# PILOT: Introducing Transformers for Probabilistic Sound Event Localization

Christopher Schymura, Benedikt Boenninghoff, Tsubasa Ochiai, Marc Delcroix, Keisuke Kinoshita, Tomohiro Nakatani, Shoko Araki and Dorothea Kolossa
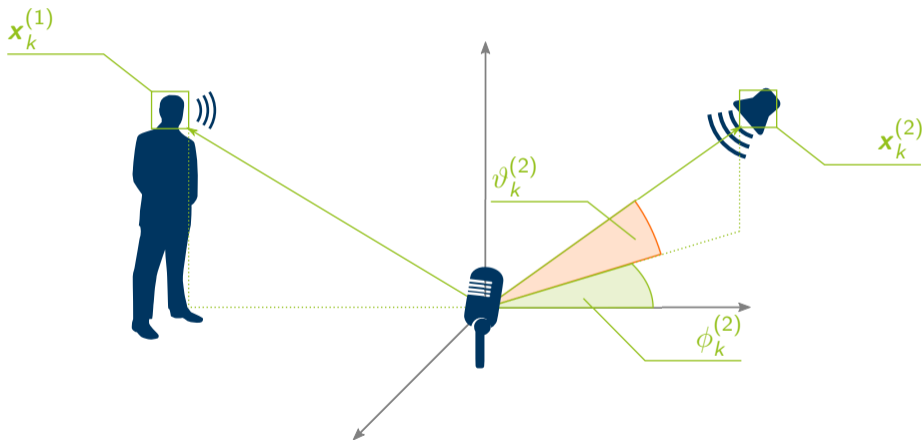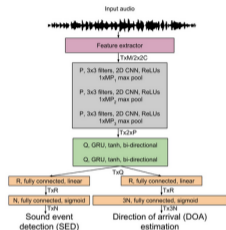
30 August - 3 September 2021

# Problem statement

Sound event localization (SEL) aims at finding the positions of *active* sound sources in the environment
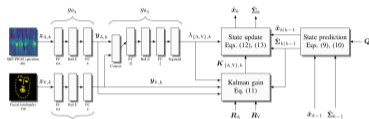
# A selection of previous approaches

Different approaches to SEL have been proposed, primarily focusing on recurrent model architectures
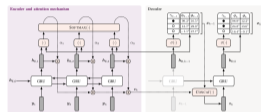


SELDNet[1]

Based on RNNs

Not probabilistic

DSW-BPKF[2]

Based on Kalman Filters

Probabilistic

ADRENALINE[3]
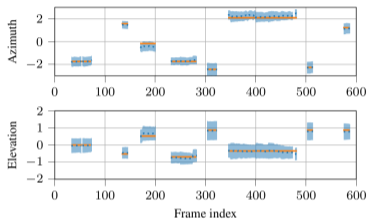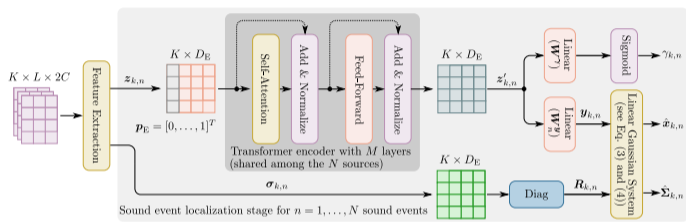
Based on RNNs+Attn

Not probabilistic

[1] Adavanne et al.: "Sound event localization and detection of overlapping sources using convolutional recurrent neural network" (JSTSP 2018)

[2] Schymura et al.: "A Dynamic Stream Weight Backprop Kalman Filter for Audiovisual Speaker Tracking" (ICASSP 2020)

[3] Schymura et al.: "Exploiting Attention-based Sequence-to-Sequence Architectures for Sound Event Localization" (EUSIPCO 2020)

# Proposed framework: PILOT[4]

PILOT is a transformer-based SEL system without recurrent structures and probabilistic output stage



- ▶ Multi-head attention instead of recurrent architectures
- ▶ Differentiable linear Gaussian system as output stage to represent uncertainty

[4] Probabilistic Localization of Sounds with Transformers

# Key results

A transformer-based architecture shows superior SEL performance over recurrent models

▶ PILOT **consistently outperforms** CNN, modified SELDNet and ADRENALINE baseline systems **with statistically significant differences in DoA error** in both simulated and recorded acoustic conditions.

▶ PILOT also shows **improved frame recall** (percentage of correctly identified sound sources) compared to SELDNet and ADRENALINE.

▶ The **probabilistic output stage** allows the model to represent the **estimation reliability** associated with individual sound source DoA estimates.