

1  
2 **Performance Benchmarking of Deep Learning Models**  
3 **for Real-Time Median Nerve Segmentation and**  
4 **Cross-Sectional Area Measurement in Ultrasound**  
5 **Imaging**  
6

7  
8 Vaddadi Venkatesh<sup>1</sup>, Lokesh Bathala<sup>2</sup>, Raji Susan Mathew<sup>3</sup>, and  
9 Phaneendra K. Yalavarthy, <sup>\*,1</sup>

10 <sup>1</sup> Department of Computational and Data Sciences, Indian Institute of Science, Bangalore-  
11 560012 India

12 <sup>2</sup> Aster CMI Hospital, Hebbal, Bangalore- 560092, India

13 <sup>3</sup> School of Data Science, Indian Institute of Science Education and Research,  
14 Thiruvananthapuram- 695551, India

15 \*e-mail: yalavarthy@iisc.ac.in

16 *RUNNING TITLE:* MNSeg-Net for Median Nerve Segmentation in Ultrasonography  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38

39 Medical Physics

40 Submitted on: March 18, 2026

## Abstract

**Background:** Median nerve, a major peripheral nerve, connects the hand to the central nervous system, facilitating upper limb motor function and sensation by transmitting sensory data from the palm and fingers. Damage to this nerve can result in motor and sensory deficits, with carpal tunnel syndrome (CTS) causing compression, leading to tingling and numbness in the thumb, index, middle, and lateral ring fingers.

**Purpose:** This study aimed to develop an accurate deep-learning-based segmentation method for measuring the cross-sectional area (CSA) of the median nerve to facilitate the diagnosis of nerve entrapment syndromes and aid in surgical planning, with a focus on CTS.

**Methods:** This study introduces MNSeg-Net, a novel lightweight multiscale feature fusion network with 2.46M parameters for median nerve segmentation in ultrasound (US) frames, specifically designed to enable a fully automated, end-to-end clinical setup supporting real-time segmentation and CSA computation. The dataset comprised 100 subjects and 30,000 ultrasound frames, which were split into training (80%), validation (10%), and testing (10%) subsets with subject-wise separation to avoid data leakage. MNSeg-Net was benchmarked against state-of-the-art segmentation models, including UNet and its variants (UNet++ and U2Net). The performance was assessed using metrics such as the Dice Similarity Coefficient (DSC) and CSA difference. The statistical significance of performance differences was evaluated using paired t-tests, effect size (Cohen’s  $d$ ), and one-way ANOVA with Tukey’s HSD correction for multiple comparisons at a  $p$ -value threshold of 0.05, while statistical equivalence between models within predefined margins was formally assessed using the two one-sided test (TOST) procedure. Following quantitative validation, the model was deployed in a real-time clinical setup utilizing an Av.io HD Epiphan frame grabber to stream ultrasound images from the ultrasound machine to a GPU-equipped system. A secondary display running parallel to the original ultrasound screen visualized the segmented median nerve and computed the CSA values in real time.

**Results:** MNSeg-Net achieved high segmentation performance, with average DSC scores of 94.7% at the wrist and 83.4% from the wrist to the elbow, and the lowest Hausdorff distance, matching the performance of the best-performing 44-million-parameter heavy U2Net model. Compared to U2Net, MNSeg-Net showed no statistically significant difference in DSC performance ( $p = 9.11 \times 10^{-1}$ ; Cohen’s  $d = -0.003$ ; mean difference =  $-0.001$ ), with formal equivalence testing confirming equivalence across all tested margins ( $\pm 0.01, \pm 0.03, \pm 0.05$ ). For CSA estimation, MNSeg-Net also showed no statistically significant difference from clinician-annotated values ( $p = 1.14 \times 10^{-1}$ ; Cohen’s  $d = -0.041$ ; mean difference =  $-0.081$ ), and equivalence was established at the  $\pm 0.5$  margin, confirming a strong alignment with expert clinical assessments. MNSeg-Net demonstrated real-time performance by processing up to 43 frames per second on a single GPU, successfully segmenting the median nerve and computing CSA from ultrasound frames.

**Conclusion:** The developed MNSeg-Net-based clinical system represents an important step toward real-time median nerve assessment, enabling a fully automated solution for CTS diagnosis. By combining a lightweight architecture, real-time processing capability, and successful clinical deployment, it represents a substantial advancement in the CTS detection and management.

87 **Keywords:** Median Nerve Segmentation, Ultrasound Video, Lightweight Network, Real-  
88 time Clinical Setup.

## 89 1. Introduction

90 The median nerve is a major peripheral nerve that serves as a critical communication pathway  
91 between the hand and central nervous system and is crucial for upper limb motor function  
92 and sensation<sup>1,2,3</sup>. It transmits sensory data from the palms and fingers to the central  
93 nervous system, which is essential for touch, temperature, and pain perception. Damage  
94 to the median nerve can result in notable impairment of movement and sensation. Carpal  
95 tunnel syndrome (CTS)<sup>1,4</sup>, the most common peripheral neuropathy, affects the thumb,  
96 index finger, middle finger, and the lateral side of the ring finger. It typically arises because  
97 of increased pressure within the enclosed carpal tunnel, which compresses the median nerve  
98 and is characterized by tingling and numbness in the affected hand.

99 High-frequency alternating currents<sup>5,6</sup> are used in peripheral nerve block therapy, re-  
100 quiring precise determination of the nerve diameter or cross-sectional area (CSA) to establish  
101 the minimum frequency for effective nerve block induction. Segmentation is crucial for di-  
102 agnosing conditions such as CTS, guiding ultrasound-guided regional anesthesia (UGRA)<sup>7</sup>,  
103 identifying nerve entrapment syndromes, and understanding nerve anatomy. Accurate seg-  
104 mentation of the median nerve in medical imaging enhances patient care by enabling precise  
105 surgical interventions. Regional anesthesia serves as an effective alternative to general anes-  
106 thesia in many surgeries. Traditionally, needle guidance to target nerves has been conducted  
107 blindly, risking nerve injury and local anesthetic toxicity. The UGRA allows real-time visu-  
108 alization of nerves and needle placement, thereby reducing these risks. Precise segmentation  
109 of the median nerve improves patient outcomes, and the safety and effectiveness of surgical  
110 procedures. Existing clinical methods often lack metrics for identifying CTS, and procedures  
111 such as nerve conduction studies may cause patient discomfort.

112 Numerous tools have been developed to aid radiologists in segmenting the median nerve  
113 in ultrasound (US) images. Hadjerci et al.<sup>8</sup> proposed an automated method for median nerve  
114 localization for UGRA, employing a machine learning framework that includes despeckling  
115 filtering, feature extraction, and selection, followed by pixel-wise classification using a sup-  
116 port vector machine with a Gaussian kernel. Similarly,<sup>9</sup> introduced a computer-aided ma-  
117 chine learning algorithm for median nerve localization. However, these techniques rely on  
118 manual feature extraction and selection, and the accuracy of the segmented images depends  
119 on the selected features.

120 In recent years, deep learning techniques have emerged as promising tools for medical ul-  
121 trasound analysis, particularly for nerve segmentation. One study introduced a convolutional  
122 neural network (CNN) with spatiotemporal consistency for nerve segmentation<sup>10</sup>, whereas  
123 another study used similarity measures to track the median nerve during wrist motion<sup>11</sup>.  
124 Huang et. al.<sup>12</sup> utilized the Attention-VGG16-UNet for median nerve segmentation at the  
125 wrist and trained the model using data from the wrist inlet. Automated frameworks, such as

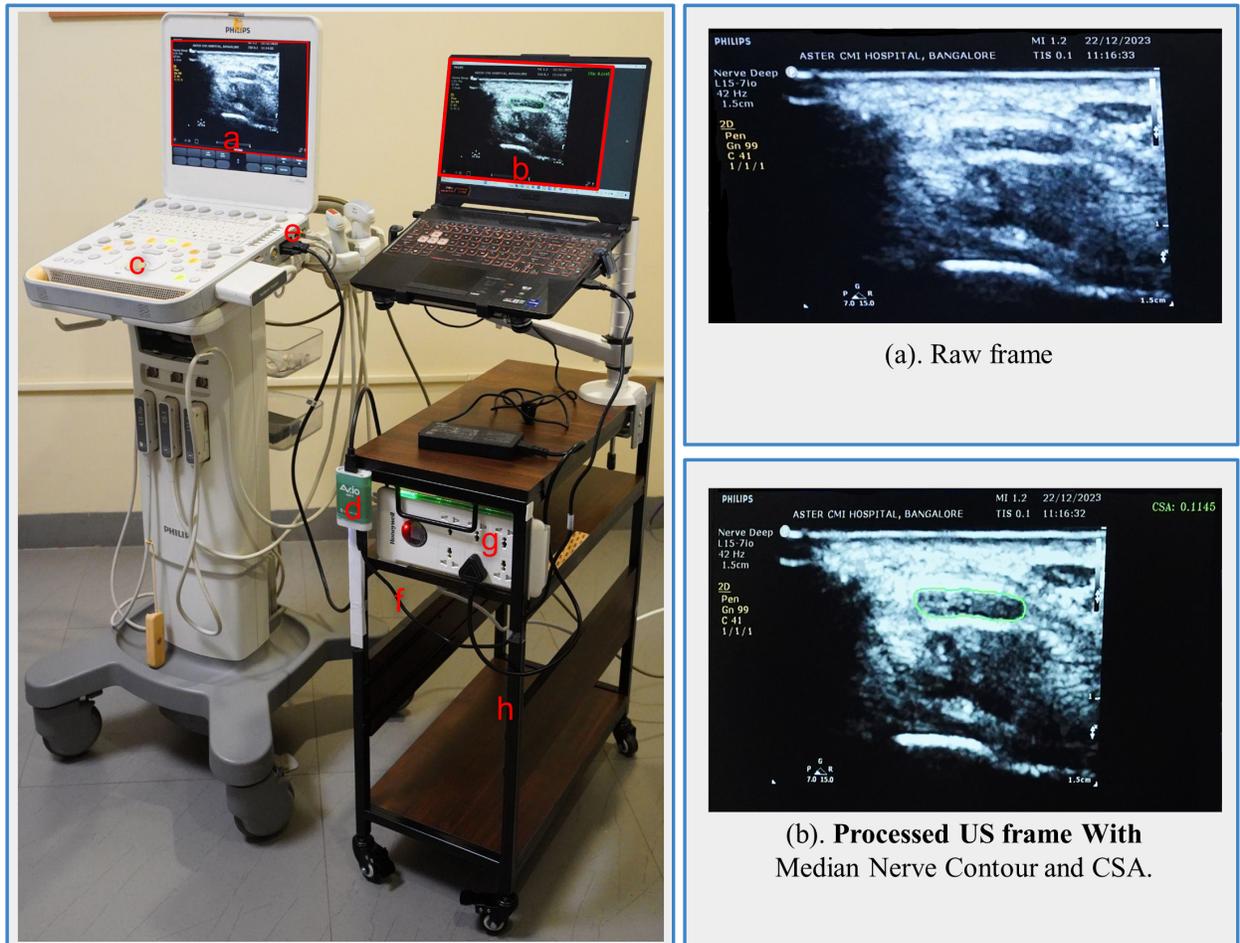


Figure 1: The physical system setup for Real-Time Median Nerve Segmentation includes the following main parts: (c). Philips cx50 ultrasound (US) machine(Input source) (d). Av.io Frame Grabber (e). DVI-to-HDMI interface (f).USB-B-to-USB-A interface (g). Switch box (power supply) (h). Portable wheeled table for the movable setup. The frame grabber was connected to the DVI output of the US machine. It is then connected via USB-B to a laptop, which serves as the computational device for executing the deep learning model.

126 a U-Net-based model<sup>13</sup> and DeepNerve, which integrates MaskTrack with convolutional long  
127 short-term memory (LSTM)<sup>14</sup>, have also been proposed. Comparative analyses evaluated  
128 the performance of pre-trained CNN models for median nerve localization/segmentation<sup>15</sup>,  
129 and a Mask R-CNN approach with additional layers to predict CSA<sup>16</sup>. A robot-assisted  
130 ultrasound imaging system was developed to enhance the monitoring of percutaneous needle  
131 insertion<sup>17</sup>, autonomously restoring visibility when misalignment occurs, and achieving  
132 precise segmentation of the needle with minimal positional and angular errors in ex vivo  
133 porcine samples. Yeh et al.<sup>18</sup> introduced a simple, real-time instance segmentation model  
134 for the median nerve at the wrist, achieving state-of-the-art speed and accuracy, making it  
135 highly suitable for diagnosing CTS with dynamic ultrasonography. Wang et al.<sup>19</sup> recently  
136 proposed a lightweight encoder-decoder network called GREnet for generalized 2D medi-  
137 cal image segmentation, including ultrasound images, that was prohibitively expensive for  
138 deployment in real-time. EU2-Net<sup>20</sup> is a novel ensemble model that enhances the U2-Net ar-  
139 chitecture for breast ultrasound image segmentation by replacing conventional convolutional  
140 layers with separable convolutions, incorporating a weighted averaging ensemble mecha-  
141 nism with learnable weights, and utilizing an attention-aided triple feature fusion technique  
142 to improve segmentation accuracy, achieving state-of-the-art results on publicly accessible  
143 datasets. Recent work by Pan et. al<sup>21</sup> a novel two-stage method combining a cascaded  
144 segmentation network and a knowledge-based classification network to accurately recognize  
145 medullary thyroid carcinoma (MTC) from ultrasound images, achieving higher accuracy  
146 than experienced clinicians. Another recent study<sup>22</sup> introduced VLSC-Net, a deep neural  
147 network-based approach for the precise segmentation and classification of dialysis access  
148 (DA) stenosis in ultrasound images, notably enhancing diagnostic accuracy and efficiency.  
149 Gujarati et al.<sup>23</sup> used a modified vision transformer (VisTR) model to segment the median  
150 nerve from wrist to elbow, achieving superior segmentation accuracy. However, the model’s  
151 high computational complexity and requirement for video input, which necessitates numer-  
152 ous frames and extensive spatial information, make it less practical for real-time clinical  
153 use.

154 Most methods have focused on segmenting the median nerve at the wrist, where localiza-  
155 tion is easier, and simple UNet-based architectures perform well. However, their performance  
156 declines from the wrist to the elbow because of the increased depth and rapid shape changes  
157 of the median nerve. UNet models struggle to capture long-range dependencies and complex  
158 spatial relationships, which are problematic in ultrasound image segmentation with rapidly  
159 changing foreground and background. In addition, many UNet-based methods lack effective  
160 metrics for identifying CTS, thereby limiting their clinical utility. Recent methods by Yeh  
161 et al.<sup>18</sup> focused on real-time instance segmentation using ultrasound for clinical use but were  
162 limited to the wrist region (not tested as a practical clinical application). Conversely, Gu-  
163 jarati et al.<sup>23</sup> aimed to enhance the segmentation accuracy; however, this was less practical  
164 for real-time clinical use with ultrasound.

165 This study introduces a novel encoder-decoder-based architecture with a unique design  
166 that captures richer contextual information through a larger receptive field. It also incorpo-  
167 rates a specialized network block to leverage the multiscale features and enhance the feature  
168 representation. These architectural changes resulted in a simple yet efficient deep learning

169 model that addressed the current limitations. The goal was to develop an end-to-end, deep  
170 learning-based automated tool to aid sonographers in segmenting the median nerve at the  
171 wrist, mid-forearm, and elbow using ultrasound images. This study proposes MNSeg-Net,  
172 a lightweight, real-time, inference-supported, and efficient CNN model for median nerve  
173 segmentation from the wrist to the elbow using a single ultrasound video frame.

174 To evaluate the effectiveness of the proposed MNSeg-Net, several well-known deep learn-  
175 ing models for medical image segmentation were selected as benchmarks. UNet<sup>24</sup> and its  
176 variants (SegNet<sup>25</sup>, ResUNet<sup>26</sup>, Attention-UNet<sup>27</sup>, UNet++<sup>28</sup>, BASNet<sup>29</sup>, and U2Net<sup>30</sup>)  
177 were included because they represent baseline encoder-decoder architectures with different  
178 enhancements, such as skip connections, attention mechanisms, and residual learning. BAS-  
179 Net and U2Net were selected because of their success in fine-grained object segmentation  
180 and hierarchical design, which aligns well with the requirements for median nerve localiza-  
181 tion in ultrasound. These benchmark models provide a comprehensive basis for evaluating  
182 the performance, computational efficiency, and clinical applicability of MNSeg-Net.

183 The specific contributions of this study, encompassing architectural innovation, statis-  
184 tical validation, and clinical integration, are as follows:

- 185 • **Proposed MNSeg-Net Architecture:** Development of MNSeg-Net, a novel  
186 lightweight deep learning architecture designed for fully automated, real-time me-  
187 dian nerve segmentation. The model integrates a computationally efficient redesigned  
188 U2Net backbone with a multi-scale feature fusion sub-network to effectively capture  
189 contextual information and enhance feature representation without adding the high  
190 computational overhead.
- 191 • **Comprehensive and Statistically Rigorous Evaluation:** A rigorous evaluation  
192 of the proposed model against state-of-the-art segmentation networks. This study  
193 establishes the effectiveness of lightweight architectures through extensive statistical  
194 analysis and demonstrated that high segmentation accuracy and cross-sectional area  
195 (CSA) estimation can be achieved with significantly reduced parameter counts.
- 196 • **Clinical Translation through Real-Time Integration:** Development and deploy-  
197 ment of an end-to-end clinical workflow that integrates the deep learning model with  
198 real-time ultrasound data. This contribution includes a custom graphical interface for  
199 real-time visualization and automatic CSA computation, validated on a diverse dataset  
200 of healthy subjects and patients with carpal tunnel syndrome (CTS) to assess clinical  
201 feasibility.

## 202 II. Methods

### 203 II.A. Data Acquisition and Preparation

204 This study used two distinct datasets: Dataset-1 for experimental evaluation and Dataset-2  
205 for clinical validation.

### 206 II.A.1. Dataset-1: Experimental Dataset

207 The primary dataset (Dataset-1) was sourced from the Aster-CMI Hospital, Bangalore,  
208 India, and acquired using a Philips CX50 ultrasound machine equipped with a Philips L15-  
209 7io hockey-stick transducer probe (frequency range: 7–15 MHz). The participants were  
210 positioned facing the examiner with their forearms and fingers extended, as shown in Fig.  
211 5. Imaging was performed at a depth of 3 cm, and each subject’s scan was recorded as an  
212 8-second video (approximately 312 frames at 39 fps), with each frame measuring  $800 \times 600$   
213 pixels. The dataset included 100 subjects, annotated by an expert sonographer. For each  
214 video, the first 300 frames were used for labeling, resulting in 30,000 annotated frames.  
215 Annotations were generated using ImageJ<sup>31</sup>, a free NIH image analysis tool, and stored  
216 as binary segmentation masks for the median nerve. The 100-subject dataset was split  
217 into training, validation, and test sets with 80, 10, and 10 subjects, respectively. Data  
218 augmentation via horizontal flipping was applied to the training set, increasing the total to  
219 48,000 images. The validation and test sets each comprised 3,000 images (300 frames per  
220 subject  $\times$  10 subjects).

### 221 II.A.2. Dataset-2: Clinical Evaluation Dataset

222 To assess clinical applicability, a separate dataset (Dataset-2) comprising 30 subjects was  
223 collected for real-time evaluation. This dataset included both NORMAL and CTS-positive  
224 cases, with 18 subjects categorized as NORMAL and 12 as CTS-positive based on clinical  
225 criteria, including a CSA threshold of 12 mm<sup>2</sup><sup>32,33,34</sup>. Each subject contributed a video  
226 recorded under the same imaging setup as Dataset-1, with six clinician annotations per  
227 subject, focused at the wrist level. These annotations were used to compare the MNSeg-  
228 Net-predicted CSA ( $CSA_{Cal}$ ) with the clinician-annotated CSA ( $CSA_{Act}$ ) for statistical and  
229 clinical reliability analyses.

230 All Ethical and experimental procedures were approved by the Internal Ethical Com-  
231 mittee of the Aster-CMI Hospital, Bangalore, India (Approval No. Aster/IEC/049/2020-21,  
232 dated June 27, 2020). Written informed consent was obtained from all participants. The  
233 study population had a male-to-female ratio of 1:3, with participants aged between 35 and  
234 65 years. Individuals with prior major nerve surgery were excluded because of the possibility  
235 of postoperative anatomical alterations.

## 236 II.B. Proposed Median Nerve Segmentation Network

237 The proposed Median Nerve Segmentation Network, referred to as MNSeg-Net, is a  
238 lightweight architecture designed for accurate, real-time segmentation of the median nerve  
239 in ultrasound images. While inspired by the U2Net framework<sup>30</sup>, MNSeg-Net introduces  
240 several important modifications to reduce computational cost and enable deployment on  
241 resource-constrained clinical hardware, without sacrificing segmentation performance.

242 The overall architecture of MNSeg-Net is organized into two components: (i) a Main-

243 Network, which is a lightweight redesign of U2Net optimized for real-time inference, and  
244 (ii) a Sub-Network, a multi-scale feature fusion (MSFF) module that efficiently aggregates  
245 encoder features to enhance contextual representation. An overview of the complete ar-  
246 chitecture, including the main encoder–decoder backbone and the MSFF sub-network, is  
247 shown in Fig. II.B.1., while the detailed UNet block configurations are provided in Fig. 1 of  
248 the supplementary materials.

### 249 II.B.1. Main-Network: Lightweight Redesign of the U2Net

250 The original U2Net<sup>30</sup> employs a fully nested U-structure with exponentially varying filter  
251 counts, resulting in a model of approximately 44M parameters that is computationally ex-  
252 pensive for real-time or edge deployment. In contrast, the proposed main network adopts a  
253 lightweight redesign in which the encoder (E1–E5) and bottleneck (E6) are mirrored by the  
254 decoder (D1–D5), with each stage implemented as a compact residual UNet block using uni-  
255 form filter allocation (64 filters in outer blocks and 32 in inner blocks). Each block includes  
256 a single residual connection to enhance the gradient flow while reducing redundancy. This  
257 modification reduces the parameter count to 2.46M and lowers the FLOPs by more than  
258 an order of magnitude, enabling real-time inference without sacrificing the representational  
259 capability. To further balance accuracy and efficiency, deeper residual U-blocks are used  
260 in the early encoder–decoder stages (E1–E3 and D1–D3), whereas the later stages (E4–E5  
261 and D4–D5) adopt shallower blocks to reduce computation. Additionally, dilated convolu-  
262 tions replace pooling and upsampling in the deeper encoder (E5), bottleneck (E6), and final  
263 decoder (D5), thereby enlarging the receptive field while preserving the fine spatial details  
264 crucial for nerve boundary localization. Together, these modifications provide a compact yet  
265 powerful backbone that retains the contextual strengths of U2Net while being optimized for  
266 clinical applications requiring efficient and accurate real-time segmentation.

### 267 II.B.2. Sub-network: Multi-Scale Feature Fusion

268 In addition to the main encoder–decoder backbone, MNSeg-Net incorporates a dedicated  
269 sub-network module for MSFF, designed to enhance segmentation performance by leveraging  
270 complementary features across encoder stages. Specifically, feature maps from all encoder  
271 blocks (E1–E5) and the bottleneck (E6) are extracted, resized to the input image resolution,  
272 and concatenated to form a multiscale representation that captures both local fine-grained  
273 details and global contextual information. This fused representation is then processed by a  
274 lightweight UNet, which transforms the multichannel input into a compact single-channel  
275 feature map at full resolution. By reusing encoder features rather than introducing new  
276 feature extraction layers, the MSFF module introduces a minimal additional computational  
277 cost while significantly improving the ability to localize nerve boundaries. Finally, the output  
278 of this sub-network is integrated with the decoder outputs (D1–D5) and bottleneck (E6)  
279 during the final prediction stage, ensuring that both multi-scale context and deep supervised  
280 features contribute to the final segmentation mask.

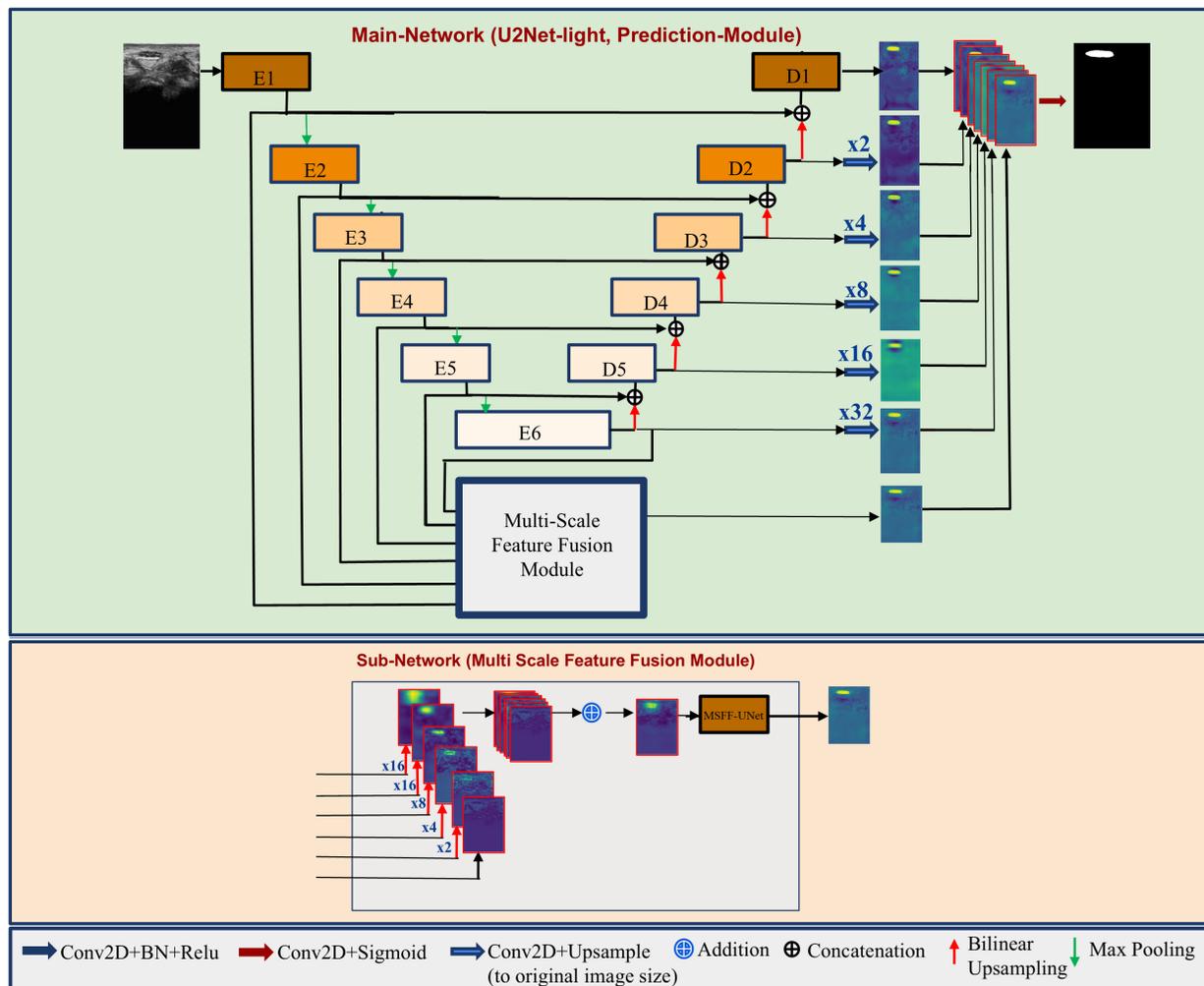


Figure 2: The proposed MNSeg-Net architecture for Median Nerve Segmentation with deep supervision and sub-network module for multi-scale feature fusion. (a) Main-network as prediction module, (b) Sub-network for multi-scale feature fusion. Detailed UNet block configurations used in MNSeg-Net were provided in Fig. 1 of supplementary information.

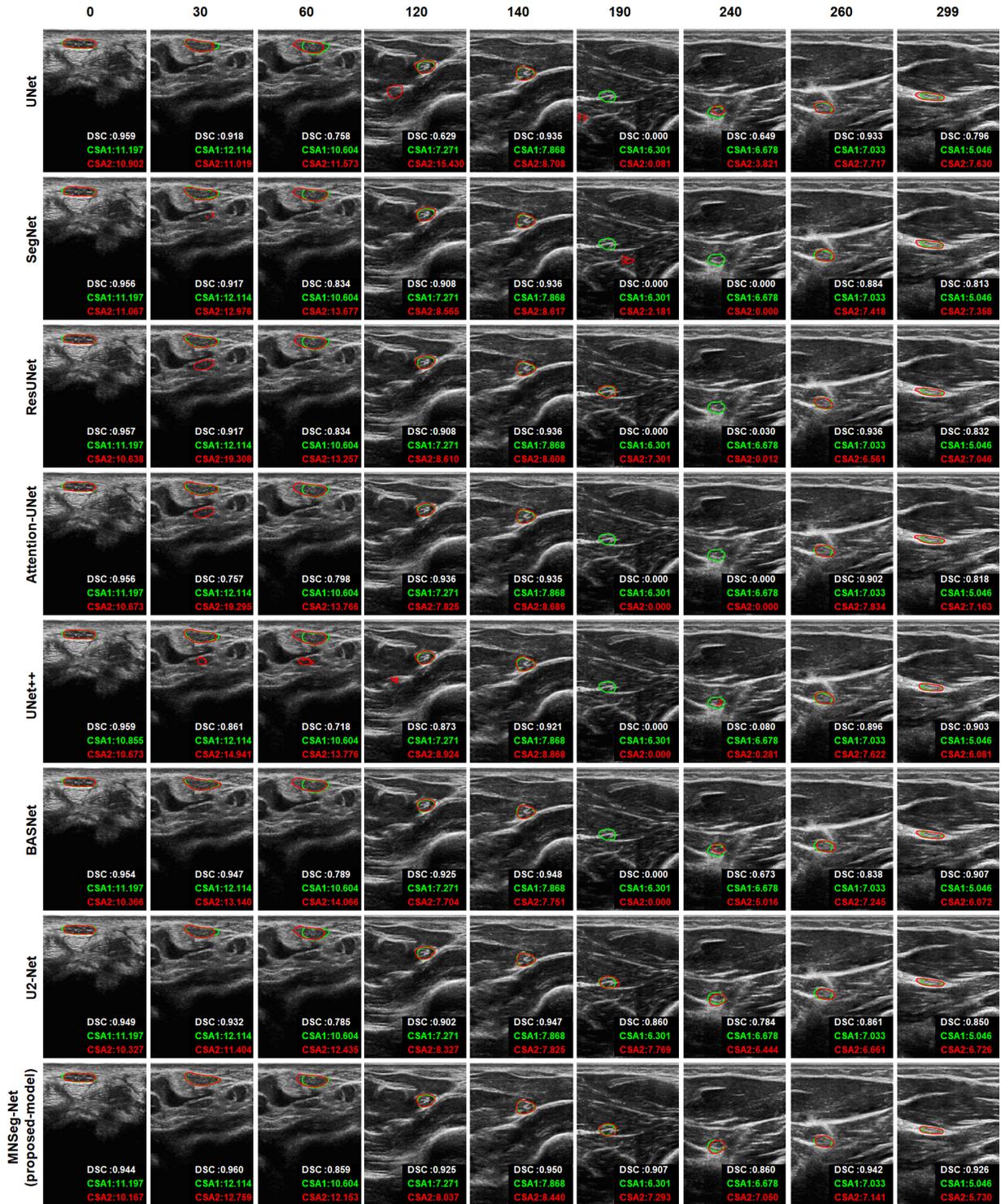


Figure 3: Example segmentation of the median nerve using the methods discussed in this work for subject-1. The green contour indicates expert annotation, and the red contour indicates the result obtained for the corresponding method, as indicated in each row. The associated frame number is given at the top of every image (0 corresponds to the start of the wrist region, and 299 corresponds to the elbow region), and the bottom of each frame has the corresponding computed cross-sectional area (CSA) from the expert and model in green and red, respectively.

281 Overall, the MSFF sub-network serves as a compact yet powerful enhancement to the  
 282 main backbone. By reusing the multi-scale features from E1–E6 and fusing them through a  
 283 lightweight UNet, the contextual representation is enriched while keeping the computation  
 284 minimal. This design ensures that MNSeg-Net effectively captures both fine local details  
 285 and global structures, making the segmentation output more accurate and reliable without  
 286 sacrificing efficiency.

Table 1: Comparison between the proposed MNSeg-Net and the full-size U2Net model.

		Proposed MNSeg-Net					U2Net (Full size)				
Type	Block	RAM Occupied (MB)	Module Size (MB)	Comput. (G. FLOPS)	Params (M)	Output (SHAPE)	RAM Occupied (MB)	Module Size (MB)	Comput. (G. FLOPS)	Params (M)	Output (SHAPE)
Encoder Blocks	E1	807.46	0.79	9.6	0.206	64x448x320	807.46	0.79	9.600	0.206	64x448x320
	E2	210.85	0.81	3.65	0.213	64x224x160	351.20	1.17	6.981	0.306	128x448x320
	E3	53.24	0.71	0.914	0.185	64x112x80	179.30	4.23	6.958	1.109	256x112x80
	E4	13.76	0.60	0.228	0.158	64x56x40	102.84	15.21	6.947	3.987	512x56x40
	E5	6.21	0.60	0.089	0.158	64x28x20	38.28	38.28	5.620	10.035	512x28x20
Bottle Neck	E6	2.00	0.60	0.022	0.158	64x14x10	49.49	38.28	1.405	10.035	512x14x10
Decoder Blocks	D5	6.48	0.74	0.109	0.194	64x28x20	93.22	47.28	6.942	12.394	512x28x20
	D4	14.45	0.74	0.311	0.194	64x56x40	75.53	16.33	7.603	4.280	256x56x40
	D3	55.57	0.85	1.245	0.222	64x112x80	122.70	4.51	7.610	1.182	128x112x80
	D2	219.74	0.95	4.98	0.250	64x224x160	237.52	1.24	7.623	0.324	64x224x160
	D1	876.09	1.06	19.922	0.278	64x448x320	753.05	0.53	14.952	0.139	64x224x160
MSFF-module	MSFF-UNet	877.05	0.92	14.646	0.241	64x448x320	-	-	-	-	-
Conv (projection) + Concat		-	-	-	-	7x448x320	-	-	-	-	6x448x320
Out Conv	Out Conv	-	-	-	-	1x448x320	-	-	-	-	1x448x320
Overall	-	3142.8	9.37	55.716	2.458	1x448x320	2810.59	167.85	82.243	43.996	1x448x320

### 287 II.B.3. MNSeg-Net Architectural Contributions

288 The proposed MNSeg-Net integrates a lightweight adaptation of U2Net as the main network  
 289 together with a novel MSFF module as a sub-network. Several architectural enhancements  
 290 have been introduced to achieve this efficient design. Specifically, the encoder (E1–E5, E6)  
 291 and decoder (D1–D5) blocks of the original U2Net were replaced with compact residual  
 292 UNet (RSU) blocks, each featuring a single residual connection to improve the gradient  
 293 flow while reducing the redundancy. Unlike the original U2Net, which employs an exponen-  
 294 tial variation in filter counts, the redesigned architecture adopts a uniform filter allocation  
 295 (64 in the outer blocks and 32 in the inner blocks), thereby reducing the number of pa-  
 296 rameters and substantially lowering the computational complexity. Side-output layers were  
 297 retained to support deep supervision and enhance the training stability. Furthermore, a novel  
 298 MSFF module was introduced, employing a half-UNet-style decoder that directly upsam-  
 299 ples and fuses multi-scale features from the encoder outputs (E1–E6), producing a compact  
 300 and context-rich representation with minimal memory and computational overhead require-  
 301 ments. These combined architectural optimisations make MNSeg-Net considerably lighter  
 302 than the original U2Net while maintaining competitive segmentation accuracy.

303 A detailed analytical derivation of the relationship between model parameters, com-  
 304 putational complexity, memory consumption, and clinical deployability, including explicit  
 305 parameter-scaling formulations for RSU blocks and the overall MNSeg-Net architecture, is  
 306 provided in the section 1 of the Supplementary Material.

307 In addition to these technical contributions, the overall design of MNSeg-Net is pur-

308 posefully aligned with clinical translation. By combining low computational demand with  
 309 accurate segmentation, the model supports deployment in fully automated, real-time clinical  
 310 setups for median nerve segmentation and CSA estimation.

311 Consequently, MNSeg-Net achieves real-time inference with only 2.46 million parameters  
 312 and approximately 56 GFLOPs per image. A detailed comparison of the computational  
 313 requirements between MNSeg-Net and the original U2Net is provided in Table 1, based on  
 314 an input size of  $448 \times 320$ .

Table 2: Deep learning models Specifications utilised in this work and their corresponding inference time

Model	Model Size (MB)	No. of Parameters (M)	Comput. (GFLOPs)	Inference Memory (MB)	Training Time (Min/Epoch)	Inference Speed (FPS)	Inference Time (ms/frame)
UNet <sup>24</sup>	96.12	25.2	72.60	1854	20	71	14.08
SegNet <sup>25</sup>	112.32	29.44	87.74	2781	30	66	15.15
ResUNet <sup>26</sup>	157.60	39.37	111.63	2472	40	38	26.32
Attention-UNet <sup>27</sup>	152.20	38.01	166.00	3702	40	29	34.48
UNet++ <sup>28</sup>	146.60	36.62	302.43	4984	40	19	52.63
BASNet <sup>29</sup>	332.11	87.06	278.91	6103	40	20	50.00
U2Net <sup>30</sup>	167.88	44.01	82.53	2978	30	40	25.00
<b>MNSeg-Net (Proposed)</b>	<b>9.40</b>	<b>2.46</b>	<b>56.00</b>	<b>3143</b>	<b>30</b>	<b>43</b>	<b>23.26</b>

## 315 Training

316 The proposed model was trained using a deep supervision strategy, in which losses were  
 317 applied not only to the final network output but also to the intermediate outputs to enhance  
 318 the gradient flow and stabilize the optimization. Specifically, eight outputs were considered:  
 319 the final output, five decoder blocks, a bottleneck, and a subnetwork module. Each output  
 320 was optimized using the same hybrid loss function, which combined three complementary  
 321 components: SSIM loss, log-cosh loss, and binary cross-entropy (BCE) loss. This  
 322 hybrid formulation leverages the perceptual sensitivity of SSIM, the robustness of log-cosh to  
 323 outliers, and the probabilistic classification strength of BCE, thereby ensuring both structural  
 324 fidelity and pixel-level accuracy. The overall training loss was obtained by summing the  
 325 hybrid loss values across all eight supervised outputs. Although this comprehensive deep  
 326 supervision improves gradient propagation and supports robust learning, it also requires a  
 327 relatively high GPU memory. Training converged within seven epochs, taking approximately  
 328 5.25 hours of computational time.

## 329 II.C. Testing and CSA Computation

330 After completing end-to-end training, the architecture was used to classify each pixel of the  
 331 test image into two categories: background and median nerve. Each pixel was assigned a  
 332 label based on the highest probability score among the two categories. In addition, the

333 proposed method automatically computes the CSA of the median nerve by calibrating the  
334 dimensions of individual pixels. Specifically, for the 3 cm depth settings, each pixel occupied  
335 an area of  $0.00432mm^2$ , and for the 1.5cm depth settings, it was  $0.00119mm^2$ . The CSA  
336 was computed by multiplying the single-pixel area by the number of median nerve pixels.

## 337 II.D. Quantitative Performance Evaluation

338 The effectiveness of the proposed deep learning model was assessed using metrics such as  
339 the DSC<sup>35</sup>, precision (Prec), recall(Rec), and the Hausdorff distance (HD)<sup>36</sup>. The DSC  
340 evaluates the overlap between expert annotations and predicted segmentation masks, with  
341 higher values indicating better alignment. Precision measures the proportion of true-positive  
342 predictions among all positive predictions, whereas recall assesses the proportion of true-  
343 positive predictions among all actual positives. HD measures the furthest distance from any  
344 point in one set to the nearest point in the other set, offering a numerical assessment of  
345 segmentation accuracy.

## 346 II.E. Implementation

347 All the deep learning models were trained using PyTorch and the Adam optimizer<sup>37</sup>. The  
348 models were trained using various loss functions, including the dice loss<sup>38</sup>, (BCE)<sup>39</sup>, logcosh  
349 loss<sup>40</sup>, and a custom hybrid loss function that combines the logcosh loss, BCE loss, and  
350 SSIM loss<sup>29,41</sup>. Computations, including model training, were performed on a Linux work-  
351 station with an Intel i9 9920X CPU, 128 GB RAM, and two NVIDIA Quadro RTX 8000  
352 GPUs with 48 GB of memory each. The hybrid loss ( $loss_{Hybrid}$ ) was selected after extensive  
353 experimentation for its superior performance, defined as

$$354 \quad loss_{Hybrid} = loss_{SSIM} + loss_{logcosh} + loss_{BCE}. \quad (1)$$

355 During the training process, the models were trained at various learning rates, including  
356  $1 \times 10^{-3}$ ,  $1 \times 10^{-4}$ ,  $5 \times 10^{-4}$ , and  $5 \times 10^{-5}$ . The UNet-based models exhibited stable training  
357 at  $1 \times 10^{-4}$ . However, the training time required was relatively high for learning rates of  
358  $1 \times 10^{-5}$  and  $5 \times 10^{-5}$ . Similarly, the proposed MNSeg-Net model exhibited stable training  
359 at  $1 \times 10^{-4}$  and comparable results at  $5 \times 10^{-4}$ , and the training time required was relatively  
360 high for a learning rate  $5 \times 10^{-5}$ . Because all the models were encoder-decoder-based, the  
361 training process was stable. All models, including the proposed model, were trained from  
362 scratch, and the best model was chosen based on the minimum validation loss. For a fair  
363 comparison, all the UNet-based models were considered to have the same depth (in terms of  
364 the encoder and decoder blocks) and the same number of filters in each stage. The Adam  
365 optimizer was used to train all models. In all experiments, training was continued for 30  
366 epochs to maintain consistency.

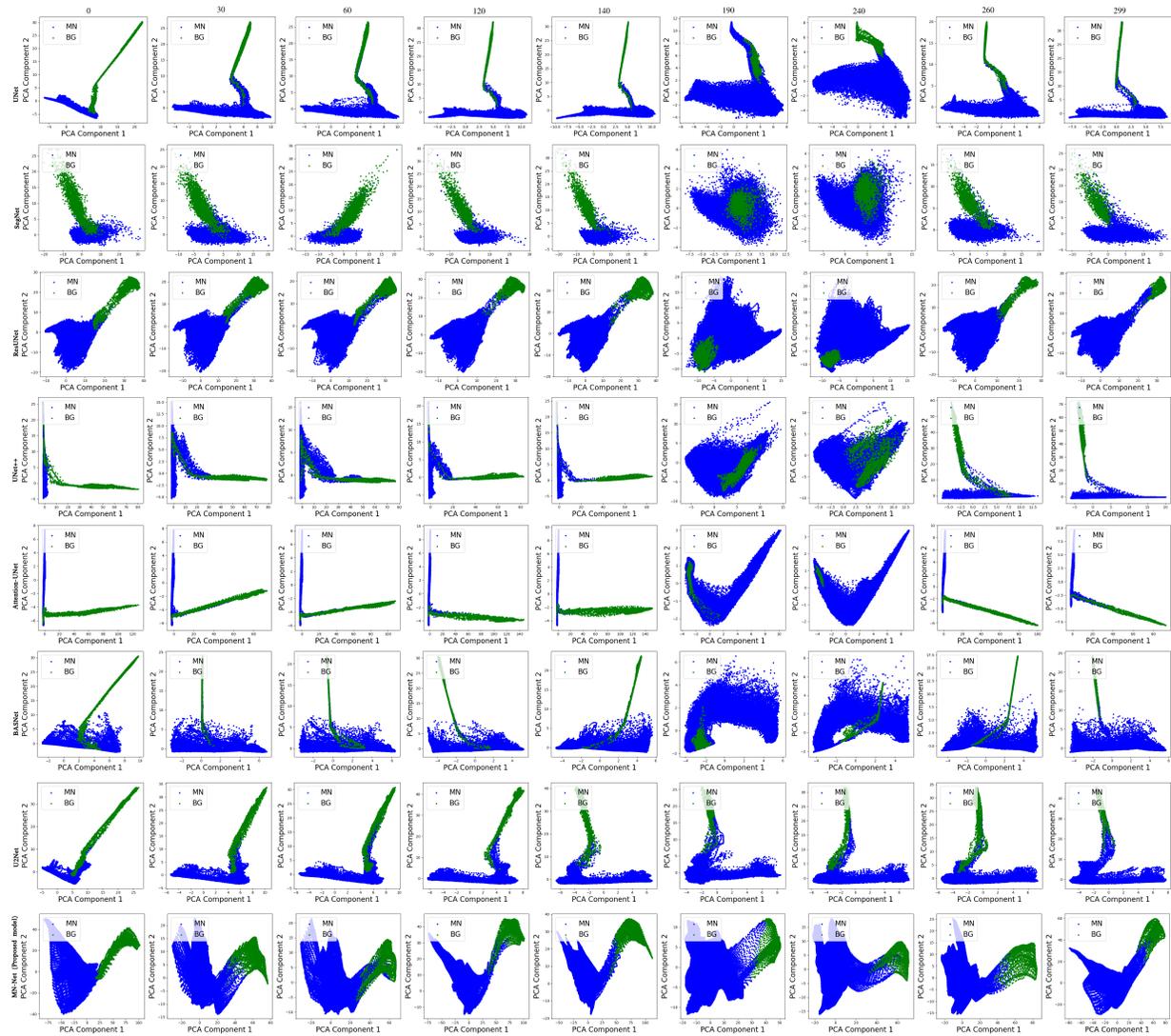


Figure 4: Qualitative visualization of the learned feature representations for patient 1. The features extracted from the final encoder layer of each model were projected into two dimensions using PCA. Each point corresponds to a pixel classified as either background (blue) or median nerve (green). Compared with baseline models, MN-Seg-Net exhibits clearer separation between the two classes, indicating more discriminative feature learning.

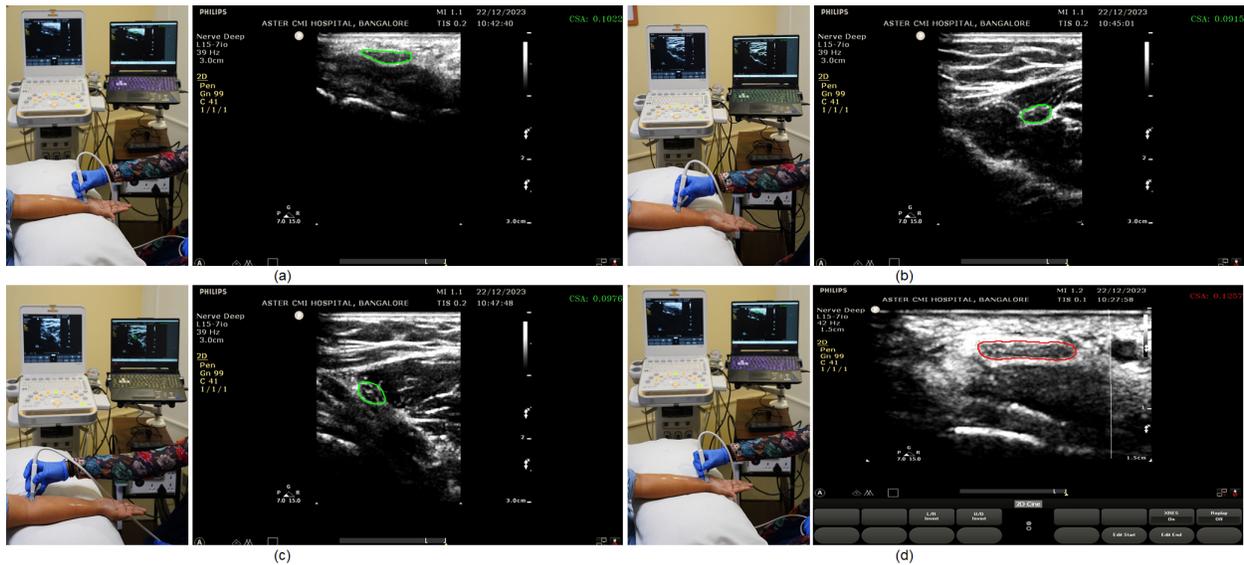


Figure 5: (a), (b), and (c) represent sample images captured during real-time testing of wrist, forearm, and elbow scanning, shown alongside their corresponding median nerve segmentations and CSA in US frames. The CSA values are presented as  $cm^2$ . In addition, (d) illustrates a sample image in which the CSA surpassed the threshold of  $0.12 cm^2$  during wrist scanning. The CSA values are displayed in the top right corner of the ultrasound frame.

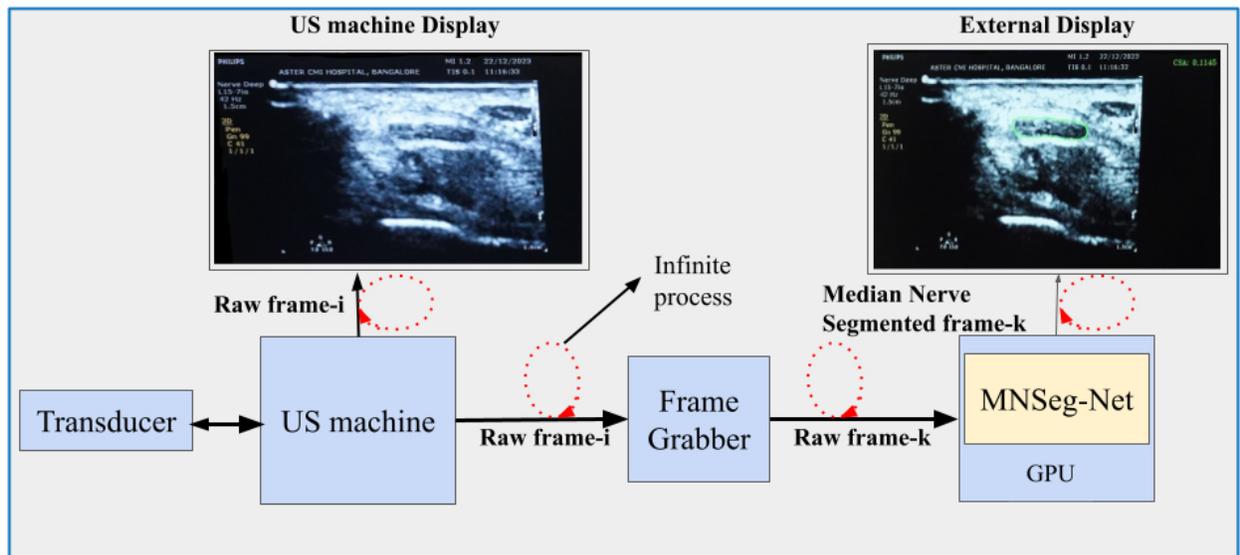


Figure 6: Timeline diagram of the clinical workflow. The US transducer generates raw frames that are displayed on the ultrasound machine while simultaneously being captured via a frame grabber. The raw frames are sent to the GPU, where the proposed MNSeg-Net model processes them to generate segmented frames with CSA computation. These are displayed on an external screen, providing real-time quantitative and visual feedback alongside the ultrasound machine's native display.

## 367 II.F. Real-Time Setup

368 The proposed median nerve segmentation model was deployed in a clinical setting using a  
369 customized desktop application. The real-time setup employed a frame grabber to capture  
370 the unprocessed video output of the ultrasound machine via a Digital Visual Interface (DVI).  
371 An AV.io frame grabber<sup>42</sup> captured the raw video output and converted it into serial frames,  
372 which were then transferred to a laptop (serving as an edge device) through a USB-B to  
373 USB-A interface. The laptop used these raw frames as input for the MNSeg-Net model,  
374 which was executed on its GPU using PyTorch inference. The model generated a segmenta-  
375 tion mask for the median nerve, which was displayed as an overlay on the input frame along  
376 with the computed CSA on the laptop screen. This setup provided feedback to the operator  
377 in addition to the main screen of the ultrasound machine. Figure 1 provides a comprehensive  
378 overview of the clinical setup, whereas Figure 6 illustrates the data flow/timeline diagram of  
379 the application. The development of this clinical setup was inspired by prior research<sup>43</sup>. To  
380 the best of our knowledge, this is the first system to provide dual-screen real-time segmenta-  
381 tion of the median nerve combined with quantitative CSA measurements during ultrasound  
382 scanning. To support this claim, a comprehensive literature review was performed. While  
383 some existing studies have addressed CSA measurement<sup>44</sup> and others have proposed real-  
384 time segmentation techniques<sup>18</sup>, these methods were typically evaluated offline and not in  
385 clinical settings. The review—performed across PubMed, IEEE Xplore, and Google Scholar  
386 on 17-04-2025 using combinations of the keywords “median nerve,” “ultrasound,” “real-time  
387 segmentation,” “dual-screen interface,” and “CSA measurement”—confirmed that no prior  
388 work has unified real-time segmentation with quantitative CSA estimation in a dual-screen  
389 clinical system

## 390 II.G. Feature Visualization

391 To qualitatively assess the discriminative capacity of the learned feature representations,  
392 feature embeddings were extracted from the final encoder layer of each model. Pixels cor-  
393 responding to the median nerve and background were projected into two dimensions using  
394 principal component analysis (PCA). PCA was selected over non-linear techniques such as  
395 t-SNE because the embedding set for a single image frame contained 1,43,360 ( $448 \times 320$ )  
396 points in 64 dimensions, making t-SNE computationally expensive, sensitive to hyper pa-  
397 rameters, and prone to overcrowding in large-scale settings.

398 Feature representations were extracted from the sample images (subject-1) of each  
399 trained model. Pixels corresponding to the median nerve and background were projected  
400 into two dimensions using principal component analysis (PCA) for data visualization. This  
401 allowed the higher-dimensional pixel embeddings in the segmentation map to be represented  
402 in a lower-dimensional space that is suitable for qualitative assessment.

## 403 II.H. Statistical Analysis

404 The performance of the proposed MNSeg-Net was evaluated using a comprehensive statisti-  
405 cal analysis based on the Dice Similarity Coefficient (DSC) and CSA. Since all segmentation  
406 models were evaluated on the same test images, paired statistical analyses were employed  
407 to account for within-sample dependence. Accordingly, paired  $t$ -tests were used to assess  
408 whether performance differences between MNSeg-Net and each benchmark model were sta-  
409 tistically significant, following the classical formulation introduced by Student<sup>45</sup>. A signifi-  
410 cance threshold of  $p < 0.05$  was adopted in accordance with widely accepted conventions in  
411 biomedical and clinical research. To quantify practical significance beyond statistical signifi-  
412 cance, effect sizes were computed using Cohen’s  $d$ . The magnitude of effects was interpreted  
413 using standardized thresholds for negligible ( $|d| < 0.2$ ), small ( $0.2 \leq |d| < 0.5$ ), moderate  
414 ( $0.5 \leq |d| < 0.8$ ), and large ( $|d| \geq 0.8$ ) effects, as proposed by Cohen<sup>46</sup>. These effect size  
415 categories provide an interpretable measure of the magnitude of observed differences. To  
416 control for the increased risk of type I error arising from multiple model comparisons, a  
417 one-way analysis of variance (ANOVA) was first conducted to assess whether overall differ-  
418 ences existed among the evaluated models<sup>47</sup>. When statistically significant differences were  
419 detected, Tukey’s Honestly Significant Difference (HSD) test was applied as a post-hoc anal-  
420 ysis to identify specific model pairs exhibiting significant differences while controlling the  
421 family-wise error rate<sup>48</sup>. In addition, the mean difference was computed separately for DSC  
422 and CSA evaluation metrics to quantify the average paired deviation between models. For  
423 DSC, the mean difference is defined as in Eq. (2), representing the average paired difference  
424 between MNSeg-Net and the reference model. For CSA, the mean difference is defined as in  
425 Eq. (3), representing the average paired difference between clinician-annotated and model-  
426 predicted values. Mean differences were not aggregated across heterogeneous metrics with  
427 opposing optimization directions and were used solely as descriptive measures of deviation.

$$428 \quad \Delta_{\text{DSC}} = \frac{1}{N} \sum_{i=1}^N (\text{DSC}_{\text{MNSeg-Net},i} - \text{DSC}_{\text{Ref},i}) \quad (2)$$

$$429 \quad \Delta_{\text{CSA}} = \frac{1}{N} \sum_{i=1}^N (\text{CSA}_{\text{Clinician},i} - \text{CSA}_{\text{Model},i}) \quad (3)$$

### 430 II.H.1. Equivalence Testing

431 In addition to conventional hypothesis testing, formal equivalence testing was performed  
432 using the two one-sided tests (TOST) procedure<sup>49</sup> to assess whether observed differences  
433 lay within predefined margins of practical relevance. Equivalence testing was applied to  
434 both DSC and CSA measurements, reflecting distinct analytical objectives. For DSC-based  
435 model-to-model comparisons, equivalence margins of  $\pm 0.01$ ,  $\pm 0.03$ , and  $\pm 0.05$  were consid-  
436 ered, representing increasingly relaxed criteria for practical similarity in segmentation per-  
437 formance commonly used in medical image analysis. For CSA evaluation, which represents

438 an agreement analysis between model-predicted measurements and clinician-annotated refer-  
 439 ence values, equivalence testing was conducted at margins of  $\pm 0.5 \text{ mm}^2$ , reflecting clinically  
 440 meaningful tolerance ranges for median nerve CSA assessment. Equivalence was concluded  
 441 only when the corresponding confidence intervals lay entirely within the specified equivalence  
 442 margins.

### 443 III. Results

Table 3: Performance comparison of deep learning models evaluated on a 10-subject test dataset across two anatomical regions: the wrist and wrist-to-elbow. Evaluation metrics include DSC, Precision (Prec), Recall (Rec), and Hausdorff Distance (HD).  $\uparrow$  indicates that higher values correspond to better performance, while  $\downarrow$  indicates that lower values correspond to better performance.

Model	At wrist region				Wrist to Elbow			
	DSC $\uparrow$	Prec $\uparrow$	Rec $\uparrow$	HD $\downarrow$	DSC $\uparrow$	Prec $\uparrow$	Rec $\uparrow$	HD $\downarrow$
UNet	0.943	0.940	0.946	4.657	0.787	0.791	0.789	14.421
SegNet	0.936	0.930	0.952	4.852	0.788	0.786	0.819	16.428
ResUNet	0.945	0.947	0.945	4.803	0.781	0.804	0.786	14.178
Attention-UNet	0.943	0.938	<b>0.953</b>	7.720	0.771	0.786	0.785	14.221
UNet++	0.942	0.935	<b>0.953</b>	7.130	0.755	0.781	0.766	16.550
BASNet	<b>0.947</b>	0.955	0.942	4.425	0.823	0.825	0.838	13.400
U2Net-full	0.945	0.957	0.937	4.605	<b>0.834</b>	<b>0.844</b>	0.840	13.301
<b>MNSeg-Net (Proposed)</b>	<b>0.947</b>	<b>0.965</b>	0.945	<b>4.177</b>	<b>0.834</b>	0.838	<b>0.846</b>	<b>12.835</b>

#### 444 III.A. Performance Comparison with Expert Annotation

445 Table 2 summarizes the models used in this study, including the parameter sizes, FLOPs per  
 446 frame, training duration, and number of frames processed per second of each model. Table  
 447 3 presents the average DSC, precision, recall, and HD across all frames for the ten test par-  
 448 ticipants. All methods exhibited high accuracy and effectiveness in segmenting the median  
 449 nerve at the wrist because of the consistent anatomy of the nerve and the lack of obstructing  
 450 tissues. However, the performance of basic UNet-based methods noticeably declined further  
 451 from the wrist towards the elbow, which is attributed to increased anatomical variability and  
 452 surrounding structures, such as muscles and blood vessels. The proposed model, with its  
 453 efficient subnetwork block design, demonstrated a superior ability to handle the complexities  
 454 of localizing the median nerve in these regions. Table 10 displays all metrics for individual  
 455 subjects (averaged across all frames) for the proposed model. Sample segmented images  
 456 for subject-1 obtained using various models are shown in Fig. 3. Each frame includes the  
 457 computed CSA of the median nerve and DSC between the segmented output and expert  
 458 annotation at the bottom. The results indicate that the proposed model outperforms exist-  
 459 ing deep learning models in segmenting the median nerve in ultrasound videos. As shown  
 460 in Fig. 4, the proposed MNSeg-Net model demonstrates a clearer separation between the  
 461 pixel-embedding representations of median nerve and background compared to those of the

462 other models.

Table 4: Comprehensive statistical comparison of segmentation models relative to the proposed MNSeg-Net (reference baseline) based on DSC. Mean differences and unadjusted p-values were computed using paired t-tests, and Cohen’s d values represent practical significance. Effect size interpretation: N = Negligible ( $d < 0.2$ ), S = Small ( $0.2 \leq d < 0.5$ ), M = Moderate ( $0.5 \leq d < 0.8$ ), and L = Large ( $d \geq 0.8$ ). An asterisk (\*) denotes a statistically significant difference from MNSeg-Net at  $p < 0.05$  after Tukey correction.

Model	Mean Difference	p-value (paired t-test)	Effect size Cohen’s d	p-value (Tukey adj)	95% CI (Lower)	95% CI (Upper)
UNet*	0.069	$8.70 \times 10^{-29}$	0.289 (S)	0.000	-0.088	-0.049
SegNet*	0.034	$6.56 \times 10^{-10}$	0.160 (N)	0.000	-0.053	-0.015
ResUNet*	0.054	$8.67 \times 10^{-19}$	0.229 (S)	0.000	-0.073	-0.034
Attention-UNet*	0.063	$3.41 \times 10^{-23}$	0.257 (S)	0.000	0.044	0.082
UNet++*	0.079	$4.06 \times 10^{-33}$	0.312 (S)	0.000	-0.098	-0.059
BASNet	0.011	$4.00 \times 10^{-2}$	0.053 (N)	0.738	-0.009	0.030
U2Net	-0.001	$9.11 \times 10^{-1}$	-0.003 (N)	1.000	-0.019	0.020

Table 5: Comprehensive statistical comparison of model-predicted CSA values relative to clinician-annotated  $CSA_{Act}$  (reference baseline). Mean differences and unadjusted p-values were computed using paired t-tests, and Cohen’s d values represent practical significance. Effect size interpretation: N = Negligible ( $d < 0.2$ ), S = Small ( $0.2 \leq d < 0.5$ ), M = Moderate ( $0.5 \leq d < 0.8$ ), and L = Large ( $d \geq 0.8$ ). An asterisk (\*) denotes a statistically significant difference from  $CSA_{Act}$  at  $p < 0.05$  after Tukey correction.

Model	Mean Diff ( $CSA_{Act} - CSA_{Cal}$ )	p-value (paired t-test)	Effect size Cohen’s d	p-value (Tukey adj)	95% CI (Lower)	95% CI (Upper)
UNet*	0.538	$1.34 \times 10^{-13}$	0.192 (N)	0.000	0.321	0.755
SegNet	-0.186	$3.80 \times 10^{-3}$	-0.075 (N)	0.168	-0.403	0.031
ResUNet*	0.659	$8.57 \times 10^{-23}$	0.255 (S)	0.000	0.442	0.876
Attention-UNet*	0.694	$7.40 \times 10^{-22}$	0.249 (S)	0.000	0.477	0.911
UNet++*	0.977	$3.99 \times 10^{-39}$	0.341 (S)	0.000	0.760	1.194
BASNet	-0.054	$3.41 \times 10^{-1}$	-0.025 (N)	0.999	-0.271	0.163
U2Net	0.105	$3.88 \times 10^{-2}$	0.053 (N)	0.882	-0.322	0.112
MNSeg-Net	-0.081	$1.14 \times 10^{-1}$	-0.041 (N)	0.976	-0.298	0.136

### 463 III.B. Statistical Analysis

#### 464 III.B.1. DSC-based Evaluation:

465 Table 4 summarizes the paired t-test results for DSC. All models, except U2Net, demon-  
 466 strated statistically significant differences compared to MNSeg-Net ( $p < 0.05$ ), with p-values  
 467 of  $8.70 \times 10^{-29}$  for UNet,  $6.56 \times 10^{-10}$  for SegNet, and  $4.06 \times 10^{-33}$  for UNet++. Despite  
 468 being statistically significant, the effect sizes were generally small (e.g., 0.160 for SegNet,  
 469 0.229 for ResUNet, and 0.312 for UNet++), suggesting a limited practical impact. U2Net  
 470 exhibited a negligible mean difference of  $-0.001$  and an effect size of  $-0.003$ , indicating that  
 471 its performance was statistically and practically similar to that of MNSeg-Net.

472 **III.B.2. CSA-based Evaluation:**

473 As shown in Table 5, MNSeg-Net produced a non-significant p-value of  $1.14 \times 10^{-1}$  and a  
 474 negligible effect size of  $-0.041$  compared with the clinician-annotated  $CSA_{Act}$ , demonstrat-  
 475 ing strong alignment with the clinical ground truth. Similarly, BASNet exhibited a non-  
 476 significant difference ( $p = 3.41 \times 10^{-1}$ ,  $d = -0.025$ ). Conversely, UNet ( $p = 1.34 \times 10^{-13}$ ,  
 477  $d = 0.192$ ), ResUNet ( $p = 8.57 \times 10^{-23}$ ,  $d = 0.255$ ), and UNet++ ( $p = 3.99 \times 10^{-39}$ ,  
 478  $d = 0.341$ ) showed significant difference and moderate effect sizes, suggesting notable devi-  
 479 ations from  $CSA_{Act}$ .

480 **III.B.3. Multiple Comparison Correction**

481 To address the risk of type I error arising from multiple comparisons, a one-way ANOVA  
 482 followed by Tukey’s HSD test was conducted. The ANOVA on the DSC yielded a highly  
 483 significant result ( $p = 1.72 \times 10^{-80} < 0.05$ ), indicating statistically significant differences  
 484 among the evaluated models. Tukey’s HSD analysis (Table 6) showed that MNSeg-Net  
 485 significantly outperformed UNet ( $p \approx 1 \times 10^{-308} < 0.05$ ), SegNet ( $p \approx 1 \times 10^{-308} < 0.05$ ),  
 486 ResUNet ( $p \approx 1 \times 10^{-308} < 0.05$ ), UNet++ ( $p \approx 1 \times 10^{-308} < 0.05$ ), and Attention-UNet  
 487 ( $p \approx 1 \times 10^{-308} < 0.05$ ). In contrast, BASNet ( $p = 0.738$ ) and U2Net ( $p = 1.000$ ) did not  
 488 exhibit statistically significant differences, confirming comparable performance.

489 For CSA, ANOVA similarly revealed significant group differences ( $p = 1.05 \times 10^{-134} <$   
 490  $0.05$ ). Tukey’s HSD results (Table 7) indicated a significant overestimation by UNet ( $p$   
 491  $\approx 1 \times 10^{-308} < 0.05$ ), ResUNet ( $p \approx 1 \times 10^{-308} < 0.05$ ), UNet++ ( $p \approx 1 \times 10^{-308} < 0.05$ ),  
 492 and Attention-UNet ( $p \approx 1 \times 10^{-308} < 0.05$ ). In contrast, MNSeg-Net ( $p = 0.9761$ ), BASNet  
 493 ( $p = 0.9988$ ), SegNet ( $p = 0.1678$ ), and U2Net ( $p = 0.8819$ ) did not differ significantly from  
 494  $CSA_{Act}$ , with MNSeg-Net exhibiting the lowest mean deviation of  $-0.0806$ .

495 Overall, MNSeg-Net demonstrated strong statistical and practical agreement with ex-  
 496 pert annotations, particularly in CSA estimation, where it showed minimal deviation from  
 497 the clinician-provided ground truth. Furthermore, it achieved consistently high DSC perfor-  
 498 mance and statistically validated improvements over several baseline models, while appro-  
 499 priately controlling for type I error through multiple-comparison correction.

Table 6: Two One-Sided Tests (TOST) equivalence test results for DSC scores of MNSeg-Net versus other models. The table reports mean differences with 95% confidence intervals and equivalence decisions at margins  $\pm 0.01$ ,  $\pm 0.03$ , and  $\pm 0.05$ . MNSeg-Net was statistically equivalent to U2Net at all tested margins, BASNet at  $\pm 0.03$  and  $\pm 0.05$ , and SegNet at  $\pm 0.05$  only. No equivalence was found with UNet, ResUNet, Attention-UNet, or UNet++.

Model	Mean Difference	95% CI Lower	95% CI Upper	Equivalent (margin=0.01)	Equivalent (margin=0.03)	Equivalent (margin=0.05)
UNet	0.0686	0.0605	0.0768	No	No	No
SegNet	0.0340	0.0269	0.0410	No	No	Yes
ResUNet	0.0537	0.0454	0.0621	No	No	No
Attention-UNet	0.0629	0.0541	0.0717	No	No	No
UNet++	0.0786	0.0694	0.0878	No	No	No
BASNet	0.0107	0.0039	0.0176	No	Yes	Yes
U2Net	-0.0006	-0.0066	0.0055	Yes	Yes	Yes

Table 7: TOST equivalence test results for CSA (actual vs. computed) at margin  $\pm 0.5$ . The proposed MNSeg-Net, SegNet, BASNet, and U2Net models met the equivalence criterion, whereas UNet, ResUNet, Attention-UNet, and UNet++ did not. At the equivalence margin of  $\pm 0.5$ , MNSeg-Net, SegNet, BASNet, and U2Net demonstrated statistical equivalence with the clinician-annotated CSA values. In contrast, UNet, ResUNet, Attention-UNet, and UNet++ did not achieve statistical equivalence under this criterion.

Model	Mean Difference	95% CI Lower	95% CI Upper	Equivalent ( $\pm 0.5$ )
UNet	0.5383	0.4178	0.6588	No
SegNet	-0.1859	-0.2840	-0.0878	Yes
ResUNet	0.6595	0.5582	0.7609	No
Attention-UNet	0.6945	0.5808	0.8083	No
UNet++	0.9778	0.8573	1.0983	No
BASNet	-0.0540	-0.1324	0.0244	Yes
U2Net	0.1046	0.0386	0.1705	Yes
MNSeg-Net	-0.0807	-0.1491	-0.0123	Yes

Table 8: TOST equivalence analysis of frame-level DSC between MNSeg-Net and U2Net. It reports the mean difference, 95% confidence interval, and equivalence decisions under predefined DSC margins.

Mean Difference	95% CI Lower	95% CI Upper	Equivalent (margin=0.01)	Equivalent (margin=0.03)	Equivalent (margin=0.05)
-0.0006	-0.0066	0.0055	Yes	Yes	Yes

Table 9: TOST equivalence analysis of frame-level CSA estimates between MNSeg-Net and U2Net. It reports the mean difference, 95% confidence interval, and equivalence decisions under clinically defined CSA margins (in  $\text{mm}^2$ ).

Mean Difference	95% CI Lower	95% CI Upper	Equivalent (margin=0.1)	Equivalent (margin=0.3)	Equivalent (margin=0.5)
0.1852	0.1359	0.2346	No	Yes	Yes

### 500 III.B.4. TOST Equivalence Analysis of DSC and CSA

501 The results, summarized in Table 6, indicated that U2Net was statistically equivalent to  
 502 MNSeg-Net even at the strictest margin of  $\pm 0.01$ , while BASNet became equivalent at the  
 503 wider margin of  $\pm 0.02$ . At the broadest margin of  $\pm 0.05$ , equivalence was observed for  
 504 U2Net, BASNet, and SegNet. In contrast, UNet consistently failed to meet equivalence  
 505 criteria across all tested margins, highlighting a significant performance gap compared to  
 506 MNSeg-Net. Similarly, Attention-UNet, ResUNet, and UNet++ did not achieve equivalence  
 507 within these margins, despite showing competitive DSC scores. These findings reinforce  
 508 that MNSeg-Net exhibits performance that is statistically indistinguishable from several  
 509 advanced architectures under practical equivalence thresholds, while maintaining superiority  
 510 over models such as UNet.

511 The results, presented in Table 7, show that at the narrowest margin of  $\pm 0.1$ , none of the  
 512 models achieved equivalence with the actual CSA, suggesting that small differences persist  
 513 at very strict thresholds. When the margin was relaxed to  $\pm 0.2$ , BASNet satisfied the equiv-  
 514 alence criteria, indicating its predictions were statistically indistinguishable from the ground

515 truth within this tolerance. At  $\pm 0.5$ , a broader set of models—including U2Net, MNISeg-Net,  
516 BASNet, and SegNet—achieved equivalence, while UNet continued to show significant deviation.  
517 At the widest tested margin of  $\pm 1.0$ , all models except UNet demonstrated equivalence  
518 with the actual CSA. These findings highlight that while certain models such as U2Net and  
519 BASNet align closely with the actual CSA even under tighter thresholds, UNet consistently  
520 fails to meet equivalence criteria, reflecting its systematic bias in CSA estimation.

Table 10: Subject-wise evaluation of the proposed MNISeg-Net model from wrist to elbow, including segmentation metrics (DSC, precision, recall, HD) and a statistical comparison between the predicted CSA values from MNISeg-Net ( $CSA_{Cal}$ ) and clinician-annotated CSA values ( $CSA_{Act}$ ). Paired t-tests, Cohen’s  $d$  effect sizes, mean differences, and confidence intervals were computed across 300 image frames per subject to assess the clinical reliability of the MNISeg-Net. Equivalence was assessed using the TOST procedure with predefined margins of  $\pm 0.5$  mm<sup>2</sup>. Asterisks (\*) denote statistically significant differences ( $p < 0.05$ ). Effect size categories: N = Negligible ( $d < 0.2$ ), S = Small ( $0.2 \leq d < 0.5$ ), M = Moderate ( $0.5 \leq d < 0.8$ ), L = Large ( $d \geq 0.8$ ).

Subject ID	DSC $\uparrow$	Prec $\uparrow$	Rec $\uparrow$	HD $\downarrow$	$CSA_{Cal}$	$CSA_{Act}$	t-test p-value	Cohen’s $d$ (N/S/M/L)	Mean Difference	95% CI Lower	95% CI Upper	Equivalent ( $\pm 0.5$ )
1	0.891	0.854	0.896	8.993	8.297	7.937	0.0015*	0.254 (S)	0.360	0.198	0.522	Yes
2	0.898	0.893	0.913	7.081	8.319	8.126	0.0124*	0.145 (N)	0.193	0.045	0.341	Yes
3	0.786	0.812	0.776	19.909	8.666	9.278	0.0000*	-0.289 (S)	-0.612	-0.808	-0.416	No
4	0.704	0.695	0.738	23.511	10.140	9.292	0.0000*	0.327 (S)	0.849	0.618	1.080	No
5	0.840	0.807	0.900	11.222	9.539	8.484	0.0000*	0.501 (M)	1.055	0.826	1.284	No
6	0.890	0.888	0.899	8.005	8.499	8.364	0.0457*	0.116 (N)	0.136	0.006	0.266	Yes
7	0.878	0.880	0.886	9.705	8.849	8.860	0.8963	-0.008 (N)	-0.011	-0.141	0.119	Yes
8	0.768	0.773	0.783	14.741	7.404	7.404	0.9987	-0.000 (N)	-0.000	-0.132	0.132	Yes
9	0.818	0.845	0.810	14.940	6.791	7.023	0.0253*	-0.130 (N)	-0.233	-0.416	-0.050	Yes
10	0.888	0.932	0.861	10.244	8.776	9.705	0.0000*	-0.501 (M)	-0.928	-1.124	-0.732	No

### 521 III.B.5. Subject-wise Statistical Analysis

522 A comprehensive subject-wise statistical analysis was performed on Dataset-1 to compare  
523 the clinician-annotated CSA ( $CSA_{Act}$ ) with the MNISeg-Net-predicted CSA ( $CSA_{Cal}$ ). For  
524 each subject, paired t-tests, mean differences, and effect sizes (Cohen’s  $d$ ) were calculated  
525 across approximately 300 images. The detailed results of this analysis are summarized in  
526 Table 10. The analysis revealed that although some subjects exhibited statistically signifi-  
527 cant differences ( $p < 0.05$ ) between the predicted and clinician-annotated CSA values, the  
528 practical impact of these differences was minimal. Based on Cohen’s  $d$  categorization, five  
529 subjects demonstrated negligible effect sizes ( $|d| < 0.2$ ), three subjects exhibited small effect  
530 sizes ( $0.2 \leq |d| < 0.5$ ), and two subjects exhibited moderate effect sizes ( $0.5 \leq |d| < 0.8$ ).  
531 Notably, no subject exhibited a large effect size ( $d \geq 0.8$ ), indicating that the MNISeg-Net  
532 CSA predictions remained clinically reliable across the entire cohort. Even in cases where  
533 statistically significant differences were observed, the associated effect sizes were predom-  
534 inantly negligible or small, confirming that these statistical deviations did not adversely  
535 affect clinical decision-making. In addition to conventional hypothesis testing, a two one-  
536 sided tests (TOST) equivalence analysis was performed at the frame level for each subject  
537 to explicitly assess the practical equivalence between  $CSA_{Cal}$  and  $CSA_{Act}$ . Equivalence mar-  
538 gins were predefined as  $\pm 0.5$  mm<sup>2</sup>, reflecting clinically acceptable deviations in median nerve  
539 CSA measurements. Based on the TOST procedure, six out of ten subjects demonstrated  
540 statistical equivalence, with the corresponding confidence intervals lying entirely within the  
541 predefined equivalence bounds, as summarized in Table 10. Overall, these findings demon-

542 strate that MNSeg-Net consistently provides robust and clinically reliable CSA estimations  
543 across diverse populations, maintaining diagnostic reliability and practical applicability de-  
544 spite minor subject-specific variability.

### 545 III.C. Performance Comparison - Model Analysis

546 The proposed MNSeg-Net model is designed to achieve competitive segmentation perfor-  
547 mance with substantially reduced computational complexity. Compared with the full-size  
548 U2Net, which contains 44M trainable parameters, MNSeg-Net requires only 2.46 M param-  
549 eters, corresponding to an approximate reduction of 94.4% in the parameter count. This  
550 reduction is at least  $12\times$  relative to other compact UNet-based architectures, such as the  
551 simple UNet (25 M), while maintaining a segmentation performance comparable to that of  
552 U2Net. A detailed quantitative comparison with the full-size U2Net is provided in Table 1,  
553 which reports both block-wise and overall metrics, including the memory footprint, model  
554 size, computational complexity (FLOPs), and parameter count across the encoder, bottle-  
555 neck, and decoder stages. As summarized in the table, the overall computational cost is  
556 reduced from 82.24 G FLOPs for U2Net to 55.72 G FLOPs for MNSeg-Net, representing a  
557 reduction of approximately 32.2%. Furthermore, the model size is reduced from 167.85 MB  
558 to 9.37 MB, yielding an approximate 94.4% reduction in storage requirements. These reduc-  
559 tions quantitatively highlight the lightweight design of the proposed architecture. Despite  
560 the substantial reduction in model complexity, MNSeg-Net remains suitable for real-time  
561 deployment. During test-time inference, the model requires a memory footprint of approxi-  
562 mately 3.14GB, primarily because of the subnetwork module, while retaining computational  
563 efficiency. The proposed model achieves a single-frame inference time of 0.023s for an in-  
564 put resolution of  $448 \times 320$ , enabling the processing of approximately 43 frames per second.  
565 These characteristics allow MNSeg-Net to operate smoothly on standard computational se-  
566 tups, such as a laptop equipped with an Intel i9 processor and an NVIDIA GeForce 3090  
567 GPU with 6 GB memory for test-time inference, as summarized in Table 2.

### 568 III.D. Effect of Training Data Size on Performance

569 The robustness and generalisation capability of the proposed MNSeg-Net model were eval-  
570 uated under varying data availability conditions. Specifically, the model was trained using  
571 only 25% and 50% of the available training data, and its results were compared against  
572 the same segmentation networks considered in the previous experiments. As shown in Ta-  
573 ble 11, MNSeg-Net consistently outperformed the competing models under both data con-  
574 ditions, achieving high DSC, precision, and recall values while maintaining relatively lower  
575 HD. Among the competing models, BASNet and U2Net demonstrated relatively strong per-  
576 formance, with BASNet achieving high DSC and recall and U2Net showing competitive  
577 precision and recall. These findings highlight the ability of MNSeg-Net to deliver stable  
578 segmentation performance even with substantially reduced training data, thereby demon-  
579 strating strong generalization potential in data-constrained scenarios.

580 To quantitatively validate these observations, a comprehensive statistical analysis in-  
581 cluding one-way ANOVA, paired  $t$ -tests, and Cohen’s  $d$  effect size estimation was conducted  
582 on the frame-wise DSC scores for the 25% and 50% training data settings. The detailed  
583 results are provided in Section 4 (Tables 1 and 2) of the Supplementary Material.

### 584 III.E. Perturbation Study

585 To assess the robustness of the proposed model, a perturbation study was conducted by  
586 introducing speckle noise of varying mean levels (0.1, 0.2, and 0.3) to the input data. The  
587 performance of MNSeg-Net was compared with several UNet-based models under these per-  
588 turbations.

589 At a low noise level (mean = 0.1), all UNet-based models performed well at the wrist  
590 (DSC  $\approx$  0.87-0.94) but exhibited a noticeable drop from the wrist to the elbow (e.g., UNet  
591 elbow DSC = 0.497). In contrast, the proposed model remained stable at the wrist (DSC =  
592 0.946) and showed only a slight decrease at the elbow (DSC = 0.793). At the medium noise  
593 level (mean = 0.2), the UNet-based models degraded severely, particularly from the wrist  
594 to the elbow where performance nearly collapsed (most elbow DSC < 0.20). The proposed  
595 model also showed a drop at the wrist (DSC = 0.868,  $\sim$ 8% reduction), but retained moderate  
596 accuracy at the elbow (DSC = 0.568). At the high noise level (mean = 0.3), most baseline  
597 models failed completely in both regions, except for Attention-UNet, which maintained some  
598 wrist performance (DSC = 0.545) but failed at the elbow. Remarkably, the proposed model  
599 remained effective, achieving 73% precision at the wrist and 33% precision from the wrist to  
600 the elbow. Average performance across the ten subjects at each noise level is summarized  
601 in Table 12.

602 Across all noise perturbation levels ( $\alpha = 0.1, 0.2, \text{ and } 0.3$ ), statistical analysis demon-  
603 strated a strong and reliable effect of model choice on segmentation performance. One-way  
604 ANOVA revealed highly significant inter-model differences at each noise level (for  $\alpha = 0.1$ :  
605  $F = 565.20$ ,  $\alpha = 0.2$ :  $F = 580.81$ , and  $\alpha = 0.3$ :  $F = 640.49$ , all with  $p \approx 1 \times 10^{-308}$ ), con-  
606 firming that performance variations are not attributable to random chance. Paired  $t$ -tests  
607 comparing MNSeg-Net against competing models showed statistically significant improve-  
608 ments ( $p < 0.05$ ) at all noise levels; however, the magnitude of these improvements varied  
609 systematically with noise severity. At moderate noise ( $\alpha = 0.1$ ), Cohen’s  $d$  values ranged  
610 from negligible to small for most baselines ( $d = 0.09\text{--}0.32$ ), with larger effects observed  
611 only for UNet and UNet++ ( $d = 0.80\text{--}1.04$ ), corresponding to mean DSC gains of approx-  
612 imately 0.02–0.41. Under severe noise ( $\alpha = 0.2$ ), effect sizes increased substantially across  
613 all comparisons, with Cohen’s  $d$  ranging from 0.67 to 1.27 and mean DSC improvements  
614 between 0.26 and 0.48, indicating large and practically meaningful robustness gains. At ex-  
615 treme noise levels ( $\alpha = 0.3$ ), MNSeg-Net maintained statistically significant improvements  
616 over all baselines with consistently moderate effect sizes ( $d = 0.52\text{--}0.70$ ), reflecting stable  
617 performance and graceful degradation rather than catastrophic failure. To control for mul-  
618 tiplicity and mitigate inflation of Type-I error arising from repeated pairwise comparisons,  
619 Tukey’s HSD post-hoc analysis was applied following the ANOVA at each noise level. After

adjustment, MNSeg-Net remained significantly superior to all baseline models at  $\alpha = 0.2$  and  $\alpha = 0.3$  (all adjusted  $p < 0.001$ ). At  $\alpha = 0.1$ , all improvements remained statistically significant after correction except for the comparison with U2Net (adjusted  $p = 0.1617$ ), confirming comparable performance between MNSeg-Net and U2Net under moderate noise. The multiplicity-adjusted analyses confirm that the majority of observed performance improvements remain statistically significant after controlling for family-wise error, thereby reinforcing the robustness and reliability of the reported gains under noisy conditions.

Overall, these results indicate that MNSeg-Net exhibits greater robustness to input perturbations than UNet-based baselines. Detailed statistical results, including repeated-measures one-way ANOVA, paired  $t$ -tests, Cohen’s  $d$  effect size estimates, and Tukey’s HSD post-hoc analyses with multiplicity-adjusted  $p$ -values and confidence intervals, are provided in Section 5 (Tables 3–5) of the Supplementary Material.

Table 11: Performance comparison of different models on 25% and 50% training data.

Model	25% data				50% data			
	DSC $\uparrow$	Prec $\uparrow$	Rec $\uparrow$	HD $\downarrow$	DSC $\uparrow$	Prec $\uparrow$	Rec $\uparrow$	HD $\downarrow$
UNet	0.685	0.720	0.686	19.158	0.776	0.808	0.777	20.237
SegNet	0.674	0.698	0.690	26.820	0.766	0.775	0.785	22.211
ResUnet	0.653	0.676	0.667	27.018	0.727	0.749	0.742	19.080
Attention-UNet	0.642	0.668	0.649	22.703	0.741	0.762	0.748	14.976
UNet++	0.623	0.652	0.627	38.567	0.702	0.727	0.714	18.236
BASNet	0.746	0.724	0.791	24.320	0.803	0.806	0.819	16.202
U2Net	0.736	0.731	0.758	25.552	0.788	0.760	0.836	17.898
MNSeg-Net	0.751	0.743	0.778	21.573	0.804	0.804	0.824	16.430

Table 12: Performance of different models under speckle noise perturbation at mean levels 0.1, 0.2, and 0.3, evaluated at the wrist and wrist-to-elbow regions.

Mean	Model	At Wrist				Wrist to Elbow			
		DSC $\uparrow$	Precision $\uparrow$	Recall $\uparrow$	HD $\downarrow$	DSC $\uparrow$	Precision $\uparrow$	Recall $\uparrow$	HD $\downarrow$
0.1	UNet	0.872	0.881	0.889	33.181	0.497	0.544	0.501	85.370
	SegNet	0.932	0.939	0.934	5.500	0.699	0.735	0.700	22.140
	ResUNet	0.944	0.941	0.950	5.363	0.734	0.755	0.747	20.075
	UNet++	0.872	0.883	0.877	18.910	0.382	0.418	0.376	54.120
	Attention-UNet	0.939	0.923	0.961	8.272	0.695	0.693	0.730	19.913
	U2Net	0.934	0.947	0.928	5.602	0.770	0.779	0.779	19.516
	BASNet	0.940	0.941	0.942	6.365	0.756	0.747	0.785	16.644
	MNSeg-Net	0.946	0.949	0.945	4.617	0.793	0.788	0.814	17.758
0.2	UNet	0.416	0.726	0.331	51.324	0.087	0.153	0.070	174.371
	SegNet	0.818	0.881	0.787	10.130	0.176	0.214	0.164	75.220
	ResUNet	0.832	0.870	0.809	9.524	0.194	0.218	0.186	113.635
	UNet++	0.480	0.788	0.398	33.184	0.101	0.165	0.085	135.286
	Attention-UNet	0.841	0.824	0.864	15.205	0.249	0.256	0.258	95.342
	U2Net	0.539	0.613	0.494	18.048	0.303	0.335	0.292	40.582
	BASNet	0.797	0.846	0.771	14.912	0.176	0.192	0.170	39.715
	MNSeg-Net	0.868	0.897	0.853	12.047	0.568	0.594	0.563	47.241
0.3	UNet	0.000	0.017	0.001	70.901	0.001	0.002	0.001	259.379
	SegNet	0.396	0.698	0.318	35.870	0.082	0.140	0.066	77.548
	ResUNet	0.235	0.609	0.167	39.594	0.054	0.119	0.041	121.031
	UNet++	0.000	0.003	0.001	73.437	0.001	0.002	0.005	163.039
	Attention-UNet	0.545	0.507	0.530	27.683	0.153	0.170	0.149	134.967
	U2Net	0.002	0.005	0.001	32.561	0.006	0.007	0.004	57.043
	BASNet	0.002	0.002	0.001	37.215	0.001	0.001	0.001	51.900
	MNSeg-Net	0.600	0.734	0.536	26.082	0.257	0.335	0.228	68.680

## IV. Clinical Evaluation and Analysis

In clinical practice, CSA is a widely accepted imaging parameter for diagnosing conditions such as CTS, inflammation, and edema in nerves, tendons, and ligaments. A CSA greater than  $12 \text{ mm}^2$  is typically used to indicate positive for CTS<sup>32,33,34</sup>. To clinically validate the proposed MNSeg-Net model, its segmentation outputs were used to compute the CSA values, which were then compared with the clinician-annotated ground truth ( $\text{CSA}_{\text{Act}}$ ) across a dataset comprising 30 patients. For each patient, six expert annotations were collected at the wrist level, with patients grouped into two clinically relevant categories based on the CSA threshold: NORMAL ( $\text{CSA} < 12 \text{ mm}^2$ ) and CTS-positive ( $\text{CSA} \geq 12 \text{ mm}^2$ ).

Table 13 summarizes the patient-wise segmentation performance, including the DSC, precision, recall, and CSA values for both  $\text{CSA}_{\text{Act}}$  and model-predicted  $\text{CSA}_{\text{Cal}}$ . MNSeg-Net demonstrated consistently high segmentation accuracy, with average DSC values exceeding 0.93 in both the NORMAL and CTS-positive groups. Fig. 5 shows sample images from real-time testing across the wrist, forearm, and elbow regions. The real-time performance of the system can also be viewed in the accompanying demonstration [video link](#).

The statistical analysis comparing the MNSeg-Net predicted CSA ( $\text{CSA}_{\text{Cal}}$ ) with clinician-annotated CSA ( $\text{CSA}_{\text{Act}}$ ) values is summarized in Table 14. In the NORMAL group, no statistically significant difference was observed ( $p = 0.208$ ), with a small negative mean difference of  $-0.206 \text{ mm}^2$  and a small effect size (Cohen’s  $d = -0.222$ ), indicating close

651 agreement between the model-predicted and clinician-annotated CSA values. MNSeg-Net  
652 also demonstrated a strong segmentation performance in this group, achieving an average  
653 DSC of 0.9381, further reinforcing its clinical reliability. For the CTS-positive group, a sta-  
654 tistically significant difference was detected ( $p = 0.0094$ ), with a positive mean difference of  
655  $+0.709 \text{ mm}^2$  and a moderate effect size (Cohen  $d = 0.529$ ). Despite the statistically signifi-  
656 cant difference, MNSeg-Net maintained high segmentation accuracy in this group, achieving  
657 an average DSC of 0.9544, thus highlighting its robustness across different clinical scenarios.  
658 The overall dataset analysis revealed no statistically significant difference ( $p = 0.296$ ), with  
659 a mean difference of  $+0.163 \text{ mm}^2$  and a negligible to small effect size (Cohen’s  $d = 0.194$ ).  
660 These results demonstrate that MNSeg-Net closely approximated clinician annotations in  
661 the NORMAL group and the overall dataset, with minor deviations. Although a moderate  
662 difference was noted in the CTS-positive group, the magnitude of the discrepancy remained  
663 within a clinically acceptable range, supporting the reliability of the model across different  
664 patient categories.

665 Overall, the strong alignment between the model-predicted and clinician-measured CSA  
666 values, combined with consistently high DSC performance, highlights the clinical applicabil-  
667 ity, diagnostic reliability, and real-world deployment potential of MNSeg-Net.

Table 13: Patient-wise segmentation performance metrics, including DSC, precision, recall, and CSA estimations, comparing clinician-annotated ( $CSA_{Act}$ ) and MNSeg-Net-predicted ( $CSA_{Cal}$ ) values for 30 patients, categorized into NORMAL and CTS-positive groups.  $\uparrow$  indicates that higher values correspond to better performance, while  $\downarrow$  indicates that lower values correspond to better performance.

S.No	Label	DSC $\uparrow$	Prec $\uparrow$	Rec $\uparrow$	$CSA_{act}$ (mm $^2$ )	$CSA_{Cal}$ (mm $^2$ )
1	NORMAL	0.933	0.829	0.991	8.413	9.036
2	NORMAL	0.952	0.949	0.892	11.749	10.568
3	NORMAL	0.963	0.916	0.946	10.910	10.493
4	NORMAL	0.960	0.931	0.927	10.080	9.500
5	NORMAL	0.917	0.819	0.973	8.857	9.612
6	NORMAL	0.963	0.910	0.947	10.078	9.735
7	NORMAL	0.961	0.888	0.976	9.900	9.963
8	NORMAL	0.897	0.772	0.987	8.471	9.748
9	NORMAL	0.934	0.840	0.972	8.578	8.959
10	NORMAL	0.896	0.763	0.996	7.132	8.213
11	NORMAL	0.946	0.854	0.986	9.538	9.954
12	NORMAL	0.929	0.828	0.987	9.848	10.640
13	NORMAL	0.932	0.841	0.969	9.301	9.798
14	NORMAL	0.893	0.757	0.990	5.382	6.178
15	NORMAL	0.945	0.862	0.974	8.604	8.850
16	NORMAL	0.947	0.869	0.968	8.589	8.766
17	NORMAL	0.965	0.908	0.953	7.770	7.527
18	NORMAL	0.952	0.930	0.911	9.361	8.692
	<b>NORMAL</b>	<b>0.938</b>	<b>0.859</b>	<b>0.964</b>	<b>9.030</b>	<b>9.240</b>
19	CTS positive	0.976	0.938	0.967	16.629	16.242
20	CTS positive	0.957	0.964	0.899	15.152	13.684
21	CTS positive	0.966	0.937	0.939	13.878	13.229
22	CTS positive	0.961	0.937	0.930	16.851	15.919
23	CTS positive	0.967	0.919	0.960	12.701	12.463
24	CTS positive	0.924	0.940	0.856	13.053	11.397
25	CTS positive	0.939	0.930	0.894	12.940	11.739
26	CTS positive	0.948	0.973	0.876	13.786	12.079
27	CTS positive	0.951	0.935	0.924	17.935	16.946
28	CTS positive	0.967	0.936	0.949	16.389	15.774
29	CTS positive	0.942	0.854	0.990	12.056	12.826
30	CTS positive	0.955	0.882	0.983	15.443	15.952
	<b>CTS positive</b>	<b>0.954</b>	<b>0.929</b>	<b>0.931</b>	<b>14.730</b>	<b>14.020</b>

Table 14: Statistical analysis on clinical dataset for comparison of MNSeg-Net predicted CSA values from clinician-annotated  $CSA_{Act}$ . The table includes mean differences, p-values from paired t-tests and Cohen’s d effect sizes. (\*) indicates a statistically significant difference from  $CSA_{Act}$  at  $p < 0.05$ . Effect size abbreviations: N = Negligible ( $d < 0.2$ ), S = Small ( $0.2 \leq d < 0.5$ ), M = Moderate ( $0.5 \leq d < 0.8$ ), L = Large ( $d \geq 0.8$ ).

Group	p-value (t-test)	Mean Diff. ( $CSA_{Act} - CSA_{Cal}$ )	Effect Size (Cohen’s d)
NORMAL	0.208	-0.206	-0.222 (S)
CTS positive	0.009*	0.709	0.529 (M)
Overall	0.296	0.163	0.194 (N)

## 668 IV.A. Ablation Study

669 To verify the effectiveness of the proposed MNSeg-Net, ablation studies were conducted on  
 670 the following four aspects: (i) loss function, (ii) deep supervision, (iii) subnetwork block, and  
 671 (iv) architecture of the main network. All the ablation studies used the same implementation  
 672 setup. Table 15 summarizes the ablation configurations and their corresponding DSC values  
 673 for the wrist-to-elbow region. Among these, Experiment 9 represents the complete MNSeg-  
 674 Net configuration, incorporating the redesigned U2Net backbone, deep supervision, MSFF  
 675 module, and the proposed hybrid loss function, and therefore serves as the reference model  
 676 for all ablation comparisons. The remaining experiments progressively disable or replace  
 677 individual components to quantify their impact on segmentation performance.

### 678 IV.A.1. Ablation on Loss Function

679 In this ablation study, the model was trained using alternative loss functions, including BCE,  
 680 Dice, and logcosh losses, to evaluate the contribution of the proposed hybrid loss function.  
 681 In the full MNSeg-Net configuration, a hybrid loss function consisting of BCE, logcosh, and  
 682 SSIM losses was utilized. Specifically, Experiments 3 and 4 replace the hybrid loss with  
 683 logcosh and BCE loss, respectively. Compared to the full MNSeg-Net configuration (Ex-  
 684 periment 9, DSC = 0.834), using logcosh loss alone (Experiment 3, DSC = 0.822) results  
 685 in an approximate 1.2% reduction in DSC, while replacing the hybrid loss with BCE loss  
 686 (Experiment 4, DSC = 0.791) leads to a more substantial performance degradation of ap-  
 687 proximately 4.4%. As shown in Table 15, substituting the hybrid loss with individual loss  
 688 functions consistently results in lower DSC values, with the most pronounced degradation  
 689 observed for BCE-only training. These findings confirm the effectiveness of the hybrid loss  
 690 formulation in improving segmentation performance.

Table 15: Ablation study results with the data from wrist to elbow region.

Exp. No.	Ablation paramters						DSC $\uparrow$
	Main-Network (U2Net redesigned)	Deep Supervision	MSFF Module	BCE loss	logcosh loss	Hybrid loss	
1	✓	X	X	-	-	✓	0.804
2	✓	✓	X	-	-	✓	0.822
3	✓	✓	✓	-	✓	-	0.822
4	✓	✓	✓	✓	-	-	0.791
5	✓	✓	X	-	✓	-	0.823
6	✓	✓	X	✓	-	-	0.791
7	X	✓	✓	-	-	✓	0.819
8	✓	X	✓	-	-	✓	0.806
<b>9</b>	✓	✓	✓	-	-	✓	<b>0.834</b>

#### 691 IV.A.2. Ablation on Deep Supervision

692 Section II.B.2 outlines the network training strategy incorporating deep supervision. To  
693 evaluate its contribution, an ablation study was conducted in which deep supervision was  
694 removed and the loss was computed solely from the final-stage output. Experiments 1 and 8  
695 correspond to configurations without deep supervision, while Experiments 2 and 9 represent  
696 their respective deep-supervised counterparts. Importantly, the paired comparisons (Experi-  
697 ment 1 vs. Experiment 2) and (Experiment 8 vs. Experiment 9) differ only in the use of deep  
698 supervision, thereby isolating its effect. As summarized in Table 15, the absence of deep su-  
699 pervision results in a consistent degradation in segmentation accuracy, with DSC reductions  
700 of approximately 1.8% when comparing Experiment 1 (DSC = 0.804) to Experiment 2 (DSC  
701 = 0.822), and approximately 2.8% when comparing Experiment 8 (DSC = 0.806) to Experi-  
702 ment 9 (DSC = 0.834). In addition to the accuracy degradation, removing deep supervision  
703 also adversely affected training dynamics, requiring substantially more epochs to converge  
704 (approximately 25–30 epochs) and resulting in longer training times of 12–15 hours, despite  
705 reduced GPU memory consumption.

#### 706 IV.A.3. Ablation on Sub-Network Block

707 In the proposed network architecture, the output feature vector of the subnetwork block  
708 is integrated into the final stage and contributes to the loss function computation as part  
709 of deep supervision. The contribution of the subnetwork block is primarily evaluated by  
710 comparing Experiment 2 with the full MNISeg-Net configuration (Experiment 9), as these two  
711 configurations differ only in the presence of the subnetwork block while all other components  
712 are kept identical. Removing the subnetwork leads to a reduction in the DSC from 0.834  
713 (Experiment 9) to 0.822 (Experiment 2), corresponding to an absolute performance drop of  
714 approximately 1.2%, as summarized in Table 15. This controlled degradation quantifies the  
715 contribution of the subnetwork block to the overall segmentation performance.

716 Additional evidence is provided by Experiments 5 and 6, which also exclude the subnet-  
717 work block but employ different loss functions. Although these configurations are not strictly  
718 comparable due to the change in loss formulation, they exhibit a consistent degradation in  
719 performance, with DSC values of 0.823 in Experiment 5 (approximately 1.1% reduction) and  
720 0.791 in Experiment 6 (approximately 4.3% reduction), further supporting the relevance of  
721 the subnetwork block.

#### 722 IV.A.4. Ablation on Main-Network Decoder Blocks

723 The proposed network consists of five encoder–decoder blocks, one bottleneck block, and one  
724 subnetwork module. In Experiment 7, the decoders of the main network were removed, such  
725 that the segmentation output and loss computation relied solely on the subnetwork. This  
726 configuration enables direct evaluation of the importance of the redesigned main-network  
727 decoder. Under this setting, the architecture effectively resembles a U2Net encoder cou-

728 pled with a Half-UNet<sup>50</sup>-style decoder, resulting in reduced representational capacity and  
729 parameter count. Consequently, a slight performance degradation of approximately 1.2% in  
730 DSC was observed compared to the full configuration (Experiment 9), confirming that the  
731 joint use of the redesigned decoder and the subnetwork is essential for achieving optimal  
732 segmentation performance.

733 To quantitatively validate the qualitative trends observed in the ablation experiments,  
734 a comprehensive statistical analysis was performed across all ablation configurations. This  
735 analysis included one-way ANOVA, paired *t*-tests, mean DSC-score differences, and corre-  
736 sponding Cohen’s *d* effect size measures, and is summarized in section 5 (Table 6) of the  
737 Supplementary Material. The statistical results confirm that the complete MNSEg-Net con-  
738 figuration (Experiment 9) consistently achieves statistically significant performance improve-  
739 ments over all ablation variants. Although the associated effect sizes are generally small,  
740 they indicate statistically supported and practically meaningful incremental gains, thereby  
741 quantitatively substantiating the importance of jointly employing the redesigned backbone,  
742 deep supervision, subnetwork block, and hybrid loss function in the proposed architecture.

## 743 V. Discussion

744 This study introduced an efficient, lightweight deep learning model for median nerve seg-  
745 mentation in ultrasound images. Ultrasound is well known for its real-time output during  
746 probing. To leverage this real-time display, the proposed model was designed to infer the  
747 output in real time. With this objective in mind, the model was constrained in its design to  
748 limit memory usage (lower footprint) and computational complexity (lower FLOPS).

### 749 V.A. Model Design for Balancing Accuracy and Complexity

750 The model design phase began with the adoption of U2Net, as its full-size model demon-  
751 strated superior performance in median nerve segmentation compared to other UNet-based  
752 models. Experiments with U2Net-p (an existing lightweight variant from<sup>30</sup>) revealed its  
753 inability to handle the complexity of median nerve segmentation owing to its limited pa-  
754 rameter capacity (4.7 MB parameter space). However, the actual full-size U2Net model  
755 requires higher computation and more weight parameters, as evidenced in Table 2. During  
756 model design, the main challenge was balancing the model accuracy with the computational  
757 complexity and weight parameters. Increasing the accuracy substantially increased the com-  
758 putational complexity and weight parameters, whereas reducing the complexity and weight  
759 decreased the accuracy. To address this challenge, a lightweight version of the U2Net model  
760 was redesigned, featuring an uniform number of channels in both the encoder and decoder,  
761 instead of the traditional exponential growth and shrinkage. This architectural modification  
762 results in a model that is heavier than U2Net-p but substantially lighter than the full-size  
763 U2Net in terms of parameter space (MNSEg-Net: 2.46M parameters vs. U2Net: 44M)  
764 and computational cost (MNSEg-Net: 56G FLOPS vs. U2Net: 82.53G), as quantitatively

demonstrated by the block-wise comparison in Table 1. Additionally, a subnetwork module was introduced that leveraged existing encoding representations without notably increasing the computational complexity or memory usage. This additional block does not require separate computations or memory for the encoder. Similarly, a half-U-Net-style decoder allows direct upsampling of all-level encoding representations without adding extra memory and computation. Thus, the additional block improves segmentation accuracy without substantial overhead in terms of memory or computational load. Since the core network is a lightweight U2Net-derived architecture (Table 1), the proposed MNSeg-Net preserves overall computational efficiency. The subnetwork module processes the multiscale and multichannel encoding representations of the main network. By integrating all features, the subnetwork generates a final feature encompassing a broader image context and applies implicit regularization by selecting important features from the combined-encoding representations. Finally, as shown in Table 2, the proposed MNSeg-Net achieves efficient performance, requiring only 56 GFLOPs per inference, occupying just 9.4 MB of model size, and consuming 3.14 GB of memory during test time while sustaining real-time processing at 43 frames per second (FPS) for inputs of size  $448 \times 320$ .

In addition to the empirical comparison of computational requirements summarized in Table 2, a detailed quantitative analysis linking architectural design parameters to model complexity is provided in Supplementary Material Section 1. Specifically, this analysis formulates how the total number of learnable parameters and FLOPs scale with the internal channel dimensions of the Residual U-blocks, thereby offering explicit insight into the accuracy–efficiency trade-off underlying the MNSeg-Net design. This quantitative characterization clarifies that the lightweight behavior of MNSeg-Net is not achieved through post hoc pruning or compression, but rather through principled architectural parameterization that directly governs memory consumption, inference cost, and deployability. By making this relationship explicit, the supplementary analysis complements the present discussion and further substantiates the suitability of MNSeg-Net for real-time, resource-constrained clinical ultrasound applications.

## V.B. Real-Time Setup

In the clinical deployment of the proposed model, inference was performed on an external computational device, introducing a slight display delay of approximately 0.023 s per frame owing to the segmentation and CSA estimation. However, this delay poses no clinical limitations because the high frame redundancy during the transducer movement ensures temporal continuity during the scanning. The current model, which operates on a single-frame basis, allows clinicians to start tracking anywhere within the hand region, from the wrist to the elbow, thus offering the necessary flexibility for effective monitoring and assessment. From the deployment perspective, several factors were considered, including model size, number of parameters, computational complexity, memory consumption, and inference speed (see Table 2). Traditional UNet-based models, such as UNet (25.2M parameters, 72.6 GFLOPs), SegNet (29.4M, 87.74 GFLOPs), and ResUNet (39.37M, 111.63 GFLOPs), are computationally manageable and provide high frame rates (71, 66, and 38 FPS, respectively), making

806 them suitable for real-time deployment. However, their segmentation accuracy remains con-  
807 strained to the wrist region and does not generalize effectively to the entire forearm region.  
808 However, models such as UNet++ and BASNet demonstrated strong segmentation capa-  
809 bilities but imposed excessive computational costs. For instance, UNet++ requires 302.43  
810 GFLOPs and 4.98 GB of inference memory, whereas BASNet requires 278.91 GFLOPs and  
811 6.1 GB of memory, making it unsuitable for real-time deployment in resource-constrained  
812 systems. In contrast, both U2Net and MNSeg-Net achieved a favorable balance between  
813 performance and efficiency. Whereas U2Net requires 44.01M parameters, 82.53 GFLOPs,  
814 and 2.97 GB inference memory, MNSeg-Net achieves comparable accuracy with only 2.46M  
815 parameters, 56 GFLOPs, and a modest memory footprint of 3.14 GB. Additionally, MNSeg-  
816 Net achieved 43 FPS, slightly outperforming U2Net’s 40 FPS while maintaining superior  
817 segmentation accuracy with minimal resource usage. This establishes MNSeg-Net as an op-  
818 timal candidate for real-time clinical applications, offering an efficient performance without  
819 compromising accuracy.

## 820 **V.C. Model Deployment**

821 The real-time system was deployed on a laptop equipped with an NVIDIA GeForce GPU  
822 (6 GB memory), Intel i7 processor, and 16 GB RAM, connected to the ultrasound machine  
823 via an AV.IO HD Epiphan frame grabber. This configuration was sufficient to achieve real-  
824 time performance of up to 43 frames per second, offering both portability and practical  
825 convenience for clinical use. A key limitation of this setup is the restricted GPU memory,  
826 which supports models with inference demands of up to approximately 6 GB (Table 2). As a  
827 result, memory-intensive architectures such as BASNet, requiring more than 6.1 GB, could  
828 not be tested or deployed in this environment. This constraint underscores the importance  
829 of lightweight designs such as the proposed MNSeg-Net, which operates efficiently within  
830 a 3.14 GB memory footprint while maintaining high segmentation accuracy and real-time  
831 inference performance. Comparable mid-range GPUs or hospital-grade workstations would  
832 also be adequate for deployment in clinical settings, ensuring flexibility and scalability.

## 833 **V.D. Interpreting Favorable Outcomes: A Statistical and Clinical** 834 **Perspective**

835 The proposed MNSeg-Net model was comprehensively evaluated on two datasets, Dataset-1  
836 for experimental analysis and Dataset-2 for independent clinical validation, considering both  
837 statistical significance and clinical relevance to assess its practical applicability.

838 In Dataset-1, MNSeg-Net demonstrated statistically significant improvements in seg-  
839 mentation accuracy over multiple baseline models, as evidenced by a higher DSC, with sta-  
840 tistical significance confirmed through both paired t-tests and small-to-moderate effect sizes  
841 (Table 4) and Tukey’s HSD post hoc comparisons (Table 4). Beyond statistical significance,  
842 the extent to which the MNSeg-Net preserves segmentation accuracy, despite its lightweight  
843 design, was further quantified using formal equivalence testing. As reported in Table 6,

844 MNSeg-Net achieved statistically equivalent DSC performance in 3 out of 7 benchmark com-  
845 parisons, including equivalence with the full-size U2Net across all tested equivalence margins  
846 and with BASNet and SegNet at relaxed margins. Importantly, equivalence with U2Net is  
847 particularly meaningful, as U2Net represents a high-capacity yet computationally intensive  
848 architecture, thereby demonstrating that MNSeg-Net achieves comparable segmentation ac-  
849 curacy while operating with substantially reduced model complexity. Overall, MNSeg-Net  
850 shows consistent and statistically significant DSC improvements over most baseline models,  
851 while maintaining no significant difference and statistical equivalence with U2Net, suggesting  
852 an effective trade-off between segmentation accuracy and lightweight design. Additionally,  
853 the CSA estimates of MNSeg-Net exhibited no statistically significant difference from the  
854 clinician-annotated ground truth, as shown by the non-significant p-values and negligible  
855 effect size in Table 5, which is further supported by the non-significant findings in Tukey’s  
856 HSD test results presented in the same Table 5. This close agreement with expert clinical  
857 measurements underscores the clinical reliability of the model in estimating the CSA for  
858 CTS assessment. Beyond significance testing, formal equivalence analysis further quantified  
859 the extent of the clinical agreement. As shown in Table 7, MNSeg-Net achieved statistical  
860 equivalence with clinician-annotated CSA values, along with only three other models (U2Net,  
861 BASNet, and SegNet). Notably, equivalence with U2Net, a high-capacity and computationally  
862 intensive architecture, as well as with clinician-derived CSA measurements, highlights  
863 that the proposed lightweight design preserves clinically meaningful measurement fidelity,  
864 despite substantial reductions in model complexity. Together, the statistically significant  
865 improvements in DSC and the lack of clinically meaningful deviation in CSA estimation  
866 demonstrate the effectiveness of the MNSeg-Net under the evaluated experimental setting  
867 from both technical and clinical perspectives. To further examine inter-model equivalence  
868 between MNSeg-Net and U2Net, a dedicated TOST analysis was conducted. For the DSC  
869 metric, equivalence was established even under the strictest  $\pm 0.01$  margin (Table 8), indi-  
870 cating that segmentation differences fall well within clinically negligible bounds. For CSA,  
871 equivalence was achieved at clinically meaningful thresholds of  $\pm 0.3$  and  $\pm 0.5$  mm<sup>2</sup>, although  
872 not at the stricter  $\pm 0.1$  mm<sup>2</sup> margin (Table 9). This reflects only a minor systematic de-  
873 viation between the two models, with magnitude substantially below diagnostic relevance.  
874 Collectively, these results confirm that MNSeg-Net achieves performance statistically equiva-  
875 lent to U2Net while maintaining a substantially more computationally efficient architecture.

876 At the subject level, MNSeg-Net demonstrated consistent agreement in CSA estimation  
877 relative to clinician-annotated reference measurements across the wrist-to-elbow region, as  
878 quantified through subject-wise statistical analysis (Table 10). Paired t-tests and effect  
879 size analysis revealed that, although statistically significant differences were observed for  
880 some subjects, the associated effect sizes were predominantly negligible to small, indicating  
881 a limited practical impact. Specifically, five subjects exhibited negligible effect sizes and  
882 three subjects showed small effect sizes, while only two subjects demonstrated moderate  
883 effect sizes and no subject exhibited a large effect size, underscoring the overall clinical  
884 reliability of the proposed method. Beyond statistical significance testing, the extent to  
885 which MNSeg-Net preserves clinically meaningful CSA agreement at the individual-subject  
886 level was further quantified using formal equivalence testing. Using a clinically motivated  
887 equivalence margin of  $\pm 0.5$  mm<sup>2</sup>, the TOST procedure established statistical equivalence

888 between the MNSeg-Net predicted CSA values and clinician annotations for six out of ten  
889 subjects. Importantly, several subjects exhibited statistically significant paired differences  
890 while still satisfying the equivalence criterion, highlighting that statistical significance alone  
891 does not necessarily imply a clinically meaningful deviation. Overall, these subject-wise  
892 findings demonstrate that MNSeg-Net maintains clinically reliable CSA estimation for the  
893 majority of patients, while residual inter-subject variability represents a limited and localized  
894 deviation rather than a systematic limitation of the proposed lightweight architecture.

895 In Dataset-2, MNSeg-Net consistently achieved high segmentation accuracy across both  
896 clinical categories, including both the NORMAL and CTS-positive groups, yielding stable  
897 and reliable DSC values, as shown in Table 13. This consistent performance indicates the  
898 potential practical value of the MNSeg-Net in clinical applications, rather than a dominant  
899 advantage over existing methods. In the NORMAL group, the difference between the pre-  
900 dicted and clinician-annotated CSA measurements was minimal (mean difference =  $-0.206$   
901  $\text{mm}^2$ ,  $p = 0.208$ , Cohen’s  $d = -0.222$ ; Table 14), indicating no statistically or practically  
902 significant deviation. This agreement with clinician-annotated CSA values, the primary  
903 clinical reference standard, indicates that MNSeg-Net preserves diagnostic measurement fi-  
904 delity in non-pathological cases despite its lightweight design. In the CTS-positive group,  
905 although MNSeg-Net exhibited a slight underestimation of CSA (mean difference =  $+0.709$   
906  $\text{mm}^2$ ), this difference was statistically significant ( $p = 0.0094$ ), with a moderate effect size  
907 (Cohen’s  $d = 0.529$ ; Table 14). However, the predicted CSA values remained well above the  
908 diagnostic threshold of  $12 \text{ mm}^2$ , thereby maintaining their clinical relevance. Overall, the  
909 combination of MNSeg-Net’s consistent segmentation performance (Table 13) and clinically  
910 reliable CSA estimation across clinical categories (Table 14) supports its practical utility  
911 in assisting clinical interpretation and decision-making, within the evaluated experimental  
912 setting.

## 913 **V.E. Interpreting Performance Variability: A Statistical and Clin-** 914 **ical Perspective**

### 915 **V.E.1. DSC-Based Equivalence with Benchmark Models**

916 Across multiple experimental settings, several DSC-based pairwise comparisons yielded neg-  
917 ligible to small Cohen’s  $d$  values, indicating that MNSeg-Net achieves segmentation perfor-  
918 mance that is statistically comparable to several benchmark models such as U2Net rather  
919 than substantially exceeding them (Table 4). Formal equivalence testing of DSC scores (Ta-  
920 ble 6) further showed that MNSeg-Net achieved statistical equivalence with selected bench-  
921 marks, most notably the full-size U2Net across all tested equivalence margins. This outcome  
922 is expected in the context of median nerve segmentation, where the anatomical structure is  
923 relatively small and well-defined and many modern deep learning models already attain high  
924 DSC values. Under such conditions, further improvements correspond to marginal absolute  
925 gains, which naturally translate into small effect sizes when quantified using standardized  
926 metrics. Importantly, the primary objective of this study was not to maximize segmenta-  
927 tion accuracy at the expense of computational cost, but to design a clinically deployable

928 model that balances accuracy, efficiency, and robustness. While several benchmark models  
929 demonstrate similar segmentation accuracy, they typically require higher parameter counts,  
930 memory usage, and computational resources. In contrast, MNSEg-Net attains statistically  
931 equivalent performance with substantially reduced computational complexity, enabling real-  
932 time inference on resource-constrained clinical hardware. From a clinical perspective, such  
933 statistical equivalence combined with improved efficiency is often preferable to marginal accu-  
934 racy gains, as it facilitates reliable real-time operation, deployment, and seamless integration  
935 into routine ultrasound workflows without specialized hardware requirements.

## 936 V.E.2. Inter-Subject Variability

937 Despite the generally favorable performance of MNSEg-Net across datasets, minor unfavor-  
938 able outcomes were observed in terms of inter-subject variability. Subject-wise analysis  
939 (Section III.B.5. and Table 10) indicates that a small subset of subjects exhibiting pro-  
940 nounced anatomical variations or imaging inconsistencies occasionally showed slightly re-  
941 duced segmentation accuracy or CSA estimation performance compared with the model’s  
942 typical behavior. Consequently, paired statistical tests sometimes identified significant dif-  
943 ferences between model predictions and clinician annotations at the subject level. However,  
944 these differences were associated with negligible-to-small effect sizes, indicating that although  
945 the differences were statistically significant, their magnitude was small to influence clinical  
946 decision-making; this interpretation is further supported by the TOST equivalence analy-  
947 sis. Importantly, these subject-specific deviations were confined to a few cases and did not  
948 influence overall performance trends or compromise the clinical reliability of the proposed  
949 method. Rather than representing a fundamental limitation, these observations reflect in-  
950 herent inter-subject variability in clinical imaging data and highlight promising directions  
951 for future improvement, including increased training diversity, adaptive learning strategies,  
952 and enhanced incorporation of anatomical context.

## 953 V.F. Robustness Analysis on a Challenging Anatomical Region

954 The zero DSC values observed for several baseline models (UNet, SegNet, ResUNet,  
955 Attention-UNet, and UNet++) on frame-190 indicate a complete segmentation failure. Upon  
956 closer inspection, frame-190 corresponded to a challenging anatomical transition region in  
957 the proximal forearm, where the median nerve exhibited a rapid reduction in cross-sectional  
958 contrast, increased speckle noise, and partial boundary ambiguity. In this frame, the nerve  
959 appears faint and partially merges with the surrounding hypoechoic tissue, significantly de-  
960 grading the local edge and intensity cues. Baseline architectures primarily rely on local  
961 convolutional features and hierarchical downsampling, making them particularly sensitive  
962 to abrupt contrast loss and spatial ambiguity. Consequently, these models failed to activate  
963 meaningful segmentation responses, resulting in empty predictions and zero DSC values. In  
964 contrast, MNSEg-Net, BASNet, and U2Net maintained non-zero DSC scores for this frame  
965 owing to their ability to aggregate multi-scale contextual information. Specifically, MNSEg-  
966 Net employs a redesigned U2Net-inspired backbone combined with a MSFF module, which

967 enhances robustness by integrating features across multiple receptive fields and preserving  
968 the global context, enabling the network to infer the presence of the median nerve even when  
969 local appearance cues are severely degraded.

## 970 **V.G. Conditional Interpretation of CSA based on DSC**

971 It is important to note that segmentation-based measures such as DSC directly evaluate  
972 the spatial alignment between predicted and ground truth segmentations, making them  
973 the most reliable indicators for assessing statistical equivalence of model performance. In  
974 contrast, CSA is a derived scalar measure that may coincide with the ground truth area  
975 even when the segmentation is spatially misaligned. Consequently, CSA equivalence is only  
976 meaningful when supported by strong segmentation equivalence on DSC. Based on this  
977 principle, conclusions emphasize DSC equivalence as primary and interpret CSA findings  
978 conditionally on segmentation reliability.

## 979 **V.H. Limitations**

980 This study has certain limitations that merit consideration. First, the work focused ex-  
981 clusively on the median nerve using ultrasound imaging, and the findings may not readily  
982 generalize to other peripheral nerves or alternative imaging modalities. Second, segmen-  
983 tation performance may be affected by ultrasound artefacts, variations in imaging depth,  
984 and differences in probe pressure, which were not explicitly controlled in this study. Third,  
985 patient-specific anatomical variability, such as differences in nerve morphology across in-  
986 dividuals, was not systematically addressed and could introduce variability in model per-  
987 formance. These factors underscore the need for further validation and adaptation of the  
988 proposed framework in broader clinical contexts, including additional nerve types, diverse  
989 imaging conditions, and larger patient populations.

## 990 **V.I. Future Work**

### 991 **V.I.1. Technical Directions**

992 Future studies that can further enhance this work can include the following: First, seg-  
993 mentation accuracy can be improved by leveraging different models to reduce processing  
994 delays. Second, optimizing the setup in terms of portability and cost-effectiveness is a viable  
995 approach. Strategies for minimizing processing delays are also presented. Currently, the pro-  
996 posed model utilizes only 50% of the GPU during computation, with an inference memory  
997 usage of 3 GB out of 6 GB available. The implementation of a thread-programming-based  
998 technique can optimize resource utilization and reduce the processing delays. Additionally,  
999 the remaining 50% GPU capacity and 3 GB of memory can be used to run parallel infer-  
1000 ences, enhance speed, and minimize delays. The proposed model achieved an approximate

1001 segmentation accuracy 83% for the test data. Thus, exploring novel deep learning archi-  
1002 tectures and advanced training techniques can yield higher accuracy rates. Furthermore,  
1003 alternative computing platforms, such as Nvidia Jetson or mobile GPU-based computation  
1004 (with an Android application interface), can be considered for compact and cost-effective  
1005 clinical setups in the future.

## 1006 V.1.2. Clinical Directions

1007 The potential clinical applications of the proposed framework extend beyond CTS diagno-  
1008 sis. While the current study focused on the median nerve, the same model design could  
1009 be adapted to segment other peripheral nerves such as the radial and ulnar nerves, pro-  
1010 vided sufficient annotated data is available for training. This generalizability would make  
1011 the framework valuable in additional clinical scenarios, including the guidance of regional  
1012 anaesthesia procedures where accurate nerve identification is essential, as well as in emerg-  
1013 ing areas such as robotic-assisted surgery where real-time nerve localization could enhance  
1014 procedural safety.

1015 In regional anaesthesia, automated nerve segmentation could serve as a valuable educational  
1016 tool by providing real-time visual feedback on nerve boundaries, thereby improving the ac-  
1017 curacy and confidence of trainees during ultrasound-guided procedures. For robotic-assisted  
1018 surgery, integration of the proposed model into the surgical workflow could support nerve  
1019 identification and tracking, reducing the risk of nerve injury and enhancing procedural safety.

## 1020 VI. Conclusion

1021 This study proposed MNSeg-Net, a novel lightweight and efficient CNN model designed  
1022 for the real-time segmentation of the median nerve in ultrasound images. MNSeg-Net fea-  
1023 tures an additional subnetwork module, achieving competitive performance with minimal  
1024 parameters. The proposed model also demonstrated high segmentation accuracy from the  
1025 wrist to the elbow compared to more complex models, such as U2Net, while being com-  
1026 putationally and memory efficient. The proposed model also showed clinical utility and  
1027 provided real-time segmentation of the median nerve while scanning the patient, making  
1028 it suitable for real-time evaluation in CTS diagnosis. Furthermore, this work included the  
1029 development of an end-to-end clinical system with a GUI for real-time median nerve seg-  
1030 mentation, serving as a practical screening tool for clinicians. The integration of this into  
1031 the clinical setup highlights its practicality and feasibility. Overall, this study contributes  
1032 toward advancing automated real-time median nerve segmentation and CTS diagnosis. The  
1033 proposed model along with other models discussed in this study were made available here:  
1034 <https://github.com/venkateshvaddadi/Real-Time-Median-Nerve-Segmentation/>

## 1035 **Conflicts of interest**

1036 There are no conflicts of interest declared by Authors.

## 1037 **Acknowledgement**

1038 The authors are thankful to Aster-CMI Hospital, Bangalore for enabling this research work.  
1039 This work was supported by S. Ramachandran-National Bioscience Award for Career De-  
1040 velopment awarded by Department of Biotechnology, Govt. of India and in part by the  
1041 Wipro-GE Collaborative Laboratory on Artificial Intelligence in Health Care and Medical  
1042 Imaging.

## 1043 **Data availability**

1044 The imaging data used in this study are not publicly available due to participant pri-  
1045 vacy constraints and ethical restrictions. The code has been made available for enthusias-  
1046 tic users as open source at [https://github.com/venkateshvaddadi/Real-Time-Median-Nerve-](https://github.com/venkateshvaddadi/Real-Time-Median-Nerve-Segmentation/)  
1047 [Segmentation/](https://github.com/venkateshvaddadi/Real-Time-Median-Nerve-Segmentation/).

## 1048 **References**

- 1049
- 1050 <sup>1</sup> Soubeyrand M, Melhem R, Protais M, Artuso M, Crézé M. Anatomy of the median  
1051 nerve and its clinical applications *Hand Surgery and Rehabilitation*. 2020;39:2–18.
  - 1052 <sup>2</sup> Agarwal Pawan, Gupta Shivkant, Yadav Prashant, Sharma D. Cadaveric study of  
1053 anatomical variations of the median nerve and persistent median artery at wrist *In-*  
1054 *Indian journal of plastic surgery: official publication of the Association of Plastic Surgeons*  
1055 *of India*. 2014;47:95.
  - 1056 <sup>3</sup> Henry Brandon Michael, Zwinczewska Helena, Roy Joyeeta, et al. The prevalence of  
1057 anatomical variations of the median nerve in the carpal tunnel: a systematic review and  
1058 meta-analysis *PLoS One*. 2015;10:e0136477.
  - 1059 <sup>4</sup> Ehler Edvard. Median nerve ultrasonography in carpal tunnel syndrome *Clinical Neuro-*  
1060 *physiology Practice*. 2017;2:186.
  - 1061 <sup>5</sup> Avendaño-Coy Juan, Serrano-Muñoz Diego, Taylor Julian, Goicoechea-García Carlos,  
1062 Gómez-Soriano Julio. Peripheral nerve conduction block by high-frequency alternating  
1063 currents: a systematic review *IEEE Transactions on Neural Systems and Rehabilitation*  
1064 *Engineering*. 2018;26:1131–1140.

- 1065 <sup>6</sup> Hopkins PM. Ultrasound guidance as a gold standard in regional anaesthesia 2007.
- 1066 <sup>7</sup> Marhofer P, Greher M, Kapral S. Ultrasound guidance in regional anaesthesia *British*  
1067 *journal of anaesthesia*. 2005;94:7–17.
- 1068 <sup>8</sup> Hadjerci Oussama, Hafiane Adel, Conte Donatello, Makris Pascal, Vieyres Pierre, Delbos  
1069 Alain. Ultrasound median nerve localization by classification based on despeckle filter-  
1070 ing and feature selection in *2015 IEEE International Conference on Image Processing*  
1071 *(ICIP)*:4155–4159IEEE 2015.
- 1072 <sup>9</sup> Hadjerci Oussama, Hafiane Adel, Conte Donatello, Makris Pascal, Vieyres Pierre, Delbos  
1073 Alain. Computer-aided detection system for nerve identification using ultrasound images:  
1074 a comparative study *Informatics in Medicine Unlocked*. 2016;3:29–43.
- 1075 <sup>10</sup> Hafiane Adel, Vieyres Pierre, Delbos Alain. Deep learning with spatiotemporal consis-  
1076 tency for nerve segmentation in ultrasound images *arXiv preprint arXiv:1706.05870*.  
1077 2017.
- 1078 <sup>11</sup> Wang You-Wei, Chang Ruey-Feng, Horng Yi-Shiung, Chen Chii-Jen. MNT-DeepSL: Me-  
1079 dian nerve tracking from carpal tunnel ultrasound images with deep similarity learning  
1080 and analysis on continuous wrist motions *Computerized Medical Imaging and Graphics*.  
1081 2020;80:101687.
- 1082 <sup>12</sup> Huang Aiyue, Jiang Li, Zhang Jiangshan, Wang Qing. Attention-VGG16-UNet: a novel  
1083 deep learning approach for automatic segmentation of the median nerve in ultrasound  
1084 images *Quantitative imaging in medicine and surgery*. 2022;12:3138.
- 1085 <sup>13</sup> Festen Raymond T, Schrier Verena JMM, Amadio Peter C. Automated Segmentation of  
1086 the Median Nerve in the Carpal Tunnel using U-Net *Ultrasound in Medicine & Biology*.  
1087 2021;47:1964–1969.
- 1088 <sup>14</sup> Horng Ming-Huwi, Yang Cheng-Wei, Sun Yung-Nien, Yang Tai-Hua. Deepnerve: A new  
1089 convolutional neural network for the localization and segmentation of the median nerve  
1090 in ultrasound image sequences *Ultrasound in Medicine & Biology*. 2020;46:2439–2452.
- 1091 <sup>15</sup> Wu Chueh-Hung, Syu Wei-Ting, Lin Meng-Ting, et al. Automated Segmentation of  
1092 Median Nerve in Dynamic Sonography Using Deep Learning: Evaluation of Model Per-  
1093 formance *Diagnostics*. 2021;11:1893.
- 1094 <sup>16</sup> Di Cosmo Mariachiara, Fiorentino Maria Chiara, Villani Francesca Pia, et al. A deep  
1095 learning approach to median nerve evaluation in ultrasound images of carpal tunnel inlet  
1096 *Medical & Biological Engineering & Computing*. 2022;60:3255–3264.
- 1097 <sup>17</sup> Jiang Zhongliang, Li Xuesong, Chu Xiangyu, et al. Needle Segmentation Using GAN:  
1098 Restoring Thin Instrument Visibility in Robotic Ultrasound *IEEE Transactions on In-*  
1099 *strumentation and Measurement*. 2024;73:1-11.

- 1100 <sup>18</sup> Yeh Cheng-Liang, Wu Chueh-Hung, Hsiao Ming-Yen, Kuo Po-Ling. Real-time auto-  
1101 mated segmentation of median nerve in dynamic ultrasonography using deep learning  
1102 *Ultrasound in Medicine & Biology*. 2023;49:1129–1136.
- 1103 <sup>19</sup> Wang Jinting, Tang Yujiao, Xiao Yang, Zhou Joey Tianyi, Fang Zhiwen, Yang Feng.  
1104 GREnet: Gradually REcurrent Network With Curriculum Learning for 2-D Medical  
1105 Image Segmentation *IEEE Transactions on Neural Networks and Learning Systems*.  
1106 2024;35:10018-10032.
- 1107 <sup>20</sup> Roy Ayush, Pramanik Payel, Sarkar Ram. EU2-Net: A Parameter Efficient Ensem-  
1108 ble Model With Attention-Aided Triple Feature Fusion for Tumor Segmentation in  
1109 Breast Ultrasound Images *IEEE Transactions on Instrumentation and Measurement*.  
1110 2024;73:1-7.
- 1111 <sup>21</sup> Pan Lin, Cai Yanjing, Lin Ning, Yang Linxin, Zheng Shaohua, Huang Liqin. A two-stage  
1112 network with prior knowledge guidance for medullary thyroid carcinoma recognition in  
1113 ultrasound images *Medical Physics*. 2022;49:2413-2426.
- 1114 <sup>22</sup> Shi Fengxin, Zhu Dongming, Zhi Jia, Hou Guocun, Cui Yaoyao, Wang Xiacong. Au-  
1115 tomatic calculation method for stenosis ratio based on dialysis access ultrasound image  
1116 segmentation *Medical Physics*. ;n/a.
- 1117 <sup>23</sup> Gujarati Karan R., Bathala Lokesh, Venkatesh Vaddadi, Mathew Raji Susan, Yalavarthy  
1118 Phaneendra K.. Transformer-Based Automated Segmentation of the Median Nerve in  
1119 Ultrasound Videos of Wrist-to-Elbow Region *IEEE Transactions on Ultrasonics, Ferro-*  
1120 *electrics, and Frequency Control*. 2024;71:56-69.
- 1121 <sup>24</sup> Ronneberger Olaf, Fischer Philipp, Brox Thomas. U-net: Convolutional networks for  
1122 biomedical image segmentation in *International Conference on Medical image computing*  
1123 *and computer-assisted intervention*:234–241Springer 2015.
- 1124 <sup>25</sup> Badrinarayanan Vijay, Kendall Alex, Cipolla Roberto. Segnet: A deep convolutional  
1125 encoder-decoder architecture for image segmentation *IEEE transactions on pattern anal-*  
1126 *ysis and machine intelligence*. 2017;39:2481–2495.
- 1127 <sup>26</sup> Diakogiannis Foivos I, Waldner François, Caccetta Peter, Wu Chen. ResUNet-a: A deep  
1128 learning framework for semantic segmentation of remotely sensed data *ISPRS Journal*  
1129 *of Photogrammetry and Remote Sensing*. 2020;162:94–114.
- 1130 <sup>27</sup> Oktay Ozan, Schlemper Jo, Folgoc Loic Le, et al. Attention u-net: Learning where to  
1131 look for the pancreas *arXiv preprint arXiv:1804.03999*. 2018.
- 1132 <sup>28</sup> Zhou Zongwei, Rahman Siddiquee Md Mahfuzur, Tajbakhsh Nima, Liang Jianming.  
1133 Unet++: A nested u-net architecture for medical image segmentation in *Deep Learn-*  
1134 *ing in Medical Image Analysis and Multimodal Learning for Clinical Decision Support:*  
1135 *4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS*  
1136 *2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018,*  
1137 *Proceedings* 4:3–11Springer 2018.

- 1138 <sup>29</sup> Qin Xuebin, Zhang Zichen, Huang Chenyang, Gao Chao, Dehghan Masood, Jager-  
1139 sand Martin. Basnet: Boundary-aware salient object detection in *Proceedings of the*  
1140 *IEEE/CVF conference on computer vision and pattern recognition*:7479–7489 2019.
- 1141 <sup>30</sup> Qin Xuebin, Zhang Zichen, Huang Chenyang, Dehghan Masood, Zaiane Osmar R, Jager-  
1142 sand Martin. U2-Net: Going deeper with nested U-structure for salient object detection  
1143 *Pattern recognition*. 2020;106:107404.
- 1144 <sup>31</sup> Ferreira Tiago, Rasband Wayne. ImageJ user guide *Imagej/fiji*. 2012;1:155–161.
- 1145 <sup>32</sup> Ghasemi Majid, Masoumi Sanaz, Ansari Behnaz, Fereidan-Esfahani Mahboobeh,  
1146 Mousavi Seyed Morteza. Determination of cut-off point of cross-sectional area of median  
1147 nerve at the wrist for diagnosing carpal tunnel syndrome *Iranian journal of neurology*.  
1148 2017;16:164.
- 1149 <sup>33</sup> Nkrumah Gideon, Blackburn Alan R, Goitz Robert J, Fowler John R. Ultrasonography  
1150 findings in severe carpal tunnel syndrome *Hand*. 2020;15:64–68.
- 1151 <sup>34</sup> Falsetti Paolo, Conticini Edoardo, Baldi Caterina, et al. A novel ultrasonographic  
1152 anthropometric-independent measurement of median nerve swelling in carpal tunnel syn-  
1153 drome: the “nerve/tendon ratio”(NTR) *Diagnostics*. 2022;12:2621.
- 1154 <sup>35</sup> Zou Kelly H, Warfield Simon K, Bharatha Aditya, et al. Statistical validation of image  
1155 segmentation quality based on a spatial overlap index1: scientific reports *Academic*  
1156 *radiology*. 2004;11:178–189.
- 1157 <sup>36</sup> Huttenlocher Daniel P, Klanderman Gregory A., Rucklidge William J. Comparing im-  
1158 ages using the Hausdorff distance *IEEE Transactions on pattern analysis and machine*  
1159 *intelligence*. 1993;15:850–863.
- 1160 <sup>37</sup> Kingma Diederik P, Ba Jimmy. Adam: A method for stochastic optimization *arXiv*  
1161 *preprint arXiv:1412.6980*. 2014.
- 1162 <sup>38</sup> Sudre Carole H, Li Wenqi, Vercauteren Tom, Ourselin Sebastien, Jorge Cardoso M.  
1163 Generalised dice overlap as a deep learning loss function for highly unbalanced segmen-  
1164 tations in *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical*  
1165 *Decision Support: Third International Workshop, DLMIA 2017, and 7th International*  
1166 *Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC,*  
1167 *Canada, September 14, Proceedings 3*:240–248Springer 2017.
- 1168 <sup>39</sup> De Boer Pieter-Tjerk, Kroese Dirk P, Mannor Shie, Rubinstein Reuven Y. A tutorial on  
1169 the cross-entropy method *Annals of operations research*. 2005;134:19–67.
- 1170 <sup>40</sup> Jadon Shruti. A survey of loss functions for semantic segmentation in *2020 IEEE*  
1171 *conference on computational intelligence in bioinformatics and computational biology*  
1172 *(CIBCB)*:1–7IEEE 2020.

- 1173 41 Wang Zhou, Bovik Alan C, Sheikh Hamid R, Simoncelli Eero P. Image quality assess-  
1174 ment: from error visibility to structural similarity *IEEE transactions on image process-*  
1175 *ing.* 2004;13:600–612.
- 1176 42 AV.Io HD Now Available: The New and Improved AV.Io HD+ Website.
- 1177 43 Van Woudenberg Nathan, Liao Zhibin, Abdi Amir H, et al. Quantitative echocardi-  
1178 ography: real-time quality estimation and view classification implemented on a mobile  
1179 android device in *Simulation, Image Processing, and Ultrasound Systems for Assisted*  
1180 *Diagnosis and Navigation: International Workshops, POCUS 2018, BIVPCS 2018, Cu-*  
1181 *RIOUS 2018, and CPM 2018, Held in Conjunction with MICCAI 2018, Granada, Spain,*  
1182 *September 16–20, 2018, Proceedings:*74–81Springer 2018.
- 1183 44 Smerilli Gianluca, Cipolletta Edoardo, Sartini Gianmarco, et al. Development of a convo-  
1184 lutional neural network for the identification and the measurement of the median nerve  
1185 on ultrasound images acquired at carpal tunnel level *Arthritis Research & Therapy.*  
1186 2022;24:38.
- 1187 45 Student . The probable error of a mean *Biometrika.* 1908;6:1–25.
- 1188 46 Cohen Jacob. *Statistical Power Analysis for the Behavioral Sciences.* Lawrence Erlbaum  
1189 Associates2nd ed. 1988.
- 1190 47 Fisher Ronald A.. *Statistical Methods for Research Workers.* Oliver and Boyd 1925.
- 1191 48 Tukey John W.. Comparing individual means in the analysis of variance *Biometrics.*  
1192 1949;5:99–114.
- 1193 49 Schuirmann Donald J. A comparison of the two one-sided tests procedure and the power  
1194 approach for assessing the equivalence of average bioavailability *Journal of pharmacoki-*  
1195 *netics and biopharmaceutics.* 1987;15:657–680.
- 1196 50 Lu Haoran, She Yifei, Tie Jun, Xu Shengzhou. Half-UNet: A simplified U-Net architec-  
1197 ture for medical image segmentation *Frontiers in Neuroinformatics.* 2022;16:911679.

1198 **List of Tables**

1199 1 Comparison between the proposed MNSeg-Net and the full-size U2Net model. 11

1200 2 Deep learning models Specifications utilised in this work and their correspond-  
1201 ing inference time . . . . . 12

1202 3 Performance comparison of deep learning models evaluated on a 10-subject  
1203 test dataset across two anatomical regions: the wrist and wrist-to-elbow. Eval-  
1204 uation metrics include DSC, Precision (Prec), Recall (Rec), and Hausdorff  
1205 Distance (HD).  $\uparrow$  indicates that higher values correspond to better perfor-  
1206 mance, while  $\downarrow$  indicates that lower values correspond to better performance. 18

1207 4 Comprehensive statistical comparison of segmentation models relative to the  
1208 proposed MNSeg-Net (reference baseline) based on DSC. Mean differences and  
1209 unadjusted p-values were computed using paired t-tests, and Cohen’s d values  
1210 represent practical significance. Effect size interpretation: N = Negligible  
1211 ( $d < 0.2$ ), S = Small ( $0.2 \leq d < 0.5$ ), M = Moderate ( $0.5 \leq d < 0.8$ ), and L  
1212 = Large ( $d \geq 0.8$ ). An asterisk (\*) denotes a statistically significant difference  
1213 from MNSeg-Net at  $p < 0.05$  after Tukey correction. . . . . 19

1214 5 Comprehensive statistical comparison of model-predicted CSA values relative  
1215 to clinician-annotated  $CSA_{Act}$  (reference baseline). Mean differences and un-  
1216 adjusted p-values were computed using paired t-tests, and Cohen’s d values  
1217 represent practical significance. Effect size interpretation: N = Negligible  
1218 ( $d < 0.2$ ), S = Small ( $0.2 \leq d < 0.5$ ), M = Moderate ( $0.5 \leq d < 0.8$ ), and L  
1219 = Large ( $d \geq 0.8$ ). An asterisk (\*) denotes a statistically significant difference  
1220 from  $CSA_{Act}$  at  $p < 0.05$  after Tukey correction. . . . . 19

1221 6 Two One-Sided Tests (TOST) equivalence test results for DSC scores of  
1222 MNSeg-Net versus other models. The table reports mean differences with  
1223 95% confidence intervals and equivalence decisions at margins  $\pm 0.01$ ,  $\pm 0.03$ ,  
1224 and  $\pm 0.05$ . MNSeg-Net was statistically equivalent to U2Net at all tested  
1225 margins, BASNet at  $\pm 0.03$  and  $\pm 0.05$ , and SegNet at  $\pm 0.05$  only. No equiv-  
1226 alence was found with UNet, ResUNet, Attention-UNet, or UNet++. . . . . 20

1227 7 TOST equivalence test results for CSA (actual vs. computed) at margin  
1228  $\pm 0.5$ . The proposed MNSeg-Net, SegNet, BASNet, and U2Net models met the  
1229 equivalence criterion, whereas UNet, ResUNet, Attention-UNet, and UNet++  
1230 did not. At the equivalence margin of  $\pm 0.5$ , MNSeg-Net, SegNet, BASNet,  
1231 and U2Net demonstrated statistical equivalence with the clinician-annotated  
1232 CSA values. In contrast, UNet, ResUNet, Attention-UNet, and UNet++ did  
1233 not achieve statistical equivalence under this criterion. . . . . 21

1234 8 TOST equivalence analysis of frame-level DSC between MNSeg-Net and  
1235 U2Net. It reports the mean difference, 95% confidence interval, and equiva-  
1236 lence decisions under predefined DSC margins. . . . . 21

1237	9	TOST equivalence analysis of frame-level CSA estimates between MNSeg-Net and U2Net. It reports the mean difference, 95% confidence interval, and equivalence decisions under clinically defined CSA margins (in mm <sup>2</sup> ). . . . .	21
1238			
1239			
1240	10	Subject-wise evaluation of the proposed MNSeg-Net model from wrist to elbow, including segmentation metrics (DSC, precision, recall, HD) and a statistical comparison between the predicted CSA values from MNSeg-Net (CSA <sub>Cal</sub> ) and clinician-annotated CSA values (CSA <sub>Act</sub> ). Paired t-tests, Cohen’s <i>d</i> effect sizes, mean differences, and confidence intervals were computed across 300 image frames per subject to assess the clinical reliability of the MNSeg-Net. Equivalence was assessed using the TOST procedure with predefined margins of ±0.5 mm <sup>2</sup> . Asterisks (*) denote statistically significant differences ( $p < 0.05$ ). Effect size categories: N = Negligible ( $d < 0.2$ ), S = Small ( $0.2 \leq d < 0.5$ ), M = Moderate ( $0.5 \leq d < 0.8$ ), L = Large ( $d \geq 0.8$ ). . . . .	22
1241			
1242			
1243			
1244			
1245			
1246			
1247			
1248			
1249			
1250	11	Performance comparison of different models on 25% and 50% training data. .	25
1251	12	Performance of different models under speckle noise perturbation at mean levels 0.1, 0.2, and 0.3, evaluated at the wrist and wrist-to-elbow regions. . .	26
1252			
1253	13	Patient-wise segmentation performance metrics, including DSC, precision, recall, and CSA estimations, comparing clinician-annotated (CSA <sub>Act</sub> ) and MNSeg-Net-predicted (CSA <sub>Cal</sub> ) values for 30 patients, categorized into NORMAL and CTS-positive groups. ↑ indicates that higher values correspond to better performance, while ↓ indicates that lower values correspond to better performance. . . . .	28
1254			
1255			
1256			
1257			
1258			
1259	14	Statistical analysis on clinical dataset for comparison of MNSeg-Net predicted CSA values from clinician-annotated CSA <sub>Act</sub> . The table includes mean differences, p-values from paired t-tests and Cohen’s <i>d</i> effect sizes. (*) indicates a statistically significant difference from CSA <sub>Act</sub> at $p < 0.05$ . Effect size abbreviations: N = Negligible ( $d < 0.2$ ), S = Small ( $0.2 \leq d < 0.5$ ), M = Moderate ( $0.5 \leq d < 0.8$ ), L = Large ( $d \geq 0.8$ ). . . . .	28
1260			
1261			
1262			
1263			
1264			
1265	15	Ablation study results with the data from wrist to elbow region. . . . .	29

1266 **List of Figures**

1267	1	The physical system setup for Real-Time Median Nerve Segmentation includes the following main parts: (c). Philips cx50 ultrasound (US) machine(Input source) (d). Av.io Frame Grabber (e). DVI-to-HDMI interface (f).USB-B-to-USB-A interface (g). Switch box (power supply) (h). Portable wheeled table for the movable setup. The frame grabber was connected to the DVI output of the US machine. It is then connected via USB-B to a laptop, which serves as the computational device for executing the deep learning model. . . . .	4
1268			
1269			
1270			
1271			
1272			
1273			

1274	2	The proposed MNSeg-Net architecture for Median Nerve Segmentation with deep supervision and sub-network module for multi-scale feature fusion. (a) Main-network as prediction module, (b) Sub-network for multi-scale feature fusion. Detailed UNet block configurations used in MNSeg-Net were provided in Fig. 1 of supplementary information. . . . .	9
1275			
1276			
1277			
1278			
1279	3	Example segmentation of the median nerve using the methods discussed in this work for subject-1. The green contour indicates expert annotation, and the red contour indicates the result obtained for the corresponding method, as indicated in each row. The associated frame number is given at the top of every image (0 corresponds to the start of the wrist region, and 299 corresponds to the elbow region), and the bottom of each frame has the corresponding computed cross-sectional area (CSA) from the expert and model in green and red, respectively. . . . .	10
1280			
1281			
1282			
1283			
1284			
1285			
1286			
1287	4	Qualitative visualization of the learned feature representations for patient 1. The features extracted from the final encoder layer of each model were projected into two dimensions using PCA. Each point corresponds to a pixel classified as either background (blue) or median nerve (green). Compared with baseline models, MNSeg-Net exhibits clearer separation between the two classes, indicating more discriminative feature learning. . . . .	14
1288			
1289			
1290			
1291			
1292			
1293	5	(a), (b), and (c) represent sample images captured during real-time testing of wrist, forearm, and elbow scanning, shown alongside their corresponding median nerve segmentations and CSA in US frames. The CSA values are presented as $cm^2$ . In addition, (d) illustrates a sample image in which the CSA surpassed the threshold of $0.12 cm^2$ during wrist scanning. The CSA values are displayed in the top right corner of the ultrasound frame. . . . .	15
1294			
1295			
1296			
1297			
1298			
1299	6	Timeline diagram of the clinical workflow. The US transducer generates raw frames that are displayed on the ultrasound machine while simultaneously being captured via a frame grabber. The raw frames are sent to the GPU, where the proposed MNSeg-Net model processes them to generate segmented frames with CSA computation. These are displayed on an external screen, providing real-time quantitative and visual feedback alongside the ultrasound machine's native display. . . . .	15
1300			
1301			
1302			
1303			
1304			
1305			