# Self-Critical Reasoning for Robust Visual Question Answering
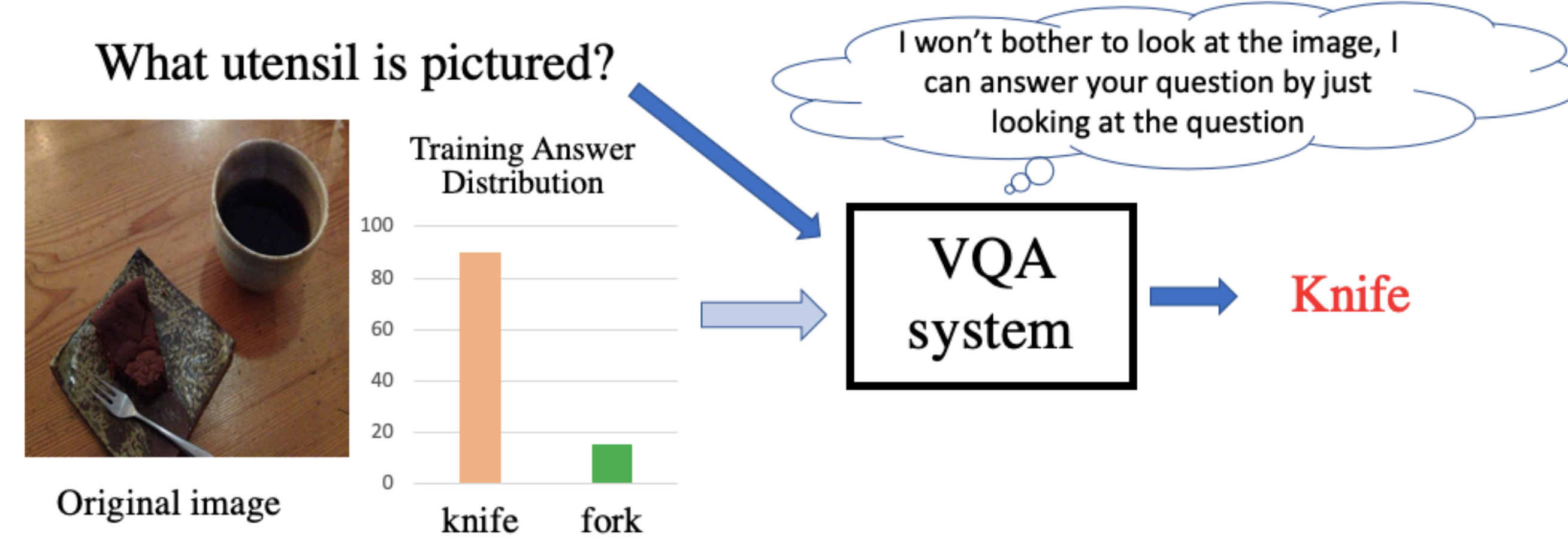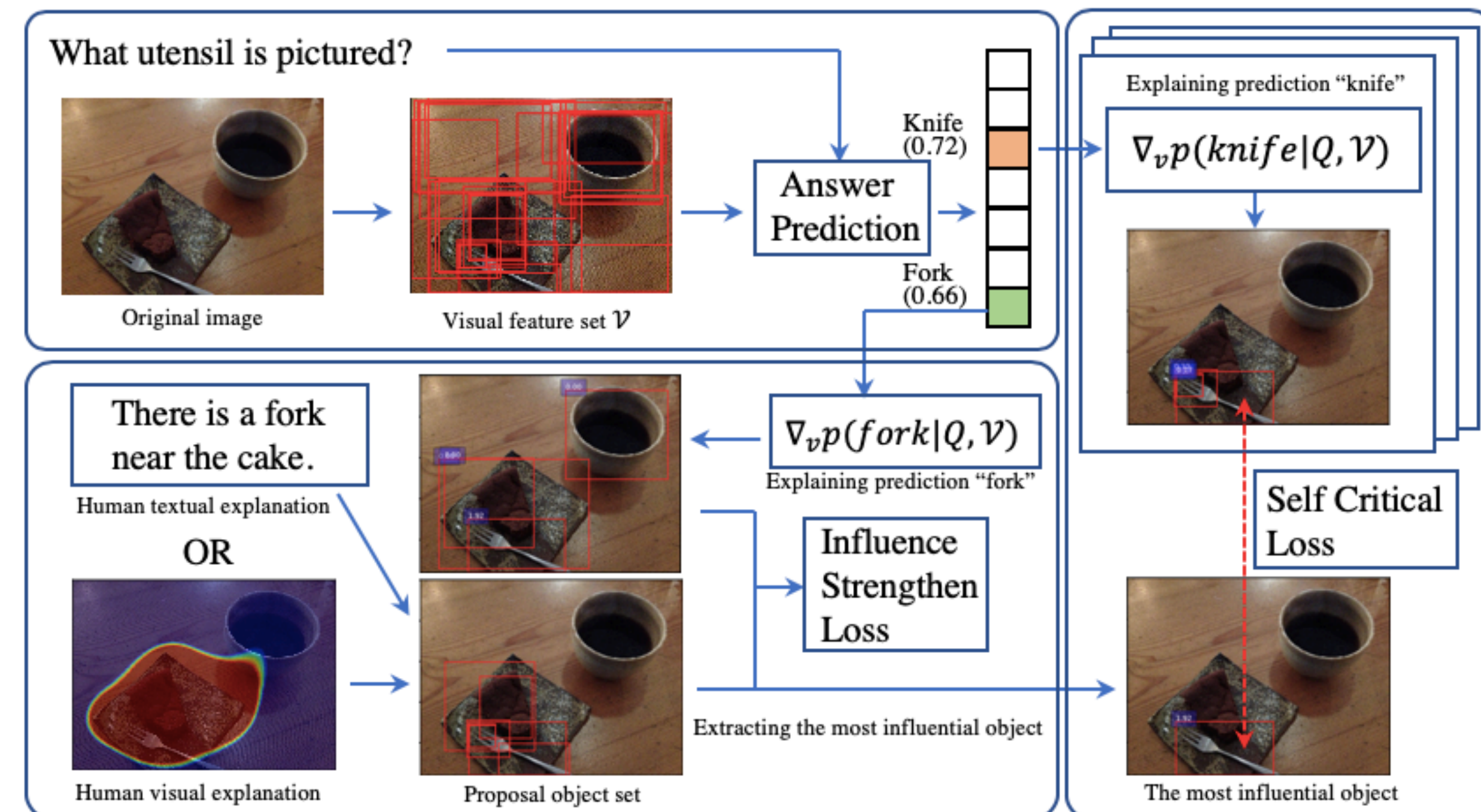## Jialin Wu, and Raymond J. Mooney The University of Texas at Austin

## Introduction

- ➤ Common VQA systems tend to only capture superficial statistical correlations between QA pairs, especially when training and test set are under different distribution.
- ➤ VQA systems should focus on the objects that human would focus.
- ➤ We propose two constraints for VQA systems that help the right objects contribute more to the right answers than to the wrong answers.



## Model Overview



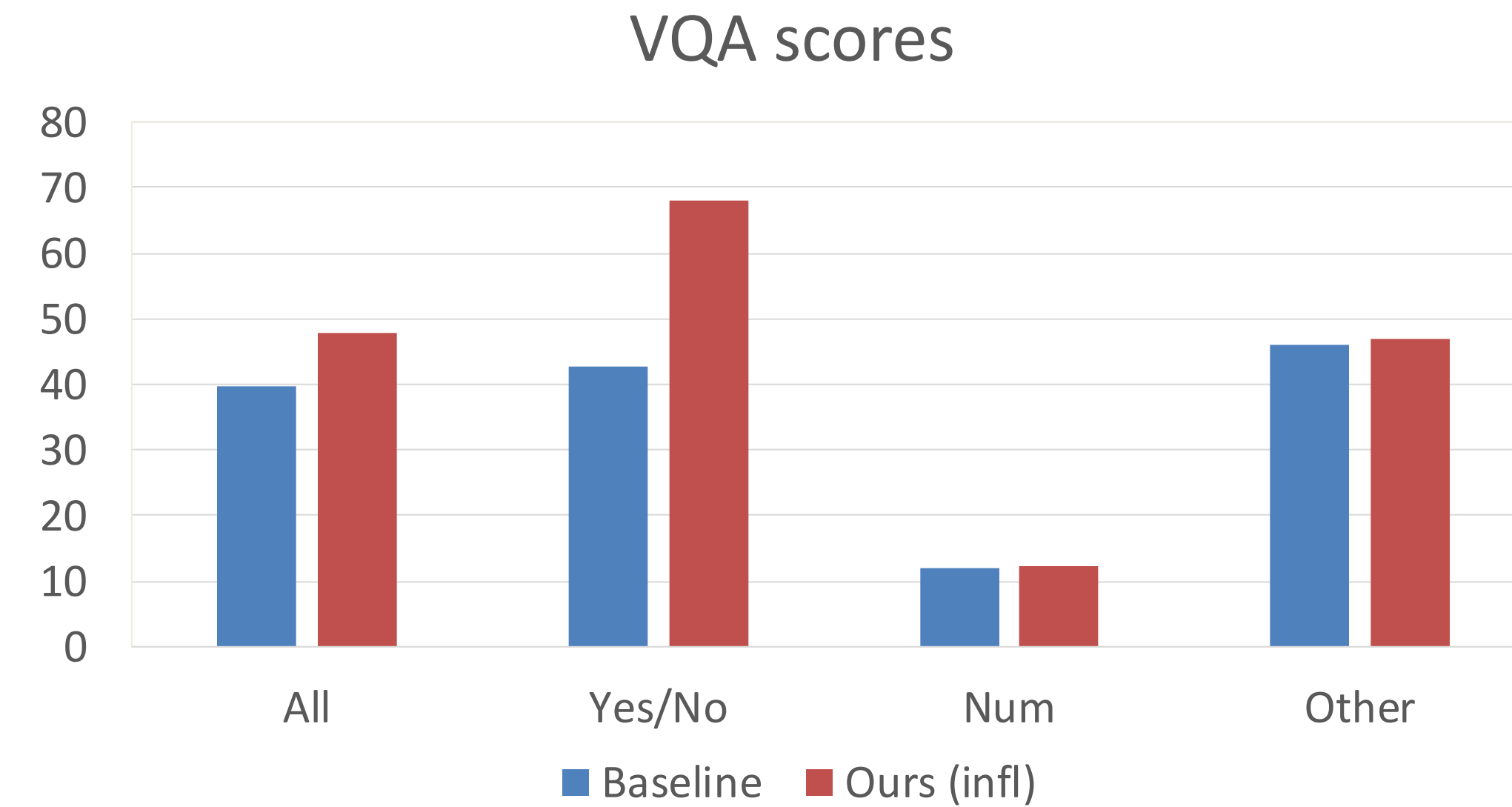## Recognizing and Strengthening Influential Objects

- ➤ Extracting a set of influential objects ($\mathcal{I}$) that humans would focus on.
- ➤ Enforcing the gradients ($\mathcal{S}(a, \mathbf{v}_i)$) from the correct answer to have the biggest value in at least one of the extracted objects.

$$\mathcal{S}(a, \mathbf{v}_i) := \left(\nabla_{\mathbf{v}_i} P(a|V, q)\right)^T \mathbf{1}$$

$$\mathcal{SV}(a, \mathbf{v}_i, \mathbf{v}_j) = \max\left(\mathcal{S}(a, \mathbf{v}_j) - \mathcal{S}(a, \mathbf{v}_i), 0\right)$$

$$\mathcal{L}_{infl} = \min_{\mathbf{v}_i \in \mathcal{I}} \left( \sum_{\mathbf{v}_j \in \mathcal{V} \backslash \mathcal{I}} \mathcal{SV}(a_{gt}, \mathbf{v}_i, \mathbf{v}_j) \right)$$

- ➤ Comparing to Up-Down VQA system.



VQA scores

■ Baseline ■ Ours (infl)

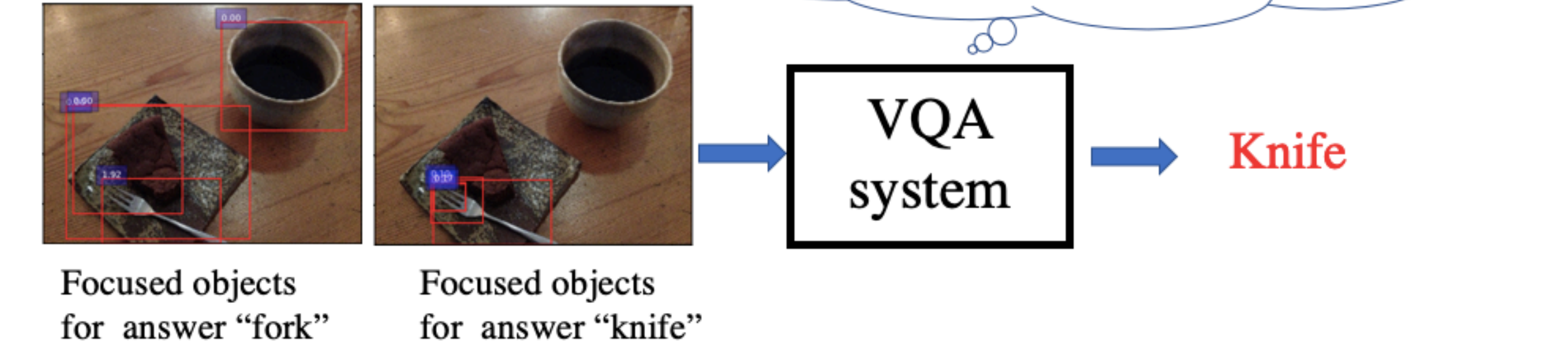## Criticizing Incorrect Dominant Answers

- ➤ Although VQA systems focus on the right object for the right answer, but the object could contribute more to the wrong answers.
- ➤ Finding the most influential object ($\mathbf{v}^*$) using gradient-based method.

$$\mathbf{v}^* = \arg\min_{\mathbf{v}_i \in \mathcal{I}} \left( \sum_{\mathbf{v}_j \in \mathcal{V} \backslash \mathcal{I}} \mathcal{SV}(a_{gt}, \mathbf{v}_i, \mathbf{v}_j) \right)$$
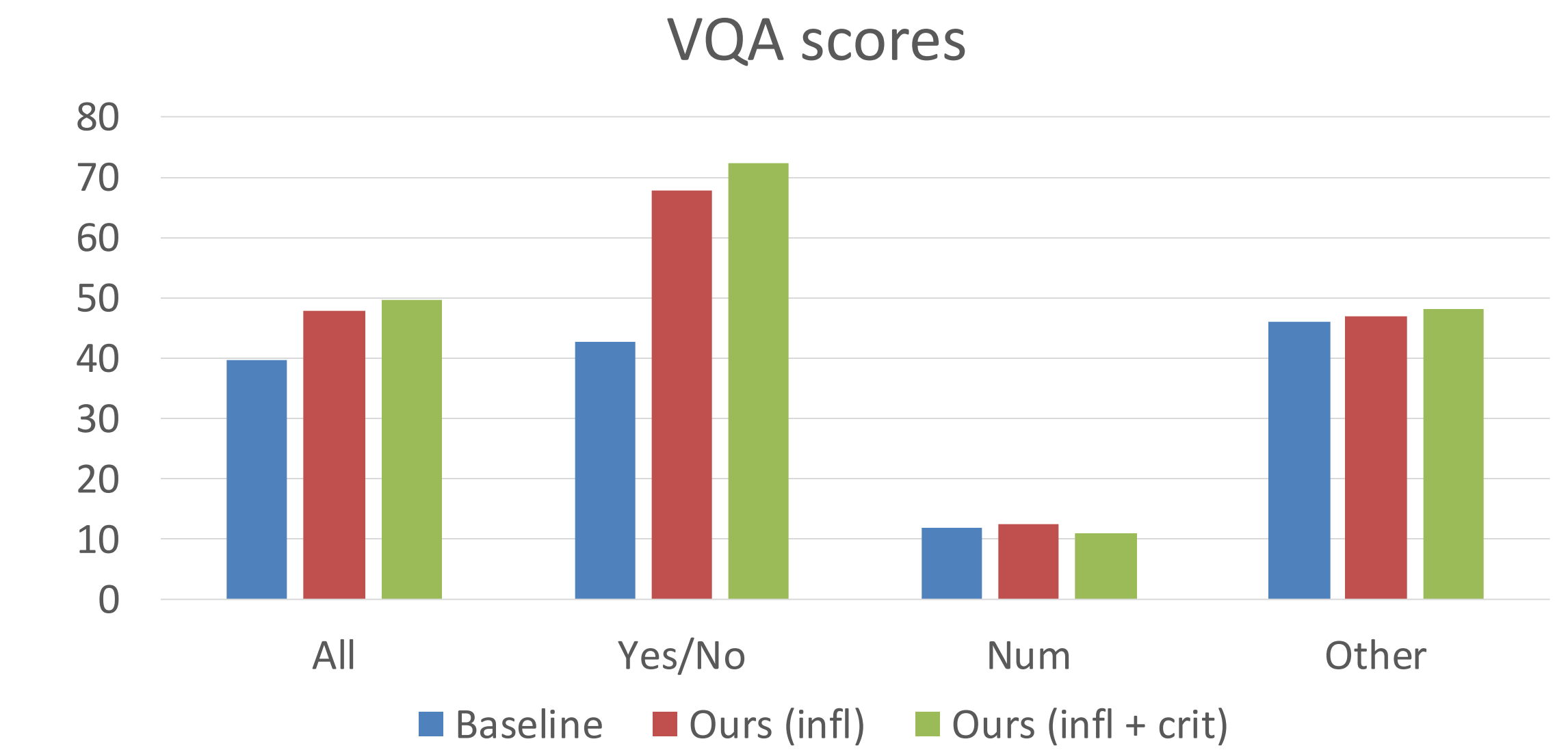
- ➤ Enforcing the object to contribute more to the correct answer.

$$\mathcal{L}_{crit} = \sum_{a \in \mathcal{B}} w(a)(\mathcal{S}(a, \mathbf{v}^*) - \mathcal{S}(a_{gt}, \mathbf{v}^*))$$

What utensil is pictured?



Focused objects for answer "fork"    Focused objects for answer "knife"

- ➤ Comparing to Up-Down VQA system.



VQA scores

■ Baseline ■ Ours (infl) ■ Ours (infl + crit)

## Conclusion

- ➤ VQA systems should be able to focus on the right set of objects as human do to predict the right answer
- ➤ It is also necessary to prevent the systems from over sensitive to the most common answers.