

Supplementary material

1 Implementation details

Self-supervised pretrain In self-supervised pretraining stage, we resize each frame to 128×171 and randomly crop images to size of 112×112 in a temporal consistent way. Each input clip has 16 frames with temporal interval of 4. Data augmentations include random color jittering, horizontal flipping and gaussian blurring. We pretrain the model for 200 epochs with an SGD optimizer with an initial learning rate of 0.003, weight decay of $1e-4$ and momentum of 0.9. Learning rate is divided by 10 at epoch 120 and 160. Batch size is set to 64. We pretrain the model on 8 Tesla V100 GPUs. Pretrained weights at epoch 190 is selected for evaluation. Size of the MoCo negative queue is set to 16,384. We pretrain on UCF101 training split in all ablation studies and provide pretraining result on large-scale Kinetics400 for performance comparison with previous works. We set τ , τ_{tc} , θ , λ_1 , λ_2 to 0.07, 0.5, 0.05, 1.0 and 1.0, respectively.

Supervised finetuning During supervised finetuning, we replace the nonlinear projection head in pretraining stage with a one-layer classification layer. The backbone is initialized with the pretrained weights while the classification head is randomly initialized. Similar as pretraining stage, each clip contains 16 frames sampled at a pace equals to 2. Each frame is first resized to 128×171 and then randomly cropped to size of 112×112 in a temporal consistent way. Data augmentations include color jittering and horizontal flipping. We finetune all layers for 150 epochs using a SGD optimizer with momentum of 0.9. Learning rate is set to 0.05 and divided by 10 at epoch 50 and 100. Batchsize is set to 64.

Finetuning testing After finetuning, we test classification accuracy on the test splits. 10 16-frame clips are temporarily and uniformly sampled in each video with a pace equals to 2. No data augmentation is applied. Each frame is first resized to 128×171 and then center cropped to size of 112×112 . Classification probabilities are averaged among the 10 clips for each video.

Video retrieval We compute average features of 10 clips uniformly sampled from each video using the pretrained backbone. Setting is the same as the *Finetuning testing* stage except that we use a pace equals to 4 to be consistent with the pretraining stage. Feature vectors are computed through spatial-temporal average pooling the last layer output of the backbone and then normalized. Cosine similarity is adopted. We conduct video retrieval on UCF101 to retrieve test videos using training videos and calculate the top- k accuracy ($k = 1, 5, 10, 20, 50$).

2 Comparison with shuffling order prediction

We compare our method with shuffling order prediction baseline in Table 1. Improvement on HMDB51 dataset is much larger than that on UCF101, showing our model has a much better transferring ability across different dataset since model is pretrained on UCF101 training split.

Table 1: Comparison between order prediction (OP) and our proposed shuffle-rank (SR). Models are pretrained on UCF101 training set.

	UCF101	HMDB51
SimCLR+OP	75.39	32.38
SimCLR+SR	76.82	40.05

3 More experiment results on inter-intra variance

Continuing our analysis in ablation study section, we provide more experiment results of inter-intra variance on HMDB51 datasets for reference. As can be seen in Table 2, the trend of inter/intra variance changes are consistent on both UCF101 and HMDB51.

Table 2: Comparison of inter-instance variance, intra-instance variance and instance discrimination factor. Results are multiplied by 100 for demonstration. Each cell’s increasing(Red) and decreasing(Green) times are compared to the cell above it.

	R3D			R(2+1)D			S3D-G		
	inter-v.	intra-v.	discrim.	inter-v.	intra-v.	discrim.	inter-v.	intra-v.	discrim.
UCF101									
SimCLR	28.8	0.8	34.3	28.2	0.8	33.6	73.1	3.0	24.7
+SR	22.6	11.3 \uparrow 13.6 \times	2.0 \downarrow 17.2 \times	23.0	11.2 \uparrow 13.4 \times	2.0 \downarrow 16.8 \times	51.7	11.7 \uparrow 4.0 \times	4.4 \downarrow 5.6 \times
+SR+TC	42.3	5.8 \downarrow 1.9 \times	7.3 \uparrow 3.7 \times	38.0	8.0 \downarrow 1.4 \times	4.8 \uparrow 2.4 \times	78.7	4.9 \downarrow 2.4 \times	16.2 \uparrow 3.7 \times
HMDB51									
SimCLR	26.6	0.3	85.8	25.8	0.3	83.3	71.1	1.2	58.3
+SR	33.7	5.8 \uparrow 18.7 \times	5.8 \downarrow 14.8 \times	32.2	5.1 \uparrow 16.4 \times	6.4 \downarrow 13.0 \times	61.8	4.8 \uparrow 3.9 \times	12.9 \downarrow 4.5 \times
+SR+TC	43.0	2.9 \downarrow 2.0 \times	14.8 \uparrow 2.6 \times	40.0	3.8 \downarrow 1.3 \times	10.6 \uparrow 1.7 \times	73.3	2.5 \downarrow 1.9 \times	29.0 \uparrow 2.2 \times

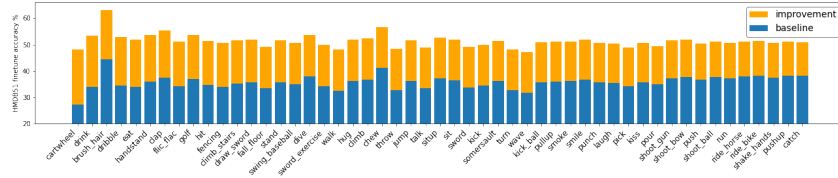


Figure 1: Per-class improvement on HMDB51 downstream supervised finetuning task with SimCLR baseline and R(2+1)D backbone. Model is unsupervisedly pretrained on UCF101 training split and θ is chosen as 0.01. The category labels along the horizontal coordinate follows decreasing order of absolute improvement magnitudes upon the baseline (height of orange bars).

39 4 Result of different number of segments

We ablated the number of segments in Table 3. Model is equipped with R(2+1)D backbone and pretrained on UCF101 training split. As input clip length is the same, increasing number of segments also increases the similarity between temporarily neighboring segments, which largely increases the difficulty for shuffle-rank task. We find increase segment number from 2 to 4 does not improve the performance much, but 8 segments make the loss not converging, which is caused by the difficulty of the task.

Number of Segment	UCF101	HMDB51
2	79.01	45.37
4	77.32	46.05
8	<i>null</i>	<i>null</i>

Table 3: Finetune task result on different number of segments with R(2+1)D backbone pretrained on UCF101 training split. *null* means loss can not converge in training time.

46 5 Per-class improvement result

Per-class Analysis We provide per-class accuracy improvement result of our method upon baseline SimCLR in Figure 1. Remind that in HMDB51 dataset, highly-improved categories such as cartwheel, drink, eat and handstand usually contain decomposable unrepertitive actions with large motion changes, while categories like running, riding horses, shake-hands consist of many repetitive actions thus are less improved by our method. All actions have been improved. These results supported our motivation that representing videos by sub-features are useful.

6 Backbone architecture

We provide architectures of R3D and R(2+1)D in Table 4. As the structure of S3D-G is too complicated to be put in a table, we ask the reader to refer to [3] for more details. We use the same implementation in our code.

Table 4: architecture of R3D and R(2+1)D we used in experiments.

layer name	output size	R3D	R(2+1)D
conv1	$L \times 56 \times 56$	$3 \times 3 \times 3, 64, \text{stride } 1 \times 2 \times 2$	
conv2_x	$L \times 56 \times 56$	$\begin{bmatrix} 3 \times 3 \times 3, 64 \\ 3 \times 3 \times 3, 64 \end{bmatrix} \times 1$	$\begin{bmatrix} 1 \times 3 \times 3, 64 \\ 3 \times 1 \times 1, 64 \end{bmatrix} \times 2$
conv3_x	$\frac{L}{2} \times 28 \times 28$	$\begin{bmatrix} 3 \times 3 \times 3, 128 \\ 3 \times 3 \times 3, 128 \end{bmatrix} \times 1$	$\begin{bmatrix} 1 \times 3 \times 3, 128 \\ 3 \times 1 \times 1, 128 \end{bmatrix} \times 2$
conv4_x	$\frac{L}{4} \times 14 \times 14$	$\begin{bmatrix} 3 \times 3 \times 3, 256 \\ 3 \times 3 \times 3, 256 \end{bmatrix} \times 1$	$\begin{bmatrix} 1 \times 3 \times 3, 256 \\ 3 \times 1 \times 1, 256 \end{bmatrix} \times 2$
conv5_x	$\frac{L}{8} \times 7 \times 7$	$\begin{bmatrix} 3 \times 3 \times 3, 512 \\ 3 \times 3 \times 3, 512 \end{bmatrix} \times 1$	$\begin{bmatrix} 1 \times 3 \times 3, 512 \\ 3 \times 1 \times 1, 512 \end{bmatrix} \times 2$
	$1 \times 1 \times 1$	spatial-temporal global average pooling	

7 Measure inter and intra variances

We provide our mathematical formulations of inter-variance and intra-variance for calculating inter and intra variances. Formally, let μ_k represent the mean embedded feature of the k -th video. We use S_k to represent the hierarchical relation between clips $\{c_i\}_{i=1}^M$ and videos $\{V_i\}_{i=1}^N$, i.e. $S_k = \{i | c_i \in V_k\}$. We define intra-variance as the average variances within clip features sampled from the same video, and inter-variance as the average pairwise distance between videos:

$$\mu_k = \frac{1}{S_k} \sum_{i \in S_k} z_i \quad (1)$$

$$\sigma_{intra} = \frac{1}{N} \sum_{k=1}^N \frac{1}{|S_k|} \sum_{i \in S_k} \|z_i - \mu_k\|_2^2 \quad (2)$$

$$\sigma_{inter} = \frac{1}{N(N-1)} \sum_{i=1}^{N-1} \sum_{j=i+1}^N \|\mu_i - \mu_j\|_2^2 \quad (3)$$

where z is the encoded feature vector as introduced in the paper. This formulation is adapted from Liu et al. [1].

8 Effect of contrastive loss on inter and intra variance

We provide theoretical analysis to our statement in preliminary section that contrastive loss suppresses the learning of σ_{inter} and encourages learning σ_{intra} . Our analysis follows discussion of Tongzhou and Phillip [2]. Following our definition in preliminary section, clip contrastive loss can be written in the form of $\mathcal{L}_c = -\frac{1}{M} \sum_{i=1}^M \log \frac{\exp(\frac{z_i \cdot z_{i^+}}{\tau})}{\sum_{k=1}^M \mathbf{1}_{[k \neq i]} \exp(\frac{z_i \cdot z_k}{\tau})}$, which can be transformed into:

$$\mathcal{L}_c = -\frac{1}{M} \sum_{i=1}^M \frac{z_i \cdot z_{i^+}}{\tau} + \frac{1}{M} \sum_{i=1}^M \left[\log \left(\exp\left(\frac{z_i \cdot z_{i^+}}{\tau}\right) + \sum_{k=1, k \notin \{i, i^+\}}^M \exp\left(\frac{z_i \cdot z_k}{\tau}\right) \right) \right] \quad (4)$$

where M is the number of sampled clips, z is normalized encoded clip feature vector and τ is the temperature parameter. As explained in [2], since $\sum_{k=1, k \notin \{i, i^+\}}^M \exp(z_i \cdot z_k / \tau)$ is always positive and bounded below, the loss favors smaller $-\frac{1}{M} \sum_{i=1}^M z_i \cdot z_{i^+} / \tau$, which is equivalent to decreasing σ_{intra} as the sample size M goes to infinity. On the other hand, when fixing σ_{intra} , minimizing the loss is equivalent to decreasing

$$\frac{1}{M} \sum_{i=1}^M \left[\log(t_i + \sum_{k=1, k \notin \{i, i^+\}}^M \exp(z_i \cdot z_k / \tau)) \right] \quad (5)$$

where $t_i = \exp(z_i \cdot z_{i^+} / \tau)$ is a random variable when σ_{intra} is fixed. Hence, equation 5 measures similarity between features from different instances. Minimizing equation 5 is thus equivalent to minimizing inter-variance σ_{inter} . Therefore, contrastive loss encourages decreasing intra-variance and increasing inter-variance.

9 Training algorithm

9.1 Training SimCLR

Algorithm 1: Training SimCLR

Input: batch size N , backbone f , clip projector f_c , dual projector f_r , videos V ;

while training not converge **do**

 sample minibatch $\{v_i\}_{i=1}^N$ from V ;

for all $i \in \{1, 2, \dots, N\}$ **do**

 sample and augment two clips c_i^1 and c_i^2 from v_i ;

$h_i^1 = f(c_i^1)$;

$h_i^2 = f(c_i^2)$;

 /* clip projection */

$z_{2i-1} = f_c(h_i^1)$;

$z_{2i} = f_c(h_i^2)$;

 /* dual projection */

$r_{2i-1} = f_r(h_i^1)$;

$r_{2i} = f_r(h_i^2)$;

 /* shuffling and dual projection */

 augment c_i^1 into s_i and randomly shuffle s_i into \hat{s}_i ;

$q_i = f_r(f(s_i))$;

$p_i = f_r(f(\hat{s}_i))$;

end

 calculate \mathcal{L}_c on $\{z\}_{i=1}^{2N}$;

 calculate \mathcal{L}_{tc} on $\{r\}_{i=1}^{2N}$;

 calculate $\mathcal{L}_{rank}^{unaug}$ on $\{r_{2i-1}, p_i\}_{i=1}^N$ and \mathcal{L}_{rank}^{aug} on $\{q_i, p_i\}_{i=1}^N$;

 update f , f_c and f_r using stochastic gradient descent;

end

9.2 Training MoCo

Unlike SimCLR, negative examples in MoCo come from the dictionary $\{k_i\}_{i=1}^m$ while positive example is the sampled augmented clip k^+ . m is set to 16384. So contrastive loss is denoted as :

$$\mathcal{L}_c = -\frac{1}{N} \sum_{i=1}^N \frac{\exp(z_i \cdot k_i^+ / \tau)}{\exp(z_i \cdot k_i^+ / \tau) + \sum_{j=1}^m \exp(z_i \cdot k_j / \tau)} \quad (6)$$

97 11 Attention visualization

98 To qualitatively measure the learned features, we draw heatmaps output from last layer of pretrained
 99 models. Specifically, we apply average pooling on last layer feature maps along both feature channel
 100 and temporal channel, i.e. from a size of (C, T, H, W) into (H, W) , to obtain a heatmap. Such
 101 heatmap is then added to each frame of the input clip for visualization. As shown in Figure 3, baseline
 102 SimCLR can be easily distracted by surrounding backgrounds, while our proposed method always
 103 focus on areas with motion semantics.

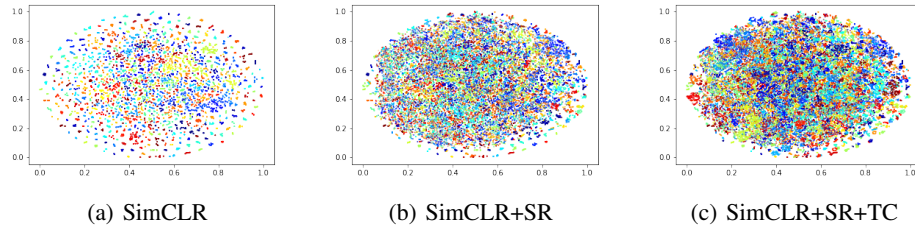


Figure 2: t-SNE visualization of features encoded by pretrained models. Models are pretrained on UCF101 train split and features are computed on UCF101 test split.

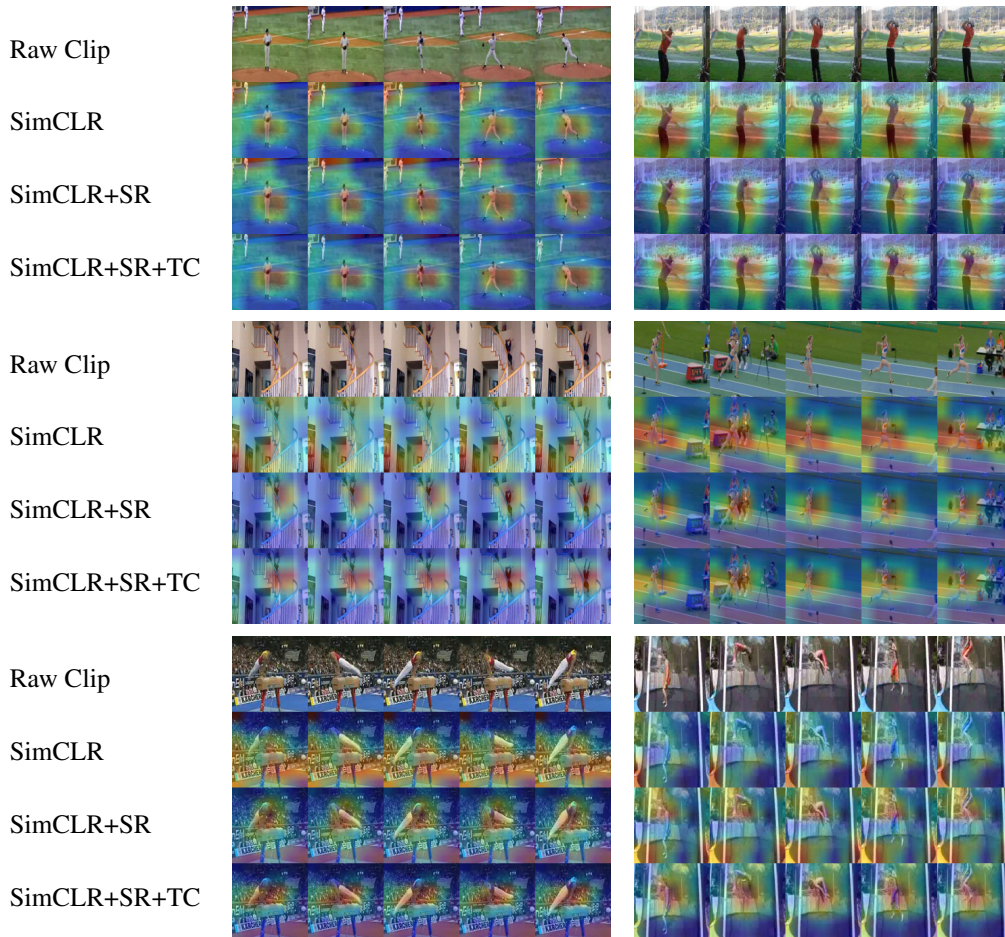


Figure 3: Heatmap visualization. The warmer the color is, the larger the response is. Zoom in for visualization.

12 Video retrieval visualization

In Figure 4, we visualize the video retrieval result of the integrated pretrained model (SimCLR+SR+TC). We pretrain the model on Kinetics-400 and conduct retrieval experiment on UCF101. It can be observed that our retrieval results are not perfect in that videos from different categories can be retrieved if having similar actions and contents, e.g. image 3 (ApplyLipstick) of ApplyEyeMakeup, image 3 (Shotput) of GolfSwing, image 3 (ParallelBars) of BalanceBeam and image 3 (BalanceBeam) of ParallelBars.

References

- [1] B. Liu, Y. Cao, Y. Lin, Q. Li, Z. Zhang, M. Long, and H. Hu. Negative margin matters: Understanding margin in few-shot classification. In *ECCV*, 2020.
- [2] W. Tongzhou and I. Phillip. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *ICML*, 2020.
- [3] S. Xie, C. Sun, J. Huang, Z. Tu, and K. Murphy. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *ECCV*, 2018.

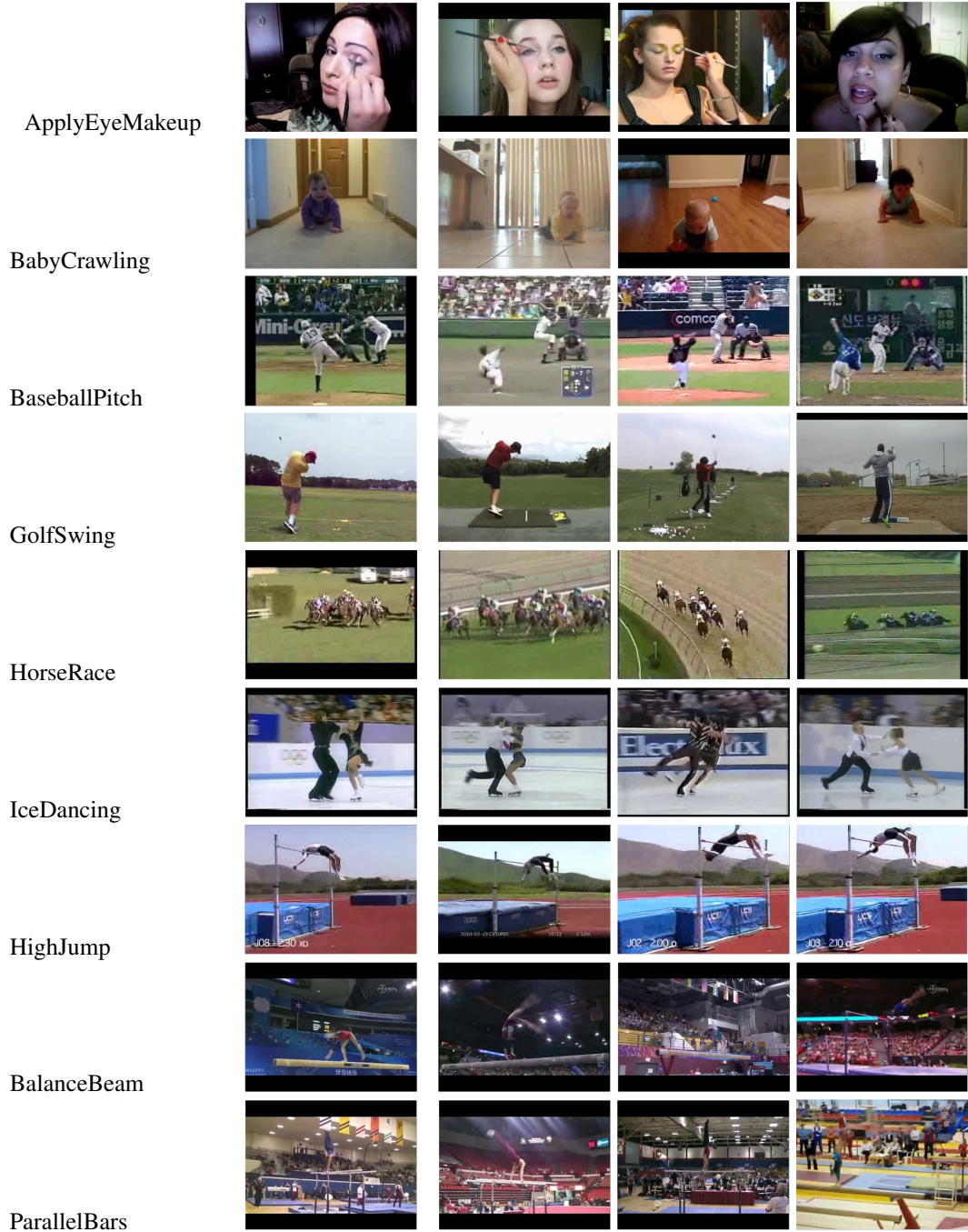


Figure 4: Video retrieval visualization. The first and second column are categories and images of test instance, respectively. The rightmost 3 columns are nearest retrieval results. We select an image from the video for demonstration.