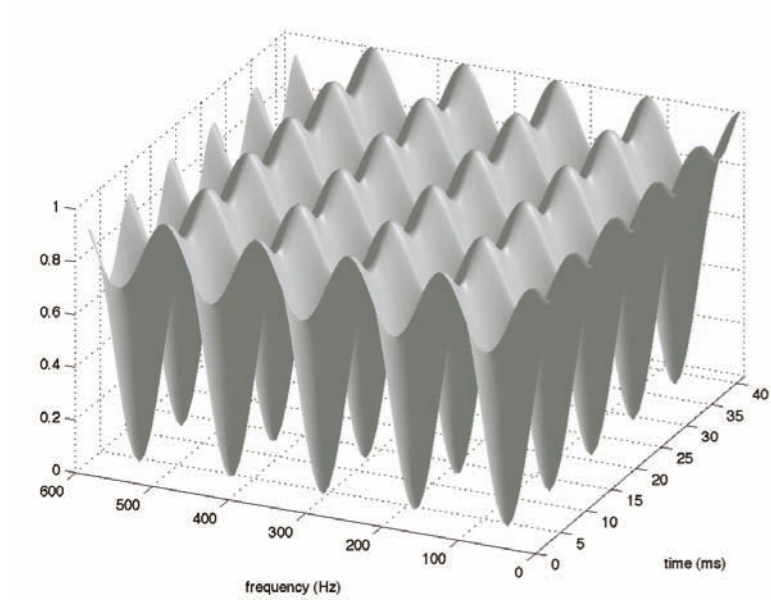
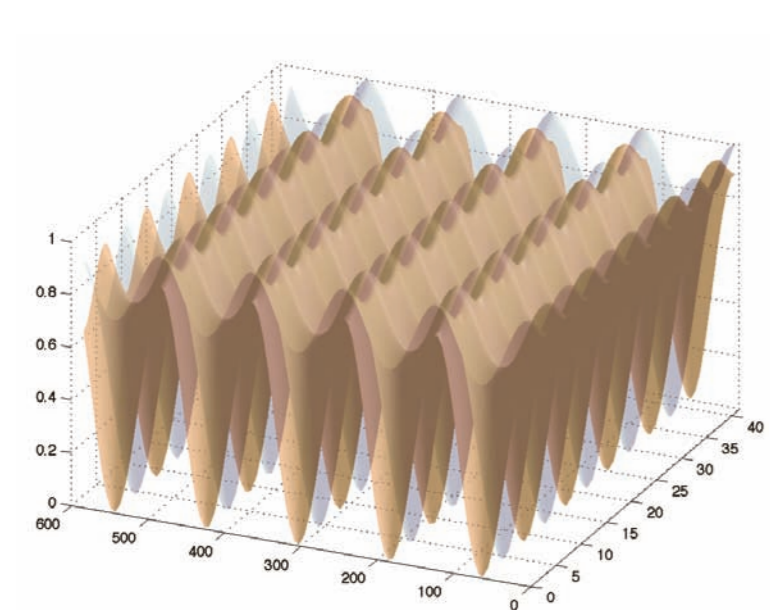


### Overview/Summary

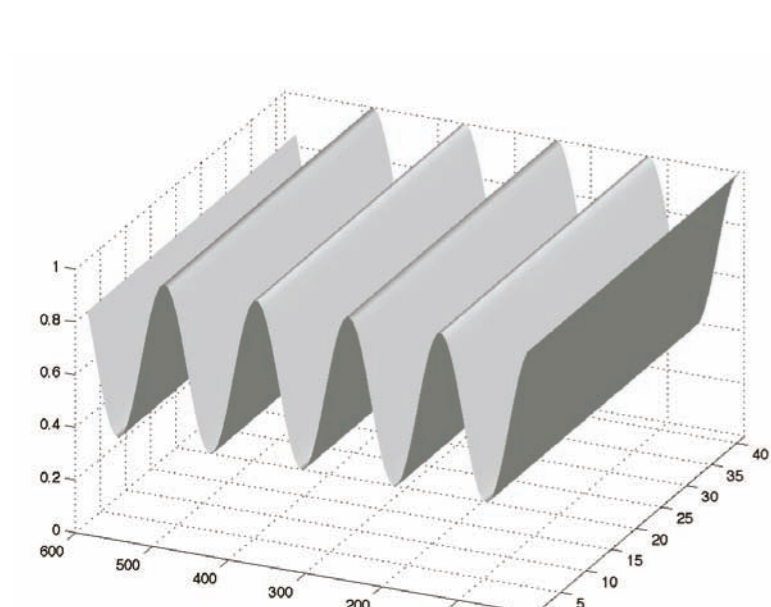
STFT power spectrum



T0/2 time shift + averaging

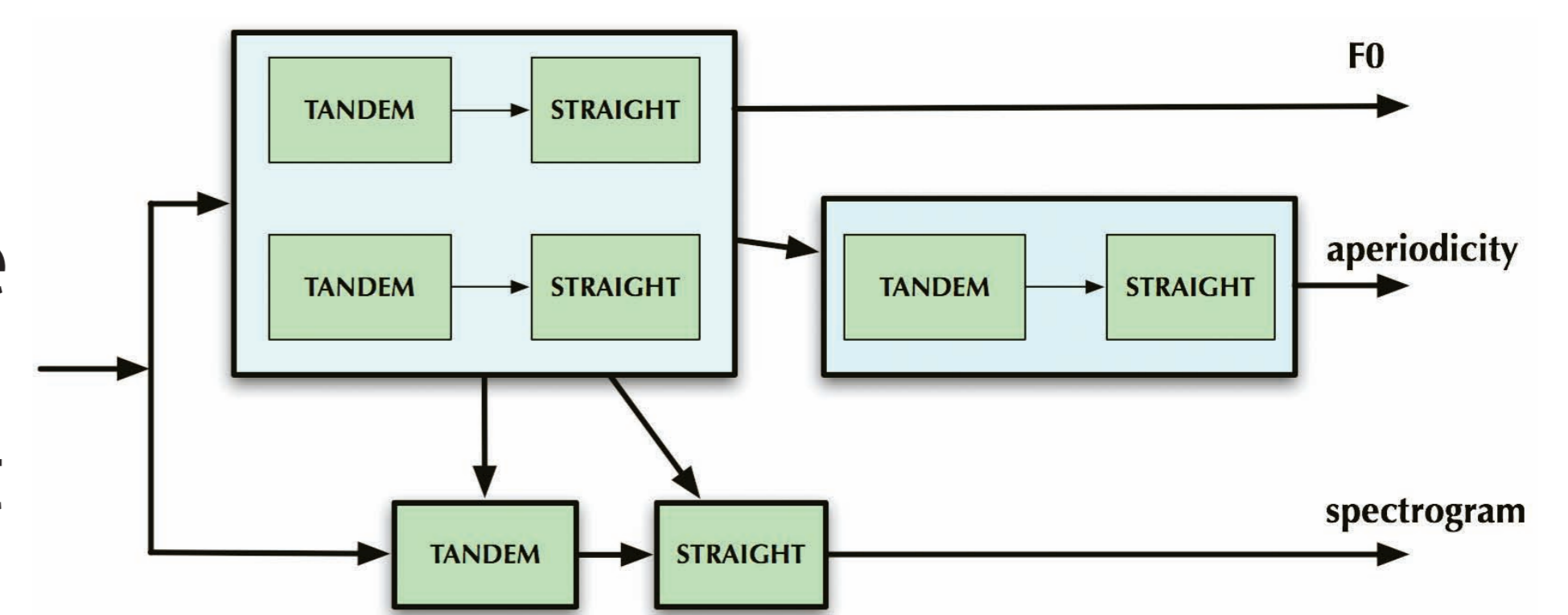


TANDEM spectrum



STRAIGHT, a speech analysis, modification synthesis system, is an extension of the classical channel VOCODER that exploits the advantages of progress in information processing technologies and a new conceptualization of the role of repetitive structures in speech sounds. STRAIGHT was designed to provide representation consistent with our auditory perception which decomposes input sounds in terms of excitation (source) and resonant (filter) characteristics. This architecture and sophisticated implementation made it the most flexible and high-quality speech manipulation system to date. However, underlying algorithms were not well formulated mathematically sound manner.

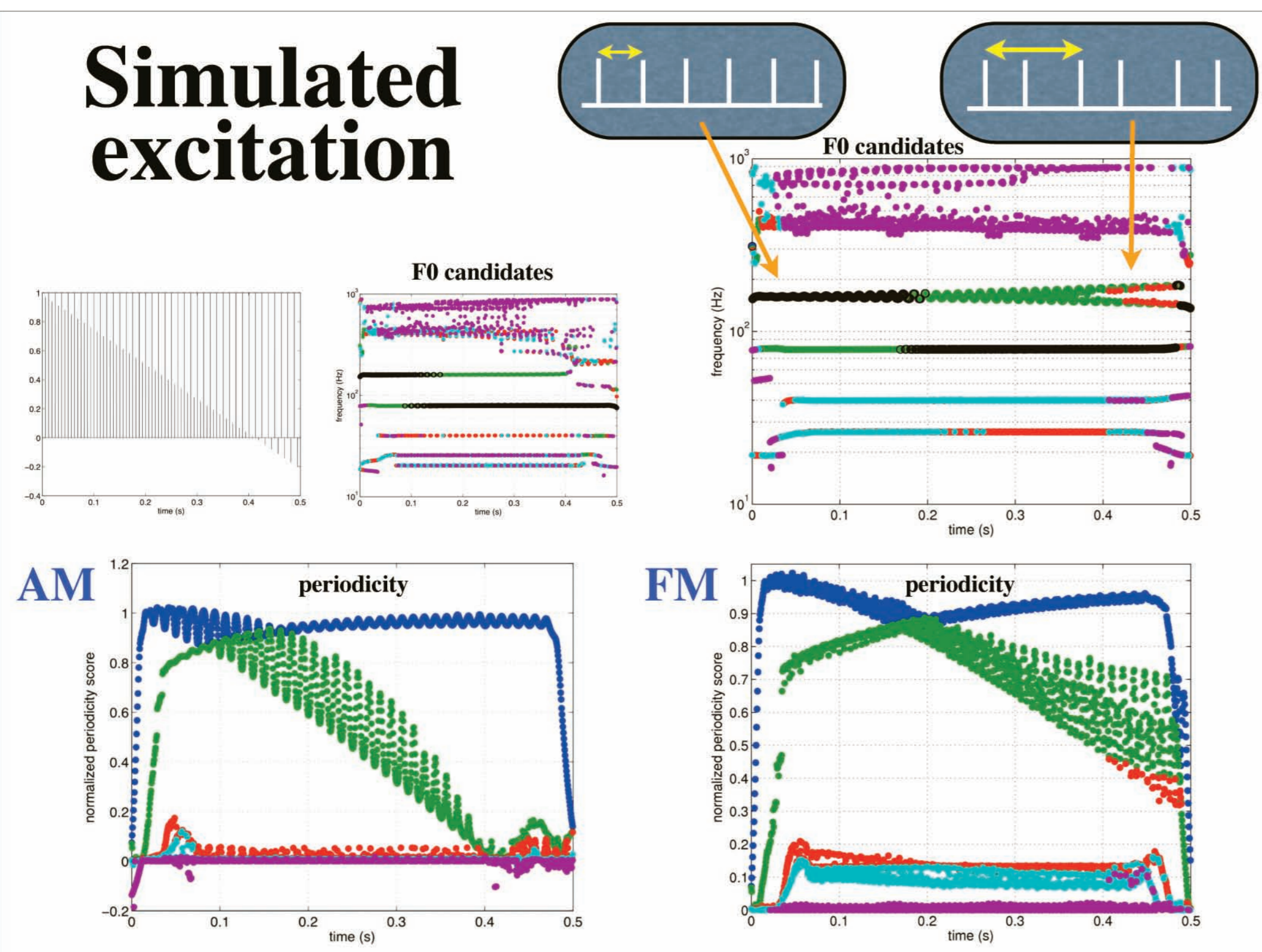
TANDEM, a method to extract temporally stable power spectral representation and consistent sampling theory, a new mathematical foundation of sampling made possible to reformulate STRAIGHT completely based on sound foundation. It also simplified codes and enhanced execution speed drastically.



TANDEM-STRAIGHT architecture

### Examples/Applications

#### Simulated excitation



TANDEM-STRAIGHT provides a unified method to extract the excitation source information with supra- and sub-harmonic structures. The baseline performance of the proposed method without any post processing is comparable to available state of the art methods.

Morphing based on STRAIGHT provides means to investigate physical correlates of perceptual attributes in singing. Preliminary tests indicated perceptual dominance of voice timbre in judgement of singers' identity. This finding and a concept of speech texture mapping are basis of a voice conversion solely based on vowel information.

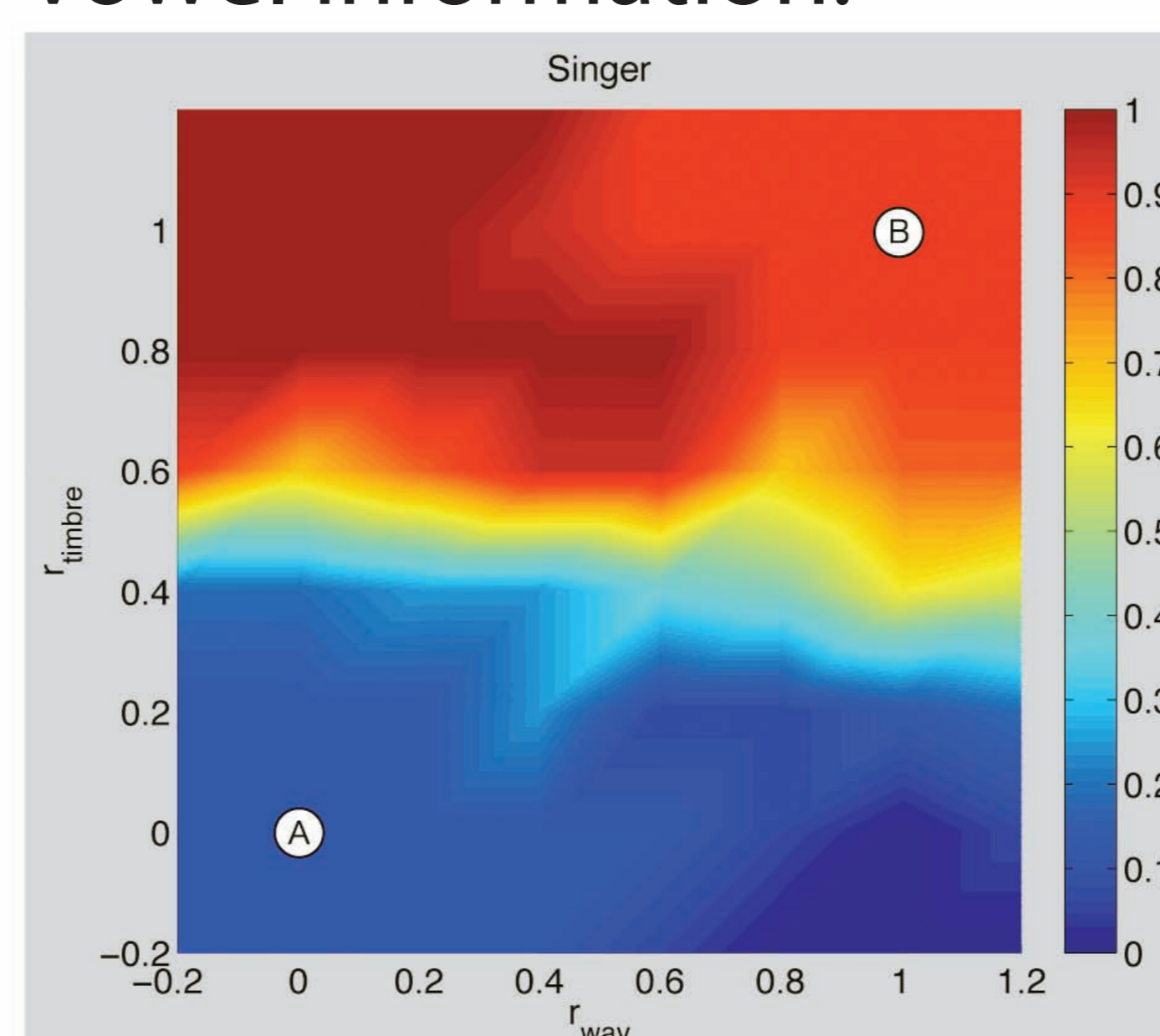
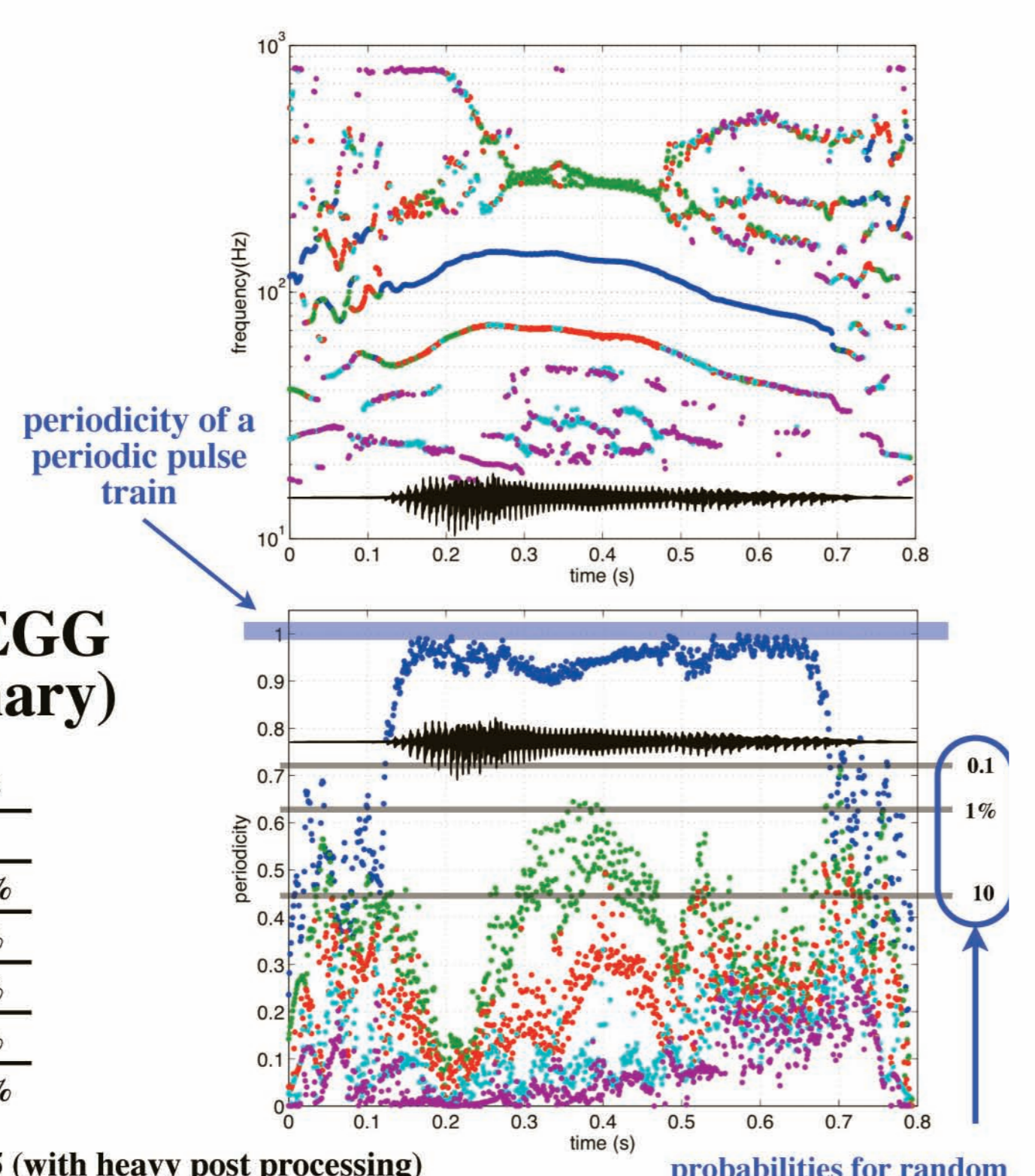
#### Natural speech

Analysis example  
Japanese vowel sequence /aiueo/ spoken by a male speaker

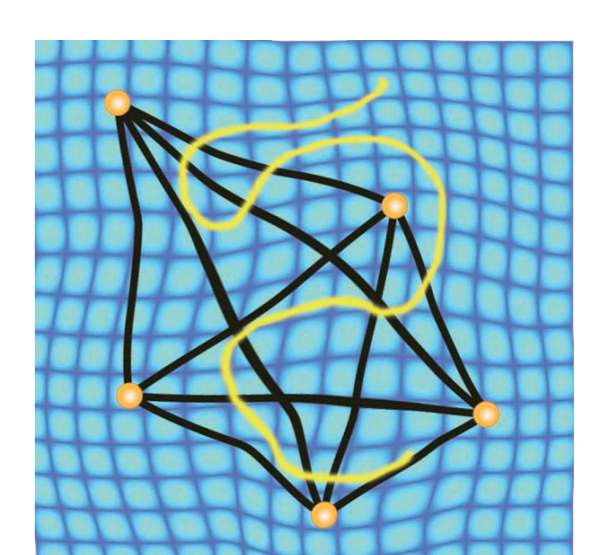
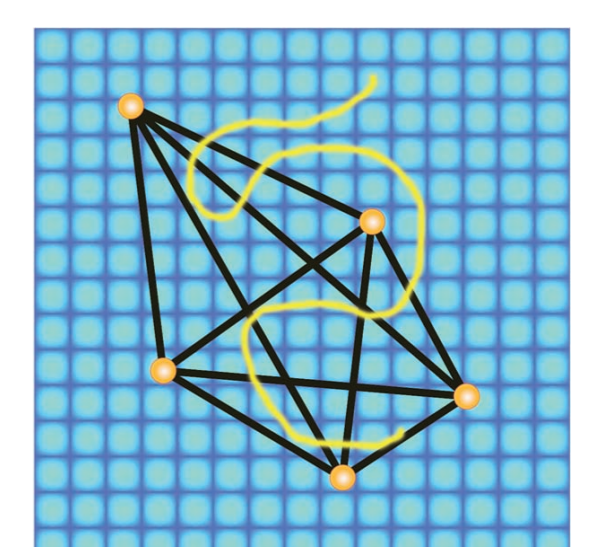
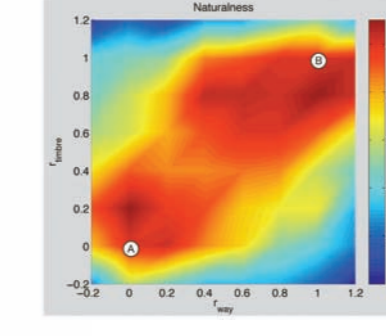
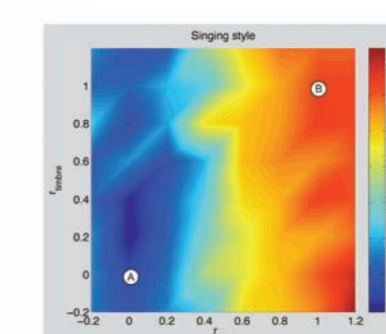
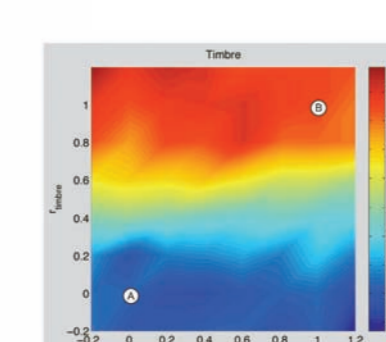
Evaluation using EGG database (preliminary)

	DB1	DB2
Proposed	0.30%	-
NDF*	0.09%	0.35%
YIN	0.29%	1.4%
TEMPO	0.77%	2.8%
ac	2.7%	4.5%
fxcep	7.3%	12.5%

\* Kawahara et.al. Interspeech 2005 (with heavy post processing)



Partial morphing for singing style, timbre and identity



Vowel-based voice conversion