

# PPDM: Parallel Point Detection and Matching for Real-time Human-Object Interaction Detection

Yue Liao, Si Liu, Fei Wang, Yanjie Chen, Chen Qian, Jiashi Feng  
School of Computer Science and Engineering, Beihang University,  
SenseTime Research, National University of Singapore

2020 CVPR

인공지능 연구실  
석사과정 구자봉

# 문제 정의 :

## HOI(Human Object Interaction)

이미지에서 오브젝트 디텍션, 인간과 상호작용이 큰 객체쌍을 선택, 술어(상관관계)를 찾는 것이 목적

### Instance Detection



### Interaction Inference



(a)

# 문제 제기 :

대부분의 HOI 모듈은

1)인간-대상 제안서 생성, 2)제안서 분류

2단계로 나뉘어 구성된다. 순차적이고 개별적인 아키텍처를 사용하여 성능이 제한된다.

## Instance Detection



## Interaction Inference

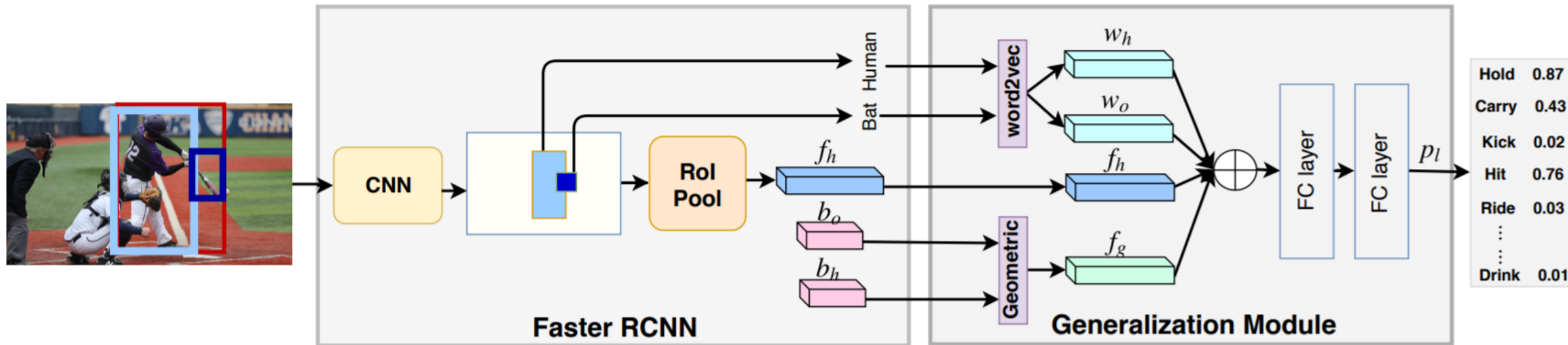


(a)

# Detecting Human-Object Interactions via Functional Generalization

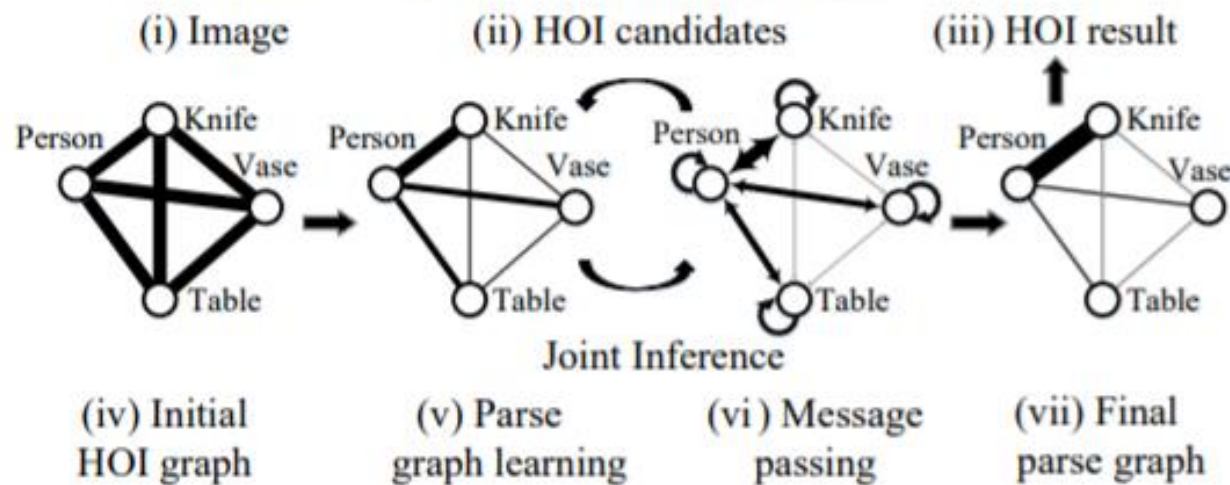
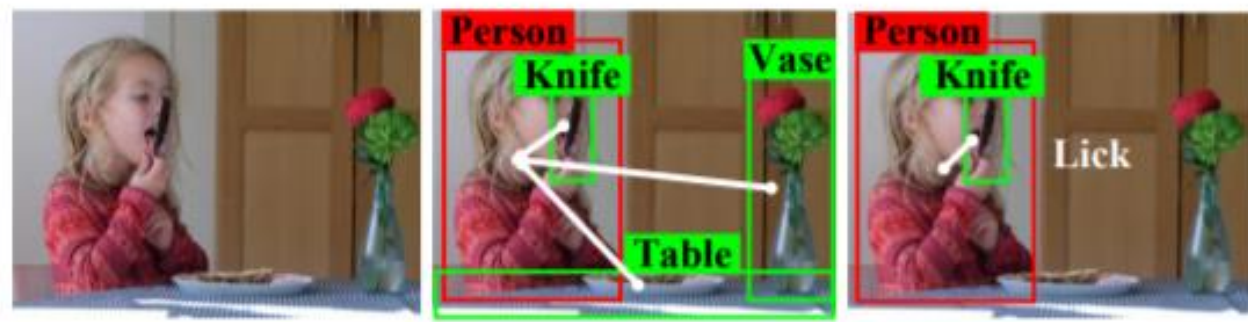


(human, eat, ...)

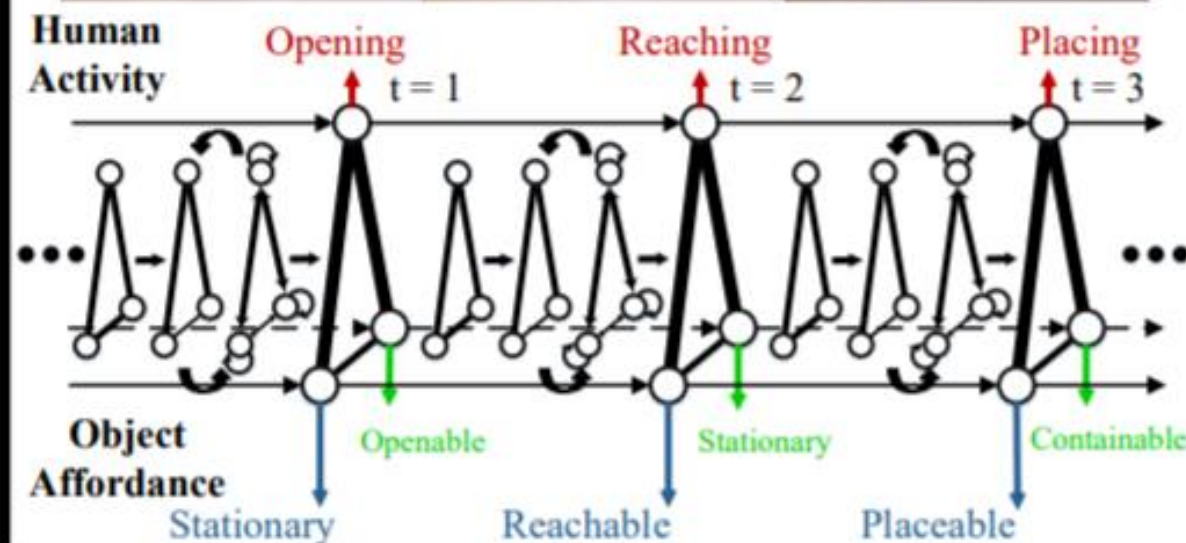
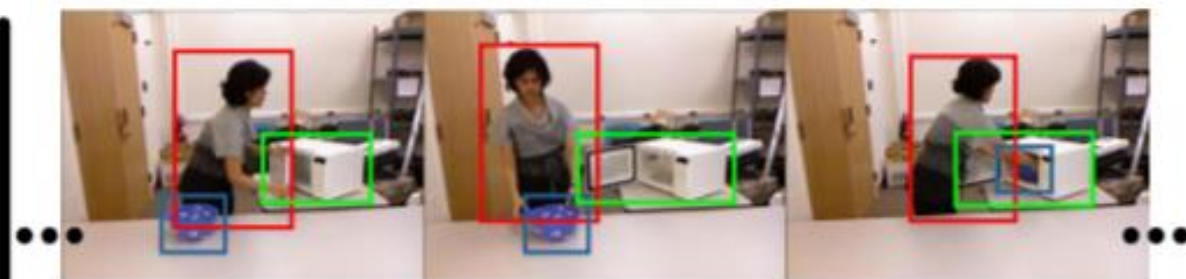




# Learning Human-Object Interactions by Graph Parsing Neural Networks (GPNN):



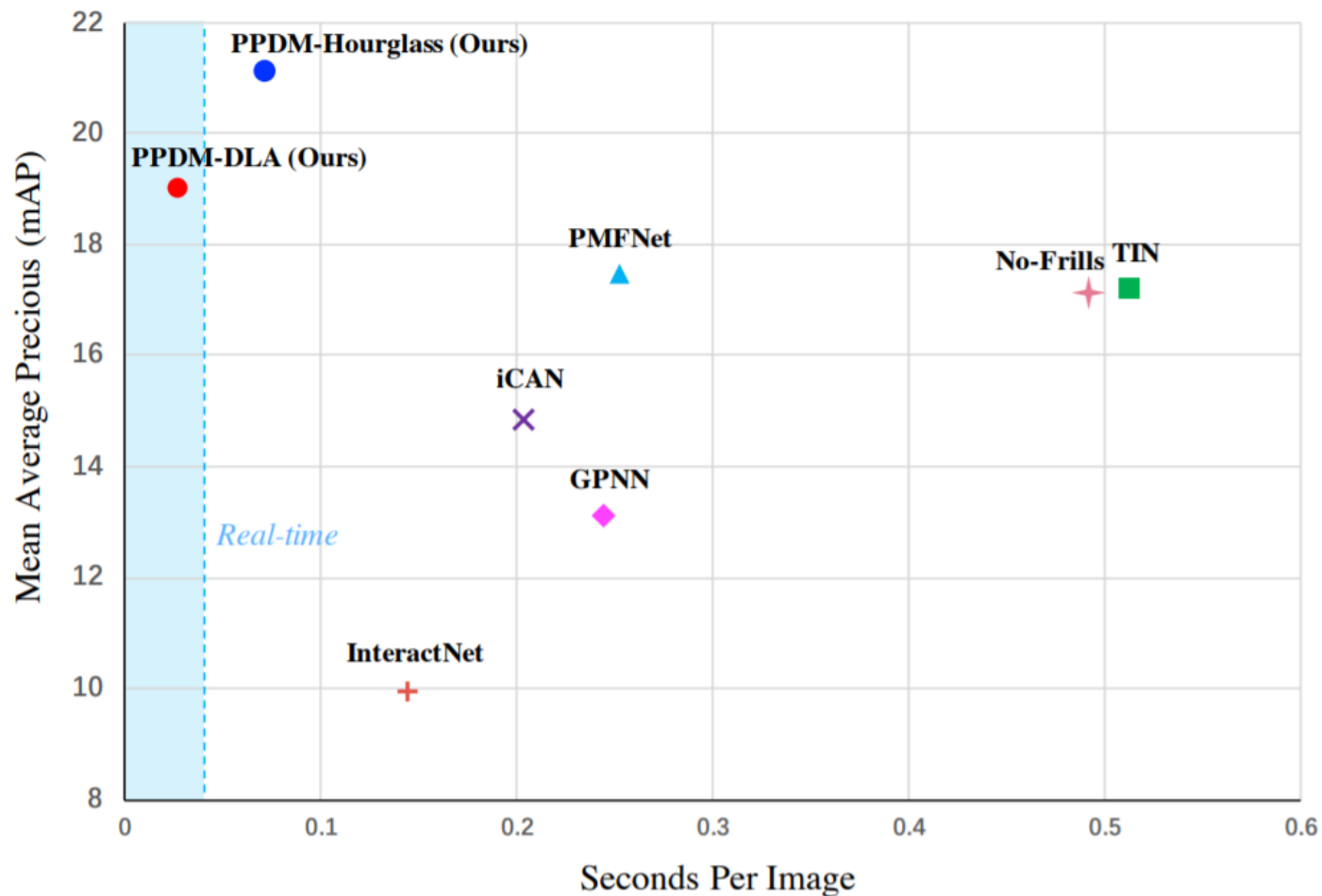
(a) Human-Object Interaction Detection in Still Images



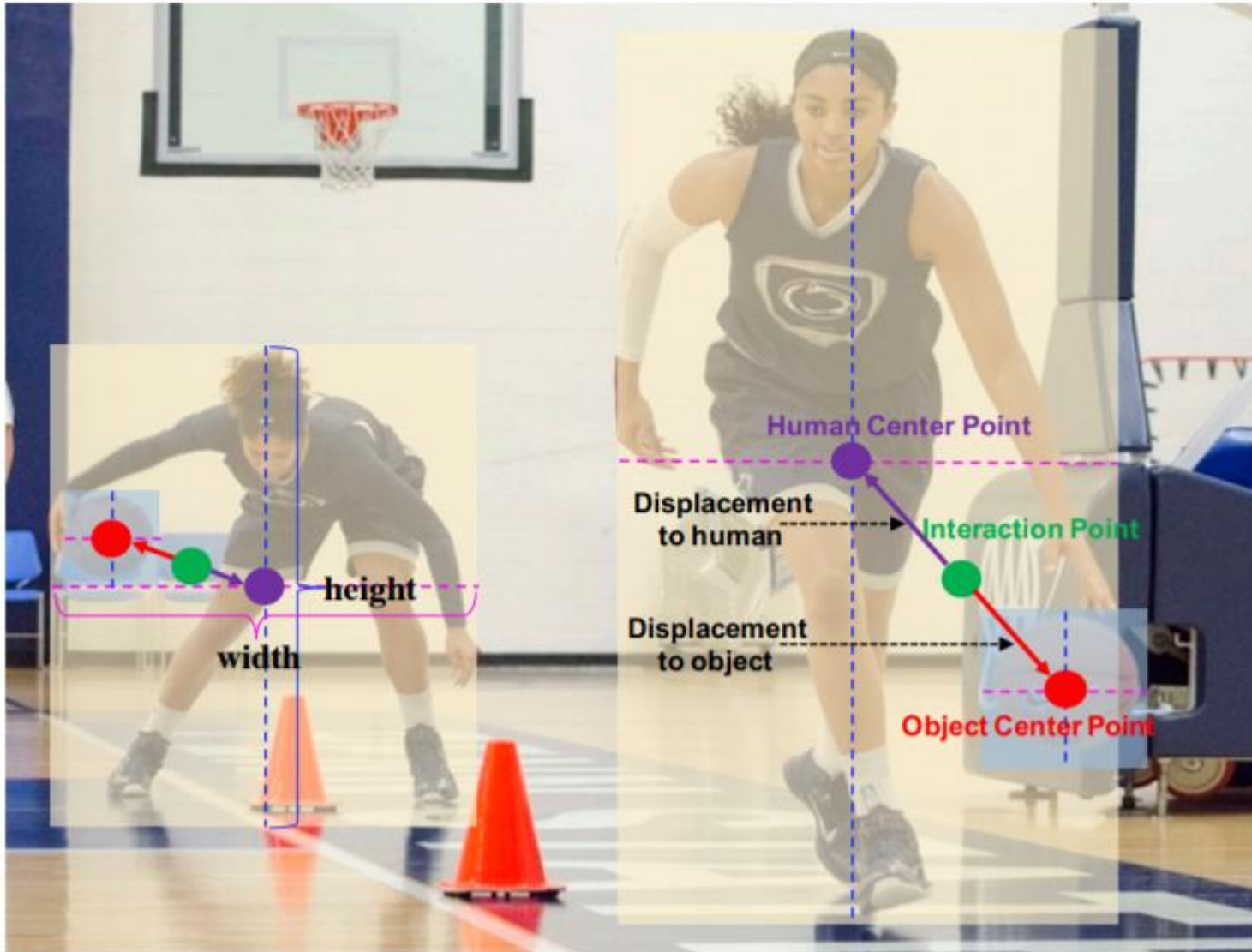
(b) Human-Object Interaction Recognition in Videos

# 목표:

HOI 문제를 **점 삼중항으로 정의**하여 포인트 감지분기와 포인트 일치 분기를 병렬 아키텍처를 사용하여 쓸모없는 삼중항을 억제함으로써 정확도를 높임과 동시에 **계산 비용**을 절감한다.



# 메인 아이디어



두개의 병렬 분기

1) 포인트 검출

- 사람 중심점, 객체 중심점, 상호작용 중심점 찾기

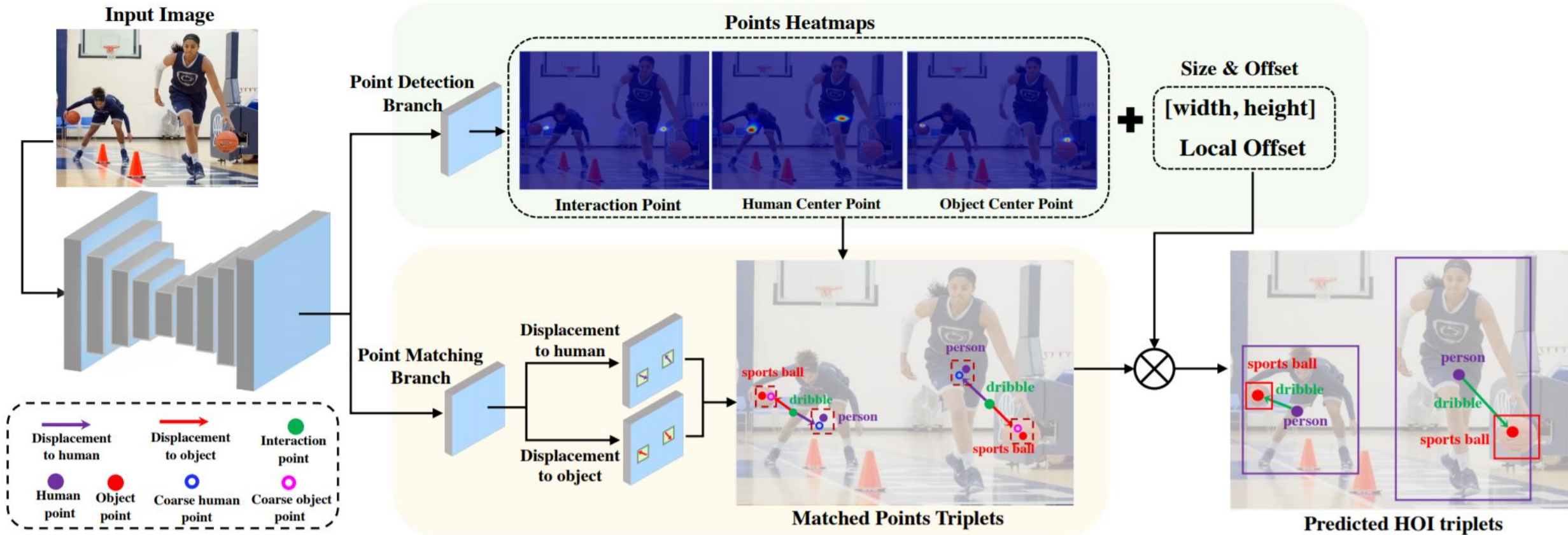
2) 포인트 매칭

- 상호작용 중심점에서 인간과 물체로의 변위 추정

➤ 동일한 상호작용 중심점에서의 변위가 일치하면 상호작용 쌍으로 간주함



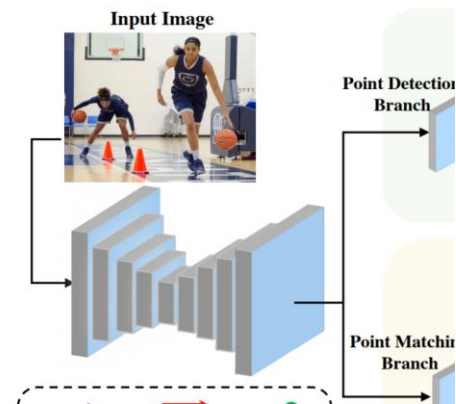
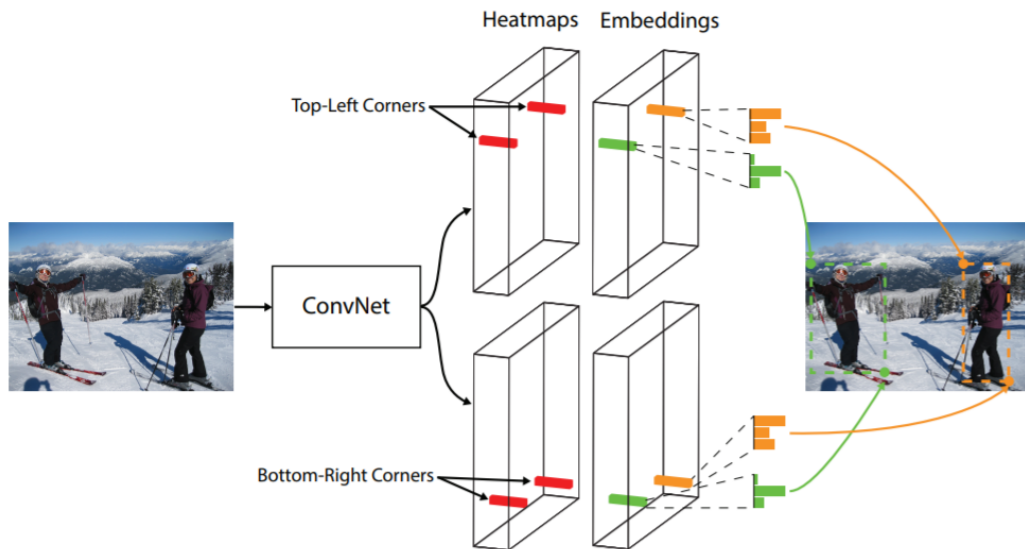
# PPDM : Parallel Point Detection and Matching for Real-time Human-Object Interaction Detection





# 피쳐 추출기(heatmap prediction networks)

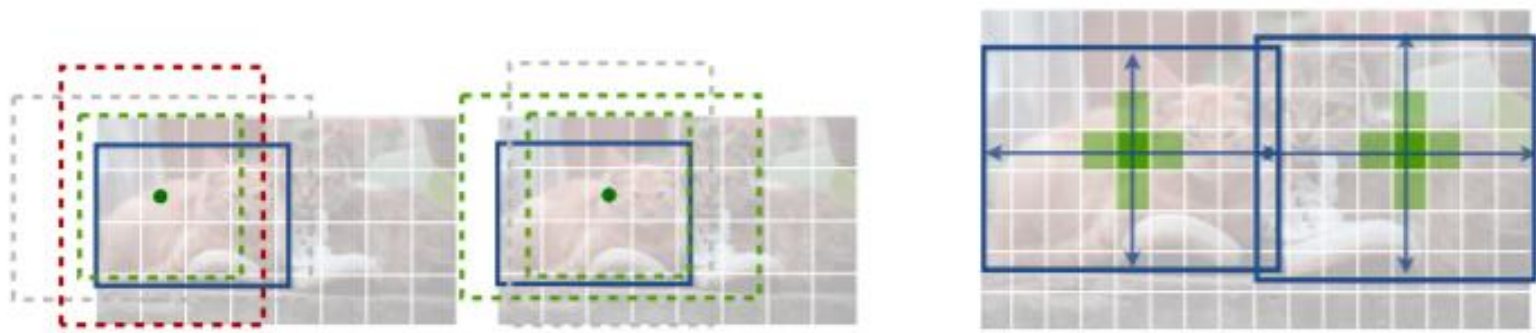
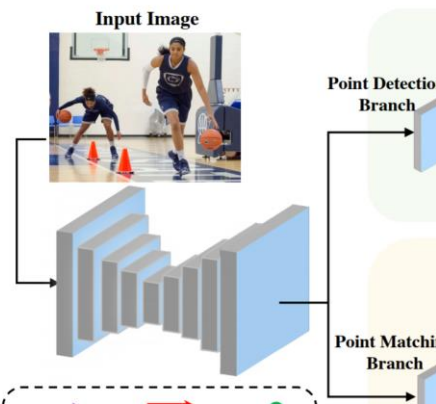
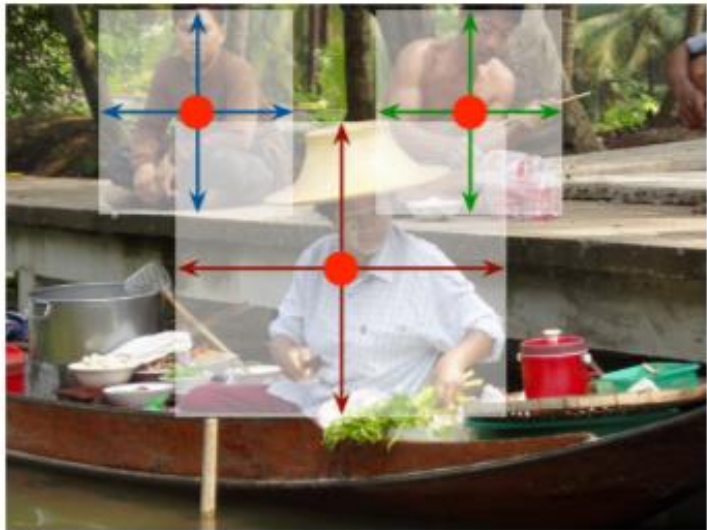
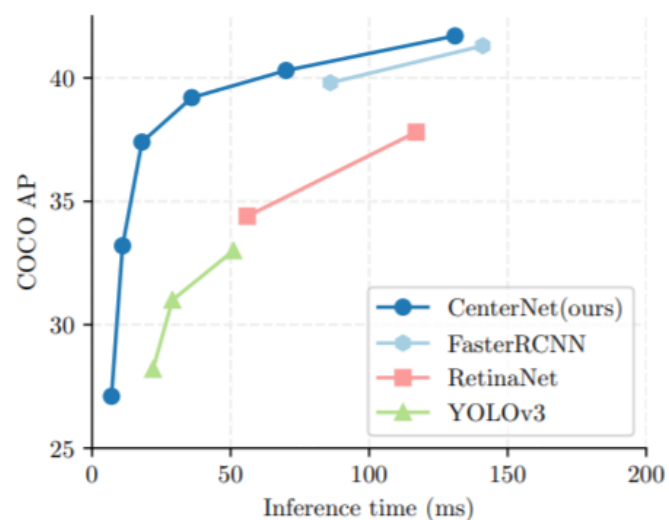
## Hourglass-104(CornerNet)



Method	Backbone	AP
<b>Two-stage detectors</b>		
DeNet (Tychsen-Smith and Petersson, 2017a)	ResNet-101	33.8
CoupleNet (Zhu et al., 2017)	ResNet-101	34.4
Faster R-CNN by G-RMI (Huang et al., 2017)	Inception-ResNet-v2 (Szegedy et al., 2017)	34.7
Faster R-CNN+++ (He et al., 2016)	ResNet-101	34.9
Faster R-CNN w/ FPN (Lin et al., 2016)	ResNet-101	36.2
Faster R-CNN w/ TDM (Shrivastava et al., 2016)	Inception-ResNet-v2	36.8
D-FCN (Dai et al., 2017)	Aligned-Inception-ResNet	37.5
Regionlets (Xu et al., 2017)	ResNet-101	39.3
Mask R-CNN (He et al., 2017)	ResNeXt-101	39.8
Soft-NMS (Bodla et al., 2017)	Aligned-Inception-ResNet	40.9
LH R-CNN (Li et al., 2017)	ResNet-101	41.5
Fitness-NMS (Tychsen-Smith and Petersson, 2017b)	ResNet-101	41.8
Cascade R-CNN (Cai and Vasconcelos, 2017)	ResNet-101	42.8
D-RFCN + SNIP (Singh and Davis, 2017)	DPN-98 (Chen et al., 2017)	45.7
<b>One-stage detectors</b>		
YOLOv2 (Redmon and Farhadi, 2016)	DarkNet-19	21.6
DSOD300 (Shen et al., 2017a)	DS/64-192-48-1	29.3
GRP-DSOD320 (Shen et al., 2017b)	DS/64-192-48-1	30.0
SSD513 (Liu et al., 2016)	ResNet-101	31.2
DSSD513 (Fu et al., 2017)	ResNet-101	33.2
RefineDet512 (single scale) (Zhang et al., 2017)	ResNet-101	36.4
RetinaNet800 (Lin et al., 2017)	ResNet-101	39.1
RefineDet512 (multi scale) (Zhang et al., 2017)	ResNet-101	41.8
CornerNet511 (single scale)	Hourglass-104	40.6
CornerNet511 (multi scale)	Hourglass-104	42.2

# 피쳐 추출기(heatmap prediction networks)

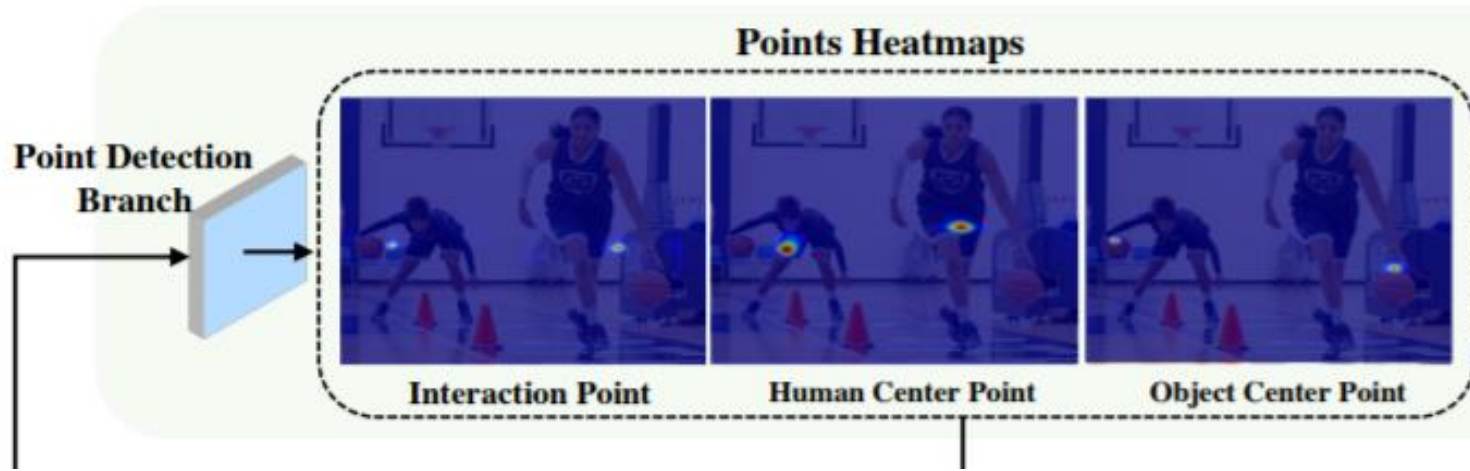
DLA-34



	AP			Time (ms)			FPS		
	N.A.	F	MS	N.A.	F	MS	N.A.	F	MS
Hourglass-104	40.3	42.2	45.1	71	129	672	14	7.8	1.4
DLA-34	37.4	39.2	41.7	19	36	248	52	28	4
ResNet-101	34.6	36.2	39.3	22	40	259	45	25	4
ResNet-18	28.1	30.0	33.2	7	14	81	142	71	12

# 포인트 디텍션

사람 상자, 객체 상자, 상호작용 지점을 추정(가우시안 커널을 사용하여 히트맵으로 표시)



$$(x^h, \hat{y}^h) \in \mathbb{R}^2$$

$$(w^h, h^h) \in \mathbb{R}^2$$

$$\delta c^h \in \mathbb{R}^2$$

$$(x^h, y^h)$$

$$\tilde{C}^h \in [0, 1]^{\frac{H}{d} \times \frac{W}{d}}$$

$$(x^a, y^a) \in \mathbb{R}^2 \quad (\tilde{x}^a, \tilde{y}^a) = (\lfloor \frac{\tilde{x}^h + \tilde{x}^o}{2} \rfloor, \lfloor \frac{\tilde{y}^h + \tilde{y}^o}{2} \rfloor).$$

$$(x^o, y^o)$$

$$\tilde{C}^o \in [0, 1]^{T \times \frac{H}{d} \times \frac{W}{d}}$$

$$(x^a, y^a)$$

$$\tilde{C}^a \in [0, 1]^{K \times \frac{H}{d} \times \frac{W}{d}}$$



# 포인트 디텍션

로스, 알파2, 베타4, 사람과 오브젝트도 동일한 식 사용하여 로스 적용

$$L_a = -\frac{1}{N} \sum_{kxy} \begin{cases} (1 - \hat{C}_{kxy}^a)^\alpha \log(\hat{C}_{kxy}^a) & \text{if } \tilde{C}_{kxy}^a = 1 \\ (1 - \tilde{C}_{kxy}^a)^\beta (\hat{C}_{kxy}^a)^\alpha & \text{otherwise} \\ \log(1 - \hat{C}_{kxy}^a), \end{cases} \quad (1)$$

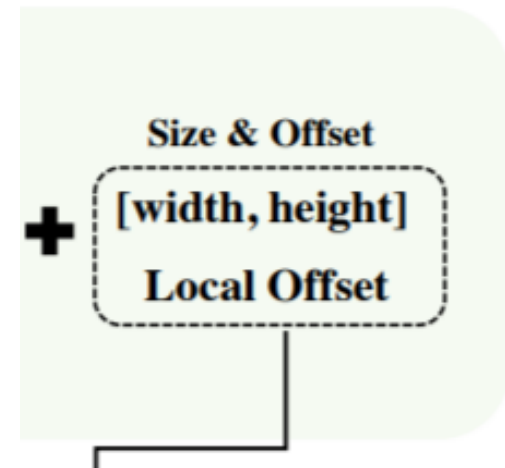
크기 및 오프셋 손실

$$(\tilde{\delta}_{(\tilde{x}^h, \tilde{y}^h)}^x, \tilde{\delta}_{(\tilde{x}^h, \tilde{y}^h)}^y) = \left( \frac{x^h}{d} - \tilde{x}^h, \frac{y^h}{d} - \tilde{y}^h \right).$$

$$L_{off}^h = \sum_{(\tilde{x}^h, \tilde{y}^h) \in \tilde{S}^h} (|\tilde{\delta}_{(\tilde{x}^h, \tilde{y}^h)}^x - \hat{\delta}_{(\tilde{x}^h, \tilde{y}^h)}^x| + |\tilde{\delta}_{(\tilde{x}^h, \tilde{y}^h)}^y - \hat{\delta}_{(\tilde{x}^h, \tilde{y}^h)}^y|), \quad (3)$$

$$L_{off} = \frac{1}{M + D} (L_{off}^h + L_{off}^o)$$

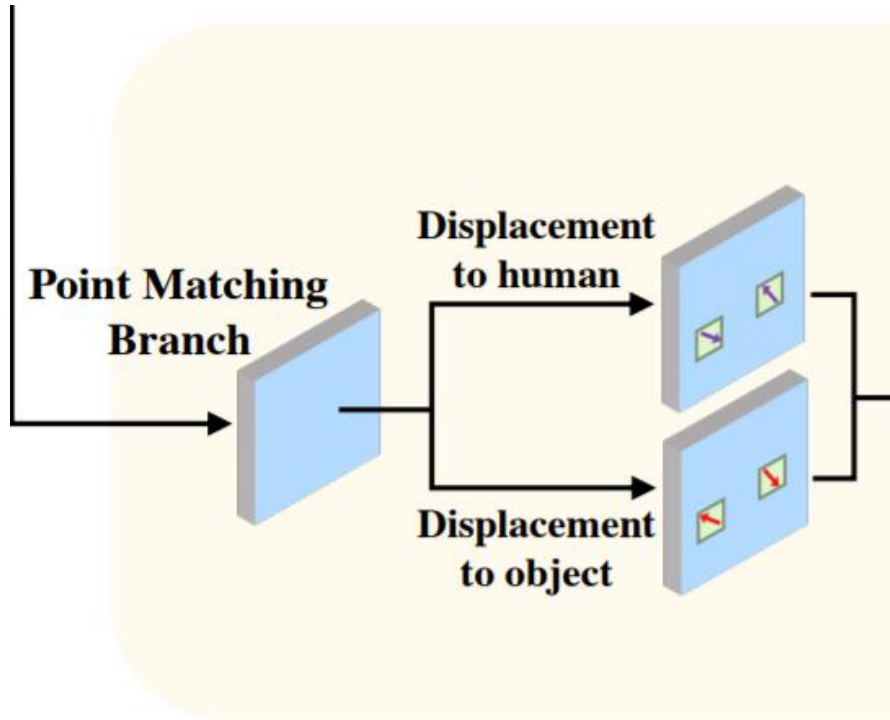
$$M = |\tilde{S}^h| \text{ and } D = |\tilde{S}^o|$$





# 포인트 매칭

상호작용 점을 앵커로 하여 사람 쪽과 객체 쪽으로 변위를 추정



로스  $\sim$ (그라운드트루),  $\wedge$ (예측된),  $L_{ao}$

$$(\tilde{d}_{(\tilde{x}^a, \tilde{y}^a)}^{hx}, \tilde{d}_{(\tilde{x}^a, \tilde{y}^a)}^{hy}) = (\tilde{x}^a - \tilde{x}^h, \tilde{y}^a - \tilde{y}^h)$$

$$L_{ah} = \frac{1}{N} \sum_{(\tilde{x}^a, \tilde{y}^a) \in \tilde{S}^a} |\hat{d}_{(\tilde{x}^a, \tilde{y}^a)}^{hx} - \tilde{d}_{(\tilde{x}^a, \tilde{y}^a)}^{hx}| + |\hat{d}_{(\tilde{x}^a, \tilde{y}^a)}^{hy} - \tilde{d}_{(\tilde{x}^a, \tilde{y}^a)}^{hy}|, \quad (4)$$

$$d^{ah} = (d_x^{ah}, d_y^{ah})$$

$$d^{ao} = (d_x^{ao}, d_y^{ao})$$

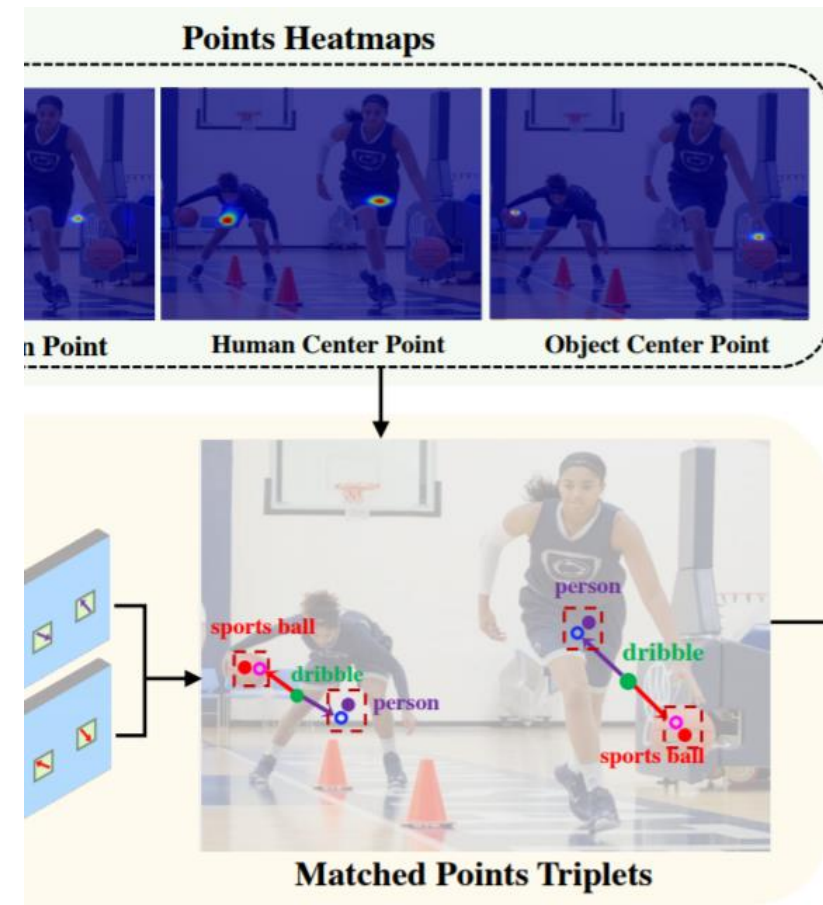
# 포인트 매칭

트리플넷 매칭

가우시안 분포인 C점수가 크고, 예측된 값의 차이가 가까운 것들중 최적의 것을 선택  
객체상자도 같은 식으로 선택

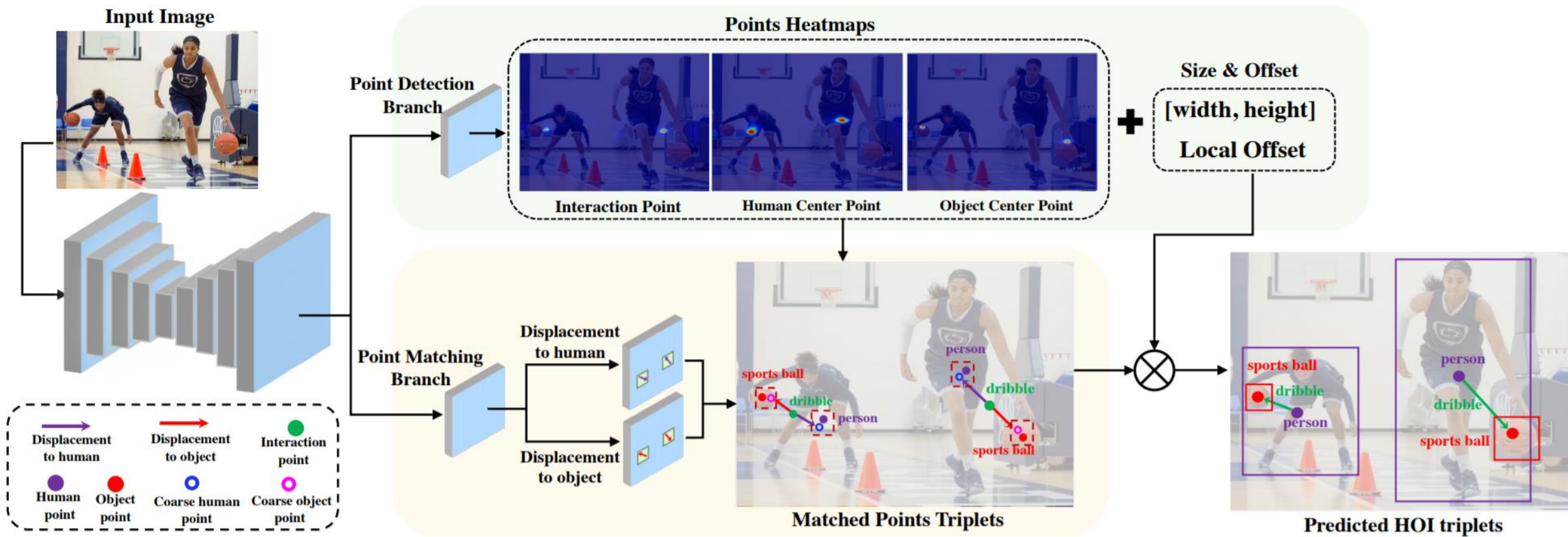
$$(\hat{x}_{opt}^h, \hat{y}_{opt}^h) = \arg \min_{(\hat{x}^h, \hat{y}^h) \in \hat{S}^h} \frac{1}{C_{(\hat{x}^h, \hat{y}^h)}^h}$$

$$(|(\hat{x}^a, \hat{y}^a) - (d_{(\hat{x}^a, \hat{y}^a)}^{hx}, d_{(\hat{x}^a, \hat{y}^a)}^{hy}) - (\hat{x}^h, \hat{y}^h)|) \quad (5)$$



# 최종 손실

$$\begin{aligned}\hat{x}_{ref}^h &= \hat{x}_{opt}^h + \hat{\delta}_{(\hat{x}_{opt}^h, \hat{y}_{opt}^h)}^x \left( \hat{x}_{ref}^h - \frac{\hat{w}_{(\hat{x}_{opt}^h, \hat{y}_{opt}^h)}}{2}, \hat{y}_{ref}^h - \frac{\hat{h}_{(\hat{x}_{opt}^h, \hat{y}_{opt}^h)}}{2}, \right. \\ \hat{y}_{ref}^h &= \hat{y}_{opt}^h + \hat{\delta}_{(\hat{x}_{opt}^h, \hat{y}_{opt}^h)}^y \left. \hat{x}_{ref}^h + \frac{\hat{w}_{(\hat{x}_{opt}^h, \hat{y}_{opt}^h)}}{2}, \hat{y}_{ref}^h + \frac{\hat{h}_{(\hat{x}_{opt}^h, \hat{y}_{opt}^h)}}{2} \right).\end{aligned}\quad (7)$$



$$L = L_a + L_h + L_o + \lambda(L_{ah} + L_{ao} + L_{wh}) + L_{off} \quad (6)$$

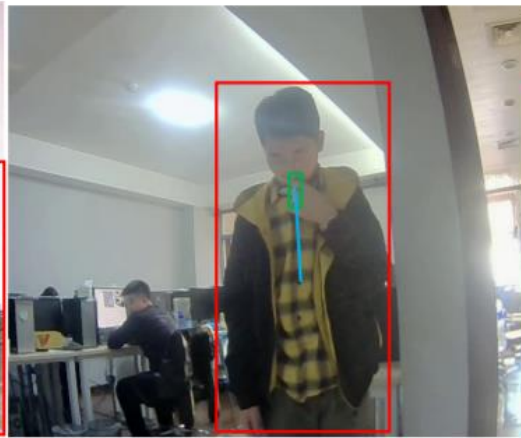
$$\hat{C}_{\hat{x}_{ref}^h \hat{y}_{ref}^h}^p \hat{C}_{\hat{x}_{ref}^o \hat{y}_{ref}^o}^o \hat{C}_{\hat{x}_{ref}^a \hat{y}_{ref}^a}^a.$$



# HOI-A Dataset(Human-Object Interaction for Application)



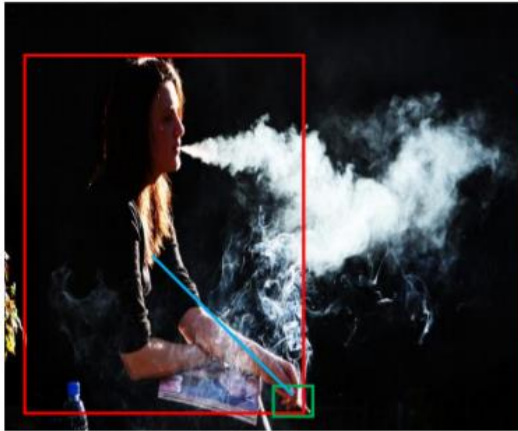
a. <human, smoke, cigarette>  
outdoor



b. <human, smoke, cigarette>  
indoor



c. <human, smoke, cigarette>  
in car & intense illumination



d. <human, smoke, cigarette>  
in dark scene



e. Attacking smoke: no cigarette  
negative sample



f. no predefined interaction  
negative sample

HOI를 위한 긍정적인 이미지와 부정적인 이미지를 모아 훈련할 때 더 효과적으로 할 수 있도록 하는 데이터셋  
38,668 이미지, 11개 객체, 10개 액션

43,820 휴먼 등장, 60,438개 오브젝트 등장, 96,160개 상호작용 등장



# Dataset

HICO-DET

Images 47,776 (38,118, 9,658)

Objects 80 (airplane, apple...)

Verbs 117 (carry, catch...)

HOI 600 (airplane – board, direct, exit, fly...)

HOI Remark  $\geq 150k$

HOI-A

## 평가지표

mAP

# 실험 (HICO-DET)

Method	Feature	Full(mAP %) ↑	Rare(mAP %) ↑	Non-Rare(mAP %) ↑	Inference Time (ms) ↓	FPS ↑
Shen <i>et. al</i> [23]	A + P	6.46	4.24	7.12	-	-
HO-RCNN [2]	A + S	7.81	5.37	8.54	-	-
VSRL [10]	A	9.09	7.02	9.71	-	-
InteractNet [8]	A	9.94	7.16	10.77	145	6.90
GPNN [21]	A	13.11	9.34	14.23	197 + 48 = 245	4.08
Xu <i>et. al</i> [27]	A + L	14.70	13.26	15.13	-	-
iCAN [6]	A + S	14.84	10.45	16.15	92 + 112 = 204	4.90
PMFNet-Base [25]	A + S	14.92	11.42	15.96	-	-
Wang <i>et. al</i> [26]	A	16.24	11.16	17.75	-	-
No-Frills [11]	A + S + P	17.18	12.17	18.68	197 + 230 + 67 = 494	2.02
TIN [15]	A + S + P	17.22	13.51	18.32	92 + 98 + 323 = 513	1.95
RPNN [32]	A + P	17.35	12.78	18.71	-	-
PMFNet [25]	A + S + P	17.46	<b>15.65</b>	18.00	92 + 98 + 63 = 253	3.95
PPDM-DLA	A	19.02	12.65	20.92	<b>27</b>	<b>37.03</b>
PPDM-Hourglass	A	<b>21.10</b>	14.46	<b>23.09</b>	71	14.08

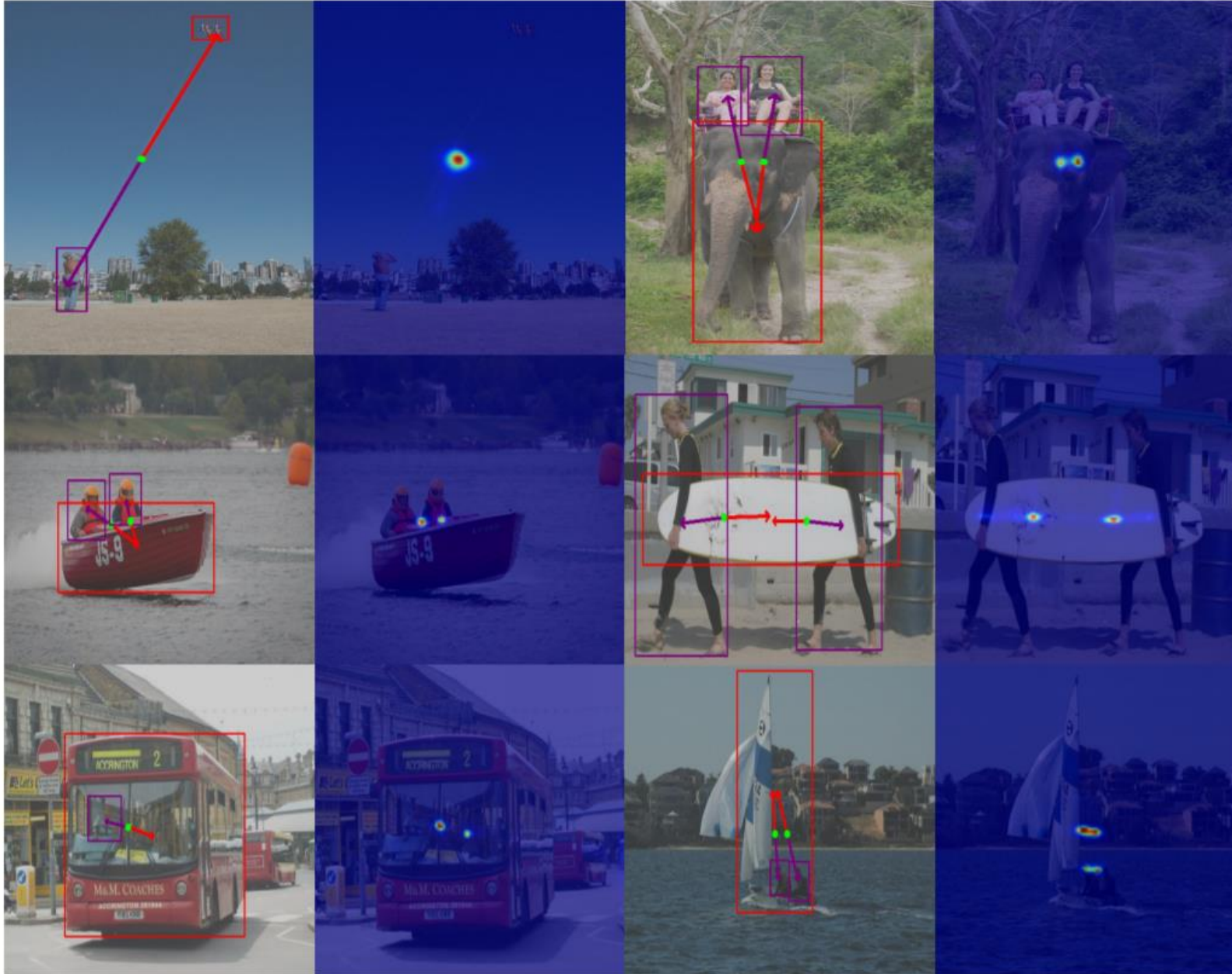
	Method	Full	Rare	Non-Rare	Time
1	Basic Model	18.46	11.97	20.40	24
2	+ Feature Fusion	18.66	11.86	20.69	26
3	+ Global Reasoning	18.63	12.61	20.42	26
4	Union Center	18.07	11.53	20.02	27
5	PPDM-DLA	19.02	12.65	20.92	27

# 실험 (HOI-A)

Method	mAP (%)	Time (ms)
Faster Interaction Net [1]	56.03	-
GMVM [1]	60.26	-
URNet [1]	66.04	-
iCAN [6]	44.23	194
TIN [15]	48.64	501
PPDM-DLA	67.03	<b>27</b>
PPDM-Hourglass	<b>71.45</b>	71

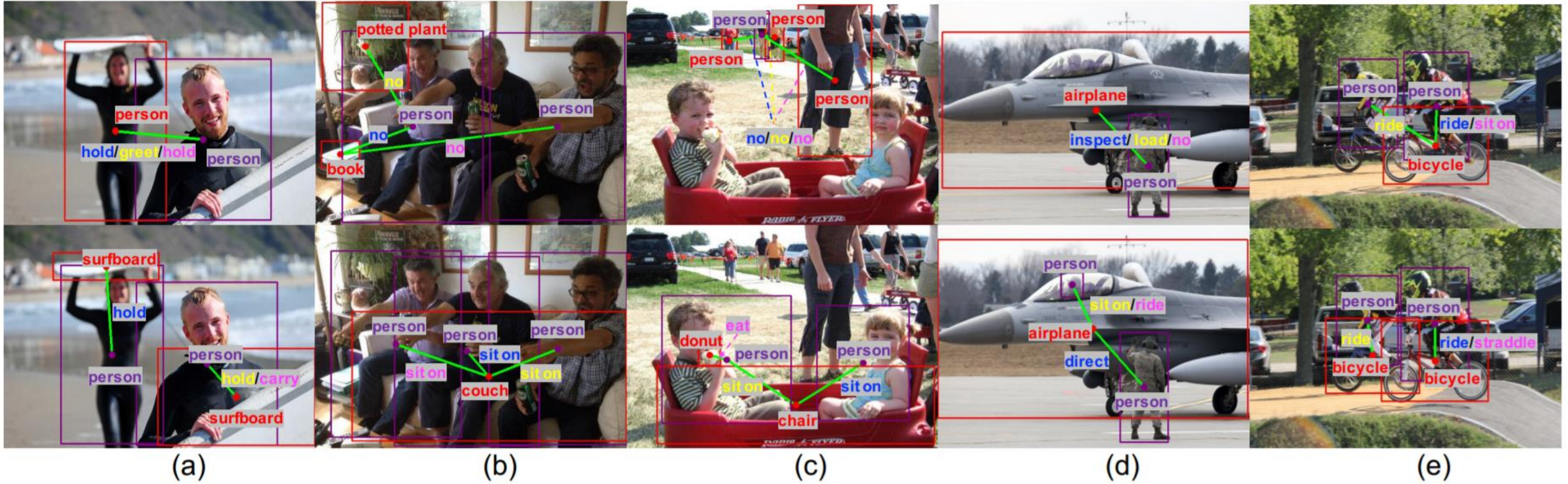
Table 3. Performance comparison on HOI-A test set.

# 실험 (인터랙션 포인트 히트맵)





# 실험 (iCAN 과 결과 비교)



Q & A