

---

# simNet: Stepwise Image-Topic Merging Network for Generating Detailed and Comprehensive Image Captions

---

Fenglin Liu<sup>1</sup>, Xuancheng Ren<sup>2\*</sup>, Yuanxin Liu<sup>1</sup>, Houfeng Wang<sup>2</sup> and Xu Sun<sup>2</sup>

<sup>1</sup> Beijing University of Posts and Telecommunications, Beijing, China

<sup>2</sup> Peking University, Beijing, China

\* Equal Contributions

# CONTENTS

**1** Introduction **3** Experiment

**2** Approach **4** Analysis



北京邮电大学 x  
Beijing University of Posts and Telecommunications



北京大学  
PEKING UNIVERSITY



# Introduction

---



北京邮电大学 X  
Beijing University of Posts and Telecommunications



北京大学  
PEKING UNIVERSITY

# simNet: Stepwise Image-Topic Merging Network for Generating Detailed and Comprehensive Image Captions



**Soft-Attention:** a open laptop computer sitting on top of a table

**ATT-FCN:** a dog sitting on a desk with a laptop computer and mouse

**simNet:** a open laptop computer and mouse sitting on a table with a dog nearby

Figure 1: Examples of using different attention mechanisms.

•**Soft-Attention:** Show, attend and tell: Neural image caption generation with visual attention. In PMLR 2015

•**ATT-FCN :** Image captioning with semantic attention. In CVPR 2016



北京邮电大学  
Beijing University of Posts and Telecommunications

x



北京大学  
PEKING UNIVERSITY

| 4

# Introduction: Soft-Attention



**Soft-Attention:** a open laptop computer sitting on top of a table

→ omitting “dog” and “mouse”

**ATT-FCN:** a dog sitting on a desk with a laptop computer and mouse

**simNet:** a open laptop computer and mouse sitting on a table with a dog nearby



•**Soft-Attention:** Show, attend and tell: Neural image caption generation with visual attention. In PMLR 2015



北京邮电大学  
Beijing University of Posts and Telecommunications

x



北京大学  
PEKING UNIVERSITY

# Introduction: ATT-FCN



**Soft-Attention:** a open laptop computer sitting on top of a table

**ATT-FCN:** a dog sitting on a desk with a laptop computer and mouse

→ missing “open” and mislocating “dog”

**simNet:** a open laptop computer and mouse sitting on a table with a dog nearby



•ATT-FCN : Image captioning with semantic attention. In CVPR 2016



北京邮电大学  
Beijing University of Posts and Telecommunications

x



北京大学  
PEKING UNIVERSITY



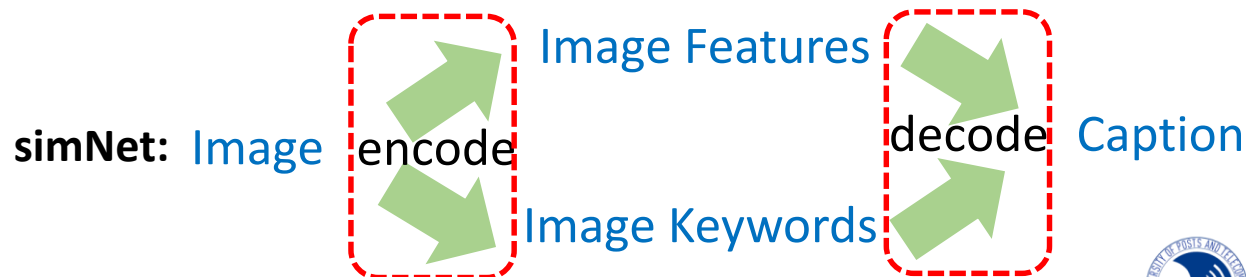
# Introduction: SimNet



**Soft-Attention:** a open laptop computer sitting on top of a table

**ATT-FCN:** a dog sitting on a desk with a laptop computer and mouse

**simNet:** a open laptop computer and mouse sitting on a table with a dog nearby



# Introduction: Main idea

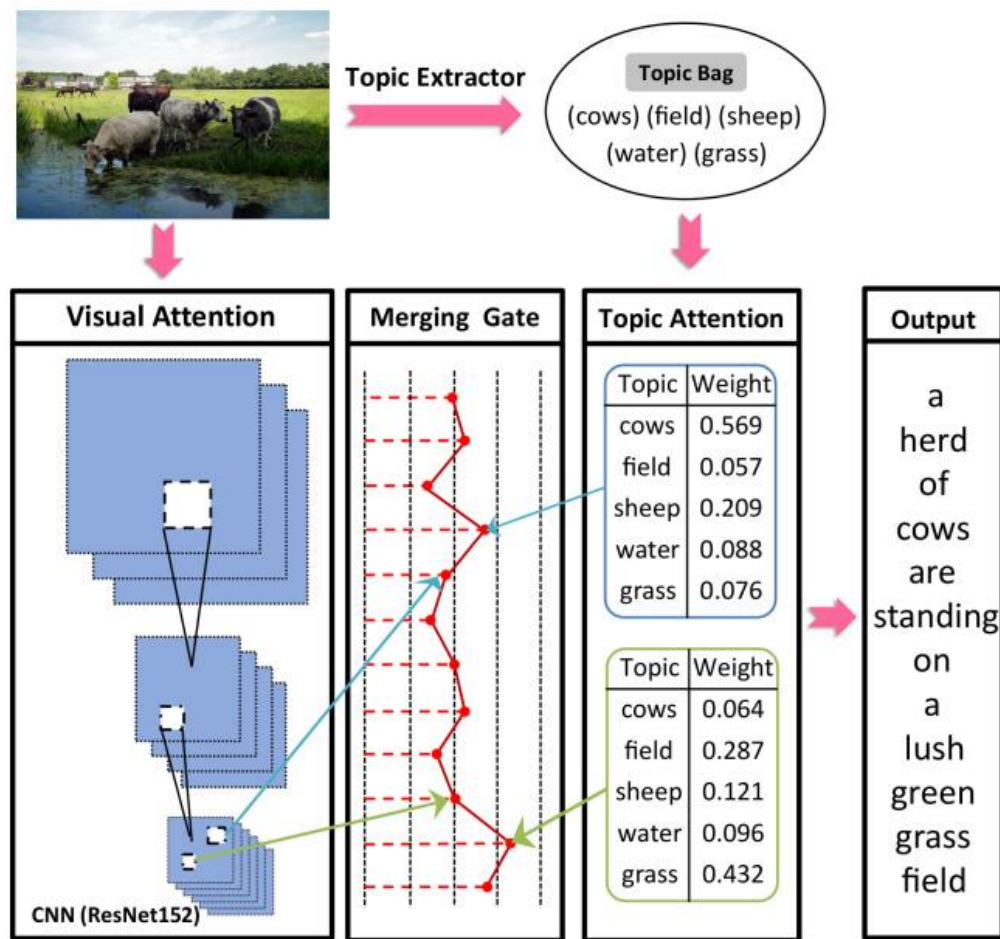


Figure 2: Illustration of the main idea.

- The visual information captured by CNN
- The topics extracted by a topic extractor
- The merging gate then **adaptively adjusts the weight** between visual attention and topic attention



# Contributions

- We propose a novel approach that can effectively merge the **information in the image** and the **topics**.
- The generated captions are both **detailed** and **comprehensive**.
- The proposed approach **outperforms** previous works in terms of SPICE, which correlates the best with human judgments.



# 2

# Approach

---



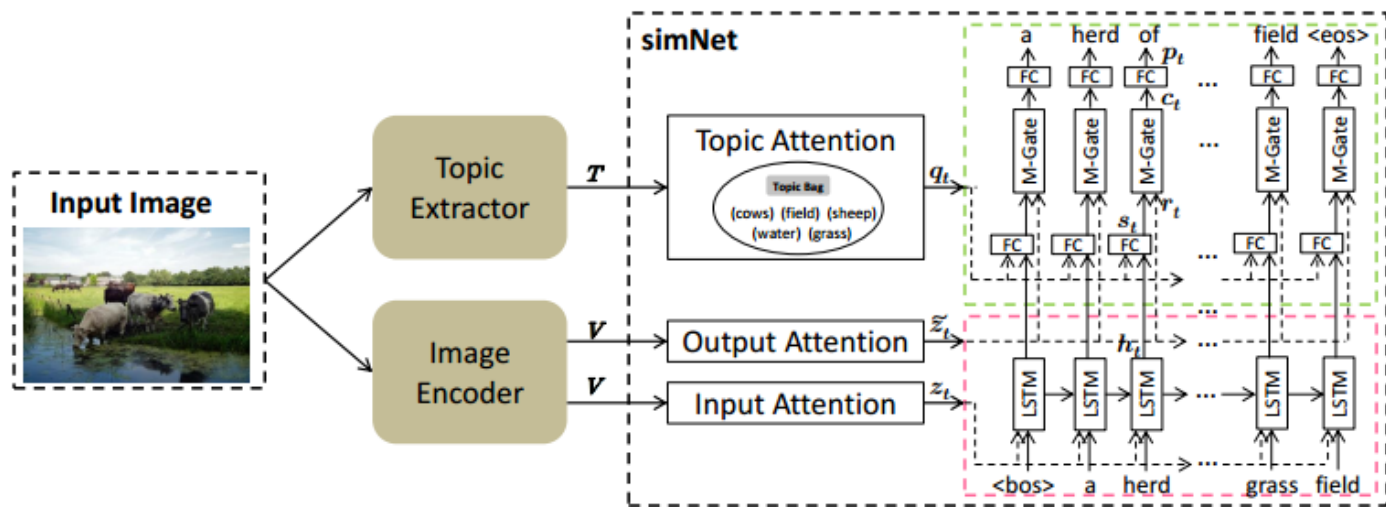
北京邮电大学  
Beijing University of Posts and Telecommunications

x

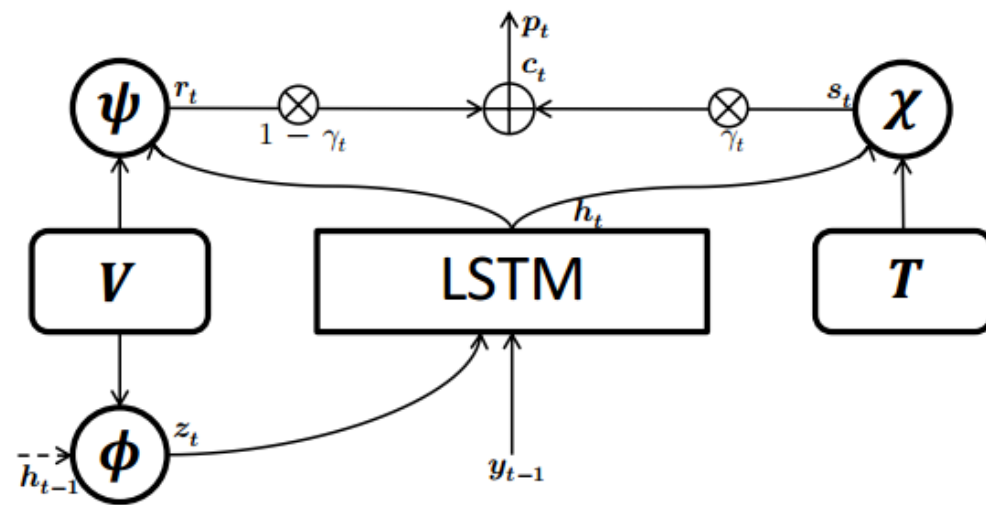


北京大学  
PEKING UNIVERSITY

# Overview



(a) The overall framework.



(b) The data flow in the proposed simNet.

Figure 3: Illustration of the proposed approach. In the right plot, we use  $\phi$ ,  $\psi$ ,  $\chi$  to denote input attention, output attention, and topic attention, respectively.

# Approach: Image Encoder

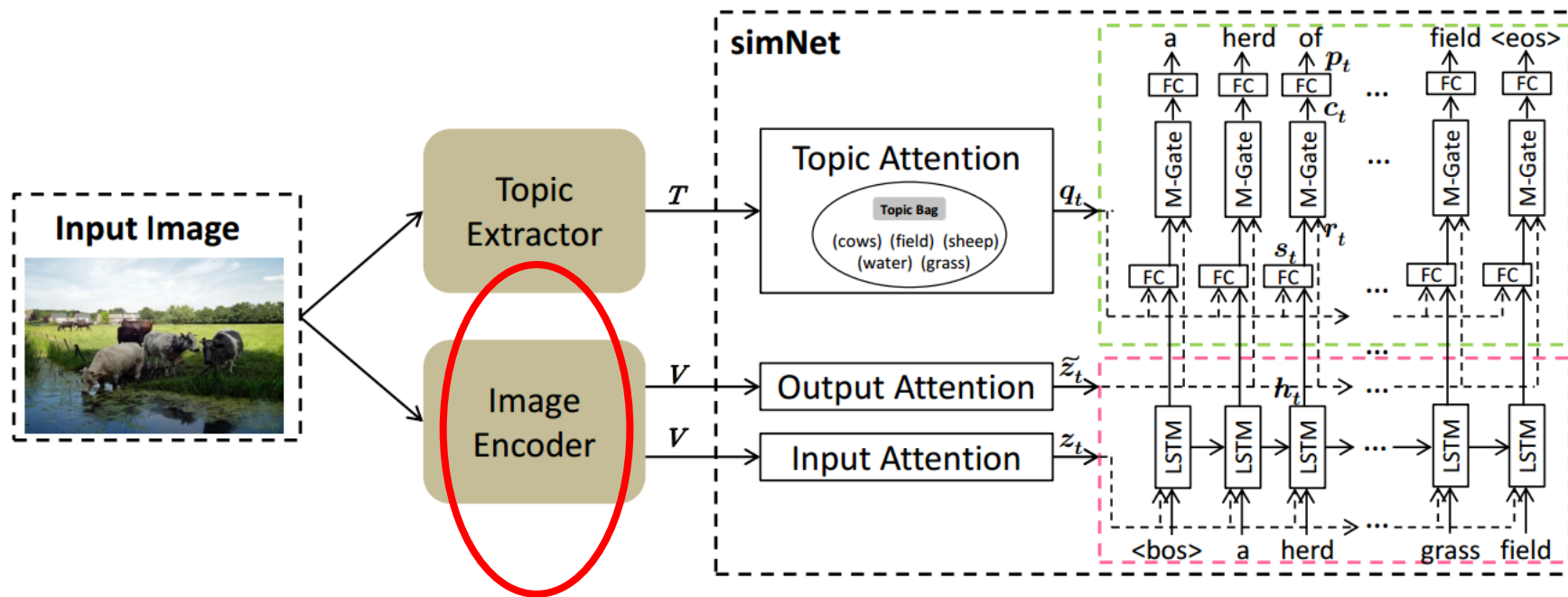


Image Encoder: ResNet152

He et al., 2016: Deep residual learning for image recognition. In CVPR 2016.

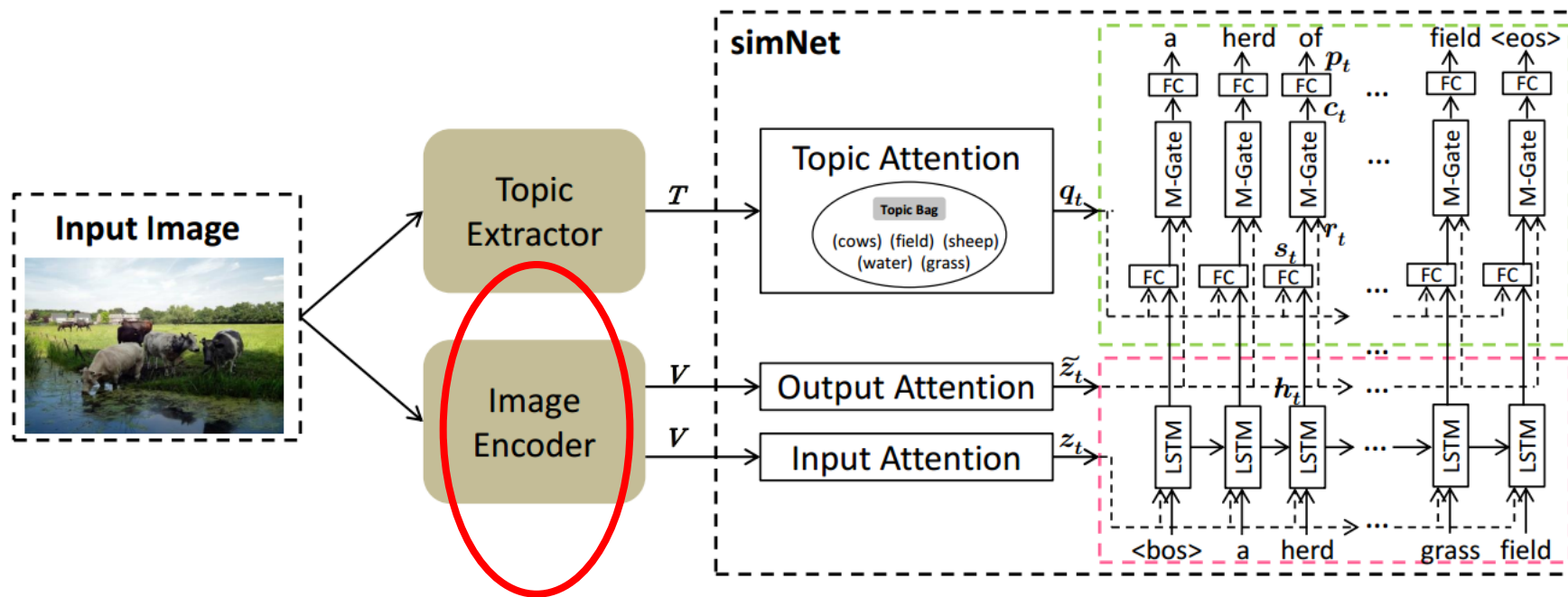


北京邮电大学  
Beijing University of Posts and Telecommunications



北京大学  
PEKING UNIVERSITY

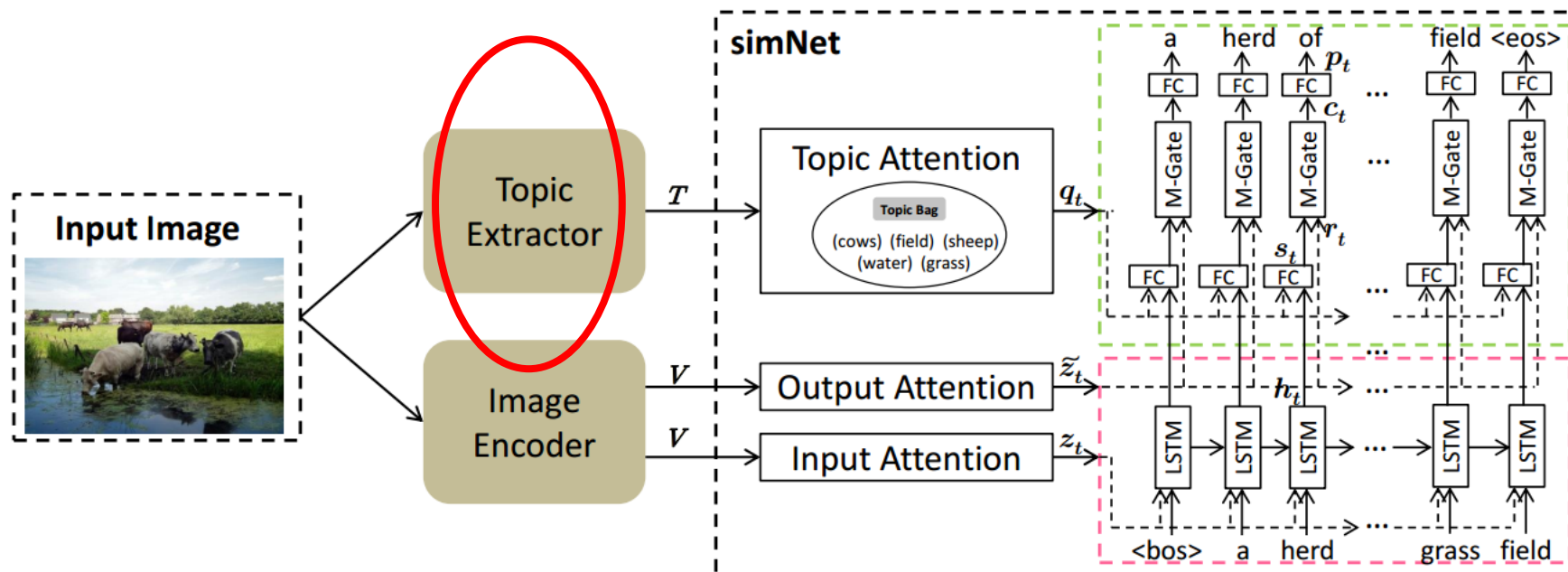
# Approach: Image Encoder



$$\text{Feature map: } V = W^{V,I} \text{CNN}(I) \quad (1)$$

where  $I$  is the input image, and  $W^{V,I}$  shrinks the last dimension of the output.

# Approach: Topic Extractor



Topic Extractor: Multiple Instance Learning

Zhang et al., 2006: Multiple instance boosting for object detection. In NIPS 2006.

Fang et al., 2015: From captions to visual concepts and back. In CVPR2015



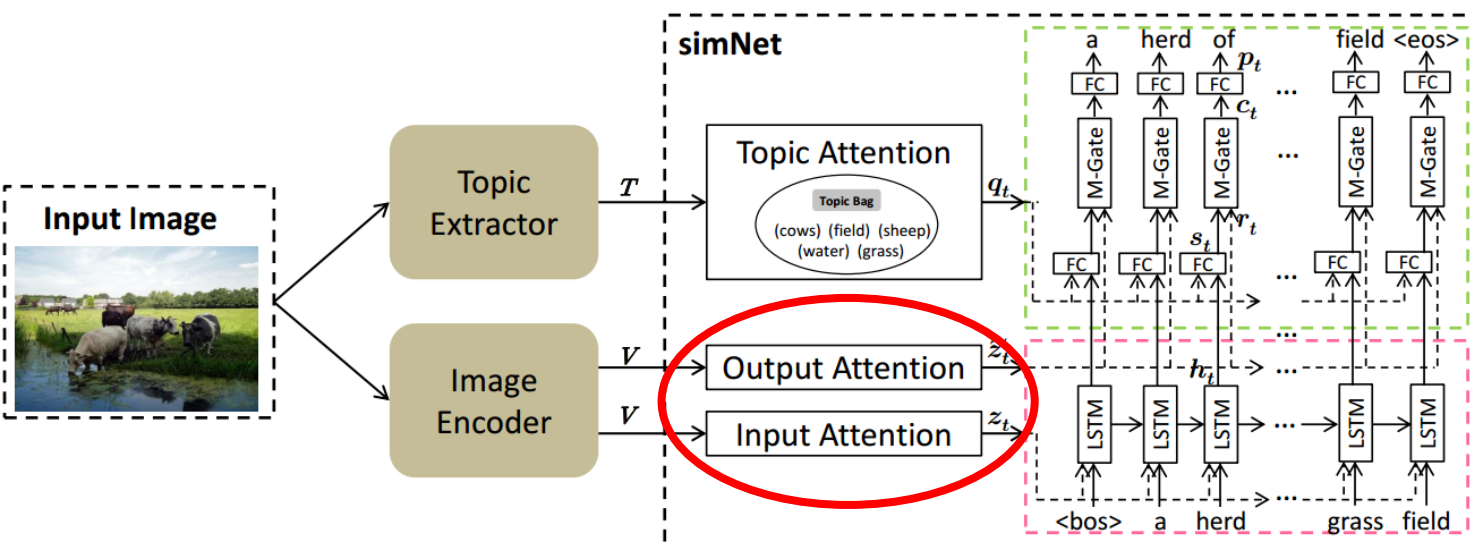
北京邮电大学  
Beijing University of Posts and Telecommunications



北京大学  
PEKING UNIVERSITY



# Approach: Input Attention



Input Attention:

$$Z_t = \tanh(W^{Z,V} V \oplus W^{Z,h} \underline{h_{t-1}}) \quad (2)$$

$$\alpha_t = \text{softmax}(Z_t w^{\alpha,Z}) \quad (3)$$

$$z_t = V \alpha_t \quad (4)$$

$$h_t = \text{LSTM}\left(\begin{bmatrix} z_t \\ y_{t-1} \end{bmatrix}, h_{t-1}\right) \quad (5)$$

Xu et al., 2015 : Show, attend and tell: Neural image caption generation with visual attention. In PMLR 2015

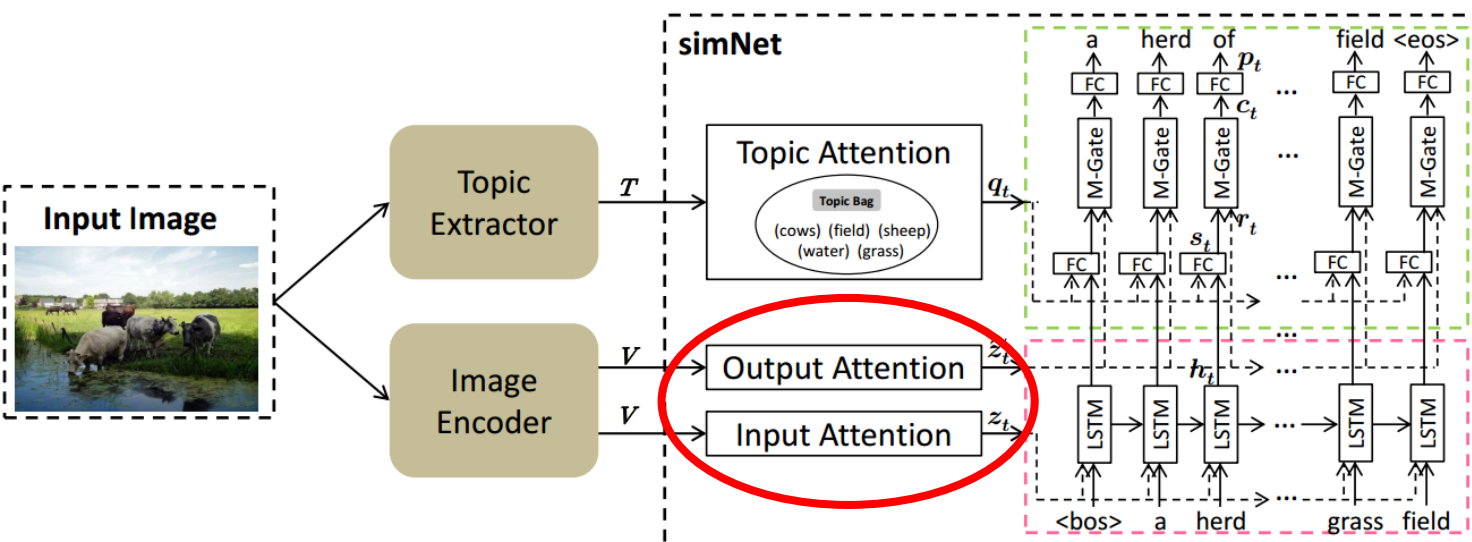


北京邮电大学 x



北京大学  
PEKING UNIVERSITY

# Approach: Output Attention



Output Attention:

$$\tilde{Z}_t = \tanh(\tilde{W}^{Z,V} V \oplus \tilde{W}^{Z,h} \underline{h_t}) \quad (6)$$

$$\tilde{\alpha}_t = \text{softmax}(\tilde{Z}_t \tilde{w}^{\alpha,Z}) \quad (7)$$

$$\tilde{z}_t = V \tilde{\alpha}_t \quad (8)$$

You et al., 2016 : Image captioning with semantic attention. In CVPR 2016

Lu et al., 2017 : Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In CVPR 2017

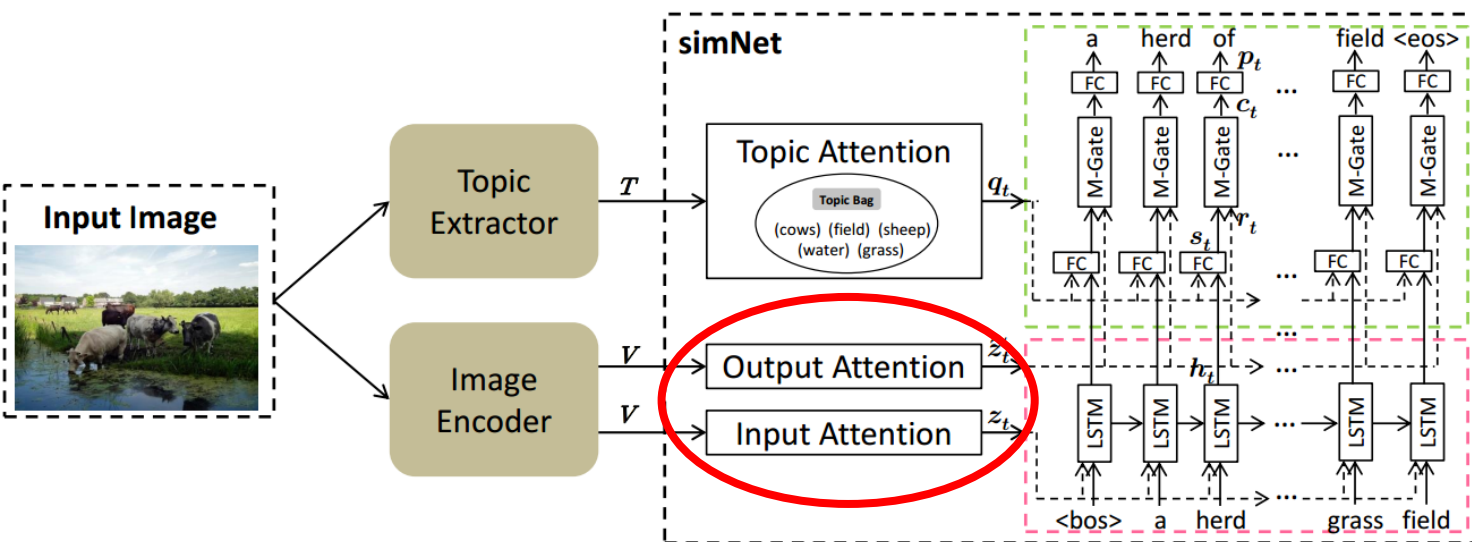


北京邮电大学 x



北京大学

# Approach: Visual Information



Output Attention:

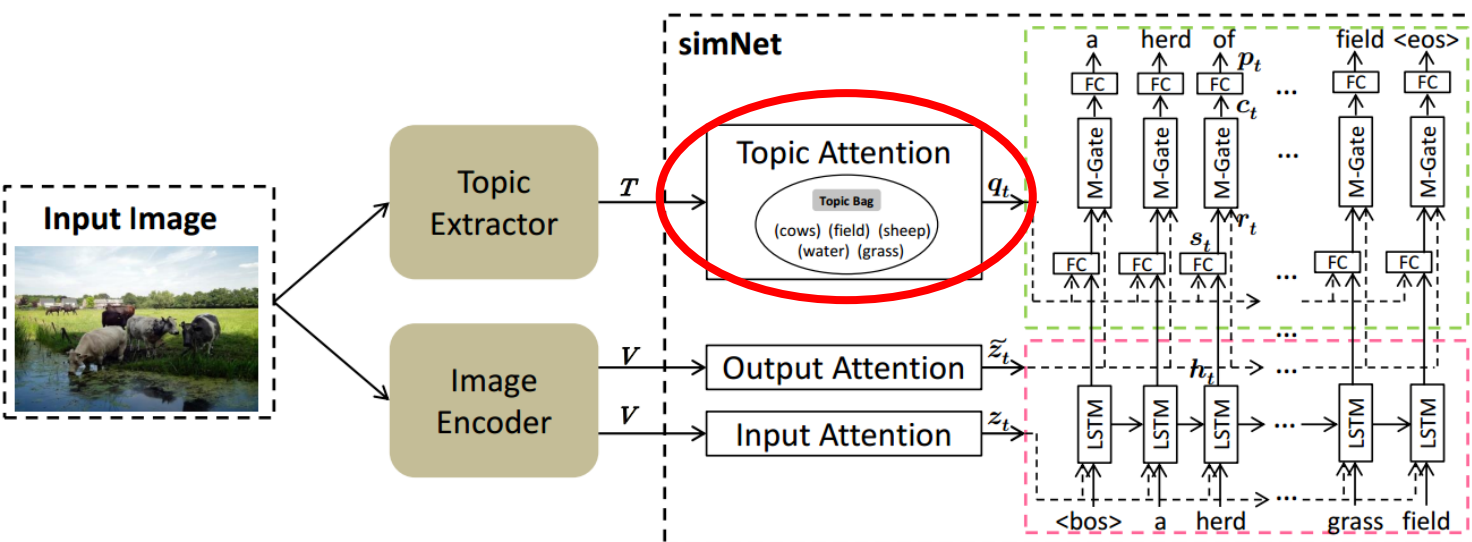
$$\tilde{Z}_t = \tanh(\tilde{W}^{Z,V} V \oplus \tilde{W}^{Z,h} h_t) \quad (6)$$

$$\tilde{\alpha}_t = \text{softmax}(\tilde{Z}_t \tilde{w}^{\alpha,Z}) \quad (7)$$

$$\tilde{z}_t = V \tilde{\alpha}_t \quad (8)$$

the visual information:  $r_t = \tanh(W^{s,z} \tilde{z}_t)$

# Approach: Previous Topic Attention



Topic Attention (Previous work):

$$\beta_t = \text{softmax}(T^T U y_{t-1}) \quad (9)$$

Lacking the attentive visual information  
when selecting topic!

You et al., 2016 : Image captioning with semantic attention. In CVPR 2016



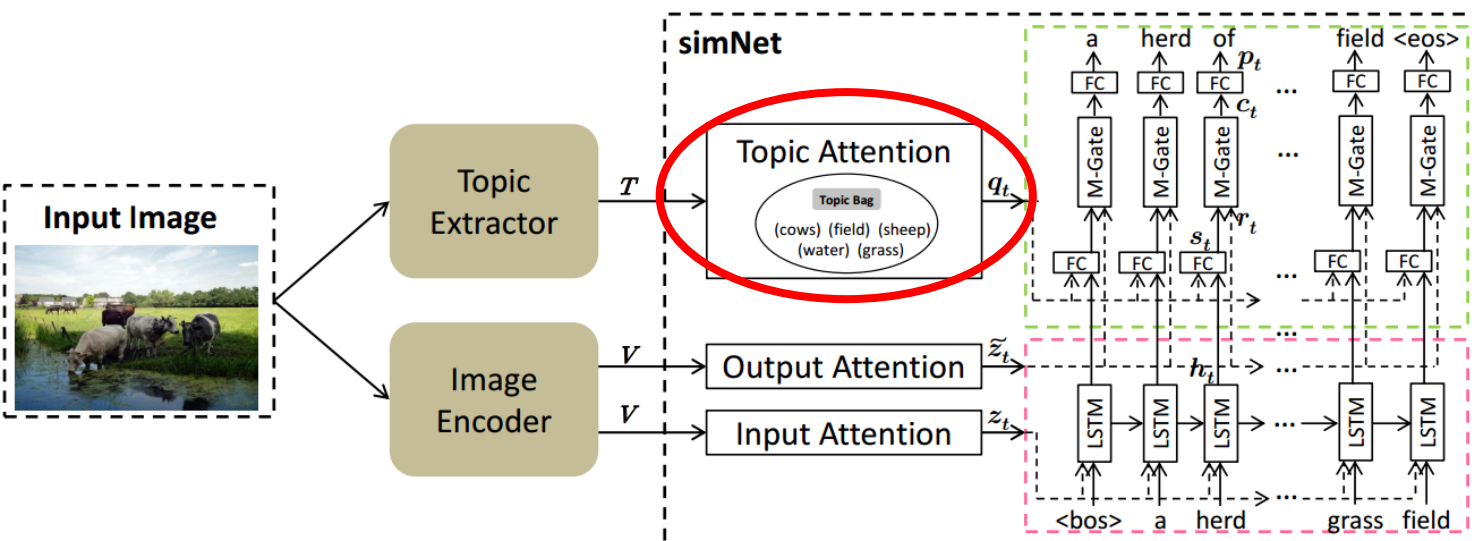
北京邮电大学 x



北京大学



# Approach: Contextual Information



Topic Attention (Our):

$$Q_t = \tanh(W^{Q,T}T \oplus W^{Q,h}h_t) \quad (10)$$

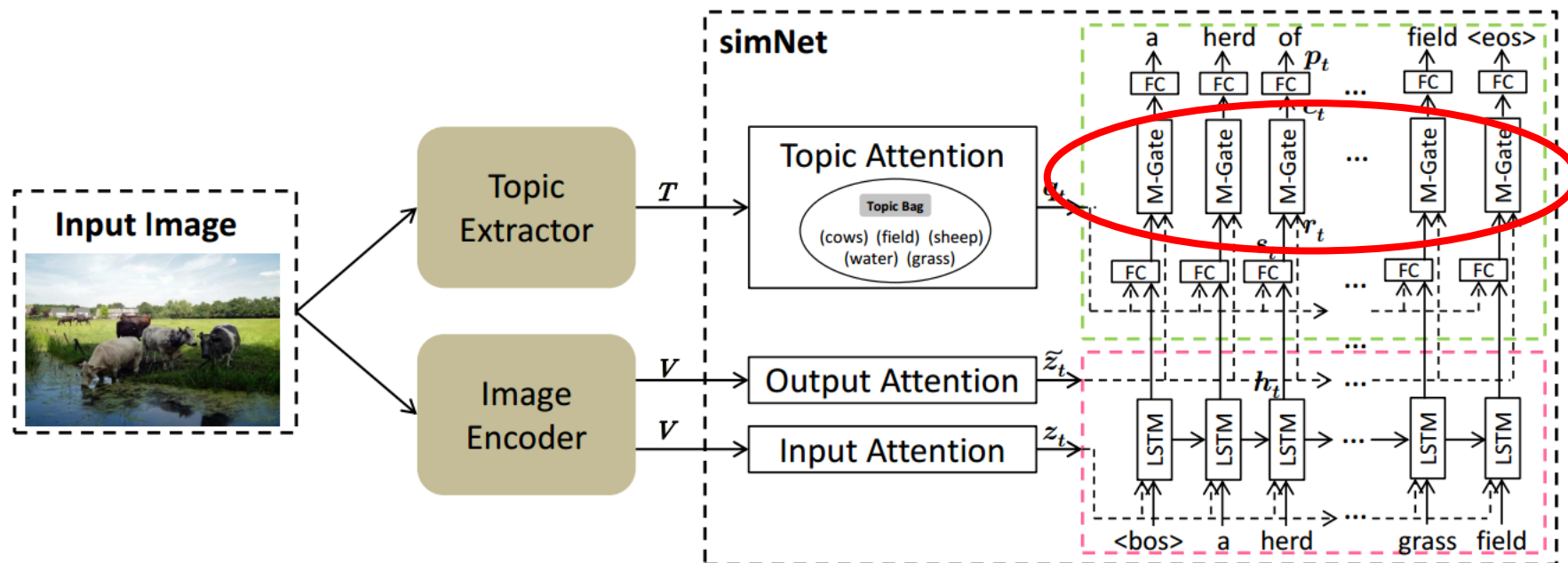
$$\beta_t = \text{softmax}(Q_t w^{\beta,Q}) \quad (11)$$

$$q_t = T\beta_t \quad (12)$$

the contextual information:  $s_t = \tanh(W^{s,q}q_t + W^{s,h}h_t)$

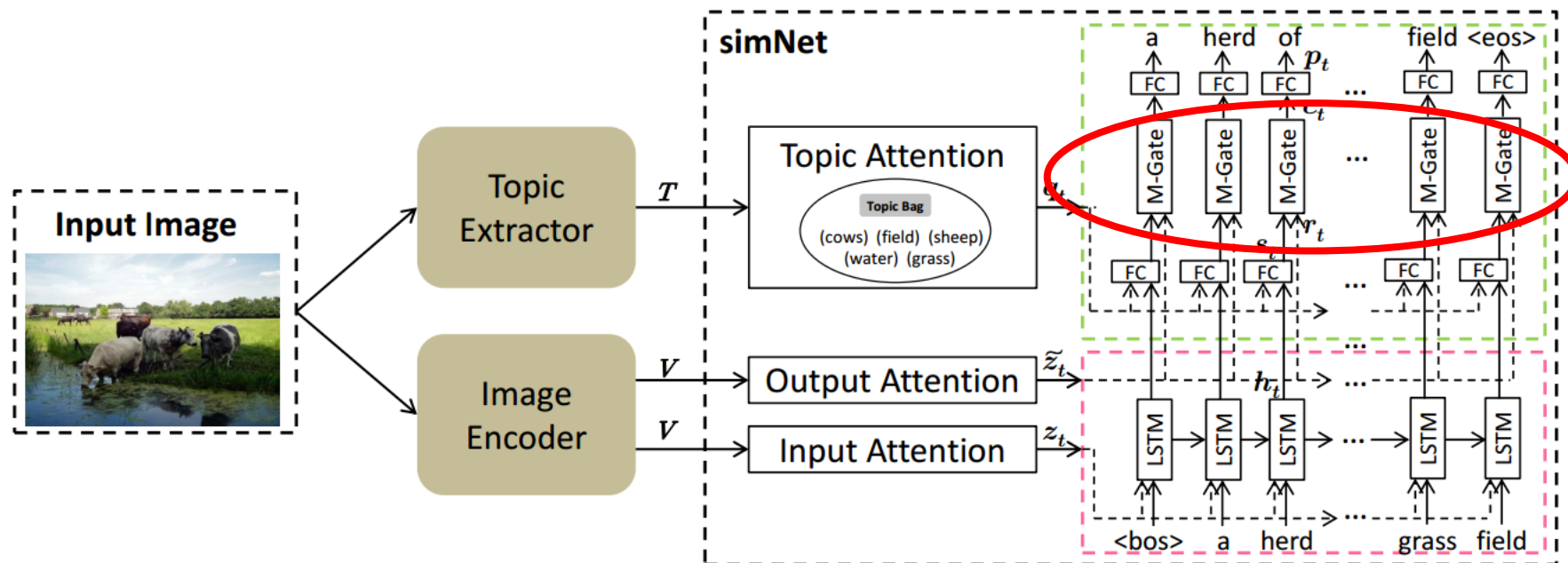


# Approach: Merging Gate



How to make full use of the visual information and the contextual information?

# Approach: Merging Gate



Visual information  
(e.g., “*behind*”, “*red*” is better)

VS

Contextual information  
(e.g., “*people*”, “*table*” is better)

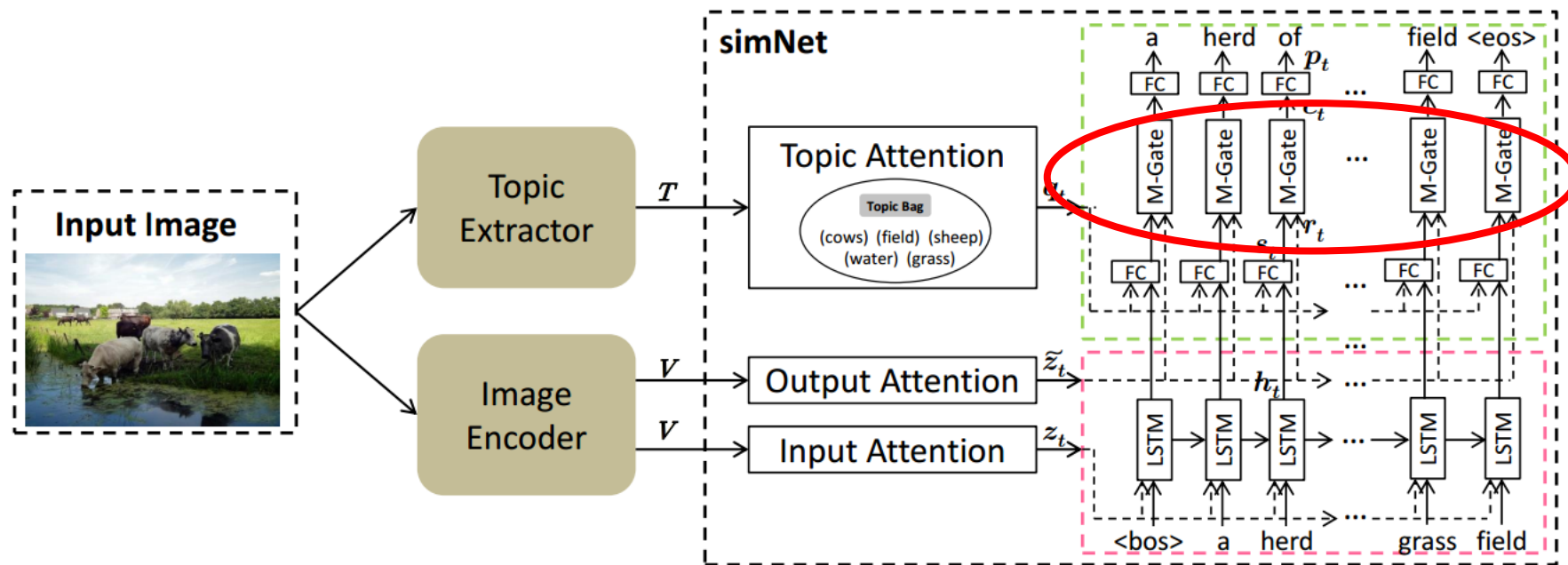


北京邮电大学  
Beijing University of Posts and Telecommunications



北京大学  
PEKING UNIVERSITY

# Approach: Merging Gate



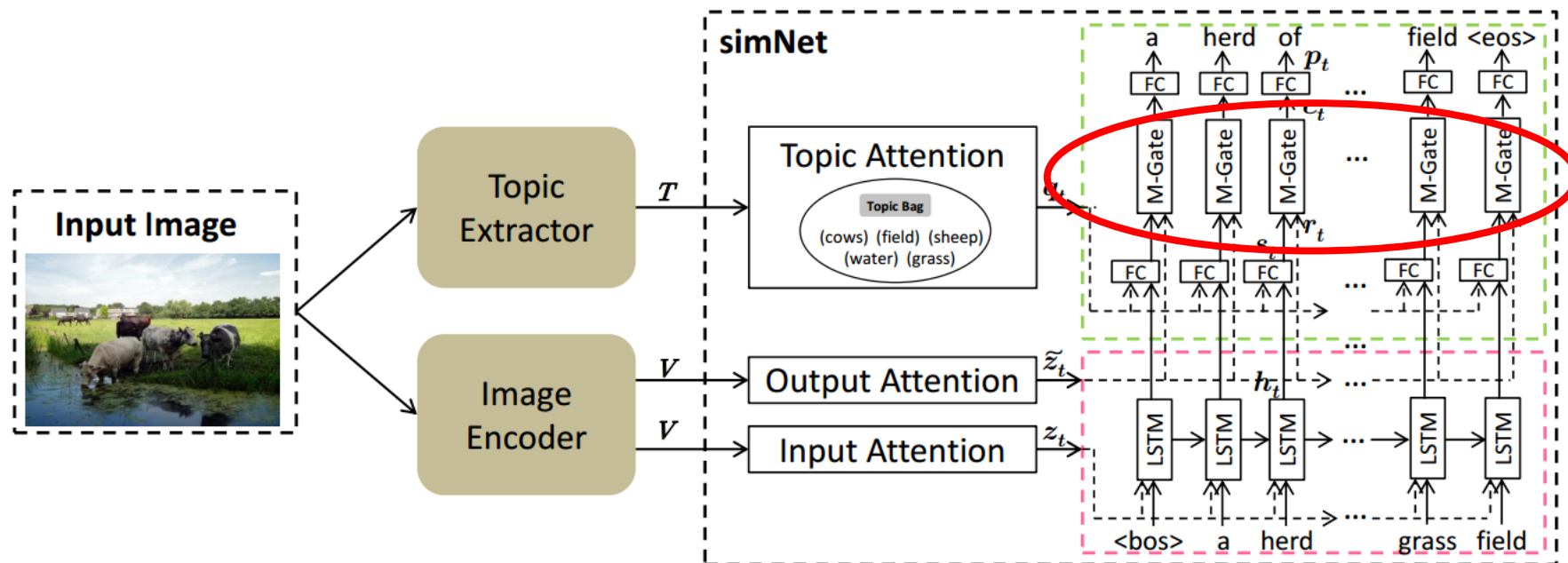
$$\gamma_t = \sigma(S(s_t) - S(r_t))$$

$$c_t = \gamma_t s_t + (1 - \gamma_t) r_t$$

( Where  $\sigma$  is the sigmoid function )



# Approach: Merging Gate



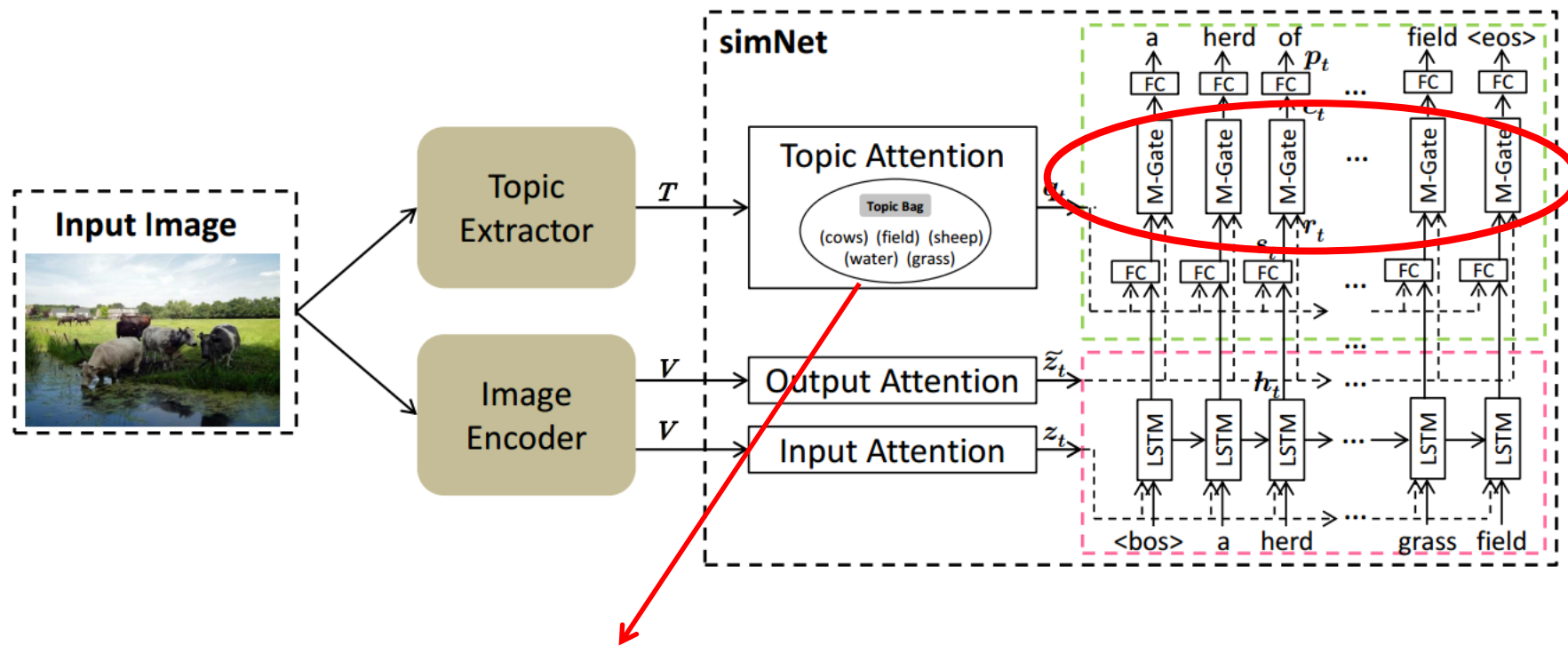
$$\gamma_t = \sigma(S(s_t) - S(r_t))$$

$$c_t = \gamma_t s_t + (1 - \gamma_t) r_t$$

The **scoring function  $S$**  is designed to evaluate the importance of the topic attention.



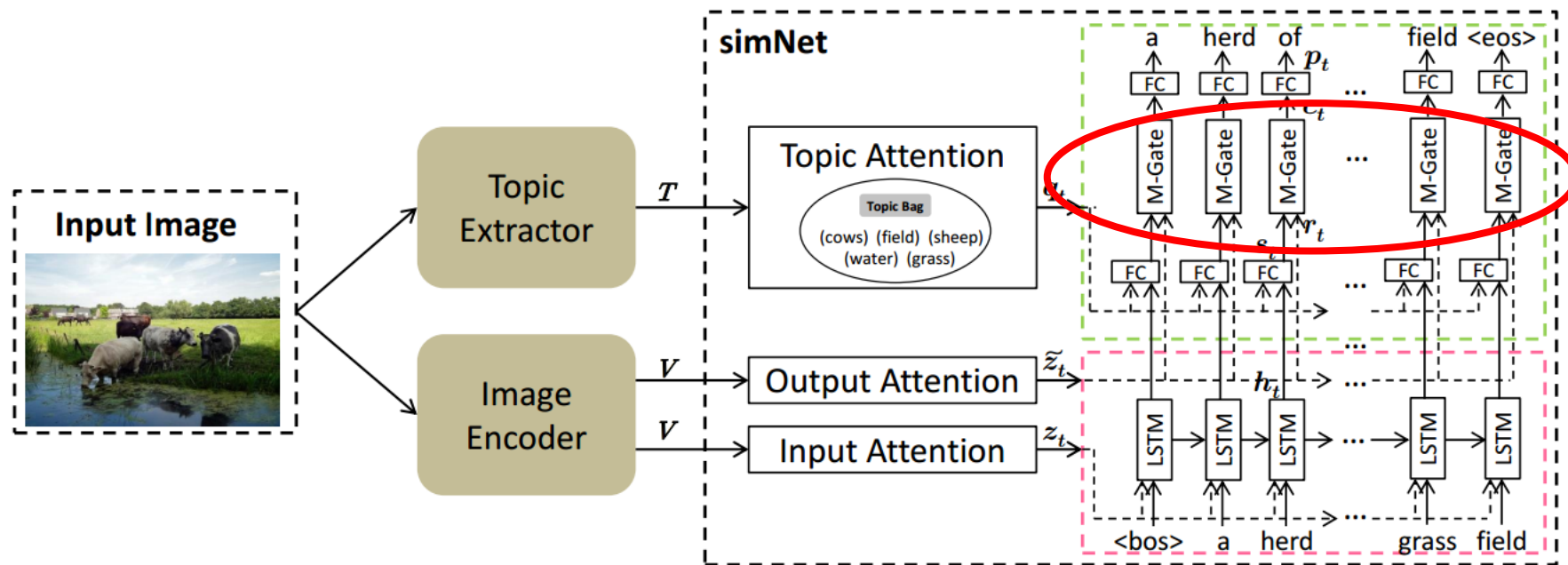
# Approach: Merging Gate



$$Q_t = \tanh(W^{Q,T}T \oplus W^{Q,h}h_t) \quad (10)$$

$$\beta_t = \text{softmax}(Q_t w^{\beta,Q}) \quad (11)$$

# Approach: Merging Gate



Share Weights

$$Q_t = \tanh(W^{Q,T}T \oplus W^{Q,h}h_t) \quad (10)$$

$$\beta_t = \text{softmax}(Q_t w^{\beta,Q}) \quad (11)$$

$$S(s_t) = \tanh(W^{S,h}h_t + W^{S,s}s_t) \cdot w^S \quad (16)$$

$$S(r_t) = \tanh(W^{S,h}h_t + W^{S,r}r_t) \cdot w^S \quad (17)$$

Share Weights



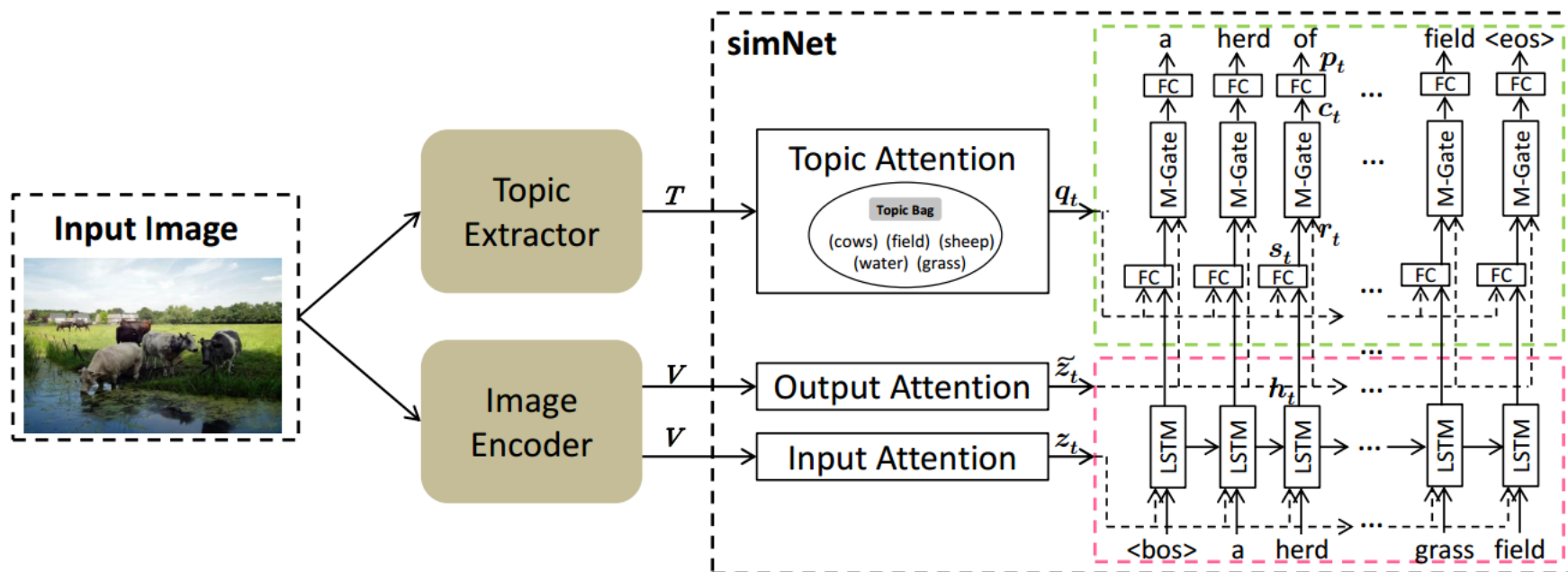
北京邮电大学 x



北京大学



# Generating Words



the contextual information:  $y_t \sim p_t = \text{softmax}(W^{p,c} c_t)$

# 3

## Experiments

---



北京邮电大学 X  
Beijing University of Posts and Telecommunications



北京大学  
PEKING UNIVERSITY

# Experiments

## Dataset

### Microsoft COCO(MSCOCO) and Flickr30k



- ✓ Sparrow bird on branch, with beak inspecting leaves on branch.
- ✓ A bird sitting on the branch of a tree near leaves.
- ✓ A bird that is sitting in a tree.
- ✓ a bird sitting on a branch of a tree.
- ✓ a bird that is on a small branch of a tree.

## Evaluation Metrics

- ✓ SPICE
- ✓ CIDEr
- ✓ BLEU
- ✓ METEOR
- ✓ ROUGE

Correlates the best with human judgments !



# Experiments: Results (MSCOCO)

	COCO	SPICE	CIDEr	METEOR	ROUGE-L	BLEU-4
Comparable Models	HardAtt (Xu et al., 2015)	-	-	0.230	-	0.250
	ATT-FCN (You et al., 2016)	-	-	0.243	-	0.304
	SCA-CNN (Chen et al., 2017)	-	0.952	0.250	0.531	0.311
	LSTM-A (Yao et al., 2017)	0.186	1.002	0.254	0.540	0.326
	SCN-LSTM (Gan et al., 2017)	-	1.012	0.257	-	0.330
	Skeleton (Wang et al., 2017)	-	1.069	0.268	0.552	0.336
	AdaAtt (Lu et al., 2017)	0.195	1.085	0.266	0.549	0.332
	NBT (Lu et al., 2018)	0.201	1.072	0.271	-	0.347
	DRL (Ren et al., 2017b)*	-	0.937	0.251	0.525	0.304
	TD-M-ATT (Chen et al., 2018)*	-	1.116	0.268	0.555	0.336
	SCST (Rennie et al., 2017)*	-	1.140	0.267	0.557	0.342
	SR-PL (Liu et al., 2018)* <sup>†</sup>	0.210	1.171	0.274	<b>0.570</b>	0.358
	Up-Down (Anderson et al., 2018)* <sup>†</sup>	0.214	<b>1.201</b>	0.277	0.569	<b>0.363</b>
	simNet	<b>0.220</b>	1.135	<b>0.283</b>	0.564	0.332



# Experiments: Results (MSCOCO)

COCO	SPICE	CIDEr	METEOR	ROUGE-L	BLEU-4
HardAtt (Xu et al., 2015)	-	-	0.230	-	0.250
ATT-FCN (You et al., 2016)	-	-	0.243	-	0.304
SCA-CNN (Chen et al., 2017)	-	0.952	0.250	0.531	0.311
LSTM-A (Yao et al., 2017)	0.186	1.002	0.254	0.540	0.326
SCN-LSTM (Gan et al., 2017)	-	1.012	0.257	-	0.330
Skeleton (Wang et al., 2017)	-	1.069	0.268	0.552	0.336
AdaAtt (Lu et al., 2017)	0.195	1.085	0.266	0.549	0.332
NBT (Lu et al., 2018)	0.201	1.072	0.271	-	0.347
DRL (Ren et al., 2017b)*	-	0.937	0.251	0.525	0.304
TD-M-ATT (Chen et al., 2018)*	-	1.116	0.268	0.555	0.336
SCST (Rennie et al., 2017)*	-	1.140	0.267	0.557	0.342
SR-PL (Liu et al., 2018)* <sup>†</sup>	0.210	1.171	0.274	<b>0.570</b>	0.358
Up-Down (Anderson et al., 2018)* <sup>†</sup>	0.214	<b>1.201</b>	0.277	0.569	<b>0.363</b>
simNet	<b>0.220</b>	1.135	<b>0.283</b>	0.564	0.332

Competitive



# 4

# Analysis

---



北京邮电大学 X  
Beijing University of Posts and Telecommunications



北京大学  
PEKING UNIVERSITY



# Analysis: The Contributions of The Sub-modules

Comprehensiveness

Detailedness

Methods	SPICE							CIDEr	METEOR	ROUGE-L	BLEU-4
	All	Objects	Attributes	Relations	Color	Count	Size				
Baseline (Plain Encoder-Decoder Network)	0.150	0.295	0.048	0.039	0.022	0.004	0.023	0.762	0.220	0.495	0.251
Up-Down (Anderson et al., 2018)* <sup>†</sup>	0.214	0.391	0.100	0.065	0.114	0.184	0.032	<b>1.201</b>	0.277	<b>0.569</b>	<b>0.363</b>
Baseline + Input Att.	0.164	0.316	0.060	0.044	0.030	0.038	0.024	0.840	0.233	0.512	0.273
Baseline + Output Att.	0.181	0.329	0.094	0.053	0.089	0.184	0.044	0.968	0.253	0.534	0.301
Baseline + Input Att. + Output Att.	0.187	0.338	0.101	0.055	<b>0.115</b>	0.161	<b>0.048</b>	1.038	0.259	0.542	0.311
Baseline + Topic Att.	0.184	0.348	0.074	0.051	0.047	0.064	0.037	0.915	0.250	0.517	0.260
Baseline + Topic Att. + MGate	0.189	0.355	0.080	0.051	0.055	0.090	0.033	0.959	0.256	0.527	0.281
Baseline + Input Att. + Output Att. + Topic Att.	0.206	0.381	0.091	0.060	0.075	0.094	0.045	1.068	0.273	0.556	0.320
simNet (Full Model)	<b>0.220</b>	<b>0.394</b>	<b>0.109</b>	<b>0.070</b>	0.088	<b>0.202</b>	0.045	1.135	<b>0.283</b>	0.564	0.332



# Analysis: Output Attention

The output attention is much more effective than the input attention

Methods	SPICE							CIDEr	METEOR	ROUGE-L	BLEU-4
	All	Objects	Attributes	Relations	Color	Count	Size				
Baseline (Plain Encoder-Decoder Network)	0.150	0.295	0.048	0.039	0.022	0.004	0.023	0.762	0.220	0.495	0.251
Up-Down (Anderson et al., 2018)* <sup>†</sup>	0.214	0.391	0.100	0.065	0.114	0.184	0.032	<b>1.201</b>	0.277	<b>0.569</b>	<b>0.363</b>
Baseline + Input Att.	0.164	0.316	0.060	0.044	0.030	0.038	0.024	0.840	0.233	0.512	0.273
Baseline + Output Att.	0.181	0.329	0.094	0.053	0.089	0.184	0.044	0.968	0.253	0.534	0.301
Baseline + Input Att. + Output Att.	0.187	0.338	0.101	0.055	<b>0.115</b>	0.161	<b>0.048</b>	1.038	0.259	0.542	0.311
Baseline + Topic Att.	0.184	0.348	0.074	0.051	0.047	0.064	0.037	0.915	0.250	0.517	0.260
Baseline + Topic Att. + MGate	0.189	0.355	0.080	0.051	0.055	0.090	0.033	0.959	0.256	0.527	0.281
Baseline + Input Att. + Output Att. + Topic Att.	0.206	0.381	0.091	0.060	0.075	0.094	0.045	1.068	0.273	0.556	0.320
simNet (Full Model)	<b>0.220</b>	<b>0.394</b>	<b>0.109</b>	<b>0.070</b>	0.088	<b>0.202</b>	0.045	1.135	<b>0.283</b>	0.564	0.332



# Analysis: Visual Attention

A combination of the input attention and the output attention makes the results even better

Methods	SPICE							CIDEr	METEOR	ROUGE-L	BLEU-4
	All	Objects	Attributes	Relations	Color	Count	Size				
Baseline (Plain Encoder-Decoder Network)	0.150	0.295	0.048	0.039	0.022	0.004	0.023	0.762	0.220	0.495	0.251
Up-Down (Anderson et al., 2018)* <sup>†</sup>	0.214	0.391	0.100	0.065	0.114	0.184	0.032	<b>1.201</b>	0.277	<b>0.569</b>	<b>0.363</b>
Baseline + Input Att.	0.164	0.316	0.060	0.044	0.030	0.038	0.024	0.840	0.233	0.512	0.273
Baseline + Output Att.	0.181	0.329	0.094	0.053	0.089	0.184	0.044	0.968	0.253	0.534	0.301
Baseline + Input Att. + Output Att.	0.187	0.338	0.101	0.055	<b>0.115</b>	0.161	<b>0.048</b>	1.038	0.259	0.542	0.311
Baseline + Topic Att.	0.184	0.348	0.074	0.051	0.047	0.064	0.037	0.915	0.250	0.517	0.260
Baseline + Topic Att. + MGate	0.189	0.355	0.080	0.051	0.055	0.090	0.033	0.959	0.256	0.527	0.281
Baseline + Input Att. + Output Att. + Topic Att.	0.206	0.381	0.091	0.060	0.075	0.094	0.045	1.068	0.273	0.556	0.320
simNet (Full Model)	<b>0.220</b>	<b>0.394</b>	<b>0.109</b>	<b>0.070</b>	0.088	<b>0.202</b>	0.045	1.135	<b>0.283</b>	0.564	0.332



# Analysis: Topic Attention

The topic attention is better at identifying objects but worse at identifying attributes.

Methods	SPICE							CIDEr	METEOR	ROUGE-L	BLEU-4
	All	Objects	Attributes	Relations	Color	Count	Size				
Baseline (Plain Encoder-Decoder Network)	0.150	0.295	0.048	0.039	0.022	0.004	0.023	0.762	0.220	0.495	0.251
Up-Down (Anderson et al., 2018)* <sup>†</sup>	0.214	0.391	0.100	0.065	0.114	0.184	0.032	<b>1.201</b>	0.277	<b>0.569</b>	<b>0.363</b>
Baseline + Input Att.	0.164	0.316	0.060	0.044	0.030	0.038	0.024	0.840	0.233	0.512	0.273
Baseline + Output Att.	0.181	0.329	0.094	0.053	0.089	0.184	0.044	0.968	0.253	0.534	0.301
Baseline + Input Att. + Output Att.	0.187	0.338	0.101	0.055	<b>0.115</b>	0.161	<b>0.048</b>	1.038	0.259	0.542	0.311
Baseline + Topic Att.	0.184	0.348	0.074	0.051	0.047	0.064	0.037	0.915	0.250	0.517	0.260
Baseline + Topic Att. + MGate	0.189	0.355	0.080	0.051	0.055	0.090	0.033	0.959	0.256	0.527	0.281
Baseline + Input Att. + Output Att. + Topic Att.	0.206	0.381	0.091	0.060	0.075	0.094	0.045	1.068	0.273	0.556	0.320
simNet (Full Model)	<b>0.220</b>	<b>0.394</b>	<b>0.109</b>	<b>0.070</b>	0.088	<b>0.202</b>	0.045	1.135	<b>0.283</b>	0.564	0.332



# Analysis: Visual Attention + Topic Attention

Combining the visual attention and the topic attention directly results in a huge boost in performance

Methods	SPICE							CIDEr	METEOR	ROUGE-L	BLEU-4
	All	Objects	Attributes	Relations	Color	Count	Size				
Baseline (Plain Encoder-Decoder Network)	0.150	0.295	0.048	0.039	0.022	0.004	0.023	0.762	0.220	0.495	0.251
Up-Down (Anderson et al., 2018)* <sup>†</sup>	0.214	0.391	0.100	0.065	0.114	0.184	0.032	<b>1.201</b>	0.277	<b>0.569</b>	<b>0.363</b>
Baseline + Input Att.	0.164	0.316	0.060	0.044	0.030	0.038	0.024	0.840	0.233	0.512	0.273
Baseline + Output Att.	0.181	0.329	0.094	0.053	0.089	0.184	0.044	0.968	0.253	0.534	0.301
Baseline + Input Att. + Output Att.	0.187	0.338	0.101	0.055	<b>0.115</b>	0.161	<b>0.048</b>	1.038	0.259	0.542	0.311
Baseline + Topic Att.	0.184	0.348	0.074	0.051	0.047	0.064	0.037	0.915	0.250	0.517	0.260
Baseline + Topic Att. + MGate	0.189	0.355	0.080	0.051	0.055	0.090	0.033	0.959	0.256	0.527	0.281
Baseline + Input Att. + Output Att. + Topic Att.	0.206	0.381	0.091	0.060	0.075	0.094	0.045	1.068	0.273	0.556	0.320
simNet (Full Model)	<b>0.220</b>	<b>0.394</b>	<b>0.109</b>	<b>0.070</b>	0.088	<b>0.202</b>	0.045	1.135	<b>0.283</b>	0.564	0.332





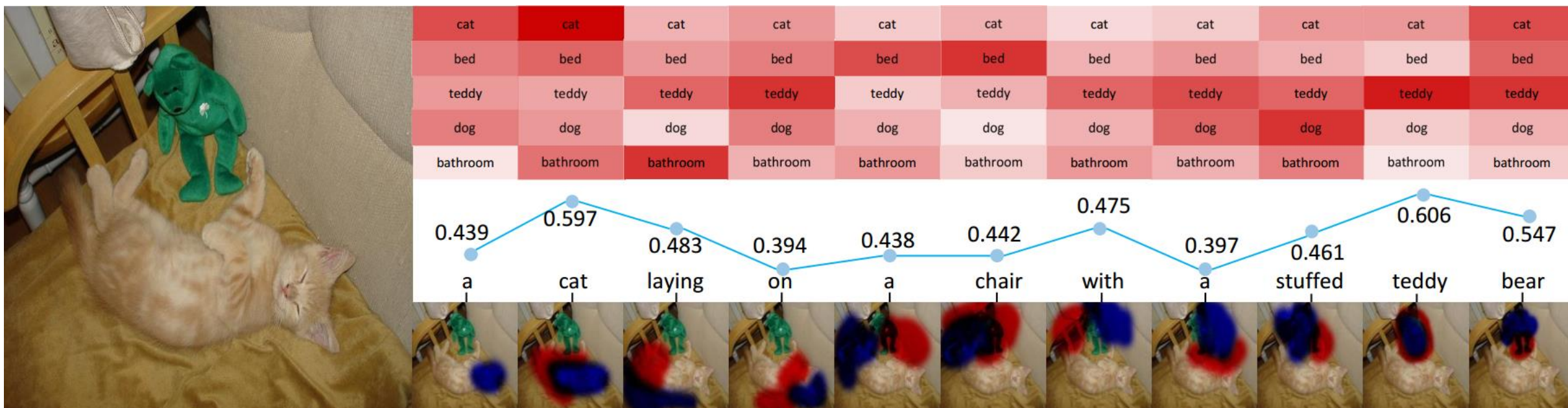
# Analysis: Full Model

Applying the merging gate is essential to the overall performance.

Methods	SPICE							CIDEr	METEOR	ROUGE-L	BLEU-4
	All	Objects	Attributes	Relations	Color	Count	Size				
Baseline (Plain Encoder-Decoder Network)	0.150	0.295	0.048	0.039	0.022	0.004	0.023	0.762	0.220	0.495	0.251
Up-Down (Anderson et al., 2018)* <sup>†</sup>	0.214	0.391	0.100	0.065	0.114	0.184	0.032	<b>1.201</b>	0.277	<b>0.569</b>	<b>0.363</b>
Baseline + Input Att.	0.164	0.316	0.060	0.044	0.030	0.038	0.024	0.840	0.233	0.512	0.273
Baseline + Output Att.	0.181	0.329	0.094	0.053	0.089	0.184	0.044	0.968	0.253	0.534	0.301
Baseline + Input Att. + Output Att.	0.187	0.338	0.101	0.055	<b>0.115</b>	0.161	<b>0.048</b>	1.038	0.259	0.542	0.311
Baseline + Topic Att.	0.184	0.348	0.074	0.051	0.047	0.064	0.037	0.915	0.250	0.517	0.260
Baseline + Topic Att. + MGate	0.189	0.355	0.080	0.051	0.055	0.090	0.033	0.959	0.256	0.527	0.281
Baseline + Input Att. + Output Att. + Topic Att.	0.206	0.381	0.091	0.060	0.075	0.094	0.045	1.068	0.273	0.556	0.320
simNet (Full Model)	<b>0.220</b>	<b>0.394</b>	<b>0.109</b>	<b>0.070</b>	0.088	<b>0.202</b>	0.045	1.135	<b>0.283</b>	0.564	0.332



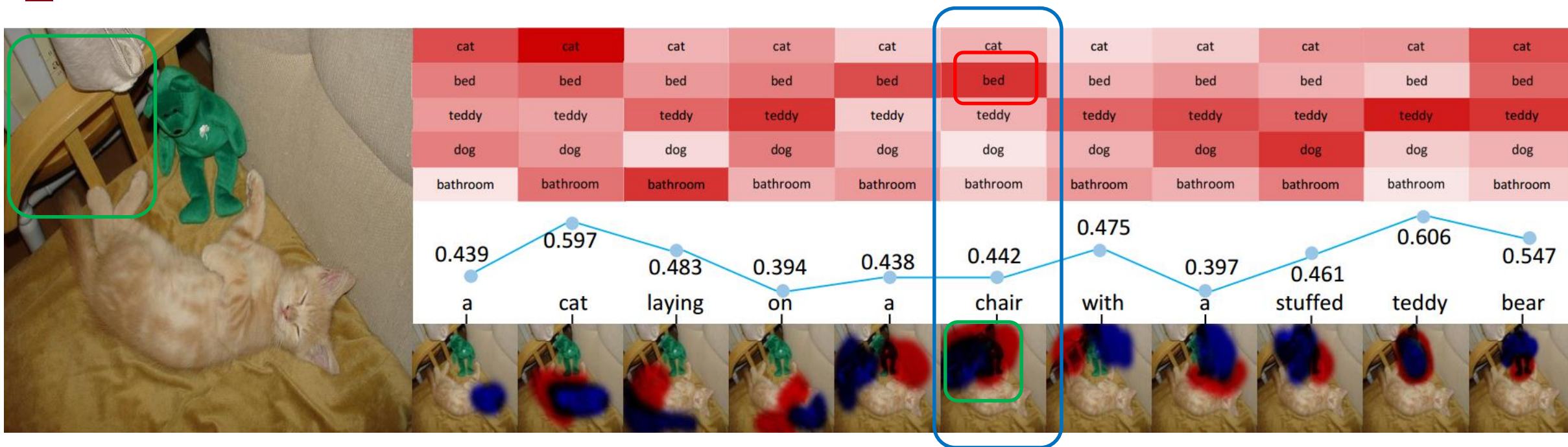
# Analysis: Visualization



- The upper part shows the attention weights of each of 5 extracted topics. —————> Deeper color means larger in value.
- The middle part shows the value of the merging gate. —————> Determines the importance of the topic attention.
- The lower part shows the visualization of visual attention. —————> The blue shade indicates the output attention. The red shade indicates the input attention.



# Analysis: Visualization



Visual information “*chair*” is **more important** than contextual information “*bed*”

# Analysis: Examples

## Comparison of Models



### Topics

woman girl  
baby bear  
kitchen

computer  
keyboard  
laptop mouse  
desk

pizza cheese  
table plate  
toppings

### Visual Attention

a girl  
and a baby  
are holding a  
stuffed animal

a computer ke  
yboard sitting  
on top of a  
wooden desk

two pizzas  
with toppings  
on a table

### Topic Attention

a woman  
holding a  
teddy bear  
in a kitchen

a computer  
keyboard and a  
mouse sitting  
on a desk

a pizza with  
a lot of  
toppings on it

### simNet

a woman  
and a baby  
are holding a  
stuffed animal

a computer  
keyboard and  
mouse on a  
wooden desk

two pizzas sitting  
on a table with  
two different ki  
nds of toppings

→ erroneous  
topic "kitchen"

→ lacking "mouse"

→ missing "wooden"

→ error count



北京  
Beijing University of P



北京大学  
PEKING UNIVERSITY

# Conclusion

- Stepwise image-topic merging network can adaptively combine the visual and the semantic attention to achieve substantial improvements.
- The generated captions are both detailed and comprehensive
- Our approach outperforms previous works in terms of SPICE on COCO and Flickr datasets.



---

# Thank you!

If you have any questions about our paper, you can send a email to [lfl@bupt.edu.cn](mailto:lfl@bupt.edu.cn)



北京邮电大学  
Beijing University of Posts and Telecommunications

x



北京大学  
PEKING UNIVERSITY

| 43

# Experiments: Results (Flickr30k)

Flickr30k	SPICE	CIDEr	METEOR	ROUGE-L	BLEU-4
HardAtt (Xu et al., 2015)	-	-	0.185	-	0.199
SCA-CNN (Chen et al., 2017)	-	-	0.195	-	0.223
ATT-FCN (You et al., 2016)	-	-	0.189	-	0.230
SCN-LSTM (Gan et al., 2017)	-	-	0.210	-	0.257
AdaAtt (Lu et al., 2017)	0.145	0.531	0.204	0.467	0.251
NBT (Lu et al., 2018)	0.156	0.575	0.217	-	0.271
SR-PL (Liu et al., 2018)* <sup>†</sup>	0.158	<b>0.650</b>	0.218	<b>0.499</b>	<b>0.293</b>
simNet	<b>0.160</b>	0.585	<b>0.221</b>	0.489	0.251

Table 1: Performance on the Flickr30k Karpathy test split. The symbol \* denotes directly optimizing CIDEr. The symbol <sup>†</sup> denotes using extra data for training, thus not directly comparable. Nonetheless, our model supersedes all existing models in SPICE, which correlates the best with human judgments.





# Analysis: Topic Extraction

Method	Precision	Recall	F1
Topics ( $m=5$ )	49.95	38.91	42.48
All words ( $m=5$ )	<b>84.01</b>	17.99	29.49
All words ( $m=10$ )	70.90	30.18	42.05
All words ( $m=20$ )	52.51	<b>44.53</b>	<b>47.80</b>

Table 4: Performance of visual word extraction.

The reason of using objects as topics is that they are easier to identify so that the generation suffers less from erroneous predictions.

Method	S	C	M	R	B
Topics ( $m=5$ )	<b>0.220</b>	<b>1.135</b>	<b>0.283</b>	<b>0.564</b>	<b>0.332</b>
All words ( $m=5$ )	0.197	1.047	0.264	0.550	0.314
All words ( $m=10$ )	0.201	1.076	0.256	0.528	0.293
All words ( $m=20$ )	0.209	1.117	0.276	0.561	0.329

Table 5: Effect of using different visual words.

It proves that for semantic attention, it is also important to limit the absolute number of incorrect visual words instead of merely the precision or the recall.



# Analysis: Merging Gate

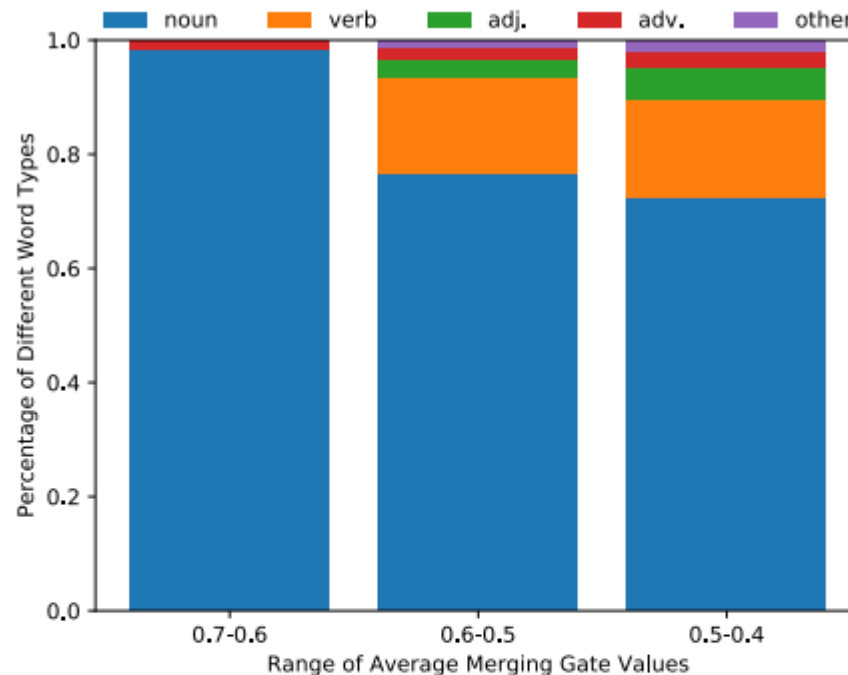





Figure 4: Average merging gate values according to word types. As we can see, object words (noun) dominate the high value range, while attribute and relation words are assigned lower values, indicating the merging gate learns to efficiently combine the information.








# Analysis: Error Analysis

<b>Error Analysis</b>			
<b>Topics</b>	clock tower building street city	people bus truck street train	garden bench park forest plants
<b>Reference</b>	a tall building that has a clock on it (near a large building)	tour buses driving down a street lined with cheering people	an old wooden bench in nature surrounded by plants
<b>simNet</b>	a large building with a clock tower in the background	a group of people standing around a parked bus at a bus stop	a wooden bench sitting in the middle of a lush green garden
<b>Error Type</b>	distance	movement	object

There are mainly three types of errors, i.e. **distance** (32, 26%), **movement** (22, 18%), and **object** (60, 49%), with 9 (7%) other errors.



# Analysis: Error Analysis

<b>Error Analysis</b>			
<b>Topics</b>	clock tower building street city	people bus truck street train	<u>garden</u> bench park forest plants
<b>Reference</b>	a tall building that has a clock on it (near a large building)	tour buses <u>driving down</u> a street lined with cheering people	an old woo den bench <u>in nature</u> surrounded by plants
<b>simNet</b>	a large building with a clock tower in the <u>background</u>	a group of people standing around a <u>parked</u> bus at a bus stop	a wooden bench sitting in the middle of a lush green <u>garden</u>
<b>Error Type</b>	distance	movement	object

There are mainly three types of errors, i.e. **distance** (32, 26%), **movement** (22, 18%), and **object** (60, 49%), with 9 (7%) other errors.



北京邮电大学  
Beijing University of Posts and Telecommunications



北京大学  
PEKING UNIVERSITY



# Experiments: Results (MSCOCO)

COCO	BLEU-1		BLEU-2		BLEU-3		BLEU-4		METEOR		ROUGE-L		CIDEr	
	c5	c40	c5	c40	c5	c40	c5	c40	c5	c40	c5	c40	c5	c40
HardAtt (Xu et al., 2015)	0.705	0.881	0.528	0.779	0.383	0.658	0.277	0.537	0.241	0.322	0.516	0.654	0.865	0.893
ATT-FCN (You et al., 2016)	0.731	0.900	0.565	0.815	0.424	0.709	0.316	0.599	0.250	0.335	0.535	0.682	0.943	0.958
SCA-CNN (Chen et al., 2017)	0.712	0.894	0.542	0.802	0.404	0.691	0.302	0.579	0.244	0.331	0.524	0.674	0.912	0.921
LSTM-A (Yao et al., 2017)	0.739	0.919	0.575	0.842	0.436	0.740	0.330	0.632	0.256	0.350	0.542	0.700	0.984	1.003
SCN-LSTM (Gan et al., 2017)	0.740	0.917	0.575	0.839	0.436	0.739	0.331	0.631	0.257	0.348	0.543	0.696	1.003	1.013
AdaAtt (Lu et al., 2017) <sup>†</sup>	0.748	0.920	0.584	0.845	0.444	0.744	0.336	0.637	0.264	0.359	0.550	0.705	1.042	1.059
TD-M-ATT (Chen et al., 2018) <sup>*†</sup>	0.757	0.913	0.591	0.836	0.441	0.726	0.324	0.609	0.259	0.342	0.547	0.689	1.059	1.090
SCST (Rennie et al., 2017) <sup>*†</sup>	0.781	0.937	0.619	0.860	0.470	0.759	0.352	0.645	0.270	0.355	0.563	0.707	1.147	1.167
Up-Down (Anderson et al., 2018) <sup>*†‡</sup>	<b>0.802</b>	<b>0.952</b>	<b>0.641</b>	<b>0.888</b>	<b>0.491</b>	<b>0.794</b>	<b>0.369</b>	<b>0.685</b>	<b>0.276</b>	<b>0.367</b>	<b>0.571</b>	<b>0.724</b>	<b>1.179</b>	<b>1.205</b>
simNet	0.766	0.941	0.605	0.874	0.462	0.778	0.350	0.671	0.267	0.362	0.558	0.716	1.087	1.111

Table 6: Performance on the online COCO evaluation server. The SPICE metric is unavailable for our model, thus not reported. c5 means evaluating against 5 references, and c40 means evaluating against 40 references. The symbol \* denotes directly optimizing CIDEr. The symbol <sup>†</sup> denotes model ensemble. The symbol <sup>‡</sup> denotes using extra data for training, thus not directly comparable. Our submission does not use the three aforementioned techniques. Nonetheless, our model is second only to Up-Down and surpasses almost all the other models in published work, especially when 40 references are considered.

# Analysis: Topic Attention

Methods	SPICE							CIDEr	METEOR	ROUGE-L	BLEU-4
	All	Objects	Attributes	Relations	Color	Count	Size				
Baseline (Plain Encoder-Decoder Network)	0.150	0.295	0.048	0.039	0.022	0.004	0.023	0.762	0.220	0.495	0.251
Up-Down (Anderson et al., 2018)* <sup>†</sup>	0.214	0.391	0.100	0.065	0.114	0.184	0.032	<b>1.201</b>	0.277	<b>0.569</b>	<b>0.363</b>
Baseline + Input Att.	0.164	0.316	0.060	0.044	0.030	0.038	0.024	0.840	0.233	0.512	0.273
Baseline + Output Att.	0.181	0.329	0.094	0.053	0.089	0.184	0.044	0.968	0.253	0.534	0.301
Baseline + Input Att. + Output Att.	0.187	0.338	0.101	0.055	<b>0.115</b>	0.161	<b>0.048</b>	1.038	0.259	0.542	0.311
Baseline + Topic Att.	0.184	0.348	0.074	0.051	0.047	0.064	0.037	0.915	0.250	0.517	0.260
Baseline + Topic Att. + MGate	0.189	0.355	0.080	0.051	0.055	0.090	0.033	0.959	0.256	0.527	0.281
Baseline + Input Att. + Output Att. + Topic Att.	0.206	0.381	0.091	0.060	0.075	0.094	0.045	1.068	0.273	0.556	0.320
simNet (Full Model)	<b>0.220</b>	<b>0.394</b>	<b>0.109</b>	<b>0.070</b>	0.088	<b>0.202</b>	0.045	1.135	<b>0.283</b>	0.564	0.332

