

---

000	055
001	056
002	057
003	058
004	059
005	060
006	061
007	062
008	063
009	064
010	065
011	066
012	067
013	068
014	069
015	070
016	071
017	072
018	073
019	074
020	075
021	076
022	077
023	078
024	079
025	080
026	081
027	082
028	083
029	084
030	085
031	086
032	087
033	088
034	089
035	090
036	091
037	092
038	093
039	094
040	095
041	096
042	097
043	098
044	099
045	100
046	101
047	102
048	103
049	104
050	105
051	106
052	107
053	108
054	109

---

## Question Answering about Images using Visual Semantic Embeddings

---

110		165
111		166
112		167
113		168
114		169
115		170
116		171
117		172
118		173
119	Q2: what is grazing in a grassy area above a body of water ?	174
120	Ground truth: sheep	175
121	2-CNN-LSTM: <b>sheep</b> (0.6904)	176
122	CNN-BOW: <b>sheep</b> (0.8393)	177
123	BLIND-LSTM: <b>zebra</b> (0.3476)	178
124	BLIND-BOW: <b>giraffe</b> (0.1673)	179
125	Image features are useful	180
126		181
127		182
128		183
129		184
130		185
131		186
132		187
133		188
134		189
135	Q1213: what is wrapped in packaging paper for shipping ?	190
136	Ground truth: toilet	191
137	2-CNN-LSTM: <b>toilet</b> (0.3934)	192
138	CNN-BOW: <b>toilet</b> (0.3243)	193
139	BLIND-LSTM: <b>sandwich</b> (0.2178)	194
140	BLIND-BOW: <b>sandwich</b> (0.2595)	195
141	Blind model has no clue from the text alone	196
142		197
143		198
144		199
145		200
146		201
147		202
148		203
149		204
150		205
151		206
152	Q4974: what is the person riding down the street ?	207
153	Ground truth: skateboard	208
154	2-CNN-LSTM: <b>skateboard</b> (0.8822)	209
155	CNN-BOW: <b>skateboard</b> (0.4611)	210
156	BLIND-LSTM: <b>motorcycle</b> (0.3368)	211
157	BLIND-BOW: <b>skateboard</b> (0.2380)	212
158	Reletively small object to focus on	213
159		214
160		215
161		216
162		217
163		218
164		219
110		165
111		166
112		167
113		168
114		169
115		170
116		171
117		172
118		173
119	Q85: what is standing in the sand alone ?	174
120	Ground truth: bird	175
121	2-CNN-LSTM: <b>bird</b> (0.6608)	176
122	CNN-BOW: <b>seagull</b> (0.3627)	177
123	BLIND-LSTM: <b>elephant</b> (0.2465)	178
124	BLIND-BOW: <b>elephant</b> (0.1757)	179
125	Image features are useful	180
126		181
127		182
128		183
129		184
130		185
131		186
132		187
133		188
134		189
135	Q223: what is the man in a helmet eating ?	190
136	Ground truth: banana	191
137	2-CNN-LSTM: <b>banana</b> (0.8961)	192
138	CNN-BOW: <b>banana</b> (0.9729)	193
139	BLIND-LSTM: <b>sandwich</b> (0.2065)	194
140	BLIND-BOW: <b>sandwich</b> (0.1627)	195
141	Occluded object are still focusable	196
142		197
143		198
144		199
145		200
146		201
147		202
148		203
149		204
150		205
151		206
152	Q1549: what are sitting on wood UNK one is cut in half ?	207
153	Ground truth: oranges	208
154	2-CNN-LSTM: <b>oranges</b> (0.8196)	209
155	CNN-BOW: <b>oranges</b> (0.3242)	210
156	BLIND-LSTM: <b>pizzas</b> (0.1006)	211
157	BLIND-BOW: <b>plate</b> (0.0424)	212
158	Both models know plural for	213
159		214
160		215
161		216
162		217
163		218
164		219
110		165
111		166
112		167
113		168
114		169
115		170
116		171
117		172
118		173
119	Q2007: what is sitting next to a cup of coffee ?	174
120	Ground truth: banana	175
121	2-CNN-LSTM: <b>banana</b> (0.4267)	176
122	CNN-BOW: <b>banana</b> (0.6397)	177
123	BLIND-LSTM: <b>sandwich</b> (0.1032)	178
124	BLIND-BOW: <b>sandwich</b> (0.1012)	179
125	When there are two objects present the model attend to the correct one	180
126		181
127		182
128		183
129		184
130		185
131		186
132		187
133		188
134		189
135		190
136		191
137		192
138		193
139		194
140		195
141		196
142		197
143		198
144		199
145		200
146		201
147		202
148		203
149		204
150		205
151		206
152	Q5429: what do two women hold with a picture on it ?	207
153	Ground truth: cake	208
154	2-CNN-LSTM: <b>cake</b> (0.5611)	209
155	CNN-BOW: <b>laptop</b> (0.1443)	210
156	BLIND-LSTM: <b>umbrellas</b> (0.1567)	211
157	BLIND-BOW: <b>phones</b> (0.1447)	212
158	Does not look like a regular cake and the model has to infer from the candles	213
159		214
160		215
161		216
162		217
163		218
164		219
110		165
111		166
112		167
113		168
114		169
115		170
116		171
117		172
118		173
119	Q1965: what did people set up on the beach next to blankets ?	174
120	Ground truth: umbrellas	175
121	2-CNN-LSTM: <b>umbrellas</b> (0.4768)	176
122	CNN-BOW: <b>umbrella</b> (0.7319)	177
123	BLIND-LSTM: <b>kites</b> (0.1585)	178
124	BLIND-BOW: <b>boat</b> (0.1258)	179
125	The CNN-LSTM model knows the plural but CNN-BOW does not	180
126		181
127		182
128		183
129		184
130		185
131		186
132		187
133		188
134		189
135		190
136		191
137		192
138		193
139		194
140		195
141		196
142		197
143		198
144		199
145		200
146		201
147		202
148		203
149		204
150		205
151		206
152		207
153		208
154		209
155		210
156		211
157		212
158		213
159		214
160		215
161		216
162		217
163		218
164		219

*Figure 1. COCO-QA Object Questions*

## Question Answering about Images using Visual Semantic Embeddings

220			275
221			276
222			277
223			278
224			279
225			280
226			281
227			282
228			283
229	Q1966: what are the man and woman holding up ? Ground truth: pizza 2-CNN-LSTM: <b>pizza (0.0775)</b> CNN-BOW: <b>pizzas (0.1807)</b> BLIND-LSTM: <b>phones (0.1421)</b> BLIND-BOW: <b>phones (0.2272)</b> The CNN-BOW again messed up the counting (mostly of the questions with LSTM win over BOW is on singular/plural)		284
230			285
231			286
232			287
233			288
234			289
235			290
236			291
237			292
238			293
239			294
240			295
241			296
242			297
243			298
244			299
245			300
246			301
247	Q15178: what are two guys riding down the street ? Ground truth: motorcycle 2-CNN-LSTM: <b>motorcycle (0.6804)</b> CNN-BOW: <b>motorcycles (0.4514)</b> BLIND-LSTM: <b>motorcycles (0.2584)</b> BLIND-BOW: <b>bikes (0.2329)</b> From the text alone it is misleading that there might exist multiple motorcycles		302
248			303
249			304
250			305
251			306
252			307
253			308
254			309
255			310
256			311
257			312
258			313
259			314
260			315
261			316
262			317
263			318
264			319
265	Q29528: what is there filled with sauce covered food ? Ground truth: tray 2-CNN-LSTM: <b>tray (0.1637)</b> CNN-BOW: <b>pizza (0.1872)</b> BLIND-LSTM: <b>plate (0.2158)</b> BLIND-BOW: <b>plate (0.2363)</b> Tray is hard to get		320
266			321
267			322
268			323
269			324
270			325
271			326
272			327
273			328
274			329
275	Q2950: what is the color of the man ? Ground truth: black 2-CNN-LSTM: <b>black (0.7063)</b> CNN-BOW: <b>black (0.2809)</b> BLIND-LSTM: <b>black (0.4410)</b> BLIND-BOW: <b>black (0.4541)</b> LSTM learns to add plural form from the verb tense		284
276			285
277			286
278			287
279			288
280			289
281			290
282			291
283			292
284	Q14930: what next to some UNK UNK a lime and a mixed drink ? Ground truth: blender 2-CNN-LSTM: <b>blender (0.3763)</b> CNN-BOW: <b>plate (0.0927)</b> BLIND-LSTM: <b>standing (0.1304)</b> BLIND-BOW: <b>plate (0.0410)</b> Many items present but the LSTM figures out the correct item		293
285			294
286			295
287			296
288			297
289			298
290			299
291			300
292			301
293			302
294	Q19533: what lays in the bed resting it 's head against a pillow ? Ground truth: bear 2-CNN-LSTM: <b>bear (0.3648)</b> CNN-BOW: <b>bed (0.3379)</b> BLIND-LSTM: <b>dog (0.5108)</b> BLIND-BOW: <b>dog (0.6639)</b> BOW is not a good blend of visual information with the question semantics		303
295			304
296			305
297			306
298			307
299			308
300			309
301			310
302	Q24952: what is the black and white cat wearing ? Ground truth: hat 2-CNN-LSTM: <b>hat (0.6349)</b> CNN-BOW: <b>cat (0.6880)</b> BLIND-LSTM: <b>tie (0.5821)</b> BLIND-BOW: <b>tie (0.3385)</b> BOW is not a good blend of visual information with the question semantics		311
303			312
304			313
305			314
306			315
307			316
308			317
309			318
310			319
311	Q30568: what stuck in an apple ? Ground truth: knife 2-CNN-LSTM: <b>knife (0.2179)</b> CNN-BOW: <b>apple (0.2666)</b> BLIND-LSTM: <b>bananas (0.0894)</b> BLIND-BOW: <b>apple (0.0838)</b> It is rare to see something stuck into something and here CNN-LSTM win over other 4 models shows a good combination of image and question semantics		320
312			321
313			322
314			323
315			324
316			325
317			326
318			327
319	Q36019: the woman wearing what UNK with a little girl as she played with her stroller UNK which also has an umbrella over it ? Ground truth: hat 2-CNN-LSTM: <b>hat (0.4144)</b> CNN-BOW: <b>UNK (0.0371)</b> BLIND-LSTM: <b>shirt (0.0738)</b> BLIND-BOW: <b>plate (0.0401)</b> LSTM remembers the beginning part of the question although the entire question is very long		328
320			329
321			329
322			329
323			329
324			329
325			329
326			329
327			329
328			329
329			329

Figure 2. COCO-QA Object Questions

## Question Answering about Images using Visual Semantic Embeddings

---

330				385
331				386
332				387
333				388
334				389
335				390
336				391
337				392
338				393
339	Q38613: what are the people watching and taking photos ?			394
340	Ground truth: elephants			395
341	2-CNN-LSTM: <b>elephants (0.4016)</b>			396
342	CNN-BOW: <b>horses (0.7849)</b>			397
343	BLIND-LSTM: <b>picture (0.1402)</b>			398
344	BLIND-BOW: <b>picture (0.2065)</b>			399
345	The object is fairly small to focus on			400
346				401
347				402
348				403
349				404
350				405
351				406
352				407
353				408
354				409
355				410
356				411
357	Q14518: what will get the lot of use by its owner ?			412
358	Ground truth: mouse			413
359	2-CNN-LSTM: <b>toothbrushes (0.1746)</b>			414
360	CNN-BOW: <b>mouse (0.2407)</b>			415
361	BLIND-LSTM: <b>dogs (0.1367)</b>			416
362	BLIND-BOW: <b>elephant (0.0866)</b>			417
363	CNN-LSTM fails			418
364				419
365				420
366				421
367				422
368				423
369				424
370				425
371				426
372				427
373				428
374				429
375	Q385: two ladies enjoying what at a roadside stand ?			430
376	Ground truth: dessert			431
377	2-CNN-LSTM: <b>bananas (0.4363)</b>			432
378	CNN-BOW: <b>banana (0.9090)</b>			433
379	BLIND-LSTM: <b>fruit (0.0820)</b>			434
380	BLIND-BOW: <b>meal (0.1226)</b>			435
381	Yellow food doesn't have to be bananas			436
382				437
383				438
384				439

*Figure 3. COCO-QA Object Questions*

440	495
441	496
442	497
443	498
444	499
445	500
446	501
447	502
448	503
449	504
450	505
451	506
452	507
453	508
454	509
455	510
456	511
457	512
458	513
459	514
460	515
461	516
462	517
463	518
464	519
465	520
466	521
467	522
468	523
469	524
470	525
471	526
472	527
473	528
474	529
475	530
476	531
477	532
478	533
479	534
480	535
481	536
482	537
483	538
484	539
485	540
486	541
487	542
488	543
489	544
490	545
491	546
492	547
493	548
494	549



Q828: what is towing the heavy load across town ?  
 Ground truth: truck  
 2-CNN-LSTM: train (0.8407)  
 CNN-BOW: train (0.5635)  
 BLIND-LSTM: truck (0.6577)  
 BLIND-BOW: truck (0.4229)  
 Very long truck looks like a train?



Q33: what is shining down on the passing airplane ?  
 Ground truth: sun  
 2-CNN-LSTM: airplane (0.4153)  
 CNN-BOW: airplane (0.2487)  
 BLIND-LSTM: sun (0.3013)  
 BLIND-BOW: car (0.1130)  
 Image provides no extra information and blind LSTM really focuses on the text to get the right answer



Q7231: what lined up along one wall with a handicap UNK one at the end ?  
 Ground truth: urinals  
 2-CNN-LSTM: suitcases (0.2665)  
 CNN-BOW: luggage (0.0518)  
 BLIND-LSTM: urinals (0.0969)  
 BLIND-BOW: plate (0.0401)  
 Interesting that the blind model knows the correct answer whereas the CNN model thinks they are suitcases



Q24432: what does the colorful market booth sell to a customer ?  
 Ground truth: bananas  
 2-CNN-LSTM: box (0.1577)  
 CNN-BOW: cake (0.1204)  
 BLIND-LSTM: bananas (0.0727)  
 BLIND-BOW: fruit (0.2503)  
 Very hard to answer bananas but the blind model made it

Figure 4. COCO-QA Object Questions

## Question Answering about Images using Visual Semantic Embeddings

550			605
551			606
552			607
553			608
554			609
555			610
556			611
557			612
558			613
559	Q3202: how many red double decker buses parked side by side ? Ground truth: three 2-CNN-LSTM: <b>three (0.3344)</b> CNN-BOW: <b>two (0.5916)</b> BLIND-LSTM: <b>two (0.3844)</b> BLIND-BOW: <b>two (0.4718)</b> Couting regarding to an object with a specific color		614
560			615
561			616
562			617
563			618
564			619
565			620
566			621
567			622
568			623
569			624
570			625
571			626
572			627
573			628
574			629
575			630
576	Q7270: how many guys in a room is playing with a wii ? Ground truth: four 2-CNN-LSTM: <b>four (0.4461)</b> CNN-BOW: <b>five (0.3932)</b> BLIND-LSTM: <b>two (0.4457)</b> BLIND-BOW: <b>two (0.6535)</b> Four again		631
577			632
578			633
579			634
580			635
581			636
582			637
583			638
584			639
585			640
586			641
587			642
588			643
589			644
590			645
591			646
592			647
593			648
594	Q32912: how many bottles of beer is sitting on a wooden table next to bananas ? Ground truth: three 2-CNN-LSTM: <b>three (0.3403)</b> CNN-BOW: <b>two (0.4664)</b> BLIND-LSTM: <b>two (0.3781)</b> BLIND-BOW: <b>four (0.3717)</b> Multiple objects but asking about bottles		649
595			650
596			651
597			652
598			653
599			654
600			655
601			656
602			657
603			658
604			659
550			
551			
552			
553			
554			
555			
556			
557			
558			
559			
560			
561			
562			
563			
564			
565			
566			
567			
568			
569			
570			
571			
572			
573			
574			
575			
576			
577			
578			
579			
580			
581			
582			
583			
584			
585			
586			
587			
588			
589			
590			
591			
592			
593			
594			
595			
596			
597			
598			
599			
600			
601			
602			
603			
604			
550			
551			
552			
553			
554			
555			
556			
557			
558			
559			
560			
561			
562			
563			
564			
565			
566			
567			
568			
569			
570			
571			
572			
573			
574			
575			
576			
577			
578			
579			
580			
581			
582			
583			
584			
585			
586			
587			
588			
589			
590			
591			
592			
593			
594			
595			
596			
597			
598			
599			
600			
601			
602			
603			
604			
550			
551			
552			
553			
554			
555			
556			
557			
558			
559			
560			
561			
562			
563			
564			
565			
566			
567			
568			
569			
570			
571			
572			
573			
574			
575			
576			
577			
578			
579			
580			
581			
582			
583			
584			
585			
586			
587			
588			
589			
590			
591			
592			
593			
594			
595			
596			
597			
598			
599			
600			
601			
602			
603			
604			
550			
551			
552			
553			
554			
555			
556			
557			
558			
559			
560			
561			
562			
563			
564			
565			
566			
567			
568			
569			
570			
571			
572			
573			
574			
575			
576			
577			
578			
579			
580			
581			
582			
583			
584			
585			
586			
587			
588			
589			
590			
591			
592			
593			
594			
595			
596			
597			
598			
599			
600			
601			
602			
603			
604			
550			
551			
552			
553			
554			
555			
556			
557			
558			
559			
560			
561			
562			
563			
564			
565			
566			
567			
568			
569			
570			
571			
572			
573			
574			
575			
576			
577			
578			
579			
580			
581			
582			
583			
584			
585			
586			
587			
588			
589			
590			
591			
592			
593			
594			
595			
596			
597			
598			
599			
600			
601			
602			
603			
604			

Figure 5. COCO-QA Number Questions

## Question Answering about Images using Visual Semantic Embeddings

660				715
661				716
662				717
663				718
664				719
665				720
666				721
667				722
668				723
669	Q9489:	how many benches surrounded by leaves near the pond ?		724
670	Ground truth:	two		725
671	2-CNN-LSTM:	four (0.3651)		726
672	CNN-BOW:	two (0.3351)		727
673	BLIND-LSTM:	three (0.3202)		728
674	BLIND-BOW:	four (0.3390)		729
675	Seems to have no reason to answer "four"			730
676				731
677				732
678				733
679				734
680				735
681				736
682				737
683				738
684				739
685				740
686				741
687				742
688	Q24263:	how many ducks are standing on the beach with water in the foreground ?		743
689	Ground truth:	five		744
690	2-CNN-LSTM:	five (0.3164)		745
691	CNN-BOW:	five (0.4691)		746
692	BLIND-LSTM:	three (0.3322)		747
693	BLIND-BOW:	four (0.2999)		748
694	Again five is not easy to count for both models			749
695				750
696				751
697				752
698				753
699				754
700				755
701				756
702				757
703				758
704				759
705	Q25671:	how many slice is there of pizza left on the plate ?		760
706	Ground truth:	one		761
707	2-CNN-LSTM:	one (0.9581)		762
708	CNN-BOW:	one (0.4745)		763
709	BLIND-LSTM:	one (0.9584)		764
710	BLIND-BOW:	one (0.8934)		765
711	Answer is revealed in the word "is"			766
712				767
713				768
714				769
38				
39				
40				
41				
42				
43		© Ryan Taylor		
44				
45				
46				
39	Q170:	how many black laptop computers in UNK on a white table ?		
40	Ground truth:	three		
41	2-CNN-LSTM:	three (0.3208)		
42	CNN-BOW:	three (0.2394)		
43	BLIND-LSTM:	two (0.4606)		
44	BLIND-BOW:	plate (0.0363)		
45	Both CNN models agree on counting			
46	Q3449:	how many men cut the cake at an event ?		
47	Ground truth:	five		
48	2-CNN-LSTM:	five (0.3079)		
49	CNN-BOW:	five (0.4365)		
50	BLIND-LSTM:	three (0.4024)		
51	BLIND-BOW:	four (0.3052)		
52	This is not easy to get firstly because "five" is fairly large number to count and secondly because cutting cakes really doesn't imply five persons			
41	Q24263:	how many ducks are standing on the beach with water in the foreground ?		
42	Ground truth:	one		
43	2-CNN-LSTM:	one (0.8434)		
44	CNN-BOW:	one (0.7174)		
45	BLIND-LSTM:	one (0.8946)		
46	BLIND-BOW:	one (0.7984)		
47	Usually signs show "one way"			
41	Q6:	how many way sign that is in front of a building ?		
42	Ground truth:	one		
43	2-CNN-LSTM:	one (0.4889)		
44	CNN-BOW:	two (0.6430)		
45	BLIND-LSTM:	two (0.4091)		
46	BLIND-BOW:	two (0.5368)		
47	Answer is revealed in the phrase "each other"			
41	Q25552:	how many trains are parked next to each other on the tracks ?		
42	Ground truth:	two		
43	2-CNN-LSTM:	two (0.4889)		
44	CNN-BOW:	two (0.6430)		
45	BLIND-LSTM:	two (0.4091)		
46	BLIND-BOW:	two (0.5368)		
47	Answer is revealed in the phrase "each other"			
41	Q25671:	how many slice is there of pizza left on the plate ?		
42	Ground truth:	one		
43	2-CNN-LSTM:	one (0.9581)		
44	CNN-BOW:	one (0.4745)		
45	BLIND-LSTM:	one (0.9584)		
46	BLIND-BOW:	one (0.8934)		
47	Answer is revealed in the word "is"			
41	Q24776:	how many small dog is running next to four sheep ?		
42	Ground truth:	one		
43	2-CNN-LSTM:	one (0.6354)		
44	CNN-BOW:	three (0.2432)		
45	BLIND-LSTM:	one (0.8920)		
46	BLIND-BOW:	two (0.2721)		
47	LSTM captures "is" whereas BOW doesn't			
41	Q811:	how many apples is sitting in a bowl on a table ?		
42	Ground truth:	four		
43	2-CNN-LSTM:	three (0.3703)		
44	CNN-BOW:	two (0.4607)		
45	BLIND-LSTM:	two (0.2551)		
46	BLIND-BOW:	four (0.2847)		
47	Easy one but only blind BOW gets it			

Figure 6. COCO-QA Number Questions

Q33482: how many surfers on a beach with their surfboards ?

Ground truth: four

2-CNN-LSTM: two (0.4726)

CNN-BOW: two (0.5776)

BLIND-LSTM: three (0.3930)

BLIND-BOW: four (0.3283)

Again blind BOW tend to guess “four” when it has no clue and other models just don’t get the correct answer

Q12255: how many engine plane is flying in lightly overcast sky ?

Ground truth: four

2-CNN-LSTM: four (0.4671)

CNN-BOW: nine (0.2623)

BLIND-LSTM: one (0.3137)

BLIND-BOW: four (0.5412)

This one is impossible to get

Q38217: how many blue jets with yellow painted on top are flying together ?

Ground truth: four

2-CNN-LSTM: four (0.3229)

CNN-BOW: six (0.1610)

BLIND-LSTM: three (0.2661)

BLIND-BOW: four (0.4019)

Objects are overlapped probably the model just guesses four planes regardless

*Figure 7. COCO-QA Number Questions*

## Question Answering about Images using Visual Semantic Embeddings

880				935	
881				936	
882				937	
883				938	
884				939	
885				940	
886				941	
887				942	
888				943	
889				944	
890		Q521: what is the color of the umbrellas ? Ground truth: yellow 2-CNN-LSTM: <b>yellow (0.9376)</b> CNN-BOW: <b>yellow (0.3751)</b> BLIND-LSTM: <b>red (0.2789)</b> BLIND-BOW: <b>blue (0.2287)</b> CNN tend to focus on bright colors even if they only occupy small portion of the image	Q1362: what is the color of the flowers ? Ground truth: red 2-CNN-LSTM: <b>red (0.4641)</b> CNN-BOW: <b>red (0.4694)</b> BLIND-LSTM: <b>purple (0.6047)</b> BLIND-BOW: <b>purple (0.5514)</b> The blue bench is more obvious to detect but the question is asking about flowers	Q2023: what is the color of the dog ? Ground truth: white 2-CNN-LSTM: <b>white (0.6819)</b> CNN-BOW: <b>white (0.9388)</b> BLIND-LSTM: <b>brown (0.5489)</b> BLIND-BOW: <b>brown (0.4795)</b> Focus on the animal not the pants	945
891				946	
892				947	
893				948	
894				949	
895				950	
896				951	
897				952	
898				953	
899				954	
900				955	
901				956	
902				957	
903				958	
904				959	
905				960	
906				961	
907		Q2683: what is the color of the table ? Ground truth: black 2-CNN-LSTM: <b>black (0.7793)</b> CNN-BOW: <b>black (0.3785)</b> BLIND-LSTM: <b>brown (0.2596)</b> BLIND-BOW: <b>white (0.3189)</b> Attend to background of the image	Q5520: what is the color of the banana ? Ground truth: green 2-CNN-LSTM: <b>green (0.4476)</b> CNN-BOW: <b>green (0.7215)</b> BLIND-LSTM: <b>yellow (0.3524)</b> BLIND-BOW: <b>yellow (0.3629)</b> Banana question 1	Q9822: what is the color of the bananas ? Ground truth: yellow 2-CNN-LSTM: <b>yellow (0.6238)</b> CNN-BOW: <b>yellow (0.8101)</b> BLIND-LSTM: <b>green (0.3822)</b> BLIND-BOW: <b>green (0.3820)</b> Banana question 2	962
908				963	
909				964	
910				965	
911				966	
912				967	
913				968	
914				969	
915				970	
916				971	
917				972	
918				973	
919				974	
920				975	
921				976	
922				977	
923				978	
924		Q34206: what is the color of the bowl ? Ground truth: red 2-CNN-LSTM: <b>red (0.4605)</b> CNN-BOW: <b>red (0.3050)</b> BLIND-LSTM: <b>white (0.1952)</b> BLIND-BOW: <b>white (0.4396)</b> Hard to look for the bowl	Q139: what is the color of the bear ? Ground truth: black 2-CNN-LSTM: <b>black (0.8752)</b> CNN-BOW: <b>green (0.3879)</b> BLIND-LSTM: <b>brown (0.5518)</b> BLIND-BOW: <b>brown (0.6222)</b> BOW now mistakenly output the dominant color whereas the blind models give a good guess	Q1404: what is the color of the tooth-brush ? Ground truth: green 2-CNN-LSTM: <b>green (0.4124)</b> CNN-BOW: <b>blue (0.3226)</b> BLIND-LSTM: <b>blue (0.1879)</b> BLIND-BOW: <b>blue (0.2653)</b> LSTM is more attentive to details of both question and image	979
925				980	
926				981	
927				982	
928				983	
929				984	
930				985	
931				986	
932				987	
933				988	
934				989	

*Figure 8. COCO-QA Color Questions*

990				1045	
991				1046	
992				1047	
993				1048	
994				1049	
995				1050	
996				1051	
997				1052	
998				1053	
999				1054	
1000				1055	
1001				1056	
1002				1057	
1003				1058	
1004				1059	
1005				1060	
1006				1061	
1007				1062	
1008				1063	
1009				1064	
1010				1065	
1011				1066	
1012				1067	
1013				1068	
1014				1069	
1015				1070	
1016				1071	
1017	Q1586: what is the color of the wheels ?	Ground truth: yellow 2-CNN-LSTM: <b>yellow (0.4101)</b> CNN-BOW: <b>brown (0.4306)</b> BLIND-LSTM: <b>orange (0.2095)</b> BLIND-BOW: <b>blue (0.1855)</b> Attend to detail	Q6358: what is the color of the boat ? Ground truth: red 2-CNN-LSTM: <b>red (0.5345)</b> CNN-BOW: <b>blue (0.2879)</b> BLIND-LSTM: <b>white (0.1832)</b> BLIND-BOW: <b>white (0.2890)</b> Not the colorful umbrella it is asking	Q10492: what is the color of the plane ? Ground truth: red 2-CNN-LSTM: <b>red (0.5145)</b> CNN-BOW: <b>blue (0.6872)</b> BLIND-LSTM: <b>white (0.2756)</b> BLIND-BOW: <b>white (0.2865)</b> BOW goes for the dominant color of the image again	1072
1018				1073	
1019				1074	
1020				1075	
1021				1076	
1022				1077	
1023				1078	
1024	Q1037: what is the color of the skies ?	Ground truth: blue 2-CNN-LSTM: <b>blue (0.4425)</b> CNN-BOW: <b>blue (0.6370)</b> BLIND-LSTM: <b>blue (0.3826)</b> BLIND-BOW: <b>blue (0.4494)</b> Sometimes colors are easy to guess	Q591: what is the color of the shorts ? Ground truth: black 2-CNN-LSTM: <b>blue (0.4518)</b> CNN-BOW: <b>blue (0.3261)</b> BLIND-LSTM: <b>black (0.1963)</b> BLIND-BOW: <b>black (0.2141)</b> CNN can't tell that shorts are worn on lower part of the body	Q2211: what is the color of the salad ? Ground truth: green 2-CNN-LSTM: <b>white (0.3428)</b> CNN-BOW: <b>white (0.3321)</b> BLIND-LSTM: <b>green (0.3224)</b> BLIND-BOW: <b>green (0.3310)</b> CNN is really looking at the background instead of the foreground	1079
1025				1080	
1026				1081	
1027				1082	
1028				1083	
1029				1084	
1030				1085	
1031				1086	
1032				1087	
1033	Q17918: what is the color of the plate ?	Ground truth: white 2-CNN-LSTM: <b>orange (0.4249)</b> CNN-BOW: <b>orange (0.4252)</b> BLIND-LSTM: <b>white (0.6865)</b> BLIND-BOW: <b>white (0.7853)</b> Not asking about the forgrond dominant object but CNN just can't help	Q2562: what is the color of the ribbon ? Ground truth: yellow 2-CNN-LSTM: <b>white (0.8117)</b> CNN-BOW: <b>white (0.8068)</b> BLIND-LSTM: <b>red (0.2475)</b> BLIND-BOW: <b>red (0.1830)</b> Very hard question and even if it gets right it might just focusing on the bananas	Q3705: what is the color of the flowers ? Ground truth: blue 2-CNN-LSTM: <b>orange (0.6924)</b> CNN-BOW: <b>orange (0.3380)</b> BLIND-LSTM: <b>purple (0.6047)</b> BLIND-BOW: <b>purple (0.5514)</b> Very unusual scene setting but again CNN doesn't focus on the right object	1088
1034				1089	
1035				1090	
1036				1091	
1037				1092	
1038				1093	
1039				1094	
1040				1095	
1041				1096	
1042				1097	
1043				1098	
1044				1099	

Figure 9. COCO-QA Color Questions

1100	1155
1101	1156
1102	1157
1103	1158
1104	1159
1105	1160
1106	1161
1107	1162
1108	1163
1109	1164
1110	1165
1111	1166
1112	1167
1113	1168
1114	1169
1115	1170
1116	1171
1117	1172
1118	1173
1119	1174
1120	1175
1121	1176
1122	1177
1123	1178
1124	1179
1125	1180
1126	1181
1127	1182
Q5826: what is the color of the sign ?	
Ground truth: blue	
2-CNN-LSTM: red (0.3877)	1183
CNN-BOW: red (0.3503)	1184
BLIND-LSTM: red (0.4757)	1185
BLIND-BOW: red (0.3888)	1186
Seems like all models are using textual information about a sign instead of the image	1187
	1188
	1189
1135	1190
1136	1191
1137	1192
1138	1193
1139	1194
1140	1195
1141	1196
1142	1197
1143	1198
1144	1199
1145	1200
1146	1201
1147	1202
1148	1203
1149	1204
1150	1205
1151	1206
1152	1207
1153	1208
1154	1209

Figure 10. COCO-QA Color Questions



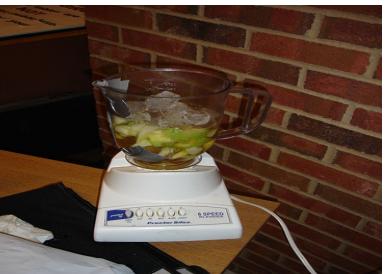
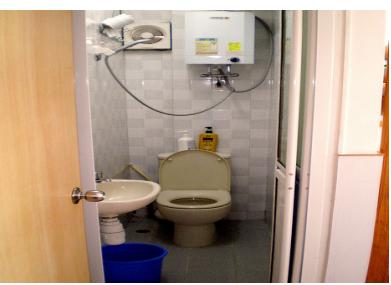
1210			1265
1211			1266
1212			1267
1213			1268
1214			1269
1215			1270
1216			1271
1217			1272
1218			1273
1219			1274
1220	Q42: where is a white tiger ? Ground truth: pool 2-CNN-LSTM: pool (0.1984) CNN-BOW: pool (0.2602) BLIND-LSTM: bathroom (0.0935) BLIND-BOW: room (0.2206) Image features are useful		1275
1221			1276
1222	Q857: where is the woman ? Ground truth: station 2-CNN-LSTM: station (0.0833) CNN-BOW: station (0.2196) BLIND-LSTM: bed (0.0644) BLIND-BOW: room (0.1590) Image features are useful		1277
1223			1278
1224			1279
1225			1280
1226			1281
1227	Q805: where is fruit and ice ? Ground truth: blender 2-CNN-LSTM: blender (0.1350) CNN-BOW: blender (0.3828) BLIND-LSTM: bowl (0.3528) BLIND-BOW: bowl (0.5479) Inclusion of object (not just room and scene)		1282
1228			1283
1229			1284
1230			1285
1231			1286
1232			1287
1233			1288
1234			1289
1235			1290
1236	Q9648: where does the man with glasses and a tie clip sit ? Ground truth: car 2-CNN-LSTM: car (0.7772) CNN-BOW: car (0.4117) BLIND-LSTM: bowl (0.1078) BLIND-BOW: room (0.3101) This one is hard to see it is in a car		1291
1237			1292
1238			1293
1239	Q15891: where did the cat curl up UNK laying ? Ground truth: sink 2-CNN-LSTM: sink (0.8603) CNN-BOW: sink (0.4068) BLIND-LSTM: bed (0.3623) BLIND-BOW: plate (0.0401) Easy one for CNN		1294
1240			1295
1241			1296
1242			1297
1243	Q25218: where are the ripe bananas sitting ? Ground truth: basket 2-CNN-LSTM: basket (0.4965) CNN-BOW: basket (0.4929) BLIND-LSTM: bowl (0.6415) BLIND-BOW: bowl (0.5170) Inclusion of object		1298
1244			1299
1245			1300
1246			1301
1247			1302
1248			1303
1249			1304
1250			1305
1251			1306
1252			1307
1253	Q25397: where are people looking are different products ? Ground truth: store 2-CNN-LSTM: store (0.2376) CNN-BOW: store (0.2212) BLIND-LSTM: kitchen (0.1704) BLIND-BOW: window (0.1182) Hard to infer		1308
1254			1309
1255	Q35349: where is an elephant approaching ? Ground truth: mirror 2-CNN-LSTM: mirror (0.1540) CNN-BOW: mirror (0.8029) BLIND-LSTM: zoo (0.1094) BLIND-BOW: zoo (0.1345) The model knows that the elephant is inside the mirror		1310
1256			1311
1257			1312
1258			1313
1259			1314
1260			1315
1261			1316
1262			1317
1263			1318
1264			1319

Figure 11.

## Question Answering about Images using Visual Semantic Embeddings

---

1320			1375
1321			1376
1322			1377
1323			1378
1324			1379
1325			1380
1326			1381
1327			1382
1328			1383
1329	Q23418: where do the girl and her dog sit ?		1384
1330	Ground truth: chair		1385
1331	2-CNN-LSTM: <b>chair</b> (0.4588)		1386
1332	CNN-BOW: <b>bed</b> (0.2899)		1387
1333	BLIND-LSTM: <b>bed</b> (0.2878)		1388
1334	BLIND-BOW: <b>bed</b> (0.2042)		1389
1335	Funny picture		1390
1336			1391
1337			1392
1338			1393
1339			1394
1340			1395
1341			1396
1342			1397
1343			1398
1344			1399
1345			1400
1346			1401
1347	Q26714: where do orange flowers fill vases ?		1402
1348	Ground truth: shop		1403
1349	2-CNN-LSTM: <b>vase</b> (0.6365)		1404
1350	CNN-BOW: <b>vase</b> (0.9500)		1405
1351	BLIND-LSTM: <b>vase</b> (0.4533)		1406
1352	BLIND-BOW: <b>vase</b> (0.4836)		1407
1353	Many location based questions are noisy in the sense that there exist many valid answer	Synonyms are also the noise of the dataset	1408
1354			1409
1355			1410
1356			1411
1357			1412
1358			1413
1359			1414
1360			1415
1361			1416
1362			1417
1363			1418
1364			1419
1365	Q172: where are the sink and toilet ?		1420
1366	Ground truth: bathroom		1421
1367	2-CNN-LSTM: <b>bathroom</b> (0.8315)		1422
1368	CNN-BOW: <b>bathroom</b> (0.6268)		1423
1369	BLIND-LSTM: <b>bathroom</b> (0.7750)		1424
1370	BLIND-BOW: <b>bathroom</b> (0.8103)		1425
1371	Some questions are fairly easy and can be answered without looking at the image	Easy one	1426
1372			1427
1373			1428
1374			1429

*Figure 12.*