# STU-Net: Scalable and Transferable Medical Image Segmentation Models Empowered by Large-Scale Supervised Pre-training

Ziyan Huang[1,2*]     Haoyu Wang[1,2*]     Zhongying Deng[2*]     Jin Ye[2*]     Yanzhou Su[2]

Hui Sun[2]     Junjun He[2]     Yun Gu[1]     Lixu Gu[1]     Shaoting Zhang[2]     Yu Qiao[2†]

[1]Shanghai Jiao Tong University
[2]Shanghai AI Laboratory

{ziyanhuang, gulixu}@sjtu.edu.cn
{hejunjun, yejin, zhangshaoting, qiaoyu}@pjlab.org.cn

## Abstract

*Large-scale models pre-trained on large-scale datasets have profoundly advanced the development of deep learning. However, the state-of-the-art models for medical image segmentation are still small-scale, with their parameters only in the tens of millions. Further scaling them up to higher orders of magnitude is rarely explored. An overarching goal of exploring large-scale models is to train them on large-scale medical segmentation datasets for better transfer capacities. In this work, we design a series of Scalable and Transferable U-Net (STU-Net) models, with parameter sizes ranging from 14 million to 1.4 billion. Notably, the 1.4B STU-Net is the largest medical image segmentation model to date. Our STU-Net is based on nnU-Net framework due to its popularity and impressive performance. We first refine the default convolutional blocks in nnU-Net to make them scalable. Then, we empirically evaluate different scaling combinations of network depth and width, discovering that it is optimal to scale model depth and width together. We train our scalable STU-Net models on a large-scale TotalSegmentator dataset and find that increasing model size brings a stronger performance gain. This observation reveals that a large model is promising in medical image segmentation. Furthermore, we evaluate the transferability of our model on 14 downstream datasets for direct inference and 3 datasets for further fine-tuning, covering various modalities and segmentation targets. We observe good performance of our pre-trained model in both direct inference and fine-tuning. The code and pre-trained models are available at https://github.com/Ziyan-Huang/STU-Net.*

*Equal contribution. This work is done when Ziyan Huang is an intern at Shanghai AI Laboratory.
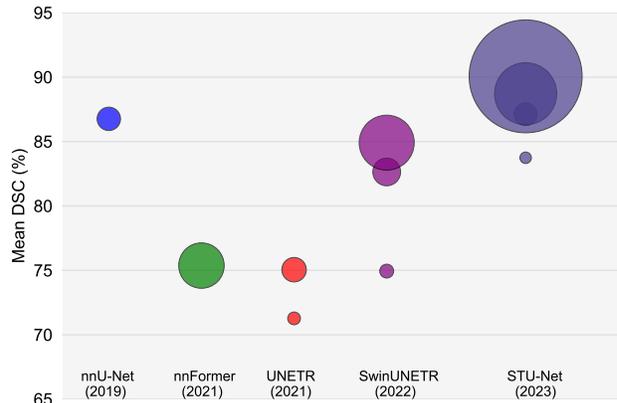
†Corresponding author



Figure 1. Segmentation performance of various models on the TotalSegmentator dataset. The area of each bubble is proportional to the FLOPs (Floating-Point Operations Per Second) of the corresponding model at different scales. Distinct colors represent different models, while multiple bubbles of the same color denote the same model with varying scales. FLOPs calculations are based on input patch sizes of $128 \times 128 \times 128$.

## 1. Introduction

Medical image segmentation, which aims to automatically annotate anatomical structures and lesions in medical images, is an important intermediate step for many downstream clinical tasks, such as medical image registration [39], quantification [24] and image-guided surgery [37]. In recent years, various specific medical image segmentation tasks have been heavily studied and great success has been achieved by many deep learning based models

[12, 13, 33, 44]. However, these models often require careful tuning to adapt to different tasks, limiting their transferability. There is a lack of work investigating a single model to simultaneously handle various medical segmentation tasks, including various input modalities (CT, MRI, PET, etc.) and various segmentation targets, such as organs and tumors.

The key to such a goal can be pre-training large-scale models on large-scale datasets to make the models **transferable**. This solution has been verified in both vision and language fields [3, 5, 7, 20, 21], and its success can be explained as follows: large-scale datasets can provide more knowledge for model training, while large-scale models with more parameters can better take advantage of this knowledge. From the perspective of datasets, some public large-scale medical image segmentation datasets are emerging. For example, AbdomenCT-1K [29], BraTS21 [2], AutoPET [11], and TotalSegmentator [38] contain more than 1000 annotated image from diverse modalities (including CT, PET, MRI) and segmentation targets (abdomen organs, brain/whole-body tumor). Notably, the large-scale TotalSegmentator [38] dataset has 1204 CT images with 104 organs annotated, which is suitable for training a large-scale model.

With the availability of large-scale datasets for pre-training, our primary focus is on designing large-scale models that are **transferable** to a variety of medical segmentation tasks. Nonetheless, large-scale models usually take much more computational costs. This can be even severer when 3-dimensional high-resolution medical images are used for training. Hence, we also hope that the large-scale model can be **scalable** to different sizes to fit different computational budgets.

To sum up, our goal is to design a large-scale medical segmentation model that is **scalable and transferable**. To this end, we propose a series of Scalable and Transferable U-Net [33], termed STU-Net, with parameter sizes ranging from 14 million to 1.4 billion. It is worth noting that the 1.4B model is the largest model in the medical image segmentation field to date. In addition to these different model sizes for scalability, we pre-train them on large-scale datasets in a supervised manner to ensure the models' transfer capacities. Specifically, we build our models based on nnU-Net framework [18] due to its state-of-the-art baseline performance and its wide use by researchers. However, there are two obstacles to developing large-scale models using this framework. 1) For the models' scalability, the basic convolutional blocks in nnU-Net may not be suitable for scaling due to gradient diffusion, and how to increase parameters remains unclear. 2) To evaluate the models' transferability, we usually need to fine-tune the models on other downstream datasets. But the model architectures in nnU-Net cannot easily be used for fine-tuning since they are treated as hyper-parameters and thus task-specific.

To tackle the first obstacle, we refine the basic convolutional blocks of nnU-Net, such as incorporating residual connection [14] to its basic blocks, to facilitate scaling model depth. Then we empirically evaluate different combinations of network depth and width. We discover that it is optimal to scale model depth and width together. Hence, we can obtain the STU-Net-L (2× depth, 2× width) with 450 million parameters and the STU-Net-H (3× depth, 3× width) with 1.4 billion parameters. We do not conduct further scaling due to the prohibitive computational costs and GPU memory consumption on a single A100 GPU.

For the second obstacle, we replace the transpose convolution in nnU-Net with the nearest interpolation followed by a $1 \times 1 \times 1$ convolution for up-sampling, which avoids the use of task-specific kernel/stride options in transpose convolution and makes up-sampling blocks transferable. We further fix hyper-parameters related to model weights, e.g., fixing the number of stages to 6 and using isotropic convolution kernels, so that the model architecture can be the same during pre-training and fine-tuning, thus feasible for transfer to other tasks.

To make our STU-Net transferable, we further adopt the large-scale TotalSegmentator dataset for supervised pre-training. TotalSegmentator is one of the largest datasets in terms of the number of training images and categories, containing 1204 volumetric CT images with 104 well-annotated structures in the whole body. The large-scale STU-Net pre-trained on the large-scale dataset excels in various downstream datasets, covering diverse modalities and segmentation targets, when doing direct inference and fine-tuning. The superior performance shows that our large-scale STU-Net has impressive transferability, benefiting from the large-scale supervised pre-training.

Our contribution can be summarized as:

- We propose scalable STU-Net models which can be scaled to different parameter sizes, with the 1.4B STU-Net as the largest medical image segmentation model to date. Based on these STU-Net models, we discover that model's performance significantly increases with respect to model size when trained on large-scale datasets.

- Our large-scale STU-Net model pre-trained on large-scale TotalSegmentator dataset demonstrates strong transfer capabilities. It can generalize well to a variety of other datasets without additional tuning or adaptation. It also excels on downstream datasets with fine-tuning.

- Our STU-Net is based on the nnU-Net framework and benefits from its task-adaptive designs, so our model can adapt to different tasks with good performance guarantee. Meanwhile, we modify the model architecture in nnU-Net to ensure that model weights trained

on the TotalSegmentator dataset can be easily transferred to downstream tasks.

## 2. Related Works

### 2.1. Medical Image Segmentation Models

Medical image segmentation is dominated by deep learning models, which can be broadly divided into two categories: CNN-based models and Transformer-based models. U-Net [33] is the pioneering CNN model proposed for medical image segmentation. Based on it, residual connection [10, 30], attention module [31] and different feature aggregation strategies [45] are applied for various tasks. Recently, vision transformers [9, 26] with self-attention mechanisms [36], which achieve success in natural image processing, are introduced in medical image segmentation tasks. Particularly, UNETR [13] and SwinUNETR [12] use Vision Transformer and Swin Transformer respectively as encoders to extract features from embedded 3D patches and positional embedding. TransUNet [6] uses transformer blocks as a bottleneck to extract global contexts. nnFormer [44] proposes an interleaved combination of transformer and convolution blocks to extract both local and global features. These medical image segmentation models have only several million parameters, thus not large enough. In addition, these models are not scalable and transferable to fit different computational budgets and handle diverse medical image segmentation tasks simultaneously.

### 2.2. Scaling-up Models

Scaling up models is widely used to boost models' performance in deep learning. The most common way is to scale the depth [14] and width [41] of models. EfficientNet [35] proposes to scale network depth, width, and resolution in a compound way. [20] and [42] present comprehensive studies of empirical scaling laws of transformers in language processing and vision recognition respectively, where the main finding shows that the relationships between computation, data size, model size, and performance fit a power law. Following the scaling laws, GPT-3 [5], which has 175 billion parameters and is pre-trained on about 45 TB text data, has near-human performance in various text processing tasks. In addition, vision Transformers [9] have been scaled up to 22 billion parameters [7] and pre-trained on approximately 4 billion images [42] [34]. With a frozen visual feature extractor, the ViT-22B model achieved a top-1 accuracy of 89.5% on the ImageNet [8] dataset. However, only a few works scale models in medical image segmentation to large scale. Following EfficientNet, [3] scales 2D U-Net for biodegradable bone implant segmentation. [17] explores task-specific scaling strategies for various medical image segmentation tasks. These works have limited scalability and are only evaluated on small datasets. In contrast, our paper successfully scales models to be an order of magnitude larger than previous works and evaluates their transfer capacities on large-scale datasets.

## 3. Methods

We build our models based on the nnU-Net framework which can configure task-specific hyper-parameters automatically and achieve state-of-the-art performance on various tasks. In this section, we first introduce our refinements on it to facilitate scalability and transferability, then present our proposed scaling method, and finally detail our large-scale supervised pre-training strategy for better transferability.

### 3.1. nnU-Net Architecture

We first briefly introduce the default nnU-Net architecture, especially its modules related to our modifications. nnU-Net adopts a symmetric encoder-decoder architecture based on skip connections. Such an architecture contains various resolution stages. Each stage comprises two convolution layers followed by instance normalization and leaky ReLU (denoted as Conv-IN-LeakyReLU). It does not contain residual connections, so simply stacking more layers in each stage may suffer from gradient diffusion, making the whole model hard to optimize. This can limit the depth of nnU-Net, further constraining the scalability of nnU-Net. On the other hand, nnU-Net determines the input patch size and input spacing according to dataset properties. Then, the dataset-sensitive patch size and spacing are used to set the hyper-parameters related to network architecture, like the number of resolution stages, convolution kernels and down-sample/up-sample ratios. Thus, these architecture-related hyper-parameters vary between tasks, leading to different network architectures for different tasks. As such, the models trained on one task cannot be transferred to the other tasks directly, which restricts the evaluation of models' transfer capacities.

### 3.2. nnU-Net Tweaks

The task-specific hyper-parameters in nnU-Net can be divided into model weight-related (e.g., convolution kernel size, number of resolution stages) and weight-unrelated (such as pool kernels, input patch size, and input spacing). We first fix weight-related hyper-parameters to keep the model architecture feasible for transfer to other tasks. Concretely, we keep the number of resolution stages in all tasks to be 6, and use isotropic kernels (3,3,3) for all convolution layers. For weight-unrelated hyper-parameters, we adopt the default setting in nnU-Net to ensure its state-of-the-art performance. We compare our settings with nnU-Net and 3D U-Net [43] in Table 1.

Table 1. Comparison of different hyper-parameters in nnU-Net, 3D U-Net and our method. The up(down)-sample ratios along different axes are (x,y,z). nnU-Net configures task-specific hyper-parameters automatically, which usually brings state-of-the-art performance. But this may make the model less transferable. The 3D-UNet uses a fixed network architecture that can easily transfer weights across different segmentation tasks but its baseline performance is usually worse than nnU-Net. Our STU-Net not only maintains good performance of nnU-Net but also can be transferred easily.

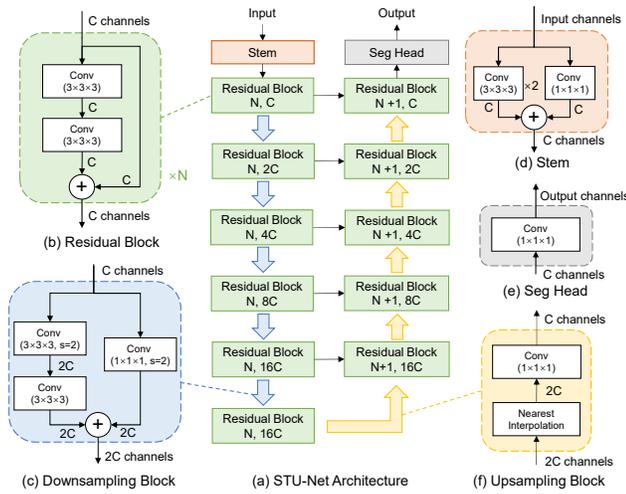| Settings | nnU-Net | 3D U-Net | STU-Net (ours) |
|---|---|---|---|
| number of resolution stages | 4-7 | 5 | 6 |
| convolution kernels | $3\times3\times3$ or $3\times3\times1$ | $3\times3\times3$ | $3\times3\times3$ |
| up(down)-sample ratios | (2,2,2) or (2,2,1) | (2,2,2) | (2,2,2) or (2,2,1) |
| input patch size | task-specific | fixed | task-specific |
| input spacing | task-specific | fixed | task-specific |
| up-sample operation | transpose convolution | transpose convolution | interpolation with (1,1,1) convolution |



Figure 2. Illustration of our STU-Net architecture which is built upon the nnU-Net architecture with several modifications to enhance its scalability and transferability. (a) An overview of the STU-Net architecture. The blue arrows denote downsampling while the yellow ones represent upsampling. (b) Residual blocks to achieve a large-scale model. (c) Downsampling in the first residual block of each encoder stage. (d-e) Stem and segmentation head for channel conversion of input and output. (f) Weight-free interpolation for upsampling, which effectively addresses the issue of weight mismatch across different tasks.

### 3.2.1 Basic Block Tweaks

In both the encoder and decoder of nnU-Net, each stage consists of a basic block. Each basic block comprises two Conv-IN-LeakyReLU layers. When increasing the number of basic blocks in each stage, optimization issues may happen due to gradient diffusion. To address these issues, we introduce residual connections in the basic block. Apart from the first resolution stage, each stage in the encoder begins with a downsampling block in Figure 2(c), followed by several residual blocks in Figure 2(b). Different from nnU-Net which uses a separate convolution for downsampling,

we integrate downsampling in the first residual block of each stage. Such a downsampling block has two branches, namely the left and right branches as in Figure 2(c). The left branch has two $3 \times 3 \times 3$ convolutions with different strides, namely a stride of 1 for the former one and 2 for the latter. The right branch uses a kernel size $1 \times 1 \times 1$ with a stride of 2 to match the output shape of the left branch. This downsampling block is with similar residual architecture to the regular residual block in Figure 2(b), making the whole architecture neat.

### 3.2.2 Upsampling Tweaks

By default, nnU-Net uses transpose convolution with stride for upsampling. However, the convolutional kernels and strides may vary between (2,2,2) and (2,2,1) for different tasks even in the same resolution stage, which causes different weight shapes in the transpose convolution for different tasks. It further leads to a weight mismatch when we want to transfer the weights from one task to another. To solve this problem, we use interpolation followed by a $1 \times 1 \times 1$ convolution with stride 1 to replace the transpose convolution, as the weight-free interpolation has no weight shape issue (see Figure 2(f)). We use nearest neighbor interpolation for up-sampling based on our experimental results (see Table 6), which show that nearest neighbor interpolation not only offers faster processing but also achieves comparable performance to cubic linear interpolation.

### 3.3. Scaling Strategy

Based on the refined modules presented above, we can obtain a new network architecture called STU-Net. We further expand the depth and width in each stage of our model to scale it up. Deeper networks usually have larger receptive fields and better representation abilities. Wider networks tend to extract richer multi-scale features in each layer. As revealed in EfficientNet [35], depth scaling and width scaling are not independent. It is desirable to scale the depth and width of a network in a compound way to achieve better ac-

Table 2. Our proposed STU-Net with different scales. Depth refers to the number of residual blocks in each resolution stage, and width denotes the channel count in each resolution stage. Parameter calculations are based on a single-channel input and a 105-channel output, which accounts for 104 foreground classes and 1 background class in the TotalSegmentator dataset. FLOPs calculations are based on input patch sizes of 128 x 128 x 128.

| Model | depth | width | Params (M) | FLOPs (T) |
|---|---|---|---|---|
| STU-Net-S | (1,1,1,1,1,1) | (16,32,64,128,256,256) | 14.60 | 0.13 |
| STU-Net-B | (1,1,1,1,1,1) | (32,64,128,256,512,512) | 58.26 | 0.51 |
| STU-Net-L | (2,2,2,2,2,2) | (64,128,256,512,1024,1024) | 440.30 | 3.81 |
| STU-Net-H | (3,3,3,3,3,3) | (96,192,384,768,1536,1536) | 1457.33 | 12.60 |

curacy and efficiency. To simplify the scaling problem, we maintain the symmetry of our model, which means that we scale the encoder and decoder simultaneously and scale the depth and width with the same ratio in each resolution stage. Table 2 displays the different scales of STU-Net, with the suffixes 'S, B, L, H' representing 'Small, Base, Large, and Huge,' respectively.

### 3.4. Large-Scale Supervised Pre-training

We pre-train our STU-Net on the TotalSegmentator dataset. In STU-Net, the final $1 \times 1 \times 1$ convolution layer for segmentation output has 105 channels, which corresponds to the total number of target annotation categories in TotalSegmentator. To make the pre-trained models more general and transferable, we do not strictly follow the standard training procedure in nnU-Net but make some modifications. Compared to the default 1000 training epochs in nnU-Net, we pre-train our models for 4000 epochs. In addition, we discover that pre-training with mirror data augmentation can improve transfer performance on downstream tasks.

Our pre-trained models can directly perform inference on downstream datasets that consist of CT images with target segmentation categories within the upstream 104 classes, without further tuning. If the downstream tasks have novel labels or different modalities, we use our trained model as initialization and randomly initialize the segmentation output layer to match the number of target output classes. For fine-tuning, the segmentation head is randomly initialized while the weights of the remaining layers are loaded from our pre-trained model. These weights are fine-tuned with a smaller learning rate ($0.1\times$) than that of the segmentation head, which leads to better results.

### 4. Experiments

**Dataset** We train our STU-Net of different scales on the TotalSegmentator [38] dataset which contains 1204 images with 104 anatomical structures (consisting of 27 organs, 59 bones, 10 muscles and 8 vessels). It covers most of the clinical segmentation targets for normal structures in the whole body. All the images are resampled to $1.5 \times 1.5 \times 1.5$ mm

isotropic resolution. We follow the original data split in [38] that uses 1081 cases for training, 57 cases for validation, and 65 cases for final testing. It is noticeable that faces have been blurred for data privacy reasons. We evaluate our trained STU-Net on 14 public datasets for direct inference and 3 public datasets for further fine-tuning to test the transferability of our trained models. The detailed properties of these downstream datasets are shown in Appendix.

**Evaluation Metric** We adopt the Dice Similarity Coefficient (DSC) as the evaluation metric and a higher DSC score indicates better segmentation performance. For a fair comparison, we report the results of models training at the last epoch instead of the best one.

**Implementation Details** We run all the experiments based on the environment of Python 3.8, CentOS 7, Pytorch 1.10, and nnU-Net 1.7.0. We roughly follow the default data pre-processing, data augmentation and training procedure in nnU-Net. We use SGD optimizer with Nestrov momentum of 0.99, and a weight decay of 1e-3. The batch size is fixed to 2 and each epoch contains 250 iterations. For all datasets, the learning rate starts at 0.01 when training from scratch, except for AutoPET [11], which begins at 0.001, following the state-of-the-art solution [40]. The learning rate is decayed following the poly learning rate policy: $(1 - epoch/1000)^{0.9}$. We adopt data augmentation of additive brightness, gamma, rotation, scaling, mirror, and elastic deformation on the fly during training. The pre-training patch size on TotalSegmentator is $128 \times 128 \times 128$. Fine-tuning patch sizes on downstream tasks are configured by nnU-Net automatically. Models are trained on NVIDIA Tesla A100 cards with 80 GB VRAM.

### 4.1. Quantitative Results on TotalSegmentator

To verify the effectiveness of our STU-Net with different scales, we compare our models with other state-of-the-art methods on the validation set of TotalSegmentator. For a fair comparison, we train all models in nnU-Net framework for 1000 epochs. To improve the performance of other methods, we use the optimizers, learning rates, and learning rate decay strategies reported in their paper. Following the approach in the original SwinUNETR [12] paper, where SwinUNETR-B has twice the feature size of SwinUNETR-S, we further scale the feature size of SwinUNETR-B by doubling it to 96, resulting in SwinUNETR-L. This allows us to observe the performance of SwinUNETR when the model is scaled up to larger size.

As shown in Figure 1 and Table 3, our STU-Net-B model surpasses both the best CNN-based model, nnU-Net, and the best transformer-based model, SwinUNETR-B, by 0.36% and 4.48% in terms of mean DSC on all classes, respectively. Further scaling our base model to large and huge sizes leads to 1.59% and 2.94% improvement in mean DSC score, respectively. We also observe that scaling up the

Table 3. Segmentation results of different methods on TotalSegmentator validation dataset. Mean DSC (% ↑) is evaluated. Due to the limited space, we divide all 104 classes into 5 groups and report the mean DSC of these 5 group classes and all classes, respectively. *: We change the maximum feature number in nnU-Net from 320 to 512 to match the parameters of our STU-Net-B model. †: SwinUNETR-L is obtained by changing the feature size in SwinUNETR-B from 48 to 96.

| Methods | Params (M) | FLOPs (T) | TotalSeg_organs | TotalSeg_vertebrae | TotalSeg_cardiac | TotalSeg_muscles | TotalSeg_ribs | TotalSeg_all |
|---|---|---|---|---|---|---|---|---|
| nnU-Net [18] | 31.28 | 0.54 | 87.45 | 86.97 | 88.70 | 85.05 | 86.11 | 86.76 |
| nnU-Net* [18] | 60.18 | 0.55 | 87.48 | 86.87 | 88.43 | 86.42 | 85.98 | 86.94 |
| nnFormer [44] | 153.97 | 2.04 | 79.26 | 73.87 | 75.96 | 74.97 | 74.03 | 75.37 |
| UNETR-S [13] | 93.01 | 0.16 | 72.59 | 71.94 | 72.67 | 66.20 | 73.06 | 71.27 |
| UNETR-B [13] | 102.02 | 0.59 | 77.11 | 74.21 | 77.57 | 73.29 | 74.06 | 75.05 |
| SwinUNETR-S [12] | 15.71 | 0.19 | 76.29 | 77.04 | 76.05 | 71.33 | 74.21 | 74.94 |
| SwinUNETR-B [12] | 62.19 | 0.76 | 84.18 | 83.27 | 83.45 | 81.07 | 81.71 | 82.64 |
| SwinUNETR-L† [12] | 248.1 | 2.99 | 85.39 | 87.27 | 85.25 | 82.97 | 83.63 | 84.91 |
| STU-Net-S | 14.60 | 0.13 | 84.72 | 83.05 | 84.80 | 82.92 | 83.67 | 83.74 |
| STU-Net-B | 58.26 | 0.51 | 87.67 | 86.46 | 88.64 | 86.46 | 86.82 | 87.12 |
| STU-Net-L | 440.30 | 3.81 | 88.92 | 88.71 | 89.66 | 87.61 | 88.79 | 88.71 |
| STU-Net-H | 1457.33 | 12.60 | **89.82** | **90.43** | **90.89** | **88.83** | **90.29** | **90.06** |

SwinUNETR model leads to significant performance improvement, but still performs worse than our STU-Net-B. Our STU-Net-H achieves the highest mean DSC across all classes and within five sub-class groups in TotalSegmentator dataset. The results show the effectiveness of our architectural refinements on nnU-Net and the scaling strategy.

## 4.2. Transferability of Trained Models

We evaluate the transferability of our trained models by 1) conducting direct inference on the downstream CT datasets, which contain a subset of the 104 classes found in the TotalSegmentator dataset; and 2) fine-tuning the trained models on three downstream datasets, which include classes (e.g., lesions) not present in TotalSegmentator dataset and modalities (e.g., MR, PET) other than CT.

### 4.2.1 Direct Inference Results

We use the pre-trained STU-Net to conduct direct inference on 14 downstream datasets that include annotation categories within the TotalSegmentator dataset. These 14 datasets contain 2494 cases in total which is a strong benchmark to evaluate models' transferability pre-trained on large-scale datasets.

The TotalSegmentator dataset provides detailed annotations for normal organs, including separate labels for left and right organs (without lesion annotations). However, downstream datasets may contain annotations for lesions within organ regions or combined annotations for left and right organs. These differences in annotation protocols can introduce extra labels that never appear in TotalSegmentator, leading to label inconsistency. To deal with this issue, we merge lesion annotations with their corresponding organ labels and, when necessary, combine left and right organ annotations to create modified labels for evaluation purposes.

Table 4 shows that with TotalSegmentator for pre-training, models of larger scales usually have higher mean DSC scores across all these 14 datasets. This conclusion generally applies to each specific dataset. The better performance supports that our elaborately designed large-scale STU-Net pre-trained on the large-scale dataset can have better transfer capacities.

Table 4. Evaluation on the transferability of models trained on TotalSegmentator. Mean DSC (%) of direct inference on various downstream datasets are reported. AMOS-CT denotes the results on the CT modality only.

| Datasets | nnU-Net [18,38] | STU-Net-B | STU-Net-L | STU-Net-H |
|---|---|---|---|---|
| MSD Liver [1] | 95.29 | 95.80 | 95.87 | 95.88 |
| MSD Pancreas [1] | 76.52 | 75.96 | 78.41 | 78.95 |
| MSD Spleen [1] | 93.97 | 95.46 | 95.38 | 95.52 |
| BTCV [23] | 80.14 | 80.96 | 83.05 | 83.83 |
| BTCV-Cervix [23] | 54.22 | 89.43 | 89.95 | 89.79 |
| AbdomenCT-1K [29] | 90.29 | 91.61 | 92.15 | 92.27 |
| KiTS2021 [16] | 86.00 | 83.45 | 84.38 | 85.44 |
| FLARE22 [28] | 85.39 | 88.27 | 89.61 | 89.87 |
| CT-ORG [32] | 69.01 | 72.09 | 71.90 | 73.14 |
| AMOS-CT [19] | 79.26 | 80.99 | 82.87 | 83.11 |
| KiPA2022 [15] | 30.72 | 53.77 | 67.04 | 78.44 |
| Verse2020 [25] | 67.82 | 62.01 | 65.35 | 66.65 |
| WORD [27] | 78.73 | 76.07 | 78.08 | 77.42 |
| SegThor [22] | 81.79 | 84.02 | 85.59 | 85.91 |
| Mean DSC (%) | 76.37 | 80.71 | 82.83 | 84.02 |

### 4.2.2 Fine-tuning Results

We fine-tune the pre-trained STU-Net with different scales on three downstream datasets that contain novel anatomical structures, modalities and domains: AutoPET22 [11] (with lesions and PET modality), AMOS22 [19] (with MR modality), and FLARE22 [28] (with multiple domains). Owing to our architectural refinement for nnU-Net, all the model weights, except these of the segmentation head, can be easily transferred to these downstream tasks. Note that these downstream datasets, such as AutoPET, may contain multi-modal inputs which can have more channels than the input data of the pre-training dataset (e.g., a single channel for TotoalSegmentator). More input channels require the first convolutional layer of the model to have more channels correspondingly, further resulting in a shape mismatch of

Table 5. Fine-tuning results on 3 downstream datasets. Mean DSC (% ↑) is evaluated. AMOS-CT (or -MR) denotes the results on the CT (or MR) modality only, otherwise, the results are on the mixed modalities. The same applies to AutoPET. The suffix 'ft' means fine-tuning.

| Methods | FLARE22 [28] | AMOS [19] | AMOS-CT [19] | AMOS-MR [19] | AutoPET [11] | AutoPET-CT [11] | AutoPET-PET [11] | Mean DSC (%) |
|---|---|---|---|---|---|---|---|---|
| nnU-Net [18] | 85.86 | 89.48 | 89.98 | 87.39 | 73.03 | 41.84 | 71.81 | 77.06 |
| STU-Net-B | 86.56 | 89.70 | 89.84 | 86.92 | 73.96 | 41.09 | 72.25 | 77.19 |
| STU-Net-L | 87.41 | 90.29 | 90.57 | 87.04 | 74.94 | 40.97 | 73.37 | 77.80 |
| STU-Net-H | 87.29 | 90.13 | 90.39 | 87.18 | 72.19 | 43.72 | 72.59 | 77.64 |
| STU-Net-B-ft | 87.11 | 89.73 | 90.25 | 87.49 | 74.68 | 50.03 | 75.37 | 79.23 |
| STU-Net-L-ft | 88.25 | **90.50** | 90.79 | 87.69 | 75.19 | 51.88 | 77.04 | 80.19 |
| STU-Net-H-ft | **88.81** | 90.49 | **91.16** | **87.98** | **76.29** | **52.80** | **77.32** | **80.69** |

convolutional weights between the pre-training and downstream fine-tuning tasks. In this case, we replicate the channels of the first layer's weights of the pre-trained model from a single channel to multiple channels to adapt to downstream tasks. On AutoPET and AMOS datasets where multiple modalities exist, we fine-tune and evaluate our models on every single modality separately (denoted with a suffix of modality name in Table 5) and on the mixed modalities.

As shown in Table 5, fine-tuning our STU-Net models, which are pre-trained on TotalSegmentator, leads to better segmentation performance than models trained from scratch on downstream datasets. Particularly, our huge model (STU-Net-H-ft) outperforms all the other models, achieving the highest mean DSC of 80.69% on these downstream datasets. This observation underscores the importance of both pre-training and large model size in enhancing segmentation performance. A qualitative comparison of the segmentation results from our models and nnU-Net is provided in Figure 3. This visual representation further highlights the advantages of our STU-Net models when fine-tuned on downstream datasets, and illustrates the superiority of our approach in various medical imaging scenarios. Notably, the improvement over the other competitors is most significant on the AutoPET dataset, probably because pre-training on TotalSegmentator provides models with the information on the normal anatomy of the whole body. Such information is complementary to (and can be the prior information as) the lesion information which is the only annotation information in AutoPET. Thus, fine-tuning the pre-trained model can effectively enhance the model's ability to segment lesions on AutoPET.

Interestingly, when trained from scratch, our huge model is slightly worse than the large one, probably because the limited training samples on these datasets cannot further boost the larger model. With sufficient training samples from TotalSegmentator for pre-training, our huge model can significantly benefit from fine-tuning and exceed the large one, usually by a clear margin.

Notably, pre-trained models fine-tuned on non-CT modalities also demonstrate significant performance improvements, such as AMOS-MR and AutoPET-PET datasets, despite TotalSegmentator's focus on CT scans. This suggests that pre-training facilitates learning funda-

Table 6. Segmentation performance of models with different architectural variants on TotalSegmentator validation dataset. Mean DSC (% ↑) is evaluated. *: We change the maximum feature number in standard nnU-Net from 320 to 512 to match the parameters of our STU-Net-B. Abbreviations: DS (Downsampling), US (Upsampling).

| Architecture | Params (M) | FLOPs (T) | Mean DSC (%) |
|---|---|---|---|
| nnU-Net | 31.28 | 0.54 | 86.76 |
| nnU-Net* | 60.18 | 0.55 | 86.94 |
| **STU-Net-B** | **58.26** | **0.51** | **87.12** |
| STU-Net-B (Replace w/ Conv DS) | 66.02 | 0.54 | 86.41 |
| STU-Net-B (Replace w/ Transpose Conv US) | 61.32 | 0.56 | 86.96 |
| STU-Net-B (Replace w/ Trilinear Interpolation US) | 58.26 | 0.51 | 86.67 |

mental features and structures that generalize across modalities, beyond modality-specific characteristics.

### 4.3. Ablation Study

#### 4.3.1 Efficacy of Architectural Refinements

Our proposed STU-Net-B model is an improved version of the default nnU-Net architecture, incorporating several refinements. Table 6 compares the segmentation performance of different STU-Net-B architectural variants on the validation set of TotalSegmentator. Our STU-Net-B model utilizes the nearest interpolation for upsampling and incorporates downsampling in the first residual block. To evaluate its performance, we conducted a comparative analysis with alternative downsampling methods employing separate convolutions, and alternative upsampling methods utilizing transpose convolutions or trilinear interpolation.

We first increase the maximum feature number of standard nnU-Net from 320 to 512 to match the parameters of our STU-Net-B, and denote it as nnU-Net*. The nnU-Net* works better than the standard nnU-Net, but performs slightly worse than our STU-Net-B. This comparison demonstrates the effectiveness of the refinements in STU-Net-B.

Then, we explore the downsampling design in STU-Net by introducing a variant that employs convolutional downsampling instead of integrating the downsampling process within the first residual block of each stage. This modification results in reduced performance and higher computational costs. We further investigate the upsampling refinements by designing two variants of STU-Net-B. The first
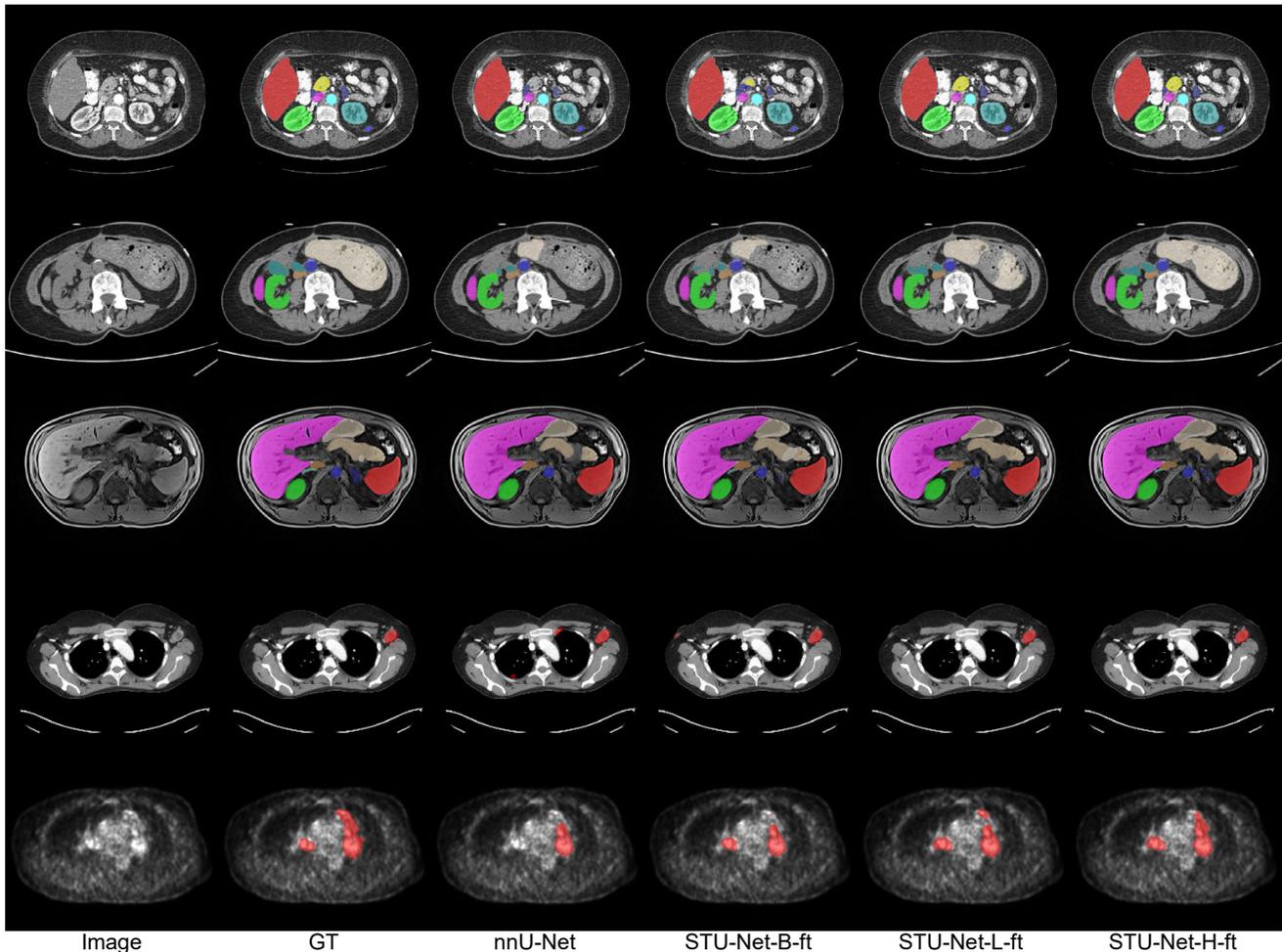
Figure 3. Qualitative visualization of our STU-Net with different scales and nnU-Net on various medical imaging datasets. The representative cases from distinct datasets are displayed in each row, including Row 1 - FLARE22 dataset, Row 2 - AMOS dataset with CT images, Row 3 - AMOS dataset with MR images, Row 4 - AutoPET dataset with CT images, and Row 5 - AutoPET dataset with PET images. The seven columns from left to right correspond to the original image, the ground truth (gt), the nnU-Net results, and our STU-Net-B-ft, STU-Net-L-ft, and STU-Net-H-ft results.

one utilizes transpose convolution to replace the default interpolation and convolution-based upsampling. This incurs a 0.16% decrease in performance, and makes the weights not transferable for downstream fine-tuning. The second one uses trilinear (or cubic linear) interpolation to replace the nearest neighbor interpolation. This change decreases the performance and slows down the running speed. Overall, the default upsampling design in STU-Net achieves better performance, faster running speed, and better transfer capacities.

Thus, our proposed refinements not only enhance the effectiveness and efficiency of nnU-Net, but also endow it with the crucial properties of weight transferability and scalability. These properties are essential for further scaling up the model and facilitating transfer learning.

### 4.3.2 Scaling Strategy

We apply three different scaling strategies to nnU-Net and our STU-Net-base, namely scaling using different depth coefficient $d \in [1.0, 2.0, 3.0, 4.0]$, width coefficient $w \in [1.0, 2.0, 3.0, 4.0]$, and simultaneously with depth and width coefficients of $[1.0, 2.0, 3.0]$. The coefficient $d$ (or $w$) means that the depth (width) of each stage is scaled to $d$ (or $w$) times. Table 7 shows the results of different scaling strategies on the TotalSegmentator dataset. Firstly, wider nnU-Net* and STU-Net consistently obtain better performance while deeper ones are not. Therefore, compared to depth scaling, width scaling is more effective in improving the model performance on the large-scale dataset. But it brings a significant increase in computational consumption. Secondly, compared to the performance drop by deeper nnU-

Table 7. Segmentation performance of models with different scaling dimensions on TotalSegmentator dataset. Mean DSC (% ↑) is evaluated. *: We change the maximum feature number in standard nnU-Net from 320 to 512 to match the parameters of our STU-Net-B.

| Methods | Params (M) | FLOPs (T) | Mean DSC (%) |
|---|---|---|---|
| nnU-Net* (d=1.0, w=1.0) | 60.18 | 0.55 | 86.94 |
| STU-Net (d=1.0, w=1.0) | 58.26 | 0.51 | 87.12 |
| scale nnU-Net* by width (w=2.0) | 240.47 | 2.19 | 87.29 |
| scale STU-Net by width (w=2.0) | 232.80 | 2.00 | 88.08 |
| scale nnU-Net* by width (w=3.0) | 540.88 | 4.92 | 87.80 |
| scale STU-Net by width (w=3.0) | 523.62 | 4.49 | 88.85 |
| scale nnU-Net* by width (w=4.0) | 961.40 | 8.73 | 88.84 |
| scale STU-Net by width (w=4.0) | 930.71 | 7.97 | 89.19 |
| scale nnU-Net* by depth (d=2.0) | 112.06 | 1.01 | 85.65 |
| scale STU-Net by depth (d=2.0) | 110.15 | 0.96 | 87.72 |
| scale nnU-Net* by depth (d=3.0) | 163.94 | 1.46 | 83.70 |
| scale STU-Net by depth (d=3.0) | 162.03 | 1.41 | 87.99 |
| scale nnU-Net* by depth (d=4.0) | 215.83 | 1.91 | 81.45 |
| scale STU-Net by depth (d=4.0) | 213.91 | 1.86 | 87.58 |
| nnU-Net* compound scale (d=2.0, w=2.0) | 447.97 | 4.00 | 87.42 |
| STU-Net compound scale (d=2.0, w=2.0) | 440.30 | 3.81 | 88.71 |
| nnU-Net* compound scale (d=3.0, w=3.0) | 1474.59 | 13.03 | 87.50 |
| STU-Net compound scale (d=3.0, w=3.0) | 1457.33 | 12.60 | 90.06 |

Net, our STU-Net can better benefit from scaling depth to a certain extent, e.g., 87.12→87.72→87.99→87.58 vs. nnU-Net's 86.94→85.65→83.70→81.45. We owe the better performance of STU-Net to its residual design. Thirdly, compound scaling, i.e., increasing depth and width simultaneously, is more effective and efficient in improving the performance of our STU-Net than the other two scaling strategies. Lastly, even with the same scaling strategy, and similar parameters and FLOPs, our STU-Net consistently outperforms nnU-Net* on all the settings, which again verifies the effectiveness of our refinement.

### 4.3.3 Pre-training and Fine-tuning Strategy

We perform empirical studies on the effectiveness of pre-training, and then on mirror augmentation used in the pre-training and fine-tuning stages. Finally, we study the training epochs for pre-training.

Firstly, we study the effectiveness of large-scale pre-training. We use STU-Net-L pre-trained on the TotalSegmentator dataset and then fine-tuned on FLARE22 and AutoPET datasets to conduct such a study. Comparing the 1st (w/o pre-training) and 3rd (w/ pre-training) rows in Table 8, we find that the model with pre-training outperforms the one without pre-training by 0.84% and 3.67%, respectively. Such better performance justifies the effectiveness of large-scale pre-training.

Secondly, we investigate the mirror augmentation used in pre-training and fine-tuning. With pre-training, the mirror is useful for improving the DSC on the downstream task as the mirror (3rd∼5th rows) improves the DSC of the 6th row (no mirror) considerably. It is also worth noting that the mirror used in both pre-training and fine-tuning stages achieves the best result.

Table 8. Segmentation performance of different pre-training settings on STU-Net-L. Mean DSC (% ↑) is evaluated. The first and second rows compare whether to use mirror data augmentation when training models from scratch (without pre-training). The third to sixth rows compare the performance of different combinations of mirror data augmentation to pre-training and fine-tuning. The last four rows investigate the effect of different pre-training epoch numbers on fine-tuning results.

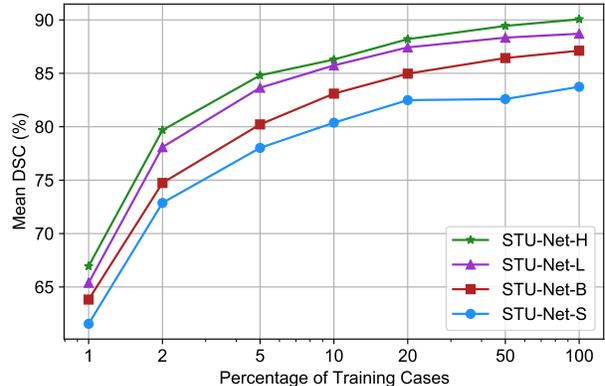| Pre-train epochs | Pre-train mirror | Fine-tune mirror | FLARE22 [28] | AutoPET-PET [11] |
|---|---|---|---|---|
| - | - | ✓ | 87.41 | 73.37 |
| - | - | ✗ | 85.79 | 73.59 |
| 4k | ✓ | ✓ | **88.25** | **77.04** |
| 4k | ✓ | ✗ | 87.65 | 76.34 |
| 4k | ✗ | ✓ | 86.16 | 75.88 |
| 4k | ✗ | ✗ | 86.41 | 73.90 |
| 1k | ✓ | ✓ | 87.06 | 74.89 |
| 2k | ✓ | ✓ | 87.63 | 76.89 |
| 4k | ✓ | ✓ | **88.25** | **77.04** |
| 8k | ✓ | ✓ | 87.71 | 76.87 |



Figure 4. Comparison of mean DSC (% ↑) performance for STU-Net models with different scales, trained on subsets of the TotalSegmentator training set with different proportions of training cases, and evaluated on the same TotalSegmentator validation set.

Finally, with the mirror incorporated in both pre-training and fine-tuning, we find that different pre-training epochs can also influence performance (see the last four rows). The best result is obtained at the 4k epochs which can ensure that the model is sufficiently trained for convergence.

### 4.3.4 Influence of Dataset Sizes on Model Performance

We investigate the impact of dataset size on the performance of our models when training them on the TotalSegmentator dataset. It is important to note that the training cases at different proportions were obtained through a stratified random selection process, ensuring that higher proportions of training cases also include the data from lower proportions.

As illustrated in Figure 4, increasing the model size leads to better segmentation performance on the TotalSegmentator subset, regardless of the number of training cases. For
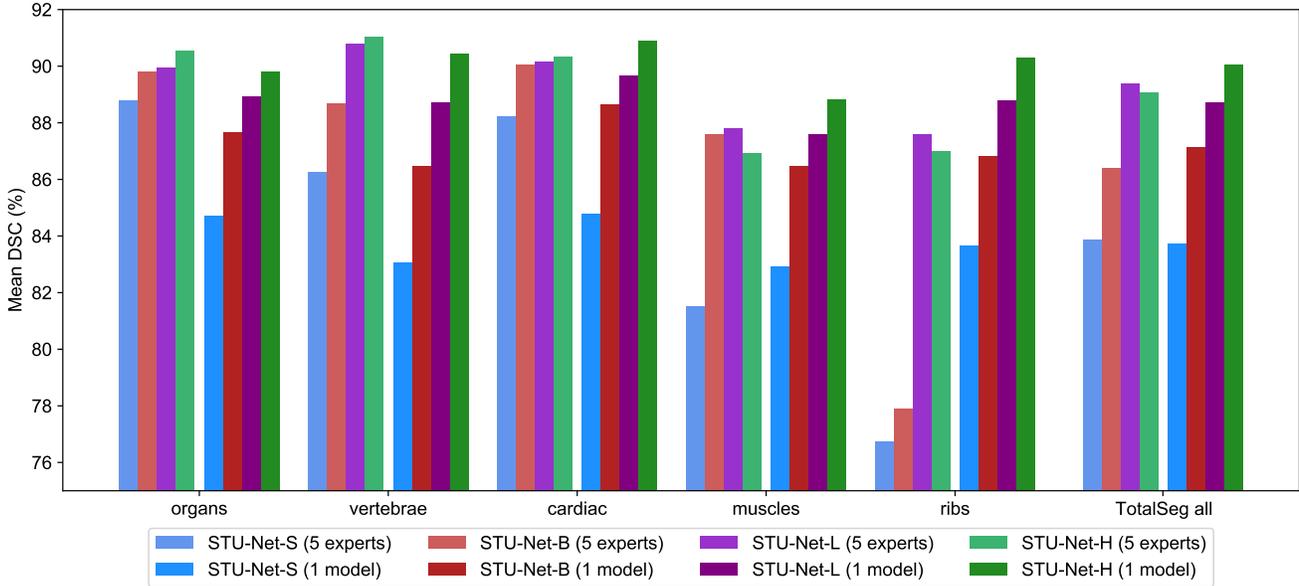
Figure 5. Comparison between five specialized expert STU-Net models and a single universal STU-Net model on the TotalSegmentator dataset. Each expert model targets one of the five subcategories (i.e., organs, vertebrae, cardiac, muscles, and ribs), while the universal model is trained on all 104 classes. The performance is measured using the mean DSC across various anatomical categories: the five subcategories and an overall performance metric for TotalSegmentator dataset. STU-Net architectures (S, B, L, and H) are depicted for both expert and universal models. Lighter colors represent expert models, and darker colors indicate universal models.

example, STU-Net-H outperforms STU-Net-S even when trained with only 5% of the cases. Similarly, STU-Net-H surpasses STU-Net-B when trained with just 20% of the cases. These results suggest that large-scale models are more data-efficient than smaller models for medical image segmentation. Moreover, the performance of different models is consistently improved with an increasing number of cases, and the trend has not yet saturated. These observations indicate that augmenting the number of training cases based on the TotalSegmentator dataset can yield further performance enhancements.

### 4.4. Universal Model Versus Expert Models

We evaluate the performance of a universal STU-Net model trained on all 104 classes in the TotalSegmentator dataset, against five expert STU-Net models, each targeting one of the five subcategories (the same as in Table 3). Furthermore, we investigate the influence of model size on the performance of both expert and universal models.

Figure 5 illustrates that as model size increases, the performance of both expert and universal models generally improves. Expert models excel in the organs, vertebrae, and cardiac subcategories, while universal models perform better in the muscles and ribs subcategories. For the largest models (STU-Net-H), the universal model surpasses expert models, achieving the highest overall mean DSC score of 90.06% on all the classes of TotalSeg dataset, compared to

the expert models' highest mean DSC score of 89.07%.

The results reveal that although expert models may outperform universal models in specific subcategories, universal models consistently deliver strong performance across different anatomical structures. The performance gap between expert and universal models varies across subcategories as the model size increases: the gap narrowing for organs and vertebrae, reversing for cardiac, and widening for muscles and ribs. The findings suggest that with increased model size, universal models are capable of concurrently segmenting numerous categories and exhibit promising performance advancements.

## 5. Outlook

Our work is an attempt toward Medical Artificial General Intelligence (MedAGI). We believe that the keys to MedAGI should be similar to that of AGI in the computer vision community – the foundation model [4, 21] and large-scale datasets. Large-scale datasets are emerging in the medical image processing field, while the foundation models, which can generalize to unseen tasks and data distributions, are less explored. The foundation models usually refer to the backbones which can be used to handle many tasks, like ResNet [14] or vision transformer [7, 9] in computer vision. To this end, they are usually large-scale models with powerful representation abilities to extract discrim-

inative features for different tasks. But in the medical image processing field, we lack such models. We thus explore how to design large-scale STU-Net models and take this problem as the initial step to the foundation models. Starting from STU-Net, researchers can build foundation models that can segment everything in medical images, like [21], or further generalize these segmentation-based foundation models to many other tasks beyond segmentation, such as detection, classification, reconstruction and registration. Following this way, we will gradually approach the goal of foundation models: MedAGI, and use MedAGI to contribute to human health.

## 6. Conclusions

In this work, we introduce a series of scalable and transferable medical image segmentation models, named STU-Net, based on the nnU-Net framework. Our STU-Net models are designed to be scalable, with the largest model consisting of 1.4 billion parameters, making it the most substantial medical image segmentation model to date. By training our STU-Net models on the large-scale TotalSegmentator dataset, we demonstrate that scaling the model size yields significant performance improvements when transferring them to various downstream tasks. It thus underscores the potential of large-scale models in the medical image segmentation domain.

Furthermore, our STU-Net-H model, trained on the TotalSegmentator dataset, exhibits robust transferability in both direct inference and fine-tuning scenarios across multiple downstream datasets. This observation emphasizes the practical value of leveraging large-scale pre-trained models for medical image segmentation tasks.

In conclusion, the development of scalable and transferable STU-Net models has the potential to advance the state-of-the-art in medical image segmentation, opening new avenues for research and innovation within the medical image segmentation community.

## References

[1] Michela Antonelli, Annika Reinke, Spyridon Bakas, Keyvan Farahani, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, Olaf Ronneberger, Ronald M Summers, et al. The medical segmentation decathlon. *Nature communications*, 13(1):4128, 2022. 6, 14

[2] Ujjwal Baid, Satyam Ghodasara, Suyash Mohan, Michel Bilello, Evan Calabrese, Errol Colak, Keyvan Farahani, Jayashree Kalpathy-Cramer, Felipe C Kitamura, Sarthak Pati, et al. The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification. *arXiv preprint arXiv:2107.02314*, 2021. 2

[3] Ivo M Baltruschat, Hanna Ćwieka, Diana Krüger, Berit Zeller-Plumhoff, Frank Schlünzen, Regine Willumeit-Römer, Julian Moosmann, and Philipp Heuser. Scaling the

u-net: segmentation of biodegradable bone implants in high-resolution synchrotron radiation microtomograms. *Scientific reports*, 11(1):24237, 2021. 2, 3

[4] 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. 10

[5] 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*, 33:1877–1901, 2020. 2, 3

[6] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image segmentation. *arXiv preprint arXiv:2102.04306*, 2021. 3

[7] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. *arXiv preprint arXiv:2302.05442*, 2023. 2, 3, 10

[8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 3

[9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021. 3, 10

[10] Michal Drozdzal, Eugene Vorontsov, Gabriel Chartrand, Samuel Kadoury, and Chris Pal. The importance of skip connections in biomedical image segmentation. In *International Workshop on Deep Learning in Medical Image Analysis, International Workshop on Large-Scale Annotation of Biomedical Data and Expert Label Synthesis*, pages 179–187. Springer, 2016. 3

[11] Sergios Gatidis, Tobias Hepp, Marcel Früh, Christian La Fougère, Konstantin Nikolaou, Christina Pfannenberg, 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. 2, 5, 6, 7, 9, 15

[12] 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 *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 7th International Workshop, BrainLes 2021, Held in Conjunction with MICCAI 2021, Virtual Event, September 27, 2021, Revised Selected Papers, Part I*, pages 272–284. Springer, 2022. 2, 3, 5, 6

[13] Ali Hatamizadeh, Yucheng Tang, Vishwesh Nath, Dong Yang, Andriy Myronenko, Bennett Landman, Holger R

Roth, and Daguang Xu. Unetr: Transformers for 3d medical image segmentation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 574–584, 2022. 2, 3, 6

[14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2, 3, 10

[15] 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. 6, 14

[16] Nicholas Heller, Fabian Isensee, Klaus H Maier-Hein, Xiaoshuai Hou, Chunmei Xie, Fengyi Li, Yang Nan, Guangrui Mu, Zhiyong Lin, Miofei Han, et al. The state of the art in kidney and kidney tumor segmentation in contrast-enhanced ct imaging: Results of the kits19 challenge. *Medical Image Analysis*, page 101821, 2020. 6, 14

[17] Ziyan Huang, Zehua Wang, Zhikai Yang, and Lixu Gu. Adwu-net: adaptive depth and width u-net for medical image segmentation by differentiable neural architecture search. In *International Conference on Medical Imaging with Deep Learning*, pages 576–589. PMLR, 2022. 3

[18] 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. 2, 6, 7

[19] 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. 6, 7, 14, 15

[20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020. 2, 3

[21] 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. 2, 10, 11

[22] Zoé Lambert, Caroline Petitjean, Bernard Dubray, and Su Kuan. Segthor: Segmentation of thoracic organs at risk in ct images. In *2020 Tenth International Conference on Image Processing Theory, Tools and Applications (IPTA)*, pages 1–6. IEEE, 2020. 6, 14

[23] 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, page 12, 2015. 6, 14

[24] Lei Li, Veronika A Zimmer, Julia A Schnabel, and Xiahai Zhuang. Medical image analysis on left atrial lge mri for atrial fibrillation studies: A review. *Medical image analysis*, page 102360, 2022. 1

[25] Hans Liebl, David Schinz, Anjany Sekuboyina, Luca Malagutti, Maximilian T Löffler, Amirhossein Bayat, Malek El Husseini, Giles Tetteh, Katharina Grau, Eva Niederreiter, et al. A computed tomography vertebral segmentation dataset with anatomical variations and multi-vendor scanner data. *Scientific Data*, 8(1):284, 2021. 6, 14

[26] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 3

[27] Xiangde Luo, Wenjun Liao, Jianghong Xiao, Jieneng Chen, Tao Song, Xiaofan Zhang, Kang Li, Dimitris N Metaxas, Guotai Wang, and Shaoting Zhang. Word: A large scale dataset, benchmark and clinical applicable study for abdominal organ segmentation from ct image. *Medical Image Analysis*, 82:102642, 2022. 6, 14

[28] Jun Ma, Yao Zhang, Song Gu, Xingle An, Zhihe Wang, Cheng Ge, Congcong Wang, Fan Zhang, Yu Wang, Yinan Xu, Shuiping Gou, Franz Thaler, Christian Payer, Darko Štern, Edward G.A. Henderson, Dónal M. McSweeney, Andrew Green, Price Jackson, Lachlan McIntosh, Quoc-Cuong Nguyen, Abdul Qayyum, Pierre-Henri Conze, Ziyan Huang, Ziqi Zhou, Deng-Ping Fan, Huan Xiong, Guoqiang Dong, Qiongjie Zhu, Jian He, and Xiaoping Yang. Fast and low-gpu-memory abdomen ct organ segmentation: The flare challenge. *Medical Image Analysis*, 82:102616, 2022. 6, 7, 9, 14, 15

[29] 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*, 44(10):6695–6714, 2021. 2, 6, 14

[30] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pages 565–571. Ieee, 2016. 3

[31] Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew Lee, Mattias Heinrich, Kazunari Misawa, Kensaku Mori, Steven McDonagh, Nils Y Hammerla, Bernhard Kainz, et al. Attention u-net: Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999*, 2018. 3

[32] 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):381, 2020. 6, 14

[33] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. 2, 3

[34] 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*, page 843–852, 2017. 3

[35] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. 3, 4

[36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3

[37] Haoyu Wang, Hongrui Yi, Jie Liu, and Lixu Gu. Integrated treatment planning in percutaneous microwave ablation of lung tumors. In *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*, pages 4974–4977. IEEE, 2022. 1

[38] 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. 2, 5, 6

[39] Wei Wu, Jingyang Zhang, Wenjia Peng, Hongzhi Xie, Shuyang Zhang, and Lixu Gu. Car-net: A deep learning-based deformation model for 3d/2d coronary artery registration. *IEEE Transactions on Medical Imaging*, 41(10):2715–2727, 2022. 1

[40] Jin Ye, Haoyu Wang, Ziyan Huang, Zhongying Deng, Yanzhou Su, Can Tu, Qian Wu, Yuncheng Yang, Meng Wei, Jingqi Niu, et al. Exploring vanilla u-net for lesion segmentation from whole-body fdg-pet/ct scans. *arXiv preprint arXiv:2210.07490*, 2022. 5, 15

[41] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. 3

[42] 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*, pages 12104–12113, 2022. 3

[43] 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*, pages 1195–1204, 2021. 3

[44] Hong-Yu Zhou, Jiansen Guo, Yinghao Zhang, Lequan Yu, Liansheng Wang, and Yizhou Yu. nnformer: Interleaved transformer for volumetric segmentation. *arXiv preprint arXiv:2109.03201*, 2021. 2, 3, 6

[45] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: Redesigning skip connections to exploit multiscale features in image segmentation. *IEEE transactions on medical imaging*, 39(6):1856–1867, 2019. 3

| Datasets | # Targets | # Annotated Scans | Overlapped Annotations w/ Totalsegmentator |
|---|---|---|---|
| 01. MSD Liver [1] | 1 | 131 | Liver* |
| 02. MSD Pancreas [1] | 1 | 281 | Pancreas* |
| 03. MSD Spleen [1] | 1 | 41 | Spleen |
| 04. BTCV [23] | 13 | 30 | Spl, RKid, LKid, Gall, Eso, Liv, Sto, Aor, IVC, R&SVeins, Pan, RAG, LAG |
| 05. BTCV-Cervix [23] | 1 | 30 | Urinary Bladder |
| 06. AbdomenCT-1K [29] | 4 | 1000 | Spleen, Kidney*, Liver, Pancreas |
| 07. KiTS2021 [16] | 1 | 300 | Kidney* |
| 08. FLARE22 [28] | 13 | 50 | Liv, Eso, Sto, Duo, LKid, RKid, Spl, Pan, Aor, IVC, RAG, LAG, Gall |
| 09. CT-ORG [32] | 4 | 140 | Lung*, Liver, Kidney* and Bladder |
| 10. AMOS-CT [19] | 15 | 200 | Spl, RKid, LKid, Gall, Eso, Liv, Sto, Aor, IVC, Pan, RAG, LAG, Duo, Bla, Pro/UTE |
| 11. KiPA2022 [15] | 1 | 70 | Kidney* |
| 12. Verse2020 [25] | 24 | 61 | Vertebrae C1-7, Vertebrae T1-12, Vertebrae L1-5 |
| 13. WORD [27] | 16 | 120 | Spl, RKid, LKid, Gall, Eso, Liv, Sto, Pan, RAG, Duo, Col, Int, Rec, Bla, LFH, RFH |
| 14. SegThor [22] | 3 | 40 | Esophagus, Trachea, Aorta |

Table 9. **Detailed information of 14 evaluation datasets for direct inference.** The table provides details on the number of overlapped targets, available annotated scans, and annotations that overlap with TotalSegmentator. *: Some conflict annotations underwent special processing to ensure consistency in targets, as described in Section A.

# Appendix

## A. Consistency of overlapped annotations

In order to maintain consistency in the annotated targets, some targets (e.g., Kidney *v.s.* Left Kidney and Right Kidney) need to be processed due to variations in annotation protocols across datasets. In some cases, we combined the inference results of related target to ensure consistency. Specifically, the following results underwent special processing: **Kidney** was created by combining the results of Totalsegmentator targets *kidney_left* and *kidney_right*; **Lung** was created by combining five Totalsegmentator targets of lung lobe, i.e., *lung_lower_lobe_left*, *lung_lower_lobe_right*, *lung_middle_lobe_right*, *lung_upper_lobe_left*, and *lung_upper_lobe_right*.

Additionally, considering the hierarchical relationships, the ground-truth annotations of organs and their inner tumors were combined. For the MSD Liver dataset, **Liver** was created by combining the ground truth of targets *liver* and *tumor*. Similarly, for the MSD Pancreas dataset, **Pancreas** was created by combining the annotations of targets *pancreas* and *tumor*. For the KiTS2021 dataset, **Kidney** was created by combining the annotations of targets *kidney*, *tumor*, and *cyst*.

## B. Dataset Details

In our evaluation, we utilized a total of 15 publicly available datasets, with 14 datasets used for direct inference and 3 datasets for fine-tuning. For direct inference evaluation, we conducted inference on all annotated cases and computed metrics for all overlapped targets between TotalSegmentator and the respective evaluation dataset, as detailed in Table 9. Details on data splitting for fine-tuning will be introduced in the following parts.

**MSD Liver** [1] contains 131 contrast-enhanced CT scans annotated with liver and tumor.

**MSD Pancreas** [1] consists of 281 portal venous phase CT scans annotated with pancreas and tumor.

**MSD Spleen** [1] includes 41 3D volumes of CT with the annotation of spleen.

**BTCV** [23] consists of 30 cases of CT with labels of 13 abdominal organs. These images are collected from metastatic liver cancer patients or post-operative ventral hernia patients.

**BTCV-Cervix** [23] includes 30 cervis CT scans annotated with 4 organs.

**AbdomenCT-1K** [29] has 1000 cases of abdominal CT scans from 12 medical centers. Each case has annotations of 4 organs.

**KiTS2021** [16] contains 300 annotated CT scans with labels of kidney, tumor and cyst. The data is collected from patients who underwent partial or radical nephrectomy for suspected renal malignancy between 2010 and 2020.

**CT-ORG** [32] is composed of 140 CT images containing 6 organ classes, which are from 8 different medical centers.

**KiPA2022** [15] is provided from the Kidney PArsing Challenge 2022. This dataset contains 70 CT images with annotations of 4 kinds of kidney-related structures.

**Verse2020** [25] is a CT spine dataset consisting of 374 scans from 355 patients. Considering the computational burden, we adopted the official train set (61 cases) for evaluation.

**WORD** [27] collects 150 CT scans from 150 patients before the radiation therapy in a single center. Each volume has 16 organs with fine pixel-level annotations and scribble-based sparse annotations. We conducted evaluation on all 120 available annotated cases.

**SegThor** [22] is the abbreviation of Segmentation of Thoracic Organs at Risk in CT images. SegThor provides 40 cases of CT scans with annotation of heart, trachea, aorta and esophagus.

**FLARE22** [28] provides 50 cases of labeled CT images and 2000 unlabeled CT images. The segmentation targets include 13 organs. We tested all 50 annotated cases for direct inference. For fine-tuning, we trained on all labeled cases and evaluated on 20 official validation cases.

**AMOS** [19] has 500 CT scans (AMOS-CT) and 100 MR scans (AMOS-MR) with voxel-level annotations of 15 abdominal organs. For our evaluation of direct inference, we utilized all 200 annotated scans from the training set of AMOS-CT. For fine-tuning, we trained on the training set (200 CT and 40 MR) and evaluated on the validation set (100 CT and 20 MR) according to official data splitting.

**AutoPET** [11] provides 1014 studies of paired 3D CT and PET scans from 900 patients where tumor lesions is annotated on each paired CT/PET images. Following the state-of-the-art solution [40], we fine-tuned our model on 398 cases and tested on 103 cases that contain lesions.