
Machine Listening for Music and Sound Analysis

Lecture 4 - Environmental Sound Analysis

Dr.-Ing. Jakob Abeßer

Fraunhofer IDMT

Jakob.abesser@idmt.fraunhofer.de

<https://machinelisting.github.io>

Overview & Learning Objectives

- Introduction
- Sound Event Detection
- Acoustic Scene Classification
- Acoustic Anomaly Detection
- Application Scenarios

Introduction

Motivation

- Sound carries information about our environment
- Challenging attempt to mimic the human's abilities
 - Environment perception
 - Context-awareness & localization of sound sources
 - Acoustic scene understanding
- Complementary sensory path to vision → multimodality
- Related to other content analysis domains (speech, music)

Introduction

Environmental Sounds (Recap)

- Sound sources
 - Nature, climate, humans, machines, etc.
- Sound characteristics
 - Structured or unstructured, stationary or non-stationary, repetitive or without any predictable nature
- Sound duration
 - From very short (gun shot, door knock, shouts) to very long and almost stationary (running machines , wind, rain)



Introduction

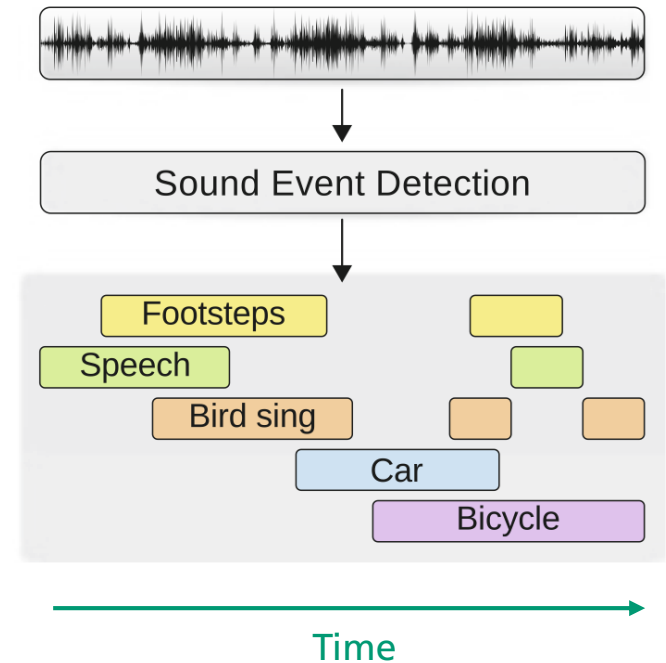
Tasks / Categories

- Sound event detection (SED)
- Acoustic scene classification (ASC)
- Acoustic anomaly detection (AAD)

Sound Event Detection

Introduction

- Sound event detection → 2 simultaneous tasks
 - Segmentation (detection of temporal boundaries)
 - Classification (type of sound)
- Sound polyphony
 - Number of simultaneous sounds
 - Depends on the acoustic scene composition & sound sources



Sound Event Detection

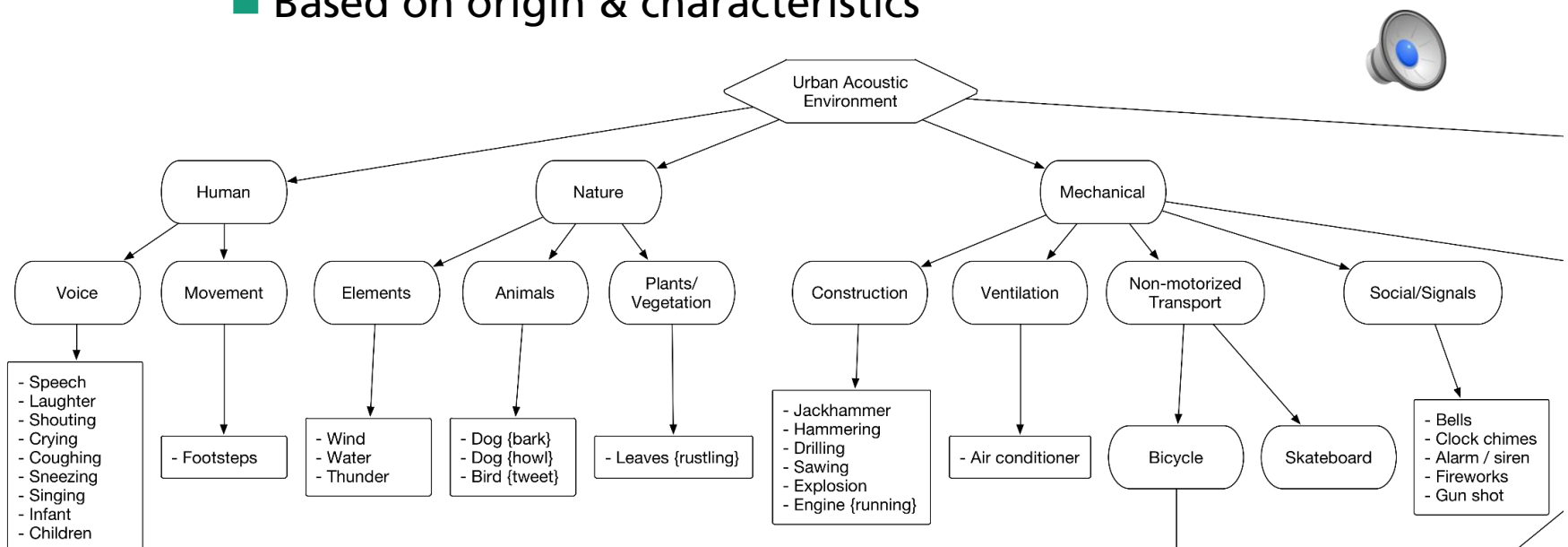
Introduction

■ Sound source categories

■ Humans, animals, vehicles, tools, machines, climate, ...

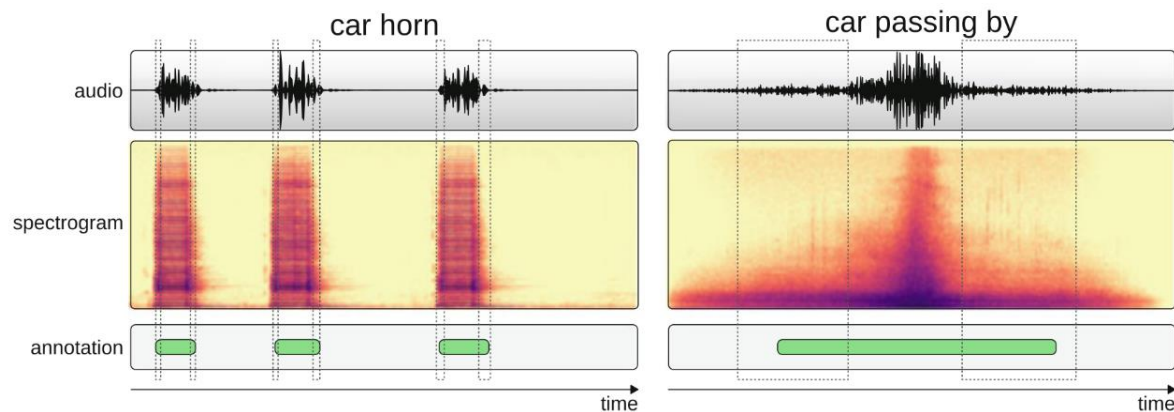
■ Hierarchies of sounds (e.g. urban sounds)

■ Based on origin & characteristics



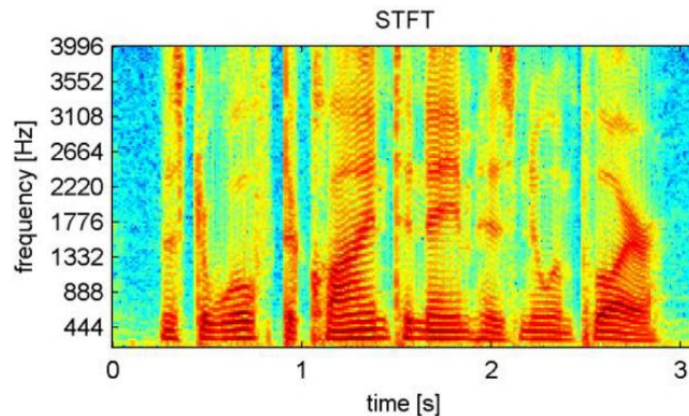
Sound Event Detection Challenges

- Large range of different timbre characteristics
 - Short transients, noise-like signals, harmonic / inharmonic signals
- Different sound durations
 - Short (gun shot, door knock) → long / stationary (machines, wind)
- Ill-defined temporal boundaries
 - Complicates annotation & detection



Sound Event Detection Challenges

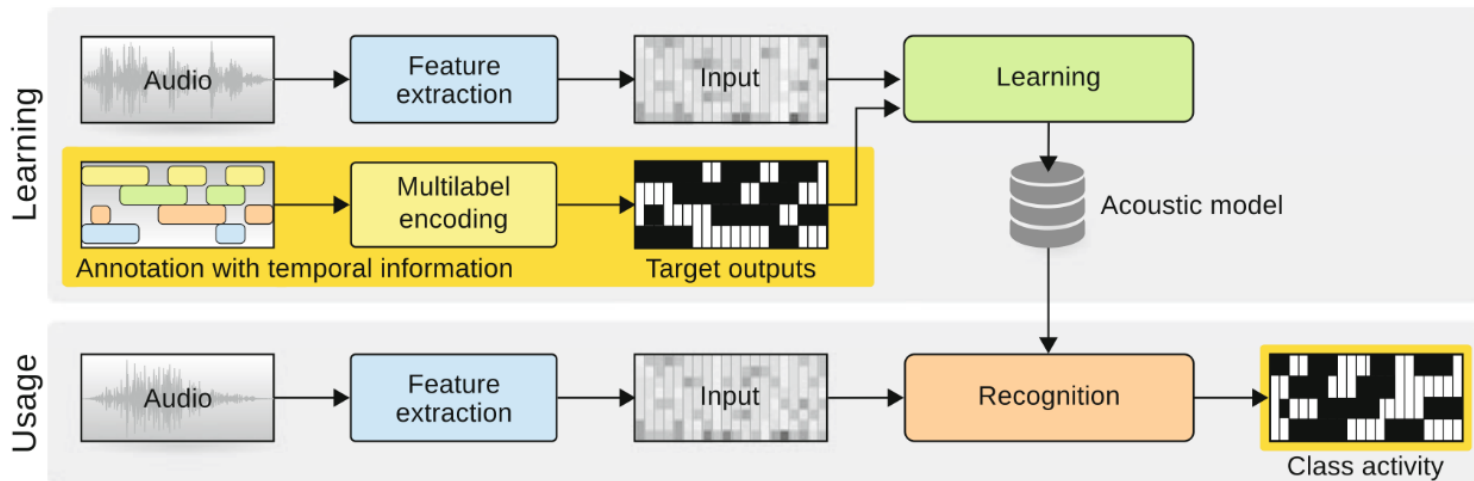
- Sound appear in the foreground & background
 - depending on relative sound source position
- Non-local / sparse energy distribution
 - Fundamental frequency & overtones



- Sounds are „transparent“
 - Phase-dependent overlap, possible cancelations

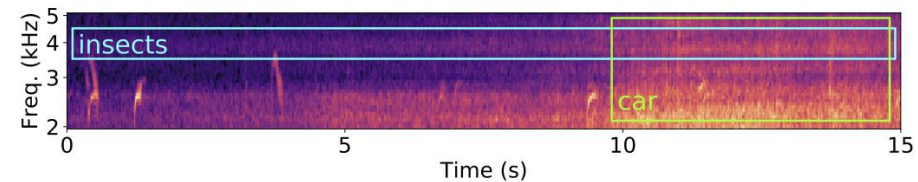
Sound Event Detection Pipeline

- Supervised learning pipeline
 - Feature extraction & pre-processing
 - Label encoding
 - Acoustic modeling

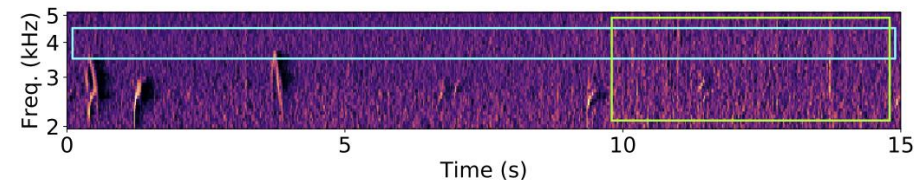


Sound Event Detection Pipeline

- Feature extraction
 - 1D features (audio samples) → “end-to-end learning”
 - 2D features (mel-spectrogram, STFT)
- Feature pre-processing
 - Log-magnitude scaling
 - Per-channel energy (PCEN)
 - Dynamic range compression
 - Adaptive gain control
 - Suppresses stationary (background) noise



(a) Logarithmic transformation.

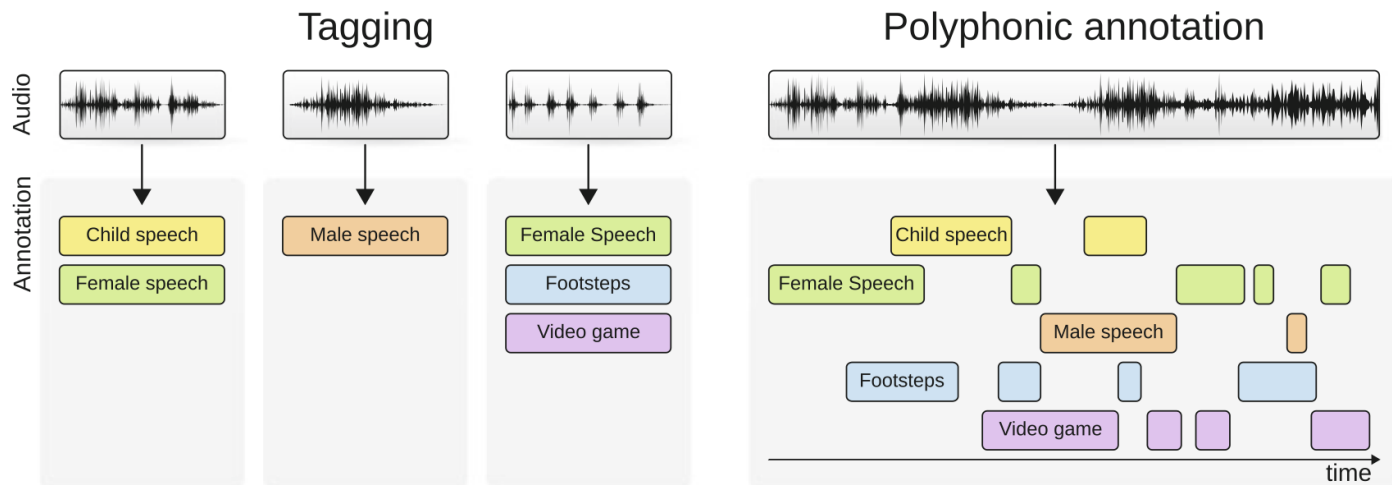


(b) Per-channel energy normalization (PCEN).

Sound Event Detection Pipeline

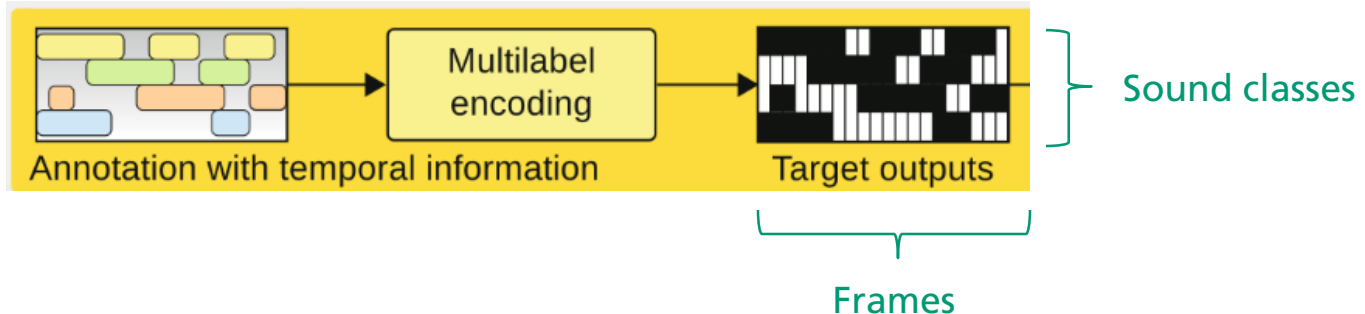
■ Annotation

- Quality of “ground truth”? (limited agreement / reliability)
- Different granularities
 - Tagging / Global level (“weak” labels) → cheap
 - Event-level (“strong” labels) → expensive



Sound Event Detection Pipeline

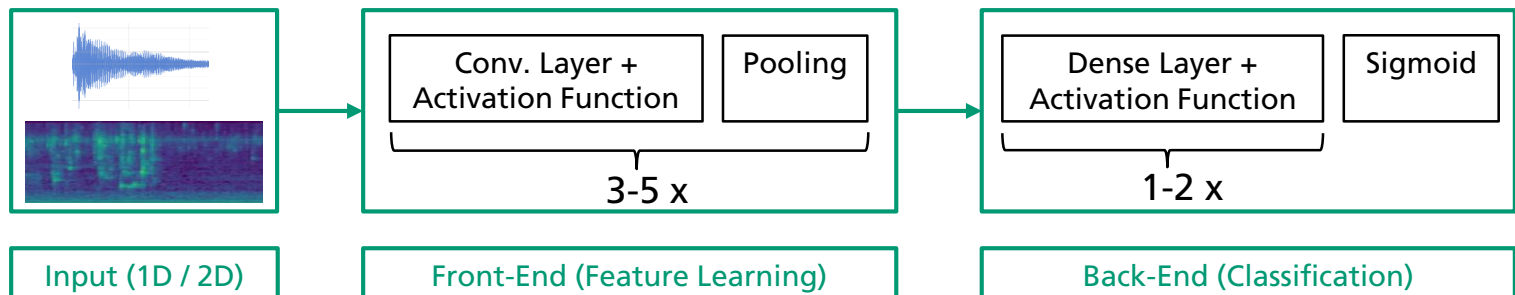
- Label encoding
 - Binarized sound activity (0/1)
 - Multilabel classification
 - 1 (independent) binary detector per class
 - Temporal resolution (duration of each annotated time frame)



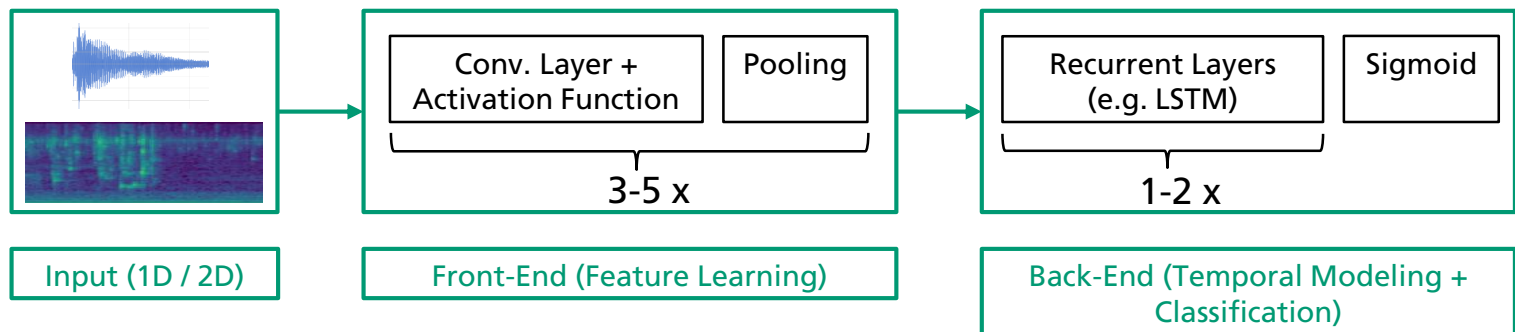
Sound Event Detection Pipeline

■ Typical neural network architectures

■ CNN



■ CRNN



Sound Event Detection Pipeline

■ Data Augmentation

- Cope with limited amount of training data
- Approach 1: Apply signal transformations
 - Adding noise
 - Time-stretching/Pitch-shifting
 - Mix-up data augmentation [Zhang]
 - Random mix of two items (audio & targets)
 - Example: $0.7 \times \text{car} + 0.3 \times \text{speaker}$
- Approach 2: Data synthesis
 - Generative Neural Network Models (e.g. Generative Adversarial Networks (GAN), SampleRNNs)

Sound Event Detection Evaluation

- Recap: Binary classification evaluation

- True/false positives (TP/FP)

- True/false negatives (TN/FN)

- Metrics

- Precision

- Recall

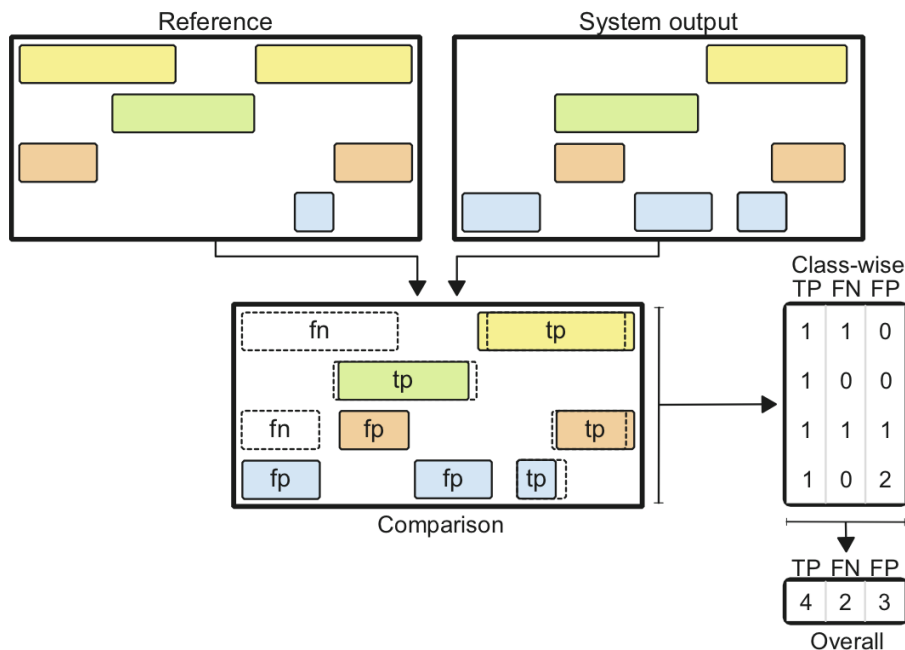
- Accuracy

- F-score

		Prediction			
		1	0		
Annotation	1	TP <i>true positives</i>	FN <i>false negatives</i>	True Positive Rate Sensitivity Recall $R = \frac{TP}{TP+FN}$	
	0	FP <i>false positives</i>	TN <i>true negatives</i>	False Positive Rate $FPR = \frac{FP}{FP+FN}$ Specificity $\text{Specificity} = \frac{TN}{FP+FN}$	
		Precision $P = \frac{TP}{TP+FP}$		Accuracy $ACC = \frac{TP+TN}{TP+TN+FP+FN}$	

Sound Event Detection Pipeline

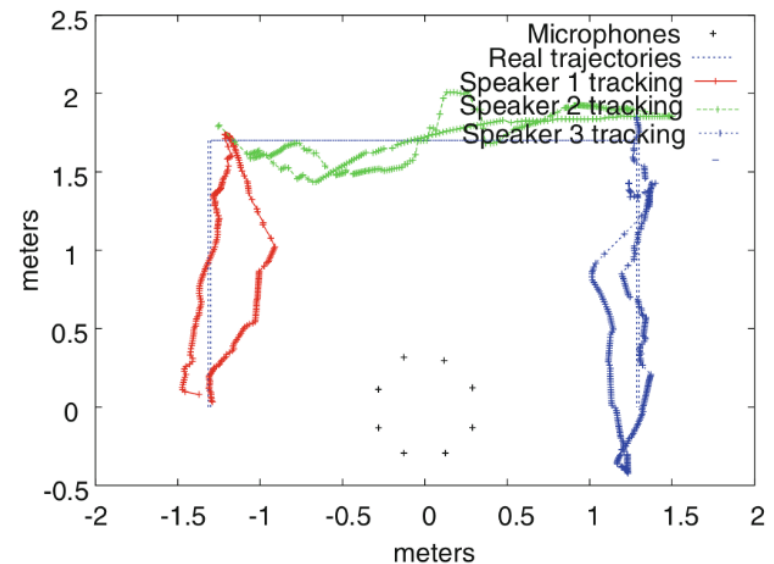
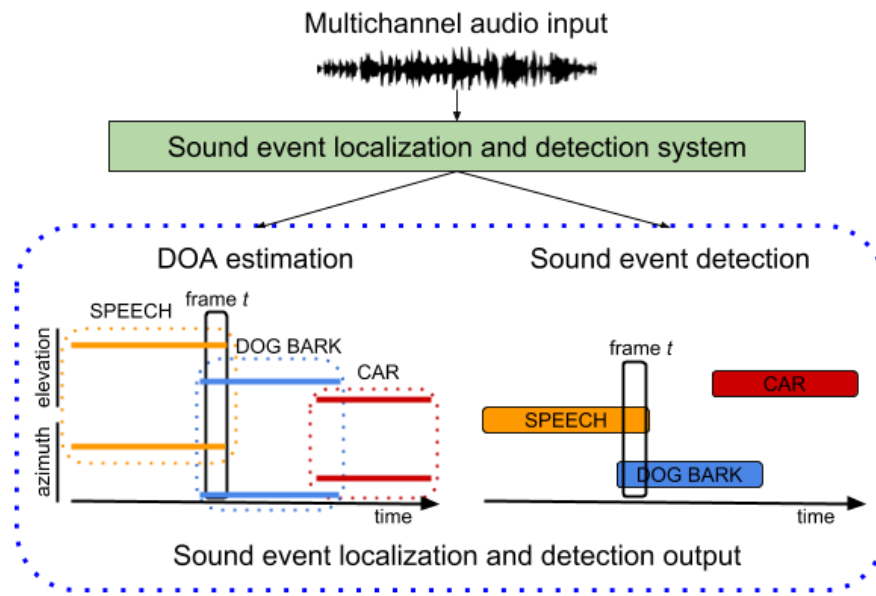
- Evaluate SED → binary classification results on a frame-level
- Compare reference with predictions
- Count TP/FN/FP → aggregate over time → compute metrics



Sound Event Detection

Related tasks

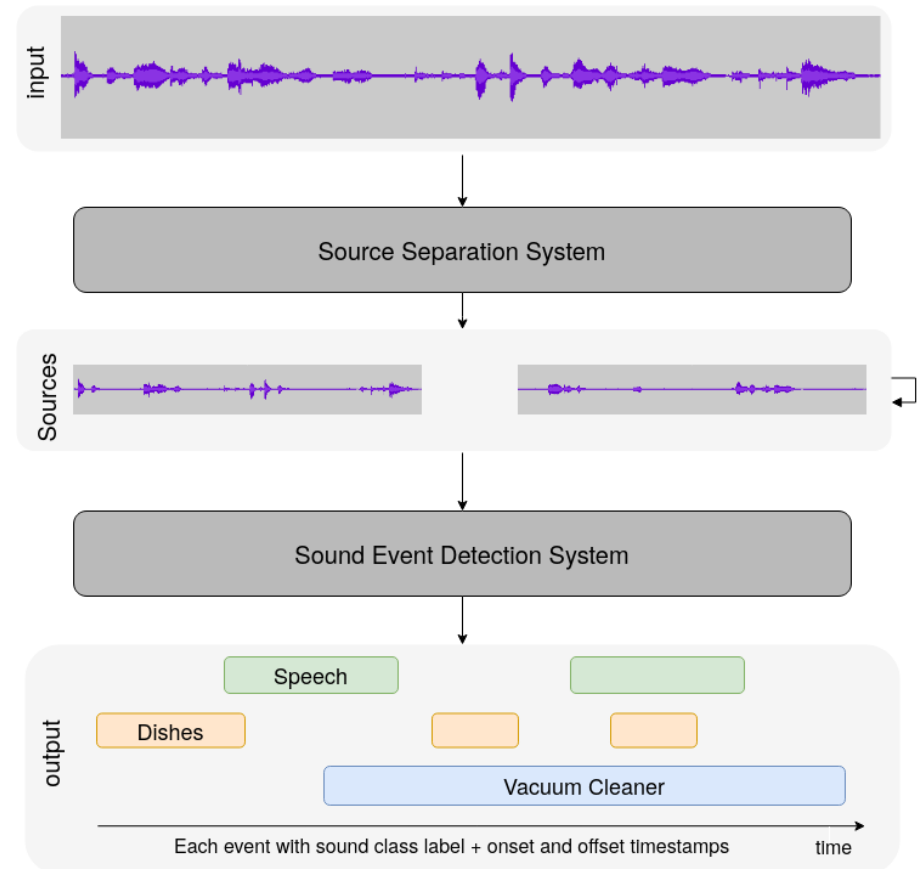
- Sound event localization & tracking
 - Multichannel audio recordings (e.g., first-order ambisonic microphones)
 - Estimate direction-of-arrival (DOA) & track source movement



Sound Event Detection

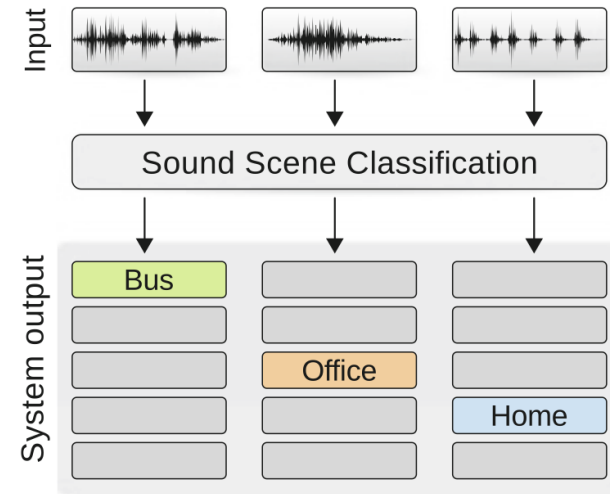
Related tasks

- Source separation
 - Prior to sound event detection
- Chicken-egg problem
 - Alternative: sound-informed source-separation



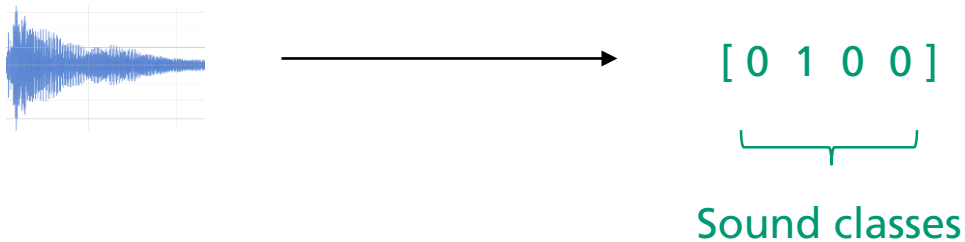
Acoustic Scene Classification Task

- Acoustic scene classification (ASC)
 - Multi-class (1 of N) classification scenario
 - Summative label (tagging)
- Acoustic scene
 - Typical set of sounds
 - Example: office
 - Keyboard clicks
 - Human conversations
 - Printer
 - Air conditioner



Acoustic Scene Classification Pipeline

- Label encoding
 - One-hot-encoded (global) target
- Example
 - 4 scene classes (bus, office, home, forest)
 - Encoding of an office recording



Acoustic Scene Classification Pipeline

- Network architectures
 - Similar to SED (CNN & CRNN)
- Differences
 - Temporal result aggregation within network
 - Dense layer / pooling
 - Final layer: softmax activation function (multiclass classification)
- Current Research Topics
 - Attention → learn to focus on spectrogram regions
 - Open-set classification → detect unknown classes
 - Transfer learning → fine-tune pre-trained models with less data

Acoustic Anomaly Detection Task

■ Goal

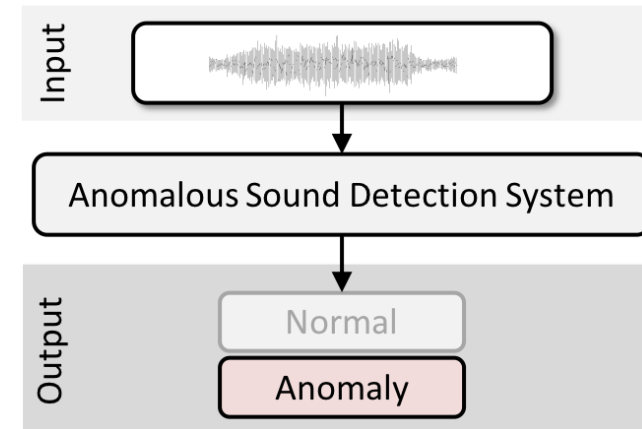
- Detect deviations from “normal” state
- Identify whether emitted sound from target object is normal or anomalous

■ Challenges

- Often only training examples for normal state available
- Acoustic anomalies are often subtle compared to louder background noise

■ Application Scenarios

- Detecting machine failures
- Intrusion detection (glass break...)



Acoustic Anomaly Detection Approaches

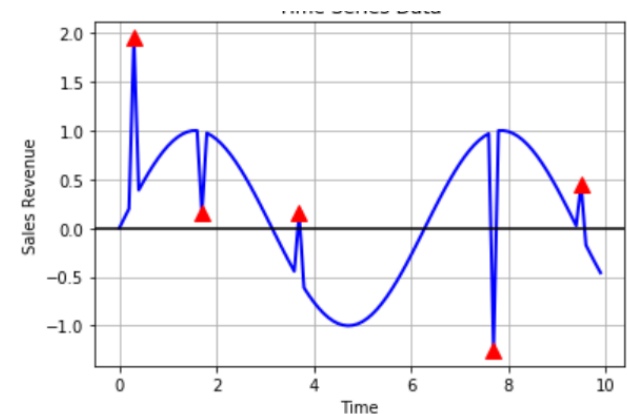
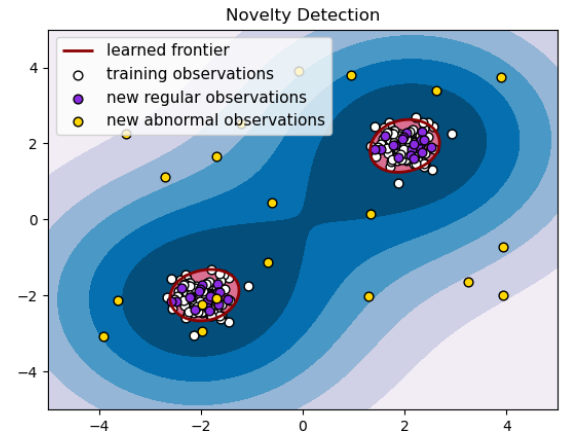
■ Traditional methods

■ Distribution outlier detection

- Modelling normal state distribution
- Detect distribution outliers
- E.g.: One-class GMM / SVM

■ Time-series analysis

- AD via local deviation from prediction
- E.g.: Autoregressive models, Hidden-Markov-Models (HMM)

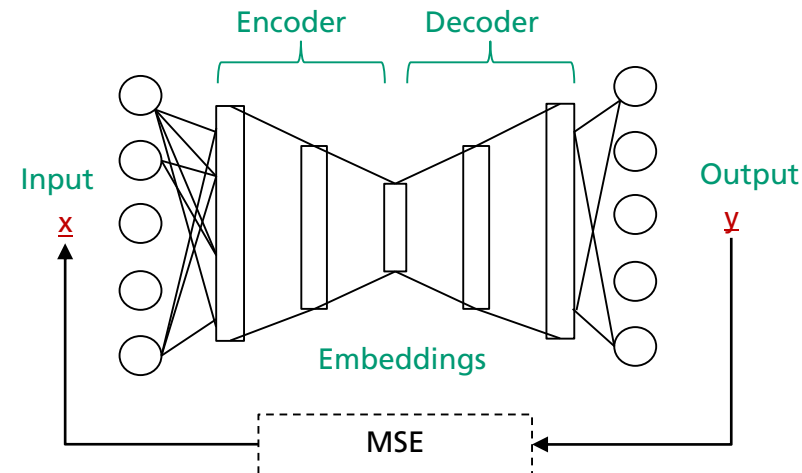


Acoustic Anomaly Detection Approaches

- Novel methods

- Autoencoder (encoder → decoder) models

- Idea: normal sounds can be better reconstructed than unknown anomalous sounds
 - Dense, convolutional, variational AE
 - Interpolation DNN
 - Interpolate spectrogram frame from surrounding frames



Application Scenarios

Urban Noise Monitoring



- Joint R&D project (2016 – 2018)
 - Fraunhofer IDMT, IMMS, SSJ GmbH, BE
- Goal
 - Develop distributed sensor network for
 - Sound level measurement
 - Sound classification
- Approach
 - Mobile sensor units
 - Raspberry Pi 3, quad-core ARM, 1GB RAM
 - Battery + MEMS microphones
 - Sensor locations (light poles)



Application Scenarios

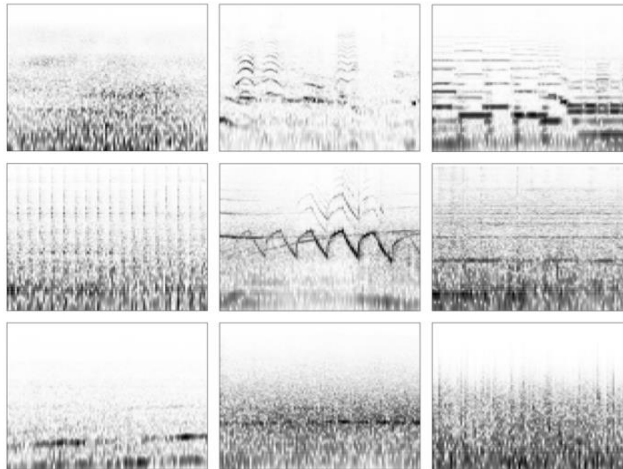
Urban Noise Monitoring



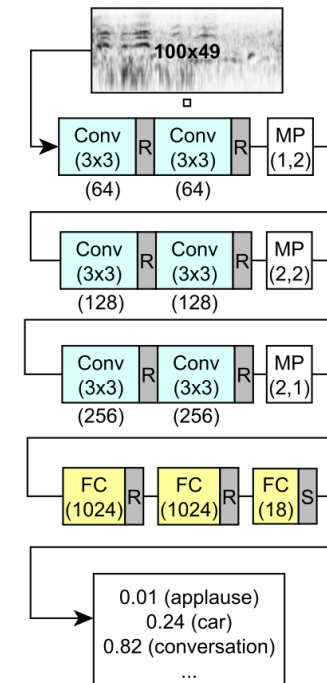
■ Measurements

- Different loudness values (8/s)
- Sound event detection (1/s)
 - 9 sound event classes (car, conversation, music, roadworks, siren, train, tram, truck, wind)

Spectrogram examples (2 s long)



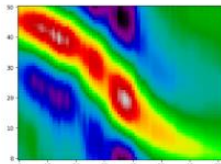
■ CNN architecture



Application Scenarios

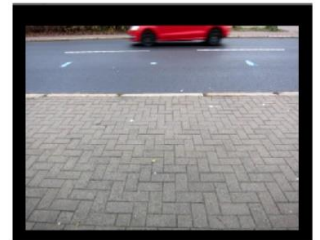
Traffic Monitoring

- Traffic monitoring
 - Vehicle detection & direction of movement
 - Stereo channel cross-correlation



Example: movement left → right

- Vehicle type classification (car, truck, bus, motorcycle)
- Challenges
 - Microphone type
 - Vehicle speed
 - Street surface quality & weather conditions

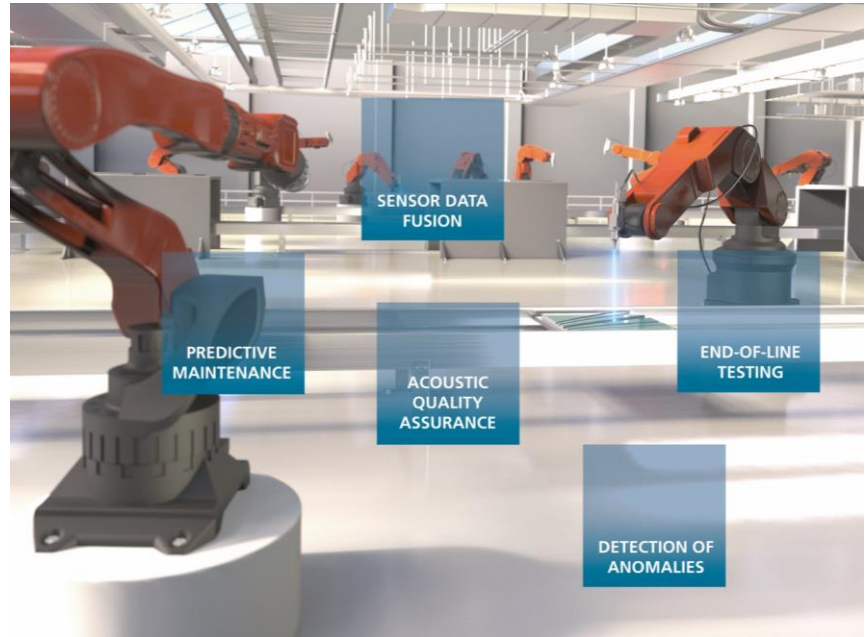


Application Scenarios

Production Chain Monitoring

■ Challenges

- Real-time analysis & classification of industrial sounds
- Energy-efficient AI algorithms
- Sound variations due to different machine states
- Acoustic anomalies subtle compared to background noises



Summary

- Introduction
- Sound Event Detection
 - Challenges
 - Pipeline
 - Evaluation
- Acoustic Scene Classification
- Acoustic Anomaly Detection
- Application Scenarios
 - Urban Noise Monitoring
 - Traffic Monitoring
 - Product Chain Monitoring