# ENHANCED SOUND EVENT LOCALIZATION AND DETECTION IN REAL 360° AUDIO-VISUAL SOUNDSCAPES

*Adrian S. Roman*\*      *Baladithya Balamurugan*      *Rithik Pothuganti*

Viterbi School of Engineering, University of Southern California, California, USA

## 1. ABSTRACT

This technical report details our work towards building an enhanced audio-visual sound event localization and detection (SELD) network. We build on top of the audio-only SELDnet23 model and adapt it to be audio-visual by merging both audio and video information prior to the gated recurrent unit (GRU) of the audio-only network. Our model leverages YOLO and DETIC object detectors. We also build a framework that implements audio-visual data augmentation and audio-visual synthetic data generation. We deliver an audio-visual SELDnet system that outperforms the existing audio-visual SELD baseline.

## 2. INTRODUCTION

A sound event localization and detection (SELD) system generates temporal detection of active sound classes with their corresponding direction of arrival (DoA) around a microphone array [1]. The spatiotemporal characterization of sound scenes generated by a SELD system has a wide range of applications, such as in audio-visual navigation systems and bio-acoustic monitoring systems [2, 3, 4].

In real-world environments, a soundscape is a mixture of one or more sound events. A sound event emerges when object interactions generate sound vibrations, e.g., a musician plays a guitar, a person opens a door. Such events virtually never occur as unimodal auditory events but almost always as audio-visual events. SELD depends on the physical characteristics of acoustic scenes to track moving sources, perform sound source distance estimation [5, 6] and account active sources even if they are occluded by another object. With such challenges in mind, it is clear that the visual modality provides useful information to mitigate ambiguities in a SELD task.

Most of the novel research in audio-visual sound source localization aims to learn audio-visual correspondence [7, 8]. While these models are robust at localizing sound on an image, they are not designed for estimating physical direction of arrival (DoA) in a soundscape. Such models are often trained with audio-visual datasets that do not contain DoA labels and only contain monoaural data[9]. For this reason, the sound localization performance strongly depends on the video content [10]. This makes models prone to erroneous SELD on frames with no audio or uncorrelated audio activity.

We introduce a visual branch into the audio-only SELDnet23 baseline from the Classification of Acoustic Scenes and Events (DCASE) Challenge 2023 task3. We equip the SELDnet system with state-of-the-art (SOTA) object detectors such as YOLO[1] and DETIC [11]. Additionally, we enable data generation for model training by implementing novel video and audio processing techniques. In summary, our contributions are:

1. An implementation of audio-visual data augmentation techniques originally proposed by Wang et al [12, 13].
2. An audio-visual synthetic data generator for spatial audio and 360° video.
3. A set of audio-visual SELD systems that outperform the existing baselines.

We provide this work as an open source framework available at https://github.com/aromanusc/SoundQ

## 3. METHODS

### 3.1. Datasets

We use the Sony-Tau Realistic Spatial Soundscapes 2023 (STARSS23) dataset[2] [14] for development and evaluation. The dataset contains two 4-channel 3-dimensional recording formats: first-order Ambisonics (FOA) and tetrahedral microphone array (tetra), each recorded at a sampling rate of 24kHz and a bit depth of 16 bits. The corresponding video is 360-degree equirectangular with 1920x960 resolution sampled at 29.97 frames-per-second. The data was strongly labeled with spatiotemporal DoA and class labels for a set of 13 target sound classes. Due to the complexities of building a real audio-visual dataset, the development set of STARSS23 is limited (3.8 hours) compared to the large scale data used in the previous iterations of the DCASE Challenge Task3. We tackle the data scarcity by applying three techniques:

---

[1]https://github.com/ultralytics/ultralytics
[2]https://zenodo.org/record/7880637

**Fig. 1**. Video pixel swapping augmentations.

**Data augmentation:** We adopt and implement the audio channel swapping (ACS) and video pixel swapping (VPS) method proposed by Wang et al. [12, 13]. The ACS method increases the amount of DoA representations by performing audio channel rotations within the recordings. In the video domain, the STARSS23 dataset includes 360° video clips with 1920x960 resolution, corresponding to an azimuth range of $[-90°, 90°]$ and an elevation range of $[-180°, 180°]$. The VPS method matches ACS by applying the same transformations at pixel level. For instance, an ACS transformation like $\phi = \phi - \pi/2, \theta = -\theta$, implies rotating the azimuth angle by $-90°$ and inverting the elevation axis. At the pixel level, the VPS transformation slides the azimuth pixels by 1440 pixel points in the negative direction and inverts the elevation axis.

Using the ACS and VPS methods we augment the original audio data by a factor of seven. Figure 1 features the total transformations (including identity) that we integrate in our training set.

**Audio-only synthetic data:** We leverage the synthetic data from the DCASE Challenge 2022 Task3[3]. The recordings were generated through convolution of isolated sound samples with real spatial room impulse responses (SRIRs) captured in nine unique spaces of Tampere University. The training data contains the same sound events classes as the STARSS23 dataset, where the sound assets are sources from the FSD50K dataset [15].

**Audio-visual synthetic data:** We build a synthetic 360° audio-visual synthetic data generator. We collected a total of 200 YouTube videos containing similar sound events to the


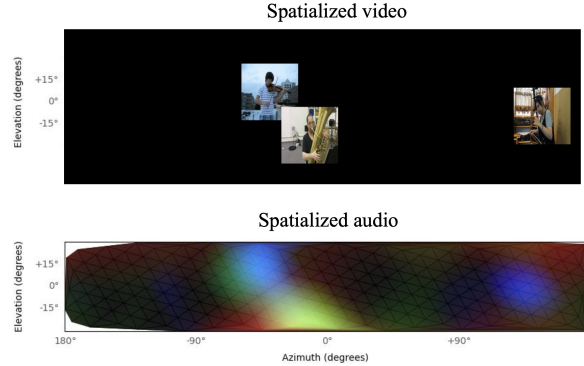
**Fig. 2**. 360° audio-visual synthetic frame (top) and its spatialized audio displayed using an acoustic camera (bottom).

STARSS23 dataset. In the audio domain, we generate spatialized sound events using room impulse responses (RIRs) from the METU-SPARG RIR dataset [16]. The spatial audio synthesizer extracts the audio from a YouTube video and convolves it with an RIR. For the simulated 360° video we generate a black video canvas of 1920x960 resolution. We resize the YouTube videos to tiles of 50x50 pixels. We use the 2D projection of the RIR coordinates to superimpose the video frames on the black canvas. The result is a synchronized audio-visual clip. Our audio-visual synthesizer generates variable length audio-visual soundscapes with a maximum of three active events per frame. The metadata format follows the same convetion as the STARSS23 dataset.

We generate a total of 100 synthetic videos each of 30 seconds duration. The top image of Figure 2 shows an example of the generated 360° video frames from our synthesizer. We validate the spatialized audio implementation using Deep-Wave's PyTorch[4] to generate acoustic maps [17, 18]. The color intensity clusters from Figure 2 correspond to the active sound sources, which match the spatial locations of the video overlays. We confirmed the sounds direction of arrival (DOA) estimated a deep acoustic imaging beamformer[5].

### 3.2. The SELD Baseline models

**Audio-only baseline:** We use the audio-only SELDnet23 baseline model from the DCASE Challenge Task3 [19, 1, 20]. The model uses multichannel audio. SELDnet23 is equipped with multi-ACCDOA and it can simultaneously infer presence, class, and spatial coordinates for up to three sound events [21]. Compared to previous SELDnet baselines, it also introduces two multi-head self-attention (MHSA) layers, which makes it more robust for SELD tasks.

**Audio-visual baseline:** We use the audio-visual SELD-net23 baseline model from the DCASE Challenge Task3
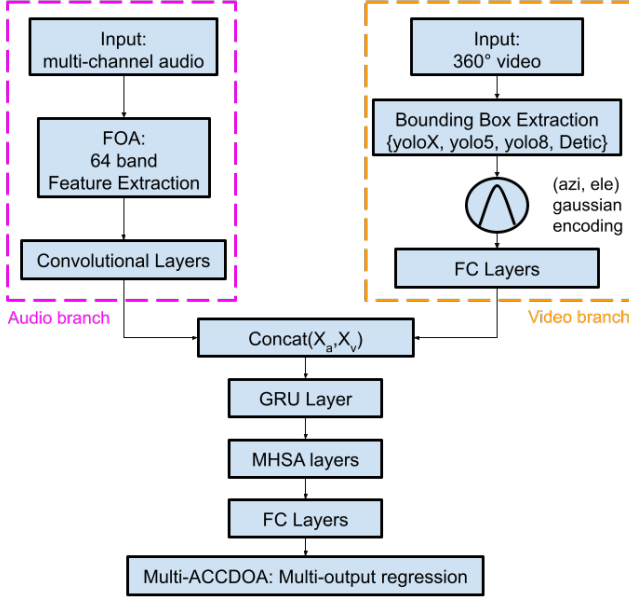
---

**Fig. 3**. Enhanced audio-visual SELD system.



**Fig. 4**. Per-class localization error performance comparing the audio-visual baseline against our proposed enhancements.

[14]. This model shared a similar architecture as the audio-only SELDnet, except that it does not contain multi-head self-attention (MHSA) layers. The audio-visual model is equipped with a visual branch that makes use of an object detector (YOLOX [22]). The detector is only used to capture "human" objects. Each bounding box is encoded into two Gaussian-like vectors, corresponding to the azimuth $\rho_{azi}(u_{ij}) \in \mathbb{R}^N$ and elevation $\rho_{ele}(v_{ij}) \in \mathbb{R}^N$. $N$ represents the encoding size, which is $N = 37$ for the baseline. Note that the baseline model only supports a maximum of $M = 6$ bounding boxes per frame (i.e. $i = 1..6$). The vectors are then concatenated in to a visual embeddings vector $\mathbf{X_v} \in \mathbb{R}^{2 \times M \times N}$. The visual embeddings $\mathbf{X_v}$ are then concatenated with audio embeddings $\mathbf{X_a}$ and sent through an audio-visual decoder network that allows to output multi-ACCDOA labels.

### 3.3. Audio-visual SELD model enhancements

The audiovisual SELDnet baseline shows limited performance compared to the audio-only baseline [14]. Such shortcomings are due to the scarce training data and the architectural limitations of its visual branch [14].

We focus our study towards architectural enhancements to the vision branch and deliver two model variants that replace the original object detector: YOLO-based and DETIC-based. This two-fold approach allows us to understand how enhancing object detection contributes to better SELD performance. Similarly, we are able to understand the contributions of fixed vocabulary detectors (i.e YOLO), compared to detectors with large and customizable vocabulary (i.e DETIC). The YOLO-based detectors attain detection fixed to the COCO dataset
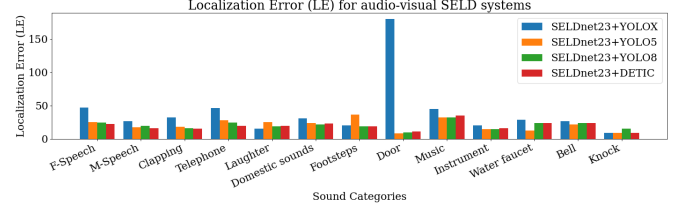
classes [23]. DETIC has the ability to identify 21,000 object classes with strong accuracy as it was trained with a variety of datasets, among those the LVIS dataset [24]. Furthermore, since DETIC employs a CLIP embedding vector, it is possible to change the model's vocabulary to a customized one. In our case, we customized DETIC's vocabulary to target the STARSS23 classes. We also equip our audio-visual SELD system with two MHSA layers inspired from the audio-only baseline [25]. Our architectural enhancements therefore integrate either YOLO5, YOLO8 or DETIC as the object detectors. The visual embeddings follow the same format as the audio-visual baseline described in section 3.2. Figure 3 displays our proposed audio-visual SELDnet23 architecture.

### 3.4. Metrics

We employ the SELD metrics proposed by the DCASE Challenge [26]. Two metrics relate to DoA estimation: F1-score ($F_{20°}$) and error rate ($ER_{20°}$). $F_{20°}$ is calculated from location-aware precision and recall. $ER_{20°}$ is the sum of insertion, deletion and substitution errors divided by the total number of inferred audio frames. The remaining two metrics relate to class-aware localization: localization error (LE) in degrees and localization recall (LR). LE is the average angular difference between each class prediction and its ground truth labels. LR is the true positive rate of instantaneous detections out of the total annotated sounds.

### 3.5. Training procedure

Our proposed SELD systems adopt the same hyper-parameters as the audio-only SELDnet23 baseline. We only augment the STARSS23 train split and our synthetic audio-visual data using Wang et al's data augmentation approach [12, 13]. We also include the audio-only synthetic data from DCASE22 discussed in section 3.1. Audio-visual models were trained with STARSS23, synthetic audio-visual data, and audio-only synthetic data. For the audio-only synthetic data a black video stream is presented as background. The audio-only model was trained with the augmented STARSS23 and audio-only synthetic data. All models were trained over 100 epochs and the best of 5 validation performances was chosen.

| Model | detector | data aug | input | $ER_{20°}$ | $F_{20°}$ | $LE° \downarrow$ | $LR \uparrow$ |
|---|---|---|---|---|---|---|---|
| Baseline AO SELDnet23 | N/A | Yes | FOA + Multi-ACCDOA | 0.57 | 29.9 | 21.6 | **47.7** |
| Baseline AV SELDnet23 | YOLOX | No | FOA + Video | 1.07 | 14.3 | 48.0 | 35.5 |
| Baseline AV SELDnet23 | YOLOX | Yes | FOA + Video | 1.37 | 15.0 | 40.62 | 40.0 |
| **Ours AV SELDnet** | YOLO5 | Yes | FOA + Multi-ACCDOA + Video | 0.64 | 27.5 | 21.0 | 41.4 |
| **Ours AV SELDnet** | YOLO8 | Yes | FOA + Multi-ACCDOA + Video | 0.63 | 30.9 | 20.3 | 46.1 |
| **Ours AV SELDnet** | DETIC | Yes | FOA + Multi-ACCDOA + Video | 0.64 | 30.6 | **19.5** | 43.7 |

**Table 1**. Comparison of the baseline SELDnet23 audio-only and audio-visual against our proposed audio-visual SELDnet architecture. Best performance in the "test-split" from the development STARSS23 dataset. AV=audio-visual, AO=audio-only.

## 4. RESULTS

Table 1 shows performance by our proposed audio-visual SELDnet enhancements and compares them against the audio-only and audio-visual SELDnet23 baselines. The results below are the evaluation metric scores using the 'test split' of the STARSS23 development set. Most models were trained using audio-visual data augmentations (see section 3.5), otherwise denoted by 'No' on the third column.

The top three rows correspond to the audio-only (AO) and the audio-visual (AV) SELDnet23 baselines. The AO SELDnet23 features more robust LE and LR metrics compared to the AV SELDnet23. This is because the AV SELDnet23 was only trained on the STARSS23 development set with no data augmentations, while the AO SELDnet23 was trained with a combination of the STARSS23 development set recordings and synthetic audio recordings from DCASE22. Additionally, the AO system is equipped with two MHSA layers that make the system more robust at localization tasks. The third row shows the benefits of data augmentations on the AV SELDnet23, giving a clear improvement on LE ~8° and LR ~5° over the original AV SELDnet23.

The models equipped with YOLO5 and YOLO8 detectors outperform in all metrics the original AV baseline that uses the tiny YOLOX version. This shows the advantage of using more robust object detectors that are resilient to the equiangular projection distortion on 360° video. The model equipped with YOLO8 outperforms other YOLO-based architectures. This is a product of YOLO8's transformer-based architecture, which has a much higher accuracy compared to earlier YOLO models. Note that the YOLO-based systems inherit object detection from the COCO dataset classes [23]; hence, because the COCO classes do not overlap with the STARSS23 classes, the systems virtually detects only the "person" class.

The DETIC-based system (last row) is the best-performing model in regards to LE. This model shows the advantage of having a detector with a vocabulary tailored to the STARSS23 classes. The slight degradation in LR performance is an indicator of the added complexity for handling multiple object class detection in the video stream.

Figure 4 breaks down the LE performance for each sound category, comparing the baseline system (SELDnet23+YOLOX) against our proposed architectural enhancements. In this plot, all models were trained with the augmented training set. It is clear that our proposed model enhancements outperform the audio-visual baseline nearly for all sound categories. In addition, the DETIC-based system generally performs better on classes where there may be clear object-human interactions.

## 5. CONCLUSION AND FUTURE WORK

We have proposed a set of audio-visual models that outperform the existing audio-only and audio-visual SELDnet23 baselines. Our experiments reveal the benefits of using audio-visual data augmentations and audio-visual synthetic data. We also shed light on the advantages of using SOTA object detectors for fixed vocabulary (e.g., YOLO) and customized vocabulary (e.g., DETIC). While our validation results show enhanced performance over baselines, there may still be overfitting or domain-shift risks.

Future work could build on our proposed audio-visual synthetic data generator. There could be large benefits from collecting 360° video samples to generate synthetic high-quality audio-visual soundscapes at scale. Improvements to our method could include superimposing video frames using 3D coordinates (including depth), which could used for training SELD systems for depth estimation. Our proposed audio-visual synthetic generation is planned to be integrated in a future release of the Spatial Scaper library [27].

## 6. ACKNOWLEDGEMENTS

# 7. REFERENCES

[1] Sharath Adavanne, Archontis Politis, Joonas Nikunen, and Tuomas Virtanen, "Sound event localization and detection of overlapping sources using convolutional recurrent neural networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 1, pp. 34–48, 2018.

[2] Pierre-Amaury Grumiaux, Srđan Kitić, Laurent Girin, and Alexandre Guérin, "A survey of sound source localization with deep learning methods," *JASA*, vol. 152, no. 1, pp. 107–151, 2022.

[3] Changan Chen, Unnat Jain, Carl Schissler, Sebastia Vicenc Amengual Gari, Ziad Al-Halah, Vamsi Krishna Ithapu, Philip Robinson, and Kristen Grauman, "Soundspaces: Audio-visual navigation in 3d environments," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VI 16*. Springer, 2020, pp. 17–36.

[4] Aurora Linh Cramer, Vincent Lostanlen, Andrew Farnsworth, Justin Salamon, and Juan Pablo Bello, "Chirping up the right tree: Incorporating biological taxonomies into deep bioacoustic classifiers," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 901–905.

[5] Saksham Singh Kushwaha, Iran R Roman, Magdalena Fuentes, and Juan Pablo Bello, "Sound source distance estimation in diverse and dynamic acoustic conditions," in *2023 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2023, pp. 1–5.

[6] Benjamin S Liang, Andrew S Liang, Iran Roman, Tomer Weiss, Budmonde Duinkharjav, Juan Pablo Bello, and Qi Sun, "Reconstructing room scales with a single sound for augmented reality displays," *Journal of Information Display*, vol. 24, no. 1, pp. 1–12, 2023.

[7] Shentong Mo and Pedro Morgado, "Localizing visual sounds the easy way," in *European Conference on Computer Vision*. Springer, 2022, pp. 218–234.

[8] Shentong Mo and Pedro Morgado, "A closer look at weakly-supervised audio-visual source localization," *Advances in Neural Information Processing Systems*, vol. 35, pp. 37524–37536, 2022.

[9] 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)*. IEEE, 2020, pp. 721–725.

[10] Julia Wilkins, Justin Salamon, Magdalena Fuentes, Juan Pablo Bello, and Oriol Nieto, "Bridging high-quality audio and video via language for sound effects retrieval from visual queries," in *2023 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2023, pp. 1–5.

[11] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra, "Detecting twenty-thousand classes using image-level supervision," in *European Conference on Computer Vision*. Springer, 2022, pp. 350–368.

[12] Qing Wang, Jun Du, Hua-Xin Wu, Jia Pan, Feng Ma, and Chin-Hui Lee, "A four-stage data augmentation approach to resnet-conformer based acoustic modeling for sound event localization and detection," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 31, pp. 1251–1264, 2023.

[13] Qing Wang, Li Chai, Huaxin Wu, Zhaoxu Nian, Shutong Niu, Siyuan Zheng, Yuyang Wang, Lei Sun, Yi Fang, Jia Pan, et al., "The nerc-slip system for sound event localization and detection of dcase2022 challenge," *DCASE2022 Challenge, Tech. Rep.*, 2022.

[14] Kazuki Shimada, Archontis Politis, et al., "Starss23: An audio-visual dataset of spatial recordings of real scenes with spatiotemporal annotations of sound events," *arXiv preprint 2306.09126*, 2023.

[15] Eduardo Fonseca, Xavier Favory, Jordi Pons, Frederic Font, and Xavier Serra, "Fsd50k: an open dataset of human-labeled sound events," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 829–852, 2021.

[16] Orhun Olgun and Huseyin Hacihabiboglu, "METU SPARG Eigenmike em32 Acoustic Impulse Response Dataset v0.1.0," Apr. 2019.

[17] Matthieu Simeoni, Sepand Kashani, Paul Hurley, and Martin Vetterli, "Deepwave: a recurrent neural-network for real-time acoustic imaging," *Advances In Neural Information Processing Systems*, vol. 32, 2019.

[18] Adrian S Roman, Iran R Roman, and Juan P Bello, "Robust doa estimation using deep acoustic imaging," *arXiv preprint arXiv:2401.08717*, 2024.

[19] Sharath Adavanne, Archontis Politis, and Tuomas Virtanen, "Localization, detection and tracking of multiple moving sound sources with a convolutional recurrent network," *arXiv preprint 1904.12769*, 2019.

[20] Sharath Adavanne, Archontis Politis, and Tuomas Virtanen, "A multi-room reverberant dataset for sound event localization and detection," in *DCASE*, 2019, pp. 10–14.

[21] Kazuki Shimada, Yuichiro Koyama, Shusuke Taka-hashi, Naoya Takahashi, Emiru Tsunoo, and Yuki Mit-sufuji, "Multi-accdoa: Localizing and detecting over-lapping sounds from the same class with auxiliary dupli-cating permutation invariant training," in *IEEE ICASSP*. IEEE, 2022, pp. 316–320.

[22] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun, "Yolox: Exceeding yolo series in 2021," *arXiv preprint arXiv:2107.08430*, 2021.

[23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, "Microsoft coco: Common ob-jects in context," in *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.

[24] Agrim Gupta, Piotr Dollar, and Ross Girshick, "Lvis: A dataset for large vocabulary instance segmentation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 5356–5364.

[25] Parthasaarathy Sudarsanam, Archontis Politis, and Kon-stantinos Drossos, "Assessment of self-attention on learned features for sound event localization and detec-tion," *arXiv preprint arXiv:2107.09388*, 2021.

[26] Archontis Politis, Annamaria Mesaros, Sharath Ada-vanne, Toni Heittola, and Tuomas Virtanen, "Overview and evaluation of sound event localization and detec-tion in dcase 2019," *IEEE/ACM Transactions on Au-dio, Speech, and Language Processing*, vol. 29, pp. 684–698, 2020.

[27] Iran R Roman, Christopher Ick, Sivan Ding, Adrian S Roman, Brian McFee, and Juan P Bello, "Spatial sca-per: a library to simulate and augment soundscapes for sound event localization and detection in realistic rooms," *arXiv preprint arXiv:2401.12238*, 2024.