAN) have been proposed to generate photo-realistic images [25]. Many SR models based on the principles of GAN have been proposed over the years [11], [4], [2], [5], [7], [3] to generate a single SR image (per λ). It is well known that GANs hallucinate HF details. It is obvious to humans that some of these hallucinations

are artifacts, while some others may look like real although they are fake. Hence, GAN-based SR models cannot offer a trustworthy solution to resolve the ambiguity problem.

One-to-Many SR Inference. To generate rich diversity, likelihood-based model training prioritizing accurate density estimation such as variational autoencoders (VAE) [26, 27] and normalizing flow based SR methods [14, 15, 28] have been introduced. They offer notable benefits compared to GAN-based methods including monotonic convergence, stable training and efficiency while generating multiple SR images, however, they exhibit low fidelity in terms of image quality. Similarly, autoregressive models (ARM) [29, 30] excel in density estimation but face high computational complexity for inference due to their sequential sampling process. In addition, pixel-based image representations require prolonged training times to learn subtle HF details [31].

Recently, significant progress has been made in one-to-many SR image generation with the advent of diffusion models [32, 33, 34, 35, 36, 16, 17, 19]. However, current diffusion models still face some challenges, including but not limited to complex two-stage pipelines, high computational requirements for training, and presence of unnatural artifacts resulting in unreliable ambiguous SR outputs. In this work, we employ a pre-trained LDM for $\times 4$ SR to avoid lengthy training and propose strategic sample selection via human feedback to address the challenge of trustworthy SR image reconstruction by integrating the strengths of diffusion models, human feedback, and image ensembling within a practical framework.

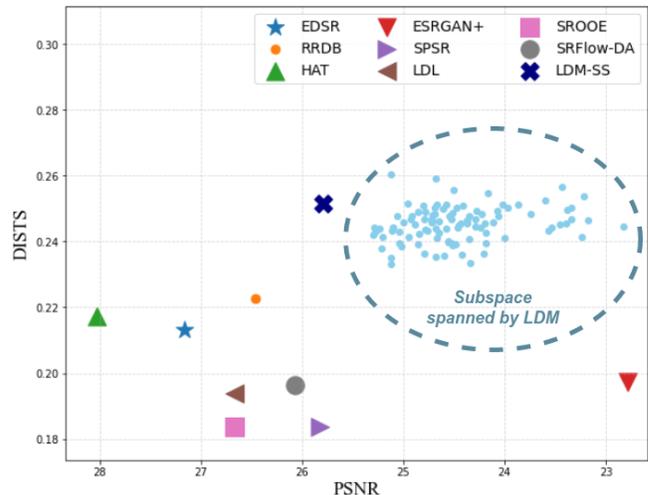


Fig. 2: Demonstration of the SR space spanned by LDM [17] samples, proposed LDM-SS and other state-of-the-art methods on the PSNR-DISTS plane. We note that perception-distortion tradeoff with respect to known metrics does not correlate well with visual quality and trustworthiness.

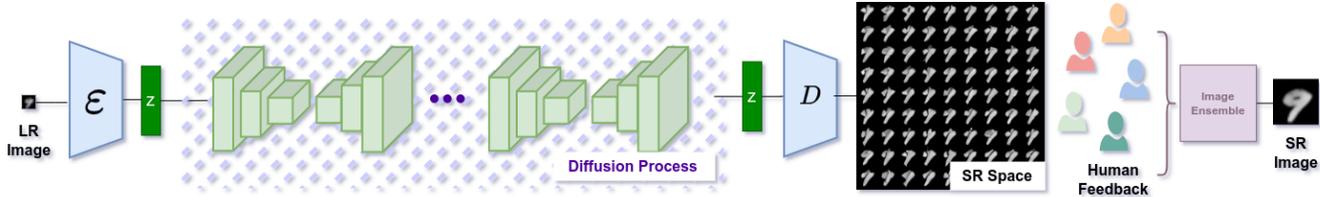


Fig. 3: Block diagram of our approach depicting sample selection from the diffusion model SR space by human feedback.

3. LDM-SS: RESOLVING AMBIGUITY THROUGH SAMPLE SELECTION IN DIFFUSION SR SPACE

3.1. SR Space Spanned by LDM

Diffusion models [32] are statistical models designed to learn a data distribution, $p(x)$, through progressively denoising a normal distributed variable. This process involves learning the inverse operation of a Markov Chain with a fixed length T . Recently proposed LDM [17] method performs diffusion process in a low dimensional latent space and provides computationally tractable and flexible SR images what we refer as subspace spanned by LDM as demonstrated in Fig. 2. On one hand all diffusion-based methods including LDM have a common problem: generation of SR image samples exhibiting rich textures but lack fidelity. On the other hand, in order to compare the performance of DMs and one-to-one SR approaches, we need to map the set of outputs to a single SR image. One can simply select a certain realization from the set. However, since there is a significant variety among all possible SR realizations produced by DMs, it is hard to obtain the well-suited result in the first realization. Therefore, generating numerous possible solutions in such information-centric applications may not result in a conclusive decision.

3.2. Resolving Ambiguity in the Diffusion SR Space

Our proposed method involves combining various diffusion samples into a unified, trustworthy image by ensembling them through human feedback. The objective is to strike a balance between high fidelity and perceptual distortion, particularly for information-centric applications. The block diagram illustrating our proposed approach is depicted in Fig. 3. Specifically, while developing the concept of strategic sample selection for trustworthy SR images, we employ LDM method to construct SR space by sampling from the learned distribution at inference time. For each LR image, we generate up to 100 images for human selection, then top 5 selected images by majority voting are ensembled by pixel-wise averaging to construct the SR image.

A hypothetical use case is demonstrated in the following example: Suppose an evidence, shown in Fig. 4-(a), is presented to a court of law. The prosecution confidently asserted that the digit in the image is unequivocally “5.” On the opposing side, the defense argued persuasively for an undeni-

Table 1: Thirty participants were asked to select two samples out of 324 generated samples that are most helpful to identify the specific digit.

	Perceived as “5”	Perceived as “6”
# of People	22 (73.3%)	8 (26.7%)

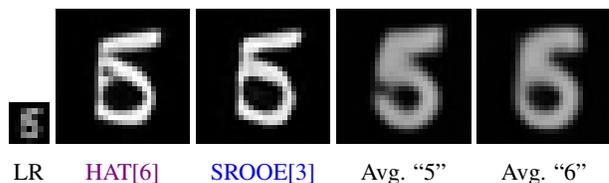


Fig. 4: Identification of the digit from LR image is impossible and results of state-of-the-art methods HAT [6] (Regressive) and SROOE [3] (GAN-SR) are ambiguous. However, the five most selected figures were combined through pixel-wise averaging, yielding single, informative SR image. The prevalence of the perception of “5” enables mitigating ambiguity.

able “6.” As experts in image processing were summoned to the witness stand, they face a tough challenge. They have employed state-of-the-art regressive and generative SR algorithms, including HAT [6] and SROOE [3]. Yet, they have not reached a consensus. This ambiguity in SR problems exemplifies the profound complexities awaiting a solution. Since the correct answer, the presented digit being 5 or 6 can be deduced from the plausible SR space spanned by diffusion samples, then direct human evaluation can be employed. Specifically, we use mean opinion score (MOS) and total of 30 participants are tasked with identifying the specific number depicted. Subsequently, they are asked to select the 2 samples among 324 generated ones that are most helpful to identify the digit. As presented in Table 1, 22 among 30 participants perceived generated samples as “5”, whereas 8 of them predicted the number as “6”. Then, 5 most selected figures for the digit is ensembled by pixel-wise averaging resulting in a single, informative SR image demonstrated in Fig. 4-(d) and (e). It is important to note that participants are not provided with information regarding the actual, ground-truth digit, relying solely on their visual preferences. Since, most people perceived the number as “5”, the ambiguity is resolved to sum extend.

EDSR [1]	RRDB [4]	HAT [6]	ESRGAN+ [2]	SPSR [5]	LDL [7]	SROOE [3]	HCFlow++ [37]	SRFlowDA [15]	LDM [17]	LDM-SS (Ours)	HR [38]
20.28	20.99	20.46	19.00	19.87	21.29	20.93	19.45	21.33	16.89	17.62	PSNR↑
0.146	0.102	0.125	0.108	0.125	0.095	0.106	0.121	0.099	0.212	0.215	DISTS↓[23]
EDSR [1]	RRDB [4]	HAT [6]	ESRGAN+ [2]	SPSR [5]	LDL [7]	SROOE [3]	HCFlow++ [37]	SRFlowDA [15]	LDM [17]	LDM-SS (Ours)	HR [38]
18.74	18.43	17.96	16.61	16.90	17.84	18.37	17.85	17.68	16.73	17.07	PSNR↑
0.136	0.146	0.136	0.122	0.145	0.118	0.125	0.204	0.151	0.190	0.151	DISTS↓[23]

Fig. 5: Visual comparison of proposed method and state-of-the-art regressive (purple), GAN-based (blue), flow-based (red) and LDM SR methods on Mnist dataset [38]. It can be seen that the proposed LDM-SS method provides more reliable SR images for information retrieval, but quantitative metrics are insufficient to capture the nuances of visual artifacts or trustworthiness.

4. EXPERIMENTAL RESULTS

Implementation Details and Benchmarks. We selected two popular datasets as benchmarks: Mnist dataset [38] and DIV2K [39]. Since, the size of Mnist images is 28x28, we obtain 7x7 LR images after downsampling by 4 in each dimension using Matlab’s bicubic kernel. Then, each LR image is upsampled by the latent diffusion model (LDM) [17]. Since LDM is designed to upsample images from 64x64 to 256x256, 7x7 LR sample is repeated 9 times both vertically and horizontally before applying the diffusion process. Since it is a stochastic process each SR image contains variety of upsampled numbers. Similarly, cropped 32x32 pixels of RGB images from DIV2K [39] are fed into the pretrained LDM model with a scaling factor of 4×.

Human Evaluation. Since none of the highly utilized evaluation metrics including PSNR, LR-Consistency [20], SSIM, LPIPS [22] and DISTS [23] do not correlate with the information that has been conveyed through the SR images. Accordingly, we employ direct human assessment for evaluation. While the mean opinion score (MOS) is a prevalent measure for assessing image quality has been found to be a more reliable method for such subjective quality assessments. In Task-1, the 30 participants are assigned the task of identifying a specific number from the Mnist [38] dataset, followed by the requirement to select the two numbers deemed most “natural” from a pool of 324 generated samples. Task-2 is similar, subjects are presented 15 natural images from DIV2K [39], focusing solely on selecting the more photo-realistic image. In both tasks, the ground-truth image is not provided, hence participants just rely entirely on their visual preferences.

4.1. Comparison with the State-of-the-Art Methods

Quantitative Comparison. Table 2 demonstrates quantitative comparison for 4× SR methods and our proposed strate-

gic sample selection approach for Task-2. The main objective of Task-2 is to have the subjects concentrate exclusively on choosing the image that appears more “photo-realistic” among 15 diffusion samples. In each round, fifteen 128x128 samples are presented simultaneously to each participant and asked to select at most 3 images with natural looking details, colors and lightning. To have a concrete evaluation, this selection process is repeated for 10 images from DIV2K [39] dataset. Also, in order not to biased the participants, the ground-truth image is kept hidden. The top 3 selected images are ensembled by pixel-wise averaging. Then, we compare the ensembled images with the existing state-of-the-art methods including EDSR [1], RRDB [4], HAT [6], ESRGAN+ [2], SPSR [5], LDL [7], SROOE [3] as well as stochastic SR methods like HCFlow++ [37], SRFlow-DA [15]. Even though we provide quantitative comparison results for our proposed approach, the efficacy of evaluating visual artifacts in SR tasks cannot solely rely on metrics such as PSNR or other quantitative perceptual scores. While these metrics provide numerical insights into image quality, they might not capture the nuances of visual artifacts effectively. In this context, it becomes crucial to explore alternative methods that go beyond quantitative assessments.

Qualitative Comparison. Visual comparisons among 4× SR approaches and LDM-SS are presented in Fig. 1, 5 and 6. These figures showcase the effectiveness of strategic sample selection by human feedback in mitigating visual artifacts. In details, we observe that all GAN-SR results including ESRGAN+ [2], SPSR [5], LDL [7], SROOE [3] as well as stochastic SR methods like HCFlow++ [37], SRFlow-DA [15] produce visible artifacts and experience excessive sharpness. For instance, the visual results presented in Fig. 5 demonstrate the effectiveness of human feedback for information-centric applications by enabling the outcome of the accurate digit. The

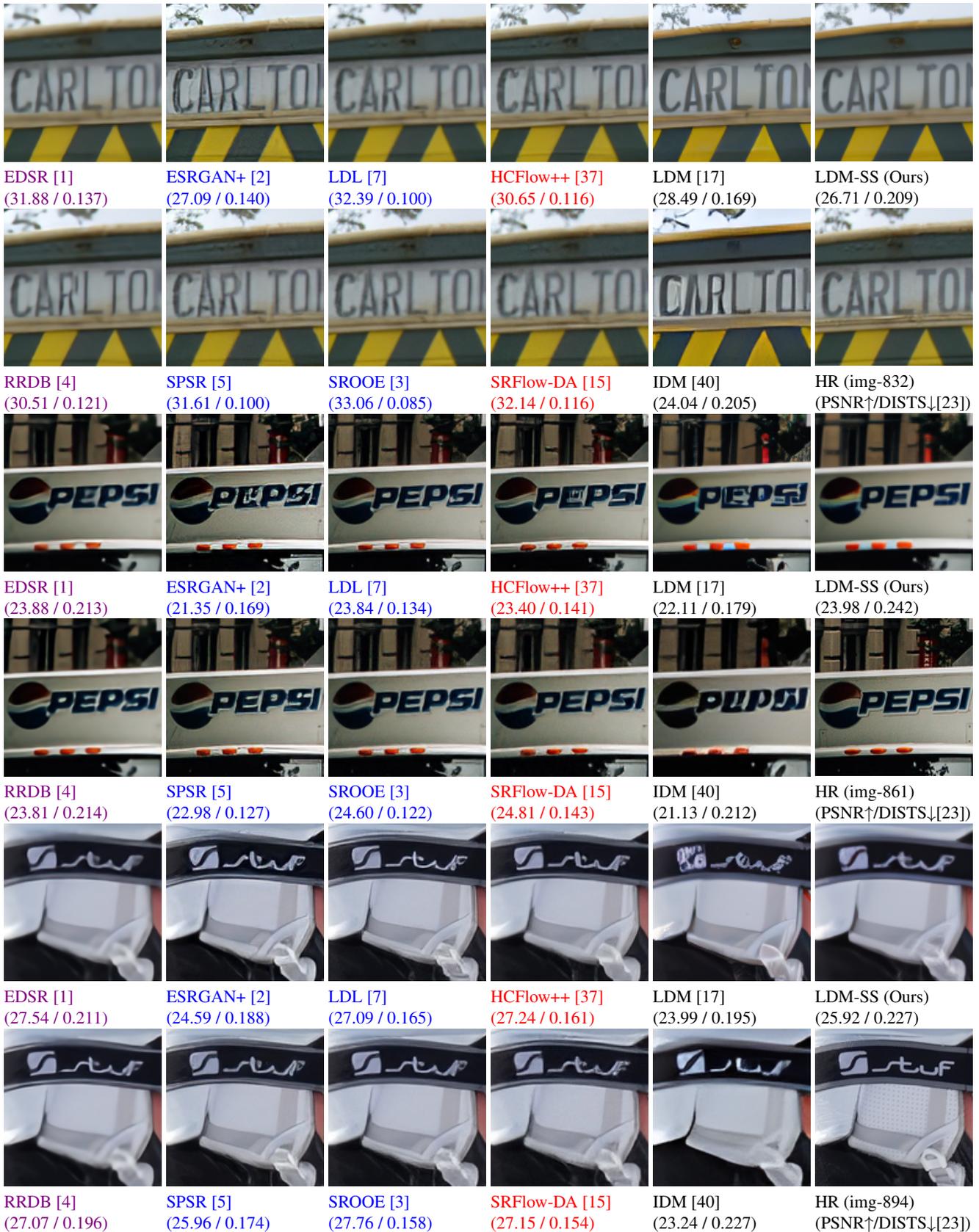


Fig. 6: Comparison of the proposed method with the state-of-the-art for $\times 4$ SR on natural images from DIV2K validation set [39]. Even though the proposed LDM-SS method with human feedback has clear advantages in reconstructing realistic high-frequency details while inhibiting artifacts, popular quantitative metrics are insufficient to reflect the visual improvements.

	SR Model	PSNR \uparrow	LR Consistency \uparrow	SSIM \uparrow	LPIPS	LPIPS _{VGG} \downarrow	PieAPP \downarrow	DISTS \downarrow	NRQM \uparrow
Regressive	EDSR	25.962	43.047	0.803	0.115	0.231	0.901	0.194	5.142
	RRDB	25.316	39.508	0.788	0.103	0.225	0.799	0.187	5.850
	HAT	27.408	44.673	0.826	0.089	0.201	0.686	0.179	5.518
GAN-based	ESRGAN+	22.666	31.718	0.716	0.083	0.224	0.292	0.168	7.757
	SPSR	24.760	36.520	0.762	0.063	0.184	0.523	0.138	7.159
	LDL	27.194	43.360	0.852	0.053	0.145	0.396	0.125	7.079
	SROOE	25.894	41.040	0.790	0.061	0.166	0.562	0.132	6.741
Flow-based	SRFlowDA	27.510	46.929	0.852	0.062	0.172	0.686	0.145	6.699
	HCFlow	25.062	43.302	0.777	0.067	0.183	0.641	0.141	6.896
Diffusion-based	SR3	21.596	25.587	0.683	0.231	0.299	2.065	0.357	6.649
	LDM	24.234	29.655	0.780	0.122	0.244	0.898	0.185	5.794
	IDM	24.573	29.526	0.716	0.149	0.294	0.651	0.227	6.496
	LDM-SS	26.047	31.447	0.823	0.141	0.227	1.120	0.194	5.195

Table 2: Performance comparison of $\times 4$ ($32 \times 32 \rightarrow 128 \times 128$) SR models on DIV2K validation set. Even though LDM-SS successfully suppresses distortions, contributing to visually enhanced and reliable outputs, there exists an intriguing divergence between the notable visual performance and the quantitative measures.

first digit presented by state-of-the-art methods can be perceived as “6.” On the contrary, LDM-SS provides the accurate digit “5.” Similarly, second digit obtained by SOTA methods is not a clear “8”, unlike the output of LDM-SS. In addition, the qualitative results from DIV2K [39] dataset demonstrated in 6 prove LDM-SS visibly suppresses unwanted distortions and enhances overall image quality in a perceptually meaningful way and enables visually on-par or better results with SOTA SR methods.

5. CONCLUSION

DMs are able to generate not one but a set of plausible SR images at their output. While this improves diversity, it brings the ambiguity of how to select the single trustworthy SR solution when the goal is to extract crucial information from LR images. In this work, we are primarily interested in obtaining a consistent and reliable image SR result within the space spanned by diffusion models. Specifically, we benefit from human feedback while selecting diverse set of diffusion samples since we found that we cannot rely on quantitative metrics to select a trustworthy result. Then, we ensemble the selected images to resolve the ambiguity in SR images and to obtain a reliable, photo-realistic solution. Our approach achieves promising results on DIV2K validation set and information-centric applications. While our method, LDM-SS, excels in delivering visually appealing results by reducing artifacts, these improvements may not be accurately reflected in quantitative scores like PSNR, MS-SSIM, LPIPS, DISTS, etc. The proposed method for selection of diffusion samples is generic in the sense that any off-the-shelf diffusion model can be easily plugged into this framework to benefit from human feedback to resolve the ambiguity.

6. REFERENCES

- [1] B. Lim, S. Son, H. Kim, S. Nah, and K. Lee, “Enhanced deep residual networks for single image super-resolution,” in *IEEE/CVF CVPR Workshops*, 2017.
- [2] N. Rakotonirina and A. Rasoanaivo, “EsrGAN+: Further improving enhanced super-resolution generative adversarial network,” in *IEEE ICASSP*, 2020, pp. 3637–3641.
- [3] S. Park, Y. Moon, and N. Cho, “Perception-oriented single image super-resolution using optimal objective estimation,” in *IEEE/CVF Conf. on Comp. Vision and Patt. Recog. (CVPR)*, 2023, pp. 1725–1735.
- [4] X. Wang, Ke Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. C. Loy, “ESRGAN: enhanced super-resolution generative adversarial networks,” in *European Conf. on Comp. Vision (ECCV) Workshops*, 2018.
- [5] C. Ma, Y. Rao, Y. Cheng, C. Chen, J. Lu, and J. Zhou, “Structure-preserving super resolution with gradient guidance,” in *IEEE/CVF Conf. Comp. Vis. and Patt. Recog. (CVPR)*, 2020.
- [6] X. Chen, X. Wang, J. Zhou, Y. Qiao, and C. Dong, “Activating more pixels in image super-resolution transformer,” in *IEEE/CVF Conf. Comp. Vision and Patt. Recog. (CVPR)*, 2023.
- [7] J. Liang, H. Zeng, and L. Zhang, “Details or artifacts: A locally discriminative learning approach to realistic image super-resolution,” in *IEEE/CVF Conf. on Comp. Vision and Patt. Recog. (CVPR)*, 2022, pp. 5657–5666.
- [8] J. Huang, A. Singh, and N. Ahuja, “Single image super-resolution from transformed self-exemplars,” in *IEEE Conf. on Comp. Vision and Patt. Recog. (CVPR)*, 2015.
- [9] C. Dong, Loy C., K. He, and X. Tang, “Image super-resolution using deep convolutional networks,” *IEEE Trans. on Pattern Analysis and Mach. Intell.*, vol. 38, pp. 295–307, 2016.
- [10] J. Kim, J. Kwon Lee, and K. Mu Lee, “Accurate image super-resolution using very deep convolutional networks,”

- IEEE/CVF Conf. on Comp. Vision and Patt. Recog. (CVPR)*, pp. 1646–1654, 2016.
- [11] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi, “Photo-realistic single image super-resolution using a generative adversarial network,” *IEEE/CVF Conf. on Comp. Vision and Patt. Recog. (CVPR)*, pp. 105–114, 2017.
- [12] G. Lin, Q. Wu, X. Huang, Q. Lida, and X. Chen, “Deep convolutional networks-based image super-resolution,” in *Int. Conf. Intelligent Computing*, 2017, pp. 338–344.
- [13] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, “Image super-resolution using very deep residual channel attention networks,” in *IEEE/CVF ECCV*, 2018.
- [14] L. Andreas, D. Martin, L. Van Gool, and R. Timofte, “Srflow: Learning the super-resolution space with normalizing flow,” in *ECCV*, 2020.
- [15] Y. Jo, S. Yang, and S. Joo Kim, “Srflow-da: Super-resolution using normalizing flow with deep convolutional block,” in *IEEE/CVF Conf. on Comp. Vision and Patt. Recog. (CVPR) Workshops*, June 2021.
- [16] C. Saharia, J. Ho, W. Chan, T. Salimans, D. Fleet, and M. Norouzi, “Image super-resolution via iterative refinement,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 4, pp. 4713–4726, 2022.
- [17] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *IEEE/CVF Conf. on Comp. Vision and Patt. Recog. (CVPR)*, 2022, pp. 10684–10695.
- [18] C. Korkmaz, A. M. Tekalp, Z. Doğan, E. Erdem, and A. Erdem, “Perception-distortion trade-off in the sr space spanned by flow models,” in *IEEE Int. Cong. on Image Processing (ICIP)*. IEEE, 2022, pp. 2396–2400.
- [19] F. Luo, J. Xiang, J. Zhang, X. Han, and W. Yang, “Image super-resolution via latent diffusion: A sampling-space mixture of experts and frequency-augmented decoder approach,” *arXiv preprint arXiv:2310.12004*, 2023.
- [20] Andreas Lugmayr and et. al., “Ntire 2021 learning the super-resolution space challenge,” in *IEEE Conf. on Comp. Vision and Patt. Recog. Workshops (CVPRW)*, 2021, pp. 596–612.
- [21] Andreas Lugmayr and et. al., “Ntire 2022 challenge on learning the super-resolution space,” in *IEEE Conf. on Comp. Vision and Patt. Recog. Workshops (CVPRW)*, 2022, pp. 785–796.
- [22] R. Zhang, P. Isola, A. Efros, E. Shechtman, and O. Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” in *IEEE/CVF Conf. Comp. Vision and Patt. Recog. (CVPR)*, 2018, pp. 586–595.
- [23] K. Ding, K. Ma, S. Wang, and E. Simoncelli, “Image quality assessment: Unifying structure and texture similarity,” *IEEE Trans. Patt. Anal. Mach. Intel.*, vol. 44, pp. 2567–2581, 2020.
- [24] Y. Zhang, B. Ji, J. Hao, and A. Yao, “Perception-distortion balanced admm optimization for single-image super-resolution,” in *European Conf. on Comp. Vision (ECCV)*, 2022.
- [25] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems*, 2014, vol. 27.
- [26] Z. Liu, W. Siu, and L. Wang, “Variational autoencoder for reference based image super-resolution,” in *IEEE/CVF Conf. on Comp. Vision and Patt. Recog. (CVPR)*, 2021, pp. 516–525.
- [27] H. Zhou, C. Huang, S. Gao, and X. Zhuang, “Vspsr: Explorable super-resolution via variational sparse representation,” in *IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2021, pp. 373–381.
- [28] K. Song, D. Shim, K. Kim, J. Lee, and Y. Kim, “FS-NCSR: Increasing diversity of the super-resolution space via frequency separation and noise-conditioned normalizing flow,” in *IEEE/CVF Conf. Comp. Vis. and Patt. Recog. Workshops (CVPRW)*, 2022, pp. 967–976.
- [29] A. Van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals, A. Graves, et al., “Conditional image generation with pixelcnn decoders,” *Advances in neural information processing systems*, vol. 29, 2016.
- [30] M. Chen, A. Radford, R. Child, J. Wu, H. Jun, D. Luan, and I. Sutskever, “Generative pretraining from pixels,” in *Int. Conf. on Mach. Learning*, 2020, pp. 1691–1703.
- [31] T. Salimans, A. Karpathy, X. Chen, and D. Kingma, “Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications,” in *Int. Conf. on Learning Representations*, 2016.
- [32] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” in *Advances in Neural Information Processing Systems*, 2020.
- [33] Y. Chen, S. Liu, and X. Wang, “Learning continuous image representation with local implicit image function,” in *IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 8628–8638.
- [34] O. Avrahami, D. Lischinski, and O. Fried, “Blended diffusion for text-driven editing of natural images,” in *IEEE/CVF Conf. Comp. Vis. and Patt. Recog. (CVPR)*, 2022, pp. 18208–18218.
- [35] S. Gu, D. Chen, J. Bao, F. Wen, B. Zhang, D. Chen, L. Yuan, and B. Guo, “Vector quantized diffusion model for text-to-image synthesis,” in *IEEE/CVF Conf. on Comp. Vision and Pattern Recog. (CVPR)*, 2022, pp. 10696–10706.
- [36] A. Blattmann, R. Rombach, K. Oktay, J. Müller, and B. Ommer, “Retrieval-augmented diffusion models,” in *Advances in Neural Info. Processing Systems*, 2022, pp. 15309–15324.
- [37] J. Liang, A. Lugmayr, K. Zhang, M. Danelljan, L. Van Gool, and R. Timofte, “Hierarchical conditional flow: A unified framework for image super-resolution and image rescaling,” in *IEEE Int. Conf. on Computer Vision*, 2021.
- [38] Li Deng, “The mnist database of handwritten digit images for machine learning research,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [39] E. Agustsson and R. Timofte, “NTIRE 2017 Challenge on single image super-resolution: Dataset and study,” in *IEEE/CVF Conf. Comp. Vis. and Patt. Recog. (CVPR) Workshops*, 2017.
- [40] S. Gao, X. Liu, B. Zeng, S. Xu, Y. Li, X. Luo, J. Liu, X. Zhen, and B. Zhang, “Implicit diffusion models for continuous super-resolution,” in *IEEE/CVF Conf. Comp. Vis. and Patt. Recog.*, 2023, pp. 10021–10030.