

Segment Anything Model (SAM) for Digital Pathology: Assess Zero-shot Segmentation on Whole Slide Imaging

Ruining Deng* (Department of Computer Science, Vanderbilt University, Nashville, TN),
Can Cui* (Department of Computer Science, Vanderbilt University, Nashville, TN),
Quan Liu* (Department of Computer Science, Vanderbilt University, Nashville, TN),
Tianyuan Yao (Department of Computer Science, Vanderbilt University, Nashville, TN),
Lucas W. Remedios (Department of Computer Science, Vanderbilt University, Nashville, TN),
Shunxing Bao (Department of Electrical and Computer Engineering, Vanderbilt University, Nashville, TN),
Bennett A. Landman (Department of Electrical and Computer Engineering, Vanderbilt University, Nashville, TN),
Lee E. Wheless (Department of Dermatology, Vanderbilt University Medical Center, Nashville, TN),
Lori A. Coburn (Department of Medicine, Vanderbilt University Medical Center, Nashville, TN),
Keith T. Wilson (Department of Medicine, Vanderbilt University Medical Center, Nashville, TN),
Yaohong Wang (Department of Anatomical Pathology, UT MD Anderson Cancer Center, Houston, TX),
Shilin Zhao (Department of Biostatistics, Vanderbilt University Medical Center, Nashville, TN),
Agnes B. Fogo (Department of Pathology, Microbiology, and Immunology, Vanderbilt University Medical Center, Nashville, TN),
Haichun Yang (Department of Pathology, Microbiology, and Immunology, Vanderbilt University Medical Center, Nashville, TN),
Yucheng Tang (NVIDIA Corporation, Redmond, WA),
Yuankai Huo (Department of Computer Science, Vanderbilt University, Nashville, TN)

*R. Deng, C. Cui, Q. Liu were contributed equally

*Y. Huo is the corresponding author, e-mail: yuankai.huo@vanderbilt.edu

Abstract

The segment anything model (SAM) was released as a foundation model for image segmentation. The promptable segmentation model was trained by over 1 billion masks on 11M licensed and privacy-respecting images. The model supports zero-shot image segmentation with various segmentation prompts (e.g., points, boxes, masks). It makes the SAM attractive for medical image analysis, especially for digital pathology where the training data are rare. In this study, we evaluate the zero-shot segmentation performance of SAM model on representative segmentation tasks on whole slide imaging (WSI), including (1) tumor segmentation, (2) non-tumor tissue segmentation, (3) cell nuclei segmentation. **Core Results:** The results suggest that the zero-shot SAM model achieves remarkable segmentation performance for large connected objects. However, it does not consistently achieve satisfying performance for dense instance object segmentation, even with 20 prompts (clicks/boxes) on each image. We also summarized the identified limitations for digital pathology: (1) image resolution, (2) multiple scales, (3) prompt selection, and (4) model fine-tuning. In the future, the few-shot fine-tuning with images from downstream pathological segmentation tasks might help the model to achieve better performance in dense object segmentation.

Introduction

Large language models (e.g., ChatGPT [6] and GPT-4 [7]), are leading a paradigm shift in natural language processing with strong zero-shot and few-shot generalization capabilities. This development has encouraged researchers to develop large-scale vision foundation models. While the first successful “foundation models” [8] in computer vision have focused on pre-training ap-

proaches (e.g., CLIP [9] and ALIGN [10]) and generative AI applications (e.g., DALL-E [13]), they have not been specifically designed for image segmentation tasks [14]. Segmenting objects (e.g., tumor, tissue, cell nuclei) for whole slide imaging (WSI) data is an essential task for digital pathology, deep learning models typically necessitate well-delineated training data. Obtaining these gold-standard data from clinical experts can be challenging due to privacy regulations, intensive manual efforts, insufficient reproducibility, and complicated annotation processes [16]. Hence, zero-shot image segmentation [20] is desired, where the model can accurately segment pathological images without prior exposure to the domain data during training.

Recently, the “Segment Anything Model” (SAM) [14] was proposed as a foundation model for image segmentation. The model has been trained on over 1 billion masks on 11 million licensed and privacy-respecting images. Furthermore, the model supports zero-shot image segmentation with various segmentation prompts (e.g., points, boxes, and masks). This feature makes it particularly attractive for pathological image analysis where the labeled training data are rare and expensive.

In this study, we assess the zero-shot segmentation performance of the SAM model on representative segmentation tasks, including (1) tumor segmentation [18], (2) tissue segmentation [19], and (3) cell nuclei segmentation [21]. Our study reveals that the SAM model has some limitations and performance gaps compared to state-of-the-art (SOTA) domain-specific models.

Experiments and Performance

We obtained the source code and the trained model from <https://segment-anything.com>. To ensure scalable assessments, all experiments were performed directly using Python,

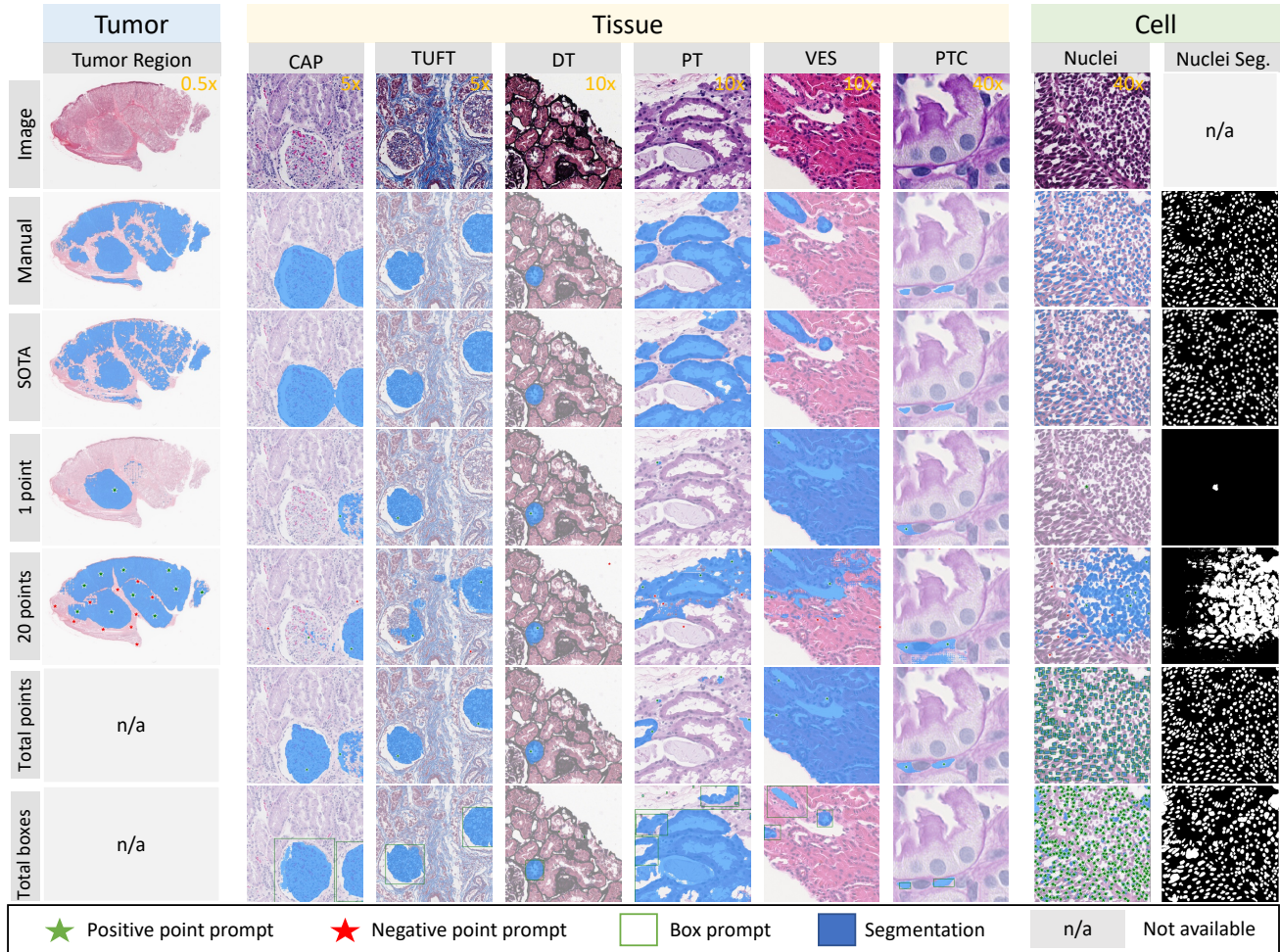


Figure 1. Qualitative segmentation results. The SOTA methods are compared with SAM method with different prompt strategies.

rather than relying on the Demo website. The results are presented in Figure 1 and Table 1.

Tumor Segmentation. The whole-slide images (WSIs) of skin cancer patients were obtained from the Cancer Genome Atlas (TCGA) datasets (TCGA Research Network: <https://www.cancer.gov/tcga>). We employed SimTriplet [18] approach as the SOTA method, with the same testing cohort to make a fair comparison. In order to be compatible with the SAM segmentation model, the WSI inputs were scaled down 80 times from a resolution of $40\times$, resulting in an average size of 860×1279 pixels. We evaluated two different scenarios: (1) SAM with a single positive point prompt, and (2) SAM with 20 point prompts (10 positive and 10 negative points). The prompts were randomly selected from manual annotations, with positive prompt points being selected from the tumor region and negative prompt points being selected from the non-tumor region.

Tissue Segmentation. A total of 1,751 regions of interest (ROIs) images were obtained from 459 WSIs belonging to 125 patients diagnosed with Minimal Change Diseases. These images were manually segmented to identify six structurally normal pathological primitives [12], using digital renal biopsies from the

NEPTUNE study [11]. To form a test cohort for multi-tissue segmentation, we captured 8,359 patches measuring 256×256 pixels. For comparison, We employed Omni-Seg [19] approach as the SOTA method. The tissue types consist of the glomerular unit (CAP), glomerular tuft (TUFT), distal tubular (DT), proximal tubular (PT), arteries (VES), and peritubular capillaries (PTC). For the SAM method, we evaluated four different scenarios: (1) SAM with a single positive point prompt, (2) SAM with 20 point prompts (10 positive and 10 negative points), and (3)/(4) SAM with all points/boxes on every single instance object, which served as a theoretical upper bound for SAM. We randomly selected point prompts from the manual annotations and eroded each connected component with a 10×10 filter to generate at most one random point. For the box prompts, we used the bounding box of each connected component.

Cell nuclei Segmentation. The dataset for nuclei segmentation was obtained from the MoNuSeg challenge [17]. It contains H&E stained images at $40\times$ magnification with 1000×1000 pixels from the TCGA dataset, along with corresponding annotations of nuclear boundaries. The MoNuSeg dataset includes 30 images for training and 14 for testing. We evaluated the performance of

* Compare SAM with state-of-the-art (SOTA) methods.(Unit: Dice score)

Method	Prompts	Tissue							
		Tumor	5×		10×			40×	Cell
		0.5×	CAP	TUFT	DT	PT	VES	PTC	40×
		Tumor							Nuclei
SOTA	no prompt	71.98	96.50	96.59	81.01	89.80	85.05	77.23	81.77
SAM	1 point	58.71	78.08	80.11	58.93	49.72	65.26	67.03	1.95
SAM	20 points	74.98	80.12	79.92	60.35	66.57	68.51	64.63	41.65
SAM	total points	n/a	88.10	89.65	70.21	73.19	67.04	67.61	69.50
SAM	total boxes	n/a	95.23	96.49	89.97	86.77	87.44	87.18	88.30

total points/boxes: we place points/boxes on every single instance object (based on the known ground truth) as a theoretical upper bound of SAM. Note that it is impractical in real applications.

SAM models against the BEDs model [21], a competitive nuclei segmentation model trained on the MoNuSeg training data. The prompt method and evaluation are as described in §Tissue Segmentation.

Limitations on Digital Pathology

The SAM models achieve remarkable performance under zero-shot learning scenarios. However, we identified several limitations during our assessment.

Image resolution. The average training image resolution of SAM is 3300×4950 pixels [14], which is significantly smaller than Giga-pixel WSI data ($> 10^9$ pixels). Moreover, analyzing WSI data at the patch level may result in an impractical number of interactions, even if only a small number of points or bounding boxes are marked for each patch.

Multiple scales. Multi-scale is a significant feature in digital pathology. Different tissue types have their optimal image resolution (as shown in Table 1). For instance, at the optimal resolution for CAP segmentation (5× scale), it is difficult to achieve good segmentation for PTC. However, zooming in (40× scale) would result in nearly 100 times more patches.

Prompt selection. Firstly, to achieve decent segmentation performance in zero-shot learning scenarios, a considerable number of prompts are still necessary. Secondly, the segmentation performance heavily depends on the quality of prompt selection. Another concern related to segmentation performance is inter-rater and intra-rater reproducibility of prompt-based segmentation.

Model fine-tuning. Currently, tedious manual prompt placements are still necessary for segmentation tasks with significant domain heterogeneities. A reasonable online/offline fine-tuning strategy is necessary to propagate the knowledge obtained from manual prompts to larger-scale automatic segmentation on Giga-pixel WSI data.

Conclusion

The zero-shot setting of SAM enables domain users to segment heterogeneous objects in digital pathology without undergoing a heavy training process. The results suggest that the zero-shot SAM model achieves remarkable segmentation performance for large connected objects. However, it does not consistently achieve satisfying performance for dense instance object segmentation, even with 20 prompts (clicks/boxes) on each image. Nonetheless, several limitations still exist and require further investigation for digital pathology.

Acknowledgment

This research was supported by NIH R01DK135597 (Huo), The Leona M. and Harry B. Helmsley Charitable Trust grant G-1903-03793 and G-2103-05128, NSF CAREER 1452485, NSF 2040462, NCRR Grant UL1 RR024975-01 (now at NCATS Grant 2 UL1 TR000445-06), NIH NIDDK DK56942 (ABF), DoD HT94252310003 (Yang), NIH R01DK128200 (Wilson), the VA grants I01CX002662, I01CX002171 and I01CX002473, VUMC Digestive Disease Research Center supported by NIH grant P30DK058404, NVIDIA hardware grant, resources of AC-CRE at Vanderbilt University. This work was supported by Integrated Training in Engineering and Diabetes, grant number T32 DK101003.

References

- [1] John Doe, Recent Progress in Digital Halftoning II, IS&T, Springfield, VA, 1999, pg. 173.
- [2] John Doe, Digital Imaging, J. Imaging. Sci. and Technol., 42, 112 (1998).
- [3] John Doe, An Inexpensive Micro-Goniophotometry You Can Build, Proc. PICS, pg. 179. (1998).
- [4] Leslie Lamport, *TEX: A Document Preparation System*, Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1986.
- [5] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean, “Distilling the Knowledge in a Neural Network,” arXiv preprint arXiv:1503.02531, 2015.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, *et al.*, “Language Models are Few-Shot Learners,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 1877–1901, 2020.
- [7] OpenAI, “GPT-4 Technical Report,” arXiv preprint arXiv:2303.08774, 2023.
- [8] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, *et al.*, “On the Opportunities and Risks of Foundation Models,” arXiv preprint arXiv:2108.07258, 2021.
- [9] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, *et al.*, “Learning Transferable Visual Models from Natural Language Supervision,” in *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 8748–8763, 2021.
- [10] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu

Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig, “Scaling Up Visual and Vision-Language Representation Learning With Noisy Text Supervision,” in *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 4904–4916, 2021.

- [11] Laura Barisoni, Cynthia C. Nast, J. Charles Jennette, Jeffrey B. Hodgins, Andrew M. Herzenberg, Kevin V. Lemley, Catherine M. Conway, Jeffrey B. Kopp, Matthias Kretzler, Christa Lienczewski, *et al.*, “Digital Pathology Evaluation in the Multicenter Nephrotic Syndrome Study Network (NEPTUNE),” *Clinical Journal of the American Society of Nephrology*, vol. 8, no. 8, pp. 1449–1459, 2013.
- [12] Catherine P. Jayapandian, Yijiang Chen, Andrew R. Janowczyk, Matthew B. Palmer, Clarissa A. Cassol, Miroslav Sekulic, Jeffrey B. Hodgins, Jarcy Zee, Stephen M. Hewitt, John O’Toole, *et al.*, “Development and Evaluation of Deep Learning-Based Segmentation of Histologic Structures in the Kidney Cortex With Multiple Histologic Stains,” *Kidney International*, vol. 99, no. 1, pp. 86–101, 2021.
- [13] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever, “Zero-Shot Text-to-Image Generation,” in *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 8821–8831, 2021.
- [14] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, *et al.*, “Segment Anything,” arXiv preprint arXiv:2304.02643, 2023.
- [15] Mohammad Hesam Hesamian, Wenjing Jia, Xiangjian He, and Paul Kennedy, “Deep Learning Techniques for Medical Image Segmentation: Achievements and Challenges,” *Journal of Digital Imaging*, vol. 32, pp. 582–596, 2019.
- [16] Yuankai Huo, Ruining Deng, Quan Liu, Agnes B. Fogo, and Haichun Yang, “AI Applications in Renal Pathology,” *Kidney International*, vol. 99, no. 6, pp. 1309–1320, 2021.
- [17] Neeraj Kumar, Ruchika Verma, Deepak Anand, Yanning Zhou, Omer Fahri Onder, Efstratios Tsougenis, Hao Chen, Pheng-Ann Heng, Jiahui Li, Zhiqiang Hu, *et al.*, “A Multi-Organ Nucleus Segmentation Challenge,” *IEEE Transactions on Medical Imaging*, vol. 39, no. 5, pp. 1380–1391, 2019.
- [18] Quan Liu, Peter C. Louis, Yuzhe Lu, Aadarsh Jha, Mengyang Zhao, Ruining Deng, Tianyuan Yao, Joseph T. Roland, Haichun Yang, Shilin Zhao, *et al.*, “SimTriplet: Simple Triplet Representation Learning With a Single GPU,” in *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021*, pp. 102–112, Springer, 2021.
- [19] Ruining Deng, Quan Liu, Can Cui, Tianyuan Yao, Jun Long, Zuhayr Asad, R. Michael Womick, Zheyu Zhu, Agnes B. Fogo, Shilin Zhao, *et al.*, “Omni-Seg: A Scale-Aware Dynamic Network for Renal Pathological Image Segmentation,” *IEEE Transactions on Biomedical Engineering*, 2023.
- [20] Wei Wang, Vincent W. Zheng, Han Yu, and Chunyan Miao, “A Survey of Zero-Shot Learning: Settings, Methods, and Applications,” *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 10, no. 2, pp. 1–37, 2019.
- [21] Xing Li, Haichun Yang, Jiabin He, Aadarsh Jha, Agnes B. Fogo, Lee E. Wheless, Shilin Zhao, and Yuankai Huo, “BEDS: Bagging Ensemble Deep Segmentation for Nucleus Segmentation With Testing Stage Stain Augmentation,” in *Proceedings of the 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, pp. 659–662, IEEE, 2021.

bilt University, working with Dr. Yuankai Huo. He received his Bachelor’s degree from China University of Mining and Technology, Beijing. In addition to being a research assistant at Vanderbilt University, Mr. Deng was a visiting scholar at the University of Notre Dame and a member of the research staff at the Guangdong Provincial Cardiovascular Institute. Recently, he also served as an imaging scientist intern at Roche Diagnostics USA.

Author Biography

Ruining Deng is a PhD candidate in Computer Science at Vander-