# Probing Contextual Language Models for Common Ground with Visual Representations

Gabriel Ilharco

Rowan Zellers

Ali Farhadi

Hannaneh Hajishirzi

UNIVERSITY of WASHINGTON

UWNLP

NAACL 2021

1

How do **text representations** relate to the visual world?

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

2

# Motivation

How do **text representations** relate to the visual world?



a dog is sleeping
on the floor

# Motivation

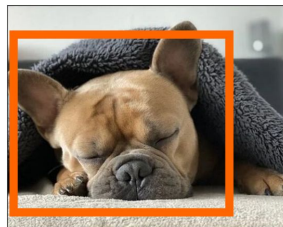How do **text representations** relate to the visual world?
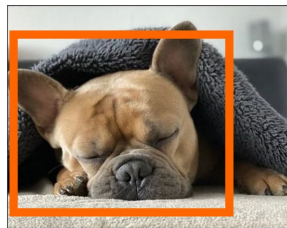
a **dog** is sleeping
on the floor

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

4

# Motivation

How do **text representations** relate to the visual world?



a **dog** is sleeping on the floor

We measure whether contextual **text representations** of concrete **objects** are effective in finding aligned image patches

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

5

# Motivation

**Context** in critical for this investigation



a **dog** is sleeping
on the floor

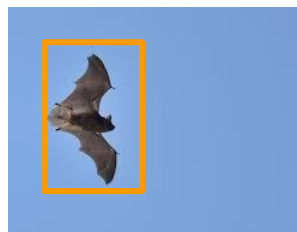Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

6

# Motivation

**Context** in critical for this investigation

A **bat** flying in the sky

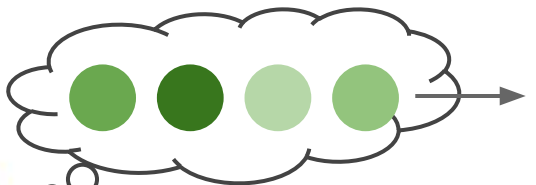Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021
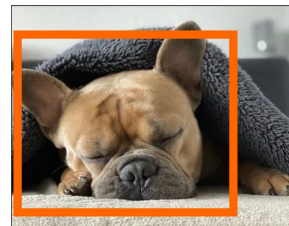
7

# Overview

Our method uses a lightweight **probe** that measures how **text** and **visual** representations are related

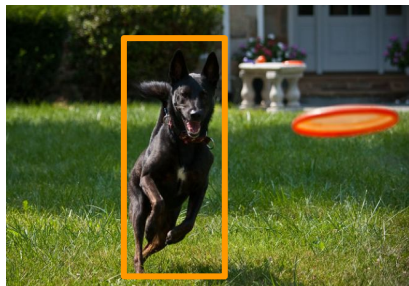a **dog** is sleeping on the floor

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021
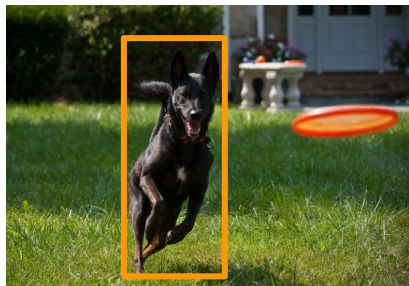
8

# Method - collecting representations

We find aligned representations of concrete **objects**

a **dog** is chasing
an orange frisbee



Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

9

# Method - collecting representations

We find aligned representations of concrete **objects**

a **dog** is chasing
an orange frisbee

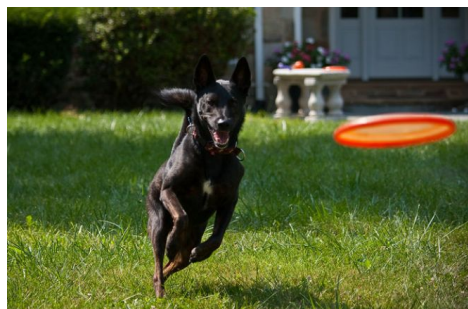Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021
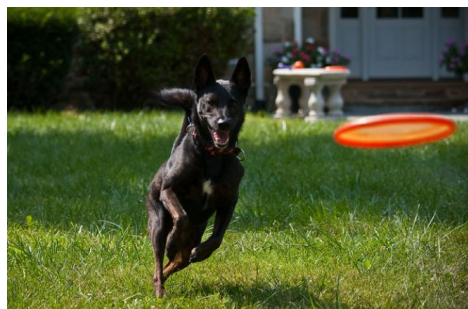
10

# Method - collecting representations

From image captioning datasets, we find aligned pairs of **instances**

using a trained object detector



a dog is chasing an
orange frisbee

# Method

From image captioning datasets, we find aligned pairs of **instances**
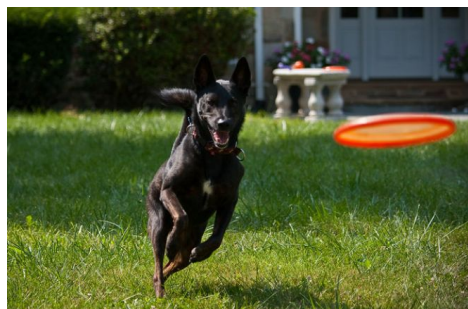
using a trained object detector



a dog is chasing an
orange frisbee



a **dog** is chasing an
orange frisbee

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021
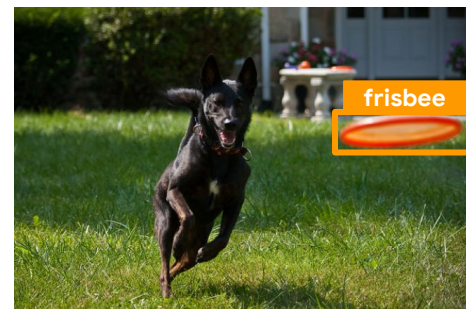
12

# Method

From image captioning datasets, we find aligned pairs of **instances**

using a trained object detector



a dog is chasing an orange frisbee
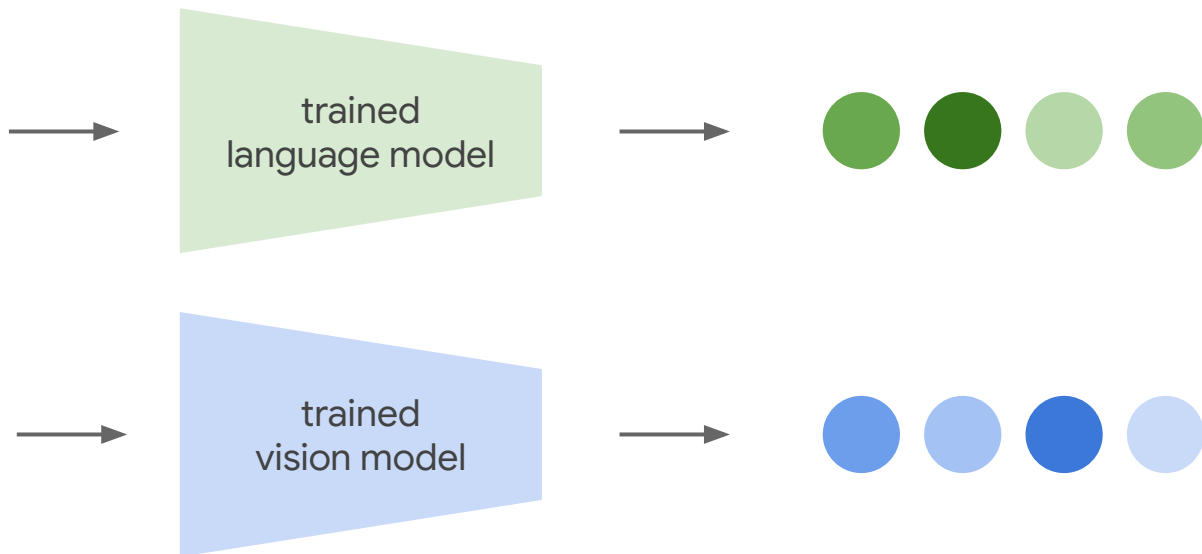
a **dog** is chasing an orange frisbee

a dog is chasing an orange **frisbee**

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

# Method - Collecting Data

**Text** and **visual** representations are extracted by trained models



a **dog** is chasing an orange frisbee → trained language model →

dog → trained vision model →

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

14

# Method - Collecting Data

**Text** and **visual** representations are extracted by trained models



a **dog** is chasing an orange frisbee

trained language model

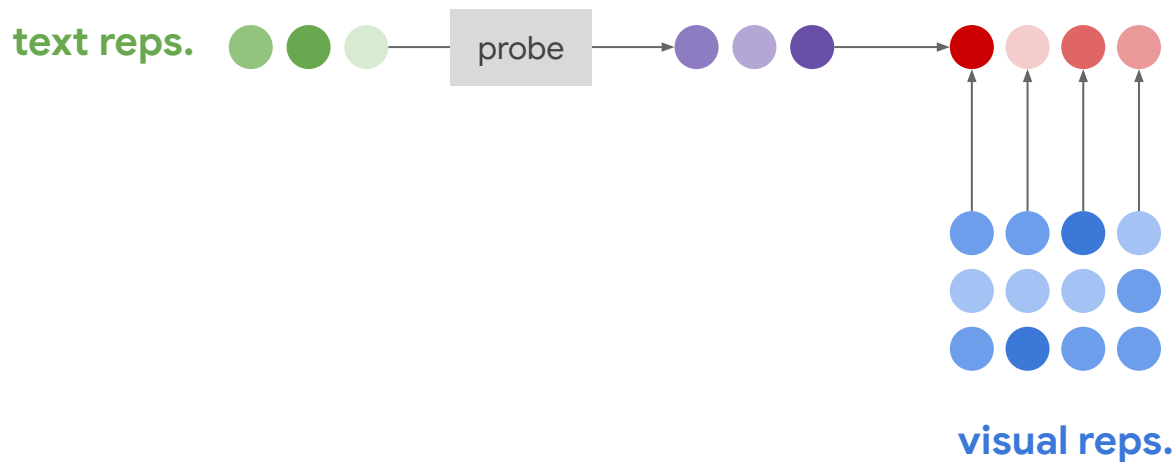trained vision model

# Method – Inspecting Text Representations

The **probe** maps **text representations** to the **visual domain.**

text reps. 〇〇〇 — probe → 〇〇〇

# Method - Inspecting Text Representations

The **probe** maps **text representations** to the **visual domain.**

We compute the **dot product** between **projected representations** and **visual representations**
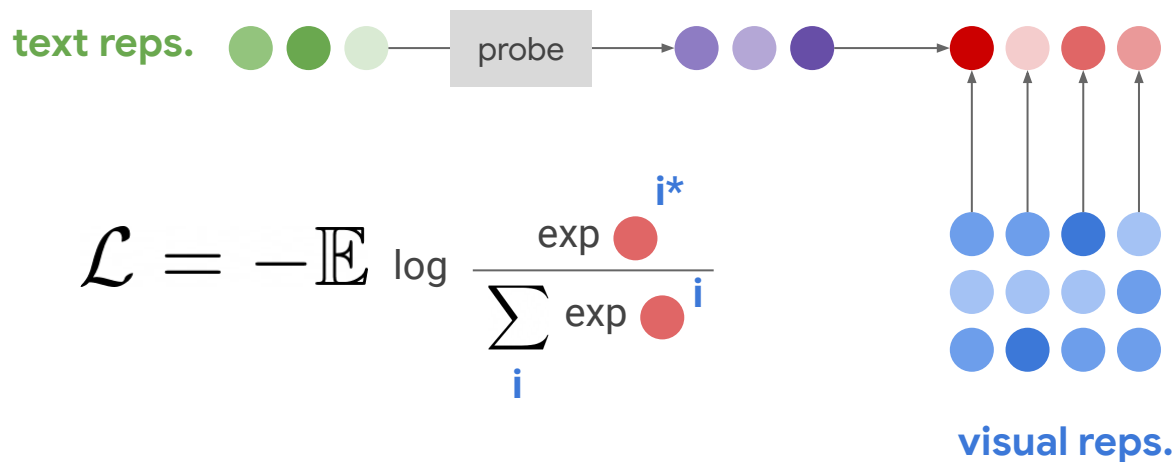


text reps.

probe

visual reps.

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

17

# Method - Inspecting Text Representations

The **probe** maps **text representations** to the **visual domain.**

We compute the **dot product** between **projected representations** and **visual representations**
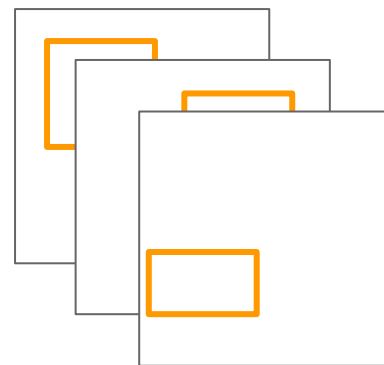
The **probe** is optimized via a **contrastive loss,** InfoNCE (Oord et al., 2018)



$$\mathcal{L} = -\mathbb{E} \log \frac{\exp \bullet^{i*}}{\sum_i \exp \bullet^{i}}$$

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

18

# Method - Evaluation

We then evaluate by retrieving image patches of **unseen object categories**

a man in the park is
flying a **kite**

probe →

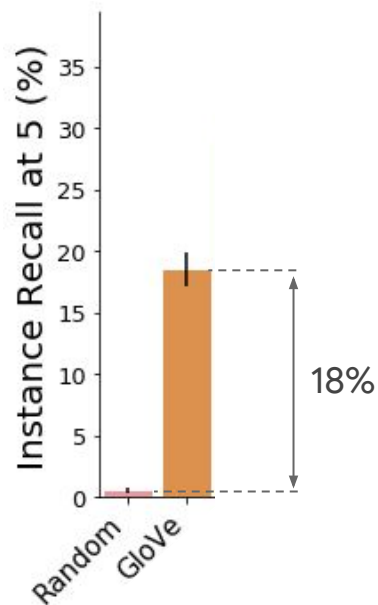top-K retrieved
image patches

# Method - Evaluation

We then evaluate by retrieving image patches of **unseen object categories**.

We report two metrics:

- **Category Recall at K**:

    - how often an image patch of the correct object category was in the top-K

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

20

# Method - Evaluation

We then evaluate by retrieving image patches of **unseen object categories**.

We report two metrics:

- **Category Recall at K**:

    - how often an image patch of the correct object category was in the top-K

- **Instance Recall at K**:

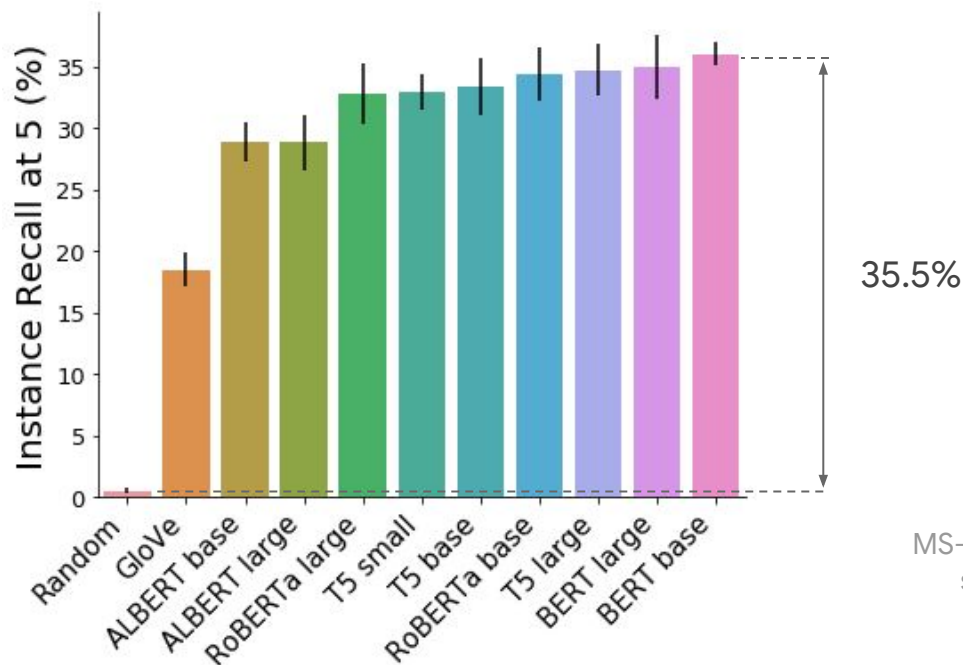    - how often the correct instance was in the top-K
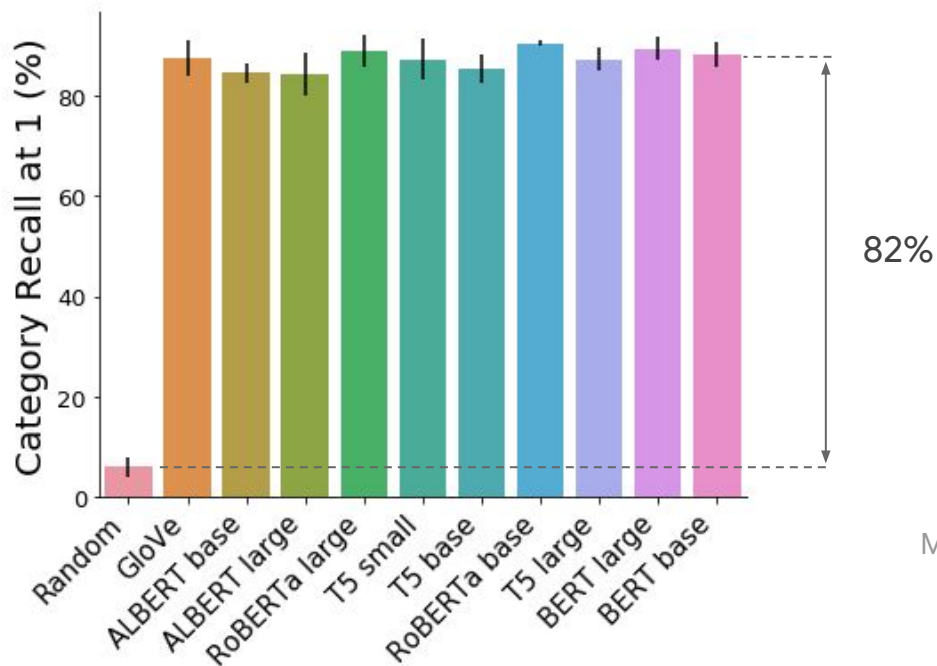
# Results

Language representations provide a strong signal for retrieval



18%

Instance retrieval results on MS-COCO using 1000 test samples spanning 200 unseen object categories
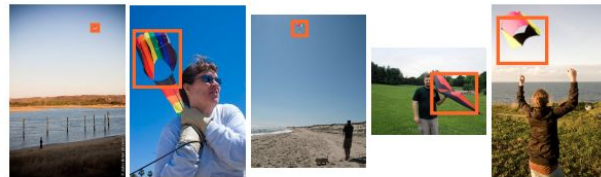
Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

22

# Results

Language representations provide a strong signal for retrieval



35.5%

Instance retrieval results on MS-COCO using 1000 test samples spanning 200 unseen object categories

# Results

Language representations provide a strong signal for retrieval



82%

Category retrieval results on MS-COCO using 1000 test samples spanning 200 unseen object categories

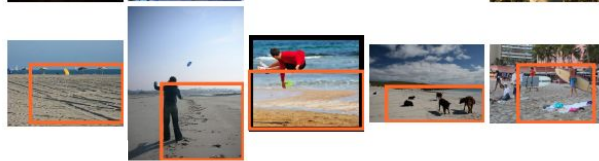Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021
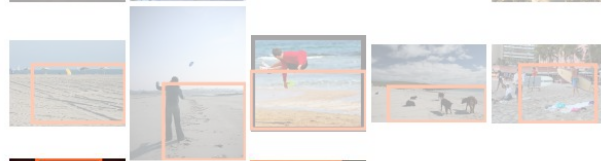
24

# Results - Qualitative Results

There is a man in the park flying a **kite**. →

A person flying a colorful kite on a **beach**. →

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

25

# Results - Qualitative Results



There is a man in the park flying a **kite**.

A person flying a colorful kite on a **beach**.

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021
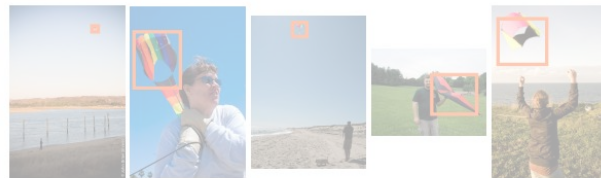
26

There is a man in the park flying a **kite**.

A person flying a colorful kite on a **beach**.

A **cat**.

A black **cat**.

A **cat** sleeping.

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

27

There is a man in the park flying a **kite**. →

A person flying a colorful kite on a **beach**. →

A **cat**. →

A black **cat**. →

A **cat** sleeping. →

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021
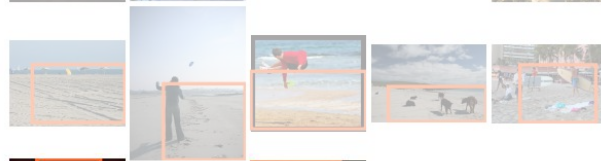
28

# Results - Qualitative Results

There is a man in the park flying a **kite**. ⟶

A person flying a colorful kite on a **beach**. ⟶
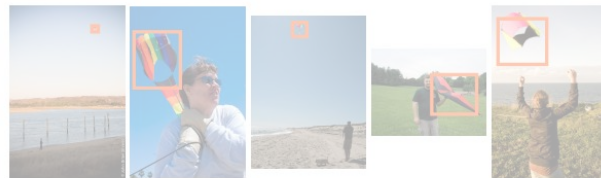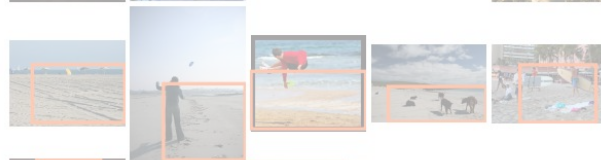
A **cat**. ⟶

**A black cat.** ⟶

A **cat** sleeping. ⟶

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

29

# Results - Qualitative Results

There is a man in the park flying a **kite**.  →

A person flying a colorful kite on a **beach**.  →
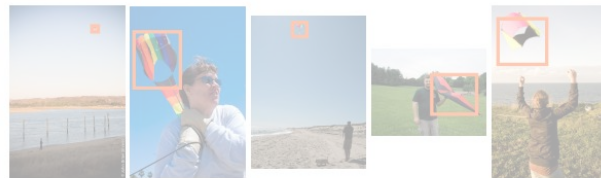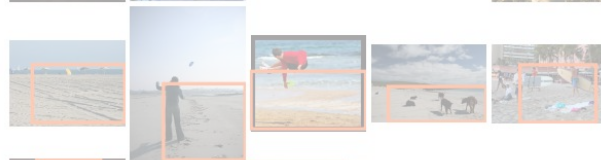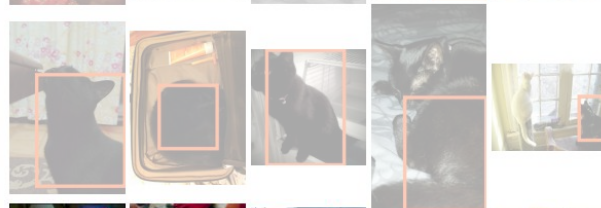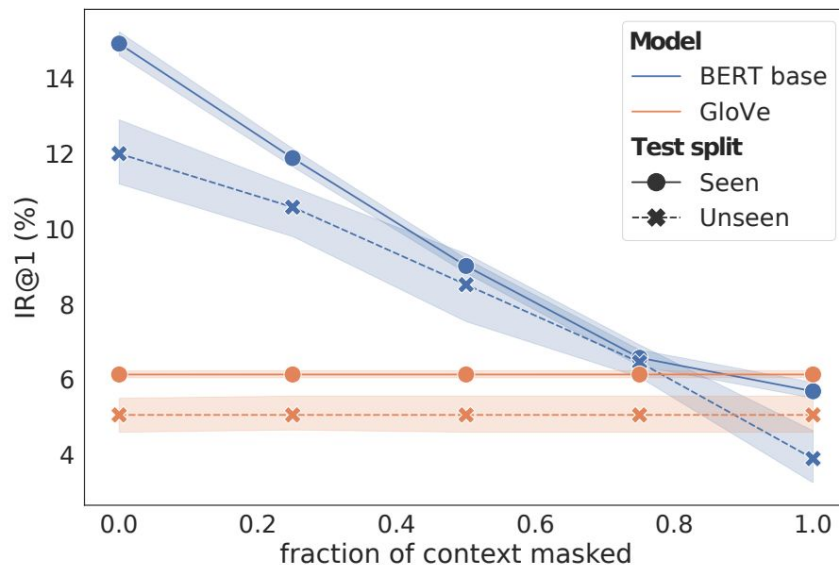
A **cat**.  →

A black **cat**.  →

A **cat** sleeping.  →

# Results - Influence of context

Performance of contextual models quickly degrades as
context tokens are progressively masked out

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

31

# Results - Influence of context

More descriptive sentences lead to better retrieval:
performance increases when objects are accompanied by at least one adjective



Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

32

# Results - Grounded Models

Grounded models slightly outperform text-only models



Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

33

All models substantially underperform humans

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

34

# Takeaways

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

35

# Takeaways

- We introduce a method for measuring similarities between text and visual representations

# Takeaways

- We introduce a method for measuring similarities between text and visual representations

- Contextual language representations are useful in finding aligned image patches

  - We explore how results are affected by variables such as context and explicit grounding

    during training

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021
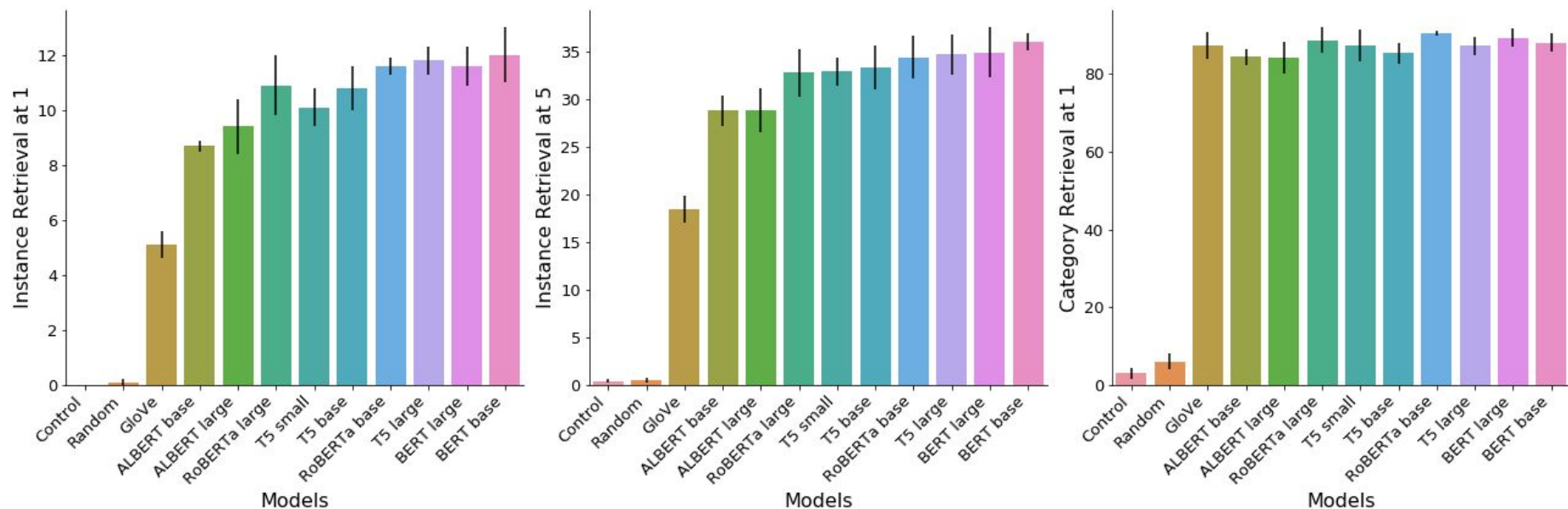
37

# Takeaways

- We introduce a method for measuring similarities between text and visual representations

- Contextual language representations are useful in finding aligned image patches

  - We explore how results are affected by variables such as context and explicit grounding

    during training

- All studied models significantly underperform humans, showing much room for future progress

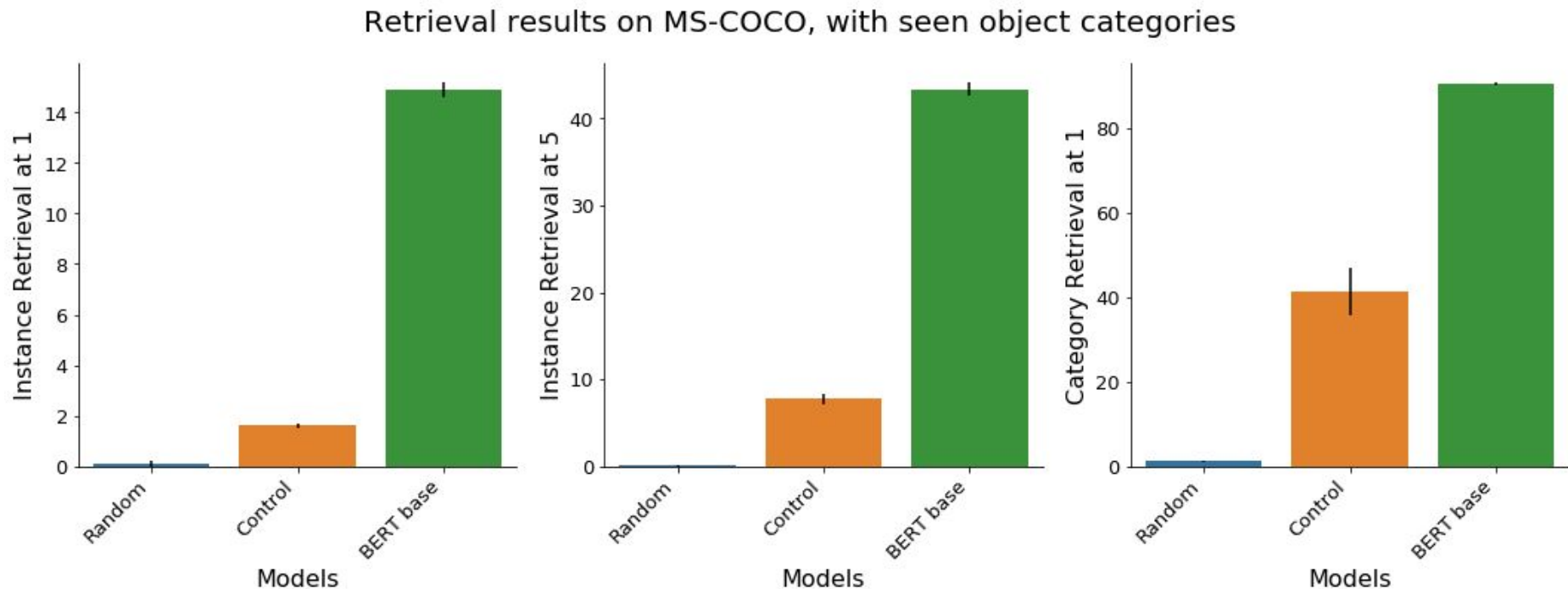Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

38

# Thank you!

# Results - Control



Retrieval results on MS-COCO, with unseen object categories

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

40

# Results - Seen object categories



Retrieval results on MS-COCO, with seen object categories

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

41

# Results - Loss ablations



Retrieval results on MS-COCO, with unseen object categories

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021
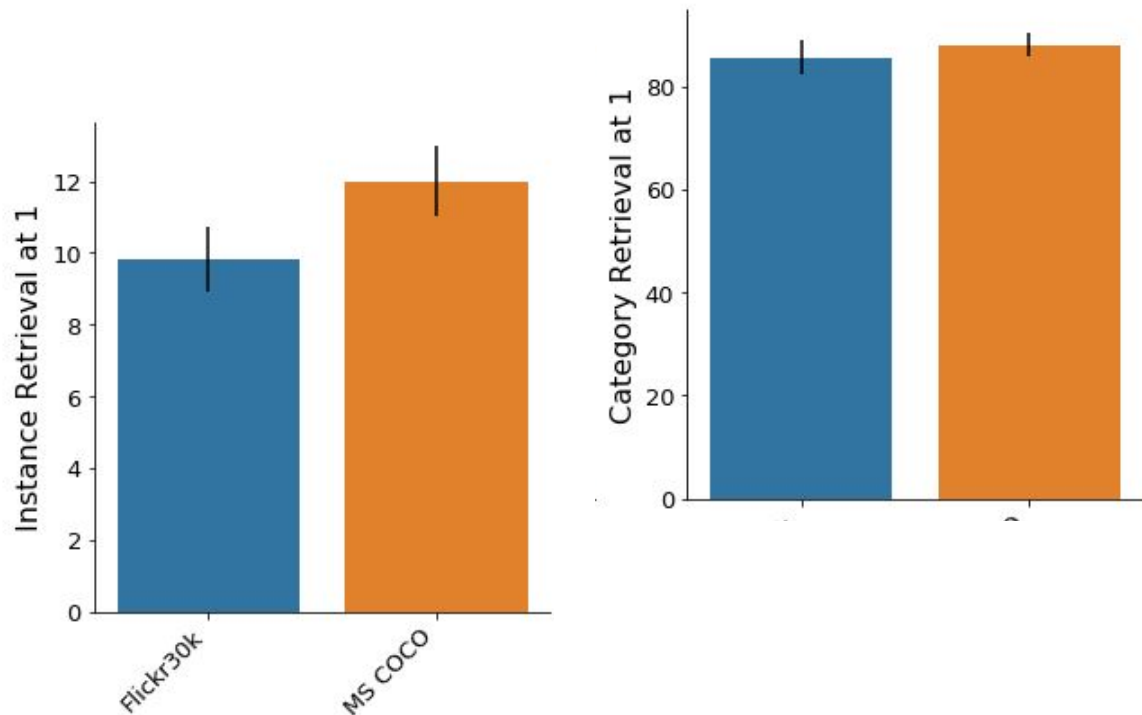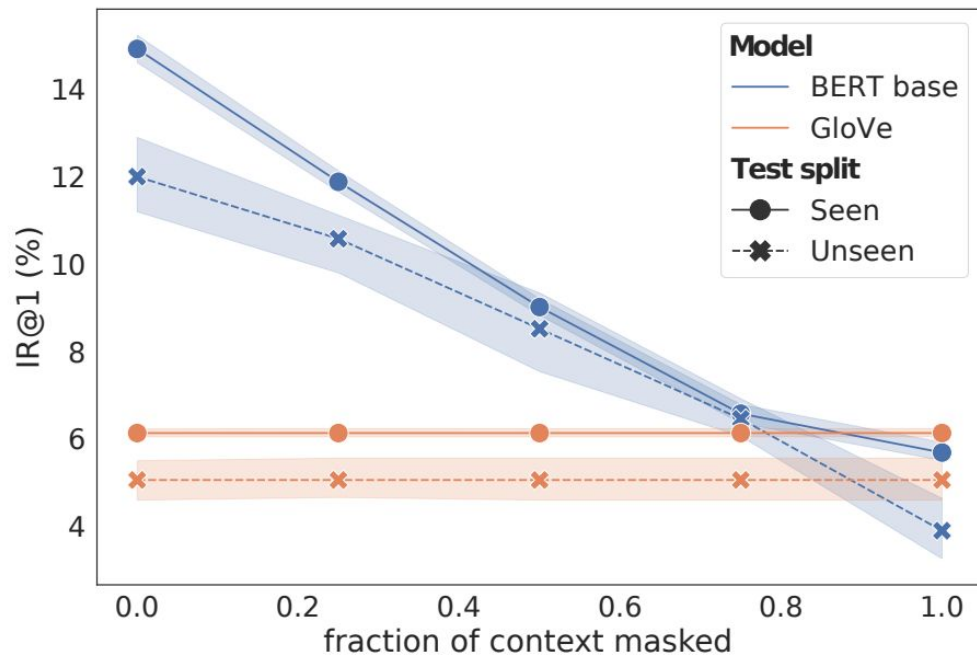
42

# Results - Data ablations



Retrieval results on multiple datasets, with unseen object categories

# Results - Data ablations

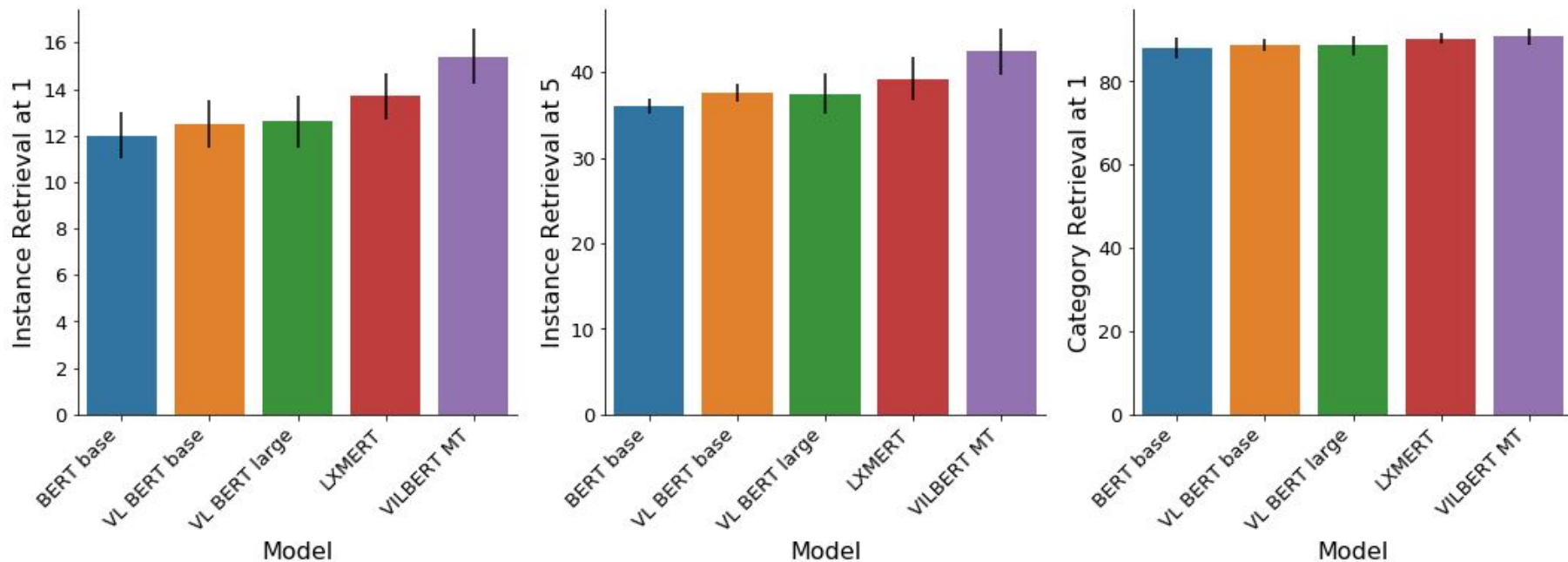# Results - Influence of context

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

45

# Results - Influence of context

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

46

# Results - Grounded Models



Retrieval results for grounded models, with unseen object categories

a **dog** is sleeping on the floor

Probing Contextual Language Models for Common Ground with Visual Representations. Ilharco et. al, 2021

48