

SEMANTIC ACOUSTIC IMAGING VIA SELD-STYLE DOA REGRESSION AND A PER-CLASS AUDIOVISUAL ENSEMBLE

Technical Report

Gwantae Kim

Samsung Electronics
kgt1103211@gmail.com

ABSTRACT

Semantic acoustic imaging aims to estimate class-conditioned spatial energy distributions from audio or audiovisual observations. This report describes our submission to DCASE 2026 Task 3, where systems predict dense acoustic-energy maps on the sphere from either 4-channel spatial audio (Track A) or paired audio and 360° video (Track B). Directly regressing dense per-class maps from recordings is difficult because the targets are spatially sparse, the losses are dominated by background regions, and real labeled recordings are limited. We therefore decompose the problem into sound event localization and detection (SELD) followed by a deterministic renderer that converts direction estimates into dense spherical-Gaussian energy fields. The audio model is trained on the provided real recordings and synthetic recordings, with 8-pattern audio-channel swapping (ACS) and SpecMix mixup used to improve spatial robustness. At inference, we explore complementary strategies across the submitted systems: 8-view inverse-aligned ACS test-time ensembling, a multi-checkpoint cluster ensemble, and per-class audiovisual model selection. On the held-out development-test split, the audio system achieves 0.0234 mask mAP / 0.358 Pearson r , and the per-class audiovisual ensemble reaches 0.0332 / 0.446 (+42% mAP, +25% Pearson).

Index Terms— SELD, acoustic imaging, ACCDOA, GCC-PHAT, audio channel swapping, audiovisual fusion, model ensemble

1. INTRODUCTION

DCASE 2026 Task 3 (SAISELD) [1] extends conventional SELD from sparse point estimates to dense semantic acoustic images. For each frame and sound-event class, the system must output an acoustic-energy map over the $360^\circ \times 180^\circ$ sphere, simultaneously encoding source direction and relative energy. Performance is evaluated using mask mAP under the COCO protocol (IoU 0.50:0.95, with masks binarized at 10% of their peak) and Pearson correlation r between matched predicted and reference energy fields. The final leaderboard score is the sum of the mAP and Pearson ranks. Track A restricts the input to the 4-channel tetrahedral audio available at evaluation time, whereas Track B additionally permits the use of 360° video.

Our approach decouples *localization* from *image rendering*. Rather than optimizing a dense pixel-wise objective whose gradients are dominated by background regions, we first train a spatial-feature Conformer for SELD-style direction regression. The predicted directions are then rendered into acoustic-energy maps using spherical Gaussian kernels. This decomposition preserves the

strong inductive bias of SELD while producing the dense output required by the task.

The submission contains two main components. First, we train the audio model with calibrated 8-pattern ACS augmentation and use the same transformations at test time to form an inverse-rotated ACS ensemble. Second, for Track B we exploit the class-independent structure of the official metric and construct a per-class ensemble between the audio renderer and an audiovisual Mask R-CNN. The resulting system combines the sharper masks of the SELD pathway with the stronger energy correlation of the audiovisual pathway.

Figure 1 summarizes the resulting unified architecture. The same audio pathway is used for Track A and as one candidate source for Track B, while the audiovisual pathway supplies complementary class-wise masks for classes where visual context improves the official metrics.

2. DATASETS

The **development** set is STAIRS26, which contains 32-channel Eigenmike recordings, synchronized 360° video, and dense JSON acoustic-image labels [1]. Because the evaluation data provide only a 4-channel tetrahedral array, all audio models are trained and validated using the corresponding 4-channel subset. We follow the official train/test partition and further track the capture-site split between Sony and TAU recordings. Unless otherwise stated, reported results use the held-out development-test recordings (Sony 30 + TAU 48 = 78 clips). The TAU-test subset is not used for model selection or hyper-parameter tuning, and is therefore treated as a generalization check.

The **evaluation** set is the blind STARSS23 set [2, 3], which provides 4-channel audio and video but no labels or 32-channel audio. Following the task rules, the final evaluation submissions are retrained on the full development set (train + test). The split-respecting models reported in Section 5 are used only for development-set analysis.

3. TRACK A: SELD-STYLE AUDIO SYSTEM

3.1. Input features

For each 4-channel microphone signal, we compute per-channel log-mel spectrograms (64 mel bands, $n_{\text{fft}}=512$, hop size 240) and all 6 pairwise GCC-PHAT cross-correlations. These features are stacked into a 13-channel tensor. The frequency/lag axis is kept uncollapsed (`flatten`) so that the Conformer [4] can model the

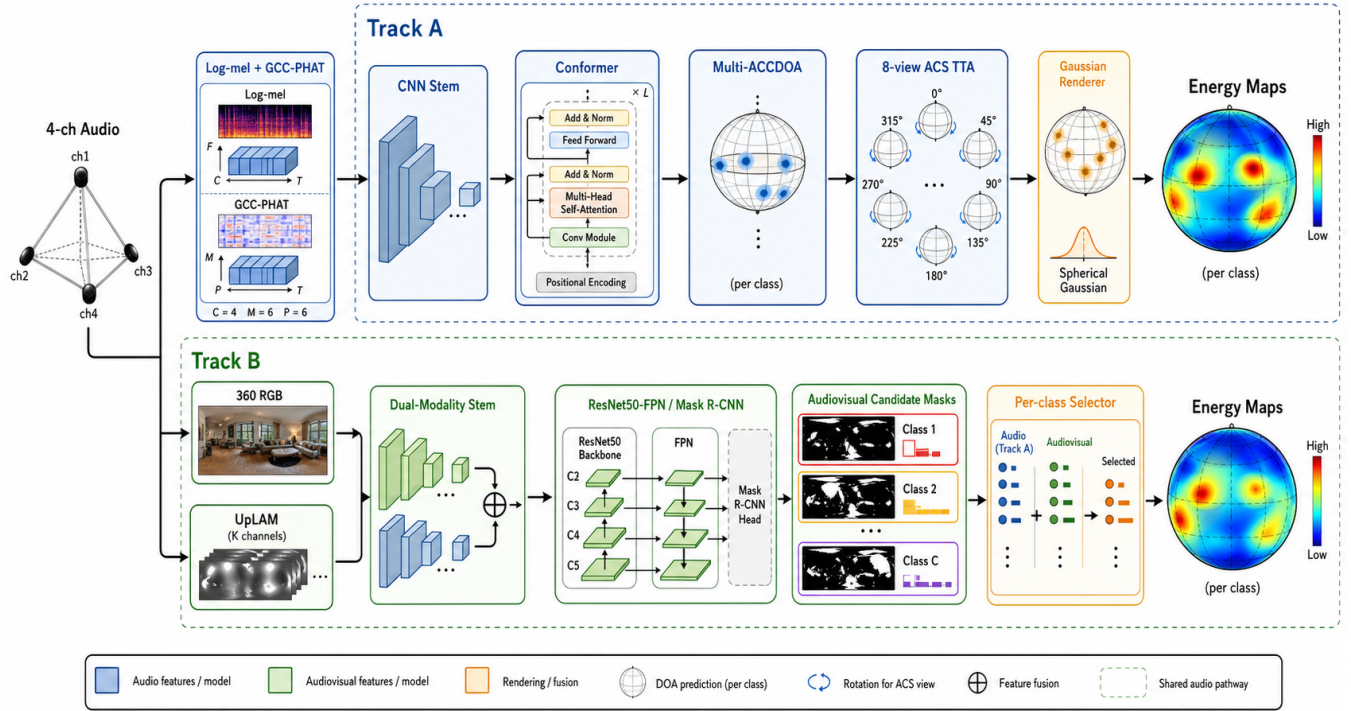


Figure 1: Unified neural-network pipeline for Track A and Track B. The audio pathway estimates class-wise directions with a Conformer-based SELD model and renders them into spherical-Gaussian energy maps. Track B additionally uses an audiovisual Mask R-CNN and selects, for each class, the stronger candidate according to development-set validation behavior.

TDOA structure exposed by GCC-PHAT. In preliminary experiments, averaging over this axis removed important spatial cues and degraded localization.

3.2. Architecture

A convolutional stem with channel widths $[64, 128, 256]$ feeds a Conformer encoder (4 blocks, $d_{\text{model}}=256$, 8 heads). The output is an ACCDOA [5] / multi-ACCDOA [6] head trained with the Auxiliary-Duplicating Permutation-Invariant Training (ADPIT) loss. The head regresses a Cartesian direction vector for each class, includes a distance branch, and supports overlapping sources through duplicated output tracks. The complete model contains approximately 6.6M parameters.

3.3. Training

The model is warm-started from a stochastic-weight-averaged baseline and optimized with AdamW using a cosine learning-rate schedule (3×10^{-4} initial learning rate), batch size 96, and gradient-norm clipping at 1.0. Each epoch samples 80k frames with a weighted sampler.

Two augmentations are used throughout training. **Audio Channel Swapping (ACS)** [7] applies one of 8 tetrahedral channel-swap patterns with probability 0.85, together with the corresponding azimuth and elevation label transformation. **SpecMix mixup** [8] is applied with probability 0.5 and $\alpha=0.4$ by convex-combining two spectra and their ADPIT losses using the same mixing coefficient

λ . ACS increases the effective spatial diversity and also makes the test-time ACS ensemble geometrically consistent.

3.4. Rendering and test-time ensemble

The predicted per-frame, per-class directions are rendered into dense equirectangular energy fields with spherical Gaussian kernels. We use a render threshold of 0.25 and a blob scale of 1.3, which slightly relaxes precision in exchange for higher IoU tolerance and improved mAP. At test time, we apply all 8 ACS transformations to the input, map the predicted direction vectors back with the inverse ACS transforms, render the aligned predictions, and average the resulting fields. Because the model is trained with the same ACS group, the views reinforce consistent directions rather than introducing rotational blur. This 8-view ensemble is the largest single contributor to the mAP gain reported in Table 1.

3.5. Multi-checkpoint ensemble

As an alternative inference strategy, we ensemble the eleven checkpoints retained from the final training run (the stochastic-weight average together with the best, last, and eight late-epoch snapshots). For each frame and class, the per-checkpoint direction estimates from the three ADPIT output tracks are pooled and grouped by great-circle proximity. Clusters supported by fewer than four checkpoints are rejected as outliers, while distinct dense clusters are emitted as separate sources, enabling the ensemble to recover overlapping same-class events. This robust pooling replaces simple averaging and is submitted as an alternative Track A system.

4. TRACK B: PER-CLASS AUDIOVISUAL ENSEMBLE

4.1. Audiovisual baseline

For the audiovisual branch, we retrain the official `EnergyInstanceModel`, a `maskrcnn_resnet50_fpn` [9] with a dual-modality stem. The input consists of a 3-channel RGB frame (360×180 , 10 fps) and 9 acoustic UpLAM channels [1]. The model is trained in its native audiovisual mode with progressive backbone unfreezing (layer4 at epoch 3 and layer3 at epoch 7), RGB dropout of 0.15, and early stopping with patience 15 under a 60-epoch cap. The validation objective plateaued by epoch 6 and did not improve for 13 subsequent epochs, which we use as the convergence point. At inference, per-class peaks are converted into dense energy masks using a score threshold of 0.3. For the multi-checkpoint variants, we further ensemble the audiovisual checkpoints by concatenating their per-frame instance detections before the shared tracking and non-maximum-suppression stage, so that consistent detections reinforce one another within a single decoding pass.

4.2. Per-class selection

The official evaluator scores each class independently: predictions are grouped by (frame, class) and matched to ground truth only within that group. Consequently, when all predictions for a class are taken from one model, that class contributes the corresponding model’s per-class AP and Pearson values to the macro averages. We therefore select the model independently for each class, using per-class AP as the primary criterion and Pearson as the tie-breaker when AP is approximately zero. This selects the audio system for classes {0, 1, 4, 5, 8, 9, 11, 12} and the Mask R-CNN audiovisual model for classes {2, 3, 6, 7, 10}. Under the evaluator’s grouping rule, this procedure forms the per-class max-envelope of the two ranked metrics.

We found that a naive union of detections does not provide the same behavior. The two models use different confidence-score ranges (0.20–0.85 for the audio renderer and 0.06–0.17 for the audiovisual model), so score-ordered greedy matching tends to collapse toward one model’s predictions. Per-class selection avoids this calibration issue and directly encodes the desired mAP–Pearson trade-off.

5. RESULTS

Table 1 reports the official evaluator on the held-out development-test split. The audio system substantially improves over our earlier imaging-only models. The Mask R-CNN audiovisual model is weaker in mAP but consistently stronger in Pearson correlation, indicating that the two architectures capture complementary properties of the target maps. The per-class ensemble improves both metrics over both standalone systems on the combined split and on each capture-site subset, including the TAU subset that is never used for selection.

Relative to the audio-only Track A system, the Track B per-class ensemble improves combined mask mAP by 42% and Pearson correlation by 25%. On the held-out TAU subset, the ensemble reaches 0.0279 mAP, corresponding to an 82% improvement over Track A on the same subset. This indicates that the gain is not solely an artifact of model selection on the development-test split.

Table 1: Official evaluator on the held-out development-test split (mask mAP / Pearson r). “Combined” is all 78 clips; TAU is never used for selection.

System	Combined	Sony	TAU
Track A (audio, TTA)	0.0234/0.358	0.0164/0.342	0.0153/0.292
Mask R-CNN AV (alone)	0.0112/0.437	0.0061/0.400	0.0143/0.452
Per-class ens. (B)	0.0332/0.446	0.0213/0.4220	0.0279/0.372

6. CONCLUSION

We formulate semantic acoustic imaging as SELD-style direction regression followed by deterministic spherical rendering. The Track A system combines a Conformer-based ACCDOA model with ACS training, inverse-aligned ACS test-time ensembling, and spherical-Gaussian map generation. The Track B system further combines this audio renderer with an audiovisual Mask R-CNN through per-class model selection, matching the class-independent structure of the official evaluation. The ACS-ensembled audio model is submitted for Track A, while the per-class audiovisual ensemble is submitted for Track B; both final systems are retrained on the full development set for evaluation.

7. REFERENCES

- [1] A. S. Roman, I. R. Roman, and J. P. Bello, “Latent acoustic mapping for direction of arrival estimation: A self-supervised approach,” in *2025 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2025, pp. 1–5.
- [2] A. Politis, K. Shimada, P. Sudarsanam, S. Adavanne, D. Krause, Y. Koyama, N. Takahashi, S. Takahashi, Y. Mitsufuji, and T. Virtanen, “STARSS22: A dataset of spatial recordings of real scenes with spatiotemporal annotations of sound events,” in *Proceedings of the 8th Detection and Classification of Acoustic Scenes and Events 2022 Workshop (DCASE2022)*, Nancy, France, November 2022, pp. 125–129. [Online]. Available: <https://dcase.community/workshop2022/proceedings>
- [3] K. Shimada, A. Politis, P. Sudarsanam, D. A. Krause, K. Uchida, S. Adavanne, A. Hakala, Y. Koyama, N. Takahashi, S. Takahashi, T. Virtanen, and Y. Mitsufuji, “STARSS23: An audio-visual dataset of spatial recordings of real scenes with spatiotemporal annotations of sound events,” in *Advances in Neural Information Processing Systems*, A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, Eds., vol. 36. Curran Associates, Inc., 2023, pp. 72 931–72 957. [Online]. Available: https://proceedings.neurips.cc/paper/_files/paper/2023/hash/e6c9671ed3b3106b71cafda3ba225c1a-Abstract-Datasets_and_Benchmarks.html
- [4] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, *et al.*, “Conformer: Convolution-augmented transformer for speech recognition,” in *Proc. Interspeech 2020*, 2020, pp. 5036–5040.
- [5] K. Shimada, Y. Koyama, N. Takahashi, S. Takahashi, and Y. Mitsufuji, “Accdoa: Activity-coupled cartesian direction of arrival representation for sound event localization and detection,” in *ICASSP 2021-2021 IEEE international conference on*

- acoustics, speech and signal processing (ICASSP)*. IEEE, 2021, pp. 915–919.
- [6] K. Shimada, Y. Koyama, S. Takahashi, N. Takahashi, E. Tsunoo, and Y. Mitsufuji, “Multi-accdoa: Localizing and detecting overlapping sounds from the same class with auxiliary duplicating permutation invariant training,” in *ICASSP 2022-2022 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2022, pp. 316–320.
- [7] L. Mazzon, M. Yasuda, Y. Koizumi, and N. Harada, “Sound event localization and detection using foa domain spatial augmentation,” DCASE2019 Challenge, Tech. Rep., June 2019.
- [8] G. Kim, D. K. Han, and H. Ko, “Specmix: A mixed sample data augmentation method for training with time-frequency domain features,” in *Proc. Interspeech 2021*, 2021, pp. 546–550.
- [9] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.