

# Supplementary Materials for “Counterfactual VQA: A Cause-Effect Look at Language Bias”

This supplementary document is organized as follows:

- Section 1 introduces that RUBi [5] and Learned-Mixin [6] can be unified into our counterfactual inference framework.
- Section 2 provides an analysis of estimating NDE using the learnable parameter.
- Section 3 describes the implementation details.
- Section 4 describes the supplementary quantitative and qualitative results.

## 1. Revisiting RUBi and Learned-Mixin

As mentioned in Section 4.3, RUBi [5] and Learned-Mixin [6] can be unified into our counterfactual inference framework, which (1) follow a simplified causal graph without the direct path  $V \rightarrow A$ , and (2) use natural indirect effect (NIE) for inference. The detailed analysis is provided as follows.

### 1.1. Cause-Effect Look

Recent works RUBi [5] and Learned-Mixin [6] apply an ensemble architecture with a vision-language branch  $\mathcal{F}_{VQ}$  and a question-only branch  $\mathcal{F}_Q$ , while the direct relation between vision and answer is not formulated. The architecture is shown in Figure 1 (a).

Note that total effect can be decomposed into natural direct effect (NDE) and total indirect effect (TIE). As introduced in the main paper, we remove language bias by subtracting the natural direct effect from the total effect. The TIE is calculated by:

$$\begin{aligned} TE &= Z_{q,k} - Z_{q^*,k^*}, \\ NDE &= Z_{q,k^*} - Z_{q^*,k^*}, \\ TIE &= TE - NDE = Z_{q,k} - Z_{q^*,k^*}, \end{aligned} \quad (1)$$

which corresponds to Eq. (4) in the main paper. An alternative option to reduce language bias is to subtract the total direct effect (TDE) of questions on answers from total effect, which is formulated as:

$$\begin{aligned} TDE &= Z_{q,k} - Z_{q^*,k}, \\ NIE &= TE - TDE = Z_{q^*,k} - Z_{q^*,k^*}. \end{aligned} \quad (2)$$

Intuitively, both TIE and NIE reflect the increase of confidence for the answer given the visual knowledge, *i.e.*, from  $k^*$  to  $k$ . The difference between TIE and NIE is the existence of question  $q$ . The question  $q$  is block to calculate NIE (*i.e.*,  $q^*$ ), while  $q$  is given to calculate TIE. We use TIE to reserve  $q$  as the language context. In addition, both TDE and NDE reflect the increase of confidence for the answer given the question, *i.e.*, from  $q^*$  to  $q$ . The difference between TDE and NDE is also the existence of question  $q$ . Note that we hope to exclude the effect directly caused by question. Therefore, the mediator knowledge should be blocked when estimating the pure language effect, which is captured by NDE.

### 1.2. Implementation

RUBi [5] and Learned-Mixin (LM) [6] use the following fusion strategies for ensemble-based training:

$$(RUBi) \quad h(Z_q, Z_k) = Z_k \cdot \sigma(Z_q) \quad (3)$$

$$(LM) \quad h(Z_q, Z_k) = \log \sigma(Z_k) + g(k) \cdot \log \sigma(Z_q) \quad (4)$$

where  $\sigma(\cdot)$  represents the sigmoid function, and  $g(\cdot)$  is a learned function  $\mathbb{R}^d \rightarrow \mathbb{R}^1$  with the knowledge representation  $k \in \mathbb{R}^d$  as input and a scalar weight as output. During the test stage, they use  $Z_k$  for inference.

Perhaps surprising As for RUBi, NIE is calculated as:

$$NIE = \underbrace{z_k \cdot \sigma(c)}_{Z_{q^*,k}} - \underbrace{c \cdot \sigma(c)}_{Z_{q^*,k^*}} \propto z_k \quad (5)$$

As for Learned-Mixin, NIE is calculated as:

$$\begin{aligned} NIE &= \underbrace{(\log \sigma(z_k) + g(k) \cdot \log \sigma(c))}_{Z_{q^*,k}} \\ &\quad - \underbrace{(\log \sigma(c) + g(k^*) \cdot \log \sigma(c))}_{Z_{q^*,k^*}} \propto z_k \end{aligned} \quad (6)$$

where  $c$ ,  $g(k)$  and  $g(k^*)$  are constants for the same sample. Therefore, we have  $NIE \propto z_k$  for both RUBi and Learned-Mixin, which is exactly the output score of the

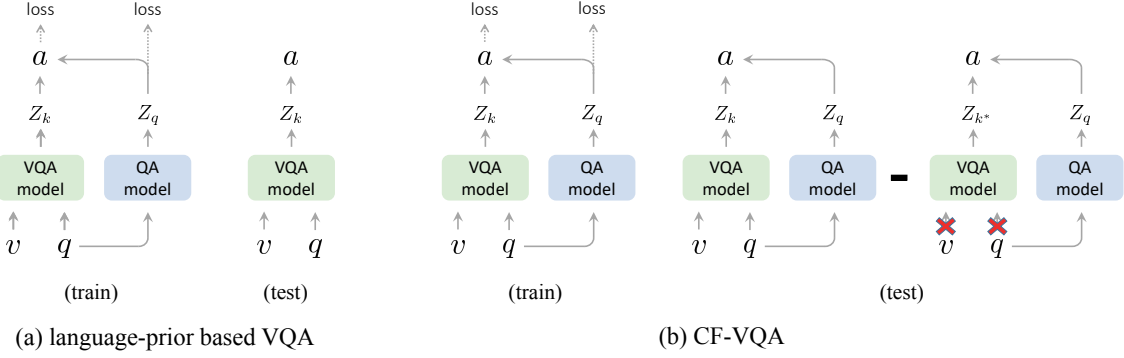


Figure 1: Comparison between our CF-VQA and language-prior based methods [5, 6] based on the simplified causal graph.

---

**Algorithm 1** Improving RUBi [5] using CF-VQA

---

```

1: function RUBi( $v, q, \text{is\_Training}; \theta, c$ )
2:    $z_q = \mathcal{F}_Q(q)$ 
3:    $z_k = \mathcal{F}_{VQ}(v, q)$ 
4:   if  $\text{is\_Training}$  then
5:      $z = z_k \cdot \sigma(z_q)$ 
6:     updating  $\theta$  according to  $\mathcal{L}_{cls}$ 
7:     updating  $c$  according to  $\mathcal{L}_{kl}$ 
8:   else
9:      $z = z_k \cdot z = (z_k - c) \cdot \sigma(z_q)$ 
10:  end if
11:  return  $z$ 
12: end function

```

---

vision-language branch  $\mathcal{F}_{VQ}$ . Note that RUBi and Learned-Mixin simply preserve the vision-language branch and uses  $z_k$  for inference. From our cause-effect view, *RUBi* and *Learned-Mixin* use natural indirect effect for inference.

### 1.3. Improving RUBi [5]

Thanks to our cause-effect look, RUBi [5] can be improved using CF-VQA, *i.e.*, using TIE for inference. Specifically, TIE for RUBi is calculated as:

$$TIE = \underbrace{z_k \cdot \sigma(z_q)}_{Z_{q,k}} - \underbrace{c \cdot \sigma(z_q)}_{Z_{q,k^*}} \quad (7)$$

where  $c$  denotes a learnable parameter. Table 5 in the main paper demonstrates that CF-VQA can outperform RUBi by 7% on VQA-CP v2. The red notes in Algorithm 1 show how RUBi is improved by changing several lines of code.

## 2. Analysis of Estimating NDE

In Section 4.2 in the main paper, we claimed that the learnable parameter  $c$  controls the sharpness of  $Z_{q,v^*,k^*}$  for estimating NDE. We give an intuitive analysis here.

For Harmonic (HM), we have:

$$(HM) \quad Z_{q,v^*,k^*} = \log \frac{\sigma(z_q) \cdot c_{HM}}{1 + \sigma(z_q) \cdot c_{HM}}, \quad (8)$$

where  $c_{HM} = (\sigma(c))^2 \in (0, 1)$ . We approximate the limits of  $Z_{q,v^*,k^*}$  and  $TIE = Z_{q,v,k} - Z_{q,v^*,k^*}$  as:

$$(HM) \quad \begin{aligned} \lim_{c_{HM} \rightarrow 0} Z_{q,v^*,k^*} &= -\infty \\ \lim_{c_{HM} \rightarrow 0} TIE &= z_{q,v,k} - C \\ &\propto z_{q,v,k}, \end{aligned} \quad (9)$$

where we use a extremely negative number  $C$  to replace  $-\infty$  for valid estimation of TIE. In this case, NDE is estimated as the same constant for all the answers, and TIE is dominated by  $z_{q,v,k}$ , which means that the language bias is not reduced. For  $c_{HM} \rightarrow 1$ , we have

$$(HM) \quad \begin{aligned} \lim_{c_{HM} \rightarrow 1} Z_{q,v^*,k^*} &= \log \frac{\sigma(z_q)}{1 + \sigma(z_q)} \\ \lim_{c_{HM} \rightarrow 1} TIE &= \log \frac{\sigma(z_v) \cdot \sigma(z_k) \cdot (1 + \sigma(z_q))}{1 + \sigma(z_q) \cdot \sigma(z_v) \cdot \sigma(z_k)}. \end{aligned} \quad (10)$$

For SUM, we have

$$(SUM) \quad Z_{q,v^*,k^*} = \log \sigma(z_q + 2c), \quad (11)$$

where  $c \in (-\infty, +\infty)$ . We approximate the limits of  $Z_{q,v^*,k^*}$  and  $TIE = Z_{q,v,k} - Z_{q,v^*,k^*}$  as:

$$(SUM) \quad \begin{aligned} \lim_{c \rightarrow \infty} Z_{q,v^*,k^*} &= -\infty \\ \lim_{c \rightarrow \infty} TIE &= z_{q,v,k} - C \\ &\propto z_{q,v,k}. \end{aligned} \quad (12)$$

Similar to HM, TIE is dominated by  $z_{q,v,k}$ . For  $c \rightarrow +\infty$ , we have:

$$(SUM) \quad \begin{aligned} \lim_{c \rightarrow +\infty} Z_{q,v^*,k^*} &= 0 \\ \lim_{c \rightarrow +\infty} TIE &= z_{q,v,k}. \end{aligned} \quad (13)$$

Table 1: **Comparison on VQA-CP v2 val set.** “Base.” indicates the VQA base model.

		VQA-CP v2					VQA v2
		val (in-domain)				test (OOD)	val (in-domain)
		Base.	All	Y/N	Num.	Other	All
GRLS [8]	–	–	56.90	69.23	42.50	49.36	42.33
GradSup[10]	–	–	62.4	77.8	43.8	53.6	–
RandImg [12]	UpDn	–	54.24	64.22	34.40	50.46	55.37
CF-VQA (HM)	UpDn	–	65.47	79.09	45.86	57.86	49.74
CF-VQA (SUM)	UpDn	–	60.29	66.32	47.48	57.96	51.27
CF-VQA (HM)	S-MRL	–	63.08	75.76	44.88	55.99	53.55
CF-VQA (SUM)	S-MRL	–	57.86	66.24	44.98	53.38	55.05

Also, TIE is dominated by  $z_{q,v,k}$ . In both cases, the language bias cannot be excluded. This analysis shows that a extremely large or small  $c$  will fail to estimate NDE and TIE, and it is necessary to control the sharpness of NDE by selecting a optimal  $c$ . In the main paper, we use a KL-divergence in Eq. (17) to force the sharpness of NDE similar to that of TE.

### 3. Implementation Details

We use the same implementation of RUBi [5] for fair comparison, including feature representation, baseline architectures, and optimization.

**Image Representation.** Following the popular bottom-up attention mechanism [2], we use a Faster R-CNN based framework to extract visual features. We select top- $K$  region proposals for each image, where  $K$  is fixed as 36.

**Question Representation.** Following [4, 5], we first lowercase all the questions and remove the punctuation, and then use the pretrained Skip-thought encoder [9] with fine-tuning. The size of final embedding is set as 4800.

**Vision-Language Branch.** The vision-language branch consists of the image representation, question representation, and a visual knowledge encoder. The baseline models for encoding visual knowledge includes SAN [13], UpDn [2], and a simplified version of the recent architecture MUREL [4] (S-MUREL) proposed in [5]. In short, S-MUREL consists of a BLOCK [3] bilinear fusion between image and question representations for each region, and a MLP classifier composed of three fully connected layers with ReLU activations. The dimension are 2,048, 2,048, and 3,000. More details can be found in [5].

**Language-Only Branch.** The language-only branch consists of the question representation and a question-only classifier. The question-only classifier is implemented by a MLP with three fully connect layers with ReLU activations. Note that this MLP has the same structure with the classifier for vision-language branch with different parameters.

**Vision-Only Branch.** The vision-only branch is composed of the question representation and a vision-only classifier. The vision-only classifier has the same structure as the

language-only classifier with different parameters.

**Optimization.** All the experiments are conducted with the Adam optimizer for 22 epochs. The learning rate linearly increases from  $1.5 \times 10^{-4}$  to  $6 \times 10^{-4}$  for the first 7 epochs, and decays after 14 epochs by multiplying 0.25 every two epochs. The batch size is set as 256.

**Datasets.** The experiments are conducted on VQA-CP [1] and VQA [7] datasets. VQA-CP v1 and v2 are created by re-organizing the train and val splits of the VQA v1 and v2 datasets, respectively [1].

## 4. Supplementary Experimental Results

We have conducted the ablation study and compared CF-VQA with state-of-the-art methods in the main paper. In this section, we show supplementary experimental results.

### 4.1. Quantitative Results

As suggested by [12, 8, 10, 11], we further hold out 8,000 instances from the training set (*i.e.*, VQA-CP v2 val) to measure the in-domain performance. Note that the results on VQA v2 val set also measure the in-domain performance. The results are given in Table 1. Compared to GRLS [8], all of our variants outperform GRLS by large margins for both in-domain and out-of-distribution (OOD) settings. Compared to GradSup [10], CF-VQA (HM) achieves better results on both VQA-CP val set and test set. Compared to RandImg [12], CF-VQA (SUM) achieves competitive results on VQA-CP v2 test set, and outperforms RandImg on in-domain settings by over 3%. These results demonstrate that CF-VQA not only effectively reduces language bias, but also performs robustly.

Table 2 shows the ablation study on VQA-CP v1 test split. As shown in Table 2, CF-VQA is general to both *base-line VQA architectures* and *fusion strategies*, which is also demonstrated by the results on VQA-CP v2. Table 3 shows the ablation study on VQA-CP v1 test split using the simplified causal graph. Similarly, CF-VQA achieves significant improvement for different baseline VQA architectures and fusion strategies.

Table 2: **Ablation of CF-VQA** on VQA-CP v1 test set. “SAN/UpDn/S-MRL” denotes the baseline VQA model. “HM/SUM” represents the strategies that train the ensemble model and test with only the vision-language branch following ensemble-based method [5, 6]. \* represents the reproduced results.

	All	Y/N	Num.	Other		All	Y/N	Num.	Other		All	Y/N	Num.	Other
SAN*	32.50	36.86	12.47	36.22	UpDn*	37.08	42.46	12.76	41.50	S-MRL*	36.68	42.72	12.59	40.35
Harmonic	49.29	72.73	<b>20.57</b>	37.51	Harmonic	<b>55.75</b>	80.65	<b>24.72</b>	43.46	Harmonic	53.55	79.38	17.39	42.38
+ CF-VQA	<b>52.06</b>	<b>80.38</b>	16.88	<b>38.04</b>	+ CF-VQA	55.16	<b>82.27</b>	16.14	<b>43.87</b>	+ CF-VQA	<b>55.26</b>	<b>82.13</b>	<b>18.03</b>	<b>43.49</b>
SUM	38.34	49.88	<b>15.82</b>	35.91	SUM	52.78	78.71	14.30	42.45	SUM	49.44	76.49	16.23	35.90
+ CF-VQA	<b>52.87</b>	<b>84.94</b>	14.85	<b>36.26</b>	+ CF-VQA	<b>57.39</b>	<b>88.46</b>	<b>14.80</b>	<b>43.61</b>	+ CF-VQA	<b>57.03</b>	<b>89.02</b>	<b>17.08</b>	<b>41.27</b>

Table 3: **Ablation of CF-VQA with the simplified causal graph** on VQA-CP v1 test set. “SAN/UpDn/S-MRL” denotes the baseline VQA model. “HM/SUM” represents the strategies that train the ensemble model and test with only the vision-language branch following ensemble-based method [5, 6]. \* represents the reproduced results.

	All	Y/N	Num.	Other		All	Y/N	Num.	Other		All	Y/N	Num.	Other
SAN*	32.50	36.86	12.47	36.22	UpDn*	37.08	42.46	12.76	41.50	S-MRL*	36.68	42.72	12.59	40.35
Harmonic	46.83	66.64	19.45	38.13	Harmonic	54.13	80.60	15.75	43.24	Harmonic	54.51	80.82	17.30	43.29
+ CF-VQA	<b>54.48</b>	<b>83.73</b>	<b>22.73</b>	<b>38.15</b>	+ CF-VQA	<b>56.19</b>	<b>85.08</b>	<b>16.00</b>	<b>43.61</b>	+ CF-VQA	<b>56.82</b>	<b>86.01</b>	<b>17.38</b>	<b>43.63</b>
SUM	40.08	54.15	15.53	<b>35.95</b>	SUM	51.20	74.70	13.61	42.94	SUM	52.54	78.42	16.77	<b>41.18</b>
+ CF-VQA	<b>52.73</b>	<b>84.64</b>	<b>16.02</b>	35.75	+ CF-VQA	<b>56.80</b>	<b>87.76</b>	<b>13.89</b>	<b>43.25</b>	+ CF-VQA	<b>57.07</b>	<b>89.28</b>	<b>17.39</b>	41.00

## 4.2. Qualitative Results


Figure 2 illustrates examples to show how CF-VQA improves RUBi by simply replacing natural indirect effect with total indirect effect for inference following Algorithm 1. The examples show that CF-VQA benefits from language context, e.g., “large or small”, “deep or shallow”, and “real or a statue” in the first row. Some failure cases are shown in the last two rows. First, CF-VQA may tend to generate broad answers, e.g., “houses” v.s “church”, and “vegetables” v.s “peas”. Second, CF-VQA may ignore visual content like traditional likelihood strategy. Therefore, there remains the challenge about how to balance visual understanding and language context.

## References

- [1] Aishwarya Agrawal, Dhruv Batra, Devi Parikh, and Anirudha Kembhavi. Don’t just assume; look and answer: Overcoming priors for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4971–4980, 2018. 3
- [2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6077–6086, 2018. 3
- [3] Hedi Ben-Younes, Remi Cadene, Nicolas Thome, and Matthieu Cord. Block: Bilinear superdiagonal fusion for visual question answering and visual relationship detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8102–8109, 2019. 3
- [4] Remi Cadene, Hedi Ben-Younes, Matthieu Cord, and Nicolas Thome. Murel: Multimodal relational reasoning for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1989–1998, 2019. 3
- [5] Remi Cadene, Corentin Dancette, Matthieu Cord, Devi Parikh, et al. Rubi: Reducing unimodal biases for visual question answering. *Advances in Neural Information Processing Systems*, 32:841–852, 2019. 1, 2, 3, 4
- [6] Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. Don’t take the easy way out: Ensemble based methods for avoiding known dataset biases. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4060–4073, 2019. 1, 2, 4
- [7] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6904–6913, 2017. 3
- [8] Gabriel Grand and Yonatan Belinkov. Adversarial regularization for visual question answering: Strengths, shortcomings, and side effects. In *Proceedings of the Second Workshop on Shortcomings in Vision and Language*, pages 1–13, 2019. 3
- [9] Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. In *Advances in neural information processing systems*, pages 3294–3302, 2015. 3
- [10] Damien Teney, Ehsan Abbasnejad, and Anton van den Hengel. Learning what makes a difference from counterfactual examples and gradient supervision. *arXiv preprint arXiv:2004.09034*, 2020. 3
- [11] Damien Teney, Ehsan Abbasnejad, and Anton van den Hengel. Unshuffling data for improved generalization. *arXiv preprint arXiv:2002.11894*, 2020. 3
- [12] Damien Teney, Kushal Kafle, Robik Shrestha, Ehsan Abbasnejad, Christopher Kanan, and Anton van den Hengel. On

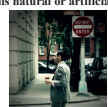


Is this area large or small?




RUBi + CF-VQA		RUBi	
<b>small</b>	94.6%	old	36.3%
large	3.5%	yes	32.8%
big	1.5%	no	18.7%
medium	0.3%	both	5.8%
huge	0.1%	<b>small</b>	3.0%

Is this natural or artificial light?




RUBi + CF-VQA		RUBi	
<b>natural</b>	50.4%	no	72.6%
yes	25.6%	yes	18.5%
both	7.7%	none	3.2%
curved	4.8%	old	1.3%
main	3.9%	unknown	1.3%

Is this fruit fresh or frozen?




RUBi + CF-VQA		RUBi	
<b>fresh</b>	63.9%	orange	25.7%
frozen	8.3%	red	15.7%
half	3.8%	none	14.3%
salt	2.2%	white	13.7%
wheat	1.6%	yellow	8.7%

What brand is the man's shirt?




RUBi + CF-VQA		RUBi	
<b>nike</b>	60.1%	billabong	42.8%
adidas	12.9%	none	19.1%
billabong	10.4%	blue	12.1%
hurley	8.8%	white	6.6%
polo	3.3%	hurley	4.1%

What brand is the racket?




RUBi + CF-VQA		RUBi	
<b>wilson</b>	50.2%	adidas	83.0%
nike	24.1%	<b>wilson</b>	6.5%
prince	9.1%	nike	4.8%
head	8.5%	w	2.8%
adidas	5.6%	white	0.8%

What type of herb is on the left?




RUBi + CF-VQA		RUBi	
<b>parsley</b>	52.1%	orange	73.7%
ginger	9.2%	red	18.6%
lily	6.0%	none	2.7%
pepper	5.7%	white	1.7%
maple	2.8%	yellow	0.9%

What type of sneakers are the players playing in?




RUBi + CF-VQA		RUBi	
<b>cleats</b>	40.7%	baseball	92.5%
baseball	10.9%	<b>cleats</b>	6.1%
tennis shoes	8.6%	baseball cap	0.6%
converse	5.9%	don't know	0.1%
giants	5.7%	yes	0.1%

What type of flower is in the vase?




RUBi + CF-VQA		RUBi	
<b>rose</b>	31.0%	pink	69.2%
daisy	28.3%	<b>rose</b>	21.8%
carnation	16.0%	red	4.6%
lily	6.2%	<b>roses</b>	2.1%
lilly	3.0%	purple	0.8%

What are these machines used for?




RUBi + CF-VQA		RUBi	
money	39.2%	<b>parking</b>	90.0%
transportation	28.2%	money	8.6%
<b>parking</b>	13.4%	picture	0.6%
<b>parking meter</b>	3.6%	<b>parking meter</b>	0.2%
riding	2.0%	driving	0.1%

Why are they wearing wetsuits?




RUBi + CF-VQA		RUBi	
safety	67.9%	<b>surfing</b>	77.5%
surf	7.7%	surf	11.9%
yes	5.1%	yes	3.9%
protection	4.0%	sunny	2.1%
<b>surfing</b>	3.5%	walking	1.0%

Is this water deep or shallow?



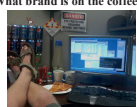
RUBi + CF-VQA		RUBi	
<b>deep</b>	97.0%	no	40.3%
<b>shallow</b>	2.5%	yes	35.7%
yes	0.3%	unknown	12.1%
ascending	0.1%	expert	2.3%
choppy	0.1%	old	1.8%

Is this bus going or coming?




RUBi + CF-VQA		RUBi	
<b>going</b>	78.4%	police	36.3%
stopped	8.6%	<b>going</b>	30.6%
<b>coming</b>	6.1%	forward	11.9%
leaving	1.1%	<b>coming</b>	6.2%
city	0.9%	no	4.3%

What brand is on the coffee cup?




RUBi + CF-VQA		RUBi	
<b>starbucks</b>	89.6%	none	37.7%
dunkin	5.1%	unknown	29.6%
donuts			
coca cola	4.7%	can't tell	7.8%
coke	0.5%	not sure	4.2%
jones	0.1%	not possible	3.5%

What brand of phone is this?




RUBi + CF-VQA		RUBi	
<b>iphone</b>	25.9%	nintendo	68.2%
motorola	24.5%	wii	26.3%
samsung	19.0%	none	4.3%
htc	10.2%	unknown	0.3%
<b>apple</b>	3.3%	<b>iphone</b>	0.3%

What brand is the box?




RUBi + CF-VQA		RUBi	
<b>hp</b>	32.2%	dell	64.2%
canon	21.5%	<b>hp</b>	26.5%
dell	15.1%	windows	6.9%
toshiba	13.4%	adidas	0.9%
head	8.0%	toshiba	0.6%

What type of seeds are stuck to the outside of the bun?




RUBi + CF-VQA		RUBi	
<b>sesame</b>	99.5%	unknown	41.7%
pepper	0.1%	none	30.1%
regular	0.1%	yellow	8.6%
sunflower	0.1%	<b>sesame</b>	6.6%
0	0.1%	white	2.9%

What type of sink is seen in the picture?




RUBi + CF-VQA		RUBi	
<b>pedestal</b>	49.6%	bathroom	48.5%
bathroom	10.3%	<b>pedestal</b>	40.0%
<b>porcelain</b>	8.0%	white	5.6%
ceramic	7.1%	ceramic	5.2%
white	4.1%	<b>porcelain</b>	1.0%

What type of building is pictured in the photo?




RUBi + CF-VQA		RUBi	
<b>clock tower</b>	28.8%	unknown	91.0%
<b>church</b>	15.1%	none	3.7%
apartment	9.5%	yellow	3.0%
school	7.9%	sesame	0.7%
tower	7.1%	white	0.5%

Where on the cow's body is there a tag?




RUBi + CF-VQA		RUBi	
<b>ear</b>	48.8%	yes	76.8%
yes	17.3%	no	6.5%
back	9.3%	left	3.0%
head	5.0%	unknown	2.9%
legs	3.2%	bowl	2.5%

What are these buildings?




RUBi + CF-VQA		RUBi	
houses	26.8%	<b>church</b>	73.1%
skyscrapers	13.8%	city	8.4%
office	6.7%	building	5.5%
tower	6.4%	london	1.8%
<b>church</b>	6.2%	castle	1.6%

What are those round green things?



RUBi + CF-VQA		RUBi	
vegetables	27.3%	<b>peas</b>	97.3%
peas	19.5%	grapes	1.2%
grapes	10.2%	vegetables	1.2%
beets	6.6%	fruit	0.1%
fruit	5.9%	blueberries	0.1%

What kind of window covering is shown?



RUBi + CF-VQA		RUBi	
blinds	96.5%	<b>curtain</b>	16.8%
shade	2.9%	canopy	15.1%
sheet	0.3%	fan	14.9%
<b>curtains</b>	0.2%	<b>curtains</b>	7.9%
cloth	0.1%	sheet	7.7%

Figure 2: Qualitative comparison of RUBi and RUBi+CF-VQA on VQA-CP v2 test split. Red bold answer denotes the ground-truth one.

the value of out-of-distribution testing: An example of good-hart's law. *arXiv preprint arXiv:2005.09241*, 2020. 3

[13] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and

Alex Smola. Stacked attention networks for image question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 21–29, 2016. 3