# GCL: Gradient-Guided Contrastive Learning for Medical Image Segmentation with Multi-Perspective Meta Labels

Yixuan Wu
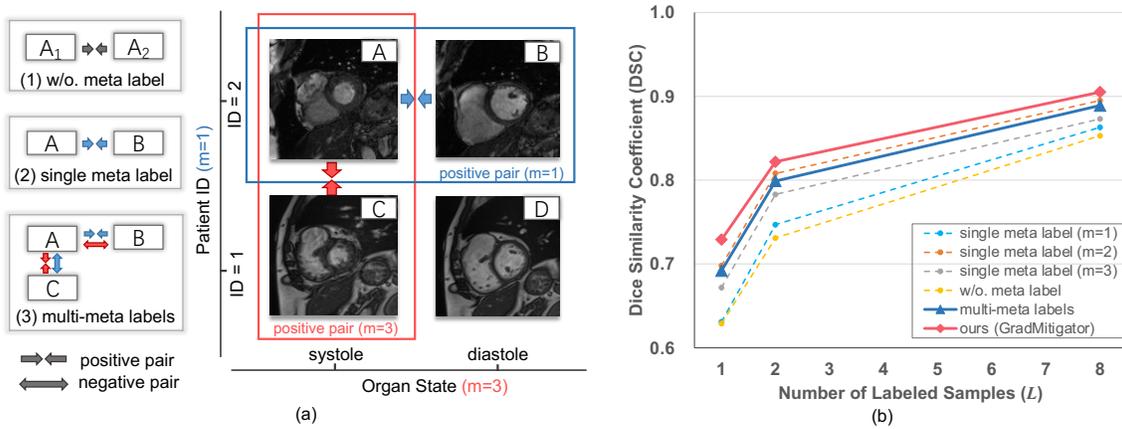Zhejiang University
Hangzhou, China
wyx_chloe@zju.edu.cn

Jintai Chen*
Zhejiang University
Hangzhou, China
jtchen721@gmail.com

Jiahuan Yan
Zhejiang University
Hangzhou, China
jyansir@zju.edu.cn

Yiheng Zhu
Zhejiang University
Hangzhou, China
zhuyiheng2020@zju.edu.cn

Danny Z. Chen
University of Notre Dame
United States
dchen@nd.edu

Jian Wu*
Zhejiang University
Hangzhou, China
wujian2000@zju.edu.cn

Figure 1: Illustrating the "semantic contradiction" problem and its negative effect. (a) Three types of common contrastive learning formulations in medical imaging scenarios: (1) vanilla contrastive learning in which a positive pair is constructed from two augmented versions $(A_1, A_2)$ of one image $A$; (2) a single meta label is used to define additional positive pairs, where images with an identical meta label $m$ are taken as positive pairs; (3) multiple meta labels are leveraged simultaneously to define positive pairs, in which "semantic contradiction" may occur (e.g., images $A$ and $B$ are regarded as both a positive pair and a negative pair simultaneously based on meta labels $m=1$ and $m=3$, respectively). Our novel gradient-guided method GradMitigator mitigates such contradiction. (b) Our preliminary experiments show that directly using multi-perspective meta labels without any additional processing can lead to worse performance (see the blue solid line). Our proposed GradMitigator enables to unify and accumulate positive effects of multi-perspective meta labels (see the red solid line).

## ABSTRACT

Since annotating medical images for segmentation tasks commonly incurs expensive costs, it is highly desirable to design an annotation-efficient method to alleviate the annotation burden. Recently, contrastive learning has exhibited a great potential in learning robust representations to boost downstream tasks with limited labels. In medical imaging scenarios, ready-made meta labels (i.e., specific attribute information of medical images) inherently reveal semantic relationships among images, which have been used to define positive pairs in previous work. However, the multi-perspective semantics revealed by various meta labels are usually incompatible and can incur intractable "semantic contradiction" when combining different meta labels. In this paper, we tackle the issue of "semantic contradiction" in a gradient-guided manner using our proposed *Gradient Mitigator* method, which systematically unifies multi-perspective meta labels to enable a pre-trained model to attain a better high-level semantic recognition ability. Moreover, we emphasize that the fine-grained discrimination ability is vital for segmentation-oriented pre-training, and develop a novel method called *Gradient Filter* to dynamically screen pixel pairs with the most discriminating power based on the magnitude of gradients. Comprehensive experiments on four medical image segmentation datasets verify that our new method GCL: (1) learns informative image representations and considerably boosts segmentation performance with limited labels, and (2) shows promising generalizability on out-of-distribution datasets.

## KEYWORDS

medical pre-training; multi-perspective meta labels; optimization

---

*Co-corresponding authors.

# 1 INTRODUCTION

Cutting-edge medical image segmentation methods usually follow the paradigm of deep learning (DL) based semantic segmentation with a pixel-wise classification process. In this paradigm, pixel-wise annotation is still a big bottleneck due to the labor-intensive and time-consuming burden on medical experts. Moreover, the semantic class of each pixel is predicted independently and pixel correlation is not explicitly specified, and thus a large amount of annotations may be needed to train a comparable model [13].

To reduce the reliance on labeled data, in this paper, we focus on contrastive learning to exploit underlying information of unlabeled data and facilitate informative model initialization for medical image segmentation with limited labels. For better segmentation, in model pre-training, we empower the model with not only recognition ability of high-level semantics (i.e., semantic similarity across the dataset) but also fine-grained discrimination ability for pixel-wise correlation.

For recognition ability, ready-made meta labels (e.g., `Patient_ID`, `Organ_state`) – specific attribute information of different images – are inherently a good source for models to identify semantic similarities between images and learn high-level semantics across a dataset. It was shown [7] that by leveraging the meta labels of slice positions as auxiliary information, contrastive learning could gain more clues to define additional positive pairs. It was verified [38] that the underlying pathology contained in meta labels helps learn image representations in pre-training. However, existing work focused only on utilizing a single meta label while the relationships between different meta labels were not systematically considered and the effects of them were not effectively unified.

When combining multi-perspective meta labels, a natural idea is to treat each meta label independently and sum up the effects of different meta labels directly. But, we observe in preliminary experiments (e.g., see Fig. 1(b)) that combining multiple meta labels without any additional processing may result in worse performance than using a single meta label ($m$=2). Based on this observation, we formulate the "semantic contradiction" problem which is caused by incompatible semantics revealed by different meta labels. For example, as shown in Fig. 1(a), images $A$ and $B$ both are from the same patient but present different organ states (i.e., $A_{\mathsf{Patient\_ID}} = B_{\mathsf{Patient\_ID}}$, $A_{\mathsf{Organ\_state}} \neq B_{\mathsf{Organ\_state}}$). Inspired by multi-objective optimization theory [12, 15, 44, 55], we hypothesize that the contradictory semantics revealed by different meta labels can lead to divergence of model optimization and also to inferior pre-trained representations. In this work, we tackle the issue of "semantic contradiction" in a gradient-guided manner using our proposed **Grad**ient **Mitigator** (GradMitigator), a gradient modifying method that systematically unifies positive effects of various meta labels and hence improves the model's optimization trajectory.

On the other hand, for a better pixel-level discrimination ability, the pre-training should be conducted to distinguish pixel-wise correlation for better detail-aware representations. Deviating from the common practice of pre-defining sub-image positive pairs based on physical coordinates or additional annotations, we utilize high-level semantics (i.e., image-wise semantic similarities) to first initialize a pool of positives for reserving potential positive pixels, from which

optimal positives are dynamically screened to update the model with our proposed **Grad**ient **Filter** (GradFilter) method. Specifically, we define *uncertainty* and *hardness* as two sampling criteria, which are both characterized based on the magnitude of gradients. In this way, it is only the reliable and discriminating pixel pairs that are included for optimizing the pre-trained model.

Based on the above key components, we develop a new overall method GCL (Gradient-guided Contrastive Learning). The main contributions of this work are as follows.

- We exploit multi-perspective meta labels to empower the model with recognition ability for high-level semantics, by mitigating the "semantic contradiction" between meta labels in a gradient-guided fashion.
- We extend the operating granularity of pre-training to the pixel level, where pixel-wise correlation is utilized to increase the model's fine-grained discrimination ability. Specifically, we develop a new GradFilter method to dynamically screen discriminating pixel pairs.
- Our experiments on various medical image segmentation tasks show that, by focusing on both high-level semantics and fine-grained details, our GCL method effectively reduces the downstream model's reliance on labeled data and outperforms known related methods.

# 2 RELATED WORKS

## 2.1 Contrastive Learning

Contrastive learning was first proposed as an instance discrimination task [45], which aims to learn a representation space where similar instances (e.g., images) are pulled closer and different instances are pushed away. In [10, 11, 14, 16, 18, 42, 56], a positive pair was constructed by two augmented versions of one image using transformations (e.g., crop, blur, and color transformations), while a negative pair was constructed by any two different images. To construct suitable positive and negative pairs for better representation learning, some recent work explored optimal combinations of transformations [35, 36, 42] for positive pairs, while other work designed interesting sampling [32] or generating [21] strategies for negative pairs. In [39], alignment and uniformity were identified as two key properties relevant to contrastive learning, and considerable work sought to optimize these two properties. However, mainstream contrastive formulations share two common drawbacks: (1) the criterion for positives and negatives is one-sided, which ignores semantic relationships between images, resulting in models' poor recognition ability for high-level semantics across the dataset; (2) the contrasting granularity is usually restricted to the image level while pixel correlations are overlooked, leading to inferior fine-grained discriminating ability of pre-trained models.

## 2.2 Contrastive Learning for Medical Data

Prior work tried to use domain-specific knowledge to construct better image representations when applying contrastive learning to medical scenarios [7–9, 17, 20, 29, 38, 40, 46, 48, 50–54]. Radiomics features were exploited as knowledge-augmentation to construct additional positive pairs for abnormality classification and localization in chest X-ray images [17]. In [7], it was shown that by leveraging 2D slice positions, contrastive learning based

pre-training could gain more clues to define additional positive pairs and the encoded image representations performed better on downstream tasks. The importance of individual images was dynamically adapted in the contrastive loss to boost performance [29]. In [38], it verified that the underlying pathology contained in meta labels helps learn pre-trained representations, and also compared the effects of different meta labels. However, the utilized domain-specific knowledge focuses only on specific features of medical datasets in a one-sided manner, while information from different perspectives is not systematically combined to characterize the overall dataset.

## 2.3 Pixel-wise Contrastive Learning

Some recent work realized that image-wise contrastive learning is classification-oriented, and extended the operating granularity from the original image level to sub-image level. Different ways to define positive pairs have been proposed. In [7], the same pixel entity after different augmentations is used to form positive pairs. Spatial transformations were leveraged as a prior to deduce location relations between two augmented views, and then matched pixel pairs were formed [46]. In [47], positive pixel pairs were selected based on spatial proximity of physical coordinates. In [31], an information-guided pixel augmentation strategy was proposed to achieve unsupervised local feature matching. Besides, fully-supervised [41] and semi-supervised [1, 8, 19, 25, 40, 52, 53, 60] settings were considered respectively, and external ground truth labels were utilized to construct positive pairs.
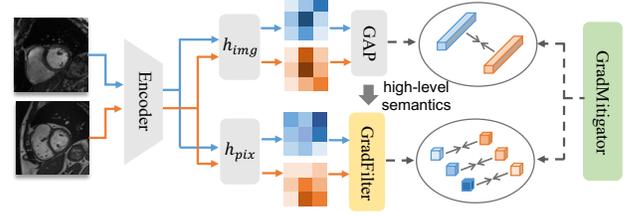
## 3 METHODOLOGY

### 3.1 Contrastive Learning with Meta Labels

Given a mini-batch of $N$ unlabeled 2D images, contrastive learning aims to learn a feature extractor $f(\cdot)$ and a projection head $h(\cdot)$ to yield an image-wise representation $z = h(f(x))$ for each 2D image $x$, by pulling the representations of similar image pairs (i.e., positive pairs) together. In vanilla contrastive learning [10], only two augmented versions of an image are regarded as a positive pair, while any two different images are taken as a negative pair even if they are semantically similar.

In order to empower a pre-trained model with recognition ability for high-level semantics, inspired by [7], we leverage the pre-specified meta labels of medical images to define additional positive pairs. **Note that such meta labels are given for free during the acquisition process of medical datasets, which reveal specific attribute information of the images** (see Sec. 4.2 for illustrations). Specifically, assume that each 2D image $x_i$ has $M$ kinds of meta labels (e.g., Patient_ID, Organ_state), denoted as $y_i^m \in \{1, \ldots, C_m\}$, where $C_m$ is the class number of the meta label $m \in \{1, \ldots, M\}$. Correspondingly, the image-wise contrastive loss guided by meta label $m$ is defined as:

$$\mathcal{L}_{\text{img}}^m = -\frac{1}{|\mathcal{P}_i^m|} \sum_{j \in \mathcal{P}_i^m} \log \frac{\exp(z_i \cdot z_j / \tau)}{\sum_{a=1}^{2N} \mathbf{I}_{i \neq a} \exp(z_i \cdot z_a / \tau)}, \quad (1)$$

where $z_i$ and $z_j$ are the representations of the anchor image and its positive respectively in the image-wise representation space ($z_i = h_{\text{img}}(f(x_i))$, $h_{\text{img}}$ projects features to image-wise representation



Figure 2: The pipeline of GCL. For two images sharing the same meta label $m$, image head $h_{img}$ and pixel head $h_{pix}$ adopt the output from the same encoder as input and project it to their own representation spaces. Image-wise branch employs global average pooling (GAP) to get global features and contrasts them to learn the high-level semantics. GradFilter in pixel-wise branch utilizes the learned high-level semantics from image-wise branch and dynamically screens discriminating pixel pairs. GradMitigator is applied for both image-wise and pixel-wise contrasts to alleviate conflicts between different meta labels (only one meta label is illustrated in figure for simplicity).

space), $\mathcal{P}_i^m$ is the set of indices $j$ of positives $z_j$, $z_a$ denotes each representation of all augmented images in the current mini-batch except the anchor itself (including positives and negatives), and $\tau$ is a temperature parameter.

Thus, the positives come from two sources: (1) the augmented versions of the same image; (2) the different images that have the same class of the meta label $m$ considered.

### 3.2 Gradient Mitigator

To combine the multi-perspective meta labels, a direct way is to simply sum up the contrastive losses guided by all the $M$ meta labels, training jointly to minimize the average loss and update model parameters $\theta$ by:

$$\theta^* = \arg \min_{\theta \in \mathbf{R}} \left( \frac{1}{M} \sum_{m=1}^{M} \mathcal{L}^m(\theta) \right). \quad (2)$$

But, this can incur the "semantic contradiction" issue since high-level semantics revealed by different meta labels may be incompatible. In Fig. 1(a), images $A$ and $B$ are taken from the same patient but present different organ states, and thus are regarded as a positive pair and a negative pair in computing contrastive losses guided by Patient_ID and Organ_state, respectively. We hypothesize that such contradiction can lead to divergence of model optimization and also to inferior pre-trained representations. Poor performance when combining multiple meta labels without any additional processing in preliminary experiments (see Fig. 1(b)) verifies our hypothesis.

**We aim to tackle the issue of "semantic contradiction" in a gradient-guided manner,** given the learning process of DL networks is dictated by gradients with respect to network parameters ($\theta$) – usually back-propagated in the network during gradient descent [34, 49]. Specifically, let $\mathbf{g}_m = \nabla_\theta \mathcal{L}^m(\theta)$ denote an individual gradient guided by meta label $m$, and $\mathbf{g} = \nabla_\theta \mathcal{L}(\theta) = \frac{1}{M} \sum_{m=1}^{M} \mathbf{g}_m$ be the average gradient. With a learning rate $\alpha$, $\theta \leftarrow \theta - \alpha \mathbf{g}$ gives the steepest descent update when optimizing Eq. (2). However, if the

individual gradient $\mathbf{g}_m$ conflicts with $\mathbf{g}$, following Eq. (2) directly will interfere the optimization trajectory guided by meta label $m$.

Thus, we propose the novel GradMitigator method to mitigate the gradient interference by modifying conflicting gradients of different meta-level contrastive losses. As shown in Fig. 3, we study three types of gradient relationships based on cosine similarity $\omega_{ij}$ between meta-level gradients $\mathbf{g}_i$ and $\mathbf{g}_j$: (a) non-conflicting (i.e., $\omega_{ij}=1$); (b) slightly-conflicting (i.e., $0 \le \omega_{ij} < 1$); (c) conflicting (i.e., $\omega_{ij} < 0$).

---

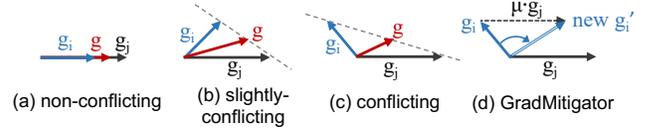**Algorithm 1** The updating process with Gradient Mitigator

---

**Require:** Model parameters $\theta$, meta labels $m \in \{1, \ldots, M\}$, loss functions $\mathcal{L}^m(\theta)$, and EMA weight $\beta$

1: Initialize time-step $t = 0$, EMA variable $\hat{\omega}_{ij}^{(0)} = 0, \forall i, j$
2: Compute $\mathbf{g}_m \leftarrow \nabla_\theta \mathcal{L}^m(\theta), \forall m$
3: **for** $i \in \{1, \ldots, M\}$ **do**
4:     Set $\mathbf{g}_i' \leftarrow \mathbf{g}_i$
5:     **for** $j \in \{1, \ldots, M\} \setminus \{i\}$ **do**
6:         Compute $\omega_{ij}^{(t)} \leftarrow \frac{\mathbf{g}_i' \cdot \mathbf{g}_j}{\|\mathbf{g}_i'\|\|\mathbf{g}_j\|}$
7:         Update $\hat{\omega}_{ij}^{(t)} \leftarrow (1-\beta)\hat{\omega}_{ij}^{(t-1)} + \beta\omega_{ij}^{(t)}$
8:         **if** $\omega_{ij}^{(t)} < \hat{\omega}_{ij}^{(t)}$ **then**
9: 
$$\mathbf{g}_i' = \mathbf{g}_i' + \frac{\|\mathbf{g}_i'\|(\hat{\omega}_{ij}^{(t)}\sqrt{1-\left(\omega_{ij}^{(t)}\right)^2} - \omega_{ij}^{(t)}\sqrt{1-\left(\hat{\omega}_{ij}^{(t)}\right)^2})}{\|\mathbf{g}_j\|\sqrt{1-\left(\hat{\omega}_{ij}^{(t)}\right)^2}} \cdot \mathbf{g}_j$$

10:         **end if**
11:     **end for**
12: **end for**
13: Update $\Delta\theta \leftarrow \mathbf{g}' = \frac{1}{M}\sum_{i=1}^{M}\mathbf{g}_i'$
14: Update time-step $t \leftarrow t + 1$

---

**The goal of GradMitigator is to softly eliminate conflicting components of gradients, seeking agreement between the individual meta-level gradients.** Alg. 1 presents the updating process. We first initialize the target cosine similarity $\hat{\omega}_{ij}^{(0)}$ as 0, and pre-compute all the gradients $\mathbf{g}_m$ of contrastive functions $\mathcal{L}^m(\theta)$ guided by different meta labels $m$. We use $i, j$ as two meta labels for illustrating an updating process: At the current time-step $t$, if the computed cosine similarity between two gradients, i.e., $\omega_{ij}^{(t)} = \frac{\mathbf{g}_i \cdot \mathbf{g}_j}{\|\mathbf{g}_i\|\|\mathbf{g}_j\|}$, is smaller than the target value $\hat{\omega}_{ij}^{(t)}$, we modify one gradient $\mathbf{g}_i$ by injecting a weighted component of the other gradient $\mathbf{g}_j$, i.e., $\mathbf{g}_i' = \mathbf{g}_i + \mu \cdot \mathbf{g}_j$, such that the resulting cosine similarity softly matches the target value $\hat{\omega}_{ij}^{(t)}$. Based on the Law of Sines, the weight $\mu$ is computed. This modifying process is described as:

$$\mathbf{g}_i' = \mathbf{g}_i + \frac{\|\mathbf{g}_i\|(\hat{\omega}_{ij}^{(t)}\sqrt{1-\left(\omega_{ij}^{(t)}\right)^2} - \omega_{ij}^{(t)}\sqrt{1-\left(\hat{\omega}_{ij}^{(t)}\right)^2})}{\|\mathbf{g}_j\|\sqrt{1-\left(\hat{\omega}_{ij}^{(t)}\right)^2}} \cdot \mathbf{g}_j. \quad (3)$$

Note that the value of the target cosine similarity is not fixed. Instead, we use the exponential moving average (EMA) for: (1) avoiding drastic change of the target value during training, and (2)



**Figure 3: Diagrams for illustrating conflicting gradients and our proposed GradMitigator method. Blue, black, and red arrows represent meta-level gradients $\mathbf{g}_i$, $\mathbf{g}_j$, and averaged gradient $\mathbf{g}$, respectively. (a)-(c) Three types of gradient relationships. (d) Our GradMitigator method aims to modify gradient $\mathbf{g}_i$ by injecting a weighted component of $\mathbf{g}_j$ to mitigate the gradient interference.**

bootstrapping for a potentially better target value in a self-adapting manner. This is why we call it 'softly', as:

$$\hat{\omega}_{ij}^{(t)} = (1-\beta)\hat{\omega}_{ij}^{(t-1)} + \beta\omega_{ij}^{(t)}. \quad (4)$$

In this manner, a new average gradient $\mathbf{g}'$ is obtained to update the model parameters $\theta$ in practice.
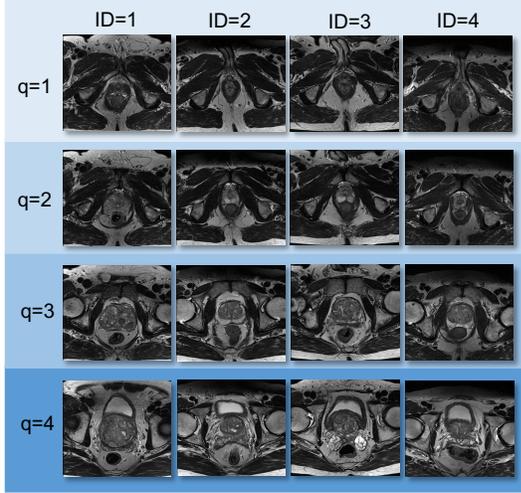
### 3.3 Gradient Filter

It is important to note that the image-level contrastive learning mentioned above is classification task-oriented, which often uses a pooling-like operation to aggregate features from all spatial locations to obtain an image-wise representation. In such a situation, pixel correlation is not explicitly concerned, restricting the pre-trained model's fine-grained discrimination ability, especially for segmentation tasks. In this work, we extend the operating granularity to the pixel level. Existing work usually defines pixel-wise positive pairs based on: (i) the same pixel entity (after different augmentations), (ii) corresponding physical coordinates, and (iii) additional ground truth labels. However, (i) restricts the source of positive pairs to the same image; (ii) largely relies on the assumption that different images in the dataset are well aligned and registered; (3) requires a large number of pixel-wise annotations. All these are not practical in real medical scenarios.

Instead, we propose to utilize the learned high-level semantics between images to pre-define a pool of positives, and dynamically introduce optimal positives from this pool to update the model with our devised Gradient Filter method.

First, for each pixel, we build its pool of positives based on pixel affinity. The pixel affinity $\mathcal{A}$ is computed based on corresponding features in the image-wise representation space (before the pooling-like operation). Suppose pixel $i(u)$ in image $x_i$ is an anchor pixel, its Top-$K$ similar pixels $j(v)$ in image $x_j$ are formed as the positive pool $\mathcal{P}_u$ for pixel $i(u)$, and all the remaining pixels in image $x_j$ are regarded as negatives $\mathcal{N}_u$. Note that the anchor pixel and its positives/negatives are not restricted to being from the same image (after various augmentations); instead, image $x_j$ provides all the image-wise positives of $x_i$ defined in Sec. 3.1 guided by a meta label $m$. Our pixel-wise contrastive loss is defined as:

$$\mathcal{L}_{\text{pix}}^m = -\frac{1}{|\mathcal{P}_u|}\sum_{u^+ \in \mathcal{P}_u}\log\frac{\exp\left(u \cdot u^+/\tau\right)}{\exp(u \cdot u^+/\tau) + \sum_{u^- \in \mathcal{N}_u}\exp(u \cdot u^-/\tau)}, \quad (5)$$

Figure 4: Examples of 2D slices taken from different quantiles ($q = 1, 2, 3, 4$) for four patients on Prostate dataset. One can see that the slices with the same quantile from different patients contain relatively similar anatomical structures.
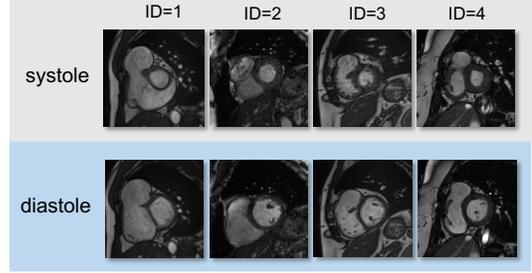
$$\mathcal{A}_{i(u)j(v)} = \frac{z_{i(u)} \cdot z_{j(v)}}{\|z_{i(u)}\| \|z_{j(v)}\|}, \tag{6}$$

where $u, u^+$, and $u^-$ denote respectively the anchor pixel, its positive and negative in the pixel-wise representation space (i.e., $u = h_{\text{Pix}}(f(x_{i(u)}))$, where $h_{\text{Pix}}$ projects features to the pixel-wise representation space), and $z_{i(u)}$ and $z_{j(v)}$ are corresponding features of pixels $i(u)$ and $j(v)$ in the image-wise representation space (e.g., $z_{i(u)} = h_{\text{img}}(f(x_{i(u)}))$).

Next, to further enhance the effectiveness of the defined positives, we consider that an "ideal" positive should both be reliable (i.e., with a low *uncertainty*) for the right optimization direction of the model and have a certain degree of *hardness* for constantly optimizing the model's decision boundary [32]. Hence, we propose GradFilter for screening positive pixels with a high discriminating power based on two criteria, *uncertainty* and *hardness*, from the pre-defined positive pool to update the model. Since DL networks are optimized using gradient-based methods, we characterize these two criteria by gradient magnitudes induced by different positives. The gradient of the pixel-wise contrastive loss w.r.t. the anchor pixel representation can be described as:

$$\frac{\partial \mathcal{L}_{\text{pix}}^m}{\partial u} = -\frac{1}{\tau|\mathcal{P}_u|} \sum_{u^+ \in \mathcal{P}_u} \left( \frac{(u^+ - u^-) \sum\limits_{u^- \in \mathcal{N}_u} \exp(u \cdot u^-/\tau)}{\exp(u \cdot u^+/\tau) + \sum\limits_{u^- \in \mathcal{N}_u} \exp(u \cdot u^-/\tau)} \right). \tag{7}$$

One may see that a harder positive usually has a smaller dot product with the anchor, which brings more gradient contribution than easier positives. Conversely, we consider that the model is more certain about a positive if a smaller gradient is induced and hence a little update is performed at the current optimization direction [3]. With these two criteria, we aim to make reconciliation inspired by Self-Pace Learning [24], **following the learning process of humans, so that the model learns better when feeding samples from easy to hard to it.**



Figure 5: Examples of 2D slices taken in different organ states (i.e., systole and diastole) from four different patients on the ACDC dataset. It can be seen that the appearances of the heart in different states are different.

Specifically, at time-step $t$, for each positive $u^+ \in \mathcal{P}_u$, we compute the corresponding gradient with respect to the parameters of the last layer of the encoder. A pace function $g(t)$ is defined to specify the positive pool size so that only positives with the $g(t)$ lowest gradients are actually used, as:

$$g(t) = \left[1 + \frac{1}{4} \log\left(\frac{t}{T} + e^{-4}\right)\right] \cdot |\mathcal{P}_u|, \tag{8}$$

where $T$ denotes the total number of training steps. This scheme schedules how the positives are introduced to the training process: At the beginning, positives with high certainty are preferred, and as the training progresses, more harder positives are introduced.

### 3.4 Training Objective

As shown in Fig. 2, we jointly perform image-wise contrastive learning (Sec. 3.1) and pixel-wise contrastive learning (Sec. 3.3) under the guidance of a meta label $m$. The effects of all the $M$ multi-perspective meta labels are dynamically unified to optimize the model with modified gradients using our GradMitigator method (Sec. 3.2), by:

$$\mathcal{L}_{\text{GCL}} = \sum_{m=1}^{M} \left( \mathcal{L}_{\text{img}}^m + \mathcal{L}_{\text{pix}}^m \right). \tag{9}$$

## 4 EXPERIMENTS

### 4.1 Experimental Setup

We conduct four sets of experiments, investigating: (1) the informativeness of learned representations compared with other pre-training methods; (2) the effectiveness to reduce downstream task's reliance on labeled data compared with other semi-supervised methods; (3) the generalizability on out-of-distribution datasets; and (4) the effects of each designed components.

### 4.2 Implementations

**Training and Evaluation.** Note that our proposed GCL is a pre-training method. Following [7, 20, 29], the performance of our GCL method is evaluated in a "pre-training and fine-tuning" paradigm: (1) GCL method is used to pre-train a U-Net encoder; and (2) the pre-trained weights are regarded as initialization for the downstream segmentation network to be fine-tuned with limited labels. The performance of our method is indicated by the segmentation accuracy.

**Table 1: Comparison of our proposed GCL method and related pre-training methods on four datasets in DSC performance. $L$ denotes the number of provided labeled samples in fine-tuning. The best results are marked in bold.**

| Method | ACDC | | | Prostate | | | MMWHS | | | ACDC $\rightarrow$ HVSMR | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $L$=1 | $L$=2 | $L$=8 | $L$=1 | $L$=2 | $L$=8 | $L$=1 | $L$=2 | $L$=8 | $L$=2 | $L$=4 | $L$=6 |
| Baseline | | | | | | | | | | | | |
| Random Init. | $0.598_{\pm.023}$ | $0.682_{\pm.012}$ | $0.847_{\pm.008}$ | $0.477_{\pm.039}$ | $0.547_{\pm.027}$ | $0.645_{\pm.013}$ | $0.443_{\pm.017}$ | $0.623_{\pm.012}$ | $0.792_{\pm.008}$ | $0.743_{\pm.037}$ | $0.817_{\pm.025}$ | $0.847_{\pm.017}$ |
| Image-level Contrastive Learning | | | | | | | | | | | | |
| SimCLR | $0.629_{\pm.037}$ | $0.731_{\pm.022}$ | $0.853_{\pm.017}$ | $0.519_{\pm.049}$ | $0.578_{\pm.033}$ | $0.663_{\pm.022}$ | $0.489_{\pm.022}$ | $0.667_{\pm.018}$ | $0.799_{\pm.009}$ | $0.737_{\pm.047}$ | $0.807_{\pm.022}$ | $0.842_{\pm.023}$ |
| MoCo | $0.623_{\pm.032}$ | $0.723_{\pm.024}$ | $0.851_{\pm.012}$ | $0.503_{\pm.047}$ | $0.577_{\pm.032}$ | $0.660_{\pm.027}$ | $0.501_{\pm.018}$ | $0.652_{\pm.020}$ | $0.787_{\pm.012}$ | $0.733_{\pm.048}$ | $0.809_{\pm.029}$ | $0.840_{\pm.025}$ |
| Pixel-level Contrastive Learning | | | | | | | | | | | | |
| PixPro | $0.642_{\pm.025}$ | $0.754_{\pm.022}$ | $0.863_{\pm.020}$ | $0.544_{\pm.029}$ | $0.597_{\pm.024}$ | $0.670_{\pm.022}$ | $0.519_{\pm.024}$ | $0.672_{\pm.017}$ | $0.798_{\pm.017}$ | $0.753_{\pm.033}$ | $0.821_{\pm.033}$ | $0.849_{\pm.018}$ |
| DenseCL | $0.633_{\pm.027}$ | $0.742_{\pm.021}$ | $0.859_{\pm.020}$ | $0.537_{\pm.029}$ | $0.589_{\pm.024}$ | $0.668_{\pm.022}$ | $0.514_{\pm.027}$ | $0.688_{\pm.018}$ | $0.797_{\pm.023}$ | $0.752_{\pm.032}$ | $0.817_{\pm.037}$ | $0.842_{\pm.022}$ |
| PointRC | $0.647_{\pm.023}$ | $0.760_{\pm.017}$ | $0.867_{\pm.013}$ | $0.552_{\pm.028}$ | $0.601_{\pm.022}$ | $0.673_{\pm.018}$ | $0.524_{\pm.025}$ | $0.677_{\pm.010}$ | $0.799_{\pm.014}$ | $0.750_{\pm.038}$ | $0.820_{\pm.031}$ | $0.848_{\pm.012}$ |
| Medical Pre-training | | | | | | | | | | | | |
| GLCL | $0.702_{\pm.020}$ | $0.783_{\pm.015}$ | $0.881_{\pm.011}$ | $0.572_{\pm.029}$ | $0.612_{\pm.018}$ | $0.687_{\pm.022}$ | $0.557_{\pm.015}$ | $0.689_{\pm.014}$ | $0.801_{\pm.005}$ | $0.778_{\pm.033}$ | $0.825_{\pm.023}$ | $0.852_{\pm.022}$ |
| PosiCL | $0.688_{\pm.021}$ | $0.803_{\pm.018}$ | $0.886_{\pm.012}$ | $0.554_{\pm.031}$ | $0.606_{\pm.017}$ | $0.689_{\pm.017}$ | $0.533_{\pm.012}$ | $0.692_{\pm.011}$ | $0.814_{\pm.003}$ | $0.780_{\pm.029}$ | $0.827_{\pm.020}$ | $0.855_{\pm.018}$ |
| SSCL | $0.697_{\pm.017}$ | $0.785_{\pm.011}$ | $0.892_{\pm.009}$ | $0.587_{\pm.023}$ | $0.621_{\pm.012}$ | $0.692_{\pm.009}$ | $0.547_{\pm.017}$ | $0.694_{\pm.009}$ | $0.812_{\pm.004}$ | $0.779_{\pm.030}$ | $0.839_{\pm.022}$ | $0.859_{\pm.019}$ |
| SPL | $0.699_{\pm.023}$ | $0.801_{\pm.012}$ | $0.889_{\pm.014}$ | $0.588_{\pm.023}$ | $0.624_{\pm.017}$ | $0.688_{\pm.023}$ | $0.560_{\pm.014}$ | $0.694_{\pm.013}$ | $0.814_{\pm.009}$ | $0.782_{\pm.030}$ | $0.827_{\pm.022}$ | $0.857_{\pm.019}$ |
| Pretext task Pre-training | | | | | | | | | | | | |
| Rotation | $0.592_{\pm.037}$ | $0.689_{\pm.022}$ | $0.852_{\pm.014}$ | $0.508_{\pm.043}$ | $0.562_{\pm.037}$ | $0.659_{\pm.029}$ | $0.447_{\pm.020}$ | $0.635_{\pm.019}$ | $0.787_{\pm.012}$ | $0.752_{\pm.044}$ | $0.819_{\pm.024}$ | $0.847_{\pm.022}$ |
| Inpainting | $0.605_{\pm.025}$ | $0.701_{\pm.018}$ | $0.863_{\pm.005}$ | $0.503_{\pm.033}$ | $0.557_{\pm.020}$ | $0.668_{\pm.013}$ | $0.463_{\pm.018}$ | $0.650_{\pm.010}$ | $0.792_{\pm.008}$ | $0.757_{\pm.023}$ | $0.820_{\pm.017}$ | $0.849_{\pm.009}$ |
| Jigsaw | $0.595_{\pm.028}$ | $0.712_{\pm.015}$ | $0.874_{\pm.010}$ | $0.501_{\pm.038}$ | $0.570_{\pm.023}$ | $0.667_{\pm.012}$ | $0.437_{\pm.044}$ | $0.642_{\pm.012}$ | $0.793_{\pm.007}$ | $0.748_{\pm.038}$ | $0.823_{\pm.020}$ | $0.852_{\pm.012}$ |
| GCL (ours) | $\mathbf{0.729}_{\pm.014}$ | $\mathbf{0.812}_{\pm.009}$ | $\mathbf{0.905}_{\pm.007}$ | $\mathbf{0.606}_{\pm.019}$ | $\mathbf{0.631}_{\pm.013}$ | $\mathbf{0.701}_{\pm.012}$ | $\mathbf{0.572}_{\pm.010}$ | $\mathbf{0.709}_{\pm.008}$ | $\mathbf{0.819}_{\pm.005}$ | $\mathbf{0.797}_{\pm.022}$ | $\mathbf{0.849}_{\pm.017}$ | $\mathbf{0.870}_{\pm.007}$ |

All the fine-tuning experiments are repeated 5 times. Segmentation results are reported in the form of mean (standard deviation) with the dice similarity coefficient (DSC).

**Data.** We evaluate the performance of our GCL method on four public medical image datasets: ACDC [5], Prostate [2], MMWHS [61, 62], and HVSMR [27]. These four datasets have different anatomical structures (i.e., cardiac, prostate), modalities (i.e., MRI, CT), resolutions, and sizes, allowing comprehensive evaluation of our method. For ACDC, we leverage the meta labels of `Patient_ID` ($m$=1), `Slice_quantile` ($m$=2) (i.e., the quantile of a 2D image along one axis), and `Organ_state` ($m$=3) (i.e., systole or diastole). For the other datasets, the meta labels of `Patient_ID` and `Slice_quantile` are used. Each dataset is split into a pre-training set $X_{ptr}$ and a fine-tuning set $X_{ft}$; the fine-tuning set $X_{ft}$ is further split into a training set $X_{tr}$, a validation set $X_{val}$, and a test set $X_{ts}$. We use $X_{ptr}$ to pre-train the GCL model without ground truth labels, and use $X_{ft}$ to fine-tune the pre-trained encoder on the downstream task and report segmentation performance. A small number of samples in $X_{tr}$ are randomly chosen as labeled samples (e.g., $L = 1, 2, 8$).

**Illustration of Meta Labels.** We show 2D images (or slices) with their meta labels in Figs. 4 and 5, in order to provide intuitive illustrations of the usage of different meta labels. `Slice_quantile` represents the quantile of a 2D image/slice along one axis (i.e., the $z$-axis for the Prostate, MMWHS, and HVSMR datasets, and the short axis for the ACDC dataset). Thus, all the slices in a 3D image are divided into four parts, and the corresponding quantile is indicated by $q$ ($q$ = 1, 2, 3, 4). In Fig. 4, we illustrate slices taken from different quantiles, for four different patients. One can see that the slices with the same quantile from different patients contain similar anatomical structures. Besides, `Organ_state` indicates the state of a target organ at the time of scanning (e.g., the systole or diastole state of the heart). It can be seen from Fig. 5 that the appearances of the heart in different states show considerable differences.

**Model Details.** Our contrastive formulation follows [10]. The architecture of the encoder follows U-Net [33]. The two projection heads $h_{\text{img}}$ and $h_{\text{pix}}$ share the same design, which consists of $1 \times 1$ convolution, ReLU, and $1 \times 1$ convolution. The hidden layer's dimension of the projection head is 512, keeping the same as its input channels, and the final output dimension is 128, the same as [10]. In the fine-tuning stage, we employ U-Net as our segmentation network.

**Training Details.** The GCL pre-training is performed on four NVIDIA GeForce RTX 3090 GPUs. We train with the SGD optimizer [6] for 300 epochs, and the cosine learning rate scheduler is adopted, with a batch size of 48 and a learning rate of 0.1. In the fine-tuning stage, we train the segmentation network with limited labels for 300 epochs. The Adam optimizer [22] and cosine learning rate scheduler are used, with a batch size of 5 and a learning rate of $10^{-4}$. The temperature $\tau$ is set to 0.1 following [10]. $K$ is set to 0.3 when defining the positive pool. The EMA weight $\beta$ is set to $10^{-2}$. All the parameters $\theta$ of the encoder are updated individually based on the modified gradients when applying the proposed GradMitigator method.

### 4.3 Comparison with Pre-training Methods

To evaluate the informativeness of learned image representations, we compare our GCL with several groups of pre-training methods. (1) Image-level contrastive learning, including SimCLR [10], MoCo [18]. (2) Pixel-level contrastive learning, including PixPro [47], DenseCL [43], and PointRC [4]. (3) Medical pre-training, including GLCL [7], PosiCL [57], SSCL [20], and SPL [29]. (4) Pretext task pre-training, including Rotation [23], Inpainting [28], and Jigsaw [26].

**Results.** Table 1 summarizes the results on four downstream segmentation tasks, which are used to indicate the pre-training performance. $L$ denotes the number of provided labeled samples in

**Table 2: Comparison of our proposed GCL method and semi-supervised methods on three datasets with limited labeled data provided in DSC performance.** $L$ **denotes the number of provided labeled samples. The best results are marked in bold.**

| Method | ACDC | | | Prostate | | | MMWHS | | |
|---|---|---|---|---|---|---|---|---|---|
| | $L$=1 | $L$=2 | $L$=8 | $L$=1 | $L$=2 | $L$=8 | $L$=1 | $L$=2 | $L$=8 |
| Baseline | $0.598_{\pm.023}$ | $0.682_{\pm.012}$ | $0.847_{\pm.008}$ | $0.477_{\pm.039}$ | $0.547_{\pm.027}$ | $0.645_{\pm.013}$ | $0.443_{\pm.017}$ | $0.623_{\pm.012}$ | $0.792_{\pm.008}$ |
| Adv. Training | $0.662_{\pm.012}$ | $0.749_{\pm.013}$ | $0.849_{\pm.007}$ | $0.544_{\pm.023}$ | $0.587_{\pm.020}$ | $0.681_{\pm.010}$ | $0.531_{\pm.018}$ | $0.679_{\pm.012}$ | $0.790_{\pm.007}$ |
| Mean Teacher | $0.674_{\pm.012}$ | $0.771_{\pm.007}$ | $0.857_{\pm.004}$ | $0.526_{\pm.014}$ | $0.557_{\pm.008}$ | $0.657_{\pm.004}$ | $0.538_{\pm.012}$ | $0.692_{\pm.008}$ | $0.809_{\pm.004}$ |
| Mixup | $0.667_{\pm.009}$ | $0.773_{\pm.010}$ | $0.862_{\pm.005}$ | $0.531_{\pm.017}$ | $0.598_{\pm.012}$ | $0.677_{\pm.008}$ | $0.549_{\pm.014}$ | $0.683_{\pm.013}$ | $0.797_{\pm.009}$ |
| GCL (ours) | $0.729_{\pm.014}$ | $0.812_{\pm.009}$ | $0.905_{\pm.007}$ | $0.606_{\pm.019}$ | $0.631_{\pm.013}$ | $0.701_{\pm.012}$ | $0.572_{\pm.010}$ | $0.709_{\pm.008}$ | $0.819_{\pm.005}$ |
| Ours + Mixup | $0.754_{\pm.007}$ | $0.833_{\pm.008}$ | $0.911_{\pm.002}$ | $0.612_{\pm.012}$ | $0.633_{\pm.014}$ | $\mathbf{0.713_{\pm.002}}$ | $0.621_{\pm.007}$ | $0.722_{\pm.008}$ | $0.822_{\pm.003}$ |
| Ours + M.T. | $\mathbf{0.762_{\pm.010}}$ | $\mathbf{0.842_{\pm.008}}$ | $\mathbf{0.913_{\pm.002}}$ | $\mathbf{0.614_{\pm.014}}$ | $\mathbf{0.643_{\pm.008}}$ | $0.711_{\pm.005}$ | $\mathbf{0.628_{\pm.008}}$ | $\mathbf{0.731_{\pm.005}}$ | $\mathbf{0.828_{\pm.005}}$ |
| Fully Superv. | $0.914_{\pm.003}$ ($L$=50) | | | $0.703_{\pm.005}$ ($L$=18) | | | $0.812_{\pm.009}$ ($L$=10) | | |

**Table 3: Ablation study with different components of our GCL method on three datasets in DSC performance.** $L$ **denotes the number of labeled samples in fine-tuning. The best results are marked in bold.**

| Method | ACDC | | | Prostate | | | MMWHS | | |
|---|---|---|---|---|---|---|---|---|---|
| | $L$=1 | $L$=2 | $L$=8 | $L$=1 | $L$=2 | $L$=8 | $L$=1 | $L$=2 | $L$=8 |
| Base | $0.629_{\pm.037}$ | $0.731_{\pm.022}$ | $0.853_{\pm.017}$ | $0.519_{\pm.049}$ | $0.578_{\pm.033}$ | $0.663_{\pm.022}$ | $0.489_{\pm.022}$ | $0.667_{\pm.018}$ | $0.799_{\pm.009}$ |
| Base + single meta label ($m$=1) | $0.631_{\pm.024}$ | $0.747_{\pm.019}$ | $0.863_{\pm.018}$ | $0.527_{\pm.032}$ | $0.589_{\pm.020}$ | $0.678_{\pm.022}$ | $0.529_{\pm.020}$ | $0.684_{\pm.012}$ | $0.799_{\pm.002}$ |
| Base + single meta label ($m$=2) | $0.688_{\pm.020}$ | $0.784_{\pm.013}$ | $0.880_{\pm.012}$ | $0.573_{\pm.027}$ | $0.617_{\pm.019}$ | $0.688_{\pm.024}$ | $0.549_{\pm.013}$ | $0.683_{\pm.014}$ | $0.801_{\pm.007}$ |
| Base + single meta label ($m$=3) | $0.672_{\pm.018}$ | $0.783_{\pm.018}$ | $0.873_{\pm.009}$ | | / | | | / | |
| Base + multi-meta labels | $0.681_{\pm.022}$ | $0.784_{\pm.020}$ | $0.878_{\pm.017}$ | $0.574_{\pm.023}$ | $0.614_{\pm.018}$ | $0.685_{\pm.012}$ | $0.548_{\pm.023}$ | $0.684_{\pm.021}$ | $0.802_{\pm.014}$ |
| Base + multi-meta labels + GradMiti. | $0.709_{\pm.019}$ | $0.804_{\pm.018}$ | $0.894_{\pm.012}$ | $0.592_{\pm.030}$ | $0.623_{\pm.028}$ | $0.692_{\pm.014}$ | $0.559_{\pm.017}$ | $0.697_{\pm.013}$ | $0.811_{\pm.008}$ |
| Base + PixCL | $0.635_{\pm.024}$ | $0.747_{\pm.022}$ | $0.860_{\pm.018}$ | $0.545_{\pm.028}$ | $0.593_{\pm.027}$ | $0.663_{\pm.028}$ | $0.513_{\pm.022}$ | $0.685_{\pm.017}$ | $0.793_{\pm.014}$ |
| Base + PixCL + GradFilter | $0.652_{\pm.020}$ | $0.775_{\pm.019}$ | $0.871_{\pm.014}$ | $0.562_{\pm.024}$ | $0.612_{\pm.028}$ | $0.679_{\pm.017}$ | $0.536_{\pm.019}$ | $0.694_{\pm.016}$ | $0.803_{\pm.009}$ |
| Full model (frozen) | $0.709_{\pm.011}$ | $0.797_{\pm.007}$ | $0.882_{\pm.006}$ | $0.593_{\pm.014}$ | $0.617_{\pm.012}$ | $0.684_{\pm.008}$ | $0.563_{\pm.011}$ | $0.694_{\pm.008}$ | $0.809_{\pm.006}$ |
| Full model (ours) | $0.729_{\pm.014}$ | $0.812_{\pm.009}$ | $0.905_{\pm.007}$ | $0.606_{\pm.019}$ | $0.631_{\pm.013}$ | $0.701_{\pm.012}$ | $0.572_{\pm.010}$ | $0.709_{\pm.008}$ | $0.819_{\pm.005}$ |
| Full model (merged data) | $\mathbf{0.734_{\pm.012}}$ | $\mathbf{0.817_{\pm.009}}$ | $\mathbf{0.908_{\pm.004}}$ | $\mathbf{0.614_{\pm.014}}$ | $\mathbf{0.634_{\pm.016}}$ | $\mathbf{0.702_{\pm.009}}$ | $\mathbf{0.579_{\pm.008}}$ | $\mathbf{0.714_{\pm.007}}$ | $\mathbf{0.822_{\pm.005}}$ |

fine-tuning. As baselines, we train the downstream segmentation network with random initialization (train from scratch). One can see that the pretext task pre-training (i.e., Rotation, Inpainting, and Jigsaw) provides less informative initialization, which shows worse performance when labeled data is extremely limited (i.e., $L = 1, 2$). Besides, the general image-level contrastive learning methods (i.e., SimCLR and MoCo) provide useful initialization to some extent. And when extending the contrasting granularity to the pixel level (i.e., PixPro, DenseCL, and PointRC), the pre-training performance gets further boosted. In addition, the pre-training methods designed in medical scenarios (i.e., GLCL, PosiCL, SSCL, and SPL) perform better than those general pre-training methods, suggesting that single-source domain-specific information in medical images provides useful clues to some extent. Yet, our GCL still boosts the downstream segmentation accuracy to a large extent (e.g., 0.131, 0.129, 0.129, and 0.054 in DSC on the ACDC, Prostate, MMWHS, and HVSMR dataset, when $L$=1, respectively). This is because our GCL method effectively unifies the information from multi-perspective meta labels. Beyond that, it can be seen that the fewer labeled samples (i.e., $L = 1, 2$) provided in fine-tuning, the more significant the superiority of our GCL pre-training method. This validates the necessity of pre-training for segmentation in medical scenarios where limited labeled data can be provided.

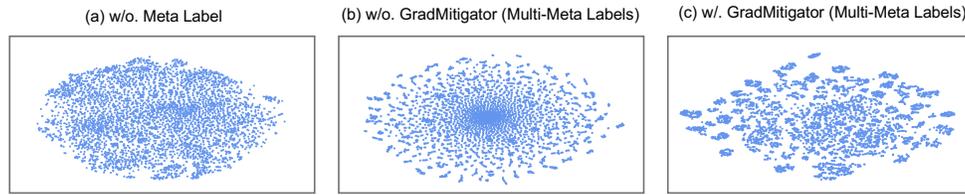## 4.4 Comparison with Semi-supervised Methods

To further investigate the effectiveness of our GCL to reduce downstream task's reliance on labeled data, we compare it with semi-supervised methods, including Adv. Training [59], Mean Teacher [30], and Mixup [58].

**Results.** In Table 2, the baseline is trained from scratch without other designs, and the segmentation results under full supervision (i.e., all labeled samples in the training set are provided). The improvements of semi-supervised methods (i.e., Adv. Training, Mean Teacher, and Mixup) are limited and largely depend on specific datasets. For example, Adv. Training performs worse on the ACDC dataset but performs better on the Prostate dataset. In contrast, our GCL method performs well on all the datasets (e.g., 0.131, 0.129, 0.129, and 0.054 in DSC on the ACDC, Prostate, MMWHS, and HVSMR dataset, when $L$=1, respectively). Moreover, it can be seen that our method shows promising compatibility with semi-supervised methods (e.g., Mixup and Mean Teacher), leading to further performance gains (e.g., 0.164, 0.137, and 0.185 in DSC on the ACDC, Prostate, and MMWHS dataset, when $L$=1, respectively), and even surpassing the results under full supervision on the Prostate and MMWHS datasets.

## 4.5 Generalizability

To further examine whether the image representations learned by our GCL method have good generalizability, we conduct pre-training and fine-tuning on different datasets. Specifically, our proposed GCL is used to pre-train an encoder on ACDC dataset, and the downstream segmentation network is fine-tuned on HVSMR. **Results.** The last column of Table 1 presents the segmentation performance on the out-of-distribution dataset. The image-level and pixel-level contrastive learning (i.e., SimCLR, MoCo, PixPro, DenseCL, and PointRC) and pretext task pre-training (i.e., Rotation, Inpainting, and Jigsaw) give little performance gain for the downstream task. In contrast, our GCL method shows superiority

**Figure 6: Some t-SNE visualization examples of the learned representations on the ACDC dataset. (a) No meta label is used, and the pre-trained model is degraded to the vanilla contrastive learning, where the learned representations are evenly distributed. (b) Multiple meta labels are simultaneously used without any additional processing, where the learned representations are excessively diffused. (c) With our proposed GradMitigator method, obvious semantic separation is formed.**

compared with the other pre-training methods. This is partially due to the multi-perspective meta labels that are unified to empower the model with recognition ability for high-level semantics, which are commonly found in different medical image datasets.

## 4.6 Ablation Study

The results of ablation study are presented in Table 3. The base model refers to the results that all designs in our GCL method are removed, with vanilla image-level contrastive learning retained.

**Effects of Introducing Meta Labels.** We compare the performance pre-trained with different meta labels one-by-one. The meta labels of `Patient_ID` ($m$=1), `Slice_quantile` ($m$=2), and `Organ_state` ($m$=3, only for the ACDC dataset) are included. One can see that by introducing the meta labels as additional clues, the segmentation performance gets well improved. In particular, the meta label of `Slice_quantile` ($m$=2) provides more gains on all the datasets.

**Effects of Our Gradient Mitigator.** We further investigate the effect of simultaneously using multiple meta labels. It can be seen that when combine all meta labels without any additional processing, the performance is poor and is even worse than using a single meta label ($m$=2). This verifies our hypothesis that the "semantic contradiction" between different meta labels can incur divergence of model optimization. Further, only by modifying the conflicting gradients with our proposed GradMitigator can multi-perspective meta labels be unified to synergistically optimize the pre-training process. Meanwhile, in Fig. 6, we utilize t-SNE [37] to visualize the learned representations on the ACDC dataset. In Fig. 6(b), when introducing multiple meta labels without any additional processing, the learned representations are excessively diffused. Instead, in Fig. 6(c), when applying our proposed GradMitigator, the obvious semantic separation is formed, suggesting the effectiveness of our GCL on exploring the semantic similarity between images and thus capturing high-level semantics across the dataset.

**Effects of Pixel-wise Contrastive Learning and Gradient Filter.** To explore the effects of fine-grained contrasting granularity, we add the pixel-wise contrastive learning (PixCL) component to our base model, where the GradFilter strategy is not applied first (remove the pace function $g(t)$ (*cf.* Eq. (8)), and use all the positives in the pre-defined positive pool to update the model). It can be seen that the PixCL boosts the downstream performance to some extent, indicating that fine-grained contrast is necessary to segmentation-oriented pre-training, which do contribute to the model's recognition ability for segmentation details. Moreover, when further adding

our proposed GradFilter method, the downstream segmentation accuracy gets largely boosted. Therefore, by considering both *uncertainty* and *hardness* to screen optimal pixel-wise positives, our pixel-wise contrastive learning exert its substantial effectiveness.

**The scope of pre-training set.** To further explore the extensibility of our GCL, we merge the pre-training sets of all used datasets (i.e., ACDC, Prostate, and MMWHS) to pre-train the encoder with our full GCL model. The results are shown in the last row of Table 3 (denoted as "merged data"). Compared with using the single dataset to pre-train (denoted as "ours"), incorporating more pre-training data from other datasets (even from different organs and modalities) does contribute to the pre-training performance.

**Effects of fine-tuning the pre-trained parameters.** During the fine-tuning process, we freeze the entire pre-trained encoder and only fine-tune the decoder. The results in Table 3 suggest that the pre-trained encoder has certain feature extraction capability, while fine-tuning can further enhance the performance. Moreover, when an extremely limited number of labeled samples ($L$=1) is provided, our GCL pre-training method plays a more significant role.

## 5 CONCLUSIONS

In this paper, we proposed to systematically unify multi-perspective meta labels without incurring the "semantic contradiction" issue by modifying their corresponding gradients. Further, when extending the contrast granularity to the pixel level, our new Gradient Filter method helps dynamically screen positive pixel pairs with the most discriminating power. Compared to other contrastive formulations, our method empowers the pre-trained model with both recognition ability for high-level semantics and discrimination ability for pixel-wise correlation in a gradient-guided manner. Extensive experiments on four public datasets verified that our GCL method not only learns informative image representations for downstream segmentation with extremely limited labels, but also shows promising generalizability on out-of-distribution datasets.

## REFERENCES

[1] Inigo Alonso, Alberto Sabater, David Ferstl, et al. 2021. Semi-supervised semantic segmentation with pixel-level contrastive learning from a class-wise memory bank. In *ICCV*.

[2] Michela Antonelli, Annika Reinke, Spyridon Bakas, et al. 2022. The medical segmentation decathlon. *Nature Communications* (2022).

[3] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, et al. 2019. Deep batch active learning by diverse, uncertain gradient lower bounds. *ArXiv Preprint ArXiv:1906.03671* (2019).

[4] Yutong Bai, Xinlei Chen, Alexander Kirillov, et al. 2022. Point-level region contrast for object detection pre-training. In *CVPR*.

[5] Olivier Bernard, Alain Lalande, Clement Zotti, et al. 2018. Deep learning techniques for automatic MRI cardiac multi-structures segmentation and diagnosis: Is the problem solved? *TMI* (2018).

[6] Léon Bottou. 2012. Stochastic gradient descent tricks. In *Neural Networks: Tricks of the Trade*.

[7] Krishna Chaitanya, Ertunc Erdil, Neerav Karani, et al. 2020. Contrastive learning of global and local features for medical image segmentation with limited annotations. In *NeurIPS*.

[8] Krishna Chaitanya, Ertunc Erdil, Neerav Karani, et al. 2023. Local contrastive loss with pseudo-label based self-training for semi-supervised medical image segmentation. *MIA* (2023).

[9] Jintai Chen, Xiangshang Zheng, Hongyun Yu, Danny Z Chen, and Jian Wu. 2021. Electrocardio panorama: Synthesizing new ECG views with self-supervision. In *IJCAI*.

[10] Ting Chen, Simon Kornblith, Mohammad Norouzi, et al. 2020. A simple framework for contrastive learning of visual representations. In *ICML*.

[11] Xinlei Chen, Haoqi Fan, Ross Girshick, et al. 2020. Improved baselines with momentum contrastive learning. *ArXiv Preprint ArXiv:2003.04297* (2020).

[12] Zhao Chen, Vijay Badrinarayanan, Chen-Yu Lee, et al. 2018. GradNorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *ICML*.

[13] Bowen Cheng, Alex Schwing, and Alexander Kirillov. 2021. Per-pixel classification is not all you need for semantic segmentation. In *NeurIPS*.

[14] Ching-Yao Chuang, Joshua Robinson, Yen-Chen Lin, et al. 2020. Debiased contrastive learning. *Advances in neural information processing systems* 33 (2020), 8765–8775.

[15] Jean-Antoine Désidéri. 2012. Multiple-gradient descent algorithm (MGDA) for multiobjective optimization. *Comptes Rendus Mathematique* (2012).

[16] Jean-Bastien Grill, Florian Strub, Florent Altché, et al. 2020. Bootstrap your own latent-a new approach to self-supervised learning. *NeurIPS* (2020).

[17] Yan Han, Chongyan Chen, Ahmed Tewfik, et al. 2022. Knowledge-augmented contrastive learning for abnormality classification and localization in chest X-rays with radiomics using a feedback loop. In *WACV*.

[18] Kaiming He, Haoqi Fan, Yuxin Wu, et al. 2020. Momentum contrast for unsupervised visual representation learning. In *CVPR*.

[19] Hanzhe Hu, Jinshi Cui, and Liwei Wang. 2021. Region-aware contrastive learning for semantic segmentation. In *ICCV*.

[20] Xinrong Hu, Dewen Zeng, Xiaowei Xu, et al. 2021. Semi-supervised contrastive learning for label-efficient medical image segmentation. In *MICCAI*.

[21] Yannis Kalantidis, Mert Bulent Sariyildiz, Noe Pion, et al. 2020. Hard negative mixing for contrastive learning. In *NeurIPS*.

[22] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *ArXiv Preprint ArXiv:1412.6980* (2014).

[23] Nikos Komodakis and Spyros Gidaris. 2018. Unsupervised representation learning by predicting image rotations. In *ICLR*.

[24] M Kumar, Benjamin Packer, and Daphne Koller. 2010. Self-paced learning for latent variable models. In *NeurIPS*.

[25] Shikun Liu, Shuaifeng Zhi, Edward Johns, et al. 2021. Bootstrapping semantic segmentation with regional contrast. *ArXiv Preprint ArXiv:2104.04465* (2021).

[26] Ishan Misra and Laurens van der Maaten. 2020. Self-supervised learning of pretext-invariant representations. In *CVPR*.

[27] Danielle F Pace, Adrian V Dalca, Tal Geva, et al. 2015. Interactive whole-heart segmentation in congenital heart disease. In *MICCAI*.

[28] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, et al. 2016. Context encoders: Feature learning by inpainting. In *CVPR*.

[29] Jizong Peng, Ping Wang, Christian Desrosiers, et al. 2021. Self-paced contrastive learning for semi-supervised medical image segmentation with meta-labels. In *NeurIPS*.

[30] Christian S Perone and Julien Cohen-Adad. 2018. Deep semi-supervised segmentation with weight-averaged consistency targets. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*.

[31] Quan Quan, Qingsong Yao, Jun Li, et al. 2022. Information-guided pixel augmentation for pixel-wise contrastive learning. *ArXiv Preprint ArXiv:2211.07118* (2022).

[32] Joshua Robinson, Ching-Yao Chuang, Suvrit Sra, et al. 2020. Contrastive learning with hard negative samples. *ArXiv Preprint ArXiv:2010.04592* (2020).

[33] Olaf Ronneberger, Philipp Fischer, Thomas Brox, et al. 2015. U-Net: Convolutional networks for biomedical image segmentation. In *MICCAI*.

[34] Karthik Abinav Sankararaman, Soham De, Zheng Xu, et al. 2020. The impact of neural network overparameterization on gradient confusion and stochastic gradient descent. In *ICML*.

[35] Alex Tamkin, Mike Wu, and Noah Goodman. 2020. Viewmaker networks: Learning views for unsupervised representation learning. *ArXiv Preprint ArXiv:2010.07432* (2020).

[36] Yonglong Tian, Chen Sun, Ben Poole, et al. 2020. What makes for good views for contrastive learning?. In *NeurIPS*.

[37] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *JMLR* (2008).

[38] Yen Nhi Truong Vu, Richard Wang, Niranjan Balachandar, et al. 2021. MedAug: Contrastive learning leveraging patient metadata improves representations for chest X-ray interpretation. In *MLHC*.

[39] Tongzhou Wang and Phillip Isola. 2020. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *ICML*.

[40] Tao Wang, Jianglin Lu, Zhihui Lai, et al. 2022. Uncertainty-Guided Pixel Contrastive Learning for Semi-Supervised Medical Image Segmentation. In *IJCAI*.

[41] Wenguan Wang, Tianfei Zhou, Fisher Yu, et al. 2021. Exploring cross-image pixel contrast for semantic segmentation. In *CVPR*.

[42] Xiao Wang and Guo-Jun Qi. 2022. Contrastive learning with stronger augmentations. *TPAMI* (2022).

[43] Xinlong Wang, Rufeng Zhang, Chunhua Shen, et al. 2021. Dense contrastive learning for self-supervised visual pre-training. In *CVPR*.

[44] Zirui Wang, Yulia Tsvetkov, Orhan Firat, et al. 2020. Gradient vaccine: Investigating and improving multi-task optimization in massively multilingual models. *ArXiv Preprint ArXiv:2010.05874* (2020).

[45] Zhirong Wu, Yuanjun Xiong, Stella X Yu, et al. 2018. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*.

[46] Yutong Xie, Jianpeng Zhang, Zehui Liao, et al. 2020. PGL: Prior-guided local self-supervised learning for 3D medical image segmentation. *ArXiv Preprint ArXiv:2011.12640* (2020).

[47] Zhenda Xie, Yutong Lin, Zheng Zhang, et al. 2021. Propagate yourself: Exploring pixel-level consistency for unsupervised visual representation learning. In *CVPR*.

[48] Ke Yan, Jinzheng Cai, Dakai Jin, et al. 2022. SAM: Self-supervised learning of pixel-wise anatomical embeddings in radiological images. *TMI* (2022).

[49] Dong Yin, Ashwin Pananjady, Max Lam, et al. 2018. Gradient diversity: A key ingredient for scalable distributed learning. In *AISTATS*.

[50] Chenyu You, Weicheng Dai, Fenglin Liu, et al. 2022. Mine your own anatomy: Revisiting medical image segmentation with extremely limited labels. *ArXiv Preprint ArXiv:2209.13476* (2022).

[51] Chenyu You, Weicheng Dai, Yifei Min, et al. 2023. Rethinking semi-supervised medical image segmentation: A variance-reduction perspective. *ArXiv Preprint ArXiv:2302.01735* (2023).

[52] Chenyu You, Weicheng Dai, Lawrence Staib, et al. 2022. Bootstrapping semi-supervised medical image segmentation with anatomical-aware contrastive distillation. *ArXiv Preprint ArXiv:2206.02307* (2022).

[53] Chenyu You, Ruihan Zhao, Lawrence H Staib, et al. 2022. Momentum contrastive voxel-wise representation learning for semi-supervised volumetric medical image segmentation. In *MICCAI*.

[54] Chenyu You, Yuan Zhou, Ruihan Zhao, et al. 2022. SimCVD: Simple contrastive voxel-wise representation distillation for semi-supervised medical image segmentation. *TMI* (2022).

[55] Tianhe Yu, Saurabh Kumar, Abhishek Gupta, et al. 2020. Gradient surgery for multi-task learning. In *NeurIPS*.

[56] Jure Zbontar, Li Jing, Ishan Misra, et al. 2021. Barlow twins: Self-supervised learning via redundancy reduction. In *ICML*.

[57] Dewen Zeng, Yawen Wu, Xinrong Hu, et al. 2021. Positional contrastive learning for volumetric medical image segmentation. In *MICCAI*.

[58] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2017. Mixup: Beyond empirical risk minimization. *ArXiv Preprint ArXiv:1710.09412* (2017).

[59] Yizhe Zhang, Lin Yang, Jianxu Chen, et al. 2017. Deep adversarial networks for biomedical image segmentation utilizing unannotated images. In *MICCAI*.

[60] Xiangyun Zhao, Raviteja Vemulapalli, Philip Andrew Mansfield, et al. 2021. Contrastive learning for label efficient semantic segmentation. In *ICCV*.

[61] Xiahai Zhuang. 2013. Challenges and methodologies of fully automatic whole heart segmentation: A review. *Journal of Healthcare Engineering* (2013).

[62] Xiahai Zhuang and Juan Shen. 2016. Multi-scale patch and multi-modality atlases for whole heart segmentation of MRI. *MIA* (2016).