



ISSN: 2454-9940



**INTERNATIONAL JOURNAL OF APPLIED
SCIENCE ENGINEERING AND MANAGEMENT**

E-Mail :
editor.ijasem@gmail.com
editor@ijasem.org

www.ijasem.org

<https://zenodo.org/records/14505535>

AN ANNOTATED CRNN AND GRU-BASED AUDIO CAPTIONING SYSTEM

Ballepu Kumar Naveen¹, Assistant Professor¹, Department of ECE, Siddhartha Institute of Technology & Sciences, Telangana, India

A. Santhosh Reddy², Assistant Professor², Department of CSE, Siddhartha Institute of Technology & Sciences, Telangana, India.

ABSTRACT

Audio captioning aims at generating a natural sentence to describe the content in an audio clip. This paper proposes the use of a powerful CRNN encoder combined with a GRU decoder to tackle this multi-modal task. In addition to standard cross-entropy, reinforcement learning is also investigated for generating richer and more accurate captions. Our approach significantly improves against the baseline model on all shown metrics achieving a relative improvement of at least 34%. Results indicate that our proposed CRNNGRU model with reinforcement learning achieves a Spider of 0.190 on the Clotho evaluation set1. With data augmentation, the performance is further boosted to 0.223. In the DCASE challenge Task 6 we ranked fourth based on Spider, second on 5 metrics including BLEU, ROUGE-L and METEOR, without ensemble or data augmentation while maintaining a small model size (only 5 million parameters). Index Terms— audio captioning, reinforcement learning, convolutional recurrent neural networks

INTRODUCTION

Automatic captioning is a challenging task that involves joint learning of different modalities. For example, image captioning requires extracting features from an image and combining them with a language model to generate reasonable sentences to describe the image. Similarly, video captioning learns features from a temporal sequence of images as well as audio to generate captions. However, audio captioning does not attract much attention [1], unlike in the image and video fields. By its nature, captioning is a novel multi-modal task that captures the fine details within an auditory scene with natural language (text). Unlike other tasks such as sound or acoustic event detection, which only focuses on narrow single-label estimation of an event, audio captioning is concerned with producing rich sentences appropriately and precisely describing an audio. Audio captioning has great potential in real-world applications, such as audio surveillance, automatic content description and content-oriented machine-to-machine interaction. Initial work in audio captioning has been done in [1], which utilized the commercial Propounds Effects [2] audio corpus as a proof of concept. The paper utilized an encoder-decoder Architecture containing a three-layer bidirectional gated recurrent unit (Bigram) encoder and a two-layer Bigram decoder. An attention pooling is added to summarize the encoder sentence. Subsequent work in [3] investigated audio

captioning within the scope of Chinese captioning, firstly proposing a public captioning corpus, focusing on dialogues within a hospital setting. Their results showed that within a limited domain, audio captions can indeed be generated by a single layer encoder-decoder GRU network successfully, but also questioned if commonly utilized metrics for machine translation can well evaluate the final performance.

The main discussion is that even though their approach achieves measurably near-human performance via objective metrics, the generated sentences are often less useful according to human evaluation. Similar to other text generation tasks like machine translation and image captioning, exposure bias also exists in audio captioning. Neural network-based models are typically trained in “teacher forcing” fashion, meaning they aim to maximize the likelihood of a future ground-truth word given the current ground-truth word. However, ground-truth annotations are only available during training, while during inference, the model can solely rely on its own predicted current word to infer the next word. This leads to an error accumulation during test-time. Another problem in text generation tasks is a mismatch between the training objective and evaluation metric. Generative models are typically evaluated by discrete metrics such as BLEU [4], ROUGE-L [5], Cider [6] or METEOR [7]. However, these non-differentiable metrics cannot be directly optimized using the standard back-propagation approach. Previous studies have shown that the application of Reinforcement Learning (RL) can partially circumvent exposure bias while optimizing the discrete evaluation metrics at the same time. RL is first proposed to train natural language generation models in [8]. It takes a generative model as an agent and treats words and context as an external environment.

The model parameters define a policy, and the choice of the current generated word corresponds to its action. The reward comes from evaluation scores (BLEU, METEOR, Cider etc.) of the sampled sentence. Policy-gradient [9] is used to estimate the gradient of the agent parameters using the reward. Work in [10] improves this method by using rewards from greedy-sampled sentences as

<https://zenodo.org/records/14505535>

the baseline to reduce the high variance of rewards. Subsequent work in [11] also adopts actor-critic methods [12] to estimate the value of generated words instead of sampling from the action space. In this paper, we explore the use of the self-critical sequence training (SCST) approach (proposed in [10]) for audio captioning. This paper is structured as follows, in Section 2 we put forth our CRNN-based encoder-decoder approach to audio captioning. Then in Section 3, the experimental setup, including front-end features and model parameters, are shown. Our results and analysis are dising-

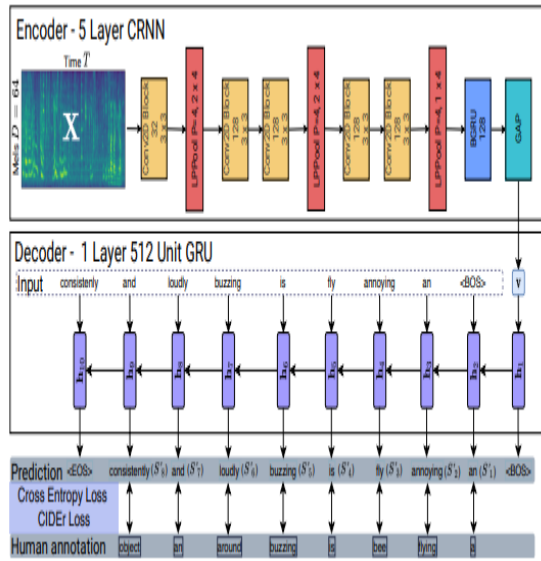


Figure 1: Our proposed encoder-decoder architecture. The encoder is a CRNN model which outputs a fixed-sized 256-dimensional embedding v after a global average pooling layer (GAP).

A convolution block refers to an initial batch normalization, then a convolution, and lastly, a Leaky (slope -0.1) activation. The numbers in each block represent the output channel size and the kernel size. For example, "32, 3×3 " means the convolution layer has 32 output channels with a kernel size of 3×3 . All convolutions use padding in order to preserve the input size. Then a GRU decoder utilizes this audio embedding v or embedding of the word $S_0 t$ at each time-step, to predict the next word $S_0 t+1$.

APPROACH

Similar to previous audio captioning frameworks [3], our approach follows a standard encoder-decoder model (see Equation (1)).

$$\begin{aligned} v &= \text{Enc}(X) \\ [S'_1, \dots, S'_T] &= \text{Dec}(v) \end{aligned} \quad (1)$$

The encoder (Enc) is fed an audio-spectrogram (X) and produces a fixed-sized vector representation v , which the decoder (Dec) uses to predict the caption sentence. Specifically, the decoder generates a single word-token $S_0 t$ for each time-step it up until an end of sentence (ϵ) token is seen (see Figure 1). In audio captioning, decoding differs between training and evaluation stages:

$$\ell_{\text{XE}}(\theta; S, v) = - \sum_{t=1}^T \log p(S_t | \theta; v) \quad (2)$$

During training, where transcriptions are available, Dec generates word-tokens given the embedding v and human-annotated data S , supervised by a cross-entropy (XE) loss (see Equation (2)). During evaluation and testing, no transcriptions are available; thus word-tokens are sampled from the decoder given the audio embedding v . From this description, it is evident that the quality of v directly affects the generated sentence quality. Thus, our approach mainly diverges from previous approaches in two ways: the encoder architecture and the loss function. Previous encoder models (GRU) might be insufficient to produce a robust vector representation, thus we replace the standard GRU encoder with a robust convolutional recurrent neural network (CRNN). Our framework can be seen in Figure 1. Moreover, standard XE training has its potential downsides. For one, the criterion only compares single word-tokens and neglects context information. Second, since each word is treated individually, syntactically incorrect sentences can be generated. Third, optimizing XE inevitably leads to monotonous sentences, because the model is required to precisely imitate a sentence word by word, instead of allowing semantically similar, but different worded sentences.

We employ reinforcement learning for audio captioning. Reinforcement learning allows us to directly back-propagate a metric (e.g., BLEU or Cider) in the form of a reward. Formally we train the model to minimize the negative reward of a single sampled sentence S_0 :

$$\ell_{\text{RL}}(\theta; v) = -r(S'), S' \sim p(S' | \theta; v) \quad (3)$$

where $S_0 = [S_0 1, S_0 2, \dots, S_0 T]$. By incorporating the policy gradient method with baseline normalization, the parameter gradients can be estimated as follows:

<https://zenodo.org/records/14505535>

$$\nabla_{\theta} \ell(\theta; \mathbf{v}) = -(\mathbf{r}(S') - b) \nabla_{\theta} \log p(S' | \theta; \mathbf{v}), S' \sim p(S' | \theta; \mathbf{v}) \quad (4)$$

here b is a pre-defined baseline normalization constant to reduce the high variance brought by sampling [12]. We set b as the greedy decoding reward because of its effectiveness in image captioning [10].

Models

Encoder Our proposed encoder is a CRNN model, which has seen success in localizing sound events [13, 14]. The architecture consists of a five-layer CNN (utilizing 3×3 convolutions), summarized into three blocks, with L4-Norm pooling after each block. The CNN blocks subsample the temporal dimension by factor of 4. A Bigram is attached after the last CNN output, enhancing our model's ability to localize sounds accurately. At last, we use a global average pooling (GAP) layer in order to remove any time-variability to a single, time-independent representation $\mathbf{v} \in \mathbb{R}^{256}$. The encoder has 679k parameters, making it comparably lightweight while only using 2.7 MB on disk. **Decoder** In the context of audio captioning, a decoder takes a fixed-sized embedding and aims to produce a sentence. We use a single-layer GRU with 512 hidden units as our decoder model.

EXPERIMENTS

Dataset

The challenge provides Clotho [2, 15] for the audio captioning task. It contains a total of 4981 audio samples, where the duration is uniformly distributed between 15 to 30 seconds. All audio samples are collected from the Freesound platform. Five native English speakers annotate each sample; thus, 24905 captions are available in total. Captions are post-processed to ensure each caption has eight to 20 words, and the caption does not contain unique words, named entities or speech transcription. The dataset is officially split into three sets, termed as development, evaluation, and testing, with a ratio of 60%-20%-20%. In the challenge, the development and evaluation sets are used for training our audio captioning model while the testing set is for evaluating the model.

Data pre-processing

We extract 64-dimensional log-Mel spectrogram (LMS) as our default input feature. Here a single frame is extracted via a 2048-point Fourier transform every 20 ms with a Hann window size of 40 ms. This results in a $\mathbf{X} \in \mathbb{R}^{T \times D}$ log-Mel spectrogram feature for each input audio, where D

$= 64$ and T is the number of frames. Moreover, the input feature is normalized by the mean and standard deviation of the development set. For each caption in the dataset, we remove punctuation and convert all letters to lowercase to reduce the vocabulary size. To mark the beginning and the end of sentences, we add special tokens “” and “” to captions. The available training data is split into a model training part, consisting of 90% of available data and a held-out 10% validation set.

Evaluation metrics

A total of eight objective metrics is utilized to evaluate our model-generated captions: BLEU@1-4 grams [4], METEOR [7], RougeL [5], Cider [6] and SPICE [16]. A further Spider metric is calculated as the mean of Cider and SPICE.

Training details

We submit predictions from four models to the challenge:

- CRNN-B (Base). This is our baseline CRNN-GRU encoder-decoder model.
- CRNN-W (Word). Here, the decoder word-embeddings are initialized from Word2Vec word-embeddings trained on the development set captions.
- CRNN-E (Ensemble). Here we fuse CRNN-B and CRNN-W results on output level.
- CRNN-R (Reinforcement). Here we finetune CRNN-W using reinforcement learning. The details for each submission are elaborated in the following. **XE training** For XE training, teacher forcing is used to accelerate the training process. We evaluate the model on the validation set at each epoch and select the best model according to the highest BLEU4 score. We train the model for 20 epochs and use Adam [17] optimizer with an initial learning rate of 5×10^{-4} . The batch size is 32. According to whether Word2Vec is used for word embedding initialization, we get CRNN-B and CRNN-W respectively. **Ensemble** In order to further enhance performance we merge the outputs of CRNN-B and CRNN-W on word-level. The encoded audio representation \mathbf{v} is fed to both CRNN-B and CRNN-W to obtain two-word probabilities p_1 and p_2 .

We ensemble the two models, which means the current word is decoded according to the mean of p_1 and p_2 . Then the current word embedding is fed to CRNN-B and CRNN-W to obtain the next word until is generated. **Reinforcement** The CRNN-R approach is first initialized by training a CRNN-W model using the standard XE criterion. This model

<https://zenodo.org/records/14505535>

is then finetuned using reinforcement learning, as seen in Section 2, by optimizing the Cider score using policy gradient with baseline normalization. Although [21] optimized Spider by policy gradient in image captioning, we choose Cider as the training objective because Cider optimized model trained by SCST achieved better performance [10]. Cider measures sentence similarity through representation by n-gram TF-IDFs while BLEU of sentences on "hard" n-gram overlaps. Such a "soft" similarity (Cider) may be a better optimization objective compared with BLEU under the condition that one audio corresponds to several semantic similar sentences, possibly composed of different n-grams. The model is trained for 25 epochs using Adam optimizer with a learning rate of 5×10^{-5} . Similar to the practice in XE training, we report the best model based on the Cider score on the validation set.

RESULTS

Results

Our results on the Clotho evaluation set are displayed in Table 1 and compared with the DCASE challenge baseline, which consists of a three-layer Bigram encoder and two-layer Bigram decoder. As it can be seen, our initial CRNN-B model largely outperforms the baseline, indicating that a potent encoder is indeed beneficial towards captioning performance. By initializing word embeddings with Word2Vec trained on the development set captions, CRNN-W gets a slight performance improvement in most metrics compared with CRNN-B, except Cider and METEOR. CRNN-E improves performance against both CRNN-B and CRNN-W. Our best performing model is CRNN-R. Interestingly, although CRNN-R is optimized towards Cider score, the relative improvement in BLEU3 and BLEU4 are more significant than Cider. The improvement in ROUGE1 and METEOR is not as significant as other metrics. However, CRNN-R does achieve the best performance in terms of all evaluation metrics, which validates the effectiveness of reinforcement learning for audio captioning with regards to the official challenge evaluation, our CRNNR achieves the fourth place in DCASE2020 task 6 on the Clotho testing set. However, there is only a slight difference between our submission and the submission ranking the third (0.194 / 0.196).

CONCLUSION

In this paper, we propose a novel audio captioning approach utilizing a CRNN encoder front-end as well as a reinforcement learning framework. Audio captioning models are trained on the Clotho dataset. The results on the Clotho evaluation set

suggest that the CRNN encoder is crucial to extract useful audio embeddings for captioning while reinforcement learning further improves the performance significantly in terms of all metrics. Our approach ranked fourth in the DCASE2020 task 6 challenge testing set with a competitive result on all metrics except Cider. Notably, our approach is the best performing non-ensemble result without data augmentation, with the least parameters (5 million). By further utilizing Specie data augmentation, we observe an additional boost in regards to the Spider score on the evaluation set from 0.190 to 0.223.

REFERENCES

- [1] K. Dross's, S. Advance, and T. Virtanen, "Automated audio captioning with recurrent neural networks," in *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, Oct 2017, pp. 374–378.
- [2] S. Lipping, K. Dross's, and T. Virtanen, "Crowdsourcing a dataset of audio captions," in *Proceedings of the Detection and Classification of Acoustic Scenes and Events Workshop (DCASE)*, Nov. 2019. [Online]. Available: <https://arxiv.org/abs/1907.09238>
- [3] M. Wu, H. Dinkel, and K. Yu, "Audio Caption: Listen and Tell," in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, vol. 2019– May. Institute of Electrical and Electronics Engineers Inc., May 2019, pp. 830–834.
- [4] K. Papini, S. Roukoops, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting on association for computational linguistics*. Association for Computational Linguistics, 2002, pp. 311–318.
- [5] C. Y. Lin, "Rouge: A package for automatic evaluation of summaries," *Proceedings of the workshop on text summarization branches out (WAS 2004)*, 2004.
- [6] R. Vedantam, C. L. Zitnick, and D. Parikh, "CIDER: Consensus-based image description evaluation," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2015.
- [7] S. Banerjee and A. Lavie, "METEOR: An automatic metric for MT evaluation with improved correlation with human judgments," in *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005.
- [8] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba, "Sequence level training with recurrent neural networks," *arXiv preprint arXiv:1511.06732*, 2015.
- [9] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *Advances in neural information processing systems*, 2000, pp. 1057–1063.
- [10] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, "Self-critical sequence training for image captioning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7008–7024.
- [11] D. Bahdanau, P. Brakel, K. Xu, A. Goyal, R. Lowe, J. Pineau, A. Courville, and Y. Bengio, "An actor-critic

<https://zenodo.org/records/14505535>

algorithm for sequence prediction,” arXiv preprint arXiv:1607.07086, 2016.

[12] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

[13] H. Dinkel and K. Yu, “Duration Robust Weakly Supervised Sound Event Detection,” in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, may 2020, pp. 311–315. [Online]. Available: <https://ieeexplore.ieee.org/document/9053459/>

[14] H. Dinkel, Y. Chen, M. Wu, and K. Yu, “GPVAD: Towards noise robust voice activity detection via weakly supervised sound event detection,” mar 2020. [Online]. Available: <http://arxiv.org/abs/2003.12222>

[15] K. Drossos, S. Lipping, and T. Virtanen, “Clotho: An audio captioning dataset,” in *45th IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Barcelona, Spain, May 2020. [Online]. Available: <https://arxiv.org/abs/1910.09387>

[16] P. Anderson, B. Fernando, M. Johnson, and S. Gould, “SPICE: semantic propositional image caption evaluation,” *CoRR*, vol. abs/1607.08822, 2016. [Online]. Available: <http://arxiv.org/abs/1607.08822>

[17] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” arXiv preprint arXiv:1412.6980, 2014.

[18] Y. Wu, K. Chen, Z. Wang, X. Zhang, F. Nian, S. Li, and X. Shao, “Audio captioning based on transformer and pre-training for 2020 DCASE audio captioning challenge,” *DCASE2020 Challenge*, Tech. Rep., June 2020.

[19] H. Wang, B. Yang, Y. Zou, and D. Chong, “Automated audio captioning with temporal attention,” *DCASE2020 Challenge*, Tech. Rep., June 2020.

[20] Y. Koizumi, D. Takeuchi, Y. Ohishi, N. Harada, and K. Kashino, “The NTT DCASE2020 challenge task 6 system: Automated audio captioning with keywords and sentence length estimation,” *DCASE2020 Challenge*, Tech. Rep., June 2020