

---

# Unfolding Videos Dynamics via Taylor Expansion

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Taking inspiration from physical motion, we present a new self-supervised dynamics  
2 learning strategy for videos: **Video Time-Differentiation for Instance**  
3 **Discrimination (ViDiDi)**. ViDiDi is a simple and data-efficient strategy, read-  
4 ily applicable to existing self-supervised video representation learning frameworks  
5 based on instance discrimination. At its core, ViDiDi observes different aspects  
6 of a video through various orders of temporal derivatives of its frame sequence.  
7 These derivatives, along with the original frames, support the Taylor series expan-  
8 sion of the underlying continuous dynamics at discrete times, where higher-order  
9 derivatives emphasize higher-order motion features. ViDiDi learns a single neural  
10 network that encodes a video and its temporal derivatives into consistent embed-  
11 dings following a balanced alternating learning algorithm. By learning consistent  
12 representations for original frames and derivatives, the encoder is steered to em-  
13 phasize motion features over static backgrounds and uncover the hidden dynamics  
14 in original frames. Hence, video representations are better separated by dynamic  
15 features. We integrate ViDiDi into existing instance discrimination frameworks  
16 (VICReg, BYOL, and SimCLR) for pretraining on UCF101 or Kinetics and test  
17 on standard benchmarks including video retrieval, action recognition, and action  
18 detection. The performances are enhanced by a significant margin without the need  
19 for large models or extensive datasets.

## 20 1 Introduction

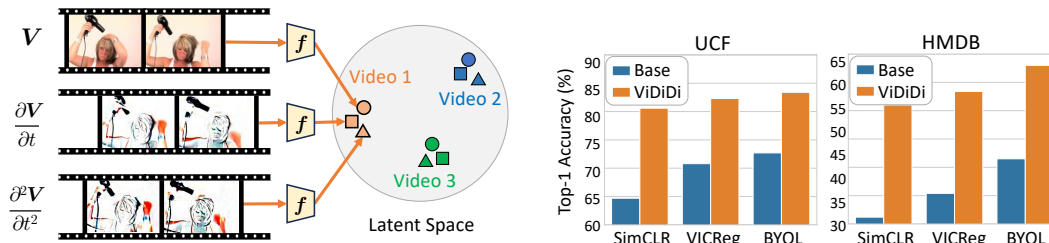


Figure 1: **Method Overview**. Left: Demonstration of ViDiDi. Right: ViDiDi enhances existing instance discrimination methods significantly on action recognition.

21 Learning video representations is central to various aspects of video understanding, such as action  
22 recognition [1, 2], video retrieval [3, 4], and action detection [5]. While supervised learning requires  
23 expensive video labeling [6], recent works highlight the strengths of self-supervised learning (SSL)  
24 from unlabeled videos [7, 8, 9] with a large number of training videos. One popular strategy for  
25 SSL on video representations uses instance discrimination objectives, such as SimCLR [10], initially  
26 demonstrated for images [11, 10, 12, 13, 14]. When adapting this approach to videos, previous  
27 methods often treat time directly as an additional spatial dimension. This may neglect the special

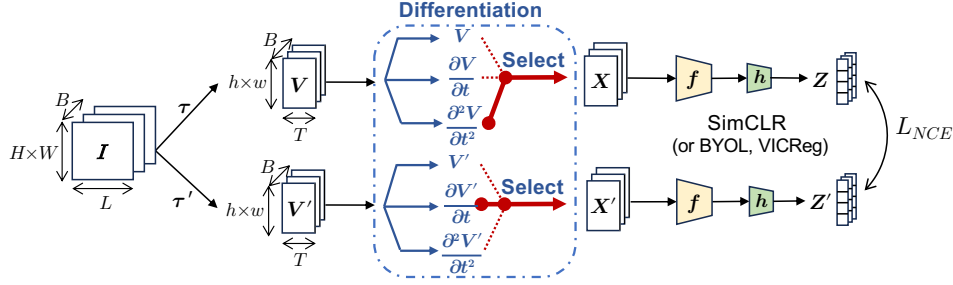


Figure 2: **Illustration of the ViDiDi framework.** For a batch of videos  $I$ , we do two spatio-temporal augmentations to obtain two batches of clips:  $V$  and  $V'$ .  $X$  and  $X'$  are the  $0^{th}$ ,  $1^{st}$ , or  $2^{nd}$  order derivatives of  $V$  and  $V'$ , decided via a balanced alternating learning strategy (alg. 1). They are the inputs to the video encoder for learning the video encoder under SimCLR, BYOL, or VICReg.

28 feature of the temporal dimension in carrying dynamic information, causing models to prioritize  
 29 static content (e.g., background scenes) over dynamic features (e.g., motion, action, and interaction),  
 30 which are often essential to video understanding [1, 2, 6, 15, 16, 17, 18, 19].

31 In contrast, we utilize the unique role of time in "unfolding" continuous real-world dynamics. We  
 32 introduce a generalizable and data-efficient method, applicable to self-supervised video representation  
 33 learning through instance discrimination including VICReg [13], BYOL [11], and SimCLR [10]. We  
 34 view a video as a continuous and dynamic process and use the Taylor series expansion to express this  
 35 continuous process as a weighted sum of its temporal derivatives at each frame. Based on this, we  
 36 evaluate temporal derivatives of videos as hidden views apart from the original frames. Following a  
 37 *balanced alternating learning strategy*, we train models to align representations for the original video  
 38 and its temporal derivatives, such that the learned representations are steered to dynamic information  
 39 in the images (section 1). Herein, we refer to this approach as Video Time-Differentiation for Instance  
 40 Discrimination (ViDiDi). Our method demonstrates excellent generalizability and data efficacy on  
 41 standard benchmarks including action recognition (section 1), video retrieval, and action detection.

## 42 2 Approach

### 43 2.1 A Thought Experiment on Physical Motion

44 Our idea is analogous to and inspired by physical motion. Consider a 1-D toy example through a  
 45 thought experiment. Imagine that we observe the free fall motion of a ball. The zeroth, first, and  
 46 second derivatives are the position  $y(t)$ , velocity  $v(t)$ , and acceleration  $a(t)$ , respectively. Given the  
 47 physical law, the free fall motion is governed by  $y(t) = y_0 + v_0 t - \frac{1}{2} g t^2$ , including three latent  
 48 factors: the initial position  $y_0$ , the initial velocity  $v_0$ , and the gravity  $g$ . Inferring the representation  
 49 shared across these different views reveals gravity  $g$ , which is the defining feature of the dynamics.

### 50 2.2 The ViDiDi Framework

51 ViDiDi involves 1) creating multiple views from videos through augmentation and differentiation, 2)  
 52 a balanced alternating learning strategy for pair-wise encoding of the different views into consistent  
 53 representations, and 3) plugging this strategy into existing instance discrimination methods, including  
 54 SimCLR [10], BYOL [11], and VICReg [13]. See fig. 2 for an overview.

55 **Creating multiple views from videos.** i) *Augmentation*: Given a batch of videos, we sample  
 56 each video by randomly cropping two clips. For each clip, we apply a random set of spatial  
 57 augmentations such as random crops to create two batches of augmented clips, denoted as  $V$   
 58 and  $V'$ . ii) *Differentiation*: We further conduct temporal differentiation on clips. For every augmented  
 59 clip within either  $V$  or  $V'$ , we evaluate its  $0^{th}$ ,  $1^{st}$ , or  $2^{nd}$  temporal derivatives. We approximate  
 60 temporal derivatives with finite forward differences. We make sure within one batch, the clips are all  
 61 original frames or derivatives in the same order. Thus, we can now denote the three possible views  
 62 for a batch of clips as  $V$ ,  $\frac{\partial V}{\partial t}$  and  $\frac{\partial^2 V}{\partial t^2}$ .

63 **Balanced alternating learning strategy.** We design a pairing schedule for learning representa-  
 64 tions among different views. Specifically, we define seven pairs:  $(V, V')$ ,  $(V, \frac{\partial V'}{\partial t})$ ,  $(\frac{\partial V}{\partial t}, V')$ ,  
 65  $(\frac{\partial V}{\partial t}, \frac{\partial V'}{\partial t})$ ,  $(\frac{\partial V}{\partial t}, \frac{\partial^2 V'}{\partial t^2})$ ,  $(\frac{\partial^2 V}{\partial t^2}, \frac{\partial V'}{\partial t})$ ,  $(\frac{\partial^2 V}{\partial t^2}, \frac{\partial^2 V'}{\partial t^2})$ .  $(X, X')$  will be chosen from one of the seven

66 pairs above, and serve as inputs to the video encoder for learning using methods such as SimCLR.  
 67 At each step,  $(\mathbf{X}, \mathbf{X}')$  is decided following alg. 1. Intuitively we choose this strategy to let learning  
 68 of derivatives guide original frames. We empirically verify that this balanced alternating learning  
 69 strategy plays an important role in the learning process in the ablation study in appendix C.

70 **Plug into existing instance discrimination methods.** We plug ViDiDi into existing instance  
 71 discrimination frameworks, namely ViDiDi-SimCLR, ViDiDi-BYOL, and ViDiDi-VIC. We describe  
 72 ViDiDi-SimCLR below and refer to the appendix for details about other models.  $(\mathbf{X}, \mathbf{X}')$  is input  
 73 to the video encoder  $f$  and then the projection head  $h$ , yielding paired embeddings  $(\mathbf{Z}, \mathbf{Z}')$  for  
 74 evaluating the loss  $\mathcal{L}_{NCE}$  and training the networks. We discuss more details in appendix D.

### 75 3 Experiments

#### 76 3.1 Experiment Setup

77 We train and evaluate ViDiDi using human action video datasets. **UCF101** [1] includes 13k videos  
 78 from 101 classes. **HMDB51** [2] contains 7k videos from 51 classes. In addition, we also use larger  
 79 and more diverse datasets, **K400** [6], aka Kinetics400, including 240k videos from 400 classes,  
 80 and **K200-40k**, including 40k videos from 200 classes, as a subset of Kinetics 400, helps verify  
 81 data-efficiency. In our experiments, we pretrain models with UCF101, K400, or K200-40k and then  
 82 test them with UCF101 or HMDB51, using split 1 for both datasets. **AVA** contains 280K videos from  
 83 60 action classes, each video is annotated with spatiotemporal localization of human actions. After  
 84 pertaining on Kinetics or UCF101, we follow the evaluation protocol in previous works [3, 7, 8],  
 85 including three types of downstream tasks. i) *video retrieval* on UCF101 and HMDB51, ii) *action*  
 86 *recognition*, and iii) *Action detection* on AVA. More details are in appendix E.

#### 87 3.2 Results

Table 1: Video retrieval using different methods.

Method	Pretrained	UCF101			HMDB51		
		1	5	10	1	5	10
SimCLR	UCF101	29.6	41.4	49.3	17.5	34.7	45.1
<b>ViDiDi-SimCLR</b>	UCF101	<b>38.3</b>	<b>54.6</b>	<b>64.5</b>	<b>17.5</b>	<b>38.9</b>	<b>52.4</b>
BYOL	UCF101	32.2	43.0	50.5	13.8	31.1	44.4
<b>ViDiDi-BYOL</b>	UCF101	<b>43.7</b>	<b>60.4</b>	<b>70.1</b>	<b>19.3</b>	<b>44.1</b>	<b>56.6</b>
VICReg	UCF101	31.1	43.6	50.9	15.7	33.7	44.5
<b>ViDiDi-VIC</b>	UCF101	<b>47.6</b>	<b>60.9</b>	<b>68.6</b>	<b>19.7</b>	<b>40.5</b>	<b>55.1</b>
VICReg	K400	41.9	56.5	64.8	21.7	44.1	56.1
<b>ViDiDi-VIC</b>	K400	<b>51.2</b>	<b>64.6</b>	<b>72.6</b>	<b>25.0</b>	<b>47.2</b>	<b>60.9</b>
<b>ViDiDi-VIC</b>	K200-40k	<b>49.5</b>	<b>63.4</b>	<b>71.0</b>	<b>24.7</b>	<b>45.4</b>	<b>56.0</b>

Table 2: Video retrieval using different backbones backbones.

Net	Method	UCF101			HMDB51		
		1	5	10	1	5	10
R(2+1)D-18	VICReg	30.2	44.1	51.4	15.8	33.6	45.5
	<b>ViDiDi-VIC</b>	<b>47.2</b>	<b>62.6</b>	<b>69.8</b>	<b>20.6</b>	<b>44.1</b>	<b>57.7</b>
MC3-18	VICReg	31.9	44.4	51.4	15.6	35.6	46.1
	<b>ViDiDi-VIC</b>	<b>44.1</b>	<b>59.8</b>	<b>68.0</b>	<b>20.3</b>	<b>40.3</b>	<b>53.4</b>
S3D	VICReg	29.2	41.9	49.2	12.8	29.8	40.9
	<b>ViDiDi-VIC</b>	<b>42.9</b>	<b>59.0</b>	<b>67.5</b>	<b>18.4</b>	<b>38.4</b>	<b>51.4</b>

Table 3: Action detection using different methods.

Method	VICReg	<b>ViDiDi-VIC</b>	SimCLR	<b>ViDiDi-SimCLR</b>	BYOL	<b>ViDiDi-BYOL</b>
mAP	0.089	<b>0.106</b>	0.079	<b>0.094</b>	0.087	<b>0.118</b>

88 **Superior performance, generalizability, and data efficiency.** ViDiDi learns effective video  
 89 representations with limited data. It can be plugged into multiple existing frameworks based on  
 90 instance discrimination, using multiple encoder architectures, and improving the performance of  
 91 various downstream tasks significantly as shown in fig. 1, table 3, table 1, and table 2. Besides, as

92 presented in table 5, table 4, pretrained on small dataset UCF101 or K200-40k, ViDiDi surpasses  
 93 prior video representation learning works pretrained on large-scale K400 or K600 dataset. Further,  
 94 ViDiDi-VIC pretrained on UCF101 or K200-40k outperforms its baseline method VICReg pretrained  
 95 on K400, and also reaches compatible performance as ViDiDi-VIC pretrained on K400, as shown in  
 96 table 1. In summary, ViDiDi is highly generalizable, and efficiently uncovers dynamic features using  
 97 limited data. To gain insights into how ViDiDi works, we further analyze the features and attention  
 98 of video encoders learned via ViDiDi.

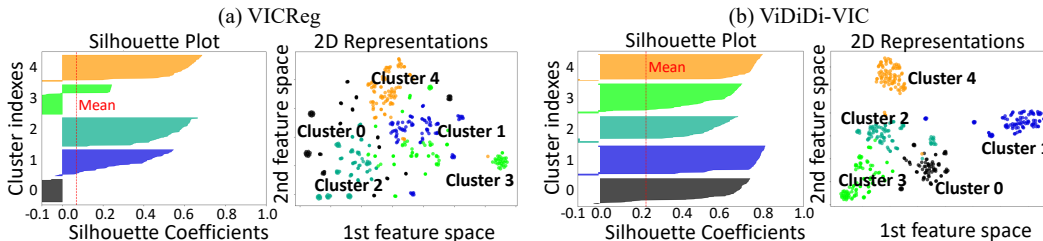


Figure 3: Silhouette scores and t-SNE of top 5 classes from VICReg (left) and ViDiDi-VIC (right).

99 **Better separation based on dynamic features.** We use t-SNE [20] to visualize the representations  
 100 from five classes of videos. As shown in fig. 3, representations learned with ViDiDi are better  
 101 clustered by action classes. We also use Silhouette score [21] in table 7 to quantify the degree of  
 102 separation. To study whether such separation results from capturing better dynamic features, we  
 103 visualize the spatiotemporal attention using Saliency Tubes [22]. ViDiDi leads the model to attend to  
 104 dynamic aspects of the video, such as motions and interactions, rather than static backgrounds as  
 105 shown in fig. 4. These results align with our intuition that ViDiDi attends to dynamic parts and avoids  
 106 learning static content as a learning shortcut, resulting in efficient utilization of data. We attach more  
 107 visualization and analysis results in appendix F.

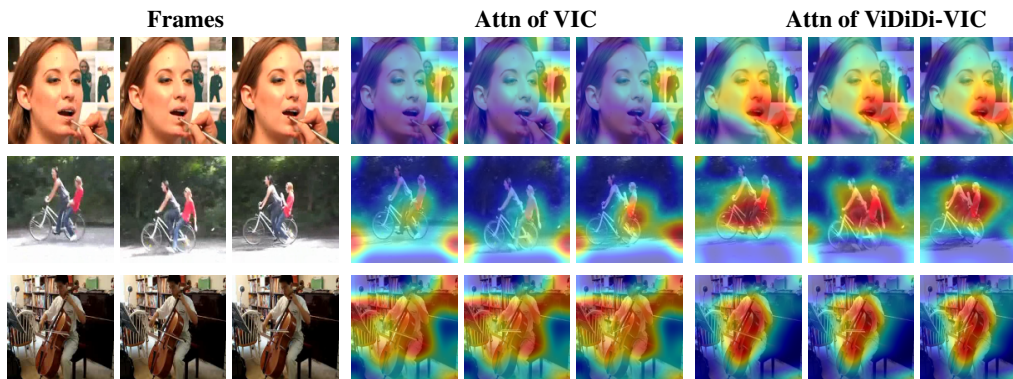


Figure 4: Spatiotemporal attention on UCF and HMDB51.

## 108 4 Conclusion

109 In this paper, we introduce ViDiDi, a novel, data-efficient, and generalizable framework for self-  
 110 supervised video representation. We utilize the Taylor series to unfold multiple views from a  
 111 video through different orders of temporal derivatives and learn representations among these views  
 112 following a balanced alternating learning strategy. ViDiDi steers the video encoder to dynamic  
 113 features instead of static shortcuts, enhancing performance on common video representation learning  
 114 tasks significantly. We identify multiple future directions as well, such as applying ViDiDi to other  
 115 modalities, and exploring other vision tasks that require a more fine-grained understanding of video  
 116 dynamics [23, 24]. Furthermore, a potentially fruitful direction is to use this approach to learn  
 117 intuitive physics, better supporting agents to understand, predict, and interact with the physical world.

## References

- 118
- 119 [1] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions  
120 classes from videos in the wild. *CoRR*, abs/1212.0402, 2012.
- 121 [2] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. Hmdb: A large video database for human  
122 motion recognition. In *ICCV*, pages 2556–2563, 2011.
- 123 [3] Sheng Guo, Zihua Xiong, Yujie Zhong, Limin Wang Wang, Xiaobo Guo, Bing Han, and Huang Weilin.  
124 Cross-architecture self-supervised video representation learning. In *CVPR*, June 2022.
- 125 [4] Zhenzhi Wang, Limin Wang, Tao Wu, Tianhao Li, and Gangshan Wu. Negative sample matters: A  
126 renaissance of metric learning for temporal grounding, 2021.
- 127 [5] Chunhui Gu, Chen Sun, Sudheendra Vijayanarasimhan, Caroline Pantofaru, David A. Ross, George  
128 Toderici, Yeqing Li, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, and Jitendra Malik. Ava: A  
129 video dataset of spatio-temporally localized atomic visual actions. *CVPR*, pages 6047–6056, 2017.
- 130 [6] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan,  
131 Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The  
132 kinetics human action video dataset, 2017.
- 133 [7] Tengda Han, Weidi Xie, and Andrew Zisserman. Self-supervised co-training for video representation  
134 learning. In *Neurips*, 2020.
- 135 [8] Rui Qian, Tianjian Meng, Boqing Gong, Ming-Hsuan Yang, Huisheng Wang, Serge Belongie, and Yin Cui.  
136 Spatiotemporal contrastive video representation learning, 2021.
- 137 [9] Ishan Rajendrakumar Dave, Rohit Gupta, Mamshad Nayeem Rizve, and Mubarak Shah. Tclr: Temporal  
138 contrastive learning for video representation. *Comput. Vis. Image Underst.*, 219:103406, 2021.
- 139 [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for  
140 contrastive learning of visual representations. In Hal Daumé III and Aarti Singh, editors, *International  
141 Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*,  
142 pages 1597–1607. PMLR, 13–18 Jul 2020.
- 143 [11] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya,  
144 Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, koray  
145 kavukcuoglu, Remi Munos, and Michal Valko. Bootstrap your own latent - a new approach to self-  
146 supervised learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *NeurIPS*,  
147 volume 33, pages 21271–21284. Curran Associates, Inc., 2020.
- 148 [12] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised  
149 visual representation learning. In *CVPR*, pages 9726–9735, 2020.
- 150 [13] Adrien Bardes, Jean Ponce, and Yann LeCun. VICReg: Variance-invariance-covariance regularization for  
151 self-supervised learning. In *ICLR*, 2022.
- 152 [14] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsuper-  
153 vised learning of visual features by contrasting cluster assignments. In *NeurIPS*, NeurIPS’20, Red Hook,  
154 NY, USA, 2020. Curran Associates Inc.
- 155 [15] Allan Jabri, Andrew Owens, and Alexei Efros. Space-time correspondence as a contrastive random walk.  
156 In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *NeurIPS*, volume 33, pages  
157 19545–19560. Curran Associates, Inc., 2020.
- 158 [16] Shengyi Qian, Linyi Jin, Chris Rockwell, Siyi Chen, and David F. Fouhey. Understanding 3d object  
159 articulation in internet videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
160 Pattern Recognition (CVPR)*, pages 1599–1609, June 2022.
- 161 [17] Zheyun Qin, Xiankai Lu, Xiushan Nie, Yilong Yin, and Jianbing Shen. Exposing the self-supervised  
162 space-time correspondence learning via graph kernels. *AAAI*, 37(2):2110–2118, Jun. 2023.
- 163 [18] Jing Tan, Jiaqi Tang, Limin Wang, and Gangshan Wu. Relaxed transformer decoders for direct action  
164 proposal generation. *ICCV*, pages 13506–13515, 2021.
- 165 [19] T. Lin, X. Liu, X. Li, E. Ding, and S. Wen. Bmn: Boundary-matching network for temporal action proposal  
166 generation. In *ICCV*, pages 3888–3897, Los Alamitos, CA, USA, nov 2019. IEEE Computer Society.

- 167 [20] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning*  
168 *Research*, 9(86):2579–2605, 2008.
- 169 [21] Peter J. Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis.  
170 *Journal of Computational and Applied Mathematics*, 20:53–65, 1987.
- 171 [22] Alexandros Stergiou, Georgios Kapidis, Grigorios Kalliatakis, Christos Chrysoulas, Remco C. Veltkamp,  
172 and Ronald Poppe. Saliency tubes: Visual explanations for spatio-temporal convolutions. *2019 IEEE*  
173 *International Conference on Image Processing (ICIP)*, pages 1830–1834, 2019.
- 174 [23] H. Zhang, D. Liu, Q. Zheng, and B. Su. Modeling video as stochastic processes for fine-grained video  
175 representation learning. In *CVPR*, pages 2225–2234, Los Alamitos, CA, USA, jun 2023. IEEE Computer  
176 Society.
- 177 [24] Hilde Kuehne, Ali Bilgin Arslan, and Thomas Serre. The language of actions: Recovering the syntax  
178 and semantics of goal-directed human activities. *2014 IEEE Conference on Computer Vision and Pattern*  
179 *Recognition*, pages 780–787, 2014.
- 180 [25] Tian Pan, Yibing Song, Tianyu Yang, Wenhao Jiang, and Wei Liu. Videomoco: Contrastive video  
181 representation learning with temporally adversarial examples. In *CVPR*, pages 11205–11214, June 2021.
- 182 [26] Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross Girshick, and Kaiming He. A large-scale study on  
183 unsupervised spatiotemporal representation learning. In *CVPR*, pages 3298–3308, 2021.
- 184 [27] Ceyuan Yang, Yinghao Xu, Bo Dai, and Bolei Zhou. Video representation learning with visual tempo  
185 consistency, 2020.
- 186 [28] S. Jenni and H. Jin. Time-equivariant contrastive video representation learning. In *ICCV*, pages 9950–9960,  
187 Los Alamitos, CA, USA, oct 2021. IEEE Computer Society.
- 188 [29] Di Yang, Yaohui Wang, Quan Kong, Antitza Dantcheva, Lorenzo Garattoni, Gianpiero Francesca, and  
189 François Brémond. Self-supervised video representation learning via latent time navigation. In *AAAI*,  
190 *AAAI’23/IAAI’23/EAAI’23*. AAAI Press, 2023.
- 191 [30] Jinpeng Wang, Yuting Gao, Ke Li, Xinyang Jiang, Xiao-Wei Guo, Rongrong Ji, and Xing Sun. Enhancing  
192 unsupervised video representation learning by decoupling the scene and the motion. In *AAAI*, 2020.
- 193 [31] Minghao Zhu, Xiao Lin, Ronghao Dang, Chengju Liu, and Qijun Chen. Fine-grained spatiotemporal motion  
194 alignment for contrastive video representation learning. In *Proceedings of the 31st ACM International*  
195 *Conference on Multimedia*, MM ’23, page 4725–4736, New York, NY, USA, 2023. Association for  
196 Computing Machinery.
- 197 [32] S. Jenni, Givi Meishvili, and Paolo Favaro. Video representation learning by recognizing temporal  
198 transformations. In *ECCV*, 2020.
- 199 [33] Longlong Jing, Xiaodong Yang, Jingen Liu, and Yingli Tian. Self-supervised spatiotemporal feature  
200 learning via video rotation prediction, 2019.
- 201 [34] Sagie Benaim, Ariel Ephrat, Oran Lang, Inbar Mosseri, William T. Freeman, Michael Rubinstein, Michal  
202 Irani, and Tali Dekel. Speednet: Learning the speediness in videos. In *CVPR*, June 2020.
- 203 [35] Peihao Chen, Deng Huang, Dongliang He, Xiang Long, Runhao Zeng, Shilei Wen, Mingkui Tan, and  
204 Chuang Gan. Rspnet: Relative speed perception for unsupervised video representation learning. In *AAAI*,  
205 2020.
- 206 [36] Jiangliu Wang, Jianbo Jiao, and Yun-Hui Liu. Self-supervised video representation learning by pace  
207 prediction. In *ECCV*, page 504–521, Berlin, Heidelberg, 2020. Springer-Verlag.
- 208 [37] Y. Yao, C. Liu, D. Luo, Y. Zhou, and Q. Ye. Video playback rate perception for self-supervised spatio-  
209 temporal representation learning. In *CVPR*, pages 6547–6556, Los Alamitos, CA, USA, jun 2020. IEEE  
210 Computer Society.
- 211 [38] Hyeon Cho, Taehoon Kim, Hyung Chang, and Wonjun Hwang. Self-supervised visual learning by variable  
212 playback speeds prediction of a video. *IEEE Access*, PP:1–1, 05 2021.
- 213 [39] Ishan Misra, C. Zitnick, and Martial Hebert. Shuffle and learn: Unsupervised learning using temporal  
214 order verification. In *ECCV*, volume 9905, pages 527–544, 10 2016.

- 215 [40] B. Fernando, H. Bilen, E. Gavves, and S. Gould. Self-supervised video representation learning with  
216 odd-one-out networks. In *CVPR*, pages 5729–5738, Los Alamitos, CA, USA, jul 2017. IEEE Computer  
217 Society.
- 218 [41] Tomoyuki Suzuki, Takahiro Itazuri, Kensho Hara, and Hirokatsu Kataoka. Learning spatiotemporal 3d  
219 convolution with video order self-supervision. In *ECCV Workshops*, 2018.
- 220 [42] Ishan Rajendrakumar Dave, Simon Jenni, and Mubarak Shah. No more shortcuts: Realizing the potential  
221 of temporal self-supervision, 2023.
- 222 [43] Hsin-Ying Lee, Jia-Bin Huang, Maneesh Kumar Singh, and Ming-Hsuan Yang. Unsupervised representa-  
223 tion learning by sorting sequence. In *ICCV*, 2017.
- 224 [44] Dejing Xu, Jun Xiao, Zhou Zhao, Jian Shao, Di Xie, and Yueting Zhuang. Self-supervised spatiotemporal  
225 learning via video clip order prediction. In *CVPR*, June 2019.
- 226 [45] Dahun Kim, Donghyeon Cho, and In So Kweon. Self-supervised video representation learning with  
227 space-time cubic puzzles. In *AAAI, AAAI’19/IAAI’19/EAAI’19*. AAAI Press, 2019.
- 228 [46] Unaiza Ahsan, Rishi Madhok, and Irfan Essa. Video jigsaw: Unsupervised learning of spatiotemporal  
229 context for video action recognition. *2019 IEEE Winter Conference on Applications of Computer Vision*  
230 (*WACV*), pages 179–189, 2018.
- 231 [47] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked autoencoders are data-efficient  
232 learners for self-supervised video pre-training. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and  
233 Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- 234 [48] X. Sun, P. Chen, L. Chen, C. Li, T. H. Li, M. Tan, and C. Gan. Masked motion encoding for self-supervised  
235 video representation learning. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
236 (*CVPR*), pages 2235–2245, Los Alamitos, CA, USA, jun 2023. IEEE Computer Society.
- 237 [49] David Fan, Jue Wang, Leo Liao, Yi Zhu, Vimal Bhat, Hector Santos, Rohith Mysore Vijaya Kumar, and  
238 Xinyu (Arthur) Li. Motion-guided masking for spatiotemporal representation learning. In *ICCV 2023*,  
239 2023.
- 240 [50] Jean-Baptiste Alayrac, Adrià Recasens, Rosalia Schneider, Relja Arandjelović, Jason Ramapuram, Jef-  
241 frey De Fauw, Lucas Smaira, Sander Dieleman, and Andrew Zisserman. Self-supervised multimodal  
242 versatile networks. In *NeurIPS, NeurIPS’20*, Red Hook, NY, USA, 2020. Curran Associates Inc.
- 243 [51] Adrià Recasens, Pauline Luc, Jean-Baptiste Alayrac, Luyu Wang, Florian Strub, Corentin Tallec, Mateusz  
244 Malinowski, Viorica Pătrăucean, Florent Alché, Michal Valko, Jean-Bastien Grill, Aaron van den Oord,  
245 and Andrew Zisserman. Broaden your views for self-supervised video learning. In *ICCV*, pages 1235–1245,  
246 2021.
- 247 [52] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from  
248 self-supervised synchronization. In *NeurIPS, NeurIPS’18*, page 7774–7785, Red Hook, NY, USA, 2018.  
249 Curran Associates Inc.
- 250 [53] Tianhao Li and Limin Wang. Learning spatiotemporal features via video and text pair discrimination.  
251 *ArXiv*, abs/2001.05691, 2020.
- 252 [54] A. Miech, J. Alayrac, L. Smaira, I. Laptev, J. Sivic, and A. Zisserman. End-to-end learning of visual  
253 representations from uncurated instructional videos. In *CVPR*, pages 9876–9886, Los Alamitos, CA, USA,  
254 jun 2020. IEEE Computer Society.
- 255 [55] Peng Jin, Jinfa Huang, Fenglin Liu, Xian Wu, Shen Ge, Guoli Song, David A. Clifton, and Jie Chen.  
256 Expectation-maximization contrastive learning for compact video-and-language representations. In Alice H.  
257 Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *NeurIPS*, 2022.
- 258 [56] Rui Li, Yiheng Zhang, Zhaofan Qiu, Ting Yao, Dong Liu, and Tao Mei. Motion-focused contrastive  
259 learning of video representations\*. *ICCV*, 2021.
- 260 [57] Jingcheng Ni, Nana Zhou, Jie Qin, Qianrun Wu, Junqi Liu, Boxun Li, and Di Huang. Motion sensitive  
261 contrastive learning for self-supervised video representation. In *ECCV*, 2022.
- 262 [58] X. Sui, S. Li, X. Geng, Y. Wu, X. Xu, Y. Liu, R. Goh, and H. Zhu. Craft: Cross-attentional flow transformer  
263 for robust optical flow. In *CVPR*, pages 17581–17590, Los Alamitos, CA, USA, jun 2022. IEEE Computer  
264 Society.

- 265 [59] Kecheng Zheng, Yang Cao, Kai Zhu, Ruijing Zhao, and Zhengjun Zha. Famlp: A frequency-aware mlp-like  
266 architecture for domain generalization. *ArXiv*, abs/2203.12893, 2022.
- 267 [60] Qinwei Xu, Ruipeng Zhang, Ya Zhang, Yanfeng Wang, and Qi Tian. A fourier-based framework for  
268 domain generalization. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
269 pages 14378–14387, 2021.
- 270 [61] Ju-Hyeon Nam and Sang-Chul Lee. Frequency filtering for data augmentation in x-ray image classification.  
271 In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 81–85, 2021.
- 272 [62] Jinhung Kim, Taeh Kim, Minh Shim, Dongyoon Han, Dongyoon Wee, and Junmo Kim. Frequency  
273 selective augmentation for video representation learning. In *AAAI*, 2022.
- 274 [63] Haoxin Li, Yuan Liu, Hanwang Zhang, and Boyang Li. Mitigating and evaluating static bias of action  
275 representations in the background and the foreground, 2023.
- 276 [64] Lei Wang, Xiuyuan Yuan, Tom Gedeon, and Liang Zheng. Taylor videos for action recognition. In  
277 *Forty-first International Conference on Machine Learning (ICML)*, 2024.
- 278 [65] Dezhao Luo, Chang Liu, Yu Zhou, Dongbao Yang, Can Ma, Qixiang Ye, and Weiping Wang. Video cloze  
279 procedure for self-supervised spatio-temporal learning. In *AAAI*, 2020.
- 280 [66] Tengda Han, Weidi Xie, and Andrew Zisserman. Video representation learning by dense predictive coding.  
281 *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 1483–1492,  
282 2019.
- 283 [67] Chengxu Zhuang, Tianwei She, Alex Andonian, Max Sobol Mark, and Daniel Yamins. Unsupervised  
284 learning from video with deep neural embeddings. In *CVPR*, June 2020.
- 285 [68] Hanwen Liang, Niamul Quader, Zhixiang Chi, Lizhe Chen, Peng Dai, Juwei Lu, and Yang Wang. Self-  
286 supervised spatiotemporal representation learning by exploiting video continuity. *AAAI*, 36(2):1564–1573,  
287 Jun. 2022.
- 288 [69] Wei Li, Dezhao Luo, Bo Fang, Yu Zhou, and Weiping Wang. Video 3d sampling for self-supervised  
289 representation learning. *ArXiv*, abs/2107.03578, 2021.
- 290 [70] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing  
291 internal covariate shift. In Francis Bach and David Blei, editors, *International Conference on Machine  
292 Learning (ICML)*, volume 37 of *Proceedings of Machine Learning Research*, pages 448–456, Lille, France,  
293 07–09 Jul 2015. PMLR.
- 294 [71] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts, 2017.
- 295 [72] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks, 2017.
- 296 [73] Herbert E. Robbins. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407,  
297 1951.
- 298 [74] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object  
299 detection with region proposal networks. *CoRR*, abs/1506.01497, 2015.
- 300 [75] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101,  
301 2017.



## Appendix

---

**Algorithm 1: Differentiation at Each Batch**


---

**Data:**  $\text{epoch} \geq 0, (V, V')$   
**Result:**  $(X, X')$   
 // Deterministic differentiation step  
 1 **if**  $\text{epoch}\%4 = 0$  **then**  $(X, X') \leftarrow \left(\frac{\partial V}{\partial t}, \frac{\partial V'}{\partial t}\right)$ ;  
 2 **else if**  $\text{epoch}\%4 = 1$  **then**  $(X, X') \leftarrow \left(\frac{\partial V}{\partial t}, V'\right)$ ;  
 3 **else if**  $\text{epoch}\%4 = 2$  **then**  $(X, X') \leftarrow \left(V, \frac{\partial V'}{\partial t}\right)$ ;  
 4 **else**  $(X, X') \leftarrow (V, V')$ ;  
 // Additional random differentiation step  
 5  $\epsilon \leftarrow \text{rand}(0, 1)$ ;  
 6 **if**  $\epsilon < 0.5$  **then**  $(X, X') \leftarrow \left(\frac{\partial X}{\partial t}, \frac{\partial X'}{\partial t}\right)$ ;

---

### 303 A Relationship to Prior Arts

304 Current methods for SSL of video representations mainly utilize instance discrimination, pretext  
 305 tasks, multimodal learning, and other ones. In the following, we discuss these results, and highlight  
 306 the advantages of our method over existing ones.

307 **Instance discrimination.** It is first applied to images [10, 11, 12, 14] and then to videos [25, 26, 8, 3,  
 308 27, 7, 28]. Given either images or videos as instances, models learn to discriminate different instances  
 309 versus different "views" of the same instance, where the views are generated by spatio-temporal  
 310 augmentations [25, 8]. The learning process is driven by contrastive learning [8], clustering [25], or  
 311 teacher-student distillation [26]. Recognizing the rich dynamics in videos, some prior works have  
 312 further modified the loss function to consider each video's temporal attributes, such as play speed [27],  
 313 time differences [8, 29, 9, 26], frame order [30], and motion diversity [9, 31]. Such modifications  
 314 do not fully capture the essence of videos as reflections of continuous real-world dynamics, and  
 315 are usually designed for a specific instance discrimination method. Apart from being applicable to  
 316 different frameworks of instance discrimination, our approach, is new for its extraction of continuous  
 317 dynamics by unfolding a video's hidden views via Taylor expansion and temporal differentiation.

318 **Pretext tasks.** Another category of methods involves creating learning tasks from videos. These  
 319 tasks have many possible variations, such as identifying transformations applied to videos [32, 33],  
 320 predicting the speed of videos [34, 35, 36, 37, 38], identifying incorrect ordering of frames or clips  
 321 [39, 40, 41, 42], resorting them in order [43, 44], and solving space-time puzzles [45, 46]. The above  
 322 methods usually require a complex combination of different tasks to learn general representations,  
 323 while some recent works utilize large transformer backbones and learn by reconstructing masked  
 324 areas [47] and further incorporating motion guidance into masking or reconstruction [48, 49]. These  
 325 tasks provide non-trivial challenges for models to learn but are unlikely to reflect the natural processes  
 326 through which humans and machines alike may learn and interpret dynamic visual information. In  
 327 contrast, ViDiDi uses simple learning objectives and models how the physical world can be intuitively  
 328 processed and understood without the need for complex tasks.

329 **Others.** In addition, prior methods also learn to align videos with other modalities, such as audio  
 330 tracks [50, 51, 52], video captions [50, 53, 54, 55], and optical flows [7, 51, 56, 57, 58]. Optical flow  
 331 also models changes between frames and is related to our method. However, our method is easy to  
 332 calculate and intuitively generalize to higher orders of motion and guides the learning of the original  
 333 frames within the same encoder in contrast to an additional encoder for optical flow [7]. Besides,  
 334 our temporal differentiation strategy may be flexibly adapted to other dynamic data such as audio  
 335 while optical flow explicitly models the movements of pixels. Incidentally, our proposed *balanced*  
 336 *alternating learning strategy* as a simple yet novel way of learning different types of data, may inspire  
 337 multimodal learning. Some other existing works manipulate frequency content to create augmented  
 338 views of images [59, 60, 61] or videos [62, 63, 56] to make models more robust to out-of-domain

Table 4: ViDiDi surpasses previous works on action recognition after finetuning.

Method	Net	Input	Pretrained	UCF	HMDB
VCOP [44]	R3D-18	16 × 112	UCF101	64.9	29.5
VCP [65]	R3D-18	16 × 112	UCF101	66.0	31.5
3D-RotNet [33]	R3D-18	16 × 112	K400	66.0	37.1
DPC [66]	R3D-18	25 × 128	K400	68.2	34.5
VideoMoCo [25]	R3D-18	32 × 112	K400	74.1	43.6
RTT [32]	R3D-18	16 × 112	K600	79.3	49.8
VIE [67]	R3D-18	16 × 112	K400	72.3	44.8
RSPNet [35]	R3D-18	16 × 112	K400	74.3	41.8
VTHCL [27]	R3D-18	8 × 224	K400	80.6	48.6
CPNet [68]	R3D-18	16 × 112	K400	80.8	52.8
CPNet [68]	R3D-18	16 × 112	UCF101	77.2	46.3
CACL [3]	T+C3D	16 × 112	K400	77.5	-
CACL [3]	T+R3D	16 × 112	UCF101	77.5	43.8
TCLR [9]	R3D-18	16 × 112	UCF101	82.4	52.9
<b>ViDiDi-BYOL</b>	R3D-18	16 × 112	UCF101	<b>83.4</b>	<b>58.0</b>
<b>ViDiDi-VIC</b>	R3D-18	16 × 112	UCF101	<b>82.3</b>	<b>53.4</b>
<b>ViDiDi-VIC</b>	R3D-18	16 × 112	K200-40k	<b>82.7</b>	<b>54.2</b>
<b>ViDiDi-VIC</b>	R3D-18	16 × 112	K400	<b>83.2</b>	<b>55.8</b>
VCOP [44]	R(2+1)D-18	16 × 112	UCF101	72.4	30.9
VCP [65]	R(2+1)D-18	16 × 112	UCF101	66.3	32.2
PacePred [36]	R(2+1)D-18	16 × 112	K400	77.1	36.6
VideoMoCo [25]	R(2+1)D-18	32 × 112	K400	78.7	49.2
V3S [69]	R(2+1)D-18	16 × 112	K400	79.2	40.4
RSPNet [35]	R(2+1)D-18	16 × 112	K400	81.1	44.6
RTT [32]	R(2+1)D-18	16 × 112	UCF101	81.6	46.4
CPNet [68]	R(2+1)D-18	16 × 112	UCF101	81.8	51.2
CACL [3]	T+R(2+1)D	16 × 112	UCF101	82.5	48.8
<b>ViDiDi-VIC</b>	R(2+1)D-18	16 × 112	UCF101	<b>83.0</b>	<b>54.9</b>

339 data. Related to but unlike these works, we seek a fundamental and computationally efficient strategy  
 340 to construct views from videos that reflect the continuous nature of real-world dynamics. Besides,  
 341 apart from a new way of processing data, we propose an alternating learning strategy that is pivotal to  
 342 boosting learning and has not been explored in previous works. Upon wrapping up our work, we  
 343 noticed a concurrent work [64] recovers views from videos inspired by the Taylor series expansion,  
 344 but from a complementary perspective, where derivatives in different orders are combined together.  
 345 Besides, they utilize the augmented input in a supervised way, while ours focuses on a generalizable  
 346 self-supervised framework including creating new views and a learning strategy.

## 347 B Comparisons with Previous Works

348 Results on both video retrieval (table 5) and action recognition (table 4) suggest that ViDiDi outper-  
 349 forms prior models. Compared to other models trained on Kinetics, ViDiDi-VIC achieves the highest  
 350 accuracy using K400, while also reaching compatible performance using UCF101 or K200-40k  
 351 subset for pretraining. Besides, ViDiDi-BYOL achieves the best performance in action recognition  
 352 on HMDB51, by a significant margin of 5.1% over the recent TCLR method [9]. Importantly, ViDiDi  
 353 supports efficient use of data. Its performance gain is most significant in the scenario of training with  
 354 smaller datasets. We discuss this with more details in the following section.

## 355 C Ablation Study

356 ViDiDi involves multiple methodological choices, including 1) the order of derivatives, 2) how to  
 357 pair different orders of derivatives as the input to two-stream video encoders, and 3) how to prescribe  
 358 the learning schedule over different pairings. We perform ablation studies to test each design choice.

359 For the order of derivatives, we consider up to the  $2^{nd}$  derivative. For pairing, we consider pairing  
 360 derivatives in the same order ( $1^{st}$  vs.  $1^{st}$ ,  $2^{nd}$  vs.  $2^{nd}$ , etc.) or between different orders ( $1^{st}$  vs.  $0^{th}$ ,  
 361  $1^{st}$  vs.  $2^{nd}$ , etc.), respectively. For scheduling, we consider either random vs. scheduled selection of  
 362 input pairs. With random selection, temporal differentiation is essentially treated as additional data  
 363 augmentation. In contrast, the scheduled selection (alg. 1) aims to provide a balanced and structured  
 364 way for the model to learn from various orders of temporal derivatives, where higher order derivatives  
 365 are intuitively used as guidance of learning the original frames.

Table 5: **ViDiDi surpasses previous SSL models on video retrieval.** T + C3D means training with an additional transformer.

Method	Net	Pretrained	UCF101			HMDB51		
			1	5	10	1	5	10
SpeedNet [34]	S3D-G	K400	13.0	28.1	37.5	-	-	-
RIT [32]	R3D-18	K600	26.1	48.5	59.1	-	-	-
RSPNet [35]	R3D-18	K400	41.1	59.4	68.4	-	-	-
CoCLR [7]	S3D	K400	46.3	62.8	69.5	20.6	43.0	54.0
CACL [3]	T+C3D	K400	44.2	63.1	71.9	-	-	-
<b>ViDiDi-VIC</b>	R3D-18	K200-40k	<b>49.5</b>	<b>63.4</b>	<b>71.0</b>	<b>24.7</b>	<b>45.4</b>	<b>56.0</b>
<b>ViDiDi-VIC</b>	R3D-18	K400	<b>51.2</b>	<b>64.6</b>	<b>72.6</b>	<b>25.0</b>	<b>47.2</b>	<b>60.9</b>
VCOP [44]	R3D-18	UCF101	14.1	30.3	40.0	7.6	22.9	34.4
VCP [65]	R3D-18	UCF101	18.6	33.6	42.5	7.6	24.4	33.6
PacePred [36]	R3D-18	UCF101	23.8	38.1	46.4	9.6	26.9	41.1
PRP [37]	R3D-18	UCF101	22.8	38.5	46.7	8.2	25.8	38.5
V3S [69]	R3D-18	UCF101	28.3	43.7	51.3	10.8	30.6	42.3
CACL [3]	T+R3D	UCF101	41.1	59.2	67.3	17.6	36.7	48.4
<b>ViDiDi-VIC</b>	R3D-18	UCF101	<b>47.6</b>	<b>60.9</b>	<b>68.6</b>	<b>19.7</b>	<b>40.5</b>	<b>55.1</b>

As shown in table 6, results demonstrate a progressive improvement in the model’s performance in video retrieval, given a higher order of derivatives (from the 1<sup>st</sup> to 2<sup>nd</sup> order), given mixed pairing, and given scheduled selection of input pairs. Therefore, temporal differentiation is not merely another data augmentation trick. Invariance to different orders of temporal derivatives is a valuable principle for SSL of video representations that lead to better performance in downstream tasks. To leverage this principle, it is beneficial to design mixed pairing and prescribe a learning schedule that provides a balanced and holistic view of different orders of temporal dynamics inherent to videos. Details about how we design the groups of models in table 6 are summarized below:

- **Base:** The direct extension of VICReg.
- **+Random 1<sup>st</sup>:** Add 1<sup>st</sup> order derivatives as random augmentation.
- **+Random 1<sup>st</sup> & 2<sup>nd</sup>:** Add 1<sup>st</sup> and 2<sup>nd</sup> order derivatives as random augmentation.
- **Reverse ViDiDi-VIC:** Reverse the order of pair alternation by epoch in ViDiDi, i.e., line 1-9 in alg. 1.
- **+Schedule 1<sup>st</sup>:** Alternate pairs across epochs in the order  $(\frac{\partial \mathbf{V}}{\partial t}, \frac{\partial \mathbf{V}'}{\partial t}) \rightarrow (\mathbf{V}, \mathbf{V}') \rightarrow (\frac{\partial \mathbf{V}}{\partial t}, \frac{\partial \mathbf{V}'}{\partial t}) \rightarrow \dots$
- **+Schedule 1<sup>st</sup> & Mix:** switch pairs by epoch in the order  $(\frac{\partial \mathbf{V}}{\partial t}, \frac{\partial \mathbf{V}'}{\partial t}) \rightarrow (\frac{\partial \mathbf{V}}{\partial t}, \mathbf{V}') \rightarrow (\mathbf{V}, \mathbf{V}') \rightarrow (\frac{\partial \mathbf{V}}{\partial t}, \frac{\partial \mathbf{V}'}{\partial t}) \rightarrow \dots$
- **+Schedule 1<sup>st</sup> & 2<sup>nd</sup> and +Schedule 1<sup>st</sup> & 2<sup>nd</sup> & Mix:** Build upon +Schedule 1<sup>st</sup> and +Schedule 1<sup>st</sup> & Mix accordingly with random differentiation at each batch to utilize 2<sup>nd</sup> order derivatives.

Table 6: **Ablation Study.** Video retrieval performance on UCF101 with different design choices.

Method	UCF101		
	1	5	10
Base (VICReg)	31.1	43.6	50.9
+Random 1 <sup>st</sup>	35.2	47.7	56.1
+Random 1 <sup>st</sup> & 2 <sup>nd</sup>	36.2	48.6	55.8
Reverse ViDiDi-VIC	39.1	54.7	62.9
+Schedule 1 <sup>st</sup>	37.1	50.3	58.2
+Schedule 1 <sup>st</sup> & Mix	39.3	53.2	60.9
+Schedule 1 <sup>st</sup> & 2 <sup>nd</sup>	40.7	56.5	64.0
+Schedule 1 <sup>st</sup> & 2 <sup>nd</sup> & Mix	43.0	59.2	66.6
<b>ViDiDi-VIC</b>	<b>47.6</b>	<b>60.9</b>	<b>68.6</b>

## 386 D Details of SimCLR, BYOL and VICReg

387 In this section, we provide more details on how we plug ViDiDi into different instance discrimination  
 388 frameworks: SimCLR [10], BYOL [11], VICReg [13].

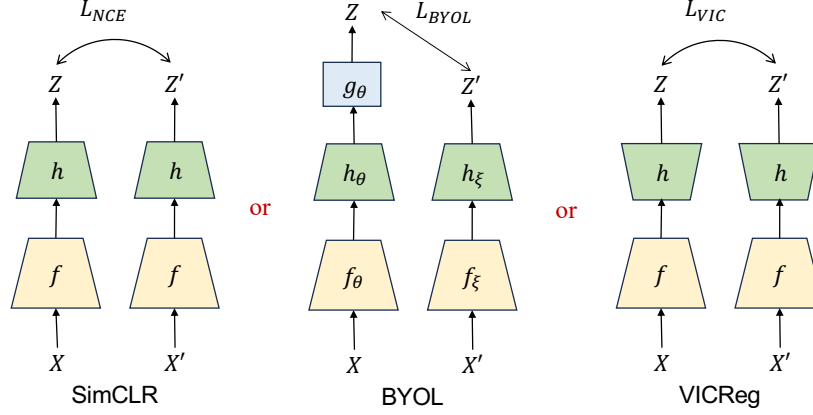


Figure 5: The SimCLR, BYOL, VICReg details.

### 389 D.1 Notation

390 The summary of the SimCLR, BYOL, and VICReg are shown in fig. 5. We begin by introducing the  
 391 notations.  $(X, X')$  represents two batches of input to the discrimination framework, both in the shape  
 392  $\mathbb{R}^{B \times C \times T \times h \times w}$ , containing  $B$  clips (or derivatives) of length  $T$ , and size  $h \times w$ .  $(Z, Z')$  denotes two  
 393 batches of latents encoded from  $(X, X')$ , in the shape of  $\mathbb{R}^{B \times D}$ , containing latents of dimension  $D$   
 394 for  $B$  clips.  $Z = [z_1, \dots, z_B]^T$  and  $Z' = [z'_1, \dots, z'_B]^T$ , expressed as collects of column vectors.  $f$   
 395 represents the encoder, which is a 3D convolutional neural network in our experiments.  $h$  serves as  
 396 the projector, either shrinking or expanding output dimensionality.  $g$  denotes the predictor.  $h$  and  
 397  $g$  are both realized as multi-layer perceptrons (MLPs). We also introduce the similarity function:  
 398  $s_{i,j} = z_i^\top z'_j / (\|z_i\| \|z'_j\|)$ .

### 399 D.2 SimCLR

400 SimCLR [10] is a contrastive learning framework, whose key idea is to contrast dissimilar instances  
 401 in the latent space. As shown in fig. 5, SimCLR uses a shared encoder  $f$  to process  $(X, X')$ , and then  
 402 project the output with an MLP projection head  $h$  into  $(Z, Z')$ .  $Z = h(f(X))$ ,  $Z' = h(f(X'))$ . The  
 403 InfoNCE loss is defined as:

$$\mathcal{L}_{NCE} = \frac{1}{2B} \sum_{i=1}^B \log \frac{\exp(s_{i,i}/\alpha)}{\sum_{j=1}^B \exp(s_{i,j}/\alpha)} + \frac{1}{2B} \sum_{i=1}^B \log \frac{\exp(s_{i,i}/\alpha)}{\sum_{j=1}^B \exp(s_{j,i}/\alpha)} \quad (1)$$

### 404 D.3 BYOL

405 BYOL [11] is a teacher-student approach. It has an online encoder  $f_\theta$ , an online projector  $h_\theta$ , and a  
 406 predictor  $g_\theta$ , learned via gradient descent. BYOL uses stop gradient for a target encoder  $f_\xi$  and a  
 407 target projector  $h_\xi$ , which are updated only by exponential moving average of the online ones  $\xi \leftarrow$   
 408  $\tau\xi + (1 - \tau)\theta$  after each training step, where  $\tau \in [0, 1]$  is the target decay rate.  $Z = g_\theta(h_\theta(f_\theta(X)))$ ,  
 409  $Z' = \text{sg}(h_\xi(f_\xi(X')))$ , here sg means stop gradient. The loss is defined as:

$$\mathcal{L}_{BYOL} = \frac{1}{2B} \sum_{i=1}^B (2 - 2s_{i,j}) \quad (2)$$

410 **D.4 VICReg**

411 VICReg [13] learns to discriminate different instances using direct variance, invariance, and co-  
 412 variance regularization in the latent space. It also has a shared encoder  $f$  and a shared projector  $h$ .  
 413  $Z = h(f(X))$ ,  $Z' = h(f(X'))$ . The invariance term is defined as:

$$s(Z, Z') = \frac{1}{B} \sum_{i=1}^B \|z_i - z'_i\|_2^2 \quad (3)$$

414 The variance term constraints variance along each dimension to be at least  $\gamma$ ,  $\gamma$  is a constant:

$$v(Z) = \frac{1}{D} \sum_{j=1}^D \max(0, \gamma - S(z^j, \epsilon)) \quad (4)$$

415 where  $S$  is the regularized standard deviation  $S(x, \epsilon) = \sqrt{\text{Var}(x) + \epsilon}$ ,  $\epsilon$  is a small constant,  $z^j$  is the  
 416  $j^{\text{th}}$  row vector of  $Z^T$ , containing the value at  $j^{\text{th}}$  dimension for all latents in  $Z$ .

417 The covariance term constraints covariance of different dimensions to be 0:

$$c(Z) = \frac{1}{D} \sum_{i \neq j} [C(Z)]_{i,j}^2 \quad (5)$$

418  $C(Z) = \frac{1}{B-1} \sum_{i=1}^B (z_i - \bar{z})(z_i - \bar{z})^T$ ,  $\bar{z} = \frac{1}{B} \sum_{i=1}^B z_i$ .

419 The total loss is a weighted sum of invariance, variance, and covariance terms:

$$\begin{aligned} \mathcal{L}_{VIC} = & \lambda s(Z, Z') + \mu [v(Z) + v(Z')] \\ & + \nu [c(Z) + c(Z')] \end{aligned} \quad (6)$$

420 **E Implementation Details**

421 **E.1 Augmentation Details**

422 We apply clipwise spatial augmentations as introduced in [8]. All the augmentations are applied  
 423 before differentiation. For example, for a clip sampled from one video, we do a random crop on  
 424 the first frame and crop all the other frames in the clip to the same area as the first frame. If a  
 425 second clip is sampled, we do random crop on its first frame and crop the other frames to the same  
 426 area. The original frames are extracted and resized to have a shorter edge of 150 pixels. The list of  
 427 augmentations is as follows:

- 428 • Random Horizontal Flip, with probability 0.5;
- 429 • Random Sized Crop, with area scale uniformly sampled in the range (0.08, 1), aspect ratio  
 430 in  $(\frac{3}{4}, \frac{4}{3})$ , BILINEAR Interpolation, and output size  $112 \times 112$ ;
- 431 • Gaussian Blur, with probability 0.5, kernal size (3, 3), sigma range (0.1, 2.0);
- 432 • Color Jitter, with probability 0.8, brightness 0.2, contrast 0.2, saturation 0.2, hue 0.05;
- 433 • Random Gray, with probability 0.5;
- 434 • Normalize, mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225].

435 **E.2 Network Architecture**

436 The output feature dimension for R3D-18, R(2+1)D-18, and MC3-18 is 512, while 1024 for S3D. In  
 437 terms of the projector architecture, we use a 2-layer MLP in BYOL, and a 3-layer MLP in SimCLR  
 438 and VICReg, as proposed by [10, 11, 13]. The output dimension of the projector is  $d_{BYOL} =$   
 439  $256$ ,  $d_{SimCLR} = 128$ ,  $d_{VICReg} = 2048$ , and the hidden dimension is  $d_{BYOL} = 4096$ ,  $d_{SimCLR} =$   
 440  $2048$ ,  $d_{VICReg} = 2048$ . The predictor for BYOL is a 2-layer MLP, with output dimension  $d = 256$ ,  
 441 and hidden dimension  $d = 4096$ . Batch normalization [70] and Rectified Linear Unit (ReLU) are  
 442 applied for all hidden layers of projectors and predictors.

### 443 E.3 Pretraining

444 UCF101, K400, or K200-40k is used as the pretraining dataset. We train the model for 400 epochs on  
445 UCF101 or K200-40k, and K400. We set  $T = 8$ , and select 1 frame every 3 frames. The learning  
446 rate follows a cosine decay schedule [71] for all frameworks. The learning rate at  $k_{th}$  iteration is  
447  $\eta \cdot 0.5 \left[ \cos\left(\frac{k}{K}\pi\right) + 1 \right]$ , where  $K$  is the maximum number of iterations and  $\eta$  is the base learning rate.  
448 A 10-epoch warmup is only employed for BYOL. Weight decay is set as  $1e - 6$ . We apply cosine-  
449 annealing of the momentum for BYOL as proposed in [11]:  $\tau = 1 - (1 - \tau_{base}) \cdot \left(\cos\left(\frac{k}{K}\pi\right) + 1\right)/2$ ,  
450 and set  $\tau_{base} = 0.99$ . The temperature for SimCLR is  $\alpha = 0.1$ , and hyper-parameters for VICReg  
451 are  $\lambda = 1.0, \mu = 1.0, \nu = 0.05$ . We train all models with the LARS optimizer [72] utilizing a batch  
452 size of 64 for UCF101 or K200-40k, batch size of 256 for K400, and a base learning rate  $\eta = 1.2$ .  
453 The pretraining can be conducted on 8 GPUs, each having at least 12 GB of memory.

### 454 E.4 Video Retrieval

455 For the pretrained model without any fine-tuning, we test its performance on video retrieval using  
456 nearest-neighborhood in the feature space [3, 7]. Specifically, given a video, we uniformly sample 10  
457 clips of length 16, apply random crop and normalization for data augmentation, encode each clip  
458 using the pretrained video encoder, and average the resulting representations into a single feature  
459 vector for encoding the given video. Through a nearest-neighborhood model that fits the training set,  
460 we use each video in the testing set as a query and retrieve the top-k ( $k = 1, 5, 10$ ) closest videos in  
461 the training set. The retrieval is successful if at least one out of the  $k$  retrieved training videos is from  
462 the same class as the query video. We report the top-k retrieval recall on UCF101 and HMDB51. The  
463 retrieval can be conducted on 1 GPU, having at least 24 GB of memory.

### 464 E.5 Action Recognition

465 We also fine-tune the pretrained model to classify human actions. For this purpose, we add a linear  
466 classification head to the pretrained model, and fine-tune it end-to-end on UCF101 or HMDB51  
467 for 100 epochs (see more details in the supplementary material). At training, we sample clips of  
468 length 16. We use the SGD optimizer [73] with a momentum value of 0.9. The model is tuned  
469 for 100 epochs. The batch size is set at 128, with an initial learning rate of 0.2 which is scaled by  
470  $\frac{1}{10}$  at the 60th and 80th epochs. We use a weight decay of  $1e - 4$ . Furthermore, a dropout rate of  
471 0.5 is applied. After fine-tuning, we sample 10 clips of length 16 from each testing video, apply  
472 random crop and normalization, feed the results as the input to the fine-tuned model, and average their  
473 resulting predictions for the final classification of the video. We report the top-1 action recognition  
474 accuracy on UCF101 and HMDB51. The finetuning can be conducted on 8 GPUs, each having at  
475 least 12 GB of memory. The testing can be conducted on 1 GPU.

### 476 E.6 Action Detection

477 We mainly follow the CVRL [8] testing pipeline, taking our pre-trained R3D-18 as the backbone and  
478 casting a Faster-RCNN [74] on top of it. To fit the time-sequential nature of the input, we extract  
479 region-of-interest (RoI) features using a 3D RoIAlign on the output from the final convolutional  
480 block. These features are then processed through temporal average pooling and spatial max pooling.  
481 The resulting feature is fed into a sigmoid-based classifier for multi-label prediction. We pretrain our  
482 R3D-18 with three different methods (VIC/BYOL/SimCLR) and two different inputs (with/without  
483 derivative). We use an AdamW[75] optimizer with a 0.01 learning rate, then shrink the learning rate  
484 to half after epoch 5. The dropout rate for Faster-RCNN is 0.5. We perform 20 epochs for our six  
485 pre-trained weights and run an evaluation after each epoch. We report the epoch with the highest  
486 mAP. Our clip length is eight frames with an interval of four frames. The finetuning can be conducted  
487 on 2 GPUs, each having at least 48 GB of memory. The testing can be conducted on 1 GPU.

## 488 F Auxiliary Results

### 489 F.1 Silhouette Score

490 Apart from visualization of clustering in the latent space, we also quantify the clustering using  
491 Silhouette Score as illustrated in 7.

Table 7: **Silhouette Score** for Base and ViDiDi with 3, 5, . . . , 101 classes. ViDiDi improves the Score, showing better clustering in the latent space.

Method	Silhouette Score					
	3	5	10	15	20	101
SimCLR	0.136	0.081	0.048	0.034	0.022	-0.026
<b>ViDiDi-SimCLR</b>	<b>0.210</b>	<b>0.132</b>	<b>0.096</b>	<b>0.078</b>	<b>0.058</b>	<b>0.003</b>
BYOL	0.038	-0.086	-0.094	0.091	-0.080	-0.186
<b>ViDiDi-BYOL</b>	<b>0.185</b>	<b>0.230</b>	<b>0.128</b>	<b>0.107</b>	<b>0.070</b>	<b>0.004</b>
VICReg	0.110	0.069	0.044	0.036	0.017	-0.038
<b>ViDiDi-VIC</b>	<b>0.235</b>	<b>0.232</b>	<b>0.150</b>	<b>0.138</b>	<b>0.098</b>	<b>0.014</b>

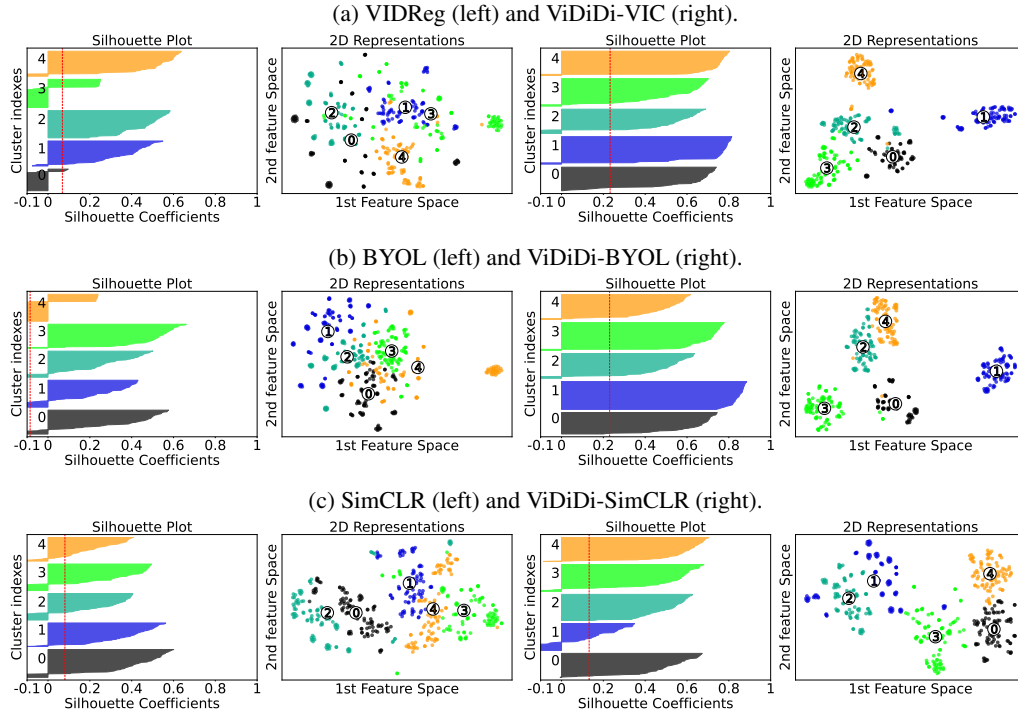


Figure 6: **Silhouette scores and t-SNE plots of top 5 classes in UCF101 train.**

## 492 F.2 Clustering of Latent Space

493 We provide more visualization of the clustering phenomenon for VICReg, BYOL, and SimCLR, with  
 494 or without the ViDiDi framework; on UCF101 train dataset or test dataset; utilizing 5 or 10 classes  
 495 of videos. Here, for each model, we choose the top 5 or 10 classes of videos that are best retrieved  
 496 during the video retrieval experiments. The results are shown in fig. 6, fig. 7, fig. 8, and fig. 9. ViDiDi  
 497 provides consistently better clustering in the latent space for both train data and test data.

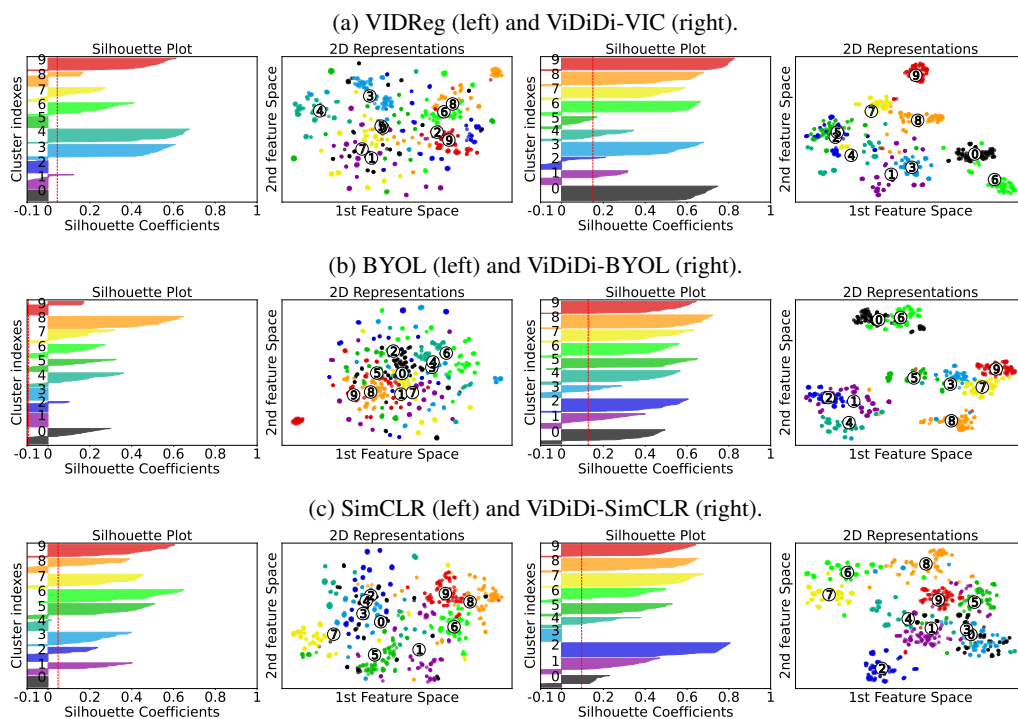


Figure 7: Silhouette scores and t-SNE plots of top 10 classes in UCF101 train.

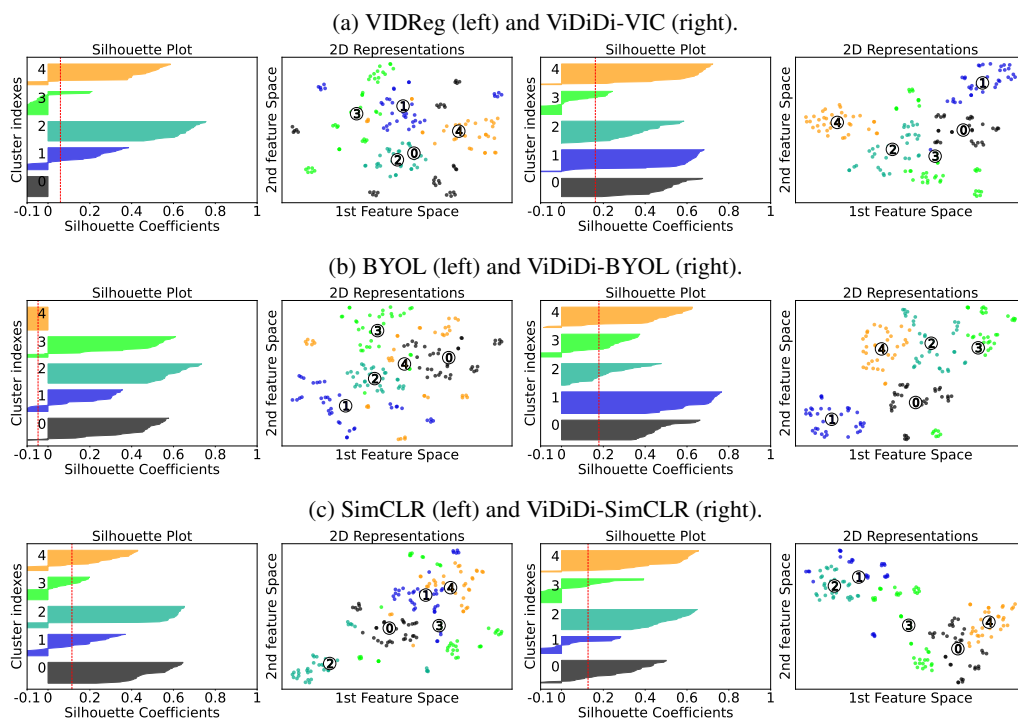


Figure 8: Silhouette scores and t-SNE plots of top 5 classes in UCF101 test.



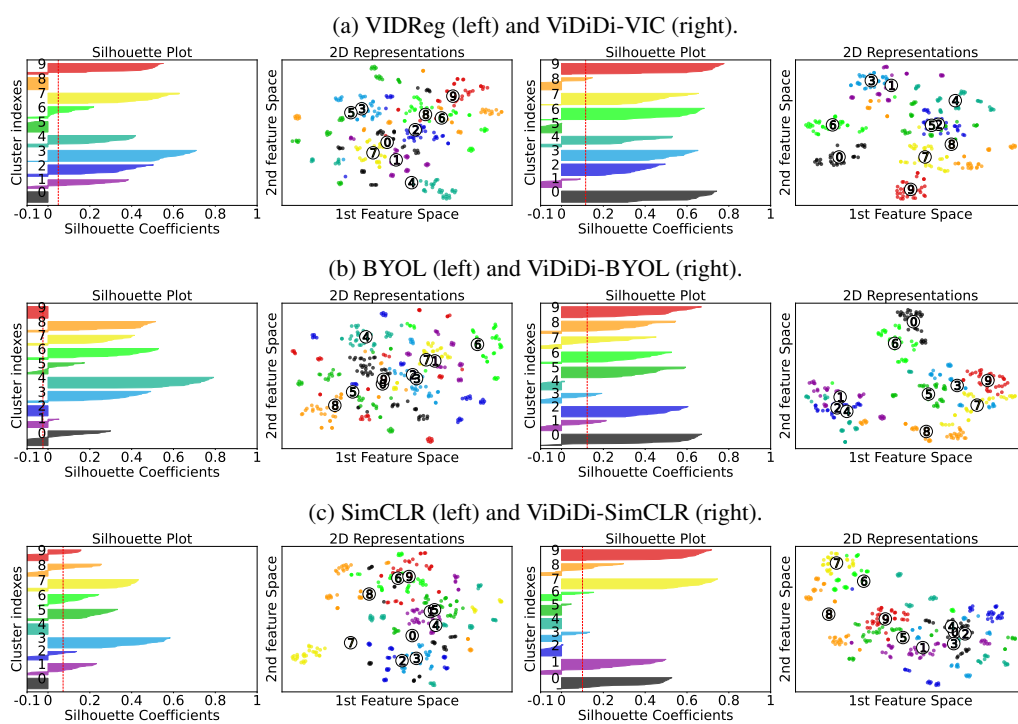


Figure 9: Silhouette scores and t-SNE plots of top 10 classes in UCF101 test.

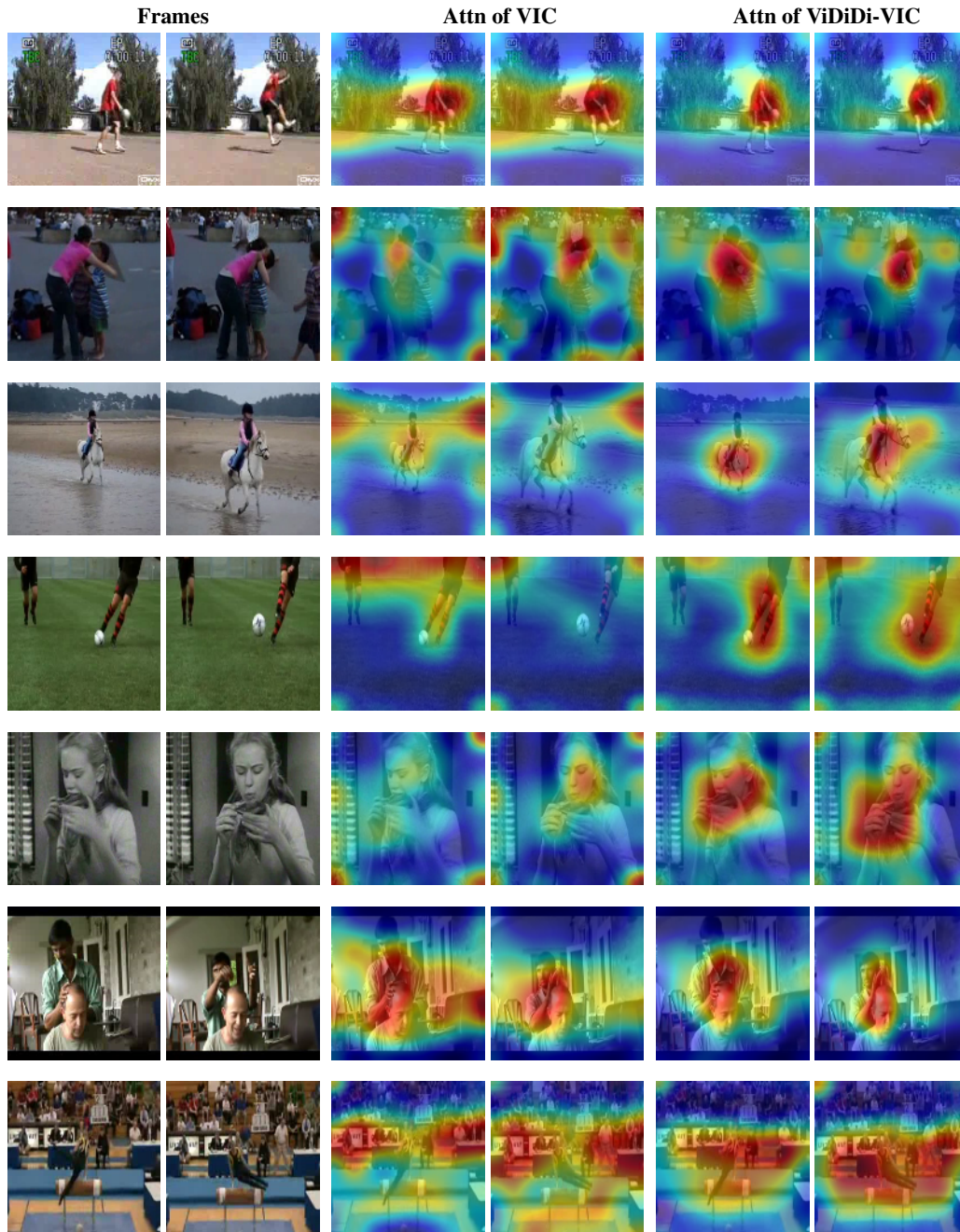


Figure 10: **More spatiotemporal attention for VICReg and ViDiDi-VIC.** Left: Original frames. Middle: Attention from VIC. Right: Attention from ViDiDi-VIC.

498 **F.3 Spatio-temporal Attention**

499 We provide more visualization of the attention for VICReg, BYOL, and SimCLR, with or without  
 500 the ViDiDi framework; on UCF101 dataset or HMDB51 dataset. The results are presented in fig. 10,  
 501 fig. 11, fig. 12, fig. 13, and fig. 14.

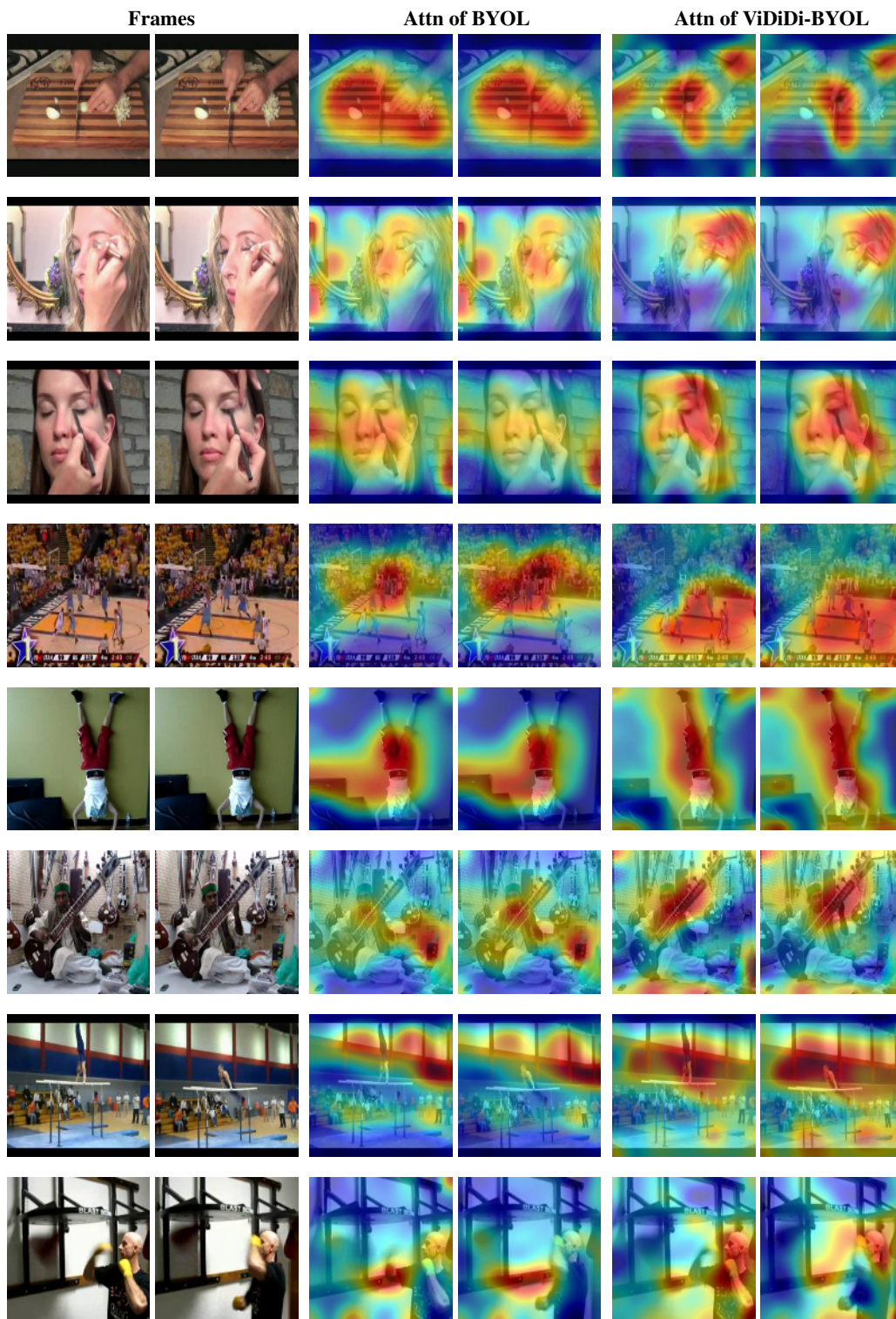


Figure 11: **Spatiotemporal attention on UCF101.** Left: Original frames. Middle: Attention from BYOL. Right: Attention from ViDiDi-BYOL.

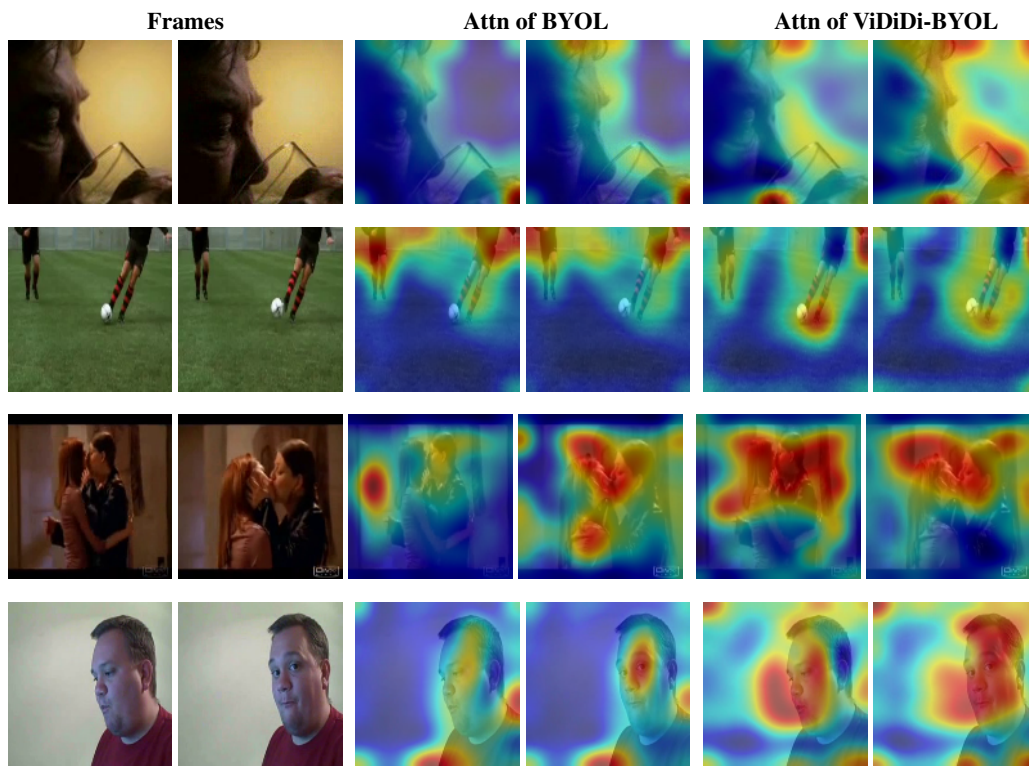


Figure 12: **Spatiotemporal attention on HMDB51**. Left: Original frames. Middle: Attention from BYOL. Right: Attention from ViDiDi-BYOL.

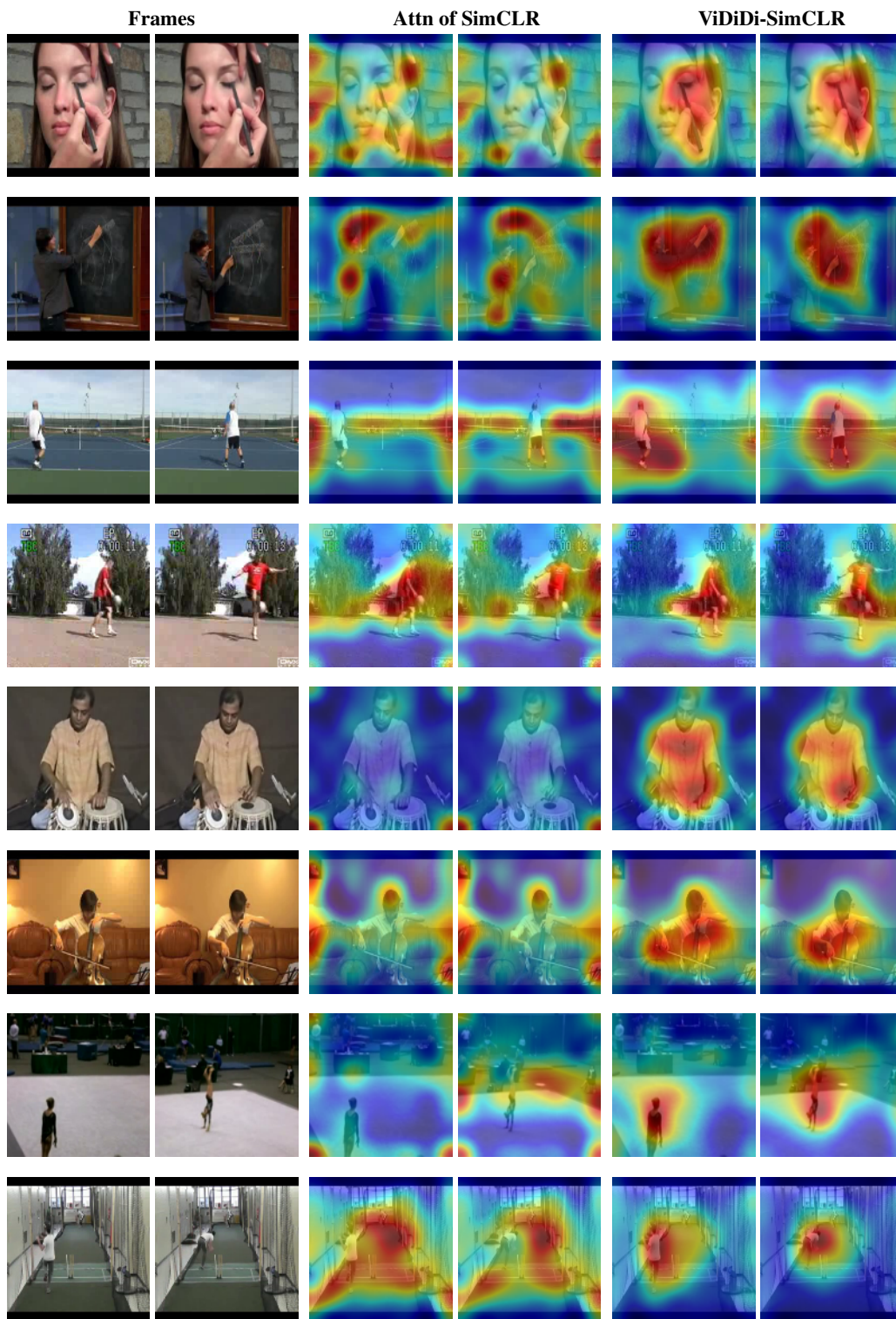


Figure 13: Spatiotemporal attention on UCF101. Left: Original frames. Middle: Attention from SimCLR. Right: Attention from ViDiDi-SimCLR.

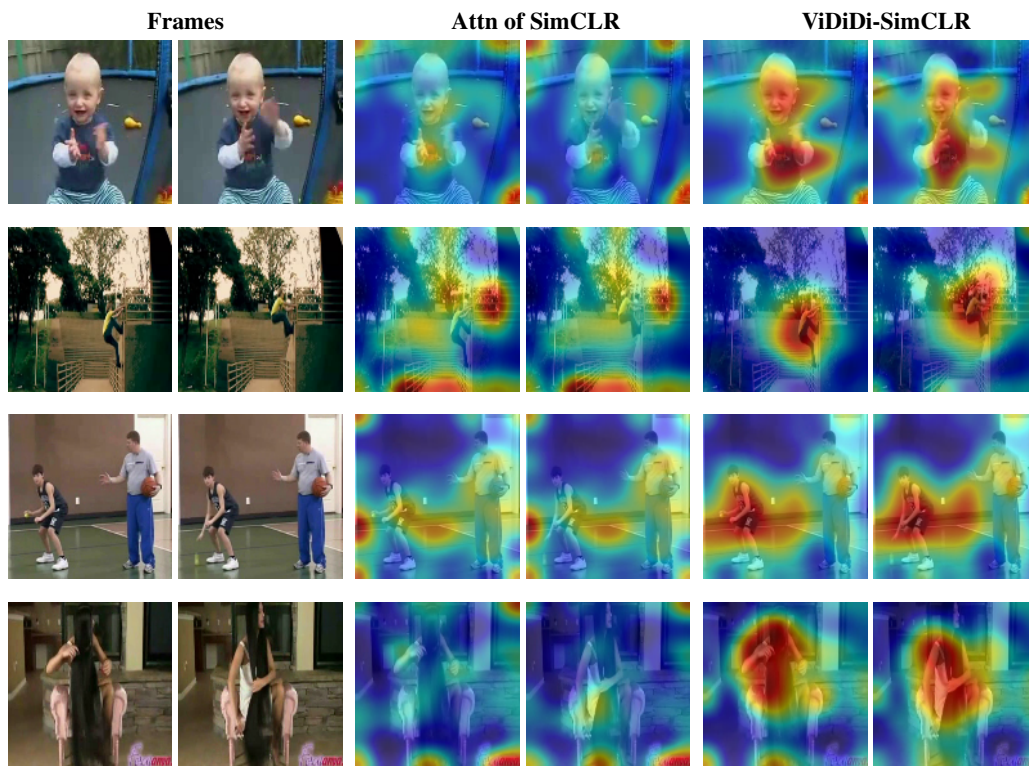


Figure 14: **Spatiotemporal attention on HMDB51**. Left: Original frames. Middle: Attention from SimCLR. Right: Attention from ViDiDi-SimCLR.