# Visual Question Answering

Project for the "Computer Vision" course
A.Y. 2022/23

**Clizia Giorgia Manganaro**  manganaro.2017897@studenti.uniroma1.it
**Chiara Giacanelli**  giacanelli.1801145@studenti.uniroma1.it
**Alessio Palma**  palma.1837493@studenti.uniroma1.it
**Davide Santoro**  santoro.1843664@studenti.uniroma1.it
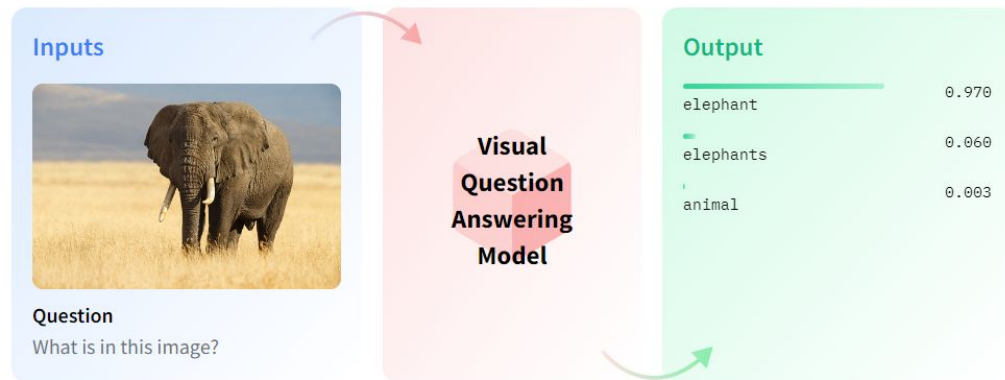
# What's Visual Question Answering?

Visual Question Answering (VQA) is a **computer vision task where a system is given a text-based question about an image, and it must infer the answer.**



## Why it's important?

- **to help blind users** to communicate through pictures;

- **to attract customers of online shopping** sites by giving "semantically" satisfying results for their search queries;

- **Visual Dialogue**, which aims to give natural language instructions to robots.
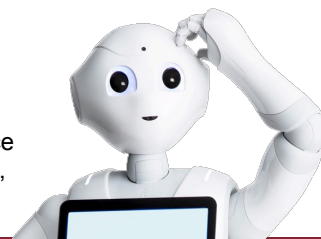
## Related works

P. Anderson et al. (2018), *"Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering."* 6077-6086. 10.1109/CVPR.2018.00636;

S. Antol et al., "*VQA: Visual Question Answering",* 2015 IEEE International Conference on Computer Vision (ICCV), 2015, pp. 2425-2433, doi: 10.1109/ICCV.2015.279.

D. Teney et al., *"Tips and Tricks for Visual Question Answering: Learnings from the 2017 Challenge"*, 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp;.

# Where did we start from?

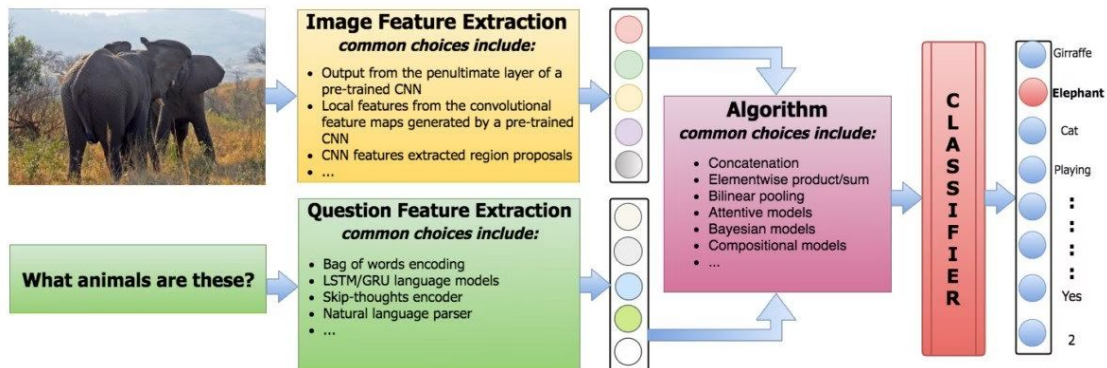We had to face the **two sides** of the project model



## Computer Vision
**Image Feature Extraction**

converting images into their feature representations for further processing

## NLP
**Question Feature Extraction**

converting natural language questions into their embeddings for further processing

🏆 The goal of our project was to **analyze and test** the VQA model in terms of the structure, hyperparameters, attention and to see if using **Transformer Networks** would positively affect the performance

VQA Annotations

Balanced Real Images [Cite]

- Training annotations 2017 v2.0*
  4,437,570 answers
- Validation annotations 2017 v2.0*
  2,143,540 answers

VQA Input Questions

- Training questions 2017 v2.0*
  443,757 questions
- Validation questions 2017 v2.0*
  214,354 questions
- Testing questions 2017 v2.0
  447,793 questions

VQA Input Images

COCO
- Training images
  82,783 images
- Validation images
  40,504 images
- Testing images
  81,434 images

an ==entry== is made of this structure

```
{'question_id': 9000,
 'image_id': 9,
 'image': 52181,
 'question': 'How many cookies can be seen?',
 'answer': {'labels': tensor([17], dtype=torch.int32), 'scores': tensor([1.])},
 'q_token': tensor([19901, 19901, 19901, 19901, 19901, 19901, 19901, 19901,   77,   78,
       949,   80,   81,   82], dtype=torch.int32)}
```

Baseline code: An efficient PyTorch implementation of the winning entry of the *2017 VQA Challenge:*
https://github.com/hengyuan-hu/bottom-up-attention-vqa
**translated by us to Python3 and a newer version of PyTorch**

# Bottom-Up and Top-Down Attention for Image Captioning and VQA[1]

**Bottom-up** + **top-down attention** mechanism = to calculate attention at the level of objects and other salient image regions.

## Bottom-Up Attention Model

We can define some stages:

1. **Region Proposal Network (RPN)** that allows to predicts object proposals
2. **Region of interest (RoI) pooling:** used to extract a small feature map for each box proposal
3. **Final output:** softmax distribution & class-specific bounding box refinements

They used a Faster R-CNN initialized with ResNet-101 weights to pretrain the model.



Figure 2. Example output from our Faster R-CNN bottom-up attention model. Each bounding box is labeled with an attribute class followed by an object class. Note however, that in captioning and VQA we utilize only the feature vectors – not the predicted labels.

[1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Tene, Mark Johnson, Stephen Gould, Lei Zhang
Australian National University, JD AI Research, Microsoft Research, University of Adelaide Macquarie University

Visual Question Answering -
Clizia Manganaro, Chiara Giacanelli, Alessio Palma, Davide Santoro

# Bottom-Up and Top-Down Attention for Image Captioning and VQA[1]

## Visual Question Answering Model

1. **Joint embedding of the question** with:
   - a gated recurrent unit (GRU), with each input word represented using a learned word embedding,
   - image features {v1, …, vk}

2. **Unnormalized attention weight** for each of the k image features

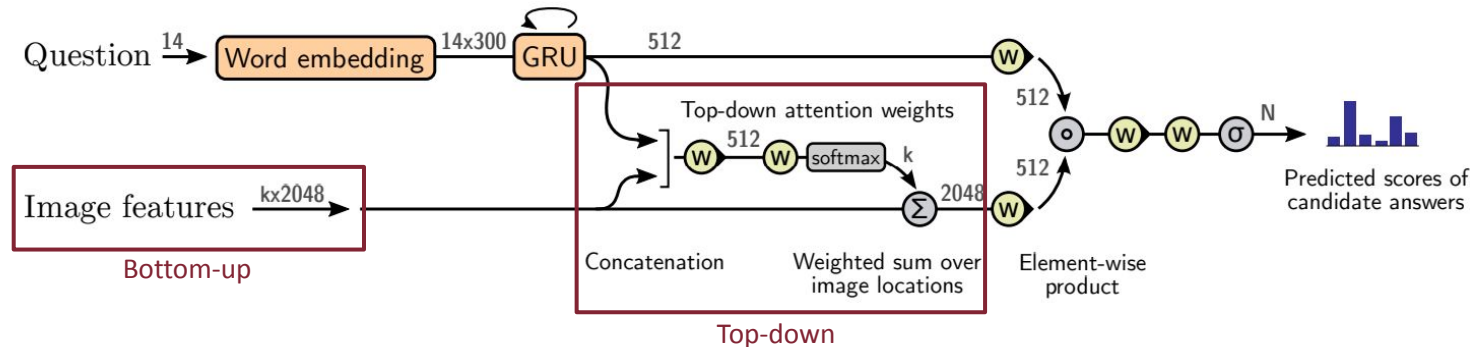Output is generated by a multi-label classifier operating over a fixed set of candidate answers



Question: What room are they in? Answer: kitchen

Figure 6. VQA example illustrating attention output. Given the question 'What room are they in?', the model focuses on the stove-top, generating the answer 'kitchen'.

[1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Tene, Mark Johnson, Stephen Gould, Lei Zhang
Australian National University, JD AI Research, Microsoft Research, University of Adelaide Macquarie University

# So…the benchmark model



- Bottom-up attention: image features are coming from a Faster R-CNN initialized with ResNet-101 weights and fine-tuned on Visual Genome;

- Top-down attention: NLP-based attention, queries are question embeddings, keys and values are image features;

- **GloVe** pre-trained embeddings;

- Fusion is done through **element-wise product**

Reference: Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, & Lei Zhang (2018). "Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering". In CVPR.

# Dataset & Evaluation Metric

**VQA v2.0** dataset: https://visualqa.org/

This dataset resolves some bias problems of the v1. It was used as the basis of the 2017 VQA Challenge, contains 1.1M questions with 11.1M answers relating to MSCOCO images.



Who is wearing glasses?
man                  woman

Where is the child sitting?
fridge               arms

Is the umbrella upside down?
yes                   no

How many children are in the bed?
2                     1

**Metric**

**Standard VQA accuracy**, which is robust to inter-human variability in phrasing the answers:

$$\text{Acc}(ans) = \min\left\{\frac{\#\text{humans that said } ans}{3}, 1\right\}$$

An answer is deemed *100% accurate* if at least 3 unique annotators provided that exact answer. Instead, it is assigned a fractional score if it produces an answer that is deemed *rare*.

# Our experiments

- **Feature Extraction** from images with:
    1. VGG16
    2. SIFT

- **Word Embedding**: varying the embedding dimension and model (we tested DistilBERT frozen and with fine-tuning); 🤗

- **Question Embedding**: training different models to produce an embedding for the question given the embeddings of single words (TCN, BiGRU, BiLSTM, …);

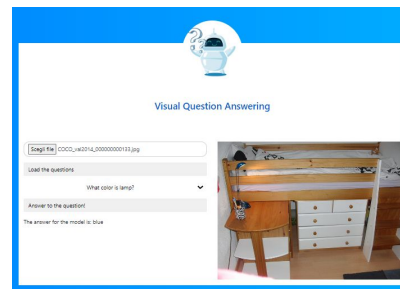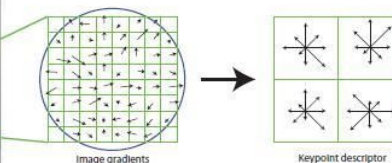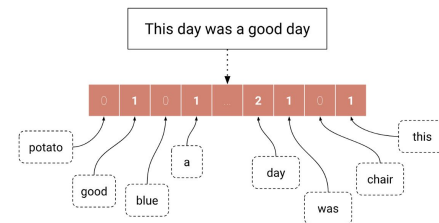- Analysis on the **gradient flow** and of the **data variance** preserved;

- **Attention analysis** and building the **bounding box** that describes where the model pays attention in the image;

- Added a **new Calibration layer** based on the "Reliable VQA" paper, according to which the model will output an answer only if deemed reliable;

- Analyzed **if the model is able to reason** and answer questions that require common sense;
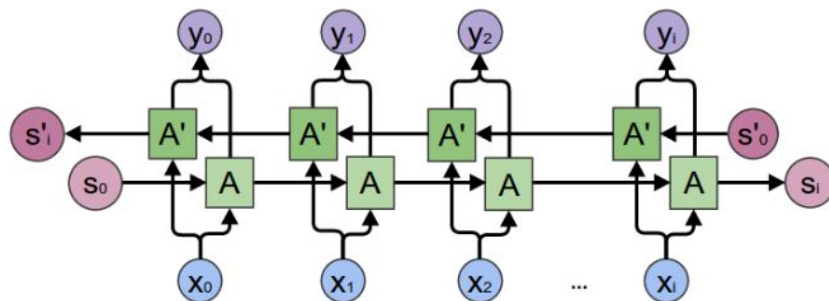
- Built a **Demo**.

# Different question embeddings



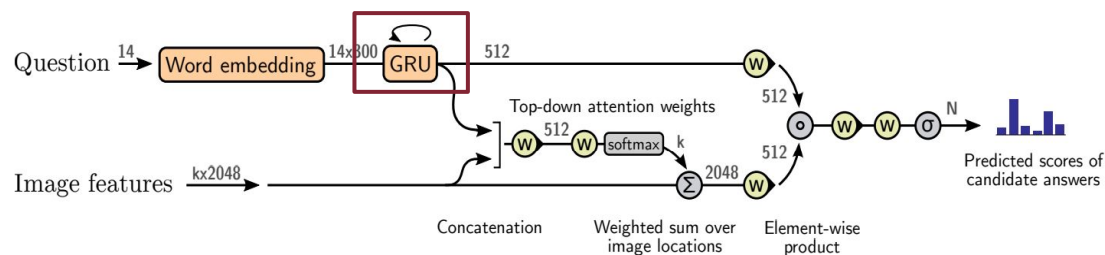| Embedding dimension | Validation VQA accuracy | Epochs of training |
|---|---|---|
| GRU (baseline model) | 62.70 | 10 |
| LSTM | 62.77 | 10 |
| BiGRU | 63.14 | 10 |
| BiLSTM | 62.74 | 10 |
| **BiGRU2x** | **63.70** | **10** |
| BiLSTM2x | 63.41 | 10 |
| TCN | 58.87 | 10 |

# Why should we go bidirectional?



Bidirectional recurrent neural networks connect two RNNs that process the input in opposite directions, and then concatenates the outputs. With this architecture, the output layer can get information from past and future states simultaneously.

E.g.: *"What color is the child's hat?"*

# Different RNNs initialization

| Parameter | TensorFlow/Keras | PyTorch |
|-----------|------------------|---------|
| weight_ih | xavier uniform | uniform $\sqrt{hidden\_size}$ |
| weight_hh | orthogonal | same as above |
| bias | 1 for forget gate, 0 otherwise | same as above |
| linear | xavier uniform | uniform $\sqrt{input\_size}$ |

} GRU

} LSTM



Obtained **+0.39 accuracy** on the GRU model and **+0.22** on the LSTM one

Reference: Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. 2015. "An empirical exploration of recurrent network architectures". In Proceedings of the 32nd International Conference on International Conference on Machine Learning.

# Features Extraction and Performance

We used three models to extract the image features:

1. **VGG16**, a deep CNN architecture with 16 layers, used for object recognition;
2. **SIFT**, a CV algorithm to detect local keypoints in images;
3. **Faster R-CNN**, region proposal algorithm + CNN feature vector for each region.

We set 10 epochs to perform our training and then we compared the performance obtained from these models.
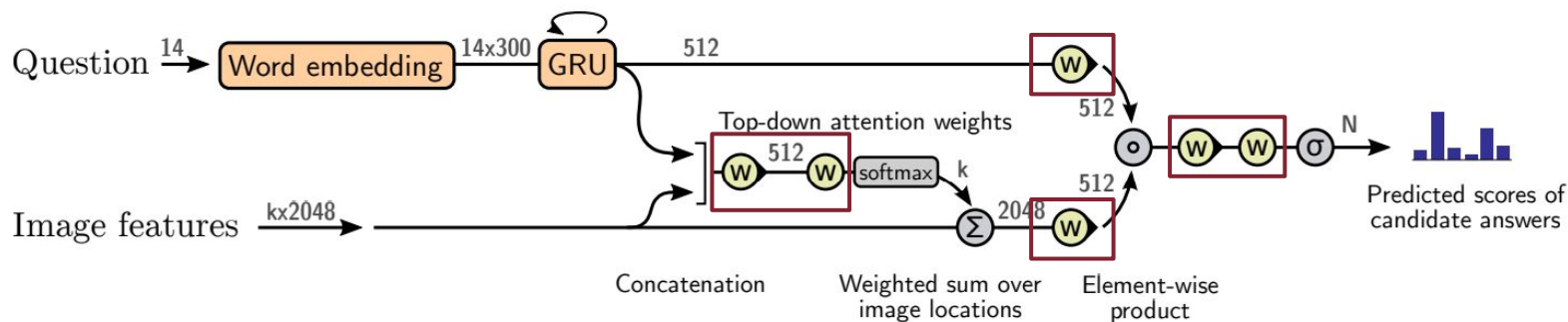


1) Region proposals with selective search

2) Feature extraction with CNN

| Image Features | Question Embedding | Validation VQA accuracy |
|---|---|---|
| VGG16 | GRU | 44.46 |
| VGG16 | BiGRU2x | 45.22 |
| SIFT | GRU | 45.01 |
| **Faster R-CNN (baseline model)** | **GRU** | **62.70** |

# Different activation functions and gradient flow

# Gradient flow

- Gradients have high variance;

- For some layers, the maximum value is much higher than the mean

**Weight normalization** is used to speed up the training

Reference: Tim Salimans and Diederik P. Kingma. 2016. "Weight normalization: a simple reparameterization to accelerate training of deep neural networks". 30th International Conference on Neural Information Processing Systems (NIPS'16).



Gradient flow

Visual Question Answering -
Clizia Manganaro, Chiara Giacanelli, Alessio Palma, Davide Santoro

# Gradient flow

- Gradients have high variance;

- For some layers, the maximum value is much higher than the mean;

- Solution: **clip the norm of the gradients**;

- Prevents exploding gradient, but…

# …is there any vanishing gradient?



Gradient flow

- Mostly depends on the chosen activation function;

- For the Sigmoid, gradients are vanishing

# …is there any vanishing gradient?


Gradient flow

- Mostly depends on the chosen activation function;

- For the Sigmoid, gradients are vanishing;

- ReLU is behaving much better;

We also tested **LeakyReLU, ELU and SELU,** but gradient flow and performances are very similar to ReLU

# Deep ReLU networks initialization and data variance

# Default PyTorch initialization

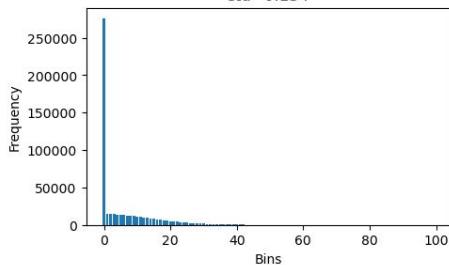$$\mathcal{W} \sim \mathcal{U}(\sqrt{\frac{2}{n}}, \sqrt{\frac{2}{n}})$$

q_net layer 0
mean=0.021
std=0.03

q_net layer 1
mean=0.011
std=0.016

q_net layer 2
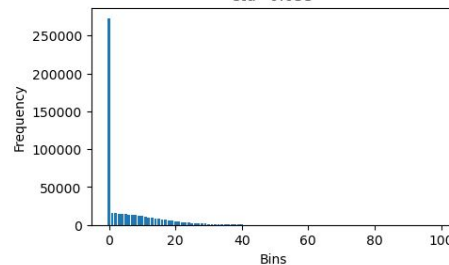mean=0.008
std=0.012

q_net layer 3
mean=0.008
std=0.011

**1**   **2**   **3**   **4**
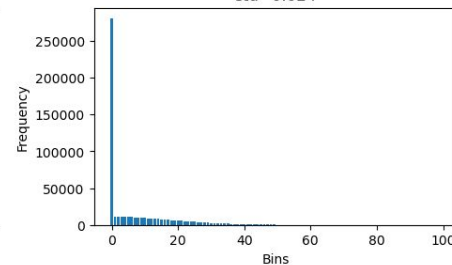
<u>layers</u>

v_net layer 0
mean=0.214
std=0.324

v_net layer 1
mean=0.09
std=0.134

v_net layer 2
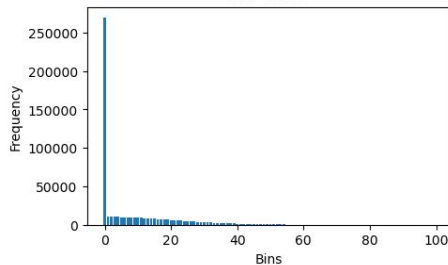mean=0.037
std=0.055

v_net layer 3
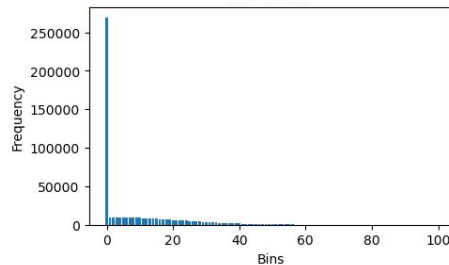mean=0.016
std=0.024

# Our initialization

$$\mathcal{W} \sim \mathcal{N}(0, \sqrt{\frac{2}{n}})$$
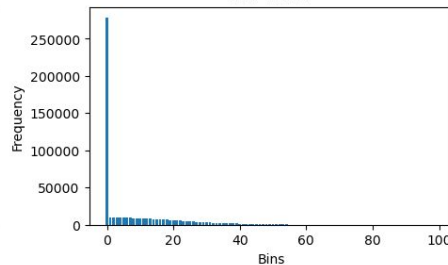
q_net layer 0
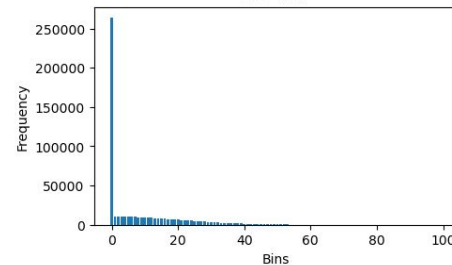mean=0.049
std=0.072

q_net layer 1
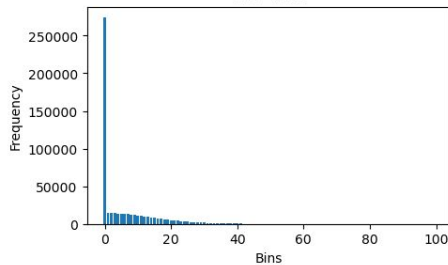mean=0.047
std=0.068

q_net layer 2
mean=0.045
std=0.067

q_net layer 3
mean=0.049
std=0.07
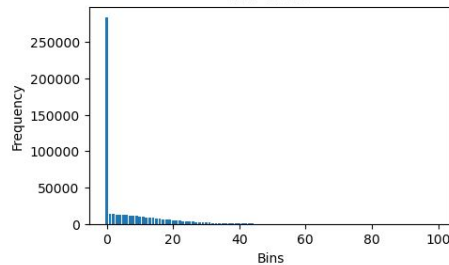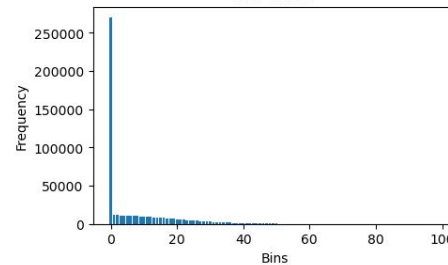
**1**     **2**     **3**     **4**
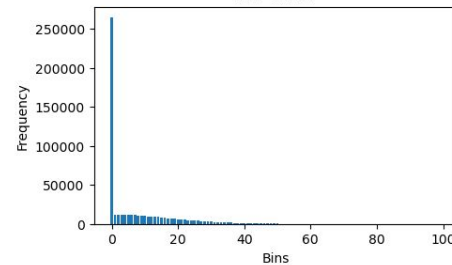
layers

v_net layer 0
mean=0.548
std=0.81

v_net layer 1
mean=0.522
std=0.803

v_net layer 2
mean=0.54
std=0.789

v_net layer 3
mean=0.551
std=0.797

# Attention on the image and reasoning

We decide to construct a **bounding box** which describes where the model pays attention to the image after processing the question



Q: What is the food?
GT: pizza    Our: pizza

The model is answering good and focusing on the right object



Q: Is that a mirror or a window?
GT: mirror    Our: window

The model can't distinguish effectively if the object is a window or a mirror, but it is still focusing it's attention on the right image patch.



Q: Is the fence made of wood?
GT: no    Our: no

The model is answering well, it is mainly focusing on a part of the image that is completely uncorrelated from the question, so it's not correctly seeing the wood structure in the very background.
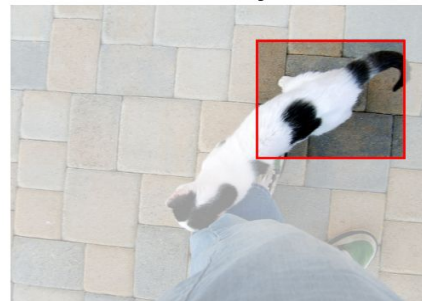
**What happens when the input question does not require simple object counting or object/material classification?**
We designed some questions on some images that require commonsense knowledge to be answered.
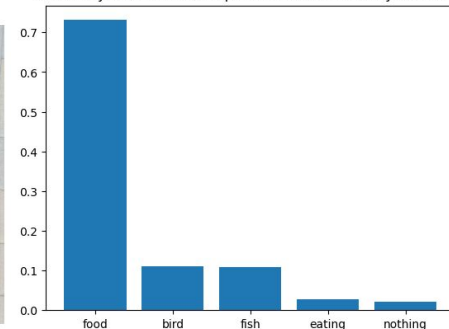
**Is the model able to reason or is it just a stochastic parrot?**

We also extracted the top-5 answers for each question and plotted their probability distribution, to check the confidence of the model and if it is coherent.
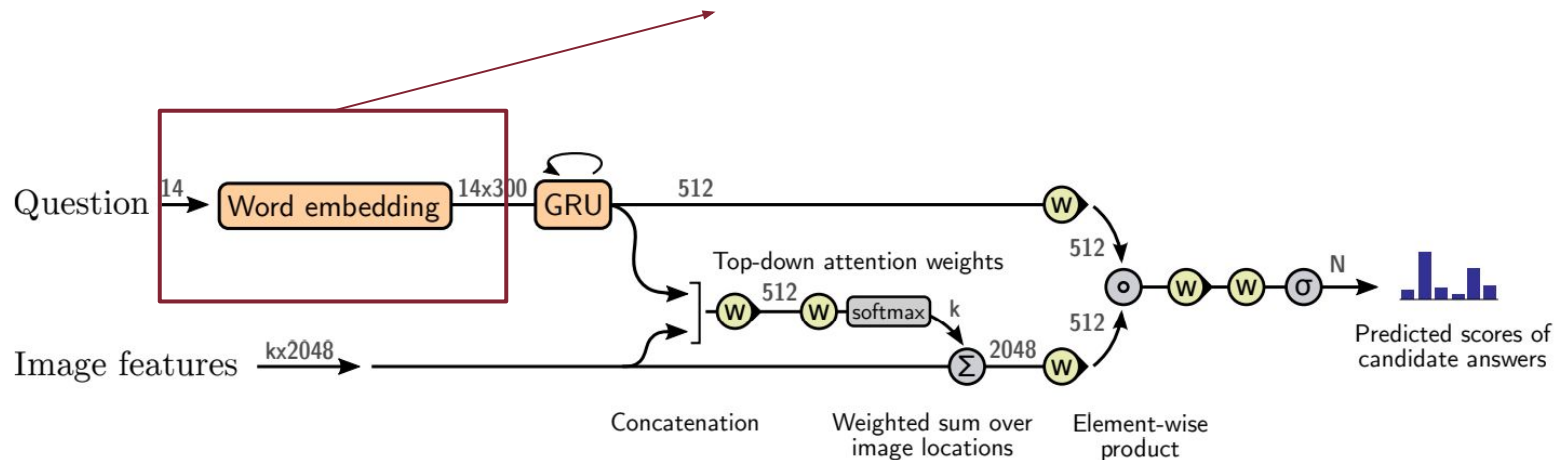


Q: What is the cat looking for?
Commonsense: cuddling    Our: food

Probability distribution of top-5 answers returned by the model

# We then worked in this area

# Word Embeddings with GloVe & DistilBert Transformers

**Standard VQA accuracy**,
which is robust to inter-human variability
in phrasing the answers:

$$\text{Acc}(ans) = \min\left\{\frac{\#\text{humans that said } ans}{3}, 1\right\}$$

There's a direct proportion between the increment of the dimension and the increment of the accuracy

| Embedding dimension | Validation VQA accuracy | Epochs of training |
|---|---|---|
| GloVe50 | 61.49 | 10 |
| GloVe100 | 62.17 | 10 |
| GloVe200 | 62.64 | 10 |
| **GloVe300 (baseline model)** | **62.70** | **10** |

# Word Embeddings with GloVe & DistilBert Transformers 🤗

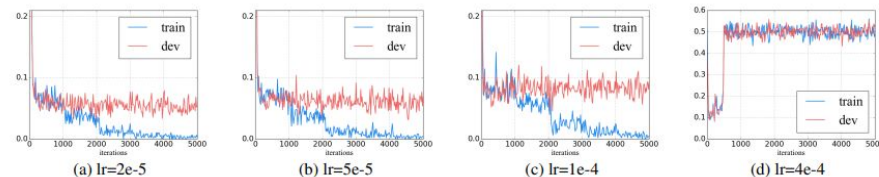| Embedding dimension | Validation VQA accuracy | Epochs of training |
|---|---|---|
| GloVe300 (baseline model) | 62.70 | 10 |
| **DistilBert Freezed Out** | **62.73** | **10** |
| DistilBert Fine-tuned (lr = 2e-3) | 25.43 | 7 |
| DistilBert Fine-tuned (lr = 2e-5) | 50.11 | 7 |

## Catastrophic forgetting!



Figure 2: Catastrophic Forgetting

Reference: "How to Fine-Tune BERT for Text Classification?" (2020), Shanghai Key Laboratory of Intelligent Information Processing, Fudan University School of Computer Science, Shanghai, China

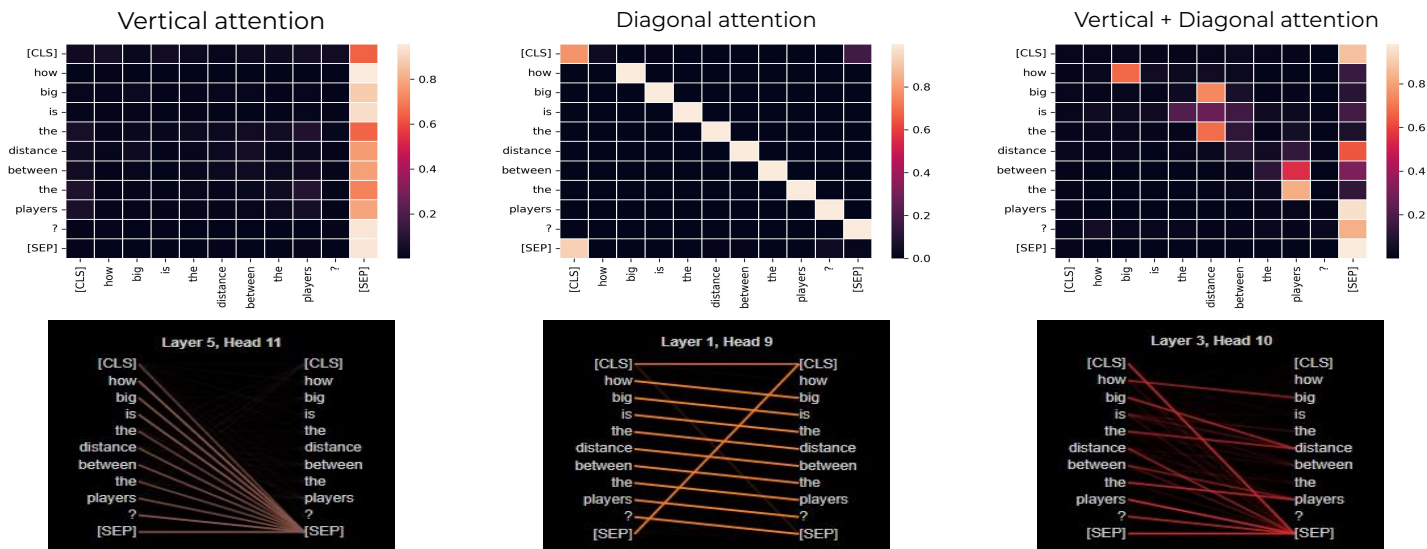# Word Embeddings with GloVe & DistilBert Transformers 🤗



Learning Rate Test

| Embedding dimension | Question Embedding | Validation VQA accuracy | Epochs of training |
|---|---|---|---|
| GloVe300 (baseline model) | GRU | 62.70 | 10 |
| DistilBert Freezed Out | GRU | 62.73 | **10** |
| DistilBert Fine-tuned (lr = 2e-3) | GRU | 25.43 | 7 |
| DistilBert Fine-tuned (lr = 2e-5) | GRU | 50.11 | 7 |
| **DistilBERT fine-tuned** (lr = 2e-4) | **BiGRU2x** | **63.51** | **10** |
| DistilBERT frozen (lr = 2e-4) | BiGRU2x | 63.11 | 10 |

# Attention Analysis

An attention analysis was carried out, in particular we analyzed both the different types of attention on the questions and how the model pays attention to the image.

We decided to analyze the **attention weights** of the pretrained DistilBERT model for some specific questions.

We can distinguish three different types of attentions:



Vertical attention

Diagonal attention

Vertical + Diagonal attention

Layer 5, Head 11

Layer 1, Head 9

Layer 3, Head 10

Visual Question Answering -
Clizia Manganaro, Chiara Giacanelli, Alessio Palma, Davide Santoro

# Some interesting (positive) examples



"What is the dog doing?"

**laying down = 0.3728**

resting = 0.1825

sitting = 0.1514

…

"What drink is shown?"

**wine = 9.9793e-01**

champagne = 8.1414e-04

alcohol = 4.5826e-04

…

"What color are the gym shoes?"

**white = 0.2967**

gray = 0.1730

red = 0.1431

…

# Some interesting (negative) examples

How big is the distance between the players?



**small = 0.1935**

big = 0.1545

large = 0.1459

…

How many dishes of food are in the picture?



**8 = 0.1628**

7 = 0.1405

6 = 0.1192

…

What time is it?
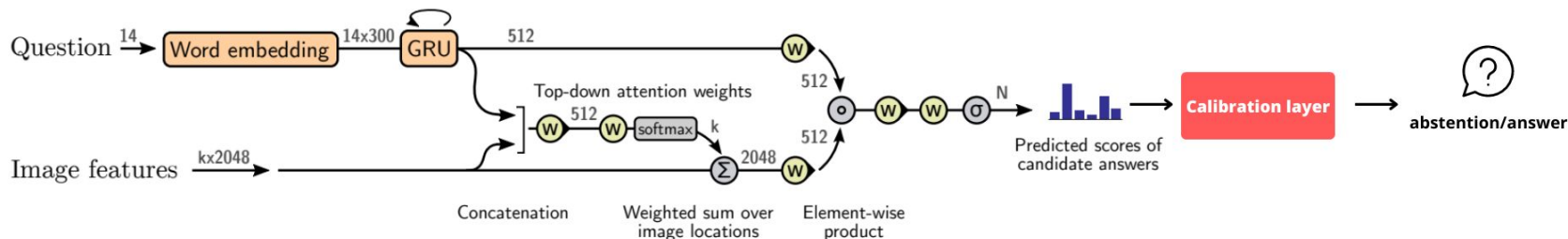


**noon = 0.0928**

daytime = 0.0425

afternoon = 0.0347

…

# Reliable Visual Question Answering

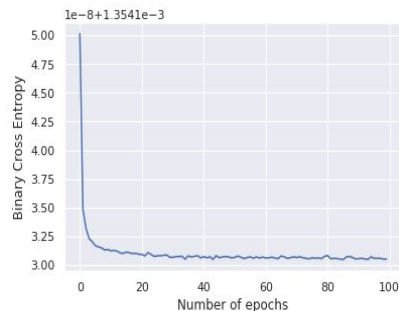What if the model responds incorrectly?



We implemented **selection function**: its output is compared with a threshold and, if it is greater, than the model is enabled to answer, otherwise it's not.

The Calibration Layer takes in input the answer logits of the VQA model and outputs the calibrated logits.

Reference: "Reliable Visual Question Answering: Abstain Rather Than Answer Incorrectly", Spencer Whitehead, Suzanne Petryk ,, Vedaad Shakib , Joseph Gonzalez , Trevor Darrell, Anna Rohrbach , and Marcus Rohrbach; Meta AI, UC Berkeley

# Reliable Visual Question Answering

Our results



**Effective Reliability metric**

$$\Phi_c(x) = \begin{cases} Acc(x) & \text{if } g(x) = 1 \text{ and } Acc(x) > 0, \\ -c & \text{if } g(x) = 1 \text{ and } Acc(x) = 0, \\ 0 & \text{if } g(x) = 0. \end{cases}$$

**Coverage**

(the percentage of questions answered by the model)
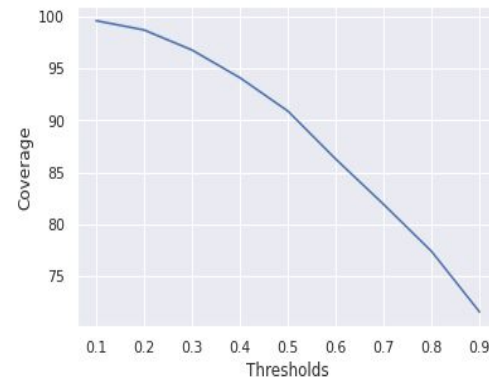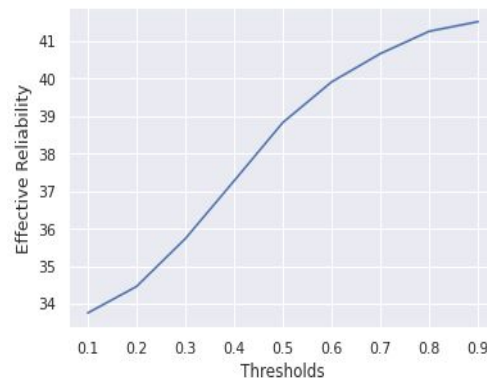
Q: How many birds are in the tree?



Answer: 7
Ground truth: 8
Selector: Abstain

Q: What is the man doing?



Answer: Talking on phone
Ground truth: Talking on phone
Selector: Answer

# Testing step…Demo!

# Demo Real Time

# Thanks for your attention!