# When Language Meets Vision:
# A Multimodal Perspective on the NLP World

Sapienza NLP Group
Bianca Scarlini
scarlini@di.uniroma1.it
*Reading Group @ Sapienza NLP*

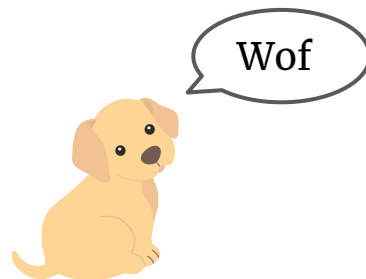# How do humans learn?

A dog is sitting on a couch with its toy.

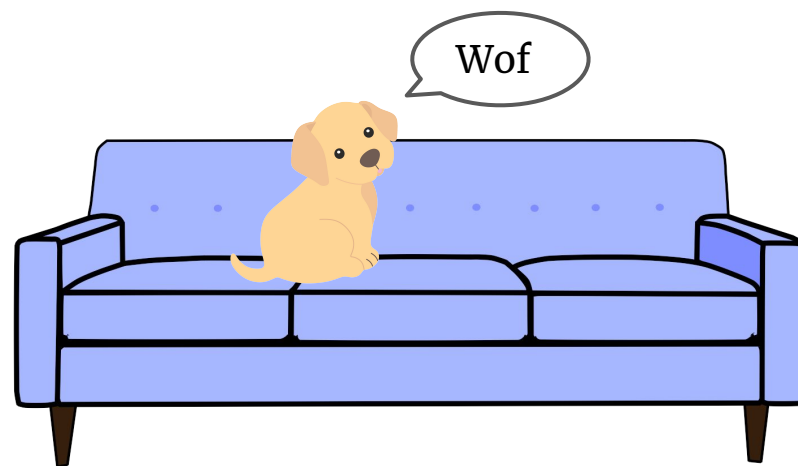Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# How do humans learn?

A **dog** is sitting on a couch with its toy.

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# How do humans learn?

A dog is sitting on a **couch** with its toy.

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# How do humans learn?

A dog is sitting on a couch with its **toy**.

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# How do text models learn?

Cat  Ball  Dog  Baby

↑

Txt

↑

A __ is sitting on a couch with its toy.

– Learn syntactic relations  *Clark et al. 2019*

✖ Not grounded in the real world  *Bender and Koller 2020*

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# How do visual models learn?

A dog is sitting on a couch with its toy.
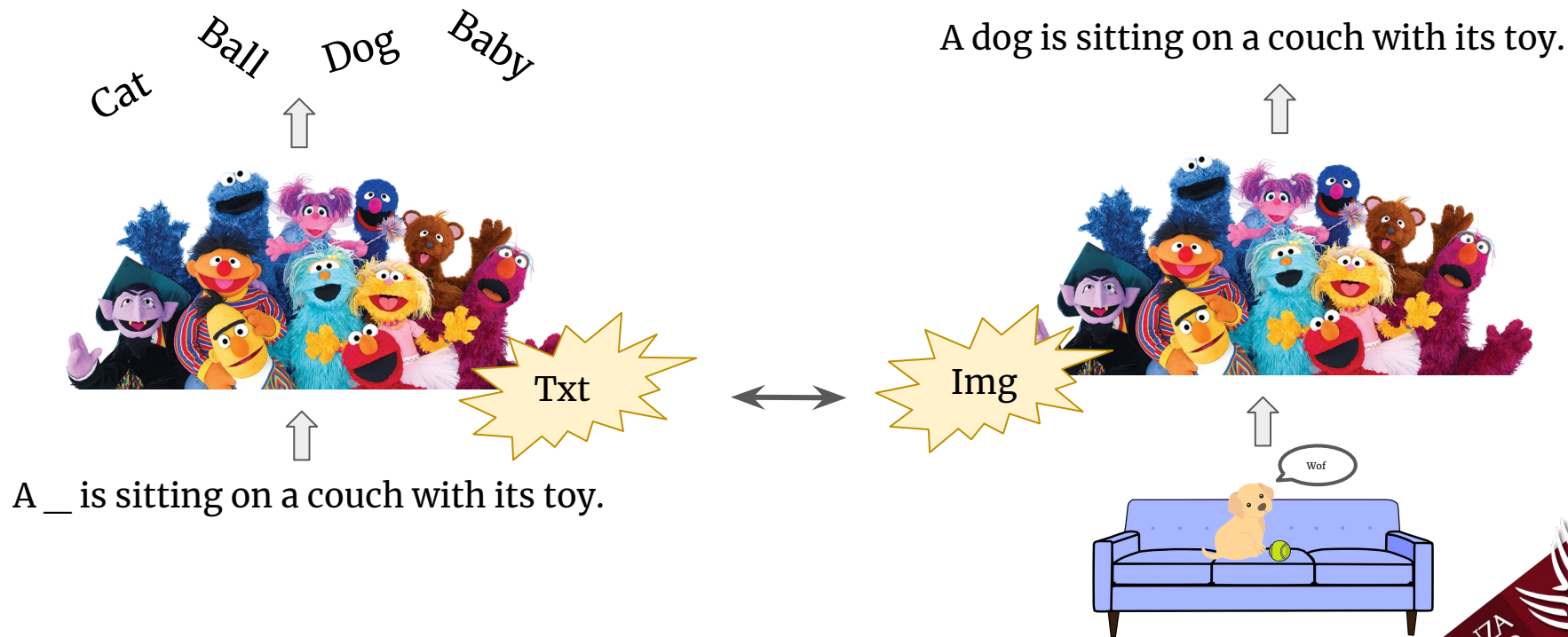
- Learn relations within objects in an image
  *Cadene et al. 2019*

- Need detailed semantics of the image for visual understanding
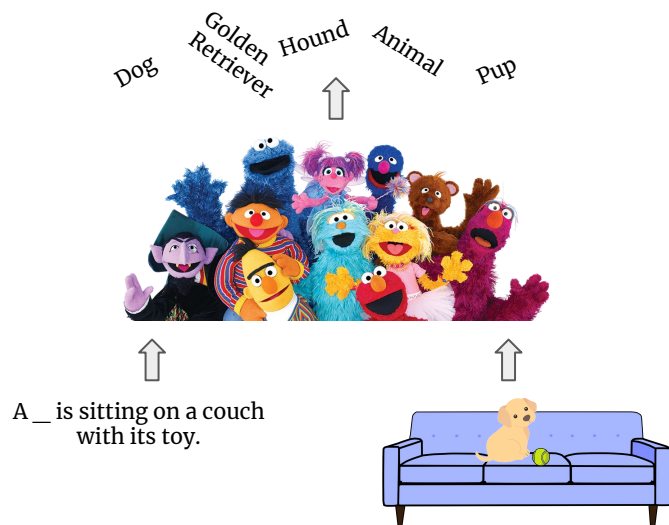  *Johnson et al. 2015*

Img

Wof

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# How do models learn?

Cat  Ball  Dog  Baby

A dog is sitting on a couch with its toy.

Txt ⟷ Img

A __ is sitting on a couch with its toy.

Wof

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

SAPIENZA NLP

# How do models learn?

Dog   Golden Retriever   Hound   Animal   Pup

Txt+Img

A __ is sitting on a couch with its toy.

Wof

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Cross-modal architectures

Single-stream architecture

Dog    Golden Retriever    Hound    Animal    Pup

A __ is sitting on a couch with its toy.

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Cross-modal architectures

## Single-stream architecture

Dog  Golden Retriever  Hound  Animal  Pup

A __ is sitting on a couch with its toy.

## Two-stream architecture

Dog  Golden Retriever  Hound  Animal  Pup

A __ is sitting on a couch with its toy.

Bianca Scarlini - When Language Meets Vision - Reading Group @ Sapienza NLP

# Cross-modal architectures

## Single-stream architecture

VL-BERT
*Su et al. 2020*

UNITER
*Chen et al. 2020*

OSCAR
*Li et al. 2020*

VisualBERT
*Li et al. 2019b*

Unicoder-VL
*Li et al. 2019a*

Golden Retriever
Dog
Pup

A __ is sitting on a couch with its toy.

## Two-stream architecture

ViLBERT
*Lu et al. 2019*

LxMERT
*Tan and Bansal 2019*

Golden Retriever
Hound
Animal
Pup

A __ is sitting on a couch
with its toy.

SAPIENZA NLP

# Cross-modal architectures

**Single-stream architecture**

Dog · Golden Retriever · Hound · Animal · Pup

VL-BERT
*Su et al. 2020*

OSCAR
*Li et al. 2020*

A __ is sitting on a couch with its toy.

**Two-stream architecture**

Dog · Golden Retriever · Hound · Animal · Pup

LxMERT
*Tan and Bansal 2019*

A __ is sitting on a couch with its toy.

Bianca Scarlini - When Language Meets Vision - Reading Group @ Sapienza NLP

SAPIENZA NLP

# VL-BERT *[Su et al. 2020]*

Visual and linguistic contents interact freely

Pretrain on visual-linguistic and text-only data

*A ___ is sitting on a couch with its toy.*

Add new visual features to BERT input embeddings

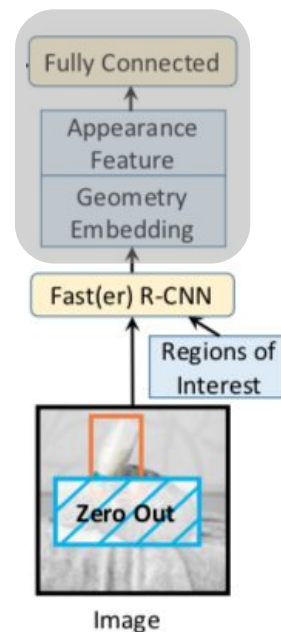Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# VL–BERT *[Su et al., 2020]*



Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# VL-BERT *[Su et al., 2020]*

Fast R-CNN *Girshick, 2015*

⟹ Object detection model

**DOG**

**BALL**

**COUCH**

**Region of Interests (RoI)**

Fully Connected

Appearance Feature

Geometry Embedding

Fast(er) R-CNN

Regions of Interest

Zero Out

Image

# VL-BERT *[Su et al., 2020]*

Fast R-CNN *Girshick, 2015*

⟹ Object detection model

**DOG**

**BALL**

**COUCH**

**Appearance Feature**
Feature vector prior to the output layer of RoI



Fully Connected

Appearance Feature

Geometry Embedding

Fast(er) R-CNN

Regions of Interest

Zero Out

Image

# VL-BERT *[Su et al., 2020]*

Fast R-CNN  *Girshick, 2015*

⟹ Object detection model



**DOG**

**BALL**

**COUCH**

**Geometry Embedding**
Sine and cosine functions in different wavelengths applied to normalized coordinates of RoI

# VL-BERT *[Su et al. 2020]*

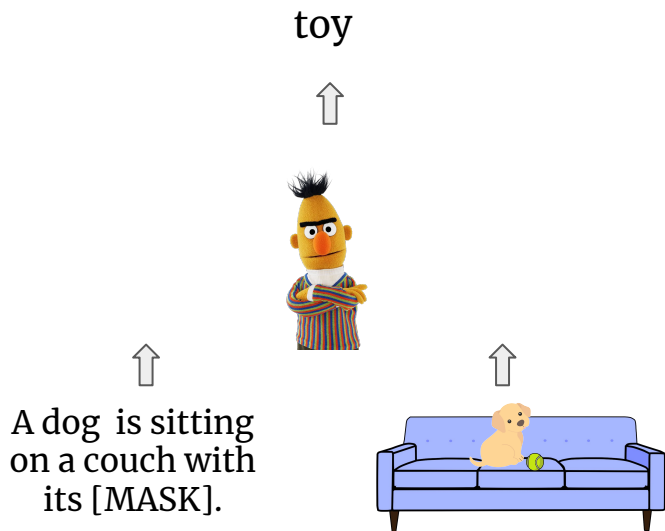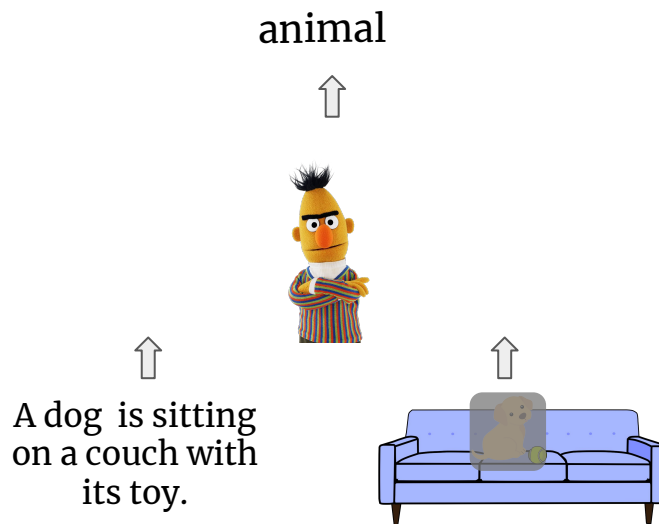# VL-BERT *[Su et al. 2020]*

# VL–BERT *[Su et al. 2020]*



**Whole Image**　　**One per RoI**

# VL-BERT *[Su et al. 2020]*

**Masked Language Modeling with Visual Clues**

**Masked RoI Classification with Linguistic Clues**

toy

animal

A dog  is sitting on a couch with its [MASK].

A dog  is sitting on a couch with its toy.

# VL-BERT *[Su et al. 2020]*

**Masked Language Modeling**

toy

⇧

⇧ ⇧

A dog is sitting
on a couch with
its [MASK].

# OSCAR *[Li et al., 2020]*

Single-stream architecture

Uses object tags in an image as anchor points

Ease the learning of image-text alignment

A dog is sitting on a couch with its toy.

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP
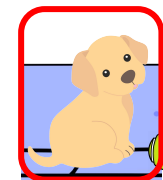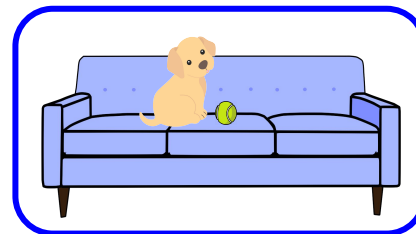
# OSCAR *[Li et al., 2020]*

A **dog** is sitting on a **couch**.

A dog is sitting on a couch.

**dog**
**couch**

# OSCAR *[Li et al., 2020]*

**Faster R-CNN,** *Ren et al., 2015*
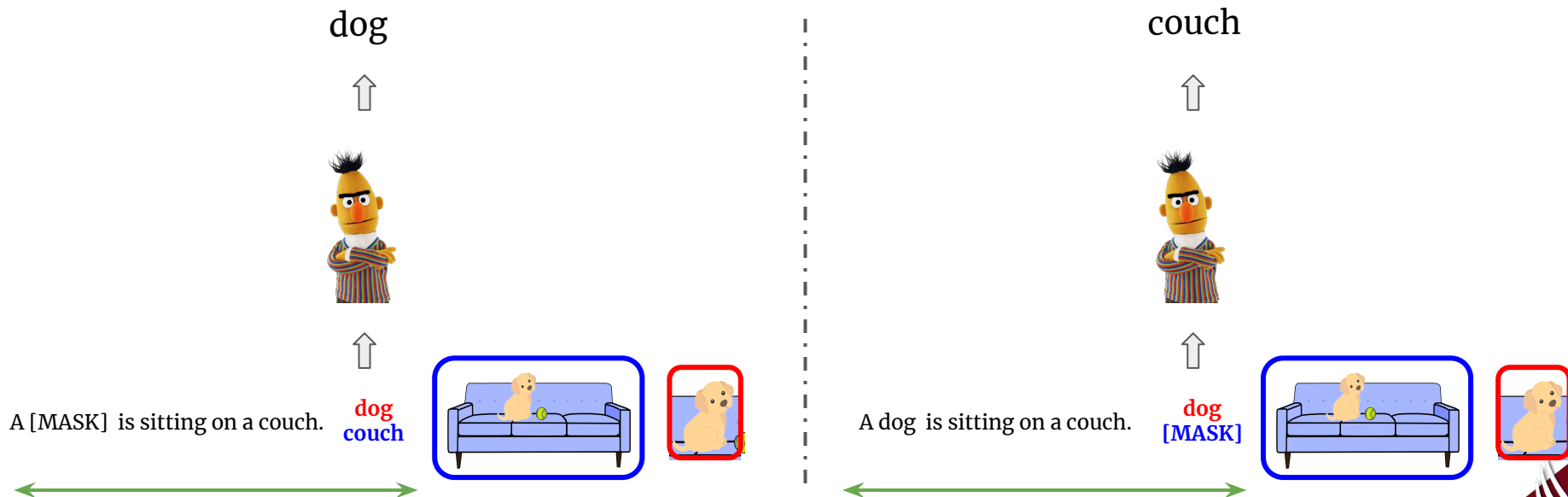


Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# OSCAR *[Li et al., 2020]*

**Masked Token Loss**

dog

couch

A [MASK] is sitting on a couch.

**dog
couch**

A dog is sitting on a couch.

**dog
[MASK]**

# OSCAR *[Li et al., 2020]*

**Contrastive loss**

yes

no

A dog is sitting on a couch.
**dog**
**couch**

A dog is sitting on a couch.
**bird**
**balloon**

# LxMERT *[Tan and Bansal, 2020]*

Double-stream architecture

Builds both intra-modality and cross-modality relations

Five diverse pre-training tasks

A [MASK] is sitting on a
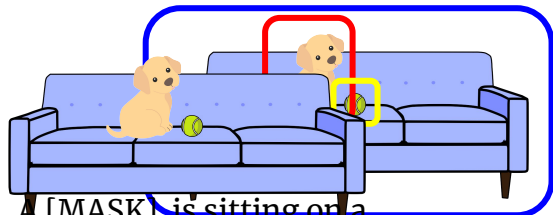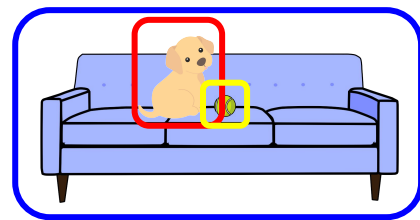couch with its toy.

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# LxMERT *[Tan and Bansal, 2020]*



Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# LxMERT *[Tan and Bansal, 2020]*



**Image Encoder**

**Language Encoder**

**Cross-modality Encoder**

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# LxMERT *[Tan and Bansal, 2020]*



Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# LxMERT *[Tan and Bansal, 2020]*



Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# LxMERT *[Tan and Bansal, 2020]*

**Cross-Modality Output [CLS]**



Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# LxMERT *[Tan and Bansal, 2020]*

## Masked Language Modeling with Visual Clues



couch

A dog  is sitting
on a [MASK].

# LxMERT *[Tan and Bansal, 2020]*

**RoI Feature Regression**

**Detected-Label Classification**

dog

A dog is sitting
on a couch.

A dog is sitting
on a couch.

# LxMERT *[Tan and Bansal, 2020]*

## Cross-Modality Matching



yes

A dog is sitting on a couch.

no

The sky is blue.

# LxMERT *[Tan and Bansal, 2020]*

## Image Question Answering

dog

⇑

⇗           ⇖

⇑                              ⇑

Who is sitting on
the couch?

# Experimental Setup

Image Text Retrieval

Visual Question Answering

Visual Commonsense Reasoning

Grounding Referring Expression

Image Captioning

Natural Language Visual Reasoning for Real

... and many more

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Visual Question Answering

dog

Who is sitting
on the couch?

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Visual Question Answering



VQA 2.0

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Natural Language Visual Reasoning for Real

yes

NLVR²
*Suhr et al., 2019*

The first image contains twice the number of dogs as the second image

# Natural Language Visual Reasoning for Real

### NLVR2



Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Visual Commonsense Reasoning

VCR
*Zellers et al., 2019*

**Image**



**Questions**

Why is [person4 🖼] pointing at [person1 🖼]?

**Answers**

a) He is telling [person3 🖼] that [person1 🖼] ordered the pancakes.
b) He just told a joke.
c) He is feeling accusatory towards [person1 🖼].
d) He is giving [person1 🖼] directions.

*I chose a) because...*

**Rationales**

a) [person1 🖼] has the pancakes in front of him.
b) [person4 🖼] is taking everyone's order and asked for clarification.
c) [person3 🖼] is looking at the pancakes and both she and [person2 🖼] are smiling slightly.
d) [person3 🖼] is delivering food to the table, and she might not know whose order is whose.

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Visual Commonsense Reasoning

Given the image and the question, return the correct answer
Q → A



Why is [person4 🖼️] pointing at [person1 🖼️]?

a) He is telling [person3 🖼️] that [person1 🖼️] ordered the pancakes.
b) He just told a joke.
c) He is feeling accusatory towards [person1 🖼️].
d) He is giving [person1 🖼️] directions.
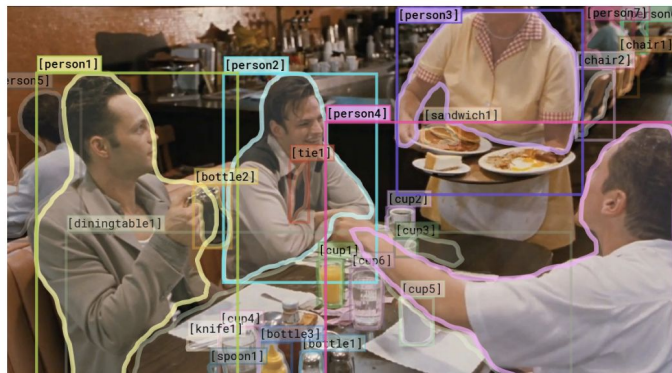
*I chose a) because...*

a) [person1 🖼️] has the pancakes in front of him.
b) [person4 🖼️] is taking everyone's order and asked for clarification.
c) [person3 🖼️] is looking at the pancakes and both she and [person2 🖼️] are smiling slightly.
d) [person3 🖼️] is delivering food to the table, and she might not know whose order is whose.

# Visual Commonsense Reasoning

Given the image, the question and the answer return the correct rationale
**QA → R**



Why is [person4 🧑] pointing at [person1 👤]?

a) He is telling [person3 👤] that [person1 👤] ordered the pancakes.
b) He just told a joke.
c) He is feeling accusatory towards [person1 👤].
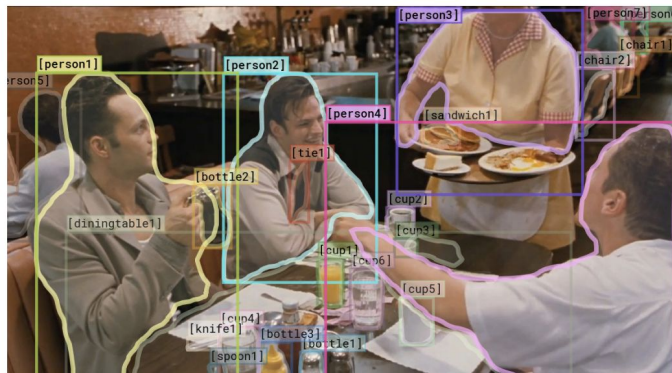d) He is giving [person1 👤] directions.

*I chose a) because...*

a) [person1 👤] has the pancakes in front of him.
b) [person4 👤] is taking everyone's order and asked for clarification.
c) [person3 👤] is looking at the pancakes and both she and [person2 👤] are smiling slightly.
d) [person3 👤] is delivering food to the table, and she might not know whose order is whose.

# Visual Commonsense Reasoning

Given the image and the question return the correct answer and rationale
**Q → AR**



Why is [person4 🧑] pointing at [person1 🧑]?

a) He is telling [person3 🧑] that [person1 🧑] ordered the pancakes.
b) He just told a joke.
c) He is feeling accusatory towards [person1 🧑].
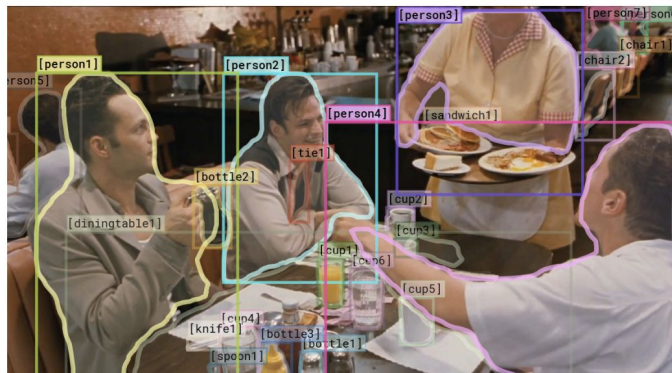d) He is giving [person1 🧑] directions.

*I chose a) because…*

a) [person1 🧑] has the pancakes in front of him.
b) [person4 🧑] is taking everyone's order and asked for clarification.
c) [person3 🧑] is looking at the pancakes and both she and [person2 🧑] are smiling slightly.
d) [person3 🧑] is delivering food to the table, and she might not know whose order is whose.

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Visual Commonsense Reasoning



VCR

■ Task-specific Model    ■ VL-BERT

Q -> A: Task-specific Model 65.1, VL-BERT 75.8
QA -> R: Task-specific Model 67.3, VL-BERT 78.4
Q -> AR: Task-specific Model 44, VL-BERT 59.7

Accuracy (y-axis), Subtasks (x-axis)

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Conclusions

Vision-and-language models gained much interest in the last couple of years

Straightforward techniques to incorporate visual features in contextualized language models

Vision-and-language models raised the bar for the state-of-the-art in many vision-and language tasks

Many directions that are still worth to be explored!

Bianca Scarlini – When Language Meets Vision – Reading Group @ Sapienza NLP

# Thanks for your attention!
*Any questions? Feel free to ask*