# DL Seminar

## MOTDT

**REAL-TIME MULTIPLE PEOPLE TRACKING WITH DEEPLY LEARNED CANDIDATE SELECTION AND PERSON RE-IDENTIFICATION**
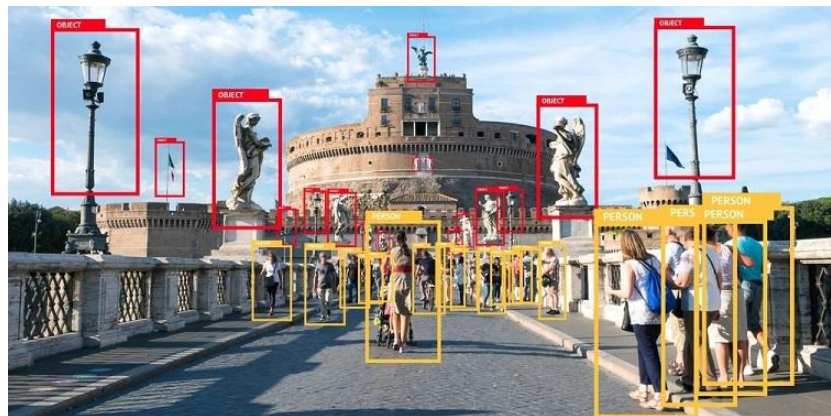
한양대학교
**HANYANG UNIVERSITY**

인공지능 연구실
김지성

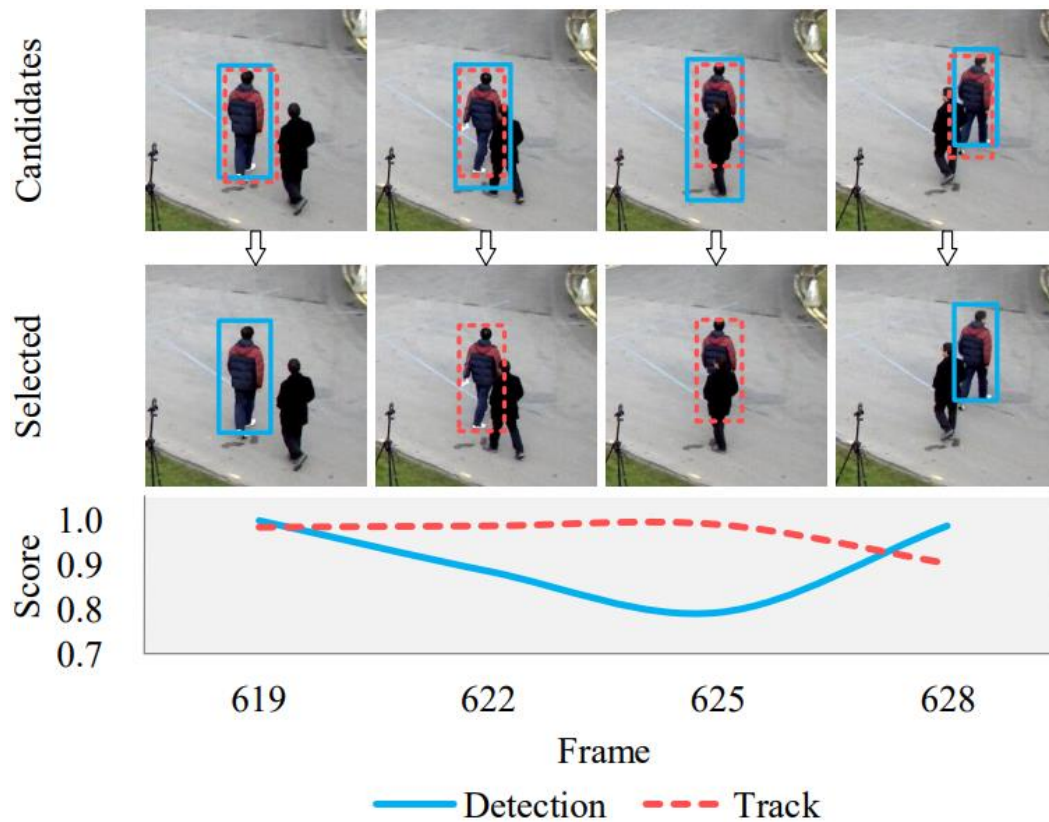# Introduction

Detector, Tracker



Detector



Frame 1

Frame 2

Frame 3

Tracker

# Introduction

Candidates

# Method

## 후보 필터링



**Algorithm 1:** The proposed tracking algorithm.

**Input:** A video sequence $v$ with $N_v$ frames and object detection $\{\mathcal{D}_k\}_{k=1}^{N_v}$

**Output:** Tracks $\mathcal{T}$ of the video

1. Initialization: $\mathcal{T} \leftarrow \emptyset$; appearance of tracks $\mathcal{F}_{trk} \leftarrow \emptyset$
2. **foreach** *frame* $f_k$ *in* $v$ **do**
3.     Estimate score maps $\mathbf{z}$ from $f$ using R-FCN
   /* collect candidates */
4.     $C_{det} \leftarrow \mathcal{D}_k; C_{trk} \leftarrow \emptyset$
5.     **foreach** $t$ *in* $\mathcal{T}$ **do**
6.         Predict new location $\mathbf{x}^*$ of $t$ using Kalman filter
7.         $C_{trk} \leftarrow C_{trk} \cup \{\mathbf{x}^*\}$
8.     **end**
   /* select candidates */
9.     $C \leftarrow C_{det} \cup C_{trk}$
10.     $S \leftarrow$ unified scores computed from Equation 3
11.     $C, S \leftarrow \text{NMS}(C, S, \tau_{nms})$
12.     $C, S \leftarrow \text{Filter}(C, S, \tau_s)$ // filter out if $s < \tau_s$
    /* extract appearance features */
13.     $\mathcal{F}_{det} \leftarrow \emptyset$
14.     **foreach** $\mathbf{x}$ *in* $C_{det}$ **do**
15.         $\mathbf{I_x} \leftarrow \text{Crop}(f_k, \mathbf{x})$
16.         $\mathcal{F}_{det} \leftarrow \mathcal{F}_{det} \cup H_{reid}(\mathbf{I_x})$
17.     **end**
    /* hierarchical data association */
18.     Associate $\mathcal{T}$ and $C_{det}$ using distances of $\mathcal{F}_{trk}$ and $\mathcal{F}_{det}$
19.     Associate remaining tracks and candidates using IoU
20.     $\mathcal{F}_{trk} \leftarrow \mathcal{F}_{trk} \cup \mathcal{F}_{det}$
    /* initialize new tracks */
21.     $C_{remain} \leftarrow$ remaining candidates from $C_{det}$
22.     $\mathcal{F}_{remain} \leftarrow$ features of $C_{remain}$
23.     $\mathcal{T}, \mathcal{F}_{trk} \leftarrow \mathcal{T} \cup C_{remain}, \mathcal{F}_{trk} \cup \mathcal{F}_{remain}$
24. **end**

Object Detection

Kalman 예측 Track

(Detection + Tracking)후보집합 구성

Score 계산 -> 후보 필터링

Body Embedding

Body, IoU Matching

Tracklet Update

Tracking 결과 출력

# Method

후보 필터링

후보 x에 대한 RoI 분류확률

$$p(y|\mathbf{z}, \mathbf{x}) = \sigma\left(\frac{1}{wh} \sum_{i=1}^{k^2} \sum_{(x,y) \in bin_i} \mathbf{z}_i(x, y)\right), \qquad (1)$$
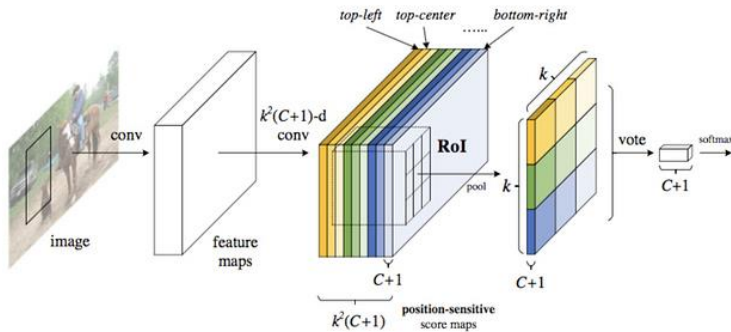


Figure 1: Key idea of **R-FCN** for object detection. In this illustration, there are $k \times k = 3 \times 3$ position-sensitive score maps generated by a fully convolutional network. For each of the $k \times k$ bins in an RoI, pooling is only performed on one of the $k^2$ maps (marked by different colors).
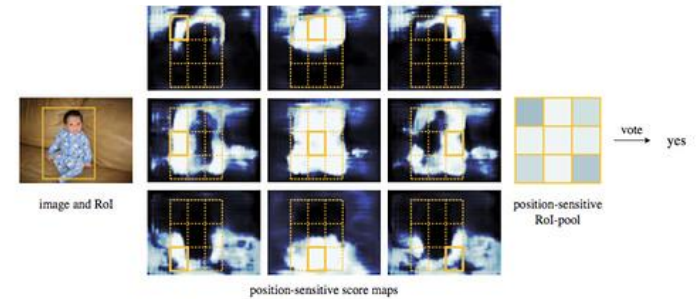


Figure 3: Visualization of R-FCN ($k \times k = 3 \times 3$) for the *person* category.



Figure 4: Visualization when an RoI does not correctly overlap the object.
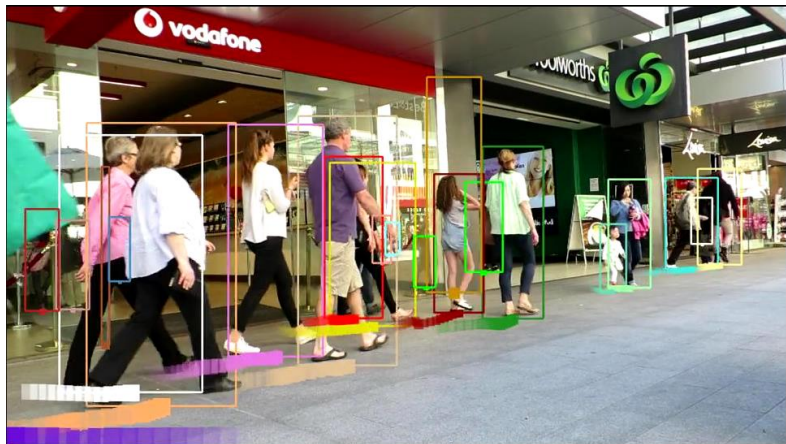
# Method

### Tracklet 신뢰도

$$s_{trk} = \max(1 - log(1 + \alpha \cdot L_{trk}), 0) \cdot \mathbb{1}(L_{det} \geq 2), \quad (2)$$

$L_{trk}$ : 최종적으로 트래킹한 결과의 수
$L_{det}$ : 최종적으로 탐지한 결과의 수

### 적어도 탐지된 결과가 2개 이상이고, 최종 트래킹 수가 적을수록 신뢰도가 높다



후보가 많을 수록 신뢰도가 낮음

후보가 적을 수록 신뢰도가 높음

# Method

후보 필터링

후보 x에 대한 score 계산

$$s = p(y|\mathbf{z}, \mathbf{x}) \cdot (\mathbb{1}(\mathbf{x} \in C_{det})) + s_{trk}\mathbb{1}(\mathbf{x} \in C_{trk})). \quad (3)$$

X : 후보
P(y|z, x) : 식 1에서 계산된 신뢰도
Strk : 식 2에서 계산된 신뢰도

R-FCN으로 계산한 점수 + Tracklet(후보) 신뢰도

만약 최종 신뢰도 s가 0.4 이하라면 후보를 버린다.

# Method

## Body Embedding

**Algorithm 1: The proposed tracking algorithm.**

**Input:** A video sequence $v$ with $N_v$ frames and object detection $\{\mathcal{D}_k\}_{k=1}^{N_v}$
**Output:** Tracks $\mathcal{T}$ of the video

1. Initialization: $\mathcal{T} \leftarrow \emptyset$; appearance of tracks $\mathcal{F}_{trk} \leftarrow \emptyset$
2. **foreach** *frame* $f_k$ *in* $v$ **do**
3.      Estimate score maps $\mathbf{z}$ from $f$ using R-FCN
     /* collect candidates */
4.      $C_{det} \leftarrow \mathcal{D}_k; C_{trk} \leftarrow \emptyset$
5.      **foreach** $t$ *in* $\mathcal{T}$ **do**
6.          Predict new location $\mathbf{x}^*$ of $t$ using Kalman filter
7.          $C_{trk} \leftarrow C_{trk} \cup \{\mathbf{x}^*\}$
8.      **end**
     /* select candidates */
9.      $C \leftarrow C_{det} \cup C_{trk}$
10.     $S \leftarrow$ unified scores computed from Equation 3
11.     $C, S \leftarrow \text{NMS}(C, S, \tau_{nms})$
12.     $C, S \leftarrow \text{Filter}(C, S, \tau_s)$ // filter out if $s < \tau_s$
     /* extract appearance features */
13.     $\mathcal{F}_{det} \leftarrow \emptyset$
14.     **foreach** $\mathbf{x}$ *in* $C_{det}$ **do**
15.         $\mathbf{I_x} \leftarrow \text{Crop}(f_k, \mathbf{x})$
16.         $\mathcal{F}_{det} \leftarrow \mathcal{F}_{det} \cup H_{reid}(\mathbf{I_x})$
17.     **end**
     /* hierarchical data association */
18.     Associate $\mathcal{T}$ and $C_{det}$ using distances of $\mathcal{F}_{trk}$ and $\mathcal{F}_{det}$
19.     Associate remaining tracks and candidates using IoU
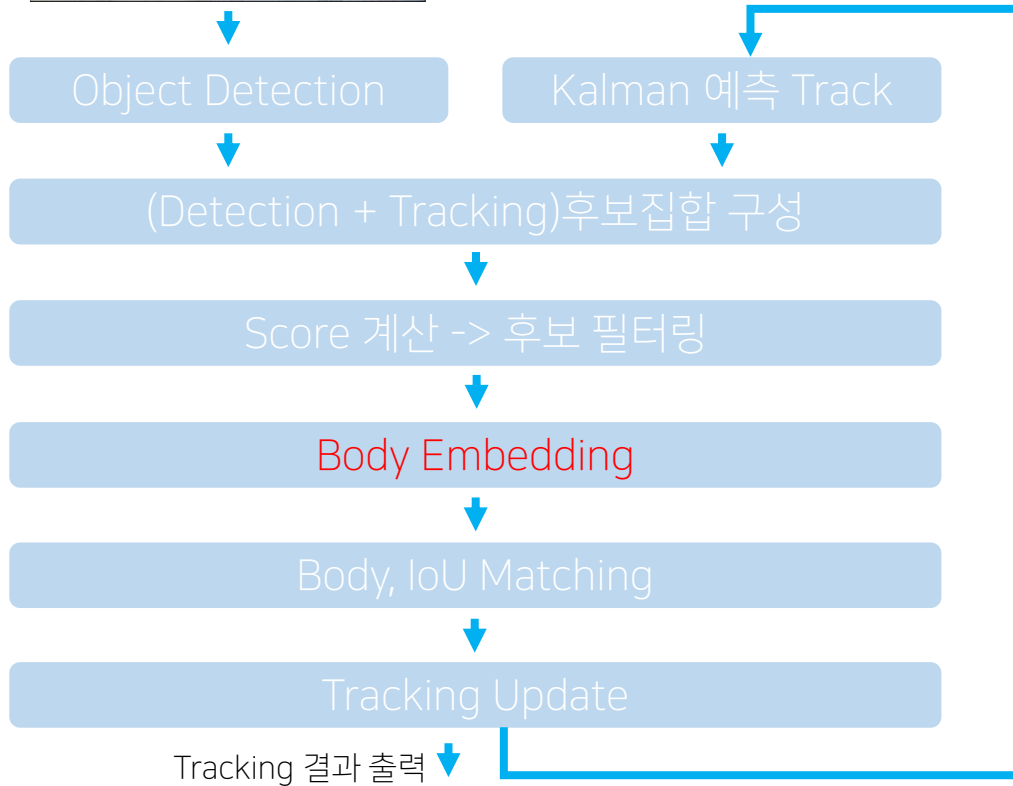20.     $\mathcal{F}_{trk} \leftarrow \mathcal{F}_{trk} \cup \mathcal{F}_{det}$
     /* initialize new tracks */
21.     $C_{remain} \leftarrow$ remaining candidates from $C_{det}$
22.     $\mathcal{F}_{remain} \leftarrow$ features of $C_{remain}$
23.     $\mathcal{T}, \mathcal{F}_{trk} \leftarrow \mathcal{T} \cup C_{remain}, \mathcal{F}_{trk} \cup \mathcal{F}_{remain}$
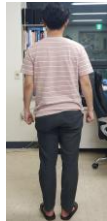24. **end**



Object Detection

Kalman 예측 Track

(Detection + Tracking)후보집합 구성

Score 계산 -> 후보 필터링

Body Embedding

Body, IoU Matching

Tracking Update

Tracking 결과 출력

# Method

Body Embedding



128d 벡터    Euclidean distance
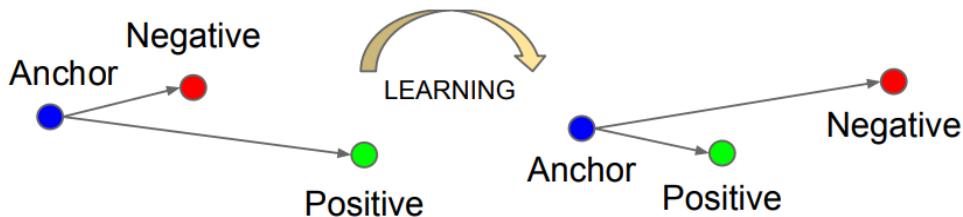
A인물의 사진1

A인물의 사진2

B인물의 사진1

0.78

1.33

1.27

# Method

Body Embedding - Triplet

$$l_{triplet} = \frac{1}{N} \sum_{\langle \mathbf{I}_i, \mathbf{I}_j, \mathbf{I}_k \rangle \in \mathbf{T}} \max(d_{ij} - d_{ik} + m, 0), \quad (4)$$

Triplet : 세 개의 데이터
- Anchor($x_i^a$) : 기준 인물의 벡터
- Positive($x_i^p$) : 기준과 같은 인물의 벡터
- Negative($x_i^n$) : 기준과 다른 인물의 벡터



기준 인물과 같은 인물은 가깝도록, 기준 인물과 다른 인물은 멀도록
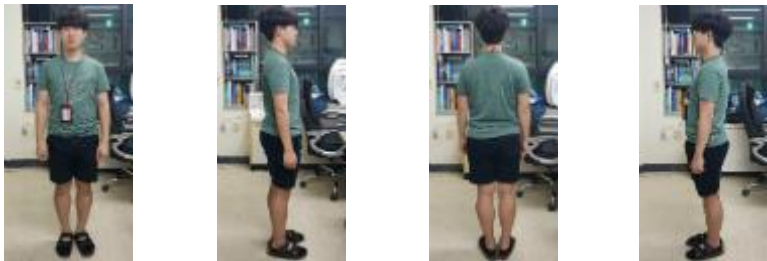
# Person Re-Identification



평균 거리

Jiseong
0    1    2    3

11.24(0제외)

gunhee
0    1    2    3

28.64

supil
0    1    2    3

27.87

```
jiseong 0, jiseong 0 : 0
jiseong 0, jiseong 1 : 14.6833314343748
jiseong 0, jiseong 2 : 7.14393362143874
jiseong 0, jiseong 3 : 13.7174603497803
jiseong 1, jiseong 0 : 14.6833314343748
jiseong 1, jiseong 1 : 0
jiseong 1, jiseong 2 : 13.3976718560914
jiseong 1, jiseong 3 : 5.4062859718527
jiseong 2, jiseong 0 : 7.14393362143874
jiseong 2, jiseong 1 : 13.3976718560914
jiseong 2, jiseong 2 : 0
jiseong 2, jiseong 3 : 13.1192839922723
jiseong 3, jiseong 0 : 13.7174603497803
jiseong 3, jiseong 1 : 5.4062859718527
jiseong 3, jiseong 2 : 13.1192839922723
jiseong 3, jiseong 3 : 0
평균거리 : 11.2446612043017

jiseong 0, gunhee0 : 33.702875406024
jiseong 0, gunhee1 : 30.6240802173837
jiseong 0, gunhee2 : 34.6539498741167
jiseong 0, gunhee3 : 30.0388592472747
jiseong 1, gunhee0 : 28.4433327086116
jiseong 1, gunhee1 : 24.1252707695199
jiseong 1, gunhee2 : 29.460226635402
jiseong 1, gunhee3 : 23.97228405069
jiseong 2, gunhee0 : 31.3361323432144
jiseong 2, gunhee1 : 28.299821729878
jiseong 2, gunhee2 : 32.1328332890316
jiseong 2, gunhee3 : 27.8813871223626
jiseong 3, gunhee0 : 27.7410707429242
jiseong 3, gunhee1 : 23.9328520321636
jiseong 3, gunhee2 : 28.771589169261
jiseong 3, gunhee3 : 23.1914684226621
평균거리 : 28.6442521100325

jiseong 0, supil0 : 32.5034984023311
jiseong 0, supil1 : 31.0123639922488
jiseong 0, supil2 : 33.1022183632668
jiseong 0, supil3 : 28.2895292724031
jiseong 1, supil0 : 26.9324479963173
jiseong 1, supil1 : 25.3354276975164
jiseong 1, supil2 : 28.4278065042272
jiseong 1, supil3 : 22.5782505170552
jiseong 2, supil0 : 31.1035040681003
jiseong 2, supil1 : 30.1407394004526
jiseong 2, supil2 : 31.6739051391745
jiseong 2, supil3 : 27.3928874045504
jiseong 3, supil0 : 25.716432805367
jiseong 3, supil1 : 23.8073005976364
jiseong 3, supil2 : 27.0416266409934
jiseong 3, supil3 : 20.8900794862151
평균거리 : 27.871751142991
```

# Method

Matching



**Algorithm 1: The proposed tracking algorithm.**

**Input:** A video sequence $v$ with $N_v$ frames and object detection $\{\mathcal{D}_k\}_{k=1}^{N_v}$
**Output:** Tracks $\mathcal{T}$ of the video

1  Initialization: $\mathcal{T} \leftarrow \emptyset$; appearance of tracks $\mathcal{F}_{trk} \leftarrow \emptyset$
2  **foreach** *frame* $f_k$ *in* $v$ **do**
3      Estimate score maps $\mathbf{z}$ from $f$ using R-FCN
    /* collect candidates */
4      $C_{det} \leftarrow \mathcal{D}_k$; $C_{trk} \leftarrow \emptyset$
5      **foreach** $t$ *in* $\mathcal{T}$ **do**
6          Predict new location $\mathbf{x}^*$ of $t$ using Kalman filter
7          $C_{trk} \leftarrow C_{trk} \cup \{\mathbf{x}^*\}$
8      **end**
    /* select candidates */
9      $C \leftarrow C_{det} \cup C_{trk}$
10     $S \leftarrow$ unified scores computed from Equation 3
11     $C, S \leftarrow \text{NMS}(C, S, \tau_{nms})$
12     $C, S \leftarrow \text{Filter}(C, S, \tau_s)$ // filter out if $s < \tau_s$
    /* extract appearance features */
13     $\mathcal{F}_{det} \leftarrow \emptyset$
14     **foreach** $\mathbf{x}$ *in* $C_{det}$ **do**
15         $\mathbf{I_x} \leftarrow \text{Crop}(f_k, \mathbf{x})$
16         $\mathcal{F}_{det} \leftarrow \mathcal{F}_{det} \cup H_{reid}(\mathbf{I_x})$
17     **end**
    /* hierarchical data association */
18     Associate $\mathcal{T}$ and $C_{det}$ using distances of $\mathcal{F}_{trk}$ and $\mathcal{F}_{det}$
19     Associate remaining tracks and candidates using IoU
20     $\mathcal{F}_{trk} \leftarrow \mathcal{F}_{trk} \cup \mathcal{F}_{det}$
    /* initialize new tracks */
21     $C_{remain} \leftarrow$ remaining candidates from $C_{det}$
22     $\mathcal{F}_{remain} \leftarrow$ features of $C_{remain}$
23     $\mathcal{T}, \mathcal{F}_{trk} \leftarrow \mathcal{T} \cup C_{remain}, \mathcal{F}_{trk} \cup \mathcal{F}_{remain}$
24 **end**

Object Detection

Kalman 예측 Track

(Detection + Tracking)후보집합 구성

Score 계산 -> 후보 필터링

Body Embedding

Body, IoU Matching

Tracking Update

Tracking 결과 출력

# Method

1. 후보와 Active tracklet을 Body Feture를 이용하여 매칭



이전 프레임        현재 프레임

☐ Active Tracklet
☐ 후보

Body Feature Threshold : 0.64
IoU Threshold : 0.4

# Method

2. 후보와 Lost Tracklet과 Body Feature를 이용하여 매칭

이전 프레임 1 　　　　　 이전 프레임 2 　　　　　　　 현재 프레임



Tracking Object를 잃어버림

☐ Active Tracklet
☐ Lost Tracklet
☐ 후보

Body Feature Threshold : 0.64
IoU Threshold : 0.4

# Method

3. 후보와 Active tracklet을 IoU를 이용하여 매칭



이전 프레임

현재 프레임

□ Active Tracklet
□ 후보

IoU Threshold : 0.7

# Method

4. 후보와 Lost Tracklet과 IoU를 이용하여 매칭

| 이전 프레임 1 | 이전 프레임 2 | 현재 프레임 |
|---|---|---|



Tracking Object를 잃어버림

□ Active Tracklet
□ Lost Tracklet
□ 후보

IoU Threshold : 0.7

# Method

MOTA


IoU = Area of Overlap / Area of Union

Ground Truth와 관심영역의 IOU가 0.5 이상일 때 True로 판단

MOTA : tracker 성능지표로 적합

$$\text{MOTA} = 1 - \frac{\sum_t (\text{FN}_t + \text{FP}_t + \text{IDSW}_t)}{\sum_t \text{GT}_t}, \qquad (1)$$

t : 프레임 인덱스
FN : 잘못탐지
FP : 놓친 객체
IDSW : id가 바뀐 횟수
GT : Ground Truth의 수

# Method

IDF1 : 얼마나 지속적으로 객체를 동일하게 판단하는가

True
Computed
time

- ID Precision $\quad P = \frac{TP}{TP+FP} = \frac{TP}{C}$

- ID Recall $\quad R = \frac{TP}{TP+FN} = \frac{TP}{T}$

- $F_1$-score $\quad F_1 = 2\frac{PR}{P+R} = \frac{TP}{\frac{T+C}{2}}$

TP : 가장 많이 나온 ID의 개수 = 12
FP : 나머지 ID의 개수 = 2
FN : True – TP = 4

False
Negatives
True
Positives
False
Positives

True
Detections
Computed
Detections

# Vision Module

## MOT Challenge

| Method | C | T | A | MOTA↑ | IDF1↑ | IDS↓ | FAF↓ |
|--------|---|---|---|-------|-------|------|------|
| Baseline | | | | 28.4 | 32.8 | 628 | 0.85 |
| | ✓ | | | 33.0 | 37.6 | 445 | 0.77 |
| | ✓ | ✓ | | 33.7 | 37.3 | 475 | 0.63 |
| | | | ✓ | 30.6 | 42.4 | 234 | 1.01 |
| Proposed | ✓ | ✓ | ✓ | **35.7** | **45.3** | **184** | **0.58** |

| Method | Length | MOTA↑ | IDF1↑ | IDS↓ | FAF↓ |
|--------|--------|-------|-------|------|------|
| None | - | 33.7 | 37.3 | 475 | 0.63 |
| Color histogram | 750 | 34.9 | 38.6 | 250 | 0.73 |
| HOG | 1152 | 34.6 | 38.5 | 317 | 0.70 |
| Color + HOG | 1902 | 34.7 | 39.3 | 307 | 0.68 |
| ReID feature | **512** | **35.7** | **45.3** | **184** | **0.58** |

| Tracker | Method | MOTA(%)↑ | IDF1(%)↑ | IDR(%)↑ | MT(%)↑ | ML(%)↓ | FP↓ | FN↓ | IDS↓ | FPS↑ |
|---------|--------|----------|----------|---------|--------|--------|-----|-----|------|------|
| LINF1 [2] | batch | 41.0 | 45.7 | 34.2 | 11.6 | 51.3 | 7,896 | 99,224 | **430** | **4.2** |
| MHT_DAM [5] | batch | 45.8 | 46.1 | 35.3 | 16.2 | 43.2 | 6,412 | 91,758 | 590 | 0.8 |
| JMC [11] | batch | 46.3 | 46.3 | 35.6 | 15.5 | **39.7** | **6,373** | 90,914 | 657 | 0.8 |
| LMP [6] | batch | **48.8** | **51.3** | **40.1** | **18.2** | 40.1 | 6,654 | **86,245** | 481 | 0.5 |
| EAMTT [7] | online | 38.8 | 42.4 | 31.5 | 7.9 | 49.1 | 8,114 | 102,452 | 965 | 11.8 |
| CDA_DDAL [1] | online | 43.9 | 45.1 | 34.1 | 10.7 | 44.4 | 6,450 | 95,175 | 676 | 0.5 |
| STAM [8] | online | 46.0 | 50.0 | 38.5 | 14.6 | 43.6 | 6,895 | 91,117 | **473** | 0.2 |
| AMIR [3] | online | 47.2 | 46.3 | 34.8 | 14.0 | 41.6 | **2,681** | 92,856 | 774 | 1.0 |
| **MOTDT (Ours)** | online | **47.6** | **50.9** | **40.3** | **15.2** | **38.3** | 9,253 | **85,431** | 792 | **20.6** |

감사합니다.