

Incremental Learning for Multi-organ Segmentation with Partially Labeled Datasets

Pengbo Liu¹, Li Xiao¹, and S. Kevin Zhou¹

Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China
{liupengbo2019,xiaoli}@ict.ac.cn
s.kevin.zhou@gmail.com

Abstract. There exists a large number of datasets for organ segmentation, which are partially annotated, and sequentially constructed. A typical dataset is constructed at a certain time by curating medical images and annotating the organs of interest. In other words, new datasets with annotations of new organ categories are built over time. To unleash the potential behind these partially labeled, sequentially-constructed datasets, we propose to learn a multi-organ segmentation model through incremental learning (IL). In each IL stage, we lose access to the previous annotations, whose knowledge is assumingly captured by the current model, and gain the access to a new dataset with annotations of new organ categories, from which we learn to update the organ segmentation model to include the new organs. We give the first attempt to conjecture that the different distribution is the key reason for ‘catastrophic forgetting’ that commonly exists in IL methods, and verify that IL has the natural adaptability to medical image scenarios. Extensive experiments on five open-sourced datasets are conducted to prove the effectiveness of our method and the conjecture mentioned above.

Keywords: Incremental learning · Partially labeled datasets · Multi-organ segmentation.

1 Introduction

As the performance of deep learning has been verified in the field of computer vision, a large number of supervised datasets have been open-sourced, which greatly accelerating the development of deep learning technology. Unlike natural image datasets [3, 4, 12, 25] that are almost completely labeled for common categories, constructing a high-quality medical image dataset requires professional knowledge of different anatomical structures, so full annotation is very difficult to achieve in medical image scenarios, especially for segmentation tasks. Multi-organ segmentation is a very important task in medical image analysis scenes [26, 27]. However, there exist now many partially labeled datasets [1, 6, 23] that only with annotation of the organs of interest to the dataset builders. Fig. 1 gives some example images in partially labeled datasets.

There exists a ‘knowledge’ conflict from these partially labeled datasets, *e.g.*, the liver is marked as foreground in Dataset 1 and background in Datasets 2-4, as

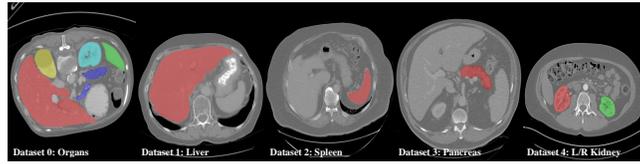


Fig. 1. Five typical partially labeled CT images from five different datasets.

shown in Fig. 1. Such a conflict prevents the direct utilization of all these datasets together, which limits their potential usefulness. So far, there has been some emerging research [5, 20, 28] on how to mix them together, and the performance of multi-organ segmentation was improved, proving that the unlabeled data in partially labeled datasets is also helpful for learning. As clinical needs increase, more categories and labeled datasets will be added, and current methods must retrain all datasets every time. When the aggregate scale of the datasets is large, there will be great pressure on storage and efficiency.

Incremental learning (IL) is a staged learning method, which learns new categories incrementally and loses access to the previous annotated images with old categories, making it an ideal choice for dealing with the above-mentioned storage and efficiency issues with better scalability in the future. And it can also solve the ethical and moral issues of sharing medical data by sharing model parameters. The main challenge in IL is the so-called ‘catastrophic forgetting’ [15]: how to keep the performance on old classes while learning new ones. IL has been studied for object recognition [7, 10, 11, 13, 19] and detection [14, 21, 22], but less in segmentation [2, 16, 18, 24]. In 2D medical image segmentation, Ozdemir and Goksel [18] made some attempts using the IL methods in natural images directly, with only two categories, and it mainly focuses on verifying the possibility of transferring the knowledge learned in the first category with more images to a second category with less images. In MiB [2], Cermelli et al. solved knowledge conflicts existing in other IL methods [11, 16] by remodeling old and new categories into background in loss functions of learning new categories and preserving old knowledge respectively, achieving a performance improvement.

However, catastrophic forgetting is still obvious even though a knowledge distillation loss is used commonly for combatting it. We make a hypothesis that the distillation loss can only be applied to the dataset at different stages *under the same distribution*. Different distributions cause the old model to output wrong responses to contents of *unseen categories* or seen categories that are quite different in appearance, violating the implicit assumption for the distillation to work. This is why ‘15-5’ setting and ‘Overlapped’ setting in MiB [2], whose distributions in different stages are closer, perform better than other comparison settings.

Compared with nature images, we believe medical images are inherently adaptable to IL due to the relatively fixed anatomical structures of the human body, e.g. liver is just close to right kidney, that old categories objects will emerge in new categories learning stage. This feature can maintain the *distribu-*

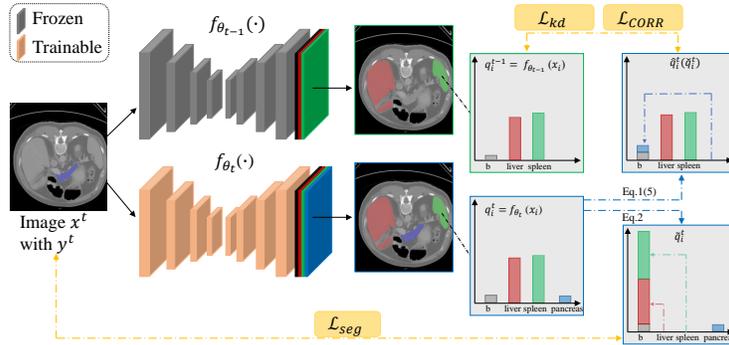


Fig. 2. Overview of the t^{th} stage of IL in multi-organ segmentation.

tion consistency of the datasets in different stages to a large extent. Then there raises an interesting question: *will the IL perform better on the medical image segmentation tasks?*

In this work, we present a novel study on incremental learning models for multi-organ segmentation task with four stages to aggregate five partially labeled datasets in medical image scene. Our main contributions can be summarized as:

- We give the first attempt to perform IL on multi-organ segmentation task, and firstly verify the effectiveness on multiple partially annotated datasets.
- Our extensive experiments on five open-sourced datasets help to prove our hypothesis that *different distributions* in different stages is the key reason for catastrophic forgetting in current IL settings.

2 Methodology

2.1 Problem Definition

The overview of t^{th} stage of IL in our method is shown in Fig. 2. Given an input image $x^t \in \mathcal{X}^t$, which is composed by a set of voxels \mathcal{I} , we firstly process it by the model in current stage, $f_{\theta_t}(\cdot)$ with trainable parameters θ_t , getting the output $q^t = f_{\theta_t}(x^t)$. For the learning of new categories (\mathcal{C}^t) in current stage, a one-hot vector y^t is the ground truth for \mathcal{C}^t in x . Label space \mathcal{Y}^t is expanded from \mathcal{Y}^{t-1} with \mathcal{C}^t , $\mathcal{Y}^t = \mathcal{Y}^{t-1} \cup \mathcal{C}^t = \mathcal{C}^0 \cup \mathcal{C}^1 \cup \dots \cup \mathcal{C}^t$. Note that the annotations of the old categories \mathcal{Y}^{t-1} will be inaccessible in the new stage under ideal IL settings. For preserving the knowledge of old categories, we process x by old model $f_{\theta_{t-1}}(\cdot)$ with frozen parameters θ_{t-1} and get $q^{t-1} = f_{\theta_{t-1}}(x^t)$ for reference. Trainable θ_t in t^{th} stage is expand with Θ_t , $\theta_t = \theta_{t-1} \cup \Theta_t$. We initialize Θ^t the same as MiB [2].

2.2 Background Remodeling

While the relatively fixed anatomical structure of human body in medical image brings help for IL, it also makes the ‘knowledge’ conflict more obvious. For

example, in Fig. 2, the label of voxel i on spleen is background in ground truth y^t of t^{th} stage. If we directly calculate the loss based on q^t and y^t , it will *punish the correct response* on spleen channel, and the same for other channels. Different from [18] based on Learning without Forgetting [11] using q^t directly, we remodel the background (b) channel of q^t based on MiB [2] by moving probabilities of new classes or old classes to background class, getting \hat{q}^t and \tilde{q}^t for the following calculation of the loss functions. Their definition is shown in Eq. 1 and Eq. 2, respectively.

$$\hat{q}_{i,c}^t = \begin{cases} \exp(q_{i,b}^t + \sum_{c \in \mathcal{C}^t} q_{i,c}^t) / \sum_{c \in \mathcal{Y}^t} \exp(q_{i,c}^t) & \text{if } c = b \\ \exp(q_{i,c}^t) / \sum_{c \in \mathcal{Y}^t} \exp(q_{i,c}^t) & \text{if } c \in \mathcal{Y}^{t-1} \& c \neq b \end{cases} \quad (1)$$

$$\tilde{q}_{i,c}^t = \begin{cases} \exp(\sum_{c \in \mathcal{Y}^{t-1}} q_{i,c}^t) / \sum_{c \in \mathcal{Y}^t} \exp(q_{i,c}^t) & \text{if } c = b \\ 0 & \text{if } c \in \mathcal{Y}^{t-1} \& c \neq b \\ \exp(q_{i,c}^t) / \sum_{c \in \mathcal{Y}^t} \exp(q_{i,c}^t) & \text{if } c \in \mathcal{C}^t \end{cases} \quad (2)$$

2.3 Loss Functions

In the IL setting, the whole loss function \mathcal{L} is composed by \mathcal{L}_{seg} for learning new knowledge of new categories and \mathcal{L}_{kd} for preserving old knowledge distilled from the previous model, $f_{\theta_{t-1}}$. For \mathcal{L}_{seg} , the cross-entropy loss is the most commonly used. We also invoke Dice loss [17] into \mathcal{L}_{seg} , which is verified useful in medical image segmentation.

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{seg}(\tilde{q}^t, y^t) + \mathcal{L}_{kd}(\hat{q}^t, \sigma(q^{t-1})) \\ &= \mathcal{L}_{CE}(\tilde{q}^t, y^t) + \mathcal{L}_{Dice}(\tilde{q}^t, y^t) + \mathcal{L}_{kd}(\hat{q}^t, \sigma(q^{t-1})) \end{aligned} \quad (3)$$

Where σ is the softmax operation.

CORR Loss We also devise a new corrective (CORR) loss to *reduce the low confident knowledge* and remove some false positive predictions, maybe caused by distribution disturbance between different datasets. CORR loss weakens voxel i with low confident response and enhances the influence of voxel i with high confident response, which is implemented by W defined in Eq. 4, where C is the target category in voxel i , \mathcal{THR} is the threshold of confidence and n is the scale exponent. We set \mathcal{THR} and n to 0.95 and 12 empirically. $y_{pseu}^{t-1} = \text{onehot}(\text{argmax}_{c \in \mathcal{Y}^{t-1}} q_c^{t-1})$.

$$W_{i,c} = \begin{cases} \left(\frac{\mathcal{THR}}{\sigma(q^{t-1})_{i,c}} \cdot y_{pseu,i,c}^{t-1} \right)^n & \text{if } c = C \\ 1 & \text{if } c \neq C \end{cases} \quad (4)$$

Then CORR loss can be calculated as shown in Eq. 6.

$$\tilde{q}_{i,c}^t = \begin{cases} q_{i,b}^t + \sum_{c \in \mathcal{C}^t} q_{i,c}^t & \text{if } c = b \\ q_{i,c}^t & \text{if } c \in \mathcal{Y}^{t-1} \& c \neq b \end{cases} \quad (5)$$

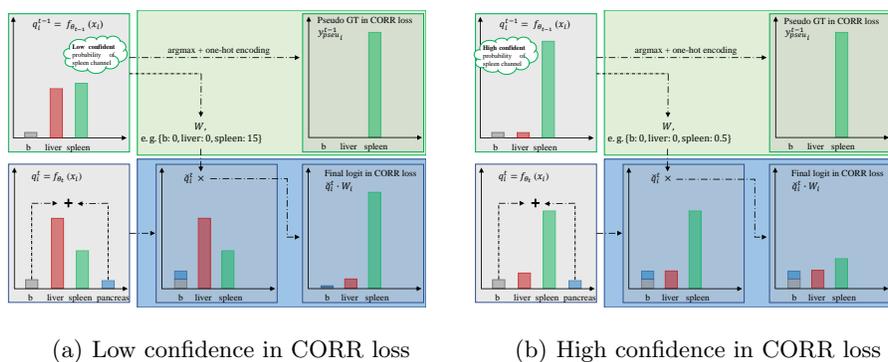


Fig. 3. Diagram of $\mathcal{L}_{CORR}(q^t, q^{t-1})$. (a) When the confidence of the pseudo GT out of model in $stage_{i-1}$ is low, we will reduce this voxel’s contribution to the CORR loss. (b) Contrary to (a) when the confidence is high.

$$\mathcal{L}_{CORR}(q^t, q^{t-1}) = \mathcal{L}_{CE}(\hat{q}^t \cdot W, y_{pseu}^{t-1}) = -\frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \sum_{c \in \mathcal{Y}^t} y_{pseu_{i,c}}^{t-1} \log(\sigma(\hat{q}_{i,c}^t * W_{i,c})) \quad (6)$$

The diagrams of CORR loss is shown in Fig. 3. In (a), When the voxel i is with a low confident response from $f_{\theta_{t-1}}$, i.e. $q_{i,C}^{t-1} < \mathcal{THR}$, $q_{i,C}^t$ will be enlarged $W_{i,C}$ times, where C is the channel of ‘Spleen’ here. The lower confidence, the higher probability in corresponding channel and the lower contribution to loss function. Vice versa at voxels with high confidence as shown in (b). So the whole loss function in our method is shown as below, and ω_1 , ω_2 and ω_3 are set to 1, 10, 1 based on [2].

$$\mathcal{L} = \omega_1 \cdot \mathcal{L}_{seg} + \omega_2 \cdot \mathcal{L}_{kd} + \omega_3 \cdot \mathcal{L}_{CORR} \quad (7)$$

3 Experiments

3.1 Implementation Details

Datasets and Preprocessing We choose five organs in our experiments, including liver, spleen, pancreas, right kidney and left kidney, and use five CT datasets as shown in Table 1. We process all datasets to a unified spacing (1.7, 0.79, 0.79) and normalize them with mean and std of 90.9 and 65.5 respectively. We split five datasets to 5 folds and select one fold, randomly, to evaluate our method. For our IL setting, five organs are collected in four stages, liver (F+P₁)→spleen (F+P₂)→pancreas (F+P₃)→R/L kidney (F+P₄). The annotations of different organs in dataset F are used separately in our IL setting.

Table 1. A summary of five benchmark datasets used in our experiments. [13] means 13 organs are in original dataset. We ignore other eight organs in our experiments. [T] means there are tumor labels in original dataset and we merge them into corresponding organs.

Datasets	Modality	# of labeled volumes	Annotated organs	Mean spacing (z, y, x)	Source
Dataset0 (F)	CT	30	Five organs [13]	(3.0, 0.76, 0.76)	Abdomen in [1]
Dataset1 (P_1)	CT	131	Liver [T]	(1.0, 0.77, 0.77)	Task03 in [23]
Dataset2 (P_2)	CT	41	Spleen	(1.6, 0.79, 0.79)	Task09 in [23]
Dataset3 (P_3)	CT	281	Pancreas [T]	(2.5, 0.80, 0.80)	Task07 in [23]
Dataset4 (P_4)	CT	210	L&R Kidneys [T]	(0.8, 0.78, 0.78)	KiTS [6]
All	CT	693	Five organs	(1.7, 0.79, 0.79)	-

Table 2. In the last stage, the 95th percentile Hausdorff distance (HD95) of the segmentation results of different methods on different datasets. The best result is shown in **bold**.

Methods/Organs	Liver∈F	Liver∈ P_1	Spleen∈F	Spleen∈ P_2	Pancreas∈F	Pancreas∈ P_3	R Kidney∈F	R Kidney∈ P_4	L Kidney∈F	L Kidney∈ P_4	Mean
ϕ_{F+P_1} (Liver)	2.39 ± 0.66	10.81 ± 25.54	-	-	-	-	-	-	-	-	-
ϕ_{F+P_2} (Spleen)	-	-	1.58 ± 0.41	24.74 ± 62.02	-	-	-	-	-	-	-
ϕ_{F+P_3} (Pancreas)	-	-	-	-	23.23 ± 33.60	6.45 ± 10.02	-	-	-	-	-
ϕ_{F+P_4} (R/L Kidney)	-	-	-	-	-	-	26.49 ± 54.34	15.15 ± 43.06	30.13 ± 61.08	6.67 ± 16.59	14.76
ϕ_F (Five organs)	1.58 ± 0.41	12.12 ± 17.81	1.00 ± 0.00	1.35 ± 0.41	5.39 ± 3.82	9.41 ± 8.98	1.36 ± 0.40	6.19 ± 4.25	2.27 ± 1.87	11.67 ± 16.45	5.23
FT	nan	nan	nan	nan	nan	nan	4.85 ± 2.33	8.16 ± 31.29	3.97 ± 1.66	3.01 ± 6.71	-
LwF [11]	2.33 ± 0.48	11.19 ± 24.66	46.11 ± 96.71	30.31 ± 76.26	4.89 ± 3.04	9.33 ± 13.54	16.03 ± 23.95	35.90 ± 57.56	25.68 ± 22.29	49.63 ± 54.46	23.14
ILT [16]	2.36 ± 0.53	11.13 ± 25.34	66.61 ± 102.64	30.59 ± 76.31	16.02 ± 19.58	10.37 ± 15.08	4.63 ± 2.21	29.34 ± 56.31	4.31 ± 1.21	21.80 ± 36.70	19.72
MiB [2]	2.56 ± 0.76	11.52 ± 25.03	1.48 ± 0.37	29.04 ± 72.38	3.59 ± 1.35	6.76 ± 9.71	4.87 ± 2.47	8.09 ± 30.75	3.63 ± 1.35	10.29 ± 28.95	8.19
MiBorgan (MiB+CORR)	2.19 ± 0.72	11.06 ± 24.24	1.96 ± 0.99	30.11 ± 75.21	3.97 ± 2.14	6.13 ± 6.04	4.44 ± 2.24	8.58 ± 33.12	3.04 ± 0.91	5.21 ± 12.16	7.45
MargExc MIA [20]	2.84 ± 1.53	4.04 ± 2.64	17.58 ± 7.27	1.00 ± 0.09	3.24 ± 0.69	3.96 ± 3.27	1.43 ± 0.14	1.28 ± 0.07	3.13 ± 0.58	1.68 ± 0.68	4.02

Code We implement our method based on open source code of 3D fullres version nnU-Net¹ [9]. We also used MONAI² and MiB³ during our algorithm development. Considering limiting the GPU memory consumption within 12Gb, our patch-size and batch-size are (80, 160, 128) and 2 in our experiments. We train the network with the same optimizer and learning rate policy as nnU-Net for about 400 epochs. The initial learning rate of the first stage and followed stages are set to 3e-4 and 15e-5.

Baselines To verify the effect of IL approach in the collection of multiple partially annotated datasets, we first construct experiments on each organ separately ($F+P_i$). Dataset F has five organs meanwhile, we also do a five-class segmentation experiment on F directly (F). To handle the datasets constructed sequentially, simple fine-tuning (FT) is the most intuitive. And we compare our proposed method with some state-of-the-art (SOTA) methods, LwF [11], ILT [16] and MiB [2]. For the results please refer to Sect. 3.2.

3.2 Results and Discussions

We use Dice coefficient (DC) and 95th percentile Hausdorff distance (HD95) for comparing different methods. The results are shown in Table 3 and Table 2.

¹ github.com/mic-dkfz/nnunet

² <https://monai.io/>

³ <https://github.com/fcd194/MiB>

Table 3. In the last stage, the Dice coefficient (DC) of the segmentation results of different methods on different datasets. The best result is shown in **bold**.

Methods\Organs	Liver∈F	Liver∈P ₁	Spleen∈F	Spleen∈P ₂	Pancreas∈F	Pancreas∈P ₃	R Kidney∈F	R Kidney∈P ₄	L Kidney∈F	L Kidney∈P ₄	Mean
ϕ_{F+P_1} (Liver)	.958 ± .017	.964 ± .030	-	-	-	-	-	-	-	-	-
ϕ_{F+P_2} (Spleen)	-	-	.951 ± .010	.955 ± .028	-	-	-	-	-	-	-
ϕ_{F+P_3} (Pancreas)	-	-	-	-	.809 ± .053	.850 ± .071	-	-	-	-	-
ϕ_{F+P_4} (R/L Kidney)	-	-	-	-	-	-	.917 ± .038	.970 ± .027	.913 ± .031	.963 ± .042	.925
ϕ_F (Five organs)	.967 ± .010	.948 ± .027	.969 ± .007	.955 ± .005	.786 ± .091	.704 ± .149	.949 ± .016	.884 ± .093	.926 ± .057	.825 ± .172	.891
FT	.000 ± .000	.000 ± .000	.000 ± .000	.000 ± .000	.000 ± .000	.000 ± .000	.917 ± .016	.978±.011	.919 ± .015	.973±.021	.379
LwF [11]	.959 ± .017	.961 ± .032	.940 ± .021	.956 ± .024	.804 ± .044	.807 ± .102	.912 ± .016	.944 ± .043	.879 ± .044	.900 ± .104	.906
ILT [16]	.958 ± .017	.962 ± .029	.937 ± .020	.949 ± .035	.795 ± .046	.807 ± .096	.913 ± .022	.955 ± .039	.912 ± .017	.919 ± .103	.911
MiB [2]	.961±.017	.959 ± .037	.953±.015	.953±.033	.817±.048	.819±.111	.918±.018	.972 ± .035	.920 ± .016	.952 ± .073	.922
MiBOrgan(MiB+CORR)	.961±.017	.960±.034	.950 ± .016	.950 ± .035	.809 ± .049	.814 ± .111	.917 ± .018	.971 ± .028	.921±.019	.953 ± .077	.921
MargExc MIA [20]	.969 ± .012	.957 ± .009	.924 ± .009	.970 ± .008	.836 ± .006	.808 ± .041	.946 ± .012	.952 ± .013	.978 ± .013	.972 ± .004	.931

Table 4. The Dice coefficient (DC) and 95th percentile Hausdorff distance (HD95) of the segmentation results in different stages. The best result is shown in **bold**. ‘-’ means Not Applicable.

Setting	Organs\Stages	DC				HD			
		S0	S1	S2	S3	S0	S1	S2	S3
FT	Liver	.963±.028	.000 ± .000	.000 ± .000	.000 ± .000	9.23±23.26	nan	nan	nan
	Spleen	-	.961±.013	.000 ± .000	.000 ± .000	-	1.35±0.34	nan	nan
	Pancreas	-	-	.844±.091	.000 ± .000	-	-	5.00±6.82	nan
	L Kidney	-	-	-	.970±.024	-	-	-	7.74±29.26
	R Kidney	-	-	-	.966±.027	-	-	-	3.13±6.30
LwF	Liver	.963±.028	.962±.028	.962±.028	.961±.030	9.23±23.26	9.43 ± 23.20	9.36±22.80	9.53 ± 22.50
	Spleen	-	.948 ± .024	.948 ± .024	.949 ± .024	-	40.02 ± 85.27	44.40 ± 86.85	37.08 ± 85.98
	Pancreas	-	-	.806 ± .094	.807 ± .098	-	-	15.59 ± 32.45	8.90 ± 12.97
	L Kidney	-	-	-	.940 ± .042	-	-	-	33.36 ± 54.84
	R Kidney	-	-	-	.897 ± .098	-	-	-	46.57 ± 52.10
MiBOrgan (MiB+CORR)	Liver	.963±.028	.961 ± .032	.961 ± .031	.961±.032	9.23±23.26	9.29±22.36	14.70 ± 36.24	9.40±22.12
	Spleen	-	.949 ± .026	.950±.025	.950±.029	-	43.13 ± 84.72	36.09±83.96	18.05±58.54
	Pancreas	-	-	.819 ± .100	.814±.107	-	-	8.80±14.28	5.92±5.82
	L Kidney	-	-	-	.964 ± .033	-	-	-	8.05 ± 30.97
	R Kidney	-	-	-	.949 ± .073	-	-	-	4.94 ± 11.38

Annotations used separately: When we do not aggregate these partially labeled data together by IL, there are some limitations in the results. Five-class segmentation model ϕ_F trained on ‘fully’ annotated dataset F, has a good performance on itself, but can not generalize well to other datasets due to the scale of the dataset F. And we train four models, ϕ_{F+P_*} , one model per organ segmentation task trained on corresponding datasets ($F+P_*$), then all datasets can be used. We can get the best performance on DC metric, but worse on HD95 metric. And this method is also poor in scalability and efficiency when the categories grow in the future.

Aggregating partially labeled datasets: When we aggregate these partially labeled datasets together, the most intuitive method FT is the worst. It has no preservation about the old knowledge because there is no restraint for it. LwF and ILT perform better than model ϕ_F , because they learn on much more data than dataset F, 554 vs 24. But wrong supervision limits the performance of LwF and ILT on these datasets, due to ‘knowledge’ conflict.

After we remodel the background of the predictions out of f_{θ_t} , MiB gets a large improvement on the performance of DC and HD all (MiB vs Lwf/ILT), ob-

taining a comparable performance on DC and an obvious improvement on HD95 compared to models trained separately. Solving ‘knowledge’ conflict between different partially annotated datasets, not only the preservation of the performance on the old categories but also the learning of the new categories has been improved. Adding CORR loss, MiBOrgan gets an exchange of 0.1% drop on DC for 9% enhancement on HD95. It shows that CORR loss removes some low confident predictions and reduces false positive results, thereby reducing HD95 to a certain extent.

Compared with partially supervised method: We also compare the result with SOTA partially supervised method, MargExc MIA [20], which have access to all partially labeled datasets and annotations in one time. The results are taken from [20], which can be regarded as our upper-bound. The performance of our method is close, but without accessing all the training data in one time.

Performance on different stages: In Table 4, we also show the performance on old and new categories of models in different stages of three typical IL settings, FT, LwF, and MiBOrgan. We can observe that FT can always get the best results on the categories learned in current stage. Because FT only needs to focus on learning fully supervised new categories with a good ‘pretrained’ base model, which is trained in former stage. MiBOrgan and LwF can preserve old knowledge and learn new knowledge meanwhile, due to the constraint from distillation loss. The more difficult task also makes the learning of new categories not as good as FT. MiBOrgan takes one step closer to the best through solving ‘knowledge’ conflict existing in LwF. And we found no obvious forgetting problem in our medical image scene, which can help to prove our hypothesis — *distribution consistency in medical image helps retain knowledge of old categories*. This implies that IL is a suitable choice for medical image analysis.

4 Conclusion

To unleash the potential from a collection of partially labeled datasets in medical image scenarios, we introduce incremental learning (IL) to aggregate them by stages, which marks the first attempt in the literature to verify for a multi-organ segmentation task the extent of the key issue associated with IL — different distributions between IL stages may mislead the direction of learning process. The introduction of CORR loss also helps to reduce the false positive predictions by removing predictions with low confidence. IL methods have natural adaptability to medical image scenarios due to the relatively fixed anatomical structure of human body, which is an inspiration to natural image scene that introducing an external dataset containing old categories of objects under the similar distribution in the new stage will give the same effect. We believe it will be a valuable research direction in the future. Further, we plan to explore a universal segmentation model [8] based on IL method, containing organs from different regions, which presents a new challenge for using IL in medical image segmentation.

References

1. Bennett, L., Zhoubing, X., Juan, Eugenio, I., Martin, S., Thomas, Robin, L., Arno, K.: 2015 miccai multi-atlas labeling beyond the cranial vault – workshop and challenge (2015). <https://doi.org/10.7303/syn3193805>
2. Cermelli, F., Mancini, M., Bulò, S.R., Ricci, E., Caputo, B.: Modeling the background for incremental learning in semantic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9233–9242 (2020)
3. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
4. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International journal of computer vision* **88**(2), 303–338 (2010)
5. Fang, X., Yan, P.: Multi-organ segmentation over partially labeled datasets with multi-scale feature abstraction. *IEEE Transactions on Medical Imaging* **39**(11), 3619–3629 (2020)
6. Heller, N., Sathianathan, N., Kalapara, A., Walczak, E., Moore, K., Kaluzniak, H., Rosenberg, J., Blake, P., Rengel, Z., Oestreich, M., Dean, J., Tradewell, M., Shah, A., Tejpaul, R., Edgerton, Z., Peterson, M., Raza, S., Regmi, S., Papanikolopoulos, N., Weight, C.: The kits19 challenge data: 300 kidney tumor cases with clinical context, CT semantic segmentations, and surgical outcomes. arXiv:1904.00445 (2019)
7. Hou, S., Pan, X., Loy, C.C., Wang, Z., Lin, D.: Learning a unified classifier incrementally via rebalancing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 831–839 (2019)
8. Huang, C., Han, H., Yao, Q., Zhu, S., Zhou, S.K.: 3d u²-net: A 3d universal u-net for multi-domain medical image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 291–299. Springer (2019)
9. Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H.: nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature Methods* **18**(2), 203–211 (2021)
10. Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A.A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al.: Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences* **114**(13), 3521–3526 (2017)
11. Li, Z., Hoiem, D.: Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence* **40**(12), 2935–2947 (2017)
12. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
13. Liu, Y., Schiele, B., Sun, Q.: Meta-aggregating networks for class-incremental learning. arXiv preprint arXiv:2010.05063 (2020)
14. Marra, F., Saltori, C., Boato, G., Verdoliva, L.: Incremental learning for the detection and classification of gan-generated images. In: 2019 IEEE International Workshop on Information Forensics and Security (WIFS). pp. 1–6 (2019). <https://doi.org/10.1109/WIFS47025.2019.9035099>
15. McCloskey, M., Cohen, N.J.: Catastrophic interference in connectionist networks: The sequential learning problem. In: *Psychology of learning and motivation*, vol. 24, pp. 109–165. Elsevier (1989)

16. Michieli, U., Zanuttigh, P.: Incremental learning techniques for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops. pp. 0–0 (2019)
17. 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)
18. Ozdemir, F., Goksel, O.: Extending pretrained segmentation networks with additional anatomical structures. *International journal of computer assisted radiology and surgery* **14**(7), 1187–1195 (2019)
19. Rebuffi, S.A., Kolesnikov, A., Sperl, G., Lampert, C.H.: icarl: Incremental classifier and representation learning. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 2001–2010 (2017)
20. Shi, G., Xiao, L., Chen, Y., Zhou, S.K.: Marginal loss and exclusion loss for partially supervised multi-organ segmentation. *Medical Image Analysis* p. 101979 (2021)
21. Shi, K., Bao, H., Ma, N.: Forward vehicle detection based on incremental learning and fast r-cnn. In: 2017 13th International Conference on Computational Intelligence and Security (CIS). pp. 73–76 (2017). <https://doi.org/10.1109/CIS.2017.00024>
22. Shmelkov, K., Schmid, C., Alahari, K.: Incremental learning of object detectors without catastrophic forgetting. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3400–3409 (2017)
23. Simpson, A.L., Antonelli, M., Bakas, S., Bilello, M., Farahani, K., Van Ginneken, B., Kopp-Schneider, A., Landman, B.A., Litjens, G., Menze, B., Ronneberger, O., Summers, R.M., Bilic, P., Christ, P.F., Do, R.K., Gollub, M., Golia-Pernicka, J., Heckers, S.H., Jarnagin, W.R., McHugo, M.K., Napel, S., Vorontsov, E., Maier-Hein, L., Cardoso, M.J.: A large annotated medical image dataset for the development and evaluation of segmentation algorithms. arXiv:1902.09063 (2019)
24. Tasar, O., Tarabalka, Y., Alliez, P.: Incremental learning for semantic segmentation of large-scale remote sensing data. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* **12**(9), 3524–3537 (2019)
25. Zhou, B., Zhao, H., Puig, X., Fidler, S., Barriuso, A., Torralba, A.: Scene parsing through ade20k dataset. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 633–641 (2017)
26. Zhou, S.K., Greenspan, H., Davatzikos, C., Duncan, J.S., van Ginneken, B., Madabhushi, A., Prince, J.L., Rueckert, D., Summers, R.M.: A review of deep learning in medical imaging: Imaging traits, technology trends, case studies with progress highlights, and future promises. *Proceedings of the IEEE* (2021)
27. Zhou, S.K., Rueckert, D., Fichtinger, G.: *Handbook of Medical Image Computing and Computer Assisted Intervention*. Academic Press (2019)
28. Zhou, Y., Li, Z., Bai, S., Wang, C., Chen, X., Han, M., Fishman, E., Yuille, A.L.: Prior-aware neural network for partially-supervised multi-organ segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10672–10681 (2019)