

# A MULTI-TASK U-NET FOR SEGMENTATION WITH LAZY LABELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The need for labour intensive pixel-wise annotation is a major limitation of many fully supervised learning methods for image segmentation. In this paper, we propose a deep convolutional neural network for multi-class segmentation that circumvents this problem by being trainable on coarse data labels combined with only a very small number of images with pixel-wise annotations. We call this new labelling strategy ‘lazy’ labels. Image segmentation is then stratified into three connected tasks: rough detection of class instances, separation of wrongly connected objects without a clear boundary, and pixel-wise segmentation to find the accurate boundaries of each object. These problems are integrated into a multi-task learning framework and the model is trained end-to-end in a semi-supervised fashion. The method is demonstrated on two segmentation datasets, including food microscopy images and histology images of tissues respectively. We show that the model gives accurate segmentation results even if exact boundary labels are missing for a majority of the annotated data. This allows more flexibility and efficiency for training deep neural networks that are data hungry in a practical setting where manual annotation is expensive, by collecting more lazy (rough) annotations than precisely segmented images.

## 1 INTRODUCTION

Image segmentation has been an active research field in the past decades. Deep learning approaches play an increasingly important role and have become state-of-the-art in various segmentation tasks (Huang et al., 2018; Khoreva et al., 2017; Tsutsui et al., 2018; Ghosh et al., 2018; Litjens et al., 2017). Though fully supervised segmentation neural networks have shown great success, one of their most challenging issues is the need for pixel-level annotations to train them. Obtaining such annotations usually requires a great amount of manual work and is therefore expensive.

In this paper, we propose a multi-class and multi-instance segmentation approach that we split into three relevant tasks: detection, separation and segmentation (cf. Figure 1). Doing so, we obtain a semi-supervised learning approach that is trained with so-called "lazy" labels, that is a lot of coarse annotations of class instances together with only a few pixel-wise annotated images that can be obtained from the coarse labels in a semi-automated way. In the following, we will refer to weak (resp. strong) annotations for coarse (resp. accurate) labels and denote them as WL and SL.

**Task 1** detects and classifies each object and roughly determines its region through an under-segmentation mask. Instance counting can be obtained as a by-product of this task. As the main objective is instance detection, exact labels for the whole object or its boundary are not necessary at this stage. We use instead weakly annotated images in which a rough region inside each object is marked, cf. the most top left part of Figure 1. For segmentation problems with a dense population of instances, such as the food components (see e.g., Figure 1), cells (Guerrero-Pena et al., 2018; Ronneberger et al., 2015), glandular tissue, or people in a crowd (Wang et al., 2018b), separating objects sharing a common boundary is a well known challenge. We can optionally perform a second task (**Task 2**) that focuses on the separation of instances that are connected without a clear boundary dividing them. Also for this task we rely on WL to reduce the burden of manual annotations: touching interfaces are specified with rough scribbles, cf. top left part of Figure 1. **Task 3** finally tackles the pixel-wise classification of the instances. It requires strong annotations that are accurate up to the boundaries of the objects. Thanks to the information brought by the weak annotations, we here just

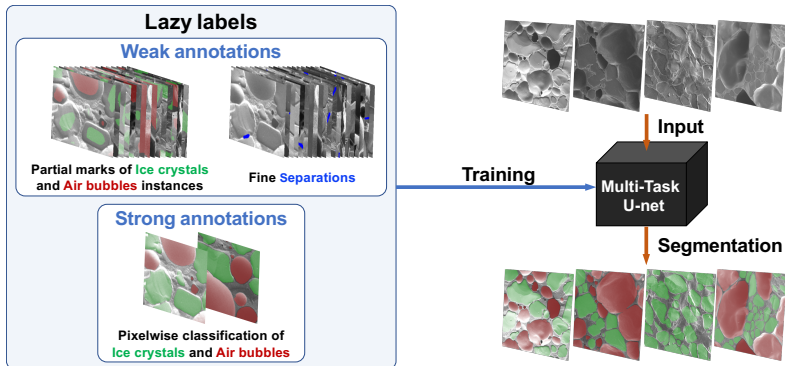


Figure 1: Multi-task learning for image segmentation with lazy labels. The figure uses Scanning Electron Microscopy (SEM) images of food microstructures as an example and demonstrates a segmentation problem of three classes, namely air bubbles (green), ice crystals (red) and background respectively. Most of the training data are weak annotations containing (i) partial marks of ice crystals and/or air bubbles instances and (ii) fine separation marks of boundaries shared by different instances. Only a few strongly annotated images are used. On the top right several SEM images are displayed. Their corresponding output, obtained with the trained network, are shown at the bottom right.

need a very small set of accurate segmentation masks, cf. bottom left part of Figure 1. To that end, we propose to refine some of the coarse labels resulting from task 1 using a semi-automatic segmentation method which requires additional manual intervention.

The three tasks are handled by a single deep neural network and are jointly optimized. The network architecture is inspired by the widely used segmentation network named U-net (Ronneberger et al., 2015). With all three tasks sharing the same contracting path, we introduce a new multi-task block for the expansive path. The network has three outputs and is fed with a combination of WL and SL described above. Since weakly and strongly annotated training data is shared between the tasks, part of the annotations are missing, especially for task 3. To accommodate for this we introduce a weighted loss function over the samples. Accurate segmentation labels for training are usually not easy to obtain, however, with our approach we demonstrate that exact labels for the whole training set are not needed for good segmentation learning performance.

We evaluate the performance of the proposed approach on two applications, namely for the segmentation of SEM images of food microstructure and stained histology images of glandular tissues.

In summary, our contributions are as follows. (1). We propose a decomposition of the segmentation problems into three tasks and a corresponding user friendly labeling strategy. (2). We develop a multi-task learning framework that learns directly from the manual labels and is trained end-to-end. Our approach outperforms state-of-the-art weakly supervised methods such as Khoreva et al. (2017) (3). The network predicts segmentation mask as well as the object inner regions and touching object interfaces. Touching objects are effectively disconnected as a side product.

## 2 RELATED WORK

In image segmentation problems, one needs to classify an image at pixel level. It is a vast topic with a diversity of algorithms being developed, including traditional unsupervised methods like  $k$ -means that splits the image into homogeneous regions according to image low level features, curve evolution based methods like snakes (Caselles et al., 1997), graph-cut based methods like Grabcut (Rother et al., 2004), just to name a few. Interactive approaches such as snakes or Grabcut enable getting involved users’ knowledge by means of initializing regions or putting constraints on the segmentation results.

Deep convolutional neural network (DCNN) approaches have been developed for segmenting complex images, especially in the semantic setting. In particular, fully convolutional networks (FCN) (Long et al., 2015) replace the last few fully connected layers of a conventional classification network by up-sampling layers and convolutional layers, to preserve spatial information. FCNs have many variants for semantic segmentation. The DeepLab (Chen et al., 2018) uses a technique called

trous convolution to handle spatial information together with a fully connected conditional random field (CRF) (Chen et al., 2014) for refining the segmentation results. Fully connected CRF can be used as post-processing or can be integrated into the network architecture, allowing for end-to-end training (Zheng et al., 2015).

Other types of FCNs are deep encoder-decoder networks (Badrinarayanan et al., 2017; Ronneberger et al., 2015). They have multiple up-sampling layers for better localizing boundary details. One of the most well-known models is the U-net (Ronneberger et al., 2015). It is a fully convolutional network made of a contracting path, which brings the input images into very low resolution features with a sequence of down-sampling, and an expansive path that has an equal amount of up-sampling layers. The higher resolution feature maps on the contracting path are merged with up-sampled layers via long skip connections to recover boundary information after down-sampling.

**Weakly/semi-supervised segmentation learning.** Motivated by the heavy cost of pixel-level annotation needed for fully supervised learning, there has been a growing interest in weakly supervised learning with coarse annotations for semantic segmentation in the last years. Common weak annotations include image-level labels (Huang et al., 2018; Papandreou et al., 2015; Lee et al., 2019), bounding boxes (Khoreva et al., 2017), scribbles (Lin et al., 2016) and points (Bearman et al., 2016). Most weakly supervised deep learning methods for segmentation are built on top of a classification network. The training of such networks may be realized using segmentation masks explicitly generated from weak annotations (Wei et al., 2017; Khoreva et al., 2017; Lee et al., 2019; Wei et al., 2018; Tsutsui et al., 2018). The segmentation masks can be improved recursively, which involves several rounds of training of the segmentation network (Wei et al., 2017; Jing et al., 2018; Ezhov et al., 2018).

Other techniques do not use any explicit pseudo segmentation masks as training data, but instead compute a composite loss from other guiding principles. For instance, the SEC method (Kolesnikov & Lampert, 2016) combines localization cues (Seed), classification coherence (Expand) and Constraint-to-boundary with a fully-connected CRF. Special modules could be designed to further exploit the information of sources, such as the Deep Seeded Region Growing (Huang et al., 2018) and saliency seeded region growing (Sun & Li, 2019).

Semi-supervised deep learning provides an alternative way to reduce the pixel-level annotation burden (Baur et al., 2017; Perone & Cohen-Adad, 2018). These methods require only a small amount of strongly annotated data and a large set of unlabelled data. Several training methods that combine weak annotations with a limited set of strong annotations have already been explored (Papandreou et al., 2015; Wei et al., 2018; Lee et al., 2019). The pixel-wise labels can be integrated through an additional loss function, and used along with their semi-supervised counterparts. Better weak supervision performance can be obtained by feeding the network with learned object localization maps (Wei et al., 2018; Lee et al., 2019). However, the networks are trained with image level labels from large scale datasets like ImageNet (Deng et al., 2009) or PASCAL VOC 2012 (Everingham et al., 2015).

The weakly and semi-supervised deep learning methods are also explored for object localization or segmentation on relatively smaller datasets, especially in medical imaging. For these methods, the weak annotations can be of image level (Mlynarski et al., 2018; Playout et al., 2019; Zhou et al., 2019; Shin et al., 2019) or in forms of bounding boxes (Shah et al., 2018).

**Multi-task learning.** Multi-task learning algorithms consider several related tasks to improve the overall performance, taking benefits from the underlying common information that may be ignored by a single task learning. In deep neural based multi-task models, the common information is conveyed by soft or hard parameter sharing (Ruder, 2017). Various multi-task deep learning methods have been developed for segmentation, for example, the stacked U-net for extracting roads from satellite imagery (Sun et al., 2018), the two stage 3D U-net framework for 3D CT or MR data segmentation (Wang et al., 2018a), encoder decoder networks for depth regression, semantic and instance segmentation (Kendall et al., 2018), and for building footprint segmentation (Bischke et al., 2017). Compared to these works, our method handles the tasks with a multi-task block at different feature resolutions and is designed upon both strong and weak notations.

Our work is more closely related to the multi-task learning methods in (Playout et al., 2018; 2019) for retinal lesions segmentation and (Mlynarski et al., 2018) for brain tumour segmentation. Three tasks are considered in (Playout et al., 2019): red lesion segmentation, bright lesions segmentation

and image level lesion detection. The proposed encoder-decoder architecture has one branch in the encoder part and two branches into decoder part specified for the red/bright lesion segmentation and the image level classification respectively. The network is weakly supervised as only image level labels are used in one of the training phases. In (Mlynarski et al., 2018), a U-net architecture is used for jointly segmenting and classifying brain tumours. An additional branch takes the last but one layer of a standard segmentation U-net as input, followed by a mean pooling layer, and outputs a score for the classification task.

Our network structure, similar to these architectures, shares an encoder part for all tasks. The learning is also supervised by a mixture of strong annotations and weak annotations. However, the proposed method does not include supervision from image level annotations and is more specialized in distinguishing the different object instances and clarifying their boundaries in every single image.

### 3 MULTI-TASK LEARNING FRAMEWORK

Fully supervised learning for segmentation approximates the conditional probability distribution of the segmentation mask given the image. Let  $\mathbf{s}^{(3)}$  be the segmentation mask and  $I$  be the image, then the segmentation task aims to estimate  $p(\mathbf{s}^{(3)} | I)$  based on a set of sample images  $\mathcal{I} = \{I_1, I_2, \dots, I_n\}$  and the corresponding labels  $\{\mathbf{s}_1^{(3)}, \mathbf{s}_2^{(3)}, \dots, \mathbf{s}_n^{(3)}\}$ . The set  $\mathcal{I}$  is randomly drawn from an unknown distribution. In our setting, having the whole set of segmentation labels  $\{\mathbf{s}_i^{(3)}\}_{1, \dots, n}$  is impractical, and we introduce two auxiliary tasks for which the labels can be more easily generated to achieve an overall small cost on labeling.

For a given image  $I \in \mathcal{I}$ , we denote as  $\mathbf{s}^{(1)}$  the rough instance detection mask,  $\mathbf{s}^{(2)}$  a map containing some interfaces shared by touching objects. To create  $\mathbf{s}^{(1)}$ , the contours of the objects are not treated carefully, resulting in a coarse mask that misses most of the boundary pixels, cf, the left of Figure 1. Let  $\mathcal{I}_k \subset \mathcal{I}$  denote the subset of images labelled for task  $k$  ( $k = 1, 2, 3$ ). As we collect a different amount of annotations for each task, the number of annotated images  $|\mathcal{I}_k|$  may not be the same for different  $k$ . Typically the number of images with strong annotations  $|\mathcal{I}_3| \ll n$ .

The set of samples in  $\mathcal{I}_3$  for segmentation being small, the computation of an accurate approximation of the true probability distribution  $p(\mathbf{s}^{(3)} | I)$  is a challenging issue. Given that much more samples of  $\mathbf{s}^{(1)}$  and  $\mathbf{s}^{(2)}$  are observed, it is relatively easier to learn the statistics of the weak labels. Therefore, in a multi-task learning setting, one aims at approximating also the conditional probability  $p(\mathbf{s}^{(1)} | I)$  and  $p(\mathbf{s}^{(2)} | I)$  for the other two tasks, or the joint probability  $p((\mathbf{s}^{(1)}, \mathbf{s}^{(2)}, \mathbf{s}^{(3)}) | I)$ . The three tasks are connected to each other. By the definition of the detection task, it is not difficult to see that  $p(\mathbf{s}^{(3)} = \mathbf{z} | \mathbf{s}^{(1)} = \mathbf{x}) = 0$  for  $\mathbf{x}$  and  $\mathbf{z}$  satisfying  $x_{i,c} = 1$  and  $z_{i,c} = 0$  for some pixel  $i$  and class  $c$  other than the background. The map of interfaces  $\mathbf{s}^{(2)}$  indicates small gaps between two connected instances, and is therefore a subset of boundary pixels of the mask  $\mathbf{s}^{(3)}$ .

Let us consider the probabilities given by the models  $p(\mathbf{s}^{(k)} | I; \theta)$  ( $k = 1, 2, 3$ ) parametrized by  $\theta$  which is determined such that the models match the desired probability distributions. The parameter  $\theta$  is shared among all the tasks. We do not optimize  $\theta$  for each individual task, but instead consider a joint probability  $p((\mathbf{s}^{(1)}, \mathbf{s}^{(2)}, \mathbf{s}^{(3)}) | I; \theta)$ . Assuming that  $\mathbf{s}^{(1)}$  (rough under-segmented instance detection) and  $\mathbf{s}^{(2)}$  (subset of shared boundaries) are conditionally independent given image  $I$ , and if the samples are i.i.d, we define the maximum likelihood (ML) estimator for  $\theta$  as

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{I \in \mathcal{I}} \left( \log p(\mathbf{s}^{(3)} | \mathbf{s}^{(1)}, \mathbf{s}^{(2)}, I; \theta) + \sum_{k=1}^2 \log p(\mathbf{s}^{(k)} | I; \theta) \right). \quad (1)$$

The set  $\mathcal{I}_3$  may not be evenly distributed across  $\mathcal{I}$ , but we assume that it is generated by a fixed distribution as well. Provided that the term  $\left\{ p(\mathbf{s}^{(3)} | \mathbf{s}^{(1)}, \mathbf{s}^{(2)}, I) \right\}_{I \in \mathcal{I}}$  can be approximated correctly by  $p(\mathbf{s}^{(3)} | \mathbf{s}^{(1)}, \mathbf{s}^{(2)}, I; \theta)$  even if  $\theta$  is computed without  $\mathbf{s}^{(3)}$  specified for  $\mathcal{I} \setminus \mathcal{I}_3$ , then

$$\sum_{I \in \mathcal{I}} \log p(\mathbf{s}^{(3)} | \mathbf{s}^{(1)}, \mathbf{s}^{(2)}, I; \theta) \propto \sum_{I \in \mathcal{I}_3} \log p(\mathbf{s}^{(3)} | \mathbf{s}^{(1)}, \mathbf{s}^{(2)}, I; \theta). \quad (2)$$



If furthermore, the segmentation mask does not depend on  $\mathbf{s}^{(1)}$  or  $\mathbf{s}^{(2)}$  given  $I \in \mathcal{I}_3$ , and if  $|\mathcal{I}_1|, |\mathcal{I}_2|$  are large enough, then from Equations (1), and (2), we approximate the ML estimator by

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{k=1}^3 \left( \alpha_k \sum_{I \in \mathcal{I}_k} \log p(\mathbf{s}^{(k)} | I; \boldsymbol{\theta}) \right) \quad (3)$$

in which  $\alpha_1, \alpha_2, \alpha_3$  are non negative constants.

### 3.1 LOSS FUNCTION

Let the output of the approximation models be denoted respectively by  $h_{\boldsymbol{\theta}}^{(1)}(I), h_{\boldsymbol{\theta}}^{(2)}(I)$ , and  $h_{\boldsymbol{\theta}}^{(3)}(I)$ , with  $[h_{\boldsymbol{\theta}}^{(k)}(I)]_{i,c}$  the estimated probability of pixel  $i$  being in class  $c$  of task  $k$ . For the label  $\mathbf{s}^{(k)}$  of  $I$ , the log likelihood function for each task is decomposed into

$$\log p(\mathbf{s}^{(k)} | I; \boldsymbol{\theta}) = \sum_i \sum_{c \in C_k} s_{i,c}^{(k)} \log [h_{\boldsymbol{\theta}}^{(k)}(I)]_{i,c}, \quad k = 1, 2, 3, \quad (4)$$

in which  $s_{i,c}^{(k)}$  denotes the element of the label  $\mathbf{s}^{(k)}$  at pixel  $i$  for class  $c$  and  $C_k$  is the set of classes for task  $k$ . For example, for ice cream images, we have three classes including air bubbles, ice crystals and the rest (background or parts of the objects ignored by the weak labels), so  $C_1, C_3 = \{1, 2, 3\}$ . For the separation task, there are only two classes for pixels (belonging or not to a touching interface):  $C_2 = \{1, 2\}$ . According to Equation (3), the network is trained by minimizing the weighted cross entropy loss:

$$L(\boldsymbol{\theta}) = - \sum_{I \in \mathcal{I}} \sum_{k=1}^3 \alpha_k \mathbb{1}_{\mathcal{I}_k}(I) \log p(\mathbf{s}^{(k)} | I; \boldsymbol{\theta}), \quad (5)$$

Here  $\mathbb{1}_{\mathcal{I}_k}(\cdot)$  is an indicator function which is 1 if  $I \in \mathcal{I}_k$  and 0 otherwise.

### 3.2 MULTI-TASK NETWORK

We follow a convolutional encoder-decoder network structure for the multi-task learning. The network architecture is illustrated in Figure 2. As an extension of the U-net structure for multiple tasks, we only have one contracting path that encodes shared features representation for all the tasks. On the expansive branch, we introduce a multi-task block at each resolution to support different learning purposes (blue blocks in Figure 2). Every multi-task block runs three paths, with three inputs and three corresponding outputs, and it consists of several sub-blocks.

In each multi-task block, the detection task (task 1) and the segmentation task (task 3) have a common path similar to the decoder part of the standard U-net. They share the same weights and use the same concatenation with feature maps from contracting path via the skip connections. However, we insert an additional residual sub-block for the segmentation task. The residual sub-block provides extra network parameters to learn information not known from the detection task, *e.g.* object boundary localization. The path for the separation task (task 2) is built on the top of detection/segmentation ones. It is also a U-net decoder block structure, but the long skip connections start from the sub-blocks of the detection/segmentation paths instead of the contracting path. The connections extract higher resolution features from the segmentation task and use them in the separation task.

To formulate the multi-task blocks, let  $\mathbf{x}_l$  and  $\mathbf{z}_l$  denote respectively the output of the detection path and segmentation path at the multi-task block  $l$ , and let  $\mathbf{c}_l$  be the feature maps received from the contracting path with the skip connections. Then for task 1 and task 3 we have

$$\mathbf{x}_{l+1} = F_{W_l}(\mathbf{x}_l, \mathbf{c}_l), \quad \mathbf{z}_{l+\frac{1}{2}} = F_{W_l}(\mathbf{z}_l, \mathbf{c}_l), \quad \mathbf{z}_{l+1} = \mathbf{z}_{l+\frac{1}{2}} + F_{W_{l+\frac{1}{2}}}(\mathbf{z}_{l+\frac{1}{2}}), \quad (6)$$

in which  $W_l, W_{l+1/2} \in \boldsymbol{\theta}$  are parameters of the network and  $F_{W_l}, F_{W_{l+\frac{1}{2}}}$  are determined respectively by a sequence of layers of the network (Cf. the small grey blocks on the right of Figure 2). For task 2 the output at  $l^{\text{th}}$  block  $\mathbf{y}_{l+1}$  is computed by  $\mathbf{y}_{l+1} = G_{\tilde{W}_l}(\mathbf{z}_{l+1}, \mathbf{y}_l)$  with additional network parameters  $\tilde{W}_l \in \boldsymbol{\theta}$ . Finally, after the last multi-task block, softmax layers are added, outputting a probability map for each task.

**Implementation details.** We implement a multi-task U-net with 6 levels of spatial resolution and input images of size  $256 \times 256$ . A sequence of down-sampling via max-pooling with pooling

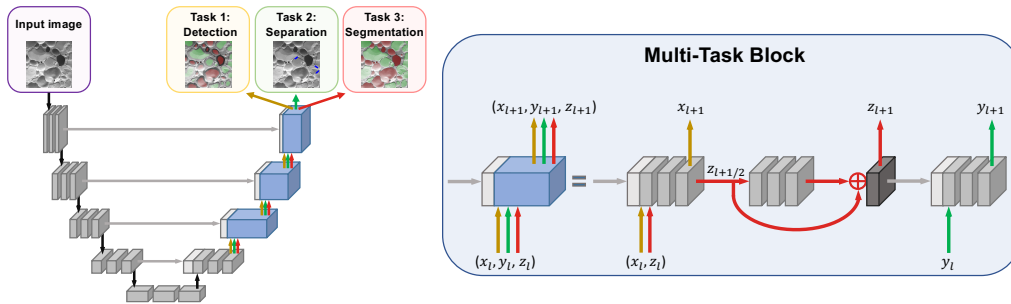


Figure 2: Architecture of the multi-task U-net. The left part of the network is a contracting path similar to the standard U-net. For multi-task learning, we construct several expansive paths with specific multi-task blocks. At each resolution, task 1 (Detection in yellow) and task 3 (Segmentation in red) run through a common sub-block, but the red path learns an additional residual to better localize object boundaries. Long skip connections with the layers from contracting path are built for yellow/red paths via concatenation. Task 2 (Separation, in green) mainly follows a separated expansive path, with its own up-sampled blocks. A link with the last layer of task 3 is added via a skip connection in order to integrate accurate boundaries in the separation task.

size  $2 \times 2$  is used for the contracting path of the network. Different from the conventional U-net (Ronneberger et al., 2015), each small gray block (see Figure 2) consists of a convolution layer and a batch normalization (Ioffe & Szegedy, 2015), followed by a leaky ReLU activation with a leakiness parameter 0.01. The same setting is also applied to gray sub-blocks of the 4 multi-task blocks. On the expansive path of the network, feature maps are up-sampled (with factor  $2 \times 2$ ) by bilinear interpolation from a low resolution multi-task block to the next one.

### 3.3 METHODS FOR LAZY LABELS GENERATION

We now explain our strategy for generating all the lazy annotations that are used for training the multi-task U-net. We introduce our method with a data set of ice cream SEM images but any other similar dataset could be used. Typical images of ice cream samples are shown in the top row of the left part of Figure 3. The segmentation problem is challenging since the images contain densely distributed small object instances (*i.e.*, air bubble and ice crystals), and poor contrast between the foreground and the background. The sizes of the objects can vary significantly in a single sample. Textures on the surfaces of objects also appear.

As a first step, scribble-based labelling is applied to obtain detection regions of air bubbles and ice crystals for task 1. This can be done in a very fast way as no effort is put on the exact object boundaries. We adopt a lazy strategy by picking out an inner region for each object in the images (see *e.g.*, the second row of the left part of Figure 3). Though one could get these rough regions as accurate as possible, we delay such refinement to task 3, for better efficiency of the global annotation process. Compared to the commonly used bounding box annotations in computer vision tasks, these labels give more confidence for a particular part of the region of interest.

In the second step, we focus on tailored labels for those instances that are close one to each other (task 2), without a clear boundary separating them. Again, we use scribbles to mark their interface. Examples for such annotations are given in Figure 3 (top line, right part) The work can be carried out efficiently especially when the target scribbles have a sparse distribution. On the other hand, as no labelling is needed for the objects that are well separated, we can collect sufficient labelled images in a limited amount of time and cover the complex ice cream sample conditions. Lazy manual labeling of tasks 1 and 2 are done independently. It follows the assumption made in Section 3 that  $s^{(1)}$  and  $s^{(2)}$  are conditionally independent given image  $I$ .

The precise labels for task 3 are created using interactive segmentation tools. Starting from the rough (inner) regions of task 1, a natural idea is to let these regions grow and stop when the boundaries are reached. This can be done with geodesic active contours (Caselles et al., 1997). Unfortunately, such a method fails to capture sharp corners and the contour evolution tends to ignore boundaries with low

contrast. The annotation then requires frequent and time consuming user interaction. Instead, we use Grabcut (Rother et al., 2004; Hong et al., 2015) a graph-cut based method. The initial labels obtained from the first step give a good guess of the whole object regions. The Grabcut works well on isolated objects. However, it gives poor results when the objects are close to each other and have boundaries with inhomogeneous colors. As corrections may be needed for each image, only a few images of the whole dataset are processed. A fully segmented example is shown in the last row of Figure 3.

## 4 EXPERIMENTS

In this section, we demonstrate the performance of our approach using two different datasets. For both datasets we use strong labels (SL) as well as weak labels (WL). We prepare the labels and design the network as described in Section 3. The overall method is summarized with the 2 procedures presented in Algorithm 1 (see the appendix).

### 4.1 SEM IMAGES OF ICE CREAM

Scanning Electron Microscopy (SEM) constitutes the state-of-the-art for analysing food microstructures as it enables the efficient acquisition of high quality images for food materials, resulting into a huge amount of image data available for analysis. However, to better delineate the microstructures and provide exact statistical information, the partition of the images into different food components and instances is needed. The structures of food, especially soft solid materials, are usually complex which makes automated segmentation a difficult task. Some SEM images of ice cream in our dataset are shown on the top right of Figure 1. A typical ice cream sample consists of air bubbles, ice crystals and a concentrated unfrozen solution. In most situations, the air bubbles and ice crystals appear as foam in the images, while the solution fills the gaps between them. We treat the solution as the background and aim at detecting and computing a pixel-wise classification for each air bubbles and ice crystals instances.

The set of ice-cream SEM dataset consists of 38 images that are split into three sets (53% for training, 16% for validation and 31% testing respectively). Each image contains a rich set of instances with an overall number of instances around 13300 for 2 classes (ice crystals and air bubbles). For comparison, the PASCAL VOC 2012 dataset has 27450 objects in total for 20 classes.

For training the network, data augmentation is applied to prevent over-fitting. The size of the raw images is  $960 \times 1280$ . They are rescaled and rotated randomly, and then cropped into an input size of  $256 \times 256$  for feeding the network. Random flipping is also performed during training. The network is trained using Adam optimizer (Kingma & Ba, 2014) with a learning rate  $r = 2 \times 10^{-4}$  and a batch size of 16.

In the inference phase, the network outputs for each patch a probability map of size  $256 \times 256$ . The patches are then aggregated to obtain a probability map for the whole image. In general, the pixels near the boundaries of each patch are harder to classify. We thus weight the spatial influence of the patches with a Gaussian kernel to emphasize the network prediction at patch center.

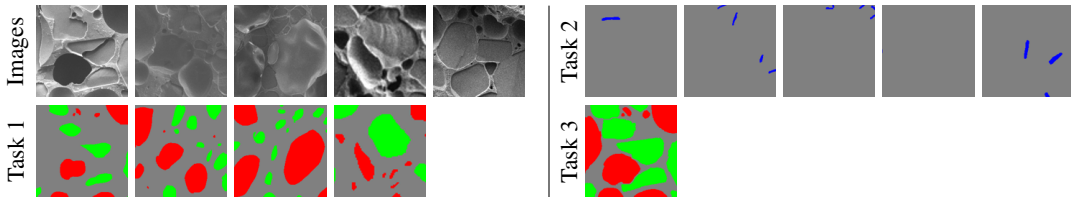


Figure 3: Example of annotated images. Some of the annotations are missing because not all images are labelled for task 1 and task 2. The marks in red are for air bubbles and the ones in green are for ice crystal instances. The blue curves on the third row are labels for interfaces of touching objects.

We now evaluate the multi-task U-net and compare it to the traditional single task U-net. The performance of each model is tested on 12 images, and average results are shown in Table 1. In the table, the dice score for a class  $c$  is defined as  $d_c = 2 \sum_i x_{i,c} y_{i,c} / (\sum_i x_{i,c} + \sum_i y_{i,c})$  where  $x$  is the computed segmentation mask and  $y$  the ground truth.

The models	air bubbles	ice crystals	Overall
U-net on WL	0.725	0.706	0.716
U-net on SL	0.837	0.794	0.818
PL approach	0.938	0.909	0.924
Multi-task U-net	<b>0.953</b>	<b>0.931</b>	<b>0.944</b>

Table 1: Dice scores of segmentation results on the test images of SEM images of ice cream dataset.

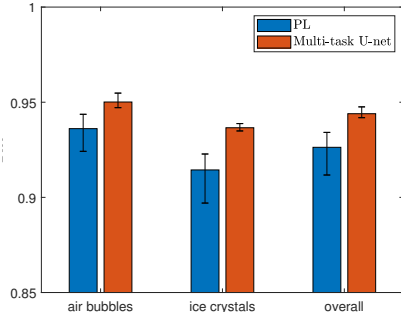


Figure 4: The error bars for the PL and multi-task U-net, each computed from 8 different experiments

We train a single task U-net (*i.e.*, without the multi-task block) on the weakly labelled set (task 1), with the 15 annotated images. The single task U-net on weak annotations gives an overall dice score at 0.72, the lowest one among the three other methods tested. One reason for the low accuracy of the single task U-net on weak (inaccurate) annotations is that in the training labels, the object boundaries are mostly ignored. Hence the U-net is not trained to recover them, leaving large parts of the object not recognized. Second, we consider strong annotations as training data, without the data of the other tasks, *i.e.* only 2 images with accurate segmentation masks are used. The score of the U-net trained on SL is only 0.82, which is significantly lower than the 0.94 obtained by the multi-task U-net.

We also compare our multi-task U-net results with one of the major weakly supervised approaches that makes use of pseudo labels (PL) (see e.g., Khoreva et al. (2017); Jing et al. (2018)). In these approaches, the pseudo segmentation masks are created from Ws and are used to feed a segmentation network. Following the work of Khoreva et al. (2017), we use the Grabcut method to create the PLs, and in our setting the partial masks are used to initialize the Grabcut. For the small subset of images that are strongly annotated, the full segmentation masks are used instead of PLs. The PL are created without human correction, and then used for feeding the segmentation network. Here we use the baseline single task U-net as the segmentation network in order to make comparison with the multi-task U-net.

Our multi-task U-net outperforms the PL approach as shown in Table 1. Figure 4 displays the error bars for the two methods with dice score collected from 8 different runs. A significant limitation of the PL method here is that its performance relies on the tools used for pseudo segmentation mask generations. Having common errors, instead of random ones, on the object boundaries in the training data, the segmentation network of PL also learns to have those patterns in the prediction. The images in the left part of Figure 5 show that the predicted label of an object tends to merge with some background pixels when there are edges of another object nearby. Similar errors from the GrabCut are illustrated in the appendix.

Besides the number of pixels that are correctly classified, the separation of touching instances is also of interest. In addition to the dice scores in Table 1, we study the learning performance of multi-task U-net on task 2, which specializes in the separation aspect. The test results on the 12 images give an overall precision of 0.70 of the detected interfaces, while 0.82 of the touching objects are recognized. We show some examples of computed separations and ground truth in the right part of Figure 5. For the detection task, the network predicts a probability map for the inner regions of the object instances. An output of the network is shown in Figure 6. With partial masks as coarse labels for this task, the network learns to identify the object instance as shown in the figure.

We finally consider the work of Bearman et al. (2016) that investigate the cost related to different types of annotations. Based on the data reported (Bearman et al., 2016) and our estimated annotation time, the collecting of WL for detection is considerably (more than 6x) faster than obtaining SL. For a fair comparison with the baseline U-net we use a larger ratio of SL for the single task learning accordingly (since no WL is used here), and the results are reported in Table 2. The proposed method still outperforms the U-net by a large margin on similar annotation time budgets, and additional SL after the first 10% do not help significantly.

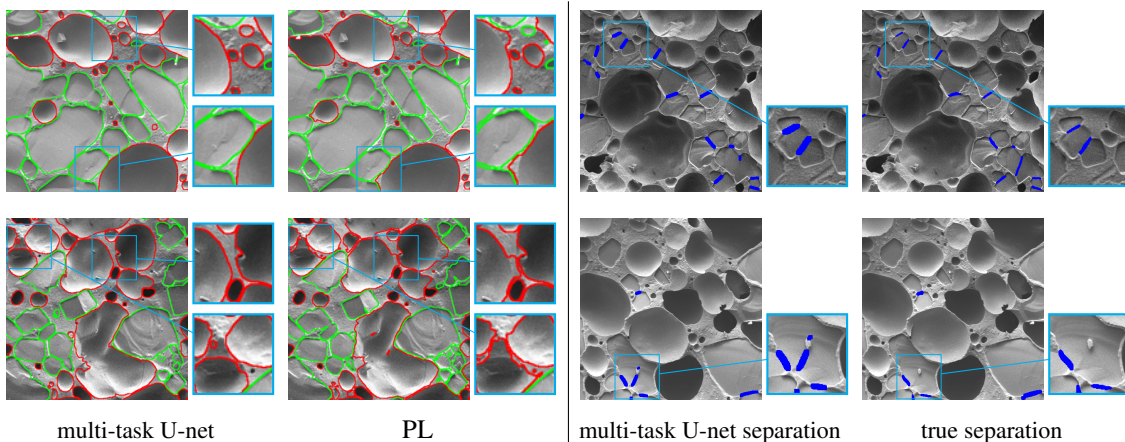


Figure 5: Segmentation and separation results (best viewed in color). Left: the computed contours are shown in red for air bubbles and green for ice crystals. While multi-task U-net and PL supervised network both have good performance, PL misclassifies the background near object boundaries. Right: Examples of separation by the multi-task U-net and the ground truth.

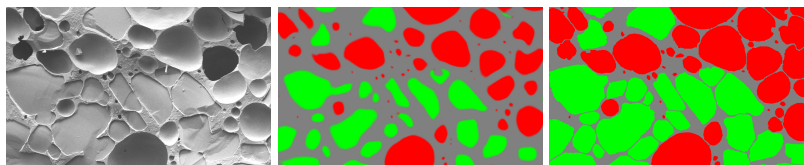


Figure 6: The image (left), the predicted probability map (middle) from the detection task, and the ground truth segmentation mask (right). The red and green on the middle and right images stand for air bubble and ice crystals respectively.)

#### 4.2 H&E-STAINED IMAGE DATASET FOR GLAND SEGMENTATION

We also apply the approach to the segmentation of tissue in histology images. In this experiment, we use the GlaS challenge dataset (Sirinukunwattana et al., 2017) that consists of 165 H&E stained images. The dataset is split into three parts, with 85 images for training, and 60 for offsite test and 20 images for onsite test (we will call the latter two sets Test part A and Test part B respectively in the following).

Apart from the SL available from the dataset, we create a set of a weak labels for the detection task and separation task (as illustrated in Appendix, Figure 11). These weak labels together with a part

	methods	labels	dice score
Annotation budget 1	multi-task	10% SL + WL	<b>0.944</b>
	single task	20% SL	0.882
Annotation budget 2	multi-task	20% SL + WL	<b>0.948</b>
	single task	30% SL	0.913
Annotation budget 3	multi-task	50% SL + WL	<b>0.949</b>
	single task	75% SL	0.940

Table 2: We compare the two methods under similar annotation time budgets. In each budget, two different combinations of SL and WL that take similar annotation time are used. The WL in this table contains 75% labels for the detection task and 100% labels for the separation task. From budget 1 to budget 3, we increase the amount of labels (that means more annotation time is needed) in the training data, and the dice score is reported for each case.

SL Ratio		2.4%	4.7%	9.4%	100%
Test Part A	Ours	<b>0.866</b>	<b>0.889</b>	<b>0.915</b>	<b>0.921</b>
	U-net	0.700	0.749	0.840	0.915
	MDU-net				0.920
Test Part B	Ours	<b>0.751</b>	<b>0.872</b>	<b>0.904</b>	<b>0.910</b>
	U-net	0.658	0.766	0.824	0.898
	MDU-net				0.871

Table 3: Average dice score for segmentation of gland. Our method uses both SL and WL. The ratio of strong labels (SL) is increased from 2.4% to 100%, and the scores of the methods are reported here for two parts A and B of the test sets, as split in Sirinukunwattana et al. (2017).

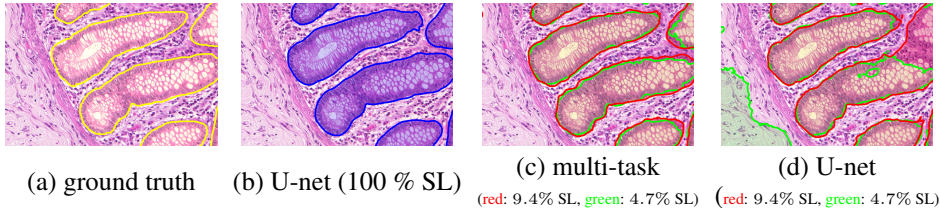


Figure 7: Segmentation results on the gland dataset (best view in color). The ground truth and the results. For (c) and (d), **Red** contour denotes the results from 9.4% strong labels; **Green** contour denotes results from 4.7% strong labels

of the strong labels are used for training the multi-task U-net as illustrated in Algorithm 1 (in the appendix).

In this experiment, we test the algorithm on different ratios of SL, and compare it with the baseline U-net, and the Multi-scale Densely Connected U-Net (DMU-net) (Zhang et al., 2018). The results on two sets of test data are reported in Table 3. As the SL ratios increase from 2.4% to 9.4%, an improvement of performance of the multi-task U-net is gained, and when it reaches 9.4% SL, the multi-task framework achieves comparable score with the fully supervised version. We emphasize that the 9.4% SL and WL are at a much lower annotation cost than that of the 100% SL used for fully supervised learning. Example of segmentation results are displayed in Figure 7.

## 5 CONCLUSION

In this paper, we develop a multi-task learning framework for image segmentation problems, which relaxes the requirement for numerous and accurate annotations to train the network. It is suitable for segmentation problem with a dense population of object instances. The model separates the segmentation problem into three smaller tasks. One of them is dedicated to the instance detection and therefore do not need exact boundary information. This gives potential flexibility as one could concentrate on the classification and rough location of the instances during data collection. The second one focuses on the separation of objects sharing a common boundary. The final task aims at extracting pixel-wise boundary information. Thanks to the information shared within the multi-task learning, this accurate segmentation can be obtained using very few annotated data.

Our model learns directly the statistics of WL as auxiliary tasks, so no further processing steps are needed before training the network. For the partial masks that ignore boundary pixels, the annotation can also be done when the boundaries of object are hard to detect. As a small amount of SL is needed and the collection of WL can be fast and cheap, the proposed framework is potentially effective for applications with growing datasets. The weakly annotated set for detection purpose could be augmented if necessary and the new images could easily be incorporated into our end-to-end framework.



## ACKNOWLEDGMENTS

## REFERENCES

- Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017.
- Christoph Baur, Shadi Albarqouni, and Nassir Navab. Semi-supervised deep learning for fully convolutional networks. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 311–319. Springer, 2017.
- Amy Bearman, Olga Russakovsky, Vittorio Ferrari, and Li Fei-Fei. What’s the point: Semantic segmentation with point supervision. In *European conference on computer vision*, pp. 549–565. Springer, 2016.
- Benjamin Bischke, Patrick Helber, Joachim Folz, Damian Borth, and Andreas Dengel. Multi-task learning for segmentation of building footprints with deep neural networks. *arXiv preprint arXiv:1709.05932*, 2017.
- Vicent Caselles, Ron Kimmel, and Guillermo Sapiro. Geodesic active contours. *International journal of computer vision*, 22(1):61–79, 1997.
- Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*, 2014.
- Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2018.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.
- Matvey Ezhov, Adel Zakirov, and Maxim Gusarev. Coarse-to-fine volumetric segmentation of teeth in cone-beam ct. *arXiv preprint arXiv:1810.10293*, 2018.
- Arthita Ghosh, Max Ehrlich, Sohil Shah, Larry Davis, and Rama Chellappa. Stacked u-nets for ground material segmentation in remote sensing imagery. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 257–261, 2018.
- Fidel A Guerrero-Pena, Pedro D Marrero Fernandez, Tsang Ing Ren, Mary Yui, Ellen Rothenberg, and Alexandre Cunha. Multiclass weighted loss for instance segmentation of cluttered cells. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pp. 2451–2455. IEEE, 2018.
- Seunghoon Hong, Hyeonwoo Noh, and Bohyung Han. Decoupled deep neural network for semi-supervised semantic segmentation. In *Advances in neural information processing systems*, pp. 1495–1503, 2015.
- Zilong Huang, Xinggang Wang, Jiasi Wang, Wenyu Liu, and Jingdong Wang. Weakly-supervised semantic segmentation network with deep seeded region growing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7014–7023, 2018.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

- Longlong Jing, Yucheng Chen, and Yingli Tian. Coarse-to-fine semantic segmentation from image-level labels. *arXiv preprint arXiv:1812.10885*, 2018.
- Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7482–7491, 2018.
- Anna Khoreva, Rodrigo Benenson, Jan Hosang, Matthias Hein, and Bernt Schiele. Simple does it: Weakly supervised instance and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 876–885, 2017.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Alexander Kolesnikov and Christoph H Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *European Conference on Computer Vision*, pp. 695–711. Springer, 2016.
- Jungbeom Lee, Eunji Kim, Sungmin Lee, Jangho Lee, and Sungroh Yoon. Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. *arXiv preprint arXiv:1902.10421*, 2019.
- Di Lin, Jifeng Dai, Jiaya Jia, Kaiming He, and Jian Sun. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3159–3167, 2016.
- Geert Litjens, Thijs Kooi, Babak Ehteshami Bejnordi, Arnaud Arindra Adiyoso Setio, Francesco Ciompi, Mohsen Ghafoorian, Jeroen Awm Van Der Laak, Bram Van Ginneken, and Clara I Sánchez. A survey on deep learning in medical image analysis. *Medical image analysis*, 42:60–88, 2017.
- Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3431–3440, 2015.
- Pawel Mlynarski, Hervé Delingette, Antonio Criminisi, and Nicholas Ayache. Deep learning with mixed supervision for brain tumor segmentation. *arXiv preprint arXiv:1812.04571*, 2018.
- George Papandreou, Liang-Chieh Chen, Kevin P Murphy, and Alan L Yuille. Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In *Proceedings of the IEEE international conference on computer vision*, pp. 1742–1750, 2015.
- Christian S Perone and Julien Cohen-Adad. Deep semi-supervised segmentation with weight-averaged consistency targets. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, pp. 12–19. Springer, 2018.
- Clément Payout, Renaud Duval, and Farida Cheriet. A multitask learning architecture for simultaneous segmentation of bright and red lesions in fundus images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 101–108. Springer, 2018.
- Clément Payout, Renaud Duval, and Farida Cheriet. A novel weakly supervised multitask architecture for retinal lesions segmentation on fundus images. *IEEE transactions on medical imaging*, 2019.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
- Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. Grabcut: Interactive foreground extraction using iterated graph cuts. In *ACM transactions on graphics (TOG)*, volume 23, pp. 309–314. ACM, 2004.
- Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*, 2017.



- Meet P Shah, SN Merchant, and Suyash P Awate. Ms-net: Mixed-supervision fully-convolutional networks for full-resolution segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 379–387. Springer, 2018.
- Seung Yeon Shin, Soochahn Lee, Il Dong Yun, Sun Mi Kim, and Kyoung Mu Lee. Joint weakly and semi-supervised deep learning for localization and classification of masses in breast ultrasound images. *IEEE transactions on medical imaging*, 38(3):762–774, 2019.
- Korsuk Sirinukunwattana, Josien PW Pluim, Hao Chen, Xiaojuan Qi, Pheng-Ann Heng, Yun Bo Guo, Li Yang Wang, Bogdan J Matuszewski, Elia Bruni, Urko Sanchez, et al. Gland segmentation in colon histology images: The glas challenge contest. *Medical image analysis*, 35:489–502, 2017.
- Fengdong Sun and Wenhui Li. Saliency guided deep network for weakly-supervised image segmentation. *Pattern Recognition Letters*, 120:62–68, 2019.
- Tao Sun, Zehui Chen, Wenxiang Yang, and Yin Wang. Stacked u-nets with multi-output for road extraction. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 187–1874. IEEE, 2018.
- Satoshi Tsutsui, Tommi Kerola, Shunta Saito, and David J Crandall. Minimizing supervision for free-space segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 988–997, 2018.
- Chengjia Wang, Tom MacGillivray, Gillian Macnaught, Guang Yang, and David Newby. A two-stage 3d unet framework for multi-class segmentation on full resolution image. *arXiv preprint arXiv:1804.04341*, 2018a.
- Xinlong Wang, Tete Xiao, Yuning Jiang, Shuai Shao, Jian Sun, and Chunhua Shen. Repulsion loss: Detecting pedestrians in a crowd. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7774–7783, 2018b.
- Yunchao Wei, Xiaodan Liang, Yunpeng Chen, Xiaohui Shen, Ming-Ming Cheng, Jiashi Feng, Yao Zhao, and Shuicheng Yan. Stc: A simple to complex framework for weakly-supervised semantic segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(11):2314–2320, 2017.
- Yunchao Wei, Huaxin Xiao, Honghui Shi, Zequn Jie, Jiashi Feng, and Thomas S Huang. Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7268–7277, 2018.
- Jiawei Zhang, Yuzhen Jin, Jilan Xu, Xiaowei Xu, and Yanchun Zhang. Mdu-net: Multi-scale densely connected u-net for biomedical image segmentation. *arXiv preprint arXiv:1812.00352*, 2018.
- Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE international conference on computer vision*, pp. 1529–1537, 2015.
- Juan Zhou, Lu-Yang Luo, Qi Dou, Hao Chen, Cheng Chen, Gong-Jie Li, Ze-Fei Jiang, and Pheng-Ann Heng. Weakly supervised 3d deep learning for breast cancer classification and localization of the lesions in mr images. *Journal of Magnetic Resonance Imaging*, 2019.

## A APPENDIX

### A.1 THE ALGORITHM

We summarize the overall procedures, including the steps of creating the labels and the multi-task training into Algorithm .

**Algorithm 1** Multi-task learning for segmentation with lazy labels

- 
- 1: **procedure** LAZY LABELS( $\mathcal{I}$ ) ▷ Choose  $\mathcal{I}_1, \mathcal{I}_2, \mathcal{I}_3 \subset \mathcal{I}$ ,  $|\mathcal{I}_3|$  can be relatively small.  
**Input:** set of images  $\mathcal{I}$
  - 2:   Select inner regions for each object in  $\mathcal{I}_1$  for the detection task
  - 3:   Indicate scribbles on images of  $\mathcal{I}_2$  for the separation task
  - 4:   Generate a few pixel-wise labels  $\mathcal{I}_3$  from  $\mathcal{I}_1$  using interactive segmentation tools (e.g., Grabcut)
  - 5:   **return**  $\mathcal{I}_1, \mathcal{I}_2, \mathcal{I}_3$  and the labels
- 
- 1: **procedure** MULTI-TASK U-NET TRAINING( $\mathcal{I}_k, \mathbf{s}_i^{(k)}, \alpha_k, r$ ) ▷  $\mathbf{s}_i^{(k)}$  denote the labels  
**Input:** labelled sets  $\mathcal{I}_k$ , loss function weights  $\alpha_k$  for  $k = 1, 2, 3$ , Adam parameters  $r$ , mini-batch size  $m$ .
  - 2:   Set the 1<sup>st</sup> and 2<sup>nd</sup> momentum vectors  $\mathbf{m}, \mathbf{v}$  as zeros.
  - 3:   Initialize the multi-task U-net parameter  $\theta$ .
  - 4:   **while** Termination criterion is not satisfied **do**
  - 5:     Obtain a mini-batch  $(I_1, \mathbf{s}_1^{(k)}), \dots, (I_m, \mathbf{s}_m^{(k)})$ . ▷  $\mathbf{s}_i^{(k)}$  can be unknown for some  $(i, k)$ .
  - 6:     Compute the gradient  $\mathbf{g} \leftarrow \nabla_{\theta} \frac{1}{m} \sum_{k=1}^3 \sum_{i=1}^m \alpha_k \mathbb{1}_{\mathcal{I}_k}(I_i) \log p(\mathbf{s}_i^{(k)} | I_i; \theta)$
  - 7:      $(\theta, \mathbf{m}, \mathbf{v}) \leftarrow \text{Adam}_r(\theta, \mathbf{m}, \mathbf{v}, \mathbf{g})$  ▷  $\text{Adam}_r(\cdot)$  is an Adam iteration
  - 8:   **return**  $\theta$
- 

## A.2 FURTHER RESULTS FOR THE ICE CREAM IMAGE DATASET

Figure 8 shows the Dice scores on the 12 test images from Figure 9. One reason for the low accuracy of the single task U-net on weak (inaccurate) annotations is that the U-net is not trained to recover object boundaries. Therefore, the corresponding dice score on every test image is low, as shown by the blue curve. On the other hand, the network trained on strong annotations has a very good performance on only a few test images, as shown by the red line. Learning from both strong and weak annotations, the multi-task U-net improves from the other two, and the dice scores are between 0.9 to 0.95 for all of the images.

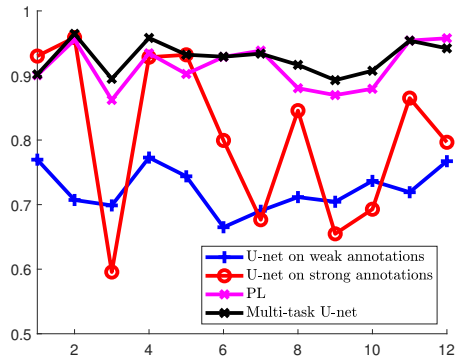


Figure 8: Overall dice scores for the air bubbles and ice crystals plotted versus the test image number. We show the scores for the U-net baseline with weak annotations, strong annotations and pseudo labels (PL) respectively, compared along with the multi-task U-net

## A.2.1 ADDITIONAL RESULTS

We compare the segmentation results given by the multi-task U-net with the ones obtained using GAC Caselles et al. (1997) and Grabcut Rother et al. (2004) (Figure 10). To apply GAC and Grabcut, we convert the problem into a single class and single object segmentation problem. We use the rough inner regions from the first step of the labelling to estimate a bounding box around each object. We then apply GAC and Grabcut inside each box. The box-wise results are aggregated to get the segmentation mask of the whole image.

Both GAC and Grabcut fail to capture the boundaries in various situations, such as intensity variations within the objects and background fragments connected to the objects. This demonstrates that a precise labelling requires many additional user interventions and thus the interest of our lazy labelling.

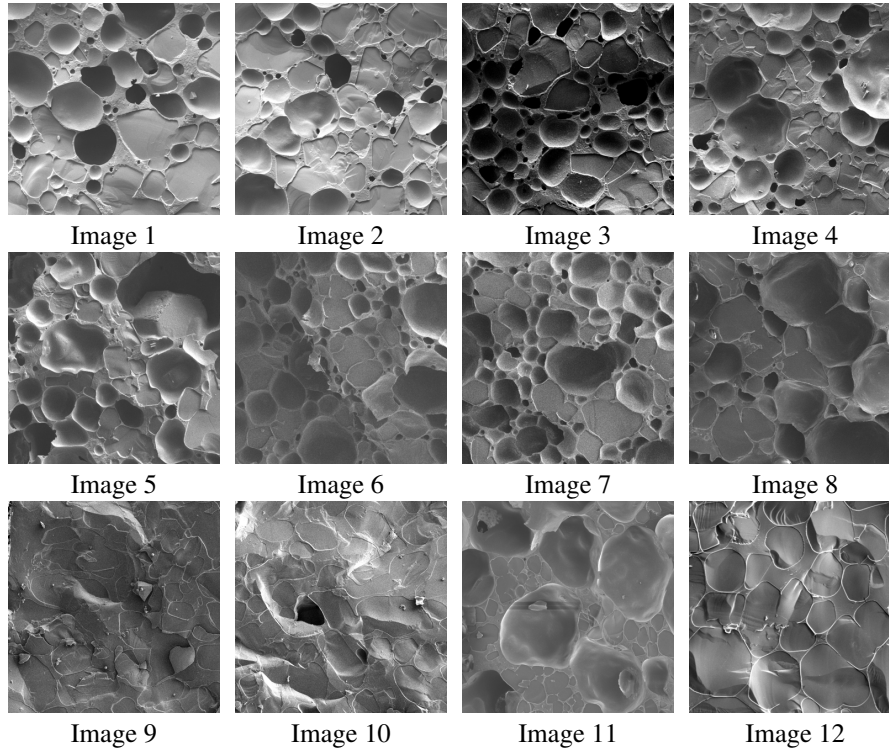


Figure 9: Twelve patches from the twelve test images respectively.

As shown in the third column of Figure 10), the Grabcut repeatedly misclassified small pieces of background as parts of the object especially when two objects are close to each other with some weak boundaries. The PL learning methods also give this kind of error in the inference phase. This means that the error has been learned by the segmentation network. The manual detection and correction of such small misclassified regions is a tedious task as well, and therefore is a limitation of the PL approach.

### A.3 THE H&E-STAINED IMAGE DATASET FOR GLAND SEGMENTATION

Figure 7 shows examples of WL and SL that are used for the segmentation of gland.

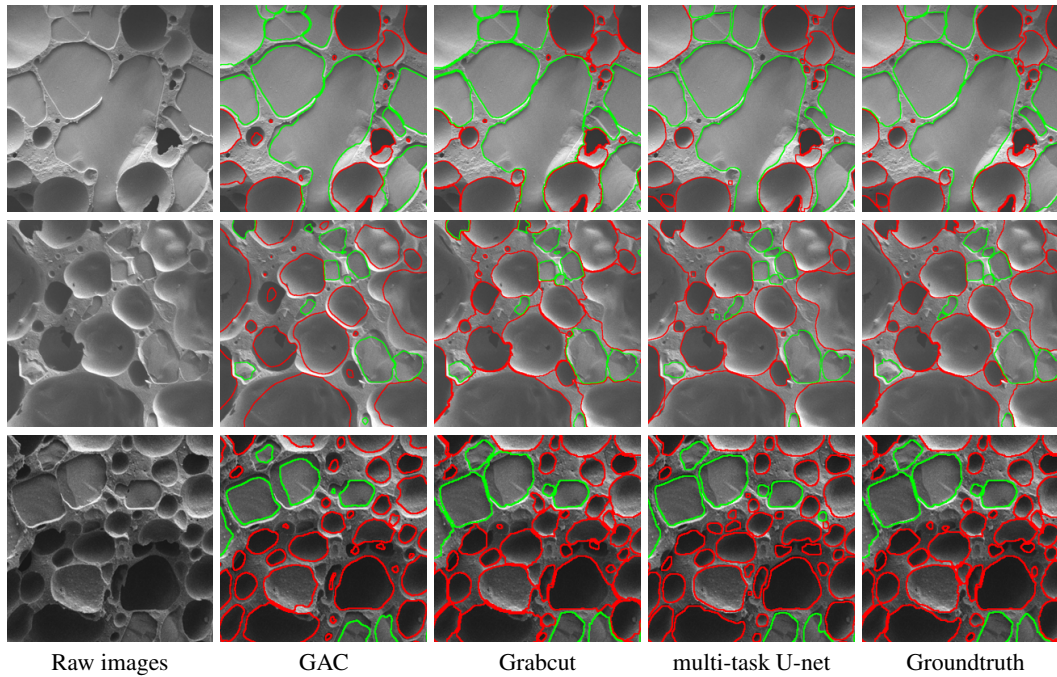


Figure 10: Comparison of multi-task U-net with the GAC and Grabcut methods using the weak labels of the validation data. Green curves: contour of the ice crystals. Red curves: air bubbles.

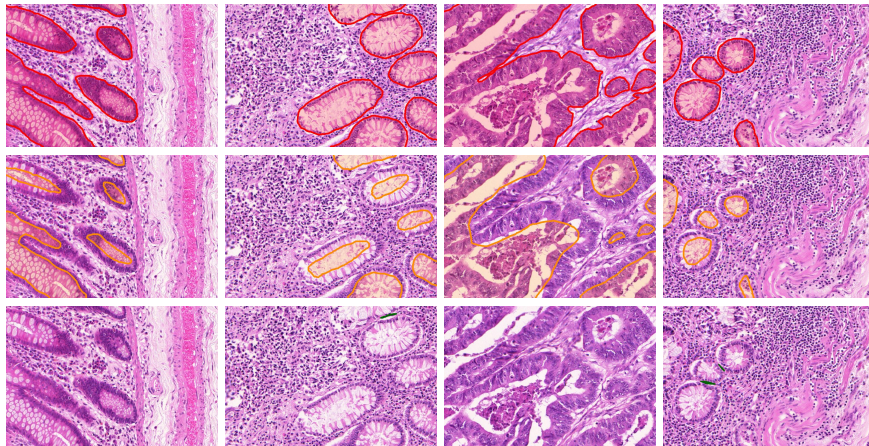


Figure 11: Gland segmentation dataset with SL in the first row (highlighted in red), WL in the second row for detection (in orange) and weak labels in the third row for separation (in dark green). The labels for separation can be sparse with some images (see the third row) having zero annotation as the instances are not touching each other.