

USCL: Pretraining Deep Ultrasound Image Diagnosis Model through Video Contrastive Representation Learning

Yixiong Chen^{1,4,*}, Chunhui Zhang^{2,*}, Li Liu^{3,4}, and Cheng Feng^{5,6}, Changfeng Dong^{5,6}, Yongfang Luo^{5,6}, Xiang Wan^{3,4}

¹ School of Data Science, Fudan university, Shanghai, China

² Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

³ Shenzhen Research Institute of Big Data, Shenzhen, China

liuli@cuhk.edu.cn

⁴ The Chinese University of Hong Kong Shenzhen, Shenzhen, China

⁵ Shenzhen Third People's Hospital, Shenzhen, China

⁶ Southern University of Science and Technology, Shenzhen, China

Abstract. Most deep neural networks (DNNs) based ultrasound (US) medical image analysis models use pretrained backbones (*e.g.*, ImageNet) for better model generalization. However, the domain gap between natural and medical images causes an inevitable performance bottleneck. To alleviate this problem, an US dataset named US-4 is constructed for direct pretraining on the same domain. It contains over 23,000 images from four US video sub-datasets. To learn robust features from US-4, we propose an US semi-supervised contrastive learning method, named USCL, for pretraining. In order to avoid high similarities between negative pairs as well as mine abundant visual features from limited US videos, USCL adopts a sample pair generation method to enrich the feature involved in a single step of contrastive optimization. Extensive experiments on several downstream tasks show the superiority of USCL pretraining against ImageNet pretraining and other state-of-the-art (SOTA) pretraining approaches. In particular, USCL pretrained backbone achieves fine-tuning accuracy of over 94% on POCUS dataset, which is 10% higher than 84% of the ImageNet pretrained model. The source codes of this work are available at <https://github.com/983632847/USCL>.

Keywords: Ultrasound · Pretrained model · Contrastive learning.

1 Introduction

Due to the low cost and portability, ultrasound (US) is a widely used medical imaging technique, leading to the common application of US images [2,24] for clinical diagnosis. To date, deep neural networks (DNNs) [9] are one of the most popular automatic US image analysis techniques. When training DNN on US

* The first two authors contributed equally. This work was done at Shenzhen Research Institute of Big Data (SRIBD).

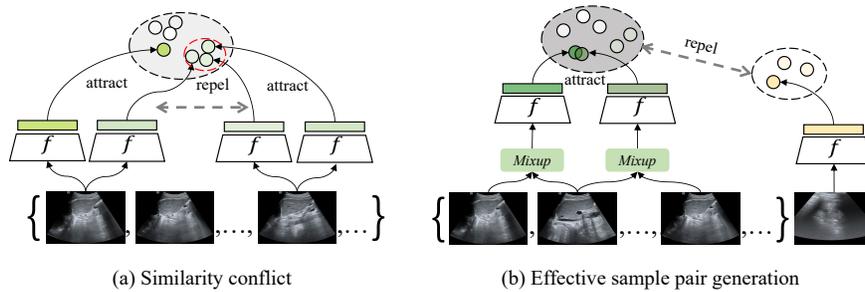


Fig. 1. Motivation of USCL. SPG tackles the harmful similarity conflict of traditional contrastive learning. (a) Similarity conflict: if a negative sample pair comes from different frames of the same video, they might be more similar than positive samples augmented from a frame, which confuses the training. (b) SPG ensures negative pairs coming from different videos, thus are dissimilar. The sample interpolation process also helps the positive pairs to have appropriate similarities, and enriches the feature involved in comparison. Representations are learned by gathering positive pairs close in representation space and pushing negative pairs apart.

images, a big challenge is the data scarcity, which is often dealt with parameters transferred from pretrained backbones (*e.g.*, ImageNet [19] pretrained VGG or ResNet). But model performance on downstream tasks suffers severely from the domain gap between *natural* and *medical* images [12]. There is a lack of public well-pretrained models specifically for US images due to the insufficient labeled pretraining US data caused by the high cost of specialized annotations, inconsistent labeling criterion and data privacy issue.

Recently, more and more literature tend to utilize unsupervised methods [3,7] to avoid medical data limitation for pretraining. The common practice is to pre-train models with pretext tasks and evaluate the representations on specific downstream tasks. Yet most existing methods can only outperform ImageNet pretraining with high-cost multi-modal data [11,14]. To get powerful pretrained models from US videos, we first build an US video dataset to alleviate data shortage. Secondly, contrastive learning[25,7,4] is also exploited to reduce the dependence on accurate annotations due to its good potential ability to learn robust visual representations without labels. However, given the fact that most US data are in video format, normal contrastive learning paradigm (*i.e.*, SimCLR [4] and MoCo [7], which considers two samples augmented from each image as a positive pair, and samples from different images as negative pairs) will cause high similarities between negative pairs sampled from the same video and mislead the training. This problem is called *similarity conflict* (Fig. 1 (a)) in this work. *Thus, is there a method which can avoid similarity conflict of contrastive learning and train a robust DNN backbone with US videos?*

To answer this question, we find that image features from the same US video can be seen as a cluster in semantic space, while features from different videos come from different clusters. We design a sample pair generation (SPG) scheme

Table 1. Statistics of the US-4 dataset containing 4 video-based sub-datasets. The total number of images is 23,231, uniformly sampled from 1051 videos. Most videos contain 10~50 similar images, which ensures the good property of semantic clusters.

Sub-dataset	Organ	Image size	Depth	Frame rate	Classes	Videos	Images
Butterfly [1]	Lung	658×738	-	23Hz	2	22	1533
CLUST [24]	Liver	434×530	-	19Hz	5	63	3150
Liver Fibrosis	Liver	600×807	~8cm	28Hz	5	296	11714
COVID19-LUSMS	Lung	747×747	~10cm	17Hz	4	670	6834

to make contrastive learning fit the natural clustering characteristics of US video (Fig. 1 (b)). Two samples from the same video act as a positive pair and two samples from different videos are regarded as a negative pair. In this process, two positive samples can naturally be seen as close feature points in the representation space, while negative samples have enough semantic differences to avoid similarity conflict. In addition, SPG does not simply choose frames as samples (*e.g.*, key frame extraction [18]), we put forward sample interpolation contrast to enrich features. Samples are generated from multiple-frame random interpolation so that richer features can be involved in positive-negative comparison. This method makes the semantic cohesion appear at the volume level of the ultrasound representation space [13] instead of the instance level. Combined with SPG, our work develops a semi-supervised contrastive learning method to train a generic model with US videos for downstream US image analysis. Here, the whole framework is called *ultrasound contrastive learning (USCL)*, which combines supervised learning to learn category-level discriminative ability, and contrastive learning to enhance instance-level discriminative ability.

2 US-4 Dataset

In this work, we construct a new US dataset named US-4, which is collected from four different convex probe [2] US datasets, involving two scan regions (*i.e.*, lung and liver). Among the four sub-datasets of US-4, *Liver Fibrosis* and *COVID19-LUSMS* datasets are collected by local sonographers [6,17], *Butterfly* [1] and *CLUST* [24] are two public sub-datasets. The first two sub-datasets are collected with *Resona 7T* ultrasound system, the frequency is FH 5.0 and the pixel size is 0.101mm - 0.127mm. All sub-datasets contain labeled images captured from videos for classification task. In order to generate a diverse and sufficiently large dataset, images are selected from original videos with a suitable sampling interval. For each video with frame rate T , we extract $n = 3$ samples per second with sampling interval $I = \frac{T}{n}$, which ensures that US-4 contains sufficient but not redundant information of videos. This results in 1051 videos and 23,231 images. The different attributes (*e.g.*, depth and frame rate) of dataset are described in Tab. 1. The US-4 dataset is relatively balanced in terms of images in each video,

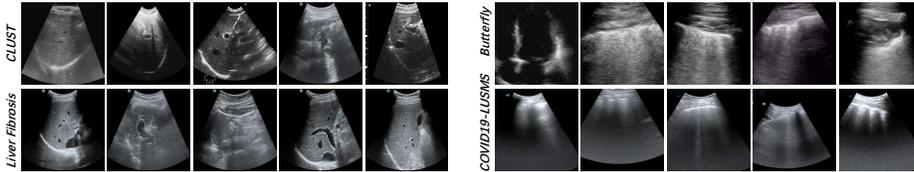


Fig. 2. Examples of US image in US-4.

where most videos contain tens of US images. Some frame examples are shown in Fig. 2.

3 Methodology

This section first formulates the proposed USCL framework (Fig. 3), then describes the details of sample pair generation. Finally, the proposed USCL will be introduced.

3.1 Problem Formulation

Given a video V_i from the US-4 dataset, USCL first extracts images to obtain a balanced distributed frame set $\mathbb{F}_i^K = \{\mathbf{f}_i^{(k)}\}_{k=1}^K$, where K is the number of extracted images. Next, a *sampler* Θ is applied to randomly sample M images, denoted as $\mathbb{F}_i^M = \{\mathbf{f}_i^{(m)}\}_{m=1}^M$ with $2 \leq M \ll K$. A following *mixed frame generator* $G : \mathbb{F}_i^M \rightarrow \mathbb{S}_i^2$ obtains two images, where $\mathbb{S}_i^2 = \{\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}\}$ is a positive pair followed by two data augmentation operations $Aug = \{Aug_i, Aug'_i\}$. These augmentations including random cropping, flipping, rotation and color jittering are used for perturbing positive pairs, making the trained backbones invariant to scale, rotation, and color style.

The objective of USCL is to train a backbone f from training samples $\{(\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}), \mathbf{y}_i\}_{i=1}^N$ by combining self-supervised contrastive learning loss \mathcal{L}_{con} and supervised cross-entropy (CE) loss \mathcal{L}_{sup} , where N is the number of videos in a training batch. Therefore, the USCL framework formulation aims to minimize following loss \mathcal{L} :

$$\mathcal{L} = \mathcal{L}_{con}(g(f(Aug(G(\mathbf{f}))); \mathbf{w}_f); \mathbf{w}_g)) + \lambda \mathcal{L}_{sup}(h(f(Aug(G(\mathbf{f}))); \mathbf{w}_f); \mathbf{w}_h); \mathbf{y}), \quad (1)$$

where λ is a hyper-parameter, $\mathbf{f} = \{\{\mathbf{f}_i^{(m)}\}_{m=1}^M\}_{i=1}^N$ are frames sampled from a batch of videos for training, $G(\mathbf{f}) = \{\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}\}_{i=1}^N$ are positive pairs, and $\mathbf{y} = \{\mathbf{y}_i\}_{i=1}^N$ are corresponding class labels. f , g and h denote backbone, projection head (two-layer MLP) and linear classifier, respectively. Different from most existing contrastive learning methods, USCL treats contrastive loss in Eq. (1) as a consistency regularization (CR) term, which improves the performance of pretraining backbone by combining supervised loss in a mutually reinforcing way.

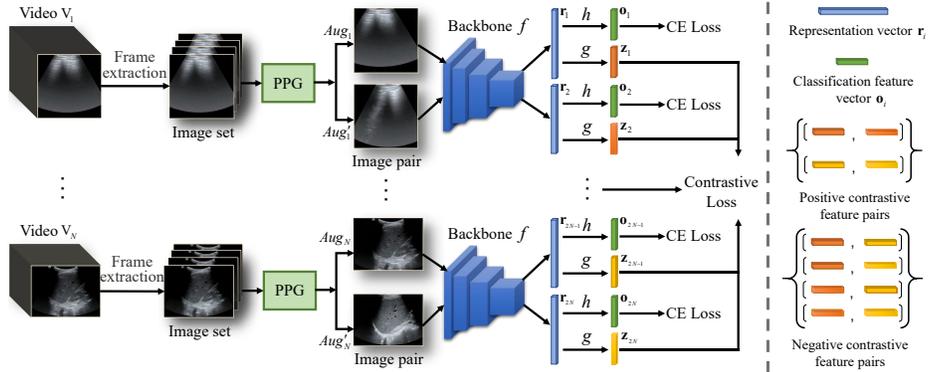


Fig. 3. System framework of the proposed USCL, which consists of sample pair generation and semi-supervised contrastive learning. (i) USCL extracts evenly distributed image sets from every US video as image dataset. (ii) The positive pair generation (PPG) module consists of a *sampler* Θ random sampling several images from an image set, and a *mixed frame generator* G obtaining two images. A generated positive pair is processed by two separate data augmentation operations. (iii) A backbone f , a projection head g and a classifier h are trained simultaneously by minimizing the self-supervised contrastive learning loss and supervised CE loss.

Here, label information instructs the model to recognize samples with different labels to be negative pairs, and contrastive process learns how US images can be semantically similar or different to assist better classification.

3.2 Sample Pair Generation

Most of the existing contrastive learning approaches construct positive pairs by applying two random data augmentations image by image. When directly applying them to the US frames, the contrastive learning fails to work normally due to the similarity conflict problem (*i.e.*, two samples coming from the same video are too similar to be a negative pair). To solve this problem, a sample pair generation (SPG) scheme is designed: it generates positive pairs with the positive pair generation (PPG) module⁷, and any two samples from different positive pairs are regarded as a negative pair.

The PPG module regards an evenly distributed image set extracted from a video as a semantic cluster, and different videos belong to different clusters. This kind of organization fits the purpose of contrastive learning properly. We expect the model can map the semantic clusters to feature clusters. Then PPG generates two images as a positive sample pair from each cluster. Note that only one positive pair is generated from a video, which can prevent the aforementioned similarity conflict problem.

⁷ For more details of PPG module, see the Supplementary Material Section 2.

In detail, firstly, a *sampler* Θ is applied to randomly sample three images $\widehat{\mathbf{x}}_i^{(1)}$, $\widehat{\mathbf{x}}_i^{(2)}$, and $\widehat{\mathbf{x}}_i^{(3)}$ in chronological order from an image set $\{\mathbf{f}_i^{(m)}\}_{m=1}^M$. Secondly, a delicate *mixed frame generator* G is performed to generate a positive sample pair. The image $\widehat{\mathbf{x}}_i^{(2)}$ is set as the anchor image, while $\widehat{\mathbf{x}}_i^{(1)}$ and $\widehat{\mathbf{x}}_i^{(3)}$ are perturbation images. In a mini-batch, G constructs positive sample pairs in interpolation manner via the mixup operation between anchor image and two perturbation images as follows.

$$\begin{cases} (\mathbf{x}_i^{(1)}, \mathbf{y}_i^{(1)}) = \xi_1(\widehat{\mathbf{x}}_i^{(2)}, \widehat{\mathbf{y}}_i^{(2)}) + (1 - \xi_1)(\widehat{\mathbf{x}}_i^{(1)}, \widehat{\mathbf{y}}_i^{(1)}) \\ (\mathbf{x}_i^{(2)}, \mathbf{y}_i^{(2)}) = \xi_2(\widehat{\mathbf{x}}_i^{(2)}, \widehat{\mathbf{y}}_i^{(2)}) + (1 - \xi_2)(\widehat{\mathbf{x}}_i^{(3)}, \widehat{\mathbf{y}}_i^{(3)}) \end{cases}, \quad (2)$$

where $\{\widehat{\mathbf{y}}_i^{(k)}\}_{k=1}^3$ are corresponding labels. $\xi_1, \xi_2 \sim \text{Beta}(\alpha, \beta)$, where α, β are parameters of *Beta* distribution.

In our contrastive learning process, sample pairs are then fed to the backbone followed by the projection head for contrastive learning task. The proposed PPG module has several benefits: 1) Interpolation makes every point in the feature convex hull enclosed by the cluster boundary possible to be sampled, making the cluster cohesive as a whole; 2) Positive pairs generated with Eq. (2) have appropriate mutual information. On the one hand, positive pairs are random offsets from the anchor image to the perturbation images, which ensures that they share the mutual information from the anchor image. On the other hand, the sampling interval $I \geq 5$ frames in US-4, resulting in low probability for SPG to sample temporarily close $\{\widehat{\mathbf{x}}_i^{(k)}\}_{k=1}^3$ which are too similar.

3.3 Ultrasound Contrastive Learning

The proposed USCL method learns representations not only by the supervision of category labels, but also by maximizing/minimizing agreement between positive/negative pairs as CR. Here, assorted DNNs can be used as backbone f to encode images, where the output representation vectors $\mathbf{r}_{2i-1} = f(\mathbf{x}_i^{(1)})$ and $\mathbf{r}_{2i} = f(\mathbf{x}_i^{(2)})$ are then fed to the following projection head and classifier.

Contrastive Branch. The contrastive branch consists of a projection head g and corresponding contrastive loss. The g is a two layer MLP which nonlinearly maps representations to other feature space for calculating contrastive regularization loss. The mapped vector $\mathbf{z}_i = g(\mathbf{r}_i) = \mathbf{w}_g^{(2)} \sigma(\mathbf{w}_g^{(1)} \mathbf{r}_i)$ is specialized for a contrast, where σ is ReLU activation function and $\mathbf{w}_g = \{\mathbf{w}_g^{(1)}, \mathbf{w}_g^{(2)}\}$ are the weights of g . The contrastive loss is proposed by Sohn [22], which aims at minimizing the distance between positive pairs $\{\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}\}_{i=1}^N$ and maximizing the distance between any negative pair $\{\mathbf{x}_i^{(1/2)}, \mathbf{x}_j^{(1/2)}\}$, $i \neq j$ for CR:

$$\mathcal{L}_{con} = \frac{1}{2N} \sum_{i=1}^N (l(2i-1, 2i) + l(2i, 2i-1)), \quad (3)$$

where

$$l(i, j) = -\log \frac{\exp(s_{i,j}/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{i \neq k} \cdot \exp(s_{i,k}/\tau)}, \quad (4)$$

and

$$s_{i,j} = \mathbf{z}_i \cdot \mathbf{z}_j / (\|\mathbf{z}_i\| \|\mathbf{z}_j\|), \quad (5)$$

where τ is a tuning temperature parameter.

Classification Branch. We use a linear classifier h with weights \mathbf{w}_c to separate sample features linearly in the representation space similar to [9,10]. The classification loss with corresponding one-hot label \mathbf{y}_i is

$$\mathcal{L}_{sup} = \frac{1}{2N} \sum_{i=1}^N (CE(\mathbf{o}_{2i-1}, \mathbf{y}_i) + CE(\mathbf{o}_{2i}, \mathbf{y}_i)), \quad (6)$$

where $\mathbf{o}_i = h(\mathbf{r}_i) = \text{softmax}(\mathbf{w}_c \mathbf{r}_i)$.

Note that USCL is a semi-supervised training method, only contrastive branch works when the framework receives unlabeled data. This semi-supervised design is intuitively simple but effective, which makes it easy to be implemented and has great potential to be applied to various pretraining scenarios.

4 Experiment

4.1 Experimental Settings

Pretraining Details. ResNet18 is chosen as a representative backbone. We use US-4 dataset (the ratio of training set to validation set is 8 to 2) with 1% labels for pretraining, and fine-tune pretrained models for various downstream tasks. During pretraining, US images are randomly cropped and resized to 224×224 pixels as the input, followed by random flipping and color jittering. We use Adam optimizer with learning rate 3×10^{-4} and weight decay rate 10^{-4} to optimize network parameters. The backbones are pretrained on US-4 for 300 epochs with batch size $N = 32$. The pretraining loss is the sum of contrastive loss and standard cross-entropy loss for classification. Like SimCLR, the backbones are used for fine-tuning on target tasks, projection head g and classifier h are discarded when the pretraining is completed. The λ in Eq. (1) is 0.2, parameters α and β in Eq. (2) are 0.5 and 0.5, respectively. The temperature parameter τ in Eq. (4) is 0.5. The experiments are implemented using PyTorch with an Intel Xeon Silver 4210R CPU@2.4GHz and a single Nvidia Tesla V100 GPU.

Fine-tuning Datasets. We fine-tuned the last 3 layers of pretrained backbones on POCUS [2] and UDIAT-B [26] datasets to testify the performance of our USCL. On POCUS and UDIAT-B datasets, the learning rates are 0.01 and 0.005, respectively. The POCUS is a widely used lung convex probe US dataset for COVID-19 consisting of 140 videos, 2116 images from three classes (*i.e.*, COVID-19, bacterial pneumonia and healthy controls). The UDIAT-B consists of 163 linear probe US breast images from different women with the mean image size of 760×570 pixels, where each of the images presents one or more lesions. Within the 163 lesion images, 53 of them are cancerous masses and other 110 are benign lesions. In this work, we use UDIAT-B dataset to perform the lesion detection and segmentation comparison experiments. 50 of 163 images are used for validation and the rest are used for training.

Table 2. Ablation study of two contrastive ingredients during pretraining: assigning a negative pair from samples of different videos to overcome similarity conflict (I_1) and using mixup operation to enrich the features of positive pairs (I_2). They both improve the model transfer ability significantly, and the classification brunch is also beneficial. All results are reported as POCUS fine-tuning accuracy.

Method	ImageNet	USCL	I_1 I_2 CE loss		✓	✓	✓	✓
Accuracy (%)	84.2			87.5	90.8	92.3	93.2	94.2

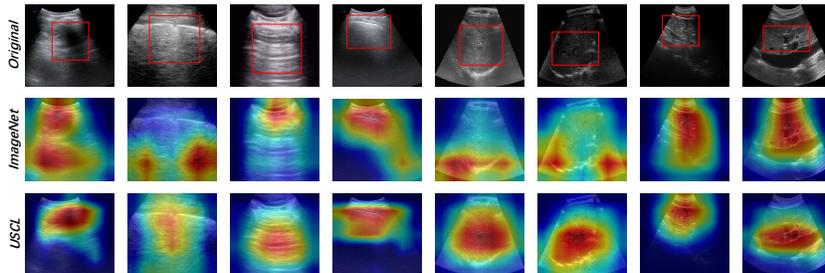


Fig. 4. The visualization results of the last Conv layer in ImageNet pretrained model and USCL model with Grad-CAM [21]. The first 4 columns are lung US images, models trained with USCL on US-4 can focus on the regions of A-line and pleural lesion instead of concentrating on regions without valid information like the ImageNet counterpart. The last 4 columns are liver US images, models trained with USCL accurately attend to the liver regions.

4.2 Ablation Studies

Here, we report the last 3 layers fine-tuning results on POCUS of US-4 pretrained backbones (ResNet18) to validate different components of USCL.

SPG & CE loss. We implement five pretraining methods considering the influence of different contrastive ingredients and the classification brunch with CE loss (Tab. 2). Compared with ImageNet, vanilla contrastive learning improves the accuracy by 3.3% due to a smaller domain gap. It is regarded as the method baseline. The negative pair assigning scheme and positive pair generation method further improves the fine-tuning performance by 3.3% and 4.8%. They can be combined to reach higher performance. In addition, CE loss improves fine-tuning accuracy by 1.0%. This indicates that extra label information is able to enhance the power of contrastive representation learning.

Visualization of Feature Representation. To illustrate the robust feature representation of pretrained backbone, we visualize the last Conv feature map of some randomly selected images produced by USCL pretrained model and ImageNet pretrained model with Grad-CAM [21] (Fig. 4). Compared with ImageNet

Table 3. Comparison of fine-tuning accuracy (%) on POCUS classification dataset and average precision (AP [16])⁸ on UDIAT-B detection (Det), segmentation (Seg) with SOTA methods.

Method	Classification				F1	Det AP	Seg AP
	COVID-19	Pneumonia	Regular	Total Acc			
ImageNet [9]	79.5	78.6	88.6	84.2	81.8	40.6	48.2
US-4 supervised	83.7	82.1	86.5	85.0	82.8	38.3	42.6
TE [20]	75.7	70.0	89.4	81.7	79.0	38.7	46.6
Π Model [20]	77.6	76.4	88.7	83.2	80.6	36.1	45.5
FixMatch [23]	83.0	77.5	85.7	83.6	81.6	39.6	46.9
MoCo v2 [5]	79.7	81.4	88.9	84.8	82.8	38.7	47.1
SimCLR [4]	83.2	89.4	87.1	86.4	86.3	43.8	51.3
USCL	90.8	97.0	95.4	94.2	94.0	45.4	52.8

pretrained backbone, attention regions given by the USCL backbone are much more centralized and more consistent with clinical observation.

4.3 Comparison with SOTA

We compare USCL with ImageNet pretrained ResNet18 [9] and other backbones pretrained on US-4 dataset with supervised method (*i.e.*, plain supervised), semi-supervised methods (*i.e.*, Temporal Ensembling (TE) [20], Π Model [20], FixMatch [23]), and self-supervised methods (*i.e.*, MoCo v2 [5], SimCLR [4]).

Results on Classification Task. On POCUS dataset, we fine-tune the last three layers to testify the representation capability of backbones on classification task (Tab. 3). USCL has consistent best performance on classification of all classes, and its total accuracy of 94.2% is also significantly better than all 7 counterparts. Compared with ImageNet pretrained backbone, USCL reaches a much higher F1 score of 94.0%, which is 12.2% higher.

Results on Detection and Segmentation Tasks. Tab. 3 shows the comparison results of detection and segmentation on UDIAT-B dataset. Mask R-CNN [8] with ResNet18-FPNs [15], whose backbones are pretrained, is used to implement this experiment. USCL generates better backbones than ImageNet and US-4 supervised learning. For detection and segmentation, it outperforms ImageNet pretraining by 4.8% and 4.6%, respectively. Importantly, the UDIAT-B images are collected with linear probe instead of convex probe like US-4, showing a superior texture encoding ability of USCL.

5 Conclusion

This work constructs a new US video-based image dataset US-4 and proposes a simple but efficient contrastive semi-supervised learning algorithm USCL for

⁸ AP is calculated as the area under the precision-recall curve drawn with different Intersection over Union (IoU) thresholds.

US analysis model pretraining. USCL achieves significantly superior performance than ImageNet pretraining by learning compact semantic clusters from US videos. Future works include adding more scan regions of US videos to US-4 dataset for a better generalization on more diseases.

6 Acknowledgement

This work is supported by the Key-Area Research and Development Program of Guangdong Province (2020B0101350001); the GuangDong Basic and Applied Basic Research Foundation (No. 2020A1515110376); Guangdong Provincial Key Laboratory of Big Data Computation Theories and Methods, The Chinese University of Hong Kong (Shenzhen).

References

1. Butterfly videos. <https://www.butterflynetwork.com/index.html>, accessed September 20, 2020
2. Born, J., Wiedemann, N., Cossio, M., Buhre, C., Brändle, G., Leidermann, K., Aujayeb, A., Moor, M., Rieck, B., Borgwardt, K.: Accelerating detection of lung pathologies with explainable ultrasound image analysis. *Applied Sciences* **11**(2), 672 (2021)
3. Celebi, M.E., Aydin, K.: *Unsupervised learning algorithms*. Springer (2016)
4. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. *arXiv:2002.05709* (2020)
5. Chen, X., Fan, H., Girshick, R., He, K.: Improved baselines with momentum contrastive learning. *arXiv:2003.04297* (2020)
6. Gao, L., Zhou, R., Dong, C., Feng, C., Li, Z., Wan, X., Liu, L.: Multi-modal active learning for automatic liver fibrosis diagnosis based on ultrasound shear wave elastography. In: *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*. pp. 410–414. IEEE (2021)
7. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: *CVPR*. pp. 9729–9738 (2020)
8. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. *IEEE TPAMI* **42**(2), 386–397 (2020)
9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *CVPR*. pp. 770–778 (2016)
10. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: *CVPR*. pp. 4700–4708 (2017)
11. Jiao, J., Cai, Y., Alsharid, M., Drukker, L., Papageorghiou, A.T., Noble, J.A.: Self-supervised contrastive video-speech representation learning for ultrasound. In: *MICCAI*. pp. 534–543. Springer (2020)
12. Ke, A., Ellsworth, W., Banerjee, O., Ng, A.Y., Rajpurkar, P.: Chextransfer: Performance and parameter efficiency of imagenet models for chest x-ray interpretation. *arXiv:2101.06871* (2021)
13. Kwitt, R., Vasconcelos, N., Razzaque, S., Aylward, S.: Localizing target structures in ultrasound video—a phantom study. *Medical image analysis* **17**(7), 712–722 (2013)

14. Li, X., Jia, M., Islam, M.T., Yu, L., Xing, L.: Self-supervised feature learning via exploiting multi-modal data for retinal disease diagnosis. *IEEE TMI* **39**(12), 4023–4033 (2020)
15. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: *CVPR*. pp. 2117–2125 (2017)
16. 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: *ECCV*. pp. 740–755 (2014)
17. Liu, L., Lei, W., Wan, X., Liu, L., Luo, Y., Feng, C.: Semi-supervised active learning for covid-19 lung ultrasound multi-symptom classification. In: *2020 IEEE 32nd International Conference on Tools with Artificial Intelligence (ICTAI)*. pp. 1268–1273. *IEEE* (2020)
18. Liu, T., Zhang, H.J., Qi, F.: A novel video key-frame-extraction algorithm based on perceived motion energy model. *IEEE TCSVT* **13**(10), 1006–1013 (2003)
19. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *IJCV* **115**(3), 211–252 (2015)
20. Samuli, L., Timo, A.: Temporal ensembling for semi-supervised learning. In: *International Conference on Learning Representations (ICLR)*. vol. 4, p. 6 (2017)
21. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: *ICCV*. pp. 618–626 (2017)
22. Sohn, K.: Improved deep metric learning with multi-class n-pair loss objective. In: *NeurIPS*. pp. 1857–1865 (2016)
23. Sohn, K., Berthelot, D., Li, C.L., Zhang, Z., Carlini, N., Cubuk, E.D., Kurakin, A., Zhang, H., Raffel, C.: Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *arXiv:2001.07685* (2020)
24. Somphone, O., Allaire, S., Mory, B., Dufour, C.: Live feature tracking in ultrasound liver sequences with sparse demons. In: *MICCAI Workshop*. pp. 53–60 (2014)
25. Vu, Y.N.T., Wang, R., Balachandar, N., Liu, C., Ng, A.Y., Rajpurkar, P.: Medaug: Contrastive learning leveraging patient metadata improves representations for chest x-ray interpretation. *arXiv:2102.10663* (2021)
26. Yap, M.H., Pons, G., Martí, J., Ganau, S., Sentís, M., Zwigelaar, R., Davison, A.K., Martí, R.: Automated breast ultrasound lesions detection using convolutional neural networks. *IEEE journal of biomedical and health informatics* **22**(4), 1218–1226 (2017)