A Journey Through Video Research at AWS

Joseph Tighe AWS Rekognition, Senior Scientist



Content

- AWS AI Video Research
- Tracking
 - Pose Understanding (CVPR 2020)
 - Tracking (Under submission)
 - Self Supervised Tracking
- Action Understanding
 - Fine-grained (ICCV 2019)
 - Spatio-temporal modeling (Under submission)
 - Zero-shot Learning (CVPR 2020)
 - Unsupervised Training (Under submission)
- Fast Video Training and Inference (On-going)
- What Comes Next



Tracking to Understand Scene Dynamics





Combining detection and tracking for human pose estimation in videos (CVPR 2020 Oral)

Manchen Wang, Joseph Tighe, Davide Modolo



Introduction







Differences and Benefits



missed person detections

people

tracklets



3. Produces more temporally consistent

2. Better predictions on highly entangled

1. Compensates for

Our approach: 3 components

Clip tracking network 1.

Performs both pose estimation and tracking simultaneously on a short video clip





Our approach: 3 components

- 1. Clip tracking network
- 2. Video Tracking Pipeline





Our approach: 3 components

- Clip tracking network 1.
- Video Tracking Pipeline 2.
- Spatial-Temporal refinement 3.

Merges multiple pose hypotheses --> picks the optimal one (spatial-temporal optimization)



Results on validation set

- SOTA on both human pose estimation (mAP) and tracking (MOTA)
- SOTA on both PoseTrack 2017 and 2018 validation
- SOTA against both top-down and bottom-up approaches

PoseTrack 2017 - validation

MOTA Method mAP Bottom-up JointFlow [9] 69.3 59.8 TML++ [13] 61.3 71.5 STAF [23] 72.6 62.7 STEmbedding [15] 71.8 77.0 Detect&Track [11] 55.2 60.6 PoseFlow [31] 58.3 66.5 Top-down FastPose [33] 70.3 63.2 65.4 FlowTrack [30] 76.7 HRNet [26] 77.3 83.8 Our approach 71.6

PoseTrack 2018 - validation

	Method	MOT
D-	STAF [23]	60.
è	TML++ [13]	65.'
	PT_CPN++ [32]	64.
Ŀ	Our approach	68.'

+6.2 MOTA +6.5 mAP





Official challenge (Test set)

- SOTA on both PoseTrack2017 and 2018 benchmarks

Challenge 3: Multi-Person Pose Tracking

No.	Entry	Additional Training Data	total AP	total MOTA
1	DetTrack	+ COCO	74.14	64.09
2	KeyTrack	+ COCO	74.04	61.15
3	PGPT	+ COCO	72.57	60.17
4	CorrTrack	+ COCO	74.21	60.01
5	POINet	+ COCO	72.49	58.41









Multi-object Tracking with Siamese Track-RCNN (Under submission)

Bing Shuai, Andrew Berneshawi, Davide Modolo, Joseph Tighe



Siamese Track-RCNN

Three functional branches with joint training











MOT Challenge 2017

Method	Year	MOTA	\uparrow IDF1 \uparrow	$\mathrm{MT}\uparrow$	$\mathrm{ML}\downarrow$	$\mathrm{FP}\downarrow$	$\mathrm{FN}\downarrow$	$\mathrm{IDsw}\downarrow$
Siamese Track-RCNN	2020	59.6	60.1	23.9 %	33.9 %	15532	210519	2068
DeepMOT-Tracktor [54]	2019	53.7	53.8	19.4%	36.6%	11731	247447	1947
Tracktor++[5]	2019	53.5	52.3	19.5%	36.6%	12201	248047	2072
DeepMOT-SiamRPN [54]	2019	52.1	47.7	16.7%	41.7%	12132	255743	2271
eHAF [44]	2018	51.8	54.7	23.4%	37.9%	33212	248047	1834
FWT [27]	2017	51.3	47.6	21.4%	35.2%	24101	247921	2648
jCC [30]	2018	51.2	54.5	20.9%	37.0%	25937	247822	1802
STRN [53]	2019	50.9	56.5	20.1%	37.0%	27532	246924	2593
MOTDT17 [7]	2018	50.9	52.7	17.5%	35.7%	24069	250768	2474
MHT_DAM [31]	2015	50.7	47.2	20.8%	36.9%	22875	252889	2314

Table 1. Results on MOT17 test set using the provided public detections.



Tracking on JTA





Improving Temporal Correspondence through Selfsupervised Visual Representation Learning (Under submission)

Daniel McKee, Bing Shuai, Davide Modolo, Joseph Tighe, Svetlana Lazebnik



Task

Given labeled starting frame, propagate dense labels through video using encoder *F* Labels are propagated based on computed affinities between feature locations

We train *F* in a fully selfsupervised manner



Predict





Previous Temporal Correspondence Methods

Self-supervised signals: Colorization Temporal cycle consistency Colorizing video frames



Vondrick et al., 2018

All these methods rely only inter-frame signals for training

Localization+colorization





Li et al., 2019

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.

Temporal cycle consistency

Wang et al., 2019



Image-based Representation Learning

Motivation: ImageNet Classifiers sets a strong baseline without training on video, what about self-supervised visual representation learning? Self-supervised methods show similarly strong performance:

Method	Supervision	\mathcal{J}	\mathcal{F}
RotNet	×	44.8	50.7
DeepCluster VGG16	×	48.2	53.2
ImageNet Classifier	1	51.3	56.6

Our goal: incorporate temporal-based objectives with image-based selfsupervised representation learning



Our Method





Results

DAVIS: video object segmentation

Method Supervision \mathcal{F} \mathcal{J} RotNet 44.850.7X TimeCycle [7] 50.046.4X JointTask [4] 57.761.3+3.6 J X +4.8 F 66.1 61.3Ours X ImageNet Classifier 1 51.356.6OSVOS [1] 11 56.663.9DMM-Net [8] 73.3 11 68.1

 \checkmark or $\checkmark \checkmark$ denote image-level or pixel supervision resp.

VIP: human part segmentation

Method	Supervision	mIoU	AP_{vol}^r	
RotNet	×	26.9	11.6	
TimeCycle ~[7]	×	28.9	15.6	
JointTask [4]	×	34.1	17.7	+
Ours	×	38.3	22.2 🚽	+
ImageNet Classifier	1	31.9	15.9	
ATEN [9]	55	37.9	24.1	

JHMDB: human keypoints

Method	Supervision	PCK@.1	PCK@.2
RotNet	×	56.0	76.1
TimeCycle [7]	×	57.3	78.1
JointTask [4]	×	58.6	79.8
Ours	×	59.8	81.3 🧹
ImageNet Classifier	1	57.3	78.5

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.

-4.2 mloU 4.5 AP

+1.2 PCK@.1 +1.5 PCK@.2



Qualitative Results

Input

JointTask

Ours

















Qualitative Results













Semantic Video Understanding













Activity Classification using Spatial-Temporal Finegrained Discriminative Filter Banks

Brais Martinez Alonso, Davide Modolo, Yuanjun Xiong, Joseph Tighe



Activity Classification using Spatial-Temporal Finegrained Discriminative Filter Banks





water skiing









- **Top**: filters that have specialized on the person's leg and the object being ridden (ski or surf)
- Bottom: filters that have specialized on the texture of the ٠ water (wake or wave)



swimming backstroke



swimming breast stroke





- Filters activate on the frame and location of the • swimmer appearing in a canonical pose
- Robust to viewpoints and scenes' changes •





Directional Temporal Modeling For Action Recognition (Under submission)

Xinyu Li, Bing Shuai, Joseph Tighe





Channel Independent Directional Convolution (CIDC)







CIDC Implementation (Directional Mask)







CIDC Network





CIDC Results – Small dataset



al modeling helps



CIDC Results – Large Dataset

Table 2: Result Kinetics-400 datas parison, we only methods that use 1 Model	com et. Fo com ResNe Conv	pariso or fair npare t-50. Top1	on on com- with Top5		Table 3: Result Something-Somethic We only compare wi use ResNet-50 as ba	comp ng V th me ckbon	arisor 2 dat thods 1e and	n on taset. s that l take	
Signal TSN [34] R2D [35] R2D-NL [35] R3D[35] R3D-NL [35] TSM [19] bLVNet [5] Slowfast (4 × 16) [6] Ours (R2D) Ours (R2D) Ours (R3D)	$\begin{array}{c} 2D\\ 2D\\ 2D\\ 3\overline{D}\\ 3D\\ 2D\\ 3D\\ 3D\\ 2D\\ 3D\\ 2D\\ 3D\\ 3D\\ 2D\\ 3\overline{D}\\ 3\overline{D}\\ \end{array}$	70.6 70.2 72.4 73.1 73.6 74.1 73.5 74.1 72.2 72.8 73.6	89.2 88.7 89.8 90.1 90.9 91.2 91.2 91.2 91.1 90.1 90.5 90.9	+2.2%	Model TSN [19, 34] MultiScale TRN [40] 2-stream TRN [40] Fine-grain [21] TSM [19] bLVNet [5] R2D R3D Ours (R2D)	Conv 2D 2D 2D 3D 2D 3D 2D 3D 2D 3D	Top1 30.0 48.8 55.5 53.4 63.4 61.7 35.5 48.6 40.2	Top5 60.5 77.6 83.1 81.1 88.5 88.1 65.4 78.2 68.6	× 17. 000 × 1. 000 × 1. 000
Ours (R3D-NL)	3D	74.1	91.4			50	00.0	00.1 <	



CIDC Visualization

Clapping (camera moves to right) Peeling Potato (no camera motion)







Spatial Temporal Separated Network (Under submission)

Xinyu Li*, Chunhui Liu*, Hao Chen, Yi Zhu, Joseph Tighe * equally contributed



Spatial Temporal Separated Network





Spatial Temporal Separated Network: Temporal Modeling

- Spatial temporal information exchange by concat along channel
- Temporal to spatial feature aggregation by apply the same attention matrix





Spatial Temporal Separated Network

Version	Spatial Branch	GFLOPs	Top1
R2D-50		64	69.9
R2D-50 NL			72.9
I3D-50		153	74.0
13D-50 NL		282	75.2
I3D-101		167	75.1
I3D-101-NL		359	76.0
IR-CSN-101		73	76.2
Slowfast-50 4x16		36	75.3
Slowfast-50 8x8		66	76.6
Slowfast-101 8x8		106	77.2
STS	R2D-50	49	74.1
STS	R2D-101	64	75.3
STS	I3D-50	34	75.5
STS	I3D-101	102	76.6
STS	I3D-152	192	78.2
STS	X3D-M*	5.7	73.1



Temporal Affinity Map

Raw frame





Visualization – More examples









Visualization – Compare with I3D





Rethinking Zero-shot Video Classification: End-to-end Training for Realistic Applications

Biagio Brattoli, Joseph Tighe, Fedor Zhdanov, Pietro Perona, Krzysztof Chalupka



Task





Related Work



- Complex models
- Hard to reproduce
- Poor performance
- Zero-shot Learning paradigm not enforced in the pre-trained network



Our Model

Out-of-the-box 3D CNN

Out-of-the-box MSE loss



Protocol: Trained once on Kinetics and tested on multiple datasets (UCF101, HMDB, Activity)

- Simple
- Easy to reproduce
- Good performance
- ZSL paradigm enforced

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Co



Enforcing ZSL paradigm

Zero-shot learning paradigm: Training and testing classes should **NOT** overlap

Previous methods: The ZSL paradigm is not enforced in the pre-trained CNN.

Example: "kayaking" (UCF101) and "canoeing and kayaking" (Kinetics) are overlapping even though they have different names



Comparison: strict ZSL protocol

We outperform state-of-the-art while using strict ZSL protocol.

Dataset	VisualFeat	UCF	HMDB	Activity
URL [64]	ResNet200	42.5	51.8	Ξ.
DataAug [60]	-	18.3	19.7	-
InfDem [39]	I3D	17.8	21.3	×
Bidirectional [55]	IDT	21.4	18.9	÷
FairZSL [40]	-	-	23.1	-
TARN [4]	C3D	19	19.5	5
Action2Vec [18]	C3D	22.1	23.5	≂.,
Ours(605classes)	C3D	41.5	25.0	24.8
Ours(664classes)	C3D	43.8	24.7	20
Ours(605classes)	R(2+1)D_18	44.1	29.8	26.6
Ours(664classes)	R(2+1)D_18	48	32.7	<u> </u>



Weakness/Strenghts study

Better more training samples or classes?





Conclusion

Contributions:

- First end-to-end model for ZSL video classification
- Strong and simple baseline for future work
- Cheap method for increasing training classes for a more robust model.
- Deep study of ZSL weaknesses and strengths



Bidirectional GANs for Unsupervised Video Representation Learning

Tom F. H. Runia, Andrew Berneshawi, Rahul Rama Varior, Davide Modolo, Joseph Tighe



BigBiGAN like Setup For Videos



© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.

0/1



BigBiGAN like Setup For Videos





Embedding Performance





Generator Examples





Cycle Examples























Cycle Examples











Training Evolution





Fast Training and Inference

Chunhui Liu, Xinyu Li, Joseph Tighe



Video Model Fast Training Effort

Cross Instance Distributed Data Parallel (DDP) for Videos

- Easy deployment on multiple AWS instances
- Support both video and image model

Multi Grid Training¹

• Speedup training by rescaling input size



1 A Multigrid Method for Efficiently Training Video Models Chao-Yuan Wu, Ross Girshick, Kaiming He, **Christoph Feichtenhofer**, Philipp Krähenbühl CVPR 2020

© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.



Iterations

Cross Instance DDP Training for Videos

Distributed Data Parallel (DDP) make it possible to use multiple instances on AWS for model training.

Each GPU has its own forward-backward process. Then, DDP uses collective functions to synchronize gradients and buffers among GPUs.

The tricky part is to fully balance GPU utilization for running data flow and CPU utilization for loading data. Decord make this possible.



© 2020, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark.

Instance n

GPU rank 8n

GPU rank 8n+1

GPU rank 8n+2

GPU rank 8n+3

...

GPU rank 8n+7



Multi Grid Training

Down-sample input data so that we can scale-up batch size, and thus speedup training process. Short-cycle Scaling Strategy happens for each iteration.

Long-cycle Scaling Strategy happens in each learning rate stage.







Video Model Fast Training: A result

Distributed Data Parallel (DDP)

- Test on Kinetics 400 dataset, with 8x p3.16xlarge instances (64x Tesla V100). \bullet
- 4 epoch/hour, 12 hour to fully train an I3D model. •
- Multi Grid Training + DDP
 - 3 times speedup using short-cycle only, •
 - 4 times speedup using long-short cycle. •
 - **3 hour** to fully train an I3D model. •

Important Learnings

- Short-cycle is more stable than long-short cycle. •
- Long-short cycle becomes unstable if temporal down-sampling is too aggressive. \bullet
- In Multi Grid training, we found that training is more stable without using • learning rate - batch scaling rule.
- Single grid finetuning can help with ~0.5% accuracy. ullet





Video Model Fast Inference Effort

Current Approaches:

- Efficient Neural Architecture Search (NAS) •
- Model Compression and Distillation •
- Model on Compressed Data •

Our effort:

Single Pass Feature Selection Model for Video Fast Inference (Ongoing) •





Motivation: From Multi Crop to One Pass

- Most state-of-the-art methods use short clips at training time and thus at test require dense • oversampling to achieve high performance numbers.
- We want to avoid dense sampling during inference time and maintains this performance \bullet



Results on R50 I3D

Method	Top-1 Acc	Top-5 Acc	FLOPS
64x30 I3D baseline	73.8	91.0	998G
256 frame randomly choose one 64 clip	62.3	78.38	33G
256 frame forward	69.2	88.2	133G
256 frame Feature Selection	71.83	89.6	154G





Future of Video Research

Datasets •







Future of Video Research

- Datasets
- Tasks











Thanks





























Links to Work Presented

This video includes descriptions of the following work:

Combining detection and tracking for human pose estimation in videos: <u>https://assets.amazon.science/9c/68/ce3ec91b41c1b6a20ee1e793709d/scipub-1326.pdf</u>

Multi-Object Tracking with Siamese Track-RCNN: <u>https://arxiv.org/pdf/2004.07786</u>

Action Recognition With Spatial-Temporal Discriminative Filter Banks: <u>http://openaccess.thecvf.com/content_ICCV_2019/papers/Martinez_Action_Recognition_With_Spatial-</u> <u>Temporal_Discriminative_Filter_Banks_ICCV_2019_paper.pdf</u>

Rethinking Zero-shot Video Classification: End-to-end Training for Realistic Applications: <u>https://assets.amazon.science/b3/28/0702b5ec441aaadcb79040b58128/scipub-1328.pdf</u>

