

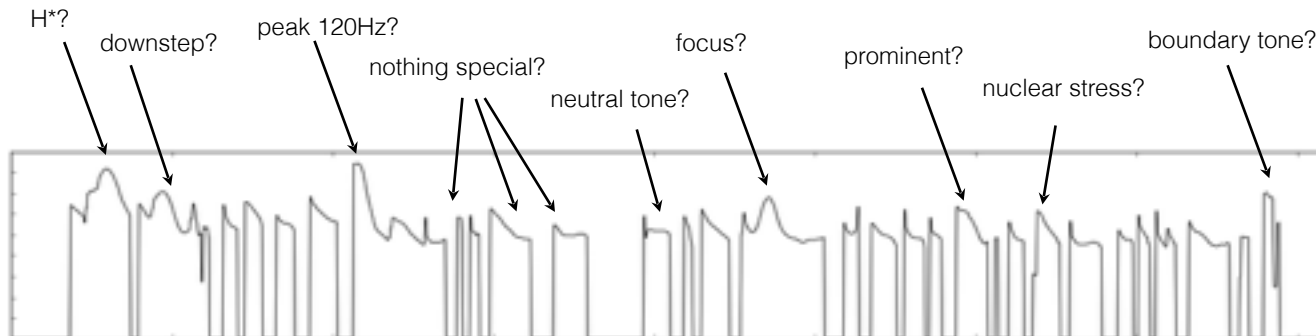


Continuous wavelet transform for speech research

Martti Vainio, Juraj Šimko, Antti Suni
Linguistic Convergence Laboratory Meeting
HSE Moscow, November 27, 2017

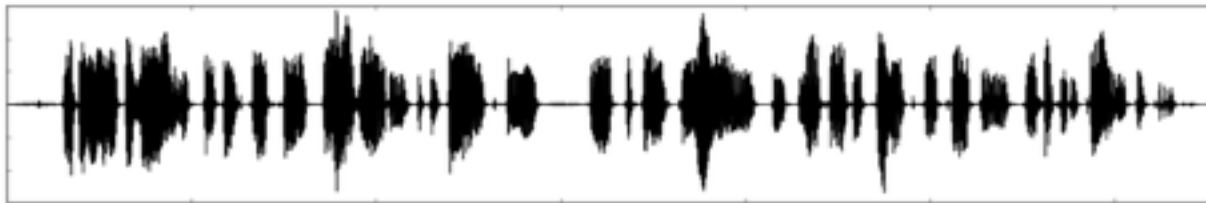
Introduction

- Prosodic signals, like f_0 , are complex, containing information on syllable, word, phrase and utterance levels, with diverse functions.
- The information is encoded *in parallel* in one dimensional signal;
 - Automatic non-trivial prosodic analysis is difficult
 - Expert analysis requires a lot of subjectivity and effort
 - No generally agreed framework for analysis

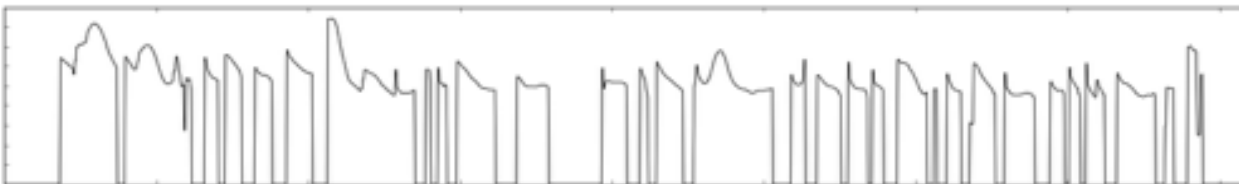
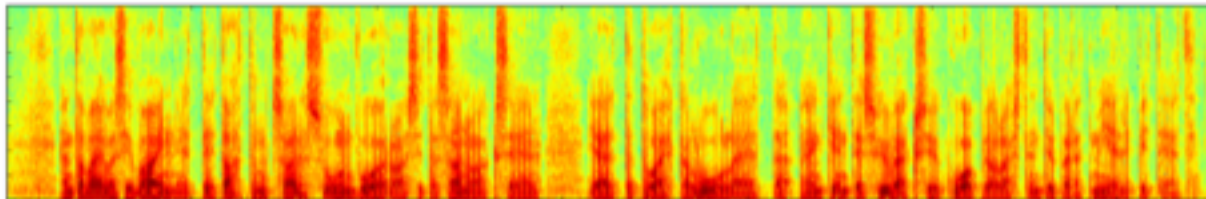


Introduction

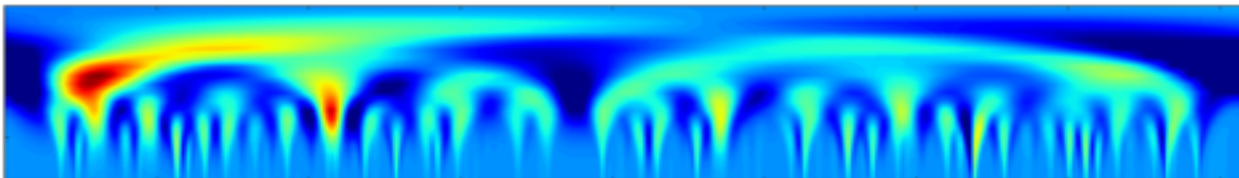
- Thus, what is sought after, is a representation for prosody, where the contribution of different phonological layers is distinguishable: **Continuous wavelet analysis**



↓ Short time fourier transform



↴ Continuous wavelet analysis



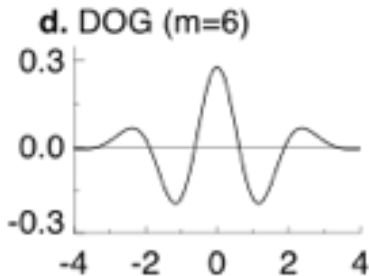
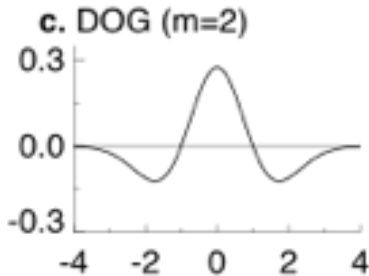
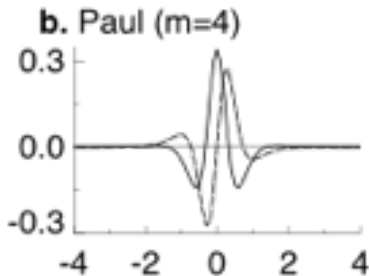
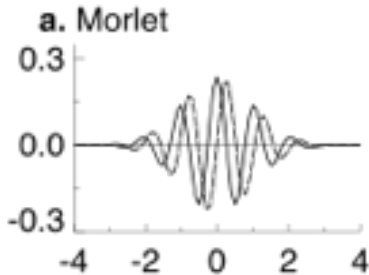
Definition

The diagram shows the wavelet transform equation with several labels and arrows pointing to specific parts of the equation:

- time*: points to the variable t in the first argument of $X(t, s)$.
- signal*: points to the function $x(\tau)$ inside the integral.
- wavelet*: points to the function $\bar{\psi}$ inside the integral.
- scale*: points to the variable s in the second argument of $X(t, s)$ and also to the denominator $|s|^{1/2}$.

$$X(t, s) = \frac{1}{|s|^{1/2}} \int_{-\infty}^{\infty} x(\tau) \bar{\psi}\left(\frac{\tau - t}{s}\right) d\tau$$

Mother wavelet



Wavelet: scaling



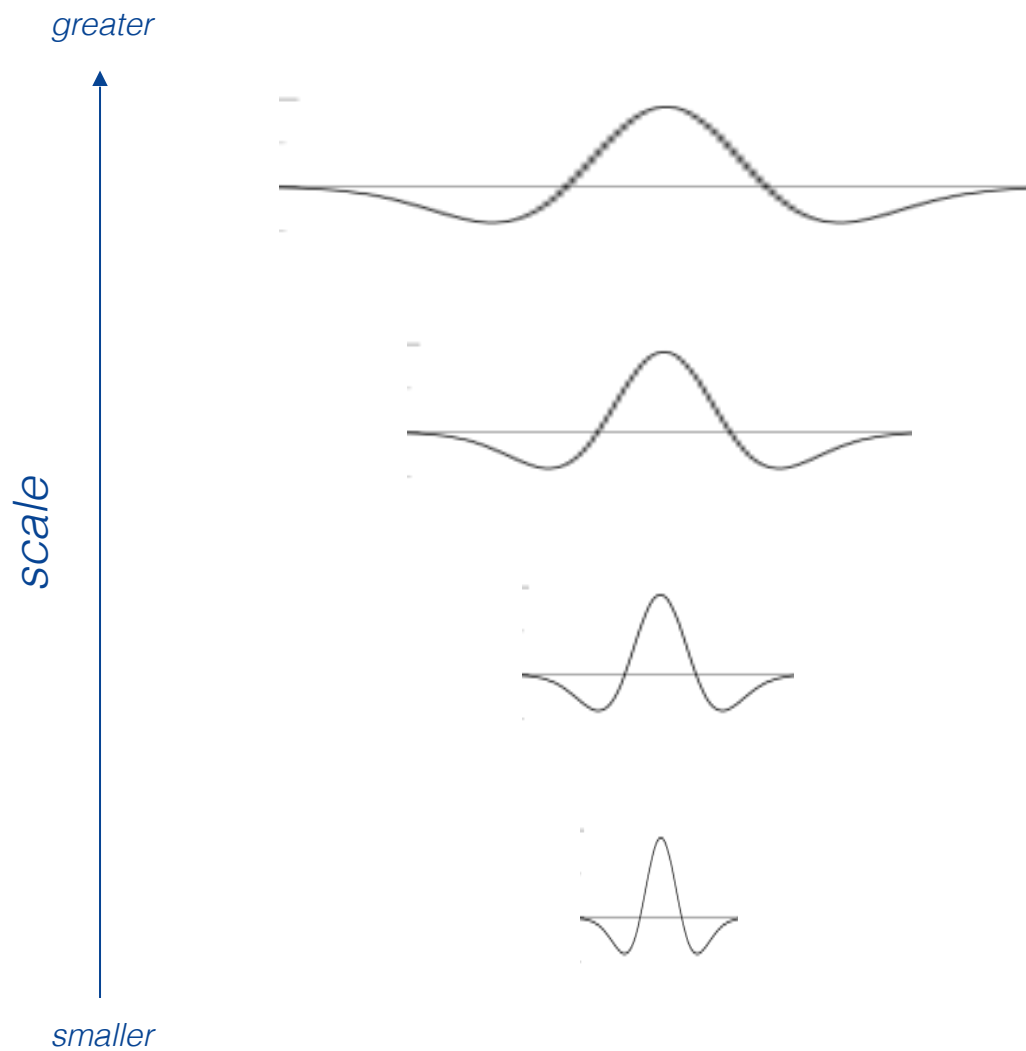
$$\psi\left(\frac{\tau - t}{s}\right)$$

s

↑

scale

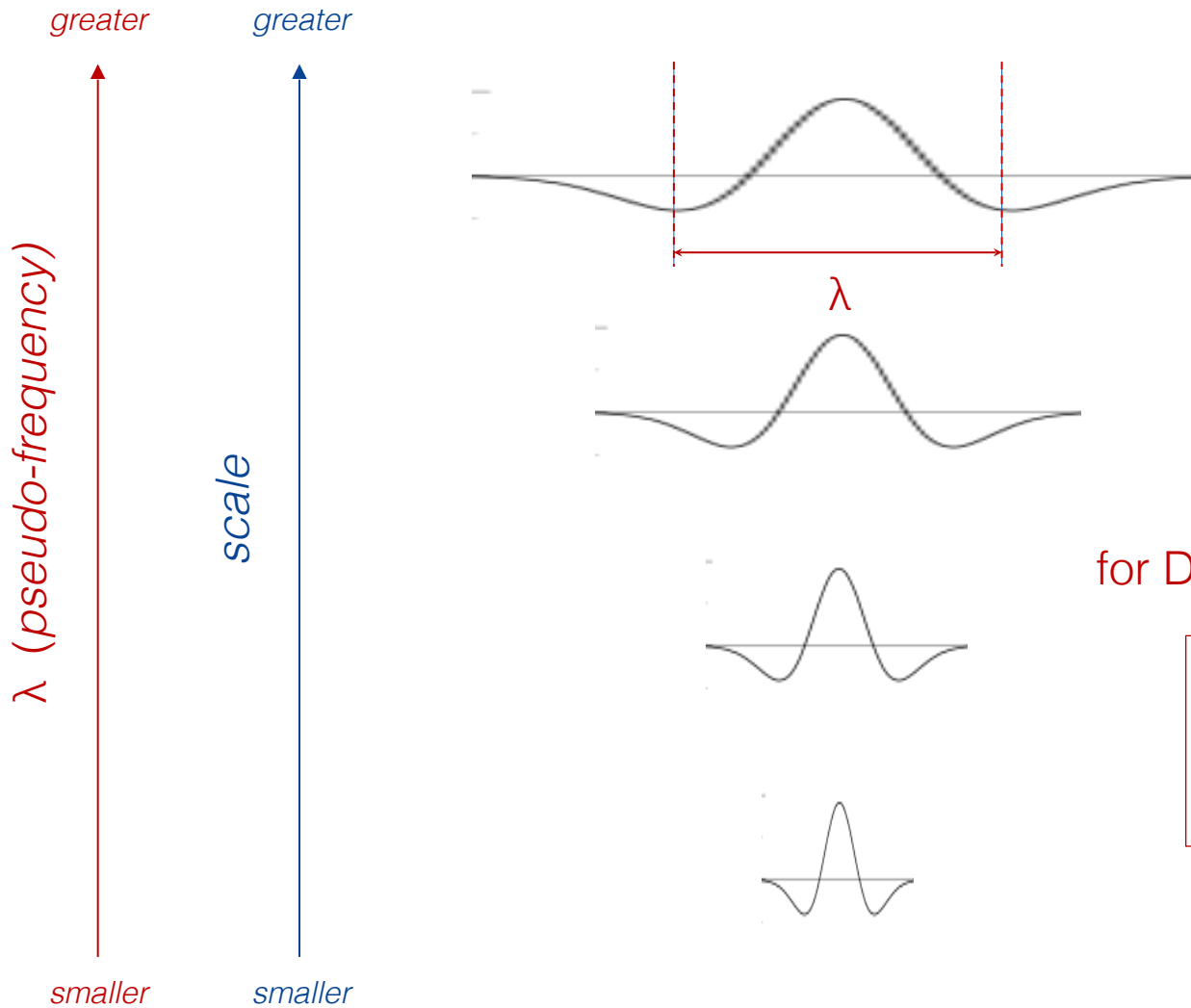
Wavelet: scaling



$$\psi\left(\frac{\tau - t}{s}\right)$$

↑
scale

Wavelet: scaling



for DOG (m=2):

$$\lambda = \frac{2\pi s}{\sqrt{5/2}}$$

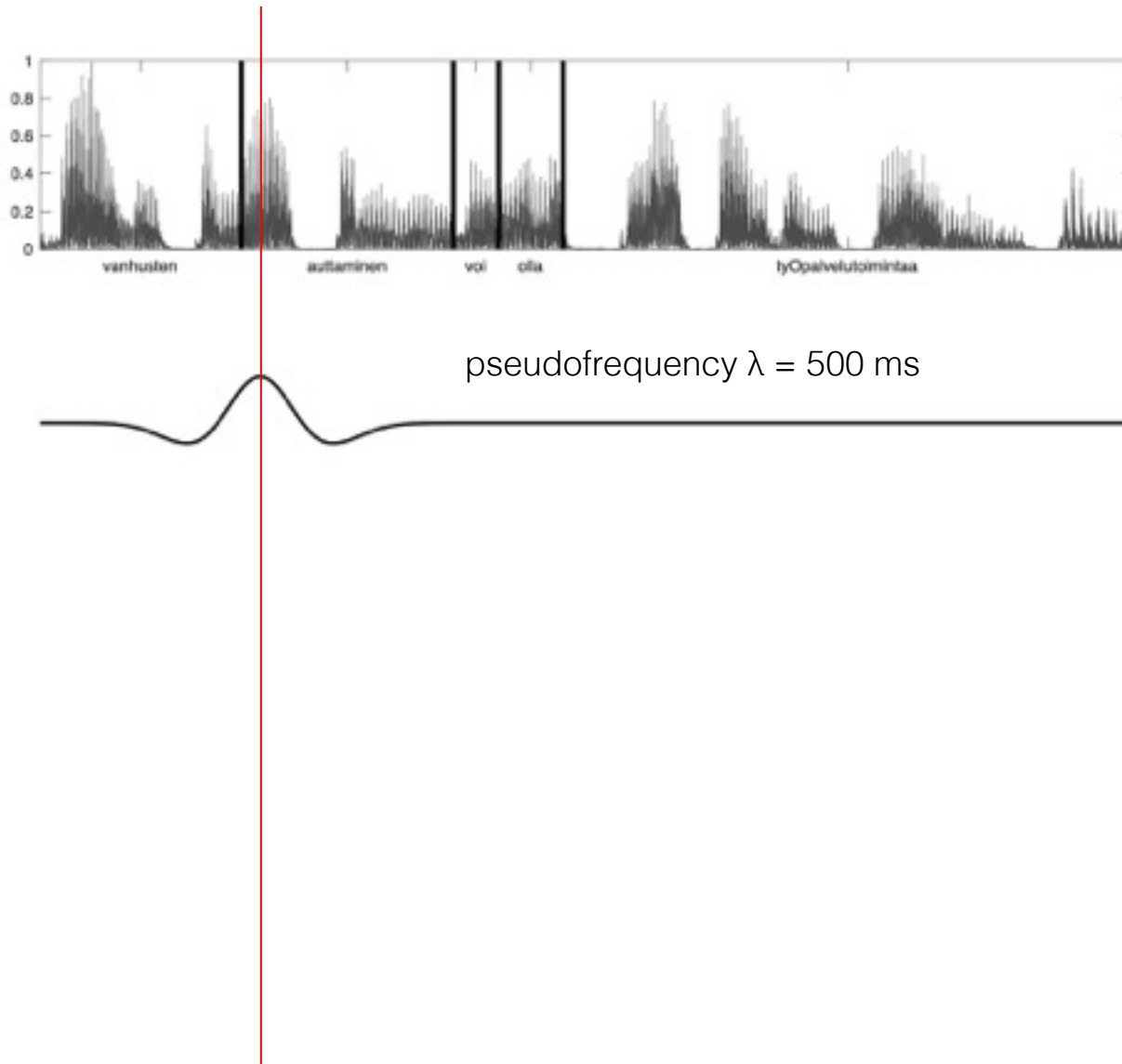
linear relationship

Convolution

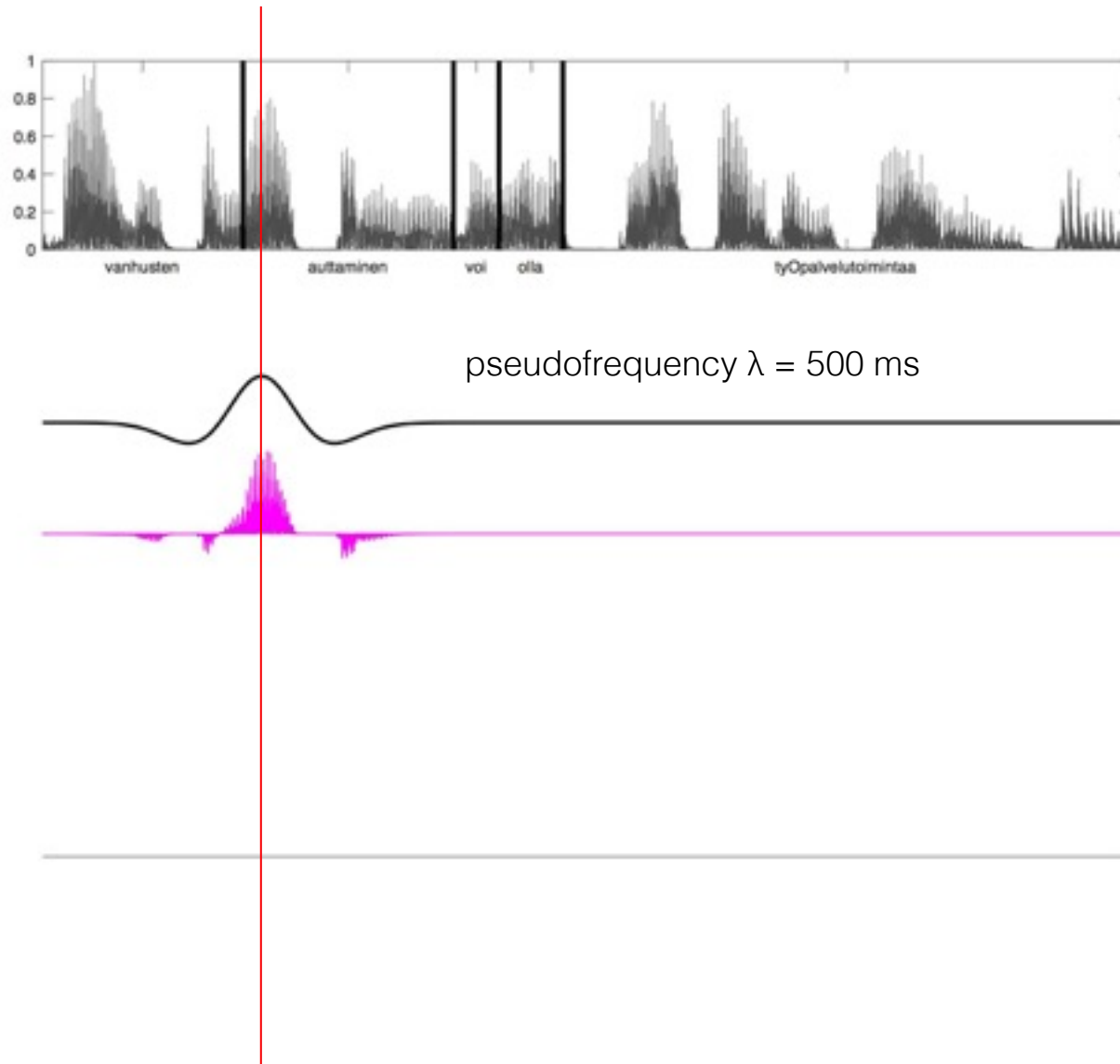
$$X(t, s) = \frac{1}{|s|^{1/2}} \int_{-\infty}^{\infty} x(\tau) \overline{\psi\left(\frac{\tau - t}{s}\right)} d\tau$$

time → t
scale → s
signal → $x(\tau)$
wavelet → $\overline{\psi\left(\frac{\tau - t}{s}\right)}$
convolution → $\int_{-\infty}^{\infty} x(\tau) \overline{\psi\left(\frac{\tau - t}{s}\right)} d\tau$

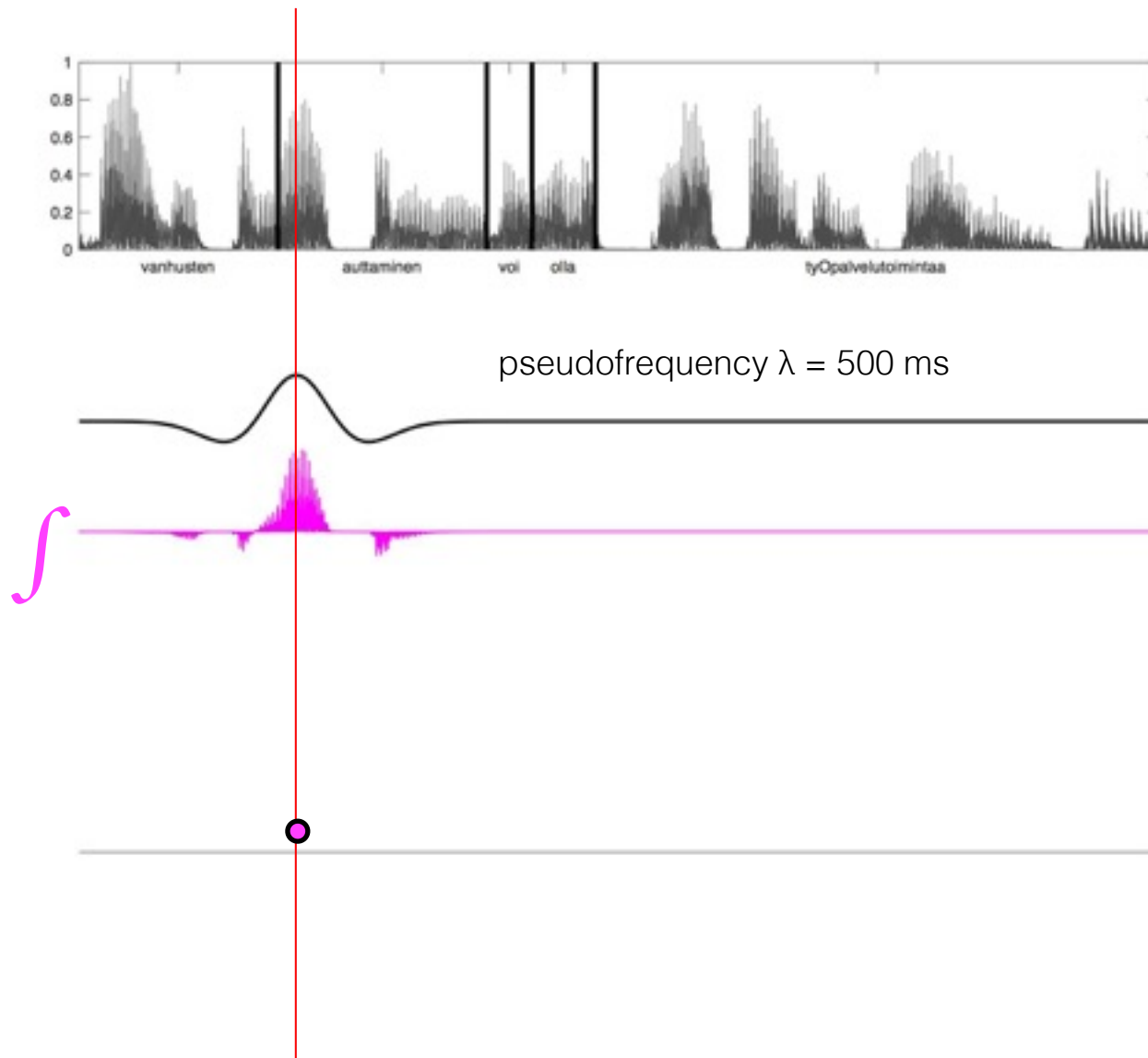
Convolution



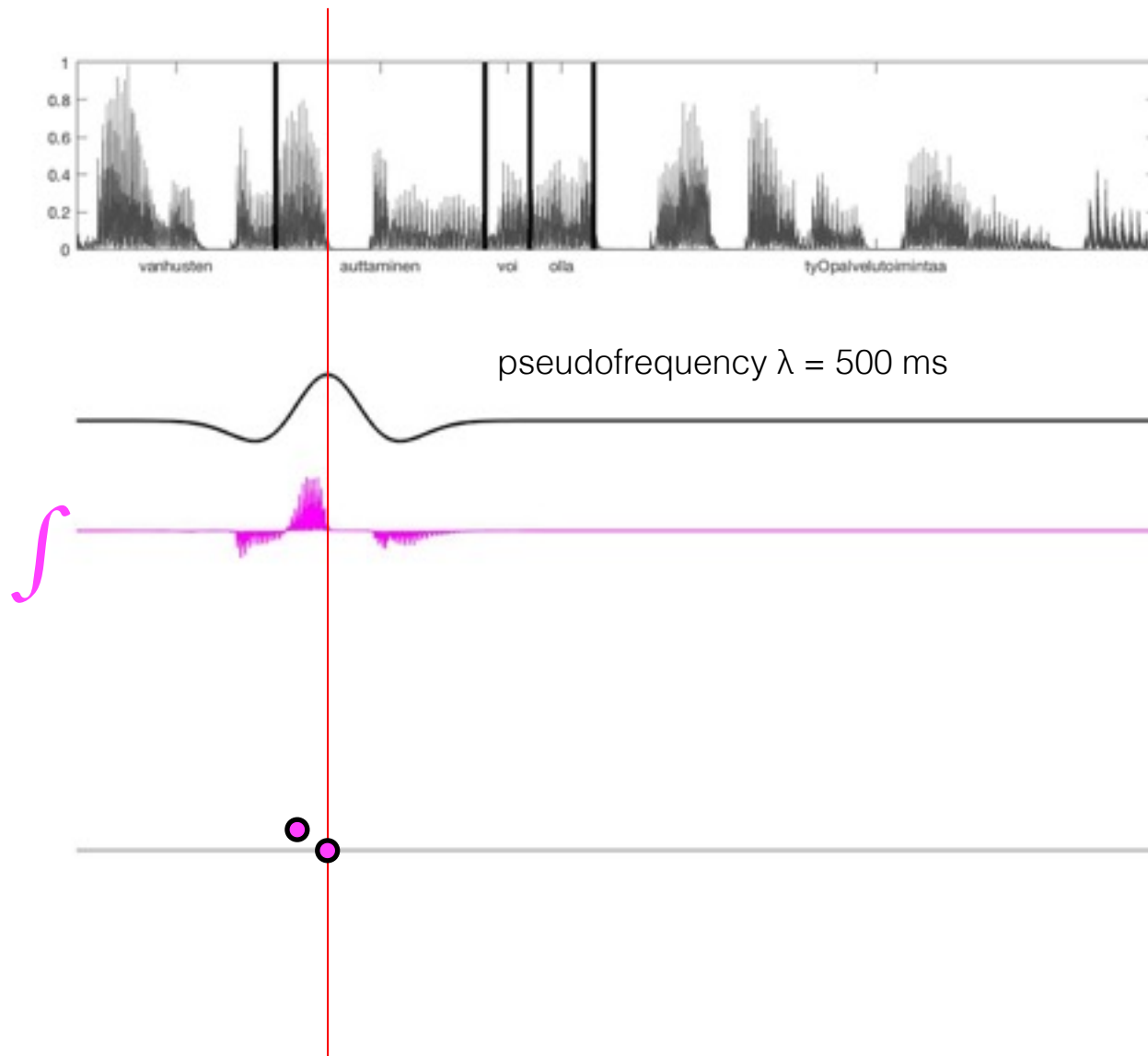
Convolution



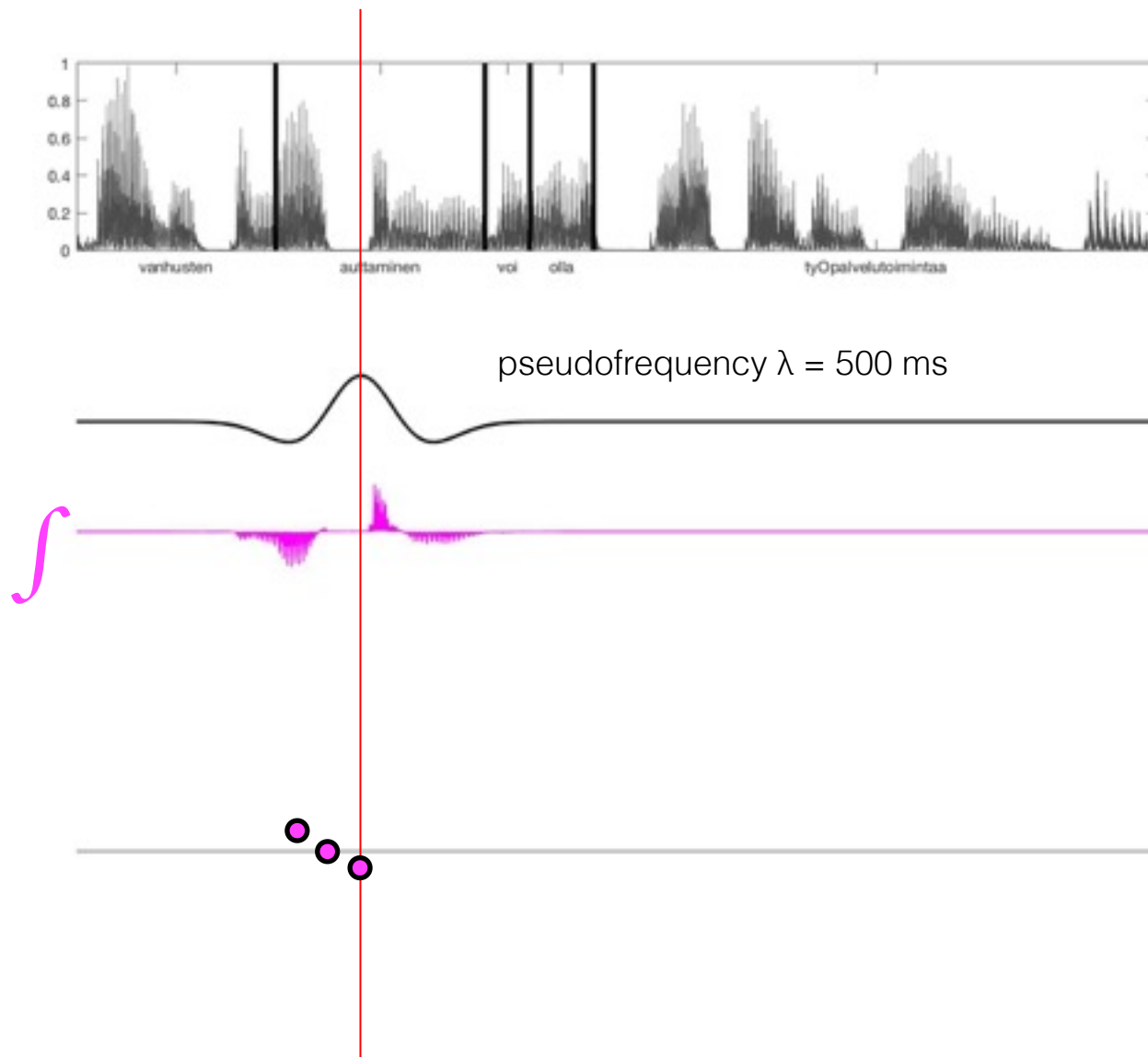
Convolution



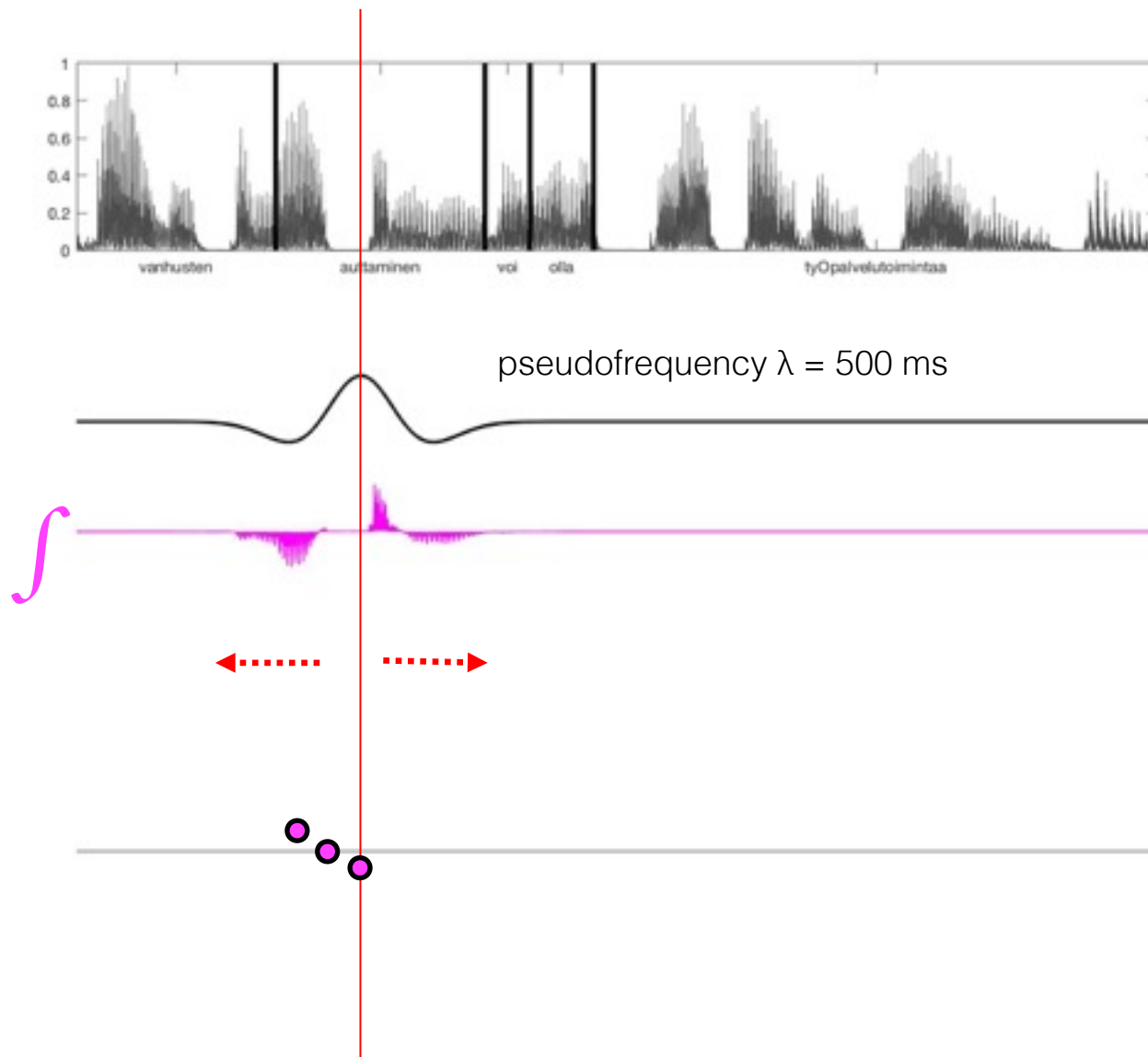
Convolution



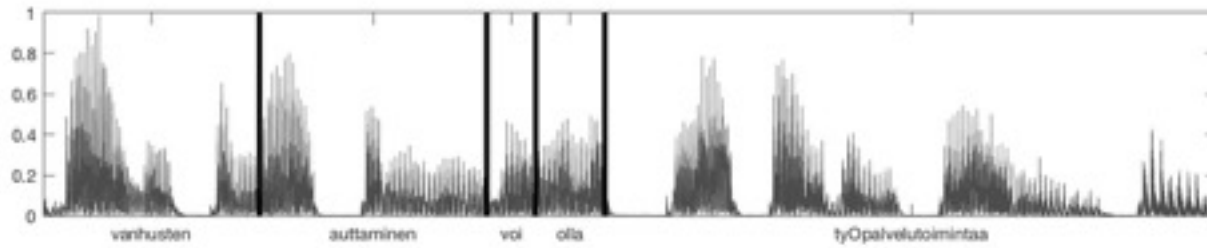
Convolution



Convolution

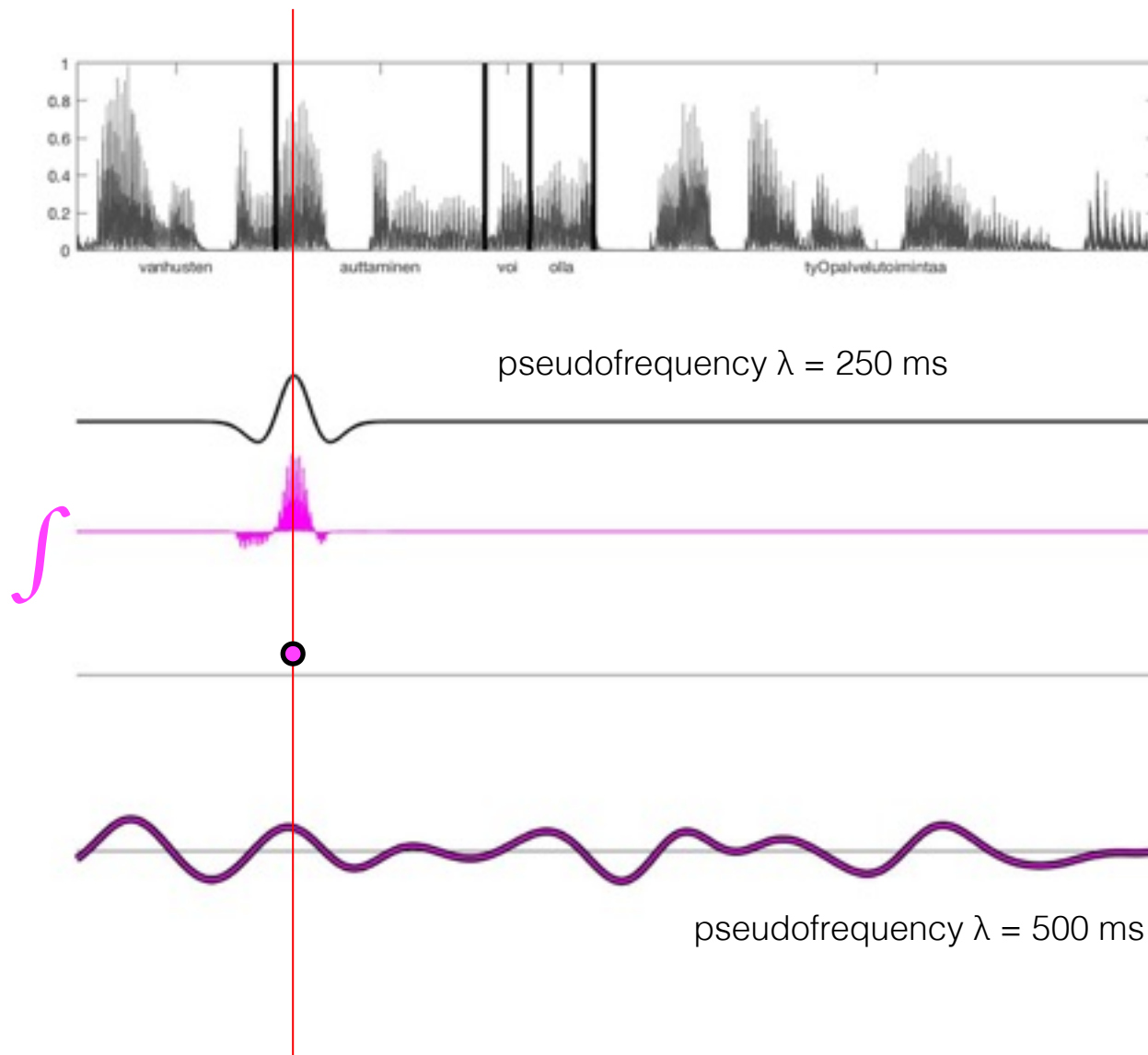


Convolution

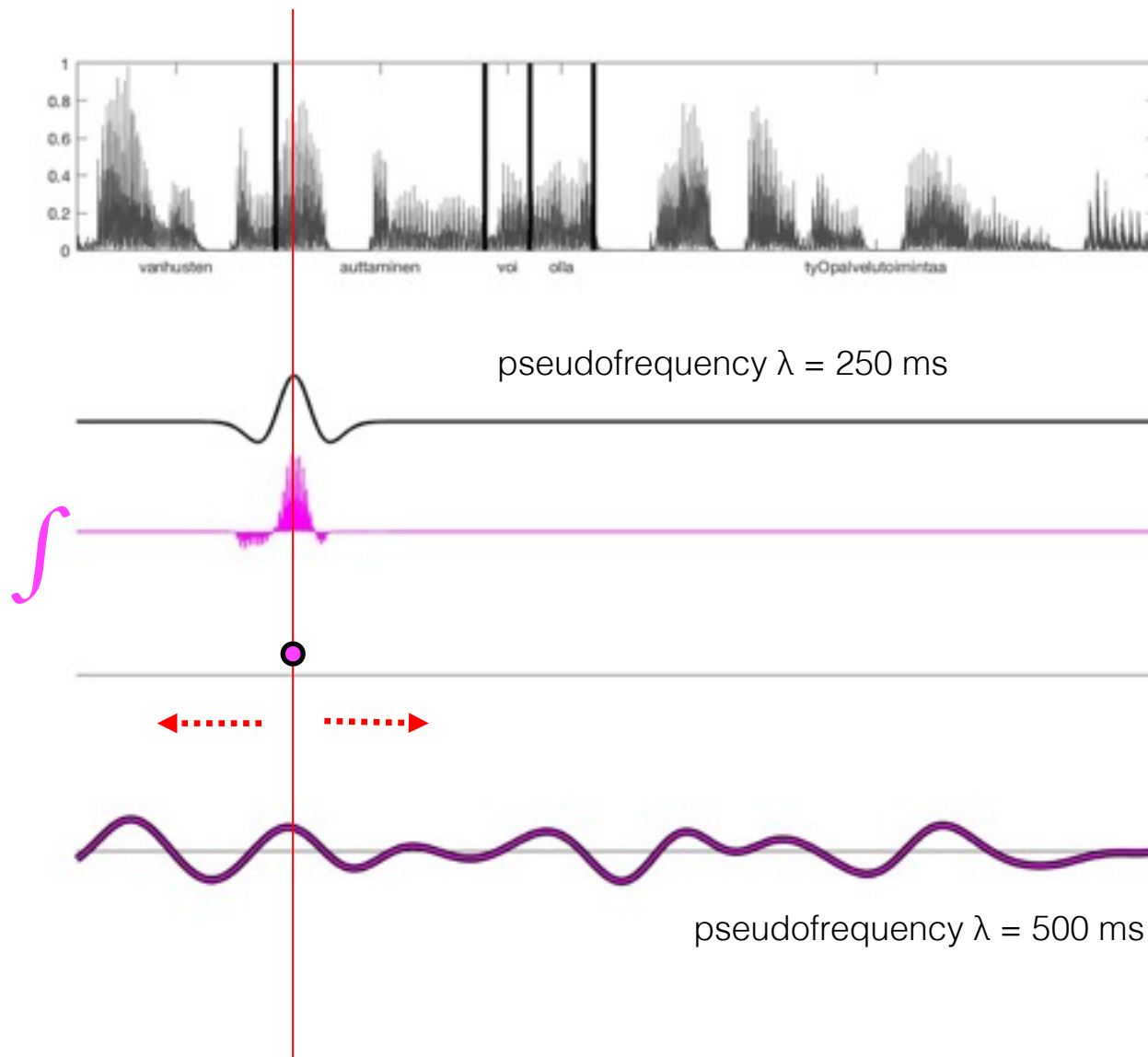


pseudofrequency $\lambda = 500$ ms

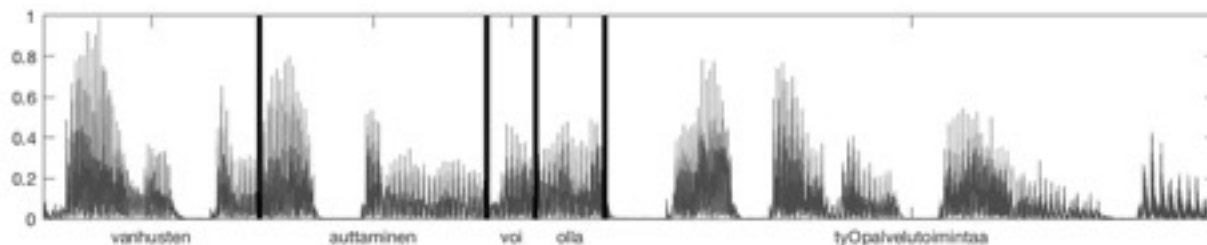
Convolution



Convolution



Convolution

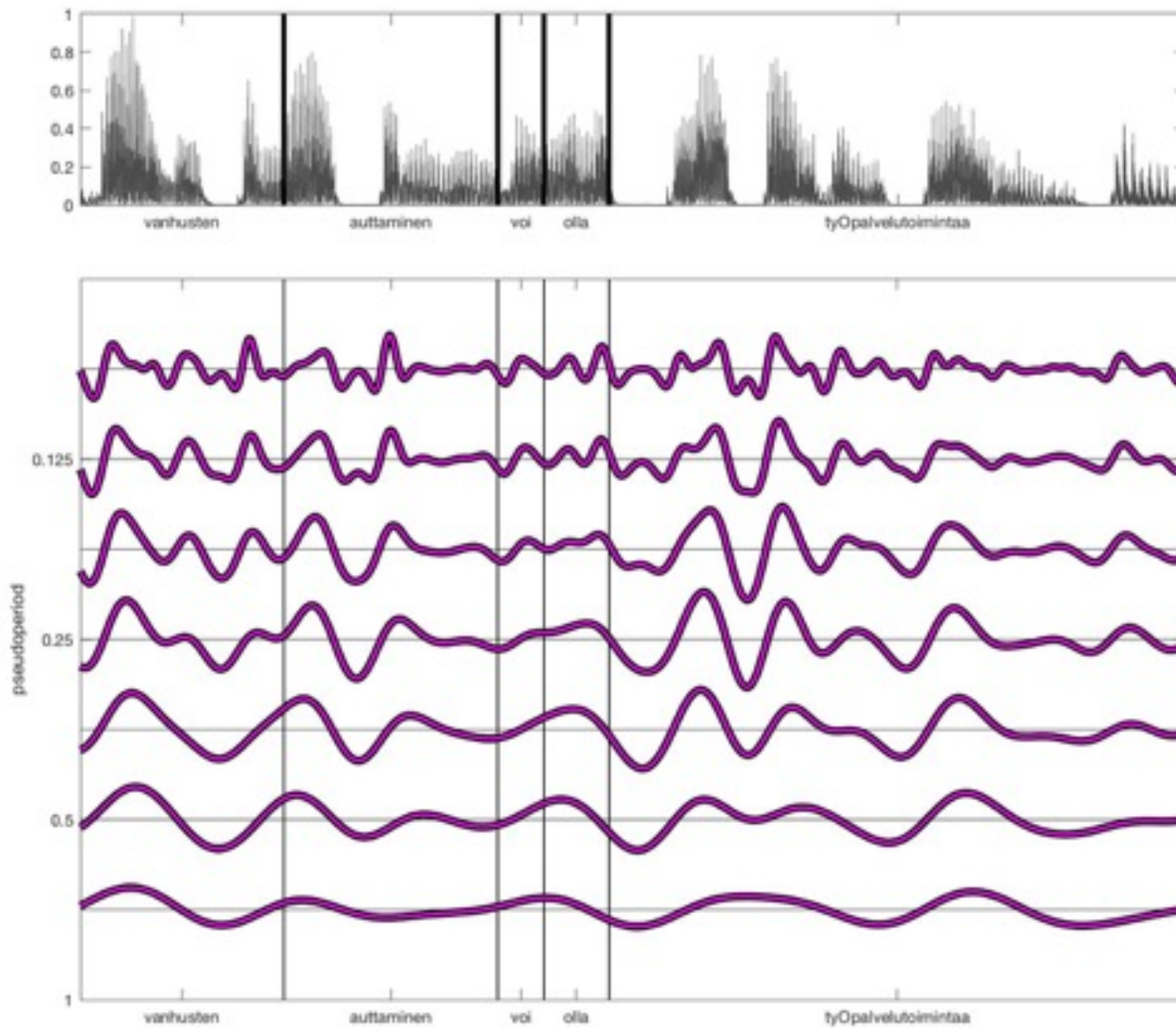


pseudofrequency $\lambda = 250$ ms

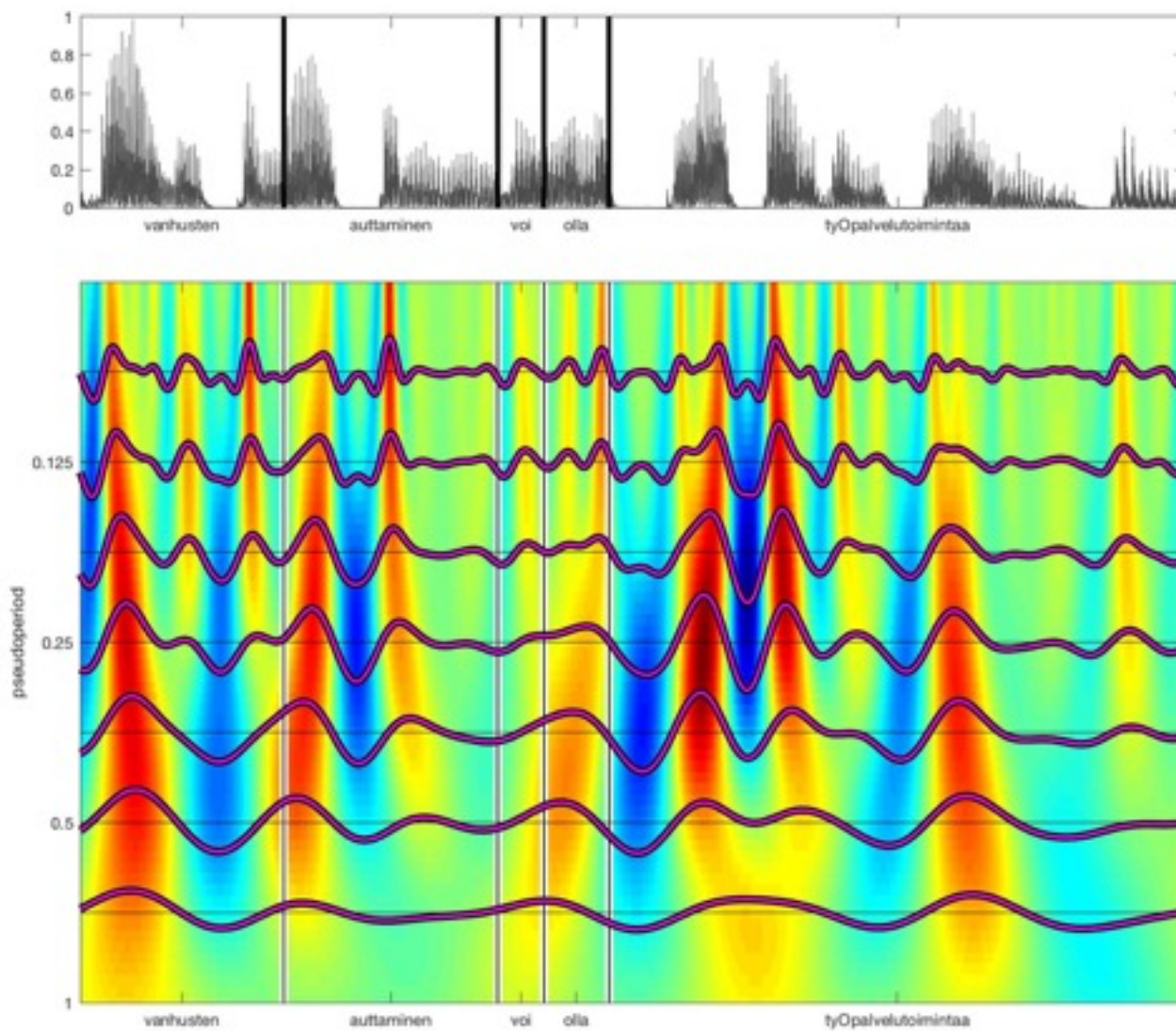


pseudofrequency $\lambda = 500$ ms

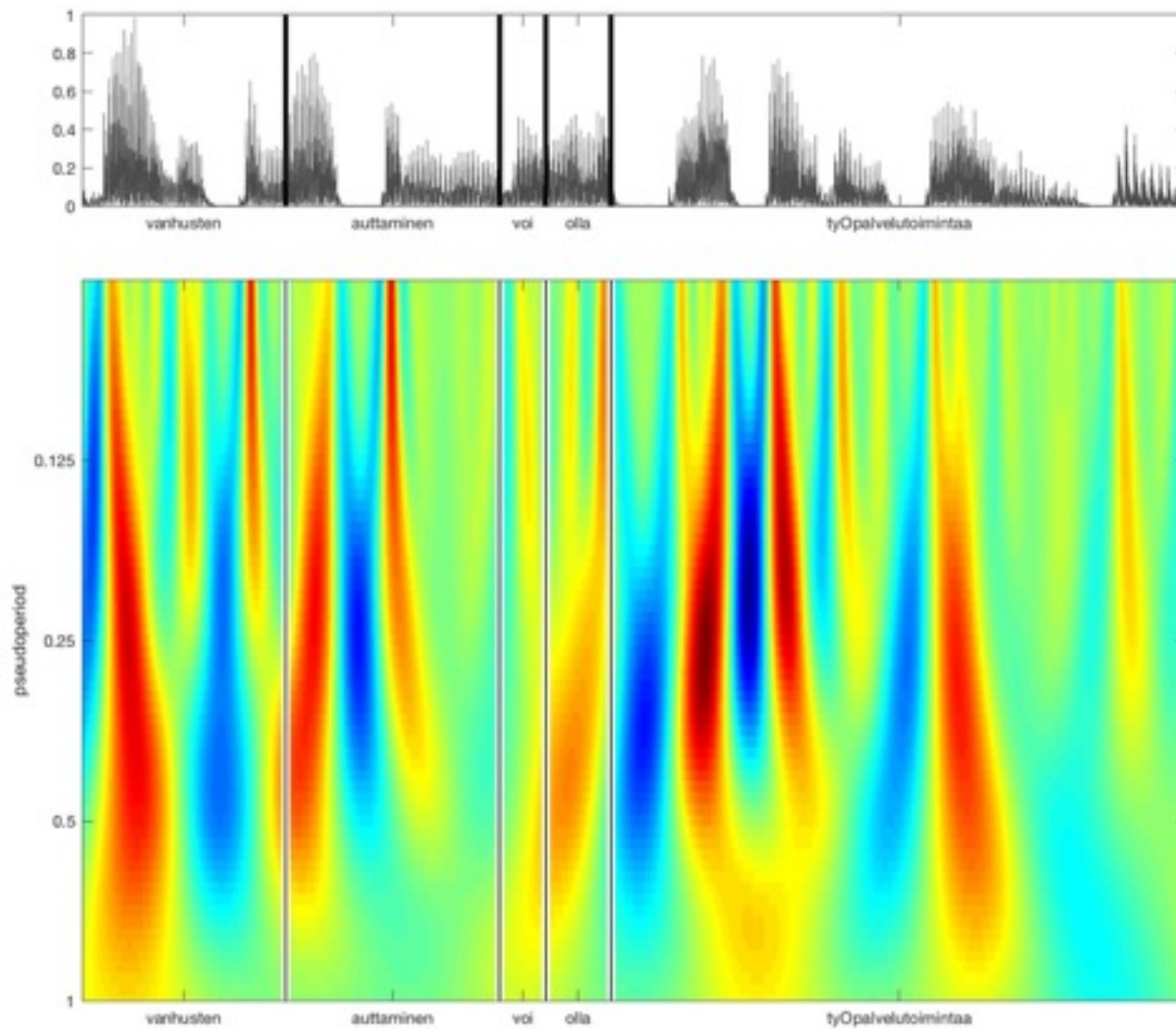
Continuous wavelet transform



Scalogram



Scalogram



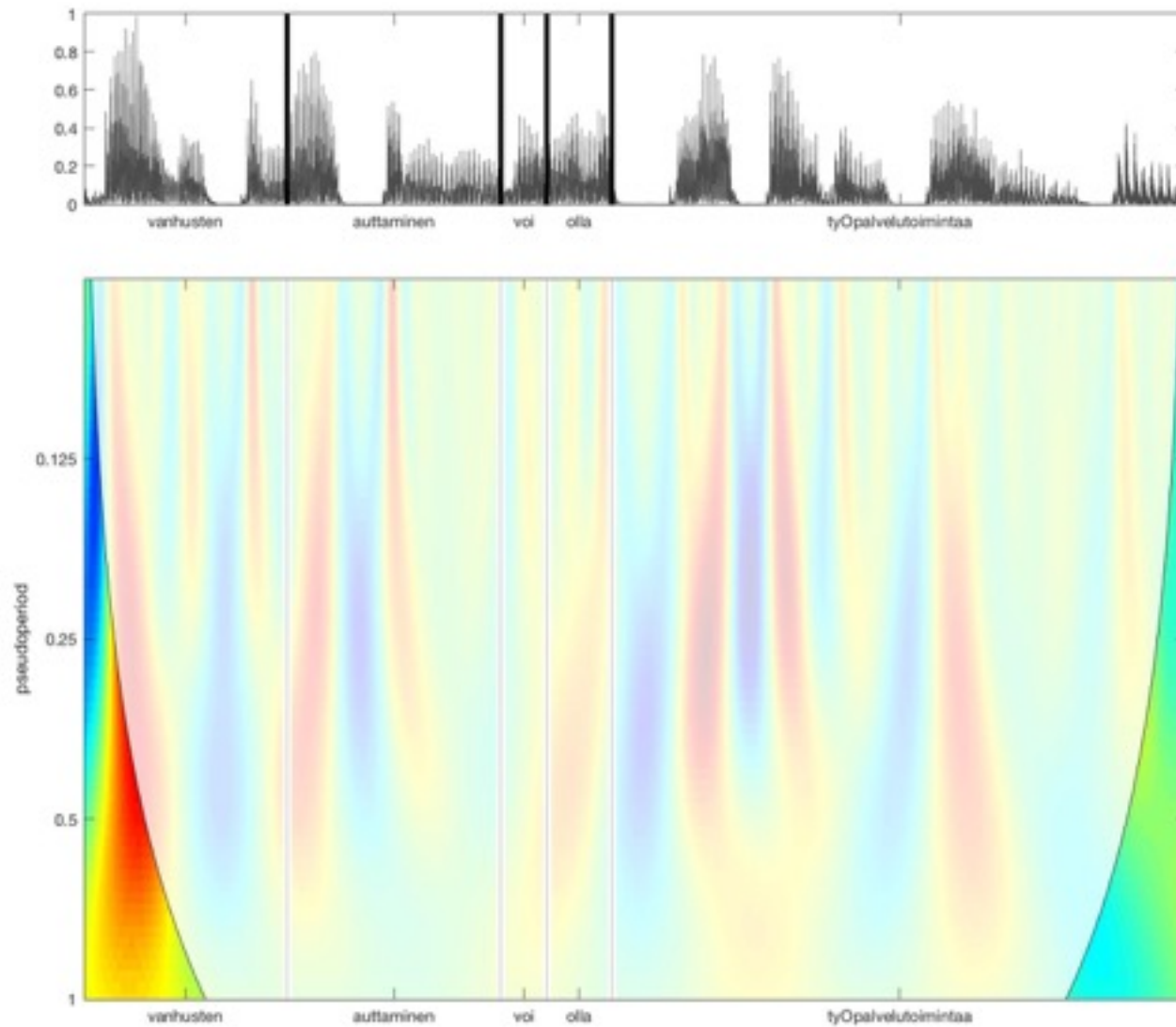
Computation of CWT

- is quite fast, uses Convolution Theorem

$$f * g = iFFT\left(c.FFT(f).FFT(g)\right)$$

- pads the signal (left and right) with zeros
- cone of influence

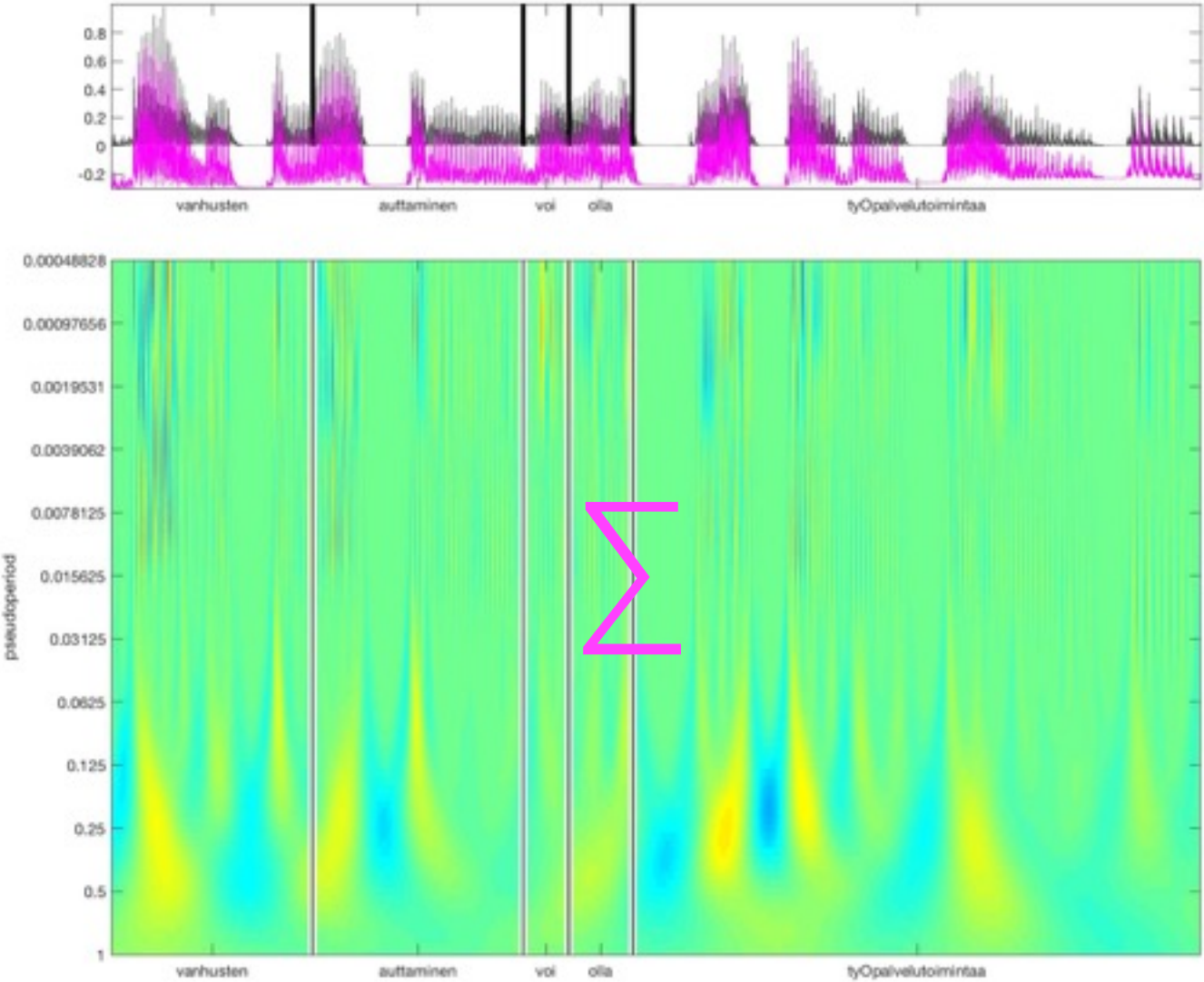
Cone of influence



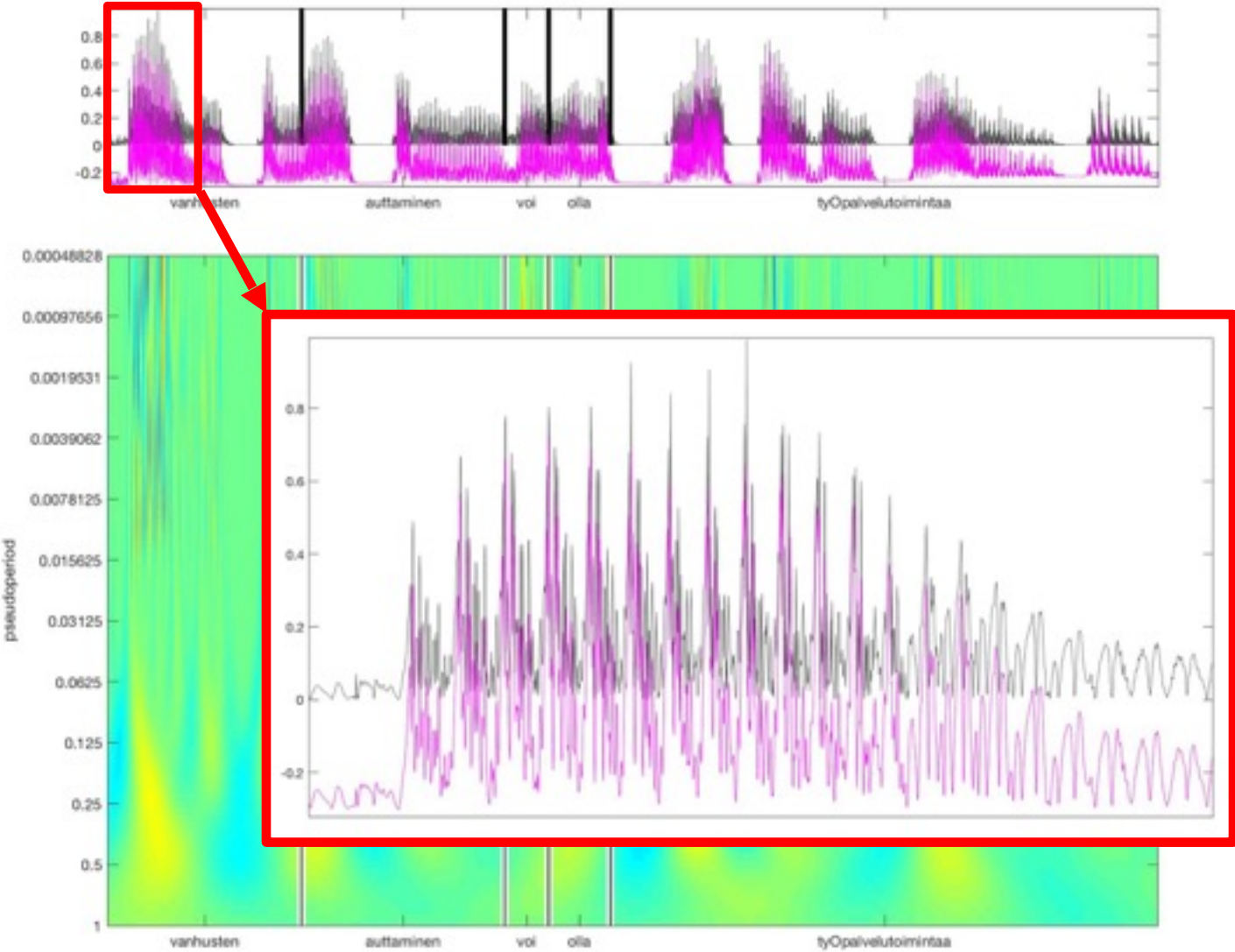
Computation of CWT

- packages in python, Matlab ...
- sources:
 - Daubechies, I. (1992). *Then Lectures on Wavelets*, Society of Industrial and Applied Mathematics.
 - Torrence & Compo (1998). A practical guide to wavelet analysis, *Bulletin of the American Meteorological Society*, 79(1), pp. 61-78
 - Grinsted, Moore and Jevrejeva. (2004). Application of the cross-wavelet transform and wavelet coherence to geographical time series. *Nonlinear processes in Geophysics*, 11(5/6), pp. 561-566.
 - *and many many more*

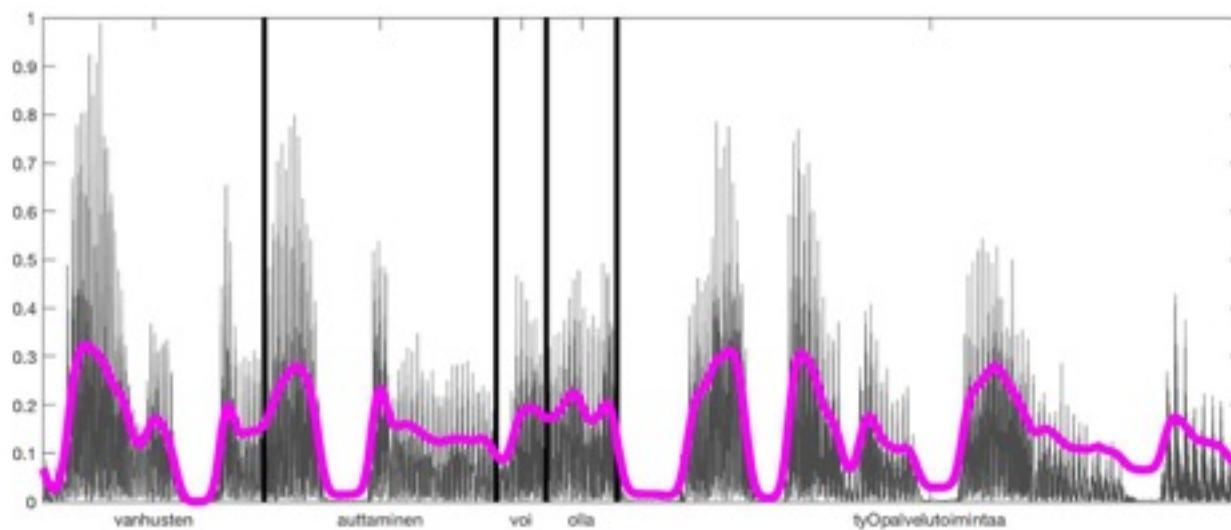
Decomposition: CWT is (almost) invertible



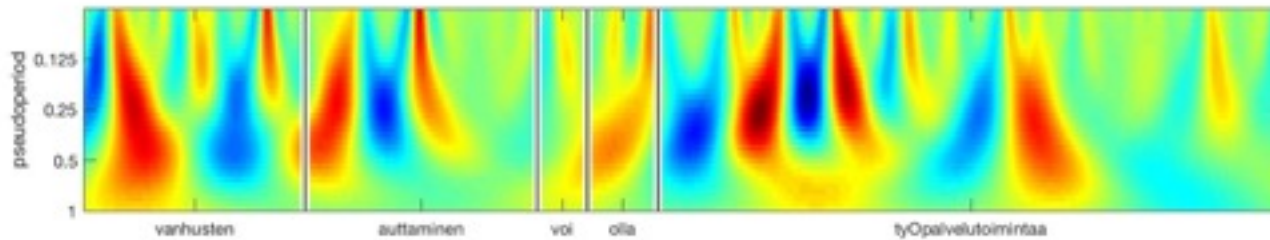
Decomposition: CWT is (almost) invertible



CWT as a filter

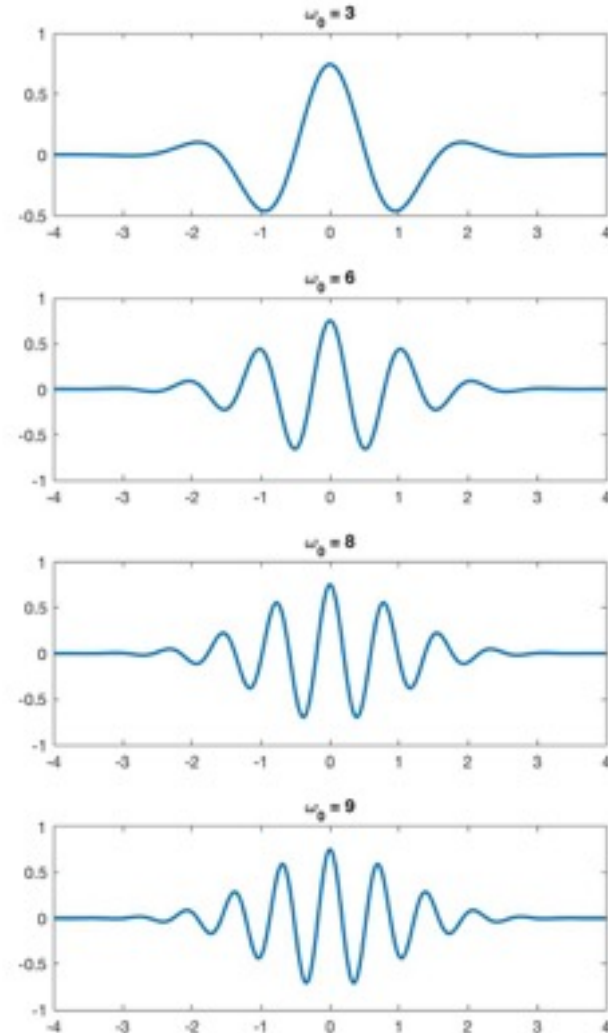
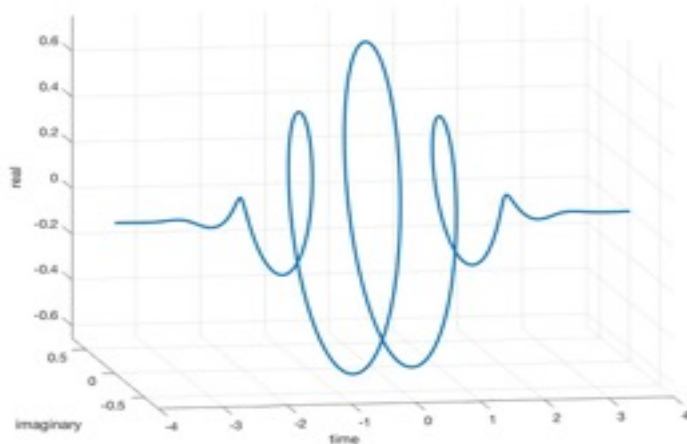


low-pass filter, approximation of signal energy



Morlet mother wavelet

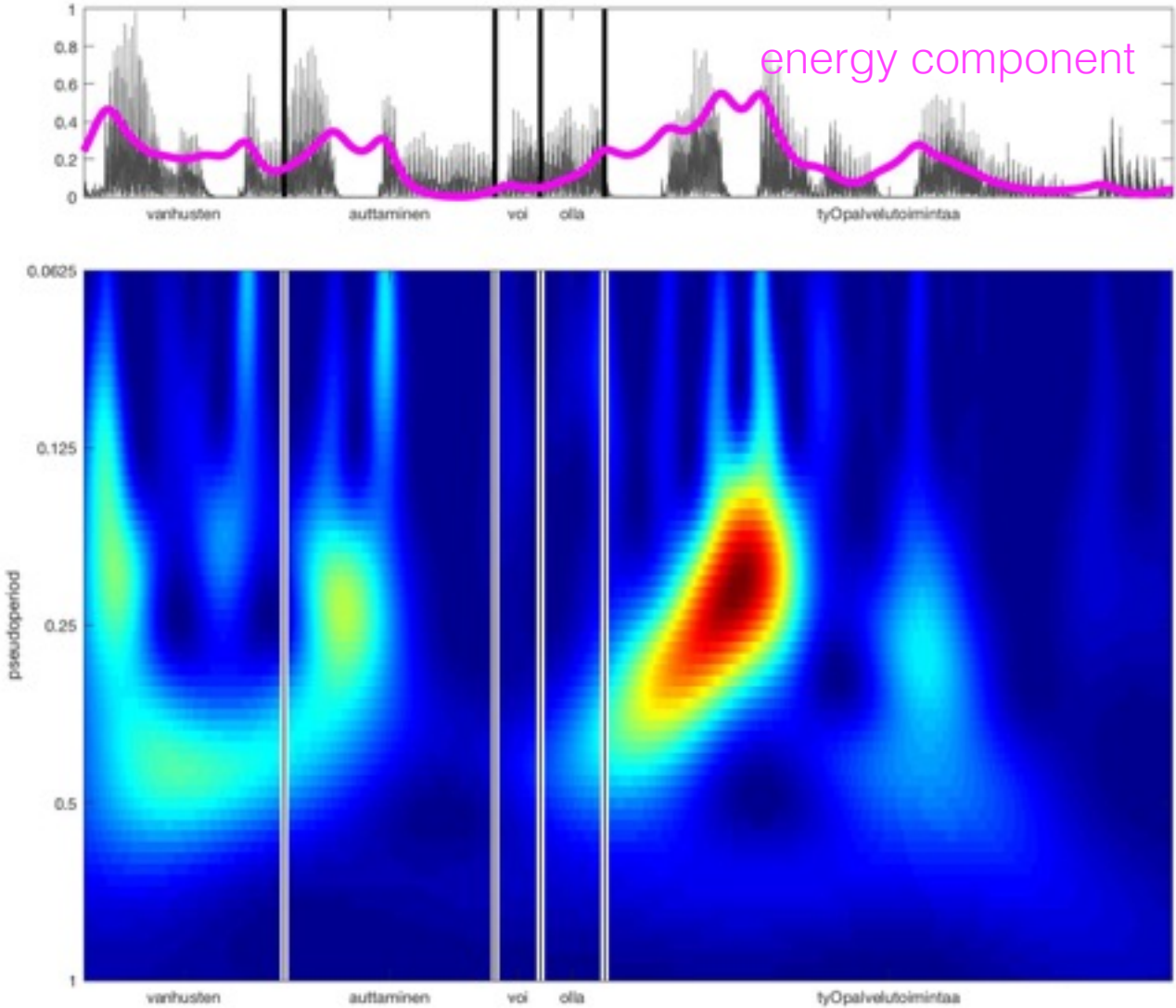
- complex wavelet
- complex exponential ($e^{i\omega t}$) within Gaussian envelope



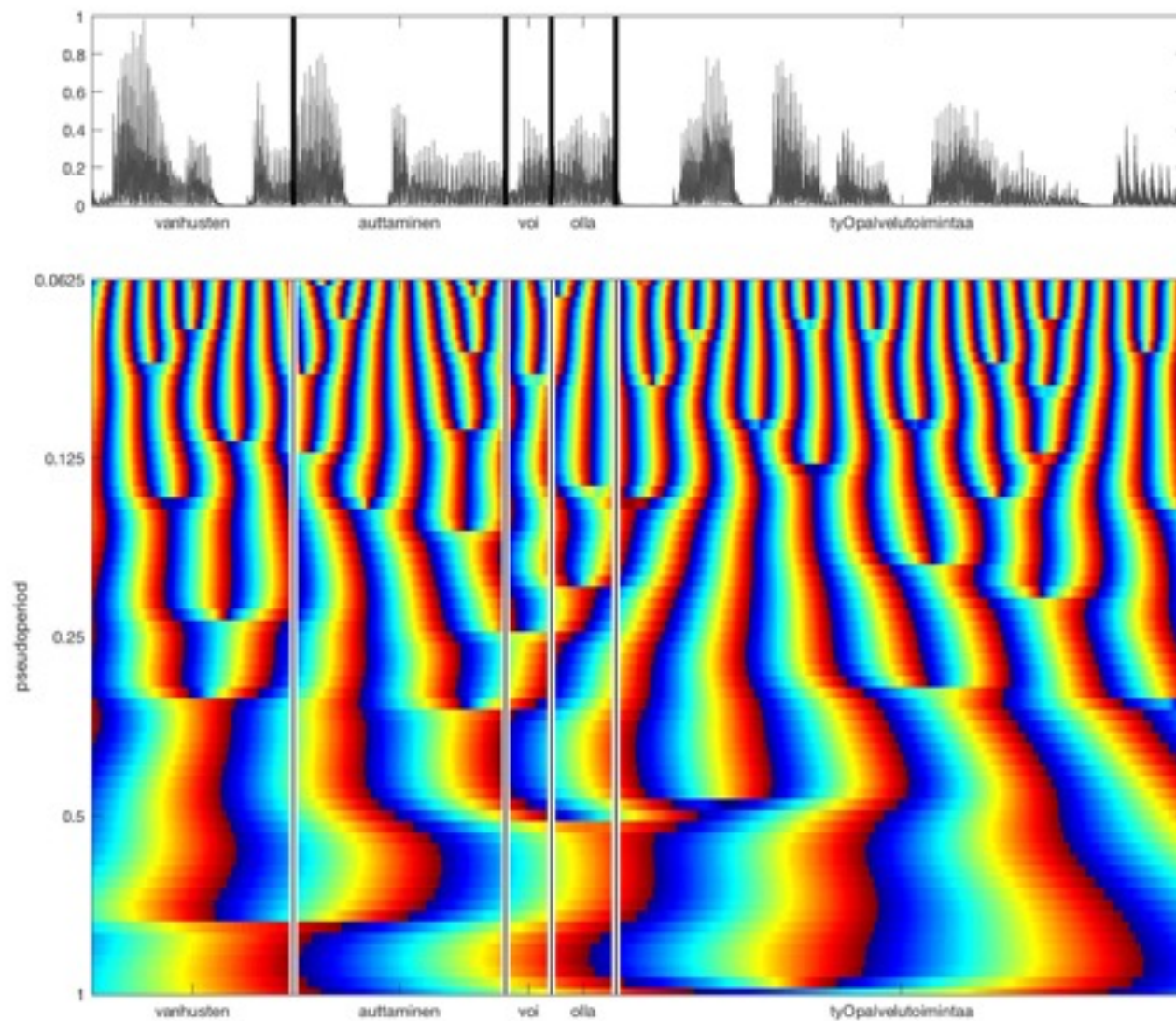
good temporal resolution

good frequency resolution

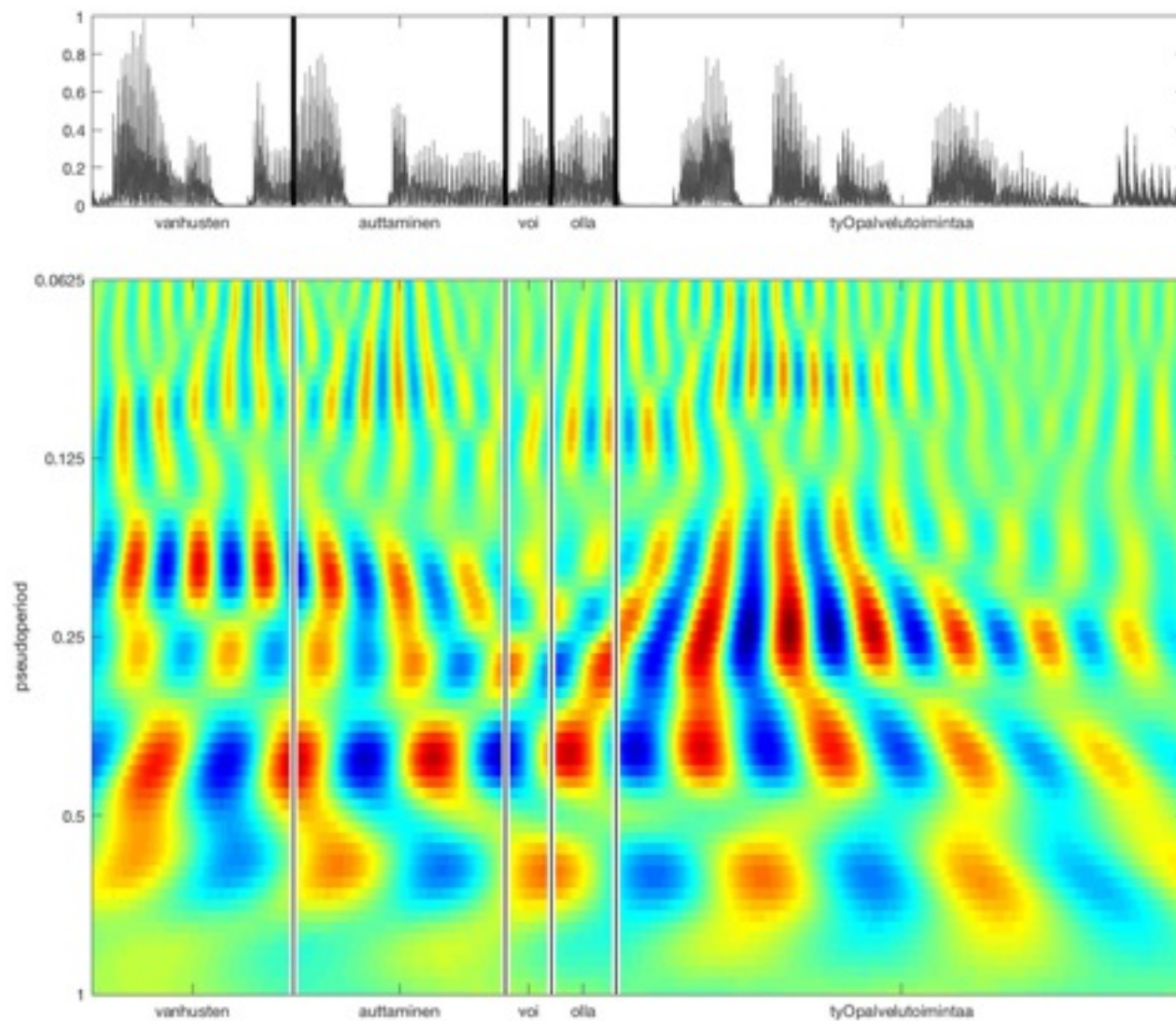
Energy – Morlet (amplitude): $\omega_0 = 3$



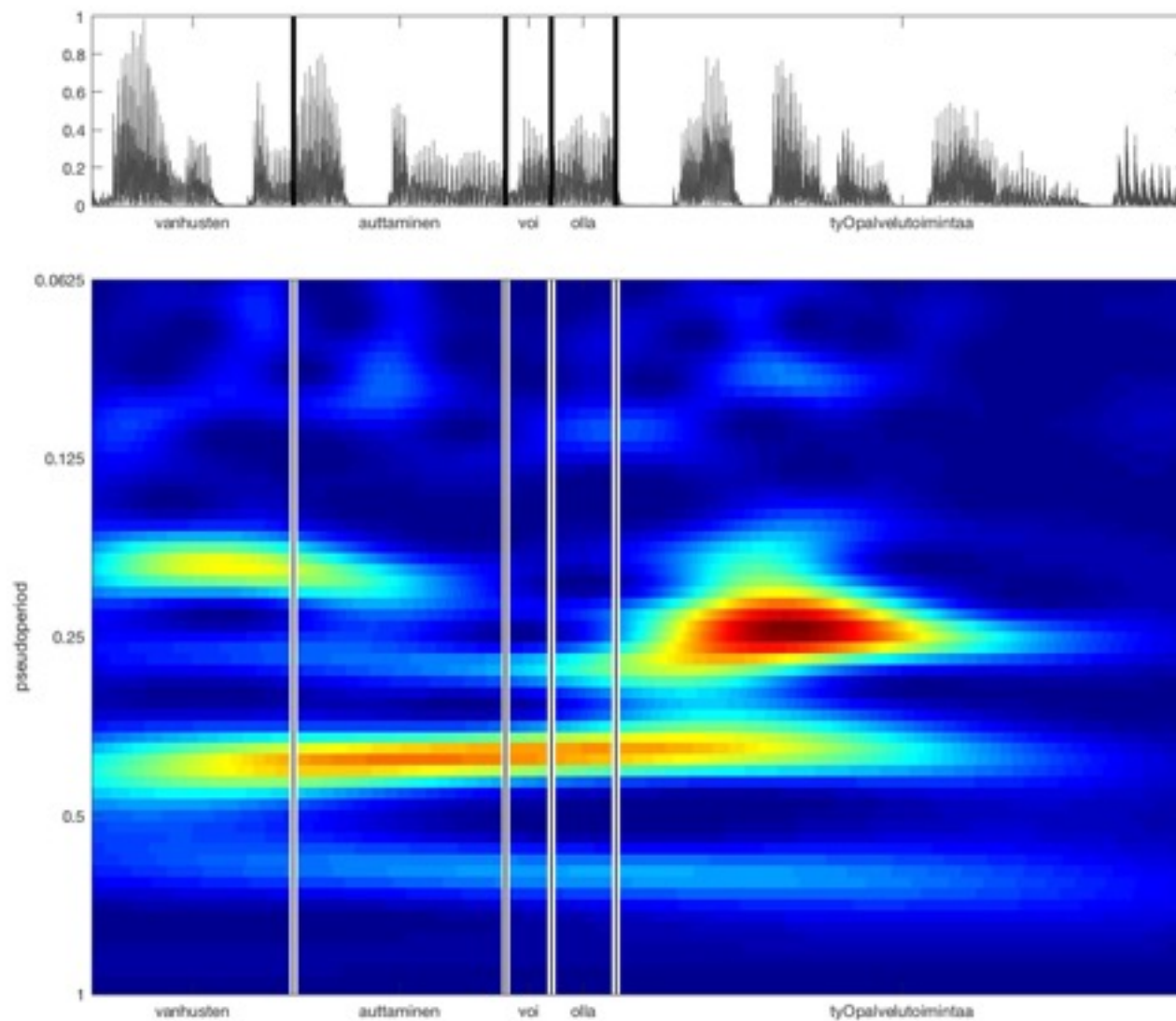
Phase – Morlet (amplitude): $\omega_0 = 10$



Morlet (real part): $\omega_0 = 10$

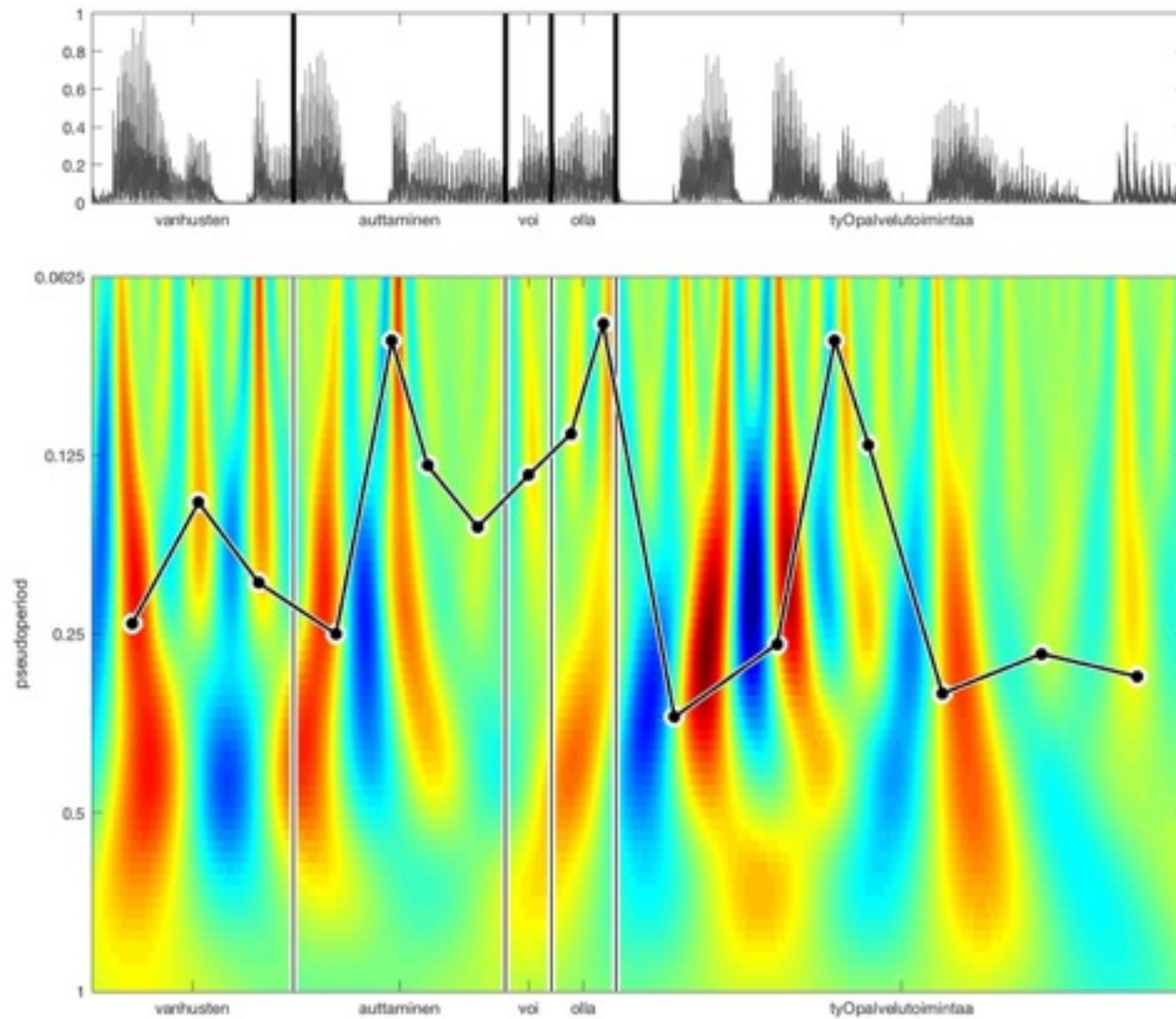


Energy – Morlet (amplitude): $\omega_0 = 10$

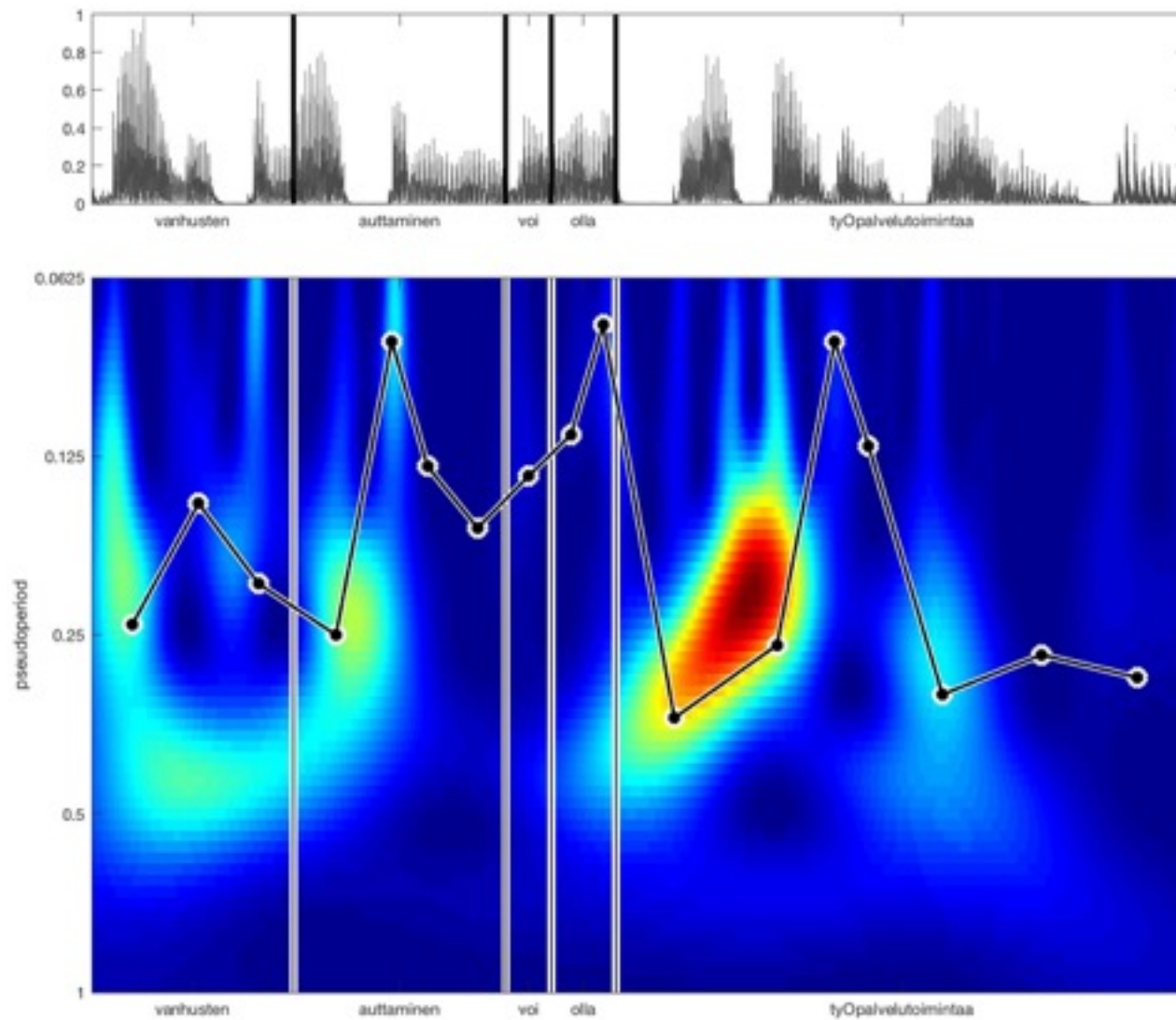


Continuous Wavelet Analysis and Speech

Syllable structure (Morlet, $\omega_0 = 3$)

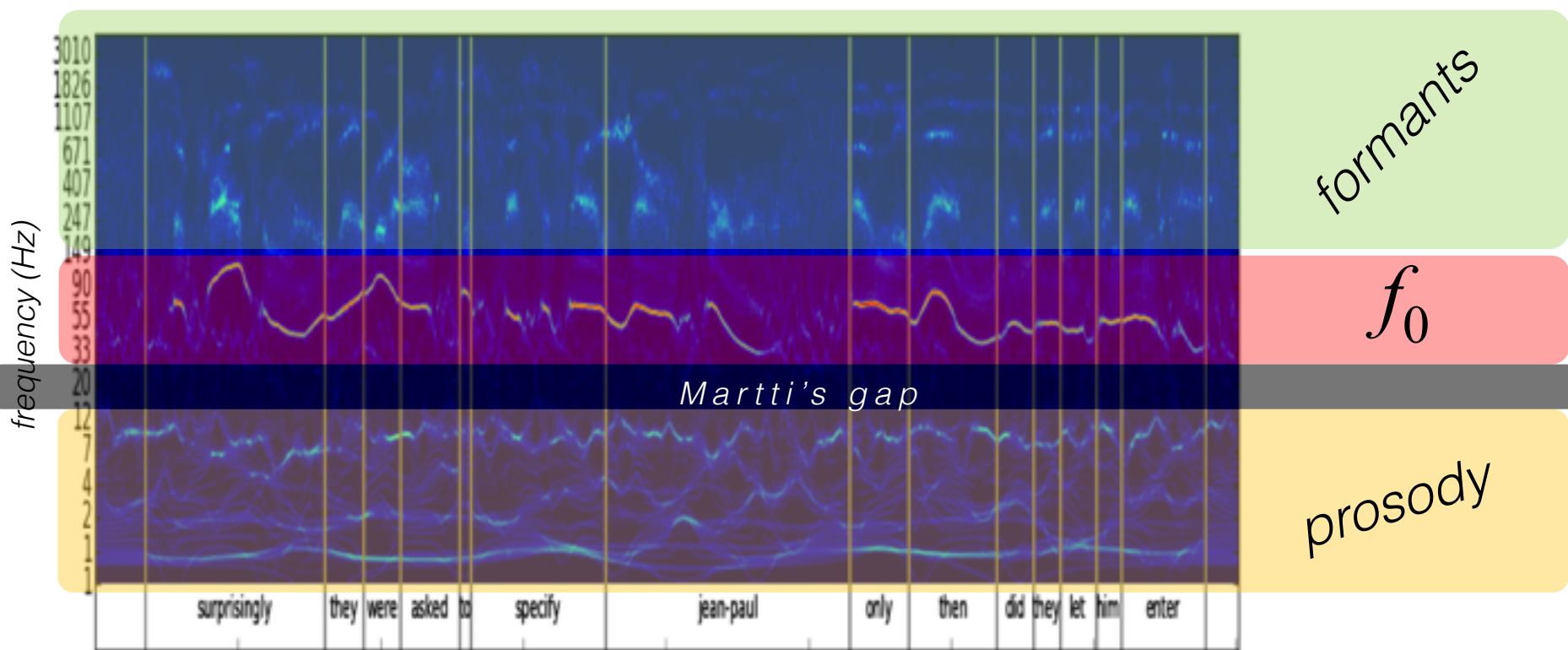


Prominence (Morlet, $\omega_0 = 3$)

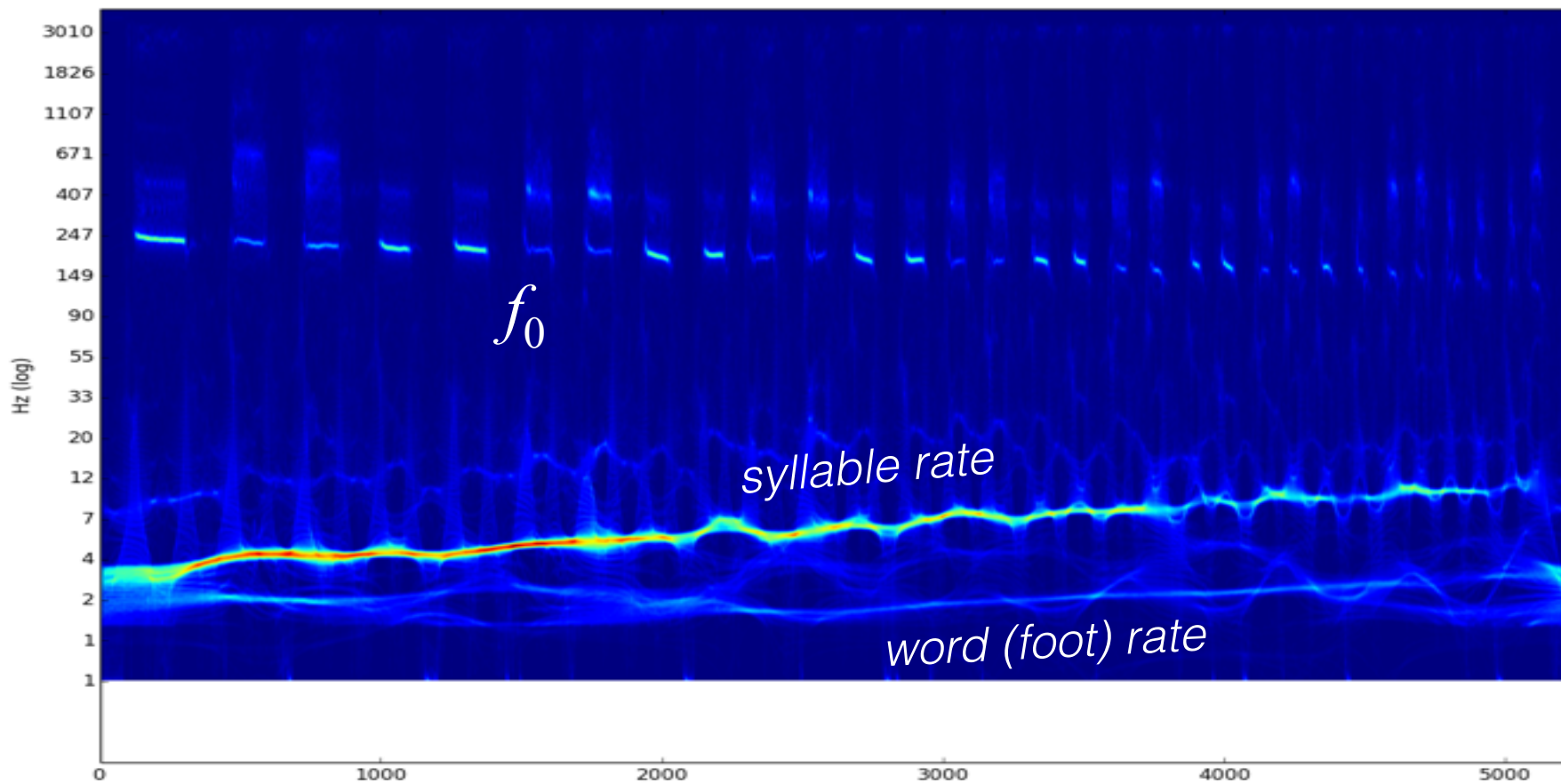


A “complete” (upside-down) scalogram

- quite heavily engineered: calculated instantaneous frequencies at each scales and then plotted in frequency domain (aka Hilbert spectrum)
- huge range of frequencies (compared to Fourier Transform)

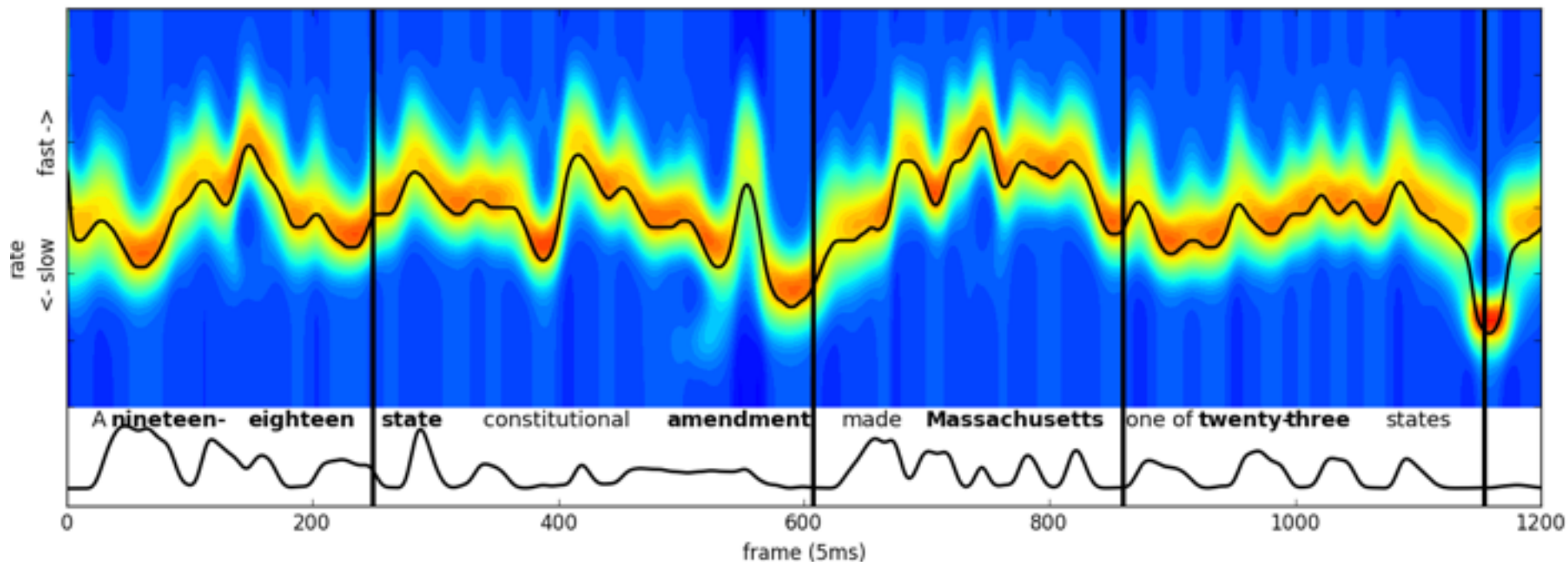


f_0 and speaking rate(s)



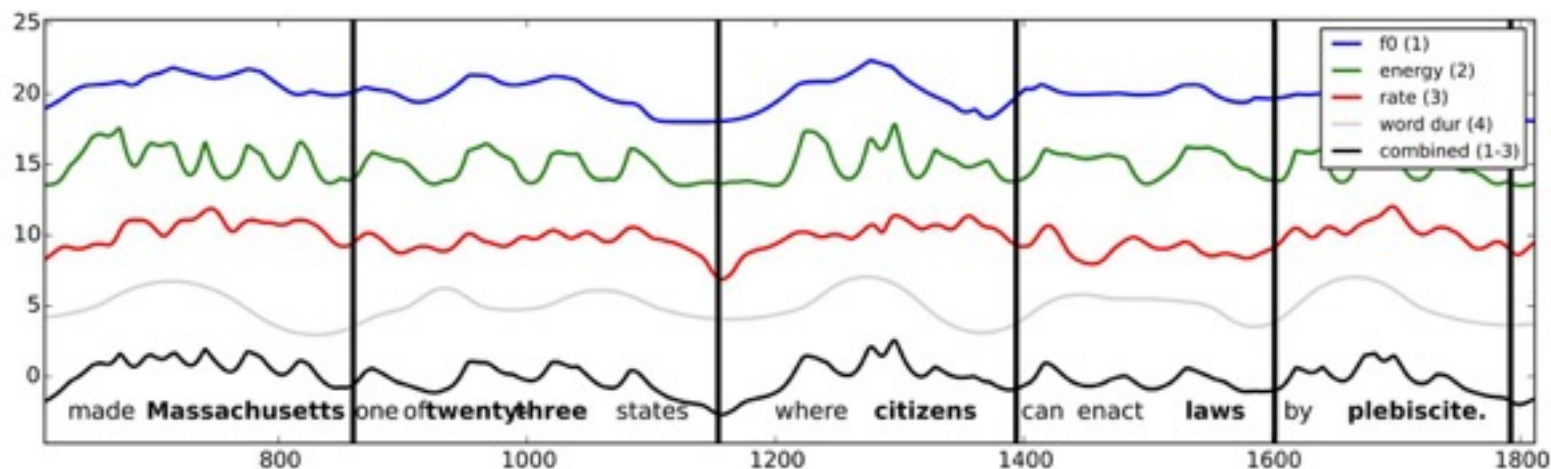
Speaking rate

- band-pass energy signal, Morlet wavelet
- amplitude scalogram (abs), normalize per frame
- follow ridge in time by viterbi

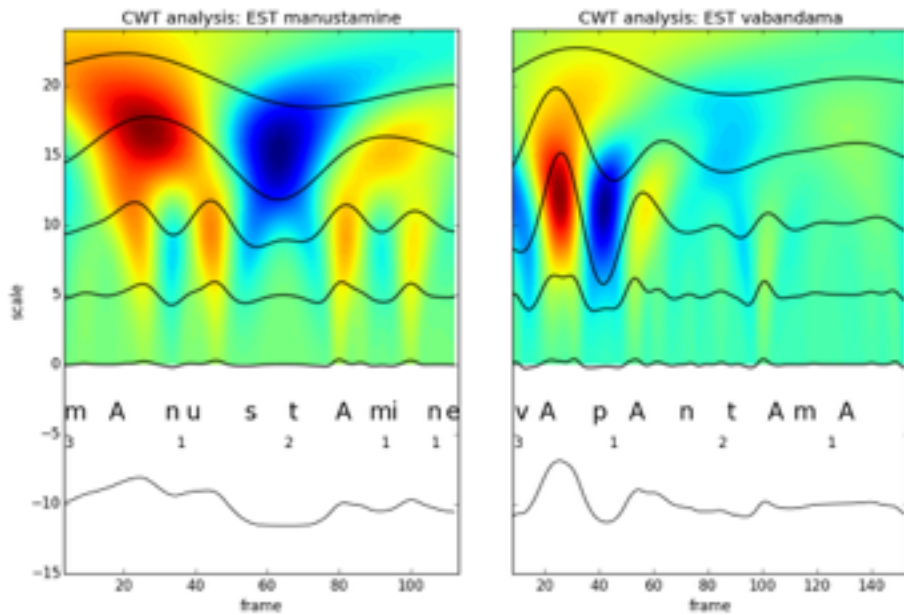


Different signals

- other speech signals can be processed by CWT:
 - (interpolated) f_0
 - (interpolated) energy (perhaps obtained via CWT)
 - duration signal (interpolated durations)
 - speaking rate (obtained via CWT)
- or even a combination thereof...

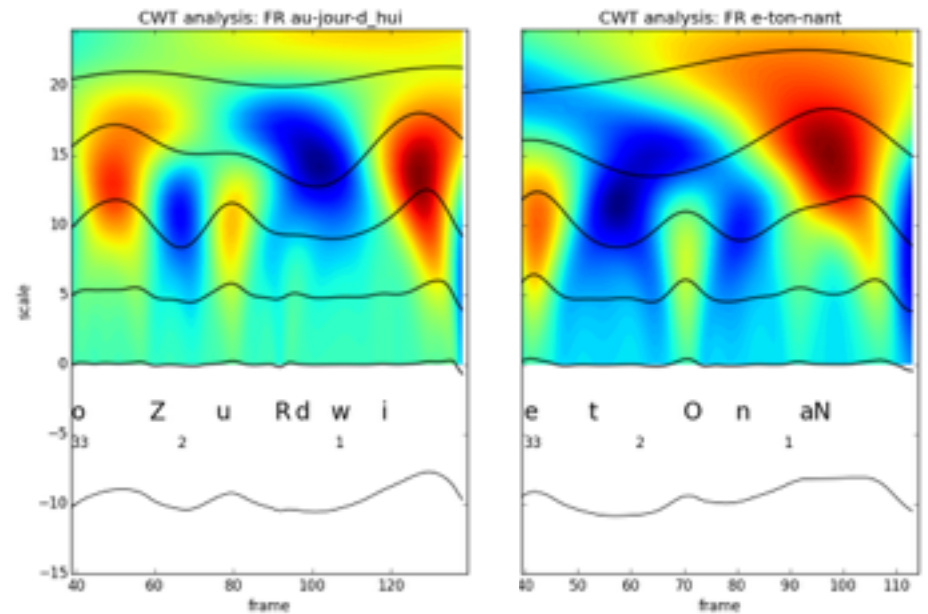


Prominence detection: lexical stress



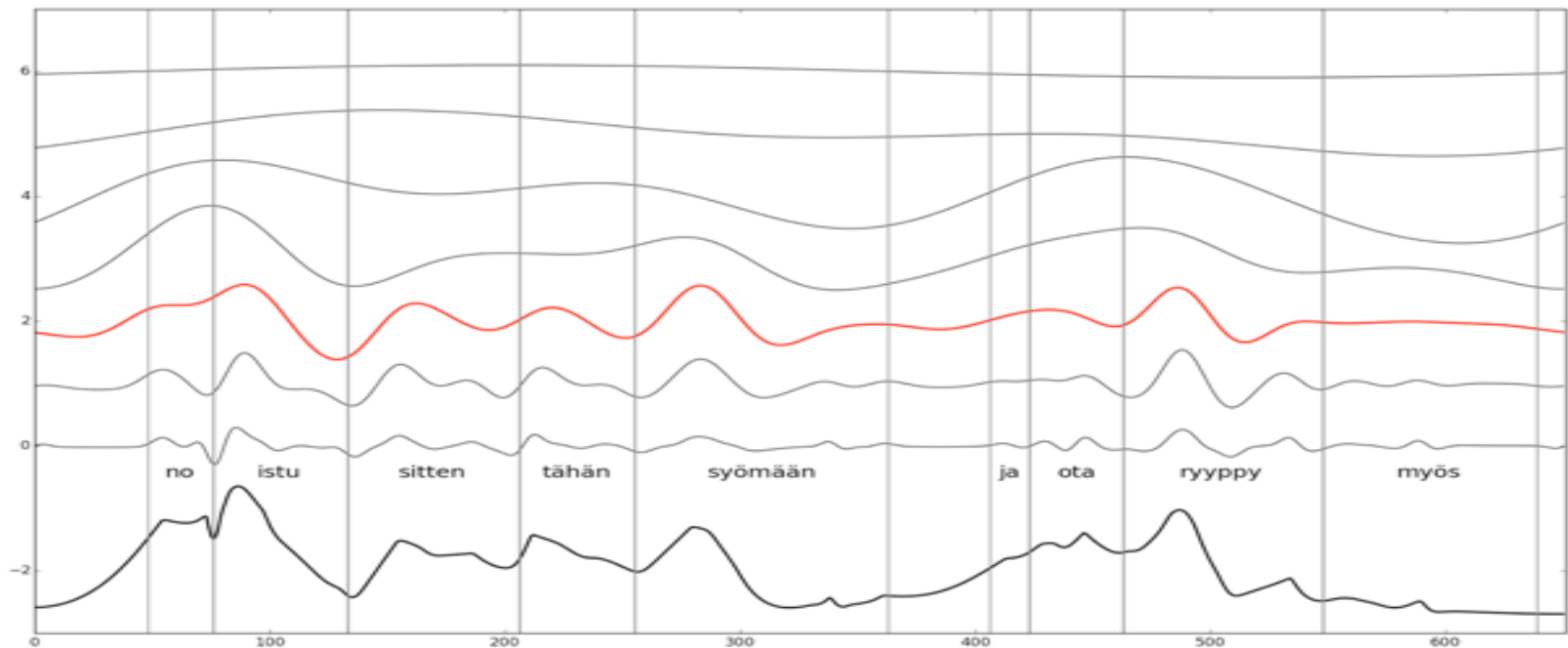
Estonian: word initial

French: word final

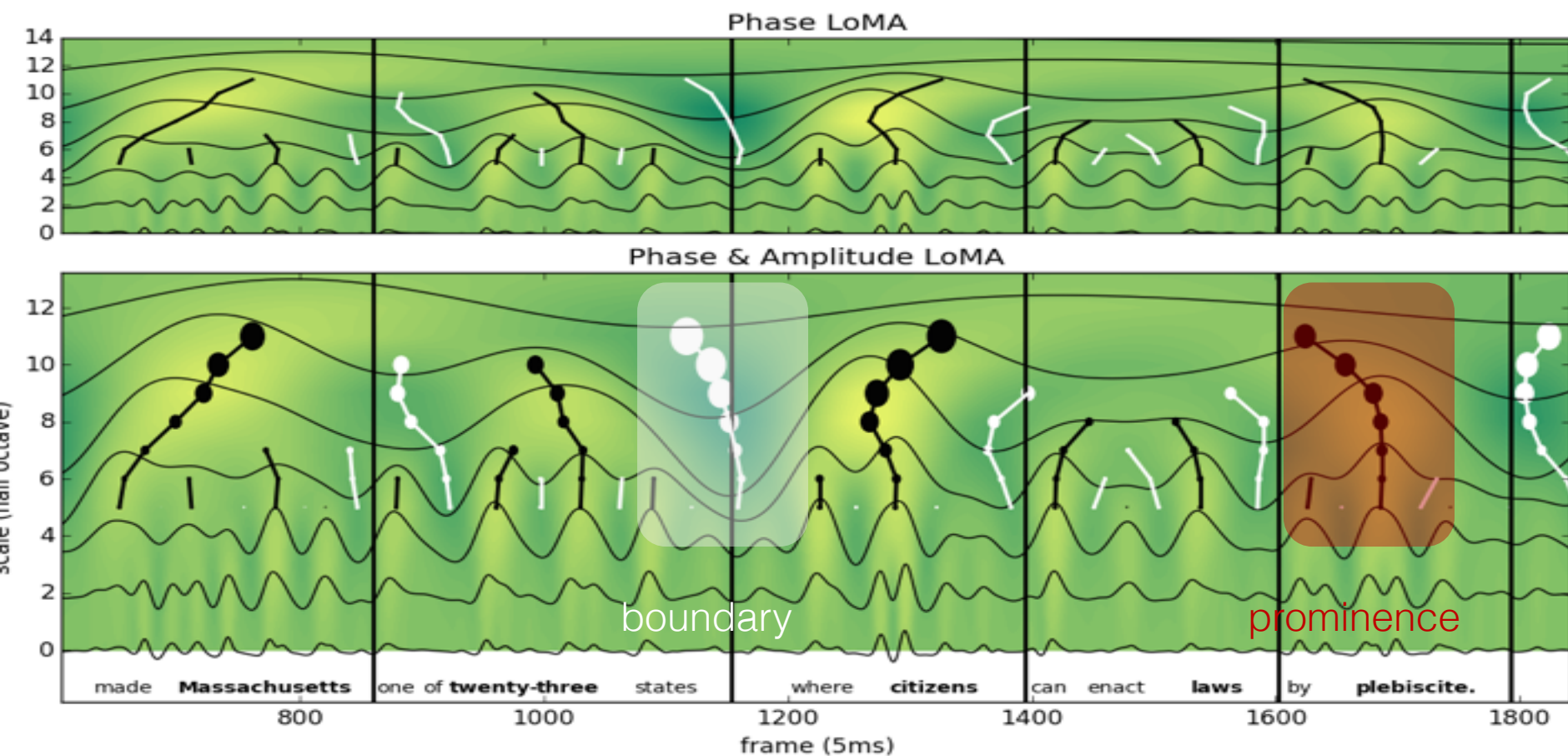


Prominence detection

1. Perform wavelet analysis on interpolated f_0
2. Select the scale with the closest match of number of peaks and number of words in the utterance
3. Prominence = the maximum peak of the word



Combining prominence and boundary detection

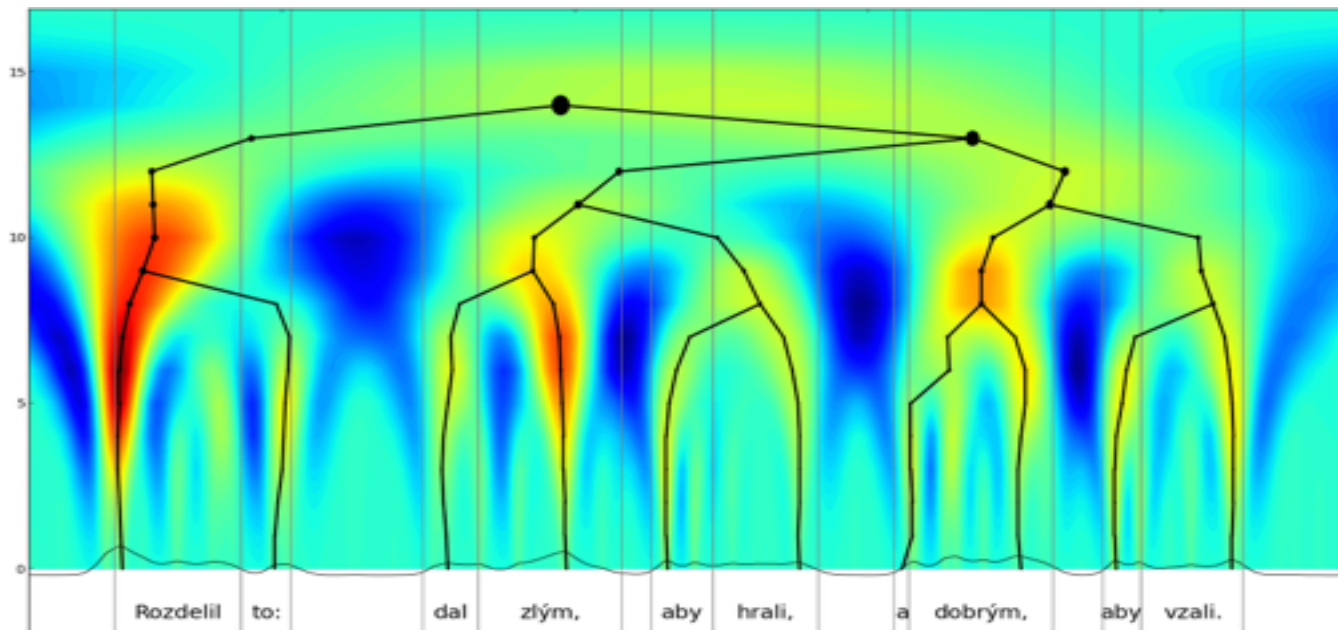


- Suni, Šimko, Aalto & Vainio (2016). Hierarchical representation and estimation of prosody using continuous wavelet transform. *Computer Speech & Language*
- Suni, Šimko & Vainio (2016). Boundary detection using Continuous Wavelet Analysis. *Proc.Speech Prosody*, Boston
- Suni & Kocharov (soon).

Prosodic hierarchy

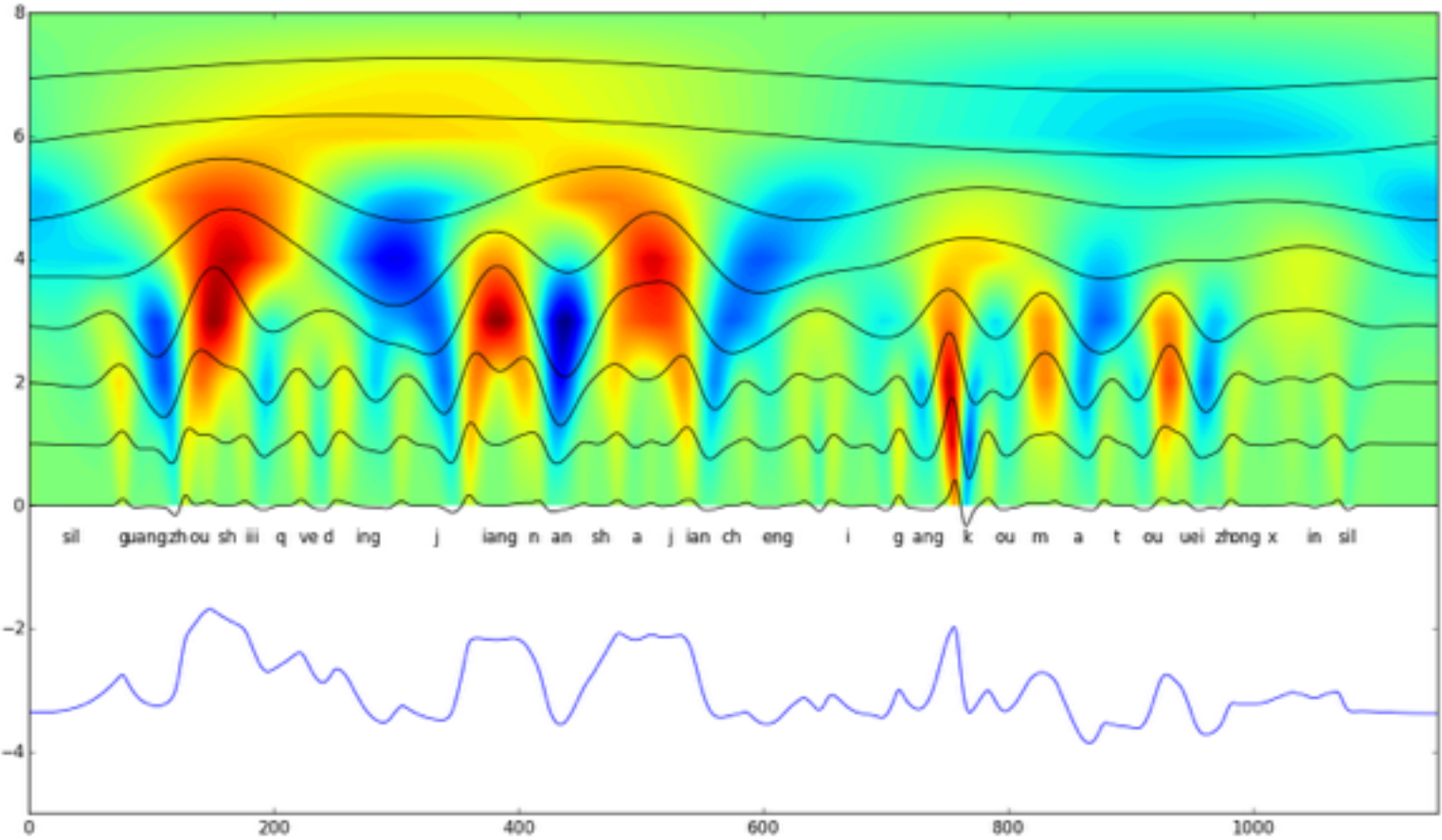
- Speech (and language) is structured hierarchically
- Utterances consist of (temporally) nested phrases > (phonological) words > syllables > speech sounds > acoustic events
- CWT allows for a time-frequency localisation of these structural components in a speech signal and can reveal how they are nested

[[[Rozdelil to:] [[[dal zlým] [aby hráli]] [a dobrým] [aby vzali.]]]]
[He] divided it: [he] gave [to] baddies to play and [to] goodies to take.

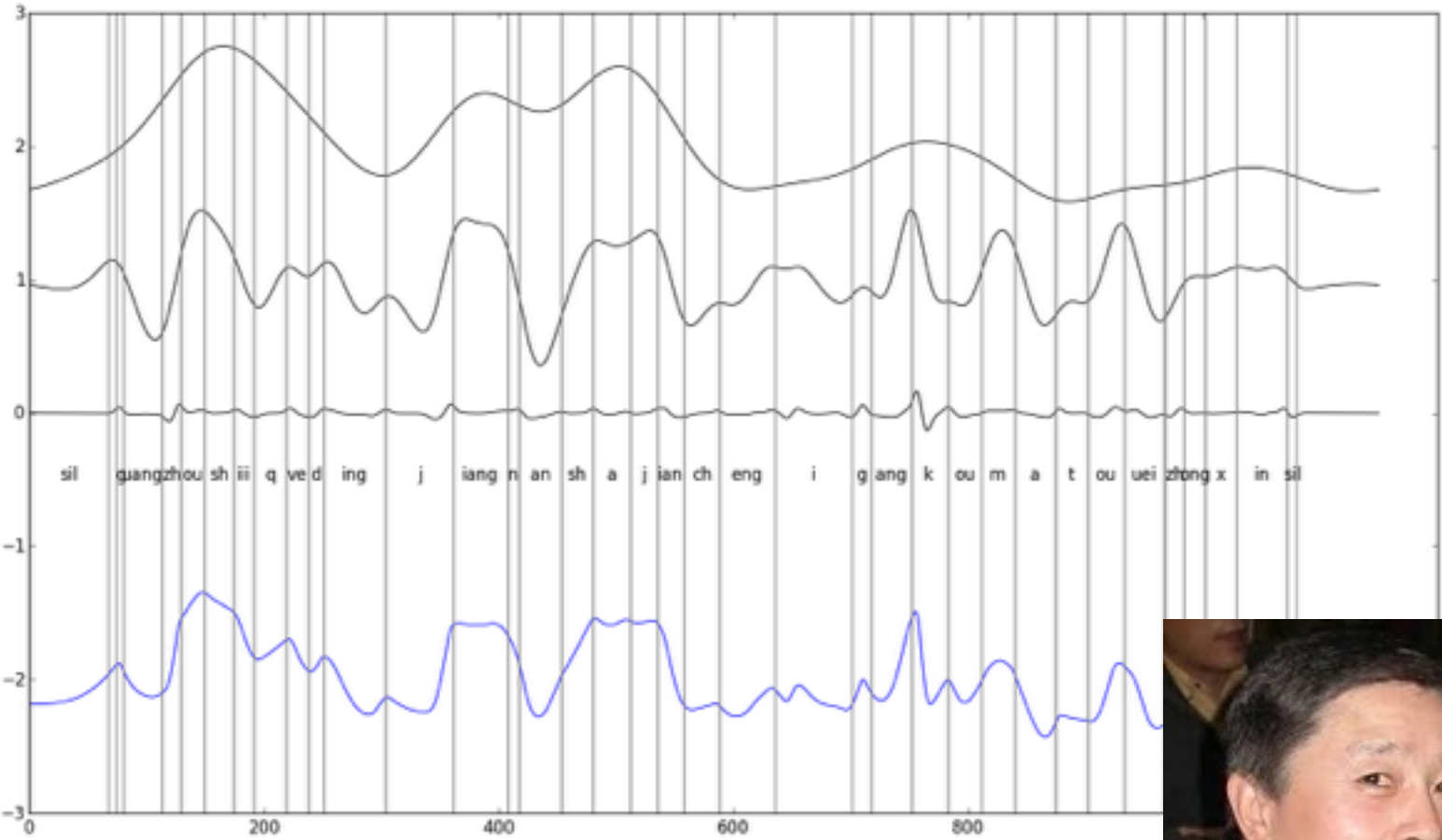


Dannenber, Vainio, Suni & Werner. (2015). Prosodic and syntactic segmentation of spontaneous speech: A preliminary study. In *Proceedings of ICPHS*, Glasgow

Prosodic hierarchy: The case of lexical tone



Prosodic hierarchy: The case of lexical tone





Thank you!



Thank you!