



Video Event Detection

Prof. Kuan-Ting Lai

2019/10/28









ets & Animals













flow To Write a Book

Books & Literature



Human Action Recognition





http://www.thumos.info/



Applications of Video Event Detection

Video Search



Drone Action Recognition



Surveillance Video Analysis



Driver Activity Detection



Complex Video Events

Basic event or action detection





push

pushup



ride bike



ride

horse





run



hands



shoot ball

Complex or high-level event detection



Wedding Ceremony

Detect Complex (High-level) Events



• Jiang et al., high level events recognition in unconstrained videos, 2012 6

Detect Complex (High-level) Events



• Jiang et al., high level events recognition in unconstrained videos, 2012 7

Detect Complex (High-level) Events



• Jiang et al., high level events recognition in unconstrained videos, 2012 ^a

TRECVID Contest

- NIST TRECVID Multimedia Event Detection (MED) contest
- Detecting complex events in around 100,000 videos





Attempting a board trick



Feeding an animal



Landing a fish



Wedding ceremony



Working on a woodworking project



MED 2011

testing





Changing a vehicle tire



Flash mob gathering



Getting a vehicle unstuck





Grooming an animal





Making a sandwich



Parade







Repairing an appliance Working on a sewing project

Challenges of Video Event Detection

- Consist of interactions between human, objects and scenes
- Actions contain spatial data and temporal data
- Temporal data are huge and noisy
- Hard to define action classes

Optical Flow for Temporal Information



https://devblogs.nvidia.com/an-introduction-to-the-nvidia-optical-flow-sdk/

Dense Trajectories

- H. Wang, et al. Dense trajectories and motion boundary descriptors for action recognition. IJCV, 2013
- Tracking dense pixels has shown better than KLT (corner) features and SIFT points



KLT trajectories



SIFT trajectories



Dense trajectories

Extracting Dense Trajectories

- 1. Sample feature points with 5-pixel step
- 2. Remove features in homogeneous areas
- 3. Track optical flows smoothed by 3x3 median filter in 15 frames
- 4. Extract HOG, HOF and MBH along each trajectory



Motion Boundary Histogram (MBH)

- First proposed by Dalal & Triggs
 - "Human Detection Using Oriented Historgrams of Flow and Appearance", ECCV, 2006
- Show best results on HMDB51 & TRECVID 2011*
- Based on oriented histogram of differential optical flows
- Can effectively cancel camera motions



 *A. Tamarakar, et al., Evaluation of Low-Level Features and their Combinations for Complex Event Detection in Open Source Videos, CVPR, 2012

Improved Trajectories

- Cancel camera motions by matching SURF points & dense optical flows between frames
- A human detector is also used



Learning Optical Flow with Deep Learning?

• Karpathy et al., "Large-scale Video Classification with Convolutional Neural Networks", 2014



http://blog.qure.ai/notes/deep-learning-for-videos-action-recognition-review

Two-stream Convolutional Neural Networks



Simonyan, K., & Zisserman, A. (2014). Two-stream convolutional networks for action recognition in videos. *NIPS*

State-of-the-Arts

• <u>LRCN</u> [15], <u>C3D</u> [16], <u>Conv3D & Attention</u> [17], <u>TwoStreamFusion</u> [18], <u>TSN</u> [19], <u>ActionVlad</u> [20], <u>HiddenTwoStream</u> [1] <u>I3D</u> [21] and <u>T3D</u> [22]



http://blog.qure.ai/notes/deep-learning-for-videos-action-recognition-review

Long-term Recurrent Convolutional Networks

- Donahue et al., "Long-term Recurrent Convolutional Networks for Visual Recognition and Description," 2014 (<u>Arxiv Link</u>)
- Key Contributions:
 - Building on previous work by using RNN as opposed to stream based designs
 - Extension of encoder-decoder architecture for video representations
 - End-to-end trainable architecture proposed for action recognition

Long-term Recurrent Convolutional Networks

Activity Recognition Sequences in the Input CNN CNN CNN LSTM LSTM LSTM Average HighJump



3D Convolutional Networks (C3D)

- Du Tran et al., "Learning Spatiotemporal Features with 3D Convolutional Networks," 2014 (<u>Arxiv Link</u>)
- Key Contributions
 - Repurposing 3D convolutional networks as feature extractors
 - Extensive search for best 3D convolutional kernel and architecture
 - Using deconvolutional layers to interpret model decision



3D Convolutional Networks (C3D)

- Extract features on 2-second clip
- C3D tends to focus on spatial appearance in first few frames and tracked the motion in the subsequent frames





Human Action Recognition using Factorized Spatio-Temporal Convolutional Networks

Conv3D + Attention

• Yao et al., Describing Videos by Exploiting Temporal Structure, 2015



Temporal Segment Networks (TSN)

- Wang et al., "Temporal Segment Networks: Towards Good Practices for Deep Action Recognition", 2016
- Sampling clips sparsely across the video to better model long range



Hidden Two Stream

- Zhu et al., "Hidden Two-Stream Convolutional Networks for Action Recognition," 2017
- Novel architecture for generating optical flow input on-the-fly using a separate network



YouTube-8M

<u>https://research.google.com/youtube8m/</u>





- <u>https://neurohive.io/en/datasets/new-datasets-for-action-recognition/</u>
- <u>http://blog.qure.ai/notes/deep-learning-for-videos-action-recognition-review</u>