

---

# Zero-shot evaluation of promptable foundation models for 3D CT segmentation

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

1 Foundation Models (FMs) have revolutionized interactive segmentation for medical  
2 imaging. However, the increasing number of promptable FMs, along with evalua-  
3 tions varying in dataset, metrics, and compared models, makes direct comparison  
4 difficult and complicates the selection of the most suitable model for specific  
5 clinical tasks. In the context of bone segmentation in CT scans, we evaluated 11  
6 promptable FMs using non-iterative 2D and 3D prompting strategies on both a  
7 private and public dataset. The models were categorized based on their prediction  
8 dimensionality (2D vs. 3D), and the Pareto-optimal models were identified.

9 **1 Introduction**

10 Promptable FMs can serve different purposes, from accelerating annotations (1) to being fine-tuned  
11 and adapted for specific datasets (2). Selecting the most suitable model depends on the task at hand,  
12 the available resources, and the required accuracy. With the growing number of models, evaluations  
13 are typically conducted in isolation, making direct comparison challenging. Although broader  
14 benchmarks have been proposed (3; 4) to demonstrate generalization across diverse datasets, clinical  
15 applications often demand solutions for specific tasks. In this work, we focus on bone segmentation  
16 in CT scans and evaluate promptable FMs, comparing them in the Pareto sense to avoid prioritizing a  
17 single metric while also incorporating computational resources.

18 **2 Method**

19 **2.1 Promptable Foundation Models**

20 The introduction of the Segment Anything Model (**SAM**) (5) made Foundation Models (FMs) popular  
21 for interactive segmentation. Building on SAM, SAM2 introduced a memory mechanism allowing  
22 video segmentation by means of prediction propagation from initial frames to whole videos. While  
23 developed for natural images and videos, initial evaluation studies showed promising results but also  
24 evident limitations for medical image segmentation, which motivated the development of dedicated  
25 medical FMs. Med-SAM (6), SAM-Med2D (7), ScribblePrompt-SAM (8) and MedicoSAM (9) are  
26 SAM-based models fine-tuned for medical data. MedicoSAM (9) also uses prompt propagation for  
27 volumetric predictions. SAM-Med3D (10) extends SAM to a 3D architecture trained from scratch.  
28 Other 3D interactive segmentation models include SegVol (11), incorporating semantic prompts  
29 alongside geometric prompts, Vista3D (12), supporting automatic closed- and interactive open-set  
30 segmentation, and nnInteractive (4), offering a well-integrated framework with extended prompting  
31 strategies. Med-SAM2 (1) is a SAM2-based model fine-tuned on medical data. All mentioned  
32 models have in common that they are promptable with sparse prompts, i.e., the object of interest  
33 is identified by a bounding box or point/click, allowing interactive segmentation. Based on the  
34 prediction dimensionality (2D vs. 3D) and evaluation manner (slice-wise or volumetric), promptable  
35 FMs can be categorized into two groups (Table 1):

Table 1: Overview of promptable FMs with model backbone architecture and the supported prompting strategies. Boxed settings are possible across different models resulting in our default settings. (✓) denotes that authors explicitly stated that the test set of (15) was excluded from training. N denotes that multiple prompts (i.e., up to N prompts) per slice were used.

Model	Arch.	Med.FT	(15)	Box (1/N)	Point (1/N)	Pt+Box (1/N)	Slice (1/N)	Vol. Limits
SAM (5)	ViT	x	x	✓/✓	✓/✓	✓/✓	-	-
SAM 2D (13)	Hiera	x	x	✓/✓	✓/✓	✓/✓	-	-
Med-SAM (6)	SAM	✓	✓	✓/✓	X/X	X/X	-	-
SAM-Med2D (7)	SAM	✓	✓	✓/✓	✓/✓	✓/✓	-	-
ScribblePrompt-U (8)	UNet	✓	(✓)	✓/✓	✓/✓	✓/✓	-	-
ScribblePrompt-SAM (8)	SAM	✓	(✓)	✓/✓	✓/✓	✓/X	-	-
MedicoSAM 2D (9)	SAM	✓	✓	✓/✓	✓/✓	✓/✓	-	-
SAM2 3D (13)	Hiera	x	x	✓/X	✓/✓	✓/X	✓/✓	✓
SAM-Med3D (10)	3D ViT	✓	(✓)	X/X	✓/✓	X/X	✓/X	x
SegVol (11)	3D ViT	✓	✓	X/X	✓/✓	X/X	✓/✓	x
MedicoSAM 3D (9)	SAM	✓	✓	✓/X	✓/✓	✓/X	✓/X	x
Vista3D (12)	SegResNet	✓	✓	X/X	✓/✓	X/X	✓/✓	x
nnInteractive (4)	CNN	✓	✓	✓/✓	✓/✓	✓/X	✓/✓	x
Med-SAM2 (1)	SAM2	✓	✓	✓/X	X/X	X/X	✓/✓	✓

36 **2D models evaluated slice-wise:** SAM Vit-B, SAM Vit-H, SAM Vit-L (5); SAM2.1 B+, SAM2.1  
37 L, SAM2.1 S, SAM2.1 T for 2D image segmentation (13); Med-SAM (6), SAM-Med2D (7),  
38 ScribblePrompt-UNet and ScribblePrompt-SAM (8), MedicoSAM2D (9)

39 **3D models evaluated slice-wise and 3D models evaluated volumetric:** SAM2.1 B+, SAM2.1  
40 L, SAM2.1 S, SAM2.1 T for 3D volume segmentation (13); Med-SAM2 (1), nnInteractive (4),  
41 MedicoSAM3D (9), SegVol (11), Vista3D (12).

## 42 2.2 Prompting strategies

43 For this study, we used non-iterative prompts, automatically extracted from the reference masks.  
44 The **2D prompting strategies** included bounding boxes, center points, or a combination of both,  
45 applied to up to five components. This choice was motivated by their effectiveness in 2D SAM-family  
46 models (14) and by the dataset characteristics, where most objects consist of up to five disconnected  
47 components. The **3D prompting strategies** extended the 2D approaches by incorporating the initial  
48 slice as a third dimension. All FMs (except SegVol (11)) rely on pseudo-3D boxes represented by 2D  
49 coordinates and slice indices, which can be similarly applied to a 3D point. The chosen settings were  
50 bounding boxes, center points or their combination extracted from the largest component in a single  
51 initial slice, since this is the minimal feasible configuration compatible with all models (Table 1).

## 52 2.3 Dataset

53 As medical FMs are usually trained with publicly available datasets (15; 16; 17), to ensure a fair  
54 comparison across models, we had to resort to a private CT test dataset, approved by the local Medical  
55 Ethics Committee. To allow reproducibility, we included 10 samples of the TotalSegmentator test set  
56 (15) to our study. The final test dataset contains four skeletal regions: (D1) 10 bilateral shoulder CT  
57 scans with labels for scapula and humerus; (D2) 10 unilateral wrist CT scans with labels for capitate,  
58 lunate, radius, scaphoid, triquetrum, and ulna; (D3) 15 unilateral lower leg CT scans with labels for  
59 tibia and fibular implant; (D4) 10 unilateral hip CT scans (15) with the original labels hip and femur,  
60 and manually added labels for femur implant.

61 To reduce computational resources, we selected random slices as *initial slices* based on the following  
62 strategy: From D1 and D2, two slices per class were extracted; from D4, two slices per original class;  
63 and from D3, three slices per class. The top and bottom 10% of each volume were excluded to avoid  
64 slices with little relevant anatomy. To ensure spatial diversity, a minimum spacing between selected  
65 slices was enforced: > 10 slices for D1, D3, and D4, and > 5 slices for D2. In total, 370 slices were  
66 extracted (D1: 80, D2: 120, D3: 90, D4: 80).

## 67 2.4 Evaluation

68 **Segmentation Performance** was measured by Dice Similarity Coefficient (DSC), 95%-percentile  
69 Hausdorff distance (HD95), and Normalized Surface Dice (NSD) ( $\tau = 1.5$  mm, based on largest

Table 2: Segmentation performance and model size for Pareto-optimal models in each category. Models highlighted in bold are Pareto-optimal and have the least model parameters.

Model	Size (M)	DSC ( $\uparrow$ ) (%)	NSD ( $\uparrow$ ) (%)	HD95 ( $\downarrow$ ) (mm)	DSC ( $\uparrow$ ) (%)	NSD ( $\uparrow$ ) (%)	HD95 ( $\downarrow$ ) (mm)	DSC ( $\uparrow$ ) (%)	NSD ( $\uparrow$ ) (%)	HD95 ( $\downarrow$ ) (mm)
2D Models										
<i>Prompt: Bounding box</i>										
<b>Med-SAM2</b>	39	-	-	-	-	-	-	<b>79.56</b> $\pm 11.1$	<b>80.25</b> $\pm 10.5$	<b>13.49</b> $\pm 11.1$
MedicoSAM2D	94	90.74 $\pm 7.7$	97.36 $\pm 3.6$	0.76 $\pm 0.9$	-	-	-	-	-	-
<b>SAM2.1 B+</b>	81	<b>90.60</b> $\pm 8.1$	<b>97.84</b> $\pm 3.5$	<b>0.82</b> $\pm 1.0$	-	-	-	-	-	-
<b>nnInteractive</b>	102	-	-	-	<b>90.80</b> $\pm 9.2$	<b>97.57</b> $\pm 4.6$	<b>1.05</b> $\pm 2.2$	-	-	-
<i>Prompt: Center point</i>										
<b>SAM B</b>	94	<b>85.43</b> $\pm 14.4$	<b>90.82</b> $\pm 13.0$	<b>4.83</b> $\pm 6.3$	-	-	-	-	-	-
<b>nnInteractive</b>	102	-	-	-	-	<b>83.47</b> $\pm 12.9$	<b>90.92</b> $\pm 9.6$	<b>2.16</b> $\pm 2.8$	<b>69.40</b> $\pm 11.2$	<b>68.23</b> $\pm 12.0$
<i>Prompt: Combination of bounding box and center point</i>										
MedicoSAM2D	94	91.27 $\pm 7.4$	97.74 $\pm 3.3$	0.69 $\pm 0.8$	-	-	-	-	-	-
<b>SAM2.1 B+</b>	81	91.98 $\pm 7.2$	98.21 $\pm 3.6$	0.73 $\pm 1.1$	-	-	-	<b>68.33</b> $\pm 9.4$	<b>67.86</b> $\pm 10.2$	<b>26.04</b> $\pm 18.2$
SAM2.1 L	224	90.90 $\pm 6.9$	98.36 $\pm 3.2$	0.69 $\pm 1.0$	-	-	-	-	-	-
SAM2.1 S	46	91.51 $\pm 7.0$	98.40 $\pm 3.3$	0.69 $\pm 0.9$	-	-	-	-	-	-
<b>SAM2.1 T</b>	39	<b>91.83</b> $\pm 6.9$	<b>98.38</b> $\pm 3.2$	<b>0.71</b> $\pm 1.0$	<b>89.27</b> $\pm 8.6$	<b>96.95</b> $\pm 5.5$	<b>1.95</b> $\pm 3.4$	-	-	-
nnInteractive	102	-	-	-	89.57 $\pm 8.6$	96.34 $\pm 4.8$	1.10 $\pm 1.6$	75.92 $\pm 9.4$	76.60 $\pm 9.6$	26.53 $\pm 10.3$

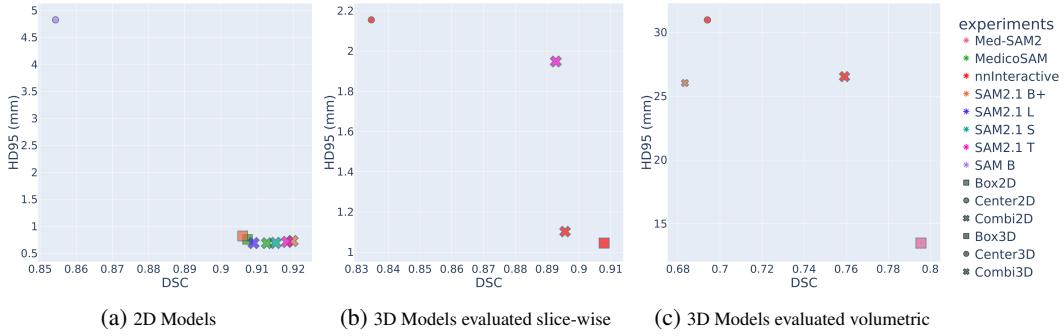


Figure 1: Segmentation performance of Pareto-front models (DSC vs. HD95 (mm)). Symbol size encodes NSD, mapped linearly to values between 2 (smallest) to 5 (largest).

spacing), following MetricsReloaded (18) and employing DisTorch (19). 2D and 3D models are evaluated on the selected slices (slice-wise), 3D models additionally on full volumes (volumetric). The **Pareto front** consists of all models that are not simultaneously outperformed across all criteria, i.e., no other model performs at least as well on every criterion and strictly better on at least one. For each of the three categories per prompt strategy, the best models in the Pareto sense are identified based on segmentation metrics (i.e., DSC, HD95, NSD). For multiple Pareto-optimal models, the model size (i.e., number of parameters) is also taken into consideration, jointly reflecting on performance and computational efficiency.

### 3 Results & Discussion

The Pareto-optimal models based on segmentation metrics for each category are illustrated in a DSC vs. HD95 (mm) scatterplot (Figure 1) and summarized in Table 2, with models that are Pareto-optimal and have the lowest model size highlighted in bold. The two main insights are: The bounding box prompt performs well in 2D and 3D, with improved results in 3D without combining with center points; Within the pool of models, medical dedicated FMs, such as MedicoSAM2D, Med-SAM2, and nnInteractive can perform on-par or outperform SAM and SAM2.1.

### 4 Conclusion

We identified the Pareto-optimal models based on segmentation metrics out of 11 promptable FMs using non-iterative 2D and 3D prompts. Although the *Pareto front* depends on the chosen criteria and different metrics may yield different selections, it avoids prioritizing a single metric. For future work, we are aiming to evaluate the prompt robustness of the selected models.

90 **References**

- 91 [1] J. Ma, Z. Yang, S. Kim, B. Chen, M. Baharoon, A. Fallahpour, R. Asakereh, H. Lyu, B. Wang,  
92 Medsam2: Segment anything in 3d medical images and videos (2025). [arXiv:2504.03600](https://arxiv.org/abs/2504.03600).  
93 URL <https://arxiv.org/abs/2504.03600>
- 94 [2] H. Gu, H. Dong, J. Yang, M. A. Mazurowski, How to build the best medical image segmen-  
95 tation algorithm using foundation models: a comprehensive empirical study with segment  
96 anything model (2025). [arXiv:2404.09957](https://doi.org/10.59275/j.melba.2025-86a6), doi:<https://doi.org/10.59275/j.melba.2025-86a6>.  
97 URL <https://arxiv.org/abs/2404.09957>
- 99 [3] C. Ulrich, T. Wald, E. Tempus, M. Rokuss, P. F. Jaeger, K. Maier-Hein, Radioactive: 3d  
100 radiological interactive segmentation benchmark (2025). [arXiv:2411.07885](https://arxiv.org/abs/2411.07885).  
101 URL <https://arxiv.org/abs/2411.07885>
- 102 [4] F. Isensee, M. Rokuss, L. Krämer, S. Dinkelacker, A. Ravindran, F. Stritzke, B. Hamm, T. Wald,  
103 M. Langenberg, C. Ulrich, J. Deissler, R. Floca, K. Maier-Hein, nninteractive: Redefining 3d  
104 promptable segmentation (2025). [arXiv:2503.08373](https://arxiv.org/abs/2503.08373).  
105 URL <https://arxiv.org/abs/2503.08373>
- 106 [5] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C.  
107 Berg, W.-Y. Lo, P. Dollár, R. Girshick, Segment anything (2023). [arXiv:2304.02643](https://arxiv.org/abs/2304.02643).  
108 URL <https://arxiv.org/abs/2304.02643>
- 109 [6] J. Ma, Y. He, F. Li, L. Han, C. You, B. Wang, Segment anything in medical images, *Nature*  
110 Communications 15 (2024) 1–9.
- 111 [7] J. Cheng, J. Ye, Z. Deng, J. Chen, T. Li, H. Wang, Y. Su, Z. Huang, J. Chen, L. J. H. Sun, J. He,  
112 S. Zhang, M. Zhu, Y. Qiao, Sam-med2d (2023). [arXiv:2308.16184](https://arxiv.org/abs/2308.16184).
- 113 [8] H. E. Wong, M. Rakic, J. Guttag, A. V. Dalca, Scribbleprompt: Fast and flexible interactive  
114 segmentation for any biomedical image, European Conference on Computer Vision (ECCV)  
115 (2024).
- 116 [9] A. Archit, L. Freckmann, C. Pape, Medicosam: Towards foundation models for medical image  
117 segmentation (2025). [arXiv:2501.11734](https://arxiv.org/abs/2501.11734).  
118 URL <https://arxiv.org/abs/2501.11734>
- 119 [10] H. Wang, S. Guo, J. Ye, Z. Deng, J. Cheng, T. Li, J. Chen, Y. Su, Z. Huang, Y. Shen, B. Fu,  
120 S. Zhang, J. He, Y. Qiao, Sam-med3d: Towards general-purpose segmentation models for  
121 volumetric medical images (2024). [arXiv:2310.15161](https://arxiv.org/abs/2310.15161).  
122 URL <https://arxiv.org/abs/2310.15161>
- 123 [11] Y. Du, F. Bai, T. Huang, B. Zhao, Segvol: Universal and interactive volumetric medical image  
124 segmentation (2025). [arXiv:2311.13385](https://arxiv.org/abs/2311.13385).  
125 URL <https://arxiv.org/abs/2311.13385>
- 126 [12] Y. He, P. Guo, Y. Tang, A. Myronenko, V. Nath, Z. Xu, D. Yang, C. Zhao, B. Simon, M. Belue,  
127 S. Harmon, B. Turkbey, D. Xu, W. Li, Vista3d: A unified segmentation foundation model for 3d  
128 medical imaging (2024). [arXiv:2406.05285](https://arxiv.org/abs/2406.05285).  
129 URL <https://arxiv.org/abs/2406.05285>
- 130 [13] N. Ravi, V. Gabeur, Y.-T. Hu, R. Hu, C. Ryali, T. Ma, H. Khedr, R. Rädle, C. Rolland,  
131 L. Gustafson, E. Mintun, J. Pan, K. V. Alwala, N. Carion, C.-Y. Wu, R. Girshick, P. Dollár, C. Fe-  
132 ichtenhofer, Sam 2: Segment anything in images and videos, arXiv preprint [arXiv:2408.00714](https://arxiv.org/abs/2408.00714)  
133 (2024).  
134 URL <https://arxiv.org/abs/2408.00714>
- 135 [14] C. Magg, C. I. Sánchez, H. Kervadec, Zero-shot capability of 2d SAM-family models for bone  
136 segmentation in CT scans, in: Medical Imaging with Deep Learning, 2025.  
137 URL <https://openreview.net/forum?id=AUv6NhK9aH>

- 138 [15] J. Wasserthal, H.-C. Breit, M. T. Meyer, M. Pradella, D. Hinck, A. W. Sauter, T. Heye, D. T.  
139 Boll, J. Cyriac, S. Yang, M. Bach, M. Segeroth, Totalsegmentator: Robust segmentation of  
140 104 anatomic structures in ct images, Radiology: Artificial Intelligence 5 (5) (2023) e230024.  
141 doi:10.1148/ryai.230024.  
142 URL <https://doi.org/10.1148/ryai.230024>
- 143 [16] P. Liu, H. Han, Y. Du, H. Zhu, Y. Li, F. Gu, H. Xiao, J. Li, C. Zhao, L. Xiao, X. Wu, S. K. Zhou,  
144 Deep learning to segment pelvic bones: Large-scale ct datasets and baseline models (2021).  
145 arXiv:2012.08721.  
146 URL <https://arxiv.org/abs/2012.08721>
- 147 [17] A. Sekuboyina, M. E. Husseini, A. Bayat, M. Löffler, H. Liebl, H. Li, G. Tetteh, J. Kukačka,  
148 C. Payer, D. Štern, M. Urschler, M. Chen, D. Cheng, N. Lessmann, Y. Hu, T. Wang, D. Yang,  
149 D. Xu, F. Ambellan, T. Amiranashvili, M. Ehlke, H. Lamecker, S. Lehnert, M. Lirio, N. P.  
150 de Olaguer, H. Ramm, M. Sahu, A. Tack, S. Zachow, T. Jiang, X. Ma, C. Angerman, X. Wang,  
151 K. Brown, A. Kirszenberg, Élodie Puybareau, D. Chen, Y. Bai, B. H. Rapazzo, T. Yeah,  
152 A. Zhang, S. Xu, F. Hou, Z. He, C. Zeng, Z. Xiangshang, X. Liming, T. J. Netherton, R. P.  
153 Mumme, L. E. Court, Z. Huang, C. He, L.-W. Wang, S. H. Ling, L. D. Huỳnh, N. Boutry,  
154 R. Jakubicek, J. Chmelik, S. Mulay, M. Sivaprakasam, J. C. Paetzold, S. Shit, I. Ezhev,  
155 B. Wiestler, B. Glocker, A. Valentinitisch, M. Rempfler, B. H. Menze, J. S. Kirschke, Verse: A  
156 vertebrae labelling and segmentation benchmark for multi-detector ct images, Medical Image  
157 Analysis 73 (2021) 102166. doi:<https://doi.org/10.1016/j.media.2021.102166>.  
158 URL <https://www.sciencedirect.com/science/article/pii/S1361841521002127>
- 159 [18] L. Maier-Hein, A. Reinke, P. Godau, M. D. Tizabi, F. Buettner, E. Christodoulou, B. Glocker,  
160 F. Isensee, J. Kleesiek, M. Kozubek, M. Reyes, M. A. Riegler, M. Wiesenfarth, A. E.  
161 Kavur, C. H. Sudre, M. Baumgartner, M. Eisenmann, D. Heckmann-Nötz, T. Rädsch,  
162 L. Acion, M. Antonelli, T. Arbel, S. Bakas, A. Benis, P. F. Jäger, et al., Metrics reloaded:  
163 recommendations for image analysis validation, Nature Methods 21 (2024) 195–212. doi:  
164 10.1038/s41592-023-02151-7.  
165 URL <https://doi.org/10.1038/s41592-023-02151-7>
- 166 [19] J. Rony, H. Kervadec, Distorch: A fast gpu implementation of 3d hausdorff distance (2025).  
167 URL <https://github.com/jeromerony/distorch>