000 TEMPORAL ENHANCEMENT CONTRASTIVE OF 001 AUDIO-LANGUAGE MODELS 002 THROUGH SELF-003 SUPERVISED POST-TRAINING TEXT-AUDIO WITH 004 PAIRS 006

Anonymous authors

Paper under double-blind review

ABSTRACT

Research on multi-modal contrastive learning strategies for audio and text has rapidly gained interest. Contrastively trained Audio-Language Models (ALMs), such as CLAP, which establish a unified representation across audio and language modalities, have enhanced the efficacy in various subsequent tasks by providing good text aligned audio encoders and vice versa. These improvements are evident in areas like zero-shot audio classification and audio retrieval, among others. However, the ability of these models to understand natural language and temporal relations is still a largely unexplored and open field for research. In this paper, we propose to equip the multi-modal ALMs with temporal understanding without loosing their inherent prior capabilities of audio-language tasks with a temporal instillation method TeminAL. We implement a two-stage training scheme TeminAL A & B, where the model first learns to differentiate between multiple sounds in TeminAL A, followed by a phase that instills a sense of time, thereby enhancing its temporal understanding in TeminAL B. This approach results in an average performance gain of 5.28% in temporal understanding on the benchmark ESC-50 dataset, while the model remains competitive in zero-shot retrieval and classification tasks on the AudioCap/Clotho datasets. We also note the lack of proper evaluation techniques for contrastive ALMs and propose a strategy for evaluating ALMs in zero-shot settings. The general-purpose Zero-Shot Temporal Evaluation (ZSTE) strategy, is used to evaluate various prior models. ZSTE demonstrates a general strategy to evaluate all ZS contrastive models. The model trained with TeminAL successfully outperforms current models on most downstream tasks.

034 035

037

008

009

010 011 012

013

015

016

017

018

019

021

025

026

027

028

029

031

032

033

1 INTRODUCTION

Audio, text, and images are among the most prevalent forms of information data. Developing models with multi-modal capabilities is well recognized as a path forward toward artificial general intelli-040 gence (Fei et al., 2022; Huang et al., 2021). In the field of multi-modal learning, contrastive learn-041 ing has emerged as an effective strategy for training models on extensive, less-structured internet-042 sourced data (Radford et al., 2021; Liang et al., 2022; Tian et al., 2020). Contrastive learning-based 043 models have demonstrated exceptional adaptability across a range of related tasks, such as image 044 classification (Chen et al., 2020; He et al., 2020a), natural language processing (Gao et al., 2021) and speech processing (Ravanelli et al., 2020), making them a crucial area of research in multi-modal machine learning. One notable early model in this domain is CLIP, developed by Radford et al. 046 (2021). CLIP learns the relationship between text and images, aligning them in a common latent 047 domain. It stands out as a groundbreaking vision-language model, enabling tasks such as generating 048 images from text (Rombach et al., 2022) and formulating image captions (Mokady et al., 2021). 049

Similar work on contrastive learning has been extended to other multi-modal domains, such as video-language (Xu et al., 2021; Fang et al., 2021; Zhao et al., 2022; Luo et al., 2022; Cheng et al., 2023; Ge et al., 2022) and audio-language models (Elizalde et al., 2023; Huang et al., 2022; Guzhov et al., 2022; Wu et al., 2023b; Deshmukh et al., 2023; Wu et al., 2023a). Contrastive models generally excel in relating different modalities through their learned embedding and performing similarity-based



Figure 1: The overview of TeminAL where we are post-training orginal CLAP encoders f_c and f_a with our **TeminAL** method to get f_c^t and f_a^t after application of the two-stage training. We only train a subset of the total weights ($f_{c_{\theta}}^t$ and $f_{a_{\phi}}^t$) in both our training stages. Mathematical formulation of the functions are elaborated in section 3.3 and section 3.4.

068

069

071

075 retrieval tasks. These multi-context encoders integrate well with other downstream models, such as 076 retrieval and open-ended generation models (Ramesh et al., 2021; Li et al., 2022; Yuan et al., 2021; 077 Singh et al., 2022). However, previous authors have shown the limitations of audio-language models in truly understanding natural language while learning the relationship between texts and audio (Wu et al., 2023a; Ghosh et al., 2023). Critical applications like medical procedures, assembly instruc-079 tions, commercial user applications, cooking instructions, and language learning may suffer from mistaken outputs in either text or audio settings. Wu et al. (2023a) highlights a critical limitation 081 in current audio-language models (ALMs): a bias towards retrieving nouns and verbs, often at the expense of understanding the complete sentence context. They illustrate this by training an ALM on 083 captions stripped of all but nouns and verbs, achieving performance comparable to or even surpass-084 ing models trained on full, non-shuffled captions. This finding questions the prevailing assumption 085 that ALMs require holistic sentence comprehension for high performance, revealing gaps in their compositional reasoning capabilities. Furthermore, studies such as Thrush et al. (2022), Ma et al. 087 (2023), and Yuksekgonul et al. (2022) have demonstrated that models like CLIP struggle with lan-088 guage reasoning despite access to extensive training datasets. These limitations arise because contrastive pre-training primarily emphasizes retrieval tasks, enabling strong benchmark performance 089 without a deep understanding of sentence composition. In response to these challenges, Ghosh et al. 090 (2023) critique existing audio-retrieval benchmarks, arguing that the perceived success of ALMs 091 often lacks true compositional understanding. They introduce CompA-CLAP, an ALM designed 092 with novel contrastive training techniques to improve both language comprehension and attribution capabilities in multiple training steps but with the same global objective of making the model di-094 rectly adapt to temporality. Although the model perform well on various downstream tasks, these approaches do not adequately address a fundamental prerequisite for compositional reasoning in au-096 dio tasks: the ability to distinguish multiple sound events before attempting to establish relationships between them. Our work emphasizes this overlooked step, proposing a framework where the model 098 first learns to recognize the existence of multiple sound events as a foundation for higher-level reasoning. Similarly Yuan et al. (2024); Wu et al. (2023a) trains a contrastive learning model without 099 requiring to address the need of multiple sounds distinction which defeats the purpose of increasing 100 the interpretability of the models. 101

In contrast, our approach achieves this advancements within a limited computational budget, training around 10% of the total trainable parameters of the base model (here CLAP) and utilizing a single dataset (ESC-50). Unlike prior works, such as those by (Ghosh et al., 2023; Yuan et al., 2024; Wu et al., 2023a), which rely on more expansive datasets and substantial computational resources. Our focus is on developing a methodology that can effectively instill a sense of time in the model within acceptable computational constraints, rather than on generalizing over large, diverse datasets. Our approach detailed in section 3 and illustrated in fig. 1, modifies the contrastive training

108 109	paradigm by introducing a multi-stage hierarchical training process. In the first stage, the model is trained to recognize and differentiate multiple sound events. In the subsequent stage, it learns the
111	temporal relationships between these events, addressing limitations of prior contrastive models that
112	training objective is based on previous works of Oord et al. (2018) on the formulation of InfoNCE
113	loss and Bagad et al. (2023) who explored temporal instillation in video-language models, however
114	we take the research forward and implement a structured, multi-step post-training process tailored
115	to complex temporal tasks in the audio-language domain. Our objective Comparative analysis in
116	section 5 demonstrates the necessity and efficacy of this approach, showing that our two-stage pro-
117	cess outperforms single-stage methods in enabling ALMs to comprehend audio-language modality
118	gaps in sound event distinction and temporal reasoning which has been overlooked in the past.
119	
120	we further critique current zero-shot evaluation methods, which predominantly rely on basic similarity based ratrieval accuracies or employ large language models (LLMs) as evaluators both
121	of which have shown inherent biases and limitations (Gao et al. 2024: Jones & Steinhardt 2023:
122	Stureborg et al., 2024; Wang et al., 2023). Although previous models have been evaluated for their
123	robustness over time (Shocher et al., 2018; Bau et al., 2019; Kundu et al., 2020; Huang et al., 2020;
124	Sun et al., 2020; Liu et al., 2021), these assessments fail to test the models' general language and
125	temporal understanding comprehensively. To bridge this gap, we propose a sequential zero-shot
127	evaluation method that poses increasingly complex tasks, aiming to create a general-purpose evalu- ation framework (details discussed in algorithm 2)
128	ation namework (details discussed in argorithm 2).
129 130	Main contributions. Here are the key contributions of our work, which, to the best of our knowl- edge, are novel and not present in current state-of-the-art models:
131	
132	• Our analysis indicates that current contrastive ALMs face challenges in accurately captur-
133	ing temporal relationships between audio and text, as shown in table 3, highlighting an area
134	for potential improvement in existing models.
135	• We propose a two step post-training within limited compute budget scheme TeminAL:
136	Temporal Instillation in Audio-Language Models for multi-modal contrastive ALMs.
137	Aimed towards developing temporally aware contrastive audio & text encoders which can
138	be employed in various close and open ended generation models as described in section 3.4.
139	• We propose ZSTE: Zero Shot Temporal Evaluation scheme for contrastively trained mod-
140	els. The sequentially complicated evaluation strategy used for evaluating our objectives of
141	temporal instillation section 4.2.
142	
143	2 BACKGROUND AND RELATED WORK
145	2 BACKOROUND AND RELATED WORK
140	

2.1 FOUNDATION MODELS AND MULTI-MODAL TEXT-AUDIO LEARNING

The expansion of Pretrained Foundation Models (PFMs) now includes auditory (Baevski et al., 148 2020), visual (Dosovitskiy et al., 2020), text-image (Ramesh et al., 2021; Radford et al., 2021), and 149 multi-modal data (Lu et al., 2019; Akbari et al., 2021), driving multi-modal integration. Recent 150 work uses audio-visual contrasts for sound localization (Chen et al., 2021; Wu et al., 2022a), cross-151 modal retrieval (Surís et al., 2022), and zero-shot classification (Wu et al., 2022b; Guzhov et al., 152 2022). Audio-text models are gaining traction, including those in the DCASE competition for audio 153 retrieval with language (Xie et al., 2022), and PFMs have been applied in music tagging (Manco 154 et al., 2022), environmental sound identification (Zhao et al., 2021; Lou et al., 2022; Mei et al., 2022; 155 Koepke et al., 2022), and zero-shot tasks (Zhao et al., 2021; Lou et al., 2022; Mei et al., 2022; Koepke 156 et al., 2022; Elizalde et al., 2023). Open-ended models (Kong et al., 2024; Chu et al., 2023; Liu et al., 157 2024; Deshmukh et al., 2023) enable QA capabilities, but our focus is on contrastive learning for 158 audio encoders. The trend is towards integrating language into auditory systems, with applications in text-to-audio (Ghosal et al., 2023; Liu et al., 2023a; Huang et al., 2023), music generation from text 159 (Agostinelli et al., 2023), and source separation (Liu et al., 2023b). Frameworks like CLAP 160 and Compa (Elizalde et al., 2023; Ghosh et al., 2023) unify auditory-linguistic domains, offering 161 strong zero-shot performance in multimodal tasks.

162 2.2 SELF-SUPERVISED LEARNING AND POST-TRAINING

Self-Supervised Learning (SSL) has revolutionized machine learning, especially in NLP and computer vision (He et al., 2020b; Bao et al., 2021). SSL involves training models to predict parts of their input using other parts, leveraging the data's inherent structure for supervision. A prominent SSL method, Contrastive Learning, learns representations by contrasting positive and negative examples, effectively distinguishing similar and dissimilar data samples (Radford et al., 2021; Liang et al., 2022; Tian et al., 2020; Chen et al., 2020; He et al., 2020a). This approach has significantly advanced representation learning, achieving state-of-the-art results across various domains (Chen et al., 2020; He et al., 2020a).

Post-training introduces an additional self-supervised phase to existing models using a limited set of data before downstream task evaluation, reducing the costs of initial large-scale training (Luo et al., 2022; Xue et al., 2022). Luo et al. (2022) employs static mean-pooling, whereas Xue et al. (2022) aligns image captions with video subtitles. In this unsupervised setting, post-training usually fine-tunes few parameters, maintaining the core strengths of the parent model.

- 177
- 178

194 195

196 197

2.3 ZERO-SHOT INFERENCE: LIMITATIONS OF CLASSICAL ZERO-SHOT RETRIEVAL

179 Zero-shot inference enables models to recognize unseen classes without relying on labeled data from each target class, unlike traditional supervised learning (Xian et al., 2018; Wang et al., 2020b). 181 While zero-shot learning facilitates generalization to unseen classes, conventional audio-retrieval 182 benchmarks often lack compositional complexity, typically involving single acoustic events without 183 proper word order (Radford et al., 2021; Baevski et al., 2020; Gemmeke et al., 2017). In traditional audio classification, models are trained on specific classes like musical genres or environmental sounds, but zero-shot audio classification requires identifying audio samples from previously un-185 seen classes. For example, a model trained on animal and vehicle sounds should also classify new categories like "machinery" or "insects" (Wang et al., 2020a). As illustrated in fig. 8, zero-shot clas-187 sification involves encoding audio and text prompts through respective encoders and using cosine 188 similarity to predict classes (Harwath & Glass, 2015; Kim & Pardo, 2018). Zero-shot audio retrieval 189 extends this concept by finding relevant audio clips from unseen classes based on queries, such as 190 retrieving "birdsong" or "ocean waves" when trained only on spoken words and ambient sounds 191 (Fonseca et al., 2021). As shown in fig. 9, the process involves encoding prompts and audio clips, 192 with cosine similarity determining the most relevant match (Chang & Yang, 2019). This approach 193 leverages class information to understand semantic relationships.

3 Methodology

3.1 PRELIMINARIES



Introduction to Fundamentals. Consider set A as the domain of audio recordings and $\mathbb C$ as the set of corresponding textual transcripts (contexts). For any two discrete and nonoverlapping audio clips $\{a_i, a_j\}$ within \mathbb{A} , let their relevant transcripts be $\{c_i, c_j\}$ in \mathbb{C} (We use 'c' for transcripts to avoid confusion with the time variable 't'). We define an integrated segment that respects the sequential order as (a_{ij}, c_{ij}) , with a_{ij} constructed by the operation $[a_i \oplus a_j]$, which concatenates the two audio clips as marked by the operator \oplus which shows the concatenation operation also shown in fig. 2. Similarly for contexts, we first introduce $\tau = \{\tau_t, \tau_o\}$ to represent a sequential relationship, where τ_t can either be preceding or

succeeding as prompted by {before or after} and we define τ_o for overlapping language prompt {while}. Following which c_{ij} is represented as $[c_i : \tau_t; c_j]$, merging the transcripts in a manner that it reflects the temporal relation $\tau = \{\tau_t, \tau_o\}$. Later in section 3.2, we relate $[a_i \oplus a_j]$ with $[a_j \oplus a_i]$ using mathematical operators. It should be noted that the arrangement of a_i and a_j within a_{ij} may vary depending on the value of τ_t . The same is applicable for overlapping sounds (a_j, a_j) , for which the overlapping texts can be represented by $[c_i : \tau_o; c_j]$, which essentially means " c_i while c_j " with overlaid audios $[a_i \wedge a_j]$. For simplicity, we will refer to the composite audio-text pair (a_{ij}, c_{ij}) as (a, c), except where additional specificity is required.

221 222

223

3.2 DATA-PROCESSING: DESIGNING OUR TRAINING DATA.

The dataset for our post-training study was meticulously curated from publicly available audio-text 224 pairs, we specifically select the ESC-50 dataset for the current study. The dataset selection and pro-225 cessing is descried in detail in appendix B.2. We introduce a temporal inversion operator 'T' and 226 temporal overlay operator ' \mathbb{O} ' to represent the transformation of audio and text training data to form 227 the temporally inverted samples and temporally overlapped samples as shown in equation 1 for the 228 temporal inversion and equation 2 for temporally overplayed samples. This function is designed to 229 operate on pairs of simultaneous audios (a_i, a_j) or transcription sequences (c_i, c_j) where sequences 230 in both these sets are initially non-overlapping. We show temporal addition/ concatenation of the 231 pair of audios by $a_i \oplus a_i$ and overlaying of the audio pair by $a_i \wedge a_i$. Meanwhile temporal addition and overlaying of texts are shown as $c_j; \tau_t; c_i$ and $c_j; \tau_o; c_i$ respectively and follows the same 232 233 convention as mentioned in section 3.1.

- 234
- 235 236

$$\mathbb{T}(a) = \mathbb{T}([a_i; a_j]) := [a_j \oplus a_i], \quad \mathbb{T}(c) = \mathbb{T}([c_i; c_j]) := [c_j; \tau_t; c_i]$$
(1)

$$\mathbb{O}(a) = \mathbb{O}([a_i; a_j]) := [a_j \wedge a_i], \quad \mathbb{O}(c) = \mathbb{O}([c_i; c_j]) := [c_j; \tau_o; c_i]$$

$$(2)$$

238 It is essential to recognize that 'T' does not literally reverse time within the audio tracks, rather it 239 rearranges the sequence of events within the compiled segments. Our goal is to cultivate a model 240 capable of distinguishing an original audio-text pair (a, c) from both of its temporally inverted coun-241 terpart $(a, \mathbb{T}(c))$, and also $(\mathbb{T}(a), c)$; furthermore to contrast all of these from the overlaid text-242 audio pair as $(\mathbb{O}(a), c)$ (which is the same as $(a, \mathbb{O}(c))$). So a typical training batch would look like $B_{a_B} = \{a, \mathbb{T}(a), \mathbb{O}(a)\}$ for the audio and $B_{t_B} = \{c, \mathbb{T}(c), \mathbb{O}(c)\}$ for the text. The details 243 for out data-preparation method is described in algorithm 3. As described earlier in section 1 we 244 245 have a hierarchical 2-stage training process TeminAL A followed by TeminAL B. The text-audio dataset $\{B_{a_B}, B_{t_B}\}$ is used to train TeminAL B. While the first pretraining TeminAL A, works on 246 learn single sounds and multiple sounds thus the input data in the batch doesn't consists of time-247 reversed data, it's made up of $B_{a_A} = \{a_i, a_i \oplus a_j \forall i, j \in \{1, N\}$ the audio and $B_{t_A} = \{c_i, c_i \oplus c_j\}$ 248 $\forall i, j \in \{1, N\}$ for the text. 249

249 250 251

3.3 PRELIMINARIES OF POST-TRAINING WITH SSL

The input texts and audios are first transformed into machine-level embeddings. Let the processed embedding for audio be x_a where $x_a \in \mathbb{R}^{F \times T}$, with F representing frequency components (e.g., Mel frequency bins) and T indicating the number of temporal segments. The corresponding textual data is denoted as x_c for a given sample. For a batch of N text-audio pairs, the audio and corresponding text are represented as $\{X_a, X_c\}_i = \{x_a^{(i)}, x_c^{(i)}\}$ for $i = 1, \ldots, N$. For simplicity, we denote the entire collection of N pairs as $\{X_a, X_c\}_i$. Each audio segment and its corresponding text description are processed through separate encoders: $f_a(.)$ for audio and $f_c(.)$ for text. For a batch of size N, we have:

260 261 262

$$m{z}_{a}^{(i)} = f_{a}(m{x}_{a}^{(i)}) \in \mathbb{R}^{d}, \quad m{z}_{c}^{(i)} = f_{c}(m{x}_{c}^{(i)}) \in \mathbb{R}^{d}, \quad i = 1, \dots, N$$

where $z_a^{(i)}$ and $z_c^{(i)}$ represent the audio and text encodings, respectively. To evaluate the similarity between embeddings $z_a^{(i)}$ and $z_c^{(i)}$, we calculate their similarity matrix as $C = \gamma \cdot (z_c z_a^{\top})$. Here, τ is a scaling constant that adjusts the logarithmic scale after applying softmax, as detailed in part D. The similarity matrix C is $\mathbb{R}^{N \times N}$, with N compatible pairs along the diagonal and $N^2 - N$ non-compatible pairs elsewhere. The overall objective function is defined as $\mathcal{L} = 0.5 \cdot (\ell_{text}(C) + \ell_{audio}(C))$.where $\ell_{text}(C)$ and $\ell_{audio}(C)$ are computed separately for the text and audio embeddings, using cross-entropy loss. Specifically, $\ell_{text}(C) = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{e^{(z_c^{(i)} \cdot z_a^{(i)}/\gamma)}}{\sum_{i=1}^{N} e^{(z_c^{(i)} \cdot z_a^{(i)}/\gamma)}}$. This promotes joint optimization of the audio and text encoders along with their respective transformations, as described in later sections.



Figure 3: The overview of TeminAL B where we are post-training orginal CLAP encoders f_c and f_a with our **TeminAL** method to get f_c^t and f_a^t . The functions as described in section 3.3, while the objective formulation for training (f_c, f_a) to achieve (f_c^t, f_a^t) has been described in section 3.4. The "Temporal contrastive loss" for TeminAL B has been elaborated in fig. 4.

3.4 OBJECTIVE FUNCTION FOR TEMINAL: WHAT ADDITION WE PROPOSE ON CLASSICAL CONTRASTIVE LEARNING

296 We propose a multi-stage training approach, outlined in fig. 1, with two stages: TeminAL A and 297 TeminAL B. In TeminAL A, the model is trained to distinguish between single and multiple sounds, 298 while in TeminAL B, the model learns to differentiate temporally distinct sounds along with their 299 corresponding text labels. Both stages use contrastive learning with a modified infoNCE loss func-300 tion (Oord et al., 2018), detailed further in this section and appendix B.3, the difference in training 301 being the training data and contrastive objective. We have already elaborated on the different train-302 ing dataset and it's formulation in section 3.2, while the detailed loss function for both stages has 303 been discussed in this section. Context (text) and audio encodings are processed through their re-304 spective encoders, producing embeddings C_e and A_e as shown in fig. 3. These embeddings are used to form a (batch \times batch) matrix to identify positive and negative pairs (see Figures 4). Similar-305 ity scores are calculated from these embeddings and used to compute the modified infoNCE loss 306 function, as described latter in the section. Logits derived from similarity scores are transformed 307 using a Softmax function to generate probabilities (equation 4 and equation 5), which are evaluated 308 with cross-entropy against true labels. The loss function is computed as the sum of text loss (L_c) 309 and audio loss (L_a) , which sum up to form $L_B = L_{c_B} + \beta(L_{a_B})$ which stands for the TeminAL 310 B loss. The text loss L_{c_B} optimizes the selection of texts from n possible options generated by 311 C_e (equation 3), while audio loss L_{a_B} does the same for the audio embeddings (equation 7). This 312 dual-component loss ensures balanced training of both context (C_e) and audio (A_e) encoders. The 313 overall methodology is schematically depicted in Figure 3. 314

Unlike traditional contrastive loss functions that primarily reinforce true positives, our approach 315 modifies the infoNCE loss to make encoders more sensitive to time-reversed and overlapping sam-316 ples, as shown in equation 4 and equation 5. For temporal alignment, we use an adapted version 317 of the InfoNCE loss function in both TeminAL A and B to distinguish the temporal sequence of 318 audio-text pairs. For a time-aligned audio-text pair (a, c), following section 3.1, we design a loss 319 function that maintains chronological order within the pair, differing only in the loss components. 320 The training batch for TeminAL A is defined as $B_{c_A} = \{B_{c_s}, B_{c_d}\}$ for texts (single and dual stitched audio) and $B_{a_A} = \{B_{a_s}, B_{a_d}\}$ for audio, following the conventions in section 3.2. For TeminAL B, 321 the batches are $B_{c_B} = \{B_{c_f}, B_{c_r}, B_{c_o}\}$ for texts and $B_{a_B} = \{B_{a_f}, B_{a_r}, B_{a_o}\}$ for audio (forward, 322 reversed, and time-overlaid); In general we represent batches of audio-text data by symbol B. These 323 batches are processed through encoders, converting them into audio and text embeddings that are

274

275 276

277

278

279

281

284

287

288

289

290

291 292 293

294

used in subsequent stages of training. For the layout of the batches of data kindly refer to fig. 5. Furthermore, our encoders are not trained from scratch; we extend our framework using a pre-existing audio-language model comprising an audio encoder $f_{c_{\theta}}$ and a text encoder $f_{a_{\phi}}$ shown in fig. 3 from CLAP by Elizalde et al. (2023). These pre-trained encoders are post-trained to enhance temporal accuracy while maintaining baseline retrieval performance, as demonstrated in table 1. Due to limited dataset size, selective refinement of specific layers within $\Theta = \{\theta, \phi\}$ is performed, as schematically shown in fig. 1 and detailed in appendix B.3.

331 332 333

334

347

348

$$L_{c_B} = \sum_{(a,c)\in B_{c_B}} (\text{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_c) + \text{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_{\mathbb{T}(c)}) + \text{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_{\mathbb{O}(c)}))$$
(3)

To complete our model construct, in the following section we have explained the details of the loss function mathematically in equations equation 3–equation 9. Kindly note that, the hyperparameters introcuded are discussed in the following section 3.5. Earlier, we had seen the discussion on text and audio losses (L_c and L_a), we now define them mathematically in the following equations. Here, TNCE stands for Temporal Noise Contrastive Estimation, a variant of the NCE loss.

$$\text{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_c) = -\log \frac{\exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_c)}{\sum_{c' \in B_{c_s}} \exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_{c'}) + C^{c_r} + C^{c_o}}$$
(4)

$$\mathsf{TNCE}(\boldsymbol{z}_{a}, \boldsymbol{z}_{\mathbb{O}(c)}) = -\log \frac{\exp(\boldsymbol{z}_{a} \cdot \boldsymbol{z}_{\mathbb{O}(c)})}{\sum_{c' \in B_{c_{a}}} \exp(\boldsymbol{z}_{a} \cdot \boldsymbol{z}_{\mathbb{O}(c')}) + C^{c_{c}}}$$
(5)

In equation 4, C^{c_r} and C^{c_o} is an accumulation of negatives fashioned via time-reversal and time-overlay respectively in equation 4, and is expressed as: $C^{c_r} = \alpha_{s_t} \exp(\mathbf{z}_a \cdot \mathbf{z}_{\mathbb{T}(c)}) + \alpha_{c_t} \sum_{c' \in B_{c_r} \setminus \{c\}} \exp(\mathbf{z}_a \cdot \mathbf{z}_{\mathbb{T}(c')})$ and $C^{c_o} = \alpha_{s_o} (\exp(\mathbf{z}_a \cdot \mathbf{z}_{\mathbb{O}(c)})) + \alpha_{t_o} \sum_{c' \in B_c \setminus \{o\}} \exp(\mathbf{z}_a \cdot \mathbf{z}_{\mathbb{O}(c')})$. While C^{c_c} from equation 5 is expressed as:

$$C^{c_c} = \exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_c) + \sum_{c' \in B_{c_f} \setminus \{c\}} \exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_{c'}) + \alpha_{s_t} \exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_{\mathbb{T}(c)}) + \alpha_{c_t} \sum_{c' \in B_{c_r} \setminus \{c\}} \exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_{\mathbb{T}(c')})$$
(6)

The loss function is constructed in such a way that it penalises the miss-classifications among the audio-text pairs. The loss formulations gives a handle on penalising the time-reversed samples and time-overlayed samples with the hyper-parameters α_{s_t} and α_{s_o} , we present a detailed analysis on effects of these hyper-parameters later in section 3.5. The total loss L_B for TeminAL B can then be written with L_{c_B} and L_{a_B} which follows the same formulation as L_{c_B} . Detailed formulation for L_{c_B} and L_{a_B} have been provided in the supplementary section. The net loss for TeminAL B is expressed as $L_B = L_{c_B} + \beta(L_{a_B})$, where L_{a_B} is as follows:

363 364

365

$$L_{a_B} = \sum_{(a,c)\in B_{a_B}} (\text{TNCE}(z_c, z_a) + \text{TNCE}(z_{\mathbb{T}(c)}, z_a) + \text{TNCE}(z_{\mathbb{O}(c)}, z_a))$$
(7)

After discussing the loss formulation of TeminAL B, we have similar formulation for TeminAL A. With necessary changes in the configuration of data within the batch $(B_{a_A} \text{ and } B_{c_A})$ as it's mentioned in the previous paragraph, kindly refer to fig. 6 for the layout of the batch. The mathematical formulation of the contrastive loss function is described as follows and schematically shown in fig. 7:

$$L_{c_A} = \sum_{(\mathbb{T}(a),\mathbb{T}(c))\in B_{c_A}} (\text{TNCE}(\boldsymbol{z}_{c_s}, \boldsymbol{z}_{a_s}) + \text{TNCE}(\boldsymbol{z}_{c_d}, \boldsymbol{z}_{a_d}))$$
(8)

372 373 374

375 376

370

$$L_{a_A} = \sum_{(\mathbb{O}(a),\mathbb{O}(c))\in B_{a_A}} (\text{TNCE}(\boldsymbol{z}_{a_s}, \boldsymbol{z}_{c_s}) + \text{TNCE}(\boldsymbol{z}_{a_d}, \boldsymbol{z}_{c_d}))$$
(9)

The loss function construction and mathematical derivation of $L_A = L_{c_A} + \beta_A(L_{a_A})$ for TeminAL A is also detailed in appendix B.3.



Figure 4: The schematic showing Temporal Contrastive Loss for TeminAL B. On the vertical axis we have the audio embeddings with batches of data corresponding to $B_{a_B} = \{B_{a_f}, B_{a_r}, B_{a_o}\}$ and text embedding batches of data corresponding to $B_{c_B} = \{B_{c_f}, B_{c_r}, B_{c_o}\}$ on the horizontal axis.

3.5 DETAILS ON HYPER-PARAMETERS OF THE LOSS FORMULATION

The loss function formulation introduces hyper-parameters that affect the temporal sensitivity of the encoders. The choice of $\{\alpha_{s_o}, \alpha_{c_o}, \alpha_{s_t}, \alpha_{c_t}\}$ and $\{\beta_A, \beta\}$, which can be either 0 or 1, significantly impacts the model's behavior. For example, setting $\alpha_{s_t} = 1$ increases sensitivity to time-reversed sound samples by adding the term $\alpha_{s_t} \exp(\mathbf{z}_a \cdot \mathbf{z}_{\mathbb{T}(c)})$ to the expression C_{tr} denominator of equation 6. This adjustment forces the encoders to ensure the sum of terms equals unity, which reduces the encoding values of non-similar pairs, enhancing sensitivity to time-reversed samples and guiding the encoders to assign lower values to dissimilar batch samples.

These coefficients also help adapt the model to different datasets. In fig. 4, these parameters extend the contrastive loss function over time: the top three sub-squares represent $\text{TNCE}(z_a, z_c)$, the middle sub-squares represent $\text{TNCE}(z_a, z_{\mathbb{T}(c)})$, and the last three sub-squares represent $\text{TNCE}(z_a, z_{\mathbb{O}(c)})$. The top-left quadrant shows contrastive loss on stitched pairs, with positive (green) and negative (red) diagonal terms crucial for temporal understanding.

The key role of β is to increase the number of training samples, while the α coefficients enhance sensitivity to time-reversal and overlapping sounds. When $\alpha = 0$, sensitivity is nullified, but higher values compel the encoders to refine recognition of temporal variations by minimizing the denominator. Without these hyper-parameters, the loss converges similarly, but encoders learn different relationships. Their inclusion ensures distinct sample treatment, enhancing temporal sensitivity.

424 425

426

428

400

401

402 403 404

405 406

4 EXPERIMENTS

427 4.1 BASE MODEL

We employ a pre-trained CLAP model Elizalde et al. (2023) using transformer-based encoders:
HTSAT Chen et al. (2022) for audio and BERT Devlin et al. (2018) for text, each with a projection layer. We focuses on the final layers and projection layers of both encoders for our training. While our TeminAL training approach is model-agnostic, we start with the CLAP model as a foundation.

432 4.2 **ZSTE** AND DOWNSTREAM TASKS

434 We construct comprehensive evaluation of our proposed method in order to satisfy our objectives 435 of temporal instillation. ZSTE (Zero Shot Temporal Evaluation) is our evaluation framework for assessing contrastive models in zero-shot tasks. The implementation is discussed in algorithm 2 and 436 algorithm 3. ZSTE begins with basic classification on unseen classes in Task 1 (fig. 10), progressing 437 to complex scenarios involving overlapping audio features and novel composite text classes in Task 438 2. The subtasks 2A and 2B involve the model being able to distinguish both classes and at-least 439 one class respectively, model which correctly understands both sounds performs well on subtask 440 A. The subtask 2C and 2D are similar tasks as 2A and 2B but on overlapping sounds rather than 441 concatenated sounds. Thus given the overlapping text-audio pair, we need compute the accuracy 442 of the model to detect both the audio classes in 2C and at-least one of the classes in 2D (fig. 11). 443 Models which are able to distinguish multiple overlapping sounds, perform better in 2C. 444

ZSTE Task 3 evaluates temporal comprehension by testing sequences of interchanged acoustic 445 events, building on the configuration from Task 2 but now using temporal texts (fig. 12). A model 446 capable of understanding temporal relationships in text will excel in this task. Task 4 (fig. 13) as-447 sesses the model's ability to maintain focus amid irrelevant class labels, which act as noise to the 448 actual audio embeddings. Task 5 (fig. 14) examines the model's generalization to out-of-distribution 449 prompts, reflecting real-world complexities. Each of these three tasks includes subtasks A and B: 450 Subtask A requires detecting all audio classes, while Subtask B involves identifying at least one 451 class. These tasks test the model's grasp of temporality and general language attribution. Models 452 with a robust understanding of both temporality and language generalisability will perform better.

This comprehensive approach ensures robust evaluation of the model's zero-shot learning capabil The primary aim is to foster model improvement rather than solely benchmark performance.
 Further details on ZSTE is shown in appendix B.4.

456 457 458

5 RESULTS

459 In this section, we present the experimental results to support the claims outlined in our objectives 460 and discuss our key findings. The results are organized around several downstream tasks, beginning 461 with audio and text retrieval and progressing to an in-depth evaluation of the models' temporal be-462 havior and finally comparing various SOTA contrastive ALM models for temporal understanding. 463 Firstly we compare the retrieval performance of closed-ended and open-ended models on benchmark 464 AudioCap and Clotho dataset, as summarized in Tables 1 and 7. Our model, T-CLAP, demonstrates 465 superior performance across retrieval tasks, surpassing most existing models in both closed-ended and open-ended categories. Notably, it achieves competetive state-of-the-art results for both text-466 to-audio (T-A) and audio-to-text (A-T) retrieval tasks. These results underscore the effectiveness 467 of our contrastive training strategy, employed both during the pre-training phase of CLAP and the 468 subsequent fine-tuning phase using TeminAL. Our approach effectively preserves the contrastive 469 knowledge acquired during pre-training, ensuring strong retrieval performance. However, as men-470 tioned previously in section 1, the retrieval metrics alone do not fully encapsulate the temporal 471 understanding capabilities of our model. To address this limitation, we conduct a rigorous Zero-472 Shot Temporal Evaluation (ZSTE), which offers deeper insights into the temporal reasoning ability 473 of T-CLAP. The method of evaluation is further elaborated in appendix B.4.

474 Firstly we try and study the model's behaviour through the hyperparameter variations, the role of 475 each hypermeter is detailed in section 3.5. The results from Table 2 convey important information 476 on how the model captures temporal behaviour in general (across the ZSTE tasks) through our mod-477 ification of the overall objective function through parametric variations and the impact of including 478 a multistage training objective. As mentioned in section 4.2 in Task 1, the model must excel in the 479 initial pre-training task, and results indicate that our training strategies prevent catastrophic forget-480 ting, although increasing temporality in the objective function we observe decrease in performance 481 in this task we still remain well above across different tasks. Simillarly, task 2 tests multi-sound 482 understanding, where the two-stage TeminAL AB training significantly improves sound distinction capabilities. Task 3 focuses on temporal reasoning, demonstrating that specific loss coefficients en-483 hance the model's ability to capture temporal relationships. Task 4 and 5 evaluates complex and 484 general text prompts, showing that our model outperforms the original CLAP in correctly mapping 485 stitched audio to appropriate text although much room for improvement is there for future models.

Model		T-A Retrieva	1	Α	A-T Retrieval	
	R@1	R@5	R@10	R@1	R@5	R@10
MMT	36.1 / 6.7	72.0/21.6	84.5 / 33.2	39.6 / 7.0	76.8 / 22.7	86.7 / 34.6
ML-ACT	33.9 / 14.4	69.7 / 36.6	82.6 / 49.9	39.4 / 16.2	72.0/37.6	83.9 / 50.2
CLAP	34.6 / 16.7	70.2 / 41.1	82.0 / 54.1	41.9 / 20.0	73.1 / 44.9	84.6 / 58.7
CLAP-LAION	36.2 / 17.2	70.3 / 42.9	82.5 / 55.4	45.0 / 24.2	76.7 / 51.1	88.0 / 66.9
CompA-CLAP	36.1 / 16.8	78.6 / 43.5	90.2 / 56.1	47.8 / 23.9	83.5 / 50.7	90.2 / 67.6
T-CLAP(ours)	35.1 / 17.0	71.2/42.2	82.1 / 54.7	49.2 / 23.1	85.1 / 52.2	87.8 / 66.4

Table 1: Comparison of models on text-to-audio (T-A) and audio-to-text (A-T) retrieval tasks. Performance of Text-to-Audio and Audio-to-Text retrieval on AudioCap/Clotho dataset. The values for other models have been taken from previous publications (Ghosh et al., 2023; Elizalde et al., 2023).

		Loss-	coeffic	ients							ZSTE					
	0.	0.	Ω.	0.	β	Task 1		Tas	sk 2		Tas	sk 3	Tas	sk 4	Tas	k 5
	αs_t	αc_t	cas _o	CC0	μ	А	А	В	С	D	А	В	А	В	А	В
	0	0	0	0	1	78.77	10.11	77.60	8.46	78.01	32.67	31.57	3.50	1.12	26.01	1.10
	1	0	1	0	1	77.67	11.02	83.20	8.71	81.21	48.50	18.20	36.28	7.10	27.07	12.7
<u> </u>	0	1	0	1	1	76.54	10.11	83.44	7.98	83.01	49.11	18.74	40.38	8.32	28.01	15.2
Γ	1	1	1	1	1	76.14	12.20	83.61	11.44	83.2	51.3	22.83	41.10	9.18	31.21	15.8
	0	0	0	0	1	77.34	38.45	80.90	49.87	79.67	34.22	33.12	34.50	32.15	27.34	2.01
8	1	0	1	0	1	76.76	39.23	86.34	50.65	83.78	52.12	42.45	39.45	38.67	28.45	14.23
-	0	1	0	1	1	75.29	43.45	85.89	59.56	84.56	50.78	24.33	52.78	41.45	29.78	16.89
Ę	1	1	1	1	1	75.11	46.78	86.12	62.34	85.45	54.56	56.78	46.23	44.89	32.45	18.34

Table 2: Hyper-parameter analysis for loss coefficients $\{\alpha, \beta\} = \{\alpha_{s_t}, \alpha_{c_t}, \alpha_{s_o}, \alpha_{c_o}, \beta\}$. Each Task is defined according to appendix B.4, kindly refer this section for details on each task. Here T-B and **T-AB** refers to models trained with only TeminAL B and TeminAL A + B respectively.

Our findings show that setting $\alpha_{c_t} = 0$ benefits ZS-tasks with fewer confusing classes, like Task 1, while $\alpha_{c_t} = 1$ improves performance in tasks with more complex classes (Tasks 4, and 5). Models trained with $\alpha_{s_0} = 0$ lack overlaid classes in the denominator, impacting how negatives are penal-ized. Table 2 illustrates our model's adaptability across ZSTE tasks, with key improvements noted when using combined TeminAL A and B training. Adjusting α_{c_t} and α_{c_o} makes the model more sensitive to time-reversed samples, enhancing performance in time-sensitive tasks. The T-CLAP model sometimes struggles with sound distinction due to training focused on temporal understand-ing, not sound separation, affecting sensitivity and overall accuracy. However, hierarchical training with both TeminAL A and B significantly improves sound distinction and general language understanding tasks. This result grounds the importance of our multi-stage training method in order to learn the temporal behavior of multiple sound as described in section 1.

Tasks	Subtasks	ML-ACT	CLAP	CLAP-LAION	CompA-CLAP	Т-0	CLAP
						TeminAL B	TeminaAL AB
1	Α	76.12	81.22	82.5	83.0	76.14	75.11
2	А	7.2	9.59	10.1	18.4	32.20	46.78
2	В	78.1	81.00	81.3	91.6	83.61	87.12
2	С	6.5	9.39	10.0	21.3	31.4	62.34
2	D	71.7	80	80.4	90.8	83.2	85.45
3	Α	28.01	33.27	34.93	54.5	51.3	54.56
3	В	27.5	34.29	34.6	49.87	22.83	56.78
4	Α	2.2	2.4	7.56	48.71	41.1	46.23
4	В	2.0	1.98	5.45	38.74	9.18	44.89
5	А	3.0	26	26.4	36.81	31	32.45
5	В	2.5	0	0.7	18.2	15.8	18.34

Table 3: Showing the comparison of various contrastive learning models on our ZSTE tasks. The details on each task is discussed in section 4.2 and further detailed in appendix B.4.

540 We next evaluate and compare the performance of various models on temporal understanding 541 through the ZSTE tasks. As shown in Table 3, T-CLAP outperforms most state-of-the-art mod-542 els across a majority of tasks. Notably, T-CLAP excels in tasks 2A and 2C, as mentioned previously 543 section 4.2, these two substasks invovle the model to distinguish multiple sounds and detecting both 544 the sounds. While remaining competitive in Task 1A we observe, demonstrating that our temporal instillation approach effectively instills the model with a sense of time without significantly degrad-545 ing its performance on the original pre-training tasks. Furthermore, results of Tasks 2B and 2C, 546 which require the model to associate multiple sounds with at least one correct class, show better 547 performance with larger models due to their capacity to handle complex associations. Following 548 the results of Task 3, we observe T-CLAP performs competitively with other models, despite those 549 models being specifically trained on audio-text pairs. Interestingly, these competing models achieve 550 strong results on Task 3 without performing as well in tasks requiring the differentiation of mul-551 tiple audio sounds, such as Tasks 2B and 2C. However, all models, including T-CLAP, encounter 552 challenges in general language understanding tasks, such as Tasks 4 and 5. This suggests that lever-553 aging a more robust pre-trained language encoder along with diversifying the dataset could further 554 enhance overall performance.

555 556

6 CONCLUSION

557 558

This research introduces the Temporal Instillation in Audio-Language Models, a post-training tech-559 nique that enhances temporal and language understanding in Audio-Language Models (ALMs). Our 560 approach employs sequential inversion and temporal augmentations, effectively improving sequen-561 tial discernment in ALMs. The hierarchical training strategy proves crucial, as seen in the perfor-562 mance comparison between TeminAL B and TeminAL AB, highlighting the need for structured 563 training in complex tasks like time instillation. Our findings also demonstrate that modifying the infoNCE loss improves model sensitivity, as shown in our parametric study (table 2). Zero Shot 565 Temporal Evaluation (ZSTE) results (section 5) confirm T-CLAP's strength in zero-shot classifi-566 cation and retrieval tasks, offering new evaluation insights for contrastive learning models. While 567 T-CLAP shows a slight decrease in traditional audio classification accuracy, it consistently outperforms baseline models in scenarios involving temporal relationships, demonstrating enhanced se-568 quential information processing. This study opens new research directions, particularly in refining 569 contrastive loss for broader task optimization, paving the way for ALMs that excel in both retrieval 570 and complex temporal-linguistic tasks. 571

572 573

574 575

576 577

578

579 580

584

585

586

587

588

7 References

References

- Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. Musiclm: Generating music from text. arXiv preprint arXiv:2301.11325, 2023.
- Hassan Akbari, Liangzhe Yuan, Rui Qian, Wei-Hong Chuang, Shih-Fu Chang, Yin Cui, and Boqing
 Gong. Vatt: Transformers for multimodal self-supervised learning from raw video, audio and
 text. Advances in Neural Information Processing Systems, 34:24206–24221, 2021.
 - Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. In Advances in Neural Information Processing Systems (NeurIPS 2020), 2020. URL https://arxiv.org/abs/2006. 11477.
- Piyush Bagad, Makarand Tapaswi, and Cees GM Snoek. Test of time: Instilling video-language
 models with a sense of time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2503–2516, 2023.
- 592
 - Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. arXiv preprint arXiv:2106.08254, 2021.

631

632

633 634

635

636

637

638

639

594	David Bau, Hendrik Strobelt, William S. Peebles, Jonas Wulff, Bolei Zhou, Jun-Yan Zhu, and
595	Antonio Torralba. Semantic photo manipulation with a generative image prior. ACM Trans-
596	actions on Graphics, 38(4), 2019. ISSN 0730-0301. doi: 10.1145/3306346.3322994. URL
597	https://doi.org/10.1145/3306346.3322994.

- H. Chang and Z. Yang. Zero-shot learning for audio-visual speech recognition. In *Proceedings of the 2019 IEEE International Conference on Multimedia and Expo (ICME)*, pp. 1341–1346. IEEE, 2019.
- Honglie Chen, Weidi Xie, Triantafyllos Afouras, Arsha Nagrani, Andrea Vedaldi, and Andrew Zisserman. Localizing visual sounds the hard way. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16867–16876, 2021.
- Ke Chen, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Hts A hierarchical token-semantic audio transformer for sound classification and detection. In
 ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 646–650. IEEE, 2022.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *Proceedings of the 37th International Conference on Machine Learning*, 119:1597–1607, 2020.
- Feng Cheng, Xizi Wang, Jie Lei, David Crandall, Mohit Bansal, and Gedas Bertasius. Vindlu: A
 recipe for effective video-and-language pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10739–10750, 2023.
- Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.
- Soham Deshmukh, Benjamin Elizalde, Rita Singh, and Huaming Wang. Pengi: An audio language
 model for audio tasks. *arXiv preprint arXiv:2305.11834*, 2023.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Alexey Dosovitskiy et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR 2021)*, 2020. URL https://arxiv.org/abs/2010.11929.
- Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning audio concepts from natural language supervision. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023.
 - Han Fang, Pengfei Xiong, Luhui Xu, and Yu Chen. Clip2video: Mastering video-text retrieval via image clip. *arXiv preprint arXiv:2106.11097*, 2021.
 - Nanyi Fei, Zhiwu Lu, Yizhao Gao, Guoxing Yang, Yuqi Huo, Jingyuan Wen, Haoyu Lu, Ruihua Song, Xin Gao, Tao Xiang, et al. Towards artificial general intelligence via a multimodal foundation model. *Nature Communications*, 13(1):3094, 2022.
 - E. Fonseca, J. Pons, and X. Serra. Unsupervised learning for large-scale zero-shot audio classification. In *Proceedings of the 2021 Conference of the International Speech Communication Association*, 2021.
- Mingqi Gao, Xinyu Hu, Jie Ruan, Xiao Pu, and Xiaojun Wan. Llm-based nlg evaluation: Current status and challenges. *arXiv preprint arXiv:2402.01383*, 2024.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*, 2021.
- Yuying Ge, Yixiao Ge, Xihui Liu, Dian Li, Ying Shan, Xiaohu Qie, and Ping Luo. Bridging video text retrieval with multiple choice questions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16167–16176, 2022.

648 Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing 649 Moore, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. 650 In 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 651 pp. 776–780. IEEE, 2017. 652 Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio gen-653 eration using instruction-tuned llm and latent diffusion model. arXiv preprint arXiv:2304.13731, 654 2023. 655 656 Sreyan Ghosh, Ashish Seth, Sonal Kumar, Utkarsh Tyagi, Chandra Kiran Evuru, S Ramaneswaran, 657 S Sakshi, Oriol Nieto, Ramani Duraiswami, and Dinesh Manocha. Compa: Addressing the gap in compositional reasoning in audio-language models. arXiv preprint arXiv:2310.08753, 2023. 658 659 Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Audioclip: Extending clip to 660 image, text and audio. In ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech 661 and Signal Processing (ICASSP), pp. 976–980. IEEE, 2022. 662 D. Harwath and J. Glass. Deep multimodal semantic embeddings for speech and images. In 2015 663 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU), pp. 237–244. 664 IEEE, 2015. 665 666 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsu-667 pervised visual representation learning. Proceedings of the IEEE/CVF Conference on Computer 668 Vision and Pattern Recognition, pp. 9729–9738, 2020a. 669 Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert 670 with disentangled attention. arXiv preprint arXiv:2006.03654, 2020b. 671 672 Qingqing Huang, Aren Jansen, Joonseok Lee, Ravi Ganti, Judith Yue Li, and Daniel PW Ellis. Mu-673 lan: A joint embedding of music audio and natural language. arXiv preprint arXiv:2208.12415, 674 2022. 675 Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin 676 Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced 677 diffusion models. arXiv preprint arXiv:2301.12661, 2023. 678 679 Yu Huang, Chenzhuang Du, Zihui Xue, Xuanyao Chen, Hang Zhao, and Longbo Huang. What 680 makes multi-modal learning better than single (provably). Advances in Neural Information Pro-681 cessing Systems, 34:10944-10956, 2021. 682 Yujia Huang, James Gornet, Sihui Dai, Zhiding Yu, Tan M. Nguyen, Doris Y. Tsao, and Anima 683 Anandkumar. Neural networks with recurrent generative feedback. In Advances in Neural Infor-684 mation Processing Systems (NeurIPS), 2020. 685 686 Sarah Jones and Jacob Steinhardt. Benchmarking cognitive biases in large language models as eval-687 uators. arXiv preprint arXiv:2309.17012, 2023. URL https://arxiv.org/abs/2309. 17012. 688 689 Y. Kim and B. Pardo. Zero-shot audio classification with transfer learning. In Proceedings of the 690 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 2261– 691 2265. IEEE, 2018. 692 693 A Sophia Koepke, Andreea-Maria Oncescu, Joao Henriques, Zeynep Akata, and Samuel Albanie. Audio retrieval with natural language queries: A benchmark study. IEEE Transactions on Multi-694 media, 2022. 696 Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio 697 flamingo: A novel audio language model with few-shot learning and dialogue abilities. arXiv 698 preprint arXiv:2402.01831, 2024. 699 Jogendra Nath Kundu, Naveen Venkat, and R. Venkatesh Babu. Universal source-free domain adap-700 tation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition 701 (CVPR), June 2020.

102	Junnan Li, Dongxu Li, Caiming Xiong, and Steven CH Hoi. Blip: Bootstrapped language-
703	image pre-training for unified vision-language understanding and generation. arXiv preprint
704	arXiv:2201.12086, 2022.
705	

- Victor Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Y Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. *Advances in Neural Information Processing Systems*, 35:17612–17625, 2022.
- Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and
 Mark D Plumbley. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503*, 2023a.
- Shansong Liu, Atin Sakkeer Hussain, Chenshuo Sun, and Ying Shan. Music understanding llama: Advancing text-to-music generation with question answering and captioning. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 286–290. IEEE, 2024.
- Xubo Liu, Qiuqiang Kong, Yan Zhao, Haohe Liu, Yi Yuan, Yuzhuo Liu, Rui Xia, Yuxuan Wang, Mark D Plumbley, and Wenwu Wang. Separate anything you describe. *arXiv preprint arXiv:2308.05037*, 2023b.
- Yuejiang Liu, Parth Kothari, Bastien van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexan dre Alahi. Ttt++: When does self-supervised test-time training fail or thrive? In Advances in
 Neural Information Processing Systems (NeurIPS), 2021.
- Siyu Lou, Xuenan Xu, Mengyue Wu, and Kai Yu. Audio-text retrieval in context. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4793–4797. IEEE, 2022.
- Jiasen Lu et al. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In Advances in Neural Information Processing Systems (NeurIPS 2019), 2019.
 URL https://arxiv.org/abs/1908.02265.
- Huaishao Luo, Lei Ji, Ming Zhong, Yang Chen, Wen Lei, Nan Duan, and Tianrui Li. Clip4clip: An empirical study of clip for end to end video clip retrieval and captioning. *Neurocomputing*, 508: 293–304, 2022.
- Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna.
 Crepe: Can vision-language foundation models reason compositionally? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10910–10921, 2023.
- ⁷³⁷ Ilaria Manco, Emmanouil Benetos, Elio Quinton, and György Fazekas. Contrastive audio-language
 ⁷³⁸ learning for music. *arXiv preprint arXiv:2208.12208*, 2022.
 ⁷³⁹
- Xinhao Mei, Xubo Liu, Jianyuan Sun, Mark D Plumbley, and Wenwu Wang. On metric learning for audio-text cross-modal retrieval. *arXiv preprint arXiv:2203.15537*, 2022.
- Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning. *arXiv* preprint arXiv:2111.09734, 2021.
- Andreea-Maria Oncescu, A Koepke, Joao F Henriques, Zeynep Akata, and Samuel Albanie. Audio retrieval with natural language queries. *arXiv preprint arXiv:2105.02192*, 2021.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pp. 8821–8831. Pmlr, 2021.

781

794

795

- Mirco Ravanelli, Jianyuan Zhong, Santiago Pascual, Pawel Swietojanski, Joao Monteiro, Jan Tr-mal, and Yoshua Bengio. Multi-task self-supervised learning for robust speech recognition. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing* (*ICASSP*), pp. 6989–6993. IEEE, 2020.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer- ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Assaf Shocher, Nadav Cohen, and Michal Irani. "zero-shot" super-resolution using deep internal learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (*CVPR*), pp. 3118–3126, June 2018.
- Amanpreet Singh, Ronghang Hu, Akhilesh Gotmare, Devi Parikh, Christoph Feichtenhofer, Stefan Lee, and Marcus Rohrbach Singh. Flava: A foundational language and vision alignment model. *arXiv preprint arXiv:2112.04482*, 2022.
- Rickard Stureborg, Dimitris Alikaniotis, and Yoshi Suhara. Large language models are inconsistent and biased evaluators. arXiv preprint arXiv:2405.01724, 2024. URL https://arxiv.org/ abs/2405.01724.
- Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei A. Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- Dídac Surís, Carl Vondrick, Bryan Russell, and Justin Salamon. It's time for artistic correspondence in music and video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10564–10574, 2022.
- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5238–5248, 2022.
- Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What
 makes for good views for contrastive learning? *Advances in neural information processing systems*, 33:6827–6839, 2020.
- Y. Wang, R. Singh, and B. Raj. Zero-shot learning for audio classification. In *Proceedings of the 2020 Conference of the International Speech Communication Association*, 2020a.
- Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples:
 A survey on few-shot learning. *ACM Computing Surveys (CSUR)*, 53(3):1–34, 2020b.
 - Zheng Wang, Yang Liu, and Wei Zhao. Large language models are not fair evaluators. *arXiv preprint arXiv:2305.17926*, 2023. URL https://arxiv.org/abs/2305.17926.
- Ho-Hsiang Wu, Magdalena Fuentes, Prem Seetharaman, and Juan Pablo Bello. How to listen?
 rethinking visual sound localization. *arXiv preprint arXiv:2204.05156*, 2022a.
- Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, and Juan Pablo Bello. Wav2clip: Learning
 robust audio representations from clip. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4563–4567. IEEE, 2022b.
- Ho-Hsiang Wu, Oriol Nieto, Juan Pablo Bello, and Justin Salomon. Audio-text models do not yet
 leverage natural language. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023a.
- Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov.
 Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption
 augmentation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. IEEE, 2023b.

810 811 812	Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a com- prehensive evaluation of the good, the bad and the ugly. <i>IEEE Transactions on Pattern Analysis</i> <i>and Machine Intelligence</i> , 41(9):2251–2265, 2018.
814 815	Huang Xie, Samuel Lipping, and Tuomas Virtanen. Language-based audio retrieval task in dcase 2022 challenge. <i>arXiv preprint arXiv:2206.06108</i> , 2022.
816 817	Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-shot
818 819	video-text understanding. arXiv preprint arXiv:2109.14084, 2021.
820 821 822	Hongwei Xue, Yuchong Sun, Bei Liu, Jianlong Fu, Ruihua Song, Houqiang Li, and Jiebo Luo. Clip- vip: Adapting pre-trained image-text model to video-language representation alignment. <i>arXiv</i> <i>preprint arXiv:2209.06430</i> , 2022.
823 824 825	Lu Yuan, Dongdong Chen, Yi-Ling Chen, Vlad Codreanu, Yandong Ge, Wenhan Guo, Yandong Guo, Jianfeng Huang, Ming Li, Ping Li, et al. Florence: A new foundation model for computer vision. <i>arXiv preprint arXiv:2111.11432</i> , 2021.
826 827 828 829	Yi Yuan, Zhuo Chen, Xubo Liu, Haohe Liu, Xuenan Xu, Dongya Jia, Yuanzhe Chen, Mark D Plumbley, and Wenwu Wang. T-clap: Temporal-enhanced contrastive language-audio pretraining. <i>arXiv preprint arXiv:2404.17806</i> , 2024.
830 831 832	Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In <i>The Eleventh International Conference on Learning Representations</i> , 2022.
833 834 835 836	Shuai Zhao, Linchao Zhu, Xiaohan Wang, and Yi Yang. Centerclip: Token clustering for efficient text-video retrieval. In <i>Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pp. 970–981, 2022.
837 838 839 840 841 842 843 844	Yanpeng Zhao, Jack Hessel, Youngjae Yu, Ximing Lu, Rowan Zellers, and Yejin Choi. Connecting the dots between audio and text without parallel data through visual knowledge transfer. <i>arXiv</i> preprint arXiv:2112.08995, 2021.
837 838 839 840 841 842 843 844 845 844 845 846 847 848	Yanpeng Zhao, Jack Hessel, Youngjae Yu, Ximing Lu, Rowan Zellers, and Yejin Choi. Connecting the dots between audio and text without parallel data through visual knowledge transfer. <i>arXiv</i> preprint arXiv:2112.08995, 2021.
837 838 839 840 841 842 843 844 845 846 845 846 847 848 849 850 851	Yanpeng Zhao, Jack Hessel, Youngjae Yu, Ximing Lu, Rowan Zellers, and Yejin Choi. Connecting the dots between audio and text without parallel data through visual knowledge transfer. <i>arXiv</i> preprint arXiv:2112.08995, 2021.
837 838 839 840 841 842 843 844 845 846 845 846 847 848 849 850 851 852 853	Yanpeng Zhao, Jack Hessel, Youngjae Yu, Ximing Lu, Rowan Zellers, and Yejin Choi. Connecting the dots between audio and text without parallel data through visual knowledge transfer. <i>arXiv</i> preprint arXiv:2112.08995, 2021.
837 838 839 840 841 842 843 844 845 844 845 846 847 848 849 850 851 852 853 854 855 856	Yanpeng Zhao, Jack Hessel, Youngjae Yu, Ximing Lu, Rowan Zellers, and Yejin Choi. Connecting the dots between audio and text without parallel data through visual knowledge transfer. <i>arXiv</i> preprint arXiv:2112.08995, 2021.
837 838 839 840 841 842 843 844 845 844 845 846 847 848 849 850 851 852 853 854 855 856 857	Yanpeng Zhao, Jack Hessel, Youngjae Yu, Ximing Lu, Rowan Zellers, and Yejin Choi. Connecting the dots between audio and text without parallel data through visual knowledge transfer. <i>arXiv</i> preprint arXiv:2112.08995, 2021.
837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 855 856 857 858 859	Yanpeng Zhao, Jack Hessel, Youngjae Yu, Ximing Lu, Rowan Zellers, and Yejin Choi. Connecting the dots between audio and text without parallel data through visual knowledge transfer. <i>arXiv preprint arXiv:2112.08995</i> , 2021.
837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 855 856 857 858 859 860 861	Yanpeng Zhao, Jack Hessel, Youngjae Yu, Ximing Lu, Rowan Zellers, and Yejin Choi. Connecting the dots between audio and text without parallel data through visual knowledge transfer. <i>arXiv preprint arXiv:2112.08995</i> , 2021.
837 838 839 840 841 842 843 844 845 844 845 846 847 848 849 850 851 852 853 854 855 855 856 857 858 859 860 861 862	Yanpeng Zhao, Jack Hessel, Youngjae Yu, Ximing Lu, Rowan Zellers, and Yejin Choi. Connecting the dots between audio and text without parallel data through visual knowledge transfer. <i>arXiv preprint arXiv:2112.08995</i> , 2021.

864 APPENDIX А 865 866 В SUPPLEMENTARY SECTION 867 868 **PROOF OF PROPOSITIONS B**.1 870 B.1.1 **PROPOSITION 1**: 871 872 Contrastive models, when used for audio-text matching, do not comprehend the semantic relationship between the audio and text, but rather operate by matching similar audios to similar texts based 873 on superficial features. 874 875 Let $f_{audio}: A \to \mathbb{R}^d$ and $f_{text}: T \to \mathbb{R}^d$ be the functions mapping audio A and text T into a 876 d-dimensional embedding space, respectively. The similarity score between an audio sample a and 877 a text sample t is given by $s(a, t) = f_{audio}(a) \cdot f_{text}(t)$. 878 1. Contrastive models can yield high similarity scores for pairs of audio and text samples that 879 share similar superficial features but lack semantic congruence. 880 2. Contrastive models, as defined, cannot inherently discern semantic relationships between audio and text but rely on the co-occurrence of similar features in their respective embed-883 dings. 884 Assume a pair of audio samples a_1, a_2 and text samples t_1, t_2 such that a_1 and t_1, a_2 and t_2 are 885 semantically congruent but share similar superficial features with a_2 and t_1 respectively. 886 887 888 According to the model, $s(a_1, t_1)$ and $s(a_2, t_2)$ should be high. However, due to the shared 889 superficial features, $s(a_1, t_2)$ and $s(a_2, t_1)$ may also be high, indicating a false positive match. 890 891 This contradiction shows that the model's high similarity score does not necessarily correspond to a 892 true semantic match, supporting the hypothesis. 893 894 **B.2** DATASET SELECTION AND CREATION 895 896 For dataset selection and creation process we chose ESC-50 dataset. Due to its high audio quality, 897 adequate pre-processing, suitable length and number of samples, and its inherent robustness. Importantly, we excluded datasets generated through crowd-sourcing to reduce labeling inaccuracies. ESC-50's assortment of 50 classes encompasses a variety of real-world sounds from natural, animal, 899 and human sources, making it versatile for different applications and particularly effective for zero-900 shot classification tasks, which require identifying items from previously unseen categories. From 901 the ESC-50 dataset we get 50 pairs of Audio, Label data, these pairs are then processed according 902 to algorithm 2 to make a training dataset. We select two distinct sounds from the possible 50 sounds 903 giving us a total of 2450 pairs (a_i, a_j) and (t_i, t_j) of sounds. For each pair we have 3 possible 904 configurations using keywords 'before', 'after' and 'while' as suggested in section 3.1. Thus our 905 total dataset thus becomes 7350 data pairs. For teminAL A, we only use 2450 pairs of data while 906 selecting either one of the audio and text from this pair. The prompt used for concatenating the texts 907 are 'single sound of t_i ' and 'combined sound of t_i and t_i '. 908 Our Sequential Inversion Approach challenges traditional contrastive learning methods, which typ-909 ically align audio segments with matching text while contrasting them against unrelated pairs. This 910 practice, akin to a bag-of-words model, often fails to capture sequential nuances as it emphasizes dis-911 tinguishing features over temporal understanding. To foster a deeper comprehension of sequences, 912 we introduced a novel technique for generating negative samples that share thematic elements, com-913 pelling the model to focus on the order of events. This method, depicted in fig. 2, utilizes two types 914 of temporal augmentations "before" and "while" to enhance the model's ability to discern sequential 915 information. The transformation aims to capture the dynamic interplay between the two arguments, allowing the model to discern the original audio-transcript pair (a, c) from its transformed versions 916

917 $(a, \mathbb{O}(c))$ and $(\mathbb{O}(a), c)$. It is applicable to concatenated audio or transcription pairs, effectuating a temporal reordering of the components.

918 Algorithm 1 Dataset Preparation and Sequential Inversion for Contrastive Learning 919 1: Input: Acoustic Events Dataset $P = \{(Audio_i, Label_i)\}$ 920 2: **Output:** Refined Dataset $(A_{pos}, A_{neg}, T_{pos}, T_{neg})$ for Contrastive Learning 921 3: Initialization: 922 4: Initialize ESC-50 dataset with selected criteria. 923 5: Organize dataset into classes $\{C_1, C_2, \ldots, C_n\}$ representing distinct sounds. 924 6: Define empty sets for positive and negative samples: $A_{pos}, A_{neg}, T_{pos}, T_{neg}$. 925 926 7: Step 1: Positive Sample Selection 927 8: for each class C_i in the dataset do for each audio-text pair $(a_i, c_i) \in C_i$ do 9: 928 10: if (a_i, c_i) meets quality standards then 929 Åppend a_j to A_{pos} , c_j to T_{pos} . 11: 930 12: end if 931 13: end for 932 14: end for 933 934 15: Step 2: Sequential Inversion and Overlay 935 16: for each $(a_j, c_j) \in (A_{pos}, T_{pos})$ do 936 17: Generate negative audio samples a'_i using inversion function \mathbb{T} . 937 18: Define $\mathbb{T}(a) = [a_j \oplus a_i], \mathbb{T}(c) = [c_j; \tau_t; c_i].$ 938 19: Generate overlapping samples using overlay function \mathbb{O} : $\mathbb{O}(a) = [a_j \wedge a_i], \mathbb{O}(c) = [c_j; \tau_o; c_i].$ 939 Append resulting negative samples a'_i, c'_i to A_{neg} and T_{neg} . 940 20: 941 21: end for 942 22: Step 3: Template-Based Caption Generation 943 23: for each positive sample $c_k \in T_{pos}$ do 944 Generate captions using Caption (c_k) and append to T_{pos} . 24: 945 25: end for 946 26: for each negative sample $c_n \in T_{neg}$ do 947 27: Generate captions using $Caption(c_n)$ and append to T_{neq} . 948 28: end for 949 950 29: **Return:** $(A_{pos}, A_{neq}, T_{pos}, T_{neq})$. 951 952 953 MATHEMATICAL DERIVATIONS: **B.3** 954 955

In this section, we derive the loss functions used in our model, specifically focusing on the Temporal Noise Contrastive Estimation (TNCE) technique. TNCE is a variant of the Noise Contrastive Estimation (NCE) loss, adapted for temporal learning tasks. This method helps in effectively distinguishing between positive and negative samples over time. (Kindly note that we have used 't' as text in the Mathematical derivation instead of 'c' as we have shown in the main paper, all the other component remains the same. For example here we have shown batch of texts $B_t = \{B_{t_f}, B_{t_r}, B_{t_o}\}$ instead of $B_c = \{B_{c_f}, B_{c_r}, B_{c_o}\}$).

963 For the loss function L_r , we define it as follows:

964 965

966

956

957

958

959

960

961

962

$$L_{t_B} = \sum_{(a,t)\in B} (\text{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_t) + \text{TNCE}(\boldsymbol{z}_{\mathbb{T}(t)}, \boldsymbol{z}_a)) + \text{TNCE}(\boldsymbol{z}_{\mathbb{O}(t)}, \boldsymbol{z}_a))$$
(10)

967 968 969

970 Here, $\text{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_t)$, $\text{TNCE}(\boldsymbol{z}_{\mathbb{T}(t)}, \boldsymbol{z}_a)$) and $\text{TNCE}(\boldsymbol{z}_{\mathbb{O}(t)}, \boldsymbol{z}_a)$) represent the temporal consistent, 971 temporally reversed and temporally overlap components of the TNCE loss, respectively. The function TNCE is calculated by the formula:



Figure 6: Schematic explanation of the terms in loss function for TeminAL A. Here we show a term (row) in the summation of L_{t_A} which is $\text{TNCE}_t(z_{a_s}, z_{t_s})$ The other term $\text{TNCE}_t(z_{a_d}, z_{t_d})$ of this loss function can be calculated in the similar way and will belong to the green block of the above schematic. Here, B_{t_s} and B_{t_d} are the batches of texts corresponding to single and concatenated (double) samples which compose the whole batch of text following the same convention as shown in section 3.4.

$$\text{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_t) := -\log \frac{\exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_t)}{\sum_{t' \in B_{t_f}} \exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_{t'}) + C^{t_r} + C^{t_o}}$$
(11)

(12)

(17)

Similarly, the overlap component $TNCE_0$ is given by:

 $\mathsf{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_t) := -\log \frac{\exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_t)}{\sum_{t' \in B_{t_o}} \exp(z_a \cdot \boldsymbol{z}_{t'}) + C^{t_c}}$

In these equations: B represents the batch of user-item pairs (a, t), where a is a user and t is a temporal context. z_a and z_t denote the latent representations of the user and the temporal context, respectively. B_{t_f} and B_{t_o} are subsets of the batch B that serve as temporal and overlap negatives, respectively. The constants C^{t_r} , C^{t_o} , and C^{t_c} are designed to account for additional temporal and contextual information, enhancing the robustness of the loss function against trivial solutions. The term C^{t_r} accounts for the influence of time-reversed negatives and is defined as:

$$C^{t_r} = \alpha_{s_t} \exp(\boldsymbol{z}_u \cdot \boldsymbol{z}_{\Pi(t)}) + \alpha_{c_t} \sum_{t' \in B_{t_r} \setminus \{t\}} \exp(\boldsymbol{z}_u \cdot \boldsymbol{z}_{\Pi(t')})$$
(13)

where: $\Pi(t)$ denotes the time-reversed representation of the context t. The coefficients α_{s_t} and α_{c_t} modulate the contribution of individual and cumulative time-reversed negatives, respectively. The term C^{t_o} captures the effect of overlapping contexts, defined as:

$$C^{t_o} = \alpha_{s_o}(\exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_t) + \exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_{\Pi(t)})) + \alpha_{c_o} \sum_{\substack{t' \in B_t \setminus \{o\}}} \exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_{\Pi(t')})$$
(14)

Here: α_{s_0} and α_{c_0} control the impact of single and multiple overlapping contexts.

Finally, C^{t_c} integrates both temporal and contextual negative sampling:

$$C^{t_{c}} = \left(\exp(z_{a} \cdot \boldsymbol{z}_{t}) + \sum_{t' \in B_{t_{f}} \setminus \{t\}} \exp(z_{a} \cdot \boldsymbol{z}_{t'}) \right) + \left(\alpha_{s} \exp(\boldsymbol{z}_{a} \cdot \boldsymbol{z}_{\Pi(t)}) + \alpha_{c} \sum_{t' \in B_{t_{r}} \setminus \{t\}} \exp(\boldsymbol{z}_{a} \cdot \boldsymbol{z}_{\Pi(t')}) \right)$$
(15)

This term combines the effect of immediate and cumulative context influences, with parameters α_s and α_c providing tunable weights.

For the loss function L_{a_B} , which deals with another set of temporal dynamics, we follow a similar structure. The formulation and constants remain analogous, ensuring consistency across different temporal modeling aspects.

$$L_{a_B} = \sum_{(a,t)\in B} (\text{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_t) + \text{TNCE}(\boldsymbol{z}_{\mathbb{T}(t)}, \boldsymbol{z}_a)) + \text{TNCE}(\boldsymbol{z}_{\mathbb{O}(t)}, \boldsymbol{z}_a))$$
(16)

 $\exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_t)$

Here, TNCE stands for Temporal Noise Contrastive Estimation, a variant of the NCE loss tailored for temporal learning, and is calculated as:

$$\text{TNCE}(\boldsymbol{z}_a, \boldsymbol{z}_t) = -\log \frac{\exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_t)}{\sum_{t' \in B_{t_f}} \exp(\boldsymbol{z}_a \cdot \boldsymbol{z}_{t'}) + C^{t_r} + C^{t_o}}$$

$$\text{TNCE}(\boldsymbol{z}_{\mathbb{O}(t)}, \boldsymbol{z}_{a})) = -\log \frac{\exp(\boldsymbol{z}_{a} \cdot \boldsymbol{z}_{t})}{\sum_{t' \in B_{t_{a}}} \exp(\boldsymbol{z}_{a} \cdot \boldsymbol{z}_{t'}) + C^{t_{c}}}$$
(18)



1133 Here, $\text{TNCE}(\boldsymbol{z}_{t_s}, \boldsymbol{z}_{a_s})$ and $\text{TNCE}(\boldsymbol{z}_{t_d}, \boldsymbol{z}_{a_d})$ represent the temporal and overlap components of the TNCE loss, respectively. $\boldsymbol{z}_{t_s}, \boldsymbol{z}_{a_s}$ represents the text and audio samples of single samples in the

batch. And z_{t_d} , z_{a_d} represents the text and audio samples of the double or concatenated batch. The function TNCE is calculated by the formula:

$$\text{TNCE}(\boldsymbol{z}_{a_s}, \boldsymbol{z}_{t_s}) = -\log \frac{\exp(\boldsymbol{z}_{a_s} \cdot \boldsymbol{z}_{t_s})}{\sum_{t' \in B_{t_s}} \exp(\boldsymbol{z}_{a_s} \cdot \boldsymbol{z}_{t'_s}) + C^{t_d}}$$
(23)

Where C^{t_d} the contribution of the concatenated samples to the above loss function.

$$C^{t_d} = \alpha_{same} \exp(\boldsymbol{z}_{a_d} \cdot \boldsymbol{z}_{\Pi(t)}) + \alpha_{diff} \sum_{t' \in B_{t_d} \setminus \{t\}} \exp(\boldsymbol{z}_{a_d} \cdot \boldsymbol{z}_{\Pi(t'_d)})$$
(24)

1150 The terms α_{same} in the above represent the concatenated samples which have one of the sounds 1151 similar to z_{a_s} , while α_{diff} is the co–efficient used for all the concatenated samples (z_{a_d}) which 1152 don't have any sound similar to z_{a_s} . Next up we have similar formulation for the other half of the 1153 TeminAL A loss function which is shown below.

$$L_{a_A} = \sum_{(\mathbb{O}(a),\mathbb{O}(t))\in B} (\text{TNCE}(\boldsymbol{z}_{a_s}, \boldsymbol{z}_{t_s}) + \text{TNCE}(\boldsymbol{z}_{a_d}, \boldsymbol{z}_{t_d}))$$
(25)

(26)

Finally the overall loss function for TeminAL A is composed of L_{t_A} and L_{a_A} shown as follows. Note, We keep all our hyper-parameters set as unity for the training of TeminAL A.

 $L_A = L_{t_A} + \beta_A(L_{a_A})$

The rest of the formulation follows the same derivation scheme as what we have detailed for TeminAL B in the above paragraphs.

B.4 ZERO SHOT DOWNSTREAM TASK AND DETAILS:





Table 4: Performance comparison on audio classification task on different datasets. For ESC-50 and US8K we have used the prompt "The sound of a {class}" over all the 50 and 10 classes respectively.
For ESC-50 the other text prompts are from the validation set of the model.

Method Wav2CLIP AudioClip CLAP CLAP-LAION-audio-630K CompA-CLAP T-CLAP (ours) n 2 ZSTE : Zero Shot Temporal Evaluation; es for General-purpose contrastive training m by is detailed in appendix B.4 also refer apport : Dataset \mathcal{D} , Contrastive Learning-based Method i: Dataset \mathcal{D} , Contrastive Learning-based Method i: Dataset \mathcal{D} and the contrastive learning-based Method i: Basic Zero-Shot Evaluation ad dataset \mathcal{D} and the contrastive learning-based 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $\lambda_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features 1: model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	ESC-50 41.4 69.4 82.6 88.0 89.1 75.1 evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat for unseer omposite fe ictions for r improven	US8K 40.4 65.3 73.2 75.8 85.7 72.2 Zero-Shot Tell models. Im for detail of \mathcal{M} tion tasks \mathcal{T}_1 a classes reference on the second	The posite instances refer
Wav2CLIP AudioClip CLAP CLAP-LAION-audio-630K CompA-CLAP T-CLAP (ours) n 2 ZSTE : Zero Shot Temporal Evaluation; es for General-purpose contrastive training m dy is detailed in appendix B.4 also refer apport i: Dataset \mathcal{D} , Contrastive Learning-based Matt: Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-based I: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $M_i \in \mathcal{U}$ $1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc1 $\mathbf{2: Zero-Shot with Overlapping Features}$ t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc2 = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	41.4 69.4 82.6 88.0 89.1 75.1 evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificate for unseer omposite for composite for the formula form	40.4 65.3 73.2 75.8 85.7 72.2 Zero-Shot Te I models. Im for detail of \mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	Temporal Classification uplementation of ZSTE n parameters. For algorithm 3: $Acc_1 =$ uposite instances refer
AudioClip CLAP CLAP-LAION-audio-630K CompA-CLAP T-CLAP (ours) n 2 ZSTE: Zero Shot Temporal Evaluation; es for General-purpose contrastive training m ly is detailed in appendix B.4 also refer apport : Dataset \mathcal{D} , Contrastive Learning-based Me ut: Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-based 1: Basic Zero-Shot Evaluation fluate model's zero-shot capabilities on basic asure accuracy by correct label identification fluate model's zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Featuress t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation of 3: Temporal Relationship Comprehension	69.4 82.6 88.0 89.1 75.1 evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat for unseer omposite for ictions for r improven	\mathcal{M} tion tasks \mathcal{T}_1 n classes refe	Temporal Classification iplementation of ZSTE n parameters. For algorithm 3: $Acc_1 =$
CLAP CLAP-LAION-audio-630K CompA-CLAP T-CLAP (ours) a 2 ZSTE : Zero Shot Temporal Evaluation; es for General-purpose contrastive training m by is detailed in appendix B.4 also refer apport b b c Dataset \mathcal{D} , Contrastive Learning-based Me ut : Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-based 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $M_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ \mathcal{C} } \sum_{j \in \mathcal{C}} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	82.6 88.0 89.1 75.1 evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat n for unseer omposite fa- ictions for r improven	73.2 75.8 85.7 72.2 Zero-Shot Ta I models. Im for detail of \mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	Temporal Classification uplementation of ZSTE n parameters. For algorithm 3: $Acc_1 =$
CLAP-LAION-audio-630K CompA-CLAP T-CLAP (ours) a 2 ZSTE : Z ero Shot Temporal Evaluation; es for General-purpose contrastive training m by is detailed in appendix B.4 also refer apport b a b b a b b c b c b c b c b c b c b c b c b c b c c c c c c c c c c	88.0 89.1 75.1 evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat n for unseer omposite fa- ictions for r improven	75.8 85.7 72.2 Zero-Shot Ta I models. Im for detail of \mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com-	Temporal Classification uplementation of ZSTE n parameters. For algorithm 3: $Acc_1 =$
CompA-CLAP T-CLAP (ours) n 2 ZSTE: Zero Shot Temporal Evaluation; es for General-purpose contrastive training m by is detailed in appendix B.4 also refer apport : Dataset \mathcal{D} , Contrastive Learning-based Me ut: Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-based 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $\sum_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat for unseer omposite fa	Zero-Shot Ta 72.2 Zero-Shot Ta models. Im for detail or \mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com-	Temporal Classification plementation of ZSTE n parameters. er algorithm 3: $Acc_1 =$ posite instances refer
T -CLAP (ours) T -CLAP (ours) n 2 ZSTE: Zero Shot Temporal Evaluation; tes for General-purpose contrastive training magnetizes for General-purpose contrastive training magnetizes for General-purpose contrastive training-based Metrix Model evaluation scores for zero-shot tallization: and dataset \mathcal{D} and the contrastive learning-based Metrix Model evaluation scores for zero-shot tallization: and dataset \mathcal{D} and the contrastive learning-based Metrix Model evaluation scores for zero-shot tallization: and dataset \mathcal{D} and the contrastive learning-based Metrix Model's zero-shot Evaluation asure accuracy by correct label identification $Y_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label predict thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation of 3: Temporal Relationship Comprehension	75.1 evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat for unseer omposite for ictions for r improven	72.2 Zero-Shot Ta models. Im for detail or \mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	Temporal Classification pplementation of ZSTE n parameters. Er algorithm 3: $Acc_1 =$ posite instances refer
n 2 ZSTE : Zero Shot Temporal Evaluation; es for General-purpose contrastive training m by is detailed in appendix B.4 also refer appo- contrastive Learning-based Me ut : Model evaluation scores for zero-shot ta lization: and dataset \mathcal{D} and the contrastive learning-based lization: and dataset \mathcal{D} and the contrastive learning-based lization: and dataset \mathcal{D} and the contrastive learning-base li Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $Y_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat for unseer omposite fa ictions for r improven	Zero-Shot To I models. Im for detail or \mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen corr	Temporal Classification aplementation of ZSTE n parameters. For algorithm 3: $Acc_1 =$ apposite instances refer
n 2 ZSTE : Zero Shot Temporal Evaluation; es for General-purpose contrastive training m dy is detailed in appendix B.4 also refer appo- : Dataset \mathcal{D} , Contrastive Learning-based Me ut : Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-base 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $\lambda_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat for unseer omposite fe ictions for r improven	Zero-Shot Ta I models. Im for detail or \mathcal{M} tion tasks \mathcal{T}_1 classes refe eatures \mathcal{T}_2 unseen com	Temporal Classification aplementation of ZSTE n parameters. Example: $Acc_1 =$ aposite instances refer
n 2 ZSTE: Zero Shot Temporal Evaluation; es for General-purpose contrastive training n dy is detailed in appendix B.4 also refer appe : Dataset \mathcal{D} , Contrastive Learning-based Me ut: Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-based 1: Basic Zero-Shot Evaluation duate model's zero-shot capabilities on basic asure accuracy by correct label identification $\sum_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	evaluating nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat for unseer omposite fa ictions for r improven	Zero-Shot Te I models. Im for detail of \mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	Temporal Classification applementation of ZSTE n parameters. Ex algorithm 3: $Acc_1 =$ apposite instances refer
es for General-purpose contrastive training n ly is detailed in appendix B.4 also refer appe : Dataset \mathcal{D} , Contrastive Learning-based Me ut: Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-base 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $f_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	nulti-modal endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classificat for unseer omposite fe ictions for r improven	I models. Im for detail of \mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	plementation of ZSTE n parameters. er algorithm 3: $Acc_1 =$ nposite instances refer
iy is detailed in appendix B.4 also refer apport : Dataset \mathcal{D} , Contrastive Learning-based Mut: Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-based 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $M_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Featuress t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ \mathcal{C} } \sum_{j \in \mathcal{C}} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation of 3: Temporal Relationship Comprehension	endix B.4.1 odel \mathcal{M} sks \mathcal{S} sed model . c classification for unseer omposite for ictions for r improven	tion tasks \mathcal{T}_1 a classes reference on tasks \mathcal{T}_2 unseen company S	n parameters. er algorithm 3: $Acc_1 =$ nposite instances refer
: Dataset \mathcal{D} , Contrastive Learning-based Met it: Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-base 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $f_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ \mathcal{C} } \sum_{j \in \mathcal{C}} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	odel <i>M</i> sks <i>S</i> sed model . c classificat for unseer omposite fo ictions for r improven	\mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	er algorithm 3: $Acc_1 =$
ut: Model evaluation scores for zero-shot ta lization: ad dataset \mathcal{D} and the contrastive learning-base 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $y_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ \mathcal{C} } \sum_{j \in \mathcal{C}} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	sks <i>S</i> sed model . c classificat for unseer omposite fo ictions for r improven	\mathcal{M} tion tasks \mathcal{T}_1 o classes refe eatures \mathcal{T}_2 unseen com	er algorithm 3: $Acc_1 =$
lization: ad dataset \mathcal{D} and the contrastive learning-base 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $y_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ \mathcal{C} } \sum_{j \in \mathcal{C}} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	sed model . c classificat for unseer omposite for ictions for r improven	\mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	er algorithm 3: $Acc_1 =$
ad dataset \mathcal{D} and the contrastive learning-bas 1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $\sum_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ \mathcal{C} } \sum_{j \in \mathcal{C}} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	sed model . c classificat n for unseer omposite for ictions for r improven	\mathcal{M} tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	er algorithm 3: $Acc_1 =$
1: Basic Zero-Shot Evaluation aluate model's zero-shot capabilities on basic asure accuracy by correct label identification $\sum_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	c classificat for unseer omposite for ictions for r improven	tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	er algorithm 3: $Acc_1 =$
Inducte model's zero-shot capabilities on basis asure accuracy by correct label identification $\sum_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	c classificat for unseer omposite for ictions for r improven	tion tasks \mathcal{T}_1 a classes refe eatures \mathcal{T}_2 unseen com	er algorithm 3: $Acc_1 =$
asure accuracy by correct label identification $\sum_{i \in \mathcal{U}} 1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	n for unseer omposite for ictions for r improven	a classes refe eatures \mathcal{T}_2 unseen com	er algorithm 3: $Acc_1 =$ nposite instances refer
$1(\hat{y}_i = y_i)$ cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	omposite fo ictions for r improven	eatures \mathcal{T}_2 unseen com	nposite instances refer
cord baseline zero-shot performance Acc ₁ 2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: Acc ₂ = $\frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	omposite fe ictions for r improven	eatures \mathcal{T}_2 unseen com	nposite instances refer
2: Zero-Shot with Overlapping Features t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: $Acc_2 = \frac{1}{ C } \sum_{j \in C} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	omposite fe ictions for r improven	eatures \mathcal{T}_2 unseen com	nposite instances refer
t model's ability to discern overlapping or c asure accuracy based on correct label pred thm 3: $\operatorname{Acc}_2 = \frac{1}{ \mathcal{C} } \sum_{j \in \mathcal{C}} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	omposite fe ictions for r improven	eatures \mathcal{T}_2 unseen com	nposite instances refer
asure accuracy based on correct label pred thm 3: $\operatorname{Acc}_2 = \frac{1}{ \mathcal{C} } \sum_{j \in \mathcal{C}} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	ictions for r improven	unseen com	nposite instances refer
thm 3: $\operatorname{Acc}_2 = \frac{1}{ \mathcal{C} } \sum_{j \in \mathcal{C}} 1(\hat{y}_j = y_j)$ cord and analyze performance degradation o 3: Temporal Relationship Comprehension	r improven	ant S.	
cord and analyze performance degradation o 3: Temporal Relationship Comprehension	r improven	pont S	
3: Temporal Relationship Comprehension	i impioven		
5. Temporal Relationship Comprehension	•	left \mathcal{O}_2	
cent model with unceen cequences () to acce	I ss tempora	l relationshi	in understanding \mathcal{T}_{r}
source accuracy in identifying the correct.	order of e	ants refer a	punderstanding 73
asure accuracy in identifying the correct $(1/2) = 1/2$		cints ieiei a	ingomulii 5. $Acc_3 =$
$k \in \mathcal{Q}$ $\mathbf{I}(O_k = O_k)$	1	1	1
iluate against known sequences to determine	e zero-snot	temporal col	mprenension \mathcal{S}_3
4: Resistance to Irrelevant Features	1		au
allenge model with unseen data \mathcal{N} that inclu	des irreleva	ant reatures	
termine model's ability to ignore noise and $\sum_{i=1}^{n} \frac{1}{2} \left(\hat{a}_{i} + i \right)$	locus on re	elevant zero-s	snot features: $Acc_4 =$
$\mathcal{L}_{l\in\mathcal{N}} \mathbf{L}(y_l = y_l)$		-	
sess confusion metrics and resilience to irrel	evant data	S_4	
5: Generalization to Novel Scenarios			_
luate model's generalization to completely	novel zero-	shot scenario	os \mathcal{T}_5
asure model's performance on tasks with ne	w contexts	or relationsl	hips refer algorithm 3:
$= \frac{1}{ \mathcal{X} } \sum_{m \in \mathcal{X}} 1(\hat{y}_m = y_m)$			
t for understanding of complex temporal see	quences and	d novel featu	re combinations \mathcal{S}_5
lusion:			
mpile and compare evaluation scores across	all tasks ${\cal S}$	$= \{\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_2, \mathcal{S}_2\}$	$\{\mathcal{S}_3, \mathcal{S}_4, \mathcal{S}_5\}$
termine model's strengths and weaknesses in	n zero-shot	learning	
vide insights into model's potential real-wor	rld applicat	oility	
n Compiled evaluation scores S , insights, ar	d potential	applications	S
	aluate against known sequences to determine 4: Resistance to Irrelevant Features allenge model with unseen data \mathcal{N} that inclu- termine model's ability to ignore noise and $\sum_{l \in \mathcal{N}} 1(\hat{y}_l = y_l)$ sess confusion metrics and resilience to irrele 5: Generalization to Novel Scenarios aluate model's generalization to completely pasure model's performance on tasks with ne $= \frac{1}{ \mathcal{X} } \sum_{m \in \mathcal{X}} 1(\hat{y}_m = y_m)$ at for understanding of complex temporal second husion: mpile and compare evaluation scores across termine model's strengths and weaknesses in wide insights into model's potential real-word n Compiled evaluation scores \mathcal{S} , insights, and Fask 1 : In our initial experiment, we aim traightforward classification task devoid of termine if T-CLAP exhibited any improvement	aluate against known sequences to determine zero-shot 4: Resistance to Irrelevant Features allenge model with unseen data \mathcal{N} that includes irrelevant termine model's ability to ignore noise and focus on re $\sum_{l \in \mathcal{N}} 1(\hat{y}_l = y_l)$ sess confusion metrics and resilience to irrelevant data of 5: Generalization to Novel Scenarios aluate model's generalization to completely novel zero- asure model's performance on tasks with new contexts $= \frac{1}{ \mathcal{X} } \sum_{m \in \mathcal{X}} 1(\hat{y}_m = y_m)$ at for understanding of complex temporal sequences and husion: mpile and compare evaluation scores across all tasks \mathcal{S} termine model's strengths and weaknesses in zero-shot wide insights into model's potential real-world applicate n Compiled evaluation scores \mathcal{S} , insights, and potential Fask 1 : In our initial experiment, we aimed to evalu- traightforward classification task devoid of a temporal ermine if T-CLAP exhibited any improvement or loss	aluate against known sequences to determine zero-shot temporal co 4: Resistance to Irrelevant Features allenge model with unseen data \mathcal{N} that includes irrelevant features termine model's ability to ignore noise and focus on relevant zero- $\sum_{l \in \mathcal{N}} 1(\hat{y}_l = y_l)$ sess confusion metrics and resilience to irrelevant data \mathcal{S}_4 5: Generalization to Novel Scenarios aluate model's generalization to completely novel zero-shot scenari asure model's performance on tasks with new contexts or relations $= \frac{1}{ \mathcal{X} } \sum_{m \in \mathcal{X}} 1(\hat{y}_m = y_m)$ at for understanding of complex temporal sequences and novel features husion: mpile and compare evaluation scores across all tasks $\mathcal{S} = \{\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_2, \mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_2, \mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_2, \mathcal{S}_1, \mathcal{S}_2, \mathcal{S}$

1296 class in the ESC dataset. We then measured accuracy by assessing how often the model 1297 correctly identified the label associated with a given audio input (refer to Figure 10). 1298 • Task 2 : Subsequently, we explored whether T-CLAP demonstrated enhanced abilities in 1299 discerning two distinct sounds within a given audio clip with one of the sounds being from 1300 the validation set. The task configuration paralleled that of Task 1, with the key difference 1301 being that the accuracy assessment was conducted on audio clips featuring either concate-1302 nated or overlapping sounds (refer to Figure 11). We measured two accuracy metrics: one 1303 based on the model correctly identifying the two highest probabilities corresponding to the correct labels and another based on the model selecting at least one correct class. 1304 1305 • **Task 3**: In contrast to the preceding task, which disregarded temporality, this new experiment focuses on assessing T-CLAP's capability to accurately discern classes with their respective temporal relationships. For this task, we presented the model with three prompts following the same format as those encountered during training: "[class label 1] before [class label 2]", "[class label 2] before [class label 1]", and "[class label 1] while [class label 2] 1309 bel 2]" (see Figure 12) while picking the 2 classes similar to Task 2. By exposing the model 1310 to an audio featuring one of these three temporal combinations, we gauged its accuracy in 1311 correctly identifying the corresponding temporal relationship within each prompt. 1312 • Task 4 : This task represents a more challenging iteration of Task 3. Here, our objec-1313 tive is to challenge the model by introducing prompts that include additional class labels not present in the audio (Figure 13), aiming to create confusion for the model during the 1315 evaluation process. 1316 • Task 5 : In our final task, we aimed to push the model's boundaries by presenting it with a 1317 temporal prompt it had not encountered during training, assessing its ability to generalize to 1318 novel temporal inputs. Our hypothesis was rooted in the nature of the text encoder, T5; if T-1319 CLAP had truly grasped the temporal nuances embedded in "before" and "while" prompts, 1320 it should demonstrate an understanding of temporality across various prompt formats. For 1321 testing its comprehension of the "before" temporal aspect, we provided the model with four prompts structured as follows: "In this concatenated sound" followed by"The first sound is [class label 1]", "The second sound is [class label 1]", "The first sound is [class label 2]" and "The second sound is [class label 2]" (refer to Figure 14). In each instance, there were two correct prompts, and we evaluated the model based on its ability to correctly identify 1326 the combination of two prompts out of the six possible options. The model received a score of 1 if it correctly identified both prompts and 0.5 if it identified only one. Regarding the "while" temporality, we presented the model with 50 diverse prompts of 1328 the form "Simultaneous sound of [class label 1] and [class label 2]." The model's task was to select the two correct prompts, considering the two correct classes in both possible 1330

1332 1333

1334

B.4.1 PARAMETER LIST FOR ALGORITHM 3

1335 \mathcal{D} : Dataset used for evaluation, \mathcal{M} : Contrastive learning-based model being evaluated, \mathcal{S} : Model evaluation scores for zero-shot tasks, \mathcal{T}_1 : Basic classification tasks for zero-shot evaluation, \mathcal{U} : 1336 Set of unseen classes in basic classification tasks, Acc_1 : Accuracy for basic zero-shot classification 1337 tasks, \hat{y}_i : Predicted label for the *i*-th unseen class, y_i : True label for the *i*-th unseen class, C: Set of 1338 unseen composite instances in overlapping features tasks, Acc₂ : Accuracy for zero-shot tasks with 1339 overlapping features, \hat{y}_i : Predicted label for the j-th composite instance, y_i : True label for the j-th 1340 composite instance, S_2 : Performance evaluation for overlapping features tasks, Q: Set of unseen 1341 sequences in temporal relationship comprehension tasks, Acc_3 : Accuracy for zero-shot temporal relationship comprehension tasks, \hat{o}_k : Predicted order for the k-th sequence, o_k : True order for the k-th sequence, S_3 : Performance evaluation for temporal relationship comprehension tasks, N1344 : Set of unseen data including irrelevant features, Acc_4 : Accuracy for tasks involving irrelevant 1345 features, \hat{y}_l : Predicted label for the *l*-th instance in irrelevant features task, y_l : True label for the *l*-th instance in irrelevant features task, S_4 : Performance evaluation for resistance to irrelevant features, \mathcal{T}_5 : Tasks for evaluating generalization to novel scenarios, \mathcal{X} : Set of instances in novel 1347 scenarios, Acc₅ : Accuracy for generalization to novel zero-shot scenarios, \hat{y}_m : Predicted label for 1348 the *m*-th instance in novel scenarios, y_m : True label for the *m*-th instance in novel scenarios, S_5 : 1349 Performance evaluation for generalization to novel scenarios

in identifying both correct prompts or one, as appropriate.

orderings. The same reward function was applied, scoring the model based on its accuracy



1404 **B.6 BASELINE MODELS** 1405

1434 1435

1436 1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1406 In evaluating retrieval tasks, specifically text-to-audio and audio-to-text, we assess CompA-CLAP 1407 alongside six other baseline models. MMT Oncescu et al. (2021) revolutionized the task of audio retrieval by introducing the use of free-form natural language queries, suggesting this method is more 1408 natural and versatile compared to traditional techniques reliant on text annotations. The research 1409 also highlights the advantages of pre-training on a variety of audio tasks. ML-ACT Mei et al. (2022) 1410 investigates the effects of distinct metric learning objectives on audio-text retrieval tasks, identifying 1411 the NT-Xent loss as a particularly effective method that consistently performs well across various 1412 datasets and training conditions, surpassing commonly-used triplet-based losses. Metric learning 1413 objectives are crucial for training cross-modal retrieval systems, as they organize data into an em-1414 bedding space where similar items cluster together and dissimilar ones are separated. CLAP Elizalde 1415 et al. (2023) presents a new framework for retrieving audio utilizing a contrastive learning objec-1416 tive along with dual audio encoders to bridge the gap between language and audio content. Lastly, 1417 CLAP-LAION Wu et al. (2023b) offers a methodology for contrastive language-audio pre-training, 1418 aiming to forge robust audio representations by marrying audio data with corresponding natural language descriptions. Their model considers various audio and text encoders and enhances the model 1419 architecture with feature fusion strategies and keyword-to-caption augmentation. 1420

Model		T-A Retrieva	1	А	-T Retrieval	
	R@1	R@5	R@10	R@1	R@5	R@10
Pengi	36.2 / 9.4	76.0 / 26.1	86.8 / 36.7	16.9 / 7.0	72.8 / 22.7	84.5 / 34.6
Qwen-Audio	39.1 / 16.2	78.9 / 45.8	87.1 / 57.2	38.0 / 16.1	73.2/23.3	85.0/35.1
Audio Flamingo	41.9 / 18.0	80.2 / 46.3	93.9 / 58.0	38.9 / 17.01	78.9 / 44.0	85.7 / 55.8
CLAP	34.6 / 16.7	70.2 / 41.1	82.0 / 54.1	41.9 / 20.0	73.1 / 44.9	84.6 / 58.7
T-CLAP(ours)	35.1 / 17.0	71.2 / 42.2	82.1 / 54.7	49.2 / 23.1	85.1 / 52.2	87.8 / 66.4

1430 Table 7: Comparison of models with open-ended generation models on Text-Audio and Audio-Text 1431 retrieval performance on the AudioCap/Clotho dataset. The results for previous models have been 1432 taken from (Deshmukh et al., 2023; Elizalde et al., 2023). For retrieval in open-ended generation 1433 models, we use a consistent prompt style as mentioned in (Deshmukh et al., 2023).

B.7 LIMITATIONS OF THE CURRENT MODEL

- 1. General-Purpose ALM: Our proposed model is not a general-purpose Audio-Language Model (ALM) capable of performing all downstream applications across all datasets. The current implementation is specifically designed and validated on the ESC-50 dataset as a proof-of-concept and to achieve our defined objectives. Consequently, generalization remains a limitation, although addressing this was beyond the scope of our work.
- 2. Temporality Beyond the Dataset: The model does not provide a general understanding of temporality beyond the ESC-50 dataset. Results from the ZSTE Task 4 and 5 confirm that neither does the model propose general temporality, nor does it achieve it. Notably, we emphasize that achieving general-purpose temporality would require larger, more comprehensive pre-trained text encoders. Specifically, open-ended text encoders (e.g., encoders from encoder-decoder models) would be more suitable than encoders trained on closed masked language modeling techniques, such as BERT.
- 1449 3. Zero-Shot Evaluation Scheme: While our zero-shot evaluation scheme is designed to be 1450 general-purpose, it is inherently limited to contrastive models. Furthermore, the evaluation has not been tested on domains beyond the ESC-50 dataset. The reported ZSTE results are 1451 restricted to this dataset because it aligns with the training data used in our model and those 1452 of prior works. To mitigate data leakage, we ensure that our evaluation dataset is separate 1453 from the training data. However, it is important to note that the primary objective of this 1454 evaluation is to validate the proof-of-concept rather than to set new benchmarks. 1455
- 4. Evaluation Dataset Overlap: A broader limitation, relevant to most models in this do-1456 main, is the overlap between evaluation datasets across various benchmarks. These over-1457 laps can occur in terms of sound events or contextual similarities. Therefore, we cau-

1458 tion against uncritical comparisons of model performance on classical benchmarks without 1459 careful consideration of dataset overlap. 1460 1461 By addressing these limitations, we hope to guide future researchers in expanding and improving 1462 upon the current work to achieve broader applicability and more generalized performance. 1463 1464 **B.8** EVALUATION METRICS 1465 1466 Our evaluation metrics are task specific, but in general they follow a similar strategy. The primary objective of the model appears to be to determine how well it can match audio clips with their 1467 corresponding textual descriptions. Here's a breakdown of key elements in the code and how they 1468 can be translated into a mathematical formulation for the evaluation section: 1469 1470 Algorithm 3 General calculation for accuracy in ZSTE tasks 1471 1: Evaluation procedure 1472 Step 1: Audio Encoding 2: 1473 3: Encode audio inputs using the Audio Encoder \mathcal{A} to get audio embeddings \mathcal{A}_i 1474 4: Ensure the embeddings are normalized to have a unit norm to maintain consistency in 1475 comparisons 1476 5: Step 2: Similarity Calculation 1477 6: Compute similarity scores between audio embeddings A_i and text embeddings T_i using a 1478 suitable similarity metric (e.g., cosine similarity) 1479 7: Generate a similarity matrix S where S_{ij} represents the similarity between the *i*-th audio embedding and the j-th text embedding 1480 **Step 3: Probability Calculation** 8: 1481 Apply the softmax function to the similarity scores to obtain probabilities p_{ij} for each class 9: 1482 $p_{ij} = \frac{e^{\mathcal{S}_{ij}}}{\sum_{k=1}^{50} e^{\mathcal{S}_{ik}}}$ 1483 10: 1484 Step 4: Classification and Accuracy Measurement 11: 1485 12: Determine the predicted class by selecting the class with the highest probability for each 1486 audio input 13: $\hat{y}_i = \arg \max_j p_{ij}$ 1487 Measure accuracy by comparing predicted labels \hat{y}_i with ground truth labels y_i : Acc₁ = 14: 1488 $\frac{1}{|\mathcal{U}|} \sum_{i \in \mathcal{U}} \mathbf{1}(\hat{y}_i = y_i)$ 1489 15: return Evaluation scores Acc1, insights, and potential improvements 1490 1491 1492 1493 1494 1495 1496 1497 1498 1499 1500 1501 1502 1503 1507 1509 1510 1511

1512		
1513		
1514		
1515		
1516		
1517		
1510		
1510		
1519		
1520		
1521		
1522		
1523	Alg	orithm 4 Audio-Text Matching Evaluation with CLAP Model
1524	1:	Initialize CLAP model with pre-trained weights
1525	2:	Load dataset $D = \{(a_i, t_i)\}_{i=1}^N$
1526	3:	Split dataset into training, validation, and test sets
1527	4:	Prepare DataLoader for batch processing
1528	5:	Load wordsList from file
1529	6:	Set prompt as 'this is a sound of '
1530	7:	Create target texts $y = [prompt + x \text{ for } x \text{ in words_list}]$
1531	8:	function ONEHOTENCODE(text, wordsList)
1532	9:	Initialize a zero vector $oneHotVector \in \{0,1\}^{ wordsList }$
1533	10:	for each $word \in wordsList$ do
1534	11:	if text starts with word then
1535	12:	Set $oneHotVector[index of word] \rightarrow 1$
1505	13:	break
1000	14:	end if
1537	15:	end for
1538	16:	return oneHotVector
1539	17:	end function
1540	18:	for each $batch \in testLoader$ do
1541	19:	Extract audio and text samples from batch
1542	20:	Compute audio embeddings $f_{\text{audio}}(a_i)$
1543	21:	One-hot encode the text samples
1544	22:	for each t_j in text samples do
1545	23:	$oneHotVector \leftarrow OneHotEncode(t_j, wordsList)$
1546	24:	Compute text embeddings $f_{\text{text}}(oneHotVector)$
1547	25:	end for
1548	26:	Compute similarity scores $s(f_{audio}(a_i), f_{text}(OneHotEncode(t_j, wordsList)))$
1549	27:	Apply softmax to get $P(t_j a_i)$
1550	28:	Record predicted and true labels
1550	29:	end for
1551	30:	Compute accuracy: $1 \sum_{n=1}^{N} \sum_{i=1}^{n} (i - i)$
1002	31:	Accuracy = $\frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(t_i = t_i)$
1553	32:	where $\hat{t}_i = \arg \max_{t \in T} s(a_i, t)$ and \mathbb{I} is used as the indicator function.
1554	33:	return accuracy
1555		
1556		
1557		
1558		
1559		
1560		
1561		
1562		
1563		
1564		
1565		
-		