

Open-Vocabulary Attribute Detection

TUE-PM-278



María A. Bravo



Sudhanshu Mittal



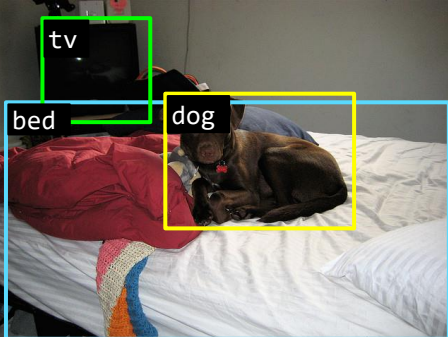
Simon Ging



Thomas Brox

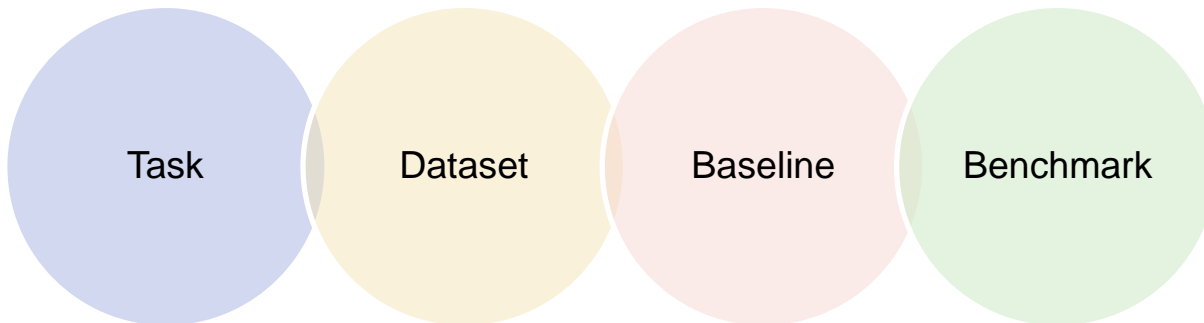
universität freiburg

OVAD: Open-vocabulary Attribute Detection

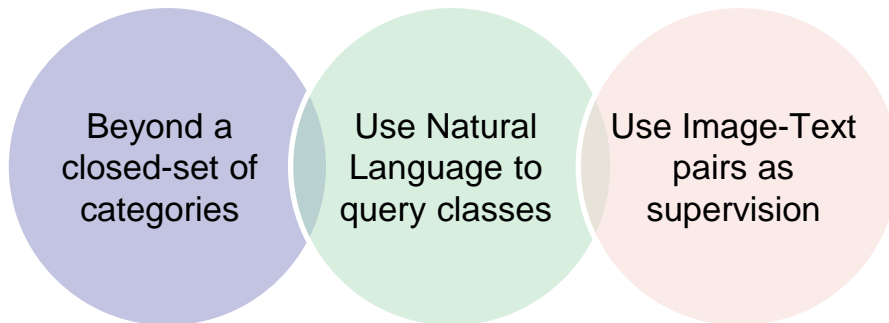
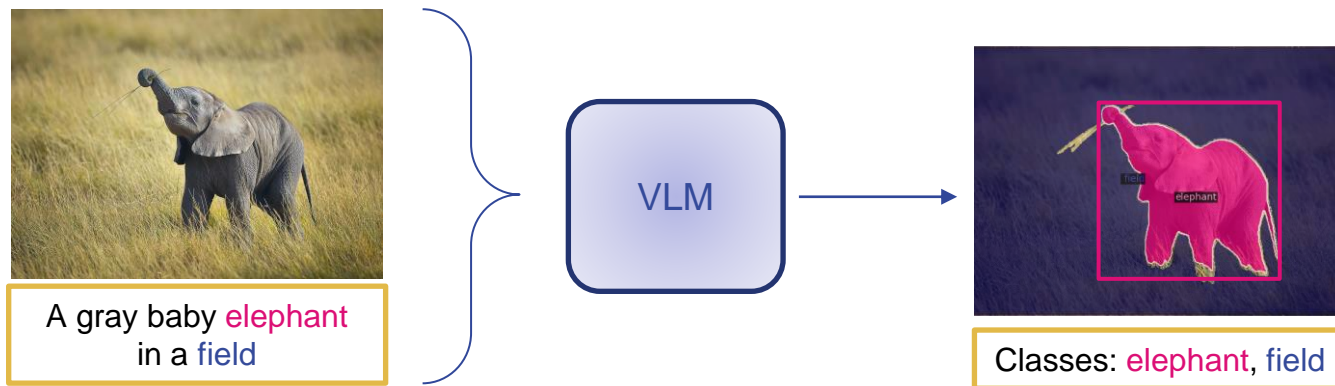


Object class: dog Attributes color quantity: single color: brown group: single maturity: adult position: lying state: dry texture: soft tone: dark	Object class: bed Attributes color quantity: single color: white group: single material: cloth optical prop: opaque size: large texture: smooth tone: light	Object class: tv Attributes color quantity: single color: black material: glass optical prop: reflective pattern: plain state: off texture: smooth tone: dark
--	---	---

■ base
■ novel



Vision Language Models for Open Vocabulary Recognition



Vision Language Models for Open Vocabulary Recognition

Noun concepts

Visual Recognition
Supervised vs Open-vocabulary

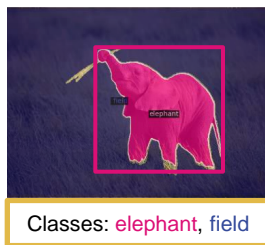
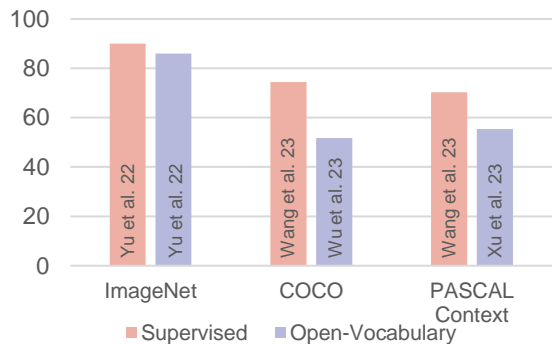


Image
Classification

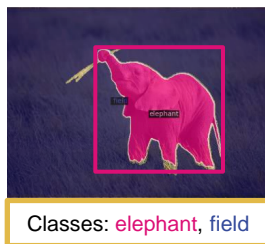
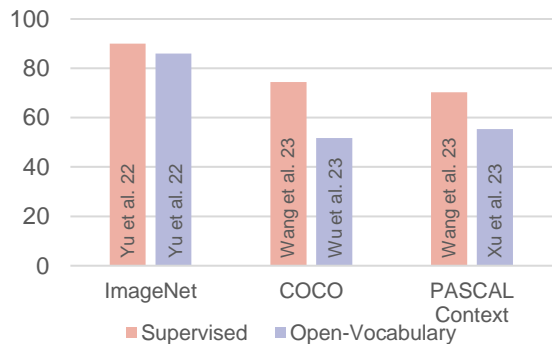
Object
Detection

Semantic
Segmentation

Vision Language Models for Open Vocabulary Recognition

Noun concepts

Visual Recognition
Supervised vs Open-vocabulary



Attribute concepts

A **gray** **baby** elephant
in a field



Object class: elephant
Attributes
color: gray
quantity: one
group: single
maturity: baby
pattern: plain
position: vertical
size: small
state: dry
texture: rough
tone: light

Image
Classification

Object
Detection

Semantic
Segmentation

Significance of attributes in an object's identity



A cup with **three** pairs of scissor, one is **blue**, and **two** are **metallic**.



A **white** and a **spotted** horse on a field of grass.



Red traffic signal in the middle of a **wide** street.

OVAD task: Open-vocabulary Attribute Detection

Object class: **bear**
Attributes
color quantity: **two**
color: **black and brown**
group: **single**
material: **wood**
maturity: **adult**
position: **upright**
size: **big**
state: **dry**
texture: **smooth**
tone: **dark**



Object class: **car**
Attributes
color quantity: **one**
color: **red**
group: **single**
material: **wood**
optical prop: **opaque**
patterns: **lettered**
state: **piece / cut**
texture: **smooth**
tone: **dark**

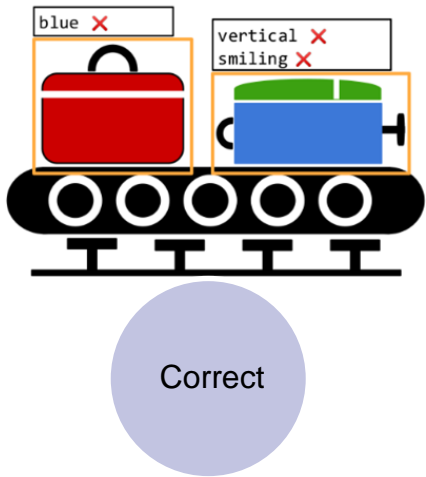
Object class: **person**
Attributes
face exp: **surprise**
group: **single**
hair color: **black**
hair length: **short**
hair tone: **dark**
hair type: **straight**
maturity: **adult**
position: **upright**
clothes color: **white**

Solving the **OVAD** task: → (1) Detect all **object classes**
(2) Detect their **attributes**.

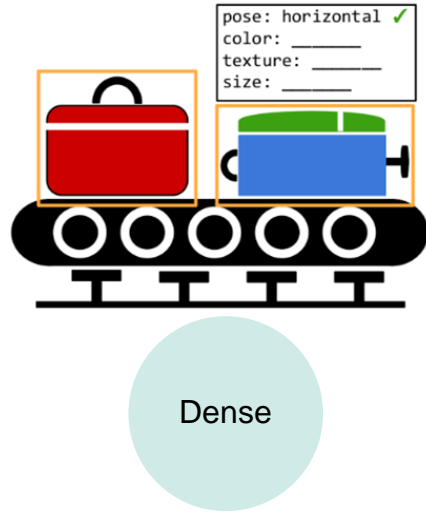
Attribute Detection Benchmarks

Type of errors:

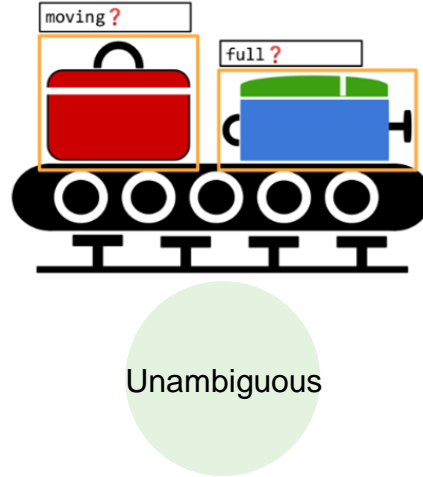
(a) Incorrect ✘



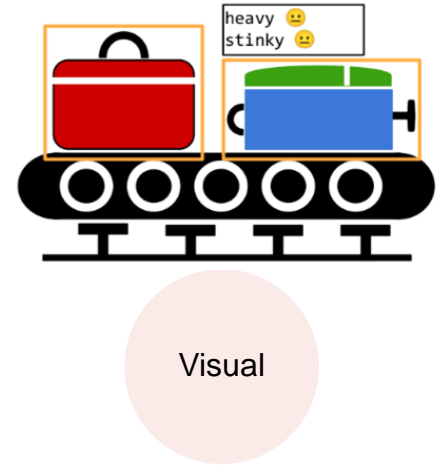
(b) Missing —



(c) Ambiguous ?



(d) Non-visual 😞



OVAD Benchmark

Open-vocabulary Attribute Detection Dataset Visualization

Dataset overview

Back to main page

Dense

Unambiguous










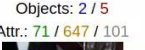


Correct

Visual

Showing page 1 of 100 (2000 images total).

First Previous Next Last

Filter with show

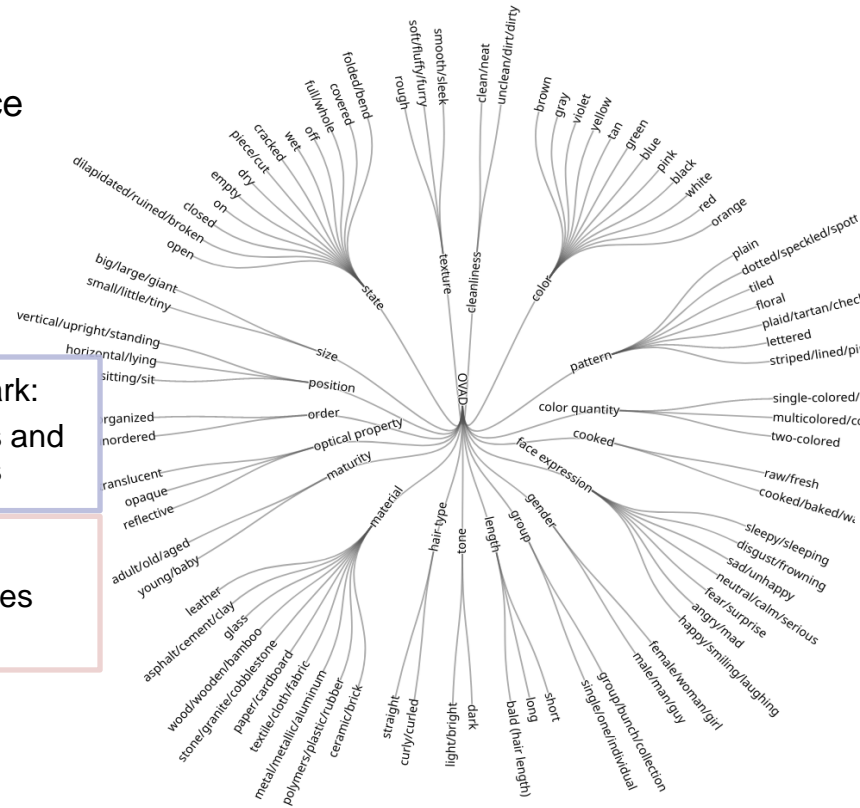
<p>2 / 251</p>  <p>Objects: 1 / 0 Attr.: 12 / 95 / 10</p>	<p>Objects: 1 / 2 Attr.: 31 / 253 / 67</p>  <p>Objects: 2 / 2 Attr.: 34 / 343 / 91</p>	<p>Objects: 8 / 1 Attr.: 46 / 700 / 307</p>  <p>Objects: 9 / 1 Attr.: 106 / 919 / 145</p>	<p>Objects: 3 / 1 Attr.: 35 / 337 / 96</p>  <p>At</p>	
<p>1 59</p>  <p>Objects: 3 / 2 Attr.: 23 / 364 / 198</p>	<p>Objects: 2 / 0 Attr.: 20 / 186 / 28</p>  <p>Objects: 3 / 1 Attr.: 45 / 349 / 74</p>	<p>Objects: 2 / 0 Attr.: 18 / 206 / 10</p>  <p>Objects: 4 / 1 Attr.: 31 / 409 / 145</p>	<p>Objects: 0 / 16 Attr.: 145 / 1358 / 369</p>  <p>A</p>	
	<p>Objects: 16 / 9 Attr.: 213 / 2131 / 581</p> 	<p>Objects: 2 / 5 Attr.: 71 / 647 / 101</p> 	<p>Objects: 5 / 0 Attr.: 30 / 407 / 148</p> 	<p>Objects: 2 / 0 Attr.: 15 / 175 / 44</p> 

OVAD Benchmark

Human annotated
~ 96.8 attributes / instance

Evaluation modes:

1. Full benchmark:
 - Detect objects and their attributes
2. Box-oracle:
 - Detect attributes given box



1.4M
Attribute
Annotations

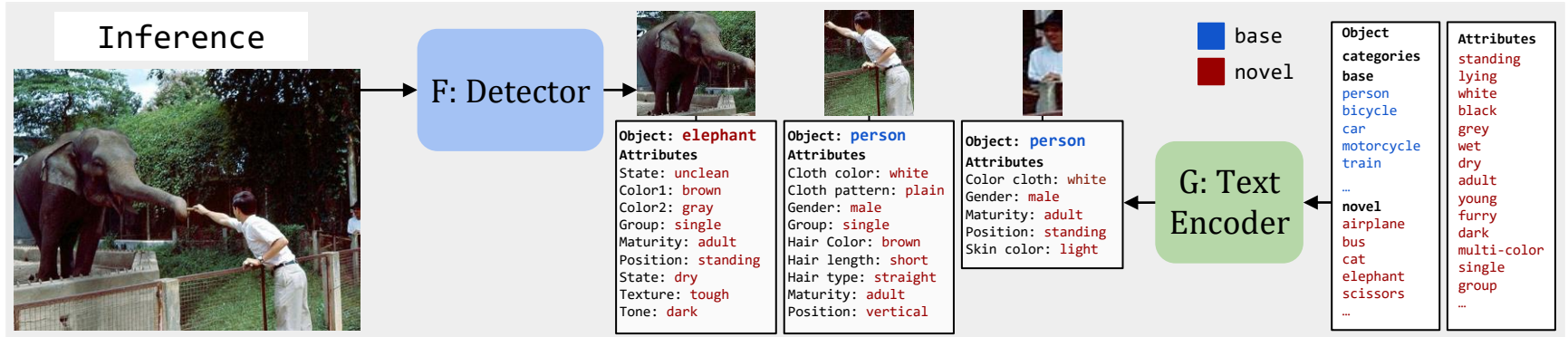
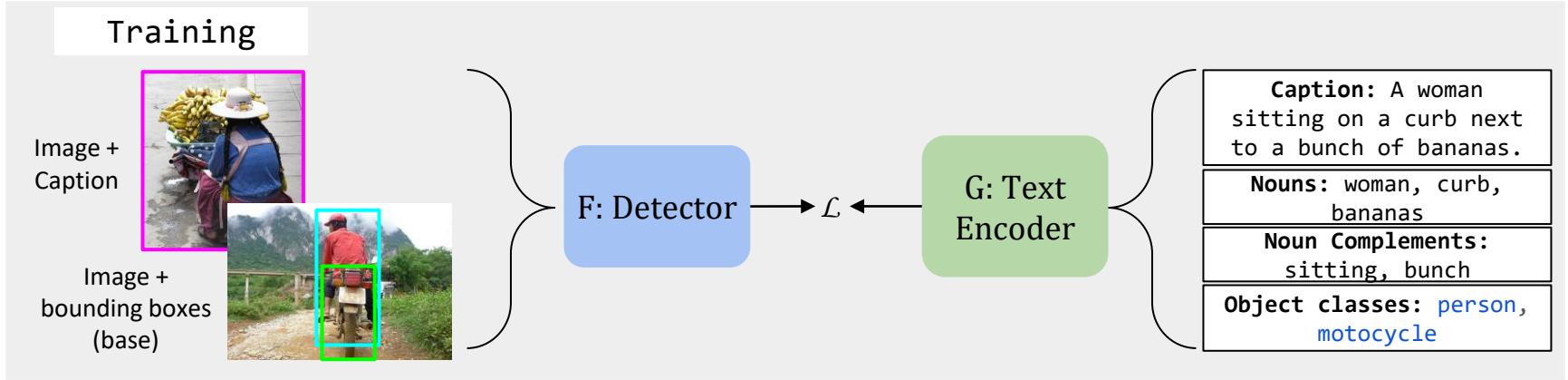
2,000 Test
Images

117
Attribute
Categories

80 Object
Categories

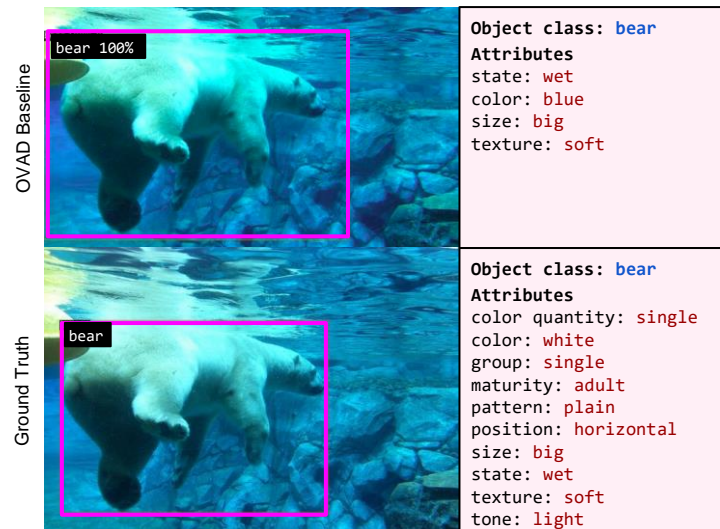
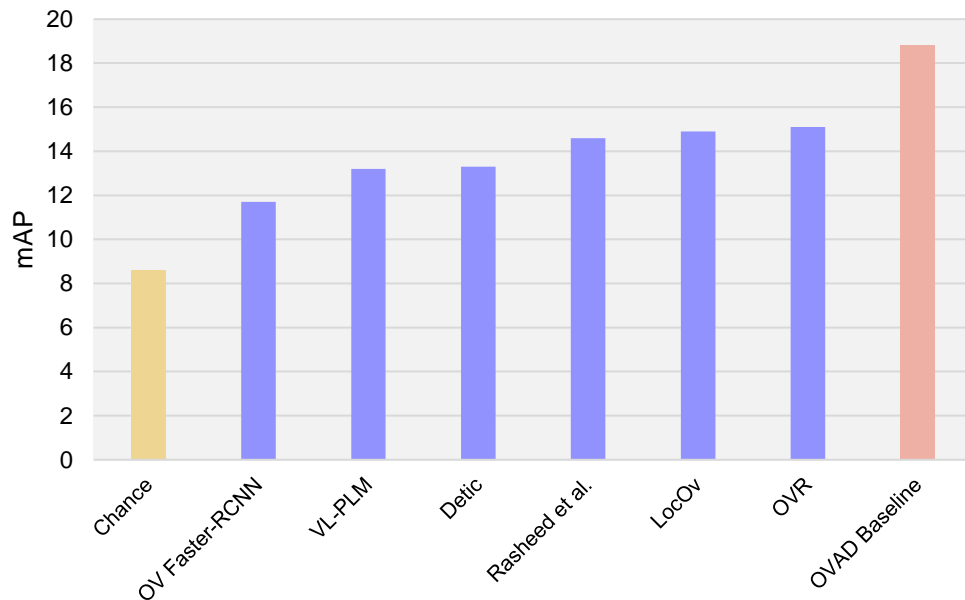
14.3k
Object
Instances

OVAD Baseline



Open-vocabulary Models on OVAD

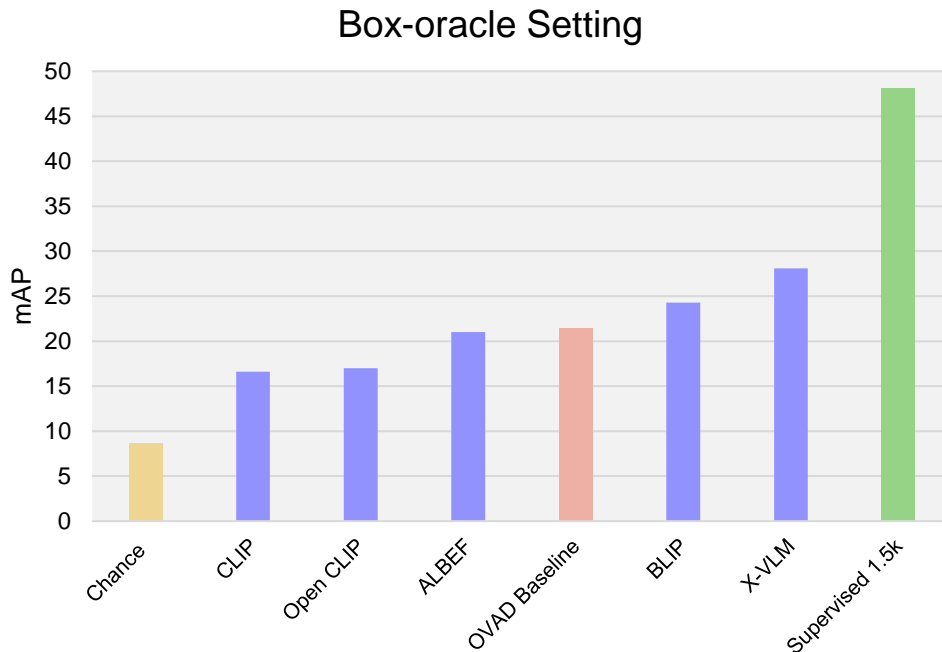
Full Evaluation Setting



➤ Our OVAD Baseline outperforms latest OVD models on the OVAD task.

Performance of VLMs on OVAD

Method	mAP OVAD-Box oracle
Chance	8.6
CLIP ViT-B16 (400M)	16.6
Open CLIP ViT-B32 (2B)	17.0
ALBEF (14M+COCOft)	21.0
OVAD Baseline (110k+COCO)	21.4
BLIP (129M+COCOft)	24.3
X-VLM (16M+COCOft)	28.1
Supervised (OVAD-1.5k)	48.2



- X-VLM model performs the on the OVAD task in the box-oracle setting.
- Large pre-trained VLMs fail to capture fine-grained information.

Open-Vocabulary Attribute Detection

TUE-PM-278

