

Robust Single-view Cone-beam X-ray Pose Estimation with Neural Tuned Tomography (NeTT) and Masked Neural Radiance Fields (mNeRF)

Chaochao Zhou^{1,*}, Syed Hasib Akhter Faruqui¹, Abhinav Patel¹, Ramez N. Abdalla¹, Michael C. Hurley^{1,2,3}, Ali Shaibani^{1,2,3}, Matthew B. Potts^{1,3}, Babak S. Jahromi^{1,3}, Leon Cho¹, Sameer A. Ansari^{1,2,3}, Donald R. Cantrell^{1,2,*}

¹Department of Radiology, ²Department of Neurology, and ³Department of Neurological Surgery, Northwestern University and Northwestern Medicine, Chicago, IL, United States

*Correspondence to: Donald R. Cantrell, MD, PhD (donald.cantrell@nm.org); Chaochao Zhou, PhD (chaochao.zhou@northwestern.edu)

Abstract

Purpose: Many tasks performed in image-guided, minimally invasive, medical procedures can be cast as pose estimation problems, where an X-ray projection is utilized to reach a target in 3D space. Recent advances in the differentiable rendering of optically reflective materials have enabled state-of-the-art performance in RGB camera view synthesis and pose estimation. Expanding on these prior works, we introduce new methods for pose estimation of radiolucent objects using X-ray projections, and we demonstrate the critical role of optimal view synthesis in performing this task.

Methods: We first develop an algorithm (DiffDRR) that efficiently computes Digitally Reconstructed Radiographs (DRRs) and leverages automatic differentiation within TensorFlow. In conjunction with classic Cone-Beam Computerized Tomography (CBCT) reconstruction algorithms, we perform pose estimation by gradient descent using a loss function that quantifies the similarity of the DRR synthesized from a randomly initialized pose and the true fluoroscopic image at the target pose. We propose two novel methods for high-fidelity view synthesis, Neural Tuned Tomography (NeTT) and masked Neural Radiance Fields (mNeRF). Both methods rely on

classic CBCT; NeTT directly optimizes the CBCT densities, while the non-zero values of mNeRF are constrained by a 3D mask of the anatomic region segmented from CBCT.

Results: We demonstrate that both NeTT and mNeRF distinctly improve pose estimation within our framework. By defining a successful pose estimate to be a 3D angle error of less than 3° , we find that NeTT and mNeRF can achieve similar results, both with overall success rates more than 93%. Furthermore, we show that a NeTT trained for a single subject can generalize to synthesize high-fidelity DRRs and ensure robust pose estimations for all other subjects.

Conclusion: Given the much lower computational cost and the cross-subject generalizability for NeTT, we suggest that NeTT is an attractive option for robust pose estimation using fluoroscopic projections.

Keywords: Differentiable Digitally Reconstructed Radiograph; Neural Tuned Tomography; Masked Neural Radiance Field; Image Registration; Pose Estimation

1. Introduction

Image-guided percutaneous and endovascular interventions are increasingly utilized in modern day healthcare as minimally-invasive alternatives to open surgical procedures. The management of neurovascular diseases has been revolutionized by these new image-guided techniques. Endovascular coiling is now the preferred treatment method for ruptured cerebral aneurysms, and endovascular aspiration thrombectomy is a proven, highly-effective, intervention for large vessel strokes [1, 2]. Both of these neurointerventions rely on X-ray fluoroscopy for intraoperative guidance. Additional image-guided neurointerventional procedures include biopsies, steroid injections for pain management, kyphoplasties, trigeminal rhizotomies, and sclerotherapy or embolization of vascular malformations and tumors.

Although 3D CT, CBCT, and Magnetic Resonance Imaging (MRI) data are readily available for the pre-procedural planning of neurointerventions, real-time *intraoperative* guidance is most often limited to 2D X-ray fluoroscopic projections of these 3D volumes. Furthermore, the patient position and orientation may change substantially between subsequent imaging procedures, or even intraoperatively as many minimally-invasive, image-guided procedures are performed

with the patient awake. Thus, many tasks in image-guided procedures can be cast as pose estimation problems where a series of 2D projections are used to reach an orientation or target in a 3D volume. This complex pose estimation and 2D/3D registration task, which often has a low tolerance for error, is currently performed in the clinical setting by the physician operators. However, machine learning aids for this critical process could greatly improve procedure efficiency and safety.

Initial efforts to leverage machine learning for this task largely utilized supervised training frameworks [3]. Due to the sizeable data requirements of this approach, as well as the limited availability of manually labeled medical datasets, many of these works relied on the generation of synthetic digitally reconstructed radiographs (DRRs) from 3D anatomic volumes [3]. For example, Miao et al. [4] proposed a method for single-view X-ray pose estimation in which a regressive convolutional neural network (CNN) was trained in a supervised manner to directly estimate the six transformation parameters of a DRR relative to a target radiograph using the difference of the two images. However, DRR synthesis methods that are based on ray-casting do not fully capture the physics and style of true X-ray fluoroscopic images, and this can be a limitation for pose estimation algorithms that rely on the comparison of DRRs to target radiographs. To address this limitation, Liao et al. [5] introduced a Siamese-like POINT network that establishes point-to-point correspondences between DRRs generated from a CBCT volume and an intraoperative X-ray and then estimates the pose of the 3D volume by triangulation. In another approach, Zhou et al. [6] explicitly incorporated style-transfer techniques into their method for landmark-based 2D/3D image registration to estimate *in vivo* skull pose using biplane fluoroscopic images. They first trained a CNN for landmark detection on DRRs that were generated from a small manually labeled 3D dataset. A cycle-consistent generative adversarial network (GAN) was then used to map real X-rays onto the DRR style domain where landmark detection and pose estimation by triangulation could be performed.

Outside of the field of medical imaging, there have recently been rapid advances in differentiable rendering, view synthesis, and pose estimation. In 2020, Mildenhall et al. [7] introduced the neural radiance field (NeRF) to represent real-world 3D scenes using a multi-layer perceptron (MLP) and to synthesize high-fidelity RGB images. The relatively small, fully connected, non-convolutional NeRF network takes the spatial coordinates (x,y,z) and the viewing angle (θ,ϕ) as input and outputs a volume density and a view-dependent RGB color. Using this

neural network scene representation, they also implemented a fully vectorized and differentiable rendering algorithm in a machine learning framework (TensorFlow) that is capable of efficiently generating 2D photorealistic images. Implementations of NeRF have now achieved state-of-the-art performance for the rendering of optically reflective materials [8, 9]. Furthermore, by leveraging the differentiability of NeRF, Yen-Chen et al. [10] subsequently introduced a gradient-based optimization framework that performs state-of-the-art 6DOF pose estimation given the NeRF scene representation and a single RGB view of the scene, which they termed “inverting” NeRF (iNeRF).

Motivated by these recent works in differentiable neural scene representation, here we introduce methods for pose estimation of radiolucent objects using 2D fluoroscopic projections, and we demonstrate the importance of optimal view synthesis for this task. First, we develop a differentiable framework of DiffDRR for generating DRRs with automatic differentiation. In conjunction with classic CBCT 3D volume reconstruction algorithms, we then perform pose estimation by iterative gradient descent using loss functions that quantify the similarity of synthesized DRRs and the true fluoroscopic image of the target pose. Next, we introduce two novel methods for neural scene representation and high-fidelity DRR view synthesis, Neural Tuned Tomography (NeTT) and masked Neural Radiance Fields (mNeRF). Each method relies on pre-computed classic CBCT reconstructions. NeTT directly optimizes the CBCT densities, while the non-zero output values of mNeRF are masked by an anatomic region segmented from the CBCT. Finally, we demonstrate that both of these techniques substantially improve pose estimation within our framework.

2. Methodology

2.1. Differentiable Digitally Reconstructed Radiograph (DiffDRR)

As illustrated in **Fig. 1**, our model consists of a movable cone-beam X-ray source and image intensifier (such as a fluoroscopy C-arm) with a 3D volume/field representing an anatomic structure fixed at the X-ray isocenter (i.e., the rotation center of the movable X-ray). The X-ray source and image intensifier translate and rotate together relative to the X-ray isocenter. Consistent with the NeRF rendering algorithm, we pad and scale our anatomic 3D CBCT volumes, mapping

them onto a normalized device coordinate (NDC) system, with X, Y, Z values in the $[-1, 1]^3$ cube and an origin set at the X-ray isocenter (O) [7]. The source-isocenter distance (SOD) and the source-intensifier distance (SID) are determined at the time of CBCT acquisition and define the geometry of the X-ray source and image intensifier relative to the anatomic volume. The SOD and SID must therefore also be mapped into the NDC so that the resulting projection geometries are maintained (*Appendix A*). We also normalize the scalar densities within the 3D volumes/fields to have values in the range $[0, 1]$. To synthesize DRRs in this framework, we adopted two algorithms, including Volume Rendering (VR) and Ray Casting (RC).

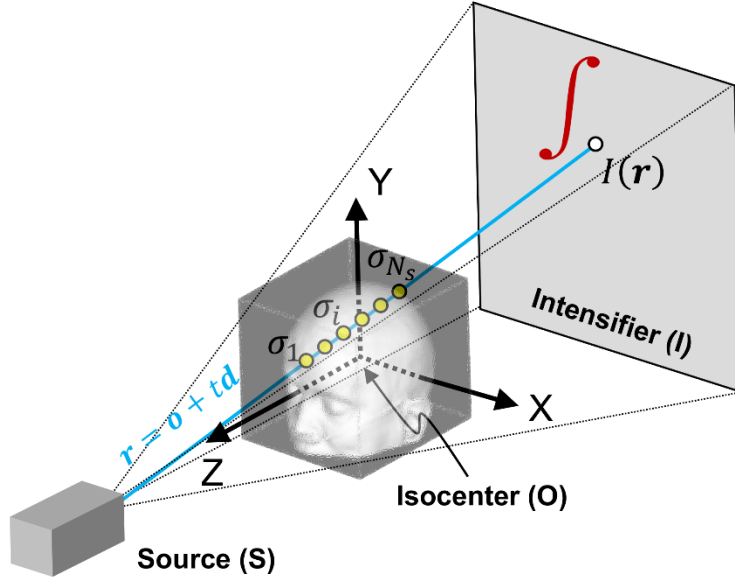


Fig. 1: Illustration of the DiffDRR geometry, consisting of an X-ray source (S) and an image intensifier (I). The NDC coordinate system origin is set at the isocenter (O). For an arbitrary ray (\mathbf{r}) cast from the source onto a single pixel of the image intensifier, the pixel intensity (I) is determined by the sampled densities ($\sigma_i, i = 1, 2, \dots, N_s$) of the 3D volume.

2.1.1. Volume Rendering (VR) Algorithm We modify the VR equations described by Mildenhall et al in their original paper on NeRF [7] to repurpose them for X-ray styled (grayscale) rendering. This method differs from the more commonly used Beer-Lambert model for generating X-ray projections [11], but it has already demonstrated success in rendering realistic, synthetic DRRs [12]. Although we describe a pixel “Intensity” in the VR equations below, this should be

interpreted only as the pixel luminance or brightness, as it does not directly correspond to the physical X-ray intensity.

We consider a single ray ($\mathbf{r} = \mathbf{o} + t\mathbf{d}$) which is emitted from the X-ray source (\mathbf{o}) along an arbitrary direction (\mathbf{d}), as shown in **Fig. 1**. The resulting image intensity (I) at the pixel on the image intensifier where the ray arrives can be formulated as a numerical integration of all intensities throughout this ray in the 3D space:

$$I(\mathbf{r}) = \sum_{i=1}^{N_s} w_i \sigma_i \quad (1)$$

where w_i are weights on each sampled volume density, σ_i , and N_s is the sample size per ray. To reduce computational cost, the intensities ($\sigma_i, i = 1, 2, \dots, N_s$) can be sampled only within the 3D volume by setting near and far bounds to correspond to the faces of the 3D volume along the Z-axis, as shown in **Fig. 1**.

The weights on each sampled volume density in **Eq. 1** can then be expressed as $w_i = \alpha_i T_i$, comprising an opacity of the sampled volume, α_i , and a transmittance, T_i , which describes the cumulative absorption of the ray by densities located between the source and the sampled volume with a nested numerical integration:

$$\alpha_i = 1 - \exp(-\sigma_i \delta_i) \quad (2)$$

$$T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right) \quad (3)$$

where $\delta_i = t_{i+1} - t_i$ is the spacing between adjacent sampled densities (t_i is the distance of σ_i to the source along the ray). The derivation of these equations and physical assumptions are well described by Tagliasacchi et al [13]. However, it is worth noting that the volume shading (RGB colors) in their formulations has been replaced with the volume luminance (densities) in our work.

2.1.2. Ray Casting (RC) Algorithm The VR equations described above are not intended to precisely describe the underlying physics. Instead, the physical RC process during X-ray imaging is most commonly approximated with the Beer-Lambert Law for generating X-ray

projections [11]. According to the Beer-Lamber Law, the resulting intensity of a pixel of the DRR along a single ray direction can be determined by the cumulative absorbance of the initially emitted X-ray intensity. With the same symbols defined above, RC Absorbance can be simply formulated below:

$$A(\mathbf{r}) = \sum_{i=1}^{N_s} \sigma_i \delta_i \quad (4)$$

The cumulative Absorbance is then related to the X-ray image intensity by $A = -\log(I/I_0)$, where I_0 is the intensity of the emitted X-ray beam, and I is the X-ray intensity at the pixel on the image intensifier. A detailed derivation of **Eq. 4** can be found in the previous work by Ruckert et al [14].

For both VR and RC (**Eq. 1 or 4**), to draw σ_i from a discretized CBCT volume, we implemented a fully-vectorized, differentiable sampling module in TensorFlow, which performed trilinear interpolation from neighboring volume densities. Moreover, we implemented two alternative methods for CBCT density sampling [7], including evenly-spaced, deterministic sampling (used for gradient-based pose optimization) and random, stochastic sampling (used for training NeTT and mNeRF). In both deterministic and stochastic sampling methods, the ray is first evenly divided into N_s segments within the anatomic volume. For deterministic sampling, the volume densities at the starting nodes of these segments were selected. For stochastic sampling, a volume density within each segment was randomly selected.

To generate all pixel intensities in a DRR image, **Eq. 1 or 4** must be evaluated for a total of $L \times W$ rays, where L and W are the number of pixels along the length and width of the DRR, respectively. In this work, the DRR image size was set to 128×128 , computed with a ray sampling size of $N_s = 128$.

In the remainder of this work, we refer to our algorithms for differentiable DRR view synthesis as DiffDRR, and we leverage DiffDRR for projection of both discretized 3D volumes and continuous neural density fields, including CBCT, NeTT, and mNeRF (*Sections 2.2 ~ 2.4*). DiffDRR also forms the backbone of our framework for gradient-based pose optimization (*Section 2.5*).

2.2. Cone-beam X-ray Imaging and CT Reconstruction

We analyzed the fluoroscopic X-ray sequences of five de-identified cerebrovascular CBCT acquisitions, as illustrated in **Fig. 2a**. Each X-ray sequence includes 133 X-ray images with a size of 960×960 covering a range of principal (axial) angles from -100° to 100° in an increment of 1.5° . We pre-processed the raw X-ray images, sequentially performing contrast enhancement (tone mapping and histogram equalization), intensity (grayscale) inversion, intensity scaling to $[0, 255]$, and image resizing to 256×256 . A sample post-processed X-ray sequence is presented in **Fig. 2b**. For traditional CBCT reconstruction, the processed X-ray sequence was imported into a GPU-based CT reconstruction toolbox, TIGRE [15]. Within TIGRE, the SOD and SID (which were unavailable in the DICOM metadata) were set to 1000 mm and 1536 mm, respectively, and the OSSART reconstruction algorithm was adopted. A sample CBCT reconstruction with a size of $256 \times 256 \times 256$ is presented in **Fig. 2c**.

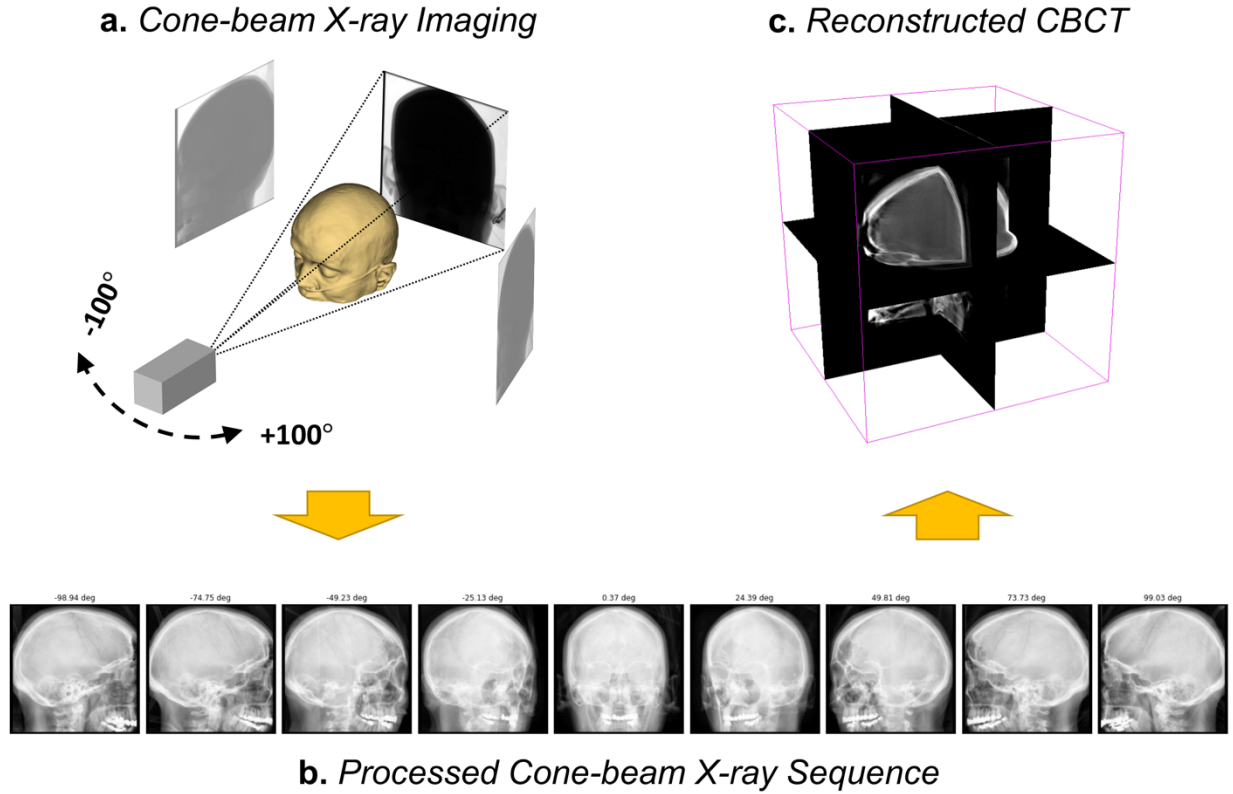


Fig. 2: Our data acquisition and pre-processing pipeline includes: **a)** Acquisition of cone-beam X-ray images with a size of 960×960 over a range of -100° to 100° ; **b)** Processed X-ray images undergoing grayscale inversion, contrast enhancement, intensity scaling, and downsampling to a size of 256×256 ; and **c)** CBCT reconstruction in TIGRE with a size of $256 \times 256 \times 256$.

2.3. Neural Tuned Tomography (NeTT)

Discrepancies in style and texture may exist between the real X-rays of the CBCT tomographic sequence and the synthetic DRR projections generated by DiffDRR. Here, we describe a deep learning method for tuning the CBCT densities to facilitate domain style transfer from synthetic DRR to real X-ray. We utilize a Multi-Layer Perceptron (MLP) consisting of fully connected layers to tune the CBCT densities sampled in DiffDRR, as illustrated in **Fig. 3**. Analogous to the positional encoding utilized in NeRF [7], we encode the CT density (σ) as a vector for input into the MLP:

$$E(\sigma) = [\sigma \quad \sin 2^i \sigma \quad \cos 2^i \sigma \quad \dots] \quad (5)$$

where a six-layer encoding is adopted, i.e., $0 \leq i \leq 5$. The MLP then outputs a new scalar density ($\hat{\sigma}$), which is used for the intensity calculation of each ray in DiffDRR (**Eq. 1 or 4**).

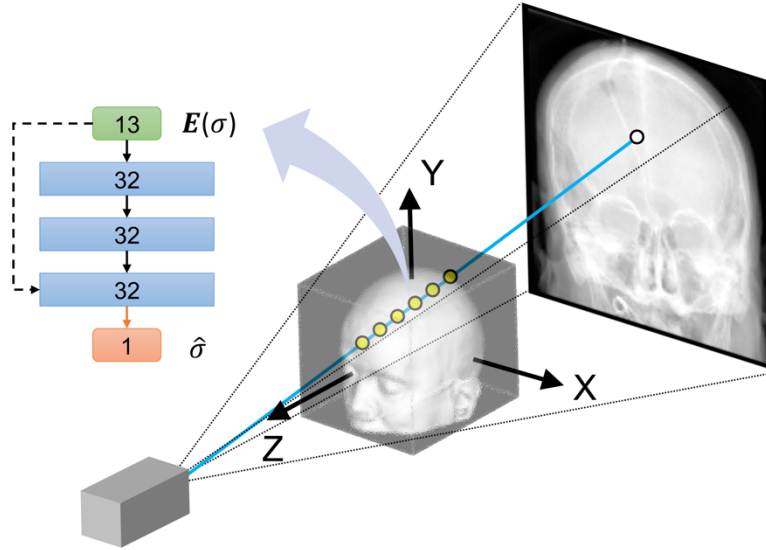


Fig. 3: Illustration of our NeTT optimization of DiffDRR. A sampled density (σ) along a ray is encoded and input into an MLP (flowchart) to output a tuned density ($\hat{\sigma}$), i.e., $\hat{\sigma} = \text{MLP}(\sigma|\mathbf{W})$. The green block represents the input density encoding, the blue blocks are hidden layers, and the orange block is the output density. The numbers in the blocks indicate the tensor sizes. The black arrows denote Dense + ReLU + Batch Normalization layers, the dashed arrow denotes a skip

connection (i.e., concatenation with the input density encoding), and the orange arrow denotes Dense + ReLU layers.

We trained the MLP using the 133 processed X-ray images of the CBCT tomographic sequence and the corresponding ground-truth poses, which were stored in the DICOM metadata (**Fig. 2b**). The X-ray images were downsized to 128×128 to match the size of the DRRs generated by DiffDRR. The SOD and SID were scaled to 7.8 and 768.0 in the NDC of DiffDRR to maintain the proper geometry, as previously described (*Section 2.2*). During MLP training, an X-ray image (I_{Xray}) and the corresponding ground-truth source pose (\mathbf{p}_t) were randomly sampled from the processed X-ray sequence. Optimization of the MLP weights (\mathbf{W}) can be formulated as:

$$\mathbf{W}_* = \underset{\mathbf{W}}{\operatorname{argmin}} \mathcal{L} \left(I_{DRR}(\mathbf{p}_t | \mathbf{W}), I_{Xray}(\mathbf{p}_t) \right) \quad (6)$$

where $I_{DRR}(\mathbf{p}_t | \mathbf{W})$ is the NeTT-optimized synthetic DRR generated by DiffDRR at the ground-truth source pose (\mathbf{p}_t) given the MLP weights (\mathbf{W}). The pose (\mathbf{p}_t) is one of the principal angles in an interval of 1.5° from -100° to 100° derived from the CBCT acquisition sequence (**Fig. 2a**). \mathcal{L} is a combined image loss function, quantifying the difference between a pair of I_{DRR} and I_{Xray} . It consists of a mean square error (MSE) loss, a focal frequency loss [16], and a structural similarity (SSIM) loss [17]. **Eq. 6** can be efficiently solved with gradient-based optimization, because the differentiability of DiffDRR enables computation of $\frac{\partial I_{DRR}}{\partial \mathbf{W}}$.

During training, we noted that feeding zero-densities (air) in the CT volume to the MLP probably resulted in DRR projections degenerating to a uniform intensity distribution (i.e., zero-intensity projection for VR and infinite-intensity projection for RC). To avoid this issue, we filtered the sampled CT densities, such that only non-zero CT densities are tuned by the MLP.

2.4. Masked Neural Radiance Field (mNeRF)

Differing from finite, discretized CBCT, CT, or MRI volumes, NeRF scene representations are unconstrained and continuous. As illustrated in **Fig. 4**, NeRF utilizes an MLP to map coordinates (input) to densities (output) [7]. To predict the density at a sampled position, the corresponding coordinates are provided to the MLP after positional encoding [7]:

$$E(x, y, z) = [x \ y \ z \ \sin 2^i x \ \sin 2^i y \ \sin 2^i z \ \cos 2^i x \ \cos 2^i y \ \cos 2^i z \ \dots] \quad (7)$$

where we adopt a six-layer encoding, i.e., $0 \leq i \leq 5$. As for synthesizing DRRs (**Eq. 1 or 4**), the densities ($\hat{\sigma}$) predicted by the NeRF MLP are sampled for each ray traced by the DiffDRR framework. The NeRF MLP is trained in a similar manner to NeTT, as described above (**Eq. 6**). Although an MSE loss was used alone to train NeRF in prior works with RGB images [7], we found that incorporation of the SSIM loss aids in convergence for our application.

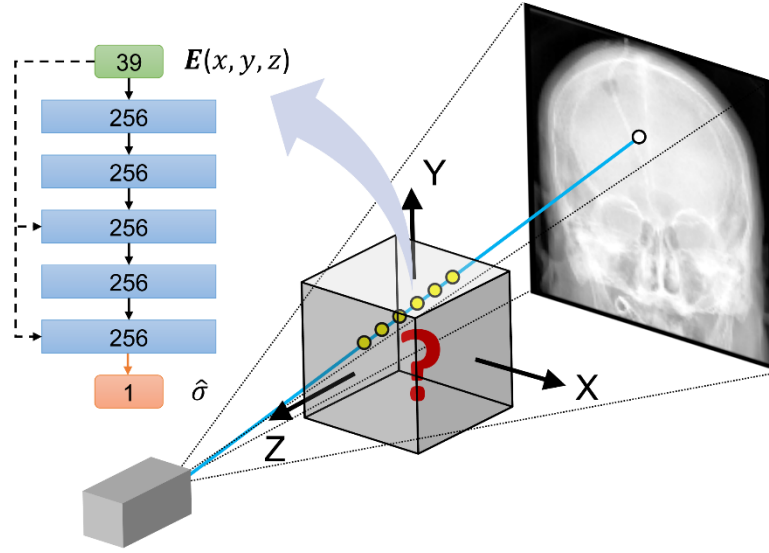


Fig. 4: Illustration of synthetic DRR generation with DiffDRR from a NeRF scene representation. The coordinates (x, y, z) at a sampled position are used as input to the NeRF MLP (shown by the flowchart) to predict the density ($\hat{\sigma}$) at the sampled position, i.e., $\hat{\sigma} = \text{MLP}(x, y, z \mid \mathbf{W})$. The green block represents the input positional encoding, the blue blocks are hidden layers, and the orange block is the output density. The numbers in the blocks indicate the tensor size. The black arrows denote Dense + ReLU + Batch Normalization layers, the dashed arrows denote skip connections (i.e., concatenation with the input positional encoding), and the orange arrow denotes Dense + ReLU layers.

Initial experiments demonstrated that the unconstrained NeRF generated artifacts outside of the CBCT reconstructed anatomic structure, as shown in **Fig. 5**. We observed that these artifacts can compensate for the style differences between DRRs and X-rays, but they also reflect

overfitting to our limited ground-truth X-ray sample, and do not yield plausible results when rendering new, out-of-sample poses. These artifacts may be exacerbated by our limited ground truth sample, which is captured over a single semi-circular tomographic pose trajectory, in contrast to the seminal NeRF work, in which the ground-truth samples covered a semi-sphere of the 3D scene [7]. To eliminate these artifacts, we introduced “masked NeRF” (mNeRF), in which the 3D region of the head, segmented from the CBCT reconstruction (**Fig. 2c**; processed using 3D Slicer [18] with a slight dilation of 3 mm applied), is used as a 3D mask to spatially constrain the non-zero values of NeRF. All densities outside of this mask are forced to zero. Our mNeRF masking module was implemented in the TensorFlow machine learning framework to ensure differentiability.

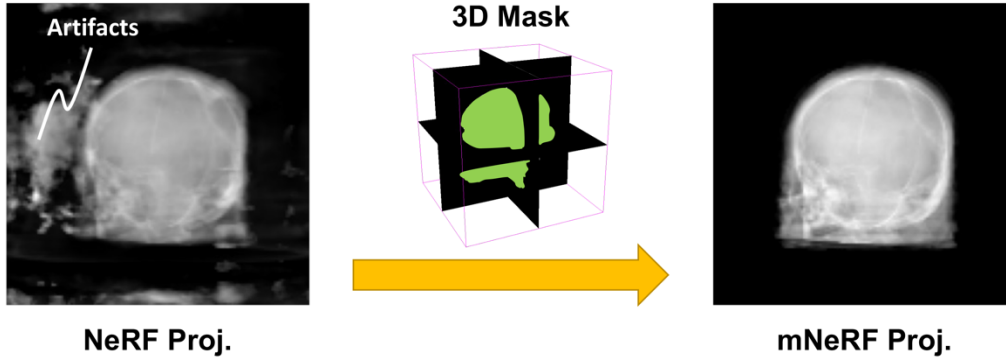


Fig. 5: Skull projections of NeRF and mNeRF scene representations after training, at an out-of-sample pose not seen during training. The application of a differentiable, spatially constrained mask results in substantial artifact reduction.

2.5. Pose Estimation by DRR/X-ray Image Registration

Intuitively, we can estimate the pose of a single X-ray projection relative to a 3D volume/field by iteratively nudging a randomly initialized pose to minimize the image difference between a rendered DRR and the target X-ray image. **Fig. 6** illustrates this process. In our model, the pose describes the location/orientation of the X-ray source and has six components, including three angular and three positional degrees of freedom (DOFs), and the anatomic volume is fixed. However, with a simple coordinate transformation, this is equivalent to pose estimation of a movable anatomic structure under a stationary X-ray source. See *Appendix B* for further details

on the equivalence of source vs. object movements, and on the coordinate transformation to move between these two configurations.

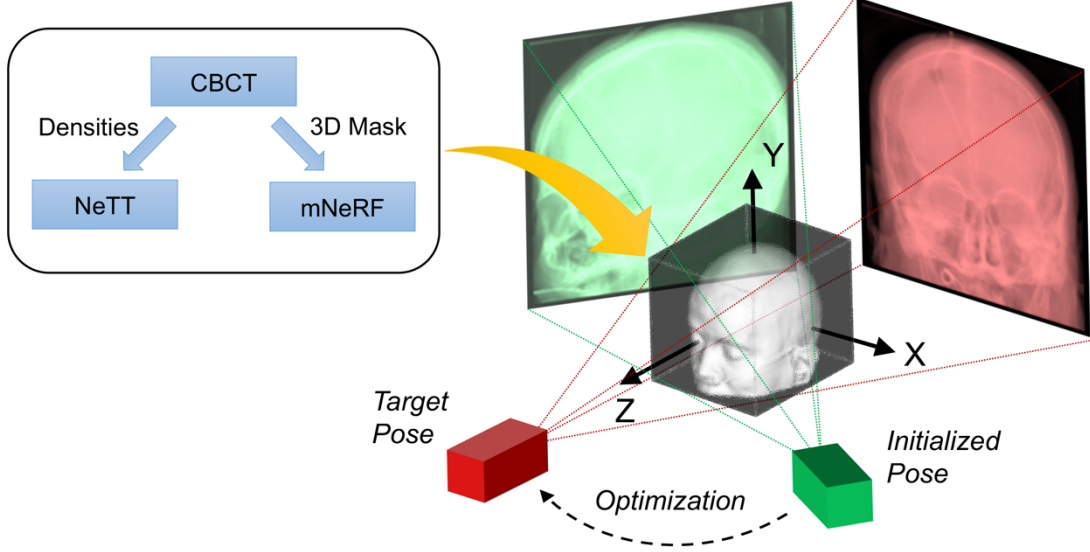


Fig. 6: Illustration of single-view X-ray pose estimation with DiffDRR. The target X-ray pose can be reached by iteratively nudging a randomly initialized DRR pose to minimize the difference between the X-ray and the rendered DRR. As a radiolucent scene representation for DiffDRR computation, CBCT, NeTT, and mNeRF can be used. Both innovative techniques of NeTT and mNeRF, that were introduced to substantially improve DRR fidelity and pose estimation, rely on classic CBCT reconstructions. As reflected in the *inset*, NeTT directly optimizes the CBCT densities, while the non-zero values of mNeRF are constrained by a 3D mask of the anatomic region segmented from CBCT.

Given an X-ray image (I_{Xray}) with an unknown pose, the estimation of a 6DoF pose (\mathbf{p}) can be described as an optimization problem:

$$\mathbf{p}_* = \underset{\mathbf{p}}{\operatorname{argmin}} \mathcal{L}(I_{DRR}(\mathbf{p}), I_{Xray}) \quad (8)$$

where $I_{DRR}(\mathbf{p})$ is the synthetic DRR generated by DiffDRR at a source pose (\mathbf{p}). $\mathbf{p} = [\boldsymbol{\theta}; \mathbf{U}]$ consists of three angular DoFs, $\boldsymbol{\theta} = [\theta_y \ \theta_x \ \theta_z]$ and three positional DoFs, $\mathbf{U} = [U_x \ U_y \ U_z]$. Here we define $\boldsymbol{\theta}$ with three *intrinsic* Euler angles applied in a “YXZ” sequence [19]. \mathcal{L} is the

image loss function used to quantify the difference between I_{DRR} and I_{Xray} . We adopted a mutual information (MI) loss [20, 21] for pose optimization. Because DiffDRR is differentiable, the gradient, $\frac{\partial I_{DRR}}{\partial \mathbf{p}}$ can be computed by backpropagation in machine learning frameworks, enabling us to efficiently solve **Eq. 8** by gradient descent.

We initialize the pose (\mathbf{p}_0) with randomly chosen angular and positional movements away from the target pose (\mathbf{p}_t):

$$\mathbf{p}_0 = \mathbf{p}_t + \Delta \mathbf{p} \quad (9)$$

where $\mathbf{p}_0 = [\boldsymbol{\theta}_0; \mathbf{U}_0]$ and $\mathbf{p}_t = [\boldsymbol{\theta}_t; \mathbf{U}_t]$. $\Delta \mathbf{p}$ is a vector consisting of six uniformly random numbers between $\min(\Delta \mathbf{p}) = [-30, -30, -30; -0.2, -0.2, -0.2]$ and $\max(\Delta \mathbf{p}) = [30, 30, 30; 0.2, 0.2, 0.2]$. The angular components are measured in degrees, and the positional components reflect a fractional distance within the full NDC space which has bounds of -1 and 1 in each dimension. This pose initialization range was chosen to simulate the starting points that physician operators would encounter during routine image-guided procedures, and the initialization range is large enough that the initial pose is most often distinctly different from the target pose, as demonstrated in **Fig. 7**. We normalized the optimization variables, \mathbf{p} , onto a regular variable space bounded by $\overline{\mathbf{lb}} = \mathbf{0}_{1 \times 6}$ and $\overline{\mathbf{ub}} = \mathbf{1}_{1 \times 6}$ for efficient optimization [22].

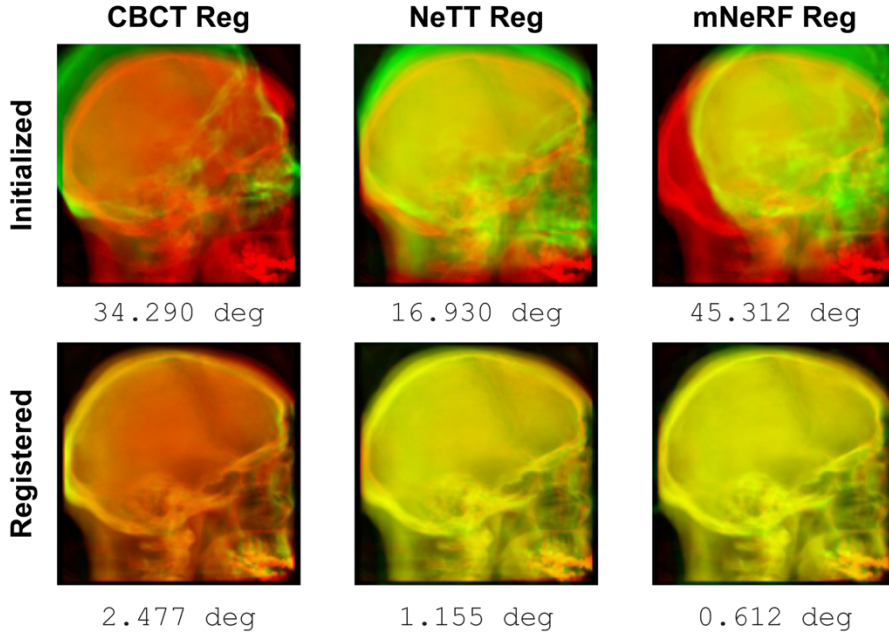


Fig. 7: Examples of successful registration of CBCT, NeTT, and mNeRF projections (*green*) onto a target X-ray image (*red*) using DiffDRR with the VR algorithm. Overlapping regions render in *yellow*. The 3D angle errors calculated between the DRR and X-ray pose estimates are listed below each plot. The *upper* row illustrates the random initial pose, and the *lower* row illustrates the optimal pose estimation.

We, and others have found that estimation of the angular pose parameters is much more challenging than estimating the positional pose parameters [23]. In our work, the angular and positional DoFs were coupled into a single optimization (**Eq. 8**). However, we often report the errors in angular DoFs, which we measure by introducing a 3D angle error (\mathcal{E}_θ). Given the angular DoFs, θ_* in the optimal pose and θ_t in the ground-truth pose, a 3D angle based on the axis-angle representation [24] can be calculated:

$$\mathcal{E}_\theta = \arccos \frac{\text{trace}(\mathbf{R}_* \cdot \mathbf{R}_t^T) - 1}{2} \quad (10)$$

where \mathbf{R}_* and $\mathbf{R}_t \in \mathbb{R}_{3 \times 3}$ are the rotation matrices calculated from the Euler angles θ_* and θ_t , respectively. The superscript T represents the matrix transpose. For this study, we define a successful image registration / pose estimation with $\mathcal{E}_\theta < 3^\circ$, which we find corresponds to good alignment on visual inspection (**Fig. 7**).

3. Experiments and Results

3.1. Reconstruction of CBCT, NeTT, and mNeRF

Tomographic X-ray sequences of the skulls of five patients were collected. Each tomographic X-ray sequence contains 133 X-ray images captured around the principal axis on an interval of 1.5° ranging from -100° to 100° . Using the X-ray sequences, CBCT, NeTT, and mNeRF were reconstructed for the skull of each patient. Prior to training of NeTT and mNeRF, reconstruction of each CBCT took about 5 min in TIGRE toolbox. For training both NeTT and mNeRF, the Adam optimizer with a learning rate of 5×10^{-4} was used. During each epoch, the 133 X-ray images and corresponding ground-truth poses from the tomographic X-ray sequence

were randomly shuffled. NeTT and mNeRF were trained for a variable number of epochs until a plateau was reached in the average MI metric calculated for the DRRs with respect to the X-ray images during an epoch. We used a stopping criterion with a patience of 10 epochs and a min delta of 1×10^{-5} .

Using DiffDRR with either VR or RC algorithms, mNeRF training was found to have a substantially larger computational cost in comparison to NeTT; the training time was less than 1 hour for NeTT and more than 4 hours for mNeRF (**Table 1**). The resulting DRRs generated by NeTT and mNeRF are similar, both by MI at training convergence (**Table 1**) and visual appearance (**Fig. 8**). DRRs generated from CBCTs with the RC algorithm are of substantially greater fidelity than those generated with VR (**Fig. 8**). This result is anticipated, as the Beer-Lambert Ray Casting algorithm more accurately reflects the underlying physics of X-ray imaging, and furthermore, the CBCT reconstruction algorithms in TIGRE also adopt the Beer-Lambert Law [25]. The suboptimal appearance of the DRR projections generated by VR of the CBCTs are primarily caused by the different physical assumptions from those in the RC algorithm. It is shown that such image domain misalignment can be compensated by either NeTT or mNeRF, thereby substantially improving the quality of the synthetic images for image registration tasks.

Table 1: Summary of the average computational cost for NeTT training, and mNeRF training for a single patient. All experiments were performed with a Tesla T4 GPU.

	NeTT		mNeRF	
	VR	RC	VR	RC
Epoch Time (sec)	78.0	74.8	169.5	166.5
Total # Epochs	45	25	90	101
Total Time (min)	58.9	31.2	253.8	281.1
Best Epoch MI	0.97	0.98	1.01	1.00

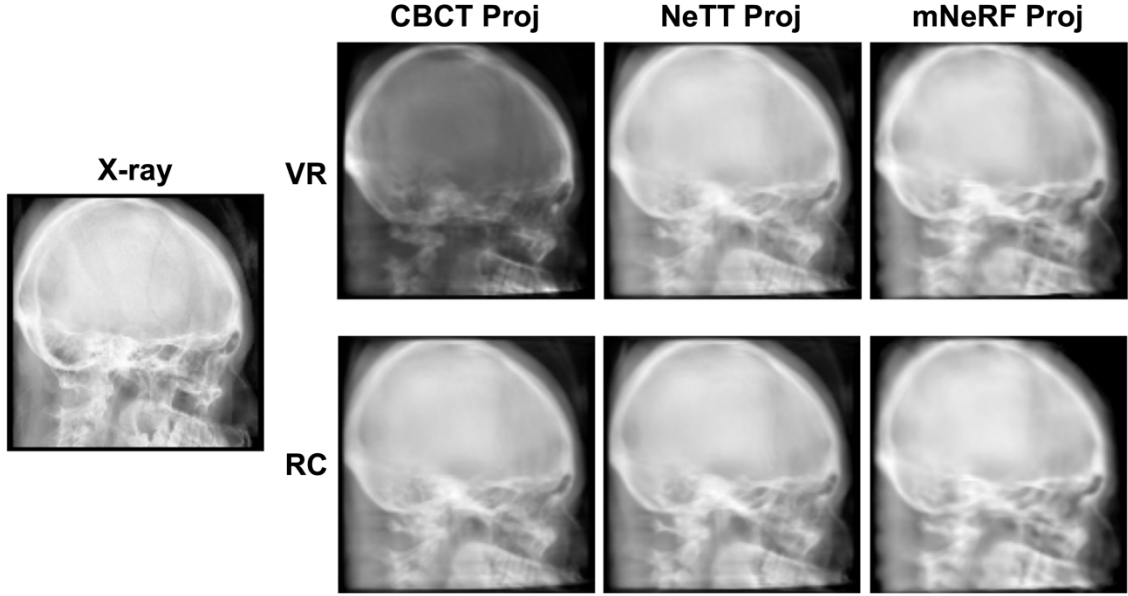


Fig. 8: Comparison of an X-ray image with DRRs generated at the same pose by projecting CBCT, NeTT, and mNeRF using volume rendering and ray casting, respectively. All images have a size of 128×128 and normalized intensities with a range of $[0, 1]$.

3.2. Pose Estimation Using CBCT, NeTT, and mNeRF

Five representative target X-ray images at poses near -90° , -45° , 0° , 45° , and 90° (**Fig. 9**) were chosen from each tomographic X-ray sequence. We estimated each target pose 20 times from a randomly initialized starting pose (**Eq. 9**) using CBCT, NeTT, and mNeRF. The Adam optimizer was used for pose estimation, with a learning rate of 0.03 and an exponential decay rate of 0.5 for the first moment estimates. Here, a large learning rate is utilized so that the gradient descent is encouraged to jump out of local optima. We terminated the pose optimization when a 50-iteration plateau was reached or at a maximum iteration number of 300. After termination, the pose with the best image loss during is considered to be the optimal solution. In all pose estimations, we adopt the MI loss as the image loss for pose optimization. Using the same optimization protocol, we discovered that the MI image loss is a superior objective function for pose estimation because it is well correlated with the pose error, and it encourages the optimization framework to be insensitive to the randomly initialized starting point (please see **Appendix C**). A summary of the computational cost for pose estimation was shown in **Table 2**. Pose estimation was nearly two times more computationally costly when using the mNeRF scene representation (~ 2 min for an

optimization), compared to the NeTT and CBCT scene representations (~1-1.5 min for an optimization).

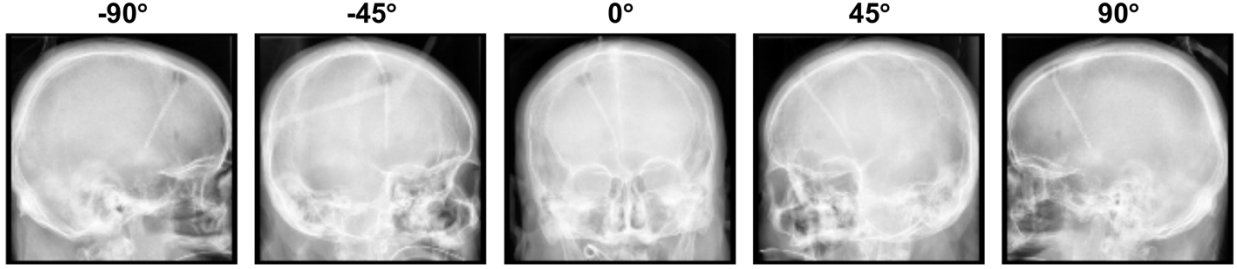


Fig. 9: Selected X-ray images at five target poses used to evaluate the performance of pose estimation using our gradient-based optimization framework.

Table 2: Summary of the average computational cost of pose estimation for a single X-ray image using CBCT, NeTT, and mNeRF. All experiments were performed with a Tesla T4 GPU.

	CBCT		NeTT		mNeRF	
	VR	RC	VR	RC	VR	RC
Iteration Time (sec)	0.50	0.51	0.60	0.68	1.01	0.97
Total # Iterations	145	121	130	126	122	132
Total Time (sec)	72.2	61.7	78.5	85.7	123.1	128.4

Given that our pose estimation framework relies upon DRR/X-ray image registration, we hypothesized that rendering methods that generate more realistic, synthetic DRRs would facilitate this process. With this in mind, we performed pose estimation with DRRs rendered from CBCT, NeTT, and mNeRF. For each method of scene representation, **Fig. 10** shows the optimal 3D angle errors for all 100 random initializations in the estimation of each target pose on all 5 patients using the MI loss. Using VR to generate DRRs (**Fig. 10a**), NeTT and mNeRF scene representations result in overall success rates of 93.2% and 95.4%, respectively, which are a substantial improvement over the success rate of 71.8% achieved using DRRs rendered directly from CBCT. Using RC to project NeTT and mNeRF scene representations (**Fig. 10b**), similar overall success rates of 93.0% and 93.4%, respectively, are achieved. However, the overall success rate for direct CBCT scene representation increases to 88.2% when using RC, since the appearance of CBCT-projected DRRs closely matches that of the X-ray images, as demonstrated in **Fig. 8**.

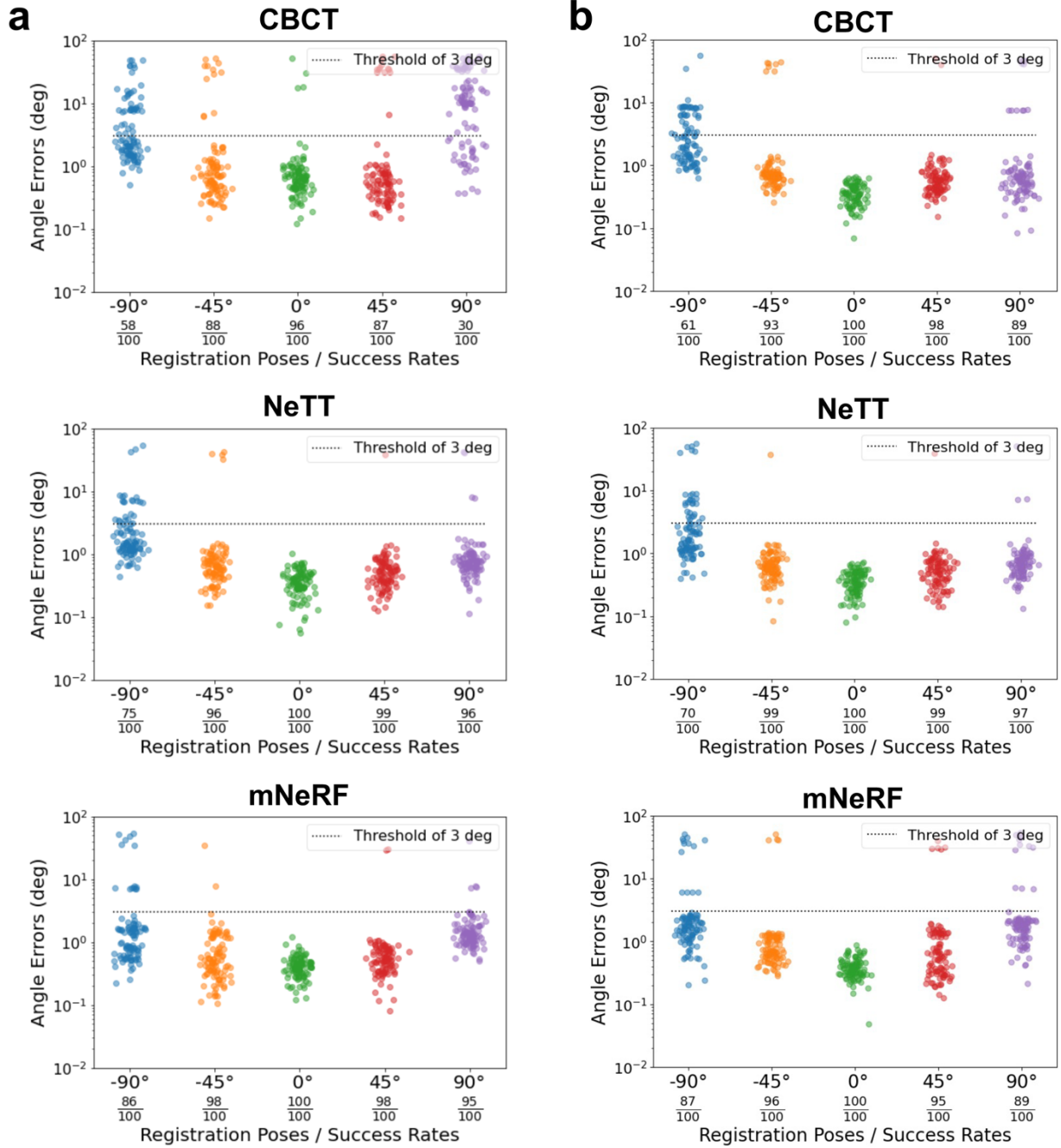


Fig. 10: Comparison of the 3D angle errors and success rates using volume rendering (a) and ray casting (b), respectively, to project CBCT, NeTT, and mNeRF scene representations for all pooled trials of all five patients. A 3D angle error less than 3° is considered to be successful.

3.3. Cross-subject Generalizability of NeTT

We further examine the generalizability of NeTT for improving DRR fidelity and pose estimations on new subjects without additional re-training. Using a CBCT for scene representation and the Volume Rendering equation for DRR view synthesis, a NeTT MLP was trained on the series of 133 X-ray images of patient #1. The trained NeTT MLP of patient #1 was then utilized for DRR view synthesis from the CBCTs of the other 4 patients (#2 ~ #5). As demonstrated in **Fig. 11**, the NeTT of a single subject can promote domain alignment of DRRs with X-ray images for all of other subjects. Furthermore, it is also shown that the performance of pose estimation is substantially improved to an overall success rate of 95.2% (**Fig. 12**), in comparison with the overall success rate of 93.2% using NeTTs trained individually (**Fig. 10a**). NeTT therefore demonstrates excellent generalizability across subjects within the X-ray modality. In contrast, mNeRF must be trained for each individual subject.

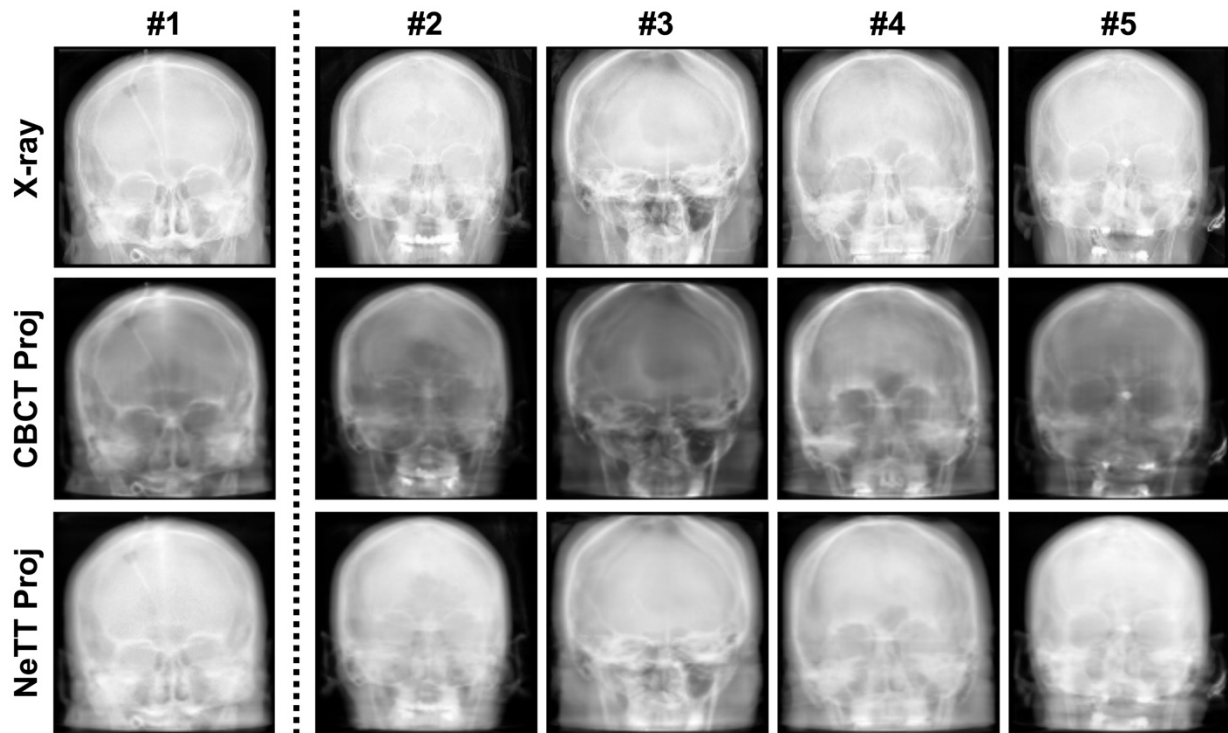


Fig. 11: Generalizability of NeTT in tuning of CBCT densities for generating DRRs using volume rendering to match the appearance of X-ray images. All images have a size of 128×128 and normalized intensities with a range of $[0, 1]$. Note that a single NeTT MLP was trained only using the series of 133 X-ray images of patient #1. The trained NeTT MLP was used to tune CBCTs of other patients (#2 ~ #5) for generating fidelitous DRRs.

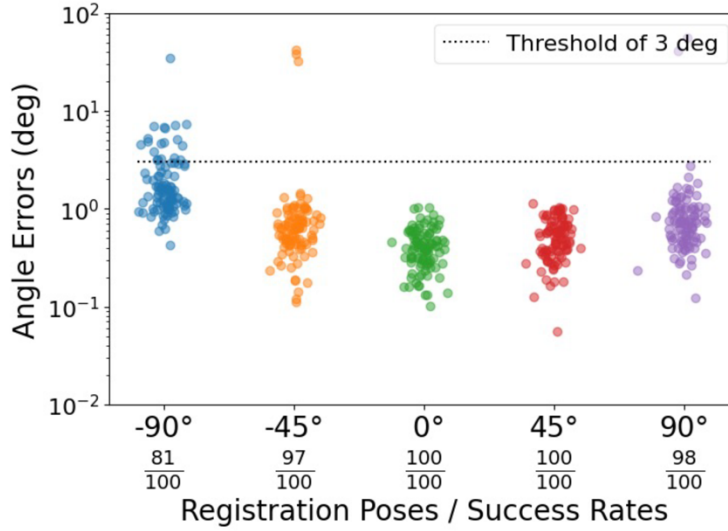


Fig. 12: Generalizability of NeTT in the improvement of success rates in pose estimation when using volume rendering. A 3D angle error less than 3° is considered to be successful. For a comparison purpose, the pose estimation results of patient #1 have been pooled together with those of other patients (#2 ~ #5). All the five CBCTs were tuned using the same NeTT MLP trained only using the series of 133 X-ray images of patient #1.

4. Discussion and Conclusion

Due to different physical assumptions, rendering algorithms, such as volume rendering and ray casting, can significantly influence the appearance and quality of synthetic radiographs. In this work, we introduced NeTT and mNeRF, and show that both techniques can effectively eliminate the discrepancy between different image domains to ensure robust DRR/X-ray image registration for pose estimation. However, the computational cost of NeTT is significantly lower than mNeRF in both training (**Table 1**) and pose estimation (**Table 2**). And furthermore, we have demonstrated that NeTT has excellent cross-subject generalizability for DRR synthesis (**Fig. 11**) and pose estimation (**Fig. 12**), eliminating the need for NeTT to be trained on each individual subject, as required by mNeRF.

Our image dataset consisted of five tomographic X-ray series from CBCT acquisitions, as these data are routinely collected in medical practice. However, a limitation of this dataset arises

from the fact that all of the ground truth tomographic images are collected on a single semi-circular arc around the patient. Consequently, real X-ray images from poses outside of this arc were not available for evaluation in this study. The five poses selected from within the CBCT tomographic series for this study do represent common projections utilized during neurointerventional procedures (AnteroPosterior, Lateral, and Oblique views), and large angles in the CranioCaudal or CaudoCranial directions are less commonly used in routine clinical practice. In addition, we performed CBCT reconstruction ($256 \times 256 \times 256$) and view synthesis (128×128) using resolutions that are below what is typically encountered in clinical scenarios. For example, our fluoroscopic images often have a size of 1024×1024 . We down-sampled images in this study due to the memory limitations of our hardware. Future optimizations may enable view synthesis and pose optimization with full resolution images, which we anticipate will improve the pose estimates.

In conclusion, we have introduced methods for pose estimation of radiolucent objects using 2D projections. We first developed a differentiable framework of DiffDRR for efficient computation of DRRs with automatic differentiation. In conjunction with classic CBCT 3D volume reconstruction algorithms, we perform pose estimation by iterative gradient descent using loss functions that quantify the similarity of the DRR synthesized from a randomly initialized pose and the true fluoroscopic image of the target pose. Next, we proposed two novel methods for high fidelity view synthesis, Neural Tuned Tomography (NeTT) and masked Neural Radiance Fields (mNeRF), and we show that both of these techniques improve pose estimation within our framework. We find that NeTT and mNeRF can achieve similar results, with overall success rates more than 93%, regardless of adopting volume rendering or ray casting for DRR synthesis. But given the much lower computational cost for NeTT in training and pose optimization as well as the cross-subject generalizability of NeTT, we suggest that NeTT is an attractive option for robust pose estimation using fluoroscopic projections.

Competing Interests: Portions of the work described in this article have been included in a related patent filed by Northwestern University, with C. Zhou, D.R. Cantrell, L. Cho, and S.A. Ansari listed as co-inventors. D.R. Cantrell, L. Cho, and S.A. Ansari are founders and have shares in Clearvoya, LLC, which aims to commercialize Computer Vision algorithms for Image-Guided

interventions, but this work was not funded or performed by Clearvoya. M.C. Hurley made contributions to this work while he was at Northwestern University, but he did not contribute from his new institution.

Ethical Statement: Patients' cone-beam CTs routinely collected during interventional treatments of neurological disorders were adopted in this study. The IRB number associated with the project is STU00212923, which was approved at Northwestern Medicine.

Acknowledgement: None

References

1. Molyneux AJ, C Kerr RS, Yu L-M, et al (2005) Articles Introduction International subarachnoid aneurysm trial (ISAT) of neurosurgical clipping versus endovascular coiling in 2143 patients with ruptured intracranial aneurysms: a randomised comparison of effects on survival, dependency, seizures, rebleeding, subgroups, and aneurysm occlusion
2. Goyal M, Menon BK, Van Zwam WH, et al (2016) Endovascular thrombectomy after large-vessel ischaemic stroke: A meta-analysis of individual patient data from five randomised trials. *The Lancet* 387:1723–1731. [https://doi.org/10.1016/S0140-6736\(16\)00163-X](https://doi.org/10.1016/S0140-6736(16)00163-X)
3. Unberath M, Gao C, Hu Y, et al (2021) The Impact of Machine Learning on 2D/3D Registration for Image-Guided Interventions: A Systematic Review and Perspective. *Front Robot AI* 8:. <https://doi.org/10.3389/frobt.2021.716007>
4. Miao S, Wang ZJ, Zheng Y, Liao R (2016) Real-time 2D/3D registration via CNN regression. In: 2016 IEEE 13th International Symposium on Biomedical Imaging (ISBI). IEEE, pp 1430–1434
5. Liao H, Lin W-A, Zhang J, et al (2019) Multiview 2D/3D Rigid Registration via a Point-Of-Interest Network for Tracking and Triangulation. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE, pp 12630–12639
6. Zhou C, Cha T, Peng Y, Li G (2021) Transfer learning from an artificial radiograph-landmark dataset for registration of the anatomic skull model to dual fluoroscopic X-ray images. *Comput Biol Med* 138:104923. <https://doi.org/10.1016/j.combiomed.2021.104923>
7. Mildenhall B, Srinivasan PP, Tancik M, et al (2022) NeRF: representing scenes as neural radiance fields for view synthesis. *Commun ACM* 65:99–106. <https://doi.org/10.1145/3503250>
8. Martin-Brualla R, Radwan N, Sajjadi MSM, et al (2021) NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* 7210–7219
9. Yu A, Ye V, Tancik M, Kanazawa A (2021) pixelNeRF: Neural Radiance Fields from One or Few Images. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* 4578–4587

10. Yen-Chen L, Florence P, Barron JT, et al (2021) iNeRF: Inverting Neural Radiance Fields for Pose Estimation. In: 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). IEEE, pp 1323–1330
11. Albarqouni S, Fotouhi J, Navab N (2017) X-ray in-depth decomposition: Revealing the latent structures. Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics) 10435 LNCS:444–452. https://doi.org/10.1007/978-3-319-66179-7_51/FIGURES/5
12. Corona-Figueroa A, Frawley J, Taylor SB-, et al (2022) 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). IEEE, pp 3843–3848
13. Tagliasacchi A, Mildenhall B (2022) Volume Rendering Digest (for NeRF)
14. Rückert D, Wang Y, Li R, et al (2022) NeAT: Neural adaptive tomography. ACM Transactions on Graphics (TOG) 41:. <https://doi.org/10.1145/3528223.3530121>
15. Biguri A, Lindroos R, Bryll R, et al (2020) Arbitrarily large tomography with iterative algorithms on multiple GPUs using the TIGRE toolbox. J Parallel Distrib Comput 146:52–63. <https://doi.org/10.1016/j.jpdc.2020.07.004>
16. Tu Z, Talebi H, Zhang H, et al (2022) MAXIM: Multi-Axis MLP for Image Processing. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) 5769–5780. <https://doi.org/https://doi.org/10.48550/arXiv.2201.02973>
17. Wang Z, Bovik AC, Sheikh HR, Simoncelli EP (2004) Image Quality Assessment: From Error Visibility to Structural Similarity. IEEE Transactions on Image Processing 13:600–612. <https://doi.org/10.1109/TIP.2003.819861>
18. Fedorov A, Beichel R, Kalpathy-Cramer J, et al (2012) 3D Slicer as an image computing platform for the Quantitative Imaging Network. Magn Reson Imaging 30:1323–41. <https://doi.org/10.1016/j.mri.2012.05.001>
19. Zhou C, Cha T, Wang W, et al (2021) Investigation of Alterations in the Lumbar Disc Biomechanics at the Adjacent Segments After Spinal Fusion Using a Combined In Vivo and In Silico Approach. Ann Biomed Eng 49:601–616. <https://doi.org/10.1007/s10439-020-02588-9>

20. Dalca A v, Guttag J, Sabuncu MR (2018) Anatomical Priors in Convolutional Networks for Unsupervised Biomedical Segmentation. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 9290–9299. <https://doi.org/https://doi.org/10.48550/arXiv.1903.03148>
21. Viola P, Wells III WM (1997) Alignment by Maximization of Mutual Information. *Int J Comput Vis* 24:137–154. <https://doi.org/10.1023/A:1007958904918>
22. Zhou C, Willing R (2020) Multiobjective Design Optimization of a Biconcave Mobile-Bearing Lumbar Total Artificial Disk Considering Spinal Kinematics, Facet Joint Loading, and Metal-on-Polyethylene Contact Mechanics. *J Biomech Eng* 142:041006. <https://doi.org/10.1115/1.4045048>
23. Hanley J, Mageras GS, Sun J, Kutcher GJ (1995) The effects of out-of-plane rotations on two dimensional portal image registration in conformal radiotherapy of the prostate. *Int J Radiat Oncol Biol Phys* 33:1331–43. [https://doi.org/10.1016/0360-3016\(95\)02062-4](https://doi.org/10.1016/0360-3016(95)02062-4)
24. Mahendran S, Ali H, Vidal R (2017) 3D Pose Regression using Convolutional Neural Networks. Proceedings of the IEEE International Conference on Computer Vision (ICCV) 2174–2182
25. Biguri A, Dosanjh M, Hancock S, Soleimani M (2016) TIGRE: a MATLAB-GPU toolbox for CBCT image reconstruction. *Biomed Phys Eng Express* 2:055010. <https://doi.org/10.1088/2057-1976/2/5/055010>

Supplementary Material

Appendix A. X-ray Configuration Parameters

In a typical X-ray configuration (e.g., C-arm fluoroscopy), the X-ray source moves relative to a coordinate system (CS) established at the isocenter (i.e., the rotational center of the X-ray), as illustrated **Fig. S1**. The geometry imaged on the intensifier depends on two parameters, the source-isocenter distance (SOD) and the source-intensifier distance (SID; also called the focal distance for a camera). The SOD and SID can be set according to the volume and image sizes, respectively.

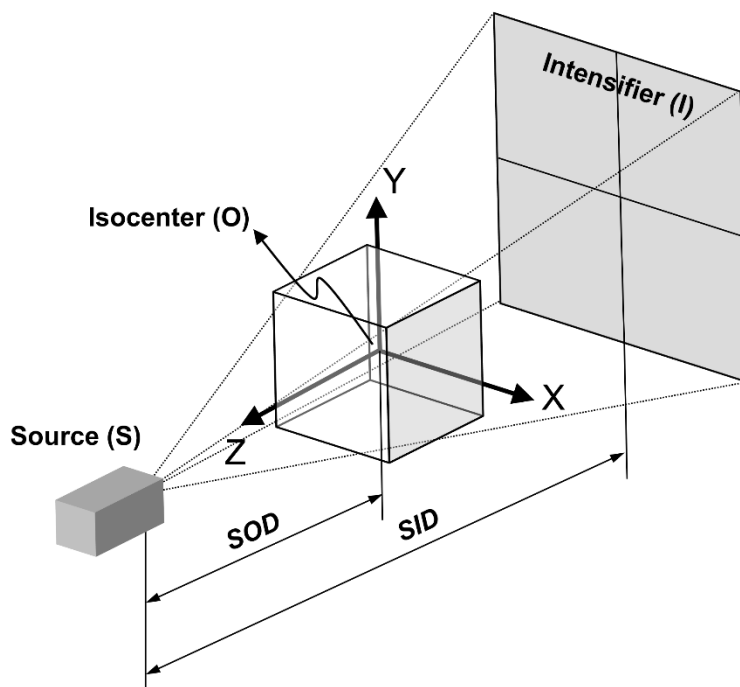


Fig. S1: The typical configuration of a movable X-ray imaging system. The imaged geometry on the intensifier is determined by two parameters, the SOD and SID.

For a cubic volume that is scaled to different sizes, l_1 and l_2 (e.g., those in the physical coordinate and NDC, respectively), the corresponding SODs, z_1 and z_2 , should be set, such that the same geometry is imaged:

$$\frac{z_1}{z_2} = \frac{l_1}{l_2} \quad (\text{S1})$$

Furthermore, the volume may be projected onto a square intensifier using different image sizes, s_1 and s_2 . To ensure that the imaged geometry remains invariant, the corresponding SIDs, f_1 and f_2 , should satisfy:

$$\frac{f_1}{f_2} = \frac{s_1}{s_2} \quad (\text{S2})$$

Appendix B. Source vs. Object Movements

In a movable X-ray configuration, the X-ray source can be moved *first* by translations and *then* by rotations with respect to a coordinate system (CS) established at the isocenter, as shown in **Fig. S2**. In particular, the angular DoFs are represented by three *intrinsic* Euler angles with a sequence of “YXZ”. Sometimes, it may also be desirable to consider a moving object relative to a stationary X-ray source (e.g., when aligning a patient’s CT to a patient’s intraoperative pose in physical 3D space), as shown in **Fig. S3**.

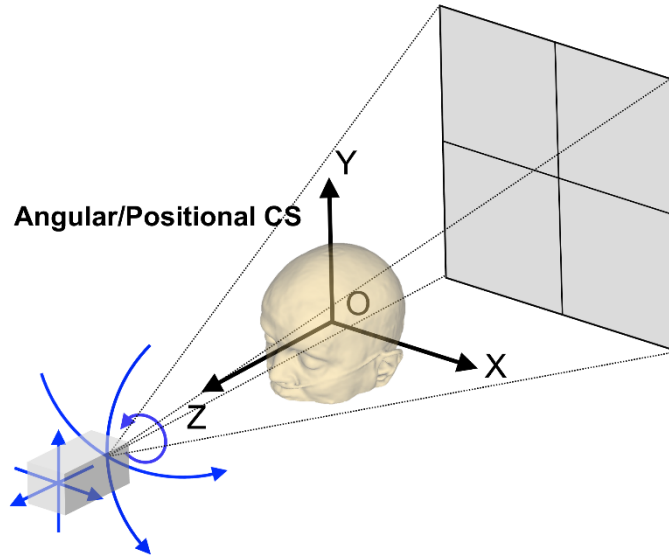


Fig. S2: The configuration of a movable X-ray, where a source moves relative to a fixed object. The positional/angular CS is set at the isocenter. The *blue* arrows indicate the positive directions of 6DoFs.

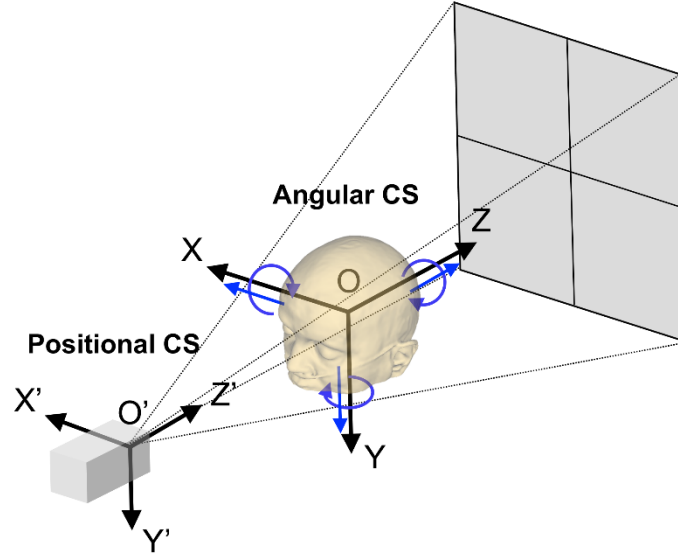


Fig. S3: The configuration of a fixed X-ray, where an object moves relative to a fixed source. The positional CS origin is set at the source, while the angular CS origin is set at the isocenter. The *blue* arrows indicate the positive directions of 6DoFs.

The 6DoF motion (*first* translations and *then* rotations) in the *movable* X-ray configuration is convertible to motion in the *fixed* X-ray configuration (*first* rotations and *then* translations). To ensure the same magnitude of the 6 parameters in both X-ray configurations, under the assumption of a fixed X-ray configuration (**Fig. S3**), the positional CS origin must be set at the X-ray source, and the angular CS origin must be set at the isocenter. In the *fixed* X-ray configuration, both the CSs have positive coordinate directions opposite to those in the *movable* X-ray configuration (i.e., the two CSs in the *fixed* X-ray configuration are left-handed). In the *fixed* X-ray configuration, the angular DoFs of the object are defined by three *extrinsic* Euler angles with the same sequence of “yxz” with respect to the angular CS, while the positional DoFs of the object are measured with respect to the positional CS.

Appendix C. Effects of Image Losses on Pose Estimation

In our pose optimization framework, the ultimate goal is to estimate the ground-truth pose, but to do so, we can only utilize an image loss as the optimization objective function to indirectly quantify the pose error. The pose error itself is unavailable to the algorithm during both training and inference. For this purpose, an optimal image loss should be highly correlated to the pose error. Here, we compare the performance of a mutual information (MI) loss (\mathcal{L}_{MI}) with a combined loss (\mathcal{L}_C) for pose estimation with DRRs rendered directly from CBCT. The combined loss for pose estimation, different from that used for NeTT and mNeRF optimization, is defined as a linear combination of an SSIM loss, a soft dice loss, and an L1 loss.

The performance of our pose estimation framework with DRRs directly rendered from CBCT using the volume rendering algorithm is shown in **Fig. S4**, when using the combined loss versus the MI loss. With the combined loss, the overall success rate was only 11%, with all 20 trials failing at the -90° , -45° , and 45° poses. There was substantial improvement when using the MI loss, which had an overall success rate of 73%, successful registrations at all poses, and success in all 20 iterations of the 0° target pose. As shown in **Fig. S5**, there were only mild correlations between the optimized 3D angle errors and the initialized 3D angle errors for both loss functions, indicating that the pose optimization framework is robust to the randomly initialized starting point. Furthermore, optimal 3D angle errors were much more highly correlated to the optimal MI image loss than the optimal combined image loss, as shown in **Fig. S6**, which is strong evidence for the superiority of the MI loss over the combined loss, as it more accurately reflects the optimal 3D angle errors.

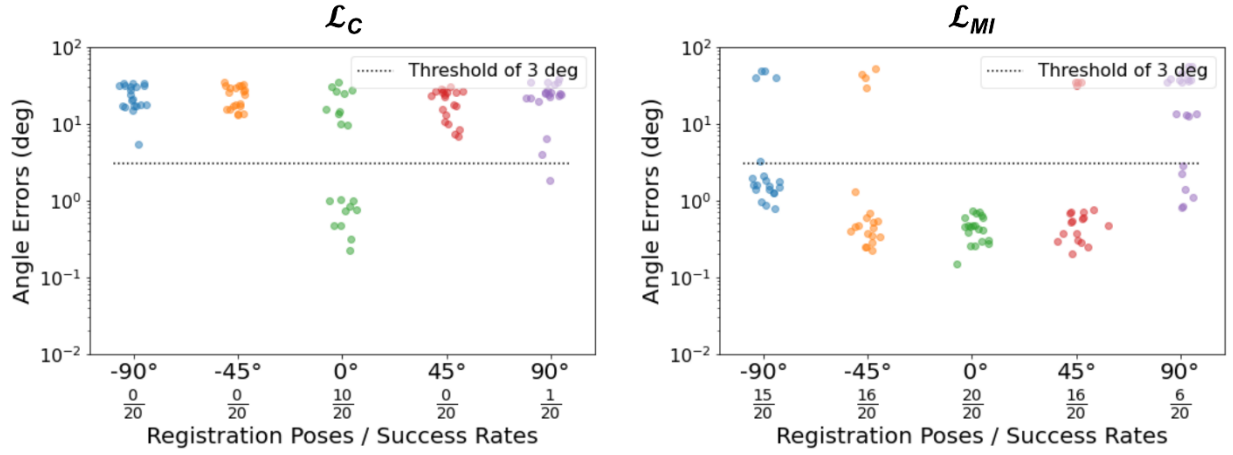


Fig. S4: Comparison of the 3D angle errors and success rates for pose estimation of five different X-ray images from a single patient using DRRs directly rendered from CBCT and the combined loss (*left*) compared to the MI loss (*right*). A 3D angle error less than 3° is considered to be successful.

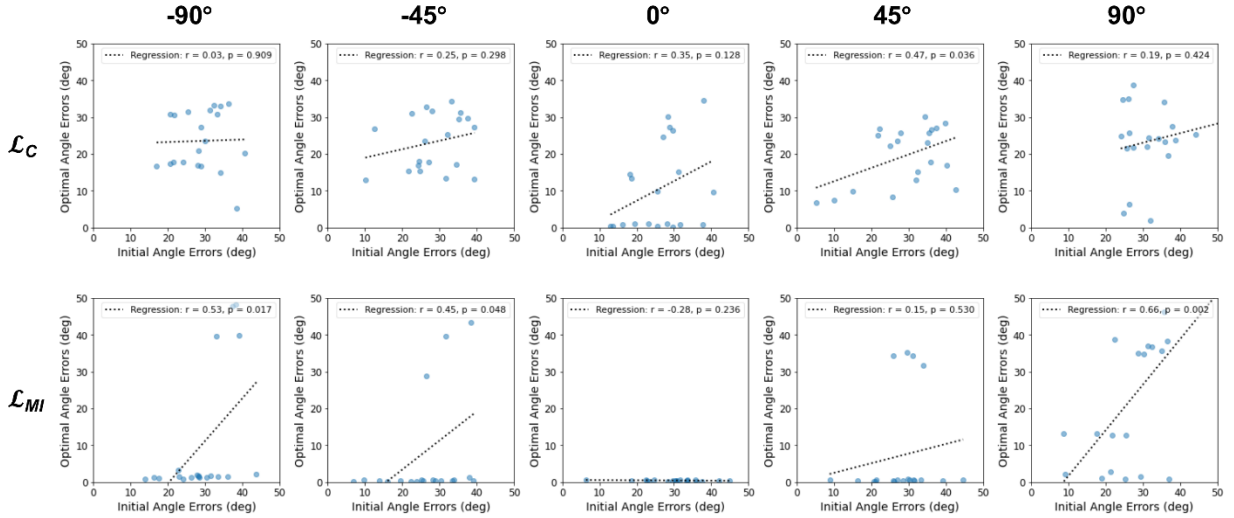


Fig. S5: Correlation between optimal 3D angle errors and initial 3D angle errors using DRRs directly rendered from CBCT with the combined loss (\mathcal{L}_C) and the MI loss (\mathcal{L}_{MI}).

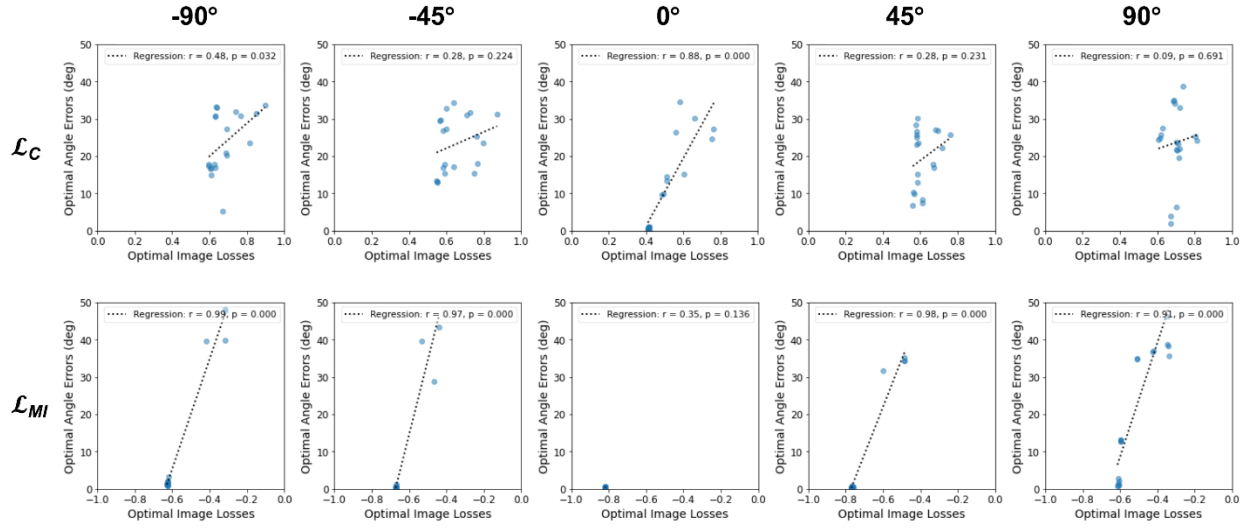


Fig. S6: Correlation between optimal 3D angle errors and optimal image losses (i.e., \mathcal{L}_C and \mathcal{L}_{MI} , respectively) using DRRs directly rendered from CBCT.