# Visual Concept Learning

## Bicycle



...

## Ride bicycle



...

🙁 Need large amount of training data
🙁 Hard to generalize beyond the training concepts

# Bongard-HOI Benchmark

Positive Examples
ride bicycle

Negative Examples
not ride bicycle

Query Images
positive

negative

# Hard Negatives in Bongard-HOI

person
ride
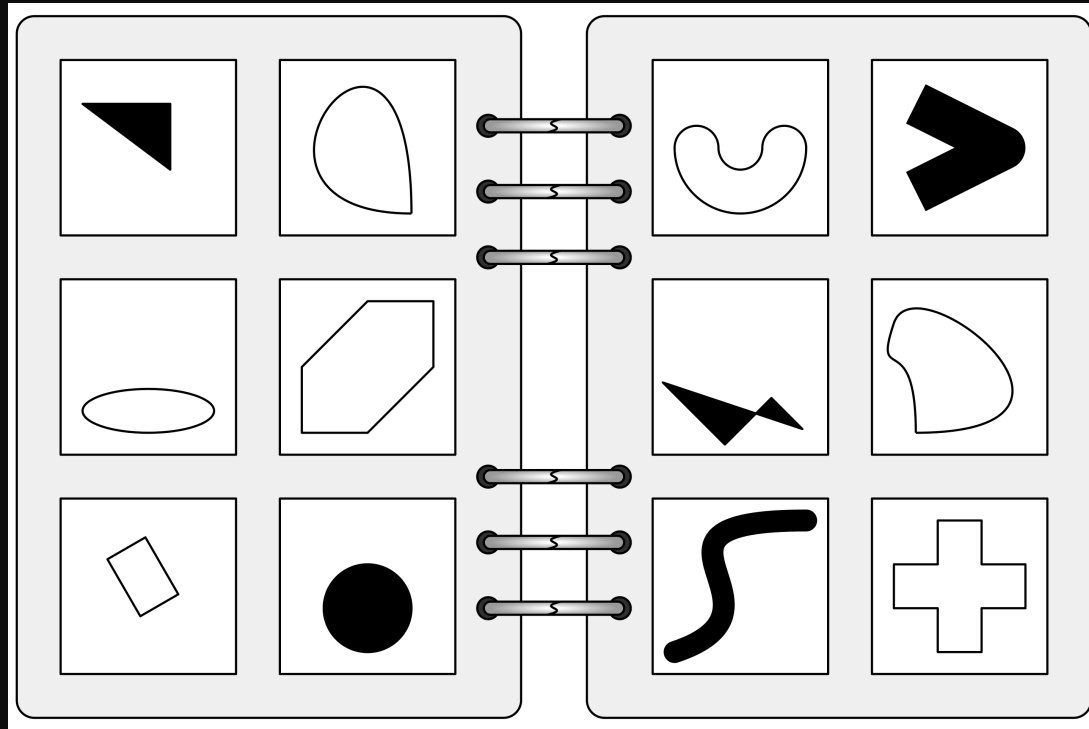bicycle

person
straddle
bicycle

person
repair
bicycle

person
walk
bicycle



Simple visual recognition is not sufficient.
Visual reasoning (e.g., few-shot learning, context reasoning) of the interactions is needed.

# Inspirations from Cognitive Science



Original Bongard problems
[Bongard, 1970.]

Bongard-LOGO
[Nie et al., NeurIPS 2021]

# Different Types of Generalization

sit_on bed     straddle bicycle     hug person     wash car

Training set

wash bicycle     sit_on bench     greet person     shear sheep

Test set

seen action, seen object     seen action, unseen object     unseen action, seen object     unseen action, unseen object

Increasing difficulty

# Context-Dependent Reasoning

## Positive Examples

drink_with cup

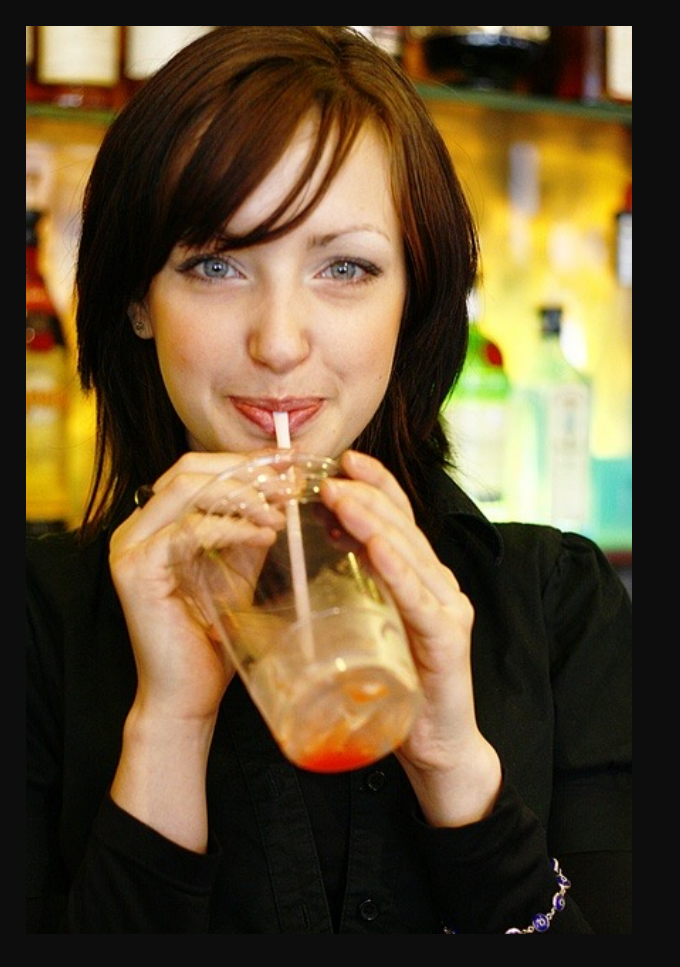








## Hard Negative Examples

not drink_with cup



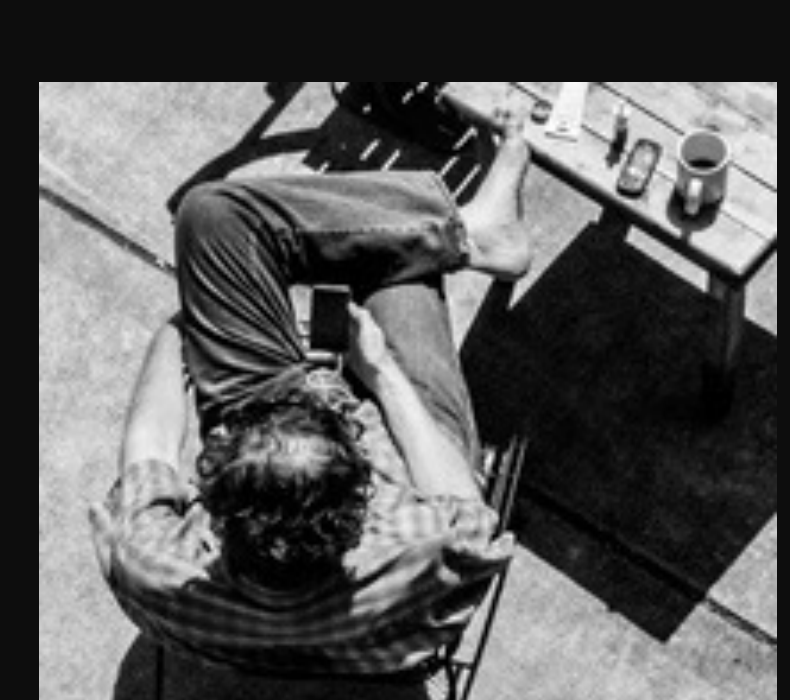



## Query Images

negative

# Context-Dependent Reasoning

Positive Examples
hold cup

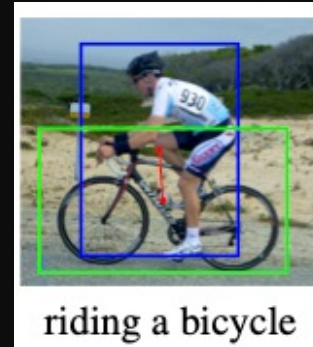<u>Hard</u> Negative Examples
not hold cup
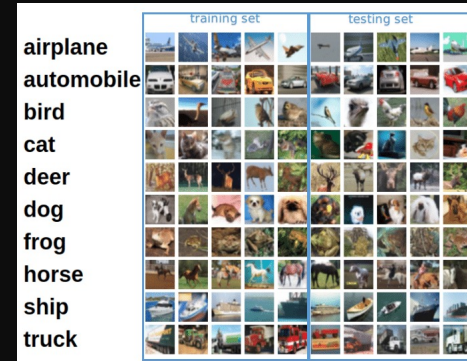
Query Images
positive
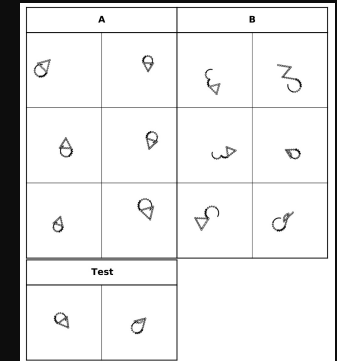
# Comparisons with Other Benchmarks

|  | Bongard-HOI | HOI detection | miniImageNet | Bongard-LOGO |
|---|---|---|---|---|
| Natural images | ✓ | ✓ | ✓ | ✗ |
| Hard negatives | ✓ | ✗ | ✗ | ✗ |
| Compositional concept | ✓ | ✓ | ✗ | ✓ |
| Few-shot learning | ✓ | ✗ | ✓ | ✓ |
| Ctx.-dependent reasoning | ✓ | ✗ | ✗ | ✓ |
| Generalization types | ✓ | ✗ | ✗ | ✓ |

# Meta-Learning for Bongard-HOI



2-way, 6-shot

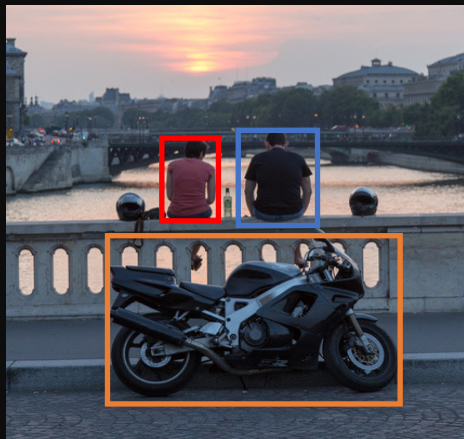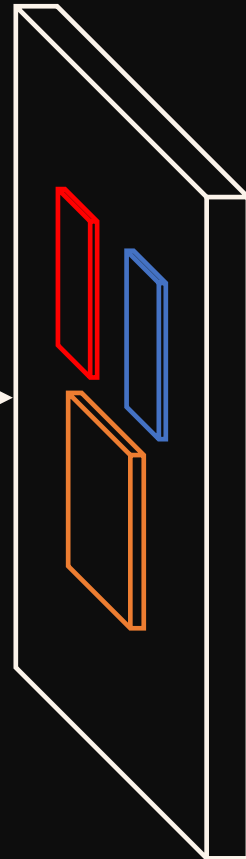mean

cos

loss

[Chen et al., Meta-Baseline. ICCV 2021]
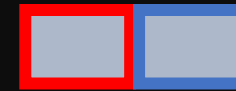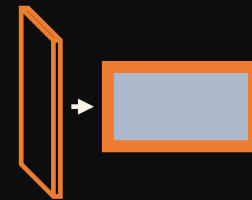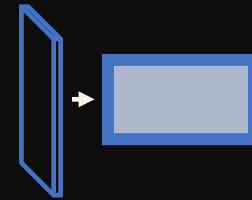
# Image Encoding with Relational Network

Objectness detection
(binary category-agnostic)

- Random initialization (scratch)
- ImageNet pre-training
- MoCo V2 [Chen et al., arXiv, 2021]
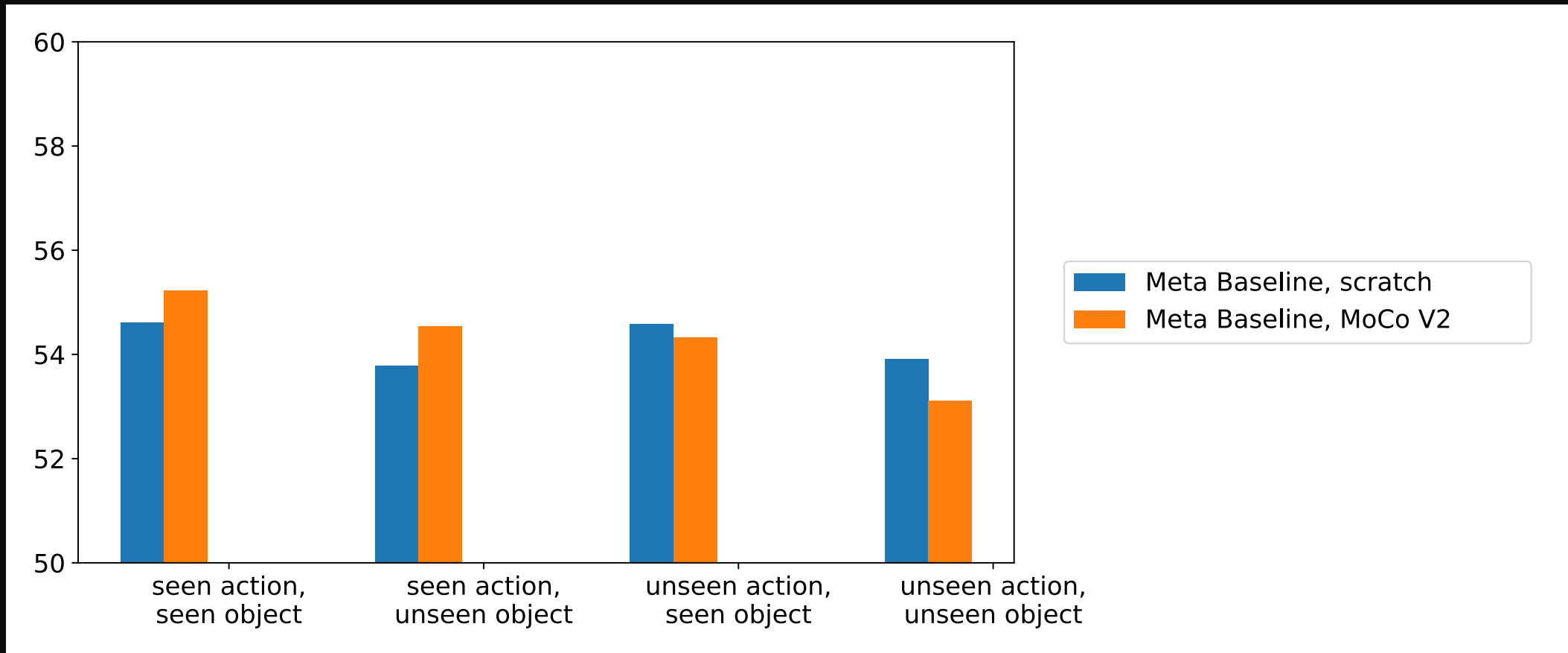
ResNet 50

RoIPool

MLP

Relational encoding for RoI features
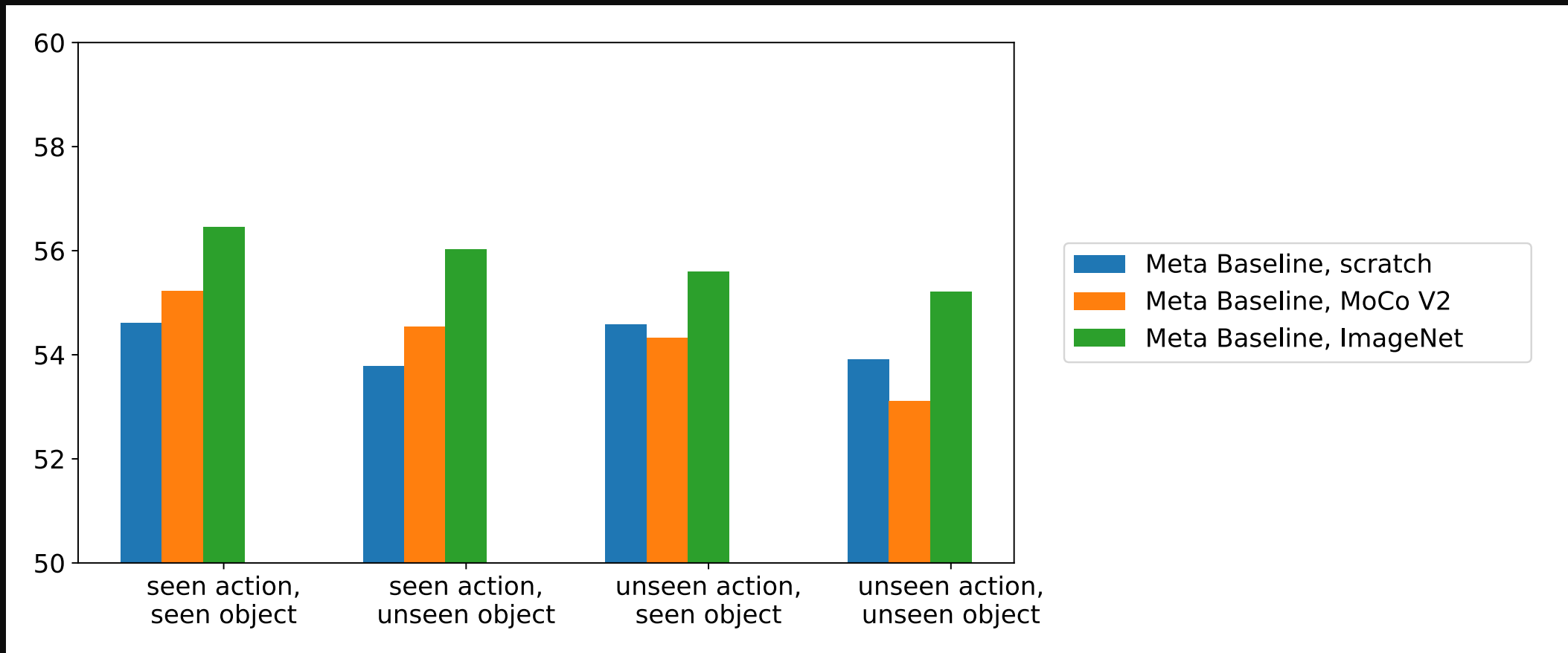
[Santoro et al., NeurIPS 2017]

# Comparisons of Meta Learning Models



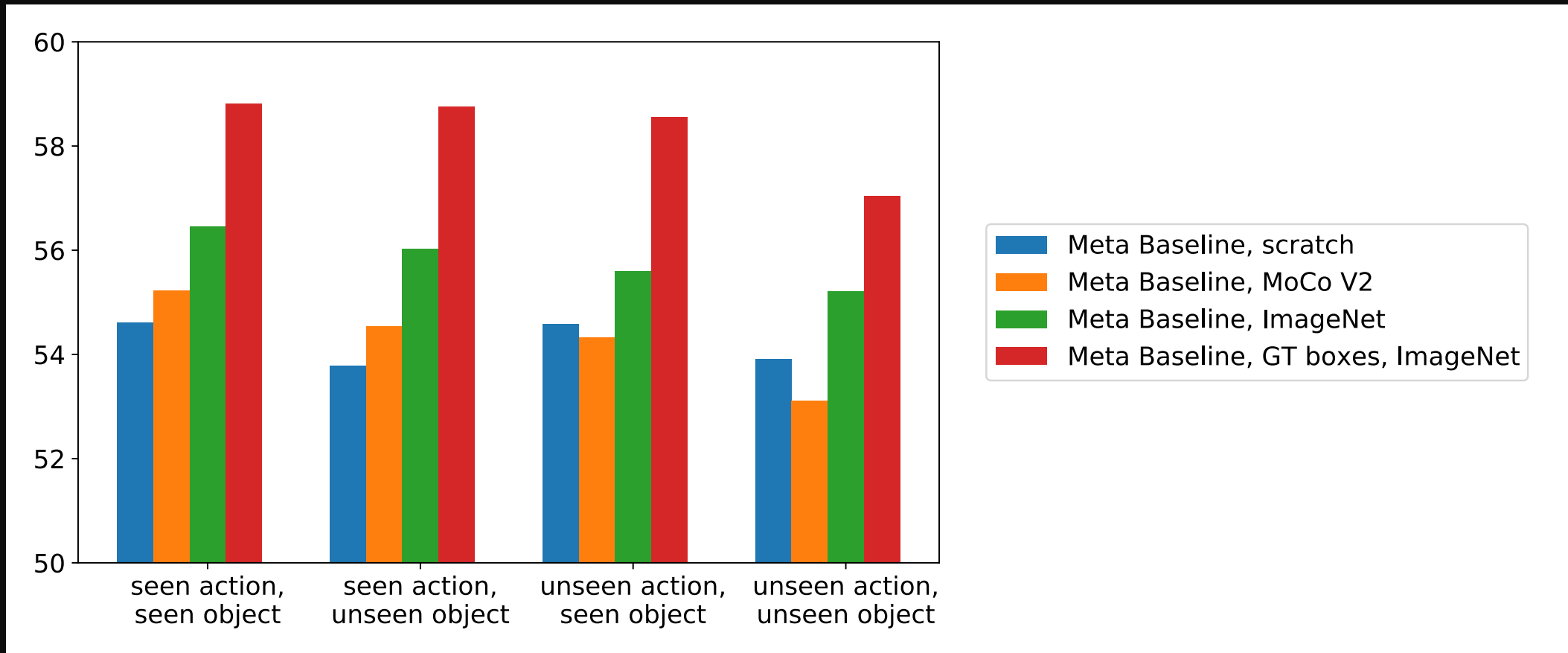Increasing difficulty

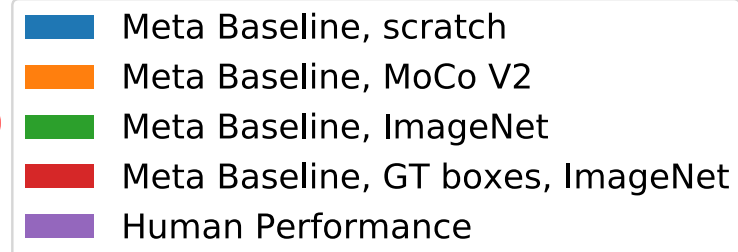# Comparisons of Meta Learning Models

# Comparisons of Meta Learning Models

# Importance of Holistic Visual Perception and Reasoning

# Comparisons with Human Accuracy

# Take-home Messages of Bongard-HOI

- A new benchmark about HOI, highlighting visual reasoning
  - Few-shot learning
  - Context reasoning
  - Generalization beyond training concepts
  - …
- Meta-learning models do not work well enough
  - Pre-training is helpful
  - Holistic visual perception and reasoning is essential
- There exit huge gap w.r.t human performance

17

# Data and Code



Poster: 36b