

Musical Analysis
by Computer
from
Audio Recordings

HANS FUGAL
PHD COMPREHENSIVE EXAM
COMPUTER SCIENCE DEPARTMENT
NEW MEXICO STATE UNIVERSITY

What? **analysis from recordings**

Why? Have the computer do what we don't want to do, or can't easily do, or do it faster. Also, to learn about our own perception of music.



CRASH COURSE

WHAT IS MUSIC?

RHYTHM

PITCH

TIMBRE

3 axes: time, frequency, texture
mention dynamics in passing

aveverum.rg - Segment Track #1 - Matrix - Rosegarden

File Edit View Composition Segment Adjust Tools Settings Help

Grid: Beat Velocity: 100 Quantize: 1/16 100%

60
4/4

1 2 3 4 5 6 7

Time: 003-02-00-00 (9.000s) D#4 (75) Click and drag to select; middle-click and drag to draw new note

two of the dimensions (time, pitch)
timbre: instrument/performer

TANGERINE

J. MERCER / 359
V. SCHERTZINGER

The musical score is written in 4/4 time and consists of ten staves. The notation includes various chord symbols and melodic lines. The chords are: Gmi, C7, F, Bb, Ami, D7(b9), Gmi, C7, Gmi, C7, F, Aφ, D7(#9), Gmi, C7, F, Bφ, E7(#9), A, Bmi, E7, A7, D7(b9), Gmi, C7, F, Bb, Ami, D7(b9), Gmi, C7, Gmi, C7, Eb7, D7(#9), Gmi, Eφ, A7(#9), Dmi.

Discuss essentials of notation: rhythm (durations, time signature), pitch (notes on staff). Discuss harmony and improvisation. Play the first two lines straight with block chords on the keyboard if possible (even better if you can transpose it down a step to match Brubeck), then play Brubeck clip.

J.S. Bach
Tocatta and Fugue in D Minor
BWV 565

Adagio



Prestissimo



Explicit notation, but still very free-flowing.

Play Herrick clip (first 3 bars).

note timing

Demonstrate timbre by playing the synthesizer (Dorsey) version of the first 3 bars, and compare with organ.

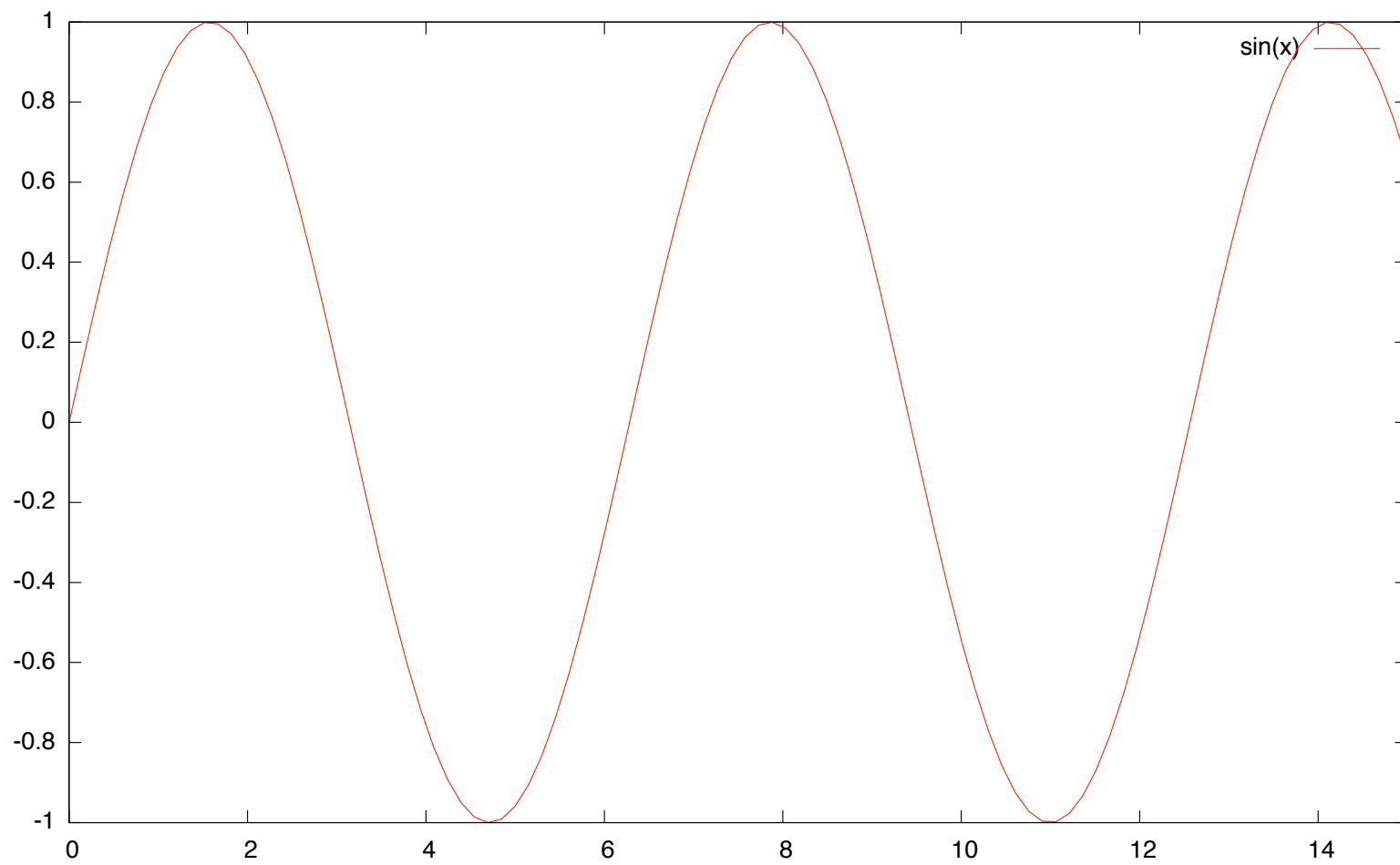
SOUND

VIBRATIONS CAUSE WAVES

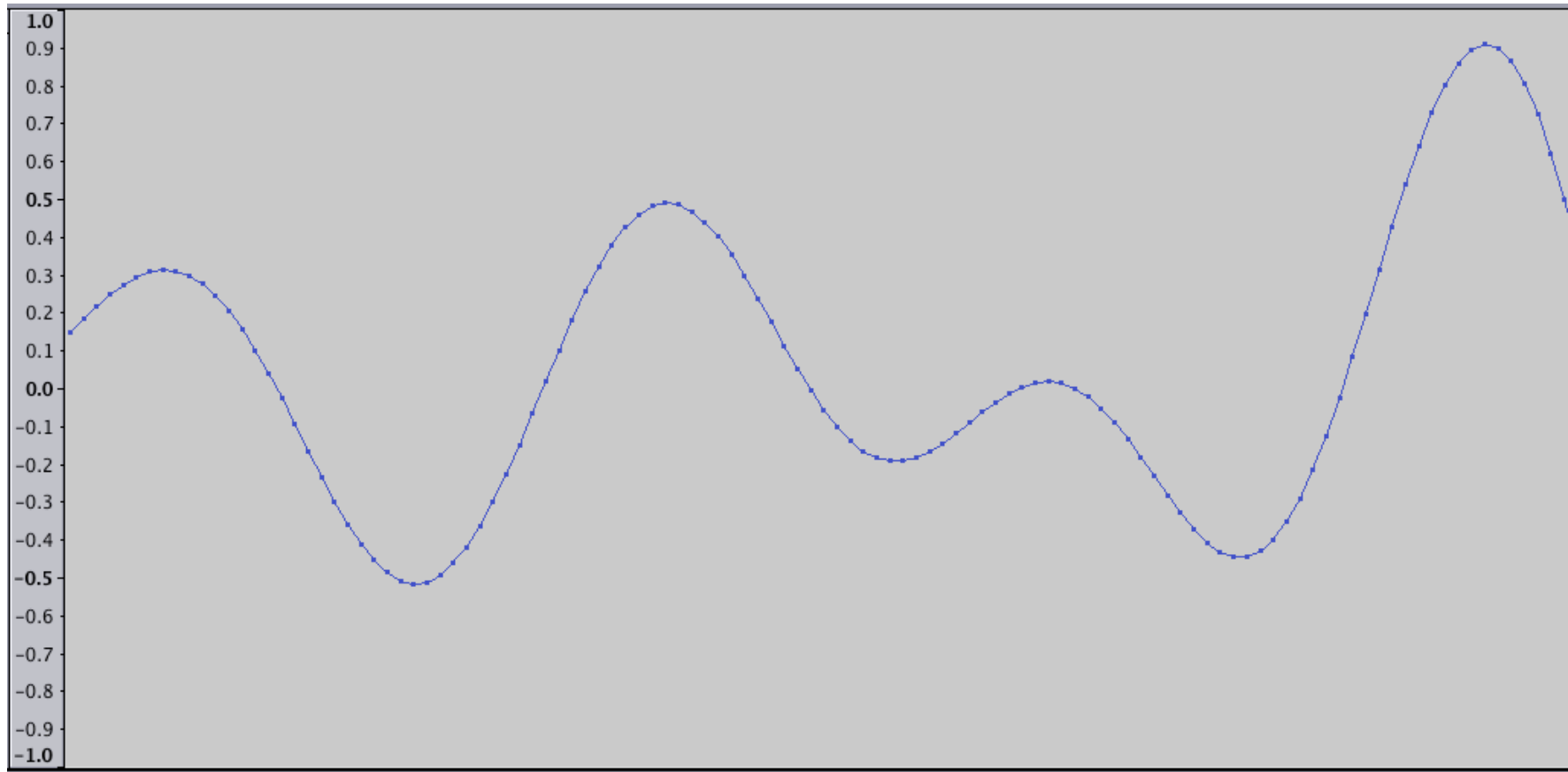
SINUSOIDAL COMPOSITION

TIME DOMAIN AND FREQUENCY DOMAIN

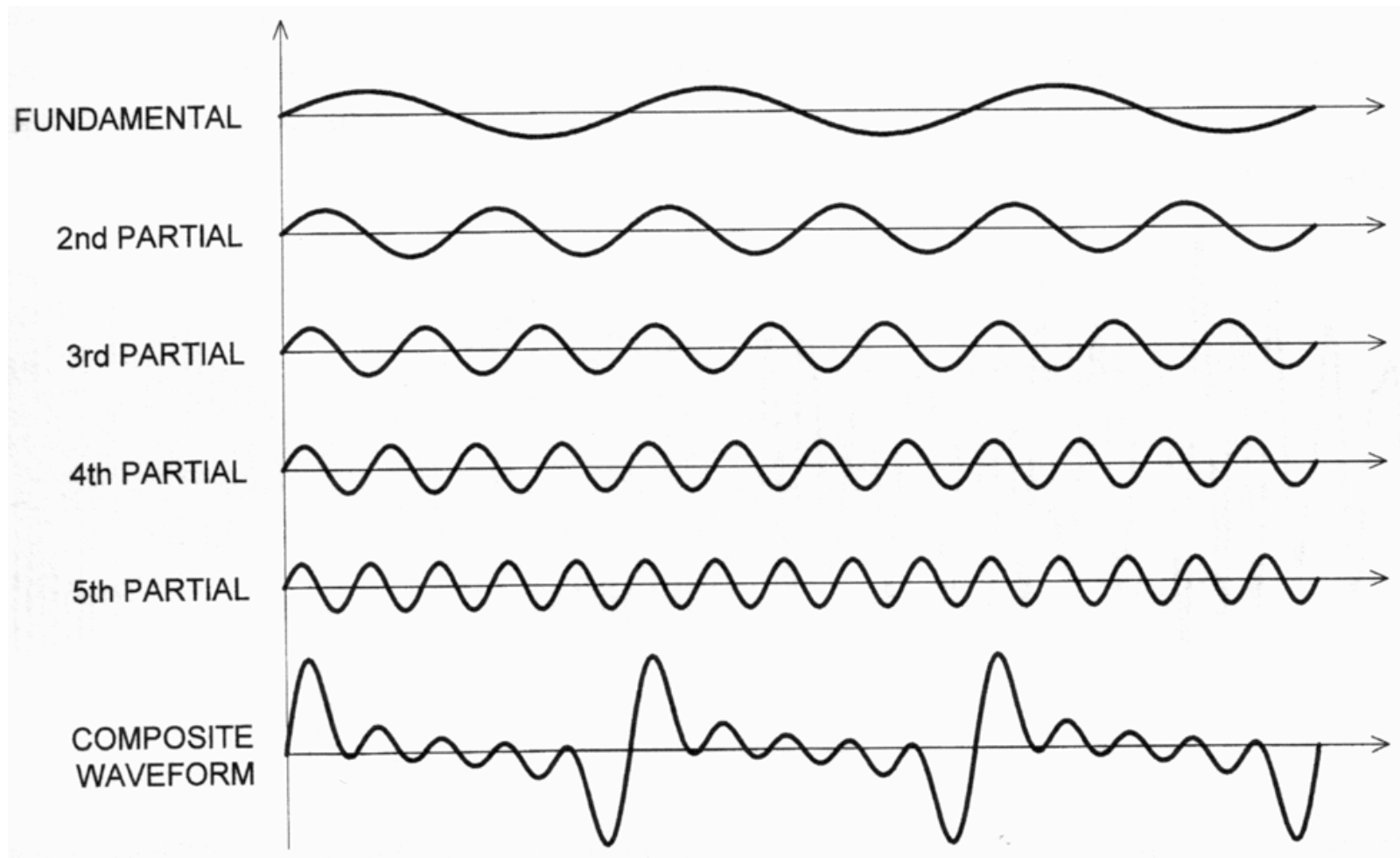
SOUND WAVES



SAMPLING



PARTIALS

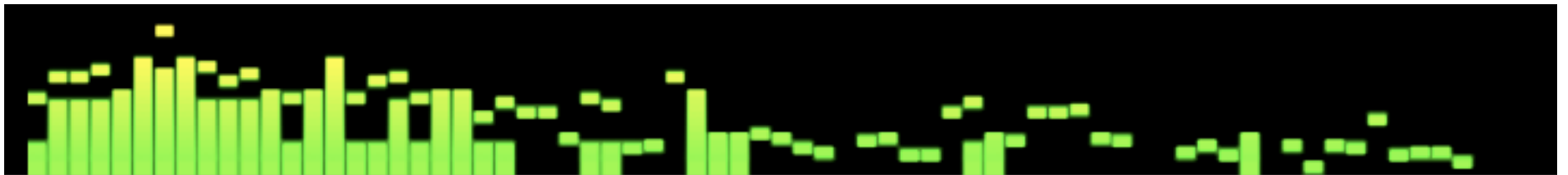


(Dodge and Jerse)

Compare with spectrogram

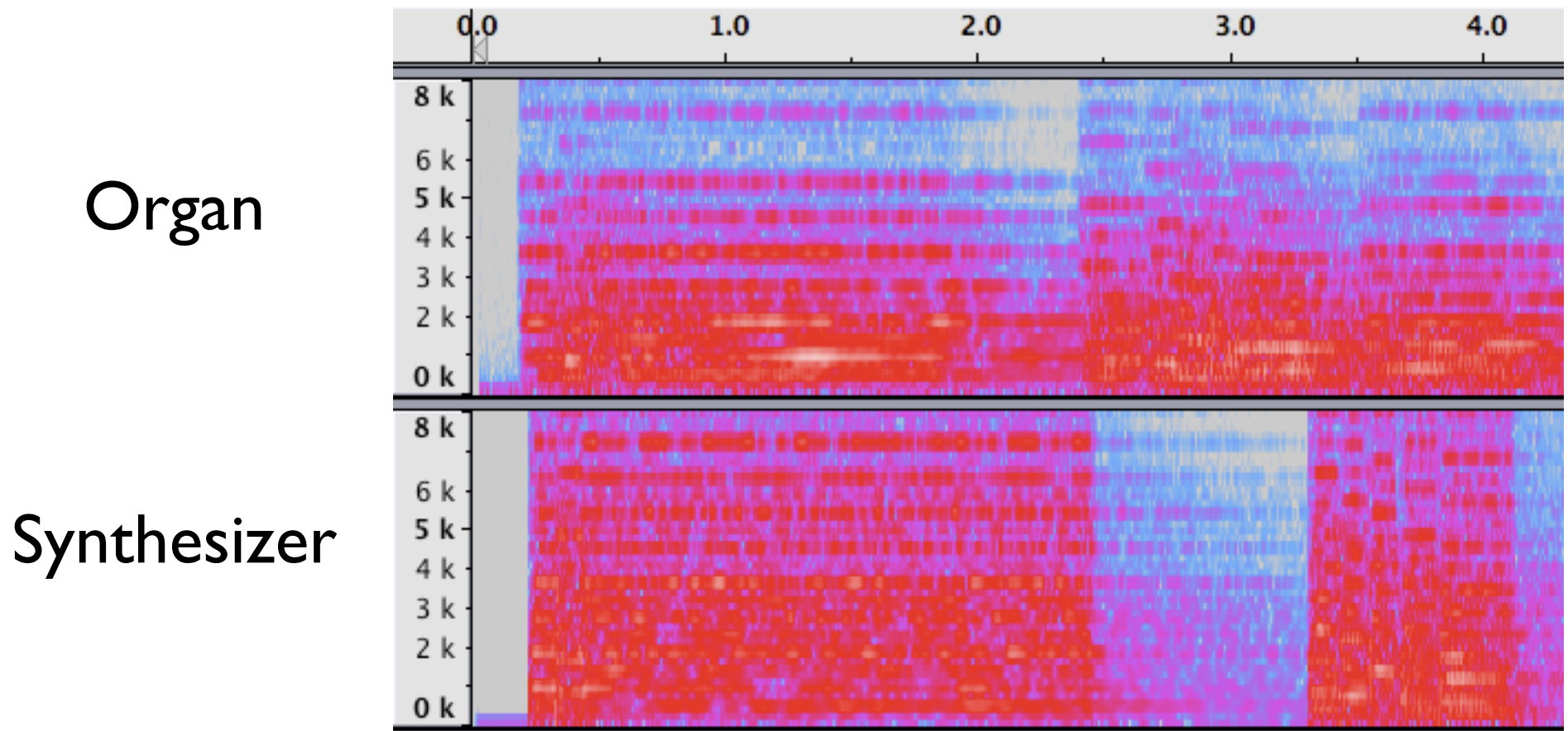
DISCRETE FOURIER TRANSFORM

$$X(\omega_k) = \sum_{n=0}^{N-1} x(t_n) e^{-j\omega_k t_n}$$



You've seen it before...
Demo Billy Joel, Movin' Out in VLC with Spectrogram.

SPECTROGRAMS



Herrick (top) and Dorsey (bottom), first few notes.

Rhythm

WHAT?

TEMPO

BEAT TRACKING

QUANTIZATION

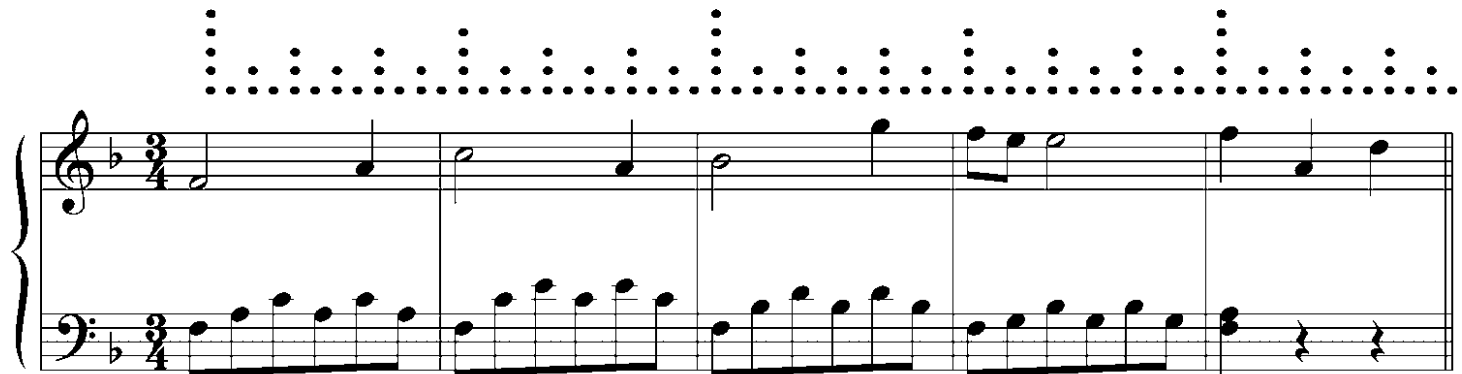
METRICAL STRUCTURE

BEAT AND TEMPO

BEAT: STEADY PULSE

TEMPO: BEATS PER MINUTE

METRICAL STRUCTURE



(1) (2) (3) (4) (5) (6)

(Temperley, 2004)

The image shows a musical score in 3/4 time. Above the staff is a dotted line with vertical stems, representing a metrical grid. The score consists of two staves: a treble clef staff and a bass clef staff. The treble staff contains a melody with notes on the first, second, and fourth beats of each measure. The bass staff contains a bass line with notes on the first, second, and fourth beats of each measure. The first six measures are numbered (1) through (6) below the staff. The notation is in 3/4 time, with a key signature of one flat (B-flat).

METRICAL LEVELS

TIME SIGNATURES

BARS, DOWNBEATS, ETC.

Note the hierarchical nature
Which is the primary metrical level? Not always clear-cut
Quantization is the alignment of events to this hierarchical structure

COMPLICATIONS

TEMPO RUBATO

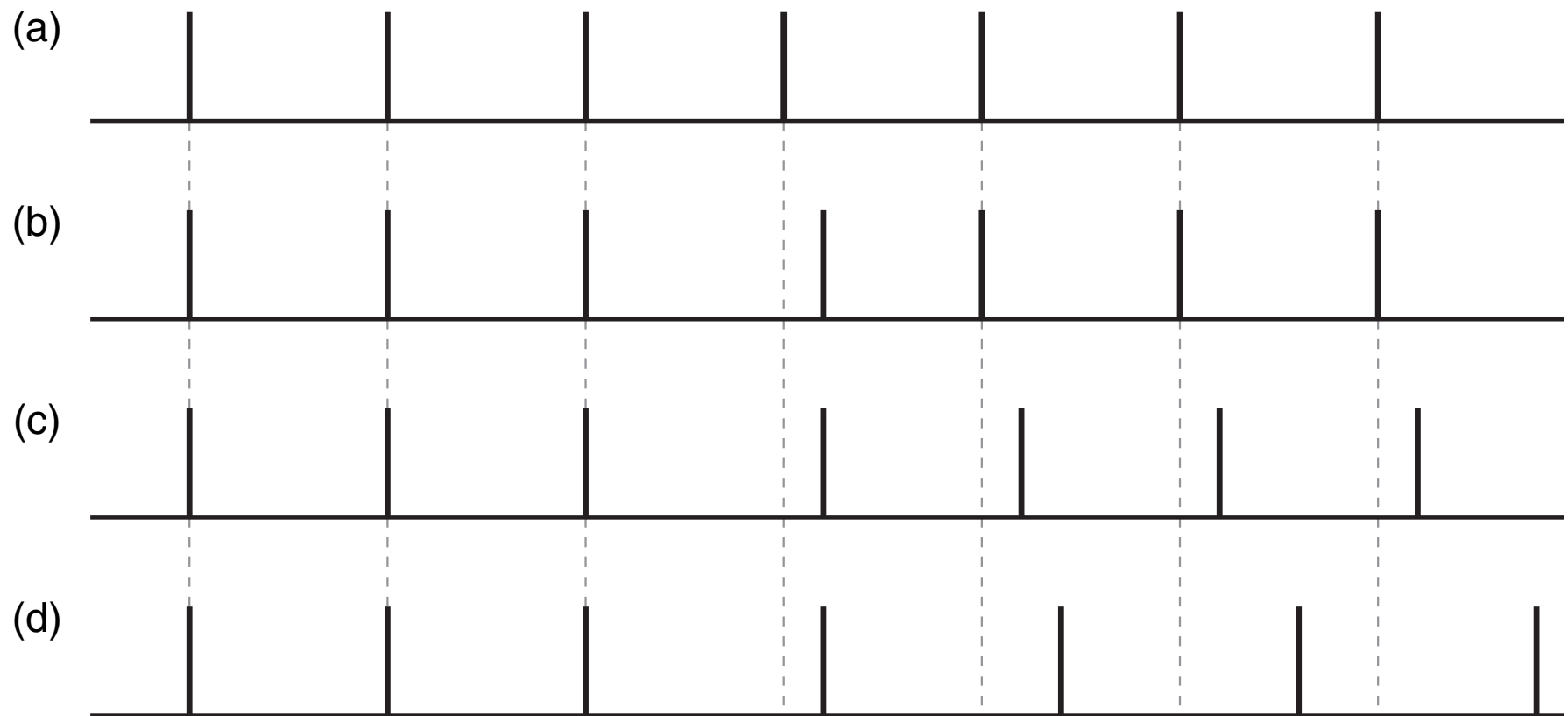
FERMATAS, ETC.

SYNCOPIATION

ARTICULATION

IMPERFECT PERFORMANCE

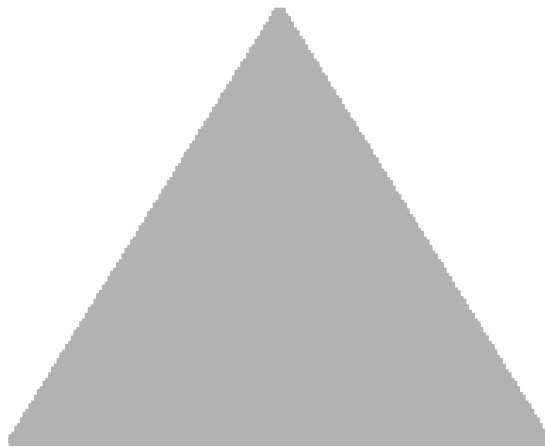
TEMPO DEVIATION



(Gouyon and Dixon 2005)

- a – no variation
- b – timing of one note
- c – offset (e.g. from a fermata)
- d – ritardando

Temporal structure



Tempo

Timing

(Honing 2001)

naïve algorithms trip up on timing
tempo changes too (rubato)
structure can be very valuable information

How?

DETECT ONSETS

CONSTRUCT A MODEL

the model has hidden state, i.e. the tempo and phase
it predicts and state is adjusted to maximize agreement

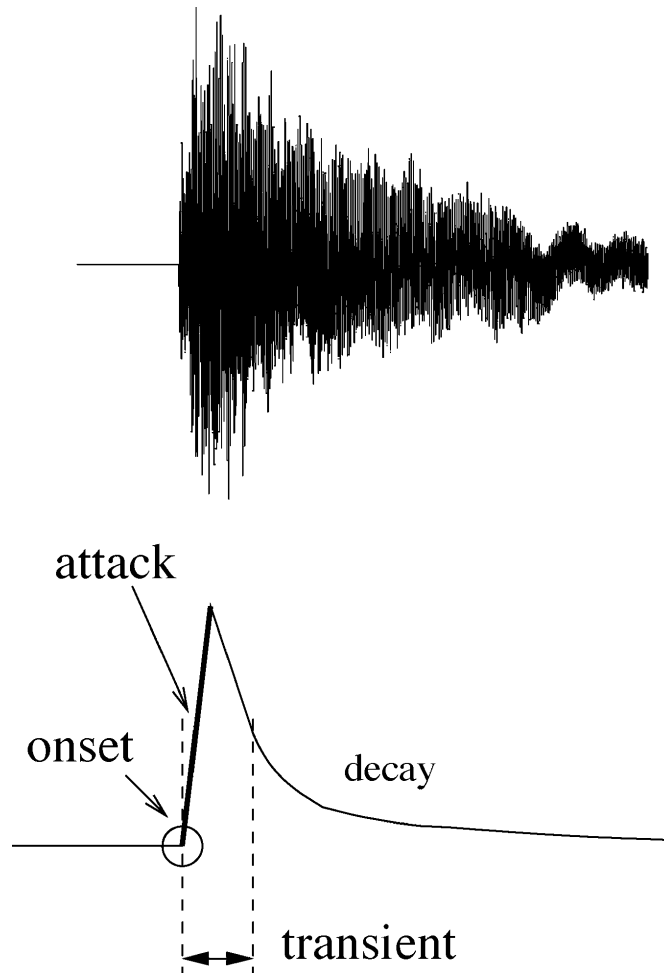
ONSETS

RHYTHMIC EVENTS

POINTS IN TIME

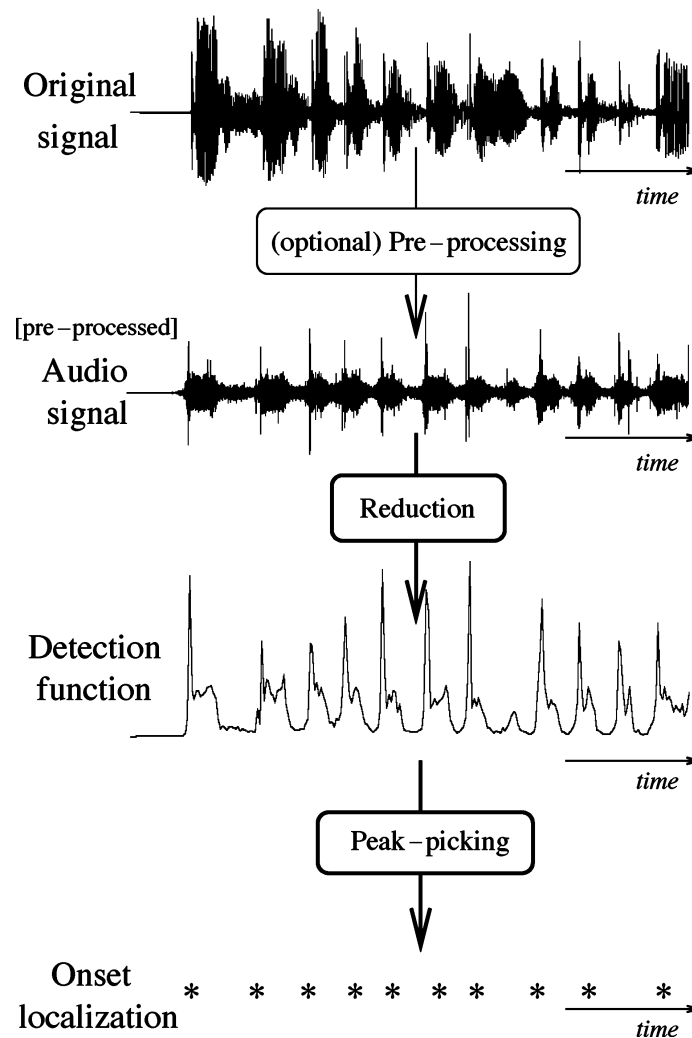
USUALLY ONSETS OF NOTES

ONSET DETECTION



(Bello et al. 2005)

ONSET DETECTION



(Bello et al. 2005)

preprocessing example: high-pass filter
detection functions aka novelty functions

DETECTON FUNCTIONS

AMPLITUDE - BEATS ARE LOUDER (OR ARE THEY?)

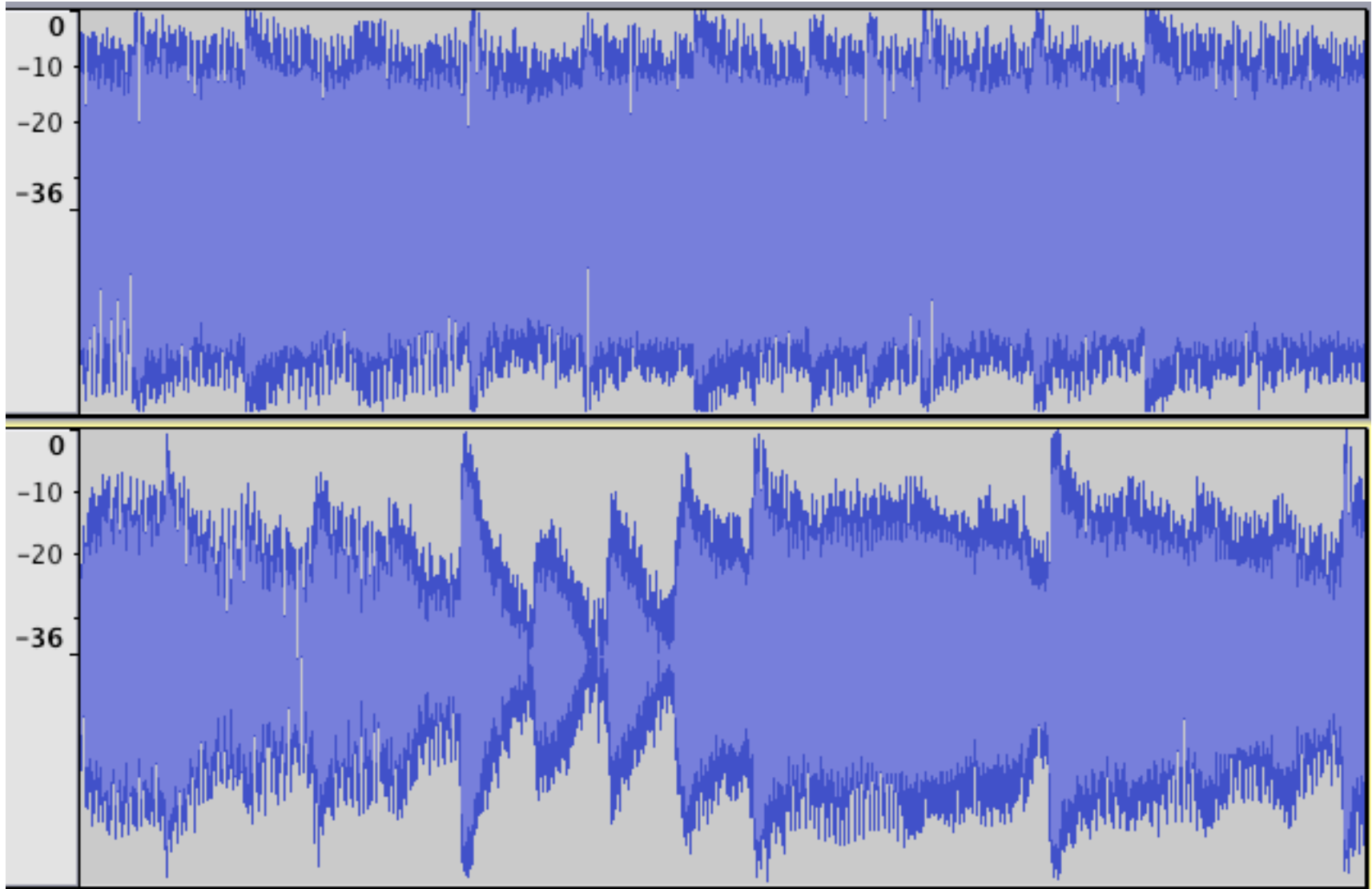
SPECTRAL FEATURES

LOOK FOR PERCUSSION IN HIGH BANDS

SPECTRAL FLUX

PROBABILISTIC MODELS

DYNAMIC RANGE COMPRESSION



5 seconds of audio
Top: Keane – Somewhere Only We Know
Bottom: Billy Joel – Movin' Out

PROBABILISTIC MODELS FOR ONSET DETECTION

MODEL THE SIGNAL

SWITCHING BETWEEN TRANSIENT/STEADY STATE

SURPRISE-DETECTING

TEMPO TRACKING

STOCHASTIC DYNAMICAL SYSTEM

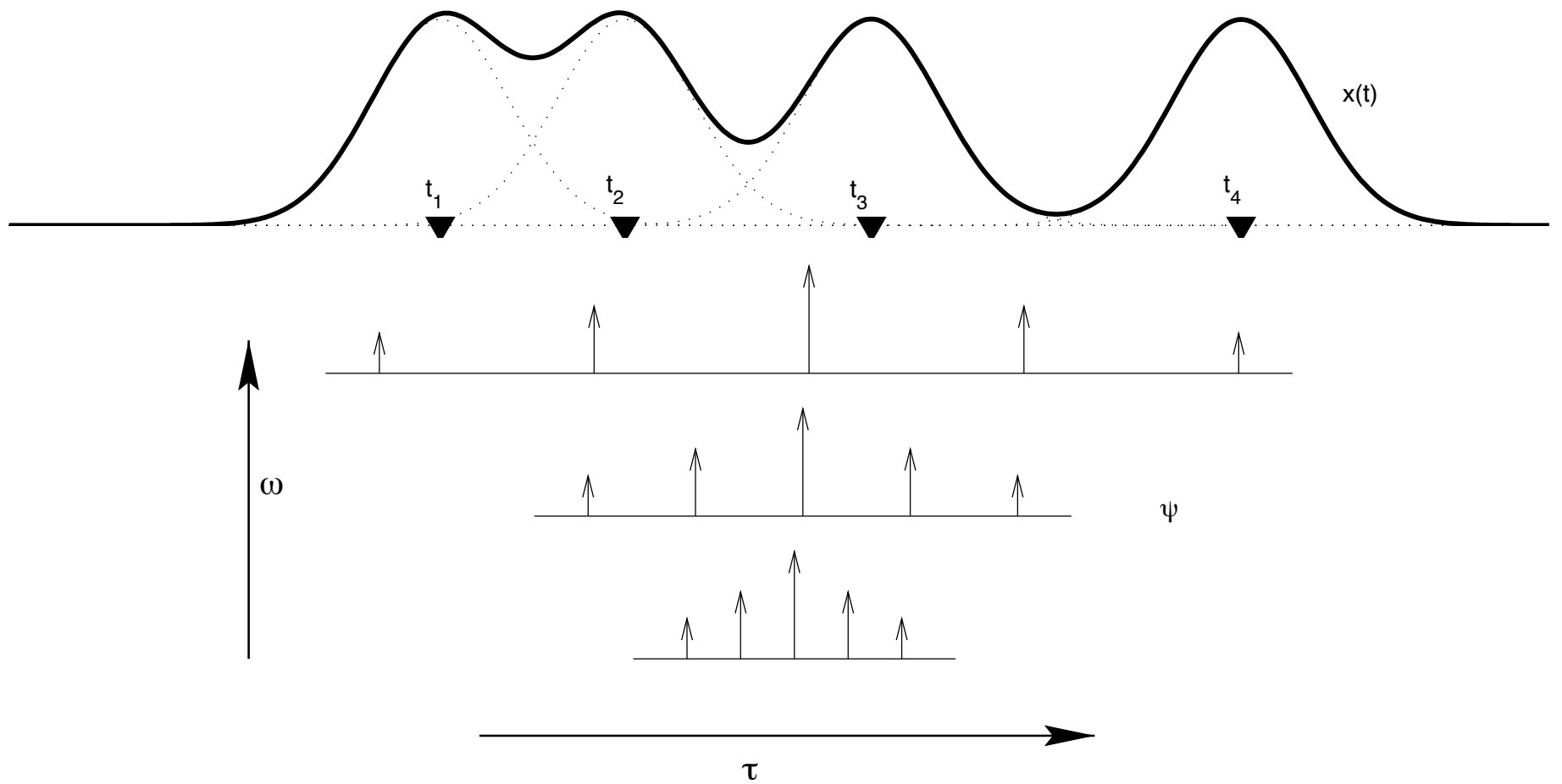
MODELS BEAT AND PERIOD - HIDDEN STATE

TEMPOGRAM

KALMAN FILTER

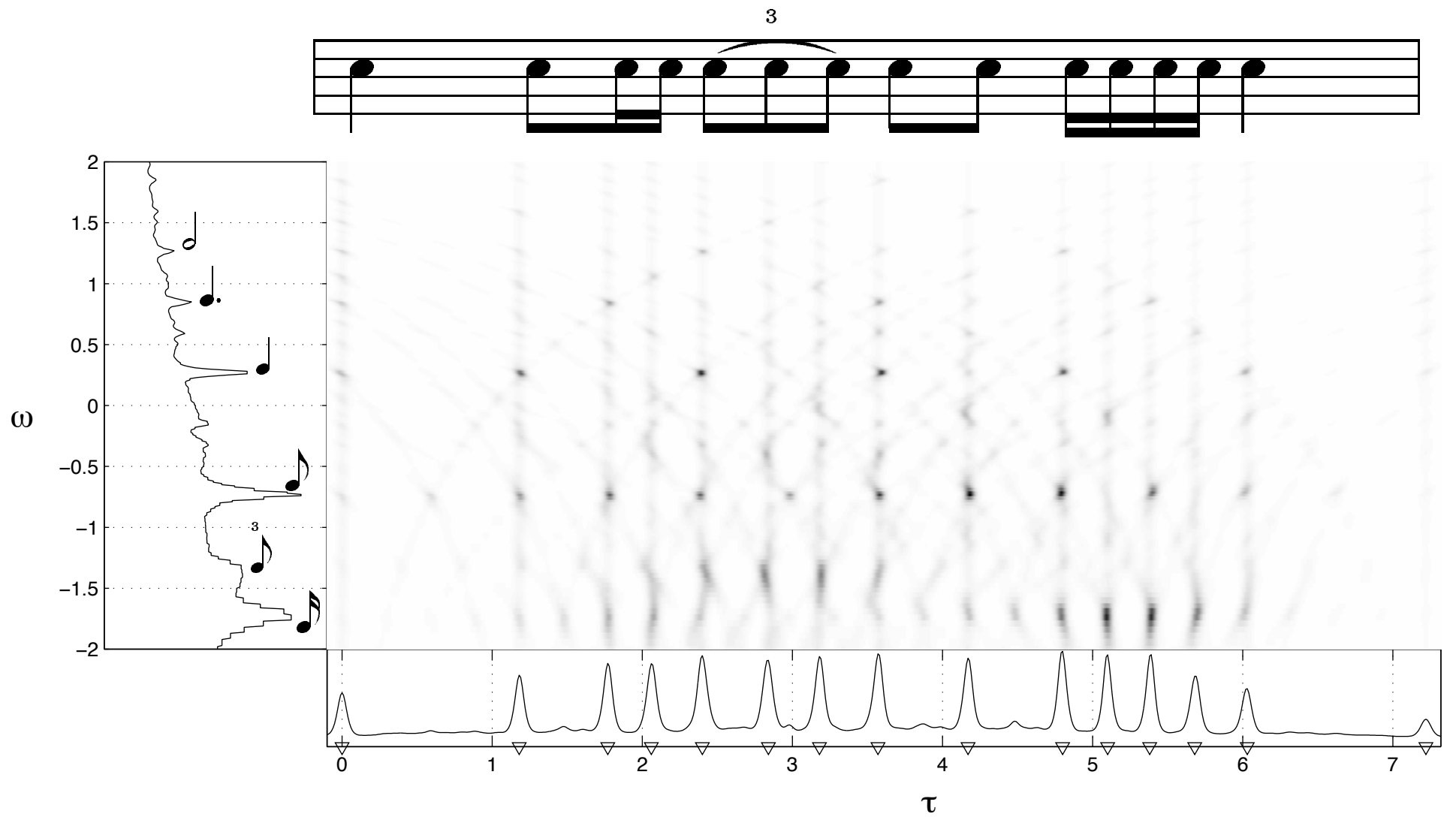
(Cemgil et al. 2001)

TEMPOGRAM



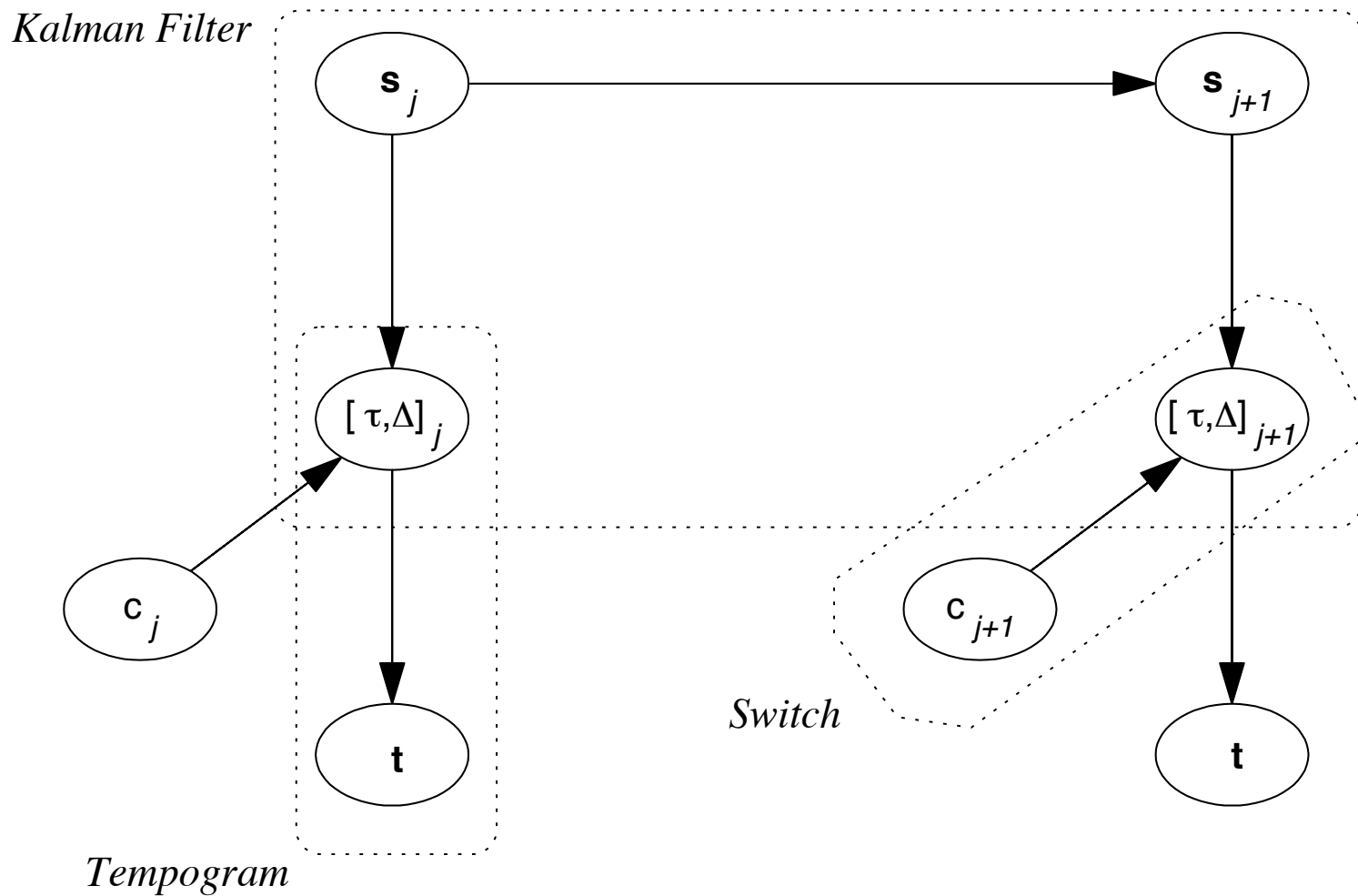
(Cemgil et al. 2001)

TEMPOGRAM



(Cemgil et al. 2001)

KALMAN FILTER



(Cemgil et al. 2001)

EXPRESSIVE PERFORMANCES

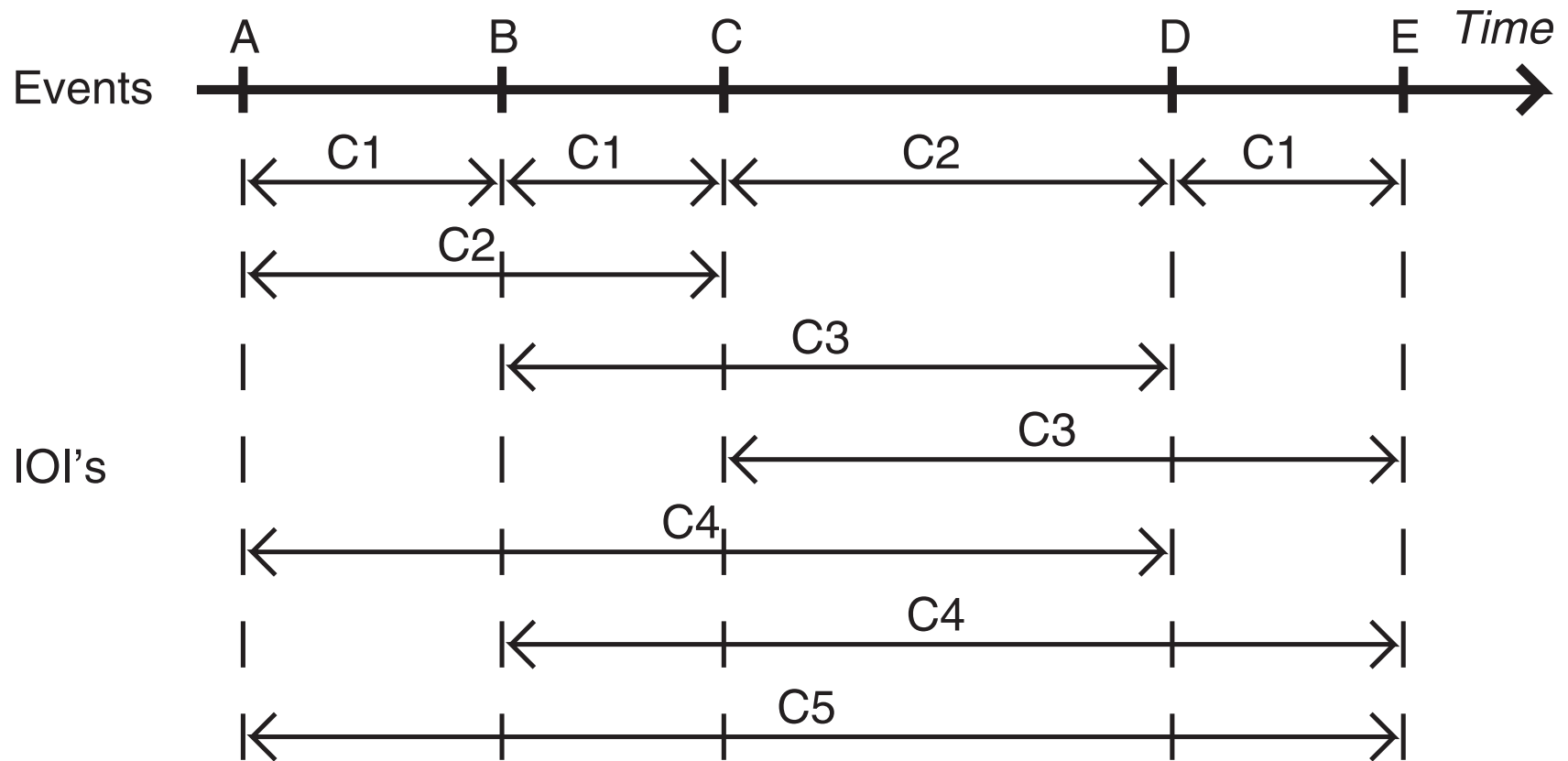
ONSET SALIENCE

ONSET CLUSTERING

MULTIPLE AGENTS

(Dixon 2001)

IOI CLUSTERING



(Dixon 2001)

Define IOI
Overlapping
Clusters C1, C2, C3, C4, C5 of similar IOI length

AGENTS

CLUSTERS BECOME HYPOTHESES

AGENTS CREATED FOR EACH HYPOTHESIS

PREDICT

EVALUATE

(Dixon 2001)

BEAT TRACKING WITH AND WITHOUT DRUMS

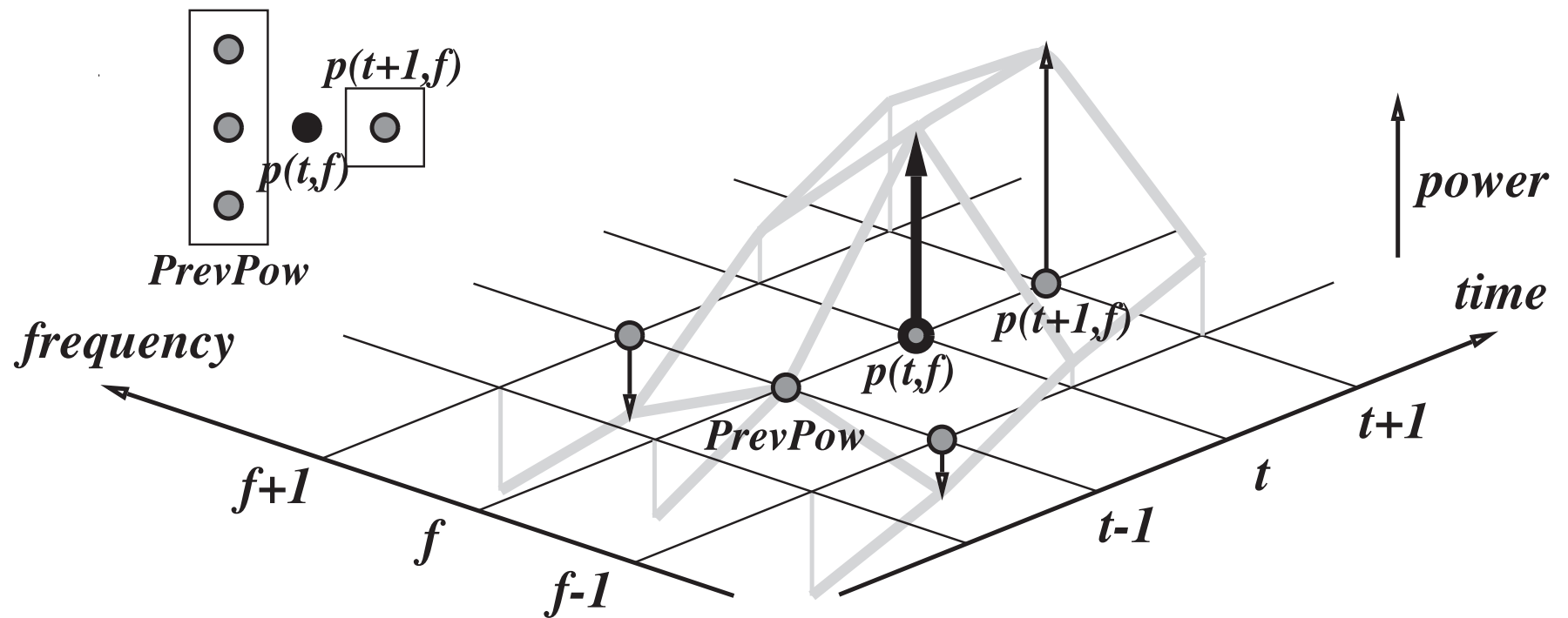
ONSETS

CHORD CHANGES

DRUM PATTERNS

(Goto 2001)

ONSET DETECTION IN CONTEXT



(Goto 2001)

f are bands

t are frames

prevpow - power spectrum of previous frame in same and neighboring bands

$p(t+1, f)$ is the same band in the next frame

DRUM PATTERNS

BASS DRUM - LOW FREQUENCY

SNARE DRUM - BROADBAND

FIND PATTERNS

METRICAL STRUCTURE

(Goto 2001)

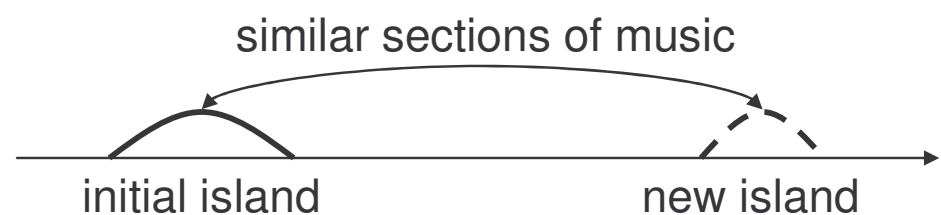
HOLISTIC BEAT TRACKING

HOLISTIC - METRICAL STRUCTURE

SIMPLE INITIAL HYPOTHESIS

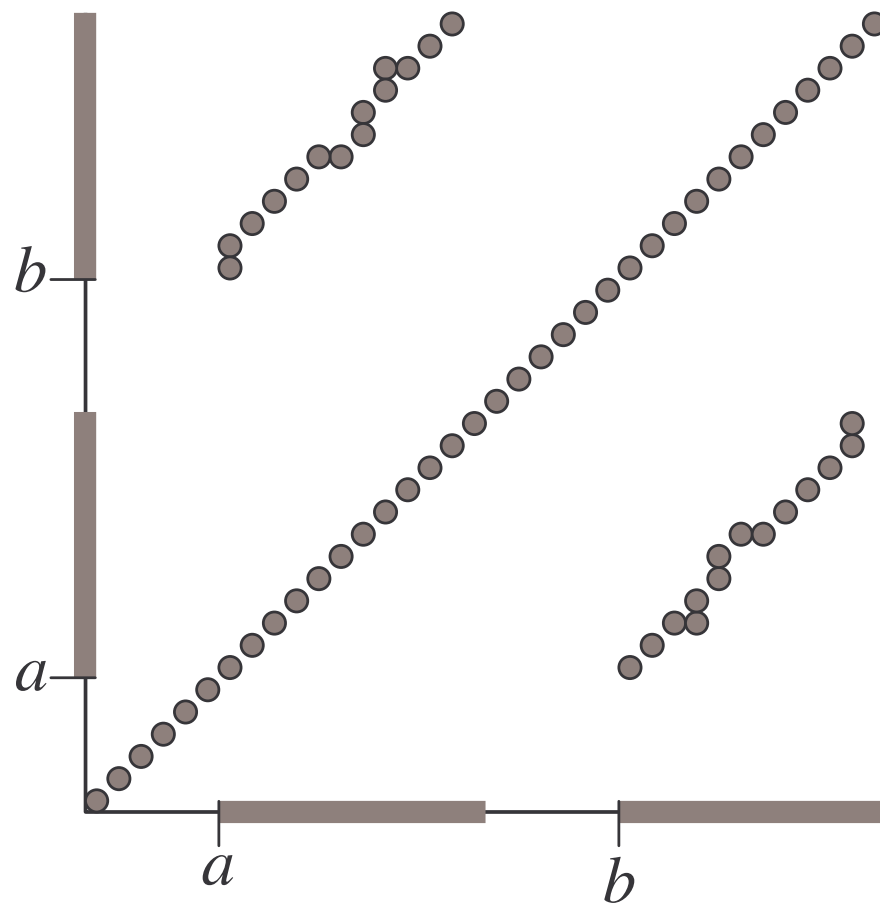
SELF-SIMILARITY MATRIX

ISLANDS OF TEMPO



(Dannenberg 2005)

SELF-SIMILARITY MATRIX



(Dannenberg 2005)

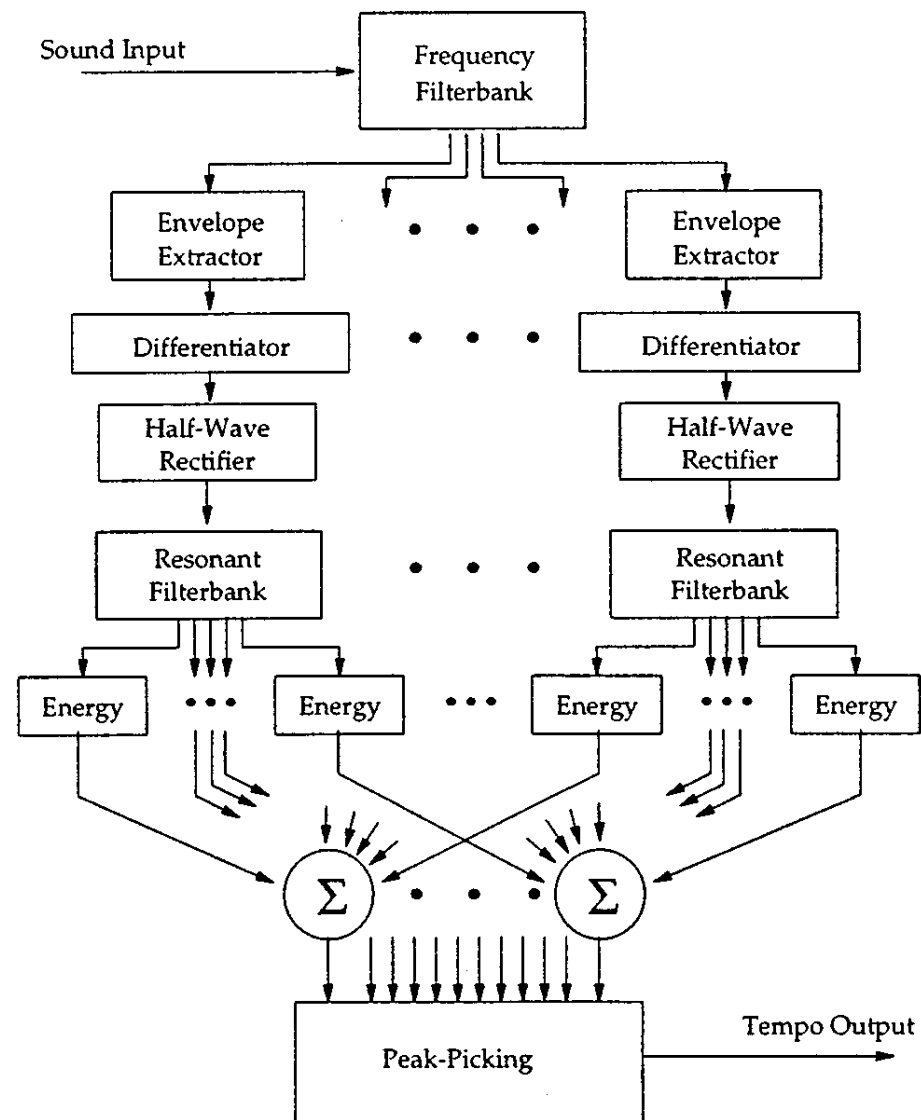
song segments a, b
diagonal where moment x is similar to itself

PERCEPTUAL VS. TRANSCRIPTIVE

ONSETS: DETECT OR NOT?

IMPLICIT RHYTHMIC EVENTS

PERCEPTUAL BEAT ANALYSIS



(Scheirer 1998)

The resonators are comb filters
6 bandpass filters

PERCEPTUAL BAYESIAN

RHYTHM TRACKS

BAYESIAN ANALYSIS OR GRADIENT DESCENT

(Sethares et al. 2005)

EVALUATION

NO COMMON CORPUS

NO GROUND TRUTH

CONFOUNDING COMPARISONS, E.G. VIOLA/VIOLIN

(Temperley 2004)

Pitch

WHAT?

TRANSCRIPTION

CHORDS

INFORM HOLISTIC APPROACHES

How?

RULE-BASED

PROBABALISTIC/LEARNING

TRANSCRIPTION

JUST PICK THE FREQUENCIES FROM THE SPECTRUM,
RIGHT?

TRANSCRIPTION

COMPLICATIONS

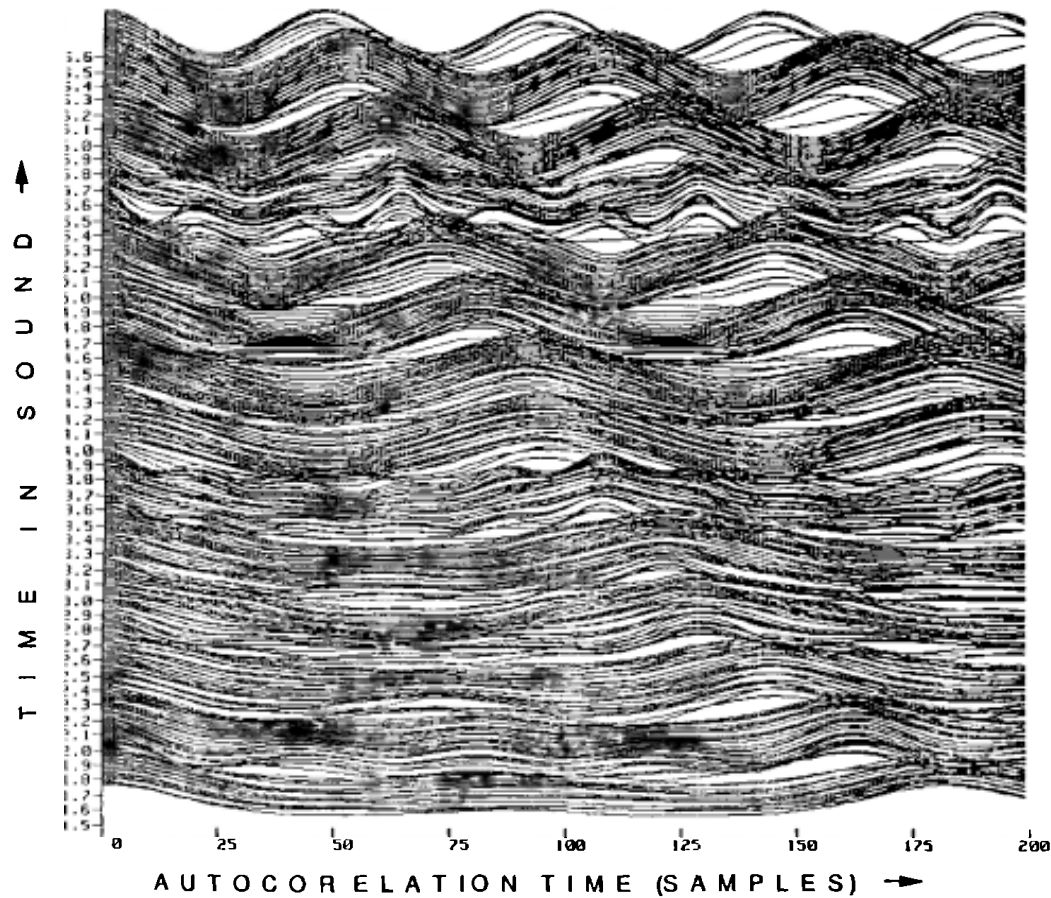
COMPLEX AND CHANGING TIMBRE

NOISE

POLYPHONIC CONFUSION

RHYTHM

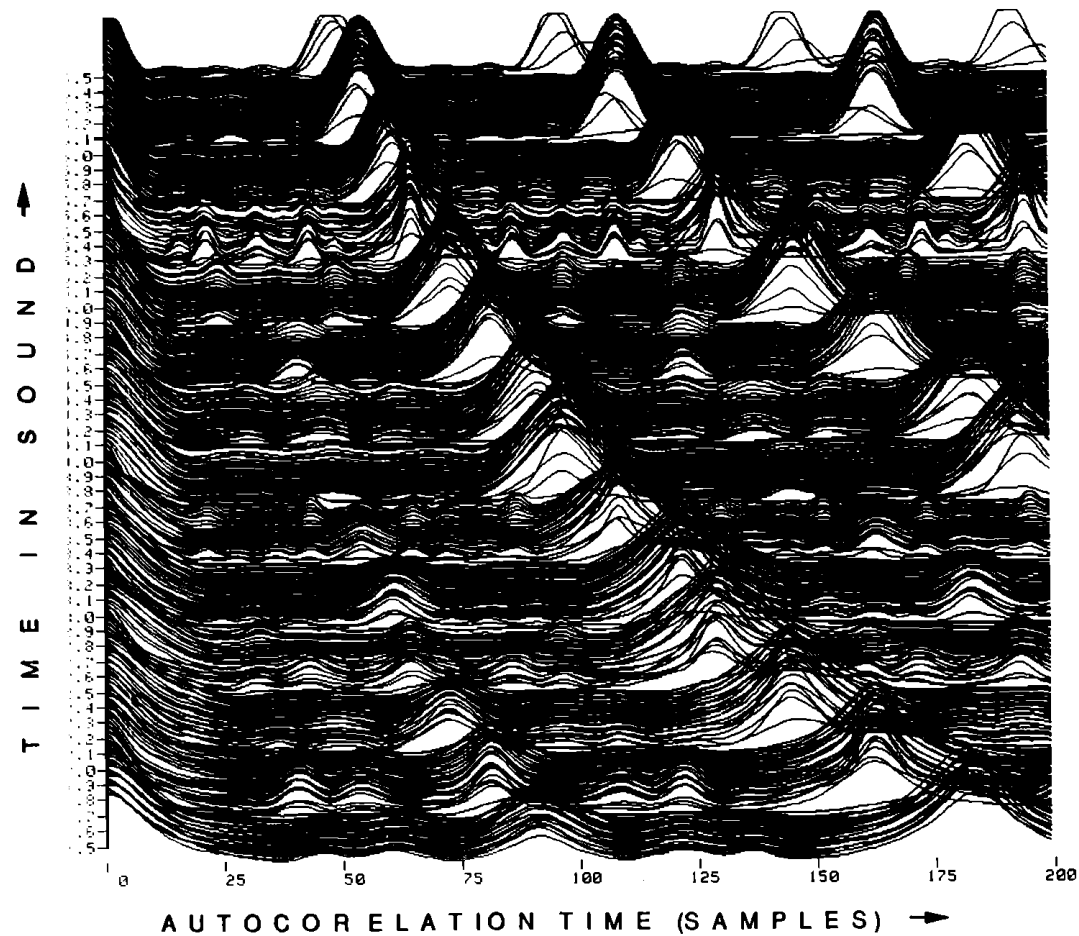
AUTOCORRELATION



$$r_{xx}(n) = \frac{1}{N} \sum_{t=1}^{N-n-1} x(t)x(t+n)$$

(Brown and Zhang 1991)

NARROWED AUTOCORRELATION



$$S_N(\tau)^2 = \langle |f(t) + f(t - \tau) + f(t - 2\tau) + \dots + f(t - (N - 1)\tau)|^2 \rangle$$

(Brown and Zhang 1991)

Where $\langle |x|^2 \rangle$ seems to mean $\langle x, x \rangle$ (inner product)
same scale
note logarithmic nature
Neither is suitable for polyphony

TRANSCRIPTION

APPROACHES:

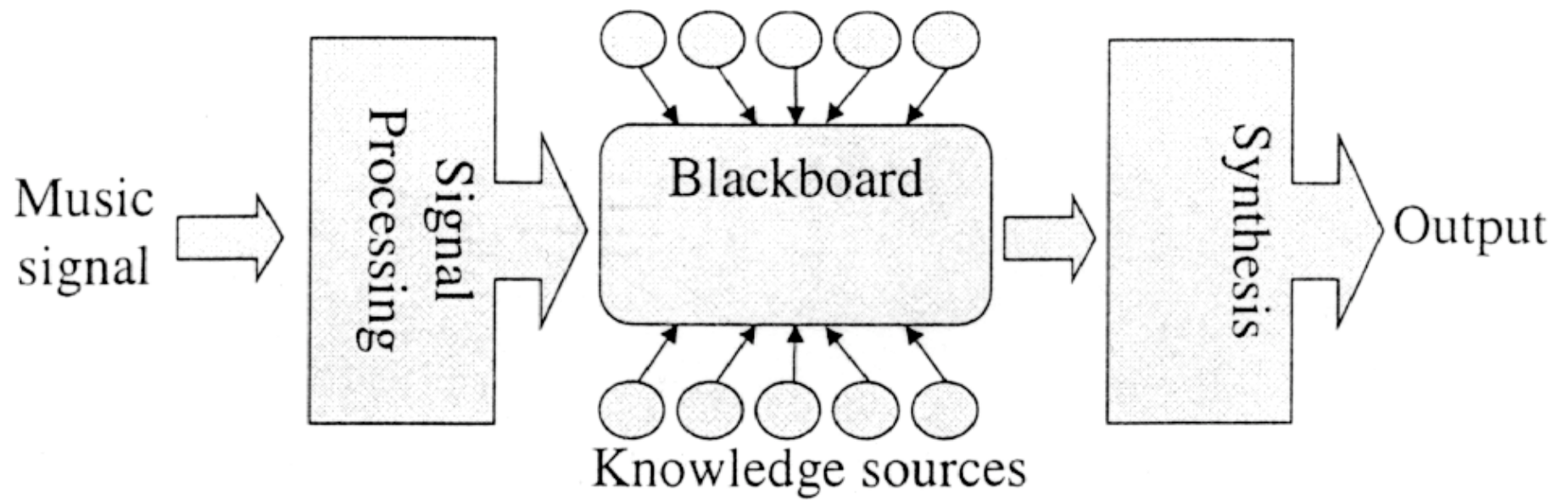
KNOWLEDGE-BASED BLACKBOARD

MULTIPLE-CAUSE MODEL

INDEPENDENT COMPONENT ANALYSIS

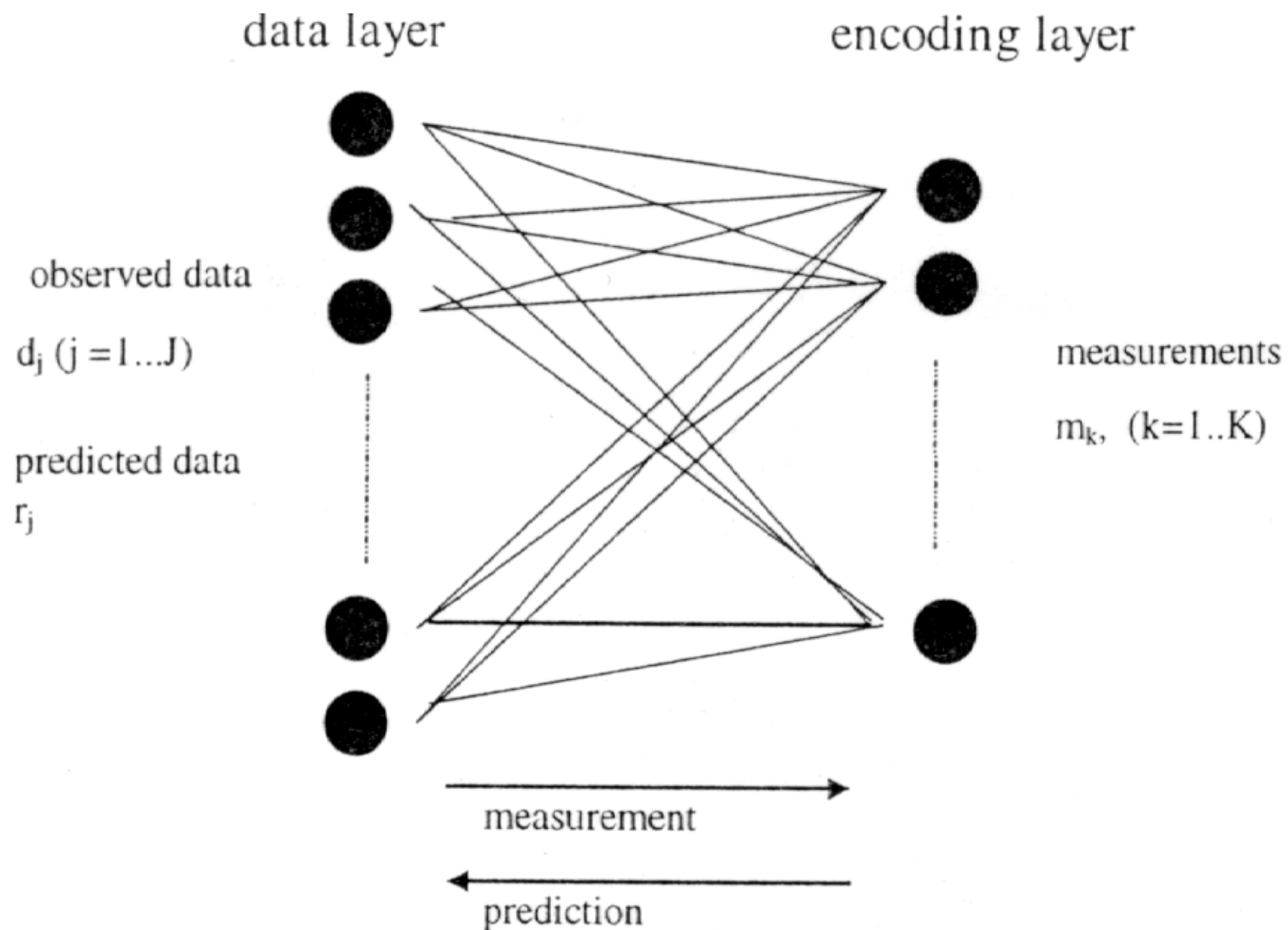
(Plumbley et al. 2002)

BLACKBOARD



(Plumbley et al. 2002)

MULTIPLE-CAUSE MODEL

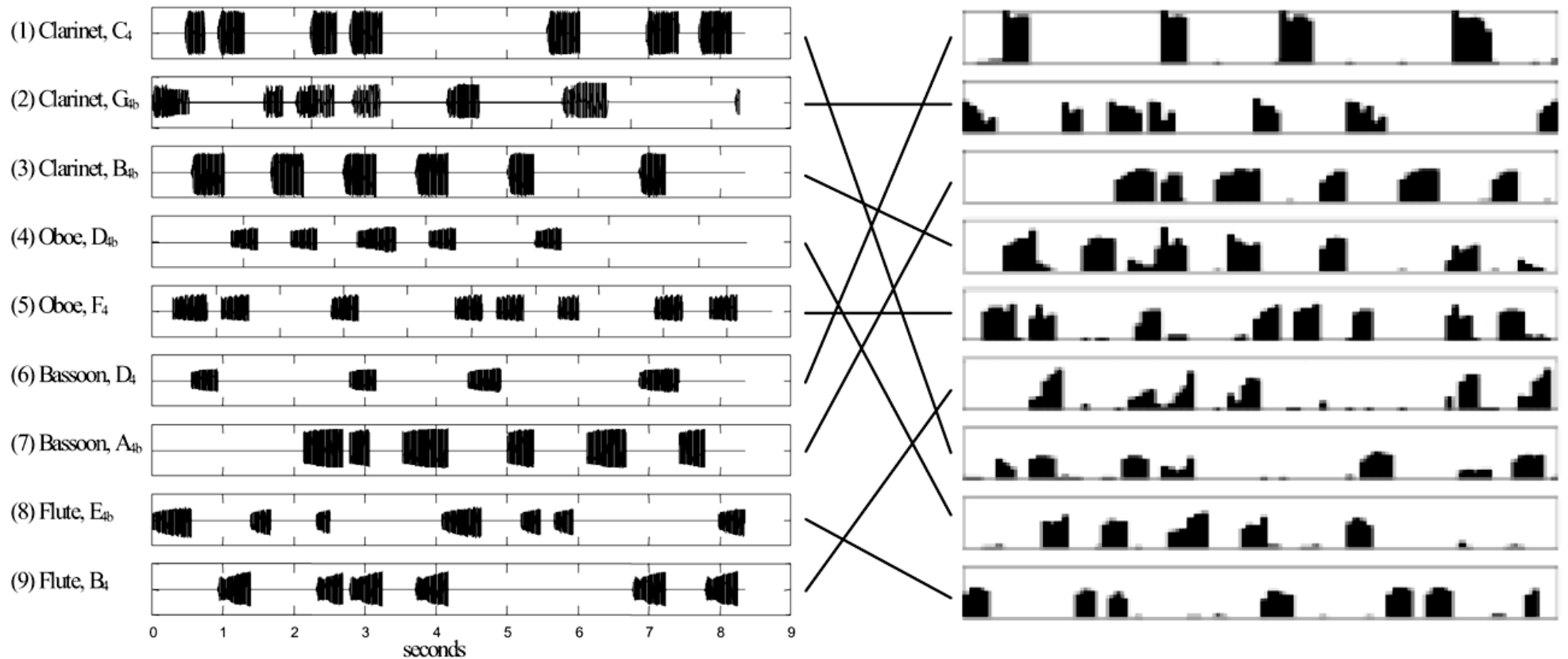


(Plumbley et al. 2002)

neural network
feeds forward to give measurements, then backward to predict
weights adjusted by minimizing prediction error

MULTIPLE-CAUSE MODEL

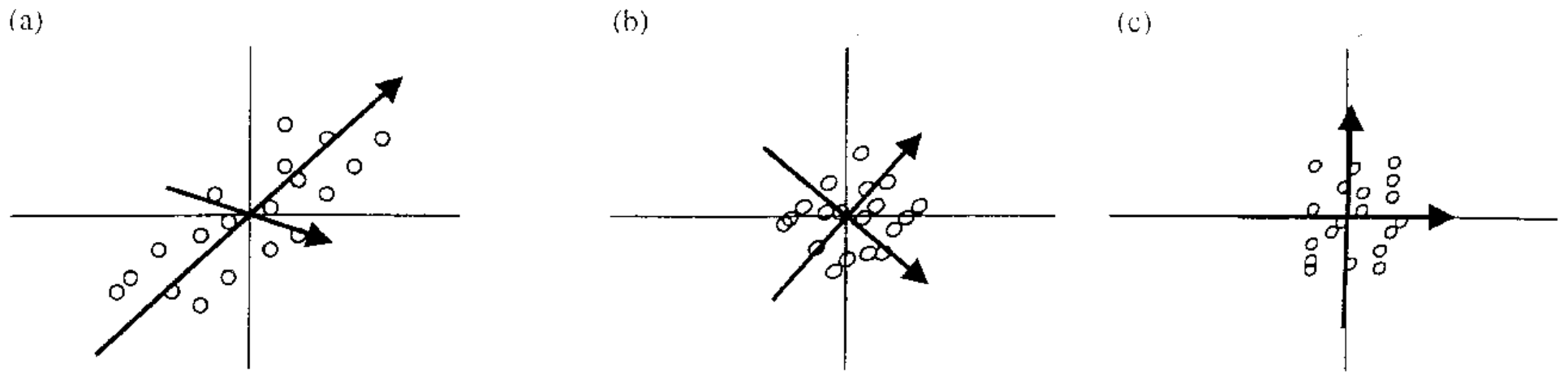
Input before mixing → Single Audio Channel → Activations after learning



(Plumbley et al. 2002)

Mapping unknown
flute/clarinet confusion, because of similarities?

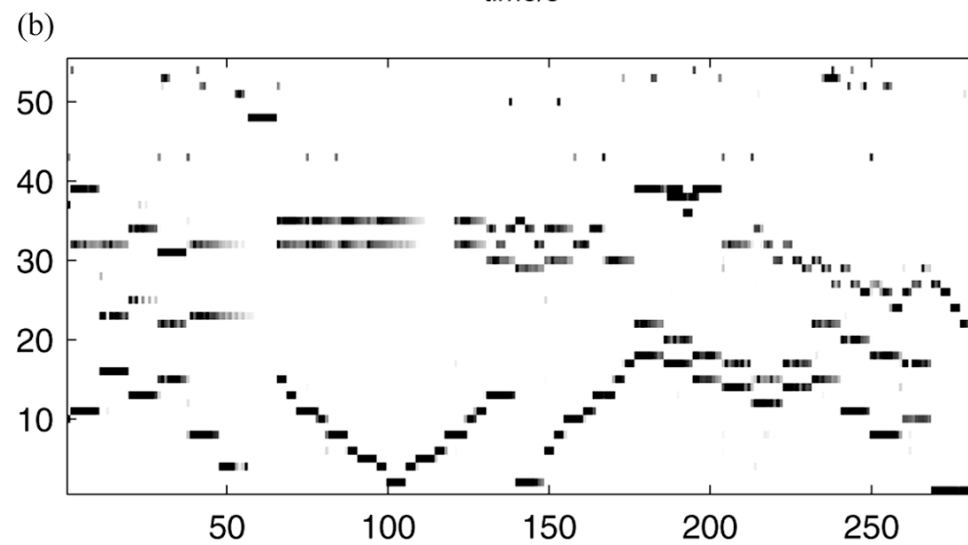
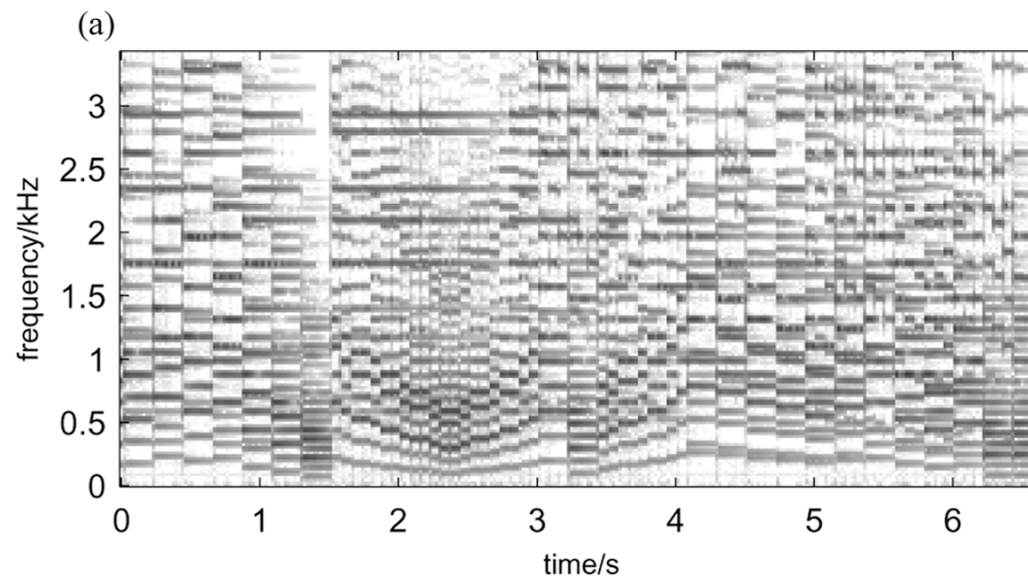
ICA



(Plumbley et al. 2002)

Independent Component Analysis
a) original data
b) after whitening
c) after rotation
Often used in source separation

ICA



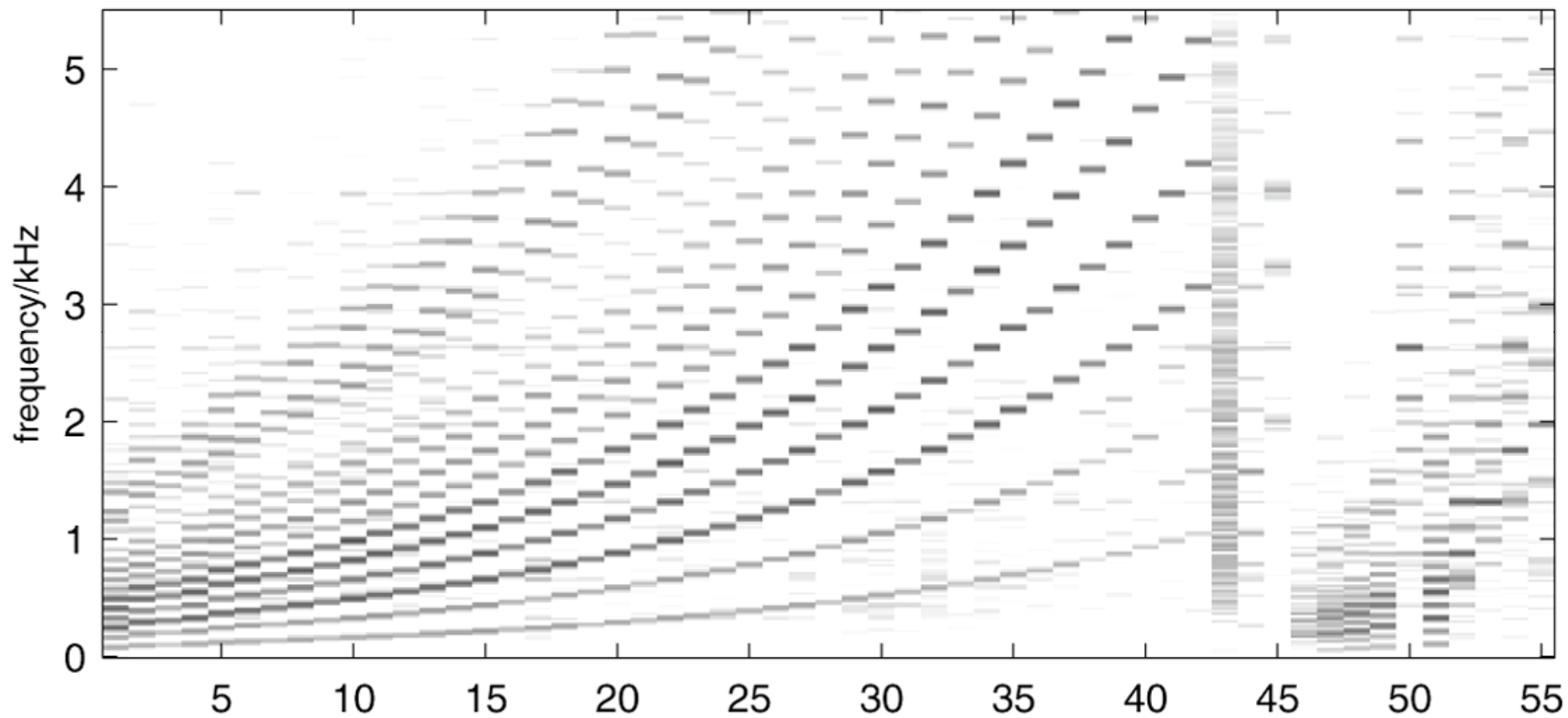
(Plumbley et al. 2002)

Plumbley et al., each note a source
Bach Partita
Spectrum and output note activations

ICA

(c)

Basis matrix



(Plumbley et al. 2002)

Note Shapes
Ordered Manually
Clear harmonic structure below 43
above 43 probably transients

GENERATIVE MODEL FOR TRANSCRIPTION

GRAPHICAL (DYNAMICAL) MODEL

BAYESIAN NETWORK

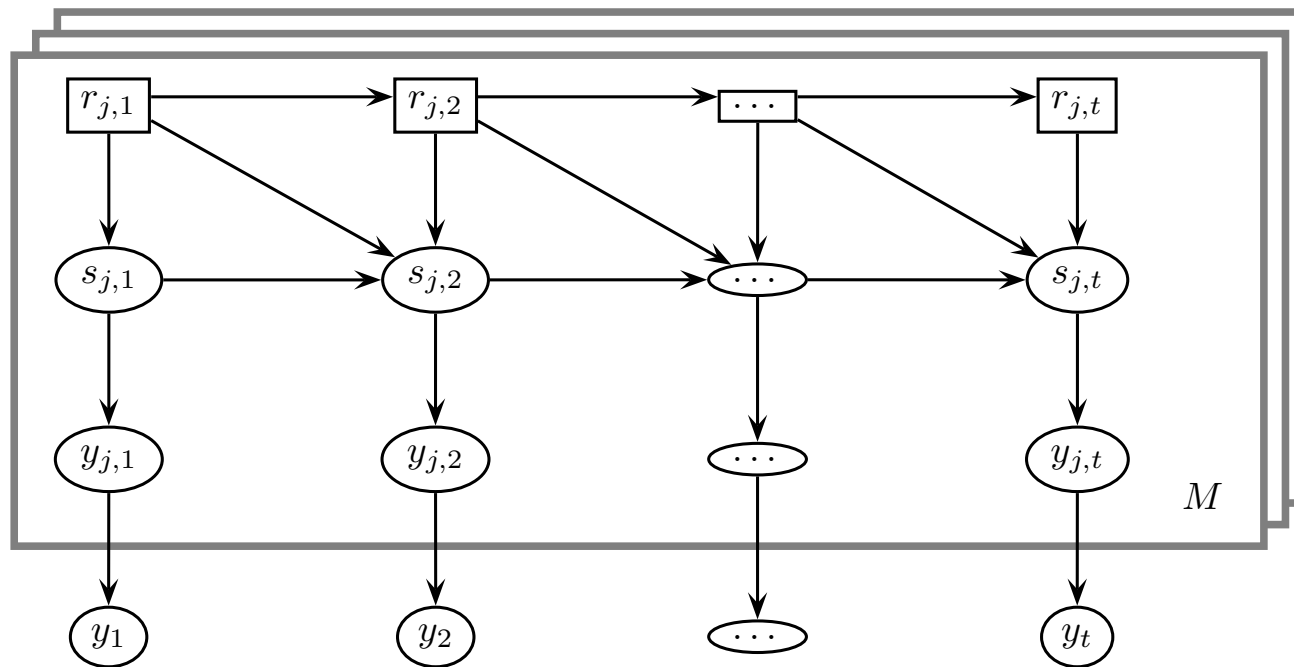
MODELS SOUND GENERATION



KALMAN

(Cemgil et al. 2004)

BAYESIAN NETWORK



$$P(r|y) \propto P(y|r)P(r)$$

(Cemgil et al. 2004)

Predict from state + piano roll
 M plates ($j=1..M$), one for each sound generator
 time t
 r for roll, s for state
 sum sound generators together at each **sample** to get y
 $P(y|r)$ from the generative model (Kalman filter used to find hidden states in box)
 $P(r)$ prior can incorporate musical knowledge

GRAPHICAL MODEL

FOR SONG MELODIES

JOINT MODEL OF PITCH, RHYTHM, TEMPO, SEGMENTATION

EVENT LIST AS A MARKOV CHAIN

OPTIMAL PARAMETERS OF MODEL VIA DYNAMIC PROGRAMMING

BRANCH AND PRUNE

(Raphael 2005)

MARKOV CHAIN

SEQUENCE OF RANDOM VARIABLES

GIVEN PRESENT, FUTURE IS INDEPENDENT OF PAST

I.E. ONLY NEED THE LAST OBSERVATION

(Raphael 2005)

AUDIO THUMBNAILS

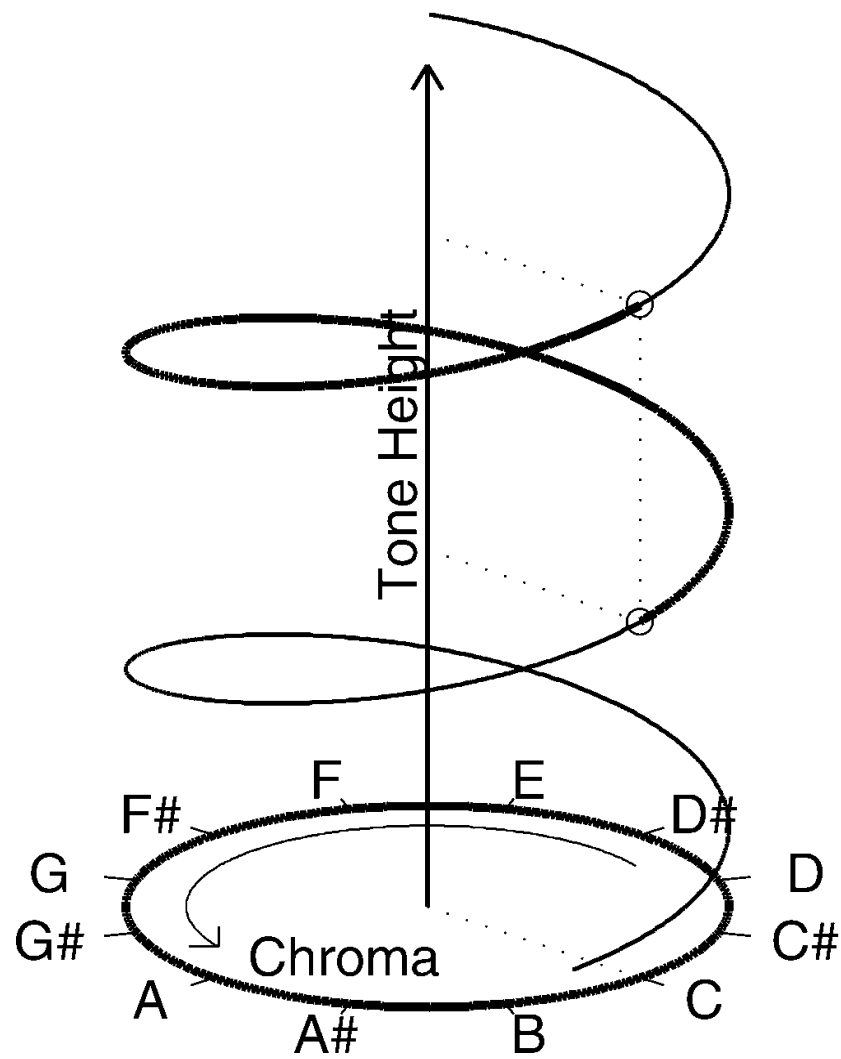
CHROMAGRAM

SIMILARITY MATRIX

MOST STRONGLY REPEATED SECTION

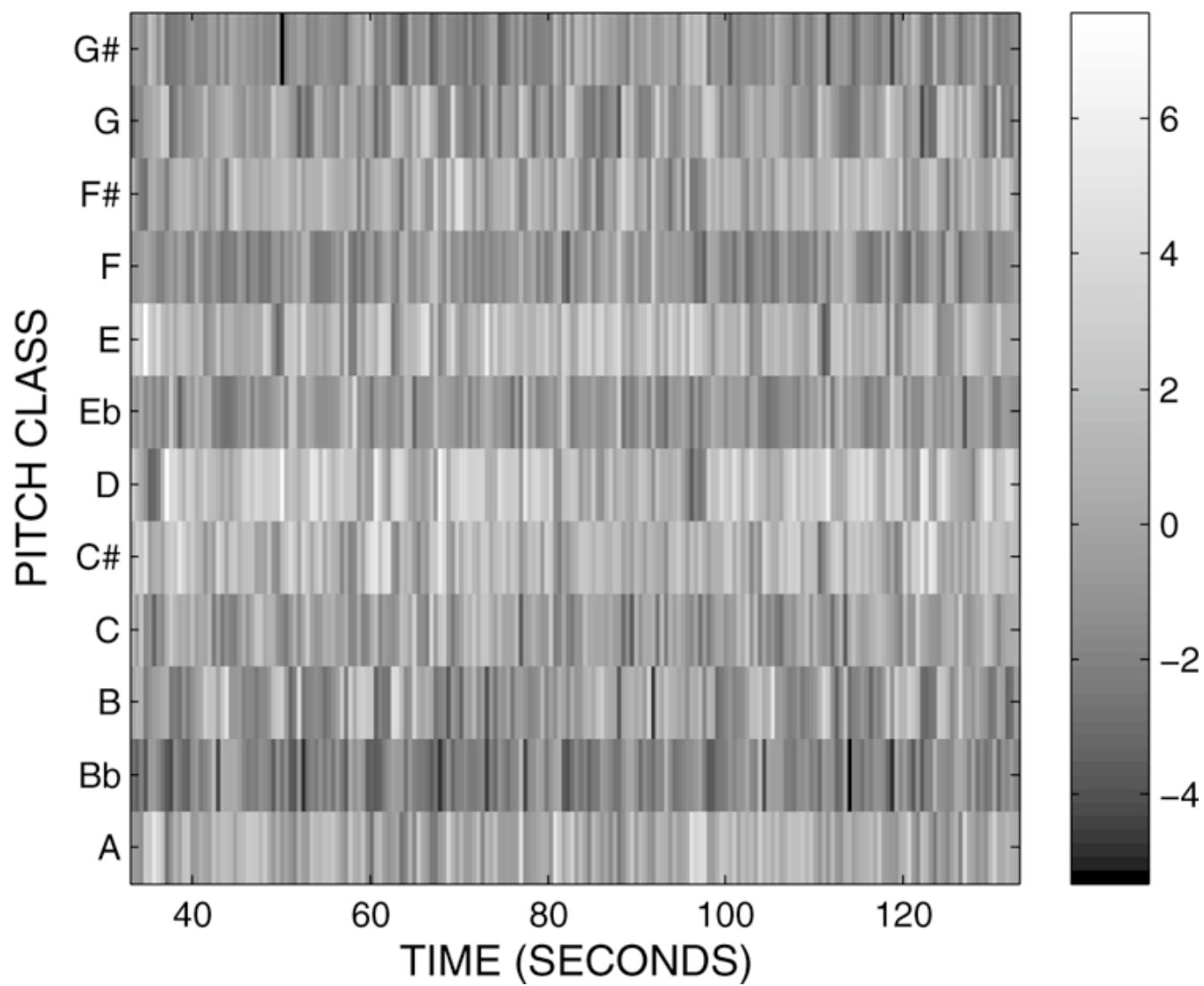
(Bartsch 2005)

CHROMA



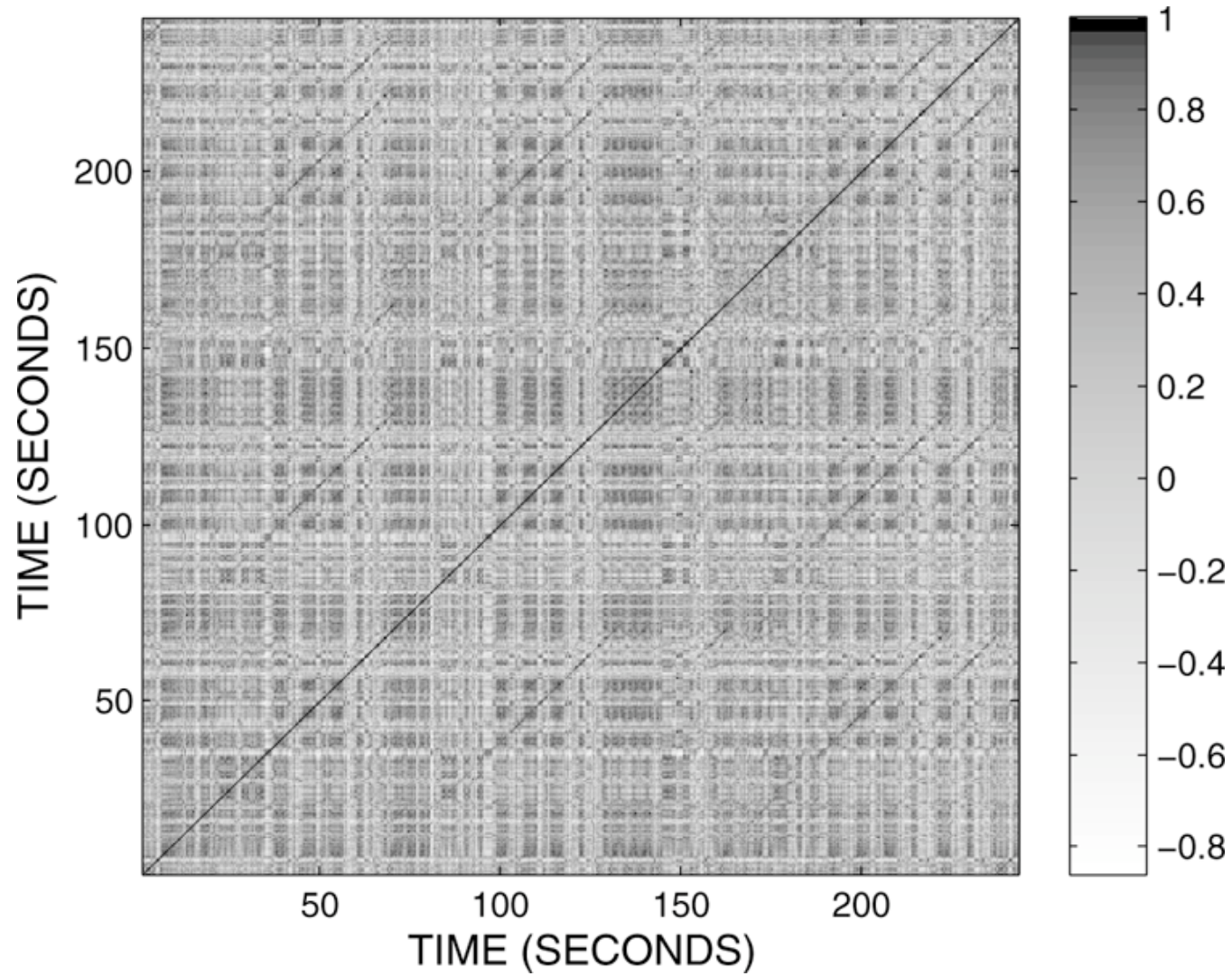
(Bartsch 2005)

CHROMAGRAM



(Bartsch 2005)

SIMILARITY MATRIX



(Bartsch 2005)

KEY, CHORD, AND RHYTHM

HOLISTIC: KEY, CHORDS, RHYTHM

RULE-BASED

(Shenoy and Wang 2005)

STEPS

BEAT TRACKING

ELIMINATE NON-KEY CHORDS

AUDIO SEGMENTATION

INCONSISTENT/MISSING
CHORDS

CHROMA FEATURES EXTRACTED

RHYTHM/METRICAL
CONSIDERATIONS

CHORD AND KEY

(Shenoy and Wang 2005)

CONCLUSION

MONOPHONIC IS WORKABLE

POLYPHONIC IS HARD

PROBABALISTIC/LEARNING METHODS MORE ROBUST

Tímbrre

WHAT?

INSTRUMENT IDENTIFICATION/CLASSIFICATION

INSTRUMENT SEPARATION

GENRE CLASSIFICATION

How?

OVERWHELMINGLY BY MACHINE LEARNING

SOME RULE-BASED SYSTEMS



MACHINE LEARNING

LEARNING CYCLE

OBSERVATION

ACTION

EVALUATION

ADJUSTMENT

PROBLEM FORMULATION

IMPORTANT!

SUFFICIENT AND REPRESENTATIVE DATA

EVALUABLE AND USEFUL TASK

PROPER EVALUATION, ENSURE GENERALIZATION

INSTRUMENT IDENTIFICATION

MONOPHONIC OR POLYPHONIC

SINGLE NOTE OR WHOLE PASSAGES

STEPS

EXTRACT FEATURES

APPLY CLASSIFIERS

SELECT RELEVANT FEATURES/CLASSIFIERS

FEATURES

AMPLITUDE

HARMONICS

ENERGY

INHARMONICS

ZERO-CROSSING RATE

CEPSTRAL COEFFICIENTS

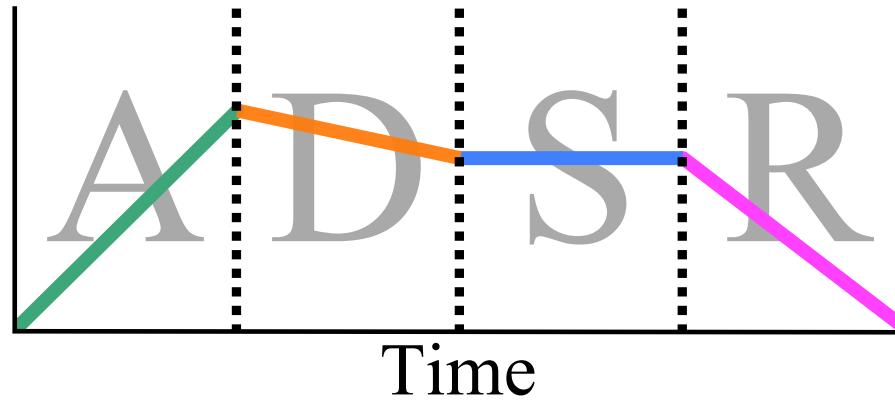
FREQUENCY

BARK, ERB

SPECTRAL SHAPE

CORRELOGRAM

AMPLITUDE



ATTACK, DECAY, SUSTAIN, RELEASE

VOLUME

FREQUENCY

FUNDAMENTAL (f_0)

RANGE

VIBRATO

SPECTRAL SHAPE

CENTROID

SPREAD

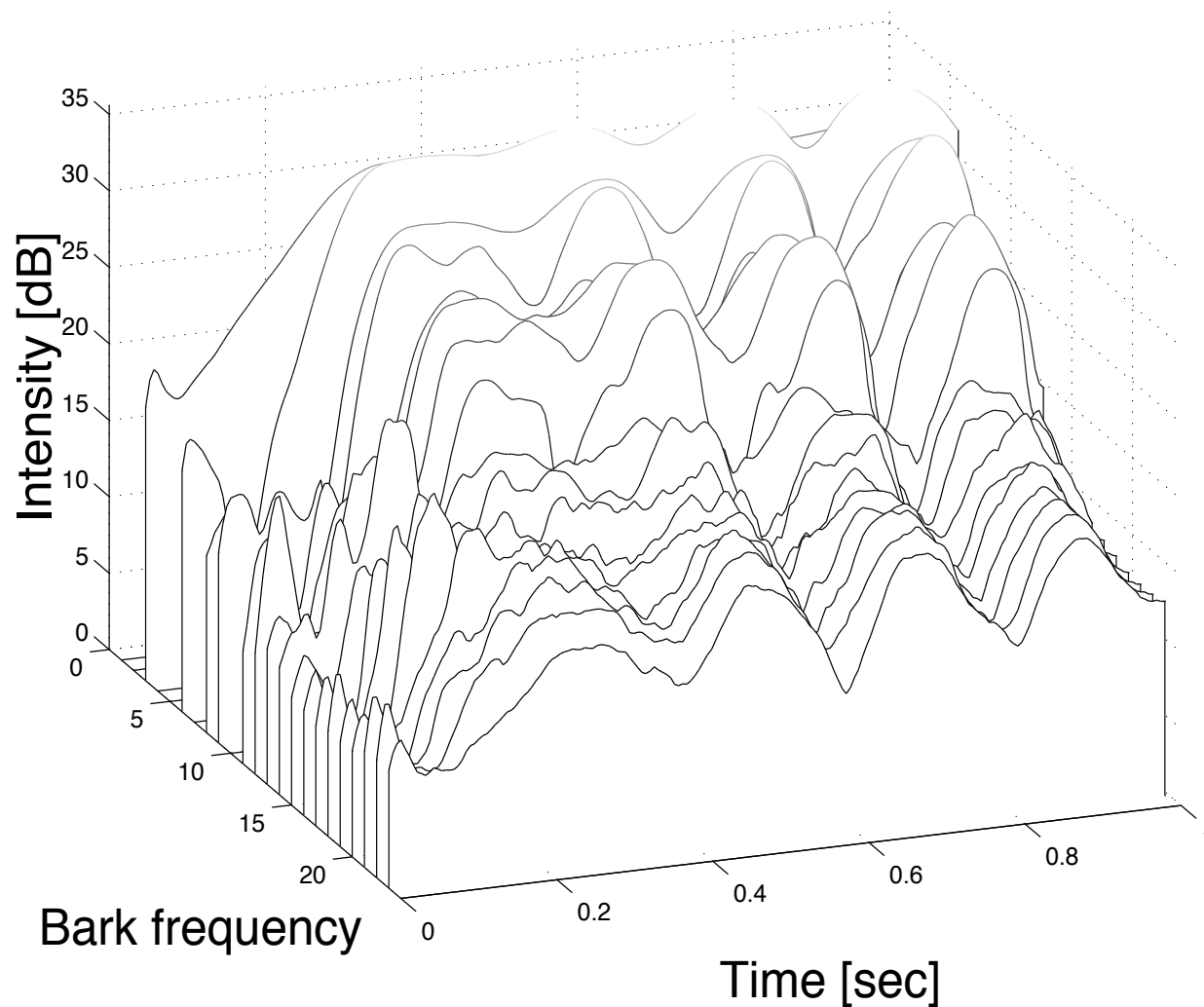
FLATNESS

FLUX

SKEW

centroid like mean of normal, most popular

SINUSOID ENVELOPES



(Eronen 2001)

Like Fourier, but bins calculated independently
sample-accurate
useful for examining spectral evolution
That's a flute

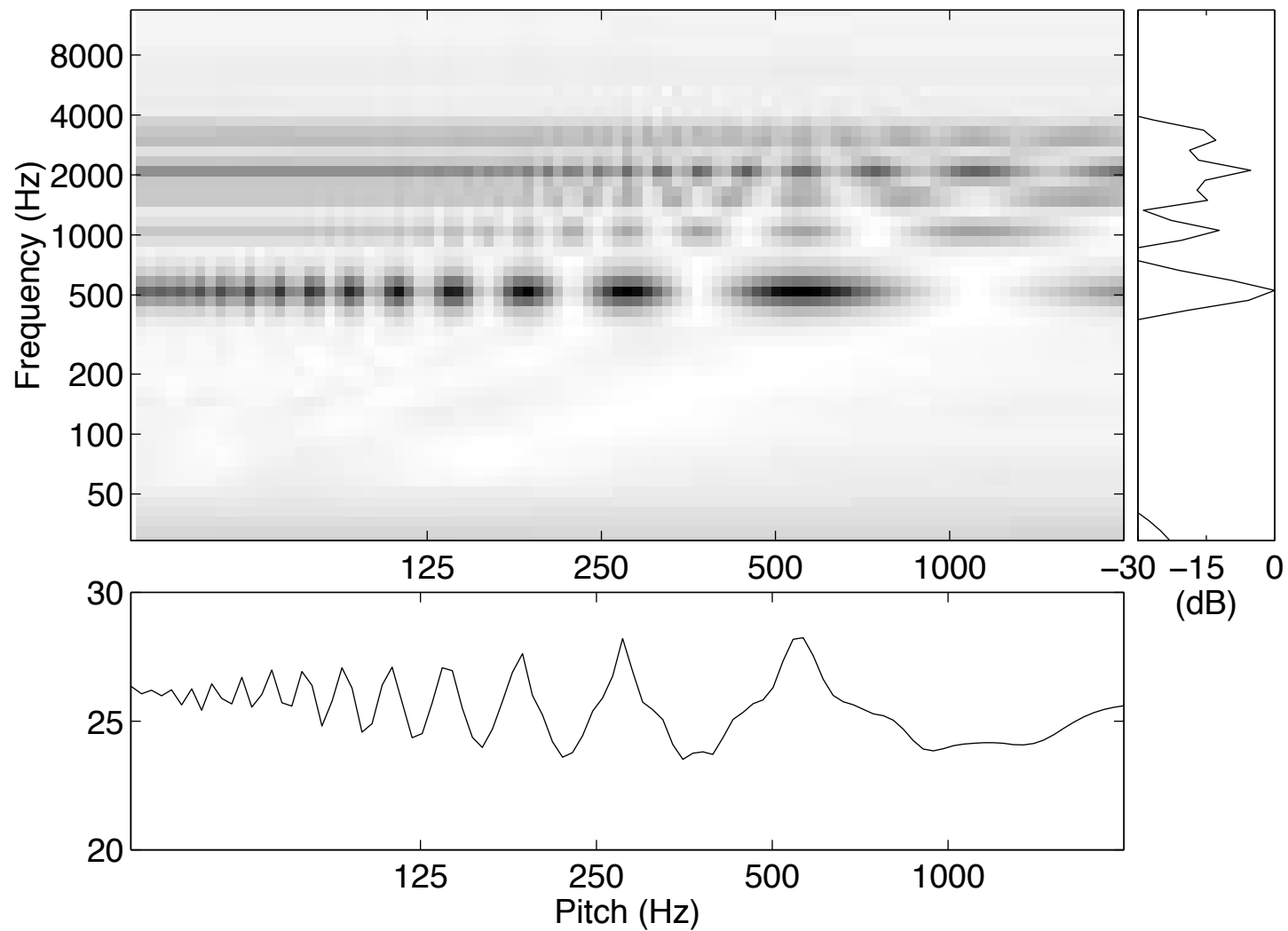
CEPSTRAL COEFFICIENTS

$$c = \mathcal{F}^{-1} (\log |\mathcal{F}(x)|)$$

USUALLY IN CONJUNCTION WITH THE MEL SCALE
(MFCC)

$$m(f) = 2595 \log(1 + f/700)$$

CORRELOGRAM



(Martin and Kim 1998)

“not based on the assumption that the signal is periodic”
modeled after cochlea and inner hair cells
filtered into channels and autocorrelation
that’s a violin

CLASSIFIERS

K-NN

ARTIFICIAL NEURAL NETWORKS

NAIVE BAYES

ROUGH SETS

DISCRIMINANT ANALYSIS

HIDDEN MARKOV MODELS

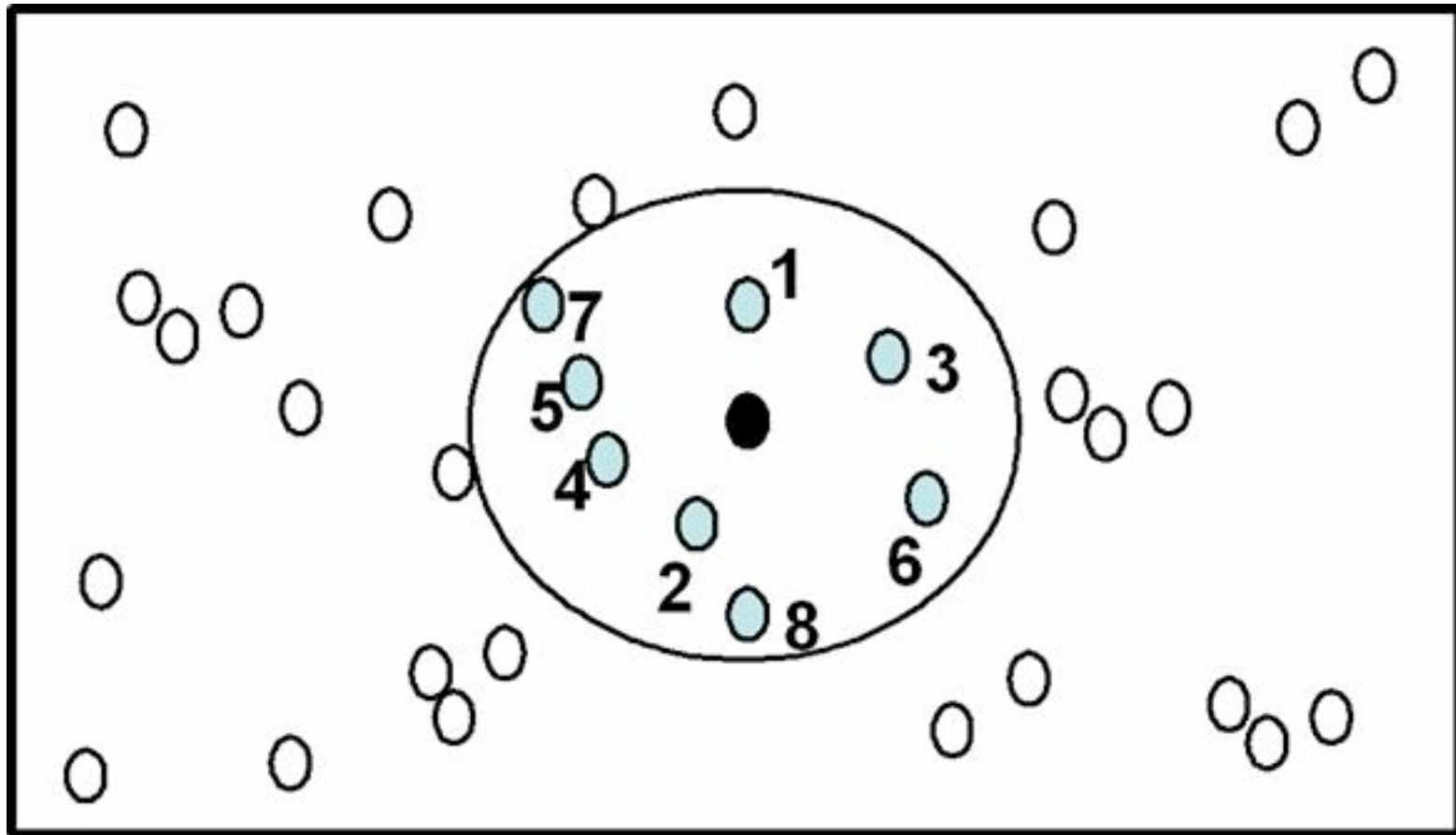
DECISION TREES

GAUSSIAN MIXTURE MODEL

SUPPORT VECTOR MACHINES

INDEPENDENT SUBSPACE
ANALYSIS

K-NEAREST NEIGHBORS



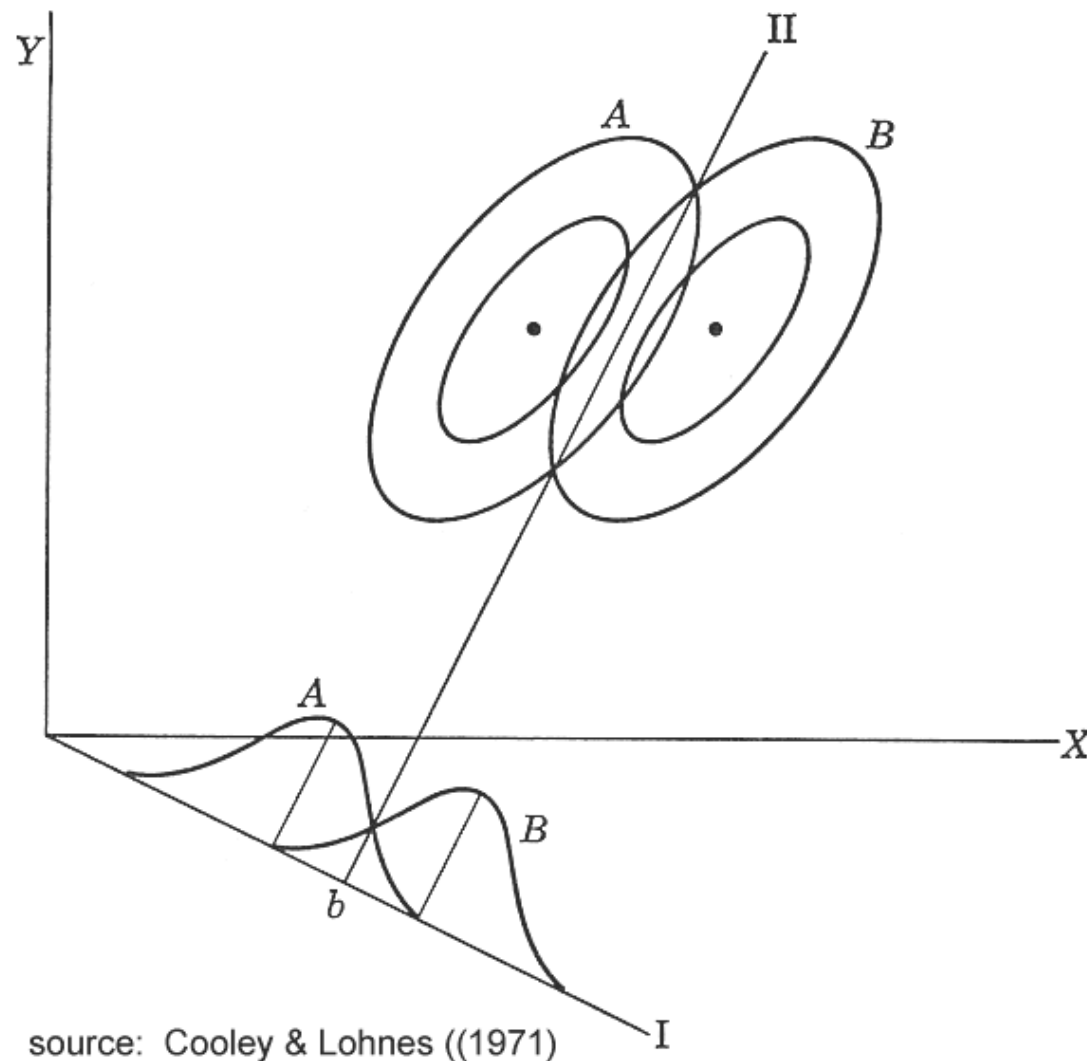
(University College Dublin)

NAIVE BAYES

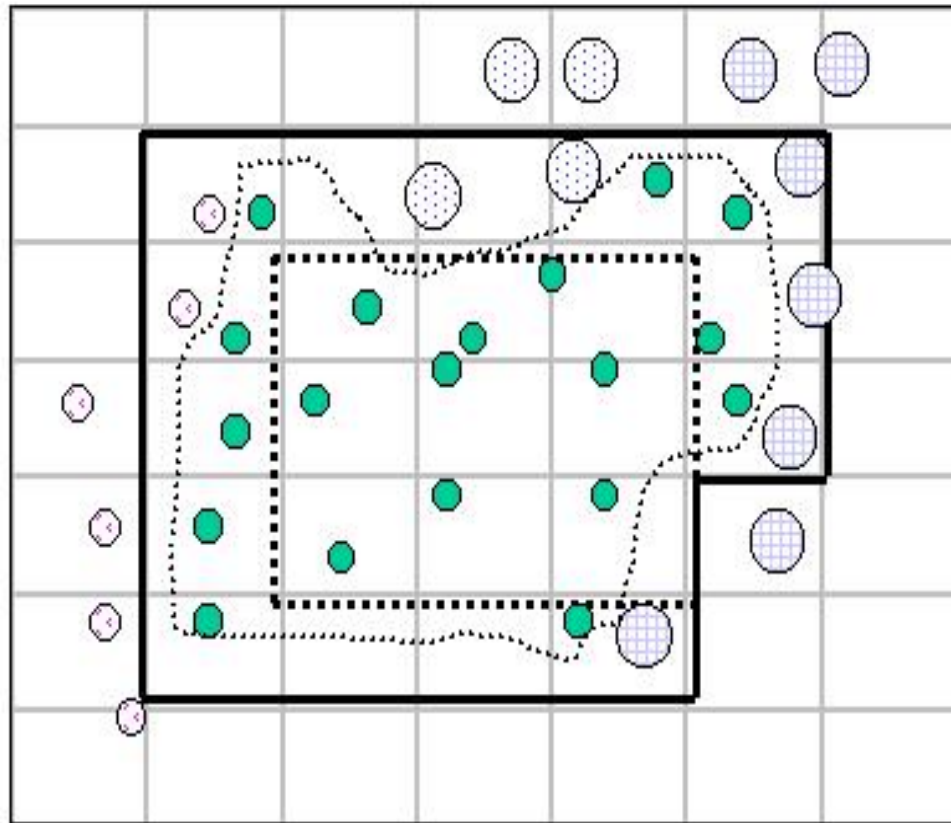
V IS SET OF CLASSES, A ARE FEATURES

$$v_{\text{NB}} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_i P(a_i | v_j)$$

DISCRIMINANT ANALYSIS



ROUGH SETS

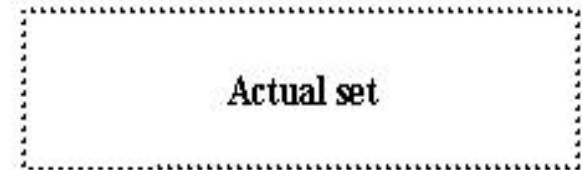
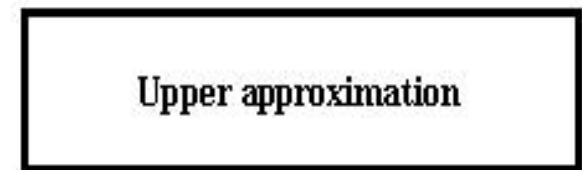


● violins



Non-violins

Equivalence class				



(Herrera-Boyer et al. 2003)

HIDDEN MARKOV MODEL

BAYESIAN NETWORK

STATE TRANSITION PROBABILITIES

OBSERVATION PROBABILITIES

STATES, TRANSITION PROBABILITIES, OBSERVATION
PROBABILITIES HIDDEN - ONLY GET OBSERVATIONS

SUPPORT VECTOR MACHINES

MAP DATA WITH KERNEL FUNCTION

FIND LINEAR HYPERPLANE TO MINIMIZE ERROR

SELECTION

PRINCIPAL COMPONENTS ANALYSIS

DISCRIMINANT ANALYSIS

SEQUENTIAL FORWARD/BACKWARD GENERATION

GRADUAL DESCRIPTOR ELIMINATION

features and/or classifiers

SFG/SBG: add/remove most/least relevant features one at a time

GDE: choose most irrelevant feature using LDA and remove, estimate at every step

STATUS

MUCH MONOPHONIC WORK DONE

MUCH POLYPHONIC WORK TO BE DONE

SHOTGUN APPROACH MOST SUCCESSFUL

Understanding

SPECTRAL ANTICIPATIONS

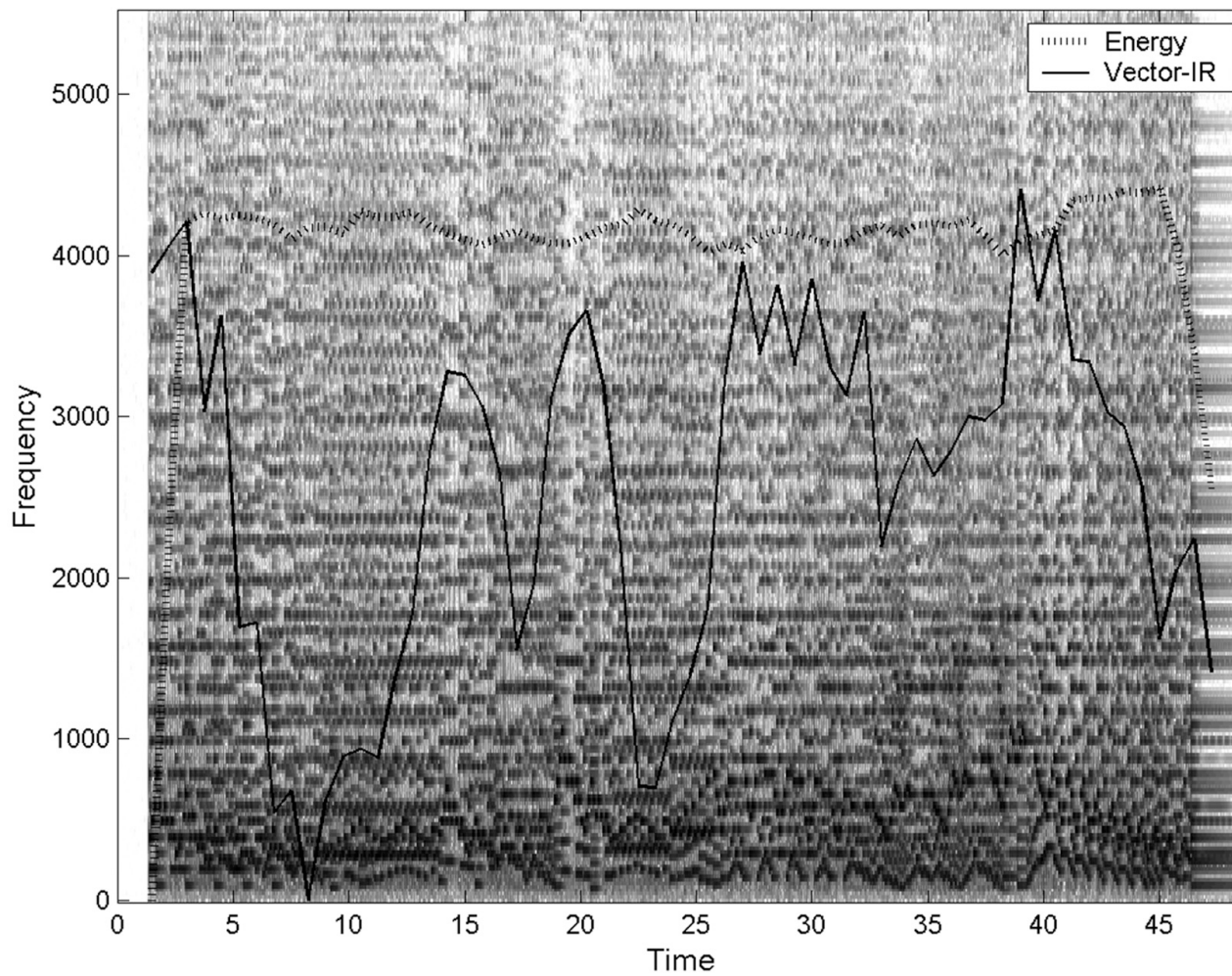
INVERTED U: SILENCE AND NOISE ARE BOTH BORING

STRUCTURE CARRIES INFORMATION

ANTICIPATION AND SURPRISE

(Dubnov 2006)

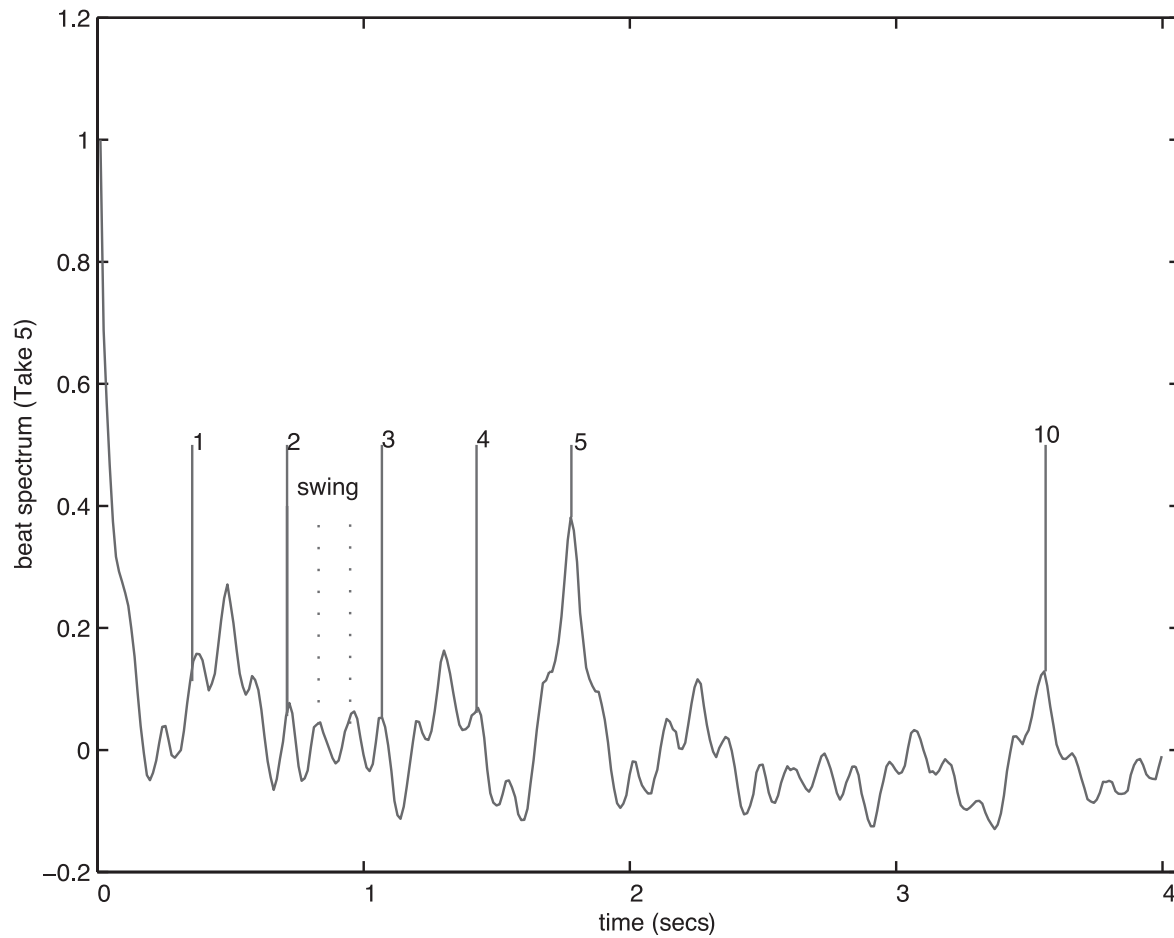
SPECTRAL ANTICIPATIONS



(Dubnov 2006)

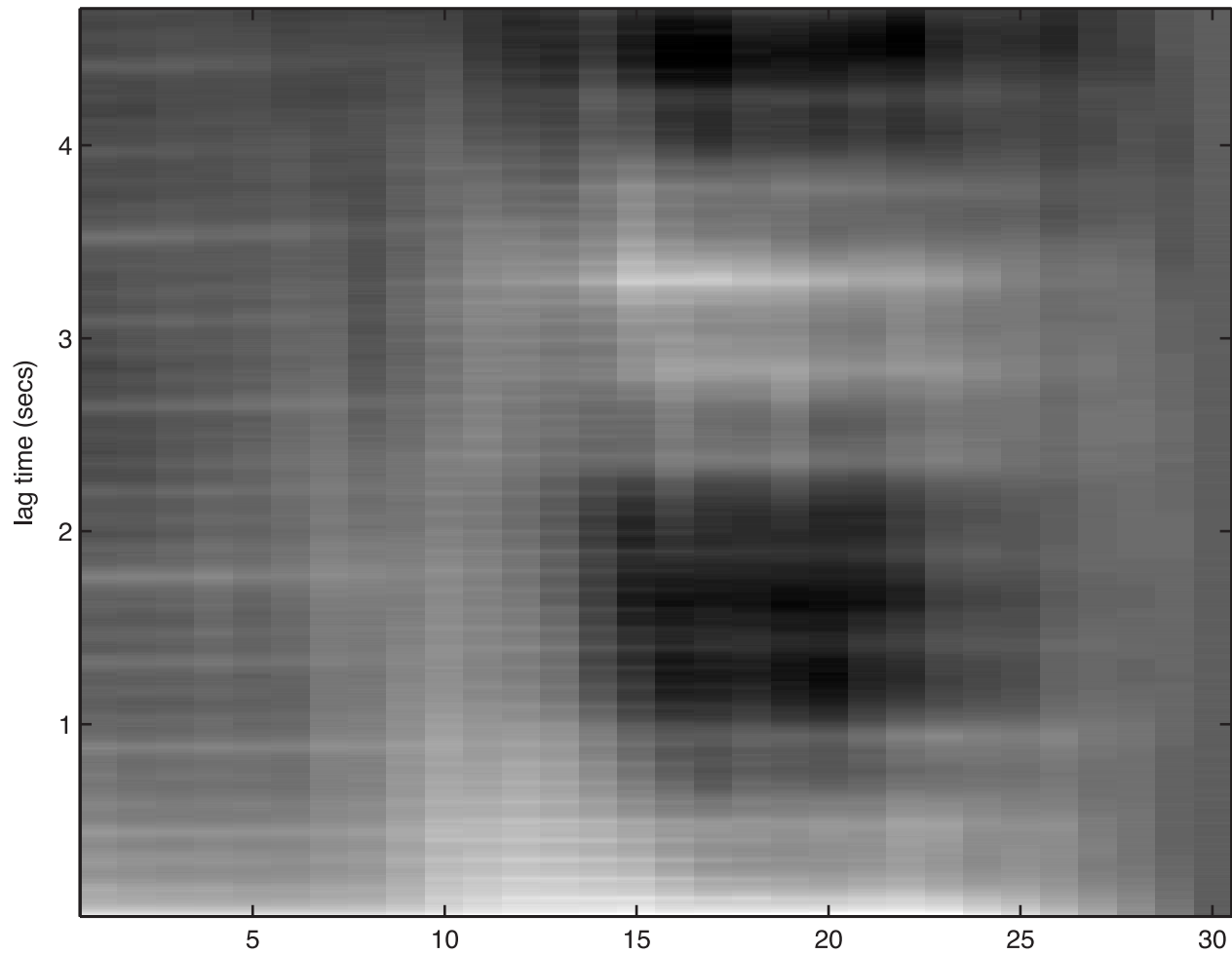
Energy doesn't say much, anticipation profile has more information

BEAT SPECTRUM



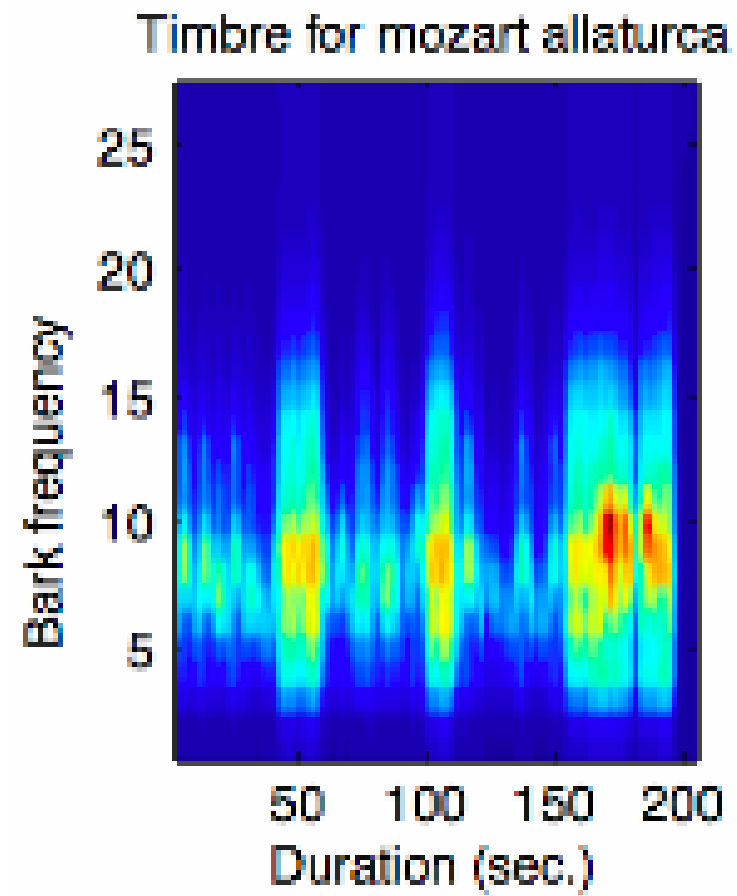
(Cooper et al. 2006)

BEAT SPECTROGRAM



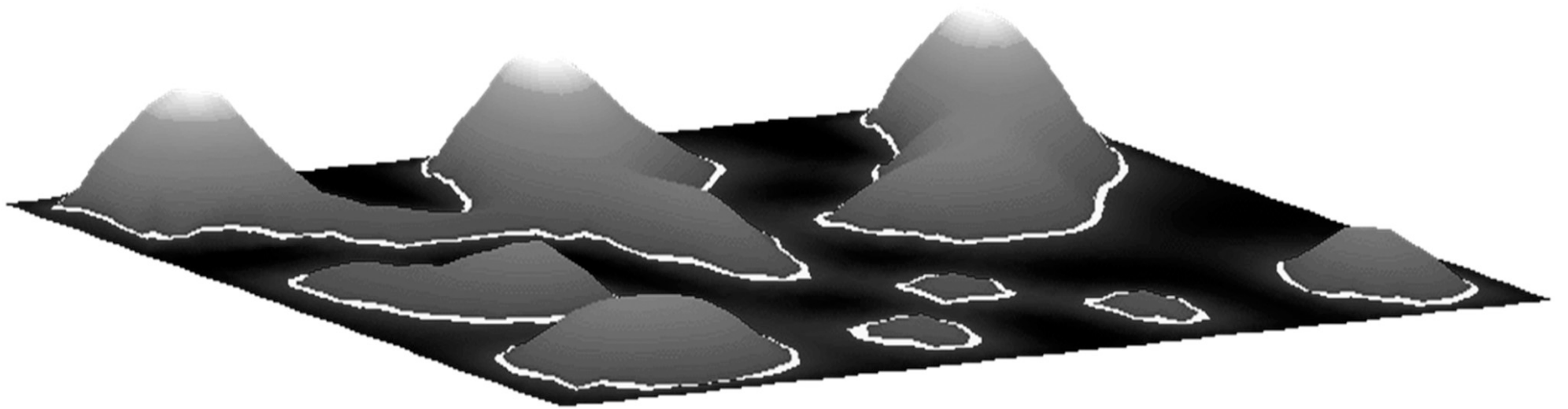
(Cooper et al. 2006)

TIMBREGRAM



(Kuhl and Jensen 2008)

ISLANDS OF MUSIC



(Cooper et al. 2006)

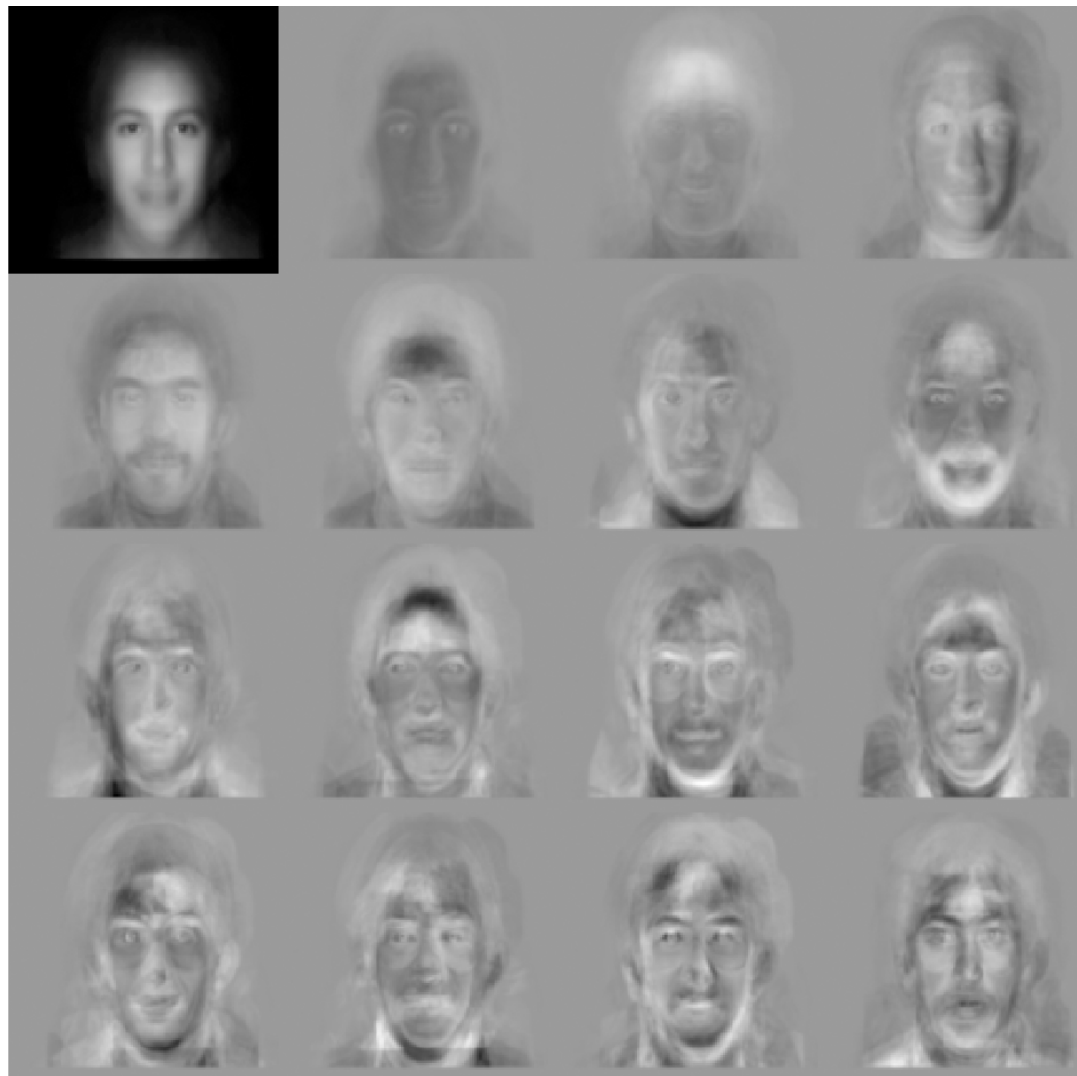
CLASSIFICATION VIA VISUALIZATION

SPECTROGRAM AND MFCC

TEXTURE-OF-TEXTURES

(Deshpande et al. 2001)

EIGENFACES



(Zhao and Chellappa 2006)

EIGENSOUND

ABSTRACT AWAY AUDITORY EVENTS

PCA

A GENERALIZED SPECTRAL TEMPLATE

WHITMAN'S EIGENRADIO

(Recht and Whitman 2003)

STYLE CLASSIFICATION

MIDI

LYRICAL OR ENERGETIC OR ...

13 MIDI-BASED FEATURES

NAIVE BAYES, GMM, ANN

(Dannenberg et al. 1997)

SPEECH/MUSIC DISCRIMINATOR

FEATURES: 4 HZ MODULATION ENERGY, LOW-ENERGY
FRAMES, SKEW, CENTROID, FLUX, ZCR, CEPSTRAL,
"RHYTHMICNESS"

GMM

(Scheirer and Slaney 1997)

Conclusion

Probabilistic and Machine Learning over direct programming
Wide application of many disciplines of computer science and signal processing
Complexity and vagueness of perception
Learn more about our own perception as we find what does and doesn't work