# POLITECNICO
## MILANO 1863

# Concept-Based Explanations for Image Classifiers Using Textual Prompts

Daniele Di Santi

**Academic year**: 2024/2025
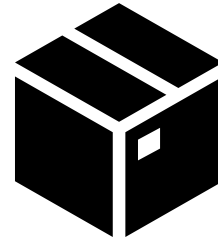
**Advisor**: Prof. Marco Brambilla

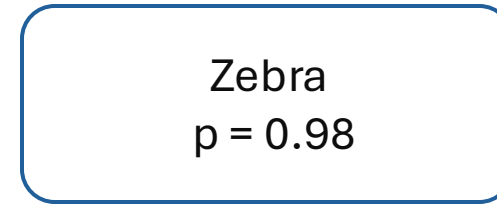**Co-advisors**: Riccardo Campi, Matteo Bianchi, Antonio De Santis

# Black-box Systems

# Explainable Artificial Intelligence (XAI)



Input

CNN

Explanation

Concept: Striped

High importance

Concept: Grass

Low Importance

# Visual TCAV



Input image      CAV      Concept map + concept attributions
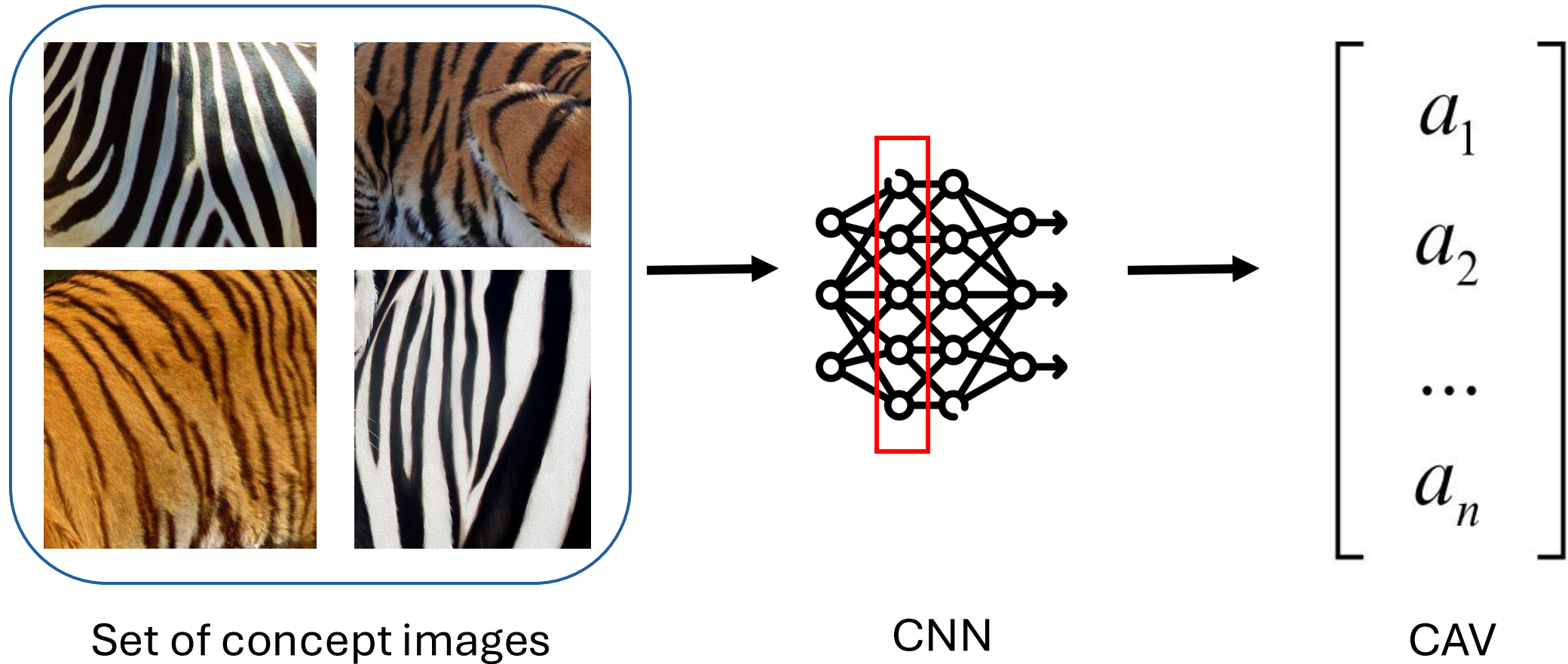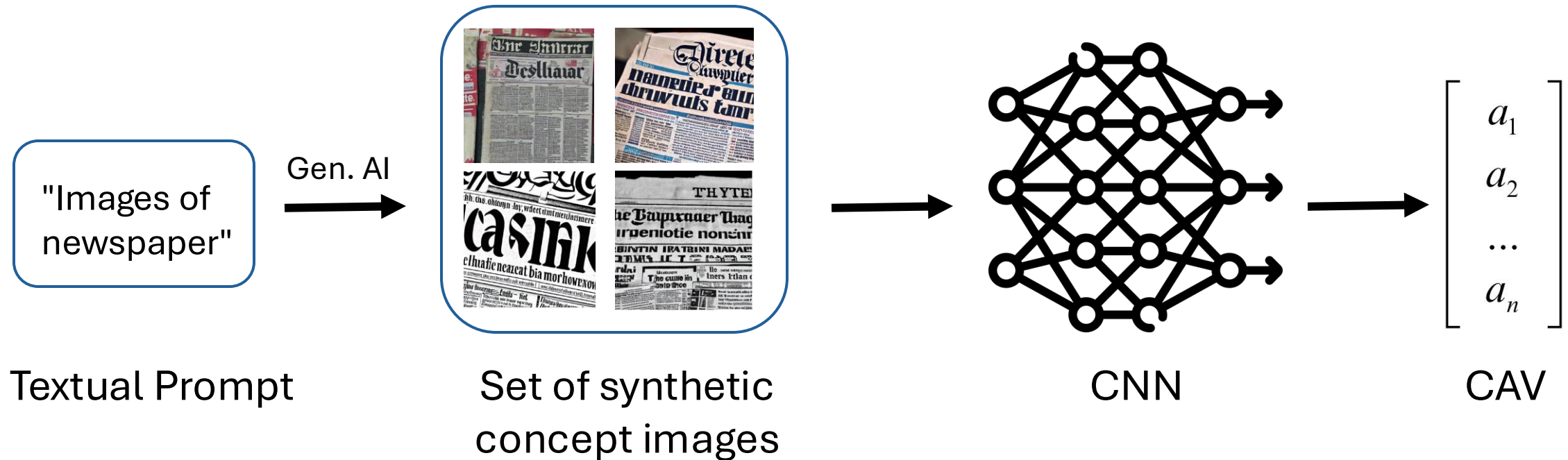
Antonio De Santis et al., Visual-tcav: Concept-based attribution and saliency maps for post-hoc explainability in image classification, 2025.
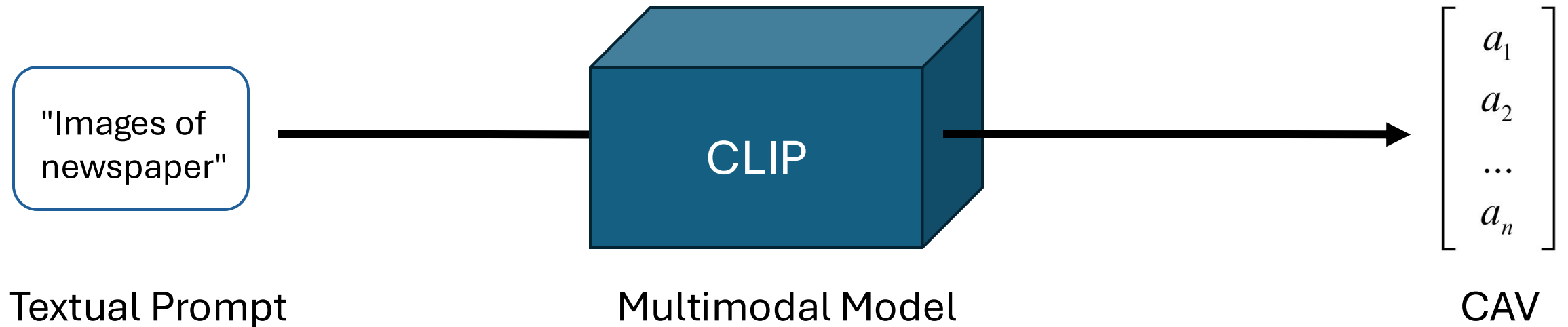
# Problem: Standard CAV Extraction



Set of concept images          CNN          CAV

POLITECNICO
MILANO 1863

# Other Solution: Synthetic Images



"Images of newspaper"

Gen. AI

$$\begin{bmatrix} a_1 \\ a_2 \\ \dots \\ a_n \end{bmatrix}$$

Textual Prompt
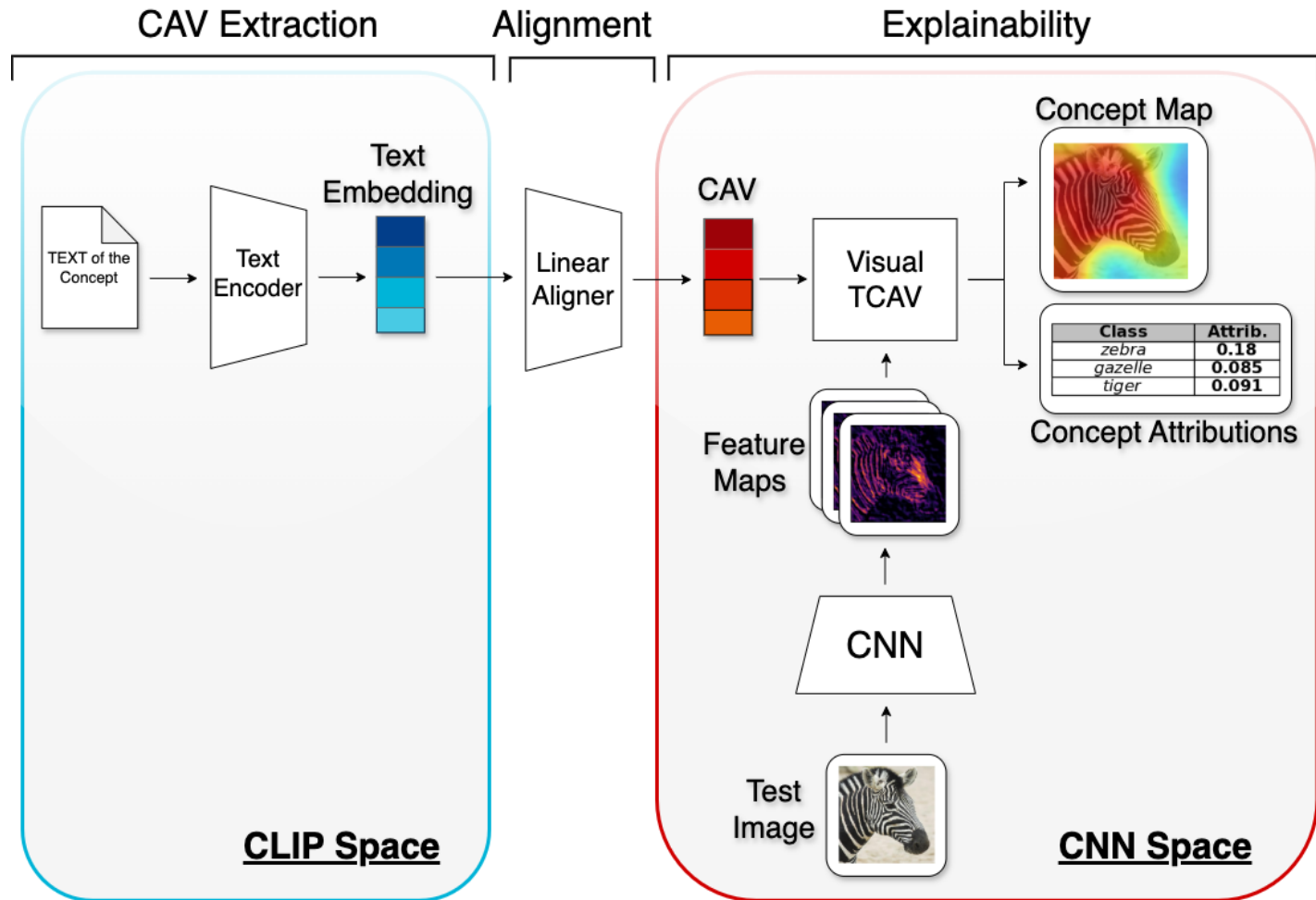
Set of synthetic concept images

CNN

CAV

Riccardo Campi et al., Towards synthetic concept activation vectors via generative models. In Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR) Workshops, pages 2720–2728, June 2025.
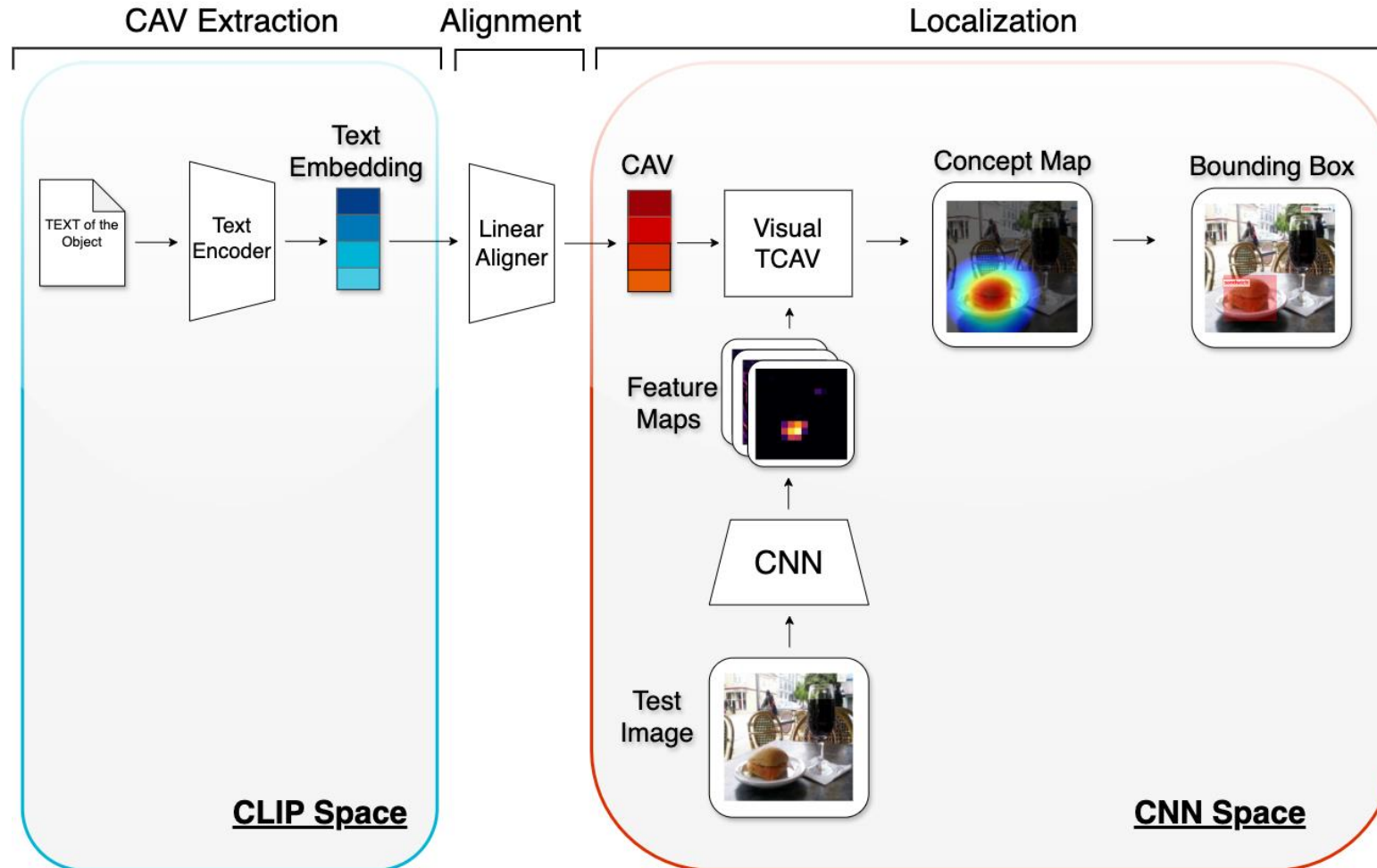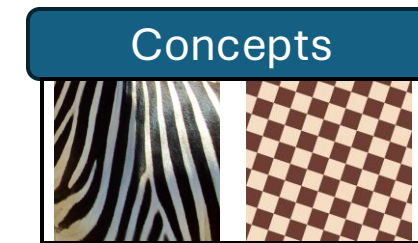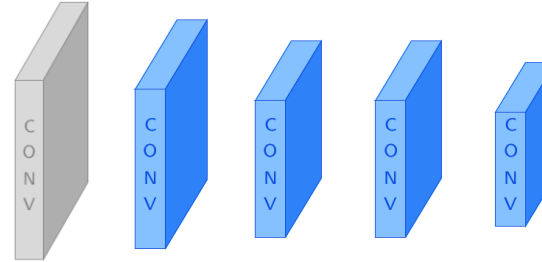
# Our idea: from text to CAV



Textual Prompt            Multimodal Model            CAV

# Our Solution

POLITECNICO MILANO 1863

# Other use: Localization

# Experiments Setup

# Linear Aligners Training



| Layer | MSE | R^2 |
|-------|-----|-----|
| layer1 | 0.492 | 0.891 |
| layer2 | 1.14 | 0.74 |
| layer3 | 1.798 | 0.6 |
| layer4 | 3.205 | 0.288 |

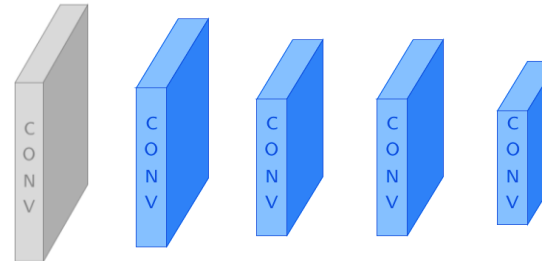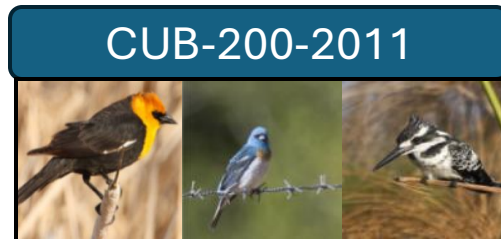| Layer | MSE | R^2 |
|-------|-----|-----|
| stage1 | 0.328 | 0.927 |
| stage2 | 0.819 | 0.818 |
| stage3 | 1.447 | 0.678 |
| stage4 | 2.853 | 0.366 |

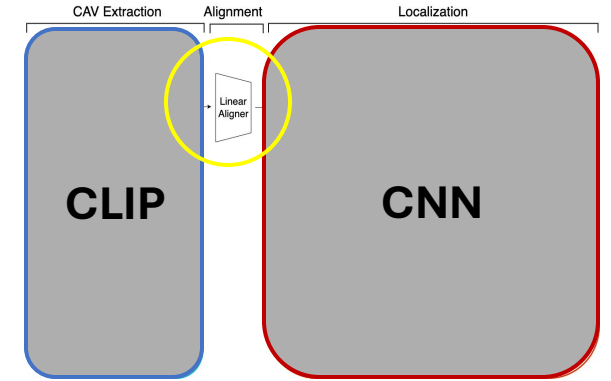Resnet50 pretrained with ImageNet-1k

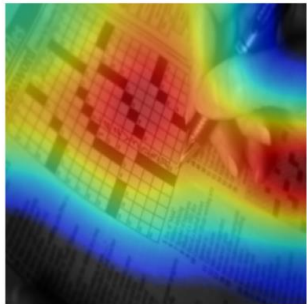Resnet18 pretrained with CUB-200-2011

# Experiments: ImageNet-1k



RESNET50 architecture
*chequered* concept

layer4 layer

| Class | Attrib. |
|---|---|
| crossword_pu.. | 0.1 |
| digital_watc.. | 0.022 |
| rule | 0.014 |

RESNET50 architecture
*hand* concept

layer4 layer

| Class | Attrib. |
|---|---|
| crossword_pu.. | 0.024 |
| digital_watc.. | 0.01 |
| rule | 0.011 |

RESNET50 architecture
*waffled* concept

layer4 layer

| Class | Attrib. |
|---|---|
| waffle_iron | 0.13 |
| mouse | 0.027 |
| computer_key.. | 0.041 |

RESNET50 architecture
*hand* concept

layer4 layer

| Class | Attrib. |
|---|---|
| waffle_iron | 0.012 |
| mouse | 0.0084 |
| computer_key.. | 0.008 |

RESNET50 architecture
*honeycombed* concept

layer4 layer

| Class | Attrib. |
|---|---|
| honeycomb | 0.2 |
| bee | 0.096 |
| apiary | 0.1 |

RESNET50 architecture
*bee* concept

layer4 layer

| Class | Attrib. |
|---|---|
| honeycomb | 0.007 |
| bee | 0.0051 |
| apiary | 0.0039 |

chequered          hand          waffled          hand          honeycombed          bee

# Experiments: ImageNet-1k



RESNET50 architecture
*legs* concept

| layer1 layer | | | layer2 layer | | | layer3 layer | | | layer4 layer | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Class** | **Attrib.** | | **Class** | **Attrib.** | | **Class** | **Attrib.** | | **Class** | **Attrib.** |
| zebra | 0.29 | | zebra | 0.035 | | zebra | 0.047 | | zebra | 0.082 |
| gazelle | 0.13 | | gazelle | 0.017 | | gazelle | 0.021 | | gazelle | 0.049 |
| impala | 0.15 | | impala | 0.017 | | impala | 0.019 | | impala | 0.039 |

RESNET50 architecture
*sky* concept

| layer1 layer | | | layer2 layer | | | layer3 layer | | | layer4 layer | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Class** | **Attrib.** | | **Class** | **Attrib.** | | **Class** | **Attrib.** | | **Class** | **Attrib.** |
| zebra | 0.019 | | zebra | 0.0e + 00 | | zebra | 0.008 | | zebra | 0.0067 |
| gazelle | 0.011 | | gazelle | 0.0e + 00 | | gazelle | 0.0033 | | gazelle | 0.0035 |
| impala | 0.015 | | impala | 0.0e + 00 | | impala | 0.0032 | | impala | 0.0041 |

RESNET50 architecture
*striped* concept

| layer1 layer | | layer2 layer | | layer3 layer | | layer4 layer | |
|---|---|---|---|---|---|---|---|
| **Class** | **Attrib.** | **Class** | **Attrib.** | **Class** | **Attrib.** | **Class** | **Attrib.** |
| zebra | 0.26 | zebra | 0.17 | zebra | 0.093 | zebra | 0.26 |
| gazelle | 0.12 | gazelle | 0.09 | gazelle | 0.032 | gazelle | 0.055 |
| impala | 0.14 | impala | 0.095 | impala | 0.027 | impala | 0.04 |

POLITECNICO
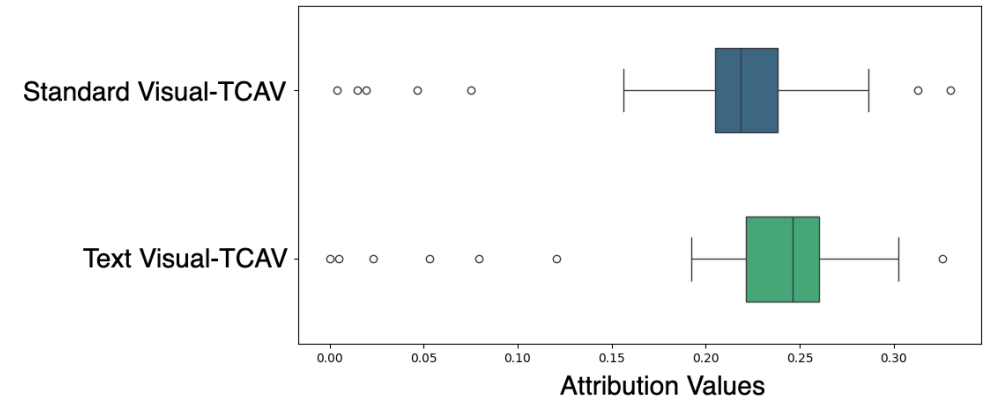MILANO 1863
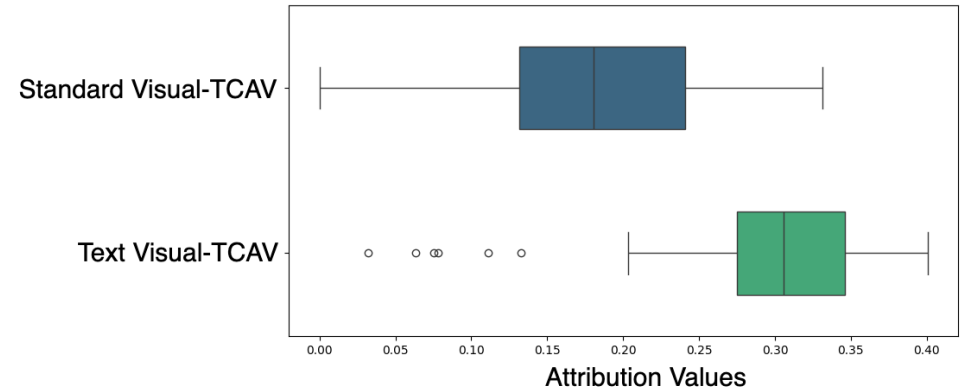
# Experiments: ImageNet-1k



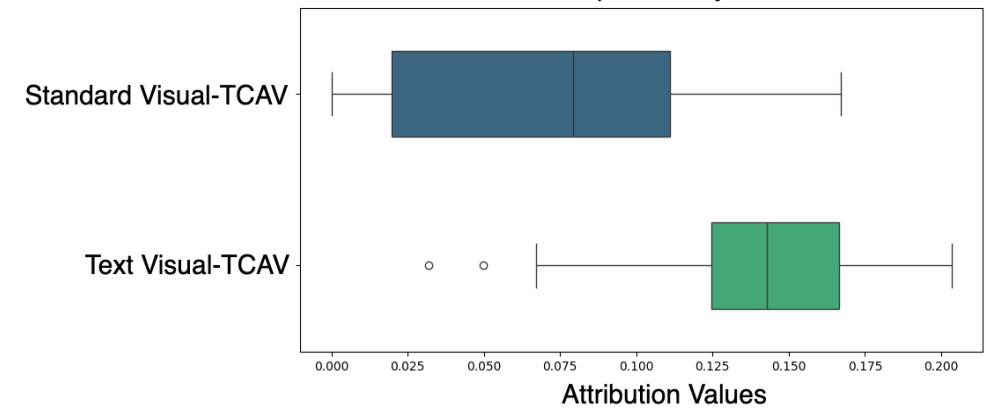Standard Visual-TCAV vs Text Visual-TCAV Attributions
waffled - layer4

Standard Visual-TCAV vs Text Visual-TCAV Attributions
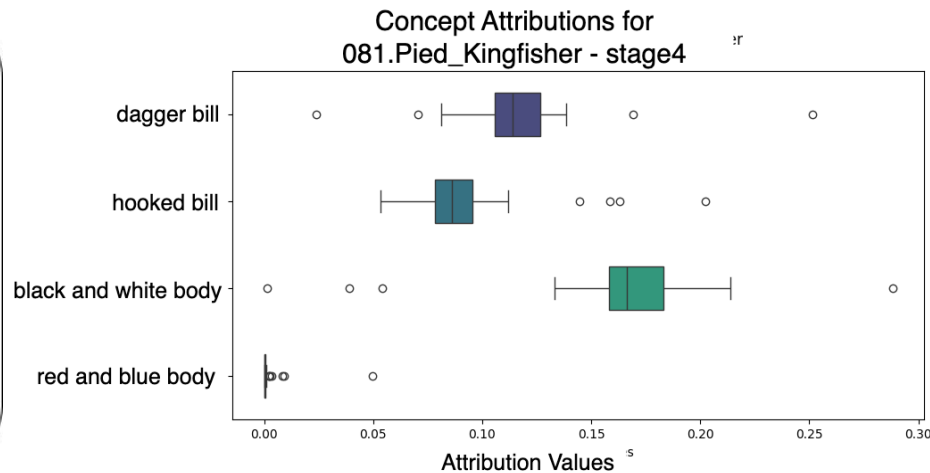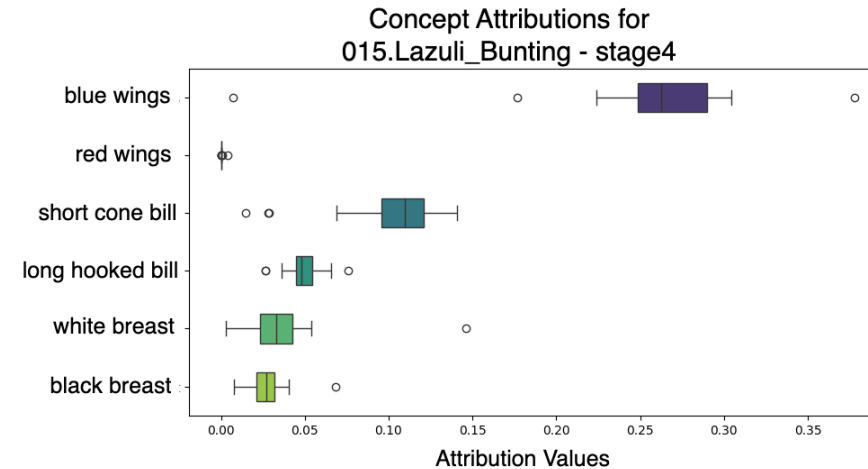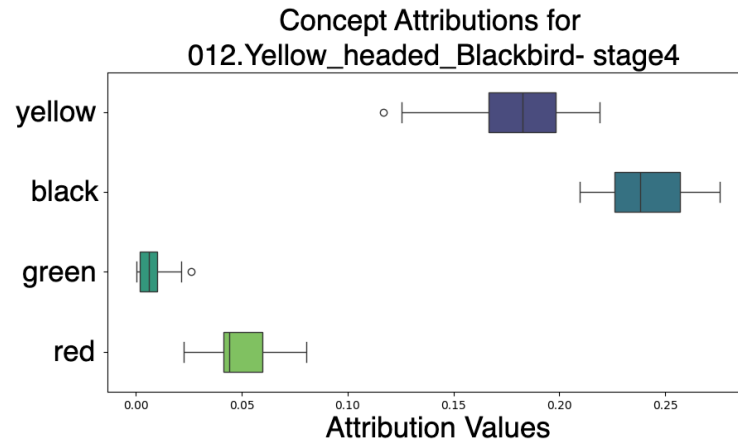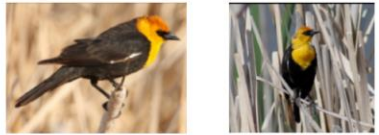striped - layer4

Standard Visual-TCAV vs Text Visual-TCAV Attributions
honeycombed - layer4

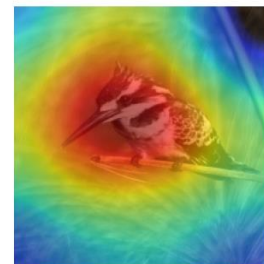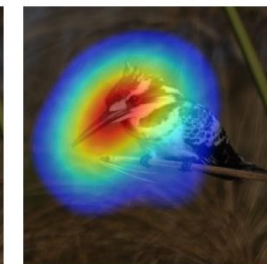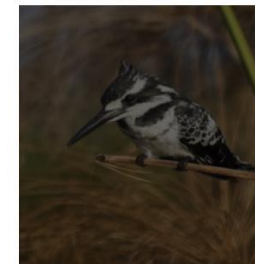Standard Visual-TCAV vs Text Visual-TCAV Attributions
chequered - layer4

POLITECNICO
MILANO 1863
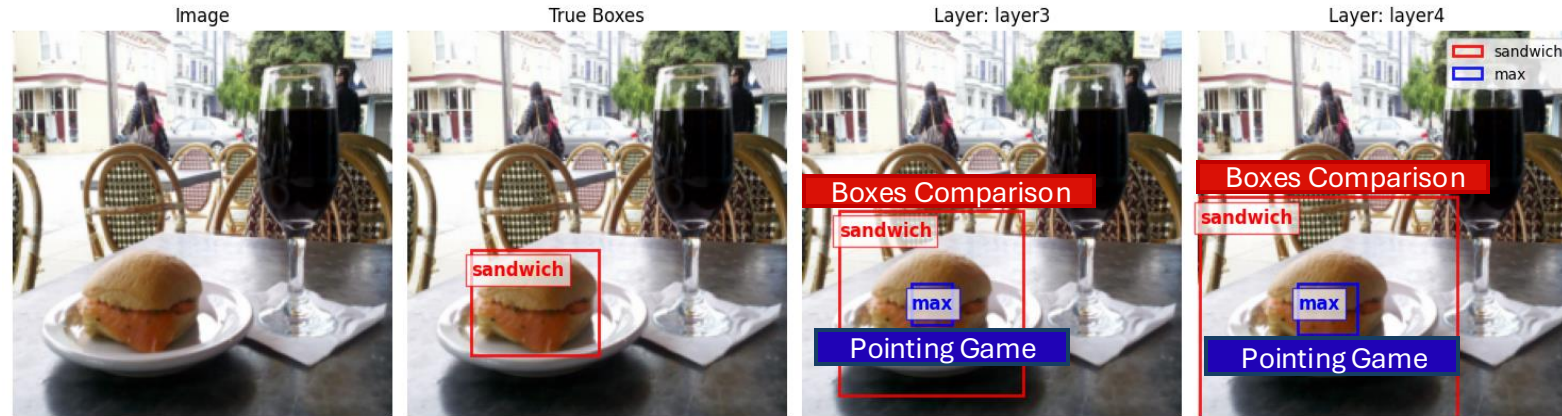
Concept Attributions for
012.Yellow_headed_Blackbird- stage4



Concept Attributions for
015.Lazuli_Bunting - stage4



Concept Attributions for
081.Pied_Kingfisher - stage4

RESNET18 architecture
daggerbill concept

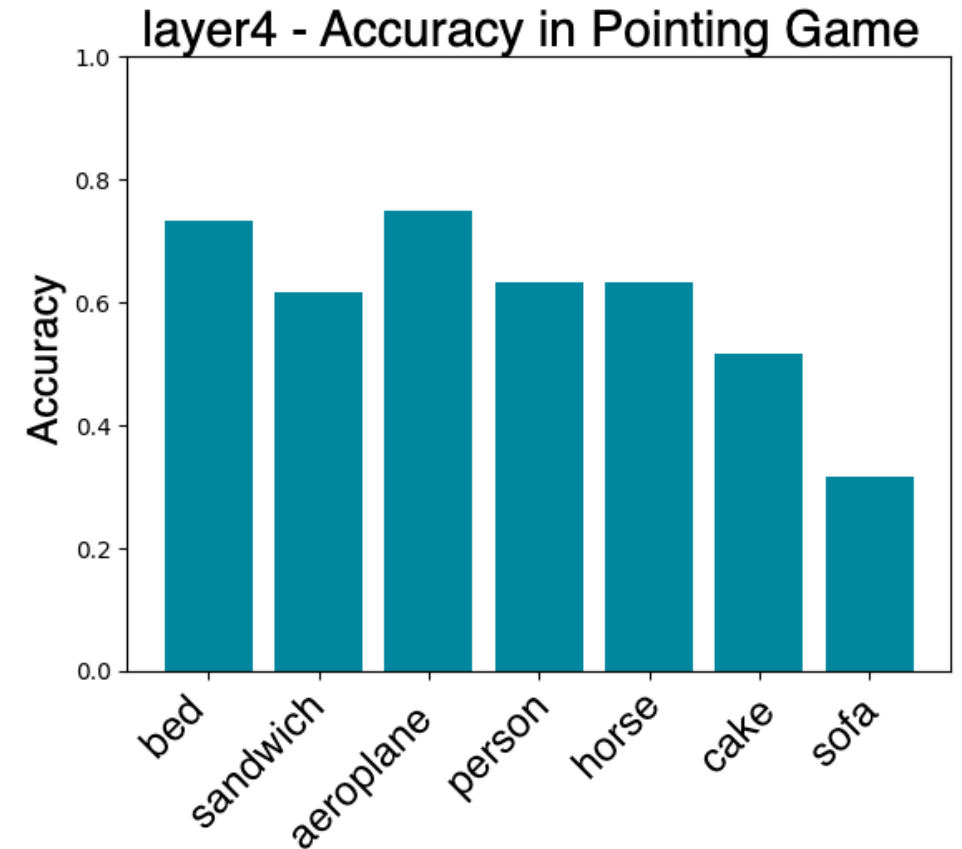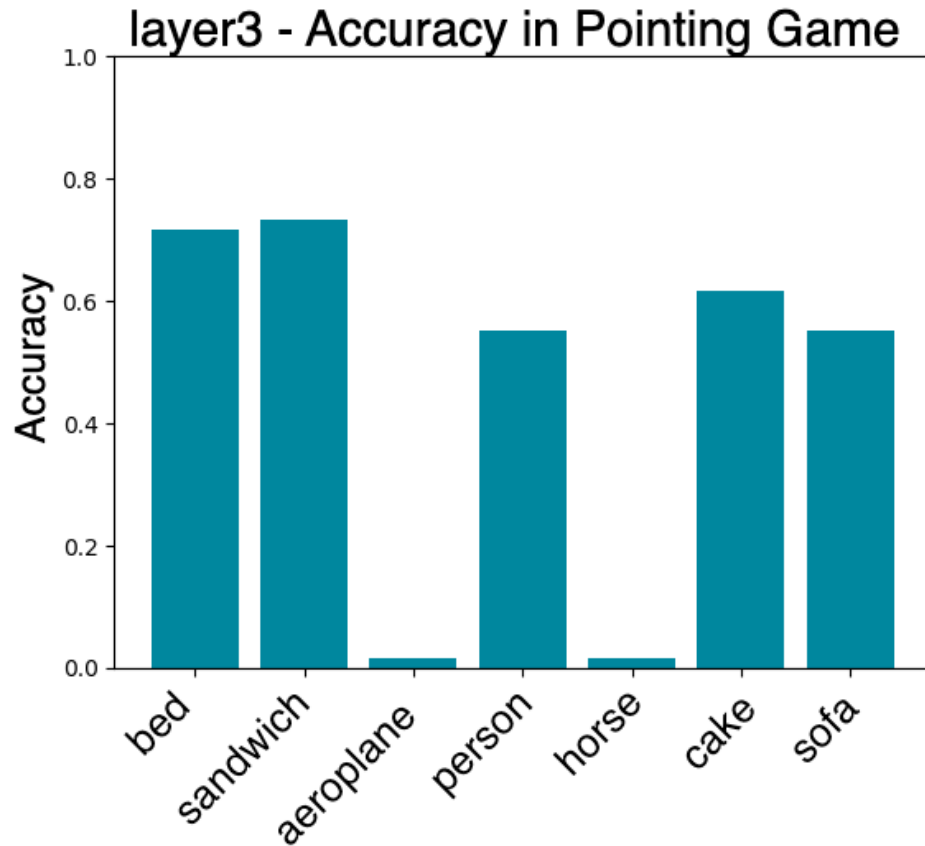| stage1 layer | | stage2 layer | | stage3 layer | | stage4 layer | |
|---|---|---|---|---|---|---|---|
| Class | Attrib. | Class | Attrib. | Class | Attrib. | Class | Attrib. |
| Pied_Kingfis.. | 0.0e+00 | Pied_Kingfis.. | 0.0e+00 | Pied_Kingfis.. | 0.012 | Pied_Kingfis.. | 0.12 |
| Black_and_wh.. | 0.0e+00 | Black_and_wh.. | 0.0e+00 | Black_and_wh.. | 0.0056 | Black_and_wh.. | 0.061 |
| Green_Kingfi.. | 0.0e+00 | Green_Kingfi.. | 0.0e+00 | Green_Kingfi.. | 0.0066 | Green_Kingfi.. | 0.07 |

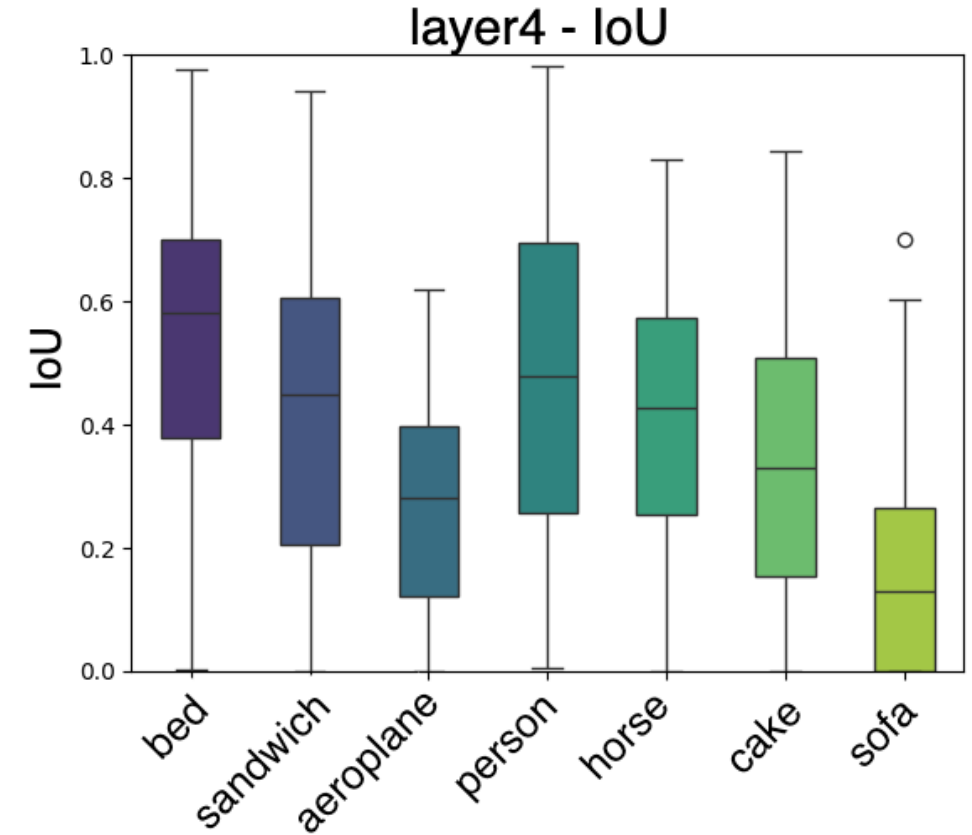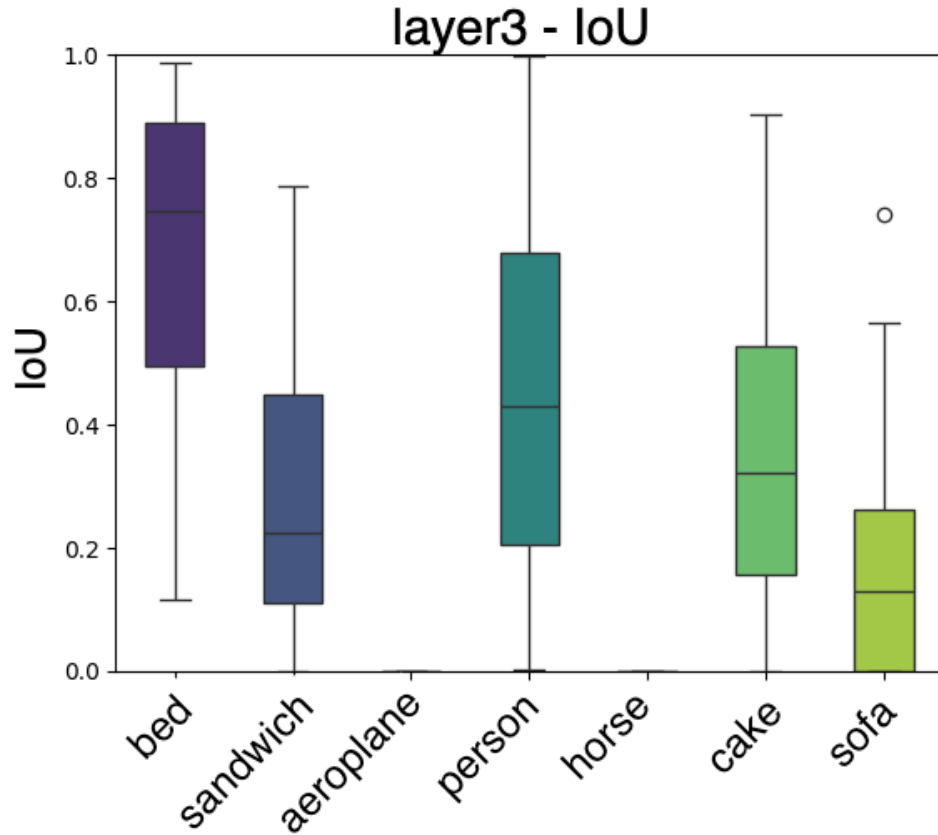# Experiments: Zero-Shot Localization with COCO



We implement two tasks to assess the performance of zero-shot localization with CAVs from text:
- Pointing Game (correct object position)
- Boxes comparison with IoU (bounding box covers most part of the object)

POLITECNICO
MILANO 1863

# Accuracy in Pointing Game



layer3 - Accuracy in Pointing Game



layer4 - Accuracy in Pointing Game

# IoU Distribution in Boxes Comparison

# Conclusions

## CAVs from text

- **High scalability** for CAV extraction

- **No concept images** for CAV extraction

- **Consistent explanations**

- **Localization** tasks without an annotated dataset

- **Zero-shot localization**

## Future work

- More **complex architectures** for the **aligners**

- More roboust ways to **refine** the **concept maps**

# Thank you for your attention