# Title

Hierarchical Audio Concept Distillation for Structured Audio Reasoning in Large Language Models

# Problem Statement

Large Language Models (LLMs) lack a structured understanding of audio concepts, making it difficult for them to reason about complex audio scenes and their relationships to language. This limitation hinders their ability to perform tasks that require detailed audio understanding and interpretation.

# Motivation

Current approaches often treat audio understanding as a flat classification or embedding problem, missing the hierarchical nature of audio concepts. Inspired by cognitive science theories of auditory scene analysis, we propose to distill a hierarchical representation of audio concepts into LLMs, enabling more structured and interpretable audio reasoning. This approach is motivated by the observation that humans understand audio in a hierarchical manner, from low-level acoustic properties to high-level semantic categories. By mimicking this cognitive process, we aim to enhance LLMs' ability to reason about audio in a more human-like and interpretable way.

# Proposed Method

We present Hierarchical Audio Concept Distillation (HACD), a framework for injecting structured audio knowledge into LLMs. HACD involves four main steps: (1) Constructing a large-scale Audio Concept Hierarchy (ACH) that organizes audio concepts from low-level acoustic properties to high-level semantic categories. (2) Training a teacher model to classify audio inputs according to the ACH. (3) Distilling the teacher's knowledge into the LLM through a novel hierarchical distillation loss that preserves the concept structure. (4) Fine-tuning the LLM on audio-language tasks while maintaining the distilled hierarchical knowledge. The resulting model can reason about audio inputs using a structured vocabulary of concepts, enabling more precise and interpretable audio understanding.

# Step-by-Step Experiment Plan

## Step 1: Construct Audio Concept Hierarchy (ACH)

Create a comprehensive hierarchy of audio concepts, ranging from low-level acoustic properties (e.g., pitch, timbre) to high-level semantic categories (e.g., music genres, environmental sounds). Use existing audio ontologies and expert knowledge to build this hierarchy. Ensure the ACH covers a wide range of audio domains and has at least 1000 concepts organized in a tree structure with 5-7 levels of depth.

## Step 2: Prepare Audio Dataset

Collect a large-scale audio dataset that covers all concepts in the ACH. Use existing datasets like AudioSet, FSD50K, and ESC-50, and supplement with additional data if needed. Ensure each audio clip is labeled with all applicable concepts from the ACH. Aim for at least 1 million labeled audio clips.

## Step 3: Train Teacher Model

Develop a multi-label hierarchical classification model (e.g., based on ResNet or Transformer architecture) that can classify audio inputs according to the ACH. Train this model on the prepared audio dataset. Use hierarchical loss functions that account for the concept structure. Evaluate the teacher model's performance using metrics like mean average precision (mAP) and hierarchical

precision/recall.

## Step 4: Implement Hierarchical Distillation

Design a novel hierarchical distillation loss that captures the structure of the ACH. This loss should encourage the LLM to learn not just individual concepts but also their relationships. Implement the distillation process by feeding audio embeddings (extracted by the teacher model) along with text to the LLM, and training it to predict the hierarchical concept probabilities.

## Step 5: Fine-tune LLM

Fine-tune the LLM on a variety of audio-language tasks while maintaining the distilled hierarchical knowledge. Tasks should include audio captioning, audio question-answering, and audio-based text generation. Use datasets like Clotho, AUDIO CAPS, and custom-created datasets based on the ACH.

## Step 6: Evaluation

Evaluate the HACD-enhanced LLM on a range of tasks that require structured audio reasoning: (1) Detailed audio scene description: Generate hierarchical descriptions of complex audio scenes. (2) Hierarchical audio classification: Classify audio inputs at different levels of the ACH. (3) Audio-based question answering: Answer questions about audio content at varying levels of abstraction. (4) Cross-modal retrieval: Retrieve relevant audio clips given textual descriptions and vice versa. Compare performance against baseline LLMs with standard audio encoders and models trained with flat audio classification. Use both standard metrics (e.g., BLEU, METEOR for captioning; accuracy, F1-score for classification) and novel metrics designed to assess the hierarchical nature of the model's understanding (e.g., concept precision at different hierarchy levels).

## Step 7: Qualitative Analysis

Conduct in-depth qualitative analyses of the model's outputs. Examine the model's ability to explain audio concepts at varying levels of abstraction. Visualize attention patterns to understand how the model leverages the hierarchical structure. Perform error analysis to identify areas for improvement.

# Test Case Examples

## Baseline Prompt Input

Describe the audio in detail: [Audio of a jazz band playing in a busy cafe]

## Baseline Prompt Expected Output

The audio contains music playing in a noisy environment. There are people talking in the background.

## Proposed Prompt Input

Describe the audio in detail, starting from general categories and moving to specific elements: [Audio of a jazz band playing in a busy cafe]

## Proposed Prompt Expected Output

3. Music details: - Genre: Jazz - Instruments: Piano (prominent), Double bass, Drums - Tempo: Moderate, swinging rhythm - Style: Bebop-influenced 4. Human activity: - Type: Conversation - Setting: Indoor, likely a cafe or restaurant - Intensity: Moderate crowd noise 5. Acoustic properties: - Reverberation: Moderate, suggesting a medium-sized room - Audio quality: High fidelity, possibly a professional recording 6. Specific elements: - Piano: Performing a complex solo with fast arpeggios - Double bass: Walking bass line, providing rhythmic foundation - Drums: Subtle brushwork on snare, occasional cymbal accents - Crowd noise: Indistinct chatter, occasional laughter, clinking of glasses

- General category: Environmental audio
- Main components: Music, Human activity
- Music details:
- Human activity:
- Acoustic properties:
- Specific elements:

# Explanation

The HACD-enhanced model provides a much more detailed and structured description of the audio, demonstrating understanding at multiple levels of abstraction. It captures both the high-level categories (environmental audio, music, human activity) and low-level details (specific instruments, playing styles, acoustic properties). The hierarchical nature of the description allows for a more comprehensive and interpretable analysis of the audio scene.

# Fallback Plan

If the proposed HACD method doesn't meet the success criteria, we can pursue several alternative directions: (1) Analyze the learned hierarchical representations to understand where the model struggles. This could involve probing tasks at different levels of the hierarchy to identify which concepts or relationships are not being captured effectively. (2) Experiment with different architectures for integrating the audio hierarchy into the LLM, such as using graph neural networks to explicitly model the concept relationships. (3) Investigate the impact of the Audio Concept Hierarchy's structure on performance, potentially refining the hierarchy based on empirical results. (4) Explore multi-task learning approaches that combine the hierarchical distillation with other audio-related tasks to improve the model's overall audio understanding. (5) If the hierarchical approach proves challenging, we could pivot to a comparative study of different audio representation methods in LLMs, analyzing the trade-offs between structured and unstructured approaches. This could yield valuable insights into the strengths and limitations of various audio understanding paradigms in the context of language models.

Ranking Score: 6