

# On the Automatic Generation of Medical Imaging Reports

Baoyu Jing<sup>†\*</sup> Pengtao Xie<sup>†\*</sup> Eric P. Xing<sup>†</sup>

<sup>†</sup>Petuum Inc, USA

<sup>\*</sup>School of Computer Science, Carnegie Mellon University, USA

{baoyu.jing, pengtao.xie, eric.xing}@petuum.com

## Abstract

Medical imaging is widely used in clinical practice for diagnosis and treatment. Report-writing can be error-prone for unexperienced physicians, and time-consuming and tedious for experienced physicians. To address these issues, we study the automatic generation of medical imaging reports. This task presents several challenges. First, a complete report contains multiple heterogeneous forms of information, including *findings* and *tags*. Second, abnormal regions in medical images are difficult to identify. Third, the reports are typically long, containing multiple sentences. To cope with these challenges, we (1) build a multi-task learning framework which jointly performs the prediction of tags and the generation of paragraphs, (2) propose a co-attention mechanism to localize regions containing abnormalities and generate narrations for them, (3) develop a hierarchical LSTM model to generate long paragraphs. We demonstrate the effectiveness of the proposed methods on two publicly available datasets.

## 1 Introduction

Medical images, such as radiology and pathology images, are widely used in hospitals for the diagnosis and treatment of many diseases, such as pneumonia and pneumothorax. The reading and interpretation of medical images are usually conducted by specialized medical professionals. For example, radiology images are read by radiologists. They write textual reports (Figure 1) to narrate the findings regarding each area of the body examined in the imaging study, specifically



**Impression:** No acute cardiopulmonary abnormality.

**Findings:** There are no focal areas of consolidation. No suspicious pulmonary opacities. Heart size within normal limits. No pleural effusions. There is no evidence of pneumothorax. Degenerative changes of the thoracic spine.

**MTI Tags:** degenerative change

Figure 1: An exemplar chest x-ray report. In the *impression* section, the radiologist provides a diagnosis. The *findings* section lists the radiology observations regarding each area of the body examined in the imaging study. The *tags* section lists the keywords which represent the critical information in the findings. These keywords are identified using the Medical Text Indexer (MTI).

whether each area was found to be normal, abnormal or potentially abnormal.

For less-experienced radiologists and pathologists, especially those working in the rural area where the quality of healthcare is relatively low, writing medical-imaging reports is demanding. For instance, to correctly read a chest x-ray image, the following skills are needed (Delrue et al., 2011): (1) thorough knowledge of the normal anatomy of the thorax, and the basic physiology of chest diseases; (2) skills of analyzing the radiograph through a fixed pattern; (3) ability of evaluating the evolution over time; (4) knowledge of clinical presentation and history; (5) knowledge of the correlation with other diagnostic results (laboratory results, electrocardiogram, and respiratory function tests).

For experienced radiologists and pathologists, writing imaging reports is tedious and time-consuming. In nations with large population such as China, a radiologist may need to read hundreds

of radiology images per day. Typing the findings of each image into computer takes about 5-10 minutes, which occupies most of their working time. In sum, for both unexperienced and experienced medical professionals, writing imaging reports is unpleasant.

This motivates us to investigate whether it is possible to automatically generate medical image reports. Several challenges need to be addressed. First, a complete diagnostic report is comprised of multiple heterogeneous forms of information. As shown in Figure 1, the report for a chest x-ray contains *impression* which is a sentence, *findings* which are a paragraph, and *tags* which are a list of keywords. Generating this heterogeneous information in a unified framework is technically demanding. We address this problem by building a multi-task framework, which treats the prediction of tags as a multi-label classification task, and treats the generation of long descriptions as a text generation task.

Second, how to localize image-regions and attach the right description to them are challenging. We solve these problems by introducing a co-attention mechanism, which simultaneously attends to images and predicted tags and explores the synergistic effects of visual and semantic information.

Third, the descriptions in imaging reports are usually long, containing multiple sentences. Generating such long text is highly nontrivial. Rather than adopting a single-layer LSTM (Hochreiter and Schmidhuber, 1997), which is less capable of modeling long word sequences, we leverage the compositional nature of the report and adopt a hierarchical LSTM to produce long texts. Combined with the co-attention mechanism, the hierarchical LSTM first generates high-level topics, and then produces fine-grained descriptions according to the topics.

Overall, the main contributions of our work are:

- We propose a multi-task learning framework which can simultaneously predict the tags and generate the text descriptions.
- We introduce a co-attention mechanism for localizing sub-regions in the image and generating the corresponding descriptions.
- We build a hierarchical LSTM to generate long paragraphs.

- We perform extensive experiments to show the effectiveness of the proposed methods.

The rest of the paper is organized as follows. Section 2 reviews related works. Section 3 introduces the method. Section 4 presents the experimental results and Section 5 concludes the paper.

## 2 Related Works

**Textual labeling of medical images** There have been several works aiming at attaching “texts” to medical images. In their settings, the target “texts” are either fully-structured or semi-structured (e.g. tags, templates), rather than natural texts. Kisilev et al. (2015) build a pipeline to predict the attributes of medical images. Shin et al. (2016) adopt a CNN-RNN based framework to predict tags (e.g. locations, severities) of chest x-ray images. The work closest to ours is recently contributed by (Zhang et al., 2017), which aims at generating semi-structured pathology reports, whose contents are restricted to 5 predefined topics.

However, in the real-world, different physicians usually have different writing habits and different x-ray images will represent different abnormalities. Therefore, collecting semi-structured reports is less practical and thus it is important to build models to learn from natural reports. To the best of our knowledge, our work represents the first one that generates truly natural reports written by physicians, which are usually long and cover diverse topics.

**Image captioning with deep learning** Image captioning aims at automatically generating text descriptions for given images. Most recent image captioning models are based on a CNN-RNN framework (Vinyals et al., 2015; Fang et al., 2015; Karpathy and Fei-Fei, 2015; Xu et al., 2015; You et al., 2016; Krause et al., 2017).

Recently, attention mechanisms have been shown to be useful for image captioning (Xu et al., 2015; You et al., 2016). Xu et al. (2015) introduce a spatial-visual attention mechanism over image features extracted from intermediate layers of the CNN. You et al. (2016) propose a semantic attention mechanism over tags of given images. To better leverage both the visual features and semantic tags, we propose a co-attention mechanism for report generation.

Instead of only generating one-sentence caption

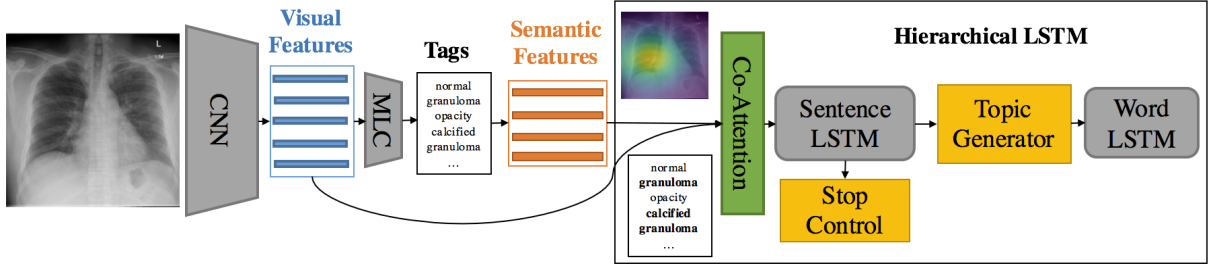


Figure 2: Illustration of the proposed model. MLC denotes a *multi-label classification* network. Semantic features are the word embeddings of the predicted tags. The boldfaced tags “calcified granuloma” and “granuloma” are attended by the co-attention network.

for images, Krause et al. (2017) and Liang et al. (2017) generate paragraph captions using a hierarchical LSTM. Our method also adopts a hierarchical LSTM for paragraph generation, but unlike Krause et al. (2017), we use a co-attention network to generate topics.

### 3 Methods

#### 3.1 Overview

A complete diagnostic report for a medical image is comprised of both text descriptions (long paragraphs) and lists of tags, as shown in Figure 1. We propose a *multi-task hierarchical model with co-attention* for automatically predicting keywords and generating long paragraphs. Given an image which is divided into regions, we use a CNN to learn visual features for these patches. Then these visual features are fed into a *multi-label classification* (MLC) network to predict the relevant tags. In the tag vocabulary, each tag is represented by a word-embedding vector. Given the predicted tags for a specific image, their word-embedding vectors serve as the semantic features of this image. Then the visual features and semantic features are fed into a *co-attention* model to generate a context vector that simultaneously captures the visual and semantic information of this image. As of now, the encoding process is completed.

Next, starting from the context vector, the decoding process generates the text descriptions. The description of a medical image usually contains multiple sentences, and each sentence focuses on one specific topic. Our model leverages this compositional structure to generate reports in a hierarchical way: it first generates a sequence of high-level topic vectors representing sentences, then generates a sentence from each topic vector. Specifically, the context vector is inputted into a

*sentence LSTM*, which unrolls for a few steps and produces a topic vector at each step. A topic vector represents the semantics of a sentence to be generated. Given a topic vector, the *word LSTM* takes it as input and generates a sequence of words to form a sentence. The termination of the unrolling process is controlled by the sentence LSTM.

#### 3.2 Tag Prediction

The first task of our model is predicting the tags of the given image. We treat the tag prediction task as a *multi-label classification* task. Specifically, given an image  $I$ , we first extract its features  $\{\mathbf{v}_n\}_{n=1}^N \in \mathbb{R}^D$  from an intermediate layer of a CNN, and then feed  $\{\mathbf{v}_n\}_{n=1}^N$  into a *multi-label classification* (MLC) network to generate a distribution over all of the  $L$  tags:

$$\mathbf{p}_{l, pred}(\mathbf{l}_i = 1 | \{\mathbf{v}_n\}_{n=1}^N) \propto \exp(\text{MLC}_i(\{\mathbf{v}_n\}_{n=1}^N)) \quad (1)$$

where  $\mathbf{l} \in \mathbb{R}^L$  is a tag vector,  $\mathbf{l}_i = 1/0$  denote the presence and absence of the  $i$ -th tag respectively, and  $\text{MLC}_i$  means the  $i$ -th output of the MLC network.

For simplicity, we extract visual features from the last convolutional layer of the VGG-19 model (Simonyan and Zisserman, 2014) and use the last two fully connected layers of VGG-19 for MLC.

Finally, the embeddings of the  $M$  most likely tags  $\{\mathbf{a}_m\}_{m=1}^M \in \mathbb{R}^E$  are used as semantic features for topic generation.

#### 3.3 Co-Attention

Previous works have shown that visual attention alone can perform fairly well for localizing objects (Ba et al., 2015) and aiding caption generation (Xu et al., 2015). However, visual attention

does not provide sufficient high level semantic information. For example, only looking at the right lower region of the chest x-ray image (Figure 1) without accounting for other areas, we might not be able to recognize what we are looking at, not to even mention detecting the abnormalities. In contrast, the tags can always provide the needed high level information. To this end, we propose a co-attention mechanism which can simultaneously attend to visual and semantic modalities.

In the sentence LSTM at time step  $s$ , the joint context vector  $\text{ctx}^{(s)} \in \mathbb{R}^C$  is generated by a co-attention network  $f_{\text{coatt}}(\{\mathbf{v}_n\}_{n=1}^N, \{\mathbf{a}_m\}_{m=1}^M, \mathbf{h}_{\text{sent}}^{(s-1)})$ , where  $\mathbf{h}_{\text{sent}}^{(s-1)} \in \mathbb{R}^H$  is the sentence LSTM hidden state at time step  $s - 1$ . The co-attention network  $f_{\text{coatt}}$  uses a single layer feed-forward network to compute the soft visual attentions and soft semantic attentions over input image features and tags:

$$\alpha_{\mathbf{v},n} \propto \exp(\mathbf{W}_{\mathbf{vatt}} \tanh(\mathbf{W}_{\mathbf{v}} \mathbf{v}_n + \mathbf{W}_{\mathbf{v,h}} \mathbf{h}_{\text{sent}}^{(s-1)})) \quad (2)$$

$$\alpha_{\mathbf{a},m} \propto \exp(\mathbf{W}_{\mathbf{aatt}} \tanh(\mathbf{W}_{\mathbf{a}} \mathbf{a}_m + \mathbf{W}_{\mathbf{a,h}} \mathbf{h}_{\text{sent}}^{(s-1)})) \quad (3)$$

where  $\mathbf{W}_{\mathbf{v}}$ ,  $\mathbf{W}_{\mathbf{v,h}}$ , and  $\mathbf{W}_{\mathbf{vatt}}$  are parameter matrices of the visual attention network.  $\mathbf{W}_{\mathbf{a}}$ ,  $\mathbf{W}_{\mathbf{a,h}}$ , and  $\mathbf{W}_{\mathbf{aatt}}$  are parameter matrices of the semantic attention network.

The visual and semantic context vectors are computed as:

$$\mathbf{v}_{\text{att}}^{(s)} = \sum_{n=1}^N \alpha_{\mathbf{v},n} \mathbf{v}_n, \quad \mathbf{a}_{\text{att}}^{(s)} = \sum_{m=1}^M \alpha_{\mathbf{a},m} \mathbf{a}_m.$$

There are many ways to combine the visual and semantic context vectors such as concatenation and element-wise operations. In this paper, we first concatenate these two vectors as  $[\mathbf{v}_{\text{att}}^{(s)}; \mathbf{a}_{\text{att}}^{(s)}]$ , and then use a fully connected layer  $\mathbf{W}_{\text{fc}}$  to obtain a joint context vector:

$$\text{ctx}^{(s)} = \mathbf{W}_{\text{fc}}[\mathbf{v}_{\text{att}}^{(s)}; \mathbf{a}_{\text{att}}^{(s)}]. \quad (4)$$

### 3.4 Sentence LSTM

The sentence LSTM is a single-layer LSTM that takes the joint context vector  $\text{ctx} \in \mathbb{R}^C$  as its input, and generates topic vector  $\mathbf{t} \in \mathbb{R}^K$  for word LSTM through topic generator and determines whether to continue or stop generating captions by a stop control component.

**Topic generator** We use a deep output layer (Pascanu et al., 2014) to strengthen the context information in topic vector  $\mathbf{t}^{(s)}$ , by combining the hidden state  $\mathbf{h}_{\text{sent}}^{(s)}$  and the joint context vector  $\text{ctx}^{(s)}$  of the current step:

$$\mathbf{t}^{(s)} = \tanh(\mathbf{W}_{\mathbf{t,h_{sent}}} \mathbf{h}_{\text{sent}}^{(s)} + \mathbf{W}_{\mathbf{t,ctx}} \text{ctx}^{(s)}) \quad (5)$$

where  $\mathbf{W}_{\mathbf{t,h_{sent}}}$  and  $\mathbf{W}_{\mathbf{t,ctx}}$  are weight parameters.

**Stop control** We also apply a deep output layer to control the continuation of the sentence LSTM. The layer takes the previous and current hidden state  $\mathbf{h}_{\text{sent}}^{(s-1)}$ ,  $\mathbf{h}_{\text{sent}}^{(s)}$  as input and produces a distribution over  $\{\text{STOP}=1, \text{CONTINUE}=0\}$ :

$$p(\text{STOP}|\mathbf{h}_{\text{sent}}^{(s-1)}, \mathbf{h}_{\text{sent}}^{(s)}) \propto \exp\{\mathbf{W}_{\text{stop}} \tanh(\mathbf{W}_{\text{stop},s-1} \mathbf{h}_{\text{sent}}^{(s-1)} + \mathbf{W}_{\text{stop},s} \mathbf{h}_{\text{sent}}^{(s)})\} \quad (6)$$

where  $\mathbf{W}_{\text{stop}}$ ,  $\mathbf{W}_{\text{stop},s-1}$  and  $\mathbf{W}_{\text{stop},s}$  are parameter matrices. If  $p(\text{STOP}|\mathbf{h}_{\text{sent}}^{(s-1)}, \mathbf{h}_{\text{sent}}^{(s)})$  is greater than a predefined threshold (e.g. 0.5), then the sentence LSTM will stop producing new topic vectors and the word LSTM will also stop producing words.

### 3.5 Word LSTM

The words of each sentence are generated by a word LSTM. Similar to (Krause et al., 2017), the topic vector  $\mathbf{t}$  produced by the sentence LSTM and the special *START* token are used as the first and second input of the word LSTM, and the subsequent inputs are the word sequence.

The hidden state  $\mathbf{h}_{\text{word}} \in \mathbb{R}^H$  of the word LSTM is directly used to predict the distribution over words:

$$p(\text{word}|\mathbf{h}_{\text{word}}) \propto \exp(\mathbf{W}_{\text{out}} \mathbf{h}_{\text{word}}) \quad (7)$$

where  $\mathbf{W}_{\text{out}}$  is the parameter matrix. After each word-LSTM has generated its word sequences, the final report is simply the concatenation of all the generated sequences.

### 3.6 Parameter Learning

Each training example is a tuple  $(I, \mathbf{l}, \mathbf{w})$  where  $I$  is an image,  $\mathbf{l}$  denotes the ground-truth tag vector, and  $\mathbf{w}$  is the diagnostic paragraph, which is comprised of  $S$  sentences and each sentence consists of  $T_s$  words.



Given a training example  $(I, \mathbf{l}, \mathbf{w})$ , our model first performs multi-label classification on  $I$  and produces a distribution  $\mathbf{p}_{\mathbf{l},pred}$  over all tags. Note that  $\mathbf{l}$  is a binary vector which encodes the presence and absence of tags. We can obtain the ground-truth tag distribution by normalizing  $\mathbf{l}$ :  $\mathbf{p}_{\mathbf{l}} = \mathbf{l} / \|\mathbf{l}\|_1$ . The training loss of this step is a cross-entropy loss  $\ell_{tag}$  between  $\mathbf{p}_{\mathbf{l}}$  and  $\mathbf{p}_{\mathbf{l},pred}$ .

Next, the sentence LSTM is unrolled for  $S$  steps to produce topic vectors and also distributions over  $\{STOP, CONTINUE\}$ :  $p_{stop}^s$ . Finally, the  $S$  topic vectors are fed into the word LSTM to generate words  $\mathbf{w}_{s,t}$ . The training loss of caption generation is the combination of two cross-entropy losses:  $\ell_{sent}$  over stop distributions  $p_{stop}^s$  and  $\ell_{word}$  over word distributions  $p_{s,t}$ . Combining the pieces together, we obtain the overall training loss:

$$\begin{aligned} \ell(I, \mathbf{l}, \mathbf{w}) = & \lambda_{tag} \ell_{tag} \\ & + \lambda_{sent} \sum_{s=1}^S \ell_{sent}(p_{stop}^s, I\{s=S\}) \\ & + \lambda_{word} \sum_{s=1}^S \sum_{t=1}^{T_s} \ell_{word}(p_{s,t}, w_{s,t}) \end{aligned} \quad (8)$$

In addition to the above training loss, there is also a regularization term for visual and semantic attentions. Similar to (Xu et al., 2015), let  $\alpha \in \mathbb{R}^{N \times S}$  and  $\beta \in \mathbb{R}^{M \times S}$  be the matrices of visual and semantic attentions respectively, then the regularization loss over  $\alpha$  and  $\beta$  is:

$$\ell_{reg} = \lambda_{reg} \left[ \sum_n^N \left(1 - \sum_s^S \alpha_{n,s}\right)^2 + \sum_m^M \left(1 - \sum_s^S \beta_{m,s}\right)^2 \right] \quad (9)$$

Such regularization encourages the model to pay equal attention over different image regions and different tags.

## 4 Experiments

In this section, we evaluate the proposed model with extensive quantitative and qualitative experiments.

### 4.1 Datasets

We used two publicly available medical image datasets to evaluate our proposed model.

**IU X-Ray** The Indiana University Chest X-Ray Collection (IU X-Ray) (Demner-Fushman et al., 2015) is a set of chest x-ray images paired with their corresponding diagnostic reports. The

dataset contains 7,470 pairs of images and reports. Each report consists of the following sections: *impression*, *findings*, *tags*<sup>1</sup>, *comparison*, and *indication*. In this paper, we treat the contents in *impression* and *findings* as the target captions<sup>2</sup> to be generated and the Medical Text Indexer (MTI) annotated tags as the target tags to be predicted (Figure 1 provides an example).

We preprocessed the data by converting all tokens to lowercases, removing all of non-alpha tokens, which resulting in 572 unique tags and 1915 unique words. On average, each image is associated with 2.2 tags, 5.7 sentences, and each sentence contains 6.5 words. Besides, we find that top 1,000 words cover 99.0% word occurrences in the dataset, therefore we only included top 1,000 words in the dictionary. Finally, we randomly selected 500 images for validation and 500 images for testing.

**PEIR Gross** The Pathology Education Informational Resource (PEIR) digital library<sup>3</sup> is a public medical image library for medical education. We collected the images together with their descriptions in the Gross sub-collection, resulting in the PEIR Gross dataset that contains 7,442 image-caption pairs from 21 different sub-categories. Different from the IU X-Ray dataset, each caption in PEIR Gross contains only one sentence. We used this dataset to evaluate our model’s ability of generating single-sentence report.

For PEIR Gross, we applied the same preprocessing as IU X-Ray, which yields 4,452 unique words. On average, each image contains 12.0 words. Besides, for each caption, we selected 5 words with the highest tf-idf scores as tags.

### 4.2 Implementation Details

We used the full VGG-19 model (Simonyan and Zisserman, 2014) for tag prediction. As for the training loss of the multi-label classification (MLC) task, since the number of tags for semantic attention is fixed as 10, we treat MLC as a multi-label retrieval task and adopt a softmax cross-entropy loss (a multi-label ranking loss), similar to (Gong et al., 2013; Guillaumin et al., 2009).

<sup>1</sup>There are two types of tags: manually generated (MeSH) and Medical Text Indexer (MTI) generated.

<sup>2</sup>The *impression* and *findings* sections are concatenated together as a long paragraph, since *impression* can be viewed as a conclusion or topic sentence of the report.

<sup>3</sup>PEIR is ©University of Alabama at Birmingham, Department of Pathology. (<http://peir.path.uab.edu/library/>)

Dataset	Methods	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE	CIDER
IU X-Ray	CNN-RNN (Vinyals et al., 2015)	0.316	0.211	0.140	0.095	0.159	0.267	0.111
	LRCN (Donahue et al., 2015)	0.369	0.229	0.149	0.099	0.155	0.278	0.190
	Soft ATT (Xu et al., 2015)	0.399	0.251	0.168	0.118	0.167	0.323	0.302
	ATT-RK (You et al., 2016)	0.369	0.226	0.151	0.108	0.171	0.323	0.155
	Ours-no-Attention	0.505	0.383	0.290	0.224	0.200	0.420	0.259
	Ours-Semantic-only	0.504	0.371	0.291	0.230	0.207	0.418	0.286
	Ours-Visual-only	0.507	0.373	0.297	0.238	0.211	0.426	0.300
	Ours-CoAttention	0.517	0.386	0.306	0.247	0.217	0.447	0.327
PEIR Gross	CNN-RNN (Vinyals et al., 2015)	0.247	0.178	0.134	0.092	0.129	0.247	0.205
	LRCN (Donahue et al., 2015)	0.261	0.184	0.136	0.088	0.135	0.254	0.203
	Soft ATT (Xu et al., 2015)	0.283	0.212	0.163	0.113	0.147	0.271	0.276
	ATT-RK (You et al., 2016)	0.274	0.201	0.154	0.104	0.141	0.264	0.279
	Ours-No-Attention	0.248	0.180	0.133	0.093	0.131	0.242	0.206
	Ours-Semantic-only	0.263	0.191	0.145	0.098	0.138	0.261	0.274
	Ours-Visual-only	0.284	0.209	0.156	0.105	0.149	0.274	0.280
	Ours-CoAttention	0.300	0.218	0.165	0.113	0.149	0.279	0.329

Table 1: Main results for paragraph generation on the IU X-Ray dataset (upper part), and single sentence generation on the PEIR Gross dataset (lower part). BLUE-n denotes the BLEU score that uses up to n-grams.

In paragraph generation, we set the dimensions of all hidden states and word embeddings as 512. For words and tags, different embedding matrices were used since a tag might contain multiple words. We utilized the embeddings of the 10 most likely tags as the semantic feature vectors  $\{\mathbf{a}_m\}_{m=1}^{M=10}$ . We extracted the visual features from the last convolutional layer of the VGG-19 network, which yields a  $14 \times 14 \times 512$  feature map.

We used the Adam (Kingma and Ba, 2014) optimizer for parameter learning. The learning rates for the CNN (VGG-19) and the hierarchical LSTM were  $1e-5$  and  $5e-4$  respectively. The weights ( $\lambda_{tag}$ ,  $\lambda_{sent}$ ,  $\lambda_{word}$  and  $\lambda_{reg}$ ) of different losses were set to 1.0. The threshold for stop control was 0.5. Early stopping was used to prevent over-fitting.

### 4.3 Baselines

We compared our method with several state-of-the-art image captioning models: CNN-RNN (Vinyals et al., 2015), LRCN (Donahue et al., 2015), Soft ATT (Xu et al., 2015), and ATT-RK (You et al., 2016). We re-implemented all of these models and adopt VGG-19 (Simonyan and Zisserman, 2014) as the CNN encoder. Considering these models are built for single sentence captions and to better show the effectiveness of the hierarchical LSTM and the attention mechanism for paragraph generation, we also implemented a hierarchical model without any attention: Ours-no-Attention. The input of Ours-no-Attention is the overall image feature of VGG-19, which has a dimension of 4096. Ours-no-Attention can be viewed as a CNN-RNN (Vinyals et al., 2015) equipped with a hierarchical LSTM decoder. To

further show the effectiveness of the proposed co-attention mechanism, we also implemented two ablated versions of our model: Ours-Semantic-only and Ours-Visual-only, which takes solely the semantic attention or visual attention context vector to produce topic vectors.

### 4.4 Quantitative Results

We report the paragraph generation (upper part of Table 1) and one sentence generation (lower part of Table 1) results using the standard image captioning evaluation tool <sup>4</sup> which provides evaluation on the following metrics: BLEU (Papineni et al., 2002), METEOR (Denkowski and Lavie, 2014), ROUGE (Lin, 2004), and CIDER (Vedantam et al., 2015).

For paragraph generation, as shown in the upper part of Table 1, it is clear that models with a single LSTM decoder perform much worse than those with a hierarchical LSTM decoder. Note that the only difference between Ours-no-Attention and CNN-RNN (Vinyals et al., 2015) is that Ours-no-Attention adopts a hierarchical LSTM decoder while CNN-RNN (Vinyals et al., 2015) adopts a single-layer LSTM. The comparison between these two models directly demonstrates the effectiveness of the hierarchical LSTM. This result is not surprising since it is well-known that a single-layer LSTM cannot effectively model long sequences (Liu et al., 2015; Martin and Cundy, 2018). Additionally, employing semantic attention alone (Ours-Semantic-only) or visual attention alone (Ours-Visual-only) to generate topic vectors does not seem to help caption generation a lot. The potential reason might be that visual at-

<sup>4</sup><https://github.com/tylin/coco-caption>


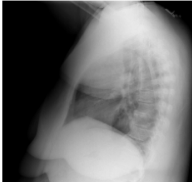


Ground Truth	Ours-CoAttention	Ours-no-Attention	Soft Attention
	No active disease. The heart and lungs have in the interval. <u>Lungs are clear and expanded.</u> Cardiomedastinal silhouette is within normal limits. No pleural effusion or pneumothorax is seen. No pleural effusion. No cavitary or pneumothorax.	The lungs are clear bilaterally. The are grossly normal. No focal lung consolidation. No acute bony abnormality. <u>cm nodule within the right lower lobe on the lateral view.</u> No pneumothorax or pleural effusion. No acute bony abnormality. The heart is not enlarged. The lungs are clear. No acute bony abnormality.	No acute cardiopulmonary abnormality. The lungs are clear bilaterally. Specifically no evidence of focal airspace consolidation pleural effusion or pneumothorax. Cardio mediastinal silhouette is unremarkable. Visualized osseous structures of the thorax are without acute abnormality.
	No evidence of active disease. The lungs are clear. There is no focal airspace consolidation. No pleural effusion or pneumothorax. Heart size and mediastinal contour are within normal limits. <u>There are multilevel degenerative changes of the spine.</u>	No acute cardiopulmonary findings. Heart size is not enlarged. No focal airspace consolidation suspicious pulmonary opacity large pleural effusion or pneumothorax. No focal areas of consolidation. <u>Degenerative changes of the spine.</u> This is moderate exam of the hydropneumothorax. Lungs are clear. There is no focal airspace consolidation pleural effusion or pneumothorax.	No acute cardiopulmonary abnormality. The lungs are clear bilaterally. There is no pleural effusion or pneumothorax. The heart and mediastinum are normal. There is no focal air space opacity to suggest a pneumonia.
	No acute cardiopulmonary abnormality. Normal heart size mediastinal contours. <u>Eventration of the right hemidiaphragm.</u> No focal airspace consolidation. No pleural effusion or pneumothorax.	No acute cardiopulmonary abnormality. Stable appearance of the thoracic aorta. <u>The right lateral lower lobe is noted in the right lower right midlung.</u> No large pleural effusion or focal airspace disease. <u>Mild interstitial opacities.</u> <u>Atherosclerotic calcifications bony structures bilaterally.</u> There is no pleural effusion or pneumothorax developed in the right lower lobe.	No acute cardiopulmonary abnormality. The lungs are clear bilaterally. There is no focal airspace consolidation. No pleural effusion or pneumothorax. Heart size and pulmonary vascularity appear within normal limits.
	No acute cardiopulmonary abnormality. Heart size appears within normal limits. Pulmonary vasculature appears within normal limits. <u>Overlying the middle cardiac silhouette representing a hiatal hernia.</u> No focal consolidation pleural effusion or pneumothorax. No acute bony abnormality.	No active disease. The heart and lungs have in the interval. <u>Nipple and lateral lucency in the lungs suggestive of focal airspace disease.</u> <u>The lungs are hyperexpanded consistent with emphysema in the left lower lobe.</u> This is most at the upper lobes. <u>This may indicate hypoventilated irregularities or effusions.</u> The lungs are otherwise grossly clear. Resolution of by normal pleural effusion.	No acute cardiopulmonary abnormality. The lungs are clear bilaterally. The are grossly normal. No focal airspace consolidation. No pneumothorax or pleural effusion. Heart size and pulmonary vascularity within normal limits. There is no pneumothorax or pleural effusion.

Figure 3: Illustration of paragraph generated by Ours-CoAttention, Ours-no-Attention, and Soft Attention models. The underlined sentences are the descriptions of detected abnormalities. The second image is a lateral x-ray image. Top two images are positive results, the third one is a partial failure case and the bottom one is failure case. These images are from test dataset.

tention can only capture the visual information of sub-regions of the image and is unable to correctly capture the semantics of the entire image. Semantic attention is inadequate of localizing small abnormal image-regions. Finally, our full model (Ours-CoAttention) achieves the best results on all of the evaluation metrics, which demonstrates the effectiveness of the proposed co-attention mechanism.

For the single-sentence generation results (shown in the lower part of Table 1), the ablated versions of our model (Ours-Semantic-only and Ours-Visual-only) achieve competitive scores compared with the state-of-the-art methods. Our full model (Ours-CoAttention) outperforms all of the baseline, which indicates the effectiveness of the proposed co-attention mechanism.

## 4.5 Qualitative Results

### 4.5.1 Paragraph Generation

An illustration of paragraph generation by three models (Ours-CoAttention, Ours-no-Attention and Soft Attention models) is shown in Figure 3.

We can find that different sentences have different topics. The first sentence is usually a high level description of the image, while each of the following sentences is associated with one area of the image (e.g. “lung”, “heart”). Soft Attention and Ours-no-Attention models detect only a few abnormalities of the images and the detected abnormalities are incorrect. In contrast, Ours-CoAttention model is able to correctly describe many true abnormalities (as shown in top three images). This comparison demonstrates that co-attention is better at capturing abnormalities.

For the third image, Ours-CoAttention model successfully detects the area (“right lower lobe”) which is abnormal (“eventration”), however, it fails to precisely describe this abnormality. In addition, the model also finds abnormalities about “interstitial opacities” and “atherosclerotic calcification”, which are not considered as true abnormality by human experts. The potential reason for this mis-description might be that this x-ray image is darker (compared with the above images), and our model might be very sensitive to this change.










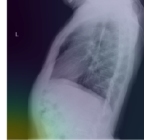
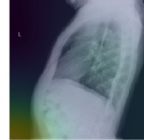
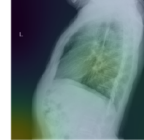
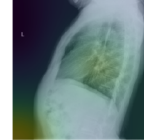
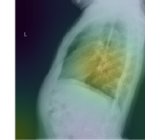




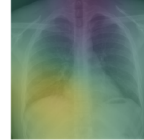
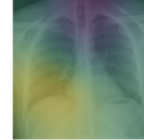

						
degenerative change; obstruction	normal; <u>degenerative change</u> ; nodule; calcified granuloma; hyper expansion; granulomatous disease; <u>pneumonia</u> ; <u>scarring</u> ; <u>sternotomy</u>	normal; <u>degenerative change</u> ; nodule; calcified granuloma; hyper expansion; granulomatous disease; <u>granuloma</u> ; pneumonia; scarring; sternotomy	normal; <u>degenerative change</u> ; nodule; calcified granuloma; hyper expansion; granulomatous disease; <u>granuloma</u> ; pneumonia; scarring; sternotomy	normal; <u>degenerative change</u> ; nodule; calcified granuloma; hyper expansion; granulomatous disease; <u>granuloma</u> ; pneumonia; scarring; sternotomy	normal; <u>degenerative change</u> ; nodule; calcified granuloma; hyper expansion; granulomatous disease; <u>granuloma</u> ; pneumonia; scarring; sternotomy	normal; <u>degenerative change</u> ; nodule; calcified granuloma; hyper expansion; granulomatous disease; <u>granuloma</u> ; pneumonia; scarring; sternotomy
No acute intrathoracic abnormality.	No bony abnormality.	The cardio mediastinal silhouette is within normal limits for appearance.	No focal areas of pulmonary consolidation.	Breast motion.	There is an age indeterminate deformity of a mid-thoracic vertebral body.	
No acute cardiopulmonary finding. The heart size and cardiopulmonary silhouette is normal. There is no focal airspace opacity pleural effusion or pneumothorax. The obstruction are intact with mild degenerative change in the thoracic spine.						
						
calcified granuloma; granulomas	normal; cardiomegaly; nodule; <u>degenerative change</u> ; <u>granulomatous disease</u> ; opacity; <u>calcified granuloma</u> ; atelectasis; lymph nodes; hiatal hernia	normal; cardiomegaly; nodule; <u>degenerative change</u> ; <u>granulomatous disease</u> ; opacity; <u>calcified granuloma</u> ; atelectasis; lymph nodes; hiatal hernia	normal; cardiomegaly; nodule; <u>degenerative change</u> ; <u>granulomatous disease</u> ; opacity; <u>calcified granuloma</u> ; atelectasis; lymph nodes; hiatal hernia	normal; cardiomegaly; nodule; <u>degenerative change</u> ; <u>granulomatous disease</u> ; opacity; <u>calcified granuloma</u> ; atelectasis; lymph nodes; hiatal hernia	normal; cardiomegaly; nodule; <u>degenerative change</u> ; <u>granulomatous disease</u> ; opacity; <u>calcified granuloma</u> ; atelectasis; lymph nodes; hiatal hernia	normal; cardiomegaly; nodule; <u>degenerative change</u> ; <u>granulomatous disease</u> ; opacity; <u>calcified granuloma</u> ; atelectasis; lymph nodes; hiatal hernia
Clear lungs.	Lungs are clear .	Left lower lung volumes is clear.	The heart is normal size.	A there is no focal airspace disease.	There is mild blunting of cost phrenic.	
No acute cardiopulmonary abnormality. The lungs are clear. There are calcified granulomas. Heart size is normal. No pneumothorax.						
						
normal	normal; <u>calcified granuloma</u> ; <u>granulomatous disease</u> ; <u>granuloma</u> ; scarring; opacity; <u>degenerative change</u> ; <u>sternotomy</u> ; <u>thoracic aorta</u> ; <u>nodule</u>	normal; calcified granuloma; granulomatous disease; granuloma; scarring; opacity; <u>degenerative change</u> ; <u>sternotomy</u> ; <u>thoracic aorta</u> ; <u>nodule</u>	normal; calcified granuloma; granulomatous disease; granuloma; scarring; opacity; <u>degenerative change</u> ; <u>sternotomy</u> ; <u>thoracic aorta</u> ; <u>nodule</u>	normal; calcified granuloma; granulomatous disease; granuloma; scarring; opacity; <u>degenerative change</u> ; <u>sternotomy</u> ; <u>thoracic aorta</u> ; <u>nodule</u>	normal; calcified granuloma; granulomatous disease; granuloma; scarring; opacity; <u>degenerative change</u> ; <u>sternotomy</u> ; <u>thoracic aorta</u> ; <u>nodule</u>	normal; calcified granuloma; granulomatous disease; granuloma; scarring; opacity; <u>degenerative change</u> ; <u>sternotomy</u> ; <u>thoracic aorta</u> ; <u>nodule</u>
Right upper lobe infiltrate.	Lungs are clear .	Stable heart size and aortic contours.	No acute displaced rib fractures.	No focal airspace opacities or consolidation.	No visualized of pneumothorax.	
No acute cardiopulmonary abnormality identified. The examination consists of frontal and lateral radiographs of the chest. The cardio mediastinal contours are within normal limits. Pulmonary vascularity is within normal limits. No focal consolidation pleural effusion or pneumothorax identified. The visualized osseous structures and upper abdomen are unremarkable.						

Figure 4: Visualization of co-attention for three examples. Each example is comprised of four things: (1) image and visual attentions; (2) ground truth tags and semantic attention on predicted tags; (3) generated descriptions; (4) ground truth descriptions. For the semantic attention, three tags with highest attention scores are highlighted. The underlined tags are those appearing in the ground truth.

The image at the bottom is a failure case of Ours-CoAttention. However, even though the model makes the wrong judgment about the major abnormalities in the image, it does find some unusual regions: “lateral lucency” and “left lower lobe”.

To further understand models’ ability of detecting abnormalities, we present the portion of sentences which describe the normalities and abnormalities in Table 2. We consider sentences which contain “no”, “normal”, “clear”, “stable” as sentences describing normalities. It is clear that Ours-CoAttention best approximates the ground truth distribution over normality and abnormality.

Method	Normality	Abnormality	Total
Soft Attention	0.510	0.490	1.0
Ours-no-Attention	0.753	0.247	1.0
Ours-CoAttention	0.471	0.529	1.0
Ground Truth	0.385	0.615	1.0

Table 2: Portion of sentences which describe the normalities and abnormalities in the image.

## 4.5.2 Co-Attention Learning

Figure 4 presents visualizations of co-attention. The first property shown by Figure 4 is that the sentence LSTM can generate different topics at different time steps since the model focuses on different image regions and tags for different sentences. The next finding is that visual attention can guide our model to concentrate on relevant re-



gions of the image. For example, the third sentence of the first example is about “cardio”, and the visual attention concentrates on regions near the heart. Similar behavior can also be found for semantic attention: for the last sentence in the first example, our model correctly concentrates on “degenerative change” which is the topic of the sentence. Finally, the first sentence of the last example presents a mis-description caused by incorrect semantic attention over tags. Such incorrect attention can be reduced by building a better tag prediction module.

## 5 Conclusion

In this paper, we study how to automatically generate textual reports for medical images, with the goal to help medical professionals produce reports more accurately and efficiently. Our proposed methods address three major challenges: (1) how to generate multiple heterogeneous forms of information within a unified framework, (2) how to localize abnormal regions and produce accurate descriptions for them, (3) how to generate long texts that contain multiple sentences or even paragraphs. To cope with these challenges, we propose a multi-task learning framework which jointly predicts tags and generates descriptions. We introduce a co-attention mechanism that can simultaneously explore visual and semantic information to accurately localize and describe abnormal regions. We develop a hierarchical LSTM network that can more effectively capture long-range semantics and produce high quality long texts. On two medical datasets containing radiology and pathology images, we demonstrate the effectiveness of the proposed methods through quantitative and qualitative studies.

## References

- Jimmy Ba, Volodymyr Mnih, and Koray Kavukcuoglu. 2015. Multiple object recognition with visual attention. *ICLR*.
- Louke Delrue, Robert Gosselin, Bart Ilse, An Van Landeghem, Johan de Mey, and Philippe Duyck. 2011. Difficulties in the interpretation of chest radiography. In *Comparative Interpretation of CT and Standard Radiography of the Chest*, pages 27–49. Springer.
- Dina Demner-Fushman, Marc D Kohli, Marc B Rosenman, Sonya E Shooshan, Laritza Rodriguez, Sameer Antani, George R Thoma, and Clement J McDonald. 2015. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310.
- Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the ninth workshop on statistical machine translation*, pages 376–380.
- Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. 2015. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2625–2634.
- Hao Fang, Saurabh Gupta, Forrest Iandola, Rupesh K Srivastava, Li Deng, Piotr Dollár, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C Platt, et al. 2015. From captions to visual concepts and back. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1473–1482.
- Yunchao Gong, Yangqing Jia, Thomas Leung, Alexander Toshev, and Sergey Ioffe. 2013. Deep convolutional ranking for multilabel image annotation. *ICLR*.
- Matthieu Guillaumin, Thomas Mensink, Jakob Verbeek, and Cordelia Schmid. 2009. Tagprop: Discriminative metric learning in nearest neighbor models for image auto-annotation. In *Computer Vision, 2009 IEEE 12th International Conference on*, pages 309–316. IEEE.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Pavel Kisilev, Eugene Walach, Ella Barkan, Boaz Ophir, Sharon Alpert, and Sharbell Y Hashoul. 2015. From medical image to automatic medical report generation. *IBM Journal of Research and Development*, 59(2/3):2–1.
- Jonathan Krause, Justin Johnson, Ranjay Krishna, and Li Fei-Fei. 2017. A hierarchical approach for generating descriptive image paragraphs. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

- Xiaodan Liang, Zhiting Hu, Hao Zhang, Chuang Gan, and Eric P. Xing. 2017. Recurrent topic-transition gan for visual paragraph generation. In *The IEEE International Conference on Computer Vision (ICCV)*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out: Proceedings of the ACL-04 workshop*, volume 8. Barcelona, Spain.
- Pengfei Liu, Xipeng Qiu, Xinchu Chen, Shiyu Wu, and Xuanjing Huang. 2015. Multi-timescale long short-term memory neural network for modelling sentences and documents. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 2326–2335.
- Eric Martin and Chris Cundy. 2018. Parallelizing linear recurrent neural nets over sequence length. *ICLR*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. 2014. How to construct deep recurrent neural networks. *ICLR*.
- Hoo-Chang Shin, Kirk Roberts, Le Lu, Dina Demner-Fushman, Jianhua Yao, and Ronald M Summers. 2016. Learning to read chest x-rays: recurrent neural cascade model for automated image annotation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2497–2506.
- K. Simonyan and A. Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575.
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3156–3164.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International Conference on Machine Learning*, pages 2048–2057.
- Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. 2016. Image captioning with semantic attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4651–4659.
- Zizhao Zhang, Yuanpu Xie, Fuyong Xing, Mason McGough, and Lin Yang. 2017. Mdnnet: A semantically and visually interpretable medical image diagnosis network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6428–6436.