

Semi-supervised Grounding Alignment for Multi-modal Feature Learning

Shih-Han Chou^{1,2}, Zicong Fan³, James J. Little¹, Leonid Sigal^{1,2,4}

¹*Department of Computer Science, University of British Columbia, Vancouver, British Columbia, Canada*

²*Vector Institute for AI, Toronto, Ontario, Canada*

³*Department of Computer Science, ETH Zürich, Rämistrasse, Zürich, Switzerland*

⁴*Canada CIFAR AI Chair*

{shchou75, little, lsigal}@cs.ubc.ca; zicong.fan@inf.ethz.ch

In this supplementary material, we present additional results to complement the main paper.

I. GROUNDING ALIGNMENT

We show examples of the pseudo ground truth of grounding alignment on the Conceptual Captions dataset in Figure 1.

II. DETAILS OF THE GROUNDING ALIGNMENT MODULE

Figure 2 summarizes the details of the grounding alignment module.

III. HYPER-PARAMETER TUNING

We show hyper-parameter tuning in Table I. The settings with SPE or at phrase level achieve higher accuracy than the settings at token level which indicates that SPE and phrase level grounding are more important.

IV. ADJUSTING LEARNING RATES

The adjustment of learning rates is shown in Table II. We train our model with 0.25x, 0.5x, and 2x of the original setting.

V. ADDITIONAL QUALITATIVE RESULTS

A. Visual Grounding

We show some qualitative results of the visual grounding task in Figure 3. The orange boxes denote the ground truth annotations. The red boxes are the baseline [1] and the blue boxes are ours. The top two rows are the results where we perform better than the baseline. As shown in the first example, the ‘silver truck’, our predicted region is more aligned with the ground truth than the baseline. We also show some failure cases in the bottom row. Although our method does not exactly match the ground truth, the predictions are more reasonable than the baseline, we predict the oven but not the table in the first example in the bottom

row. We also show some examples in Figure 4 for the model pre-trained on 1/8 Conceptual Captions dataset and fine-tuned on 1/8 RefCOCO+ dataset.

B. Visual Question Answering

We show some qualitative results of the VQA task in Figure 5. The green color means the prediction is the same as the ground truth and the red color represents the wrong answer. The top row shows the positive examples where our method predicts the correct answers but the baseline doesn’t. We also show some failure cases in the bottom row. In the failure case, ‘How many vegetables are there?’, our model is confused by the beans so that the prediction is ‘3’ but not ‘2’. Similarly, in the last example, the color of the curtains is too dark to distinguish whether they are blue or black. We also show some examples in Figure 6 for the model pre-trained on 1/8 Conceptual Captions dataset and fine-tuned on 1/8 VQA dataset.

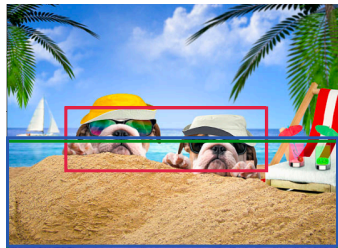
C. Visual Commonsense Reasoning

We show some qualitative results of the VCR task in Figure 7. The green color denotes the correct answers and the red color denotes the wrong answers. We show the probabilities among the answer choices. The examples show that our method is able to put the right attention on the correct answers. For the last example, although the answer is (b) in the QA \rightarrow R task, we argue that the answer choices are somewhat ambiguous. Answers (a) and (b) are both suitable for this question. We also show some examples in Figure 8 for the model pre-trained on 1/16 Conceptual Captions dataset and fine-tuned on 1/4 VCR dataset.

REFERENCES

- [1] Lu, Jiasen and Batra, Dhruv and Parikh, Devi and Lee, Stefan, “Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks,” in *NeurIPS*, 2019. 1, 3, 4, 5

Caption: happy dogs in the sand on the beach
NPs: [happy dogs, the sand, the beach]



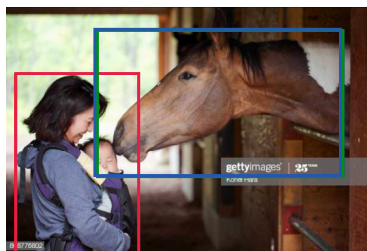
Caption: young women loading shopping bags in a car trunk
NPs: [young women loading shopping bags, a car trunk]



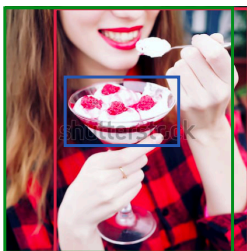
Caption: buddhist prayer flags in the barren region
NPs: [buddhist prayer flags, the barren region]



Caption: mother and baby communicating with a horse in the horse stable
NPs: [mother and baby, a horse, the horse]



Caption: close up portrait of young beautiful woman eating a dessert
NPs: [portrait, young beautiful woman, a dessert]



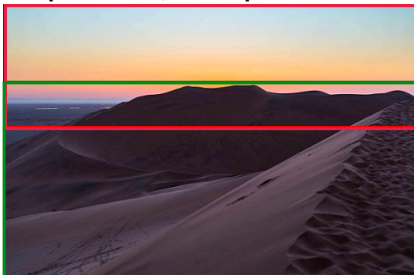
Caption: glasses on the buffet table in the restaurant
NPs: [glasses, the buffet table, the restaurant]



Caption: a woman carries her belongings
NPs: [a woman, belongings]



Caption: colorful sunset over the desert
NPs: [colorful sunset, the desert]



Caption: cowboys drive cattle down the road
NPs: [cowboys drive cattle, the road]



Figure 1: Pseudo Ground Truth of Grounding Alignment.

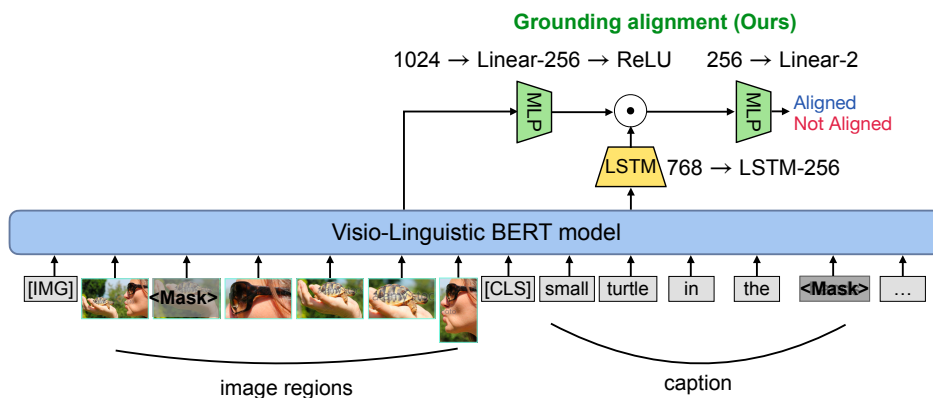


Figure 2: Detailed architecture of Grounding Alignment.

Table I: **Hyper-parameter Tuning.** We do hyper-parameter tuning using 1/8 amount of the Conceptual Caption Dataset. SPE denotes spatial positional encoding and SS Ground represents semi-supervised grounding alignment. Best results are highlighted in bold.

| Settings | SPE | λ_{align} | λ_{gnd} | token | phrase | Visual Grounding | VQA | VCR (Q→A) | VCR (QA→R) |
|----------------------|-----|-------------------|-----------------|-------|--------|------------------|--------------|--------------|--------------|
| ViLBERT [1] | - | 1 | - | - | - | 70.92 | 67.85 | 70.83 | 72.47 |
| + SS Ground (token) | | 0 | 1 | ✓ | | - | 67.38 | 71.28 | 72.73 |
| + SS Ground (token) | | 1 | 1 | ✓ | | - | 67.79 | 71.77 | 73.12 |
| + SS Ground (token) | | 1 | 20 | ✓ | | 70.98 | 67.84 | 71.91 | 73.36 |
| + SS Ground (phrase) | | 1 | 20 | | ✓ | - | 67.39 | 71.52 | 72.55 |
| + SPE | ✓ | 1 | 0 | | | 71.19 | 68.12 | 71.67 | 73.58 |
| + SS Ground (token) | ✓ | 1 | 20 | ✓ | | 69.02 | 67.64 | 71.72 | 73.06 |
| + SS Ground (phrase) | ✓ | 1 | 20 | | ✓ | 72.23 | 68.98 | 71.88 | 73.62 |

Table II: **Adjustment of Learning Rates.** We conduct the learning adjustment on the full dataset and fine-tune to visual grounding and VQA tasks. The setting is using our final model which includes SPE and semi-supervised grounding at the phrase level. Best results for each task are bold.

| Tasks / Learning Rate | lr=2.5e-05 | lr=5e-05 | lr=1e-04 (original) | lr=2e-04 |
|-----------------------|------------|----------|---------------------|----------|
| Visual Grounding | 71.74 | 71.83 | 72.47 | 71.45 |
| VQA | 68.82 | 69.19 | 69.63 | 68.02 |

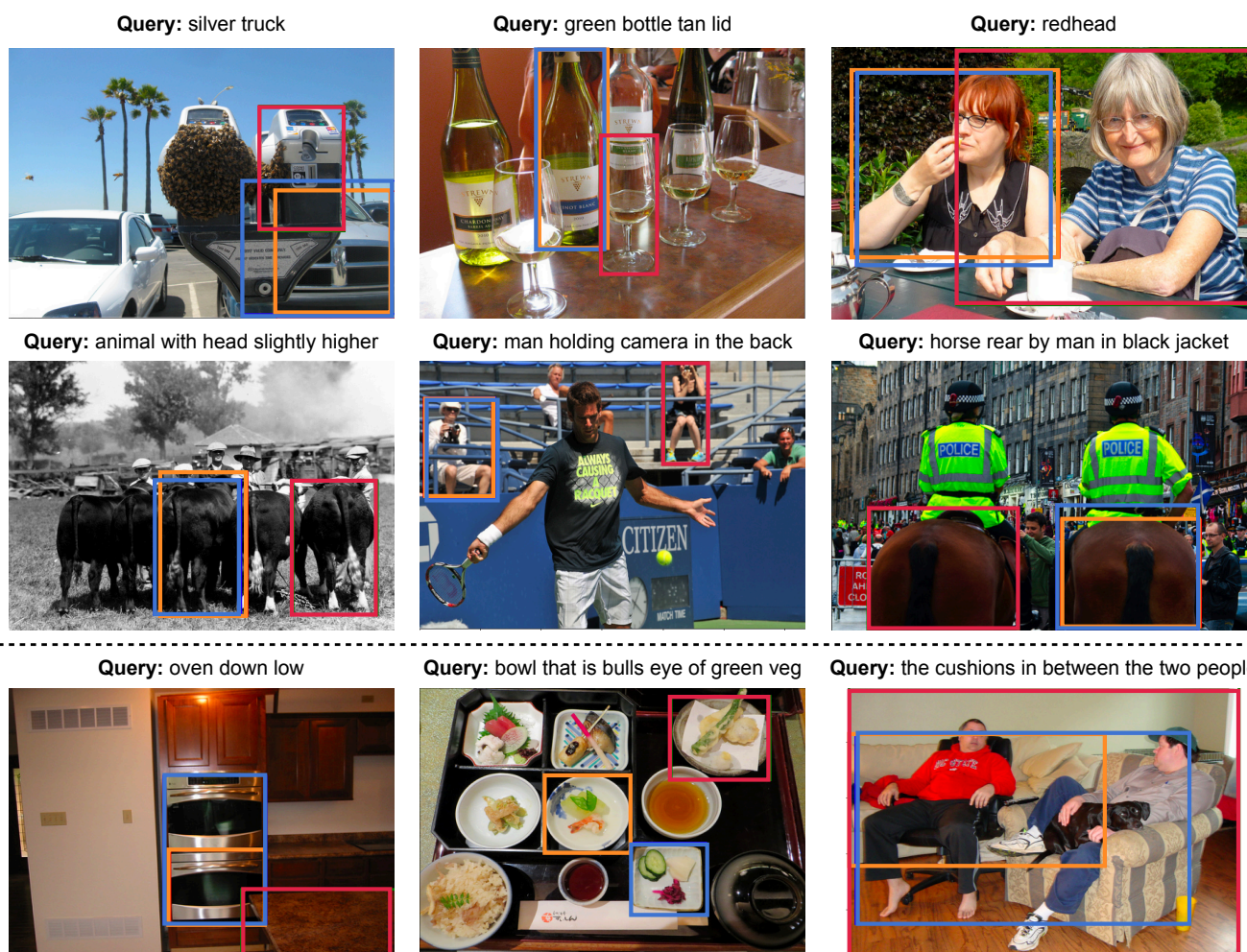
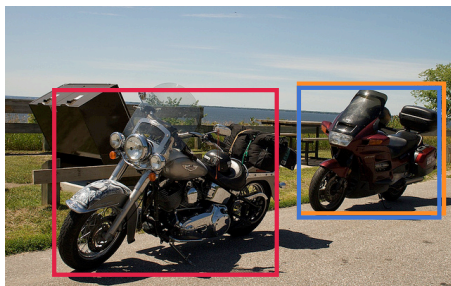


Figure 3: **Qualitative Results for visual grounding task pre-trained on full Conceptual Captions dataset and fine-tuned on full RefCOCO+ dataset.** Orange box: Ground truth. Red box: Baseline (ViLBERT [1]). Blue box: Ours. The top two rows are the positive examples and the bottom row are the failure examples.

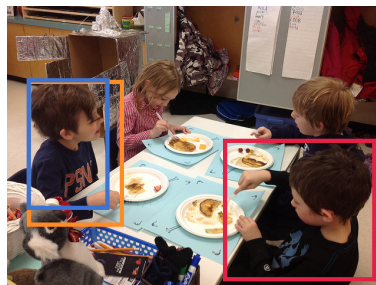
Query: Black cat under piped



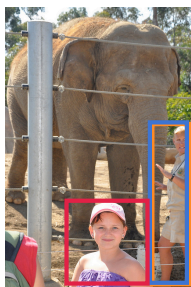
Query: black motorcycle behind blue motorcycle



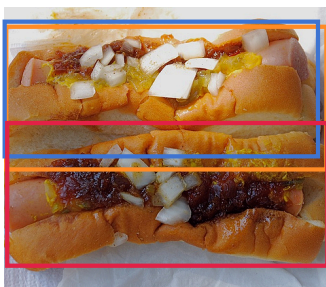
Query: boy with psn shirt



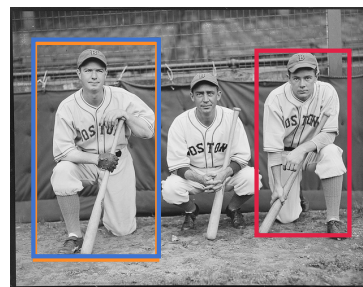
Query: woman in white closer to elephant



Query: hot dog with less chili



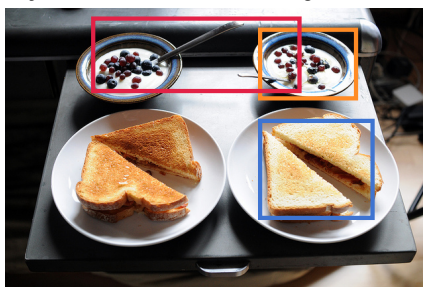
Query: Taller guy



Query: phantom magazine



Query: bowl of bubble in back of lightest sandwich



Query: peron next to man inn glasses

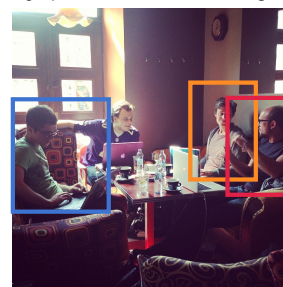


Figure 4: **Qualitative Results for visual grounding task pre-trained on 1/8 Conceptual Captions dataset and fine-tuned on 1/8 RefCOCO+ dataset.** Orange box: Ground truth. Red box: Baseline (ViLBERT [1]). Blue box: Ours. The top two rows are the positive examples and the bottom row are the failure examples.

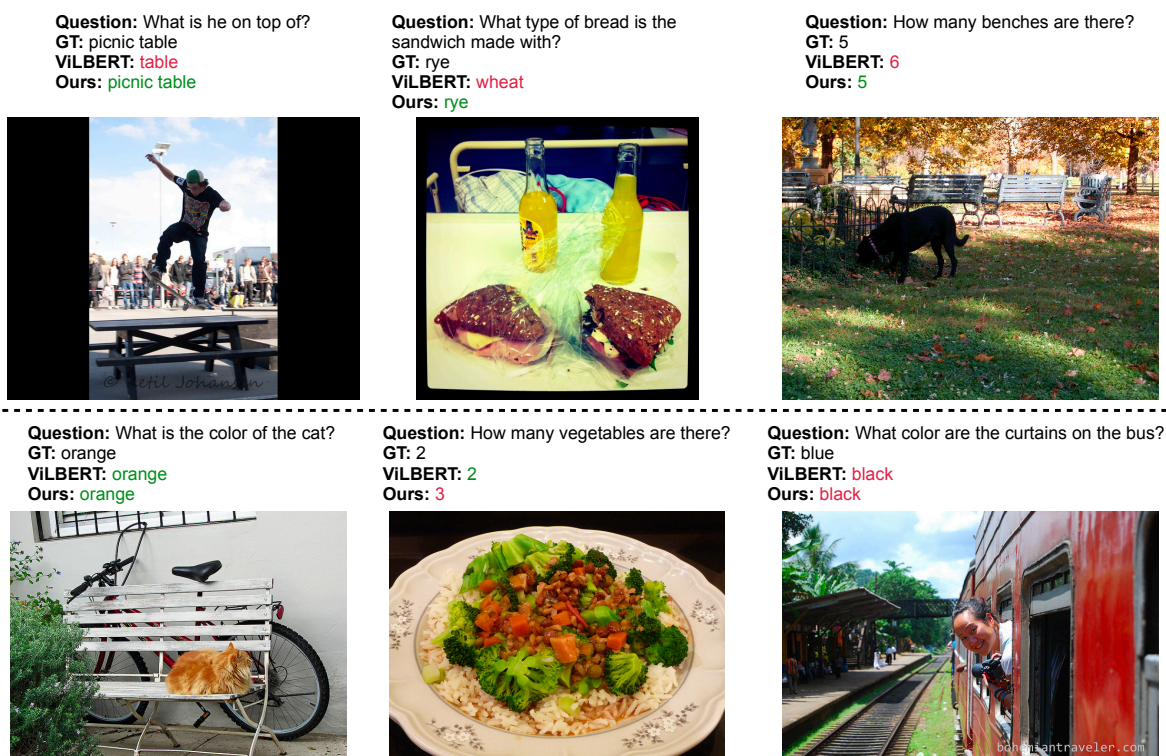


Figure 5: Qualitative Results for VQA task pre-trained on full Conceptual Captions dataset and fine-tuned on full VQA dataset. The top row is the positive examples and the two on the right of the bottom row are failure cases.

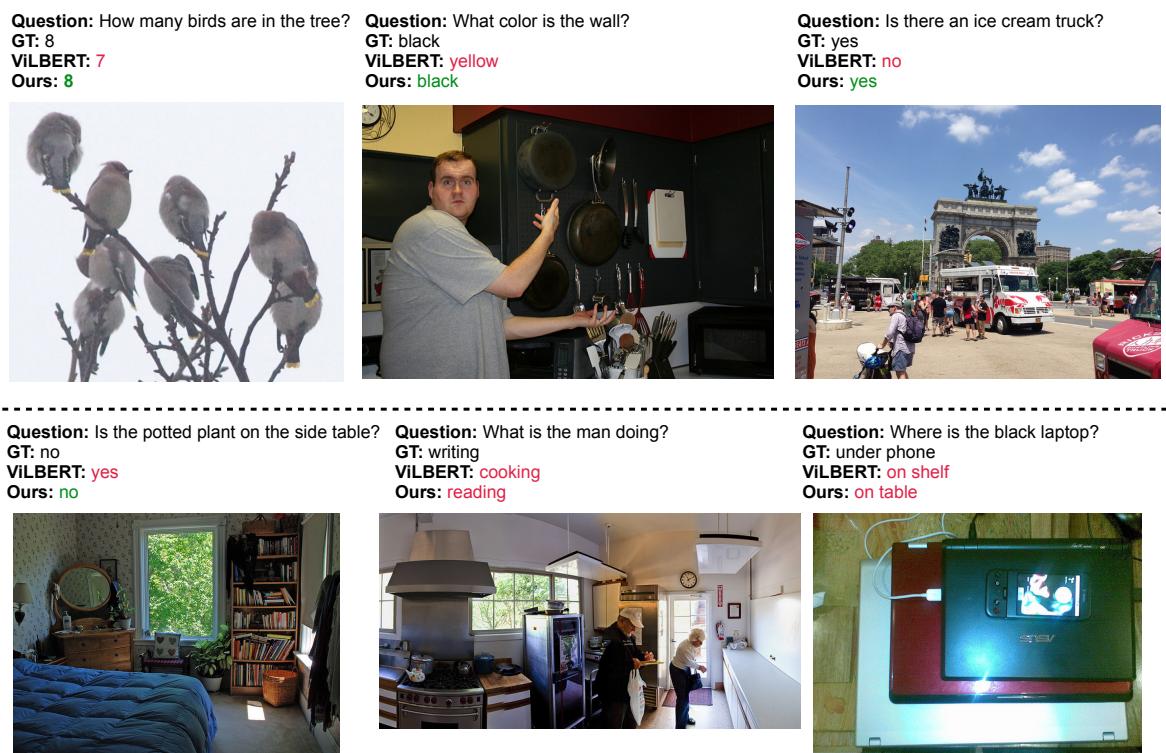
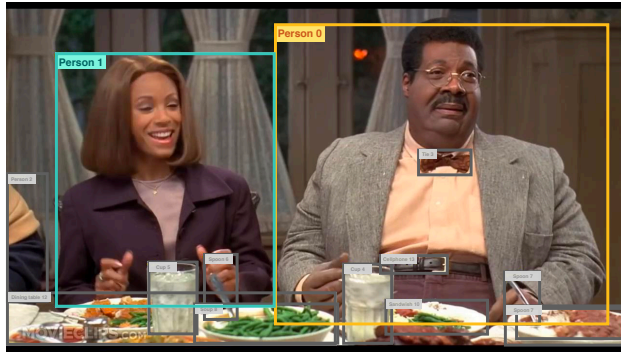


Figure 6: Qualitative Results for VQA task pre-trained on 1/8 Conceptual Captions dataset and fine-tuned on 1/8 VQA dataset. The top row is the positive examples and the two on the right of the bottom row are failure cases.



Q → A task

What are [person0, person1], looking at?

- (a) They are looking at something at the bottom of the cliff.
- (b) They are looking at the structure in front of them.
- (c) They are watching something out of the window of the plane.
- (d) [person0, person1] are looking at someone who is joking about [person0].

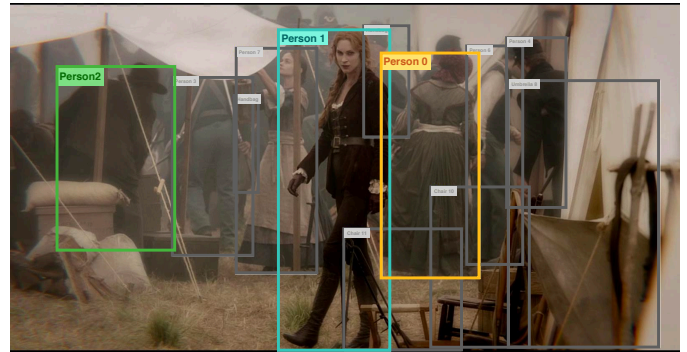
| Baseline | Ours |
|----------|------|
| ~0 | ~0 |
| ~0 | ~0 |
| ~0 | ~0 |
| ~1 | ~1 |

QA → R task

I think so because...

- (a) [person1] is smiling and opening her mouth as if about to speak to someone she is looking at.
- (b) [person1] is amused by what she is seeing while [person0], is embarrassed and uncomfortable.
- (c) [person2] is wearing a silly outfit and has her hair tied back to avoid getting hair in the food, the table is also looking at her and [person1] is speaking to [person2] like she is giving an order.
- (d) Looking away from [person0]'s eyes are closed and her mouth is pursed as if she is trying to contain her laughter.

| Baseline | Ours |
|----------|-------|
| ~0 | ~0 |
| 0.4397 | 0.999 |
| ~0 | ~0 |
| 0.5602 | ~0 |



Q → A task

Why is [person1] looking so defiant?

- (a) [person1] has tried something in this shop she does not like.
- (b) [person1] may not think that a teenage girl can do what he can.
- (c) [person1] is just in a stance that shows she is very excited about [person2, person0]'s dance moves.
- (d) Someone is threatening her.

| Baseline | Ours |
|----------|------|
| ~0 | ~0 |
| ~0 | ~0 |
| 0.011 | ~0 |
| 0.9889 | ~1 |

QA → R task

I think so because...

- (a) She has a look of fear and apprehension.
- (b) Someone could be forcing her to leave and she is looking back at her family.
- (c) She is gagged and tied up in order not to make any noise.
- (d) She has a daring expression on her face as if she is saying I dare you to come after me.

| Baseline | Ours |
|----------|--------|
| 0.2665 | 0.205 |
| 0.721 | 0.0014 |
| ~0 | ~0 |
| 0.012 | 0.7939 |



Q → A task

Why does [person0] keep this safe in a coat closet?

- (a) [person0] expects bad weather outside.
- (b) It is winter outside and [person0] needs to wear warm clothes.
- (c) This keeps [person0] safe.
- (d) So that people will not know where to find it.

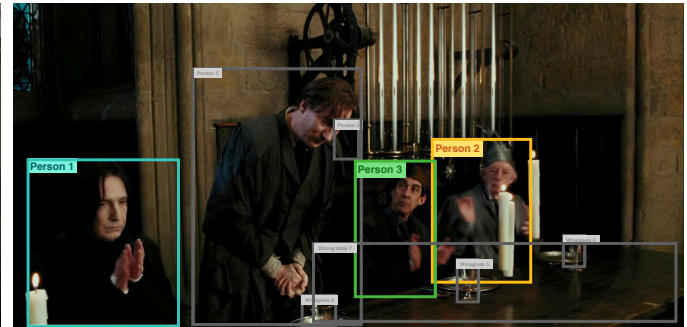
| Baseline | Ours |
|----------|--------|
| 0.945 | 0.0001 |
| ~0 | ~0 |
| 0.0001 | ~0 |
| 0.055 | 0.999 |

QA → R task

I think so because...

- (a) They are in this wooded like they're hiding so they would not want people to discover them.
- (b) Safes are designed to stay closed if you do not put in the right combination.
- (c) People ask for the gun to be kicked away so that it can't be immediately picked up again.
- (d) This is an unusual place to keep a safe, but might help to conceal it.

| Baseline | Ours |
|----------|--------|
| ~0 | ~0 |
| ~0 | ~0 |
| ~0 | ~0 |
| 0.9999 | 0.9999 |



Q → A task

What are [person1, person2], and [person3] doing ?

- (a) [person1, person2], and [person3] are having a group meeting.
- (b) [person1, person2], and [person3] are exiting the room with door open.
- (c) Clapping their hands.
- (d) [person1, person2], and [person3] are participating in a riot.

| Baseline | Ours |
|----------|-------|
| 0.954 | 0.038 |
| 0.0001 | 0.001 |
| 0.046 | 0.96 |
| ~0 | ~0 |

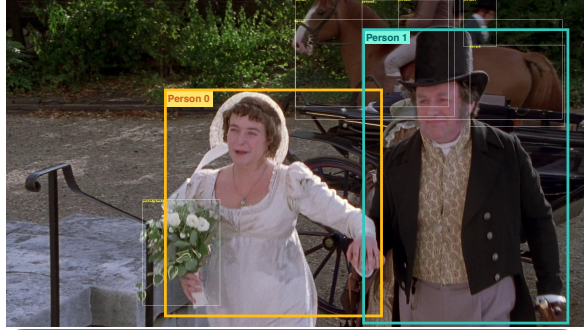
QA → R task

I think so because...

- (a) Clapping one's hands is a common sign of approval or support.
- (b) They are holding their hands in front of themselves and moving them back and forth.
- (c) People clap when they are happy or excited to see something or someone.
- (d) When you like the beat it is common to clap along.

| Baseline | Ours |
|----------|-------|
| 0.990 | 0.999 |
| 0.001 | ~0 |
| 0.008 | ~0 |
| ~0 | ~0 |

Figure 7: **Qualitative Results for VCR task pre-trained on full Conceptual Captions dataset and fine-tuned on full VCR dataset.** We show the ground truth answers with green color on the multiple choices side. The green/red colors on the probability side mean the correct/incorrect answers. The top row is the positive examples and the one on the right of the bottom row is the failure case.



Q → A task

Are [person0] and [person1] happy to get married?

- (a) Yes [person0, person1] are in love.
- (b) No, [person0, person1] are not discussing something happy.
- (c) No, they are not.
- (d) **Yes, they re both very happy today.**

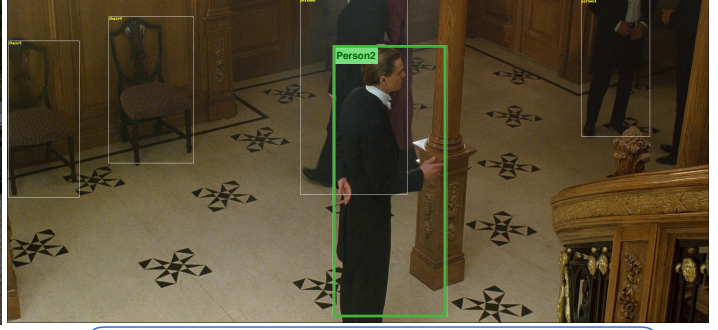
| Baseline | Ours |
|--------------|--------------|
| ~0 | ~0 |
| ~0 | ~0 |
| 0.999 | 0.068 |
| 0.0008 | 0.932 |

QA → R task

I think so because...

- (a) They re facing each other as they take their vows, while dressed in wedding attire.
- (b) [person0, person1] are dressed formally, [person3] has on a wedding dress and there is draping above them.
- (c) Both of them raise there arms up and slap hands together in a sign of celebration.
- (d) **They are both smiling and seem delighted.**

| Baseline | Ours |
|--------------|--------------|
| ~0 | ~0 |
| ~0 | ~0 |
| ~0 | ~0 |
| 0.999 | 0.999 |



Q → A task

What is [person2] doing next?

- (a) [person2] tells [person1] he is sneaking up on the staff to perform an audit.
- (b) He is going to attack [person1].
- (c) **He is going to greet someone.**
- (d) [person0] is going to retrieve his hat from the hat rack.

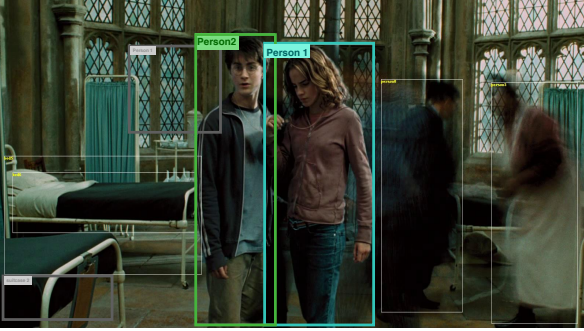
| Baseline | Ours |
|---------------|---------------|
| 0.7556 | 0.00016 |
| ~0 | ~0 |
| 0.2445 | 0.9998 |
| ~0 | ~0 |

QA → R task

I think so because...

- (a) He's standing next to the helicopter with an open door and someone is looking at him.
- (b) [person2] is standing near the door and wearing a store uniform, indicating that he works at the store. store greeters stand near the entrance of their store.
- (c) He is walking through a door and holding it open.
- (d) **He has his hand extended waiting on someone to arrive.**

| Baseline | Ours |
|----------|--------------|
| ~0 | ~0 |
| ~0 | ~0 |
| ~0 | ~0 |
| ~1 | 0.999 |



Q → A task

How does [person2] feel about [person1]?

- (a) [person2] is romantically interested in [person1].
- (b) [person2] is angry with [person1].
- (c) [person0] is worried and doesn't want [person1] to go, but can not do much else.
- (d) **[person2] trusts [person1] quite a lot.**

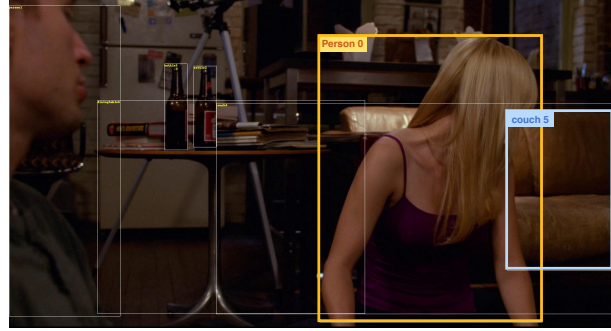
| Baseline | Ours |
|---------------|---------------|
| 0.8259 | 0.0028 |
| 0.0002 | 0.3915 |
| ~0 | ~0 |
| 0.1739 | 0.6057 |

QA → R task

I think so because...

- (a) He is looking at him keenly.
- (b) [person2] is making no move to step away from [person1] and he has his eyes closed as he continues to interact with him.
- (c) [person2] has his arms around [person1].
- (d) **He is standing close to her, and she is sharing something with him.**

| Baseline | Ours |
|---------------|---------------|
| ~0 | ~0 |
| ~0 | ~0 |
| ~0 | ~0 |
| 0.9999 | 0.9999 |



Q → A task

What is [person0] looking at?

- (a) [person0] is looking at some drawings.
- (b) **At something that is just out of view.**
- (c) [person0] is looking down.
- (d) [person1] is looking towards [person0].

| Baseline | Ours |
|---------------|----------------|
| 0.00102 | 0.99964 |
| ~0 | ~0 |
| 0.9981 | 0.00035 |
| ~0 | ~0 |

QA → R task

I think so because...

- (a) She is staring intently at something along with everyone else.
- (b) **She seems to have noticed something but because of [couch5] we can not see it.**
- (c) She is turned around looking over her shoulder.
- (d) She is looking straight ahead.

| Baseline | Ours |
|----------------|----------------|
| ~0 | ~0 |
| ~0 | 0.00012 |
| 0.99993 | 0.99988 |
| ~0 | ~0 |

Figure 8: **Qualitative Results for VCR task pre-trained on 1/16 Conceptual Captions dataset and fine-tuned on 1/4 VCR dataset.** We show the ground truth answers with green color on the multiple choices side. The green/red colors on the probability side mean the correct/incorrect answers. The top row is the positive examples and the one on the right of the bottom row is the failure case.