



Audio Flamingo: A Novel Audio Language Model with Few-Shot Learning and Dialogue Abilities



Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, Bryan Catanzaro
{zkong, wping, rafaelvalle}@nvidia.com

Summary

- We propose **Audio Flamingo**: a Flamingo-based audio language model for audio understanding.
- Audio Flamingo achieves SOTA results on several close-ended and open-ended audio understanding tasks.
- We design a series of methodologies for efficient use of in-context learning and retrieval, which lead to SOTA few-shot learning results.
- Audio Flamingo has strong multi-turn dialogue ability.
- Code and checkpoints at <https://github.com/NVIDIA/audio-flamingo>
- Sound demos at <https://audioflamingo.github.io/>

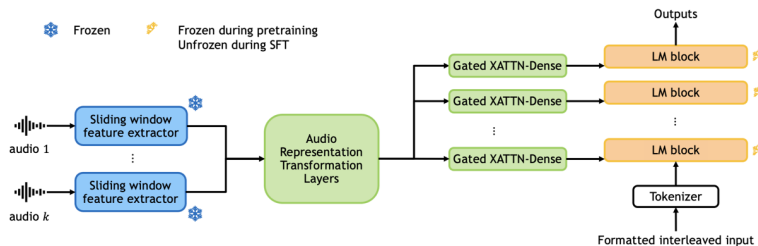
Motivation

We aim to build an audio language model that can understand sound beyond speech transcriptions.

Tasks

- ✓ Audio Captioning
- ✓ Audio Question Answering
- ✓ Audio Classification
- ✓ Retrieval-augmented few-shot learning
- ✓ Multi-turn dialogues

Architecture



Training

Our training objective combines non-interleaved and interleaved samples.

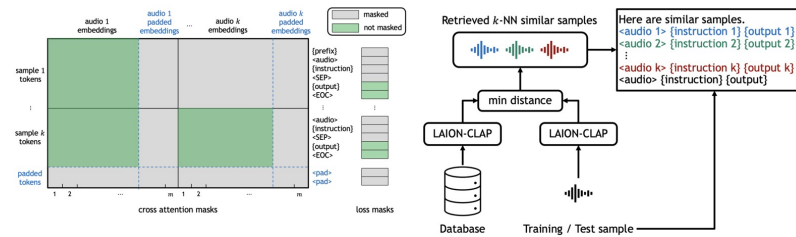
$$\mathcal{L}(z) = \sum_{t=1}^{|y_{out}|} \log p_{\theta}((y_{out})_t | x, y_{ins}, (y_{out})_{<t})$$

$$\mathcal{L}_{int}(z_{int} = \{z^1, \dots, z^J\}) = \sum_{j=1}^J \sum_{t=1}^{|y_{out}^j|} \log p_{\theta}((y_{out}^j)_t | z^{<j}, x^j, y_{ins}^j, (y_{out}^j)_{<t})$$

$$L = - \sum_{i \in \mathcal{I}} \lambda_i \mathbb{E}_{z \sim \mathcal{D}^i} \mathcal{L}(z) - \sum_{i' \in \mathcal{I}_{int}} \lambda_{i'} \mathbb{E}_{z_{int} \sim \mathcal{D}_{int}^{i'}} \mathcal{L}_{int}(z_{int})$$

Retrieval Augmented Generation and In-Context Learning

We use a block-triangular cross-attention mask for interleaved data (left), and a retrieval method to construct interleaved training samples (right).



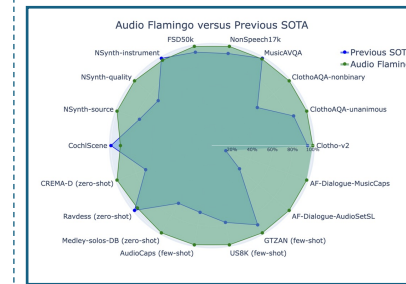
Templates

`<s>[task description]<audio>[instruction]
Options: \n- option_1\n- ... option_n
<SEP>[output]<EOC></s>`

`<s>The task is dialogue.<audio>
user: [instruction]
assistant: <SEP>[output]<EOC>
...
user: [instruction]
assistant: <SEP>[output_s]<EOC></s>`

`<s>[task description]Here are similar samples.
<audio>[instruction_1]<SEP>[output_1]<EOC>
...
<audio>[instruction_k]<SEP>[output_k]<EOC>
<audio>[instruction]
Options: \n- option_1\n- ... option_n
<SEP>[output]<EOC></s>`

Results



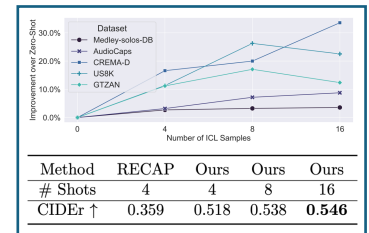
Dataset	Task	Metric	Previous SOTA ↑	Ours ↑
Clotho-v2	CAP	CIDEr	0.441 (Chu et al., 2023)	0.465
ClothoAQA _{unimodal}	AQA	ACC	74.9% (Chu et al., 2023)	86.9%
ClothoAQA _{non-binary}	AQA	ACC	29.1% (Deshmukh et al., 2023)	49.5%
ClothoAQA _{unimodal}	AQA	ACC	26.2% (Deshmukh et al., 2023)	36.4%
MusicAQA _{audio-only}	AQA	ACC	72.1% (Chu et al., 2023)	71.6%
CochiScene	CLS	ACC	91.6% (Deshmukh et al., 2023)	83.0%
NonSpeech7k	CLS	ACC	79.0% (Rashed et al., 2023)	85.1%
PSD50k	CLS	F _{approx}	65.6% (Deshmukh et al., 2023)	69.7%
NS _{instrument}	CLS	ACC	78.8% (Chu et al., 2023)	77.1%
NS _{quality}	CLS	F1	46.3% (Deshmukh et al., 2023)	66.7%
NS _{source}	CLS	ACC	60.1% (Deshmukh et al., 2023)	78.7%

Dataset	Task	Metric	Previous SOTA (0-shot) ↑	Ours (0-shot) ↑
AudioCaps (Kim et al., 2019)	CAP	CIDEr	0.281 (Salowski et al., 2023)	0.502
CREMA-D (Cui et al., 2014)	CLS	ACC	18.5% (Deshmukh et al., 2023)	26.5%
Rawsets (Livingstone & Rame, 2018)	CLS	ACC	21.7% (Elizalde et al., 2023)	20.9%
US8K (Salamon et al., 2014)	CLS	ACC	71.9% (Deshmukh et al., 2023)	75.0%
GTZAN (Storpe, 2011)	CLS	ACC	71.0% (Rao et al., 2023)	67.9%
Medley-solos-DB (Lesteläinen et al., 2019)	CLS	ACC	61.3% (Deshmukh et al., 2023)	92.7%

Audio Flamingo achieves SOTA generation quality on several audio understanding tasks. **Left:** overview of results, where 100% means the best of all baseline and our models. **Upper right:** in-domain tasks (evaluated on test splits). **Lower right:** 0-shot evaluation results.

Testset	Method	CIDEr ↑	Bleu4 ↑	R-L ↑
A	Qwen-Audio	0.507	0.060	0.292
A	LTU [†]	0.823	0.153	0.403
A	Ours [†]	1.622	0.237	0.473
M	MU-LLaMA	0.585	0.083	0.348
M	LTU [†]	0.419	0.108	0.336
M	Ours [†]	1.143	0.142	0.417

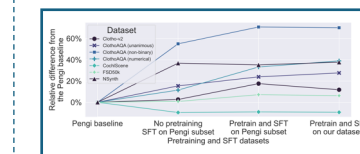
Dialogue evaluation results on our Audio-Dialogues dataset. † indicates the model is finetuned on Audio-Dialogues. A is the audio subset and M is the music subset.



Retrieval augmented in-context learning improves Audio Flamingo's generation quality. **Upper:** relative improvements with respect to #few-shot samples. **Lower:** retrieval-augmented audio captioning results.

References

- [1] Alayrac, Jean-Baptiste, et al. Flamingo: a visual language model for few-shot learning. NeurIPS 2022.
- [2] Iyer, Srinivasan, et al. OPT-ML: Scaling language model instruction meta learning through the lens of generalization. arXiv 2022.
- [3] Elizalde, Benjamin, Soham Deshmukh, and Huaming Wang. Natural language supervision for general-purpose audio representations. ICASSP 2024.
- [4] Wu, Yulong, et al. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. ICASSP 2023.
- [5] Deshmukh, Soham, et al. Pengi: An audio language model for audio tasks. NeurIPS 2023.
- [6] Chu, Yunfei, et al. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. arXiv 2023.
- [7] Tang, Changli, et al. SALMONN: Towards Generic Hearing Abilities for Large Language Models. ICLR 2024.
- [8] Gong, Yuan, et al. Listen, Think, and Understand. ICLR 2024.
- [9] Liu, Shansong, et al. Music understanding LLaMA: Advancing text-to-music generation with question answering and captioning. ICASSP 2024.
- [10] Ghosh, Sreyan, et al. RECAP: retrieval-augmented audio captioning. ICASSP 2024.



Effects of data scaling.