# CoLoC: Conditioned Localizer and Classifier for Sound Event Localization and Detection

Sławomir Kapka, Jakub Tkaczuk
s.kapka@samsung.com, j.tkaczuk@samsung.com

Samsung R&D Institute Poland
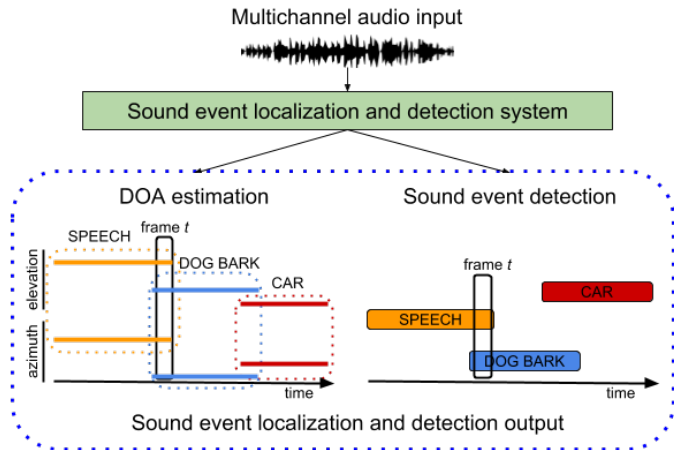
**SAMSUNG**

# SELD problem overview



Figure: https://dcase.community/challenge2022/task-sound-event-localization-and-detection-evaluated-in-real-spatial-sound-scenes

# SELD is about sets

In SELD we predict **sets** of events

$$P(\{\text{class}_i \wedge \text{location}_i\}_{i=1..k}|\text{audio}),$$

where $k \leq N$ and $N$ is the max number of overlapping events.

| Usual approach | Our approach |
| --- | --- |
| Permutation Invariant Training (PIT) | **Sequential Set Generation (SSG)** |

- ▶ Start from empty set $\emptyset$
- ▶ Step 1: Get location of first event $l_1 = \mathbb{E}(l|X, \emptyset)$
- ▶ Step k: Get location of k'th event by conditioning on previously predicted locations $l_k = \mathbb{E}(l|X, \{l_i\}_{i=1..k-1})$
- ▶ Terminate when given a special token $\tau = \mathbb{E}(l|X, \{l_i\}_{i=1..n})$

Based on the output from the localizer we then classify events corresponding to predicted DOAs.

$$P(c_i \wedge l_i | X) = \boldsymbol{P(c_i | X, l_i)} \cdot P(l_i | X).$$

In summary, given:

- ► SSG localizer $\mathbb{E}(l | X, \{l_i\}_i)$,
- ► Location-conditioned classifier $P(c | X, l)$,

we can resolve SELD task.

# CoLoC: Localizer

# CoLoC: Classifier

# Results

We report our results on STARSS22 developement test dataset

Table: Official DCASE metrics; the **boldface** denotes the best scores.

|  | $ER_{20°}$ | $F_{20°}$ | $LE_{CD}$ | $LR_{CD}$ |
|---|---|---|---|---|
| Baseline | **0.71** | 21% | 29.3° | 46% |
| max-ov3 | 0.85 | 32% | 24.7° | **51**% |
| max-ov2 | 0.76 | **33**% | **24.6°** | 49% |

Thank You!

Questions, comments? Contact me:
s.kapka@samsung.com