# SIDAR: Synthetic Image Dataset for Alignment & Restoration

**Monika Kwiatkowski**
Computer Vision & Remote Sensing
Technische Universität Berlin
m.kwiatkowski@tu-berlin.de

**Simon Matern**
Computer Vision & Remote Sensing
Technische Universität Berlin
s.matern@tu-berlin.de

**Olaf Hellwich**
Computer Vision & Remote Sensing
Technische Universität Berlin
olaf.hellwich@tu-berlin.de

## ABSTRACT

Image alignment and image restoration are classical computer vision tasks. However, there is still a lack of datasets that provide enough data to train and evaluate end-to-end deep learning models. Obtaining ground-truth data for image alignment requires sophisticated structure-from-motion methods or optical flow systems that often do not provide enough data variance, i.e., typically providing a high number of image correspondences, while only introducing few changes of scenery within the underlying image sequences. Alternative approaches utilize random perspective distortions on existing image data. However, this only provides trivial distortions, lacking the complexity and variance of real-world scenarios. Instead, our proposed data augmentation helps to overcome the issue of data scarcity by using 3D rendering: images are added as textures onto a plane, then varying lighting conditions, shadows, and occlusions are added to the scene. The scene is rendered from multiple viewpoints, generating perspective distortions more consistent with real-world scenarios, with homographies closely resembling those of camera projections rather than randomized homographies. For each scene, we provide a sequence of distorted images with corresponding occlusion masks, homographies, and ground-truth labels. The resulting dataset can serve as a training and evaluation set for a multitude of tasks involving image alignment and artifact removal, such as deep homography estimation, dense image matching, 2D bundle adjustment, inpainting, shadow removal, denoising, content retrieval, and background subtraction. Our data generation pipeline is customizable and can be applied to any existing dataset, serving as a data augmentation to further improve the feature learning of any existing method.

## 1 Introduction

Many classical computer vision tasks deal with the problem of image alignment and homography estimation. This usually requires the detection of sparse key points and the computation of correspondences across multiple images. In recent years, increasingly more methods have begun to use neural networks for feature extraction and matching of key points [1, 2]. However, there is a lack of datasets that provide sufficient data and variety to train models. In order to train end-to-end image alignment models, large datasets of high-resolution images are necessary. Existing end-to-end deep learning methods, therefore, utilize datasets containing sparse image patches [3], dense correspondences from structure-from-motion (SfM) datasets [4, 5] or optical flow datasets [6]. Image patches only provide sparse correspondences between images, SfM datasets often lack a variety of scenes, and optical flow data has a high correlation between images with small displacements.

A synthetic data generation for image alignment and restoration (SIDAR) is proposed. A planar object is generated, and a texture is added to it to simulate an artificial painting. Randomized geometric objects, lights, and cameras are added to the scene. By rendering these randomized scenes changing illumination, specular highlights, occlusions, and shadows are added to the original image. The data generation uses images from the WikiArt dataset [7] to create a large variety of different content. However, the data generation can take any image dataset as input. For each image, several distorted images are created with corresponding occlusion masks and pairwise homographies. The datasets can be used for multiple objectives, such as homography estimation, image restoration, dense image matching, and robust feature learning. The rendering pipeline can be configured to generate specific artifacts or any combination of them. It can be used as a data augmentation to any existing datasets to further improve the performance and robustness of existing methods.
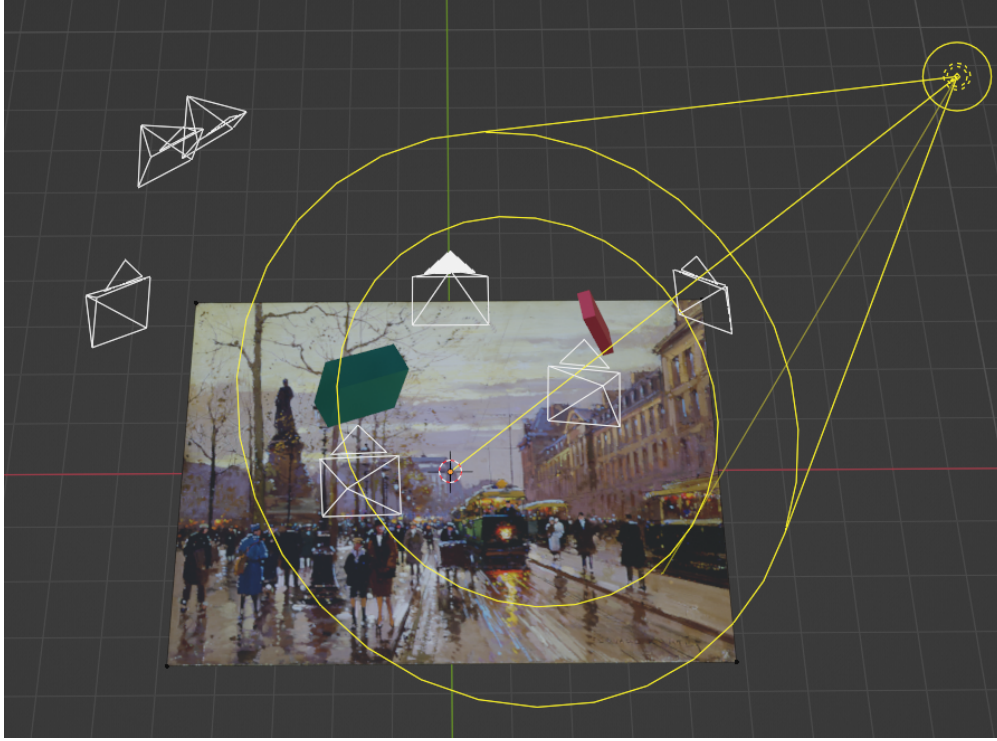


Figure 1: Illustration of a randomly generated scene using Blender. The plane shows a painting. The white pyramids describe randomly generated cameras; the yellow cone describes a spotlight. Geometric objects serve as occlusions and cast shadows onto the plane.

Our dataset and rendering pipeline provide the following contribution:

- Data quantity: An arbitrary number of distortions can be generated for each input image. Arbitrarily many scenes can be generated with an arbitrarily long image sequence each.
- Data variety: The input images and the image distortions can create a large variation in content and artifacts.
- Customizability: The rendering pipeline can be configured to control the amount and type of artifacts. For example, custom datasets can be created that only contain shadows or occlusions.
- Extendibility: The rendering pipeline can be extended to other datasets. It is also possible to change the rendering pipeline in order to generate new artifacts and data modalities, such as videos of moving objects.

One dataset is generated that only consists of aligned image sequences with various distortions, such as changing illumination, shadows, and occlusions. Another dataset is generated that, in addition to the aforementioned distortions, also contains perspective distortions. In order to reconstruct the underlying signal from a sequence of distortions, the information needs to be aligned first. Shadows and illumination distort the original signal, but they do not remove the information completely. Occlusions cover the underlying content in parts of the image. However, for all of these cases, using additional images allows us to identify and remove distortions using pixel-wise or patchwise comparisons along the temporal dimension. Perspective distortions create misalignment that makes reconstruction more complicated. For this reason, we treat misaligned images as a categorically different problem.

In the following chapters, we discuss related datasets and compare their advantages and shortcomings. The rendering pipeline is described in detail. Finally, the generated datasets are presented, and several use cases for end-to-end learning are discussed.

## 2 Related Work

In this section, we review some existing datasets and compare them with our dataset. We discuss the application of our dataset to various tasks.

### 2.1 Image Restoration & Background Subtraction

Our data generation allows the creation of image sequences containing various distortions. In the case of a static camera, this is similar to a task such as background subtraction or change detection. Existing dataset for background subtraction use videos of a static scene [8, 9, 10, 11, 12]. Each individual frame can contain deviations from an ideal background image. The variance can be due to changes in the background, such as illumination changes, weather conditions, or small movement of objects. Additionally, distortions can be caused by foreground objects or camera noise. The goal is to extract a background model of the scene, which can be further used for segmentation into background and foreground.

When training deep learning models, one can differentiate between scene-dependent models and scene-independent models [13]. Scene-dependent models learn a background model on a specific scene and must generalize on new images of the same scene, whereas scene-independent models learn are evaluated on new scenes. Many existing datasets, such as CDNet [9] or SBMNet[10], contain many individual frames but only a few different scenes. Training a model end-to-end on these datasets allows a robust training of scene-dependent models, but the low variation in scenes limits the generalization of models. Our data generation can create arbitrarily many variations within a specific scene and across different scenes. This can be useful to train and evaluate scene-independent models. Compared to video sequences, SIDAR has much more variation between each image, and there is no correlation between each frame. Methods that rely on optical flow cannot be used here.

The proposed SIDAR dataset can be customized to generate specific artifacts, such as changes in illumination, occlusions, and shadows. It can also create any combination of these artifacts. One can apply our dataset to train models for various tasks, such as the detection and removal of shadows, specular highlights, or occlusions. Existing datasets containing shadows [14, 15, 16] or illumination changes [17, 18] often deal with each artifact individually and provide less control over their combinations.

The main shortcoming of our dataset compared to existing datasets are that it is limited to a planar background and relies on synthetic data generation. Certain artifacts caused by the geometry of the scene or due to noise in the image formation from real cameras cannot be synthetically replicated.

### 2.2 Homography Estimation

Homography estimation describes a fundamental task in computer vision. However, there are few datasets providing enough ground truth homographies between image pairs to train a deep-learning model. A common approach to generate labeled data is to use any image dataset, such as MS COCO [19], and apply perspective distortion [20, 21, 22]. These methods do not introduce any additional challenges or artifacts. Other methods apply perspective distortion to video sequences with static cameras [23, 24]. The video sequences introduce dynamic objects that can create outliers when computing correspondences between images. In both cases, the perspective distortion is created artificially and does not follow a realistic planar projection.

Datasets that rely on structure-from-motion or otherwise estimate homographies from real images, such as *HPatches* [3], *Oxford Affine* [25] or *AdelaideRMF*[26] only provide enough image data to train deep learning models on smaller image patches for feature detection. *NYU-VP* and *YUD+* use detected planar objects in 3D scenes [27] to develop a large dataset that is used for self-supervised training. *NYU-VP* and *YUD+* only provide sparse correspondences of line segments.

For all of these datasets, there is little variation in scenes, and each individual image does not contain many distortions. The proposed SIDAR dataset tries to overcome these shortcomings by generating strong distortions within each scene. Homographies between images can be computed regardless of the complexity of the scene and the presence of artifacts. As described in section 3.3, a homography can be computed between any image pair if the relative position of cameras and plane are known. Furthermore, we provide homographies between all images. This provides a unique setting where

the relative orientation of all images can be jointly estimated under various distortions. It is possible to evaluate Bundle Adjustment methods under perturbations and develop end-to-end deep learning methods [28, 29].

## 2.3 Descriptor Learning & Dense Correspondences

Keypoint detection and image descriptors are fundamental methods in many computer vision tasks. Structure-from-motion and other photogrammetric methods rely on point correspondences computed from sparse key points [30]. A requirement for local feature detectors is to identify the location of distinct image points and compute a robust feature representation. The features should be invariant to various distortions, such as noise, changing illumination, scale, and perspective distortions. Recent developments in feature detectors use a data-driven approach to learn robust feature representations with deep learning [31, 2, 32]. The descriptors are either trained on sparse image correspondences from structure-from-motion methods [4] or by using perspective transformations on any image dataset, such as MS COCO [19], as described in section 2.2.

SIDAR provides the ground truth homographies between image pairs with arbitrarily many artifacts. The dataset allows explicitly adding illumination changes, shadows, specular highlights, etc., as data augmentations to train more robust descriptors. The homographies also provide dense correspondences for each pixel. This allows us to train models for dense image matching [33, 34]. Structure-from-motion datasets [4, 5] also provide dense correspondences, but they often contain a limited amount of scenes, distortion, and only a few occlusions. By changing the texture of the image plane, SIDAR has a larger variety of patches and key points.

## 3 Rendering Pipeline

Figure 1 illustrates the data generation process. We use paintings from the Wiki Art dataset [7] as our ground-truth labels. Any other image dataset could also be used, but Wiki Art contains a large variety of artworks from various periods and art styles. We believe that the diversity of paintings makes the reconstruction more challenging and reduces biases towards a specific type of image. We take an image from the dataset and use it as a texture on a plane in 3D.

Furthermore, we generate geometric objects and put them approximately in between the plane and the camera's positions. We utilize Blender's ability to apply different materials to textures. We apply randomized materials to the image texture and occluding objects. The appearance of an occluding object can be diffuse, shiny, reflective, and transparent. The material properties also change the effect lighting has on the plane. It changes the appearance of specularities, shadows, and overall brightness.

Finally, we iterate over the cameras and render the images. Blender's physically-based path tracer, Cycles, is used for rendering the final image. Path tracing enables more realistic effects compared to rasterization. It allows the simulation of effects, such as reflections, retractions, and soft shadows.

### 3.1 Virtual Painting

We first generate a 2D image plane. The plane lies on the $xy$-plane, i.e., the plane is described as:

$$0 \cdot x + 0 \cdot y + z = 0 \tag{1}$$

The center of the plane also lies exactly in the origin $(0, 0, 0)$. The plane is described by its four corner points. Let $w, h$ be the width and height of the plane, then the corners are defined as:

$$\mathbf{X}_1 = \begin{pmatrix} w/2 \\ h/2 \\ 0 \end{pmatrix}, \mathbf{X}_2 = \begin{pmatrix} -w/2 \\ h/2 \\ 0 \end{pmatrix}, \mathbf{X}_3 = \begin{pmatrix} w/2 \\ -h/2 \\ 0 \end{pmatrix}, \mathbf{X}_4 = \begin{pmatrix} -w/2 \\ -h/2 \\ 0 \end{pmatrix} \tag{2}$$

We scale the plane along the $x$ and $y$ direction to fit the image's aspect ratio. We apply the given image as a texture to the plane.

### 3.2 Aligned Camera

To render the scene, we add virtual cameras to the scene. We differentiate between a camera that is aligned with the painting's plane and a setup that adds perspective distortions. To enforce the alignment, we use a single static camera that perfectly fits the image plane. The camera's viewing direction is set to be perpendicular to the image plane and centered on the image plane. Our goal is to adjust the vertical and horizontal field of view such that only the image can be seen. Figure 2 illustrates this problem.
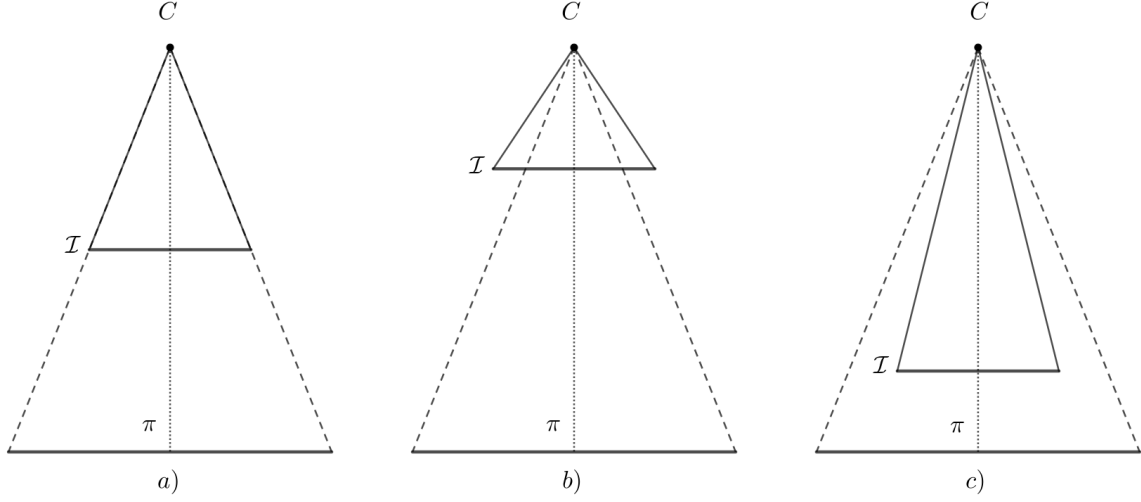
Figure 2: An Illustration of changing the principal distance between the projection center $C$ and the image plane $\mathcal{I}$. The position of the plane $\pi$ and the projection center are fixated. In a) the image plane $\mathcal{I}$ captures all of the content from the painting plane $\pi$.

We fixate on the camera's projection center, and we also fixate on the size of the sensor. We set the resolution and aspect ratio of the sensor equal to the painting. As can be seen in figure 2, the alignment only depends on the principal distance.

Let $f$ be the principal distance, $w$ the image width, $w'$ the sensor width, and $d$ the distance between the projection center and image as illustrated in figure 3. The optimal $f$ can be computed from the intercept theorem:

$$\frac{f}{d} = \frac{w'}{w} \tag{3}$$

$$\Leftrightarrow f = d\frac{w'}{w} \tag{4}$$

Note that the previous formulas and graphics only describe the alignment with regard to a single dimension, i.e., image width. Let $h, h'$ be the image height and sensor height, respectively.

$$\frac{w}{h} = \frac{w'}{h'} \tag{5}$$

$$\Leftrightarrow \frac{h'}{h} = \frac{w'}{w} \tag{6}$$

We can use the intercept theorem to find the optimal focal length $f'$ along the other dimension:

$$\frac{f'}{d} = \frac{h'}{h} \tag{7}$$

$$\Leftrightarrow f' = d\frac{h'}{h} = d\frac{w'}{w} = f \tag{8}$$

$$\Rightarrow f' = f \tag{9}$$

Since we set the sensor's aspect ratio equal to the image's aspect ratio, both dimensions share the same optimal focal length.

### 3.3 Perspective Distortions

To create perspective distortions, we randomize the generation of cameras by sampling from a range of 3D positions. The field of view of the camera is also randomly sampled to create varying zoom effects. This also creates a variety of
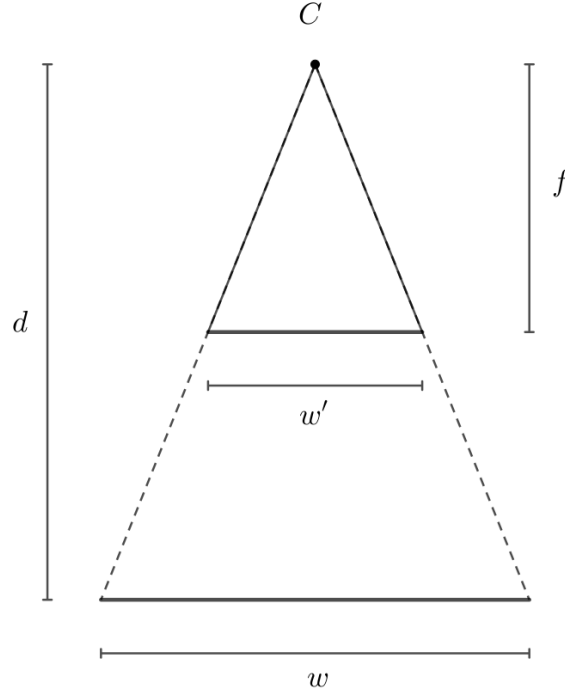
Figure 3: Illustration of a camera sensor with width $w'$ that is aligned with the image with width $w$. The principal distance is given as $f$, and the distance between the camera and image is given as $d$.

intrinsic camera parameters. Additionally, we center the viewing direction of the camera into the image center. This is done to guarantee that the image is seen by the camera. Otherwise, the camera might only see empty space.

Since we project a planar object onto image planes, all images are related by 2D homographies. This relationship is illustrated in figure 4.
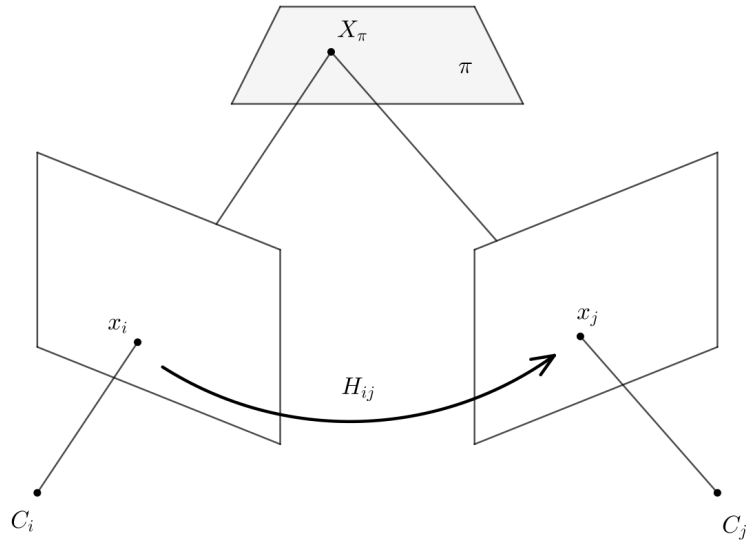


Figure 4: Illustration of a homography induced by a plane.

A point with pixel coordinates $(x, y)$ in image $i$ is projected onto the coordinates $(x', y')$ in image $j$ with the homography $H_{ij}$:

$$\lambda \begin{pmatrix} x' \\ y' \\ 1 \end{pmatrix} = H_{ij} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \tag{10}$$

The following chapters explain two different methods for calculating the homography. Unlike in real datasets, these methods can always compute the corresponding homographies regardless of overlap and distortions. Real datasets rely on more complicated photogrammetric methods that try to find point correspondences between images and optimize the orientation of cameras using bundle adjustment [3, 30, 4]. Regardless of the complexity of the scene, the homographies can be precisely computed.

### 3.3.1 Homography from Projection Matrices and Plane

We can derive the homography by projecting the point $x_i$ onto a point on the plane $X_\pi$ and then project it into the second image to derive $x_j$ [30]. The resulting transformation describes the homography between both cameras. We assume that the origin lies in the projection center of the first camera $C_i$.

**Theorem 1** *Given two projection matrices $P_i = [I|0]$, $P_j = [A|a]$ and a plane defined as $\pi^T X = 0$ with homogeneous coordinates $\pi = (v^T, 1)^T$, the homography induced by the plane is:*

$$H_{ij} = A - av^T \tag{11}$$

We can extend the previous theorem to a calibrated setting.

**Theorem 2** *Given two calibrated cameras described by projection matrices $P_i = K_i[I|0]$, $P_j = K_j[R|t]$ and a plane $\pi_E = (n^T, d)^T$, the homography induced by the plane is:*

$$H_{ij} = K_j \left( R - \frac{1}{d} t n^T \right) K_i^{-1} \tag{12}$$

In order to compute the projective transformation between two arbitrary cameras, the relative orientation $R, t$ has to be extracted first by changing the origin into the projection center of the first camera. Furthermore, the plane has to be computed with regard to the new coordinate system.

### 3.3.2 Direct Linear Transform

The homographies can be calculated using (12) with Blender's scene graph. However, we decided to apply the Direct Linear Transform to point correspondences[30] since it does not require traversing the scene graph and transforming the objects relative to the new coordinate system. We can compute the homographies from the projection matrices and correspondences directly.
Homographies can be estimated from four point-correspondences:

**Theorem 3** *Let $\mathbf{x}_i' = (x_i, y_i, w_i)^T$ and $\mathbf{x}_i' = (x_i, y_i, w_i)^T$ be the homogeneous coordinates of two corresponding points $\mathbf{x}_i \leftrightarrow \mathbf{x}_i'$. The points are related by a homography $\mathbf{x}_i' = H\mathbf{x}_i$. Then the correspondence can be described by an under-determined homogeneous equation:*

$$A_i h = \begin{pmatrix} 0 & -w_i' \mathbf{x}_i^T & y_i' \mathbf{x}_i^T \\ -w_i' \mathbf{x}_i^T & 0 & x_i' \mathbf{x}_i^T \end{pmatrix} \begin{pmatrix} h_1 \\ \vdots \\ h_9 \end{pmatrix} = 0$$

*where $A_i$ is a design matrix containing the measured point correspondences $\mathbf{x}_i \leftrightarrow \mathbf{x}_i'$ and $h = vec(H)$ is vectorized version of $H$. Using at least four correspondences, the homography can be estimated by concatenating into one larger homogeneous equation system:*

$$Ah = \begin{pmatrix} A_1 \\ \vdots \\ A_4 \end{pmatrix} \begin{pmatrix} h_1 \\ \vdots \\ h_9 \end{pmatrix} = 0 \tag{13}$$

*$h$ can be exactly determined by finding the null space of $A$.*

Since Blender provides the camera's positions in space and their intrinsic parameters, we can calculate the projection points explicitly. We project each of the painting's corners $X_k, \quad k = 1, .., 4$ into the $i$-th image plane $\mathbf{x}_i^{(k)}$ using their projection matrices:

$$\mathbf{x}_i^{(k)} = P_i X_k \tag{14}$$

$$= K_i [I|0] \begin{bmatrix} R_i & t_i \\ 0 & 1 \end{bmatrix} X_k \quad k = 1, .., 4 \tag{15}$$

$(R_i, t_i)$ describes the global rotation and translation of the camera, $K_i$ describes the intrinsic camera parameters. Using the four-point pairs $\mathbf{x}_k^{(i)} \leftrightarrow \mathbf{x}_k^{(j)} \quad k = 1, .., 4$, we can compute the homography by solving the homogeneous equation (13). We compute the Singular Value Decomposition of A:

$$A = U \Sigma V^T = U [\text{diag}(\sigma_1, \dots, \sigma_8)|0] V^T \tag{16}$$

$$= U \begin{pmatrix} \sigma_1 & & \cdots & & 0 \\ & \sigma_2 & \cdots & & 0 \\ & & \ddots & & 0 \\ & & & \sigma_8 & 0 \end{pmatrix} \begin{pmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_8^T \\ v_9^T \end{pmatrix} \tag{17}$$

Since $A \in \mathbb{R}^{8 \times 9}$ it has $\text{rank}(A) = 8$. We assume sorted singular values: $\sigma_1 > \sigma_2 > \cdots > \sigma_8$. The solution to (13) is then found in the corresponding row $v_9$ of $V^T$.

$$A v_9 = U \begin{pmatrix} \sigma_1 & & \cdots & & 0 \\ & \sigma_2 & \cdots & & 0 \\ & & \ddots & & 0 \\ & & & \sigma_8 & 0 \end{pmatrix} V^T v_9 \tag{18}$$

$$= U \begin{pmatrix} \sigma_1 & & \cdots & & 0 \\ & \sigma_2 & \cdots & & 0 \\ & & \ddots & & 0 \\ & & & \sigma_8 & 0 \end{pmatrix} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \tag{19}$$

$$= U \cdot 0 = 0 \tag{20}$$

In summary, the homography estimation is described by the pseudocode 1.

---

**Algorithm 1** Direct Linear Transform

---

**Input:** $P_i, P_j, X_1, \cdots, X_4$
**Output:** $H$
   **for** $k \in \{1, 2, 3, 4\}$ **do**
      $\mathbf{x}_k \leftarrow P_i X_k = (x_k, y_k, w_k)$
      $\mathbf{x}_k' \leftarrow P_j X_k = (x_k', y_k', w_k')$
      $A_k \leftarrow \begin{pmatrix} 0 & -w_k' \mathbf{x}_k^T & y_k' \mathbf{x}_k^T \\ -w_k' \mathbf{x}_k T & 0 & x_k' \mathbf{x}_k^T \end{pmatrix}$
   **end for**
$A \leftarrow \begin{pmatrix} A_1 \\ \vdots \\ A_4 \end{pmatrix}$
$U, \Sigma, V^T \leftarrow \mathbf{SVD}(A)$
$v_1, \cdots, v_9 \leftarrow V$
$H_{ij} \leftarrow \mathbf{reshape}(v_9)$
**return** $H_{ij}$

---

## 3.4 Illumination

We add varying illumination by randomizing the light sources in the scene. The randomization can be described by the pseudocode 2. The light is randomly sampled from Blender's predefined light sources: spotlight, point light, and

area light. We do not use a light source that simulates the sun because we want to simulate indoor environments. We randomize the intensity and color of the light. This randomization can create a large variety in the appearance of shadows, specularities, and the overall color scheme of the image.

The color is sampled using the HSV color space. Let $c = (H, S, V)$ be the color of the light. We sample $c$ using $H \sim \mathcal{U}[0, 1], S \sim \mathcal{U}[0, \epsilon], V = 1$. $V$ is fixed because the intensity of the color is already affected by the light source itself. The hue is completely randomized to create random colors. However, the range of the saturation is limited to a small $\epsilon$ to generate lights that are closer to white. Having too high saturation creates unrealistic and extreme lighting conditions that make the reconstruction of the painting's color very ambiguous, while low saturation still introduces enough variance in the appearance of the painting.

Furthermore, we set the orientation of the light such that its direction is centered on the images. By default, Blender sets the orientation of the lights along the $z$-axis. As can be seen in figure 5, the correct orientation $\theta$ can be computed from the angle between the location vector $v$ and the z-axis $(0, 0, 1)$.
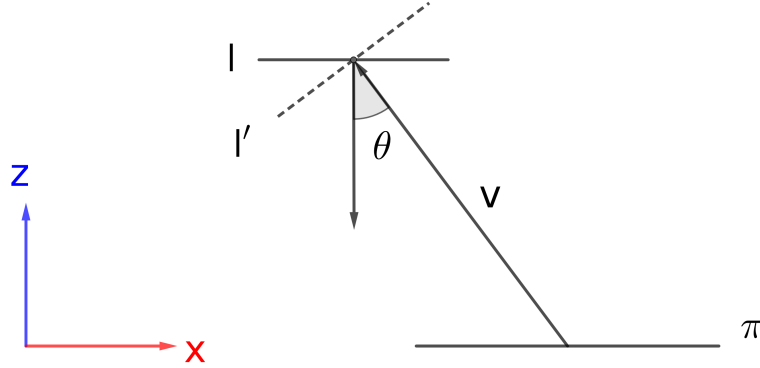


Figure 5: An illustration for the centering of an area light $l$. By default, Blender aligns new light sources along the $z$-axis. The correct orientation is described by the dotted line $l'$. $\theta$ describes the rotation offset between both configurations. $\theta$ is the angle between the (inverse) location vector $-v$ of the light and the vector $(0, 0, -1)$

### 3.5 Add Geometric Objects

Randomized geometric objects are added in the space between light sources and the image plane. The purpose of these objects is to create occlusions and shadows. The objects obfuscate image content from the cameras and block light from reaching the image plane. The material of the object is also chosen from Blender's shaders. Depending on the material, light is either completely blocked, refracted, or color filtered. As can be seen in figure 6, the material affects the appearance of the image content and casts a shadow. Materials in figure 6a and 6d create hard shadows, while figures 6b and 6c create softer shadows and even change the shadow's color. Transparent materials also do not fully obfuscate the image content.

---

**Algorithm 2** Add Illumination

---

**Input:** $n$
**Output:** $\{l_1, \cdots, l_n\}$
  lights $\leftarrow \emptyset$
  **for** $k \in \{1, \cdots, n\}$ **do**
      type $\sim$ **sample**{"Spot","Point","Area"}
      $l_k \leftarrow$ new Light(type)                                      $\triangleright$ create light object

      $l_k$.intensity $\sim \mathcal{U}(I_{min}, I_{max})$                              $\triangleright$ random intensity

      $x, y, z \sim \mathcal{U}(x_{min}, x_{max}) \times \mathcal{U}(y_{min}, y_{max}) \times \mathcal{U}(z_{min}, z_{max})$       $\triangleright$ random position
      $l_k$.position $\leftarrow x, y, z$

      h,s,v $\leftarrow \mathcal{U}(0,1) \times \mathcal{U}(0, \epsilon) \times 1$                      $\triangleright$ sample color in HSV space
      $l_k$.color $\leftarrow HSV(h, s, v)$

      $v \leftarrow l_k$.position                                    $\triangleright$ get location vector of light
      $l_k$.orientation $\leftarrow \measuredangle(v, (0,0,1)^T)$     $\triangleright$ center the direction by computing the angle between $v$ and z-axis.

      lights.append($l_k$)
  **end for**
  **return** lights

---

---

**Algorithm 3** Add Geometric Objects

---

**Input:** $n$
**Output:** $\{o_1, \cdots, o_n\}$
  objects $\leftarrow \emptyset$
  **for** $k \in \{1, \cdots, n\}$ **do**
      geometry $\sim$ **sample**{"Sphere","Cube","Cylinder", ... }
      $o_k \leftarrow$ new Object(type)

      $x, y, z \sim \mathcal{U}(x_{min}, x_{max}) \times \mathcal{U}(y_{min}, y_{max}) \times \mathcal{U}(z_{min}, z_{max})$       $\triangleright$ random position
      $o_k$.positions $\leftarrow (x, y, z)$

      $\alpha_x, \alpha_y, \alpha_z \sim \mathcal{U}(0, 2\pi) \times \mathcal{U}(0, 2\pi) \times \mathcal{U}(0, 2\pi)$           $\triangleright$ random orientation
      $o_k$.rotation $\leftarrow (\alpha_x, \alpha_y, \alpha_z)$

      $s_x, s_y, s_z \sim \mathcal{U}(s_{x_{min}}, s_{xmax}) \times \mathcal{U}(s_{ymin}, s_{ymax}) \times \mathcal{U}(s_{zmin}, s_{zmax})$    $\triangleright$ random scaling
      $o_k$.scale $\leftarrow (s_x, s_y, s_z)$

      material $\sim$ **sample**{"Refraction","Transparent","Metallic", ... }     $\triangleright$ sample random material type
      $o_k$.material $\leftarrow$ new Material(material)

      objects.append($o_k$)
  **end for**
  **return** objects

---

### 3.6 Rendering

After a scene is configured with randomized lighting, occlusions, and cameras, an image is rendered using path tracing. Blender's physics-based path tracer is used to render the image from the perspective of a specific camera. Path tracing allows the creation of more realistic shadows and lighting effects compared to rasterization. Especially effects such as transparency, reflection, and refraction can be realistically modeled using path tracing by sampling light rays. Path tracing simulates the physical image formation process much closer than rasterization techniques. Rasterization requires the use of various techniques, such as texture maps and shadow maps, to approximate the same effects.The number of rays cast to estimate the light distribution of the scene and the resulting image can be limited using a time constraint. A

(a) Gloss

(b) Refraction

(c) Transparency

(d) Solid Material

Figure 6: A randomly generated torus under the same lighting conditions but with varying materials.

time limit is set to balance the image quality and the amount of data that can be efficiently generated. Figure 7 illustrates the effect of rendering time on the resulting image. The results do not differ too much. Shadows, illumination, and occlusions are visualized correctly, even when using fewer rays. A low amount of sampling can create aliasing effects, as can be seen in the glassy material. In comparison, more sampling creates more realistic effects.

Each rendered image describes a distorted data point of the original image. Path tracing is also used to generate the ground truth label. It is possible to use the original image as a ground truth label. However, the rendering pipeline can create bias and small misalignment. For this reason, the label is created under similar conditions as the distorted images using ambient illumination without occlusion. A virtual camera is used that is aligned with the image as described in chapter 3.2. The resulting image is free of artifacts. Figure 8 compares the original images and the rendered image. Although both images look very similar, there is a small difference between both images around edges, as can be seen in figure 8c. The difference might be caused by small misalignments or texturing artifacts.

In order to differentiate between distorted parts of the image (caused by lighting, shadows, and specular highlights) and occlusions, a segmentation mask is computed that separates foreground and background pixels. The mask should be aligned with the corresponding camera. After rendering any scene, a corresponding occlusion mask is rendered using the same camera and geometric objects. For any given scene, the material of the plane is changed to a diffuse black. All objects are changed to a diffuse white material. We also use ambient illumination. The segmentation mask is rendered using regular rasterization. Figure 9 shows a mask generated from a given scene. The resulting image is binary, with black pixels describing parts of the painting and white pixels describing occlusions or background.

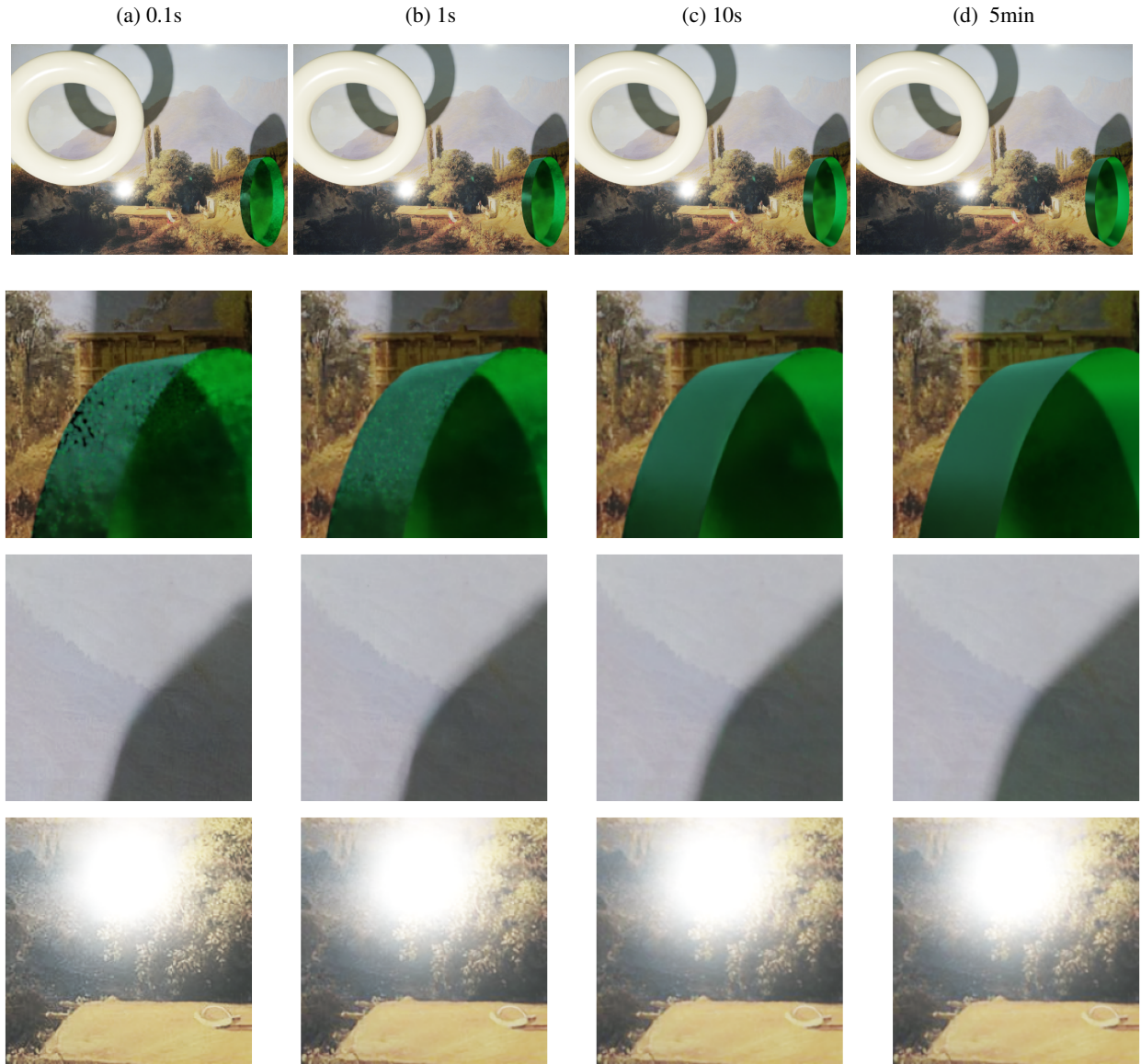(a) 0.1s  (b) 1s  (c) 10s  (d) 5min

Figure 7: The images show the same scene rendered with different time limits. The first row shows the whole images, while the second to fourth rows show specific parts of the image.



(a) Original image: $I_O$  (b) Ambient lighting: $I_A$  (c) Difference: $|I_O - I_A|$

Figure 8: The images show the original image and the output image generated by the rendering pipeline using an aligned camera and ambient lighting. c) shows the absolute difference of each pixel.
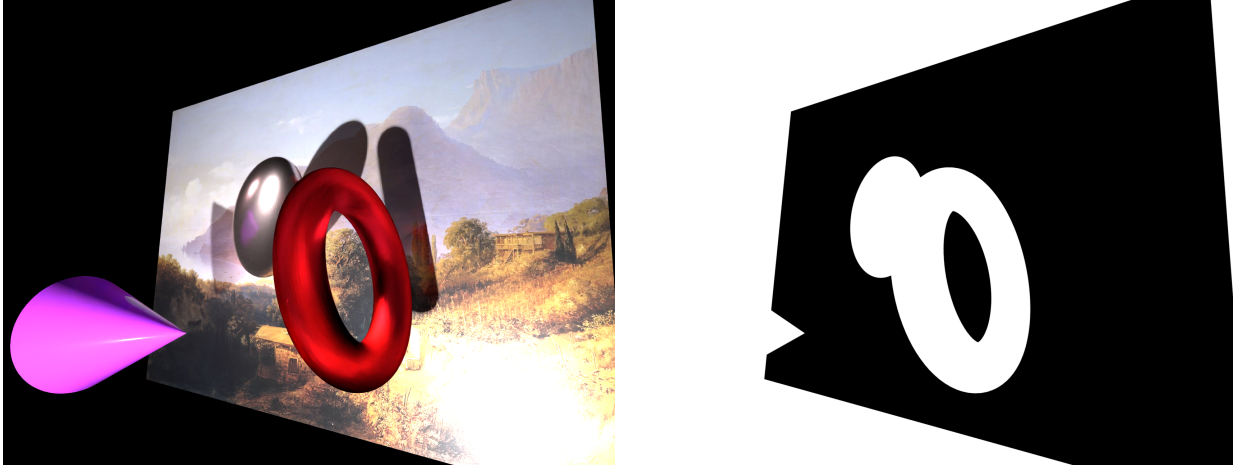
Figure 9: The left image shows a randomized scene. The right image shows the same scene rendered after changing the plane to a diffuse black and all geometric objects to a diffuse white color.

## 4 Datasets

Using the rendering pipeline from chapter 3, we create two datasets. By aligning the camera with the image, we create a dataset of aligned sequences without perspective distortion. Another dataset is created by adding perspective distortions.

### 4.1 Aligned Dataset

Figure 10 shows an example sequence of distorted images with the corresponding ground truth label. We generate image sequences, each containing multiple distorted images as our dataset. For each distorted image, an additional occlusion mask is created. The dataset can be used for multiple tasks. Since the amount of data can be increased indefinitely, it is especially useful for deep learning models.
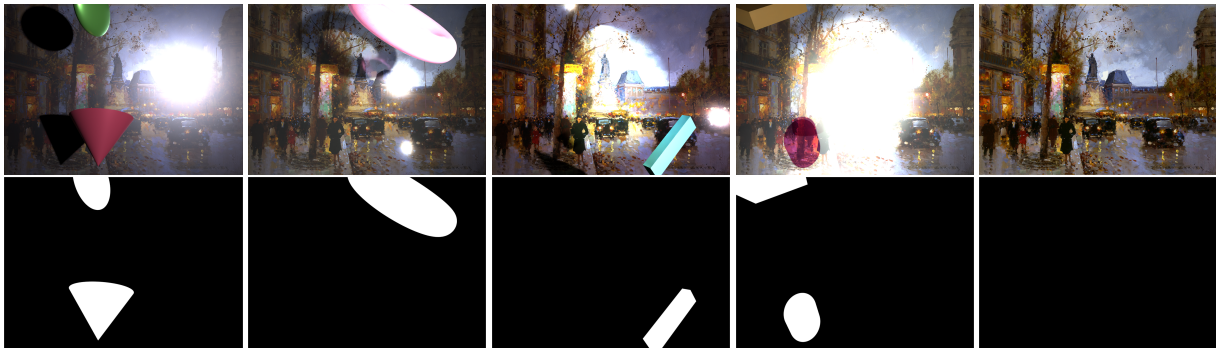


Figure 10: The top row shows four generated images without perspective distortions and with the corresponding label. The bottom row shows the occlusion mask for each image.

**Image Restoration**

Using the aligned dataset, a model can be trained to remove artifacts from images. This can be done by a sequential model that learns to aggregate information from multiple images, such as Deep Sets [35, 36] Alternatively, a model can be trained to remove artifacts from single images, e.g., single-image shadow removal [16]. Existing methods often deal with each type of artifact individually, whereas this dataset allows combining multiple artifacts simultaneously. The data generation can be configured to generate specific artifacts or combinations of them. Using the occlusion mask, inpainting methods can be trained to reconstruct the content of occluded regions.

**Background Subtraction**

The dataset can be used for background subtraction tasks. The artifact-free image describes the underlying background, while lighting and occlusion create distortions. A model can be trained to detect occlusions and variations from the underlying background model. The generated occlusion masks can be used for learning foreground-background segmentation. Existing datasets consist of videos with a few scenes, such as CDNet [9]. Although the videos provide a lot of data, there is little variation within each scene. Our dataset generation allows the creation of a large variety of different scenes and also increases the variation within each scene. This is useful to enable models to generalize over different backgrounds and make the background modeling more robust to changing lighting conditions.

**Representation Learning**

Data augmentations, such as noise, blurring, random cropping, and geometric transformations, are used to create more robust representation learning or self-supervised training [37, 38]. Using our dataset, a representation can be learned that is invariant to illumination and occlusions. Alternatively, a representation can be learned that disentangles the image content from lighting, shadows, occlusions, etc.

### 4.2 Misaligned Dataset

In addition to the aligned dataset, we generate images with perspective distortions. For each randomly generated camera, we render an image. In order to evaluate the image alignment with reconstruction, we also generate a single ground truth image as described in section 3.2 under ambient lighting conditions. Figure 11 shows a sequence of distorted images. The last image contains no distortions and is aligned with the camera's field of view.

For each image pair $(I_i, I_j)$, the corresponding homography $H_{ij}$ is computed using DLT as described in section 3.3.2. It is possible to warp any image $I_j : \mathbb{R}^2 \mapsto \mathbb{R}^3$ into the reference frame of any other image $I_j : \mathbb{R}^2 \mapsto \mathbb{R}^3$ using a warp function $\mathcal{W}_{ij} : \mathbb{R}^2 \mapsto \mathbb{R}^2$. Figure 12 shows a sequence of images under perspective distortion. Using the estimated homographies, all images can be aligned with the first image.



Figure 11: The top row shows four generated images with perspective distortions and their corresponding label. The bottom row shows the occlusion mask for each image.

Homography Adaption is a technique that applies random homographies to images from MS COCO [19] for interest point detection and descriptor learning [2]. Although this technique is successfully applied in various image-matching tasks, it does not explicitly enforce invariance to illumination changes or occlusions. Our dataset allows to overcome these shortcomings. Although our misaligned dataset also creates randomized homographies, it generates homographies from camera projections, which better represent realistic 2D projections.

The dataset containing perspective distortions extends all tasks mentioned in section 4.1, but it also creates new challenges and applications. The following chapters discuss some potential applications of our dataset.

**Homography Estimation**

Our dataset provides the ground-truth homography for any image pair within a sequence, making it possible to train and evaluate deep homography estimation methods. Let $I_i$ and $I_j$ be two images, let $H_{ij}$ be their corresponding perspective transformation, and let $f_\theta(I_i, I_j) \in \mathbb{R}^9$ be a deep learning model that estimates a homography from two images. The
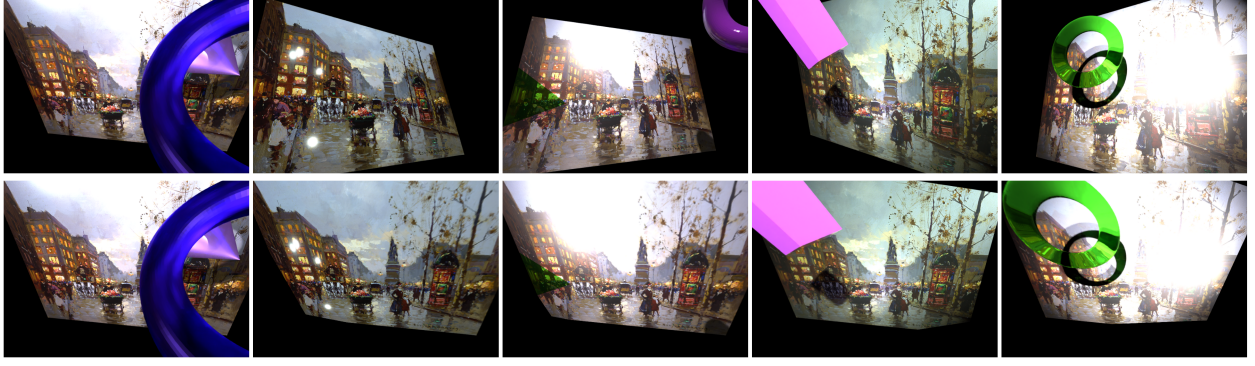
Figure 12: The first row shows a sequence of images with perspective distortions. The second row shows the images warped into the reference frame of the first image.

learning objective can be described by a regression loss $L(\theta)$:

$$\hat{H}_{ij} := f_\theta(I_i, I_j)$$

$$L(\theta) = \left\| \frac{1}{\sqrt[3]{|\hat{H}_{ij}|}} \hat{H}_{ij} - \frac{1}{\sqrt[3]{|H_{ij}|}} H_{ij} \right\|_F \tag{21}$$

Here $|\cdot|$ describes the determinant, and $\|\cdot\|_F$ describes the Frobenius norm. Both matrices have to be normalized since homographies are equivalent under a scaling factor. The given normalization fixes the scale of both matrices such that their determinants are equal to one. The objective functions allow the training of an end-to-end deep learning model. The data generation provides an indefinite amount of data. The changes in illumination and occlusion provide additional challenges and require a more robust estimation of features.

### Bundle Adjustment

Many existing homography estimation methods compute the alignment from image pairs only. This can be extended to sets of images. The problem can be described as a bundle adjustment problem. The SIDAR dataset can be used for neural bundle adjustment methods, such as BARF[28]. It could be possible to learn neural priors for bundle adjustment. Larger image sets also enforce more consistency across images compared to image pairs.

### Descriptor Learning

Given the correspondences between images, local descriptors can be learned. Correspondences exist even under very strong distortions, which allows the development of descriptors that are invariant or equivariant to the given perturbations. The methodology of HPatches [3] could also be extended to the SIDAR dataset to add more variety in image patches.

### Dense Correspondences

The ground-truth homographies also provide dense correspondences between each image point with high accuracy. Correspondences can be estimated not only for sparse key points but for every pixel. The neighborhood of points remains mostly unchanged under perspective distortions. This puts additional constraints on image-matching tasks. The occlusion masks also provide regions of outliers, while the other distortions add robustness. Image matching models can be trained to densely detect image regions under various perturbations and also detect outliers as points with no matches.

## 5    Conclusion

In this work, we propose a data generation with a corresponding dataset based on 3D rendering that introduces various disturbances, such as shadows, illumination changes, specular highlights, occlusions, and perspective distortions, to any given input image. Although it is a synthetic dataset, the data augmentations are not trivial, and they are customizable. To the best of our knowledge, we provide the first large-scale dataset containing ground-truth homographies of image sequences for deep homography estimation, which does not consist of trivial perspective distortions. Our rendering pipeline allows us to both generate new datasets and augment existing data. Models trained on these augmentations can learn a representation that is robust to these perturbations. It can be used as an extension of Homographic Adaption [2] or contrastive learning [37]. We discuss a range of computer vision applications for which this dataset can be used. It can contribute to the training of end-to-end deep learning models that solve image alignment and restoration tasks such as deep homography estimation, dense image matching, descriptor learning, 2D bundle adjustment, inpainting, shadow removal, denoising, content retrieval, and background subtraction.

The limitation of most synthetic datasets lies in their deviation from real data. This can result in biased models and limit generalization and, thus, model performance. Compared to existing augmentation methods that apply randomized homographies [2] to images, SIDAR adds additional complexity. Especially adding illumination changes, shadows, and occlusions can improve the robustness of learned descriptors and feature matching.

Future work could focus on developing benchmarks to provide specific evaluation metrics for the discussed tasks and compare the generalization across different datasets. Additionally, future work could further improve the realism of the data generation and add new data modalities, such as videos. The data generation can be adapted to include other distortions, such as reflective surfaces, translucent occlusions, or camera lens distortions. Many of these effects are common artifacts in real imaging systems, but it is difficult to create large-scale datasets for these cases. Our data generation can provide an effective way to approximate these artifacts and provide a large enough dataset for training and evaluation. It can serve as a baseline to study these effects in a more controlled environment.

## References

[1]  Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4938–4947, 2020.

[2]  Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 224–236, 2018.

[3]  Vassileios Balntas, Karel Lenc, Andrea Vedaldi, and Krystian Mikolajczyk. Hpatches: A benchmark and evaluation of handcrafted and learned local descriptors. In *CVPR*, 2017.

[4]  Zhengqi Li and Noah Snavely. Megadepth: Learning single-view depth prediction from internet photos. In *Computer Vision and Pattern Recognition (CVPR)*, 2018.

[5]  Thomas Schops, Johannes L Schonberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3260–3269, 2017.

[6]  D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black. A naturalistic open source movie for optical flow evaluation. In A. Fitzgibbon et al. (Eds.), editor, *European Conf. on Computer Vision (ECCV)*, Part IV, LNCS 7577, pages 611–625. Springer-Verlag, October 2012.

[7]  Michael Danielczuk, Matthew Matl, Saurabh Gupta, Andrew Li, Andrew Lee, Jeffrey Mahler, and Ken Goldberg. Segmenting unknown 3d objects from real depth images using mask r-cnn trained on synthetic data. In *Proc. IEEE Int. Conf. Robotics and Automation (ICRA)*, 2019.

[8]  Antoine Vacavant, Thierry Chateau, Alexis Wilhelm, and Laurent Lequievre. A benchmark dataset for outdoor foreground/background extraction. In *Computer Vision-ACCV 2012 Workshops: ACCV 2012 International Workshops, Daejeon, Korea, November 5-6, 2012, Revised Selected Papers, Part I 11*, pages 291–300. Springer, 2013.

[9]  Nil Goyette, Pierre-Marc Jodoin, Fatih Porikli, Janusz Konrad, and Prakash Ishwar. Changedetection. net: A new change detection benchmark dataset. In *2012 IEEE computer society conference on computer vision and pattern recognition workshops*, pages 1–8. IEEE, 2012.

[10] Pierre-Marc Jodoin, Lucia Maddalena, Alfredo Petrosino, and Yi Wang. Extensive benchmark and survey of modeling methods for scene background initialization. *IEEE Transactions on Image Processing*, 26(11):5244–5256, 2017.

[11] Kentaro Toyama, John Krumm, Barry Brumitt, and Brian Meyers. Wallflower: Principles and practice of background maintenance. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 1, pages 255–261. IEEE, 1999.

[12] Rudrika Kalsotra and Sakshi Arora. A comprehensive survey of video datasets for background subtraction. *IEEE Access*, 7:59143–59171, 2019.

[13] Murari Mandal and Santosh Kumar Vipparthi. An empirical review of deep learning frameworks for change detection: Model design, experimental frameworks, challenges and research needs. *IEEE Transactions on Intelligent Transportation Systems*, 2021.

[14] Jifeng Wang, Xiang Li, Le Hui, and Jian Yang. Stacked conditional generative adversarial networks for jointly learning shadow detection and shadow removal, 2017.

[15] Netanel Kligler, Sagi Katz, and Ayellet Tal. Document enhancement using visibility detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2374–2382, 2018.

[16] Liangqiong Qu, Jiandong Tian, Shengfeng He, Yandong Tang, and Rynson WH Lau. Deshadownet: A multi-context embedding deep network for shadow removal. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4067–4075, 2017.

[17] Daniel J Butler, Jonas Wulff, Garrett B Stanley, and Michael J Black. A naturalistic open source movie for optical flow evaluation. In *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part VI 12*, pages 611–625. Springer, 2012.

[18] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M. Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *International Conference on Computer Vision (ICCV) 2021*, 2021.

[19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.

[20] Che-Han Chang, Chun-Nan Chou, and Edward Y Chang. Clkn: Cascaded lucas-kanade networks for image alignment. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2213–2221, 2017.

[21] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Deep image homography estimation. *arXiv preprint arXiv:1606.03798*, 2016.

[22] Farzan Erlik Nowruzi, Robert Laganiere, and Nathalie Japkowicz. Homography estimation from image pairs with hierarchical convolutional networks. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 913–920, 2017.

[23] Hoang Le, Feng Liu, Shu Zhang, and Aseem Agarwala. Deep homography estimation for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[24] Si-Yuan Cao, Jianxin Hu, Zehua Sheng, and Hui-Liang Shen. Iterative deep homography estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1879–1888, June 2022.

[25] Krystian Mikolajczyk and Cordelia Schmid. A performance evaluation of local descriptors. *IEEE transactions on pattern analysis and machine intelligence*, 27(10):1615–1630, 2005.

[26] Hoi Sim Wong, Tat-Jun Chin, Jin Yu, and David Suter. Dynamic and hierarchical multi-structure geometric model fitting. In *2011 International Conference on Computer Vision*, pages 1044–1051. IEEE, 2011.

[27] Florian Kluger, Eric Brachmann, Hanno Ackermann, Carsten Rother, Michael Ying Yang, and Bodo Rosenhahn. Consac: Robust multi-model fitting by conditional sample consensus. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[28] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021.

[29] Philipp Lindenberger, Paul-Edouard Sarlin, Viktor Larsson, and Marc Pollefeys. Pixel-Perfect Structure-from-Motion with Featuremetric Refinement. In *ICCV*, 2021.

[30] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003.

[31] Anastasiia Mishchuk, Dmytro Mishkin, Filip Radenovic, and Jiri Matas. Working hard to know your neighbor's margins: Local descriptor learning loss. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[32] Xi Shen, François Darmon, Alexei A Efros, and Mathieu Aubry. Ransac-flow: generic two-stage image alignment. In *16th European Conference on Computer Vision*, 2020.

[33] Prune Truong, Martin Danelljan, and Radu Timofte. GLU-Net: Global-local universal network for dense flow and correspondences. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2020*, 2020.

[34] Prune Truong, Martin Danelljan, Radu Timofte, and Luc Van Gool. Pdc-net+: Enhanced probabilistic dense correspondence network. In *Preprint*, 2021.

[35] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30, 2017.

[36] Monika Kwiatkowski and Olaf Hellwich. Specularity, shadow, and occlusion removal from image sequences using deep residual sets. In *VISIGRAPP (4: VISAPP)*, pages 118–125, 2022.

[37] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations, 2020.

[38] Arpit Bansal, Eitan Borgnia, Hong-Min Chu, Jie S. Li, Hamid Kazemi, Furong Huang, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Cold diffusion: Inverting arbitrary image transforms without noise, 2022.