# Geometry-Aware Attenuation Learning for Sparse-View CBCT Reconstruction

Zhentao Liu, Yu Fang, Changjian Li, Han Wu, Yuan Liu, Dinggang Shen, *Fellow, IEEE*, Zhiming Cui

*Abstract*—**Cone Beam Computed Tomography (CBCT) plays a vital role in clinical imaging. Traditional methods typically require hundreds of 2D X-ray projections to reconstruct a high-quality 3D CBCT image, leading to considerable radiation exposure. This has led to a growing interest in sparse-view CBCT reconstruction to reduce radiation doses. While recent advances, including deep learning and neural rendering algorithms, have made strides in this area, these methods either produce unsatisfactory results or suffer from time inefficiency of individual optimization. In this paper, we introduce a novel geometry-aware encoder-decoder framework to solve this problem. Our framework starts by encoding multi-view 2D features from various 2D X-ray projections with a 2D CNN encoder. Leveraging the geometry of CBCT scanning, it then back-projects the multi-view 2D features into the 3D space to formulate a comprehensive volumetric feature map, followed by a 3D CNN decoder to recover 3D CBCT image. Importantly, our approach respects the geometric relationship between 3D CBCT image and its 2D X-ray projections during feature back projection stage, and enjoys the prior knowledge learned from the data population. This ensures its adaptability in dealing with extremely sparse view inputs without individual training, such as scenarios with only 5 or 10 X-ray projections. Extensive evaluations on two simulated datasets and one real-world dataset demonstrate exceptional reconstruction quality and time efficiency of our method.**

*Index Terms*—**Sparse-view CBCT reconstruction, geometry awareness, multi-view consistence, prior knowledge.**

## I. INTRODUCTION

CONE Beam Computed Tomography (CBCT), a specialized form of CT scanning, is extensively utilized in clinical settings for diagnostic purposes, including dental [1], [2], spinal [3], [4], and vascular disease diagnosis [5], [6]. Compared to the traditional Fan Beam CT (FBCT) [2], CBCT

Zhentao Liu, Yu Fang, Han Wu, Dinggang Shen, and Zhiming Cui are with the School of Biomedical Engineering & State Key Laboratory of Advanced Medical Materials and Devices, ShanghaiTech Univerisity, Shanghai, 201210, China. Dinggang Shen is also with Shanghai United Imaging Intelligence Co., Ltd., Shanghai, 200230, China, and Shanghai Clinical Research and Trial Center, Shanghai, 201210, China. (e-mail: {liuzht2022, fangyu, wuhan2022, dgshen, cuizhm}@shanghaitech.edu.cn). (*Corresponding author: Zhiming Cui.*)

Changjian Li is with the School of Informatics, The University of Edinburgh, Edinburgh, UK (e-mail: chjili2011@gmail.com).

Yuan Liu is with the Department of Computer Science, The University of Hong Kong, Hong Kong, China (e-mail: yuanly@connect.hku.hk).

offers high-resolution images in a shorter scanning time. A typical CBCT imaging system is illustrated in Fig. 1. During CBCT scanning, an X-ray source moves uniformly along an arc-shaped trajectory, emitting a cone-shaped beam at each angular step towards the human body, such as the oral cavity. A detector positioned on the patient's opposite side records the 2D projections. Therefore, the CBCT reconstruction essentially transforms into an inverse problem, with the goal of recovering 3D anatomical information from these 2D X-ray projections. However, traditional methods usually require hundreds of projections to produce a high-quality CBCT image [8], raising concerns about radiation exposure. Hence, sparse-view CBCT reconstruction, which reduces the number of projection views to mitigate radiation exposure, has received widespread attention in the research field.

Traditionally, CBCT reconstruction has mainly employed analytical methods like the Filtered Back Projection (FBP) algorithm [7], [8]. While these methods provide rapid reconstruction solutions, they rely heavily on numerous projection views. To address these limitations, several iterative optimization-based methods [9]–[12] have been proposed. However, despite producing improved reconstruction results under sparse input, they are not time efficient and often lack fine details. Recently, with the advent of deep learning, some researchers have explored deep learning techniques to learn the mapping between multi-view projections and CBCT image in an end-to-end manner from extensive datasets [13]–[15]. But these approaches directly concatenate multi-view image information, neglecting the inherent geometric properties of the CBCT system, leading to structurally erroneous reconstruction results. Concurrently, in the 3D vision society, the technique of neural rendering (e.g., NeRF [16]) has garnered significant interest for novel view synthesis and multi-view reconstruction. These techniques represent a 3D scene by the neural representation, and employ differentiable rendering for optimization. It is also a typical inverse problem that tries to recover the 3D scene from multi-view image captures, similar to the concept of CBCT reconstruction. Our goal is to reconstruct 3D anatomical information from multi-view X-ray projections. Building on this, many methods, e.g., IntroTomo [17], NeAT [18], NAF [19], SNAF [20], SAX-NeRF [21], have been developed to incorporate neural rendering into CBCT reconstruction, and achieved impressive results. However, a significant drawback of these methods is that they must be individually optimized for each CBCT scan, making the reconstruction process extremely time-consuming,

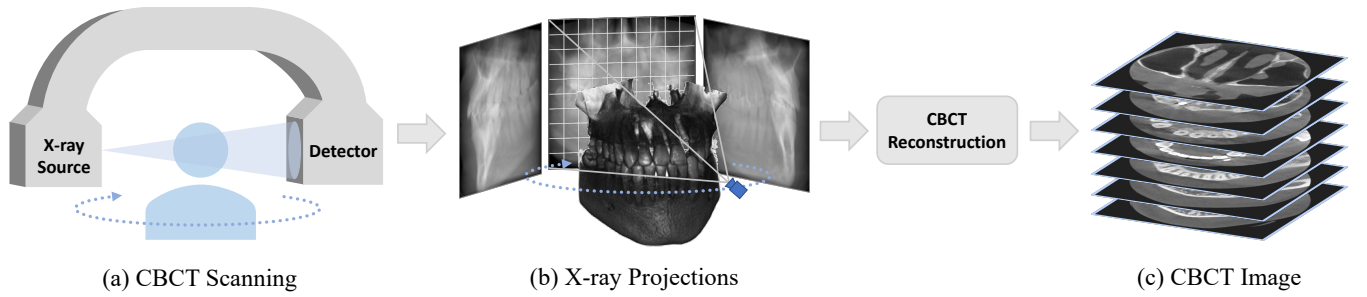(a) CBCT Scanning      (b) X-ray Projections      (c) CBCT Image

Fig. 1. CBCT scanning and reconstruction. In the CBCT imaging process, CBCT scanning (a) would generate a sequence of 2D X-ray projections (b). These projections are utilized to reconstruct 3D CBCT image (c).

i.e., using tens of minutes for one patient. More importantly, the quality of reconstructed images diminishes considerably with using the extremely sparse inputs, such as 5 or 10 views.

In this way, sparse-view CBCT reconstruction is a highly ill-posed problem with twofold challenges: (1) How to bridge the dimension gap between multi-view 2D X-ray projections and the CBCT image; (2) How to solve information insufficiency introduced by extremely sparse-view input. In this study, we introduce a geometry-aware encoder-decoder framework to solve this task efficiently. It seamlessly integrates the multi-view consistency of neural rendering and the generalization ability of deep learning, effectively addressing the challenges mentioned above. Specifically, we first adopt a 2D convolutional neural network (CNN) encoder to extract multi-view 2D features from different X-ray projections. Then, in aligning with the geometry of CBCT scanning, we back-project multi-view 2D features into 3D space, which properly bridges the dimension gap with multi-view consistency. Particularly, as different views offer varying degrees of information, an adaptive feature fusion strategy is further introduced to aggregate these multi-view features. Consequently, a 3D volumetric feature is constructed and then decoded into 3D CBCT image with a 3D CNN decoder. Our framework's inherent geometry awareness ensures accurate information retrieval from multi-view X-ray projections. Moreover, by capturing prior knowledge from populations in extensive datasets, our method can generalize well across different patients without individual optimization, even with extremely sparse input views, such as 5 or 10 views. Experiments on two simulated datasets and one real-world dataset underscore the effectiveness and time efficiency of our method.

## II. RELATED WORKS

Sparse-view CBCT reconstruction aims to reduce the number of projections while preserving the quality of the reconstructed CBCT image. Existing methods predominantly fall into three categories: traditional methods, learning-based methods, and neural rendering-based methods, as briefed below.

### A. Traditional CBCT Reconstruction

The Filtered Back Projection (FBP) [7] and its 3D cone-beam variant FDK (Feldkamp, Davis, and Kress) [8] stand out as the most prevalent CBCT reconstruction algorithm in the commercial CBCT systems. FDK back-projects the filtered

projection values from different viewpoints into the original 3D space to reconstruct the CBCT image. But, to avoid streaks in the CBCT image, it needs hundreds of X-ray projections. This can be harmful because of the high radiation. To solve this, several iterative methods [9]–[12] have been proposed to leverage the consistency between acquired projections and the reconstructed image. For example, SART [9] introduces an iterative approach to diminish the disparity between 2D projections and their estimates, updating the reconstructed image at the same time. Other methods [10]–[12] further incorporate regularization terms to improve reconstruction quality. However, while these methods work well in noisy and view-limited scenarios, they usually lose image details and demand more computational resources.

### B. Learning-Based CBCT Reconstruction

Deep learning has been increasingly applied to CT/CBCT reconstruction, and could be roughly divided into three classes: data restoration reconstruction [22]–[24], model-based iterative reconstruction [25]–[27], and end-to-end reconstruction [13]–[15]. First, data restoration reconstruction methods mainly include image domain restoration [22], sinogram domain restoration [23], and dual-domain restoration solutions [24]. FBPConvNet [22] utilizes CNN for image domain restoration to improve upon an initial reconstruction from FBP algorithm. SinoConvNet [23] operates in sinogram domain, employing a CNN to inpaint sparse view sinogram to a full view one before performing FBP reconstruction. DRONE [24] combines image and sinogram restoration together as a dual-domain approach. While effective for 2D CT slices, extending these methods to 3D CBCT reconstruction can lead to significant GPU memory overhead. Additionally, processing 3D CBCT image slice by slice leads to inter-slice inconsistency. Second, some model-based iterative methods [25], [26] have been developed based on the learned iterative reconstruction technique. Their basic idea is to unroll an iterative scheme with a deep neural network architecture, explicitly integrating forward projection and back projection into model architecture. These methods would produce satisfactory results but they may suffer from high-memory footprints and time inefficiency. IRON [27] is a notable iterative method for limited-angle 2D CT reconstruction. It adopts a training strategy like image domain restoration [22] and then deploys the trained sub-network as a regularization prior in the residual image domain for iterative image reconstruction. Compared to traditional

unrolled iterative methods [25], [26], IRON is more memory efficient for avoiding a large number of projection and back projection operations in training. Its underlying principles could be adapted for sparse-view 3D CBCT reconstruction. Third, there are also some end-to-end methods that directly learn the mapping from muti-view projections to CBCT image, as described below. PatRecon [13] can produce 3D CBCT reconstruction results using just a single X-ray projection. This framework primarily uses a 2D-to-3D autoencoder architecture. A 2D CNN encoder first extracts features from the X-ray projection, which are then reshaped in the feature space to transition into 3D. This is followed by a 3D CNN decoder to produce volumetric image. For multiple input views, PatRecon directly concatenates them along the channel dimension, which neglects intrinsic geometric relationships between different projections. Similarly, Bi-Recon [14] and X2CTGAN [15] operate on similar principle but are limited to handling only two orthogonal projections due to their specific model designs. They usually produce blurry reconstruction and inaccurate anatomical structures. These issues arise from their ignorance of the geometric relationship between X-ray projections and CBCT image when combing X-ray projections or their feature representations.

Recently, diffusion models have been widely employed for sparse-view or limited angle CT reconstruction in 2D [28]–[32] or 3D format [33], [34]. Their basic idea is to incorporate data consistency term to guide the diffusion model sampling process, aligning the forward projections from the generated image with the input X-ray projections. MOGM [30] establishes a multi-channel fusion module for sparse-view CT reconstruction, enhancing the efficiency of data consistency module and providing the diffusion model with more accurate guidance. DCDS [31] proposes a collaborative diffusion mechanism for sparse-view CT reconstruction that combines both sinogram and image diffusion model to simultaneously consider dual-domain prior distribution. TIFA [32] presents a novel rapid-sampling technique for limited-angle CT reconstruction that incorporates jump-sampling and time-reversion with re-sampling. This scheme not only accelerates the sampling but also enhances reconstruction outcomes. TIFA further integrates Diagonal Total Variation (DTV) to mitigate directional artifacts arising from limited-angle input. DiffusionMBIR [33] has been utilized for sparse-view 3D CT reconstruction, and DDS [34] offers an accelerated solution with decomposed sampling strategy. They generate 3D CT volume slice by slice along $z$-axis and use $z$-axis Total Variation (TV) to enforce inter-slice consistency. However, inter-slice inconsistency still exists. Besides, their performance has only been validated on parallel beam geometry, and their effectiveness on cone beam geometry requires further verification.

### C. Neural Rendering-Based CBCT Reconstruction

Neural rendering, including NeRF [16] and its derivatives [35]–[39] has rapidly progressed in the field of 3D vision. It stands out for its ability to generate new viewpoints and perform multi-view 3D reconstruction, demonstrating significant improvement in maintaining consistency across multiple views. In X-ray imaging and tomographic reconstruction, many methods have adopted neural rendering for multi-view X-ray synthesizing [40], 2D CT reconstruction [41], [42], rigid motion correction [43], metal artifact removal [44], DSA reconstruction [45], [46], and 3D CBCT reconstruction [17]–[21]. And in this study, we mainly focus on the methods designed for CBCT reconstruction. For example, Intratomo [17] employs Multilayer Perceptrons (MLPs) with Fourier feature encoding to predict novel sinograms from sparse-view inputs for tomographic reconstruction. It then uses geometrical priors to regularize the results. However, the outputs are often blurry due to the limitations of its encoding module. NeAT [18] uses an octree structure [35] to represent 3D anatomical structures. By combining this with a differentiable rendering algorithm, NeAT achieves excellent reconstruction results. Similarly, NAF [19] incorporates a multi-resolution hash table [36] for sparse-view CBCT reconstruction, showing impressive results even with as few as 50 projection views. Its efficient hash data structure helps preserve detailed features, even with limited input views. SNAF [20], an improvement on NAF, tackles more sparse view conditions (like 20 views) by introducing a new technique for view augmentation. It creates additional views between existing ones, which are then used to enhance reconstruction quality. SAX-NeRF [21] achieves plausible results with a line segment-based transformer and masked local-global ray sampling strategy. However, all these models require individual training for every subject, taking up to tens of minutes for a single CBCT image reconstruction. Moreover, they do not fully overcome the challenges posed by extremely sparse views. Their performance tends to drop in quality with extremely limited projections, such as 5 or 10 views. A notable deep learning variation of NeRF called PixelNeRF [38] was introduced, aiming to tackle the sparse-view challenge by using scene knowledge learned from large datasets. Another concurrent approach, DIF-Net [47], applies this concept to medical imaging and shows promising results on a knee CBCT dataset. However, it uses MLPs for point-wise decoding, which overlooks the interactions between neighboring points of CBCT image. And the point-wise 3D supervision makes it difficult to capture the global structure information of CBCT image. As a result, it exhibits streaky artifacts when dealing with extremely sparse input, like 5 views.

Our approach combines the advantages of deep learning and neural rendering, providing both generalization capability and multi-view consistency. Remarkably, our framework is capable of delivering reliable reconstruction quality using 20 input views in a second. This highlights superior performance and time efficiency of our method.

## III. METHOD

In this section, we first introduce data preparation in Sec. III-A. Following that, we elaborate our methodology in detail in Sec. III-B.

### A. Data Preparation

As shown in Fig. 1, the X-ray source moves around the patient along a set arc-shaped path. Opposite this source,
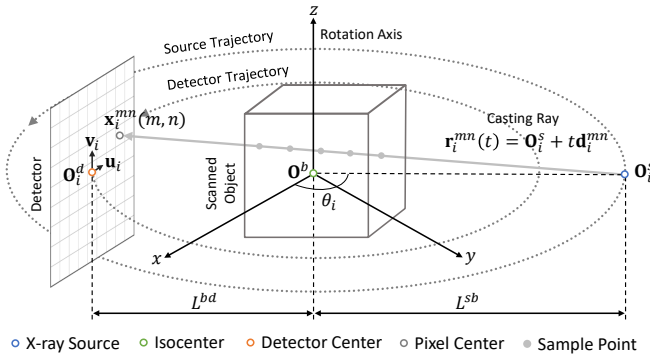
**Fig. 2.** Geometric configuration of CBCT scanning and X-ray projection simulation.

there is a 2D detector that captures X-ray projections of the body, like the oral cavity. In real-world collection, initial X-ray photons emitted from the source are attenuated differently by various tissues in the body. The detector recorded X-ray photon count would be converted into X-ray projection value after flat-field and dark-field corrections [48]. The projection value represents the line integral of attenuation along the X-ray path. For a detailed derivation of X-ray attenuation process, please refer to Sec. III-A.2. Throughout this paper, the term **X-ray projection** specifically refers to the attenuation ray integral for clarity. The attenuation value is shown as voxel intensity value we see in CBCT image. Therefore, our objective is to reconstruct the 3D CBCT image, using a limited set of X-ray projections. In our experiments, we adopt two simulated datasets of dental and spine [49] and one real-world dataset of walnut [48] to verify our method performance. In the following, we introduce how we prepare our X-ray and CBCT paired data. The detailed specifics of dataset will be provided in Sec. IV-A.1.

*1) Simulated X-ray and CBCT Pairs Collection:* For simulated datasets, we use the Digitally Reconstructed Radiography (DRR) technique to simulate multiple X-ray projections from a given CBCT image.

A typical geometric configuration of CBCT scanning is illustrated in Fig. 2. The isocenter $\mathbf{O}^b \in \mathbb{R}^3$ is the center point of the scanned object. We establish the world coordinate system with $\mathbf{O}^b$ as the origin, i.e., $\mathbf{O}^b = [0,0,0]^T$, aligning the axes with the bounding box of the scanned object. Given a volumetric CBCT image $\mathbf{V}_{gt} \in \mathbb{R}^{1 \times W \times H \times D}$ as the scanned object with a spatial resolution of $W \times H \times D$ and a voxel spacing of $v_x \times v_y \times v_z$, we define the volumetric coordinate matrix $\mathcal{V} \in \mathbb{R}^{3 \times W \times H \times D}$ that bounds the CBCT image in world space. Each element from $\mathcal{V}$ represents the center coordinate of voxel grid from CBCT image. The X-ray source and the detector plane rotate around $z$-axis simultaneously within the range of $[0°, 360°)$, completing a full scan. Using a consistent angular step $\Delta\theta = 18°$, this setup creates 20 X-ray projections. For other configurations, like 10 or 5 views, we use larger angular steps of $36°$ and $72°$, respectively. The distance between X-ray source and object is denoted as $L^{sb}$ and the distance between object and detector is denoted as $L^{bd}$. In the following, we introduce the simulation of the $i$-th X-ray projection, where $i \in \{1, 2, \dots, N\}$ is the view index and $N$ is the total number of views.

For the $i$-th viewpoint, our objective is to simulate X-ray projection $\mathbf{P}_i \in \mathbb{R}^{1 \times w \times h}$ with a spatial resolution of $w \times h$ and a pixel spacing of $p_u \times p_v$ from the given CBCT image. The rotation angle is defined as $\theta_i = \Delta\theta(i-1)$. Then, the position of the X-ray source $\mathbf{O}_i^s \in \mathbb{R}^3$ and the position of the detector center $\mathbf{O}_i^d \in \mathbb{R}^3$ could be given by:

$$\mathbf{O}_i^s = \begin{bmatrix} L^{sb} \cos\theta_i \\ L^{sb} \sin\theta_i \\ 0 \end{bmatrix}, \quad \mathbf{O}_i^d = \begin{bmatrix} -L^{bd} \cos\theta_i \\ -L^{bd} \sin\theta_i \\ 0 \end{bmatrix} \quad (1)$$

Eq. 1 is derived under the assumption that $\mathbf{O}_i^s$, $\mathbf{O}_i^d$, and $\mathbf{O}^b$ are colinear in the $xy$-plane, with the $z$-axis values of $\mathbf{O}_i^s$ and $\mathbf{O}_i^d$ both being zeros [50]. The detector 3D basis vector between adjacent pixels in horizontal direction $\mathbf{u}_i \in \mathbb{R}^3$ and the basis vector in vertical direction $\mathbf{v}_i \in \mathbb{R}^3$ could be given by:

$$\mathbf{u}_i = \begin{bmatrix} -p_u \sin\theta_i \\ p_u \cos\theta_i \\ 0 \end{bmatrix}, \quad \mathbf{v}_i = \begin{bmatrix} 0 \\ 0 \\ p_v \end{bmatrix} \quad (2)$$

Given a pixel $(m,n)$ on detector plane, where $m \in \{0, 1, \dots, w-1\}, n \in \{0, 1, \dots, h-1\}$, we aim to derive its pixel center coordinate $\mathbf{x}_i^{mn} \in \mathbb{R}^3$ in world space. And we first derive the pixel center coordinate $\mathbf{x}_i^{00} \in \mathbb{R}^3$ of pixel $(0,0)$ as follows, suitable for both odd and even image resolution:

$$\mathbf{x}_i^{00} = \mathbf{O}_i^d - \left( \left\lfloor \frac{w}{2} \right\rfloor + \left\lfloor \frac{w+1}{2} \right\rfloor - \frac{w+1}{2} \right) \mathbf{u}_i$$
$$- \left( \left\lfloor \frac{h}{2} \right\rfloor + \left\lfloor \frac{h+1}{2} \right\rfloor - \frac{h+1}{2} \right) \mathbf{v}_i \quad (3)$$

where $\lfloor \cdot \rfloor$ denotes the flooring operation. Then $\mathbf{x}_i^{mn}$ could be given by:

$$\mathbf{x}_i^{mn} = \mathbf{x}_i^{00} + m\mathbf{u}_i + n\mathbf{v}_i \quad (4)$$

Consider an incident X-ray path, $\mathbf{r}_i^{mn}(t) = \mathbf{O}_i^s + t\mathbf{d}_i^{mn} \in \mathbb{R}^3$, radiating from the X-ray source $\mathbf{O}_i^s$ towards the $\mathbf{x}_i^{mn}$ on the detector, and $t \in \mathbb{R}_0^+$ is the scaling parameter. The direction vector $\mathbf{d}_i^{mn} \in \mathbb{R}^3$ could be formulated as:

$$\mathbf{d}_i^{mn} = \mathbf{x}_i^{mn} - \mathbf{O}_i^s \quad (5)$$

This ray intersects the bounding box of the CBCT volume at entry point $\mathbf{r}_i^{mn}(t_n)$ and exit point $\mathbf{r}_i^{mn}(t_f)$ according to Ray-AABB algorithm [51]. Here, $t_n$ and $t_f$ denote the near and far bounds of the volume, respectively. We then simulate the ray integral of the attenuation value along the ray:

$$\hat{P}(\mathbf{r}_i^{mn}|\mathbf{V}_{gt}) = \int_{t_n}^{t_f} \hat{\mu}(\mathbf{r}_i^{mn}(t)|\mathbf{V}_{gt}) \mathrm{d}t \quad (6)$$

where $\hat{P}(\mathbf{r}_i^{mn}|\mathbf{V}_{gt}) \in \mathbb{R}_0^+$ is the simulated projection value. $\hat{\mu} : \left( \mathbb{R}^3 | \mathbb{R}^{1 \times W \times H \times D} \right) \to \mathbb{R}_0^+$ is a function that describes the attenuation distribution of the given CBCT image, and $\hat{\mu}(\mathbf{r}_i^{mn}(t)|\mathbf{V}_{gt})$ is the attenuation value for point $\mathbf{r}_i^{mn}(t)$, defined as follows:

$$\hat{\mu}(\mathbf{r}_i^{mn}(t)|\mathbf{V}_{gt}) = \text{Interp}_k \left( \mathbf{V}_{gt}, \mathbf{r}_i^{mn}(t) \right), k = 3 \quad (7)$$

where $\text{Interp}_k : \left( \mathbb{R}^{C_1 \times D_1 \times D_2 \times \dots \times D_k}, \mathbb{R}^k \right) \to \mathbb{R}^{C_1}$ is $k$-linear interpolation. Here $C_1 = 1$ represents the attenuation value, and we use trilinear interpolation in the above equation. For

computational feasibility, we discretize Eq. 6 by uniformly sampling points between the near and far bounds as shown in Fig. 2, deriving the following equation:

$$\hat{P}(\mathbf{r}_i^{mn}|\mathbf{V}_{gt}) = \sum_j \hat{\mu}(\mathbf{r}_i^{mn}(t_j)|\mathbf{V}_{gt})\delta_j \qquad (8)$$

where $\mathbf{r}_i^{mn}(t_j)$ represents the $j$-th sampling point, and $\delta_j = 0.5\min(v_x, v_y, v_z)$ is the distance between adjacent points.

By combining the projection values of all pixels on the detector, we can finally simulate the X-ray projection $\mathbf{P}_i$ for the $i$-th viewpoint. Similarly, we can synthesize X-ray projections from all viewpoints $\{\mathbf{P}_i\}_{i=1}^N$ for the given CBCT volume $\mathbf{V}_{gt}$. In this way, we obtain an X-ray and CBCT paired data. For each of dental and spine dataset, we collect 130 CBCT/CT volumes, and simulate the X-ray projections as mentioned above.

*2) Real-World X-ray and CBCT Pairs Collection:* However, the noise-free X-ray projections from our simulated dataset are always too ideal. In real-world collection, X-ray projections inevitably include Possion noise, and electronic noise introduced by the imaging system. To verify our robustness, we further collect another real-world captured CBCT walnut dataset [48]. In the following, we provide a detailed formulation of X-ray attenuation process in real-world collection, and the details of the walnut dataset.

Similarly, imaging an X-ray path $\mathbf{r}_i^{mn}(t)$ as we defined in Sec. III-A.1, the X-ray attenuation process along this path can be expressed according to Beer's Law [7]:

$$I(\mathbf{r}_i^{mn}) = (I_0(\mathbf{r}_i^{mn}) - I_1(\mathbf{r}_i^{mn}))\exp\left(-\int \mu(\mathbf{r}_i^{mn}(t))\,\mathrm{d}t\right) + I_1(\mathbf{r}_i^{mn}) \qquad (9)$$

here $I(\mathbf{r}_i^{mn}) \in \mathbb{Z}_0^+$ denotes the X-ray photon count recorded by the detector, $I_0(\mathbf{r}_i^{mn}) \in \mathbb{Z}_0^+$ denotes the X-ray source emitted X-ray photon count (flat-field image), and $I_1(\mathbf{r}_i^{mn}) \in \mathbb{Z}_0^+$ denotes the detector offset count (dark-field image). All these quantities are pixel-dependent. $\mu : \mathbb{R}^3 \to \mathbb{R}_0^+$ is a function that describes the attenuation distribution of the scanned scene, and $\mu(\mathbf{r}_i^{mn}(t))$ is attenuation value for point $\mathbf{r}_i^{mn}(t)$. By conducting flat-field and dark-field corrections [48], we would get the X-ray projection value as follows:

$$P(\mathbf{r}_i^{mn}) = -\ln\left(\frac{I(\mathbf{r}_i^{mn}) - I_1(\mathbf{r}_i^{mn})}{I_0(\mathbf{r}_i^{mn}) - I_1(\mathbf{r}_i^{mn})}\right) = \int \mu(\mathbf{r}_i^{mn}(t))\mathrm{d}t \qquad (10)$$

In this dataset, each walnut was scanned three times with the X-ray source set at three different heights: high, middle, and low position. Each scan comprises 1201 views with a uniform increment of $0.3°$ within a range of $[0°, 360°]$, where the first and the last views are captured at the same location. For each scan, flat-field and dark-field images were also recorded to pre-process the raw photon count data, as defined in Eq. 10. The processed X-ray projections from three scans were combined together to obtain a reference ground-truth CBCT reconstruction that is free from cone angle artifacts. In our experiments, we utilize $N$ uniformly spaced X-ray projections $\{\mathbf{P}_i\}_{i=1}^N$ selected from middle scan and the reference CBCT volume $\mathbf{V}_{gt}$ as paired data. We use $N = 20, 10, 5$, resulting in angular

increments of $18°$, $36°$, and $72°$, respectively. Note that in this dataset, CBCT scanning geometry parameters including $\mathbf{O}_i^s$, $\mathbf{O}_i^d$, $\mathbf{u}_i$, and $\mathbf{v}_i$ are not calculated by Eqs. 1–2, but are recorded in geometry description files provided by [48].

### B. Geometry-Aware Attenuation Learning

Fig. 3 shows the design of our proposed method. Given a set of sparse-view X-ray projections $\{\mathbf{P}_i\}_{i=1}^N$, our objective is to reconstruct the CBCT image $\mathbf{V}_{pred} \in \mathbb{R}^{1\times W\times H\times D}$ to be as close as possible to the ground truth one $\mathbf{V}_{gt}$. We first extract their features using a shared 2D CNN encoder. Then, we build the 3D feature map by combining feature back projection module with an adaptive feature fusion process. After this, the 3D feature map is processed by a 3D CNN decoder to recover the target CBCT image. In particular, one key step of our method is the feature back projection, which effectively bridges the dimensional gap between 2D X-ray projections and 3D CBCT image. It is crucial for achieving anatomically precise reconstructions. Secondly, the prior knowledge learned from datasets by deep learning, equips our model with the capability to adapt to new patients without individual optimization. It is particularly beneficial in handling sparse-view inputs with limited information, especially when dealing with 5 or 10 views. These two factors are the primary reasons why our model could produce satisfactory results with sparse input. Furthermore, our volume-wise CNN decoder acts as a learnable filter to reduce noise and extract more robust feature representations. It helps us better capture the global structure information of the target CBCT image, mitigating streaky artifacts. We will delve into the details of our method in the following.

*1) 2D Feature Extraction:* Given a set of X-ray projections $\{\mathbf{P}_i\}_{i=1}^N$ that we defined in Sec. III-A, we employ a shared 2D CNN encoder to extract a 2D feature map for each projection. This process allows us to capture view-specific feature representations, expressed as $\{\mathbf{F}_i\}_{i=1}^N \subset \mathbb{R}^{C\times w\times h}$, where $C$ is the feature channel size.

*2) Feature Back Projection:* The key of our framework is feature back projection, crafted to align with the geometry of CBCT imaging system. The idea is to retrieve the view-specific pixel-aligned feature representation $\mathbf{f}_i \in \mathbb{R}^C$ from $\mathbf{F}_i$ given any 3D query point $\mathbf{x} \in \mathbb{R}^3$ in world space, where $\mathbf{x}$ is sampled from volumetric coordinate matrix $\mathcal{V}$, i.e., $\mathbf{x} \in \mathcal{V}$. We first need to project $\mathbf{x}$ onto detector plane as view-specific projected point $\mathbf{x}_i \in \mathbb{R}^3$ and then transfer it into pixel $\mathbf{x}_i' \in \mathbb{R}^2$ for bilinear interpolation from $\mathbf{F}_i$. This process leverages the CBCT scanning geometry parameters, including the source position $\mathbf{O}_i^s$, the detector center $\mathbf{O}_i^d$, and the detector plane basis vectors $\mathbf{u}_i$ and $\mathbf{v}_i$ that we mentioned in Sec. III-A.

In the following, we provide mathematically detailed process of feature back projection. A visual illustration is given in Fig. 4. If an X-ray path $\mathbf{r}_i(t) = \mathbf{O}_i^s + t\mathbf{d}_i \in \mathbb{R}^3$ radiating from X-ray source position $\mathbf{O}_i^s$ and passing through $\mathbf{x}$, it will intersect the detector plane with projected point $\mathbf{x}_i$ for the $i$-th view. The ray direction $\mathbf{d}_i \in \mathbb{R}^3$ could be denoted as:

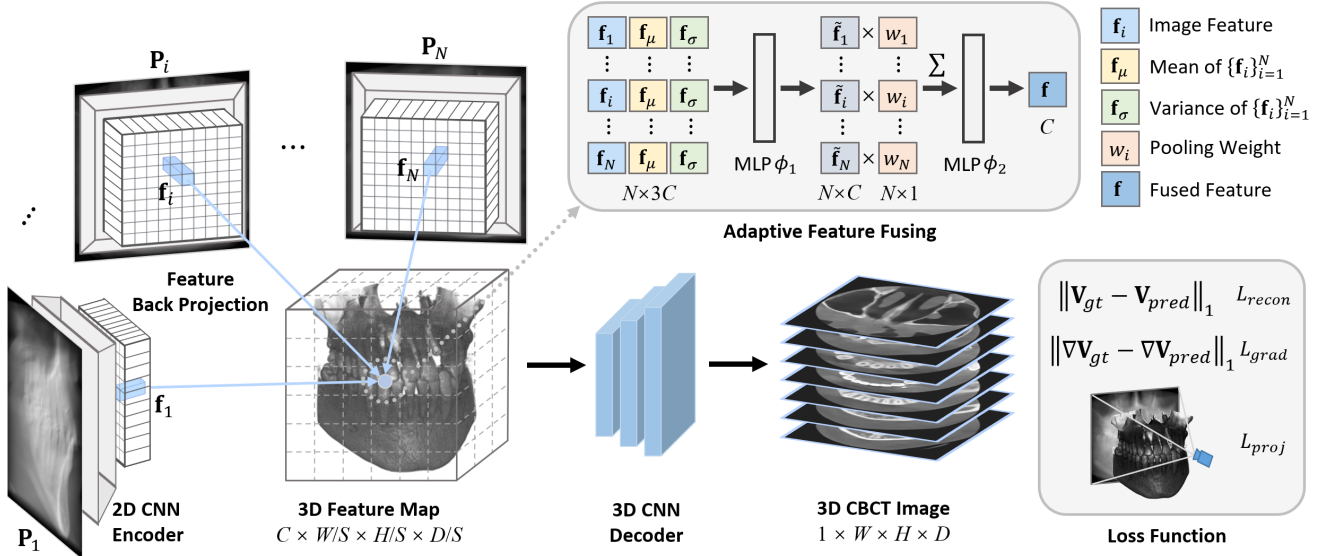$$\mathbf{d}_i = \mathbf{x} - \mathbf{O}_i^s \qquad (11)$$

Fig. 3.    Overview of our proposed method. A 2D CNN encoder first extracts feature representations from multi-view X-ray projections. Then, we build a 3D feature map by feature back projection and adaptive feature fusing. Finally, this 3D feature map is fed into a 3D CNN decoder to produce the final CBCT image.
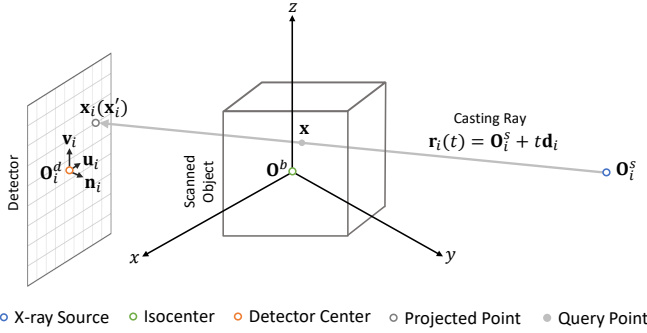


Fig. 4.    Coordinate transformation of query point for feature back projection.

The normal vector $\mathbf{n}_i$ of detector plane could be denoted as:

$$\mathbf{n}_i = \mathbf{u}_i \times \mathbf{v}_i \tag{12}$$

$\mathbf{x}_i$ would satisfy the following plane function as it belongs to the detector plane:

$$\mathbf{n}_i \cdot (\mathbf{x}_i - \mathbf{O}_i^d) = 0 \tag{13}$$

and $\mathbf{x}_i$ also satisfy the ray function:

$$\mathbf{x}_i = \mathbf{O}_i^s + t_{\mathbf{x}_i}\mathbf{d}_i \tag{14}$$

Substitute Eq. 14 into Eq. 13, we could get:

$$\mathbf{n}_i \cdot (\mathbf{O}_i^s - \mathbf{O}_i^d) + t_{\mathbf{x}_i}(\mathbf{n}_i \cdot \mathbf{d}_i) = 0 \tag{15}$$

In this way, we could get the scaling parameter $t_{\mathbf{x}_i}$ for point $\mathbf{x}_i$ as:

$$t_{\mathbf{x}_i} = -\frac{\mathbf{n}_i \cdot (\mathbf{O}_i^s - \mathbf{O}_i^d)}{\mathbf{n}_i \cdot \mathbf{d}_i} \tag{16}$$

Then $\mathbf{x}_i$ could be represented using $\mathbf{x}$, $\mathbf{O}_i^s$, $\mathbf{O}_i^d$, $\mathbf{u}_i$, and $\mathbf{v}_i$:

$$\mathbf{x}_i = \mathbf{O}_i^s - \frac{(\mathbf{u}_i \times \mathbf{v}_i) \cdot (\mathbf{O}_i^s - \mathbf{O}_i^d)}{(\mathbf{u}_i \times \mathbf{v}_i) \cdot (\mathbf{x} - \mathbf{O}_i^s)} (\mathbf{x} - \mathbf{O}_i^s) \tag{17}$$

Following that, we transform $\mathbf{x}_i$ into pixel $\mathbf{x}_i'$:

$$\mathbf{x}_i' = \begin{bmatrix} (\mathbf{x}_i - \mathbf{x}_i^{00}) \cdot \mathbf{u}_i/\|\mathbf{u}_i\|_2^2 \\ (\mathbf{x}_i - \mathbf{x}_i^{00}) \cdot \mathbf{v}_i/\|\mathbf{v}_i\|_2^2 \end{bmatrix} \tag{18}$$

Then, we could get view-specific pixel-aligned features $\mathbf{f}_i$ from $\mathbf{F}_i$ by bilinear interpolation:

$$\mathbf{f}_i = \text{Interp}_k(\mathbf{F}_i, \mathbf{x}_i'), k = 2 \tag{19}$$

Similarly, we obtain multi-view feature vectors $\{\mathbf{f}_i\}_{i=1}^N$ for a query point $\mathbf{x}$ from different 2D feature maps $\{\mathbf{F}_i\}_{i=1}^N$.

This design bridges the dimension gap between 2D X-ray projections and 3D CBCT image and facilitates accurate information retrieval from multi-view X-ray projections for 3D spatial queries, which is crucial for achieving anatomically precise reconstructions. It distinguishes our method from previous end-to-end learning-based methods [13]–[15] that brutally concatenate information of different views.

*3) Adaptive Feature Fusing:* After gathering the multi-view feature vectors $\{\mathbf{f}_i\}_{i=1}^N$ for a query point $\mathbf{x}$ through feature back projection, our goal is to merge these feature vectors into a point-wise feature vector $\mathbf{f} \in \mathbb{R}^C$. Different view X-ray projections present different information due to the variations in X-ray attenuation process introduced by varying viewpoints. To account for this property during the feature merging process, we draw inspiration from [52] and employ an adaptive feature fusion mechanism to effectively integrate these feature vectors.

For the multi-view feature vectors $\{\mathbf{f}_i\}_{i=1}^N$ associated with the query point $\mathbf{x}$, we start by calculating an element-wise average vector $\mathbf{f}_\mu \in \mathbb{R}^C$ and a variance vector $\mathbf{f}_\sigma \in \mathbb{R}^C$ to capture global information across $N$ views:

$$\mathbf{f}_\mu = \frac{1}{N}\sum_{i=1}^N \mathbf{f}_i, \ \ \mathbf{f}_\sigma = \frac{1}{N}\sum_{i=1}^N (\mathbf{f}_i - \mathbf{f}_\mu)^2 \tag{20}$$

Each $\mathbf{f}_i$ represents local information retrieved from the $i$-th X-ray projection, while $\mathbf{f}_\mu$ and $\mathbf{f}_\sigma$ capture the common

characteristics and variabilities across different viewpoints. These local features are concatenated with $\mathbf{f}_\mu$ and $\mathbf{f}_\sigma$, and then fed into the first MLP $\phi_1 : \mathbb{R}^{3C} \to \mathbb{R}^{C+1}$ to merge both local and global information:

$$\tilde{\mathbf{f}}_i, \tilde{w}_i = \phi_1 \left( \mathbf{f}_i \oplus \mathbf{f}_\mu \oplus \mathbf{f}_\sigma \right), \ i = 1, 2, \ldots, N \qquad (21)$$

here $\oplus$ denotes concatenation operation. Note that in this context, the local or global information represents the single-view or cross-view information for this specific query point $\mathbf{x}$. In this way, local feature $\mathbf{f}_i$ could interact with other views information through global features $\mathbf{f}_\mu, \mathbf{f}_\sigma$ with the help of $\phi_1$. This process produces a refined feature vector $\tilde{\mathbf{f}}_i \in \mathbb{R}^C$ and an unnormalized weight $\tilde{w}_i \in \mathbb{R}$ for each view. And $\tilde{w}_i$ further undergoes Softmax normalization to obtain a normalized weight $w_i \in [0, 1]$ as follows:

$$w_i = \frac{\exp(\tilde{w}_i)}{\sum_{l=1}^{N} \exp(\tilde{w}_l)}, \ i = 1, 2, \ldots, N \qquad (22)$$

Finally, $\{\tilde{\mathbf{f}}_i\}_{i=1}^N$ are summed together using their respective weights $\{w_i\}_{i=1}^N$ and passed through the second MLP $\phi_2 : \mathbb{R}^C \to \mathbb{R}^C$, resulting in the well-integrated feature vector $\mathbf{f}$:

$$\mathbf{f} = \phi_2 \left( \sum_{i=1}^{N} w_i \tilde{\mathbf{f}}_i \right) \qquad (23)$$

The weight $w_i$ reflects the relative importance of the $i$-th view, and both $\phi_1$ and $\phi_2$ contain one neuron layer with GELU output activation.

This fusion technique excels at integrating features from various views, enhancing the ability to capture fine details in regions of CBCT images with typically low contrast.

*4) Attenuation Decoding:* After obtaining the feature vector $\mathbf{f}$ for each query point $\mathbf{x} \in \mathcal{V}$, we compose these vectors to form the volumetric feature map. To mitigate GPU memory overhead and speed up computations, we employ a sparse sampling technique on $\mathcal{V}$ with a downsampling rate $S \in \{2^a | a \in \mathbb{Z}^+\}$. This approach enables us to generate a lower-resolution 3D feature map $\mathcal{F} \in \mathbb{R}^{C \times \frac{W}{S} \times \frac{H}{S} \times \frac{D}{S}}$. This assembled 3D feature map serves as the feature representation for the target CBCT image. We then feed it into our 3D CNN decoder to recover the CBCT image $\mathbf{V}_{pred}$, decoding the feature representations as attenuation values.

Our CNN-based decoder considers interactions among neighboring query points. This approach effectively acts as a learnable filter to mitigate noise and extracts more robust feature representations. And the volume-wise 3D supervision will help our model capture the global structural information of the target CBCT image. Consequently, our reconstructions exhibit high quality with less noises and streaky artifacts.

*5) Model Optimization:* To effectively train our framework, we incorporate several loss terms to guide and supervise the training process. First, we define the reconstruction loss $L_{recon}$ to enforce voxel-wise similarity between the ground-truth $\mathbf{V}_{gt}$ and the prediction $\mathbf{V}_{pred}$:

$$L_{recon} = \|\mathbf{V}_{gt} - \mathbf{V}_{pred}\|_1 \qquad (24)$$

To capture finer details, we introduce a gradient loss $L_{grad}$:

$$L_{grad} = \|\nabla \mathbf{V}_{gt} - \nabla \mathbf{V}_{pred}\|_1 \qquad (25)$$

here $\nabla$ denotes the first derivative. To make the model output align with the input X-ray projections, we further incorporate projection loss $L_{proj}$. For simulated dataset, $L_{proj}$ is defined as follows:

$$L_{proj} = \sum_{\mathbf{r}_i^{mn} \in \mathbf{B}} \|\hat{P}(\mathbf{r}_i^{mn}|\mathbf{V}_{gt}) - \hat{P}(\mathbf{r}_i^{mn}|\mathbf{V}_{pred})\|_1, \quad (26)$$

where $\hat{P}(\mathbf{r}_i^{mn}|\mathbf{V}_{pred})$ could be easily derived from Eq. 8, and $\mathbf{B} \subset \mathbb{R}^3$ is the random sampled ray batch set from input views. As for real-world dataset, $L_{proj}$ is defined as follows:

$$L_{proj} = \sum_{\mathbf{r}_i^{mn} \in \mathbf{B}} \|P(\mathbf{r}_i^{mn}) - \hat{P}(\mathbf{r}_i^{mn}|\mathbf{V}_{pred})\|_1, \qquad (27)$$

Therefore, our final objective function is defined as:

$$L = L_{recon} + \lambda_{grad} L_{grad} + \lambda_{proj} L_{proj}, \qquad (28)$$

where $\lambda_{grad}$ and $\lambda_{proj}$ determine the relative importance of the gradient and projection loss terms, respectively.

## IV. EXPERIMENT

### A. Experimental Settings

*1) Dataset:* To evaluate the effectiveness of our framework, we conducted experiments using both simulated (dental and spine [49]) and real-world dataset (walnut [48]). Below are the specifics of these datasets.

For the dental dataset, we collected 130 dental CBCT images from various patients. Each image is characterized by a resolution of $256 \times 256 \times 256$ and a voxel size of $0.3133\text{mm} \times 0.3133\text{mm} \times 0.3133\text{mm}$. Of these images, 100 images are used for training, 10 for validation, and the remaining 20 for testing. In our experiments, we test the reconstruction of CBCT images under three different input scenarios, i.e., 5 views, 10 views, and 20 views, respectively. The corresponding X-ray projections, following the procedure described in Sec. III-A.1, have a resolution of $256 \times 256$ and a pixel size of $0.4386\text{mm} \times 0.4386\text{mm}$. The source-to-object distance $L^{sb} = 500\text{mm}$, and the object-to-detector distance $L^{bd} = 200\text{mm}$.

Our second dataset includes 130 spinal CT images sourced from CTSpine1K [49]. It is worth noting that spinal images typically encompass a variety of organs and soft tissues, often with limited contrast in intensity. This aspect poses significant challenges for sparse-view reconstruction, a task notably more challenging than dental images. Hence, including spinal images in our study enables us to verify robustness and versatility of our framework across different anatomical areas. To ensure uniformity across cases, we resampled, cropped, and padded each spinal CT image to a resolution of $256 \times 256 \times 256$ and a voxel size of $2\text{mm} \times 2\text{mm} \times 2\text{mm}$. We split the spinal images into three groups: 100 images for training, 10 for validation, and 20 for testing. We also replicated the same input view configurations (namely, 5, 10, and 20 projections) for each case, as detailed in Sec. III-A.1. Each projection has a resolution of $256 \times 256$ with a pixel size of $3\text{mm} \times 3\text{mm}$. The source-to-object distance $L^{sb} = 1000\text{mm}$, and the object-to-detector distance $L^{bd} = 500\text{mm}$.

To further verify our robustness in real-world setting, we collect the third dataset includes 42 real-world walnut CBCT scanning collections from [48]. We split this dataset into 32 cases for training, 5 cases for validation, and 5 cases for testing. The original X-ray projection has a resolution of $768 \times 972$ and a pixel size of $0.1496\text{mm} \times 0.1496\text{mm}$. We downsample the projection to a resolution of $256 \times 324$ and a pixel size of $0.4488\text{mm} \times 0.4488\text{mm}$. For each case, we use the downsampled X-ray projections from three positioned scans to reconstruct the reference CBCT image based on the iterative algorithm from [48]. The reference CBCT image has a resolution of $256 \times 256 \times 256$ and a voxel size of $0.1961\text{mm} \times 0.1961\text{mm} \times 0.1961\text{mm}$. Uniformly spaced 5, 10, 20 X-ray projections from middle positioned scan are used as model input as mentioned in Sec. III-A.2. The source-to-object distance $L^{sb} = 66\text{mm}$, and the object-to-detector distance $L^{bd} = 133\text{mm}$.

*2) Implementation Details:* We adopt ResNet34 [53] as the backbone of the 2D CNN encoder. The encoder consists of an initial CNN layer followed by four subsequent residual layers. For each query point, we concatenate the feature vectors interpolated from each layer's feature map to create the output feature vector. As for the 3D CNN decoder, we adopt the generator network structure from SRGAN [54], [55]. This decoder is composed of six residual blocks at the beginning and a variable number of upsampling blocks. Each upsampling block increases the 3D feature map resolution by a factor of two. For instance, when $S = 4$, there are two upsampling blocks. And we would include an additional upsampling block if $S$ doubled.

In our experiments, we empirically set $\lambda_{grad} = 1$, $\lambda_{proj} = 0.01$, downsampling rate $S = 4$, channel size $C = 256$, random ray batch size $|\mathbf{B}| = 1024$. We employ the Adam optimizer with an initial learning rate of $1 \times 10^{-4}$, which decays by a factor of 0.5 every 50 epochs. We train the model for 200 epochs with a training batch size of 1. All experiments are conducted on a single A100 GPU.

*3) Competing Methods and Evaluation Metrics:* In our study, we evaluate our framework against a diverse set of methods, including conventional techniques like FDK and SART, neural rendering-based approaches like NAF, SCOPE3D [42], and SNAF, deep learning-based approaches like PatRecon, PixelNeRF, DIF-Net, and DDS3D [34]. FDK and SART are implemented based on ASTRA-toolbox [56]. The other methods are directly adopted from their source codes, except the following ones. SCOPE [42] incorporates hash-based neural rendering technique and re-projection strategy for sparse-view 2D CT reconstruction specified for parallel/fan beam geometry. We have extended SCOPE to a cone beam setting for CBCT reconstruction based on NAF codebase, named SCOPE3D. The primary difference between SCOPE3D and NAF is that SCOPE3D employs a re-projection strategy on its well-trained model to synthesize dense-view projections. Following that, it replaces the projections at the corresponding views in the synthesized dense-view projections with the given sparse input. These swapped projections are then used for reconstruction with traditional methods such as FDK. In our experiments, we use SART for the re-projection reconstruction

instead. PixelNeRF is originally designed for natural scene, we modify its forward projection equation to make it suitable for X-ray imaging. DDS adopts 2D pretrained diffusion model to generate 3D CT image slice by slice. It employs parallel beam forward projection and $z$-axis TV regularization to guide its inference sampling process. We modify it to use cone beam forward projection for CBCT reconstruction, named DDS3D. To assess performance, we utilize the Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) [57] as our primary metrics. Furthermore, the efficiency analysis in terms of both time and memory would be presented in Sec. IV-B.4.

## B. Reconstruction Results

*1) Qualitative Results:* Fig. 5 presents a side-by-side comparison of 3D CBCT images reconstructed in axial slices from case #10 of dental dataset. It is evident that the FDK approach struggles with sparse-view input, leading to significant streaky artifacts because of the limited number of input views. Although SART reduces these artifacts, it often loses fine details in the process. When it comes to neural rendering-based methods, NAF achieves decent results with 20 views by incorporating neural rendering with hash encoding. Yet, its performance greatly diminishes with very few input views (such as 5 or 10), as it is optimized for individual objects and lacks prior knowledge learned from the data population. SCOPE3D shows similar performance, as the re-projection strategy offers negligible new information. SNAF demonsrates improvements due to its view augmentation strategy but still struggles with 10 or 5 views. PatRecon ignores geometric relationships among multi-view projections, which results in blurry reconstructions with erroneous structures. Benefiting from CNN layers, PixelNeRF enjoys prior knowledge and maintains multi-view consistency. But it tends to produce noticeable streaky artifacts due to its point-wise MLP decoding and 2D supervision. DIF-Net builds upon the principles of PixelNeRF, achieving better results due to its 3D supervision. The results with 20 views input are comparable to ours, with slight blurriness and noise as highlighted in the orange box. However, its performance degrades with sparser inputs, such as 5 views, exhibiting streaky artifacts due to its point-wise MLP decoding approach. This is because the point-wise MLP independently decodes the attenuation of each query point, disregarding the spatial relationships among neighboring voxel points in CBCT image. MLP decoder is also unable to capture the global structure information of CBCT image with the point-wise 3D supervision. As a result, it would deliver streaky artifacts, especially when facing extremely sparse input like 5 views. In contrast, our CNN-based decoding module considers interactions among neighboring points, effectively acting as a learnable filter to mitigate noise and extract more robust feature representations. Moreover, 3D CNN decoder is capable of capturing the global structure information with the volume-wise 3D supervision. Consequently, our reconstructed CBCT images exhibit higher quality with less streaky artifacts. Notably, our approach surpasses all other methods, providing reconstruction quality comparable to the ground truth with 20
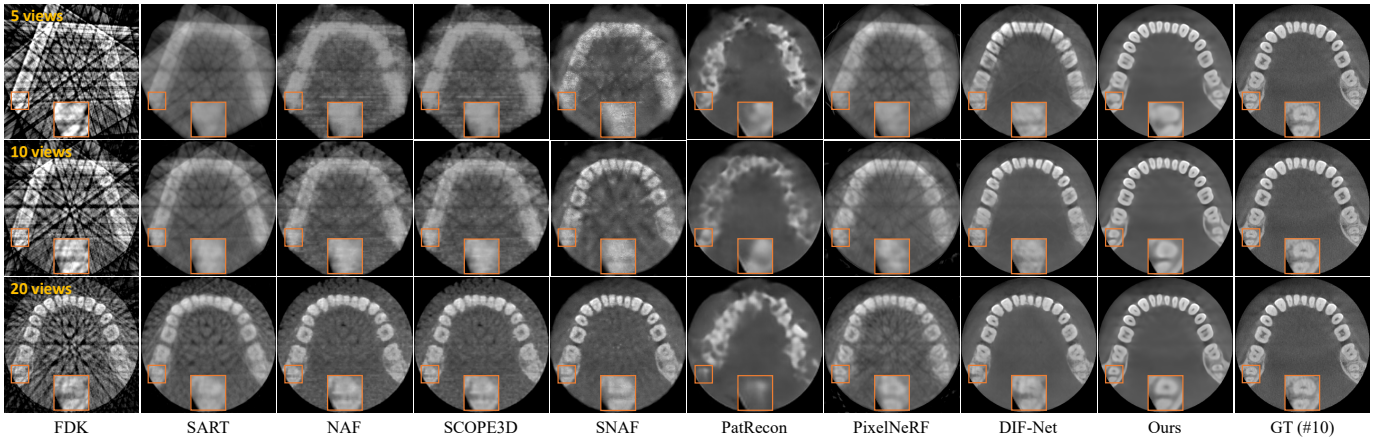
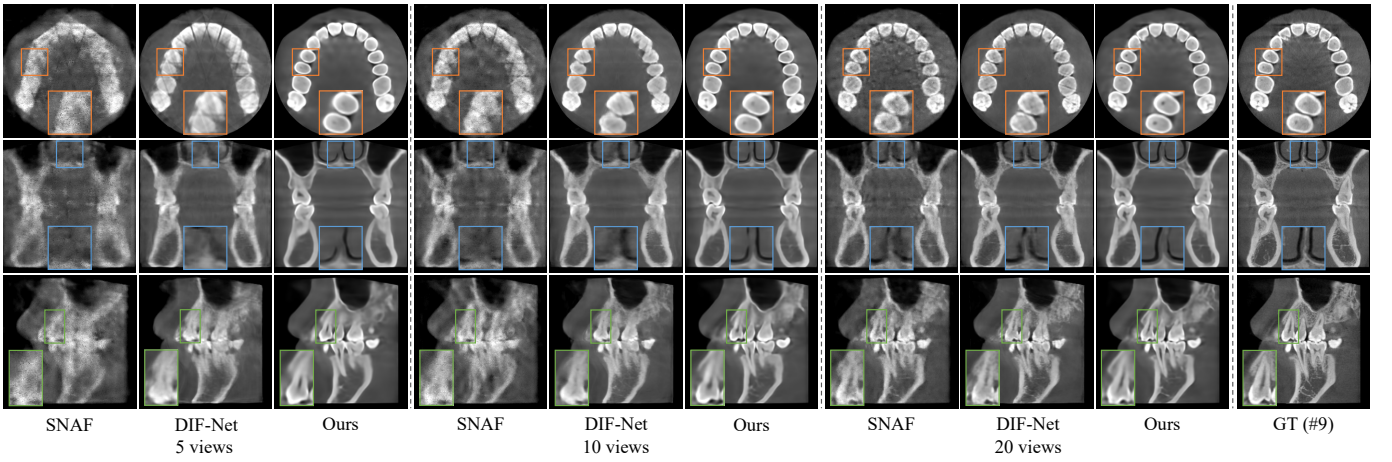Fig. 5. Qualitative comparison on case #10 from dental dataset (axial slice). Window: [-1000, 2000] HU.



Fig. 6. Qualitative comparison with SNAF and DIF-Net on case #9 from dental dataset. From top to bottom: axial, coronal, and sagittal slices. Window: [-1000, 2000] HU.
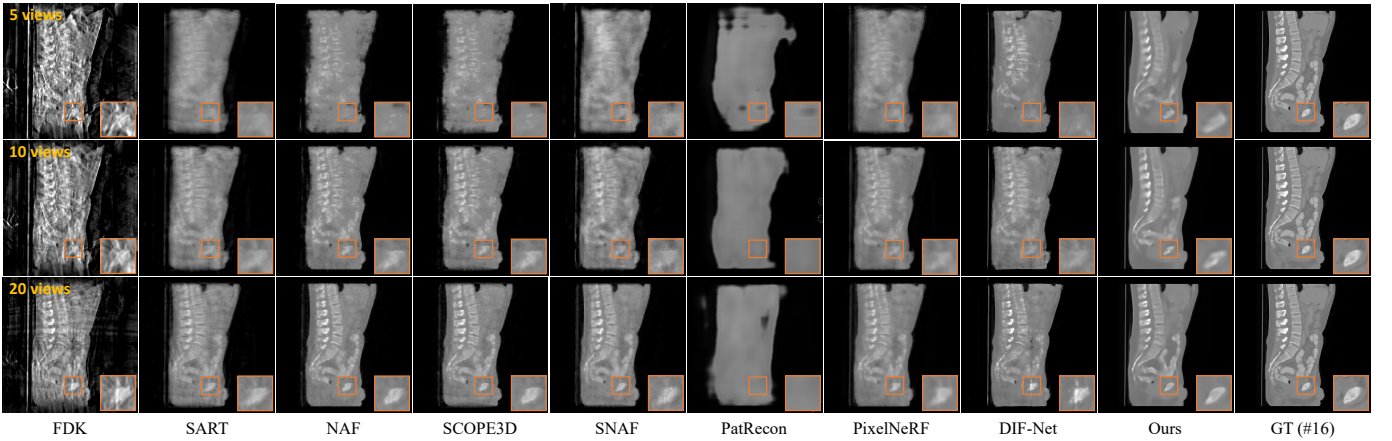


Fig. 7. Qualitative comparison on case #16 from spine dataset (sagittal slice). Window: [-1000, 1000] HU.

input views. However, recovering details with high fidelity becomes challenging for our method when facing 10 or 5 views. Despite this limitation, our method still maintains a clear advantage over the competition, showing less streaky artifacts and preserving a better global structure.

Additionally, we provide a detailed comparison with two current state-of-the-art methods, SNAF and DIF-Net, using case #9 from the dental dataset. This comparison includes axial, coronal, and sagittal slices, as shown in Fig. 6. It is clear that SNAF performs the worst, especially when limited to 10 or 5 views. This is because its lack of leveraging prior knowledge from data populations, unlike DIF-Net and our method. DIF-Net shows some noises and streaky artifacts, particularly with 5 views, and its detail recovery is not as good
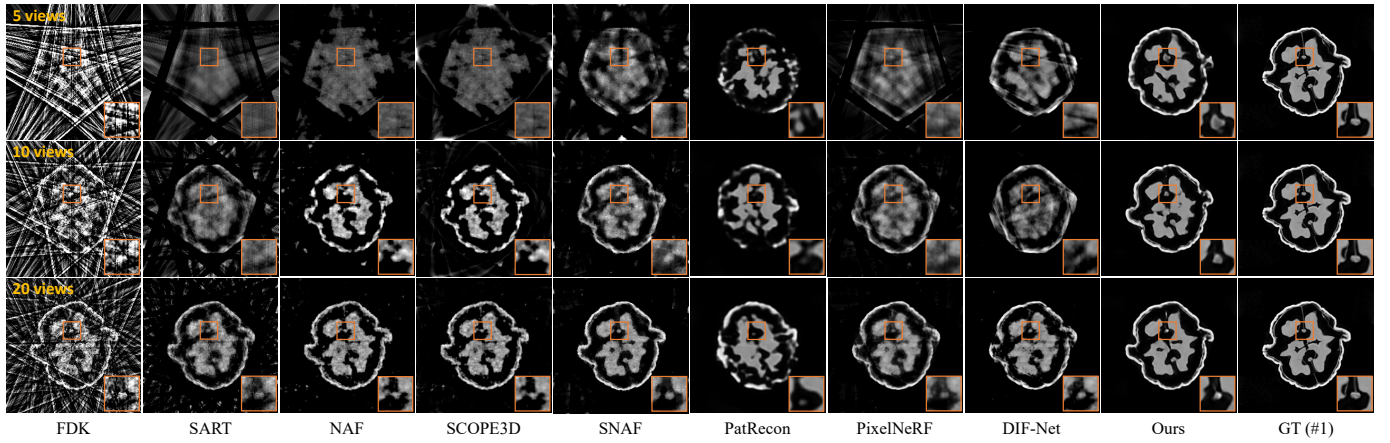
Fig. 8. Qualitative comparison on case #1 from walnut dataset (axial slice). Window: [-1000, 2000] HU.

as ours. Overall, our method achieves superior performance compared to the others. However, we also acknowledge that there can be some blurriness and minor structural errors with our method when limited to 10 or 5 views.

Fig. 7 provides visual comparison of reconstructed images in sagittal slices using case #16 from spinal dataset. We notice a consistent trend, with our method outperforming the others in performance. However, a common challenge emerges across all methods, including ours, in accurately reconstructing soft tissues and organs in the abdomen. This difficulty is due to the nature of spinal scans, which encompass a variety of organs and soft tissues. Many of these soft tissues exhibit low contrast differences, making it hard to achieve clear reconstructions with limited X-ray views. In comparison, dental scans primarily focus on teeth and jawbones, which present more distinct contrasts. Therefore, reconstructing sparse-view CBCT images is particularly challenging for spinal scans due to their intricate content. Despite these hurdles, the successful application of our method to both dental and spinal datasets showcases its versatility and robustness across different body parts.

Fig. 8 provides visual comparison of reconstructed CBCT image in axial slices using case #1 from walnut dataset. When dealing with 20 views input, SNAF and our method stand out with satisfactory results, albeit with some missing details. When reduced to 10 or 5 views, all methods, including ours, exhibit structural errors and do not perform well. Our method's results are relatively clear, with less noise and streaky artifacts, and the overall structure and quality are better compared to other methods. We have to admit that the walnut dataset is particularly difficult for reconstruction. This is partially due to the presence of noises from real-world captures. More importantly, the extremely low contrast between different regions within the walnut image poses another significant challenge. Despite such challenges, our method successfully adapts to real-world walnut dataset, further verifying our robustness. Based on the results from both the spine and walnut datasets, improving reconstruction quality for low-contrast images remains a challenging yet meaningful research direction.

*2) Quantitative Results:* To evaluate our method statistically, we present quantitative comparisons in Tab. I, Tab. II, and

## TABLE I
QUANTITATIVE COMPARISON ON DENTAL DATASET. THE BEST PERFORMANCE IS SHOWN IN BOLD.

| Method | 5 views | | 10 views | | 20 views | |
|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| FDK | 16.31±0.43 | 0.221±0.014 | 18.53±0.47 | 0.301±0.017 | 22.56±0.56 | 0.422±0.021 |
| SART | 20.22±1.82 | 0.621±0.024 | 24.17±1.13 | 0.704±0.018 | 27.93±0.91 | 0.784±0.015 |
| NAF | 21.55±1.21 | 0.578±0.036 | 23.89±1.14 | 0.673±0.024 | 28.77±0.86 | 0.793±0.020 |
| SCOPE3D | 21.74±0.99 | 0.588±0.024 | 24.05±1.12 | 0.684±0.023 | 29.39±0.87 | 0.807±0.019 |
| SNAF | 23.46±0.47 | 0.608±0.023 | 25.97±0.51 | 0.706±0.019 | 30.93±0.51 | 0.844±0.015 |
| PatRecon | 19.89±0.89 | 0.573±0.046 | 19.91±0.66 | 0.574±0.030 | 19.95±0.85 | 0.569±0.038 |
| PixelNeRF | 22.12±1.36 | 0.643±0.023 | 24.03±0.91 | 0.710±0.018 | 26.85±0.57 | 0.775±0.014 |
| DIF-Net | 25.80±0.93 | 0.759±0.027 | 27.52±0.84 | 0.818±0.021 | 30.48±0.80 | 0.870±0.015 |
| Ours | **27.48±1.12** | **0.823±0.028** | **28.83±1.19** | **0.850±0.025** | **31.44±1.00** | **0.891±0.016** |

## TABLE II
QUANTITATIVE COMPARISON ON SPINE DATASET. THE BEST PERFORMANCE IS SHOWN IN BOLD.

| Method | 5 views | | 10 views | | 20 views | |
|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| FDK | 17.06±1.27 | 0.258±0.051 | 20.27±1.30 | 0.286±0.044 | 22.87±1.37 | 0.352±0.041 |
| SART | 19.65±1.96 | 0.743±0.067 | 22.60±2.27 | 0.794±0.056 | 27.08±1.96 | 0.846±0.039 |
| NAF | 20.58±2.92 | 0.781±0.063 | 26.02±2.50 | 0.861±0.036 | 30.80±1.93 | 0.912±0.024 |
| SCOPE3D | 20.77±2.78 | 0.779±0.062 | 26.24±2.38 | 0.858±0.035 | 30.91±2.00 | 0.908±0.024 |
| SNAF | 22.05±1.48 | 0.711±0.075 | 26.33±1.85 | 0.806±0.062 | 31.69±1.39 | 0.902±0.035 |
| PatRecon | 18.30±1.62 | 0.686±0.054 | 18.29±2.14 | 0.681±0.061 | 18.57±1.75 | 0.641±0.047 |
| PixelNeRF | 20.81±1.98 | 0.772±0.075 | 25.41±1.67 | 0.848±0.043 | 28.68±1.64 | 0.890±0.033 |
| DIF-Net | 25.75±2.46 | 0.825±0.061 | 28.65±2.00 | 0.858±0.052 | 32.28±1.75 | 0.901±0.034 |
| DDS3D | 22.31±1.62 | 0.405±0.083 | 23.94±2.52 | 0.505±0.102 | 26.25±1.33 | 0.602±0.072 |
| Ours | **28.32±1.80** | **0.884±0.035** | **31.50±1.94** | **0.921±0.025** | **33.14±1.88** | **0.938±0.020** |

## TABLE III
QUANTITATIVE COMPARISON ON WALNUT DATASET. THE BEST PERFORMANCE IS SHOWN IN BOLD.

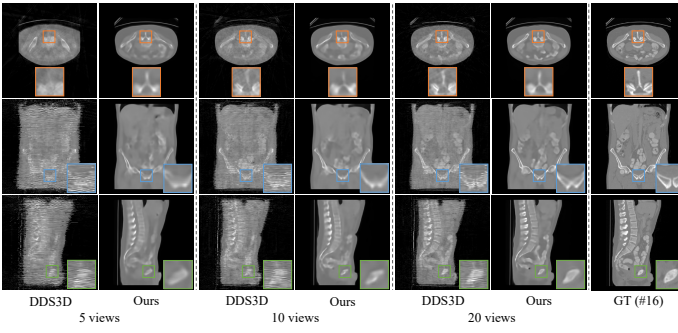| Method | 5 views | | 10 views | | 20 views | |
|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| FDK | 11.08±0.18 | 0.114±0.003 | 13.25±0.18 | 0.141±0.004 | 16.09±0.19 | 0.191±0.003 |
| SART | 17.25±0.29 | 0.378±0.013 | 19.55±0.29 | 0.487±0.016 | 23.71±0.40 | 0.588±0.012 |
| NAF | 18.15±1.17 | 0.522±0.047 | 21.72±0.82 | 0.612±0.027 | 25.36±0.22 | 0.680±0.006 |
| SCOPE3D | 19.08±0.64 | 0.523±0.046 | 20.75±1.27 | 0.535±0.046 | 25.92±0.19 | 0.714±0.009 |
| SNAF | 19.51±0.33 | 0.542±0.015 | 23.29±0.37 | 0.651±0.009 | **27.43±0.29** | 0.748±0.006 |
| PatRecon | 17.54±0.41 | 0.713±0.024 | 17.51±0.43 | 0.712±0.022 | 17.30±0.55 | 0.699±0.033 |
| PixelNeRF | 19.87±0.19 | 0.482±0.009 | 23.41±0.43 | 0.682±0.009 | 24.79±0.15 | 0.732±0.004 |
| DIF-Net | 21.30±0.34 | 0.640±0.006 | 21.76±0.45 | 0.640±0.008 | 25.99±0.14 | 0.707±0.004 |
| Ours | **21.64±0.53** | **0.811±0.012** | **24.78±0.35** | **0.867±0.010** | 26.84±0.37 | **0.895±0.008** |

Fig. 9. Qualitative comparison with DDS3D on case #16 from spine dataset. From top to bottom: axial, coronal, and sagittal slices. Window: [-1000, 1000] HU.

TABLE IV

EFFICIENCY ANALYSIS. THE BEST-PERFORMING IS SHOWN IN BOLD. UNIT: TIME (S); MEM. (GB): GPU MEMORY CONSUMPTION IN TRAINING; SIZE (MB).

| Method | 5 views | | 10 views | | 20 views | | Size |
|---|---|---|---|---|---|---|---|
| | Time | Mem. | Time | Mem. | Time | Mem. | |
| FDK | 0.48 | - | 0.51 | - | 0.50 | - | - |
| SART | 102.15 | - | 127.18 | - | 111.65 | - | - |
| NAF | 212.29 | 12.74 | 403.38 | 12.74 | 787.26 | 12.74 | 54.28 |
| SCOPE3D | 718.46 | 12.74 | 933.47 | 12.74 | 1303.97 | 12.74 | 54.28 |
| SNAF | 1681.20 | 7.03 | 1685.98 | **7.03** | 1806.09 | **7.03** | 134.83 |
| PatRecon | **0.007** | 21.25 | **0.007** | 21.25 | **0.012** | 21.25 | 1557.40 |
| PixelNeRF | 5.38 | 17.35 | 7.84 | 24.56 | 12.77 | 35.39 | **26.35** |
| DIF-Net | 1.60 | **5.05** | 3.35 | 7.82 | 7.62 | 14.45 | 118.66 |
| DDS3D | 1073.11 | 8.83 | 2072.44 | 8.83 | 4218.67 | 8.97 | 108.51 |
| Ours | 0.53 | 32.86 | 0.65 | 39.95 | 0.93 | 68.69 | 104.83 |

Tab. III for the dental, spine, walnut dataset, respectively. The performance trend is similar to qualitative results. For the dental and spine datasets, our method consistently outperforms others, as measured by PSNR and SSIM metrics. On the walnut dataset, our method achieves the best performance for both the 10 views and 5 views input conditions. Under the 20 views input condition, our method ranks second in PSNR, slightly lower than SNAF, but achieves the highest SSIM. A possible explanation is that methods based on per-scene optimization may be more robust to noisy real-world data. The noises vary among different cases, which may hinder the learning of generalization methods. Another reason is that the available training data is too limited on walnut dataset (32 cases), preventing our method from fully extracting prior knowledge from such a small data population. Due to these two factors, our method's PSNR metric is slightly lower than that of SNAF. Overall, considering all datasets, all view settings, and two metrics, our method achieves the best performance.

*3) Compare with Diffusion Model:* In this section, we compare our method with diffusion-based method DDS3D on spine dataset. The quantitative results are presented in Tab. II, and a visual example is shown in Fig. 9. It is clear that our method outperforms DDS3D by a large margin in both PSNR and SSIM metrics as well as in visual quality. When dealing with 20 views, DDS3D provides decent result on axial slices, although there are still some streaky artifacts and noises. If we

TABLE V

QUANTITATIVE RESULTS OF ABLATION STUDY ON FEATURE FUSING STRATEGY ON DENTAL DATASET.

| Fusing | PSNR | SSIM |
|---|---|---|
| Max | 29.87±1.17 | 0.870±0.020 |
| Average | 30.74±1.34 | 0.883±0.017 |
| $\mathbf{f}_i$ | 30.70±0.76 | 0.881±0.015 |
| $\mathbf{f}_i \oplus \mathbf{f}_\sigma$ | 31.19±0.86 | 0.884±0.015 |
| $\mathbf{f}_i \oplus \mathbf{f}_\mu$ | 31.27±0.96 | 0.889±0.015 |
| Adaptive | 31.44±1.00 | 0.891±0.016 |

observe the coronal and sagittal slices, we find that DDS3D exhibits severe inter-slice inconsistency due to its slice-wise diffusion sampling process. This inconsistency is particularly noticeable at the upper and lower ends, attributed to cone beam forward projection guided sampling. The sparse nature of cone beam ray casting at these ends leads to insufficient supervisory signals. This issue becomes more pronounced with fewer input views, consequently degrading the reconstruction quality. Overall, despite DDS [34] performing well in CT reconstruction with parallel beam geometry, its performance with cone beam geometry is still far from satisfactory.

*4) Efficiency Analysis:* We further evaluate each method's efficiency in terms of both time and memory. Time efficiency includes reconstruction time for one case, while memory efficiency encompasses model size and GPU memory consumption during the training stage. The quantitative results are presented in Tab. IV. Note that the metrics of SNAF are reported on its optimization stage, and those for DDS3D are reported on its diffusion sampling process.

Our method could complete reconstruction within a second, demonstrating exceptional time efficiency. It is much faster than individual optimization methods (e.g., SART, NAF, SCOPE3D, SNAF) and diffusion-based method (e.g., DDS3D), which are hindered by time-consuming iterative calculation and slice-wise diffusion sampling process, respectively. PixelNeRF and DIF-Net, while also time efficient, face delays because they need to query full-resolution features from different views during reconstruction. In contrast, our method samples query points at a lower resolution, making it even faster. What's more, our model size is also manageable compared with other methods. However, it is important to note that our model is not GPU memory efficient, as it requires considerable GPU memory during the training phase. And it will be included as one of our limitations.

*C. Ablation Study*

In this section, we study the influence of our different model components, including feature fusing strategy, loss term, loss weight, and downsampling rate. Note that, our model is trained on dental dataset under the setting of downsampling rate $S = 4$, 20 input views, and uniformly sampled angles within the range of $[0°, 360°)$, unless otherwise specified.

*1) Feature Fusing Strategy:* We first study the influence of different feature fusing strategies, including max pooling, average pooling, and our designed adaptive pooling. Besides,
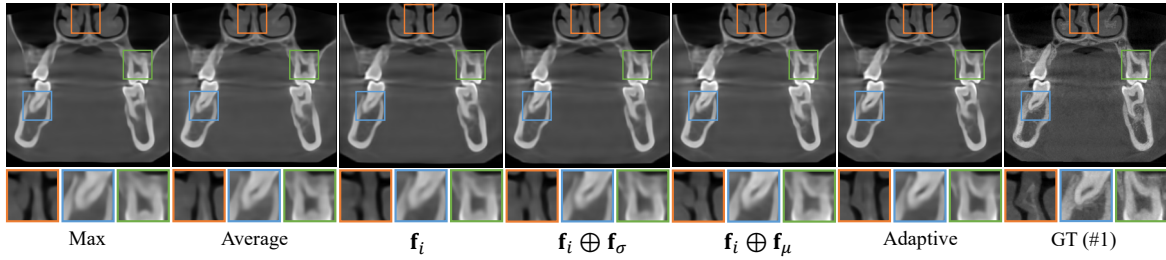
Fig. 10.  Qualitative results of ablation study on feature fusing strategy on case #1 from dental dataset (coronal slice). Window: [-1000, 2000] HU.

TABLE VI
QUANTITATIVE RESULTS OF ABLATION STUDY ON LOSS TERM ON
DENTAL DATASET.

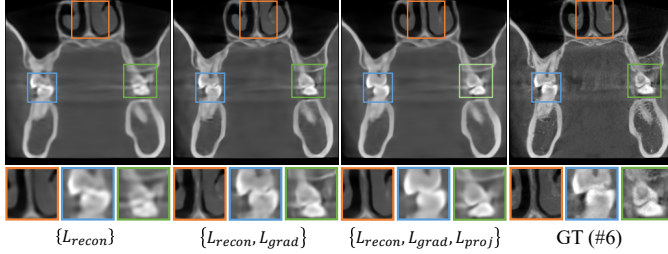| $L_{recon}$ | $L_{grad}$ | $L_{proj}$ | PSNR | SSIM |
|:---:|:---:|:---:|:---:|:---:|
| ✓ | | | 30.33±1.05 | 0.870±0.016 |
| ✓ | ✓ | | 30.93±0.94 | 0.887±0.016 |
| ✓ | ✓ | ✓ | 31.44±1.00 | 0.891±0.016 |



Fig. 11.  Qualitative results of ablation study on loss term on case #6 from dental dataset (coronal slice). Window: [-1000, 2000] HU.

we ablate each component $\mathbf{f}_i$, $\mathbf{f}_\mu$, and $\mathbf{f}_\sigma$ in our adaptive fusion module to study their importance. The quantitative results are delivered in Tab. V, and a visual example is shown in Fig. 10. Take the setting of $\mathbf{f}_i$ for example, the input of the first MLP block $\phi_1$ in Eq. 21 is replaced by $\mathbf{f}_i$.

It's obvious that max pooling performs the worst due to significant information loss. Average pooling retains all the information from different views, leading to performance improvement. However, the feature from each view lacks interaction, limiting performance gains. For instance, the low-contrast regions highlighted in Fig. 10, such as the nasal cavity and tooth pulp, remain poorly reconstructed. While the setting of $\mathbf{f}_i$ also retains information from various views, it still lacks interaction among different views. Both of $\mathbf{f}_i \oplus \mathbf{f}_\sigma$ and $\mathbf{f}_i \oplus \mathbf{f}_\mu$ incorporate global feature term $\mathbf{f}_\sigma$ or $\mathbf{f}_\mu$, enabling $\mathbf{f}_i$ to interact with features from other views through the first MLP block $\phi_1$, resulting in more effective feature fusion and enhanced performance. Moreover, incorporating $\mathbf{f}_\mu$ yields a slight better performance than $\mathbf{f}_\sigma$, indicating higher importance of $\mathbf{f}_\mu$. Finally, our full adaptive feature fusing design achieves the best results, with visual details closely matching the ground truth and yielding the highest metrics. Overall, the adaptive feature fusion strategy demonstrates superior reconstruction performance, capturing finer details, especially in regions with low contrast.

*2) Loss Term:* Second, we also study the effectiveness of our different loss terms, including the 3D reconstruction
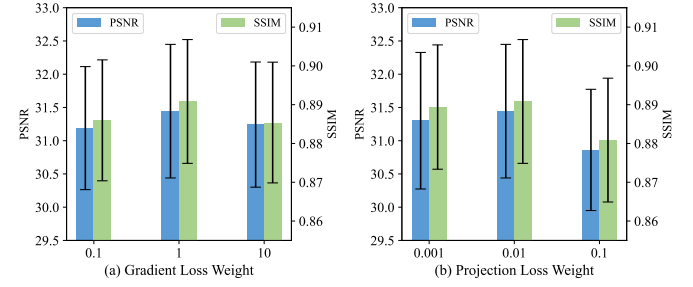


Fig. 12.  Quantitative results of ablation study on (a) gradient loss weight and (b) projection loss weight on dental dataset.

TABLE VII
ABLATION STUDY ON DOWNSAMPLING RATE $S$. UNIT: TIME (S); SIZE
(MB); MEM. (GB): GPU MEMORY CONSUMPTION IN TRAINING.

| $N$ | $S$ | PNSR | SSIM | Time | Size | Mem. |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | 16 | 24.44 ± 0.71 | 0.748 ± 0.027 | 0.21 | 104.94 | 11.43 |
| 5 | 8 | 26.35 ± 0.92 | 0.793 ± 0.028 | 0.35 | 104.88 | 13.85 |
| | 4 | 27.48 ± 0.71 | 0.823 ± 0.029 | 0.53 | 104.83 | 32.86 |
| | 16 | 24.98 ± 0.75 | 0.760 ± 0.026 | 0.33 | 104.94 | 12.49 |
| 10 | 8 | 27.21 ± 0.88 | 0.816 ± 0.026 | 0.43 | 104.88 | 15.91 |
| | 4 | 27.21 ± 0.88 | 0.842 ± 0.021 | 0.65 | 104.83 | 39.95 |
| | 16 | 25.68 ± 0.67 | 0.776 ± 0.025 | 0.45 | 104.94 | 14.59 |
| 20 | 8 | 28.66 ± 0.91 | 0.844 ± 0.021 | 0.90 | 104.88 | 19.47 |
| | 4 | 31.44 ± 1.00 | 0.891 ± 0.016 | 0.95 | 104.83 | 68.69 |

loss ($L_{recon}$), gradient loss ($L_{grad}$), and 2D projection loss ($L_{proj}$). Quantitative results are presented in Tab. VI. We use $L_{recon}$ as our baseline, and each symbol ✓ indicates the inclusion of the particular loss term in the training process, thereby offering an alternative solution. Our baseline model, using just the 3D reconstruction loss $L_{recon}$, already shows impressive performance, achieving over 30 dB for the dental dataset. This highlights the importance of 3D supervision. With the addition of $L_{grad}$ and $L_{proj}$, the PSNR and SSIM values increase gradually. A visual comparison is available in Fig. 11. It is noticeable that $L_{grad}$ assists in recovering sharper details around teeth pulp areas. However, it can also introduce some noise. $L_{proj}$, acting as a regularization term, effectively reduces these noises and contributes to more accurate recovery of anatomical structures.

*3) Loss Weight:* Third, we verify our selection of the loss weights $\lambda_{grad}$ and $\lambda_{proj}$, which determine the relative importance of the gradient loss and projection loss, respectively. Gradient loss helps model more effectively recover details, and

TABLE VIII
QUANTITATIVE RESULTS OF ROBUSTNESS ANALYSIS ON
RECONSTRUCTION RESOLUTION ON DENTAL DATASET.

| Test Resolution | PSNR | SSIM |
|---|---|---|
| 128 res | 24.85±0.43 | 0.783±0.017 |
| 256 res | 31.44±1.00 | 0.891±0.016 |

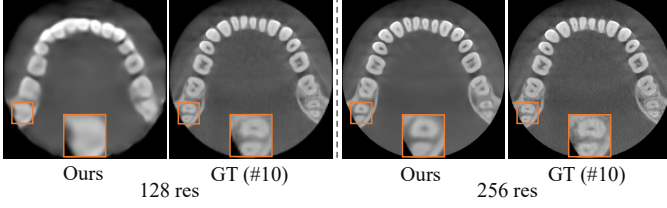

Ours 128 res    GT (#10)    Ours 256 res    GT (#10)

Fig. 13. Qualitative results of robustness analysis on reconstruction resolution on case #10 from dental dataset (axial slice). Window: [-1000, 2000] HU.

TABLE IX
QUANTITATIVE RESULTS OF ROBUSTNESS ANALYSIS ON ANGLE
SAMPLING ON DENTAL DATASET.

| Test Angle | PSNR | SSIM |
|---|---|---|
| Uniform $[0°, 360°)$ | 31.44±1.00 | 0.891±0.016 |
| Uniform $[5°, 360°)$ | 30.97±0.90 | 0.884±0.016 |
| Uniform $[10°, 360°)$ | 30.28±0.76 | 0.873±0.016 |
| Uniform $[20°, 360°)$ | 29.67±0.74 | 0.863±0.016 |
| Random $[0°, 360°)$ | 25.04±1.76 | 0.766±0.046 |



GT (#10)   Uniform [0°,360°)   Uniform [5°,360°)   Uniform [10°,360°)   Uniform [20°,360°)   Random [0°,360°)

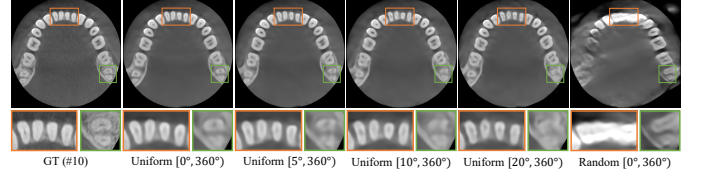Fig. 14. Qualitative results of robustness analysis on angle sampling on case #10 from dental dataset (axial slice). Window: [-1000, 2000] HU.

projection loss enforces the model output aligns with the input X-ray projections. In our work, we consider gradient loss is more important that projection loss.

Determining the optimal weights $\lambda_{grad}$ and $\lambda_{proj}$ through grid search can be experimentally laborious. Therefore, we perform a rough verification to support our choice. The quantitative results with both mean and standard deviation are presented in Fig. 12. For the ablation study on the gradient loss weight $\lambda_{grad}$, we fixed $\lambda_{proj} = 0.01$. Conversely, for the ablation study on the projection loss weight $\lambda_{proj}$, we fixed $\lambda_{grad} = 1$. The results clearly demonstrate that our chosen values $\lambda_{grad} = 1$ and $\lambda_{proj} = 0.01$ yield the best performance in terms of both PSNR and SSIM metrics.

*4) Downsampling Rate:* At last, we assess the effect of the downsampling rate $S$. The quantitative results are presented in Tab. VII. Note that, $N$ represents the number of input views in this table. The performance metrics (PSNR, SSIM) are reported on dental dataset. The trend is clear that, as the downsampling rate $S$ decreases, there is an improvement in reconstruction performance, as well as an increase in reconstruction time, and GPU memory usage in training stage. The model size slightly decreases for less unsampling blocks in decoder. A smaller $S$ means denser sampling on the voxel grid, which allows for capturing of more X-ray projection information, thereby enhancing the reconstruction quality. However, this benefit comes with the trade-off of higher computational demands and longer sampling time. $S = 4$ offers the best performance among 5, 10, and 20 views, and the computation cost is also affordable. Therefore, we choose $S = 4$ for our experiments.

### D. Robustness Analysis

In this section, we analyze our well-trained framework robustness to varying test conditions that differ from the training settings, including reconstruction resolution, angle sampling, and the number of input views. Note that, our model is trained on dental dataset under the setting of downsampling rate $S = 4$, 20 input views, and uniformly sampled angles within the range of $[0°, 360°)$, unless otherwise specified.

*1) Reconstruction Resolution:* We first verify whether our well-trained framework can perform reconstructions at different resolutions. We train our model with a resolution of $256^3$ using a downsampling rate $S = 4$, and test our model with a resolution of $128^3$ using $S = 8$. Note that the ground truth CBCT image is also downsampled to $128^3$ resolution for statistical analysis. The quantitative results are delivered in Tab. VIII, and a visual example is shown in Fig. 13. The downsampling rate $S = 8$ retains only $\frac{1}{8}$ of the sample points from $S = 4$ when constructing 3D feature map. Testing with $128^3$ resolution leads to deteriorated results due to insufficient input information. In conclusion, our current framework faces challenge in performing reconstructions at varying resolutions. We would modify our training strategy by considering different resolutions to overcome it in the future.

*2) Angle Sampling:* Second, we examine whether our framework is robust to the different viewing angles when testing. Our model is trained with uniformly sampled angles within a range of $[0°, 360°)$ and a consistent angular step. We first test with different starting angles for uniform sampling, e.g., $5°$, $10°$, and $20°$. Second, we test with randomly sampled angles within a range of $[0°, 360°)$ for inconsistent angular step. The quantitative results are delivered in Tab. IX, and a visual example is shown in Fig. 14. As the starting angle increases from $0°$ to $5°$, $10°$ and then to $20°$, the reconstruction quality gradually decreases. This decline occurs because the input view sampling angles progressively deviate from the training setting, making adaption more challenging for our model. When testing with random sampling, the reconstruction quality drops significantly and the metric variances increase. This is because the randomly sampled angles deviate considerably from the training setting. And the angle deviation varies among different test samples, leading to increased variances. In conclusion, our current framework is not robust to different starting angles or inconsistent angular steps, especially for random sampling angles. We would solve it by using different angle sampling settings during training stage in the future.

*3) Number of Input Views:* At last, we examine whether our framework is robust to the number of input views when testing.

TABLE X

QUANTITATIVE RESULTS OF ROBUSTNESS ANALYSIS ON NUMBER OF INPUT VIEWS ON DENTAL DATASET.

| Test Views | PSNR | SSIM |
|---|---|---|
| 5 views | 23.95±0.55 | 0.715±0.029 |
| 10 views | 28.83±1.19 | 0.850±0.025 |
| 20 views | 28.54±0.78 | 0.842±0.021 |

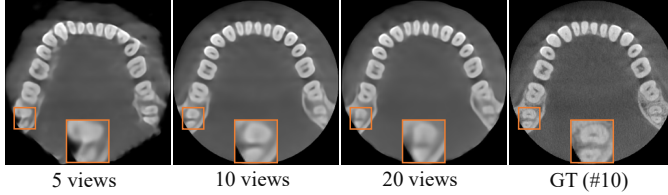

5 views    10 views    20 views    GT (#10)

Fig. 15. Qualitative results of robustness analysis on number of input views on case #10 from dental dataset (axial slice). Window: [-1000, 2000] HU.

We train our model with 10 input views, and test it with 5 or 20 views. The quantitative results are delivered in Tab. X, and a visual example is shown in Fig. 15. According to Sec. III-A, the input of 20 views retains all the sampling angels from 10 views. In contrast, the input of 5 views retains only a half of the sampling angles from 10 views. When training with 10 views, testing with 5 views results in degradation due to insufficient sampling angles, but testing with 20 views is acceptable as 20 views contain sufficient sampling angles. Hence, our model can still achieve satisfactory reconstruction results when the sampling angles of test input include the training ones. Otherwise, the reconstruction quality may degrade due to insufficient input information. Thus, our model is not robust to the number of input views. In the future, we plan to take different number of views as input during the training stage to address it.

## V. DISCUSSION

While our framework has shown promising results, it does have its limitations. Firstly, one significant challenge is that the method is resource-intensive and not computational memory efficient. Secondly, our reconstruction results for low contrast images, such as spine and walnut images, are not yet up to the mark. The reconstructed results often turn out too smooth, lacking finer details, And there can even be structural errors in walnut images with 10 or 5 input views. Thirdly, our current framework may be not robust to varying test conditions that differ from the training ones, including reconstruction resolution, angle sampling, and the number of input views. At last, real-world experiments on walnut dataset do not fully reflect the real-life scenario as walnut scans lack significant scattering effects. Throughout our paper, we do not take scattering effects into consideration and only account for the primary rays. However, in real-life thorax and pelvic scans, scattering effects are much more severe, even producing more scatter rays than primary rays. This usually leads to increased image artifacts and noises. Our method's robustness in real-life clinical scenario requires further verification.

In the future, we aim to improve our reconstruction performance by integrating explicit shape or boundary priors to help correct structural errors. Additionally, we plan to design robust training strategies that enable the model to adapt to varying test conditions. What's more, it would be valuable to verify our method's robustness using real-life projection data with significant scattering effects, such as scans of thorax and pelvic. We could also incorporate scattering effects to improve our algorithm for more accurate image reconstruction in clinical practice.

## VI. CONCLUSION

In this paper, we introduced a novel framework for sparse-view CBCT reconstruction. Our method respects the inherent nature of X-ray perspective projection during the feature back projection, ensuring accurate information retrieval from multiple X-ray projections. Moreover, by leveraging the prior knowledge learned from our extensive dataset, our framework efficiently tackles the challenges posed by sparse-view inputs, delivering high-quality reconstructions. The effectiveness and time efficiency are thoroughly validated through extensive testing on both simulated and real-world datasets.

## REFERENCES

[1] T. Kaasalainen, M. Ekholm, T. Siiskonen, and M. Kortesniemi, "Dental cone beam ct: An updated review," *Physica Medica*, vol. 88, pp. 193–217, 2021.

[2] G. P. John, T. E. Joy, J. Mathew, and V. R. Kumar, "Fundamentals of cone beam computed tomography for a prosthodontist," *The Journal of the Indian Prosthodontic Society*, vol. 15, no. 1, p. 8, 2015.

[3] W. Xu, G. Chang, E. Truong, S. Babu, E. Gan, S. Wang, A. Dayo, A. Le, and C. Rajapakse, "Cone beam computed tomography assessment of cervical spine bone density," *Journal of Clinical Densitometry*, vol. 25, no. 2, pp. 279–280, 2022.

[4] M. F. Powell, D. DiNobile, and A. S. Reddy, "C-arm fluoroscopic cone beam ct for guidance of minimally invasive spine interventions," *Pain Physician*, vol. 13, no. 1, p. 51, 2010.

[5] K. Shi, W. Xiao, G. Wu, Y. Xiao, Y. Lei, J. Yu, and Y. Gu, "Temporal-spatial feature extraction of dsa video and its application in avm diagnosis," *Frontiers in Neurology*, vol. 12, p. 655523, 2021.

[6] K. Ruedinger, S. Schafer, M. Speidel, and C. Strother, "4d-dsa: development and current neurovascular applications," *American Journal of Neuroradiology*, vol. 42, no. 2, pp. 214–220, 2021.

[7] A. C. Kak and M. Slaney, *Principles of computerized tomographic imaging*. SIAM, 2001.

[8] L. A. Feldkamp, L. C. Davis, and J. W. Kress, "Practical cone-beam algorithm," *Josa a*, vol. 1, no. 6, pp. 612–619, 1984.

[9] A. H. Andersen and A. C. Kak, "Simultaneous algebraic reconstruction technique (sart): a superior implementation of the art algorithm," *Ultrasonic imaging*, vol. 6, no. 1, pp. 81–94, 1984.

[10] E. Y. Sidky and X. Pan, "Image reconstruction in circular cone-beam computed tomography by constrained, total-variation minimization," *Physics in Medicine & Biology*, vol. 53, no. 17, p. 4777, 2008.

[11] E. Y. Sidky, C.-M. Kao, and X. Pan, "Accurate image reconstruction from few-views and limited-angle data in divergent-beam ct," *Journal of X-ray Science and Technology*, vol. 14, no. 2, pp. 119–139, 2006.

[12] Y. Zhu, Y. Liu, Q. Zhang, C. Zhang, and X. Gao, "A fast iteration approach to undersampled cone-beam ct reconstruction," *Journal of X-ray science and technology*, vol. 27, no. 1, pp. 111–129, 2019.

[13] L. Shen, W. Zhao, and L. Xing, "Patient-specific reconstruction of volumetric computed tomography images from a single projection view via deep learning," *Nature biomedical engineering*, vol. 3, no. 11, pp. 880–888, 2019.

[14] Y. Kasten, D. Doktofsky, and I. Kovler, "End-to-end convolutional neural network for 3d reconstruction of knee bones from bi-planar x-ray images," in *Machine Learning for Medical Image Reconstruction: Third International Workshop, MLMIR 2020, Held in Conjunction with MICCAI 2020, Lima, Peru, October 8, 2020, Proceedings 3*, pp. 123–133, Springer, 2020.

[15] X. Ying, H. Guo, K. Ma, J. Wu, Z. Weng, and Y. Zheng, "X2ct-gan: reconstructing ct from biplanar x-rays with generative adversarial networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10619–10628, 2019.

[16] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[17] G. Zang, R. Idoughi, R. Li, P. Wonka, and W. Heidrich, "Intratomo: self-supervised learning-based tomography via sinogram synthesis and prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1960–1970, 2021.

[18] D. Rückert, Y. Wang, R. Li, R. Idoughi, and W. Heidrich, "Neat: Neural adaptive tomography," *ACM Transactions on Graphics (TOG)*, vol. 41, no. 4, pp. 1–13, 2022.

[19] R. Zha, Y. Zhang, and H. Li, "Naf: neural attenuation fields for sparse-view cbct reconstruction," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 442–452, Springer, 2022.

[20] Y. Fang, L. Mei, C. Li, Y. Liu, W. Wang, Z. Cui, and D. Shen, "Snaf: Sparse-view cbct reconstruction with neural attenuation fields," *arXiv preprint arXiv:2211.17048*, 2022.

[21] Y. Cai, J. Wang, A. Yuille, Z. Zhou, and A. Wang, "Structure-aware sparse-view x-ray 3d reconstruction," in *CVPR*, 2024.

[22] K. H. Jin, M. T. McCann, E. Froustey, and M. Unser, "Deep convolutional neural network for inverse problems in imaging," *IEEE transactions on image processing*, vol. 26, no. 9, pp. 4509–4522, 2017.

[23] H. Lee, J. Lee, H. Kim, B. Cho, and S. Cho, "Deep-neural-network-based sinogram synthesis for sparse-view ct image reconstruction," *IEEE Transactions on Radiation and Plasma Medical Sciences*, vol. 3, no. 2, pp. 109–119, 2018.

[24] W. Wu, D. Hu, C. Niu, H. Yu, V. Vardhanabhuti, and G. Wang, "Drone: Dual-domain residual-based optimization network for sparse-view ct reconstruction," *IEEE Transactions on Medical Imaging*, vol. 40, no. 11, pp. 3002–3014, 2021.

[25] A. Hauptmann, J. Adler, S. Arridge, and O. Öktem, "Multi-scale learned iterative reconstruction," *IEEE transactions on computational imaging*, vol. 6, pp. 843–856, 2020.

[26] N. Moriakov, J.-J. Sonke, and J. Teuwen, "End-to-end memory-efficient reconstruction for cone beam ct," *Medical Physics*, vol. 50, no. 12, pp. 7579–7593, 2023.

[27] J. Pan, H. Yu, Z. Gao, S. Wang, H. Zhang, and W. Wu, "Iterative residual optimization network for limited-angle tomographic reconstruction," *IEEE Transactions on Image Processing*, 2024.

[28] H. Chung, J. Kim, M. T. Mccann, M. L. Klasky, and J. C. Ye, "Diffusion posterior sampling for general noisy inverse problems," *arXiv preprint arXiv:2209.14687*, 2022.

[29] J. Liu, R. Anirudh, J. J. Thiagarajan, S. He, K. A. Mohan, U. S. Kamilov, and H. Kim, "Dolce: A model-based probabilistic diffusion framework for limited-angle ct reconstruction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10498–10508, 2023.

[30] W. Wu, J. Pan, Y. Wang, S. Wang, and J. Zhang, "Multi-channel optimization generative model for stable ultra-sparse-view ct reconstruction," *IEEE Transactions on Medical Imaging*, 2024.

[31] Z. Li, D. Chang, Z. Zhang, F. Luo, Q. Liu, J. Zhang, G. Yang, and W. Wu, "Dual-domain collaborative diffusion sampling for multi-source stationary computed tomography reconstruction," *IEEE Transactions on Medical Imaging*, 2024.

[32] Y. Wang, Z. Li, and W. Wu, "Time-reversion fast-sampling score-based model for limited-angle ct reconstruction," *IEEE Transactions on Medical Imaging*, 2024.

[33] H. Chung, D. Ryu, M. T. McCann, M. L. Klasky, and J. C. Ye, "Solving 3d inverse problems using pre-trained 2d diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22542–22551, 2023.

[34] H. Chung, S. Lee, and J. C. Ye, "Decomposed diffusion sampler for accelerating large-scale inverse problems," *arXiv preprint arXiv:2303.05754*, 2024.

[35] J. N. Martel, D. B. Lindell, C. Z. Lin, E. R. Chan, M. Monteiro, and G. Wetzstein, "Acorn: adaptive coordinate networks for neural scene representation," *ACM Transactions on Graphics (TOG)*, vol. 40, no. 4, pp. 1–13, 2021.

[36] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Transactions on Graphics (ToG)*, vol. 41, no. 4, pp. 1–15, 2022.

[37] M. Niemeyer, J. T. Barron, B. Mildenhall, M. S. Sajjadi, A. Geiger, and N. Radwan, "Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5480–5490, 2022.

[38] A. Yu, V. Ye, M. Tancik, and A. Kanazawa, "pixelnerf: Neural radiance fields from one or few images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4578–4587, 2021.

[39] A. Chen, Z. Xu, F. Zhao, X. Zhang, F. Xiang, J. Yu, and H. Su, "Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 14124–14133, 2021.

[40] A. Corona-Figueroa, J. Frawley, S. Bond-Taylor, S. Bethapudi, H. P. Shum, and C. G. Willcocks, "Mednerf: Medical neural radiance fields for reconstructing 3d-aware ct-projections from a single x-ray," in *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pp. 3843–3848, IEEE, 2022.

[41] L. Shen, J. Pauly, and L. Xing, "Nerp: implicit neural representation learning with prior embedding for sparsely sampled image reconstruction," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[42] Q. Wu, R. Feng, H. Wei, J. Yu, and Y. Zhang, "Self-supervised coordinate projection network for sparse-view computed tomography," *IEEE Transactions on Computational Imaging*, 2023.

[43] Q. Wu, X. Li, H. Wei, J. Yu, and Y. Zhang, "Joint rigid motion correction and sparse-view ct via self-calibrating neural field," in *2023 IEEE 20th International Symposium on Biomedical Imaging (ISBI)*, pp. 1–5, IEEE, 2023.

[44] Q. Wu, L. Chen, C. Wang, H. Wei, S. K. Zhou, J. Yu, and Y. Zhang, "Unsupervised polychromatic neural representation for ct metal artifact reduction," *Advances in Neural Information Processing Systems*, vol. 36, pp. 69605–69624, 2023.

[45] Z. Zhou, H. Zhao, J. Fang, D. Xiang, L. Chen, L. Wu, F. Wu, W. Liu, C. Zheng, and X. Wang, "Tiavox: Time-aware attenuation voxels for sparse-view 4d dsa reconstruction," *arXiv preprint arXiv:2309.02318*, 2023.

[46] Z. Liu, H. Zhao, W. Qin, Z. Zhou, X. Wang, W. Wang, X. Lai, C. Zheng, D. Shen, and Z. Cui, "3d vessel reconstruction from sparse-view dynamic dsa images via vessel probability guided attenuation learning," *arXiv preprint arXiv:2405.10705*, 2024.

[47] Y. Lin, Z. Luo, W. Zhao, and X. Li, "Learning deep intensity field for extremely sparse-view cbct reconstruction," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 13–23, Springer, 2023.

[48] H. Der Sarkissian, F. Lucka, M. van Eijnatten, G. Colacicco, S. B. Coban, and K. J. Batenburg, "A cone-beam x-ray computed tomography data collection designed for machine learning," *Scientific data*, vol. 6, no. 1, p. 215, 2019.

[49] Y. Deng, C. Wang, Y. Hui, Q. Li, J. Li, S. Luo, M. Sun, Q. Quan, S. Yang, Y. Hao, *et al.*, "Ctspine1k: A large-scale dataset for spinal vertebrae segmentation in computed tomography," *arXiv preprint arXiv:2105.14711*, 2021.

[50] A. Brahme, *Comprehensive biomedical physics*. Newnes, 2014.

[51] J. Mahovsky and B. Wyvill, "Fast ray-axis aligned bounding box overlap tests with plucker coordinates," *Journal of Graphics Tools*, vol. 9, no. 1, pp. 35–46, 2004.

[52] Q. Wang, Z. Wang, K. Genova, P. P. Srinivasan, H. Zhou, J. T. Barron, R. Martin-Brualla, N. Snavely, and T. Funkhouser, "Ibrnet: Learning multi-view image-based rendering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4690–4699, 2021.

[53] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

[54] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, *et al.*, "Photo-realistic single image super-resolution using a generative adversarial network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4681–4690, 2017.

[55] I. Sánchez and V. Vilaplana, "Brain mri super-resolution using 3d generative adversarial networks," *arXiv preprint arXiv:1812.11440*, 2018.

[56] W. Van Aarle, W. J. Palenstijn, J. Cant, E. Janssens, F. Bleichrodt, A. Dabravolski, J. De Beenhouwer, K. J. Batenburg, and J. Sijbers, "Fast and flexible x-ray tomography using the astra toolbox," *Optics express*, vol. 24, no. 22, pp. 25129–25147, 2016.

[57] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.