

---

# Neuro-Symbolic Visual Reasoning: Disentangling “Visual” from “Reasoning”

---

Saeed Amizadeh<sup>1</sup> Hamid Palangi<sup>\*2</sup> Oleksandr Polozov<sup>\*2</sup> Yichen Huang<sup>2</sup> Kazuhito Koishida<sup>1</sup>

## Appendix A: Proofs

*Proof. Lemma 3.1:* Let  $X$  be the left most variable appearing in formula  $\mathcal{F}(X, \dots)$ , then depending on the quantifier  $q$  of  $X$ , we will have:

$$\begin{aligned}
 \text{If } q = \forall: \alpha(\mathcal{F}) &= \Pr(\mathcal{F} \Leftrightarrow \top \mid \mathcal{V}) \\
 &= \Pr\left(\bigwedge_{i=1}^N \mathcal{F}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V}\right) \\
 &= \prod_{i=1}^N \Pr(\mathcal{F}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V}) \\
 &= \prod_{i=1}^N \alpha(\mathcal{F} \mid x_i) = \mathcal{A}_{\forall}(\alpha(\mathcal{F} \mid X))
 \end{aligned}$$

$$\begin{aligned}
 \text{If } q = \exists: \alpha(\mathcal{F}) &= \Pr(\mathcal{F} \Leftrightarrow \top \mid \mathcal{V}) \\
 &= \Pr\left(\bigvee_{i=1}^N \mathcal{F}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V}\right) \\
 &= 1 - \prod_{i=1}^N \Pr(\mathcal{F}_{X=x_i} \Leftrightarrow \perp \mid \mathcal{V}) \\
 &= 1 - \prod_{i=1}^N (1 - \alpha(\mathcal{F} \mid x_i)) = \mathcal{A}_{\exists}(\alpha(\mathcal{F} \mid X))
 \end{aligned}$$

$$\begin{aligned}
 \text{If } q = \# : \alpha(\mathcal{F}) &= \Pr(\mathcal{F} \Leftrightarrow \top \mid \mathcal{V}) \\
 &= \Pr\left(\bigwedge_{i=1}^N \mathcal{F}_{X=x_i} \Leftrightarrow \perp \mid \mathcal{V}\right) \\
 &= \prod_{i=1}^N \Pr(\mathcal{F}_{X=x_i} \Leftrightarrow \perp \mid \mathcal{V}) \\
 &= \prod_{i=1}^N (1 - \alpha(\mathcal{F} \mid x_i)) = \mathcal{A}_{\#}(\alpha(\mathcal{F} \mid X))
 \end{aligned}$$

---

<sup>\*</sup>Equal contribution <sup>1</sup>Microsoft Applied Sciences Group (ASG), Redmond WA, USA <sup>2</sup>Microsoft Research AI, Redmond WA, USA. Correspondence to: Saeed Amizadeh <saamizad@microsoft.com>.

Note that the key underlying assumption in deriving the above proofs is that the binary logical statements  $\mathcal{F}_{X=x_i}$  for all objects  $x_i$  are independent random variables *given* the visual featurization of the scene, which is a viable assumption.  $\square$

$$\text{Proof. Lemma 3.2: } \alpha(\mathcal{F} \mid X) = [\Pr(\mathcal{F}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V})]_{i=1}^N = [\Pr(\top \Leftrightarrow \top \mid \mathcal{V})]_{i=1}^N = \mathbf{1} \quad \square$$

*Proof. Lemma 3.3:*

(A) If  $\mathcal{F}(X, Y, Z, \dots) = \neg\mathcal{G}(X, Y, Z, \dots)$ :

$$\begin{aligned}
 \alpha(\mathcal{F} \mid X) &= [\Pr(\mathcal{F}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V})]_{i=1}^N \\
 &= [\Pr(\mathcal{G}_{X=x_i} \Leftrightarrow \perp \mid \mathcal{V})]_{i=1}^N \\
 &= [1 - \alpha(\mathcal{G} \mid x_i)]_{i=1}^N = \mathbf{1} - \alpha(\mathcal{G} \mid X)
 \end{aligned}$$

(B) If  $\mathcal{F}(X, Y, Z, \dots) = \pi(X) \wedge \mathcal{G}(X, Y, Z, \dots)$  where  $\pi(\cdot)$  is a unary predicate:

$$\begin{aligned}
 \alpha(\mathcal{F} \mid X) &= [\Pr(\mathcal{F}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V})]_{i=1}^N \\
 &= [\Pr(\pi(x_i) \wedge \mathcal{G}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V})]_{i=1}^N \\
 &= [\Pr(\pi(x_i) \Leftrightarrow \top \wedge \mathcal{G}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V})]_{i=1}^N \\
 &= [\Pr(\pi(x_i) \Leftrightarrow \top \mid \mathcal{V}) \cdot \Pr(\mathcal{G}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V})]_{i=1}^N \\
 &= [\alpha(\pi \mid x_i) \cdot \alpha(\mathcal{G} \mid x_i)]_{i=1}^N \\
 &= \alpha(\pi \mid X) \odot \alpha(\mathcal{G} \mid X)
 \end{aligned}$$

(C) If  $\mathcal{F}(X, Y, Z, \dots) = [\bigwedge_{\pi \in \Pi_{XY}} \pi(X, Y)] \wedge \mathcal{G}(Y, Z, \dots)$  where  $\Pi_{XY}$  is the set of all binary predicates defined on variables  $X$  and  $Y$  in  $\mathcal{F}$  and  $Y$  is the left most variable in  $\mathcal{G}$  with quantifier  $q$ :

$$\begin{aligned}
 \alpha(\mathcal{F} \mid X) &= [\Pr(\mathcal{F}_{X=x_i} \Leftrightarrow \top \mid \mathcal{V})]_{i=1}^N \\
 &= \left[ \Pr\left(\underbrace{\left[\bigwedge_{\pi \in \Pi_{XY}} \pi(x_i, Y)\right]}_{\mathcal{R}_{x_i}(Y)} \wedge \mathcal{G} \Leftrightarrow \top \mid \mathcal{V}\right) \right]_{i=1}^N
 \end{aligned}$$

$$\stackrel{\text{L3.1}}{=} [\mathcal{A}_q(\alpha(\mathcal{R}_{x_i} \wedge \mathcal{G} \mid Y))]_{i=1}^N$$

$$\stackrel{\text{L3.3B}}{=} [\mathcal{A}_q(\alpha(\mathcal{R}_{x_i} \mid Y) \odot \alpha(\mathcal{G} \mid Y))]_{i=1}^N$$

$$\begin{aligned}
 &\stackrel{\text{L3.3B}}{=} \left[ \mathcal{A}_q \left( \left[ \bigodot_{\pi \in \Pi_{XY}} \alpha(\pi_{X=x_i} | Y) \right] \odot \alpha(\mathcal{G} | Y) \right) \right]_{i=1}^N \\
 &= \left[ \mathcal{A}_q \left( \left[ \bigodot_{\pi \in \Pi_{XY}} \alpha(\pi | x_i, Y) \right] \odot \alpha(\mathcal{G} | Y) \right) \right]_{i=1}^N \\
 &= \left[ \bigodot_{\pi \in \Pi_{XY}} \alpha(\pi | X, Y) \right] \times_q \alpha(\mathcal{G} | Y)
 \end{aligned}$$

Note that the key underlying assumption in deriving the above proofs is that all the unary and binary predicates  $\pi(x_i)$  and  $\pi(x_i, y_i)$  for all objects  $x_i$  and  $y_j$  are independent binary random variables *given* the visual featurization of the scene, which is a viable assumption.  $\square$

## Appendix B: The Language System

Our language system defines the pipeline to translate the questions in the natural language (NL) all the way to the DFOL language which we can then run to find the answer to the question. However, as opposed to many similar frameworks in the literature, our translation process takes place in two steps. First, we *parse* the NL question into the *task-dependent*, high-level, domain-specific language (DSL) of the target task. We then *compile* the resulted DSL program into the *task-independent*, low-level DFOL language. This separation is important because the  $\nabla$ -FOL core reasoning engine executes the task-independent, four basic operators of the DFOL language (i.e. **Filter**, **Relate**, **Neg** and  $\mathcal{A}_{\{\forall, \exists, \#\}}$ ) and *not* the task specific DSL operators. This distinguishes  $\nabla$ -FOL from similar frameworks in the literature as a *general-purpose* formalism; that is,  $\nabla$ -FOL can cover any reasoning task that is representable via first-order logic, and not just a specific DSL. This is mainly due to the fact that DFOL programs are equivalent to FOL formulas (up to reordering) as shown in Section 3.3. Figure 1 shows the proposed language system along with its different levels of abstraction. For more details, please refer to our PyTorch code base: <https://github.com/microsoft/DFOL-VQA>.

For the GQA task, we train a neural semantic parser using the annotated programs in the dataset to accomplish the first step of translation. For the second step, we simply use a *compiler*, which converts each high-level GQA operator into a composition of DFOL basic operators. Table 1 shows this (fixed) conversion along with the equivalent FOL formula for each GQA operator.

Most operators in the GQA DSL are parameterized by a set of NL tokens that specify the arguments of the operation (e.g. “*attr*” in **GFilter** specifies the attribute that the operator is expected to filter the objects based upon). In addition to the NL arguments, both terminal and non-terminal operators take as input the attention vector(s) on the objects present in the scene (except for **GSelect** which does not take

any input attention vector). However, in terms of their outputs, terminal and non-terminal operators are fundamentally different. A terminal operator produces a scalar likelihood or a list of scalar likelihoods (for “query” type operators). Because they are “terminal”, terminal operators have logical quantifiers in their FOL description; this, in turn, prompts the aggregation operator  $\mathcal{A}_{\{\forall, \exists, \#\}}$  in their equivalent DFOL translation. Non-terminal operators, on the other hand, produce attention vectors on the objects in the scene without calculating the aggregated likelihood.

## Appendix C: Some Examples from the Hard and the Easy Sets

In this appendix, we visually demonstrate a few examples from the hard and the easy subsets of the GQA Test-Dev split. Figures 2,3,4 show a few examples from the hard set with their corresponding questions, while Figures 5,6 show a few examples from the easy set. In these examples, the green rectangles represent where in the image the model is attending according to the attention vector  $\alpha(\mathcal{F} | X)$ . Here the formula  $\mathcal{F}$  represents either the entire question for the easy set examples or the partial question up until the point where the visual system failed to produce correct likelihoods for the hard set examples. We have included the exact nature of the visual system’s failure for the hard set examples in the captions. As illustrated in the paper, the visually hard-easy division here is with respect to the original Faster-RCNN featurization. This means that the “hard” examples presented here are *not* necessarily impossible in general, but are hard with respect to this specific featurization.

Furthermore, in Figure 7, we have demonstrated two examples from the hard set for which taking into the consideration the context of the question via the calibration process helped to overcome the imperfectness of the visual system and find the correct answer. Please refer to the caption for the details.

GQA OP	T	Equivalent FOL Description	Equivalent DFOL Program
<b>GSelect</b> ( <i>name</i> )[]	N	$name(X)$	$\mathbf{Filter}_{name}[1]$
<b>GFilter</b> ( <i>attr</i> ) $[\alpha_X]$	N	$attr(X)$	$\mathbf{Filter}_{attr}[\alpha_X]$
<b>GRelate</b> ( <i>name, rel</i> ) $[\alpha_X]$	N	$name(Y) \wedge rel(X, Y)$	$\mathbf{Filter}_{name}[\mathbf{Relate}_{rel, \exists}[\alpha_X]]$
<b>GVerifyAttr</b> ( <i>attr</i> ) $[\alpha_X]$	Y	$\exists X : attr(X)$	$\mathcal{A}_{\exists}(\mathbf{Filter}_{attr}[\alpha_X])$
<b>GVerifyRel</b> ( <i>name, rel</i> ) $[\alpha_X]$	Y	$\exists Y \exists X : name(Y) \wedge rel(X, Y)$	$\mathcal{A}_{\exists}(\mathbf{Filter}_{name}[\mathbf{Relate}_{rel, \exists}[\alpha_X]])$
<b>GQuery</b> ( <i>category</i> ) $[\alpha_X]$	Y	$[\exists X : c(X) \text{ for } c \text{ in } category]$	$[\mathcal{A}_{\exists}(\mathbf{Filter}_c[\alpha_X]) \text{ for } c \text{ in } category]$
<b>GChooseAttr</b> ( <i>a1, a2</i> ) $[\alpha_X]$	Y	$[\exists X : a(X) \text{ for } a \text{ in } [a_1, a_2]]$	$[\mathcal{A}_{\exists}(\mathbf{Filter}_a[\alpha_X]) \text{ for } a \text{ in } [a_1, a_2]]$
<b>GChooseRel</b> ( <i>n, r1, r2</i> ) $[\alpha_X]$	Y	$[\exists Y \exists X : n(Y) \wedge r(X, Y) \text{ for } r \text{ in } [r_1, r_2]]$	$[\mathcal{A}_{\exists}(\mathbf{Filter}_{name}[\mathbf{Relate}_{rel, \exists}[\alpha_X]]) \text{ for } r \text{ in } [r_1, r_2]]$
<b>GExists</b> () $[\alpha_X]$	Y	$\exists X \dots$	$\mathcal{A}_{\exists}(\alpha_X)$
<b>GAnd</b> () $[\alpha_X, \alpha_Y]$	Y	$\exists X \dots \wedge \exists Y \dots$	$\mathcal{A}_{\exists}(\alpha_X) \cdot \mathcal{A}_{\exists}(\alpha_Y)$
<b>GOr</b> () $[\alpha_X, \alpha_Y]$	Y	$\exists X \dots \vee \exists Y \dots$	$1 - (1 - \mathcal{A}_{\exists}(\alpha_X)) \cdot (1 - \mathcal{A}_{\exists}(\alpha_Y))$
<b>GTwoSame</b> ( <i>category</i> ) $[\alpha_X, \alpha_Y]$	Y	$\exists X \exists Y \bigvee_{c \in category} (c(X) \wedge c(Y))$	$\mathcal{A}_{\exists}([\mathcal{A}_{\exists}(\mathbf{Filter}_c[\alpha_X]) \cdot \mathcal{A}_{\exists}(\mathbf{Filter}_c[\alpha_Y]) \text{ for } c \text{ in } category])$
<b>GTwoDifferent</b> ( <i>category</i> ) $[\alpha_X, \alpha_Y]$	Y	$\exists X \exists Y \bigwedge_{c \in category} (\neg c(X) \vee \neg c(Y))$	$1 - \mathbf{GTwoSame}(category)[\alpha_X, \alpha_Y]$
<b>GAllSame</b> ( <i>category</i> ) $[\alpha_X]$	Y	$\bigvee_{c \in category} \forall X : \dots \rightarrow c(X)$	$1 - \prod_{c \in category} \mathcal{A}_{\exists}(\alpha_X \odot \mathbf{Neg}[\mathbf{Filter}_c[\alpha_X]])$

Table 1. The GQA operators translated to our FOL formalism. Here the notation  $\alpha_X$  is the short form for the attention vector  $\alpha(\mathcal{F} | X)$  where  $\mathcal{F}$  represents the formula the system has already processed up until the current operator. For the sake of simplicity, we have not included all of our GQA DSL here but the most frequent ones. Also the “Relate”-related operators are only shown for the case where the input variable  $X$  is the “subject” of the relation. The formalism is the same for the “object” role case except that the order of  $X$  and  $Y$  are swapped in the relation. The column **T** in the table indicates whether the operator is terminal or not. The full DSL can be found at our code base: <https://github.com/microsoft/DFOL-VQA>.

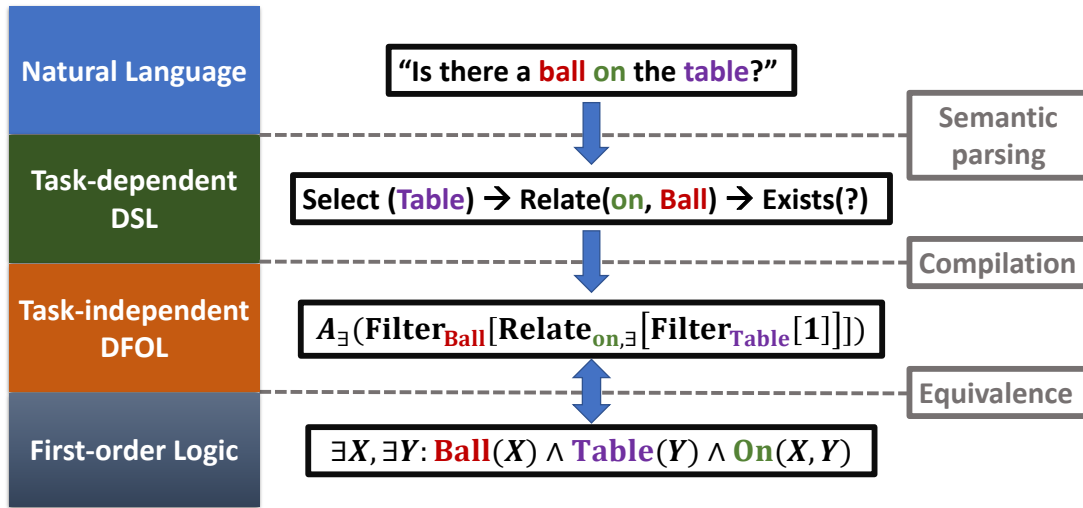


Figure 1. The language system: natural language question  $\xrightarrow{\text{semantic parser}}$  DSL program  $\xrightarrow{\text{compiler}}$  DFOL program  $\Leftrightarrow$  FOL formula.



(a)



(b)

Figure 2. **Hard Set:** (a) Q: “What are the rackets are lying on the top of?” As the attention bounding boxes show, the visual system has a hard time detecting the rackets in the first place and as a result is not able to reason about the rest of the question. (b) Q: “Does the boy’s hair have short length and white color?” In this example, the boy’s hair are not even visible, so even though the model can detect the boy, it cannot detect his hair and therefore answer the question correctly.

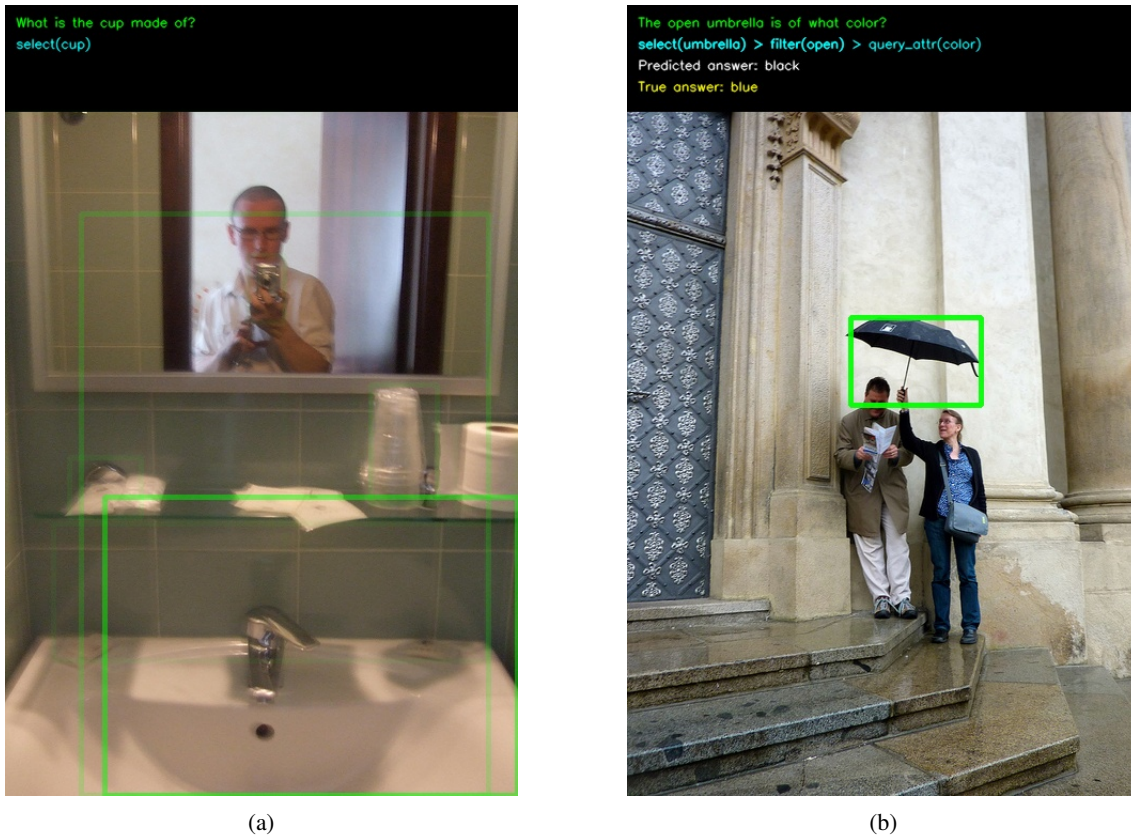


Figure 3. **Hard Set:** (a) Q: "What is the cup made of?" As the attention bounding boxes show, the visual system has a hard time finding the actual cups in the first place as they are pretty blurry. (b) Q: "The open umbrella is of what color?" In this example, the visual system was in fact able to detect an object that is both "umbrella" and "open" but its color is ambiguous and can be classified as "black" even by the human eye. However, the ground truth answer is "blue" which is hard to see visually.

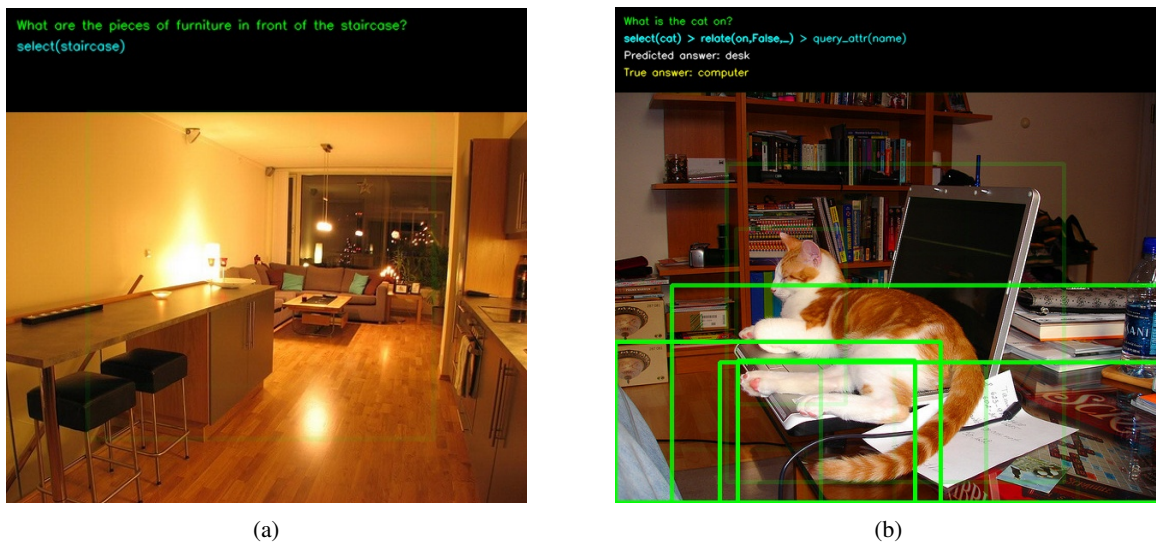
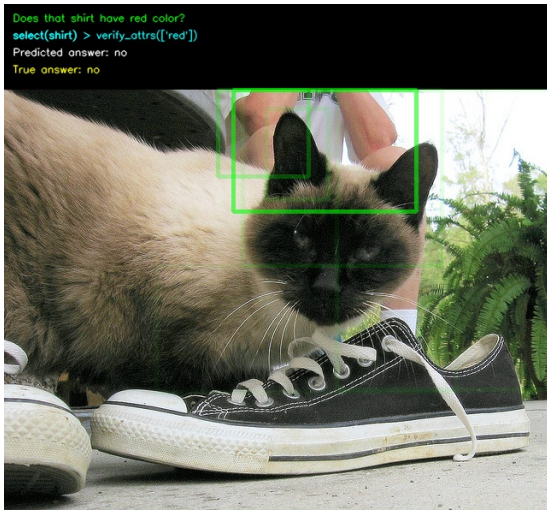


Figure 4. **Hard Set:** (a) Q: "What are the pieces of furniture in front of the staircase?" In this case, the model has a hard time detecting the staircase in the scene in the first place and therefore cannot find the correct answer. (b) Q: "What's the cat on?" In this example, the visual system can in fact detect the cat and supposedly the object that cat is "on"; however, it cannot infer the fact that there is actually a laptop keyboard invisible between the cat and the desk.

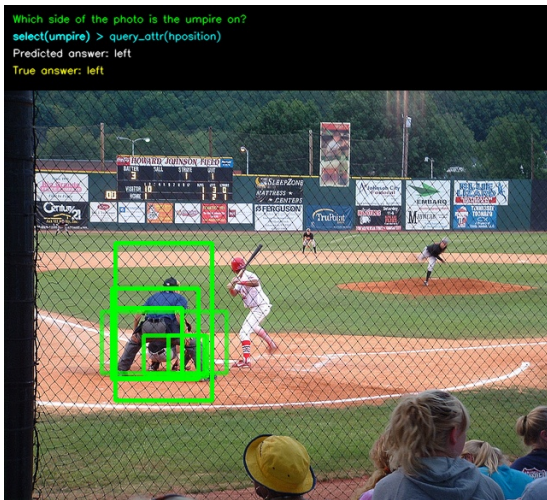


(a)

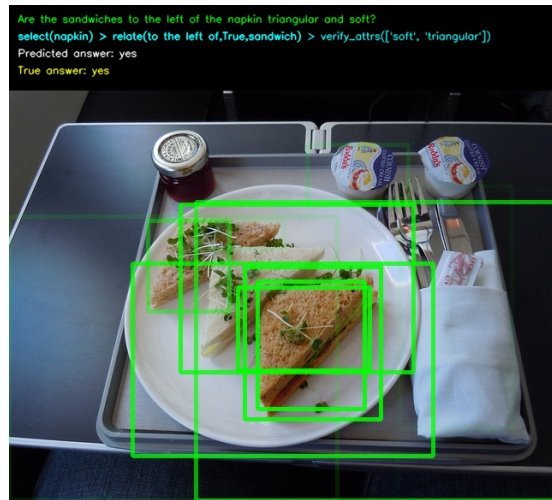


(b)

Figure 5. Easy Set: (a) Q: "Does that shirt have red color?" (b) Q: "Are the glass windows round and dark?"



(a)



(b)

Figure 6. Easy Set: (a) Q: "What side of the photo is umpire on?" (b) Q: "Are the sandwiches to the left of the napkin triangular and soft?"



(a)



(b)

Figure 7. **(a)** Q: "Are there any lamps next to the books on the right?" Due to the similar color of the lamp with its background, the visual oracle assigned a low probability for the predicate 'lamp' which in turn pushes the answer likelihood below 0.5. The calibration, however, was able to correct this by considering the context of 'books' in the image. **(b)** Q: "Is the mustard on the cooked meat?" In this case, the visual oracle had a hard time recognizing the concept of 'cooked' which in turn pushes the answer likelihood below 0.5. The calibration, however, was able to alleviate this by considering the context of 'mustard' and 'meat' in the visual input and boosts the overall likelihood.