GabriellaV2: UCF DIVA System Center for Research in Computer Vision (CRCV) at the University of Central Florida (UCF)

Zacchaeus Scheffer, Ishan Dave, Mubarak Shah, Yogesh Rawat

December 8, 2021



UCF CENTER FOR RESEARCH IN COMPUTER VISION

◆□ ▶ ◆昼 ▶ ◆臣 ▶ ◆臣 ▶ ○ ● ○ ○ ○

Problem Statement





Problem StatementNew in GabriellaV2



Problem Statement
 New in GabriellaV2
 Localization



Problem Statement
 New in GabriellaV2
 Localization

Action Classification

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●



Problem Statement

- New in GabriellaV2
 - Localization
 - Action Classification

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Post Processing



Problem Statement

- New in GabriellaV2
 - Localization
 - Action Classification

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Post Processing





Problem Statement

- New in GabriellaV2
 - Localization
 - Action Classification

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

- Post Processing
- Scores
- Qualitative Results



- Problem Statement
- New in GabriellaV2
 - Localization
 - Action Classification

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

- Post Processing
- Scores
- Qualitative Results

Questions

UCF CENTER FOR RESEARCH IN COMPUTER VISION



▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへの

Untrimmed Video



(ロ) (型) (E) (E) (E) (O)

Untrimmed Video

► No special processing



Untrimmed Video

 No special processing
 Should be comparable to "real-world"

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●



Untrimmed Video

 No special processing
 Should be comparable to "real-world"

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Detect activities



Untrimmed Video

 No special processing
 Should be comparable to "real-world"

Detect activities

Human and vehicle

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □



Untrimmed Video

 No special processing
 Should be comparable to "real-world"

Detect activities

- Human and vehicle
- Indoor and outdoor

UCF CENTER FOR RESEARCH IN COMPUTER VISION

Untrimmed Video

Types

 No special processing
 Should be comparable to "real-world"

Detect activities

- Human and vehicle
- Indoor and outdoor

UCF CENTER FOR RESEARCH IN COMPUTER VISION

Untrimmed Video

 No special processing
 Should be comparable to "real-world"

Detect activities

- Human and vehicle
- Indoor and outdoor



Untrimmed Video

 No special processing
 Should be comparable to "real-world"

Detect activities

- Human and vehicle
- Indoor and outdoor

Types



Multi-actor

Untrimmed Video

 No special processing
 Should be comparable to "real-world"

Detect activities

- Human and vehicle
- Indoor and outdoor

Types

- Single-actor
- Multi-actor
- Actor-object

Untrimmed Video

 No special processing
 Should be comparable to "real-world"

Detect activities

- Human and vehicle
- Indoor and outdoor

Types

- Single-actor
- Multi-actor
- Actor-object

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ めへで

Output

Untrimmed Video

 No special processing
 Should be comparable to "real-world"

Detect activities

- Human and vehicle
- Indoor and outdoor

Types

- Single-actor
- Multi-actor
- Actor-object

Output

Start/end times

Untrimmed Video

 No special processing
 Should be comparable to "real-world"

Detect activities

- Human and vehicle
- Indoor and outdoor

Types

- Single-actor
- Multi-actor
- Actor-object

Output

- Start/end times
- Spacial location

▲□▶ ▲□▶ ▲ 臣▶ ▲ 臣▶ ― 臣 … のへで

GabriellaV2 System



▲□▶▲圖▶★≧▶★≧▶ ≧ の�@





First, we split the video using a sliding window approach



First, we split the video using a sliding window approach





First, we split the video using a sliding window approach

- ▶ 0-16
- ▶ 8-24



First, we split the video using a sliding window approach

- ▶ 0-16
- ▶ 8-24▶ 16-32



- First, we split the video using a sliding window approach
 - ▶ 0-16
 - ▶ 8-24▶ 16-32
- Every 4th frame used for tracklet generation



- First, we split the video using a sliding window approach
 - 0-16
 8-24
 - ▶ 16-32
- Every 4th frame used for tracklet generation
 Full clip is used in tracklet extraction/classification





Tracklet Generation happens in 3 steps



◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ ○ ♥ ♥

Tracklet Generation happens in 3 steps▶ Object Detector (YOLOv5)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00



Tracklet Generation happens in 3 steps Object Detector (YOLOv5) spacially localizes objects


Tracklet Generation happens in 3 steps Object Detector (YOLOv5) spacially localizes objects Background Subtractor (MOG)



Tracklet Generation happens in 3 steps Object Detector (YOLOv5) spacially localizes objects Background Subtractor (MOG) Removes Stationary Objects



Tracklet Generation happens in 3 steps
Object Detector (YOLOv5)
spacially localizes objects
Background Subtractor (MOG)
Removes Stationary Objects
Object Tracker (SORT)



Tracklet Generation happens in 3 steps Object Detector (YOLOv5) spacially localizes objects Background Subtractor (MOG) Removes Stationary Objects Object Tracker (SORT) groups detections of the same object



For each Clip





くしゃ (四)・(日)・(日)・(日)・

For each Clip

send 4 frames to
 Object Detector





For each Clip

 send 4 frames to Object Detector
 Place bounding box around potential actor





For each Clip

- send 4 frames to
 Object Detector
- Place bounding box around potential actor
 - person





For each Clip

- send 4 frames to
 Object Detector
- Place bounding box around potential actor
 - personvehicle

CENTER FOR RESEARCH



Gaussian Blur Gray scale Conversion	MOG2 background subtraction	Morphological operation		
--	-----------------------------------	-------------------------	--	--





▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

We use the MOG2 background subtractor





▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ めへで

We use the MOG2 background subtractor
simple to use





▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00

- We use the MOG2 background subtractor
- simple to use
- drastically reduce localization false alarm





▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ めへで

We use the MOG2 background subtractor

- ► simple to use
- drastically reduce localization false alarm
 - 40% FA reduction





▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00

- We use the MOG2 background subtractor
- simple to use
- drastically reduce localization false alarm
 - 40% FA reduction
 - Runtime reduced by half

UCF CENTER FOR RESEARCH IN COMPUTER VISION

Background Subtractor Results





◆□▶ ◆□▶ ◆三▶ ◆三▶ → □ ◆ ○へ⊙



UCF CENTER FOR RESEARCH IN COMPUTER VISION

► We use the SORT tracker



We use the SORT trackerIOU-based matching

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで



▶ We use the SORT tracker

- ► IOU-based matching
- Same object gets same ID in subsequent frames

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00



- We use the SORT tracker
- ► IOU-based matching
- Same object gets same ID in subsequent frames
 Can "remember" 2 frames prior



- We use the SORT tracker
- ► IOU-based matching
- Same object gets same ID in subsequent frames
- Can "remember" 2 frames prior
 - Track continues if one bounding box is missing



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●



Center for Research in Computer Vision



 Legacy method used action localization instead of object detection/tracking

▲□▶▲□▶▲目▶▲目▶ 目 のへぐ



Legacy method used action localization instead of object detection/tracking New method attains higher activity recall than legacy method

・ 日本 《四本 《日本 《日本 《日本



Legacy method used action localization instead of object detection/tracking New method attains higher activity recall than legacy method ▶ 22% higher recall using 0.8 IOU threshold





For each object ID produced by the SORT tracker in a given clip:

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで



For each object ID produced by the SORT tracker in a given clip:

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00

collect all bounding boxes with that object ID



For each object ID produced by the SORT tracker in a given clip:

- collect all bounding boxes with that object ID
- Combine them by taking the smallest-bounding bounding box



For each object ID produced by the SORT tracker in a given clip:

- collect all bounding boxes with that object ID
- Combine them by taking the smallest-bounding bounding box
- Extend the resulting bounding box in the shorter dimension to obtain a square



For each object ID produced by the SORT tracker in a given clip:

- collect all bounding boxes with that object ID
- Combine them by taking the smallest-bounding bounding box
- Extend the resulting bounding box in the shorter dimension to obtain a square
 - classifier gets consistent aspect ratio



For each object ID produced by the SORT tracker in a given clip:

- collect all bounding boxes with that object ID
- Combine them by taking the smallest-bounding bounding box
- Extend the resulting bounding box in the shorter dimension to obtain a square
 - classifier gets consistent aspect ratio

► Take spatial crop of clip. <- Extracted Tracklet UCF CENTER FOR RESEARCH IN COMPUTER VISION

Action Classification



Action Classification

Surveillance video activities are multi-label



Action Classification

Surveillance video activities are multi-label Ground truth tracks to train 3D-CNN backbone + sigmoid activation with standard BCE loss

◆□▶ ◆□▶ ◆三▶ ◆三▶ → □ ◆ ○へ⊙


Action Classification

Surveillance video activities are multi-label Ground truth tracks to train 3D-CNN backbone + sigmoid activation with standard BCE loss

$$\mathcal{L}_{\mathsf{BCE}}(y,\hat{y}) = -rac{1}{N}\sum_{i=1}^{N}\left[y_i\log(\hat{y}_i) - (1-y_i)\log(1-\hat{y}_i)
ight]$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ → □ ◆ ○へ⊙



Action Classification

Surveillance video activities are multi-label Ground truth tracks to train 3D-CNN backbone + sigmoid activation with standard BCE loss

$$\mathcal{L}_{\mathsf{BCE}}(y,\hat{y}) = -rac{1}{N}\sum_{i=1}^{N}\left[y_i\log(\hat{y}_i) - (1-y_i)\log(1-\hat{y}_i)
ight]$$

Each action class is independent of each other



Adapting to Yolo Tracklets

Need to train classifier to take YOLO tracklets



Adapting to Yolo Tracklets

Need to train classifier to take YOLO tracklets



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00



Number of Samples per Class





Number of Samples per Class



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00

Some classes are very frequent, some very rare



Number of Samples per Class



Some classes are very frequent, some very rare In the most extreme, 1:1000 difference



Number of Samples per Class



Some classes are very frequent, some very rare In the most extreme, 1:1000 difference Used Partial Label Masking (PLM) for class imbalance UCF CENTER FOR RESEARCH



◆□ ▶ ◆□ ▶ ◆三 ▶ ◆三 ▶ ○ ● ● ●

Log-Sum-Exponential Pairwise (LSEP) loss for multilabel action recognition.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00



Log-Sum-Exponential Pairwise (LSEP) loss for multilabel action recognition.

$$\mathcal{L}_{\mathsf{LSEP}} = \mathsf{log}\left(1 + \sum_{i \in y} \sum_{j
otin y} e^{x_j - x_i}
ight)$$



Log-Sum-Exponential Pairwise (LSEP) loss for multilabel action recognition.

$$\mathcal{L}_{\mathsf{LSEP}} = \mathsf{log}\left(1 + \sum_{i \in y} \sum_{j
otin y} e^{x_j - x_i}
ight)$$

Standard LSEP is based on BCE loss



Log-Sum-Exponential Pairwise (LSEP) loss for multilabel action recognition.

$$\mathcal{L}_{\mathsf{LSEP}} = \mathsf{log}\left(1 + \sum_{i \in y} \sum_{j
otin y} e^{x_j - x_i}
ight)$$

Standard LSEP is based on BCE loss Use LSEP loss with PLM to reweight samples to class balance and learn class correlations together

We perform Knowledge Distillation in two stages:



We perform Knowledge Distillation in two stages:standard phase



We perform Knowledge Distillation in two stages:

- standard phase
 - reduce ensemble size



We perform Knowledge Distillation in two stages:

- standard phase
 - reduce ensemble size
- compression phase



We perform Knowledge Distillation in two stages:

- standard phase
 - reduce ensemble size
- compression phase
 - reduce model size



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00





Train student model with ensemble

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00







 Train student model with ensemble
 loss based on hidden

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00

layers





- Train student model with ensemble
- loss based on hidden layers
- loss on raw outputs





- Train student model with ensemble
- loss based on hidden layers
- loss on raw outputs
- loss on label prediction





- Train student model with ensemble
- loss based on hidden layers
- loss on raw outputs
- loss on label prediction
 val set mAP >4% improvment







- Train student model with ensemble
- loss based on hidden layers
- loss on raw outputs
- loss on label prediction
 val set mAP >4%
 - improvment
 - same mAP as fulla ensemble

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00





 Student model from before as Teacher

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 善 - のへ⊙



- Student model from before as Teacher
- New student is slimmer



- Student model from before as Teacher
- New student is slimmer

loss on raw outputs



- Student model from before as Teacher
- New student is slimmer
- loss on raw outputs
- ▶ loss on label prediciton



- Student model from before as Teacher
- New student is slimmer
- loss on raw outputs
- loss on label prediciton

 val set mAP >7% improvement



- Student model from before as Teacher
- New student is slimmer
- loss on raw outputs
- loss on label prediciton
- ▶ val set mAP >7%
 - improvement
- same mAP as full ensemble

Post-Processing



Post-Processing





▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Post-Processing

- ► TMAS
 - smooth detection temporally



Post-Processing

- ► TMAS
 - smooth detection temporally
 - convert actor tracks into action tubes

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00


Post-Processing

Post-Processing has two steps:

- ► TMAS
 - smooth detection temporally
 - convert actor tracks into action tubes

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ めへで

NMS



Post-Processing

Post-Processing has two steps:

- TMAS
 - smooth detection temporally
 - convert actor tracks into action tubes

- ► NMS
 - Remove duplicate predictions





UCF CENTER FOR RESEARCH IN COMPUTER VISION



UCF CENTER FOR RESEARCH IN COMPUTER VISION

Chain tracklets with same ID into one actor track



into one actor track Smooth classwise detection temporally

Chain tracklets with same ID

UCF CENTER FOR RESEARCH IN COMPUTER VISION



Chain tracklets with same ID into one actor track Smooth classwise detection temporally Create action tubes for each class using (connected) regions of actor tracks where class-scores is high

< ロ > < 酒 > < 直 > < 直 > < 三 > < 三 > < 三 > < ○へ ()











We square our object detections

◆□▶ ◆□▶ ◆三▶ ◆三▶ → □ ◆ ○へ⊙





 We square our object detections

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Very large overlaps





- We square our object detections
- Very large overlaps
- Actions contained in multiple bounding boxes





- We square our object detections
- Very large overlaps
- Actions contained in multiple bounding boxes
- multi-actor actions





- We square our object detections
- Very large overlaps
- Actions contained in multiple bounding boxes
- multi-actor actions
- Need to remove duplicate predictions





Need method to work regardless of actor size













$$d(A,B) = rac{d_e(A,B)}{\sqrt[4]{AreaA * Area(B)}}$$





$$d(A,B) = rac{d_e(A,B)}{\sqrt[4]{AreaA*Area(B)}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ → □ ◆ ○へ⊙

Where d_e is the euclidian distance.





$$d(A,B) = rac{d_e(A,B)}{\sqrt[4]{AreaA*Area(B)}}$$

Where d_e is the euclidian distance. If actors are close, treat as same action, if VERY close, same actor





$$d(A,B) = rac{d_e(A,B)}{\sqrt[4]{AreaA*Area(B)}}$$

Where d_e is the euclidian distance. If actors are close, treat as same action, if VERY close, same actor False Negatives very rare

Rank	team_name	team_abbrev	nAUDC@0.2tf a	p_miss@0.15t fa
1	BUPT-MCPRL	BUPT-MC_26542	0.4085	0.3249
2	UCF	UCF_26546	0.4306	0.3408
3	INF	INF_26532	0.4444	0.3508
4	M4D_2021	M4D_202_2646 7	0.8466	0.7941
5	TokyoTech_AIST	TOKYOTE_2650 8	0.8516	0.8197
6	Team UEC	TEAMUE_26530	0.964	0.9503

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ○ ○ ○ ○ ○



Rank	team_name	team_abbrev	nAUDC@0.2tf a	p_miss@0.15t fa
1	BUPT-MCPRL	BUPT-MC_26542	0.4085	0.3249
2	UCF	UCF_26546	0.4306	0.3408
3	INF	INF_26532	0.4444	0.3508
4	M4D_2021	M4D_202_2646 7	0.8466	0.7941
5	TokyoTech_AIST	TOKYOTE_2650 8	0.8516	0.8197
6	Team UEC	TEAMUE_26530	0.964	0.9503



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで



Rank	team_name	team_abbrev	nAUDC@0.2tf a	p_miss@0.15t fa
1	BUPT-MCPRL	BUPT-MC_26542	0.4085	0.3249
2	UCF	UCF_26546	0.4306	0.3408
3	INF	INF_26532	0.4444	0.3508
4	M4D_2021	M4D_202_2646 7	0.8466	0.7941
5	TokyoTech_AIST	TOKYOTE_2650 8	0.8516	0.8197
6	Team UEC	TEAMUE_26530	0.964	0.9503





Rank	team_name	team_abbrev	nAUDC@0.2tf a	p_miss@0.15t fa
1	BUPT-MCPRL	BUPT-MC_26542	0.4085	0.3249
2	UCF	UCF_26546	0.4306	0.3408
3	INF	INF_26532	0.4444	0.3508
4	M4D_2021	M4D_202_2646 7	0.8466	0.7941
5	TokyoTech_AIST	TOKYOTE_2650 8	0.8516	0.8197
6	Team UEC	TEAMUE_26530	0.964	0.9503

► based on VIRAT dataset has known cameras ► 2nd place



Rank	Team Name	sub_id	mean p_miss@0.01tf a	mean p_miss@0.02 tfa	mean nAUDC@0.2tfa	relative_ processi ng_time
1	UCF	25908	0.62	0.5372	0.3518	0.684
2	CMU-DIVA	26095	0.65	0.5438	0.333	0.776
3	IBM-Purdue	26113	0.65	0.5531	0.3533	0.575
4	UMD	26619	0.68	0.5938	0.3898	0.515
5	UMD- Columbia	25031	0.68	0.5975	0.4002	0.52
6	UMCMU	25576	0.75	0.6861	0.4922	0.614
7	Purdue	25782	0.8	0.7294	0.4942	0.239
8	MINDS_JHU	24666	0.84	0.7791	0.6343	0.898



Rank	Team Name	sub_id	mean p_miss@0.01tf a	mean p_miss@0.02 tfa	mean nAUDC@0.2tfa	relative_ processi ng_time
1	UCF	25908	0.62	0.5372	0.3518	0.684
2	CMU-DIVA	26095	0.65	0.5438	0.333	0.776
3	IBM-Purdue	26113	0.65	0.5531	0.3533	0.575
4	UMD	26619	0.68	0.5938	0.3898	0.515
5	UMD- Columbia	25031	0.68	0.5975	0.4002	0.52
6	UMCMU	25576	0.75	0.6861	0.4922	0.614
7	Purdue	25782	0.8	0.7294	0.4942	0.239
8	MINDS_JHU	24666	0.84	0.7791	0.6343	0.898

 based on MEVA dataset



Rank	Team Name	sub_id	mean p_miss@0.01tf a	mean p_miss@0.02 tfa	mean nAUDC@0.2tfa	relative_ processi ng_time
1	UCF	25908	0.62	0.5372	0.3518	0.684
2	CMU-DIVA	26095	0.65	0.5438	0.333	0.776
3	IBM-Purdue	26113	0.65	0.5531	0.3533	0.575
4	UMD	26619	0.68	0.5938	0.3898	0.515
5	UMD- Columbia	25031	0.68	0.5975	0.4002	0.52
6	UMCMU	25576	0.75	0.6861	0.4922	0.614
7	Purdue	25782	0.8	0.7294	0.4942	0.239
8	MINDS_JHU	24666	0.84	0.7791	0.6343	0.898





Rank	Team Name	sub_id	mean p_miss@0.01tf a	mean p_miss@0.02 tfa	mean nAUDC@0.2tfa	relative_ processi ng_time
1	UCF	25908	0.62	0.5372	0.3518	0.684
2	CMU-DIVA	26095	0.65	0.5438	0.333	0.776
3	IBM-Purdue	26113	0.65	0.5531	0.3533	0.575
4	UMD	26619	0.68	0.5938	0.3898	0.515
5	UMD- Columbia	25031	0.68	0.5975	0.4002	0.52
6	UMCMU	25576	0.75	0.6861	0.4922	0.614
7	Purdue	25782	0.8	0.7294	0.4942	0.239
8	MINDS_JHU	24666	0.84	0.7791	0.6343	0.898

- based on MEVA dataset
- has unknown

cameras

 1st place in UF



Generalization



UCF CMU UMD



Generalization



Drop in performance going from Known Facility to Unknown



Generalization



Drop in performance going from Known Facility to Unknown

Our system gets best generalization
UCF CENTER FOR RESEARCH
IN COMPUTER VISION

Qualitative Results

Many common situations that hurt performance



Qualitative Results

Many common situations that hurt performance

Distant actors



Qualitative Results

Many common situations that hurt performance

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ めへで

Distant actors

Actor changes distance over time



Qualitative Results

Many common situations that hurt performance

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00

Distant actors

- Actor changes distance over time
- Temporal variability



Qualitative Results

Many common situations that hurt performance

Distant actors

- Actor changes distance over time
- Temporal variability
- Crowded Scenes



Conclusion

GabriellaV2 is a real-time action detection system which can generalize very well to the unknown facility cameras.



Conclusion

GabriellaV2 is a real-time action detection system which can generalize very well to the unknown facility cameras.Built upon our Gabriella system by

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の00


GabriellaV2 is a real-time action detection system which can generalize very well to the unknown facility cameras.Built upon our Gabriella system by

Replacing localization by tracklet generation,



GabriellaV2 is a real-time action detection system which can generalize very well to the unknown facility cameras.Built upon our Gabriella system by

- Replacing localization by tracklet generation,
- Better classifiers: class imbalance, multi-label class correlation, Knowledge distillation



GabriellaV2 is a real-time action detection system which can generalize very well to the unknown facility cameras.Built upon our Gabriella system by

- Replacing localization by tracklet generation,
- Better classifiers: class imbalance, multi-label class correlation, Knowledge distillation
- Improved Post processing using Spatio-temporal deduplication



GabriellaV2 is a real-time action detection system which can generalize very well to the unknown facility cameras.Built upon our Gabriella system by

- Replacing localization by tracklet generation,
- Better classifiers: class imbalance, multi-label class correlation, Knowledge distillation
- Improved Post processing using Spatio-temporal deduplication

Top place for ActEV-SDL21 UF leaderboard and Runners-up for TRECVID21

References

Monfort, Mathew, et al. "Multi-moments in time: Learning and interpreting models for multi-action video understanding." arXiv preprint arXiv:1911.00232 (2019).

You, Shan, et al. "Learning from multiple teacher networks." Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2017.

Hinton, Geoffrey, Oriol Vinyals, and Jeff Dean. "Distilling the knowledge in a neural network." arXiv preprint arXiv:1503.02531 (2015).

Duarte, Kevin, Yogesh Rawat, and Mubarak Shah. "PLM: Partial Label Masking for Imbalanced Multi-label Classification." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.

UCF CENTER FOR RESEARCH IN COMPUTER VISION