



Anomaly Detection-Inspired Few-Shot Medical Image Segmentation Through Self-Supervision With Supervoxels^{*}

Stine Hansen^{a,*}, Srishti Gautam^a, Robert Jenssen^a, Michael Kampffmeyer^a

^aDepartment of Physics and Technology, UiT The Arctic University of Norway, NO-9037 Tromsø, Norway

ARTICLE INFO

Article history:

Received 22 June 2021

Accepted 1 February 2022

Available online 11 February 2022

Keywords: Organ Segmentation, Cardiac Segmentation, Few-Shot Learning, Anomaly Detection, Self-Supervision, Supervoxels

ABSTRACT

Recent work has shown that label-efficient few-shot learning through self-supervision can achieve promising medical image segmentation results. However, few-shot segmentation models typically rely on prototype representations of the semantic classes, resulting in a loss of local information that can degrade performance. This is particularly problematic for the typically large and highly heterogeneous background class in medical image segmentation problems. Previous works have attempted to address this issue by learning additional prototypes for each class, but since the prototypes are based on a limited number of slices, we argue that this ad-hoc solution is insufficient to capture the background properties. Motivated by this, and the observation that the foreground class (e.g., one organ) is relatively homogeneous, we propose a novel anomaly detection-inspired approach to few-shot medical image segmentation in which we refrain from modeling the background explicitly. Instead, we rely solely on a single foreground prototype to compute anomaly scores for all query pixels. The segmentation is then performed by thresholding these anomaly scores using a learned threshold. Assisted by a novel self-supervision task that exploits the 3D structure of medical images through supervoxels, our proposed anomaly detection-inspired few-shot medical image segmentation model outperforms previous state-of-the-art approaches on two representative MRI datasets for the tasks of abdominal organ segmentation and cardiac segmentation.

© 2022 Elsevier B. V. All rights reserved.

1. Introduction

Many applications in medical image analysis, such as diagnosis (Tsochatzidis et al., 2021), treatment planning (Chen et al., 2021), and quantification of tissue volumes (Abdeltawab et al., 2020) rely heavily on semantic segmentation. To lessen

^{*}All the authors are with the UiT Machine Learning Group (machine-learning.uit.no) and with *Visual Intelligence*, a Norwegian Centre for Research-based Innovation (visual-intelligence.no).

^{*}Corresponding author

e-mail: s.hansen@uit.no (Stine Hansen), srishti.gautam@uit.no (Srishti Gautam), robert.jenssen@uit.no (Robert Jenssen), michael.c.kampffmeyer@uit.no (Michael Kampffmeyer)

the burden on the medical practitioners performing these manual, slice-by-slice segmentations, the use of deep learning for automatic segmentation has a great potential. Unfortunately, existing segmentation frameworks (Ronneberger et al., 2015; Li et al., 2018; Isensee et al., 2021) depend on supervised training and large amounts of densely labeled data, which are often unavailable in the medical domain. Moreover, their generalization properties to previously unseen classes are typically poor, necessitating the collection and labeling of new data to re-train for new tasks. Due to the huge number of potential segmentation tasks in medical images, this makes these models impractical to use.

Inspired by how humans learn from only a handful of in-

stances, few-shot learning has emerged as a learning paradigm to foster models that can easily adapt to new concepts when exposed to just a few new, labeled samples. These models typically follow an episodic framework (Vinyals *et al.*, 2016) where, in each episode, k labeled samples, called the support set, are used to segment the unlabeled query image(s). The models are trained on one set of classes and learn to, with only a few annotated examples, segment objects from new classes. A *trained* few-shot segmentation (FSS) model is thus able to segment an unseen organ class based on just a few labeled instances. However, in order to avoid over-fitting, typical FSS models rely on training data containing a large set of labeled training classes, generally not available in the medical domain.

In a recent work, Ouyang *et al.* (2020) proposed a label-efficient approach to medical image segmentation, building on metric-learning based prototypical FSS (Liu *et al.*, 2020b; Wang *et al.*, 2019). They suggest a model that follows the traditional few-shot episodic framework, where class-wise prototypes are extracted from the labeled support set and used to reduce the segmentation of the unlabeled query image to a pixel-wise prototype matching in the embedding space. Whereas traditional few-shot learning models require a set of annotated training classes, Ouyang *et al.* (2020) propose a clever way to bypass this need by employing self-supervised training (Jing and Tian, 2020). Instead of sampling labeled support and query images, they construct the training episodes based on *one* unlabeled image slice and its corresponding superpixel (Ren and Malik, 2003) segmentation: One randomly sampled superpixel serves as foreground mask, and together with the original image slice, these form the support image-label pair. The query pair is then constructed by applying random transformations to the support pair. In this way, they enable training of the network without using annotations, i.e. the model is trained unsupervised. Finally, in the inference phase, they only need a few labeled image slices to perform segmentation on new classes.

However, a general problem with prototypical FSS is the loss of local information caused by average pooling of features during prototype extraction. This is particularly problematic for spatially heterogeneous classes like the background class in medical image segmentation problems, which can contain any semantic class other than the foreground class. Previous metric-learning based works have addressed this issue by computing additional prototypes per class to capture more diverse features. Liu *et al.* (2020b) clustered the features within each class to obtain *part-aware* prototypes and in the current state-of-the-art method, Ouyang *et al.* (2020) computed additional *local* prototypes on a regular grid.

We argue that it is insufficient to model the entire background volume with prototypes estimated from a few support slices and propose a conceptually different approach where we do *not* increase the number of background prototypes but remove the need for these altogether. Inspired by the anomaly detection literature (Chandola *et al.*, 2009; Ruff *et al.*, 2021), we propose to only model the relatively homogeneous foreground class with a single prototype and introduce an anomaly score that measures the dissimilarity between this foreground prototype and all query pixels. Segmentation is then performed

by thresholding the anomaly scores using a learned threshold that encourages compact foreground representations. For direct comparison of our novel anomaly detection-inspired few-shot medical image segmentation method to that of Ouyang *et al.* (2020) and other representative works, our baseline setup follows their approach, working with 2D image slices. Within the existing 2D setup, we, as an added contribution, propose a new self-supervision task by extending the superpixel-based self-supervision scheme by Ouyang *et al.* (2020) to 3D in order to utilize the volumetric nature of the data. As a natural extension, facilitated by the new self-supervision task, we further indicate potential benefits beyond this 2D setup by exploring a direct 3D treatment of the problem by employing a 3D convolutional neural network (CNN) as embedding network.

By only explicitly modeling the foreground class, we argue that our proposed approach is more robust to background outside the support slices, compared to current state-of-the-art methods (Ouyang *et al.*, 2020; Roy *et al.*, 2020). To further illustrate this, we introduce a new evaluation protocol where we, based on labeled slices from the support image, segment the entire query image, thus being more exposed to background effects. Previous works, on the other hand, limit the evaluation of the query image only to the slices containing the class of interest. However, this approach requires additional weak labels in the form of information about the location of the class in the query image, which is unrealistic and cumbersome, especially in the medical setting.

In summary, the main contributions of this work are threefold. We propose:

- (1) A simple but effective anomaly detection-inspired approach to FSS that outperforms prior state-of-the-art methods and removes the need to learn a large number of prototypes.
- (2) A novel self-supervision task that exploits the 3D structural information in medical images within the 2D setup and indicate the potential of training 3D CNNs for direct volume segmentation.
- (3) A new evaluation protocol for few-shot medical image segmentation that does not rely on weak-labels and therefore is more applicable in practical scenarios.

2. Related Work

2.1. Few-Shot Meta-learning

As opposed to classical supervised learning that specializes a model to perform one specific task by optimizing over training samples, few-shot meta-learning optimizes over a set of training tasks, with the goal of obtaining a model that can quickly adapt to new, unseen tasks. There exist various approaches to few-shot learning, including i) learning to fine-tune (Finn *et al.*, 2017; Ravi and Larochelle, 2017), ii) sequence based (Mishra *et al.*, 2018; Santoro *et al.*, 2016), and iii) metric-learning based approaches (Vinyals *et al.*, 2016; Snell *et al.*, 2017; Nguyen *et al.*, 2020). Due to its simplicity and efficiency, the latter category has recently received a lot of attention, and the models relevant for this paper build on this principle. Vinyals *et al.* (2016) combined deep feature learning with non-parametric methods

in the Matching Network, by performing weighted nearest-neighbor classification in the embedding space. They proposed to train the model in episodes where a small labeled support set and an unlabeled query image are mapped to the query label, making the model able to adapt to unseen classes without the need for fine-tuning. Whereas the Matching Network only performed one-shot image classification, Snell *et al.* (2017) later proposed the Prototypical Network, which extended the problem to include few-shot classification. Based on the idea that there exists an embedding space, in which samples cluster around their class prototype representation, they proposed a simpler model with a shared encoder between the support and query set, and a nearest-neighbor prototype matching in the embedding space.

2.2. Few-Shot Semantic Segmentation

Few-shot semantic segmentation extends few-shot image classification (Vinyals *et al.*, 2016; Snell *et al.*, 2017; Nguyen *et al.*, 2020) to pixel-level classifications (Shaban *et al.*, 2017; Rakelly *et al.*, 2018; Zhang *et al.*, 2020; Wang *et al.*, 2019), and the goal is to, based on a few densely labeled samples from one (or more) new class(es), segment the class(es) in a new image. A recent line of work builds on the ideas from the Prototypical Network by Snell *et al.* (2017), and can be roughly split into two groups: models where predictions are based directly on the cosine similarity between query features and prototypes in the embedding space (Wang *et al.*, 2019; Liu *et al.*, 2020b; Ouyang *et al.*, 2020), and models that find the correlation between query features and prototypes by employing decoding networks to get the final prediction (Dong and Xing, 2018; Zhang *et al.*, 2019; Liu *et al.*, 2020a; Li *et al.*, 2021a; Zhang *et al.*, 2021; Tian *et al.*, 2020).

Dong and Xing (2018) first adopted the idea of metric-learning based prototypical networks to perform few-shot semantic segmentation. They proposed a two-branched model: a prototype learner, learning class-wise prototypes from the labeled support set, and a segmentation network where the prototypes were used to guide the segmentation of the query image. Most relevant for this work, Wang *et al.* (2019) argued that parametric segmentation generalizes poorly, and proposed the Prototype Alignment Network (PANet), a simpler model where the knowledge extraction and segmentation process is separated. By exploiting prototypes extracted from the semantic classes of the support set, they reduced the segmentation of the query image to a non-parametric pixel-wise nearest-neighbor prototype matching, thereby creating a new branch of FSS models. Building on PANet, (Liu *et al.*, 2020b) addressed the limitation of reducing semantic classes to a simple prototype and proposed the Part-aware Prototype Network (PPNet), where each semantic class is represented by multiple prototypes to capture more diverse features. Liu *et al.* (2020b) further adopted a semantic branch for parametric segmentation during training to learn better representations. Ouyang *et al.* (2020) adapted ideas from PANet to perform FSS in the medical domain. They addressed the major restricting factor preventing medical FSS, e.g. the dependency on a large a set of annotated training classes. This barrier was overcome by the introduction of a superpixel-based self-supervised learning scheme, enabling the training of

FSS networks without the need for labeled data. Ouyang *et al.* (2020) further introduced the Adaptive Local Prototype pooling empowered prototypical Network (ALPNet) where additional *local* prototypes are computed on a regular grid to preserve local information and enhance segmentation performance.

A different approach to medical FSS was suggested by Roy *et al.* (2020), and was the first FSS model for medical image segmentation. Their proposed SE-Net employs squeeze and excite blocks (Hu *et al.*, 2018) in a two-armed architecture consisting of one conditioner arm, processing the support set, and one segmenter arm, interacting with the conditioner arm to segment the query image. However, this model is trained supervised, requiring a set of labeled classes for training.

Based on our experience, training a decoder in a self-supervised setting, where the training task (superpixel segmentation) differs from the inference task (organ segmentation), is challenging and leads to performance degradation. In this paper, we thus, partially inspired by the state-of-the-art model (Ouyang *et al.*, 2020), build further on the branch initiated by Wang *et al.* (2019) to perform FSS in the medical domain. We propose a novel FSS model that, unlike previous approaches in this branch (Wang *et al.*, 2019; Liu *et al.*, 2020b; Ouyang *et al.*, 2020), does *not* explicitly model the complex background class, but relies solely on one foreground prototype.

2.3. Self-Supervised Learning

When large labeled datasets are not available, self-supervision can be used to learn representations by training the deep learning model on an auxiliary task that is defined such that the label is *implicitly* available from the data. A good auxiliary task should require high-level image understanding to be solved, thereby encouraging the network to encode this type of information. Commonly used auxiliary tasks include image inpainting (Larsson *et al.*, 2016; Pathak *et al.*, 2016; Zhang *et al.*, 2016), contrastive learning (Chen *et al.*, 2020; Misra and Maaten, 2020), rotation prediction (Komodakis and Gidaris, 2018), solving jigsaw puzzles (Noroozi and Favaro, 2016), and relative patch location prediction (Doersch *et al.*, 2015).

In the medical domain, self-supervised learning (SSL) has been used to improve performance on other (main) tasks by exploiting unlabeled data in a multi-task learning setting (Chen *et al.*, 2019; Li *et al.*, 2021b) and to pre-train models before transferring them to new (main) tasks (Bai *et al.*, 2019; Zhu *et al.*, 2020; Dong *et al.*, 2021; Lu *et al.*, 2021). In Ouyang *et al.* (2020), SSL was used to train a FSS model completely unsupervised using a novel superpixel-based auxiliary task, removing the need for labeled data during training. We build on this work by extending the proposed self-supervision scheme to 3D supervoxels.

2.4. Supervoxel Segmentation

Supervoxels and superpixels are groupings of local voxels/pixels in an image that share similar characteristics. The boundaries of a supervoxel/superpixel therefore tend to follow the boundaries of the structures in the image, providing natural sub-regions. Supervoxel and superpixel segmentation has

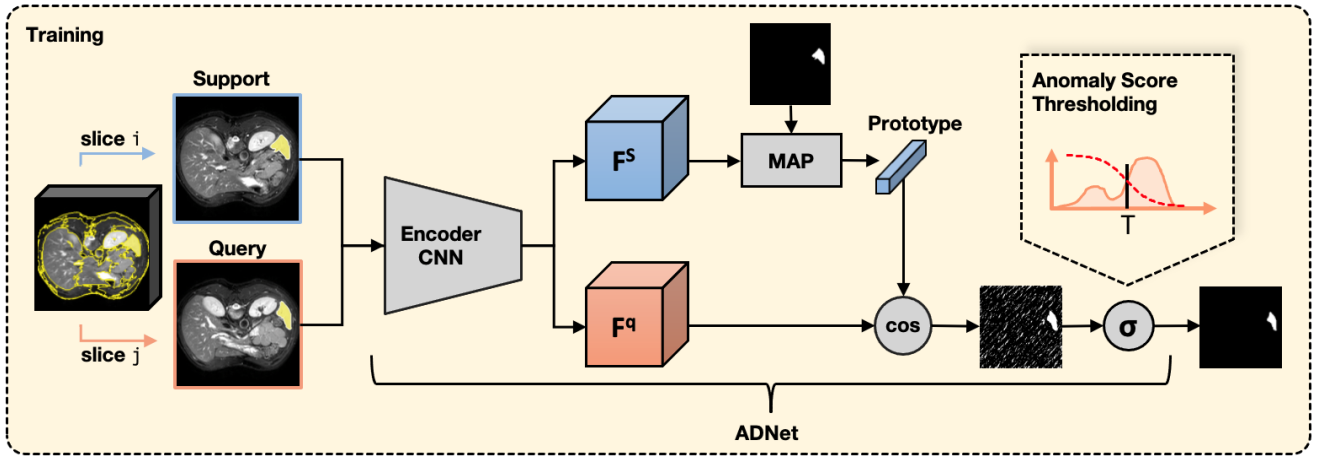


Fig. 1. Illustration of the model during training. Support and query slices are obtained from the same image volume as two different 2D slices containing a randomly sampled supervoxel. A shared feature encoder encodes the query and the support images into deep feature maps. The support features are then resized to the mask size and masked average pooling is applied to compute the foreground prototype. For each query feature vector, an anomaly score is computed based on the cosine similarity to the prototype. Finally, the segmentation of the query image is performed by thresholding the anomaly scores using a learned anomaly threshold.

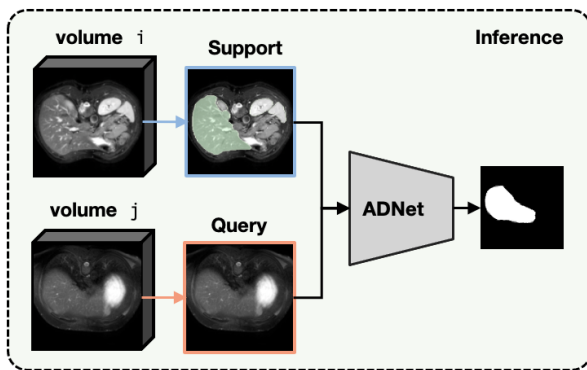


Fig. 2. Illustration of the model during inference. Based on labeled slices from the support volume, the query volume is segmented slice by slice, one class at a time.

become a common tool in computer vision, also in the medical domain (Huang et al., 2020; Irving et al., 2016). For a detailed comparison of available superpixel segmentation algorithms, we refer the reader to (Stutz et al., 2018).

3. Problem Definition

Given a labeled dataset with classes C_{train} (here: $C_{train} = \{supervoxel_1, supervoxel_2, \dots\}$), FSS models aim to learn a quick adaption to new classes C_{test} (e.g. $C_{test} = \{liver, kidney, spleen\}$) when exposed to only a few labeled samples. The training and testing are performed in an episodic manner (Vinyals et al., 2016) where, in each episode, N classes are sampled from C to create a support set and a query set. The input to an episode is the support image(s) (with annotations) and a query image, and the output is the predicted query mask. In an N -way k -shot setting, the support set $S = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_{N \times k}, \mathbf{y}_{N \times k})\}$ consists of k image slices $\mathbf{x} \in \mathbb{R}^{H \times W}$

(with annotations $\mathbf{y} \in \mathbb{R}^{H \times W}$ indicating the class of each pixel) from each of the N classes, whereas the query set consists of one query image $Q = \{(\mathbf{x}_1^*, \mathbf{y}_1^*)\}$ containing one or more of the N classes.

4. Methods

In this work, we propose an anomaly detection-inspired network (ADNet) for prototypical FSS¹. We employ a shared feature extractor between the support and query images and perform metric learning-based segmentation in the embedding space. Unlike prior approaches that obtain prototypes for both foreground *and* background classes (Liu et al., 2020b; Ouyang et al., 2020; Wang et al., 2019), we only consider foreground prototypes to avoid the aforementioned problems related to explicitly modeling the large and heterogeneous background class. Based on *one* foreground prototype, we compute anomaly scores for all query feature vectors. The segmentation of the query image is then based on these anomaly scores and a learned anomaly threshold. To train our model, we take inspiration from Ouyang et al. (2020) and propose a new supervoxel-based self-supervision pipeline. Fig. 1 and Fig. 2 provide an overview of the model during training and inference, respectively.

4.1. Anomaly Detection-Inspired Few-Shot Segmentation

We denote the encoding network as f_θ and start by embedding the support and query images into deep features, $f_\theta(\mathbf{x}) = F^S$ and $f_\theta(\mathbf{x}^*) = F^Q$, respectively. As opposed to previous works, we are only interested in explicitly modeling the foreground in each episode. We do this by employing the segmentation mask to perform masked average pooling (MAP), but only

¹By "anomaly" we refer to abnormalities compared to our defined normal class (foreground), and not necessarily something that occurs infrequently.

for the foreground class c . We resize the support feature map F^s to the mask size (H, W) and compute one foreground prototype $p \in \mathbb{R}^d$, where d is the dimension of the embedding space:

$$p = \frac{\sum_{x,y} F^s(x, y) \odot \mathbf{y}^{fg}(x, y)}{\sum_{x,y} \mathbf{y}^{fg}(x, y)}, \quad (1)$$

where \odot denotes the Hadamard product and $\mathbf{y}^{fg} = \mathbb{1}(\mathbf{y} = c)$ is the binary foreground mask of class c .

To segment the query image based on this *one* class-prototype, we design a threshold-based metric learning approach to the segmentation. We first obtain an anomaly score S for each query feature vector $F^q(x, y)$ by calculating the (negative) cosine similarity to the foreground prototype p of the episode:

$$S(x, y) = -\alpha \frac{F^q(x, y) \cdot p}{\|F^q(x, y)\| \|p\|}, \quad (2)$$

where $\alpha = 20$ is a scaling factor introduced by Oreshkin *et al.* (2018). In this way, query feature vectors that are identical to the prototype will get an anomaly score of $-\alpha$ (minimum), whereas query feature vectors that are pointing in the opposite direction, relative to the prototype, get an anomaly score of α (maximum). The predicted foreground mask is then found by thresholding these anomaly scores with a learned parameter T . To make the process differentiable, we perform soft thresholding by applying a shifted Sigmoid:

$$\hat{\mathbf{y}}_{fg}^q(x, y) = 1 - \sigma(S(x, y) - T), \quad (3)$$

where $\sigma(\cdot)$ denotes the Sigmoid function with a steepness parameter $\kappa = 0.5$. The impact of the steepness parameter is examined in Section 5.3.4. In this way, query feature vectors with an anomaly score below T (similar to the prototype) get a foreground probability above 0.5, whereas query feature vectors with an anomaly score above T (dissimilar to the prototype) get a foreground probability below 0.5. The predicted background mask is finally found as $\hat{\mathbf{y}}_{bg}^q = 1 - \hat{\mathbf{y}}_{fg}^q$.

The predicted foreground and background masks for the query image are then upsampled to the image size (H, W) and we compute the binary cross-entropy segmentation loss:

$$\mathcal{L}_S = -\frac{1}{HW} \sum_{x,y} \mathbf{y}_{bg}^q(x, y) \log(\hat{\mathbf{y}}_{bg}^q(x, y)) + \mathbf{y}_{fg}^q(x, y) \log(\hat{\mathbf{y}}_{fg}^q(x, y)). \quad (4)$$

In order to encourage a compact embedding of the foreground classes, we construct an additional loss term $\mathcal{L}_T = T/\alpha$ that minimizes the learned threshold. The effect of this loss component is examined in Section 5.3.2.

Following common practice (Liu *et al.*, 2020b; Ouyang *et al.*, 2020; Wang *et al.*, 2019), we also add a prototype alignment regularization loss where the roles of support and query are reversed. The *predicted* query mask is used to compute a proto-

type that segments the support image:

$$\mathcal{L}_{PAR} = -\frac{1}{HW} \sum_{x,y} \mathbf{y}_{bg}^s(x, y) \log(\hat{\mathbf{y}}_{bg}^s(x, y)) + \mathbf{y}_{fg}^s(x, y) \log(\hat{\mathbf{y}}_{fg}^s(x, y)). \quad (5)$$

This gives us the overall loss function

$$\mathcal{L} = \mathcal{L}_S + \mathcal{L}_T + \mathcal{L}_{PAR}. \quad (6)$$

4.2. Supervoxel-Based Self-Supervision

The ADNet is parameterized by $\mathcal{P} = \{\theta, T\}$ and trained self-supervised (unsupervised) end-to-end in an episodic manner. For ease of comparison to previous approaches, our baseline setup follows a 2D approach, where volumes are segmented slice-by-slice. However, to better utilize the volumetric nature of the medical images, we propose a new self-supervision task that exploits 3D supervoxels during the model's training phase. As supervoxels are sub-volumes of the image, representing groups of similar voxels in local regions of the image volume, this allows us to sample 3D pseudo-segmentation masks for semantically uniform regions in the image.

In the training phase, each episode is constructed based on one unlabeled image volume and its supervoxel segmentation: First, one random supervoxel is sampled to represent the foreground class, resulting in a binary 3D segmentation mask. Then, we sample two 2D slices from the image containing this "class"/supervoxel to serve as support and query images. By exploiting the relations across slices, we are able to increase the amount of information that can be extracted in the self-supervision task compared to prior approaches. Following Ouyang *et al.* (2020), we additionally apply random transformations to one of the images (query or support) to encourage invariance to shape and intensity differences.

The supervoxels for all image volumes are computed offline using a 3D extension of the same unsupervised segmentation algorithm (Felzenszwalb and Huttenlocher, 2004) as in (Ouyang *et al.*, 2020). This is an efficient graph-based image segmentation algorithm building on euclidean distances between neighboring pixels. In the 3D extension, this corresponds to the distances from each voxel to its 26 nearest neighbours. In medical images, the resolution in z -direction (slice thickness) is typically different from the in-plane (x, y) resolution. To account for this anisotropic voxel resolution, we re-weight all distances along the z -direction (xz -, yz - and xyz -direction) according to the spatial ratios.

The supervoxel generation has one hyper-parameter ρ that controls the minimum supervoxel size, where a larger ρ corresponds to larger and fewer supervoxels. The effect of this parameter on the final segmentation result is examined in Section 5.3.3.

4.3. Implementation Details

The implementation is based on the PyTorch (v1.7.1) implementation of SSL-ALPNet (Ouyang *et al.*, 2020). The encoder network used in all the 2D experiments is a ResNet-101 pretrained on MS-COCO, where the classifier is replaced by a

² $\mathbb{1}(\cdot)$ is the indicator function, returning 1 if the argument is true and 0 otherwise.

1 × 1 convolutional layer to reduce the feature dimension from 2048 to 256. Following ALPNet, we optimize the loss using stochastic gradient descent with momentum 0.9, a learning rate of 1e-3 with a decay rate of 0.98 per 1k epochs, and a weight decay of 5e-4 over 50k iterations. To address the class imbalance, we follow previous work and weigh the foreground and background class in the cross-entropy loss (1.0 and 0.1, respectively). To further stabilize training, we set a minimum threshold of 200 pixels on the supervoxel size in the slices sampled as support/query. Supervoxel generation is done offline (once per image volume) and is relatively computationally efficient³. Training takes 1.8h on a Nvidia RTX 2080Ti GPU.

5. Experiments

5.1. Setup

5.1.1. Data

We assess the proposed method on representative publicly available datasets⁴:

- (1) **MS-CMRSeg** (bSSFP fold), from the MICCAI 2019 Multi-sequence Cardiac MRI Segmentation Challenge, containing 35 3D cardiac MRI scans with on average 13 slices (Zhuang, 2018, 2016).
- (2) **CHAOS**, from the ISBI 2019 Combined Healthy Abdominal Organ Segmentation Challenge (task 5), containing 20 3D T2-SPIR MRI scans with on average 36 slices (Kavur *et al.*, 2021, 2019, 2020).

To compare our results to Ouyang *et al.* (2020), we follow the same pre-processing scheme: 1) Cut the top 0.5% intensities. 2) Re-sample image slices (short-axis slices for the cardiac images and axial slices for the abdominal images) to the same spatial resolution. 3) Crop slices to unify size (256×256 pixels). Further, to fit into the pretrained network, each slice is repeated three times along the channel dimension.

In all experiments, the models are trained self-supervised (unsupervised) and evaluated in a five-fold cross-validation manner, where, in each fold, the support images are sampled from *one* of the patients and the remaining patients are treated as query (see Fig. 3). Furthermore, to account for the stochasticity in the model and optimization, we repeat each fold three times. In the cardiac MRI scans we segment three classes: Left-ventricle blood pool (LV-BP), left-ventricle myocardium (LV-MYO) and right-ventricle (RV). In the abdominal MRI scans, we segment four classes: left kidney (L. kid.), right kidney (R. kid.), liver, and spleen. Following previous methods (Ouyang *et al.*, 2020; Roy *et al.*, 2020), each class is segmented separately in binary foreground/background segmentation problems⁵. Since the models are trained self-supervised, we do *not* exclude image slices that contain the target classes.

³The compute time for generating all supervoxels for the MS-CMRSeg dataset is less than 3 minutes using a Quad-Core Intel Core i7 processor.

⁴Links to public datasets: [MS-CMRSeg](#) and [CHAOS](#)

⁵As the segmentation only relies on the computation of the cosine similarity to a class-specific prototype and a threshold which is shared among classes, the proposed method may be extended to account for multi-class scenarios. A detailed analysis of this is left for future work.

Split 0 :	1	2	3	4	5	6	7	8
Split 1 :	8	9	10	11	12	13	14	15
Split 2 :	15	16	17	18	19	20	21	22
Split 3 :	22	23	24	25	26	27	28	29
Split 4 :	29	30	31	32	33	34	35	1

■ Query ■ Support

Fig. 3. Setup for the five-fold cross-validation. This illustrates how the patient IDs are distributed among the splits and how the support/query volumes are selected for the cardiac MRI dataset. For each fold, a model is trained on all images *not* present in that fold. During inference, the left-out fold is used exclusively, where the labeled support image is exploited to segment the query images slice by slice, class by class. The CHAOS dataset is split into five folds in a similar manner.

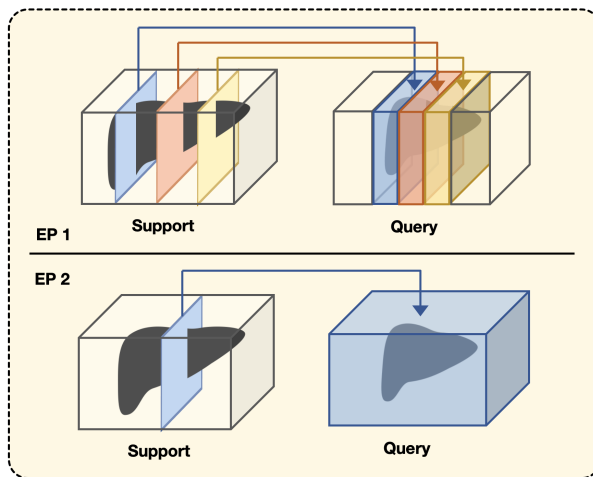


Fig. 4. Illustration of EP1 (top) and EP2 (bottom). In EP1, the support and query volumes are divided into three succeeding sub-chunks. The middle slice in each sub-chunk of the support volume is labeled and used to segment all the slices in the corresponding sub-chunk in the query volume. This means that the protocol requires weak labels indicating where the class of interest is located in the query volume. In EP2, the middle slice of the support volume is labeled and used to segment *all* slices in the query volume, avoiding the need for additional weak labels.

5.1.2. Evaluation metric

Following common practice (Ouyang *et al.*, 2020; Roy *et al.*, 2020) we employ the mean dice score to compare the model predictions to the ground truth segmentations. The dice score, D , between two segmentations A and B is given by

$$D(A, B) = 2 \frac{|A \cap B|}{|A| + |B|} \cdot 100\%, \quad (7)$$

meaning that a dice score of 100% corresponds to a perfect match between the segmentations.

5.1.3. Evaluation protocols

During inference, the query volumes are segmented episode-wise, slice-by-slice, based on labeled support slices. For this reason, it is necessary to define an evaluation protocol that describes how to construct the episodes during inference, i.e. how

Table 1. Mean dice score and standard deviation over three runs per split under EP1.

Method	Cardiac MRI				Abdominal MRI				
	LV-BP	LV-MYO	RV	Mean	L kid.	R kid.	Spleen	Liver	Mean
pSSL-PANet	80.20 ± 4.39	45.67 ± 2.58	66.95 ± 4.65	64.27 ± 14.23	63.09 ± 9.31	66.09 ± 8.73	63.93 ± 8.65	72.08 ± 3.83	66.30 ± 3.51
pSSL-ALPNet	87.54 ± 1.63	60.19 ± 4.55	76.08 ± 4.72	74.60 ± 11.21	81.00 ± 4.01	84.66 ± 2.40	72.32 ± 7.69	75.89 ± 3.02	78.46 ± 4.72
vSSL-PPNet	67.78 ± 8.31	42.61 ± 6.16	60.80 ± 6.44	57.06 ± 10.61	62.13 ± 7.85	71.78 ± 11.04	66.57 ± 9.04	73.12 ± 2.51	68.40 ± 4.37
vSSL-CANet	78.99 ± 4.72	43.61 ± 3.38	61.10 ± 3.60	61.07 ± 14.64	69.53 ± 12.05	77.15 ± 10.71	67.05 ± 6.87	72.88 ± 3.27	71.65 ± 3.79
vSSL-ADNet	87.53 ± 2.03	62.43 ± 3.98	77.31 ± 3.48	75.76 ± 10.31	75.28 ± 14.80	83.28 ± 13.36	75.92 ± 8.90	80.81 ± 2.36	78.82 ± 3.35

Table 2. Mean dice score and standard deviation over three runs per split under EP2. * indicates that the increase in mean dice score for the best performing model is statistically significant ($p < 0.05$).

Method	Cardiac MRI				Abdominal MRI				
	LV-BP	LV-MYO	RV	Mean	L kid.	R kid.	Spleen	Liver	Mean
pSSL-PANet	68.28 ± 5.67	38.60 ± 3.72	55.22 ± 5.18	54.03 ± 12.15	32.85 ± 6.74	30.18 ± 4.85	34.82 ± 8.52	53.89 ± 3.15	37.94 ± 9.36
pSSL-ALPNet	80.65 ± 3.93	53.31 ± 6.31	69.25 ± 2.80	67.74 ± 11.21	56.42 ± 5.74	50.37 ± 7.77	44.70 ± 7.77	56.73 ± 3.07	52.05 ± 4.94
vSSL-PPNet	56.69 ± 8.35	34.78 ± 7.30	47.60 ± 6.07	46.35 ± 8.99	43.36 ± 10.44	56.94 ± 14.39	43.06 ± 9.73	56.32 ± 7.57	49.92 ± 6.71
vSSL-CANet	74.54 ± 4.20	35.08 ± 4.30	47.65 ± 5.73	52.42 ± 16.46	50.18 ± 13.02	69.91 ± 12.84	48.84 ± 8.61	64.00 ± 3.44	58.23 ± 8.98
vSSL-ADNet	82.81 ± 3.20	59.46 ± 2.97	66.58 ± 4.74	69.62 ± 9.77*	62.33 ± 9.70	86.46 ± 2.74	63.73 ± 11.66	77.12 ± 3.41	72.41 ± 9.96*

to pair support and query images in episodes. In the experiments, we evaluate all models under two different evaluation protocols (EPs), illustrated in Fig. 4.

Evaluation protocol 1 (EP1). Previous works (Ouyang et al., 2020; Roy et al., 2020) follow an evaluation protocol that requires weak labels for all query images, i.e. there is a need to indicate (label) in which slices the foreground class is located. For a given class to be segmented, the chunk of slices in both the support and query volumes containing this class is divided into three succeeding sub-chunks. The middle slice in each sub-chunk of the support volume is used to segment all the slices in the corresponding sub-chunk in the query. In practice, this requires manual and time-consuming input from medical experts during the inference phase, where they have to scroll through each query image volume to mark the slices containing the class(es) of interest.

Evaluation protocol 2 (EP2). To avoid the need for weak query labels during inference, we introduce a new evaluation protocol that does not depend on the position of the target volume, and thus is more applicable in practical situations. Here, we simply sample $k = 1$ slices from the support foreground volume and use this information to segment the entire query volume. To limit boundary effects, we choose the middle slice of the support foreground volume.

5.2. Comparison to state-of-the-art

We compare our model to three modern FSS models: PANet (Wang et al., 2019), ALPNet Ouyang et al. (2020), and PPNet (Liu et al., 2020b) with five (default) prototypes per class. Additionally, to compare our one-prototype anomaly approach to a one-prototype decoder approach, we adopt the dense comparison module proposed in (Zhang et al., 2019) as

Table 3. Summarized information about the models. *The number of prototypes in ALPNet is adaptive and we report the average number over all classes during inference.

Method	Backbone	Decoder	# Foreground prototypes	# Background prototypes
pSSL-PANet		✗	1	1
pSSL-ALPNet*		✗	4	246
vSSL-PPNet	ResNet-101	✗	5	5
vSSL-CANet		✓	1	0
vSSL-ADNet		✗	1	0

a decoder on top of the backbone network and refer to this network as CANet⁶.

The current state-of-the-art method for medical FSS, Ouyang et al. (2020), showed that training PANet and ALPNet in a self-supervised manner improved the dice scores of the segmentation results considerably, compared to classical supervised FSS. Specifically, the dice scores on the MS-CMRSeg and CHAOS datasets increased by an average of 17.9 and 26.1 percentage points, respectively. Here, we are thus only focusing on SSL approaches. pSSL refers to the superpixel SSL approach presented in Ouyang et al. (2020), whereas vSSL refers to our proposed supervoxel-based approach.

Table 1 and Table 2 present the results under EP1 and EP2, respectively, as mean and standard deviations over three runs (over all splits). Summarized details about the models can be found in Table 3.

In Table 1 we can see that our proposed model under EP1 performs similarly to the state-of-the-art on both datasets, while using significantly fewer prototypes compared to the closest competitors. We can also observe that the models that use just a few

⁶Code available: PANet, ALPNet, PPNet, and CANet.

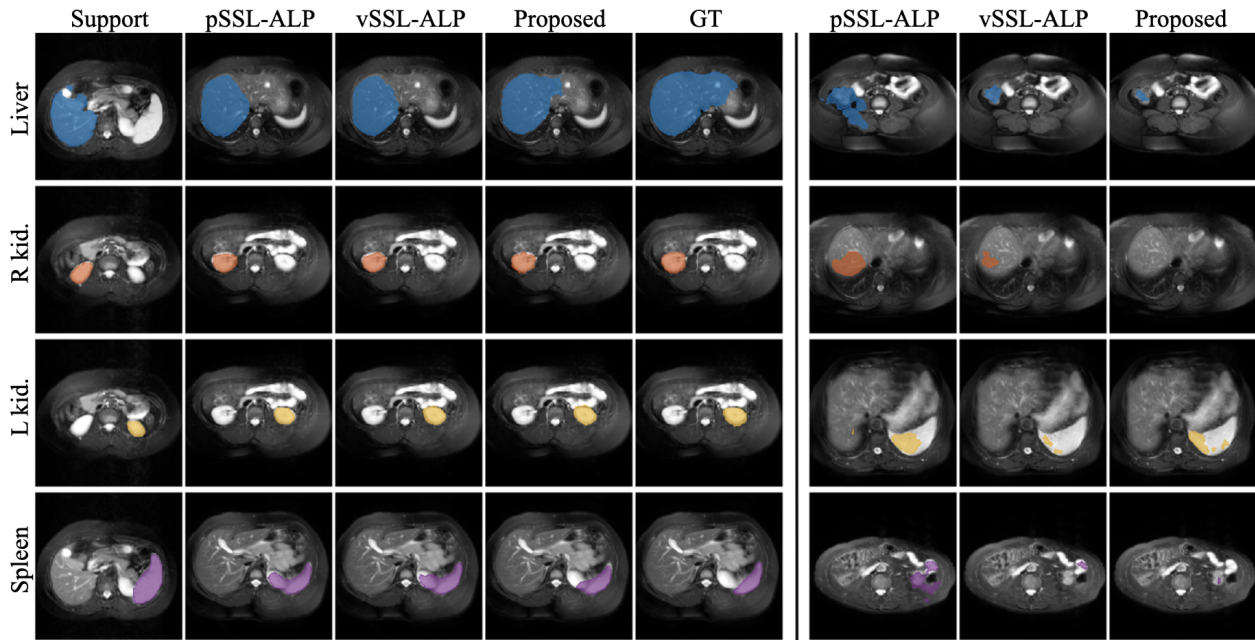


Fig. 5. Qualitative comparisons for the abdominal MRI dataset. To the left of the solid line, we see (left to right) the support image, the segmentation results of a query slice containing the foreground class, and the ground truth segmentation of this query image. To the right, we see segmentation results for query slices *not* containing the foreground class. Top to bottom: liver, right kidney, left kidney, and spleen. The proposed method is more robust to background outside the support slice, resulting in less over-segmentation.

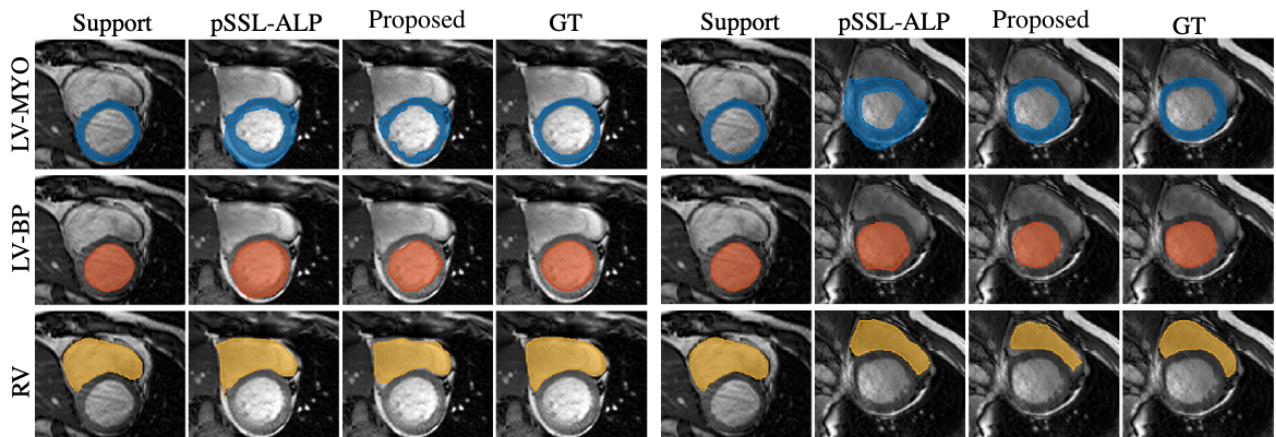


Fig. 6. Qualitative comparisons for two episodes with the same support volume from the cardiac MRI dataset. Left to right: Support image, segmentation results of a query slice, and ground truth segmentation of this query image. The segmentation results are quite similar but the proposed method captures the left-ventricle myocardium and left ventricle blood pool better, with less over-segmentation.

prototypes to model the background (PANet, PPNet) perform poorly and are among the three worst performing models for both datasets. Furthermore, by only modeling the foreground class and segmenting the query image using a decoding network, CANet results in the lowest (overall) dice score on the cardiac dataset.

In a more realistic scenario, information about the location of the foreground volume in the query images is typically not available. We therefore evaluate the models under EP2 (Table 2) and we observe that our proposed approach outperforms the state-of-the-art. One-sided Wilcoxon signed rank

tests (Wilcoxon, 1992) on the mean dice scores across all runs indicate a significant difference between the segmentation results obtained from vSSL-ADNet and pSSL-ALPNet for both datasets under EP2 ($p < 0.05$). For the abdominal data, our model improves the segmentation results by more than 20 percentage points compared to pSSL-ALPNet. The main reason for this large improvement is that we now have to consider *all* the query slices (not only the slices containing the organ to be segmented), meaning that the background class is much larger and much more diverse. This again complicates the task of modeling the background with prototypes, whereas our

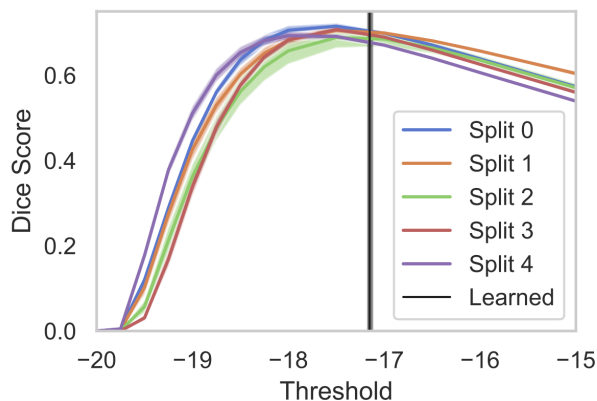


Fig. 7. Analysis of the precision of the learned threshold. The plot shows the mean dice score (with standard deviation) obtained for a range of thresholds during inference on the MS-CMRSeg dataset. The learned threshold is indicated by the black vertical line.

Table 4. Ablation study showing how the loss function components affect the results under EP1. * indicates that the increase in mean dice score for the best performing model is statistically significant ($p < 0.05$).

\mathcal{L}_S	\mathcal{L}_T	\mathcal{L}_{PAR}	Cardiac MRI			
			LV-BP	LV-MYO	RV	Mean
✓	✓	✓	87.53 ± 2.03	62.43 ± 3.98	77.31 ± 3.48	75.76 ± 10.31*
✓		✓	87.41 ± 2.08	58.48 ± 3.17	74.95 ± 3.33	73.61 ± 11.85
✓	✓		82.80 ± 3.26	57.70 ± 3.05	72.63 ± 2.43	71.05 ± 10.31
✓			83.62 ± 2.51	51.38 ± 2.52	68.36 ± 3.02	67.77 ± 13.17

anomaly detection-inspired model without background prototypes is less affected. The somewhat lower performance and high standard deviation for left-kidney and spleen are related to the weak boundaries between these organs (see discussion in Section 6). Furthermore, we obtain considerable, but smaller, improvements on the cardiac dataset under EP2. This is related to the lower number of slices and the less diverse background in these images, making the task of modeling the background with prototypes less complicated. Qualitative comparisons are provided in Fig. 5 and Fig. 6, where we can see that our approach is less prone to over-segmentation.

5.3. Model analysis

5.3.1. Analysis of learned threshold

To evaluate the learned threshold’s precision on the unseen test data, we have conducted a line search where we, in the inference phase, evaluate the dice score obtained using a range of different thresholds between -20 and -15. The experiment was performed on three runs for each split and the mean dice score and standard deviation (shaded region) are reported in Fig. 7. The learned threshold is averaged over all runs and represented by the vertical black line⁷. From the plot, we see that

⁷The small, gray shaded region indicates the range of learned threshold values.

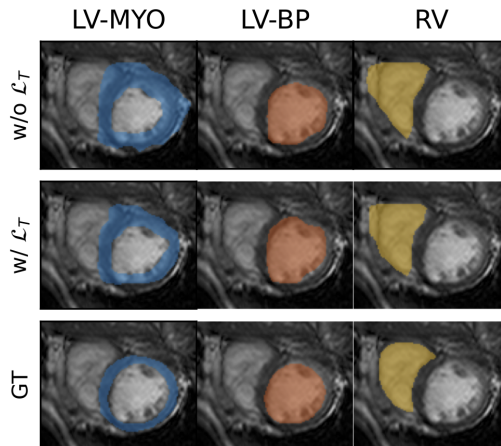


Fig. 8. Qualitative (zoomed in) segmentation results for one slice in the MS-CMRSeg dataset obtained from a model trained with (middle) and without (top) \mathcal{L}_T in the total loss. The lower row shows the ground truth, and it is evident that the threshold loss reduces the over-segmentation, especially for the left-ventricle myocardium.

the threshold optimized for the training data is close to the ideal threshold for the test data, with little to gain in terms of increased dice score.

5.3.2. Ablation study

To evaluate the effect of the three components of our loss function, we conduct an ablation study on the cardiac dataset. Table 4 illustrates that \mathcal{L}_T and \mathcal{L}_{PAR} improve the dice score across all classes. Further, Fig. 8 shows qualitatively the effect of \mathcal{L}_T on the segmentation of one image slice from the MS-CMRSeg dataset. Here, it can be seen how the encouraging of a more compact foreground embedding via \mathcal{L}_T reduces the over-segmentation, especially for the left-ventricle myocardium.

5.3.3. Sensitivity of supervoxel size

A sensitivity analysis of the parameter ρ , controlling the supervoxel size, is conducted on the MS-CMRSeg dataset and the results are presented in Table 5. As shown by these results, the final segmentation performance is relatively robust for a range of minimum size values from $\rho = 1000$ to $\rho = 2000$. However, if we allow the sizes to become too small ($\rho = 500$) or too large ($\rho = 5000$), we see that the performance is negatively affected. Examples of 2D slices from the 3D supervoxel segmentations for the different values of ρ are shown in Fig. 9.

According to the sensitivity study, a reasonable value is $\rho = 1000$, and all the reported vSSL results are obtained with this value for the MS-CMRSeg dataset and $\rho = 5000$ for the CHAOS dataset, unless otherwise stated. The difference in value of ρ reflects the differences in volume size.

5.3.4. Influence of steepness parameter

The steepness of the sigmoid function controls how soft the threshold operation performed is. If the steepness is high

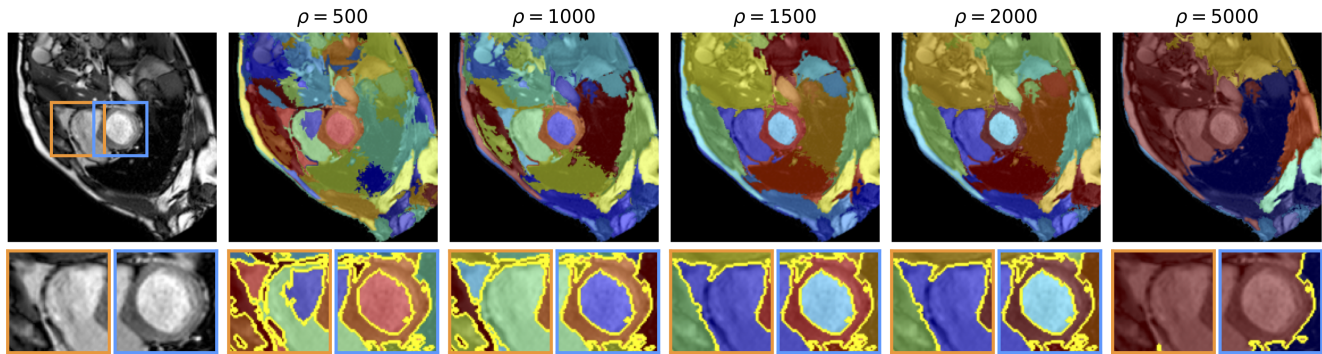


Fig. 9. Examples of supervoxel segmentation results in one slice from the MS-CMRSeg dataset for different values of ρ . The parameter ρ controls the minimum size of a supervoxel for it not to be joined with an adjacent supervoxel. A larger ρ corresponds to larger and fewer supervoxels.

Table 5. Supervoxel parameter sensitivity. Analysis of the parameter controlling the minimum supervoxel size (n.o. voxels), on the cardiac MRI dataset under EP1.

ρ	Cardiac MRI			
	LV-BP	LV-MYO	RV	Mean
500	84.64 \pm 1.51	43.48 \pm 6.25	66.87 \pm 5.00	65.00 \pm 16.86
1000	87.53 \pm 2.03	62.43 \pm 3.98	77.31 \pm 3.48	75.76 \pm 10.31
1500	86.91 \pm 2.47	62.60 \pm 3.44	75.30 \pm 1.91	74.94 \pm 9.93
2000	87.30 \pm 1.80	61.21 \pm 3.33	73.92 \pm 2.51	74.14 \pm 10.65
5000	77.84 \pm 8.49	49.30 \pm 7.93	66.44 \pm 10.1	64.53 \pm 14.72

Table 6. Steepness parameter sensitivity. Analysis of the parameter controlling the sigmoid steepness parameter, on the cardiac MRI dataset under EP1.

κ	Cardiac MRI			
	LV-BP	LV-MYO	RV	Mean
0.1	80.69 \pm 4.04	17.10 \pm 5.11	52.69 \pm 10.35	50.16 \pm 26.02
0.3	87.95 \pm 1.35	56.79 \pm 6.32	78.44 \pm 2.48	74.39 \pm 13.04
0.5	87.54 \pm 2.03	62.44 \pm 3.98	77.32 \pm 3.48	75.76 \pm 10.31
0.7	87.22 \pm 2.13	62.21 \pm 2.71	76.90 \pm 3.40	75.45 \pm 10.26
0.9	85.41 \pm 2.90	60.67 \pm 3.57	76.06 \pm 3.33	74.05 \pm 10.20
1.0	85.68 \pm 2.81	59.40 \pm 4.04	75.22 \pm 4.00	73.43 \pm 10.80

(harder thresholding), the class assignments of samples becomes harder, also close to the threshold. To examine the influence of the steepness parameter, κ , on the final segmentation results, we have conducted six experiments with different values of κ , from $\kappa = 0.1$ to $\kappa = 1.0$ on the MS-CMRSeg dataset⁸. The results presented in Table 6 indicate the model’s robustness with respect to this parameter, and we can observe a gain of more than two percentage points in the dice score by decreasing the steepness from 1.0 to 0.5.

⁸Note that this is equivalent to changing the scaling between $\alpha = 0.2$ and $\alpha = 20$ in Eq. (2).

5.3.5. vSSL vs pSSL

To disentangle and isolate the effect from the proposed extension of the self-supervision task, we have conducted additional experiments where we train our proposed model (ADNet), and the closest competing model (ALPNet) with the two different self-supervision tasks. From the results in Table 7, we see that the supervoxels overall yield better or comparable results for both models. For our proposed ADNet, there is a significant improvement ($p < 0.05$) in dice score from pSSL to vSSL for both datasets. Moreover, the improvements appear most prominent for the abdominal dataset, which is assumed to be related to the nature of the image volumes: In the abdominal dataset, the image volumes contain more slices and more potential information to utilize when the self-supervision task is extended to 3D, compared to the cardiac dataset.

A different implication of the proposed extension to supervoxel-based self-supervision is the enabling of training 3D CNNs for direct volume segmentation, as discussed in the next section.

5.4. Extension to one-step volume segmentation

Thus far, we have adopted a hybrid strategy to 3D segmentation, following Ouyang *et al.* (2020), where the 3D image volumes are segmented slice by slice, independently. However, a natural extension that is facilitated by the new self-supervision task is to adopt a 3D CNN as backbone to process the volumes in one step, thereby fully exploiting the potentially useful information along the third axis. Unfortunately, the high memory consumption and computational cost of 3D CNNs has limited their use to smaller images (in number of voxels), often obtained by down-sampling the original images (Çiçek *et al.*, 2016) or by patch-based approaches (Huo *et al.*, 2019).

To investigate the potential of utilizing 3D convolutions to do one-step 3D segmentations within our proposed framework, we employ a 3D ResNeXt-101 (Hara *et al.*, 2018), which is the 3D extension of ResNeXt (Xie *et al.*, 2017), pretrained on the Kinetics-600 dataset (Kay *et al.*, 2017), as our encoder network. The 3D ResNeXt-101 is a more resource efficient network, compared to the 3D ResNet-101, with approximately half

Table 7. Mean dice score and standard deviation over three runs per split for ADNet and ALPNet with superpixel-based and supervoxel-based self-supervision. * indicates that the increase in mean dice score for the best performing model is statistically significant ($p < 0.05$).

Model	pSSL	vSSL	Cardiac MRI				Abdominal MRI				
			LV-BP	LV-MYO	RV	Mean	L kid.	R kid.	Spleen	Liver	Mean
ALPNet	✓		80.65 ± 3.93	53.31 ± 6.31	69.25 ± 2.80	67.74 ± 11.21	56.42 ± 5.74	50.37 ± 7.77	44.70 ± 7.77	56.73 ± 3.07	52.05 ± 4.94
		✓	79.44 ± 2.79	57.64 ± 3.96	61.22 ± 4.22	66.10 ± 9.59	68.19 ± 12.30	82.45 ± 6.27	55.39 ± 10.09	66.38 ± 5.20	68.10 ± 9.62*
ADNet	✓		78.25 ± 9.68	54.59 ± 6.42	66.37 ± 4.73	66.40 ± 12.07	49.65 ± 7.59	59.00 ± 13.77	52.47 ± 9.56	54.78 ± 3.87	53.97 ± 10.00
		✓	82.81 ± 3.20	59.46 ± 2.97	66.58 ± 4.74	69.62 ± 9.77*	62.33 ± 9.70	86.46 ± 2.74	63.73 ± 11.66	77.12 ± 3.41	72.41 ± 9.96*

Table 8. Mean dice score and standard deviation over three runs per split for vSSL-ADNet with 2D ResNet-101 as backbone and 3D ResNeXt-101 as backbone (under EP2). * indicates that the increase in mean dice score for the best performing model is statistically significant ($p < 0.05$).

Backbone	Params	Labeled slices, k	Cardiac MRI				Abdominal MRI				
			LV-BP	LV-MYO	RV	Mean	L kid.	R kid.	Spleen	Liver	Mean
2D ResNet-101	42.50M	One	82.81 ± 3.20	59.46 ± 2.97	66.58 ± 4.74	69.62 ± 9.77	62.33 ± 9.70	86.46 ± 2.74	63.73 ± 11.66	77.12 ± 3.41	72.41 ± 9.96
3D ResNeXt-101	47.52M	One	81.28 ± 2.51	56.47 ± 0.75	66.22 ± 4.24	67.99 ± 10.60	77.95 ± 16.57	73.55 ± 28.69	75.04 ± 8.55	75.48 ± 8.58	75.50 ± 17.71
3D ResNeXt-101	47.52M	All	82.87 ± 1.15	56.30 ± 0.76	67.93 ± 4.04	69.03 ± 11.15	81.06 ± 4.20	84.88 ± 4.22	75.18 ± 8.40	77.17 ± 8.60	79.58 ± 7.67*

as many trainable parameters in total. The number of parameters is comparable to the 2D ResNet-101 (see Table 8).

To retain the same spatial resolution in the embedding space as for our 2D backbone, we modify the network by *i*) removing the maxpooling in z-direction and *ii*) changing the strides in conv 3, conv 4, and conv 5 to (1, 2, 2), (1, 1, 1), and (1, 1, 1), respectively (see architecture details in Table 9). Similarly to the 2D ResNet-101, we replace the classifier with $1 \times 1 \times 1$ convolutions to reduce the feature dimension from 2048 to 256. Each voxel is repeated three times along the channel dimension in the input to fit into the pretrained network. The network is trained self-supervised end-to-end on 3D patches of size (10, 215, 215), and the loss is optimized according to Section 4.3. During inference, we evaluate the performances under EP2 with two different levels of supervision: *i*) Only labeling the middle slice of the target class in the support volume ($k = one$), as is done in the 2D experiments. *ii*) Labeling all the support slices ($k = all$) and computing one prototype for the entire support volume, which is enabled by the volume-wise embedding.

Table 8 provides a summary of the performance of vSSL-ADNet with 3D ResNeXt-101 and 2D ResNet-101 backbones. Though it is difficult to directly compare 2D CNNs and 3D CNNs for many different reasons, such as difference in pre-training datasets and the number of weights modelling relations within slices and between slices, the results are meant to indicate the potential of using 3D convolutions in our framework to perform one-step 3D segmentation.

From the results on the cardiac dataset, we see that the differences between 2D and 3D are relatively small, which agrees with observations in previous work (Vesal *et al.*, 2019). In the abdominal dataset, on the other hand, there appears to be a greater potential for utilizing the 3D structure via 3D convolutions. This mirrors our results from Section 5.3.5, where we found that the abdominal dataset benefited more from extending the self-supervision task from superpixels to supervoxels.

The largest performance difference between the backbones

Table 9. Modified 3D ResNeXt-101 architecture with cardinality $C = 32$ used as backbone in the 3D experiments.

Layer name	Output size	Architecture
conv 1	$(1, \frac{1}{2}, \frac{1}{2})$	$7 \times 7 \times 7$, 64, stride 1, 2, 2
conv 2	$(1, \frac{1}{4}, \frac{1}{4})$	$1 \times 3 \times 3$ max pool, stride 1, 2, 2
		$1 \times 1 \times 1$, 128, stride 1, 1, 1
		$3 \times 3 \times 3$, 128, stride 1, 1, 1, $C = 32$ $\times 3$
conv 3	$(1, \frac{1}{8}, \frac{1}{8})$	$1 \times 1 \times 1$, 256, stride 1, 1, 1
		$3 \times 3 \times 3$, 128, stride 1, 2, 2, $C = 32$ $\times 4$
		$1 \times 1 \times 1$, 512, stride 1, 1, 1
conv 4	$(1, \frac{1}{8}, \frac{1}{8})$	$1 \times 1 \times 1$, 512, stride 1, 1, 1
		$3 \times 3 \times 3$, 512, stride 1, 1, 1, $C = 32$ $\times 23$
		$1 \times 1 \times 1$, 1024, stride 1, 1, 1
conv 5	$(1, \frac{1}{8}, \frac{1}{8})$	$1 \times 1 \times 1$, 1024, stride 1, 1, 1
		$3 \times 3 \times 3$, 1024, stride 1, 1, 1, $C = 32$ $\times 3$
		$1 \times 1 \times 1$, 2048, stride 1, 1, 1
conv 6	$(1, \frac{1}{8}, \frac{1}{8})$	$1 \times 1 \times 1$, 256, stride 1, 1, 1

can be observed for the left kidney and spleen classes. While the 2D CNN results in a segmentation where these classes are confused, the 3D CNN leads to a better separation between the classes, as illustrated in Fig. 10. We further observe a drop in performance on the right kidney class for the 3D CNN with $k = 1$, which demonstrates the importance of having good support features to achieve robust results with the 3D backbone.

6. Limitations and Outlook

The key observation leading to our anomaly-detection inspired few-shot medical image segmentation is that the foreground class typically is relatively homogeneous. By only modeling the foreground class with a single prototype, we avoid having to model the large and highly inhomogeneous back-

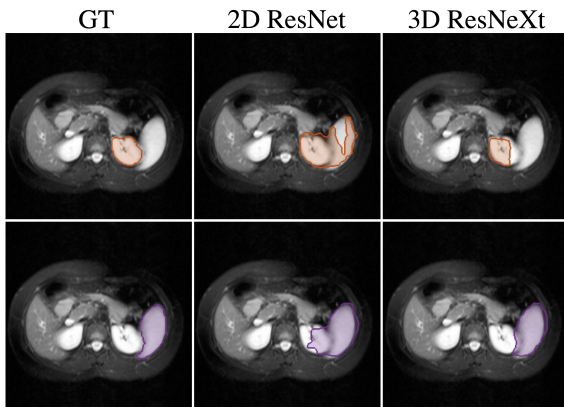


Fig. 10. Comparison of the segmentation results for the left kidney (orange, top) and spleen (purple, bottom) classes for vSSL-ADNet with 2D ResNet-101 and 3D ResNeXt-101 as backbone. The 3D CNN leads to a better separation between the classes.

ground, which we believe is the main challenge in prototypical few-shot medical image segmentation. However, if our assumption of a relatively homogeneous foreground class is not met, and the foreground consists of multiple distinct regions with strong edges, e.g. combining left-ventricle blood pool and left-ventricle myocardium into one foreground class (left-ventricle), modeling the foreground with one prototype might not be sufficient. This is related to the nature of the supervoxels, which tend to follow the boundaries of the structures in the image; Left-ventricle blood pool and left-ventricle myocardium will typically belong to different supervoxels during training and the network therefore learns to separate their feature representations into different clusters. To be able to capture this combined foreground class during inference, one option could be to take inspiration from PpNet (Liu *et al.*, 2020b) and cluster the features into multiple foreground prototypes and then merge the results.

Both the superpixel-based and the supervoxel-based self-supervision tasks are inevitably vulnerable to merging different classes during training *if* the boundaries between them are weak: If the boundaries are weak, the classes will end up in the same superpixel/voxel and the network learns to embed the classes into the same cluster, which makes them difficult to separate during inference. Moreover, in the supervoxel case, it is enough for *one* slice to contain a weak boundary between the classes before they leak into the same supervoxel. This is something that happens between the left-kidney and the spleen in the abdominal dataset, and leads to confusion between these two classes during inference, thereby resulting in lower dice scores and high standard deviations. Taking into account this weak/noisy nature of the supervoxel pseudo-labels is a promising direction for future research.

7. Conclusion

In this work, we proposed a novel and end-to-end trainable anomaly detection-inspired FSS network for medical image segmentation. By approaching the segmentation task as an

anomaly detection problem, our model eliminates the need to explicitly model the large and heterogeneous background class. Moreover, to train the model in an unsupervised manner, we introduced a new self-supervision task that captures the 3D nature of the data by utilizing supervoxels. We assessed our proposed model on representative datasets for cardiac segmentation and abdominal organ segmentation, and showed that it improves segmentation performance and robustness, especially in the realistic scenario where no weak labels for the query images are assumed. Furthermore, we demonstrated how the proposed model, together with the new self-supervision task, has the potential to perform one-step 3D segmentation of the entire image volumes. We believe that fully exploiting the 3D nature of the medical images in this manner for few-shot segmentation represents an interesting line of research for future work.

Acknowledgements

This work was supported by The Research Council of Norway (RCN), through its Centre for Research-based Innovation funding scheme [grant number 309439] and Consortium Partners; RCN FRIPRO [grant number 315029]; RCN IKTPLUSS [grant number 303514]; and the UiT Thematic Initiative.

References

- Abdeltawab, H., Khalifa, F., Taher, F., Alghamdi, N.S., Ghazal, M., Beache, G., Mohamed, T., Keynton, R., El-Baz, A., 2020. A deep learning-based approach for automatic segmentation and quantification of the left ventricle from cardiac cine mr images. *Computerized Medical Imaging and Graphics* 81, 101717.
- Bai, W., Chen, C., Tarroni, G., Duan, J., Guitton, F., Petersen, S.E., Guo, Y., Matthews, P.M., Rueckert, D., 2019. Self-supervised learning for cardiac mr image segmentation by anatomical position prediction, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 541–549.
- Chandola, V., Banerjee, A., Kumar, V., 2009. Anomaly detection: A survey. *ACM computing surveys (CSUR)* 41, 1–58.
- Chen, L., Bentley, P., Mori, K., Misawa, K., Fujiwara, M., Rueckert, D., 2019. Self-supervised learning for medical image analysis using image context restoration. *Medical image analysis* 58, 101539.
- Chen, T., Kornblith, S., Norouzi, M., Hinton, G., 2020. A simple framework for contrastive learning of visual representations, in: *International conference on machine learning*, PMLR. pp. 1597–1607.
- Chen, X., Sun, S., Bai, N., Tang, H., Liu, Q., Yao, S., Han, K., Zhang, C., Lu, Z., Huang, Q., *et al.*, 2021. A deep learning-based auto-segmentation system for organs-at-risk on whole-body computed tomography images for radiation therapy. *Radiotherapy and Oncology*.
- Çiçek, Ö., Abdulkadir, A., Lienkamp, S.S., Brox, T., Ronneberger, O., 2016. 3d u-net: learning dense volumetric segmentation from sparse annotation, in: *International conference on medical image computing and computer-assisted intervention*, Springer. pp. 424–432.
- Doersch, C., Gupta, A., Efros, A.A., 2015. Unsupervised visual representation learning by context prediction, in: *Proceedings of the IEEE international conference on computer vision*, pp. 1422–1430.
- Dong, N., Kampffmeyer, M., Voiculescu, I., 2021. Self-supervised multi-task representation learning for sequential medical images, in: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, Springer. pp. 779–794.
- Dong, N., Xing, E.P., 2018. Few-shot semantic segmentation with prototype learning., in: *British Machine Vision Conference*, British Machine Vision Association.
- Felzenszwalb, P.F., Huttenlocher, D.P., 2004. Efficient graph-based image segmentation. *International journal of computer vision* 59, 167–181.

- Finn, C., Abbeel, P., Levine, S., 2017. Model-agnostic meta-learning for fast adaptation of deep networks, in: International Conference on Machine Learning, Proceedings of Machine Learning Research, pp. 1126–1135.
- Hara, K., Kataoka, H., Satoh, Y., 2018. Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet?, in: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, pp. 6546–6555.
- Hu, J., Shen, L., Sun, G., 2018. Squeeze-and-excitation networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7132–7141.
- Huang, Q., Huang, Y., Luo, Y., Yuan, F., Li, X., 2020. Segmentation of breast ultrasound image with semantic classification of superpixels. *Medical image analysis* 61, 101657.
- Huo, Y., Xu, Z., Xiong, Y., Aboud, K., Parvathaneni, P., Bao, S., Bermudez, C., Resnick, S.M., Cutting, L.E., Landman, B.A., 2019. 3d whole brain segmentation using spatially localized atlas network tiles. *NeuroImage* 194, 105–119.
- Irving, B., Franklin, J.M., Papież, B.W., Anderson, E.M., Sharma, R.A., Gleeson, F.V., Brady, M., Schnabel, J.A., 2016. Pieces-of-parts for supervoxel segmentation with global context: Application to dce-mri tumour delineation. *Medical image analysis* 32, 69–83.
- Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H., 2021. nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods* 18, 203–211.
- Jing, L., Tian, Y., 2020. Self-supervised visual feature learning with deep neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Kavur, A.E., Gezer, N.S., Barış, M., Aslan, S., Conze, P.H., Groza, V., Pham, D.D., Chatterjee, S., Ernst, P., Özkan, S., Baydar, B., Lachinov, D., Han, S., Pauli, J., Isensee, F., Perkonig, M., Sathish, R., Rajan, R., Sheet, D., Dvletov, G., Speck, O., Nürnberger, A., Maier-Hein, K.H., Bozdağı Akar, G., Ünal, G., Dicle, O., Selver, M.A., 2021. CHAOS Challenge - combined (CT-MR) healthy abdominal organ segmentation. *Medical Image Analysis* 69, 101950. URL: <http://www.sciencedirect.com/science/article/pii/S1361841520303145>, doi:<https://doi.org/10.1016/j.media.2020.101950>.
- Kavur, A.E., Gezer, N.S., Barış, M., Şahin, Y., Özkan, S., Baydar, B., Yüksel, U., Kilikçier, Ç., Olut, Ş., Bozdağı Akar, G., Ünal, G., Dicle, O., Selver, M.A., 2020. Comparison of semi-automatic and deep learning based automatic methods for liver segmentation in living liver transplant donors. *Diagnostic and Interventional Radiology* 26, 11–21. URL: <https://doi.org/10.5152/dir.2019.19025>, doi:10.5152/dir.2019.19.
- Kavur, A.E., Selver, M.A., Dicle, O., Barış, M., Gezer, N.S., 2019. CHAOS - Combined (CT-MR) Healthy Abdominal Organ Segmentation Challenge Data. URL: <https://doi.org/10.5281/zenodo.3362844>, doi:10.5281/zenodo.3362844.
- Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al., 2017. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*.
- Komodakis, N., Gidaris, S., 2018. Unsupervised representation learning by predicting image rotations, in: International Conference on Learning Representations (ICLR).
- Larsson, G., Maire, M., Shakhnarovich, G., 2016. Learning representations for automatic colorization, in: European conference on computer vision, Springer, pp. 577–593.
- Li, G., Jampani, V., Sevilla-Lara, L., Sun, D., Kim, J., Kim, J., 2021a. Adaptive prototype learning and allocation for few-shot segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition.
- Li, X., Chen, H., Qi, X., Dou, Q., Fu, C.W., Heng, P.A., 2018. H-denseunet: hybrid densely connected unet for liver and tumor segmentation from ct volumes. *IEEE transactions on medical imaging* 37, 2663–2674.
- Li, Z., Zhao, W., Shi, F., Qi, L., Xie, X., Wei, Y., Ding, Z., Gao, Y., Wu, S., Liu, J., et al., 2021b. A novel multiple instance learning framework for covid-19 severity assessment via data augmentation and self-supervised learning. *Medical Image Analysis* 69, 101978.
- Liu, W., Zhang, C., Lin, G., Liu, F., 2020a. Crnet: Cross-reference networks for few-shot segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4165–4173.
- Liu, Y., Zhang, X., Zhang, S., He, X., 2020b. Part-aware prototype network for few-shot semantic segmentation, in: European Conference on Computer Vision, Springer, pp. 142–158.
- Lu, Q., Li, Y., Ye, C., 2021. Volumetric white matter tract segmentation with nested self-supervised learning using sequential pretext tasks. *Medical Image Analysis*, 102094.
- Mishra, N., Rohaninejad, M., Chen, X., Abbeel, P., 2018. A simple neural attentive meta-learner, in: International Conference on Learning Representations (ICLR).
- Misra, I., Maaten, L.v.d., 2020. Self-supervised learning of pretext-invariant representations, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6707–6717.
- Nguyen, V.N., Løkse, S., Wickstrøm, K., Kampffmeyer, M., Roverso, D., Jenssen, R., 2020. Sen: A novel feature normalization dissimilarity measure for prototypical few-shot learning networks, Springer, pp. 118–134.
- Noroosi, M., Favaro, P., 2016. Unsupervised learning of visual representations by solving jigsaw puzzles, in: European conference on computer vision, Springer, pp. 69–84.
- Oreshkin, B.N., Rodriguez, P., Lacoste, A., 2018. Tadam: task dependent adaptive metric for improved few-shot learning, in: Proceedings of the 32nd International Conference on Neural Information Processing Systems, pp. 719–729.
- Ouyang, C., Biffi, C., Chen, C., Kart, T., Qiu, H., Rueckert, D., 2020. Self-supervision with superpixels: Training few-shot medical image segmentation without annotation, in: European Conference on Computer Vision, Springer, pp. 762–780.
- Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., Efros, A.A., 2016. Context encoders: Feature learning by inpainting, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2536–2544.
- Rakelly, K., Shelhamer, E., Darrell, T., Efros, A., Levine, S., 2018. Conditional networks for few-shot semantic segmentation, in: International Conference on Learning Representations (ICLR).
- Ravi, S., Larochelle, H., 2017. Optimization as a model for few-shot learning, in: International Conference on Learning Representations (ICLR).
- Ren, X., Malik, J., 2003. Learning a classification model for segmentation, in: Computer Vision, IEEE International Conference on, IEEE Computer Society, pp. 10–10.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical image computing and computer-assisted intervention, Springer, pp. 234–241.
- Roy, A.G., Siddiqui, S., Pölsterl, S., Navab, N., Wachinger, C., 2020. “squeeze & excite” guided few-shot segmentation of volumetric images. *Medical image analysis* 59, 101587.
- Ruff, L., Kauffmann, J.R., Vandermeulen, R.A., Montavon, G., Samek, W., Kloft, M., Dietterich, T.G., Müller, K.R., 2021. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*.
- Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., Lillicrap, T., 2016. Meta-learning with memory-augmented neural networks, in: International conference on machine learning, pp. 1842–1850.
- Shaban, A., Bansal, S., Liu, Z., Essa, I., Boots, B., 2017. One-shot learning for semantic segmentation, in: Proceedings of the British Machine Vision Conference (BMVC), pp. 167.1–167.13.
- Snell, J., Swersky, K., Zemel, R., 2017. Prototypical networks for few-shot learning, in: Advances in neural information processing systems, pp. 4077–4087.
- Stutz, D., Hermans, A., Leibe, B., 2018. Superpixels: An evaluation of the state-of-the-art. *Computer Vision and Image Understanding* 166, 1–27.
- Tian, Z., Zhao, H., Shu, M., Yang, Z., Li, R., Jia, J., 2020. Prior guided feature enrichment network for few-shot segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Tsochatzidis, L., Koutla, P., Costaridou, L., Pratikakis, I., 2021. Integrating segmentation information into cnn for breast cancer diagnosis of mammographic masses. *Computer Methods and Programs in Biomedicine* 200, 105913.
- Vesal, S., Ravikumar, N., Maier, A., 2019. Automated multi-sequence cardiac mri segmentation using supervised domain adaptation, in: International workshop on statistical atlases and computational models of the heart, Springer, pp. 300–308.
- Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al., 2016. Matching networks for one shot learning, in: Advances in neural information processing systems, pp. 3630–3638.
- Wang, K., Liew, J.H., Zou, Y., Zhou, D., Feng, J., 2019. Panet: Few-shot image semantic segmentation with prototype alignment, in: Proceedings of the IEEE International Conference on Computer Vision, pp. 9197–9206.
- Wilcoxon, F., 1992. Individual comparisons by ranking methods, in: Breakthroughs in statistics. Springer, pp. 196–202.

- Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K., 2017. Aggregated residual transformations for deep neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1492–1500.
- Zhang, B., Xiao, J., Qin, T., 2021. Self-guided and cross-guided learning for few-shot segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition.
- Zhang, C., Lin, G., Liu, F., Yao, R., Shen, C., 2019. Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5217–5226.
- Zhang, R., Isola, P., Efros, A.A., 2016. Colorful image colorization, in: European conference on computer vision, Springer. pp. 649–666.
- Zhang, X., Wei, Y., Yang, Y., Huang, T.S., 2020. Sg-one: Similarity guidance network for one-shot semantic segmentation. *IEEE Transactions on Cybernetics* 50, 3855–3865.
- Zhu, J., Li, Y., Hu, Y., Ma, K., Zhou, S.K., Zheng, Y., 2020. Rubik’s cube+: A self-supervised feature learning framework for 3d medical image analysis. *Medical Image Analysis* 64, 101746.
- Zhuang, X., 2016. Multivariate mixture model for cardiac segmentation from multi-sequence mri, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 581–588.
- Zhuang, X., 2018. Multivariate mixture model for myocardial segmentation combining multi-source images. *IEEE transactions on pattern analysis and machine intelligence* 41, 2933–2946.