# Neural Annotation Refinement: Development of a New 3D Dataset for Adrenal Gland Analysis

Jiancheng Yang[1,2,⋆], Rui Shi[1,⋆], Udaranga Wickramasinghe[2],
Qikui Zhu[3,4], Bingbing Ni[1⋆⋆], and Pascal Fua[2]

[1] Shanghai Jiao Tong University, Shanghai, China
`nibingbing@sjtu.edu.cn`
[2] EPFL, Lausanne, Switzerland
[3] Dept. of Computer and Data Science, Case Western Reserve University, OH, USA
[4] Dept. of Biomedical Engineering, Case Western Reserve University, OH, USA

**Abstract.** The human annotations are imperfect, especially when produced by junior practitioners. Multi-expert consensus is usually regarded as golden standard, while this annotation protocol is too expensive to implement in many real-world projects. In this study, we propose a method to refine human annotation, named *Neural Annotation Refinement (NeAR)*. It is based on a learnable implicit function, which decodes a latent vector into represented shape. By integrating the appearance as an input of implicit functions, the appearance-aware NeAR fixes the annotation artefacts. Our method is demonstrated on the application of adrenal gland analysis. We first show that the NeAR can repair distorted golden standards on a public adrenal gland segmentation dataset. Besides, we develop a new Adrenal gLand ANalysis (ALAN) dataset with the proposed NeAR, where each case consists of a 3D shape of adrenal gland and its diagnosis label (normal vs. abnormal) assigned by experts. We show that models trained on the shapes repaired by the NeAR can diagnose adrenal glands better than the original ones. The ALAN dataset will be open-source, with 1,584 shapes for adrenal gland diagnosis, which serves as a new benchmark for medical shape analysis. Code and dataset are available at https://github.com/M3DV/NeAR.

**Keywords:** neural annotation refinement · adrenal gland · ALAN dataset · geometric deep learning · shape analysis.

## 1 Introduction

Deep learning has enjoyed a great success in medical image analysis, but large annotated datasets are required to achieve this [7,6,12,1,5]. Unfortunately, such datasets are difficult to obtain in part because human annotations are known to be imperfect [11,24]. In medical image segmentation, multi-expert consensus is employed as golden standard, where the agreement of multiple annotators is

---

⋆ These authors have contributed equally: Jiancheng Yang and Rui Shi.
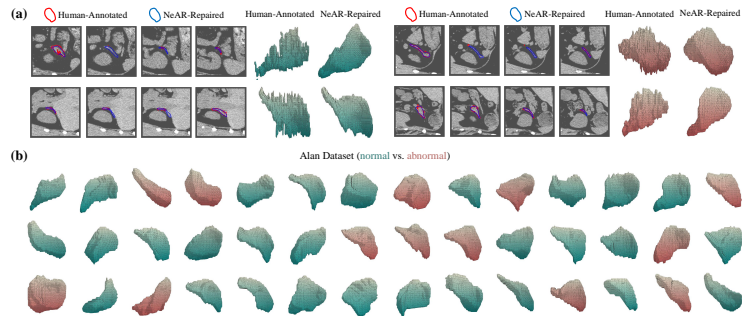⋆⋆ Corresponding author: Bingbing Ni (nibingbing@sjtu.edu.cn).

Fig. 1: **The ALAN Dataset.** It features 1,584 adrenal glands, each one of which has been tagged by human experts as normal or abnormal. The normal ones are shown in green and the abnormal in red. Best viewed on screen. **(a)** Images with human and NeAR-repaired annotations, in red and blue contours respectively, and the corresponding 3D visualization. **(b)** The repaired ALAN Dataset.

regarded as ground truth. Nevertheless, this protocol involving multiple medical experts is often too time-consuming and expensive to achieve in practice. In those cases, there are often high-frequency artefacts and false positive/negative in human segmentation. Please refer to Fig. 1 (a) for illustration.

In this paper, we introduce *Neural Annotation Refinement (NeAR)*, an approach to automatically correcting human-annotated segmentation databases, so that networks trained using the corrected database perform better than those trained using the original one. Our method is developed based on the fact that a neural network with appropriate inductive bias could serve as a deep prior [25,8]. By leveraging the recent advance in implicit surface modeling [2,16,20] that uses a neural network (*e.g.*, MLP and CNN [22]) as a mapping from spatial coordinates to a shape representation, the NeAR learns data-efficient implicit functions as a shape prior of the target annotations, which can be used to repair human annotations. To make the repaired segmentation appearance-aware, we integrate the appearance as an input of the implicit function. As illustrated in Fig. 1 (a), the repaired segmentation by the proposed NeAR is visually appealing. We will further show that the repairing can be used to improve downstream applications.

Our method is demonstrated on the application of adrenal gland analysis. We first show that the NeAR can repair distorted golden standards on a public adrenal gland segmentation dataset, consisting of 100 cases. The NeAR outperforms standard segmentation methods quantitatively in terms of the repaired annotation quality. Furthermore, we apply the NeAR to repair a new Adrenal gLand ANalysis (ALAN) dataset of 1,584 cases, where each adrenal gland is segmented by 1 clinician and diagnosed—as normal or abnormal—by 2 clinicians and 1 senior endocrinologist. In other words, the diagnosis label is quite reliable whereas the segmentation exhibits problems as shown in Fig. 1 (a). As shown in Fig. 1 (b), NeAR can effectively repair these segmentations, as evidenced by

the fact that models trained on the shapes repaired by the NeAR can better diagnose adrenal glands (normal vs. abnormal) than the original ones.

As an independent contribution, the ALAN dataset will be open-source, with NeAR-repaired shapes of adrenal glands and the corresponding diagnosis labels (normal vs. abnormal), as illustrated in Fig. 1 (b). This shape classification benchmark with 1,584 high-quality 3D models will be of interest for medical image analysis and geometric deep learning research community.[1]

## 2 Method

In this section, we first briefly review deep implicit surface, an emerging technique in 3D vision. We then introduce how this technique can be applied to repair human annotated segmentation labels, and propose the Neural Annotation Refinement (NeAR) based on appearance-aware implicit surface model.

### 2.1 Deep Implicit Surfaces

Implicit surface modeling [2,16,20] maps spatial coordinates to shape representations with a neural network. Typically, the shape representation could be either binary occupancy or signed / unsigned distance. For simplicity, we use occupancy fields [16] in this study, while the whole framework can be easily applied on distance functions [20]. In implicit surface modeling, a 3D shape is first encoded with a $c$-dimensional latent vector $\mathbf{z} \in \mathbb{R}^c$, and a continuous representation of the shape is then obtained by learning a mapping:

$$\mathcal{F}(\mathbf{z}, \mathbf{p}) = o : \mathbb{R}^c \times \mathbb{R}^3 \to [0, 1]. \tag{1}$$

Here, a $c$-dimensional latent vector $\mathbf{z} \in \mathbb{R}^c$ and coordinates of a query point $\mathbf{p} \in \mathbb{R}^3$ are inputted into a neural network $\mathcal{F}$–typically multi-layer perceptron (MLP)–to classify whether the query point is inside or outside the represented shape, with the occupancy probability $o$ close to 1 for $\mathbf{p}$ inside the shape and 0 otherwise. With a thresholding parameter $t$, the underlying surface is implicitly represented by the decision boundary $\mathcal{F}(\mathbf{z}, \mathbf{p}) = t$. For model training, we apply auto-decoding [20], an encoder-free approach where a learnable latent vector $\mathbf{z}$ of each shape is directly taken as input, jointly optimized with the parameters of $\mathcal{F}$ through back-propagation. The number of latent vectors is equal to the number of training shape samples.

The deep implicit models have achieved a great success in a wide range of applications, *e.g.*, shape modeling [2,16,20,10], 3D reconstruction [3,26] and differentiable rendering [17,18]. The deep implicit surface serves as a deep prior [25,8] for shape modeling, thus can be used as a tool to refine annotations, as the implicit reconstructions tend to remove high-frequency artefacts introduced by human annotators. However, standard implicit surface methods are not aware of the appearance, thus the reconstructed surfaces could be misaligned with

---

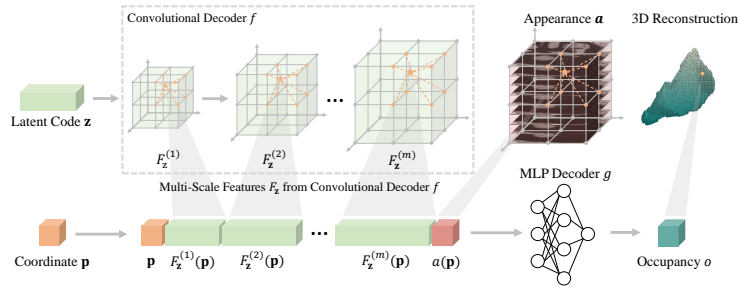[1] Code and dataset are available at https://github.com/M3DV/NeAR.

Fig. 2: **Neural Annotation Refinement (NeAR).** Given a learnable latent vector, it builds multi-scale feature maps $F_{\mathbf{z}}$ by a convolutional decoder $f$. A query coordinate $\mathbf{p}$ aggregates global and local features $F_{\mathbf{z}}^{(1)}, F_{\mathbf{z}}^{(2)}, ..., F_{\mathbf{z}}^{(m)}$, with its appearance $a$ from image. Finally, these point-wise features are fed into a light MLP $g$ for occupancy prediction $o$ to reconstruct appearance-aware surface.

the actual boundaries. It motivates us to propose the appearance-aware implicit surface model for annotation refinement. Moreover, as the MLP-based implicit functions tend to be data-hungry, which is hard to be satisfied in medical imaging scenario, we introduce the convolutional architecture with multi-scale features to reconstruct the shapes.

## 2.2 Neural Annotation Refinement

*Appearance-Aware Annotation Refinement.* The standard deep implicit surface takes spatial coordinates as input; Although the learned shape prior is able to reconstruct a high-quality surface, the reconstructed surface is possible to be misaligned with the actual boundaries. To make the implicit model appearance-aware, we employ a simple strategy by changing the input of the implicit function from spatial coordinates $\mathbf{p}$ to $\mathbf{p}$ with its appearance $a$, *i.e.*,

$$\mathcal{F}(\mathbf{z}, \mathbf{p}, a) = o : \mathbb{R}^c \times \mathbb{R}^3 \times \mathbb{R} \to [0, 1], \tag{2}$$

where $a$ short for $a(\mathbf{p})$ denotes the image appearance at the position $\mathbf{p}$, *e.g.*, Hounsfield Units in computed tomography. As will be shown, this simple modification leads to significant improvement over shape-only implicit models in both annotation refinement and downstream applications.

*Network Architecture.* As the data size is generally small in medical imaging applications, standard MLP-based implicit functions can be hard to train. To improve the data efficiency of MLP-based implicit functions, we introduce a convolutional decoder with multi-scale feature aggregation into the deep pipeline, inspired by [22,3,30]. As illustrated in Fig. 2, a latent vector $\mathbf{z}$ is first transformed by a convolutional decoder $f$ into multi-scale feature maps,

$$F_{\mathbf{z}} = [F_{\mathbf{z}}^{(1)}, F_{\mathbf{z}}^{(2)}, \cdots, F_{\mathbf{z}}^{(m)}] = f(\mathbf{z}). \tag{3}$$

In our experiments, the resolution of largest feature map $F_{\mathbf{z}}^{(m)}$ is $32 \times 32 \times 32$.

For a query point $\mathbf{p}$, it obtains its features $F_{\mathbf{z}}^{(1)}(\mathbf{p}), F_{\mathbf{z}}^{(2)}(\mathbf{p}), \cdots, F_{\mathbf{z}}^{(m)}(\mathbf{p})$ by trilinear interpolation from the multi-scale maps. To make the model appearance-aware, we further integrate the appearance $a = a(\mathbf{p})$ as an input of the implicit function. To be concrete, the coordinates and the point-wise features, are concatenated and transformed into the occupancy $o$ through a light MLP $g$,

$$o = \mathcal{F}(\mathbf{z}, \mathbf{p}, a) = g(\mathbf{p}, F_{\mathbf{z}}^{(1)}(\mathbf{p}), F_{\mathbf{z}}^{(2)}(\mathbf{p}), \cdots, F_{\mathbf{z}}^{(m)}(\mathbf{p}), a(\mathbf{p})). \tag{4}$$

*Model Training and Inference.* We treat the shape implicitly as occupancy field and train the model via auto-decoding [20]. The shape loss is measured by binary cross entropy ($BCE$) between predicted occupancy $o$ and ground-truth occupancy $\hat{o}$. Different from standard auto-encoding technique, the auto-decoding is encoder-free. We thus add regularization loss, the $l_2$-norm of the latent code. In total, the training loss is weighted by $\lambda$ ( $\lambda = 0.01$ in our experiments),

$$\mathcal{L} = BCE(o, \hat{o}) + \lambda \cdot ||\mathbf{z}||_2. \tag{5}$$

The flexibility of implicit functions enables different training resolution from actual resolution. At training stage, we sample $64 \times 64 \times 64$ meshgrid coordinates from full-resolution $128 \times 128 \times 128$ input volumes to reduce training cost. The training meshgrid is added by a random Gaussian noise $\mathcal{N}(0, 0.01^2)$, whose occupancy labels are sampled from the full-resolution ground truth. We utilize an Adam optimizer [13] with an initial learning rate of 0.001 and train the model for $1,500$ epochs. At inference stage, we sample full-resolution uniform meshgrid to reconstruct surfaces, *i.e.*, repaired annotations in this study.

*Counterpart.* Segmentation models [23] with image as input to refine the source masks, *e.g.*, 3D ResNet-based FCN [14,9] (Seg-FCN) and 3D UNet [4] (Seg-UNet), are used as counterparts. They are trained with human-annotated segmentation masks, and try to output refined annotations. Note that the label refinement counterparts are trained with the same datasets.

## 3   Datasets

### 3.1   Distorting a Golden Standard Segmentation Dataset

In order to quantitatively analyze the performance of annotation refinement, we synthesize distorted segmentation masks from golden standard. Here, we use the public AbdomenCT-1K [15], an abdominal CT organ segmentation dataset, which is annotated under multi-expert consensus protocol and thus can be regarded as golden standard. We use the adrenal gland subset[2], containing 100 adrenal glands from 50 patients. For each case, we calculate the center of left and right adrenal gland respectively, and center-crop the left and right adrenal

---

[2] https://github.com/JunMa11/AbdomenCT-1K

gland into $128 \times 128 \times 128$ volumes with a normalized spacing of $1mm^3$. The resulting dataset has 100 cases with golden standard segmentation of adrenal glands. For image pre-processing, we clip the Hounsfield Units using soft-organ window [-60, 140] and then normalize to [0,1].

To synthesize distorted segmentation masks, we randomly add or cut out cubes on the boundary. We then apply random dilation or erosion operation to the shape followed by adding small salt-and-pepper noises. The resulting distorted masks are demonstrated in Fig. 3, which imitate imperfect human annotations, including high-frequency artefacts (unsmooth boundaries) and false positive/negative. The average Dice between distorted masks and ground truth is 0.71, with a lower bound of 0.65 and an upper bound of 0.75.

### 3.2 ALAN Dataset: A New 3D Dataset for Adrenal Gland Analysis

In this study, we introduce a new 3D Adrenal gLand ANalysis dataset, named ALAN. It consists of computed tomography (CT) scans from 792 patients (*i.e.*, 1,584 left and right adrenal glands). Each case is annotated with a segmentation mask and a binary diagnosis label (normal vs. abnormal) [29]. The segmentation mask is annotated by a single clinician using 3D Slicer software. As the boundary of adrenal gland–soft organ–is difficult to identify, and the segmentation is made slice-by-slice, the resulting 3D segmentation is imperfect with potential errors, *e.g.*, inconsistent cross-slice segmentation, high-frequency artefacts and human mistakes. Different from the segmentation masks, the diagnosis labels are independently made by 2 clinicians, and confirmed by 1 senior endocrinologist when diagnoses of the 2 clinicians disagree. We pre-process the dataset from raw 792 CT scans into 1,584 3D image cubes of $128 \times 128 \times 128$ following Sec. 3.1.

As the segmentation mask of adrenal glands is imperfect, we repair the annotation with the proposed NeAR. To demonstrate the usefulness of repairing, we run 3D convolutional networks with the shapes of adrenal glands as inputs, to output the diagnosis labels. The networks are trained with the ALAN dataset, with training / validation / test split of 1,188 / 98 / 298 on a patient level. As will be shown in Sec. 4.2, the adrenal gland shapes repaired by NeAR could be better classified than the human annotated ones.

The ALAN dataset will be open-source, with 1,584 cases of NeAR-repaired, high-quality 3D models of adrenal glands together with the corresponding diagnosis labels. As there are only a few publicly available medical shape datasets [31,27] (see supplementary materials for a comparison), our dataset will be a valuable addition to the medical image analysis and geometric deep learning community.

## 4   Experiments

### 4.1   Quantitative Experiments on Distorted Golden Standards

*Experiment Setting.*  All our experiments are implemented with PyTorch 1.8 [21]. To quantitatively analyze the performance of the proposed NeAR method on

Table 1: **Segmentation Repairing on the Distorted Golden Standard Dataset.** We compare the counterparts (Seg-FCN and Seg-UNet), NeAR w/ shape only (S), and NeAR w/ shape and appearance (S+A) in Dice Similarity Coefficient (DSC) and Normalized Surface Dice (NSD) over 5 trials.

| Metrics | Seg-FCN | Seg-UNet | NeAR (S) | NeAR (S+A) |
|---|---|---|---|---|
| DSC (%, ↑) | $79.56 \pm 0.45$ | $78.70 \pm 0.45$ | $78.79 \pm 0.45$ | $\mathbf{81.07 \pm 0.22}$ |
| NSD (%, ↑) | $89.54 \pm 0.33$ | $87.71 \pm 0.90$ | $87.96 \pm 0.51$ | $\mathbf{91.22 \pm 0.12}$ |

repairing segmentation annotations, we implement several methods on the distorted golden standard segmentation dataset, including

- **NeAR (S+A)**: The full proposed method;
- **NeAR (S)**: The shape-only NeAR model without appearance $a$ as input;
- **Seg-FCN / Seg-UNet**: Segmentation counterparts, see Sec. 2.2.

All these methods are trained with all 100 cases, consisting of image and distorted segmentation mask, and evaluated by comparing the similarity between the model-predicted and golden standard segmentation masks. Best models are selected with the lowest training loss. The evaluation is based on volume-based Dice Similarity Coefficient (DSC), and surface-based Normalized Surface Dice (NSD) [19] with a distance tolerance of 1.0.

The segmentation models need to maintain 3D feature maps in the encoder, and thus take much larger memory and computation under the same training resolution as the NeAR. We use a slightly different training schedule. We observed that the number of training iterations of these segmentation models is smaller than the NeAR. Therefore, we utilize an Adam optimizer [13] with an initial learning rate of 0.001 for 100 epochs, delaying the learning rate by 0.1 after 50 and 75 epochs. Longer training schedule does not lead to higher performance.

*Results.* As depicted in Tab. 1, the NeAR (S+A) surpasses all the other methods on both DSC and NSD, especially in surface-based method (NSD). NeAR (S) underperforms NeAR (S+A) as well as the standard segmentation method Seg-FCN and Seg-UNet, indicating that appearance-awareness boosts the repair performance significantly. Fig. 3 shows contours of adrenal gland on image slices and 3D visualization of repaired annotations for each method. As shown by the contours on image slices, NeAR (S+A) can repair the distorted annotations more accurately and fix distortions that other methods fail to repair.

Moreover, we add manual smoothing as a baseline, including morphological closing and connected components filtering. We tried several settings, and the highest Dice is 76.90%, much lower than neural methods.

## 4.2   Adrenal Diagnosis on the Repaired ALAN Dataset

*Experiment Setting.* As described in Sec. 3.2, the segmentation masks in the ALAN dataset are imperfect. We utilize 4 different methods to repair the 3D
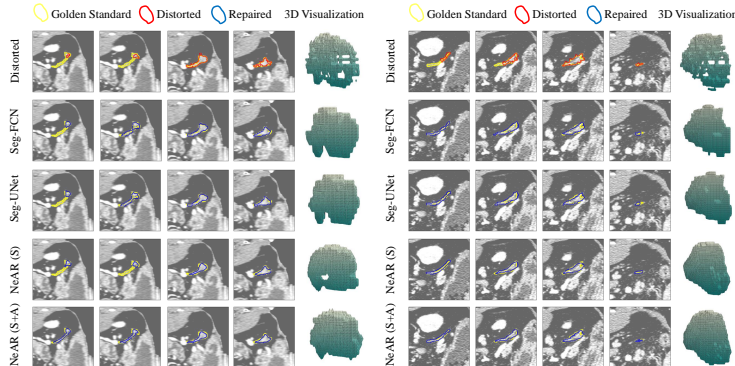
Fig. 3: **Visualization of Repaired Annotations.** Contours of adrenal glands are shown on image slices. Red are distorted, yellow are golden standard, and blue are repaired. The 3D visualization is shown on the right side.

adrenal gland analysis dataset, standard segmentation methods Seg-FCN and Seg-UNet as the baseline methods, NeAR (S) and NeAR (S+A). These methods are conducted in the same settings as that in Sec. 4.1 respectively. As there are no golden standard segmentation masks in the ALAN dataset, we only provide the qualitative results by visualizing the repaired segmentation masks in the supplementary materials.

To quantitatively analyze the annotation repairing quality, we conduct shape classification experiments on the human-annotated and model-repaired ALAN dataset. ResNet [9] variants with 2D / 3D / ACS [28] convolutions are implemented to classify the shapes into binary diagnosis labels (normal vs. abnormal). The shapes of $128 \times 128 \times 128$ are resized to the size of $48 \times 48 \times 48$ as model inputs. For model training, we utilize an Adam optimizer [13] with an initial learning rate of 0.001 for 50 epochs, delaying the learning rate by 0.1 after 25 and 40 epochs. We use cross-entropy loss, and report area under ROC curve (AUC) as the evaluation metric. Best models are selected with lowest validation loss. We repeat experiments for 5 trials for each setting.

*Results.* As depicted in Tab. 2, models trained with shapes repaired by the NeAR (S+A) can diagnose the adrenal glands better than other methods, as well as the human annotated imperfect ones. Notably, standard segmentation methods (Seg-FCN and Seg-UNet) deliver worse shape classification results than NeAR (S) and raw human annotation, though the segmentation methods produce better shape repairing results than NeAR (S) in Sec. 4.1. This implies that the learned prior of deep implicit surfaces can be particularly useful for downstream applications.

## 5    Conclusion

This study addresses a practical problem in medical image analysis: how to repair the imperfect segmentation. We propose Neural Annotation Refinement,

Table 2: **Shape Classification on the ALAN Dataset.** We repair the 3D shapes of adrenal glands using standard segmentation (Seg-FCN and Seg-UNet), NeAR w/ shape only (S), and NeAR w/ shape and appearance (S+A). ResNet-18 and ResNet-50 variants are trained to classify the 3D shapes of adrenal glands (normal vs. abnormal) on the human-annotated and the repaired datasets. We report mean and standard deviation of AUC (%, ↑) on the test set over 5 trials.

| Networks | Human-Annotated | Seg-FCN | Seg-UNet | NeAR (S) | NeAR (S+A) |
|---|---|---|---|---|---|
| ResNet-18 [9] (2D) | $68.35 \pm 2.53$ | $65.17 \pm 2.21$ | $64.68 \pm 2.45$ | $65.95 \pm 0.83$ | $\mathbf{69.69 \pm 1.44}$ |
| ResNet-18 [9] (3D) | $89.77 \pm 1.20$ | $86.27 \pm 1.91$ | $86.55 \pm 0.89$ | $88.64 \pm 1.24$ | $\mathbf{90.38 \pm 0.57}$ |
| ResNet-18 [9] (ACS [28]) | $90.10 \pm 0.90$ | $87.02 \pm 2.20$ | $86.83 \pm 2.00$ | $89.22 \pm 1.08$ | $\mathbf{91.11 \pm 0.35}$ |
| ResNet-50 [9] (2D) | $66.36 \pm 2.56$ | $64.88 \pm 3.47$ | $66.06 \pm 2.79$ | $69.04 \pm 3.11$ | $\mathbf{69.94 \pm 1.18}$ |
| ResNet-50 [9] (3D) | $89.72 \pm 1.33$ | $85.64 \pm 1.59$ | $84.92 \pm 1.08$ | $88.76 \pm 0.91$ | $\mathbf{89.78 \pm 0.79}$ |
| ResNet-50 [9] (ACS [28]) | $90.13 \pm 0.40$ | $85.91 \pm 2.11$ | $82.28 \pm 2.34$ | $89.42 \pm 1.35$ | $\mathbf{90.72 \pm 0.72}$ |

an appearance-aware implicit method, whose values are validated in repairing segmentation and downstream applications. Moreover, the ALAN dataset for 3D shape classification will be an addition for the research community. There are limitations in the current study, *e.g.*, validated on adrenal glands only. We will test the NeAR on sparse annotations and small objects in the future research.

# References

1. Ardila, D., Kiraly, A.P., Bharadwaj, S., Choi, B., Reicher, J.J., Peng, L., Tse, D., Etemadi, M., Ye, W., Corrado, G., et al.: End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. Nature medicine **25**(6), 954–961 (2019)
2. Chen, Z., Zhang, H.: Learning implicit fields for generative shape modeling. In: Conference on Computer Vision and Pattern Recognition. pp. 5939–5948 (2019)
3. Chibane, J., Alldieck, T., Pons-Moll, G.: Implicit functions in feature space for 3d shape reconstruction and completion. In: Conference on Computer Vision and Pattern Recognition. pp. 6970–6981 (2020)
4. Çiçek, Ö., Abdulkadir, A., Lienkamp, S.S., Brox, T., Ronneberger, O.: 3d u-net: learning dense volumetric segmentation from sparse annotation. In: Conference on Medical Image Computing and Computer Assisted Intervention. pp. 424–432. Springer (2016)
5. Esteva, A., Chou, K., Yeung, S., Naik, N., Madani, A., Mottaghi, A., Liu, Y., Topol, E., Dean, J., Socher, R.: Deep learning-enabled medical computer vision. NPJ digital medicine **4**(1), 1–9 (2021)
6. Esteva, A., Kuprel, B., Novoa, R.A., Ko, J., Swetter, S.M., Blau, H.M., Thrun, S.: Dermatologist-level classification of skin cancer with deep neural networks. Nature **542**(7639), 115–118 (2017)

7. Gulshan, V., Peng, L., Coram, M., Stumpe, M.C., Wu, D., Narayanaswamy, A., Venugopalan, S., Widner, K., Madams, T., Cuadros, J., et al.: Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. Jama **316**(22), 2402–2410 (2016)

8. Hanocka, R., Metzer, G., Giryes, R., Cohen-Or, D.: Point2mesh: A self-prior for deformable meshes. ACM SIGGRAPH (2020)

9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Conference on Computer Vision and Pattern Recognition. pp. 770–778 (2016)

10. Huang, X., Yang, J., Wang, Y., Chen, Z., Li, L., Li, T., Ni, B., Zhang, W.: Representation-agnostic shape fields. In: International Conference on Learning Representations (2022)

11. Karimi, D., Dou, H., Warfield, S.K., Gholipour, A.: Deep learning with noisy labels: Exploring techniques and remedies in medical image analysis. Medical Image Analysis **65**, 101759 (2020)

12. Kermany, D.S., Goldbaum, M., Cai, W., Valentim, C.C., Liang, H., Baxter, S.L., McKeown, A., Yang, G., Wu, X., Yan, F., et al.: Identifying medical diagnoses and treatable diseases by image-based deep learning. Cell **172**(5), 1122–1131 (2018)

13. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv Preprint (2014)

14. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Conference on Computer Vision and Pattern Recognition. pp. 3431–3440 (2015)

15. Ma, J., Zhang, Y., Gu, S., Zhu, C., Ge, C., Zhang, Y., An, X., Wang, C., Wang, Q., Liu, X., Cao, S., Zhang, Q., Liu, S., Wang, Y., Li, Y., He, J., Yang, X.: Abdomenct-1k: Is abdominal organ segmentation a solved problem? IEEE Transactions on Pattern Analysis and Machine Intelligence (2021)

16. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. In: Conference on Computer Vision and Pattern Recognition. pp. 4460–4470 (2019)

17. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: European Conference on Computer Vision. pp. 405–421. Springer (2020)

18. Niemeyer, M., Geiger, A.: Giraffe: Representing scenes as compositional generative neural feature fields. In: International Conference on Computer Vision. pp. 11453–11464 (2021)

19. Nikolov, S., Blackwell, S., Zverovitch, A., Mendes, R., Livne, M., De Fauw, J., Patel, Y., Meyer, C., Askham, H., Romera-Paredes, B., et al.: Deep learning to achieve clinically applicable segmentation of head and neck anatomy for radiotherapy. arXiv Preprint (2018)

20. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. In: Conference on Computer Vision and Pattern Recognition. pp. 165–174 (2019)

21. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. Advances in Neural Information Processing Systems **32** (2019)

22. Peng, S., Niemeyer, M., Mescheder, L., Pollefeys, M., Geiger, A.: Convolutional occupancy networks. In: European Conference on Computer Vision. pp. 523–540. Springer (2020)

23. Rajchl, M., Lee, M.C., Oktay, O., Kamnitsas, K., Passerat-Palmbach, J., Bai, W., Damodaram, M., Rutherford, M.A., Hajnal, J.V., Kainz, B., et al.: Deepcut: Object segmentation from bounding box annotations using convolutional neural networks. IEEE Transactions on Medical Imaging **36**(2), 674–683 (2016)
24. Tajbakhsh, N., Jeyaseelan, L., Li, Q., Chiang, J.N., Wu, Z., Ding, X.: Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation. Medical Image Analysis **63**, 101693 (2020)
25. Ulyanov, D., Vedaldi, A., Lempitsky, V.: Deep image prior. In: Conference on Computer Vision and Pattern Recognition (2018)
26. Xu, Q., Wang, W., Ceylan, D., Mech, R., Neumann, U.: Disn: Deep implicit surface network for high-quality single-view 3d reconstruction. Advances in Neural Information Processing Systems **32** (2019)
27. Yang, J., Gu, S., Wei, D., Pfister, H., Ni, B.: Ribseg dataset and strong point cloud baselines for rib segmentation from ct scans. In: Conference on Medical Image Computing and Computer Assisted Intervention. pp. 611–621. Springer (2021)
28. Yang, J., Huang, X., He, Y., Xu, J., Yang, C., Xu, G., Ni, B.: Reinventing 2d convolutions for 3d images. IEEE Journal of Biomedical and Health Informatics **25**(8), 3009–3018 (2021)
29. Yang, J., Shi, R., Wei, D., Liu, Z., Zhao, L., Ke, B., Pfister, H., Ni, B.: Medmnist v2: A large-scale lightweight benchmark for 2d and 3d biomedical image classification. arXiv Preprint (2021)
30. Yang, J., Wickramasinghe, U., Ni, B., Fua, P.: Implicitatlas: Learning deformable shape templates in medical imaging. In: Conference on Computer Vision and Pattern Recognition. pp. 15861–15871 (2022)
31. Yang, X., Xia, D., Kin, T., Igarashi, T.: Intra: 3d intracranial aneurysm dataset for deep learning. In: Conference on Computer Vision and Pattern Recognition. pp. 2656–2666 (2020)

Table A1: **A Comparison of Public Available Medical Shape Datasets.**

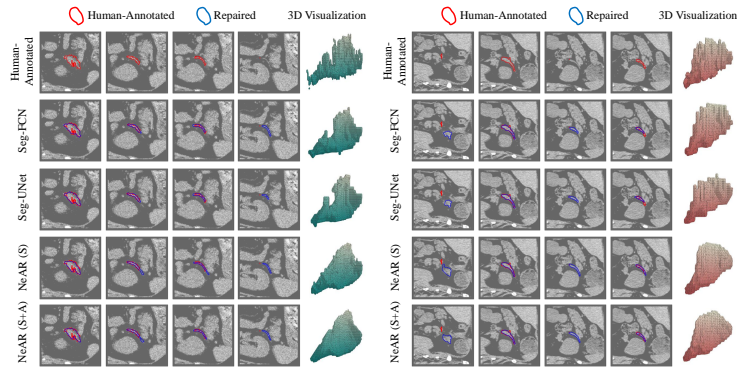|                 | IntrA [31]                    | RibSeg [27]   | ALAN           |
|-----------------|-------------------------------|---------------|----------------|
| Task            | Classification & Segmentation | Segmentation  | Classification |
| Body Part       | Vessel                        | Rib           | Adrenal Gland  |
| Characteristics | Tubular                       | Tubular       | Soft Organ     |
| No. Cases       | 1,909                         | 490           | 1,584          |



Fig. A1: **Visualization of Two Samples in the Repaired ALAN Dataset.**
Contours of adrenal glands are shown on image slices. Red are human-annotated,
and blue are repaired. The 3D visualization is shown on the right side, green/red
denotes normal/abnormal adrenal glands.