

Sprachverarbeitung I / 11 HS 2016

Spracherkennung: Merkmalsextraktion

Buch: Kapitel 11.7 und 11.8

Beat Pfister



Programm heute:

- Vorlesung:
- Aufnehmen einer Äusserung
 - Sprachmerkmale für die Spracherkennung
 - Mel-Cepstrum als Sprachmerkmal
- Übung:
- ★ Sprachmerkmale zur Lautunterscheidung

Spracherkennung

Aufgabe: Ermitteln der Aussage aus einem Sprachsignal

Annahme: Sprachsignal enthält eine (ganze) Äusserung!

Frage: Wie kommt ein Sprachsignal mit einer ganzen Äusserung in den Computer bzw. in den Spracherkenner?

→ Äusserung muss aus dem vom Benutzer aufgenommenen Signal detektiert werden!

Detektieren einer Äusserung

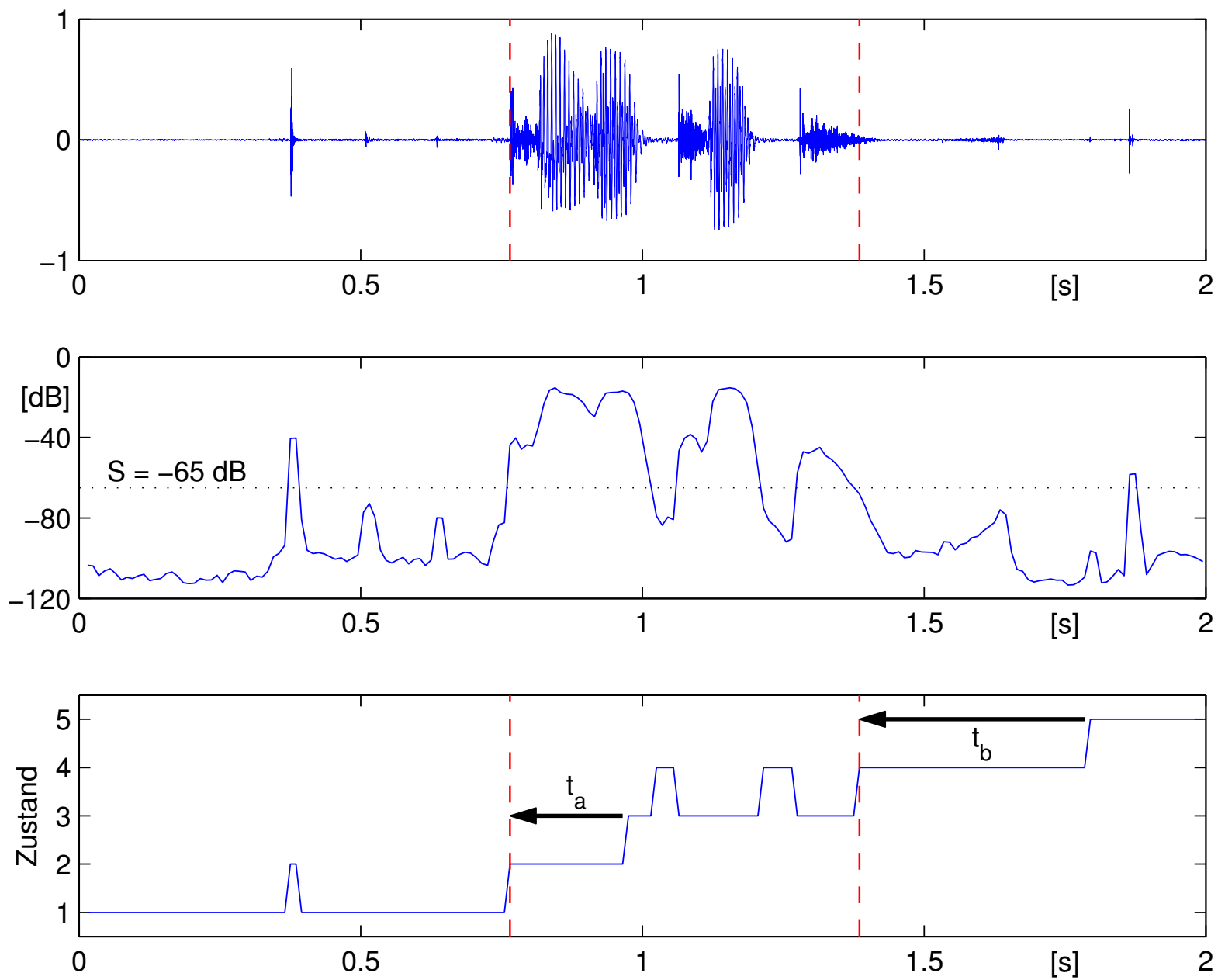
(Buch: Kapitel 11.8)

manuell:	Sprechtaste (push-to-talk)	<ul style="list-style-type: none">– robust bzgl. Umgebungslärm– eine Hand muss frei sein– Äusserung wird oft abgeschnitten
halbautomatisch:	Starttaste (tap-and-talk)	<ul style="list-style-type: none">– ziemlich robust– eine Hand muss frei sein
automatisch:	Anfangs- und Endpunktdetektion	

Automatische Anfangs- und Endpunktdetektion

Anforderungen bei Dialogsystemen:

- Flexible Antwortzeit des Benutzers
 - Zeitfenster für Eingabe genügend lang
- Kurze Reaktionszeit des Systems
 - Aufnahmedauer darf nicht fix sein
- Robust für leise oder kurze Störungen
 - Geräuschpegel viel tiefer als Leistung des Sprachsignals
 - starke Störungen kürzer als jede gültige Äusserung



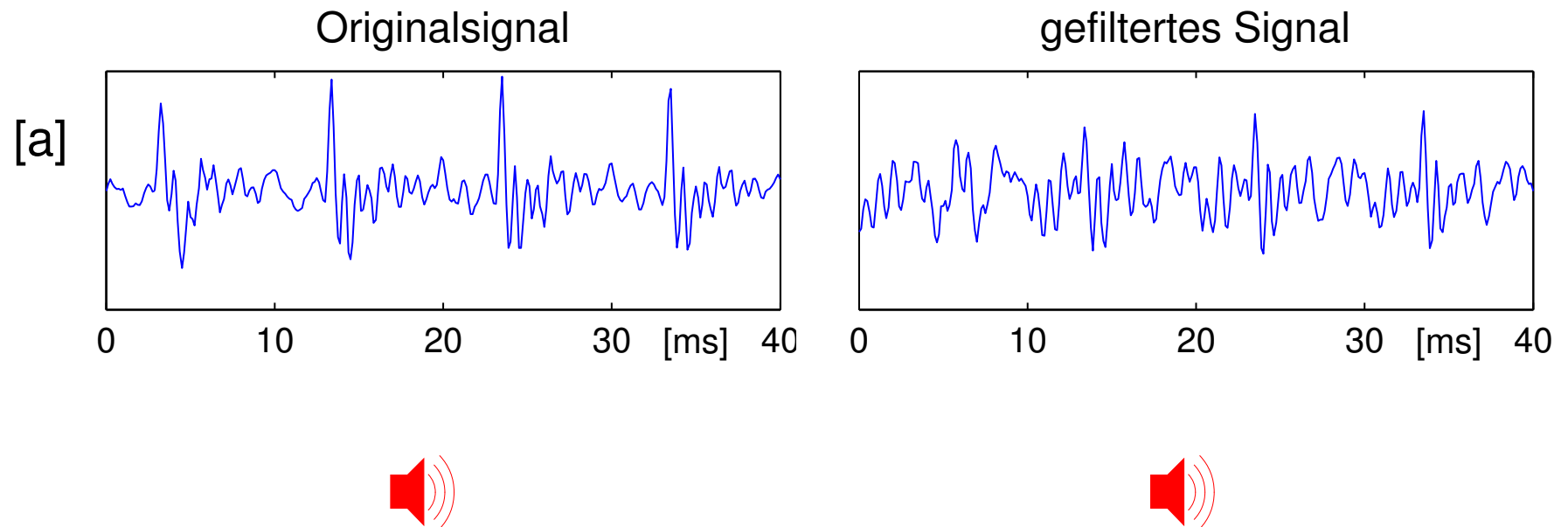
Spracherkennung

Resultat der Anfangs- und Endpunktdetektion:

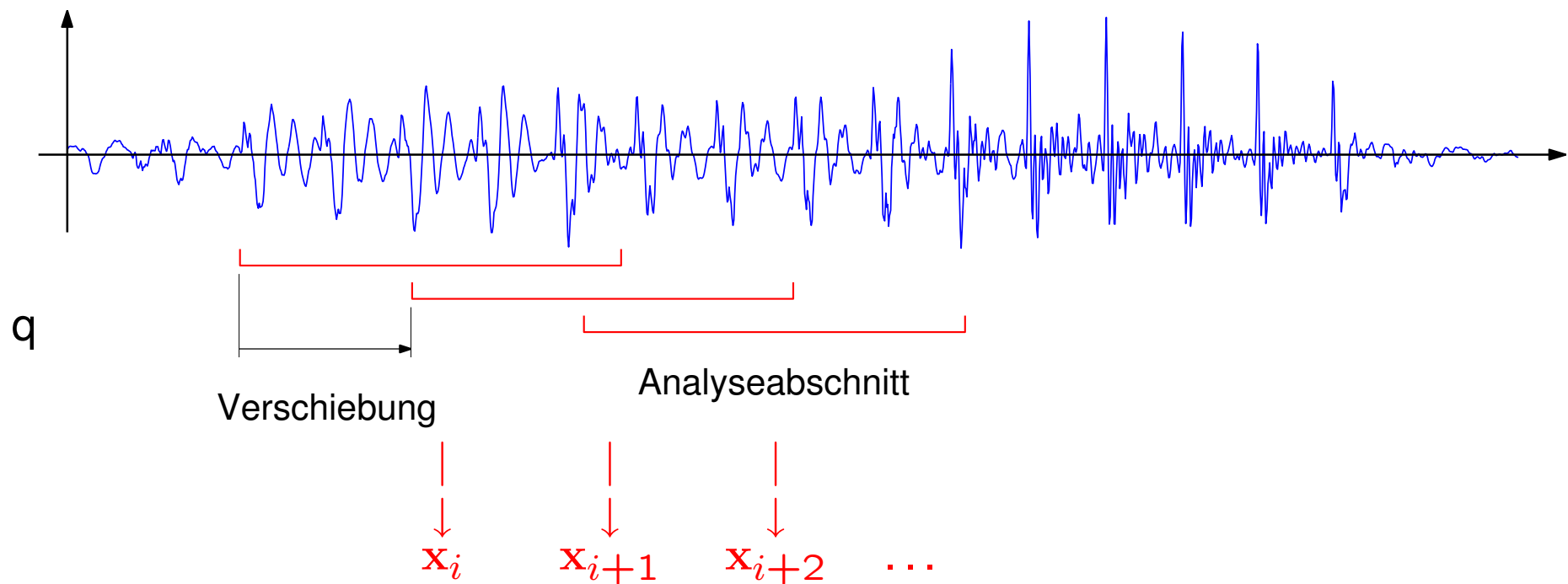
—→ Sprachsignal mit einer Äusserung
Input des Spracherkenners

Sprachsignale im Zeitbereich

Signal kann nicht direkt für die Erkennung verwendet werden!



Merkmalsextraktion Kurzzeitanalyse des Sprachsignals

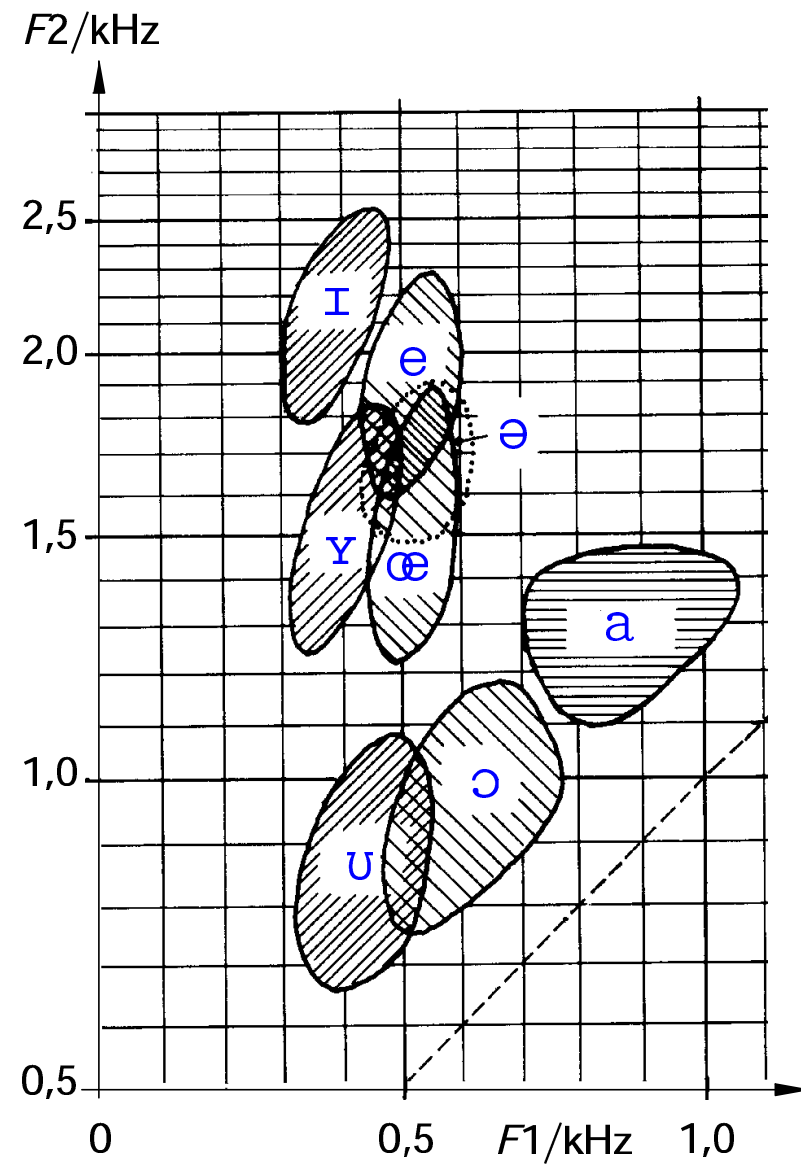
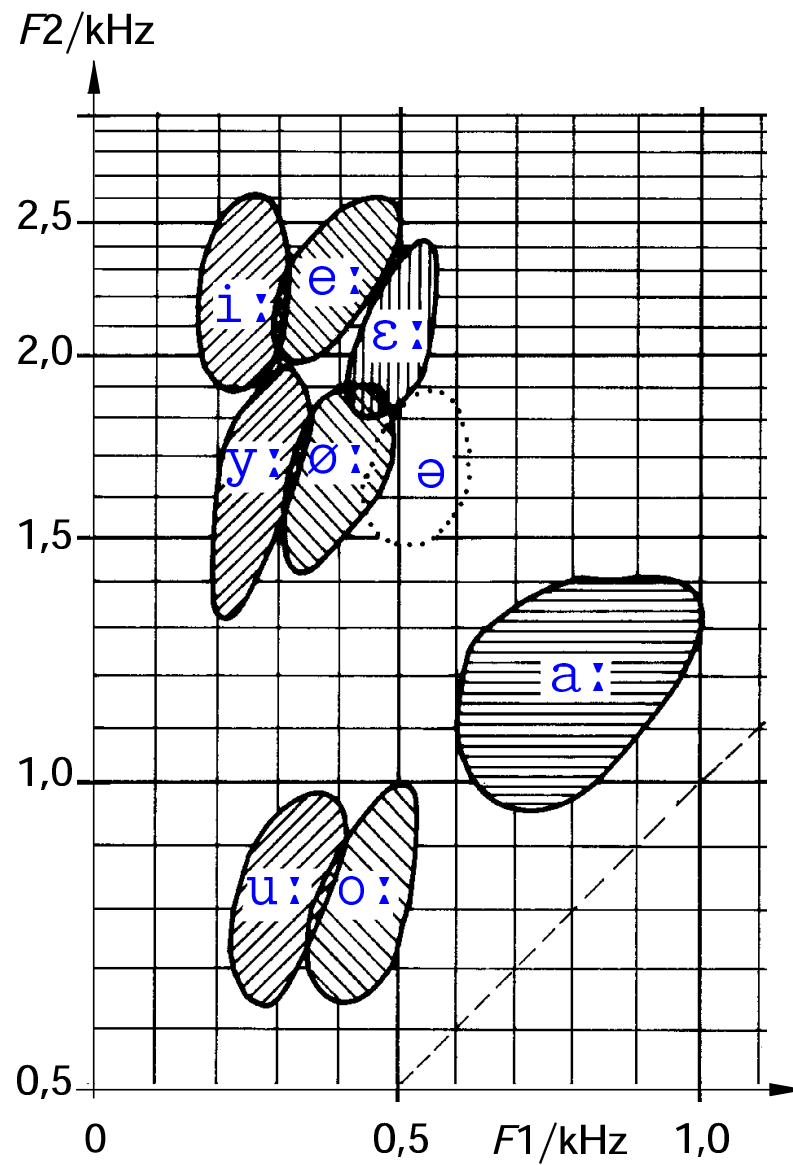


Merkmalssequenz: $\mathbf{X} = x_1 x_2 x_3 \dots x_T$

Fragen

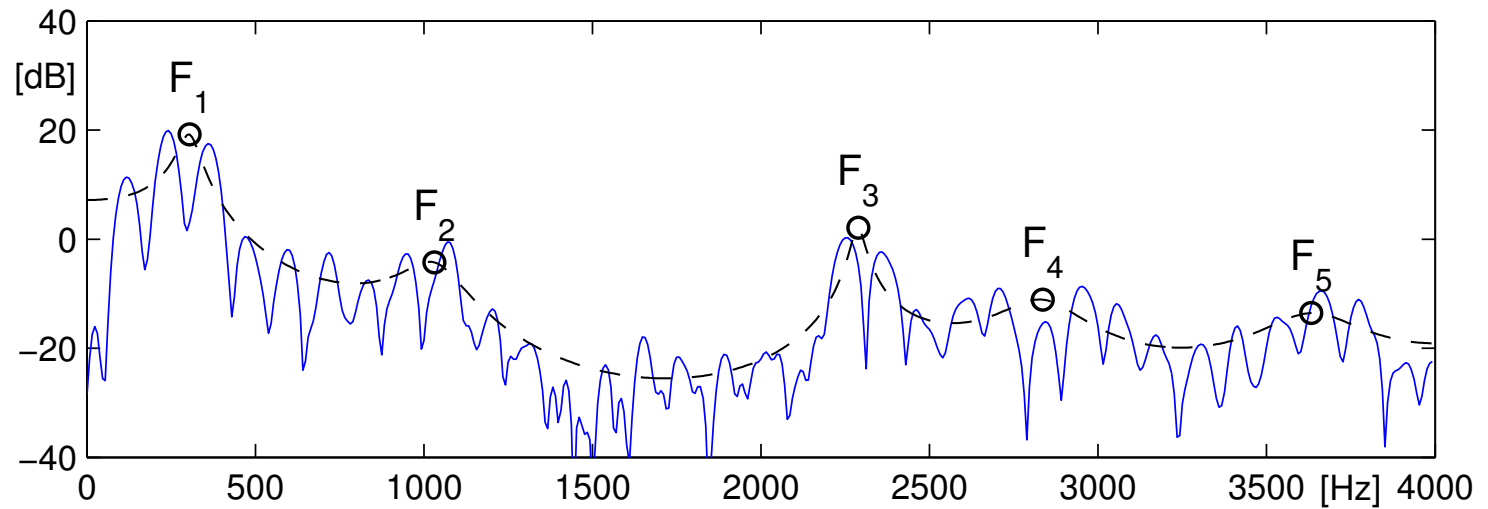
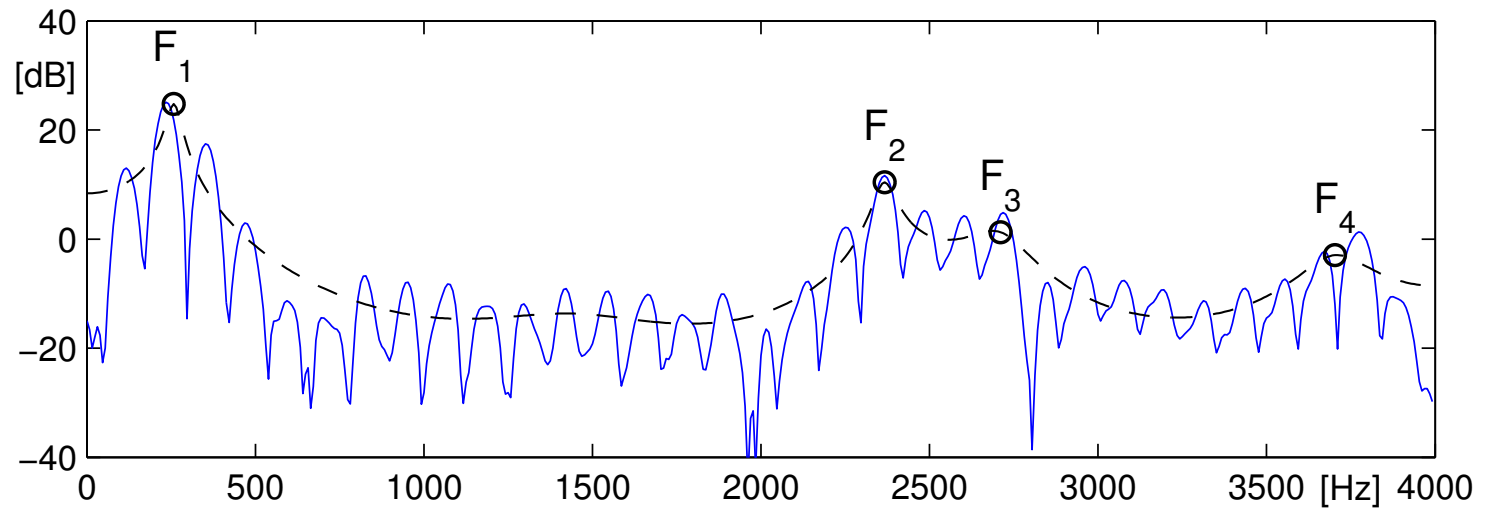
Welche Merkmale sind für die Spracherkennung geeignet?

Gibt uns die Phonetik nützliche Hinweise?



Formanten

Laut: [i:]



Anforderungen an Sprachmerkmale

Grundsatz: Ein Merkmal \mathbf{x} soll Laute $L_1, L_2, L_3 \dots$ unterscheiden

Problem: Sprachsignale eines Lautes L_i stark verschieden

Frage: Wie kann man ein Merkmal \mathbf{x} beurteilen?

Kriterium: Eigendistanzen d_E vs. Kreuzdistanzen d_K

- $d_E = d(\mathbf{x}\{L_i\}, \mathbf{x}\{L_j\})$ mit $L_i = L_j$ \longrightarrow klein
- $d_K = d(\mathbf{x}\{L_i\}, \mathbf{x}\{L_j\})$ mit $L_i \neq L_j$ \longrightarrow gross

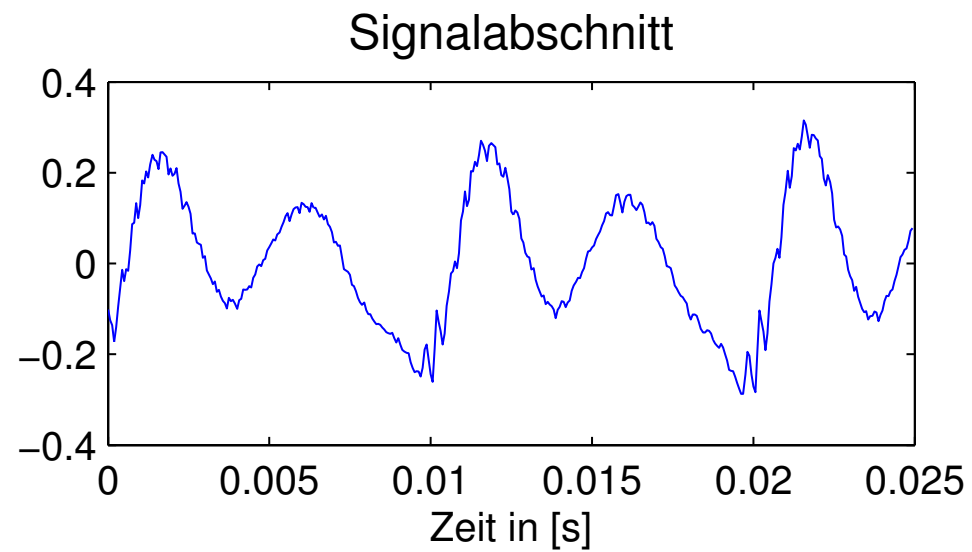
Merkmale für die Spracherkennung: MFCC

(MFCCs: mel frequency cepstral coefficients)

Ermittlung des Mel-Cepstrums:



Analyseabschnitt



Merkmale für die Spracherkennung: MFCC

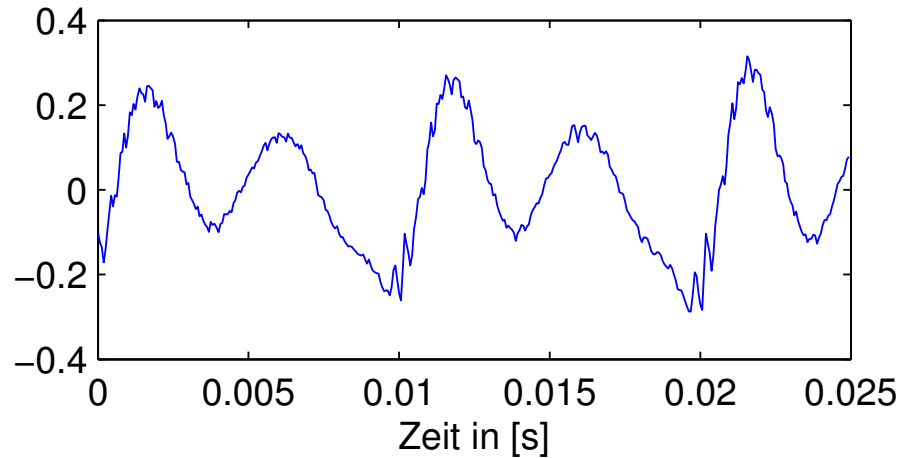
(MFCCs: mel frequency cepstral coefficients)

Ermittlung des Mel-Cepstrums:

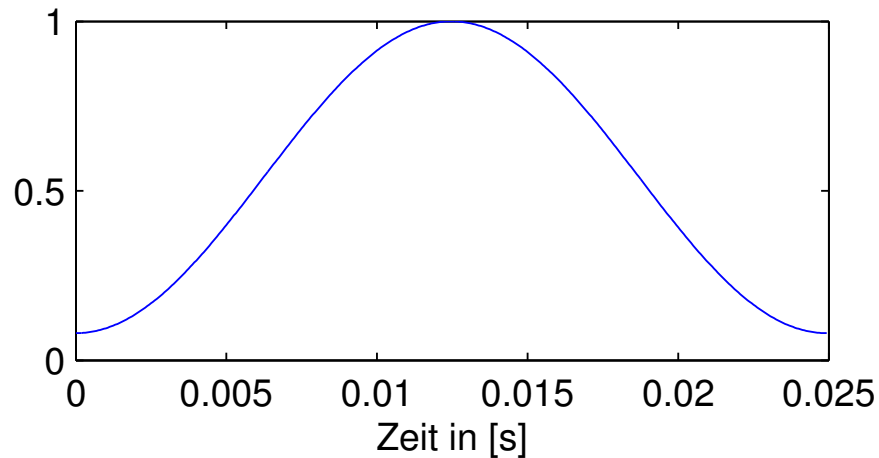


Multiplikation mit Fensterfunktion

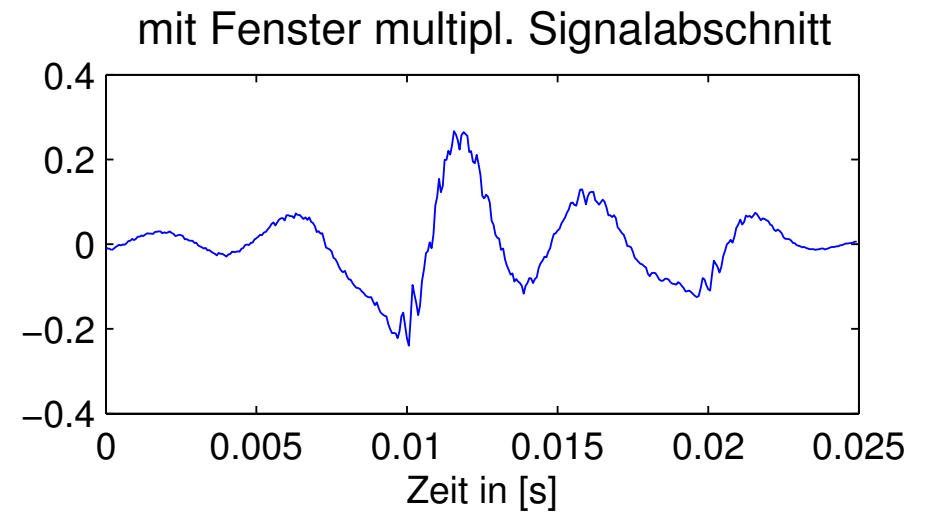
Signalabschnitt



Hamming-Fenster



Signalabschnitt · Fensterfunktion



Merkmals für die Spracherkennung: MFCC

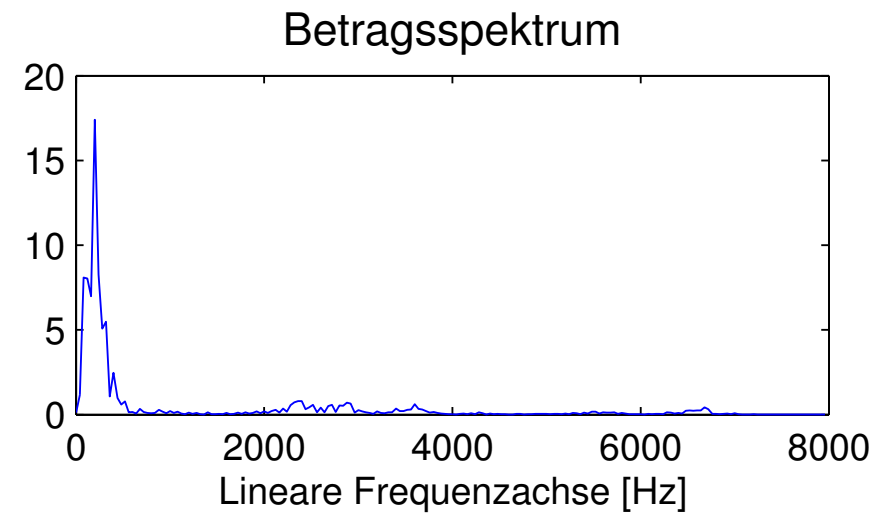
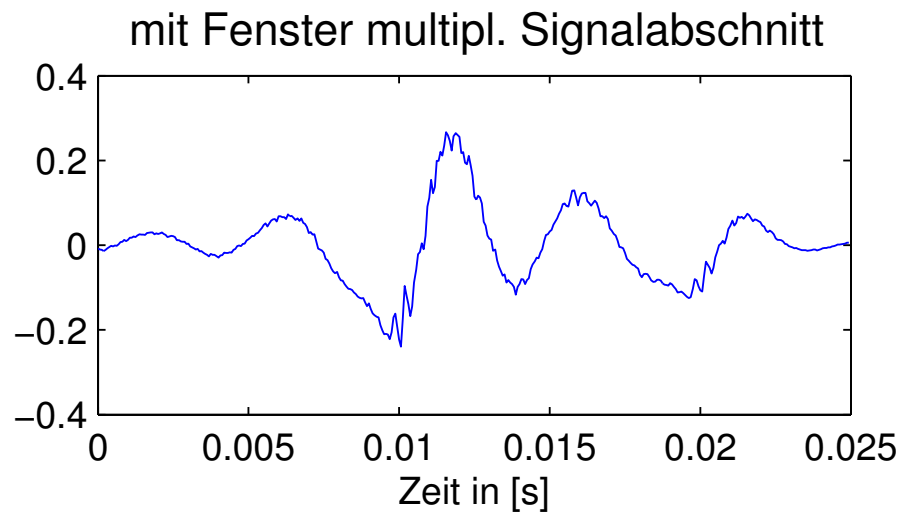
(MFCCs: mel frequency cepstral coefficients)

Ermittlung des Mel-Cepstrums:



Betragsspektrum

DFT → Betrag



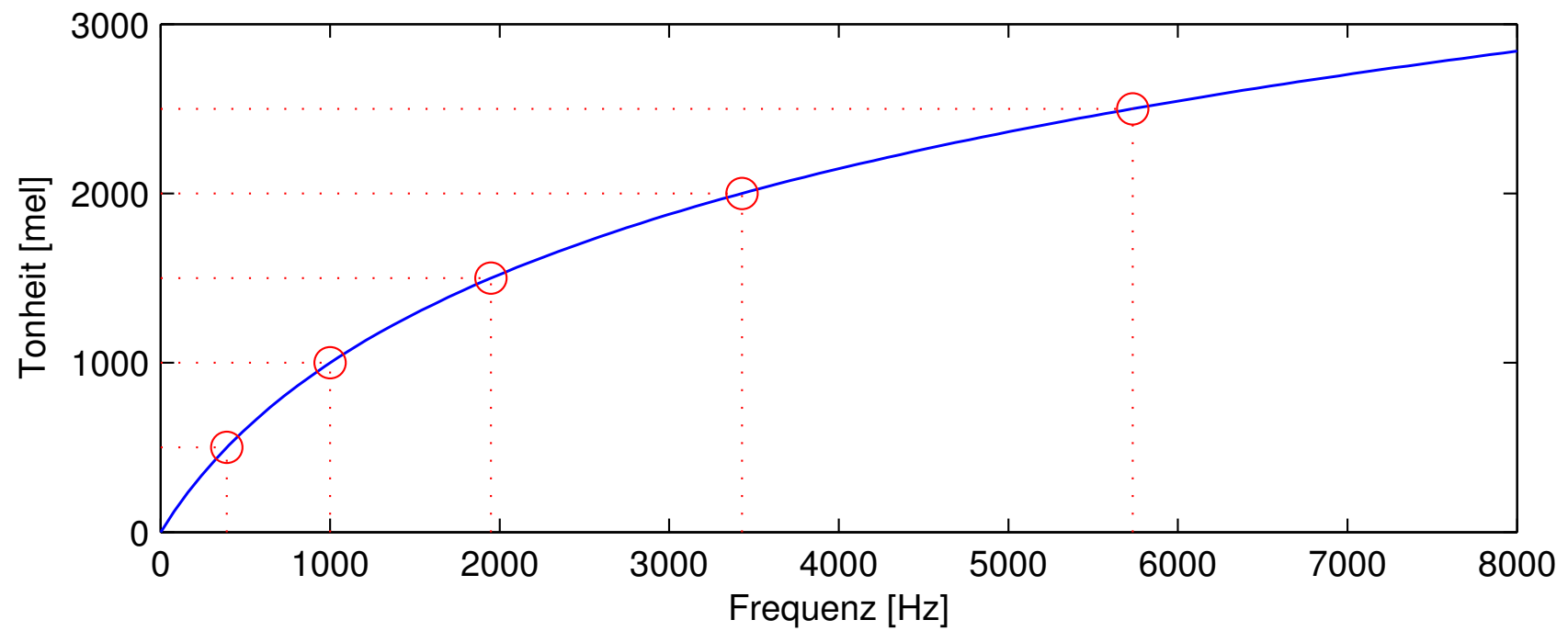
Merkmale für die Spracherkennung: MFCC

(MFCCs: mel frequency cepstral coefficients)

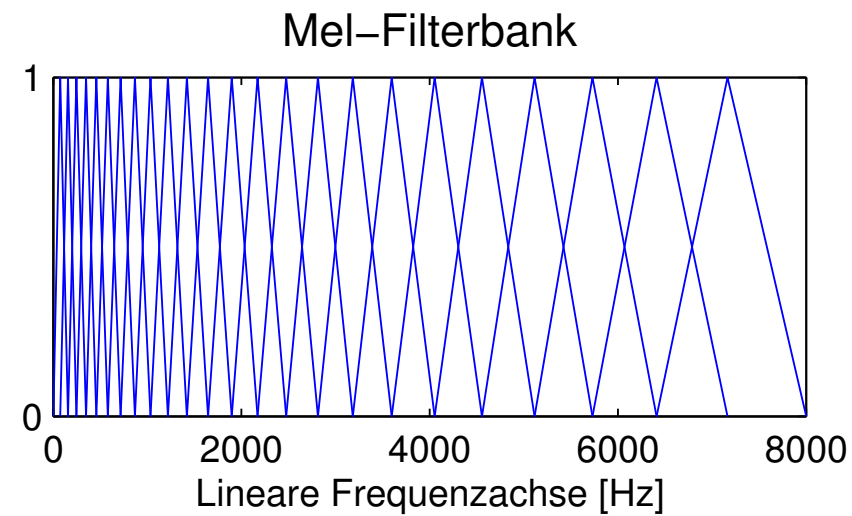
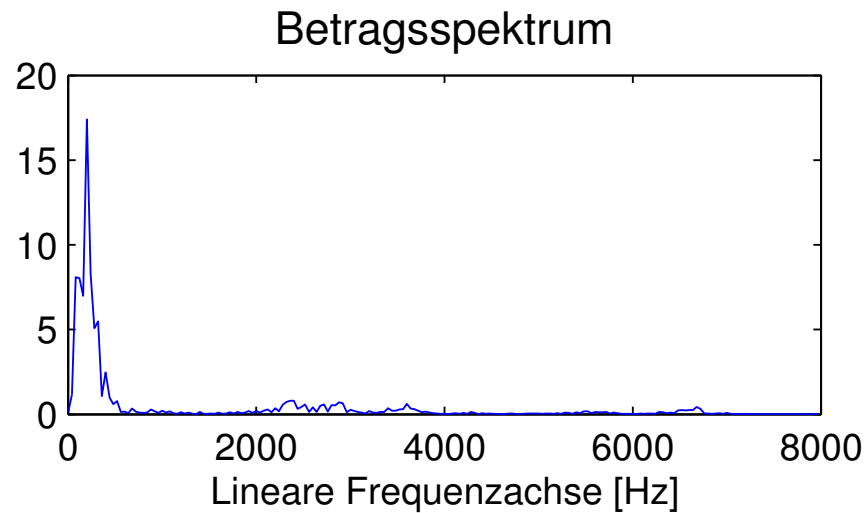
Ermittlung des Mel-Cepstrums:



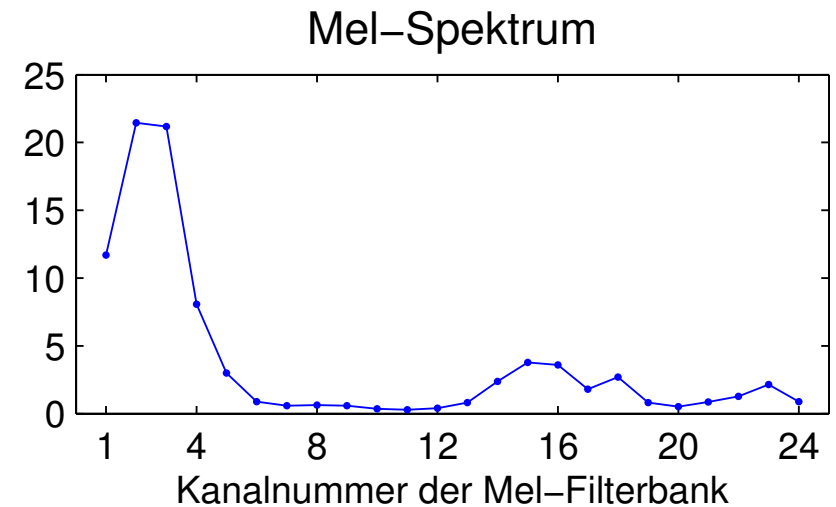
Mel-Skala



Mel-Spektrum



Betragsspektrum · Mel-Filterbank



Merkmale für die Spracherkennung: MFCC

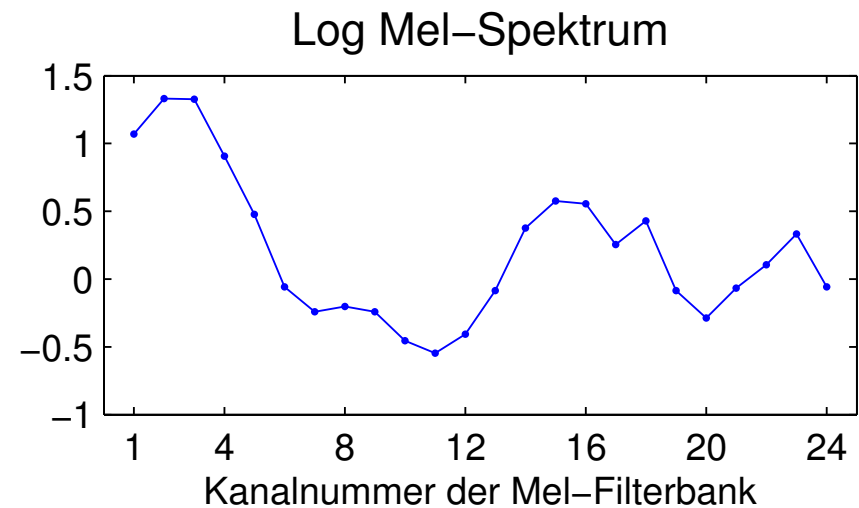
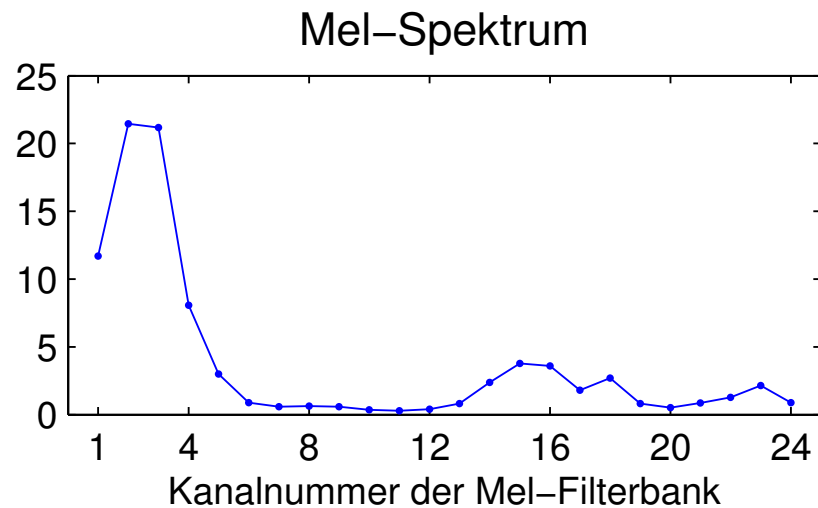
(MFCCs: mel frequency cepstral coefficients)

Ermittlung des Mel-Cepstrums:



Log-Mel-Spektrum

Mel-Spektrum \rightarrow Log



Merkmals für die Spracherkennung: MFCC

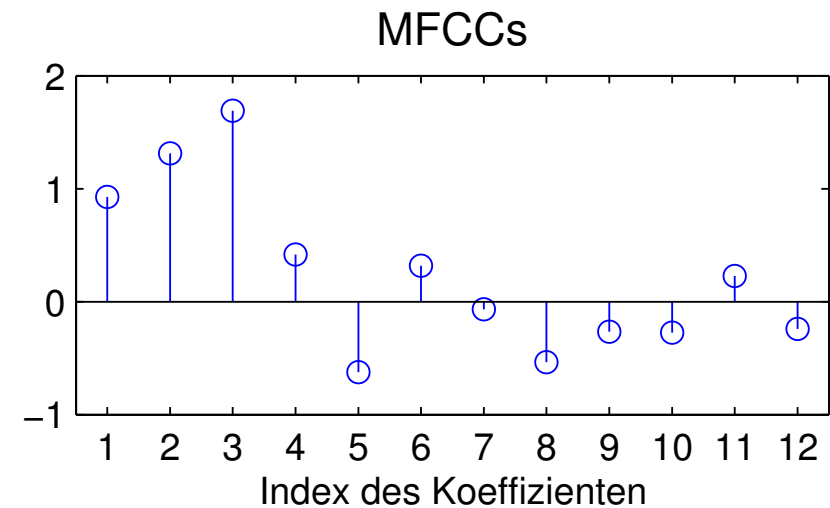
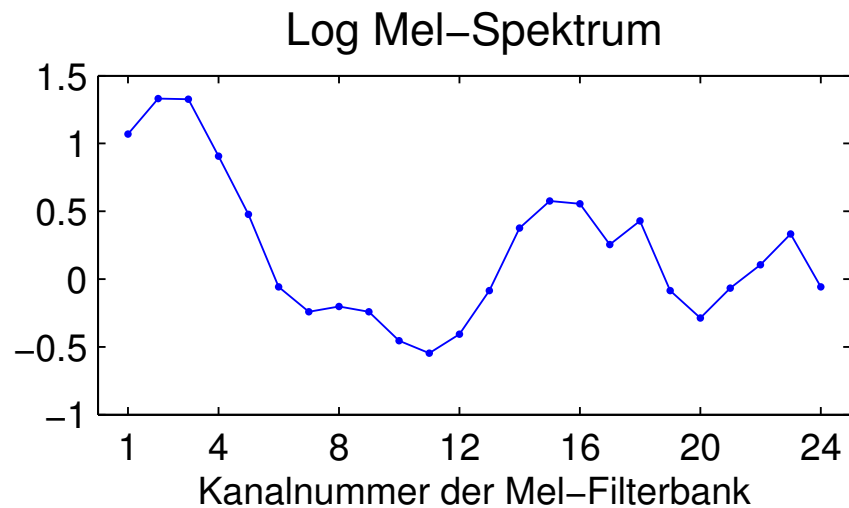
(MFCCs: mel frequency cepstral coefficients)

Ermittlung des Mel-Cepstrums:



Mel-Cepstrum

Log-Mel-Spektrum \rightarrow DCT



Merkmale für die Spracherkennung: MFCC

(MFCCs: mel frequency cepstral coefficients)

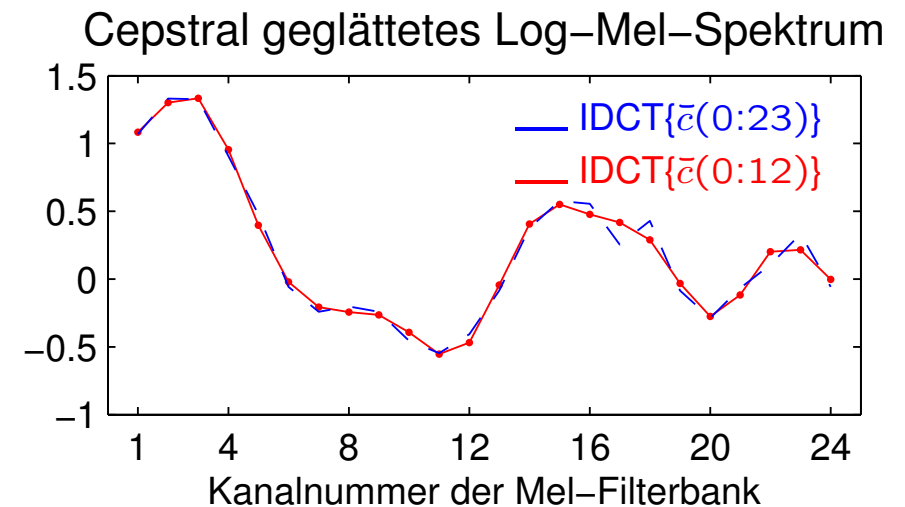
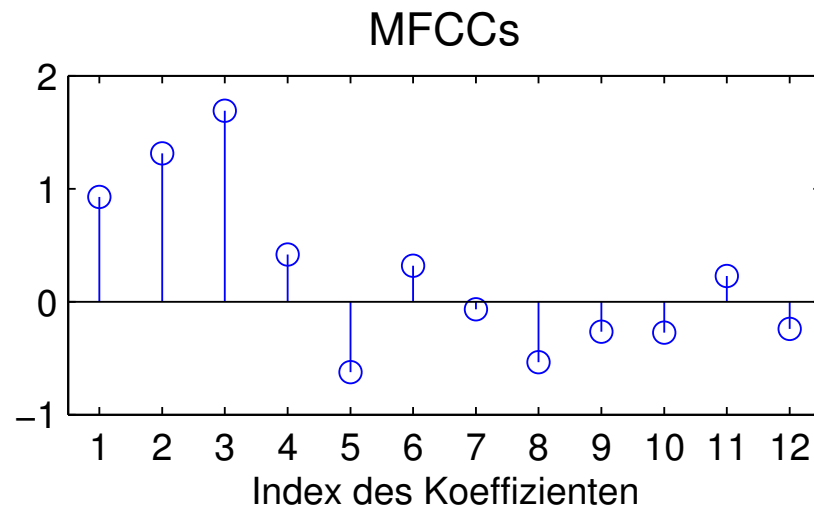
Ermittlung des Mel-Cepstrums:



Was beschreiben die MFCCs?

IDCT(MFCCs)

>>>

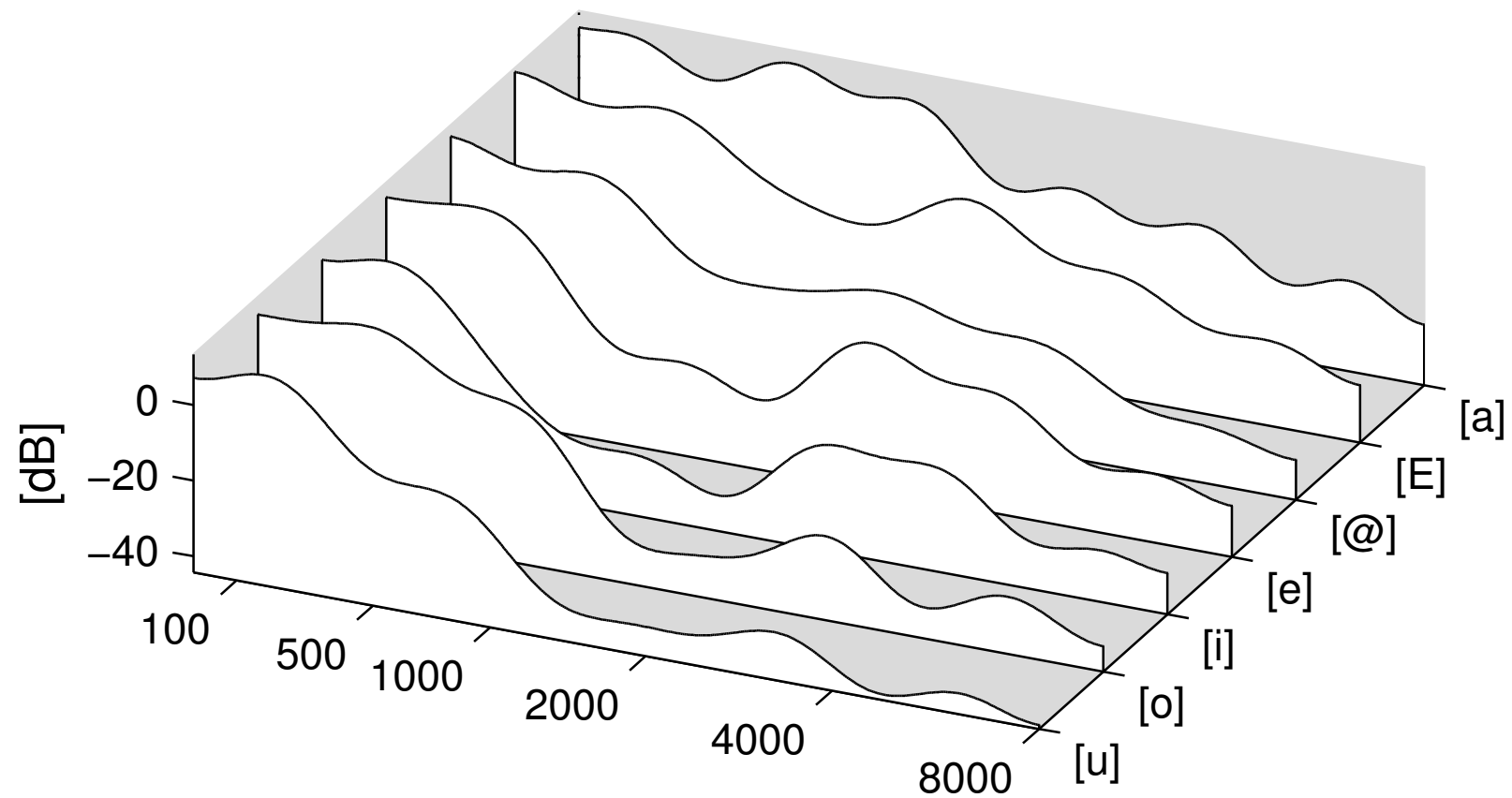


Der Merkmalsvektor repräsentiert das cepstral geglättete Log-Mel-Spektrum

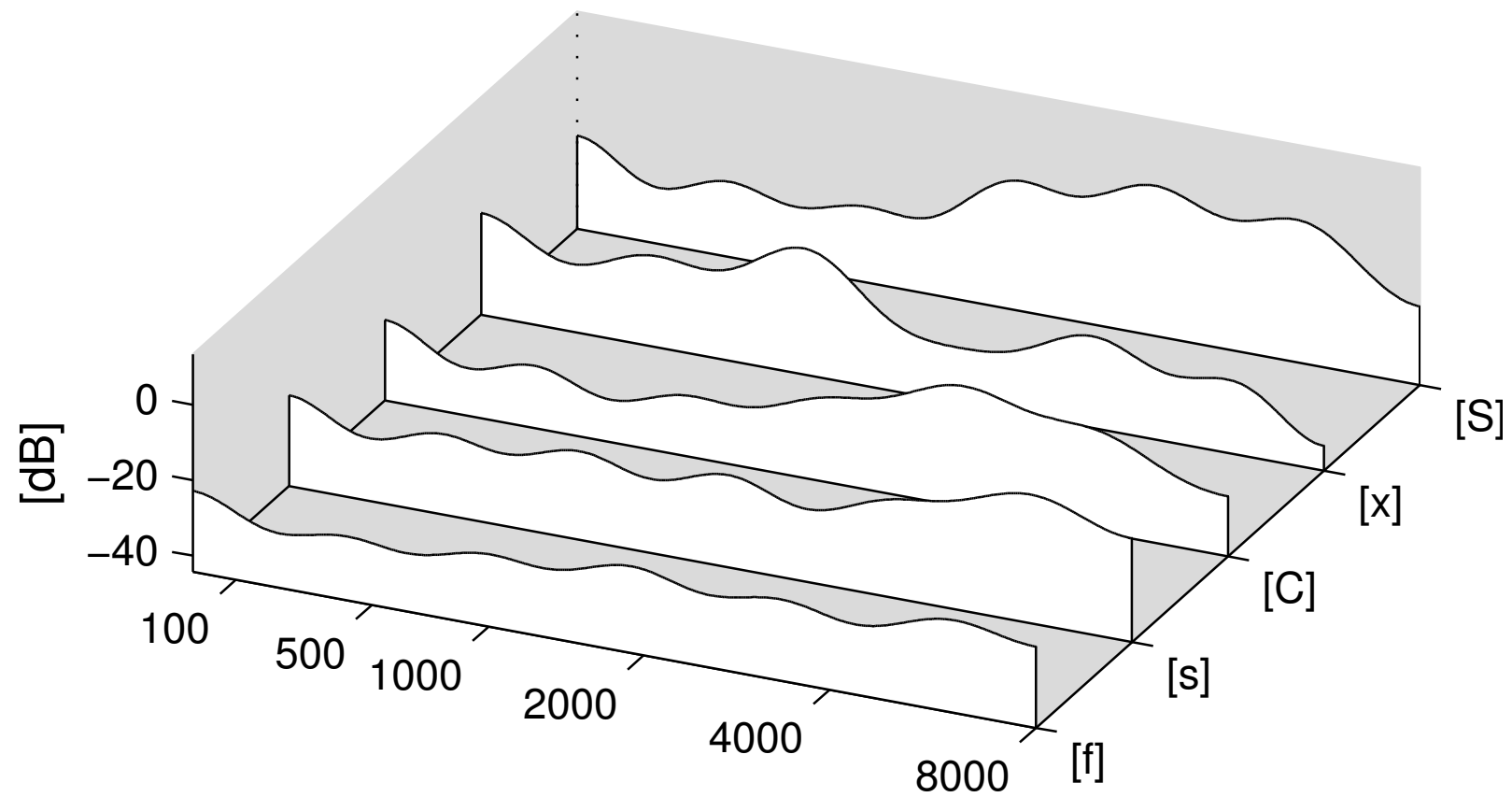
Wahl der Parameter

Abtastfrequenz	8/16 kHz
Länge des Analysefensters	25 ms
Verschiebung des Analysefensters	10 ms
Anzahl Dreiecksfilter	24
Cepstrale Koeffizienten	$\bar{c}(0) - \bar{c}(12)$

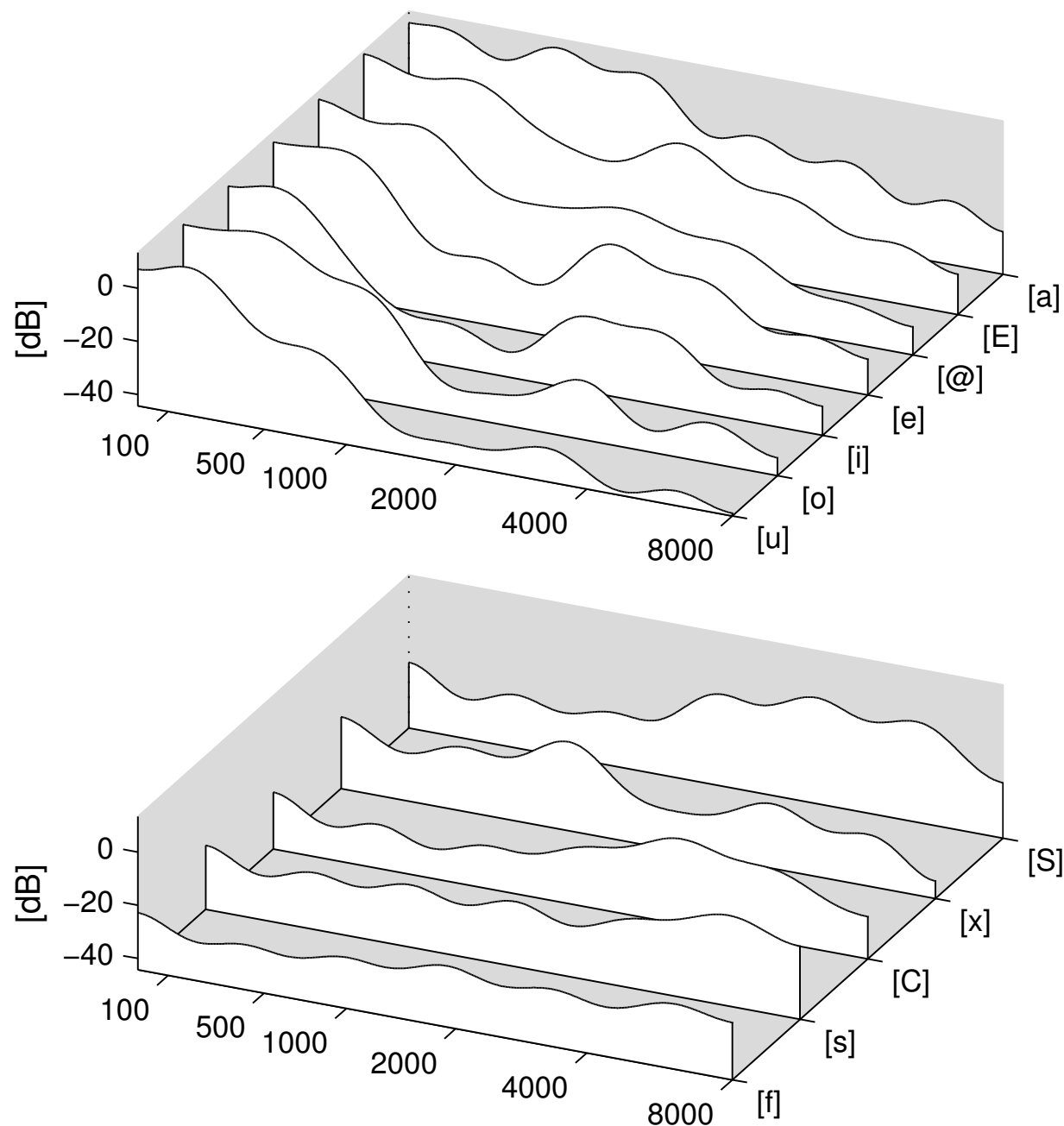
Vokale



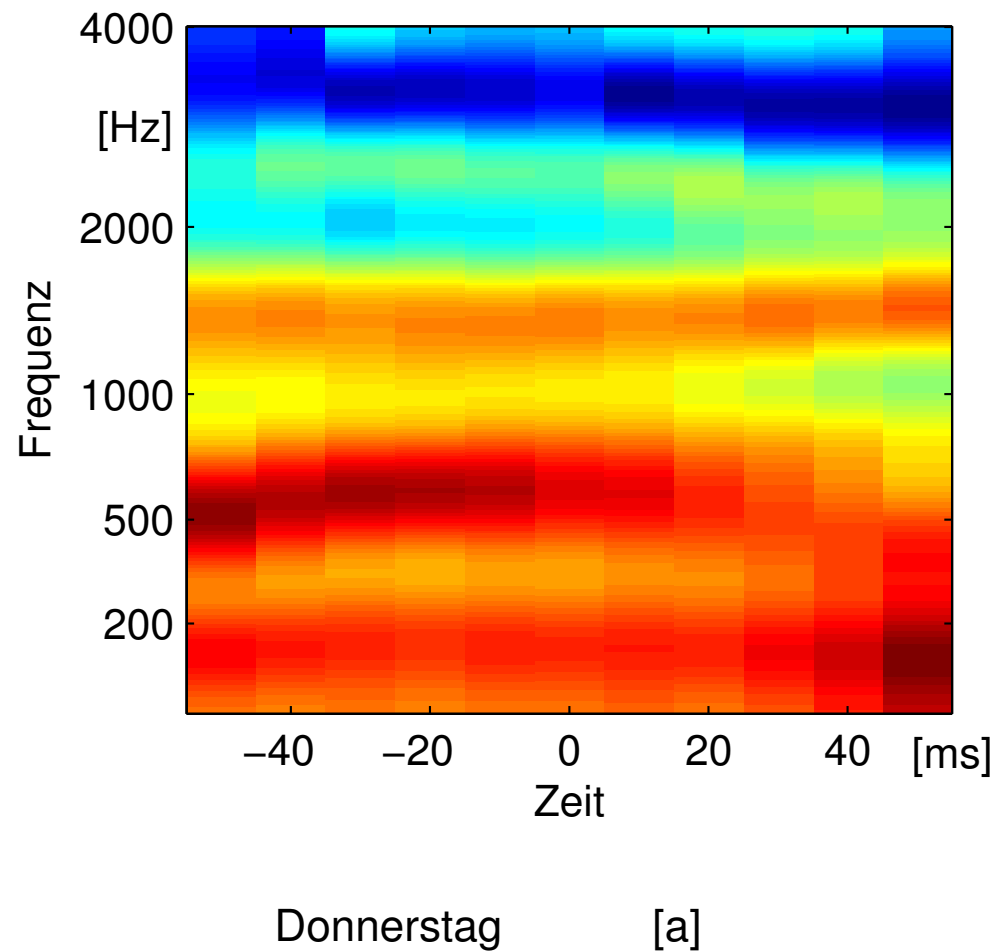
Frikative



Vokale und Frikative im Vergleich

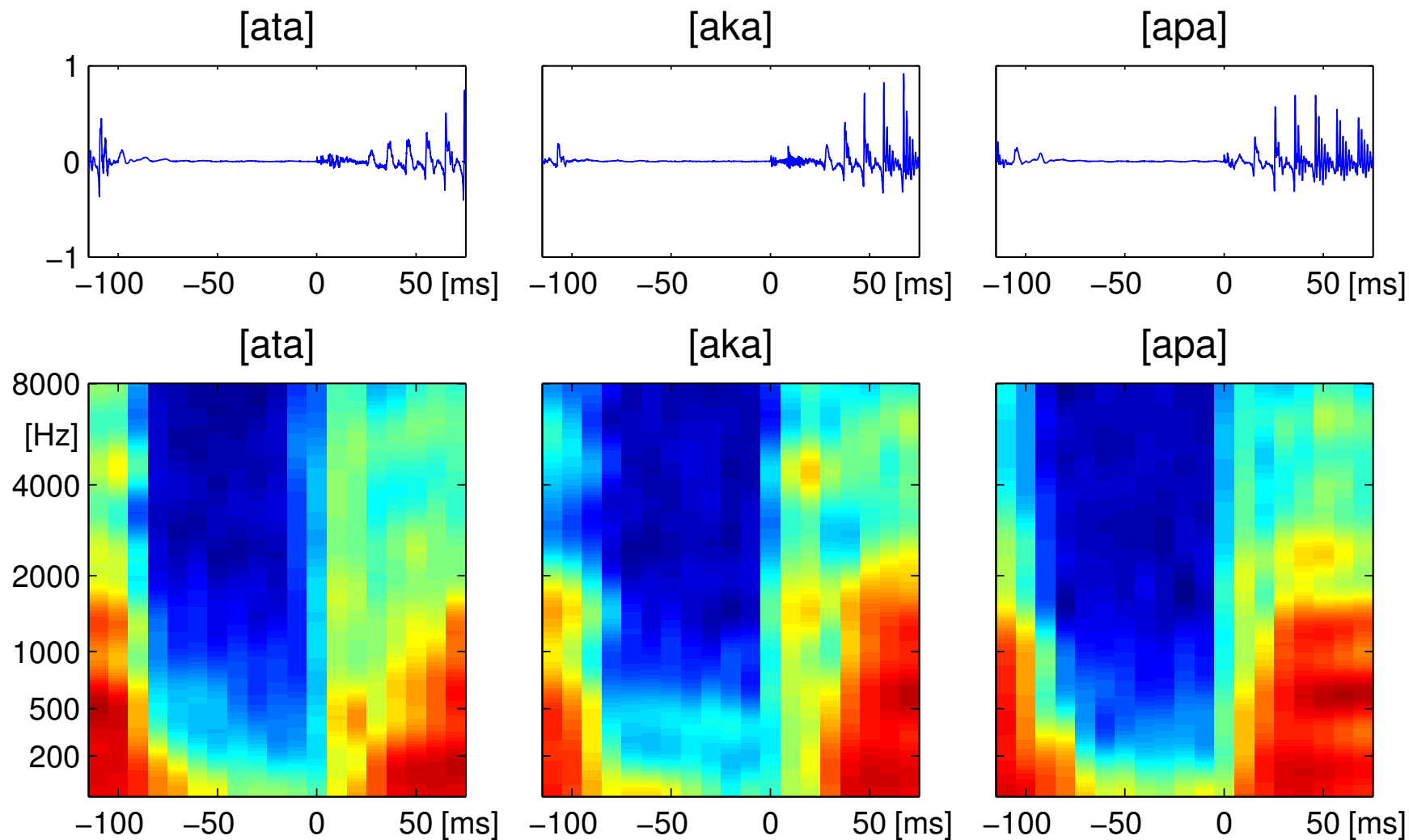


Zeitlicher Verlauf des Spektrums



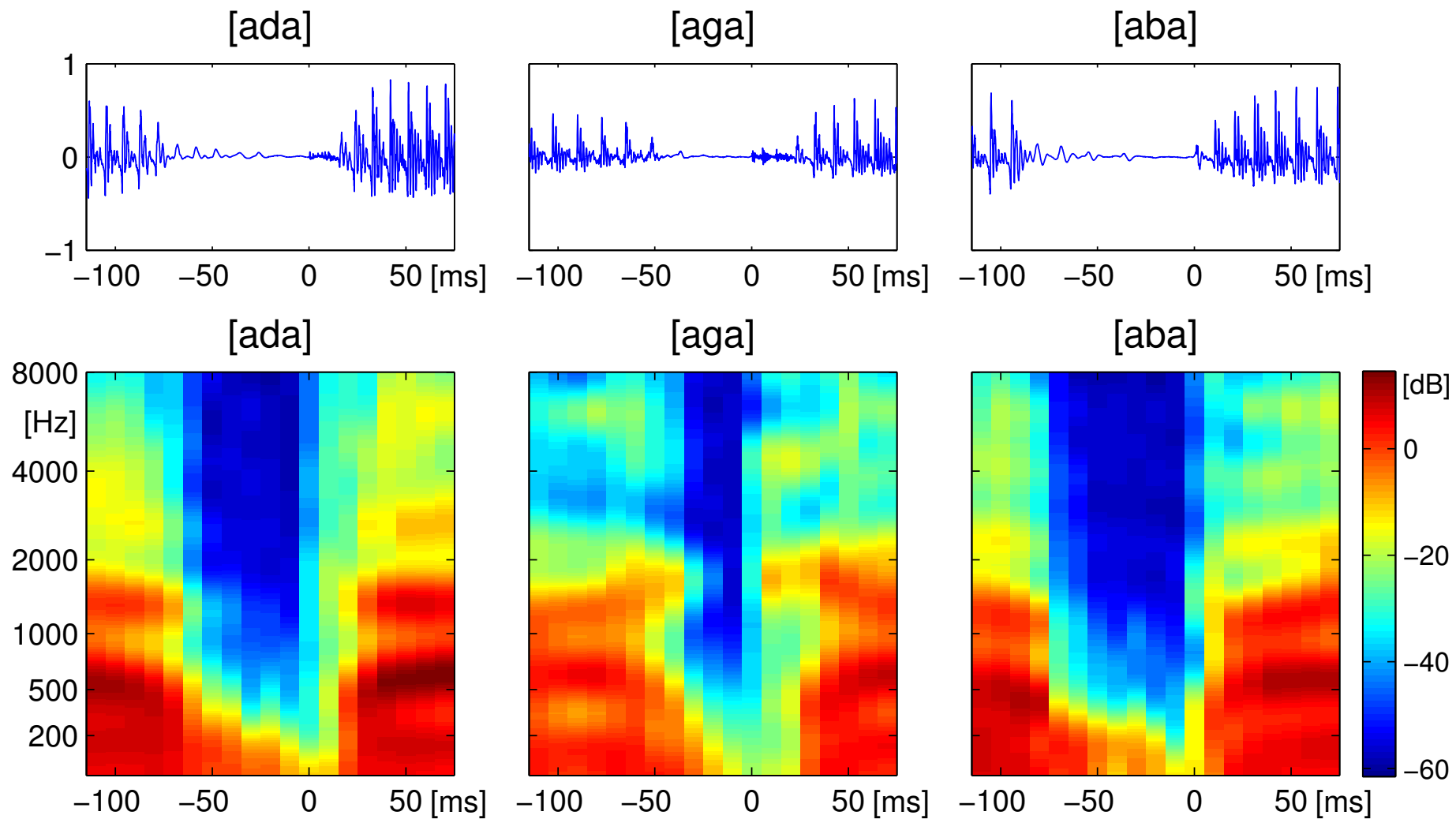
Stimmlose Plosive

Verschlusslaute [t] [k] [p]

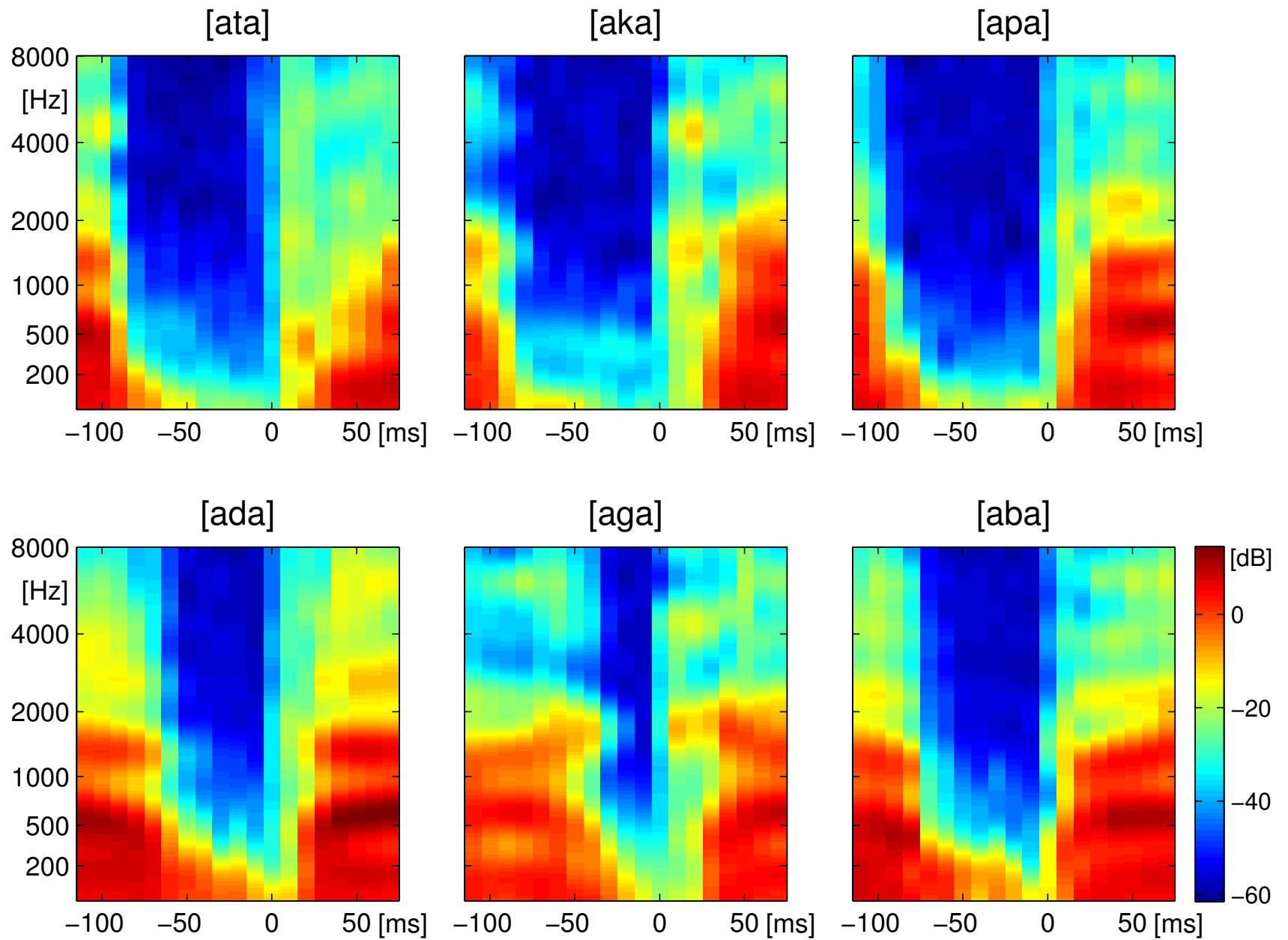


Stimmhafte Plosive

Verschlusslaute [d] [g] [b]



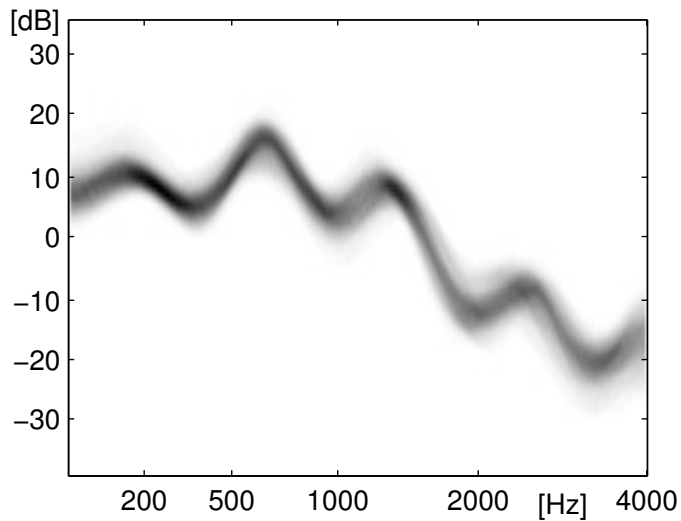
Plosive



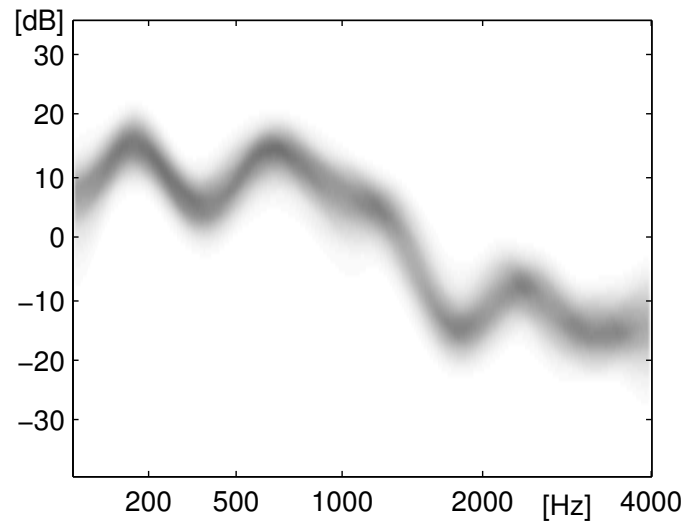
Spektrale Variabilität

Laut [a:]

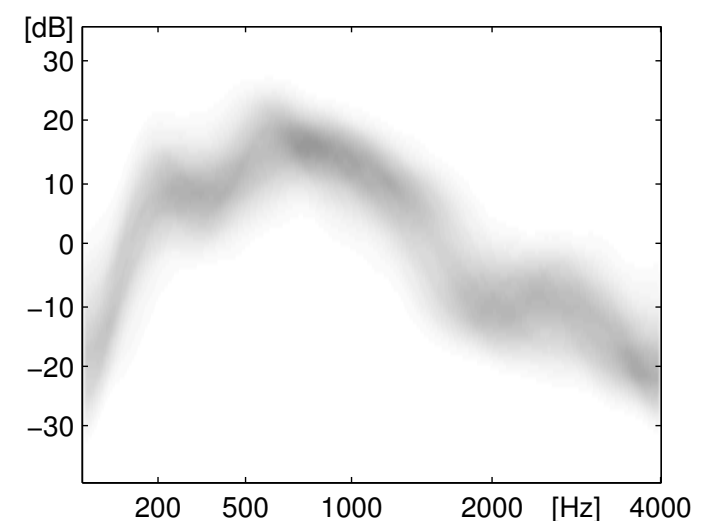
professioneller Sprecher



normaler Sprecher

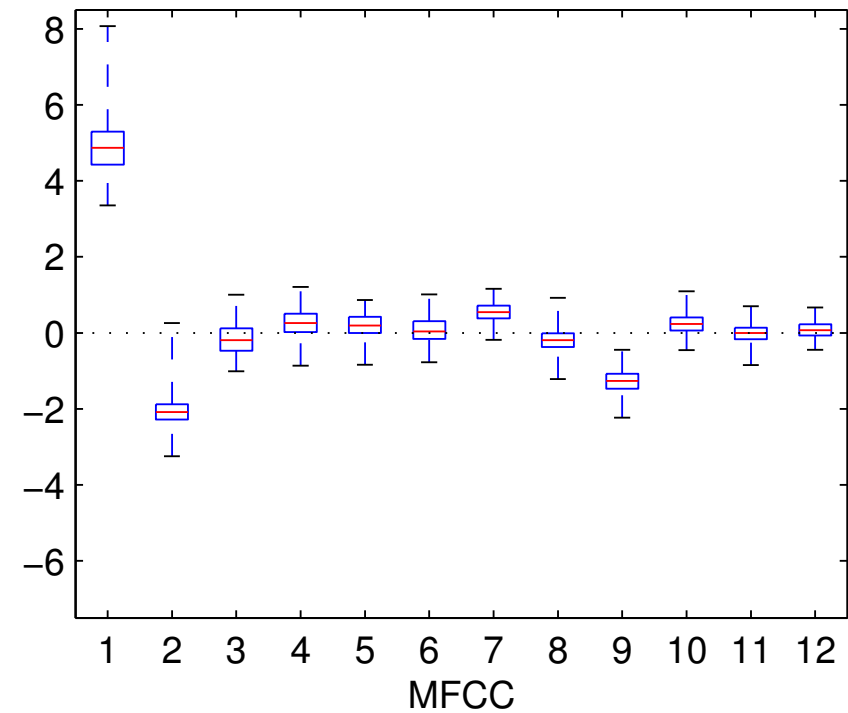
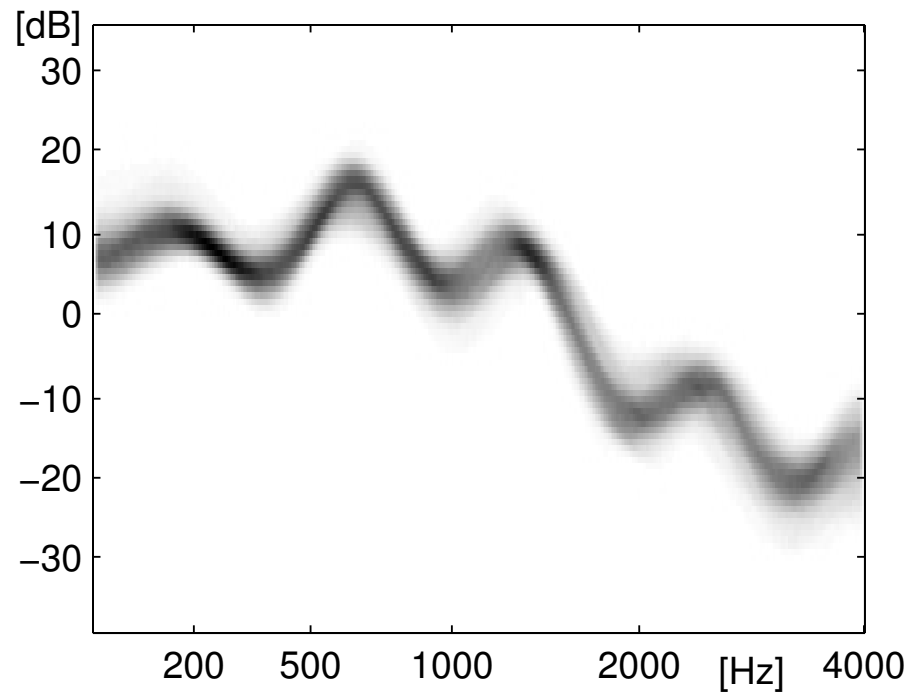


1000 Sprecher



Spektrum vs. Cepstrum

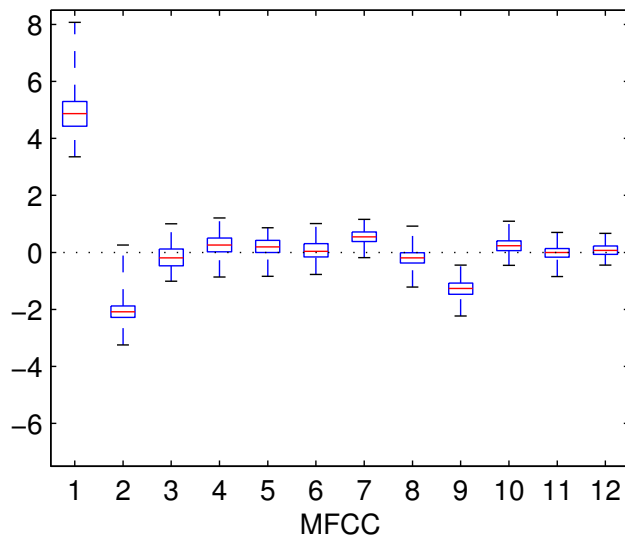
professioneller Sprecher



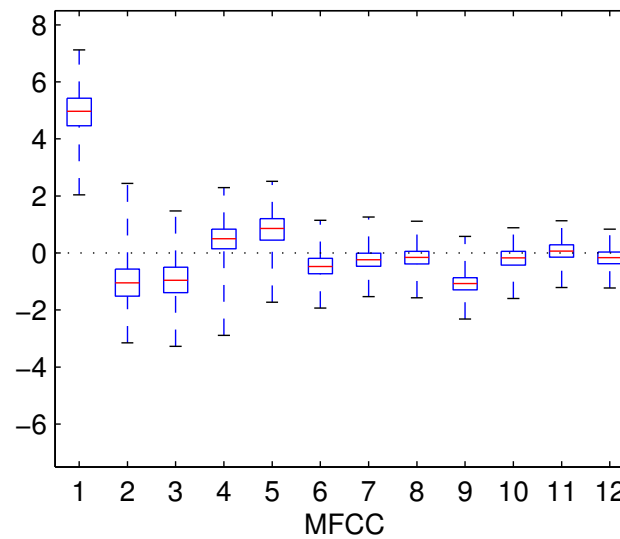
Cepstrale Variabilität

Laut [a:]

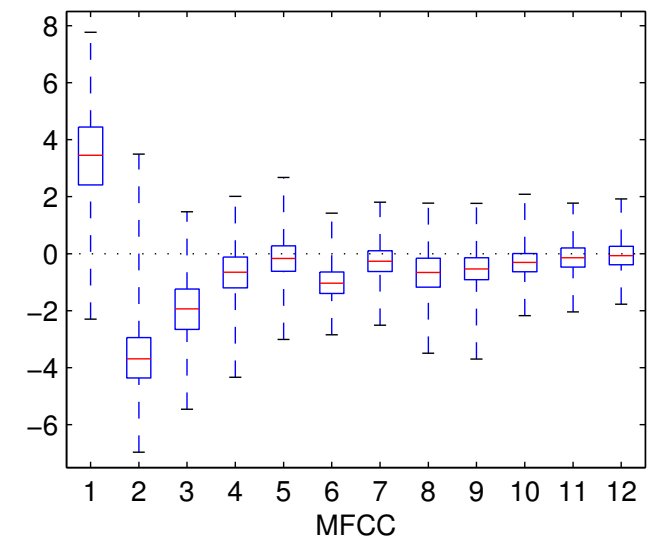
professioneller Sprecher



normaler Sprecher



1000 Sprecher



Welche Eigenschaften des Sprachsignals
sind in den MFCC **nicht** enthalten?

Welche Eigenschaften des Sprachsignals sind nicht in den MFCC enthalten?

- die Phase
- die Feinstruktur des Spektrums
- die Grundfrequenz
- die Periodizität (stimmhaft/stimmlos)

Rekonstruktion des Signals aus den MFCC

Rekonstruktion



Originalsignal

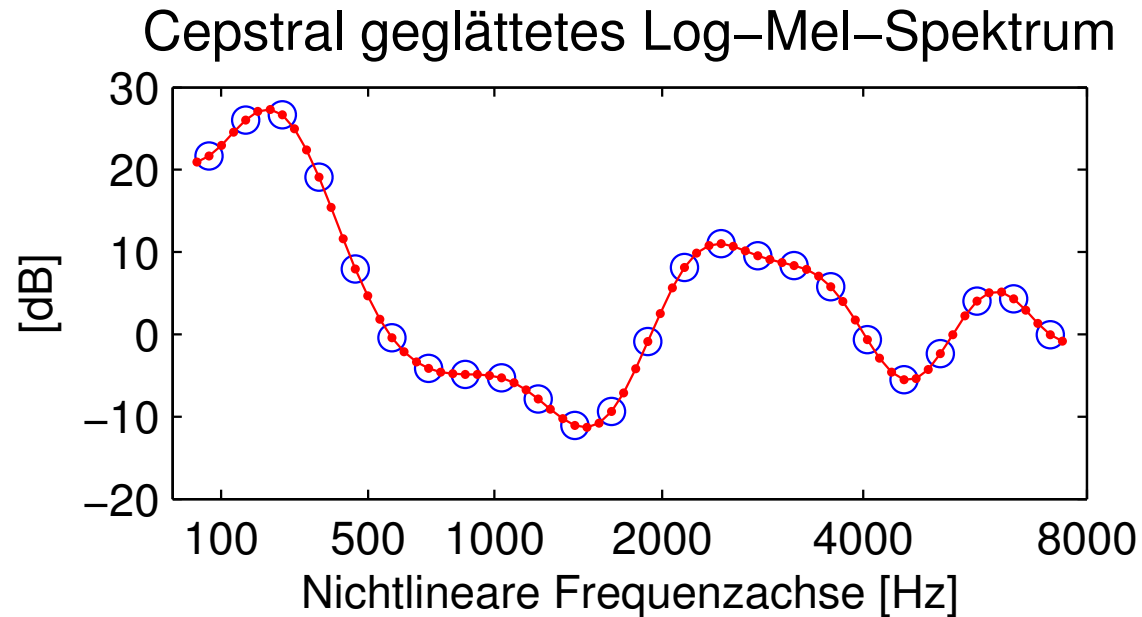


Thema der nächsten Lektion:

Spracherkennung mittels Mustervergleich

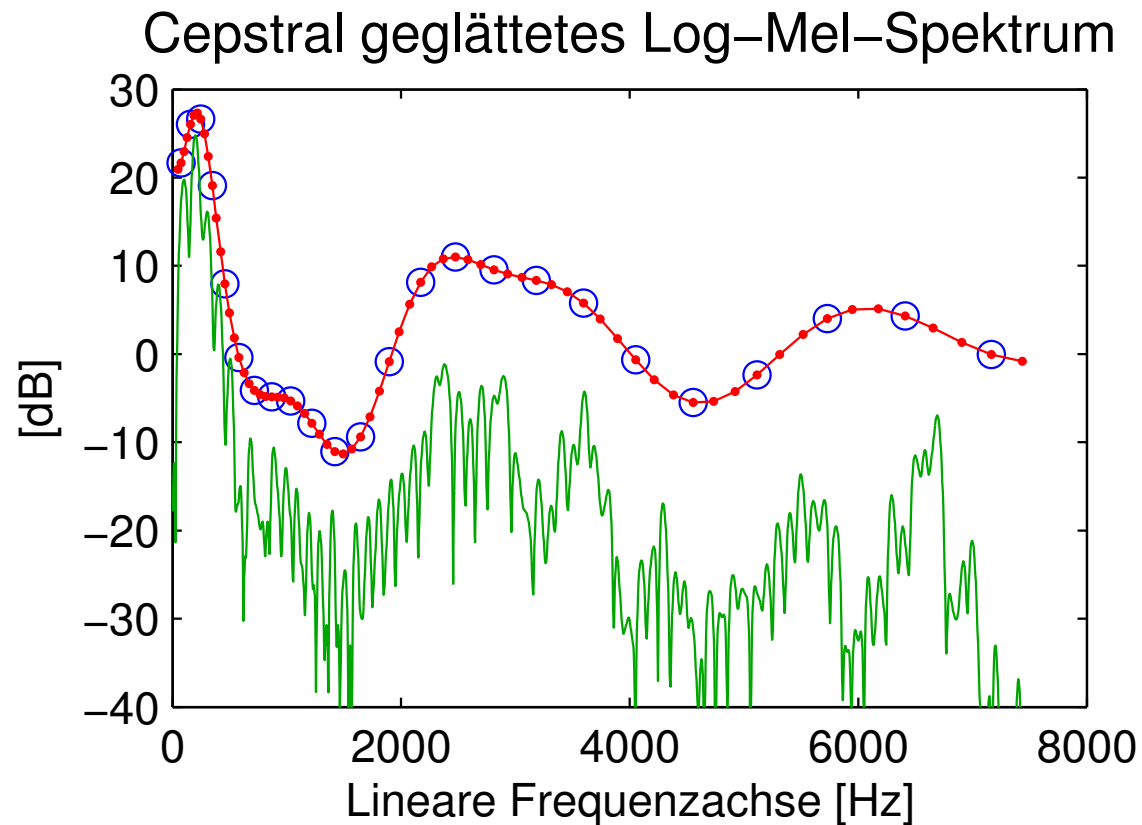
Zur Übersicht der Vorlesung *Sprachverarbeitung I* >>>

Rücktransformation: MFCC \rightarrow Log-Mel-Spektrum



- Skala der Achse in Hertz (nicht linear)
- Interpolation (IDCT mit Zero-Padding)
- Vergleich mit Betragsspektrum des Signals?

Vergleich mit Betragsspektrum des Signals



Frage: Woher rührt der Unterschied zwischen den Spektren?

<<<

Mel-Filterbank

