

Neural Image Caption Generation with Visual Attention

Anton Karazeev
493 group

MIPT, 2017

Plan

- Image Caption Generation with Attention Mechanism
- “Soft” vs “Hard” Attention
- Experiments

<https://arxiv.org/pdf/1502.03044.pdf>

**Show, Attend and Tell: Neural Image Caption
Generation with Visual Attention**

Kelvin Xu

Jimmy Lei Ba

Ryan Kiros

Kyunghyun Cho

Aaron Courville

Ruslan Salakhutdinov

Richard S. Zemel

Yoshua Bengio

KELVIN.XU@UMONTREAL.CA

JIMMY@PSI.UTORONTO.CA

RKIROS@CS.TORONTO.EDU

KYUNGHYUN.CHO@UMONTREAL.CA

AARON.COURVILLE@UMONTREAL.CA

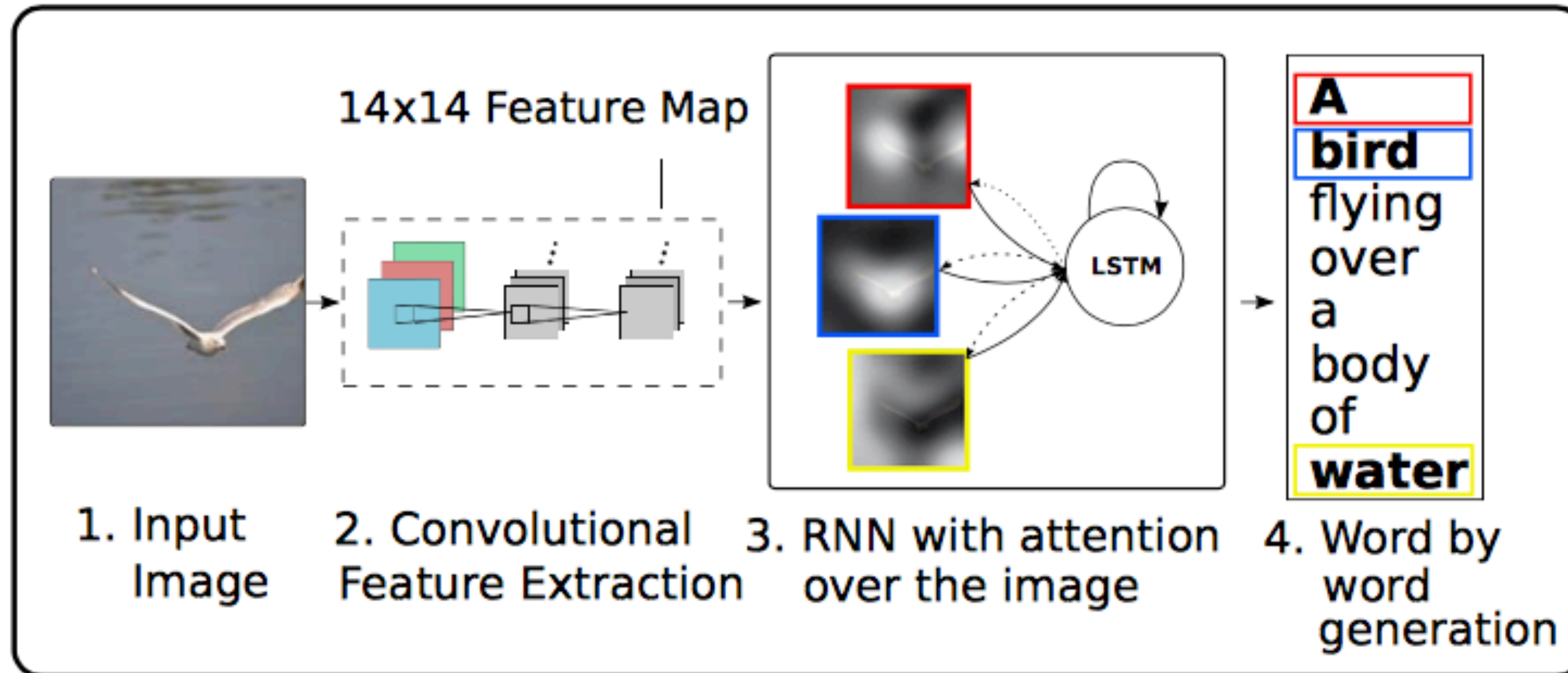
RSALAKHU@CS.TORONTO.EDU

ZEMEL@CS.TORONTO.EDU

FIND-ME@THE.WEB

Main Idea

Figure 1. Our model learns a words/image alignment. The visualized attentional maps (3) are explained in section 3.1 & 5.4



Attention

- The ability to visualize what the model “sees”



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



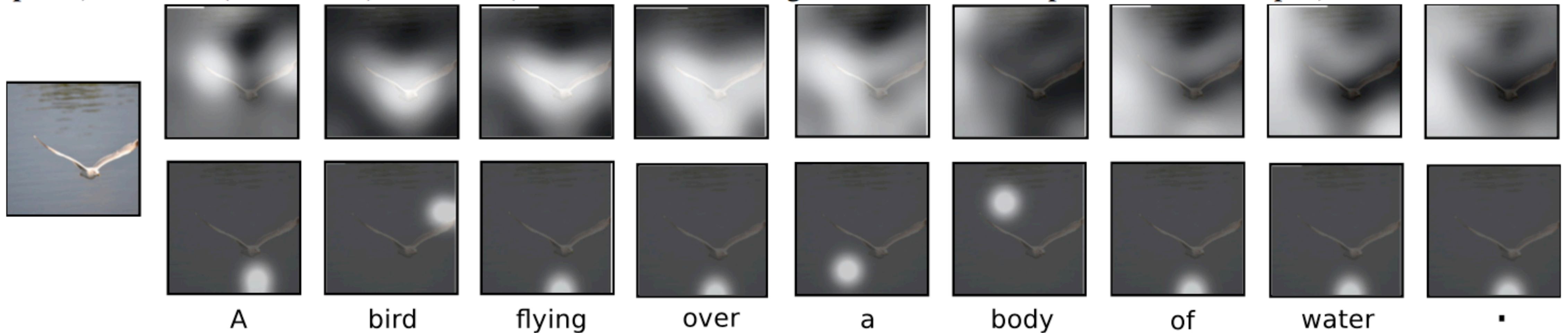
A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

“Soft” vs “Hard” Attention

Figure 2. Attention over time. As the model generates each word, its attention changes to reflect the relevant parts of the image. “soft” (top row) vs “hard” (bottom row) attention. (Note that both models generated the same captions in this example.)



Metrics

- BLEU
- METEOR

Experiments

Table 1. BLEU-1,2,3,4/METEOR metrics compared to other methods, † indicates a different split, (—) indicates an unknown metric, ◦ indicates the authors kindly provided missing metrics by personal communication, Σ indicates an ensemble, a indicates using AlexNet

Dataset	Model	BLEU				METEOR
		BLEU-1	BLEU-2	BLEU-3	BLEU-4	
Flickr8k	Google NIC(Vinyals et al., 2014) ^{†Σ}	63	41	27	—	—
	Log Bilinear (Kiros et al., 2014a) [◦]	65.6	42.4	27.7	17.7	17.31
	Soft-Attention	67	44.8	29.9	19.5	18.93
	Hard-Attention	67	45.7	31.4	21.3	20.30
Flickr30k	Google NIC ^{†$\circ\Sigma$}	66.3	42.3	27.7	18.3	—
	Log Bilinear	60.0	38	25.4	17.1	16.88
	Soft-Attention	66.7	43.4	28.8	19.1	18.49
	Hard-Attention	66.9	43.9	29.6	19.9	18.46
COCO	CMU/MS Research (Chen & Zitnick, 2014) ^a	—	—	—	—	20.41
	MS Research (Fang et al., 2014) ^{†a}	—	—	—	—	20.71
	BRNN (Karpathy & Li, 2014) [◦]	64.2	45.1	30.4	20.3	—
	Google NIC ^{†$\circ\Sigma$}	66.6	46.1	32.9	24.6	—
	Log Bilinear [◦]	70.8	48.9	34.4	24.3	20.03
	Soft-Attention	70.7	49.2	34.4	24.3	23.90
	Hard-Attention	71.8	50.4	35.7	25.0	23.04

Summary

- Image Caption Generation with Attention Mechanism
- “Soft” vs “Hard” Attention
- Experiments

“The best way to predict the future is to create it.”

– Abraham Lincoln