

# Human-Object Interaction Detection

KCCV 2022 Tutorial  
2022.02.24

김현우 (MLV Lab)  
고려대학교 기계학습 및 비전 연구실  
(hyunwoojkim@korea.ac.kr)

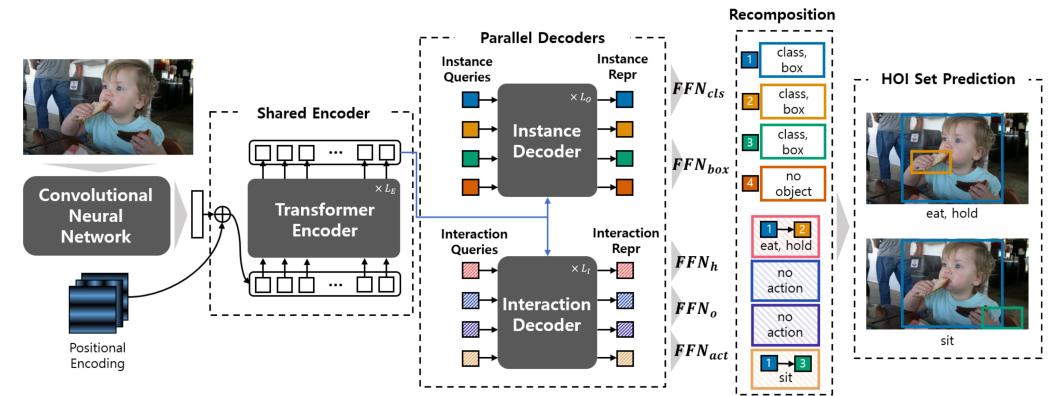
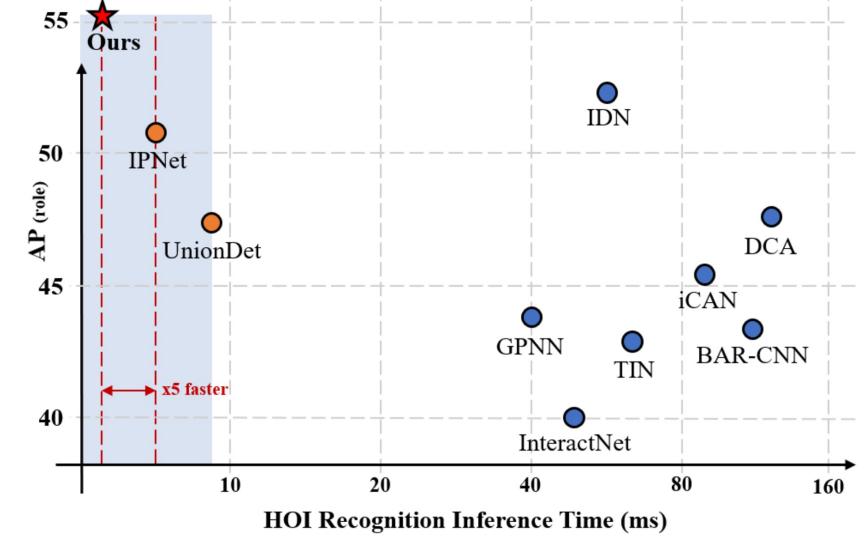
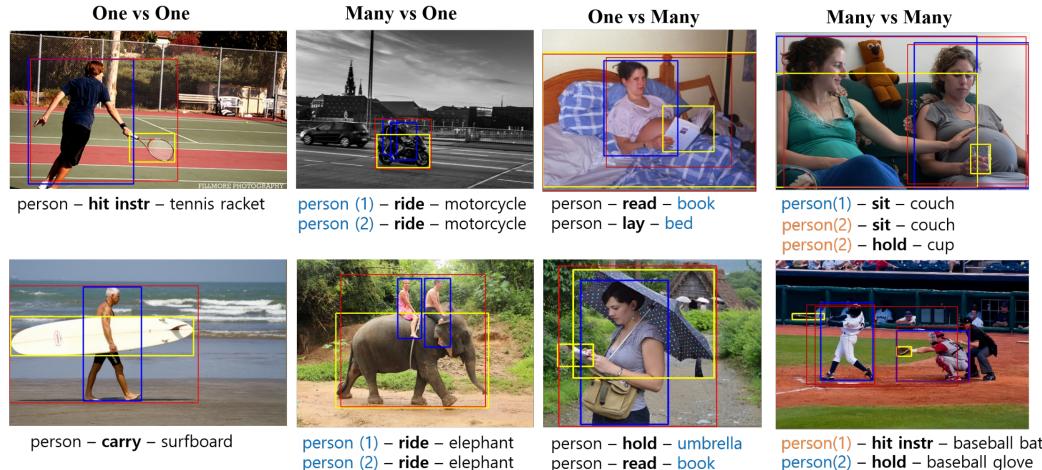
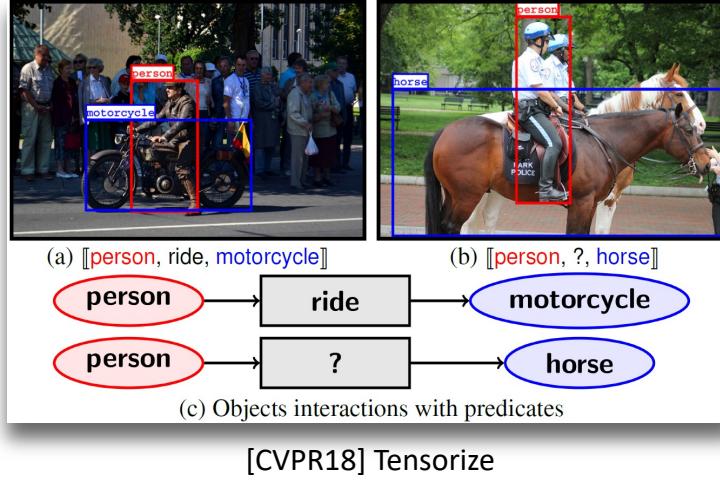
# 목차

---

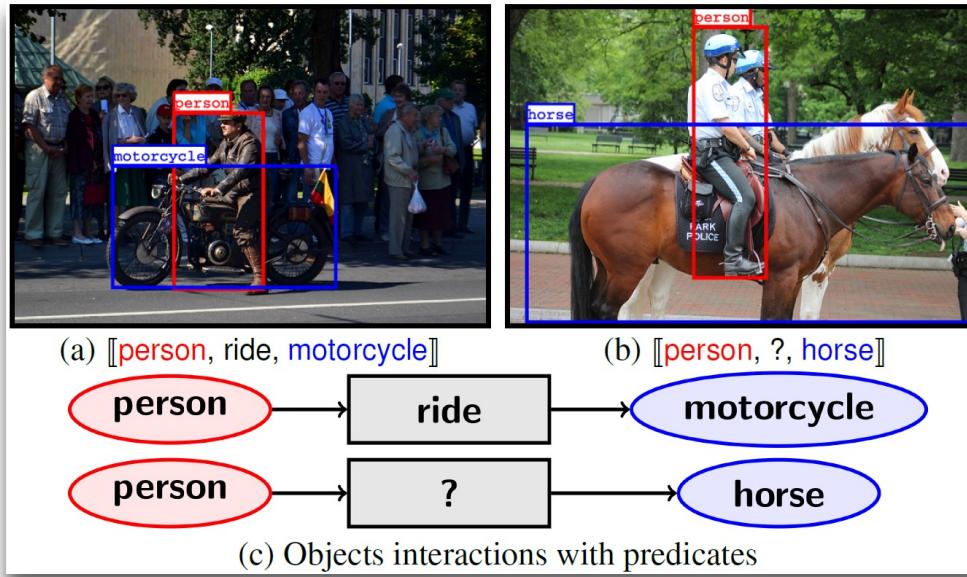
## 1. 고차원 장면이해 (Scene understanding)

- [CVPR 2018] Tensorization. 사물-사물 관계 탐지. 장면 그래프 생성
- [ECCV 2020] UnionDet. 사람-사물 상호작용 탐지
- [CVPR 2021(Oral)] HOTR. 사람-사물 상호작용 탐지

# 고차원 장면 이해 및 물체 탐지

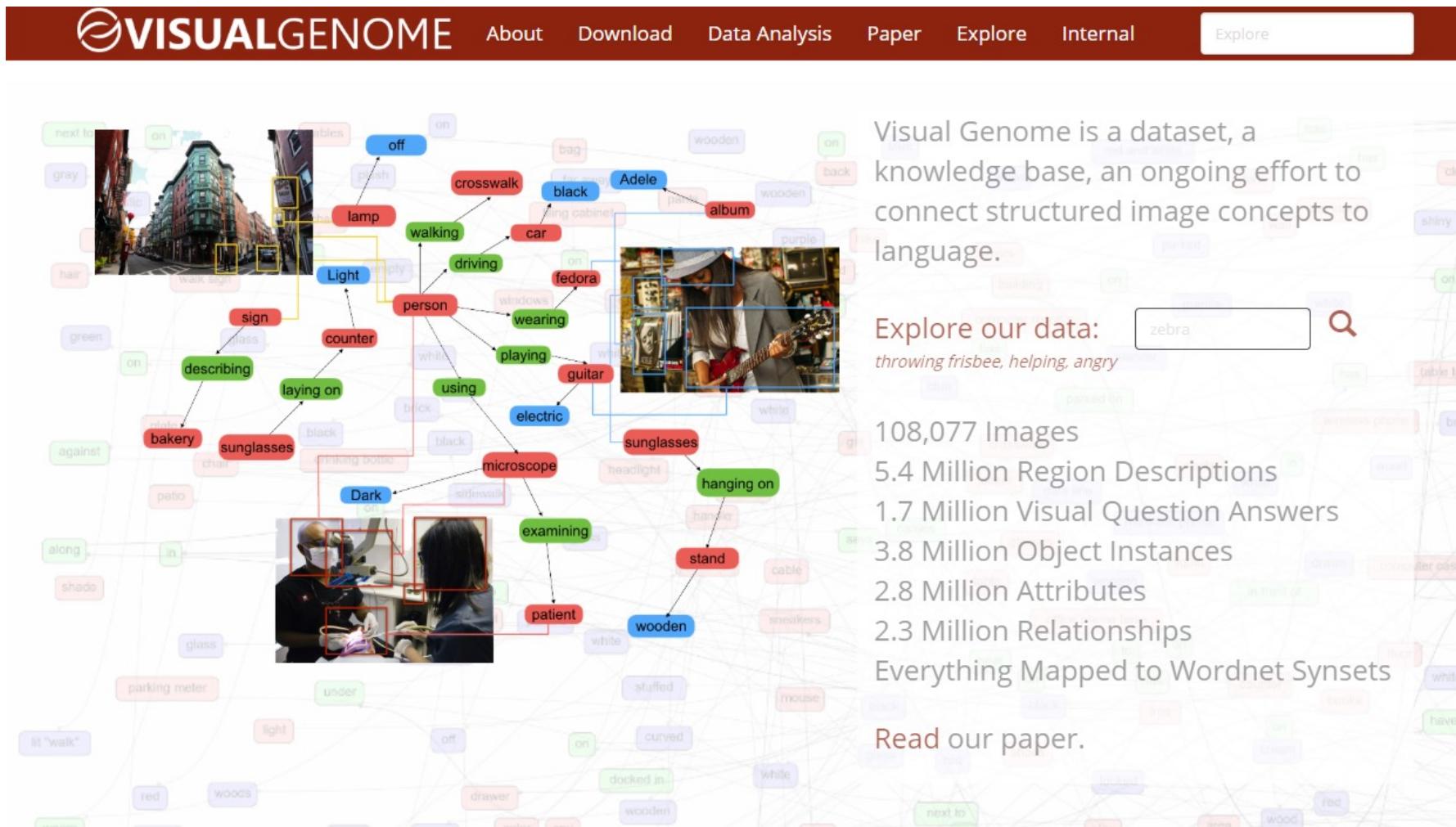


# 물체간 관계 추론 (사물-사물 상호작용 탐지)

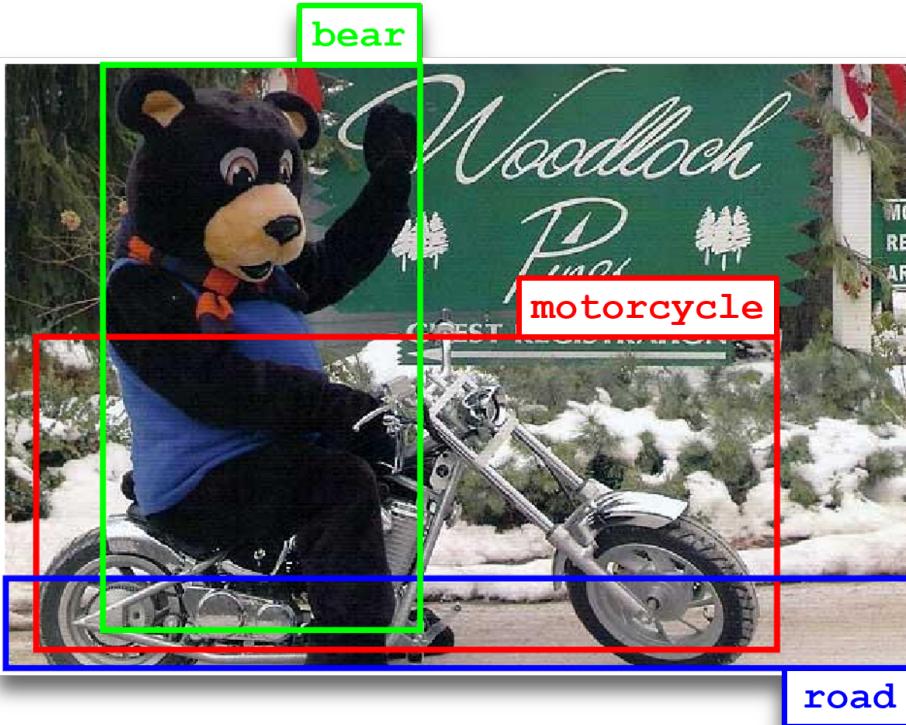


- **태스크 정의:**
  1. 주어진 이미지에서 물체를 탐지
  2. 탐지된 물체들의 관계를 추론
- **Goal:**
- **관계 추론 <물체1, 관계, 물체2>**

# 비주얼 게놈 (Visual Genome)



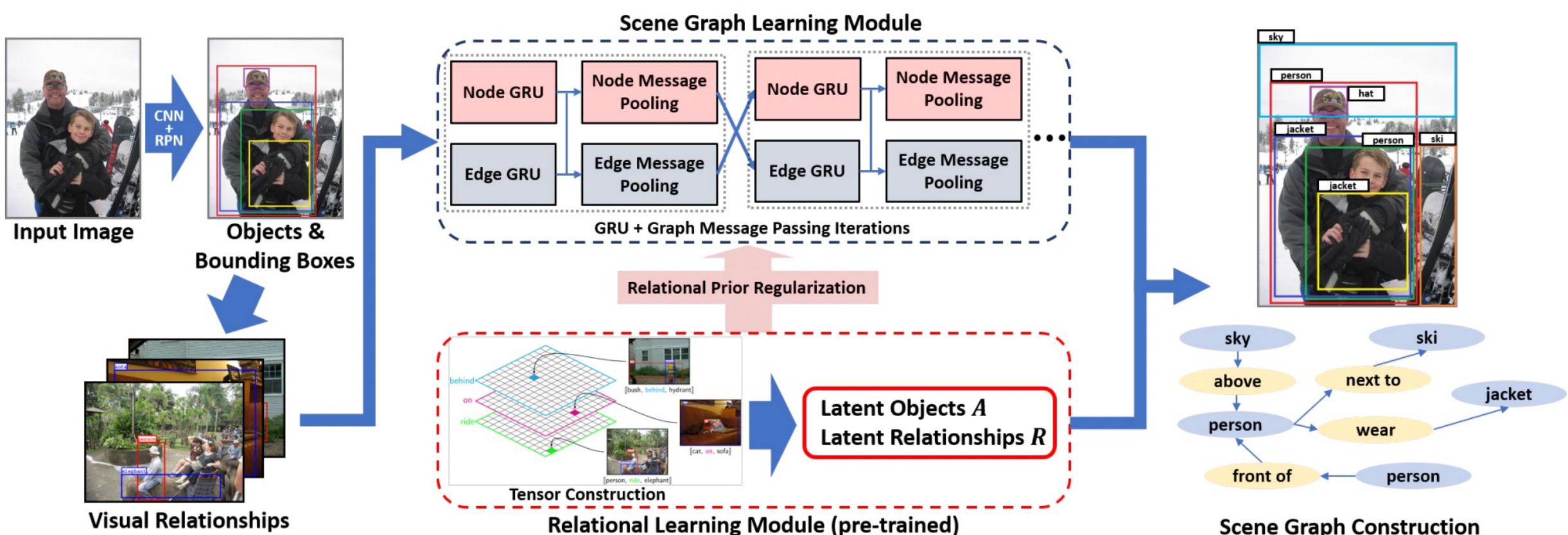
# 장면 그래프 생성 도전성



- 의미 추론의 어려움
- 데이터 의존성
- 물체관계의 다양성 및 큰 관계공간  
**<물체1, 관계, 물체2>**  
 $100 \times 100 \times 100 \sim 1M$
- 제한된 데이터 (Visual Genome)
  - >1M 이미지
  - 가능한 조합의 2% 만 커버
- 제로샷 학습 (Zero-shot Learning):  
한번도 본 적 없는 관계 예측 필요

해결 방안: 다중관계 텐서 분해. 관계에 대한 사전 지식 생성/습득

# 장면 그래프 생성 과정



# 실험 결과

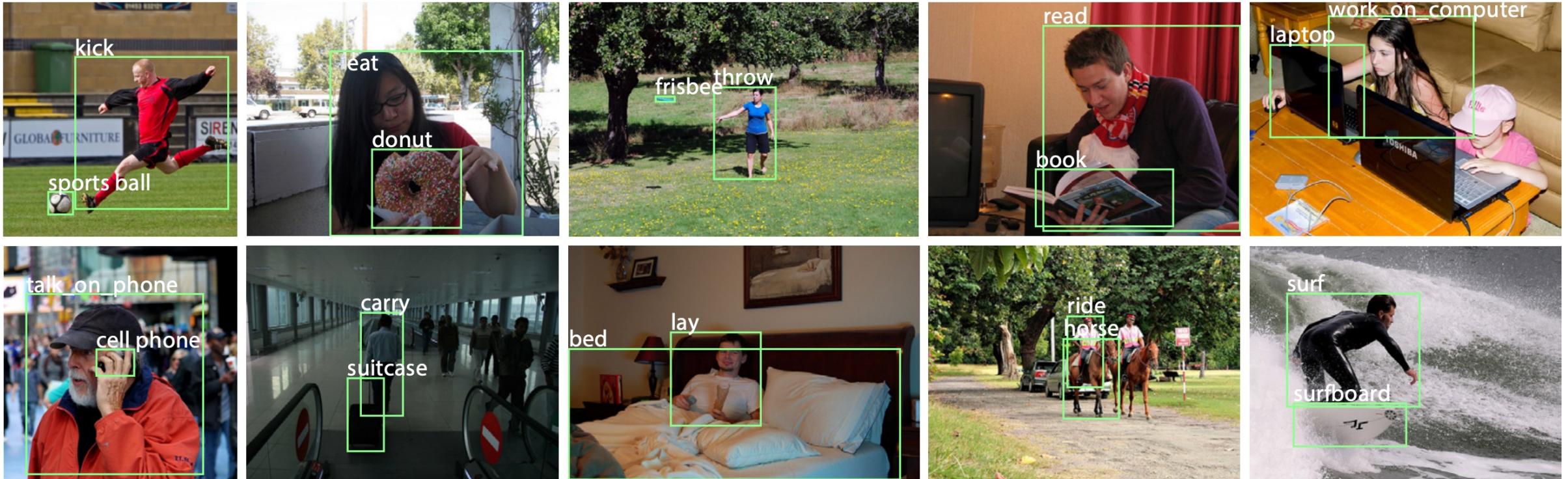
Total Relationship Detection Results: Ours (top caption) and [29] (bottom caption)



Zero-shot Relationship Detection Results: Ours (top caption) and [29] (bottom caption)



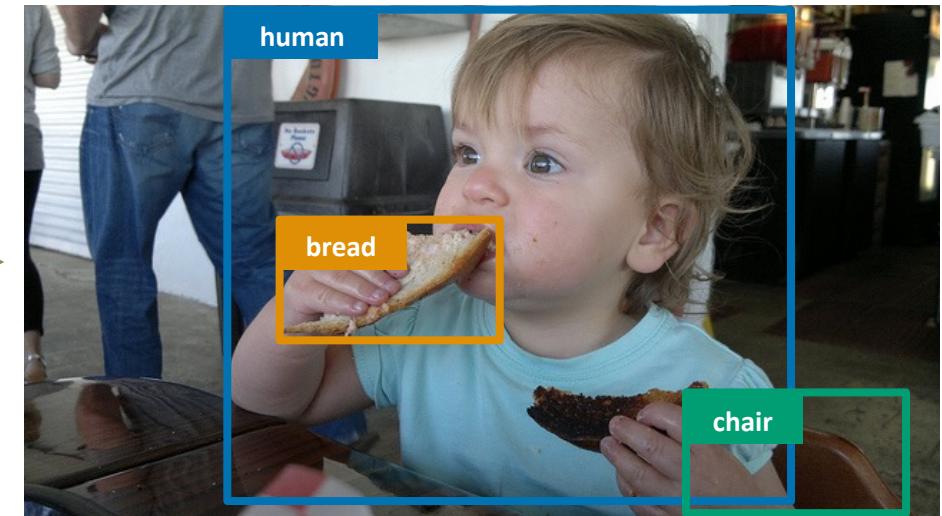
# 사람-사물 상호작용 탐지



[CVPR '18] Gkioxari, Georgia, et al. "Detecting and recognizing human-object interactions."

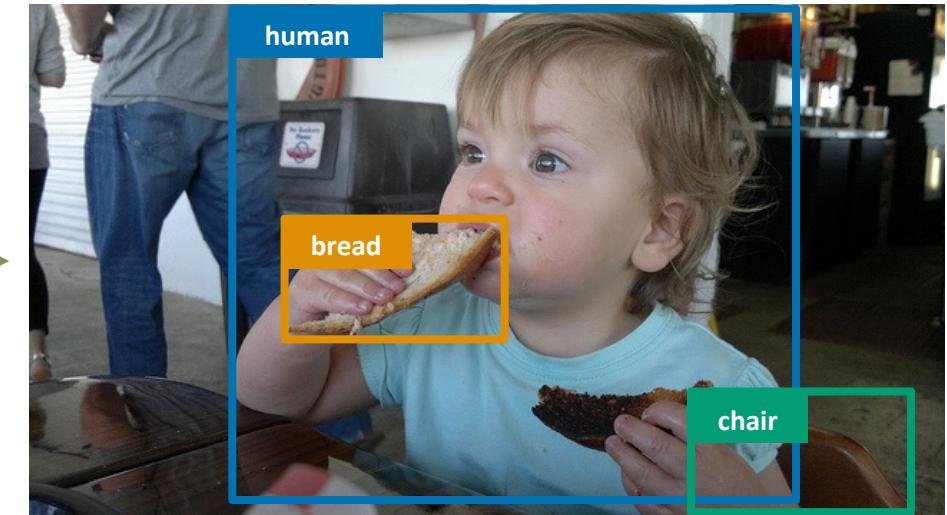
# Human-Object Interaction (HOI) Detection

- Set  $\{(\text{bbox}_1^h, \text{bbox}_1^o, [\text{eat, hold}]), (\text{bbox}_2^h, \text{bbox}_2^o, [\text{sit}])\}$



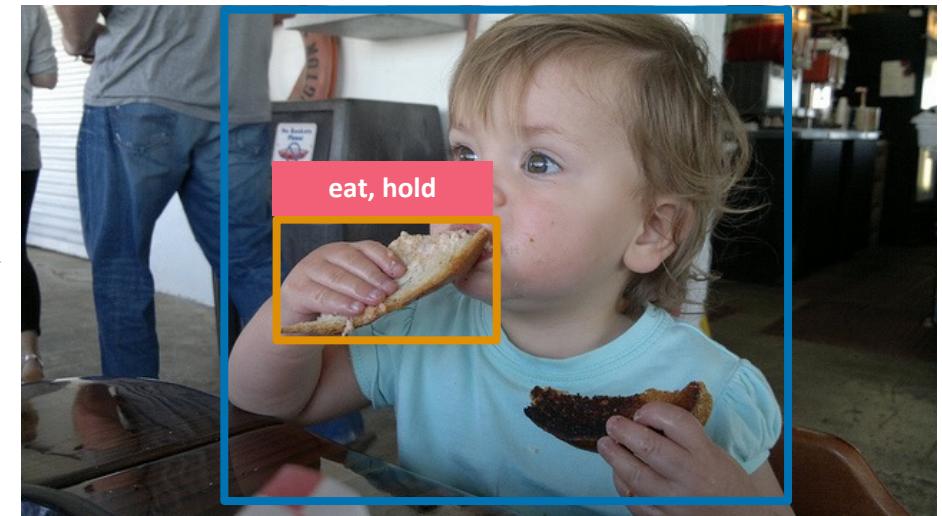
# Human-Object Interaction (HOI) Detection

- Set  $\{(\text{bbox}_1^h, \text{bbox}_1^o, (\text{bbox}_2^h, \text{bbox}_2^o))\}$



# Human-Object Interaction (HOI) Detection

- Set  $\{(\text{bbox}_1^h, \text{bbox}_1^o, [\text{eat, hold}]), (\text{bbox}_2^h, \text{bbox}_2^o, \dots)\}$



# Human-Object Interaction (HOI) Detection

- Set  $\{(bbox_1^h, bbox_1^o, [\text{eat, hold}]), (bbox_2^h, bbox_2^o, [\text{sit}])\}$



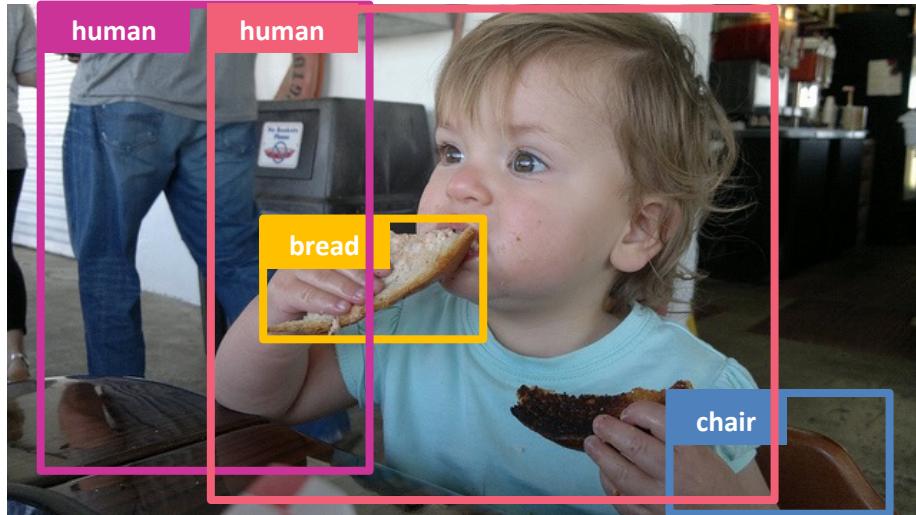
# Previous approach in HOI detection

- Sequential HOI Detectors
- 



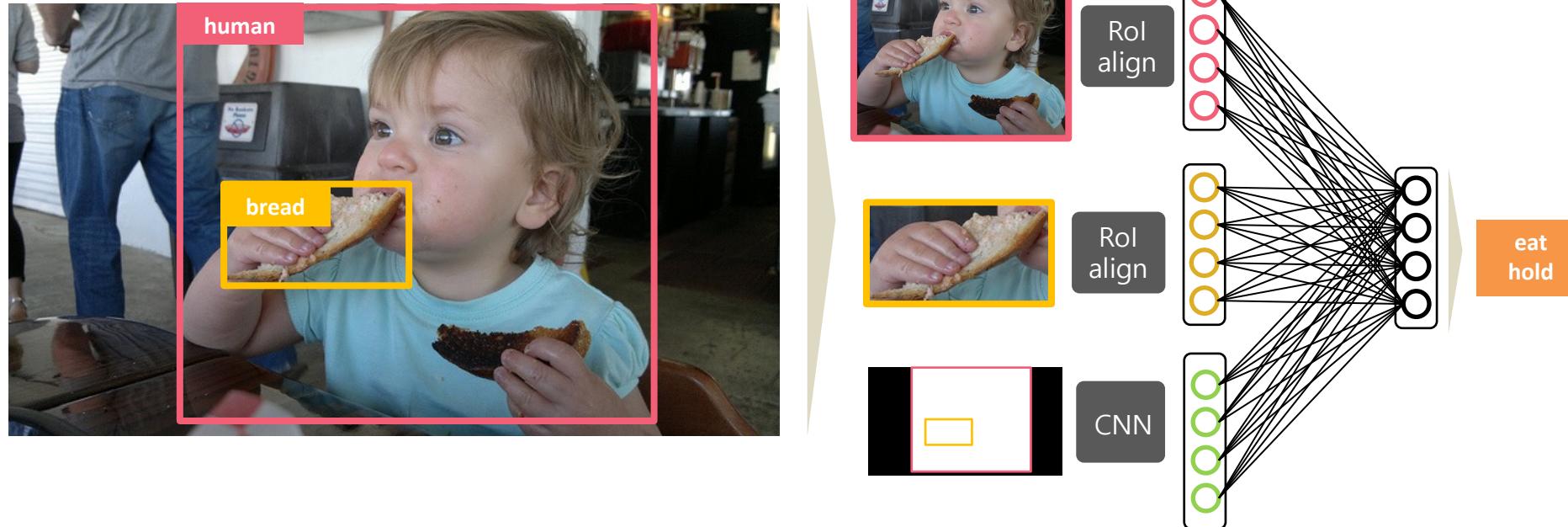
# Previous approach in HOI detection

- Sequential HOI Detectors
- 



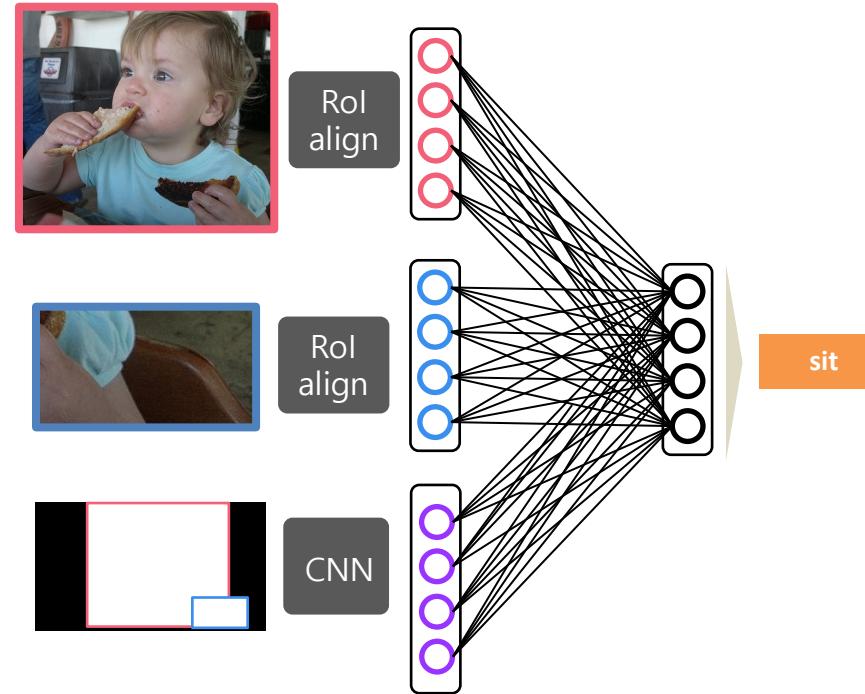
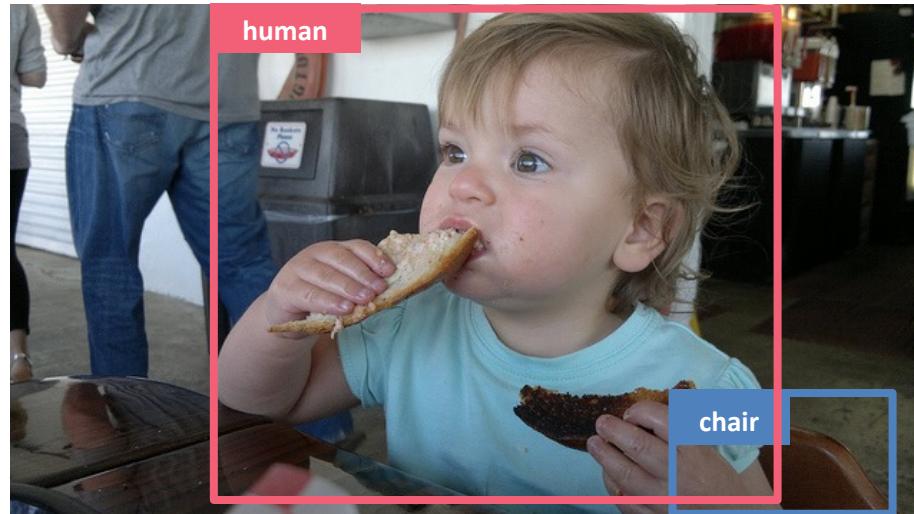
# Previous approach in HOI detection

- Sequential HOI Detectors



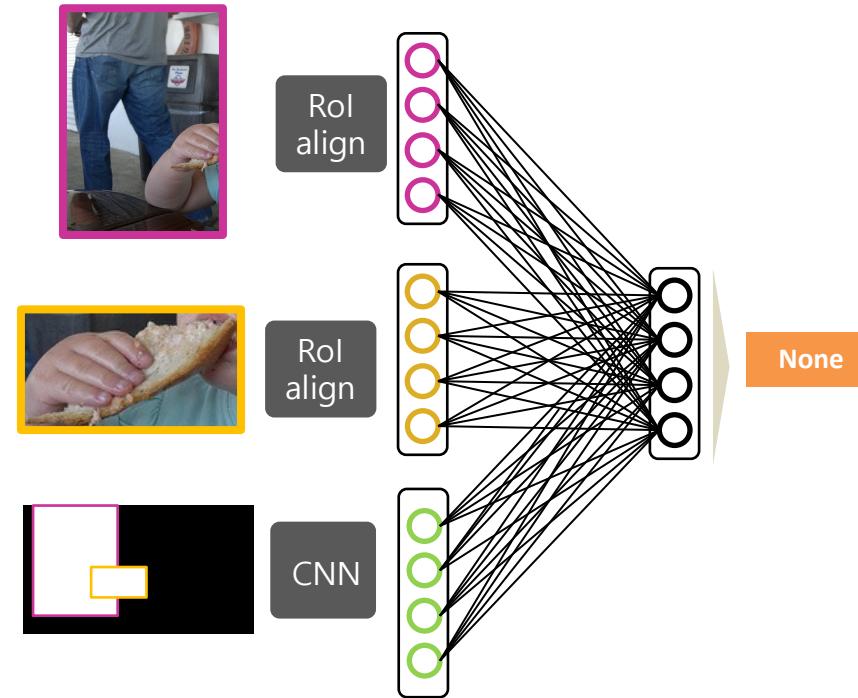
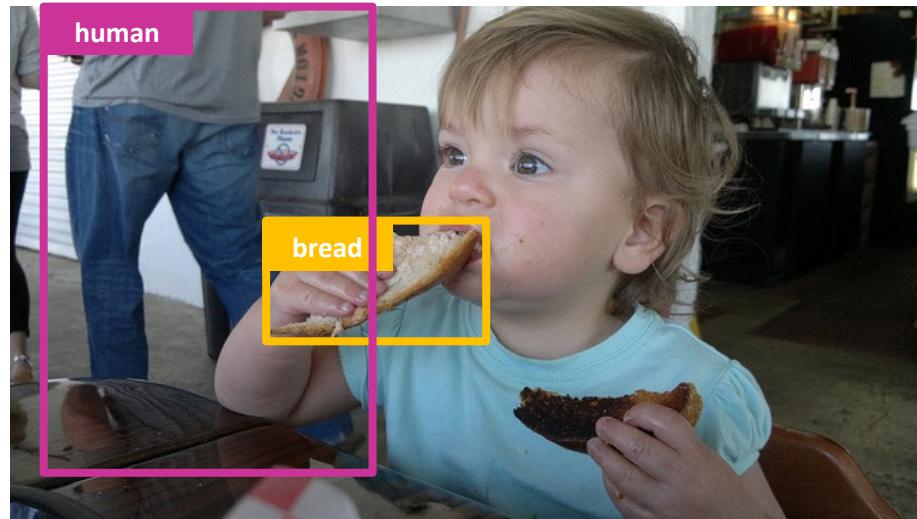
# Previous approach in HOI detection

- Sequential HOI Detectors



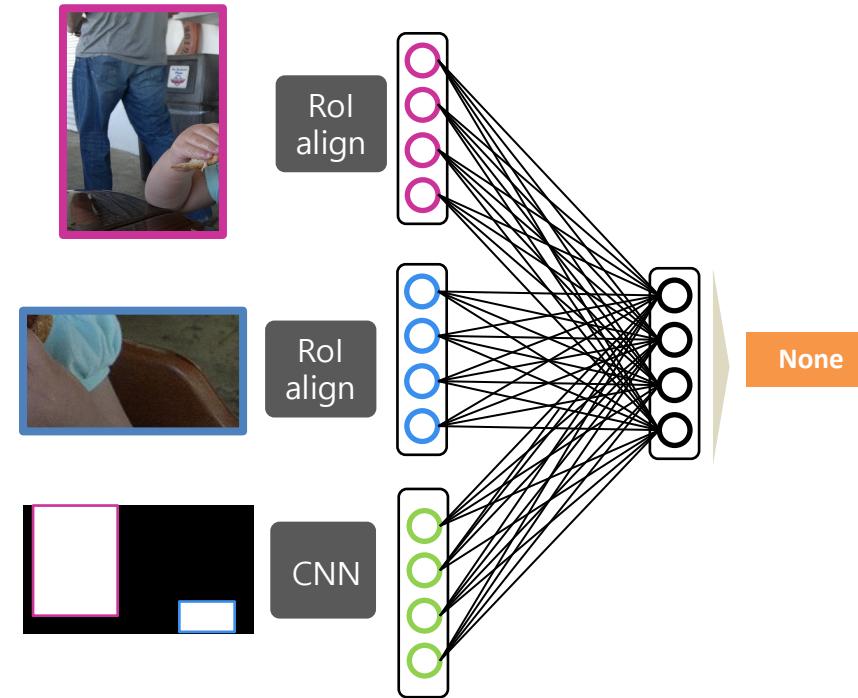
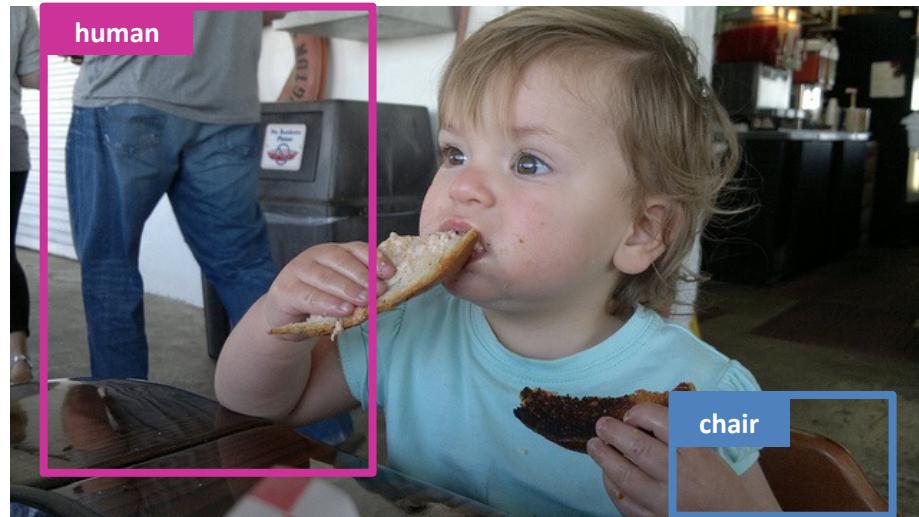
# Previous approach in HOI detection

- Sequential HOI Detectors



# Previous approach in HOI detection

- Sequential HOI Detectors



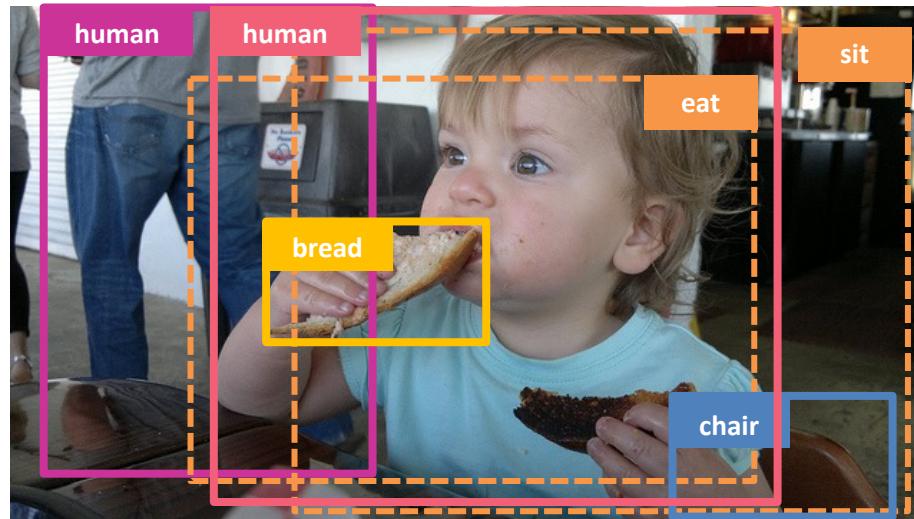
# Previous approach in HOI detection

---

- Sequential HOI Detectors
  - Intuitive Pipeline
  - Pairwise Neural Network Inference : Slow

# Previous approach in HOI detection

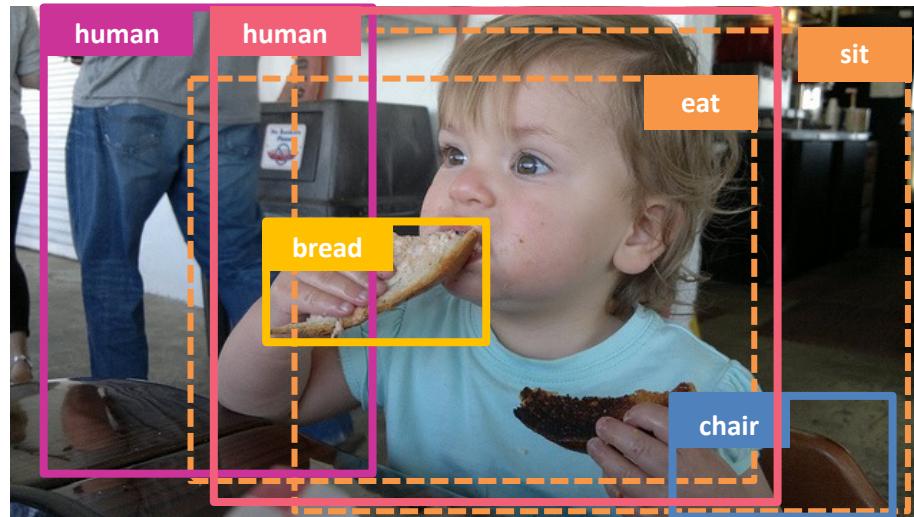
- Parallel HOI Detectors
- 



[ECCV20] UnionDet: Union-Level Detector Towards Real-Time Human-Object Interaction Detection

# Previous approach in HOI detection

- Parallel HOI Detectors

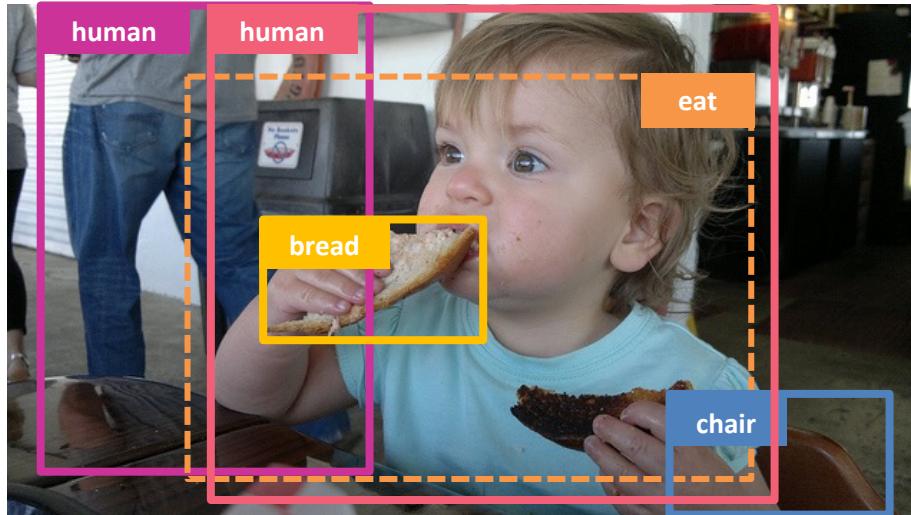


Region of Interaction

[ECCV20] UnionDet: Union-Level Detector Towards Real-Time Human-Object Interaction Detection

# Previous approach in HOI detection

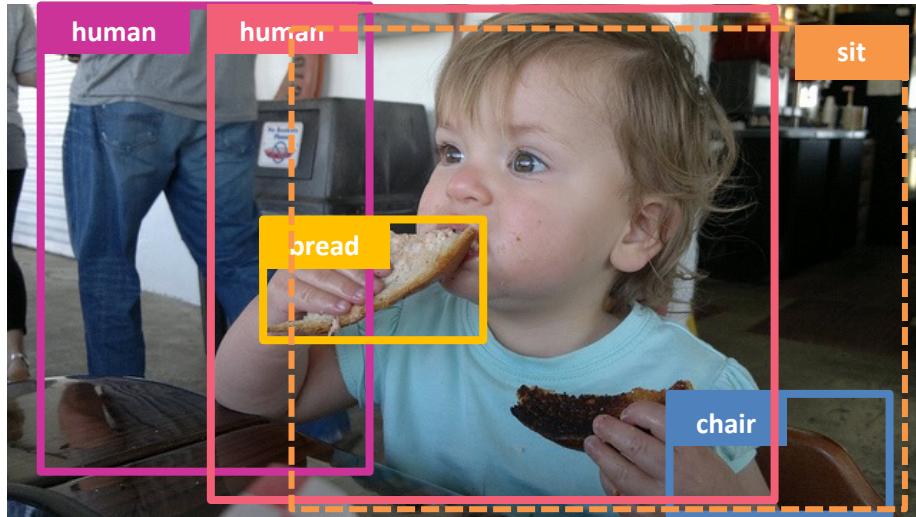
- Parallel HOI Detectors



[ECCV20] UnionDet: Union-Level Detector Towards Real-Time Human-Object Interaction Detection

# Previous approach in HOI detection

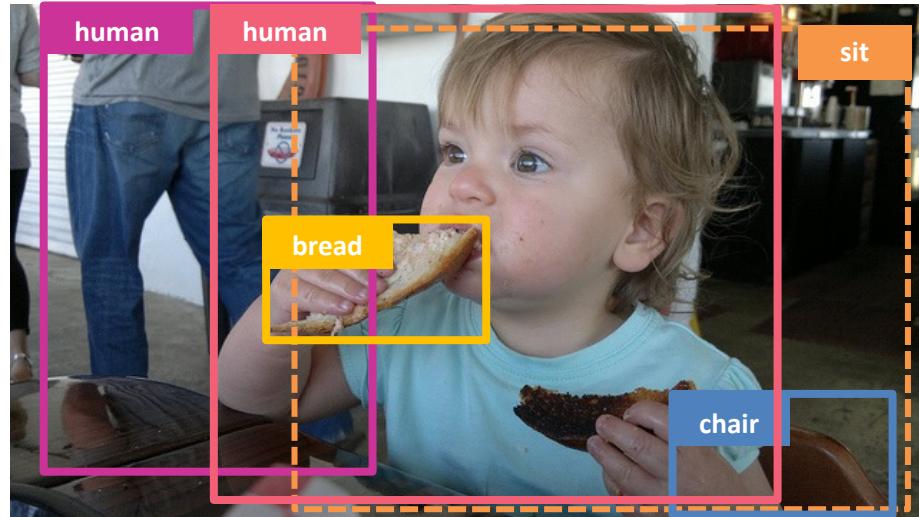
- Parallel HOI Detectors
- 



[ECCV20] UnionDet: Union-Level Detector Towards Real-Time Human-Object Interaction Detection

# Previous approach in HOI detection

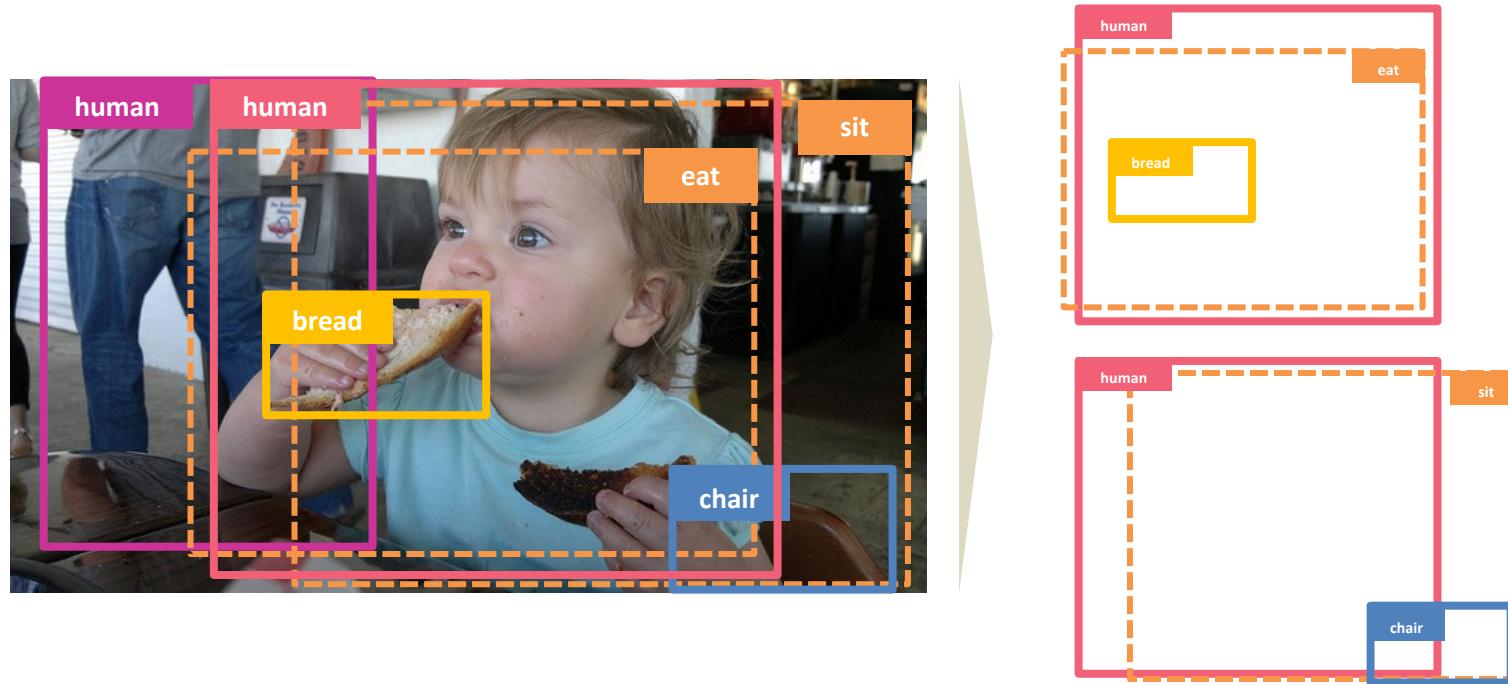
- Parallel HOI Detectors
- 



[ECCV20] UnionDet: Union-Level Detector Towards Real-Time Human-Object Interaction Detection

# Previous approach in HOI detection

- Parallel HOI Detectors



[ECCV20] UnionDet: Union-Level Detector Towards Real-Time Human-Object Interaction Detection

# Previous approach in HOI detection

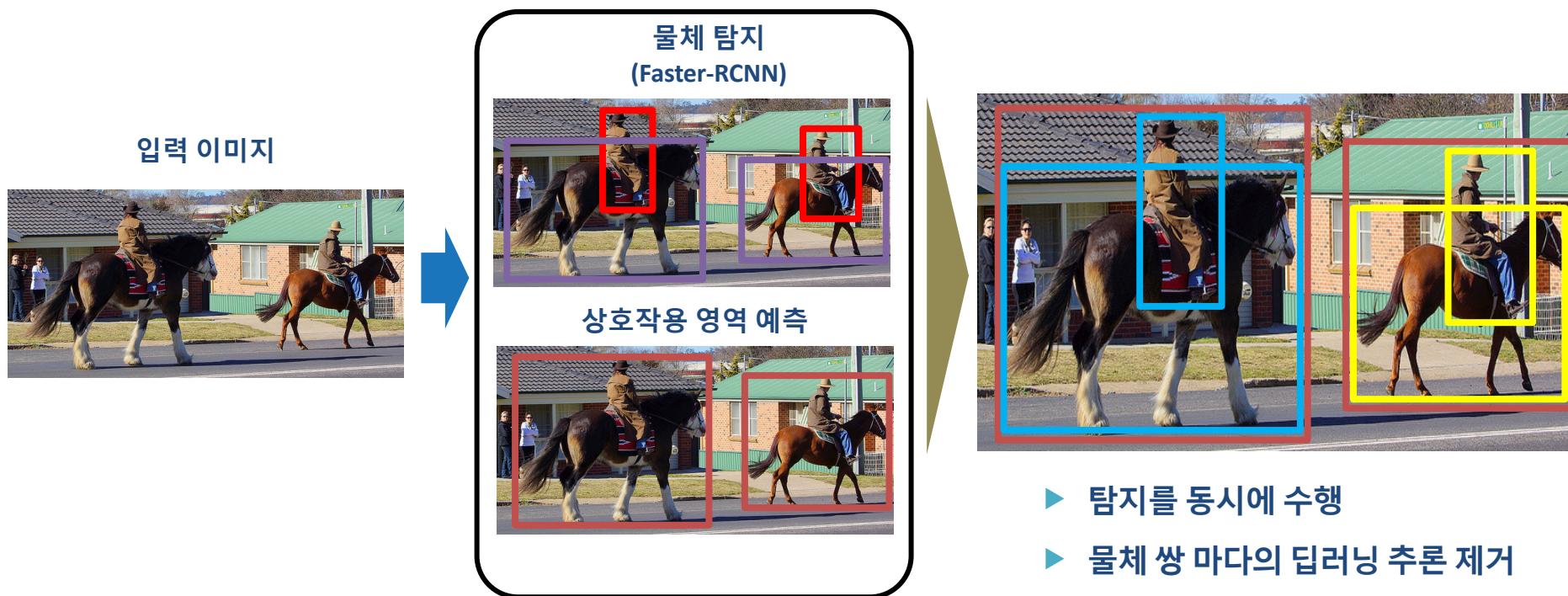
---

- Sequential HOI Detectors
  - Intuitive Pipeline
  - Pairwise Neural Network Inference : Slow
- Parallel HOI Detectors: UnionDet [ECCV20]
  - Define “region of interaction” : Union / Interaction
  - Speed-up in HOI inference time
  - However, the triplet search is still a bottleneck and has room for improvement

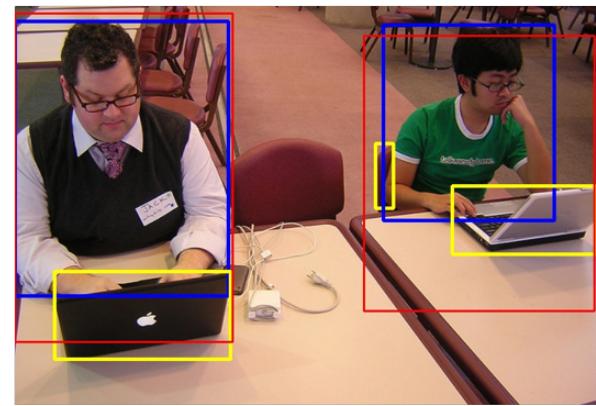
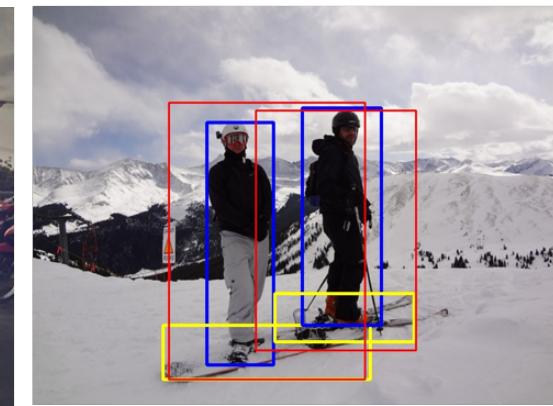
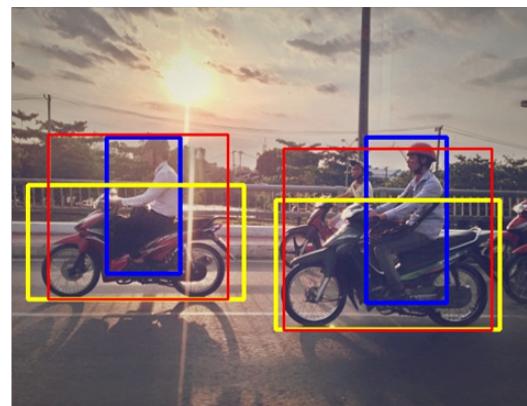
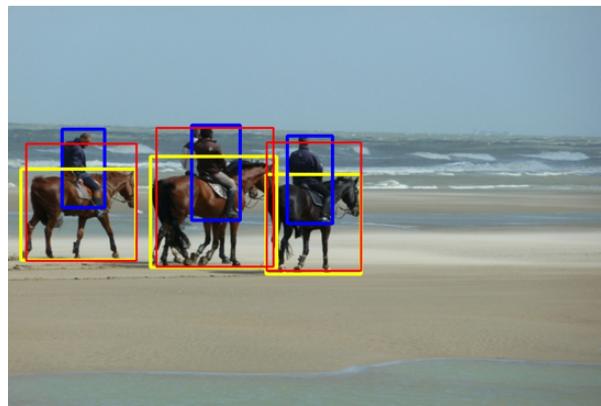
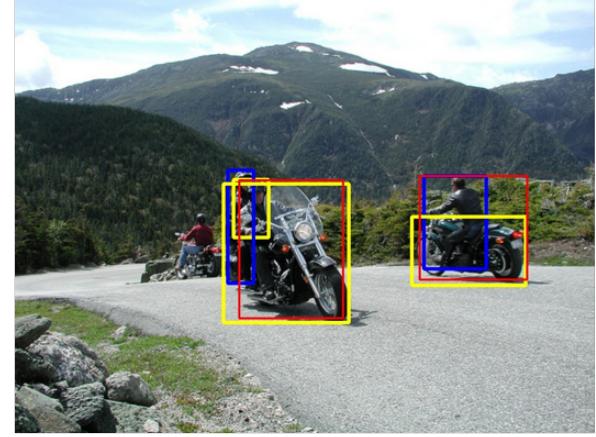
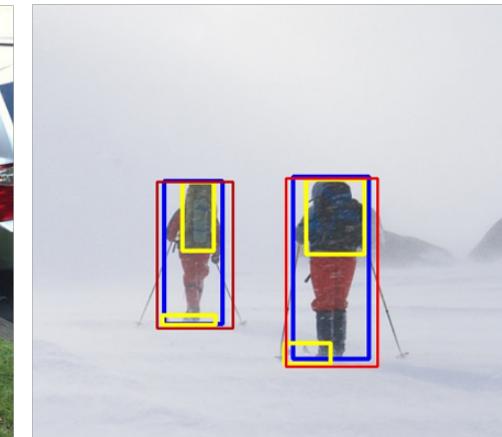
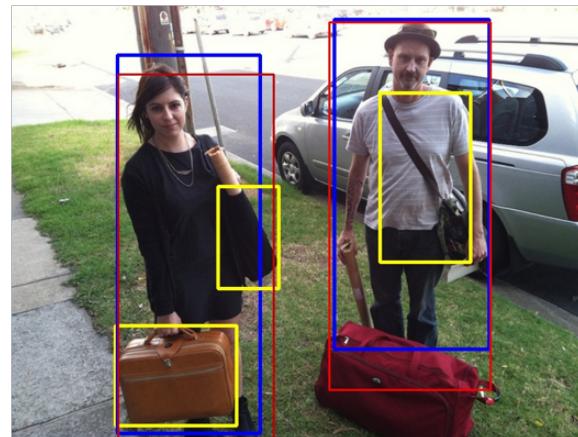
# UnionDet

- UnionDet

상호 작용 예측  
(모든 물체 쌍에 딥러닝 추론)

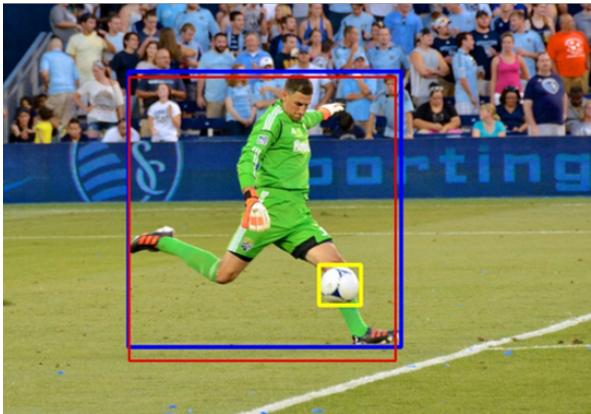
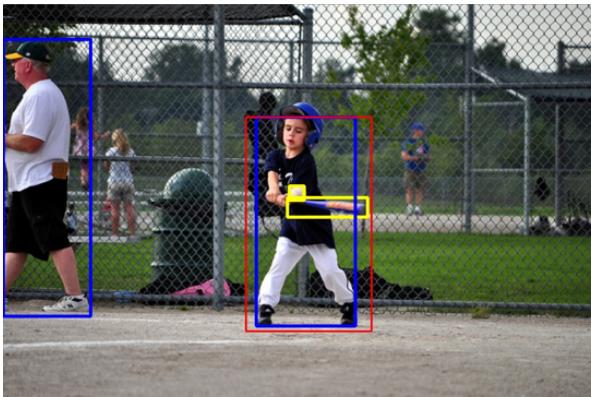


# 동일 관계 독립적 탐지 가능

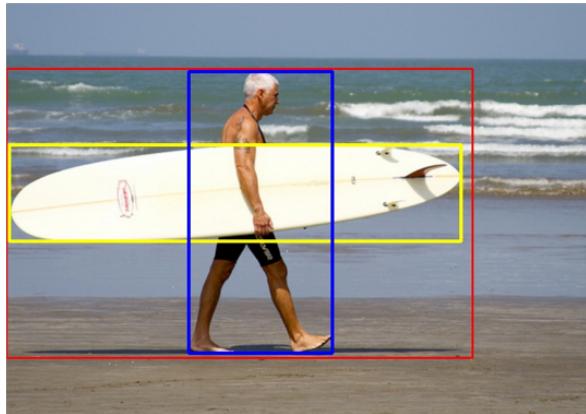
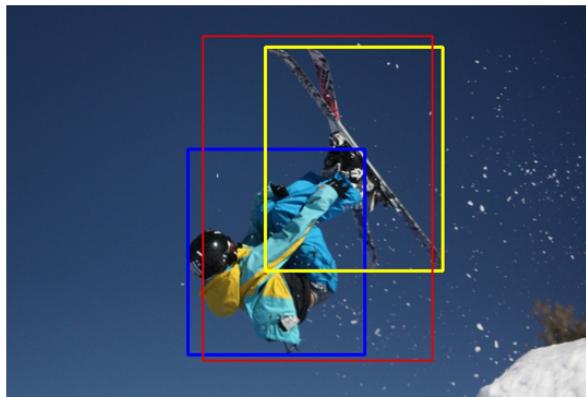


# 다양한 거리의 상호작용 탐지

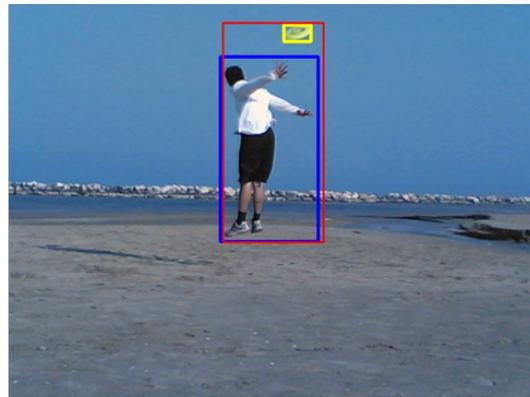
Included



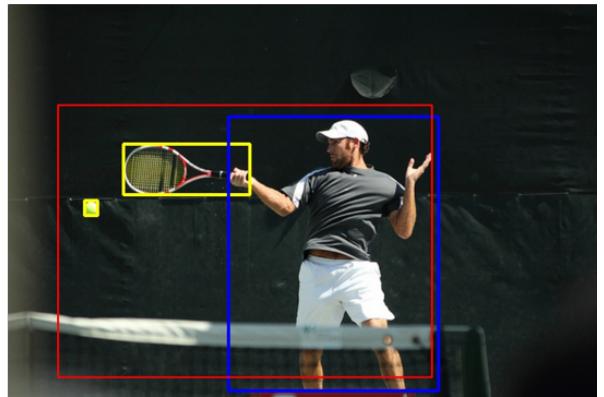
Adjacent



Distant

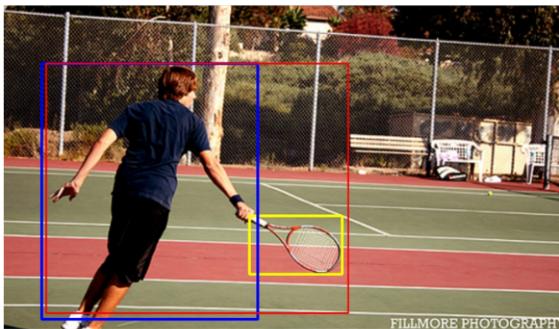


Remote



# 다대다 상호작용 탐지

One vs One



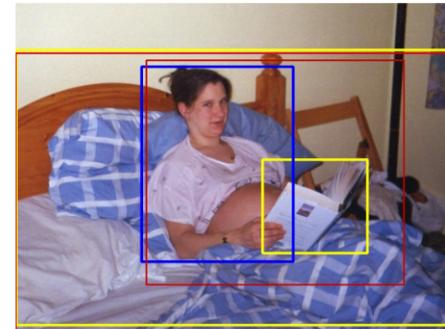
person – **hit instr** – tennis racket

Many vs One



person (1) – **ride** – motorcycle  
person (2) – **ride** – motorcycle

One vs Many

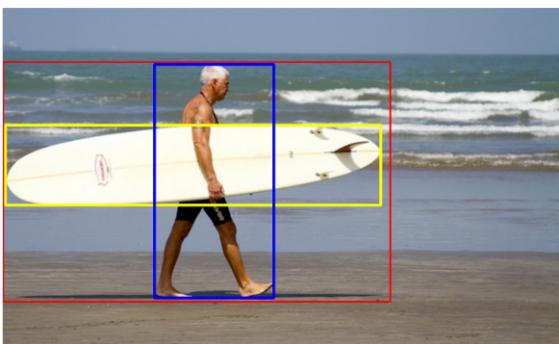


person – **read** – book  
person – **lay** – bed

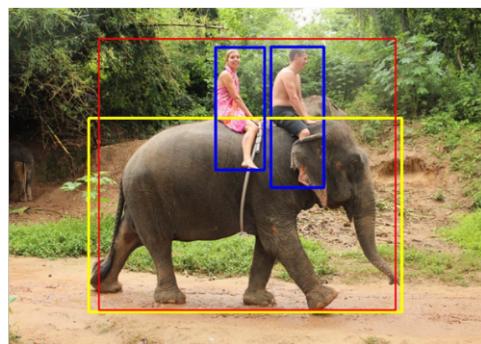
Many vs Many



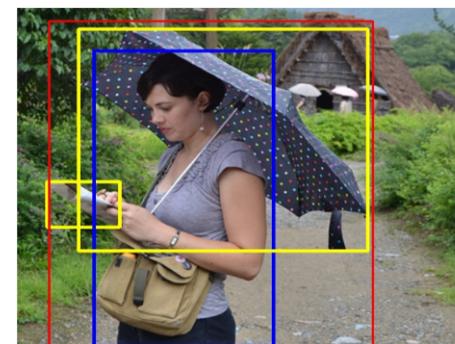
person(1) – **sit** – couch  
person(2) – **sit** – couch  
person(2) – **hold** – cup



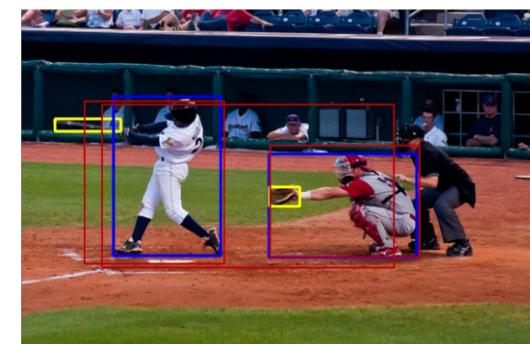
person – **carry** – surfboard



person (1) – **ride** – elephant  
person (2) – **ride** – elephant



person – **hold** – umbrella  
person – **read** – book



person(1) – **hit instr** – baseball bat  
person(2) – **hold** – baseball glove  
person(2) – **look** – person(1)

# HOTR: End-to-End Human-Object Interaction Detection with Transformers



**Bumsoo Kim**  
Kakao Brain  
Korea Univ.



**Junhyun Lee**  
Korea Univ.



**Jaewoo Kang**  
Korea Univ.



**Eun-Sol Kim**  
Kakao Brain



**Hyunwoo J. Kim**  
Korea Univ.

kakao**brain**



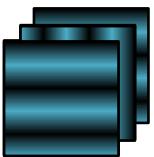
KOREA  
UNIVERSITY

# Architecture Overview

---



Convolutional  
Neural  
Network

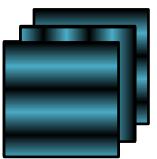


Positional  
Encoding

# Architecture Overview



Convolutional  
Neural  
Network



Positional  
Encoding

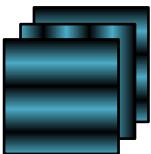
# Architecture Overview



Convolutional  
Neural  
Network



Transformer  
Encoder  $\times L_E$

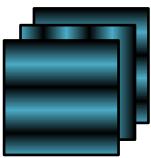


Positional  
Encoding

# Architecture Overview



Convolutional  
Neural  
Network



Positional  
Encoding

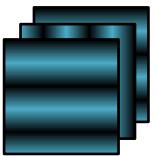
Transformer  
Encoder  $\times L_E$



# Architecture Overview

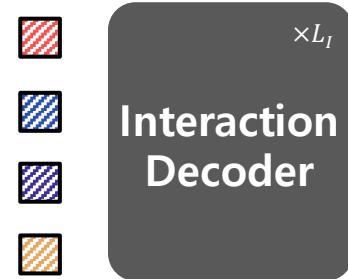
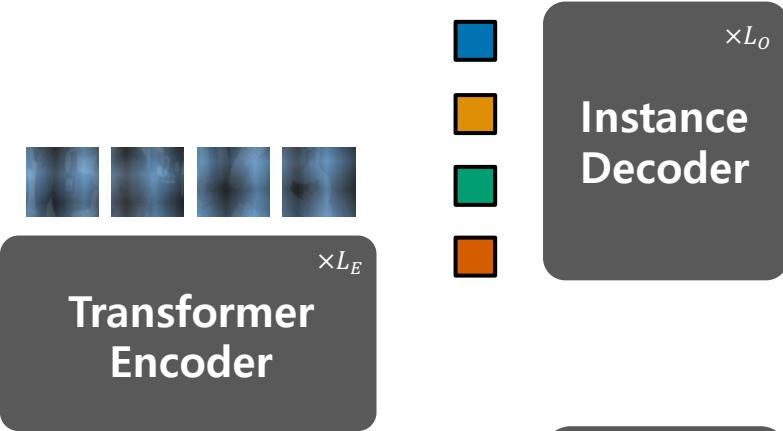


Convolutional  
Neural  
Network



Positional  
Encoding

Instance Queries

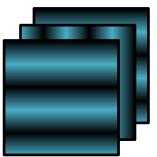


Interaction Queries

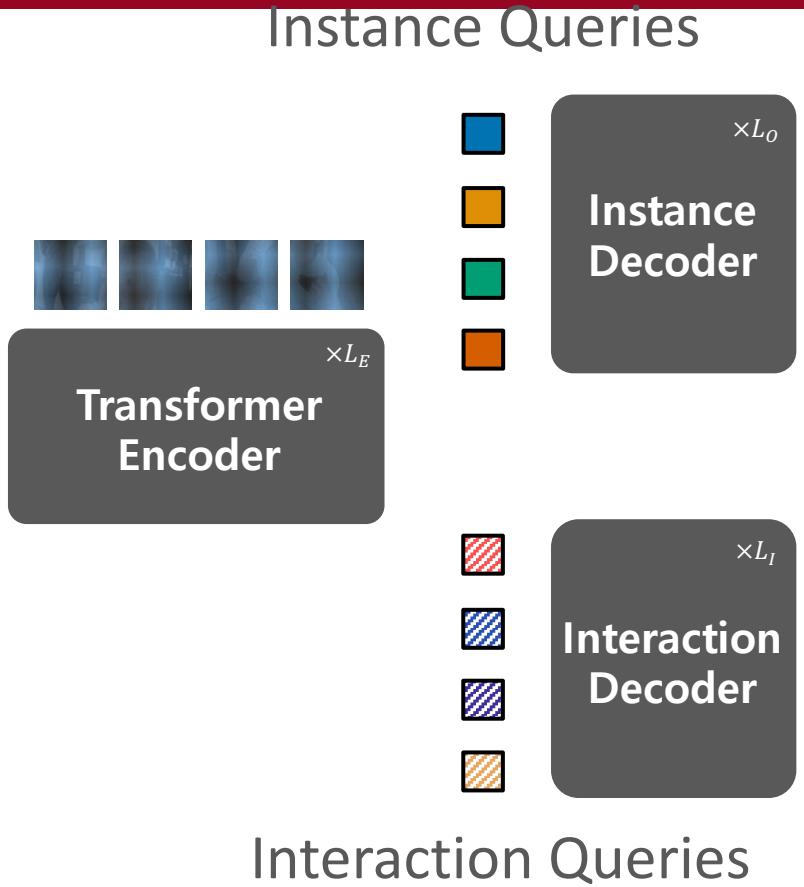
# Architecture Overview



Convolutional  
Neural  
Network



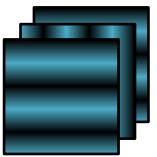
Positional  
Encoding



# Architecture Overview



Convolutional  
Neural  
Network

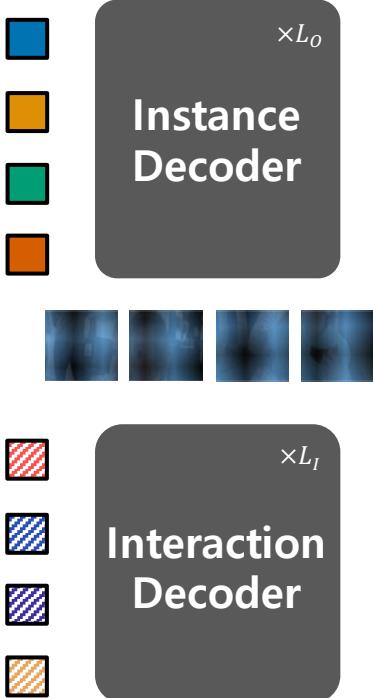


Positional  
Encoding

Transformer  
Encoder  $\times L_E$

Interaction Queries

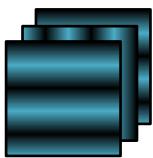
Instance Queries



# Architecture Overview



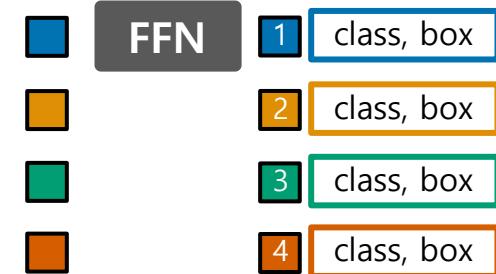
Convolutional  
Neural  
Network



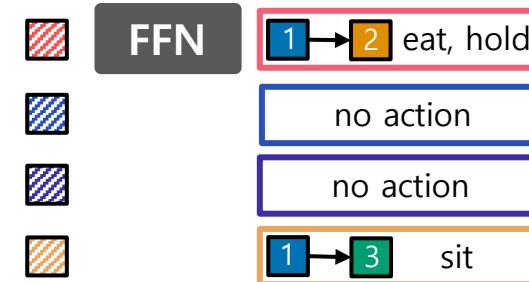
Positional  
Encoding

Transformer  
Encoder  $\times L_E$

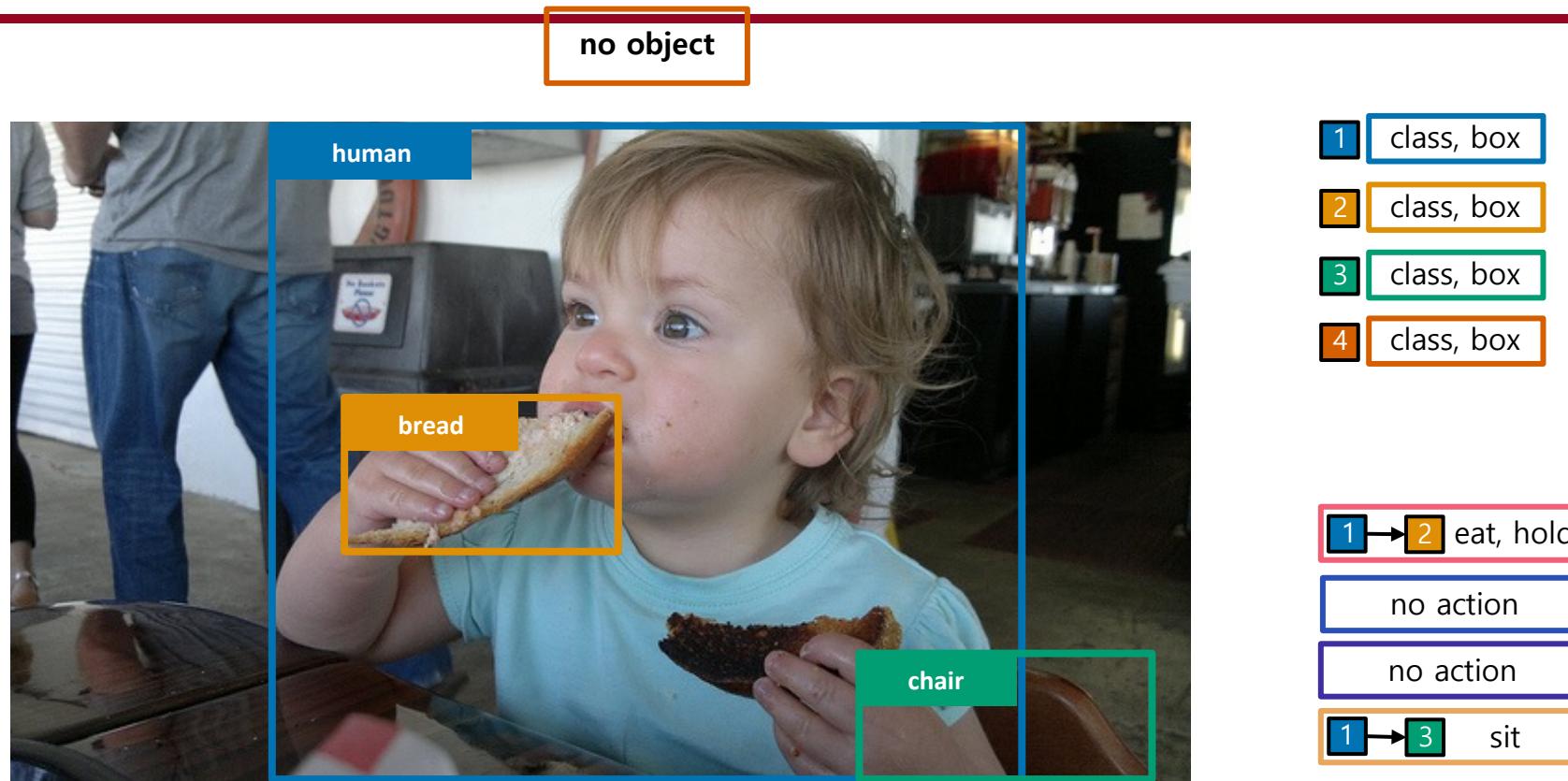
Instance  
Decoder  $\times L_O$



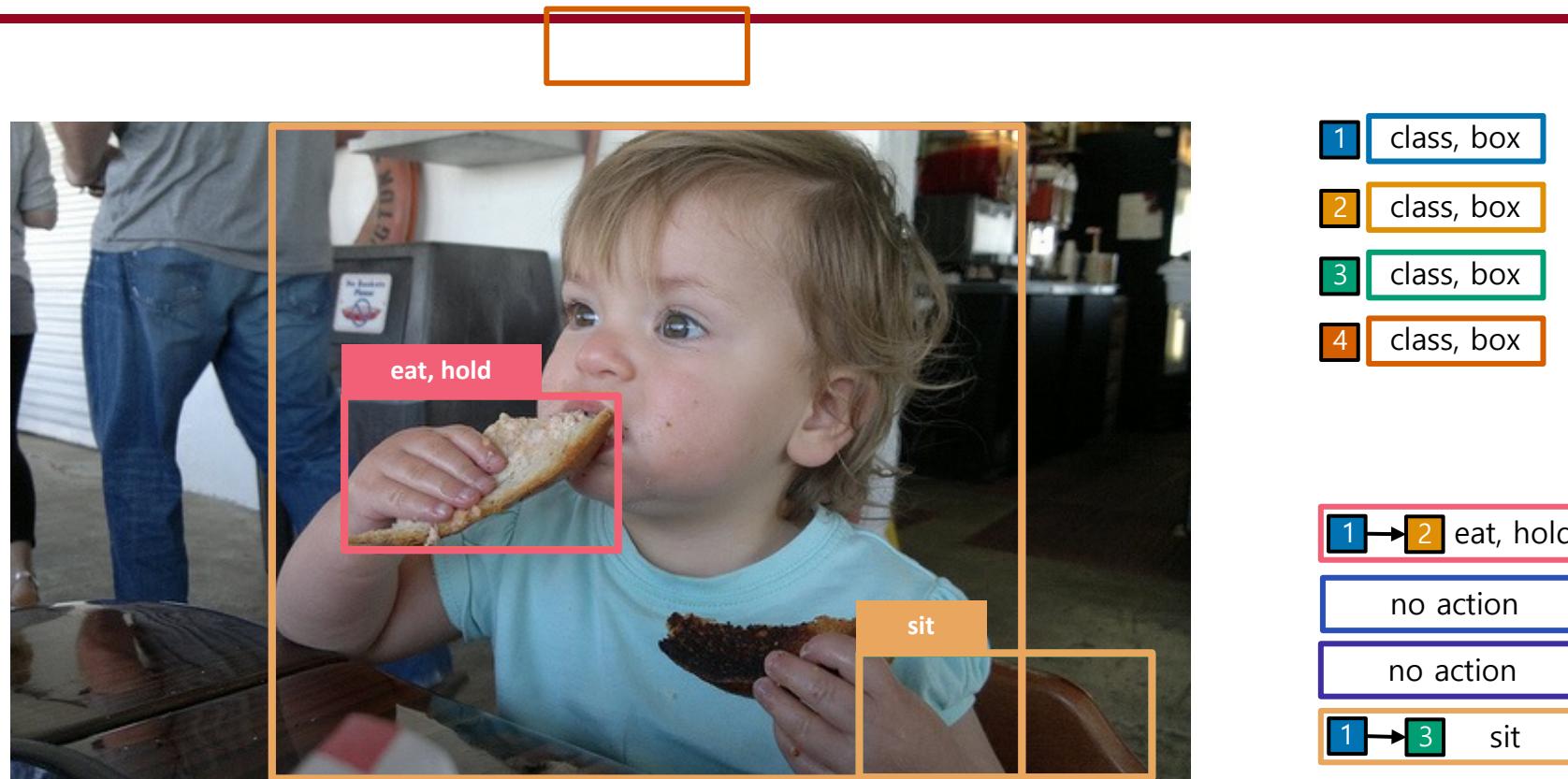
Interaction  
Decoder  $\times L_I$



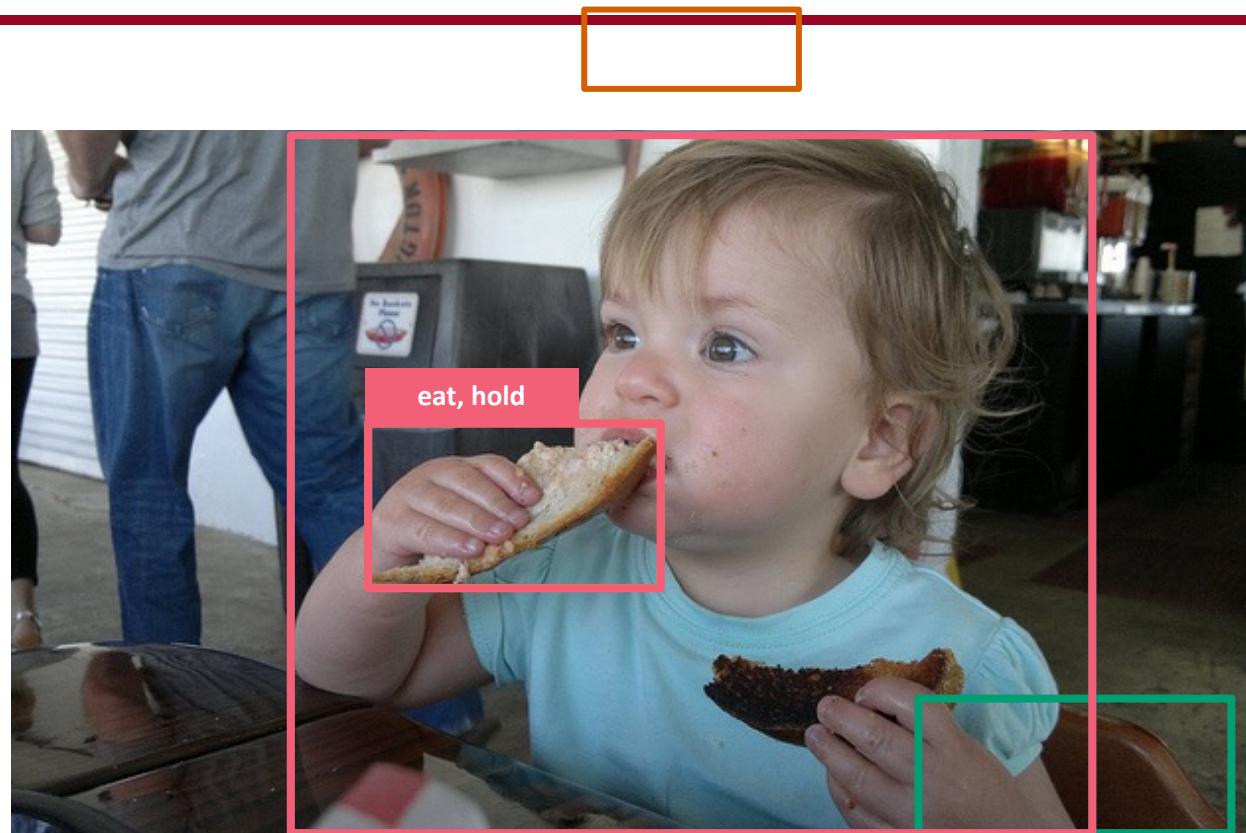
# Architecture Overview



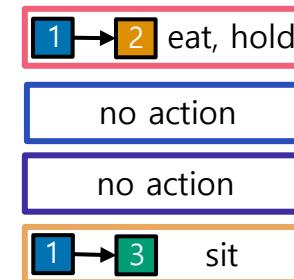
# Architecture Overview



# Architecture Overview

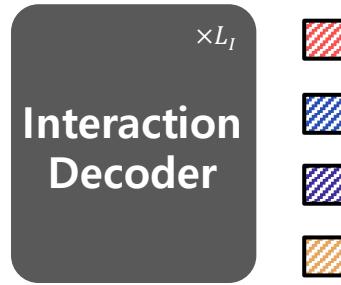
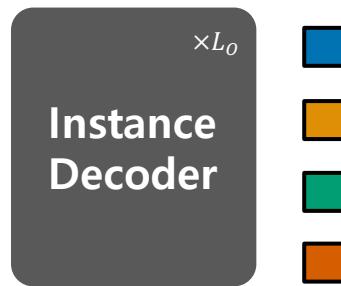


- 1 class, box
- 2 class, box
- 3 class, box
- 4 class, box

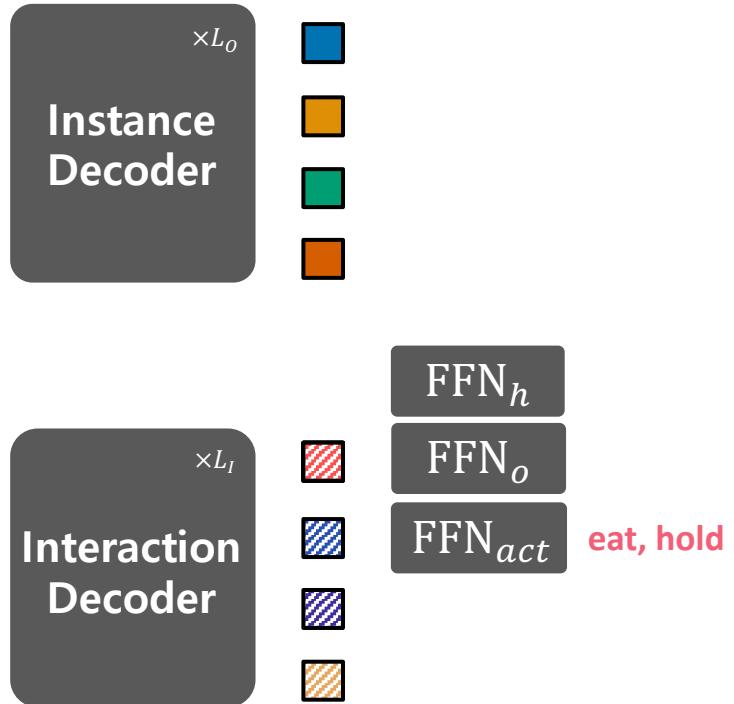


# HO Pointers

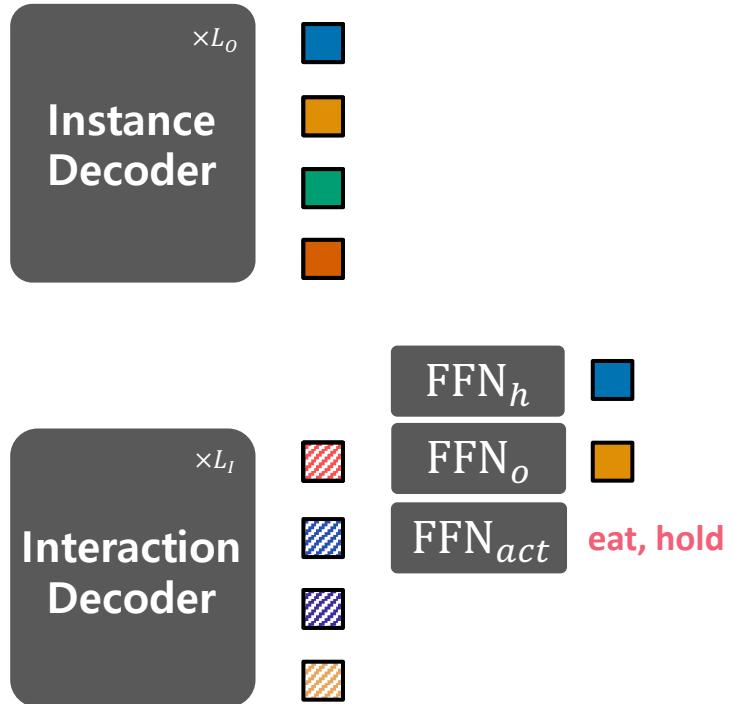
---



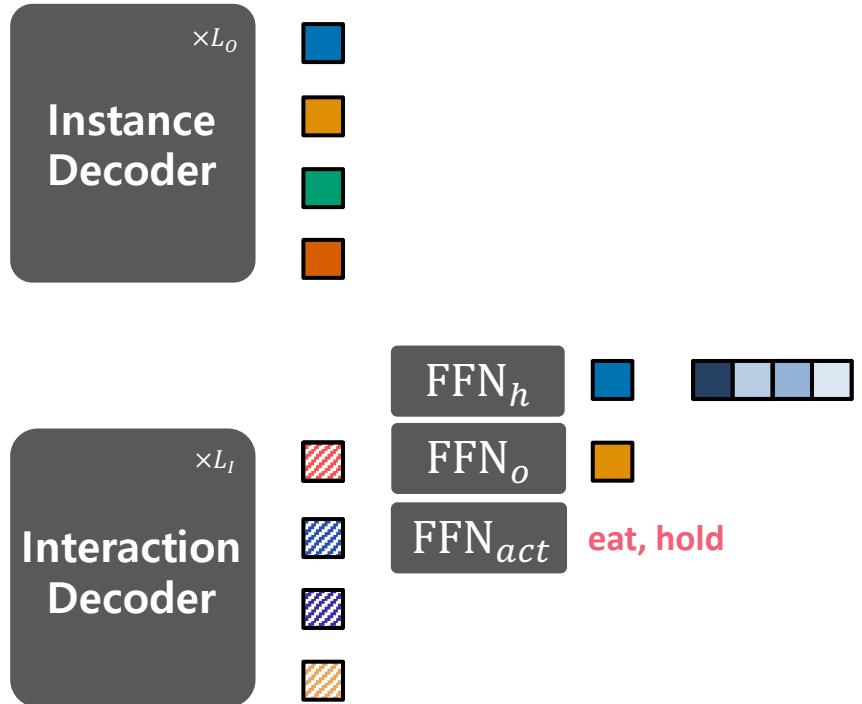
# HO Pointers



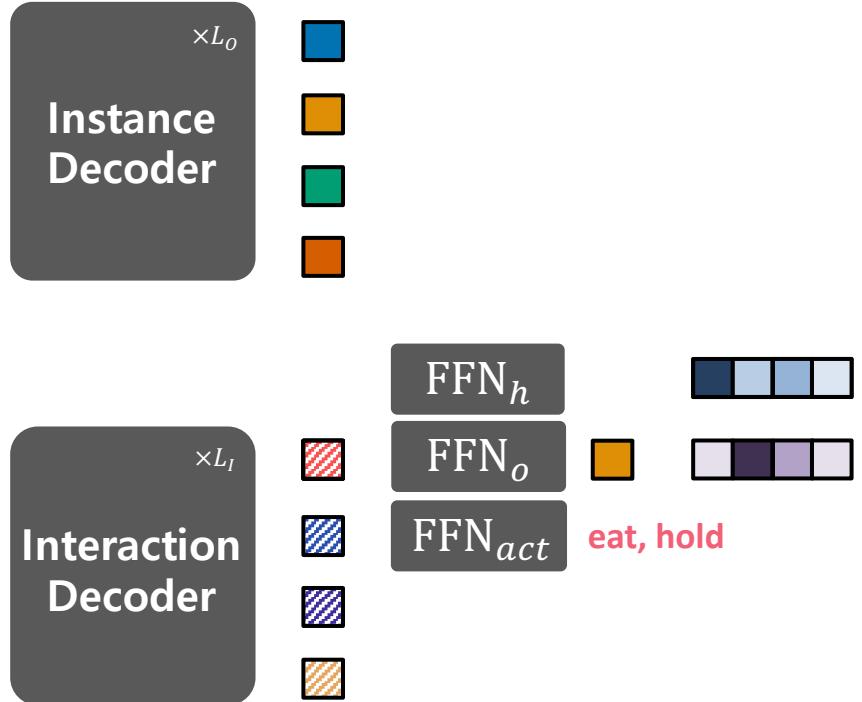
# HO Pointers



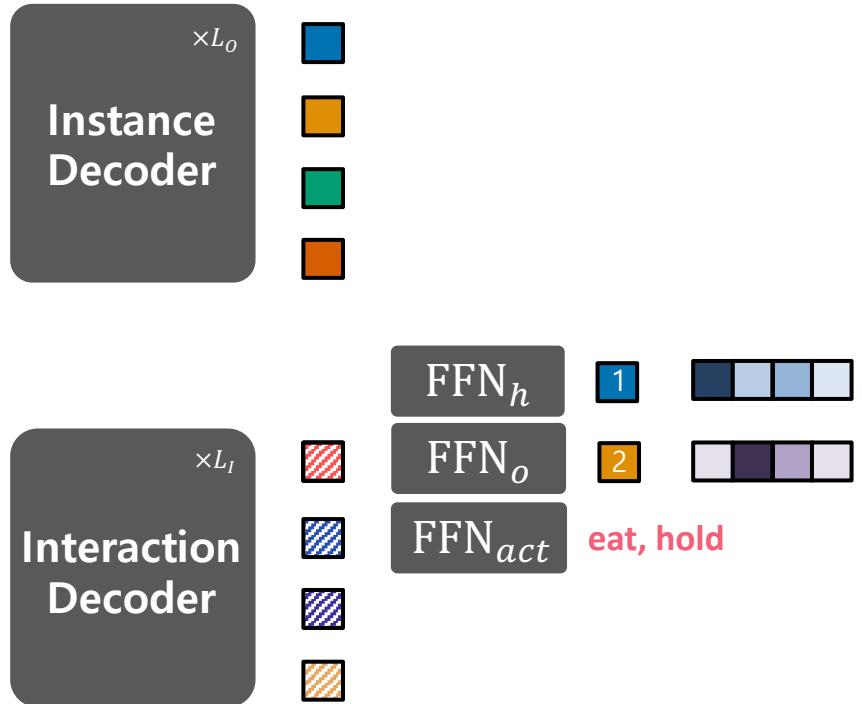
# HO Pointers



# HO Pointers



# HO Pointers



# HO Pointers



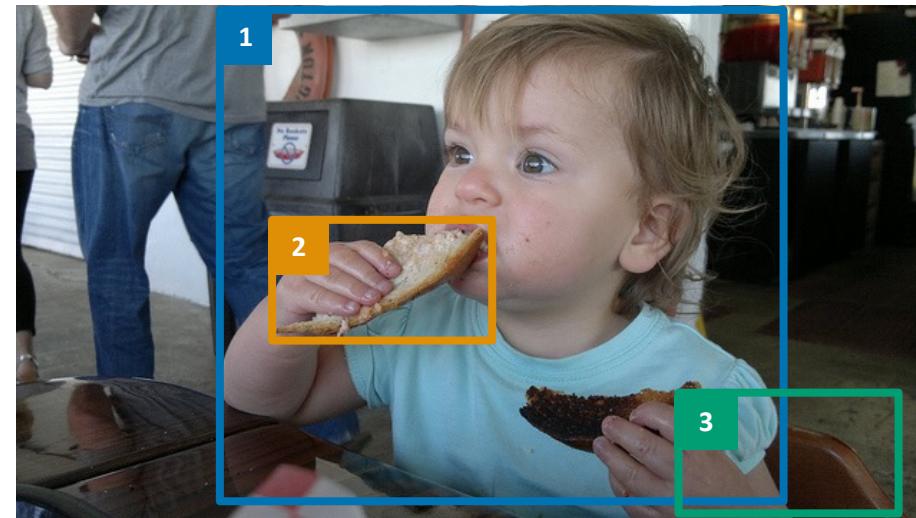
$\text{FFN}_h$

1

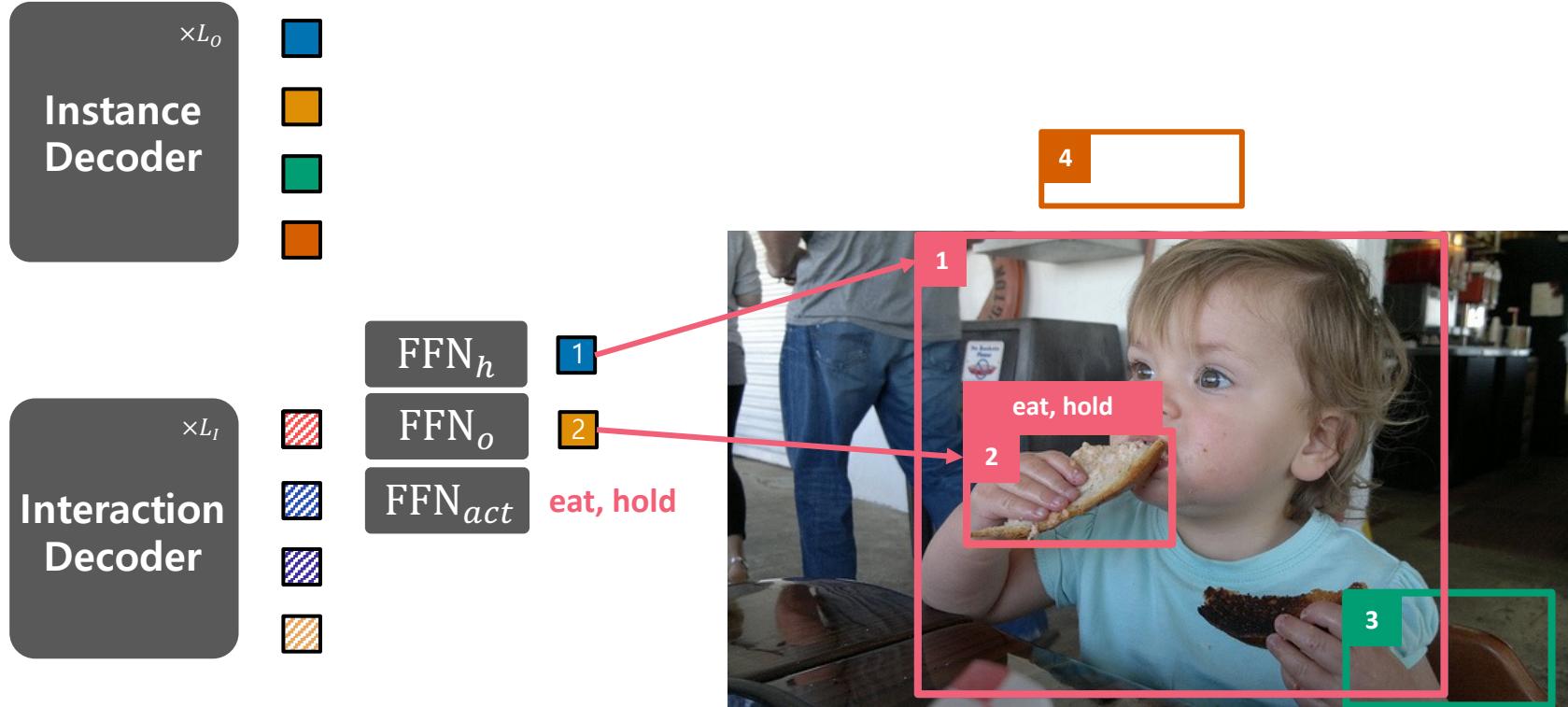
$\text{FFN}_o$

2

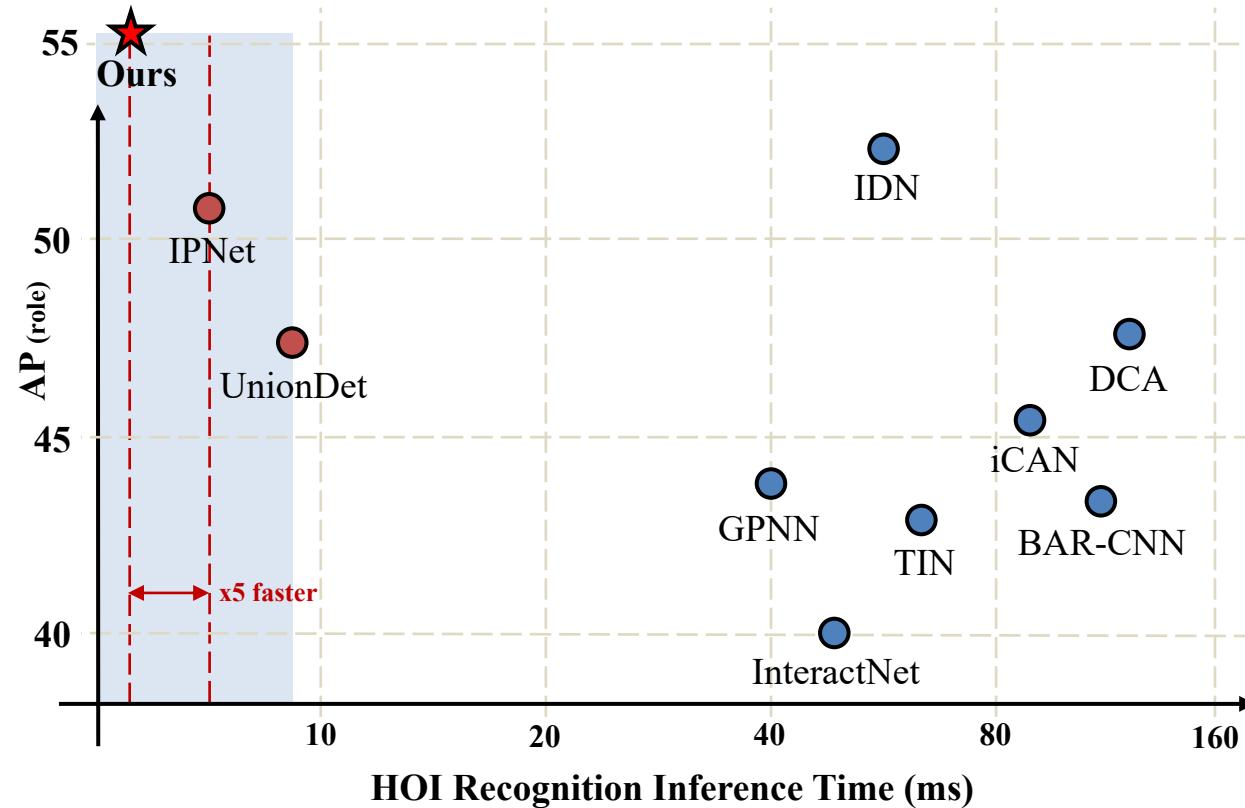
$\text{FFN}_{act}$  eat, hold



# HO Pointers



# Experimental Results (HOI Inference Time)



WBSC

3초

이스라엘 0

대한민국 3

25

1-1

TOKYO 2020

hold  
hit\_objhold  
hit\_inst

야구 본선 2라운드 연장 중계  
김나진·허구연·김선우

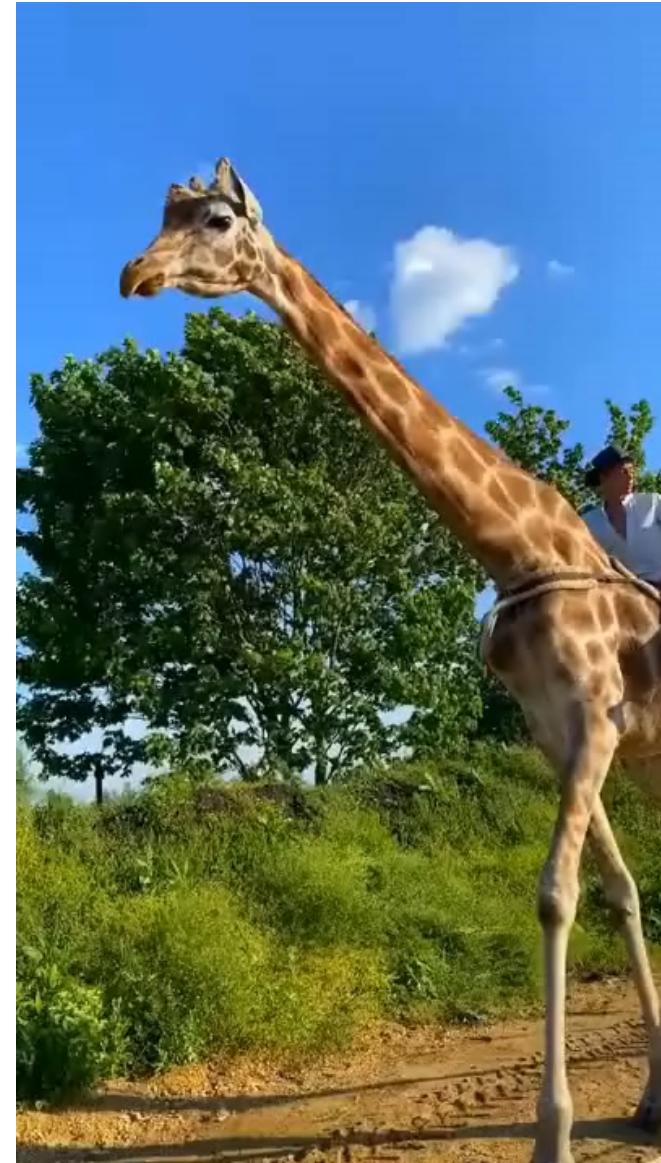


- Ride Giraffe

Work in progress



BCE



Focal+CB

# 참고 문헌

---

- [CVPR 21 (ORAL)] HOTR: End-to-End Human-Object Interaction Detection with Transformers
- [CVPRW 19] ANTNets: Mobile Convolutional Neural Networks for Resource Efficient Image Classification
- [CVPR 18] Tensorize, Factorize and Regularize: Robust Visual Relationship Learning
- [ECCV 20] UnionDet: Union-Level Detector Towards Real-Time Human-Object Interaction Detection

# Questions?

---