

APISR: Anime Production Inspired Real-World Anime Super-Resolution

Boyang Wang¹ Fengyu Yang^{1,2*} Xihang Yu¹ Chao Zhang³ Hanbin Zhao^{3†}

¹University of Michigan ²Yale University ³Zhejiang University

Abstract

While real-world anime super-resolution (SR) has gained increasing attention in the SR community, existing methods still adopt techniques from the photorealistic domain. In this paper, we analyze the anime production workflow and rethink how to use characteristics of it for the sake of the real-world anime SR. First, we argue that video networks and datasets are not necessary for anime SR due to the repetition use of hand-drawing frames. Instead, we propose an anime image collection pipeline by choosing the least compressed and the most informative frames from the video sources. Based on this pipeline, we introduce the Anime Production-oriented Image (API) dataset. In addition, we identify two anime-specific challenges of distorted and faint hand-drawn lines and unwanted color artifacts. We address the first issue by introducing a prediction-oriented compression module in the image degradation model and a pseudo-ground truth preparation with enhanced hand-drawn lines. In addition, we introduce the balanced twin perceptual loss combining both anime and photorealistic high-level features to mitigate unwanted color artifacts and increase visual clarity. We evaluate our method through extensive experiments on the public benchmark, showing our method outperforms state-of-the-art anime dataset-trained approaches. The code is available at <https://github.com/Kiteretsu77/APISR>.

1. Introduction

As an important subdiscipline of real-world super-resolution (SR), anime SR focuses on restoring and enhancing low-quality low-resolution (LR) anime visual art images and videos to high-quality high-resolution (HR) forms. It has demonstrated significant practical impacts in the fields of entertainment and commerce [46, 48, 54, 56, 61]. An emerging line of work has addressed the problem by extending SR networks to capture multi-scale information or learning an adaptive degradation model [46, 56]. We argue these methods lack understanding of the anime domain as

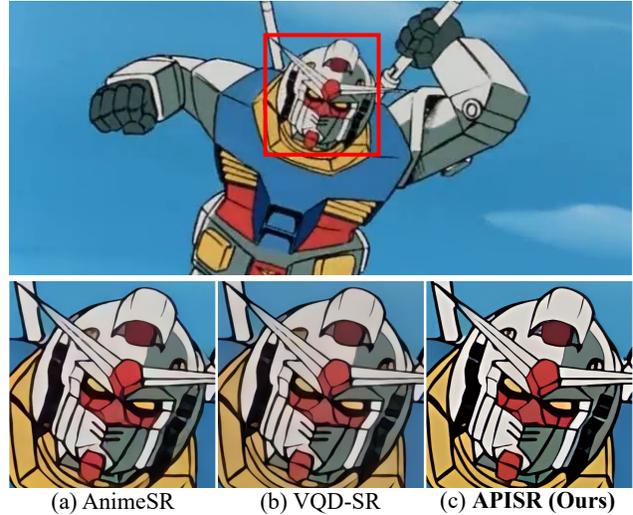


Figure 1. **Comparisons between proposed APISR and other SOTA anime SR methods.** Ours present clearer and sharper hand-drawn lines, better restoration with more natural details, and do not present unwanted color artifacts. **Zoom in for best view.**

their techniques are directly transplanted from the photorealistic SR approach.

In this paper, we thoroughly analyze the anime production process, exploring ways to leverage its unique aspects for practical applications in anime SR. The production workflow first starts with hand-drawing sketches on paper, which are then colored and enhanced by computer-generated imagery (CGI) processing [69]. Then, these processed sketches are concatenated into a video. Due to the fact that the drawing process is extremely labor-intensive and human eyes are not sensitive to motions [10, 37], it is a standard practice to reuse a single image across multiple consecutive frames when forming the video. This procedure in production motivates us to rethink whether it is necessary and efficient to use video networks and video datasets to train SR networks in the anime domain.

To this end, we explore the use of image-based methods and datasets as a unified super-resolution and restoration framework for both anime images and videos. Creating an image dataset allows us more flexibility to exclusively choose the least-compressed video frames as our potential

† Corresponding author

* works done at University of Michigan

dataset pool, rather than gathering sequential frames that contain temporal distortions to create a video dataset. Furthermore, by forming an image dataset, we can selectively focus on the most informative frames, as anime videos typically possess less information than photorealistic videos. If we randomly crop a patch from an anime image, there is a high probability that it is a monochromatic area signifying a lack of information. In light of these phenomena, we introduce an anime image collection pipeline that focuses on keyframes in video, along with an image complexity assessment-based selection criteria. This method is designed to identify and select the least-compressed and the most informative images from video sources. Using our pipeline, we propose **Anime Production-oriented Image (API)** dataset for SR training.

In addition, we identify two new anime-specific challenges for real-world SR tasks. First, in anime production, the clarity of hand-drawn lines is a highly emphasized detail [6, 26, 56] as shown in Fig. 2 a, but hand-drawn lines are easily weakened due to compression in internet transmission and physical aging in production. This deterioration at the edges of lines exerts a substantial negative impact on the visual effects. To address this, we start from the perspective of restoration and enhancement. Concretely, we propose a prediction-oriented compression module in the image degradation model to simulate compression in internet transmission such that the model trained with this self-supervised method can restore hand-drawn line distortions. In addition, we propose a ground-truth (GT) enhancement approach to enhance faint, aging hand-drawn lines, by merging hand-drawn lines extracted from the overly sharpened GT images.

Second, we realize an issue of unwanted color artifacts in anime images, which is a consequence of employing the GAN-based SR networks [17] (see Fig. 2 b). These artifacts are presented as irregularly shaped colored spots with varying intensities that are scattered randomly across generated images, which significantly undermines visual perception. We attribute this issue to the reason that image features of perceptual loss are trained on the photorealistic image datasets, which is inconsistent in the anime domain. To mitigate this issue, we conduct a comprehensive study of perceptual loss and introduce balanced twin perceptual loss, which assembles perceptual features from both the photorealistic domain and the anime domain by a balanced layer scaling distribution.

Thus, we summarize our contributions as follows:

- We propose a novel anime dataset curation pipeline that is capable of collecting the least compressed and the most informative anime images from video sources.
- We propose an image degradation model to deal with harder compression restoration challenges, especially for hand-drawn line distortions, and the first methodologies

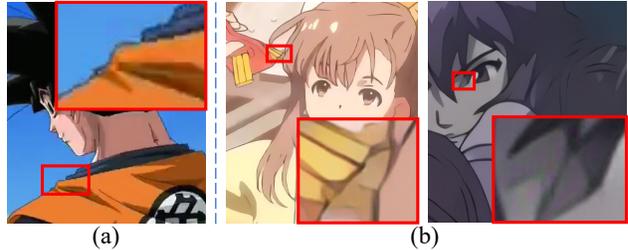


Figure 2. **We identify two new anime-specific challenges:** (a) Distorted and faint hand-drawn lines frequently appear in real-world anime images. (b) Unwanted color artifacts in AnimeSR [56] and VQD-SR [46]. **Zoom in for the best view.**

in the anime domain to attentively enhance faint hand-drawn lines.

- We realize and address the unwanted color artifacts in GAN-based SR network training caused by the domain inconsistency of the perceptual loss.
- We thoroughly evaluate our method on the real-world anime SR dataset and show that our method outperforms state-of-the-art anime dataset-trained SR approaches by a large margin with only 13.3% training sample complexity of the prior work.

2. Related Work

Real-World Super-Resolution. Classical SR methods [7, 15, 57, 58] typically employs a straightforward approach, using a single bicubic downsampling operation to convert high-resolution (HR) ground-truth (GT) images into their low-resolution (LR) counterparts. Classical image restoration methods [30, 31, 59, 60] train different weights for different tasks. In contrast, real-world SR is dedicated to implementing a sophisticated degradation model by one model weight to restore the diverse degradations found in the real-world scenario, such as blurring, noise, and compression [24, 30, 31, 48, 54, 65].

Generic degradation model design can be broadly classified into two categories: explicit models [24, 30, 48, 54, 65] and implicit models [46, 56, 64]. Explicit degradation models employ kernels and mathematical formulas to simulate real-world degradation processes. On the other hand, the implicit degradation models focus on training neural networks to capture the distribution of real-world degradations. Nevertheless, implicit models face challenges of interpretability and scalability. The efficacy of implicit models lacks a clear rationale, and adapting them to new domains requires the creation of bespoke datasets and extra training complexity.

Anime Processing. Anime represents a distinctive form of visual art, often characterized by exaggerated visual representation. Creators of anime typically start by sketching line art, followed by 2D and 3D animation techniques, which include elements like colorization, CGI effects, and

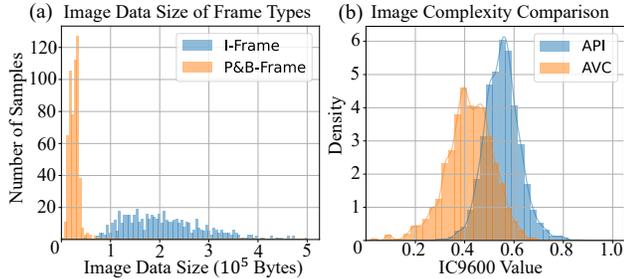


Figure 3. Histogram of (a) the average image data size comparison between I-Frames and Non-I-Frames (P and B-Frames) in collected video sources and (b) image complexity [16] comparison between proposed API and AVC [56] dataset.

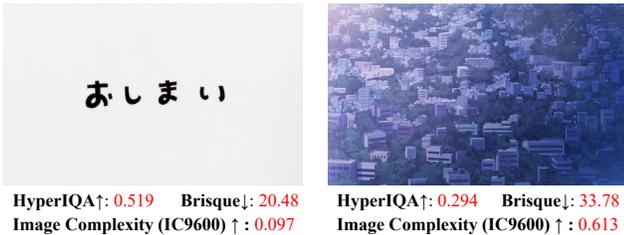


Figure 4. Image Quality Assessment (IQA) with HyperIQA [43] and Brisque [33] vs. Image Complexity Assessment (ICA) with IC9600 [16]. IQA favors simple scenes and gives low scores to images with strong CGI. However, ICA is the opposite.

frame interpolation. Notably, recent research in the realm of anime has garnered substantial attention, *e.g.*, AI painting with anime content [18, 22, 67], vectorization of anime images [63, 68], anime interpolation and inbetweening [10, 40, 42], anime sketch colorization [5, 6, 13, 50, 66], 3D representation [11, 41], and anime domain adaptation [26].

AnimeSR (NeurIPS 2022) [56] and VQD-SR (ICCV 2023) [46] are two recent representative studies in the domain of real-world anime super-resolution tasks. However, they have not fully addressed the unique challenges of low-level anime restoration. This includes the faint hand-drawn lines and domain inconsistency in the training of GAN-based networks. This paper conducts a comprehensive exploration of several meticulously crafted approaches to the anime SR domain.

3. Proposed Method

3.1. Anime Production-Oriented Image SR Dataset

In this section, we present the **API (Anime Production-oriented Image) SR dataset** and its curation workflow. This curation leverages the characteristics of anime videos to select the least compressed and the most informative frames.

I-Frame-based Image Collection. AnimeSR introduces AVC-Train, the first video-based anime SR dataset, but they overlook the impact of compression during the collection process, which leads VQD-SR to propose a post-processing technique to enhance the dataset. Instead, we propose a



Figure 5. **Samples of API Super-Resolution Dataset.** API includes versatile CGI effects scenes (*e.g.*, different lightning and special effects) and presents high image complexity.

novel method to select the least compressed frames from the source level with minimum effort.

All videos on the internet are compressed and encapsulated with a video compression standard (*e.g.*, H.264 [36] and H.265 [44]) for a trade-off between the quality and the data size. There are numerous video compression standards, each with a complex engineering system, but they share a similar backbone design. This characteristic motivates us to find the pattern that the compression quality assigned to each frame is different. Video compression designates some keyframes, known as I-Frames, as individual units for compression. Empirically, I-Frames are the first frame of scene-changing scenarios. These I-Frames are allocated with a high data size budget. On the contrary, a higher compression ratio requires non-I-Frames, namely P-Frames and B-Frames, to take I-Frames as the reference during compression, which introduces temporal distortions. As shown in Fig. 3 a, among the anime videos we collect, I-Frames on average have a much higher data size than other non-I-Frames, which genuinely stand for higher quality. Thus, we use *ffmpeg*, a video processing tool, to extract all I-Frames from the video source as an initial pool.

Image Complexity-based Selection. To further select idealistic images from the I-Frames pool, we need some criteria. A straightforward method involves following AVC-Train to use the Image Quality Assessment (IQA) to rank and choose frames with better scores. However, IQA ranking does not prefer anime images with CGI effects but favors simple scenes with little information (see Fig. 4). Thus, we argue that image complexity assessment (ICA) is a better option in the anime domain.

ICA evaluates the level of intricacy in an image by scoring the amount and variety of details present. Compared to IQA, ICA demonstrates greater robustness against changes in saturation, lightning, contrast, and motion blurring. The ICA metric we use is a recent rising analysis network, IC9600 [16]. In the anime domain, employing ICA presents two primary advantages. First, many scenes in anime videos are characteristically monotonous (as exemplified in Fig. 4 left), where the majority of pixels lack significant information in training. IQA favors these simple images and gives higher score compared to other images, but ICA enables the

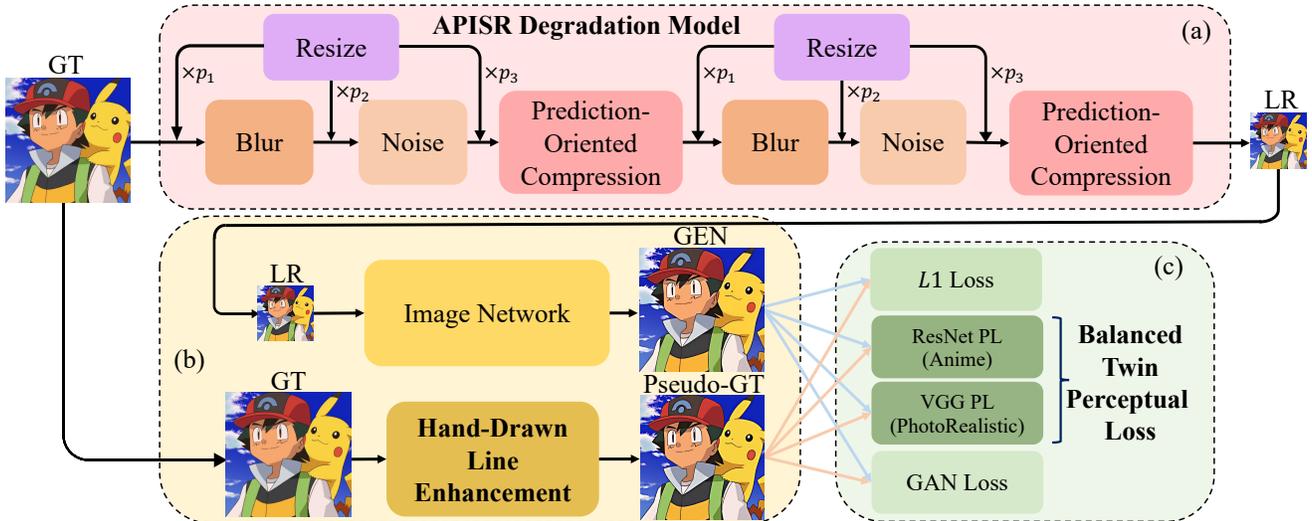


Figure 6. **The overview of our proposed methods.** (a) We proposed a prediction-oriented compression module in the degradation model to simulate versatile compression degradations for a single image input (detailed in Sec. 3.2). Proposed shuffled resize module is randomly positioned to augment the representation of the degradation model. (b) GT images are augmented with proposed hand-drawn line enhancement to promote the generation of images with sharpened line edge details in training (detailed in Sec. 3.3). (c) Proposed balanced twin perceptual loss avoids unwanted color artifacts in GAN network training (detailed in Sec. 3.4).

exclusion of these scenes, which, in turn, contributes to a reduced training sample complexity. Second, ICA is more adept at identifying meaningful scenes within anime production, especially those featuring CGI effects, such as the dark scene in Fig. 4 right. These are scenarios where IQA methods typically falter. By collecting versatile scenes, the network training can become more robust in handling complex real-world anime inputs.

API Dataset. We began by manually sourcing 562 high-quality anime videos. From these, we extracted all I-Frames as an initial selection pool. Utilizing the image complexity assessment method mentioned above, we then selected the top 10 highest-scoring frames from the I-Frames pool of each video. After discarding inappropriate images (*e.g.*, nudity, violence, abnormality, and anime images mixed with photorealistic content), 3,740 high-quality images are obtained as our proposed dataset. Example images are shown in Fig. 5. Moreover, as shown in Fig.3 b, the density of high image complexity scored frames of our API dataset is remarkably superior to that of AVC-Train. More analysis and data can be found in the supplementary materials.

720P Back-to-Original Production Resolution. While studying the anime production pipeline, we observed that most anime productions follow a 720P format (with an image height of 720 pixels). However, in real-world scenarios, anime is often falsely upscaled to 1080P or other formats, for the sake of standardizing multimedia formats. We empirically find that rescaling all anime images back to the original 720P can provide feature density envisioned by the creators with more compact anime hand-drawn lines and CGI information.



Figure 7. H.264 [36] compression of regular multi-frame video compression and our proposed single-frame compression. They exert similar degradations (*e.g.*, distortion to hand-drawn lines).

3.2. An Anime Practical Degradation model

In the real-world SR, the design of the degradation model is of great importance. Based on the high-order degradation model [54] and a recent image-based video compression restoration model [48], we propose two improvements to restore distorted hand-drawn lines and versatile compression artifacts and to augment the representation of the degradation model. Our degradation model is shown in Fig. 6 a.

Prediction-Oriented Compression. Utilizing the image degradation model presents a challenge in the anime restoration of video compression artifacts. This is because previous real-world image SR methods employ JPEG, an old but widely-used image compression standard, as the sole compression module in the image degradation model. JPEG performs repetitive and independent compression on all encoding units, without considering the existence of other units. However, video compression algorithms, for higher compression ratios, apply prediction algorithms to search for a reference with similar pixel content and only compress their differences (residual), thereby reducing in-

formation entropy. Prediction algorithms can search their reference spatially (intra-prediction) or temporally (inter-prediction). Regardless of the category, the intrinsic cause of distortion comes from the misalignment in residual due to prediction limitation.

Hence, we argue that artifacts equivalent to real-world video compression artifacts can be synthesized using a single image input in conjunction with a prediction-oriented compression algorithm (*e.g.*, WebP [38] and H.264). The need for genuinely sequential frames is not necessary. To this end, we design a prediction-oriented compression module within the image degradation model. This module requires video compression algorithms to compress inputs on a single-frame basis. Compared to VCISR [48], we don't need multiple frames for one turn of execution of compression. This methodology is theoretically reasonable and practically viable from an engineering perspective. With a single-frame input, video compression trivially applies intra-prediction to compress the frame without using its inter-prediction functionality. Utilizing this approach, the image degradation model is capable of synthesizing compression artifacts akin to those observed in conventional multi-frame video compression as shown in Fig. 7. Subsequently, by feeding these synthesized images into the image SR network, the system can effectively learn the patterns of versatile compression artifacts and engage in the restoration.

Shuffled Resize Module. Degradation models in the real-world SR domain consider blurring, resize, noise, and compression modules. Blurring, noise, and compression are real-world artifacts that can be synthesized with clear mathematical models or algorithms. However, the logic of the resize module is entirely different. Resize is not a part of natural image generation but is introduced solely for SR-paired dataset purposes. Given this notion, we believe that previous fixed resize module is not very suitable. We propose a more robust and effective solution, which involves randomly placing resize operations at various orders in the degradation model.

3.3. Anime Hand-Drawn Lines Enhancement

To enhance faint hand-drawn lines, directly employing global methods, such as modifying the degradation model or sharpening the entire GT, is not an ideal approach, as the network cannot learn with attention to hand-drawn line changes. Thus, we choose to extract sharpened hand-drawn line information and merge it back with GT to form pseudo-GT. By introducing this attentively enhanced pseudo-GT to SR training, the network can generate sharpened hand-drawn lines without the need to introduce additional neural network modules or separate post-processing networks.

To extract hand-drawn lines, a direct approach is to apply a sketch extraction model. However, current learning-based

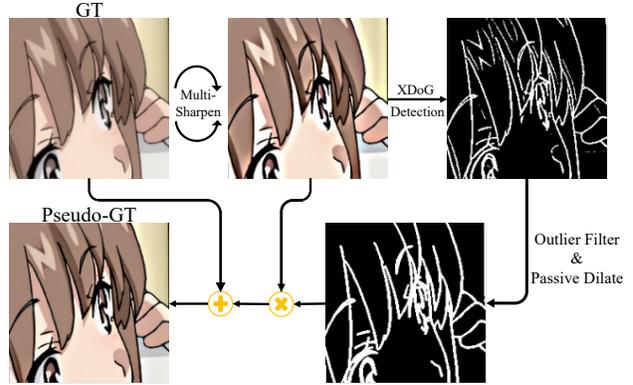


Figure 8. **Anime Hand-Drawn Lines Enhancement Pipeline.**

sketch extraction is often characterized by a style transfer to the reference image, which distorts hand-drawn line details and encompasses unrelated pixel content (*e.g.*, shadows and edges of CGI effects). Consequently, we need a more granular, pixel-by-pixel methodology to extract hand-drawn lines. Thus, we utilize XDoG [55], a pixel-by-pixel Gaussian kernel-based sketch extraction algorithm, to extract edge maps from the sharpened GT. Nevertheless, XDoG edge maps are marred by excessive noise, containing outlier pixels and fragmented line representations. To address this ill-posed issue, we propose an outlier filtering technique coupled with a custom-designed passive dilation method (detailed in the supplementary materials). In this way, we yield a more coherent and undisturbed representation of hand-drawn lines.

We empirically find that overly sharpened pre-processed GT makes the hand-drawn line margins more noticeable than other unrelated shadow edge details, which makes the outlier filter easier to distinguish their differences. Thus, we propose three rounds of unsharp masking to the GT first. To sum up, the formula is as follows:

$$I_{\text{Sharp}} = f^n(I_{\text{GT}}), \quad (1)$$

$$I_{\text{Map}} = h(g(I_{\text{Sharp}})), \quad (2)$$

$$I_{\text{pseudo-GT}} = I_{\text{Sharp}} \cdot I_{\text{Map}} + I_{\text{GT}} \cdot (1 - I_{\text{Map}}), \quad (3)$$

where f is the sharpening function that recursively executes n times, g denotes XDoG edge detection and h stands for post-processing techniques of passive dilation with outlier filtering. I_{Map} is a binary value map. The visual pipeline is shown in Fig. 8.

3.4. Balanced Twin Perceptual Loss for Anime

The existence of unwanted color artifacts is attributed to the inconsistent dataset domain in training between the generator and perceptual loss. Currently, most SR models trained with GAN, including AnimeSR and VQD-SR, use the same ImageNet [14] pre-trained VGG [39] network as the perceptual loss. However, anime content, particularly those

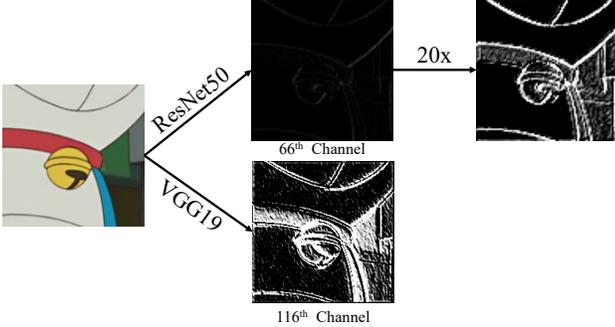


Figure 9. The second middle-layer feature outputs comparison between VGG19 used by photo-realistic perceptual loss [27] and ResNet50 used by anime recognition task [3, 4]. With scaling, ResNet50 presents a similar intensity as the VGG outputs.

mixed with CGI and extensive illustrations, differs significantly from photorealistic features in ImageNet. To tackle this problem, we investigate perceptual loss and the subsequent improvements made in their following works.

The core idea behind perceptual loss is to utilize high-level features (*e.g.*, segmentation, classification, recognition) to complement low-level pixel features by comparing middle-layer feature outputs. In this regard, we employ a pre-trained ResNet50 [3, 21] on anime object classification task with Danbooru [4] dataset, a substantial and rich tagging anime illustration database. Since the pre-trained network is ResNet50 instead of VGG, we propose a similar middle-layer comparison (detailed in the supplementary material). Overall, the formula is as follows:

$$L_{ResNet}^{\phi}(\hat{y}, y) = \sum_j \frac{w_j}{C_j H_j W_j} |\phi_j(\hat{y}) - \phi_j(y)|, \quad (4)$$

where y and \hat{y} are the pseudo-GT by Sec. 3.3 and the generated images. ϕ_j represents the perceptual function that returns j th layer output of ResNet50. C_j , H_j , and W_j are dimensions of the layer output and w_j is the scaling factor for each layer. There are 5 middle-layer feature outputs, which is the same quantity as VGG-based perceptual loss. We also observe that the intensity of shallow feature layers in ResNet50 is very weak (see Fig. 9). To resemble a similar intensity balance as the VGG, we apply a high w_j to the early layers, which leads to stable training.

Notably, introducing the ResNet-based perceptual loss as the sole perceptual loss can solve unwanted color artifacts and lead to quantitative improvements. However, there may be instances of poor visual results. This is attributed to the inherent bias in the Danbooru dataset, where most images are character faces or relatively simple illustrations. Hence, we seek a tradeoff by using real-world features as an auxiliary primer to guide the ResNet-based perceptual loss in training. This approach results in visually appealing images and also resolves the unwanted color issue. The overall loss

function for our GAN training is defined as follows:

$$L = \alpha L_1 + \beta L_{per} + \gamma L_{adv}, \quad (5)$$

$$L_{per} = L_{ResNet} + \delta L_{VGG}, \quad (6)$$

where L_1 , L_{VGG} , and L_{adv} are L1 pixel loss, photorealistic VGG-based perceptual loss, and the adversarial loss. α , β , γ and δ are weight parameters.

4. Experiment

4.1. Implementation Details

In our experiment, we employ our proposed API dataset as the training dataset for the image network. The image network we utilize is a tiny version of GRL [30] with the nearest convolution upsample module (detailed in the supplementary).

To train the GAN, we follow the same two-stage training approach as prior works [8, 53, 54, 56, 65]. In the first stage, we train the network with L1 pixel loss for 300K iterations. In the second stage, we introduce our balanced twin perceptual loss and the adversarial loss, conducting an additional 300K iterations. The weights of $\{\alpha, \beta, \gamma, \delta\}$ are $\{1, 0.5, 0.2, 1\}$ respectively. The layer weight of perceptual loss is $\{0.1, 20, 25, 1, 1\}$ for ResNet and $\{0.1, 1, 1, 1, 1\}$ for VGG. Our discriminator is the same three-scale patch discriminator [23, 35, 51] as in AnimeSR [56] and VQD-SR [46]. We use the Adam optimizer [28] with a learning rate of 2×10^{-4} in the first stage and 1×10^{-4} in the second stage. A learning rate decay is applied every 100K iterations in both stages. The entire training process was carried out on one Nvidia RTX 4090, with HR patch sizes set at 256x256 and a batch size of 32.

As for the degradation model, we perform degradation on the whole HR image first rather than directly on a cropped patch as in previous works [30, 48, 54, 65]. Within the degradation model, noise and blurring are configured identically to Real-ESRGAN [54], and the first prediction-oriented compression is implemented with JPEG [47] and WebP [38]. The second prediction-oriented compression includes AVIF [19], JPEG [47], WebP [38], and single-frame compression of MPEG2 [32], MPEG4 [2], H.264 [36], and H.265 [44]. The probability of placing the resize module is equally divided among all positions. Specific parameter settings can be found in our supplementary materials.

4.2. Comparisons with State-of-the-art Methods

We compare our APISR quantitatively and qualitatively with other SOTA real-world image and video SR methods, which include Real-ESRGAN [54], BSRGAN [65], Real-BasicVSR [8], AnimeSR [56], and VQD-SR [46].

Quantitative Comparison. Following previous real-world SR works [8, 25, 46, 54, 56], we conduct inference on low-quality LR datasets to generate high-quality HR images

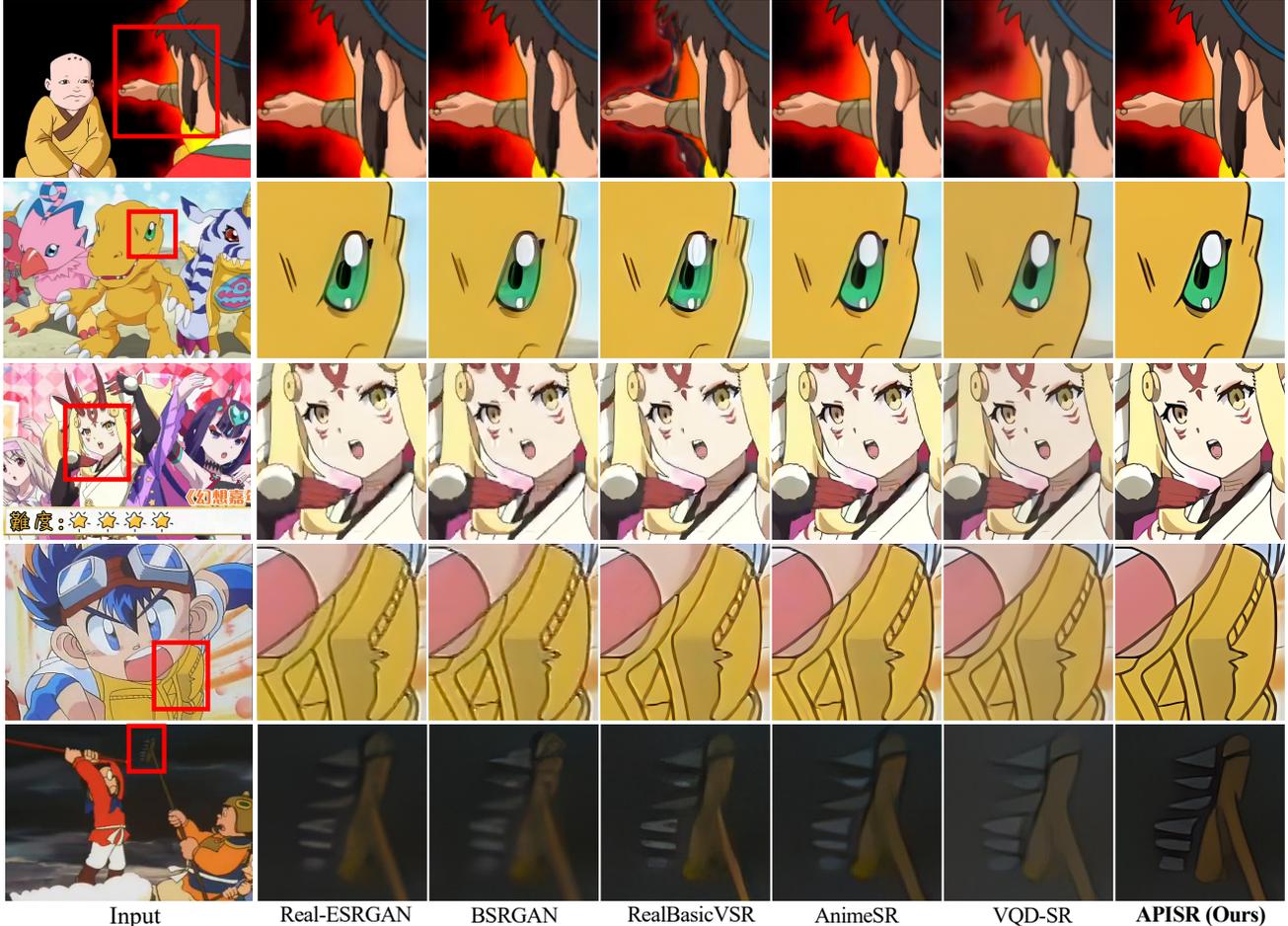


Figure 10. Qualitative comparisons on AVC-ReallQ [56] for 4× scaling. **Zoom in for the best view.**

Table 1. Quantitative comparisons on AVC-ReallQ [56]. **Bold** text indicates the best performance. (*’ denotes fine-tune on animation videos from [56])

Methods	Params ↓	NIQE ↓	MANIQA ↑	CLIPQA ↑
Real-ESRGAN* [54]	16.70	8.281	0.381	-
BSRGAN* [65]	16.70	8.632	0.376	-
RealBasicVSR* [8]	6.30	8.621	0.362	-
AnimeSR [56]	1.50	8.109	0.462	0.539
VQD-SR [46]	1.47	8.202	0.464	0.567
APISR (Ours)	1.03	6.719	0.514	0.711

and evaluate them using no-reference metrics. The scaling factor is 4 for all methods. To validate the effectiveness of our approach, our evaluation is based on AVC-ReallQ [56], which has 46 video clips each with 100 frames. This dataset is the only known dataset designed for real-world anime SR testing. For no-reference metrics, we employ the same metrics used in VQD-SR and AnimeSR, which are NIQE [34] and MANIQA [62]. We also incorporate other SOTA learning-based image quality assessment metrics like CLIPQA [49]. All metrics are based on pyiqa[9] library.

As shown in Tab. 1, our model has the smallest network size, 1.03M parameters, but has SOTA performance in all metrics among all image and video-based methods. Apart

from the various proposed methods that contribute to our success, special acknowledgment is due to the design of the prediction-oriented compression model, which enables us to train image datasets and image networks to restore video compression degradations. Meanwhile, it is worth mentioning that we achieved the result with only 13.3% and 25% of the training sample complexity of AnimeSR [56] and VQD-SR [46]. This is especially thanks to the introduction of image complexity assessment in dataset curation which selects informative images to increase the efficacy of learning the representation of anime images. Further, we require zero training on the degradation model due to the explicit degradation model we design.

Qualitative Comparison. As shown in Fig. 10, APISR greatly improves the visual quality than other methods. In restoring heavily compressed images, our model exhibits exceptional proficiency than all other methods, as exemplified in the first row, where we have much fewer ringing artifacts. Moreover, owing to the proposed hand-drawn lines enhancement, our generated images manifest increased line density and clarity as observed in the second row. In addressing various twisted lines and shadow arti-

Table 2. Ablation study results of different training datasets. IQA stands for image quality assessment. ICA stands for image complexity assessment.

Dataset	NIQE↓	MANIQA↑	CLIQQA↑
AVC-Train [56]	7.681	0.476	0.658
Random Select	8.006	0.446	0.625
I-Frame + IQA Select	7.876	0.493	0.675
I-Frame + ICA Select	6.912	0.499	0.683
I-Frame + ICA Select + 720P Rescale	6.719	0.514	0.711

Table 3. Ablation study results of different degradation model.

Degradation Model	NIQE↓	MANIQA↑	CLIQQA↑
High-Order [54]	6.667	0.483	0.663
Random Order [65]	6.975	0.491	0.674
Prediction-Oriented Compression	7.133	0.506	0.709
Compression + Shuffled Resize	6.719	0.514	0.711

Table 4. Ablation study results of hand-drawn lines enhancement denoted as **Sharpen** and twin perceptual loss denoted as **APL**.

	NIQE↓	MANIQA↑	CLIQQA↑
Plain	7.351	0.501	0.689
Plain + Sharpen	7.182	0.504	0.707
Plain + Sharpen + APL	6.835	0.512	0.708
Plain + Sharpen + APL + Balanced Scale	6.719	0.514	0.711

facts, our model outperforms others in effective restoration, evidenced by the third and fourth rows. This is thanks to our improvement to the image degradation model where we provide a robust restoration capability on compression and resize functionality. Meanwhile, due to our proposed balanced twin perceptual loss, images generated by our GAN network do not show unwanted color artifacts as in AnimeSR and VQD-SR, which can be seen in the fifth row. Further, thanks to the versatile scenes collected in our proposed dataset, we are capable of achieving effective restoration in dark scenes. More visual results can be found in the supplementary materials.

4.3. Ablation Study

In this section, we conduct ablation studies to evaluate the substantial impact of our proposed dataset, degradation model, and hand-drawn lines enhancement with balanced twin perceptual loss. The inference dataset is still AVC-RealLQ [56]. Visual comparisons are presented in the supplementary materials.

Impact of the Dataset. As shown in Tab. 2, we substitute our API training dataset with several alternatives for comparative analysis: AVC-Train [56], frames randomly selected from the same video source as our API, a collection of I-Frames with IQA selection, and a collection of I-Frames with ICA selection. For a fair comparison, we keep a similar intensity of the training dataset size. If we take the AVC-Train video training dataset as an image dataset to train, we include temporal distorted images and less infor-

mative frames, which makes the performance hard to compete with the model trained with API in all metrics. Randomly selected image datasets perform poorly because they lack attention to high-quality frames in videos. With our I-Frame collection, we take off temporally distorted frames and choose the least compressed frames, but IQA-based selection limits the performance. With the same training iterations and conditions, the dataset selected by ICA-based criteria leads to an improvement over the dataset by IQA-based selection. With the 720P rescaling method, anime images have more compact hand-drawn lines and CGI information than falsely upscaled versions, and this back-to-original thinking boosts the performance in all metrics.

Degradation Model. As shown in Tab. 3, to validate the superiority of our degradation model, we replace our proposed degradation model with the high-order degradation model from the Real-ESRGAN [54] and random order degradation model from BSRGAN [65], which share certain similarity as our methods. Our degradation model with prediction-oriented compression model reaches an outstanding improvement in MANIQA [62] and CLIQQA [49] metrics. With our shuffled resize design, our network becomes more robust to versatile real-world SR scenarios and the performance can move one step further, especially the NIQE [34] metrics.

Benefits of proposed Enhancement and Perceptual Loss.

As shown in Tab. 4, we compare our model with the plain version that is not trained with proposed hand-drawn lines enhancement and balanced twin perceptual loss. The introduction of our hand-drawn lines enhancement presents a significant improvement on CLIQQA [49]. When we append ResNet perceptual loss in GAN training, it shows outstanding improvement in NIQE [34]. Further, with the proposed scaling on the early layers of the ResNet perceptual loss part, two perceptual losses have reached a stable balance and the performance moves one step further. This proves that a perceptual loss that is compatible with the anime domain is very insightful and instructive.

5. Conclusion

In this paper, we thoroughly utilize the characteristics of anime production knowledge and fully leverage it to enrich and enhance anime SR. To be specific, we propose a high-quality and informative anime production-oriented image (API) SR dataset with a novel dataset curation design. To restore and enhance hand-drawn lines, we propose an image degradation model to restore video compression artifacts and a pseudo-GT enhancement strategy. We further address unwanted color artifacts by introducing a network trained with high-level anime tasks to construct a balanced twin perceptual loss. Extensive experiment results demonstrate our superiority over existing SOTA methods, where we can restore harder real-world low-quality anime images.

References

- [1] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 126–135, 2017. **1**
- [2] Olivier Avaro, Alexandros Eleftheriadis, Carsten Herpel, Ganesh Rajan, and Liam Ward. Mpeg-4 systems: overview. *Signal Processing: Image Communication*, 15(4-5):281–298, 2000. **6, 2**
- [3] Matthew Baas. Danbooru2018 pretrained resnet models for pytorch. <https://rf5.github.io>, 2019. Accessed: DATE. **6, 3**
- [4] Gwern Branwen and Aaron Gokaslan. Danbooru2019: A large-scale crowdsourced and tagged anime illustration dataset. *Danbooru2017*, 2019. **6**
- [5] Yu Cao, Xiangqiao Meng, PY Mok, Xueting Liu, Tong-Yee Lee, and Ping Li. Animediffusion: Anime face line drawing colorization via diffusion models. *arXiv preprint arXiv:2303.11137*, 2023. **3, 1**
- [6] Hernan Carrillo, Michaël Clément, Aurélie Bugeau, and Edgar Simo-Serra. Diffusart: Enhancing line art colorization with conditional diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3485–3489, 2023. **2, 3, 1**
- [7] Kelvin CK Chan, Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Basicvsr: The search for essential components in video super-resolution and beyond. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4947–4956, 2021. **2**
- [8] Kelvin CK Chan, Shangchen Zhou, Xiangyu Xu, and Chen Change Loy. Investigating tradeoffs in real-world video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5962–5971, 2022. **6, 7, 2**
- [9] Chaofeng Chen and Jiadi Mo. IQA-PyTorch: Pytorch toolbox for image quality assessment. [Online]. Available: <https://github.com/chaofengc/IQA-PyTorch>, 2022. **7**
- [10] Shuhong Chen and Matthias Zwicker. Improving the perceptual quality of 2d animation interpolation. In *European Conference on Computer Vision*, pages 271–287. Springer, 2022. **1, 3**
- [11] Shuhong Chen, Kevin Zhang, Yichun Shi, Heng Wang, Yiheng Zhu, Guoxian Song, Sizhe An, Janus Kristjansson, Xiao Yang, and Matthias Zwicker. Panic-3d: Stylized single-view 3d reconstruction from portraits of anime characters. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21068–21077, 2023. **3**
- [12] Yuanzheng Ci, Xinzhu Ma, Zhihui Wang, Haojie Li, and Zhongxuan Luo. User-guided deep anime line art colorization with conditional adversarial networks. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1536–1544, 2018. **1**
- [13] Yuekun Dai, Shangchen Zhou, Qinyue Li, Chongyi Li, and Chen Change Loy. Learning inclusion matching for animation paint bucket colorization, 2024. **3**
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. **5**
- [15] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307, 2015. **2**
- [16] Tinglei Feng, Yingjie Zhai, Jufeng Yang, Jie Liang, Deng-Ping Fan, Jing Zhang, Ling Shao, and Dacheng Tao. Ic9600: A benchmark dataset for automatic image complexity assessment. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. **3**
- [17] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. **2**
- [18] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. **3**
- [19] Jingning Han, Bohan Li, Debargha Mukherjee, Ching-Han Chiang, Adrian Grange, Cheng Chen, Hui Su, Sarah Parker, Sai Deng, Urvang Joshi, et al. A technical overview of av1. *Proceedings of the IEEE*, 109(9):1435–1462, 2021. **6, 1, 2**
- [20] Ylies Hat, Gregor Jouet, Francis Rousseaux, and Clément Duhart. Paintstorch: a user-guided anime line art colorization tool with double generator conditional adversarial network. In *Proceedings of the 16th ACM SIGGRAPH European Conference on Visual Media Production*, pages 1–10, 2019. **1**
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 630–645. Springer, 2016. **6, 1, 3**
- [22] Zhengyu Huang, Haoran Xie, Tsukasa Fukusato, and Kazunori Miyata. Anifacedrawing: Anime portrait exploration during your sketching. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11, 2023. **3, 1**
- [23] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. **6**
- [24] Xiaozhong Ji, Yun Cao, Ying Tai, Chengjie Wang, Jilin Li, and Feiyue Huang. Real-world super-resolution via kernel estimation and noise injection. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2020. **2**
- [25] Xiaozhong Ji, Yun Cao, Ying Tai, Chengjie Wang, Jilin Li, and Feiyue Huang. Real-world super-resolution via kernel estimation and noise injection. In *proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 466–467, 2020. **6**
- [26] Yuxin Jiang, Liming Jiang, Shuai Yang, and Chen Change Loy. Scenimefy: Learning to craft anime scene via semi-supervised image-to-image translation. In *Proceedings of*

- the *IEEE/CVF International Conference on Computer Vision*, pages 7357–7367, 2023. 2, 3
- [27] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pages 694–711. Springer, 2016. 6, 1
- [28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [29] Yeongseop Lee and Seongjin Lee. Automatic colorization of anime style illustrations using a two-stage generator. *Applied Sciences*, 10(23):8699, 2020. 1
- [30] Yawei Li, Yuchen Fan, Xiaoyu Xiang, Denis Demandolx, Rakesh Ranjan, Radu Timofte, and Luc Van Gool. Efficient and explicit modelling of image hierarchies for image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18278–18289, 2023. 2, 6, 1
- [31] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1833–1844, 2021. 2
- [32] Joan L Mitchell, William B Pennebaker, Chad E Fogg, Didier J LeGall, Joan L Mitchell, William B Pennebaker, Chad E Fogg, and Didier J LeGall. Mpeg-2 overview. *MPEG Video Compression Standard*, pages 171–186, 1996. 6, 2
- [33] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. No-reference image quality assessment in the spatial domain. *IEEE Transactions on image processing*, 21(12):4695–4708, 2012. 3
- [34] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a “completely blind” image quality analyzer. *IEEE Signal processing letters*, 20(3):209–212, 2012. 7, 8
- [35] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018. 6
- [36] Heiko Schwarz, Detlev Marpe, and Thomas Wiegand. Overview of the scalable video coding extension of the h. 264/avc standard. *IEEE Transactions on circuits and systems for video technology*, 17(9):1103–1120, 2007. 3, 4, 6, 2
- [37] Wang Shen, Cheng Ming, Wenbo Bao, Guangtao Zhai, Li Chenn, and Zhiyong Gao. Enhanced deep animation video interpolation. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 31–35. IEEE, 2022. 1
- [38] Zhanjun Si and Ke Shen. Research on the webp image format. In *Advanced graphic communications, packaging technology and materials*, pages 271–277. Springer, 2016. 5, 6, 1, 2
- [39] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5, 3
- [40] Li Siyao, Shiyu Zhao, Weijiang Yu, Wenxiu Sun, Dimitris Metaxas, Chen Change Loy, and Ziwei Liu. Deep animation video interpolation in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6587–6595, 2021. 3
- [41] Li Siyao, Yuhang Li, Bo Li, Chao Dong, Ziwei Liu, and Chen Change Loy. Animerun: 2d animation visual correspondence from open source 3d movies. *Advances in Neural Information Processing Systems*, 35:18996–19007, 2022. 3
- [42] Li Siyao, Tianpei Gu, Weiye Xiao, Henghui Ding, Ziwei Liu, and Chen Change Loy. Deep geometrized cartoon line inbetweening. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7291–7300, 2023. 3
- [43] Shaolin Su, Qingsen Yan, Yu Zhu, Cheng Zhang, Xin Ge, Jinqiu Sun, and Yanning Zhang. Blindly assess image quality in the wild guided by a self-adaptive hyper network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3667–3676, 2020. 3
- [44] Gary J Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand. Overview of the high efficiency video coding (hevc) standard. *IEEE Transactions on circuits and systems for video technology*, 22(12):1649–1668, 2012. 3, 6, 2
- [45] Radu Timofte, Eirikur Agustsson, Luc Van Gool, Ming-Hsuan Yang, and Lei Zhang. Ntire 2017 challenge on single image super-resolution: Methods and results. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 114–125, 2017. 1
- [46] Zixi Tuo, Huan Yang, Jianlong Fu, Yujie Dun, and Xueming Qian. Learning data-driven vector-quantized degradation model for animation video super-resolution. *arXiv preprint arXiv:2303.09826*, 2023. 1, 2, 3, 6, 7
- [47] Gregory K Wallace. The jpeg still picture compression standard. *IEEE transactions on consumer electronics*, 38(1): xviii–xxxiv, 1992. 6, 2
- [48] Boyang Wang, Bowen Liu, Shiyu Liu, and Fengyu Yang. Vcizr: Blind single image super-resolution with video compression synthetic data. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4302–4312, 2024. 1, 2, 4, 5, 6
- [49] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring clip for assessing the look and feel of images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2555–2563, 2023. 7, 8
- [50] Ning Wang, Muyao Niu, Zhi Dou, Zhihui Wang, Zhiyong Wang, Zhaoyan Ming, Bin Liu, and Haojie Li. Coloring anime line art videos with transformation region enhancement network. *Pattern Recognition*, 141:109562, 2023. 3, 1
- [51] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018. 6
- [52] Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Recovering realistic texture in image super-resolution by deep spatial feature transform. In *Proceedings of the IEEE*

- conference on computer vision and pattern recognition, pages 606–615, 2018. 1
- [53] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European conference on computer vision (ECCV) workshops*, pages 0–0, 2018. 6, 1
- [54] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1905–1914, 2021. 1, 2, 4, 6, 7, 8
- [55] Holger Winnemöller, Jan Eric Kyprianidis, and Sven C Olsen. Xdog: An extended difference-of-gaussians compendium including advanced image stylization. *Computers & Graphics*, 36(6):740–753, 2012. 5, 1
- [56] Yanze Wu, Xintao Wang, Gen Li, and Ying Shan. Animesr: Learning real-world super-resolution models for animation videos. *arXiv preprint arXiv:2206.07038*, 2022. 1, 2, 3, 6, 7, 8, 4, 5
- [57] Zeyu Xiao, Zhiwei Xiong, Xueyang Fu, Dong Liu, and Zheng-Jun Zha. Space-time video super-resolution using temporal profiles. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 664–672, 2020. 2
- [58] Zeyu Xiao, Xueyang Fu, Jie Huang, Zhen Cheng, and Zhiwei Xiong. Space-time distillation for video super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2113–2122, 2021. 2
- [59] Zeyu Xiao, Jiawang Bai, Zhihe Lu, and Zhiwei Xiong. A dive into sam prior in image restoration. *arXiv preprint arXiv:2305.13620*, 2023. 2
- [60] Zeyu Xiao, Yutong Liu, Ruisheng Gao, and Zhiwei Xiong. Cutmib: Boosting light field super-resolution via multi-view image blending. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1672–1682, 2023. 2
- [61] Shizhuo Xu, Vibekananda Dutta, Xin He, and Takafumi Matsumaru. A transformer-based model for super-resolution of anime image. *Sensors*, 22(21):8126, 2022. 1
- [62] Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and Yujiu Yang. Maniqa: Multi-dimension attention network for no-reference image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1191–1200, 2022. 7, 8
- [63] Chih-Yuan Yao, Shih-Hsuan Hung, Guo-Wei Li, I-Yu Chen, Reza Adhitya, and Yu-Chi Lai. Manga vectorization and manipulation with procedural simple screentone. *IEEE transactions on visualization and computer graphics*, 23(2):1070–1084, 2016. 3
- [64] Yuan Yuan, Siyuan Liu, Jiawei Zhang, Yongbing Zhang, Chao Dong, and Liang Lin. Unsupervised image super-resolution using cycle-in-cycle generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 701–710, 2018. 2
- [65] Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4791–4800, 2021. 2, 6, 7, 8, 1
- [66] Lvmin Zhang, Chengze Li, Edgar Simo-Serra, Yi Ji, Tien-Tsin Wong, and Chunping Liu. User-guided line art flat filling with split filling mechanism. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9889–9898, 2021. 3
- [67] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 3
- [68] Song-Hai Zhang, Tao Chen, Yi-Fei Zhang, Shi-Min Hu, and Ralph R Martin. Vectorizing cartoon animations. *IEEE Transactions on Visualization and Computer Graphics*, 15(4):618–629, 2009. 3
- [69] Yang Zhao, Diya Ren, Yuan Chen, Wei Jia, Ronggang Wang, and Xiaoping Liu. Cartoon image processing: A survey. *IJCV*, 2022. 1

APISR: Anime Production Inspired Real-World Anime Super-Resolution

Supplementary Material

In this supplementary material, Sec. A first presents more statistics and details of our proposed anime image SR training dataset. Then, Sec. B shows details about our implementations in super-resolution (SR) network training. Specifically, Sec. B.1 presents the image SR network we used in our training. Sec. B.2 presents details of post-processing techniques we use on the pseudo-GT preparation for hand-drawn line enhancement. Sec. B.3 presents figures and details of the ResNet50 [21] perceptual loss for our proposed balanced twin perceptual loss. Sec. B.4 provides the hyperparameter setting for our proposed prediction-oriented compression and shuffled resize module in the degradation model. Finally, Sec. C provides more visual results of comparisons among SOTA methods and ablation studies.

A. API Dataset Details

Our Anime Production-oriented Image (API) SR dataset contains 3,740 high-quality and informative images. This quantity is roughly the same quantity as the previous photorealistic SR training dataset size [54, 65], which includes DIV2K [1], Flickr2K [45], and OutdoorSceneTraining [52]. The aspect ratio and resolution information before scaling are shown in Fig. 11.

B. Implementation Details

B.1. Training Network Details

The generator network we deploy is GRL [30], a SOTA image SR network (CVPR 2023). GRL leverages interconnected relationships within various layers of image structures through a Transformer-based framework, attaining improvement in multiple tasks of SR and image restoration. The model we chose is its tiny version, which has 0.91M parameters. To better adapt the real-world SR task, we changed its upsampler module from the default pixel shuffle strategy to the nearest neighbor interpolation with the convolution layer approach, which is used for the base model version but not for the tiny version in their proposed methods. We change the upsampler because the nearest neighbor interpolation with the convolution layer is claimed to show fewer artifacts in the upsampling process than the pixel shuffle strategy. The final network parameter is 1.03M, which is the smallest network among all image and video-based SOTA methods that we compare.

B.2. Hand-drawn line enhancement Details

In the hand-drawn line enhancement, we have proposed outlier filter and passive dilate techniques to obtain a clean XDoG-extracted [55] hand-drawn line edge map. XDoG is widely used in paired dataset preparation in anime colorization [5, 6, 22, 50]. The extracted edge map by XDoG is a binary output, where the white pixel stands for the active edge map region and the black pixel stands for the unrelated region.

For the outlier filter, we use breadth-first search in eight directions to recursively detect the surrounding pixels of all white pixels and turn white pixel regions into black pixels if the total quantity of connected white pixels is less than the threshold. We empirically set the threshold as 32.

For the dilation, we passively replace the black pixel with the white pixel if it has more than 3 white pixel neighbors, which is different from independent kernel-based active dilation methods in [12, 20, 29] that directly spread the surrounding neighbors to be white pixels if the central pixel is white. Compared to active dilation methods, our proposed passive dilation is more concentrated on the hand-drawn lines region instead of covering unrelated pixel information (see Fig. 13). Thus, we name our methods as passive dilata-tion.

In the implementation, we will do an unsharp mask for the whole image first to increase overall visualization sharpness and then apply two extra turns of sharpening to the hand-drawn lines specifically based on the pipeline design mentioned above. More implementation details can be found in our released code.

B.3. Balanced Twin Perceptual Loss Details

As shown in Fig. 12, our proposed middle-layer output comparisons for ResNet50 [21] follow the idea proposed by ESRGAN [53] which compares feature map outputs before the activation layer. Following VGG-based perceptual loss [27], we compare the last convolution layer of each stage. There are five middle-layer output comparisons, which are the same quantity as VGG-based perceptual loss [27]. Thus, our proposed twin perpetual loss reaches a mutual balance in training.

B.4. Degradation Details

For the prediction-oriented compression module of the degradation model, we deploy both the image compression with prediction mechanism (*i.e.*, WebP [38] and AVIF [19]) and single-frame video compression. Meanwhile, for the robustness of the degradation model, we keep the

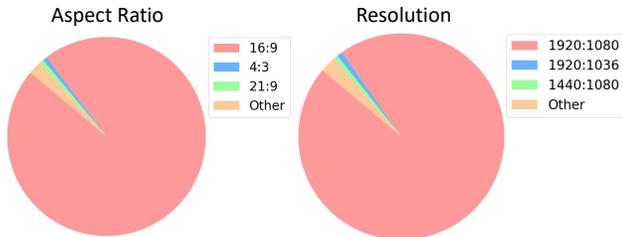


Figure 11. API dataset extra statistics.

JPEG [47]. The quality factor range of JPEG, WebP, and AVIF is [20, 95] with encoding speed in the range of [0, 6] for WebP and AVIF. The probability of fetching the value in the range is equal.

For the stability of video compression processing, we choose the widely-used video processing tools, *ffmpeg*, to perform the proposed single-frame compression of MPEG2 [32], MPEG4 [2], H.264 [36], and H.265 [44]. In *ffmpeg*, CRF is an engineering system to control the quantization level, and preset is a speed control mechanism whose setting is directly related to compression distortions. For MPEG2 and MPEG4, we empirically find that the quality factor control (*-qscale:v*) is a better way to control single-frame compression, but for H.264 and H.265, CRF is a better way to control. For MPEG2 and MPEG4, we set the quality factor in the range [8, 31]. For H.264 and H.265, we set the CRF in the range [23, 38] and [28, 42] respectively. The preset for all of them is $\{slow, medium, fast, faster, super\ fast\}$ with probability $\{0.05, 0.35, 0.3, 0.2, 0.1\}$.

The first prediction-oriented compression includes JPEG [47] and WebP [38] with a probability of $\{0.4, 0.6\}$ respectively. The second prediction-oriented compression includes JPEG [47], WebP [38], AVIF [19], and single-frame compression of MPEG2 [32], MPEG4 [2], H.264 [36], and H.265 [44] with probability of $\{0.06, 0.1, 0.1, 0.12, 0.12, 0.3, 0.2\}$ respectively. For the first resize module, we set the scaling in the range of [0.1, 1.2] with probability $\{0.2, 0.7, 0.1\}$ to scale up, scale down, or remain current resolution. For the second resize module, we choose the range of [0.15, 1.2] with probability $\{0.2, 0.7, 0.1\}$. More implementation details can be found in our released code.

C. More Qualitative Comparisons

In this section, we present more qualitative results to verify the effectiveness of our APISR among SOTA methods. Moreover, we provide visual comparisons for the ablation studies.

Extra Qualitative Comparisons with SOTA methods.

Fig. 14 and Fig. 15 show extra qualitative comparisons on AVC-RealLQ [56] datasets for $4\times$ scaling. This includes

image-based Real-ESRGAN [54] and BSRGAN [65], and video-based RealBasicVSR [8], AnimeSR [56], and VQD-SR [46]. Our APISR presents clearer and sharper hand-drawn lines (first example of Fig. 14, first and second examples of Fig. 15, and third example of Fig. 16), better restoration with more natural details (second and third examples of Fig. 14, and third example of Fig. 15), and does not present unwanted color artifacts (first and second examples of Fig. 16).

Qualitative Comparisons of Ablation Studies. Fig. 17, Fig. 18, and Fig. 19 shows the qualitative comparisons of ablations studies.

As shown in Fig. 17, the network trained with AVC-Train [56] over-sharpens the grid texture and produces annoying artifacts as denoted by the arrows in the figure. Similarly, the network trained with the random sampled or IQA-based sampled dataset can alleviate this artifact but is still hard to completely remove it. However, when we introduce the ICA-based selection method with I-Frame dataset collection, this artifact is greatly removed and the generated image shows more natural details. This is thanks to versatile complex scenes included in the dataset due to ICA-based selection. With 720P rescaling, fewer ringing artifacts appear.

As shown in Fig. 18, the network trained with high-order [54] and random order [65] degradation model presents ringing artifacts, rainbow effects, and color distortions as denoted by the arrows in the figure. Nevertheless, introducing our proposed prediction-oriented compression module in the degradation model promotes the network to greatly restore these problems and generate more natural details with less distorted hand-drawn lines. Moreover, with the shuffled resize module in the degradation model, more distortions are restored and present natural shadow details.

As shown in Fig. 19, the network trained with the plain version presents unwanted color pixel artifacts and sparse hand-drawn line information as denoted by the arrows in the figure. With the hand-drawn line enhancement, the hand-drawn line around the eyes of the character is greatly intensified and more details are generated. However, the unwanted color pixels still exist and they are presented as an annoying artifact. With the twin perceptual loss, the unwanted color pixels are greatly alleviated. Further, with the scaling to early layers in ResNet perceptual loss, more shadow artifacts and distortions are restored.

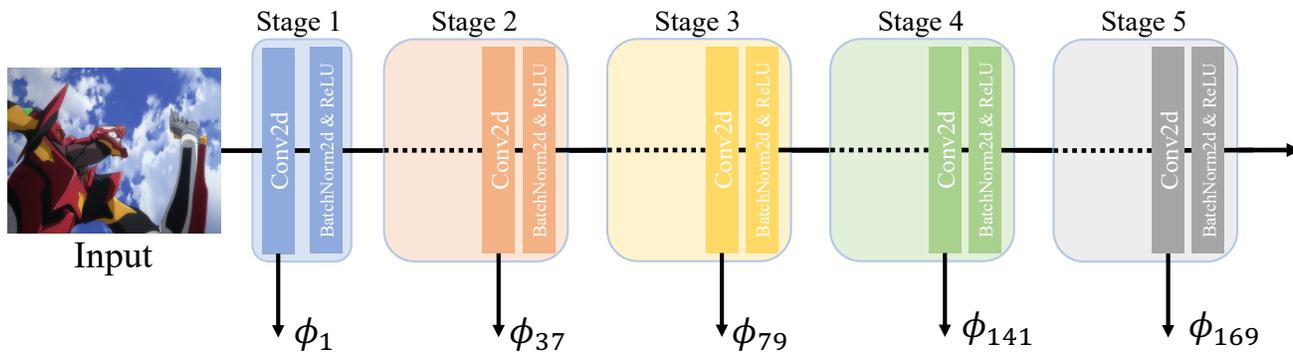


Figure 12. The overview of our proposed middle-layer outputs of ResNet50 [21] perceptual loss trained by Danbooru dataset [3]. Overall, ResNet50 can be summarized into five stages which is similar to VGG [39]. ϕ_j represents the perceptual function that returns j th layer output of ResNet50.



Figure 13. **Comparisons between active and passive dilation.** Our proposed passive dilation is more concentrated on the hand-drawn line region without producing over-sharpened pseudo-GT images as in active dilation methods.

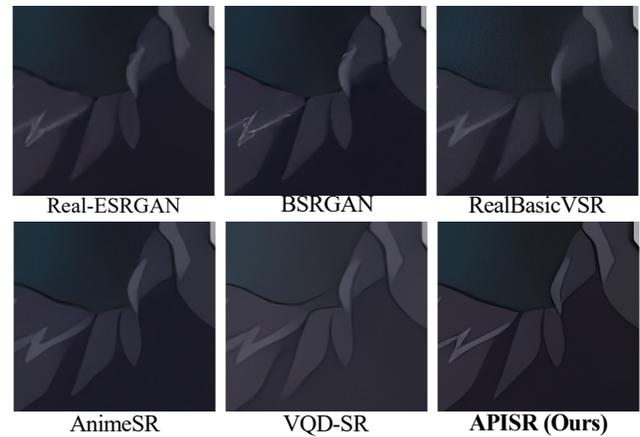
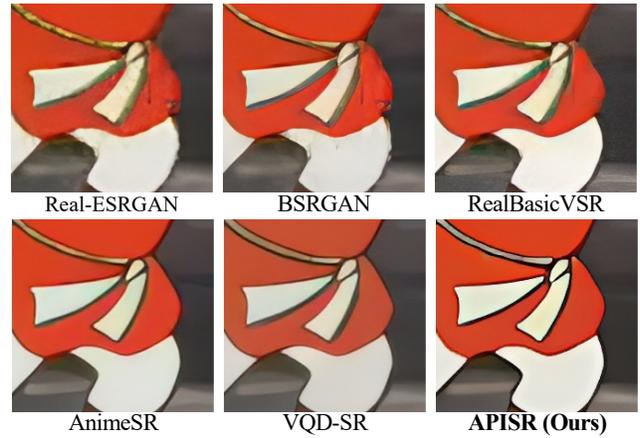


Figure 14. Qualitative comparisons on AVC-RealLQ [56] for 4× scaling. Our APISR presents clearer and sharper hand-drawn lines, better restoration with more natural details, and does not present unwanted color artifacts. **Zoom in for the best view.**

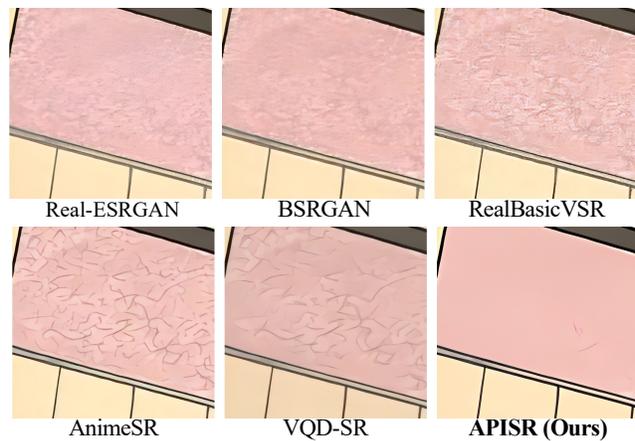
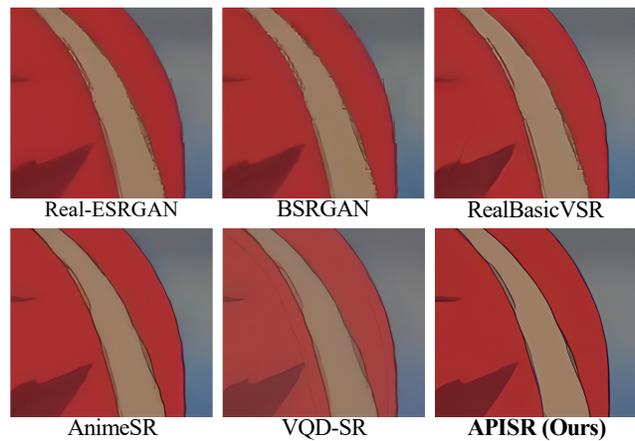
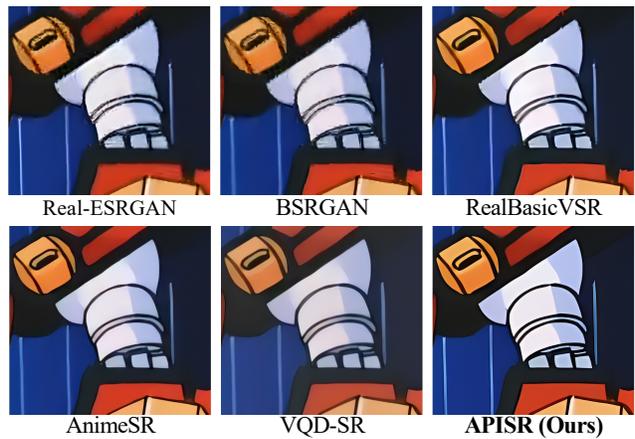
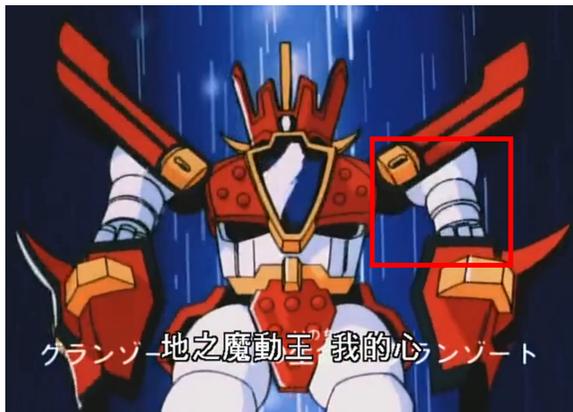


Figure 15. Qualitative comparisons on AVC-ReallQ [56] for 4× scaling. Our APISR presents clearer and sharper hand-drawn lines, better restoration with more natural details, and does not present unwanted color artifacts. **Zoom in for the best view.**

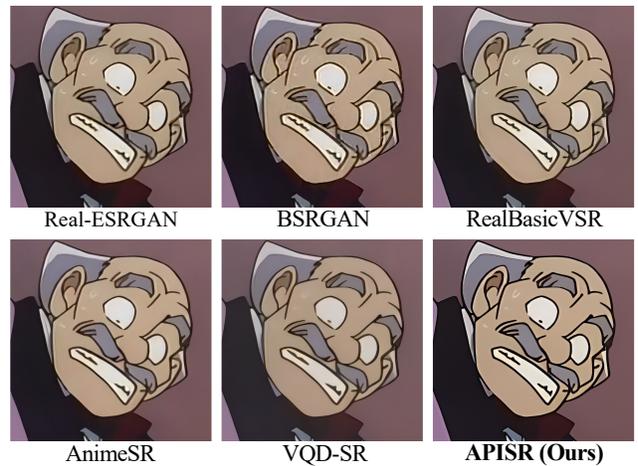
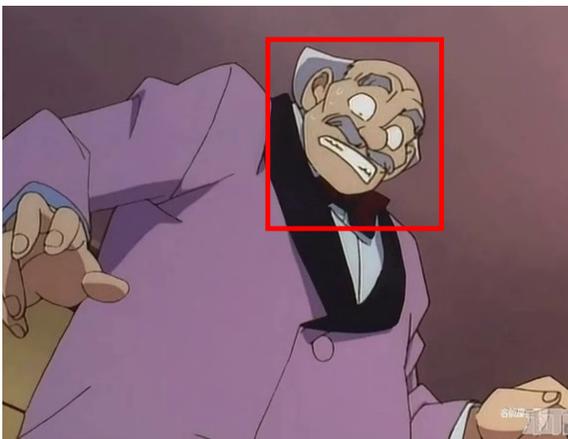
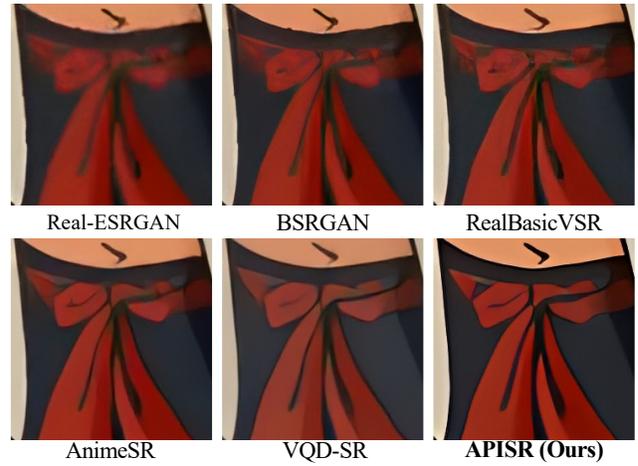
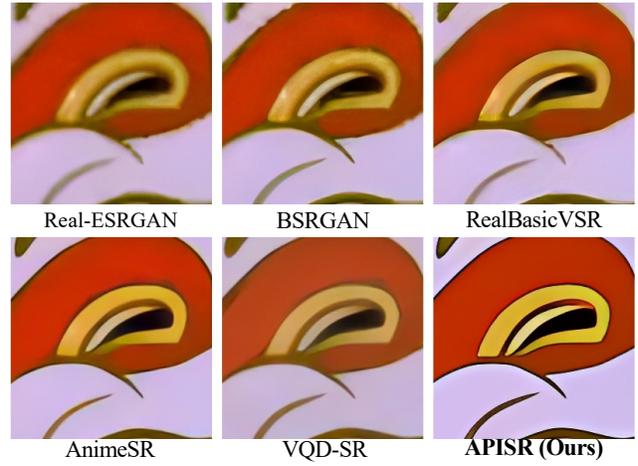


Figure 16. Qualitative comparisons on AVC-RealLQ [56] for $4\times$ scaling. Our APISR presents clearer and sharper hand-drawn lines, better restoration with more natural details, and does not present unwanted color artifacts. **Zoom in for the best view.**

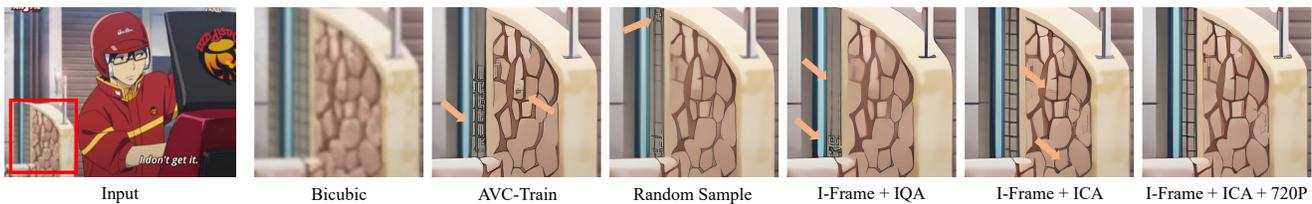


Figure 17. **Qualitative comparisons of the first ablation study.** IQA stands for image quality assessment. ICA stands for image complexity assessment. 720P stands for our proposed 720P rescaling. **Zoom in for the best view.**

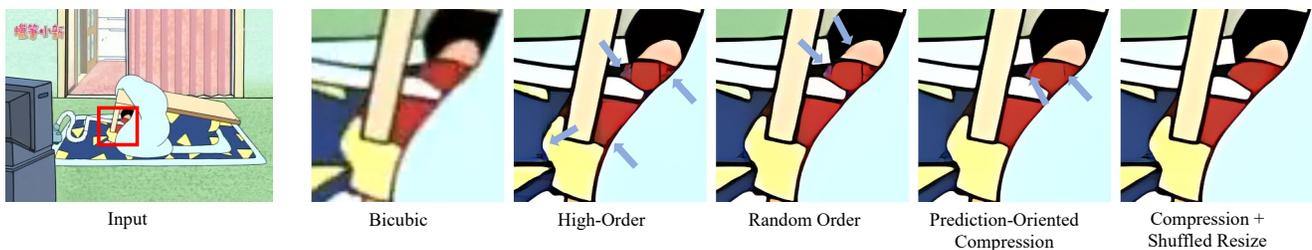


Figure 18. **Qualitative comparisons of the second ablation study.** **Zoom in for the best view.**

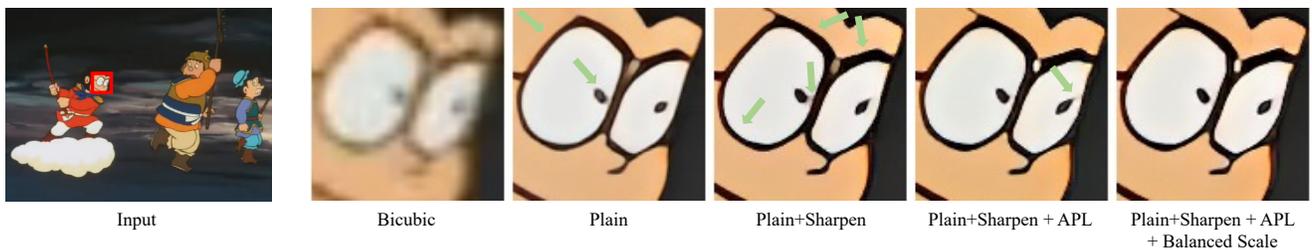


Figure 19. **Qualitative comparisons of the third ablation study.** Hand-drawn lines enhancement is denoted as **Sharpen** and twin perceptual loss is denoted as **APL**. **Balanced Scale** presents the early layer scaling to ResNet perceptual loss. **Zoom in for the best view.**