

A Comprehensive Benchmarking of A U-Net based Model for Midbrain Auto-segmentation on Transcranial Sonography

1 **Hong-yu KANG^{1, #}, Wei ZHANG^{2, #, *}, Shuai LI¹, Xinyi WANG¹, Yu SUN², Xin SUN², Fang-**
2 **xian LI², Chao HOU², Sai-kit LAM^{1, 3}, Yong-ping ZHENG^{1, 3, *}**

3 ¹Department of Biomedical Engineering, The Hong Kong Polytechnic University, HKSAR, China.

4 ²Department of Ultrasound, Beijing Tiantan Hospital, Capital Medical University, Beijing, China.

5 ³Research Institute of Smart Ageing, The Hong Kong Polytechnic University, HKSAR, China.

6

7 [#]The two authors contributed equally to this work.

8

9 ***Correspondence:**

10 Yongping ZHENG, PhD

11 Department of Biomedical Engineering

12 The Hong Kong Polytechnic University

13 Hung Hom, Kowloon

14 Hong Kong, China

15 Email: yongping.zheng@polyu.edu.hk

16

17 Wei ZHANG, MD

18 Department of Ultrasound

19 Beijing Tiantan Hospital

20 Capital Medical University

21 NO.119, South 4th Ring West Road

22 Fengtai District

23 Beijing 100070, China

24 Email: ultrazhangwei@126.com

25

26 **Highlights**

- 27 • The nnU-Net model yielded the best midbrain segmentation agreement, achieving the top-ranked
28 averaged scores compared U-Net and U-Net+++ on original TCS Images.
- 29 • The nnU-Net demonstrated the best-performing performance in aspects of segmentation
30 agreement, model stability, and time efficiency.
- 31 • The segmentation power of the nnU-Net network remained robust against ultrasound imaging
32 noise.
- 33 • This large cohort study is the first of its kind to benchmark the best-performing state-of-the-art
34 deep neural network for TCS-based midbrain auto-segmentation in a wide spectrum of application
35 perspectives.

36

37 **Abstract**

38 **Background and Objective:** Transcranial sonography-based grading of Parkinson's Disease has
39 gained increasing attention in recent years, and it is currently used for assistive differential diagnosis
40 in some specialized centers. To this end, accurate midbrain segmentation is considered an important
41 initial step. However, current practice is manual, time-consuming, and bias-prone due to the subjective
42 nature. Relevant studies in the literature are scarce and lacks comprehensive model evaluations from
43 application perspectives. Herein, we aimed to benchmark the best-performing U-Net model for
44 objective, stable and robust midbrain auto-segmentation using transcranial sonography images.

45 **Methods:** A total of 584 patients who were suspected of Parkinson's Disease were retrospectively
46 enrolled from Beijing Tiantan Hospital. The dataset was divided into training (n=416), validation
47 (n=104), and testing (n=64) sets. Three state-of-the-art deep-learning networks (U-Net, U-Net+++, and
48 nnU-Net) were utilized to develop segmentation models, under 5-fold cross-validation and three
49 randomization seeds for safeguarding model validity and stability. Model evaluation was conducted in
50 testing set in three key aspects: (i) segmentation agreement using DICE coefficients (DICE),
51 Intersection over Union (IoU), and Hausdorff Distance (HD); (ii) model stability using standard
52 deviations of segmentation agreement metrics; (iii) prediction time efficiency, and (iv) model
53 robustness against various degrees of ultrasound imaging noise produced by the salt-and-pepper noise
54 and Gaussian noise.

55 **Results:** The nnU-Net achieved the best segmentation agreement (averaged DICE: 0.910, IoU: 0.836,
56 HD: 2.793-mm) and time efficiency (1.456-s). Under mild noise corruption, the nnU-Net outperformed
57 others with averaged scores of DICE (0.904), IoU (0.827), HD (2.941 mm) in the salt-and-pepper noise
58 (signal-to-noise ratio, SNR = 0.95), and DICE (0.906), IoU (0.830), HD (2.967 mm) in the Gaussian
59 noise (sigma value, $\sigma = 0.1$); by contrast, intriguingly, performance of the U-Net and U-Net+++ models
60 were remarkably degraded. Under increasing levels of simulated noise corruption (SNR decreased
61 from 0.95 to 0.75; σ increased from 0.1 to 0.5), the nnU-Net network exhibited marginal decline in
62 segmentation agreement meanwhile yielding decent performance as if there were absence of noise
63 corruption.

64 **Conclusions:** The nnU-Net model was the best-performing midbrain segmentation model in terms of
65 segmentation agreement, stability, time efficiency and robustness, providing the community with an
66 objective, effective and automated alternative. Moving forward, a multi-center multi-vendor study is
67 warranted when it comes to clinical implementation.

68

69 **Keywords:** Parkinson's Disease, Midbrain, Transcranial Sonography, Auto-segmentation, Deep
70 Learning.

71

72 1 Introduction

73 Parkinson's Disease (PD) is a chronic, progressive and devastating neurodegenerative disorder,
74 accounting for over 9.4 million of sufferers across the globe in 2020 [1] and it has been considered as
75 the fastest-growing neurological illness [2, 3]. Early diagnosis and frequent surveillance are the key
76 antidote to this crippling disease. Among various imaging approaches, transcranial sonography (TCS)
77 possesses unique particularities of being non-invasive, radiation-free, rapid, low cost, highly
78 accessible, high patient compliance, easy-to-operate, as well as proven to offer the correlations of its
79 hyper echoic pattern in brain substantia nigra with incidence and severity of PD [4-8]. In light of this,
80 several international guidelines on TCS-based PD assessment and grading have been well-documented
81 to aid neurologists in offering enhanced healthcare delivery to PD sufferers in aspects of early diagnosis
82 and treatment outcome assessment [9-11]. Its popularity is on a rapid rise around the world [12, 13].

83 In clinical practice, an accurate segmentation of the midbrain structure is considered a crucial
84 initial task. [14, 15] in TCS-based PD grading and severity assessment [16, 17]. However, traditional
85 practice of midbrain segmentation is entirely a manual process, which suffers from two major
86 drawbacks. First, the time-consuming nature of the practice has not only posed significant clinical
87 burdens, especially in busy clinics, but also called for experienced clinicians for accurate midbrain
88 segmentation which may not be easily overcome in resource-demanding regions, such as developing
89 and under-developed countries. Second, the issue of subjectiveness is intrinsically embedded in the
90 nature of the manual segmentation process, which may in turn incur bias in subsequent PD grading and
91 assessment. In the contemporary era of AI in medicine, there is a pressing demand for an automated,
92 efficient, and objective technique for midbrain segmentation in order to lay a solid foundation toward
93 revolutionizing the clinical practice of TCS-based PD assessment in the long run.

94 With this regard, the U-Net deep-neural network and its variants, including the U-Net+++ and
95 the nnU-Net, have been gaining increasing popularity in the community for medical image
96 segmentation. In 2015, the U-Net neural network was first introduced; its design presents a U-shaped
97 architecture and it utilizes an encoder-decoder structure with skip connections to capture contextual
98 information for achieving precise localization [18]. In 2020, an extension of the U-Net network called
99 U-Net +++ was introduced. It is coupled with nested pathways and deeply supervised full-scale skip
100 connections, enabling multi-scale contextual information [19]. In 2021, the nnU-Net model, an
101 adaptive, unified extensible training framework, was introduced for segmenting medical images.
102 Instead of altering the network architecture, the authors proposed a comprehensive training workflow
103 structure. [20]. The segmentation capabilities of these three neural networks on various imaging
104 modalities have been well-documented, including computed tomography [21-23], magnetic resonance
105 imaging (MRI) [24-26] and ultrasound [27-30]. In spite of their promising segmentation power, the
106 application of AI in midbrain segmentation is scarce in the literature.

107 TCS-based midbrain segmentation can be broadly classified into two categories according to
108 the current body of literature: traditional and AI-based approaches. For traditional method, Sakalauskas
109 et al. revealed a modified shape-based active contour segmentation algorithm for midbrain using a
110 dataset of 40 TCS images in 2013, which achieved a DICE of $73.1 \pm 7.5\%$ [17]. Later in 2016,
111 Sakalauskas et al. integrated an experience-based statistical shape model with an intensity-amplitude
112 invariant edge detector, for extracting the fuzzy boundaries of the midbrain on TCS images using a
113 dataset of 130 TCS images, which achieved a mean DICE and Hausdorff Distance (HD) between 87.6-
114 89.6% and 3.60-6.09 mm, respectively [31]. For AI-based strategies, the only relevant research work
115 found in the literature was conducted by Milletari et al. (2017), they enrolled 34 subjects and developed
116 a Hough-CNN network for midbrain segmentation, which achieved a DICE between 0.77-0.85,

117 depending on different configurations of the convolutional layers [32]. These abovementioned studies
118 have reflected the importance and laid a foundation of TCS-based automated midbrain segmentation
119 over the past decade.

120 Nevertheless, the contemporary era of AI has placed a growing emphasis on the sample size and
121 a boarder scope of model evaluations from application perspectives, including and beyond the
122 segmentation accuracy. The only relevant research in the literature by Milletari et al. may not be
123 capable of providing sufficient values of their developed models from clinical implementation
124 perspectives, on top of the inadequacy of sample size [32]. To address these challenges, this study
125 aimed to conduct a large-cohort study and benchmark a U-Net based network from U-Net, U-Net+++,
126 and nnU-Net for midbrain auto-segmentation on TCS images. Efficacy of the developed models were
127 comprehensively evaluated in multiple application facets, including segmentation agreements, model
128 stability, time efficiency, as well as model robustness against noise-corrupted TCS images. Findings
129 of this study are anticipated to provide the community with enhanced understanding regarding the
130 clinical usefulness of the developed auto-segmentation models, and to pave the way for further
131 investigations on substantia nigra auto-segmentation towards objective, effective and fully automated
132 PD assessment using TCS images in the long run.

133

134 **2 Methods**

135 *2.1 Data Acquisition, Annotation and Image Pre-processing*

136 TCS images of 584 patients, who were suspected or diagnosed with PD, were retrospectively enrolled
137 from Beijing Tiantan Hospital between 2021 and 2023. Patient consent was waived because of the
138 retrospective nature of this study. Ethical approval for this study was obtained from the Human Subject
139 Ethics Sub-committee (HSESC) of the Hong Kong Polytechnic University (HSEARS20231102004).
140 All the original TCS images were acquired by a Canon Aplio i900 i-series ultrasound system (Canon
141 Medical Systems Corporation, Otawara, Tochigi, Japan) with i8CX1 convex array (center frequency =
142 2.6 MHz). An experienced physician maneuvered the ultrasound probe over the temporal region,
143 followed by capturing the ultrasound imaging frames that display the butterfly-shaped midbrain
144 structure.

145 The ground-truth midbrain annotation was manually delineated on the captured imaging frames using
146 the Canon Aplio i900 i-series ultrasound system by experienced physicians with over 10 years of
147 experience.

148 In this study, several pre-processing procedures were performed on the original TCS images prior to
149 downstream analyses. First, all the TCS images together with the corresponding ground-truth midbrain
150 annotations were cropped from a resolution of 1280x960 to 400x320, in order to effectively localize
151 the midbrain region within the region-of-interest (ROI) meanwhile minimizing irrelevant information
152 from other structures for achieving effective learning during subsequent model development. Second,
153 a binary mask was generated by extracting the region within the ground-truth midbrain annotations for
154 each patient; this procedure enabled an accurate representation of the segmentation target for
155 downstream model development and evaluation.

156 *2.2 Deep-learning Neural Networks*

157 In the present study, three U-Net based neural networks, including U-Net, U-Net+++, and nnU-Net,
158 were deployed for developing auto-segmentation models. Architectures of these three networks are
159 illustrated in **Supplementary Figure 1A-1D**. Details on settings of each of these networks were
160 presented as follows:

161 2.2.1 *U-Net*

162 The U-Net model is a state-of-the-art architecture designed dedicatedly for medical image
163 segmentation [18]. It is one of the most widely adopted segmentation network within the medical
164 imaging community. The U-Net network exhibits a U-shaped architecture, it strategically combines
165 encoders and decoders, enabling information fusion throughout the network, for achieving precise
166 segmentation outcomes. In this present experiment, the U-Net model contains four layers with 64, 128,
167 256 and 512 feature channels, respectively. In the encoder path, the input undergoes four layers
168 convolution blocks for feature extraction, followed by four max pooling operations; each layer contains
169 two 3x3 convolutions and same padding, followed by a 2x2 max pooling with strides of two in each
170 dimension. In the decoder path, the extracted features are subsequently processed through four layers
171 of convolution blocks and undergo four rounds of deconvolutions; each layer includes a 2x2 up-
172 convolution with strides of two in each dimension, followed by two 3x3 convolutions and the same
173 padding. For model development, a maximum of 100 epochs was specified, starting with an initial
174 learning rate of 1e-4. The momentum was set to 0.9 to aid in faster convergence and stability, while
175 the batch size of 4 was chosen to determine the number of samples processed before updating the model
176 parameters.

177 2.2.2 *U-Net+++*

178 The U-Net+++ model is a full-scale skip connection and deep supervision based on the U-Net
179 architecture [19]. The U-Net+++ introduces a novel full-size skip connection, which strengthens the
180 interconnections between the encoder and decoder, as well as the intra-connections among decoder
181 subnetworks, allowing each decoder layer to effectively incorporate both small-scale and same-scale
182 feature maps from the encoder. U-Net+++ is designed to aggregate features across all scales,
183 incorporating skip connections between the top and bottom layers of both the encoder and decoder. In
184 this study, the U-Net+++ network contains four layers with 64, 128, 256 and 512 feature channels,
185 respectively. In the encoder path, each layer comprises two 3x3 convolutions and the same padding,
186 followed by a 2x2 max pooling operation with strides of two in each dimension. In the decoder path,
187 each layer of the U-Net+++ model comprises a 2x2 up-convolution with strides of two in each
188 dimension, followed by two 3x3 convolutions and the same padding. The model includes skip
189 connections at the same feature map level within the encoder, skip connections from higher to lower
190 feature map levels, and decoder skip connections. The parameters of the U-Net+++ model remained
191 the same as those of the U-Net network, including a maximum of 100 epochs, an initial learning rate
192 of 0.0001, a momentum value of 0.9, and a batch size of 4.

193 2.2.3 *nnU-Net*

194 The nnU-Net model, an adaptive and unified extensible training framework, was introduced for the
195 segmentation of medical images. [20]. It is an automatic deep learning-based segmentation method that
196 configures itself for new tasks, encompassing preprocessing, network architecture, training, and post-
197 processing. It is built on an adaptive framework based on 2D and 3D U-Net, and automatically tunes
198 hyperparameters based on the specific properties of the given dataset, such as the exact patch size,
199 batch size, and inference settings based on the size and characteristics of the input images. In this
200 research work, the 2D version of nnU-Net was utilized. The model consists of an encoder path and a
201 decoder path. Within the nnU-Net architecture used in this study, there are seven layers, each with 32,

64, 128, 256, 512, 512 and 512 feature channels, respectively. In the encoder path, each layer employs two 3x3 convolutions to extract hierarchical features from the input data. This process allows for progressive capture of intricate details and patterns presented in the midbrain of the TCS images. Conversely, in the decoder path, each layer consists of a 2x2 up-convolution with strides of two in each dimension, followed by two 3x3 convolutions. The nnU-Net model was trained with a maximum of 100 epochs. Several other parameters were set to their initial values in the experiment. The learning rate was set to 0.0001. A momentum value of 0.9 was used to enhance the stability of the training process. The training was performed using a batch size of 4.

2.3 Model Development and Evaluation

To comprehensively compare the performance of the three models, the model Development involves dataset partitioning, image pre-processing, and training with 5-fold cross-validation using 3 random seeds. The evaluation of the models includes two essential aspects: assessing model segmentation performance and evaluating model robustness under various noise conditions.

2.3.1 Model Development

Figure 1 showcases a schematic diagram of the entire experimental workflow. First, the entire dataset (n=584) into was randomly into training, validation, and testing sets via a randomization seed, at an approximate ratio of 70% (n=416), 20% (n=104), 10% (n=64), respectively, taking reference from previous literature [33]. All the analyzed images at this stage were processed according to the image pre-processing procedures depicted in **Section 2.1.3**.

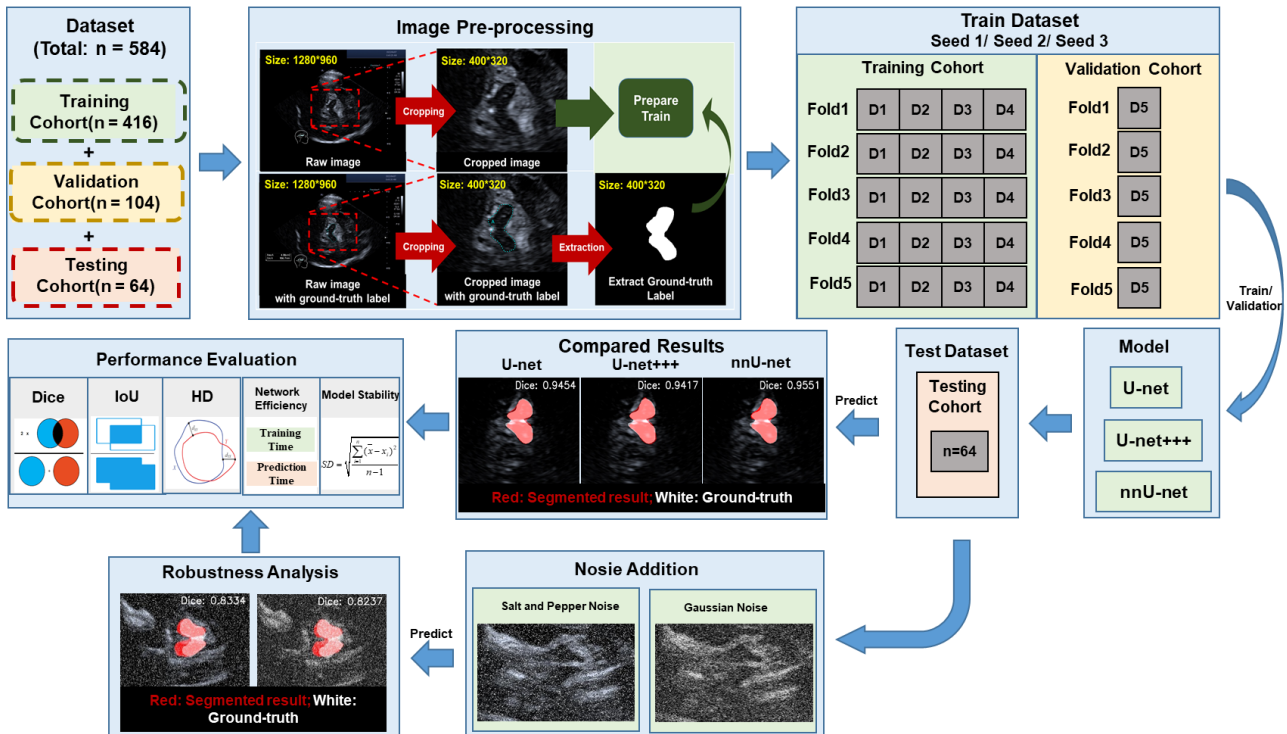


Fig. 1. Overall study workflow from dataset randomization, image pre-processing, to model training, robustness analysis after noise addition on the original images, and model performance evaluation

In the training set, three deep-learning networks (U-Net, U-Net+++, and nnU-Net) were separately utilized to develop three segmentation models under 5-fold cross-validation. In this experiment, three independent randomization seeds were deployed to obtain three sets of results for

227 the sake of harvesting a fair model performance across different patient sub-populations. The approach
 228 of multiple randomized stratification was commonly adopted [16, 34, 35]. An averaged performance
 229 of each model was reported in this study to assess model stability and to improve validity of the findings.
 230 The training process was carried out on an in-house service machine equipped with an Nvidia RTX
 231 A6000 GPU card. Training time for each model was recorded.

232 2.3.2 Model Evaluation

233 The developed models were evaluated in the testing set in two scenarios: (a) Model segmentation
 234 performance on original images; and (b) Model robustness assessment on images corrupted by two
 235 types of typical ultrasound imaging noise (after implementing noise addition on the original images).
 236 Details of the model evaluation procedures are described below:

237 (a) Model Segmentation Performance on Original TCS Images

238 All the developed models were evaluated in the testing set in aspects of model segmentation agreement,
 239 model stability and prediction time efficiency. For model segmentation agreement, DICE coefficient
 240 (DICE), Intersection over Union (IoU), and Hausdorff Distance (HD), were calculated for evaluation.
 241 These evaluating metrics have been widely adopted in the field of medical image segmentation [36,
 242 37]. Relevant equations and definitions are shown below:

243 **DICE** is a measure of similarity between two segmentations. It provides a value between 0 and 1,
 244 where a higher value indicates better, X and Y represent the area of predict image and ground truth
 245 respectively:

$$246 \quad DICE = \frac{2|X \cap Y|}{|X| + |Y|} \quad (\text{Eq. 1})$$

247 **IoU** is a metric that measures the overlap between two regions. It is calculated by dividing the
 248 intersection area of the predicted and ground truth by their union area. The IoU value ranges from 0 to
 249 1, where a higher value indicates better, X and Y represent the area of predict image and ground truth
 250 respectively:

$$251 \quad IoU = \frac{|X \cap Y|}{|X \cup Y|} \quad (\text{Eq. 2})$$

252 **HD** measures the maximum dissimilarity between two regions. A smaller HD indicates better
 253 agreement between segmentations, quantifying the maximum discrepancy between predicted and
 254 ground truth boundaries, A and B are point sets, a and b are points in A and B respectively, and $d(a,$
 255 $b)$ represents the distance between point a and point b:

$$256 \quad HD(A, B) = \max(\max_{a \in A}(\min_{b \in B}(d(a, b))), \max_{b \in B}(\min_{a \in A}(d(b, a)))) \quad (\text{Eq. 3})$$

257
 258
 259 **For model stability**, standard deviations (SD) of each of the above segmentation evaluating
 260 metrics (DICE, IoU, HD) under 5-fold cross-validation across the three randomization seeds were
 261 adopted for evaluation.

$$SD = \sqrt{\frac{\sum_{i=1}^n (\bar{x} - x_i)^2}{n-1}} \quad (\text{Eq. 4})$$

For prediction time efficiency, the time (in second) required for the developed models to execute prediction on the testing set was calculated for assessment.

(b) Model Robustness Assessment on Noise-corrupted TCS Images

In addition to segmentation agreement, model stability and time efficiency, the developed models were evaluated in aspects of their robustness under a varying degree of noise-corrupted ultrasound images. Salt-and-pepper noise and Gaussian noise are two common types of ultrasound imaging noise [38], and they have been adopted in AI studies of ultrasound-based auto-segmentation [39]. Salt-and-pepper noise is a random impulse noise where certain pixels have different intensities compared to their neighbors. Gaussian noise, on the other hand, is a statistical noise caused by random signal fluctuations and follows a Gaussian distribution. In this study, the quality of the images in the testing set of 64 testing images was corrupted by addition of the salt-and-pepper noise and the Gaussian noise.

For the salt-and-pepper noise, three noise levels reflected by the signal-to-noise ratio (SNR) of 0.95, 0.85, and 0.75 were introduced and analyzed. For the Gaussian noise, three noise levels in terms of the sigma values (σ) of 0.1, 0.3, and 0.5, representing the standard deviation of a Gaussian distribution, were employed and analyzed.

To analyze model robustness, segmentation agreements of the three comparing networks (U-Net, U-Net+++, and nnU-Net) were first assessed in terms of DICE, IoU, and HD, under a mild level of noise corruption on the salt-and-pepper noise corrupted images (SNR=0.95) and the Gaussian noise corrupted images ($\sigma=0.1$). Following the above analysis, the best-performing network was benchmarked. Finally, robustness of the benchmarked model was subsequently further analyzed on images with increasing degrees of noise corruption on the salt-and-pepper noise corrupted images (SNR=0.95, 0.85, 0.75) and the Gaussian noise corrupted images ($\sigma=0.1, 0.3, 0.5$).

3 Results

After completing the comprehensive experiments mentioned above, this study primarily presents results in two parts: the model segmentation performance, and the robustness of the model under noise-corrupted conditions.

3.1 Model Segmentation Performance on Original TCS Images

The section mainly discusses the model segmentation agreement, model stability, and prediction time efficiency. Furthermore, extensive quantitative comparisons were conducted.

Model Segmentation Agreement and Model Stability

From quantitative perspectives, **Table 1A** summarizes the quantitative comparisons among the U-Net, U-Net+++, and nnU-Net networks for midbrain segmentation agreement (in terms of averaged DICE, IoU, and HD) and model stability (in terms of SD) on the testing dataset (relevant results for the

299 validation set can be found in **Supplementary Table 1A**). It showcases the detailed results of the 3
 300 randomization seeds under 5-fold cross validation. It is worth noting that the nnU-Net model yielded
 301 the best segmentation agreement, achieving the top-ranked averaged scores of DICE (0.910, 0.907,
 302 0.910 in randomization seed 1 to 3, respectively), IoU (0.838, 0.832, 0.837 in randomization seed 1 to
 303 3, respectively), HD (2.727 mm, 2.801 mm, 2.852 mm in randomization seed 1 to 3, respectively).
 304 With respect to model stability, the nnU-Net model outperformed the other two comparing networks,
 305 achieving the top-ranked values of SD in DICE (0.001 in all the three randomization seeds), IoU (0.002,
 306 0.001, 0.002 in randomization seed 1 to 3, respectively), and HD (0.095 mm, 0.065 mm, 0.206 mm in
 307 randomization seed 1 to 3, respectively).

308 **Table 1A.** Results of segmentation agreement and stability of the three comparing neural networks in each of the three randomization
 309 seeds under 5 cross-validation on the testing dataset. Averaged (AVG) scores of the segmentation agreement metrics in DICE, IoU,
 310 and HD are presented for each model in each randomization seed. Standard deviation (SD) was calculated to illuminate the underlying
 311 model stability. The top-ranked scores are bolded.

		Randomization Seed 1			Randomization Seed 2			Randomization Seed 3		
Network	Fold	DICE	IoU	HD (mm)	DICE	IoU	HD (mm)	DICE	IoU	HD (mm)
U-Net	Fold1	0.905	0.829	3.375	0.899	0.820	3.394	0.903	0.825	3.359
	Fold2	0.901	0.823	3.321	0.901	0.823	3.207	0.904	0.826	3.194
	Fold3	0.904	0.828	3.223	0.892	0.808	3.504	0.899	0.818	3.541
	Fold4	0.900	0.822	3.337	0.899	0.820	3.271	0.903	0.824	3.556
	Fold5	0.903	0.826	3.704	0.894	0.811	3.566	0.898	0.817	3.830
	AVG	0.903	0.826	3.392	0.897	0.816	3.389	0.901	0.822	3.496
	SD	0.002	0.003	0.184	0.004	0.006	0.152	0.003	0.004	0.239
U-Net+++	Fold1	0.889	0.806	4.385	0.898	0.817	3.680	0.898	0.817	3.452
	Fold2	0.888	0.806	3.889	0.898	0.817	3.899	0.893	0.811	3.417
	Fold3	0.894	0.813	3.520	0.892	0.808	4.571	0.894	0.811	3.887
	Fold4	0.893	0.811	3.776	0.893	0.810	3.987	0.894	0.810	3.637
	Fold5	0.886	0.802	3.804	0.895	0.813	4.034	0.893	0.810	3.734
	AVG	0.890	0.807	3.875	0.895	0.813	4.034	0.894	0.812	3.626
	SD	0.003	0.004	0.317	0.003	0.004	0.329	0.002	0.003	0.197
nnU-Net	Fold1	0.909	0.836	2.882	0.908	0.833	2.791	0.910	0.837	2.723
	Fold2	0.912	0.841	2.646	0.907	0.831	2.743	0.910	0.836	2.723
	Fold3	0.911	0.838	2.654	0.908	0.833	2.832	0.913	0.841	2.669
	Fold4	0.909	0.836	2.731	0.908	0.832	2.896	0.910	0.837	3.026
	Fold5	0.911	0.839	2.721	0.906	0.830	2.743	0.909	0.835	3.120
	AVG	0.910	0.838	2.727	0.907	0.832	2.801	0.910	0.837	2.852
	SD	0.001	0.002	0.095	0.001	0.001	0.065	0.001	0.002	0.206

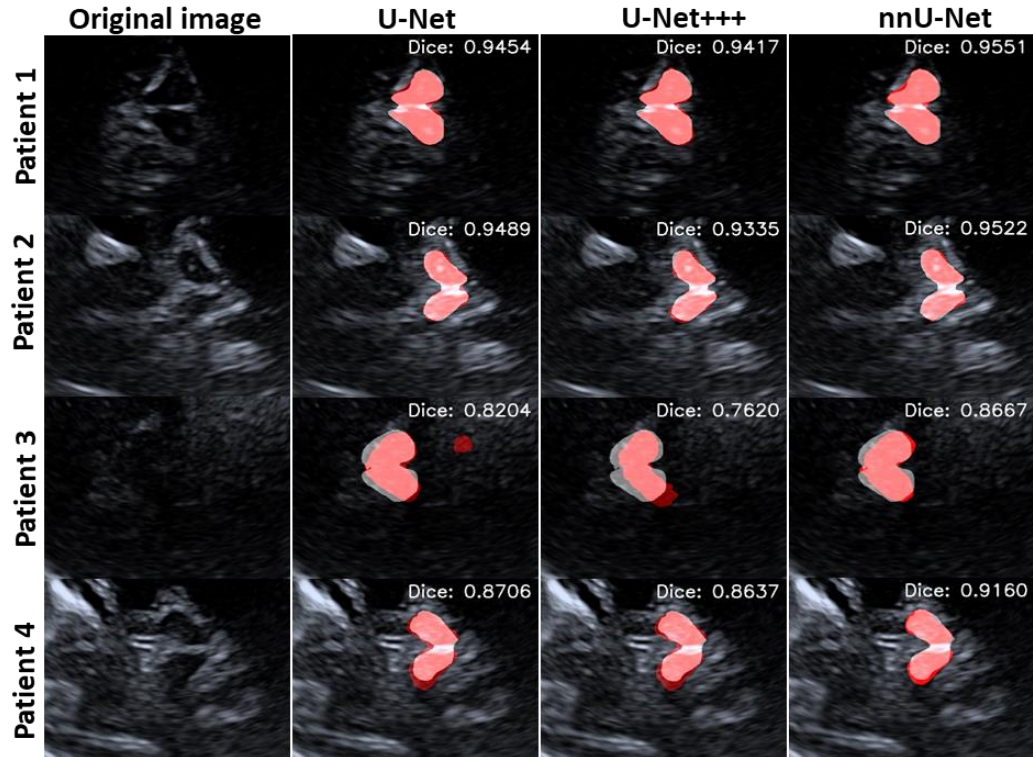
312

313 **Table 1B.** Results of segmentation agreement (in terms of DICE, IoU, and HD) and stability (in terms of SD) of the three deep-
 314 learning networks across all the three randomization seeds under 5-fold cross validation on the testing dataset. The top-ranked scores
 315 are bolded.

Network	DICE (SD)	IoU (SD)	HD (mm) (SD)
U-Net	0.900 (0.004)	0.821 (0.006)	3.425 (0.187)
U-Net+++	0.893 (0.003)	0.811 (0.004)	3.845 (0.318)
nnU-Net	0.910 (0.002)	0.836 (0.003)	2.793 (0.137)

316 **Table 1B** presents the averaged results of the three comparing networks across all the
 317 randomization seeds under 5-fold cross-validation on the testing dataset (relevant results for the
 318 validation set can be found in **Supplementary Table 1B**). The nnU-Net network outperformed U-Net
 319 and U-Net+++ networks, achieving the best segmentation agreement in terms of DICE (0.910), IoU
 320 (0.836), and HD (2.793 mm), and model stability in terms of SD in DICE (0.002), IoU (0.003), and
 321 HD (0.137 mm).

322 From qualitative aspects, **Figure 2** shows the segmentation agreement in terms of the DICE
 323 score in 4 representative patients, with the red regions representing the predicted segments and the
 324 white regions indicating the ground-truth annotation. Notably, the nnU-Net network consistently
 325 outperformed the other two networks in these 4 patients. By contrast, the U-Net and U-Net+++
 326 networks were only capable of performing comparable segmentation agreements compared to the nnU-
 327 Net model in Patient 1 and 2; while their capabilities were considerably under-performed in Patient 3
 328 and 4.



329 Fig 2. Qualitative visualization of the segmentation agreement in terms of the averaged DICE score in four representative patients, with the red region
 330 representing the predicted segments and the white region indicating the ground-truth annotation.
 331

332

333 Prediction Time Efficiency

334 **Table 2** presents the time efficiency of the three studied networks, in terms of training and prediction
 335 time in the training and testing sets, respectively. It is worth noting that the nnU-Net network was
 336 determined to be the most efficient model, requiring the least amount of time (1.456 second) in
 337 generating the predictions in the testing set, despite that it took the longest duration of time during
 338 model training (114 minute).

339 **Table 2.** Results of the training time and prediction time of the three studied networks. The shortest time are bolded.

Network	Training Time	Prediction Time
U-Net	34 minutes	2.427 second
U-Net+++	36 minutes	3.083 second
nnU-Net	114 minutes	1.456 second

340 Based on the above evaluations, the nnU-Net demonstrated the best-performing performance
 341 among the three neural networks in aspects of segmentation agreement, model stability, and time
 342 efficiency. Details of the segmentation agreement of the three comparing networks for each testing

343 case under randomization seed-1 are presented in DICE and visualized in **Supplementary Figure 2A-**
344 **2C**.

345 3.2 Model Robustness Assessment on Noise-corrupted TCS Images

346 The section primarily discusses the impact of noise addition on the performance of the benchmarked
347 model and the impact of incremental noise addition on the performance of the nnU-Net Model.

348 Impact of Noise Addition on Performance of the Benchmarked Model

349 **Table 3** summarizes the segmentation agreement of the three studied networks in both original and
350 noise-corrupted ultrasound images under mild degree of noise (i.e., noise level of SNR=0.95 for the
351 salt-and-pepper noise, and $\sigma=0.1$ for the Gaussian noise). The nnU-Net was the best-performing model,
352 demonstrating a dramatically greater capacity of midbrain segmentation in both noise-corrupted
353 images compared to the other two comparing networks, and the underlying segmentation agreements
354 were approximate to that on the original images. By contrast, the U-Net+++ network was the most
355 underperforming model in the salt-and-pepper noise corrupted images (SNR=0.95: DICE=0.032,
356 IoU=0.017, HD=37.827-mm) and in the Gaussian noise corrupted images ($\sigma=0.1$: DICE=0.157,
357 IoU=0.089, HD=41.802-mm); and the underlying segmentation agreements on the noise-corrupted
358 images were far lower than those of the nnU-Net model, incurring up to 28 times degradation in DICE,
359 48 folds in IoU and 14 times in HD.

360 Of note, it is intriguing to point out that both the U-Net and U-Net+++ segmentation models
361 were highly sensitive to the two studied types of ultrasound imaging noise, even at a low intensity level
362 of noise corruption, resulting in tremendous drops in the segmentation agreements. On the contrary,
363 the segmentation power of the nnUNet network remained robust in the same noise addition setting.

364 **Figure 3 and Figure 4** visualize the segmentation capacity of the three neural networks in three
365 representative patients when the images were corrupted with the salt-and-pepper noise (SNR=0.95)
366 and the Gaussian noise ($\sigma=0.1$), respectively. The results were found consistent with those illustrated
367 in **Table 3**.

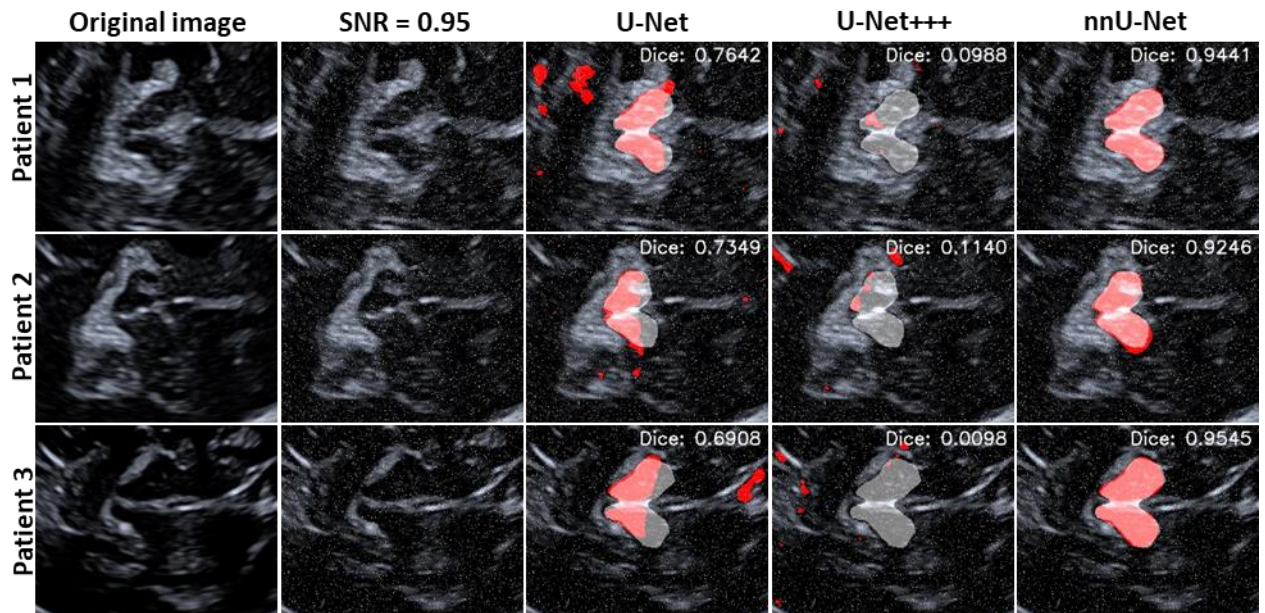
368 **Table 3.** Segmentation agreement of the three studied networks (in terms of averaged DICE, IoU and HD) in both original and noise-
369 corrupted ultrasound images (noise level of SNR=0.95 for the salt-and-pepper noise, and $\sigma=0.1$ for the Gaussian noise), across the
370 three randomization seeds under 5-fold cross validation on the testing set.

Network	Original image			Salt-and-pepper noise corrupted image (SNR = 0.95)			Gaussian noise corrupted image ($\sigma = 0.1$)		
	DICE	IoU	HD (mm)	DICE	IoU	HD (mm)	DICE	IoU	HD (mm)
nnU-Net	0.909	0.836	2.793	0.904	0.827	2.941	0.906	0.830	2.967
U-Net	0.900	0.821	3.425	0.389	0.268	41.946	0.535	0.389	38.241
U-Net+++	0.893	0.811	3.845	0.032	0.017	37.837	0.157	0.089	41.802

371

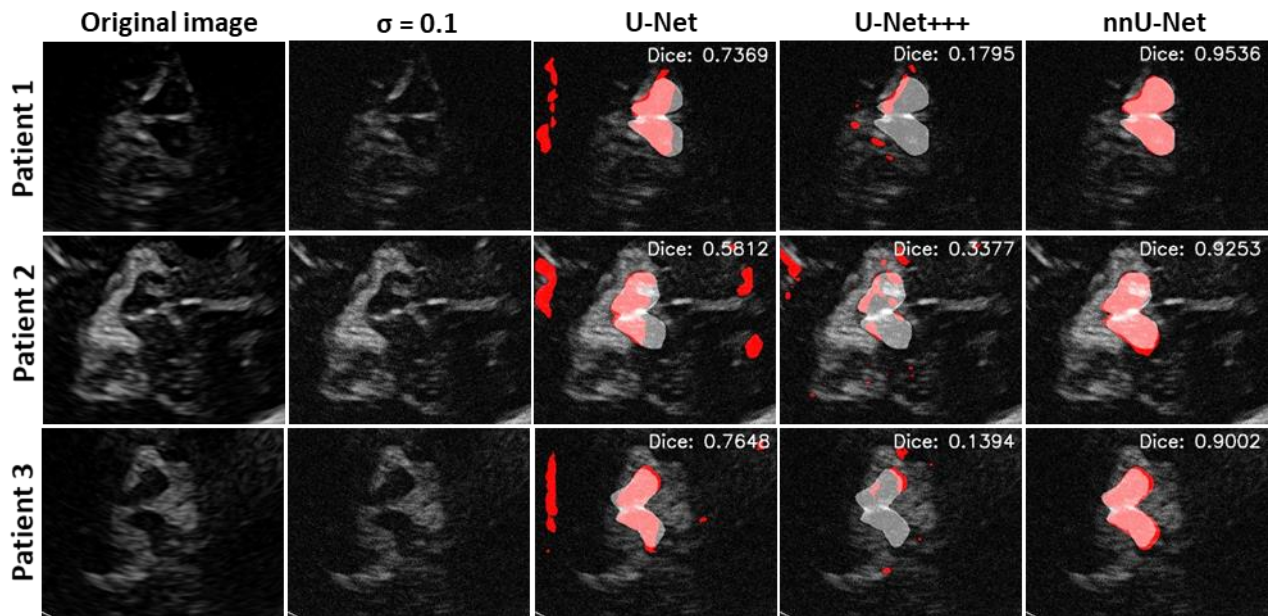
372 The nnU-Net network illuminated the greatest segmentation agreement and robustness both
373 qualitatively in terms of the matching between the red region (i.e., the predicted segment) and the white
374 region (i.e., the ground-truth annotation), and quantitatively in aspects of averaged DICE ranging from
375 0.900 to 0.955 (**Figure 3 and Figure 4**). Consistent with the findings in **Table 3**, the UNet+++
376 exhibited the worst segmentation agreement and robustness on the noise-corrupted images in both

377 qualitative and quantitative perspectives (**Figure 3 and Figure 4**), with averaged DICE scores ranging
 378 from 0.010 to 0.338. It may also be worth noting that the predicted segments in the U-Net and U-
 379 Net+++ models were of small size and were sparsely distributed inside and outside the midbrain region
 380 on the noise-corrupted images.



381

382 **Fig. 3.** Qualitative visualization of the segmentation agreement in terms of the averaged DICE score in three representative patients, on the salt-and-
 383 pepper noise corrupted images (SNR=0.95). The red region representing the predicted segments and the white region indicating the
 384 ground-truth annotation.



385

386 **Fig. 4.** Qualitative visualization of the segmentation agreement in terms of the averaged DICE score in three representative patients, on
 387 the Gaussian noise corrupted images ($\sigma=0.1$). The red region representing the predicted segments and the white region indicating the
 388 ground-truth annotation.

389

390

391

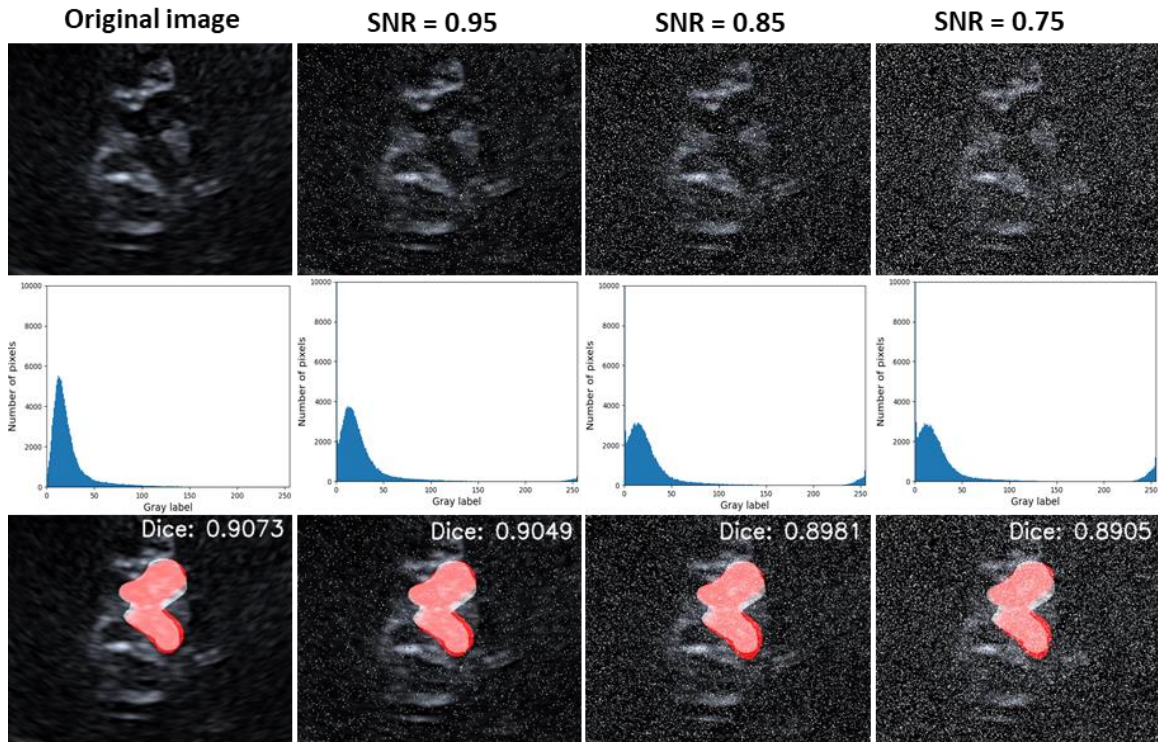
392 Impact of Incremental Noise Addition on Performance of nnU-Net Model

393 **Table 4** summarizes the segmentation power of the benchmarked nnU-Net model on the original
394 images and corrupted-images of varying degrees of noise intensity (SNR ranged from 0.75 to 0.95
395 for the salt-and-pepper noise; σ values ranged from 0.10 to 0.50 for the Gaussian noise); the results were
396 averaged across the three randomization seeds under 5-fold cross validation on the testing set.

397 **Table 4** summarizes the segmentation power of the benchmarked nnU-Net model on the original images and corrupted-images of
398 varying degree of noise intensity; the results were averaged across the three randomization seeds under 5-fold cross validation on the
399 testing set.

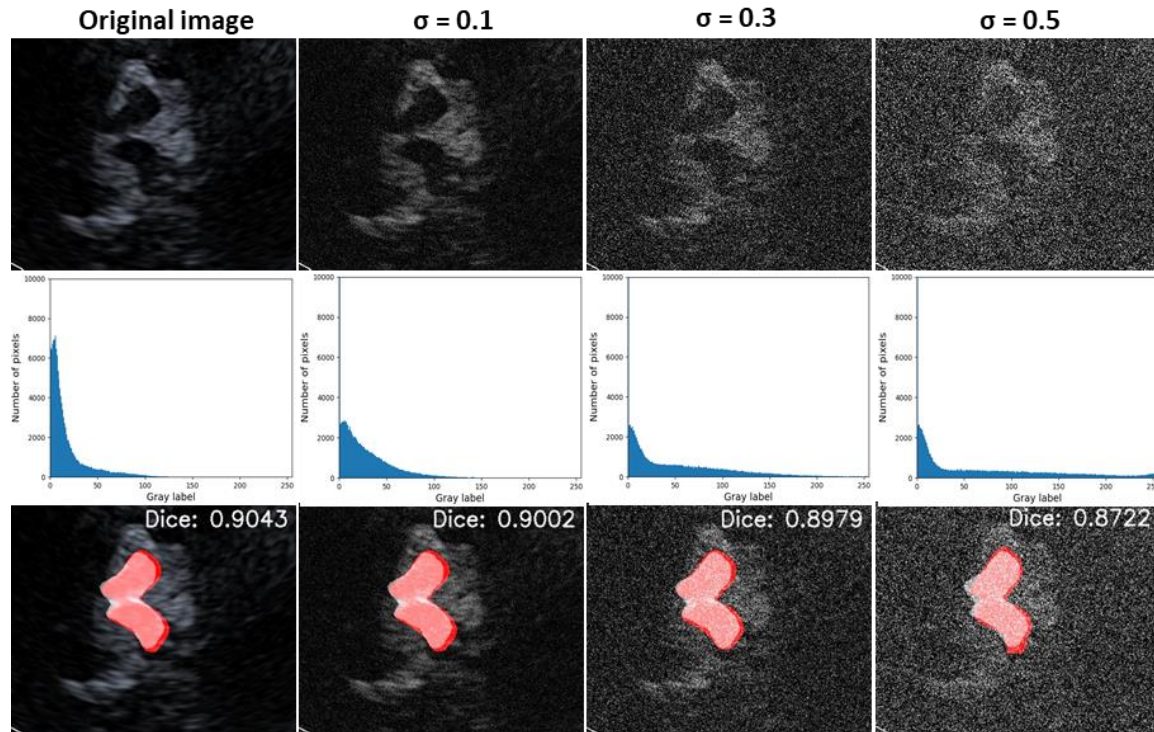
Noise level	Segmentation Agreement		
	DICE	IoU	HD (mm)
Original Image	0.909	0.836	2.793
SNR = 0.95	0.904	0.827	2.941
SNR = 0.85	0.888	0.805	3.552
SNR = 0.75	0.882	0.796	4.040
$\sigma = 0.10$	0.906	0.830	2.967
$\sigma = 0.30$	0.891	0.809	3.501
$\sigma = 0.50$	0.874	0.784	4.183

400 **Figure 5** and **Figure 6** demonstrate a representative case for visually illuminating the impact
401 of varying intensity levels of the salt-and-pepper noise and Gaussian noise, respectively, on both the
402 image appearance and the segmentation agreement of the nnU-Net network; additionally, histograms
403 are presented for each of the studied noise intensity levels.



404 **Fig. 5.** demonstrate a representative case for visually illuminating the impact of varying intensity levels of the salt-and-pepper noise
405 (SNR= 0.75, 0.85, and 0.95), on both the image appearance (first row), the underlying histograms (middle row), and the segmentation
406 agreement of the nnU-Net network (bottom row), where the red regions and white regions representing the predicted and ground-truth
407 annotation, respectively.
408

409 Two important findings were noted. Firstly, additions of imaging noise and increasing degree
 410 of noise corruption to the original ultrasound images rendered a gradual drop of the segmentation
 411 power of the nnU-Net network in aspects of DICE, IoU and HD, irrespective of the types of imaging
 412 noise studied (**Table 4, Figure 5 and Figure 6**). Secondly, in spite of the reduction in segmentation
 413 agreement under the increasing intensity of imaging noise, the nnU-Net network was still capable of
 414 producing satisfactory segmentation power, even at the highest level of noise intensity (**Table 4, Figure
 415 5 and Figure 6**). For instance, in the setting SNR of 0.75 resulting from the addition of the salt-and-
 416 pepper noise, the averaged scores of DICE, IoU and HD were 0.882, 0.796, and 4.040-mm respectively,
 417 corresponding to approximately 2%, 4.5%, and 45% reductions compared to those on the original
 418 images (**Table 4**).



419
 420 **Fig. 6.** demonstrate a representative case for visually illuminating the impact of varying intensity levels of the salt-and-pepper noise (σ
 421 = 0.1, 0.3, and 0.5), on both the image appearance (first row), the underlying histograms (middle row), and the segmentation agreement
 422 of the nnU-Net network (bottom row), where the red regions and white regions representing the predicted and ground-truth annotation,
 423 respectively.

424 4 Discussion

425 Midbrain segmentation is a key initial step in TCS-based PD grading for purposes of diagnosis, severity
 426 assessment, and treatment response surveillance. Yet, the current practice of midbrain segmentation is
 427 a manual, time-consuming and subjective procedure, putting a huge demand for automated, effective
 428 and objective alternatives in the contemporary paradigm of AI. Despite the superiority of TCS in
 429 offering a non-invasive, radiation-free, rapid, low cost, highly accessible, high patient compliance,
 430 easy-to-operate tool in a stark contrast to other imaging modalities [8], relevant studies on AI for
 431 midbrain segmentation are desperately wanting. The only relevant research in the literature by Milletari
 432 et al. was deficient in providing sufficient value of their developed models from clinical
 433 implementation perspectives, on top of the inadequacy of sample size of 34 subjects [32]. Given the
 434 challenges outlined above, our objective was to develop deep learning methods and benchmark best-
 435 performing state-of-the-art U-Net based neural networks using a substantial cohort of 584 subjects.

436 The evaluation focused on segmentation agreement, model stability, time efficiency, and model
437 robustness. Results indicated that the nnU-Net model outperformed the other two networks in all these
438 evaluating aspects, potentially providing the community with a favorable and automated alternative for
439 midbrain segmentation towards TCS-based PD assessment in the future.

440 Aligned with the present study, all the abovementioned studies reflected the importance of
441 TCS-based automated midbrain segmentation over the past decade, which are all valuable to the
442 research community. However, there exists a growing demand on auto-segmentation works for an
443 adoption of large sample size and comprehensive model evaluation beyond segmentation accuracy. In
444 contrast to these previous works, the present study benchmarked the state-of-the-art nnU-Net network
445 as the best-performing midbrain segmentation model in aspects of segmentation agreement (in DICE,
446 IoU, HD, and AD) (**Table 1A-1B, Figure 2**), model stability (**Table 1A-1B**), time efficiency (**Table**
447 **2**), as well as model robustness against noise-corrupted TCS images (**Table 3-4, Figure 5-6**), when
448 compared to the U-Net and U-Net+++ models; these findings were obtained using a larger cohort of
449 584 subjects. For segmentation agreement, the present nnU-Net model yielded a higher averaged scores
450 of DICE (91.0%), IoU (83.6%), HD (2.793 mm), and AD (79%) as shown in **Table 1A-1B**, which
451 were superior over the abovementioned previous works [32]. For model stability, the nnU-Net model
452 also achieved the top-ranked values of SD in DICE, IoU, HD, and AD, compared to the U-Net and U-
453 Net+++ models (**Table 1A-1B**). For time efficiency, the nnU-Net network required the least amount
454 of time (1.456 second) in generating the predictions in the testing set (**Table 2**). For model robustness,
455 the nnU-Net demonstrated a dramatically greater capacity of midbrain segmentation in images
456 corrupted by the salt-and-pepper noise (SNR=0.95) and the Gaussian noise ($\sigma=0.10$), and the
457 underlying segmentation agreements were approximate to that on the original images (**Table 3**).
458 Notably, such assessments of model robustness against noise corruption have not been previously
459 revealed in the current body of literature. The overall performance of the benchmarked nnU-Net
460 network outperformed the traditional and AI-based approaches reported in previous studies [16, 20,
461 40, 41]. Based on the findings in the present study, it is believed that the above superiorities of the
462 nnU-Net network would further bring added values to the community for midbrain segmentation, as
463 well as TCS-based PD assessment, in the future.

464 For our clinical objectives, we have meticulously selected these three models (U-Net, U-
465 Net+++ and nnU-Net) to conduct a thorough and comprehensive benchmarking study [18-20, 42, 43].
466 First, for our specific clinical task, we selected the U-Net model for its classic, well-established, and
467 effective performance in supervised segmentation tasks, as acknowledged in literature. Given our fully
468 supervised segmentation task and the meticulous annotation by experienced clinicians with over a
469 decade of practice, we possess a valuable dataset. Furthermore, the U-Net+++ model was chosen for
470 its extensive skip connections and deep supervision, enhancing supervision quality. Additionally, the
471 nnU-Net integrates the complete training pipeline, from pre-processing to post-processing, based on
472 the U-Net architecture. This choice optimizes the training process for our benchmark tasks. Secondly,
473 the three selected networks are relatively more mature and readily applicable in clinics than some
474 emerging networks, considering the length of time they have been developed, evaluated, and applied
475 for a wide spectrum of segmentation tasks over the past decade. We believe that performing a deeper
476 and more comprehensive analysis of their performance via the submitted work, from the perspectives
477 of segmentation agreement, model stability and robustness analysis etc, would provide enlightening
478 insights for their applications when it comes to clinical implementation.

479 Intriguingly, results of this study indicated that although both the U-Net and U-Net+++ models
480 achieved good segmentation agreements even they were inferior to the nnU-Net model (**Table 1A-1B**),
481 their segmentation capabilities dramatically degraded when it comes to the images corrupted by the

482 salt-and-pepper noise (SNR=0.95) and the Gaussian noise ($\sigma=0.1$) even at a low intensity level of noise
483 addition (**Table 3**). Conversely, the nnU-Net network not only outperformed on TCS images at the
484 same level of noise corruption (**Table 3**), its segmentation power remained even at greater degree of
485 noise corruption (SNR=0.85, 0.75; $\sigma=0.30, 0.50$), as illustrated in **Table 4**. We speculated that the
486 exceptional performance of the nnU-Net model in midbrain segmentation may be attributed to its
487 unique and comprehensive training architecture, distinguishing it from the U-Net and U-Net+++
488 models. Although the network structure of the nnU-Net was not modified, its adaptive nature enables
489 nnU-Net to automatically adjust the underlying hyper-parameters, such as the exact patch size, batch
490 size, and inference settings, to adapt to the input dataset, while integrating image preprocessing
491 strategies and image augmentation techniques seamlessly into the training workflow, crucially
492 contribute to nnU-Net's superior performance.

493 In addition, it is worth noting that the present work contained multiple layers of strengths. First,
494 we enrolled a remarkably larger cohort of data (n=584) compared to the previous studies (n=34-130).
495 The inclusion of a larger sample size enhanced the validity of the findings from this study. Second, the
496 findings were presented in a wide spectrum of application scopes beyond segmentation agreement.
497 Particularly, the noise-corruption test was applied for the first time on TCS-based midbrain auto-
498 segmentation for securing model robustness; this would hopefully provide the community with more
499 valuable insights into clinical contextualization of the benchmarked model from application point-of-
500 view. Third, three independent randomization seeds for train-test splits and 5-fold cross-validation
501 were applied in this study in the hope of generating more reliable findings of the developed networks;
502 these techniques have been proven to be effective in previous studies [34, 35, 44], and is believed to
503 enhance the validity of the findings in this study.

504 From a public health standpoint, TCS exhibits a great application potential for cost-effective
505 population-wide screening and frequent surveillance for timely PD management owing to its
506 advantages of being radiation-free, high accessibility and high affordability [9, 45]. Several research
507 groups have also explored the utility of methods based on other imaging modalities as valuable tools
508 for early diagnosis of PD [46, 47]. However, such as MRI due to their low accessibility and
509 affordability to the general public, particularly in developing or under-developed countries [48-50], as
510 well as Positron Emission Tomography and Single-Photon Emission Computed Tomography due to
511 their low accessibility, low affordability, and the potential radiation hazards [51, 52]. As such, we hope
512 that results of this study would provide enlightening insights and stimulate researchers in the field to
513 steer more focuses on AI-assisted TCS-based PD assessment in the long run, providing a cost-effective
514 alternative to both the sufferers and medical practitioners in PD management in the smart-ageing era.

515 In spite of the encouraging results, the study presents several shortcomings that worth further
516 investigations in the future. First, results of this study were generated with the dataset collected from a
517 single institution and the same scanner vendor, which may limit the generalizability of the
518 benchmarked model in real-world clinical settings. A multi-center study would improve and validate
519 model generalizability, so a multi-center study with TCS data acquired from different vendors are
520 warranted in the future in order to further validate findings of this work. Second, although the sample
521 size of 584 in this study was considerably larger than previous studies where the sample sizes ranged
522 from 40 to 130, further investigations using a larger cohort are preferred in the context of deep learning
523 to improve and validate model generalizability. Third, for the robustness analysis, we chose salt-and-
524 pepper noise and Gaussian noise in the study, the two most common types of noise in medical imaging,
525 effectively simulating the majority of ultrasound noise scenarios. Additionally, other speckle noise and
526 poisson noise are two types of noises that medical professionals may encounter. We anticipate further
527 detailed discussions on noise impact in future research. Finally, due to the retrospective nature of this

528 study, although this study recruited highly experienced physicians (>10 years) in ultrasound imaging
529 for ground-truth annotation generation, there may exist intra- or inter-rater variabilities in the dataset.
530 Therefore, addressing this issue in a prospective study is warranted in the future to further facilitate
531 widespread adoption of the developed model in clinical practice.

532

533 **5 Conclusions**

534 The nnU-Net model was benchmarked to be the best-performing model for midbrain segmentation in
535 terms of segmentation agreement, model stability, prediction time efficiency, and model robustness,
536 compared to the U-Net and U-Net+++ models. Despite the discussed limitations, results of this study
537 presented multiple layers of superiorities in comparison to previous works, in terms of sample size,
538 comprehensive scope of model evaluation, as well as the enhanced reliability owing to the adoption of
539 multiple independent randomization seeds together with 5-fold cross-validation. Results of this study
540 would provide enlightening insights and stimulate researchers in the field to steer more focuses on
541 cost-effective AI-assisted TCS-based PD assessment in the long run, potentially benefiting millions of
542 sufferers as well as medical practitioners worldwide for PD management in the smart-ageing era.
543 Moving forward, a multi-center multi-vendor study is warranted in the future when it comes to clinical
544 implementation.

545

546 **6 Acknowledgement**

547 This study was partially supported by Research Institute for Smart Ageing of The Hong Kong
548 Polytechnic University (1-CD5B).

549

550 **7 Ethical statement**

551 Ethical approval for this study was obtained from the Human Subject Ethics Sub-committee (HSESC)
552 of the Hong Kong Polytechnic University (HSEARS20231102004).

553

554 **8 Declaration of interests**

555 All the authors declare that they have no known competing financial interests or personal relationships
556 that could have appeared to influence the work reported in this paper.

557 **References**

- 558 [1] N. Maserejian, L. Vinikoor-Imler, A. Dilley, Estimation of the 2020 Global Population of
559 Parkinson's Disease (PD), *Movement Disord*, 35 (2020) S79-S80.
- 560 [2] G.B.D.N. Collaborators, Global, regional, and national burden of neurological disorders, 1990-
561 2016: a systematic analysis for the Global Burden of Disease Study 2016, *Lancet Neurol*, 18 (2019)
562 459-480.
- 563 [3] B.R. Bloem, M.S. Okun, C. Klein, Parkinson's disease, *Lancet*, 397 (2021) 2284-2303.
- 564 [4] S. Behnke, D. Berg, M. Naumann, G. Becker, Differentiation of Parkinson's disease and atypical
565 parkinsonian syndromes by transcranial ultrasound, *J Neurol Neurosurg Psychiatry*, 76 (2005) 423-
566 425.
- 567 [5] L. Zecca, D. Berg, T. Arzberger, P. Ruprecht, W.D. Rausch, M. Musicco, D. Tampellini, P.
568 Riederer, M. Gerlach, G. Becker, In vivo detection of iron and neuromelanin by transcranial
569 sonography: A new of substantia approach for early detection nigra damage, *Movement Disord*, 20
570 (2005) 1278-1285.
- 571 [6] Y.J. Bae, J.M. Kim, C.H. Sohn, J.H. Choi, B.S. Choi, Y.S. Song, Y. Nam, S.J. Cho, B. Jeon, J.H.
572 Kim, Imaging the Substantia Nigra in Parkinson Disease and Other Parkinsonian Syndromes,
573 *Radiology*, 300 (2021) 260-278.
- 574 [7] D. Berg, K. Seppi, S. Behnke, I. Liepelt, K. Schweitzer, H. Stockner, F. Wollenweber, A. Gaenslen,
575 P. Mahlknecht, J. Spiegel, J. Godau, H. Huber, K. Srulijes, S. Kiechl, M. Bentele, A. Gasperi, T.
576 Schubert, T. Hiry, M. Probst, V. Schneider, J. Klenk, M. Sawires, J. Willeit, W. Maetzler, K.
577 Fassbender, T. Gasser, W. Poewe, Enlarged Substantia Nigra Hyperechogenicity and Risk for
578 Parkinson Disease A 37-Month 3-Center Study of 1847 Older Persons, *Arch Neurol-Chicago*, 68
579 (2011) 932-937.
- 580 [8] Y.L. Mei, J. Yang, Z.R. Wu, Y. Yang, Y.M. Xu, Transcranial Sonography of the Substantia Nigra
581 for the Differential Diagnosis of Parkinson's Disease and Other Movement Disorders: A Meta-
582 Analysis, *Parkinsons Dis-Us*, 2021 (2021).
- 583 [9] P. Bartova, D. Skoloudik, M. Bar, P. Ressner, P. Hlustik, R. Herzig, P. Kanovsky, Transcranial
584 Sonography in Movement Disorders, *Biomed Pap*, 152 (2008) 251-258.
- 585 [10] D. Skoloudik, T. Fadrna, P. Bartova, K. Langova, P. Ressner, O. Zapletalova, P. Hlustik, R.
586 Herzig, P. Kannovsky, Reproducibility of sonographic measurement of the substantia nigra,
587 *Ultrasound Med Biol*, 33 (2007) 1347-1352.
- 588 [11] A. Berardelli, G.K. Wenning, A. Antonini, D. Berg, B.R. Bloem, V. Bonifati, D. Brooks, D.J.
589 Burn, C. Colosimo, A. Fanciulli, J. Ferreira, T. Gasser, F. Grandas, P. Kanovsky, V. Kostic, J.
590 Kulisevsky, W. Oertel, W. Poewe, J.P. Reese, M. Relja, E. Ruzicka, A. Schrag, K. Seppi, P. Taba, M.
591 Vidailhet, EFNS/MDS-ES recommendations for the diagnosis of Parkinson's disease, *Eur J Neurol*, 20
592 (2013) 16-+.
- 593 [12] D. Monaco, D. Berg, A. Thomas, V. Di Stefano, F. Barbone, M. Vitale, C. Ferrante, L. Bonanni,
594 M. Di Nicola, T. Garzarella, L.P. Marchionno, G. Malferrari, R. Di Mascio, M. Onofrj, R. Franciotti,
595 The predictive power of transcranial sonography in movement disorders: a longitudinal cohort study,
596 *Neurol Sci*, 39 (2018) 1887-1894.
- 597 [13] S. Zhu, Y.X. Wang, Y.Y. Jiang, R.X. Gu, M. Zhong, X. Jiang, B. Shen, J. Zhu, J. Yan, Y. Pan, L.
598 Zhang, Clinical Features in Parkinson's Disease Patients with Hyperechogenicity in Substantia Nigra:
599 A Cross-Sectional Study, *Neuropsych Dis Treat*, 18 (2022) 1593-1601.

- 600 [14] P. Singh, A neutrosophic-entropy based clustering algorithm (NEBCA) with HSV color system:
601 A special application in segmentation of Parkinson's disease (PD) MR images, *Comput Meth Prog Bio*,
602 189 (2020).
- 603 [15] D. Basukala, R. Mukundan, A. Lim, M.A. Hurrell, R.J. Keenan, J.C. Dalrymple-Alford, T.J.
604 Anderson, D.J. Myall, T.R. Melzer, Automated segmentation of substantia nigra and red nucleus using
605 quantitative susceptibility mapping images: Application to Parkinson's disease, *Comput Electr Eng*, 91
606 (2021).
- 607 [16] S.A. Ahmadi, M. Baust, A. Karamalis, A. Plate, K. Boetzel, T. Klein, N. Navab, Midbrain
608 Segmentation in Transcranial 3D Ultrasound for Parkinson Diagnosis, *Medical Image Computing and*
609 *Computer-Assisted Intervention, Miccai 2011, Pt Iii*, 6893 (2011) 362-+.
- 610 [17] A. Sakalauskas, A. Lukosevicius, K. Lauckaite, D. Jegelevicius, S. Rutkauskas, Automated
611 segmentation of transcranial sonographic images in the diagnostics of Parkinson's disease, *Ultrasonics*,
612 53 (2013) 111-121.
- 613 [18] O. Ronneberger, P. Fischer, T. Brox, U-Net: Convolutional Networks for Biomedical Image
614 Segmentation, *Lect Notes Comput Sc*, 9351 (2015) 234-241.
- 615 [19] H.M. Huang, L.F. Lin, R.F. Tong, H.J. Hu, Q.W. Zhang, Y. Iwamoto, X.H. Han, Y.W. Chen, J.
616 Wu, Unet 3+: A Full-Scale Connected Unet for Medical Image Segmentation, *Int Conf Acoust Spee*,
617 (2020) 1055-1059.
- 618 [20] F. Isensee, P.F. Jaeger, S.A.A. Kohl, J. Petersen, K.H. Maier-Hein, nnU-Net: a self-configuring
619 method for deep learning-based biomedical image segmentation, *Nat Methods*, 18 (2021) 203-211.
- 620 [21] Y. Hiasa, Y. Otake, M. Takao, T. Ogawa, N. Sugano, Y. Sato, Automated Muscle Segmentation
621 from Clinical CT Using Bayesian U-Net for Personalized Musculoskeletal Modeling, *Ieee T Med*
622 *Imaging*, 39 (2020) 1030-1040.
- 623 [22] K.B. Chen, Y. Xuan, A.J. Lin, S.H. Guo, Lung computed tomography image segmentation based
624 on U-Net network fused with dilated convolution, *Comput Meth Prog Bio*, 207 (2021).
- 625 [23] S.Q. Li, J.Y. Zheng, D.J. Li, Precise segmentation of non-enhanced computed tomography in
626 patients with ischemic stroke based on multi-scale U-Net deep network model, *Comput Meth Prog*
627 *Bio*, 208 (2021).
- 628 [24] H. Dong, G. Yang, F.D. Liu, Y.H. Mo, Y.K. Guo, Automatic Brain Tumor Detection and
629 Segmentation Using U-Net Based Fully Convolutional Networks, *Comm Com Inf Sc*, 723 (2017) 506-
630 517.
- 631 [25] H.F. Cui, Y.W. Chang, L. Jiang, Y. Xia, Y.N. Zhang, Multiscale attention guided U-Net
632 architecture for cardiac segmentation in short-axis MRI images, *Comput Meth Prog Bio*, 206 (2021).
- 633 [26] A.M.G. Allah, A.M. Sarhan, N.M. Elshennawy, Edge U-Net: Brain tumor segmentation using
634 MRI based on deep U-Net model with boundary information, *Expert Syst Appl*, 213 (2023).
- 635 [27] J. Yang, M. Faraji, A. Basu, Robust segmentation of arterial walls in intravascular ultrasound
636 images using Dual Path U-Net, *Ultrasonics*, 96 (2019) 24-33.
- 637 [28] M. Amiri, R. Brooks, B. Behboodi, H. Rivaz, Two-stage ultrasound image segmentation using U-
638 Net and test time augmentation, *Int J Comput Ass Rad*, 15 (2020) 981-988.
- 639 [29] M. Amiri, R. Brooks, H. Rivaz, Fine-Tuning U-Net for Ultrasound Image Segmentation: Different
640 Layers, Different Outcomes, *Ieee T Ultrason Ferr*, 67 (2020) 2510-2518.

- 641 [30] G.P. Chen, L. Li, Y. Dai, J.X. Zhang, M.H. Yap, AAU-Net: An Adaptive Attention U-Net for
642 Breast Lesions Segmentation in Ultrasound Images, *Ieee T Med Imaging*, 42 (2023) 1289-1300.
- 643 [31] A. Sakalauskas, K. Lauckaite, A. Lukosevicius, D. Rastenyte, Computer-Aided Segmentation of
644 the Mid-Brain in Trans-Cranial Ultrasound Images, *Ultrasound in Medicine and Biology*, 42 (2016)
645 322-332.
- 646 [32] F. Milletari, S.A. Ahmadi, C. Kroll, A. Plate, V. Rozanski, J. Maiostre, J. Levin, O. Dietrich, B.
647 Ertl-Wagner, K. Botzel, N. Navab, Hough-CNN: Deep learning for segmentation of deep brain regions
648 in MRI and ultrasound, *Comput Vis Image Und*, 164 (2017) 92-102.
- 649 [33] M. Weinreich, J.J. Chudow, B. Weinreich, T. Krumerman, T. Nag, K. Rahgozar, E. Shulman, J.
650 Fisher, K.J. Ferrick, Development of an Artificially Intelligent Mobile Phone Application to Identify
651 Cardiac Devices on Chest Radiography, *Jacc-Clin Electrophys*, 5 (2019) 1094-1095.
- 652 [34] S.K. Lam, Y. Zhang, J. Zhang, B. Li, J.C. Sun, C.Y. Liu, P.H. Chou, X. Teng, Z.R. Ma, R.Y. Ni,
653 T. Zhou, T. Peng, H.N. Xiao, T. Li, G. Ren, A.L. Cheung, F.K. Lee, C.W. Yip, K.H. Au, V.H. Lee,
654 A.T. Chang, L.W. Chan, J. Cai, Multi-Organ Omics-Based Prediction for Adaptive Radiation Therapy
655 Eligibility in Nasopharyngeal Carcinoma Patients Undergoing Concurrent Chemoradiotherapy, *Front
656 Oncol*, 11 (2021) 792024.
- 657 [35] J. Zhang, S.K. Lam, X.Z. Teng, Z.R. Ma, X.Y. Han, Y.P. Zhang, A.L.Y. Cheung, T.C. Chau,
658 S.C.Y. Ng, F.K.H. Lee, K.H. Au, C.W.Y. Yip, V.H.F. Lee, Y. Han, J. Cai, Radiomic feature
659 repeatability and its impact on prognostic model generalizability: A multi-institutional study on
660 nasopharyngeal carcinoma patients, *Radiother Oncol*, 183 (2023).
- 661 [36] M.V. Sherer, D.N. Lin, S. Elguindi, S. Duke, L.T. Tan, J. Cacicedo, M. Dahele, E.F. Gillespie,
662 Metrics to evaluate the performance of auto-segmentation for radiation treatment planning: A critical
663 review, *Radiother Oncol*, 160 (2021) 185-191.
- 664 [37] K. Harrison, H. Pullen, C. Welsh, O. Oktay, J. Alvarez-Valle, R. Jena, Machine Learning for
665 Auto-Segmentation in Radiotherapy Planning, *Clin Oncol-Uk*, 34 (2022) 74-88.
- 666 [38] A.C. Bovik, Handbook of image and video processing, Academic press 2010.
- 667 [39] T. Peng, C.Y. Tang, Y.Y. Wu, J. Cai, H-SegMed: A Hybrid Method for Prostate Segmentation in
668 TRUS Images via Improved Closed Principal Curve and Improved Enhanced Machine Learning, *Int J
669 Comput Vision*, 130 (2022) 1896-1919.
- 670 [40] C. Gonzalez, K. Gotkowski, A. Bucher, R. Fischbach, I. Kaltenborn, A. Mukhopadhyay, Detecting
671 When Pre-trained nnU-Net Models Fail Silently for Covid-19 Lung Lesion Segmentation, *Medical
672 Image Computing and Computer Assisted Intervention - Miccai 2021, Pt Vii*, 12907 (2021) 304-314.
- 673 [41] L. Huo, X.X. Hu, Q. Xiao, Y.J. Gu, X. Chu, L. Jiang, Segmentation of whole breast and
674 fibroglandular tissue using nnU-Net in dynamic contrast enhanced MR images, *Magn Reson Imaging*,
675 82 (2021) 31-41.
- 676 [42] D. Gut, Z. Tabor, M. Szymkowski, M. Rozynek, I. Kucybala, W. Wojciechowski, Benchmarking
677 of Deep Architectures for Segmentation of Medical Images, *IEEE Trans Med Imaging*, 41 (2022) 3231-
678 3241.
- 679 [43] R. Azad, E.K. Aghdam, A. Rauland, Y. Jia, A.H. Avval, A. Bozorgpour, S. Karimijafarbigloo,
680 J.P. Cohen, E. Adeli, D. Merhof, Medical image segmentation review: The success of u-net, *IEEE
681 Transactions on Pattern Analysis and Machine Intelligence*, (2024).
- 682 [44] T. Fushiki, Estimation of prediction error by using K-fold cross-validation, *Stat Comput*, 21 (2011)
683 137-146.

684 [45] D. Skoloudik, M. Jelinkova, J. Blahuta, P. Cermak, T. Soukup, P. Bartova, K. Langova, R. Herzig,
685 Transcranial Sonography of the Substantia Nigra: Digital Image Analysis, *Am J Neuroradiol*, 35 (2014)
686 2273-2278.

687 [46] N. Amoroso, M. La Rocca, A. Monaco, R. Bellotti, S. Tangaro, Complex networks reveal early
688 MRI markers of Parkinson's disease, *Med Image Anal*, 48 (2018) 12-24.

689 [47] J. Zhang, Mining imaging and clinical data with machine learning approaches for the diagnosis
690 and early detection of Parkinson's disease, *Npj Parkinsons Dis*, 8 (2022).

691 [48] B. Heim, F. Krismer, R. De Marzi, K. Seppi, Magnetic resonance imaging for the diagnosis of
692 Parkinson's disease, *J Neural Transm*, 124 (2017) 915-964.

693 [49] S. Lehericy, M.A. Sharman, C.L. Dos Santos, R. Paquin, C. Gallea, Magnetic resonance imaging
694 of the substantia nigra in Parkinson's disease, *Movement Disord*, 27 (2012) 822-830.

695 [50] M. Brammerloh, M. Morawski, I. Friedrich, T. Reinert, C. Lange, P. Pelicon, P. Vavpetic, S.
696 Jankuhn, C. Jager, A. Alkemade, R. Balesar, K. Pine, F. Gavriilidis, R. Trampel, E. Reimer, T. Arendt,
697 N. Weiskopf, E. Kirilina, Measuring the iron content of dopaminergic neurons in substantia nigra with
698 MRI relaxometry, *Neuroimage*, 239 (2021).

699 [51] T. Brucke, S. Djamshidian, G. Bencsits, W. Pirker, S. Asenbaum, I. Podreka, SPECT and PET
700 imaging of the dopaminergic system in Parkinson's disease, *J Neurol*, 247 (2000) 2-7.

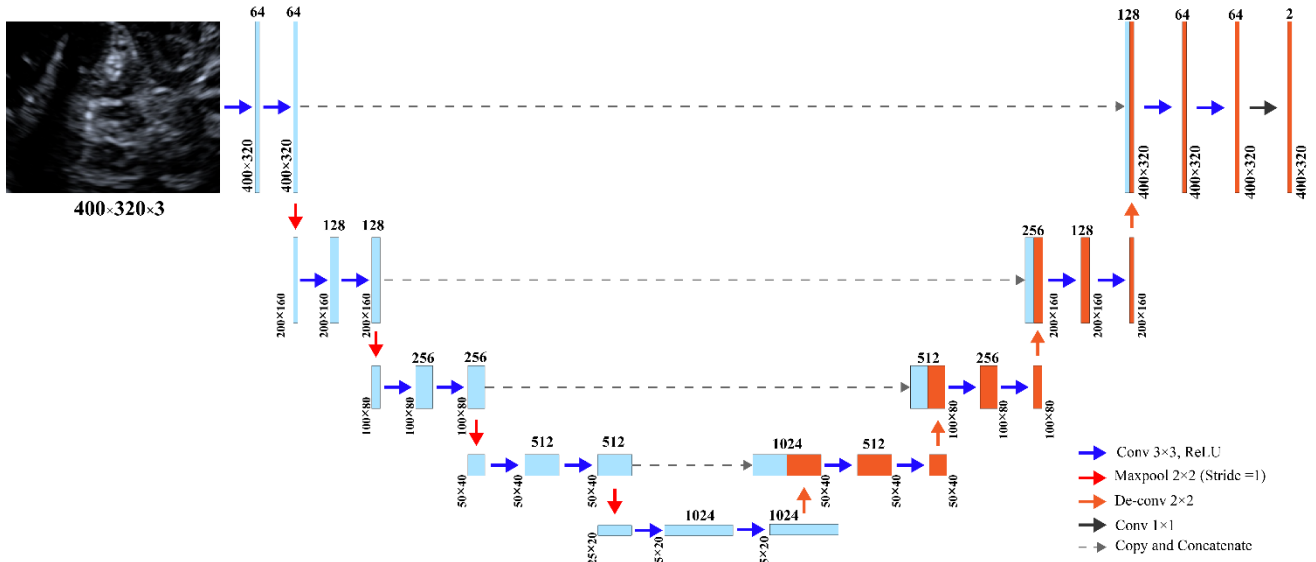
701 [52] G. Pagano, F. Niccolini, M. Politis, Imaging in Parkinson's disease, *Clin Med*, 16 (2016) 371-375.

702

703

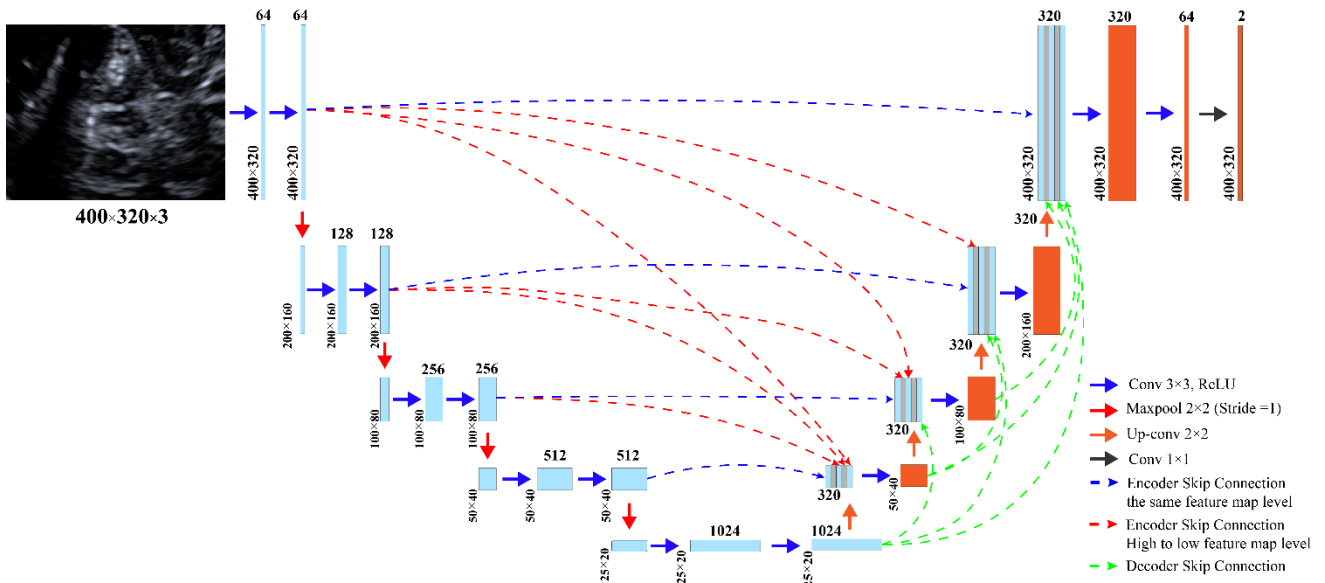
704
705
706

Supplementary Materials



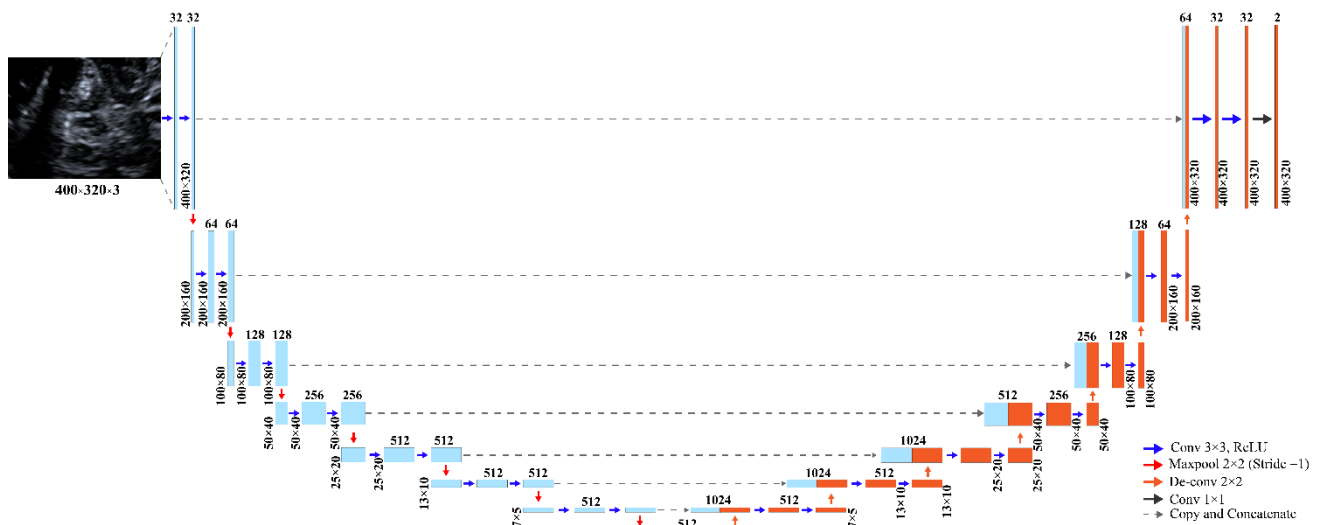
707
708
709
710

Fig. 1A. U-Net architecture



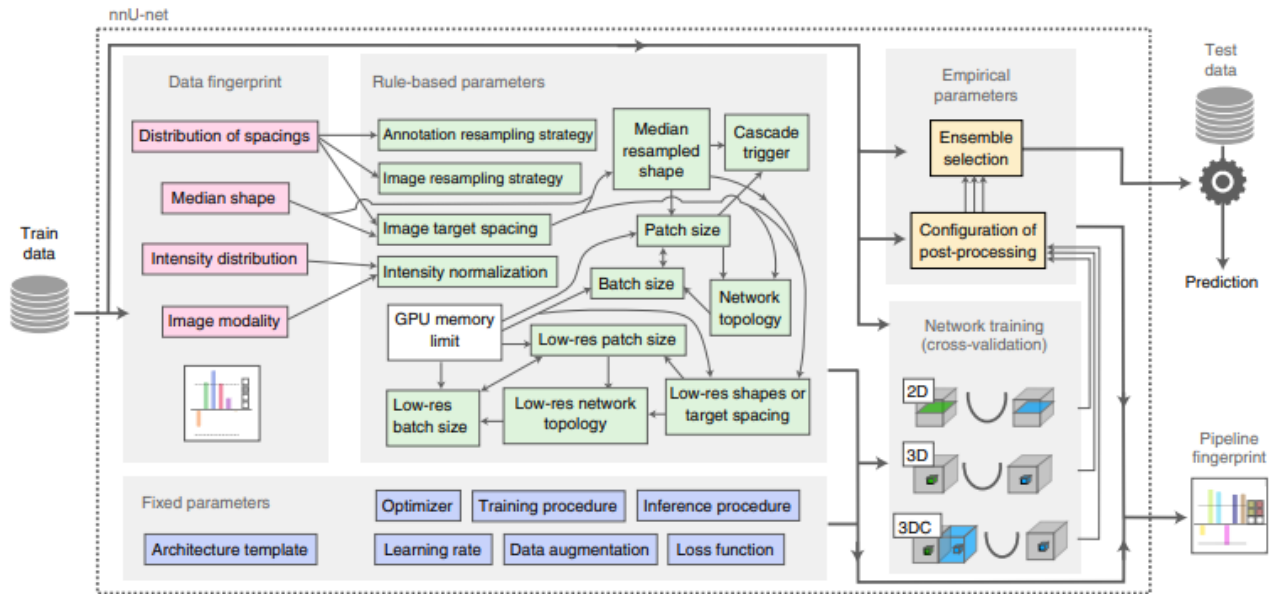
711
712
713

Fig. 1B. U-Net+++ architecture



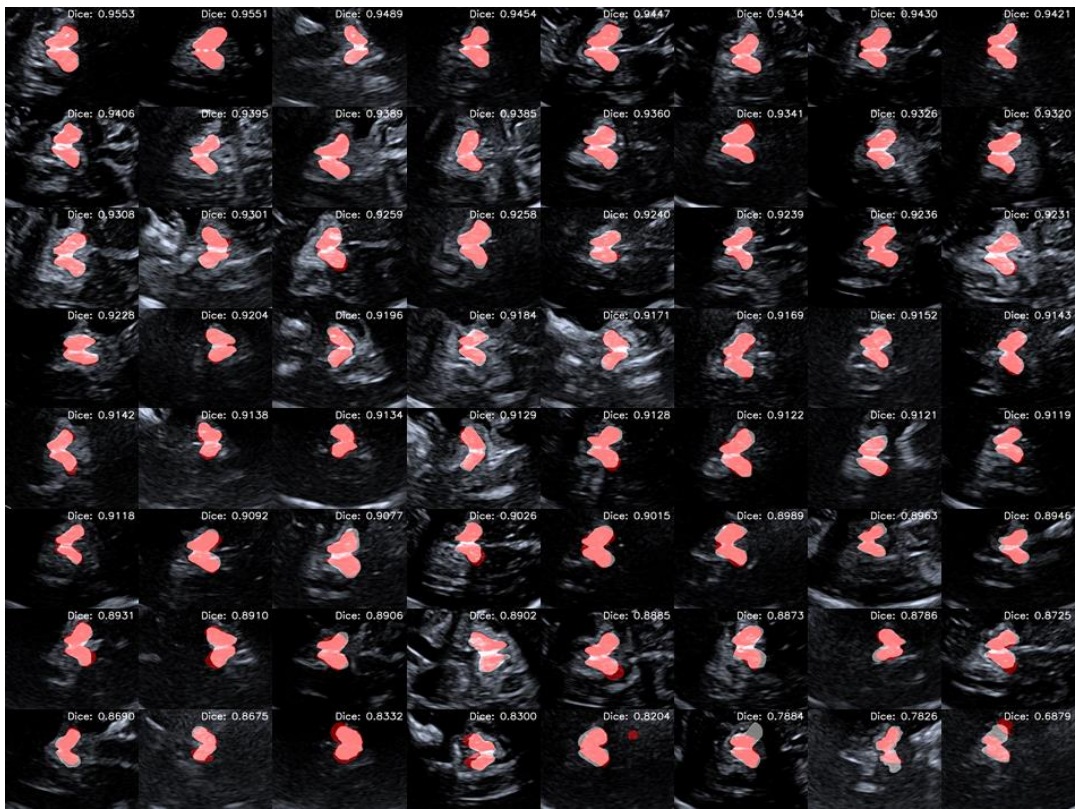
714
715
716

Fig. 1C. nnU-Net architecture



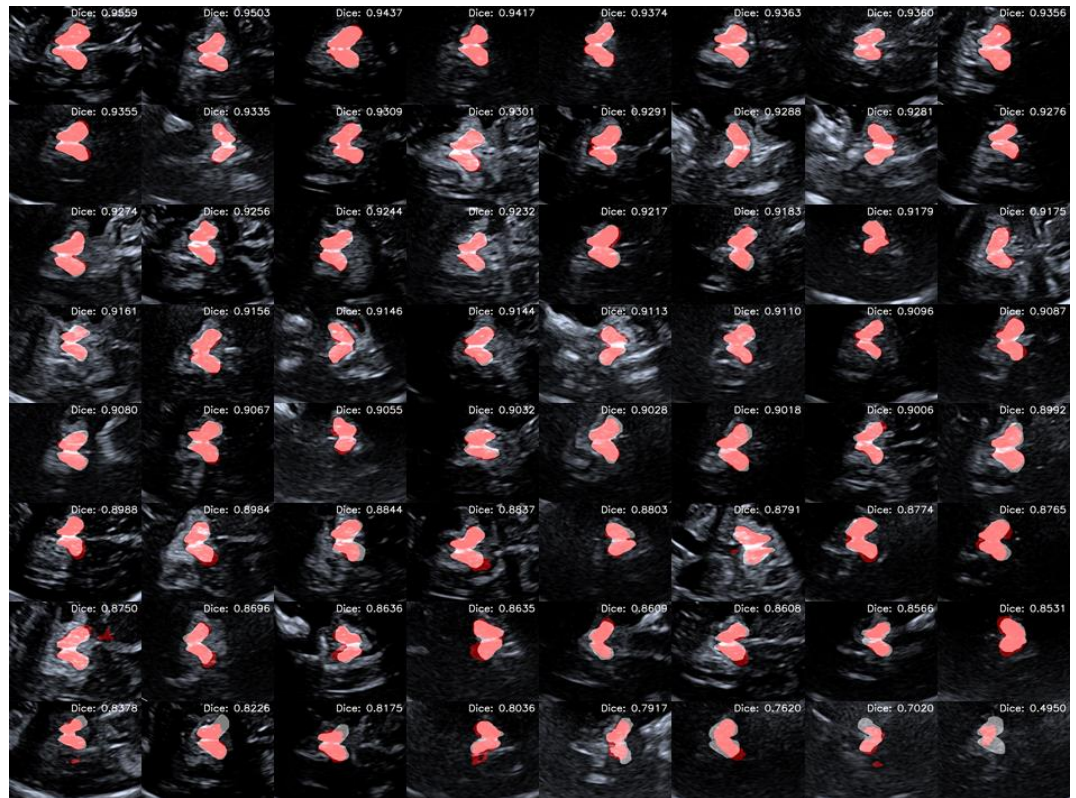
717
718

Fig. 1D. Flowchart of the nnU-Net architecture [20]



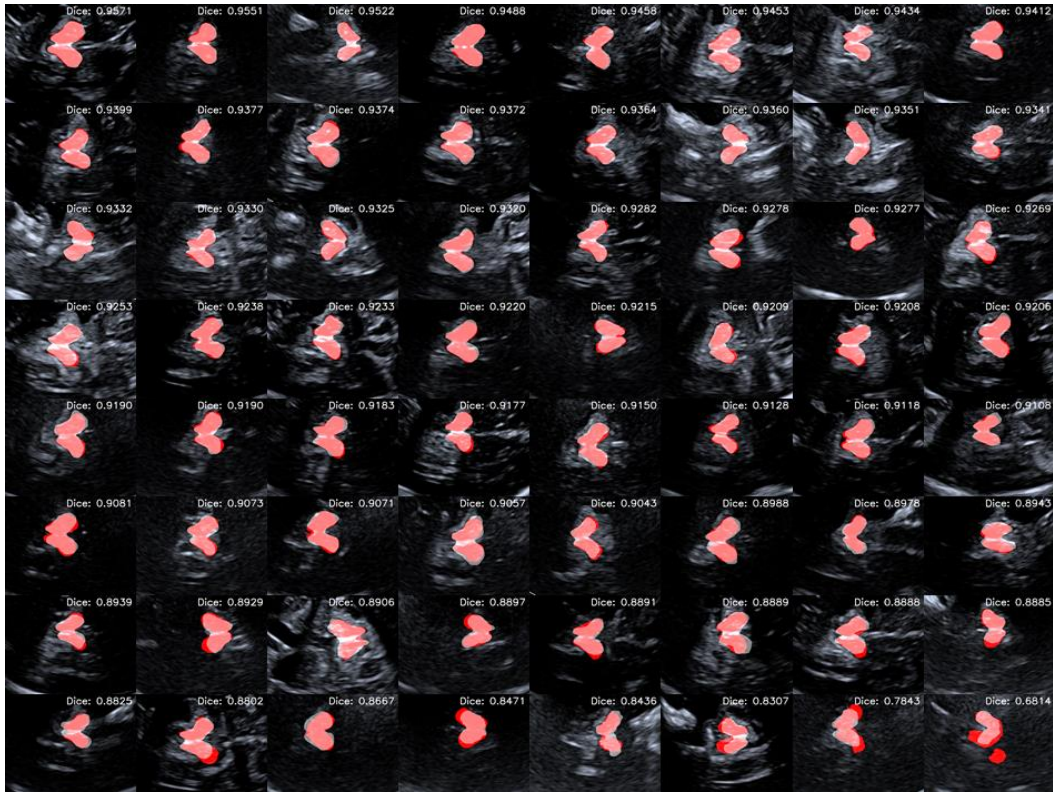
719
720
721
722

Fig. 2A. Qualitative visualization of the U-Net segmentation agreement with quantitative DICE score for all the 64 testing cases under the randomization seed 1, where the red regions represent the predicted segments and the white regions indicate the ground-truth annotations.



723
724
725
726

Fig. 2B. Qualitative visualization of the U-Net+++ segmentation agreement with quantitative DICE score for all the 64 testing cases under the randomization seed 1, where the red regions represent the predicted segments and the white regions indicate the ground-truth annotations.



727
728
729
730
731

Fig. 2C. Qualitative visualization of the nnU-Net segmentation agreement with quantitative DICE score for all the 64 testing cases under the randomization seed 1, where the red regions represent the predicted segments, and the white regions indicate the ground-truth annotations.

732

733
734
735
736

Table 1A. Results of segmentation agreement and stability of the three comparing neural networks in each of the three randomization seeds under 5 cross-validation on validation dataset. Averaged (AVG) scores of the segmentation agreement metrics in DICE, IoU, and HD are presented for each model in each randomization seed. Standard deviation (SD) was calculated to illuminate the underlying model stability. The top-ranked scores are bolded.

		Randomization Seed 1			Randomization Seed 2			Randomization Seed 3		
Network	Fold	DICE	IoU	HD (mm)	DICE	IoU	HD (mm)	DICE	IoU	HD (mm)
U-Net	Fold1	0.915	0.848	2.710	0.912	0.840	2.804	0.919	0.850	2.499
	Fold2	0.905	0.828	2.567	0.899	0.818	2.921	0.900	0.819	2.854
	Fold3	0.914	0.842	2.775	0.909	0.834	2.766	0.913	0.845	2.834
	Fold4	0.918	0.849	2.719	0.920	0.854	2.515	0.918	0.850	2.555
	Fold5	0.922	0.855	2.534	0.926	0.862	2.407	0.915	0.845	2.783
	AVG	0.915	0.845	2.661	0.913	0.842	2.683	0.913	0.842	2.705
	SD	0.006	0.011	0.105	0.010	0.017	0.214	0.008	0.013	0.166
U-Net+++	Fold1	0.909	0.835	3.126	0.908	0.832	3.198	0.909	0.836	2.898
	Fold2	0.892	0.807	3.069	0.893	0.810	3.097	0.893	0.810	3.189
	Fold3	0.910	0.836	2.906	0.901	0.824	3.419	0.906	0.831	3.016
	Fold4	0.914	0.843	3.084	0.913	0.840	2.961	0.914	0.842	2.958
	Fold5	0.915	0.845	2.890	0.918	0.848	2.859	0.918	0.849	2.881
	AVG	0.908	0.833	3.015	0.906	0.831	3.107	0.908	0.834	2.988
	SD	0.010	0.015	0.109	0.010	0.015	0.217	0.009	0.015	0.124
nnU-Net	Fold1	0.967	0.939	1.158	0.963	0.934	1.308	0.972	0.949	0.932
	Fold2	0.975	0.954	0.887	0.965	0.936	1.051	0.969	0.943	1.121
	Fold3	0.969	0.943	1.053	0.963	0.934	1.124	0.968	0.941	1.151
	Fold4	0.975	0.954	0.940	0.966	0.939	1.092	0.970	0.946	1.130
	Fold5	0.966	0.939	1.123	0.961	0.931	1.727	0.964	0.935	1.287
	AVG	0.970	0.945	1.091	0.964	0.935	1.261	0.971	0.944	1.124

AVG	0.970	0.946	1.032	0.964	0.935	1.261	0.969	0.943	1.124
SD	0.004	0.007	0.116	0.002	0.003	0.279	0.003	0.005	0.127

737

738

739

740

Table 1B. Results of segmentation agreement (in terms of DICE, IoU, and HD) and stability (in terms of SD) of the three deep-learning networks across all the three randomization seeds under 5-fold cross validation on validation dataset. The top-ranked scores are bolded.

Network	DICE (SD)	IoU (SD)	HD (mm) (SD)
U-Net	0.914 (0.008)	0.843 (0.013)	2.683 (0.156)
U-Net+++	0.908 (0.009)	0.833 (0.014)	3.037 (0.155)
nnU-Net	0.968 (0.004)	0.941 (0.007)	1.139 (0.200)

741