

Human Centered Evaluation of CLAP-Based Emotion Annotations for Therapeutic Music Using MTurk

Liam Stapley, Abhishek Karwankar, Daniel Stevens, & Matthew Louis Mauriello

Motivation

- CLAP (Contrastive Language–Audio Pretraining) generates emotion tags by aligning audio & text embeddings.
- Enables large-scale music annotation without manual labeling.
- Outputs are weakly supervised and may not align with human emotional perception.
- Our work:** Develops an MTurk-based framework to validate & refine CLAP-generated annotations.
 - Can help recommendations for music based on desired emotional profiles

Research Questions

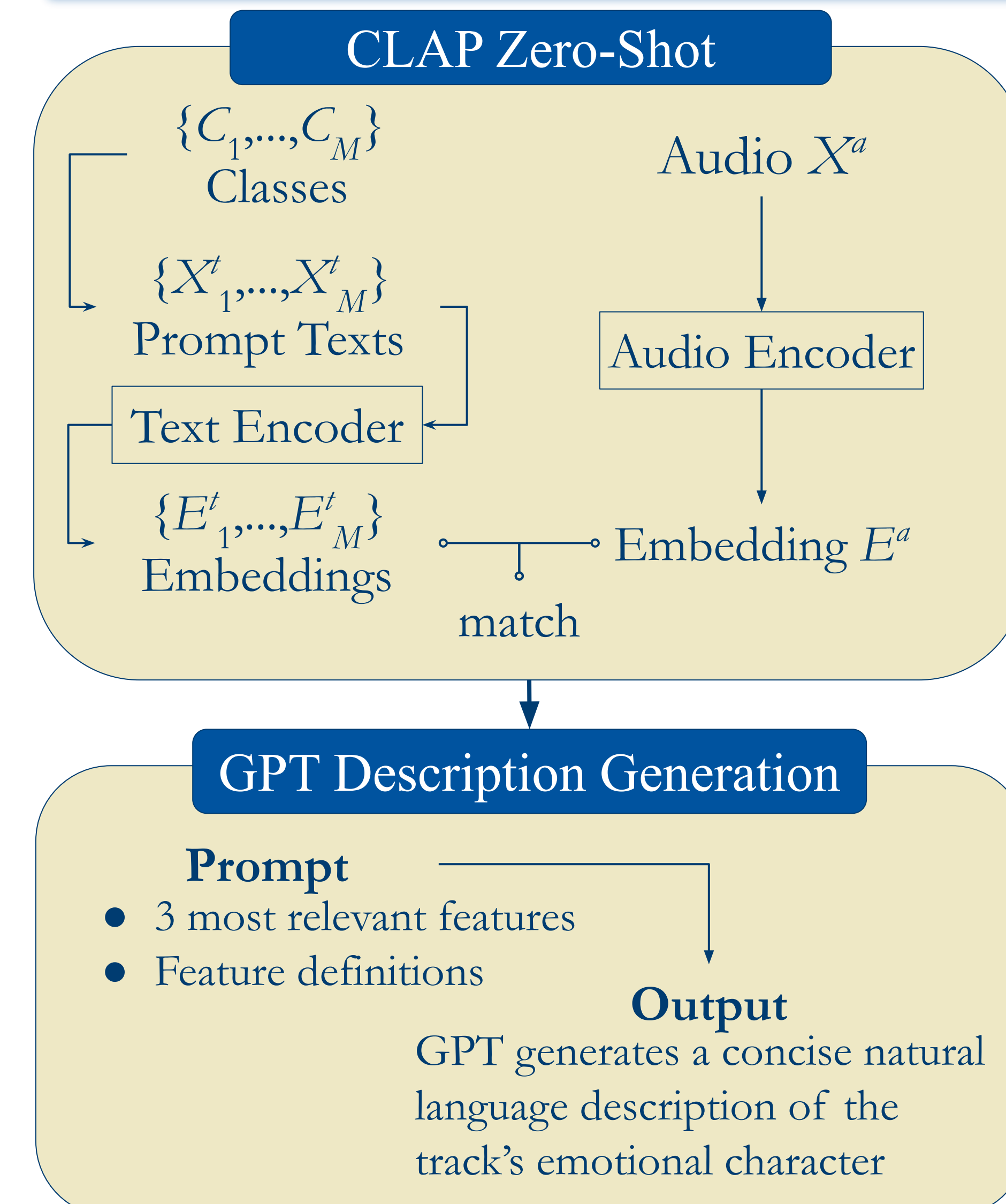
- 🎵 **Can** contrastive learning models like CLAP reliably generate emotionally coherent tags for therapeutic music without explicit supervision?
- 👁️ **How** closely do CLAP-generated emotion annotations align with human emotional perception and interpretation?
- 🔄 **Can** iterative refinement techniques, such as pseudo-labeling based on human feedback, enhance the validity and interpretability of weakly-labeled datasets?

Methodology

Dataset Preparation

- Source:** uCue's music library containing 10k+ possible song combinations.
- Selection criteria:** Chose 432 tracks for diversity in emotion features.
- Structure:** Modular arrangements with separate layers that can be added/removed.
- Metadata:** Includes song IDs, modular layer info, and predefined emotion feature definit
- Purpose:** Provide a representative subset of the full library for CLAP annotation and human validation.

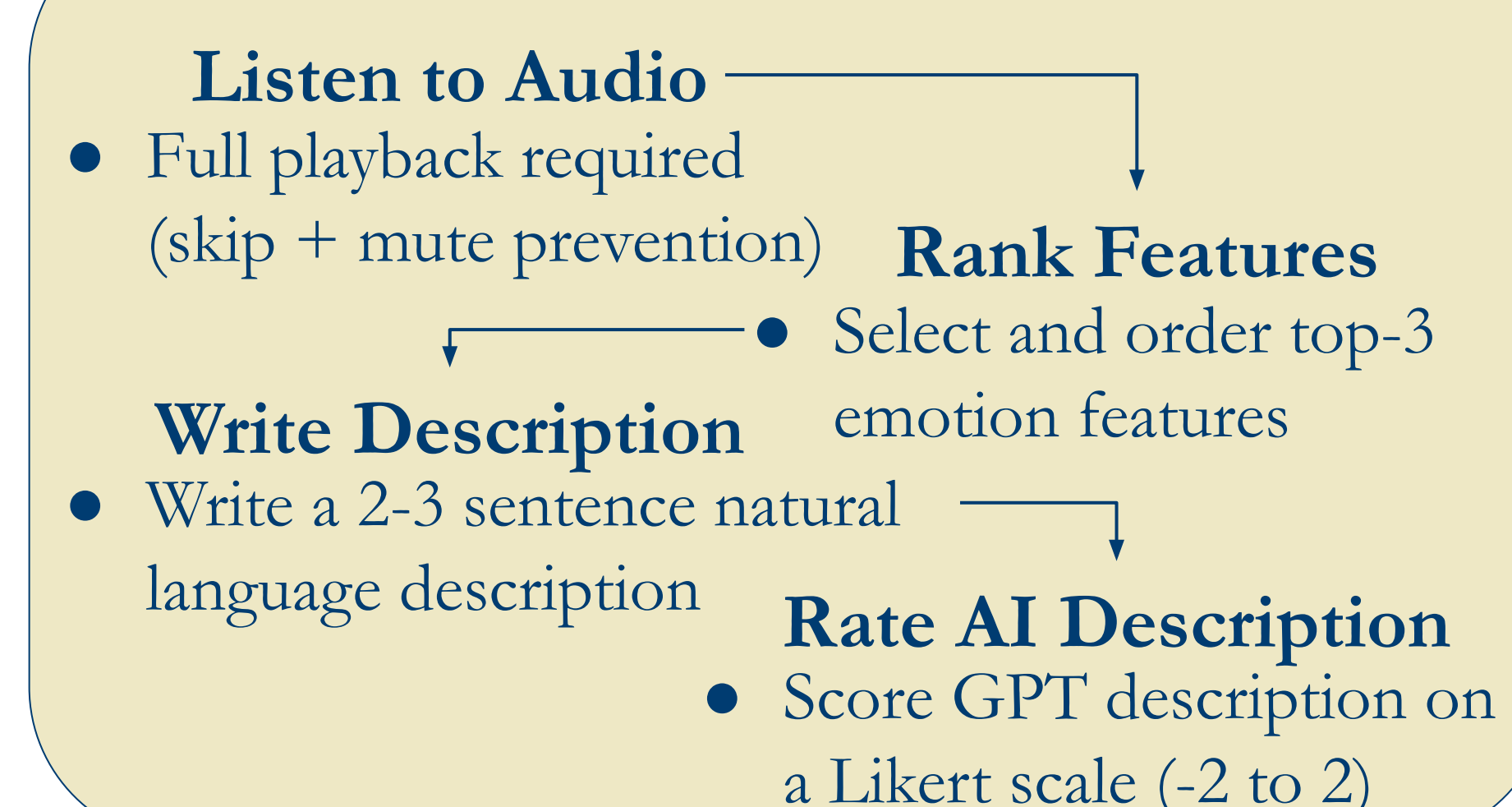
Methodology cont.



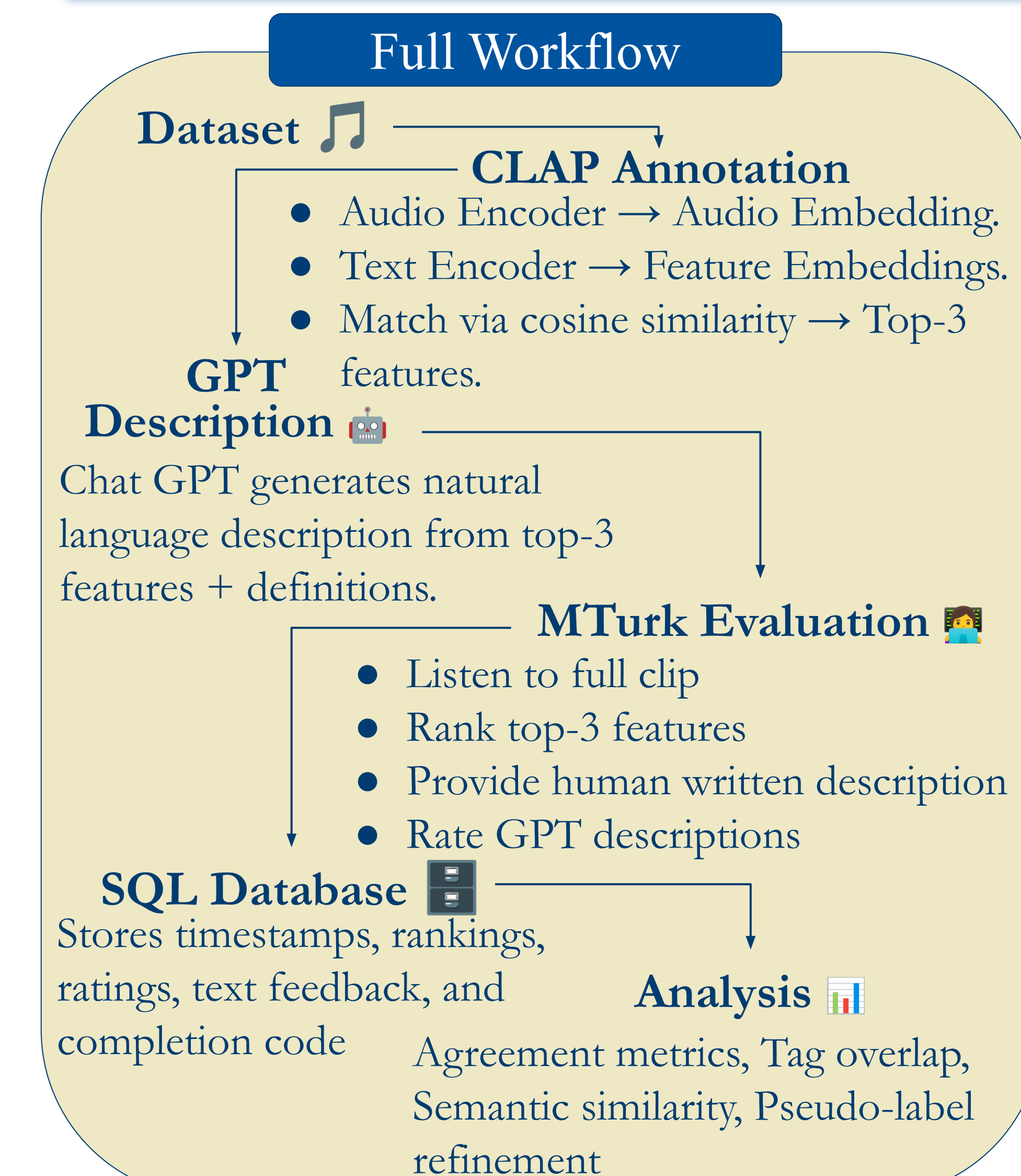
Crowdsourced Evaluation (MTurk)

- Participants: Pre-qualified workers who passed a gold-standard clearance test.
- Task Flow:**
 - Listen to the full audio clip.
 - Rank the top 3 emotion features.
 - 8 features given with definition drop-down available
 - Give 2-3 sentence description.
 - Rate agreement with generated description.
 - Only shown after “Next” button is clicked to prevent bias

MTurk Interface Workflow

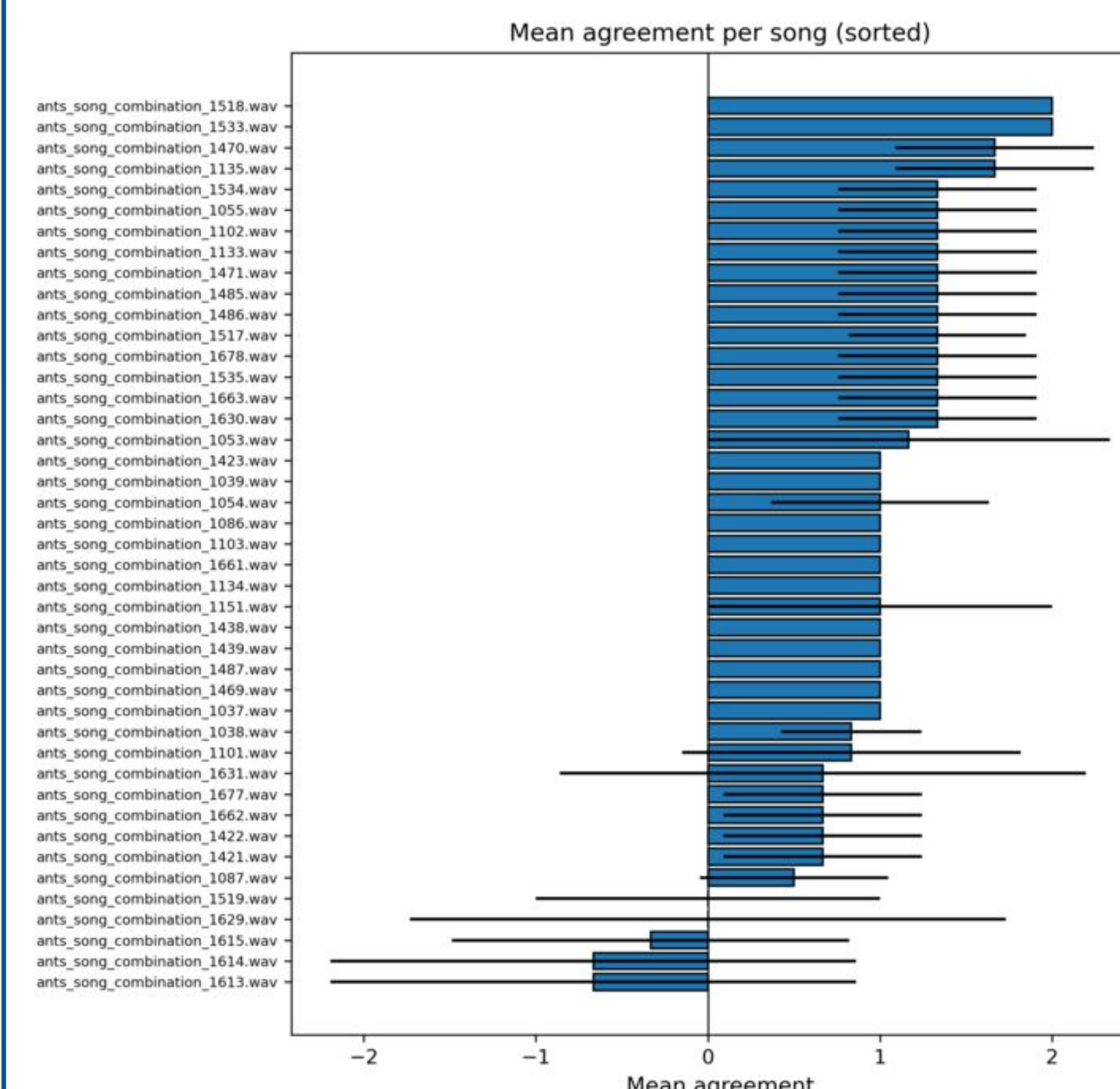


Methodology cont.



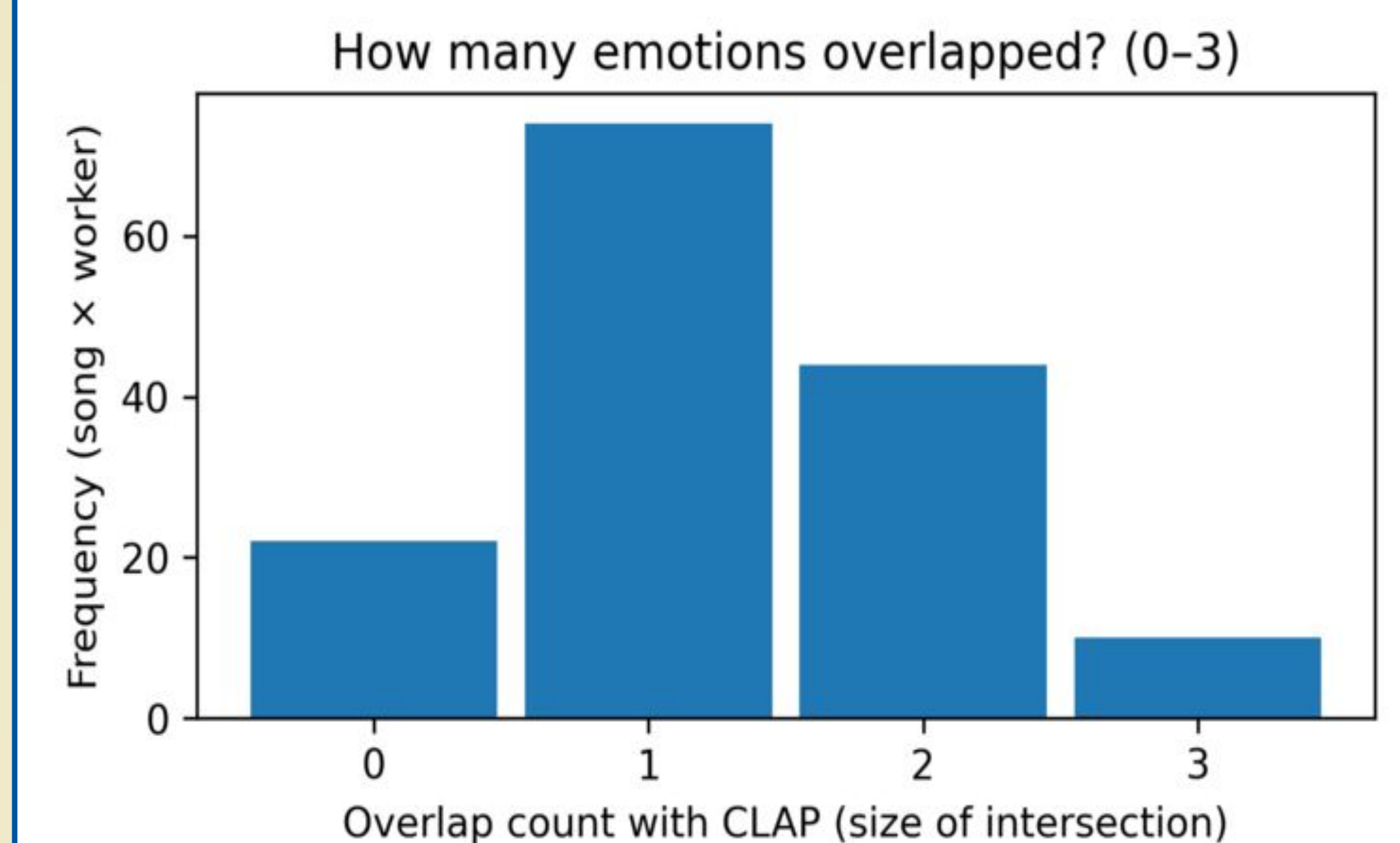
Preliminary Results

- Agreement Score Analysis:** Mean & variance per song; identify high/low consensus cases.



Preliminary Results cont.

- Top-3 Overlap:** Jaccard index & exact match counts between CLAP and human rankings.



Takeaway: CLAP typically shares exactly one label with a given worker.

Future Work

- Complete full dataset evaluation with MTurk.
- Apply pseudo-label refinement to improve annotation accuracy.
- Conduct demographic and contextual analyses of perception patterns.
 - Description Similarity:** Embedding-based semantic similarity between GPT and human descriptions.
 - Cluster Analysis:** Identify patterns of disagreement by emotion type or complexity.
 - Pseudo-Label Refinement:** Use high-agreement samples to improve dataset accuracy.
- Publish validated therapeutic music dataset and open-source evaluation platform.
- Extend methodology to other domains, including environmental audio and affective speech.