# Dense-Localizing Audio-Visual Events in Untrimmed Videos:
# A Large-Scale Benchmark and Baseline

Tiantian Geng[1,2], Teng Wang[1,3], Jinming Duan[2], Runmin Cong[4], Feng Zheng[1,5*]

[1]Southern University of Science and Technology  [2]University of Birmingham
[3]The University of Hong Kong  [4]Shandong University  [5]Peng Cheng Laboratory

gengtiantian97@gmail.com tengwang@connect.hku.hk j.duan@bham.ac.uk

rmcong@sdu.edu.cn    f.zheng@ieee.org

## Abstract

*Existing audio-visual event localization (AVE) handles manually trimmed videos with only a single instance in each of them. However, this setting is unrealistic as natural videos often contain numerous audio-visual events with different categories. To better adapt to real-life applications, in this paper we focus on the task of dense-localizing audio-visual events, which aims to jointly localize and recognize all audio-visual events occurring in an untrimmed video. The problem is challenging as it requires fine-grained audio-visual scene and context understanding. To tackle this problem, we introduce the first Untrimmed Audio-Visual (UnAV-100) dataset, which contains 10K untrimmed videos with over 30K audio-visual events. Each video has 2.8 audio-visual events on average, and the events are usually related to each other and might co-occur as in real-life scenes. Next, we formulate the task using a new learning-based framework, which is capable of fully integrating audio and visual modalities to localize audio-visual events with various lengths and capture dependencies between them in a single pass. Extensive experiments demonstrate the effectiveness of our method as well as the significance of multi-scale cross-modal perception and dependency modeling for this task.*

## 1. Introduction

Understanding real-world scenes and events is inherently a multisensory perception process for humans [16,34]. However, for machines, how to integrate multi-modal information, especially audio and visual ones, to facilitate comprehensive video understanding is still a challenging problem. In recent years, with the introduction of many audio-visual datasets [7,8,11,39], we have seen progress in
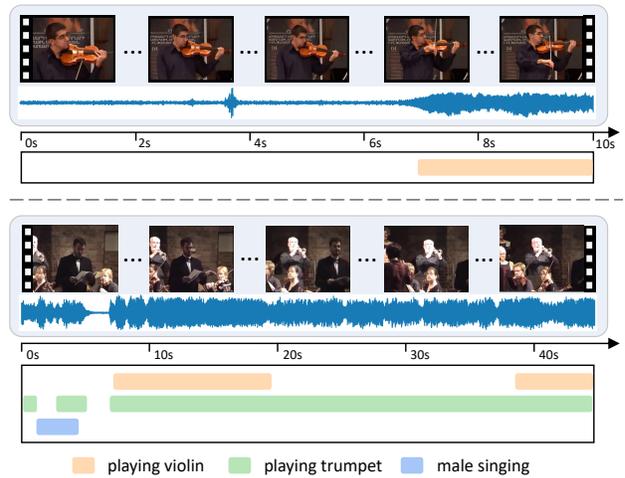
---

∗ Corresponding author



Figure 1. Different from the previous AVE task, dense-localizing audio-visual events involves localizing and recognizing all audio-visual events occurring in an untrimmed video. In real-life audio-visual scenes, there are often multiple audio-visual events that might be very short or long, and occur concurrently. The top and bottom examples are from the current AVE dataset [39] and our UnAV-100 dataset, respectively.

learning joint audio-visual representations [1, 29, 30], spatially localizing visible sound sources [7, 24] and temporally localizing audio-visual events [42, 43, 50], *etc.* While the success of these algorithms is encouraging, they all focus on manually trimmed videos that often just contain a single audio-visual instance/object in each of them. In particular, audio-visual event localization (AVE) [39] aims to localize a single event that is both audible and visible at the same time in a short, trimmed video, as shown in the upper part of Fig. 1. The task setting is impractical as a real-life video is usually long, untrimmed and associated to multiple audio-visual events from different categories, and these events might have various duration and occur simultaneously. For example, as illustrated at the bottom of Fig. 1, a man starts singing and other people accompany him on

1

trumpet and violin, and they pause several times along with the music. Therefore, we argue that it is necessary to re-examine and re-define the AVE task to better adapt to real-life audio-visual scenarios.

In this work, we conduct in-depth research starting from dataset construction to technical solutions. On the one hand, different from the existing AVE dataset [39] that only contains a single audio-visual event in each 10s trimmed video, we introduce a large-scale *Untrimmed Audio-Visual* (UnAV-100) dataset. It consists of more than 10K untrimmed videos with over 30K audio-visual events covering 100 different event categories. Our dataset spans a wide range of domains, including human activities, music performances, and sounds from animals, vehicles, tools, nature, *etc*. As the first audio-visual dataset built on untrimmed videos, UnAV-100 is quite challenging for many reasons. For instance, each video contains 2.8 audio-visual events on average (23 events maximum), and around 25% of videos have concurrent events. Besides, the length of audio-visual events varies greatly from 0.2s to 60s. There are also rich temporal dependencies among events occurring in a video, *e.g*., people often clap when cheering, and rain is usually with thunder, *etc*. We believe that the UnAV-100 dataset, with its realistic complexity, can promote the exploration on comprehensive audio-visual video understanding.

On the other hand, facing such a complex real-life scene, current methods [39, 42, 43, 45, 50] formulate the AVE task as a single-label segment-level classification problem and can only identify one audio-visual event for each segment in a trimmed video. They fail to locate concurrent events and provide an exact temporal extent for each event in untrimmed videos. To address the above issues, we re-define AVE as an instance-level localization problem, called *dense-localizing audio-visual events*. We also present a new framework to flexibly recognize all audio-visual events in an untrimmed video and meanwhile regress their temporal boundaries in a single pass. Firstly, the sound and its visual information are both critical to identify an audio-visual event, and the events can range across multiple time scales. Hence, we propose a cross-modal pyramid transformer encoder that enables the model to fully integrate informative audio and visual signals and capture both very short as well as long audio-visual events. Secondly, with the observation that the events in a video are usually related to one another, we conduct temporal dependency modeling to learn such correlations, allowing the model to use context to localize events more correctly. Finally, we design a class-aware regression head for decoding temporal boundaries of overlapping events, together with a classification head to obtain the final localization results. Extensive experiments demonstrate the effectiveness of our method, and show that it outperforms related state-of-the-art methods for untrimmed videos by a large margin.

Our contributions can be summarized as follows:
- We introduce a large-scale UnAV-100 dataset, as the first audio-visual benchmark based on untrimmed videos. There exist multiple audio-visual events in each video, and these events are usually related to one another and co-occur as in real-life scenes.
- We shift the AVE task to a more realistic setup of *dense-localizing audio-visual events*, and propose a new framework, allowing to flexibly recognize all audio-visual events in an untrimmed video and regress their temporal boundaries in a single pass.
- Extensive experiments demonstrate the significance of multi-scale cross-modal perception and dependency modeling for the task. Our method can achieve superior performance over related state-of-the-art methods for untrimmed videos by a large margin.

## 2. Related Work

### 2.1. Uni-Modal Temporal Localization Tasks

Deep learning methods have achieved promising performance in temporally localizing target instances using one modality as input, including temporal action localization (TAL) and sound event detection (SED) tasks. **Temporal action localization (TAL)** aims to detect and classify actions in untrimmed videos. It can be divided into two-stage and single-stage approaches. A two-stage TAL approach first generates action boundaries with confidence scores, and then classifies their corresponding segments into action categories and refines the generated temporal boundaries [2, 20, 21, 46]. By contrast, single-stage TAL localizes actions in a single shot without using pre-generated proposals, including anchor-based [26] and anchor-free methods [19, 47]. Besides, Transformers [41], with its powerful ability of long-range relation modeling, are recently also considered in some single-stage TAL methods [25, 36, 48]. **Sound event detection (SED)** focuses on recognizing and locating audio events in pure acoustic environments [27]. Approaches [5, 28, 31] cast it as a classification problem to classify the sound category for each temporal unit. Overall, both of them belong to uni-modal temporal localization tasks, *i.e*., TAL detects visual actions, ignoring the auditory information, while SED only considers sound tracks without utilizing visual content. Thus, they are both not beneficial for joint audio-visual scene understanding.

### 2.2. Audio-Visual Event Localization

Tian *et al*. [39] first proposed the audio-visual event localization task and introduced the AVE dataset. Afterward, Wu *et al*. [42] presented a dual attention matching module for better high-level event information modeling and also attaining local temporal cues. Xu *et al*. [45] designed a relation-aware module to build connections between visual

| Dataset | Videos | Classes | Avg. Length | Annotations | TB | ME |
|---|---|---|---|---|---|---|
| AudioSet [11] | 2.1M | 527 | 10s | A | ✗ | ✗ |
| Kinetics-Sound [1] | 19K | 34 | 10s | V | ✗ | ✗ |
| VGGSound [8] | 200K | 300 | 10s | AV | ✗ | ✗ |
| ACAV100M [18] | 100M | - | 10s | weak AV | ✗ | ✗ |
| AVE [39] | 4,143 | 28 | 10s | AV | ✓ | ✗ |
| LLP [38] | 11,849 | 25 | 10s | weak A, V | ✓ | ✓ |
| UnAV-100 (Ours) | 10,790 | 100 | 42.1s | AV | ✓ | ✓ |

Table 1. Comparison with related audio-visual datasets. A: audio events; V: visual events; AV: audio-visual events; TB: temporal boundaries; ME: multiple events.

and audio modalities. Besides, a positive sample propagation module was proposed by Zhou *et al*. [50] to adaptively aggregate positive audio-visual pairs and avoid interference of irrelevant pairs. Yan *et al*. [43] devised a background suppression scheme to suppress cross-modal asynchronous information and uni-modal noise. However, all these methods treat the task as a single-label segment classification problem, which fails to localize multiple, concurrent events in an untrimmed video. In addition, AVE [39] dataset is based on manually trimmed, short videos that only contain a single audio-visual event in each of them, which is inconsistent with real-life audio-visual scenes. On the other hand, the recent audio-visual video parsing task [38] aims to identify multiple audio, visual and audio-visual events occurring in videos. However, the methods [23, 38, 50] are also based on simple, trimmed videos in LLP dataset [38], and can be only deployed in a weekly-supervised manner.

## 3. The UnAV-100 Dataset

### 3.1. Overview

To explore audio-visual event localization in more practical scenes, we build a large-scale UnAV-100 dataset, as the first audio-visual dataset for untrimmed videos. Each video usually contains multiple audio-visual events annotated with categories and accurate temporal boundaries. The events can be very short as well as long and even overlap in time. Besides, the dataset covers a wide range of domains, including human activities, music performances, animals/vehicles/tools/natural sounds, *etc*.

The comparison with other related audio-visual datasets is shown in Tab. 1. The datasets in the top rows are mainly designed for audio-visual representation learning. They all consist of 10s short clips, and there is only a single video-level label provided for each clip. Among them, AudioSet [11] annotates videos only based on their sound without considering visual information. Kinetics-Sound [1], as a subset of Kinetics [15] for action recognition, is annotated based on visual actions, resulting that many videos contain sound tracks unrelated to the visual content (*e.g*., background music, offscreen voice). Besides, the videos in

VGGSound [8] and ACAV100M [18] have relatively good audio-visual correspondence, while they are curated using automatic algorithms leading to numerous noisy data. And ACAV100M [18] just provides weak labels obtained from pre-trained uni-modal classifiers. AVE [39], the existing dataset for audio-visual event localization, just contains 4K samples with limited 28 event classes. Each video is a trimmed 10s clip containing only one audio-visual event, and most events span over the entire video, which is seriously inconsistent with real-world scenarios. LLP [38] is designed for weakly-supervised audio-visual video parsing, only providing video-level weak labels for all training data. By contrast, our UnAV-100 dataset is based on untrimmed videos, containing over 10K samples with 100 audio-visual event categories. Moreover, there are usually multiple audio-visual events annotated in a video with their categories and precise temporal boundaries. In the following, we provide detailed descriptions about the dataset construction and statistical analysis of our UnAV-100 dataset.

### 3.2. Dataset Construction

**Collection.** We select VGGSound [8] as our data collection source for its relatively high audio-visual correspondence in videos. Specifically, we first chose the categories that are common in our daily life from 300 classes in VGGSound covering diverse domains. Then, we downloaded raw videos rather than 10s trimmed clips using the provided YouTube URLs. Since the lengths of raw videos usually span several hours, we randomly cut them within one minute to ensure reasonable duration, meanwhile keeping the videos containing the original 10s clips. Afterward, we manually verified the presence of audio-visual events in each obtained video. We found that, since VGGSound was collected in an automated manner, there exist numerous videos that do not contain any audio-visual events. For instance, some videos have correct visual content with unrelated sounds like background music and narrations. Additionally, it also contains low-quality and animated videos with unrecognizable events. Finally, by filtering the above cases, we selected around 10K from downloaded 15K videos for annotation.

**Annotation.** We annotated videos via an open-source annotation tool VIA [10] by crowdsourcing. Specifically, we provided expert annotators with a category list for reference, and all audio-visual events occurring in videos are required to be annotated with their categories and independent start and end timestamps. Different from the temporal boundaries of visual contents that are usually ambiguous [4], the start and end time points of an audio-visual event are usually clearer and can be easily identified by judging if the event occurs in both audio and visual channels. Thus, there is usually high agreement among annotators in labeling temporal boundaries. In order to ensure
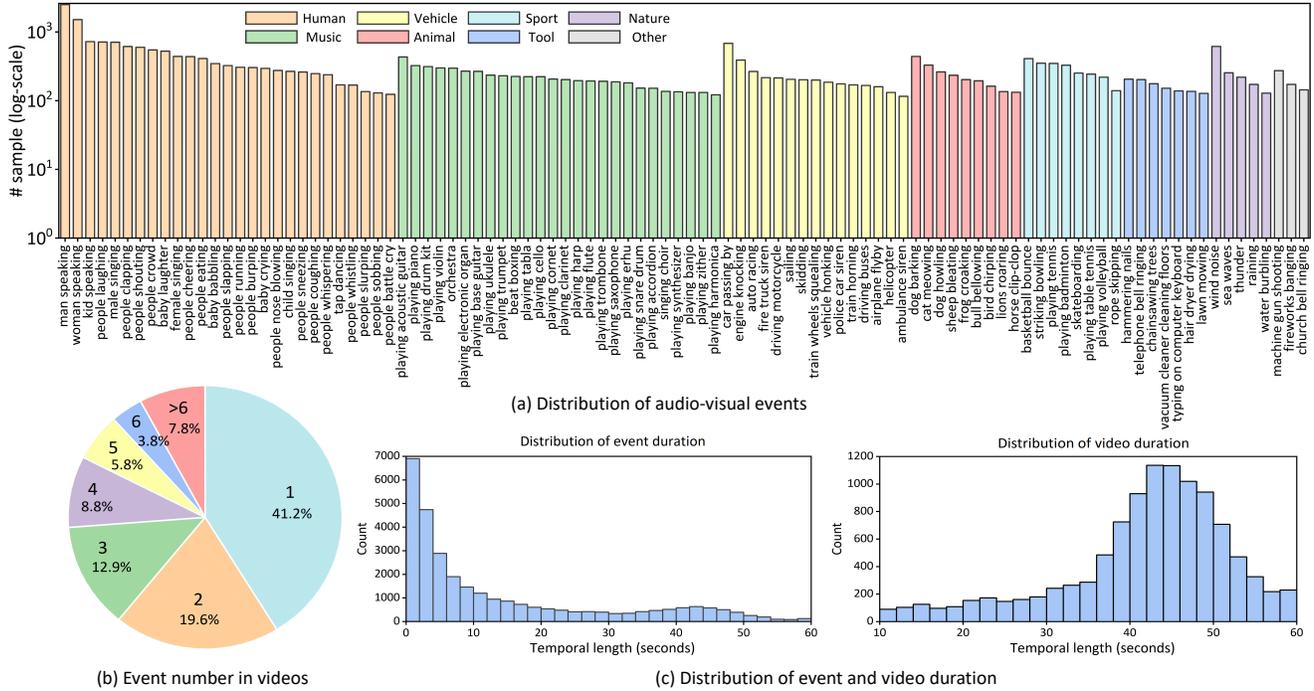
Figure 2. Illustrations of statistics on our UnAV-100 dataset. (a) Distribution of audio-visual events. Bars are grouped by domains, and different colors mean that event categories belong to different domains. (b) The number of audio-visual events in videos. (c) Distribution of event (left) and video duration (right).

annotation quality high, the labeling team is required to check all annotated data carefully. We also employed another group of crowdworkers to manually check it again, resulting in a very time-consuming process.

### 3.3. Statistical Analysis

Overall, our UnAV-100 dataset contains 30,059 audio-visual events of 100 categories, distributed in 10,790 untrimmed videos for over 126 video hours. The dataset is split into training, validation, and testing sets with a ratio of 3:1:1, where a multi-label split strategy [35] is applied to ensure a well-balanced data distribution in subsets. Besides, in order to alleviate the effect of long tails, we make sure that there are more than 116 audio-visual events for each category. Fig. 2 provides the statistics of our dataset, and the challenges of UnAV-100 include the following:

**1) Multiple events in videos.** As shown in Fig. 2(b), around 60% of videos contain more than one audio-visual event. Each video has 2.8 audio-visual events on average (1.6 for distinct ones), and the maximum number is 23. Besides, about 25% of videos have concurrent events (the details are in the *Supp. Materials*), which means that there is more than one visible sound source at the same time.

**2) Various lengths of events and videos.** Fig. 2(c) shows that a large number of events have very short duration, with the shortest being only 0.2s. Short events are often difficult to detect, but it aligns with real-life scenes. For example,
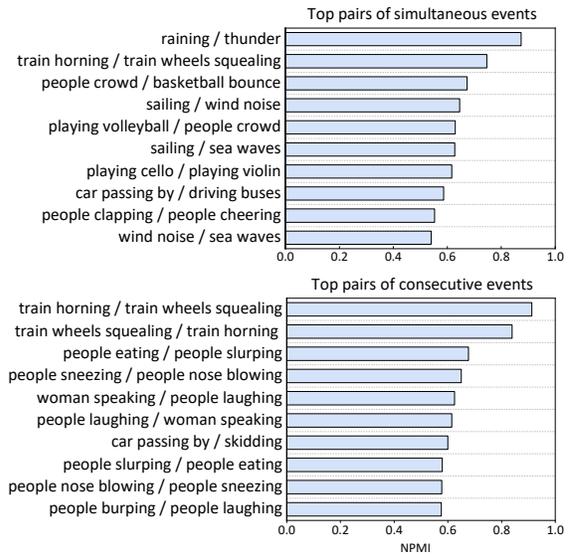


Figure 3. Top pairs of simultaneous and consecutive audio-visual events computed by NPMI falling in the range (-1, 1].

*dog barking, basketball bounce*, and *fireworks banging* are normally very short audio-visual events. Besides, the average lengths of audio-visual events and videos are 13.9s and 42.1s, respectively.

**3) Rich temporal dependencies between events.** The related audio-visual events usually occur simultaneously or consecutively in a video. In Fig. 3, we show the pairs of
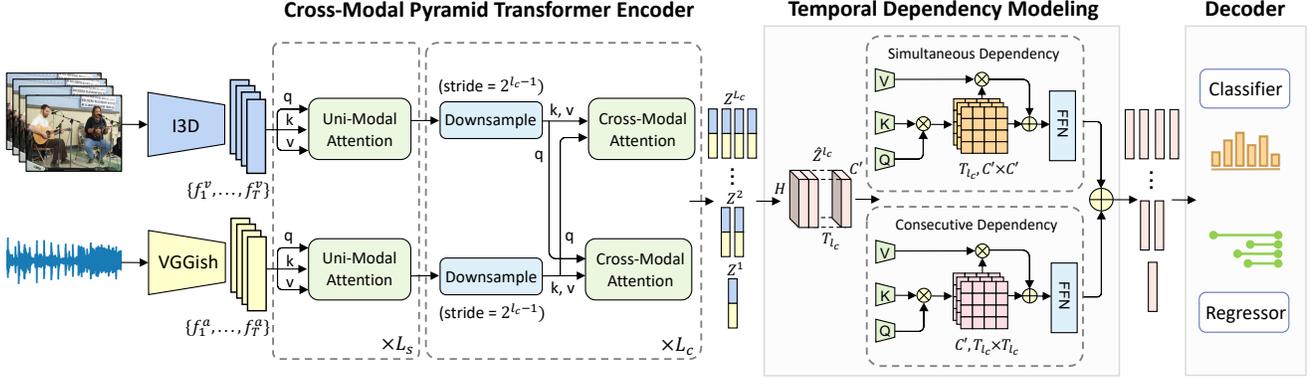
Figure 4. The overview of our proposed framework for dense-localizing audio-visual events. The model takes pre-trained CNNs to extract audio and visual features, and uses a cross-modal pyramid transformer encoder to encode and fuse cross-modal features at various temporal scales, which consists of $L_s$ uni-modal and $L_c$ cross-modal transformer blocks. Then, the temporal dependency modeling is conducted to capture correlations between events. Finally, a classification and a regression head are used to predict the categories and temporal boundaries of events in an end-to-end manner.

simultaneous and consecutive events with the top 10 Normalized Pointwise Mutual Information (NPMI) [9], respectively. We can see, frequently, rain is accompanied by thunder, violin and cello are played together, people clap when cheering, *etc*. Such event dependencies are very similar to real-world intuition and reflect human behavior. We note that UnAV-100 is the first audio-visual dataset with such context information, which provides excellent data for building many complex models for audio-visual event dependency modeling.

## 4. Method

To solve the problem of dense-localizing audio-visual events, we design an architecture to jointly recognize and localize multiple, concurrent audio-visual events with various lengths, and meanwhile capture event dependencies in an untrimmed video. An overview of the proposed framework is illustrated in Fig. 4.

### 4.1. Preliminaries

**Problem Statement.** Different from previous AVE methods, we formulate the task of dense-localizing audio-visual events as a joint classification and regression problem. Formally, given an input video sequence containing both visual and audio tracks, we first divide it into $T$ visual and audio segment pairs $\{V_t, A_t\}_{t=1}^T$, where $T$ varies across videos. The groundtruth event set for each video is denoted as $Y = \{y_n = (t_{s,n}, t_{e,n}, c_n)\}_{n=1}^N$, where $t_{s,n}, t_{e,n}$ are the start and end timestamp of $n$-th event, $c_n \in \{1, \cdots, C\}$ is the event category, and $N$ is the total number of audio-visual events in the video. Then, the model is required to predict $\hat{Y} = \{\hat{y}_t = (d_{s,t}, d_{e,t}, p(c_t))\}_{t=1}^T$ during inference, where $p(c_t) \in \mathbb{R}^{1 \times C}$ is the probabilities of $C$ event categories at moment $t$, $d_{s,t}$ and $d_{e,t}$ are the distances between the moment $t$ to the event's start and end timestamp. Note

that $d_{s,t}$ and $d_{e,t}$ are only defined when an event presents at moment $t$. Thus, the final localization results can be obtained by:

$$c_t = \arg\max p(c_t), \quad t_{s,t} = t - d_{s,t}, \quad t_{e,t} = t + d_{e,t}. \quad (1)$$

**Audio and Visual Representations.** We extract audio feature vectors using the VGGish model [13] pre-trained on AudioSet [11]. And visual feature vectors are extracted by the two-stream I3D [6] pre-trained on Kinetics-400 [15]. Then, we apply two convolutional layers with ReLU to project features from two modalities into a shared embedding space, resulting $F_V = \{f_t^v\}_{t=1}^T$, $F_A = \{f_t^a\}_{t=1}^T \in \mathbb{R}^{T \times D}$, where $D$ is the dimension of the embedding space.

### 4.2. Architecture

**Cross-Modal Pyramid Transformer Encoder.** We consider that the sound and its corresponding visual information are both crucial to identify an audio-visual event. However, the audio and visual tracks of an untrimmed video often contain a lot of irrelevant information (*e.g.*, background music and off-screen voice), and their content might be misaligned with each other (*e.g.*, a dog appears without barking). Besides, the events occurring in untrimmed videos usually range across multiple time scales. Thus, how to appropriately integrate the two modalities and capture very short as well as long events are both significant for this task. Here, a cross-modal pyramid transformer encoder is proposed to address the above challenges.

Specifically, in order to capture long-term temporal relations among uni-modal segments and filter out noise in each modality, the feature sequences from two modalities are first fed into $L_s$ stacked uni-modal transformer blocks separately. Each block regularly contains a multiheaded self-attention (MSA) and a feed-forward network (FFN) with LayerNorm (LN) and residual connections. And position

embeddings $E_{pos} \in \mathbb{R}^{T \times D}$ as in [41] are also added in input sequences. By doing this, the model can focus more on event-related information in each modality. Afterward, the obtained feature sequences are further encoded into a cross-modal pyramid transformer to integrate informative signals from two modalities at different temporal resolutions. The module consists of $L_c$ stacked blocks. In each block, as shown in Fig. 4, we first temporally downsample the feature sequence of each modality with the stride $2^{l_c-1}$, where $l_c$ is the index of the current block, and the longer strides are able to capture longer events. Then, we assign downsampled features in the current modality as the key and value vectors, and the features of another modality as the query vector in multiheaded cross-attention (MCA), followed by FFN and LN layers. Thus, the audio-guided visual feature $F_{Va}$ and visual-guided audio feature $F_{Av}$ from $l_c$-th block can be denoted as:

$$
\begin{aligned}
F_{Va}^{l_c} &= \text{MCA}(\hat{F}_{Av}^{l_c-1}W_q, \hat{F}_{Va}^{l_c-1}W_k, \hat{F}_{Va}^{l_c-1}W_v), \\
F_{Av}^{l_c} &= \text{MCA}(\hat{F}_{Va}^{l_c-1}W_q, \hat{F}_{Av}^{l_c-1}W_k, \hat{F}_{Av}^{l_c-1}W_v),
\end{aligned}
\quad (2)
$$

where $l_c = \{1, \cdots, L_c\}$, $F_{Va}^{l_c}, F_{Av}^{l_c} \in \mathbb{R}^{T_{l_c} \times D}$ ($T_{l_c} = T/2^{l_c-1}$), $\hat{F}_{Va}^{l_c-1}$ and $\hat{F}_{Av}^{l_c-1}$ are the features after downsampling, $W_q, W_k, W_v \in \mathbb{R}^{D \times D_m}$ are learnable parameters and $D_m = D$ is the dimension of learned query, key and value vectors. After cross-modal interactions at various temporal scales, we concatenate the enhanced audio and visual features at the same pyramid level, getting a cross-modal feature pyramid $Z = \{Z^{l_c}\}_{l_c=1}^{L_c}$, where $Z^{l_c} = \text{Concat}(F_{Va}^{l_c}, F_{Av}^{l_c}) \in \mathbb{R}^{T_{l_c} \times 2D}$.

**Temporal Dependency Modeling.** The key characteristic of real-life audio-visual scenes is that the related events usually occur simultaneously or consecutively. For example, people are used to clapping when cheering, and cars often honk when passing by. Here, inspired by the method [40] for action dependency modeling in the TAL task, we implicitly capture such simultaneous and consecutive dependencies among audio-visual events at the obtained cross-modal feature pyramid. Concretely, for each cross-modal feature sequence $Z^{l_c}$, we first transform and expand the feature dimension to $\hat{Z}^{l_c} \in \mathbb{R}^{T_{l_c} \times C' \times H}$, splitting it into $C'$ groups, where $C'$ represents the number of hidden classes and $H$ is the transformed feature dimension. We suppose each hidden class is learned to carry a group of distinctive features for event classification. For simultaneous dependency modeling, the self-attention is performed along the $C'$ dimension of $\hat{Z}^{l_c}$, which means a $C' \times C'$ attention matrix that denotes the relevance among hidden classes at each time step can be obtained. For consecutive dependency modeling, the self-attention is performed along $T_{l_c}$ dimension, getting a $T_{l_c} \times T_{l_c}$ attention matrix to indicate the correlations among all time steps for the classification of the given class. Then the output of the two branches followed by FFN

and LN layers with residual connections are simply merged by element-wise summation to enable the model to capture both types of dependencies. Note that we share the parameters of dependency modeling across all pyramid levels.

**Decoder.** Next, a decoder, consisting of a classification head and a regression head, is applied to decode the enhanced feature pyramid into prediction results in a single pass. Specifically, the classification head predicts the probability $p(c_t)$ of events at every moment $t$ of all pyramid levels. It consists of three layers of 1D convolutions following a sigmoid function as in [48]. Besides, the regression head outputs the distances to the start and end timestamp of an event $(d_{s,t}, d_{e,t})$ at time step $t$ if the event exists. We highlight that the regression head is designed to be class-aware, which allows the model to regress temporal boundaries for the overlapping events with different categories. It is realized by using three 1D convolutions attached with a ReLU, getting the output with the shape of $[2, C, T_{l_c}]$ for each pyramid level. Here, the pyramid architecture enables the regression head to predict temporal boundaries at different temporal scales, allowing the model to capture the events with various lengths. Note that the parameters of both two heads are shared across all pyramid levels.

### 4.3. Training and Inference

**Loss Function.** We use two losses to train our model in an end-to-end manner, *i.e.*, a focal loss [22] $\mathcal{L}_{cls}$ for classification and a generalized IoU loss [32] $\mathcal{L}_{reg}$ for distance regression, as in the TAL method [48]. For each video, the loss function is denoted as:

$$
\mathcal{L} = \frac{1}{\mathcal{T}} \sum_t \mathcal{L}_{cls} + \frac{\lambda}{\mathcal{N}} \sum_t \mathbb{I}_t \mathcal{L}_{reg}, \quad (3)
$$

where $\mathcal{T}$ is the total segment number of all levels, $\mathbb{I}_t$ is an indicator function denoting if a timestamp contains events, $\mathcal{N}$ is the number of positive segments that contain events across all levels. Here, we weight the contribution of $\mathcal{L}_{reg}$ with $\lambda = 1$ by default.

**Inference.** During inference, the outputs of the model are as in Eq. (1) for every timestamp $t$ across all levels. Then the obtained event candidates are post-processed by a multi-class version of Soft-NMS [3] to suppress redundant temporal boundaries with high overlaps within the same class.

## 5. Experiments

### 5.1. Experimental Settings

**Implementation Details.** For each video, we sample frames at 25 fps, and feed 24 consecutive RGB and optical flow frames into two-stream I3D [6], using a sliding window with stride 8. Then, the two-stream features are concatenated (2048-d) as a visual segment. Here, the optical flow is extracted by RAFT [37]. Besides, we extract 128-d

| Modality | Method | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | Avg. |
|---|---|---|---|---|---|---|---|
| A | VSGN [49] | 18.0 | 14.2 | 10.8 | 8.2 | 5.3 | 17.8 |
| | TadTR [25] | 23.0 | 20.5 | 17.6 | 14.4 | 10.4 | 22.8 |
| | ActionFormer [48] | 37.7 | 32.8 | 27.3 | 22.5 | 15.6 | 36.0 |
| | Ours | 39.0 | 34.5 | 29.1 | 23.3 | 12.5 | 37.1 |
| V | VSGN [49] | 14.8 | 11.5 | 8.5 | 6.0 | 4.1 | 15.5 |
| | TadTR [25] | 23.1 | 20.5 | 17.8 | 15.3 | 12.0 | 23.0 |
| | ActionFormer [48] | 36.3 | 31.9 | 27.4 | 21.8 | 14.8 | 35.4 |
| | Ours | 37.3 | 32.6 | 28.3 | 22.9 | 14.7 | 35.9 |
| A&V | VSGN [49] | 24.5 | 20.2 | 15.9 | 11.4 | 6.8 | 24.1 |
| | TadTR [25] | 30.4 | 27.1 | 23.3 | 19.4 | 14.3 | 29.4 |
| | ActionFormer [48] | 43.5 | 39.4 | 33.4 | 27.3 | 17.9 | 42.2 |
| | Ours | **50.6** | **45.8** | **39.8** | **32.4** | **21.1** | **47.8** |

Table 2. Comparison of the results on the test set of UnAV-100 dataset. A: only audio modality; V: only visual modality; A&V: both audio and visual modalities.

audio features by VGGish [13] for each 0.96s segment with a sliding window (stride=0.32s) to temporally align with the visual ones. Since the input sequences vary in length, we pad or crop them to the maximum length $T = 224$, and add masks for all operations in the model. The dimensions of the embedding space in the encoder and temporal dependency modeling are $D = 512$ and $H = 128$, respectively. The number of hidden classes $C' = 100$. Our model is trained with the Adam optimizer, and the number of epochs is 40 with a linear warmup of 5 epochs. The initial learning rate is 1e-4 and a cosine learning rate decay is used. The mini-batch size is 16 and the weight decay is 1e-4.

**Evaluation Metrics.** As a temporal localization task for untrimmed videos, we use mean Average Precision (mAP) to evaluate results. Specifically, we report mAPs at the tIoU thresholds [0.5:0.1:0.9] and the average mAP at the thresholds [0.1:0.1:0.9].

**Baseline Models.** Since previous AVE and SED methods are limited to localizing a single event on trimmed videos with the same duration and cannot be applied on untrimmed videos, we only compare our model with recent state-of-the-art TAL models, as shown in Tab. 2. It includes the two-stage model VSGN [49] and single-stage models (TadTR [25] and ActionFormer [48]). Here, $L_s = 2$ and $L_c = 6$ in the pyramid transformer encoder of our model. Note that all compared approaches use the same input features as ours to keep a fair comparison.

## 5.2. Results and Analysis

To validate the effectiveness of the proposed model, we compare it with recent TAL methods using different modalities as input, and also conduct extensive ablation studies.

**Comparison Results.** As shown in Tab. 2, when using one modality as input, our model variants that only apply self-attention in the encoder outperform all compared TAL methods, where TadTR [25] and ActionFormer [48] also use an end-to-end transformer-based architecture. When using both audio and visual modalities, the performance of

| $L_s$ | $L_c$ | TD | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | Avg. |
|---|---|---|---|---|---|---|---|---|
| 2 | 0 | | 36.8 | 29.3 | 21.8 | 13.8 | 4.9 | 35.5 |
| 2 | 1 | | 37.6 | 29.6 | 22.2 | 14.0 | 5.1 | 35.4 |
| 2 | 2 | | 37.0 | 29.3 | 20.9 | 12.5 | 3.8 | 35.3 |
| 2 | 2 | ✓ | 41.0 | 33.1 | 25.7 | 18.0 | 8.1 | 39.4 |
| 2 | 4 | ✓ | 49.8 | 43.0 | 35.4 | 25.5 | 11.2 | 45.0 |
| 2 | 6 | ✓ | 50.6 | 45.8 | 39.8 | 32.4 | 21.1 | **47.8** |
| 2 | 7 | ✓ | 49.1 | 44.8 | 39.5 | 32.4 | **21.8** | 46.8 |
| 0 | 6 | ✓ | 49.4 | 45.2 | 39.2 | **32.5** | 21.6 | 46.7 |
| 1 | 6 | ✓ | 49.6 | 45.3 | 39.5 | **32.5** | 21.3 | 47.0 |
| 3 | 6 | ✓ | 48.8 | 44.4 | 39.2 | 32.2 | **21.8** | 46.4 |

Table 3. Ablation study on cross-modal fusion strategies and the design of feature pyramid. TD: temporal downsampling.

| DM | CA | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | Avg. |
|---|---|---|---|---|---|---|---|
| | | 48.2 | 42.3 | 35.5 | 28.0 | 18.1 | 45.2 |
| ✓ | | 48.5 | 44.2 | 38.7 | **32.6** | 21.0 | 46.1 |
| | ✓ | 48.5 | 43.4 | 36.9 | 29.9 | 20.2 | 45.8 |
| ✓ | ✓ | **50.6** | **45.8** | **39.8** | 32.4 | **21.1** | **47.8** |

Table 4. Ablation study on dependency modeling (DM) and class-aware regression (CA).

| Model | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | Avg. |
|---|---|---|---|---|---|---|
| ResNet50 [12] (RGB) | 46.6 | 42.2 | 37.2 | 30.8 | 20.1 | 44.3 |
| I3D [6] (RGB) | 49.1 | 44.8 | 39.0 | 32.0 | **21.3** | 46.7 |
| I3D [6] (RGB + Flow) | **50.6** | **45.8** | **39.8** | **32.4** | 21.1 | **47.8** |

Table 5. Ablation study on different visual features.

our model boosts significantly, e.g., +11.9% and +10.7% at the average mAP compared with our visual-only and audio-only variants, respectively. These results clearly indicate that both modalities are equally crucial for this task. Besides, our model surpasses the compared TAL methods by a large margin, even though they also benefit greatly from multi-modal input. Here, we simply concatenate audio and visual features as input of these methods.

**Cross-Modal Fusion and Pyramid Levels.** We explore the cross-modal fusion strategies and the design of the cross-modal feature pyramid. In Tab. 3, we can see that using only two uni-modal transformer blocks ($L_s = 2$ and $L_c = 0$) for each modality separately decreases the performance dramatically. Later, adding one or two cross-modal blocks at the original temporal resolution can just slightly increase mAP scores. Instead, applying temporal downsampling in cross-modal blocks boosts the performance, indicating that the cross-modal fusion at multiple temporal resolutions is essential for our model. Then, the performance gradually increases by further adding cross-modal pyramid levels, and yet is saturated when $L_c = 6$. In addition, we found that the appropriate number of uni-modal blocks is also important, which reveals that applying self-attention before cross-modal interaction can help the model to focus on informative signals and eliminate noise from each modality.
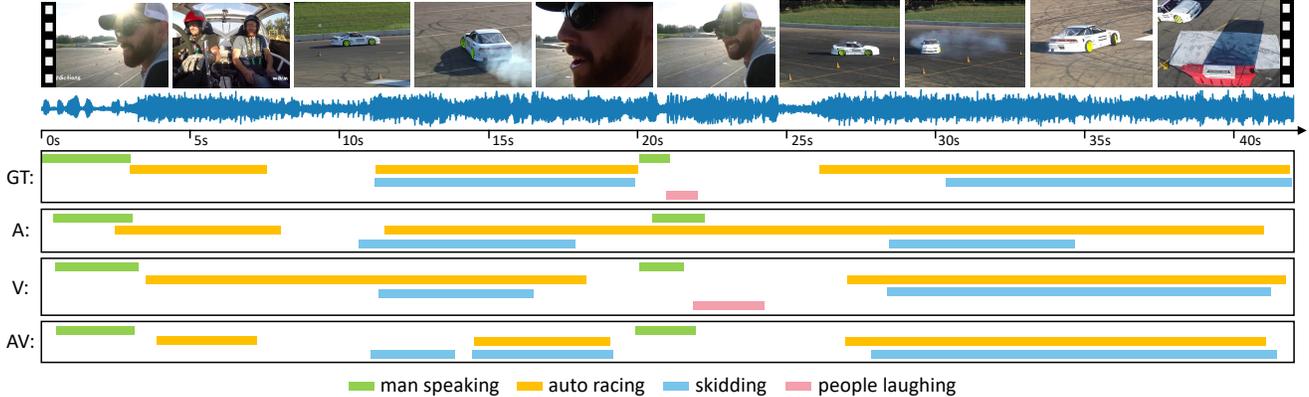
Figure 5. Qualitative results on the UnAV-100 test set. GT: ground truth, A: the prediction of our audio-only variant, V: the prediction of our visual-only variant, AV: the prediction of our audio-visual model. We show boundaries with the highest overlap with ground truth.
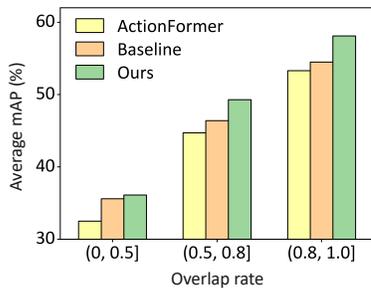


Figure 6. Performance comparison of our models and the TAL method (ActionFormer [48]) on the videos from the UnAV-100 test set, containing concurrent events with different overlap rates.

**Dependency Modeling and Class-Aware Regression.** As shown in Tab. 4, applying temporal dependency modeling and class-aware regression separately can both achieve higher performances than the base model that just contains our transformer encoder with a class-agnostic regression head in the decoder. Besides, we found that when using both of them, they can promote each other and achieve a further significant performance boost, which clearly demonstrates their effectiveness.

**The Impact of Motion Features.** In Tab. 5, we observe that utilizing both RGB and optical flow features extracted by I3D [6] achieves the best performance. It outperforms the model that uses visual features extracted by ResNet50 [12] pre-trained on ImageNet by a large margin (+3.5% at the average mAP). Even though it is proved in [39] that motion features are useless for audio-visual event localization, we argue that our experiment clearly demonstrates their significance for dense-localizing audio-visual events.

**The Capability of Localizing Concurrent Events.** We further evaluate our models and the state-of-the-art TAL method [48] on the videos that contain concurrent events with different overlap rates in Fig. 6. We observe that our model equipped with dependency modeling and class-aware regression obviously gains more performance improvement on the videos with higher overlap rates, compared with

our baseline and ActionFormer [48]. It suggests that our model has a better ability to localize overlapping audio-visual events in untrimmed videos.

**Qualitative Results.** In Fig. 5, we present the qualitative results of our model variants that utilize different modalities as input. We observe that the model using both modalities can localize audio-visual events more correctly, even though some events occur simultaneously or have short duration. By contrast, since the sound of *auto racing* almost spans the whole video, the audio-only model gets the wrong boundaries of the event without the help of visual information. And similar errors also occur when using the visual-only model. Overall, it demonstrates again that audio and visual modalities complement each other and are equally significant for dense-localizing audio-visual events. More ablation studies and qualitative results can be found in the *Supp. Materials*.

## 6. Conclusion

In this work, we investigate the dense-localizing audio-visual events problem, which aims to recognize and localize all audio-visual events occurring in an untrimmed video. To facilitate this research, we build a large-scale UnAV-100 dataset consisting of more than 10K untrimmed videos with over 30K audio-visual events covering 100 categories. We also propose a new framework, formulating the task as a joint classification and regression problem, which is capable of localizing audio-visual events that have various lengths and overlap in time, and capturing the dependencies between them in a video. Our results demonstrate the superiority of our model, indicating the significance of cross-modal perception and dependency modeling for this task.

8

# References

[1] Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 609–617, 2017. 1, 3

[2] Yueran Bai, Yingying Wang, Yunhai Tong, Yang Yang, Qiyue Liu, and Junhui Liu. Boundary content graph neural network for temporal action proposal generation. In *European Conference on Computer Vision*, pages 121–137. Springer, 2020. 2

[3] Navaneeth Bodla, Bharat Singh, Rama Chellappa, and Larry S Davis. Soft-nms–improving object detection with one line of code. In *Proceedings of the IEEE international conference on computer vision*, pages 5561–5569, 2017. 6

[4] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the ieee conference on computer vision and pattern recognition*, pages 961–970, 2015. 3, 11

[5] Emre Cakir, Toni Heittola, Heikki Huttunen, and Tuomas Virtanen. Polyphonic sound event detection using multi label deep neural networks. In *2015 international joint conference on neural networks (IJCNN)*, pages 1–7. IEEE, 2015. 2

[6] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017. 5, 6, 7, 8, 11

[7] 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 (CVPR)*, pages 16867–16876, June 2021. 1

[8] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 721–725. IEEE, 2020. 1, 3, 13

[9] Kenneth Church and Patrick Hanks. Word association norms, mutual information, and lexicography. *Computational linguistics*, 16(1):22–29, 1990. 5, 11

[10] Abhishek Dutta and Andrew Zisserman. The VIA annotation software for images, audio and video. In *Proceedings of the 27th ACM International Conference on Multimedia*, MM '19, New York, NY, USA, 2019. ACM. 3

[11] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 776–780. IEEE, 2017. 1, 3, 5

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 7, 8

[13] Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. Cnn architectures for large-scale audio classification. In *2017 ieee international conference on acoustics, speech and signal processing (icassp)*, pages 131–135. IEEE, 2017. 5, 7, 11

[14] Haroon Idrees, Amir R Zamir, Yu-Gang Jiang, Alex Gorban, Ivan Laptev, Rahul Sukthankar, and Mubarak Shah. The thumos challenge on action recognition for videos "in the wild". *Computer Vision and Image Understanding*, 155:1–23, 2017. 11, 12

[15] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 3, 5

[16] Christoph Kayser and Ladan Shams. Multisensory causal inference in the brain. *PLoS biology*, 13(2):e1002075, 2015. 1

[17] Hilde Kuehne, Ali Arslan, and Thomas Serre. The language of actions: Recovering the syntax and semantics of goal-directed human activities. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 780–787, 2014. 11

[18] Sangho Lee, Jiwan Chung, Youngjae Yu, Gunhee Kim, Thomas Breuel, Gal Chechik, and Yale Song. Acav100m: Automatic curation of large-scale datasets for audio-visual video representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10274–10284, 2021. 3

[19] Chuming Lin, Chengming Xu, Donghao Luo, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Yanwei Fu. Learning salient boundary feature for anchor-free temporal action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3320–3329, 2021. 2

[20] Tianwei Lin, Xiao Liu, Xin Li, Errui Ding, and Shilei Wen. Bmn: Boundary-matching network for temporal action proposal generation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3889–3898, 2019. 2

[21] Tianwei Lin, Xu Zhao, Haisheng Su, Chongjing Wang, and Ming Yang. Bsn: Boundary sensitive network for temporal action proposal generation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018. 2

[22] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 6

[23] Yan-Bo Lin, Hung-Yu Tseng, Hsin-Ying Lee, Yen-Yu Lin, and Ming-Hsuan Yang. Exploring cross-video and cross-modality signals for weakly-supervised audio-visual video parsing. *Advances in Neural Information Processing Systems*, 34:11449–11461, 2021. 3

[24] Xian Liu, Rui Qian, Hang Zhou, Di Hu, Weiyao Lin, Ziwei Liu, Bolei Zhou, and Xiaowei Zhou. Visual sound localization in the wild by cross-modal interference erasing. *arXiv preprint arXiv:2202.06406*, 2, 2022. 1

[25] Xiaolong Liu, Qimeng Wang, Yao Hu, Xu Tang, Shiwei Zhang, Song Bai, and Xiang Bai. End-to-end temporal action detection with transformer. *IEEE Transactions on Image Processing*, 31:5427–5441, 2022. 2, 7

[26] Fuchen Long, Ting Yao, Zhaofan Qiu, Xinmei Tian, Jiebo Luo, and Tao Mei. Gaussian temporal awareness networks for action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 344–353, 2019. 2

[27] Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen. Metrics for polyphonic sound event detection. *Applied Sciences*, 6(6):162, 2016. 2

[28] Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen. Tut database for acoustic scene classification and sound event detection. In *2016 24th European Signal Processing Conference (EUSIPCO)*, pages 1128–1132. IEEE, 2016. 2

[29] Arsha Nagrani, Shan Yang, Anurag Arnab, Aren Jansen, Cordelia Schmid, and Chen Sun. Attention bottlenecks for multimodal fusion. *Advances in Neural Information Processing Systems*, 34:14200–14213, 2021. 1

[30] Andrew Owens and Alexei A. Efros. Audio-visual scene analysis with self-supervised multisensory features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. 1

[31] Giambattista Parascandolo, Heikki Huttunen, and Tuomas Virtanen. Recurrent neural networks for polyphonic sound event detection in real life recordings. In *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 6440–6444. IEEE, 2016. 2

[32] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 658–666, 2019. 6

[33] Gunnar A Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pages 510–526. Springer, 2016. 11

[34] Charles Spence. Audiovisual multisensory integration. *Acoustical science and technology*, 28(2):61–70, 2007. 1

[35] P. Szymański and T. Kajdanowicz. A scikit-based Python environment for performing multi-label classification. *ArXiv e-prints*, Feb. 2017. 4

[36] Jing Tan, Jiaqi Tang, Limin Wang, and Gangshan Wu. Relaxed transformer decoders for direct action proposal generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13526–13535, 2021. 2

[37] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *European conference on computer vision*, pages 402–419. Springer, 2020. 6

[38] Yapeng Tian, Dingzeyu Li, and Chenliang Xu. Unified multisensory perception: Weakly-supervised audio-visual video parsing. In *European Conference on Computer Vision*, pages 436–454. Springer, 2020. 3

[39] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization in unconstrained videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 247–263, 2018. 1, 2, 3, 8

[40] Praveen Tirupattur, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. Modeling multi-label action dependencies for temporal action localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1460–1470, 2021. 6

[41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2, 6

[42] Yu Wu, Linchao Zhu, Yan Yan, and Yi Yang. Dual attention matching for audio-visual event localization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6292–6300, 2019. 1, 2

[43] Yan Xia and Zhou Zhao. Cross-modal background suppression for audio-visual event localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19989–19998, 2022. 1, 2, 3

[44] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021. 12

[45] Haoming Xu, Runhao Zeng, Qingyao Wu, Mingkui Tan, and Chuang Gan. Cross-modal relation-aware networks for audio-visual event localization. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 3893–3901, 2020. 2

[46] Mengmeng Xu, Chen Zhao, David S Rojas, Ali Thabet, and Bernard Ghanem. G-tad: Sub-graph localization for temporal action detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10156–10165, 2020. 2

[47] Le Yang, Houwen Peng, Dingwen Zhang, Jianlong Fu, and Junwei Han. Revisiting anchor mechanisms for temporal action localization. *IEEE Transactions on Image Processing*, 29:8535–8548, 2020. 2

[48] Chenlin Zhang, Jianxin Wu, and Yin Li. Actionformer: Localizing moments of actions with transformers. In *European Conference on Computer Vision*, pages 492–510. Springer, 2022. 2, 6, 7, 8, 12

[49] Chen Zhao, Ali K Thabet, and Bernard Ghanem. Video self-stitching graph network for temporal action localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13658–13667, 2021. 7

[50] Jinxing Zhou, Liang Zheng, Yiran Zhong, Shijie Hao, and Meng Wang. Positive sample propagation along the audio-visual event line. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8436–8444, 2021. 1, 2, 3

# Appendix

## A. More Statistical Analysis

**Concurrent Events.** There are usually multiple audio-visual events occurring simultaneously in UnAV-100 dataset as in real-life scenes. Here, we define the overlap rate $\mathcal{O}$ of each video as:

$$\mathcal{O} = \frac{U_o}{U_e}, \tag{4}$$

where $U_o$ is the temporal union of overlapping intervals, and $U_e$ is the temporal union of the intervals of all audio-visual events in the video. Totally, there are around 25% of videos (2,651) containing concurrent audio-visual events ($\mathcal{O} > 0.01$, considering annotation errors) in our UnAV-100 dataset. The overlap rate distribution of these videos is illustrated in Fig. 7. We can see that the videos with low and high overlap rates both have high proportions. Higher overlap rates might indicate that the events have higher correlations and usually occur at the same time, which requires the model to have a strong ability of dependency modeling.

**Temporal Dependencies between Events.** We show NPMI (Normalized Pointwise Mutual Information) [9] of the pairs of simultaneous and consecutive audio-visual events for all 100 event categories in Fig. 8(a) and Fig 8(b), respectively. NPMI is commonly used in linguistics to represent the co-occurrence between two words. Firstly, in Fig. 8(a), we can observe that the event categories from the same domains are more likely to occur concurrently, *e.g.*, the events of human activities, music performances, and the sounds of vehicles/natural. Besides, the events from various domains are usually accompanied by human activities, *e.g.*, *playing acoustic guitar* with *male singing*, *basketball bounce* with *people crowd*, *etc*. Secondly, in Fig 8(b), in addition to the NPMI for consecutive occurrences of different audio-visual events, we also compute the values for the events from the same categories, which might be larger than 1. It can be observed that the same events tend to occur repetitively in a video, especially for some events that usually happen in a short period of time, such as *people nose blowing*, *people sneezing* and *basketball bounce*, *etc*. Moreover, diverse consecutive dependencies also exist between different audio-visual events.

**Comparison with Existing TAL Datasets.** In Tab. 6, we compare our UnAV-100 dataset with four popular benchmarks for temporal action localization. All these datasets are based on untrimmed videos and have relatively small scales, since annotating temporal boundaries for all instances in videos is labor-intensive and time-consuming. Our UnAV-100 is the only dataset that combines both audio and visual signals to annotate instances, while others just utilize visual content in videos. Their audio tracks are usually very noisy and unrelated to the visual information,
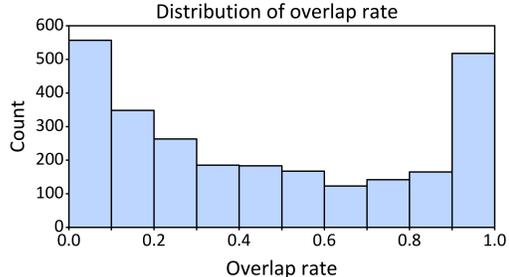


Figure 7. Overlap rate distribution of the videos that contain concurrent events in our UnAV-100 dataset.

| Dataset | Videos | Classes | Avg. Length | Avg. Instances | Domains |
|---|---|---|---|---|---|
| Breakfast [17] | 1,712 | 48 | 162s | 6 | Cooking |
| THUMOS14 [14] | 413 | 20 | 212s | 15.5 | Sports |
| ActivityNet [4] | 19,994 | 200 | 115s | 1.5 | Human Activities |
| Charades [33] | 9,848 | 157 | 30s | 6.8 | Daily Activities |
| UnAV-100 (ours) | 10,790 | 100 | 42s | 2.8 | Unconstrained |

Table 6. Comparison with temporal action localization datasets based on untrimmed videos.

*e.g.*, background music and narrations, thus these datasets are not suitable for joint audio-visual video understanding. Besides, these benchmarks all focus on specific domains, such as human activities, sports, cooking, *etc*. By contrast, our UnAV-100 covers many different domains including human/music/sport/animal/nature, *etc*., which helps machines to understand more diverse audio-visual scenes in the wild.

## B. Implementation Details

**Feature Extraction.** The visual features are extracted using two-stream I3D [6], which inputs a set of 24 RGB and optical flow frames extracted at 25 fps. Each frame is first resized such that the shortest side is 256 pixels, and then the center region is cropped to $224 \times 224$. A 1024-d RGB or flow feature vector is obtained from the final convolutional layer of the corresponding branch of I3D. Then, the two vectors are concatenated producing 2048-d features for each stack of 24 frames. The audio features are extracted using VGGish [13]. The input is a $96 \times 64$ log mel-scaled spectrogram extracted for each 0.96s segment, which is obtained by applying *Short-Time Fourier Transform* on a 16 kHz mono audio track. Then, a 128-d feature vector can be obtained after an activation function and before a classification layer. Here, we use 24 frames for each visual segment to temporally match with the input of the audio modality as $\frac{24}{25} = 0.96$.

**Network Architecture.** In the cross-modal pyramid transformer encoder, the number of attention heads is 4 in both uni-modal and cross-modal blocks. The temporal downsampling operation is realized by using a single depth-wise

| PE | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | Avg. |
|---|---|---|---|---|---|---|
| ✓ | **50.6** | 44.8 | **39.8** | 32.4 | 21.1 | **47.8** |
|  | 49.5 | **45.1** | 39.7 | **32.8** | **21.9** | 47.0 |

Table 7. Ablation study on position encoding (PE).

| $\lambda$ | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | Avg. |
|---|---|---|---|---|---|---|
| 0.2 | 49.9 | 45.0 | 39.6 | 32.2 | 20.7 | 46.9 |
| 0.5 | 50.1 | 45.4 | 39.8 | 32.3 | 21.2 | 47.3 |
| 1 | **50.6** | **45.8** | 39.8 | 32.4 | 21.1 | **47.8** |
| 2 | 49.8 | 45.3 | **40.2** | **33.0** | **22.4** | 47.2 |
| 5 | 49.0 | 44.7 | 39.2 | 32.3 | 22.2 | 46.4 |

Table 8. Ablation study on loss weight $\lambda$.

| Stride | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | Avg. |
|---|---|---|---|---|---|---|
| 8 | **50.6** | **44.8** | **39.8** | 32.4 | 21.1 | **47.8** |
| 16 | 48.9 | 44.6 | 39.0 | **32.9** | **21.8** | 46.7 |
| 24 | 49.7 | 44.7 | 38.5 | 31.0 | 20.9 | 47.0 |

Table 9. Ablation study on temporal feature stride.

| $T_{max}$ | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | Avg. |
|---|---|---|---|---|---|---|
| 192 | 49.9 | 45.2 | 39.7 | 32.6 | 21.7 | 47.0 |
| 224 | **50.6** | **45.8** | 39.8 | 32.4 | 21.1 | **47.8** |
| 256 | 49.6 | 45.3 | **39.9** | **33.1** | **22.3** | 47.2 |

Table 10. Ablation study on maximum input sequence length.

| SD | CD | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | Avg. |
|---|---|---|---|---|---|---|---|
|  |  | 48.5 | 43.4 | 36.9 | 29.9 | 20.2 | 45.8 |
| ✓ |  | 49.5 | **45.3** | 39.7 | 32.6 | 21.2 | 46.9 |
|  | ✓ | 49.5 | 44.8 | 39.7 | **32.7** | **21.7** | 46.8 |
| ✓ | ✓ | **50.6** | 44.8 | **39.8** | 32.4 | 21.1 | **47.8** |

Table 11. Ablation study on dependency modeling. SD: simultaneous dependency branch; CD: consecutive dependency branch.

1D convolution as in [48]. For temporal dependency modeling, the output dimension is converted as the shape of input to formulate it as a plug-and-play operation, and we just apply this operation once in our model.

**Reproducibility.** All our models are trained on a single 32GB NVIDIA Tesla V100 GPU and implemented in PyTorch deep-learning framework. During inference, we evaluate the performances of our method on the test set of our UnAV-100 and use the best models on the validation set.

## C. Ablation Study

**Position Encoding.** We explore the impact of position encoding in our transformer encoder. As shown in Tab. 7, adding position embeddings can improve the performance

| Method | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | Avg. |
|---|---|---|---|---|---|---|
| ActionFormer [48] | 73.4 | 67.5 | 57.6 | 47.6 | 33.7 | 56.0 |
| Ours | **74.8** | **70.1** | **60.7** | **48.1** | **34.0** | **57.5** |

Table 12. Experiments on THUMOS14 dataset with only visual modality as input (mAP@[0.3:0.1:0.7] is reported).

by $0.8\%$ in average mAP, even though the temporal convolutions (*i.e.*, the projection layer and downsampling operations) already leak the location information as pointed out in [44, 48].

**Loss Weight.** We also provide the ablation study on the loss weight $\lambda$ in our loss function. We train the model using different loss weights $\lambda \in [0.2, 0.5, 1, 2, 5]$, and report the results in Tab. 8. It can be seen that the default value $\lambda = 1$ can yield the best performance.

**Feature Stride.** In our experiments, we use stride=8 with a sliding window of 24 frames by default when extracting visual and audio features. Here, we study the performance variation using different feature strides in Tab. 9. Reducing the temporal feature resolution (*i.e.*, larger strides, 16/24) leads to obvious performance degradation, which is intuitively reasonable since the model might fail to detect many short audio-visual events at a low temporal resolution.

**Maximum Input Sequence Length.** Furthermore, we vary the length of the maximum input sequences of our model, and the results are provided in Tab. 10. We can observe that our model has quite stable results when using different $T_{max}$, and $T_{max} = 224$ gets the best results.

**Dependency Modeling.** Since the two branches of temporal dependency modeling aim to capture different correlations between events within a video, we run an ablation by removing each of the branches and show the results in Tab. 11. It indicates that applying each branch separately also leads to improvement, and the best result can be achieved by combing both branches to model simultaneous and consecutive dependencies at the same time.

## D. Experiments on Existing TAL Dataset

We also conduct experiments on THUMOS14 dataset [14], a widely-used dataset for temporal action localization. The evaluation results on THUMOS14 test set using only visual input are provided in Tab. 12. We use the same strategy to extract features on THUMOS14 as used on UnAV-100 for both methods to keep a fair comparison. We can see that our model outperforms ActionFormer [48] by a large margin (+3.1% mAP at tIoU=0.5), even without the cross-modal fusion strategy. Besides, we tried to only use the audio modality in THUMOS14 to locate actions, but got very bad results (just 4.3% average mAP) on both models, which indicates that the audio tracks in THUMOS14 are quite noisy and cannot provide useful information.
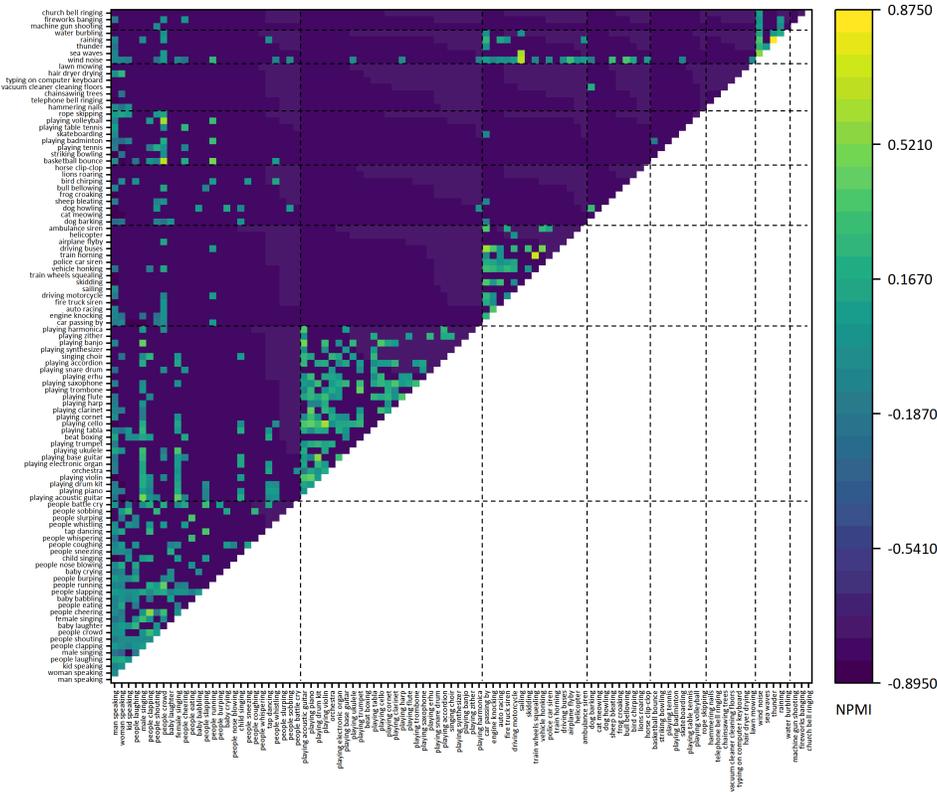
## E. More Qualitative Results

More qualitative results are presented in Fig 9, which includes the prediction results of our model variants using different modalities as input. Generally speaking, cross-modal perception encourages the model to obtain more correct localization results. For example, Fig. 9(a) refers to the relatively constant visual information versus dramatically changing audio signals. By integrating both modalities, the model can better judge the event boundaries. Besides, our audio-visual model can also get promising performance in some complex audio-visual scenarios, as in Fig. 9(c) and Fig. 9(d), where many audio-visual events occur concurrently or over very short periods of time.

## F. Discussion

**Limitations.** There is still a wide scope for exploration and improvement on the basis of our work. For instance, our dataset is limited to a temporal localization task. We will explore other audio-visual problems, such as representation learning and sound source localization in real-life and complex scenarios in our subsequent study. Besides, although our model can obtain a promising performance, as a baseline, its capability is still limited in some complex situations. For example, in Fig. 9(c), the model gets an incorrect boundary of the *dog barking* event when the barking brown dog is out of the screen and a non-barking black one can be seen. This indicates that our model might fail to effectively filter out interference information for such a difficult case. And the model might also fail to predict precise boundaries when one modality persists while another disappears for a short period of time (*e.g.*, the event of *vacuum cleaner cleaning floors* in Fig. 9(c)). In addition, for some instant events with very short duration (*e.g.*, *basketball bounce* in Fig. 9(d)), our model might get unsatisfactory results. Overall, dense-localizing audio-visual events is inherently a very challenging task, and it requires the model to have a strong fine-grained cross-modal understanding ability. Therefore, more advanced models that could better solve the above difficulties are expected to boost performance further. We hope our work as the first attempt at untrimmed audio-visual video understanding can inspire more people to explore the field.

**Ethic concerns and biases.** Our UnAV-100 is sourced from VGGSound dataset [8] that has already tried to mitigate ethical issues. During data collection, we made further efforts to manually check all videos to avoid mature, sensitive, or offensive content. Besides, our UnAV-100 follows the natural distribution of instances present on the website, which may reflect some biases in topics. For example, there are more *man/woman speaking* events than other categories. Efforts have been made to mitigate such imbalance.
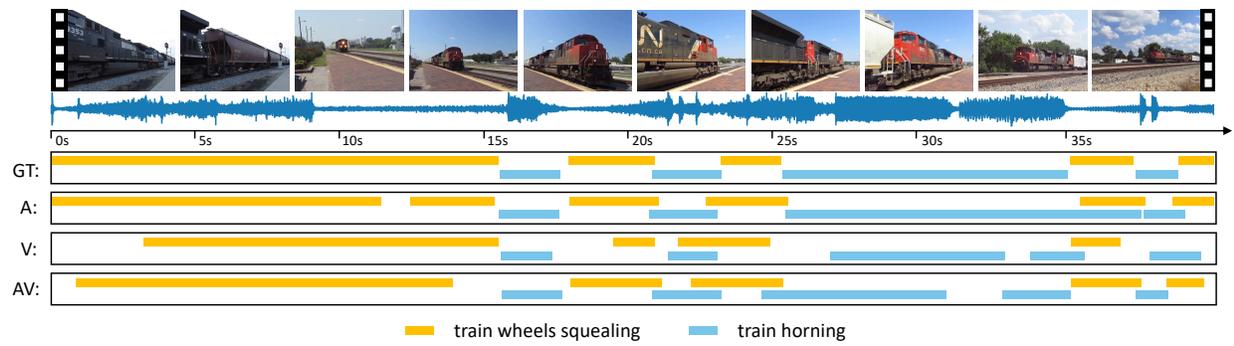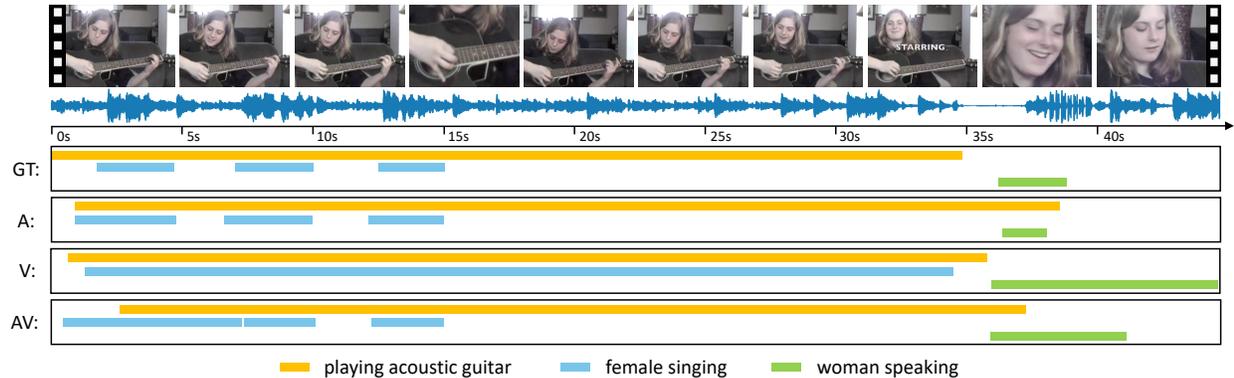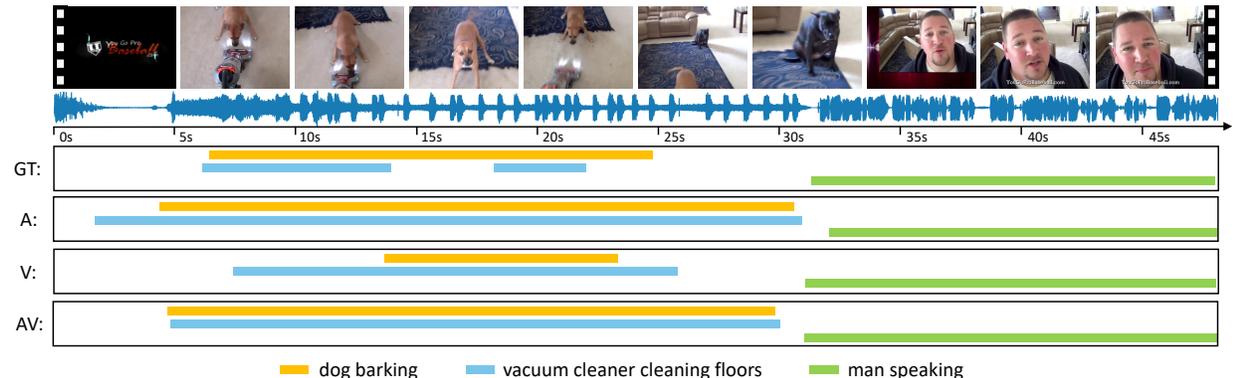
Figure 8. NPMI of the pairs of simultaneous (a) and consecutive (b) audio-visual events in our UnAV-100 dataset. In (b), the horizontal axis shows the first event, and the vertical axis shows the second subsequent event. The event categories are grouped by domains.
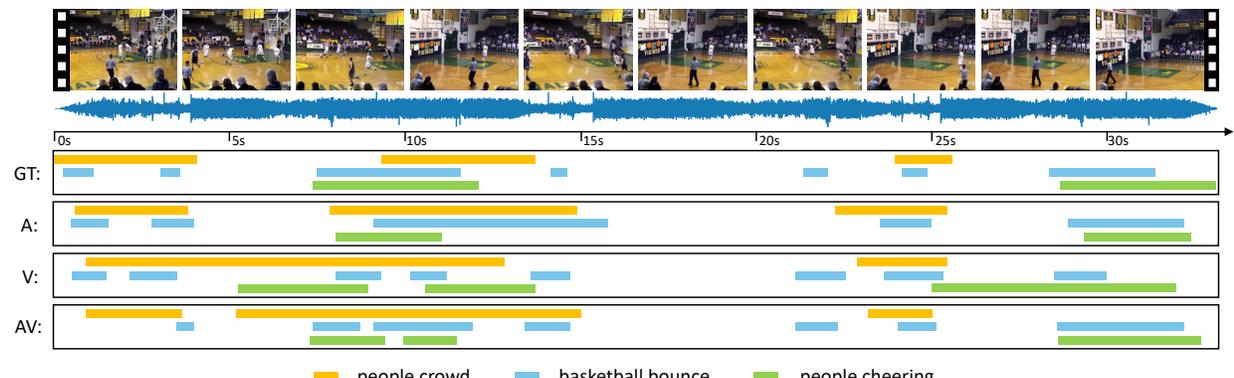
Figure 9. More qualitative results on the UnAV-100 test set. GT: ground truth, A: the prediction of the audio-only variant, V: the prediction of the visual-only variant, AV: the prediction of our audio-visual model. We show boundaries with the highest overlap with ground truth.