# UniverSeg: Universal Medical Image Segmentation

Victor Ion Butoi*
MIT CSAIL
vbutoi@mit.edu

Jose Javier Gonzalez Ortiz*
MIT CSAIL
josejg@mit.edu

Tianyu Ma
Cornell University
tm544@cornell.edu

Mert R. Sabuncu
Cornell University
msabuncu@cornell.edu

John Guttag
MIT CSAIL
guttag@mit.edu

Adrian V. Dalca
MIT CSAIL & MGH, HMS
adalca@mit.edu

## Abstract

*While deep learning models have become the predominant method for medical image segmentation, they are typically not capable of generalizing to unseen segmentation tasks involving new anatomies, image modalities, or labels. Given a new segmentation task, researchers generally have to train or fine-tune models, which is time-consuming and poses a substantial barrier for clinical researchers, who often lack the resources and expertise to train neural networks. We present UniverSeg, a method for solving unseen medical segmentation tasks without additional training. Given a query image and example set of image-label pairs that define a new segmentation task, UniverSeg employs a new Cross-Block mechanism to produce accurate segmentation maps without the need for additional training. To achieve generalization to new tasks, we have gathered and standardized a collection of 53 open-access medical segmentation datasets with over 22,000 scans, which we refer to as MegaMedical. We used this collection to train UniverSeg on a diverse set of anatomies and imaging modalities. We demonstrate that*
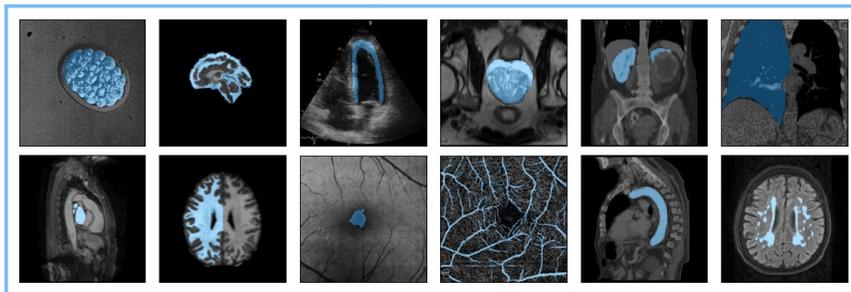
*UniverSeg substantially outperforms several related methods on unseen tasks, and thoroughly analyze and draw insights about important aspects of the proposed system. The UniverSeg source code and model weights are freely available at https://universeg.csail.mit.edu*

## 1. Introduction

Image segmentation is a widely studied problem in computer vision and a central challenge in medical image analysis. Medical segmentation tasks can involve diverse imaging modalities, such as magnetic resonance imaging (MRI), X-ray, computerized tomography (CT), and microscopy; different biomedical domains, such as the abdomen, chest, brain, retina, or individual cells; and different labels within a region, such as heart valves or chambers (Figure 1). This diversity has inspired a wide array of segmentation tools, each usually tackling one task or a small set of closely related tasks [14, 20, 38, 39, 83, 90]. In recent years, deep-learning models have become the predominant strategy for medical image segmentation [42, 71, 83].

A key problem in image segmentation is *domain shift*,
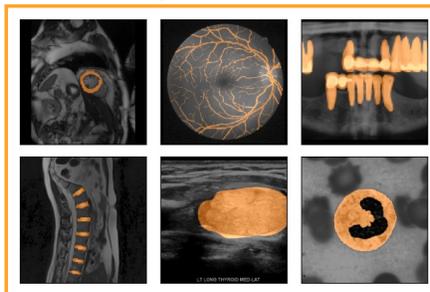
Figure 1: Medical segmentation involves many imaging types, biomedical domains, and target labels. We employ a large diverse set of training tasks **(blue)** to build a model that can segment unseen tasks **(orange)** without additional training.
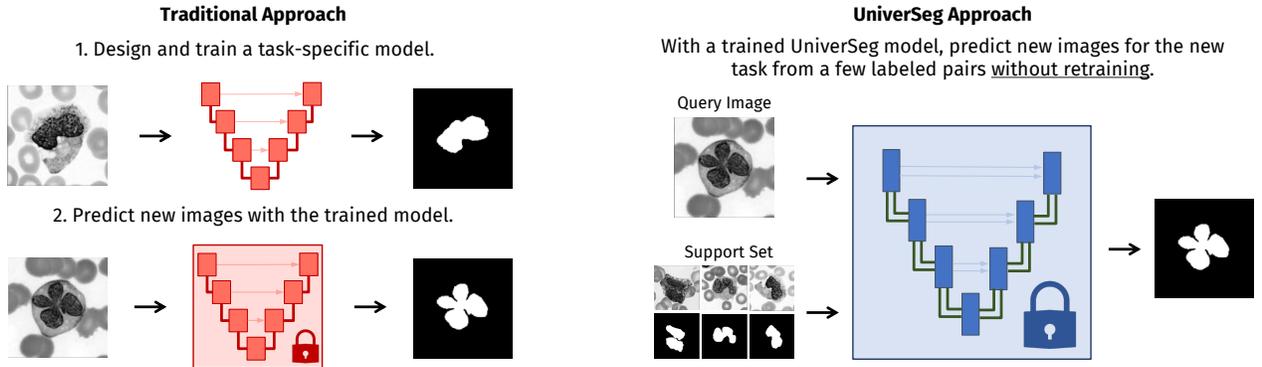
Figure 2: **Workflow for inference on a new task, from an unseen dataset.** Given a new task, traditional models **(left)** are trained before making predictions. UniverSeg **(right)** employs a *single* trained model which can make predictions for images (queries) from the new task with a few labeled examples as input (support set), without additional fine-tuning.

where models often perform poorly given out-of-distribution examples. This is especially problematic in the medical domain where clinical researchers or other scientists are constantly defining new segmentation tasks driven by evolving populations, and scientific and clinical goals. To solve these problems they need to either train models from scratch or fine-tune existing models. Unfortunately, training neural networks requires machine learning expertise, computational resources, and human labor. This is infeasible for most clinical researchers or other scientists, who do not possess the expertise or resources to train models. In practice, this substantially slows scientific development. We therefore focus on avoiding the need to do *any* training given a new segmentation tasks.

Fine-tuning models trained on the natural image domain can be unhelpful in the medical domain [82], likely due to the differences in data sizes, features, and task specifications between domains, and importantly still requires substantial retraining. Some few-shot semantic segmentation approaches attempt to predict novel classes without fine-tuning in limited data regimes, but mostly focus on classification tasks, or segmentation of new classes within the same input domain, and do not generalize across anatomies or imaging modalities.

In this paper, we present UniverSeg – an approach to learning a *single* general medical-image segmentation model that performs well on a variety of tasks without any retraining, including tasks that are substantially different from those seen at training time. UniverSeg learns how to exploit an input set of labeled examples that specify the segmentation task, to segment a new biomedical image in one forward pass. We make the following contributions.

- We propose UniverSeg – a framework that enables solving new segmentation tasks without retraining, using a novel flexible CrossBlock mechanism that transfers information from the example set to the new image.

- We demonstrate that UniverSeg substantially outperforms several models across diverse held-out segmentation tasks involving unseen anatomies, and even approaches the performance of fully-supervised networks trained specifically for those tasks.
- In extensive analysis, we show that the generalization capabilities of UniverSeg are linked to task diversity during training and image diversity during inference.

UniverSeg source code and model weights are available at https://universeg.csail.mit.edu

## 2. Related Works

**Medical Image Segmentation.** Medical image segmentation has been widely studied, with state-of-the-art methods training convolutional neural networks in a supervised fashion, predicting a label map for a given input image [20, 38, 39, 43, 83]. For a new segmentation problem, models are typically trained from scratch, requiring substantial design and tuning.

Recent strategies, such as the nnUNet [39], automate some design decisions such as data processing or model architecture, but still incur substantial overhead from training. In contrast to these methods, UniverSeg generalizes to new medical segmentation tasks without training or fine-tuning.

**Multi-task Learning.** Multi-Task Learning (MTL) frameworks learn several tasks simultaneously [13, 21, 86]. For medical imaging, this can involve multiple modalities [72], population centers [61], or anatomies [73]. However, the tasks are always pre-determined by design: once trained, each network can only solve tasks presented during training. UniverSeg overcomes this limitation, enabling tasks to be dynamically specified during inference.

**Transfer Learning.** Transfer learning strategies involve fine-tuning pre-trained models, often from a different domain [63, 97]. This is used in medical image segmentation
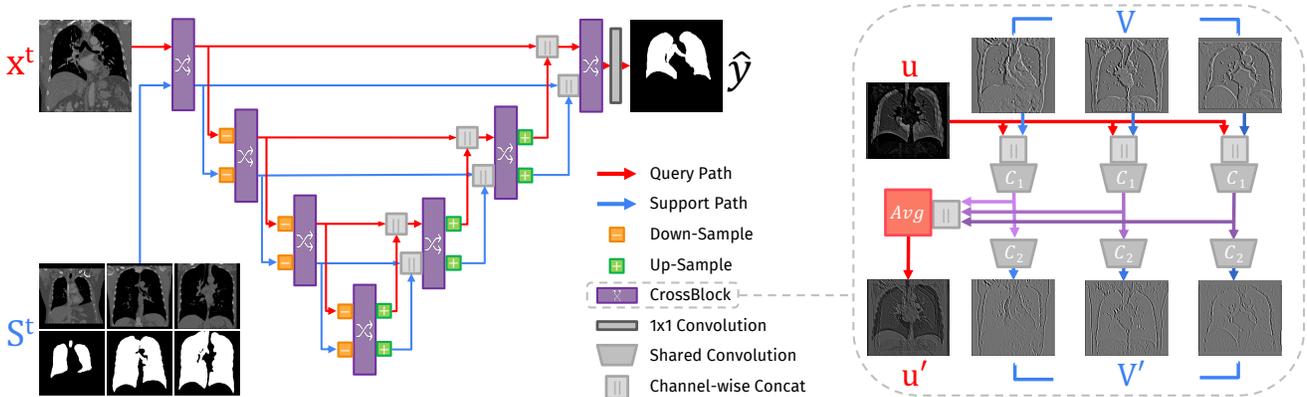
Figure 3: A UniverSeg network **(left)** takes as input a query image and a support set of image and label-maps (pairwise concatenated in the channel dimension) and employs multi-scale CrossBlock features. A CrossBlock **(right)** takes as input representations of the query $u$ and support set $V = \{v_i\}$, and interacts $u$ with each support entry $v_i$ to produce $u'$ and $V'$.

starting with models trained on natural images [3, 24, 41, 106, 109], where the amount of data far exceeds the amount in the target biomedical domain. However, this technique still involves substantial training for each new task, which UniverSeg avoids. Additionally, the differences between medical and natural images often make transfer learning from large pre-trained models unhelpful [82].

**Optimization-based Meta-Learning.** Optimization-based meta-learning techniques often learn representations that minimize downstream fine-tuning steps by using a few examples per task, sometimes referred to as few-shot learning [22, 75, 100, 94]. Meta-learning via fine-tuning has been studied in medical image segmentation to handle multiple image modalities [105], anatomies [103], and generalization to different targets [48, 49, 93]. While these strategies reduce the amount of data and training required for downstream tasks [30], fine-tuning these models nevertheless requires machine learning expertise and computational resources, which are often not available to medical researchers.

**Few-shot Semantic Segmentation.** Few-shot (FS) methods adapt to new tasks from few training examples, often by fine-tuning pretrained networks [22, 75, 100, 94]. Some few-shot semantic segmentation models generate predictions for new images (queries) containing unseen classes from just a few labeled examples (support) without additional retraining. One strategy prevalent in both natural image [74, 87, 102] and medical image [19, 59, 77, 91] FS segmentation methods is to employ large pre-trained models to extract deep features from the query and support images. These methods often involve learning meaningful prototypical representations for each label [101]. Another medical FS segmentation strategy uses self-supervised learning to make up for the lack of training data and tasks [29, 76]. In contrast to UniverSeg, these methods, focused on limited

data regimes, tackle specific tasks involving generalizing to new classes in a particular subdomain, like abdominal CT or MRI scans [29, 76, 84, 98].

In our work, we focus on avoiding *any* fine-tuning, even when given many examples for a new task, to avoid requiring the clinical or scientific user to have machine learning expertise and compute resources. Our proposed framework draws inspiration from ideas from some few-shot learning solutions, but aims to generalize to a universally broad set of anatomies, modalities, and datasets – even those completely unseen during training.

## 3. UniverSeg Method

Let $t$ be a segmentation task comprised of a set of image-label pairs $\{(x_i^t, y_i^t)\}_{i=1}^N$. Common segmentation strategies learn parametric functions $\hat{y} = f_\theta^t(x)$, where $f_\theta^t$ is most often modeled using a convolutional neural network that estimates a label map $\hat{y}$ given an input image $x$. By construction, $f_\theta^t$ only learns to predict segmentations for task $t$.

In contrast, we learn a universal function $\hat{y} = f_\theta(x^t, S^t)$ that predicts a label map for input $x^t$ of task $t$, according to the task-specifying support $S^t = \{(x_j^t, y_j^t)\}_{j=1}^n$ comprised of example image-label pairs available for $t$.

### 3.1. Model

We implement $f_\theta$ using a fully convolutional neural network illustrated in Figure 3. We first introduce the proposed building blocks: the *cross-convolution* layer and the CrossBlock module. We then specify how we combine these blocks into a complete segmentation network.

**CrossBlock.** To transfer information between the support set and query image, we introduce a *cross-convolution* layer that interacts a query feature map $u$ with a set of support

feature maps $V = \{v_i\}_{i=1}^n$:

$$\text{CrossConv}(u, V; \theta_z) = \{z_i\}_{i=1}^n, \quad (1)$$
$$\text{for } z_i = \text{Conv}(u || v_i; \theta_z),$$

where $||$ is the concatenation operation along the feature dimension and $\text{Conv}(x; \theta_z)$ is a convolutional layer with learnable parameters $\theta_z$. Due to the weight reuse of $\theta_z$, cross-convolution operations are permutation invariant with respect to $V$. From this layer, we design a higher-level building block that produces updated versions of query representation $u$ and support $V$ at each step in the network:

$$\text{CrossBlock}(u, V; \theta_z, \theta_v) = (u', V'), \text{where:} \quad (2)$$
$$z_i = A(\text{CrossConv}(u, v_i; \theta_z)) \quad \text{for } i = 1, 2, \ldots, n$$
$$u' = 1/n \sum_{i=1}^n z_i$$
$$v_i' = A(\text{Conv}(z_i; \theta_v)) \quad \text{for } i = 1, 2, \ldots, n,$$

where $A(x)$ is a non-linear activation function. This strategy enables the representations of each support set entry and query to interact with the others through their average representation, and facilitates variably sized support sets.

**Network.** To integrate information across spatial scales, we compose the CrossBlock modules in an encoder-decoder structure with residual connections, similarly to the popular UNet architecture (Figure 3). The network takes as input the query image $x^t$ and support set $S^t = \{(x_i^t, y_i^t)\}_{i=1}^n$ of image and label-map pairs, each concatenated channel-wise, and outputs the segmentation prediction map $\hat{y}^t$.

Each level in the encoder path consists of a CrossBlock followed by a spatial down-sampling operation of both query and support set representations. Each level in the expansive path consists of up-sampling both representations, which double their spatial resolutions, concatenating them with the equivalently-sized representation in the encoding path, followed by a CrossBlock. We perform a single 1x1 convolution to map the final query representation to a prediction.

### 3.2. Training

Algorithm 1 describes UniverSeg training using a large and varied set of training tasks $\mathcal{T}$ and the loss

$$\mathcal{L}(\theta; \mathcal{T}) = \mathbb{E}_{t \in \mathcal{T}} \mathbb{E}_{(x^t, y^t), S^t} \left[ \mathcal{L}_{\text{seg}}(f_\theta(x^t, S^t), y^t) \right], \quad (3)$$

where $x^t \notin S^t$, and $\mathcal{L}_{\text{seg}}(\hat{y}, y^t)$ is a standard segmentation loss like cross-entropy or soft Dice [71], capturing the agreement between the predicted $\hat{y}$ and ground truth $y_t$.

**Data Augmentation.** We employ data augmentation to grow the diversity of training tasks and increase the number of effective training examples belonging to any particular task.

*In-Task Augmentation ($Aug_t(x, y)$).* To reduce overfitting to individual subjects, we perform standard data augmentation operations, like affine transformations, elastic deformation, or adding image noise to the query image and *each entry* of the support set independently.

---

**Algorithm 1** UniverSeg Training Loop using SGD with learning rate $\eta$ over tasks $\mathcal{T}$, main architecture $f_\theta$, in-task augmentations $\text{Aug}_t$ and task augmentations $\text{Aug}_T$

---

**for** $k = 1, \ldots, \text{NumTrainSteps}$ **do**
   $t \sim \mathcal{T}$          ▷ Sample Task
   $(x_i^t, y_i^t) \sim t$        ▷ Sample Query
   $S^t \leftarrow \{(x_j^t, y_j^t)\}_{j \neq i}^n$   ▷ Sample Support
   $x_i^t, y_i^t \leftarrow \text{Aug}_t(x_i^t, y_i^t)$   ▷ Augment Query
   $S^t \leftarrow \{\text{Aug}_t(x_j^t, y_j^t)\}_j^n$  ▷ Augment Support
   $x_i^t, y_i^t, S^t \leftarrow \text{Aug}_T(x_i^t, y_i^t, S^t)$   ▷ Task Aug
   $\hat{y}_i \leftarrow f_\theta(x_i^t, S^t)$   ▷ Predict label map
   $\ell \leftarrow \mathcal{L}_{\text{seg}}(\hat{y}_i, y_i^t)$   ▷ Compute loss
   $\theta \leftarrow \theta - \eta \nabla_\theta \ell$   ▷ Gradient step
**end for**

---

*Task Augmentation ($Aug_T(x, y, S)$).* Similar to standard data augmentation that reduces overfitting to training examples, augmenting the training *tasks* is useful for generalizing to *new tasks*, especially those far from the training task distribution. We introduce task augmentation – alterations that modify all query and support images, and/or all segmentation maps, with the same type of task-changing transformation. Example task augmentations include edge detection of the segmentation maps or a horizontal flip to all images and labels. We provide a list of all augmentations and the parameters we used in the supplemental Section C.

### 3.3. Inference

For a given query image $x^t$, UniverSeg predicts segmentation $\hat{y} = f_\theta(x^t, S^t)$ given a support set $S^t$, where the prediction quality depends on the choice of the support set $S^t$. To reduce this dependence, and to take advantage of more data when memory constraints limit the support set size at inference, we combine predictions from an ensemble of $K$ independently sampled support sets $\{S_i^t\}_{i=1}^K$ as their the pixel-wise average to produce the prediction $\hat{y} = \frac{1}{K} \sum_{k=1}^K f_\theta(x, S_k^t)$.

## 4. MegaMedical Dataset

To train our universal model $f_\theta$, we employ a set of segmentation tasks that is large and diverse, so that it is able to generalize to new tasks. We compiled MegaMedical – an extensive collection of open-access medical segmentation datasets with diverse anatomies, imaging modalities, and labels. It is constructed from 53 datasets encompassing 26 medical domains and 16 imaging modalities.

We standardize data across the wildly diverse formats of original datasets, processed images, and label maps. We also expand the training data using synthetic segmentation tasks to further increase the training task diversity. Because of individual dataset agreements, we are prohibited from re-

releasing our processed version of the datasets. Instead, we will provide data processing code to construct MegaMedical from its source datasets.

**Datasets.** MegaMedical features a wide array of biomedical domains, such as eyes [37, 58, 66, 80, 95], lungs [85, 89, 92], spine vertebrae [107], white blood cells [108], abdominal [9, 11, 32, 40, 46, 54, 55, 57, 60, 64, 65, 81, 92], and brain [4, 25, 33, 52, 53, 67, 68, 69, 92], among others. Supplemental Table 3 provides a detailed list of MegaMedical datasets. Acquisition details, subject age ranges, and health conditions are different for each dataset. We provide preprocessing and data normalization details in supplemental Section A.

**Medical Image Task Creation.** While datasets in MegaMedical feature a variety of imaging tasks and label protocols, in this work we focus on the general problem of 2D binary segmentation. For datasets featuring 3D data, for each subject, we extract the 2D mid-slice of the volume along all the major axes. When multiple modalities are present, we include each modality as a new task. For datasets containing multiple segmentation labels, we create as many binary segmentation tasks as available labels. All images are resized to $128 \times 128$ pixels and intensities are normalized to the range [0,1].

**Synthetic Task Generation.** We adapt the image generation procedure involving random synthetic shapes described in SynthMorph [34] to produce a thousand synthetic tasks to be used alongside the medical tasks during training. We detail the generation process and include examples of synthetic tasks in supplemental Section D.

## 5. Experiments

We start by describing experimental details. The first set of experiments compares the performance of UniverSeg in the held-out datasets against several single-pass methods used in few-shot learning. We then report on a variety of analyses, including ablations of modeling decisions, and the effect of training task diversity, support set size, and number of examples available for a new task.

### 5.1. Experimental Setup

**Model.** We implement the network in UniverSeg (Figure 3) using an encoder with 5 CrossBlock stages and a decoder with 4 stages, with 64 output features per stage and LeakyReLU non-linearities after each convolution. We use bilinear interpolation when downsampling or upsampling.

**Data.** For each dataset $d$, we construct three disjoint splits $d = \{d_{\text{support}}, d_{\text{dev}}, d_{\text{test}}\}$ with 60%, 20%, and 20% of the subjects, respectively. Similar to dataset generalization [99], we divide the available datasets into a training set $\mathcal{D}^T$ and a held-out test set $\mathcal{D}^H$. We train models using the support and development splits of the training datasets $\{d_{\text{support}} | d \in \mathcal{D}^T\}$.

We performed model selection and hyper-parameter tuning using the development split of held-out dataset WBC, and trained models until they stopped improving in the $d_{\text{dev}}$ split, averaged across the held-out datasets. We report results using the unseen test split of the held-out datasets $\{d_{\text{test}} | d \in \mathcal{D}^H\}$. Support set image-label pairs are sampled with replacement from each dataset's support split.

For held-out datasets, we evaluated three datasets containing anatomies represented in the training datasets (ACDC [8] and SCD [81] (heart), and STARE[37] (retinal blood vessels)), and three datasets of anatomies not covered by the rest of MegaMedical (PanDental [2] (mandible), SpineWeb [107] (vertebrae), and WBC [108] (white blood cells).

**Few-Shot Baselines.** We compare UniverSeg models to three segmentation methods from the few-shot (FS) literature, since these approaches also predict the segmentation of a query image given a support set of image-label pairs, although they were designed for the low-data regime. SE-net [84] features a fully-convolutional network, squeeze-excitation blocks, and a UNet-like model architecture. ALP-Net [76] and PANet [101], employ prototypical networks that extract prototypes from their inputs to match the given query with the support set. While ALPNet also employs a self-supervised method to generate additional label maps in settings with few tasks, we omit this step since MegaMedical includes a large collection of tasks.

Unlike UniverSeg, these methods were designed to generalize to similar tasks, such as different labels in the same anatomy and image type, or different modalities for the same anatomy. To make the comparison to UniverSeg fair, we make several additions to the training and inference procedures of these baselines as described below, and chose the best performing variant of each baseline.

**Supervised Task-Specific Models.** While it is often impractical for clinical researchers to train individual networks for each task, for evaluation we train a set of task-specific networks to serve as an upper bound of supervised performance on the held-out datasets. We employ the widely-used nnUNet [39], which automatically configures the model and training pipeline based on data properties. Each model is task-specific, using the support and development splits for training and model selection, respectively. We report results on the test split.

**Evaluation.** We evaluate models on the held-out datasets $\mathcal{D}^H$ using the test split for query images and the support split for support-sets. For all methods, unless specified otherwise, we perform 5 independent predictions per test subject using randomly drawn support sets, and ensemble the predictions. We enforce that the same random support sets are used for all methods. We evaluate predictions using the Dice score [18] (0 - 100, 0=no overlap, 100=perfect match), which quantifies the region overlap between two regions and is widely used
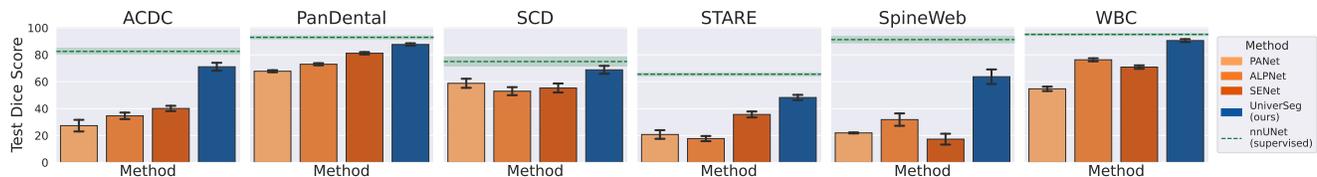
Figure 4: **Average Dice score per each held out dataset**. Performance of UniverSeg and several few-shot baselines, and the upper bound of each dataset determined by the individual fully-trained networks. For each of the unseen datasets, we average across tasks and subjects, and show the bootstrap variability in the error bars.
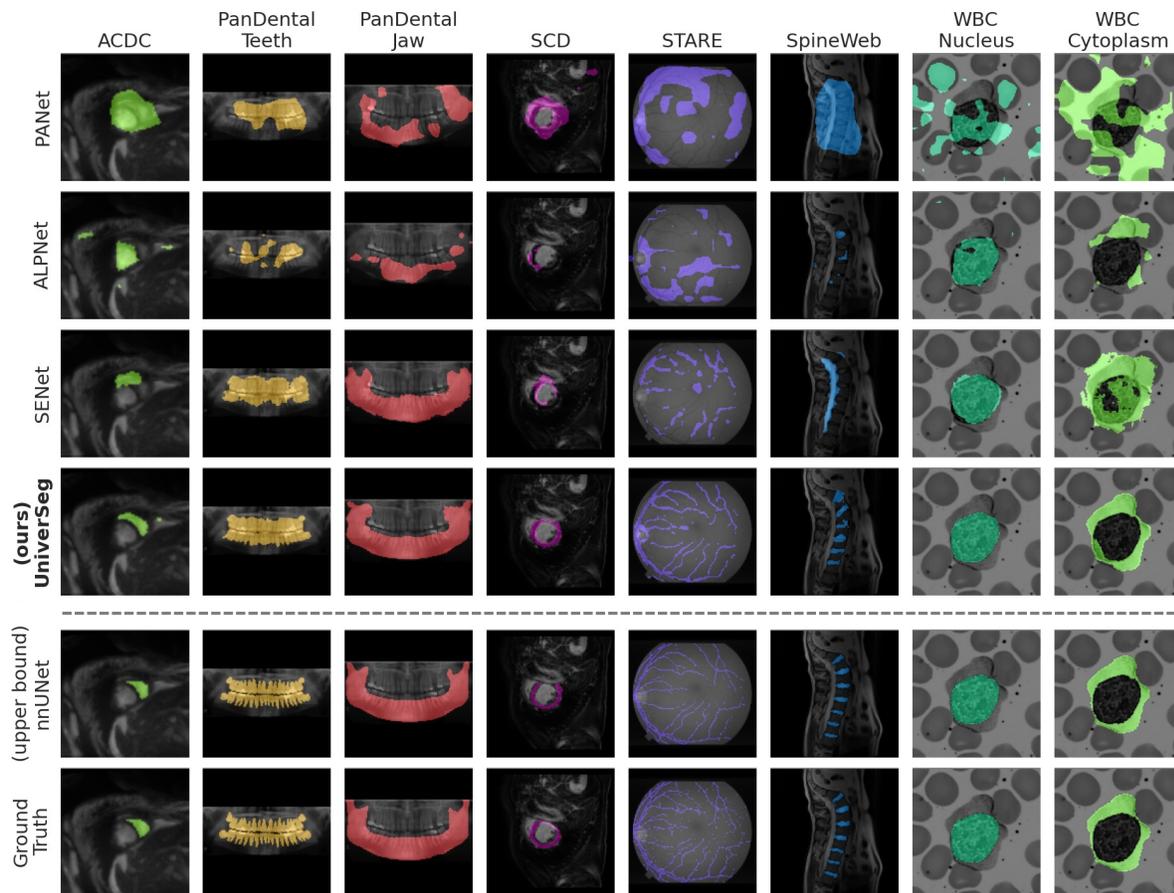


Figure 5: **Example model predictions for unseen tasks**. For a randomly sampled image per held-out task, we visualize the predictions of UniverSeg, few-shot baselines, and individually trained nnUNet models, along with ground truth maps.

in medical segmentation. For tasks with more than one label, we average Dice across all labels. For datasets with multiple tasks, we average performance across all tasks. We estimate prediction variability using subject bootstrapping, with 1,000 independent repetitions. At each repetition, we treat each task independently, sampling subjects with replacement, and report the standard deviation across bootstrapped estimates.

**Training**. We train networks with the Adam optimizer [50] and soft Dice loss [71, 96]. For the ALPNet and PANet baselines, we add an additional prototypical loss term as

described in their original works. Models trained with cross-entropy performed substantially worse than soft Dice.

While the original baseline methods were not introduced with significant data augmentation, we trained all UniverSeg and FS models with and without the proposed augmentation transformations, and report results on the best-performing setting. Unless specified otherwise, models are trained using a support size of 64. While the baselines were originally designed with small support sizes (1 or 5) as they tackled the few-shot setting, we found that training and evaluating them with larger support sizes improved their performance.

| Model | #Params | Runtime ms | Dice Score |
|---|---|---|---|
| PANet | 14.71 | $240.0 \pm 1.8$ | $41.8 \pm 1.3$ |
| ALPNet | 43.02 | $527.7 \pm 8.7$ | $47.8 \pm 1.1$ |
| SENet | 0.92 | $4.1 \pm 0.8$ | $50.1 \pm 1.3$ |
| UniverSeg (ours) | 1.18 | $142.0 \pm 0.4$ | $\mathbf{71.8 \pm 0.9}$ |
| nnUNet (sup.) | $17\times\ 1.87$ | $17\times\ 1.4{\cdot}10^{7}$ | $84.4 \pm 1.0$ |

Table 1: **Performance Summary**. For UniverSeg and each FS baseline we report model size (in millions), inference run-time, and average held-out Dice score (with bootstrapping standard deviation) . As an upper bound, we include the set of 17 individually trained task-specific nnUNets for the 6 held-out datasets, where their run-time is their cumulative required training time.

**Implementation**. We provide additional implementation and experimental details in supplemental Section B. Code and pre-trained model weights for UniverSeg are available at https://universeg.csail.mit.edu.

## 5.2. Task Generalization Results

First, we compare the segmentation quality of UniverSeg with FS baselines and the task-specific upper bounds. Our primary goal is to assess the effectiveness of UniverSeg in solving tasks from unseen datasets. Figure 4 presents the average Dice scores per dataset for each method, and Figure 5 presents example segmentation results for each method and dataset.

**Few-shot methods**. UniverSeg significantly outperforms all FS methods in all held-out datasets. For each FS method, we report the best-performing model, which involved adding components of the UniverSeg training pipeline. In the supplemental material, we show that few-shot methods perform worse when trained with a support set size of 1 and without ensembling, as they were originally introduced.

UniverSeg outperforms the highest performing baseline for all datasets with Dice improvements ranging from 7.3 to 34.9. Figure 5 also shows clear qualitative improvements in the predicted segmentations. Given the similarities between SENet and UniverSeg (fully convolutional UNet-like structure), these results suggest that the proposed CrossBlock is better suited to transferring spatial information from the support set to the dquery. Table 1 shows that UniverSeg also requires fewer model parameters than PANet, ALPNet, and the nnUNets, and a similar number to SENet.

**Task-specific networks**. For some datasets like PanDental or WBC, UniverSeg performs competitively with the supervised task-specific networks, which were extensively trained on each of the held-out tasks, and are unfeasible to run in many clinical research settings. Moreover, from the qualitative results of Figure 5, we observe that segmenta-
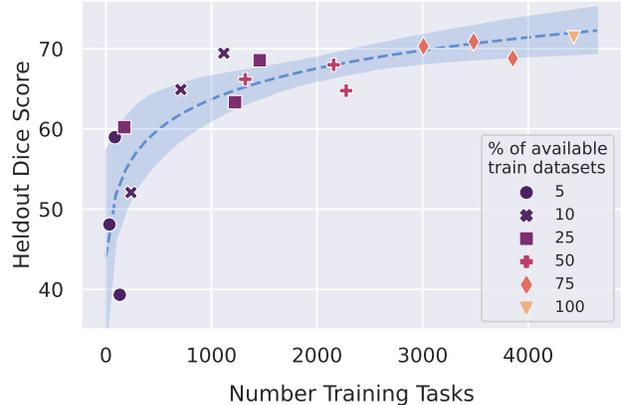


Figure 6: **Average held-out Dice versus the number of training tasks**. Points represent individual UniverSeg networks trained on a percentage of available training datasets and shown in terms of the number of underlying training tasks. In blue, we report a logarithmic fit to the data and 95% confidence intervals obtained by bootstrapped fits.

tions produced by UniverSeg more closely match those of the supervised baselines than those of any other few-shot segmentation task, especially in challenging datasets like SpineWeb or STARE.

## 5.3. Analysis

We analyze how several of the data, model, and training decisions affect the performance of UniverSeg.

**Task Quantity and Diversity.** We study the effect of the number of datasets and individual tasks used for training UniverSeg. We leave out synthetic tasks for this experiment, and train models on random subsets of the MegaMedical training datasets.

Figure 6 presents performance on the held-out datasets for different random subsets of training datasets. We find that having more training tasks improves the performance on held-out tasks. In some scenarios, the *choice* of datasets has a substantial effect. For instance, for models trained with 10% of the datasets, the best model outperforms the worst one by 17.3 Dice points, and comparing those subsets we find that the best performing one was trained on a broad set of anatomies including heart, abdomen, brain, and eyes; while the least accurate model was trained on less common lesion tasks, leading to worse generalization.

**Ablation of Training Strategies**. We perform an ablation study over the three main techniques we employ for increasing data and task diversity during training: in-task augmentation, task augmentation, and synthetic tasks.

Table 2 shows that all proposed strategies lead to improvements in model performance, with the best results achieved

| Synth | Medical | In-Task | Task | Dice Score |
|:---:|:---:|:---:|:---:|:---:|
| ✓ | | | | 61.7 ± 1.5 |
| | ✓ | | | 62.7 ± 1.1 |
| ✓ | ✓ | | | 64.5 ± 1.0 |
| | ✓ | ✓ | | 67.0 ± 0.9 |
| | ✓ | | ✓ | 70.4 ± 1.3 |
| | ✓ | ✓ | ✓ | 70.0 ± 1.5 |
| ✓ | ✓ | ✓ | ✓ | **71.8 ± 0.9** |

Table 2: **Training Strategies Ablation**. Average held-out Dice for UniverSeg models trained with different combinations of the proposed techniques to increase task diversity: in-task augmentation, task augmentation, and synthetic tasks.
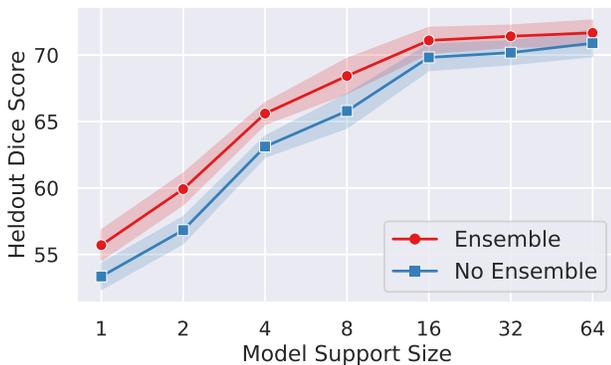


Figure 7: **Effects of support size.** Relationship between models trained at certain support sizes and their average held-out Dice score. Results improve with higher support size, with ensembling consistently helping.

when using all strategies jointly, providing a boost of 9 Dice points over no augmentations or synthetic tasks. Incorporating task augmentation leads to the largest individual improvement of 7.7 Dice points. Remarkably, the model trained using only synthetic data performs surprisingly well on the medical held-out tasks despite having never been exposed to medical training data. These results suggest that increasing image and task diversity during training, even artificially, has a substantial effect on how the model generalizes to unseen segmentation tasks.

**Support Set Size.** We study the effect of support size on models trained with support sizes $N$ from 1 to 64.

Figure 7 shows that the best results are achieved with large training support set sizes, with the average held-out Dice rapidly improving from 53.7 to 69.9 for supports sizes from 1 to 16, and then providing diminishing returns at greater support sizes, with a maximum of 71 Dice at support size 64. We find that ensembling predictions leads to consistent improvements in all cases, with greater improvements of 2.4-3.1 Dice points for small support sets ($N < 16$).
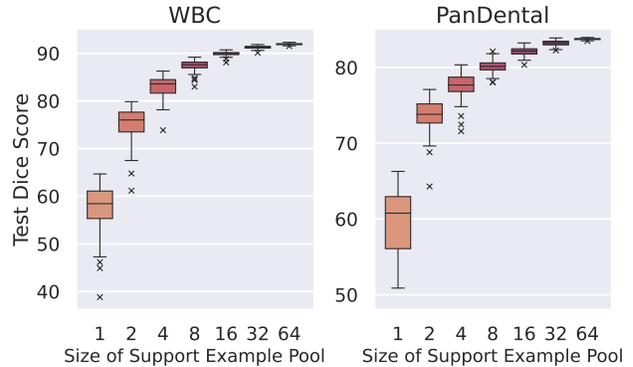


Figure 8: **Effect of available data at inference.** UniverSeg predictions using a limited $d_{support}$ example pool on the held-out WBC and PanDental datasets. For each size, we perform 100 repetitions using different random subsets.

**Limited Example Data**. Since manually annotating examples from new tasks is expensive for medical data, we investigate how the number of labeled images affects the performance of UniverSeg. We study UniverSeg when using a limited amount of labeled examples $N$ at inference, for $N = 1, 2, \ldots, 64$. We perform 100 repetitions for each size, each corresponding to an independent random subset of the data. Here, the support set contains all available data for inference, and thus we do not perform ensembling.

Figure 8 presents results for the WBC and PanDental held-out datasets, which have 108 and 116 examples in their $d_{support}$ splits respectively. For small values of support size $N$, we observe a large variance caused by very diverse support sets. As $N$ increases, we observe that average segmentation quality monotonically improves and the variance from the sample of available data examples is greatly reduced. We include analogous figures in the supplement for the other held-out datasets, where we find similar trends.

**Support Set Ensembling.** We study the effect of varying the support size $N$ at inference, and number $K$ of predictions being ensembled. We first sample 100 independent support sets for each inference support size $N$. Then, for each ensembling amount $K$, we compute ensembled predictions by averaging $K$ independently drawn predictions.

Figure 9 shows that given a certain support size, increasing the ensemble size leads to monotonic improvements and reduced variance, likely by being less dependent on the specific examples in the support set. The performance also monotonically improves with increased support size $N$, which has a significantly larger effect on segmentation accuracy than increasing the ensemble size. For instance, non-ensembled predictions with support size 64 ($N = 64$, $K = 1$) are better than heavily ensembled predictions with smaller support sizes ($N = 2, 4, 8$ and $K = 64$), even though the latter uses more support examples. This suggests
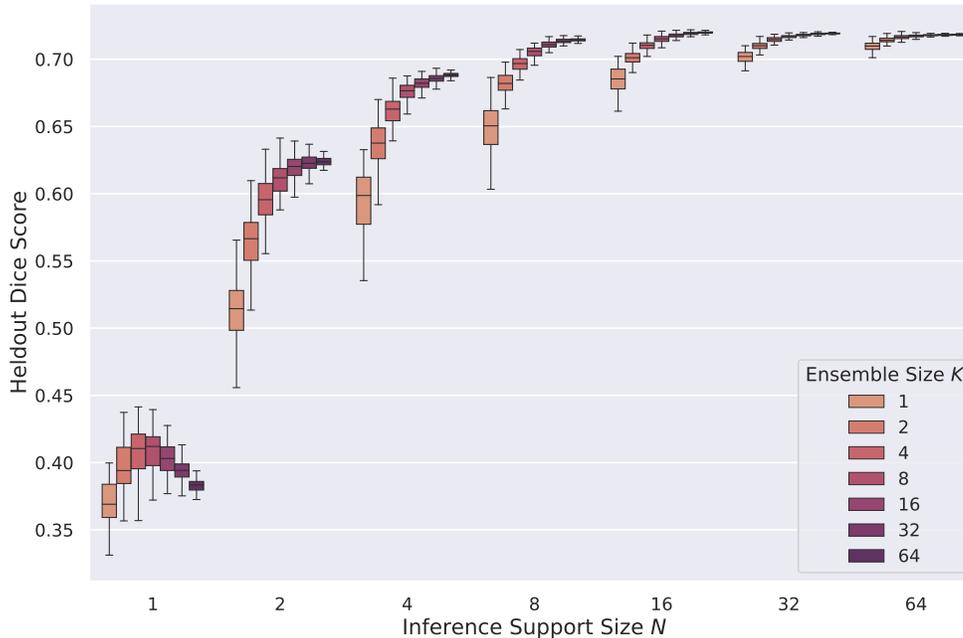
Figure 9: **Ensembling predictions at different inference support sizes.** Average held-out test Dice Score for different settings of ensembling and support size. For each inference support size $N$, we report the results (in average held-out Dice Score) of taking 100 predictions ($K = 1$) and ensembling by averaging in groups of size $K$, performing 100 repetitions for each $K$. The value boxes report quantiles over the 100 values for each setting and find that increasing either $K$ or $N$ leads to improved model performance, with $N$ having a significantly larger effect than $K$.

that UniverSeg models exploit information coming from the support examples in a fundamentally different way than existing ensembling techniques used in FS learning.

## 6. Discussion and Conclusion

We introduce UniverSeg, an approach for learning a *single* task-agnostic model for medical image segmentation. We use a large and diverse collection of open-access medical segmentation datasets to train UniverSeg, which is capable of generalizing to unseen anatomies and tasks. We introduce a novel *cross-convolution* operation that interacts the query and support representations at different scales.

In our experiments, UniverSeg substantially outperforms existing few-shot methods in all held-out datasets. Through extensive ablation studies, we conclude that UniverSeg performance is strongly dependent on task diversity during training and support set diversity during inference. This highlights the utility of UniverSeg facilitating variably-sized support sets, enabling flexibility to potential users' datasets.

**Limitations.** In this work, we focused on demonstrating and thoroughly analyzing the core idea of UniverSeg, using 2D data and single labels. We are excited by future extensions to segment 3D volumes using 2.5D or 3D models and multi-label maps, and further closing the gap with the upper bounds.

**Outlook.** UniverSeg promises to easily adapt to new segmentation tasks determined by scientists and clinical researchers, without model retraining that is often impractical for them.

## References

[1] Thyroid ultrasound cine-clip, Oct 2021. 16

[2] Amir Hossein Abdi, Shohreh Kasaei, and Mojdeh Mehdizadeh. Automatic segmentation of mandible in panoramic x-ray. *Journal of Medical Imaging*, 2(4):044003, 2015. 5, 16

[3] Laith Alzubaidi, Muthana Al-Amidie, Ahmed Al-Asadi, Amjad J. Humaidi, Omran Al-Shamma, Mohammed A. Fadhel, Jinglan Zhang, J. Santamaría, and Ye Duan. Novel transfer learning approach for medical imaging with limited labeled data. *Cancers*, 13(7), 2021. 3

[4] 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. 5, 16

[5] Spyridon Bakas, Hamed Akbari, Aristeidis Sotiras, Michel Bilello, Martin Rozycki, Justin S Kirby, John B Freymann, Keyvan Farahani, and Christos Davatzikos. Advancing the cancer genome atlas glioma mri collections with expert

segmentation labels and radiomic features. *Scientific data*, 4(1):1–13, 2017. 16

[6] Guha Balakrishnan, Amy Zhao, Mert R Sabuncu, John Guttag, and Adrian V Dalca. Voxelmorph: a learning framework for deformable medical image registration. *IEEE transactions on medical imaging*, 38(8):1788–1800, 2019. 15

[7] Sophia Bano, Francisco Vasconcelos, Luke M Shepherd, Emmanuel Vander Poorten, Tom Vercauteren, Sebastien Ourselin, Anna L David, Jan Deprest, and Danail Stoyanov. Deep placental vessel segmentation for fetoscopic mosaicking. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part III 23*, pages 763–773. Springer, 2020. 16

[8] Olivier Bernard, Alain Lalande, Clement Zotti, Frederick Cervenansky, Xin Yang, Pheng-Ann Heng, Irem Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, et al. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE transactions on medical imaging*, 37(11):2514–2525, 2018. 5, 16

[9] 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. 5, 16

[10] Benjamin Billot, Douglas Greve, Koen Van Leemput, Bruce Fischl, Juan Eugenio Iglesias, and Adrian V Dalca. A learning strategy for contrast-agnostic mri segmentation. *arXiv preprint arXiv:2003.01995*, 2020. 20

[11] Nicholas Bloch, Anant Madabhushi, Henkjan Huisman, John Freymann, Justin Kirby, Michael Grauer, Andinet Enquobahrie, Carl Jaffe, Larry Clarke, and Keyvan Farahani. Nci-isbi 2013 challenge: automated segmentation of prostate structures. *The Cancer Imaging Archive*, 370(6):5, 2015. 5, 16

[12] Mateusz Buda, Ashirbani Saha, and Maciej A Mazurowski. Association of genomic subtypes of lower-grade gliomas with shape features automatically extracted by a deep learning algorithm. *Computers in biology and medicine*, 109:218–225, 2019. 16

[13] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997. 2

[14] 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. 1

[15] Adrian V Dalca, John Guttag, and Mert R Sabuncu. Anatomical priors in convolutional networks for unsupervised biomedical segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9290–9299, 2018. 16

[16] Etienne Decenciere, Guy Cazuguel, Xiwei Zhang, Guillaume Thibault, J-C Klein, Fernand Meyer, Beatriz Marcotegui, Gwénolé Quellec, Mathieu Lamard, Ronan Danno, et al. Teleophta: Machine learning and image processing methods for teleophthalmology. *Irbm*, 34(2):196–203, 2013. 16

[17] Aysen Degerli, Morteza Zabihi, Serkan Kiranyaz, Tahir Hamid, Rashid Mazhar, Ridha Hamila, and Moncef Gabbouj. Early detection of myocardial infarction in low-quality echocardiography. *IEEE Access*, 9:34442–34453, 2021. 16

[18] Lee R Dice. Measures of the amount of ecologic association between species. *Ecology*, 26(3):297–302, 1945. 5, 17

[19] Hao Ding, Changchang Sun, Hao Tang, Dawen Cai, and Yan Yan. Few-shot medical image segmentation with cycle-resemblance attention. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2488–2497, 2023. 3

[20] Jose Dolz, Karthik Gopinath, Jing Yuan, Herve Lombaert, Christian Desrosiers, and Ismail Ben Ayed. Hyperdensenet: A hyper-densely connected cnn for multi-modal image segmentation. *IEEE Transactions on Medical Imaging*, 38(5):1116–1126, 2019. 1, 2

[21] Theodoros Evgeniou and Massimiliano Pontil. Regularized multi–task learning. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 109–117, 2004. 2

[22] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017. 3

[23] Bruce Fischl. Freesurfer. *Neuroimage*, 62(2):774–781, 2012. 15

[24] Mohsen Ghafoorian, Alireza Mehrtash, Tina Kapur, Nico Karssemeijer, Elena Marchiori, Mehran Pesteie, Charles RG Guttmann, Frank-Erik de Leeuw, Clare M Tempany, Bram van Ginneken, et al. Transfer learning for domain adaptation in mri: Application in brain lesion segmentation. In *International conference on medical image computing and computer-assisted intervention*, pages 516–524. Springer, 2017. 3

[25] Randy L Gollub, Jody M Shoemaker, Margaret D King, Tonya White, Stefan Ehrlich, Scott R Sponheim, Vincent P Clark, Jessica A Turner, Bryon A Mueller, Vince Magnotta, et al. The mcic collection: a shared repository of multi-modal, multi-site brain image data from a clinical investigation of schizophrenia. *Neuroinformatics*, 11:367–388, 2013. 5, 15, 16

[26] Ioannis S Gousias, A David Edwards, Mary A Rutherford, Serena J Counsell, Jo V Hajnal, Daniel Rueckert, and Alexander Hammers. Magnetic resonance imaging of the newborn brain: manual segmentation of labelled atlases in term-born and preterm infants. *Neuroimage*, 62(3):1499–1509, 2012. 16

[27] Ioannis S Gousias, Daniel Rueckert, Rolf A Heckemann, Leigh E Dyet, James P Boardman, A David Edwards, and Alexander Hammers. Automatic segmentation of brain mris of 2-year-olds into 83 regions of interest. *Neuroimage*, 40(2):672–684, 2008. 16

[28] Simon Graham, Quoc Dang Vu, Shan E Ahmed Raza, Ayesha Azam, Yee Wah Tsang, Jin Tae Kwak, and Nasir Rajpoot. Hover-net: Simultaneous segmentation and classi-

fication of nuclei in multi-tissue histology images. *Medical Image Analysis*, 58:101563, 2019. 16

[29] Stine Hansen, Srishti Gautam, Robert Jenssen, and Michael Kampffmeyer. Anomaly detection-inspired few-shot medical image segmentation through self-supervision with supervoxels. *Medical Image Analysis*, 78:102385, 2022. 3

[30] Kaiming He, Ross Girshick, and Piotr Dollár. Rethinking imagenet pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4918–4927, 2019. 3

[31] 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. 17

[32] 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. 5, 16

[33] Moritz R Hernandez Petzsche, Ezequiel de la Rosa, Uta Hanning, Roland Wiest, Waldo Valenzuela, Mauricio Reyes, Maria Meyer, Sook-Lei Liew, Florian Kofler, Ivan Ezhov, et al. Isles 2022: A multi-center magnetic resonance imaging stroke lesion segmentation dataset. *Scientific data*, 9(1):762, 2022. 5, 16

[34] Malte Hoffmann, Benjamin Billot, Douglas N Greve, Juan Eugenio Iglesias, Bruce Fischl, and Adrian V Dalca. Synthmorph: learning contrast-invariant registration without acquired images. *IEEE Transactions on Medical Imaging*, 41(3):543–558, 2022. 5, 20

[35] Andrew Hoopes, Malte Hoffmann, Douglas N. Greve, Bruce Fischl, John Guttag, and Adrian V. Dalca. Learning the effect of registration hyperparameters with hypermorph. volume 1, pages 1–30, 2022. 15, 16

[36] Andrew Hoopes, Jocelyn S Mora, Adrian V Dalca, Bruce Fischl, and Malte Hoffmann. Synthstrip: skull-stripping for any brain image. *NeuroImage*, 260:119474, 2022. 20

[37] AD Hoover, Valentina Kouznetsova, and Michael Goldbaum. Locating blood vessels in retinal images by piecewise threshold probing of a matched filter response. *IEEE Transactions on Medical imaging*, 19(3):203–210, 2000. 5, 16

[38] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2261–2269, 2017. 1, 2

[39] 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. 1, 2, 5

[40] 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. 5, 16

[41] Zhexin Jiang, Hao Zhang, Yi Wang, and Seok-Bum Ko. Retinal blood vessel segmentation using fully convolutional

[42] Konstantinos Kamnitsas, Enzo Ferrante, Sarah Parisot, Christian Ledig, Aditya V Nori, Antonio Criminisi, Daniel Rueckert, and Ben Glocker. Deepmedic for brain tumor segmentation. In *International workshop on Brainlesion: Glioma, multiple sclerosis, stroke and traumatic brain injuries*, pages 138–149. Springer, 2016. 1

[43] Konstantinos Kamnitsas, Christian Ledig, Virginia FJ Newcombe, Joanna P Simpson, Andrew D Kane, David K Menon, Daniel Rueckert, and Ben Glocker. Efficient multiscale 3d cnn with fully connected crf for accurate brain lesion segmentation. *Medical image analysis*, 36:61–78, 2017. 2

[44] Rashed Karim, R James Housden, Mayuragoban Balasubramaniam, Zhong Chen, Daniel Perry, Ayesha Uddin, Yosra Al-Beyatti, Ebrahim Palkhi, Prince Acheampong, Samantha Obom, et al. Evaluation of current algorithms for segmentation of scar tissue from late gadolinium enhancement cardiovascular magnetic resonance of the left atrium: an open-access grand challenge. *Journal of Cardiovascular Magnetic Resonance*, 15(1):1–17, 2013. 16

[45] A. Emre Kavur, N. Sinem Gezer, Mustafa Barış, Sinem Aslan, Pierre-Henri Conze, Vladimir Groza, Duc Duy Pham, Soumick Chatterjee, Philipp Ernst, Savaş Özkan, Bora Baydar, Dmitry Lachinov, Shuo Han, Josef Pauli, Fabian Isensee, Matthias Perkonigg, Rachana Sathish, Ronnie Rajan, Debdoot Sheet, Gurbandurdy Dovletov, Oliver Speck, Andreas Nürnberger, Klaus H. Maier-Hein, Gözde Bozdağı Akar, Gözde Ünal, Oğuz Dicle, and M. Alper Selver. CHAOS Challenge - combined (CT-MR) healthy abdominal organ segmentation. *Medical Image Analysis*, 69:101950, Apr. 2021. 16

[46] Ali Emre Kavur, M. Alper Selver, Oğuz Dicle, Mustafa Barış, and N. Sinem Gezer. CHAOS - Combined (CT-MR) Healthy Abdominal Organ Segmentation Challenge Data. Apr. 2019. 5, 15

[47] Ali Emre Kavur, M. Alper Selver, Oğuz Dicle, Mustafa Barış, and N. Sinem Gezer. CHAOS - Combined (CT-MR) Healthy Abdominal Organ Segmentation Challenge Data, Apr. 2019. 16

[48] Rabindra Khadka, Debesh Jha, Steven Hicks, Vajira Thambawita, Michael A Riegler, Sharib Ali, and Pål Halvorsen. Meta-learning with implicit gradients in a few-shot setting for medical image segmentation. *Computers in Biology and Medicine*, 143:105227, 2022. 3

[49] Pulkit Khandelwal and Paul Yushkevich. Domain generalizer: A few-shot meta learning framework for domain generalization in medical imaging. In *Domain Adaptation and Representation Transfer, and Distributed and Collaborative Learning*, pages 73–84. Springer, 2020. 3

[50] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6, 17

[51] Serkan Kiranyaz, Aysen Degerli, Tahir Hamid, Rashid Mazhar, Rayyan El Fadil Ahmed, Rayaan Abouhasera, Morteza Zabihi, Junaid Malik, Ridha Hamila, and Moncef Gabbouj. Left ventricular wall motion estimation by

active polynomials for acute myocardial infarction detection. *IEEE Access*, 8:210301–210317, 2020. 16

[52] Hugo J Kuijf, J Matthijs Biesbroek, Jeroen De Bresser, Rutger Heinen, Simon Andermatt, Mariana Bento, Matt Berseth, Mikhail Belyaev, M Jorge Cardoso, Adria Casamitjana, et al. Standardized assessment of automatic segmentation of white matter hyperintensities and results of the wmh segmentation challenge. *IEEE transactions on medical imaging*, 38(11):2556–2568, 2019. 5, 16

[53] Maria Kuklisova-Murgasova, Paul Aljabar, Latha Srinivasan, Serena J Counsell, Valentina Doria, Ahmed Serag, Ioannis S Gousias, James P Boardman, Mary A Rutherford, A David Edwards, et al. A dynamic 4d probabilistic atlas of the developing brain. *NeuroImage*, 54(4):2750–2763, 2011. 5, 16

[54] 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. 5, 16

[55] 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. 5, 16

[56] Sarah Leclerc, Erik Smistad, Joao Pedrosa, Andreas Østvik, Frederic Cervenansky, Florian Espinosa, Torvald Espeland, Erik Andreas Rye Berg, Pierre-Marc Jodoin, Thomas Grenier, et al. Deep learning for segmentation using an open large-scale dataset in 2d echocardiography. *IEEE transactions on medical imaging*, 38(9):2198–2210, 2019. 16

[57] Guillaume Lemaître, Robert Martí, Jordi Freixenet, Joan C Vilanova, Paul M Walker, and Fabrice Meriaudeau. Computer-aided detection and diagnosis for prostate cancer based on mono and multi-parametric mri: a review. *Computers in biology and medicine*, 60:8–31, 2015. 5, 16

[58] Mingchao Li, Yuhan Zhang, Zexuan Ji, Keren Xie, Songtao Yuan, Qinghuai Liu, and Qiang Chen. Ipn-v2 and octa-500: Methodology and dataset for retinal image segmentation. *arXiv preprint arXiv:2012.07261*, 2020. 5, 16

[59] Yiwen Li, Yunguan Fu, Iani Gayo, Qianye Yang, Zhe Min, Shaheer Saeed, Wen Yan, Yipei Wang, J Alison Noble, Mark Emberton, et al. Prototypical few-shot segmentation for cross-institution male pelvic structures with spatial registration. *arXiv preprint arXiv:2209.05160*, 2022. 3

[60] Geert Litjens, Robert Toth, Wendy van de Ven, Caroline Hoeks, Sjoerd Kerkstra, Bram van Ginneken, Graham Vincent, Gwenael Guillard, Neil Birbeck, Jindang Zhang, et al. Evaluation of prostate segmentation algorithms for mri: the promise12 challenge. *Medical image analysis*, 18(2):359–373, 2014. 5, 16

[61] Quande Liu, Qi Dou, Lequan Yu, and Pheng Ann Heng. Msnet: Multi-site network for improving prostate segmentation with heterogeneous mri data. *IEEE Transactions on Medical Imaging*, 2020. 2

[62] Vebjorn Ljosa, Katherine L Sokolnicki, and Anne E Carpenter. Annotated high-throughput microscopy image sets for validation. *Nature methods*, 9(7):637–637, 2012. 16

[63] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Unsupervised domain adaptation with residual transfer networks. *Advances in neural information processing systems*, 29, 2016. 2

[64] 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. 5, 16

[65] Jun Ma, Yao Zhang, Song Gu, Cheng Zhu, Cheng Ge, Yichi Zhang, Xingle An, Congcong Wang, Qiyuan Wang, Xin Liu, Shucheng Cao, Qi Zhang, Shangqing Liu, Yunpeng Wang, Yuhui Li, Jian He, and Xiaoping Yang. Abdomenct-1k: Is abdominal organ segmentation a solved problem? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 5, 16

[66] Yuhui Ma, Huaying Hao, Jianyang Xie, Huazhu Fu, Jiong Zhang, Jianlong Yang, Zhen Wang, Jiang Liu, Yalin Zheng, and Yitian Zhao. Rose: a retinal oct-angiography vessel segmentation dataset and new model. *IEEE Transactions on Medical Imaging*, 40(3):928–939, 2021. 5, 16

[67] Daniel S Marcus, Tracy H Wang, Jamie Parker, John G Csernansky, John C Morris, and Randy L Buckner. Open access series of imaging studies (oasis): cross-sectional mri data in young, middle aged, nondemented, and demented older adults. *Journal of cognitive neuroscience*, 19(9):1498–1507, 2007. 5, 15, 16

[68] Kenneth Marek, Danna Jennings, Shirley Lasch, Andrew Siderowf, Caroline Tanner, Tanya Simuni, Chris Coffey, Karl Kieburtz, Emily Flagg, Sohini Chowdhury, et al. The parkinson progression marker initiative (ppmi). *Progress in neurobiology*, 95(4):629–635, 2011. 5, 15, 16

[69] Maciej A Mazurowski, Kal Clark, Nicholas M Czarnek, Parisa Shamsesfandabadi, Katherine B Peters, and Ashirbani Saha. Radiogenomics of lower-grade glioma: algorithmically-assessed tumor shape is associated with tumor genomic subtypes and patient outcomes in a multi-institutional study with the cancer genome atlas data. *Journal of neuro-oncology*, 133:27–35, 2017. 5, 16

[70] Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions on medical imaging*, 34(10):1993–2024, 2014. 16

[71] 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. 1, 4, 6

[72] Pim Moeskops, Jelmer M Wolterink, Bas HM van der Velden, Kenneth GA Gilhuijs, Tim Leiner, Max A Viergever, and Ivana Išgum. Deep learning for multi-task medical image segmentation in multiple modalities. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 478–486. Springer, 2016. 2

[73] Fernando Navarro, Suprosanna Shit, Ivan Ezhov, Johannes Paetzold, Andrei Gafita, Jan C Peeken, Stephanie E Combs,

and Bjoern H Menze. Shape-aware complementary-task learning for multi-organ segmentation. In *International Workshop on Machine Learning in Medical Imaging*, pages 620–627. Springer, 2019. 2

[74] Khoi Nguyen and Sinisa Todorovic. Feature weighting and boosting for few-shot segmentation. *CoRR*, abs/1909.13140, 2019. 3

[75] Alex Nichol and John Schulman. Reptile: a scalable met-alearning algorithm. *arXiv preprint arXiv:1803.02999*, 2(3):4, 2018. 3

[76] Cheng Ouyang, Carlo Biffi, Chen Chen, Turkay Kart, Huaqi Qiu, and Daniel Rueckert. Self-supervision with superpix-els: Training few-shot medical image segmentation without annotation. In *European Conference on Computer Vision*, pages 762–780. Springer, 2020. 3, 5

[77] Prashant Pandey, Mustafa Chasmai, Tanuj Sur, and Brejesh Lall. Robust prototypical few-shot organ segmentation with regularized neural-odes. *arXiv preprint arXiv:2208.12428*, 2022. 3

[78] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 17

[79] Kelly Payette, Priscille de Dumast, Hamza Kebiri, Ivan Ezhov, Johannes C Paetzold, Suprosanna Shit, Asim Iqbal, Romesa Khan, Raimund Kottke, Patrice Grehten, et al. An automatic multi-tissue human fetal brain segmentation benchmark using the fetal tissue annotation dataset. *Scientific Data*, 8(1):1–14, 2021. 16

[80] Prasanna Porwal, Samiksha Pachade, Ravi Kamble, Manesh Kokare, Girish Deshmukh, Vivek Sahasrabuddhe, and Fabrice Meriaudeau. Indian diabetic retinopathy image dataset (idrid), 2018. 5, 16

[81] Perry Radau, Yingli Lu, Kim Connelly, Gideon Paul, AJWG Dick, and Graham Wright. Evaluation framework for algorithms segmenting short axis cardiac mri. *The MIDAS Journal-Cardiac MR Left Ventricle Segmentation Challenge*, 49, 2009. 5, 16

[82] Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding transfer learning for medical imaging. *Advances in neural information processing systems*, 32, 2019. 2, 3

[83] 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*, pages 234–241. Springer, 2015. 1, 2

[84] Abhijit Guha Roy, Shayan Siddiqui, Sebastian Pölsterl, Nassir Navab, and Christian Wachinger. 'squeeze & ex-cite'guided few-shot segmentation of volumetric images. *Medical image analysis*, 59:101587, 2020. 3, 5

[85] Adriel Saporta, Xiaotong Gui, Ashwin Agrawal, Anuj Pareek, SQ Truong, CD Nguyen, Van-Doan Ngo, Jayne Seekins, Francis G Blankenberg, AY Ng, et al. Deep learning saliency maps do not accurately highlight diagnostically relevant regions for medical image interpretation. *MedRxiv*, 2021. 5, 16

[86] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018. 2

[87] Jun Seo, Young-Hyun Park, Sung Whan Yoon, and Jaekyun Moon. Task-adaptive feature transformer with semantic enrichment for few-shot segmentation. *arXiv preprint arXiv:2202.06498*, 2022. 3

[88] Ahmed Serag, Paul Aljabar, Gareth Ball, Serena J Counsell, James P Boardman, Mary A Rutherford, A David Edwards, Joseph V Hajnal, and Daniel Rueckert. Construction of a consistent high-definition spatio-temporal atlas of the developing brain using adaptive kernel regression. *Neuroimage*, 59(3):2255–2265, 2012. 16

[89] Arnaud Arindra Adiyoso Setio, Alberto Traverso, Thomas De Bel, Moira SN Berens, Cas Van Den Bogaard, Piergiorgio Cerello, Hao Chen, Qi Dou, Maria Evelina Fantacci, Bram Geurts, et al. Validation, comparison, and combination of algorithms for automatic detection of pulmonary nodules in computed tomography images: the luna16 challenge. *Medical image analysis*, 42:1–13, 2017. 5, 16

[90] Neeraj Sharma and Lalit M Aggarwal. Automated medical image segmentation techniques. *Journal of medical physics/Association of Medical Physicists of India*, 35(1):3, 2010. 1

[91] Qianqian Shen, Yanan Li, Jiyong Jin, and Bin Liu. Q-net: Query-informed few-shot medical image segmentation. *arXiv preprint arXiv:2208.11451*, 2022. 3

[92] Amber L Simpson, Michela Antonelli, Spyridon Bakas, Michel Bilello, Keyvan Farahani, Bram Van Ginneken, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, et al. A large annotated medical image dataset for the development and evaluation of segmentation algorithms. *arXiv preprint arXiv:1902.09063*, 2019. 5, 16

[93] Rishav Singh, Vandana Bharti, Vishal Purohit, Abhinav Kumar, Amit Kumar Singh, and Sanjay Kumar Singh. Metamed: Few-shot medical image classification using gradient-based meta-learning. *Pattern Recognition*, 120:108111, 2021. 3

[94] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017. 3

[95] Joes Staal, Michael D Abràmoff, Meindert Niemeijer, Max A Viergever, and Bram Van Ginneken. Ridge-based vessel segmentation in color images of the retina. *IEEE transactions on medical imaging*, 23(4):501–509, 2004. 5, 15, 16

[96] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: Third International Workshop, DLMIA 2017, and 7th International Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, Proceedings 3*, pages 240–248. Springer, 2017. 6

[97] Baochen Sun, Jiashi Feng, and Kate Saenko. Return of frustratingly easy domain adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, 2016. 2

[98] Hao Tang, Xingwei Liu, Shanlin Sun, Xiangyi Yan, and Xiaohui Xie. Recurrent mask refinement for few-shot medical image segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3918–3928, 2021. 3

[99] Eleni Triantafillou, Hugo Larochelle, Richard Zemel, and Vincent Dumoulin. Learning a universal template for few-shot dataset generalization. In *International Conference on Machine Learning*, pages 10424–10433. PMLR, 2021. 5

[100] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29, 2016. 3

[101] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. *CoRR*, abs/1908.06391, 2019. 3, 5

[102] Chi Zhang, Guosheng Lin, Fayao Liu, Rui Yao, and Chunhua Shen. Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5217–5226, 2019. 3

[103] Penghao Zhang, Jiayue Li, Yining Wang, and Judong Pan. Domain adaptation for medical image segmentation: a meta-learning method. *Journal of Imaging*, 7(2):31, 2021. 3

[104] Yingtao Zhang, Min Xian, Heng-Da Cheng, Bryar Shareef, Jianrui Ding, Fei Xu, Kuan Huang, Boyu Zhang, Chunping Ning, and Ying Wang. Busis: A benchmark for breast ultrasound image segmentation. In *Healthcare*, volume 10, page 729. MDPI, 2022. 16

[105] Yao Zhang, Jiawei Yang, Jiang Tian, Zhongchao Shi, Cheng Zhong, Yang Zhang, and Zhiqiang He. Modality-aware mutual learning for multi-modal medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 589–599. Springer, 2021. 3

[106] Amy Zhao, Guha Balakrishnan, Fredo Durand, John V Guttag, and Adrian V Dalca. Data augmentation using learned transformations for one-shot medical image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8543–8553, 2019. 3

[107] Guoyan Zheng, Chengwen Chu, Daniel L Belavỳ, Bulat Ibragimov, Robert Korez, Tomaž Vrtovec, Hugo Hutt, Richard Everson, Judith Meakin, Isabel Lŏpez Andrade, et al. Evaluation and comparison of 3d intervertebral disc localization and segmentation methods for 3d t2 mr data: A grand challenge. *Medical image analysis*, 35:327–344, 2017. 5, 15, 16

[108] Xin Zheng, Yong Wang, Guoyou Wang, and Jianguo Liu. Fast and robust segmentation of white blood cell images by self-supervised learning. *Micron*, 107:55–71, 2018. 5, 15, 16

[109] Zongwei Zhou, Vatsal Sodha, Jiaxuan Pang, Michael B Gotway, and Jianming Liang. Models genesis. *Medical image analysis*, 67:101840, 2021. 3

# A. MegaMedical

**Preprocessing.** Medical images involve large variations of voxel or pixel values. For example, MRI intensities in MegaMedical range from [0, 800], CT intensities range from [-2000, 2000], while other modalities might already be in the $[0, 1]$ range.

To normalize data across the diverse datasets, we apply several preprocessing steps for each modality. For MRI datasets, we clip the intensity to $[0.5, 99.5]$ percentiles for non-zero voxels. For CT images, we clip intensity values to the range $[-500, 1000]$. We min-max normalize all resulting volumes to $[0, 1]$ and resize them to $128 \times 128 \times 128$. From any 3D volumes, we extract two different kinds of slices: *mid-slices* and *max-slices*.

**Slicing.** For *mid-slices*, from any 3D image and label volumes we extract the middle slice along each axis, resulting in a representative $128 \times 128$ slice. This strategy avoids biasing the data toward knowing the location of labels in the scans. This is especially important for inference, where the location of the foreground label would not be known in a 3D volume.

For training, we also extract *max slices*. For each label $l$ of a dataset, we find the slice (along each axis) in each volume that contains the most voxels with that label. We extract this slice from both volume and label map and repeat this for all labels in the dataset. These slices provide additional training data, and we do not use them during evaluation.

**Label Maps.** Most datasets include label maps that were either manually obtained, or manually curated after being obtained using an automatic tool. For adult brain datasets [25, 67, 68], we follow recent large-scale analyses [6, 35] and obtain semantic sub-cortical segmentations using FreeSurfer [23].

Datasets can often contain multiple tasks – such as segmenting both lesions and anatomy – and the same task can appear in different datasets – like segmenting the hippocampus in different MRI collections. Certain labels can sometimes tackle the same anatomical region of interest but be defined differently in two different datasets. In this work, we focus on single-label, single-modality, and 2D segmentation.

**Medical Task Creation.** To create a task, the subjects of dataset $d$ can contain labels for either a particular biomedical target (e.g. eye-vessels [95], vertebrae [107], white blood cells [108]) or a set of targets (e.g. abdominal organs [46], brain regions [67]), and an imaging modality $m \in M_d$ (e.g. CT, MRI, X-Ray). If $d$ is multi-class, we split it into several single-label $l \in L_d$ tasks. If $d$ is a 3D dataset, we extract different axes $a \in A_d$ as different tasks. Following this construction, each task can be described using a unique tuple $t = (d, m, l, a)$.

Table 3: We assembled the following set of datasets to train UniverSeg. For the relative size of datasets, we have included the number of unique scans (subject and modality pairs) that each dataset has.

| Dataset Name | Description | # of Scans | Image Modalities |
|---|---|---|---|
| AbdomenCT-1K [65] | Abdominal organ segmentation (overlap with KiTS, MSD) | 361 | CT |
| ACDC [8] | Left and right ventricular endocardium | 99 | cine-MRI |
| AMOS [40] | Abdominal organ segmentation | 240 | CT, MRI |
| BBBC003 [62] | Mouse embryos | 15 | Microscopy |
| BrainDevelopment [26, 27, 53, 88] | Adult and Neonatal Brain Atlases | 53 | multi-modal MRI |
| BRATS [4, 5, 70] | Brain tumors | 6,096 | multi-modal MRI |
| BTCV [55] | Abdominal Organs | 30 | CT |
| BUS [104] | Breast tumor | 163 | Ultrasound |
| CAMUS [56] | Four-chamber and Apical two-chamber heart | 500 | Ultrasound |
| CDemris [44] | Human Left Atrial Wall | 60 | CMR |
| CHAOS [45, 47] | Abdominal organs (liver, kidneys, spleen) | 40 | CT, T2-weighted MRI |
| CheXplanation [85] | Chest X-Ray observations | 170 | X-Ray |
| CoNSeP [28] | Histopathology Nuclei | 27 | Microscopy |
| DRIVE [95] | Blood vessels in retinal images | 20 | Optical camera |
| EOphtha [16] | Eye Microaneurysms and Diabetic Retinopathy | 102 | Optical camera |
| FeTA [79] | Fetal brain structures | 80 | Fetal MRI |
| FetoPlac [7] | Placenta vessel | 6 | Fetoscopic optical camera |
| HMC-QU [17, 51] | 4-chamber (A4C) and apical 2-chamber (A2C) left wall | 292 | Ultrasound |
| I2CVB [57] | Prostate (peripheral zone, central gland) | 19 | T2-weighted MRI |
| IDRID [80] | Diabetic Retinopathy | 54 | Optical camera |
| ISLES [33] | Ischemic stroke lesion | 180 | multi-modal MRI |
| KiTS [32] | Kidney and kidney tumor | 210 | CT |
| LGGFlair [12, 69] | TCIA lower-grade glioma brain tumor | 110 | MRI |
| LiTS [9] | Liver Tumor | 131 | CT |
| LUNA [89] | Lungs | 888 | CT |
| MCIC [25] | Multi-site Brain regions of Schizophrenic patients | 390 | T1-weighted MRI |
| MSD [92] | Large-scale collection of 10 Medical Segmentation Datasets | 3,225 | CT, multi-modal MRI |
| NCI-ISBI [11] | Prostate | 30 | T2-weighted MRI |
| OASIS [35, 67] | Brain anatomy | 414 | T1-weighted MRI |
| OCTA500 [58] | Retinal vascular | 500 | OCT/OCTA |
| PanDental [2] | Mandible and Teeth | 215 | X-Ray |
| PROMISE12 [60] | Prostate | 37 | T2-weighted MRI |
| PPMI [68, 15] | Brain regions of Parkinson patients | 1,130 | T1-weighted MRI |
| ROSE [66] | Retinal vessel | 117 | OCT/OCTA |
| SCD [81] | Sunnybrook Cardiac Multi-Dataset Collection | 100 | cine-MRI |
| SegTHOR [54] | Thoracic organs (heart, trachea, esophagus) | 40 | CT |
| SpineWeb [107] | Vertebrae | 15 | T2-weighted MRI |
| STARE [37] | Blood vessels in retinal images | 20 | Optical camera |
| TUCC [1] | Thyroid nodules | 167 | Ultrasound cine-clip |
| WBC [108] | White blood cell and nucleus | 400 | Microscopy |
| WMH [52] | White matter hyper-intensities | 60 | multi-modal MRI |
| WORD [64] | Organ segmentation | 120 | CT |

# B. Additional Implementation Details

**Data Storage**. For each gradient step, a UniverSeg model needs to load $B \times (N + 1) \times 2$ images, where $B$ is the batch size, $N$ is the support size, and the factor of 2 corresponds to the combination of the image and label map. This can pose a serious challenge for traditional data loading strategies, especially as $N$ increases. Therefore, we store data samples in a highly optimized way to ensure that I/O does not bottleneck the training process, using LMDB data stores that are optimized for read-only access. Within the database, data is encoded using msgpack and compressed with the LZ4 codec for fast decompression. We find that this setup exceeds regular file-system random access by over two orders of magnitude.

**Task Sampling**. To ensure task and data heterogeneity during training, we do not sample all tasks equally. Some datasets contain substantially more tasks than others, and we aim to avoid overfitting medical domains where tasks are abundant (such as neuroimaging tasks). Instead, we perform hierarchical uniform sampling with multiple stages: dataset, subject group, acquisition modality, axis, and label. We first sample the dataset uniformly from all datasets, then sample a task among the tasks from that dataset, and so on.

**Model**. We implemented UniverSeg in PyTorch [78] and used the official implementations for the baselines (ALPNet, PANet, and SENet) and supervised network nnUNet. Based on the experimental details in the ALPNet work, we used an off-the-shelf ResNet101 [31] for both the pre-trained encoder for ALPNet and PANet. For these two methods, because their feature encoder expects three-channel inputs, we duplicate the input dimension $1 \times 128 \times 128$ three times channel-wise to get inputs of dimension $3 \times 128 \times 128$.

We efficiently perform the CrossConvolution operation by exploiting the batch dimension. Instead of performing $N$ convolutions with the same learnable parameters, we perform a single convolution by tiling the inputs along the batch dimension. We use the same strategy for the convolutions predicting the CrossBlock outputs $V'$.

**Optimization**. For all models during training, we minimize the soft Dice loss:

$$\mathcal{L}_{\text{Dice}}(y_t, \hat{y}) = \frac{2 \sum y_t \odot \hat{y}}{\sum y_t^2 + \sum \hat{y}^2}, \tag{4}$$

using a learning rate of $\eta = 10^{-4}$, the Adam optimizer[50], and a batch size of 1. We searched learning rates over the range $[10^{-5}, 10^{-2}]$ and found the best results on the validation split of the training datasets with learning rates around $10^{-4}$ and set on $10^{-4}$ for comparison and reproducibility purposes.

**Evaluation**. We evaluate predicted label maps $\hat{y}$ using the Dice score [18], which quantifies the overlap between two regions and is widely used in the segmentation literature:

$$\text{Dice}(y_t, \hat{y}) = 100 * \frac{2|y_t \cap \hat{y}|}{|y_t|^2 + |\hat{y}|^2} \tag{5}$$

where $y$ is the ground truth segmentation map and $\hat{y}$ is the predicted segmentation map. A Dice score of 100 indicates perfectly overlapping regions, while 0 indicates no overlap.

**Task-Specific Networks** The nnUNet framework trains 5 networks per task using multiple folds of the support data for training, and ensemble their predictions at inference. We apply the nnUNet framework independently for each held-out task, which corresponds to a set of subjects and the segmentation labels for a particular binary task.

We also designed and trained additional individual U-Net networks. For the majority of the tasks, we found the best results after searching batch sizes and augmentation policies. We omitted these as we found that the nnUNets performed very similarly.

# C. Data Augmentation

During UniverSeg training, we found that using substantial data augmentation was important. Augmentation techniques enable UniverSeg to see effectively both a greater diversity of tasks as well as a greater number of examples of each. We separate these two kinds of augmentations into *Task* and *In-Task*.

In Table 4, we detail included augmentations. During model development, we experimented to find the hyperparameters which worked best for each kind of augmentation. Several augmentations are repeated (although with different parameters) across task and in-task sections of Table 4.

Table 4: **List of augmentations used in model training.**

| Augmentation | Aug Type | Parameter Details |
|---|---|---|
| Flip Intensities | Task | $\mathbf{p} = 0.50$ |
| Flip Labels | Task | $\mathbf{p} = 0.50$ |
| Horizontal/Vertical Flip | Task | $\mathbf{p} = 0.50$ |
| Sobel-Edge Label | Task | $\mathbf{p} = 0.50$ |
| Task Affine Shift | Task | $\mathbf{p} = 0.50$ **degrees** $= [0, 360]$ **translate** $= [0, 0.2]$ **scale** $= [0.8, 1.1]$ |
| Task Brightness Contrast Change | Task | $\mathbf{p} = 0.50$ **brightness** $= [-0.1, 0.1]$ **contrast** $= [0.8, 1.2]$ |
| Task Elastic Warp | Task | $\mathbf{p} = 0.25$ $\boldsymbol{\alpha} = [1, 2]$ $\boldsymbol{\sigma} = [6, 8]$ |
| Task Gaussian Blur | Task | $\mathbf{p} = 0.50$ **k-size** $= 5$ $\boldsymbol{\sigma} = [0.1, 1.1]$ |
| Task Gaussian Noise | Task | $\mathbf{p} = 0.50$ $\boldsymbol{\mu} = [0, 0.05]$ $\boldsymbol{\sigma^2} = [0, 0.05]$ |
| Task Sharpness Change | Task | $\mathbf{p} = 0.50$ **sharpness** $= 5$ |
| Example Affine Shift | In-Task | $\mathbf{p} = 0.50$ **degrees** $= [0, 360]$ **translate** $= [0, 0.2]$ **scale** $= [0.8, 1.1]$ |
| Example Brightness Contrast Change | In-Task | $\mathbf{p} = 0.25$ **brightness** $= [-0.1, 0.1]$ **contrast** $= [0.5, 1.5]$ |
| Example Gaussian Blur | In-Task | $\mathbf{p} = 0.25$ **k-size** $= 5$ $\boldsymbol{\sigma} = [0.1, 1.1]$ |
| Example Gaussian Noise | In-Task | $\mathbf{p} = 0.25$ $\boldsymbol{\mu} = [0, 0.05]$ $\boldsymbol{\sigma^2} = [0, 0.05]$ |
| Example Sharpness Change | In-Task | $\mathbf{p} = 0.25$ **sharpness** $= 5$ |
| Example Variable Elastic Warp | In-Task | $\mathbf{p} = 0.80$ $\boldsymbol{\alpha} = [1, 2.5]$ $\boldsymbol{\sigma} = [7, 8]$ |

We briefly describe each augmentation and its parameters. Each augmentation also has a parameter **p** which controls the probability that augmentation is applied at each iteration. For in-task augmentation, this probability controls whether or not **all** of the support set entries are individually augmented or not. For operations that we developed, we include examples in Figure 10.

- *Flip Intensities* (Task): Flip the intensity values for all images (query and support), but not the label maps, using 1 - image for each.

- *Flip Labels* (Task): Reverse the foreground and background in the segmentation maps.

- *Horizontal/Vertical Flip* (Task): Flip all entries in the support horizontally or vertically (all flipped in the same way).

- *Sobel-Edge Label* (Task): We propose an operation that increases the number of tasks with thin segmentation structures. We apply a Sobel filter to each label map in the x and y directions, compute the squared norm, which becomes our new label map.

- *Affine Shift* (Task, In-Task): Apply a consistent random affine transformation to all entries in the support set; **degrees** controls how much to randomly rotate, **translate** controls how far the images and labels can shift, and **scale** controls the amount of zoom.

- *Brightness Contrast Change* (Task, In-Task): Apply a random brightness and contrast change to all images; how much brightness can change is controlled by **brightness** and contrast is controlled by the parameter **contrast**.

- *Elastic Warp* (Task, In-Task): Apply a consistent elastic deformation warp to all entries in the support and to the query; $\alpha$ controls the strength of the warp and $\sigma$ controls the smoothness of the warp.

- *Gaussian Blur* (Task, In-Task): Apply a convolutional Gaussian blur to each image in the support set and the query with a certain kernel size, **k-size**, and standard-deviation $\sigma$.

Figure 10: **Example augmentation operations applied to the WBC Dataset.** We visualize several examples of unique task augmentations we apply during training.

- *Gaussian Noise* (Task, In-Task): Apply Gaussian noise to all images in the support set and query with mean $\mu$ and variance $\sigma^2$.

- *Sharpness Change* (Task, In-Task): Apply a sharpness filter to the images (query and support), where the sharpness strength is controlled by **sharpness**.

## D. Synthetic Tasks

We found improvement in held-out performance by introducing synthetic tasks during training, building on recent methods that use synthetic medical images to solve specific tasks [10, 34, 36], especially the synthetic shapes in SynthMorph [34]. We generate 1,000 new tasks with high diversity (Figure 12). As shown in Figure 11, for each task, we first synthesize a label map of 16 random shapes, representing 16 regions of interest. We deform this label map with 100 random smooth deformation fields, representing 100 subjects with the same simulated anatomy. We then add texture to the resulting images by filling in each region of interest with slightly varied intensities around a sampled mean and adding Gaussian and Perlin noise.



Figure 11: **Generation process for synthetic tasks.** For a new synthetic task, we first generate random shapes to obtain a label map, then synthesize 100 spatial variations on this label map, and finally synthesize resulting intensity images. We repeat this process for 1000 tasks.
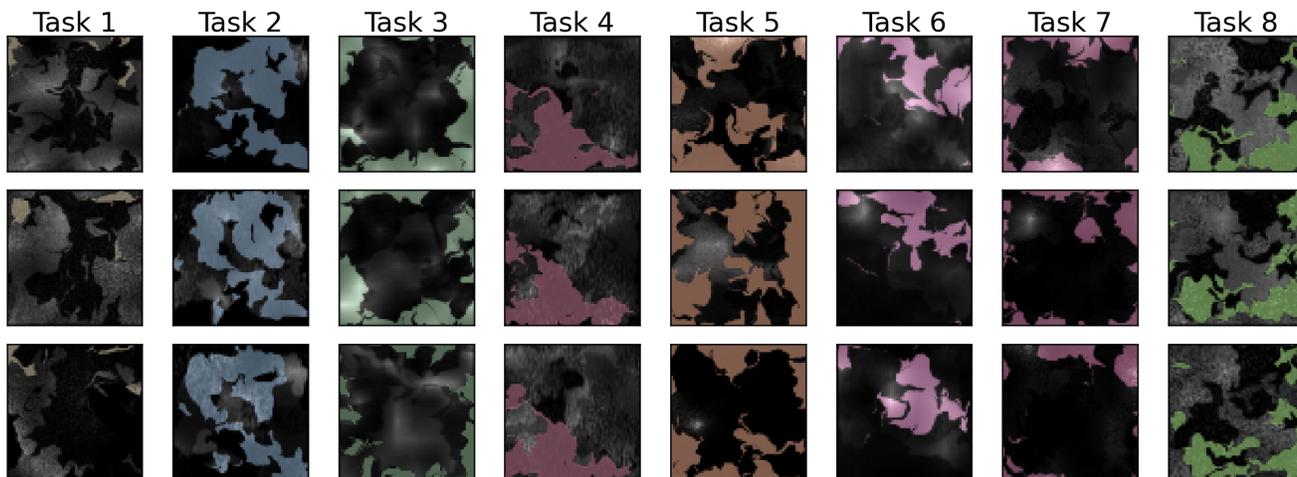


Figure 12: **Examples of Synthetically Generated Tasks.** We visualize 10 of the 1000 synthetically generated tasks, involving varying shapes, textures, and label shapes.

# E. Extended Results

## E.1. Main Results

We include detailed numbers corresponding to figures in the main body of the paper.

- **Method Comparison**. Table 5 reports test performance numbers of the results from Figure 1 and Table 1, comparing the segmentation results of UniverSeg to the FS baselines and the supervised nnUNet upper bounds.

- **Training Strategies Ablation**. Table 6 reports per-dataset test performance numbers for the results of Table 2 comparing several ways of augmenting the task diversity artificially. While the overall trend holds for most datasets, we find that the increase in task diversity has a detrimental effect on the STARE eye vessel segmentation task.

- **Model Support Size**. Table 7 reports held-out test performance numbers of the results from Figure 7 along with per-dataset breakdowns. We find that the global trend holds for each individual dataset, with larger support sizes achieving better results and ensembling (with $K = 10$) consistently improving predictions.

- **Available Data for Inference Ablation**. Table 8 reports extended results from Figure 8 with per-task results as we change the size of the support example pool. All tasks showcase the same trend with consistent improvements as more support examples are used during inference and with a reduced variance across random subsets of the support split.

- **Support Set Ensembling**. Table 9 reports results for the support set ensembling experiment. We observe a clear difference between $N = 1$ and $N > 1$ for ensembled predictions. For $N = 1$ ensembling leads to small improvements that eventually decline as $K$ grows. In contrast for $N > 1$, ensembling leads to substantial improvements that also reduce the variance of the distribution, limiting the dependence on the specific subset used for the support set.

- **Number of Tasks Ablation**. Table 10 reports the per-dataset and global dice numbers for the models trained with a subset of the training datasets.

Table 5: **Method Comparison**. Test Dice Score for the baselines, UniverSeg, and the nnUNet upper bounds in each of the held-out datasets. Standard deviation is computed by bootstrapping subjects before hierarchically averaging the data.

| Model | ACDC | PanDental | SCD | STARE | SpineWeb | WBC | All (avg.) |
|---|---|---|---|---|---|---|---|
| ALPNet | $34.6 \pm 2.4$ | $72.9 \pm 0.8$ | $53.4 \pm 3.0$ | $17.8 \pm 1.9$ | $31.6 \pm 4.6$ | $76.2 \pm 1.1$ | $47.8 \pm 1.1$ |
| PANet | $27.8 \pm 4.3$ | $67.7 \pm 0.8$ | $58.9 \pm 3.4$ | $20.1 \pm 3.2$ | $21.8 \pm 0.4$ | $54.7 \pm 1.6$ | $41.8 \pm 1.3$ |
| SENet | $40.1 \pm 2.0$ | $81.1 \pm 0.9$ | $55.4 \pm 3.3$ | $35.2 \pm 2.2$ | $18.3 \pm 4.0$ | $70.8 \pm 1.3$ | $50.1 \pm 1.3$ |
| UniverSeg (ours) | $\mathbf{70.9 \pm 2.9}$ | $\mathbf{87.5 \pm 0.9}$ | $\mathbf{69.0 \pm 2.9}$ | $\mathbf{48.1 \pm 2.0}$ | $\mathbf{64.6 \pm 5.4}$ | $\mathbf{90.6 \pm 1.1}$ | $\mathbf{71.8 \pm 0.9}$ |
| nnUNet (sup.) | $82.5 \pm 2.3$ | $92.9 \pm 1.1$ | $75.0 \pm 3.4$ | $65.5 \pm 1.1$ | $91.2 \pm 2.3$ | $95.1 \pm 0.7$ | $84.4 \pm 1.0$ |

Table 6: **Training Stategies Ablation**. Per dataset held-out Dice for UniverSeg models trained with different combinations of the proposed techniques to increase task diversity: in-task augmentation, task augmentation, and synthetic tasks.

| Synth | Medical | In-Task | Task | ACDC | PanDental | SCD | STARE | SpineWeb | WBC |
|---|---|---|---|---|---|---|---|---|---|
| ✓ | | | | $55.4 \pm 3.4$ | $80.6 \pm 1.3$ | $55.7 \pm 2.4$ | $42.6 \pm 2.5$ | $50.1 \pm 6.5$ | $86.0 \pm 1.4$ |
| | ✓ | | | $44.9 \pm 1.8$ | $85.3 \pm 0.9$ | $59.9 \pm 1.9$ | $\mathbf{63.8 \pm 0.9}$ | $40.3 \pm 6.0$ | $82.0 \pm 1.6$ |
| ✓ | ✓ | | | $50.6 \pm 2.9$ | $85.7 \pm 0.9$ | $59.0 \pm 1.9$ | $61.9 \pm 1.6$ | $45.6 \pm 4.8$ | $84.2 \pm 1.4$ |
| | ✓ | ✓ | | $52.3 \pm 4.3$ | $86.5 \pm 0.9$ | $64.9 \pm 2.7$ | $56.0 \pm 2.3$ | $57.2 \pm 3.7$ | $85.1 \pm 1.4$ |
| | ✓ | | ✓ | $68.0 \pm 3.0$ | $\mathbf{87.5 \pm 1.0}$ | $63.5 \pm 2.3$ | $56.6 \pm 2.1$ | $57.8 \pm 6.6$ | $89.2 \pm 1.3$ |
| | ✓ | ✓ | ✓ | $\mathbf{70.0 \pm 2.8}$ | $\mathbf{88.0 \pm 0.9}$ | $\mathbf{71.2 \pm 3.1}$ | $42.2 \pm 2.1$ | $58.4 \pm 8.5$ | $\mathbf{90.3 \pm 1.2}$ |
| ✓ | ✓ | ✓ | ✓ | $\mathbf{70.9 \pm 2.9}$ | $87.5 \pm 0.9$ | $69.0 \pm 2.9$ | $48.1 \pm 2.0$ | $\mathbf{64.6 \pm 5.4}$ | $\mathbf{90.6 \pm 1.1}$ |

Table 7: **Model Support Size**. Comparison of predictions for models trained with various of support sizes $N$ and evaluated with and without ensembling $K = 10$ predictions. We report results on each held-out dataset as well as the global average. Standard deviation is computed by bootstrapping subjects before hierarchically averaging the data. For all datasets, we find that increasing the support size leads to better predictions, with diminishing returns after $N > 16$. Ensembling predictions significantly improve performance in the majority of settings (paired t-test).

| N | K | ACDC | PanDental | SCD | STARE | SpineWeb | WBC | All (avg.) |
|---|---|------|-----------|-----|-------|----------|-----|------------|
| 1 | 1 | $41.3 \pm 1.3$ | $76.3 \pm 0.9$ | $60.2 \pm 1.8$ | $37.4 \pm 3.8$ | $30.4 \pm 5.5$ | $74.0 \pm 1.2$ | $53.3 \pm 1.0$ |
|   | 10 | $44.5 \pm 2.4$ | $79.1 \pm 1.0$ | $60.0 \pm 1.9$ | $38.5 \pm 4.0$ | $32.4 \pm 6.6$ | $79.4 \pm 1.4$ | $55.7 \pm 1.1$ |
| 2 | 1 | $41.3 \pm 2.6$ | $80.0 \pm 1.0$ | $63.5 \pm 2.0$ | $40.4 \pm 2.1$ | $38.0 \pm 4.3$ | $77.6 \pm 1.1$ | $56.8 \pm 1.0$ |
|   | 10 | $42.8 \pm 3.2$ | $82.4 \pm 1.1$ | $68.0 \pm 2.5$ | $40.7 \pm 2.3$ | $43.4 \pm 4.1$ | $82.3 \pm 1.4$ | $60.0 \pm 1.2$ |
| 4 | 1 | $53.9 \pm 1.9$ | $83.9 \pm 1.0$ | $64.7 \pm 1.7$ | $47.9 \pm 2.9$ | $45.5 \pm 4.0$ | $82.7 \pm 1.4$ | $63.1 \pm 0.8$ |
|   | 10 | $57.0 \pm 2.6$ | $84.6 \pm 1.1$ | $66.4 \pm 2.8$ | $48.6 \pm 2.9$ | $50.8 \pm 4.1$ | $85.7 \pm 1.5$ | $65.5 \pm 0.8$ |
| 8 | 1 | $57.0 \pm 2.5$ | $85.0 \pm 0.9$ | $66.9 \pm 3.2$ | $45.9 \pm 3.5$ | $57.3 \pm 6.5$ | $83.7 \pm 1.5$ | $66.0 \pm 1.3$ |
|   | 10 | $61.6 \pm 3.3$ | $86.1 \pm 0.9$ | $69.0 \pm 4.1$ | $47.1 \pm 3.5$ | $62.3 \pm 6.0$ | $85.9 \pm 1.5$ | $68.6 \pm 1.3$ |
| 16 | 1 | $64.1 \pm 2.4$ | $86.1 \pm 0.9$ | $69.1 \pm 3.1$ | $48.8 \pm 3.0$ | $64.4 \pm 5.8$ | $86.9 \pm 1.4$ | $69.9 \pm 1.0$ |
|    | 10 | $66.8 \pm 2.5$ | $86.7 \pm 0.9$ | $68.7 \pm 3.5$ | $\mathbf{49.7 \pm 2.8}$ | $\mathbf{66.8 \pm 5.7}$ | $88.3 \pm 1.5$ | $71.2 \pm 1.0$ |
| 32 | 1 | $65.6 \pm 3.0$ | $87.1 \pm 0.9$ | $69.0 \pm 2.0$ | $45.7 \pm 2.2$ | $65.8 \pm 4.6$ | $87.6 \pm 1.3$ | $70.1 \pm 0.9$ |
|    | 10 | $\mathbf{69.3 \pm 2.9}$ | $\mathbf{87.6 \pm 0.9}$ | $69.5 \pm 1.9$ | $\mathbf{46.4 \pm 2.1}$ | $\mathbf{66.4 \pm 4.3}$ | $88.9 \pm 1.4$ | $\mathbf{71.4 \pm 0.8}$ |
| 64 | 1 | $69.0 \pm 2.9$ | $87.2 \pm 0.9$ | $68.7 \pm 2.9$ | $47.2 \pm 2.2$ | $64.2 \pm 5.5$ | $89.7 \pm 1.1$ | $71.0 \pm 1.0$ |
|    | 10 | $\mathbf{70.9 \pm 2.9}$ | $\mathbf{87.5 \pm 0.9}$ | $69.0 \pm 2.9$ | $\mathbf{48.1 \pm 2.0}$ | $\mathbf{64.6 \pm 5.4}$ | $90.6 \pm 1.1$ | $\mathbf{71.8 \pm 0.9}$ |

Table 8: **Limited Example Data**. UniverSeg predictions using a limited $d_{\text{support}}$ example pool for each held-out task. For each size, we perform 100 repetitions using different random subsets, reporting the mean and standard deviation across them. Since some tasks do not have enough subjects to be evaluated for all values of $N$, we report $\min(N, |d_{\text{support}}|)$ and omit repeated settings where $N > |d_{\text{support}}|$.

| Task | N = 1 | N = 2 | N = 4 | N = 8 | N = 16 | N = 32 | N = 64 |
|------|-------|-------|-------|-------|--------|--------|--------|
| ACDC | $22.9 \pm 5.5$ | $38.5 \pm 6.9$ | $51.4 \pm 4.7$ | $59.1 \pm 3.0$ | $64.4 \pm 2.2$ | $68.6 \pm 1.4$ | $71.0 \pm 0.0$ |
| $\text{PanDental}_0$ | $59.1 \pm 7.4$ | $73.3 \pm 4.2$ | $77.6 \pm 1.6$ | $80.1 \pm 0.8$ | $82.1 \pm 0.5$ | $83.2 \pm 0.3$ | $83.7 \pm 0.1$ |
| $\text{PanDental}_1$ | $65.5 \pm 3.9$ | $84.1 \pm 2.2$ | $87.5 \pm 2.7$ | $89.5 \pm 1.2$ | $90.6 \pm 0.5$ | $91.1 \pm 0.3$ | $91.3 \pm 0.0$ |
| $\text{SCD}_0$ | $34.3 \pm 9.2$ | $63.1 \pm 5.1$ | $70.8 \pm 2.5$ | $73.1 \pm 1.2$ | $74.3 \pm 0.4$ | $74.2 \pm 0.0$ | |
| $\text{SCD}_1$ | $33.0 \pm 10.4$ | $61.8 \pm 9.4$ | $72.8 \pm 5.0$ | $76.7 \pm 2.4$ | $78.5 \pm 0.8$ | $78.6 \pm 0.0$ | |
| $\text{SCD}_2$ | $45.0 \pm 11.4$ | $71.5 \pm 12.7$ | $80.6 \pm 7.2$ | $84.8 \pm 0.0$ | | | |
| $\text{SCD}_3$ | $30.5 \pm 9.3$ | $47.1 \pm 6.5$ | $54.9 \pm 4.5$ | $63.0 \pm 2.3$ | $64.2 \pm 0.0$ | | |
| $\text{SCD}_4$ | $9.2 \pm 4.4$ | $13.3 \pm 8.3$ | $25.9 \pm 7.8$ | $39.0 \pm 3.1$ | $41.0 \pm 0.0$ | | |
| STARE | $25.5 \pm 3.5$ | $33.5 \pm 2.0$ | $40.2 \pm 1.2$ | $45.2 \pm 0.5$ | $47.7 \pm 0.0$ | | |
| SpineWeb | $28.1 \pm 2.1$ | $39.3 \pm 6.1$ | $52.1 \pm 7.0$ | $63.1 \pm 3.3$ | $64.7 \pm 0.0$ | | |
| $\text{WBC}_0$ | $49.4 \pm 4.5$ | $65.0 \pm 4.3$ | $74.8 \pm 3.0$ | $81.0 \pm 1.9$ | $85.0 \pm 1.2$ | $87.5 \pm 0.8$ | $88.8 \pm 0.0$ |
| $\text{WBC}_1$ | $57.4 \pm 4.8$ | $75.2 \pm 3.5$ | $83.0 \pm 2.1$ | $87.4 \pm 1.0$ | $89.9 \pm 0.4$ | $91.3 \pm 0.3$ | $91.9 \pm 0.2$ |

Table 9: **Ensembling predictions at different inference support sizes.** For each inference support size $N$, we report the results (in average held-out Dice Score) of taking 100 predictions ($K = 1$) and ensembling by averaging in groups of size $K$, performing 100 repetitions for each $K$. We report the mean and standard deviation across the 100 values for each setting and find that increasing either $K$ or $N$ leads to improved model performance, with $N$ having a significantly larger effect than $K$.

| N | K = 1 | K = 2 | K = 4 | K = 8 | K = 16 | K = 32 | K = 64 |
|---|---|---|---|---|---|---|---|
| 1 | $36.9 \pm 2.0$ | $39.4 \pm 2.3$ | $40.7 \pm 2.0$ | $40.9 \pm 1.6$ | $40.3 \pm 1.1$ | $39.4 \pm 0.7$ | $38.3 \pm 0.4$ |
| 2 | $51.0 \pm 3.2$ | $56.3 \pm 2.2$ | $59.5 \pm 1.6$ | $61.0 \pm 1.2$ | $61.9 \pm 0.9$ | $62.3 \pm 0.6$ | $62.4 \pm 0.3$ |
| 4 | $59.4 \pm 2.3$ | $63.7 \pm 1.5$ | $66.2 \pm 1.0$ | $67.5 \pm 0.7$ | $68.2 \pm 0.4$ | $68.6 \pm 0.3$ | $68.8 \pm 0.2$ |
| 8 | $64.8 \pm 1.9$ | $68.0 \pm 1.1$ | $69.6 \pm 0.6$ | $70.5 \pm 0.4$ | $71.1 \pm 0.2$ | $71.3 \pm 0.2$ | $71.4 \pm 0.1$ |
| 16 | $68.4 \pm 1.1$ | $70.1 \pm 0.5$ | $71.0 \pm 0.4$ | $71.5 \pm 0.3$ | $71.8 \pm 0.2$ | $71.9 \pm 0.1$ | $72.0 \pm 0.1$ |
| 32 | $70.1 \pm 0.6$ | $71.0 \pm 0.3$ | $71.5 \pm 0.2$ | $71.7 \pm 0.1$ | $71.8 \pm 0.1$ | $71.9 \pm 0.1$ | $71.9 \pm 0.0$ |
| 64 | $71.0 \pm 0.3$ | $71.4 \pm 0.2$ | $71.6 \pm 0.2$ | $71.7 \pm 0.1$ | $71.8 \pm 0.1$ | $71.8 \pm 0.1$ | $71.8 \pm 0.0$ |

Table 10: **Number of Training Datasets and Tasks**. Test Dice score results for models trained with $N_D$ datasets comprising $N_T$ tasks. The subsets of the training datasets are chosen randomly so we report three realizations for each $N_D$, except for the case where all datasets are included. Each row corresponds to a separate UniverSeg model.

| $N_D$ | $N_T$ | ACDC | PanDental | SCD | STARE | SpineWeb | WBC | All (avg) |
|---|---|---|---|---|---|---|---|---|
| | 25 | $24.7 \pm 2.8$ | $82.2 \pm 0.8$ | $43.7 \pm 3.3$ | $7.2 \pm 2.5$ | $0.2 \pm 0.2$ | $61.1 \pm 1.3$ | $36.5 \pm 0.8$ |
| 1 | 29 | $3.3 \pm 2.4$ | $18.4 \pm 0.7$ | $0.2 \pm 0.1$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $0.0 \pm 0.0$ | $3.7 \pm 0.4$ |
| | 156 | $63.1 \pm 2.7$ | $80.0 \pm 1.8$ | $53.5 \pm 5.3$ | $17.5 \pm 2.3$ | $46.2 \pm 1.6$ | $86.2 \pm 1.3$ | $57.8 \pm 1.1$ |
| | 33 | $41.7 \pm 3.7$ | $83.8 \pm 1.0$ | $43.9 \pm 2.3$ | $16.9 \pm 1.4$ | $22.2 \pm 3.3$ | $80.0 \pm 1.5$ | $48.1 \pm 1.1$ |
| 2 | 85 | $49.5 \pm 3.5$ | $83.1 \pm 1.1$ | $59.9 \pm 2.5$ | $22.9 \pm 2.1$ | $52.9 \pm 2.1$ | $85.5 \pm 1.4$ | $59.0 \pm 1.0$ |
| | 131 | $31.9 \pm 2.5$ | $82.0 \pm 0.8$ | $42.3 \pm 3.4$ | $0.0 \pm 0.0$ | $10.2 \pm 3.2$ | $69.5 \pm 1.3$ | $39.3 \pm 0.6$ |
| | 237 | $43.6 \pm 5.5$ | $75.3 \pm 1.3$ | $52.1 \pm 2.8$ | $26.8 \pm 3.4$ | $28.3 \pm 10.1$ | $86.3 \pm 1.2$ | $52.1 \pm 2.2$ |
| 5 | 710 | $63.9 \pm 3.6$ | $86.6 \pm 1.0$ | $63.5 \pm 2.2$ | $27.2 \pm 3.3$ | $61.1 \pm 5.0$ | $87.4 \pm 1.6$ | $64.9 \pm 1.3$ |
| | 1117 | $67.4 \pm 3.1$ | $86.5 \pm 1.0$ | $62.9 \pm 4.5$ | $51.6 \pm 2.4$ | $58.4 \pm 8.2$ | $89.9 \pm 1.0$ | $69.4 \pm 2.0$ |
| | 174 | $51.2 \pm 2.8$ | $84.0 \pm 0.9$ | $66.5 \pm 2.7$ | $22.8 \pm 1.3$ | $52.2 \pm 0.4$ | $84.5 \pm 1.1$ | $60.2 \pm 0.8$ |
| 11 | 1223 | $60.6 \pm 3.6$ | $86.8 \pm 1.0$ | $59.0 \pm 4.9$ | $29.7 \pm 2.4$ | $58.4 \pm 5.6$ | $85.5 \pm 1.3$ | $63.3 \pm 1.7$ |
| | 1457 | $69.7 \pm 2.8$ | $87.5 \pm 0.9$ | $65.4 \pm 3.3$ | $40.5 \pm 3.0$ | $59.6 \pm 7.0$ | $88.6 \pm 1.1$ | $68.6 \pm 1.7$ |
| | 1320 | $66.9 \pm 3.5$ | $85.3 \pm 0.9$ | $68.1 \pm 2.4$ | $31.2 \pm 0.7$ | $57.2 \pm 5.9$ | $88.4 \pm 1.3$ | $66.2 \pm 1.1$ |
| 23 | 2157 | $66.5 \pm 3.0$ | $86.1 \pm 1.0$ | $64.0 \pm 2.3$ | $36.9 \pm 1.1$ | $64.9 \pm 5.4$ | $89.6 \pm 1.3$ | $68.0 \pm 1.1$ |
| | 2276 | $66.5 \pm 3.8$ | $86.7 \pm 0.9$ | $65.7 \pm 2.8$ | $28.1 \pm 2.4$ | $52.4 \pm 6.6$ | $89.1 \pm 1.2$ | $64.8 \pm 1.4$ |
| | 3008 | $69.8 \pm 3.0$ | $88.8 \pm 0.9$ | $68.3 \pm 2.6$ | $42.5 \pm 3.2$ | $62.4 \pm 6.1$ | $90.0 \pm 1.2$ | $70.3 \pm 1.4$ |
| 34 | 3483 | $70.7 \pm 2.9$ | $87.1 \pm 1.0$ | $67.2 \pm 3.5$ | $46.7 \pm 3.6$ | $63.0 \pm 5.0$ | $90.7 \pm 1.4$ | $70.9 \pm 1.0$ |
| | 3854 | $65.3 \pm 3.6$ | $88.2 \pm 1.0$ | $65.1 \pm 1.6$ | $43.1 \pm 3.1$ | $62.2 \pm 5.5$ | $89.0 \pm 1.1$ | $68.8 \pm 1.0$ |
| 46 | 4432 | $71.3 \pm 2.6$ | $87.9 \pm 0.9$ | $67.9 \pm 2.5$ | $44.9 \pm 2.9$ | $65.5 \pm 4.7$ | $91.0 \pm 1.1$ | $71.4 \pm 1.1$ |

## E.2. Additional Results

**Few-shot Baseline Model Variants**. The FS baselines (ALPNet, PANet, and SENet) were introduced in a few-shot setting where the underlying assumption is that any new task can only have very few examples, rather than our setting where we avoid re-training due to the limitations of the clinical settings. These baselines were therefore presented with a support size of 1 example. They also involved no data or task augmentation. Our ablations show that UniverSeg models performed best with large support set sizes and increased data and task diversity from augmenting examples. Consequently, we test whether incorporating these changes to the baseline methods leads to improved performance in the held-out datasets in our setting, where more data *might* be available for some datasets. Similarly, we also test whether ensembling predictions from several support sets lead to better predictions, as we do for UniverSeg.

Table 11 and Table 12 report results of the hyperparameter grid search for all the few-shot baseline models and UniverSeg. Table 11 shows that ensembling ($K = 10$) and an increased support size ($N = 64$) leads to held-out improvements for all methods. In contrast, augmentation strategies do not benefit all methods. While UniverSeg and SENet improve when using augmentation strategies, PANet and ALPNet experience a decrease in performance. Table 12 shows that the best hyperparameter setting is not consistent across held-out datasets for the baseline methods.

Table 11: **FS baseline hyperparameter search**. For each method, we report results for models trained with a support size $N$, ensemble size $K$, and with and without data and task augmentation. Dice scores are averaged across all datasets and the standard deviation is computed via subject-level bootstrapping.

| Model | N | No Aug | | Aug | |
|---|---|---|---|---|---|
| | | K=1 | K=10 | K=1 | K=10 |
| ALPNet | 1 | $40.2 \pm 0.9$ | $42.3 \pm 1.3$ | $35.4 \pm 0.6$ | $37.0 \pm 0.8$ |
| | 64 | $46.3 \pm 1.3$ | $\mathbf{47.8 \pm 1.1}$ | $42.3 \pm 1.0$ | $45.2 \pm 1.2$ |
| PANet | 1 | $37.4 \pm 0.7$ | $39.3 \pm 0.8$ | $33.2 \pm 1.3$ | $34.3 \pm 1.4$ |
| | 64 | $\mathbf{41.6 \pm 1.3}$ | $\mathbf{41.8 \pm 1.3}$ | $38.7 \pm 0.9$ | $40.8 \pm 0.8$ |
| SENet | 1 | $40.0 \pm 0.9$ | $41.2 \pm 0.9$ | $40.1 \pm 1.2$ | $41.1 \pm 1.4$ |
| | 64 | $42.1 \pm 0.7$ | $42.4 \pm 0.8$ | $\mathbf{50.2 \pm 1.1}$ | $50.1 \pm 1.3$ |
| UniverSeg (ours) | 1 | $49.7 \pm 0.9$ | $53.4 \pm 1.1$ | $51.9 \pm 0.8$ | $54.0 \pm 1.0$ |
| | 64 | $64.0 \pm 1.1$ | $64.5 \pm 1.0$ | $71.0 \pm 1.0$ | $\mathbf{71.8 \pm 0.9}$ |

Table 12: **Few-shot baseline hyperparameter search per dataset**. For each method, we report results for models trained with a support size $N$, ensemble size $K = 10$, and with and without data and task augmentation. Dice score values are averaged across all datasets and the standard deviation is computed via subject-level bootstrapping. For each dataset and model, we highlight the setting with the best performance

| Model | N | Aug | ACDC | PanDental | SCD | STARE | SpineWeb | WBC |
|---|---|---|---|---|---|---|---|---|
| ALPNet | 1 | No | $22.1 \pm 3.2$ | $66.8 \pm 1.0$ | $49.1 \pm 3.8$ | $\mathbf{22.7 \pm 2.0}$ | $29.7 \pm 3.7$ | $63.2 \pm 0.9$ |
| | | Yes | $26.7 \pm 2.9$ | $51.8 \pm 1.5$ | $41.5 \pm 1.9$ | $11.0 \pm 2.8$ | $19.7 \pm 4.9$ | $71.0 \pm 1.6$ |
| | 64 | No | $34.6 \pm 2.4$ | $\mathbf{72.9 \pm 0.8}$ | $53.4 \pm 3.0$ | $17.8 \pm 1.9$ | $\mathbf{31.6 \pm 4.6}$ | $\mathbf{76.2 \pm 1.1}$ |
| | | Yes | $\mathbf{38.3 \pm 2.5}$ | $71.1 \pm 1.0$ | $\mathbf{56.1 \pm 1.6}$ | $6.3 \pm 2.2$ | $25.5 \pm 6.4$ | $73.9 \pm 1.2$ |
| PANet | 1 | No | $\mathbf{33.4 \pm 2.5}$ | $\mathbf{69.8 \pm 1.3}$ | $48.7 \pm 3.5$ | $17.4 \pm 4.3$ | $25.4 \pm 3.9$ | $40.9 \pm 1.8$ |
| | | Yes | $30.3 \pm 3.0$ | $63.2 \pm 1.3$ | $48.4 \pm 3.3$ | $4.6 \pm 2.9$ | $\mathbf{28.6 \pm 5.6}$ | $31.0 \pm 2.1$ |
| | 64 | No | $27.8 \pm 4.3$ | $67.7 \pm 0.8$ | $\mathbf{58.9 \pm 3.4}$ | $\mathbf{20.1 \pm 3.2}$ | $21.8 \pm 0.4$ | $54.7 \pm 1.6$ |
| | | Yes | $29.6 \pm 2.3$ | $66.4 \pm 1.4$ | $46.8 \pm 2.3$ | $15.1 \pm 2.1$ | $27.9 \pm 5.8$ | $\mathbf{58.8 \pm 1.5}$ |
| SENet | 1 | No | $17.0 \pm 2.9$ | $61.7 \pm 1.1$ | $47.5 \pm 2.3$ | $\mathbf{41.3 \pm 2.7}$ | $\mathbf{21.7 \pm 3.7}$ | $58.1 \pm 0.9$ |
| | | Yes | $32.2 \pm 2.8$ | $62.4 \pm 1.3$ | $48.2 \pm 2.9$ | $31.4 \pm 2.1$ | $16.8 \pm 7.8$ | $55.5 \pm 1.3$ |
| | 64 | No | $32.0 \pm 2.4$ | $79.1 \pm 0.8$ | $43.8 \pm 2.9$ | $37.5 \pm 2.9$ | $3.2 \pm 2.4$ | $58.7 \pm 1.1$ |
| | | Yes | $\mathbf{40.1 \pm 2.0}$ | $\mathbf{81.1 \pm 0.9}$ | $\mathbf{55.4 \pm 3.3}$ | $35.2 \pm 2.2$ | $18.3 \pm 4.0$ | $\mathbf{70.8 \pm 1.3}$ |
| UniverSeg | 1 | No | $29.1 \pm 2.0$ | $76.1 \pm 0.9$ | $58.0 \pm 2.1$ | $54.5 \pm 2.7$ | $31.6 \pm 6.4$ | $70.9 \pm 1.7$ |
| | | Yes | $37.5 \pm 2.0$ | $76.8 \pm 1.1$ | $62.9 \pm 2.6$ | $33.9 \pm 3.7$ | $36.0 \pm 5.0$ | $76.8 \pm 1.6$ |
| | 64 | No | $50.6 \pm 2.9$ | $85.7 \pm 0.9$ | $59.0 \pm 1.9$ | $\mathbf{61.9 \pm 1.6}$ | $45.6 \pm 4.8$ | $84.2 \pm 1.4$ |
| | | Yes | $\mathbf{70.9 \pm 2.9}$ | $\mathbf{87.5 \pm 0.9}$ | $\mathbf{69.0 \pm 2.9}$ | $48.1 \pm 2.0$ | $\mathbf{64.6 \pm 5.4}$ | $\mathbf{90.6 \pm 1.1}$ |

**Using different training and inference support sizes**. In Figure 13, we report dataset-level results of performing inference with a support size of $M$ using a UniverSeg model trained with a support size of $N$ examples. We find that using support sets larger than those seen in training ($M > N$, lower quadrant of heat-maps) leads to improvements for $N \geq 2$, which demonstrates the model is learning to interact the elements of the support set and benefits from larger amounts of examples.
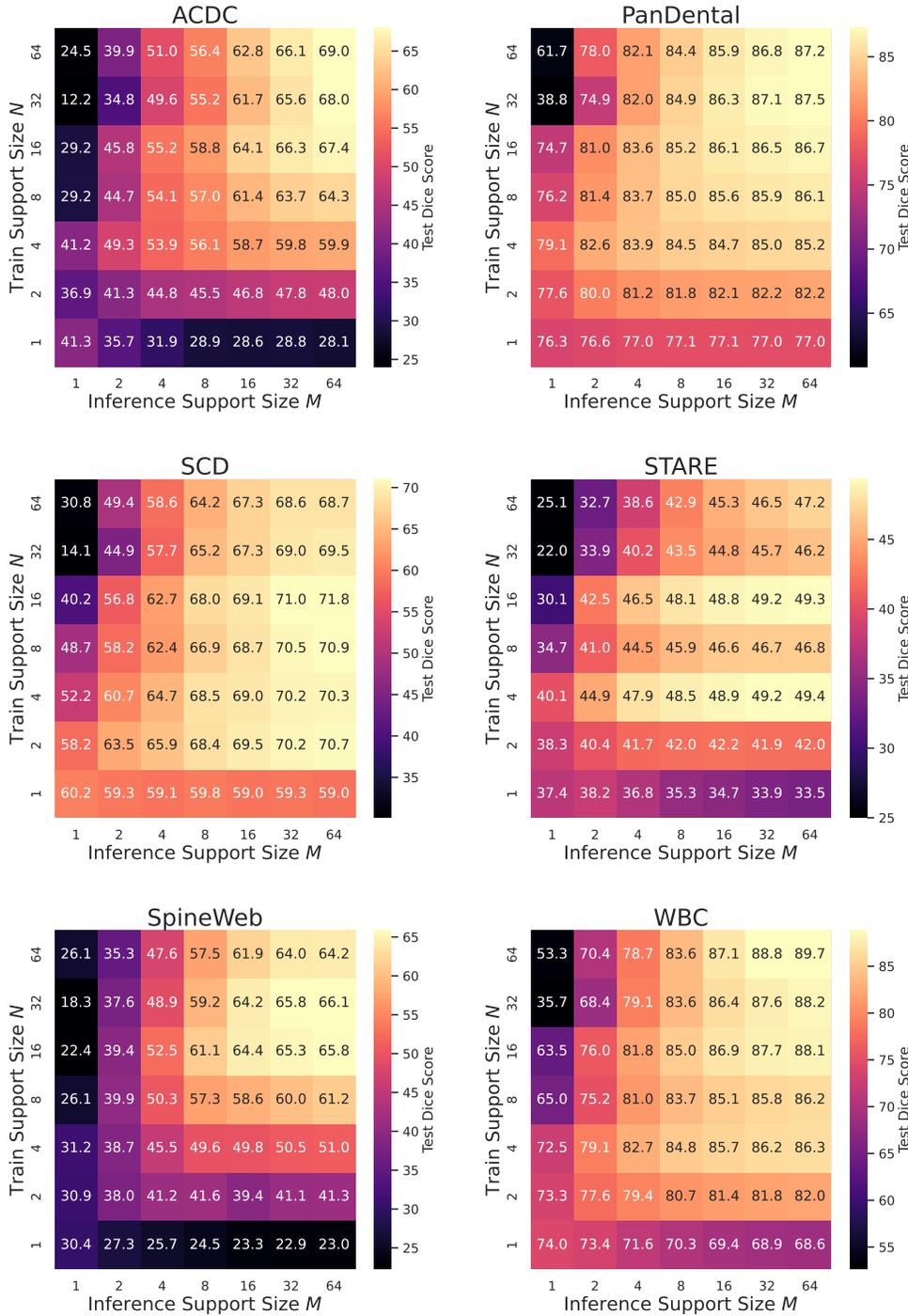


Figure 13: **Cartesian Product of Training and Inference Support Sizes.** Test results for using a UniverSeg trained with a support size of $N$ examples and performing inference with a support size of $M$ examples. No ensembling is performed ($K = 1$), but we perform 10 repetitions with varying support sets and report the average.

# F. Additional Visualizations

**Visualization of Held-Out Support Sets.** UniverSeg networks take advantage of large support sets of (image, label-map) pairs, which can be very diverse. In Figure 14, we visualize a random subset of 10 pairs for each held-out dataset. The diversity of subjects amongst support sets differs between tasks, which likely plays a role in the number of examples required to perform well.



Figure 14: **Example Support Sets for Held-Out Datasets.**

**Visualization of Soft Predictions.** Thresholding segmentation predictions (at 0.5) provides a binary segmentation and enables computation of well-known metrics such as the (hard) Dice score. However, for certain regions of interest, like thin structures, thresholding can hide network performance. In Figure 15, we show this effect visually. For example, focusing on STARE, we see that UniverSeg networks can capture the thin structures very well, which is lost when thresholding the predictions to create a binary segmentation.



Figure 15: **Visualization of Soft (Non-Thresholded) Predictions for All Methods.**

**WBC task visualizations**. We include some visualizations of UniverSeg's capability to adapt based on the support set specification. We use the WBC dataset, which presents substantial variability between support set examples.

- Figure 16 presents support set examples for the WBC Cytoplasm label as well as held-out predictions, showing that UniverSeg closely matches the ground truth.

- Figure 17 shows how UniverSeg is equivariant with respect to the support set labels. Given the same images as in Figure 16 but different labels, UniverSeg adapts its predictions to the nucleus label.

- Figure 18 showcases UniverSeg's invariance to image transformations. Using the same images and label examples from Figure 16, we invert the image data (i.e. $1 - x$) for both the query and support set images. UniverSeg correctly segments the label regardless of the image transformation.

- Figure 19 shows that while UniverSeg is trained on binary segmentation tasks, it can adequately perform multi-label segmentation. To produce multi-label predictions, we treat each label independently, and then combine the predictions for each label using a softmax operation.

- Figure 20 shows the effect of the support set size $N$ in the prediction results. We observe that segmentation mask quality substantially improves as we increase the number of support set image-label pairs.

- Figure 21 shows prediction variability for predictions performed with support size $N = 8$ along with an ensembled prediction.

(a) **Support Set Examples - Cytoplasm Label**

(b) **Predictions - Cytoplasm Label**

Figure 16: Visualization of support set examples (a) and predictions (b) for the WBC Cytoplasm label

(a) **Support Set Examples - Nucleus Label**
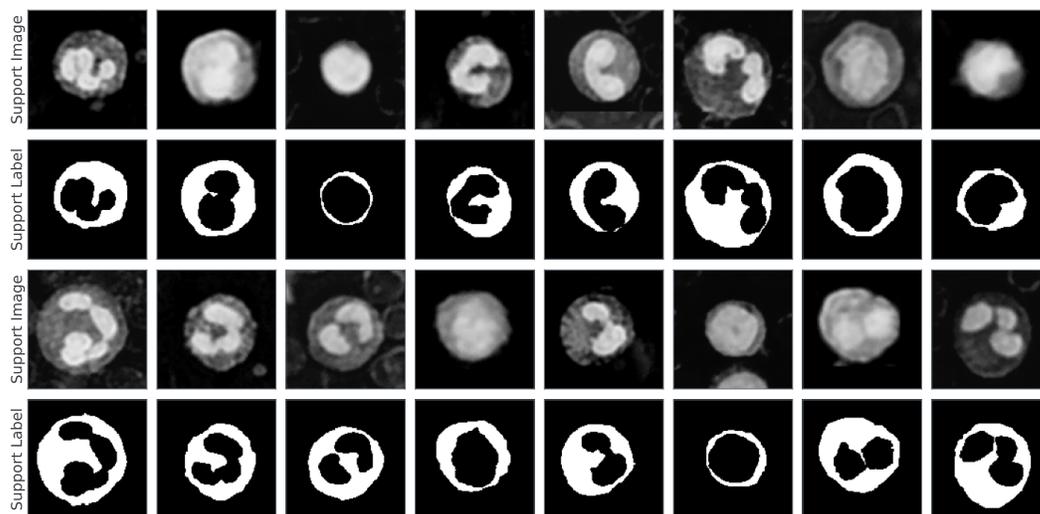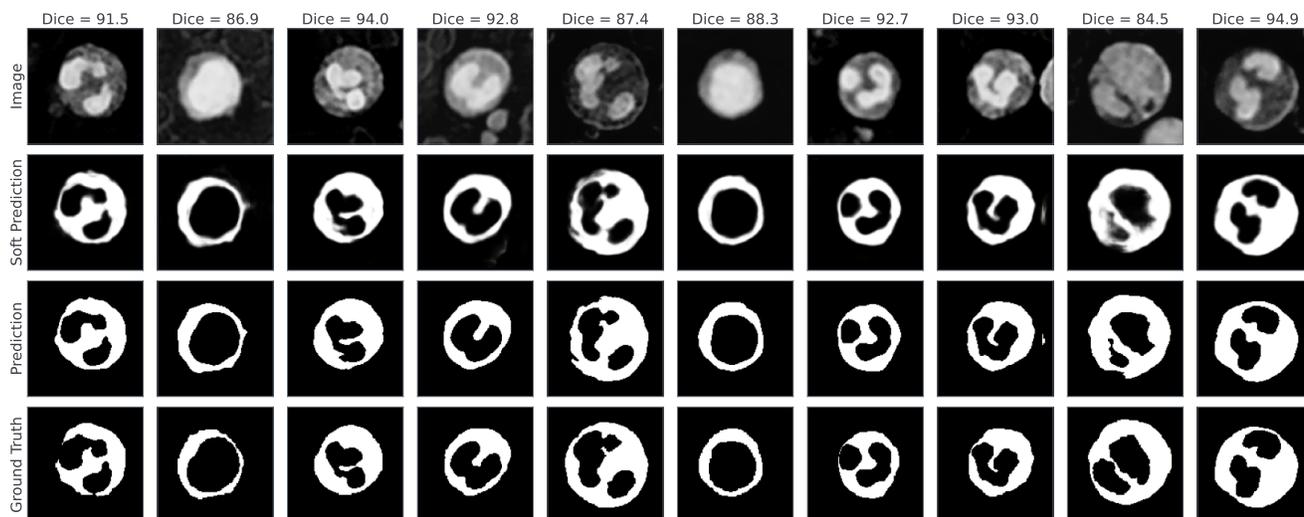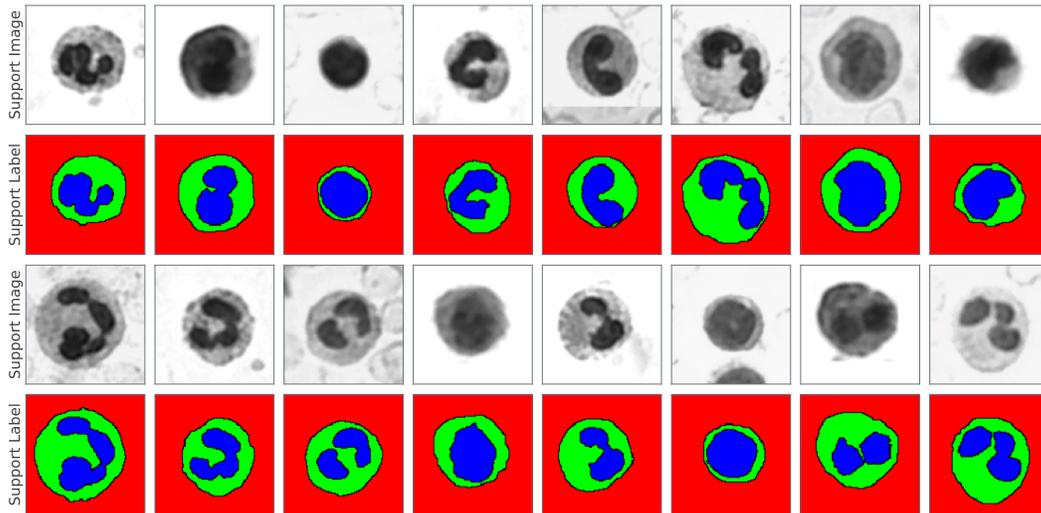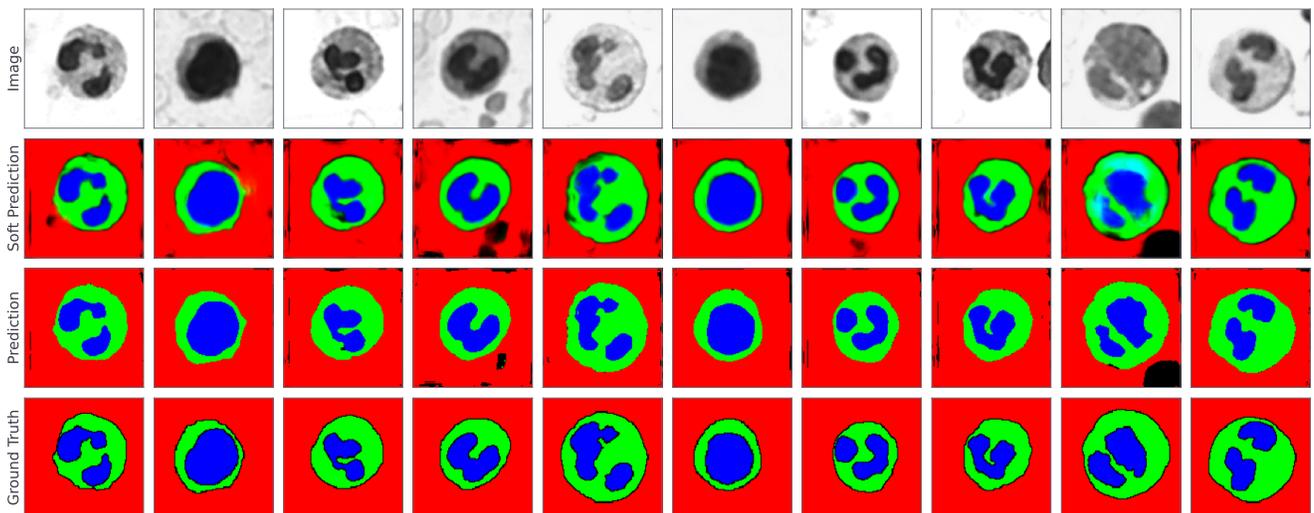


(b) **Predictions - Nucleus Label**



Figure 17: Visualization of support set examples (a) and predictions (b) for the WBC Nucleus label

Figure 18: Visualization of support set examples (a) and predictions (b) for the WBC Cytoplasm label with inverted images

Figure 19: Visualization of support set examples (a) and predictions (b) for the WBC task with multiple labels being predicted independently. Each label is encoded using a RGB channel (Red=backgroud, Green=Cytoplasm, Blue=Nuclues), we only see some mild nucleus-cytoplasm overlaps in cyan for one example.
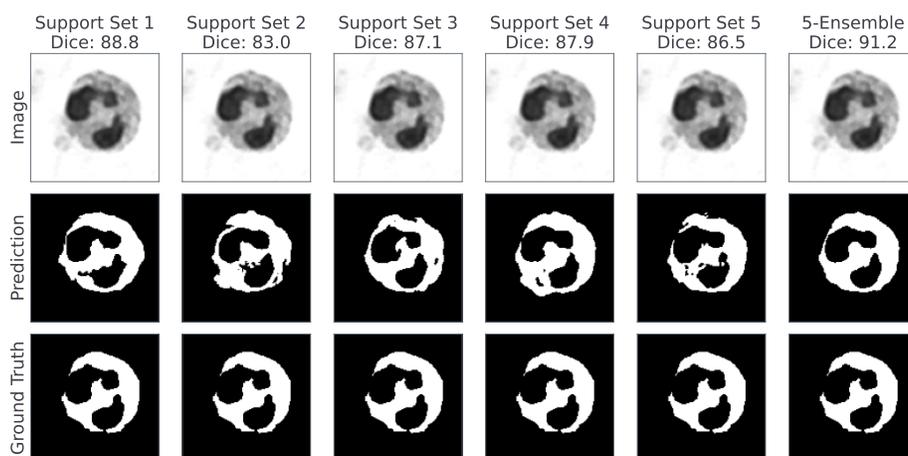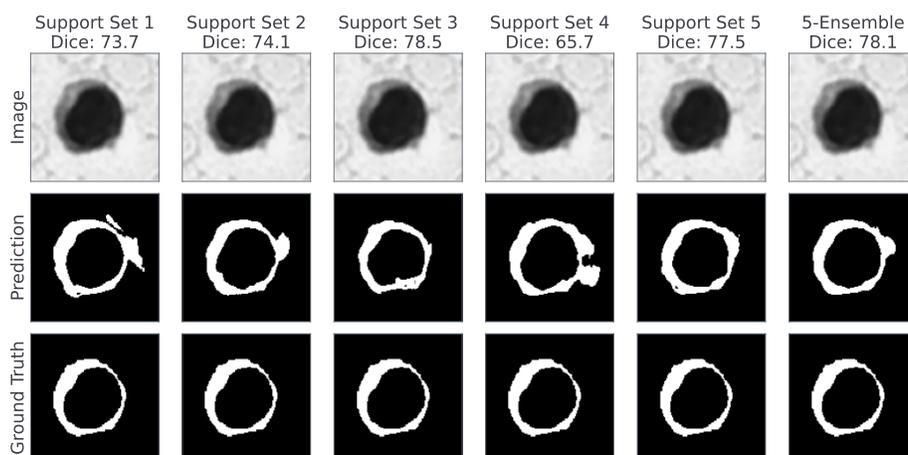
Figure 20: Visualization of predictions for the WBC Cytoplasm task with varying number of support set examples $N$. Larger support sets lead to better segmentation masks.
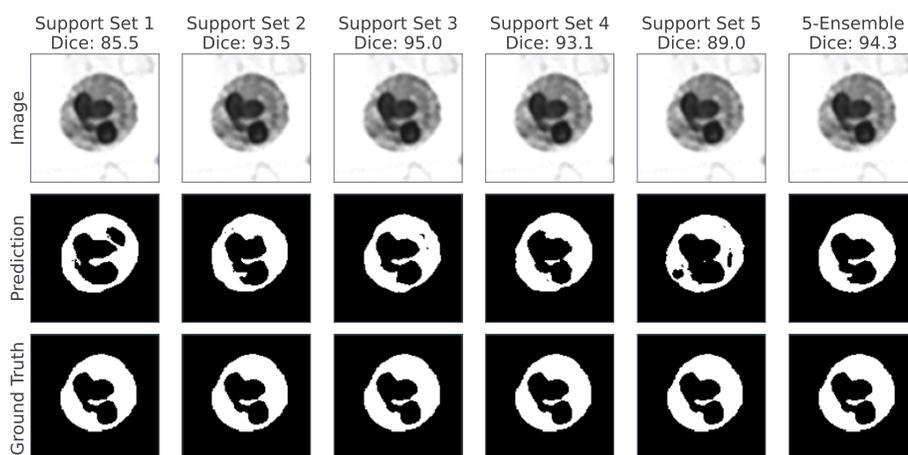
Figure 21: Visualization of predictions for the WBC Cytoplasm task with various choices of support set ($N = 8$) as well as the ensembled prediction (last column). Ensembling reduces the variance of predictions.