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

*Fenglin Liu*[1], Xuancheng Ren[2]*, Yuanxin Liu[1], Houfeng Wang[2] and Xu Sun[2]

[1] Beijing University of Posts and Telecommunications, Beijing, China

[2] Peking University, Beijing, China

* Equal Contributions

# CONTENTS

北京邮电大学
Beijing University of Posts and Telecommunications

X

北京大学
PEKING UNIVERSITY

# 1 Introduction

# 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

# Introduction: Soft-Attention

**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

omitting "dog" and "mouse"

encode    decode

**Soft-Attention:**    Image    →    Image Features    →    Caption

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

# 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

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

missing "open" and mislocating "dog"

ATT-FCN:   Image → (encode) → Image Keywords → (decode) → Caption

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

# 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

**simNet:** Image → encode → Image Features / Image Keywords → decode → Caption
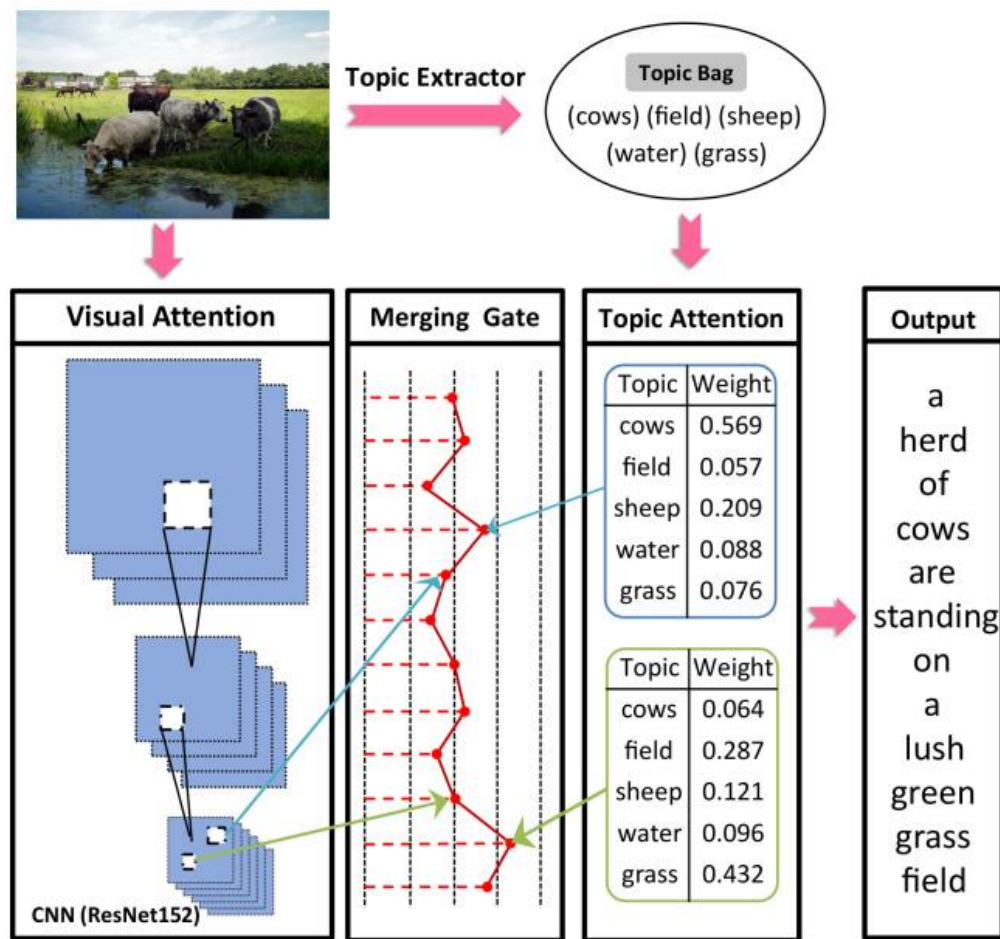
# Introduction: Main idea



**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

# Overview
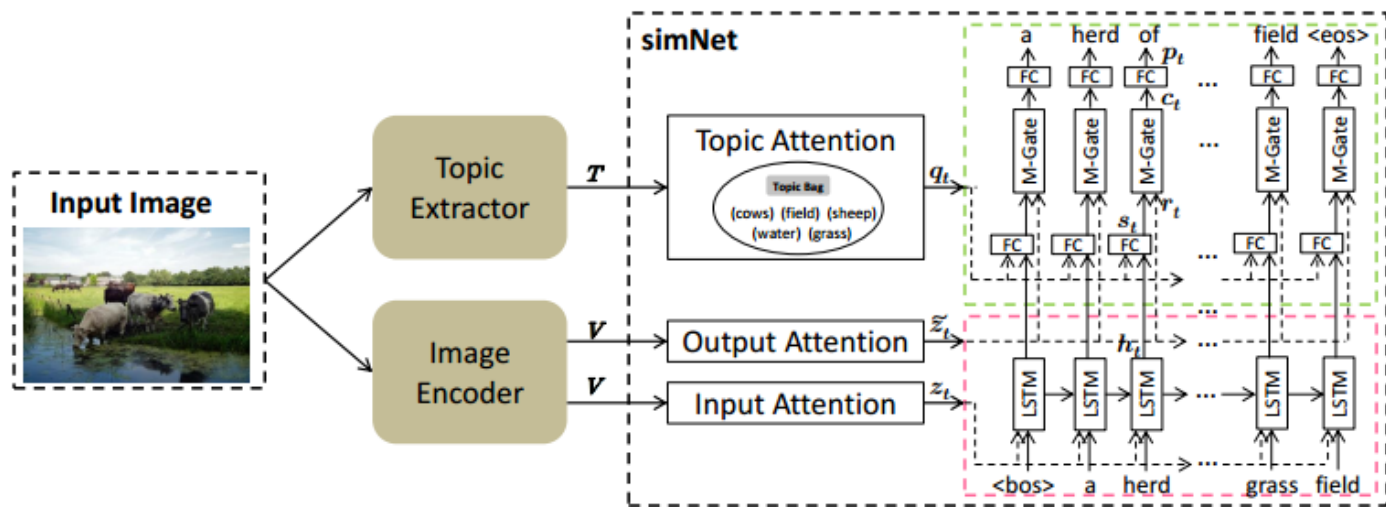


(a) The overall framework.

(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.

# Approach: Image Encoder



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

# Approach: Input Attention



Input Attention:

$$\boldsymbol{Z_t} = \tanh(\boldsymbol{W}^{Z,V}\boldsymbol{V} \oplus \boldsymbol{W}^{Z,h}\boldsymbol{h}_{t-1}) \qquad (2)$$

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

$$\boldsymbol{z_t} = \boldsymbol{V}\boldsymbol{\alpha_t} \qquad (4)$$

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

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

# Approach: Output Attention



Output Attention:

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

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

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

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

# Approach: Visual Information



Output Attention:

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

$$\widetilde{\alpha}_t = \mathrm{softmax}(\widetilde{Z}_t \widetilde{w}^{\alpha,Z}) \qquad (7)$$

$$\widetilde{z}_t = V \widetilde{\alpha}_t \qquad (8)$$

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

# Approach: Previous Topic Attention



Topic Attention (Previous work):

$$\boldsymbol{\beta}_t = \mathrm{softmax}(\boldsymbol{T}^\top \boldsymbol{U} \boldsymbol{y}_{t-1}) \qquad (9)$$

**Lacking the attentive visual information when selecting topic!**

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

# Approach: Our Topic Attention



Topic Attention (Our):

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

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

$$\boldsymbol{q}_t = \boldsymbol{T}\beta_t \quad (12)$$

# Approach: Contextual Information



Topic Attention (Our):

$$Q_t = \tanh(\boldsymbol{W}^{Q,T}\boldsymbol{T} \oplus \boldsymbol{W}^{Q,h}\boldsymbol{h}_t) \qquad (10)$$

$$\beta_t = \mathrm{softmax}(\boldsymbol{Q}_t\boldsymbol{w}^{\beta,Q}) \qquad (11)$$

$$\boldsymbol{q}_t = \boldsymbol{T}\boldsymbol{\beta}_t \qquad (12)$$

the contextual information: $\boldsymbol{s}_t = \tanh(\boldsymbol{W}^{s,q}\boldsymbol{q}_t + \boldsymbol{W}^{s,h}\boldsymbol{h}_t)$

# 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)

# Approach: Merging Gate



$$\gamma_t = \sigma(S(\boldsymbol{s}_t) - S(\boldsymbol{r}_t))$$

$$\boldsymbol{c}_t = \gamma_t \boldsymbol{s}_t + (1 - \gamma_t)\boldsymbol{r}_t$$

( Where $\sigma$ is the sigmoid function )

# Approach: Merging Gate



$$\gamma_t = \sigma(S(\boldsymbol{s}_t) - S(\boldsymbol{r}_t))$$

$$\boldsymbol{c}_t = \gamma_t \boldsymbol{s}_t + (1 - \gamma_t)\boldsymbol{r}_t$$

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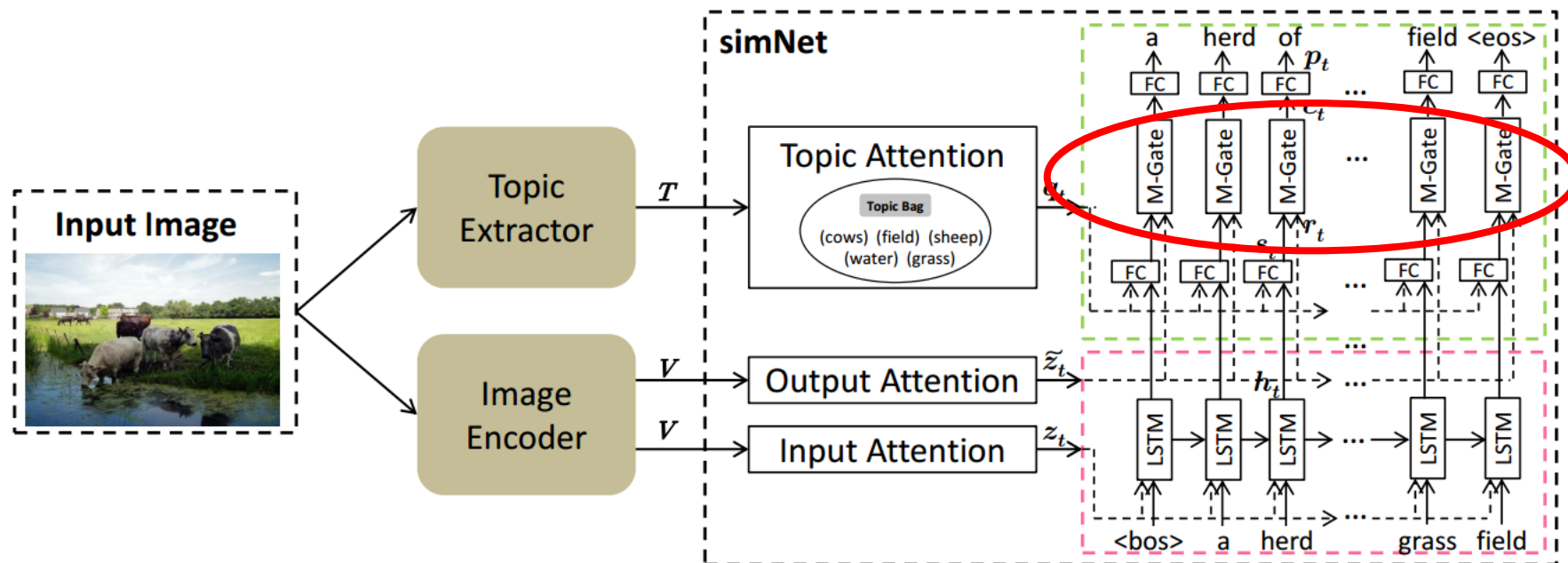the importance of the topic attention.

# Approach: Merging Gate



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

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

# Approach: Merging Gate



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

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

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

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

**Share Weights**

**Share Weights**

# Generating Words



the contextual information: $y_t \sim \boldsymbol{p}_t = \mathrm{softmax}(\boldsymbol{W}^{p,c}\boldsymbol{c}_t)$

# 3

# Experiments

# Experiments

Dataset

**Microsoft COCO(MSCOCO) and Flickr30k**

Evaluation Metrics



- ✓ 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.

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

Correlates the best with human Judgments !

# 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)[*†] | 0.210 | 1.171 | 0.274 | **0.570** | 0.358 |
| Up-Down (Anderson et al., 2018)[*†] | 0.214 | **1.201** | 0.277 | 0.569 | **0.363** |
| simNet | **0.220** | 1.135 | **0.283** | 0.564 | 0.332 |

**Comparable Models**

# 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)*† | 0.210 | 1.171 | 0.274 | **0.570** | 0.358 |
| Up-Down (Anderson et al., 2018)*† | 0.214 | **1.201** | 0.277 | 0.569 | **0.363** |
| simNet | **0.220** | 1.135 | **0.283** | 0.564 | 0.332 |

Competitive

# 4

Analysis

# 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)*† | 0.214 | 0.391 | 0.100 | 0.065 | 0.114 | 0.184 | 0.032 | **1.201** | 0.277 | **0.569** | **0.363** |
| 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 | **0.115** | 0.161 | **0.048** | 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) | **0.220** | **0.394** | **0.109** | **0.070** | 0.088 | **0.202** | 0.045 | 1.135 | **0.283** | 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)[*†] | 0.214 | 0.391 | 0.100 | 0.065 | 0.114 | 0.184 | 0.032 | **1.201** | 0.277 | **0.569** | **0.363** |
| 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 | **0.115** | 0.161 | **0.048** | 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) | **0.220** | **0.394** | **0.109** | **0.070** | 0.088 | **0.202** | 0.045 | 1.135 | **0.283** | 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)* † | 0.214 | 0.391 | 0.100 | 0.065 | 0.114 | 0.184 | 0.032 | **1.201** | 0.277 | **0.569** | **0.363** |
| 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 | **0.115** | 0.161 | **0.048** | 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) | **0.220** | **0.394** | **0.109** | **0.070** | 0.088 | **0.202** | 0.045 | 1.135 | **0.283** | 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)[*][†] | 0.214 | 0.391 | 0.100 | 0.065 | 0.114 | 0.184 | 0.032 | **1.201** | 0.277 | **0.569** | **0.363** |
| 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 | **0.115** | 0.161 | **0.048** | 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) | **0.220** | **0.394** | **0.109** | **0.070** | 0.088 | **0.202** | 0.045 | 1.135 | **0.283** | 0.564 | 0.332 |

# Analysis: Visual Attention + Topic Attention

**Combing 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)[*][†] | 0.214 | 0.391 | 0.100 | 0.065 | 0.114 | 0.184 | 0.032 | **1.201** | 0.277 | **0.569** | **0.363** |
| 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 | **0.115** | 0.161 | **0.048** | 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) | **0.220** | **0.394** | **0.109** | **0.070** | 0.088 | **0.202** | 0.045 | 1.135 | **0.283** | 0.564 | 0.332 |

# Analysis: Full Model

**Applying the merging gate is essential to the overall performance.**

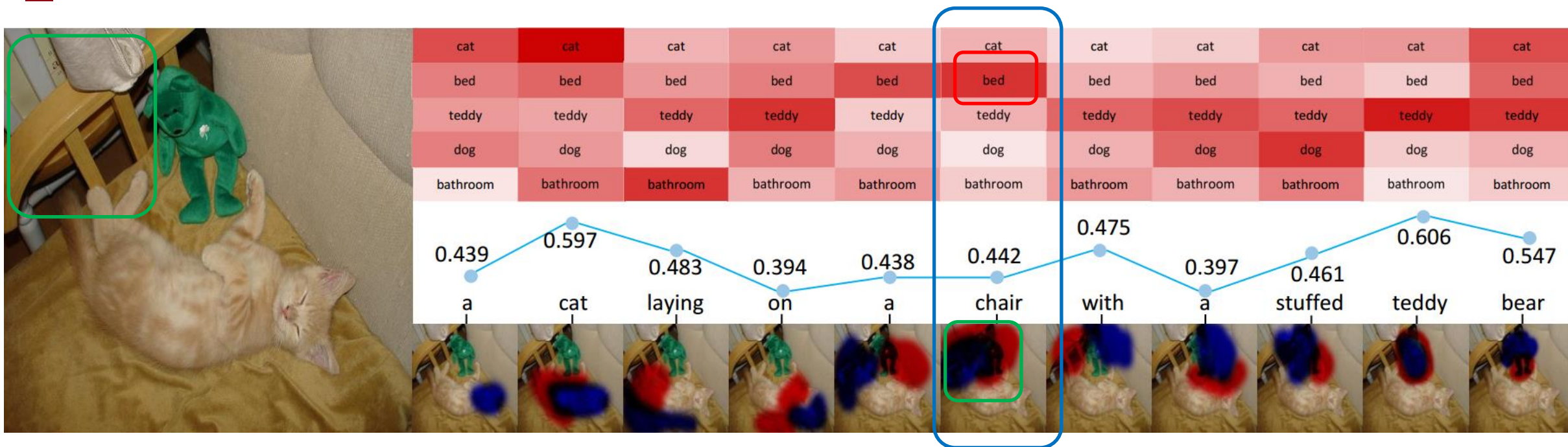| Methods | SPICE | | | | | | | CIDEr | METEOR | ROUGE-L | BLEU-4 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | All | Objects | Attributes | Relations | Color | Count | Size | | | | |
| Baseline (Plain Encoder-Decoder Network) | 0.150 | 0.295 | 0.048 | 0.039 | 0.022 | 0.004 | 0.023 | 0.762 | 0.220 | 0.495 | 0.251 |
| Up-Down (Anderson et al., 2018)[*][†] | 0.214 | 0.391 | 0.100 | 0.065 | 0.114 | 0.184 | 0.032 | **1.201** | 0.277 | **0.569** | **0.363** |
| 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 | **0.115** | 0.161 | **0.048** | 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) | **0.220** | **0.394** | **0.109** | **0.070** | 0.088 | **0.202** | 0.045 | 1.135 | **0.283** | 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 | Visual Attention | Topic Attention | simNet |
|---|---|---|---|
| woman girl baby bear kitchen | a girl and a baby are holding a stuffed animal | a woman holding a teddy bear in a kitchen → erroneous topic "kitchen" | a woman and a baby are holding a stuffed animal |
| computer keyboard laptop mouse desk | a computer keyboard sitting on top of a wooden desk → lacking "mouse" | a computer keyboard and a mouse sitting on a desk → missing "wooden" | a computer keyboard and mouse on a wooden desk |
| pizza cheese table plate toppings | two pizzas with toppings on a table | a pizza with a lot of toppings on it → error count | two pizzas sitting on a table with two different kinds of toppings |

# 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