UNIVERSITÄT BONN

Juergen Gall

An Introduction to Temporal Action Segmentation From Fully Supervised Learning to Weakly Supervised Learning

Action Recognition



• Large annotated datasets



- UCF101 (98.2%), HMDB (82.5%), Kinetics-400 (82.8%), Epic-Kitchens (36.7%)
- http://actionrecognition.net
- Continuous data streams



Action Segmentation vs. Action Detection

• Action Detection (THUMOS, ActivityNet)



Ground Truth:

Action Segmentation (Breakfast, 50 Salads, GTEA)







Action Segmentation vs. Action Detection

Action Detection (Object Detection)



Ground Truth:



Action Segmentation (Semantic Segmentation)





Action Segmentation vs. Action Detection

• Action Detection (THUMOS, ActivityNet)



Ground Truth:

• Action Segmentation (Breakfast, 50Salads, GTEA)





Why Action Segmentation?







Datasets

Breakfast

https://serre-lab.clps.brown.edu/resource/breakfastactions-dataset/

• 50 Salads

https://cvip.computing.dundee.ac.uk/datasets/foodpre paration/50salads/

- GTEA <u>http://cbs.ic.gatech.edu/fpv/#gtea</u>
- COIN <u>https://coin-dataset.github.io/</u>

Let's build a baseline...





Hidden Markov Model





[S. Prince. Computer Vision: Models, Learning, and Inference. Cambridge University Press]

Juergen Gall – Institute of Computer Science III – Computer Vision Group

Simon J.D. Prince

Inference



HMM:

$$Pr(\mathbf{x}_{1...N}, w_{1...N}) = \left(\prod_{n=1}^{N} Pr(\mathbf{x}_n | w_n)\right) \left(\prod_{n=2}^{N} Pr(w_n | w_{n-1})\right)$$

MAP inference:

$$\hat{w}_{1...N} = \operatorname*{argmax}_{w_{1...N}} [Pr(\mathbf{x}_{1...N}, w_{1...N})]$$
$$= \operatorname*{argmin}_{w_{1...N}} [-\log [Pr(\mathbf{x}_{1...N}, w_{1...N})]]$$

Substituting:

$$\hat{w}_{1...N} = \underset{w_{1...N}}{\operatorname{argmin}} \left[-\sum_{n=1}^{N} \log \left[Pr(\mathbf{x}_n | w_n) \right] - \sum_{n=2}^{N} \log \left[Pr(w_n | w_{n-1}) \right] \right]$$

Global minimum by dynamic programming

[S. Prince. Computer Vision: Models, Learning, and Inference. Cambridge University Press] simon J

Features: Dense Trajectories



- Dense sampling of features
- Feature tracking

Dense sampling in each spatial scale



[H. Wang et al. Dense Trajectories and Motion Boundary Descriptors for Action Recognition. International Journal of Computer Vision 2013]



KLT

Dense trajectories

Hidden Markov Model



Hidden Markov Model (HMM) for each activity







• HMM + GMM (IDT)







[H. Kuehne et al. An end-to-end generative framework for video segmentation and recognition. WACV 2016]





• HMM + GMM (IDT)





[H. Kuehne et al. An end-to-end generative framework for video segmentation and recognition. WACV 2016]





• HMM + GMM (IDT)



[H. Kuehne et al. An end-to-end generative framework for video segmentation and recognition. WACV 2016]



Grammar

 Transitions between activity HMMs are modeled by context free grammar



- SIL: start and end points
- Transition probability is 1 if connection exists otherwise 0
- [H. Kuehne et al. An end-to-end generative framework for video segmentation and recognition. WACV 2016]





Breakfast dataset (~65 hours)

Method	Frame-wise Accuracy (%)
Kuehne et al. 2016 (HMM+GMM)	56.3

[H. Kuehne et al. An end-to-end generative framework for video segmentation and recognition. WACV 2016]



Hybrid RNN-HMM

HMM + RNN with Gated Recurrent Units (GRU)



Gated Recurrent Units (GRU)



 Similar to LSTM, but it does not need an additional memory cell



 [J. Chung et al. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. NIPS Workshop 2014]
[K. Cho et al. On the Properties of Neural Machine Translation: Encoder-Decoder Approaches. Workshop SSST 2014]

Hybrid RNN-HMM



Breakfast dataset (~65 hours)

Method	Frame-wise Accuracy (%)
Kuehne et al. 2016 (HMM+GMM)	56.3
Richard et al. 2017 (HMM+RNN)	60.6
Kuehne et al. 2020 (HMM+RNN)	61.3

[A. Richard et al. Weakly Supervised Action Learning with RNN based Fine-to-Coarse Modeling. CVPR 2017] [H. Kuehne et al. A Hybrid RNN-HMM Approach for Weakly Supervised Temporal Action Segmentation. PAMI 2020]

Temporal Convolutional Neural Network UNIVERSITÄT BONN



[C. Lea et al. Temporal Convolutional Networks for Action Segmentation and Detection. CVPR 2017] **Temporal Convolutional Network**



Breakfast dataset (~65 hours)

Method	Frame-wise Accuracy (%)
Lea et al. 2017 (ED-TCN)*	43.3
Kuehne et al. 2016 (HMM+GMM)	56.3
Richard et al. 2017 (HMM+RNN)	60.6
Kuehne et al. 2020 (HMM+RNN)	61.3

*[L. Ding and C. Xu. Weakly-supervised action segmentation with iterative soft boundary assignment. CVPR 2018]

[C. Lea et al. Temporal Convolutional Networks for Action Segmentation and Detection. CVPR 2017] Temporal Convolutional Neural Network

UNIVERSITÄT BONN

Dilated convolutions for audio

Output	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	
Hidden Layer	\bigcirc	0	0	\bigcirc	0	\bigcirc	0	0	\bigcirc	0	\bigcirc	0	\bigcirc	\bigcirc	0	
Hidden Layer	\bigcirc	0	\bigcirc	\bigcirc	\bigcirc	\bigcirc	0	\circ	\bigcirc	0	\bigcirc	\bigcirc	\bigcirc	\bigcirc	0	
Hidden Layer	\bigcirc	0	0	0	0	\bigcirc	0	0	0	0	0	\bigcirc	\bigcirc	\bigcirc	0	
Input	•	0	•	0	•	•	0	0	0	0	0	•	0	•	0	С

van den Oord et al. WaveNet: A Generative Model for Raw Audio. SSW 2016]

Temporal Convolutional Neural Network

Dilated convolutions capture long temporal receptive field

Causal convolutions: Input for t depends only on previous observations



[C. Lea et al. Temporal Convolutional Networks for Action Segmentation and Detection. CVPR 2017]

UNIVERSITÄT BONN

Temporal Convolutional Network



• 50 Salads

Method	Frame- wise Accuracy (%)
Lea et al. 2017 (ED-TCN)	64.7
Lea et al. 2017 (Dilated TCN)	59.3

[C. Lea et al. Temporal Convolutional Networks for Action Segmentation and Detection. CVPR 2017]

Temporal Convolutional Network



• 50 Salads

Method	1 – Norm. Edit Distance (%)	Frame- wise Accuracy (%)
Lea et al. 2017 (ED-TCN)	59.8	64.7
Lea et al. 2017 (Dilated TCN)	43.1	59.3

• Edit distance (sensitive to oversegmentation):





UNIVERSITÄT BONN

[Y. Abu Farha and J. Gall. MS-TCN: Multi-Stage Temporal Convolutional Network for Action Segmentation. CVPR 2019]

03.08.2020



[Y. Abu Farha and J. Gall. MS-TCN: Multi-Stage Temporal Convolutional Network for Action Segmentation. CVPR 2019]

UNIVERSITÄT BONN



[Y. Abu Farha and J. Gall. MS-TCN: Multi-Stage Temporal Convolutional Network for Action Segmentation. CVPR 2019]

UNIVERSITÄT BONN

Over-segmentation



Frame-wise classification loss:

$$\mathcal{L}_{cls} = \frac{1}{T} \sum_{t} -log(y_{t,c})$$

 Additional loss is required to avoid oversegmentation:



[Y. Abu Farha and J. Gall. MS-TCN: Multi-Stage Temporal Convolutional Network for Action Segmentation. CVPR 2019]



Loss

- Frame-wise classification loss *L*_{cls}
- Additional loss is required to avoid oversegmentation:

$$\mathcal{L}_{T-MSE} = \frac{1}{TC} \sum_{t,c} \tilde{\Delta}_{t,c}^2$$
$$\tilde{\Delta}_{t,c} = \begin{cases} \Delta_{t,c} & : \Delta_{t,c} \leq \tau \\ \tau & : otherwise \end{cases}$$

$$\Delta_{t,c} = |\log y_{t,c} - \log y_{t-1,c}|$$

Loss functions of all stages s:

$$\mathcal{L} = \sum \mathcal{L}_s \qquad \mathcal{L}_s = \mathcal{L}_{cls} + \lambda \mathcal{L}_{T-MSE}$$

[Y. Abu Farha and J. Gall. MS-TCN: Multi-Stage Temporal Convolutional Network for Action Segmentation. CVPR 2019]







Additional loss is required to avoid oversegmentation





Impact of stages

Impact of stages (50 Salads)

	Edit	Acc
SS-TCN	20.5	78.2
MS-TCN (2 stages)	47.9	79.8
MS-TCN (3 stages)	64.0	78.6
MS-TCN (4 stages)	67.9	80.7
MS-TCN (5 stages)	69.2	79.5

	Edit	Acc
SS-TCN (48 layers)	40.7	78.0
MS-TCN	67.9	80.7

[Y. Abu Farha and J. Gall. MS-TCN: Multi-Stage Temporal Convolutional Network for Action Segmentation. CVPR 2019]



Breakfast dataset (~65 hours)

Method	Frame-wise Accuracy (%)
Lea et al. 2017 (ED-TCN)*	43.3
Kuehne et al. 2016 (HMM+GMM)	56.3
Richard et al. 2017 (HMM+RNN)	60.6
Kuehne et al. 2020 (HMM+RNN)	61.3
MS-TCN (TCN)	65.1
MS-TCN (TCN+I3D)	66.3

[J. Carreira and A. Zisserman. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. CVPR 2017] [Y. Abu Farha and J. Gall. MS-TCN: Multi-Stage Temporal Convolutional Network

for Action Segmentation. CVPR 2019]

Temporal Action Segmentation





03.08.2020

MS-TCN







[Y. Abu Farha and J. Gall. MS-TCN: Multi-Stage Temporal Convolutional Network for Action Segmentation CVPR 2019]

MS-TCN++





MS-TCN++



Breakfast dataset

Method	Frame-wise Accuracy (%)
Lea et al. 2017 (TCN)*	43.3
Kuehne et al. 2016 (HMM+GMM)	56.3
Richard et al. 2017 (HMM+RNN)	60.6
Kuehne et al. 2020 (HMM+RNN)	61.3
MS-TCN (TCN)	65.1
MS-TCN (TCN+I3D)	66.3
MS-TCN++ (TCN+I3D)	67.6

[S. Li et al. MS-TCN++: Multi-Stage Temporal Convolutional Network for Action Segmentation. arXiv]

MS-TCN++ vs. MS-TCN







Training video



• Fully supervised:



• Weakly supervised (transcripts) $A \rightarrow C \rightarrow F \rightarrow D \rightarrow A \rightarrow E \rightarrow H$

> [A. Richard et al. Weakly Supervised Action Learning with RNN based Fine-to-Coarse Modeling. CVPR 2017]

Recall: Hybrid RNN-HMM



HMM + RNN with Gated Recurrent Units (GRU)



















The transcripts define the order of activities:



Action transcript: action_1 action_2 action_3









The transcripts define the order of activities:



Action transcript: action_1 action_2 action_3









The transcripts define the order of activities:



Action transcript: action_1 action_2 action_3





(Initialization)

Action transcript:

action_1 action_2 action_3

linear segmentation

[A. Richard et al. Weakly Supervised Action Learning with RNN based Fine-to-Coarse Modeling. CVPR 2017]

03.08.2020



Action transcript:

action_1 action_2 action_3



[A. Richard et al. Weakly Supervised Action Learning with RNN based Fine-to-Coarse Modeling. CVPR 2017]





[A. Richard et al. Weakly Supervised Action Learning with RNN based Fine-to-Coarse Modeling. CVPR 2017]

03.08.2020





Action transcript:

[A. Richard et al. Weakly Supervised Action Learning with RNN based Fine-to-Coarse Modeling. CVPR 2017]

03.08.2020





Action transcript:

[A. Richard et al. Weakly Supervised Action Learning with RNN based Fine-to-Coarse Modeling. CVPR 2017]





	Breakfast frame accuracy (%)
pseudo-GT (HMM+RNN) Richard et al. 2017	33.3
pseudo-GT (HMM+RNN) Kuehne et al. 2020	36.7
<i>Fully supervised (HMM+RNN)</i> Kuehne et al. 2020	61.3

Disadvantage: Offline and sensitive to initialization

[H. Kuehne et al. A Hybrid RNN-HMM Approach for Weakly Supervised Temporal Action Segmentation. PAMI 2020]

[A. Richard et al. Weakly Supervised Action Learning with RNN based Fine-to-Coarse Modeling. CVPR 2017]



Incremental learning



[A. Richard et al. NeuralNetwork-Viterbi: A Framework for Weakly Supervised Video Learning. CVPR 2018]

03.08.2020





	Breakfast frame accuracy (%)
pseudo-GT (HMM+RNN) Richard et al. 2017	33.3
pseudo-GT (HMM+RNN) Kuehne et al. 2020	36.7
NN-Viterbi (HMM+RNN) Richard et al. 2018	43.0
<i>Fully supervised (HMM+RNN)</i> Kuehne et al. 2020	61.3

[A. Richard et al. NeuralNetwork-Viterbi: A Framework for Weakly Supervised Video Learning. CVPR 2018]

03.08.2020

Pseudo GT vs. NN-Viterbi







Evaluation Issues

 Weakly supervised approaches are sensitive to initialization

Model	MoF	MoF	MoF	MoF
	Reported	Avg (\pm Std)	Max	Min
ISBA [4]	38.4	36.4 (± 1.0)	37.6	35.1
NNV [13]	43.0	39.7 (± 2.4)	43.5	37.5
CDFL [11]	50.2	48.1 (± 2.5)	50.9	44.6

[L. Ding and C. Xu. Weakly-supervised action segmentation with iterative soft boundary assignment. CVPR 2018]

- [J. Li et al. Weakly supervised energy-based learning for action segmentation. ICCV 2019]
- [Y. Souri et al. On Evaluating Weakly Supervised Action Segmentation Methods. arXiv]



Features

- Some approaches struggle with pre-trained features (I3D)
- Dimensionality is just one issue

Approach	Features	Average MoF
NNV	IDT	40.6
NNV	I3D	11.4
NNV	PCA-I3D	23.2
CDFL	IDT	48.9
CDFL	I3D	34.9
CDFL	PCA-I3D	38.0

[Y. Souri et al. On Evaluating Weakly Supervised Action Segmentation Methods. arXiv]



Training video



• Fully supervised:



• Weakly supervised (transcripts) $A \rightarrow C \rightarrow F \rightarrow D \rightarrow A \rightarrow E \rightarrow H$

> [A. Richard et al. Weakly Supervised Action Learning with RNN based Fine-to-Coarse Modeling. CVPR 2017]



• Fully supervised:



- Weakly supervised (transcripts) $A \rightarrow C \rightarrow F \rightarrow D \rightarrow A \rightarrow E \rightarrow H$
- Weakly supervised (action set) {A, C, D, E, F, H}
 - Order unknown
 - Number of occurrence unknown

[M. Fayyaz and J. Gall. SCT: Set Constrained Temporal Transformer for Set Supervised Action Segmentation. CVPR 2020]





	Supervision	frame accuracy (%)
SCT (TCN+I3D) Fayyaz and Gall 2020	Action set	30.4
pseudo-GT (HMM+RNN) Kuehne et al. 2020	Transcript	36.7
NN-Viterbi (HMM+RNN) Richard et al. 2018	Transcript	43.0
HMM+RNN Kuehne et al. 2020	Full	61.3
MS-TCN++ (TCN+I3D) Li et al. arXiv	Full	67.6

Source Code



- MS-TCN: <u>https://github.com/yabufarha/ms-tcn</u>
- ISBA: https://github.com/Zephyr-D/TCFPN-ISBA
- NN-Viterbi: <u>https://github.com/alexanderrichard/NeuralNetwork-</u> <u>Viterbi</u>
- CDFL: https://github.com/JunLi-Galios/CDFL
- Action sets: <u>https://github.com/alexanderrichard/action-sets</u>
- SCT: <u>https://github.com/MohsenFayyaz89/SCT</u> (Codes not uploaded yet)
- Unsupervised learning: <u>https://github.com/Annusha/unsup_temp_embed</u>

Thank you for your attention.



European Research Council







03.08.2020