

# Recurrent Multimodal Interaction for Referring Image Segmentation (Supplementary Material)

Chenxi Liu<sup>1</sup> Zhe Lin<sup>2</sup> Xiaohui Shen<sup>2</sup> Jimei Yang<sup>2</sup> Xin Lu<sup>2</sup> Alan Yuille<sup>1</sup>  
Johns Hopkins University<sup>1</sup> Adobe Research<sup>2</sup>

{cxliu, alan.yuille}@jhu.edu {zlin, xshen, jimyang, xinl}@adobe.com

In section 1 of this supplementary material, we provide the full table of segmentation performance on the four datasets. In addition to the IOU reported in the main paper, we also report Precision@X ( $X \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$ ) which means the percentage of images with IOU higher than X. This is consistent with previous work [1, 2] to allow for comparison. We show that the improvement of our RMI model over the baseline is solid.

In section 2 we provide more visualization of intermediate segmentation results. We show from the progression how our RMI model can better memorize long-term information instead of being distracted by later words.

In section 3 we provide more visualization of final segmentation results. We show that our RMI model can usually generate better segmentation masks, and also DenseCRF [3] offers more visually pleasing results through refinement.

## References

- [1] R. Hu, M. Rohrbach, and T. Darrell. Segmentation from natural language expressions. In *ECCV (1)*, volume 9905 of *Lecture Notes in Computer Science*, pages 108–124. Springer, 2016. 1, 3
- [2] R. Hu, M. Rohrbach, S. Venugopalan, and T. Darrell. Utilizing large scale vision and text datasets for image segmentation from referring expressions. *CoRR*, abs/1608.08305, 2016. 1
- [3] P. Krähenbühl and V. Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *NIPS*, pages 109–117, 2011. 1

## 1. Full Table of Segmentation Performance

Table 1.1: Comparison of segmentation performance on Google-Ref.

Model	Set	Pr@0.5	Pr@0.6	Pr@0.7	Pr@0.8	Pr@0.9	IOU
[2]	val	15.25	8.37	3.75	1.29	0.06	28.14
R+LSTM	val	17.57	10.37	4.52	<b>1.26</b>	0.05	28.60
R+RMI	val	<b>20.32</b>	<b>11.62</b>	<b>5.02</b>	1.13	<b>0.09</b>	<b>32.06</b>
R+LSTM+DCRF	val	20.50	14.16	7.86	3.01	0.38	28.94
R+RMI+DCRF	val	<b>23.87</b>	<b>16.12</b>	<b>9.04</b>	<b>3.03</b>	<b>0.42</b>	<b>32.85</b>
D+LSTM	val	25.66	18.23	<b>10.82</b>	4.17	0.64	33.08
D+RMI	val	<b>26.19</b>	<b>18.46</b>	10.68	<b>4.28</b>	<b>0.73</b>	<b>34.40</b>
D+LSTM+DCRF	val	26.71	20.50	13.69	<b>7.11</b>	1.32	33.11
D+RMI+DCRF	val	<b>27.77</b>	<b>21.06</b>	<b>13.92</b>	6.83	<b>1.43</b>	<b>34.52</b>

Table 1.2: Comparison of segmentation performance on UNC.

Model	Set	Pr@0.5	Pr@0.6	Pr@0.7	Pr@0.8	Pr@0.9	IOU
R+LSTM	val	30.98	18.65	8.23	2.12	<b>0.10</b>	38.74
R+RMI	val	<b>32.12</b>	<b>19.49</b>	<b>8.60</b>	<b>2.12</b>	0.03	<b>39.74</b>
R+LSTM+DCRF	val	34.83	24.12	13.97	5.24	<b>0.59</b>	39.88
R+RMI+DCRF	val	<b>36.38</b>	<b>25.36</b>	<b>14.53</b>	<b>5.58</b>	0.56	<b>41.17</b>
D+LSTM	val	39.93	28.19	18.04	<b>7.45</b>	<b>0.91</b>	43.27
D+RMI	val	<b>41.27</b>	<b>29.71</b>	<b>18.41</b>	7.37	0.76	<b>44.33</b>
D+LSTM+DCRF	val	41.42	31.86	22.15	11.90	<b>2.36</b>	43.97
D+RMI+DCRF	val	<b>42.99</b>	<b>33.24</b>	<b>22.75</b>	<b>12.11</b>	2.23	<b>45.18</b>
R+LSTM	testA	30.49	18.63	<b>8.71</b>	1.96	<b>0.09</b>	39.18
R+RMI	testA	<b>31.36</b>	<b>18.95</b>	8.34	<b>2.02</b>	0.04	<b>39.99</b>
R+LSTM+DCRF	testA	34.45	24.62	14.39	5.43	<b>0.50</b>	40.44
R+RMI+DCRF	testA	<b>35.92</b>	<b>24.94</b>	<b>14.50</b>	<b>5.73</b>	0.44	<b>41.35</b>
D+LSTM	testA	39.77	29.61	18.83	8.01	0.88	43.60
D+RMI	testA	<b>40.68</b>	<b>30.14</b>	<b>18.99</b>	<b>8.03</b>	<b>0.88</b>	<b>44.74</b>
D+LSTM+DCRF	testA	41.52	32.61	23.63	12.59	2.19	44.25
D+RMI+DCRF	testA	<b>42.99</b>	<b>33.59</b>	<b>23.69</b>	<b>12.94</b>	<b>2.44</b>	<b>45.69</b>
R+LSTM	testB	32.54	20.80	10.21	2.12	<b>0.16</b>	39.01
R+RMI	testB	<b>35.00</b>	<b>22.30</b>	<b>10.40</b>	<b>2.67</b>	0.12	<b>40.44</b>
R+LSTM+DCRF	testB	35.53	25.28	15.76	6.08	<b>0.65</b>	40.07
R+RMI+DCRF	testB	<b>38.31</b>	<b>27.34</b>	<b>16.96</b>	<b>7.01</b>	0.55	<b>41.87</b>
D+LSTM	testB	40.00	29.03	17.41	7.56	<b>0.92</b>	43.31
D+RMI	testB	<b>42.75</b>	<b>30.40</b>	<b>18.19</b>	<b>7.83</b>	0.86	<b>44.63</b>
D+LSTM+DCRF	testB	41.86	31.66	21.55	11.25	2.30	44.07
D+RMI+DCRF	testB	<b>44.99</b>	<b>34.21</b>	<b>22.69</b>	<b>11.84</b>	<b>2.65</b>	<b>45.57</b>

Table 1.3: Comparison of segmentation performance on UNC+.

Model	Set	Pr@0.5	Pr@0.6	Pr@0.7	Pr@0.8	Pr@0.9	IOU
R+LSTM	val	13.16	6.80	2.40	0.51	0.02	26.25
R+RMI	val	<b>14.49</b>	<b>7.57</b>	<b>2.56</b>	<b>0.59</b>	<b>0.06</b>	<b>27.85</b>
R+LSTM+DCRF	val	16.17	9.76	4.82	1.57	0.11	26.29
R+RMI+DCRF	val	<b>17.88</b>	<b>11.12</b>	<b>5.34</b>	<b>1.62</b>	<b>0.15</b>	<b>28.26</b>
D+LSTM	val	16.58	10.64	5.30	1.78	0.17	28.42
D+RMI	val	<b>18.39</b>	<b>11.50</b>	<b>5.86</b>	<b>1.85</b>	<b>0.20</b>	<b>29.91</b>
D+LSTM+DCRF	val	18.36	12.88	7.85	3.37	0.62	28.07
D+RMI+DCRF	val	<b>20.52</b>	<b>14.02</b>	<b>8.46</b>	<b>3.77</b>	<b>0.62</b>	<b>29.86</b>
R+LSTM	testA	14.01	6.97	2.41	0.40	<b>0.02</b>	26.95
R+RMI	testA	<b>15.33</b>	<b>7.75</b>	<b>2.50</b>	<b>0.40</b>	0.00	<b>28.69</b>
R+LSTM+DCRF	testA	17.22	10.76	5.31	<b>1.82</b>	0.10	27.03
R+RMI+DCRF	testA	<b>19.21</b>	<b>12.10</b>	<b>5.80</b>	1.73	<b>0.10</b>	<b>29.16</b>
D+LSTM	testA	17.31	11.11	5.78	1.59	0.19	28.57
D+RMI	testA	<b>18.76</b>	<b>11.67</b>	<b>6.08</b>	<b>1.78</b>	<b>0.26</b>	<b>30.37</b>
D+LSTM+DCRF	testA	18.97	13.46	8.61	3.63	<b>0.51</b>	28.29
D+RMI+DCRF	testA	<b>21.22</b>	<b>14.43</b>	<b>8.99</b>	<b>3.91</b>	0.49	<b>30.48</b>
R+LSTM	testB	13.15	6.91	2.95	0.51	<b>0.14</b>	24.57
R+RMI	testB	<b>14.15</b>	<b>7.92</b>	<b>3.46</b>	<b>0.82</b>	0.08	<b>26.65</b>
R+LSTM+DCRF	testB	14.99	9.47	5.22	1.76	0.20	24.44
R+RMI+DCRF	testB	<b>16.69</b>	<b>10.62</b>	<b>5.89</b>	<b>2.25</b>	<b>0.27</b>	<b>26.86</b>
D+LSTM	testB	17.26	10.88	6.05	2.25	0.27	27.70
D+RMI	testB	<b>19.08</b>	<b>12.11</b>	<b>6.44</b>	<b>2.70</b>	<b>0.31</b>	<b>29.43</b>
D+LSTM+DCRF	testB	18.41	13.17	8.14	3.95	0.59	27.44
D+RMI+DCRF	testB	<b>20.78</b>	<b>14.56</b>	<b>8.80</b>	<b>4.58</b>	<b>0.80</b>	<b>29.50</b>

Table 1.4: Comparison of segmentation performance on ReferItGame.

Model	Set	Pr@0.5	Pr@0.6	Pr@0.7	Pr@0.8	Pr@0.9	IOU
[1]	test	34.02	26.71	19.32	11.63	3.92	48.03
R+LSTM	test	38.17	29.73	20.91	<b>11.99</b>	<b>3.81</b>	54.01
R+RMI	test	<b>39.44</b>	<b>30.50</b>	<b>21.20</b>	11.95	3.64	<b>54.55</b>
R+LSTM+DCRF	test	40.56	33.40	25.27	16.08	<b>6.09</b>	55.90
R+RMI+DCRF	test	<b>42.21</b>	<b>34.53</b>	<b>25.72</b>	<b>16.09</b>	6.05	<b>56.61</b>
D+LSTM	test	43.86	35.75	26.65	16.75	<b>6.47</b>	56.83
D+RMI	test	<b>44.33</b>	<b>36.13</b>	<b>27.20</b>	<b>16.99</b>	6.43	<b>57.34</b>
D+LSTM+DCRF	test	45.26	38.37	30.20	20.41	<b>8.56</b>	58.20
D+RMI+DCRF	test	<b>46.08</b>	<b>38.90</b>	<b>30.77</b>	<b>20.62</b>	8.54	<b>58.73</b>

## 2. Visualization of Intermediate Segmentation

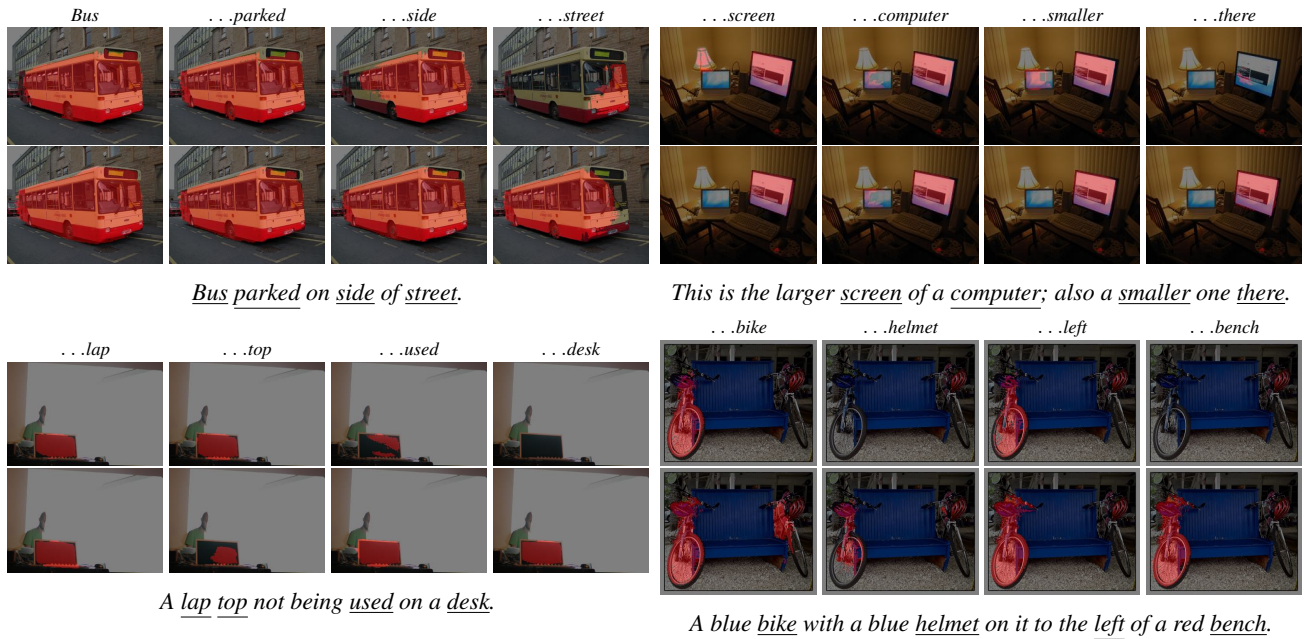


Figure 2.1: Comparison of D+LSTM+DCRF (first row) and D+RMI+DCRF (second row) on Google-Ref. Each column shows segmentation result until after reading the underlined word.



Figure 2.2: Comparison of D+LSTM+DCRF (first row) and D+RMI+DCRF (second row) on UNC. Each column shows segmentation result until after reading the underlined word.

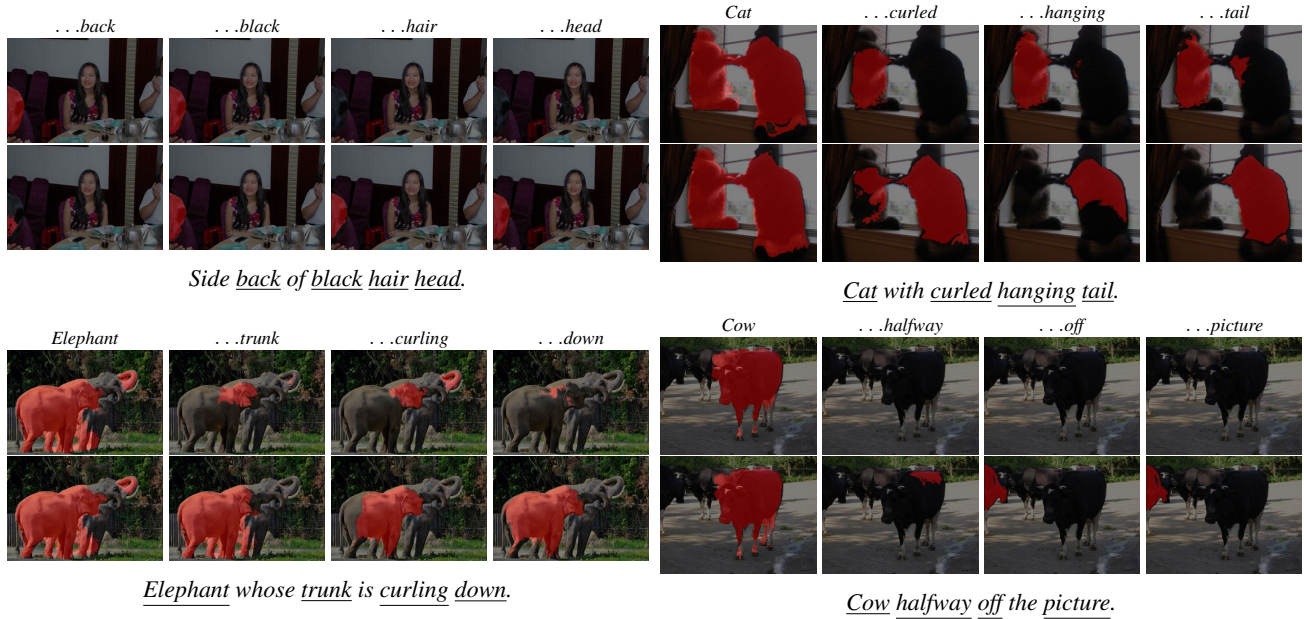


Figure 2.3: Comparison of D+LSTM+DCRF (first row) and D+RMI+DCRF (second row) on UNC+. Each column shows segmentation result until after reading the underlined word.

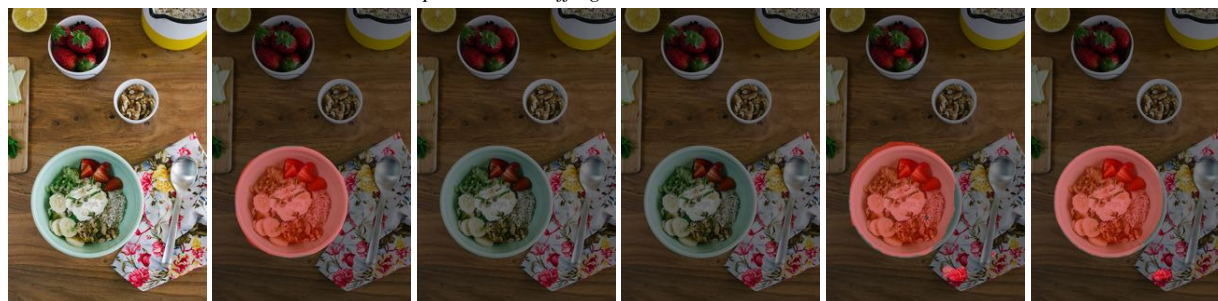


Figure 2.4: Comparison of D+LSTM+DCRF (first row) and D+RMI+DCRF (second row) on ReferItGame. Each column shows segmentation result until after reading the underlined word.

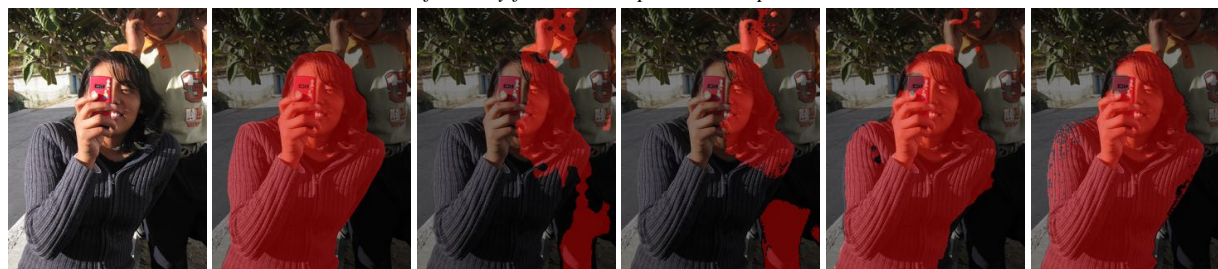
### 3. Visualization of Final Segmentation



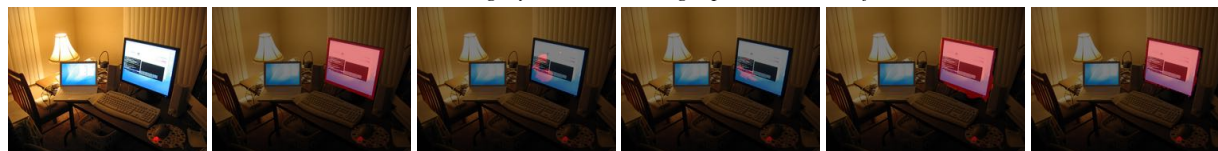
*A polar bear sniffing a slanted concrete slab.*



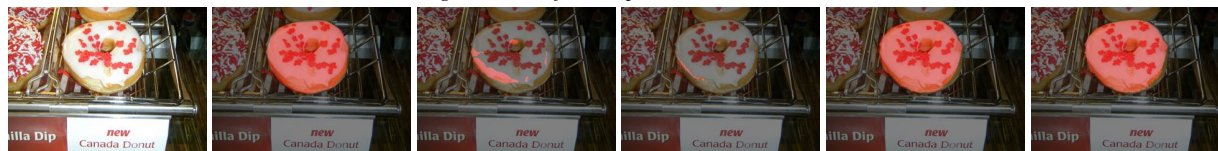
*A bowl of healthy food with a spoon and napkin next to it.*



*A woman with a gray sweater holding a phone near her face.*



*This is the larger screen of a computer; also a smaller one there.*



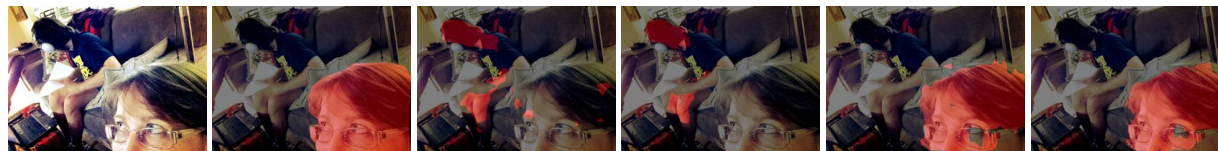
*A new canada donut with white frosting.*



*A green and blue bed.*



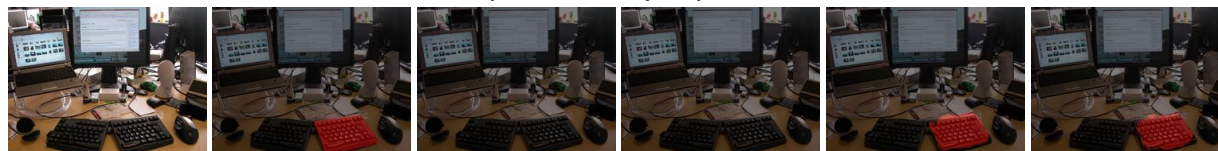
*The bottom two luggage cases being rolled.*



*woman with glasses looking up.*



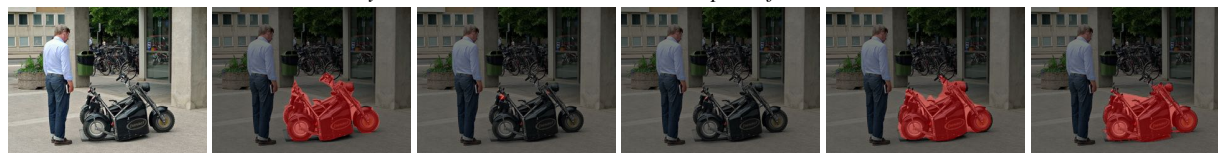
*A yellow and blue fire hydrant.*



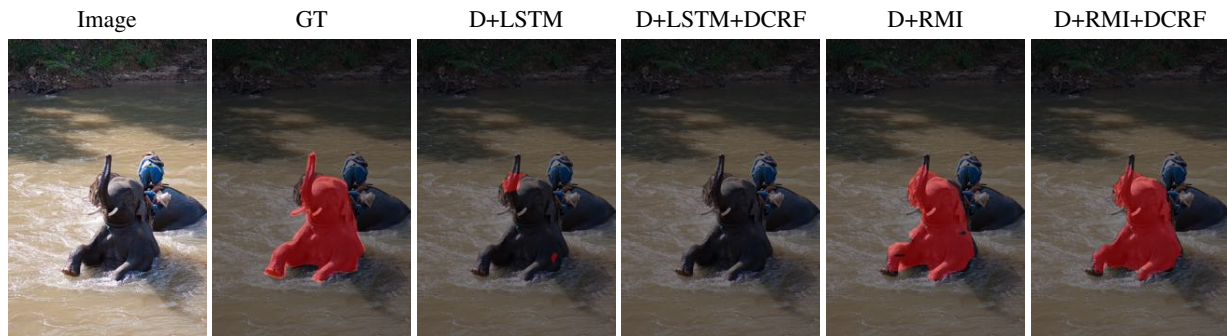
*right hand side of the split keyboard.*



*yellow lemon which is on the lower part of the bark.*



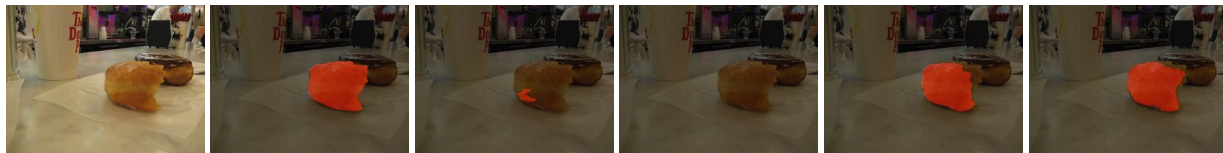
*a weird black bike that sticks to the ground and there are tires for it.*



*The elephant is in the water and has his trunk up while he is being rode by a person.*



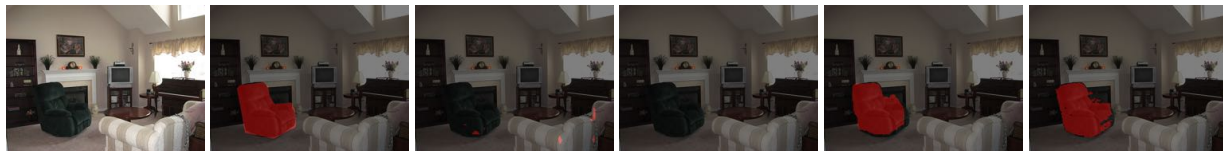
*a woman in a gray sweater standing next to a yellow line.*



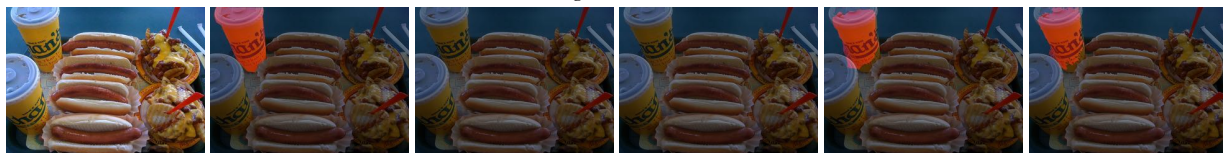
*A half of a glazed donut.*



*The dark pants leg of someone standing to the left of the white sheep carcass.*



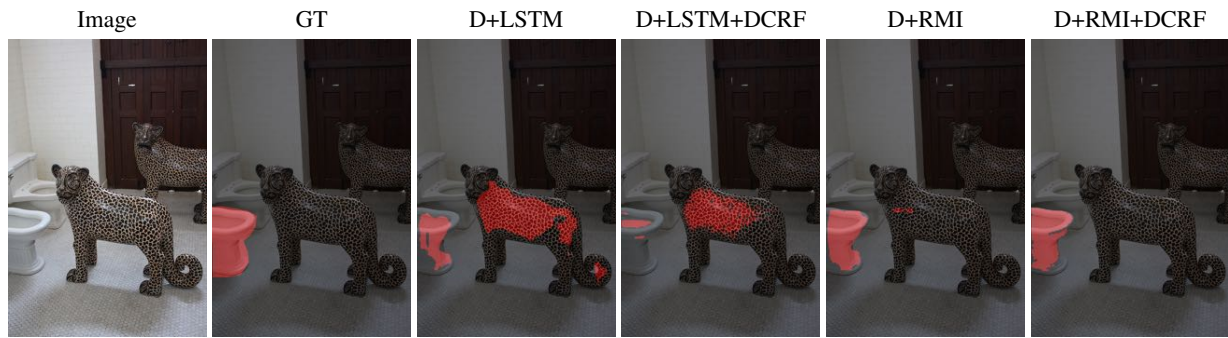
*A dark green arm chair.*



*Yellow drink cup in the back of photo.*

Figure 3.1: Qualitative results of referring image segmentation on Google-Ref.





*first toilet.*



*white truck.*



*man's body above hotdog.*



*person in back in white.*



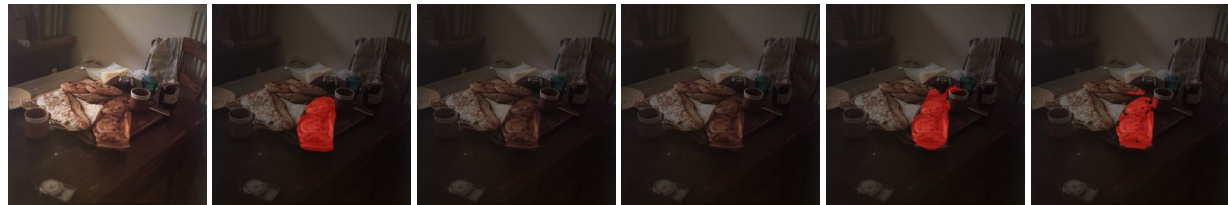
*blk bike.*



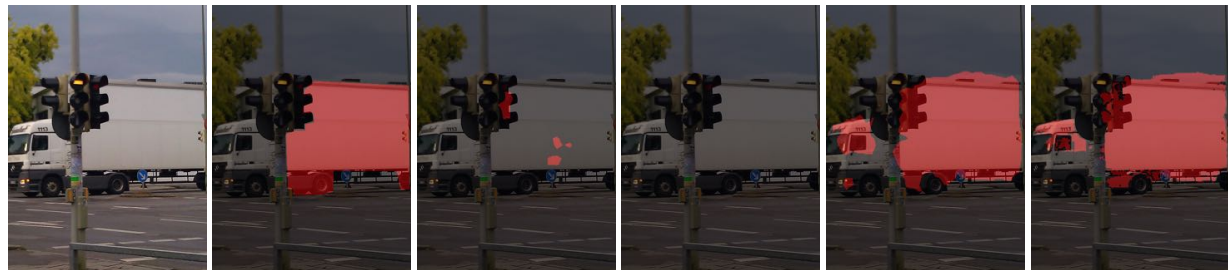
*bear arm top edge of photo.*



*all the way right in front.*



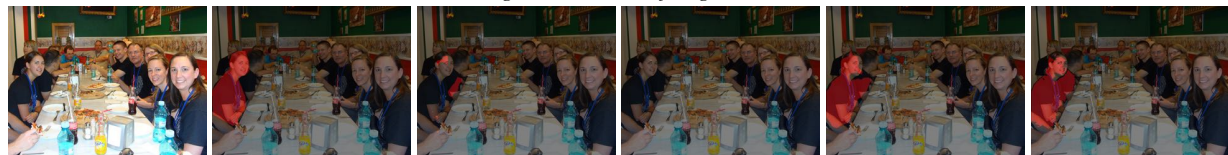
*pastry next to right hand coffee mug.*



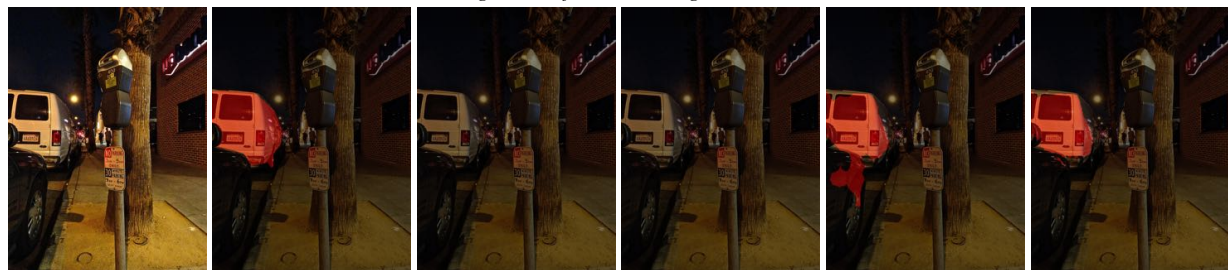
*back part of truck not cab.*



*sitting boarder in the foregrounds.*



*girl on left side looking at us.*



*van.*

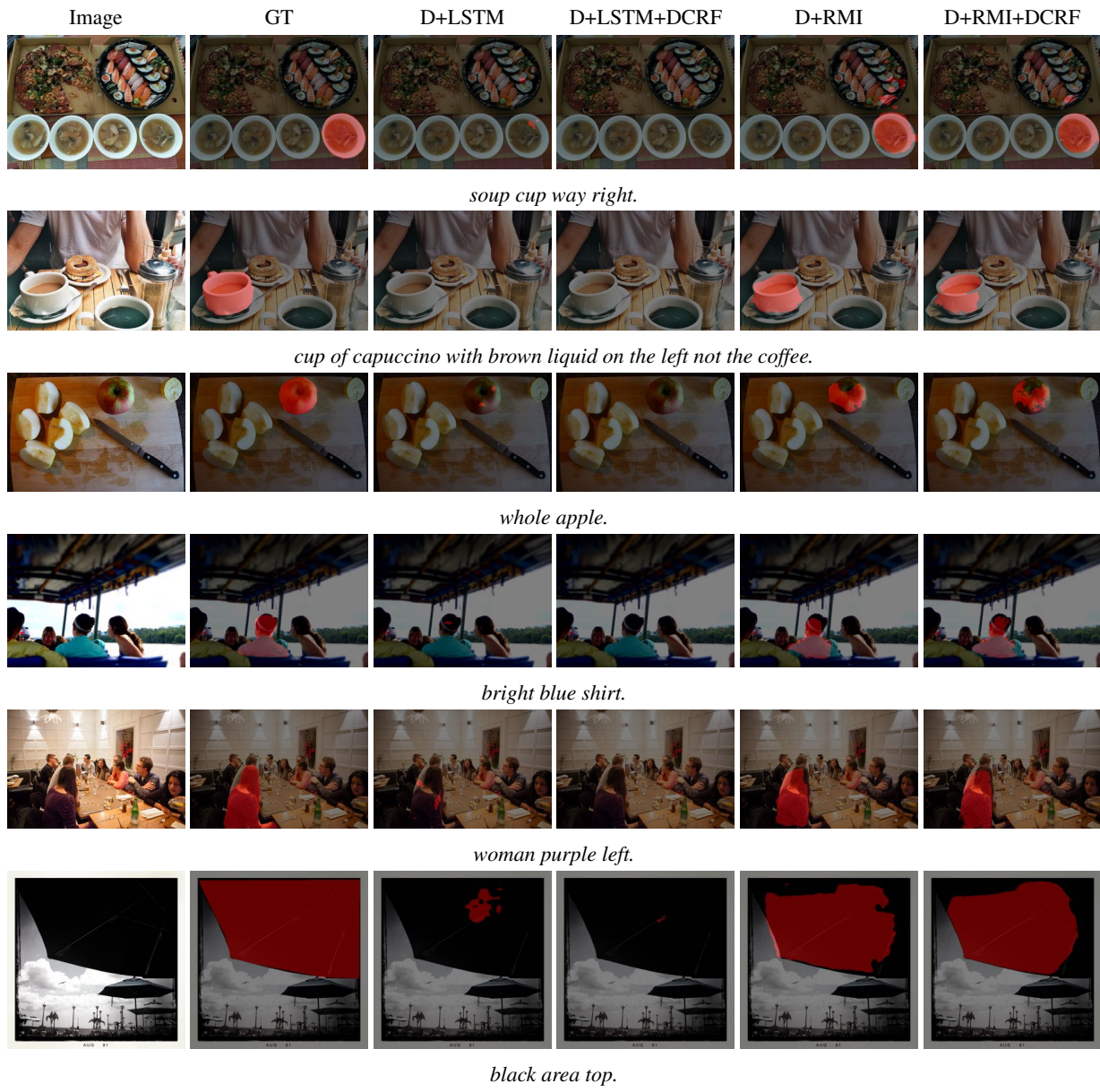


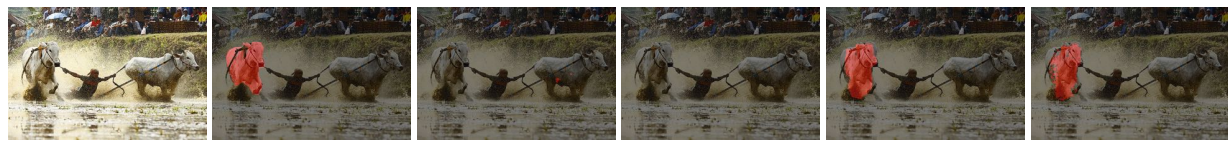
Figure 3.2: Qualitative results of referring image segmentation on UNC.



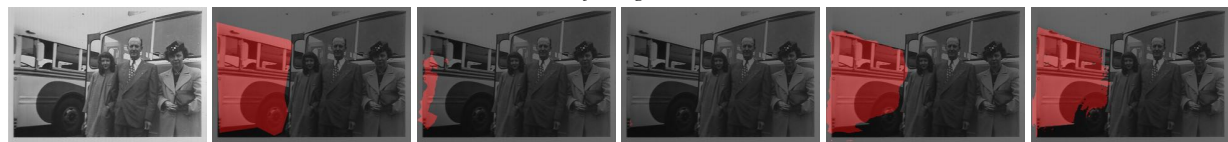
*carrying water.*



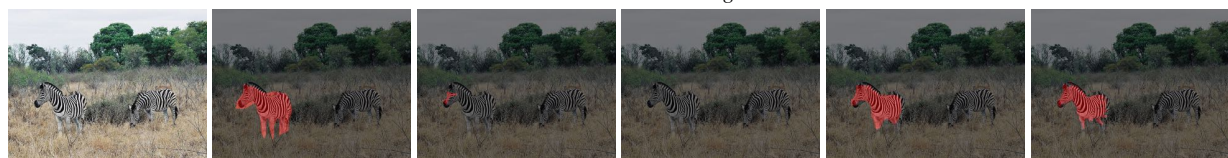
*white car with back window.*



*ox facing us.*



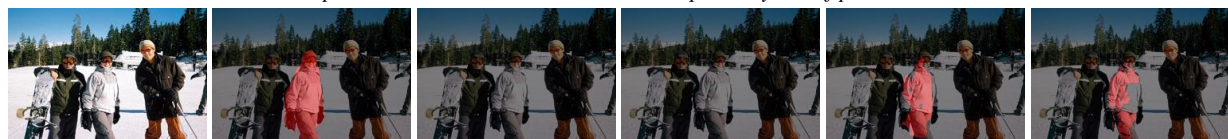
*bus with the tire showing.*



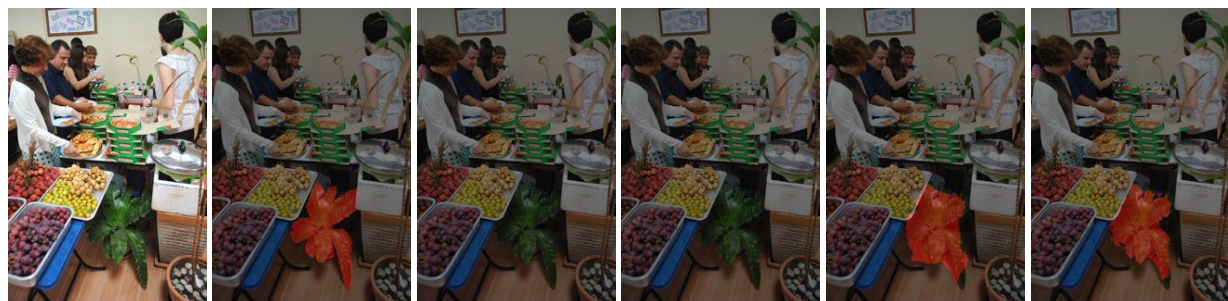
*zebra not eating.*



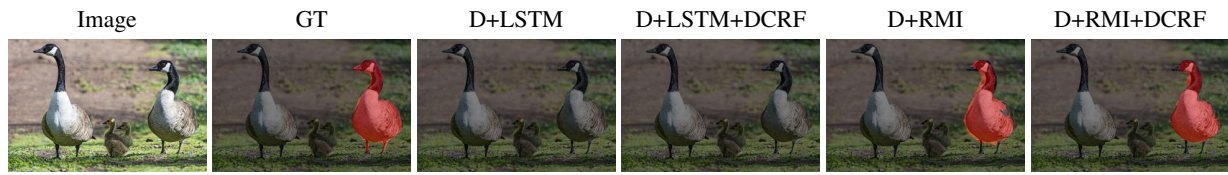
*person in red dress whose umbrella is partially out of picture.*



*between others.*



*green veg close to clock.*



*full size goose with neck bent.*



*half guy.*



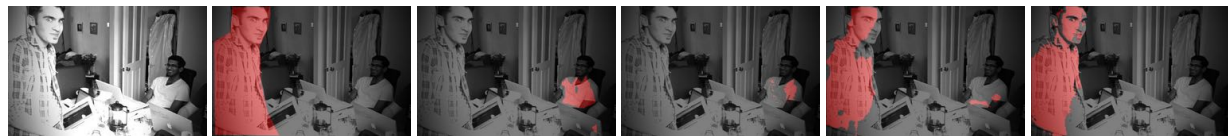
*white umbrella with birds on it.*



*the tv screen by the dude in the hat.*



*tv.*



*creepy eyes.*



*guy with tan hat and fur collar.*

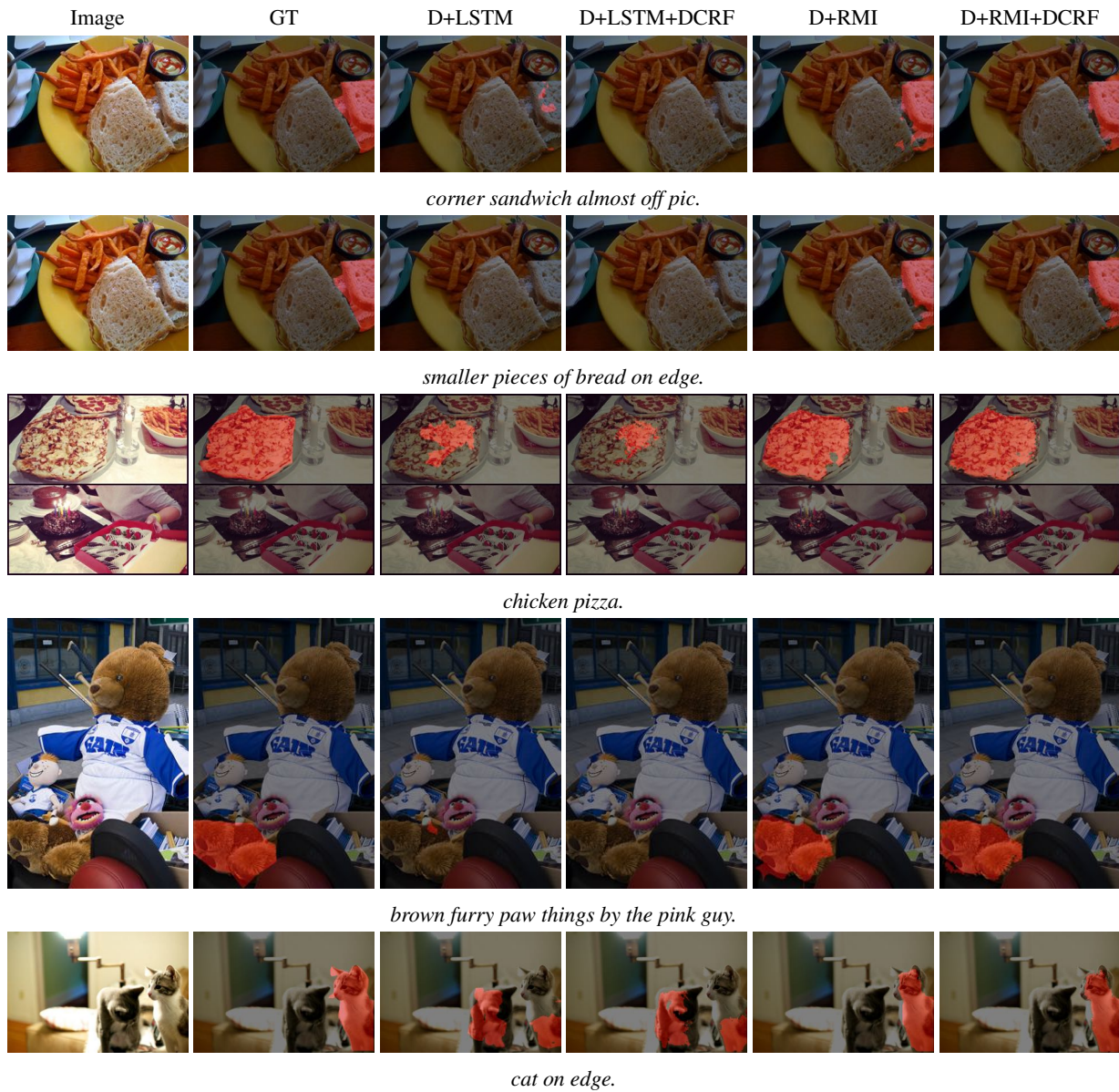


Figure 3.3: Qualitative results of referring image segmentation on UNC+.

Image

GT

D+LSTM

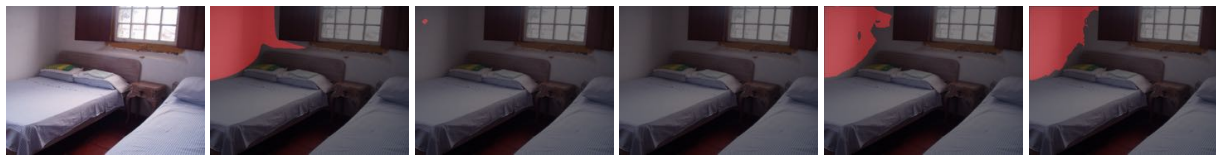
D+LSTM+DCRF

D+RMI

D+RMI+DCRF



*water.*



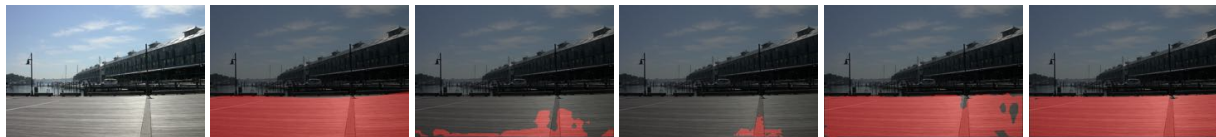
*wall on far left.*



*bush to right of yellow boy.*



*the white van on the left.*



*sidewalk.*



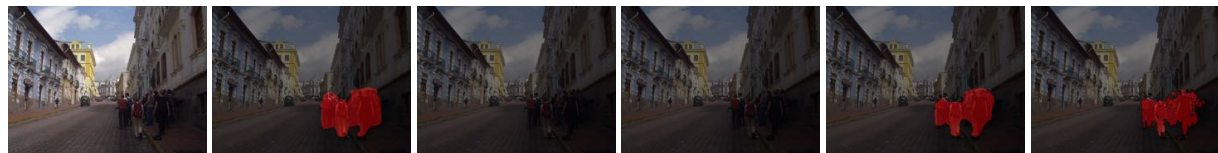
*The painting photo on the wall.*



*the white pillar thing.*



*I would like you to click on the sky please.*



*cluster of ppl.*



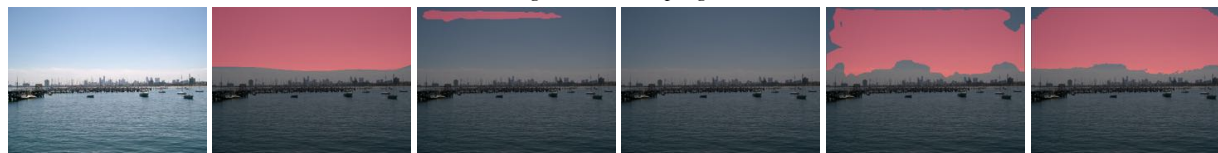
*The sky on the right side of the hut.*



*orange flag on right.*



*car on the right that's creeping in the back.*



*top part of sky.*

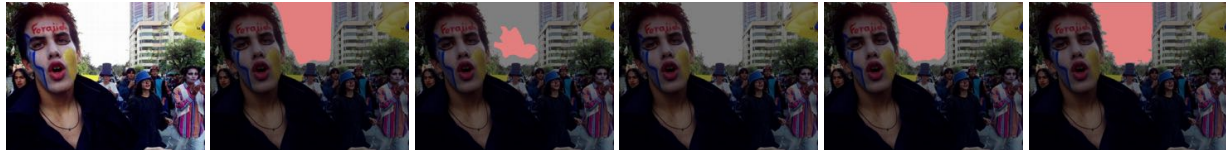


*climber.*

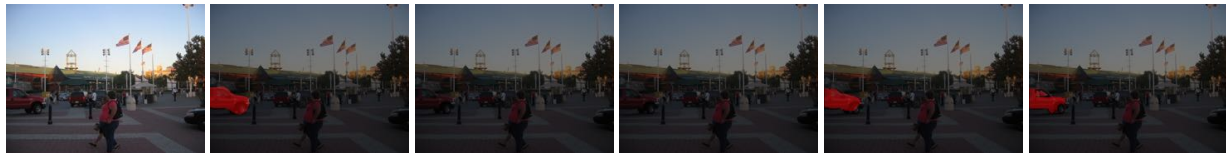




*yes water.*



*the empty sky in the middle.*



*Red truck left.*



*dog.*



*The ground below the wall.*



*white paper on wall on right.*

Figure 3.4: Qualitative results of referring image segmentation on ReferItGame.