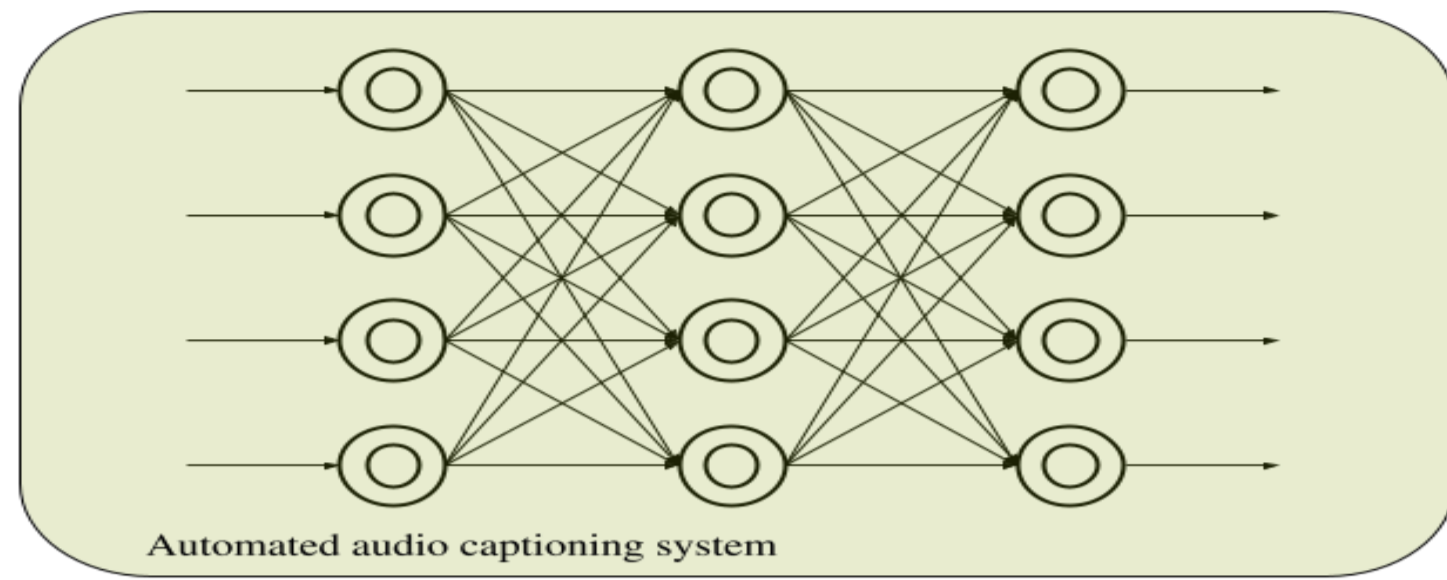




Introduction



The typewriter clacks along, slowly and quickly and when the bell rings the carriage is moved to the next line.

Figure 1: An example of an AAC system and process.

Automated Audio Captioning (AAC) automatically creates captions that can explain the given audio sound data using machine learning techniques.

Researchers investigate the solutions for AAC on DCASE 2021 audio captioning challenge(task 6). In the challenge, a model is required to generate natural language descriptions of a given audio signal.

- We propose transfer learning followed by transformer architecture. With the transfer learning, our proposed model takes two pre-trained networks, 14-layers CNN (CNN14) and 54- layers ResNet (ResNet54), trained on AudioSet as the encoder part.
- Audio feature preprocessing method is log-mel spectrogram.
- Data augmentation method is spec augmentation

<http://dcase.community/challenge2021/task-automatic-audio-captioning>

Proposed Method

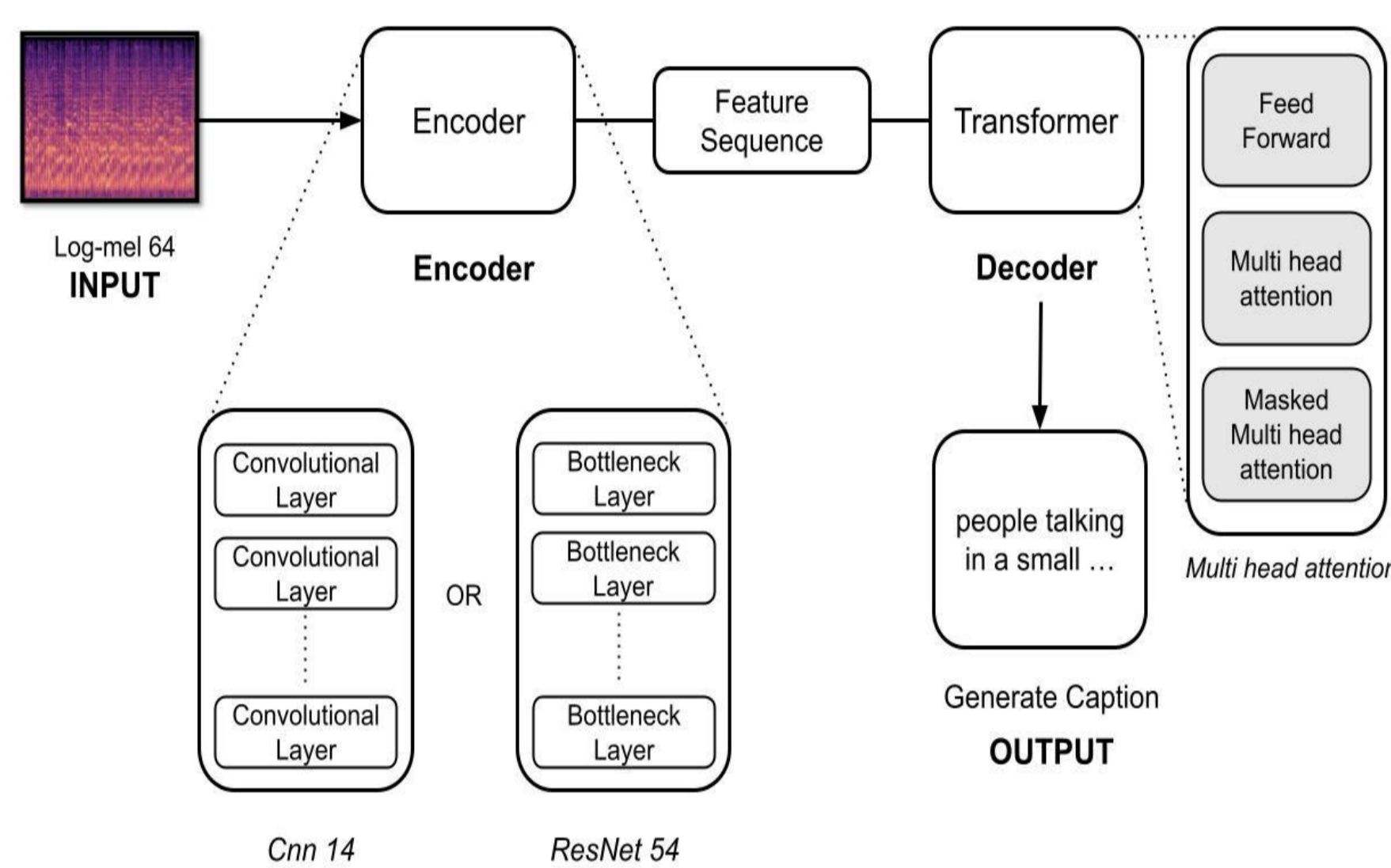


Figure 2. Proposed system architecture

Our model uses CNN14 or ResNet54 as an encoder and Transformer Decoder for natural language generation.

Encoder

- CNN14: The 14-layer CNN consists of four convolution blocks, each having two 3×3 convolution layers with ReLU activation function and batch normalization, with an 2×2 average pooling layer between the blocks.

- ResNet54: We had three bottleneck blocks with 64 filters, four bottleneck blocks with 128 filters, 6 bottleneck blocks with 256 filters, and 3 bottleneck blocks with 512 filters. Finally, two 3×3 convolutions are applied.

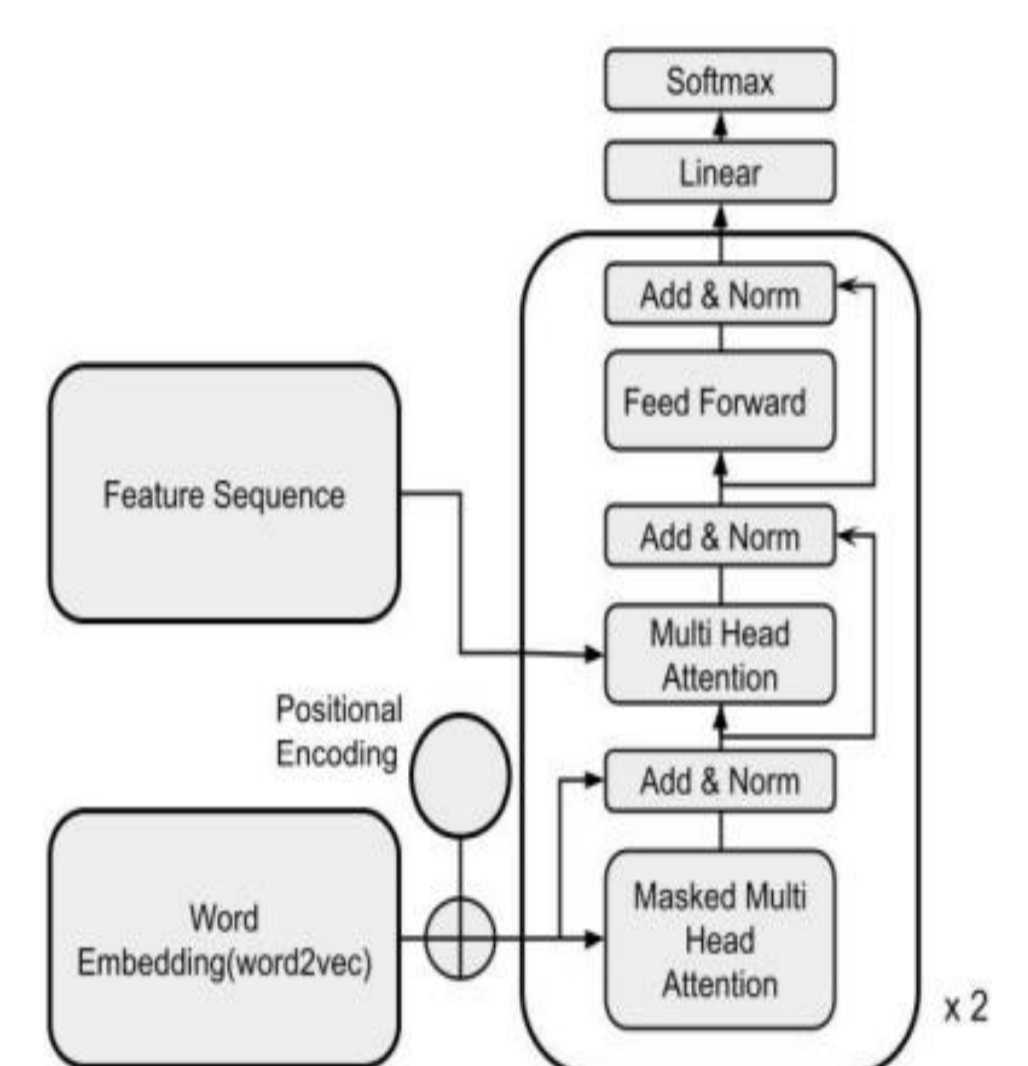


Figure 3: Multi head attention(Decoder)

Decoder

The decoder uses a standard transformer decoder consisting of multi-head self-attention.

The decoder uses a 2-layers transformer with a hidden dimension of 192 and 4 heads.

CNN14
Log-mel spectrogram 64 mel bins
$(3 \times 3 @64, \text{BN}, \text{ReLU}) \times 2$
Pooling 2×2
$(3 \times 3 @128, \text{BN}, \text{ReLU}) \times 2$
Pooling 2×2
$(3 \times 3 @256, \text{BN}, \text{ReLU}) \times 2$
Pooling 2×2
$(3 \times 3 @512, \text{BN}, \text{ReLU}) \times 2$
Pooling 2×2
$(3 \times 3 @1024, \text{BN}, \text{ReLU}) \times 2$
Pooling 2×2
$(3 \times 3 @2048, \text{BN}, \text{ReLU}) \times 2$

Table 1: CNN14(Encoder)

ResNet54
Log-mel spectrogram 64 mel bins
$(3 \times 3 @512, \text{BN}, \text{ReLU}) \times 2$
Pooling 2×2
$(\text{bottleneckB}@64) \times 3$
Pooling 2×2
$(\text{bottleneckB}@128) \times 4$
Pooling 2×2
$(\text{bottleneckB}@256) \times 6$
Pooling 2×2
$(\text{bottleneckB}@512) \times 3$
Pooling 2×2
$(3 \times 3 @512, \text{BN}, \text{ReLU}) \times 2$

Table 2: ResNet54(Encoder)

Result

Model	BLEU ₁	BLEU ₂	BLEU ₃	BLEU ₄	ROUGE _L	METEOR	CIDE _r	SPICE	SPIDE _r
Baseline Model	0.378	0.119	0.050	0.017	0.263	0.078	0.075	0.028	0.051
CNN14 + Transformer (From Scratch)	0.466	0.262	0.156	0.092	0.309	0.137	0.208	0.087	0.148
ResNet54 + Transformer (From Scratch)	0.459	0.253	0.152	0.084	0.312	0.131	0.182	0.085	0.133
CNN14 + Transformer (Transfer-learning with fine-tuning)	0.552	0.364	0.244	0.159	0.378	0.168	0.395	0.118	0.257
ResNet54 + Transformer (Transfer-learning with fine-tuning)	0.546	0.358	0.239	0.156	0.373	0.166	0.379	0.113	0.246
CNN14 + Transformer (Transfer-learning)	0.564	0.376	0.254	0.163	0.388	0.177	0.441	0.128	0.285
ResNet54 + Transformer (Transfer-learning)	0.540	0.345	0.230	0.152	0.361	0.161	0.383	0.109	0.246

Table 3: Score for model performance on evaluation data

All scenarios of CNN14 Encoder+Transformer Decoder and ResNet54 Encoder+Transformer Decoder models have higher scores than the baseline model in all evaluation metrics. Also, All Transfer learning scenarios have better performance than the training from the scratch scenario. We also experiment with fine-tuning the last convolution block of encoder networks (CNN14 and ResNet54 models). For CNN14 Encoder + Transformer Decoder, the transfer-learning model without fine-tuning works better than with the fine-tuning scenario in all evaluation metrics.

Conclusion

- In DCASE 2021, the use of external data is allowed. Thus, we propose transfer learning followed by a transformer approach. We adopt CNN14 and ResNet54 pre-trained on AudioSet data because it achieves state-of-the-art performance on audio pattern recognition.
- The pre-trained CNN14 or ResNet54 models are taken as encoder networks for informative audio feature extraction. With the transferred encoder and a transformer decoder, our proposed systems outperform the baseline system with all evaluation metrics.
- Further, we experiment with three training scenarios, 1) from scratch, 2) transfer learning, and 3) transfer learning with fine-tuning.
- Among them, the transfer learning of CNN14 encoder without fine-tuning works the best, achieving a SPIDE_r score of 0.285. Our proposed system ranked 6th in this competition for DCASE 2021 task 6.