

---

# BYOLMed3D: SELF-SUPERVISED REPRESENTATION LEARNING OF MEDICAL VIDEOS USING GRADIENT ACCUMULATION ASSISTED 3D BYOL FRAMEWORK

---

**Siladittya Manna**

Computer Vision and Pattern Recognition Unit,  
Indian Statistical Institute  
Kolkata, India  
siladittya\_r@isical.ac.in

**Souvik Chakraborty**

Computer Vision Researcher  
souvik.chakraborty207@gmail.com

## ABSTRACT

Applications on Medical Image Analysis suffer from acute shortage of large volume of data properly annotated by medical experts. Supervised Learning algorithms require a large volumes of balanced data to learn robust representations. Often supervised learning algorithms require various techniques to deal with imbalanced data. Self-supervised learning algorithms on the other hand are robust to imbalance in the data and are capable of learning robust representations. In this work, we train a 3D BYOL self-supervised model using gradient accumulation technique to deal with the large number of samples in a batch generally required in a self-supervised algorithm. To the best of our knowledge, this work is one of the first of its kind in this domain. We compare the results obtained through our experiments in the downstream task of ACL Tear Injury detection with the contemporary self-supervised pre-training methods and also with ResNet3D-18 initialized with the Kinetics-400 pre-trained weights. From the downstream task experiments, it is evident that the proposed framework outperforms the existing baselines.

**Keywords** Self-Supervised · Medical · Gradient Accumulation · BYOL.

## 1 Introduction

Self supervised learning is a technique of building generalist models that can learn from the data itself by recognizing underlying patterns in the data without any supervision or ground truth labels. Many real world applications suffer from a dearth of labelled data which acts as an impediment to training models sufficiently. Annotating huge volumes of unlabelled data not only requires ample time but also necessitates skilled domain knowledge. In this regard, self-supervised learning evolved as a machine learning method that learns from the data without explicit requirement of labels and is widely considered as an intermediate form between supervised and unsupervised learning. It relies on the correlation in the data, metadata embedded in the data in the stage termed as the pretext task. It is analogous to learning an approximate form of sense in humans which enables them to approach various problems and tasks based on previously acquired background of the scheme of things.

The pretext task often has the objective of predicting either some hidden data or some feature from the available data through finding patterns in the data. The downstream task follows the pretext task where the network trained in the pretext task helps in the actual classification or detection task and the limitation of having less annotated examples is henceforth overcome.

MoCo [1] is a self-supervised learning (SSL) algorithm which involves building huge dynamic dictionaries whose keys are sampled from data and represented by the encoder network. For every query, a corresponding matching key is looked up. There is an existent upgrade on this work [2] which uses a non-linear projector, larger batch size and color augmentations inspired by SimCLR [3], which is also a SSL framework. SimCLR uses large batch sizes instead of a memory bank like MoCo. The major drawback of both these works is the need to maintain a huge memory bank or large batch size.

Another novel approach for the pretext task is the Barlow Twins redundancy reduction method [4]. The method aims to compute loss from two standpoints - to make the network robust to noise and at the same time, minimize the redundancy between embeddings of distorted versions of the same image, which is achieved by trying to make the cross-correlation matrix between the outputs of two identical networks as close to the identity matrix as possible. However, the work has been found to be not immune to all kinds of data augmentations.

Self supervised learning methods that maximize similarity between augmented views of the same image achieved high performance using a Siamese network [?]. But the major bottleneck of such networks is the requirement of higher batch sizes and in turn, higher computational resources, failing which the accuracy drops significantly. Hence, there was a need to develop computationally efficient methods that did not rely so desperately on high batch sizes. SimCLR [5] is one such contemporary work that learns data representations through augmented views of the same example, employing a contrastive loss in the latent space. However, the method suffers slightly when using random crops as image augmentations and to mitigate this drawback, BYOL [6] offers more robustness by not relying on negative pairs and directly bootstrapping the learnt representations. Triplet BYOL [7] improves accuracy further by compensating for much lower batch sizes through the usage of a triple-view loss.

In all the aforementioned works, [1, 2, 4, ?, 5, 6, 7], the main area of application is natural images. Medical images, however, need to be treated differently as they are grayscale and there is less spatio-temporal variation in medical data. They are a class of images that unfortunately suffer from the common problem, the lack of sufficiently large, well-annotated datasets. Fully supervised deep neural networks are data hungry and so the usage of the former on medical images becomes restricted. Hence, various methods to learn representations from non-annotated medical data (pretext task) and utilize the same in the downstream task have emerged to fill the void.

In the domain of medical images, [8] uses a 3D self-supervised network based on contrastive SimCLR [3] and Monte Carlo dropout to enhance results in downstream task. In [9], temporal order prediction and geometric transformation prediction have been fused for ultrasound videos for representation learning. The learnt representations were employed for standard plane detection as the downstream task. However, there has not been much focus on learning features from a class-imbalanced multi-label dataset in these works. [10] uses momentum contrast method [1] against chest X-ray images and achieved better results than supervised learning algorithms. [11] proposes a CNN architecture for solving a jigsaw puzzle as the pretext task, to learn explainable visual representations. It also addresses memory constraint issues in the downstream task through a Divide-and-Teach approach. [12] proposes a semi-parallel architecture for the pretext task. The learnt spatial anatomical representations from the frames of magnetic resonance (MR) video clips are used for the diagnosis of knee medical conditions.

In this paper, we propose a self-supervised framework to learn features from knee MR videos. These features would aid in detection of knee ACL tear in the ultimate downstream task through learning accurate feature representations in the pretext task. Similar to BYOL, the pretext task relies on two neural networks - the online and target networks - that interact and learn mutually to output a representation that is pivotal in the downstream task, but the method makes considerable progress in extrapolating the original ResNet-50[13] encoder with a 3D version to leverage the usage of the same on medical videos and includes optimizations to account for improved performance despite low batch sizes.

The contributions of this paper are as follows:

- We propose a 3D variant of BYOL, a self-supervised learning algorithm for medical videos.
- We also apply gradient accumulation technique in self-supervised learning to accommodate an effectively large batch size, which is the first of its kind to the best of our knowledge.

The rest of the paper is organized as follows: Section 2 describes the issue of imbalanced data in modern machine learning or deep learning tasks. Section 3 consists of 3 subsections. In Section 3.1, we describe the pre-training framework proposed in this paper. In Sec. 3.2, we explain the gradient accumulation technique used in this work, and in Sec. 3.3 we describe the modified momentum update strategy. The details of the downstream task of detecting ACL tear injury is described in Sec. 3.4. In Section 4, we have described the dataset used in our work, the results obtained and comparative studies with the contemporary self-supervised algorithms and the supervised baseline models in Sec. 4.1, 4.3, 4.4, respectively. Finally, we conclude the work in Section 5.

## 2 The problem of data imbalance

A main deterrent in dealing with medical images is class imbalance. Most machine learning problems work best when the number of examples in each class are approximately comparable. Otherwise, there are high chances of most models to achieve good prediction results only with the majority class examples, whereas the minority class is incorrectly captured.

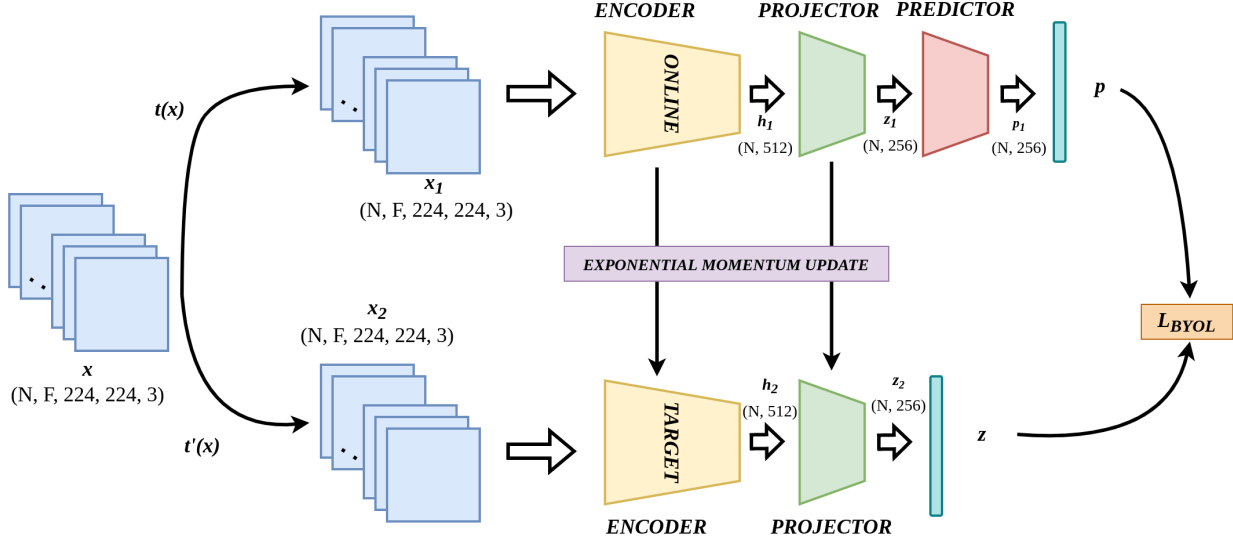


Figure 1: Pretraining Model Architecture : BYOLMed3D

Most frequently, images pertaining to a specific disease type are very less in number when compared to the size of the entire dataset. So most datasets available are highly skewed towards disease-free examples and hence, analysis of medical datasets with such class imbalance becomes a caveat. Any machine learning problem applied to the domain of medical images must attach consideration to overcome this limitation to ensure success.

In our work, we aim to use self-supervised learning techniques to overcome the first limitation - that of sufficient annotated data. Secondly, the algorithm does not rely on negative data pairs - a feature that, in fact, contributes to its increased robustness.

### 3 Methodology

#### 3.1 Pretraining

The basis of our framework in this paper is based on the BYOL [6], a non-contrastive self-supervised learning framework, which utilizes only positive pairs for representation learning. However, BYOL is used only for representation learning of images, whereas in this work, we use BYOL for representation of medical videos. To this end, we replace the ResNet-50 [13] encoder backbone with a 3D ResNet-18 [5] encoder backbone. The pretraining model architecture with the 3D Resnet-18 backbone is shown in Figure 1. The projector and predictor consists of a Multi Layered Perceptron with a single hidden layer with 4096 hidden units. The output from both projector and predictor consists of 256 dimensions.

The output from the online network is taken from the predictor, whereas the output from the target network is taken from the projector. For the rest of the paper, we will denote the online network parameters by  $\theta$  and the target network parameters using  $\xi$ . Furthermore, we will denote the input as  $x$ , the encoder as  $f_\theta/f_\xi$ , the output from the encoder  $f_\theta/f_\xi$  as  $h_\theta/h_\xi$ , the projector as  $g_\theta/g_\xi$  the output from the projector  $g_\theta/g_\xi$  as  $z_\theta/z_\xi$ , the predictor as  $q_\theta$  and the output from the predictor  $q_\theta$  as  $p_\theta$ . At the time of weight initialization, the online network is initialized and then the weights of the encoder and the projector are copied to the target network. The target network is updated only using an exponential momentum update routine. Thus at any iteration, the parameters of the target network is updated following the following equation,  $\xi = \alpha\xi + (1 - \alpha)\theta$ , where  $\alpha$  is the momentum update hyper-parameter representing the degree of weighting.

The input consists of 16 frames chosen randomly from a pool consisting of all the frames of a particular video. In this paper, we have used Knee MRI videos from the MRNet [14] dataset, but our framework will work for any medical video dataset, be it Magnetic Resonance Imaging, Computerized Tomography or Ultrasound. The 16 frames consist of a clip and  $N$  number of clips constitute a single batch. However, due to memory constraint on a single 8GB NVIDIA 2080Ti GPU, allocating a large number of clips with each frame having dimensions of  $224 \times 224$ , we have resorted to gradient accumulation technique for increasing the effective batch size to 256. Once the output from the predictor of the online network ( $p_\theta$ ) and the projector of the target network ( $z_\xi$ ) is obtained, we pass both of them to the loss function, which is same as used in [6], as shown in Equation 1.

$$\mathcal{L}_{BYOL} = 2 - 2 \langle \frac{p}{\|p\|_2}, \frac{z}{\|z\|_2} \rangle \quad (1)$$

,where  $p$  represents the output from the predictor of the online network and  $z$  represents the output from the projector of the target network, respectively.  $\langle p, z \rangle$  denotes the cosine similarity value calculated between the feature vectors  $p$  and  $z$ .

However, for feeding the input to online and target network, we take two different augmented versions of the input  $x$ . The augmented version  $t(x)$  is fed to the online network and the augmented version  $t'(x)$  is fed to the target network. Let, the loss obtained be  $\mathcal{L}_{BYOL}(p_\theta, z_\xi)$ . To maintain symmetry, the augmented versions  $t(x)$  and  $t'(x)$  are fed to the target and online networks, respectively as well. Let, the loss obtained in this case be  $\mathcal{L}_{BYOL}(p'_\theta, z'_\xi)$ . The final loss is the sum of the two losses, i.e.,  $\mathcal{L} = \mathcal{L}_{BYOL}(p_\theta, z_\xi) + \mathcal{L}_{BYOL}(p'_\theta, z'_\xi)$ .

The optimal representations are learnt by optimizing this loss  $\mathcal{L}_{BYOL}$  using LARS [15] optimizer with a learning rate of 0.2 for 200 epochs decayed following a cosine annealing schedule. In each iteration of any epoch, a batch size of 4 is taken. But the gradient accumulation is done only when the total number of samples for which the gradients are calculated, reaches a value of 256. This work is one of the first to apply gradient accumulation to lessen the ill-effects of low batch size in self-supervised representation learning scenario. The rest of the pretraining methodology is similar to BYOL [6].

### 3.2 Gradient Accumulation

Deep Learning applications often require large batch size for better performance. In a general gradient descent scenario, batch size is a hyper-parameter which helps in stabilizing the optimization process. In the neighbourhood of a minima, which we can consider to be a convex hyperplane, lower batch sizes are often associated with large fluctuations on the loss surface. Whereas under the same conditions, when using large batch size, the fluctuations are smoothed out the optimization process converges faster. However, when the optimization process is stuck in a local minima, it is preferable to use a large learning rate to get enough fluctuations to walk out of the local minima. A similar phenomenon can be replicated by using a low batch size.

However, since self-supervised algorithms like BYOL [6, 3, 1, 2], generally require large batch size, we resort to using gradient accumulation under the constraint of limited GPU memory. Gradient accumulation is generally used to increase the effective batch size to a large number even though only a small number of samples can be used as input in a single batch. The gradients for each parameter are added up for each batch for a fixed number of steps such that the total batch size for all such steps equals the required batch size. The effective gradient is then back propagated. This technique helps us to deal with large batch size requirements.

In this work, the self-supervised algorithm we are using is BYOL [6]. BYOL uses only positive pairs for learning representations. Hence, it is possible to use gradient accumulation with BYOL as no loss term for each such positive pair depends on the other positive pairs or samples. Hence, it is possible to divide a large batch of 256 into 64 steps each with a batch size of 4. After 64 steps, the summed up gradients are back propagated. The gradient accumulation strategy is described using pseudocode in [algorithm 1](#).

### 3.3 Modified Momentum Update Strategy

As mentioned in the previous section, the gradient updates are implemented after 64 iterations, to obtain an effective batch size of 256 using gradient accumulation technique. To maintain the same parameter for those 64 iterations for the target encoder, the momentum update of the target encoder is also implemented at an interval of 64 iterations as well. This ensures that the parameters of the target encoder remains same for all the 256 samples in the effective batch. Once, the gradient update is done after 64 iterations and the parameters of the online network are updated the momentum update of the target network is also implemented and the parameters of the target network are updated. The momentum update strategy is described using pseudocode in [algorithm 1](#).

### 3.4 Downstream

The downstream task involves predicting an appropriate medical diagnosis from the medical videos. When using the MRNet [14] dataset, the objective is to predict ACL Tear injury in the Knee MR videos. Thus, the downstream task in this case is a binary classification problem. However, due to the imbalanced nature of the dataset, we resort to random oversampling to deal with the imbalance.

**Algorithm 1:** Gradient Accumulation and Momentum Update Strategy

---

```

1  $T_B \leftarrow 256 \leftarrow$  Total Batch Size for which gradient accumulation is done;
2  $n_B \leftarrow 4 \leftarrow$  Actual Batch Size in each iteration;
3  $aiter \leftarrow \frac{T_B}{n_B} \leftarrow$  Number of iterations after which gradient will be applied;
4  $bnum \leftarrow$  Batch Number;
5  $numepoch \leftarrow$  Training Epoch Number;
6  $compute\_grad \leftarrow$  Computes Gradients by Back Propagation;
7 while  $numepoch \neq 200$  do
8    $bnum = 0$ ;
9   while DataLoader not empty do
10     $p_\theta = q_\theta(g_\theta(f_\theta(t(x))))$ 
11     $z_\xi = g_\xi(f_\xi(t'(x)))$ 
12     $p'_\theta = q_\theta(g_\theta(f_\theta(t'(x))))$ 
13     $z'_\xi = g_\xi(f_\xi(t(x)))$ 
14     $\mathcal{L} = \mathcal{L}_{BYOL}(p_\theta, z_\xi) + \mathcal{L}_{BYOL}(p'_\theta, z'_\xi)$ 
15     $\delta_G = \delta_G + compute\_grad(\mathcal{L})$ 
16    if  $bnum \% aiter == 0$  then
17       $\theta = \theta - \eta \delta_G$ ;
18       $\xi = \alpha \xi + (1 - \alpha) \theta$ ;
19     $bnum = bnum + 1$ ;
20   $numepoch = numepoch + 1$ ;

```

---

In this task, we discard the predictor and projector from the online network. We only use the encoder backbone 3D ResNet18 [5] of the online network as the feature extractor. Similar to the pretraining phase, we randomly select 16 frames from a MR video and pass it through the 3D ResNet18 encoder to obtain feature vectors. These feature vectors are then used to train a linear classifier for the final predictions. We use Adam optimizer for fine-tuning the 3D ResNet18 encoder. We fine tune the encoder and the classifier until the loss flattens and use the model with the lowest validation accuracy for inference on the test set.

## 4 Experiments and Results

In this section, we present the results obtained with the proposed framework on MRNet dataset. We also present comparative studies with state-of-the-art self-supervised learning algorithms, as well as, supervised baseline methods.

### 4.1 Dataset

In our experiments, we use the MRNet [14] dataset as our reference dataset. The MRNet dataset contains 1370 knee MR video clips in total. Out of 1370 clips, 1130 MR video clips are included in the training set and 120 MR video clips are considered as the tuning or validation set. The rest 120 are used for external validation. Out of the 1,130 training examples, only 208 videos contain ACL tear. It is evident that the dataset we are using for this work is highly imbalanced, with an imbalance ratio greater than 4. This gives us an opportunity to explore the effects of self-supervised pretraining representation learning techniques on imbalanced datasets.

### 4.2 Experimental Details

In this subsection we present the configuration of the model used in our experiments. The model was pre-trained with an encoder backbone of 3D ResNet18 and trained with LARS [15] optimizer. The number of samples in a single batch was 4 and the effective accumulated batch size was set to 256. So, in the pre-training phase, the gradient update was done after 64 iterations, propagating the gradient obtained from 256 samples in one update. The pre-training model was trained for 200 epochs with LARS optimizer with a learning rate of 0.2. The learning rate was decayed following a cosine annealing schedule.

For the downstream task, we used the weights learnt in the pre-training phase as the initializing weights. We randomly selected 16 frames from each MR video and used them as input. The downstream model was fine-tuned by optimizing

the binary cross-entropy loss using Adam optimizer. As the dataset is highly imbalanced, we oversample the dataset using random oversampling strategy to deal with the imbalance.

### 4.3 Results

In this subsection, we present the results obtained in the downstream experiments with the pre-trained model for detecting ACL Tear injury from MR data. In Table 1, we present the accuracy, AUC, Sensitivity and specificity for different configuration of the downstream experiments. For the experiments using Multistep decay, the learning rate was decayed by 0.1 times at the  $0.6 \times T_E$  and  $0.8 \times T_E$ -th epochs, where  $T_E$  is the total number of epochs of downstream fine-tuning. For the experiments using step decay, the learning rate was decayed by 0.98 times at every epoch.

Table 1: Experimental Results of the proposed framework on different sets of hyperparameters

Model	Epochs	Optim.	LR	Decay	Accuracy(%)	AUC	Sensitivity	Specificity
1	30	Adam	1e-4	Multistep	70.8	0.81	0.463	0.909
2	50	Adam	1e-4	Multistep	<b>75.8</b>	0.824	<b>0.574</b>	0.909
3	100	Adam	1e-4	Multistep	74.2	<b>0.825</b>	0.5	<b>0.939</b>

It is to be noted, that the accuracy scores were obtained on the validation set of MRNet [14] as the test set is hidden. For conducting fair experiment, the validation set was treated as the test set in the experiments. A separate validation set was constructed by using a 90 : 10 split from the training set consisting of 1130 samples.

### 4.4 Comparison with Existing Self-Supervised Methods

In this subsection, we compare our work with [11] and [12], which are two self-supervised learning algorithms on MRNet [14] dataset for detection of ACL tear injury from MR data. The comparative results are mentioned in Table 2.

Table 2: Comparison of the proposed method with contemporary self-supervised learning algorithms

Method	No. of Parameters	Accuracy(%)
S. Manna et al.[11]	77 Million	76.72
SSLM [12]	103 Million	74.0
BYOLMed3D (Best)	33.4 Million	<b>75.8</b>

The models used in our experiments ResNet3D-18 contains 33.4M parameters, whereas, the best performing model used in [11] has 77M parameters. The model used in SSLM [12] has 103M parameters in the downstream fine-tuning. Thus, we can see that the proposed framework performs at par or better than the existing state-of-the-art methods on ACL Tear detection on MRNet dataset with less than half number of parameters used in the models in the other methods.

It is to be noted that, to the best of our knowledge only S. Manna et al. [11] and SSLM [12] has used MRNet for their experiments. Thus, it has not been possible to compare with other self-supervised methods. Although the work [8] also presents a 3D adapted version of the self-supervised algorithm SimCLR [3], it is not possible to implement gradient accumulation on that algorithm as it requires negative pairs to be accommodated in a batch along with a positive pair. Furthermore, SimCLR requires a large batch size to give satisfactory performance. Hence, it will not be possible to compare with the work [8] here under identical circumstances.

### 4.5 Comparison with Kinetics-400 pretrained and randomly initialized weights

In this subsection, we present comparison results of the downstream task fine-tuned with pre-trained weights obtained by training on Kinetics-400 [16] dataset. We also trained the downstream model from scratch, that is with randomly initialized weights. The comparative results are presented in Table 3. All the models were fine-tuned by optimizing binary cross-entropy loss with Adam optimizer with an initial learning rate of  $10^{-4}$  and multistep decay by 0.1 times at the  $0.6 \times T_E$  and  $0.8 \times T_E$ -th epochs, where  $T_E$  is the total number of epochs of downstream fine-tuning. As the number of classes in the Kinetics-400 dataset is 400, the final fully-connected layer in the model initialized with Kinetics-400 pre-trained weights, was changed to a linear layer with 1 output and randomly initialized prior to fine-tuning.

We can observe that the self-supervised pre-trained weights outperforms the model initialized with weights obtained by training on Kinetics-400 dataset for the task of Action Recognition. The weights learnt from the Kinetics-400 dataset were tailored to the dataset and as the Kinetics-400 dataset is quite clearly different from MR data, the weights do not



provide a good initialization point in the loss landscape for the downstream task. This is primarily because the features learnt by training the ResNet3D-18 model on Kinetics-400 dataset are specific to the objects present in the dataset sample.

We can also observe that the randomly initialized weight provides a better starting point in the loss landscape and consequently a better end result in our experiments. This phenomenon may be because the optimizing the weights from a random start point is easier than to change the weights which are already specific to some dataset.

Table 3: Comparison of the Self-supervised pre-trained model with Kinetics 400 pre-trained and randomly initialized model

Downstream initialization	Accuracy(%)	AUC	Sensitivity	Specificity
Kinetics 400 pre-trained	68.3	0.902	0.296	1.0
Randomly Initialized	74.2	0.804	0.537	0.909
BYOLMed3D pre-trained (Best)	<b>75.8</b>	0.825	0.574	0.909

## 5 Conclusion

The primary objective of this work is to expand the horizon of self-supervised learning in the domain of medical imaging. The pretext task has been shown to extract structural features which are aiding the downstream task. Gradient accumulation helps the task to use small batch sizes but save the gradients and update the weights of the network after intervals of batches, thereby increasing the effective batch size. It serves the same purpose as having a mini-batch with a higher number of images. BYOL learns its representation by predicting previous versions of its outputs, without using negative pairs. In class imbalanced datasets, the lack of dependence on negative pairs is an obvious advantage. This is the first work of its kind that incorporates learning of medical videos using the mentioned gradient accumulation technique and data adapted BYOL, the downstream task ultimately being a classification problem. The results encourage further research in this direction, especially for annotation-starved medical images, with the objective of learning more robust and reliable representational features in the pretext task.

## References

- [1] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick, “Momentum contrast for unsupervised visual representation learning,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pp. 9726–9735, IEEE, 2020. 1, 2, 4
- [2] X. Chen, H. Fan, R. B. Girshick, and K. He, “Improved baselines with momentum contrastive learning,” *CoRR*, vol. abs/2003.04297, 2020. 1, 2, 4
- [3] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, “A simple framework for contrastive learning of visual representations,” in *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, vol. 119 of *Proceedings of Machine Learning Research*, pp. 1597–1607, PMLR, 2020. 1, 2, 4, 6
- [4] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, “Barlow twins: Self-supervised learning via redundancy reduction,” in *ICML*, 2021. 2
- [5] K. Hara, H. Kataoka, and Y. Satoh, “Learning spatio-temporal features with 3d residual networks for action recognition,” *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, pp. 3154–3160, 2017. 2, 3, 5
- [6] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, B. Piot, k. kavukcuoglu, R. Munos, and M. Valko, “Bootstrap your own latent - a new approach to self-supervised learning,” in *Advances in Neural Information Processing Systems* (H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, eds.), vol. 33, pp. 21271–21284, Curran Associates, Inc., 2020. 2, 3, 4
- [7] G. Li, R. Togo, T. Ogawa, and M. Haseyama, “Tribyol: Triplet byol for self-supervised representation learning,” in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3458–3462, 2022. 2
- [8] Y. Ali, A. Taleb, M. M.-C. Hohne, and C. Lippert, “Self-supervised learning for 3d medical image analysis using 3d simclr and monte carlo dropout,” *ArXiv*, vol. abs/2109.14288, 2021. 2, 6
- [9] J. Jiao, R. Droste, L. Drukker, A. T. Papageorgiou, and J. A. Noble, “Self-supervised representation learning for ultrasound video,” in *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*, pp. 1847–1850, 2020. 2
- [10] H. Sowrirajan, J. Yang, A. Ng, and P. Rajpurkar, “Moco pretraining improves representation and transferability of chest x-ray models,” *ArXiv*, vol. abs/2010.05352, 2020. 2
- [11] S. Manna, S. Bhattacharya, and U. Pal, “Self-supervised representation learning for detection of acl tear injury in knee mr videos,” *Pattern Recogn. Lett.*, vol. 154, p. 37–43, feb 2022. 2, 6

- [12] S. Manna, S. Bhattacharya, and U. Pal, “Sslm: Self-supervised learning for medical diagnosis from mr video,” *ArXiv*, vol. abs/2104.10481, 2021. 2, 6
- [13] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016. 2, 3
- [14] N. Bien, P. Rajpurkar, R. Ball, J. Irvin, A. Park, E. Jones, M. Bereket, B. Patel, K. Yeom, K. Shpanskaya, S. Halabi, E. Zucker, G. Fanton, D. Amanatullah, C. Beaulieu, G. M. Riley, R. Stewart, F. Blankenberg, D. Larson, R. Jones, C. Langlotz, A. Ng, and M. Lungren, “Deep-learning-assisted diagnosis for knee magnetic resonance imaging: Development and retrospective validation of mrnet,” *PLoS Medicine*, vol. 15, 2018. 3, 4, 5, 6
- [15] Y. You, I. Gitman, and B. Ginsburg, “Large batch training of convolutional networks,” 2017. 4, 5
- [16] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman, “The kinetics human action video dataset,” *CoRR*, vol. abs/1705.06950, 2017. 6