

# SPEEDNET: LEARNING THE SPEEDINESS IN VIDEOS

Sagie Benaim, Ariel Ephrat, Oran Lang, Inbar Mosseri, William T. Freeman,  
Michael Rubinstein, Michal Irani, Tali Dekel

Google Research, Tel Aviv University, Weizmann Institute

CVPR, 2020

# SPEEDINESS

Slower



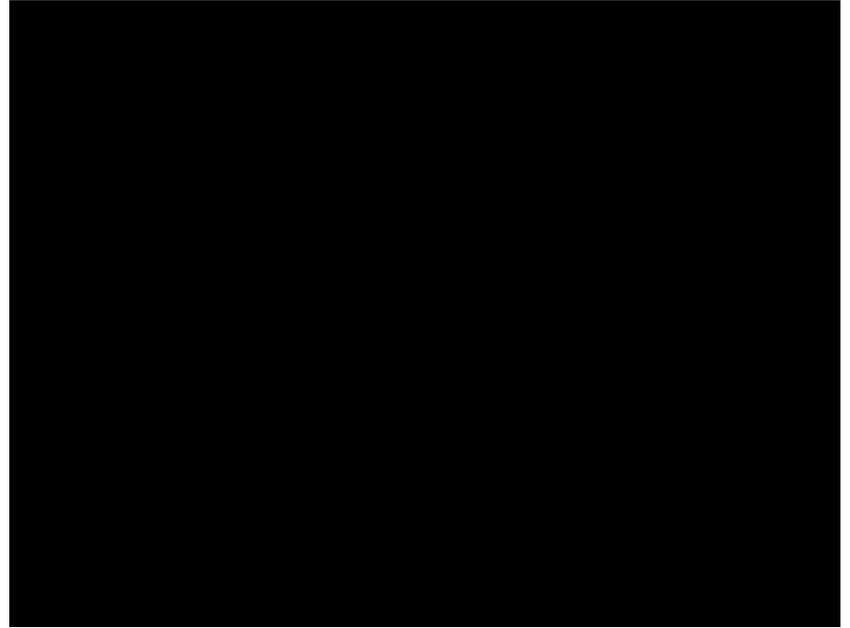
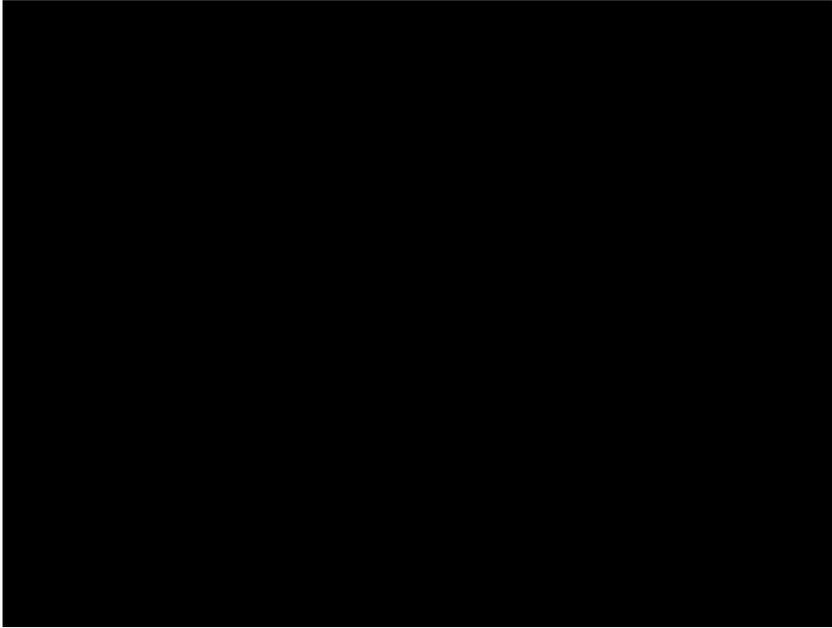
Normal



Faster



# SPEEDINESS

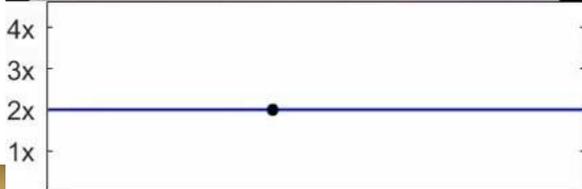


# MOTIVATION

- Video playback speed classification
  - Detecting replays in sports

# MOTIVATION

- Video playback speed classification
  - Detecting replays in sports
- Video time remapping
  - Generate adaptive speedup videos, depending on the speediness score, so that when sped up, their motion looks more natural to the viewer



# MOTIVATION

- Video playback speed classification
  - Detecting replays in sports
- Video time remapping
  - Generate adaptive speedup videos, depending on the speediness score, so that when sped up, their motion looks more natural to the viewer
- Self-supervised learning from videos
  - SpeedNet learns a powerful spacetime representation that can be used for self-supervised action recognition and for video retrieval

# CHALLENGES

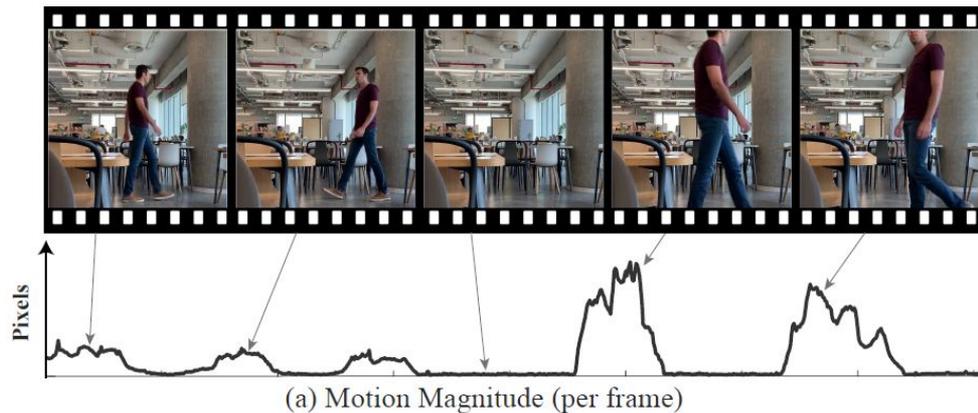
- high-level reasoning is necessary for solving this task

# CHALLENGES

- high-level reasoning is necessary for solving this task
- Avoid tendency to detect easy shortcuts specific to the dataset
  - Rely on artificial, low-level cues, such as compression artifacts

# CHALLENGES

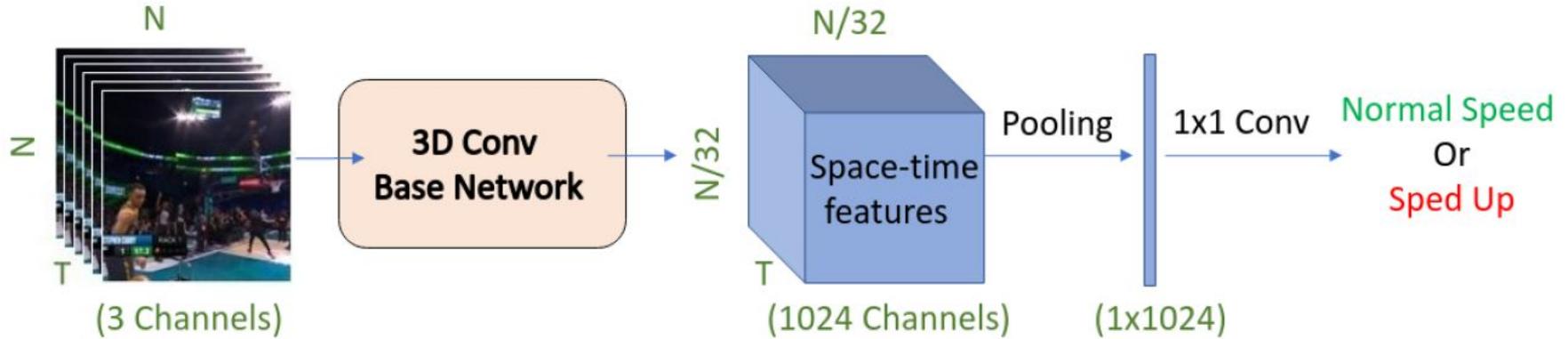
- high-level reasoning is necessary for solving this task
- Avoid tendency to detect easy shortcuts specific to the dataset
  - Rely on artificial, low-level cues, such as compression artifacts
- Trivial case of motion magnitude => speediness (optical flow)
  - Two people walking normally at two different distances from the camera
  - Normal video of hare vs Sped up video of turtle



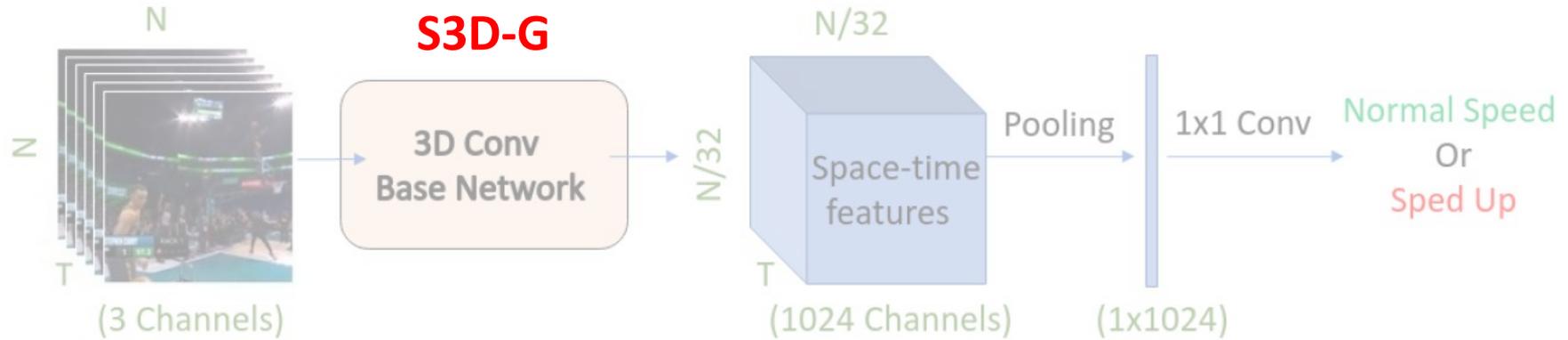
# SPEEDNET

- Classify whether an object in an input video sequence is moving at its normal speed, or faster than the normal speed
  - Given a set of  $L$  frames in an  $L$ -fps video, predict whether those frames depict 1 second of the object's motion (normal), or more than 1 second (sped up) [specifically 2 seconds]
- To determine whether motion in a video is natural or not, a regression objective may be unnecessarily difficult to learn

# SPEEDNET ARCHITECTURE



# SPEEDNET ARCHITECTURE



– most dominant spatially moving object

▪ spatially apply global max pooling

temporally avoid sensitivity to instantaneous “spiky” motions

▪ temporally apply global average pooling.

# SPEEDNET - AVOIDING ARTIFICIAL CUES

## Spatial augmentations

- randomly resize the input video clip to a spatial dimension  $N$ 
  - blurring during resize process can help mitigate potential pixel intensity jitter caused by compression
  - Since the input is of variable size, space-time features correspond to differently sized regions in the unresized input. This forces network not to rely only on size-dependent factors, such as motion magnitude

# SPEEDNET - AVOIDING ARTIFICIAL CUES

Spatial augmentations

Temporal augmentations

- Introduce variability in the time domain
- Normal speed => rate of 1-1.2
- Spedup version => rate of 1.7-2.2

# SPEEDNET - AVOIDING ARTIFICIAL CUES

Spatial augmentations

Temporal augmentations

Same-batch training

- Each batch contains both normal-speed and sped-up versions of each video clip

# EXPERIMENTS ON SPEEDNET

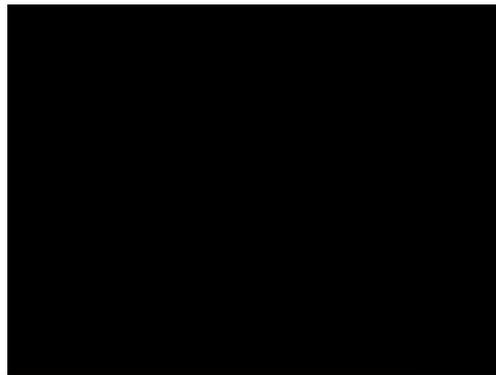
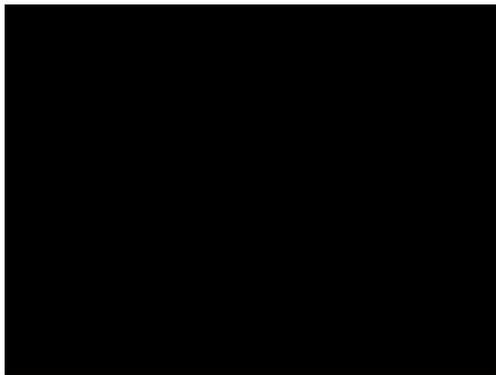
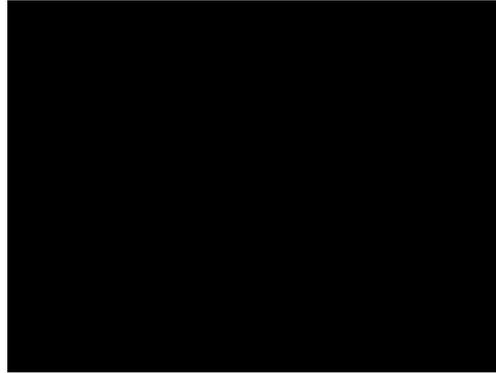
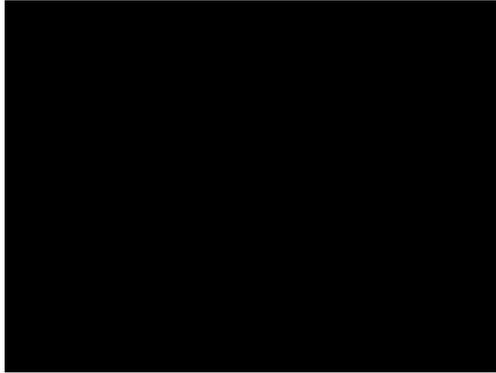
## Train

- **Kinetics train set** (Normal = 1X | Spedup = 2X)

## Test

- **Kinetics test set** (Normal = 1X | Spedup = 2X)
- **Need for Speed dataset (NFS)** (Normal = 10X | Spedup = 20X)

# EXPERIMENTS ON SPEEDNET



## EXPERIMENTS ON SPEEDNET

<b>Batch</b>	<b>Model Type</b>		<b>Accuracy</b>	
	<b>Temporal</b>	<b>Spatial</b>	<b>Kinetics</b>	<b>NFS</b>
Yes	Yes	Yes	75.6%	73.6%
No	Yes	Yes	88.2%	59.3%
No	No	Yes	90.0%	57.7%
No	No	No	96.9%	57.4%
Mean Flow			55.8%	55.0%

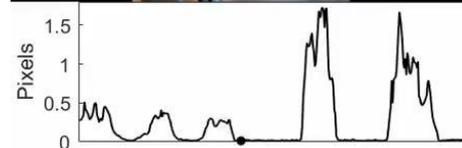
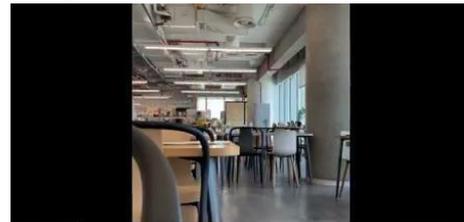
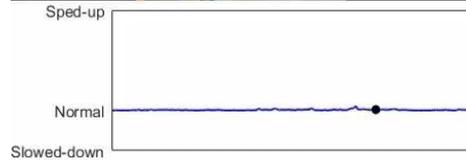
## EXPERIMENTS ON SPEEDNET

Batch	Model Type		Accuracy	
	Temporal	Spatial	Kinetics	NFS
Yes	Yes	Yes	75.6%	73.6%
No	Yes	Yes	88.2%	59.3%
No	No	Yes	90.0%	57.7%
No	No	No	96.9%	57.4%
Mean Flow			55.8%	55.0%

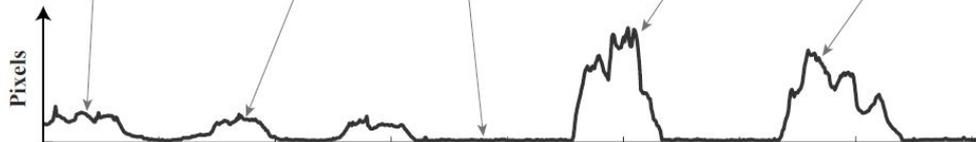
## EXPERIMENTS ON SPEEDNET

Batch	Model Type		Accuracy	
	Temporal	Spatial	Kinetics	NFS
Yes	Yes	Yes	75.6%	73.6%
No	Yes	Yes	88.2%	59.3%
No	No	Yes	90.0%	57.7%
No	No	No	96.9%	57.4%
Mean Flow			55.8%	55.0%

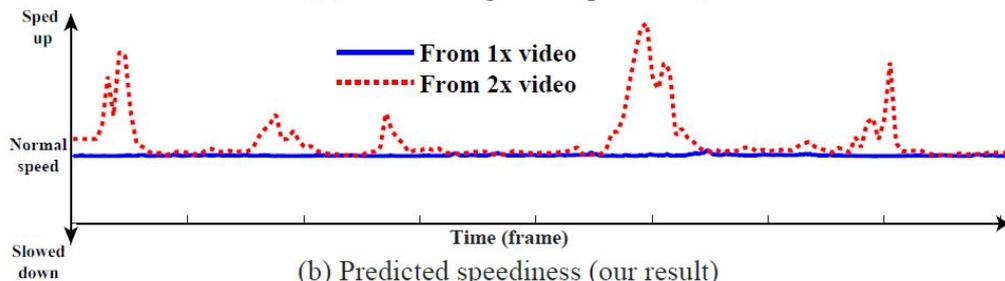
# SPEEDNET - MOTION MAGNITUDE



# SPEEDNET - MOTION MAGNITUDE



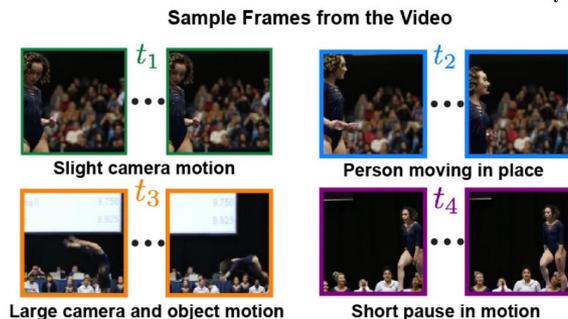
(a) Motion Magnitude (per frame)



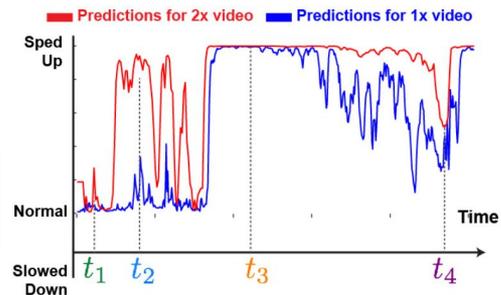
(b) Predicted speediness (our result)

# SPEEDNET - MOTION MAGNITUDE

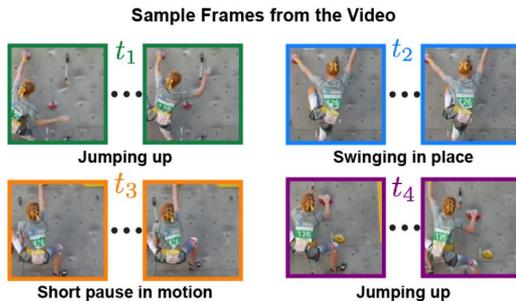
"Gymnast"



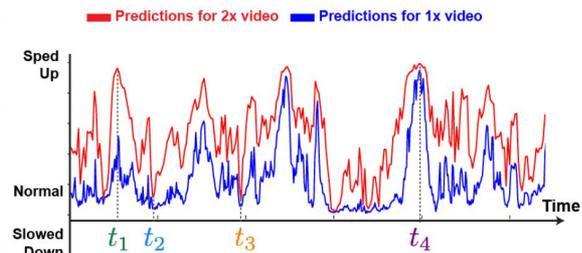
Video sped up probability (probability of x2 speed)



"Rock Climbing"



Video sped up probability (probability of x2 speed)



## ADAPTIVE VIDEO SPEEDUP

- $\{V_i\}$  => Videos of different speeds
  - $i \rightarrow 0$  to  $k$
- Speed of  $V_i = X^i$ 
  - $X = 1.25$

# ADAPTIVE VIDEO SPEEDUP

- $\{V_i\}$
- $\{P_{v_i}\} \Rightarrow$  probability of normal speed
- $\{P_{v_i}^\delta\} \Rightarrow$  threshold probability to get speed label
  - 0  $\Rightarrow$  spedup
  - 1  $\Rightarrow$  normal speed

# ADAPTIVE VIDEO SPEEDUP

- $\{V_i\}$
- $\{P_{v_i}\}$
- $\{V_i(t)\} = \{P_{v_i} \delta\} * X$
- $V(t) = \max(\{V_i(t)\})$  along  $i$ 
  - contains the maximum possible speedup for each timestep that was still classified as not sped-up

# ADAPTIVE VIDEO SPEEDUP

Optimal speedup  $S^*$

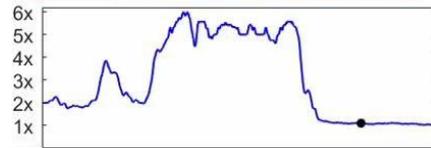
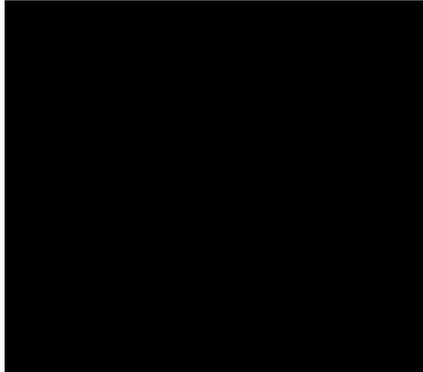
Match overall speedup

$$\arg \min_S E_{\text{speed}}(S, V) + \beta E_{\text{rate}}(S, R_o) + \alpha E_{\text{smooth}}(S')$$

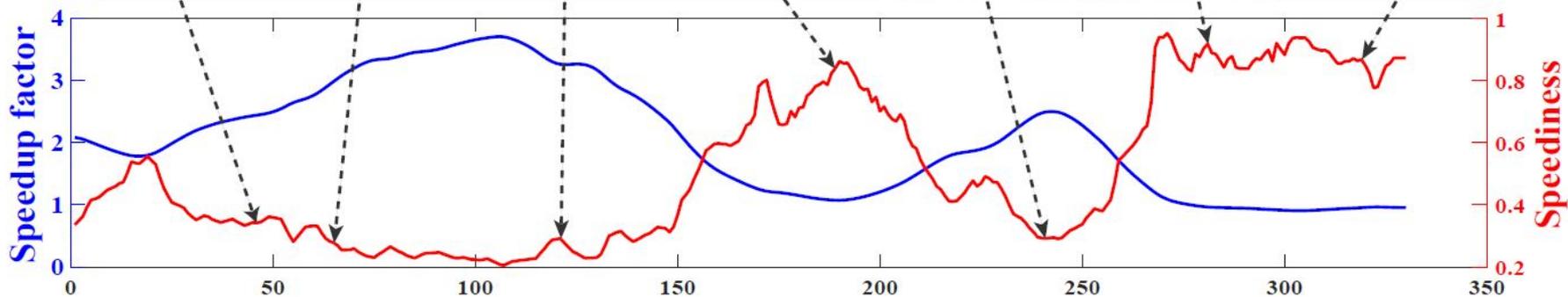
Match speedup  
given by  $V$

smoothness regularizer

# ADAPTIVE VIDEO SPEEDUP



# ADAPTIVE VIDEO SPEEDUP



## EXPERIMENTS ON ADAPTIVE VIDEO SPEEDUP

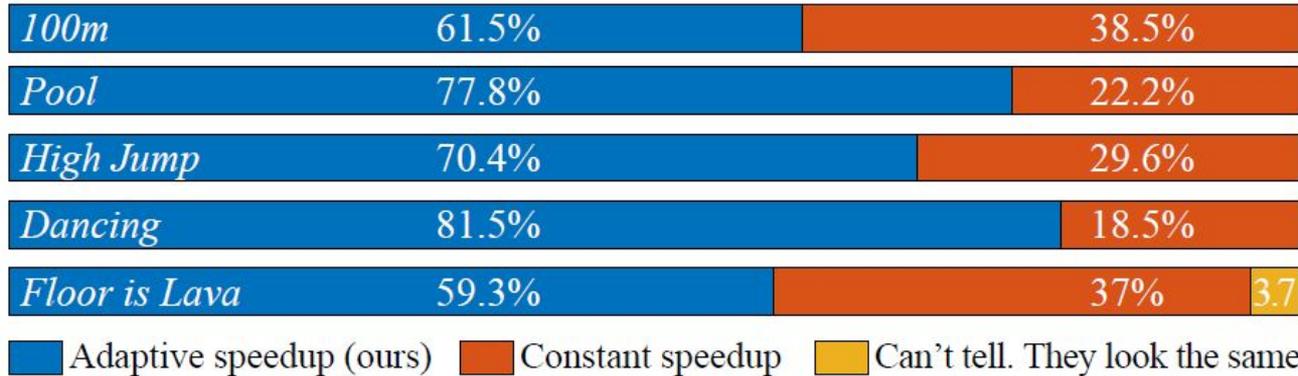
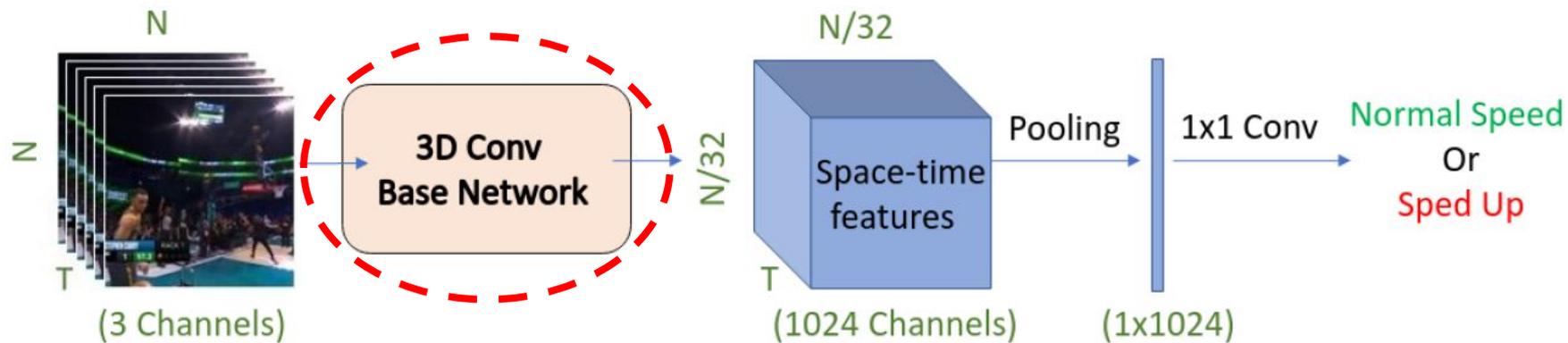


Figure 6. **Adaptive video speedup user study.** We asked 30 participants to compare our adaptive speedup results with constant uniform speedup for 5 videos (without saying which is which), and select the one they liked better. Our adaptive speedup results were consistently (and clearly) preferred over uniform speedup.

# ACTION RECOGNITION



**Pretraining**

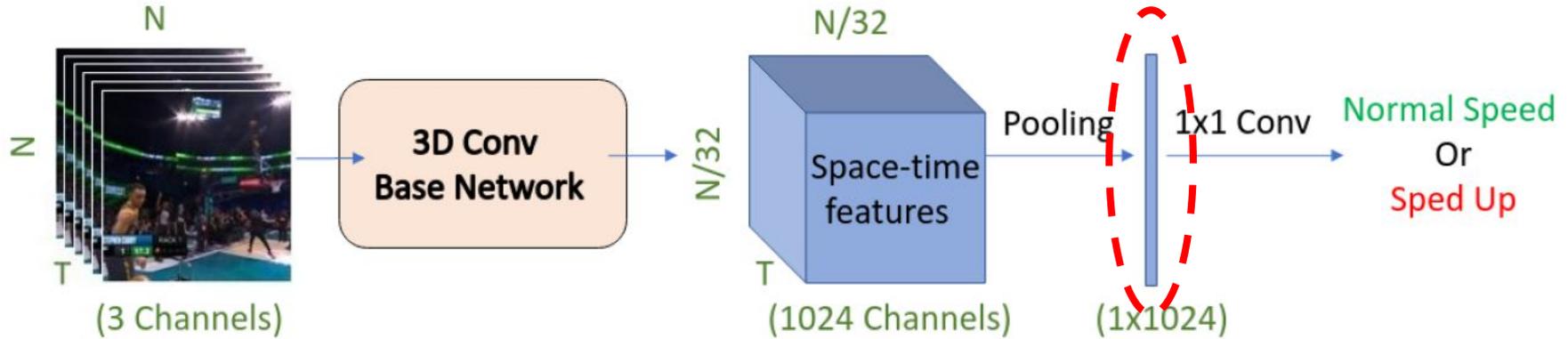
# ACTION RECOGNITION

Method	Initialization	Supervised accuracy	
	Architecture	UCF101	HMDB51
Random init	S3D-G	73.8	46.4
ImageNet inflated	S3D-G	86.6	57.7
Kinetics supervised	S3D-G	96.8	74.5
CubicPuzzle [19]	3D-ResNet18	65.8	33.7
Order [40]	R(2+1)D	72.4	30.9
DPC [13]	3D-ResNet34	75.7	35.7
AoT [38]	T-CAM	79.4	-
<b>SpeedNet (Ours)</b>	<b>S3D-G</b>	<b>81.1</b>	<b>48.8</b>
Random init	I3D	47.9	29.6
SpeedNet (Ours)	I3D	66.7	43.7

# ACTION RECOGNITION

Method	Initialization	Supervised accuracy	
	Architecture	UCF101	HMDB51
Random init	S3D-G	73.8	46.4
ImageNet inflated	S3D-G	86.6	57.7
Kinetics supervised	S3D-G	96.8	74.5
CubicPuzzle [19]	3D-ResNet18	65.8	33.7
Order [40]	R(2+1)D	72.4	30.9
DPC [13]	3D-ResNet34	75.7	35.7
AoT [38]	T-CAM	79.4	-
SpeedNet (Ours)	S3D-G	<b>81.1</b>	<b>48.8</b>
Random init	I3D	47.9	29.6
SpeedNet (Ours)	I3D	66.7	43.7

# NEAREST NEIGHBOR RETRIEVAL



**Video representation**

# NEAREST NEIGHBOR RETRIEVAL



(a) Within a video



(b) Across videos

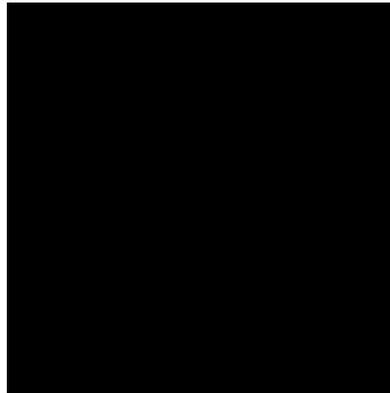
**Not necessarily  
same class**

## NEAREST NEIGHBOR CLASS CORRELATION

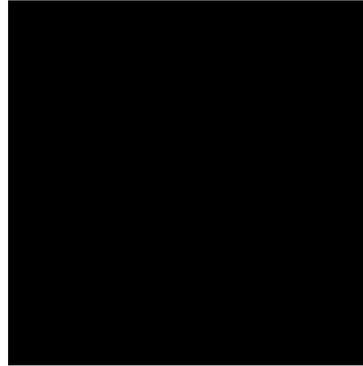
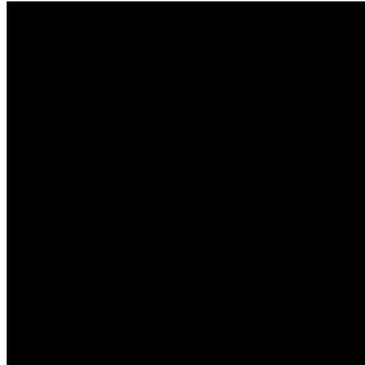
Method	Architecture	1	5	10	20	50
Jigsaw [28]	CFN	19.7	28.5	33.5	40.0	49.4
OPN [22]	OPN	19.9	28.7	34.0	40.6	51.6
Buchler [3]	CaffeNet	25.7	36.2	42.2	49.2	59.5
Order [40]	C3D	12.5	29.0	39.0	50.6	66.9
Order [40]	R(2+1)D	10.7	25.9	35.4	47.3	63.9
Order [40]	R3D	14.1	30.3	40.0	51.1	66.5
Ours	S3D-G	13.0	28.1	37.5	49.5	65.0

Table 3. **Recall-at-topK.** Top-K accuracy for different values of K for UCF101.

# VISUALIZING SALIENT SPACETIME REGIONS



# SPATIALLY-VARYING SPEEDINESS



## STRENGTHS

- Novel idea to detect video speed
- Theory and details behind the method well explained in the paper
- Great experimentation and ablation
- Learns general representations useful for multiple applications
- Great project page with details and examples -  
<https://speednet-cvpr20.github.io/>

## WEAKNESSES

- Could have used backbones other than S3D to compare with the action recognition task baselines
- NFS dataset framerate adjusted so that it's similar to that of Kinetics. Could have tried with original framerate
- Adaptive video speedup user study is human based and subjective

# DISCUSSION