# On Temporal Granularity in Self-Supervised Video **Representation Learning**

Rui Qian, Yeqing Li, Liangzhe Yuan, Boqing Gong, Ting Liu, Matthew Brown, Serge Belongie, Ming-Hsuan Yang, Hartwig Adam, Yin Cui



### **Research Motivation**

- Common video representation learning methods mainly focus on temporally persistent learning:
  - The features learned are persistent across all temporal locations of the whole video.
- This objective is not always preferable:
  - The strong performances are only on common video benchmarks of action recognition.
  - The feature learned ignores the changing nature in videos, may be not suitable for all video tasks.

### Contribution

### • We propose the concept of temporal granularity.



#### Given the same short video:

- Event boundary detection task calls for temporally fine-grained features to be aware of the temporal content shifts.
- Video-level recognition task requires the model to robustly predict the target label based on some sampled clips from the video, preferring coarse-grained features.



## Evaluation and benchmarks

- Two representation setting:
  - TeG-PS: a = 0.0, only persistent
  - TeG-FG: a = 0.9, mostly fine-grained
- Benchmarks:
  - Prefer fine-grained feature:
  - Event classification
  - Event boundary detection
  - Prefer coarse-grained feature:
  - Scene-heavy action recognition
  - Interesting to explore:
  - Temporal-heavy action recognition
  - Spatio-temporal action localization on **AVA-Kinetics**

### **Proposed Method**

- Temporal sampling:
  - Long-short sampling of two clips.
- Two embedding spaces:
  - Fine-grained space using dense contrastive learning for temporal discriminative features.
  - Coarse-grained space using global contrastive learning for coarse features.

# Results

method		backbone	acc.	method	external data	finetuning	F1 score	
Supervised		121)	21.9	BMN [50]	IN + THUMOS	×	18.6	
		13D	31.2	SceneDet [12]	-	×	27.5	
	Train from	I3D + NL	30.2	PA [65]	IN	×	39.6	
	scratch [62]	R3D-50 + NL	33.1	BMN-SE [50]	IN + THUMOS	×	49.1	
		SlowFast + NL	32.6	TCN [44]	IN	×	58.8	
nsupervised	CVDI [59]	P2D 50	28.2	PC [65]	IN	1	62.5	
	C 1 m [00]	1630-50	20.0	CVRL [58]	-	1	69.1	
	TeG-PS	R3D-50	28.3	TeG-PS	-	1	69.9	
	TeG-FG	R3D-50	31.1	TeG-FG	-	1	71.4	
Event classification				Event houndary detection				

	backbone	nro train				
method		pre-train				acc.
		dataset	epocn	Irame	FLOPS (G)	
SimCLR [58]	R3D-50	K400	-	-	-	46.8
VINCE [28]	R-50	K400	200	-	-	49.1
ImageNet [58]	R3D-50	IN	-	-	-	53.5
SeCo [83]	R-50	$IN^{\dagger} + K400$	400	-	-	61.9
$\rho SwAV (\rho=2)$ [21]	R3D-50	K400	800	8+8	83.5	63.2
CVRL [58]	R3D-50↓	K400	800	$16{+}16 \rightarrow 8{+}8$	91.2	66.1
$\rho BYOL (\rho=2)$ [21]	R3D-50	K400	800	8+8	83.5	66.2
$\rho$ MoCo ( $\rho$ =2) [21]	R3D-50	K400	800	8+8	83.5	67.4
$\rho$ MoCo ( $\rho$ =2) [21]	R3D-50	K400	200	16+16	167.0	67.6
$\rho BYOL (\rho = 4) [21]$	R3D-50	K400	200	16+16+16+16	334.0	71.4
TeG-FG	R3D-50↓	K400	800	$16+32 \rightarrow 8+16$	136.4	65.0
TeG-PS	R3D-50↓	K400	800	$16+32 \rightarrow 8+16$	136.4	67.8

Scene-heavy action recognition