



Attention-Based Abnormal-Aware Fusion Network for Radiology Report Generation

Xiancheng Xie¹, Yun Xiong^{1,2(✉)}, Philip S. Yu^{1,2,3}, Kangan Li⁴,
Suhua Zhang⁵, and Yangyong Zhu^{1,2}

¹ Shanghai Key Laboratory of Data Science, School of Computer Science,
Fudan University, Shanghai, China
{17212010043,yunx}@fudan.edu.cn

² Shanghai Institute for Advanced Communication and Data Science,
Fudan University, Shanghai, China

³ Computer Science Department, University of Illinois at Chicago,
Chicago, IL 60607, USA

⁴ Shanghai General Hospital, Shanghai, China

⁵ Department of Nephrology, Suzhou Kowloon Hospital, Shanghai Jiaotong
University, School of Medicine, Suzhou 2015028, Jiangsu Province, China

Abstract. Radiology report writing is error-prone, time-consuming and tedious for radiologists. Medical reports are usually dominated by a large number of normal findings, and the abnormal findings are few but more important. Current report generation methods often fail to depict these prominent abnormal findings. In this paper, we propose a model named Attention-based Abnormal-Aware Fusion Network (A3FN). We break down sentence generation into abnormal and normal sentence generation through a high level gate module. We also adopt a topic guide attention mechanism for better capturing visual details and develop a context-aware topic vector for model cross-sentence topic coherence. Experiments on real radiology image datasets demonstrate the effectiveness of our proposed method.

Keywords: Radiology report generation · Attention · Abnormal-aware

1 Introduction

In recent years, radiology images are playing a vital role in the auxiliary diagnosis. Most of the existing captioning methods [3, 4, 6] perform poorly since report usually consists of multiple long, structural and informative sentences. Besides generating short captions, pioneering radiology report generation method [2] adopts a hierarchical LSTM framework combined with a co-attention mechanism to generate paragraph on IU X-Ray dataset [1]. However, their generated reports tend to describe normal findings with some repetitions and are incapable of capturing rare but prominent abnormality. Xue et al. [5] adopt a feedback

mechanism to learn long sequence dependency. However, their method is prone to generate fluent but general report without prominent abnormal narratives.

In this paper, we propose an Attention-based Abnormal-Aware Fusion Network (A3FN). Our method can effectively capture these rare but prominent abnormal observations. For more accurate abnormal descriptions, we adopt a topic guide attention mechanism to support abnormal findings with its detailed visual context (e.g., location, size, and severity etc.). We also develop a context-aware topic vector to control topic coherence and descriptive completeness. Our proposed A3FN method achieves the highest detection accuracy of positive abnormality terminologies.

2 The Proposed A3FN Method

This section describes our proposed A3FN in detail. Figure 1 shows the architecture of A3FN. We will go through vital parts of the model in following section.

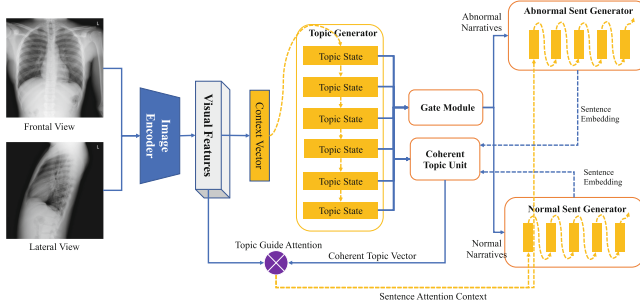


Fig. 1. The architecture of the proposed A3FN method.

Topic Generator and Coherent Topic Unit. Topic generator is a single-layer LSTM which generates a sequence of high-level topic vectors \mathbf{q} , one for each sentence. Generally, topic generator can be written as:

$$\mathbf{c}_i^t = F_{attn}^t(\mathbf{h}^c, \mathbf{h}_{i-1}^t), \mathbf{h}_i^t = F_{RNN}^t(\mathbf{c}_i^t, \mathbf{h}_{i-1}^t),$$

$$\mathbf{q}_i = \tanh(\mathbf{W}_t \mathbf{h}_i^t + \mathbf{b}_t), \mathbf{y}_i = \text{softmax}(\mathbf{W}_{stop} \mathbf{h}_i^t + \mathbf{b}_{stop}).$$

where F_{attn}^t is the attention mechanism same as [4], and F_{RNN}^t is a single-layer LSTM. The probability \mathbf{y}_i controls whether to stop generating next topic vector. After generate all topic states, we enhance current topic state \mathbf{q}_i with global topic vector \mathbf{q}_g and last generated sentence embedding \mathbf{r}_{i-1} for topic coherence as follows:

$$\mathbf{q}_g = \sum_{i=1}^M \alpha_i \mathbf{q}_i, \text{ where } \alpha_i = \frac{\|\mathbf{q}_i\|_2}{\sum_i \|\mathbf{q}_i\|_2}.$$

We first transform sentence embedding \mathbf{r}_{i-1} to \mathbf{q}_{i-1}^h through a 2-layer fully-connected unit. Then we merge \mathbf{q}_{i-1}^h and \mathbf{q}_i as follows:

$$\mathbf{q}_i' = \beta_i \mathbf{q}_i + (1 - \beta_i) \mathbf{q}_{i-1}^h, \beta_i = \text{sigmoid}(\mathbf{W}_{c,i} \mathbf{q}_i - \mathbf{W}_{c,i-1} \mathbf{q}_{i-1}^h).$$

The obtained vector \mathbf{q}_i' is then coupled with global topic vector \mathbf{q}_g using a gating function implemented by a single-layer GRU.

$$\mathbf{q}_i^c = F_{GRU}^c(\mathbf{q}_i', \mathbf{q}_g).$$

Gate Module. Given the i -th topic vector \mathbf{q}_i , previous topic vector \mathbf{q}_{i-1} , the gate module generate a distribution \mathbf{u}_i over $\{abnormal = 0, normal = 1\}$, that is:

$$\mathbf{u}_i = \text{softmax}(\mathbf{W}_u \tanh(\mathbf{W}_{u,i} \mathbf{q}_i + \mathbf{W}_{u,i-1} \mathbf{q}_{i-1})).$$

Attentional Sentence Generator. Different from [2] that directly feed topic vector into generator, we enhance our sentence generator with a *topic guide attention module*. More specifically, coherent topic vector \mathbf{q}_i^c , regional visual feature \mathbf{v} and previous hidden state $\mathbf{h}_{i,j-1}^s$ are fed into a fully connected layer, followed by a softmax to get the attention distribution over k regions as follows:

$$\mathcal{A} = \mathbf{W}_{att} \tanh(\mathbf{W}_v^{att} \mathbf{v} + \mathbf{W}_t^{att} \mathbf{q}_i^c \mathbf{1}^k + \mathbf{W}_s^{att} \mathbf{h}_{i,j-1}^s \mathbf{1}^k).$$

$$\alpha_n = \frac{\exp(\mathcal{A}_n)}{\sum_n \exp(\mathcal{A}_n)}, \mathbf{c}_{i,j}^s = \sum_{n=1}^k \alpha_n \mathbf{v}_n.$$

The sentence generator is a single-layer LSTM which takes $\mathbf{e}_{i,j-1}, \mathbf{q}_i^c, \mathbf{c}_{i,j}^s$ as inputs:

$$\begin{aligned} \mathbf{h}_{i,j}^s &= F_{RNN}^s(\mathbf{e}_{i,j-1}, [\mathbf{q}_i^c, \mathbf{c}_{i,j}^s], \mathbf{h}_{i,j-1}^s), \\ \mathbf{a}_{i,j} &= \text{softmax}(\mathbf{W}_s \mathbf{h}_{i,j}^s + \mathbf{b}_s), s_{i,j} = \text{argmax}(\mathbf{a}_{i,j}), \mathbf{e}_{i,j} = \mathbf{E} \mathbf{1}_{s_{i,j}}. \end{aligned}$$

where \mathbf{E} is learnable embedding matrix and $\mathbf{1}_{y_{i,j}}$ is one-hot vector.

3 Experiments and Analysis

Settings. We evaluate the proposed A3FN method on public IU X-Ray dataset [1]. There are 2914 reports associated with 5828 images after filtering out reports without two complete image views. We randomly pick 10% reports as testing set. We compare the proposed A3FN with state-of-the-art report generation methods including CNN-RNN [3], Soft-ATT [4], ATT-RK [6], Co-ATT [2] and Multi-Modal [5]. For comparison, we reimplement baseline models [2, 5] using Pytorch since the codes of these models are not available. We report the performance of all models on frequently used metrics BLEU- $\{1,2,3,4\}$, CIDEr and ROUGE in Table 1. We compute the keywords accuracy (KA) used in [5] and positive abnormality terminology detection accuracy (Acc) as a measurement of abnormal detection. Table 2 shows evaluation results of KA and Acc.

Table 1. Automatic evaluation results.

Dataset	Model	BLEU-1	BLEU-2	BLEU-3	BLEU-4	CIDEr	ROUGE
IU X-Ray	CNN-RNN [3]	0.311	0.218	0.137	0.092	0.124	0.262
	Soft-ATT [4]	0.351	0.237	0.161	0.12	0.278	0.314
	ATT-RK [6]	0.341	0.221	0.153	0.106	0.187	0.302
	Co-ATT [2]	0.421	0.324	0.225	0.174	0.331	0.341
	Multi-Modal [5]	0.434	0.331	0.234	0.177	0.312	0.346
	A3FN	0.443	0.337	0.236	0.181	0.374	0.347

Table 2. Positive abnormality terminology detection accuracy and KA

Dataset	Model	KA (%)	Acc (%)
IU X-Ray	Multi-Modal [5]	58.7	10.1
	Co-ATT [2]	57.3	11.37
	A3FN	61.0	13.25

Results and Analysis. The proposed A3FN method outperforms all other baselines. Compared with these hierarchical models [2, 5], our A3FN method outperforms Co-ATT [2] and Multi-Modal [5] by a large margin, with respectively **11.3%** and **19.9%** relative improvement on CIDEr score. Furthermore, the proposed A3FN model achieves best KA score and best positive abnormality detection accuracy score, which means our generated reports detect more positive abnormality and cover more topics.

4 Conclusion

In this paper, we propose a novel model named Attention-based Abnormal-Aware Attention Network (A3FN) which aims to generate structured, detailed, topic coherent, and abnormal-aware radiology reports. Our model achieves the competitive results compared with all state-of-the-art models.

Acknowledgements. This work is supported in part by the National Natural Science Foundation of China Projects No. U1636207, No. 91546105, the Shanghai Science and Technology Development Fund No. 16JC1400801, No. 17511105502, No. 17511101702, Suzhou Science and Technology Bureau Technology Demonstration Project (SS201712, SS201812).

References

1. Demner-Fushman, D., et al.: Preparing a collection of radiology examinations for distribution and retrieval. *JAMIA* **23**, 304–310 (2015)
2. Jing, B., Xie, P., Xing, E.P.: On the automatic generation of medical imaging reports. In: *Proceedings of ACL* (2018)

3. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: a neural image caption generator. In: Proceedings of CVPR (2015)
4. Xu, K., et al.: Show, attend and tell: neural image caption generation with visual attention. In: Proceedings of ICML (2015)
5. Xue, Y., et al.: Multimodal recurrent model with attention for automated radiology report generation. In: Frangi, A.F., Schnabel, J.A., Davatzikos, C., Alberola-López, C., Fichtinger, G. (eds.) MICCAI 2018. LNCS, vol. 11070, pp. 457–466. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-00928-1_52
6. You, Q., Jin, H., Wang, Z., Fang, C., Luo, J.: Image captioning with semantic attention. In: Proceedings of CVPR (2016)