# Audio Flamingo: Advancing Audio Intelligence with Open and Powerful Large Audio-Language Models

Sreyan Ghosh*, Arushi Goel*, Jaehyeon Kim, Sonal Kumar, Zhifeng Kong, Sang-gil Lee, Chao-Han Huck Yang, Ramani Duraiswami, Dinesh Manocha, Rafael Valle, Bryan Catanzaro
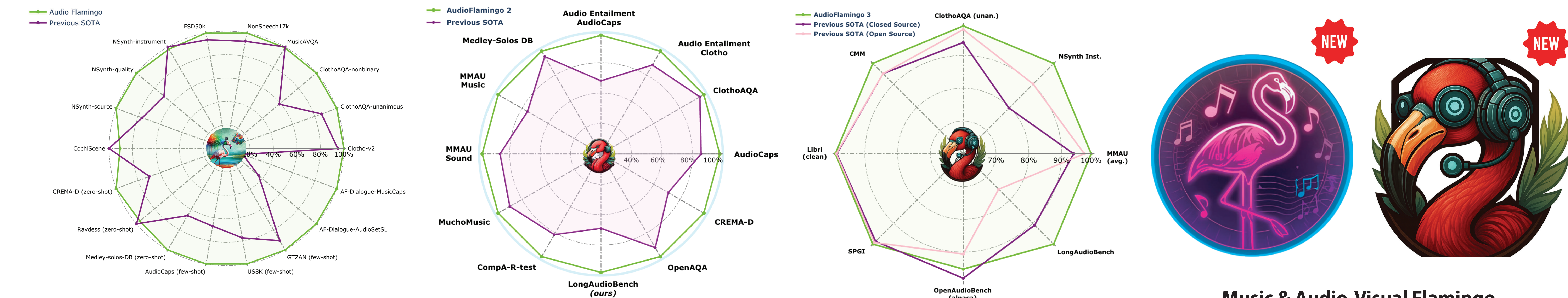
NVIDIA · UNIVERSITY OF MARYLAND

## The Audio Flamingo Series



**Audio Flamingo 1**
- Focuses on foundational audio processing tasks (e.g., captioning)
- Introduced RAG and multi-turn dialogue abilities
- ~5.9 M training instances

**Audio Flamingo 2**
- Focuses on improving performance on expert reasoning tasks
- Introduced long-audio understanding and reasoning
- ~10 M training instances and synthetic data

**Audio Flamingo 3**
- Focuses on advancing comprehensive audio understanding and reasoning
- Introduces speech in the Audio Flamingo series
- Introduces long-speech, thinking abilities and multi-turn, multi-audio chat
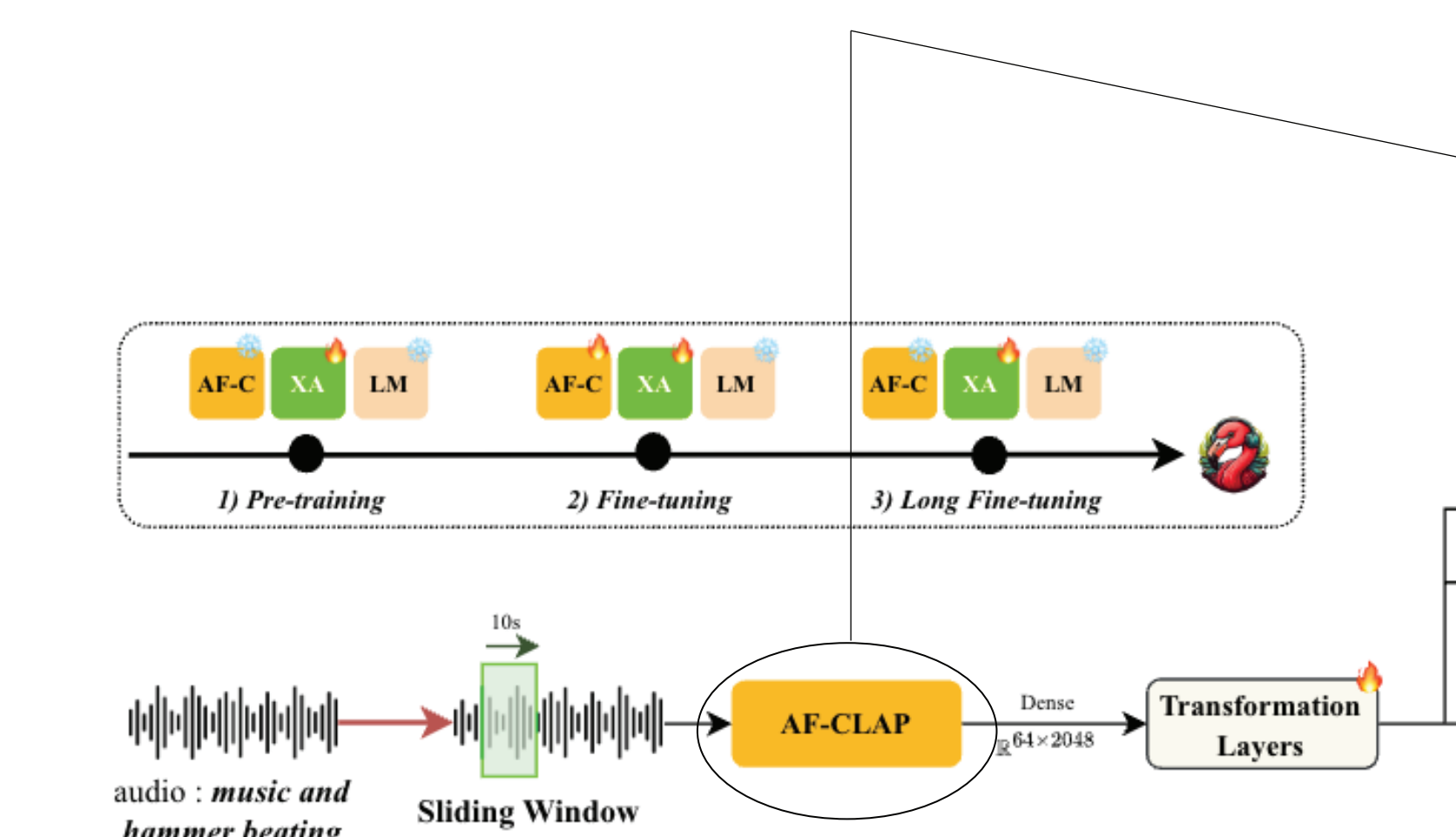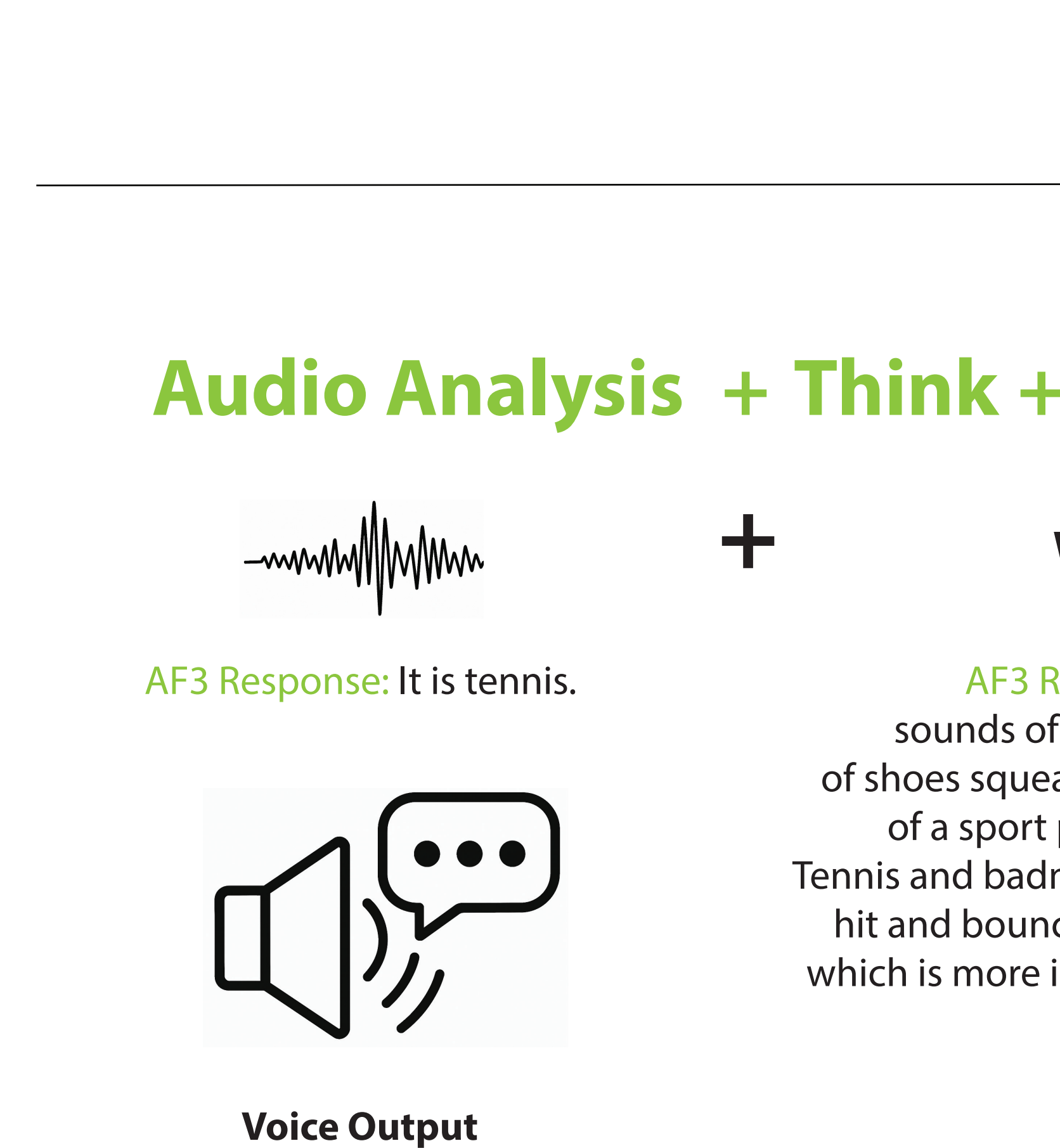- ~50 M training instances

**Music & Audio-Visual Flamingo**
- Focused on deep understanding of full-length songs.
- With Omni Flamingo, we bring audio-intelligence to visual understanding models

**Improved Perception = Improved Performance!!**

Cross-Attention for Audio Flamingo 1 & 2

Prefix for Music & Audio Flamingo 3

Audio-Visual Flamingo

## Data is All You Need!

**We present AudioSkills (XL) with 10M+ AQAs across 12+ skills.**



## Open datasets focused on Long Audios and Full-length Songs



## Audio Analysis + Think + Voice Chat (streaming)

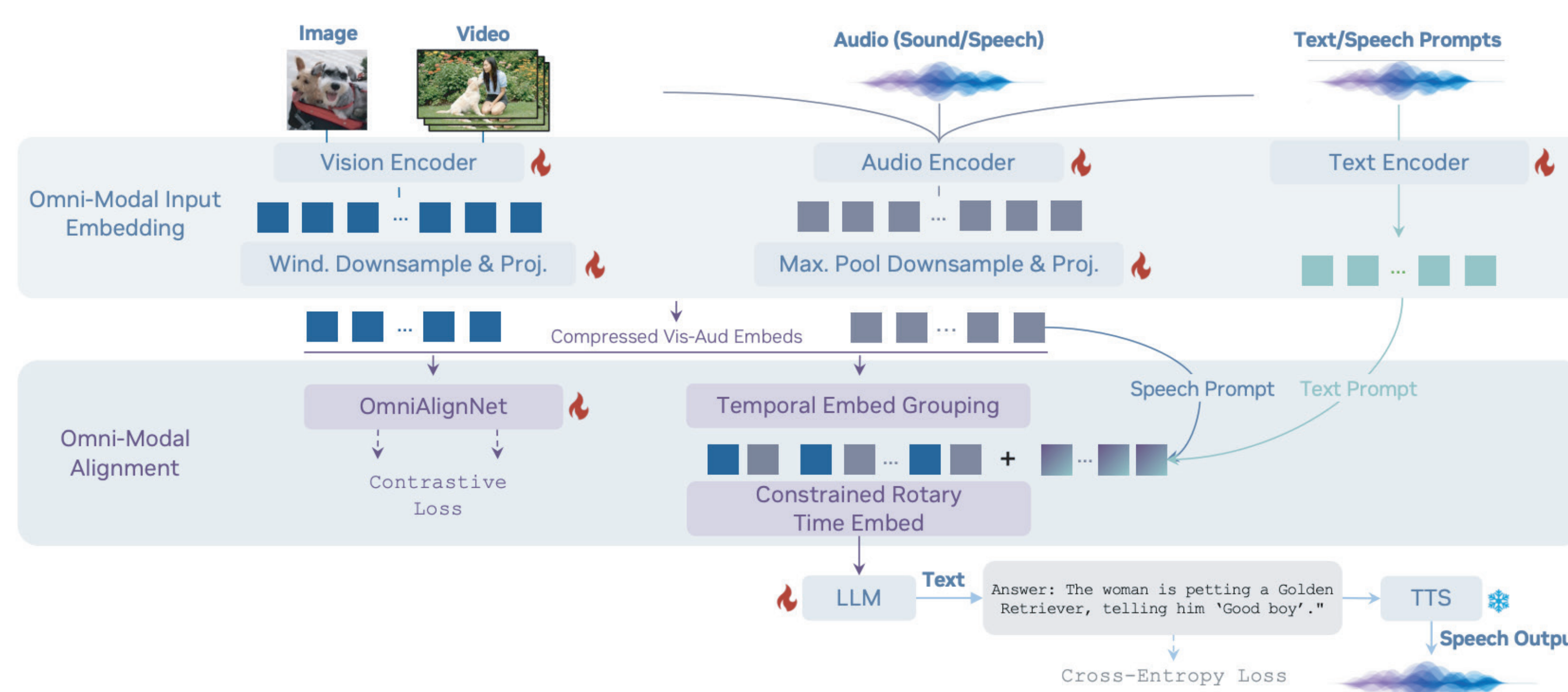**What sport is being played in the audio?**

AF3 Response: It is tennis.

AF3 Response in Thinking Mode: The audio contains sounds of a ball being hit and bounced, along with the sound of shoes squeaking on a hard surface. These sounds are characteristic of a sport played on a hard court with a racket or shuttlecock. Tennis and badminton are both such sports. The presence of a ball being hit and bounced suggests a sport where a ball, which is more indicative of tennis than badminton. Therefore, the most likely sport is tennis.
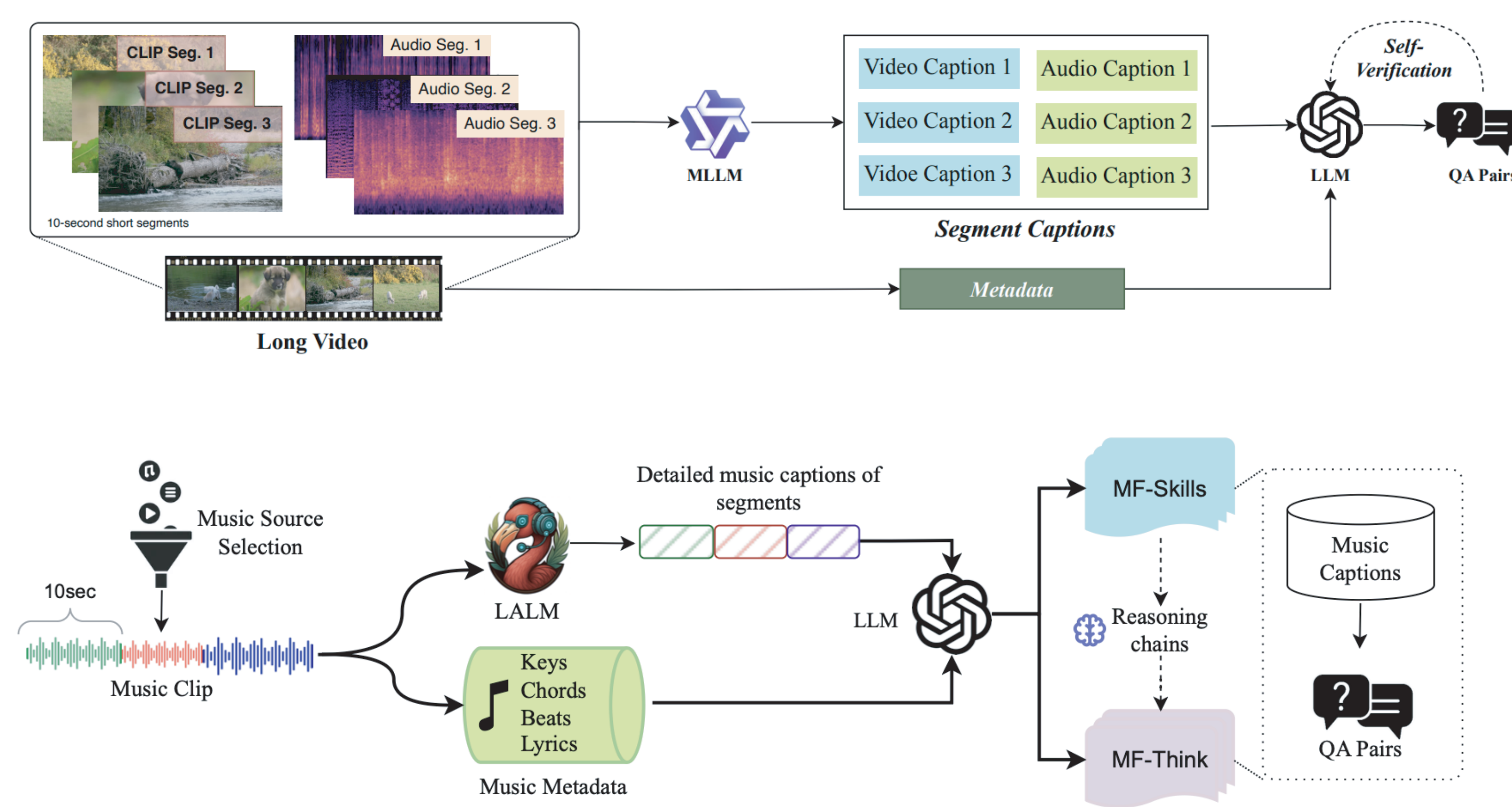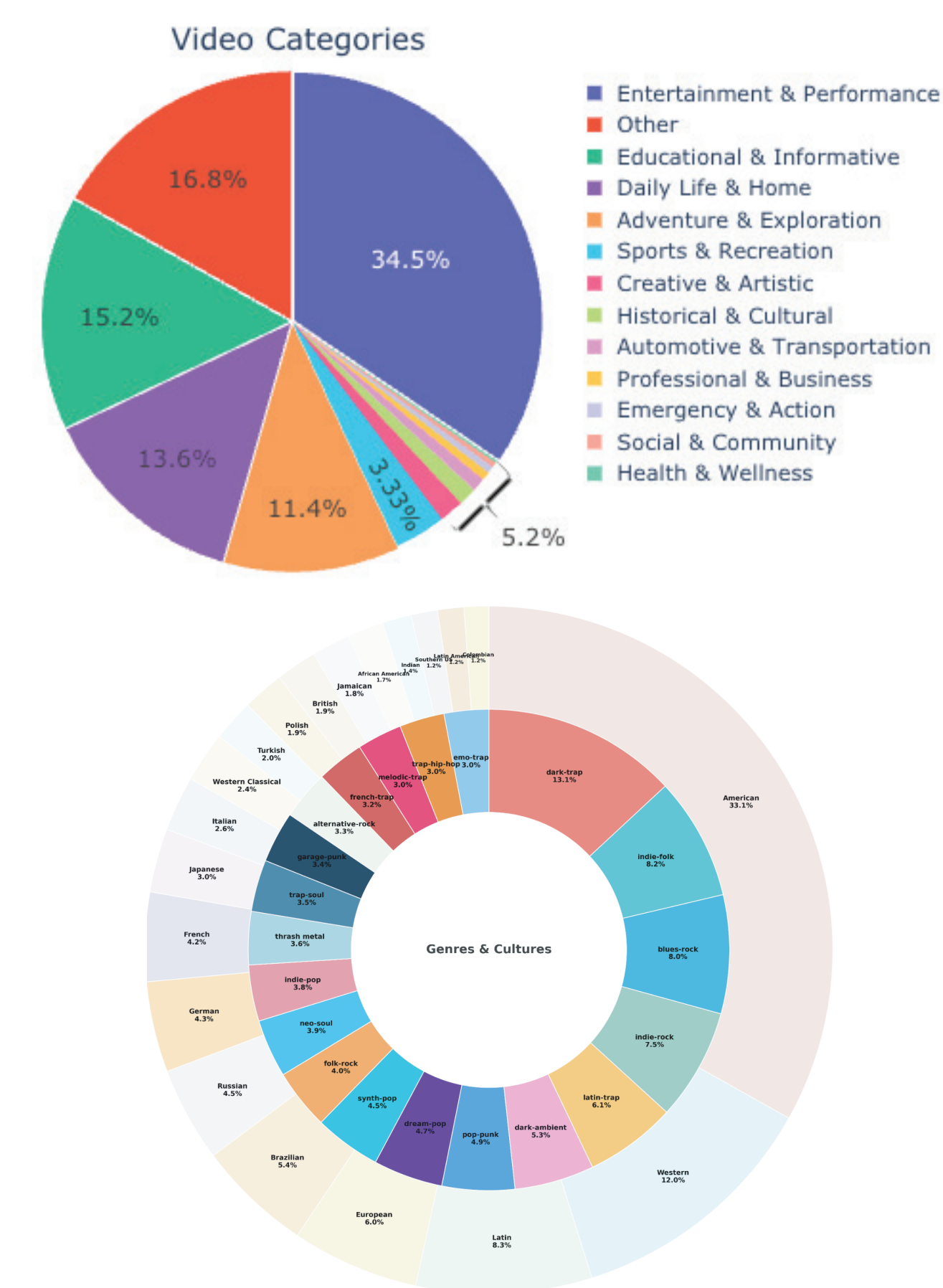
Voice Output

## Multi-turn Multi-Audio Chat

**What elements of the track make it feel so calming and meditative?**

AF3 Response: The calming effect comes from several factors: the gentle acoustic instrumentation, the slow and steady tempo, and the soft, expressive vocal delivery.

**If I wanted to remix the calming folk piece into something closer to this, what would you suggest?**
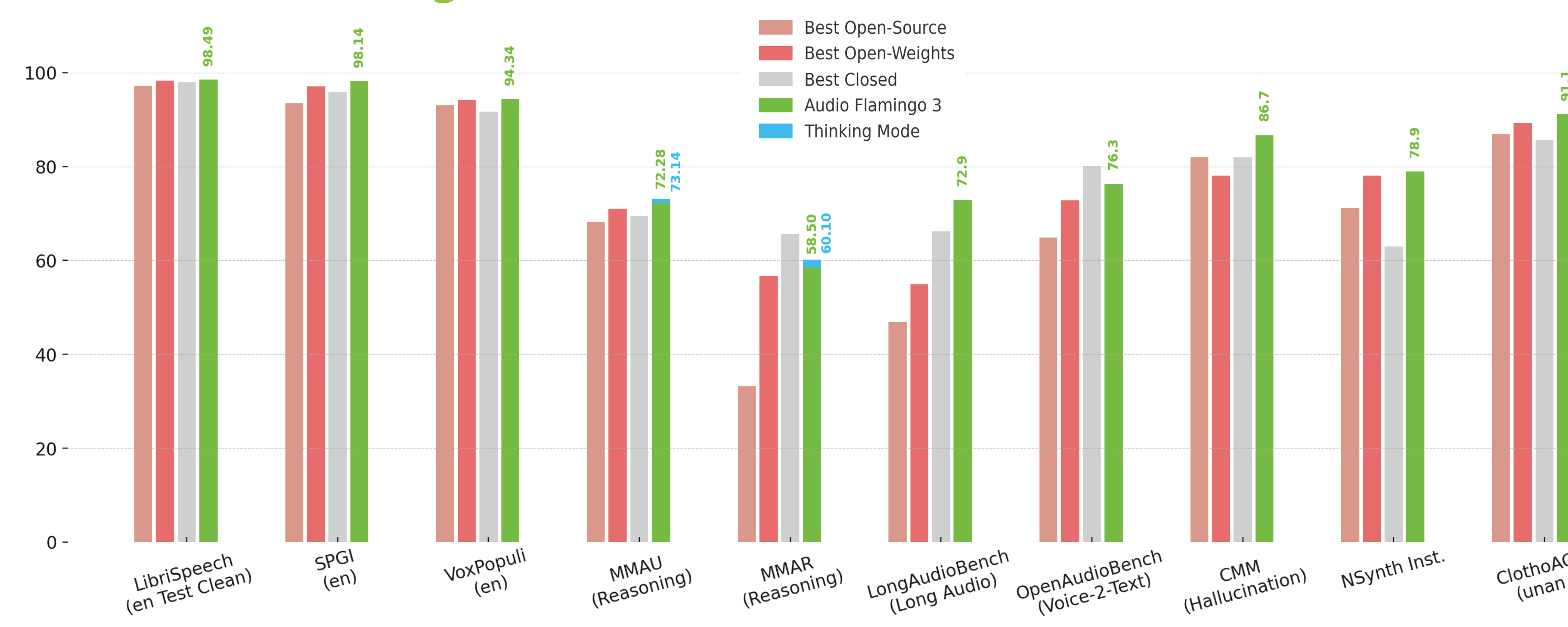
AF3 Response: To transform the folk track into a more energetic, synth-driven style, start by increasing the tempo and introducing electronic drums with a punchy, danceable groove.
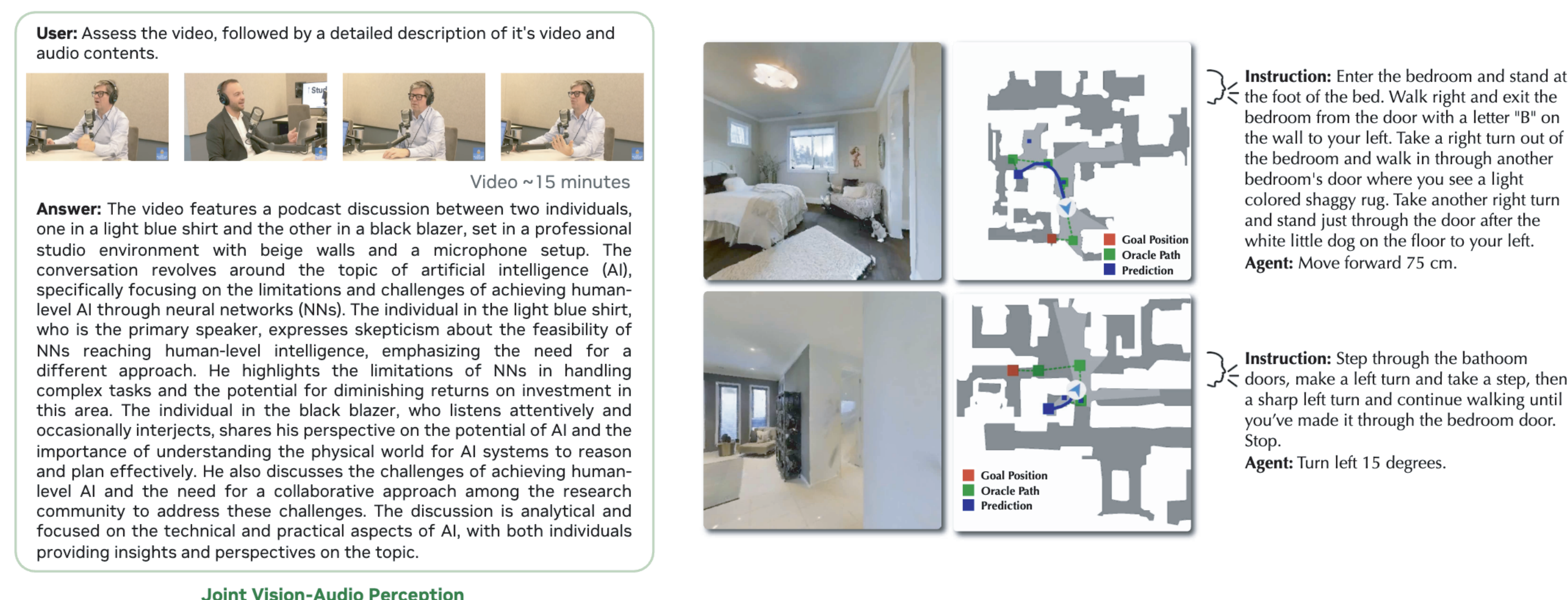
## A Snapshot of MF and AVF capabilities!



## Flamingo Models are SOTA across Benchmarks!