

**Erstellung einer  
Sprachdatenbank**

**sowie**

**eines Programms zu deren  
Analyse im Kontext einer  
Sprachsynthese  
mit spektralen Modellen**

Tobias Platen

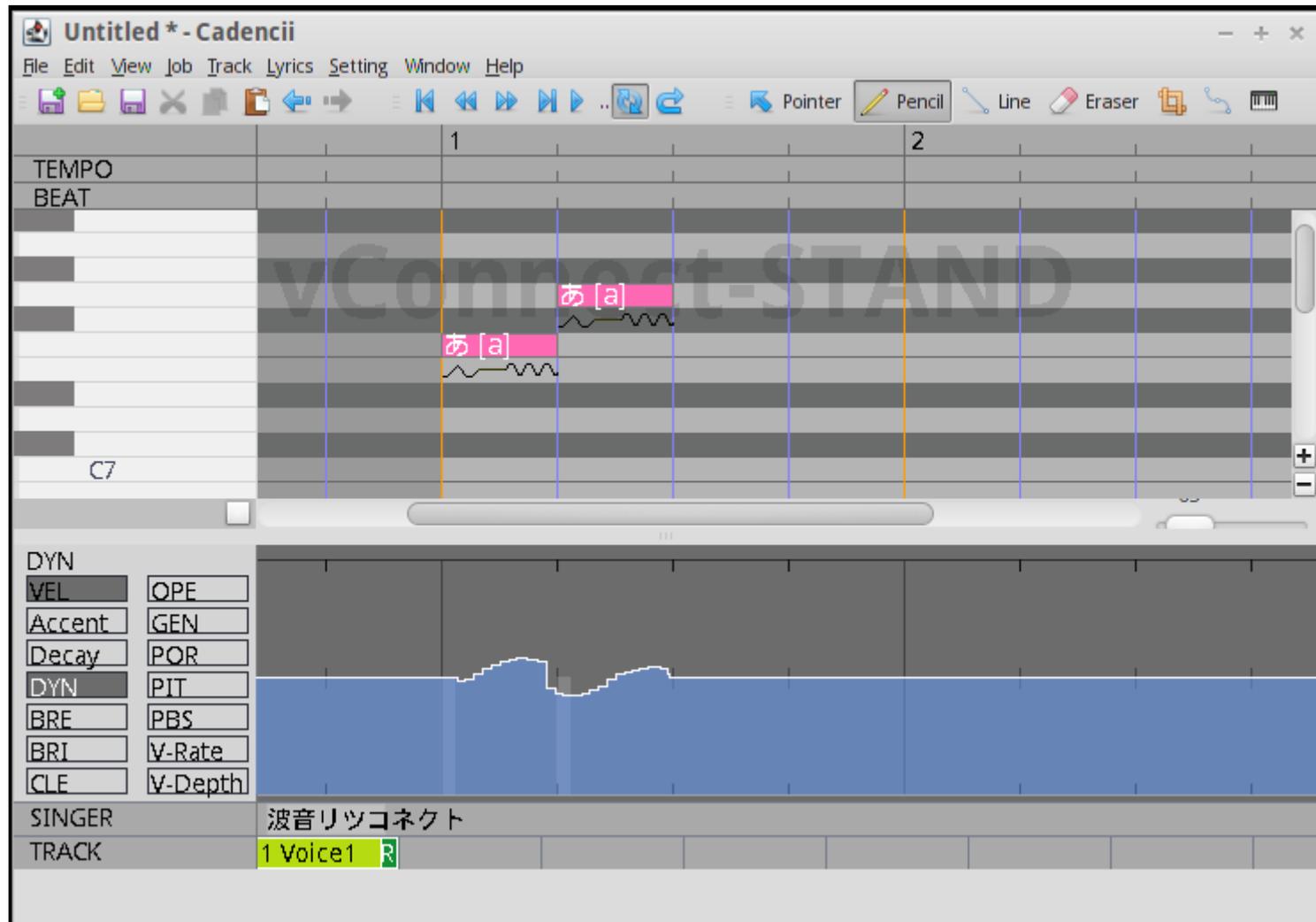
# Gliederung

1. Einleitung
2. Untersuchung spektraler Modelle im Hinblick auf Eignung für Sprachsynthese
3. Erstellung der Sprachdatenbank
4. Ergebnisse
5. Ausblick

# 1. Einleitung

- **Motivation**
- **Ziele**
- **Wege und Methoden zum Ziel**

# Cadencii



# Spectral Model Synthesis

Funktioniert durch H/S Zerlegung von Signalen

Vorteile des Modells:

- ermöglicht Transformationen der Klangfarbe
- effiziente Resynthese mit IFFT-Methode

Nachteile:

- verwirft Phaseninformation
- Referenzimplementierung fehlerhaft

# Das WORLD Speech Toolkit

Funktioniert durch Zerlegung in drei Teile:

- DIO: fundamentale Frequenz
- STAR: spektrale Hüllkurve
- PLATINUM: Anregungssignal

Vorteile:

- Sehr gute Sprachqualität
- Freie Software unter der BSD-Lizenz

# Verbesserung von eSpeak durch WORLD

- Integration in den speech-dispatcher
- Generierung der Prosody mit eSpeak
- Schreiben des Parsers für MBROLA pho-files
- Entnahme von Samples aus einer Datenbank
- Implementierung der Echtzeitsynthese
- Test der konkatenativen Sprachausgabe

# Konkrete Schritte zur Erstellung der Sprachdatenbank

- Corpus-Design kompatibel zu MBROLA de2
- Rendern der Prompts
- Aufnahme der Sprachsamples
- Beschriften der Samples
- Spektrale Analyse
- Verpacken der Datenbank
- Produktivtest durch Endnutzer

# Die Sekai Speech Tools

selbstgeschriebene grafische Programme:

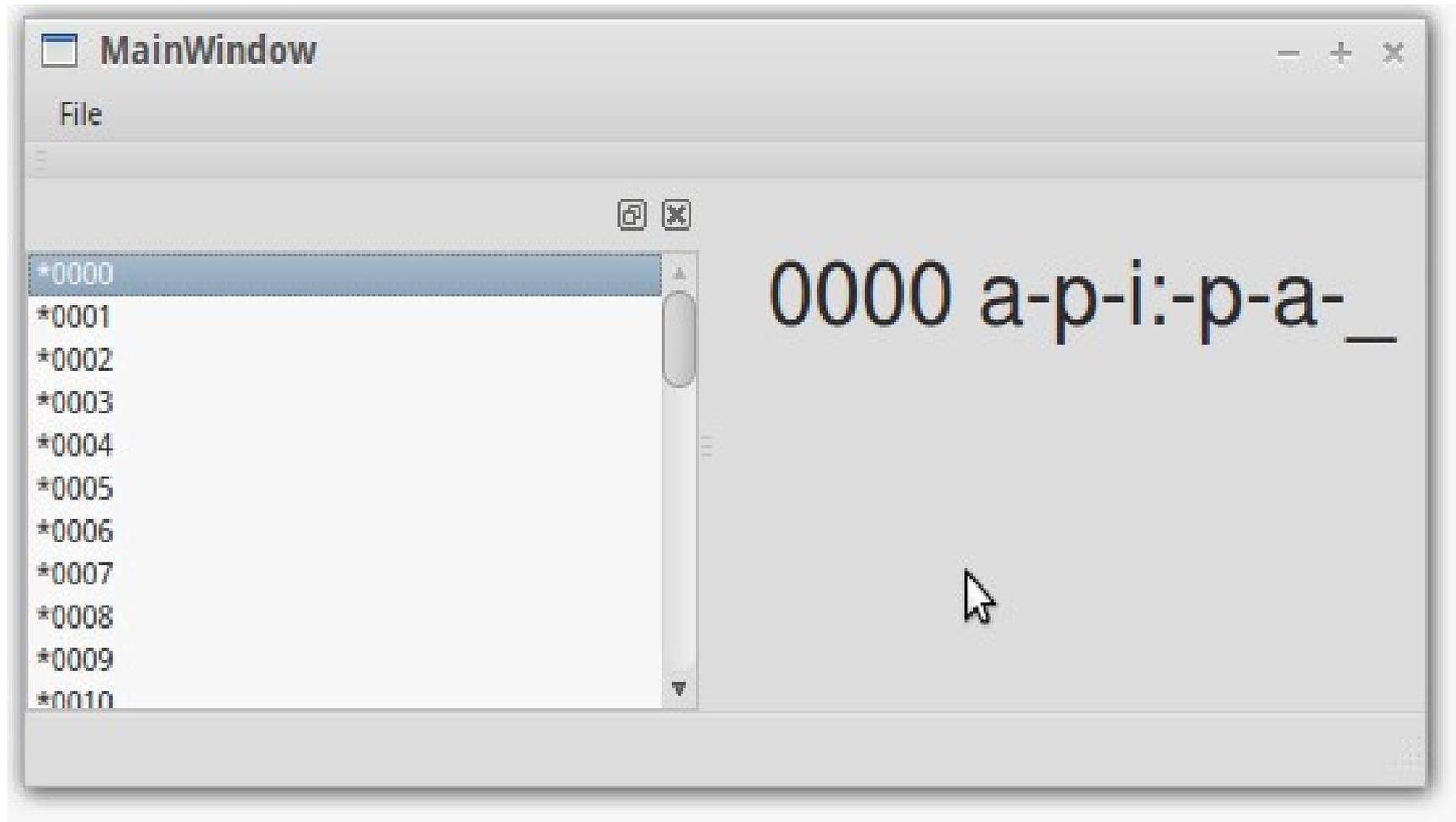
- Oremo2: Echtzeitrecording mit Jack
- OtoEdit: Visuelle Sample-Segmentierung
- SekaiViewer: Visualisierung der Sprachdaten

Batch-processing der einzelnen Samples

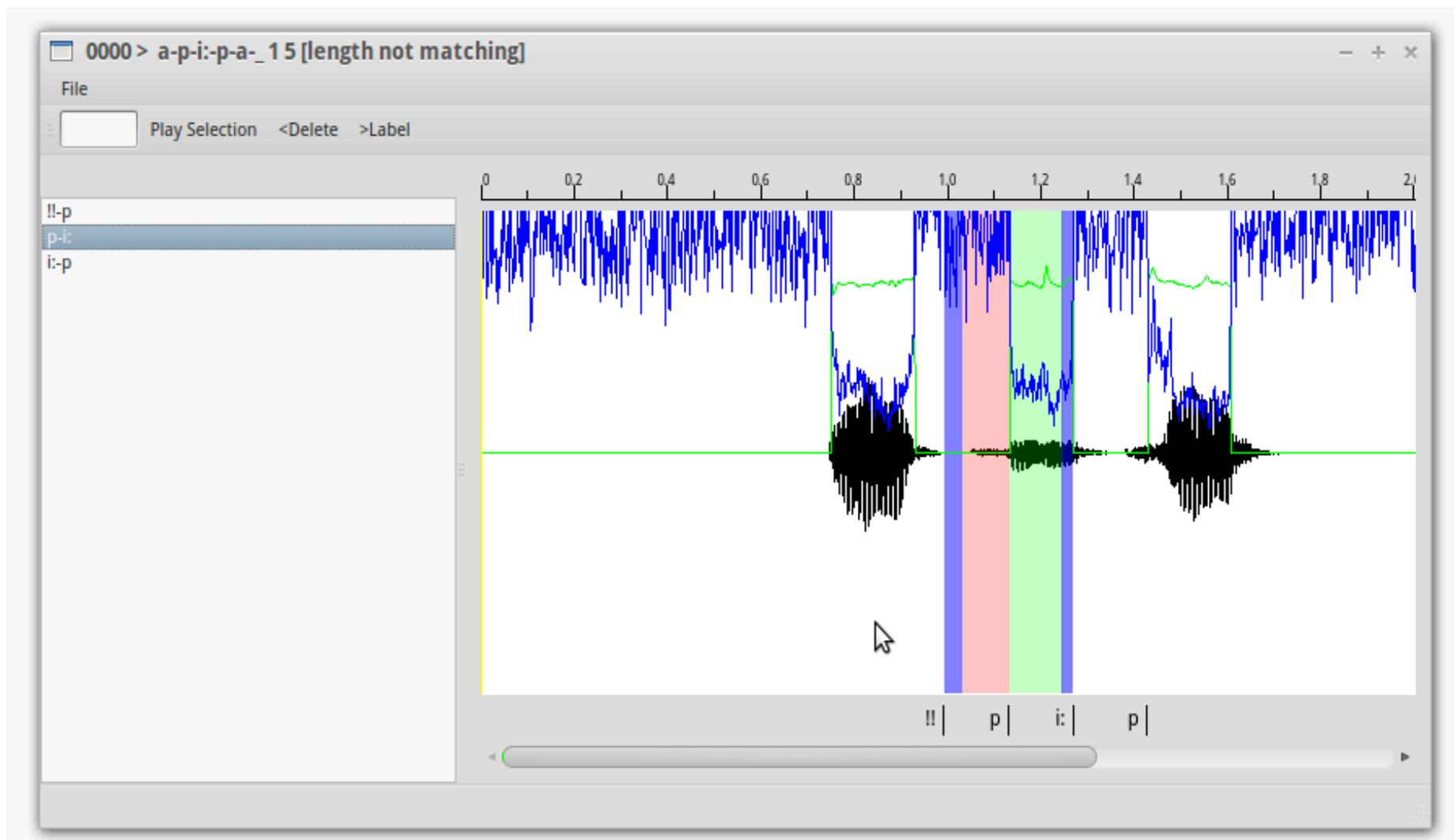
- temporäre Dateien sehr groß
- Vorbis und MFCC zur Kompression

Test der Konkatenation und Kompression

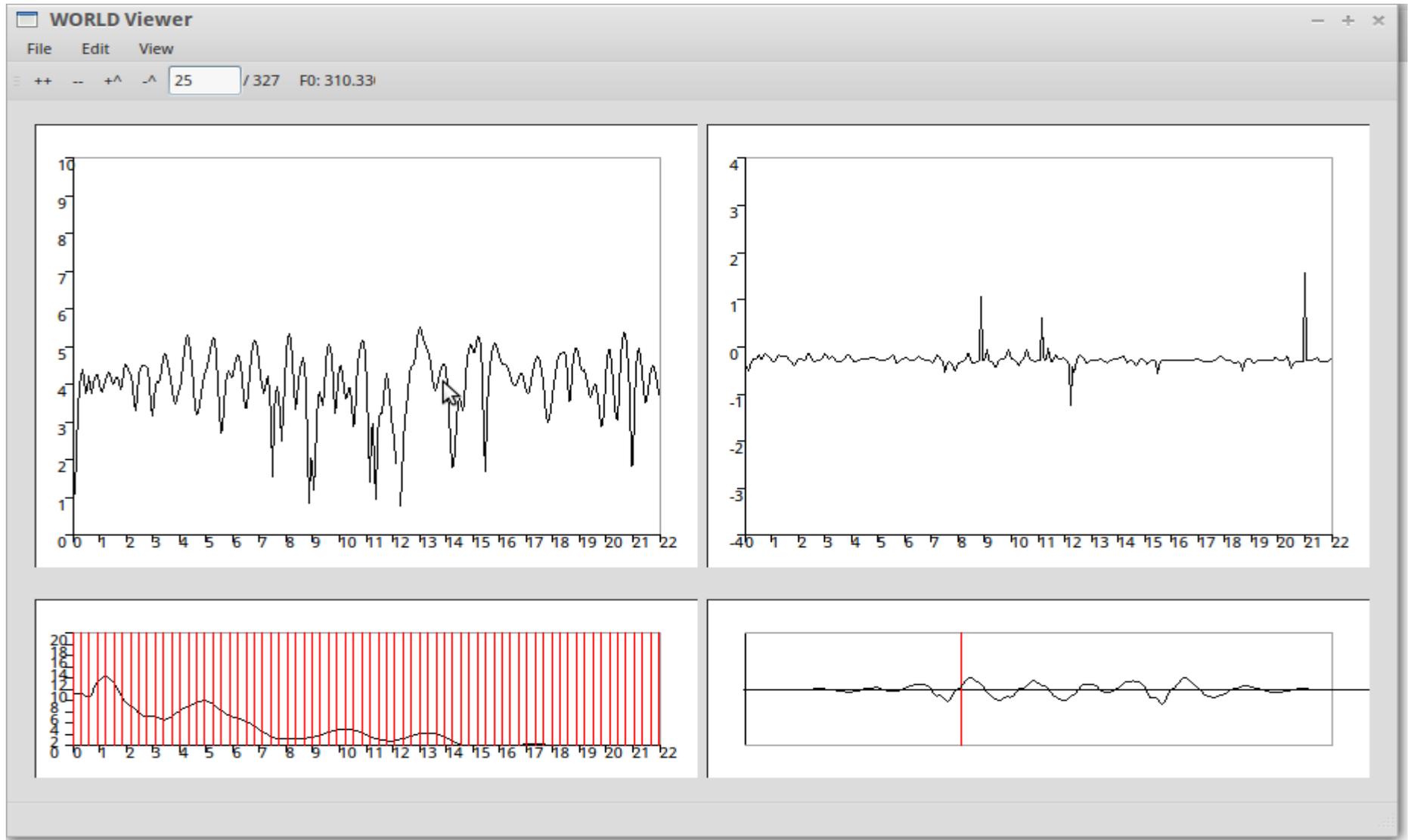
# Oremo2



# OtoEdit



# SekaiViewer



# Ergebnisse

- Vorteile von Vocoder und TD-PSOLA nutzbar
- automatische Bestimmung der GCIs gelöst
- Nutzung von spektralen Stimmenmodellen sinnvoll
- realistische Formantbewegungen jetzt möglich
- Tests zeigen Eignung der Algorithmen
- besserer Klang auch gegenüber Festival

# Ausblick

nötig:

- weitere Tests der Sprachdatenbank

sinnvoll:

- Refactoring von eSpeak
- Implementation von EpR und IFFT-Methode
- Ergänzung weiterer Phoneme
- Spracherkennung mit HMMs und DTW
- Integration der HTS-Sprachausgabe
- Bessere Modellierung des Anregungssignals
- Verbesserung der Echtzeitsynthese

**Fazit**

Vielen Dank für Ihre  
Aufmerksamkeit

Noch Fragen?