

What can we extract from explicit detections for image captioning?

		Person	Bicycle		Train			Banana			Bench		Sno	Hood
Frequency	[3	1			0		0			3		()
Normalized Frequency	[0.43	0.1	4		0		0		().43	3	()
Presence	[1	1			0		0			1		()
Most central instance	[0.09	0.1	5		0		0			0.81		()
Size of biggest instance	[0.18	0.2	27		0		0		().20)	()
Size centrality	[0.18 0.2	27 0	0 0.2	20 (0		0.0	09 (0.15	00	0 0	.81 ()
Frequency centrality	[3 1	0 () 3	()		0.0	09 (0.15	00	0 0	.81 ()
Frequency size	[3 1	0 () 3)		0.1	8 0	.27	0 0	0.	.20 ()
Freq. centrality size	[310	03	0	C	.09	0.15	5	0		0.18	0.	27	. 0
		person	хy	wł	ı a	X	y w	h	a >	(y	W	h	a	
Full spatial info	[bicycle train	x y 0 0	w ł 0 (1 a) 0	00	0 C 0 C	0 0	0 (0 0 (0	0 0 0	0 0	0 0	0 0	-
		bench	хy	wł	ı a	X	yw	h	a>	(y	W	h	a]

(x, y): coordinates of object center, (w, h): width/height of object bounding box, (a): area of object segment



Rank correlation between Δ CIDEr versus:

- frequency of object being annotated
- $\circ f(v_c)$: Kendal's au=0.093, Spearman's ho=0.137
- probability of object being mentioned given that it is depicted $p(t_c|v_c)$: Kendal's au=0.153, Spearman's ho=0.227





What happens when information is removed? Mask N% of object frequencies with different heuristics:



Summary:

- Importance: Frequency > Size > Centrality
- All three features are complementary
- A more frequent category is not necessarily more important
- Larger objects are more important than smaller objects
- Objects near the image center > away from the center (but size is more important)



Frequency: Centrality: All 3 feats:

 (\mathbf{x},\mathbf{y}) : (w,h): (a): (x,y,w,h,a): CNN: person removed:

Example output

person, clock

a large clock tower with a large clock on it . a clock tower with a large clock on it 's face . a man standing in front of a clock tower. a clock tower with people standing in the middle of the water.

a large clock tower with a clock on the front . a clock on a pole in front of a building a large clock tower with people walking around it a group of people standing around a clock tower. a large building with a clock tower in the middle of it. a clock tower with a weather vane on top of it.

Josiah Wang, Pranava Madhyastha, Lucia Specia

Why does bag of objects work so well?

Object distribution in train vs. test is almost identical! A retrieval machine maybe?



Is neural image captioning merely a multimodal retrieval machine?

If yes, then the model should generate similar captions for the nearest neighbour(s) of an image!





