

# Vision-language Models

Internship Presentation

10 May 2021 – 9 Jul 2021

Akkapaka Saikiran

CSE, IIT Bombay

Mentors: Swati Tiwari, Neelesh Khanna, Niranjana Paramkusum

# Outline

Problem statement

Current pipeline

Baseline scores

Multimodal learning

- Pre-Oscar
- Oscar
- VinVL

Results

- Object detection examples
- VinVL scores

Future work

# The Problem Statement


➤ Product Ads (PA) classification – adult, weapon

**\$74.95**  
Cricket Best Buy  
GN Batting Cricket Shoes Players  
Rubber Sole US 12 / White

**\$59.95**  
Cricket Best Buy  
GM Maestro Multi Function Cricket  
Shoes Metal & Rubber Spikes US 10....

**\$68.03**  
Amazon.com  
Free shipping  
Payntr V Pimple - White & Blue Crick...

**\$57.95**  
Cricket Best Buy  
Payntr X Cricket Shoes White & Yellow  
Pimple Rubber Sole US 12



A white and blue cricket shoe with metal spikes on the sole. The shoe is shown from a side profile, with the sole visible. The sole is blue with white spikes. The shoe has a white upper with blue accents and the GM logo.

### GM Maestro Multi Function Cricket Shoes Metal & Rubber Spikes US 10.5 / White

**\$59.95** Cricket Best Buy

Visit site

Product details

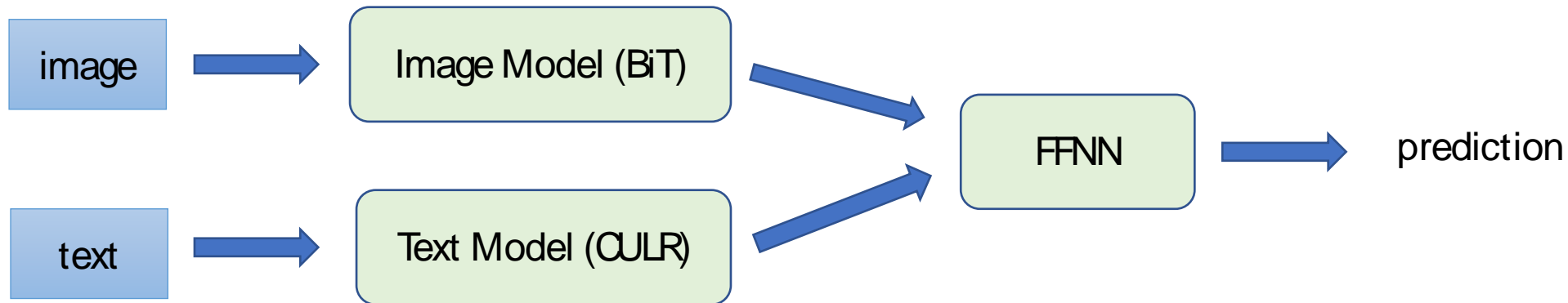
Cricket Shoes Maestro Multi Function Cricket Shoe GM. LIGHTWEIGHT UPPER XLO air mesh provides supreme vamp ventilation Three dimensional comfort lining for superior, long lasting fit and feel Ergonomic slip lasted construction for consistent sock-like fit time after time Moulded TPR heel cradle locks foot in place maintaini... +

Track price ☐

Receive email and browser notifications if the price drops

# The Problem Statement

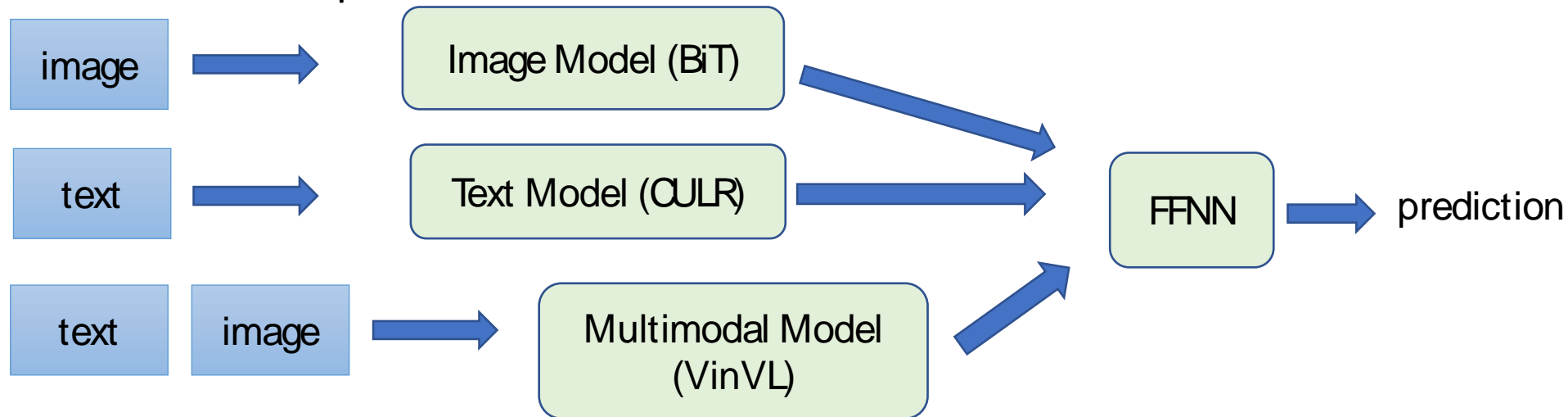
- Product Ads (PA) classification – adult, weapon
- PA - Image (raw and thumbnail)
  - Text (Product name, Merchant name, Description, etc.)
  - Labels (TextDisallowed, ImageDisallowed, OverallAdDisallowed)
- Unimodal pipeline



# The Problem Statement

- Product Ads (PA) classification – adult, weapon
- PA - Image (raw and thumbnail)
  - Text (Product name, Merchant name, Description, etc.)
  - Labels (TextDisallowed, ImageDisallowed, OverallAdDisallowed)

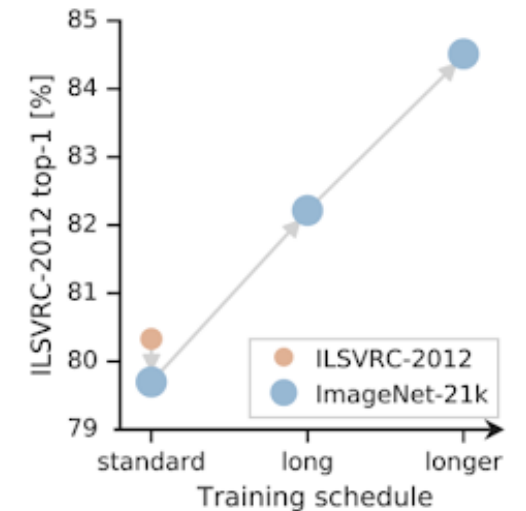
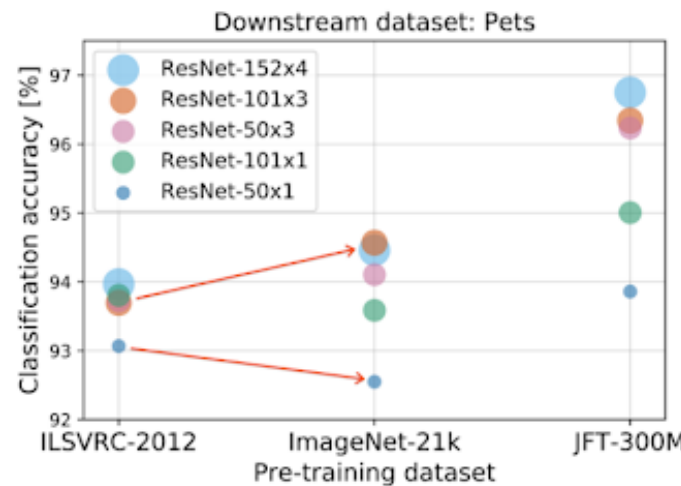
## ➤ Multimodal Pipeline:



# BiT (Big Transfer)<sup>1</sup>

- A pre-training recipe
- Large datasets, large models
- Replace Batch Normalisation with Group Normalisation and Weight Standardization
- We use BiT-M-R152x2

Model	ILSVRC2012	ImageNet Real
SOTA	90.45	91.12
BiT-L	87.54	90.54
BiT-M	85.39	89.02
BiT-S	81.30	86.21



<sup>1</sup> Big Transfer (BiT): General Visual Representation Learning <https://arxiv.org/pdf/1912.11370.pdf>  
Images source: <https://ai.googleblog.com/2020/05/open-sourcing-bit-exploring-large-scale.html>

# Baseline scores

Model	Test Target	Adult PRAUC	Weapon PRAUC
Text	Text label	0.9200	0.9477
Text	Overall label	0.7539	0.9462
Image	Image label	0.9328	0.8813
Image	Overall label	0.8755	0.8651
Combined	Overall label	0.9522	0.9381

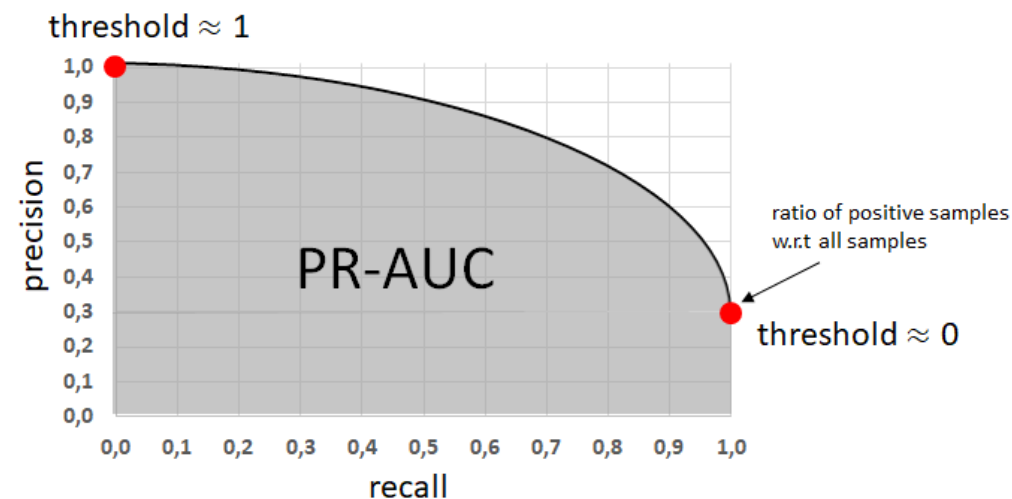


Image source: <https://towardsdatascience.com/gaining-an-intuitive-understanding-of-precision-and-recall-3b9df37804a7>

## Details

- Text model: CULR-Large v1, a multilingual BERT-like model
- Data: ~8.6% adult, ~6.3% weapon

# Multimodal Learning

## Image captioning



The man at bat readies to swing at the pitch while the umpire looks on.

Image source: COCO Captions 2015, <https://cocodataset.org/#captions-2015>

## Visual Question Answering



What color are her eyes?  
What is the mustache made of?



How many slices of pizza are there?  
Is this a vegetarian pizza?

Image source: VQA Dataset, <https://visualqa.org/>

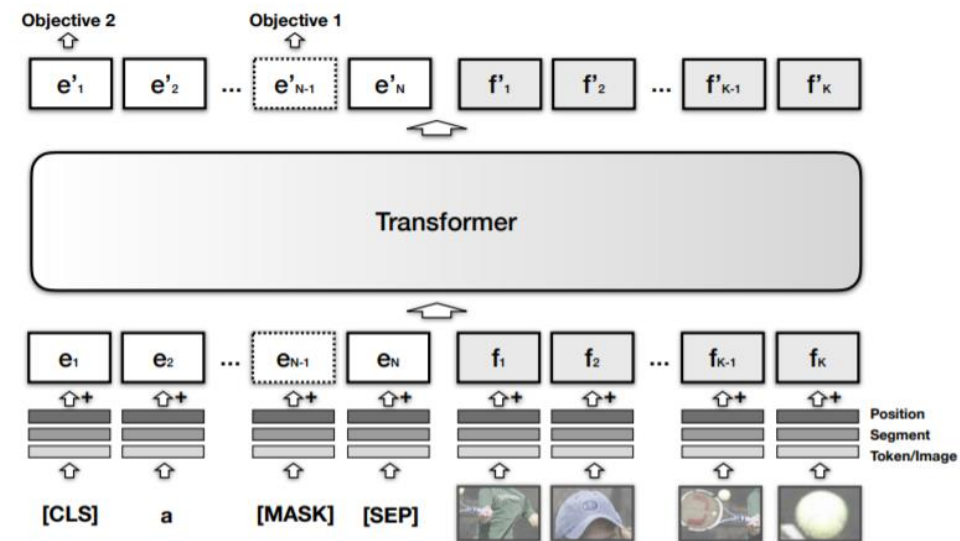


# Multimodal Learning

- Learn cross-modal representations by large-scale pre-training
- Concatenate image region features and text features to form input
- Use self-attention to implicitly learn image-text semantic alignments
- Image region features extracted from pre-trained Object Detection (OD) models
- Self-supervised pre-training objectives
  - Masked Language Modeling (MLM)
  - Image-caption matching

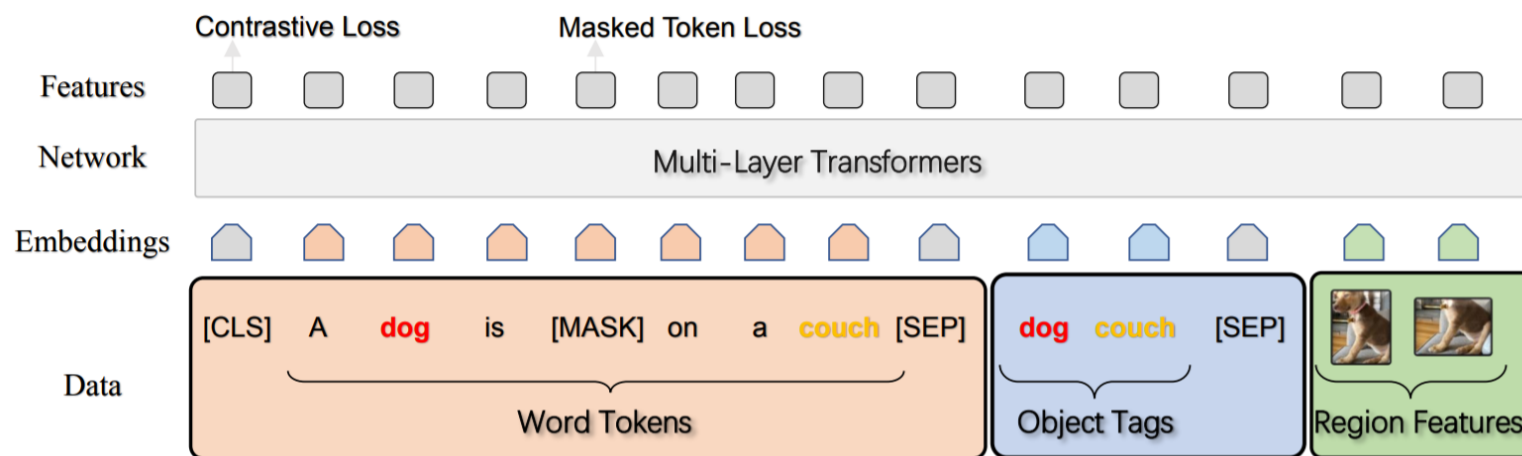
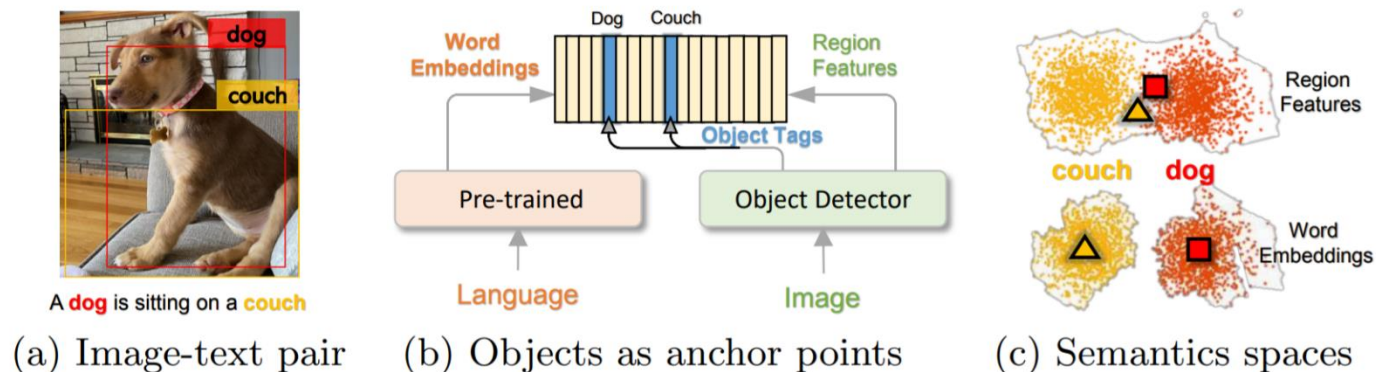


A person hits a ball with a tennis racket



# Oscar<sup>1</sup>

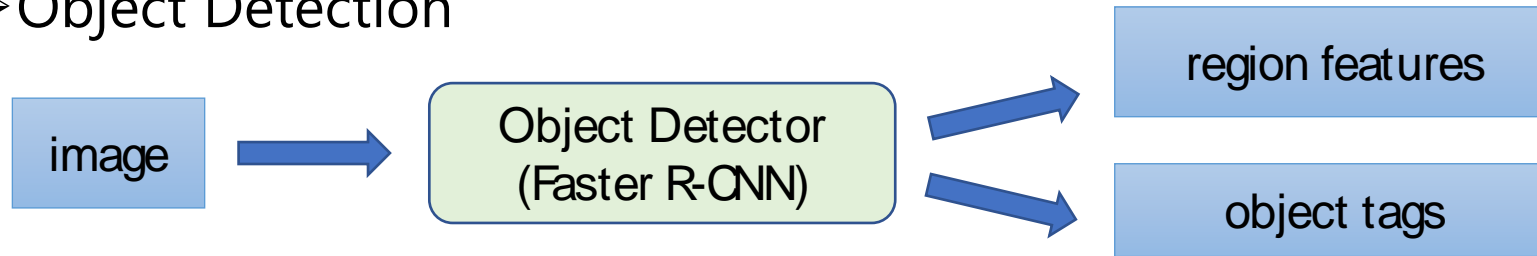
- Key idea – use object tags as *anchor points* to ease learning of alignment
- Motivation – salient objects detected in images are often mentioned in the paired text



<sup>1</sup> Oscar: Object-Semantics Aligned Pre-training for Vision-Language Tasks, <https://arxiv.org/pdf/2004.06165.pdf>

# Oscar: Input Representation

## ➤ Object Detection



## ➤ Input – each image-text pair is represented as a tuple $(w, q, v)$

- $w$  :- sequence of word embeddings of the text
- $q$  :- sequence of word embeddings of the object tags
- $v$  :- sequence of region features

## ➤ Idea (to summarize): Both $q$ and $w$ share the same semantic space (BERT initialisation), so their alignments are easy to identify

## ➤ So image regions $v$ corresponding to relevant object tags are likely to have higher attention weights when queried by related words in $q$

# Oscar: Pre-training Objectives

## Masked Token Loss (Dictionary view) - $L_M$

- Mask each input token in  $h = \text{concat}(w, q)$  with 15% probability
- Goal of training – predict masked tokens using surrounding (partial) text and image context

## Contrastive loss (Modality view) - $L_C$

- Sample a set of “polluted” image representations by replacing  $q$  with 50% probability with a randomly sampled tag sequence  $q'$
- Idea – Utilize object tags as a proxy for images
- Goal of training – predict whether triplet is polluted (using the [CLS] token)

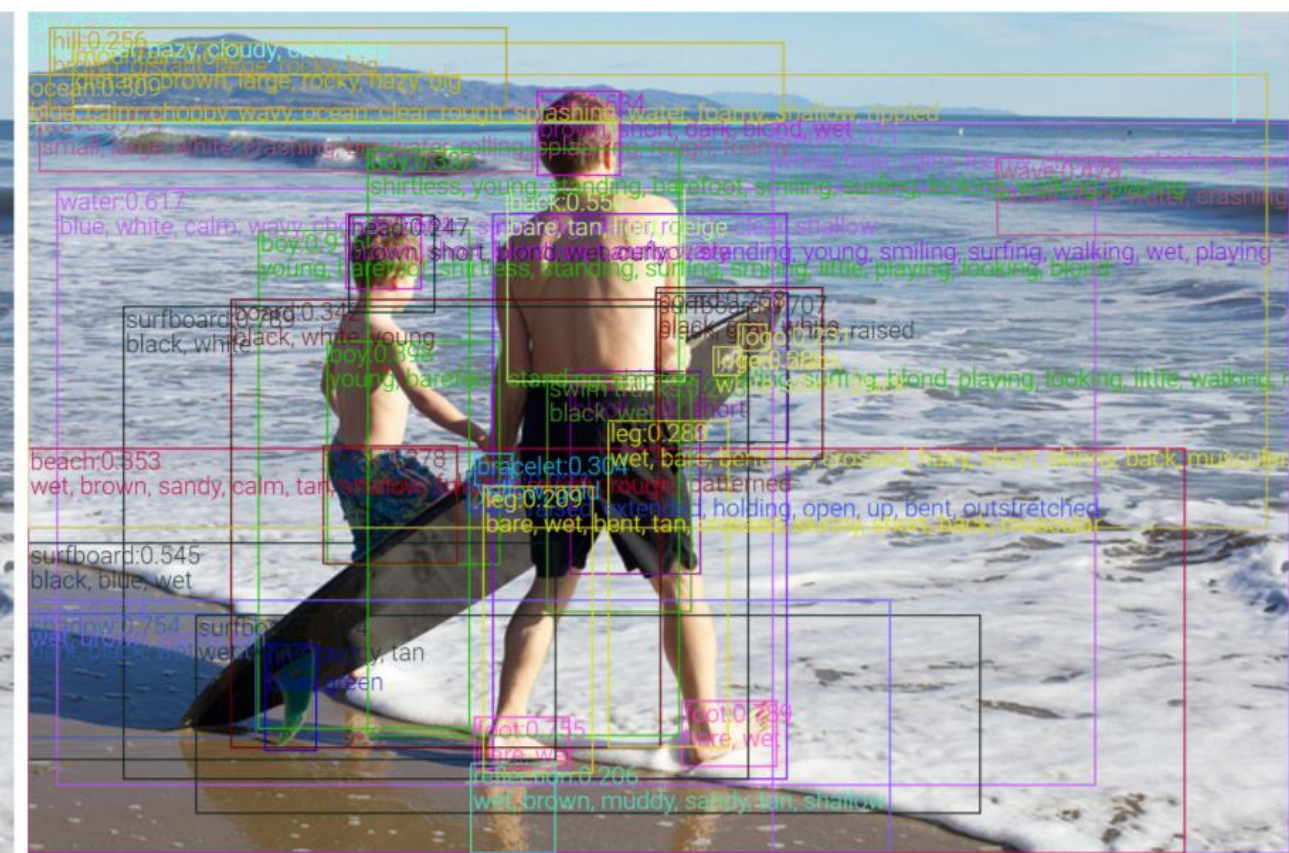
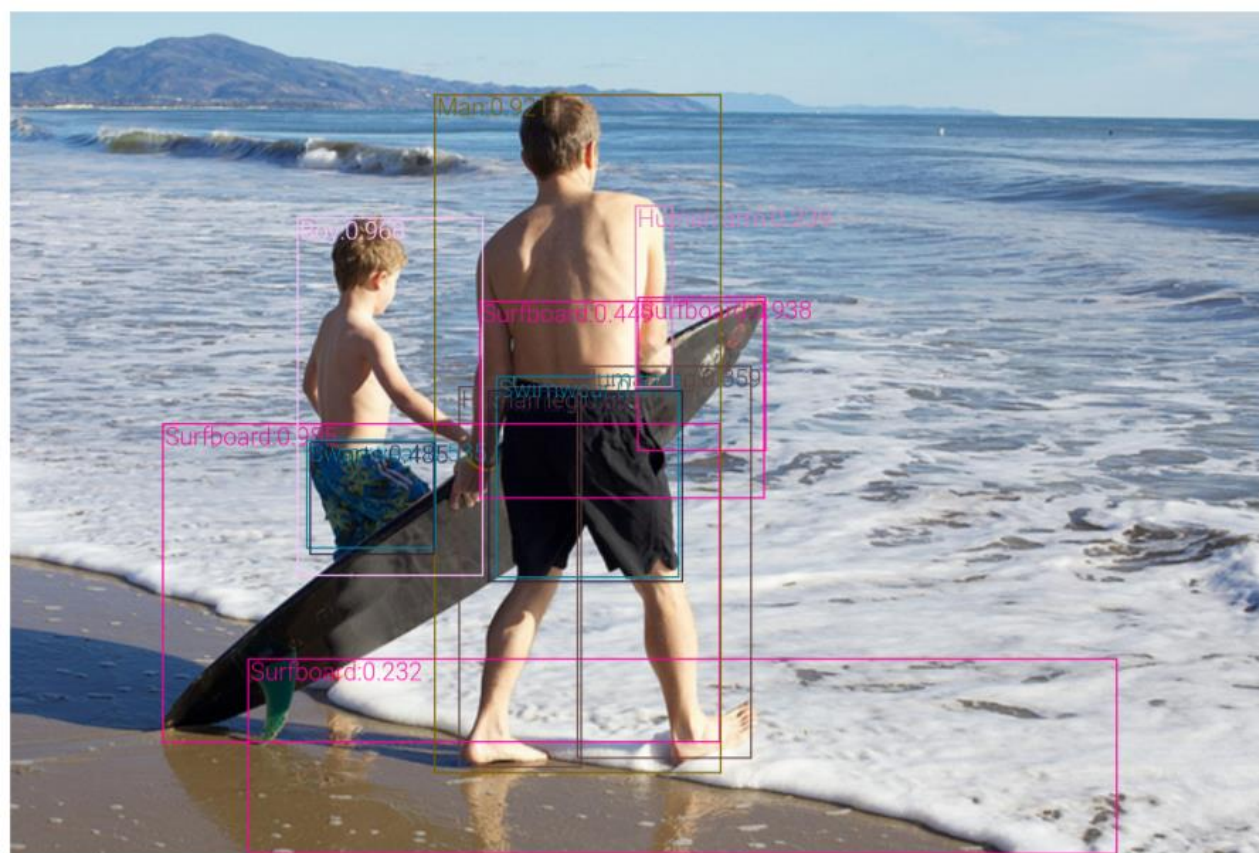
## Full pre-training objective

- $L_{\text{pre-training}} = L_M + L_C$

## Oscar+ novel 3-way contrastive loss – $L_{\text{CL3}}$

- Create two types of polluted training samples –  $(w', q, v)$  and  $(w, q', v)$
- Idea –  $(w', q, v)$  is polluted captions,  $(w, q', v)$  is polluted answers
- Goal of training – predict type of pollution in triplet (using the [CLS] token)





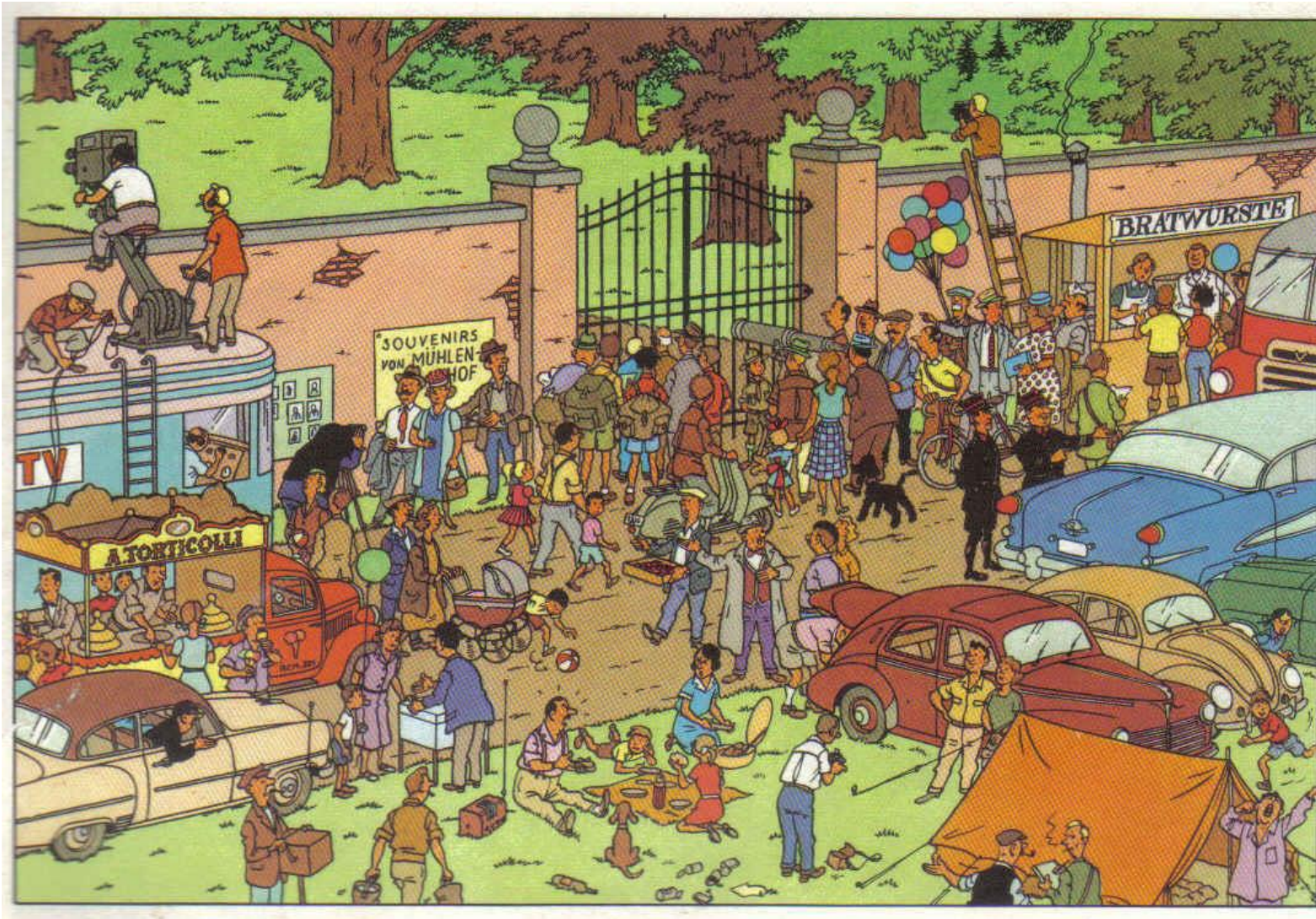
# VinVL<sup>1</sup>

- An improved OD model for VL
- Model based on ResNeXt-152 C4 architecture
- Trained on a much larger composite dataset – 1848 classes
- Datasets employed – COCO, OpenImages, Objects365, VisualGenome

<sup>1</sup> VinVL: Revisiting Visual Representations in Vision-Language Models, <https://arxiv.org/pdf/2101.00529.pdf>

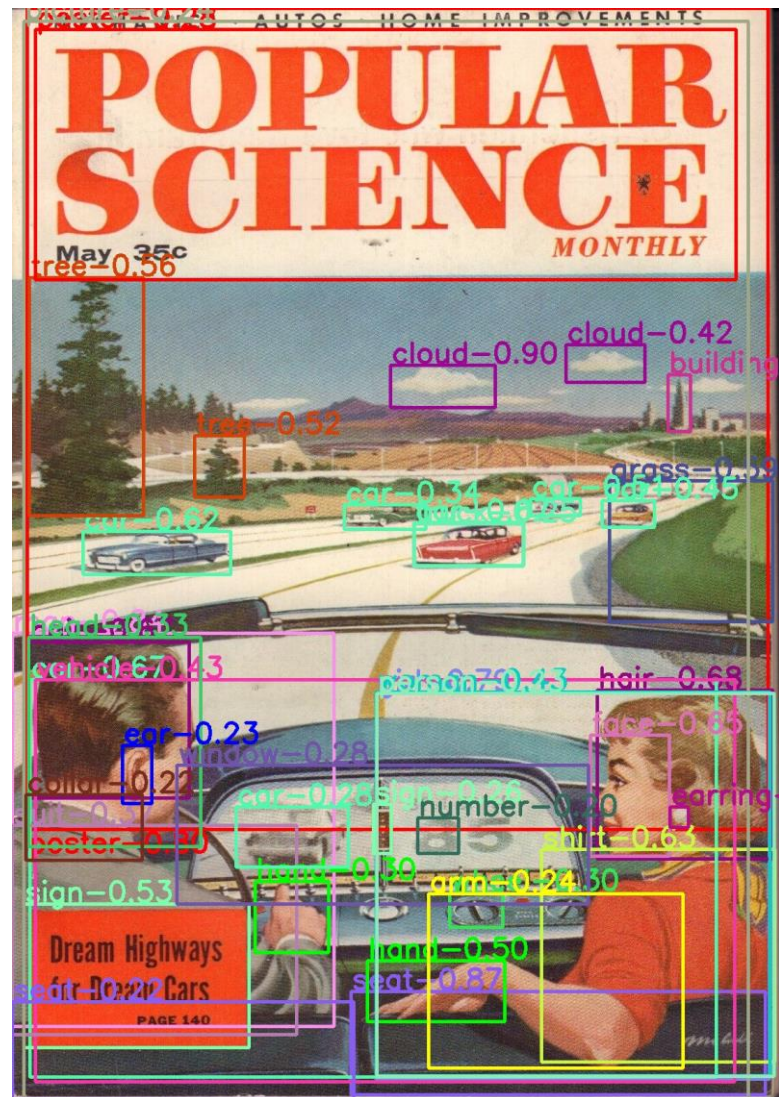
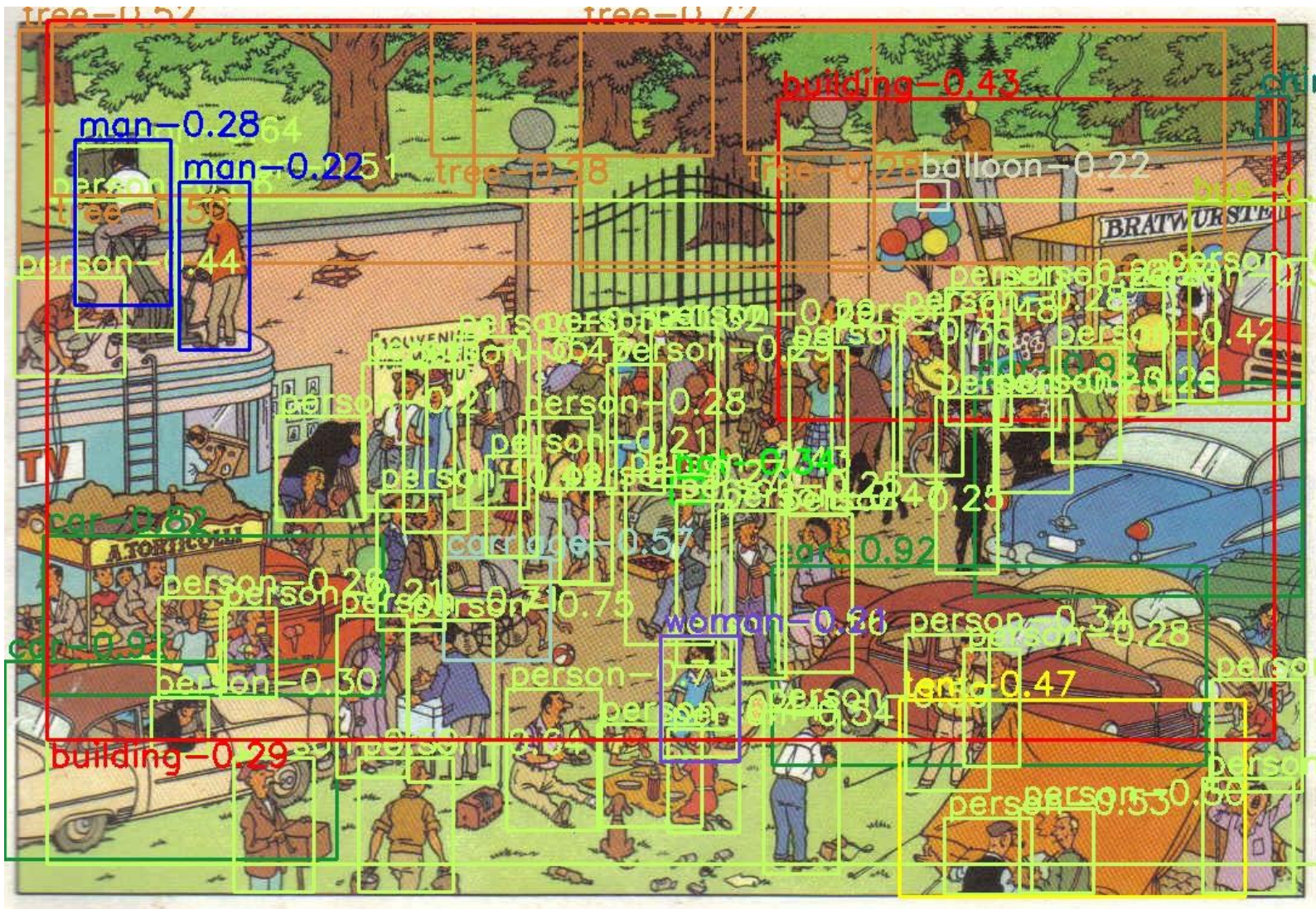


# Object Detection Examples





# Object Detection Examples





# Object Detection Examples

- Most ads have gun accessories – OD model not built to detect



screen light hole device button handle box

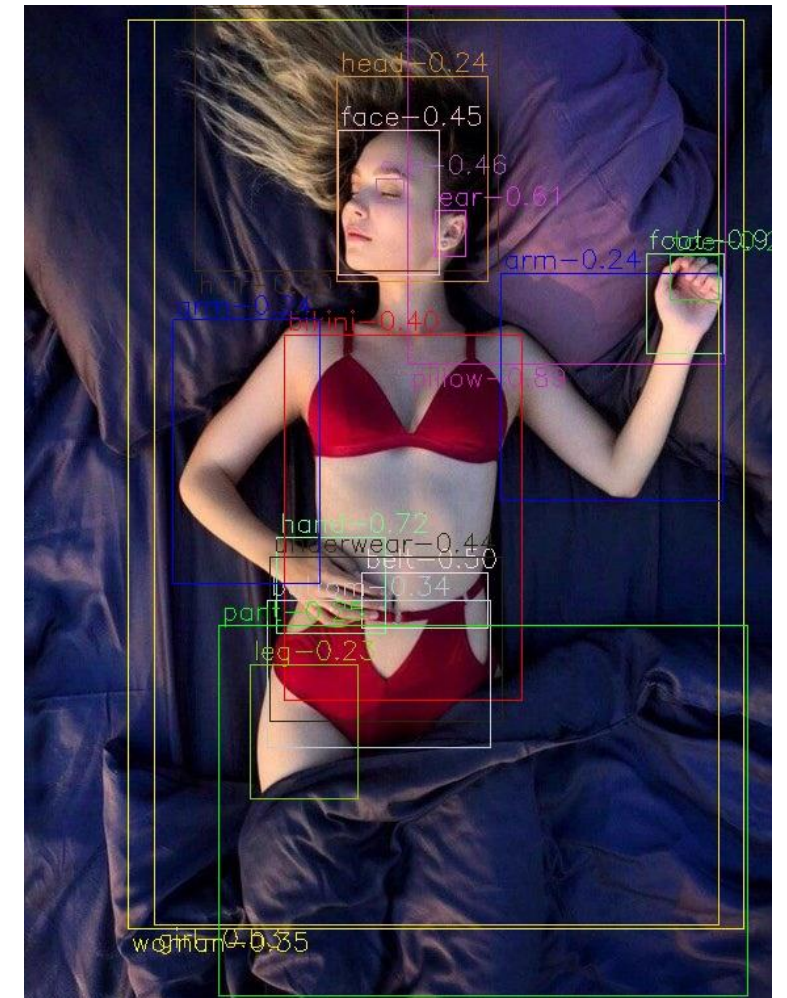


box cap pen word sticker word letter logo gun bracket cap table screw toy handle light





# Object Detection Examples



# PA Editorial Experiments

Model	OD tag len	LR (With Linear Schedule)	Adult PRAUC	Weapon PRAUC
VinVL base	0	2e-5	0.9088	0.9319
VinVL base	10	2e-5	0.9071	0.9335
VinVL base	20	2e-5	0.9078	0.9331
VinVL base	20	2.5e-5	0.9161	0.9310
VinVL large	20	2e-5	0.9110	0.9322

# PA Editorial Results

Model	Target : Overall Disallowed	
	Adult PRAUC	Weapon PRAUC
Unimodal NN baseline	0.9521	0.9380
VinVL base	0.9078	0.9331
Multimodal base	0.9518	0.9364
VinVL large	0.9110	0.9322
Multimodal large	0.9516	0.9365

- Weapon Class: obtaining similar results
- Adult Class: significant dip
- Adult is heavily dependent on Image Modality and there is scope of improvement there.
  - Objects tags are oversampled and not refined for Ads data.
  - Object Detector is not finetuned on Ads data.
- VinVL uses BERT initialization and our baselines are using CULR as the text model

# Future Work

Do Oscar+ pre-training on newer models (XLM-R, CULR, etc)

Explore adding an image token representing whole image (BiT penultimate layer)

Policy-based domain-specific object tags

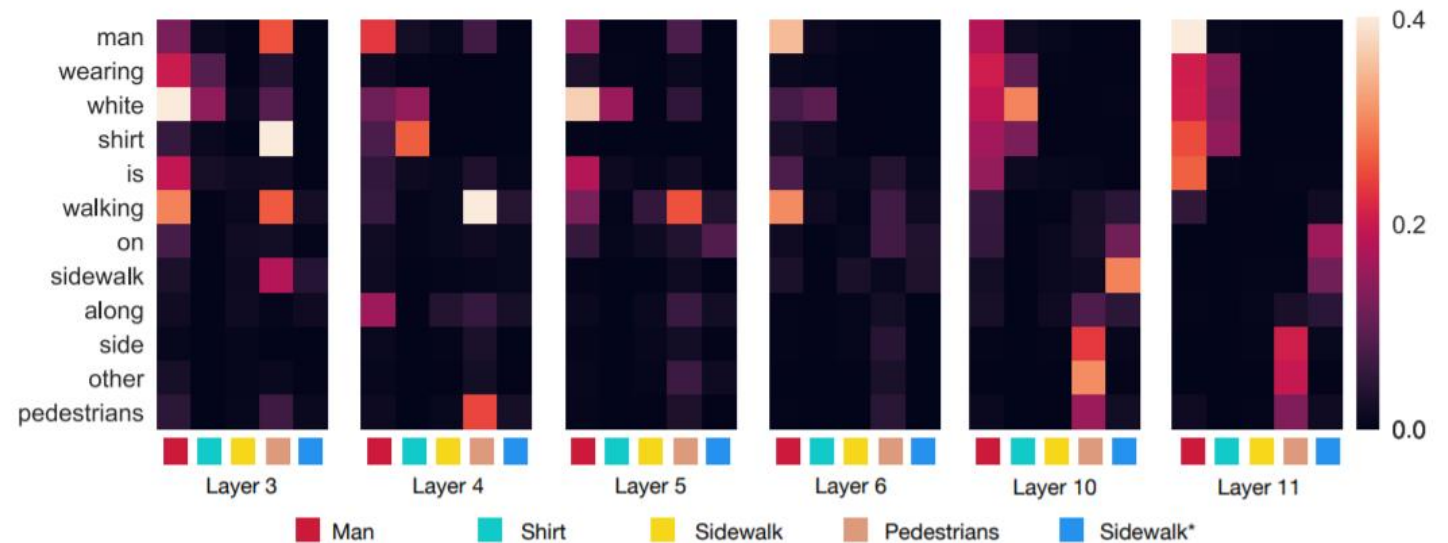
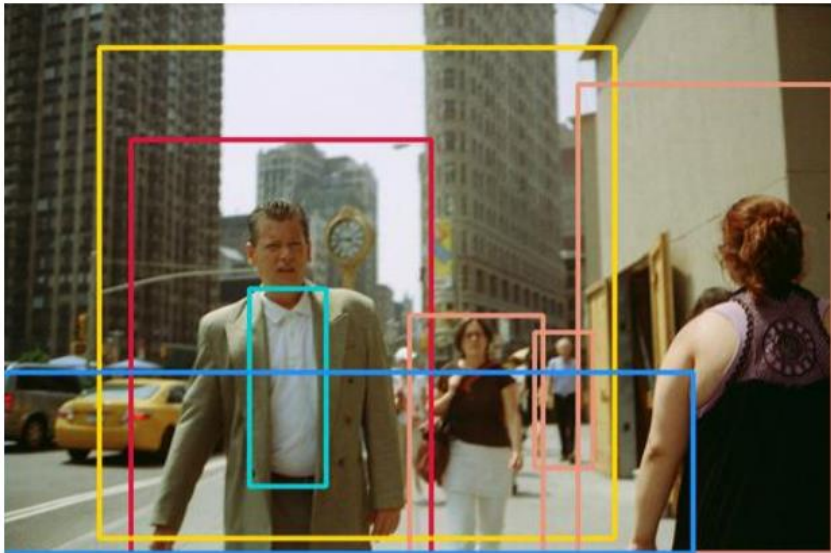
**Thanks**

**Open To Questions**

# Appendix

# Multimodal Learning

- Attention weights of some selected heads in VisualBERT
  - Implicit grounding of visual concepts in higher layers (eg "man wearing white shirt")
  - Refinement of understanding (eg "walking"-*Pedestrians*)
  - Correction of alignment (eg "shirt"-*Man*)



# Data Distribution

Overall AdDisallowed	Train	Val	GoldenTest
Compliance	950K	230K	59k
Adult	148K	27K	6.3k
Weapon	72K	16K	5k