



Co-heterogeneous and Adaptive Segmentation from Multi-source and Multi-phase CT Imaging Data: A Study on Pathological Liver and Lesion Segmentation

Ashwin Raju^{1,2}, Chi-Tung Cheng³, Yuankai Huo¹, Jinzheng Cai¹,
Junzhou Huang², Jing Xiao⁴, Le Lu¹, ChienHung Liao³,
and Adam P. Harrison¹(✉)

¹ PAII Inc., Bethesda, MD, USA

adampharrison070@piai-labs.com

² The University of Texas at Arlington, Arlington, TX, USA

³ Chang Gung Memorial Hospital, Linkou, Taiwan, ROC

⁴ PingAn Technology, Shenzhen, China

Abstract. Within medical imaging, organ/pathology segmentation models trained on current publicly available and fully-annotated datasets usually do not well-represent the heterogeneous modalities, phases, pathologies, and clinical scenarios encountered in real environments. On the other hand, there are tremendous amounts of unlabelled patient imaging scans stored by many modern clinical centers. In this work, we present a novel segmentation strategy, co-heterogeneous and adaptive segmentation (CHASE), which only requires a small labeled cohort of *single* phase data to adapt to any unlabeled cohort of heterogeneous *multi-phase* data with possibly new clinical scenarios and pathologies. To do this, we develop a versatile framework that fuses appearance-based semi-supervision, mask-based adversarial domain adaptation, and pseudo-labeling. We also introduce co-heterogeneous training, which is a novel integration of co-training and hetero-modality learning. We evaluate CHASE using a clinically comprehensive and challenging dataset of multi-phase computed tomography (CT) imaging studies (1147 patients and 4577 3D volumes), with a test set of 100 patients. Compared to previous state-of-the-art baselines, CHASE can further improve pathological liver mask Dice-Sørensen coefficients by ranges of 4.2% to 9.4%, depending on the phase combinations, *e.g.*, from 84.6% to 94.0% on non-contrast CTs.

Keywords: Multi-phase segmentation · Semi-supervised learning · Co-training · Domain adaptation · Liver and lesion segmentation

Electronic supplementary material The online version of this chapter (https://doi.org/10.1007/978-3-030-58592-1_27) contains supplementary material, which is available to authorized users.

1 Introduction

Delineating anatomical structures is an important task within medical imaging, *e.g.*, to generate biomarkers, quantify or track disease progression, or to plan radiation therapy. Manual delineation is prohibitively expensive, which has led to a considerable body of work on automatic segmentation. However, a perennial problem is that models trained on available image/mask pairs, *e.g.*, publicly available data, do not always reflect clinical conditions upon deployment, *e.g.*, the present pathologies, patient characteristics, scanners, and imaging protocols. This leads to potentially drastic performance gaps. When multi-modal or multi-phase imagery is present these challenges are further compounded, as datasets may differ in their composition of available modalities or consist of heterogeneous combinations of modalities. The challenges then are in both managing new patient/disease variations and in handling heterogeneous multi-modal data. Ideally segmentation models can be deployed without first annotating large swathes of additional data matching deployment scenarios. This is our goal, where we introduce co-heterogeneous and adaptive segmentation (CHASe), which can adapt models trained on single-modal and public data to produce state-of-the-art results on *multi-phase and multi-source* clinical data *with no extra annotation cost*.

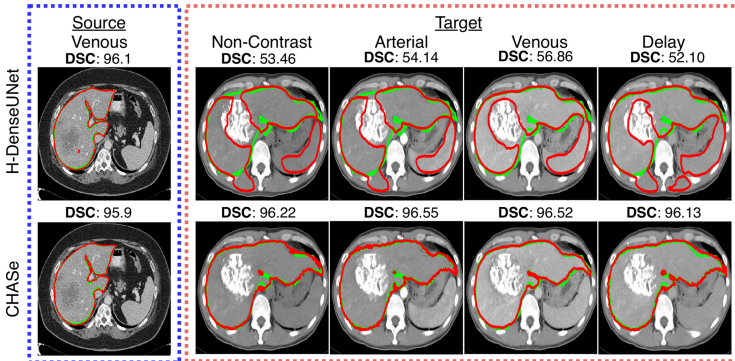


Fig. 1. Ground truth and predictions are rendered in green and red, respectively. Despite performing excellently on labeled source data, state-of-the-art fully-supervised models [25] can struggle on cohorts of multi-phase data with novel conditions, *e.g.*, the patient shown here with splenomegaly and a TACE-treated tumor. CHASe can adapt such models to perform on new data. (Color figure online)

Our motivation stems from challenges in handling the wildly variable data found in picture archiving and communication systems (PACSs), which closely follows deployment scenarios. In this study, we focus on segmenting pathological livers and lesions from dynamic computed tomography (CT). Liver disease and cancer are major morbidities, driving efforts toward better detection and

characterization methods. The dynamic contrast CT protocol images a patient under multiple time-points after a contrast agent is injected, which is critical for characterizing liver lesions [30]. Because accurate segmentation produces important volumetric biomarkers [3, 12], there is a rich body of work on automatic segmentation [10, 20, 23, 25, 28, 36, 41, 46, 48], particularly for CT. Despite this, all publicly available data [3, 7, 11, 18, 38] is limited to venous-phase (single-channel) CTs. Moreover, when lesions are present, they are typically limited to hepatocellular carcinoma (HCC) or metastasized tumors, lacking representation of intrahepatic cholangiocellular carcinoma (ICC) or the large bevy of benign lesion types. Additionally, public data may not represent other important scenarios, *e.g.*, the transarterial chemoembolization (TACE) of lesions or splenomegaly, which produce highly distinct imaging patterns. As Fig. 1 illustrates, even impressive leading entries [25] within the public LiTS challenge [3], can struggle on clinical PACS data, particularly when applied to non-venous contrast phases.

To meet this challenge, we integrate together powerful, but complementary, strategies: hetero-modality learning, appearance-based consistency constraints, mask-based adversarial domain adaptation (ADA), and pseudo-labeling. The result is a semi-supervised model trained on smaller-scale supervised public *venous-phase* data [3, 7, 11, 18, 38] and large-scale unsupervised *multi-phase* data. Crucially, we articulate non-obvious innovations that avoid serious problems from a naive integration. A key component is co-training [4], but unlike recent deep approaches [44, 49], we do not need artificial views, instead treating each phase as a view. We show how co-training can be adopted with a minimal increase of parameters. Second, since CT studies from clinical datasets may exhibit any combination of phases, ideally liver segmentation should also be able to accept whatever combination is available, with performance topping out as more phases are available. To accomplish this, we fuse hetero-modality learning [15] together with co-training, calling this *co-heterogeneous training*. Apart from creating a natural hetero-phase model, this has the added advantage of combinatorially increasing the number of views for co-training from 4 to 15, boosting even single-phase performance. To complement these appearance-based semi-supervision strategies, we also apply pixel-wise ADA [40], guiding the network to predict masks that follow a proper shape distribution. Importantly, we show how ADA can be applied to co-heterogeneous training with no extra computational cost over adapting a single phase. Finally, we address edge cases using a principled pseudo-labelling technique specific to pathological organ segmentation. These innovations combine to produce a powerful approach we call CHASe.

We apply CHASe to a challenging unlabelled dataset of 1147 dynamic-contrast CT studies of patients, all with liver lesions. The dataset, extracted directly from a hospital PACS, exhibits many features not seen in public single-phase data and consists of a heterogeneous mixture of non-contrast (NC), arterial (A), venous (V), and delay (D) phase CTs. With a test set of 100 studies, *this is the largest, and arguably the most challenging, evaluation of multi-phase pathological liver segmentation to date*. Compared to strong fully-supervised

baselines [14,25] only trained on public data, CHASE can dramatically boost segmentation performance on non-venous phases, *e.g.*, moving the mean Dice-Sørensen coefficient (DSC) from 84.5 to 94.0 on NC CTs. Importantly, performance is also boosted on V phases, *i.e.*, from 90.7 mean DSC to 94.9. Inputting all available phases to CHASE maximizes performance, matching desired behavior. Importantly, CHASE also significantly improves robustness, operating with much greater reliability and without deleterious outliers. Since CHASE is general-purpose, it can be applied to other datasets or even other organs with ease.

2 Related Work

Liver and Liver Lesion Segmentation. In the past decade, several works addressed liver and lesion segmentation using traditional texture and statistical features [10,23,41]. With the advent of deep learning, fully convolutional networks (FCNs) have quickly become dominant. These include 2D [2,36], 2.5D [13,43], 3D [20,28,45,48], and hybrid [25,46] FCN-like architectures. Some reported results show that 3D models can improve over 2D ones, but these improvements are sometimes marginal [20].

Like related works, we also use FCNs. However, all prior works are trained and evaluated on venous-phase CTs in a fully-supervised manner. In contrast, we aim to robustly segment a large cohort of multi-phase CT PACS data in a semi-supervised manner. As such, our work is orthogonal to much of the state-of-the-art, and can, in principle, incorporate any future fully-supervised solution as a starting point.

Semi-supervised Learning. Annotating medical volumetric data is time consuming, spurring research on semi-supervised solutions [39]. In *co-training*, predictions of different “views” of the same unlabelled sample are enforced to be consistent [4]. Recent works integrate co-training to deep-learning [31,44,49]. While CHASE is related to these works, it uses different contrast phases as views and therefore has no need for artificial view creation [31,44,49]. More importantly, CHASE effects a stronger appearance-based semi-supervision by fusing co-training with hetero-modality learning (co-heterogenous learning). In addition, CHASE complements this appearance-based strategy via prediction-based ADA, resulting in significantly increased performance.

Adversarial domain adaptation (ADA) for semantic segmentation has also received recent attention in medical imaging [39]. The main idea is to align the distribution of predictions or features between the source and target domains, with many successful approaches [6,9,40,42]. Like Tsai *et al.*, we align the distribution of *mask shapes* of source and target predictions [40]. Unlike Tsai *et al.*, we use ADA in conjunction with appearance-based semi-supervision. In doing so, we show how ADA can effectively adapt 15 different hetero-phase predictions at the same computational cost as a single-view variant. Moreover, we demonstrate the need for complementary semi-supervised strategies in order to create a robust and practical medical segmentation system.

In *self-learning*, a network is first trained with labelled data. The trained model then produces pseudo labels for the unlabelled data, which is then added to the labelled dataset via some scheme [24]. This approach has seen success within medical imaging [1, 29, 47], but it is important to guard against “confirmation bias” [26], which can compound initial misclassifications [39]. While we thus avoid this approach, we do show later that co-training can also be cast as self-learning with consensus-based pseudo-labels. Finally, like more standard self-learning, we also generate pseudo-labels to finetune our model. But, these are designed to deduce and correct likely mistakes, so they do not follow the regular “confirmation”-based self-learning framework.

Hetero-Modality Learning. In medical image acquisition, multiple phases or modalities are common, *e.g.*, dynamic CT. It is also common to encounter hetero-modality data, *i.e.*, data with possible missing modalities [15]. Ideally a segmentation model can use whatever phases/modalities are present in the study, with performance improving the more phases are available. CHASe uses hetero-modal fusion for fully-supervised FCNs [15]; however it fuses it with co-training, thereby using multi-modal learning to perform appearance-based learning from *unlabelled* data. Additionally, this *co-heterogeneous training* combinatorially increases the number of views for co-training, significantly boosting even single-phase performance by augmenting the training data. To the best of our knowledge, we are the first to propose co-heterogenous training.

3 Method

Although CHASe is not specific to any organ, we will assume the liver for this work. We assume we are given a curated and labelled dataset of CTs and masks, *e.g.*, from public data sources. We denote this dataset $\mathcal{D}_\ell = \{\mathcal{X}_i, Y_i\}_{i=1}^{N_\ell}$, with \mathcal{X}_i denoting the set of available phases and $Y_i(k) \in \{0, 1, 2\}$ indicating background, liver, and lesion for all pixel/voxel indices k . Here, without loss of generality, we assume the CTs are all V-phase, *i.e.*, $\mathcal{X}_i = V_i \forall \mathcal{X}_i \in \mathcal{D}_\ell$. We also assume we are given a large cohort of unlabelled multi-phase CTs from a challenging and uncurated clinical source, *e.g.*, a PACS. We denote this dataset $\mathcal{D}_u = \{\mathcal{X}_i\}_{i=1}^{N_u}$, where $\mathcal{X}_i = \{NC_i, A_i, V_i, D_i\}$ for instances with all contrast phases. When appropriate, we drop the i for simplicity. Our goal is to learn a segmentation model which can accept any combination of phases from the target domain and robustly delineate liver or liver lesions, despite possible differences in morbidities, patient characteristics, and contrast phases between the two datasets.

Figure 2 CHASe, which integrates supervised learning, co-heterogenous training, ADA, and specialized pseudo-labelling. We first start by training a standard fully-supervised segmentation model using the labelled data. Then under CHASe, we finetune the model using consistency and ADA losses:

$$\mathcal{L} = \mathcal{L}_{seg} + \mathcal{L}_{cons} + \lambda_{adv} \mathcal{L}_{adv}, \quad (1)$$

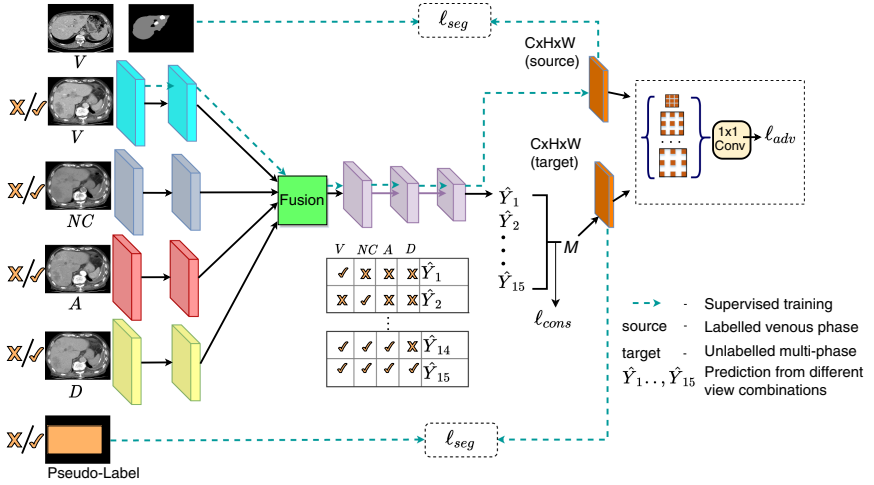


Fig. 2. Overview of CHASe. As shown at the top, labelled data, V-phase in our experiments, are trained using standard segmentation losses. Multi-phase unlabelled data are inputted into an efficient co-heterogenous pipeline, that minimizes divergence between mask predictions across all 15 phase combinations. ADA and specialized pseudo-labelling are also applied. Not shown are the deeply-supervised intermediate outputs of the PHNN backbone.

where \mathcal{L} , \mathcal{L}_{seg} , \mathcal{L}_{cons} , \mathcal{L}_{adv} are the overall, supervised, co-heterogenous, and ADA losses, respectively. For adversarial optimization, a discriminator loss, \mathcal{L}_d , is also deployed in competition with (1). We elaborate on each of the above losses below. Throughout, we minimize hyper-parameters by using standard architectures and little loss-balancing.

3.1 Backbone

CHASe relies on an FCN backbone, $f(\cdot)$, for its functionality and can accommodate any leading choice, including encoder/decoder backbones [28, 32]. Instead, we adopt the deeply-supervised progressive holistically nested network (PHNN) framework [14], which has demonstrated leading segmentation performance for many anatomical structures [5, 14, 21, 22, 34], sometimes even outperforming U-Net [21, 22]. Importantly, PHNN has roughly half the parameters and activation maps of an equivalent encoder/decoder. Since we aim to include additional components for semi-supervised learning, this lightweightness is important.

In brief, PHNN relies on deep supervision in lieu of a decoder and assumes the FCN can be broken into stages based on pooling layers. Here, we assume there are five FCN stages, which matches popular FCN configurations. PHNN produces a sequence of logits, $a^{(m)}$, using 1×1 convolutions and upsamplings operating on the terminal backbone activations of each stage. Sharing similarities to residual connections [16], predictions are generated for each stage using a progressive scheme that adds to the previous stage’s activations:

$$\hat{Y}^{(1)} = \sigma(a^{(1)}), \quad (2)$$

$$\hat{Y}^{(m)} = \sigma(a^{(m)} + a^{(m-1)}) \quad \forall m > 1, \quad (3)$$

where $\sigma(\cdot)$ denotes the softmax and $\hat{Y}^{(\cdot)}$ represents the predictions, with the final stage's predictions acting as the actual segmentation output, \hat{Y} . Being deeply supervised, PHNN optimizes a loss at each stage, with higher weights for later stages:

$$\ell_{seg}(f(V), Y) = \sum_{j=1}^5 \frac{m}{5} \ell_{ce}(\hat{Y}^m, Y), \quad (4)$$

where we use pixel-wise cross-entropy loss (weighted by prevalence), $\ell_{ce}(\cdot, \cdot)$. Prior to any semi-supervised learning, this backbone is pre-trained using \mathcal{D}_ℓ :

$$\mathcal{L}_{seg} = \frac{1}{N_\ell} \sum_{V, Y \in \mathcal{D}_\ell} \ell_{seg}(f(V), Y). \quad (5)$$

Our supplementary visually depicts the PHNN structure.

3.2 Co-training

With a pretrained fully-supervised backbone, the task is now to leverage the unlabelled cohort of dynamic CT data, \mathcal{D}_u . We employ the ubiquitous strategy of enforcing consistency. Because dynamic CT consists of the four NC, A, V, and D phases, each of which is matched to same mask, Y , they can be regarded as different views of the same data. This provides for a natural co-training setup [4] to penalize inconsistencies of the mask predictions across different phases.

To do this, we must create predictions for each phase. As Fig. 2 illustrates, we accomplish this by using phase-specific FCN stages, *i.e.*, the first two low-level stages, and then use a shared set of weights for the later semantic stages. Because convolutional weights are more numerous at later stages, this allows for an efficient multi-phase setup. All layer weights are initialized using the fully-supervised V-phase weights from Sect. 3.1, including the phase-specific layers. Note that activations across phases remain distinct. Despite the distinct activations, for convenience we abuse notation and use $\hat{Y} = f(\mathcal{X})$ to denote the generation of all phase predictions for one data instance. When all four phases are available in \mathcal{X} , then \hat{Y} corresponds to $\{\hat{Y}^{NC}, \hat{Y}^A, \hat{Y}^V, \hat{Y}^D\}$.

Like Qiao *et al.* [31] we use the Jensen-Shannon divergence (JSD) [27] to penalize inconsistencies. However, because it will be useful later we devise the JSD by first deriving a consensus prediction:

$$M = \frac{1}{|\hat{\mathcal{Y}}|} \sum_{\hat{Y} \in \hat{\mathcal{Y}}} \hat{Y}. \quad (6)$$

The JSD is then the divergence between the consensus and each phase prediction:

$$\ell_{cons}(f(\mathcal{X})) = \frac{1}{|\hat{\mathcal{Y}}|} \sum_{\hat{Y} \in \hat{\mathcal{Y}}} \sum_{k \in \Omega} KL(\hat{Y}(k) \parallel M(k)), \quad (7)$$

$$\mathcal{L}_{cons} = \frac{1}{N_u} \sum_{\mathcal{X} \in \mathcal{D}_u} \ell_{cons}(f(\mathcal{X})), \quad (8)$$

where Ω denotes the spatial domain and $KL(\cdot \parallel \cdot)$ is the Kullback-Leibler divergence across label classes. (8) thus casts co-training as a form of self-learning, where pseudo-labels correspond to the consensus prediction in (6). For the deeply-supervised PHNN, we only calculate the JSD across the final prediction.

3.3 Co-heterogeneous Training

While the loss in (8) effectively incorporates multiple phases of the unlabelled data, it is not completely satisfactory. Namely, each phase must still be inputted separately into the network, and there is no guarantee of a consistent output. Despite only having single-phase *labelled* data, ideally, the network should be adapted for multi-phase operation on \mathcal{D}_u , meaning it should be able to consume whatever contrast phases are available and output a unified prediction that is stronger as more phases are available.

To do this, we use hetero-modality image segmentation (HeMIS)-style feature fusion [15], which can predict masks given any arbitrary combination of contrast phases. To do this, a set of phase-specific layers produce a set of phase-specific activations, \mathcal{A} , whose cardinality depends on the number of inputs. The activations are then fused together using first- and second-order statistics, which are flexible enough to handle any number of inputs:

$$\mathbf{a}_{fuse} = \text{concat}(\mu(\mathcal{A}), \text{var}(\mathcal{A})), \quad (9)$$

where \mathbf{a}_{fuse} denotes the fused feature, and the mean and variance are taken across the available phases. When only one phase is available, the variance features are set to 0. To fuse intermediate predictions, an additional necessity for deeply-supervised networks, we simply take their mean.

For choosing a fusion point, the co-training setup of Sect. 3.2, with its phase-specific layers, already offers a natural fusion location. We can then readily combine hetero-phase learning with co-training, re-defining a “view” to mean any possible combination of the four contrast phases. This has the added benefit of combinatorially increasing the number of co-training views. More formally, we use $\mathcal{X}^* = \mathcal{P}(\mathcal{X}) \setminus \{\emptyset\}$ to denote all possible contrast-phase combinations, where $\mathcal{P}(\cdot)$ is the powerset operator. The corresponding predictions we denote as $\hat{\mathcal{Y}}^*$. When a data instance has all four phases, then the cardinality of \mathcal{X}^* and $\hat{\mathcal{Y}}^*$ is 15, which is a drastic increase in views. With hetero-modality fusion in place, the consensus prediction and co-training loss of (6) and (7), respectively, can be supplanted by ones that use $\hat{\mathcal{Y}}^*$:

$$M = \frac{1}{|\hat{\mathcal{Y}}^*|} \sum_{\hat{Y} \in \hat{\mathcal{Y}}^*} \hat{Y}, \quad (10)$$

$$\ell_{cons}(f(\mathcal{X})) = \frac{1}{|\hat{\mathcal{Y}}^*|} \sum_{\hat{Y} \in \hat{\mathcal{Y}}^*} \sum_{k \in \Omega} KL(\hat{Y}(k) \parallel M(k)). \quad (11)$$

When only single-phase combinations are used, (10) and (11) reduce to standard co-training. To the best of our knowledge we are the first to combine co-training with hetero-modal learning. This combined workflow is graphically depicted in Fig. 2.

3.4 Adversarial Domain Adaptation

The co-heterogeneous training of Sect. 3.3 is highly effective. Yet, it relies on accurate consensus predictions, which may struggle to handle significant appearance variations in \mathcal{D}_u that are not represented in \mathcal{D}_ℓ . Mask-based ADA offers an complementary approach that trains a network to output masks that follow a *prediction-based* distribution learned from labelled data [40]. Since liver shapes between \mathcal{D}_u and \mathcal{D}_ℓ should follow similar distributions, this provides an effective learning strategy that is not as confounded by differences in appearance. Following Tsai *et al.* [40], we can train a discriminator to classify whether a softmax output originates from a labelled- or unlabelled-dataset prediction. However, because we have a combinatorial number (15) of possible input phase combinations, *i.e.*, $\hat{\mathcal{X}}^*$, naively domain-adapting all corresponding predictions is prohibitively expensive. Fortunately, the formulations of (7) and (11) offer an effective and efficient solution. Namely, we can train the discriminator on the consensus prediction, M . This adapts the combinatorial number of possible predictions *at the same computational cost as performing ADA on only a single prediction*.

More formally, let $d(\cdot)$ be defined as an FCN discriminator, then the discriminator loss can be expressed as

$$\mathcal{L}_d = \frac{1}{N_\ell} \sum_{\mathcal{D}_\ell} \ell_{ce}(d(\hat{Y}^V), \mathbf{1}) + \frac{1}{N_u} \sum_{\mathcal{D}_u} \ell_{ce}(d(M, \mathbf{0})), \quad (12)$$

where ℓ_{ce} represents a pixel-wise cross-entropy loss. The opposing labels pushes the discriminator to differentiate semi-supervised consensus predictions from fully-supervised variants. Unlike natural image ADA [40], we do not wish to naively train the discriminator on all output classes, as it not reasonable to expect similar distributions of liver *lesion* shapes across datasets. Instead we train the discriminator on the *liver region*, *i.e.*, the union of healthy liver and lesion tissue predictions. Finally, when minimizing (12), we only optimize the discriminator weights. The segmentation network can now be tasked with fooling the discriminator, through the addition of an adversarial loss:

$$\mathcal{L}_{adv} = \frac{1}{N_u} \sum_{\mathcal{D}_u} \ell_{ce}(d(M, \mathbf{1})), \quad (13)$$

where the ground-truth labels for ℓ_{ce} have been flipped from (12). Note that here we use single-level ADA as we found the multi-level variant [40] failed to offer significant enough improvements to offset the added complexity. When minimizing (13), or (1) for that matter, the discriminator weights are frozen. We empirically set λ_{adv} to 0.001.

3.5 Pseudo-Labeling

By integrating co-heterogeneous training and ADA, CHASe can robustly segment challenging multi-phase unlabelled data. However some scenarios still present challenging edge cases, *e.g.*, lesions treated with TACE. See the supplementary for some visualizations. To manage these cases, we use a simple, but effective, domain-specific pseudo-labelling.

First, after convergence of (1), we produce predictions on \mathcal{D}_u using all available phases and extract any resulting 3D holes in the liver region (healthy tissue plus lesion) greater than 100 voxels. Since there should never be 3D holes, these are mistakes. Under the assumption that healthy tissue in both datasets should equally represented, we treat these holes as missing “lesion” predictions. We can then create a pseudo-label, Y_h , that indicates lesion at the hole, *with all others regions being ignored*. This produces a new “holes” dataset, $\mathcal{D}_h = \{\mathcal{X}, Y_h\}_{i=1}^{N_h}$, using image sets extracted from \mathcal{D}_u . We then finetune the model using (1), but replace the segmentation loss of (5) by

$$\begin{aligned} \mathcal{L}_{seg} &= \frac{1}{N_\ell} \sum_{V, Y \in \mathcal{D}_\ell} \ell_{seg}(f(V), Y) \\ &+ \frac{\lambda_h}{N_h} \sum_{\mathcal{X}, Y_h \in \mathcal{D}_h} \sum_{X \in \mathcal{X}^*} \ell_{seg}(f(X), Y_h), \end{aligned} \quad (14)$$

where we empirically set λ_h to 0.01 for all experiments. We found results were not sensitive to this value. While the hole-based pseudo-labels do not capture all errors, they only have to capture enough of missing appearances to guide CHASe’s training to better handle recalcitrant edge cases.

4 Results

Datasets. To execute CHASe, we require datasets of single-phase labelled and multi-phase unlabelled studies, \mathcal{D}_u and \mathcal{D}_ℓ , respectively. The goal is to robustly segment patient studies from \mathcal{D}_u while only having training mask labels from the less representative \mathcal{D}_ℓ dataset. 1) For \mathcal{D}_u , we collected 1147 multi-phase dynamic CT studies (4577 volumes in total) directly from the PACS of Chang Gung Memorial Hospital (CGMH). The only selection criteria were patients with biopsied or resected liver lesions, with dynamic contrast CTs taken within one month

before the procedure. Patients may have ICC, HCC, benign or metastasized tumors, along with co-occurring maladies, such as liver fibrosis, splenomegaly, or TACE-treated tumors. Thus, \mathcal{D}_u directly reflects the variability found in clinical scenarios. We used the DEEDS algorithm [19] to correct any misalignments. 2) For \mathcal{D}_ℓ , we collected 235 V-phase CT studies collected from as many public sources as we could locate [3, 7, 11, 18]. This is a superset of the LiTS training data [3], and consists of a mixture of healthy and pathological livers, with only HCC and metastasis represented.

Evaluation Protocols. 1) To evaluate performance on \mathcal{D}_u , 47 and 100 studies were randomly selected for validation and testing, respectively, with 90 test studies having all four phases. Given the extreme care required for lesion annotation, *e.g.*, the four readers used in the LiTS dataset [3], only the liver region, *i.e.*, union of healthy liver and lesion tissue, of the \mathcal{D}_u evaluation sets was annotated by a clinician. For each patient study, this was performed independently for each phase, with a final *study-wise* mask generated via majority voting. 2) We also evaluate whether the CHASe strategy of learning from unlabelled data can also improve performance on \mathcal{D}_ℓ . To do this, we split \mathcal{D}_ℓ , with 70%/10%/20% for training, validation, and testing, respectively, resulting in 47 test volumes. To measure CHASe’s impact on *lesion* segmentation, we use the \mathcal{D}_ℓ test set.

Backbone Network. We used a 2D segmentation backbone, an effective choice for many organs [14, 25, 33, 35], due to its simplicity and efficiency. We opt for the popular ResNet-50 [17]-based DeepLabv2 [8] network with PHNN-style deep supervision. We also tried VGG-16 [37], which also performed well and its results can be found in the supplementary. To create 3D masks we simply stack 2D predictions. **CHASe training.** We randomly sample multi-phase slices and, from them, randomly sample four out of the 15 phase combinations from \mathcal{X}^* to stochastically minimize (11) and (13). For standard co-training baselines, we sample all available phases to minimize (7). **Discriminator Network.** We use an atrous spatial pyramid pooling (ASPP) layer, employing dilation rates of 1,2,3,4 with a kernel size of 3 and a leaky ReLU with negative slope 0.2 as our activation function. After a 1×1 convolution, a sigmoid layer classifies whether a pixel belongs to the labelled or unlabelled dataset. *Specific details on data pre-processing, learning rates and schedules can be found in the supplementary material.*

4.1 Pathological Liver Segmentation

We first measure the performance of CHASe on segmenting pathological livers from the unlabeled CGMH PACS dataset, *i.e.*, \mathcal{D}_u . We use PHNN trained only on \mathcal{D}_ℓ as a baseline, testing against different unlabeled learning baselines, *i.e.*, co-training [4], co-heterogeneous training, ADA [40], and hole-based pseudo-labelling. We measure the mean DSC and average symmetric surface distance (ASSD). For non hetero-modality variants, we use majority voting across each single-phase prediction to produce a multi-phase output. We also test against

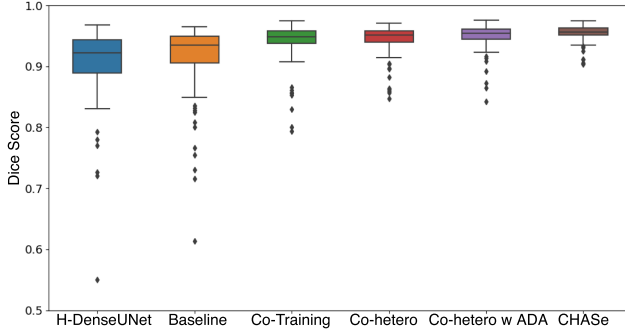


Fig. 3. Box and whisker plot. Shown is the distribution of DSCs of pathological liver segmentation on the CGMH PACS when using *all available* phases for inference.

the publicly available hybrid H-DenseUNet model [25], one of the best published models. It uses a cascade of 2D and 3D networks.

Table 1. Pathological liver segmentation. Mean DSC and ASSD scores on the CGMH PACS dataset are tabulated across different contrast phase inputs. “All” means all available phases are used as input. Number of samples are in parentheses.

Models	NC (96)		A (98)		V (97)		D (98)		All (100)	
	DSC	ASSD	DSC	ASSD	DSC	ASSD	DSC	ASSD	DSC	ASSD
HDenseUNet [25]	85.2	3.25	90.1	2.19	90.7	2.61	85.2	2.91	89.9	2.59
Baseline [14]	84.6	2.97	90.3	1.23	90.7	1.18	86.7	2.12	91.4	1.21
Baseline w pseudo	89.4	1.97	90.5	1.34	90.9	1.29	90.6	2.03	91.9	1.27
Baseline w ADA [40]	90.9	1.34	91.9	1.13	91.5	1.14	90.9	1.65	92.6	1.03
Co-training [31]	92.8	0.95	93.4	0.84	93.4	0.83	92.4	0.99	94.0	0.92
Co-hetero	93.4	0.81	93.7	0.77	94.5	0.79	93.6	0.86	94.7	0.89
Co-hetero w ADA	93.8	0.81	93.9	0.79	94.8	0.66	93.9	0.81	95.0	0.68
CHASE	94.0	0.79	94.2	0.74	94.9	0.66	94.1	0.80	95.4	0.63

As Table 1 indicates, despite being only a single 2D network, our PHNN baseline is strong, comparing similarly to the cascaded 2D/3D H-DenseUNet on our dataset¹. However, both H-DenseUNet and our PHNN baseline still struggle to perform well on the CGMH dataset, particularly on non *V*-phases, indicating that training on public data alone is not sufficient. In contrast, through its principled semi-supervised approach, CHASE is able to dramatically increase performance, producing boosts of 9.4%, 3.9%, 4.2%, 7.4%, and 4.0% in mean DSCs for inputs of NC, A, V, D, and all phases, respectively. As can also be seen,

¹ A caveat is that the public H-DenseUNet model was only trained on the LiTS subset of \mathcal{D}_ℓ .

all components contribute to these improvements, indicating the importance of each to the final result. Compared to established baselines of co-training and ADA, CHASe garners marked improvements. In addition, CHASe performs more strongly as more phases are available, something the baseline models are not always able to do. Results across all 15 possible combinations, found in our supplementary material, also demonstrate this trend.

More compelling results can be found in Fig. 3’s box and whisker plots. As can be seen, each component is not only able to reduce variability, but more importantly significantly improves worst-case results. These same trends are seen across all possible phase combinations. Compared to improvements in mean DSCs, these worst-case reductions, with commensurate boosts in reliability, can often be more impactful for clinical applications. Unlike CHASe, most prior work on pathological liver segmentation is fully-supervised. Wang *et al.* report 96.4% DSC on 26 LiTS volumes and Yang *et al.* [45] report 95% DSC on 50 test volumes with unclear healthy vs pathological status. We achieve comparable, or better, DSCs on 100 pathological multi-phase test studies. As such, we articulate a versatile strategy to use and learn from the vast amounts of uncured multi-phase clinical data housed within hospitals.

These quantitative results are supported by qualitative examples in Fig. 4. As the first two rows demonstrate, H-DenseUNet [25] and our baseline can perform inconsistently across contrast phases, with both being confused by the splenomegaly (overly large spleen) of the patient. The CHASe components are able to correct these issues. The third row in Fig. 4 depicts an example of a TACE-treated lesion, not seen in the public dataset and demonstrates how CHASe’s components can progressively correct the under-segmentation. Finally, the last row depicts the *worst-case* performance of CHASe. Despite this unfavorable selection, CHASe is still able to predict better masks than the alternatives. Of note, CHASe is able to provide tangible improvements in consistency and reliability, robustly predicting even when presented with image features not seen in \mathcal{D}_ℓ . More qualitative examples can be found in our supplementary material.

4.2 Liver and Lesion Segmentation

We also investigate whether using CHASe on unlabelled data can boost performance on the labelled data, *i.e.*, the \mathcal{D}_ℓ test set of 47 single-phase V volumes. Note, we include results from the public H-DenseUNet [25] implementation, even though it was only trained on LiTS and included some of our test instances originating from LiTS in its training set.

As Table 2 indicates, each CHASe component progressively boosts performance, with lesion scores being the most dramatic. The one exception is that the holes-based pseudo-labelling produces a small decrease in mean lesion scores. Yet, box and whisker plots, included in our supplementary material, indicate that holes-based pseudo-labelling boosts *median* values while reducing variability. Direct comparisons against other works, all typically using the LiTS challenge, are not possible, given the differences in evaluation data. nnUNet [20], the

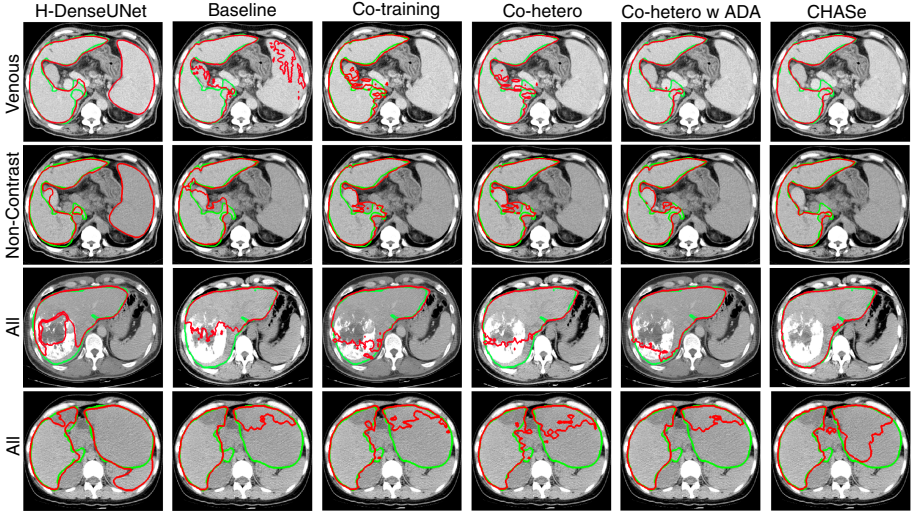


Fig. 4. Qualitative results. The first two rows are captured from the *same patient* across different contrast phases. The third and fourth row shows the performance when all available phases are included. **Green** and **red** curves depict the ground truth and segmentation predictions, respectively. (Color figure online)

Table 2. Ablation study on public data. Presented are test set DSC scores with their standard deviation of healthy liver, lesion, and liver region.

Model	Liver	Lesion	Liver region
HDenseUNet [25]	96.5 ± 2.0	51.7 ± 19.4	96.8 ± 1.8
Baseline [14]	96.3 ± 2.2	47.5 ± 24.1	96.6 ± 2.1
Co-training	96.3 ± 1.8	51.9 ± 20.5	96.7 ± 1.7
Co-hetero	96.4 ± 1.5	53.2 ± 19.1	96.7 ± 1.4
Co-hetero w ADA	96.5 ± 1.5	61.0 ± 17.2	97.0 ± 1.3
CHASe	96.8 ± 1.3	60.3 ± 18.0	97.1 ± 1.1

winner of the Medical Decathlon, reported 61% and 74% DSCs for their own validation and challenge test set, respectively. However, 57% of the patients in our test set are healthy, compared to the 3% in LiTS. More healthy cases will tend to make it a harder lesion evaluation set, as any amount of false positives will produce DSC scores of zero. For unhealthy cases, CHASe’s lesion mean DSC is 61.9% compared to 53.2% for PHNN. CHASe allows a standard backbone, with no bells or whistles, to achieve dramatic boosts in lesion segmentation performance. As such, these results broaden the applicability of CHASe, suggesting it can even improve the *source*-domain performance of fully-supervised models.

5 Conclusion

We presented CHASe, a powerful semi-supervised approach to organ segmentation. Clinical datasets often comprise multi-phase data and image features not represented in single-phase public datasets. Designed to manage this challenging domain shift, CHASe can adapt publicly trained models to robustly segment multi-phase clinical datasets *with no extra annotation*. To do this, we integrate co-training and hetero-modality into a co-heterogeneous training framework. Additionally, we propose a highly computationally efficient ADA for multi-view setups and a principled holes-based pseudo-labeling. To validate our approach, we apply CHASe to a highly challenging dataset of 1147 multi-phase dynamic contrast CT volumes of patients, all with liver lesions. Compared to strong fully-supervised baselines, CHASe dramatically boosts mean performance (>9% in NC DSCs), while also drastically improving worst-case scores. Future work should investigate 2.5D/3D backbones and apply this approach to other medical organs. Even so, these results indicate that CHASe provides a powerful means to adapt publicly-trained models to challenging clinical datasets found “in-the-wild”.

References

1. Bai, W., et al.: Semi-supervised learning for network-based cardiac MR image segmentation. In: Descoteaux, M., Maier-Hein, L., Franz, A., Jannin, P., Collins, D.L., Duchesne, S. (eds.) MICCAI 2017. LNCS, vol. 10434, pp. 253–260. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-66185-8_29
2. Ben-Cohen, A., Diamant, I., Klang, E., Amitai, M., Greenspan, H.: Fully convolutional network for liver segmentation and lesions detection. In: Carneiro, G., et al. (eds.) LABELS/DLMIA -2016. LNCS, vol. 10008, pp. 77–85. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-46976-8_9
3. Bilic, P., et al.: The liver tumor segmentation benchmark (LiTS). [arXiv:1901.04056](https://arxiv.org/abs/1901.04056) (2019). <http://arxiv.org/abs/1901.04056>, [arXiv: 1901.04056](https://arxiv.org/abs/1901.04056)
4. Blum, A., Mitchell, T.: Combining labeled and unlabeled data with co-training. In: Proceedings of the eleventh annual conference on Computational learning theory, pp. 92–100. Citeseer (1998)
5. Cai, J., et al.: Accurate weakly-supervised deep lesion segmentation using large-scale clinical annotations: slice-propagated 3D mask generation from 2D RECIST. In: Frangi, A.F., Schnabel, J.A., Davatzikos, C., Alberola-López, C., Fichtinger, G. (eds.) MICCAI 2018. LNCS, vol. 11073, pp. 396–404. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-00937-3_46
6. Chang, W.L., Wang, H.P., Peng, W.H., Chiu, W.C.: All about structure: adapting structural information across domains for boosting semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1900–1909 (2019)
7. Chaos: Chaos - combined (CT-MR) healthy abdominal organ segmentation (2019). https://chaos.grand-challenge.org/Combined_Healthy_Abdominal_Organ_Segmentation

8. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: DeepLab: semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. [arXiv:1606.00915](https://arxiv.org/abs/1606.00915) (2016)
9. Chen, Y.C., Lin, Y.Y., Yang, M.H., Huang, J.B.: CrDoCo: pixel-level domain transfer with cross-domain consistency. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 1791–1800 (2019)
10. Conze, P.H., et al.: Scale-adaptive supervoxel-based random forests for liver tumor segmentation in dynamic contrast-enhanced CT scans. *Int. J. Comput. Assist. Radiol. Surg.* **12**(2), 223–233 (2017). <https://doi.org/10.1007/s11548-016-1493-1>
11. Gibson, E., et al.: Multi-organ abdominal CT reference standard segmentations (2018). <https://doi.org/10.5281/zenodo.1169361>. This data set was developed as part of independent research supported by Cancer Research UK (Multidisciplinary C28070/A19985) and the National Institute for Health Research UCL/UCL Hospitals Biomedical Research Centre
12. Gotra, A., et al.: Liver segmentation: indications, techniques and future directions. *Insights Imaging* **8**(4), 377–392 (2017). <https://doi.org/10.1007/s13244-017-0558-1>
13. Han, X.: Automatic liver lesion segmentation using a deep convolutional neural network method. *arXiv preprint arXiv:1704.07239* (2017)
14. Harrison, A.P., Xu, Z., George, K., Lu, L., Summers, R.M., Mollura, D.J.: Progressive and multi-path holistically nested neural networks for pathological lung segmentation from CT images. In: Descoteaux, M., Maier-Hein, L., Franz, A., Janin, P., Collins, D.L., Duchesne, S. (eds.) MICCAI 2017. LNCS, vol. 10435, pp. 621–629. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-66179-7_71
15. Havaei, M., Guizard, N., Chapados, N., Bengio, Y.: HeMIS: hetero-modal image segmentation. In: Ourselin, S., Joskowicz, L., Sabuncu, M.R., Unal, G., Wells, W. (eds.) MICCAI 2016. LNCS, vol. 9901, pp. 469–477. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-46723-8_54
16. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–778 (2015)
17. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778 (2016)
18. Heimann, T., et al.: Comparison and evaluation of methods for liver segmentation from CT datasets. *IEEE Trans. Med. Imaging* **28**(8), 1251–1265 (2009). <https://doi.org/10.1109/TMI.2009.2013851>
19. Heinrich, M.P., Jenkinson, M., Brady, M., Schnabel, J.A.: MRF-based deformable registration and ventilation estimation of lung CT. *IEEE Trans. Med. Imaging* **32**, 1239–1248 (2013)
20. Isensee, F., et al.: nnU-Net: self-adapting framework for u-net-based medical image segmentation. *arXiv preprint arXiv:1809.10486* (2018)
21. Jin, D., et al.: Accurate esophageal gross tumor volume segmentation in PET/CT using two-stream chained 3D deep network fusion. In: Shen, D., et al. (eds.) MICCAI 2019. LNCS, vol. 11765, pp. 182–191. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-32245-8_21
22. Jin, D., et al.: Deep esophageal clinical target volume delineation using encoded 3D spatial context of tumors, lymph nodes, and organs at risk. In: Shen, D., et al. (eds.) MICCAI 2019. LNCS, vol. 11769, pp. 603–612. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-32226-7_67

23. Kuo, C., Cheng, S., Lin, C., Hsiao, K., Lee, S.: Texture-based treatment prediction by automatic liver tumor segmentation on computed tomography. In: 2017 International Conference on Computer, Information and Telecommunication Systems (CITS), pp. 128–132 (2017). <https://doi.org/10.1109/CITS.2017.8035318>
24. Lee, D.H.: Pseudo-label : the simple and efficient semi-supervised learning method for deep neural networks. In: ICML 2013 Workshop : Challenges in Representation Learning (WREPL) (2013)
25. Li, X., Chen, H., Qi, X., Dou, Q., Fu, C.W., Heng, P.A.: H-DenseUNet: hybrid densely connected UNet for liver and tumor segmentation from CT volumes. *IEEE Trans. Med. Imaging* **37**(12), 2663–2674 (2018)
26. Li, Y., Liu, L., Tan, R.T.: Decoupled certainty-driven consistency loss for semi-supervised learning (2019)
27. Lin, J.: Divergence measures based on the Shannon entropy. *IEEE Trans. Inf. Theory* **37**(1), 145–151 (1991)
28. Milletari, F., Navab, N., Ahmadi, S.A.: V-Net: fully convolutional neural networks for volumetric medical image segmentation. In: 2016 Fourth International Conference on 3D Vision (3DV), pp. 565–571. IEEE (2016)
29. Min, S., Chen, X., Zha, Z.J., Wu, F., Zhang, Y.: A two-stream mutual attention network for semi-supervised biomedical segmentation with noisy labels (2019)
30. Oliva, M., Saini, S.: Liver cancer imaging: role of CT, MRI, US and PET. *Cancer Imaging Official Publ. Int. Cancer Imaging Soc.* **4**(Spec No A), S42–6 (2004). <https://doi.org/10.1102/1470-7330.2004.0011>
31. Qiao, S., Shen, W., Zhang, Z., Wang, B., Yuille, A.: Deep co-training for semi-supervised image recognition. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 135–152 (2018)
32. Ronneberger, O., Fischer, P., Brox, T.: U-Net: convolutional networks for biomedical image segmentation. In: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (eds.) MICCAI 2015. LNCS, vol. 9351, pp. 234–241. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-24574-4_28
33. Roth, H.R., et al.: DeepOrgan: multi-level deep convolutional networks for automated pancreas segmentation. In: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (eds.) MICCAI 2015. LNCS, vol. 9349, pp. 556–564. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-24553-9_68
34. Roth, H.R., et al.: Spatial aggregation of holistically-nested convolutional neural networks for automated pancreas localization and segmentation. *Med. Image Anal.* **45**, 94–107 (2018)
35. Roth, H.R., et al.: Spatial aggregation of holistically-nested convolutional neural networks for automated pancreas localization and segmentation. *Med. Image Anal.* **45**, 94–107 (2018). <https://doi.org/10.1016/j.media.2018.01.006>, <http://www.sciencedirect.com/science/article/pii/S1361841518300215>
36. Roth, K., Konopczyński, T., Hesser, J.: Liver lesion segmentation with slice-wise 2D Tiramisu and Tversky loss function. arXiv preprint [arXiv:1905.03639](https://arxiv.org/abs/1905.03639) (2019)
37. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint [arXiv:1409.1556](https://arxiv.org/abs/1409.1556) (2014)
38. Soler, L., et al.: 3D image reconstruction for comparison of algorithm database: a patient specific anatomical and medical image database. Technical report, IRCAD, Strasbourg, France (2010)
39. Tajbakhsh, N., Jeyaseelan, L., Li, Q., Chiang, J., Wu, Z., Ding, X.: Embracing imperfect datasets: a review of deep learning solutions for medical image segmentation. *Med. Image Anal.* **63**, 101693 (2019)

40. Tsai, Y.H., Hung, W.C., Schultze, S., Sohn, K., Yang, M.H., Chandraker, M.: Learning to adapt structured output space for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7472–7481 (2018)
41. Vorontsov, E., Abi-Jaoudeh, N., Kadoury, S.: Metastatic liver tumor segmentation using texture-based omni-directional deformable surface models. In: Yoshida, H., Nappi, J., Saini, S. (eds.) Abdominal Imaging. Computational and Clinical Applications. ABD-MICCAI 2014. Lecture Notes in Computer Science, vol. 8676, pp. 74–83. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-13692-9_7
42. Vu, T.H., Jain, H., Bucher, M., Cord, M., Pérez, P.: Advent: adversarial entropy minimization for domain adaptation in semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2517–2526 (2019)
43. Wang, R., Cao, S., Ma, K., Meng, D., Zheng, Y.: Pairwise semantic segmentation via conjugate fully convolutional network. In: Shen, D., et al. (eds.) MICCAI 2019. LNCS, vol. 11769, pp. 157–165. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-32226-7_18
44. Xia, Y., et al.: 3D semi-supervised learning with uncertainty-aware multi-view co-training. arXiv preprint [arXiv:1811.12506](https://arxiv.org/abs/1811.12506) (2018)
45. Yang, D., et al.: Automatic liver segmentation using an adversarial image-to-image network. In: Descoteaux, M., Maier-Hein, L., Franz, A., Jannin, P., Collins, D.L., Duchesne, S. (eds.) MICCAI 2017. LNCS, vol. 10435, pp. 507–515. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-66179-7_58
46. Zhang, J., Xie, Y., Zhang, P., Chen, H., Xia, Y., Shen, C.: Light-weight hybrid convolutional network for liver tumor segmentation. In: Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19, pp. 4271–4277. International Joint Conferences on Artificial Intelligence Organization (2019). <https://doi.org/10.24963/ijcai.2019/593>
47. Zhang, L., Gopalakrishnan, V., Lu, L., Summers, R.M., Moss, J., Yao, J.: Self-learning to detect and segment cysts in lung CT images without manual annotation. In: 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018), pp. 1100–1103. IEEE (2018)
48. Zhang, Q., Fan, Y., Wan, J., Liu, Y.: An efficient and clinical-oriented 3D liver segmentation method. IEEE Access **5**, 18737–18744 (2017). <https://doi.org/10.1109/ACCESS.2017.2754298>
49. Zhou, Y., et al.: Semi-supervised 3D abdominal multi-organ segmentation via deep multi-planar co-training. In: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV), pp. 121–140. IEEE (2019)