

Competence-based Multimodal Curriculum Learning for Medical Report Generation

Fenglin Liu¹, Shen Ge², Xian Wu²

¹School of ECE, Peking University

²Tencent Medical AI Lab, Beijing, China

fenglinliu98@pku.edu.cn; {shenge, kevinxwu}@tencent.com

Abstract

Medical report generation task, which targets to produce long and coherent descriptions of medical images, has attracted growing research interests recently. Different from the general image captioning tasks, medical report generation is more challenging for data-driven neural models. This is mainly due to 1) the serious data bias and 2) the limited medical data. To alleviate the data bias and make best use of available data, we propose a Competence-based Multimodal Curriculum Learning framework (CMCL). Specifically, CMCL simulates the learning process of radiologists and optimizes the model in a step by step manner. Firstly, CMCL estimates the difficulty of each training instance and evaluates the competence of current model; Secondly, CMCL selects the most suitable batch of training instances considering current model competence. By iterating above two steps, CMCL can gradually improve the model’s performance. The experiments on the public IU-Xray and MIMIC-CXR datasets show that CMCL can be incorporated into existing models to improve their performance.

1 Introduction

Medical images, e.g., radiology and pathology images, and their corresponding reports, which describe the observations in details of both normal and abnormal regions, are widely-used for diagnosis and treatment (Delrue et al., 2011; Goergen et al., 2013). In clinical practice, writing a medical report can be time-consuming and tedious for experienced radiologists, and error-prone for inexperienced radiologists. Therefore, automatically generating medical reports can assist radiologists in clinical decision-making and emerge as a prominent attractive research direction in both artificial intelligence and clinical medicine (Jing et al., 2018,

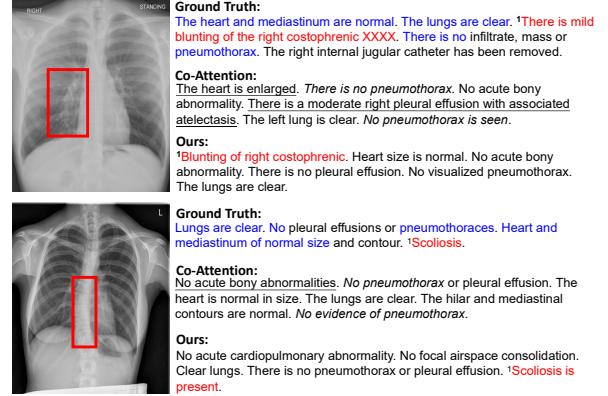


Figure 1: Two examples of ground truth reports and reports generated by a state-of-the-art approach Co-Attention (Jing et al., 2018) and our approach. The Red bounding boxes and Red colored text indicate the abnormalities in images and reports, respectively. The Blue colored text stands for the similar sentences used to describe the normalities in ground truth reports. There are notable visual and textual data biases and the Co-Attention (Jing et al., 2018) fails to depict the rare but important abnormalities and generates some error sentences (Underlined text) and repeated sentences (*Italic* text).

2019; Li et al., 2018, 2019; Wang et al., 2018; Xue et al., 2018; Yuan et al., 2019; Zhang et al., 2020a; Chen et al., 2020; Liu et al., 2021a,b, 2019c).

Many existing medical report generation models adopt the standard image captioning approaches: a CNN-based image encoder followed by a LSTM-based report decoder, e.g., CNN-HLSTM (Jing et al., 2018; Liang et al., 2017). However, directly applying image captioning approaches to medical images has the following problems: 1) **Visual data bias**: the normal images dominate the dataset over the abnormal ones (Shin et al., 2016). Furthermore, for each abnormal image, the normal regions dominate the image over the abnormal ones. As shown in Figure 1, abnormal regions (Red bounding boxes) only occupy a small part of the entire

image; 2) **Textual data bias**: as shown in Figure 1, in a medical report, radiologists tend to describe all the items in an image, making the descriptions of normal regions dominate the entire report. Besides, many similar sentences are used to describe the same normal regions. 3) **Training efficiency**: during training, most existing works treat all the samples equally without considering their difficulties. As a result, the visual and textual biases could mislead the model training (Jing et al., 2019; Xue et al., 2018; Yuan et al., 2019; Liu et al., 2021a,b; Li et al., 2018). As shown in Figure 1, even a state-of-the-art model (Jing et al., 2018) still generates some repeated sentences of normalities and fails to depict the rare but important abnormalities.

To this end, we propose a novel Competence-based Multimodal Curriculum Learning framework (CMCL) which progressively learns medical reports following an easy-to-hard fashion. Such a step by step process is similar to the learning curve of radiologists: (1) first start from simple and easy-written reports; (2) and then attempt to consume harder reports, which consist of rare and diverse abnormalities. In order to model the above gradual working patterns, CMCL first assesses the difficulty of each training instance from multiple perspectives (i.e., the Visual Complexity and Textual Complexity) and then automatically selects the most rewarding training samples according to the current competence of the model. In this way, once the easy and simple samples are well-learned, CMCL increases the chance of learning difficult and complex samples, preventing the models from getting stuck in bad local optima¹, which is obviously a better solution than the common approaches of uniformly sampling training examples from the limited medical data. As a result, CMCL could better utilize the limited medical data to alleviate the data bias. We evaluate the effectiveness of the proposed CMCL on IU-Xray (Demner-Fushman et al., 2016) and MIMIC-CXR (Johnson et al., 2019).

Overall, the main contributions of this work are:

- We introduce the curriculum learning in medical report generation, which enables the models to gradually proceed from easy samples to more complex ones in training, helping existing models better utilize the limited medical data to alleviate the data bias.

¹Current models tend to generate plausible general reports with no prominent abnormal narratives (Jing et al., 2019; Li et al., 2018; Yuan et al., 2019; Liu et al., 2021a,b)

- We assess the difficulty of each training instance from multiple perspectives and propose a competence-based multimodal curriculum learning framework (CMCL) to consider multiple difficulties simultaneously.
- We evaluate our proposed approach on two public datasets. After equipping our proposed CMCL, which doesn't introduce additional parameters and only requires a small modification to the training data pipelines, performances of the existing baseline models can be improved on most metrics. Moreover, we conduct human evaluations to measure the effectiveness in terms of its usefulness for clinical practice.

2 Related Work

In this section, we will introduce the related works from three aspects: 1) Image Captioning and Paragraph Generation; 2) Medical Report Generation and 3) Curriculum Learning.

2.1 Image Captioning and Paragraph Generation

The task of image captioning (Chen et al., 2015; Vinyals et al., 2015), which aims to generate a sentence to describe the given image, has received extensive research interests (Rennie et al., 2017; Liu et al., 2018; Anderson et al., 2018; Liu et al., 2019a,b, 2020a). These approaches mainly adopt the encoder-decoder framework which translates the image to a *single* descriptive sentence. Such an encoder-decoder framework have achieved great success in advancing the state-of-the-arts (Vinyals et al., 2015; Lu et al., 2017; Xu et al., 2015; Liu et al., 2018, 2019b; Cornia et al., 2020; Pan et al., 2020; Liu et al., 2020a). Specifically, the encoder network (Krizhevsky et al., 2012; He et al., 2016) computes visual representations for the visual contents and the decoder network (Hochreiter and Schmidhuber, 1997; Vaswani et al., 2017) generates a target sentence based on the visual representations. In contrast to the image captioning, image paragraph generation, which aims to produce a long and semantic-coherent paragraph to describe the input image, has recently attracted growing research interests (Krause et al., 2017; Liang et al., 2017; Yu et al., 2016). To perform the image paragraph generation, a hierarchical LSTM (HLSTM) (Krause et al., 2017; Liang et al., 2017) is proposed as the decoder to well generate long paragraphs.

2.2 Medical Report Generation

The medical reports are expected to 1) cover contents of key medical findings such as heart size, lung opacity, and bone structure; 2) correctly capture any abnormalities and support with details such as the location and shape of the abnormality; 3) correctly describe potential diseases such as effusion, pneumothorax and consolidation (Delrue et al., 2011; Goergen et al., 2013; Li et al., 2018; Liu et al., 2021a,b). Therefore, correctly describing the abnormalities become the most urgent goal and the core value of this task. Similar to image paragraph generation, most existing medical report generation works (Jing et al., 2018, 2019; Li et al., 2018; Wang et al., 2018; Xue et al., 2018; Yuan et al., 2019; Zhang et al., 2020a,b; Miura et al., 2021; Lovelace and Mortazavi, 2020; Liu et al., 2021b, 2019c) attempt to adopt a CNN-HLSTM based model to automatically generate a fluent report. However, due to the data bias and the limited medical data, these models are biased towards generating plausible but general reports without prominent abnormal narratives (Jing et al., 2019; Li et al., 2018; Yuan et al., 2019; Liu et al., 2021a,b).

2.3 Curriculum Learning

In recent years, curriculum learning (Bengio et al., 2009), which enables the models to gradually proceed from easy samples to more complex ones in training (Elman, 1993), has received growing research interests in natural language processing field, e.g., neural machine translation (Platanios et al., 2019; Kumar et al., 2019; Zhao et al., 2020; Liu et al., 2020b; Zhang et al., 2018; Kocmi and Bojar, 2017; Xu et al., 2020) and computer vision field, e.g., image classification (Weinshall et al., 2018), human attribute analysis (Wang et al., 2019) and visual question answering (Li et al., 2020). For example, in neural machine translation, Platanios et al. (2019) proposed to utilize the training samples in order of easy-to-hard and to describe the “difficulty” of a training sample using the sentence length or the rarity of the words appearing in it (Zhao et al., 2020). However, these methods (Platanios et al., 2019; Liu et al., 2020b; Xu et al., 2020) are single difficulty-based and unimodal curriculum learning approaches. It is obviously not applicable to medical report generation task, which involves multi-modal data, i.e., visual medical images and textual reports, resulting in multi-modal complexities, i.e., the visual complexity and the

textual complexity. Therefore, it is hard to design one single metric to estimate the overall difficulty of medical report generation. To this end, based on the work of Platanios et al. (2019), we propose a competence-based multimodal curriculum learning approach with multiple difficulty metrics.

3 Framework

In this section, we briefly describe typical medical report generation approaches and introduce the proposed Competence-based Multimodal Curriculum Learning (CMCL).

As shown in the top of Figure 2, many medical report generation models adopt the encoder-decoder manner. Firstly, the visual features are extracted from the input medical image via a CNN model. Then the visual features are fed into a sequence generation model, like LSTM to produce the medical report. In the training phase, all training instances are randomly shuffled and grouped into batches for training. In other words, all training instances are treated equally. Different from typical medical report generation models, CMCL builds the training batch in a selective manner. The middle part of Figure 2 displays the framework of CMCL equipped with one single difficulty metric. CMCL first ranks all training instances according to this difficulty metric and then gradually enlarges the range of training instances that the batch is selected. In this manner, CMCL can train the models from easy to difficult instances.

Since medical report generation involves multimodal data, like visual medical images and textual reports, it is hard to design one single metric to estimate the overall difficulty. Therefore, we also propose a CMCL with multiple difficulty metrics. As shown in the bottom of Figure 2, the training instances are ranked by multiple metrics independently. At each step, CMCL generates one batch for each difficulty metric and then calculates the perplexity of each batch based on current model. The batch with highest perplexity is selected to train the model. It can be understood that CMCL sets multiple syllabus in parallel, and the model is optimized towards the one with lowest competence.

4 Difficulty Metrics

In this section, we define the difficulty metrics used by CMCL. As stated in Section 2, the key challenge of medical report generation is to accurately capture and describe the abnormalities (Delrue et al.,

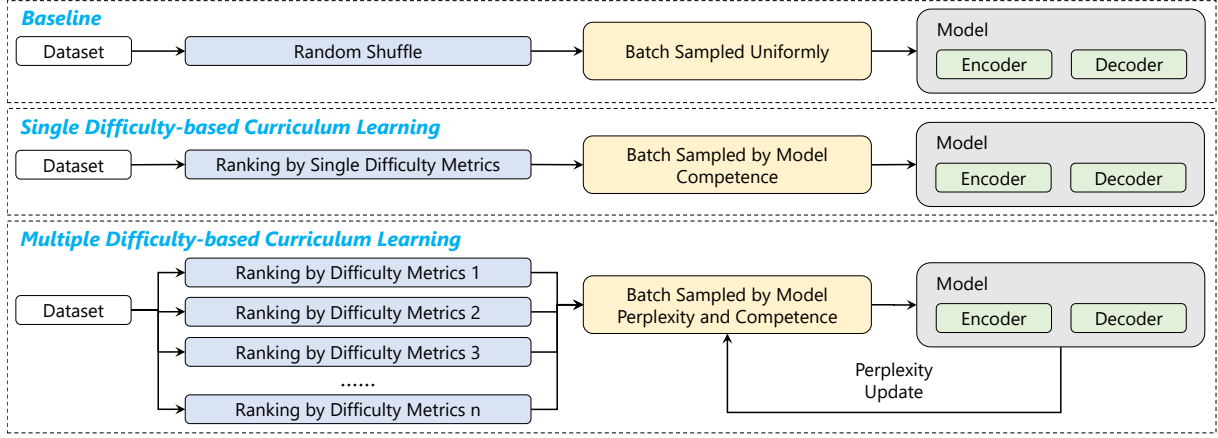


Figure 2: The top illustrates the typical encoder-decoder approach; The middle illustrates the Single Difficulty-based Curriculum Learning, where only one difficulty metric is used; The bottom illustrates the Multiple Difficulty-based Curriculum Learning, where multiple difficulty metrics are introduced.

2011; Goergen et al., 2013; Li et al., 2018). Therefore, we assess the difficulty of instances based on the difficulty of accurately capturing and describing the abnormalities.

4.1 Visual Difficulty

We define both a heuristic metric and a model-based metric to estimate the visual difficulty.

Heuristic Metric d_1 If a medical image contains complex visual contents, it is more likely to contain more abnormalities, which increases the difficulty to accurately capture them. To measure such visual difficulty, we adopt the widely-used ResNet-50 (He et al., 2016) pre-trained on ImageNet (Deng et al., 2009) and fine-tuned on CheXpert dataset (Irvin et al., 2019), which consists of 224,316 X-ray images with each image labeled with occurrences of 14 common radiographic observations. Specifically, we first extract the normal image embeddings of all normal training images from the last average pooling layer of ResNet-50. Then, given an input image, we again use the ResNet-50 to obtain the image embedding. At last, the average cosine similarity between the input image and normal images is adopted as the heuristic metric of visual difficulty.

Model Confidence d_2 We also introduce a model-based metric. We adopt the above ResNet-50 to conduct the abnormality classification task. We first adopt the ResNet-50 to acquire the classification probability distribution $P(I) = \{p_1(I), p_2(I), \dots, p_{14}(I)\}$ among the 14 common diseases for each image I in the training dataset, where $p_n(I) \in [0, 1]$. Then, we employ the entropy

value $H(I)$ of the probability distribution, defined as follows:

$$H(I) = - \sum_{n=1}^{14} (p_n(I) \log(p_n(I)) + (1 - p_n(I)) \log(1 - p_n(I))) \quad (1)$$

We employ the entropy value $H(I)$ as the model confidence measure, indicating whether an image is easy to be classified or not.

4.2 Textual Difficulty

We also define a heuristic metric and a model-based metric to estimate the textual difficulty.

Heuristic Metric d_3 A serious problem for medical report generation models is the tendency to generate plausible general reports with no prominent abnormal narratives (Jing et al., 2019; Li et al., 2018; Yuan et al., 2019). The normal sentences are easy to learn, but are less informative, while most abnormal sentences, consisting of more rare and diverse abnormalities, are relatively more difficult to learn, especially at the initial learning stage. To this end, we adopt the number of abnormal sentences in a report to define the difficulty of a report. Following Jing et al. (2018), we consider sentences which contain “no”, “normal”, “clear”, “stable” as normal sentences, the rest sentences are consider as abnormal sentences.

Model Confidence d_4 Similar to visual difficulty, we further introduce a model confidence as a metric. To this end, we define the difficulty using the negative log-likelihood loss values (Xu et al., 2020; Zhang et al., 2018) of training samples. To

Algorithm 1 Single Difficulty-based Curriculum Learning (Platanios et al., 2019).

Input: The training set D^{train} .

Output: A model with single difficulty-based curriculum learning.

- 1: Compute difficulty d for each training sample in D^{train} ;
 - 2: Sort D^{train} based d to acquire D_1^{train} ;
 - 3: At $t = 0$, initialize the model competence $c(0)$ by Eq. (2); Uniformly sample a data batch, $B(0)$, from the top $c(0)$ portions of D_1^{train} ;
 - 4: **repeat**
 - 5: Train the model with the $B(t)$;
 - 6: $t \leftarrow t + 1$;
 - 7: Estimate the model competence, $c(t)$, by Eq. (2); Uniformly sample a data batch, $B(t)$, from the top $c(t)$ portions of D_1^{train} ;
 - 8: **until** Model converge.
-

acquire the negative log-likelihood loss values, we adopt the widely-used and classic CNN-HLSTM (Jing et al., 2018), in which the CNN is implemented with ResNet-50, trained on the downstream dataset used for evaluation with a cross-entropy loss.

It is worth noting that since we focus on the medical report generation and design the metrics based on the difficulty of accurately capturing and describing the abnormalities, we do not consider some language difficulty metrics used in neural machine translation, e.g., the sentence length (Platanios et al., 2019), the n-gram rarity together with Named Entity Recognition (NER) and Parts of Speech (POS) taggings (Zhao et al., 2020).

5 Approach

In this section, we first briefly introduce the conventional single difficulty-based curriculum (Platanios et al., 2019). Then we propose the multiple difficulty-based curriculum learning for medical report generation.

5.1 Single Difficulty-based Curriculum Learning

Platanios et al. (2019) proposed a competence-based and single difficulty-based curriculum learning framework (see Algorithm 1), which first sorts each instance in the training dataset D^{train} according to a single difficulty metric d , and then defines the model competence $c(t) \in (0, 1]$ at training step t by following functional forms:

$$c(t) = \min \left(1, \sqrt[p]{t \frac{1 - c(0)^p}{T} + c(0)^p} \right) \quad (2)$$

where $c(0)$ is the initial competence and usually set to 0.01, p is the coefficient to control the cur-

Algorithm 2 Multiple Difficulty-based Curriculum Learning. The **Red** colored text denotes the differences from Algorithm 1.

Input: The training set D^{train} , $i \in \{1, 2, 3, 4\}$.

Output: A model with multiple difficulty-based curriculum learning.

- 1: Compute four difficulties, d_i , for each training sample in D^{train} ;
 - 2: Sort D^{train} based each difficulty of every sample, resulting in D_i^{train} (i.e., D_1^{train} , D_2^{train} , D_3^{train} , D_4^{train});
 - 3: **for** $i = 1, 2, 3, 4$ **do**
 - 4: $t_i = 0$; Initialize the model competence from i^{th} perspective, $c_i(0)$, by Eq. (2); Uniformly sample a data batch, $B_i(0)$, from the top $c_i(0)$ portions of D_i^{train} ;
 - 5: **Compute the perplexity (PPL) on $B_i(0)$, $\text{PPL}(B_i(0))$;**
 - 6: **end for**
 - 7: **repeat**
 - 8: $j = \arg \max_i (\text{PPL}(B_i(t_i)))$;
 - 9: Train the model with the $B_j(t_j)$;
 - 10: $t_j \leftarrow t_j + 1$;
 - 11: Estimate the model competence from j^{th} perspective, $c_j(t_j)$, by Eq. (2); Uniformly sample a data batch, $B_j(t_j)$, from the top $c_j(t_j)$ portions of D_j^{train} ;
 - 12: **Compute the perplexity (PPL) of model on $B_j(t_j)$, $\text{PPL}(B_j(t_j))$;**
 - 13: **until** Model converge.
-

riculum schedule and is usually set to 2, and T is the duration of curriculum learning and determines the length of the curriculum. In implementations, at training time step t , the top $c(t)$ portions of the sorted training dataset are selected to sample a training batch to train the model. In this way, the model is able to gradually proceed from easy samples to more complex ones in training, resulting in first starting to utilize the simple and easy-written reports for training, and then attempting to utilize harder reports for training.

5.2 Multiple Difficulty-based Curriculum Learning

The training instances of medical report generation task are pairs of medical images and corresponding reports which is a multi-modal data. It's hard to estimate the difficulty with only one metric. In addition, the experimental results (see Table 4) show that directly fusing multiple difficulty metrics as one ($d_1 + d_2 + d_3 + d_4$) is obviously inappropriate, which is also verified in Platanios et al. (2019). To this end, we extend the single difficulty-based curriculum learning into the multiple difficulty-based curriculum learning, where we provide the medical report generation models with four different difficulty metrics, i.e., d_1, d_2, d_3, d_4 (see Section 4).

A simple and natural way is to randomly or sequentially choose a curricula to train the model,

i.e., $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 1$. However, a better approach is to adaptively select the most appropriate curricula for each training step, which follows the common practice of human learning behavior: When we have learned some curricula well, we tend to choose the under-learned curricula to learn. Algorithm 2 summarizes the overall learning process of the proposed framework and Figure 3 illustrates the process of Algorithm 2. In implementations, similarly, we first sort the training dataset based on the four difficulty metrics and acquire four sorted training datasets in line 1-2. Then, based on the model competence, we acquire the training samples for each curricula, in line 4. In line 5, we further estimate the perplexity (PPL) of model on different training samples $B_i(t_i)$ corresponding to different curricula, defined as:

$$\text{PPL}(B_i(t_i)) = \sum_{R^k \in B_i(t_i)} \sqrt[N]{\prod_{m=1}^N \frac{1}{P(w_m^k | w_1^k, \dots, w_{m-1}^k)}}$$

where $R^k = \{w_1^k, w_2^k, \dots, w_N^k\}$ denotes the k -th report in $B_i(t_i)$. The perplexity (PPL) measures how many bits on average would be needed to encode each word of the report given the model, so the current curricula with higher PPL means that the model is not well-learned for this curricula and need to be improved. Therefore, the PPL can be used to determine the curricula at each training step dynamically. Specifically, in line 8-9, we select the under-learned curricula, i.e., the curricula with maximum PPL, to train the current model. After that, we again estimate the model competence in the selected curricula in line 11 and compute the PPL of model on the training samples corresponding to the selected curricula in line 12.

6 Experiment

We firstly describe two public datasets as well as the widely-used metrics, baselines and settings. Then we present the evaluation of our CMCL.

6.1 Datasets

We conduct experiments on two public datasets, i.e., a widely-used benchmark IU-Xray (Demner-Fushman et al., 2016) and a recently released large-scale MIMIC-CXR (Johnson et al., 2019).

- **IU-Xray**² is collected by Indiana University and is widely-used to evaluate the performance of medical report generation methods.

²<https://openi.nlm.nih.gov/>

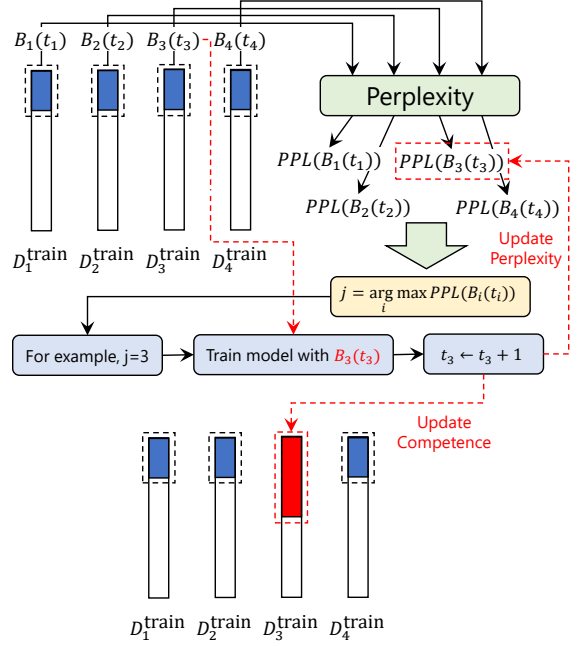


Figure 3: Illustration of Algorithm 2.

It contains 7,470 chest X-ray images associated with 3,955 radiology reports sourced from Indiana Network for Patient Care.

- **MIMIC-CXR**³ is the recently released largest dataset to date and consists of 377,110 chest X-ray images and 227,835 radiology reports from 64,588 patients of the Beth Israel Deaconess Medical Center.

For IU-Xray dataset, following previous works (Chen et al., 2020; Jing et al., 2019; Li et al., 2019, 2018), we randomly split the dataset into 70%-10%-20% training-validation-testing splits. At last, we preprocess the reports by tokenizing, converting to lower-cases and removing non-alpha tokens. For MIMIC-CXR, following Chen et al. (2020); Liu et al. (2021a,b), we use the official splits to report our results, resulting in 368,960 samples in the training set, 2,991 samples in the validation set and 5,159 samples in the test set. We convert all tokens of reports to lower-cases and filter tokens that occur less than 10 times in the corpus, resulting in a vocabulary of around 4,000 tokens.

6.2 Baselines

We tested three representative baselines that were originally designed for image captioning and three

³<https://physionet.org/content/mimic-cxr/2.0.0/>

competitive baselines that were originally designed for medical report generation.

6.2.1 Image Captioning Baselines

- **NIC:** Vinyals et al. (2015) proposed the encoder-decoder network, which employs a CNN-based encoder to extract image features and a RNN-based decoder to generate the target sentence, for image captioning.
- **Spatial-Attention:** Lu et al. (2017) proposed the visual attention, which is calculated on the hidden states, to help the model to focus on the most relevant image regions instead of the whole image.
- **Adaptive-Attention:** Considering that the decoder tends to require little or no visual information from the image to predict the non-visual words such as “the” and “of”, Lu et al. (2017) designed an adaptive attention model to decide when to employ the visual attention.

6.2.2 Medical Report Generation Baselines

- **CNN-HLSTM:** Jing et al. (2018) introduced the Hierarchical LSTM structure (HLSTM), which contains the paragraph LSTM and the sentence LSTM. HLSTM first uses the paragraph LSTM to generate a series of high-level topic vectors representing the sentences, and then utilizes the sentence LSTM to generate a sentence based on each topic vector.
- **HLSTM+att+Dual:** Harzig et al. (2019) proposed a hierarchical LSTM with the attention mechanism and further introduced two LSTMs, i.e., Normal LSTM and Abnormal LSTM, to help the model to generate more accurate normal and abnormal sentences.
- **Co-Attention:** Jing et al. (2018) proposed the co-attention model, which combines the merits of visual attention and semantic attention, to attend to both images and predicted semantic tags⁴ simultaneously, exploring the synergistic effects of visual and semantic information.

6.3 Metrics and Settings

We adopt the widely-used BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005) and ROUGE-L (Lin, 2004), which are reported by the

evaluation toolkit (Chen et al., 2015)⁵, to test the performance. Specifically, ROUGE-L is proposed for automatic evaluation of the extracted text summarization. METEOR and BLEU are originally designed for machine translation evaluation.

For all baselines, since our focus is to change the training paradigm, which improves existing baselines by efficiently utilizing the limited medical data, we keep the inner structure of the baselines untouched and preserve the original parameter setting. For our curriculum learning framework, following previous work (Platanios et al., 2019), the $c(0)$ and p are set to 0.01 and 2, respectively. For different baselines, we first re-implement the baselines without using any curriculum. When equipping baselines with curriculum, following Platanios et al. (2019), we set T in Eq.(2) to a quarter of the number of training steps that the baseline model takes to reach approximately 90% of its final BLEU-4 score. To boost the performance, we further incorporate the Batching method (Xu et al., 2020), which batches the samples with similar difficulty in the curriculum learning framework. To re-implement the baselines and our approach, following common practice (Jing et al., 2019; Li et al., 2019, 2018; Liu et al., 2021a,b), we extract image features for both dataset used for evaluation from a ResNet-50 (He et al., 2016), which is pretrained on ImageNet (Deng et al., 2009) and fine-tuned on public available CheXpert dataset (Irvin et al., 2019). To ensure consistency with the experiment settings of previous works (Chen et al., 2020), for IU-Xray, we utilize paired images of a patient as the input; for MIMIC-CXR, we use single image as the input. For parameter optimization, we use Adam optimizer (Kingma and Ba, 2014) with a batch size of 16 and a learning rate of 1e-4.

6.4 Automatic Evaluation

As shown in Table 1, for two datasets, all baselines equipped with our approach receive performance gains over most metrics. The results prove the effectiveness and the compatibility of our CMCL in promoting the performance of existing models by better utilizing the limited medical data. Besides, in Table 2, we further select six existing state-of-the-art models, i.e., HRGR-Agent (Li et al., 2018), CMAS-RL (Jing et al., 2019), SentSAT + KG (Zhang et al., 2020a), Up-Down (Anderson et al., 2018), Transformer (Chen et al., 2020) and

⁴<https://ii.nlm.nih.gov/MTI/>

⁵<https://github.com/tylin/coco-caption>

Methods	Dataset: MIMIC-CXR (Johnson et al., 2019)						Dataset: IU-Xray (Demner-Fushman et al., 2016)					
	B-1	B-2	B-3	B-4	M	R-L	B-1	B-2	B-3	B-4	M	R-L
NIC (Vinyals et al., 2015) [†] w/ CMCL	0.290 0.301	0.182 0.189	0.119 0.123	0.081 0.085	0.112 0.119	0.249 0.241	0.352 0.358	0.227 0.223	0.154 0.160	0.109 0.114	0.133 0.137	0.313 0.317
Spatial-Attention (Lu et al., 2017) [†] w/ CMCL	0.302 0.312	0.189 0.200	0.122 0.125	0.082 0.087	0.118 0.118	0.259 0.258	0.374 0.381	0.235 0.246	0.158 0.164	0.120 0.123	0.146 0.153	0.322 0.327
Adaptive-Attention (Lu et al., 2017) [†] w/ CMCL	0.307 0.302	0.192 0.192	0.124 0.129	0.084 0.091	0.119 0.125	0.262 0.264	0.433 0.437	0.285 0.281	0.194 0.196	0.137 0.140	0.166 0.174	0.349 0.338
CNN-HLSTM (Krause et al., 2017) [†] w/ CMCL	0.321 0.337	0.203 0.210	0.129 0.136	0.092 0.097	0.125 0.131	0.270 0.274	0.435 0.462	0.280 0.293	0.187 0.207	0.131 0.155	0.173 0.179	0.346 0.360
HLSTM+att+Dual (Harzig et al., 2019) [†] w/ CMCL	0.328 0.330	0.204 0.206	0.127 0.133	0.090 0.088	0.122 0.119	0.267 0.272	0.447 0.461	0.289 0.298	0.192 0.201	0.144 0.150	0.175 0.173	0.358 0.359
Co-Attention (Jing et al., 2018) [†] w/ CMCL	0.329 0.344	0.206 0.217	0.133 0.140	0.095 0.097	0.129 0.133	0.273 0.281	0.463 0.473	0.293 0.305	0.207 0.217	0.155 0.162	0.178 0.186	0.365 0.378

Table 1: Performance of automatic evaluations on the test sets of the MIMIC-CXR and the IU-Xray datasets. CMCL denotes the Competence-based Multimodal Curriculum Learning framework. B-n, M and R-L are short for BLEU-n, METEOR and ROUGE-L, respectively. Higher is better in all columns. [†] denotes our re-implementation. As we can see, all baseline models enjoy comfortable improvements in most metrics with our CMCL.

Methods	Dataset: MIMIC-CXR (Johnson et al., 2019)						Dataset: IU-Xray (Demner-Fushman et al., 2016)					
	B-1	B-2	B-3	B-4	M	R-L	B-1	B-2	B-3	B-4	M	R-L
HRGR-Agent (Li et al., 2018)	-	-	-	-	-	-	0.438	0.298	0.208	0.151	-	0.322
CMAS-RL (Jing et al., 2019)	-	-	-	-	-	-	0.464	0.301	0.210	0.154	-	0.362
SentSAT + KG (Zhang et al., 2020a)	-	-	-	-	-	-	0.441	0.291	0.203	0.147	-	0.367
Up-Down (Anderson et al., 2018)	0.317	0.195	0.130	0.092	0.128	0.267	-	-	-	-	-	-
Transformer (Chen et al., 2020)	0.314	0.192	0.127	0.090	0.125	0.265	0.396	0.254	0.179	0.135	0.164	0.342
R2Gen (Chen et al., 2020)	0.353	0.218	0.145	0.103	0.142	0.277	0.470	0.304	0.219	0.165	0.187	0.371
CMCL (Ours)	0.344	0.217	0.140	0.097	0.133	0.281	0.473	0.305	0.217	0.162	0.186	0.378

Table 2: Comparison with existing state-of-the-art methods on the test set of the MIMIC-CXR dataset and the IU-Xray dataset. CMCL is taken from the “Co-Attention w/ CMCL” in Table 1. In this table, the Red and Blue colored numbers denote the best and second best results across all approaches, respectively.

vs. Models	Baseline wins	Tie	‘w/ CMCL’ wins
CNN-HLSTM (Jing et al., 2018) [†]	15	28	57
Co-Attention (Jing et al., 2018) [†]	24	35	41

Table 3: We invite 2 professional clinicians to conduct the human evaluation for comparing our method with baselines. All values are reported in percentage (%).

R2Gen (Chen et al., 2020), for comparison. For these selected models, we directly quote the results from the original paper for IU-Xray, and from Chen et al. (2020) for MIMIC-CXR. As we can see, based on the Co-Attention (Chen et al., 2020), our approach CMCL achieves results competitive with these state-of-the-art models on major metrics, which further demonstrate the effectiveness of the proposed approach.

6.5 Human Evaluation

In this section, to verify the effectiveness of our approach in clinical practice, we invite two professional clinicians to evaluate the perceptual quality of 100 randomly selected reports generated by “Baselines” and “Baselines w/ CMCL”. For the baselines, we choose a representative model: CNN-HLSTM and a state-of-the-art model: Co-Attention. The clinicians are unaware of which

model generates these reports. In particular, to have more documents examined, we did not use the same documents for both clinicians and check the agreements between them. That is to say, the documents for different clinicians do not overlap. The results in Table 3 show that our approach is better than baselines in clinical practice with winning pick-up percentages. In particular, all invited professional clinicians found that our approach can generate fluent reports with more accurate descriptions of abnormalities than baselines. It indicates that our approach can help baselines to efficiently alleviate the data bias problem, which also can be verified in Section 6.7.

6.6 Quantitative Analysis

Analysis on the Difficulty Metrics In this section, we conduct an ablation study by only using a single difficulty metric during the curriculum learning, i.e., single difficulty-based curriculum learning, to investigate the contribution of each difficulty metric in our framework and the results are shown in Table 4. Settings (a-d) show that every difficulty metric can boost the performance of baselines, which verify the effectiveness of our designed difficulty metrics. In particular, 1) the

Settings	Visual Difficulty		Textual Difficulty		Route Strategy	Dataset: IU-Xray (Demner-Fushman et al., 2016)					
	Heuristic Metric	Model Confidence	Heuristic Metric	Model Confidence		Baseline: CNN-HLSTM (Jing et al., 2018)					
						B-1	B-2	B-3	B-4	M	R-L
Baseline	-	-	-	-	-	0.435	0.280	0.187	0.131	0.173	0.346
(a)	✓	-	-	-	-	0.438	0.283	0.188	0.132	0.173	0.348
(b)	-	✓	-	-	-	0.447	0.288	0.195	0.143	0.175	0.354
(c)	-	-	✓	-	-	0.443	0.287	0.192	0.135	0.175	0.351
(d)	-	-	-	✓	-	0.454	0.290	0.201	0.148	0.177	0.357
(e)	✓	✓	-	-	Dynamically	0.450	0.289	0.196	0.144	0.176	0.355
(f)	✓	✓	✓	-	Dynamically	0.455	0.290	0.199	0.145	0.176	0.357
(g)	✓	✓	✓	✓	Dynamically	0.462	0.293	0.207	0.155	0.179	0.360
(h)	✓	✓	✓	✓	Fuse	0.440	0.282	0.190	0.134	0.174	0.349
(i)	✓	✓	✓	✓	Randomly	0.457	0.291	0.199	0.146	0.178	0.358
(j)	✓	✓	✓	✓	Sequentially	0.459	0.290	0.203	0.150	0.176	0.354

Table 4: Quantitative analysis of our approach, which includes four designed difficulty metrics (see Section 4) and the route strategy (see Section 5.2). We conduct the analysis on the widely-used baseline model CNN-HLSTM (Jing et al., 2018). The setting (g) also denotes our full proposed approach.

model confidence in both visual and textual difficulties achieves better performance than the heuristic metrics. It shows that the model confidence is the more critical in neural models. 2) Both the model confidence and heuristic metrics in the textual difficulty achieve better performance than their counterparts in the visual difficulty, which indicates that the textual data bias is the more critical in textual report generation task. When progressively incorporate each difficulty metric, the performance will increase continuously (see settings (e-g)), showing that integrating different difficulty metrics can bring the improvements from different aspects, and the advantages of all difficulty metrics can be united as an overall improvement.

Analysis on the Route Strategy As stated in Section 5.2, to implement the multiple difficulty-based curriculum learning, three simple and natural ways is to: 1) **Fuse** multiple difficulty metrics directly as a single mixed difficulty metric, $d_1 + d_2 + d_3 + d_4$; 2) **Randomly** choose a curricula and 3) **Sequentially** choose a curricula (i.e., $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 1$) to train the model. Table 4 (h-j) show the results of the three implementations. As we can see, all route strategies are viable in practice with improved performance of medical report generation, which proves the effectiveness and robustness of our CMCL framework. Besides, all of them perform worse than our approach (Setting (g)), which confirms the effectiveness of dynamically learning strategy at each training step.

6.7 Qualitative Analysis

In Figure 1, we give two intuitive examples to better understand our approach. As we can see, our approach generates structured and robust reports, which show significant alignment with ground truth

reports and are supported by accurate abnormal descriptions. For example, the generated report correctly describes “*Blunting of right costophrenic*” in the first example and “*Scoliosis is present*” in the second example. The results prove our arguments and verify the effectiveness of our proposed CMCL in alleviating the data bias problem by enabling the model to gradually proceed from easy to more complex instances in training.

7 Conclusion

In this paper, we propose the novel competence-based multimodal curriculum learning framework (CMCL) to alleviate the data bias by efficiently utilizing the limited medical data for medical report generation. To this end, considering the difficulty of accurately capturing and describing the abnormalities, we first assess four sample difficulties of training data from the visual complexity and the textual complexity, resulting in four different curricula. Next, CMCL enables the model to be trained with the appropriate curricula and gradually proceed from easy samples to more complex ones in training. Experimental results demonstrate the effectiveness and the generalization capabilities of CMCL, which consistently boosts the performance of the baselines under most metrics.

Acknowledgments

This work is partly supported by Tencent Medical AI Lab, Beijing, China. We would like to sincerely thank the clinicians Xiaoxia Xie, Jing Zhang and Minghui Shao of the Harbin Chest Hospital in China for providing the human evaluation. We sincerely thank all the anonymous reviewers for their constructive comments and suggestions. Xian Wu is the corresponding author of this paper.

Ethical Considerations

In this work, we focus on helping a wide range of existing medical report generation systems alleviate the data bias by efficiently utilizing the limited medical data for medical report generation. Our work can enable the existing systems to gradually proceed from easy samples to more complex ones in training, which is similar to the learning curve of radiologist: (1) first start from simple and easy-written reports; (2) and then attempt to consume harder reports, which consist of rare and diverse abnormalities. As a result, our work can promote the usefulness of existing medical report generation systems in better **assisting** radiologists in clinical decision-makings and reducing their workload. In particular, for radiologists, given a large amount of medical images, the systems can automatically generate medical reports, the radiologists only need to make revisions rather than write a new report from scratch. We conduct the experiments on the public MIMIC-CXR and IU-Xray datasets. All protected health information was de-identified. De-identification was performed in compliance with Health Insurance Portability and Accountability Act (HIPAA) standards in order to facilitate public access to the datasets. Deletion of protected health information (PHI) from structured data sources (e.g., database fields that provide patient name or date of birth) was straightforward. All necessary patient/participant consent has been obtained and the appropriate institutional forms have been archived.

References

- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and VQA. In *CVPR*.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: an automatic metric for MT evaluation with improved correlation with human judgments. In *IEEE-valuation@ACL*.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *ICML*.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C. Lawrence Zitnick. 2015. Microsoft COCO captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*.
- Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang Wan. 2020. Generating radiology reports via memory-driven transformer. In *EMNLP*.
- Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. 2020. Meshed-memory transformer for image captioning. In *CVPR*.
- Louke Delrue, Robert Gosselin, Bart Ilsen, An Van Landeghem, Johan de Mey, and Philippe Duyck. 2011. Difficulties in the interpretation of chest radiography. *Comparative Interpretation of CT and Standard Radiography of the Chest*, pages 27–49.
- Dina Demner-Fushman, Marc D. Kohli, Marc B. Rosenman, Sonya E. Shooshan, Laritza Rodriguez, Sameer K. Antani, George R. Thoma, and Clement J. McDonald. 2016. Preparing a collection of radiology examinations for distribution and retrieval. *J. Am. Medical Informatics Assoc.*, 23(2):304–310.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*.
- Jeffrey L Elman. 1993. Learning and development in neural networks: The importance of starting small. *Cognition*, 48(1):71–99.
- Stacy K Goergen, Felicity J Pool, Tari J Turner, Jane E Grimm, Mark N Appleyard, Carmel Crock, Michael C Fahey, Michael F Fay, Nicholas J Ferris, Susan M Liew, et al. 2013. Evidence-based guideline for the written radiology report: Methods, recommendations and implementation challenges. *Journal of medical imaging and radiation oncology*, 57(1):1–7.
- Philipp Harzig, Yan-Ying Chen, Francine Chen, and Rainer Lienhart. 2019. Addressing data bias problems for chest x-ray image report generation. In *BMVC*.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Comput.*, 9(8):1735–1780.
- Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghighi, Robyn L. Ball, Katie S. Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. 2019. Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *AAAI*.
- Baoyu Jing, Zeya Wang, and Eric P. Xing. 2019. Show, describe and conclude: On exploiting the structure information of chest x-ray reports. In *ACL*.

- Baoyu Jing, Pengtao Xie, and Eric P. Xing. 2018. On the automatic generation of medical imaging reports. In *ACL*.
- Alistair E. W. Johnson, Tom J. Pollard, Seth J. Berkowitz, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih-ying Deng, Roger G. Mark, and Steven Horng. 2019. MIMIC-CXR: A large publicly available database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042*.
- Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. In *ICLR*.
- Tom Kocmi and Ondrej Bojar. 2017. Curriculum learning and minibatch bucketing in neural machine translation. In *RANLP*.
- Jonathan Krause, Justin Johnson, Ranjay Krishna, and Li Fei-Fei. 2017. A hierarchical approach for generating descriptive image paragraphs. In *CVPR*.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *NIPS*.
- Gaurav Kumar, George F. Foster, Colin Cherry, and Maxim Krikun. 2019. Reinforcement learning based curriculum optimization for neural machine translation. In *NAACL-HLT*.
- Christy Y. Li, Xiaodan Liang, Zhiting Hu, and Eric P. Xing. 2019. Knowledge-driven encode, retrieve, paraphrase for medical image report generation. In *AAAI*.
- Qing Li, Siyuan Huang, Yining Hong, and Song-Chun Zhu. 2020. A competence-aware curriculum for visual concepts learning via question answering. In *ECCV*.
- Yuan Li, Xiaodan Liang, Zhiting Hu, and Eric P. Xing. 2018. Hybrid retrieval-generation reinforced agent for medical image report generation. In *NeurIPS*.
- Xiaodan Liang, Zhiting Hu, Hao Zhang, Chuang Gan, and Eric P. Xing. 2017. Recurrent topic-transition GAN for visual paragraph generation. In *ICCV*.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *ACL*.
- Fenglin Liu, Yuanxin Liu, Xuancheng Ren, Xiaodong He, and Xu Sun. 2019a. Aligning visual regions and textual concepts for semantic-grounded image representations. In *NeurIPS*.
- Fenglin Liu, Xuancheng Ren, Yuanxin Liu, Kai Lei, and Xu Sun. 2019b. Exploring and distilling cross-modal information for image captioning. In *IJCAI*.
- Fenglin Liu, Xuancheng Ren, Yuanxin Liu, Houfeng Wang, and Xu Sun. 2018. simnet: Stepwise image-topic merging network for generating detailed and comprehensive image captions. In *EMNLP*.
- Fenglin Liu, Xuancheng Ren, Xian Wu, Shen Ge, Wei Fan, Yuexian Zou, and Xu Sun. 2020a. Prophet attention: Predicting attention with future attention. In *NeurIPS*.
- Fenglin Liu, Xian Wu, Shen Ge, Wei Fan, and Yuexian Zou. 2021a. Exploring and distilling posterior and prior knowledge for radiology report generation. In *CVPR*.
- Fenglin Liu, Changchang Yin, Xian Wu, Shen Ge, Ping Zhang, and Xu Sun. 2021b. Contrastive attention for automatic chest x-ray report generation. In *ACL (Findings)*.
- Guanxiong Liu, Tzu-Ming Harry Hsu, Matthew B. A. McDermott, Willie Boag, Wei-Hung Weng, Peter Szolovits, and Marzyeh Ghassemi. 2019c. Clinically accurate chest x-ray report generation. In *MLHC*.
- Xuebo Liu, Houtim Lai, Derek F. Wong, and Lidia S. Chao. 2020b. Norm-based curriculum learning for neural machine translation. In *ACL*.
- Justin R. Lovelace and Bobak Mortazavi. 2020. Learning to generate clinically coherent chest x-ray reports. In *EMNLP (Findings)*.
- Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. 2017. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *CVPR*.
- Yasuhide Miura, Yuhao Zhang, Emily Bao Tsai, Curtis P. Langlotz, and Dan Jurafsky. 2021. Improving factual completeness and consistency of image-to-text radiology report generation. In *NAACL-HLT*.
- Yingwei Pan, Ting Yao, Yehao Li, and Tao Mei. 2020. X-linear attention networks for image captioning. In *CVPR*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a Method for automatic evaluation of machine translation. In *ACL*.
- Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neubig, Barnabás Póczos, and Tom M. Mitchell. 2019. Competence-based curriculum learning for neural machine translation. In *NAACL-HLT*.
- Steven J. Rennie, Etienne Marcheret, Youssef Mroueh, Jarret Ross, and Vaibhava Goel. 2017. Self-critical sequence training for image captioning. In *CVPR*.
- 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 *CVPR*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.

- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. In *CVPR*.
- Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, and Ronald M. Summers. 2018. Tienet: Text-image embedding network for common thorax disease classification and reporting in chest x-rays. In *CVPR*.
- Yiru Wang, Weihao Gan, Jie Yang, Wei Wu, and Junjie Yan. 2019. Dynamic curriculum learning for imbalanced data classification. In *ICCV*.
- Daphna Weinshall, Gad Cohen, and Dan Amir. 2018. Curriculum learning by transfer learning: Theory and experiments with deep networks. In *ICML*.
- Chen Xu, Bojie Hu, Yufan Jiang, Kai Feng, Zeyang Wang, Shen Huang, Qi Ju, Tong Xiao, and Jingbo Zhu. 2020. Dynamic curriculum learning for low-resource neural machine translation. In *COLING*.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron C. Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*.
- Yuan Xue, Tao Xu, L. Rodney Long, Zhiyun Xue, Sameer K. Antani, George R. Thoma, and Xiaolei Huang. 2018. Multimodal recurrent model with attention for automated radiology report generation. In *MICCAI*.
- Haonan Yu, Jiang Wang, Zhiheng Huang, Yi Yang, and Wei Xu. 2016. Video paragraph captioning using hierarchical recurrent neural networks. In *CVPR*.
- Jianbo Yuan, Haofu Liao, Rui Luo, and Jiebo Luo. 2019. Automatic radiology report generation based on multi-view image fusion and medical concept enrichment. In *MICCAI*.
- Xuan Zhang, Gaurav Kumar, Huda Khayrallah, Kenton Murray, Jeremy Gwinnup, Marianna J. Martindale, Paul McNamee, Kevin Duh, and Marine Carpuat. 2018. An empirical exploration of curriculum learning for neural machine translation. *arXiv preprint arXiv:1811.00739*.
- Yixiao Zhang, Xiaosong Wang, Ziyue Xu, Qihang Yu, Alan L. Yuille, and Daguang Xu. 2020a. When radiology report generation meets knowledge graph. In *AAAI*.
- Yuhao Zhang, Derek Merck, Emily Bao Tsai, Christopher D. Manning, and Curtis Langlotz. 2020b. Optimizing the factual correctness of a summary: A study of summarizing radiology reports. In *ACL*.
- Mingjun Zhao, Haijiang Wu, Di Niu, and Xiaoli Wang. 2020. Reinforced curriculum learning on pre-trained neural machine translation models. In *AAAI*.