

Procedural Generation of Videos to Train Deep Action Recognition Networks



Introduction

- **Problem:** labeled video data bottleneck
- Cost: manual video annotation is expensive
- Bias: lack of diversity and variations
- **Solution:** synthetic video generation
- Modern computer graphics and game engines
- Realistic, efficient, and scriptable simulation

39,982

1,142.34

2d 07h 31m

5,996,286

1-10s

35

- Contributions
- A procedural generative model of human action videos
- Physically plausible variations of scenes and actions
- Large synthetic dataset with real and procedural actions
- Experimental validation as a data supplement for training deep action recognition models

Procedural Human Action Videos (PHAV)

- Action clips:
- Categories:
- Clips / class:
- Clip duration
- Total duration:
- Frames / modality:





Scene composition

Protagonis





Actions

- **One person synthetic (10):** car hit, crawl, dive floor, flee, hop, leg split, limp, moonwalk, stagger, surrender
- Two people synthetic (4): walking hug, walk hold hands, walk the line, bump into each other

Adding variations

Limited-time physics simulations for human body (ragdoll physics)

Random perturbations

- Muscle weakening
- Interaction with object

Synthetic action generation

- Library of atomic actions, e.g.:
- Left arm movement during walk
- Upper right leg movement during jump

New actions by programmatic composition, e.g.:

- Flee

Limp

Hold hands

Procedural Generation with an Interpretable Parametric Generative Model



César Roberto de Souza¹², Adrien Gaidon³, Yohann Cabon¹, Antonio Manuel López² ¹NAVER LABS Europe, France; ² Centre de Visió per Computador & UAB, Spain, ³ Toyota Research Institute, USA

Generating a Large Synthetic Dataset for Action Recognition

• "A human action video contains a protagonist performing an action in an environment, under particular weather conditions at a specific period of the day. There can be one or more background actors in the scene, as well as one or more supporting characters. The scene is filmed using a camera behaviour."

sub-HMDB (21): brush hair, catch, clap, climb stairs, golf, jump, kick ball, push, pick, pour, pull up, run, shoot ball, shoot bow, shoot gun, sit, stand, swing baseball, throw, walk, wave



 \rightarrow Place a protagonist with all movements from walk except left hand, plus one supporting character with all movements from walk except right hand, then tie puppet hands together \rightarrow Place a protagonist with leg movements from run, raise puppet hands

 \rightarrow Place a **protagonist** with **leg** movements from **walk**, increase physical weight of left leg

Right lower leg Human ragdoll with 15 parts



Cool mixed-source mini-batches

- Virtual-world and real-world frames mixed together in the same minibatch
- Two paths after last feature layer:
- Prediction of **virtual** dataset class labels
- Prediction of **real** dataset class labels
- Different losses for each data path

RGB







Random perturbations

- Animation motion blending
- Ragdoll muscle weakening
- Objects / inv. kinematics
- Annotations
- Extrinsic and intrinsic camera p
- Actor location in camera coordi
- 2D bounding boxes in screen of a screen
- 3D bounding boxes in world cod

Cool Temporal Segment Networks } real data classes ual data classes **Cool** Temporal Segment Networks Multi-task loss • *z* indexes the source dataset of the video $\mathcal{L}(y, \boldsymbol{G}) =$ $\delta_{\{y \in C_z\}} w_z \mathcal{L}_z(y, G)$

$$\mathcal{L}_{z}(y, \mathbf{G}) = -\sum_{i \in C_{z}} y_{i} \left(G_{i} - \log \sum_{j \in C_{z}} \exp G_{j} \right)$$

z∈{real,virtual}

- w_z is a loss weight (e.g. relative
- proportion of *z* in the mini-batch) • C_{τ} is the set of action
- categories for dataset z
- $\delta_{\{y \in C_z\}}$ is the indicator function that returns one when label *y* belongs to C_z , and zero otherwise

Target PHAV UCF-101 UCF-101 UCF-101 UCF-101 HMDB-51 HMDB-51 HMDB-51 HMDB-51

	Data modalities	s and annotations						
ç	Instance Segmentation	Depth Map	Vertical Optical Flow	Horizontal Optical Flow		Method	UCF-101 %mAcc	HMDB-51 %mAcc
					ONE SOURCE	iDT+FV [73] iDT+StackFV [43] iDT+SFV+STP [72] iDT+MIFS [30] VideoDarwin [13]	84.8 - 86.0 89.1 -	57.2 66.8 60.1 65.1 63.7
	A A A A				PLE SOURCES	2S-CNN [53] TDD [74] TDD+iDT [74] C3D+iDT[62] Actions~Trans [76]	88.0 90.3 91.5 90.4 92.0	59.4 63.2 65.9 - 62.0
	_				MULTIN	2S-Fusion [12] Hybrid-iDT [9] TSN-3M [75] VGAN [70]	93.5 92.5 94.2 52.1 94.2	69.2 70.4 69.4 -
					9. C. R. de Sou and Hybrid	za, A. Gaidon, E. Vig, and A. Classification Architectures fo	M. López. Sym r Action Recogn	pathy for the Details: Dense nition. In ECCV, pages 1–27
oarameters inates coordinates oordinates	 Body joint loca Physical prope Detailed inform environment, v 	ations in screen coordinates (perties of body parts (weight, s nation about the generation p weather, day phase, location,		 C. Feichten Action Reco 13. B. Fernando evolution fo Z. Lan, M. L stacking for X. Peng, C. 2014 K. Simonya In NIDS 20 	norer, A. Pinz, and A. Zisserrr gnition. In CVPR, 2016. b, E. Gavves, M. Oramas, A. G r action recognition. In CVPR, .in, X. Li, A. G. Hauptmann, a action recognition. In CVPR, Zou, Y. Qiao, and Q. Peng. A n and A. Zisserman. Two-stre	an. Convolutio Ghodrati, T. Tu 2015. nd B. Raj. Beyc 2015. action recognition am convolution	vtelaars, and K. U. Leuven. ond gaussian pyramid: Multi on with stacked fisher vector al networks for action recog	

http://adas.cvc.uab.es/phav

A P R

July 21-26 2017

Experiments

Model	Spatial	Temporal	Full
TSN	65.9	81.5	82.3
[75]	85.1	89.7	94.0
TSN	84.2	89.3	93.6
TSN-FT	86.1	89.7	94.1
Cool-TSN	86.3	89.9	94.2
[75]	51.0	64.2	68.5
TSN	50.4	61.2	66.6
TSN-FT	51.0	63.0	68.9
Cool-TSN	53.0	63.9	69.5

Procedural generation can act as Strong prior during model initialization: fine-tuning Regularizer during network training: multi-task Surrogate for missing data samples: **fractioning**

Fine-tuning from PHAV to the real world Improvements for all modalities (RGB and flow) Does not require storing virtual training data

Multi-task loss mixing virtual and real data Improvements for all modalities (RGB and flow) Better results for smaller datasets (HMDB-51)

Fractioning real world data sets into smaller sets What performance improvements we could expect if we had only a fraction of the real world training data? Improvements of up to 6.6 percent points on small data

(342,7K) (84,8K) (34.7K) (17.6K) Cool-TSN UCF-101



Conclusion

- Procedural generation is useful for action recognition
- Strong physical priors present in game and physics engines can be leveraged to train deep neural networks
- Procedural generation can be controlled by a generative probabilistic graphical model of human action videos
- Quantitative evidence that virtual samples can act as drop-in complement for small datasets



- Not necessary to generate particular action categories to obtain performance improvements
- New perspectives for video modelling and understanding

References

- Fraiectories , 2016.
- i-skip feature

- 62. D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3D convolutional networks. In CVPR, 2014. usion for Video **70.** C. Vondrick, H. Pirsiavash, and A. Torralba. Generating videos with scene dynamics. In NIPS, 2016.
 - 72. H. Wang, D. Oneata, J. Verbeek, and C. Schmid. A robust and efficient video representation for action recognition. IJCV, pages 1–20, July 2015.
 - 73. H. Wang and C. Schmid. Action recognition with improved trajectories. In ICCV, 2013.
 - 74. L. Wang, Y. Qiao, and X. Tang. Action recognition with trajectory-pooled deep-convolutional descriptors. In CVPR, 2015
 - 75. L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. V. Gool. Temporal Segment Networks: Towards Good Practices for Deep Action Recognition. In ECCV, 2016.
- gnition in videos. **76.** X. Wang, A. Farhadi, and A. Gupta. Actions ~ Transformations. In CVPR, 2015.

Impact of PHAV on training with fractions of real data (RGB + Flow