

# LookingGlass: Generative Anamorphoses via Laplacian Pyramid Warping

Pascal Chang<sup>1,2</sup>Sergio Sancho<sup>1,2</sup>Jingwei Tang<sup>2</sup>Markus Gross<sup>1,2</sup>Vinicius Azevedo<sup>2</sup><sup>1</sup>ETH Zürich<sup>2</sup>DisneyResearch|Studios

Figure 1. We propose a method to generate *ambiguous anamorphoses*—images that reveal a hidden image when viewed through a mirror or lens. In the examples above, a conic mirror viewed from the top reveals a turtle hidden in an Earth image; a garden, seen through a lens, shows a bunny, and rotating the lens slightly reveals a gnome; a cylindrical mirror reflects a village painting into the face of an old man.

## Abstract

*Anamorphosis refers to a category of images that are intentionally distorted, making them unrecognizable when viewed directly. Their true form only reveals itself when seen from a specific viewpoint, which can be through some catadioptric device like a mirror or a lens. While the construction of these mathematical devices can be traced back to as early as the 17th century [29], they are only interpretable when viewed from a specific vantage point and tend to lose meaning when seen normally. In this paper, we revisit these famous optical illusions with a generative twist. With the help of latent rectified flow models, we propose a method to create anamorphic images that still retain a valid interpretation when viewed directly. To this end, we introduce Laplacian Pyramid Warping, a frequency-aware image warping technique key to generating high-quality visuals. Our work extends Visual Anagrams [18] to latent space models and to a wider range of spatial transforms, enabling the creation of novel generative perceptual illusions.*

## 1. Introduction

Anamorphosis, derived from the Greek *ana* (“back” or “again”) and *morphe* (“form”), refers to a category of images that are deliberately distorted, rendering them unrecognizable when viewed directly. These optical illusions reveal their true form only when observed from a precise vantage point or through reflective or refractive surfaces, such as mirrors or lenses—objects collectively known as *anamorphoscopes* [24]. These mathematical curiosities became more popular since the 17th century, when the pioneering treatise by the French mathematician J.-F. Nicéron, *La Perspective Curieuse*, laid the foundation for their rigorous construction [29]. However, these images are typically interpretable only from specific angles, losing their meaning when viewed normally.

In this paper, we propose a method for creating anamorphic images using latent text-to-image models. We focus on setups in which the image has a valid interpretation when viewed as-is without distortions. Our work is similar to the recent framework of Visual Anagrams proposed by Geng *et*

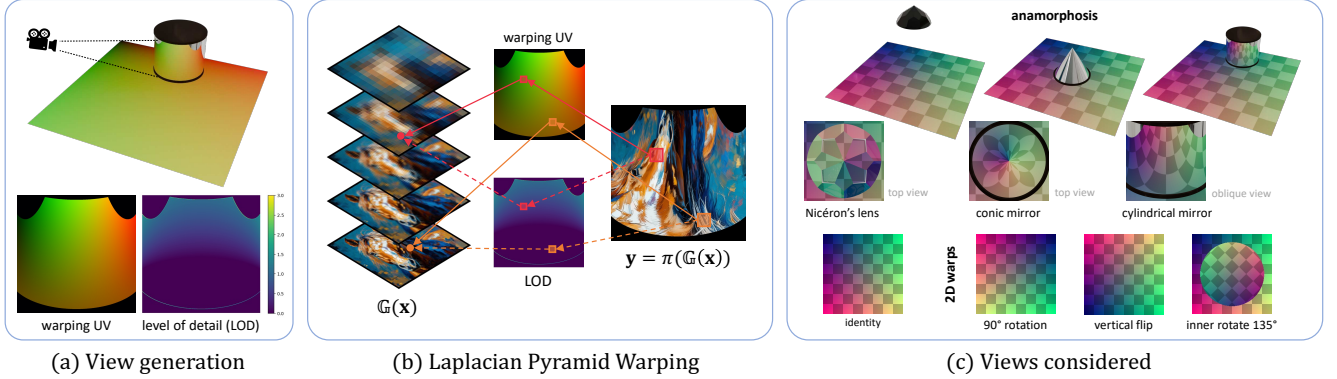


Figure 2. **Laplacian Pyramid Warping.** (a) The view mappings are generated using a ray tracer and the Level of Detail (LOD) map is computed. (b) For each pixel, our forward warping algorithm looks up in the warping UV mapping and LOD to fetch the corresponding value. (c) We consider different views, from 2D transformations like vertical flip and arbitrary angle rotation to complex 3D projections.

*al.* [18], which generates ambiguous images by synchronizing diffusion paths across multiple views. However, their method is constrained to pixel-space diffusion models and limited to orthogonal transformations of image pixels. We address these two limitations in this paper. First, we enable the use of latent diffusion and flow models in an artifact-free manner, improving the generation quality. We believe this will render the generation of these illusions more accessible. Second, we introduce *Laplacian Pyramid Warping*, a robust image-warping technique that handles complex image transformations while preserving high-frequency details. This enables the generation of intricate anamorphoses involving complex reflective and refractive surfaces, with minimal sacrifice to image quality. Compared to previous work, our method demonstrates a significant boost in both the quality and expressiveness of the generated results.

## 2. Related Work

**Computational optical illusions.** Anamorphosis can be dated back to around the 16th century [22, 29, 39], when artists either hand drew the illusions on paper or used grids to create them systematically. Since then, the generation of optical illusions has seen significant progress, especially with the advent of computational methods in recent years. Earlier work focuses on creating illusion with 2D images, such as revealing an image by stacking transparent sheets of images [28], achieving appearance change of images at different viewing distances [30], creating static images that appear to move [11], and designing a refractive lens for revealing a hidden image from dots [31]. Beyond image manipulation, several works explored 3D illusions. Hsiao *et al.* [21] introduced multi-view wire art, where a single 3D wireframe produces different projected images from various perspectives. Perroni-Scharf and Rusinkiewicz [32] extended this idea to 3D-printed view-dependent surfaces.

Apart from illusion based on 3D geometries, Chu *et al.* [12] explored camouflaging objects by retexturing them, while Chandra *et al.* [8] developed models that shift in perception based on lighting changes. In contrast, we focus on 2D illusions that require 3D objects to reveal the hidden views.

**Illusions with diffusion models.** Recent work has revealed the potential of diffusion models in creating optical illusions. Burgert *et al.* [4] employ score distillation sampling (SDS) to generate images that align with multiple prompts from different viewpoints. Although their optimization-based method can theoretically produce anamorphoses, it suffers from lower image quality and long inference times. Visual Anagrams [18] introduces a formal framework for illusion generation in a single diffusion pass. However, their approach is limited to orthogonal transformations, making it unsuitable for generating the complex deformations needed for anamorphoses. Subsequent studies have also explored various types of illusions, such as visually meaningful spectrograms [10] and generative hybrid images [17]. Our proposed method is most similar to Visual Anagrams. Key differences, however, are that we extend to latent space models and a broader range of transformations. A concurrent work, Illusion3D [16], builds on [4] to generate 3D anamorphic illusions, but appears constrained in quality and artistic flexibility. We outline the key differences with our method in the supplementary material.

Beyond academic research, the artistic community has also explored diffusion models for optical illusions. Notably, an anonymous artist known as MrUgleh [38] repurposed a model fine-tuned for generating QR codes [25, 44] to create images that subtly mimic the global structure of a specified template image. Our focus is on generating ambiguous images and anamorphoses based on text prompts, which does not require an image template.



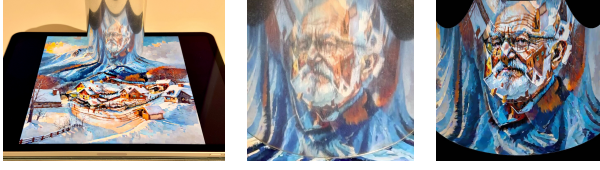


Figure 3. A real life demo of the cylindrical mirror illusion.

**Synchronized diffusion.** In Visual Anagrams [18], diffusion paths from different viewpoints are synchronized by averaging the predicted noise at each timestep. Numerous studies have explored merging diffusion paths, often in the context of controlled image generation. MultiDiffusion [3] proposes a least-squares formulation for merging views, which simplifies to averaging in the special case of equal-size crops—a setup they apply to panorama generation. DiffCollage [45] synthesizes large-scale content by merging outputs from diffusion models trained on segments of the larger composition. SyncTweedies [23] thoroughly examines synchronization techniques, finding that averaging the predicted clean images yields the best quality. Closer to our approach, Generative Powers of Ten [40] creates infinite zoom videos by merging concentric views at different resolutions using Laplacian pyramids. But their method is tailored to the specific use case of zooming. One of our contributions, *Laplacian Pyramid Warping*, generalizes this approach to arbitrary views.

**Image pyramids in vision and graphics.** Image pyramids, particularly Gaussian and Laplacian pyramids, are widely used in computer vision for their multi-scale representation capabilities [5, 13, 42]. By decomposing images hierarchically, pyramids enable efficient compression, progressive image reconstruction, and seamless blending—essential in applications like panorama stitching and HDR imaging [6]. Beyond blending, Gaussian pyramids are central to scale-invariant object detection and recognition, where they assist in feature detection for algorithms like SIFT [26]. They are also valuable in texture analysis and synthesis [19], and optical flow estimation [34], where multi-scale representations enhance accuracy and reduce artifacts.

In computer graphics, pyramids relate closely to techniques like texture MIP-mapping [15] and antialiasing, which address the challenges of rendering textures at varying distances and viewing angles. MIP-maps, essentially a Gaussian pyramid form, allow graphics engines to select the appropriate level of detail (LOD) based on screen space, minimizing artifacts like flickering and enhancing both quality and efficiency. We repurpose these texture MIP-mapping techniques in our proposed method for frequency-aware image warping.

### 3. Preliminaries

#### 3.1. Text-conditioned Rectified Flows

In Rectified Flows (RFs), a noise sample  $\mathbf{z}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is mapped to an image  $\mathbf{z}_1 \sim p_1$  through the ODE:

$$d\mathbf{z}_t = \mathbf{u}_t(\mathbf{z}_t, y)dt, \quad (1)$$

where  $t \in [0, 1]$ ,  $y$  is an optional text prompt conditioning, and the velocity field is typically parameterized with a neural network, *i.e.*  $\mathbf{u}_t(\mathbf{z}_t, y) = \mathbf{u}_\theta(\mathbf{z}_t; t, y)$ . At inference, the ODE is discretized, and solved with classical integration schemes such as forward Euler:

$$\mathbf{z}_{t+\Delta t} = \mathbf{z}_t + \mathbf{u}_\theta(\mathbf{z}_t; t, y)\Delta t, \quad (2)$$

**Classifier-free guidance (CFG).** As in diffusion models, classifier-free guidance [20] can be used to improve sample quality in RFs. The final velocity interpolates between a text-conditioned and an unconditional prediction:

$$\hat{\mathbf{u}}_t = (1 + \omega)\mathbf{u}_\theta(\mathbf{z}_t; t, y) - \omega\mathbf{u}_\theta(\mathbf{z}_t; t, \emptyset), \quad (3)$$

where  $\omega$  is the classifier-free *guidance scale*. Higher guidance scales typically improve sample quality at the expense of diversity, but also tend to produce over-saturated images.

**Predicted clean image.** At any intermediate timestep  $t$ , an estimate of the clean image, denoted  $\mathbf{z}_{1|t}$ , can be obtained by a single Euler step to  $t = 1$  using the current velocity estimate:

$$\mathbf{z}_{1|t} = \mathbf{z}_t + \mathbf{u}_\theta(\mathbf{z}_t; t, y)(1 - t). \quad (4)$$

Equation (4) can be seen as the flow matching equivalent of the Tweedie’s formula [35] in diffusion models.

#### 3.2. Gaussian & Laplacian Pyramid

A Gaussian pyramid is a multi-scale representation of an image obtained by iteratively applying a Gaussian blur kernel  $\kappa$  and downsampling  $\mathbf{D}(\cdot)$ . Given an image  $\mathbf{x}$ , the image  $\mathbb{G}_l$  at level  $l$  is computed from the previous level as

$$\mathbb{G}_l(\mathbf{x}) = \mathbf{D}(\kappa(\mathbb{G}_{l-1}(\mathbf{x}))),$$

where  $\mathbb{G}_0(\mathbf{x}) = \mathbf{x}$  is the original image.

A Laplacian pyramid stores the high-frequency details between each level of a Gaussian pyramid. Each Laplacian level  $\mathbb{L}_l$  is defined as the difference between a Gaussian level and the upsampled version of the next level

$$\mathbb{L}_l(\mathbf{x}) = \mathbb{G}_l(\mathbf{x}) - \mathbf{U}(\mathbb{G}_{l+1}(\mathbf{x})),$$

with  $\mathbb{L}_{L-1}(\mathbf{x}) = \mathbb{G}_{L-1}(\mathbf{x})$  for a pyramid of depth  $L$  and  $\mathbf{U}(\cdot)$  being the upsample operator. To reconstruct the image, we recursively add each Laplacian level back to the upsampled version of the next level.

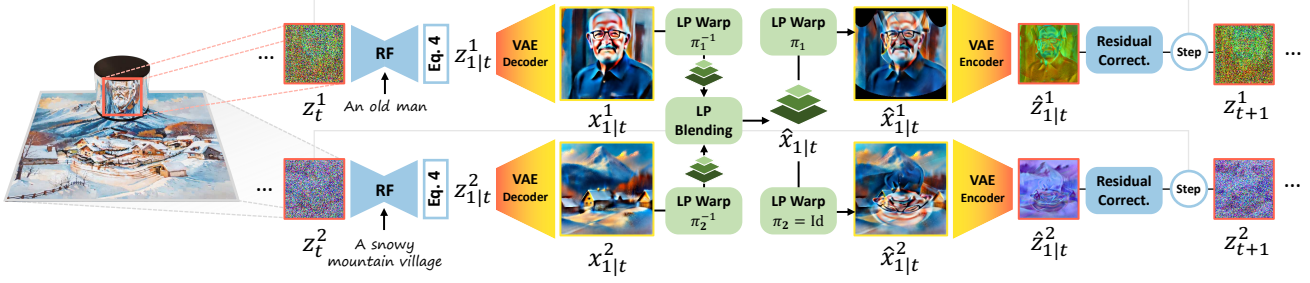


Figure 4. **Our Proposed Pipeline.** At each denoising step, the estimated final image is computed from the network velocity estimate and decoded into image space. Image warping and view aggregation is performed in image space using Laplacian pyramids, before encoding back into latent space for the diffusion step.

### 3.3. Visual Anagrams

The work of Geng *et al.* [18] proposes creating multi-view images by using a text-to-image generative model to simultaneously denoise multiple views of an image. The original paper utilizes diffusion models. We summarize the method here through the terminologies of RFs for simplicity.

A canonical space  $\mathcal{C}$  is defined for an image. A set of prompts  $y_i$  are associated with different view functions  $\pi_i$ , which transform the image from the canonical space to the target space  $\mathcal{T}$  where rectified flow models are applied. At each timestep of the inference,  $\mathbf{z}_t^i$  of each view is transformed into the canonical space, and averaged together with the other views. After averaging, the noisy images in the canonical space are transformed back to the target space as  $\hat{\mathbf{z}}_t^i$ , replacing the original  $\mathbf{z}_t^i$  as

$$\hat{\mathbf{z}}_t^i = \pi_i \left( \frac{1}{N} \sum_j \pi_j^{-1} (\mathbf{z}_t^j) \right). \quad (5)$$

Transforming noisy samples, as in Geng *et al.* [18], limits possible mappings between the canonical and target spaces to be orthogonal transformations such as flipping, rotation and permutation of pixels. This happens since arbitrarily warping a noise sample is generally more difficult than warping images, as bilinear and bicubic interpolations can destroy the Gaussian noise properties [9]. Moreover, SyncTweedies [23] showed that averaging the predicted clean image  $\mathbf{z}_{1|t}^i$ , instead of noisy samples, produces higher quality results. Thus, to allow more general transformations and improve generation quality, we modify Equation (5) to average predicted clean image using Tweedie’s formula as

$$\hat{\mathbf{z}}_{1|t}^i = \pi_i \left( \frac{1}{N} \sum_j \pi_j^{-1} (\mathbf{z}_{1|t}^j) \right). \quad (6)$$

## 4. Generative Anamorphosis

Our goal is to generate high-quality anamorphoses from text prompts. Anamorphoses involve view functions beyond simple transformations such as rotation, flipping and pixel permutations where no analytical transformations can be defined. We opt for a more general representation:  $\pi$  is now a 2-channel image of UV coordinates indicating where to fetch values in the canonical view for each pixel  $\pi(x, y) = (u, v)$ . We implement the transformation with a simple raytracer by placing a UV coordinate texture on the main image plane, and rendering the result when viewing through mirrors or lenses (see Figure 2).

Simply adopting SyncTweedies [23], however, is not enough when considering latent diffusion models. Thus, in Section 4.1, we present a generalization of previous approaches to latent diffusion models. Naively averaging arbitrary transformations with highly distorted regions, typical when looking through curved mirrors or lenses, results in visual artifacts. We introduce a novel Laplacian Image Warping method that utilizes a multi-level texture structure inspired by classic works in computer graphics [13, 42] in Section 4.2 to alleviate this problem. Additional design choices are then discussed in Section 4.3. Figure 4 and Algorithm 1 show an overview of our method.

### 4.1. Latent Visual Anagrams

Tancik [37] generates multi-view illusions with Stable Diffusion 1.5 [36] using a similar pipeline to Geng *et al.* [18]. Images are transformed to canonical space through views in the latent space, and decoded to the clean image after the denoising process is finished. The results contain visual artifacts due to the fact that VAEs for latent diffusion and flow models are generally not trained to be equivariant. When the latent images are deformed, the corresponding decoded image does not necessarily share the same transformation. While these artifacts are less pronounced in recent models with larger latent spaces (*e.g.* Stable Diffusion 3 [14]), they still persist and create undesirable strokes (see Figure 5).

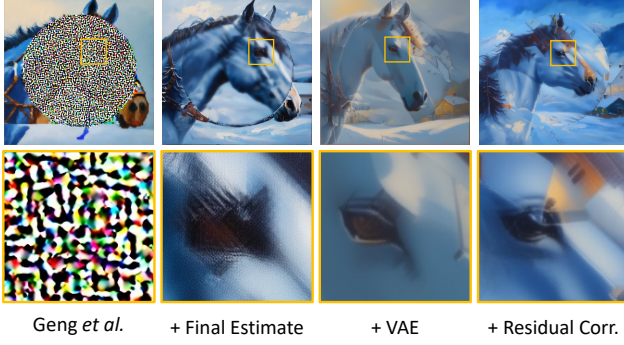


Figure 5. **Latent Visual Anagrams.** In this 135° rotation example, we demonstrate that contributions from Sec. 4.1 improve the generation of visual anagrams with latent models. While the final estimate from SyncTweedies [23] partially addresses noise issues, artifacts from the VAE persist. Our VAE encoding/decoding process and residual correction further enhance image quality.

**VAE encoding-decoding.** To tackle this issue, we propose to keep all image transformations in image space. This is done by decoding the estimated clean image latent to pixel space, applying the transformation, and re-encoding it to latent space. Denoting  $\mathcal{E}$ ,  $\mathcal{D}$  the encoder and decoder of the VAE respectively, our estimated clean latent at timestep  $t$ , in view  $i$ , becomes

$$\hat{\mathbf{z}}_{1|t}^i = \mathcal{E} \circ \pi_i \left( \frac{1}{N} \sum_j \pi_j^{-1} \circ \mathcal{D} \left( \mathbf{z}_{1|t}^j \right) \right). \quad (7)$$

**Residual correction.** While Eq. (7) proved to be effective in removing the artifacts, we observed that it can be sensitive to the VAE reconstruction quality. In particular, when the input latent is far from the training distribution of the VAE (e.g. predicted clean image from early steps of denoising), the reconstruction typically fails to match the value range of the input. This creates washed-out colors in the final image. We correct the reconstruction failure through a first-order term  $\Delta \hat{\mathbf{z}}_{1|t}^i$ . This is done by first computing the *residual* between the latent and the decoded-encoded latent. This residual is then transformed to a target view as

$$\Delta \hat{\mathbf{z}}_{1|t}^i = \pi_i \left( \frac{1}{N} \sum_j \pi_j^{-1} \left( \mathbf{z}_{1|t}^j - \mathcal{E} \circ \mathcal{D} \left( \mathbf{z}_{1|t}^j \right) \right) \right). \quad (8)$$

This correction term is then added to Eq. (7) as the final estimated clean latent. The correction has two nice properties. First, if the VAE reconstruction is perfect, then  $\mathcal{E}(\mathcal{D}(\mathbf{z})) = \mathbf{z}$ , and  $\Delta \mathbf{z} = 0$ . Therefore, no correction is done. Second, when there is only one single identity view ( $N = 1$ ), we recover  $\hat{\mathbf{z}}_{1|t}^j = \mathbf{z}_{1|t}^j$ , the original predicted clean latent, as expected. Figure 5 compares different steps presented in this section, and shows the effectiveness of the residual correction in maintaining the vibrant colors.

## 4.2. Laplacian Pyramid Warping (LPW)

While the above modifications enable using latent diffusion and flow models for creating ambiguous images with arbitrary 2D transformations, applying it to the generation of anamorphoses still presents a few challenges. First, a view in anamorphosis rarely covers the entirety of the main image. As such, boundaries can create visible artifacts and seams that are undesirable. Second, having views that have varying degrees of stretching can lead to degraded high frequency details (see Figure 6).

We identify the key problem being using the averaging operation as aggregation of views. When images of different views are transformed to the canonical view, they can have different frequency components. One view may map to a small pixel subset in the canonical view or undergo large stretching. Averaging pixels of these views with another view not stretched so much ignores the frequency mismatching problem. To solve this problem, we propose to use Laplacian pyramid [7] for blending the views.

After decoding the predicted clean latent, *Inverse Laplacian Warping*  $\pi^{-1}$  is used to map the image to a Laplacian pyramid in the canonical view. The canonical views are then aggregated through Laplacian Pyramid Blending. Afterwards, the canonical views are transformed back to images through *Forward Laplacian Warping*  $\pi$ .

**Forward Laplacian Warping.** Given an image  $\mathbf{x}$  and a view projection  $\pi$ , we propose a Level-of-Detail-Aware (LOD-aware) method to compute  $\mathbf{y} = \pi(\mathbf{x})$ . First, we build a Gaussian pyramid from our image  $\mathbb{G}(\mathbf{x}) = \{\mathbb{G}_0(\mathbf{x}), \dots, \mathbb{G}_{L-1}(\mathbf{x})\}$ . Then, we compute a LOD level map using  $\pi$  and image-space derivatives. For a given pixel  $(x, y)$ , the LOD level is given by:

$$l = \log_2 \left( \max \left( \sqrt{\left( \frac{\partial u}{\partial x} \right)^2 + \left( \frac{\partial u}{\partial y} \right)^2}, \sqrt{\left( \frac{\partial v}{\partial x} \right)^2 + \left( \frac{\partial v}{\partial y} \right)^2} \right) \right). \quad (9)$$

The transformed image is obtained by sampling the pyramid  $\mathbb{G}(\mathbf{x})$  with the LOD map using nearest or trilinear interpolation. This approach is commonly used in computer graphics when rendering textures to avoid aliasing. Our contribution lies in connecting the method to our problem and repurposing the idea for image warping. For clarity, we simplify the equation as  $\mathbf{y} = \pi(\mathbb{G}(\mathbf{x}))$ .

**Inverse Laplacian Warping.** To transform an image  $\mathbf{y}$ , into a Laplacian pyramid in the canonical view  $\mathbb{G}(\mathbf{x}) = \pi^{-1}(\mathbf{y})$ , we define the inverse Laplacian warping operation. As inverting an arbitrarily complex image transformation is infeasible, we make use of the gradient of Forward Laplacian Warping. Consider a dummy zero image  $\mathbf{x}^0 = \mathbf{0}$ , our inverted image pyramid is given by:

$$\mathbb{G}(\mathbf{x}) = -\nabla_{\mathbf{x}^0} \left[ \frac{1}{2} \|\pi(\mathbf{x}^0) - \mathbf{y}\|^2 \right]. \quad (10)$$



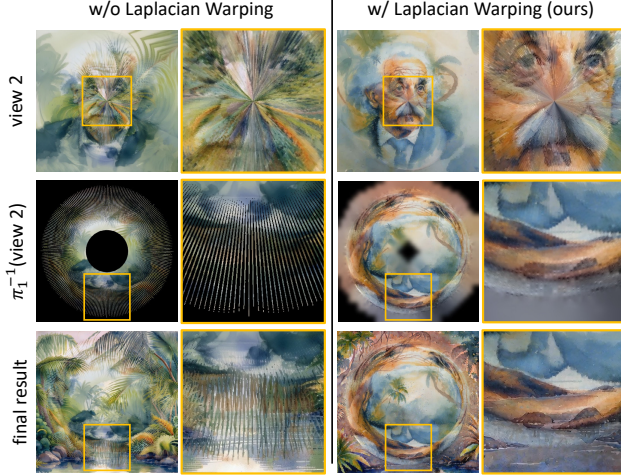


Figure 6. **Inverse Laplacian Warping.** In this conic mirror example, we use a challenging pair of prompts to demonstrate our inverse warping. The main view shows “a jungle,” while the mirror view reveals “a portrait of Einstein.” Without inverse Laplacian warping, gaps from the inverse transformation cause striped artifacts and distortions. Our Inverse Laplacian Warping (Sec. 4.2) correctly assigns values at the right frequencies, eliminating artifacts and making the mirror view more recognizable.

This effectively transports the pixels in  $\mathbf{y}$  to the corresponding location and level in the pyramid. We then propagate the changes of lower levels to higher ones, and extract a Laplacian pyramid out of the resulting pyramid. More details can be found in the supplementary material.

**Laplacian Pyramid Blending** is used for aggregating different canonical views to obtain a synchronized image for the given denoising step. Each level is averaged and the final image is reconstructed from the resulting Laplacian pyramid. In the case of anamorphic illusions, most views will only cover the identity space partially, so the blending is weighted by a set of masks at each level. Special care need to be taken at the boundary of the masks, as well as during averaging, which we further discuss in the supplementary material.

### 4.3. Design Choices & Further Improvements

**Model choice.** Our method is designed to work with latent rectified flow and diffusion models. However, we observed that different models behave differently to our synchronization scheme, which we discuss in the supplementary. For the majority of our experiments, we use Stable Diffusion 3.5 because of their higher visual quality.

**Improved consistency with time travel.** Similar to DDNM [41], RePaint [27] and Bansal *et al.* [2], we found

---

#### Algorithm 1: LookingGlass

---

**Input:**  $\forall i \in 0, \dots, N-1$ : view transformations  $\pi_i$  and text prompts  $y_i$ , with  $N$  being the number of views. Pretrained RF model  $\mathbf{u}_\theta$ .

**Output:** Final images in each view  $\mathbf{x}_1^0, \dots, \mathbf{x}_1^{N-1}$

```

1  $\mathbf{z}_0^0, \dots, \mathbf{z}_0^{N-1} \sim \mathcal{N}(0, I)$ 
2 for  $t \leftarrow 0 : T$  do
3   for  $i \leftarrow 0 : N-1$  do
4      $\mathbf{z}_{1|t}^i \leftarrow \mathbf{z}_t^i + \mathbf{u}_\theta(\mathbf{z}_t^i; t, y_i)(1-t)$   $\triangleright$  Eq. (4)
5      $\mathbf{x}_{1|t}^i \leftarrow \text{VAE\_DECODE}(\mathbf{z}_{1|t}^i)$ 
6      $\Delta \mathbf{z}_{1|t}^i \leftarrow \mathbf{z}_{1|t}^i - \text{VAE\_ENCODE}(\mathbf{x}_{1|t}^i)$ 
7   end
8    $\hat{\mathbf{x}}_{1|t} \leftarrow \text{LAPLACIAN\_BLENDING}$ 
      $(\pi_0^{-1}(\mathbf{x}_{1|t}^0), \dots, \pi_{N-1}^{-1}(\mathbf{x}_{1|t}^{N-1}))$ 
9    $\Delta \hat{\mathbf{z}}_{1|t} \leftarrow \text{LAPLACIAN\_BLENDING}$ 
      $(\pi_0^{-1}(\Delta \mathbf{z}_{1|t}^0), \dots, \pi_{N-1}^{-1}(\Delta \mathbf{z}_{1|t}^{N-1}))$ 
10  for  $i \leftarrow 0 : N-1$  do
11     $\hat{\mathbf{x}}_{1|t}^i \leftarrow \pi_i(\hat{\mathbf{x}}_{1|t})$ 
12     $\hat{\mathbf{z}}_{1|t}^i \leftarrow \text{VAE\_ENCODE}(\hat{\mathbf{x}}_{1|t}^i)$   $\triangleright$  Eq. (7)
13     $\Delta \hat{\mathbf{z}}_{1|t}^i \leftarrow \pi_i(\Delta \hat{\mathbf{z}}_{1|t})$   $\triangleright$  Eq. (8)
14     $\hat{\mathbf{z}}_{1|t}^i \leftarrow \hat{\mathbf{z}}_{1|t}^i + \Delta \hat{\mathbf{z}}_{1|t}^i$ 
15     $\mathbf{z}_{t+1}^i \leftarrow \text{DENOISING\_STEP}(\hat{\mathbf{z}}_{1|t}^i, \mathbf{z}_t^i)$ 
16  end
17 end
18  $\{\mathbf{x}_1^i\}_i \leftarrow \text{VAE\_DECODE}(\{\mathbf{z}_T^i\}_i)$ 
19 return  $\mathbf{x}_1^0, \dots, \mathbf{x}_1^{N-1}$ 
```

---

repeating segments of the denoising process allows the model to blend different views better. To keep the inference efficient, similar to FreeDoM [43], we only apply time traveling at intermediate timesteps between 20% and 80% of the denoising process and use segments of size 1.

**Prioritizing a single view.** We observed that a key component to the success of previous works [4, 18, 37] lies in the fact that generated images generally lack detail. This makes it easier for human imagination to interpret image features differently based on different prompts. With latent flow models, our method generates highly detailed images at 1K resolution. This poses a new challenge, as high-frequency details are rarely compatible between the views and easily give away hidden views. As this is an inherent problem with high resolution images, we propose to prioritize one of the views, which is defined by the user. We set a portion of the last timesteps of the denoising process to be solely denoising for the chosen view. This encourages the model to create coherent details towards the end, while hiding the remaining views better.

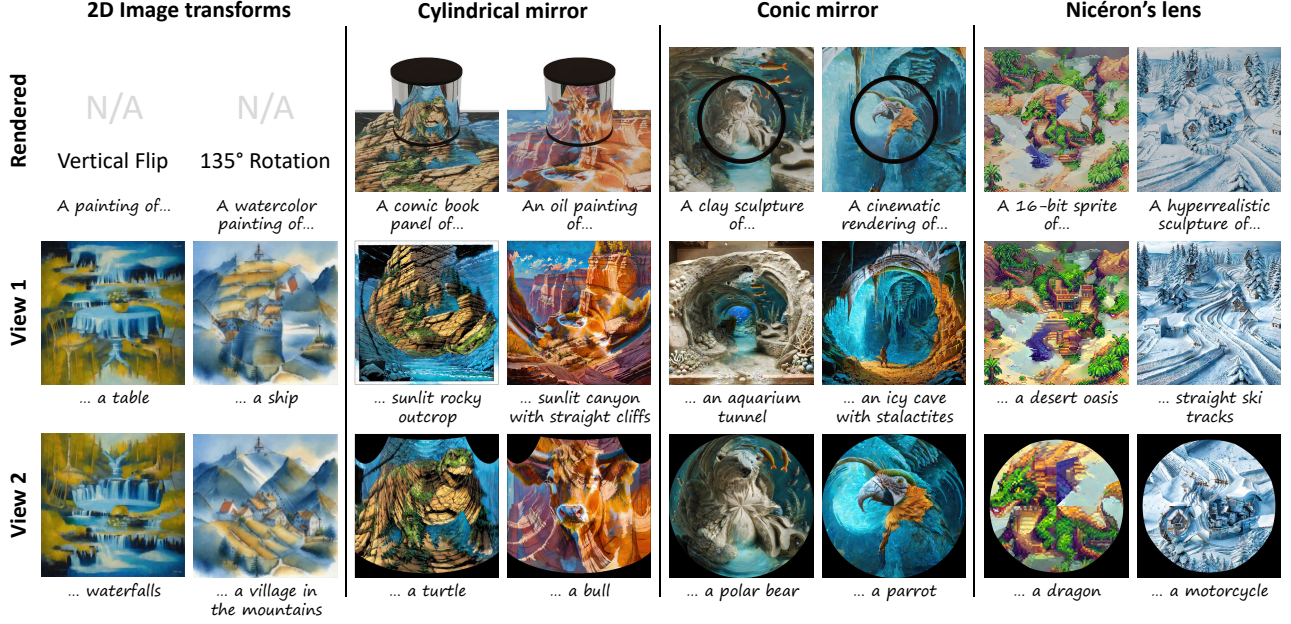


Figure 7. **Generative Anamorphoses.** Generated results with our approach for 2D transformations, and the three types of anamorphoses: a cylindrical mirror, a conic mirror, and Nicéron’s lens. A rendering of the physical setup is shown in the top row when applicable.

## 5. Experiments

### 5.1. Views Considered

In this section we briefly describe the views considered for generating illusions, which are shown in Figure 2 (c).

- **2D transforms.** As in Visual Anagrams [18], we generate vertical flip and 90° rotation illusions. Since our method can handle arbitrary view projections, we also create optical illusions involving 135° rotations.
- **Mirror cone.** A conic mirror placed at the center of an image reveals a hidden picture when viewed from the top.
- **Mirror cylinder.** A cylindrical mirror located over an image shows a hidden picture when viewed from an angle.
- **Nicéron’s lens.** We replicate the setting described by Nicéron [29]. Given an image, and looking at it through a lens sculpted with polygonal faces, the irregular refraction will generate novel images.

Note that in our examples, we always set one view as the identity view (*i.e.* the canonical view), but this is not required. Transforms  $\pi_i$  can be defined relative to a canonical space  $\mathcal{C}$  without explicitly specifying the latter.

### 5.2. Quantitative Results

We compare our method quantitatively with Visual Anagrams [18], Diffusion Illusions [4], SyncTweedies [23], and Tancik [37]. Following Geng *et al.*, CLIP is used [33] to measure the **alignment score**  $\mathcal{A}$  and the **concealment score**  $\mathcal{C}$ . We generate 50 pairs of prompts and create images for these prompts using standard denoising (no optical

illusion), as well as with our method. The results are reported in Table 1 for three tasks: simple vertical flip, a more complex 135° rotation, and the cylindrical mirror anamorphosis. Our results are generated with Stable Diffusion 3.5 Medium on a Nvidia GeForce RTX 4090 GPU. Using 30 inference steps, with time-traveling between 20% and 80% of the diffusion process repeated twice, we generate an image pair in approximately 80 seconds.

	Method	$\mathcal{A} \uparrow$	$\mathcal{A}_{0.9} \uparrow$	$\mathcal{C} \uparrow$	$\mathcal{C}_{0.9} \uparrow$	FID $\downarrow$	KID $\downarrow$
Vertical Flip	Geng <i>et al.</i> [18]	0.306	0.340	0.695	0.786	149.24	0.057
	Tancik SD 3.5 [37]	0.306	0.349	0.693	0.806	132.52	0.049
	Burgert <i>et al.</i> [4]	0.281	0.324	0.679	0.778	219.84	0.115
	SyncTweedies [23]	0.302	0.341	0.707	0.801	132.62	0.054
	<b>LookingGlass (ours)</b>	<b>0.297</b>	<b>0.338</b>	<b>0.680</b>	<b>0.779</b>	<b>124.67</b>	<b>0.049</b>
135° Rotation	Geng <i>et al.</i> [18]	0.262	0.308	0.563	0.652	293.00	0.254
	Tancik SD 3.5 [37]	0.194	0.216	0.498	0.509	439.35	0.545
	Burgert <i>et al.</i> [4]	0.280	0.326	0.654	0.760	223.21	0.120
	SyncTweedies [23]	0.283	0.335	0.647	0.753	166.03	0.083
	<b>LookingGlass (ours)</b>	<b>0.295</b>	<b>0.338</b>	<b>0.666</b>	<b>0.767</b>	<b>129.74</b>	<b>0.055</b>
Cylindrical Mirror	Geng <i>et al.</i> [18]	0.171	0.198	0.506	0.546	285.23	0.216
	Tancik SD 3.5 [37]	0.171	0.198	0.505	0.547	284.97	0.215
	Burgert <i>et al.</i> [4]	0.261	0.304	0.706	0.795	229.65	0.138
	SyncTweedies [23]	0.241	0.284	0.673	0.763	138.69	0.082
	<b>LookingGlass (ours)</b>	<b>0.272</b>	<b>0.318</b>	<b>0.698</b>	<b>0.810</b>	<b>130.27</b>	<b>0.070</b>

Table 1. **Quantitative Comparison.** Sample quality is assessed with FID/KID against a reference dataset of 3.2k images generated from the same set of prompts. Image-prompt alignment is assessed using CLIP alignment score  $\mathcal{A}$ , and concealment score  $\mathcal{C}$  introduced in [18]. While all methods achieve comparable results for the vertical flip, LookingGlass surpasses previous approaches on more complex transformations, including anamorphoses. Please see the supplementary material for more quantitative evaluations.



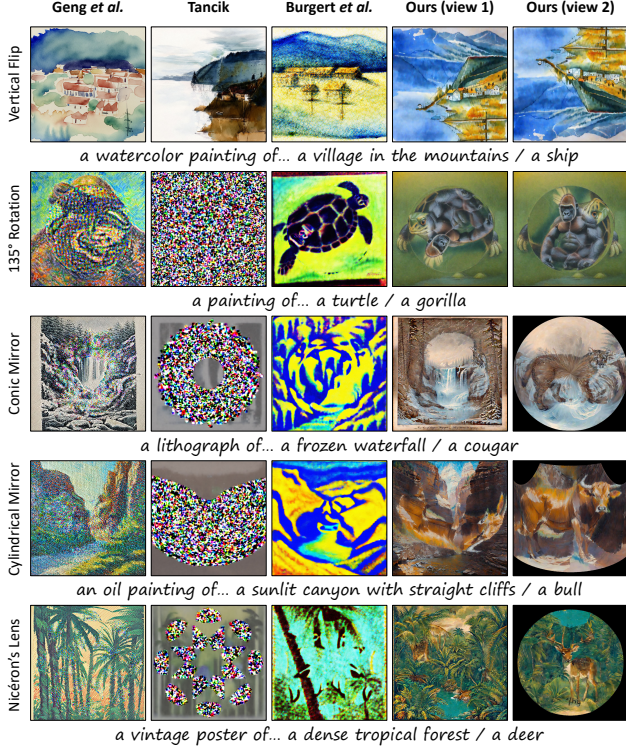


Figure 8. **Qualitative Comparison.** We compare our method against prior work for the considered views (Sec. 5.1). Note that all transformations except the vertical flip are not supported by Geng *et al.* [18] and Tancik [37] due to their inherent limitations.

### 5.3. Ablations

We perform an ablation of the proposed approach. Figure 5 shows that warping the clean image estimate significantly improves image quality compared to warping the predicted noise. Additionally, the proposed VAE encoding/decoding and latent residual correction enhance detail preservation and reduce reconstruction errors. Table 2 suggests that time traveling improves visual quality, as reflected by the FID metric, but may lead to reduced prompt alignment. Further qualitative ablations on time traveling and the effects of prioritizing a single view can be found in the Appendix.

	$\mathcal{A} \uparrow$	$\mathcal{A}_{0.9} \uparrow$	$\mathcal{C} \uparrow$	$\mathcal{C}_{0.9} \uparrow$	FID $\downarrow$	KID $\downarrow$
Geng <i>et al.</i> SD 3 [18]	0.219	0.249	0.516	0.571	335.24	0.297
+ Final Estimate (Eq. 4)	0.273	0.320	0.657	0.757	171.61	0.106
+ VAE (Sec. 4.1)	0.293	0.333	0.717	0.814	160.17	0.083
+ LPW (Sec. 4.2)	<b>0.300</b>	<b>0.336</b>	<b>0.723</b>	0.814	153.27	<b>0.074</b>
+ Time Travel (Sec. 4.3)	0.295	0.331	0.716	<b>0.816</b>	<b>150.01</b>	<b>0.074</b>

Table 2. **Ablation Study.** Starting from the baseline of Geng *et al.* [18], we ablate the contributions introduced by our approach. We also show the 90th-percentile of the CLIP-based metrics, as we are interested in the best case performance.

### 5.4. Qualitative Results

Figures 8 and 7 show example images generated using our method in combination with Stable Diffusion 3.5. Figure 3 features a real-life demonstration of the cylindrical example, confirming that the generated results function as intended in practice. Additional results for all three types of anamorphoses are provided in the Appendix. These results are selected from a set of curated prompts that worked best, as not all prompt combinations are expected to work well.

### 5.5. User Study

We conducted a user study with 27 participants. A total of 10 prompt pairs, the same for all three types of illusions, were selected to generate results using different methods, with the same random seed (no hand-picked samples). Participants ranked the samples from 1 (best) to 5 (worst) based on prompt fidelity, style adherence and overall visual quality. The results are reported in Figure 9.

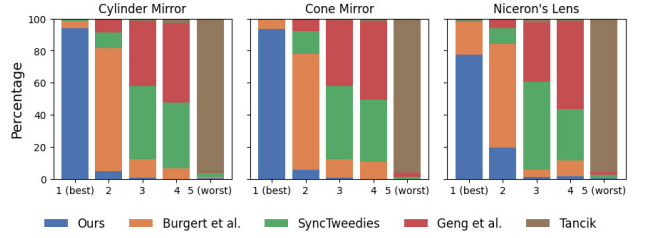


Figure 9. **User Study.** A survey of 27 participants shows that our method (blue) is consistently preferred to prior works.

## 6. Conclusion and Discussion

We presented LookingGlass, a method for generating high-quality, ambiguous anamorphic illusions with latent RF models. Our contributions are twofold. First, we extended Geng *et al.* [18] to latent space models without introducing artifacts, enabling the first feed-forward generation of high-quality illusions with common models. This makes illusion generation more accessible and allows for possible integration with modern generative techniques like ControlNet and DreamBooth. Second, we introduced *Laplacian Pyramid Warping* (LPW), a warping method that preserves fine details while handling extreme distortions. While essential for illusion generation, LPW is also compatible with pixel diffusion models and has potential applications in generative mesh texturing and panorama synthesis.

Despite these advances, selecting effective prompts remains a challenge, as not all prompts lead to high-quality illusions. Additionally, while our method is significantly more efficient than optimization-based approaches like Burgert *et al.* [4], it is computationally more expensive than Geng *et al.* [18] due to the VAE intermediate steps. We leave this to future work.



## References

- [1] DeepFloyd Lab at StabilityAI. DeepFloyd IF: a novel state-of-the-art open-source text-to-image model with a high degree of photorealism and language understanding. <https://www.deepfloyd.ai/deepfloyd-if>, 2023. Retrieved on 2023-11-08. [15](#)
- [2] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 843–852, 2023. [6](#)
- [3] Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: Fusing diffusion paths for controlled image generation. *arXiv preprint arXiv:2302.08113*, 2023. [3](#)
- [4] Ryan Burgert, Xiang Li, Abe Leite, Kanchana Ranasinghe, and Michael Ryoo. Diffusion illusions: Hiding images in plain sight. In *ACM SIGGRAPH 2024 Conference Papers*, New York, NY, USA, 2024. Association for Computing Machinery. [2](#), [6](#), [7](#), [8](#), [14](#), [15](#), [17](#), [18](#)
- [5] P. Burt and E. Adelson. The laplacian pyramid as a compact image code. *IEEE Transactions on Communications*, 31(4): 532–540, 1983. [3](#)
- [6] Peter J. Burt and Edward H. Adelson. A multiresolution spline with application to image mosaics. *ACM Trans. Graph.*, 2(4):217–236, 1983. [3](#)
- [7] Peter J Burt and Edward H Adelson. The laplacian pyramid as a compact image code. In *Readings in computer vision*, pages 671–679. Elsevier, 1987. [5](#)
- [8] Kartik Chandra, Tzu-Mao Li, Joshua Tenenbaum, and Jonathan Ragan-Kelley. Designing perceptual puzzles by differentiating probabilistic programs. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Proceedings (SIGGRAPH ’22 Conference Proceedings)*, 2022. [2](#)
- [9] Pascal Chang, Jingwei Tang, Markus Gross, and Vinicius C. Azevedo. How i warped your noise: a temporally-correlated noise prior for diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024. [4](#)
- [10] Ziyang Chen, Daniel Geng, and Andrew Owens. Images that sound: Composing images and sounds on a single canvas. *Neural Information Processing Systems (NeurIPS)*, 2024. [2](#)
- [11] Ming-Te Chi, Tong-Yee Lee, Yingge Qu, and Tien-Tsin Wong. Self-animating images: Illusory motion using repeated asymmetric patterns. *ACM Transactions on Graphics (SIGGRAPH 2008 issue)*, 27(3):62:1–62:8, 2008. [2](#)
- [12] Hung-Kuo Chu, Wei-Hsin Hsu, Niloy J. Mitra, Daniel Cohen-Or, Tien-Tsin Wong, and Tong-Yee Lee. Camouflage images. *ACM Trans. Graph.*, 29(4), 2010. [2](#)
- [13] James H. Clark. Hierarchical geometric models for visible surface algorithms. *Commun. ACM*, 19(10):547–554, 1976. [3](#), [4](#)
- [14] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. [4](#)
- [15] Jon P Ewins, Marcus D Waller, Martin White, and Paul F Lister. Mip-map level selection for texture mapping. *IEEE Transactions on Visualization and Computer Graphics*, 4(4): 317–329, 1998. [3](#)
- [16] Yue Feng, Vaibhav Sanjay, Spencer Lutz, Badour AlBahar, Songwei Ge, and Jia-Bin Huang. Illusion3d: 3d multiview illusion with 2d diffusion priors. 2024. [2](#), [17](#)
- [17] Daniel Geng, Inbum Park, and Andrew Owens. Factorized diffusion: Perceptual illusions by noise decomposition. In *European Conference on Computer Vision (ECCV)*, 2024. [2](#)
- [18] Daniel Geng, Inbum Park, and Andrew Owens. Visual anagrams: Generating multi-view optical illusions with diffusion models. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#), [8](#), [15](#), [16](#), [17](#), [18](#)
- [19] David J. Heeger and James R. Bergen. Pyramid-based texture analysis/synthesis. In *Proceedings of the 22nd Annual Conference on Computer Graphics and Interactive Techniques*, page 229–238, New York, NY, USA, 1995. Association for Computing Machinery. [3](#)
- [20] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. [3](#)
- [21] Kai-Wen Hsiao, Jia-Bin Huang, and Hung-Kuo Chu. Multi-view wire art. *ACM Trans. Graph.*, 37(6), 2018. [2](#)
- [22] Baltrusaitis Jurgis and WJ Strachan. *Anamorphic Art*. Abrams, New York, 1977. [2](#)
- [23] Jaihoon Kim, Juil Koo, Kyeongmin Yeo, and Minhyuk Sung. Synctweedies: A general generative framework based on synchronized diffusions. *arXiv:2403.14370*, 2024. [3](#), [4](#), [5](#), [7](#), [15](#), [18](#)
- [24] Philip Kuchel. Anamorphoscopes: A visual aid for circle inversion. *The Mathematical Gazette*, 63:82, 1979. [1](#)
- [25] Monster Laabs. Controlnet qr code monster v2 for sd-1.5, 2023. [2](#)
- [26] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91–110, 2004. [3](#)
- [27] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11461–11471, 2022. [6](#)
- [28] Mizuho Nakajima and Yasushi Yamaguchi. Picture illusion by overlap. In *ACM SIGGRAPH 2004 Posters*, page 56, New York, NY, USA, 2004. Association for Computing Machinery. [2](#)
- [29] Jean François Nicéron. *La perspective curieuse*. P. Billaine, Paris, 1638. [1](#), [2](#), [7](#)
- [30] Aude Oliva, Antonio Torralba, and Philippe G. Schyns. Hybrid images. In *ACM SIGGRAPH 2006 Papers*, page 527–532, New York, NY, USA, 2006. Association for Computing Machinery. [2](#)
- [31] Marios Papas, Thomas Houit, Derek Nowrouzezahrai, Markus Gross, and Wojciech Jarosz. The magic lens: Refractive steganography. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia)*, 31(6), 2012. [2](#)

- [32] Maxine Perroni-Scharf and Szymon Rusinkiewicz. Constructing printable surfaces with view-dependent appearance. In *ACM SIGGRAPH 2023 Conference Proceedings*, New York, NY, USA, 2023. Association for Computing Machinery. 2
- [33] 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. 7
- [34] Anurag Ranjan and Michael J Black. Optical flow estimation using a spatial pyramid network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4161–4170, 2017. 3
- [35] Herbert E Robbins. An empirical bayes approach to statistics. In *Breakthroughs in Statistics: Foundations and basic theory*, pages 388–394. Springer, 1992. 3
- [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 4
- [37] Matthew Tancik. Illusion diffusion: optical illusions using stable diffusion. <https://github.com/tancik/Illusion-Diffusion>, 2023. 4, 6, 7, 8, 15, 17, 18
- [38] Ugleh. [https://www.reddit.com/r/StableDiffusion/comments/16ew9fz/spiral\\_town\\_different\\_approach\\_to\\_qr\\_monster/](https://www.reddit.com/r/StableDiffusion/comments/16ew9fz/spiral_town_different_approach_to_qr_monster/), 2023. 2
- [39] Jean-Louis Vaulezard. *Perspective cylindrique et conique, concave et convexe ou traité des apparences vueus par le moyen des miroirs*. J. Jacquin, Paris, 1630. 2
- [40] Xiaojuan Wang, Janne Kontkanen, Brian Curless, Steve Seitz, Ira Kemelmacher, Ben Mildenhall, Pratul Srinivasan, Dor Verbin, and Aleksander Holynski. Generative powers of ten. *arXiv preprint arXiv:2312.02149*, 2023. 3
- [41] Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion null-space model. *The Eleventh International Conference on Learning Representations*, 2023. 6
- [42] Lance Williams. Pyramidal parametrics. *SIGGRAPH Comput. Graph.*, 17(3):1–11, 1983. 3, 4
- [43] Jiwen Yu, Yinhuai Wang, Chen Zhao, Bernard Ghanem, and Jian Zhang. Freedom: Training-free energy-guided conditional diffusion model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23174–23184, 2023. 6
- [44] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. 2
- [45] Qinsheng Zhang, Jiaming Song, Xun Huang, Yongxin Chen, and Ming yu Liu. Diffcollage: Parallel generation of large content with diffusion models. In *CVPR*, 2023. 3

# LookingGlass: Generative Anamorphoses via Laplacian Pyramid Warping

## Supplementary Material

### A. Laplacian Pyramid Warping

In this section, we provide additional details about Sec. 4.2. In Section A.1 and A.2, we illustrate the forward and inverse warping with the example of the conic mirror, and cover some implementation subtleties for the inverse operation. Section A.3 explains how we properly blend pyramid levels in the case of partial views. Lastly, we provide pseudo-code for Laplacian Pyramid Warping in Section A.4.

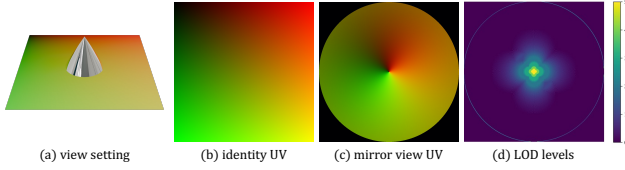


Figure 10. **Conic mirror view.** Rendering the scene (a) with a top down camera yields the UV map (c) for the conic mirror case. Using Eq. (9), the associated level-of-detail map (d) is computed.

#### A.1. Forward Warping

We consider the conic mirror example and its corresponding warping function (Fig. 10). The forward warping operation was explained in Sec. 4.2. Figure 11.b) illustrates the result of applying forward warping to an image. In Figure 14.b), we visualize the pyramid levels in the identity view, highlighting in white the pixels used in the forward warping. For instance, pixels closer to the cone’s center are mapped to an outer ring in the identity view. Because the view function is locally more compressed for these pixels, they are sampled from a higher level of the pyramid (*i.e.*, lower resolution).

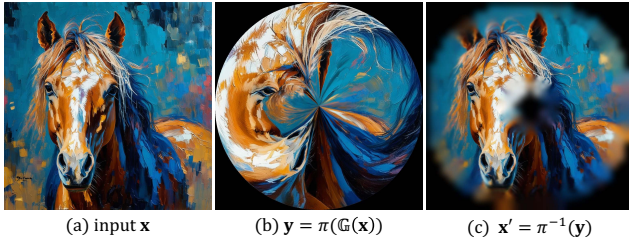


Figure 11. **Warp example.** Given an image (a), we warp it to the conic mirror view (b), and warp back to the original view (c). Our multi-scale approach warps values to different frequency bands based on the local compression of the view function  $\pi$ .

#### A.2. Inverse Warping

Figure 11.c) shows the result of warping the transformed image back to the identity view: pixels closer to the cen-

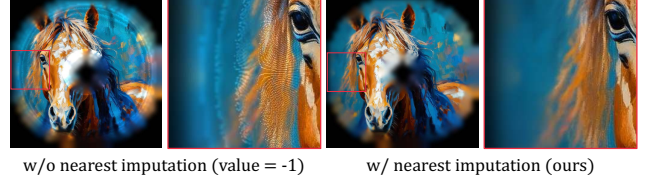


Figure 12. **Imputation.** We use nearest neighbor imputation for pixels in  $\mathbb{P}(\mathbf{x})$  that has no colors (*i.e.* did not receive gradients). We consider images with values in  $[-1, 1]$ .

ter of the cone in the warped view naturally map back to a lower frequency band in the identity view, occupying larger regions at the boundary in the reconstructed image. Rigorously speaking, the output of the inverse warping is a Laplacian pyramid (see Fig. 14.g), not an image.

**Implementation details.** We provide a more detailed illustration of the inverse operation in Figure 16. The subfigures show:

- Starting from a dummy pyramid  $\mathbb{P}(\mathbf{x}^0)$ , we perform a Laplacian-to-Gaussian pyramid conversion to get  $\mathbb{G}(\mathbf{x}^0)$  (*i.e.* *reparameterization*). Forward warping is then applied to obtain a warped dummy image  $\tilde{\mathbf{y}}$ .
- An  $L_2$  loss is computed between  $\tilde{\mathbf{y}}$  and  $\mathbf{y}$ , the image we wish to warp back. The gradients populate each level of the dummy pyramid, yielding  $\mathbb{P}(\mathbf{x})$ . The reparameterization ensures that gradients from each level flows to all levels above it.
- Lastly, we extract the final Laplacian pyramid  $\mathbb{L}(\mathbf{x})$  from  $\mathbb{P}(\mathbf{x})$  following

$$\mathbb{L}(\mathbf{x}) = \mathbb{M}(\mathbf{x}) \odot (\mathbb{P}^*(\mathbf{x}) - \mathbf{U}(\mathbf{D}(\kappa(\mathbb{P}^*(\mathbf{x}))))), \quad (11)$$

where  $\mathbb{M}$  is a pyramid of binary masks indicating pixels that have gradients,  $\mathbb{P}^*$  is the result of imputing missing values in  $\mathbb{P}$  with nearest color (more details below).

**Imputation.** It is important to impute the pixels that received no gradients with meaningful colors. Otherwise, if left black, the pyramid computation will pick up these discontinuities at the mask boundaries as high-frequency details, leading to incorrect values in those regions. We opt for a simple nearest neighbor imputation that fills pixels that have no gradients with the nearest pixel value. Figure 12 illustrates the difference between no imputation (default value 0) and our nearest imputation.



**Comparison with baselines.** In Figure 13, we compare our Laplacian-based inverse warping with standard warping using nearest or bilinear interpolation. Because the view function has extreme compression around the cone center, this translates to missing values on the outer ring when warping back with nearest or bilinear, as some pixels in the identity view are not sampled during forward warping due to the compression. We showed in Fig. 6 that this can lead to artifacts in the generated image.

It is worth noting that another way to avoid missing values could be to do a *forward* warping with the *inverse* UV map. However, the inverse map is usually not trivial to compute. Additional problems arise when the mapping is not bijective or has discontinuities, which can be quite tricky to solve in a robust way. In contrast, our method automatically handles these cases, and only requires the forward UV map, which is easily obtainable through simple rendering of the desired view.

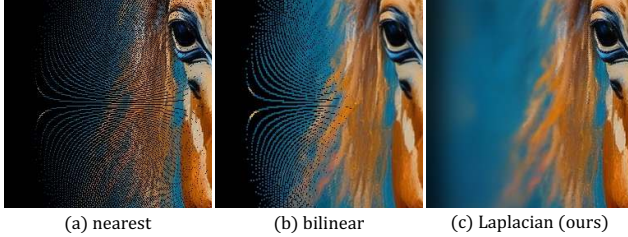


Figure 13. **Backward warping baseline comparison.** In regions where the view function is compressing, standard interpolation like nearest (a) or bilinear interpolation (b) creates holes, while our Laplacian Pyramid Warping ensures smooth results (c).

### A.3. Pyramid Blending

As explained in Sec. 4.2, special care needs to be taken when blending the pyramids.

**Blending with partial views.** In standard Laplacian Pyramid Blending, the blended pyramid is obtained by averaging each level of the input pyramids. Given two pyramids  $\mathbb{L}^0, \mathbb{L}^1$ , the level  $k$  of the blended pyramid  $\mathbb{L}$  is given by:

$$\mathbb{L}_k = \frac{1}{2} (\mathbb{L}_k^0 + \mathbb{L}_k^1). \quad (12)$$

In our case, Laplacian pyramids come from inverse warping of views, and might look like Fig. 11.c), with missing values. In this case, the average should only be computed over defined pixels. Assuming each level  $\mathbb{L}_k$  is associated with a binary mask  $\mathbb{M}_k$ , the blending is:

$$\mathbb{L}_k = \frac{1}{\mathbb{M}_k^0 + \mathbb{M}_k^1} (\mathbb{M}_k^0 \odot \mathbb{L}_k^0 + \mathbb{M}_k^1 \odot \mathbb{L}_k^1). \quad (13)$$

In practice, we map missing values to `torch.nan`, and use `torch.nanmean()` to perform averaging.

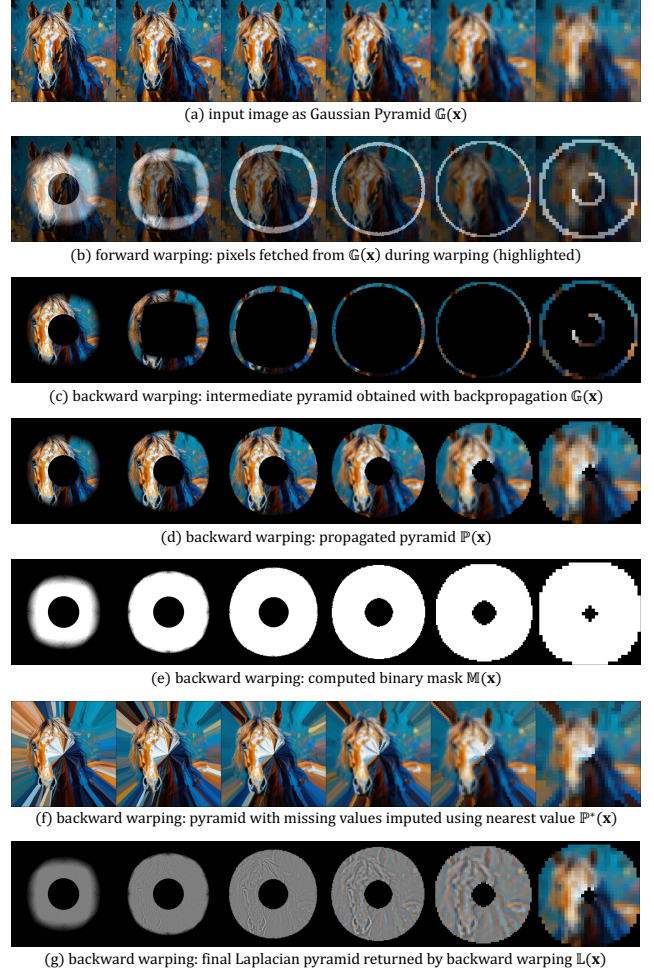


Figure 14. **Intermediate pyramids.** We visualize some intermediate pyramids that appear in forward (a, b) and backward warping (c-g). Refer to the text and Fig. 16 for more context.

**Detail-preserving averaging.** Another issue with averaging in general is that it reduces variance and washes out details. While this is already improved with the Laplacian pyramid, averaging still leads to the loss of sharp details. Given two pixel values  $x, y \in [-1, 1]$ , we define a normal averaging `avg` and a value-weighted averaging `vavg`:

$$\text{avg}(x, y) = \frac{1}{2}(x + y), \quad \text{vavg}(x, y) = \frac{|x|x + |y|y}{|x| + |y|} \quad (14)$$

The value-weighted average will give more weights to extreme pixel values, hence preserving better the value range. In the results shown, we linearly interpolate between the two types of averaging through a parameter  $\alpha \in [0, 1]$ :

$$z = \text{avg}(x, y) + \alpha(\text{vavg}(x, y) - \text{avg}(x, y)). \quad (15)$$

Figure 15 shows the effect of  $\alpha$  for two examples with the cylinder mirror. When standard averaging is used ( $\alpha =$

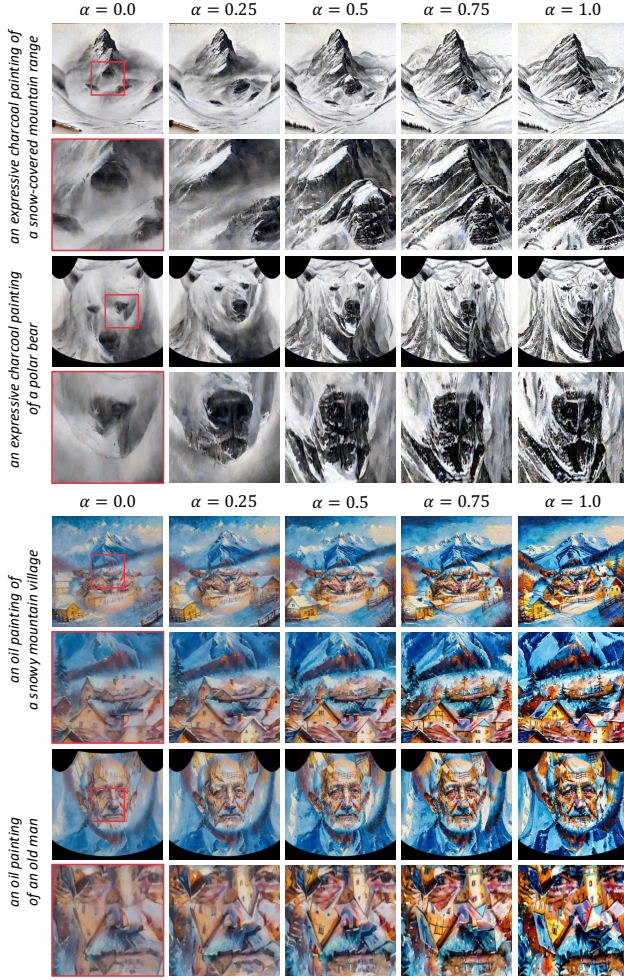


Figure 15. **Effects of parameter  $\alpha$ .** The parameter affects how views are merged together at each level of the pyramid.  $\alpha = 0$  corresponds to standard average, while  $\alpha = 1$  gives more weights to pixels with extreme values, preserving the details in all the views.

0), the image looks blurrier. However, when value-weighted average is used ( $\alpha = 1$ ), all details from both views are preserved, creating very saturated, over-sharpened results. In most of our results, we opt for  $\alpha \in [0.25, 0.5]$ .

#### A.4. Pseudo-code

We provide a pseudo-code for our novel Laplacian Pyramid warping. Some notes:

- 1.8: `_compute_lod_level` returns the LOD level map as a  $(H, W)$  tensor using Equation (9);
- 1.29: `pyrStack` takes a pyramid, upsamples all levels to highest resolution, and stacks them along a dimension;
- 1.73: `impute_with_nearest` fills missing values with nearest ones in the image given the mask.

```

1 def _get_grid(warp, maxLOD):
2     """
3     Takes in numpy array warp of shape (1, 3, 1024, 1024)
4     """
5     # compute lod and normalize to (-1, 1)
6     mapping = 1024 * warp[:, :2]
7     lod = _compute_lod_level(mapping, maxLOD=maxLOD)
8     lod = 2 * lod / maxLOD - 1.0 # (h, w)
9     lod = lod.unsqueeze(0).unsqueeze(-1)
10
11     # normalize uv to (-1, 1)
12     grid = torch.tensor(warp[:, :2]).permute(0, 2, 3, 1)
13     mask = torch.all(grid == 0.0, dim=-1, keepdim=True)
14     mask = mask.expand(-1, -1, -1, 2)
15     grid[mask] = torch.nan # replace undefined with NaN
16     grid = 2 * grid - 1
17
18     # combine into a 3D coordinate grid (lod, u, v)
19     grid = torch.cat([lod, grid], -1).unsqueeze(1)
20     return grid # (1, 1, dim, dim, 3)
21
22
23
24 def view(lp, warp, leveln):
25     # get 3d coordinate grid
26     grid = _get_grid(warp, maxLOD=leveln-1)
27
28     # sample values
29     layers = pyrStack(lp, dim=-1)
30     new_im = F.grid_sample(
31         layers,
32         grid,
33         mode='nearest',
34         padding_mode="zeros",
35         align_corners=True,
36     ).squeeze(2)
37
38     # replace by NaN where image is 0
39     new_im[new_im == 0.0] = torch.nan
40     new_im = torch.nanmean(new_im, 0).half()
41     return new_im
42
43
44 def inverse_view(im, warp, leveln):
45     c, h, w = im.shape
46     grid = _get_grid(warp, maxLOD=leveln - 1)
47     with torch.enable_grad():
48         # create an empty pyramid
49         opt_var = torch.zeros(1, c + 1, h, w)
50         opt_var = LaplacianPyramid(opt_var, leveln)
51         for lvl in opt_var:
52             lvl.requires_grad_()
53
54         # convert to Gaussian pyramid
55         opt_gp = Laplacian2Gaussian(opt_var)
56         target = torch.cat([im, torch.ones_like(im[:1])])
57         layers = pyrStack(opt_gp, -1).float()
58         warped = F.grid_sample(
59             layers,
60             grid,
61             mode='nearest',
62             padding_mode="zeros",
63             align_corners=True,
64         ).squeeze(2)
65         loss = 0.5 * ((warped - target) ** 2).sum()
66         loss.backward()
67         result = [-1.grad.detach() for l in opt_var]
68         result = [(r[:, :c] / r[:, -1:]) for r in result]
69
70     # extract laplacian pyramid
71     for k in range(leveln - 1):
72         mask = torch.isnan(result[k])
73         imputed = impute_with_nearest(result[k], ~mask)
74         result[k] = imputed - pyrDown(pyrUp(imputed))
75         result[k][mask] = torch.nan
76
77     return result
78

```



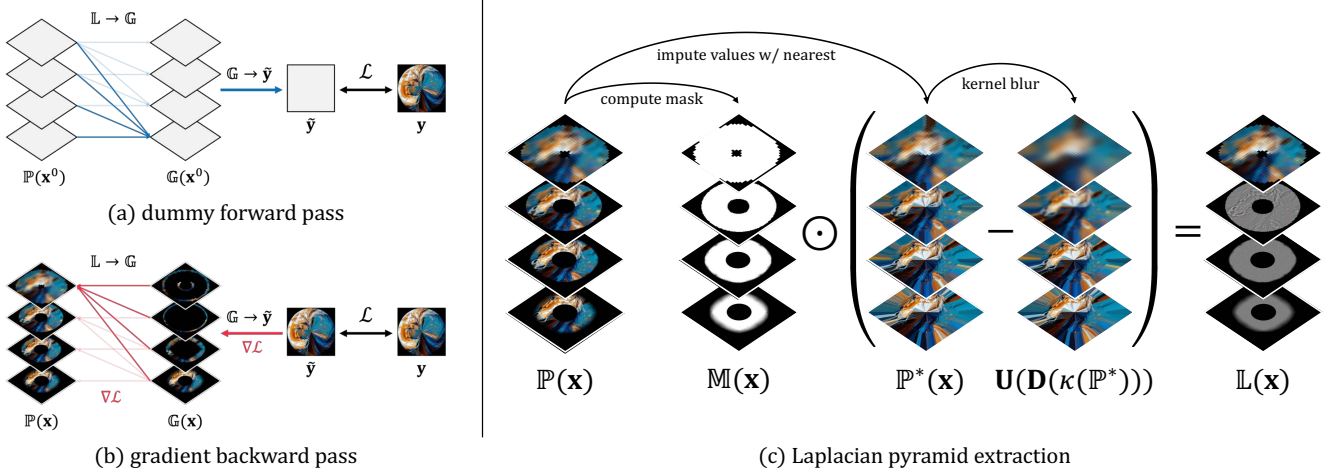


Figure 16. **Detailed look at inverse warping.** We provide an illustration to Equation (10), explaining how to obtain the final Laplacian pyramid using back-propagation, as mentioned in the main paper.

## B. Additional Ablations

We provide more qualitative ablations to show the effects of Laplacian warping, view prioritization, and time travel.

### B.1. Prioritizing a Single View

In Section 4.3, we proposed to let the last  $x\%$  of the denoising process only denoise a single one of the views. Figure 17 shows the effect of varying values of  $x$  for an example with  $90^\circ$  rotation. As the ratio increases, details from the first view (“a snowy mountain village”) dominate over the second view (“a horse”). This is a useful parameter to control the trade-off between views.

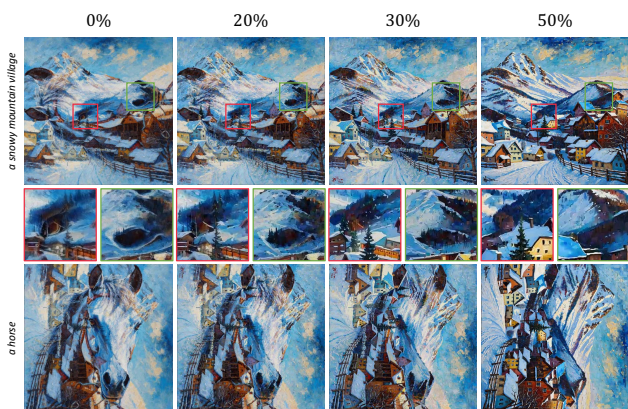


Figure 17. **Prioritizing a single view.** We show the effect of prioritizing the first view (“a snowy mountain village”) over the second (“a horse”) in the example of  $90^\circ$  rotation.

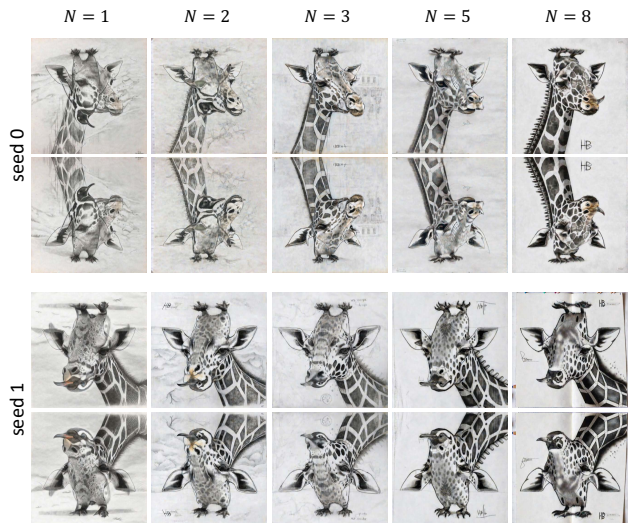


Figure 18. **Time travel.** We show the effect of time traveling for the flip example (“a giraffe head” / “a penguin”) with two different sampled seeds. Time traveling is only applied for timesteps between 20% and 80%, by repeating each step  $N$  times.

### B.2. Time travel

We show in Figure 18 that time traveling is an effective strategy for improving image consistency and blending between views. Without time traveling ( $N = 1$ ), the *penguin* and the *giraffe head* seem like two independent entities overlaid on top of each other in the image. As the repeating number  $N$  increases, the two views align better, resulting in a sharper image with more coherent details.

Obviously, runtime scales linearly with the repeating number. In Burgert et al. [4], image quality and blend-



ing quality are entangled, and improving with the number of optimization steps. Here, time traveling is only responsible for blending quality, but is independent of the generated image quality. We believe this provides a more intuitive control parameter than the number of optimization steps.

Naturally, better blending still yields better overall quality, which explains the lower FID in Table 2. However, it seems to come at the cost of slightly worse prompt alignment. We hypothesize that blending better means making more compromise between views, which in turn makes it harder to match the prompts as well.

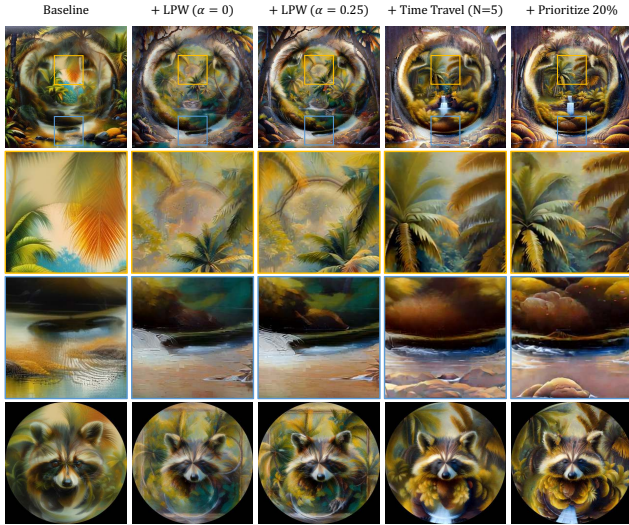


Figure 19. **Baseline ablation.** We consider the cone mirror example with two views (“a tropical jungle forest”-“a raccoon”). We show successively the effect of Laplacian Pyramid Warping (LPW), value-weighted average ( $\alpha$ ), time traveling, and single-view prioritization.

### B.3. Comparison with Baseline

Lastly, we show in Figure 19 the effects of these different features. Laplacian Pyramid Warping addresses the artifacts at the edge of the ring, visible in the baseline. Value-weighted average recovers sharp details lost in standard averaging. Time traveling ensures a smoother blending, while single-view prioritization adds back the final details in the identity view without destroying the second view.

### C. Quantitative Evaluation Details

We provide additional details regarding the quantitative evaluation in this section.

**Dataset generation.** Similar to Geng *et al.* [18], we use a custom list of 50 prompt pairs in the form of [`<style>`, `<prompt1>`, `<prompt2>`] for our quantitative evaluation. These are a mix between prompt pairs from Visual

Anagrams’ paper [18], and from querying ChatGPT. For each method, we generate 10 samples per prompt pair. This results in 500 pairs of images, *i.e.* 1k images per method. For FID/KID computation, we generated a reference dataset comprised of 3.2k images from SD3 and 3.2k images from SD3.5 following the same prompts (only single view).

**Comparing with prior work.** For Visual Anagrams [18], we used the original codebase from Geng *et al.*: [https://github.com/dangeng/visual\\_anagrams](https://github.com/dangeng/visual_anagrams). Results are generated with DeepFloyd IF [1] in two stages, and subsequently up-sampled to  $1024 \times 1024$  using Stable Diffusion x4 Upscaler, as provided by the code.

For Burgert *et al.* [4], we used the available codebase at (<https://github.com/RyannDaGreat/Diffusion-Illusions>) and added a function to rotate the inner circle of an image by  $135^\circ$ . Each image is generated with 1000 optimization steps using Stable Diffusion v1.4, the default model from their codebase.

Lastly, we set the hyperparameters of our pipeline to replicate the original implementation of Tancik [37] and SyncTweedies [23]. The results are generated with Stable Diffusion 3.5 Medium, as in our method.

**Runtime.** We generate our results with Stable Diffusion 3.5 Medium on a Nvidia GeForce RTX 4090 GPU. With 30 steps of inference and time traveling between 20% and 80% repeating 2 times, we can generate an image pair in  $\sim 80$ s.

Geng <i>et al.</i> [18]	Tancik <i>et al.</i> [37]	SyncTweedies [23]	Burgert <i>et al.</i> [4]	Ours SD 3.5
18.6s	17.2s	18.8s	176.0s	79.4s

Table 3. **Inference time comparison.**

### D. Using Other Latent Models

Most of our results were generated using Stable Diffusion 3.5. However, we also experimented with other models such as SD2.1 and SDXL.

Here, we show a simple experiment with *two identity views* associated with two distinct prompts. This way, we abstract out the VAE, as well as other problems that come from warping. The results are shown in Figure 20 for two pairs of prompts. Interestingly, even in such a simple setting, not all models behave the same: SD2.1 tend to generate images with poor quality, while SD3+ attempts to blend the two concepts by compositing them as much as possible. Curiously, SDXL seems to blend the concepts *semantically*: “a snowy mountain village” and “a horse” give horses *in* a village, while “people at a campfire” and “an old man” produce “old men at a campfire”.

While intriguing, this is not ideal for ambiguous images, as the goal is more to blend *spatially* the different views. Further investigation is needed to explain the discrepancy between SDXL and SD3+, possibly due to the switch from diffusion to flow matching or from U-Net to a transformer backbone. A deeper study is left for future work.

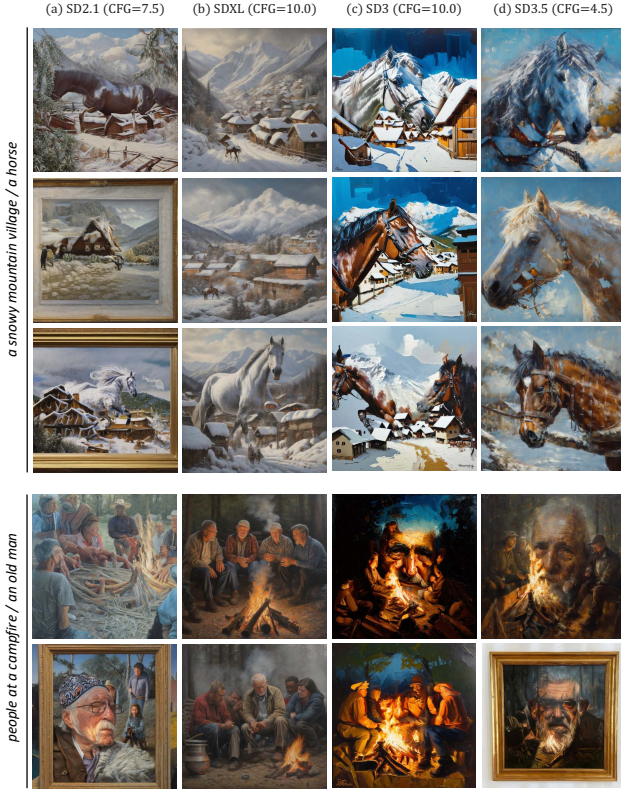


Figure 20. **Comparing different models for two-view setting.** We consider the case of two identity views, associated with two distinct prompts. While most models blend the two concepts *spatially*, SDXL seems to blend them *semantically*.

## E. User Study

As mentioned in Section 5.5, we conducted a user study over 27 participants to compare human preferences between our proposed method and existing prior work.

**Study structure.** The study consists of three sections, each corresponding to a different type of anamorphosis: cylindrical mirror, conic mirror, and Niceron’s lens. Each section begins with an example image and video demonstrating how the illusion works. Participants are then shown 10 different samples generated from 10 pairs of prompts. These prompts remain the same across all three sections to avoid bias. Additionally, the order in which the methods are presented is randomized for each sample. To ensure a fair

comparison, we first display the results in high resolution before asking users to rank them (see Fig. 21).

**Ranking criteria.** Participants are asked to provide a ranking of the five methods from 1 (best) to 5 (worst). At the beginning of the study, they are instructed to evaluate the images based on the following criteria:

- **Match both text prompts:** When viewed directly (resp. through a mirror or lens), the image should correspond to the first (resp. second) prompt.
- **Maintain the specified style:** The image should adhere to the given style prompt (e.g. painting, photograph).
- **Be high-quality:** The final image should be sharp, detailed, and visually appealing.

## F. Failed Experiments

**Relative negative prompting.** In Visual Anagrams [18], the authors note that including other views in the negative prompt does not improve substantially the visual quality. Moreover, prompt conflicts can arise. For example, the prompts “an oil painting of a village” and “an oil painting of a horse” both share the style descriptor “an oil painting of”, which should not be included in the negative prompt. Similarly, prompts like “a cat” and “a dog” share features such as fur. If these shared features appear in both the negative and positive prompts, it can lead to suboptimal results.

We experimented with a variant of this, which we dub *relative negative prompting*. The idea is to put in the negative prompt only the relative direction between the positive prompt (from the current view) and the negative prompt (from the other view). Thus, we subtract the positive prompt embedding from the negative one. Our experiments showed that this removes the need to manually select which part to keep for the negative prompt. In the example of the shared style, our difference vector cancels it out automatically, and the model is thus still able to generate the correct style. However, similar to Geng *et al.* [18], we did not observe substantial improvement using this method.

## G. Choosing Prompts & Failure Cases

The quality of the generation relies on choosing a good style and pair of prompts. Here are some tips we found for generating good anamorphoses:

1. a place or location (e.g. jungle, desert, library etc.) gives a lot of freedom to the composition and generally works well for the identity view;
2. the second view is generally seen through some mirror or a lens, which is smaller than the main image. For this view, easily recognizable subjects like animals or faces are good prompts in most cases;
3. artistic styles are more prone to produce good results than photorealistic styles;



Start of the section:  
example image pair  
and video showing  
the illusion type

Image sample pairs  
from five methods in  
high resolution for  
better viewing

Users are asked to  
rank the five samples  
based on **text-  
prompt adherence,**  
**style matching,** and  
**overall visual quality**

Part 1: Cylindrical Mirror (~5 min.)

In this first part, the hidden image is revealed when looking through a cylindrical mirror at an angle.

Each method consists of two images. The example video below shows how the two images are physically related to each other:

A paper collage of a mountain pass

A paper collage of a lion

Sample 1/10  
Take a look at this first sample in high resolution and answer below. Look at the text prompts, and the images from different methods A-E. A smaller version of the images will be available afterwards.

An oil painting of a sunset canyon with straight cliffs

An oil painting of a bull

A

B

C

D

E

Look at the text and images above and rank the methods from best to worst based on **text prompt adherence, style matching and overall visual quality.**

An oil painting of a sunset canyon with straight cliffs

An oil painting of a bull

A

B

C

D

E

1 (best)

2

3

4

5 (worst)

☐

☐

☐

☐

☐

Figure 21. Sample from the user study.

- styles with no colors (e.g. sketches, ink, marble sculpture) will generate better results when the two prompts have very different color palettes.

Our method is still prone to fail in certain cases. For example, the model can still cheat and put all the views in the image without properly blending them (see Figure 22).



Figure 22. **Failure case.** Similar to Geng *et al.*, our method can cheat and put both views in the image without properly mixing them. In this example, the shark from the cylinder mirror view can be seen in the sky behind the wind turbines.

## H. Concurrent Work

After our initial submission, a preprint titled “*Illusion3D: 3D Multiview Illusion with 2D Diffusion Priors*” [16] appeared on arXiv, addressing a similar problem. While their code is not available at the time of writing, we identified a few key differences between our method and theirs.

First, Illusion3D builds on optimization-based approaches like Burgert et al. [4], whereas our method improves upon feedforward techniques [18, 37], generating images in a single inference pass. Their approach replaces Score Distillation Sampling (SDS) [4] with Variational Score Distillation (VSD), which requires training a LoRA module during optimization. Due to the inherent limitations of score distillation methods, we believe our approach produces higher-quality images while supporting a broader range of styles, from artistic to photorealistic, as demonstrated in Section I.

A key advantage of Illusion3D, however, is its ability to generate full 3D structures, which our method does not support. That said, our approach can still generate 2D textures for mapping onto 3D surfaces, similar to their sphere or cube illusions. Lastly, we expect our method to have significantly lower generation time and memory requirements compared to Illusion3D.

## I. Additional Results

In the next pages, we show additional qualitative results for the three anamorphic views: cylinder mirror, conic mirror, and Nicéron’s lens. Please refer to the supplementary videos to see these anamorphoses in action. Table 4 shows additional quantitative evaluations.

	Method	$\mathcal{S} \uparrow$	$\mathcal{S}_{0.9} \uparrow$	$\mathcal{A} \uparrow$	$\mathcal{A}_{0.9} \uparrow$	$\mathcal{C} \uparrow$	$\mathcal{C}_{0.9} \uparrow$	FID $\downarrow$	KID $\downarrow$
Vertical Flip	Geng <i>et al.</i> [18]	0.325	0.362	0.306	0.340	0.695	0.786	149.24	0.057
	Tancik SD 3.5 [37]	0.328	0.367	0.306	0.349	0.693	0.806	132.52	0.049
	Burgert <i>et al.</i> [4]	0.303	0.347	0.281	0.324	0.679	0.778	219.84	0.115
	SyncTweedies [23]	0.323	0.360	0.302	0.341	0.707	0.801	132.62	0.054
	<b>LookingGlass (ours)</b>	<b>0.320</b>	<b>0.358</b>	<b>0.297</b>	<b>0.338</b>	<b>0.680</b>	<b>0.779</b>	<b>124.67</b>	<b>0.049</b>
135° Rotation	Geng <i>et al.</i> [18]	0.284	0.340	0.262	0.308	0.563	0.652	293.00	0.254
	Tancik SD 3.5 [37]	0.203	0.225	0.194	0.216	0.498	0.509	439.35	0.545
	Burgert <i>et al.</i> [4]	0.301	0.347	0.280	0.326	0.654	0.760	223.21	0.120
	SyncTweedies [23]	0.308	0.354	0.283	0.335	0.647	0.753	166.03	0.083
	<b>LookingGlass (ours)</b>	<b>0.319</b>	<b>0.357</b>	<b>0.295</b>	<b>0.338</b>	<b>0.666</b>	<b>0.767</b>	<b>129.74</b>	<b>0.055</b>
Cylindrical Mirror	Geng <i>et al.</i> [18]	0.190	0.228	0.171	0.198	0.506	0.546	285.23	0.216
	Tancik SD 3.5 [37]	0.189	0.225	0.171	0.198	0.505	0.547	284.97	0.215
	Burgert <i>et al.</i> [4]	0.285	0.334	0.261	0.304	0.706	0.795	229.65	0.138
	SyncTweedies [23]	0.285	0.348	0.241	0.284	0.673	0.763	138.69	0.082
	<b>LookingGlass (ours)</b>	<b>0.307</b>	<b>0.360</b>	<b>0.272</b>	<b>0.318</b>	<b>0.698</b>	<b>0.810</b>	<b>130.27</b>	<b>0.070</b>

Table 4. **Additional quantitative comparison.** We additionally assess image-prompt alignment using CLIP similarity score  $\mathcal{S}$  on all three transformations evaluated in the main paper. While all methods achieve comparable results for the vertical flip, LookingGlass surpasses previous approaches on more complex transformations, including anamorphoses.

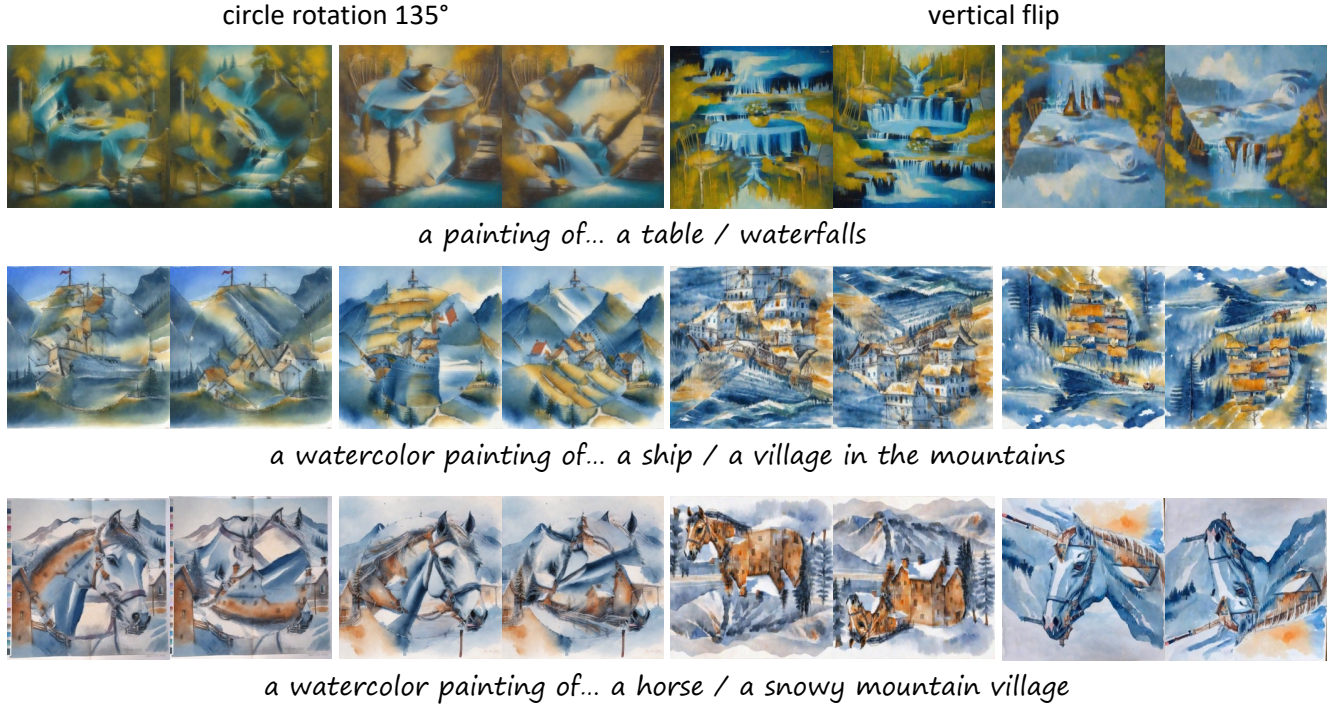
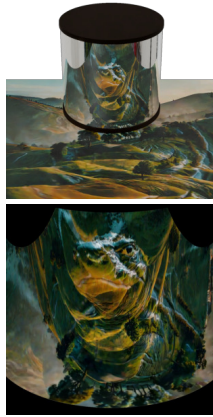
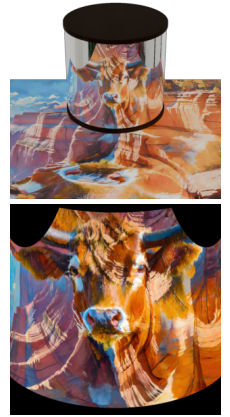
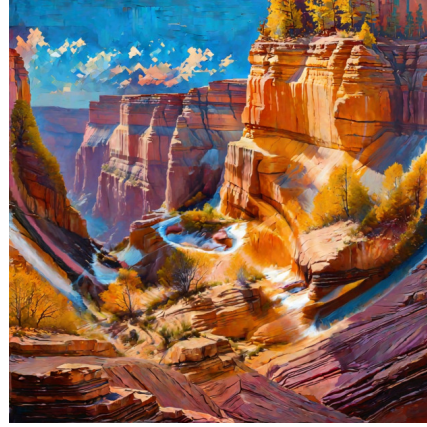


Figure 23. **2D transform results.** Here are some generated results for the two 2D transforms: vertical flip, and 135° rotation (not supported by Geng *et al.* [18]).





*a cinematic rendering of  
rolling hills in golden light / turtle*



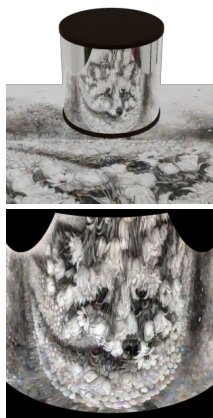
*an oil painting of  
sunlit canyon with straight cliffs / bull*



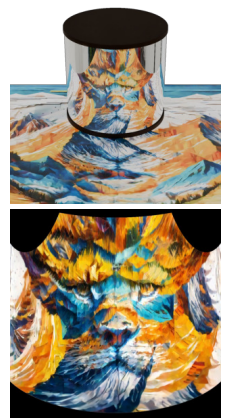
*an oil painting of  
desert dunes at sunset / jellyfish*



*a clay sculpture of  
cobblestone street / cheetah*



*a charcoal drawing of  
flower garden with rows of tulips / fox*



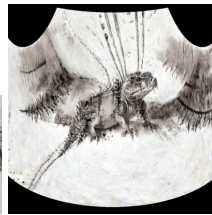
*a paper collage of  
mountain pass / lion*

Figure 24. **Cylinder mirror anamorphosis.** In this figure and the two following ones, we show additional results for the cylinder mirror example. Each example contains the identity view, the mirror view as predicted by the flow model, and a rendering of the actual physical setting to validate our examples. Kindly refer to the supplementary videos to see these results in action.





*an ink wash drawing of straight ski tracks*



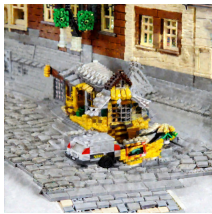
*an ink wash drawing of iguana*



*a watercolor painting of desert plateau*



*a watercolor painting of seal*



*a LEGO model of cobblestone street*



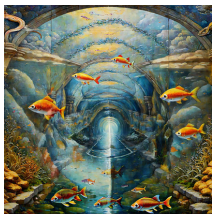
*a LEGO model of cheetah*



*a clay sculpture of aquarium tunnel*



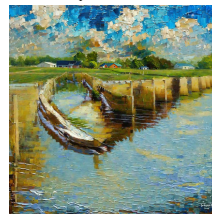
*a clay sculpture of polar bear*



*a fresco painting of aquarium tunnel*



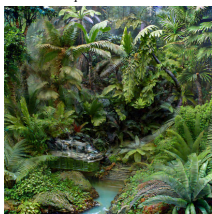
*a fresco painting of knight's helmet*



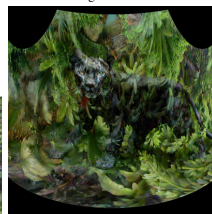
*an oil painting of straight levee dividing water*



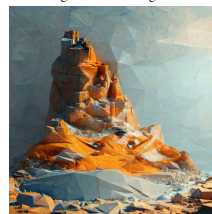
*an oil painting of otter*



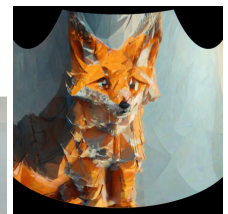
*a diorama of dense tropical rainforest*



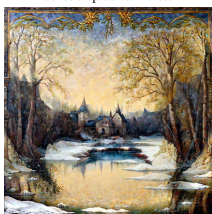
*a diorama of panther*



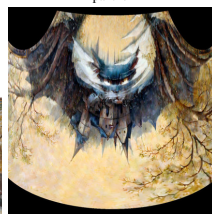
*a low-poly model of sunlit rocky outcrop*



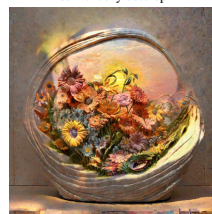
*a low-poly model of fox*



*a fresco painting of mirror-like frozen pond*



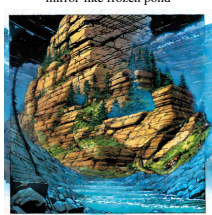
*a fresco painting of bat*



*a clay sculpture of flower meadow at sunrise*



*a clay sculpture of iguana*



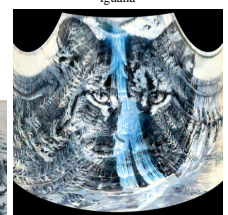
*a comic book panel of sunlit rocky outcrop*



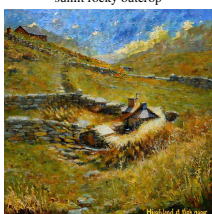
*a comic book panel of turtle*



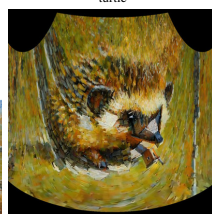
*a lithograph of frozen waterfall*



*a lithograph of cougar*



*an oil painting of highland moor with stone walls*



*an oil painting of hedgehog*

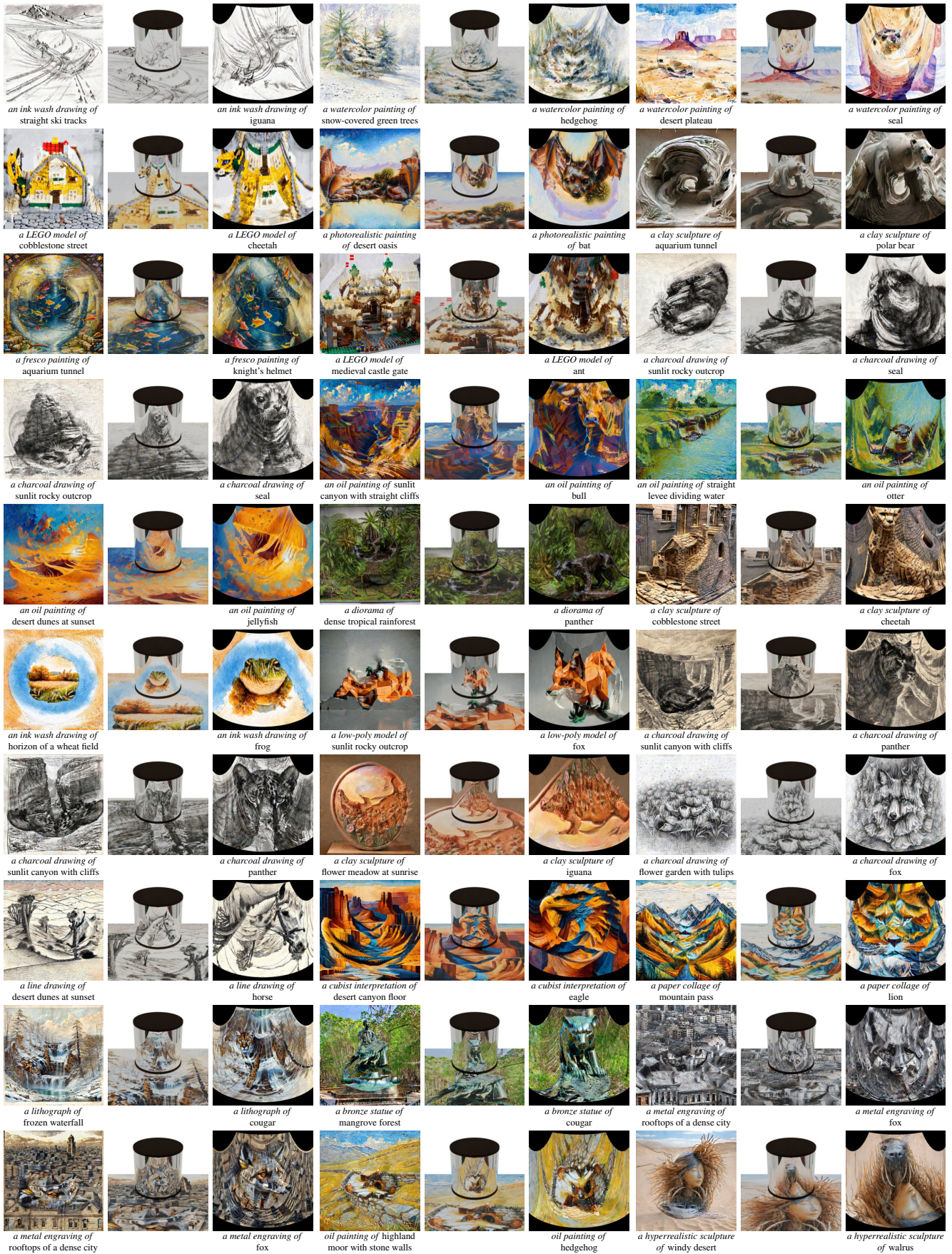


*a hyperrealistic sculpture of windy desert*



*a hyperrealistic sculpture of walrus*





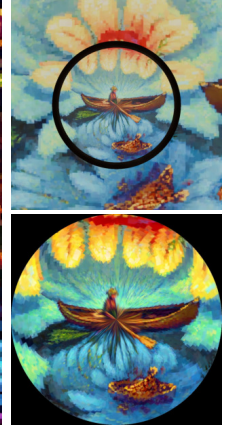




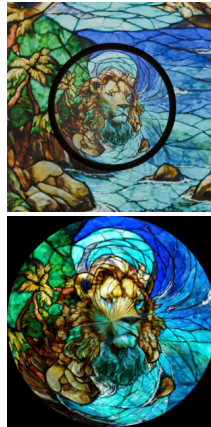
*a clay sculpture of  
aquarium tunnel / polar bear*



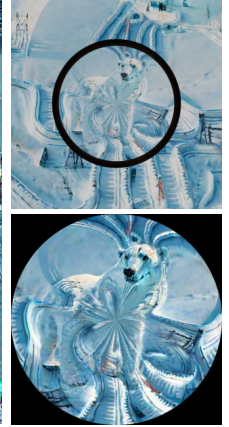
*a pixel art version of  
flower petals close-up / canoe*



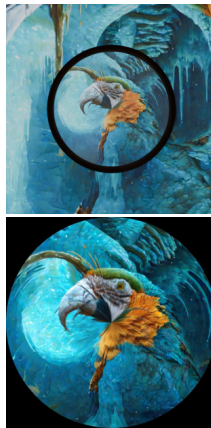
*a stained glass depiction of  
straight coastline / lion*



*a 3D rendering of  
straight ski tracks / polar bear*



*a cinematic rendering of  
icy cave with stalactites / parrot*



*a pointillism painting of  
straight canal lined with trees / butterfly*

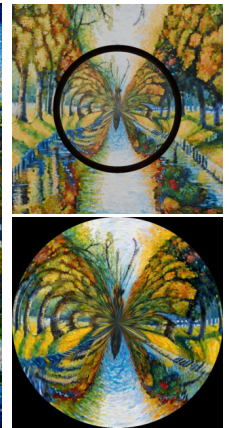


Figure 27. **Conic mirror anamorphosis.** In this figure and the two following ones, we show additional results for the conic mirror example. Each example contains the identity view, the mirror view as predicted by the flow model, and a rendering of the actual physical setting from the top to validate our examples. Kindly refer to the supplementary videos to see these results in action.

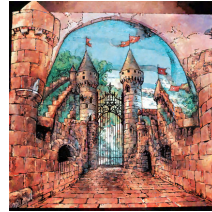




*a clay sculpture of aquarium tunnel*



*a clay sculpture of polar bear*



*a comic book panel of medieval castle gate*



*a comic book panel of octopus*



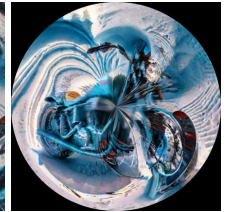
*a clay sculpture of cobblestone street*



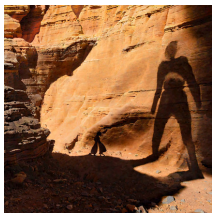
*a clay sculpture of cheetah*



*a hyperrealistic sculpture of straight ski tracks*



*a hyperrealistic sculpture of motorcycle*



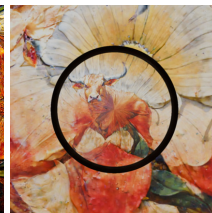
*a shadow puppet silhouette of desert canyon floor*



*a shadow puppet silhouette of otter*



*a fresco painting of flower petals close-up*



*a fresco painting of bull*



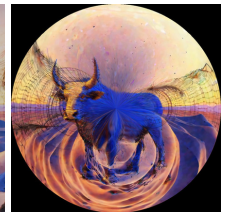
*a hyperrealistic sculpture of icy cave with stalactites*



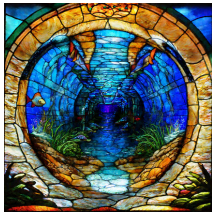
*a hyperrealistic sculpture of fox*



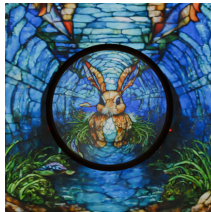
*a wireframe rendering of desert dunes at sunset*



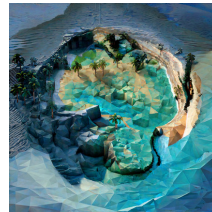
*a wireframe rendering of bull*



*a stained glass depiction of aquarium tunnel*



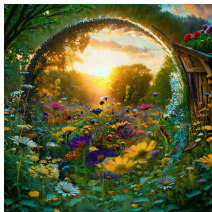
*a stained glass depiction of rabbit*



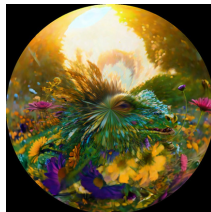
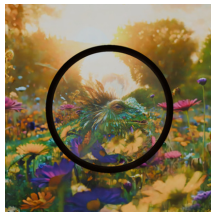
*a low-poly model of straight coastline*



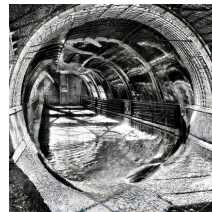
*a low-poly model of wolf*



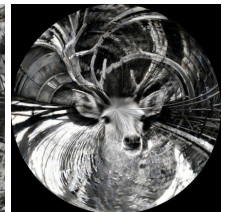
*a cinematic rendering of flower meadow at sunrise*



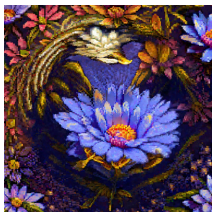
*a cinematic rendering of iguana*



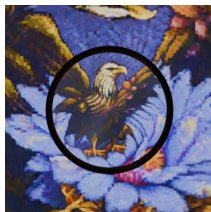
*a black-and-white photo of aquarium tunnel*



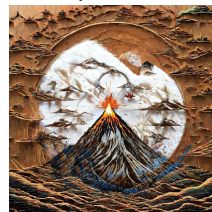
*a black-and-white photo of reindeer*



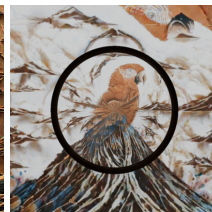
*a 16-bit sprite of flower petals close-up*



*a 16-bit sprite of eagle*



*a cut-paper silhouette of volcanic crater*

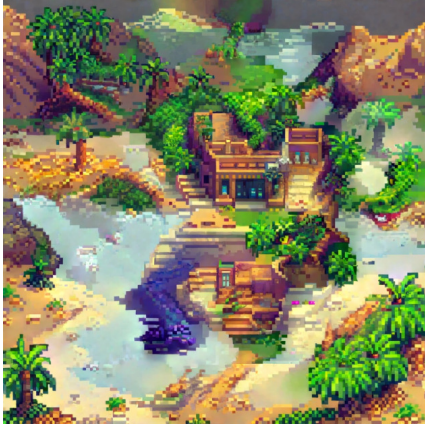


*a cut-paper silhouette of macaw*

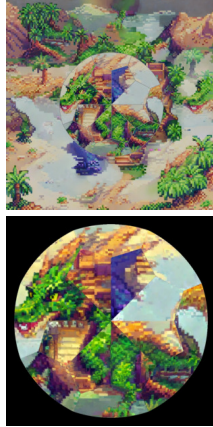




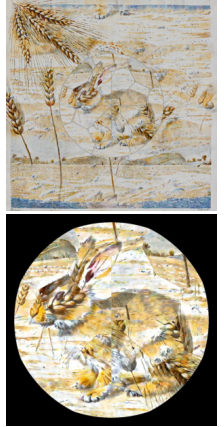




*a 16-bit sprite of  
desert oasis / dragon*



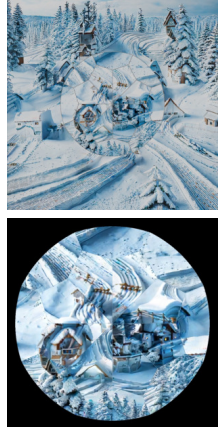
*a lithograph of  
horizon of a wheat field / rabbit*



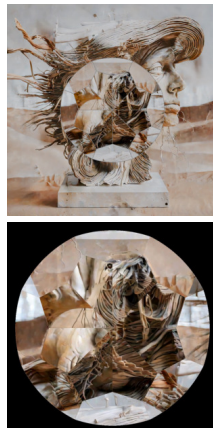
*a vintage poster of  
dense tropical rainforest / deer*



*a hyperrealistic sculpture of  
straight ski tracks / motorcycle*



*a hyperrealistic sculpture of  
windy desert / walrus*



*a pencil sketch of  
harbor pier / lion*

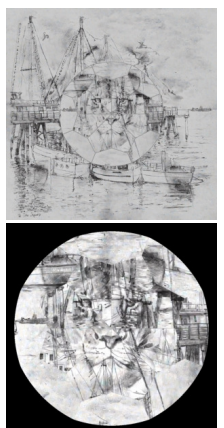
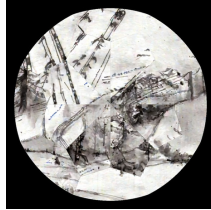


Figure 30. **Nicéron's lens anamorphosis.** In this figure and the two following ones, we show additional results for the lens example. Each example contains the identity view, the lens view as predicted by the flow model, and a rendering of the actual image through the lens to validate our examples. Kindly refer to the supplementary videos to see these results in action.

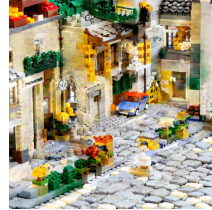




*an ink wash drawing of straight ski tracks*



*an ink wash drawing of iguana*



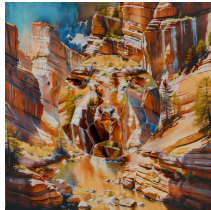
*a LEGO model of cobblestone street*



*a LEGO model of cheetah*



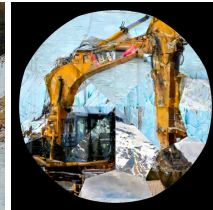
*an oil painting of sunlit canyon with straight cliffs*



*an oil painting of bull*



*a photo of glacier valley*



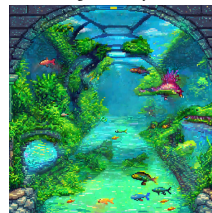
*a photo of excavator*



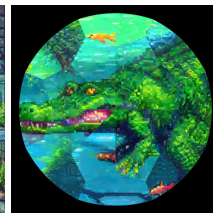
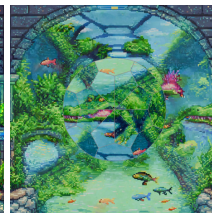
*a low-poly model of icebergs in the ocean*



*a low-poly model of dragon*



*a pixel art version of aquarium tunnel*



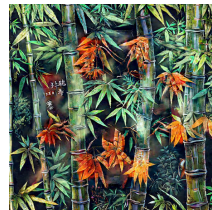
*a pixel art version of alligator*



*a digital illustration of beach with straight shoreline & waves*



*a digital illustration of fox*



*a chalkboard drawing of bamboo forest*



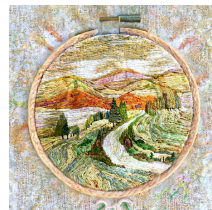
*a chalkboard drawing of lobster*



*a low-poly model of straight coastline*



*a low-poly model of wolf*



*an embroidered version of rolling hills in golden light*



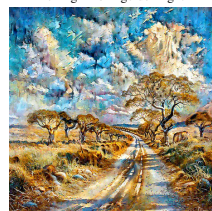
*an embroidered version of frog*



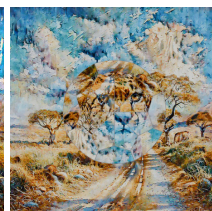
*a pencil sketch of castle walls with battlements*



*a pencil sketch of deer*



*a pastel artwork of straight dirt road through a savannah*



*a pastel artwork of cheetah*



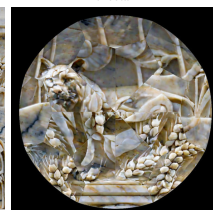
*a cinematic rendering of ocean waves*



*a cinematic rendering of horse*



*a marble carving of horizon of a wheat field*



*a marble carving of cougar*



