

Segmentation with mixed supervision: Confidence maximization helps knowledge distillation

Bingyuan Liu^{a,*}, Christian Desrosiers^a, Ismail Ben Ayed^{a,b}, Jose Dolz^{a,b,**}

^aÉTS Montréal

^bCentre de recherche du Centre hospitalier de l'Université de Montréal (CRCHUM)

Abstract

Despite achieving promising results in a breadth of medical image segmentation tasks, deep neural networks (DNNs) require large training datasets with pixel-wise annotations. Obtaining these curated datasets is a cumbersome process which limits the applicability of DNNs in scenarios where annotated images are scarce. Mixed supervision is an appealing alternative for mitigating this obstacle. In this setting, only a small fraction of the data contains complete pixel-wise annotations and other images have a weaker form of supervision, e.g., only a handful of pixels are labeled. In this work, we propose a dual-branch architecture, where the upper branch (teacher) receives strong annotations, while the bottom one (student) is driven by limited supervision and guided by the upper branch. Combined with a standard cross-entropy loss over the labeled pixels, our novel formulation integrates two important terms: (i) a Shannon entropy loss defined over the less-supervised images, which encourages confident student predictions in the bottom branch; and (ii) a Kullback-Leibler (KL) divergence term, which transfers the knowledge (i.e., predictions) of the strongly supervised branch to the less-supervised branch and guides the entropy (student-confidence) term to avoid trivial solutions. We show that the synergy between the entropy and KL divergence yields substantial improvements in performance. We also discuss an interesting link between Shannon-entropy minimization and standard pseudo-mask generation, and argue that the former should be preferred over the latter for leveraging information from unlabeled pixels. We evaluate the effectiveness of the proposed formulation through a series of quantitative and qualitative experiments using two publicly available datasets. Results demonstrate that our method significantly outperforms other strategies for semantic segmentation within a mixed-supervision framework, as well as recent semi-supervised approaches. Moreover, in line with recent observations in classification, we show that the branch trained with reduced supervision and guided by the top branch largely outperforms the latter. Our code is publicly available: <https://github.com/by-liu/ConfKD>.

Keywords: CNN, image segmentation, mixed-supervision, semi supervision,

1. Introduction

The advent of deep learning has led to the emergence of high-performance models which currently

dominate the medical image segmentation literature (Litjens et al., 2017; Dolz et al., 2018; Ronneberger et al., 2015). The availability of large training datasets with high-quality segmentation ground-truth has been a key factor for these advances. Nevertheless, obtaining such annotations is a cumbersome process prone to observer variability, which is further magnified when volumetric data is involved. To alle-

*Corresponding author: bingyuan.Liu@etsmtl.ca

**Corresponding author: jose.dolz@etsmtl.ca

viate the need for large labeled datasets, weakly supervised learning has recently emerged as an appealing alternative. In this scenario, one has access to a large amount of weakly labeled data that can come in the form of bounding boxes (Kervadec et al., 2020; Rajchl et al., 2016), scribbles (Lin et al., 2016), image tags (Lee et al., 2019) or anatomical priors (Kervadec et al., 2019b; Peng et al., 2020b). However, even though numerous attempts have been made to train segmentation models from weak supervision, most of them still fall behind their supervised counterparts, limiting their applicability in real-world settings.

Another promising learning scenario is mixed supervision, where only a small fraction of data is densely annotated and a larger dataset contains less-supervised images. In this setting, which helps keeping the annotation budget under control, strongly-labeled data – where all pixels are annotated – can be combined with images presenting weaker forms of supervision. Prior literature (Lee et al., 2019; Rajchl et al., 2016) has focused mainly on leveraging weak annotations to generate accurate initial pixel-wise annotations, or *pseudo-masks*, which are then combined with strong types of supervision to augment the training dataset. The resulting dataset is employed to train a segmentation network, mimicking fully supervised training. Nevertheless, we argue that treating both equally in a single branch may result in limited improvements, as the less-supervised data is underused. Other approaches resort to multi-task learning (Mlynarski et al., 2019; Shah et al., 2018; Wang et al., 2019), where the mainstream task (i.e., segmentation) is assisted by auxiliary objectives that are typically integrated in the form of localization or classification losses. While multi-task learning might enhance the common representation for both tasks in the feature space, this strategy has some drawbacks. First, the learning of relevant features is driven by commonalities between the multiple tasks, which may generate suboptimal representations for the mainstream task. Secondly, having distinct task-objectives ignores the direct interaction between the multi-stream outputs, for example, explicitly enforcing consistency between the predictions of multiple branches. As we show in our experiments, consider-

ing such interaction significantly improves the results.

Motivated by these observations, we propose a novel formulation for learning with mixed supervision in medical image segmentation. Particularly, our dual-branch network imposes a separate processing of the strong and weak annotations, which prevents direct interference of different supervision cues. This is supported by the recent findings in (Luo and Yang, 2020), who demonstrated empirically that bundling different forms of supervision together to train a segmentation network is problematic, and argued that joint treatment of different supervision under-exploits less supervised samples, introducing limited improvement. In particular, authors pointed out two key issues related to equal treatment of different levels of supervision: *sample imbalance* (commonly much more less-supervised samples than fully supervised ones) and *supervision inconsistency* (less-supervised samples provide lower quality annotations). The former introduces a high risk of overfitting towards less supervised data, whereas the later induces inconsistencies in the supervisory signals. Authors validated these hypothesis in their experiments, and showed that by decoupling branches receiving different types of labels the supervision inconsistency and biases from class imbalance can be eliminated.

As the uncertainty of predictions for unlabeled pixels can be high, the proposed model includes a loss term based Shannon entropy that enforces high-confidence predictions over the whole image. Moreover, in contrast to prior works in mixed-supervision (Mlynarski et al., 2019; Shah et al., 2018; Wang et al., 2019), which have overlooked the co-operation between multiple branches by considering independent multi-task objectives, we introduce a Kullback-Leibler (KL) divergence term. The benefits of the latter are two-fold. First, it transfers the predictions generated by the strongly supervised branch (teacher) to the less-supervised branch (student). Second, it guides the entropy (student-confidence) term to avoid trivial solutions. Interestingly, we show that the synergy between the entropy and KL term yields substantial improvements in performances. Furthermore, we discuss an interesting link between Shannon-entropy minimization and pseudo-

mask generation, and argue that the former should be preferred over the latter for leveraging information from unlabeled pixels. We report comprehensive experiments and comparisons with other strategies for learning with mixed supervision, which show the effectiveness of our novel formulation. An interesting finding is that the branch receiving weaker supervision considerably outperforms the strongly supervised branch. This phenomenon, where **the student surpasses the teacher’s performance**, is in line with recent observations in the context of image classification (Furlanello et al., 2018; Yim et al., 2017).

A preliminary conference version of this work has appeared at IPMI’21 (Dolz et al., 2021). Nevertheless, this journal version provides a substantial extension. First, we further discuss the current literature in semi-supervised segmentation, which is closely related to the proposed methodology. Furthermore, we have performed several additional experiments to demonstrate the robustness and usability of our approach. In particular, new main experiments include: 1) benchmark against well-known and recent semi-supervised segmentation methods, 2) evaluation of our model in the publicly available Left Atrium (LA) segmentation challenge, 3) assessing the impact of several components in the methodology and 4) studying the impact of alternative divergence functionals as consistency terms in our formulation. In addition to the theoretical insights regarding the preference of directly minimizing the entropy of the predictions over using pseudo-labels given in the conference version, we provide empirical evidence that employing pseudo-labels has indeed a strong pushing effect on uncertain predictions at the beginning of the training.

2. Related work

Mixed-supervised segmentation An appealing alternative to training CNNs with large labeled datasets is to combine a reduced number of fully-labeled images with a larger set of images with reduced annotations. These annotations can come in the form of bounding boxes, scribbles or image tags,

for example¹. A large body of the literature in this learning paradigm addresses the problem from a multi-task objective perspective (Hong et al., 2015; Bhalgat et al., 2018; Shah et al., 2018; Mlynarski et al., 2019; Wang et al., 2019), which might hinder their capabilities to fully leverage joint information for the mainstream objective. Furthermore, these methods typically require carefully-designed task-specific architectures, which also integrate task-dependent auxiliary losses, limiting the applicability to a wider range of annotations. For example, the architecture designed in (Shah et al., 2018) requires, among others, landmark annotations, which might be difficult to obtain in many applications. More recently, Luo et al. (Luo and Yang, 2020) promoted the use of a dual-branch architecture to deal separately with strongly and weakly labeled data. Particularly, while the strongly supervised branch is governed by available fully annotated masks, the weakly supervised branch receives supervision from a proxy ground-truth generator, which requires some extra information, such as class labels. While we advocate the use of independent branches to process naturally different kinds of supervision, we believe that this alone is insufficient, and may lead to suboptimal results. Thus, our work differs from (Luo and Yang, 2020) in several aspects. First, we make a better use of the labeled images by enforcing consistent segmentations between the strongly and weakly supervised branches on these images. Furthermore, we enforce confident predictions at the weakly supervised branch by minimizing the Shannon entropy of the softmax predictions.

Semi-supervised segmentation in medical images Semi-supervised learning is closely related to the proposed methodology. In this scenario, a small number of labeled images are leveraged with a much larger set of unlabeled images. In recent years, a breadth of semi-supervised approaches have been proposed for medical image segmentation, includ-

¹Note that this type of supervision differs from semi-supervised methods, which leverage a small set of labeled images and a much larger set of unlabeled images.

ing techniques based on adversarial learning (Zhang et al., 2017), self-training (Bai et al., 2017), manifold learning (Baur et al., 2017), co-training (Peng et al., 2020a; Zhou et al., 2019), temporal ensembling (Perone and Cohen-Adad, 2018), data augmentation (Chaitanya et al., 2019), consistency regularization (Bortsova et al., 2019) and mutual information maximization (Peng et al., 2021). The common element of these approaches is adding an unsupervised loss computed on unlabeled images, which regularizes the learning. In contrast, our model exploits images that can be fully or partly annotated, processing each type in a separate branch of the proposed network.

Distilling knowledge in semantic segmentation Transferring knowledge from one model to another has recently gained attention in segmentation tasks. For example, the teacher-student strategy has been employed in model compression (Bar et al., 2019), to distil knowledge from multi-modal to mono-modal segmentation networks (Hu et al., 2020), or in domain adaptation (Xu et al., 2019). Semi-supervised segmentation has also benefited from teacher-student architectures (Cui et al., 2019; Sedai et al., 2019). In these approaches, however, the segmentation loss evaluating the consistency between the teacher and student models is computed on the unannotated data. A common practice, for example, is to add additive Gaussian noise to the unlabeled images, and enforce similar predictions for the original and noised images. This contrasts with our method, which enforces consistency only on the strongly labeled data, thereby requiring less additional images to close the gap with full supervision.

3. Methodology

We first define the set of training images as $\mathcal{D} = \{(\mathbf{X}_n, \mathbf{Y}_n)\}_n$, where $\mathbf{X}_i \in \mathbb{R}^{\Omega_i}$ represents the i^{th} image and $\mathbf{Y}_i \in \{0, 1\}^{\Omega_i \times C}$ its corresponding ground-truth segmentation mask. Ω_i denotes the spatial image domain and C the number of segmentation classes (or regions). We assume the dataset has two subsets: $\mathcal{D}_s =$

$\{(\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_m, \mathbf{Y}_m)\}$, which contains complete pixel-level annotations of the associated C categories, and $\mathcal{D}_w = \{(\mathbf{X}_{m+1}, \mathbf{Y}_{m+1}), \dots, (\mathbf{X}_n, \mathbf{Y}_n)\}$, whose labels can take the form of semi- or weakly-supervised annotations (e.g., scribbles, points, bounding boxes or image-tags). Furthermore, for each image \mathbf{X}_i in $\mathcal{D} = \mathcal{D}_s \cup \mathcal{D}_w$, $\mathbf{P}_i \in [0, 1]^{\Omega_i \times C}$ denotes the softmax probability outputs of the network, i.e., the matrix containing a simplex column vector $\mathbf{p}_i^l = (p_i^{l,1}, \dots, p_i^{l,C})^T \in [0, 1]^C$ for each pixel $l \in \Omega_i$. Note that we omit the parameters of the network here to simplify notation.

3.1. Multi-branch architecture

The proposed architecture is composed of multiple branches, each dedicated to a specific type of supervision (see Fig. 1). It can be divided in two components: a shared feature extractor and independent but identical decoding networks (one per type of supervision), which differ in the type of annotations received. It is worth mentioning that the initialisation of the weights in the decoders and the different gradients received by each branch ensure that the parameters from both decoders will not have the same values during training. Even though the proposed multi-branch architecture has similarities with the recent work in (Luo and Yang, 2020), there are significant differences, particularly in the loss functions, which leads to different optimization scenarios.

3.2. Supervised learning

The top-branch is trained under the fully-supervised paradigm, where a set of training images containing pixel-level annotations for all the pixels is given, i.e., \mathcal{D}_s . The problem amounts to minimizing with respect to the network parameters a standard full-supervision loss, which typically takes the form of a cross-entropy:

$$\mathcal{L}_s = - \sum_{i=1}^m \sum_{l \in \Omega_i} (\mathbf{y}_i^l)^T \log (\mathbf{p}_i^l)_{\text{top}} \quad (1)$$

where vector $\mathbf{y}_i^l = (y_i^{l,1}, \dots, y_i^{l,C}) \in \{0, 1\}^C$ describes the ground-truth annotation for pixel $l \in \Omega_i$.

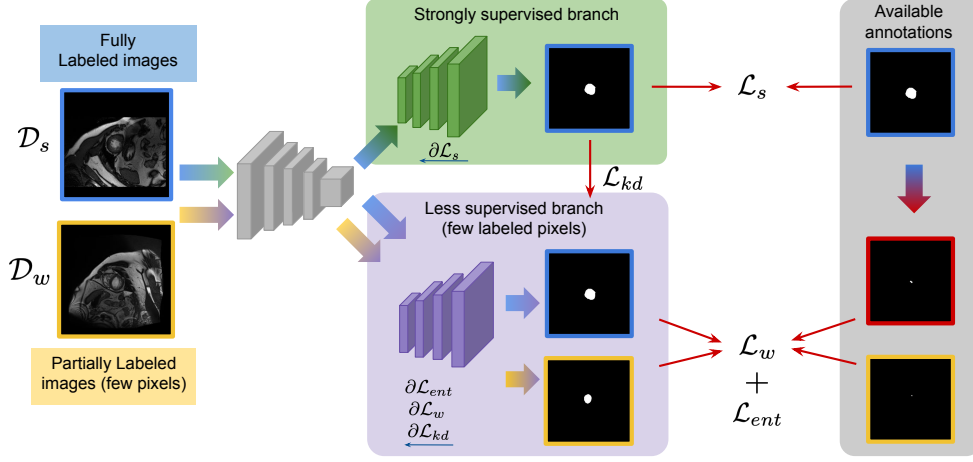


Figure 1: Overview of the proposed method. **Training:** Both fully and partial labeled images are fed to the network. The top branch generates predictions for fully labeled images, whereas the bottom branch generates the outputs for partially labeled images. Furthermore, the bottom branch also generates segmentations for the fully labeled images, which are guided by the KL term between the two branches. **Inference:** Once the model is trained, we can remove the strongly supervised branch (*Top*), and get the final segmentation result from the bottom stream. The gradients for each loss term are highlighted in the figure. Note that, similarly to UNet-like architectures, the proposed model has skip connections between the encoder and the decoders.

Here, notation $(\cdot)_{\text{top}}$ refers to the softmax outputs of the *top* branch of the network. Note that all the losses are normalized by the cardinality of the training dataset, which has been omitted to simplify the notation.

3.3. Not so-supervised branch

We consider the scenario where only the labels for a handful of pixels are known, i.e., scribbles or points. Particularly, we use the dataset \mathcal{D}_w whose pixel-level labels are partially provided. Furthermore, for each image on the labeled training set, \mathcal{D}_s , we generate partially supervised labels (more details in the experiments’ section), which are added to augment the dataset \mathcal{D}_w . Then, for the partially-labeled set of pixels, denoted as $\Omega_i^{\text{partial}}$ for each image $i \in \{1, \dots, n\}$, we can resort to the following partial-supervision loss, which takes the form of a cross-entropy on the fraction of labeled pixels:

$$\mathcal{L}_w = - \sum_{i=m+1}^n \sum_{l \in \Omega_i^{\text{partial}}} (\mathbf{y}_i^l)^T \log (\mathbf{p}_i^l)_{\text{bottom}} \quad (2)$$

where notation $(\cdot)_{\text{bottom}}$ refers to the softmax outputs of the *bottom* branch of the network.

3.4. Distilling strong knowledge

In addition to the specific supervision available at each branch, we transfer the knowledge from the teacher (top branch) to the student (bottom branch). This is done by forcing the softmax distributions from the bottom branch to mimic the probability predictions generated by the top branch for the fully labeled images in \mathcal{D}_s . This knowledge-distillation regularizer takes the form of a Kullback-Leibler divergence (\mathcal{D}_{KL}) between both distributions:

$$\mathcal{L}_{kd} = \sum_{i=1}^m \sum_{l \in \Omega_i} \mathcal{D}_{\text{KL}} \left((\mathbf{p}_i^l)_{\text{top}} \parallel (\mathbf{p}_i^l)_{\text{bottom}} \right) \quad (3)$$

where $\mathcal{D}_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) = \mathbf{p}^T \log \frac{\mathbf{p}}{\mathbf{q}}$, with T denoting the transpose operator.

3.5. Shannon-Entropy minimization

Finally, we encourage high confidence in the student softmax predictions for the partially labeled images by minimizing the Shannon entropy of the predictions on the bottom branch:

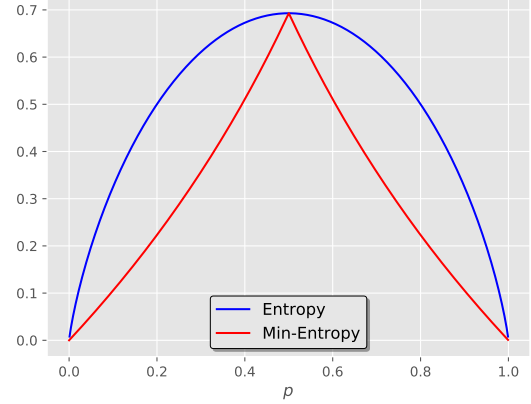
$$\mathcal{L}_{ent} = \sum_{i=m+1}^n \sum_{l \in \Omega_i} \mathcal{H}(\mathbf{p}_i^l) \quad (4)$$

where $\mathcal{H}(\mathbf{p}) = -\mathbf{p}^T \log \mathbf{p}$ is the Shannon entropy of distribution \mathbf{p} .

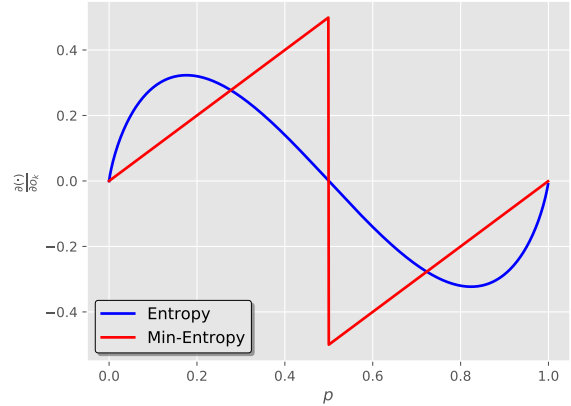
Entropy minimization is widely used in semi-supervised learning (SSL) and transductive classification (Grandvalet and Bengio, 2005; Berthelot et al., 2019; Dhillon et al., 2019; Boudiaf et al., 2020) to encourage confident predictions at unlabeled data points. Fig. 2a plots the entropy in the case of a two-class distribution $(p, 1 - p)$, showing how the minimum is reached at the vertices of the simplex, i.e., when $p = 0$ or $p = 1$. However, surprisingly, in segmentation, entropy is not commonly used, except a few recent works in the different contexts of SSL and domain adaptation (Peng et al., 2020c; Bateson et al., 2020; Vu et al., 2019). As we will see in our experiments, we found that the synergy between the entropy term for confident students, \mathcal{L}_{ent} , and the student-teacher knowledge transfer term, \mathcal{L}_{kd} , yields substantial increases in performances. Furthermore, in the following, we discuss an interesting link between *pseudo-mask generation*, common in the segmentation literature, and entropy minimization, showing that the former could be viewed as a proxy for minimizing the latter. We further provide insights as to why entropy minimization should be preferred for leveraging information from the set of unlabeled pixels.

3.6. Link between entropy and pseudo-mask supervision

In the weakly- and semi-supervised segmentation literature, a very dominant technique to leverage information from unlabeled pixels is to generate pseudo-masks and use these as supervision in a cross-entropy training, in an alternating way (Lin et al.,



(a) Shannon entropy (blue) and min-entropy (red) for a two-class distribution $(p, 1 - p)$, with $p \in [0, 1]$.



(b) Derivatives of both the entropy and Min-Entropy with respect to the input logit of class k .

2016; Khoreva et al., 2017; Papandreou et al., 2015). This self-supervision principle is also well known in classification (Lee, 2013). Given pixel-wise predictions $\mathbf{p}_i^l = (p_i^{l,1}, \dots, p_i^{l,C})$, pseudo-masks $q_i^{l,k}$ are generated as follows: $q_i^{l,k} = 1$ if $p_i^{l,k} = \max_c p_i^{l,c}$ and 0 otherwise. By plugging these pseudo-labels in a cross-entropy loss, it is easy to see that this corresponds to minimizing the *min-entropy*, $\mathcal{H}_{\min}(\mathbf{p}_i^l) = -\log(\max_c p_i^{l,c})$, which is a lower bound on the Shannon entropy; see the red curve in Figure 2a. Fig.

2a and 2b provide a good insight as to why entropy should be preferred over min-entropy (pseudo-masks) as a training loss for unlabeled data points, and our experiments confirm this. With entropy, the gradients of low-confidence predictions at the middle of the simplex are small and, therefore, dominated by the other terms at the beginning of training. However, with min-entropy, the inaccuracies resulting from uncertain predictions are reinforced (pushed towards the simplex vertices), yielding early/unrecoverable errors in the predictions, which might mislead training. This is a well-known limitation of self-supervision in the SSL literature (Chapelle et al., 2009).

3.7. Joint objective

Our final loss function takes the following form:

$$\mathcal{L}_t = \mathcal{L}_s + \lambda_w \mathcal{L}_w + \lambda_{kd} \mathcal{L}_{kd} + \lambda_{ent} \mathcal{L}_{ent} \quad (5)$$

where λ_w , λ_{kd} and λ_{ent} balance the importance of each term.

4. Experimental setting

Benchmark dataset To evaluate the proposed model, we employ two public segmentation benchmarks. First, we focus on the task of left ventricular (LV) endocardium segmentation on cine MRI images. Particularly, we used the training set from the publicly available data of the 2017 ACDC Challenge (Bernard et al., 2018), which consists of 100 cine magnetic resonance (MR) exams covering several well defined pathologies: dilated cardiomyopathy, hypertrophic cardiomyopathy, myocardial infarction with altered left ventricular ejection fraction and abnormal right ventricle. Following prior literature, e.g., (Peng et al., 2020a; Wang et al., 2021), slices contained in each 3D-MRI scan were considered as 2D images, whose spatial resolution was 256×256 . We split this dataset into 80 exams for training, 5 for validation and the remaining 15 for testing. Then, to demonstrate the broad applicability of our method,

we use images from the Left Atrium (LA) segmentation challenge², which has been widely used in the context of semi-supervised segmentation. It includes 100 gadolinium-enhanced magnetic resonance imaging (GE-MRI) scans, with aligned LA segmentation masks. These 3D GE-MRI scans have isotropic resolution of $0.625 \times 0.625 \times 0.625 \text{mm}^3$. Following the setting of (Yu et al., 2019), all the scans are cropped centering at the heart region for better comparison of the segmentation performance and normalized as zero mean and unit variance. Contrary to the ACDC dataset, the inputs of the network in this setting was 3D volumetric data. In our experiments, we randomly divided them into 80 for training, 5 for validation and the remaining 15 for testing.

Generating partially labeled images The training exams are divided into a small set of fully labeled images, \mathcal{D}_s , and a larger set of images with reduced supervision, \mathcal{D}_w , where only a handful of pixels are labeled. Concretely, we employ the same partial labels as in (Kervadec et al., 2019a,b), which are obtained by eroding iteratively the full pixel-wise masks with a kernel of size 10×10 , until the smallest possible contour is obtained. To evaluate how increasing the amount of both fully and partially labeled affects the performance, we evaluated the proposed models in three settings, referred to as *Set-3*, *Set-5*, and *Set-10*. In these settings, the number of fully labeled images is 3, 5 and 10, respectively, while the number of images with partial labels is $\times 5$ times the number of labeled images.

Evaluation metrics For evaluation purposes we employ two well-known metrics in medical image segmentation: the Dice similarity score (DSC) and the modified Hausdorff-Distance (MHD). Particularly, the MHD represents the 95th percentile of the symmetric HD between the binary objects in two images.

Baseline methods To demonstrate the efficiency of the proposed model, we compared it to several baselines. First, we include full-supervised baselines that

²<http://atriaseg2018.cardiacatlas.org/>

Table 1: Results on ACDC (Left-ventricle) on the testing set for the *top* and *bottom* branches (when applicable). Results are averaged over three runs. Best results for non-fully supervised methods are highlighted in bold.

					Top		Bottom		Ensemble	
Setting		Model	FS	PS	DSC	HD-95	DSC	HD-95	DSC	HD-95
Set-3	Single Branch	Lower bound	✓	–	54.66	80.05	–	–	–	–
		Single	✓	✓	57.42	78.80	–	–	–	–
		Single + Ent	✓	✓	43.01	83.98	–	–	–	–
		Upper Bound (Set-3)	✓	–	87.17	5.34	–	–	–	–
	Dual Branch	Decoupled (Luo and Yang, 2020)	✓	✓	56.61	74.95	5.01	120.06	–	–
		Ours (KL)	✓	✓	68.25	63.15	71.49	63.51	–	–
Ours (KL+Ent)		✓	✓	78.38	46.73	86.94	8.84	86.35	11.97	
Set-5	Single Branch	Lower bound	✓	–	69.71	51.75	–	–	–	–
		Single	✓	✓	70.73	51.34	–	–	–	–
		Single + Ent	✓	✓	74.92	55.24	–	–	–	–
		Upper Bound (Set-5)	✓	–	87.69	4.93	–	–	–	–
	Dual Branch	Decoupled (Luo and Yang, 2020)	✓	✓	70.96	54.42	4.29	127.68	–	–
		Ours (KL)	✓	✓	80.64	23.25	79.06	34.83	–	–
Ours (KL+Ent)		✓	✓	85.57	20.68	88.77	5.40	88.54	5.49	
Set-10	Single Branch	Lower bound	✓	–	78.28	44.16	–	–	–	–
		Single	✓	✓	78.17	42.99	–	–	–	–
		Single + Ent	✓	✓	80.63	37.75	–	–	–	–
		Upper Bound (Set-10)	✓	–	91.18	3.71	–	–	–	–
	Dual Branch	Decoupled (Luo and Yang, 2020)	✓	✓	77.53	32.23	4.58	125.36	–	–
		Ours (KL)	✓	✓	88.29	12.47	88.68	11.93	–	–
Ours (KL+Ent)		✓	✓	86.53	5.64	90.92	1.39	90.75	1.55	
All images	Single Branch	Upper bound	✓	–	93.31	3.46	–	–		

FS and PS indicate fully or partially supervised images.

will act as lower and upper bounds. The lower bound employs only a small set of fully labeled images (either 3, 5 or 10, depending on the setting), whereas the upper bound considers all the available training images. The fully labeled images used to train the lower bound baselines are exactly the same images employed in the other models. Then, we consider a single-branch network, referred to as *Single*, which receives both fully and partial labeled images without making distinction between them. In order to have a fair comparison with the proposed method, we also include a version of the *Single* model where the entropy of the predictions is minimized during training, i.e., *Single+Ent*. To assess the impact of decoupling the branches without further supervision, similar to (Luo and Yang, 2020), we modify the baseline network by integrating two independent decoders, while the encoder remains the same. This model, which we refer to as *Decoupled*, is governed by different types of supervision at each branch and is equivalent to our model without the proposed KL and entropy min-

imization terms. Then, our first model, which we refer to as *KL*, integrates the KL divergence term presented in Eq. (3), whereas *KL+Ent* corresponds to the whole proposed model, which couples the two important terms in Eq. (3) and Eq. (4) in the formulation. Last, to evaluate the benefits of the proposed method compared to a equivalent fully supervised model when the same amount of images are available, we include the results when the *Single* model is trained with all the images on each setting, and their corresponding pixel-level mask. We refer to this model as to the *Upper Bound (Set-N)*, where *N* indicates the setting.

Implementation details On the LV segmentation task, we employed UNet (Ronneberger et al., 2015) as backbone architecture for the single branch models, whereas VNet (Milletari et al., 2016) was utilized in the LA segmentation task. The reason behind these choices is that segmentation was performed in a 2D manner in ACDC, whereas we employ volumetric in-

puts for the LA dataset, following the literature. Regarding the dual-branch architectures, we modified the decoding path of the standard UNet and VNet to accommodate two separate branches. Note that the encoders remain the same for both single and dual-branch architectures. All the networks on the LV segmentation task are trained during 500 epochs by using Adam optimizer, with a batch size equal to 24 (i.e., 8 labeled and 16 partially labeled images), and the learning rate was initialized to 1×10^{-4} . We empirically set the values of λ_w , λ_{kd} and λ_{ent} to 0.001, 50 and 1, respectively. For the LA segmentation task, we followed the setting in (Yu et al., 2019) by training the network with SGD optimizer and batch size equal to 12 (i.e., 4 labeled and 8 partially labeled images). We found that our formulation provided the best results when the input distributions to the KL term in eq. (3) were very smooth, which was achieved by applying softmax over the softmax predictions. All the hyperparameters, for all the baselines and models, were fixed by using the independent validation set. In particular, we first found the best value for λ_w (0.001) for the *Single* model, which was then fixed to find the other hyperparameters of the proposed model. Then, the optimal λ_{kd} weight of the KL term was found (Detailed results can be found in Supplemental Materials, section Appendix A). Last, with λ_w and λ_{kd} fixed, we found the best value of λ_{ent} , whose value is the same in both the proposed model and the baseline *Single+Ent*. Furthermore, we perform 3 runs for each model and report the average values. The code was implemented in PyTorch and all the experiments were performed in a server equipped with a NVIDIA Titan RTX GPU.

4.1. Results

Main results. Table 1 reports the quantitative evaluation of the proposed method compared to the different baselines on the ACDC dataset.

The first thing we observe is that, across all the settings, simply adding partial annotations to the training set does not considerably improve the segmentation performance, i.e., *Single* model. Furthermore, integrating the entropy minimization term in this baseline results in a performance degradation in

the less-supervised setting (i.e., Set-3), whereas the performance typically increases in the other two settings. Indeed, the lower performance observed in the setting Set-3 might be due to the entropy minimization term pushing towards trivial solutions. This is particularly important on this setting: as the number of fully-labeled images is low the information derived from these images might not be enough to serve as a strong prior to avoid these trivial solutions.

On the other hand, by integrating the guidance from the upper branch, the network is capable of leveraging additional partially-labeled images more efficiently through the bottom branch. Furthermore, if we couple the KL divergence term with an objective based on minimizing the entropy of the predictions on the partial labeled images, the segmentation performance substantially increases. Particularly, the gain obtained by the complete model is consistent across the several settings, improving the DSC by 6-12% compared to the *KL* model, and reducing the MHD by nearly 30%. Compared to the baseline dual-branch model, i.e., *Decoupled*, our approach brings improvements of 10-20% in terms of DSC and reduces the MHD values by 30-40%. Last, the results obtained by the proposed model in each setting are on par with their individual upperbound counterparts (i.e., those trained with the same images as our model), particularly in terms of DSC. Furthermore, as the number of labeled images increases, the gap in the HD metric is decreased between the two models, with our model outperforming the upper bound in the *Set-10* setting.

These results demonstrate the strong capabilities of the proposed model to leverage fully and partially labeled images during training. It is noteworthy to mention that findings on these results, where **the student excels the teacher**, align with recent observations in classification (Furlanello et al., 2018; Yim et al., 2017). Note that the top branch is also improved in our formulation. This can be explained from the fact that even though the supervision the teacher (*top*) receives remains unchanged, changes in the student (*bottom*) also affect the encoder, which is shared among both. Thus, the integration of the proposed objective also results in an improvement on

the latent representation of the model.

The predictions from top and bottom streams can be seen as independent model outputs. In this scenario, ensemble learning has often demonstrated to be an efficient solution to boost the performance of single models (Dolz et al., 2020). Nevertheless, as only two different predictions are available, and there exists a significant gap in performance between both, combining both outputs hampers the performance compared to the bottom branch, particularly in lower data regimes. This supports our choice of using only the weakly supervised stream at inference.

Comparison with proposals. As mentioned previously, a popular paradigm in weakly and semi-supervised segmentation is to resort to pseudo-masks generated by a trained model, which are used to re-train the network mimicking full supervision. To demonstrate that our model leverages more efficiently the available data, we train a network with an augmented dataset composed by the available labeled images and the proposals generated on the unlabeled images by the *Lower bound* and *KL* models, whose results are reported in Table 2. We can observe that despite typically improving the base model, minimizing the cross-entropy over proposals does not outperform directly minimizing the entropy on the predictions of the partially labeled images.

Table 2: Results obtained by training on an augmented dataset composed by fully labeled images and proposals generated from the *Lower bound* and *KL* models on the test set. Results on the ACDC dataset.

Setting	Proposals (Lower bound)		Proposals (KL)		Ours (KL+Ent)	
	DSC	HD-95	DSC	HD-95	DSC	HD-95
Set-3	63.11	49.99	74.27	45.44	86.94	8.84
Set-5	73.91	45.54	81.35	20.28	88.77	5.40
Set-10	81.31	29.95	89.26	7.98	90.92	1.39

In addition to the quantitative results reported below, we depict the performance evolution on the validation set, for the setting Set-3, in Fig. 3. An interesting observation is that, while training with pseudo-labels converges faster, our method needs more iterations to reach convergence which is achieved at a slower pace. Nevertheless, minimizing the entropy

on the predictions results in the model outperforming the pseudo-label based approach. We argue that this behaviour can be explained from a gradient dynamics perspective. At the beginning of the training, there exist low-confidence predictions which might result in inaccurate predictions (e.g., in Fig 2a we can observe that these predictions lie within the middle of the simplex). As we showed in Section 3.6, employing pseudo-masks in a cross-entropy loss corresponds to minimizing the *min-entropy*, which quickly pushes low-confidence predictions towards the simplex vertices at the beginning of the training. On the other hand, if we employ an entropy term, the gradients of low-confidence predictions in the same region (i.e., middle of the simplex) are small compared to the other terms at the beginning. However, as the other terms start to be satisfied, the scale of their gradients becomes comparable to the entropy gradients term, and hence this term begins its regularization role.

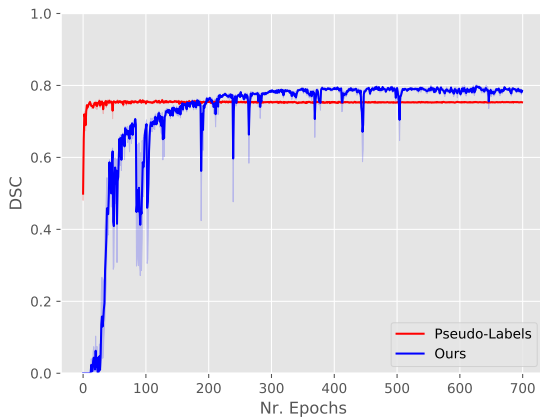


Figure 3: DSC evolution across the training when trained with pseudo-labels (*red*) and the proposed formulation (*blue*) on 3 subjects fully labeled on the ACDC validation dataset.

Several failure cases from proposals are depicted in Fig 4. With this strategy, these errors are propagated during training, which might explain the low performance compared to our method. In addition to lower performances, this approach requires to fully train a first model, generate pseudo-labels and then re-train a second model with the generated masks.

This contrasts to our method, which is trained in an end-to-end manner.

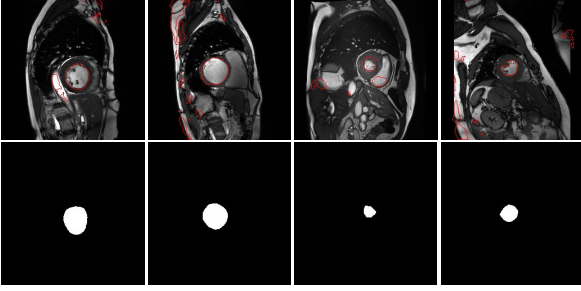


Figure 4: Failure cases which are employed as pseudo-labels in the *proposals-based* approach (top), whose errors are reinforced during training, and their corresponding ground truth. Best viewed in colours.

Are errors actually propagated through the training? To illustrate the weaknesses of pseudo labels based approaches, particularly in reinforcing the errors, we perform the following experiment. We start with an initial model, i.e., *Ours (KL)*, generates the pseudo labels to train the model at iteration \mathcal{I}_1 , i.e., *Proposals (KL)*. Similarly, once the model at iteration \mathcal{I}_t is trained, it is used to generate the pseudo-labels for training the next model at iteration \mathcal{I}_{t+1} . Furthermore, at each iteration, the model parameters are initialized randomly (note that this is similar to the so-called self-training strategy). We see in Figure 5 that, despite having a performance improvement in early iterations, this improvement quickly saturates. This degradation of results over time suggests that the model is unable to correct noisy pseudo-labels and accumulates these errors across iterations. As explained recently in (Huo et al., 2021), this *trapping* effect can be explained from an optimization perspective. If we use \mathbf{X}_i to denote a training image and \mathbf{Y}_i its corresponding ground truth, we can assume that the predicted segmentation is $\hat{\mathbf{Y}}_i = \mathbf{Y}_i + \epsilon_i$, where ϵ_i denotes the prediction error. If \mathbf{X}_i belongs to the fully labeled dataset \mathcal{D}_s , \mathbf{Y}_i is known. Thus, as the optimization objective involves minimizing $|\epsilon_i|$ on \mathcal{D}_s , it follows a zero mean distribution. In contrast, \mathbf{Y}_i is unknown on the weakly supervised dataset \mathcal{D}_w , which allows ϵ_i to follow a distribution with non-zero

mean. In this scenario, the prediction $\hat{\mathbf{Y}}_i$ (used later as pseudo-label) might integrate this noise, which can be propagated to the model in subsequent iterations.

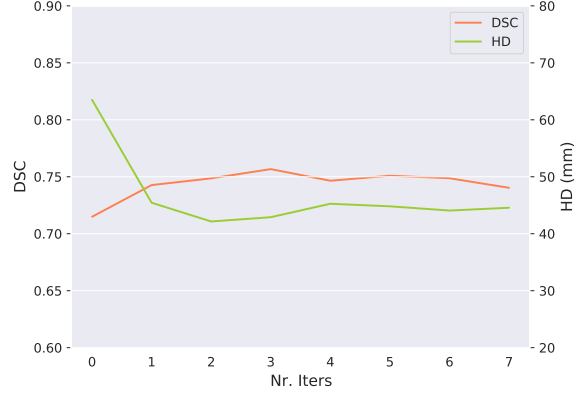


Figure 5: Evolution of the results from the *pseudo-labels* after several iterations (on the test set of ACDC for the Setting-3).

Comparison with semi-supervised methods.

We now compare the proposed approach with a series of state-of-the-art semi-supervised segmentation approaches including: UA-MT (Yu et al., 2019), GLMI (Peng et al., 2021) and SSCO (Wang et al., 2021). UA-MT is an established method to benchmark semi-supervised segmentation approaches, whereas GLMI and SSCO have recently demonstrated superior performance over the existing literature. We carefully tune the hyper-parameters and report the best performance for each method. For each semi-supervised method, we fixed the number of labeled and unlabeled training examples in each mini-batch to 4 as this yield consistent results. For UA-MT, the consistency weight is set to 0.1 and increased by a Gaussian ramp-up function during the first 100 epochs. For GLMI (Peng et al., 2021), the three balancing weights for controlling the relative contributions of global mutual information (MI), local MI and consistency loss are set to 1.0, 0.1 and 10.0, respectively. To evaluate SSCO, the light-weight architecture ENet used in the original implementation (Wang et al., 2021) is replaced by UNet for a fair comparison in our experiments. The values for the hyper-parameters are the

same as those reported in (Wang et al., 2021), as it empirically worked well in our experiment. For detailed parameter search for each method, please refer to Section Appendix B in the Supplemental materials. Table 3 reports the results from this study. We can observe that, under the same conditions, our approach significantly outperforms SSL state-of-the-art methods, particularly in the scenarios where less labeled images are available. For example, when only 3 scans are fully labeled, our method significantly outperforms the recent method in (Wang et al., 2021), with a gap of nearly 10% in terms of DSC. Even though one can argue that our method employs more supervision than these approaches, the cost of it is negligible, as extra labeled pixels could mimic the human behaviour of quickly drawing scribbles in a volumetric scan.

Table 3: Comparison to semi-supervised segmentation approaches on the ACDC test dataset. Results are averaged over three runs.

Setting	Model	DSC	HD-95
Set-3	Lower bound	54.66	80.05
	UA-MT (Yu et al., 2019) *	70.62	39.06
	GLMI (Peng et al., 2021) *	76.27	11.21
	SSCO (Wang et al., 2021) *	77.16	10.10
	Ours (KL+Ent) *	86.94	8.84
Set-5	Lower bound	69.71	51.75
	UA-MT (Yu et al., 2019) *	74.71	30.36
	GLMI (Peng et al., 2021) *	80.58	8.14
	SSCO (Wang et al., 2021) *	81.17	8.15
	Ours (KL+Ent) *	88.77	5.40
Set-10	Lower bound	78.28	44.16
	UA-MT (Yu et al., 2019) *	82.67	28.04
	GLMI (Peng et al., 2021) *	88.43	8.04
	SSCO (Wang et al., 2021) *	89.25	2.58
	Ours (KL+Ent) *	90.92	1.39

*Ours uses partially supervised images.

Ablation study on the importance of the KL term. The objective of this ablation study (Table 4) is to assess the effect of balancing the importance of the KL term in our formulation. Particularly, the KL term plays a crucial role in the proposed formulation, as it guides the entropy term during training to avoid degenerate solutions. We note that the value of the KL term is typically 2 orders of magnitude smaller than the entropy objective. Therefore, by setting its weight (λ_{KL}) to 1, we demonstrate empir-

Table 4: Impact of λ_{KL} on the proposed formulation. Results on the ACDC dataset.

	Set-3		Set-5		Set-10	
	DSC	HD-95	DSC	HD-95	DSC	HD-95
$\lambda_K = 0.1$	66.36	51.86	77.26	31.30	81.95	30.11
$\lambda_K = 1$	71.31	39.67	83.88	21.68	89.73	7.34
$\lambda_K = 10$	85.87	12.38	86.52	7.86	90.55	2.84
$\lambda_K = 50$	86.94	8.84	88.77	5.40	90.92	1.39
$\lambda_K = 100$	83.92	18.17	87.31	9.34	89.38	1.62
$\lambda_K = 1000$	76.20	29.99	85.59	13.93	88.90	4.46

ically its crucial role during training when coupled with the entropy term, as in this setting the latter strongly dominates the training. In this scenario, we observe that the model is negatively impacted, particularly when fully-labeled images are scarce, i.e., Set-3, significantly outperforming the lower bound model. This confirms our hypothesis that minimizing the entropy alone results in degenerated solutions. Increasing the weight of the KL term typically alleviates this issue. However, if much importance is given to this objective the performance also degrades. This is likely due to the fact that the bottom branch is strongly encouraged to follow the behaviour of the top branch, and the effect of the entropy term is diminished.

Sensitivity to λ_w . This ablation study quantifies the contribution of the partially labeled cross-entropy term in Eq. (5) by evaluating the performance across several λ_w values. In this study, the number of partially labeled images is 5 times larger than the number of fully labeled images. Table 5 reports these results, from which we can see that a value of $\lambda_w = 0.001$ consistently brings the best performance across all the settings. Differences are larger in the Set-3 case, which corresponds to the lowest amount of extra information, in terms of partial labels. We believe that, as the amount of additional supervision increases, the performance is less sensitive to the weight of this term.

On the divergence terms. In addition to the widely well-known KL-divergence, we study a series of additional divergences for the constraining term in Eq. (3). In particular, we first resort to the Bhattacharyya distance (Bhattacharyya, 1946). For two discrete distributions $\mathbf{p} = (p_k)_{k=1}^K$ and $\mathbf{q} = (q_k)_{k=1}^K$

Table 5: Impact of λ_w on the proposed formulation. Results on the ACDC dataset.

	<i>Set-3</i>		<i>Set-5</i>		<i>Set-10</i>	
	DSC	HD-95	DSC	HD-95	DSC	HD-95
$\lambda_w = 1$	80.55	20.35	84.09	18.78	90.26	6.99
$\lambda_w = 0.1$	84.09	11.46	87.09	18.03	90.17	4.84
$\lambda_w = 0.01$	85.40	12.05	87.39	9.87	90.14	2.70
$\lambda_w = 0.001$	86.94	8.84	88.77	5.40	90.92	1.39
$\lambda_w = 0.0001$	81.24	19.77	83.68	15.88	87.64	4.39

this term takes the following form:

$$\mathcal{D}_{BC}(\mathbf{p}, \mathbf{q}) = -\log \sum_{k=1}^K (p_k q_k)^{\frac{1}{2}} \quad (6)$$

Furthermore, we also investigate the Tsallis’s formulation of α -divergence (Cichocki and Amari, 2010; Tsallis, 1988), which generalizes the KL:

$$\begin{aligned} \mathcal{D}_\alpha(\mathbf{p} \parallel \mathbf{q}) &= -\sum_{k=1}^K p_k \log_\alpha \left(\frac{q_k}{p_k} \right) \\ &= \frac{1}{1-\alpha} \left(1 - \sum_{k=1}^K p_k^\alpha q_k^{1-\alpha} \right) \end{aligned} \quad (7)$$

Table 6 reports the results generated by the different divergence functionals, evaluated on the *Set-3* setting. We see that there are two cases, i.e., Bhattacharyya distance and α -divergence with $\alpha = 2.0$, that outperforms the model integrating the KL-divergence. This suggests that our model can be further improved by replacing the KL term by alternative divergence functions as consistency losses.

Table 6: Comparison of several divergence functions for the third term (i.e., \mathcal{L}_{kd} term in eq. 5), on the *Set-3* setting. Results on the ACDC dataset.

Setting	Model	DSC	HD-95
<i>Set-3</i>	Kullback-Leibler	86.94	8.84
	Bhattacharyya	86.98	5.26
	α -Divergence ($\alpha = 2.0$)	88.04	6.43
	α -Divergence ($\alpha = 3.0$)	86.89	10.24
	α -Divergence ($\alpha = 5.0$)	85.57	11.77

Figure 6 depicts the evolution, in terms of DSC, on the validation set for the different divergences. Despite showing a similar performance, employing the α -divergence with $\alpha = 2$ (green line) stands out from

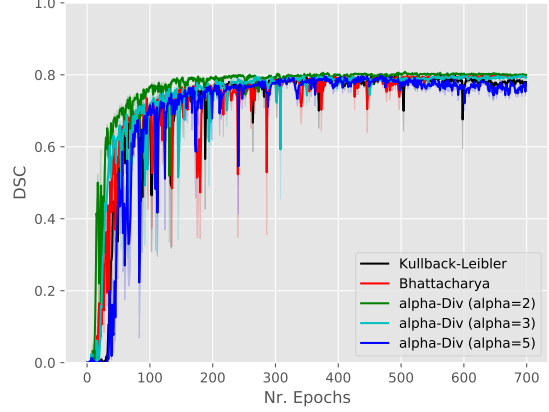


Figure 6: Comparison of different divergences in terms of DSC in the validation set (on the *Set-3* setting).

the others, with a faster convergence at the beginning of the training.

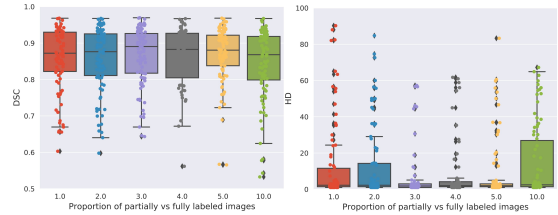


Figure 7: Ablation experiments on the effect of increasing the number of partially labeled images for a fix set of labeled images (on the *Set-3* setting for ACDC). The value in the x axis represents the amount of partially labeled images with respect the labeled images, e.g., 2.0 indicates that there are 2 times more partially labeled than fully labeled images.

Impact on the number of partially labeled images. We now evaluate the impact of training the proposed model with a diverse amount of partially labeled images. These results, which are depicted in Fig 7, show that by having a ratio between labeled and partially labeled data ranging from 3 to 5 typically brings the best performance, both in terms of DSC and HD distance. It is important to note that these results represent the lowest supervised scenario, where only 3 fully labeled images are available. Nev-

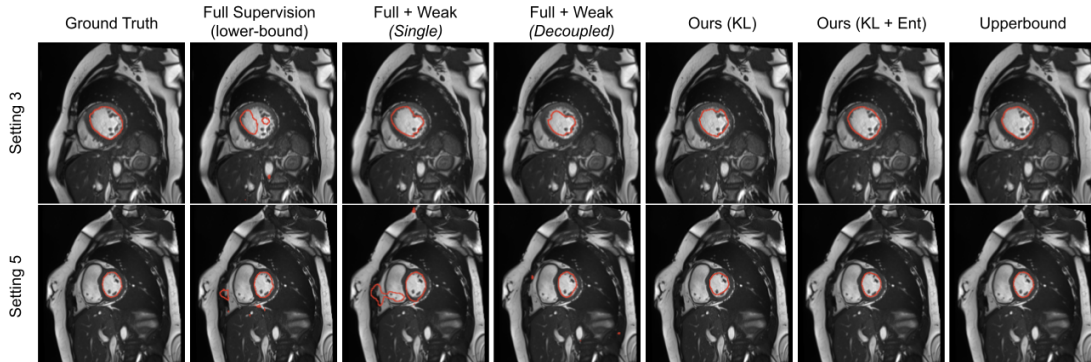


Figure 8: Qualitative results for the analyzed models under two different settings.

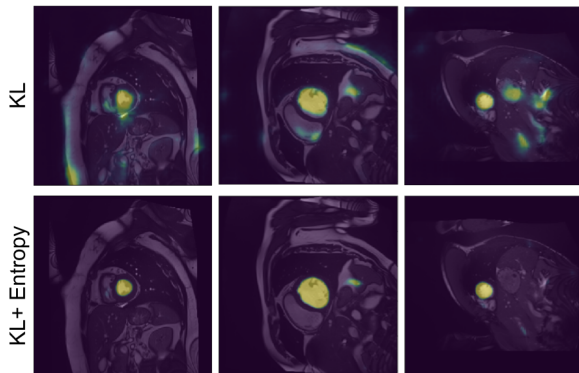


Figure 9: Probability maps obtained by the proposed KL and $KL + Ent$ models.

ertheless, as seen in previous ablation, as the number of fully labeled cases becomes larger, the impact of several elements is reduced.

Qualitative results. In addition to the numerical results presented before, we also depict qualitative results in Fig. 8 and Fig. 9. Particularly, Fig. 8 depicts the segmentation results for the models evaluated in Table 1. We see that results obtained by models with a single network typically under-segment the object of interest (*first row*) or generate many false positives (*second row*). Decoupling the decoding branches might reduce the false positive rate, however, it also tends to under-segment the target. Finally, we ob-

serve that both of our formulations achieve qualitatively better segmentation results, with the $KL+Ent$ model yielding segmentations similar to those generated by the upper bound model. Furthermore, in Fig. 9, we illustrate additional qualitative results of our models. We observe that without the entropy term our model produces less confident predictions, which results in more noisy segmentations.

Results on left-atrium (LA) segmentation. Beyond the ACDC dataset, we performed experiments on the more challenging LA segmentation task, whose results are reported in Table 7. These results align with the observations in the ACDC dataset (Table 1). In particular, the proposed approach outperforms consistently the different baselines across all the settings, with a significant gap in less supervised scenarios. For example, in *Set-3*, our model brings a gain of 15% in terms of DSC compared to the recent work in (Luo and Yang, 2020), whereas the difference amounts to approximately 10 mm in terms of HD. On the other hand, there is a noticeable decline in this gap as the number of fully supervised samples increases. Nevertheless, the differences between our method and the best performing baseline are still remarkable, with nearly 4% in terms of DSC. Furthermore, and similarly to the ACDC dataset, simply adding an entropy minimization term to the Single model does not translate into similar performances to those observed by the proposed model. This empirically *i)* demonstrates that our method is not equiva-

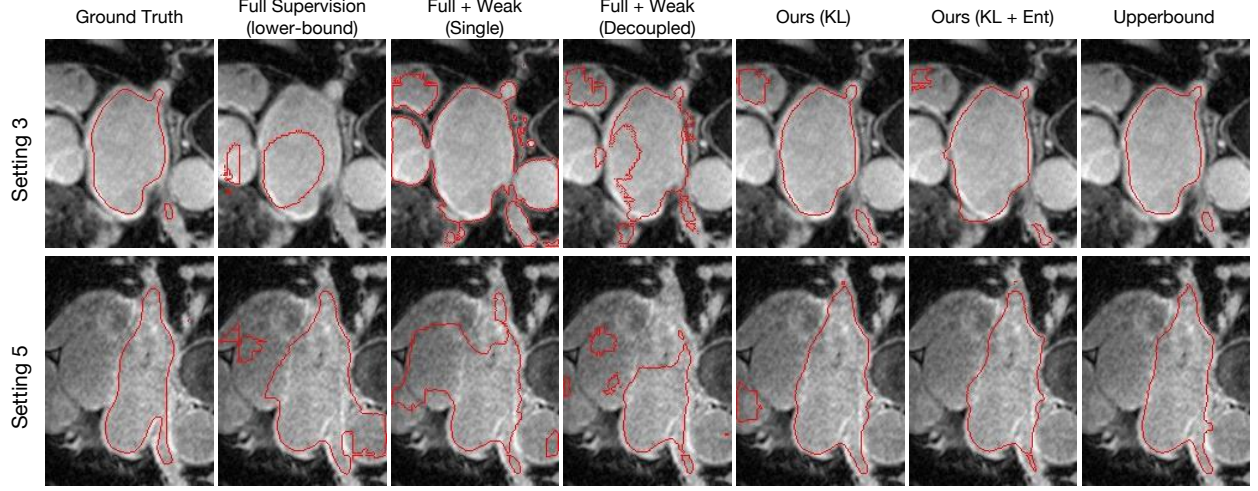


Figure 10: Qualitative results on LA segmentation for the analyzed models under two different settings.

Table 7: Results on the test set of LA segmentation for the *top* and *bottom* branches (when applicable). Results are averaged over three runs. Best results highlighted in bold and second best results are underlined.

				Top		Bottom		Ensemble		
Setting		Model	FS	PS	DSC	HD-95	DSC	HD-95	DSC	HD-95
Set-3	Single Branch	Lower bound	✓	–	37.76	42.30	–	–	–	–
		Single	✓	✓	38.11	34.28	–	–	–	–
		Single + Ent	✓	✓	41.82	36.54	–	–	–	–
		Upper Bound (Set-3)	✓	–	85.91	15.10	–	–	–	–
	Dual Branch	Decoupled (Luo and Yang, 2020)	✓	✓	57.85	42.04	18.11	61.53	–	–
		Ours (KL)	✓	✓	61.41	33.63	66.75	33.71	–	–
Ours (KL+Ent)		✓	✓	69.94	33.05	72.09	32.33	71.29	32.16	
Set-5	Single Branch	Lower bound	✓	–	64.86	35.97	–	–	–	–
		Single	✓	✓	72.06	28.97	–	–	–	–
		Single + Ent	✓	✓	71.79	28.62	–	–	–	–
		Upper Bound (Set-5)	✓	–	86.21	14.01	–	–	–	–
	Dual Branch	Decoupled (Luo and Yang, 2020)	✓	✓	75.33	29.06	17.58	61.18	–	–
		Ours (KL)	✓	✓	76.21	26.62	77.19	28.64	–	–
Ours (KL+Ent)		✓	✓	78.10	23.64	78.50	24.39	78.45	23.34	
Set-10	Single Branch	Lower bound	✓	–	77.65	22.45	–	–	–	–
		Single	✓	✓	79.14	18.22	–	–	–	–
		Single + Ent	✓	✓	81.72	22.99	–	–	–	–
		Upper Bound (Set-10)	✓	–	87.15	11.08	–	–	–	–
	Dual Branch	Decoupled (Luo and Yang, 2020)	✓	✓	79.18	20.16	18.11	61.53	–	–
		Ours (KL)	✓	✓	81.50	21.99	81.01	23.05	–	–
Ours (KL+Ent)		✓	✓	81.79	21.13	83.05	18.93	82.36	20.08	
All images	Single Branch	Upper bound	✓	–	91.30	5.01	–	–	–	–

lent to a model with a single architecture integrating an entropy term and *ii*) supports our hypothesis that decoupling the decoder network to avoid supervision interference yields better segmentation results.

Qualitative evaluation is visually assessed in Fig. 10, which depicts the segmentation results across models on the *Set-3* and *Set-5* settings. Similarly to the visual examples in Fig. 8, single models generate

inconsistent segmentations, which result in both large under and over-segmentations. Even though decoupling single models in dual-stream architectures seem to reduce the amount of false positives, it typically comes at the price of failing to identify target regions. In contrast, both of our models provide a substantial improvement on the segmentation quality, with the model integrating the KL and the entropy terms providing the closest results to the ground-truth.

Model complexity. Last, we evaluate the model complexity, measured in number of parameters for the different analyzed methods (Table 8). Several methods (Peng et al., 2021) employ a single model, whereas other approaches (Yu et al., 2019) need to duplicate this into a teacher-student architecture. It is important to realize that some of these models, e.g., (Yu et al., 2019), have only half of reported parameters to be learned by gradient descent, since the teacher parameters are updated via exponential average moving. Nevertheless, as the parameter values need to be also stored, we have included them in our calculation. In terms of complexity, our model lies in between these two strategies, as despite integrating two decoupled decoders, the encoder is shared among the two branches. On the other hand, the closest method in terms of performance, i.e., (Wang et al., 2019), comes at the price of significant complexity increase, which may hinder its deployment in realistic scenarios.

Table 8: Model complexity in terms of parameters during training and inference time.

Model	#Params	Time (per sample) (ms)
Single Branch	31,042,434	4.9
GLMI (Peng et al., 2021)	31,042,434	4.9
Dual Branch (Ours)	41,137,220	7.1
UA-MT (Yu et al., 2019)	62,084,868	9.9
SSCO (Wang et al., 2021)	248,339,472	37.6

5. Conclusion

In this work we have presented a novel formulation for semantic segmentation under the mixed-supervised paradigm. In addition to the standard

cross-entropy loss over the labeled pixels, we integrate two important terms in the global learning objective, which have demonstrated to bring a significant boost in performance. First, a Shannon entropy loss defined over the less-supervised images encourages confident predictions on these images. Secondly, a KL divergence transfers the knowledge from the predictions generated by the strongly supervised branch to the less-supervised branch. As shown in the experiments, the latter term plays a crucial role in our global learning objective, as it serves as a strong prior for the bottom branch, avoiding trivial solutions resulting from the entropy term.

Furthermore, we have discussed an interesting link between Shannon-entropy minimization and standard pseudo-mask generation, typically used in semi and weakly supervised semantic segmentation. Motivated from a gradient dynamics perspective, we further argue that the former should be preferred over the latter to leverage information from unlabeled pixels. In particular, we show that plugging pseudo-masks in a cross-entropy loss is equivalent to minimizing the min-entropy (Fig 2a). Hence, for uncertain predictions, i.e., in the middle of the simplex, the gradient is much higher than with the entropy, pushing the predictions at the beginning of the training towards the vertex. In addition, we provide empirical evidence that employing pseudo-labels has this undesired effect.

Through extensive experiments, we have rigorously assessed the impact of the different elements of the proposed formulation. Our experiments have further confirmed the usefulness of our method on two publicly available segmentation benchmarks. We have also demonstrated the significant superiority of our approach compared, not only to existing literature in mixed-supervised segmentation, but also to well-known recent semi-supervised methods. It is worth mentioning that, even though our method requires slightly more supervision, the cost of obtaining it is negligible. Furthermore, compared to similar performing methods that require training multiple models, our approach is substantially less complex in terms of number of parameters.

The proposed framework is straightforward to use, does not incur in significant computational costs and can be used with any segmentation network architecture or segmentation loss. Future work will address the integration of other types of supervision in the bottom branch, for example in the form of image tags or anatomical priors.

Acknowledgements

This work is supported by the National Science and Engineering Research Council of Canada (NSERC), via its Discovery Grant program. We also thank Calcul Quebec and Compute Canada.

References

- Bai, W., Oktay, O., Sinclair, M., Suzuki, H., Rajchl, M., Tarroni, G., Glocker, B., King, A., Matthews, P.M., Rueckert, D., 2017. Semi-supervised learning for network-based cardiac MR image segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 253–260.
- Bar, A., Huger, F., Schlicht, P., Fingscheidt, T., 2019. On the robustness of redundant teacher-student frameworks for semantic segmentation, in: CVPRW.
- Bateson, M., Kervadec, H., Dolz, J., Lombaert, H., Ben Ayed, I., 2020. Source-relaxed domain adaptation for image segmentation, in: MICCAI.
- Baur, C., Albarqouni, S., Navab, N., 2017. Semi-supervised deep learning for fully convolutional networks, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 311–319.
- Bernard, O., Lalande, A., Zotti, C., Cervenansky, F., Yang, X., Heng, P.A., Cetin, I., Lekadir, K., Camara, O., Ballester, M.A.G., et al., 2018. Deep learning techniques for automatic MRI cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE TMI* 37, 2514–2525.
- Berthelot, D., Carlini, N., Goodfellow, I., Papernot, N., Oliver, A., Raffel, C.A., 2019. Mixmatch: A holistic approach to semi-supervised learning, in: NeurIPS.
- Bhargat, Y., Shah, M., Awate, S., 2018. Annotation-cost minimization for medical image segmentation using suggestive mixed supervision fully convolutional networks, in: Medical Imaging meets NeurIPS Workshop.
- Bhattacharyya, A., 1946. On some analogues of the amount of information and their use in statistical estimation. *Sankhyā: The Indian Journal of Statistics*, 1–14.
- Bortsova, G., Dubost, F., Hogeweg, L., Katramados, I., de Bruijne, M., 2019. Semi-supervised medical image segmentation via learning consistency under transformations, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 810–818.
- Boudiaf, M., Ziko, I., Rony, J., Dolz, J., Piantanida, P., Ben Ayed, I., 2020. Information maximization for few-shot learning. *NeurIPS* 33.
- Chaitanya, K., Karani, N., Baumgartner, C.F., Becker, A., Donati, O., Konukoglu, E., 2019. Semi-supervised and task-driven data augmentation, in: International Conference on Information Processing in Medical Imaging, Springer. pp. 29–41.
- Chapelle, O., Scholkopf, B., Zien, A., 2009. Semi-supervised learning. *IEEE Transactions on Neural Networks* 20, 542–542.
- Cichocki, A., Amari, S.i., 2010. Families of alpha-beta-and gamma-divergences: Flexible and robust measures of similarities. *Entropy* 12, 1532–1568.
- Cui, W., Liu, Y., Li, Y., Guo, M., Li, Y., Li, X., Wang, T., Zeng, X., Ye, C., 2019. Semi-supervised brain lesion segmentation with an adapted mean teacher model, in: IPMI, pp. 554–565.
- Dhillon, G.S., Chaudhari, P., Ravichandran, A., Soatto, S., 2019. A baseline for few-shot image classification, in: ICLR.

- Dolz, J., Desrosiers, C., Ayed, I.B., 2018. 3D fully convolutional networks for subcortical segmentation in MRI: A large-scale study. *NeuroImage* 170, 456–470.
- Dolz, J., Desrosiers, C., Ben Ayed, I., 2021. Teach me to segment with mixed supervision: Confident students become masters, in: *International Conference on Information Processing in Medical Imaging*, Springer. pp. 517–529.
- Dolz, J., Desrosiers, C., Wang, L., Yuan, J., Shen, D., Ben Ayed, I., 2020. Deep cnn ensembles and suggestive annotations for infant brain mri segmentation. *Computerized Medical Imaging and Graphics* 79, 101660.
- Furlanello, T., Lipton, Z., Tschannen, M., Itti, L., Anandkumar, A., 2018. Born again neural networks, in: *ICML*.
- Grandvalet, Y., Bengio, Y., 2005. Semi-supervised learning by entropy minimization, in: *NeurIPS*.
- Hong, S., Noh, H., Han, B., 2015. Decoupled deep neural network for semi-supervised semantic segmentation, in: *NeurIPS*.
- Hu, M., Maillard, M., Zhang, Y., Ciceri, T., La Barbera, G., Bloch, I., Gori, P., 2020. Knowledge distillation from multi-modal to mono-modal segmentation networks, in: *MICCAI*.
- Huo, X., Xie, L., He, J., Yang, Z., Zhou, W., Li, H., Tian, Q., 2021. ATSO: Asynchronous teacher-student optimization for semi-supervised image segmentation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1235–1244.
- Kervadec, H., Dolz, J., Granger, É., Ben Ayed, I., 2019a. Curriculum semi-supervised segmentation, in: *MICCAI*.
- Kervadec, H., Dolz, J., Tang, M., Granger, E., Boykov, Y., Ayed, I.B., 2019b. Constrained-CNN losses for weakly supervised segmentation. *MedIA* 54, 88–99.
- Kervadec, H., Dolz, J., Wang, S., Granger, E., Ben Ayed, I., 2020. Bounding boxes for weakly supervised segmentation: Global constraints get close to full supervision, in: *MIDL*.
- Khoreva, A., Benenson, R., Hosang, J., Hein, M., Schiele, B., 2017. Simple does it: Weakly supervised instance and semantic segmentation, in: *CVPR*.
- Lee, D.H., 2013. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks, in: *Workshop on challenges in representation learning*, *ICML*.
- Lee, J., Kim, E., Lee, S., Lee, J., Yoon, S., 2019. Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference, in: *CVPR*.
- Lin, D., Dai, J., Jia, J., He, K., Sun, J., 2016. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation, in: *CVPR*.
- Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., Van Der Laak, J.A., Van Ginneken, B., Sánchez, C.I., 2017. A survey on deep learning in medical image analysis. *MedIA* 42, 60–88.
- Luo, W., Yang, M., 2020. Semi-supervised semantic segmentation via strong-weak dual-branch network, in: *ECCV*.
- Milletari, F., Navab, N., Ahmadi, S., 2016. V-net: Fully convolutional neural networks for volumetric medical image segmentation, in: *3DV*.
- Mlynarski, P., Delingette, H., Criminisi, A., Ayache, N., 2019. Deep learning with mixed supervision for brain tumor segmentation. *J. Med. Imaging* 6, 034002.
- Papandreou, G., Chen, L.C., Murphy, K.P., Yuille, A.L., 2015. Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation, in: *ICCV*.

- Peng, J., Estrada, G., Pedersoli, M., Desrosiers, C., 2020a. Deep co-training for semi-supervised image segmentation. *Pattern Recognition* 107, 107269.
- Peng, J., Kervadec, H., Dolz, J., Ben Ayed, I., Pedersoli, M., Desrosiers, C., 2020b. Discretely-constrained deep network for weakly supervised segmentation. *Neural Networks* 130, 297–308.
- Peng, J., Pedersoli, M., Desrosiers, C., 2020c. Mutual information deep regularization for semi-supervised segmentation, in: *MIDL*.
- Peng, J., Pedersoli, M., Desrosiers, C., 2021. Boosting semi-supervised image segmentation with global and local mutual information regularization. *arXiv preprint arXiv:2103.04813*.
- Perone, C.S., Cohen-Adad, J., 2018. Deep semi-supervised segmentation with weight-averaged consistency targets, in: *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*. Springer, pp. 12–19.
- Rajchl, M., Lee, M.C., Oktay, O., Kamnitsas, K., Passerat-Palmbach, J., Bai, W., Damodaram, M., Rutherford, M.A., Hajnal, J.V., Kainz, B., et al., 2016. Deepcut: Object segmentation from bounding box annotations using convolutional neural networks. *IEEE TMI* 36, 674–683.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: *MICCAI*.
- Sedai, S., Antony, B., Rai, R., Jones, K., Ishikawa, H., Schuman, J., Gadi, W., Garnavi, R., 2019. Uncertainty guided semi-supervised segmentation of retinal layers in OCT images, in: *MICCAI*.
- Shah, M.P., Merchant, S., Awate, S., 2018. MS-Net:mixed-supervision fully-convolutional networks for full-resolution segmentation, in: *MICCAI*.
- Tsallis, C., 1988. Possible generalization of boltzmann-gibbs statistics. *Journal of statistical physics* 52, 479–487.
- Vu, T.H., Jain, H., Bucher, M., Cord, M., Pérez, P., 2019. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation, in: *CVPR*.
- Wang, D., Li, M., Ben-Shlomo, N., Corrales, C.E., Cheng, Y., Zhang, T., Jayender, J., 2019. Mixed-supervised dual-network for medical image segmentation, in: *MICCAI*.
- Wang, P., Peng, J., Pedersoli, M., Zhou, Y., Zhang, C., Desrosiers, C., 2021. Self-paced and self-consistent co-training for semi-supervised image segmentation. *Medical Image Analysis* 73, 102146.
- Xu, Y., Du, B., Zhang, L., Zhang, Q., Wang, G., Zhang, L., 2019. Self-ensembling attention networks:addressing domain shift for semantic segmentation, in: *AAAI*.
- Yim, J., Joo, D., Bae, J., Kim, J., 2017. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning, in: *CVPR*.
- Yu, L., Wang, S., Li, X., Fu, C.W., Heng, P.A., 2019. Uncertainty-aware self-ensembling model for semi-supervised 3D left atrium segmentation, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 605–613.
- Zhang, Y., Yang, L., Chen, J., Fredericksen, M., Hughes, D.P., Chen, D.Z., 2017. Deep adversarial networks for biomedical image segmentation utilizing unannotated images, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 408–416.
- Zhou, Y., Wang, Y., Tang, P., Bai, S., Shen, W., Fishman, E., Yuille, A., 2019. Semi-supervised 3D abdominal multi-organ segmentation via deep multi-planar co-training, in: *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE. pp. 121–140.

Appendix A. Parameters search

The proposed model contains several objective terms, each balanced with a different weighting factor. Thus, we sequentially found the best value for each balancing term, and with this term fixed we moved to the next term. We report below the intermediate steps to find the best set of hyperparameters, and their corresponding results.

Appendix A.1. Sensitivity of λ_w on the Single model.

To find the best set of hyperparameters, we first searched the optimal value of λ_w in the *Single* model, which was then fixed to find the remaining hyperparameters. The results from the conducted experiments are reported in Table A.1. We can observe that, similarly to the model trained with the whole learning objective, setting $\lambda_w = 0.001$ provided the best results consistently across the three settings.

Table A.1: Impact of λ_w on the *Single* model. Results on the ACDC dataset.

	Set-3		Set-5		Set-10	
	DSC	HD-95	DSC	HD-95	DSC	HD-95
$\lambda_w = 1$	22.93	125.61	64.21	64.31	69.55	53.10
$\lambda_w = 0.1$	45.20	92.37	63.74	52.01	75.35	50.52
$\lambda_w = 0.01$	55.04	78.15	65.38	52.59	76.94	48.53
$\lambda_w = 0.001$	57.42	78.80	70.73	51.34	78.17	42.99
$\lambda_w = 0.0001$	41.82	91.17	60.93	51.93	74.18	50.98

Appendix A.2. Sensitivity of λ_{KL} on the KL model.

Once the optimal value for λ_w ($\lambda_w = 0.001$) is found, we optimize the λ_{KL} hyperparameter. In particular, with λ_w fixed, we evaluate the performance of the *KL* model across different values of λ_{KL} . The optimal value found will be used in our final model, which also includes the entropy term into the learning objective. We can observe that both $\lambda_{KL} = 50$ and $\lambda_{KL} = 100$ achieve similar results in terms of DSC. Nevertheless, when the value of λ_{KL} is fixed to 50, the values of the HD metric decrease, suggesting that it is a better value for the *KL* term. Hence, we will fix this hyperparameter to 50 in the whole proposed model.

Table A.2: Impact of λ_{KL} on the *KL*-based model. Results on the ACDC dataset.

	Set-3		Set-5		Set-10	
	DSC	HD-95	DSC	HD-95	DSC	HD-95
$\lambda_{KL} = 0.1$	63.37	87.63	74.26	44.63	84.12	32.09
$\lambda_{KL} = 1$	65.70	79.77	75.21	52.29	87.74	21.10
$\lambda_{KL} = 10$	64.12	66.47	75.56	37.45	87.04	27.61
$\lambda_{KL} = 50$	71.41	61.37	81.79	24.07	89.20	9.78
$\lambda_{KL} = 100$	71.23	61.92	81.52	27.71	88.66	12.01
$\lambda_{KL} = 1000$	63.64	77.40	78.60	26.71	86.65	14.05

Appendix A.3. Single vs Double Branch models

We now perform further experiments on the impact of the entropy weighting factor on the single model. These results, which are reported in Table A.3, show that the value selected for the entropy term is actually a good compromise across settings. As this value increases, the performance is degraded due to the entropy term dominating the learning and driving the results towards trivial solutions. On the other hand, if the balancing weight is small, it resembles to the baseline approach. Note that, however, the results obtained by this model (Single + Entropy) are not comparable to the performance of the proposed formulation, as suggested by the reviewer.

Table A.3: Impact of λ_{Ent} on the Single-based model that integrates only the entropy as additional term, whose overall learning objective is $\mathcal{L}_s + \lambda_w \mathcal{L}_w + \lambda_{ent} \mathcal{H}(\mathbf{p})$. Results on the ACDC dataset.

	Set-3		Set-5	
	DSC	HD-95	DSC	HD-95
$\lambda_{Ent} = 0.0$	57.42	78.80	70.73	51.34
$\lambda_{Ent} = 0.1$	51.32	84.45	69.62	48.53
$\lambda_{Ent} = 1$	43.01	83.98	74.92	55.24
$\lambda_{Ent} = 10$	38.61	100.86	62.66	54.34

Appendix B. Parameter search for semi-supervised methods

During the validation of our model we found that the default hyperparameters of compared methods were not optimal in our setting. Thus, for a fair comparison, we searched the optimal hyperparameters on ACDC dataset for each semi-supervised method as well, and reported the best scores (shown in Table 3

of the main text). Here, we include the intermediate steps to find the best set of hyperparameters for each method and their corresponding results. Note that we tune the parameters under the Set-3 setting and use the best findings for the other two settings, i.e. Set-5 and Set-10.

Appendix B.1. UA-MT

The formulation of UA-MT (Yu et al., 2019) consists of two terms, i.e., the supervised loss and the unsupervised consistency loss for measuring the consistency between the prediction of the teacher and the student model, with a balancing consistency weight λ_c controlling their relative contribution. We searched the optimal λ_c and reported the corresponding results in Table B.4, where we can see that $\lambda_c = 0.1$ provided the best results.

Table B.4: Impact of the consistency weight (λ_c) on UA-MT (Yu et al., 2019). The study is conducted on the ACDC dataset (Set-3 setting).

	DSC	HD-95
$\lambda_c = 0.01$	67.90	42.57
$\lambda_c = 0.1$	70.62	38.06
$\lambda_c = 1$	68.14	45.01
$\lambda_c = 5$	67.31	54.96
$\lambda_c = 10$	67.45	43.06

Appendix B.2. GLMI

The formulation of GLMI (Peng et al., 2021) includes a composite loss with four terms, i.e., supervised loss, the mutual information loss on global feature embedding and local feature embedding, and the consistency loss on different transformation of the same given input. Hence, three hyper-parameters are introduced for tuning the relative contributions of the last three terms, i.e., λ_{MI}^{global} , λ_{MI}^{local} and λ_{cons} . In our implementations, the best hyper-parameter setting was slightly different from that suggested in (Peng et al., 2021). The reason for this could be that we focus on the task of left ventricular (LV) endocardium segmentation and the potential difference on the experimental environment. Specifically, we empirically fixed λ_{cons} to 10 and report the results of different λ_{MI}^{global} and λ_{MI}^{local} , as shown in Table B.5. The best

result is achieved by setting λ_{MI}^{global} and λ_{MI}^{local} to 1 and 0.1 respectively.

Table B.5: Impact of λ_{MI}^{global} and λ_{MI}^{local} on GLMI (Peng et al., 2021). Dice score is reported for each setting and the study is conducted on ACDC dataset (Set-3).

	$\lambda_{MI}^{local} = 0.05$	$\lambda_{MI}^{local} = 0.1$	$\lambda_{MI}^{local} = 1$
$\lambda_{MI}^{global} = 0.1$	70.61	70.07	69.99
$\lambda_{MI}^{global} = 1$	74.31	76.27	72.85
$\lambda_{MI}^{global} = 2$	73.76	75.79	74.70

Appendix B.3. SSCO

In SSCO (Wang et al., 2021), two hyper-parameters, i.e., λ_1 and λ_2 , are introduced to balance the relative contributions of the self-paced co-training loss and the self-consistency loss. We tune the two hyper-parameters based on the values reported in the original paper (Wang et al., 2021) and found the values reported in the paper ($\lambda_1 = 0.5$, $\lambda_2 = 4$) provided best scores in our implementation as well. Detailed comparison results are shown in Table B.6

Table B.6: Impact of λ_1 and λ_2 on SSCO (Wang et al., 2021). Dice score is reported for each setting, and the study is conducted on ACDC dataset (Set-3).

	$\lambda_2 = 1$	$\lambda_2 = 4$	$\lambda_2 = 8$
$\lambda_1 = 0.1$	75.08	74.80	73.38
$\lambda_1 = 0.5$	74.71	77.16	75.95
$\lambda_1 = 1$	74.01	75.68	74.83

Appendix C. Additional failure cases from Pseudo-Labels

In this section we show additional failure cases obtained from the *Pseudo-labels* approach. As shown in Figure C.1, It is important to mention that these pseudo-labels are used in subsequent iterations to train the deep network, which brings a high risk of propagating these errors through the training. These visual results, which are supported by the quantitative evaluation in the main paper, support our hypothesis that minimizing entropy should be preferred over the use of pseudo-labels, as the errors are propagated through the training.

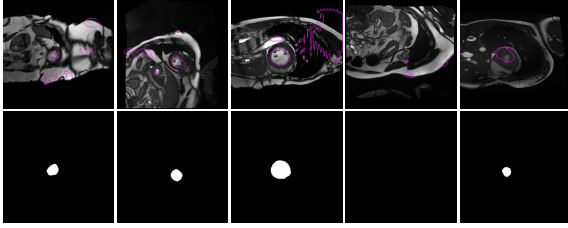


Figure C.1: Additional failure cases (*Set-3*) which are employed as pseudo-labels in the *proposals-based* approach (top), whose errors are reinforced during training, and their corresponding ground truth. Best viewed in colours.

Appendix D. Qualitative comparison with the upperbound

We depict in Figures D.2 and D.3 several visual results on the *Set-3* and *Set-5* from our method compared to the individual upperbounds. In this scenario, both models are trained with the same images, being the only difference the level of supervision provided. We can observe that despite the slight quantitative differences in Table 1 of the main paper, visual results indicate that the proposed method can achieve very similar results to the upperbound model, sometimes generating more similar segmentations to the ground truth.

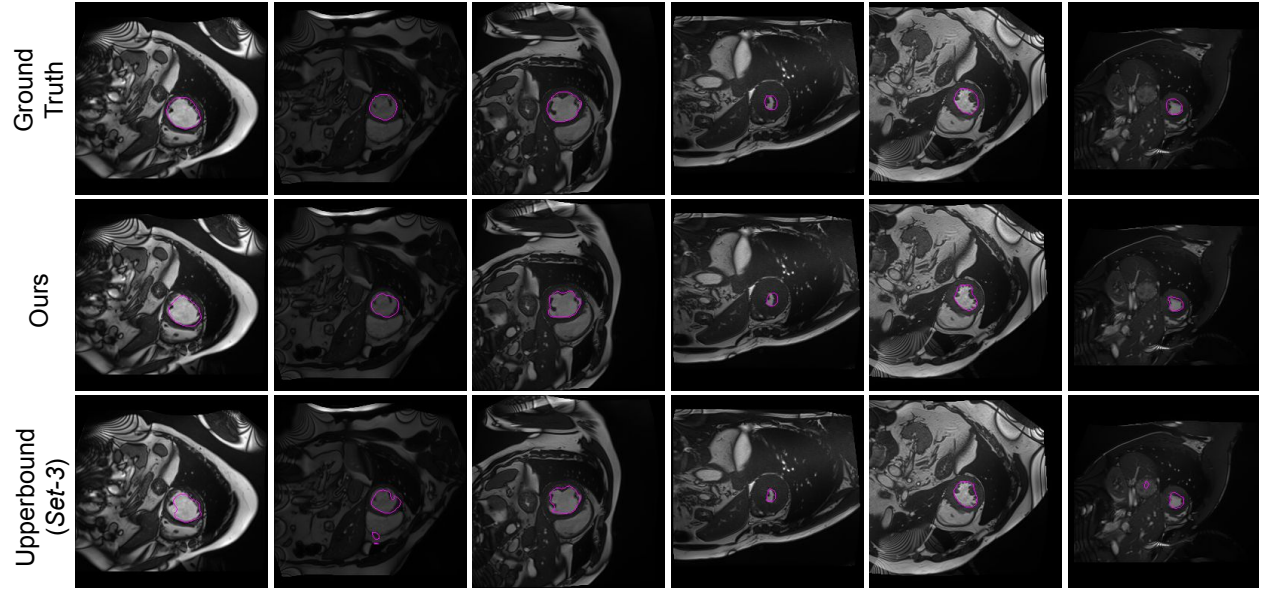


Figure D.2: Qualitative results on the Test dataset of *Set-3* obtained by our model (*middle*) and the upperbound model (*top*) on the same dataset. The corresponding ground truth is depicted in the (*top*) row.

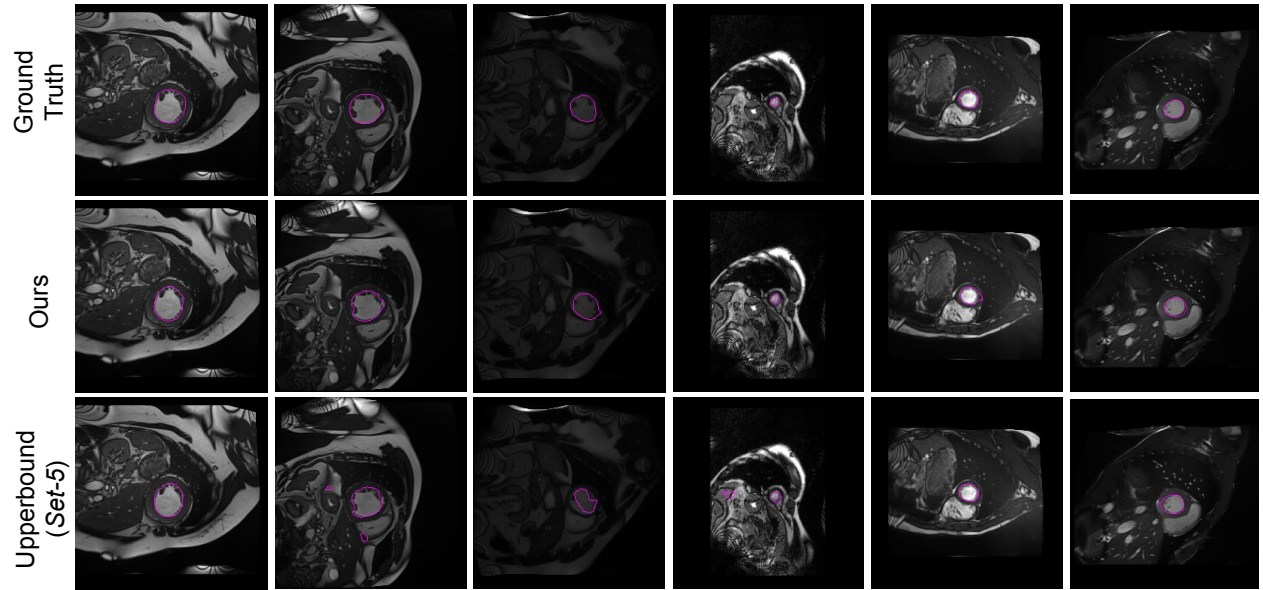


Figure D.3: Qualitative results on the Test dataset of *Set-5* obtained by our model (*middle*) and the upperbound model (*top*) on the same dataset. The corresponding ground truth is depicted in the (*top*) row.