# S$^3$D: **S**ingle **S**hot multi-**S**pan **D**etector via Fully 3D Convolutional Network

**Da Zhang**[1], Xiyang Dai[2], Xin Wang[1], and Yuan-Fang Wang[1]

dazhang@cs.ucsb.edu

[1]UC Santa Barbara & [2]University of Maryland

BMVC 2018
Newcastle upon Tyne

# Task: Temporal Activity Detection

Input: untrimmed videos



1. *Localization*: when do activities start/end?

2. *Classification*: what are the activities?
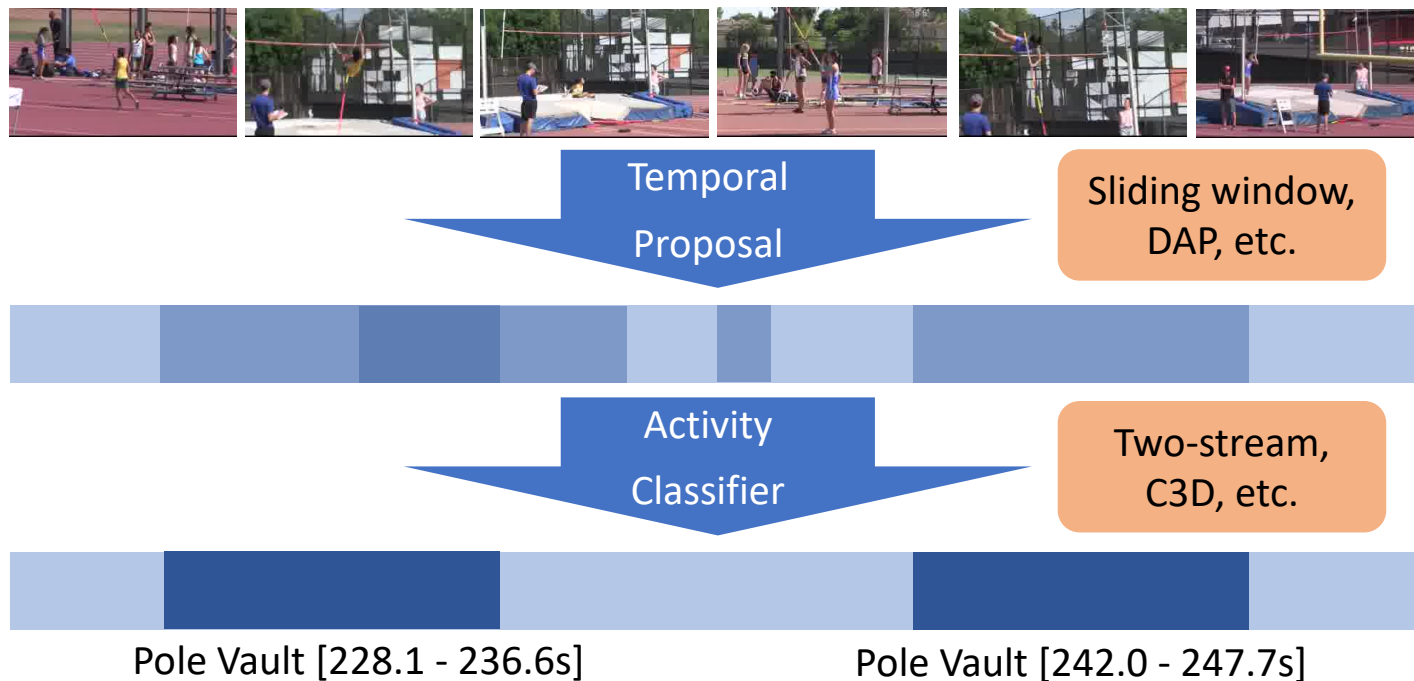
Detection Results



Pole Vault

[228.1 - 236.6s]
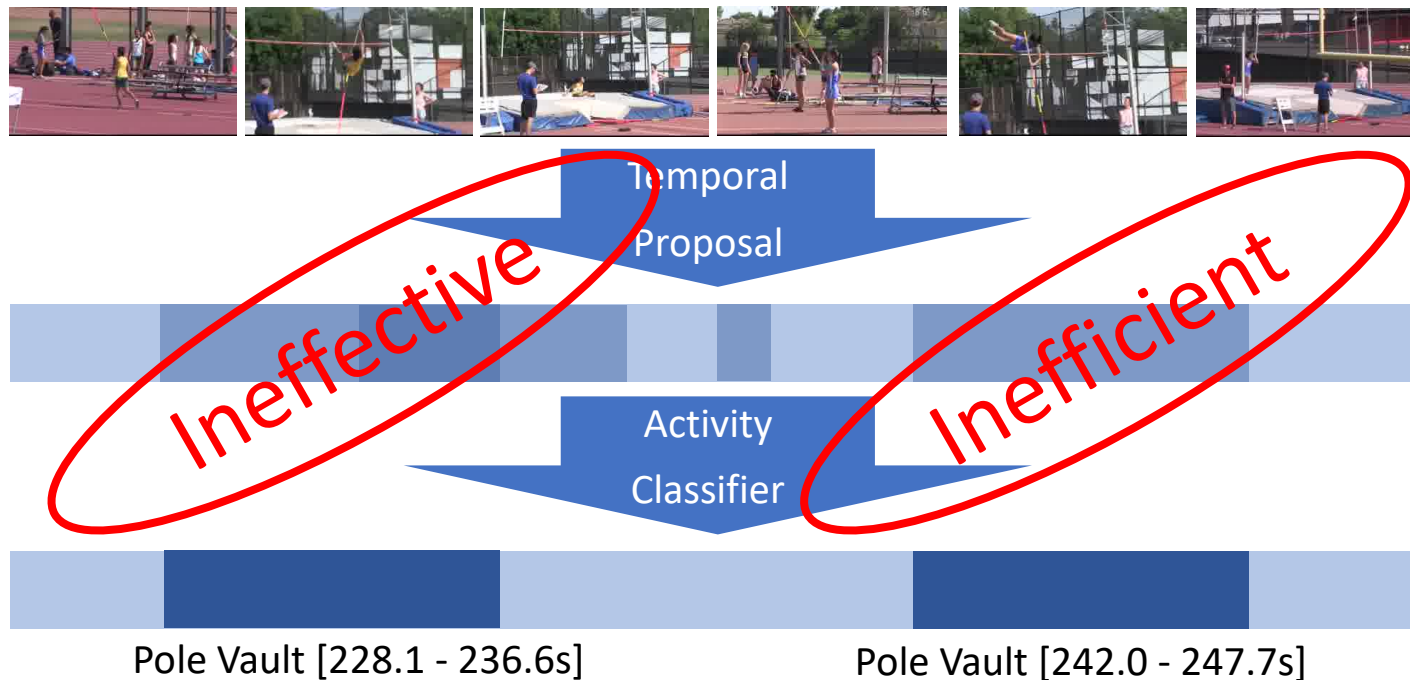
Pole Vault

[242.0 - 247.7s]

# Related Works

**Conventional two-stage approach: Proposal + Classification**



S-CNN (CVPR 2016), CDC (CVPR 2017), TSN (ICCV 2017), R-C3D (ICCV 2017), SSN (ICCV 2017)

# Related Works

**Current limitations:**



Temporal Proposal

*Ineffective*   *Inefficient*

Activity Classifier

Pole Vault [228.1 - 236.6s]          Pole Vault [242.0 - 247.7s]

S-CNN (CVPR 2016), CDC (CVPR 2017), TSN (ICCV 2017), R-C3D (ICCV 2017), SSN (ICCV 2017)

# Motivation

**Can we do better?**



Single-shot
End-to-end

Pole Vault [228.1 - 236.6s]          Pole Vault [242.0 - 247.7s]

Introducing a novel **S**ingle **S**hot multi-**S**pan **D**etector (S$^3$D)

# Motivation

**Quick Summary**





Single-shot
End-to-end

Pole Vault [228.1 - 236.6s]            Pole Vault [242.0 - 247.7s]
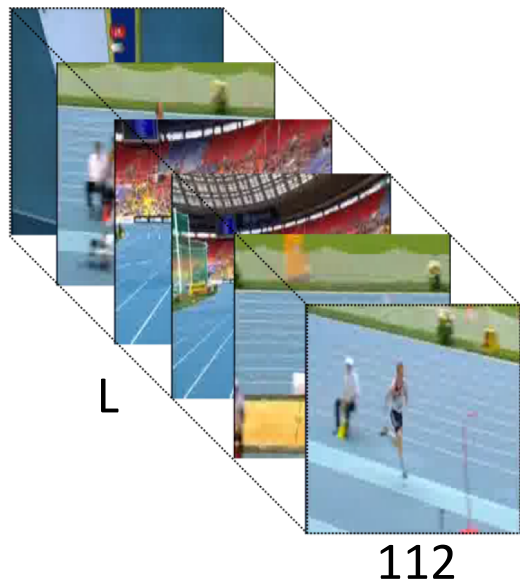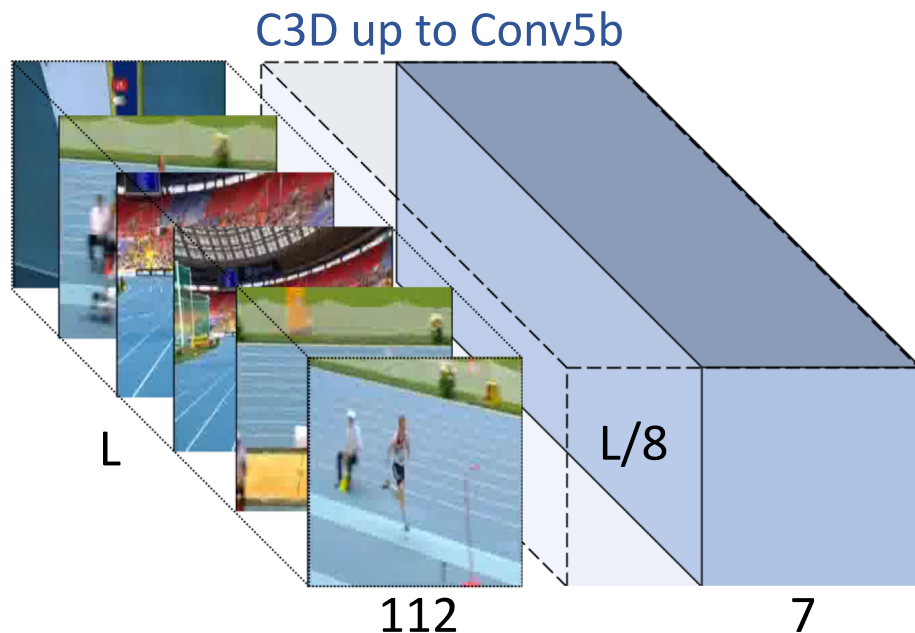
❑ Directly encode entire input video with Conv3D kernels

❑ Multi-scale default spans associated to temporal feature maps

❑ End-to-end trainable and single forward-pass inference

# S³D: Input Video



L

112

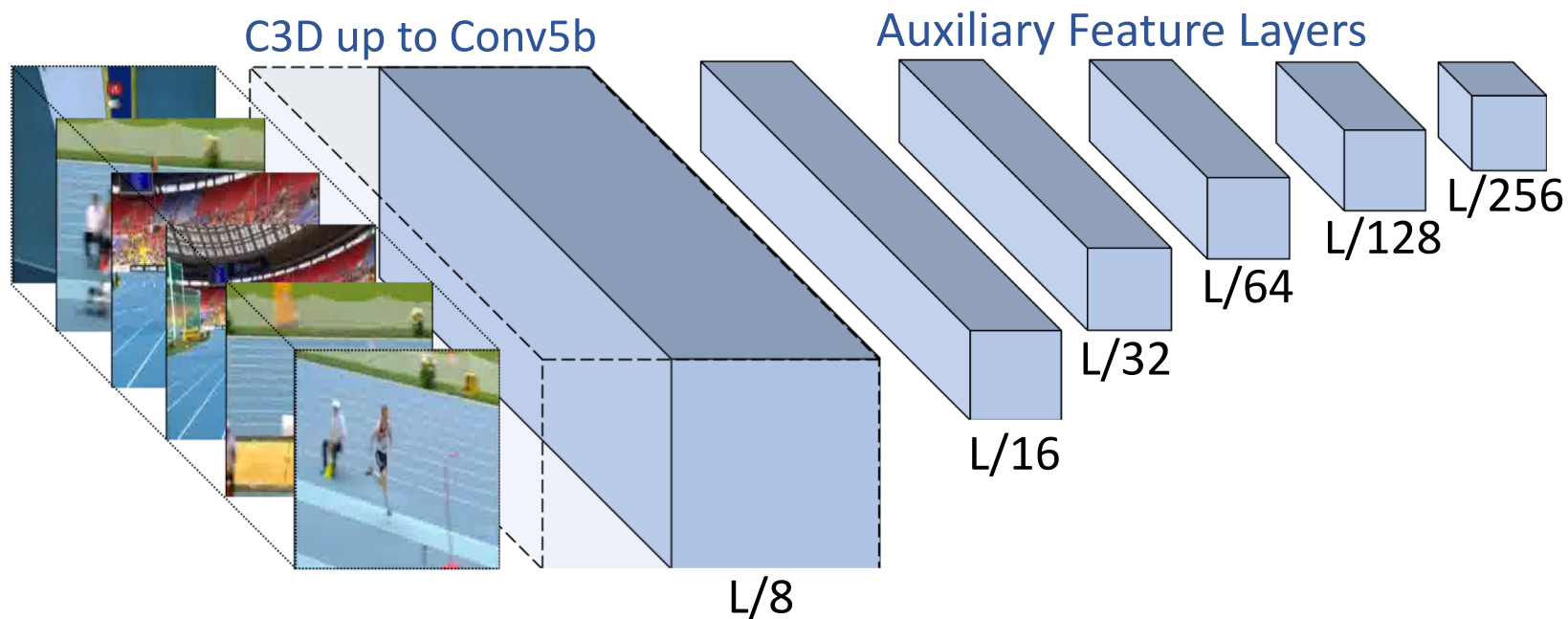Our model takes the whole video stream as input (L frames)

# S³D: Base Feature Layers



C3D up to Conv5b

L

112

L/8

7

We apply the standard C3D network to extract spatial-temporal features.

D. Tran, L. Bourdev, R. Fergus, L. Torresani and M. Paluri. Learning spatiotemporal features with 3D convolutional networks. In CVPR, 2015.
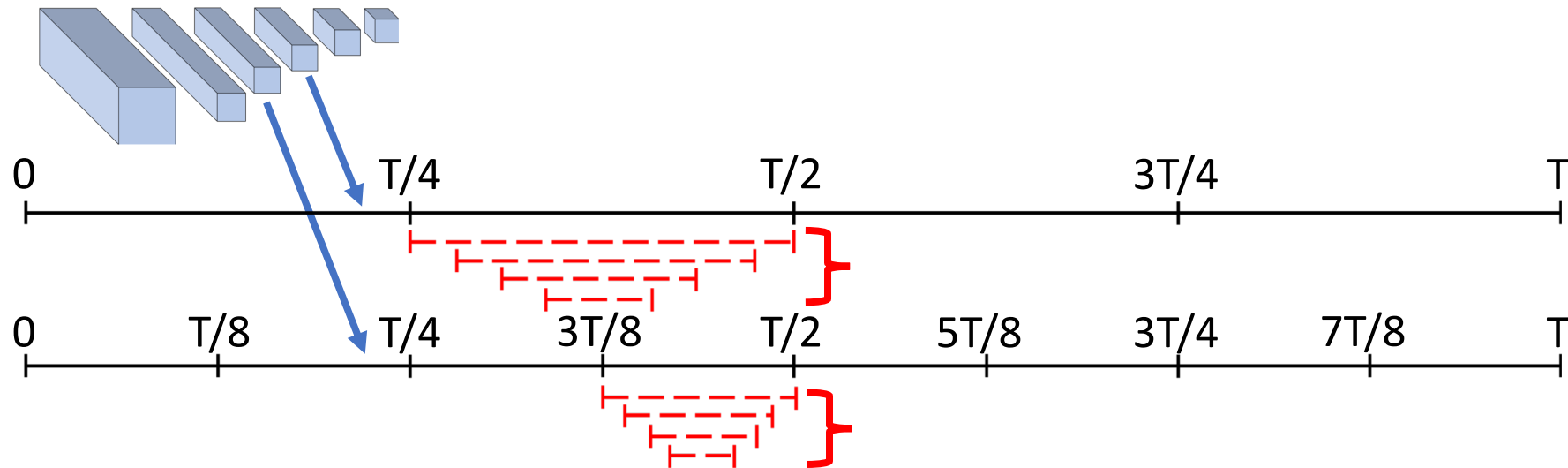
# S³D: Auxiliary Feature Layers

C3D up to Conv5b

Auxiliary Feature Layers



L/8

L/16

L/32

L/64

L/128

L/256

We produce a sequence of feature maps that progressively decrease in temporal dimension.

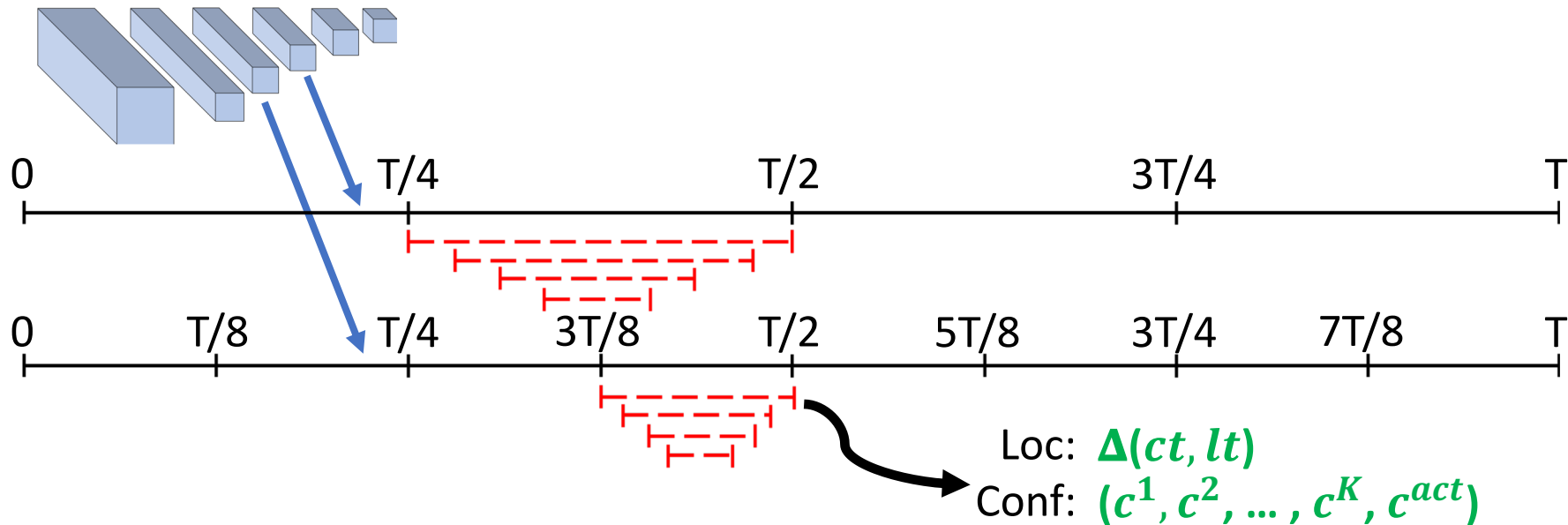# S³D: Multi-scale Default Spans

Temporal Feature Layers



Multi-scale default spans are associated to each temporal feature map
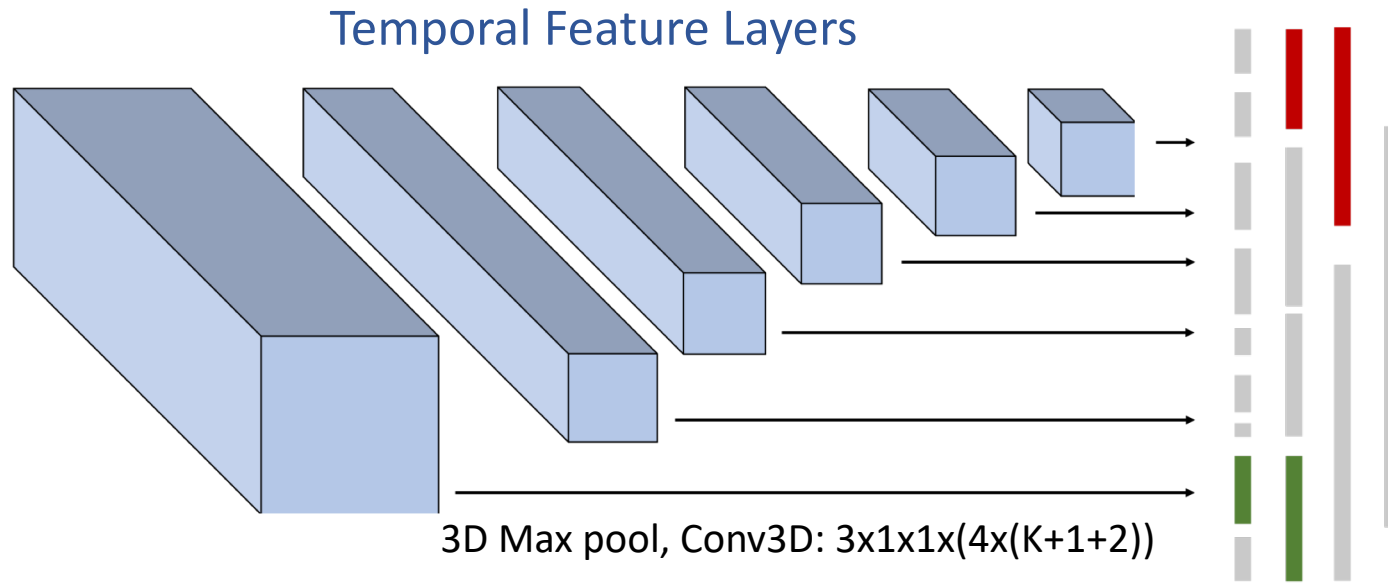
# S³D: Multi-scale Default Spans

**Temporal Feature Layers**



Loc: $\boldsymbol{\Delta(ct, lt)}$
Conf: $\boldsymbol{(c^1, c^2, \ldots, c^K, c^{act})}$

Localization and classification results are predicted at each default span.

# S³D: Convolutional Predictors

Temporal Feature Layers



3D Max pool, Conv3D: 3x1x1x(4x(K+1+2))
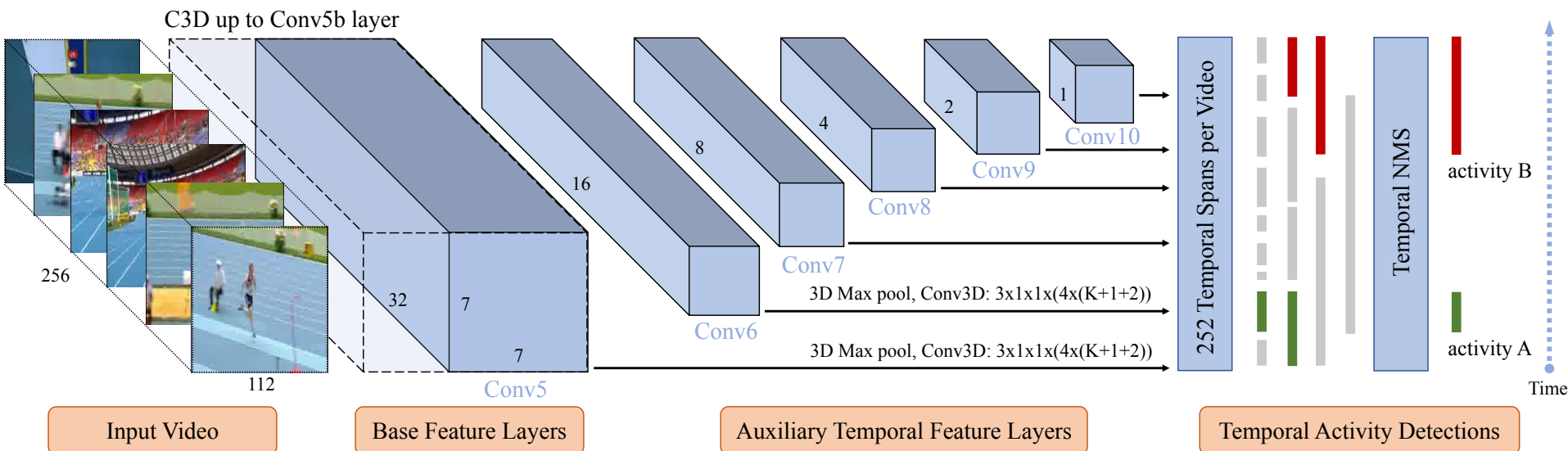
We apply on top of each feature map a Conv3D filter to produce the results.

# S³D: Convolutional Predictors



Temporal Feature Layers

3D Max pool, Conv3D: 3x1x1x(4x(K+1+2))

Kernel size    # of scales    Classes + BG    Localization offsets

$(c^1, c^2, \ldots, c^K, c^{act})$    $\Delta(ct, lt)$

# Single Shot multi-Span Detector



Training of S³D:

$$Loss = \boxed{L_{loc}(x,t,g)} + \alpha \boxed{L_{conf}(x,c)} + \beta \boxed{L_{act}(s,c)}$$

Smooth L1        Softmax Cross Entropy        Sigmoid Cross Entropy

# Quantitative Results

Evaluation: mean Average Precision over 20 activities on THUMOS'14

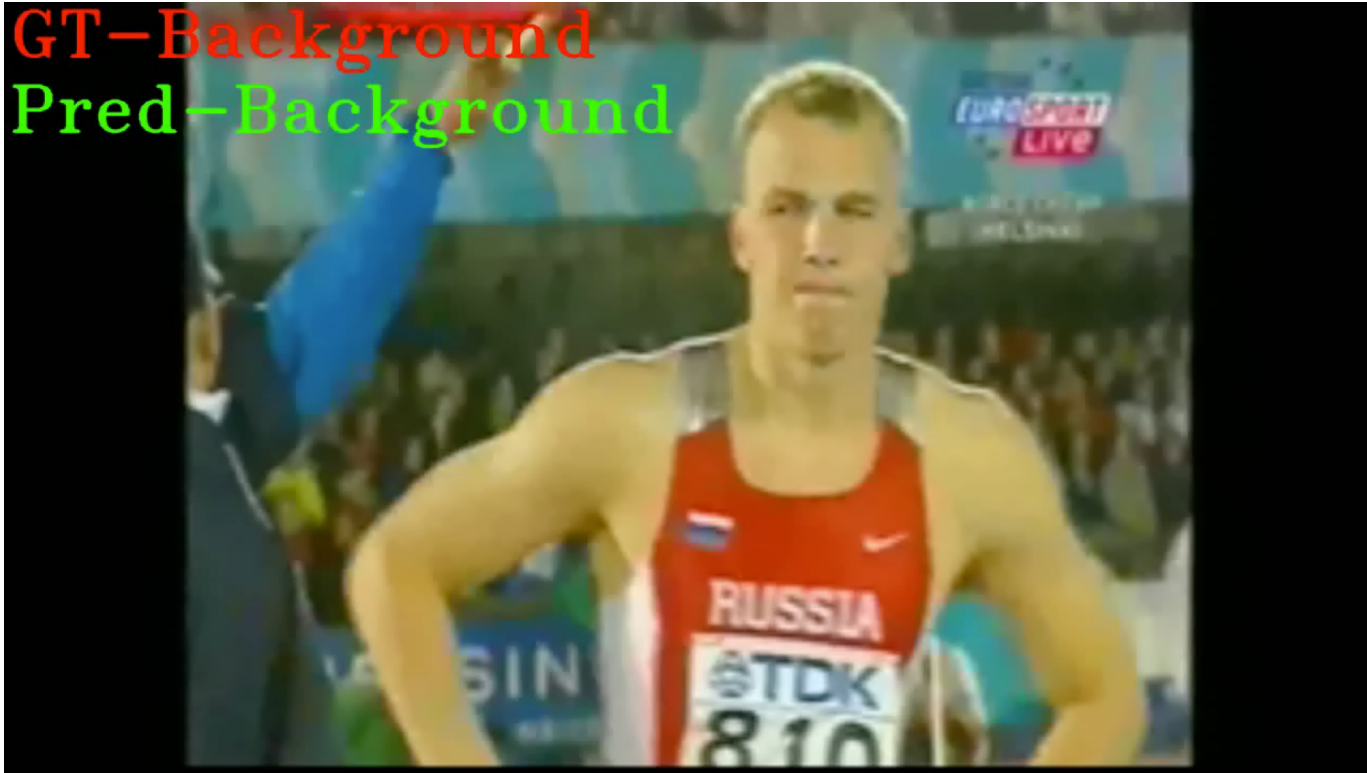| IoU threshold | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 |
|---|---|---|---|---|---|
| S-CNN (CVPR 2016) | 36.3 | 28.7 | 19.0 | 10.3 | 5.3 |
| CDC (CVPR 2017) | 40.1 | 29.4 | 23.3 | 13.1 | 7.9 |
| SSAD (MM 2017) | 43.0 | 35.0 | 24.6 | - | - |
| TCN (ICCV 2017) | - | 33.3 | 25.6 | 15.9 | 9.0 |
| R-C3D (ICCV 2017) | 44.8 | 35.6 | 28.9 | - | - |
| SSN (ICCV 2017) | **50.6** | 40.8 | 29.1 | - | - |
| SS-TAD (BMVC 2017) | 40.1 | - | 29.2 | - | 9.6 |
| $S^3D$ (ours) | 47.9 | **41.2** | **32.6** | **23.3** | **14.3** |

1271 FPS on a single GTX 1080 Ti GPU

# Qualitative Results

THUMOS'14 segment: Pole Vault

# Qualitative Results

THUMOS'14 segment: Javelin Throw

# Qualitative Results

THUMOS'14 segment: Shotput

# Qualitative Results

THUMOS'14 segment: Clean and Jerk

# Conclusions

Introduced **S³D**:

- ❑ A novel single-shot end-to-end model for Temporal Activity Detection.
- ❑ *Simple*: completely based on Conv3D kernels.
- ❑ *Strong*: state-of-the-art performance on THUMOS'14 benchmark.
- ❑ *Speed*: operates at 1271 FPS on a single GeForce GTX 1080 Ti GPU.

TensorFlow code coming soon at https://github.com/dazhang-cv/S3D

# Thank you!



C3D up to Conv5b layer

256 · 112

32 · 7 · 7 — Conv5

16 — Conv6

8 — Conv7

4 — Conv8

2 — Conv9

1 — Conv10

3D Max pool, Conv3D: 3x1x1x(4x(K+1+2))

3D Max pool, Conv3D: 3x1x1x(4x(K+1+2))

252 Temporal Spans per Video

Temporal NMS

activity B

activity A

Time

**Input Video**  **Base Feature Layers**  **Auxiliary Temporal Feature Layers**  **Temporal Activity Detections**