

Doubly-Attentive Decoder for Multi-modal Neural Machine Translation

Iacer Calixto¹ and Qun Liu¹

August 14, 2018

¹ADAPT Centre

School of Computing

Dublin City University

{*FirstName.LastName*}@adaptcentre.ie



Introduction

Model Architecture

Data sets

Experiments

Conclusions

Introduction

Introduction

- **Machine Translation (MT)**: the task in which we wish to learn a model to translate text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - text-only task;
 - model is trained on parallel source/target sentence pairs.
- **Image description generation (IDG)**: the task in which we wish to learn a model to describe an image using natural language (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on image/target sentence pairs.

Introduction

- **Machine Translation (MT)**: the task in which we wish to learn a model to translate text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - text-only task;
 - model is trained on parallel source/target sentence pairs.
- **Image description generation (IDG)**: the task in which we wish to learn a model to describe an image using natural language (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on image/target sentence pairs.

Introduction

- **Machine Translation (MT)**: the task in which we wish to learn a model to translate text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - text-only task;
 - model is trained on parallel source/target sentence pairs.
- **Image description generation (IDG)**: the task in which we wish to learn a model to describe an image using natural language (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on image/target sentence pairs.

Introduction

- **Machine Translation (MT)**: the task in which we wish to learn a model to translate text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - text-only task;
 - model is trained on parallel source/target sentence pairs.
- **Image description generation (IDG)**: the task in which we wish to learn a model to describe an image using natural language (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on image/target sentence pairs.

Introduction

- **Machine Translation (MT)**: the task in which we wish to learn a model to translate text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - text-only task;
 - model is trained on parallel source/target sentence pairs.
- **Image description generation (IDG)**: the task in which we wish to learn a model to describe an image using natural language (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on image/target sentence pairs.

Introduction

- **Machine Translation (MT)**: the task in which we wish to learn a model to translate text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - text-only task;
 - model is trained on parallel source/target sentence pairs.
- **Image description generation (IDG)**: the task in which we wish to learn a model to describe an image using natural language (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on image/target sentence pairs.

- **Multi-Modal Machine Translation (MMT):** learn a model to translate text and an image that illustrates this text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on source/image/target triplets;
 - can be seen as a form of augmented MT or augmented image description generation.

- **Multi-Modal Machine Translation (MMT)**: learn a model to translate text and an image that illustrates this text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on source/image/target triplets;
 - can be seen as a form of augmented MT or augmented image description generation.

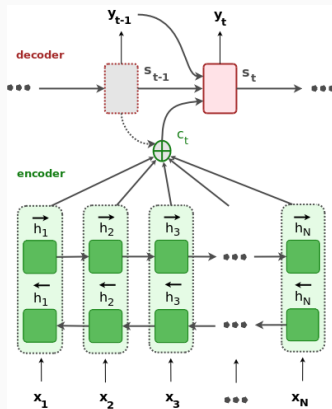
- **Multi-Modal Machine Translation (MMT)**: learn a model to translate text and an image that illustrates this text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on source/image/target triplets;
 - can be seen as a form of augmented MT or augmented image description generation.

- **Multi-Modal Machine Translation (MMT)**: learn a model to translate text and an image that illustrates this text from one natural language (e.g., English) into another (e.g., Brazilian Portuguese).
 - multi-modal task (text and vision);
 - model is trained on source/image/target triplets;
 - can be seen as a form of **augmented MT** or **augmented image description generation**.

Model Architecture

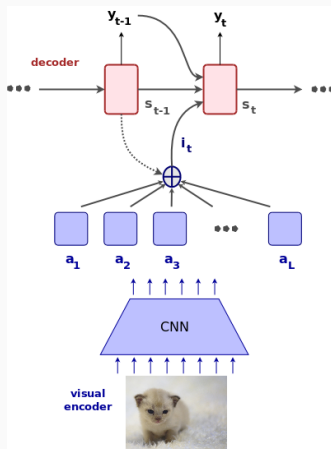
Attentional Neural Machine Translation

The **attention mechanism** lets the decoder **search for** the best source words to generate each target word, e.g. Bahdanau et al. (2015).



Attentional Neural Image Description Generation

The **attention mechanism** lets the decoder look at or **attend to** specific parts of the image when generating each target word, e.g. Xu et al. (2015).



Doubly-Attentive Multi-Modal NMT Model

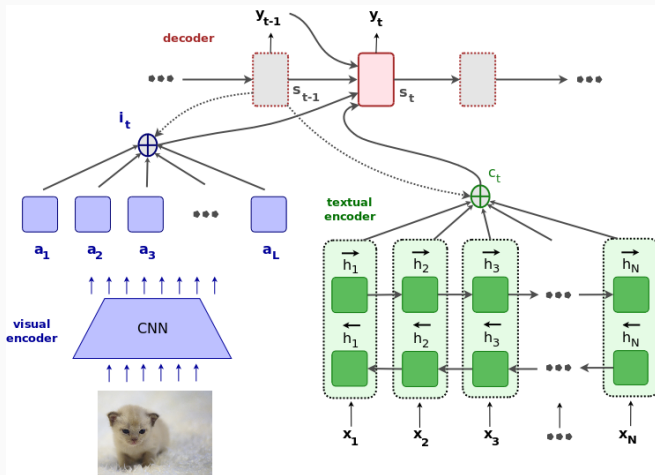


Figure 1: Doubly-Attentive Multi-modal NMT (paper accepted in ACL 2017)

image gating

Doubly-Attentive Multi-Modal NMT Model

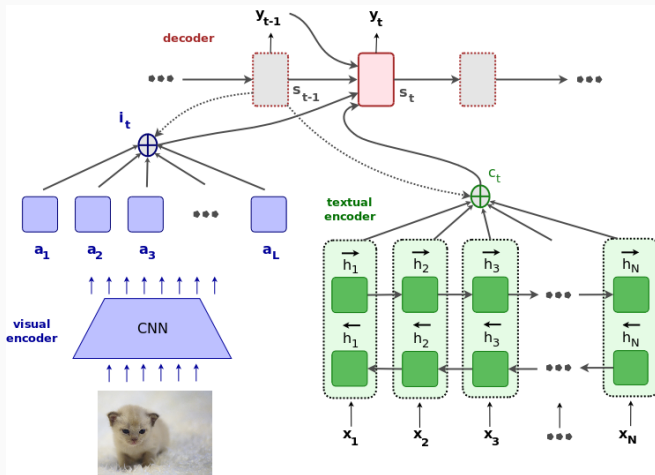


Figure 1: Doubly-Attentive Multi-modal NMT (paper accepted in ACL 2017)

image gating

Data sets

- **29K** images;
- **Translated Multi30k** – **29K** English–German **parallel descriptions** (1 per image);
- **Comparable Multi30k** – **145K** English and **145K** German **comparable descriptions** (5 English and 5 German per image);
- **~2M** English and **~1.6M** German words.

- **29K** images;
- **Translated Multi30k** – **29K** English–German **parallel descriptions** (1 per image);
- **Comparable Multi30k** – **145K** English and **145K** German **comparable descriptions** (5 English and 5 German per image);
- **~2M** English and **~1.6M** German words.

- **29K** images;
- **Translated Multi30k** – **29K** English–German **parallel descriptions** (1 per image);
- **Comparable Multi30k** – **145K** English and **145K** German **comparable descriptions** (5 English and 5 German per image);
- **~2M** English and **~1.6M** German words.

- **29K** images;
- **Translated Multi30k** – **29K** English–German **parallel descriptions** (1 per image);
- **Comparable Multi30k** – **145K** English and **145K** German **comparable descriptions** (5 English and 5 German per image);
- **~2M** English and **~1.6M** German words.

- Corpora:
 - Europarl;
 - News Commentary;
 - Common Crawl;
- **~4.3M sentence pairs** in total;
- **~103M** English and **~103M** German words.

- Corpora:
 - Europarl;
 - News Commentary;
 - Common Crawl;
- ~4.3M sentence pairs in total;
- ~103M English and ~103M German words.

- Corpora:
 - Europarl;
 - News Commentary;
 - Common Crawl;
- ~4.3M sentence pairs in total;
- ~103M English and ~103M German words.

- Corpora:
 - Europarl;
 - News Commentary;
 - Common Crawl;
- **~4.3M sentence pairs** in total;
- **~103M** English and **~103M** German words.

- Corpora:
 - Europarl;
 - News Commentary;
 - Common Crawl;
- **~4.3M sentence pairs** in total;
- **~103M** English and **~103M** German words.

- **Training:**
 - **Translated Multi30k** (29k triples);
- **Pre-training:**
 - **Translated Multi30k** (29k triples) +
back-translated Comparable Multi30k (145k triples);
 - **WMT 2015** (4.3M sentence pairs);

Experiments

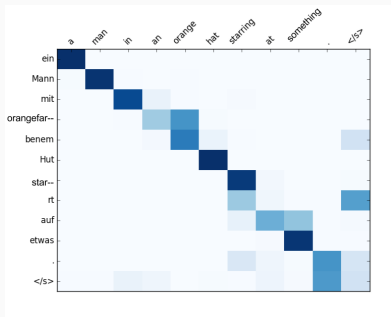
Model	Training data	BLEU4 \uparrow	METEOR \uparrow	TER \downarrow	chrF3 \uparrow (prec. / recall)	
NMT	M30k _T	<u>33.7</u>	52.3	46.7	65.2	(67.7 / 65.0)
PBSMT	M30k _T	32.9	<u>54.3</u> \dagger	<u>45.1</u> \dagger	67.4	(66.5 / 67.5)
Huang et al. (2016)	M30k _T	35.1 (\uparrow 1.4)	52.2 (\downarrow 2.1)	—	—	—
	+ RCNN	36.5 (\uparrow 2.8)	54.1 (\downarrow 0.2)	—	—	—
NMT _{SRC+IMG}	M30k _T	36.5 $\dagger\dagger$	55.0 \dagger	43.7 $\dagger\dagger$	67.3	(66.8 / 67.4)
Improvements						
NMT _{SRC+IMG} vs. NMT		\uparrow 2.8	\uparrow 2.7	\downarrow 3.0	\uparrow 2.1	\downarrow 0.9 / \uparrow 2.4
NMT _{SRC+IMG} vs. PBSMT		\uparrow 3.6	\uparrow 0.7	\downarrow 1.4	\downarrow 0.1	\uparrow 0.3 / \downarrow 0.1
NMT _{SRC+IMG} vs. Huang		\uparrow 1.4	\uparrow 2.8	—	—	—
NMT _{SRC+IMG} vs. Huang (+RCNN)		\uparrow 0.0	\uparrow 0.9	—	—	—
Pre-training data set: back-translated M30k_C (in-domain)						
PBSMT (LM)	M30k _T	34.0	55.0 \dagger	44.7	68.0	(66.8 / 68.1)
NMT	M30k _T	<u>35.5</u> \dagger	53.4	<u>43.3</u> \dagger	65.2	(67.7 / 65.0)
NMT _{SRC+IMG}	M30k _T	37.1 $\dagger\dagger$	54.5 \dagger	42.8 $\dagger\dagger$	66.6	(67.2 / 66.5)
NMT_{SRC+IMG} vs. best PBSMT		\uparrow 3.1	\downarrow 0.5	\downarrow 1.9	\downarrow 1.4	\uparrow 0.4 / \downarrow 1.6
NMT_{SRC+IMG} vs. NMT		\uparrow 1.6	\uparrow 1.1	\downarrow 0.5	\uparrow 1.4	\downarrow 0.5 / \uparrow 1.5
Pre-training data set: WMT'15 English-German corpora (general domain)						
PBSMT (concat)	M30k _T	32.6	53.9	46.1	67.3	(66.3 / 67.4)
PBSMT (LM)	M30k _T	32.5	54.1	46.0	67.3	(66.0 / 67.4)
NMT	M30k _T	<u>37.8</u>	<u>56.7</u>	<u>41.0</u>	<u>69.2</u>	(69.7 / 69.1)
NMT _{SRC+IMG}	M30k _T	39.0 $\dagger\dagger$	56.8 \dagger	40.6 \dagger	69.6	(69.6 / 69.6)
NMT_{SRC+IMG} vs. best PBSMT		\uparrow 6.4	\uparrow 2.7	\downarrow 5.4	\uparrow 2.3	\uparrow 3.3 / \uparrow 2.2
NMT_{SRC+IMG} vs. NMT		\uparrow 1.2	\uparrow 0.1	\downarrow 0.4	\uparrow 0.4	\downarrow 0.1 / \uparrow 0.5

Model	BLEU4 [↑]	METEOR [↑]	TER [↓]	chrF3 [↑]
PBSMT	32.8	34.8	43.9	61.8
NMT	<u>38.2</u>	<u>35.8</u>	<u>40.2</u>	<u>62.8</u>
NMT _{SRC+IMG}	40.6^{†‡}	37.5^{†‡}	37.7^{†‡}	65.2
Improvements				
Ours vs. NMT	↑ 2.4	↑ 1.7	↓ 2.5	↑ 2.4
Ours vs. PBSMT	↑ 7.8	↑ 2.7	↓ 6.2	↑ 3.4
Pre-training data set: back-translated M30k_C (in-domain)				
PBSMT	36.8	36.4	40.8	64.5
NMT	<u>42.6</u>	<u>38.9</u>	<u>36.1</u>	<u>67.6</u>
NMT _{SRC+IMG}	43.2[‡]	39.0[‡]	35.5[‡]	67.7
Improvements				
Ours vs. PBSMT	↑ 6.4	↑ 2.6	↓ 5.3	↑ 3.2
Ours vs. NMT	↑ 0.6	↑ 0.1	↓ 0.6	↑ 0.1

Example



(a) Image-target word alignments.



(b) Source-target word alignments.

Conclusions

Conclusions

- **multi-modal neural MT** model with **two separate attention mechanisms**;
- **consistent improvements** when translating from English into German and vice-versa;
- model can **efficiently exploit additional data** in pre-training, e.g. **back-translated** and **text-only**;
- **visual attention** seems to learn to **focus on one important aspect** of the image, and the model chooses when to use it in generating a word via the **image gate**;

Conclusions

- **multi-modal neural MT** model with **two separate attention mechanisms**;
- **consistent improvements** when translating from English into German and vice-versa;
- model can **efficiently exploit additional data** in pre-training, e.g. **back-translated** and **text-only**;
- **visual attention** seems to learn to **focus on one important aspect** of the image, and the model chooses when to use it in generating a word via the **image gate**;

Conclusions

- **multi-modal neural MT** model with **two separate attention mechanisms**;
- **consistent improvements** when translating from English into German and vice-versa;
- model can **efficiently exploit additional data** in pre-training, e.g. **back-translated** and **text-only**;
- **visual attention** seems to learn to **focus on one important aspect** of the image, and the model chooses when to use it in generating a word via the **image gate**;

Conclusions

- **multi-modal neural MT** model with **two separate attention mechanisms**;
- **consistent improvements** when translating from English into German and vice-versa;
- model can **efficiently exploit additional data** in pre-training, e.g. **back-translated** and **text-only**;
- **visual attention** seems to learn to **focus on one important aspect** of the image, and the model chooses when to use it in generating a word via the **image gate**;

References

- Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. In *International Conference on Learning Representations. ICLR 2015*.
- Huang, P.-Y., Liu, F., Shiang, S.-R., Oh, J., and Dyer, C. (2016). Attention-based multimodal neural machine translation. In *Proceedings of the First Conference on Machine Translation*, pages 639–645, Berlin, Germany. Association for Computational Linguistics.

Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., and Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. In Blei, D. and Bach, F., editors, *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 2048–2057. JMLR Workshop and Conference Proceedings.

Thank you!