

# Bringing back simplicity and lightness into neural image captioning

Jean-Benoit Delbrouck and Stéphane Dupont

TCTS Lab, University of Mons, Belgium

{jean-benoit.delbrouck, stephane.dupont}@umons.ac.be

## Abstract

Neural Image Captioning (NIC) or neural caption generation has attracted a lot of attention over the last few years. Describing an image with a natural language has been an emerging challenge in both fields of computer vision and language processing. Therefore a lot of research has focused on driving this task forward with new creative ideas. So far, the goal has been to maximize scores on automated metric and to do so, one has to come up with a plurality of new modules and techniques. Once these add up, the models become complex and resource-hungry. In this paper, we take a small step backwards in order to study an architecture with interesting trade-off between performance and computational complexity. To do so, we tackle every component of a neural captioning model and propose one or more solution that lightens the model overall. Our ideas are inspired by two related tasks: Multimodal and Monomodal Neural Machine Translation.

## 1 Introduction

Problems combining vision and natural language processing such as image captioning (Chen et al. 2015) is viewed as an extremely challenging task. It requires to grasp and express low to high-level aspects of local and global areas in an image as well as their relationships. Over the years, it continues to inspire considerable research. Visual attention-based neural decoder models (Xu et al. 2015; Karpathy and Li 2015) have shown gigantic success and are now widely adopted for the NIC task. These recent advances are inspired from the neural encoder-decoder framework (Sutskever, Vinyals, and Le 2014; Bahdanau, Cho, and Bengio 2014)—or sequence to sequence model (seq2seq)—used for Neural Machine Translation (NMT). In that approach, Recurrent Neural Networks (RNN, Mikolov et al. 2010) map a source sequence of words (encoder) to a target sequence (decoder). An attention mechanism is learned to focus on different parts of the source sentence while decoding. The same mechanism applies for a visual input; the attention module learns to attend the salient parts of an image while decoding the caption.

These two fields, NIC and NMT, led to a Multimodal Neural Machine Translation (MNMT, Specia et al. 2016) task where the sentence to be translated is supported by the information from an image. Interestingly, NIC and MNMT

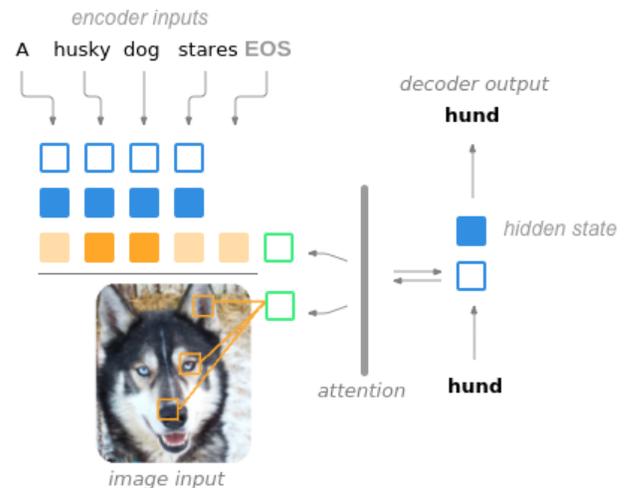


Figure 1: This image depicts a decoder timestep of the MNMT architecture. At time  $t$ , the decoder attend both a visual and textual representations. In NIC, the decoder only attends an image. This shows how both tasks are related.

share a very similar decoder: they are both required to generate a meaningful natural language description or translation with the help of a visual input. However, both tasks differ in the amount of annotated data. MNMT has  $\approx 19$  times less unique training examples, reducing the amount of learnable parameters and potential complexity of a model. Yet, over the years, the challenge has brought up very clever and elegant ideas that could be transferred to the NIC task. The aim of this paper is to propose such an architecture for NIC in a straightforward manner. Indeed, our proposed models work with less data, less parameters and require less computation time. More precisely, this paper intends to:

- Work only with in-domain data. No additional data besides proposed captioning datasets are involved in the learning process;
- Lighten as much as possible the training data used, i.e. the visual and linguistic inputs of the model;
- Propose a subjectively light and straightforward yet efficient NIC architecture with high training speed.

## 2 Captioning Model

As quickly mentioned in section 1, a neural captioning model is a RNN-decoder (Bahdanau, Cho, and Bengio 2014) that uses an attention mechanism over an image  $I$  to generate a word  $y_t$  of the caption at each time-step  $t$ . The following equations depict what a baseline time-step  $t$  looks like (Xu et al. 2015):

$$\mathbf{x}_t = \mathbf{W}^x \mathbf{E} y_{t-1} \quad (1)$$

$$\mathbf{c}_t = f_{\text{att}}(\mathbf{h}_{t-1}, I) \quad (2)$$

$$\mathbf{h}_t = f_{\text{rnn}}(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_t) \quad (3)$$

$$y_t \sim \mathbf{p}_t = \mathbf{W}^y [\mathbf{y}_{t-1}, \mathbf{h}_t, \mathbf{c}_t] \quad (4)$$

where equation 1 maps the previous embedded word generated to the RNN hidden state size with matrix  $\mathbf{W}^x$ , equation 2 is the attention module over the image  $I$ , equation 3 is the RNN cell computation and equation 4 is the probability distribution  $\mathbf{p}_t$  over the vocabulary (matrix  $\mathbf{W}^y$  is also called the projection matrix).

If we denote  $\theta$  as the model parameters, then  $\theta$  is learned by maximizing the likelihood of the observed sequence  $\mathbf{Y} = (\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n)$  or in other words by minimizing the cross entropy loss. The objective function is given by:

$$\mathcal{L}(\theta) = - \sum_{t=1}^n \log p_{\theta}(\mathbf{y}_t | \mathbf{y} < t, I) \quad (5)$$

The paper is structured so that each section tackles an equation (i.e. a main component of the captioning model) in the following manner: section 2.1 for equation 1 (embeddings), section 2.2 for equation 3 ( $f_{\text{rnn}}$ ), section 2.3 for equation 2 ( $f_{\text{att}}$ ), section 2.4 for equation 4 (projection) and section 2.5 for equation 5 (objective function).

### 2.1 Embeddings

Recall that:  $\mathbf{x}_t = \mathbf{W}^x \mathbf{E} y_{t-1}$

The total size of the embeddings matrix  $\mathbf{E}$  depends on the vocabulary size  $|\mathcal{Y}|$  and the embedding dimension  $d$  such that  $\mathbf{E} \in \mathbb{R}^{|\mathcal{Y}| \times d}$ . The mapping matrix  $\mathbf{W}^x$  also depends on the embedding dimension because  $\mathbf{W}^x \in \mathbb{R}^{d \times |\mathcal{X}_t|}$ .

Many previous researches (Karpathy and Li 2015; You et al. 2016; Yao et al. 2017; Anderson et al. 2018) uses pretrained embeddings such as Glove and word2vec or one-hot-vectors. Both word2vec and glove provide distributed representation of words. These models are pre-trained on 30 and 42 billions words respectively (Mikolov et al. 2013; Pennington, Socher, and Manning 2014), weights several gigabytes and work with  $d = 300$ .

For our experiments, each word  $y_i$  is a column index in an embedding matrix  $\mathbf{E}_y$  learned along with the model and initialized using some random distribution. Whilst the usual allocation is from 512 to 1024 dimensions per embedding (Xu et al. 2015; Lu et al. 2017; Mun, Cho, and Han 2017;

Rennie et al. 2017) we show that a small embedding size of  $d = 128$  is sufficient to learn a strong vocabulary representation. The solution of an jointly-learned embedding matrix also tackles the high-dimensionality and sparsity problem of one-hot vectors. For example, (Anderson et al. 2018) works with a vocabulary of 10,010 and a hidden size of 1000. As a result, the mapping matrix  $\mathbf{W}^x$  of equation 1 has 10 millions parameters.

Working with a small vocabulary, besides reducing the size of embedding matrix  $\mathbf{E}$ , presents two major advantages: it lightens the projection module (as explained further in section 2.4) and reduces the action space in a Reinforcement Learning setup (detailed in section 2.5). To marginally reduce our vocabulary size (of  $\approx 50\%$ ), we use the byte pair encoding (BPE) algorithm on the train set to convert space-separated tokens into subwords (Sennrich, Haddow, and Birch 2016). Applied originally for NMT, BPE is based on the intuition that various word classes are made of smaller units than words such as compounds and loanwords. In addition of making the vocabulary smaller and the sentences length shorter, the subword model is able to productively generate new words that were not seen at training time.

### 2.2 Conditional GRU

Recall that:  $\mathbf{h}_t = f_{\text{rnn}}(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_t)$

Most previous researches in captioning (Karpathy and Li 2015; You et al. 2016; Rennie et al. 2017; Mun, Cho, and Han 2017; Lu et al. 2017; Yao et al. 2017) used an LSTM (Hochreiter and Schmidhuber 1997) for their  $f_{\text{rnn}}$  function. Our recurrent model is a pair of two Gated Recurrent Units (GRU (Cho et al. 2014)), called conditional GRU (cGRU), as previously investigated in NMT<sup>1</sup>. A GRU is a lighter gating mechanism than LSTM since it doesn't use a forget gate and lead to similar results in our experiments.

The cGRU also addresses the encoding problem of the  $f_{\text{att}}$  mechanism. As shown in equation 2 the context vector  $\mathbf{c}_t$  takes the previous hidden state  $\mathbf{h}_{t-1}$  as input which is outside information of the current time-step. This could be tackled by using the current hidden  $\mathbf{h}_t$ , but then context vector  $\mathbf{c}_t$  is not an input of  $f_{\text{rnn}}$  anymore. A conditional GRU is an efficient way to both build and encode the result of the  $f_{\text{att}}$  module.

Mathematically, a first independent GRU encodes an intermediate hidden state proposal  $\mathbf{h}'_t$  based on the previous hidden state  $\mathbf{h}_{t-1}$  and input  $\mathbf{x}_t$  at each time-step  $t$ :

<sup>1</sup><https://github.com/nyu-dl/dl4mt-tutorial/blob/master/docs/cgru.pdf>

$$\begin{aligned}
z'_t &= \sigma(x_t + U'_z h_{t-1}) \\
r'_t &= \sigma(x_t + U'_r h_{t-1}) \\
\underline{h}'_t &= \tanh(x_t + r'_t \odot (U' h_{t-1})) \\
h'_t &= (1 - z'_t) \odot \underline{h}'_t + z'_t \odot h_{t-1}
\end{aligned} \tag{6}$$

Then, the attention mechanism computes  $c_t$  over the source sentence using the image  $I$  and the intermediate hidden state proposal  $h'_t$  similar to 3:

$$c_t = f_{\text{att}}(h'_t, I)$$

Finally, a second independent GRU computes the hidden state  $h_t$  of the cGRU by looking at the intermediate representation  $h'_t$  and context vector  $c_t$ :

$$\begin{aligned}
z_t &= \sigma(W_z c_t + U_z h'_t) \\
r_t &= \sigma(W_r c_t + U_r h'_t) \\
\underline{h}_t &= \tanh(W c_t + r_t \odot (U h'_t)) \\
h_t &= (1 - z_t) \odot \underline{h}_t + z_t \odot h'_t
\end{aligned} \tag{7}$$

We see that both problem are addressed: context vector  $c_t$  is computed according to the intermediate representation  $h'_t$  and the final hidden state  $h_t$  is computed according to the context vector  $c_t$ . Again, the size of the hidden state  $|h_t|$  in the literature varies between 512 and 1024, we pick  $|h_t| = 256$ .

The most similar approach to ours is the Top-Down Attention of (Anderson et al. 2018) that encodes the context vector the same way but with LSTM and a different hidden state layout.

### 2.3 Attention model

Recall that:  $c_t = f_{\text{att}}(h_{t-1}, I)$

Since the image is the only input to a captioning model, the attention module is crucial but also very diverse amongst different researches. For example, (You et al. 2016) use a semantic attention where, in addition of image features, they run a set of attribute detectors to get a list of visual attributes or concepts that are most likely to appear in the image. (Anderson et al. 2018) uses the Visual Genome dataset to pre-train his bottom-up attention model. This dataset contains 47,000 out-of-domain images of the captioning dataset densely annotated with scene graphs containing objects, attribute and relationships. (Yang et al. 2016) proposes a review network which is an extension to the decoder. The review network performs a given number of review steps on the hidden states and outputs a compact vector representation available for the attention mechanism.

Yet everyone seems to agree on using a Convolutional Neural Network (CNN) to extract features of the image  $I$ . The trend is to select features matrices, at the convolutional layers, of size  $14 \times 14 \times 1024$  (Resnet, He et al. 2016, res4f layer) or  $14 \times 14 \times 512$  (VGGNet Simonyan and Zisserman

2014, conv5 layer). Other attributes can be extracted in the last fully connected layer of a CNN and has shown to bring useful information (Yang et al. 2016; Yao et al. 2017; You et al. 2016) Some models also finetune the CNN during training (Yang et al. 2016; Mun, Cho, and Han 2017; Lu et al. 2017) stacking even more trainable parameters.

Our attention model  $f_{\text{att}}$  is guided by a unique vector with global 2048-dimensional visual representation  $V$  of image  $I$  extracted at the pool5 layers of a ResNet-50. Our attention vector is computed so:

$$c_t = h'_t \odot \tanh(W^{\text{img}} V_I) \tag{8}$$

Recall that following the cGRU presented in section 2.2, we work with  $h'_t$  and not  $h_{t-1}$ . Even though pooled features have less information than convolutional features ( $\approx 50$  to 100 times less), pooled features have shown great success in combination with cGRU in MNMT (Caglayan et al. 2017a). Hence, our attention model is only the single matrix  $W^{\text{img}} \in \mathbb{R}^{2048 \times |h'_t|}$

### 2.4 Projection

Recall that:  $y_t \sim p_t = W^y[y_{t-1}, h_t, c_t]$

The projection also accounts for a lot of trainable parameters in the captioning model, especially if the vocabulary is large. Indeed, in equation 4 the projection matrix is  $\in \mathbb{R}^{|y_{t-1}, h_t, c_t| \times |\mathcal{Y}|}$ . To reduce the number of parameters, we use a bottleneck function:

$$\begin{aligned}
b_t &= f_{\text{bot}}(y_{t-1}, h_t, c_t) = W^{\text{bot}}[y_{t-1}, h_t, c_t] \\
y_t \sim p_t &= W^{y.\text{bot}} b_t
\end{aligned} \tag{9}$$

where  $|b_t| < |[y_{t-1}, h_t, c_t]|$  so that  $|W^{\text{bot}}| + |W^{y.\text{bot}}| < |W^y|$ . Interestingly enough, if  $|b_t| = d$  (embedding size), then  $|W^{y.\text{bot}}| = |E|$ . We can share the weights between the two matrices (i.e.  $W^{y.\text{bot}} = E$ ) to marginally reduce the number of learned parameters. Moreover, doing so doesn't negatively impact the captioning results.

We push our projection further and use a deep-GRU, used originally in MNMT (Delbrouck and Dupont 2018), so that our bottleneck function  $f_{\text{bot}}$  is now a third GRU as described by equations 7:

$$b_t = f_{\text{bot}}(y_{t-1}, h'_t, c_t) = \text{cGRU}([y_{t-1}, h'_t, c_t], h_t) \tag{11}$$

Because we work with small dimension, adding a new GRU block on top barely increases the model size.

### 2.5 Objective function

To directly optimize a automated metric, we can see the captioning generator as a Reinforcement Learning (RL) problem. The introduced  $f_{\text{mn}}$  function is viewed as an agent that interact with an environment composed of words and image features. The agent interacts with the environment by

taking actions that are the prediction of the next word of the caption. An action is the result of the policy  $p_{\theta}$  where  $\theta$  are the parameters of the network. Whilst very effective to boost the automatic metric scores, porting the captioning problem into a RL setup significantly reduce the training speed.

Ranzato et al. 2015 proposed a method (MIXER), based on the REINFORCE method, combined with a baseline reward estimator. However, they implicitly assume each intermediate action (word) in a partial sequence has the same reward as the sequence-level reward, which is not true in general. To compensate for this, they introduce a form of training that mixes together the MLE objective and the REINFORCE objective. Liu et al. 2017 also addresses the delayed reward problem by estimating at each time-step the future rewards based on Monte Carlo rollouts. Rennie et al. 2017 utilizes the output of its own test-time inference model to normalize the rewards it experiences. Only samples from the model that outperform the current test-time system are given positive weight.

To keep it simple, and because our reduced vocabulary allows us to do so, we follow the work of Ranzato et al. 2015 and use the naive variant of the policy gradient with REINFORCE. The loss function in equation 5 is now given by:

$$\mathcal{L}(\theta) = -\mathbb{E}_{\mathbf{Y} \sim p_{\theta}} r(\mathbf{Y}) \quad (12)$$

where  $r(\mathbf{Y})$  is the reward (here the score given by an automatic metric scorer) of the outputted caption  $\mathbf{Y} = (y_1, y_2, \dots, y_n)$ .

We use the REINFORCE algorithm based on the observation that the expected gradient of a non-differentiable reward function is computed as follows:

$$\nabla_{\theta} \mathcal{L}(\theta) = -\mathbb{E}_{\mathbf{Y} \sim p_{\theta}} [r(\mathbf{Y}) \nabla_{\theta} \log_{p_{\theta}}(\mathbf{Y})] \quad (13)$$

The expected gradient can be approximated using  $N$  Monte-Carlo sample  $\mathbf{Y}$  for each training example in the batch:

$$\nabla_{\theta} \mathcal{L}(\theta) = \nabla - \left[ \frac{1}{N} \sum_{i=1}^N [r_i(\mathbf{Y}_i) \log_{p_{\theta}}(\mathbf{Y}_i)] \right] \quad (14)$$

In practice, we can approximate with one sample:

$$\nabla_{\theta} \mathcal{L}(\theta) \approx -r(\mathbf{Y}) \nabla_{\theta} \log_{p_{\theta}}(\mathbf{Y}) \quad (15)$$

The policy gradient can be generalized to compute the reward associated with an action value relative to a baseline  $b$ . This baseline either encourages a word choice  $y_t$  if  $r_t > b_t$  or discourages it  $r_t < b_t$ . If the baseline is an arbitrary function that does not depend on the actions  $y_1, y_2, \dots, y_n \in \mathbf{Y}$  then baseline does not change the expected gradient, and importantly, reduces the variance of the gradient estimate. The final expression is given by:

$$\nabla_{\theta} \mathcal{L}(\theta) \approx -(r(\mathbf{Y}) - b) \nabla_{\theta} \log_{p_{\theta}}(\mathbf{Y}) \quad (16)$$

### 3 Settings

Our decoder is a cGRU where each GRU is of size  $|\mathbf{h}_t| = 256$ . Word embedding matrix  $E$  allocates  $d = 128$  features per word. To create the image annotations used by our decoder, we used a ResNet-50 and extracted the features of size 1024 at the pool-5 layer. As regularization method, we apply dropout with a probability of 0.5 on bottleneck  $b_t$  and we early stop the training if the validation set CIDER metric does not improve for 10 epochs. All variants of our models are trained with ADAM optimizer (Kingma and Ba 2014) with a learning rate of  $4e^{-4}$  and mini-batch size of 256. We decode with a beam-search of size 3. In the RL setting, the baseline is a linear projection of  $\mathbf{h}_t$ .

We evaluate our models on MSCOCO (Lin et al. 2014), the most popular benchmark for image captioning which contains 82,783 training images and 40,504 validation images. There are 5 human-annotated descriptions per image. As the annotations of the official testing set are not publicly available, we follow the settings in prior work (or "Kaparthys splits"<sup>2</sup>) that takes 82,783 images for training, 5,000 for validation and 5,000 for testing. On the training captions, we use the byte pair encoding algorithm on the train set to convert space-separated tokens into subwords (Sennrich, Haddow, and Birch 2016, 5000 symbols), reducing our vocabulary size to 5066 english tokens. For the online-evaluation, all images are used for training except for the validation set.

### 4 Results

Our models performance are evaluated according to the following automated metrics: BLEU-4 (Papineni et al. 2002), METEOR (Vedantam, Lawrence Zitnick, and Parikh 2015) and CIDER-D (Lavie and Agarwal 2007). Results shown in table 1 are using cross-entropy (XE) loss (cfr. equation 5). Reinforced learning optimization results are compared in table 3.

#### 4.1 XE scores

We sort the different works in table 1 by CIDER score. For every of them, we detail the trainable weights involved in the learning process (Wt.), the number of visual features used for the attention module (Att. Feat), the amount of out-of-domain data (O.O.D) and the convergence speed (epoch).

As we see, our model has the third best METEOR and CIDER scores across the board. Yet our BLEU metric is quiet low, we postulate two potential causes. Either our model has not enough parameters to learn the correct precision for a set of n-grams as the metric would require or it is a direct drawback from using subwords. Nevertheless, the CIDER and METEOR metric show that the main concepts are presents our captions. Our models are also the lightest in regards to trainable parameters and attention features number. As far as convergence in epochs were reported in

<sup>2</sup><https://github.com/kaparthys/neuraltalk2/tree/master/coco>

Table 1: Table sorted per CIDER-D score of models being optimized with cross-entropy loss only (cfr. equation 5).  
 ◦ pool features, • conv features, \* FC features, § means glove or word2vec embeddings, † CNN finetuning in-domain, ‡ using in-domain CNN, II CNN finetuning OOD

	<b>B4</b>	<b>M</b>	<b>C</b>	<b>Wt. (in M)</b>	<b>Att. feat. (in K)</b>	<b>O.O.D. (in M)</b>	<b>epoch</b>
<i>This work</i>							
cGRU ◦	0.302	0.258	1.018	2.46	2	-	9
<i>Comparable work</i>							
Adaptive(Lu et al. 2017) •†	0.332	0.266	1.085	17.4	100	-	42
Top-down (Anderson et al. 2018) •II	0.334	0.261	1.054	≈ 25	204	-	-
Boosting (Yao et al. 2017) ◦*	0.325	0.251	0.986	≈ 28.4	2	-	123
Review (Yang et al. 2016) •* †	0.290	0.237	0.886	≈ 12.3	101	-	100
SAT (Xu et al. 2015) •	0.250	0.230	-	≈ 18	100	-	-
<i>O.O.D work</i>							
Top-down (Anderson et al. 2018) •II	0.362	0.27	1.135	≈ 25	920	920 <sub>RCNN</sub>	-
T-G att (Mun et al. 2017) •†	0.326	0.257	1.024	≈ 12.8	200	14 <sub>TSV</sub>	-
Semantic (You et al. 2016) ◦* §	0.304	0.243	-	5.5	2	3 <sub>glove</sub>	-
NT(Karpathy and Li 2015)•§	0.230	0.195	0.660	-	-	3 <sub>glove</sub>	-

previous works, our cGRU model is by far the fastest to train.

The following table 2 concerns the online evaluation of the official MSCOCO test-set 2014 split. Scores of our model are an ensemble of 5 runs with different initialization.

Table 2: Published Ranking image captioning results on the online MSCOCO test server

	<b>B4(c5)</b>	<b>M (c5)</b>	<b>C(c5)</b>
cGru	0.326	0.253	0.973
<i>Comparable work</i>			
(Lu et al. 2017)	0.336	0.264	1.042
(Yao et al. 2017)	0.330	0.256	0.984
(Yang et al. 2016)	0.313	0.256	0.965
(You et al. 2016)	0.316	0.250	0.943
(Wu et al. 2016)	0.306	0.246	0.911
(Xu et al. 2015)	0.277	0.251	0.865

We see that our model suffers a minor setback on this test-set, especially in term of CIDER score whilst the adaptive (Lu et al. 2017) and boosting method (Yao et al. 2017) yields to stable results for both test-sets.

## 4.2 RL scores

The table 3 depicts the different papers using direct metric optimization. Rennie et al. 2017 used the SCST method, the most effective one according to the metrics boost (+23, +3

an +123 points respectively) but also the most sophisticated. Liu et al. 2017 used a similar approach than ours (MIXER) but with Monte-Carlo roll-outs (i.e. sampling until the end at every time-step  $t$ ). Without using this technique, two of our metrics improvement (METEOR and CIDER) surpasses the MC roll-outs variant (+0 against -2 and +63 against +48 respectively).

Table 3: All optimization are on the CIDER metric.

	<b>B4</b>	<b>M</b>	<b>C</b>
<i>Renn. et al. 2017</i>			
XE	0.296	0.252	0.940
RL-SCST	0.319 ↑23	0.255 ↑3	1.063 ↑123
<i>Liu et al. 2017</i>			
XE	0.294	0.251	0.947
RL-PG	0.333 ↑39	0.249 ↓2	0.995 ↑48
<i>Ours</i>			
XE	0.302	0.258	1.018
RL-PG	0.315 ↑13	0.258	1.071 ↑63

## 4.3 Scalability

An interesting investigation would be to leverage the architecture with more parameters to see how it scales. We showed our model performs well with few parameters, but we would like to show that it could be used as a base for more complex posterior researches.

We propose two variants to effectively do so :

- **cGRUx2** The first intuition is to double the width of the

model, i.e. the embedding size  $d$  and the hidden state size  $|\mathbf{h}_t|$ . Unfortunately, this setup is not ideal with a deep-GRU because the recurrent matrices of equations 6 and 7 for the bottleneck GRU gets large. We can still use the classic bottleneck function (equation 9).

- **MHA** We trade our attention model described in section 2.3 for a standard multi-head attention (MHA) to see how convolutional features could improve our CIDEr-D metrics. Multi-head attention (Vaswani et al. 2017) computes a weighted sum of some values, where the weight assigned to each value is computed by a compatibility function of the query with the corresponding key. This process is repeated  $n$  multiple times. The compatibility function is given by:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_q}}\right)V$$

where the query  $Q$  is  $\mathbf{h}'_t$ , the keys and values are a set of 196 vectors of dimension 1024 (from the layer res4f\_relu of ResNet-50 CNN). Authors found it beneficial to linearly project the queries, keys and values  $n$  times with different learned linear projections to dimension  $d_q$ . The output of the multi-head attention is the concatenation of the  $n$  number of  $d_q$  values linearly projected to  $d_q$  again. We pick  $d_q = |\mathbf{h}'_t|$  and  $n = 3$ . The multi-head attention model adds up 0.92M parameters if  $|\mathbf{h}'_t| = 256$  and 2.63M parameters if  $|\mathbf{h}'_t| = 512$  (in the case of cGRUx2).

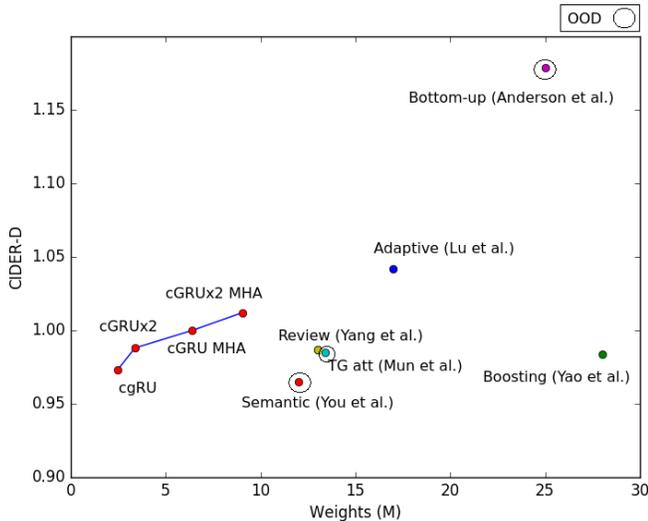


Figure 2: The figure shows a new set of results on the online MSCOCO test server and tries to put those in perspective

We have hence proposed an image captioning architecture that, compared to previous work, cover a different area in the performance-complexity trade-off plane. We hope these will be of interest and will fuel more research in this direction.

## 5 Related work

As mentioned in the introduction, our model is largely inspired by the work carried out in NMT and MNMT.

Component such as attention models like multi-head, encoder-based and pooled attention (2017a; 2017b; 2017); reinforcement learning in NMT and MNMT (2016; 2017a); embeddings (2016; 2017) are well investigated.

In captioning, Anderson et al. used a very similar approach where two LSTMs build and encode the visual features. Two other works, Yao et al. and You et al. used pooled features as described in this paper. However, they both used an additional vector taken from the fully connected layer of a CNN.

## 6 Conclusion

We presented a novel and light architecture composed of a cGRU that showed interesting performance. The model builds and encodes the context vector from pooled features in an efficient manner. The attention model presented in section 2.3 is really straightforward and seems to bring the necessary visual information in order to output complete captions. Also, we empirically showed that the model can easily scale with more sought-after modules or simple with more parameters. In the future, it would be interesting to use different attention features, like VGG or GoogleNet (that have only 1024 dimensions) or different attention models to see how far this architecture can get.

## 7 Acknowledgements

This work was partly supported by the Chist-Era project IGLU with contribution from the Belgian Fonds de la Recherche Scientifique (FNRS), contract no. R.50.11.15.F, and by the FSO project VCYCLE with contribution from the Belgian Walloon Region, contract no. 1510501.

We also thank the authors of nmtpytorch<sup>3</sup> (Caglayan et al. 2017b) that we used as framework for our experiments. Our code is made available for posterior research<sup>4</sup>.

<sup>3</sup><https://github.com/lium-lst/nmtpytorch>

<sup>4</sup>[https://github.com/jbdel/light\\_captioning](https://github.com/jbdel/light_captioning)

## References

- Anderson, P.; He, X.; Buehler, C.; Teney, D.; Johnson, M.; Gould, S.; and Zhang, L. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*.
- Bahdanau, D.; Cho, K.; and Bengio, Y. 2014. Neural machine translation by jointly learning to align and translate. *arXiv e-prints* abs/1409.0473.
- Caglayan, O.; Aransa, W.; Bardet, A.; García-Martínez, M.; Bougares, F.; Barrault, L.; Masana, M.; Herranz, L.; and van de Weijer, J. 2017a. Lium-cvc submissions for wmt17 multimodal translation task. In *Proceedings of the Second Conference on Machine Translation*, 432–439. Association for Computational Linguistics.
- Caglayan, O.; García-Martínez, M.; Bardet, A.; Aransa, W.; Bougares, F.; and Barrault, L. 2017b. Nmtpy: A flexible toolkit for advanced neural machine translation systems. *Prague Bull. Math. Linguistics* 109:15–28.
- Chen, X.; Fang, H.; Lin, T.; Vedantam, R.; Gupta, S.; Dollár, P.; and Zitnick, C. L. 2015. Microsoft COCO captions: Data collection and evaluation server. *CoRR* abs/1504.00325.
- Cho, K.; van Merriënboer, B.; Gülçehre, Ç.; Bahdanau, D.; Bougares, F.; Schwenk, H.; and Bengio, Y. 2014. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1724–1734. Doha, Qatar: Association for Computational Linguistics.
- Delbrouck, J.-B., and Dupont, S. 2017a. An empirical study on the effectiveness of images in multimodal neural machine translation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 910–919. Association for Computational Linguistics.
- Delbrouck, J., and Dupont, S. 2017b. Modulating and attending the source image during encoding improves multimodal translation. *CoRR* abs/1712.03449.
- Delbrouck, J.-B., and Dupont, S. 2018. Umons submission for wmt18 multimodal translation task. In *Proceedings of the First Conference on Machine Translation*. Brussels, Belgium: Association for Computational Linguistics.
- Delbrouck, J.-B.; Dupont, S.; and Seddati, O. 2017. Visually grounded word embeddings and richer visual features for improving multimodal neural machine translation. In *Proc. GLU 2017 International Workshop on Grounding Language Understanding*, 62–67.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 770–778.
- Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural Comput.* 9(8):1735–1780.
- Karpathy, A., and Li, F.-F. 2015. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 3128–3137. IEEE Computer Society.
- Kingma, D. P., and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Lavie, A., and Agarwal, A. 2007. Meteor: An automatic metric for mt evaluation with high levels of correlation with human judgments. In *Proceedings of the Second Workshop on Statistical Machine Translation, StatMT '07*, 228–231. Stroudsburg, PA, USA: Association for Computational Linguistics.
- Lin, T.-Y.; Maire, M.; Belongie, S.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; and Zitnick, C. L. 2014. Microsoft coco: Common objects in context. In Fleet, D.; Pajdla, T.; Schiele, B.; and Tuytelaars, T., eds., *Computer Vision – ECCV 2014*, 740–755. Cham: Springer International Publishing.
- Liu, S.; Zhu, Z.; Ye, N.; Guadarrama, S.; and Murphy, K. 2017. Improved image captioning via policy gradient optimization of spider. *2017 IEEE International Conference on Computer Vision (ICCV)* 873–881.
- Lu, J.; Xiong, C.; Parikh, D.; and Socher, R. 2017. Knowing when to look: Adaptive attention via a visual sentinel for image captioning.
- Mikolov, T.; Karafit, M.; Burget, L.; Cernock, J.; and Khudanpur, S. 2010. Recurrent neural network based language model. In Kobayashi, T.; Hirose, K.; and Nakamura, S., eds., *INTERSPEECH*, 1045–1048. ISCA.
- Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G. S.; and Dean, J. 2013. Distributed representations of words and phrases and their compositionality. In Burges, C. J. C.; Bottou, L.; Welling, M.; Ghahramani, Z.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems 26*. Curran Associates, Inc. 3111–3119.
- Mun, J.; Cho, M.; and Han, B. 2017. Text-guided attention model for image captioning. In *AAAI*.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL '02*, 311–318. Stroudsburg, PA, USA: Association for Computational Linguistics.
- Pennington, J.; Socher, R.; and Manning, C. D. 2014. Glove: Global vectors for word representation. In *EMNLP*, volume 14, 1532–1543.
- Ranzato, M.; Chopra, S.; Auli, M.; and Zaremba, W. 2015. Sequence level training with recurrent neural networks. *CoRR* abs/1511.06732.
- Rennie, S. J.; Marcheret, E.; Mroueh, Y.; Ross, J.; and Goel, V. 2017. Self-critical sequence training for image captioning. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 1179–1195.
- Sennrich, R.; Haddow, B.; and Birch, A. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1715–1725. Association for Computational Linguistics.
- Shen, S.; Cheng, Y.; He, Z.; He, W.; Wu, H.; Sun, M.; and Liu, Y. 2016. Minimum risk training for neural machine translation. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1:*

*Long Papers*), 1683–1692. Association for Computational Linguistics.

Simonyan, K., and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *CoRR* abs/1409.1556.

Specia, L.; Frank, S.; Sima'an, K.; and Elliott, D. 2016. A shared task on multimodal machine translation and crosslingual image description. In *Proceedings of the First Conference on Machine Translation*, 543–553. Berlin, Germany: Association for Computational Linguistics.

Sutskever, I.; Vinyals, O.; and Le, Q. V. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, 3104–3112.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L. u.; and Polosukhin, I. 2017. Attention is all you need. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc. 5998–6008.

Vedantam, R.; Lawrence Zitnick, C.; and Parikh, D. 2015. Cider: Consensus-based image description evaluation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Wu, Q.; Shen, C.; Liu, L.; Dick, A.; and van den Hengel, A. 2016. What value do explicit high level concepts have in vision to language problems? In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

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 Bach, F., and Blei, D., eds., *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, 2048–2057. Lille, France: PMLR.

Yang, Z.; Yuan, Y.; Wu, Y.; Cohen, W. W.; and Salakhutdinov, R. R. 2016. Review networks for caption generation. In Lee, D. D.; Sugiyama, M.; Luxburg, U. V.; Guyon, I.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 29*. Curran Associates, Inc. 2361–2369.

Yao, T.; Pan, Y.; Li, Y.; Qiu, Z.; and Mei, T. 2017. Boosting image captioning with attributes. In *ICCV*.

You, Q.; Jin, H.; Wang, Z.; Fang, C.; and Luo, J. 2016. Image captioning with semantic attention. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 4651–4659.