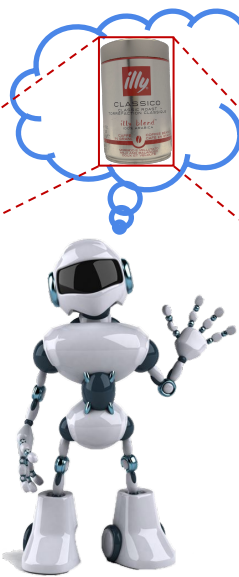# A High-Resolution Dataset for **Ins**tance **Det**ection with Multi-View Object Capture

Qianqian Shen, Yunhan Zhao, Nahyun Kwon, Jeeeun Kim, Yanan Li, Shu Kong

# Instance Detection



**Assistive robots**

locating the **wanted** object at distance!

# Object Detection (ObjDet) *vs.* Instance Detection (InsDet)



coffee-bean,
bottle,
cup,
…

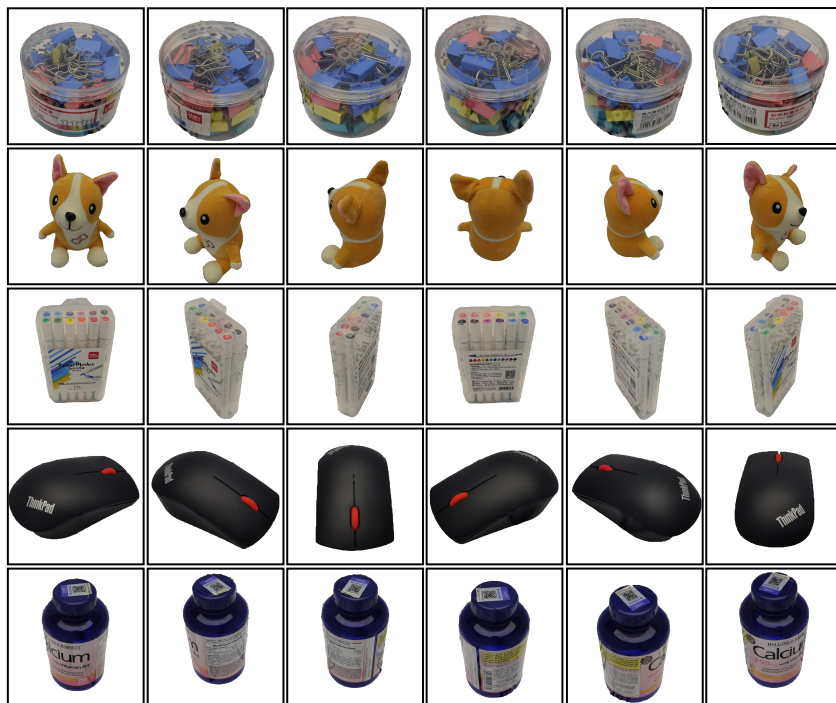**ObjDet aims to detect all objects belonging to some predefined classes.**

**InsDet requires detecting specific object instances defined by some visual examples.**

# Our dataset: **InsDet**
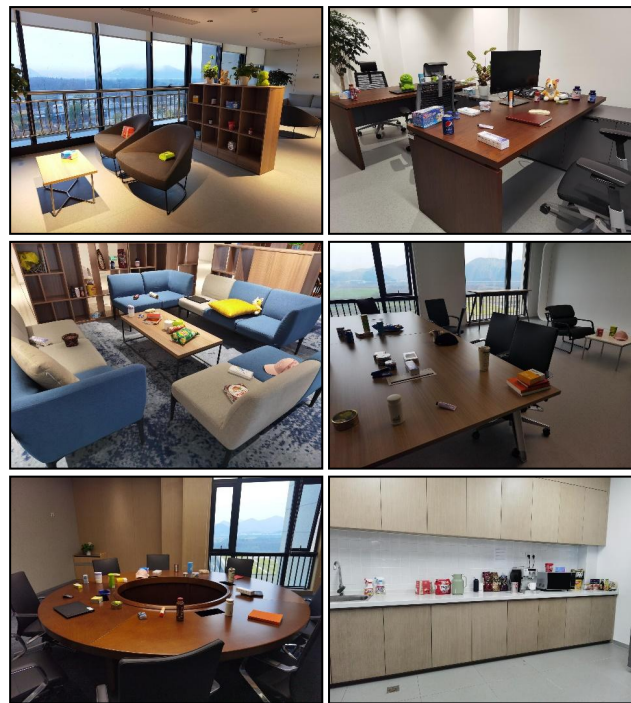
## Instance's profile images

- **100** object instances
- **24** samples per instance



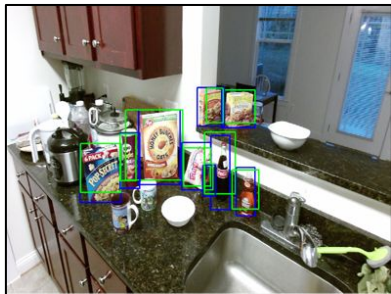**3072×3072**

## Real-World Scenes

- **Diverse** scenes
- **High-resolution** images



**6144×8192**

# Comparison against existing datasets

- **23** instances
- **9** scenes
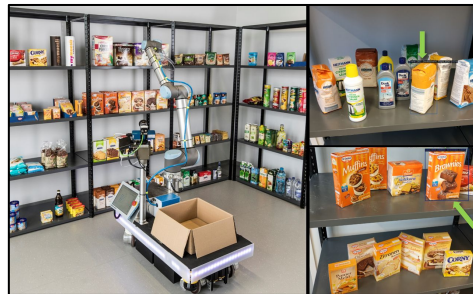- resolution: **1080×1920**
- publicly available



**GMU dataset**

[Georgakis et al. 2016]

- **33** instances
- **9** scenes
- resolution: **1080×1920**
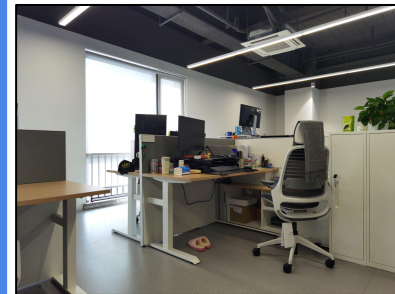- publicly available



**AVD dataset**

[Ammirato et al. 2017]

- **100** instances
- **10** scenes
- resolution: **unknown**
- publicly unavailable



**Grocery dataset**

[Bormann et al. 2021]

- **100** instances
- **14** scenes
- resolution: **6144×8192**
- publicly available



**InsDet dataset**

[Shen et al. 2023]

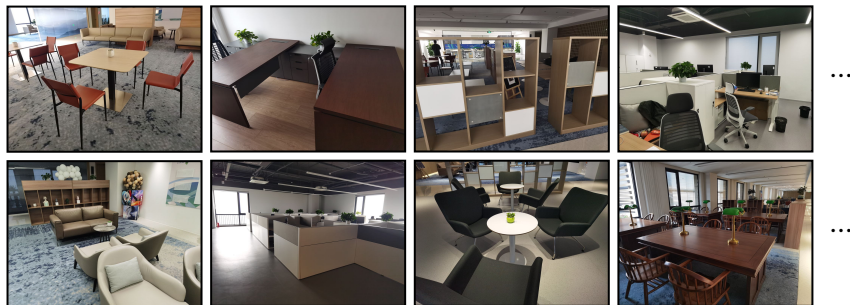# A unified **InsDet** protocol



**Training**

objects captured in various views

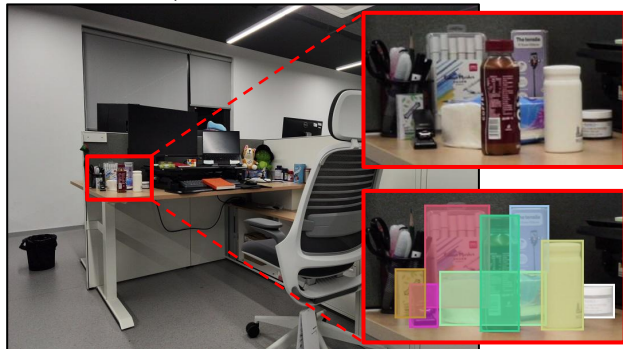indoor scene images (not containing instances of interest)

**Testing**

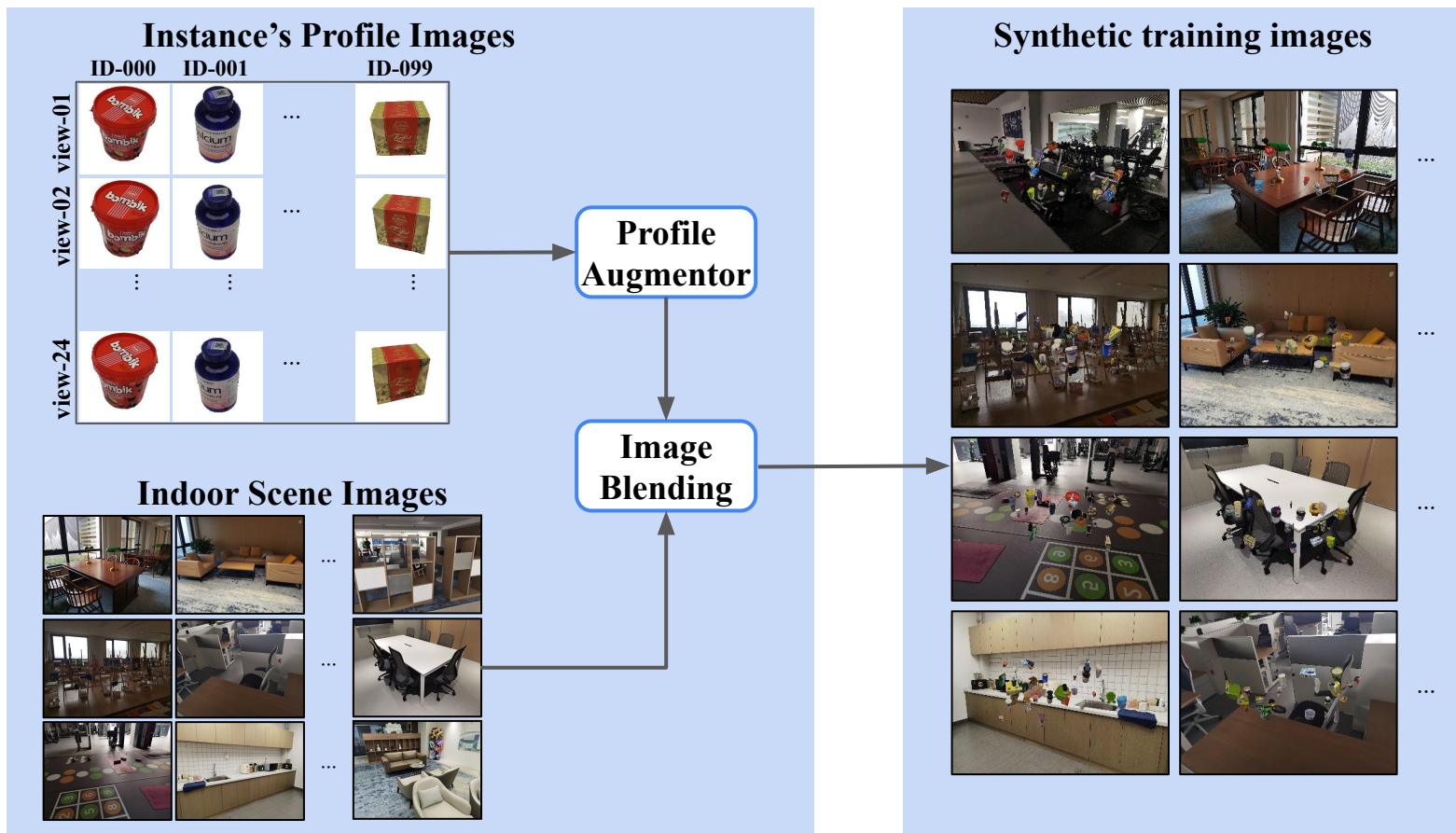*hard* scene (w/ more cluttered environments)

*easy* scene (w/ sparse placement of objects)

# Baseline: Cut-Paste-Learn



## Instance's Profile Images

## Synthetic training images

## Indoor Scene Images

**Profile Augmentor**

**Image Blending**

[1] Cut, paste and learn: Surprisingly easy synthesis for instance detection. In ICCV, 2017.

# A simple, non-learned method



## Proposal Generation

## Feature Extraction

ID-000  ID-001  ID-099

view-01  view-02  ...  view-24

**SAM** ❄️

**DINOv2** ❄️

## Proposal Matching & Selection

**Instances' Features**

# instances

# views

feature vector

| | | | |
|---|---|---|---|
| 0.124 | 0.591 | 0.106 | 0.720 |
| 0.187 | 0.659 | 0.230 | 0.514 |
| 0.084 | 0.157 | 0.820 | 0.092 |
| 0.124 | 0.199 | 0.228 | 0.194 |
| 0.058 | 0.208 | 0.175 | 0.208 |

**Cosine Similarity**

**Stable Matching**

0.659
0.720
0.820

**Proposals' Features**

# proposals
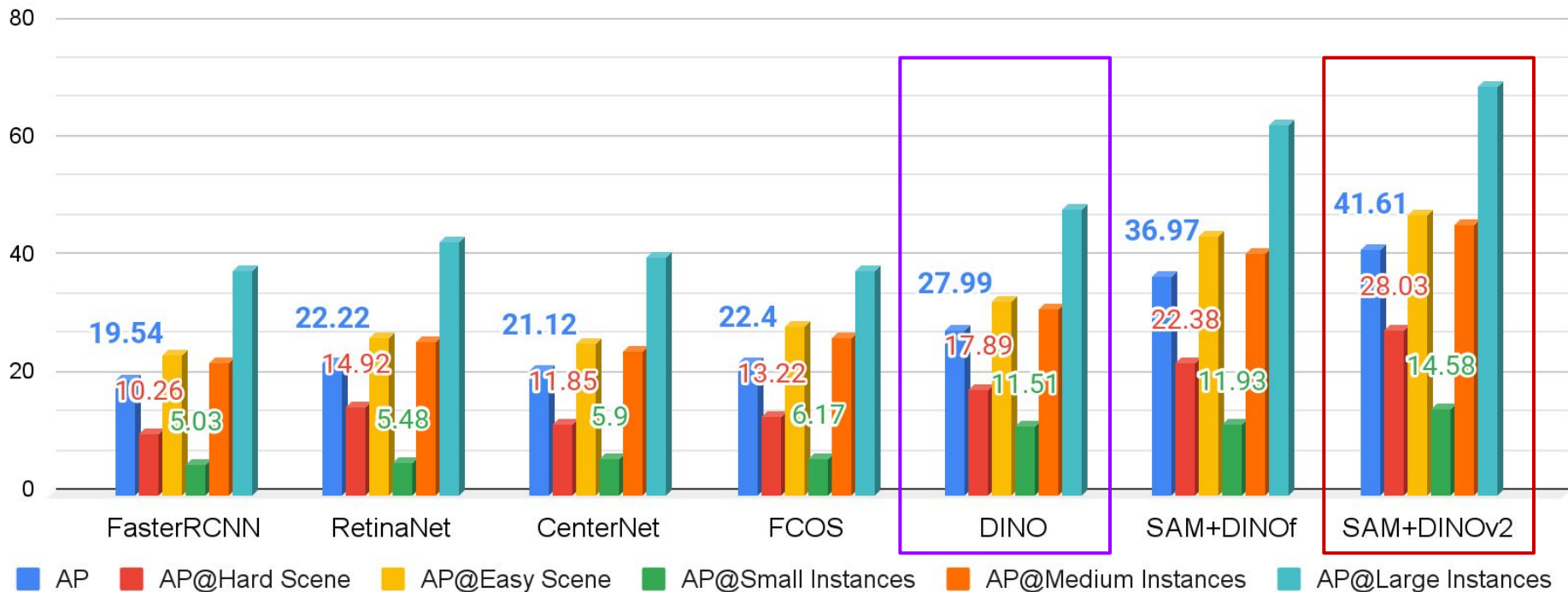
feature vector

[1] Segment anything. In ICCV, 2023.

[2] DINOv2: Learning robust visual features without supervision. In arXiv, 2023

**Benchmarking results**
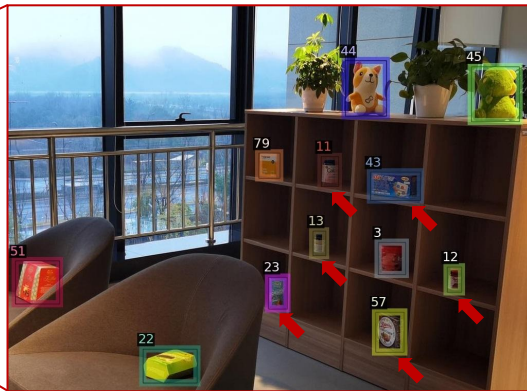
Legend: AP · AP@Hard Scene · AP@Easy Scene · AP@Small Instances · AP@Medium Instances · AP@Large Instances

FasterRCNN: 19.54, 10.26, 5.03
RetinaNet: 22.22, 14.92, 5.48
CenterNet: 21.12, 11.85, 5.9
FCOS: 22.4, 13.22, 6.17
DINO: 27.99, 17.89, 11.51
SAM+DINOf: 36.97, 22.38, 11.93
SAM+DINOv2: 41.61, 28.03, 14.58

# Qualitative evaluations on *easy* scenes

**Easy
(sparse)**

**GroundTruth**



**DINO**



**FasterRCNN**



**SAM + DINOv2**

# Qualitative evaluations on *hard* scenes



**GroundTruth**

**DINO**

**Hard (cluttered)**

**FasterRCNN**

**SAM + DINOv2**

# Suggested further directions

- **Exploring high-resolution images**

  Leverage high-resolution visual signals to help detect small objects.

- **Exploring faster algorithms**

  Building multi-scale detectors for more efficient processing.

- **Exploring more foundational models**

  Learn lightweight adaptors to bridge pretrained foundational models for better performance.

# Thank You!



https://github.com/insdet