# Highlights

**InfraDiffusion: zero-shot depth map restoration with diffusion models and prompted segmentation from sparse infrastructure point clouds**

Yixiong Jing, Cheng Zhang, Haibing Wu *, Guangming Wang, Olaf Wysocki, Brian Sheil

- Introduce a zero-shot framework, InfraDiffusion, to restore depth maps from masonry point clouds.

- Propose a virtual camera projection for depth map generation from point clouds.

- Adapt DDNM with boundary masking for depth image restorations using pre-trained diffusion models.

- Improve segmentation metrics across five datasets using SAM segmentation.

# InfraDiffusion: zero-shot depth map restoration with diffusion models and prompted segmentation from sparse infrastructure point clouds

Yixiong Jing[a], Cheng Zhang[b], Haibing Wu *[a], Guangming Wang[a], Olaf Wysocki[a], Brian Sheil[a]

[a]*Construction Engineering, University of Cambridge, Trumpington Street, Cambridge, CB2 1PZ, Cambridge, UK*
[b]*College of Civil Engineering, Hunan University, Yuelu South Road, Changsha, 410082, Hunan Province, China*

## Abstract

Point clouds are widely used for infrastructure monitoring by providing geometric information, where segmentation is required for downstream tasks such as defect detection. Existing research has automated semantic segmentation of structural components, while brick-level segmentation (identifying defects such as spalling and mortar loss) has been primarily conducted from RGB images. However, acquiring high-resolution images is impractical in low-light environments like masonry tunnels. Point clouds, though robust to dim lighting, are typically unstructured, sparse, and noisy, limiting fine-grained segmentation. We present InfraDiffusion, a zero-shot framework that projects masonry point clouds into depth maps using virtual cameras and restores them by adapting the Denoising Diffusion Null-space Model (DDNM). Without task-specific training, InfraDiffusion enhances visual clarity and geometric consistency of depth maps. Experiments on masonry bridge and tunnel point cloud datasets show significant improvements in brick-level segmentation using the Segment Anything Model (SAM), underscoring its potential for automated inspection of masonry assets. Our code and data is available at `https://github.com/Jingyixiong/InfraDiffusion-official-implement`.

*Keywords:* InfraDiffusion, Diffusion models, Image restoration, Point clouds, Masonry structures, Depth maps, Semantic segmentation, Structural health monitoring

## 1. Introduction

Masonry infrastructure forms a significant part of the transport network and civil engineering assets in the UK (Orbán, 2004). These structures require regular inspection to ensure long-term safety and functionality (Acikgoz et al., 2018). However, these inspections are often time-consuming and reliant on visual judgment of engineers, which introduces both subjectivity and variability into the assessment process (Brackenbury, 2022).

LiDAR scan point clouds are increasingly adopted to support infrastructure assessment, enabling efficient acquisition of the as-is geometry of built environments (Wang and Cho, 2015; Dai and Lu, 2013; Lubowiecka et al., 2009; Shanoer and Abed, 2018). To perform structural analysis from the point cloud, segmentation is required as a pre-processing step to identify structural elements. Recent advances in deep learning (DL) have enabled the automated segmentation of structural 'components' (such as arches, spandrels, piers, or walls) in masonry bridge point clouds (Jing et al., 2022, 2024a). This component-level segmentation enables a range of downstream tasks, including geometric modelling (Jing et al., 2023; Han et al., 2025), deformation tracking (Ye et al., 2018; Acikgoz et al., 2017), and defect detection (Jing et al., 2024b; Han et al., 2025).

Whilst previous work has predominantly focused on component-level segmentation, finer segmentation of individual bricks (i.e., "brick-level segmentation") from point clouds remains unexplored. Brick-level segmentation is important for automating early-stage defect detection Dais et al. (2021); Loverdos and Sarhosis (2022); Hallee et al. (2021); Ye et al. (2024), particularly in large-scale masonry structures. Different from steel and concrete structures, masonry structures do not behave as continuous materials but as assemblages of discrete units where the brick–mortar interaction governs both stiffness and failure mechanisms (Lourenço, 2013; Lemos, 2007). Different deterioration mechanisms can be introduced at this scale: mortar joints typically experience erosion and volume loss, while bricks may undergo spalling, cracking, and surface detachment (Giordano et al., 2002; Milani et al., 2006). Identifying individual bricks within point clouds is therefore crucial for distinguishing between brick- and mortar-related damage, enabling more accurate structural diagnosis (Wu et al., 2019b; Katsigiannis et al., 2023) and supporting advanced analysis approaches such as discrete element modelling (Lemos, 2007).

Previous studies (Dais et al., 2021; Ye et al., 2024) have typically per-

formed brick-level segmentation on images, which provide dense visual information for annotation. However, acquiring high-quality images is often impractical in low-light conditions common in tunnels (Cheng et al., 2019), indoor facilities (Luo et al., 2023), and underground pipes (You et al., 2025). In contrast, point clouds collected by active sensors such as LiDAR are robust to poor lighting (Qu et al., 2014), but their sparsity and noise significantly limit fine-grained segmentation. These limitations highlight the need for methods that enable brick-level segmentation directly from point clouds.

However, it is nontrivial to automate segmentation directly on point clouds using DL models due to two significant drawbacks:

(i) **Lack of large annotated datasets:** 3D DL models designed for point cloud segmentation typically require large and well-annotated datasets for extensive training, which are time-consuming and costly to produce.

(ii) **Lack of generalisation:** Unlike 2D vision (Kirillov et al., 2023) and NLP (Guo et al., 2025), the point cloud domain lacks robust foundation models. As a result, existing 3D DL models trained on narrow datasets often fail to generalise across different types of infrastructure.

Unlike 2D vision (Kirillov et al., 2023) and NLP (Guo et al., 2025), the point cloud domain lacks robust foundation models. As a result, existing 3D DL models trained on narrow datasets often fail to generalise across different types of infrastructure.

To overcome these limitations, recent studies have proposed project-based segmentation methods, which project point clouds into structured 2D depth maps (Zhang et al., 2022; Chen et al., 2024). The 2D depth maps can leverage powerful image-based foundation models that provide greater robustness and generalisation across diverse scenes and segmentation tasks. For example, Ye et al. (2025) recently demonstrated the use of the Segment Anything Model (SAM) (Kirillov et al., 2023) to perform zero-shot instance segmentation of tunnel lining segments, which achieves accurate and robust results without training.

However, the generalisation of current projection-based methods remains limited due to two challenges:

(i) These methods have been predominantly applied to tunnel environments, where cylindrical projection is used for unwrapping the geometry into depth images. Cylindrical projection does not generalise well to other types of infrastructure with more complex or irregular topologies, such as masonry bridges.

3

(ii) The quality of the resulting depth maps heavily depends on the equipment noise, registration errors, and density of the original point cloud. Therefore, the projected depth maps are always incomplete and noisy in real-world scenarios, which makes them unsuitable for brick-level segmentation tasks.

To address challenge (i), the autonomous driving domain has introduced virtual camera projection to allow flexible viewpoint selection by simulating the optics of real cameras and adapting to any complex 3D scene. This approach has been applied across a range of tasks, including segmentation (Wu et al., 2019a; Lyu et al., 2020; Krawciw et al., 2024), 3D object detection (Chen et al., 2017; Meyer et al., 2019), and multi-view data fusion (Massa and Grobler, 2024). Unlike cylindrical or planar projections (Chen et al., 2024; Ye et al., 2025), virtual camera projection enables depth maps generation for infrastructure with irregular or complex forms (such as masonry bridges), where global unfolding techniques often fail to preserve structural continuity.

To address challenge (ii) related to noisy and sparse depth maps, image restoration (IR) technologies provide a feasible solution for enhancing image quality. Researchers have explored DL models to guide the restoration process, which can be classified into non-generative, GAN-based, and diffusion-based methods. Non-generative methods explicitly estimate the degradation parameters of an image (such as noise levels (Chen et al., 2019), blur kernels for super-resolution (Dong et al., 2014, 2015), or missing regions for inpainting (Cao et al., 2022)) and then restore high-quality images using the predicted information. However, these models often struggle to generalise to complex real-world degradation patterns due to oversimplified degradation assumptions. To address this, Generative Adversarial Networks (GANs) have been used to implicitly learn both the underlying data distribution and the degradation process, achieving promising results in tasks such as super-resolution (Fritsche et al., 2019; Wang et al., 2021a,b). Nonetheless, GAN-based methods are limited by training stability due to the use of dual-network architectures and intricate adversarial loss functions.

Recent advancements in Denoising Diffusion Probabilistic Models (DDPMs) have shown remarkable performance in synthesising visually compelling images (Ho et al., 2020; Song et al., 2020; Rombach et al., 2022). Leveraging these generative strengths, researchers have extended DDPMs to IR tasks, either by fine-tuning pre-trained DDPMs (Lin et al., 2024; Saharia et al., 2022; Zhu et al., 2024) or by training new DDPMs from scratch (Wang et al., 2023;
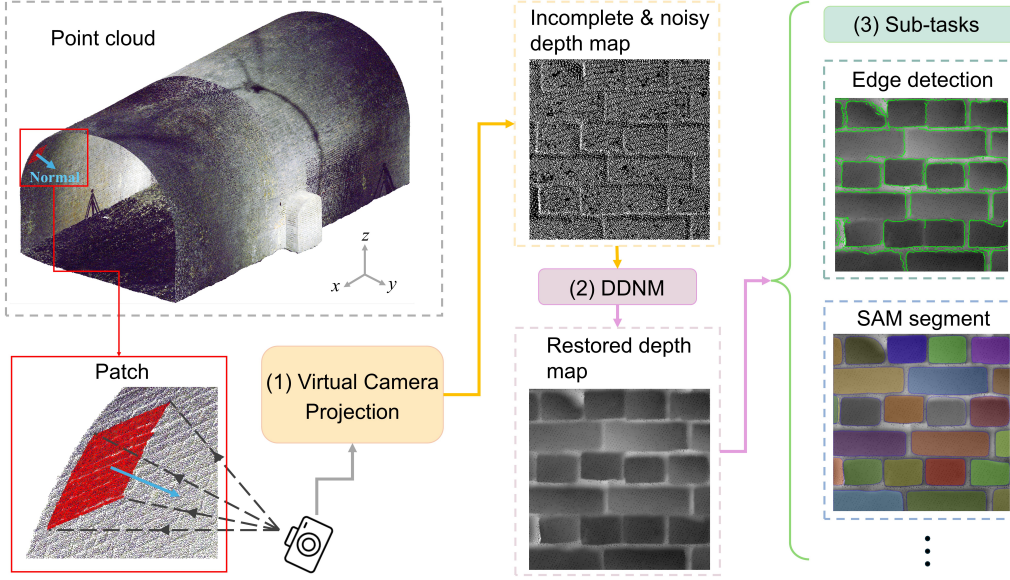
Figure 1: Overview of the InfraDiffusion pipeline.

Lo et al., 2024). However, most of the existing diffusion-based methods have focused on restoring RGB images, with limited attention on depth maps or geometric modalities. Recently, Lo et al. (2024) proposed RoofDiffusion to effectively restore roof height maps from severely incomplete depth maps by conditioning on building footprints to address sparsity and occlusion. While DDPMs achieve robust and effective IR tasks, they require paired training data comprising clean and corrupted images as additional conditions in the loss functions. This reliance limits their applicability to infrastructure point clouds, where ground-truth (GT) depth maps are rarely available due to the scarcity of dense geometric information.

To overcome the need for GT depth maps, we adapt the Denoising Diffusion Null-space Model (DDNM) (Wang et al., 2022), a zero-shot IR framework that uses pre-trained DDPMs to provide strong image priors. DDNM decomposes the restoration process into range-space and null-space components, ensuring strict data consistency by fixing the known degraded part and refining only the unknown null-space through sampling images from pre-trained DDPMs. This approach eliminates the need for task-specific training and enables effective restoration in the absence of clean depth maps. In our work, we adapt DDNM to restore sparse and noisy depth maps projected

from infrastructure point clouds to enable intricate downstream tasks such as brick-level segmentation.

To overcome the need for GT depth maps, we adapt the Denoising Diffusion Null-space Model (DDNM) (Wang et al., 2022), a zero-shot image restoration framework that leverages pre-trained DDPMs as strong image priors. DDNM decomposes the restoration process into range-space and null-space components, ensuring strict data consistency by fixing the known degraded part and refining only the unknown null-space through sampling from pre-trained DDPMs. In our work, we extend DDNM by introducing boundary masks that constrain the diffusion generation to regions where point cloud projections are valid. Our modification addresses a common challenge in utilising virtual camera projection in infrastructure point cloud datasets, where virtual camera viewpoints often lead to incomplete depth maps near the boundaries of point clouds.

By preventing spurious generation in empty regions, we introduce a complete zero-shot pipeline, i.e., InfraDiffusion, that restores sparse and noisy depth maps projected from infrastructure point clouds without task-specific training, thereby enabling fine-grained tasks such as brick-level segmentation. As illustrated in Figure 1, the pipeline begins by extracting surface patches from the point cloud and projecting them into 2D depth maps using virtual cameras. These incomplete and noisy depth maps are then restored using InfraDiffusion, which enhances depth quality without requiring GT images and fine-tuning. We validate InfraDiffusion on two representative datasets: a masonry tunnel with three different sampled sections and three masonry bridges, each presenting different challenges in terms of point cloud sparsity and geometric complexity. To evaluate the effectiveness of the restoration, we perform brick-level semantic segmentation using the Segment Anything Model (SAM) (Ye et al., 2024) with coordinate prompts to isolate individual bricks. The results demonstrate that the restored depth maps significantly improve segmentation accuracy. As a fully zero-shot framework, InfraDiffusion offers a scalable and adaptable solution for fine-grained geometric analysis of infrastructure point clouds in real-world inspection workflows.

The main contributions of this work can be summarised as follows:

(i) We propose InfraDiffusion, a zero-shot pipeline that leverages virtual camera projection and image restoration to improve the quality of depth maps derived from masonry infrastructure point clouds.

(ii) We extend the DDNM by introducing boundary masks that constrain

generation to valid projection regions, addressing boundary effects common in infrastructure point clouds and preventing spurious content.

(iii) We evaluate InfraDiffusion on two representative datasets (one masonry tunnel and two masonry bridges) featured by sparsity and geometric irregularity. By combining InfraDiffusion with the SAM for prompt-based zero-shot segmentation, we demonstrate substantial improvements in brick-level segmentation accuracy.

The paper is organised as follows. Section 2 introduces the InfraDiffusion pipeline in detail, including the virtual camera-based depth map projection and the adaptation of DDNM in InfraDiffusion for zero-shot depth restoration. Section 3 describes the two datasets used in this study: a masonry tunnel and three masonry bridges, both of which present real-world challenges in point cloud sparsity and geometric irregularity. Section 4 evaluates the performance of InfraDiffusion and prompt-based segmentation for assessing restoration quality. Finally, Section 5 concludes the paper with a discussion of our findings and potential directions for future research.

## 2. Methodology

### 2.1. Virtual camera projection of point clouds

As illustrated in Figure 2, the general framework which transfers infrastructure point clouds into depth maps is achieved in three steps: (1) normal estimation, (2) patch cropping, and (3) virtual camera projection (pinhole camera).

***Normal estimation.*** To improve computational efficiency, voxel-based downsampling with a voxel size of 0.2 m is first applied to the original point cloud $\mathbf{P} = \{\mathbf{x}_i \in \mathbb{R}^3\}_{i=1}^N$ to obtain a sparse representation, where $N$ represents the total point number. Surface normals are then estimated on the downsampled point cloud using a ball query with a radius of 0.4 m, implemented via the Open3D library (Zhou et al., 2018). This process generates a set of patch centre points $\mathbf{c}_j \in \mathbb{R}^3$ and their corresponding normal vectors $\mathbf{n}_j$, where $j$ represents the index of the centre points. These parameter values were selected based on empirical trials on real masonry infrastructure point clouds, which balance computational efficiency with robust normal estimation.
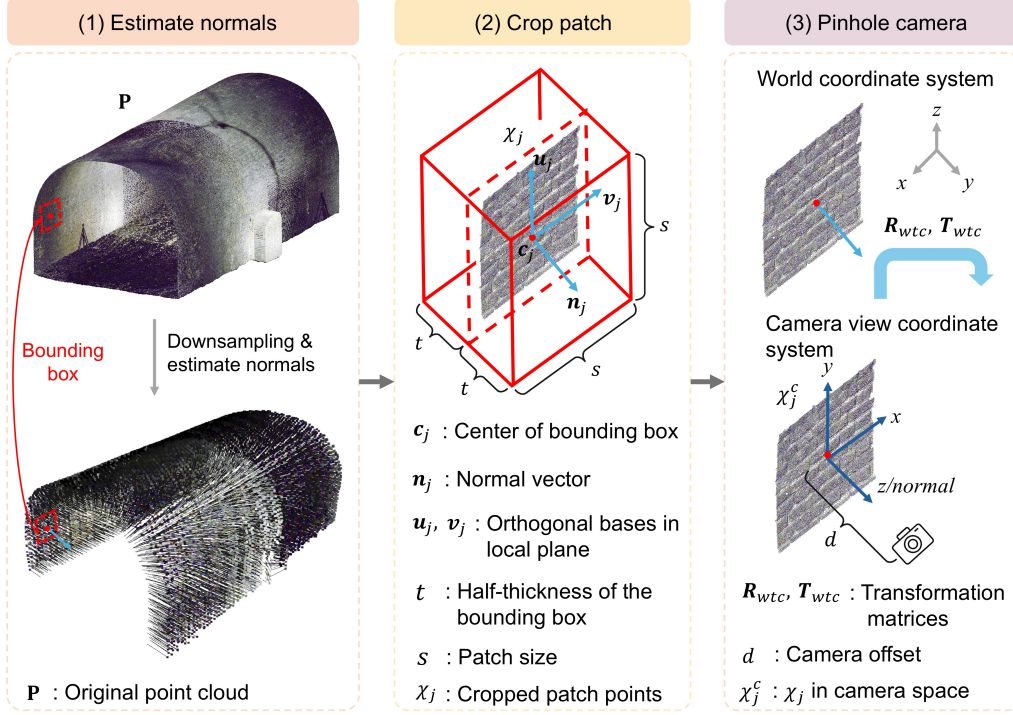
Figure 2: Overview of the patch extraction and projection using the pinhole camera model.

***Patch cropping.*** For each $\mathbf{c}_j$, we define a bounding box centered at $\mathbf{c}_j$, as illustrated in the red box in Figure 2. Each patch has a square face with side length $s$ and thickness $2t$. The specific values of these parameters are listed in Table 1. The added thickness accounts for geometric variation and ensures that sufficient points are captured near corners or regions where the surface is not shell-like. The cropped points $\mathcal{X}_j$ in the bounding box is defined as:

$$
\mathcal{X}_j = \left\{ \mathbf{x}_{ij} \in \mathbf{P} \;\middle|\; \left|(\mathbf{x}_{ij} - \mathbf{c}_j)^\top \mathbf{u}_j\right| \leq \frac{s}{2}, \; \left|(\mathbf{x}_{ij} - \mathbf{c}_j)^\top \mathbf{v}_j\right| \leq \frac{s}{2}, \right.
$$
$$
\left. \left|(\mathbf{x}_{ij} - \mathbf{c}_j)^\top \mathbf{n}_j\right| \leq t \right\}
\tag{1}
$$

where $\{\mathbf{u}_j, \mathbf{v}_j\}$ form an orthonormal basis spanning the local tangent plane orthogonal to $\mathbf{n}_j$ as shown in Figure 2. The parameter choices were determined through trial-and-error on real masonry infrastructure datasets to leverage between local completeness and boundary clarity of brick and mortar.

8

Table 1: Patch extraction and virtual camera parameters.

| Module | Parameter | Value |
|---|---|---|
| Crop patch | $s$ | 0.8 m |
| | $t$ | 0.25 m |
| Pinhole camera | $d$ | 0.8 m |
| | $(f_x, f_y)$ | (400, 400) pixels |
| | $(H, W)$ | $(256 \times 256)$ pixels |
| | $(c_x, c_y)$ | (128, 128) pixels |

***Pinhole camera***. Each $\mathcal{X}_j$ is projected onto a sparse and noisy 2D depth map, i.e., $\tilde{\mathbf{y}}_j \in \mathbb{R}^{H \times W}$, using a virtual pinhole camera model. The virtual camera is positioned a fixed distance $d$ from $\mathbf{c}_j$, with its optical axis aligned with $\mathbf{n}_j$, as shown in Figure 2. The camera origin $\mathbf{o}_j$ can thus be computed as:

$$\mathbf{o}_j = \mathbf{c}_j + d \cdot \mathbf{n}_j \tag{2}$$

The camera view coordinate system is constructed using $\{\mathbf{n}_j, \mathbf{u}_j, \mathbf{v}_j\}$. This defines the rotation matrix $\mathbf{R}_{wtc} \in \mathbb{R}^{3 \times 3}$ and the translation vector $\mathbf{T}_{wtc} = \mathbf{o}_j$, such that the world-to-camera transformation is:

$$\mathbf{x}_{ij}^c = \mathbf{R}_{wtc}(\mathbf{x}_{ij} - \mathbf{T}_{wtc}) \tag{3}$$

where $\mathbf{x}_{ij}^c = (x_{ij}^c, y_{ij}^c, z_{ij}^c)$ is the transformed point in the camera view coordinate system to form local patch $\mathcal{X}_j^c$.

$\mathcal{X}_j^c$ is then projected onto the image plane using the pinhole camera model with focal lengths $f_x, f_y$ and principal point $(c_x, c_y)$:

$$\begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} f_x \cdot \frac{x_{ij}^c}{z_{ij}^c} + c_x \\ f_y \cdot \frac{y_{ij}^c}{z_{ij}^c} + c_y \end{bmatrix} \tag{4}$$

where the depth value at pixel $(u, v)$ is taken as the $z$-coordinate of the point in the camera frame. The image resolution $(H, W)$, $d$, $f_x, f_y$, and $(c_x, c_y)$ are listed in Table 1. These values are selected to maximise coverage of $\mathcal{X}_j^c$ whilst preserving the necessary details of individual bricks.
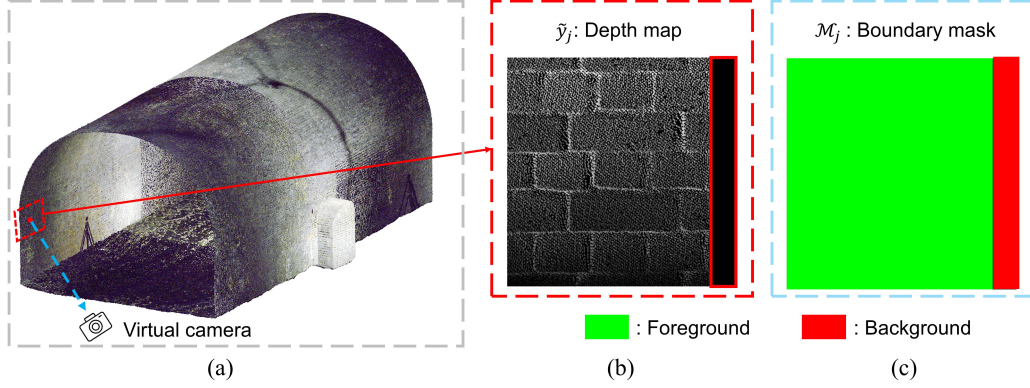
Figure 3: Boundary mask $\mathcal{M}_j$ extraction. (a) represents a virtual camera projection near the infrastructure boundary; (b) shows depth mask $\tilde{\mathbf{y}}_j$ with missing regions shown by the solid red box; (c) visualises the boundary mask $\mathcal{M}_j$ where the foreground and background are shown in green and red respectively.

Due to the fixed field of view of the virtual camera, the resulting image plane may include regions not covered by any 3D points from $\mathcal{X}_j^c$, particularly when the patch lies near the boundary of the structure. As illustrated in Figure 3a, the projection produces $\tilde{\mathbf{y}}_j$ with partially empty regions, as shown in Figure 3b (red solid box). These empty regions arise naturally from projecting irregular masonry geometries, and if left unaddressed, subsequent IR algorithms (e.g., DDNM) may introduce spurious bricks in empty regions or blur true structural boundaries.

To address this, we introduce a boundary mask $\mathcal{M}_j \subset \mathbb{R}^2$, defined as the axis-aligned 2D bounding box enclosing all projected pixels of $\mathcal{X}_j^c$:

$$\mathcal{M}_j = \text{BBox}\left(\left\{(u_i, v_i) \,\middle|\, \mathbf{x}_{ij}^c \in \mathcal{X}_j\right\}\right) \tag{5}$$

As shown in Figure 3c, pixels within $\mathcal{M}_j$ that correspond to valid projections are designated as the foreground (green), while pixels inside the mask but not covered by any projected point are labelled as the background (red). By explicitly separating valid and invalid regions, the boundary mask constrains restoration to physically meaningful areas, preventing the generation of non-existing bricks around structural edges.

### 2.2. Zero-shot depth map restoration with DDNM

To enable depth map restoration without task-specific training or clean GT supervision, we adapt the DDNM (Wang et al., 2022) conditioned on the

10

boundary mask $\mathcal{M}_j$. DDNM is a zero-shot IR framework that leverages the generative capacity of pre-trained DDPMs. DDPMs can learn statistical dependencies among pixels from large image datasets, thereby capturing strong contextual priors that can be exploited to reconstruct missing or degraded regions in depth maps.

This section begins with a brief overview of the classical IR problem. We then introduce the range–null space decomposition framework. Following this, we review the theoretical foundations of DDPM (Ho et al., 2020) and DDIM (Song et al., 2020), where DDIM eliminates the Markov assumption to improve sampling efficiency. Finally, we describe the DDNM algorithm and detail how it is adapted herein to restore sparse and noisy depth maps projected from masonry infrastructure point clouds.

### 2.2.1. Background of IR problem

IR is a long-standing inverse problem involving the recovery of a clean image $\mathbf{y}_j \in \mathbb{R}^{H \times W}$ from a degraded observation (Richardson, 1972; Andrews, 1974) (denoted as $\tilde{\mathbf{y}}_j$ in this paper). The degradation is typically modelled by a known linear operator $\mathbf{A}_j$, which represents corruption such as downsampling, missing pixels, blur, additive noise, or combinations thereof. The general observation model is expressed as:

$$\tilde{\mathbf{y}}_j = \mathbf{A}_j \mathbf{y}_j + \mathbf{n} \tag{6}$$

where $\mathbf{n} \in \mathbb{R}^{H \times W} \sim \mathcal{N}(\mathbf{0}, \sigma_y^2 \mathbf{I})$ denotes additive noise. Here, we address two canonical subproblems of IR, namely inpainting (data sparsity) and noise, by restoring $\mathbf{y}_j$ through InfraDiffusion. The degradation operator $\mathbf{A}_j$ is instantiated as a binary spatial mask derived from the projection of the point cloud, indicating which pixels in $\tilde{\mathbf{y}}_j$ correspond to valid 3D points and which are unobserved.

Restoring $\mathbf{y}_j$ from $\tilde{\mathbf{y}}_j$ is inherently ill-posed due to severe sparsity and noise during projection and point cloud data collection. Traditional non-generative IR approaches (Dong et al., 2014) often rely on task-specific training with paired data (e.g., $\mathbf{y}_j$ and $\tilde{\mathbf{y}}_j$), which is impractical for infrastructure due to the dearth of $\mathbf{y}_j$, and often yields poor results in real-world scenarios. To overcome this, we adapt the zero-shot DDNM restoration method guided by the strong generative priors of DDPM, enabling robust and accurate restoration without supervised training.

*2.2.2. Range-null space decomposition in noise-free IR problem*

We begin our analysis with the noise-free IR setting, where the simplification allows for a more interpretable formulation of the range-null space decomposition, which is represented as:

$$\tilde{\mathbf{y}}_j = \mathbf{A}_j \mathbf{y}_j \tag{7}$$

The objective is to reconstruct a plausible depth image $\hat{\mathbf{y}}_j$ that satisfies two essential conditions:

$$\text{Consistency:} \quad \mathbf{A}_j \hat{\mathbf{y}}_j = \tilde{\mathbf{y}}_j \tag{8}$$
$$\text{Realness:} \quad \hat{\mathbf{y}}_j \sim q(\mathbf{y}) \tag{9}$$

where *Consistency* ensures $\hat{\mathbf{y}}_j$ satisfies the degradation relationship, and *Realness* determines whether $\hat{\mathbf{y}}_j$ is sampled from the clean depth image distribution, i.e., $q(\mathbf{y})$.

To analyse the solution space, DDNM (Wang et al., 2022) decomposes $\hat{\mathbf{y}}_j$ into its range and null spaces of $\mathbf{A}_j$:

$$\hat{\mathbf{y}}_j = \mathbf{A}_j^\dagger \mathbf{A}_j \hat{\mathbf{y}}_j + (\mathbf{I} - \mathbf{A}_j^\dagger \mathbf{A}_j)\hat{\mathbf{y}}_j \tag{10}$$

where $\mathbf{A}_j^\dagger$ is the Moore–Penrose pseudoinverse of $\mathbf{A}_j$, which can be derived via singular value decomposition (Golub and Reinsch, 1971). $\mathbf{A}_j^\dagger$ for the IR problem has been properly introduced in DDPM (Wang et al., 2022), which is omitted in this paper for simplicity. By letting:

$$\mathbf{y}_j^{\text{range}} := \mathbf{A}_j^\dagger \mathbf{A}_j \hat{\mathbf{y}}_j, \quad \mathbf{y}_j^{\text{null}} := (\mathbf{I} - \mathbf{A}_j^\dagger \mathbf{A}_j)\bar{\mathbf{y}}_j \tag{11}$$

where the $\hat{\mathbf{y}}_j$ in the second item of Equation 10 is changed to $\bar{\mathbf{y}}_j$ as it is generated from DDPM. The Equation 10 is re-written as:

$$\hat{\mathbf{y}}_j = \mathbf{y}_j^{\text{range}} + \mathbf{y}_j^{\text{null}} \tag{12}$$

Applying the degradation operator $\mathbf{A}_j$ to both sides, we obtain:

$$\mathbf{A}_j \hat{\mathbf{y}}_j = \mathbf{A}_j \mathbf{y}_j^{\text{range}} + \mathbf{A}_j \mathbf{y}_j^{\text{null}} = \mathbf{A}_j \mathbf{y}_j^{\text{range}} = \mathbf{A}_j \hat{\mathbf{y}}_j \tag{13}$$

since $\mathbf{A}_j \mathbf{y}_j^{\text{null}} = \mathbf{A}_j(\mathbf{I} - \mathbf{A}_j^\dagger \mathbf{A}_j)\bar{\mathbf{y}}_j = \mathbf{0}$, $\hat{\mathbf{y}}_i$ is not influenced by any $\bar{\mathbf{y}}_j$ after being applied by $\mathbf{A}_j$; all observable content is confined to $\mathbf{y}_j^{\text{range}}$, which guarantees the *Consistency* constraint. However, $\mathbf{y}_j^{\text{null}}$ determines whether it is perceptually plausible, i.e., whether it satisfies the *Realness* criterion by resembling a sample from $\hat{\mathbf{y}}_j \sim q(\mathbf{y})$. To achieve *Realness*, DDNM leverages the generative capability of DDPM to sample $\bar{\mathbf{y}}_j$.

### 2.2.3. Review of DDPM and DDIM

DDPMs (Ho et al., 2020) are a class of generative models that reconstruct data by reversing a gradual noising process. The coral idea is to learn how to transform random noise into realistic image samples by exploiting the statistical dependencies observed in training data. Let $\mathbf{y}_0 \in \mathbb{R}^{H \times W}$ denote the original clean image. The forward diffusion process progressively adds Gaussian noise to $\mathbf{y}_0$ over $T$ steps, eventually converting it into approximately pure noise drawn from a standard Gaussian distribution:

$$q(\mathbf{y}_t \mid \mathbf{y}_{t-1}) = \mathcal{N}(\mathbf{y}_t; \sqrt{1 - \beta_t}\, \mathbf{y}_{t-1},\, \beta_t \mathbf{I}) \tag{14}$$

where $\{\beta_t\}_{t=1}^T$ is a predefined variance schedule. By using the reparameterization trick, noisy samples at step $t$ can be written as:

$$\mathbf{y}_t = \sqrt{\bar{\alpha}_t}\, \mathbf{y}_0 + \sqrt{1 - \bar{\alpha}_t}\, \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{15}$$

where $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$. The generative process aims to reverse the diffusion process by sampling $\mathbf{y}_{t-1}$ from $\mathbf{y}_t$ via Bayes' rule:

$$p(\mathbf{y}_{t-1} \mid \mathbf{y}_t, \mathbf{y}_0) = q(\mathbf{y}_t \mid \mathbf{y}_{t-1}, \mathbf{y}_0) \cdot \frac{q(\mathbf{y}_{t-1} \mid \mathbf{y}_0)}{q(\mathbf{y}_t \mid \mathbf{y}_0)} \tag{16}$$

$$= \mathcal{N}(\mathbf{y}_{t-1}; \boldsymbol{\mu}_t(\mathbf{y}_t, \mathbf{y}_0), \sigma_t^2 \mathbf{I}) \tag{17}$$

where $p(\mathbf{y}_{t-1} \mid \mathbf{y}_t, \mathbf{y}_0)$ still represents a Gaussian distribution. The mean $\boldsymbol{\mu}_t(\mathbf{y}_t, \mathbf{y}_0)$ and variance $\sigma_t^2$ are derived from the forward process and can be analytically computed given $\{\beta_t\}_{t=1}^T$. In practice, rather than learning the full reverse distribution directly, DDPM trains a neural network $\boldsymbol{\epsilon}_\theta(\mathbf{y}_t, t)$ to predict the added noise $\boldsymbol{\epsilon}$. The training objective then becomes:

$$\mathcal{L}_{DDPM} = \mathbb{E}_{\mathbf{y}_0, \boldsymbol{\epsilon}, t} \left[ \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t}\, \mathbf{y}_0 + \sqrt{1 - \bar{\alpha}_t}\, \boldsymbol{\epsilon}, t) \right\|^2 \right] \tag{18}$$

Equation 18 enables the network to learn to remove noise at arbitrary timesteps $t$, and forms the foundation for image generation and restoration via iterative denoising.

While DDPM achieves high-quality generative performance, its sampling process is slow due to the need for thousands of denoising steps. This inefficiency arises from its reliance on a Markovian reverse process, i.e., $q(\mathbf{y}_t \mid \mathbf{y}_{t-1}, \mathbf{y}_0) = q(\mathbf{y}_t \mid \mathbf{y}_{t-1})$, which requires step-wise resampling from Gaussian distributions.

To accelerate reverse sampling, DDIM (Song et al., 2020) removes the Markov assumption and reformulates the reverse process as a non-stochastic transformation. This allows the use of an arbitrary set of decreasing timesteps $\{\tau_1, \ldots, \tau_K\} \subset \{1, \ldots, T\}$, enabling larger step sizes in sampling. The DDIM sampling rule is given by:

$$\mathbf{y}_{\tau_{k-1}} = \sqrt{\bar{\alpha}_{\tau_{k-1}}} \left( \frac{\mathbf{y}_{\tau_k} - \sqrt{1 - \bar{\alpha}_{\tau_k}} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{y}_{\tau_k}, \tau_k)}{\sqrt{\bar{\alpha}_{\tau_k}}} \right)$$
$$+ \sqrt{1 - \bar{\alpha}_{\tau_{k-1}} - \sigma_{\tau_k}^2} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{y}_{\tau_k}, \tau_k) + \sigma_{\tau_k} \cdot \boldsymbol{\epsilon} \qquad (19)$$

where the setting of variance $\sigma_{\tau_k} \in [0, 1]$ controls the amount of stochasticity. DDIM preserves the forward process and training objective (gradient of $\mathcal{L}_{DDPM}$) of DDPM but replaces the stochastic differential equation (SDE) with an equivalent ordinary differential equation (ODE), enabling faster and controllable sampling.

### 2.2.4. InfraDiffusion with boundary-constrained DDNM in noise IR problem

We now consider the IR problem in the presence of noise, which is given by Equation 6. During the diffusion process, we denote $\mathbf{y}_{0|t}$ as the predicted clean image at timestep $t$ from DDPM. Following the DDNM framework (Wang et al., 2022), the range–null space decomposition can be expressed as:

$$\hat{\mathbf{y}}_{j,0|t} = \mathbf{A}_j^\dagger \tilde{\mathbf{y}}_j + (\mathbf{I} - \mathbf{A}_j^\dagger \mathbf{A}_j)\mathbf{y}_{0|t} = \mathbf{y}_{0|t} - \mathbf{A}_j^\dagger(\mathbf{A}_j \mathbf{y}_{0|t} - \mathbf{A}_j \mathbf{y}_j) + \mathbf{A}_j^\dagger \mathbf{n} \qquad (20)$$

where the second term, i.e., $\mathbf{A}_j^\dagger(\mathbf{A}_j \mathbf{y}_{0|t} - \mathbf{A}_j \mathbf{y}_j)$, is a correction in the range-space that ensures *Consistency*. $\mathbf{A}_j^\dagger \mathbf{n} \in \mathbb{R}^{H \times W}$ propagates the observation noise into $\hat{\mathbf{y}}_{j,0|t}$ and subsequently into $\mathbf{y}_{t-1}$ during sampling.

To align with the noise scale required by DDPM (Ho et al., 2020), we adopt the DDNM formulation and reformulate the update as:

$$\hat{\mathbf{y}}_{j,0|t} = \mathbf{y}_{0|t} - \boldsymbol{\Sigma}_t \mathbf{A}_j^\dagger(\mathbf{A}_j \mathbf{y}_{0|t} - \tilde{\mathbf{y}}_j) \qquad (21)$$

where $\boldsymbol{\Sigma}_t$ is a time-dependent scaling matrix determined by $\sigma_y$ and the DDPM variance schedule (see (Wang et al., 2022) for full derivation).

Different from the original DDNM, InfraDiffusion incorporate $\mathbf{M}_j$ as an extra condition to constrain the correction region. This restricts the correction region to valid projections, thereby preventing spurious generations in

---

**Algorithm 1** InfraDiffusion sampling with boundary masks

---

1: **Input:** Noisy image $\mathbf{y}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, known data $\tilde{\mathbf{y}}_j$, degradation matrix $\mathbf{A}_j$, mask $\mathcal{M}_j$, step schedule $\{\tau_k\}_{k=1}^K$

2: **for** $k = K, \ldots, 1$ **do**

3: $\quad \mathbf{y}_{0|\tau_k} = \frac{1}{\sqrt{\bar{\alpha}_{\tau_k}}} \left( \mathbf{y}_{\tau_k} - \sqrt{1 - \bar{\alpha}_{\tau_k}} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{y}_{\tau_k}, \tau_k) \right)$

4: $\quad \hat{\mathbf{y}}_{j,0|\tau_k} = \mathbf{M}_j \odot \left( \mathbf{y}_{0|\tau_k} - \mathbf{\Sigma}_{\tau_k} \mathbf{A}_j^\dagger (\mathbf{A}_j \mathbf{y}_{0|\tau_k} - \tilde{\mathbf{y}}_j) \right)$

5: $\quad \mathbf{y}_{\tau_{k-1}} = \sqrt{\bar{\alpha}_{\tau_{k-1}}} \cdot \hat{\mathbf{y}}_{j,0|\tau_k}$ $\qquad\qquad\qquad\qquad \triangleright\ \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$

$\qquad\qquad + \sqrt{1 - \bar{\alpha}_{\tau_{k-1}} - \sigma_{\tau_k}^2} \cdot \boldsymbol{\epsilon}_\theta(\mathbf{y}_{\tau_k}, \tau_k)$

$\qquad\qquad + \sigma_{\tau_k} \cdot \boldsymbol{\epsilon}$

6: **end for**

7: **Return:** Reconstructed image $\mathbf{y}_0$

---

empty areas and ensuring that denoising and inpainting occur only within physically meaningful regions:

$$\hat{\mathbf{y}}_{j,0|t} = \mathbf{M}_j \odot \left( \mathbf{y}_{0|t} - \mathbf{\Sigma}_t \mathbf{A}_j^\dagger (\mathbf{A}_j \mathbf{y}_{0|t} - \tilde{\mathbf{y}}_j) \right) \tag{22}$$

During the sampling, $\hat{\mathbf{y}}_{j,0|t}$ is used to produce the $\mathbf{y}_{t-1}$, where any hallucinated content outside the valid support at step $t$ is propagated forward by both the noise prediction network and the range-space correction without $\mathbf{M}_j$. Pixels outside $\mathcal{M}_j$ lie in the null space of $\mathbf{A}_j$ and, without constraints, are filled by the generative prior at every step. Applying $\mathcal{M}_j$ *within* the update suppresses these contributions at each iteration, preserves sharp boundaries, and prevents the accumulation of boundary artefacts that a post-hoc mask cannot remove.

To perform efficient sampling, we adopt the DDIM framework (Song et al., 2020) instead of DDPM. The complete InfraDiffusion sampling procedure is summarised and explained in Algorithm 1.

## 3. Dataset

We evaluate the proposed InfraDiffusion framework on two types of masonry infrastructure: a historic railway tunnel and two masonry bridges. All point clouds were captured via terrestrial laser scanning (TLS) and semantically segmented into structural components. To avoid extensive computation while demonstrating our method, we sample a subset of representative
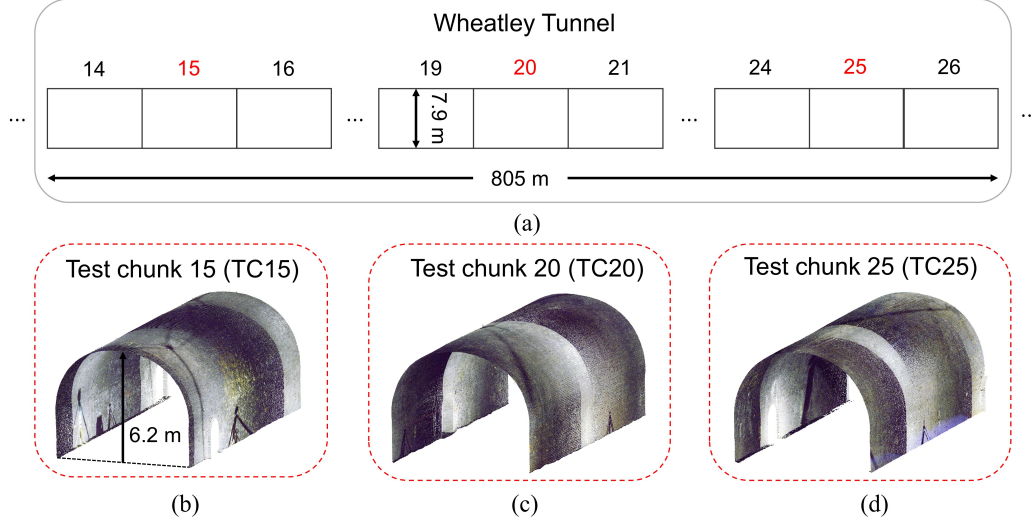
Figure 4: Wheatley Tunnel overview and selected test chunks. (a) Schematic plan view of the 805 m tunnel divided into 51 chunks, where chunks 15, 20, and 25 are selected for testing. (b)–(d) represent point cloud views of the selected chunks (e.g., TC15, TC20, and TC25, respectively).

chunks from the tunnel and randomly select patches across components in all datasets.

### 3.1. Wheatley Tunnel

The Wheatley Tunnel is a historic masonry railway tunnel located on the Halifax High Level line, spanning a total length of 805 m. The original point cloud was collected by National Highways as part of a structural health monitoring initiative. According to inspection records, the tunnel remains in overall fair condition but exhibits widespread dampness, soft mortar loss, and calcite deposits throughout the bore lining. Signs of spalling are also observed in the arch and walls.

As shown in Figure 4a-d, the tunnel was divided into 51 contiguous chunks, each approximately 15 m in length, due to the substantial volume of point cloud data. For evaluation, we selected three representative chunks (e.g., TC15, TC20, and TC25), covering structurally diverse regions of the tunnel. The typical cross-sectional dimensions of the tunnel are approximately 7.9 m in width and 6.2 m in height (Figure 4a and b).

To ensure unbiased coverage across different structural components, $\mathbf{c}_j$ (e.g., patch centre as shown in Figure 2) is randomly sampled from the full

16

Table 2: Summary of selected Wheatley Tunnel chunks with point counts, structural composition, and number of extracted patches.

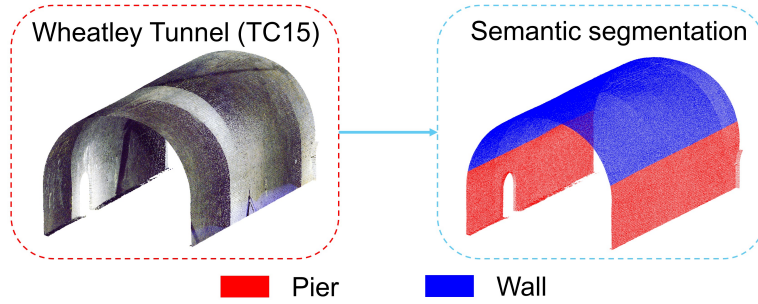| Chunk | Points | Arch | Wall | Patch |
|-------|--------|------|------|-------|
| TC15 | 55,821,396 | 1 | 2 | 36 |
| TC20 | 32,278,956 | 1 | 2 | 25 |
| TC25 | 37,756,228 | 1 | 2 | 39 |



Figure 5: Semantic segmentation of the Wheatley Tunnel (TC15) point cloud. The arch and two walls are labelled in blue and red, respectively.

tunnel point cloud, without preference toward any specific region. Due to the large number of available patches, we sample a total of 100 patches from all three chunks for evaluation, given that they belong to the same tunnel. The distribution of selected patches across the three chunks is shown in Table 2. Each selected chunk includes a typical masonry composition of one arch barrel and two supporting walls, which are semantically segmented into distinct regions for analytical purposes. As shown in Figure 5, the segmentation allows us to assess how projection quality and IR performance vary across these structural elements. Such differences are often introduced by the varying distance between the LiDAR scanner and the scanned surface; regions farther from the scanner (e.g., crown or upper arch) tend to produce noisier and more diffuse depth images than closer surfaces (e.g., walls).

*3.2. Masonry bridges*

We also evaluate our method on two representative masonry bridge point clouds (Figure 6): (1) an anonymised UK single-span masonry bridge acquired from an external industry partner, and (2) Hertford Viaduct, which is a multi-span masonry bridge (Jing et al., 2022). Both bridges were segmented
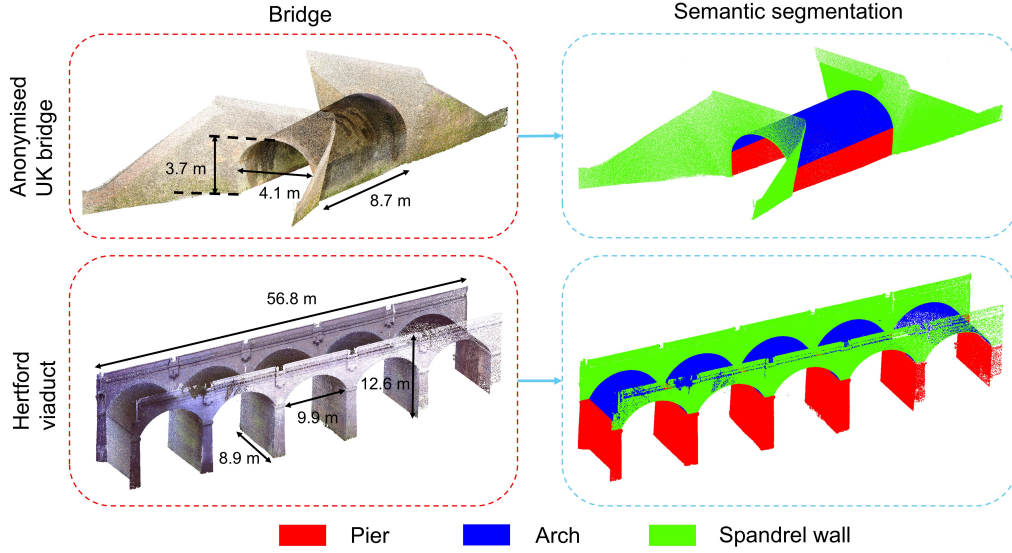
Figure 6: Point clouds and semantic segmentation of anonymised UK bridge and Hertford Viaduct. Dimensions annotated in the raw point cloud aid understanding of the structural scale and data volume. The arch, pier, and spandrel wall are coloured red, blue, and green, respectively.

into structural components (arches, piers, and spandrel walls) to compare the depth image quality in different components. Figure 6 also shows the overall geometry and scale of each bridge, which includes key structural dimensions. The semantic segmentation uses blue for arches, red for piers, and green for spandrel walls to distinguish structural components across the bridge geometry.

Table 3 summarises the number of points, structural components, and sampled patches for each bridge. Patch selection was randomly conducted across all components to ensure structural diversity while maintaining a tractable evaluation size.

## 4. Experiments and results

### 4.1. Experiment Settings

We employ a pre-trained, unconditional DM released by OpenAI (Dhariwal and Nichol, 2021) as the generative prior for InfraDiffusion. All experiments are conducted in a zero-shot manner, without any task-specific fine-tuning or additional training.

18

Table 3: Summary of masonry bridge datasets with point counts, component counts, and number of patches sampled.

| Bridge | Points | Arch | Pier | Spandrel wall | Patch |
|---|---|---|---|---|---|
| Anonymised UK bridge | 74,183,011 | 1 | 2 | 2 | 60 |
| Hertford Viaduct | 98,629,976 | 5 | 6 | 2 | 140 |

Table 4: Patch indices selected for qualitative IR evaluation across arch and wall/pier components. Arch patches are visualised in Figure 7a-j and wall/pier patches in Figure 8a-j.

| Infrastructure | Arch patch indices | Wall/Pier patch indices |
|---|---|---|
| Wheatley Tunnel (TC15) | 1828 | 969 |
| Wheatley Tunnel (TC20) | 807 | 270 |
| Wheatley Tunnel (TC25) | 2352 | 986 |
| Anonymised UK bridge | 509 | 269 |
| Hertford Viaduct | 703 | 145 |

Input depth images with a fixed resolution of $256 \times 256$ (see Table 1) are normalised to the $[0, 1]$ range, which was empirically found to yield better restoration quality than normalising to $[-1, 1]$. The measurement noise in $\mathbf{y}$ is modelled using a standard deviation of $\sigma_y = 0.16$, which provides a stable assumption for subsequent restoration. Inference is performed on a laptop equipped with an NVIDIA RTX 4090 GPU (16 GB VRAM).

The model architecture consists of a U-Net backbone with a channel depth of 256, two residual blocks per level, and self-attention mechanisms applied at spatial resolutions of 32, 16, and 8, where more details can be found at Dhariwal and Nichol (2021). The diffusion process utilizes a linear $\beta$ schedule ranging from $\beta_{\text{start}} = 10^{-4}$ to $\beta_{\text{end}} = 0.02$ over 1000 steps. For sampling, we adopt the DDIM inversion strategy, reducing the number of denoising steps to 100 using 10-step intervals to accelerate inference.

### 4.2. Depth map restoration results

Our evaluation of depth map restoration by InfraDiffusion is conducted in a qualitative manner, since quantitative GT supervision (dense depth maps) is unavailable. Instead, we demonstrate the effectiveness of IR through the
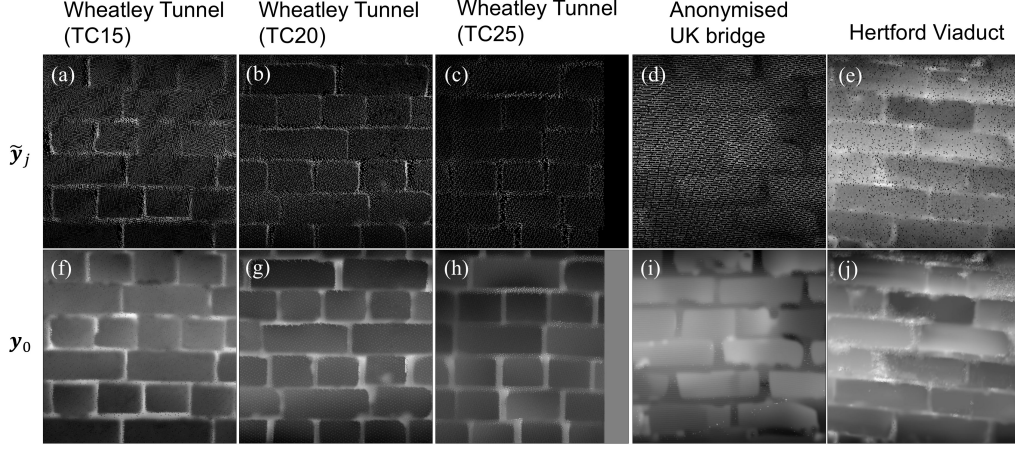
Figure 7: IR results for arch components across five infrastructures with $\sigma_y = 0.16$. (a)–(e) show degraded inputs $\tilde{\mathbf{y}}_j$, and (f)–(j) present restored images $\mathbf{y}_0$. (a, f), (b, g), and (c, h) are randomly sampled from Wheatley Tunnel (TC15, TC20, and TC25). (d, i) and (e, j) correspond to anonymised UK bridge and Hertford Viaduct.
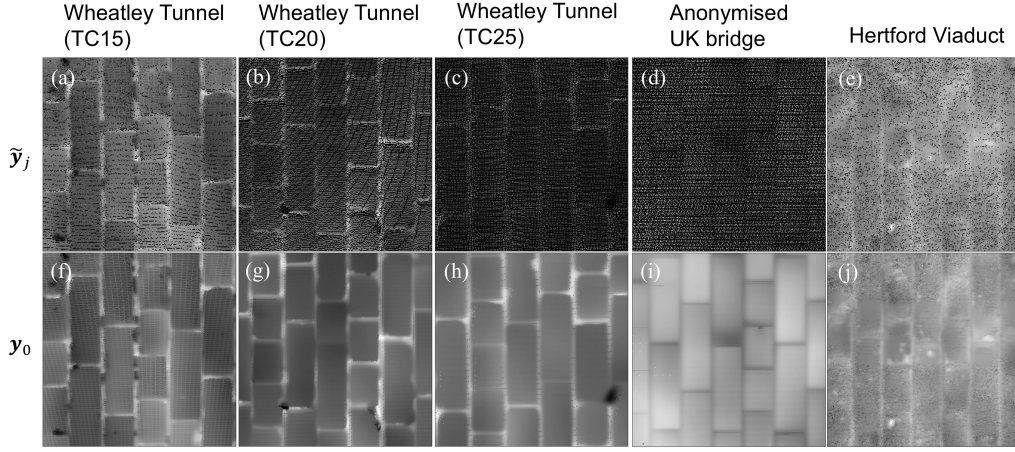


Figure 8: IR results for wall/pier components across five infrastructures with $\sigma_y = 0.16$. (a)–(e) show degraded inputs $\tilde{\mathbf{y}}_j$, and (f)–(j) present restored images $\mathbf{y}_0$. (a, f), (b, g), and (c, h) are randomly sampled from Wheatley Tunnel (TC15, TC20, and TC25). (d, i) and (e, j) correspond to anonymised UK bridge and Hertford Viaduct.

semantic segmentation by leveraging zero-shot and prompt-based SAM. The IR results are demonstrated in Figure 7a-j and Figure 8a-j across the five cases: the three selected Wheatley Tunnel chunks (TC15, TC20, and TC25),

20

the anonymised UK bridge, and the Hertford Viaduct, where $\sigma_y = 0.16$. The selected patch indices are summarised in Table 4, which correspond to Figures 7–9.

The different orientations of $\tilde{\mathbf{y}}_j$ between arch and wall/pier arise from the default generatrix direction $\mathbf{v}_j$ (as shown in Figure 2) used during virtual camera projection. The different orientations of $\tilde{\mathbf{y}}_j$ between arch and wall/pier arise from the default generatrix direction $\mathbf{v}_j$ (as shown in Figure 2) used during virtual camera projection. While this choice influences the visual orientation and the density distribution of projected points, we observed that InfraDiffusion produces consistently plausible restorations across viewpoints given robust normal estimations.

By comparing results for the anonymised UK bridge (Figure 7d, e, i, and j) and Hertford Viaduct (Figure 8d, e, i, and j), the arch surfaces exhibit blurrier and less distinct brick boundaries than the pier surfaces. This discrepancy is mainly due to the TLS acquisition geometry discussed in Section 3.1, where arches are scanned from greater distances than piers, resulting in lower point density and higher noise.

We further examine the influence of varied $\sigma_y$ on restoration quality in Figure 9a–t. Results are illustrated for both the Wheatley Tunnel (arch patch 1828 and wall patch 969) and the anonymised UK bridge (arch patch 703 and pier patch 145), which indicate that the assumed noise level in datasets strongly influences the denoising strength. For $\sigma_y = 0$, the restored images still exhibit residual speckle and blurred mortar boundaries (see Figure 9b, g, l, and q). As $\sigma_y$ increases to 0.16, the background noise is effectively suppressed while brick boundaries remain clear and continuous, producing the most visually coherent results (Figure 9c, h, m, and r). At $\sigma_y = 0.32$ and $\sigma_y = 0.64$, the denoising starts to blur the brick edges within the red boxes in Figure 9b–e progressively, and fine boundary details are lost despite smoother textures.

Based on this observation, we select $\sigma_y = 0.16$ for all subsequent experiments. This choice represents a practical balance between noise suppression and boundary preservation. We note, however, that the optimal value is empirically and visually determined and may vary across datasets due to differences in point density, registration errors, and sensor noise. For consistency and ease of comparison, we apply the same $\sigma_y$ across all cases in this study.
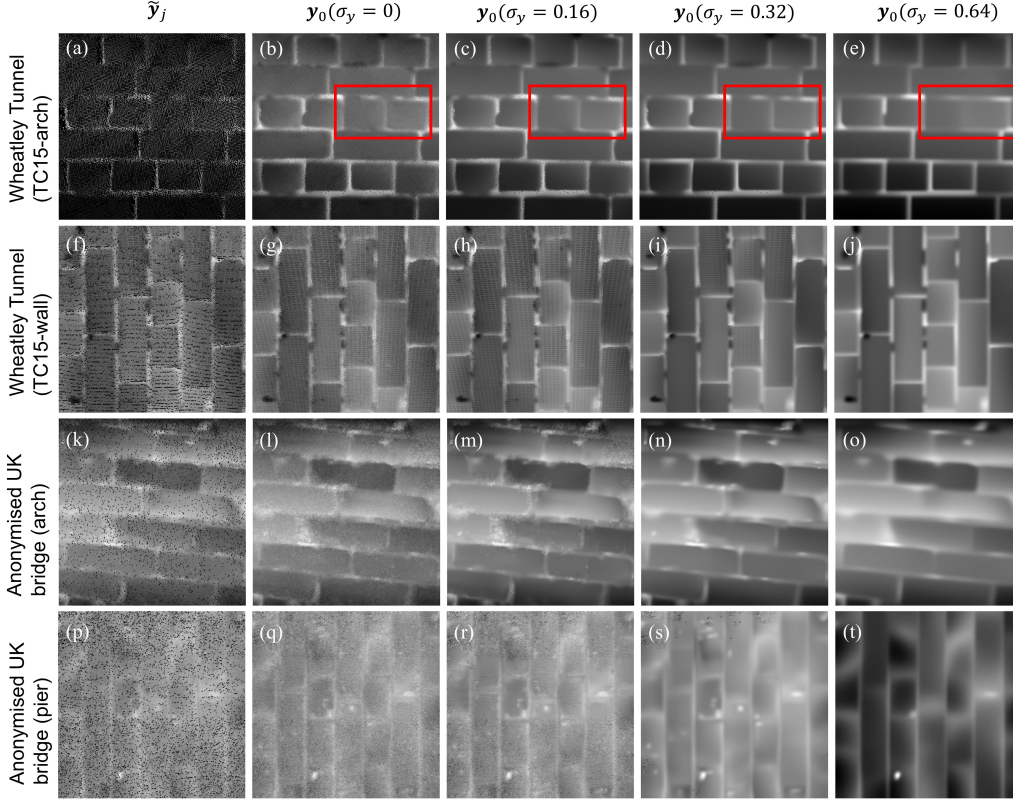
Figure 9: Effect of varying $\sigma_y \in \{0, 0.16, 0.32, 0.64\}$ on IR quality. (a)–(e) and (f)–(j) show $\tilde{\mathbf{y}}_j$ and $\mathbf{y}_0$ for the arch and walls of Wheatley Tunnel (TC15) (patch 1828 and patch 969, respectively). (k)–(o) and (p)–(t) show corresponding results for the arch and pier of the anonymised UK bridge (patch 703 and patch 145, respectively). The red boxes from (b)-(e) indicate the blurring of brick boundaries by increasing $\sigma_y$.

## 4.3. Ablation test on InfraDiffusion

To demonstrate the effectiveness of the proposed InfraDiffusion conditioned on $\mathcal{M}_j$, we compare $\mathbf{y}_0$ from the vanilla DDNM and our InfraDiffusion framework. All selected patches are taken from walls/piers, where structural boundaries are common, such as wall/pier edges or regions connected to the ground. Figure 10a-p shows four representative cases: Wheatley Tunnel (patches 897 and 1194 of TC15) and the anonymised UK bridge (patches 265 and 476). For each case, we present the $\mathcal{M}_j$, $\tilde{\mathbf{y}}_j$, and the corresponding IR results (e.g., $\mathbf{y}_0$) using DDNM and InfraDiffusion.

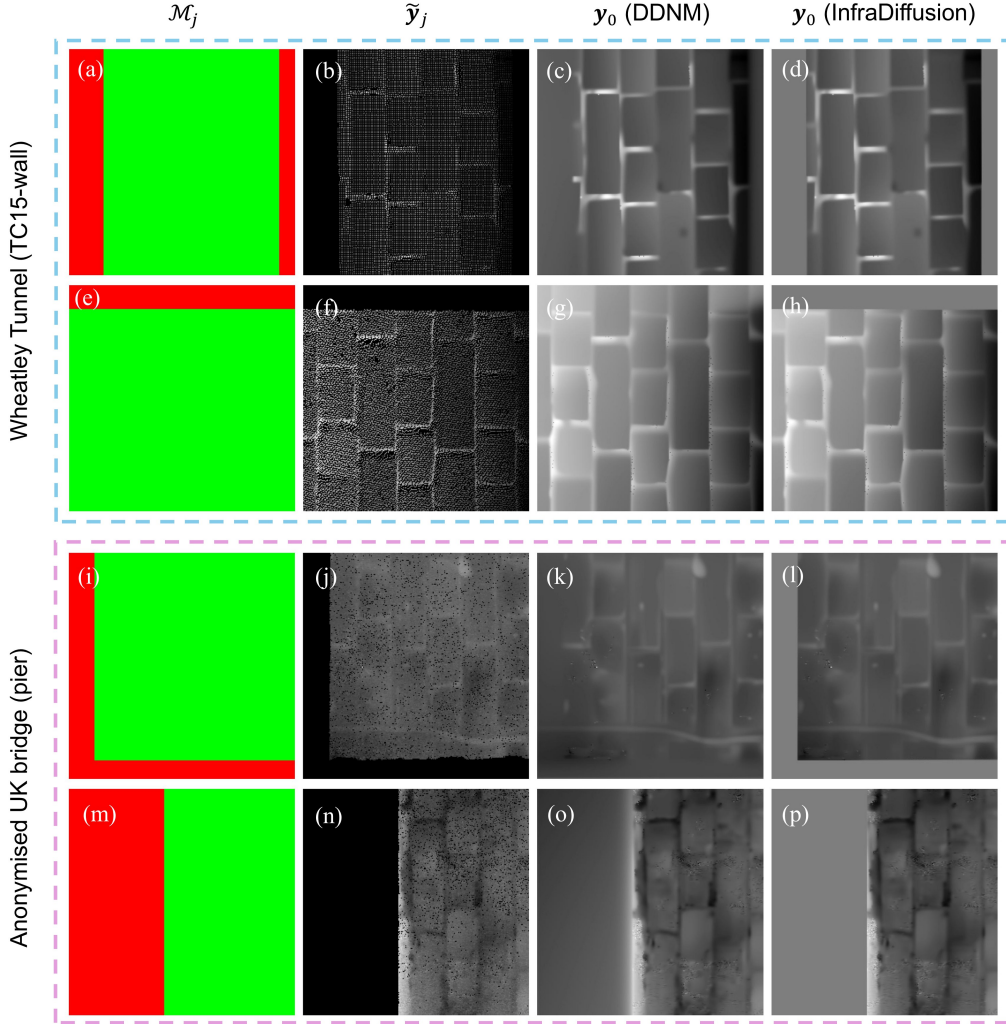The results indicate that DDNM, while achieving the same results in

Figure 10: Comparison between DDNM and InfraDiffusion on wall/pier patches. (a)–(d) and (e)–(h) represent Wheatley Tunnel (TC15) wall patches 897 and 1194, respectively. (i)–(l) and (m)–(p) correspond to anonymised UK bridge pier patches 265 and 479. For each patch, the four columns show: boundary mask $\mathcal{M}_j$ (colours are defined in Figure 3), projected depth map $\tilde{\mathbf{y}}_j$, restoration $\mathbf{y}_0$ by DDNM, and restoration by InfraDiffusion.

foreground regions (defined in Figure 3), often produces spurious structures beyond the valid projection region (e.g., background regions). This is particularly evident in Figures 10g and 10k, where non-existent bricks are generated outside the wall/pier boundary. By contrast, InfraDiffusion con-
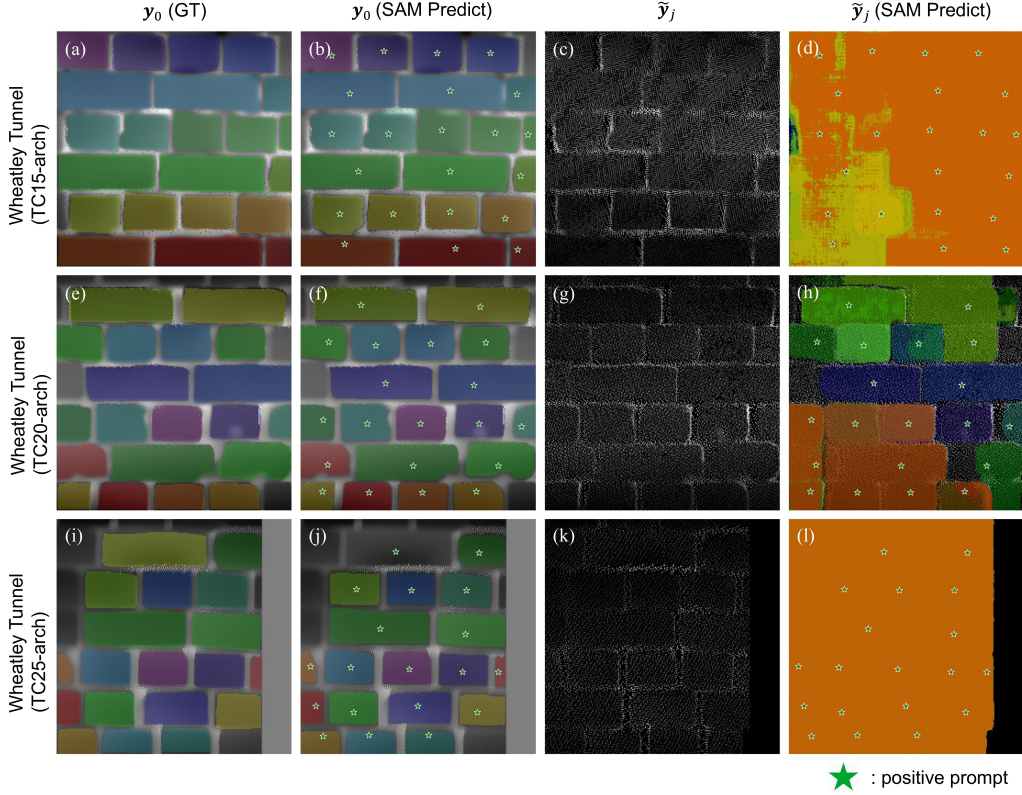
Figure 11: Zero-shot segmentation of arch patches from Wheatley Tunnel using SAM with ViT-H. (a)–(d) show GT ($\mathbf{y}_0$), SAM predictions on $\mathbf{y}_0$, $\tilde{\mathbf{y}}_j$, and SAM predictions on $\tilde{\mathbf{y}}_j$ for test chunk 15. (e)–(h) and (i)–(l) correspond to test chunks 20 and 25, respectively. Green pentagrams indicate positive prompts.

strains restoration strictly within the masked region, preserving true edges and avoiding hallucinated patterns (Figures 10d, h, l, and p). Across all cases, InfraDiffusion yields more physically consistent restorations by preserving the original spatial structures of point clouds.

## 4.4. Zero-shot segmentation with SAM

We further evaluate the effectiveness of the InfraDiffusion for downstream segmentation tasks using SAM (Ye et al., 2024). Specifically, we adopt the ViT-H SAM model weights to generate segmentation masks for masonry brick boundaries. GT annotations are created only for bricks with clear and well-defined boundaries, while instances with blurred, indistinct, and incomplete
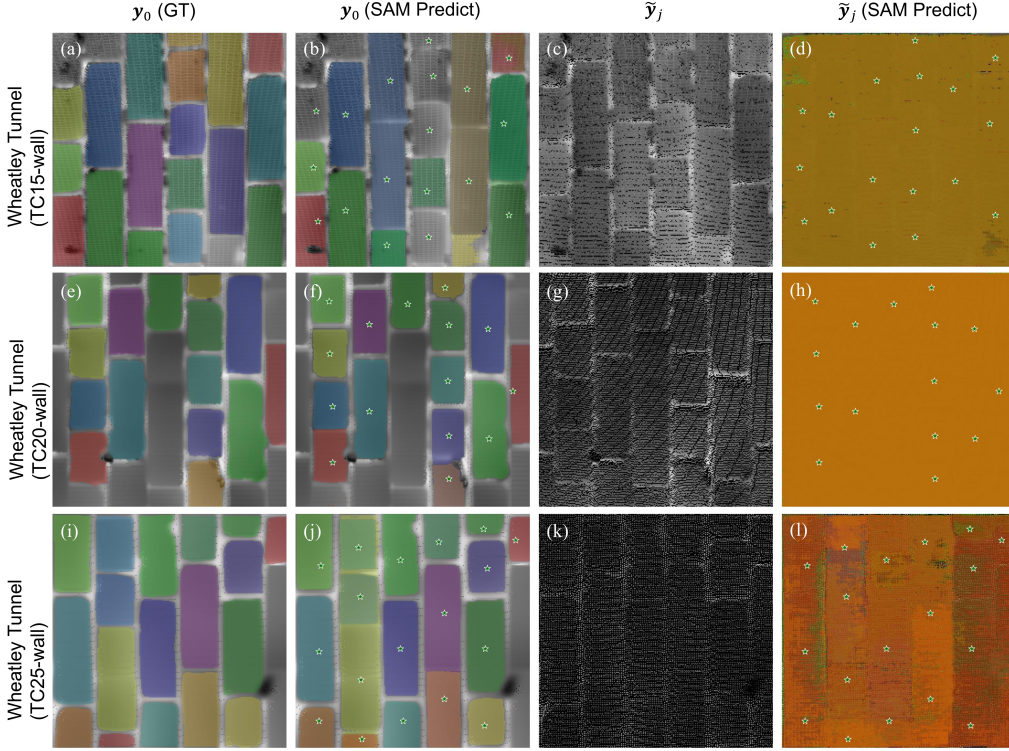
Figure 12: Zero-shot segmentation of wall patches from Wheatley Tunnel using SAM with ViT-H. (a)–(d) represent GT ($\mathbf{y}_0$), SAM predictions on $\mathbf{y}_0$, $\tilde{\mathbf{y}}_j$, and SAM predictions on $\tilde{\mathbf{y}}_j$ of test chunk 15. (e)–(h) and (i)–(l) correspond to test chunks 20 and 25. Green pentagrams indicate positive prompts.

outlines are omitted for rigorous performance. The labelling is performed on the sampled patches introduced in the Dataset section.

We adopt prompt-based SAM rather than training/fine-tuning DL-based segmentation methods such as YOLO (Khanam and Hussain, 2024), as the latter requires large labelled datasets that are often unavailable for masonry structures. In contrast, SAM, which supports zero-shot operation, can better reflect realistic annotation scenarios, where practitioners rely on AI-assisted tools by providing prompts (e.g., X-AnyLabeling (Wang, 2023)) to speed up labelling while retaining user control.

The experiments use the same patch indices listed in Table 4 of the depth map restoration. To simulate the interaction between humans and AI-assisted annotation tools, prompts are placed at representative brick regions.

Table 5: mIoU values of SAM segmentation on $\tilde{\mathbf{y}}_j$ and $\mathbf{y}_0$ using different prompting strategies.

| | Wheatley Tunnel (TC15) | Wheatley Tunnel (TC20) | Wheatley Tunnel (TC20) | Anonymised UK bridge | Hertford viaduct |
|---|---|---|---|---|---|
| $\tilde{\mathbf{y}}_j$ (1 pos) | 0.064 | 0.216 | 0.181 | 0.047 | 0.064 |
| $\tilde{\mathbf{y}}_j$ (1 pos + 1 neg) | / | / | / | 0.049 | 0.074 |
| $\mathbf{y}_0$ (1 pos) | 0.708 | 0.875 | 0.780 | 0.436 | 0.708 |
| $\mathbf{y}_0$ (1 pos + 1 neg) | / | / | / | 0.460 | 0.729 |

A positive prompt is generated by averaging the centroids of GT polygons, while the negative prompt is a randomly sampled point outside the annotated regions. Negative prompts are used only for the masonry bridges to mitigate the influence of lower data quality, but not for the tunnel experiments (only positive prompts are used). Since random negative prompts can scatter across the image plane and introduce confusion, they are omitted from the visualisations but retained in the quantitative analysis.

Figure 11a–l and Figure 12a–l illustrate zero-shot segmentation results on the arch and wall components of Wheatley Tunnel (TC15, TC20, and TC25) using SAM (ViT-H). As shown in Figure 11a, e, and i and Figure 12a, e, and i, the transparent coloured overlays denote the GT masks of masonry bricks. The corresponding SAM predictions on $\mathbf{y}_0$ with positive prompts only are shown in Figure 11b, f, and j and Figure 12b, f, and j. For visual clarity, the colours of the predicted masks are matched to the GT labels. Only brick instances with an Intersection-over-Union (IoU) score greater than 0.3 are displayed to prevent overlapping of incorrect masks. The IoU of a single masonry brick is defined as:

$$\text{IoU} = \frac{|\mathcal{M}_{\text{GT}} \cap \mathcal{M}_{\text{Pred}}|}{|\mathcal{M}_{\text{GT}} \cup \mathcal{M}_{\text{Pred}}|} \quad (23)$$

where $\mathcal{M}_{\text{GT}}$ and $\mathcal{M}_{\text{Pred}}$ denote the GT and predicted masks of brick instances, respectively.

In contrast, Figure 11d, h, and l and Figure 12d, h, and l present SAM predictions on $\tilde{\mathbf{y}}_j$, visualised using the same positive prompts. To make the comparison consistent, we show the same set of masks selected for $\mathbf{y}_0$, but without IoU-based filtering. The predictions on $\tilde{\mathbf{y}}_j$ collapse into large and overlapping masks that fail to delineate brick boundaries.

To quantify segmentation performance, we compute the mean Intersection-over-Union (mIoU) by first evaluating the IoU for each annotated brick in-
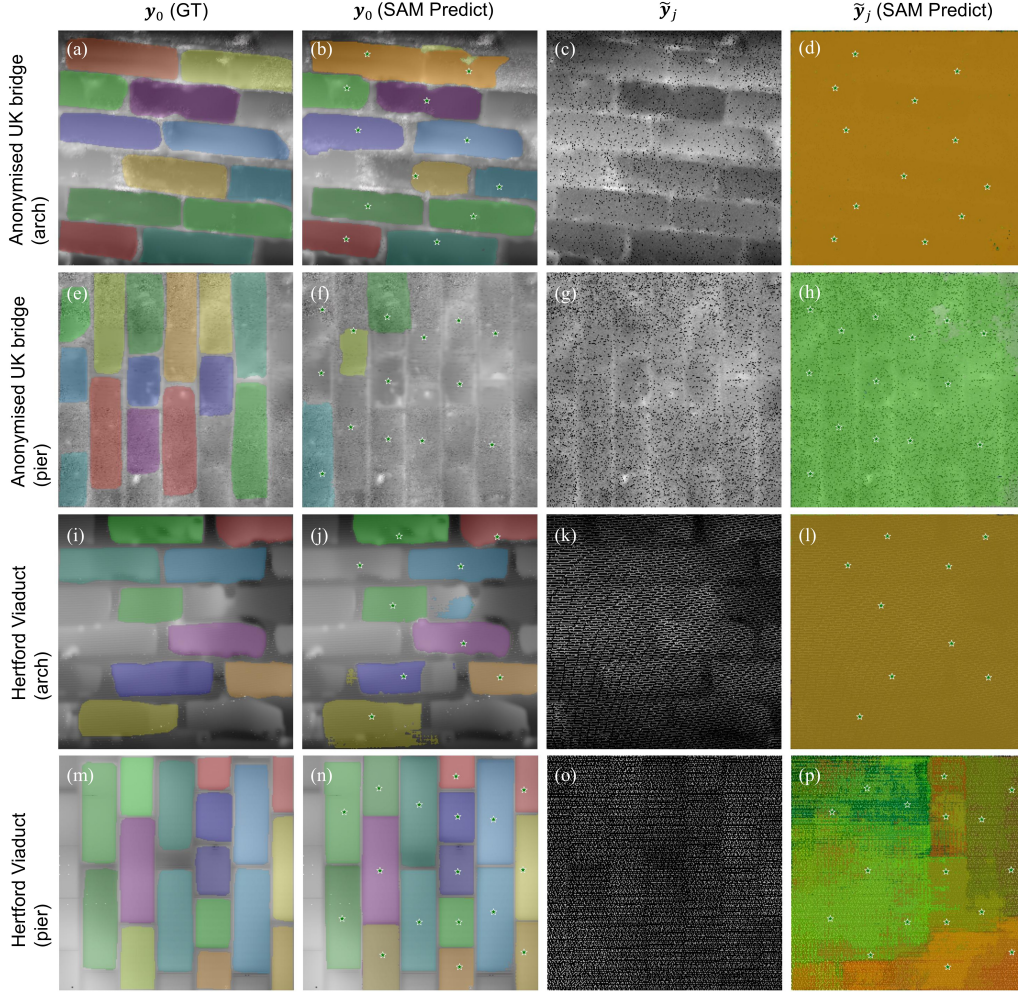
Figure 13: Zero-shot segmentation of masonry bridge patches using SAM with ViT-H by providing both positive and negative prompts. (a)–(d) show results for the arch of the anonymised UK bridge, (e)–(h) for the pier of the anonymised UK bridge, (i)–(l) for the arch of the Hertford Viaduct, and (m)–(p) for the pier of the Hertford Viaduct.

stance, then averaging across all bricks within an infrastructure, and finally averaging across all patches in the dataset. As summarised in Table 5, the predictions on $\mathbf{y}_0$ achieve substantially higher scores than those on $\tilde{\mathbf{y}}_j$. For TC15, TC20, and TC25 of Wheatley Tunnel, the mIoU improves from 0.064, 0.216, and 0.181 on $\tilde{\mathbf{y}}_j$ to 0.708, 0.875, and 0.780 on $\mathbf{y}_0$, respectively, which demonstrates the effectiveness of InfraDiffusion in enabling accurate brick-

level segmentation.

For the bridge datasets, we evaluate two prompting strategies. The first uses a single positive prompt, consistent with the tunnel experiments, while the second combines a positive–negative prompt pair to increase robustness as the masonry datasets are noisier. As shown in Table 5, using a single positive prompt yields mIoU scores of 0.436 for the anonymised UK bridge and 0.708 for the Hertford viaduct. Adding a negative prompt raises the mIoU to 0.460 and 0.729, respectively.

The anonymised UK bridge presents a more challenging case due to the lower quality of $\tilde{\mathbf{y}}_j$ (Figure 13c and g), in that the boundaries are more blurred compared to the Wheatley Tunnel shown in Figure 11e and g). Therefore, SAM predictions fail to capture meaningful masonry boundaries and degenerate into large spreading masks (Figure 13d and h). By contrast, InfraDiffusion substantially improves segmentation on $\mathbf{y}_0$ (Figure 13b and f), where SAM (guided by positive–negative prompts) recovers coherent brick-level boundaries. The improvement is reflected in the quantitative results, where the mIoU increases from 0.049 on $\tilde{\mathbf{y}}_j$ to 0.460 on $\mathbf{y}_0$ (Table 5).

For the Hertford Viaduct, both arch and pier patches benefit from InfraDiffusion, as shown in Figure 13i–l and m–p, respectively. However, the improvement is more evident for the pier regions. $\mathbf{y}_0$ of the pier yields sharper and more continuous brick boundaries, enabling SAM to recover detailed masks, whilst arch predictions are still influenced by blurred boundaries even after InfraDiffusion, as shown in Figure 13i and j.

Such a discrepancy between arch and pier performance is primarily attributable to the relative distance between the TLS and the structural component. As shown in Figure 6, the height of the Hertford Viaduct arch is 12.6 m, which results in noisier $\tilde{\mathbf{y}}_j$. In contrast, the piers are located at a lateral distance of approximately half the bridge width (4.95 m) from the scanner when positioned at the centreline, leading to denser and cleaner projections. This trend is not observed in the anonymised UK bridge, where a smaller height of 3.7 m ensures that both arches and piers remain within a similar scanning range, which results in more uniform (albeit consistently lower-quality) inputs.

As shown in Table 5, the mIoU on the Hertford Viaduct improves from 0.064 on $\tilde{\mathbf{y}}_j$ to 0.708 on $\mathbf{y}_0$ when using a single positive prompt. With the positive–negative prompting strategy, the score further rises to 0.729 (Table 5). These results confirm the effectiveness of InfraDiffusion in restoring depth map quality and demonstrate its effectiveness and robustness in supporting

28

downstream zero-shot segmentation tasks.

## Conclusions

This work introduced **InfraDiffusion**, a zero-shot diffusion framework that restored sparse and noisy depth maps derived from masonry point clouds, enabling zero-shot brick-level image segmentation for structural assessment. The pipeline first applied a virtual camera projection to systematically generate depth maps from large-scale masonry infrastructure. It then adapted the vanilla DDNM framework with boundary masking to restore these maps, addressing the boundary effects of point cloud projection in civil engineering applications. By constraining the restoration to valid regions only, InfraDiffusion prevented spurious content outside structural edges. InfraDiffusion relied solely on pre-trained DDPMs to provide strong generative priors, operating entirely in a zero-shot setting without requiring task-specific training or ground-truth supervision.

Extensive experiments were conducted on three representative chunks of the Wheatley Tunnel (TC15, TC20, and TC25 with 100 patches) and two masonry bridges (an anonymised UK bridge with 60 patches and the Hertford Viaduct with 140 patches). InfraDiffusion consistently produced clean and geometrically coherent depth maps across these datasets. The ablation study further demonstrated the effectiveness of InfraDiffusion conditioned on extra boundary masks. While vanilla DDNM produces the same results in foreground regions, it also generates spurious bricks outside structural edges when applied to projected depth maps. InfraDiffusion suppresses such artefacts by constraining restoration to masked regions, yielding more physically consistent outputs for infrastructure point clouds.

Extensive experiments were conducted on three representative chunks of the Wheatley Tunnel (TC15, TC20, and TC25 with 100 patches) and two masonry bridges (an anonymised UK bridge with 60 patches and the Hertford Viaduct with 140 patches). InfraDiffusion consistently produced clean and geometrically coherent depth maps across these datasets. The ablation study further demonstrated the effectiveness of conditioning DDNM with boundary masks. While vanilla DDNM produced comparable results in foreground regions, it also generated spurious bricks outside structural edges when applied to projected depth maps. InfraDiffusion suppressed such artefacts by constraining restoration to masked regions, yielding more physically consistent outputs for infrastructure point clouds.

Quantitative validation was carried out through zero-shot semantic segmentation using SAM. With only a single positive prompt inside bricks, InfraDiffusion significantly improved segmentation performance: across all datasets, mIoU scores rose from as low as 0.064 on Wheatley Tunnel TC15 to 0.708 after restoration. Additional robustness was achieved by introducing a negative prompt for the two noisier bridge datasets, where mIoU further improved from 0.436 to 0.460 on the anonymised UK bridge and from 0.708 to 0.729 on the Hertford Viaduct. These results showed that InfraDiffusion effectively transformed noisy depth maps into representations that downstream segmentation models could exploit.

Our results demonstrated the robustness and effectiveness of InfraDiffusion across both tunnels and bridges. They highlighted the potential of diffusion models for civil engineering: despite being trained solely on natural images, pre-trained diffusion models transferred effectively to infrastructure depth maps because their strong generative priors captured local pixel dependencies that were equally relevant in masonry. This pointed to promising avenues for applying diffusion models to other structural assessment tasks where sparse and noisy sensing data must be transformed into complete and analysable representations.

**Limitations and future work** While InfraDiffusion demonstrates robustness and effectiveness, several limitations remain. First, although DDIM accelerates sampling compared to DDPM, the restoration process is still computationally intensive when applied to large numbers of structural patches, limiting scalability for large-scale masonry infrastructures.

Second, the current framework is not end-to-end. The IR of InfraDiffusion and brick-level segmentation are performed sequentially, whereas future work should aim to integrate these steps into a unified pipeline for improved efficiency and consistency.

Third, our method restores only geometric depth information, without considering generating synthetic RGB colour, which is often unavailable in LiDAR scans taken from low-light conditions. Incorporating colour restoration could further enhance the point cloud quality and support downstream tasks such as defect detection.

## Data availability statement

The data and the code are now available at `https://github.com/Jingyixiong/InfraDiffusion-official-implement`.

## References

Acikgoz, S., DeJong, M.J., Soga, K., 2018. Sensing dynamic displacements in masonry rail bridges using 2d digital image correlation. Structural Control and Health Monitoring 25, e2187. URL: `https://doi.org/10.1002/stc.2187`.

Acikgoz, S., Soga, K., Woodhams, J., 2017. Evaluation of the response of a vaulted masonry structure to differential settlements using point cloud data and limit analyses. Construction and Building Materials 150, 916–931. URL: `https://doi.org/10.1016/j.conbuildmat.2017.05.075`.

Andrews, H.C., 1974. Digital image restoration: A survey. Computer 7, 36–45. URL: `10.1109/MC.1974.6323527`.

Brackenbury, D., 2022. Automated Image-Based Inspection of Masonry Arch Bridges. Ph.D. thesis. Apollo - University of Cambridge Repository. URL: `https://www.repository.cam.ac.uk/handle/1810/338321`, doi:`10.17863/CAM.85731`.

Cao, C., Dong, Q., Fu, Y., 2022. Learning prior feature and attention enhanced image inpainting, in: European conference on computer vision, Springer. pp. 306–322. URL: `10.1007/978-3-031-19784-0_18`.

Chen, C., Xiong, Z., Tian, X., Zha, Z.J., Wu, F., 2019. Real-world image denoising with deep boosting. IEEE Transactions on Pattern Analysis and Machine Intelligence 42, 3071–3087. URL: `10.1109/TPAMI.2019.2921548`.

Chen, Q., Kang, Z., Cao, Z., Xie, X., Guan, B., Pan, Y., Chang, J., 2024. Combining cylindrical voxel and mask r-cnn for automatic detection of water leakages in shield tunnel point clouds. Remote Sensing 16, 896. URL: `https://doi.org/10.3390/rs16050896`.

Chen, X., Ma, H., Wan, J., Li, B., Xia, T., 2017. Multi-view 3d object detection network for autonomous driving, in: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, pp. 1907–1915. URL: `10.1109/CVPR.2017.691`.

Cheng, Y.J., Qiu, W.G., Duan, D.Y., 2019. Automatic creation of as-is building information model from single-track railway tunnel point clouds. Automation in Construction 106, 102911.

Dai, F., Lu, M., 2013. Three-dimensional modeling of site elements by analytically processing image data contained in site photos. Journal of construction engineering and management 139, 881–894. URL: `https://doi.org/10.1061/(ASCE)CO.1943-7862.0000655`.

Dais, D., Bal, I.E., Smyrou, E., Sarhosis, V., 2021. Automatic crack classification and segmentation on masonry surfaces using convolutional neural networks and transfer learning. Automation in Construction 125, 103606. URL: `https://doi.org/10.1016/j.autcon.2021.103606`.

Dhariwal, P., Nichol, A., 2021. Diffusion models beat gans on image synthesis. Advances in neural information processing systems 34, 8780–8794.

Dong, C., Loy, C.C., He, K., Tang, X., 2014. Learning a deep convolutional network for image super-resolution, in: European conference on computer vision, Springer. pp. 184–199. URL: `10.1007/978-3-319-10593-2_13`.

Dong, C., Loy, C.C., He, K., Tang, X., 2015. Image super-resolution using deep convolutional networks. IEEE transactions on pattern analysis and machine intelligence 38, 295–307. URL: `10.1109/TPAMI.2015.2439281`.

Fritsche, M., Gu, S., Timofte, R., 2019. Frequency separation for real-world super-resolution, in: 2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW), IEEE. pp. 3599–3608. URL: `10.1109/TCSVT.2024.3367876`.

Giordano, A., Mele, E., De Luca, A., 2002. Modelling of historical masonry structures: comparison of different approaches through a case study. Engineering Structures 24, 1057–1069. URL: `https://doi.org/10.1016/S0141-0296(02)00033-0`.

Golub, G.H., Reinsch, C., 1971. Singular value decomposition and least squares solutions, in: Linear algebra. Springer, pp. 134–151. URL: `https://link.springer.com/content/pdf/10.1007/978-3-662-39778-7_10.pdf`.

Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., et al., 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. arXiv preprint arXiv:2501.12948 URL: `https://doi.org/10.48550/arXiv.2501.12948`.

Hallee, M.J., Napolitano, R.K., Reinhart, W.F., Glisic, B., 2021. Crack detection in images of masonry using cnns. Sensors 21, 4929. URL: `https://doi.org/10.3390/s21144929`.

Han, J., Jang, M., Han, H., Shin, D.H., 2025. Automatic extraction of bridge dimensional information using 3d point cloud data. KSCE Journal of Civil Engineering , 100312URL: `https://doi.org/10.1016/j.kscej.2025.100312`.

Ho, J., Jain, A., Abbeel, P., 2020. Denoising diffusion probabilistic models. Advances in neural information processing systems 33, 6840–6851. URL: `10.5555/3495724.3496298`.

Jing, Y., Sheil, B., Acikgoz, S., 2022. Segmentation of large-scale masonry arch bridge point clouds with a synthetic simulator and the bridgenet neural network. Automation in Construction 142, 104459. URL: `https://doi.org/10.1016/j.autcon.2022.104459`.

Jing, Y., Sheil, B., Acikgoz, S., 2023. A method to generate realistic synthetic point clouds of damaged single-span masonry arch bridges, in: International Conference on Structural Analysis of Historical Constructions, Springer. pp. 436–448. URL: `https://doi.org/10.1007/978-3-031-39603-8_36`.

Jing, Y., Sheil, B., Acikgoz, S., 2024a. A lightweight transformer-based neural network for large-scale masonry arch bridge point cloud segmentation. Computer-Aided Civil and Infrastructure Engineering URL: `https://doi.org/10.1111/mice.13201`.

Jing, Y., Zhong, J.X., Sheil, B., Acikgoz, S., 2024b. Anomaly detection of cracks in synthetic masonry arch bridge point clouds using fast point feature histograms and patchcore. Automation in Construction 168, 105766. URL: `https://doi.org/10.1016/j.autcon.2024.105766`.

Katsigiannis, S., Seyedzadeh, S., Agapiou, A., Ramzan, N., 2023. Deep learning for crack detection on masonry façades using limited data and transfer learning. Journal of Building Engineering 76, 107105. URL: `https://doi.org/10.1016/j.jobe.2023.107105`.

Khanam, R., Hussain, M., 2024. Yolov11: An overview of the key architectural enhancements. arXiv preprint arXiv:2410.17725 URL: https://arxiv.org/abs/2410.17725.

Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al., 2023. Segment anything, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 4015–4026. URL: https://doi.org/10.48550/arXiv.2304.02643.

Krawciw, A., Lilge, S., Barfoot, T.D., 2024. Lasersam: Zero-shot change detection using visual segmentation of spinning lidar. arXiv preprint arXiv:2402.10321 URL: https://doi.org/10.48550/arXiv.2402.10321.

Lemos, J.V., 2007. Discrete element modeling of masonry structures. International Journal of Architectural Heritage 1, 190–213. URL: https://doi.org/10.1080/15583050601176868.

Lin, X., He, J., Chen, Z., Lyu, Z., Dai, B., Yu, F., Qiao, Y., Ouyang, W., Dong, C., 2024. Diffbir: Toward blind image restoration with generative diffusion prior, in: European conference on computer vision, Springer. pp. 430–448. URL: 10.1007/978-3-031-73202-7_25.

Lo, K.S.H., Peters, J., Spellman, E., 2024. Roofdiffusion: Constructing roofs from severely corrupted point data via diffusion, in: European Conference on Computer Vision, Springer. pp. 38–57. URL: https://doi.org/10.48550/arXiv.2404.09290.

Lourenço, P.B., 2013. Computational strategies for masonry structures : multi-scale modeling, dynamics, engineering applications and other challenges. URL: https://api.semanticscholar.org/CorpusID:111285063.

Loverdos, D., Sarhosis, V., 2022. Automatic image-based brick segmentation and crack detection of masonry walls using machine learning. Automation in Construction 140, 104389. URL: https://doi.org/10.1016/j.autcon.2022.104389.

Lubowiecka, I., Armesto, J., Arias, P., Lorenzo, H., 2009. Historic bridge modelling using laser scanning, ground penetrating radar and finite ele-

ment methods in the context of structural dynamics. Engineering Structures 31, 2667–2676. URL: `https://doi.org/10.1016/j.engstruct.2009.06.018`.

Luo, J., Ye, Q., Zhang, S., Yang, Z., 2023. Indoor mapping using low-cost mls point clouds and architectural skeleton constraints. Automation in Construction 150, 104837. URL: `https://doi.org/10.1016/j.autcon.2023.104837`.

Lyu, Y., Huang, X., Zhang, Z., 2020. Learning to segment 3d point clouds in 2d image space, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12255–12264. URL: `10.1109/CVPR42600.2020.01227`.

Massa, K.J., Grobler, H., 2024. Adapting projection-based lidar semantic segmentation to natural domains. Journal of Visual Communication and Image Representation 100, 104111. URL: `https://doi.org/10.1016/j.jvcir.2024.104111`.

Meyer, G.P., Laddha, A., Kee, E., Vallespi-Gonzalez, C., Wellington, C.K., 2019. Lasernet: An efficient probabilistic 3d object detector for autonomous driving, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 12677–12686. URL: `10.1109/CVPR.2019.01296`.

Milani, G., Lourenço, P.B., Tralli, A., 2006. Homogenised limit analysis of masonry walls, part i: Failure surfaces. Computers & structures 84, 166–180. URL: `https://doi.org/10.1016/j.compstruc.2005.09.005`.

Orbán, Z., 2004. Assessment, reliability and maintenance of masonry arch railway bridges in europe. Arch Bridges IV–Advances in Assessment, Structural Design and Construction. Eds: P. Roca and C. Molins, Barcelona 2004, 152–161.

Qu, T., Coco, J., Rönnäng, M., Sun, W., 2014. Challenges and Trends of Implementation of 3D Point Cloud Technologies in Building Information Modeling (BIM): Case Studies. pp. 809–816. URL: `https://ascelibrary.org/doi/abs/10.1061/9780784413616.101`, `arXiv:https://ascelibrary.org/doi/pdf/10.1061/9780784413616.101`.

Richardson, W.H., 1972. Bayesian-based iterative method of image restoration. Journal of the optical society of America 62, 55–59. URL: `https://doi.org/10.1364/JOSA.62.000055`.

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B., 2022. High-resolution image synthesis with latent diffusion models, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 10684–10695. URL: `10.1109/CVPR52688.2022.01042`.

Saharia, C., Ho, J., Chan, W., Salimans, T., Fleet, D.J., Norouzi, M., 2022. Image super-resolution via iterative refinement. IEEE transactions on pattern analysis and machine intelligence 45, 4713–4726. URL: `10.1109/TPAMI.2022.3204461`.

Shanoer, M.M., Abed, F.M., 2018. Evaluate 3d laser point clouds registration for cultural heritage documentation. The Egyptian Journal of Remote Sensing and Space Science 21, 295–304. URL: `https://doi.org/10.1016/j.ejrs.2017.11.007`.

Song, J., Meng, C., Ermon, S., 2020. Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502 URL: `https://doi.org/10.48550/arXiv.2010.02502`.

Wang, C., Cho, Y.K., 2015. Smart scanning and near real-time 3d surface modeling of dynamic construction equipment from a point cloud. Automation in Construction 49, 239–249. URL: `https://doi.org/10.1016/j.autcon.2014.06.003`.

Wang, L., Wang, Y., Dong, X., Xu, Q., Yang, J., An, W., Guo, Y., 2021a. Unsupervised degradation representation learning for blind super-resolution, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 10581–10590. URL: `10.1109/CVPR46437.2021.01044`.

Wang, W., 2023. Advanced auto labeling solution with added features. URL: `https://github.com/CVHub520/X-AnyLabeling`.

Wang, X., Xie, L., Dong, C., Shan, Y., 2021b. Real-esrgan: Training real-world blind super-resolution with pure synthetic data, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 1905–1914. URL: `10.1109/ICCVW54120.2021.00217`.

Wang, Y., Yu, J., Zhang, J., 2022. Zero-shot image restoration using denoising diffusion null-space model. arXiv preprint arXiv:2212.00490 URL: https://doi.org/10.48550/arXiv.2212.00490.

Wang, Y., Yu, Y., Yang, W., Guo, L., Chau, L.P., Kot, A.C., Wen, B., 2023. Exposurediffusion: Learning to expose for low-light image enhancement, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 12438–12448. URL: https://doi.org/10.48550/arXiv.2307.07710.

Wu, B., Zhou, X., Zhao, S., Yue, X., Keutzer, K., 2019a. Squeezesegv2: Improved model structure and unsupervised domain adaptation for road-object segmentation from a lidar point cloud, in: 2019 international conference on robotics and automation (ICRA), IEEE. pp. 4376–4382.

Wu, H., Ao, X., Chen, Z., Liu, C., Xu, Z., Yu, P., 2019b. Concrete spalling detection for metro tunnel from point cloud based on roughness descriptor. Journal of Sensors 2019, 8574750. URL: https://doi.org/10.1155/2019/8574750.

Ye, C., Acikgoz, S., Pendrigh, S., Riley, E., DeJong, M., 2018. Mapping deformations and inferring movements of masonry arch bridges using point cloud data. Engineering Structures 173, 530–545. URL: https://doi.org/10.1016/j.engstruct.2018.06.094.

Ye, Z., Lin, W., Faramarzi, A., Xie, X., Ninić, J., 2025. Sam4tun: No-training model for tunnel lining point cloud component segmentation. Tunnelling and Underground Space Technology 158, 106401. URL: https://doi.org/10.1016/j.tust.2025.106401.

Ye, Z., Lovell, L., Faramarzi, A., Ninić, J., 2024. Sam-based instance segmentation models for the automation of structural damage detection. Advanced Engineering Informatics 62, 102826. URL: https://doi.org/10.1016/j.aei.2024.102826.

You, K., Zhou, C., Ding, L., Wang, Y., 2025. Construction robotics in extreme environments: From earth to space. Engineering URL: https://doi.org/10.1016/j.eng.2024.11.037.

Zhang, Z., Ji, A., Wang, K., Zhang, L., 2022. Unrollingnet: An attention-based deep learning approach for the segmentation of large-scale point

clouds of tunnels. Automation in Construction 142, 104456. URL: `https://doi.org/10.1016/j.autcon.2022.104456`.

Zhou, Q.Y., Park, J., Koltun, V., 2018. Open3d: A modern library for 3d data processing. arXiv preprint arXiv:1801.09847 URL: `https://doi.org/10.48550/arXiv.1801.09847`.

Zhu, Y., Zhao, W., Li, A., Tang, Y., Zhou, J., Lu, J., 2024. Flowie: Efficient image enhancement via rectified flow, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 13–22. URL: `https://doi.org/10.48550/arXiv.2406.00508`.