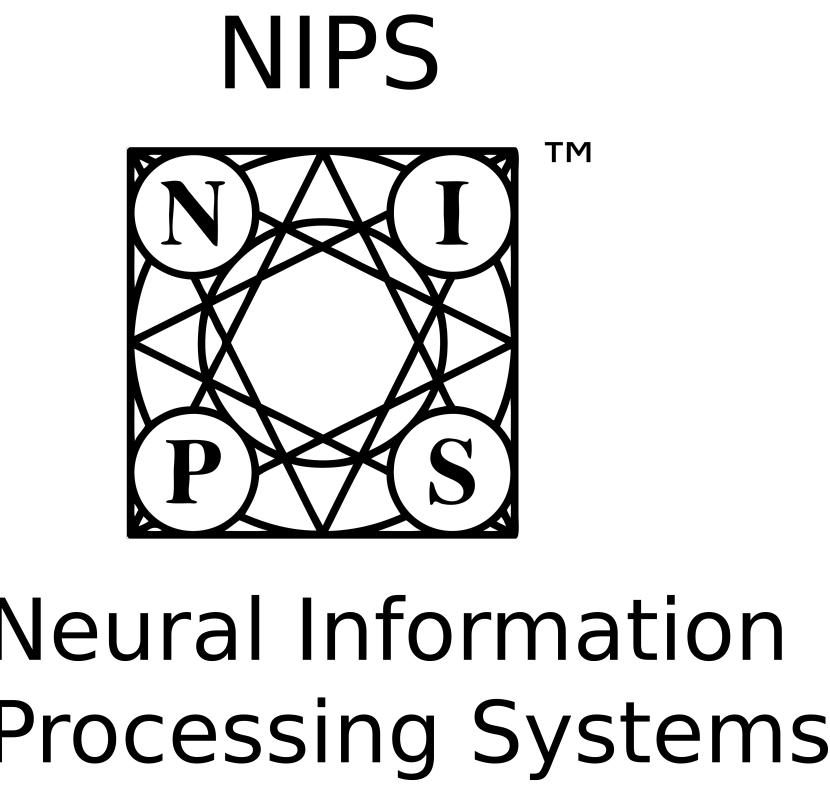


# Learning Conditioned Graph Structures for Interpretable Visual Question Answering

Will Norcliffe-Brown, Stathis Vafeias and Sarah Parisot



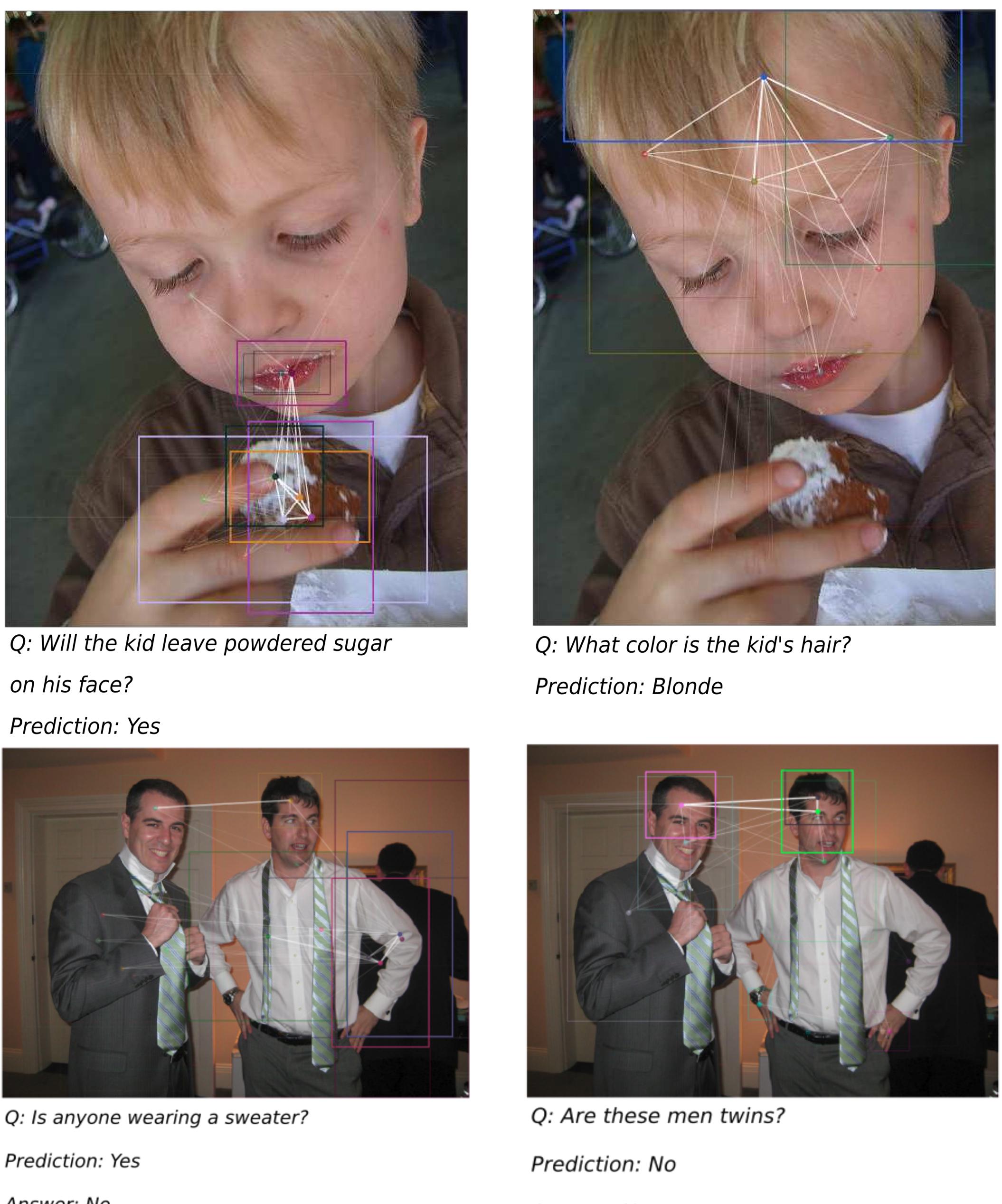
## Introduction

### Context

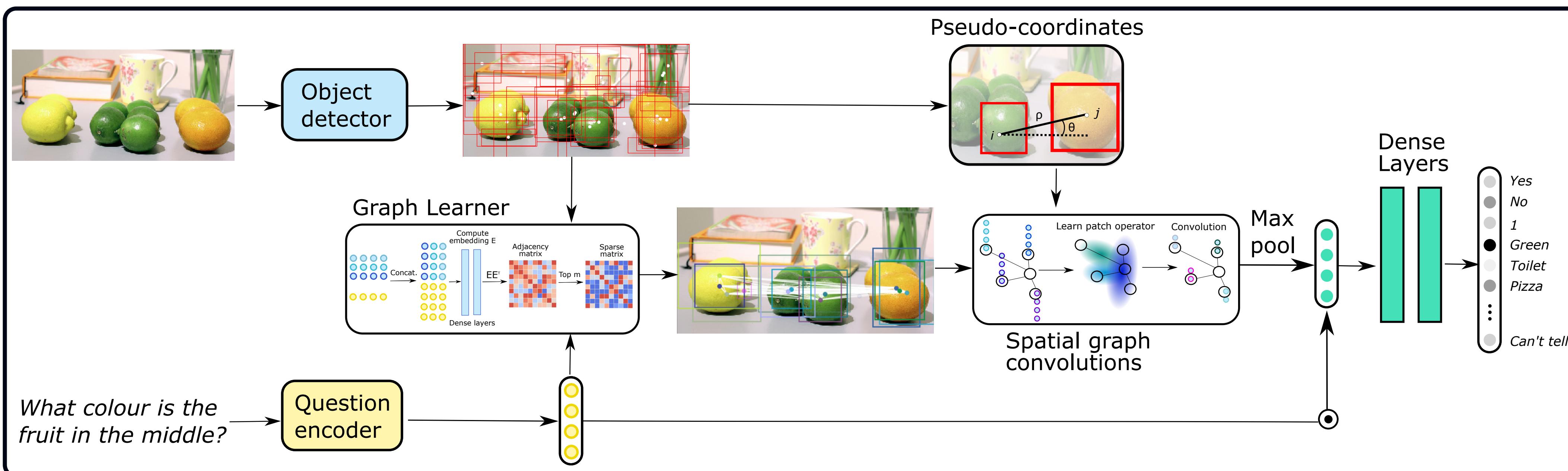
- VQA is a challenging topic requiring a combination of natural language and image processing techniques
- An effective VQA system requires **modelling the interactions between objects** in an image
- Previous approaches have either modelled objects as independent or use engineered graph structures

### Contributions

- A novel graph-based approach for VQA
- Learn a graph representation of the image, **conditioned on the question at hand**, that identify the most relevant objects and relationships for answering a question
- Learned graph structures provide **interpretability**
- Promising results with **minimal architecture tuning**

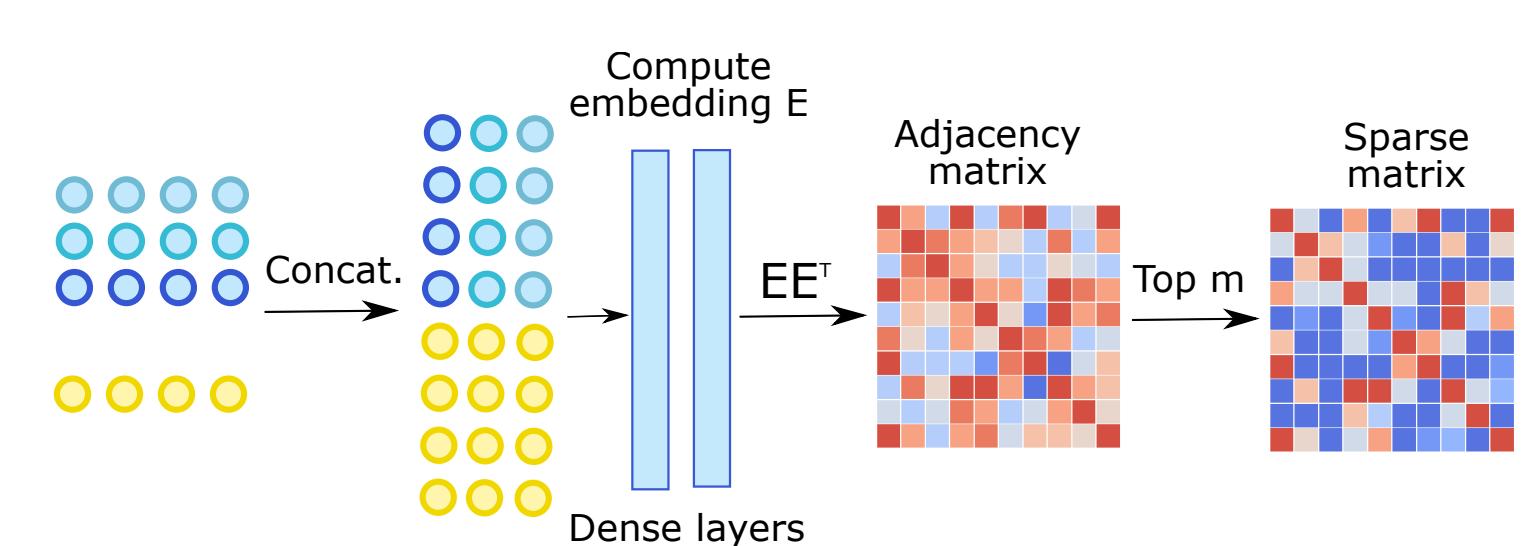


## Method Overview



## Graph Learner Module

- Define an undirected graph where nodes are object bounding boxes and **edges are learned and conditioned on the question**
- Compute a joint embedding from concatenated visual and question features:  $\mathbf{e}_n = F([\mathbf{v}_n \parallel \mathbf{q}])$
- Compute an adjacency matrix as  $A_{i,j} = \mathbf{e}_i^T \mathbf{e}_j$
- **Sparsity** imposed by defining each node's neighbourhood system  $\mathcal{N}(i)$  as the nodes with the **top m** strongest edge weights.



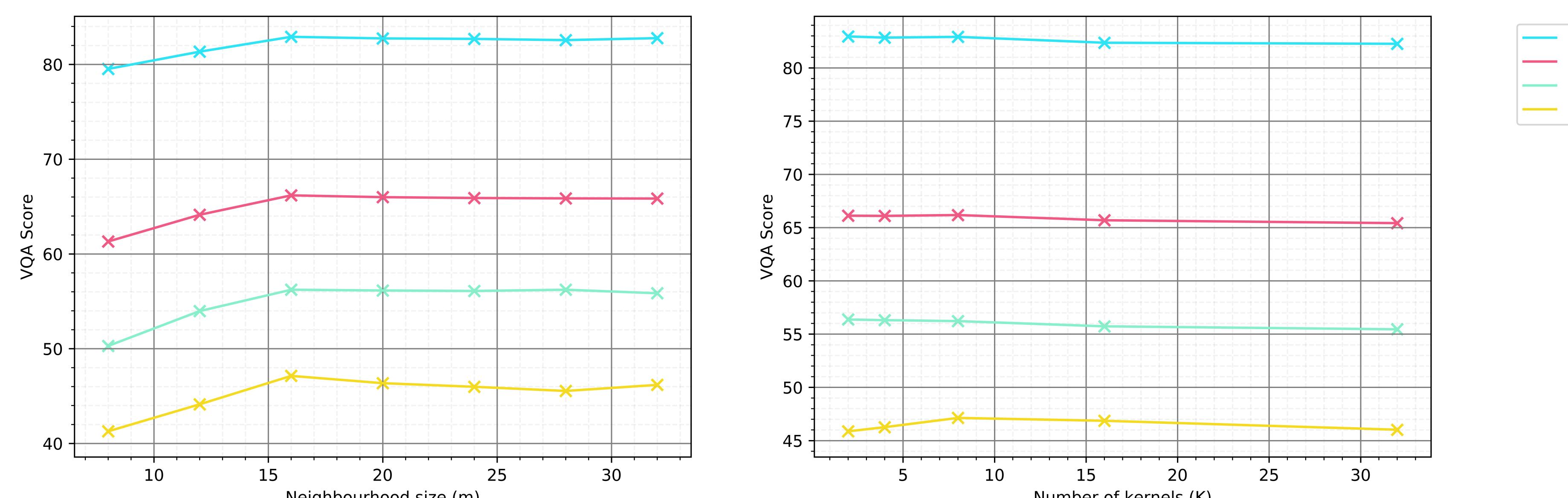
## Spatial Graph Convolutions

- Learn object representations informed by graph neighbours
- A patch operator describes the influence of each node in  $\mathcal{N}(i)$  as a mixture of Gaussian kernels[2]. Compute updated visual features as

$$\mathbf{f}_k(i) = \sum_{j \in \mathcal{N}(i)} \underbrace{w_k(\mathbf{u}(i, j))}_{\substack{\text{Gaussian kernel with} \\ \text{learnable weights}}} \underbrace{\mathbf{v}_j}_{\substack{\text{Object/bounding box} \\ \text{visual features}}} \underbrace{\alpha_{ij}}_{\substack{\text{Graph edge weights}}} \underbrace{\mathbf{v}_j}_{\substack{\text{Bounding box} \\ \text{Pseudo-} \\ \text{coordinates}}} \underbrace{\alpha_{ij}}_{\substack{\text{Graph edge weights}}}$$

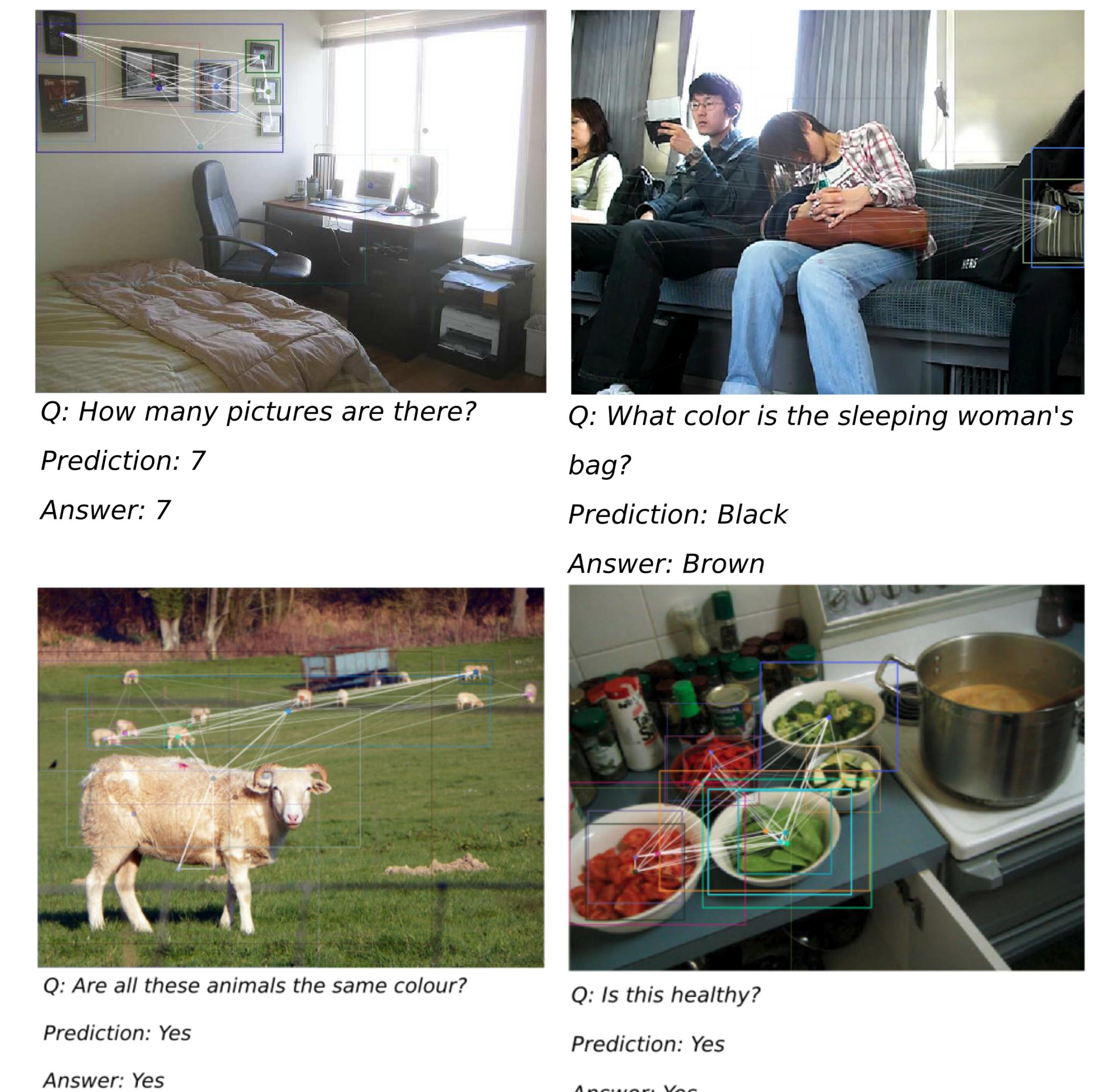
- Output at vertex  $i$ :  $\mathbf{h}_i = \parallel_{k=1}^K \mathbf{G}_k \mathbf{f}_k(i)$  where  $\mathbf{G}_k$  is a matrix of learnable weights.

## Parameter Analysis



## Results

Answer type	All	Y/N	Num.	Other
ReasonNet	64.61	78.86	41.98	57.39
Bottom-Up [1]	65.67	82.20	43.90	56.26
Counting module [3]	68.41	83.56	51.39	59.11
kNN graph	61.00	79.35	41.63	49.70
Attention	61.90	79.87	42.48	50.95
<b>Ours</b>	66.18	82.91	47.13	56.22



## References

- [1] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. Bottom-up and top-down attention for image captioning and visual question answering (2017).
- [2] F. Monti, D. Boscaini, J. Masci, E. Rodolà, J. Svoboda, and M. M. Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns (2017).
- [3] Y. Zhang, J. S. Hare, and A. Prügel-Bennett. Learning to count objects in natural images for visual question (2018).

## Links

### Code

[github.com/aimbrain/vqa-project](https://github.com/aimbrain/vqa-project)

### Contact

- [will.norcliffe@aimbrain.com](mailto:will.norcliffe@aimbrain.com)
- [stathis@aimbrain.com](mailto:stathis@aimbrain.com)
- [sarah.parisot@huawei.com](mailto:sarah.parisot@huawei.com)