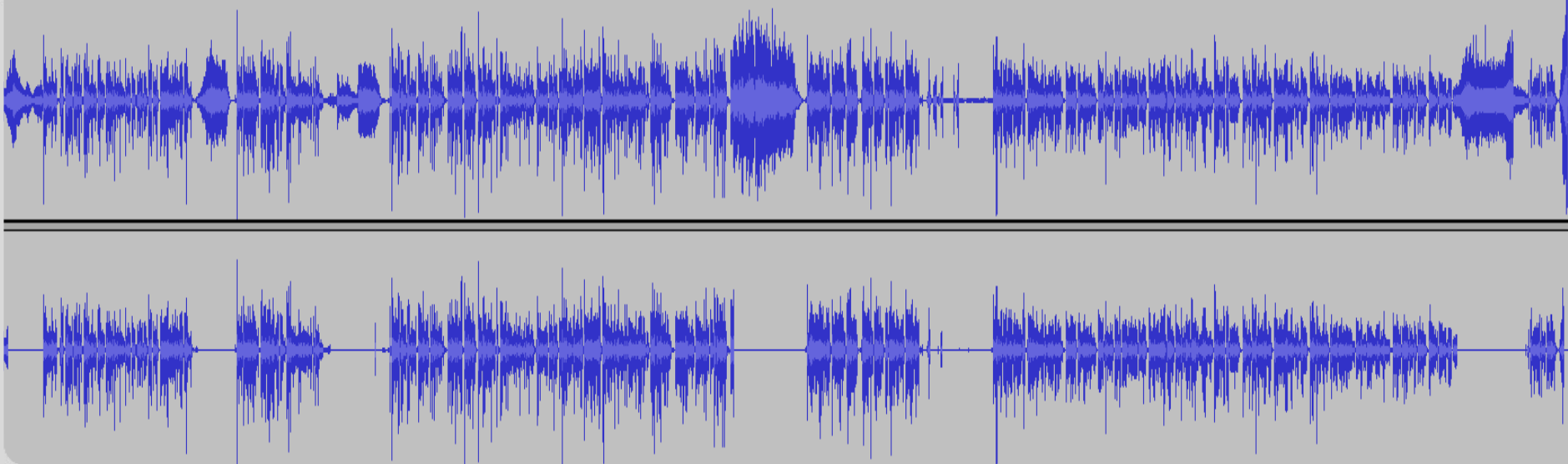# Audio Segmentation For Robust Real-Time Speech Recognition Based on Neural Networks

**Micha Wetzel, Matthias Sperber, Alex Waibel**

Institute for Anthropomatics and Robotics,
Karlsruhe Institute of Technology, Germany
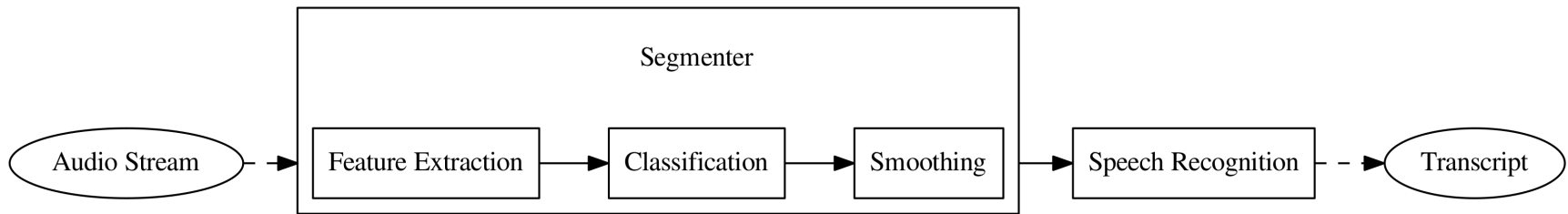
# Problems in Speech Recognition

- Garbage in, garbage out
  - „ Thank you you You Can you you you ..."
  - „ you see you may and you see that them ..."
  - „ if if if if if if f F f.."
- Causes high latency

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Challenges

- Real-time
  - Limited temporal information
  - Fast algorithm needed
- Classifying speech correctly
- No clear segment borders
- Humans need ~200ms to classify reliably[1]

[1] H. Harb et al. Signal Processing and Its Applications Vol 2. 2003

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Approach

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Feature Extraction

- 10 ms frames
- MFCC and ZCR
- 13 frames stacked
- Dimensionality reduction



$$\bar{x}, \sigma, \sigma^2$$

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for Robust Real-Time Speech Recognition Based on Neural Networks
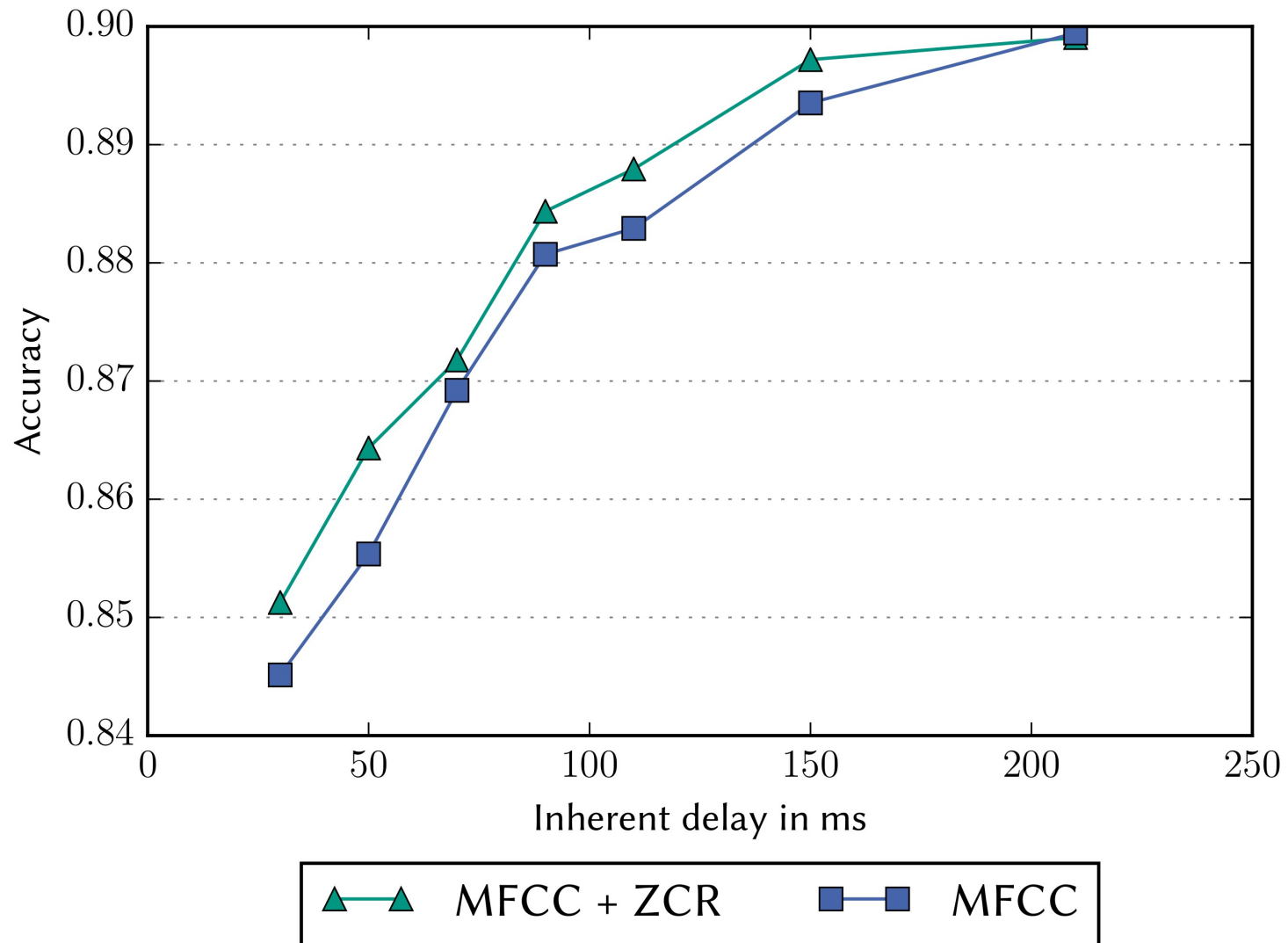
Institute for Anthropomatics and Robotics

# Classification

- Multilayer perceptron
- < 1ms computational latency
- ~70 ms inherent latency
- 87 % accuracy
- 2Class: 95 % accuracy

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Classification Results



Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Classification - Problems

- Class fluctuation
- Misclassifications
  - → ASR performs worse

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for Robust Real-Time Speech Recognition Based on Neural Networks
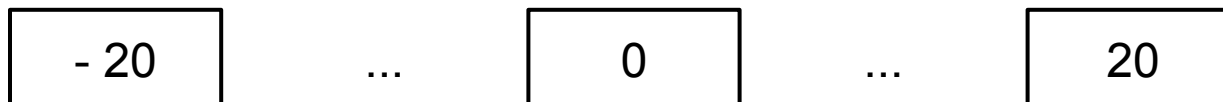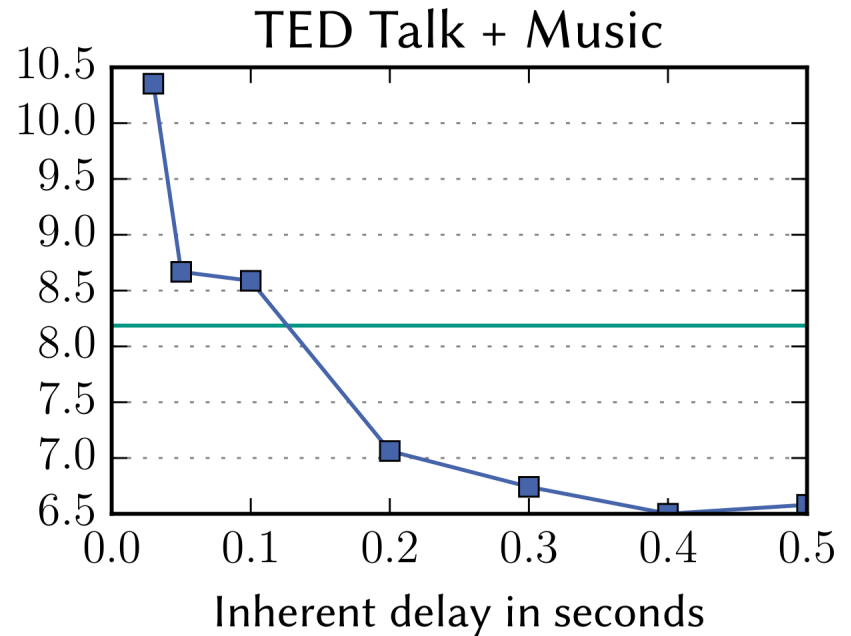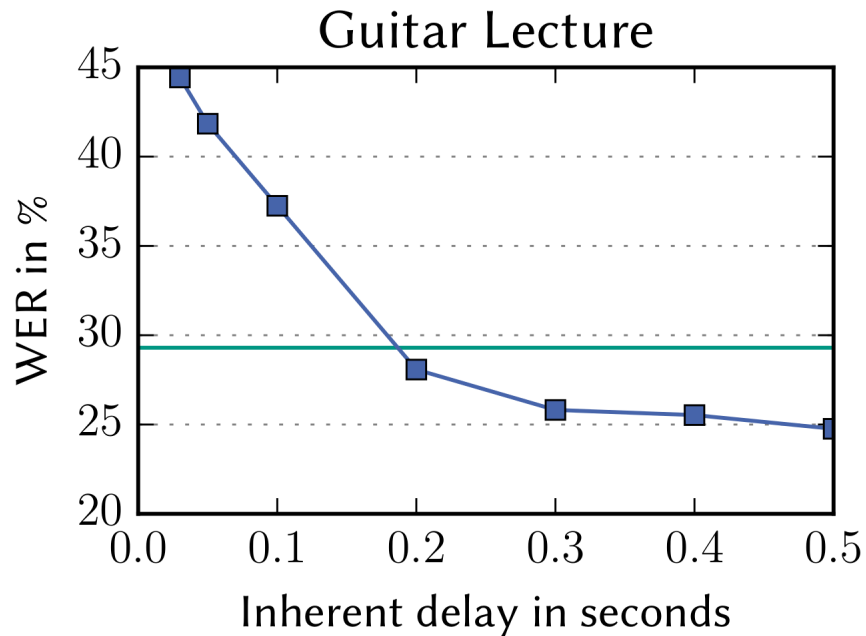
Institute for Anthropomatics and Robotics

# Smoothing

- Need more temporal information
- Class off adjacent frames is correlated
- Make use of knowledge about the past

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Mode Smoothing

- Mode of adjacent frames
- Remove misclassifications
- Additional temporal information
    $\rightarrow$ additional 200 ms latency

- Erosion & dilation tested but not needed

| - 20 | ... | 0 | ... | 20 |
|---|---|---|---|---|

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Mode - Results



- Tradeoff latency ↔ accuracy

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics
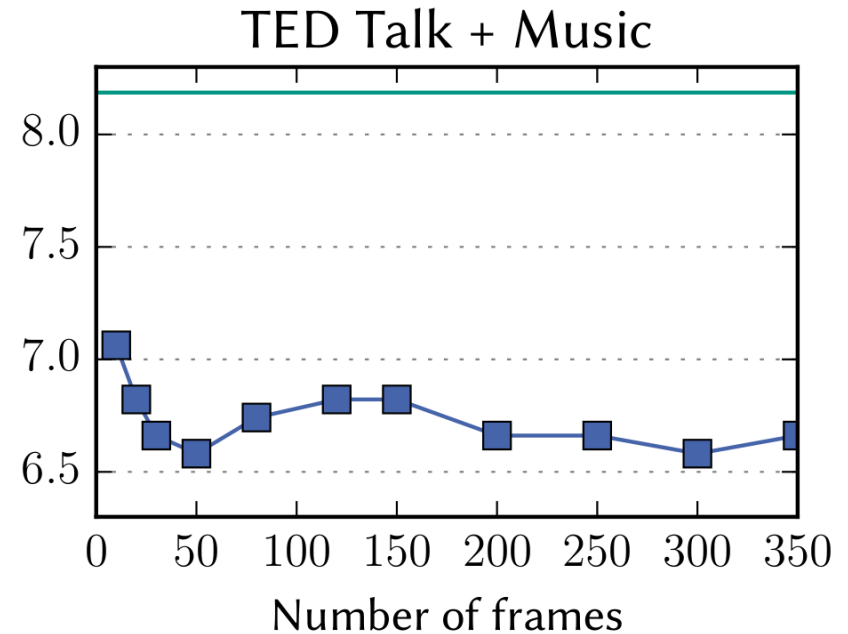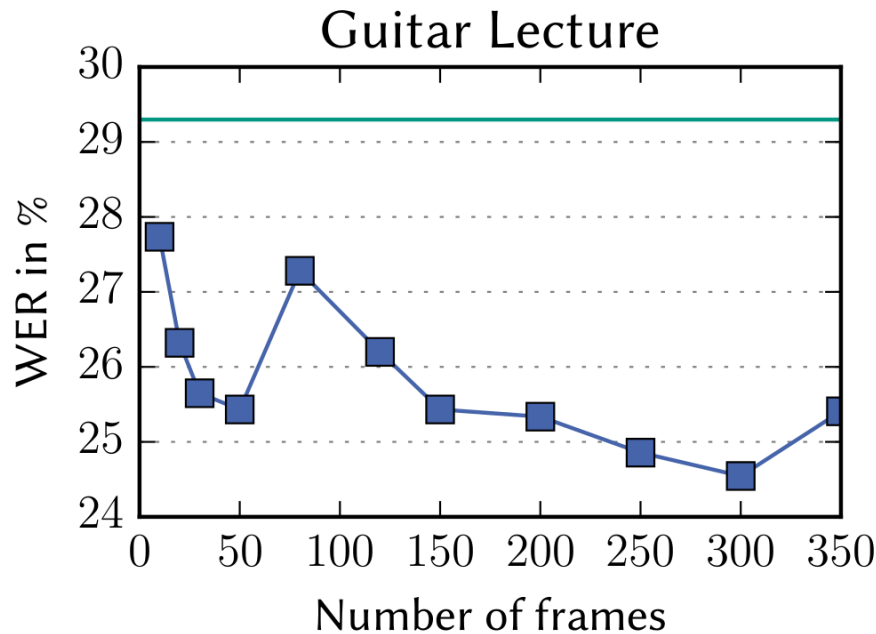
# Minimum Change Support Smoothing

- Use knowledge about past
- Create big segments
- Favour speech

```
                                        ┌──────────┐
                                        │ current  │
                                        │ segment  │
                                        │  class   │
                                        └──────────┘
                                             ↕  = ?
┌─────────┐              ┌─────────┐    ┌─────────┐
│  - 300  │      ...     │   - 1   │    │    0    │
└─────────┘              └─────────┘    └─────────┘
```

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Minimum Change Support - Results



- With mode smoothing
- No additional latency

Institute for Anthropomatics and Robotics

# Segmentation Results

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Summary

- Fast computation
- ~270 ms latency
- Removes 39 % of segmentation based errors
- Language independent

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# The End

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Classification Results

- 3Class accuracy: 87%
- 2Class accuracy: 95%

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Segmentation Results



Rate of Resolved Segmentation Errors = 1 – (S – M) / (N - M)

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics

# Music Example

Micha Wetzel, Matthias Sperber, Alexander Waibel - Audio Segmentation for
Robust Real-Time Speech Recognition Based on Neural Networks

Institute for Anthropomatics and Robotics