# PAniC-3D: Stylized Single-view 3D Reconstruction from Portraits of Anime Characters

Shuhong Chen[*†]    Kevin Zhang[*]    Yichun Shi[†]    Heng Wang[†]    Yiheng Zhu[†]

Guoxian Song[†]    Sizhe An[†]    Janus Kristjansson[*]    Xiao Yang[†]    Matthias Zwicker[*]

University of Maryland - College Park, MD, USA[*]
ByteDance[†]

## Abstract

*We propose PAniC-3D, a system to reconstruct stylized 3D character heads directly from illustrated (p)ortraits of (ani)me (c)haracters. Our anime-style domain poses unique challenges to single-view reconstruction; compared to natural images of human heads, character portrait illustrations have hair and accessories with more complex and diverse geometry, and are shaded with non-photorealistic contour lines. In addition, there is a lack of both 3D model and portrait illustration data suitable to train and evaluate this ambiguous stylized reconstruction task. Facing these challenges, our proposed PAniC-3D architecture crosses the illustration-to-3D domain gap with a line-filling model, and represents sophisticated geometries with a volumetric radiance field. We train our system with two large new datasets (11.2k Vroid 3D models, 1k Vtuber portrait illustrations), and evaluate on a novel AnimeRecon benchmark of illustration-to-3D pairs. PAniC-3D significantly outperforms baseline methods, and provides data to establish the task of stylized reconstruction from portrait illustrations.*

## 1. Introduction & Related Work

With the rise of AR/VR applications, there is increased demand for not only high-fidelity human avatars, but also non-photorealistic 3D characters, especially in the "anime" style. Most character designers typically create concept illustrations first, allowing them to express complex and highly diverse characteristics like hair, accessories, eyes, skins, headshapes, etc. Unfortunately, the process of developing illustrated concept art into an AR/VR-ready 3D asset is expensive, requiring professional 3D artists trained to use expert modeling software. While template-based creators democratize 3D avatars to an extent, they are often restricted to 3D assets compatible with a specific body model.

We propose PAniC-3D, a system to automatically recon-

struct a stylized 3D character head directly from illustrated (p)ortraits of (ani)me (c)haracters. We formulate our problem in two parts: 1) implicit single-view head reconstruction, 2) from across an illustration-3D domain gap. To summarize our contributions:

- **PAniC-3D**: a system to reconstruct the 3D radiance field of a stylized character head from a single line-based portrait illustration.

- The **Vroid 3D dataset** of 11.2k character models and renders, the first such dataset in the anime-style domain to provide 3D assets with multiview renders.

- The **Vtuber dataset** of 1.0k reconstruction-friendly portraits (aligned, front-facing, neutral-expression) that bridges the illustration-render domain gap through the novel task of line removal from drawings.

- The **AnimeRecon benchmark** with 68 pairs of aligned 3D models and corresponding illustrations, enabling quantitative evaluation of both image and geometry metrics for stylized reconstruction.

### 1.1. Implicit 3D Reconstruction

While there has been much work on mesh-based reconstruction from images [23], these systems are not expressive enough to capture the extreme complexity and diversity of topology of our 3D characters. Inspired by the recent successes in generating high-quality 3D radiance fields [4, 5, 25, 39], we instead turn to implicit representations. However, to achieve high-quality results, recent implicit reconstruction work such as PixelNerf [40] tend to operate solely from 2D images, due to the lack of publicly-available high-quality 3D data. Some implicit reconstruction systems like Pifu [31] employing complex 3D assets have shown reasonable success using point-based supervision, but require careful point sampling techniques and loss balancing.
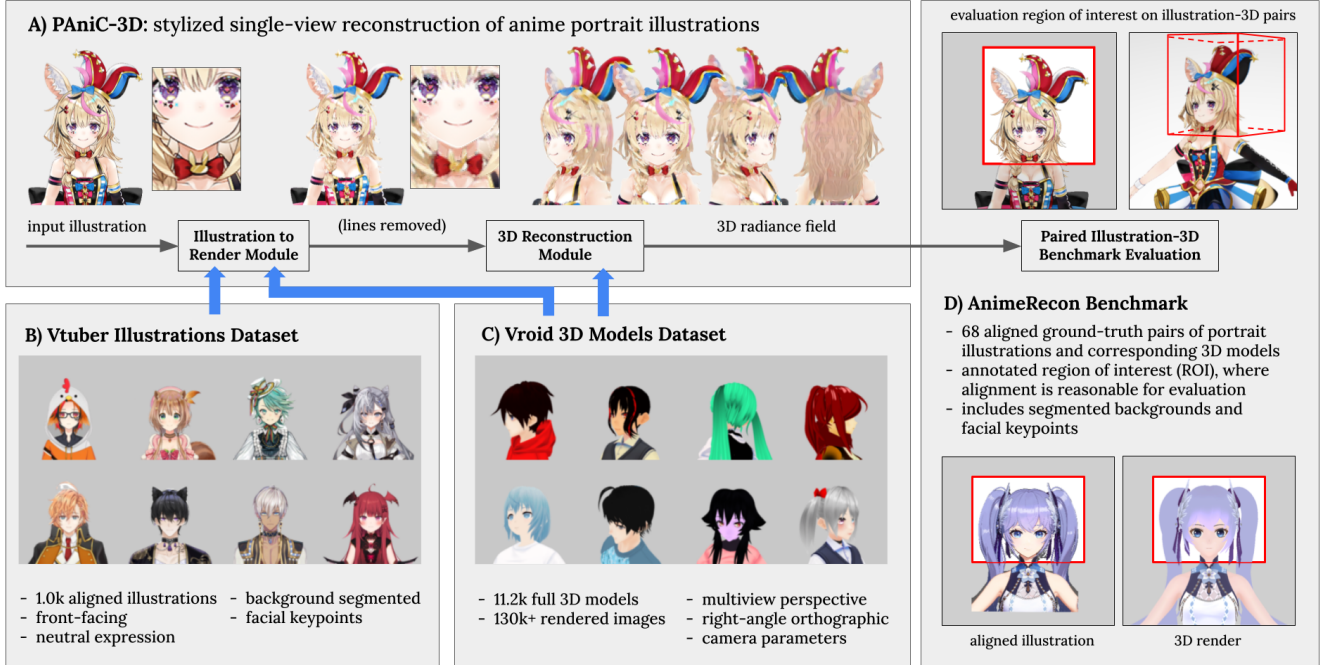
Figure 1. Overview of contributions. Our (A) PAniC-3D system is able to reconstruct a 3D radiance field directly from a line-based portrait illustration. We gather a new (B) Vtuber illustration dataset and (C) Vroid 3D models dataset in order to cross the illustration-render domain gap and supervise reconstruction. To evaluate, we provide a new (D) AnimeRecon benchmark of paired illustrations and 3D models, establishing the novel task of stylized single-view reconstruction of anime characters. (Art attributions in suppl.)

There is also a body of work on sketch-based modeling, where 3D representations are recovered from contour images. For example, Rui et al. [24] use a multi-view decoder to predict sketch-to-depth and normals, which are then used for surface reconstruction. Song et al. [44] additionally try to compensate multi-view drawing discrepancies by learning to realign the inputs. While related to our single-view portrait reconstruction problem, these methods require multi-view sketches that are difficult for character artists to draw consistently, and cannot handle color input.

For our case with complex high-quality 3D assets, we demonstrate the superiority of differentiable volumetric rendering for reconstruction. We build off of recent unconditional generative work (EG3D [4]), formulating the problem of reconstruction as a conditional generation, proposing several architecture improvements, and applying direct 2.5D supervision signal as afforded by our 3D dataset.

## 1.2. Anime-style 3D Avatars and Illustrations

It is a fairly common task for 3D character artists to produce a 3D model from a portrait illustration; however from the computer graphics standpoint, this stylized reconstruction setup adds additional ambiguity to an already ill-posed problem. In addition, while there's work in the popular anime/manga domain using 3D character assets (for pose estimation [18], re-targetting [17, 20], and reposing [22],

etc.), there's a lack of publicly-available 3D character assets with multi-view renders that allow scalable training (Tab. 1). In light of these issues, we propose AnimeRecon (Fig. 1d) to formalize the stylization task with a paired illustration-to-3D benchmark, and provide the Vroid dataset (Fig. 1c) of 3D assets to enable large-scale training.

Within the problem of stylized reconstruction, we solve the task of contour removal from illustrations. There is much work on line extraction [21, 38], sketch simplification [33, 34], reconstruction from lines [11, 24], line exploits for artistic imagery [6, 41], and scratch-line removal [8, 27, 29, 32, 35]; however, the removal of lines from line-based illustration has seen little focus. We examine this contour deletion task in the context of adapting drawings to render-like images more conducive to 3D reconstruction; we find that naive image-to-image translation [19,45] is unsuited to the task, and propose a simple yet effective adversarial training setup with facial feature awareness. Lastly, we provide a Vtuber dataset (Fig. 1b) of portraits to train and evaluate contour removal for 3D reconstruction.

## 2. Methodology

PAniC-3D is composed of two major components (Fig. 1a): a 3D reconstructor directly supervised to predict a radiance field from a given front render, and an illustration-

| 3D datasets | Vroid | AniRec. | AC | CoNR | ADD |
|---|---|---|---|---|---|
| 3D avail. | Y | Y | - | - | - |
| renders avail. | Y | Y | Y | - | - |
| multiview | Y | Y | - | Y | Y |
| paired illustr. | - | Y | - | - | - |

Table 1. 3D anime datasets comparison. Our new Vroid dataset is the first to make 3D anime models and multiview renders available. The AniRecon benchmark allows quantitative evaluation of both 2D image and 3D geometry metrics. Others left-to-right: AnimeCeleb [20], CoNR [22], Anime Drawings Dataset [18].

| Portraits | Vtuber | AC | AP | iCF | BP |
|---|---|---|---|---|---|
| illustrations | Y | - | Y | Y | Y |
| aligned face | Y | Y | Y | Y | - |
| front-facing | Y | Y | - | - | - |
| neutral expr. | Y | Y | - | - | - |
| segmented | Y | Y | - | - | Y |
| face kpts. | Y | - | - | - | Y |

Table 2. Anime image datasets comparison. Our new Vtuber dataset allows us to examine the domain gap between line-based illustrations and 3D renders, and is filtered/standardized specifically for characteristics desirable in 3D head reconstruction. Others left-to-right: AnimeCeleb [20], AnimePortraits [3], iCartoonFace [43], BizarrePose [7].

to-render module that translates images to the reconstructor's training distribution. The two parts are trained independently, but are used sequentially at inference.

## 2.1. 3D Reconstruction Module

The 3D reconstruction module Fig. 3 is trained with direct supervision to invert a frontal render into a volumetric radiance field. We build off of recent unconditional generative work (EG3D [4]), formulating the problem of reconstruction as that of conditional generation, proposing several architecture improvements, and applying direct 2.5D supervision signal as afforded by our 3D dataset.

**Conditional input:** The given front orthographic view to reconstruct is resized and appended to the intermediate feature maps of the Stylegan2 backbone used in EG3D [4]. In addition, at the earliest feature map we give the model high-level domain-specific semantic information about the input by concatenating the penultimate features of a pre-trained Resnet-50 anime tagger. The tagger provides high-level semantic features appropriate for conditioning the generator; it was pretrained by prior work [7] on 1062 relevant classes such as blue_hair, cat_ears, twin_braids, etc.

**Feature pooling:** As the spatial feature maps are to be reshaped into a 3D triplane as in EG3D [4], we found it beneficial to pool a fraction of each feature map's channels

along the image axes (see Fig. 3 left). This simple technique helps distribute information along common triplane axes, improving performance on geometry metrics.

**Multi-layer triplane:** As proposed in concurrent work [1], we improve the EG3D triplane by stacking more channels along each plane (see Fig. 3 center). The method may be interpreted as a hybrid between a triplane and a voxel grid (they are equivalent if the number of layers equals the spatial size). Setting three layers per plane allows better spatial disambiguation when bilinearly sampling the volume, and particularly helps our model generate more plausible backs of heads (a challenge not faced by EG3D).

**Losses:** We take full advantage of the ground-truth 2.5D representations afforded to us by our available 3D assets. Our reconstruction losses include: RGB $L_1$, LPIPS [42], silhouette $L_1$, and depth $L_2$; these are applied to the front, back, right, and left orthographic views, as shown in Fig. 3. A discriminative loss is applied to improve the detail quality, in addition to maintaining the generation orientation. We also keep the R1 and density regularization losses from EG3D training. Our 2.5D representations and adversarial setup allow us to surpass similar single-view reconstructors such as PixelNerf [40] which work only with color losses.

**Post-processing:** We leverage our assumption that front-orthographic views are given as input, by stitching the given input onto the generated radiance field at inference. The xy-coordinates of each pixel's intersection within the volume are used to sample the input as a uv-texture map; we cast few additional rays from each intersection to test for visibility from the front, and apply the retexturing accordingly. This simple yet effective method improves detail preservation from the input at negligible cost.

## 2.2. Illustration-to-Render Module

In order to remove non-realistic contour lines present in the input illustration, but absent in a diffusely-lit radiance field, we design an illustration-to-render module (Fig. 4). Assuming access to unpaired illustrations and renders (our Vtuber and Vroid datasets, respectively), the shallow network re-generates pixel colors near lines in the drawing in order to adversarially match the render image distribution.

Similar to unpaired image-to-image models like Cycle-GAN and UGATIT [19, 45], we also impose a small identity loss; while this may seem counter-productive for our infilling case where identity is preserved in non-generated regions, we found that this stabilizes the GAN training. Note that our setup also differs from other infilling models, in that we inpaint to match a distribution different from the input.

Following prior work extracting sketches from line-based animations [6], we use the simple difference of gaussians (DoG) operator in order to prevent double-line extraction around each stroke.

While most lines present in the drawing should be re-

Figure 2. (a) No-line diffuse render, (b) real-lined illustration, (c) Blender Freestyle [14], (d) RTSC suggestive contours [9]. Toon shaders over-draw (c, cheeks) or miss lines (d, bowtie); it is non-trivial to model the artistic line placement in real drawings (b). Thus, we train our Illustration-to-Render module on real lines drawn by artists instead of synthetically-added lines.

moved, certain lines around key facial features must be preserved as they indeed appear in renderings (eyes, mouth, nose, etc.). We employ an off-the-shelf anime facial landmark detector [16] to create convex hulls around critical structures, where infilling is disallowed.

We show that this line removal module indeed achieves a more render-like look; it performs image translation more accurately than baseline methods when evaluated over our AnimeRecon pairs (Tab. 4), and removes line artifacts from the ultimate radiance field renders (Fig. 6).

## 3. Data

Unless otherwise mentioned, we use 80-10-10 splits for training, validation, and testing.

### 3.1. Vroid 3D Dataset

We collect a large dataset of 11.2k 3D anime characters from VroidHub, in order to train both the reconstruction and image-to-image translation modules of PAniC-3D. Unlike previous work in the 3D anime domain using MikuMiku-Dance PMD/PMX models [17, 18, 20], Vroid VRM models conform to the GLTF2 standard [15] with several extensions [37], allowing us to render using ModernGL [10,12]. As the data is crowd-sourced, we filter against a variety of undesirable properties, such as: texture corruption, too much or too little transparency, extreme character sizes, missing bones, etc. The 11.2k renders we use represent 70% retension of our original 16.0k scraped models.

All our image data are rendered with unit ambient lighting (*e.g.* using diffuse surface color only), unit distance away from the "neck" bone common to VRMs [37]. Our choice to model only diffuse lighting is motivated by the frequent artistic decision to paint speculars as textures (Tab. 4 row 2, hair and eye highlights on diffuse renders). Adding specular renderering to painted highlights would introduce inconsistent lighting effects. Linear blend skinning is used to lower characters' arms 60-degrees from their resting T-pose position. We supersample at 1024px resolution, before bilinear downsampling to 512px for training and testing.

The dataset for PAniC-3D's reconstruction module consists of both random perspective views for adversarial training, and fixed orthographic views for the input and reconstruction losses. Each 3D model is rendered from 8 random perspective views (with uniformly-sampled 360-degree azimuth, normally-sampled elevation of 20-degree standard deviation, and fixed 30-degree full field-of-view); this yields a total of 89.6k perspective images, each with known camera parameters. The four orthographic views are taken at fixed 90-degree angles from the front, sides, and back. The unpaired 3D data for our image-to-image translation module is simply the front orthographic view of each character identity.

### 3.2. Vtuber Portraits Dataset

We gather 1004 portraits from the VirtualYoutuber Fandom Wiki as the unpaired illustration dataset for our image-to-image module. These portraits were manually filtered from 15.2k scraped images for desired properties, *e.g.* high-resolution, front-facing, neutral-expression, uncropped, etc. In order to mimic the image distribution of the orthographic front-view Vroid renders, we white out the backgrounds using the character segmenter from Chen *et al.* [7] and select the largest connected component. In addition, we align the facial keypoints of each illustration (extracted with a pre-trained YOLOv3 [16,28]) to match the height and scale distributions of keypoint detections on the Vroid renders.

While we could add lines artificially, we found toon shaders [9, 14] to poorly model artistic judgements needed to place lines on intricate characters (Fig. 2). To avoid these artifacts, we designed a data-driven Illustr2Render module to train on real lines drawn by artists.

### 3.3. AnimeRecon Benchmark

We collect a benchmark set for stylized reconstruction by finding 68 characters with both 3D models and closely-aligned illustrations. Specifically, we source from the 3D mobile game Genshin Impact (for which we can match character portraits from the Fandom Wiki) and the virtual talent agency Hololive (from which several members have both 3D avatars and a Fandom Wiki portrait). As the raw 3D data from both sources comes in MMD format, we convert to VRM using the DanSingSing converter; the rendering process is the same as that of Vroid.

The portraits are aligned to their corresponding front-view orthographic renders by a manually-decided mixture of YOLOv3 facial keypoints [16, 28] and ORB detections. We segment out backgrounds using the same model [7] and procedure as with Vtuber portraits. In order to maintain separation from training data, we remove all Hololive identities present in this benchmark set from the Vtuber portraits set.

Inevitably, the illustrations do not all align perfectly with their corresponding renders, and many Genshin illustrations

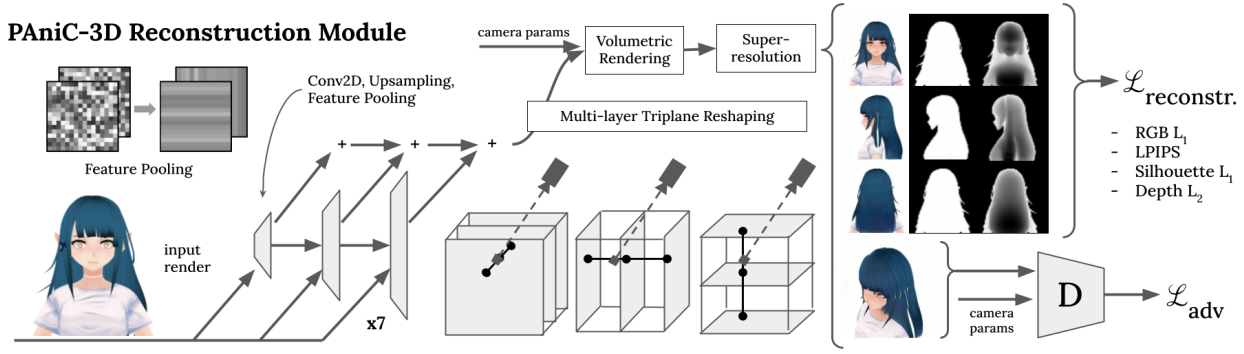**PAniC-3D Reconstruction Module**

Figure 3. Schematic of the 3D reconstruction module. The front-orthographic input rendering is fed to a series of upsampling convolutions, with intermediate feature pooling to help distribute information along common triplane axes. The final feature stack is reshaped into a multi-layer triplane, which is volumetrically rendered and super-resolutioned to the final output. Reconstruction losses are applied to front, left, right, and back views (right view omitted in figure), and adversarial loss is applied to a random perspective view. (Art attributions in suppl.)



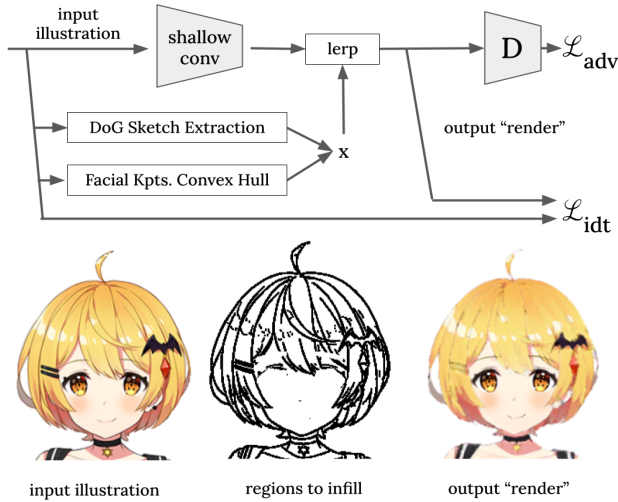**PAniC-3D Illustration-to-Render Module**

Figure 4. Schematic of the illustration-to-render module. We design a simple yet effective network to cross the domain gap by removing illustration contour lines absent in a diffuse 3D render, while retaining lines present around facial features like the eyes and mouth. (Art attributions in suppl.)

are cropped when aligned. In order to perform more reasonable evaluation, we additionally label a rectangular region of interest (ROI) for each aligned pair, within which the alignment is appropriate; all 2D image and 3D geometry metrics reported are restricted to the ROI when calculated.

# 4. Results & Evaluation

In this section, we provide a breakdown of reconstruction and image-to-image translation performance, with both qualitative and quantitative comparison to other baselines.

## 4.1. 3D Reconstruction Results

### 4.1.1 Metrics

As shown in Tab. 3, we evaluate over our new AnimeRecon benchmark using both 2D image and 3D geometry metrics. All the illustration inputs went through our illustration-to-render module before being fed to the respective method; we believe this is a fairer comparison, as all the reconstructors were trained on renders.

Image metrics are measured by comparing the predicted radiance field's integrated image with the ground-truth render from the same camera viewpoint. We show such measures for the front orthographic view (which should match the input), the back orthographic view, and an average of 12 perspective cameras circling the character at 30-degree intervals. For the front and back views, we restrict evaluation to our AnimeRecon ROI (Fig. 1d); for the 12 circling views, we crop to the horizontal bounding box strip.

We show standard color metrics like PSNR for completeness, but as there are inevitably imperfections on the AnimeRecon illustration-render pairs (even within the ROI), perceptual metrics like CLIP image cosine similarity [26] and LPIPS [42] are generally more relevant to quality.

The geometry metrics are on the right of Tab. 3. We extract meshes from both the ground-truth 3D asset and the predicted radiance field (through marching cubes), and delete faces with vertices outside the rectangular prism defined by the ROI annotation (see Fig. 1d top). We then sample 10k points randomly from each mesh subset, to compute the point-cloud chamfer distance and F-1 scores at 5cm and 10cm [13]. To put F-1 in perspective, the average Vroid head width is 25.5cm (real heads are 14cm). The F-1 are low overall, due to larger proportions and protruding hair/accessories with large surface area.

Figure 5. Qualitative comparison of baselines. Illustrations with lines removed by our illustration-to-render module are fed to various reconstruction frameworks; our PAniC-3D system delivers plausible reconstructions, while other methods struggle to preserve identity and predict reasonable geometry. Note that metrics (Tab. 3) between the displayed ground-truth and prediction are restricted to an ROI bounding box/rectangular prism (unshown) during evaluation. (Art attributions in suppl.)

| | front | | | back | | | 360 | | | geom. | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CLIP | LPIPS | PSNR | CLIP | LPIPS | PSNR | CLIP | LPIPS | PSNR | CD | F-1@5 | F-1@10 |
| PAniC-3D (full) | **94.66** | **19.37** | 16.91 | **85.08** | **30.02** | 15.51 | **84.68** | 25.25 | **15.98** | **1.33** | **37.73** | 65.50 |
| no feature pooling | **94.64** | **19.26** | **16.95** | 84.06 | **30.23** | **15.53** | 84.30 | **25.42** | 15.96 | 1.38 | 35.26 | 65.51 |
| no multi-layer triplane | 94.39 | 19.99 | 16.94 | 84.28 | 30.37 | 15.41 | 84.49 | 26.13 | 15.82 | 1.44 | 34.55 | **65.95** |
| no side/back loss | 94.54 | 19.93 | 16.86 | **85.50** | 32.18 | 14.73 | 83.98 | 27.34 | 15.39 | 1.56 | 36.82 | 58.67 |
| no 2.5D loss | 94.10 | 20.76 | **17.72** | 83.27 | 32.51 | 15.11 | **84.58** | 26.19 | 15.65 | 1.37 | 36.42 | 63.88 |
| PixelNerf [40] | 91.07 | 22.96 | 16.76 | 81.02 | 32.01 | **15.74** | 80.42 | **24.86** | **16.31** | 1.45 | 35.07 | 65.50 |
| EG3D+Img2SG [2] | 85.90 | 30.78 | 13.92 | 77.78 | 39.84 | 12.68 | 79.78 | 31.10 | 13.86 | 2.05 | 11.88 | 23.55 |
| EG3D+PTI [30] | 89.90 | 25.93 | 15.50 | 77.37 | 39.16 | 13.37 | 79.78 | 32.62 | 14.32 | 2.19 | 11.38 | 23.54 |
| Pifu [31] | 75.12 | 41.62 | 11.94 | 75.01 | 43.63 | 12.51 | 73.62 | 32.45 | 13.81 | **1.35** | **37.37** | **68.13** |

Table 3. Ablations and baselines. We evaluate and compare our system using our new AnimeRecon benchmark of illustrations with their ground-truth 3D reconstructions. Image reconstruction metrics are listed for the front (reconstruction of the given input) and back orthographic views, as well as averaged over twelve perspective spin-around views. On the right, chamfer distances and F-1 scores [13] capture the correctness of reconstructed meshes. For fair comparison, lines are preprocessed out of all input images using our illustration-to-render model, and evaluation is performed within our annotated ROIs. The top two of each category are bolded.

input illustration
(original vs. lines removed)

radiance field
(given original vs. lines removed)

Figure 6. The effect of our illustration-to-render module on reconstruction. Without our proposed line-infilling method, the reconstruction module trained on contour-less renders is unable to cross the illustration domain gap. Notice the line artifacts along the shoulders, and improper contours along the chins. (Art attributions in suppl.)



Input
Illustration

Ground Truth
Reconstruction

Our Full
Model

No feature
pooling

No multi-layer
triplane

No side/back
losses

Figure 7. 3D reconstruction module ablations. Row 1: feature pooling helps distribute information across the triplane, and correctly structures the hair as short. Rows 2+3: the multi-layer triplane helps disambiguate the front/back, and side/back-view losses significantly improve both geometry and texture. (Art attributions in suppl.)

Though Mesh R-CNN also proposes the normal consistency metric [13], it is difficult to obtain for our data, as there are assets with broken ground truth normals (zero values, inverted meshes). We can calculate rough consistency by re-estimating all normals, resulting in: Ours 75.6, Pixel-

Nerf 77.5, Img2SG 73.0, PTI 72.7, Pifu 77.3. We observed that these approximate measures don't match well with visual quality, but include them here for completeness.

### 4.1.2 Baselines

Comparisons are shown between our method and several other implicit reconstruction methods (Fig. 5). Of the methods shown, only Pifu [31] receives point-wise signals for optimization; we see that this works to its detriment, as the near-surface point sampling strategies bias the model towards certain geometric structures (such as black eyelashes that often protrude from the face). The other methods using volumetric rendering inherently weigh the supervision signal such that the final rendered product is consistent.

Image2stylegan [2] and Pivotal Tuning Inversion [30] are optimization-based methods of performing single-view reconstruction, the latter of which was used in EG3D to demonstrate reconstruction by projection [4]. We train an EG3D with similar hyperparameters as the Stylegan2 backbone used in our reconstruction module, and allow the two respective baselines to optimize features/weights in the prior to match the given input. Unfortunately, the EG3D prior is not sufficient to regularize the optimization, leading to implausible results; while PTI worked reasonably well for EG3D, it may not work as well on our setup where the back of the head must also be projected.

Lastly we compare to our model with the PixelNerf [40] single-view setup. The key difference in this comparison is that PixelNerf does not use adversarial nor 2.5D reconstruction losses (for fairer comparison with us, we trained PixelNerf with LPIPS in addition to $L_1$, resulting in significantly less blurry results). We see that without these signals provided by our available 3D assets, the quality of details and geometry does not match up to PAniC-3D.

### 4.1.3 Ablations

From Tab. 3 and Fig. 7, we conclude that our architecture decisions improve both the qualitative and quantitative performance of our reconstruction system. It is shown that feature pooling is able to propagate information along triplane axes to improve generated geometries, the multi-layer triplane adds additional model capacity to further disambiguate locations for features, and the addition of fixed side/back-view losses significantly improves geometry and texture. Expectedly, removing 2.5D supervision also deteriorates performance.

### 4.2. Illustration-to-Render Results

As shown in Tab. 4 and Fig. 6, our illustration-to-render model was able to effectively remove contours that would not appear in renders. We evaluate in Tab. 4 how closely the translated illustrations match their rendered counterpart

from the AnimeRecon benchmark (restricted to the ROI annotation), and find that we are able to significantly outperform naive inpainting [36] as well as off-the-shelf image-to-image translation systems [19, 45]. The others struggled to retain key semantic structures like eyes/mouths, and often failed to retain the original identity.

In Fig. 6, it is shown that removing lines indeed alleviates the effects of domain transfer between illustrations and their output predicted radiance field; artifacts like extra contours and bands along shoulders are significantly reduced using our line removal strategy.

## 5. Limitations & Future Work

From the figures comparing the ground-truth 3D model renders and generated radiance fields, it is evident that this task is incredibly difficult. Although PAniC-3D performs better than other baseline methods, there is still a large gap in quality between our reconstructions and real character assets; this is still particularly evident for the hind occluded regions, as well as the areas around the ears connecting the front and back. Our model is usually able to prevent faces from appearing on the backside, but still largely relies on copying the visible portions of accessories to cover the occluded parts. In the future, we may look into ways of decomposing the radiance field in an object-centric manner, and generate accessories through a separate more expressive process. An obvious extension of this work would be to model full-body characters and exploit more opportunities given by the Vroid dataset, which supports blendshape facial expressions, hair physics, full-body rigging, and more.

In conclusion, we propose PAniC-3D, a system to reconstruct stylized 3D character heads directly from illustrated portraits of anime characters. Our anime-style domain poses unique challenges to single-view reconstruction when compared to natural images of human heads, such as hair and accessories with more complex and diverse geometry, and non-photorealistic contour lines. Furthermore, there is a lack of both 3D model and portrait illustration data suitable to train and evaluate this ambiguous stylized reconstruction task. Facing these challenges, our proposed PAniC-3D architecture crosses the illustration-to-3D domain gap with a line-filling model, and represents sophisticated geometries with a volumetric radiance field. We train our system with two large new datasets (11.2k Vroid 3D models, 1k Vtuber portrait illustrations), and evaluate on a novel AnimeRecon benchmark of illustration-to-3D pairs. PAniC-3D significantly outperforms baseline methods, and can generate plausible fully-textured geometries from a single input drawing. We hope that our proposed system and provided datasets may help establish the stylized reconstruction task.



|  | LPIPS↓ | CLIP↑ | PSNR↑ |
|---|---|---|---|
| Ours | **18.26** | **94.97** | **16.96** |
| Telea [36] | 23.91 | 93.88 | 16.67 |
| CycleGAN [45] | 21.39 | 93.81 | 15.23 |
| UGATIT [19] | 39.64 | 85.48 | 13.18 |

Table 4. Illustration-to-render results. Over our AnimeRecon benchmark, we calculate perceptual and color metrics between translated illustrations and their corresponding 3D renders, restricted to the ROI annotations. Our method crosses the illustration-render domain gap better than both naive inpainting (which struggles to retain facial features) and deep image2image models (which fail to retain identity). LPIPS is scaled by 1e2.

## Acknowledgements

# References

[1] Panohead: Geometry-aware 3d full-head generative adversarial networks. 2022. 3

[2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4432–4441, 2019. 6, 7

[3] Gwern Branwen, Anonymous, and Danbooru Community. Danbooru2019 portraits: A large-scale anime head illustration dataset. https://www.gwern.net/Crops#danbooru2019-portraits, March 2019. 3

[4] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. *CVPR*, 2022. 1, 2, 3, 7

[5] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *CVPR*, 2021. 1

[6] Shuhong Chen and Matthias Zwicker. Improving the perceptual quality of 2d animation interpolation. *arXiv preprint arXiv:2111.12792*, 2021. 2, 3

[7] Shuhong Chen and Matthias Zwicker. Transfer learning for pose estimation of illustrated characters. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 793–802, 2022. 3, 4

[8] Zhiguo Cheng and Yuncai Liu. A graph-based method to remove interferential curve from text image. *Machine Vision and Applications*, 17:219–228, 2006. 2

[9] Doug DeCarlo et al. Suggestive contours for conveying shape. In *SIGGRAPH*, pages 848–855. 2003. 4

[10] Szabolcs Dombi. Moderngl, high performance python bindings for opengl 3.3+. *GitHub repository*. 4

[11] Marek Dvorožňák, Daniel Sỳkora, Cassidy Curtis, Brian Curless, Olga Sorkine-Hornung, and David Salesin. Monster mash: a single-view approach to casual 3d modeling and animation. *ACM Transactions on Graphics (TOG)*, 39(6):1–12, 2020. 2

[12] Einar Forselv. moderngl-window, a cross-platform windowing/utility library for moderngl. *GitHub repository*. 4

[13] Georgia Gkioxari and Jitendra Malik. Jj: Mesh r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 5, 6, 7

[14] Stéphane Grabli et. al. Programmable style for npr line drawing. In *Eurographics*, 2004. 4

[15] The Khronos Group. Gltf specification. 2022. 4

[16] hysts. Anime face detector. https://github.com/hysts/anime-face-detector, 2021. 4

[17] Pramook Khungurn. Talking head anime from a single image 2: More expressive. 2021. 2, 4

[18] Pramook Khungurn and Derek Chou. Pose estimation of anime/manga characters: a case for synthetic data. pages 1–6, 2016. 2, 3, 4

[19] Junho Kim, Minjae Kim, Hyeonwoo Kang, and Kwanghee Lee. U-gat-it: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. *arXiv preprint arXiv:1907.10830*, 2019. 2, 3, 8

[20] Kangyeol Kim, Sunghyun Park, Jaeseong Lee, Sunghyo Chung, Junsoo Lee, and Jaegul Choo. Animeceleb: Large-scale animation celebheads dataset for head reenactment. 2022. 2, 3, 4

[21] Chengze Li, Xueting Liu, and Tien-Tsin Wong. Deep extraction of manga structural lines. *ACM Transactions on Graphics (TOG)*, 36(4):1–12, 2017. 2

[22] Zuzeng Lin, Ailin Huang, Zhewei Huang, Chen Hu, and Shuchang Zhou. Collaborative neural rendering using anime character sheets. *arXiv preprint arXiv:2207.05378*, 2022. 2, 3

[23] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7708–7717, 2019. 1

[24] Zhaoliang Lun, Matheus Gadelha, Evangelos Kalogerakis, Subhransu Maji, and Rui Wang. 3d shape reconstruction from sketches via multi-view convolutional networks. In *2017 International Conference on 3D Vision (3DV)*, pages 67–77. IEEE, 2017. 2

[25] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. Stylesdf: High-resolution 3d-consistent image and geometry generation. *CVPR*, 2022. 1

[26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 5

[27] N Shobha Rani, P Vineeth, and Deeptha Ajith. Detection and removal of graphical components in pre-printed documents. *International Journal of Applied Engineering Research*, 11(7):4849–4856, 2016. 2

[28] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 4

[29] A Rehman, F Kurniawan, and T Saba. An automatic approach for line detection and removal without smash-up characters. *The Imaging Science Journal*, 59(3):177–182, 2011. 2

[30] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Transactions on Graphics (TOG)*, 42(1):1–13, 2022. 6, 7

[31] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2304–2314, 2019. 1, 6, 7

[32] Timothy K Shih, Louis H Lin, and Wonjun Lee. Detection and removal of long scratch lines in aged films. In *2006 IEEE International Conference on Multimedia and Expo*, pages 477–480. IEEE, 2006. 2

[33] Edgar Simo-Serra, Satoshi Iizuka, and Hiroshi Ishikawa. Mastering sketching: adversarial augmentation for structured prediction. *ACM Transactions on Graphics (TOG)*, 37(1):1–13, 2018. 2

[34] Edgar Simo-Serra, Satoshi Iizuka, and Hiroshi Ishikawa. Real-time data-driven interactive rough sketch inking. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. 2

[35] Domenico Tegolo and Francesco Isgro. Scratch detection and removal from static images using simple statistics and genetic algorithms. In *Proceedings 2001 International Conference on Image Processing (Cat. No. 01CH37205)*, volume 1, pages 265–268. IEEE, 2001. 2

[36] Alexandru Telea. An image inpainting technique based on the fast marching method. *Journal of graphics tools*, 9(1):23–34, 2004. 8

[37] VRoid. Vrm specification. 2022. 4

[38] 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. 2

[39] Yang Xue, Yuheng Li, Krishna Kumar Singh, and Yong Jae Lee. Giraffe hd: A high-resolution 3d-aware generative model. In *CVPR*, 2022. 1

[40] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 1, 3, 6, 7

[41] Lvmin Zhang, Yi Ji, and Chunping Liu. Danbooregion: An illustration region dataset. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIII 16*, pages 137–154. Springer, 2020. 2

[42] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 3, 5

[43] Yi Zheng, Yifan Zhao, Mengyuan Ren, He Yan, Xiangju Lu, Junhui Liu, and Jia Li. Cartoon face recognition: A benchmark dataset. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2264–2272, 2020. 3

[44] Yue Zhong, Yonggang Qi, Yulia Gryaditskaya, Honggang Zhang, and Yi-Zhe Song. Towards practical sketch-based 3d shape generation: The role of professional sketches. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(9):3518–3528, 2020. 2

[45] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 2, 3, 8