

Visual Question Answering with DeepProbLog

Jorrit Willaert¹

¹Catholic University of Leuven, Leuven – Belgium
e-mail: jorrit.willaert@student.kuleuven.be

Abstract – TODO

Keywords – Neuro Symbolic AI, Visual Question Answering, DeepProbLog, Problog, Convolutional Neural Networks

I. INTRODUCTION

The Neuro Symbolic AI field is interested in building a bridge between the robustness of probabilistic knowledge, with the well-known popularity and proven strengths of deep neural networks. DeepProbLog [1] offers this ability, by using both the strengths of neural networks (i.e. system 1, typical subconscious tasks such as visual recognition, the processing of languages, ...), along with the strengths of rule-based probabilistic systems (i.e. system 2, slow, sequential thinking such as the derivation of a proof).

This paper elaborates on an application that requires both systems to be used, namely Visual Question Answering. System 1 will be required in order to gain an understanding of the image under investigation, with in particular their shapes and colors. System 2, on the other hand, will use this extracted information for deriving certain properties of objects¹, or even for capturing the relations² between the objects.

II. ORGANIZATION OF THE PAPER

This paper will first provide some necessary background in Section III on certain types of datasets that are typically used for Visual Question Answering purposes. Then, Section IV will dive deeper in the different components constituting the overall system. Section V, furthermore, provides some results of conducted experiments. A main focus in this paper is to outline the advantages of using a Neuro-Symbolic AI approach (offered by the DeepProbLog framework), instead of a purely neural network based approach. These specific experiments are listed in Subsection A. Finally, the main findings of this paper are repeated in Section VI.

III. LITERATURE SURVEY

The application focuses on Visual Question Answering (VQA), for which huge datasets are present, along with very sophisticated methods. The best known dataset for VQA is CLEVR [2], which contains 100k images with one million questions. An example image is given in Figure 1, while example questions are:

- Are there an equal number of large things and metal spheres?

- What size is the cylinder that is left of the brown metal thing that is left of the big sphere?
- How many objects are either small cylinders or metal things?

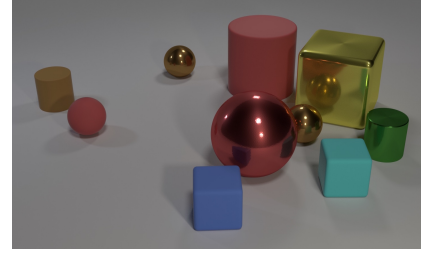


Fig. 1. A sample image from the CLEVR dataset [2]

Clearly, both system 1 and system 2 are actively used when answering these questions. One could wonder if neural networks alone could answer these questions without having an explicit system 2 encoding (i.e. the rule based knowledge base). Intuitively, it makes sense that if certain facts of the world are known³, learning can proceed much more quickly⁴. Seen from an optimization viewpoint, errors made during prediction in this setup can be targeted exactly, which makes the optimization process more targeted as well, and hence more efficient. Finally, this paper also provides evidence for these statements, since in Subsection A the comparison between a VQA implementation with DeepProbLog is made with a purely CNN based approach.

This paper is inspired on the CLEVR dataset, but uses however a much more simplified version. In essence, it is almost like the Sort-Of-CLEVR dataset [3]. This Sort-Of-CLEVR dataset contains images such as in Figure 2, while asking questions such as:

- Non-relational questions: the shape, horizontal or vertical location of an object.
- Relational questions: shape of the closest/furthest object to the object under investigation, or the number of objects with the same shape.

As illustrated earlier, both system 1 and system 2 are required for these types of VQA's.

Finally, since this application uses DeepProbLog, quite some time was spent in digesting the DeepProbLog paper [1], along with understanding the examples provided in the code repository [4].

¹For example, finding the shape of the green object, or deriving if it is located on the left hand side of the image.

²Here, one could think of counting the number of circles in the image.

³Facts can be encoded, e.g. counting the number of spheres is simply a matter of detecting all the spheres in the image, after which a mathematical summation is a statement in the knowledge base.

⁴Not to say that learning might even be impossible if a lot of background knowledge is required.

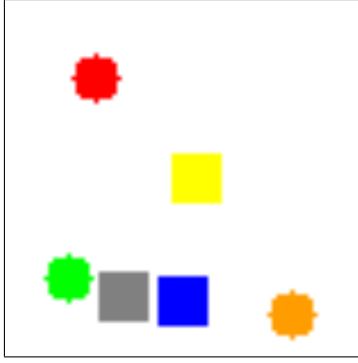


Fig. 2. A sample image from the CLEVR dataset [3]

IV. APPROACH

The implementation process involved three main parts:

1. Generation of data.
2. Linking the data and controlling the training process in pure Python code.
3. Creation of the logical part with DeepProbLog statements.

A. Generation of data

As mentioned in Section III, the data used in this application is based on the Sort-Of-CLEVR dataset, with one extra simplification. Given that the logical part will have to decide whether an object is for example located on the left side of an image, the neural network will have to convey positional information to the logical part. Hence, each discrete position will have to be encoded by a possible outcome of the neural network. Therefore, objects may only be located at certain places in a grid. In this paper, results on a grid of 2x2 and 6x6 are discussed.

The data generator that was used for the creation of the Sort-Of-CLEVR dataset has been modified in order to place objects in the mentioned grid positions [3]. An example of a generated image is given in Figure 3, where the difference with Figure 2 is the grid-layout.

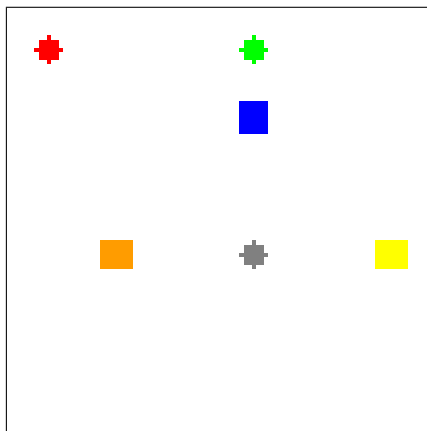


Fig. 3. A sample image from the dataset that has been used for this application

Each specified color will have an object located somewhere in the grid, of which the shape can be a square or a circle.

These images are accompanied with a question about a random object, which can be one of the following:

- Non-binary - What is the shape of this object ⁵?
- Non-binary - Is this object located on the left side of the image?
- Non-binary - Is this object located on the bottom side of the image?
- Binary - How many objects have the same shape as this object?

These questions are encoded in a one-hot encoding, after which they are stored in a CSV file, along with the expected answers. To make the training process more efficient, a training and test dataset has been generated beforehand.

B. Controlling the training process

The overall training process is controlled via the Python API of DeepProbLog, along with general PyTorch implementations of the CNN's. First of all, CNN's are defined with PyTorch. A relatively simple network is used, where the input is given as a square RGB image of 100 pixels wide, which is transformed by the CNN into 72 output features for the 6x6 grid⁶. Each color that is present in the image has its accompanied CNN network, hence the 72 output features encode the possible positions of the object with that color, along with their shape, which can be either square or circular ($6 \cdot 6 \cdot 2 = 72$).

The final thing (besides the logical rule encodings) required before commencing the training process, are the data loaders. The most challenging part here is the transformation from the generated data to specific query mappings and their outcome.

C. Logical rule encodings

Once the CNN belonging to a specific color has determined the position and the shape of that object, logical rules can deduce whether this object is located on the left hand side of the image, on the bottom side, and how many objects have the same shape. The logical rule program has been listed in Appendix A.

V. EXPERIMENTS

A. COMPARISONS WITH PURE SYSTEM 1 APPROACHES

The network based on pure neural predicates is able to recognize the questions quickly, however, seems to experience difficulties when having to decide for the correct answer. This can clearly be seen in the confusion matrix (Figure 4).

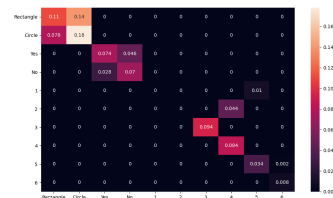


Fig. 4. The confusion matrix of the 6x6 dataset on a purely Neural Network based model

⁵I.e. the shape will be either a square or a circle.

⁶For the 2x2 grid example, 8 output features are required.

It is clear that the network trained with DeepProbLog is able to learn way faster then the purely NN based approach, if measured in the number of iterations ⁷. However, the differences in training time are much less significant, due to the high computational cost associated with the arithmetic parts of the circuit (i.e. the semirings).

VI. CONCLUSIONS

VII. APPENDIX

A. Logical rule encodings

```
width(6).
size(N) :-
    width(Width),
    N is Width * Width.
states(N) :-
    size(Size),
    N is Size * 2.

nn(cnn_red,[X],Y_red, [0, 1, 2, 3, 4, 5, 6, 7,
    ↪ 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18,
    ↪ 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29,
    ↪ 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40,
    ↪ 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51,
    ↪ 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62,
    ↪ 63, 64, 65, 66, 67, 68, 69, 70, 71]) ::
    ↪ detect_state(red, X, Y_red).
nn(cnn_green,[X], Y_green, [0, 1, 2, 3, 4, 5,
    ↪ 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,
    ↪ 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28,
    ↪ 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39,
    ↪ 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50,
    ↪ 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61,
    ↪ 62, 63, 64, 65, 66, 67, 68, 69, 70, 71]) ::
    ↪ detect_state(green, X, Y_green).
nn(cnn_blue,[X],Y_blue, [0, 1, 2, 3, 4, 5, 6,
    ↪ 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,
    ↪ 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28,
    ↪ 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39,
    ↪ 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50,
    ↪ 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61,
    ↪ 62, 63, 64, 65, 66, 67, 68, 69, 70, 71]) ::
    ↪ detect_state(blue, X, Y_blue).
nn(cnn_orange,[X], Y_orange, [0, 1, 2, 3, 4, 5,
    ↪ 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,
    ↪ 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28,
    ↪ 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39,
    ↪ 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50,
    ↪ 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61,
    ↪ 62, 63, 64, 65, 66, 67, 68, 69, 70, 71]) ::
    ↪ detect_state(orange, X, Y_orange).
```

```
nn(cnn_grey,[X], Y_grey, [0, 1, 2, 3, 4, 5, 6,
    ↪ 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,
    ↪ 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28,
    ↪ 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39,
    ↪ 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50,
    ↪ 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61,
    ↪ 62, 63, 64, 65, 66, 67, 68, 69, 70, 71]) ::
    ↪ detect_state(grey, X, Y_grey).
nn(cnn_yellow,[X], Y_yellow, [0, 1, 2, 3, 4, 5,
    ↪ 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,
    ↪ 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28,
    ↪ 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39,
    ↪ 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50,
    ↪ 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61,
    ↪ 62, 63, 64, 65, 66, 67, 68, 69, 70, 71]) ::
    ↪ detect_state(yellow, X, Y_yellow).
```

```
detect_same_states(Shape, Img, Count) :-
    check_state(red, Shape, Img, Return_red),
    check_state(green, Shape, Img,
        ↪ Return_green),
    check_state(blue, Shape, Img, Return_blue),
    check_state(orange, Shape, Img,
        ↪ Return_orange),
    check_state(grey, Shape, Img, Return_grey),
    check_state(yellow, Shape, Img,
        ↪ Return_yellow),
    Count is Return_red + Return_green +
        ↪ Return_blue + Return_orange +
        ↪ Return_grey + Return_yellow.
```

```
check_state(Color, Shape, Img, Return) :-
    detect_state(Color, Img, Y_color),
    to_shape(Y_color, Shape_color),
    Shape == Shape_color,
    Return = 1.
check_state(Color, Shape, Img, Return) :-
    detect_state(Color, Img, Y_color),
    to_shape(Y_color, Shape_color),
    Shape \== Shape_color,
    Return = 0.
```

% NON-BINARY QUESTIONS

% Question about the shape

```
shape(Img, Color, Rectangle) :-
    shape_is_rectangle(Img, Color),
    Rectangle = 1.
shape(Img, Color, Rectangle) :-
    \+ shape_is_rectangle(Img, Color),
    Rectangle = 0.
```

```
shape_is_rectangle(Img, Color) :-
    detect_state(Color, Img, Y_color),
    size(Size),
    Y_color < Size.
```

% Question about the horizontal side

```
horizontal_side(Img, Color, Left) :-
    position_is_left(Img, Color),
```

⁷By the number of iterations, the number of training steps on a batch is meant.

```

Left = 1.
horizontal_side(Img, Color, Left) :-
    \+ position_is_left(Img, Color),
    Left = 0.

position_is_left(Img, Color) :-
    detect_state(Color, Img, Y_color),
    position(Y_color, Pos),
    left_side(Pos).

% Question about the vertical side
vertical_side(Img, Color, Bottom) :-
    position_is_bottom(Img, Color),
    Bottom = 1.
vertical_side(Img, Color, Bottom) :-
    \+ position_is_bottom(Img, Color),
    Bottom = 0.

position_is_bottom(Img, Color) :-
    detect_state(Color, Img, Y_color),
    position(Y_color, Pos),
    bottom_side(Pos).

% BINARY QUESTION
% Question about the number of same shapes like
% → that one
number_of_shapes(Img, Color, Count) :-
    detect_state(Color, Img, Y_red),
    to_shape(Y_red, Shape),
    detect_same_states(Shape, Img, Count).

to_shape(Y_color, Shape) :-
    size(Size),
    Y_color < Size,
    Shape = 0. % Shape = 0 --> Rectangle
to_shape(Y_color, Shape) :-
    size(Size),
    Y_color >= Size,
    Shape = 1. % Shape = 1 --> Circle

position(Out, Pos) :-
    size(Size),
    Pos is Out mod Size.

left_side(Out) :-
    to_coor(Out, X, Y),
    width(Width),
    X < (Width // 2).

bottom_side(Out) :-
    to_coor(Out, X, Y),
    width(Width),
    Y >= (Width // 2). % Be careful. Y axis
    % → runs down!

to_coor(Out, X, Y) :-
    width(Width),

```

```

X is Out mod Width,
Y is Out // Width.

```

REFERENCES

- [1] R. Manhaeve, A. Kimmig, S. Dumančić, T. Demeester, L. De Raedt, “Deepproblog: Neural probabilistic logic programming”, in *Advances in Neural Information Processing Systems*, vol. 2018-Decem, pp. 3749–3759, jul 2018, doi:10.48550/arxiv.1907.08194, URL: <https://arxiv.org/abs/1907.08194v2>, 1805.10872.
- [2] J. Johnson, L. Fei-Fei, B. Hariharan, C. L. Zitnick, L. Van Der Maaten, R. Girshick, “CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning”, *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, vol. 2017-Janua, pp. 1988–1997, 2017, doi: 10.1109/CVPR.2017.215, 1612.06890.
- [3] K. Heecheol, “kimhc6028/relational-networks: Pytorch implementation of "A simple neural network module for relational reasoning" (Relational Networks)”, URL: <https://github.com/kimhc6028/relational-networks>.
- [4] R. Manhaeve, “ML-KULeuven/deepproblog: DeepProbLog is an extension of ProbLog that integrates Probabilistic Logic Programming with deep learning by introducing the neural predicate.”, URL: <https://github.com/ML-KULeuven/deepproblog>.