

HOW WELL DO SUPERVISED 3D MODELS TRANSFER TO MEDICAL IMAGING TASKS?

Wenxuan Li Alan Yuille Zongwei Zhou*

Johns Hopkins University

<https://github.com/MrGiovanni/SuPreM>

ABSTRACT

The pre-training and fine-tuning paradigm has become prominent in transfer learning. For example, if the model is pre-trained on ImageNet and then fine-tuned to PASCAL, it can significantly outperform that trained on PASCAL from scratch. While ImageNet pre-training has shown enormous success, it is formed in 2D, and the learned features are for classification tasks; when transferring to more diverse tasks, like 3D image segmentation, its performance is inevitably compromised due to the deviation from the original ImageNet context. A significant challenge lies in the lack of large, annotated 3D datasets rivaling the scale of ImageNet for model pre-training. To overcome this challenge, we make two contributions. Firstly, we construct AbdomenAtlas 1.1 that comprises 9,262 three-dimensional computed tomography (CT) volumes with high-quality, per-voxel annotations of 25 anatomical structures and pseudo annotations of seven tumor types. Secondly, we develop a suite of models that are pre-trained on our AbdomenAtlas 1.1 for transfer learning. Our preliminary analyses indicate that the model trained only with 21 CT volumes, 672 masks, and 40 GPU hours has a transfer learning ability similar to the model trained with 5,050 (unlabeled) CT volumes and 1,152 GPU hours. More importantly, the transfer learning ability of supervised models can further scale up with larger annotated datasets, achieving significantly better performance than preexisting pre-trained models, irrespective of their pre-training methodologies or data sources. We hope this study can facilitate collective efforts in constructing larger 3D medical datasets and more releases of supervised pre-trained models.

1 INTRODUCTION

Pre-training and fine-tuning is a widely adopted transfer learning paradigm (Zoph et al., 2020). Given the relationship across different vision tasks, a model pre-trained on one dataset is expected to benefit another. Over the past few decades, pre-training has been important in AI development (Kumar, 2017; Radford et al., 2021). For 2D vision tasks, there are two available options: (i) supervised pre-training and (ii) self-supervised pre-training, but for 3D vision tasks, option (i) is often not available simply due to the lack of large, annotated 3D volumetric datasets (Wang et al., 2022).

Supervised pre-training can learn image features that are transferable to many target tasks. It has been common practice to pre-train models using ImageNet and then fine-tune the model on target tasks that often have less training data, e.g., PASCAL. However, two challenges arise in ImageNet pre-training. Firstly, ImageNet predominantly comprises 2D images, leaving a palpable void in large-scale 3D datasets and investigation in 3D transfer learning (Huang et al., 2023). Secondly, ImageNet is intended for image classification, so the benefit for segmentation (and other vision tasks) can be somewhat compromised (He et al., 2019). If such an ImageNet-like dataset exists—formed in 3D and annotated per voxel—supervised pre-trained models are expected to transfer better to 3D image segmentation than self-supervised ones for two reasons.

1. **Supervised pre-training is more efficient in data and computation because of its explicit learning objective.** While self-supervised pre-training can learn features without manual annotation, it often requires a large corpus of datasets (Xiao et al., 2022). Extracting meaningful

*Correspondence to Zongwei Zhou (ZZHOU82@JH.EDU).

features directly from raw, unlabeled data is inherently challenging. Unlabeled data have a high degree of redundancy (Haghighi et al., 2020; 2021) and noise (Mahajan et al., 2018), which can complicate the learning process. Therefore, self-supervised pre-training often calls for greater computational resources and time to match the outcomes achieved by supervised pre-training (Chen et al., 2020a; Tang et al., 2022). We have quantified the improved data and computational efficiency from perspectives of both pre-training (Figure 2a; 99.6% fewer data) and fine-tuning (Figure 2b; 66% less computation). Specifically, the model trained with 21 CT volumes, 672 masks, and 40 GPU hours shows transfer learning ability similar to that trained with 5,050 CT volumes and 1,152 GPU hours, highlighting the remarkable efficiency of supervised pre-training.

2. **Supervised pre-training enables the model to learn image features that are relevant to image segmentation.** Self-supervised pre-training must extract image features from raw, unlabeled data using pretext tasks such as mask image modeling (Chen et al., 2019a; Tao et al., 2020; Zhou et al., 2021b; He et al., 2022), instance discrimination (Xie et al., 2020; Chaitanya et al., 2020; Shekoofeh et al., 2021), etc. Despite their efficacy in pre-training, these pretext tasks share no obvious relation to the target image segmentation. In contrast, supervised pre-training uses semantically meaningful annotations (e.g., organ/tumor segmentation) as supervision, with which the model can mimic the behavior of medical professionals—identifying the edge and boundary of specific anatomical structures. As a result, the pre-training is interpretable, and the learned features are expected to be relevant to image segmentation tasks (Zamir et al., 2018; Ilharco et al., 2022; You et al., 2022). We have demonstrated that the learned features can be *direct inference* for organ segmentation on CT volumes collected from hospitals worldwide (Table 3; evaluated on three novel hospitals). The features learned by supervision can also be *fine-tuned* to perform novel class segmentation (unseen in the pre-training) with higher accuracy and less annotated data than the features learned by self-supervision (Table 4; evaluated on 63 novel classes).

This paper seeks to answer the question *how well the model transfers to 3D medical imaging tasks* IF it is pre-trained on large, annotated 3D datasets. Naturally, we start with creating an *IF* dataset at a massive scale. **Firstly**, we construct a dataset (termed AbdomenAtlas 1.1¹) of 9,262 CT volumes with per-voxel annotations of 25 anatomical structures and pseudo annotations of seven types of tumors. This large-scale, fully-annotated dataset enables us to train models in a fully supervised manner using multi-organ segmentation as the pretext task. As reviewed in Table 1, this dataset is much more extensive (considering both the number of CT volumes and annotated classes) than public datasets (Wasserthal et al., 2022; Ma et al., 2022; Qu et al., 2023). Scaling experiments in §3.1 suggested that pre-training models on more annotated datasets can further improve the transfer learning ability. **Secondly**, we develop a suite of **Supervised Pre-trained Models**, termed SuPreM, that combined the good of large-scale datasets and per-voxel annotations, demonstrating the efficacy across a range of target segmentation tasks. As reported in §3.2, some of the dominant segmentation backbones have been pre-trained and will be available to the public. Current pre-trained backbones are U-Net (CNN-type) (Ronneberger et al., 2015), SegResNet (CNN-type) (Myronenko, 2019), and Swin UNETR (Transformer-type) (Tang et al., 2022), and more backbones will be added along time.

In prospective endeavors, we anticipate that the expansion of datasets and annotations will not only enhance feature learning, as demonstrated in this study, but also promote the development of advanced AI algorithms and benchmark the state of the art in terms of segmentation performance, inference efficiency, and domain generalizability.

2 BRIEF HISTORY: SUPERVISED PRE-TRAINING

In a major initiative aimed at developing widely transferable AI models—known as Foundation Models in the medical domain (Moor et al., 2023; Butoi et al., 2023; Ma & Wang, 2023a)—one faces a critical decision: *should the focus of pre-training be supervised or self-supervised?* While human annotations undeniably improve task-specific performance, such as semantic segmentation, the best strategy for learning generic image features that can be transferable across a spectrum of tasks has yet to be determined. For 2D vision tasks, the advent of ImageNet (Deng et al., 2009) makes it possible to debate the merits and limitations of supervised pre-trained models for transfer learning compared

¹Segmentation is fundamental in the medical domain (Ma & Wang, 2023b). It can be viewed as a per-voxel classification task. Therefore, the per-voxel supervision used in our pre-training (272.7B annotated voxels) is much stronger than the per-image supervision used in ImageNet pre-training (14M images).

with self-supervised ones. We refer the readers to [Yang et al. \(2020\)](#) and [Tendle & Hasan \(2021\)](#) for a plethora of viewpoints from either side. In essence, the debates are about clarifying the learning objective (loss function) of emulating human vision ([Zhou, 2021](#)).

The learning objective of supervised pre-training is to minimize the discrepancy between AI predictions and semantic labels annotated by humans. Over the years, supervised pre-training on ImageNet has shown marked success in transfer learning ([Yosinski et al., 2014](#)). Moreover, the transfer learning ability can be further enhanced when models are trained on increasingly expansive datasets, such as ImageNet-21K ([Kolesnikov et al., 2020](#)), Instagram ([Mahajan et al., 2018](#)), JFT-300M ([Sun et al., 2017](#)), and JFT-3B ([Zhai et al., 2022](#)). In general, supervised pre-training exhibits clear advantages over self-supervised pre-training when sizable annotated datasets are available ([Steiner et al., 2021](#); [Ridnik et al., 2021](#)). However, acquiring millions of manual annotations is labor-intensive, time-consuming, and challenging to scale—but certainly not impossible—evidenced by several recent influential endeavors ([Kuznetsova et al., 2020](#); [Mei et al., 2022](#); [Kirillov et al., 2023](#); [Bai et al., 2023](#)).

On the other hand, self-supervised pre-training offers an alternative by enabling AI models to learn from raw, unlabeled data ([Jing & Tian, 2020](#); [Zoph et al., 2020](#); [Ren et al., 2022](#); [2023](#)), thus reducing the need for manual annotation. Self-supervised pre-training has historically lagged behind the state-of-the-art supervised pre-training in ImageNet benchmarks ([Pathak et al., 2016](#); [Noroozi & Favaro, 2016](#)). The recent pace of progress in self-supervised pre-training has yielded models whose performance not only matches but, at times, surpasses those achieved by supervised pre-training ([Chen et al., 2020a](#); [Grill et al., 2020](#); [Chen et al., 2020b](#); [Zhou et al., 2021a](#); [Wei et al., 2022](#)). This has raised hopes that self-supervised pre-training could indeed replace the ubiquitous supervised pre-training in advanced computer vision going forward. The caveat, however, is the significant demand for both data and computational power, often exceeding the resources available in academic settings. For example, [He et al. \(2020\)](#) have demonstrated that self-supervised features trained on 1B images (a factor of $714\times$ larger) can transfer comparably or better than ImageNet features.

Supervised pre-training on ImageNet has demonstrated benefit for 2D medical image tasks after transfer learning ([Tajbakhsh et al., 2016](#); [Shin et al., 2016](#); [Zhou et al., 2017](#)). Unfortunately, it has been constrained for 3D medical imaging tasks due to the lack of a 3D counterpart to ImageNet. Although there are a great number of raw, unlabeled medical images available ([Team, 2011](#); [Baxter et al., 2023](#); [Zhao et al., 2023](#); [Saenz et al., 2024](#)), annotating these images is a labor-intensive undertaking for professionals. Our contribution to a large, annotated 3D dataset could spark the debate of whether self-supervised or supervised pre-training leads to better performance and data/computational efficiency, which would not be possible without the invention of a dataset of such a scale.

3 MATERIAL & METHOD

We constructed an AbdomenAtlas 1.1 dataset comprising **9,262** three-dimensional CT volumes and over **251,323** masks spanning **25** anatomical structures and **7** types of tumors. In addition, we released a suite of supervised pre-trained models (SuPreM) to benefit 3D medical imaging tasks.

3.1 EXTENSIVE DATASET: ABDOMENATLAS 1.1

Interactive segmentation, an integration of AI algorithms and human expertise, was used to create AbdomenAtlas 1.1 in a semi-automatic procedure. We recruited a team of ten radiologists to perform manual annotations to ensure the annotation quality². Given the complexity of 3D data, rather than annotating the entire dataset voxel by voxel, we asked the radiologists to focus on the most important CT volumes and regions therein. In doing so, an importance score for each volume was computed, derived from the uncertainty, consistency, and overlap ([Qu et al., 2023](#)). Six junior radiologists revised the annotations predicted by AI under the supervision of four senior radiologists, and in turn, AI improved its predictions by learning from these revised annotations. This interactive procedure continued to enhance the quality of annotations until no major revision was required from the radiologists. Subsequently, four senior radiologists went through the final visualizations for all the annotations, detecting and revising major errors as needed before the dataset was released. Annotation tools employed included a licensed version from [Pair](#) and an open-source [MONAI Label](#).

²Ensuring high-quality annotations is costly and time-consuming, yet it is critical for transfer learning and reducing ambiguity when training AI models for image segmentation.

Table 1: Contribution #1: An extensive dataset of 9,262 CT volumes with per-voxel annotations of 25 anatomical structures. This dataset is unprecedented in terms of data and annotation scales, providing over 251,323 organ/tumor masks and 2,789,975 annotated images that are taken from 88 hospitals worldwide. In 2009, before the advent of ImageNet (Deng et al., 2009), it was challenging to empower an AI model with generalized image representation using a small or even medium size of labeled data, the same situation, we believe, that presents in 3D medical image analysis today. As seen in the table, the annotations of public datasets are limited, partial, and incomplete, and the CT volumes in these datasets are often biased toward specific populations, medical centers, and countries. Our constructed dataset mitigates these gaps, representing a significant leap forward in the field. The CT volumes in datasets 1–17 are used to construct AbdomenAtlas 1.1. The domain gap across these datasets is illustrated in Appendix A.1.

dataset (year) [source]	# of organ	# of [†] volume	# of center	dataset (year) [source]	# of organ	# of [†] volume	# of center
1. Pancreas-CT (2015) [link]	1	42	1	2. CHAOS (2018) [link]	4	20	1
3. CT-ORG (2020) [link]	5	140	8	4. BTCV (2015) [link]	12	47	1
5. AMOS22 (2022) [link]	15	200	2	6. WORD (2021) [link]	16	120	1
7-12. MSD CT Tasks (2021) [link]	9	945	1	13. LiTS (2019) [link]	1	131	7
14. AbdomenCT-1K (2021) [link]	4	1,050	12	15. KiTS (2020) [link]	1	489	1
16. FLARE'23 (2022) [link]	13	4,100	30	17. Trauma Det. (2023) [link]	0	4,711	23
18. AbdomenAtlas 1.0 (2023) [link]	9	5,195	26	19. AbdomenAtlas 1.1	25	9,262 [‡]	88

[†]Our reported number of CT volumes may differ from original publications, as some CT volumes are reserved for validation purposes.

[‡]The number of CT volumes in AbdomenAtlas 1.1 is lower than the sum of datasets 1–17 due to overlaps within these public datasets.

AbdomenAtlas 1.1 is a composite dataset that unifies CT volumes from public datasets 1–17 as summarized in Table 1. AbdomenAtlas 1.1 presents a level of diversity because the CT volumes are sourced from 88 hospitals worldwide, including pre, portal, arterial, and delayed phases. The gap between these CT volumes includes changes in image quality due to different acquisition parameters, reconstruction kernels, and contrast enhancement, shown in Appendix A.1. Moreover, we provide per-voxel annotations for 25 anatomical structures, including 16 abdominal organs, two thorax organs, five vascular structures, and two skeletal structures. We also provide pseudo annotations for seven types of tumors, namely liver, kidneys, pancreatic, hepatic vessel, lung, colon tumors, and kidney cysts. In total, more than 272.7B voxels are annotated in AbdomenAtlas 1.1, marking a significant leap compared with the 4.3B voxels annotated in the public datasets, amplifying the annotations by a factor of $63.4\times$ (shown in Appendix Figure 4). The high annotation quality is due to the uniform annotation standards described in Appendix A.2. *We commit to releasing AbdomenAtlas 1.1 to the public.* However, this dataset, the largest public per-voxel annotated CT collection by far, accounts for around 0.01% of the CT volumes annually acquired in the United States (Papanicolas et al., 2018). Therefore, cross-institutional collaboration is crucial for accelerating data sharing, annotation, and AI development (Saenz et al., 2024).

3.2 A SUITE OF PRE-TRAINED MODELS: SUPREM

The magnitude of our AbdomenAtlas 1.1 is unprecedented in terms of data and annotations. One of the advantages is that it enables us to train AI models in both a supervised and self-supervised manner. At the time this paper is written, neither supervised nor self-supervised pre-training has been performed on this scale of dataset (9,262 volumetric data)³. We have developed models (termed SuPreM) pre-trained on data and annotations in AbdomenAtlas 1.1, which leverage established CNN backbones, such as U-Net and SegResNet, as well as Transformer backbones, such as Swin UNETR. With the growing trend of using pre-trained models, we have maintained a standardized, accessible **model repository** for sharing public model weights as well as a suite of supervised pre-trained models (SuPreM) released by us. Releasing pre-trained models should be considered a marked contribution as they offer an alternative way of knowledge sharing while protecting patient privacy (Sellersgren et al., 2022; Zhang & Metaxas, 2023; Ma et al., 2023a). In this study, all of the models in SuPreM follow pre-training and fine-tuning configurations as below.

³For supervised pre-training, the largest study to date was by Liu et al. (2023), which was developed on 3,410 (2,100 for training and 1,310 for validation) annotated CT volumes. For self-supervised pre-training, the largest one was by Tang et al. (2022), which was trained on 5,050 unannotated CT volumes. Concurrently, Valanarasu et al. (2023) pre-trained a model on 50K volumes of CT and MRI using self-supervised learning.

Table 2: Contribution #2: A suite of pre-trained models (termed SuPreM) comprising several widely recognized AI models. We provide pre-trained AI models based on CNN, Transformer, and their hybrid versions, and more AI models will be added. Each model was supervised pre-trained on large datasets and per-voxel annotations from AbdomenAtlas 1.1. Compared with learning from scratch and publicly available models, fine-tuning the models in SuPreM consistently achieves state-of-the-art organ and tumor segmentation performance on two datasets. All of the results, including the mean and standard deviation (mean \pm s.d.) across ten trials. In addition, we have further performed an independent two-sample t -test between learning from scratch and fine-tuning models in our SuPreM. The performance gain is statistically significant at the $P = 0.05$ level, with highlighting in **light red**.

model (# of param)	pre-training	TotalSegmentator v1			proprietary dataset		
		organ	muscle	cardiac	organ	gastro	cardiac
U-Net (2015) family (19.08M)	scratch	88.9 \pm 0.6	92.9 \pm 0.4	88.8 \pm 0.7	85.6 \pm 0.5	69.8 \pm 1.2	38.1 \pm 1.1
	Zhou et al. (2019b)	87.8	90.1	86.3	80.1	65.5	36.9
	Chen et al. (2019b)	86.9	91.4	87.4	79.0	66.2	36.7
	Xie et al. (2022)	88.5	92.9	89.0	-	-	-
	Zhang et al. (2021)	89.3	93.8	89.1	85.7	72.7	38.3
	SuPreM	92.1\pm0.3	95.4\pm0.1	92.2\pm0.3	90.8\pm0.2	76.2\pm0.8	70.5\pm0.5
Swin UNETR (2021) (62.19M)	scratch	86.4 \pm 0.5	88.8 \pm 0.5	84.5 \pm 0.6	77.3 \pm 0.9	65.9 \pm 1.7	35.5 \pm 1.4
	Tang et al. (2022)	89.3	93.8	88.3	87.9	72.5	38.9
	Liu et al. (2023)	89.7	94.1	89.4	89.1	74.6	67.6
	SuPreM	91.3\pm0.3	94.6\pm0.2	90.3\pm0.3	90.4\pm0.7	75.9\pm1.2	69.8\pm0.9
SegResNet (2019) (4.7M)	scratch	88.6 \pm 0.5	91.3 \pm 0.4	89.8 \pm 0.4	80.6 \pm 0.8	67.0 \pm 1.4	36.0 \pm 1.3
	SuPreM	91.3\pm0.5	94.0\pm0.1	91.3\pm0.5	86.6\pm0.3	73.7\pm1.0	67.9\pm0.8

To perform a fair and rigorous comparison, we benchmarked with public pre-training methods by pre-training SuPreM using 2,100 CT volumes (same as Liu et al. (2023) and fewer than Tang et al. (2022)) in Tables 2, 4 and Figures 1, 2b, 3. Then, we scaled up the number of CT volumes for pre-training to 9,262 CT volumes to perform direct inference in Table 3. Lastly, we scaled down the number of CT volumes to 21 to explore the edge of our SuPreM in Figure 2a. All these pre-trained models and configurations have been summarized in Appendix Table 6. The best-performing model was selected based on the highest average DSC score over 32 classes on a validation set of 1,310 CT volumes. Implementation details of both pre-training and fine-tuning can be found in Appendix B.1.

The transfer learning ability is assessed by segmentation performance on two datasets, i.e., TotalSegmentator v1 and a proprietary dataset. Benchmarking results in Table 2 indicate that, in comparison with learning from scratch and with existing public models, those fine-tuned from our SuPreM consistently attain superior organ, muscle, cardiac, and gastro segmentation performance on both datasets. U-Net, as a simple and lightweight segmentation backbone, still performs competitively compared with alternative choices like Swin UNETR. This observation is aligned with the majority of the medical imaging community (Isensee et al., 2021; Eisenmann et al., 2023), suggesting that more exploration is needed for advancing segmentation backbones. Moreover, in the scenarios of either small data regimes shown in Figure 1 or large data regimes shown in Appendix Figure 7a–d, supervised models transfer better than their self-supervised counterparts. In summary, our SuPreM surpasses all existing 3D pre-trained models by a large margin in transfer learning performance, irrespective of their pre-training methodologies or data sources.

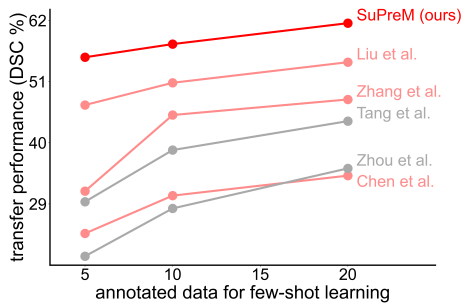


Figure 1: We present the transfer performance on a proprietary dataset with few-shot examples ($N = 5, 10, 20$). The transfer performance (Y-axis) stands for the average DSC score across 20-class organ segmentation and 3-class tumor segmentation. Generally speaking, in a few-shot learning setting, supervised pre-trained models (in red) transfer better than self-supervised pre-trained models (in gray). Notably, our SuPreM achieves the best transfer performance over other well-known publicly available models.

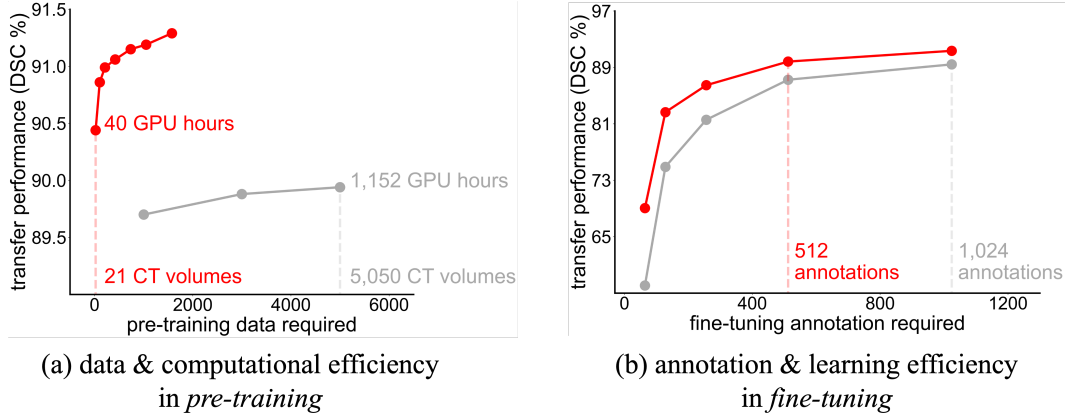


Figure 2: **Analysis of pre-training and fine-tuning efficiency.** For a fair comparison, both supervised (in red) and self-supervised (in gray) models use Swin UNETR as the backbone, and the compared self-supervised pre-training is the current state of the art (Tang et al., 2022). The target task was on TotalSegmentator v1. (a) scales the model transfer learning ability when pre-trained on varying numbers of images. The results indicate a consistent improvement in transfer learning ability when pre-training on more images. The model trained with 21 CT volumes, 672 masks, and 40 GPU hours shows a transfer learning ability similar to that trained with 5,050 CT volumes and 1,152 GPU hours. Specifically, supervised pre-training is more efficient, requiring 99.6% fewer data and 96.5% less computation. (b) assesses the annotation & learning efficiency by fine-tuning models on different numbers of annotated CT volumes from TotalSegmentator. Specifically, SuPreM, fine-tuned on 512 per-voxel annotated CT volumes, can achieve a segmentation performance on par with self-supervised models fine-tuned on 1,024 volumes, reducing 50% manual annotation cost for target tasks.

4 EXPERIMENT & ANALYSIS

4.1 DATA, ANNOTATION, AND COMPUTATIONAL EFFICIENCY

Summary. We demonstrate the remarkable efficiency: (1) SuPreM trained with 21 CT volumes, 672 masks, and 40 GPU hours shows transfer learning ability similar to that trained with 5,050 CT volumes and 1,152 GPU hours. (2) SuPreM requires 50% fewer manual annotations for organ/tumor segmentation than self-supervised pre-training.

Data efficiency for pre-training. As shown in Figure 2a, supervised pre-training requires less data (21 vs. 5,050 CT volumes) for the pretext task than self-supervised pre-training. This discrepancy arises from the inherent differences in their learning objectives and the information they leverage. Supervised pre-training benefits from explicit annotations, which provide direct guidance for the task, i.e., segmentation in this study. The model learns features from both data and annotations, which offer strong and precise supervision. On the other hand, self-supervised learning relies on pretext tasks derived from the raw data, which may offer a more ambiguous learning signal, therefore requiring more examples to capture meaningful features. Importantly, our finding suggests that supervised pre-training is more scalable with increased data. When data are increased from 21 to 1,575 volumes, the transfer learning performance on TotalSegmentator improves from 90.4% to 91.3%. In comparison, for self-supervised pre-training, an increase in data from 1,000 to 5,050 volumes only marginally improves performance from 89.7% to 89.9%. Therefore, supervised pre-training requires significantly less data than self-supervised and is more scalable and effective with increased data.

Annotation efficiency for fine-tuning. We have assessed the annotation efficiency by fine-tuning SuPreM and self-supervised models (Tang et al., 2022) on the TotalSegmentator dataset. Figure 2b suggests that fine-tuning SuPreM can reduce annotation costs for the segmentation task by 50%, averaged over the classes that were not used for pre-training (per-class performance can be found in Appendix Figure 8a–d). Specifically, SuPreM fine-tuned on 512 per-voxel annotated CT volumes can achieve segmentation performance similar to Tang et al. (2022) fine-tuned on 1,024 annotated CT volumes. The fine-tuning performance improvement gets bigger when the number of annotated CT volumes is limited in the target task (e.g., 64, 128, 256). In addition, similar levels of annotation

Table 3: Direct inference on three external datasets. We conduct external validation across four hospitals worldwide. Specifically, our SuPreM—trained on 9,262 CT volumes—is directly inferred on three external datasets, i.e., TotalSegmentator (representing the Central European population from Switzerland; one hospital), DAP Atlas (the Central European population from Germany; two hospitals), and the proprietary dataset (the North American population from the United States; one hospital) measured by DSC scores. For every dataset, we compare the *out-of-distribution* (OOD) performance obtained by SuPreM with *independently and identically distributed* (IID) performance obtained by AI models directly trained on that specific dataset, which are often considered as upper bound performance in domain transfer literature. We find that SuPreM can be generalized well across external datasets without additional fine-tuning, yielding comparable or even superior performance to the IID counterparts, evidenced by the one-sample *t*-test results. Appendix D.1 provides visual examples of anatomical structure segmentation.

class	TotalSegmentator v1		DAP Atlas		our proprietary dataset	
	SuPreM (OOD)	Liu et al. (IID)	SuPreM (OOD)	Jaus et al. (IID)	SuPreM (OOD)	Wang et al. (IID)
spleen	96.0±0.0 ****	93.6	96.8±0.0 ^{ns}	96.8	95.0±0.0 ****	89.6
kidney right	93.3±0.1 *	94.1	96.3±0.1 ****	95.3	92.2±0.0 ****	88.0
kidney left	91.2±0.2 ****	87.7	96.4±0.1 ****	97.4	91.6±0.1 ****	83.9
gall bladder	81.8±0.3 ****	73.9	87.6±0.4 ****	71.2	83.6±0.2 ^{ns}	85.4
liver	96.4±0.1 ^{ns}	96.8	97.3±0.1 ****	98.5	95.0±0.3 ****	91.4
stomach	87.3±0.3 ^{ns}	89.2	95.3±0.2 ****	96.1	92.2±0.1 *	93.6
aorta	80.8±0.4 ****	90.7	90.7±0.5 ****	97.7	73.9±0.3 ****	87.0
postcava	77.9±0.3 ****	82.1	89.1±0.4 ****	95.9	77.7±0.4 **	80.8
pancreas	84.6±0.2 ****	80.8	90.6±0.2 ****	93.7	79.0±0.3 ^{ns}	79.3
average	87.7±0.2 ^{ns}	87.6	93.3±0.2 ****	93.6	86.7±0.2 ^{ns}	86.1

^{ns} $P > 0.05$ * $P \leq 0.05$ ** $P \leq 0.01$ *** $P \leq 0.001$ **** $P \leq 0.0001$

efficiency (reduced 50% cost) are observed when fine-tuning SuPreM on the three-class tumor segmentation task using the proprietary dataset, as presented in Appendix Figure 8e–g.

Computational efficiency for both pre-training and fine-tuning. This efficiency stems, in part, from the reduced data requirements inherent to supervised pre-training, as discussed above. As shown in Figure 2a, supervised pre-training only needs 40 GPU hours to achieve a transfer learning performance comparable to that of self-supervised pre-training, which requires 1,152 GPU hours—a factor increase of $28.8\times$. When fine-tuning on target tasks, such as on a 10% subset of TotalSegmentator in Appendix Figure 9, the supervised pre-trained model converges much faster than the self-supervised one, reducing the GPU hours needed from 60 to 20. This implies that image features learned by supervised pre-training are intrinsically more expressive, enabling the model to seamlessly adapt across a myriad of 3D image segmentation tasks with minimal annotated data for fine-tuning. This computational efficiency makes supervised pre-training a compelling choice for 3D image segmentation without compromising model performance, especially when the large, annotated dataset is available.

4.2 ENHANCED FEATURES FOR NOVEL DATASETS, CLASSES, AND TASKS

Summary. The learned features manifest considerable generalizability and adaptability. The features can *direct inference* for organ segmentation on external datasets of CT volumes taken from different hospitals. The features can also be *fine-tuned* to segment novel organ/tumor classes and classify tumor sub-types with higher accuracy and less annotated data than those learned by self-supervision.

Direct inference on external datasets. AI models trained on a specific dataset often encounter challenges in generalizing to novel datasets when a marked difference—referred to as a *domain gap*—exists between them (Zhang & Metaxas, 2023). While domain adaptation and generalization are prevalent research strategies to mitigate this challenge (Guan & Liu, 2021; Zhou et al., 2022a), we choose to address this issue by training a model on an expansive and diverse dataset (elaborated in Appendix A.1). We assume the domain gap between CT volumes from different hospitals is not as pronounced as those in computer vision. This is because of the relatively standardized nature of computer tomography as an imaging modality, where pixel intensity conveys consistent anatomical significance (Zhou et al., 2022b). AbdomenAtlas 1.1 presents impressive diversity, covering CT volumes with variations in contrast enhancement, reconstruction kernels, CT scanner types, and acquisition parameters. This breadth and diversity are imperative for developing an AI model with the robustness required to accommodate the variations present in novel datasets. We conduct external

Table 4: **Fine-tuning SuPreM on 66 novel classes.** Following the standard transfer learning paradigm, we fine-tune our SuPreM on the segmentation task of novel classes. These tasks include segmenting 19 muscles, 15 cardiac structures, 5 organs, and 24 vertebrae from TotalSegmentator, as well as three fine-grained pancreatic tumor types from the proprietary dataset. It is important to note that these classes were not part of the pre-training of SuPreM. We observe that SuPreM, supervised pre-trained on only a few classes, can transfer better than those self-supervised pre-trained on raw, unlabeled data measured by DSC scores. In other words, it is the task of segmentation itself that can enhance the model’s capability of segmenting novel-class objects. This benefit is much more straightforward and understandable than such self-supervised tasks as contextual prediction, mask image modeling, and instance discrimination in the context of transfer learning. We hypothesize that it is because the model learns to understand the concept of *objectness* in a broader sense through full supervision, as suggested by Kirillov et al. (2023), but this certainly deserves further exploration. In addition, an independent two-sample *t*-test was performed between the self-supervised pre-trained model and the supervised pre-trained model.

novel class	self-super.	super.	Δ	novel class	self-super.	super.	Δ
humerus left	92.8 \pm 0.7	93.2 \pm 0.3 ^{ns}	0.4	vertebrae L5	94.1 \pm 0.2	95.7 \pm 0.3 ^{****}	1.6
humerus right	87.5 \pm 1.0	95.0 \pm 0.5 ^{****}	7.6	vertebrae L4	90.4 \pm 0.6	93.0 \pm 0.5 ^{****}	2.6
... (15 more classes)				... (20 more classes)			
iliopsoas left	84.4 \pm 0.3	85.7 \pm 0.3 ^{****}	1.3	vertebrae C2	86.8 \pm 2.0	91.8 \pm 0.2 ^{****}	5.1
iliopsoas right	87.4 \pm 0.3	88.7 \pm 0.2 ^{****}	1.3	vertebrae C1	87.1 \pm 0.8	87.4 \pm 0.8 ^{ns}	0.3
average (muscle)	93.9 \pm 0.1	94.3 \pm 0.1 ^{****}	0.4	average (vertebrae)	86.4 \pm 0.3	89.2 \pm 0.2 ^{****}	2.7
trachea	93.4 \pm 0.1	93.4 \pm 0.1 ^{ns}	0.0				
heart myocardium	88.9 \pm 0.2	89.8 \pm 0.2 ^{****}	0.9				
... (11 more classes)							
urinary bladder	90.5 \pm 0.9	91.5 \pm 0.9 [*]	1.0	PDAC	53.3 \pm 0.4	53.6 \pm 0.3 [*]	0.3
pulmonary artery	89.0 \pm 0.9	92.0 \pm 0.2 ^{****}	3.0	Cyst	41.5 \pm 0.3	49.4 \pm 0.3 ^{****}	7.9
average (cardiac)	88.9 \pm 0.1	90.7 \pm 0.1 ^{****}	1.8	PanNet	35.5 \pm 0.8	46.0 \pm 0.5 ^{****}	10.5
				average (tumor)	43.4 \pm 0.3	49.7 \pm 0.2 ^{****}	6.2

^{ns} $P > 0.05$ ^{*} $P \leq 0.05$ ^{**} $P \leq 0.01$ ^{***} $P \leq 0.001$ ^{****} $P \leq 0.0001$

validation on several novel datasets sourced from Switzerland and East Asia to challenge the AI model on the data distribution that it has not encountered during the training. This result is referred to as *out-of-distribution* (OOD) performance. For comparison, we also collect the result achieved by dataset-specific AI models—those individually trained on the specific datasets—referred to as *independently and identically distributed* (IID) performance. As shown in Table 3, our SuPreM can be generalized well to novel data distribution without the need for further fine-tuning or adaptation, consistently offering OOD performance that matches or even exceeds that of its IID counterparts.

Fine-tuning on novel classes. The value of transfer learning lies in fine-tuning the pre-trained models on novel scenarios (Zhou et al., 2021b), such as novel classes, image modalities, and vision tasks that are completely unseen during the pre-training. In this study, we evaluate the proficiency of SuPreM when transferred to a wide variety of novel classes for 3D image segmentation tasks⁴. These novel classes include 19 muscles, 15 cardiac structures, 5 organs, and 24 vertebrae from the TotalSegmentator dataset, as well as three fine-grained pancreatic tumor types from the proprietary dataset. As shown in Table 4, our SuPreM, supervised pre-trained on 25 classes, can transfer better to novel classes than those self-supervised models pre-trained on raw, unlabeled data. We find that the pretext task of segmentation itself can enhance the model capability of segmenting novel classes. The benefit of same-task transfer learning, i.e., segmentation as pretext and target task, is much more straightforward and understandable than other pretext tasks such as contextual prediction, mask image modeling, and instance discrimination. Through full supervision in segmentation tasks, the model learns to understand the concept of *objectness*⁵, wherein the model gains a more profound understanding of what characterizes an object. The model does not just recognize predefined objects but begins to understand the foundational factors of objects in general. Such factors include texture, boundary, shape, size, and other low-level visual cues that are often deemed essential for image segmentation. This resonates with our assertion in the introduction: just as classification-based features from ImageNet transfer optimally for classification tasks (Huh et al., 2016; He et al., 2019;

⁴The fine-tuning performance of 17 seen classes is promising, but this is expected because the model is exposed to more examples of these classes in both pre-training and fine-tuning phases.

⁵Objectness refers to the inherent attributes that distinguish something as an object within an image, differentiating it from the background or other entities.

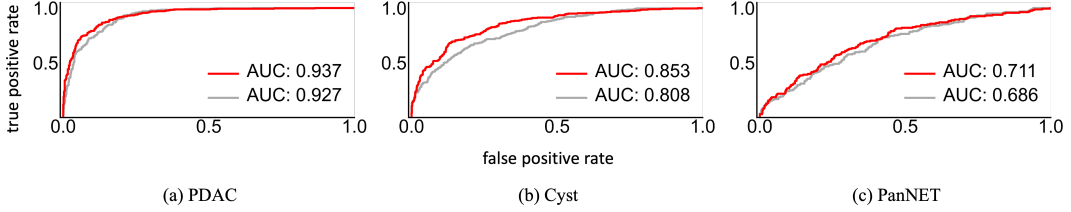


Figure 3: **Fine-tuning SuPreM on fine-grained tumor classification.** We plot receiver operating characteristic (ROC) curves to evaluate the transfer learning performance of tumor classification. Detecting Cysts and PanNETs raises additional challenges for AI because these lesions exhibit a greater variety of texture patterns than PDACs. This diversity in texture patterns is reflected in the values of the Area Under the Curve (AUC) that we obtained. For all three sub-types of pancreatic tumors, SuPreM (in red) demonstrates superior performance over the self-supervised model (Tang et al., 2022) (in gray), showcasing its effectiveness in fine-grained tumor classification.

Zoph et al., 2020; Ridnik et al., 2021), segmentation-based features are optimal for segmentation tasks. Our findings do not negate the value of self-supervised pre-training. With 9,262 CT volumes, should self-supervised pre-training outperforms supervised pre-training in model transferability in the future, its value will be further highlighted by eliminating the need for manual annotations.

Fine-tuning on novel tasks. We have investigated the cross-task transfer learning ability of SuPreM between organ segmentation and fine-grained tumor classification. The distance between the two tasks is much larger than transferring among segmentation tasks. It is challenging to benchmark fine-grained tumor classification, particularly due to the scarcity of annotations in public datasets (often limited to hundreds of tumors). To overcome this limitation, we employed our proprietary dataset (Xia et al., 2022), which contains 1,869 and 1,073 CT volumes (i.e., 1,174 and 684 patients) for training and testing, respectively. These volumes present 3,577 annotated pancreatic tumors, including detailed sub-types: 1,704 PDACs, 945 Cysts, and 928 PanNETs. This extensive dataset enabled us to thoroughly assess the transfer learning ability of SuPreM in tumor-related tasks. Figure 3 shows that supervised models (SuPreM) transfer better to target classification tasks than self-supervised models (Tang et al., 2022), leading to improved Area Under the Curve (AUC) for identifying each tumor type. Notably, the transfer learning results detailed in Appendix D.2 reveal a sensitivity of 86.1% and specificity of 95.4% for PDAC detection. This performance surpasses the average radiologist’s performance in PDAC identification by 27.6% in sensitivity and 4.4% in specificity, as reported in Cao et al. (2023). Moreover, Appendix Figure 8 shows that SuPreM requires 50% fewer manual annotations for fine-grained tumor classification than self-supervised pre-training. This is particularly critical for tumor imaging tasks because annotating tumors requires much more effort and often relies on the availability of pathology reports.

5 CONCLUSION AND DISCUSSION

This study examines the transfer learning ability of supervised models that are pre-trained on 3D annotated datasets and fine-tuned on 3D image segmentation tasks. We start by constructing AbdomenAtlas 1.1, an extensive collection of **9,262** three-dimensional CT volumes with high-quality, per-voxel annotations. The magnitude of this dataset is unprecedented regarding data volume (**2,789,975 images**), granularity of annotations (**251,323 masks**), and inclusive diversity (**88 hospitals**). This dataset facilitates the development of a suite of pre-trained models, termed SuPreM, that can be effectively transferred to a broad spectrum of 3D image segmentation tasks. Notably, SuPreM transfers better than all existing 3D models by a large margin, especially when transferred to under-annotated datasets. The model trained with 21 CT volumes, 672 masks, and 40 GPU hours shows a transfer learning ability similar to that trained with 5,050 CT volumes and 1,152 GPU hours, highlighting the remarkable efficiency of supervised pre-training. We also demonstrate that the learned features can *direct inference* effectively on external datasets and *fine-tune* to segment novel classes and classify multiple types of tumors with higher accuracy and less annotated data than those learned by self-supervision.

ACKNOWLEDGMENTS

This work was supported by the Lustgarten Foundation for Pancreatic Cancer Research and the Patrick J. McGovern Foundation Award. This work has utilized the GPUs provided partially by ASU Research Computing and NVIDIA. We appreciate the effort of the MONAI Team to provide open-source code for the community. We thank Chongyu Qu, Yixiong Chen, Junfei Xiao, Jie Liu, Yucheng Tang, Tiezheng Zhang, Yaoyao Liu, Chen Wei, Fengrui Tian, Yu-Cheng Chou, Angtian Wang, and Dora Zhiyu Yang for their constructive suggestions at several stages of the project.

REFERENCES

- Michela Antonelli, Annika Reinke, Spyridon Bakas, Keyvan Farahani, Bennett A Landman, Geert Litjens, Bjoern Menze, Olaf Ronneberger, Ronald M Summers, Bram van Ginneken, et al. The medical segmentation decathlon. *arXiv preprint arXiv:2106.05735*, 2021.
- Yutong Bai, Xinyang Geng, Karttikeya Mangalam, Amir Bar, Alan Yuille, Trevor Darrell, Jitendra Malik, and Alexei A Efros. Sequential modeling enables scalable learning for large vision models. *arXiv preprint arXiv:2312.00785*, 2023.
- Rob Baxter, Thomas Nind, James Sutherland, Gordon McAllister, Douglas Hardy, Ally Hume, Ruairidh MacLeod, Jacqueline Caldwell, Susan Krueger, Leandro Tramma, et al. The scottish medical imaging archive: 57.3 million radiology studies linked to their medical records. *Radiology: Artificial Intelligence*, pp. e220266, 2023.
- Patrick Bilic, Patrick Ferdinand Christ, Eugene Vorontsov, Grzegorz Chlebus, Hao Chen, Qi Dou, Chi-Wing Fu, Xiao Han, Pheng-Ann Heng, Jürgen Hesser, et al. The liver tumor segmentation benchmark (lits). *arXiv preprint arXiv:1901.04056*, 2019.
- Victor Ion Butoi, Jose Javier Gonzalez Ortiz, Tianyu Ma, Mert R Sabuncu, John Guttag, and Adrian V Dalca. Universeg: Universal medical image segmentation. *arXiv preprint arXiv:2304.06131*, 2023.
- Kai Cao, Yingda Xia, Jiawen Yao, Xu Han, Lukas Lambert, Tingting Zhang, Wei Tang, Gang Jin, Hui Jiang, Xu Fang, et al. Large-scale pancreatic cancer detection via non-contrast ct and deep learning. *Nature Medicine*, pp. 1–11, 2023.
- Krishna Chaitanya, Ertunc Erdil, Neerav Karani, and Ender Konukoglu. Contrastive learning of global and local features for medical image segmentation with limited annotations. *arXiv preprint arXiv:2006.10511*, 2020.
- Liang Chen, Paul Bentley, Kensaku Mori, Kazunari Misawa, Michitaka Fujiwara, and Daniel Rueckert. Self-supervised learning for medical image analysis using image context restoration. *Medical image analysis*, 58:101539, 2019a.
- Sihong Chen, Kai Ma, and Yefeng Zheng. Med3d: Transfer learning for 3d medical image analysis. *arXiv preprint arXiv:1904.00625*, 2019b.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020a.
- Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big self-supervised models are strong semi-supervised learners. *Advances in neural information processing systems*, 33:22243–22255, 2020b.
- Errol Colak, Hui-Ming Lin, Robyn Ball, Melissa Davis, Adam Flanders, Sabeena Jalal, Kirti Magudia, Brett Marinelli, Savvas Nicolaou, Luciano Prevedello, Jeff Rudie, George Shih, Maryam Vazirabad, and John Mongan. Rsn2023 abdominal trauma detection, 2023. URL <https://kaggle.com/competitions/rsna-2023-abdominal-trauma-detection>.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255. IEEE, 2009.

- Yang Deng, Ce Wang, Yuan Hui, Qian Li, Jun Li, Shiwei Luo, Mengke Sun, Quan Quan, Shuxin Yang, You Hao, et al. Ctspine1k: A large-scale dataset for spinal vertebrae segmentation in computed tomography. *arXiv preprint arXiv:2105.14711*, 2021.
- Matthias Eisenmann, Annika Reinke, Vivienn Weru, Minu D Tizabi, Fabian Isensee, Tim J Adler, Sharib Ali, Vincent Andrearczyk, Marc Aubreville, Ujjwal Baid, et al. Why is the winner the best? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19955–19966, 2023.
- Sergios Gatidis, Tobias Hepp, Marcel Früh, Christian La Fougère, Konstantin Nikolaou, Christina Pfannenbergl, Bernhard Schölkopf, Thomas Küstner, Clemens Cyran, and Daniel Rubin. A whole-body fdg-pet/ct dataset with manually annotated tumor lesions. *Scientific Data*, 9(1):601, 2022.
- Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- Hao Guan and Mingxia Liu. Domain adaptation for medical image analysis: a survey. *IEEE Transactions on Biomedical Engineering*, 69(3):1173–1185, 2021.
- Fatemeh Haghighi, Mohammad Reza Hosseinzadeh Taher, Zongwei Zhou, Michael B Gotway, and Jianming Liang. Learning semantics-enriched representation via self-discovery, self-classification, and self-restoration. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 137–147. Springer, 2020. URL <https://github.com/fhaghighi/SemanticGenesis>.
- Fatemeh Haghighi, Mohammad Reza Hosseinzadeh Taher, Zongwei Zhou, Michael B Gotway, and Jianming Liang. Transferable visual words: Exploiting the semantics of anatomical patterns for self-supervised learning. *IEEE Transactions on Medical Imaging*, 2021. URL <https://github.com/fhaghighi/SemanticGenesis>.
- Ali Hatamizadeh, Vishwesh Nath, Yucheng Tang, Dong Yang, Holger R Roth, and Daguang Xu. Swin unetr: Swin transformers for semantic segmentation of brain tumors in mri images. In *International MICCAI Brainlesion Workshop*, pp. 272–284. Springer, 2021.
- Kaiming He, Ross Girshick, and Piotr Dollár. Rethinking imagenet pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4918–4927, 2019.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9729–9738, 2020.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16000–16009, 2022.
- Yuting He, Guanyu Yang, Jian Yang, Rongjun Ge, Youyong Kong, Xiaomei Zhu, Shaobo Zhang, Pengfei Shao, Huazhong Shu, Jean-Louis Dillenseger, et al. Meta grayscale adaptive network for 3d integrated renal structures segmentation. *Medical image analysis*, 71:102055, 2021.
- Nicholas Heller, Sean McSweeney, Matthew Thomas Peterson, Sarah Peterson, Jack Rickman, Bethany Stai, Resha Tejpal, Makinna Oestreich, Paul Blake, Joel Rosenberg, et al. An international challenge to use artificial intelligence to define the state-of-the-art in kidney and kidney tumor segmentation in ct imaging., 2020.
- Ziyan Huang, Haoyu Wang, Zhongying Deng, Jin Ye, Yanzhou Su, Hui Sun, Junjun He, Yun Gu, Lixu Gu, Shaoting Zhang, et al. Stu-net: Scalable and transferable medical image segmentation models empowered by large-scale supervised pre-training. *arXiv preprint arXiv:2304.06716*, 2023.
- Minyoung Huh, Pulkit Agrawal, and Alexei A Efros. What makes imagenet good for transfer learning? *arXiv preprint arXiv:1608.08614*, 2016.

- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089*, 2022.
- Fabian Isensee, Paul F Jaeger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein. nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature Methods*, 18(2):203–211, 2021.
- Alexander Jaus, Constantin Seibold, Kelsey Hermann, Alexandra Walter, Kristina Giske, Johannes Haubold, Jens Kleesiek, and Rainer Stiefelhagen. Towards unifying anatomy segmentation: automated generation of a full-body ct dataset via knowledge aggregation and anatomical guidelines. *arXiv preprint arXiv:2307.13375*, 2023.
- Yuanfeng Ji, Haotian Bai, Jie Yang, Chongjian Ge, Ye Zhu, Ruimao Zhang, Zhen Li, Lingyan Zhang, Wanling Ma, Xiang Wan, et al. Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. *arXiv preprint arXiv:2206.08023*, 2022.
- Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):4037–4058, 2020.
- 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.
- Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pp. 491–507. Springer, 2020.
- Siddharth Krishna Kumar. On weight initialization in deep neural networks. *arXiv preprint arXiv:1704.08863*, 2017.
- Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *International Journal of Computer Vision*, 128(7):1956–1981, 2020.
- Bennett Landman, Zhoubing Xu, J Igelsias, Martin Styner, T Langerak, and Arno Klein. Miccai multi-atlas labeling beyond the cranial vault—workshop and challenge. In *Proc. MICCAI Multi-Atlas Labeling Beyond Cranial Vault—Workshop Challenge*, volume 5, pp. 12, 2015.
- Bowen Li, Zongwei Zhou, Alan Yuille, Max Allan, and Jonathan McLeod. Ultra-transunet: ultrasound segmentation framework with spatial-temporal context feature fusion. In *Medical Imaging 2024: Ultrasonic Imaging and Tomography*, volume 12932, pp. 8–15. SPIE, 2024.
- Jie Liu, Yixiao Zhang, Jie-Neng Chen, Junfei Xiao, Yongyi Lu, Bennett A Landman, Yixuan Yuan, Alan Yuille, Yucheng Tang, and Zongwei Zhou. Clip-driven universal model for organ segmentation and tumor detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 21152–21164, 2023. URL <https://github.com/ljwztc/CLIP-Driven-Universal-Model>.
- Xiangde Luo, Wenjun Liao, Jianghong Xiao, Tao Song, Xiaofan Zhang, Kang Li, Guotai Wang, and Shaoting Zhang. Word: Revisiting organs segmentation in the whole abdominal region. *arXiv preprint arXiv:2111.02403*, 2021.
- DongAo Ma, Jiaxuan Pang, Michael B Gotway, and Jianming Liang. Foundation ark: Accruing and reusing knowledge for superior and robust performance. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 651–662. Springer, 2023a.
- Jun Ma and Bo Wang. Segment anything in medical images. *arXiv preprint arXiv:2304.12306*, 2023a.
- Jun Ma and Bo Wang. Towards foundation models of biological image segmentation. *Nature Methods*, 20(7):953–955, 2023b.

- Jun Ma, Yao Zhang, Song Gu, Cheng Zhu, Cheng Ge, Yichi Zhang, Xingle An, Congcong Wang, Qiyuan Wang, Xin Liu, et al. Abdomenct-1k: Is abdominal organ segmentation a solved problem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- Jun Ma, Yao Zhang, Song Gu, Xingle An, Zhihe Wang, Cheng Ge, Congcong Wang, Fan Zhang, Yu Wang, Yinan Xu, et al. Fast and low-gpu-memory abdomen ct organ segmentation: the flare challenge. *Medical Image Analysis*, 82:102616, 2022.
- Jun Ma, Yao Zhang, Song Gu, Cheng Ge, Shihao Ma, Adamo Young, Cheng Zhu, Kangkang Meng, Xin Yang, Ziyang Huang, et al. Unleashing the strengths of unlabeled data in pan-cancer abdominal organ quantification: the flare22 challenge. *arXiv preprint arXiv:2308.05862*, 2023b.
- Zhiyu Ma, Chen Li, Tianming Du, Le Zhang, Dechao Tang, Deguo Ma, Shanchuan Huang, Yan Liu, Yihao Sun, Zhihao Chen, et al. Aatct-ids: A benchmark abdominal adipose tissue ct image dataset for image denoising, semantic segmentation, and radiomics evaluation. *arXiv preprint arXiv:2308.08172*, 2023c.
- Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 181–196, 2018.
- Mojtaba Masoudi, Hamid-Reza Pourreza, Mahdi Saadatmand-Tarzjan, Noushin Eftekhari, Fateme Shafiee Zargar, and Masoud Pezeshki Rad. A new dataset of computed-tomography angiography images for computer-aided detection of pulmonary embolism. *Scientific data*, 5(1): 1–9, 2018.
- Xueyan Mei, Zelong Liu, Philip M Robson, Brett Marinelli, Mingqian Huang, Amish Doshi, Adam Jacobi, Chendi Cao, Katherine E Link, Thomas Yang, et al. Radimagenet: an open radiologic deep learning research dataset for effective transfer learning. *Radiology: Artificial Intelligence*, 4(5): e210315, 2022.
- Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein Abad, Harlan M Krumholz, Jure Leskovec, Eric J Topol, and Pranav Rajpurkar. Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265, 2023.
- Andriy Myronenko. 3d mri brain tumor segmentation using autoencoder regularization. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 4th International Workshop, BrainLes 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Revised Selected Papers, Part II 4*, pp. 311–320. Springer, 2019.
- Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pp. 69–84. Springer, 2016.
- Irene Papanicolas, Liana R Woskie, and Ashish K Jha. Health care spending in the united states and other high-income countries. *Jama*, 319(10):1024–1039, 2018.
- S Park, LC Chu, EK Fishman, AL Yuille, B Vogelstein, KW Kinzler, KM Horton, RH Hruban, ES Zinreich, D Fadaei Fouladi, et al. Annotated normal ct data of the abdomen for deep learning: Challenges and strategies for implementation. *Diagnostic and interventional imaging*, 101(1): 35–44, 2020.
- Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2536–2544, 2016.
- Chongyu Qu, Tiezheng Zhang, Hualin Qiao, Jie Liu, Yucheng Tang, Alan Yuille, and Zongwei Zhou. Abdomenatlas-8k: Annotating 8,000 abdominal ct volumes for multi-organ segmentation in three weeks. In *Conference on Neural Information Processing Systems*, volume 21, 2023. URL <https://github.com/MrGiovanni/AbdomenAtlas>.

- 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 *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Sucheng Ren, Huiyu Wang, Zhengqi Gao, Shengfeng He, Alan Yuille, Yuyin Zhou, and Cihang Xie. A simple data mixing prior for improving self-supervised learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 14595–14604, 2022.
- Sucheng Ren, Fangyun Wei, Zheng Zhang, and Han Hu. Tinymim: An empirical study of distilling mim pre-trained models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3687–3697, 2023.
- Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. *arXiv preprint arXiv:2104.10972*, 2021.
- Blaine Rister, Darvin Yi, Kaushik Shivakumar, Tomomi Nobashi, and Daniel L Rubin. Ct-org, a new dataset for multiple organ segmentation in computed tomography. *Scientific Data*, 7(1):1–9, 2020.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 234–241. Springer, 2015.
- Holger R Roth, Le Lu, Amal Farag, Hoo-Chang Shin, Jiamin Liu, Evrim B Turkbey, and Ronald M Summers. Deeporgan: Multi-level deep convolutional networks for automated pancreas segmentation. In *International conference on medical image computing and computer-assisted intervention*, pp. 556–564. Springer, 2015.
- Agustina Saenz, Emma Chen, Henrik Marklund, and Pranav Rajpurkar. The maida initiative: establishing a framework for global medical-imaging data sharing. *The Lancet Digital Health*, 6(1):e6–e8, 2024.
- Andrew B SELLERGRN, Christina Chen, Zaid Nabulsi, Yuanzhen Li, Aaron Maschinot, Aaron Sarna, Jenny Huang, Charles Lau, Sreenivasa Raju Kalidindi, Mozziyar Etemadi, et al. Simplified transfer learning for chest radiography models using less data. *Radiology*, 305(2):454–465, 2022.
- Azizi Shekoofeh, Mustafa Basil, Ryan Fiona, Beaver Zachary, Freyberg Jan, Deaton Jonathan, Loh Aaron, Karthikesalingam Alan, Kornblith Simon, Chen Ting, et al. Big self-supervised models advance medical image classification. *arXiv preprint arXiv:2101.05224*, 2021.
- Hoo-Chang Shin, Holger R Roth, Mingchen Gao, Le Lu, Ziyue Xu, Isabella Nogues, Jianhua Yao, Daniel Mollura, and Ronald M Summers. Deep convolutional neural networks for computer-aided detection: CNN architectures, dataset characteristics and transfer learning. *IEEE transactions on medical imaging*, 35(5):1285–1298, 2016.
- Nahian Siddique, Sidike Paheding, Colin P Elkin, and Vijay Devabhaktuni. U-net and its variants for medical image segmentation: A review of theory and applications. *Ieee Access*, 9:82031–82057, 2021.
- Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *arXiv preprint arXiv:2106.10270*, 2021.
- Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, pp. 843–852, 2017.
- Nima Tajbakhsh, Jae Y Shin, Suryakanth R Gurudu, R Todd Hurst, Christopher B Kendall, Michael B Gotway, and Jianming Liang. Convolutional neural networks for medical image analysis: Full training or fine tuning? *IEEE transactions on medical imaging*, 35(5):1299–1312, 2016.

- Yucheng Tang, Dong Yang, Wenqi Li, Holger R Roth, Bennett Landman, Daguang Xu, Vishwesh Nath, and Ali Hatamizadeh. Self-supervised pre-training of swin transformers for 3d medical image analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20730–20740, 2022.
- Xing Tao, Yuexiang Li, Wenhui Zhou, Kai Ma, and Yefeng Zheng. Revisiting rubik’s cube: Self-supervised learning with volume-wise transformation for 3d medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 238–248. Springer, 2020.
- National Lung Screening Trial Research Team. The national lung screening trial: overview and study design. *Radiology*, 258(1):243–253, 2011.
- Atharva Tendle and Mohammad Rashedul Hasan. A study of the generalizability of self-supervised representations. *Machine Learning with Applications*, 6:100124, 2021.
- Jeya Maria Jose Valanarasu, Yucheng Tang, Dong Yang, Ziyue Xu, Can Zhao, Wenqi Li, Vishal M Patel, Bennett Landman, Daguang Xu, Yufan He, et al. Disruptive autoencoders: Leveraging low-level features for 3d medical image pre-training. *arXiv preprint arXiv:2307.16896*, 2023.
- Vanya V Valindria, Nick Pawlowski, Martin Rajchl, Ioannis Lavdas, Eric O Aboagye, Andrea G Rockall, Daniel Rueckert, and Ben Glocker. Multi-modal learning from unpaired images: Application to multi-organ segmentation in ct and mri. In *2018 IEEE winter conference on applications of computer vision (WACV)*, pp. 547–556. IEEE, 2018.
- Yan Wang, Yuyin Zhou, Wei Shen, Seyoun Park, Elliot K Fishman, and Alan L Yuille. Abdominal multi-organ segmentation with organ-attention networks and statistical fusion. *Medical image analysis*, 55:88–102, 2019.
- Ziyi Wang, Xumin Yu, Yongming Rao, Jie Zhou, and Jiwen Lu. P2p: Tuning pre-trained image models for point cloud analysis with point-to-pixel prompting. *Advances in neural information processing systems*, 35:14388–14402, 2022.
- Jakob Wasserthal, Manfred Meyer, Hanns-Christian Breit, Joshy Cyriac, Shan Yang, and Martin Segeroth. Totalsegmentator: robust segmentation of 104 anatomical structures in ct images. *arXiv preprint arXiv:2208.05868*, 2022.
- Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14668–14678, 2022.
- Yingda Xia, Qihang Yu, Linda Chu, Satomi Kawamoto, Seyoun Park, Fengze Liu, Jieneng Chen, Zhuotun Zhu, Bowen Li, Zongwei Zhou, et al. The felix project: Deep networks to detect pancreatic neoplasms. *medRxiv*, 2022.
- Junfei Xiao, Yutong Bai, Alan Yuille, and Zongwei Zhou. Delving into masked autoencoders for multi-label thorax disease classification. *IEEE Winter Conference on Applications of Computer Vision*, 2022. URL https://github.com/lambert-x/medical_mae.
- Yutong Xie, Jianpeng Zhang, Zehui Liao, Yong Xia, and Chunhua Shen. Pgl: Prior-guided local self-supervised learning for 3d medical image segmentation. *arXiv preprint arXiv:2011.12640*, 2020.
- Yutong Xie, Jianpeng Zhang, Yong Xia, and Qi Wu. Unimiss: Universal medical self-supervised learning via breaking dimensionality barrier. In *European Conference on Computer Vision*, pp. 558–575. Springer, 2022.
- Xingyi Yang, Xuehai He, Yuxiao Liang, Yue Yang, Shanghang Zhang, and Pengtao Xie. Transfer learning or self-supervised learning? a tale of two pretraining paradigms. *arXiv preprint arXiv:2007.04234*, 2020.
- Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pp. 3320–3328, 2014.

- Chenyu You, Ruihan Zhao, Fenglin Liu, Siyuan Dong, Sandeep Chinchali, Ufuk Topcu, Lawrence Staib, and James Duncan. Class-aware adversarial transformers for medical image segmentation. *Advances in Neural Information Processing Systems*, 35:29582–29596, 2022.
- Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3712–3722, 2018.
- Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12104–12113, 2022.
- Jianpeng Zhang, Yutong Xie, Yong Xia, and Chunhua Shen. Dodnet: Learning to segment multi-organ and tumors from multiple partially labeled datasets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1195–1204, 2021.
- Shaoting Zhang and Dimitris Metaxas. On the challenges and perspectives of foundation models for medical image analysis. *arXiv preprint arXiv:2306.05705*, 2023.
- Ziheng Zhao, Yao Zhang, Chaoyi Wu, Xiaoman Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. One model to rule them all: Towards universal segmentation for medical images with text prompts. *arXiv preprint arXiv:2312.17183*, 2023.
- Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021a.
- Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022a.
- Zongwei Zhou. *Towards Annotation-Efficient Deep Learning for Computer-Aided Diagnosis*. PhD thesis, Arizona State University, 2021. URL <https://github.com/MrGiovanni/Dissertation>.
- Zongwei Zhou, Jae Shin, Lei Zhang, Suryakanth Gurudu, Michael Gotway, and Jianming Liang. Fine-tuning convolutional neural networks for biomedical image analysis: actively and incrementally. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7340–7351, 2017. URL <https://github.com/MrGiovanni/Active-Learning>.
- Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: A nested u-net architecture for medical image segmentation. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, pp. 3–11. Springer, 2018. URL <https://github.com/MrGiovanni/UNetPlusPlus>.
- Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: Re-designing skip connections to exploit multiscale features in image segmentation. *IEEE Transactions on Medical Imaging*, 39(6):1856–1867, 2019a. URL <https://github.com/MrGiovanni/UNetPlusPlus>.
- Zongwei Zhou, Vatsal Sodha, Md Mahfuzur Rahman Siddiquee, Ruibin Feng, Nima Tajbakhsh, Michael B Gotway, and Jianming Liang. Models genesis: Generic autodidactic models for 3d medical image analysis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 384–393. Springer, 2019b. URL <https://github.com/MrGiovanni/ModelsGenesis>.
- Zongwei Zhou, Vatsal Sodha, Jiaxuan Pang, Michael B Gotway, and Jianming Liang. Models genesis. *Medical Image Analysis*, 67:101840, 2021b. URL <https://github.com/MrGiovanni/ModelsGenesis>.
- Zongwei Zhou, Michael B Gotway, and Jianming Liang. Interpreting medical images. In *Intelligent Systems in Medicine and Health*, pp. 343–371. Springer, 2022b.
- Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin Dogus Cubuk, and Quoc Le. Rethinking pre-training and self-training. *Advances in neural information processing systems*, 33: 3833–3845, 2020.

Appendix

Table of Contents

A An Extensive Dataset: AbdomenAtlas 1.1	18
A.1 Domain Transfer Across Datasets	20
A.2 Uniform Annotation Standards	21
B A Suite of Pre-trained Models: SuPreM	22
B.1 Supervised and Self-supervised Benchmarking	22
C Data, Annotation, and Computational Efficiency	24
C.1 Annotation Efficiency in Fine-tuning	24
C.2 Convergence and Learning Efficiency in Fine-tuning	25
D Enhanced Features for Novel Datasets, Classes, and Tasks	26
D.1 Direct Inference on Three External Datasets	26
D.2 Fine-tuning SuPreM on Fine-grained Tumor Classification	28

A AN EXTENSIVE DATASET: ABDOMENATLAS 1.1

Table 5: **An extensive dataset of 9,262 CT volumes with per-voxel annotations of 25 anatomical structures.** AbdomenAtlas 1.1 marks a breakthrough in data and annotation scales, encompassing 251,323 organ and tumor masks and 2,789,975 annotated images sourced from 88 global hospitals in 19 countries. In 2009, prior to the creation of ImageNet (Deng et al., 2009), AI models struggled with general image representation due to limited data availability, a challenge still prevalent in 3D medical image analysis today. As demonstrated in the table, existing public datasets often suffer from limited, partial, and incomplete annotations, and exhibit biases towards certain populations, medical centers, and countries. Our dataset not only addresses these shortcomings but has also been refined to eliminate redundancy, detailing the count of unique CT volumes sourced from each existing dataset incorporated into ours. AbdomenAtlas 1.1 offers a diverse and extensive range of annotated data, thus marking a significant advancement in the field.

dataset (year) [source]	# of organ	# of volume	# of center	source countries	license
1. Pancreas-CT (2015) [link]	1	82	1	US	CC BY 3.0
2. CHAOS (2018) [link]	4	40	1	TR	CC BY-SA 4.0
3. CT-ORG (2020) [link]	5	140	8	DE, NL, CA, FR, IL, US	CC BY 3.0
4. BTCV (2015) [link]	12	50	1	US	CC BY 4.0
5. AMOS22 (2022) [link]	15	500	2	CN	CC BY-NC-SA
6. WORD (2021) [link]	16	150	1	CN	GNU GPL 3.0
7-12. MSD CT Tasks (2021) [link]	9	947	1	US	CC BY-SA 4.0
13. LiTS (2019) [link]	1	201	7	DE, NL, CA, FR, IL	CC BY-SA 4.0
14. AbdomenCT-1K (2021) [link]	4	1,050	12	DE, NL, CA, FR, IL, US, CN	CC BY-NC-SA
15. KiTS (2020) [link]	1	300	1	US	CC BY-NC-SA 4.0
16. FLARE'23 (2022) [link]	13	4,000	30	-	CC BY-NC-ND 4.0
17. Trauma Det. (2023) [link]	0	4,711	23	CL, DE, ES, TR, AUS, TH, MA, MT, CA, IE, BR, BA	-
18. FUMPE (2018) [link]	1	35	1	IR	CC BY 4.0
19. KiPA22 (2021) [link]	4	100	1	CN	CC BY-NC-ND 3.0
20. AATTCT-IDS (2023c) [link]	0	300	1	CN	-
21. CTSpine1K (2021) [link]	26	1,005	-	-	CC BY 4.0
22. AutoPET (2022) [link]	0	1,014	2	DE	TCIA Restricted
23. TotalSegmentator (2022) [link]	104	1,204	1	CH	CC BY 4.0
24. AbdomenAtlas 1.0 (2023) [link]	9	5,195	26	US, DE, NL, CA, FR, IL, CN, TR, MT, IE, BR, BA, AUS, TH, CA, TR, CH, CL, ES, MA, US, DE, NL, FR, IL, CN	CC BY-NC-SA 4.0
25. AbdomenAtlas 1.1	25	9,262	88	US, DE, NL, CA, FR, IL, CN	CC BY-NC-SA 4.0

US: United States DE: Germany NL: Netherlands CA: Canada FR: France IL: Israel IR: Iran
 CN: China TR: Turkey CH: Switzerland AUS: Australia TH: Thailand ES: Spain CL: Chile
 MA: Morocco MT: Malta IE: Ireland BR: Brazil BA: Bosnia and Herzegovina



Figure 4: **Evolution from a combination of public data to AbdomenAtlas 1.1.** AbdomenAtlas 1.1 is not a simple combination of existing datasets. The 9,262 CT volumes in the combination of public datasets only contain a total of **39K** annotated organ masks while our AbdomenAtlas 1.1 provides over **251,323** annotated organ/tumor masks for these CT volumes, substantially increasing the number of masks by **6.4** times. Creating 251,323 high-quality organ/tumor masks for 9,262 CT volumes requires extensive medical knowledge and annotation costs (much more difficult than annotating natural images). Based on our experience and those reported in [Park et al. \(2020\)](#), trained radiologists annotate abdominal organs at a rate of 30–60 minutes per organ per CT volume. This translates to **247K** human hours for completing AbdomenAtlas 1.1. We employed a highly efficient annotation method, combining AI with the expertise of ten radiologists using active learning (details in Appendix A.2), to overcome this challenge and produce the largest annotated dataset to date.

A.1 DOMAIN TRANSFER ACROSS DATASETS

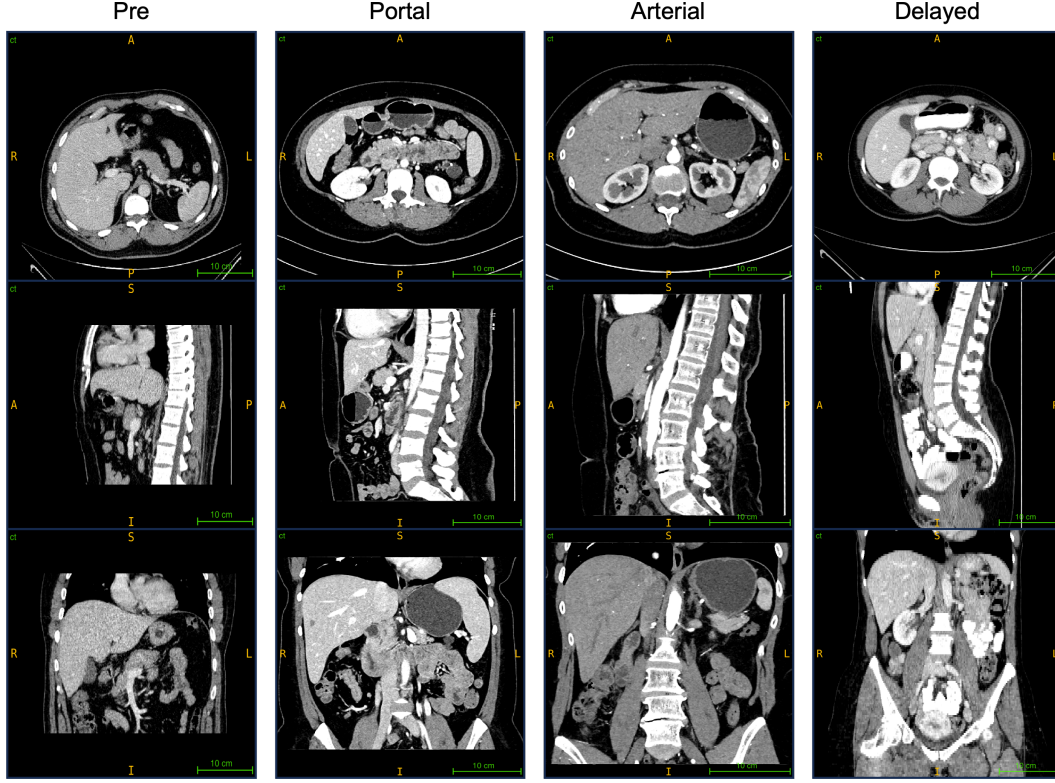


Figure 5: **Domain gaps.** Examples of CT volumes from different domains (e.g., hospitals and countries) illustrate the variability in images. AbdomenAtlas 1.1 are created by a large variety of CT scanners, imaging protocols, and acquired from numerous hospitals worldwide (Table 1). We note that substantial differences in CT volumes occur in image quality and technical display, originating from different acquisition parameters, reconstruction kernels, and contrast enhancements.

Table 3 shows that SuPreM is pretty robust because our AbdomenAtlas 1.1 covers a variety of domains (i.e., 88 hospitals with different scanners and protocols), as shown in Figure 5; models pre-trained on this dataset are expected to be generalizable for novel domains. Therefore, domain transfer becomes less important if the model is pre-trained on large and diverse datasets, elaborating on the two points below.

1. The domain transfer problem could be solved by methodology innovation, and also by training AI models on enormous datasets. This point has been more clear recently demonstrated by large language models (ChatGPT) and vision foundation models (SAM), which show incredible performance in the “novel domain”. However, this achievement may not be directly attributed to method-driven solutions for domain transfer, but simply because the AI might have been trained on similar sentences or images. This was also pointed out by Yann Lecun—*beware of testing on the training set*—in response to the incredible results achieved by ChatGPT.
2. In some sense, our paper explores dataset-driven solutions for domain transfer. The robust performance of our models when direct inference on multiple domains could also be attributed to our large-scale, fully-annotated medical dataset—as one of our major contributions. The release of AbdomenAtlas 1.1 can foster AI models that are more robust than the majority of existing models that are only trained on a few hundred CT volumes from limited domains. In addition, existing domain transfer methods could also be supplemented with direct inference and fine-tuning to further improve AI performance.

A.2 UNIFORM ANNOTATION STANDARDS

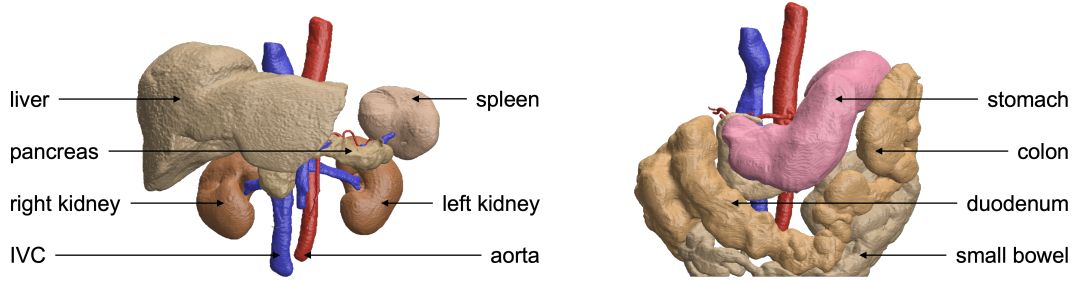


Figure 6: **Automated organ annotations.** Our annotation pipeline involved an interactive segmentation approach, a synergy of AI algorithms and human expertise, which promises to improve efficiency while upholding high-quality annotations. *Four senior radiologists* revised the annotations predicted by our AI models, and in turn, the AI models improved their predictions by learning from these revised annotations. This interactive procedure continued to enhance the quality of annotations until no major revision was needed. Subsequently, *Six junior radiologists* examine the final visualizations for accuracy (examples of the rendered images are shown above). The junior radiologists were responsible for reviewing the correctness of the annotations and marking the patient ID for any major discrepancies. Such cases are then reviewed by senior radiologists. Our uniform annotation standards, largely overlapping with those in [Ma et al. \(2023b\)](#), require trained radiologists to spend approximately 30–60 minutes annotating each organ in a three-dimensional CT volume.

Automated (pseudo) tumor annotations. We have established uniform annotation standards for tumors, with both senior and junior radiologists actively refining and adhering to these guidelines.

- **Liver tumors:** Liver tumors include primary tumor lesions and metastases in the liver. Annotations should encompass the entire tumor, including any invasive parts, necrosis, hemorrhage, fibrous scars, and calcifications. Healthy areas or unrelated lesions are not included.
- **Kidney tumors:** Kidney tumors include both benign and malignant tumor lesions growing in the kidneys. The entire tumor and its invasive parts to surrounding areas, plus internal changes like necrosis and calcification, should be annotated. Exclude healthy structures.
- **Pancreatic tumors:** Pancreatic tumors include all benign and malignant tumor lesions growing in the pancreas. Annotations cover the whole tumor and its invasive growth into adjacent areas, including changes like cysts, necrosis, and calcification. Exclude healthy structures.
- **Colon tumors:** Colon tumors include all benign and malignant tumor lesions developing from the colon wall. The entire tumor and its invasion into nearby structures, along with internal changes like necrosis, should be annotated, excluding healthy areas.
- **Hepatic vessel tumors:** Hepatic vessel tumors include all primary tumor lesions developing from the intrahepatic vessel wall and tumor thrombus in intrahepatic vessels. Annotations should include the tumor within the vessels, excluding external parts and unrelated lesions.

Overall, AbdomenAtlas 1.1 offers 51.8K pseudo tumor masks visually inspected by radiologists, though without biopsy confirmation. While these masks lack pathological validation, we anticipate they will serve as a valuable foundation for expanding precise tumor annotations in future research.

B A SUITE OF PRE-TRAINED MODELS: SUPReM

B.1 SUPERVISED AND SELF-SUPERVISED BENCHMARKING

B.1.1 BACKGROUND & STATEMENT

The goal of Table 2 and Appendix Table 6 is to provide a practical benchmark for the transfer learning ability of readily available pre-trained models. Our intent is not to compare the specific pre-training methodologies of each model for two primary reasons.

1. The majority of researchers tend to fine-tune pre-existing models rather than retrain them from scratch due to convenience and accessibility.
2. Reproducing these models would require specialized hyper-parameter tuning and varied computational resources. For example, models like Swin UNETR (Tang et al., 2022) were pre-trained using large-scale GPU clusters at NVIDIA, making them challenging for us to faithfully retrain.

Considering both practical user scenarios and computational constraints, we decided to directly use their released models and fine-tune them with consistent settings on the same datasets.

However, using existing pre-trained models can inevitably lead to certain problems. For example, the U-Net family has seen numerous variations over the years (Zhou et al., 2018; 2019a; Siddique et al., 2021; Li et al., 2024). Pre-trained models released before 2021 typically employed a basic version of U-Net (Zhou et al., 2019b; Chen et al., 2019b). On the other hand, our U-Net benefits from a more advanced code base, thanks to the MONAI platform at NVIDIA, which includes enhanced architectures and advanced training optimization strategies. Consequently, our U-Net, even trained from scratch, is capable of surpassing the performance of these older baseline models.

B.1.2 IMPLEMENTATION DETAILS OF PRE-TRAINING

For benchmark purposes (Tables 2, 4 and Figures 1, 2b, 3), we pre-trained U-Net, Swin UNETR, and SegResNet on 2,100 fully annotated CT volumes with 25 anatomical structures and pseudo annotations of seven tumors. The best model was selected based on the largest average DSC over the 32 classes on 310 CT volumes as the validation set. We randomly crop sub-volumes, sized $96 \times 96 \times 96$ voxels, from the original CT volumes. Our SuPreM is pre-trained with AdamW using $\beta_1 = 0.9$ and $\beta_2 = 0.999$ with a batch size of 2 per GPU and a cosine learning rate schedule with a warm-up for the first 100 epochs. We start with an initial learning rate of $1e^{-4}$ and a decay of $1e^{-5}$. The pre-training has been conducted on four NVIDIA A100 using multi-GPU (4) with distributed data parallel (DDP), implemented in MONAI 0.9.0., with a maximum of 800 epochs. We use the binary cross-entropy and Dice Similarity Coefficient (DSC) losses as the objective function for pre-training.

B.1.3 IMPLEMENTATION DETAILS OF FINE-TUNING

We fine-tune the pre-trained models using TotalSegmentator and the proprietary dataset datasets. During fine-tuning, configurations from pre-training persist, but we adjust the warm-up scheduler to 20 epochs, set a maximum of 200 epochs, and use a single GPU.

Table 6: **Benchmarking all the self-supervised and supervised models.** All the publicly available pre-trained models can be downloaded from our [model repository](#) at Huggingface. We will continue to include more 3D pre-trained models when they are available.

	name	backbone	params	pre-trained data	performance [†]
self-supervised	Models Genesis (Zhou et al., 2019b)	U-Net	19.08M	623 CT volumes	90.1
	UniMiSS (Xie et al., 2022)	nnU-Net	61.79M	5,022 CT&MRI volumes	92.9
	NV*	Swin UNETR	62.19M	1,000 CT volumes	93.2
	NV*	Swin UNETR	62.19M	3,000 CT volumes	93.4
	NV (Tang et al., 2022)	Swin UNETR	62.19M	5,050 CT volumes	93.8
	NV*	Swin UNETR	62.19M	5,050 CT volumes	94.2
supervised	NV*	Swin UNETR	62.19M	9,262 CT volumes	94.3
	Med3D (Chen et al., 2019b)	Residual U-Net	85.75M	1,638 CT volumes	91.4
	DoDNet (Zhang et al., 2021)	U-Net	17.29M	920 CT volumes	93.8
	DoDNet*	U-Net	17.29M	920 CT volumes	94.4
	Universal Model (Liu et al., 2023)	Swin UNETR	62.19M	2,100 CT volumes	94.1
	SuPreM*	U-Net	19.08M	2,100 CT volumes	95.4
	SuPreM*	Swin UNETR	62.19M	2,100 CT volumes	94.6
	SuPreM*	SegResNet	4.7M	2,100 CT volumes	94.0

[†] We report the transfer learning performance of muscle segmentation on TotalSegmentator.

The name with a star () denotes it is implemented by us and pre-trained using our AbdomenAtlas 1.1.

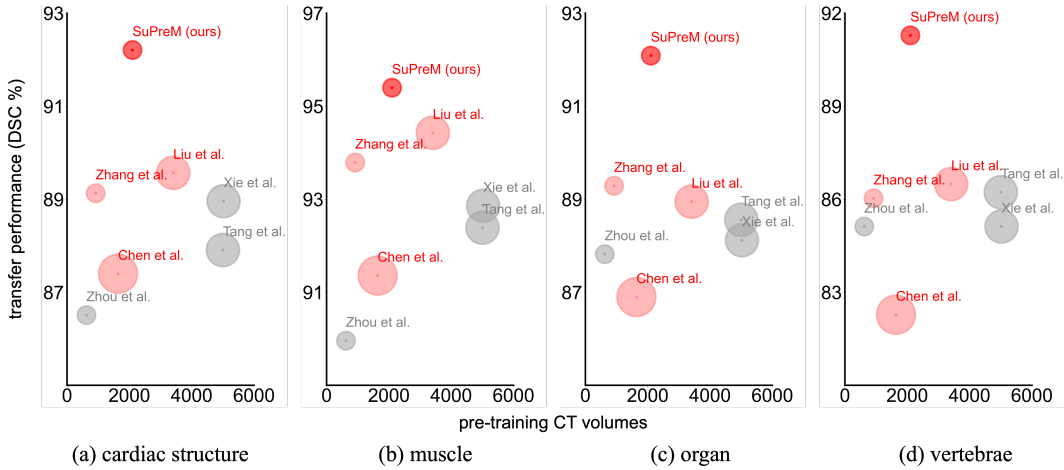


Figure 7: **A comprehensive benchmark on supervised and self-supervised models.** We present the segmentation performance achieved by fine-tuning models using the entire TotalSegmentator training set ($N = 1,081$ annotated CT volumes) as target tasks. A larger circle size denotes a greater number of model parameters. Overall, for target tasks, supervised models (in red) transfer better for pre-training in comparison with self-supervised models (in gray).

C DATA, ANNOTATION, AND COMPUTATIONAL EFFICIENCY

C.1 ANNOTATION EFFICIENCY IN FINE-TUNING

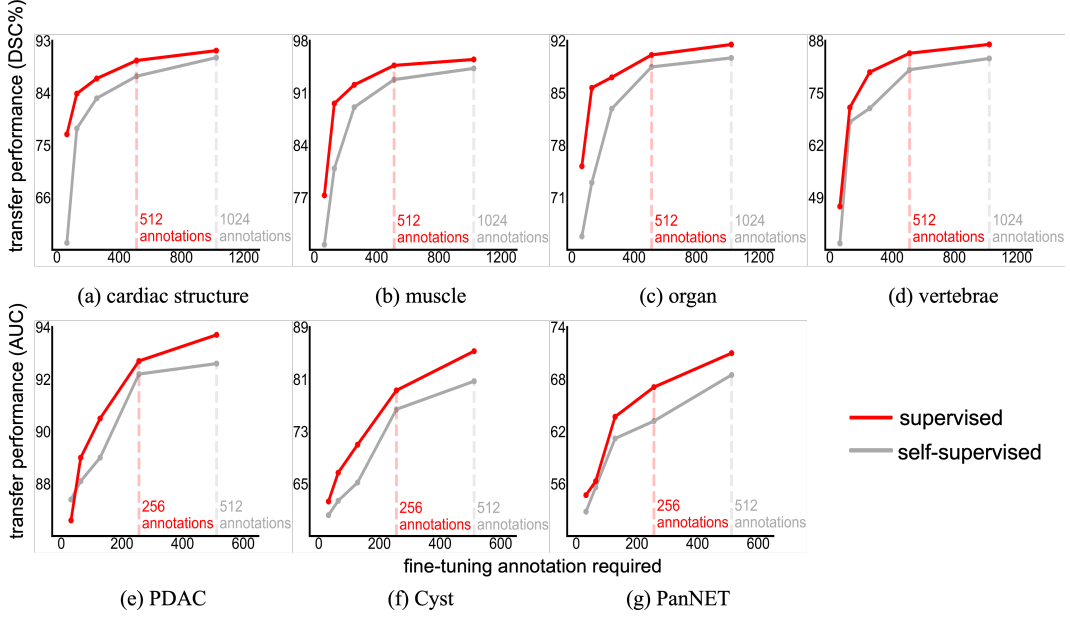


Figure 8: **SuPreM is annotation efficient when transferred to novel class segmentation tasks.** We assess the annotation & learning efficiency by fine-tuning models on different numbers of annotated CT volumes from TotalSegmentator and the proprietary dataset of a total of 66 novel classes. Specifically, TotalSegmentator provides 19 muscles, 15 cardiac structures, 5 organs, and 24 vertebrae; the proprietary dataset offers three sub-types of pancreatic tumors, including pancreatic ductal adenocarcinoma (PDAC), pancreatic cysts, and pancreatic neuroendocrine tumors (PanNET).

C.2 CONVERGENCE AND LEARNING EFFICIENCY IN FINE-TUNING

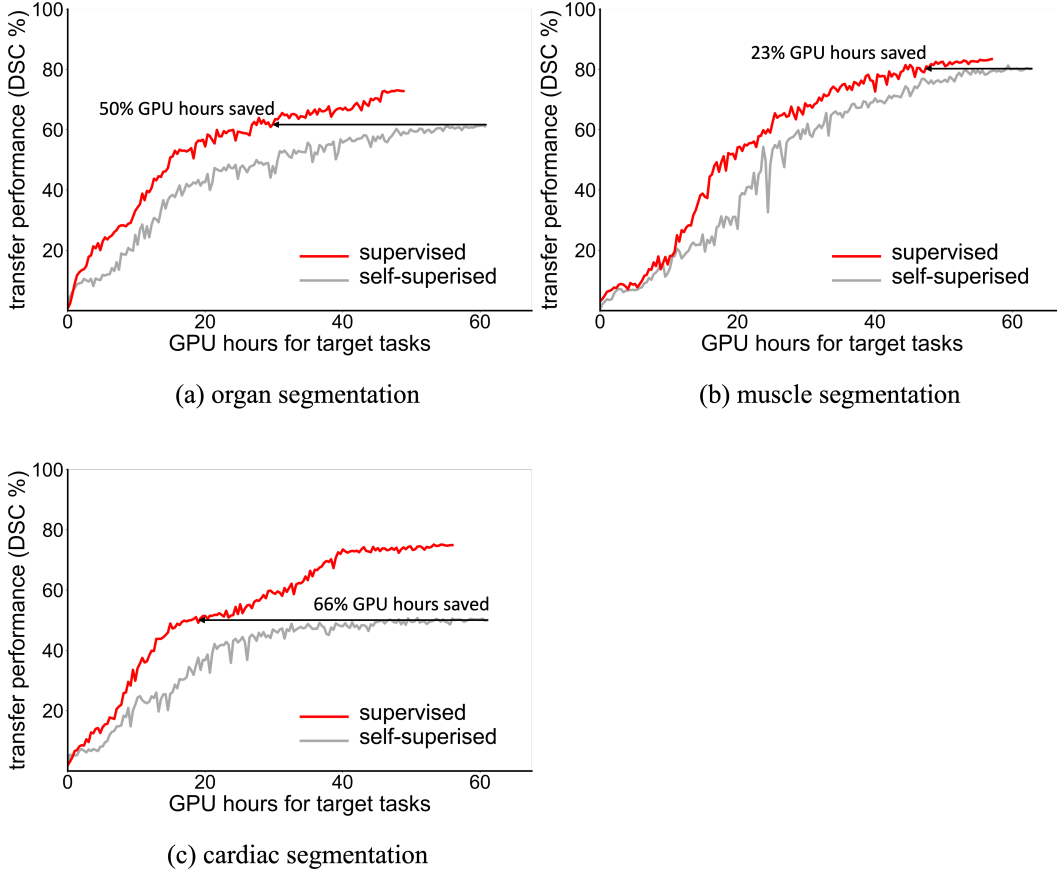


Figure 9: **Convergence & learning efficiency in fine-tuning.** We present the learning curves for fine-tuning supervised and self-supervised models for three target tasks using 10% of the training set. Supervised models achieve markedly better performance and converge faster than self-supervised counterparts by 50%, 23%, and 66% for the tasks of organ, muscle, and cardiac segmentation, respectively.

D ENHANCED FEATURES FOR NOVEL DATASETS, CLASSES, AND TASKS

D.1 DIRECT INFERENCE ON THREE EXTERNAL DATASETS

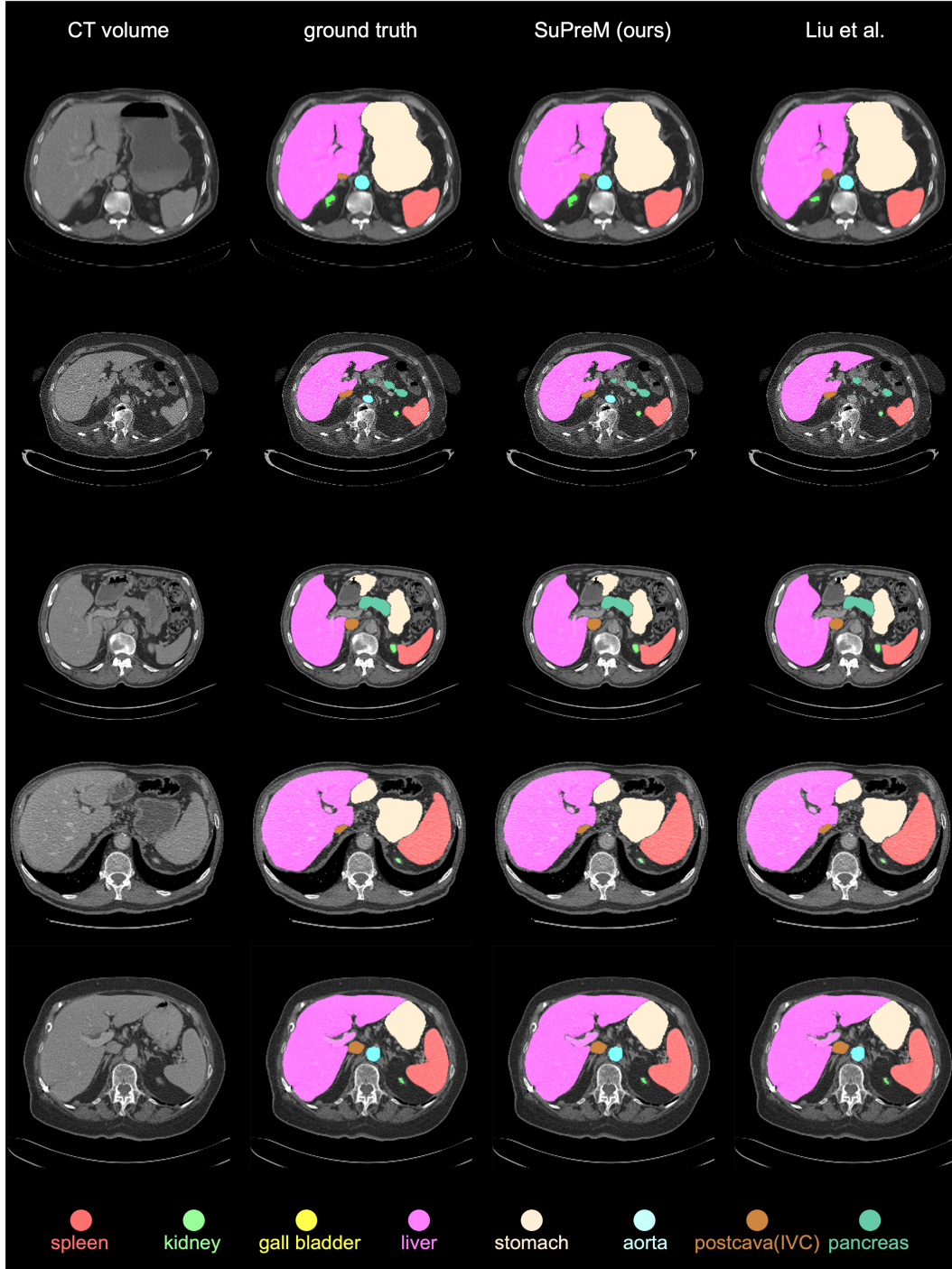


Figure 10: **Direct inference on TotalSegmentator.** We performed direct inference using TotalSegmentator, covering 104 classes. Here only 25 out of the 104 classes were visualized.

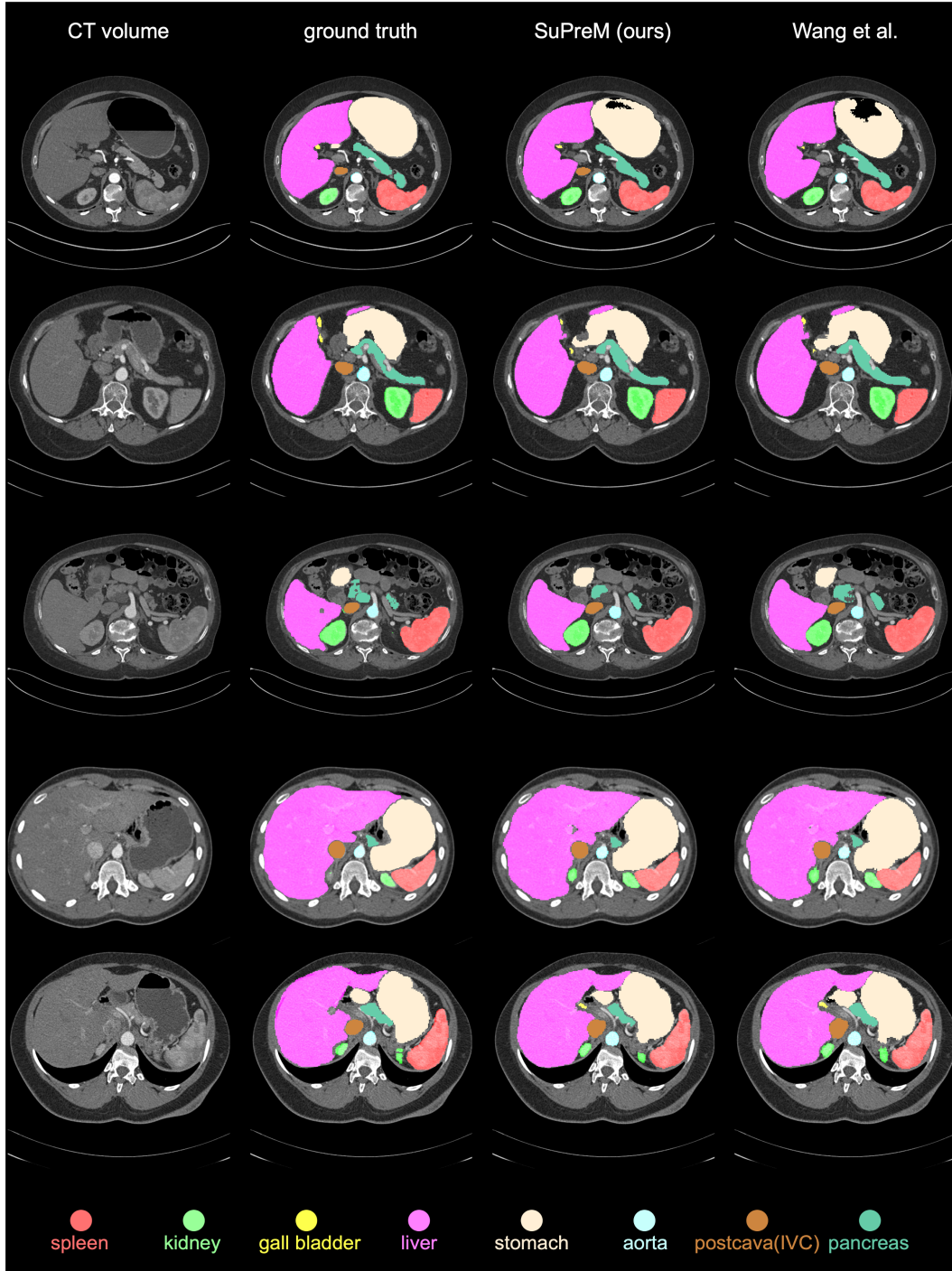


Figure 11: **Direct inference on the proprietary dataset.** We performed direct inference using a proprietary dataset at JHU, which covers 20 organ classes and three sub-types of pancreatic tumors. These include the aorta, adrenal gland, common bile duct, celiac abdominal aorta, colon, duodenum, gallbladder, postcava, left kidney, right kidney, liver, pancreas, pancreatic duct, superior mesenteric artery, small bowel, spleen, stomach, veins, left renal vein and right renal vein. The pancreatic tumor classes are pancreatic ductal adenocarcinoma (PDAC), pancreatic cysts, and pancreatic neuroendocrine tumors (PanNET).

D.2 FINE-TUNING SUPREM ON FINE-GRAINED TUMOR CLASSIFICATION

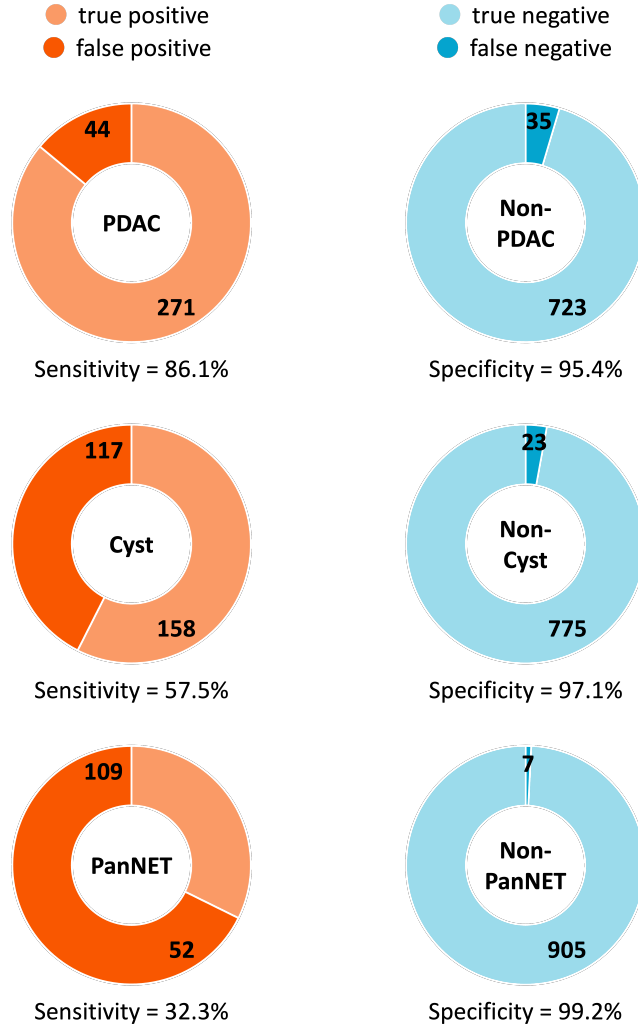


Figure 12: **Fine-grained pancreatic tumor classification.** We would like to stress the challenges in benchmarking tumor segmentation/classification, particularly due to the scarcity of annotations in publicly available datasets (often limited to hundreds of tumors). To overcome this limitation, we employed our proprietary dataset, which comprises 3,577 annotated pancreatic tumors, including detailed sub-types: 1,704 PDACs, 945 Cysts, and 928 PanNETs. The proprietary dataset contains CT scans taken by a variety of vendors, e.g., Philips, Siemens, GE, and Toshiba. This extensive dataset enabled us to thoroughly assess the transfer learning ability of our pre-trained models in tumor-related tasks. Notably, the transfer learning results detailed here demonstrate a sensitivity of 86.1% and specificity of 95.4% for PDAC detection. This performance surpasses the average radiologist’s performance in PDAC identification by 27.6% in sensitivity and 4.4% in specificity, as reported in [Cao et al. \(2023\)](#). This is one of the demonstrations of how our pre-trained models could be deployed for clinical applications.