

# Perception Matters: Detecting Perception Failures of VQA Models Using Metamorphic Testing

## Supplementary Materials

Yuan Yuan<sup>1</sup> Shuai Wang<sup>1</sup> Mingyue Jiang<sup>2</sup> Tsong Yueh Chen<sup>3</sup>  
<sup>1</sup>HKUST, <sup>2</sup>Zhejiang Sci-Tech University, <sup>3</sup>Swinburne University of Technology  
 {yyuanaq, shuaiw}@cse.ust.hk, mjiang@stu.edu.cn, tychen@swin.edu.au

### A. Errors found Using Question-Oriented MRs

This section provides some errors that we found using question-oriented MRs.



Q: # of pillows and hands? A: 2  
 Q<sub>1</sub>: # of pillows ? A<sub>1</sub>: 1  
 Q<sub>2</sub>: # of hands ? A<sub>2</sub>: 2  
 -----  
 Q: # of girls and boards? A: 2  
 Q<sub>1</sub>: # of girls? A<sub>1</sub>: 1  
 Q<sub>2</sub>: # of boards? A<sub>2</sub>: 0



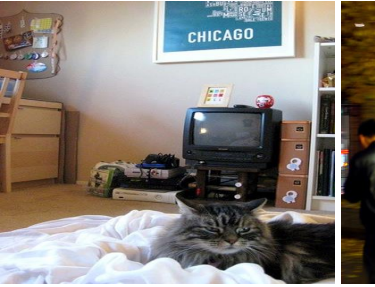
Q: # of lines and roads? A: 2  
 Q<sub>1</sub>: # of lines? A<sub>1</sub>: 3  
 Q<sub>2</sub>: # of roads? A<sub>2</sub>: 1  
 -----  
 Q: # of bikes and cyclists? A: 4  
 Q<sub>1</sub>: # of bikes? A<sub>1</sub>: 3  
 Q<sub>2</sub>: # of cyclists? A<sub>2</sub>: 4



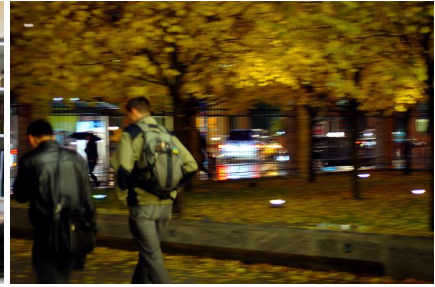
Q: # of trees and frisbees? A: 2  
 Q<sub>1</sub>: # of trees? A<sub>1</sub>: 5  
 Q<sub>2</sub>: # of frisbees? A<sub>2</sub>: 1  
 -----  
 Q: # of horses and trees? A: 2  
 Q<sub>1</sub>: # of horses? A<sub>1</sub>: 0  
 Q<sub>2</sub>: # of trees? A<sub>2</sub>: 5



Q: # of towels and windows? A: 2  
 Q<sub>1</sub>: # of towels? A<sub>1</sub>: 3  
 Q<sub>2</sub>: # of windows? A<sub>2</sub>: 2  
 -----  
 Q: # of lights and towels? A: 5  
 Q<sub>1</sub>: # of lights? A<sub>1</sub>: 3  
 Q<sub>2</sub>: # of towels? A<sub>2</sub>: 3



Q: # of televisions and floors? A: 1  
 Q<sub>1</sub>: # of televisions ? A<sub>1</sub>: 1  
 Q<sub>2</sub>: # of floors? A<sub>2</sub>: 1  
 -----  
 Q: # of blankets and cats? A: 1  
 Q<sub>1</sub>: # of blankets? A<sub>1</sub>: 1  
 Q<sub>2</sub>: # of cats? A<sub>2</sub>: 2



Q: # of persons and lights? A: 2  
 Q<sub>1</sub>: # of persons ? A<sub>1</sub>: 2  
 Q<sub>2</sub>: # of lights? A<sub>2</sub>: 2  
 -----  
 Q: # of black jackets and green jackets? A: 2  
 Q<sub>1</sub>: # of black jackets ? A<sub>1</sub>: 2  
 Q<sub>2</sub>: # of green jackets ? A<sub>2</sub>: 2

Figure 1: Errors found on Oscar<sub>large</sub><sup>+</sup> using object-/property-oriented partitioning MRs. Due to the limited space, we replace phrase “what is the total number” with “#”.



Q: Is there any bed? A: Yes  
 ¬Q: Is there no bed? A: Yes  
 -----  
 Q: Is there any yellow bed? A: No  
 ¬Q: Is there no yellow bed? A: No



Q: Is there any purple dog? A: No  
 ¬Q: Is there no purple dog? A: No  
 -----  
 Q: Is there any brick building? A: No  
 ¬Q: Is there no brick building? A: No



Q: Is there any microwave? A: Yes  
 ¬Q: Is there no microwave? A: Yes  
 -----  
 Q: Is there any golden chair? A: No  
 ¬Q: Is there no golden chair? A: No



Q: Is there any black spoon? A: Yes  
 ¬Q: Is there no black spoon? A: Yes  
 -----  
 Q: Is there any white bowl? A: Yes  
 ¬Q: Is there no white bowl? A: Yes



Q: Is there any sink? A: Yes  
 ¬Q: Is there no sink? A: Yes  
 -----  
 Q: Is there any drain? A: Yes  
 ¬Q: Is there no drain? A: Yes



Q: Is there any car? A: Yes  
 ¬Q: Is there no car? A: Yes  
 -----  
 Q: Is there any yellow statue? A: No  
 ¬Q: Is there no yellow statue? A: No

Figure 2: Errors found on Oscar<sub>large</sub><sup>+</sup> using object-/property-oriented reversion MRs. Due to the limited space, we replace phrase “*what is the total number*” with “#”.





Q: # of horses and signs? A: 1  
 Q<sub>3</sub>: # of signs and horses? A<sub>3</sub>: 2  
 -----  
 Q: # of men and trees? A: 2  
 Q<sub>3</sub>: # of trees and men? A<sub>3</sub>: 10



Q: # of women and giraffes? A: 2  
 Q<sub>3</sub>: # of giraffes and women? A<sub>3</sub>: 4  
 -----  
 Q: # of trees and ears? A: 10  
 Q<sub>3</sub>: # of ears and trees? A<sub>3</sub>: 2



Q: # of umbrellas and chairs? A: 2  
 Q<sub>3</sub>: # of chairs and umbrellas? A<sub>3</sub>: 3  
 -----  
 Q: # of kites and persons? A: 2  
 Q<sub>3</sub>: # of persons and kites? A<sub>3</sub>: 1



Q: # of watches and bottles? A: 2  
 Q<sub>3</sub>: # of bottles and watches? A<sub>3</sub>: 1  
 -----  
 Q: # of tables and men? A: 1  
 Q<sub>3</sub>: # of men and tables? A<sub>3</sub>: 2



Q: # of dogs and benches? A: 1  
 Q<sub>3</sub>: # of benches and dogs? A<sub>3</sub>: 4  
 -----  
 Q: # of sidewalks and teeth? A: 2  
 Q<sub>3</sub>: # of teeth and sidewalks? A<sub>3</sub>: 0



Q: # of players and pants? A: 2  
 Q<sub>3</sub>: # of pants and players? A<sub>3</sub>: 4  
 -----  
 Q: # of men and hats? A: 2  
 Q<sub>3</sub>: # of hats and men? A<sub>3</sub>: 3

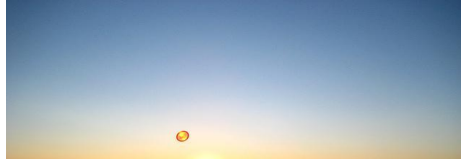
Figure 3: Errors found on Oscar<sup>+</sup><sub>large</sub> using object-/property-oriented reordering MRs. Due to the limited space, we replace phrase “*what is the total number*” with “#”.

## B. Errors found Using Image-Oriented MRs

This section provides some errors that we found using image-oriented MRs.



Q: # of ovens? A: 1  
A<sub>1</sub>: 1    A<sub>2</sub>: 1  
-----  
Q: # of apples? A: 2  
A<sub>1</sub>: 0    A<sub>2</sub>: 0



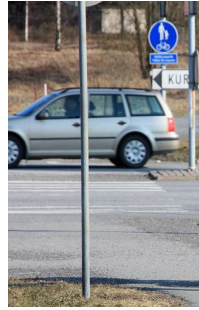
Q: # of persons? A: 1  
A<sub>1</sub>: 0    A<sub>2</sub>: 2  
-----  
Q: # of walking persons and standing persons? A: 1  
A<sub>1</sub>: 0    A<sub>2</sub>: 2



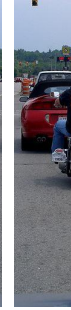
Q: # of persons? A: 1  
A<sub>1</sub>: 0    A<sub>2</sub>: 0    A<sub>3</sub>: 2    A<sub>4</sub>: 4  
-----  
Q: # of kites? A: 4  
A<sub>1</sub>: 6    A<sub>2</sub>: 6    A<sub>3</sub>: 0    A<sub>4</sub>: 0



Q: # of standing persons? A: 1  
A<sub>1</sub>: 0    A<sub>2</sub>: 2    A<sub>3</sub>: 1  
-----  
Q: # of sinks? A: 2  
A<sub>1</sub>: 0    A<sub>2</sub>: 2    A<sub>3</sub>: 2



Q: # of backpacks? A: 0  
A<sub>1</sub>: 0    A<sub>2</sub>: 1  
-----  
Q: Is there any rolling person? A: No  
A<sub>1</sub>: No    A<sub>2</sub>: Yes



Q: # of persons? A: 2  
A<sub>1</sub>: 2    A<sub>2</sub>: 0    A<sub>3</sub>: 0    A<sub>4</sub>: 2    A<sub>5</sub>: 0  
-----  
Q: # of trunks? A: 3  
A<sub>1</sub>: 1    A<sub>2</sub>: 0    A<sub>3</sub>: 0    A<sub>4</sub>: 2    A<sub>5</sub>: 2

Figure 4: Errors found on Oscar<sub>large</sub><sup>+</sup> using image-cutting MRs. Due to the limited space, we replace phrase “*what is the total number*” with “#”.





Q: # of blue baseball bats? A: 0  
Q: # of brown benches? A: 0



Q: # of red traffic lights? A: 1  
Q: # of yellow traffic lights and green traffic lights? A: 1



Q: # of benches? A: 2  
Q: Is there no green bench? A: Yes



Q: # blue baseball bats? A: 1  
Q: # of brown benches? A: 2



Q: # of red traffic lights? A: 2  
Q: # of yellow traffic lights and green traffic lights? A: 0



Q: # of benches? A: 0  
Q: Is there no green bench? A: No



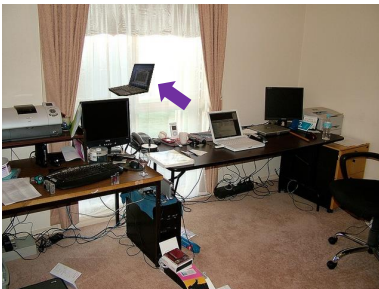
Q: # number of TVs? A: 1  
Q: # of mice? A: 2



Q: # of white frisbees? A: 1  
Q: # of yellow frisbees? A: 1



Q: # of lying persons? A: 0  
Q: # of black handbags? A: 1



Q: # number of TVs? A: 2  
Q: # of mice? A: 1



Q: # of white frisbees? A: 2  
Q: # of yellow frisbees? A: 0



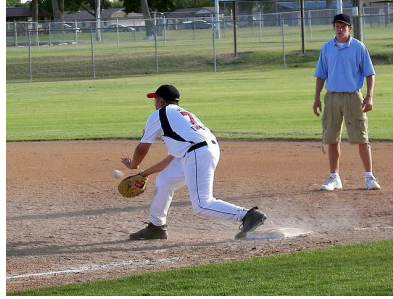
Q: # of lying persons? A: 2  
Q: # of black handbags? A: 2

Figure 5: Errors found on Oscar<sup>+</sup><sub>large</sub> using object-insertion MRs. The inserted objects are pinpointed in the figures. Due to the limited space, we replace phrase “what is the total number” with “#”.





Q: # of gray tennis rackets? A: 0



Q: # of standing persons? A: 2



Q: # of cars and birds? A: 2



Q: # of gray tennis rackets? A: 2



Q: # of standing persons? A: 1



Q: # of cars and birds? A: 1



Q: # of tennis rackets? A: 2



Q: # of standing persons? A: 2



Q: # of walking persons? A: 2



Q: # of tennis rackets? A: 0



Q: # of standing persons? A: 0



Q: # of walking persons? A: 3

Figure 6: Errors found on Oscar<sub>large</sub><sup>+</sup> using object-removal MRs. The removed objects are pinpointed in the figures. Due to the limited space, we replace phrase “what is the total number” with “#”.





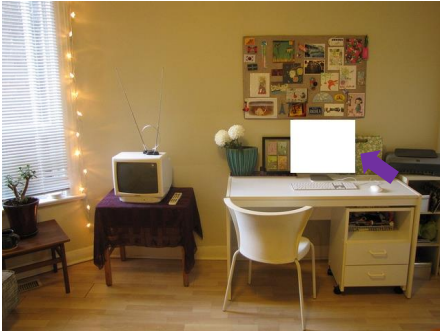
Q: # of potted plants? A: 2



Q: # of potted plants? A: 0



Q: # of golden sinks? A: 3



Q: # of potted plants? A: 1



Q: # of potted plants? A: 1



Q: # of golden sinks? A: 2



Q: # of skating persons? A: 1



Q: # black refrigerators? A: 0



Q: # of blue mice? A: 0



Q: # of skating persons? A: 0



Q: # black refrigerators? A: 1



Q: # of blue mice? A: 1

Figure 7: Errors found on Oscar<sup>+</sup><sub>large</sub> using object-removal MRs. The removed objects are pinpointed in the figures. Due to the limited space, we replace phrase “what is the total number” with “#”.





Q: # of toilets? A: 1



Q: # of couches? A: 2



Q: # of bottles? A: 2



Q: # of toilets? A: 1



Q: # of couches? A: 2



Q: # of bottles? A: 2



Q: # of boats? A: 10



Q: # of clocks? A: 1



Q: # of persons? A: 2



Q: # of boats? A: 4



Q: # of clocks? A: 2



Q: # of persons? A: 2

Figure 8: Errors found on GridFeat+MoViE<sup>+</sup> using object-removal (**Removal<sup>+</sup>**) MRs. The removed objects are **pinpointed** in the figures. Due to the limited space, we replace phrase “*what is the total number*” with “#”.





Q: # of chairs? A: 2



Q: # of cars? A: 2



Q: # of mice? A: 1



Q: # of chairs? A: 2



Q: # of cars? A: 2



Q: # of mice? A: 1



Q: # of persons? A: 3



Q: # of microwaves? A: 1



Q: # of cars? A: 3



Q: # of persons? A: 3



Q: # of microwaves? A: 1



Q: # of cars? A: 3

Figure 9: Errors found on GridFeat+MoViE<sup>+</sup> using object-removal (**Removal<sup>+</sup>**) MRs. The removed objects are pinpointed in the figures. Due to the limited space, we replace phrase “what is the total number” with “#”.



## C. False Positive Cases

This section provides some false positive cases we found on BERT-like models.



Q: What is the total number of cats and couches? A: 0  
Q<sub>1</sub>: What is the total number of cats? A<sub>1</sub>: 0  
Q<sub>2</sub>: What is the total number of couches? A<sub>2</sub>: 0  
Q<sub>3</sub>: What is the total number of couches and cats? A<sub>3</sub>: 0

Q: What is the total number of heads and plants? A: 0  
Q<sub>1</sub>: What is the total number of heads? A<sub>1</sub>: 0  
Q<sub>2</sub>: What is the total number of plants? A<sub>2</sub>: 0  
Q<sub>3</sub>: What is the total number of plants and heads? A<sub>3</sub>: 0



Q: What is the total number of birds and dogs? A: 0  
Q<sub>1</sub>: What is the total number of birds? A<sub>1</sub>: 0  
Q<sub>2</sub>: What is the total number of dogs? A<sub>2</sub>: 0  
Q<sub>3</sub>: What is the total number of dogs and birds? A<sub>3</sub>: 0

Q: What is the total number of birds and heads? A: 0  
Q<sub>1</sub>: What is the total number of birds? A<sub>1</sub>: 0  
Q<sub>2</sub>: What is the total number of heads? A<sub>2</sub>: 0  
Q<sub>3</sub>: What is the total number of heads and birds? A<sub>3</sub>: 0



Q: What is the total number of helmets and players? A: 0  
Q<sub>1</sub>: What is the total number of helmets? A<sub>1</sub>: 0  
Q<sub>2</sub>: What is the total number of players? A<sub>2</sub>: 0  
Q<sub>3</sub>: What is the total number of players and helmets? A<sub>3</sub>: 0

Q: What is the total number of bats and men? A: 0  
Q<sub>1</sub>: What is the total number of bats? A<sub>1</sub>: 0  
Q<sub>2</sub>: What is the total number of men? A<sub>2</sub>: 0  
Q<sub>3</sub>: What is the total number of men and bats? A<sub>3</sub>: 0



Q: What is the total number of buses and women? A: 0  
Q<sub>1</sub>: What is the total number of buses? A<sub>1</sub>: 0  
Q<sub>2</sub>: What is the total number of women? A<sub>2</sub>: 0  
Q<sub>3</sub>: What is the total number of women and buses? A<sub>3</sub>: 0

Q: What is the total number of men and buildings? A: 0  
Q<sub>1</sub>: What is the total number of men? A<sub>1</sub>: 0  
Q<sub>2</sub>: What is the total number of buildings? A<sub>2</sub>: 0  
Q<sub>3</sub>: What is the total number of buildings and men? A<sub>3</sub>: 0

Figure 10: False positive cases found on VisualBERT.



#### D. Definition of $P_{\text{human}}$ and $P_{\text{color}}$

This section provides the definition of  $P_{\text{human}}$  and  $P_{\text{color}}$ . We collect  $P_{\text{human}}$  and  $P_{\text{color}}$  from the attributes of Visual G.

$P_{\text{color}}$	$P_{\text{human}}$
black, blue, brown, dark, golden, gray, green, orange, pink, purple, red, white, yellow	throwing, blowing, cooking, flying, skiing, surfing, sleeping, jumping, swimming, crouching, waving, rolling, grazing, walking, smiling, laying, batting, sitting, running, moving, eating, skating, playing, standing, driving, watching, parking, bending, hanging, squatting, riding, landing, resting, looking, holding, racing, kneeling, sliding, lying, serving, hitting, pointing, posing, wearing, swinging, laughing, talking, skateboarding

Figure 11: Definition of  $P_{\text{human}}$  and  $P_{\text{color}}$ , which are used in the property-oriented question partitioning scheme.