# An Engine for Comparative Time-Series Analysis

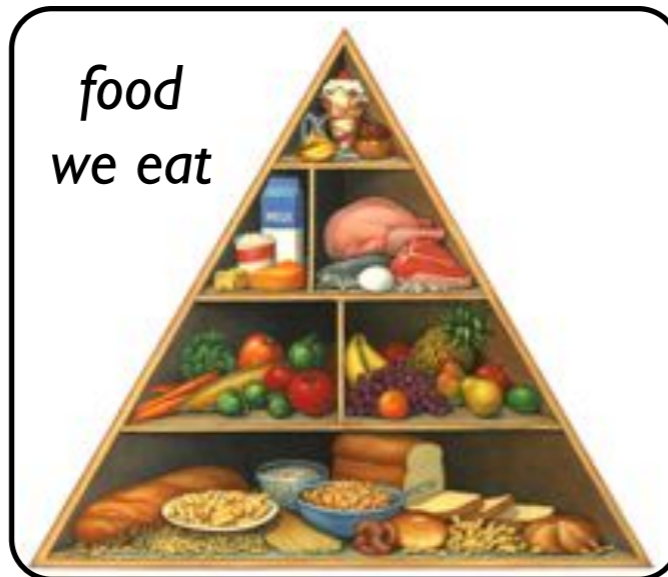*"the taming of the zoo"*

Ben Fulcher,

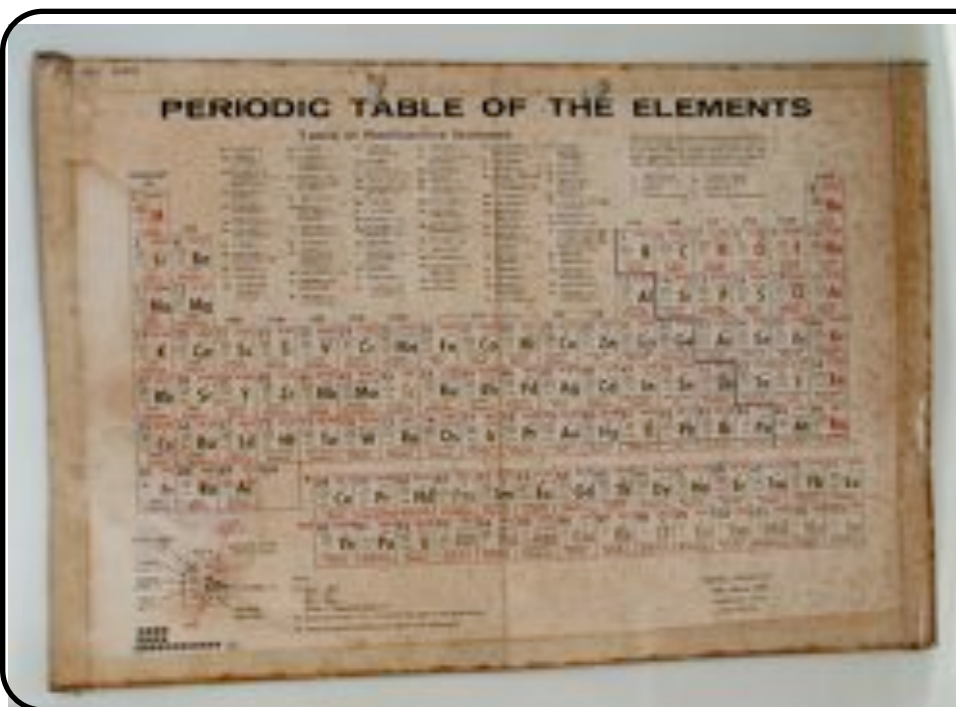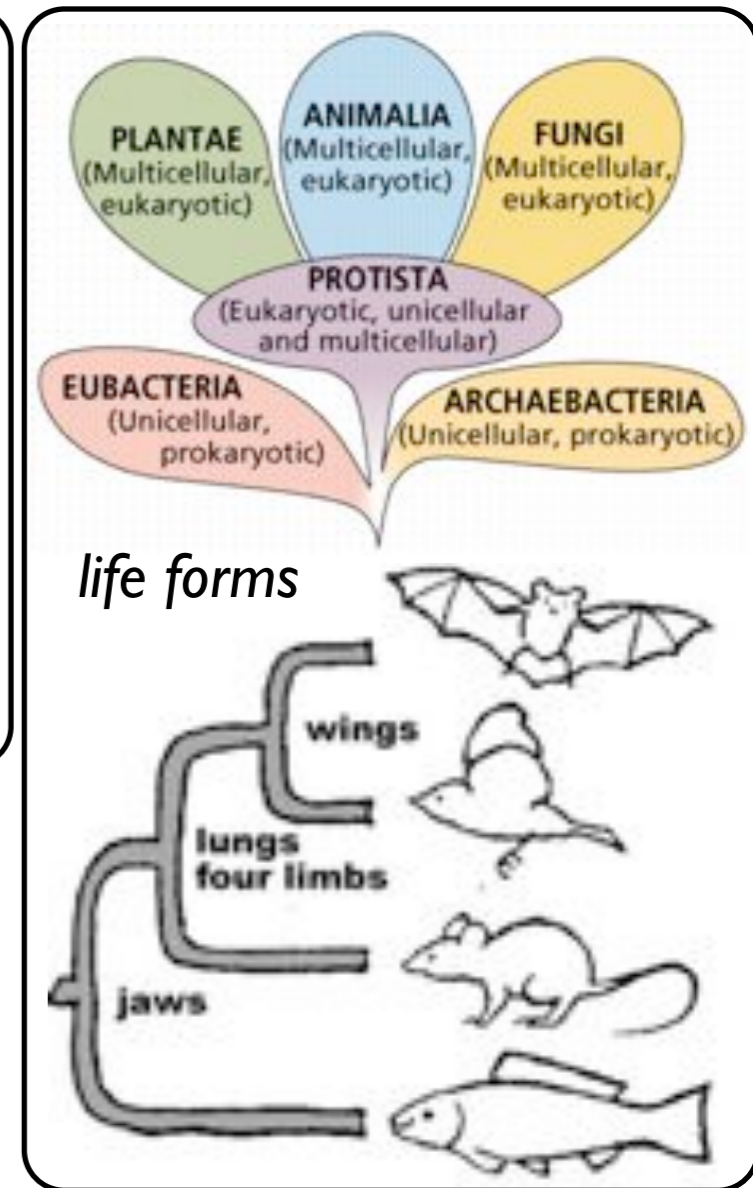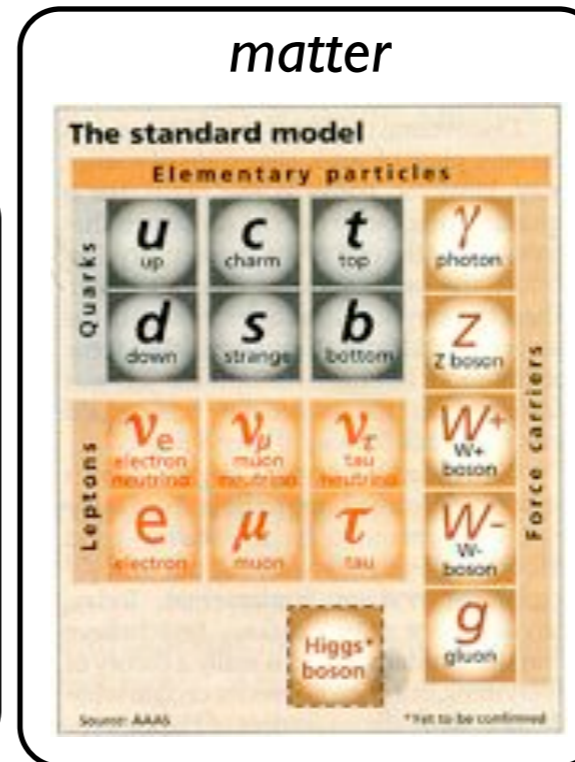Max Little, Nick Jones

May 2011

# Organizing

Scientific endeavors often focus on structuring libraries of collected information.


personalities


food we eat


matter
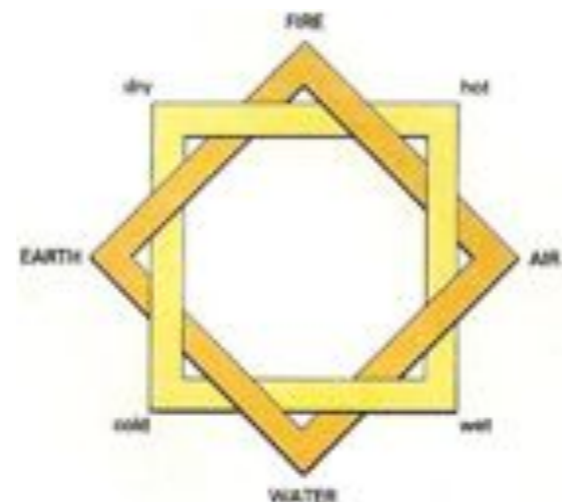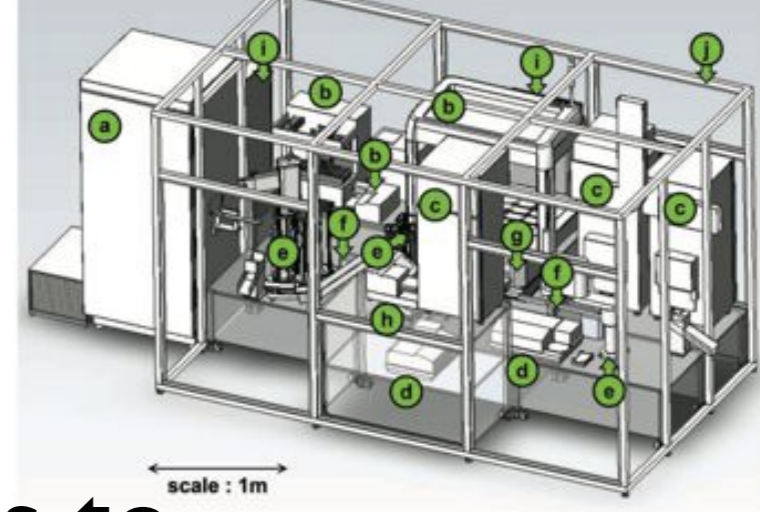

life forms


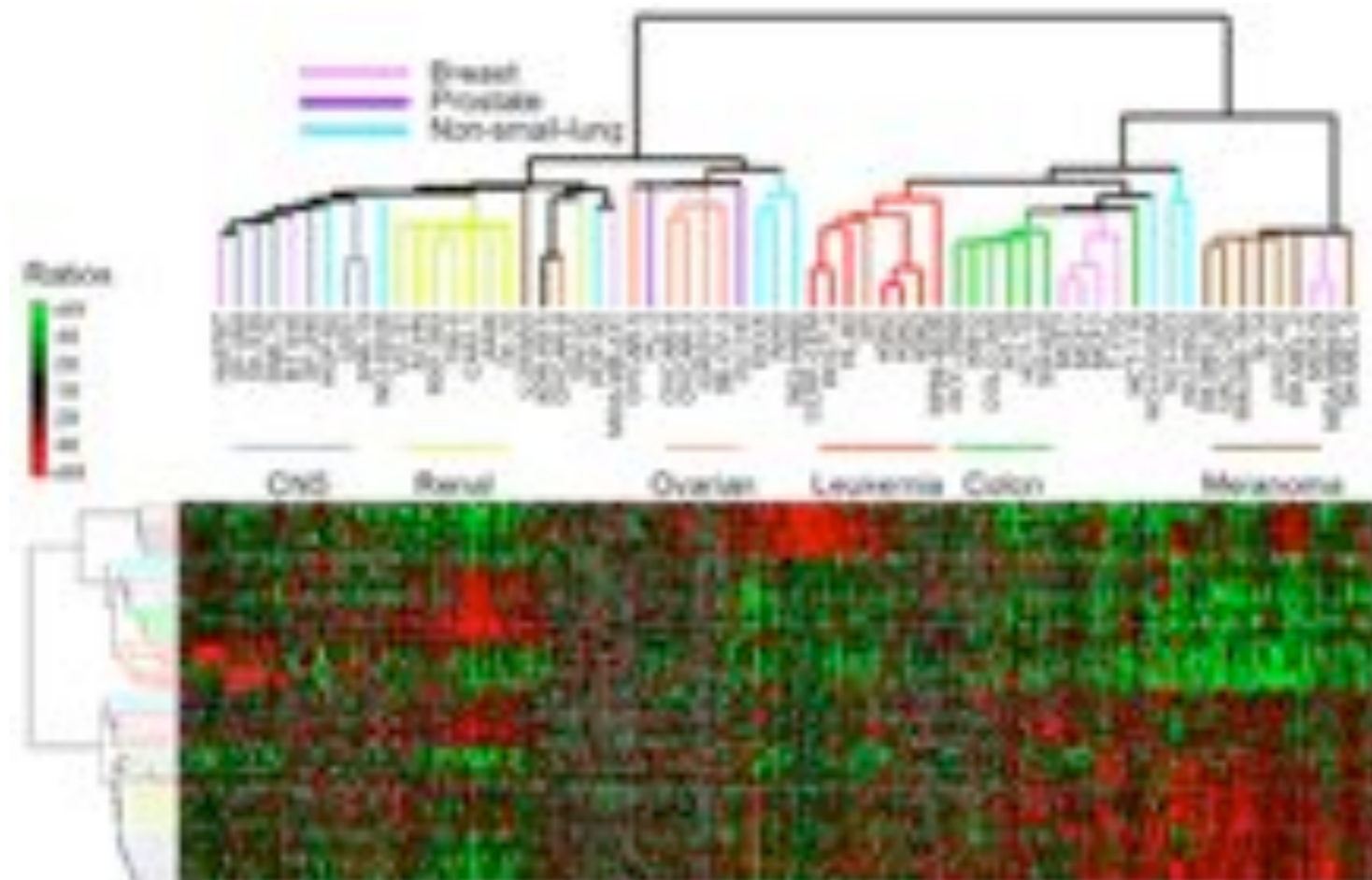macroscopic substances we encounter


color

*this helps us understand the complexity in the world*

# Organizing

Statistical learning techniques allow us to organize and understand things on a greater scale than ever before.

*e.g., the genetic microarray*

# Organizing Science

- Our data and our methods are also objects that require organization.

- How do we make sense of the time-series data that we observe in the world?

We construct a comparative framework for time-series analysis

# Challenges

Time Series Analysis:

- Pervasive importance in science

- Huge quantities of data

- Vast and growing quantity of methods
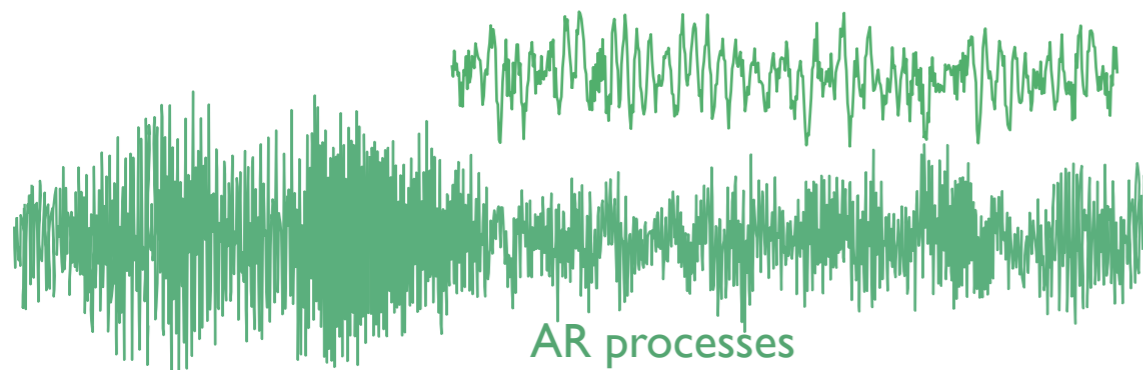
- Interdisciplinary boundaries

# Structure

- Framework

- Structure of methods for time-series analysis

- Structure of empirical time series

- Utility for specific applications
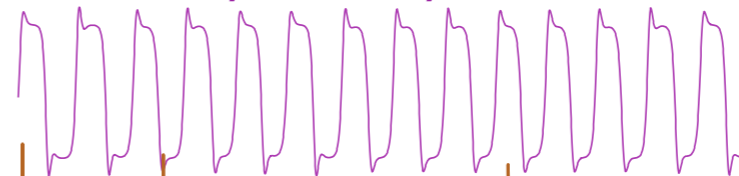
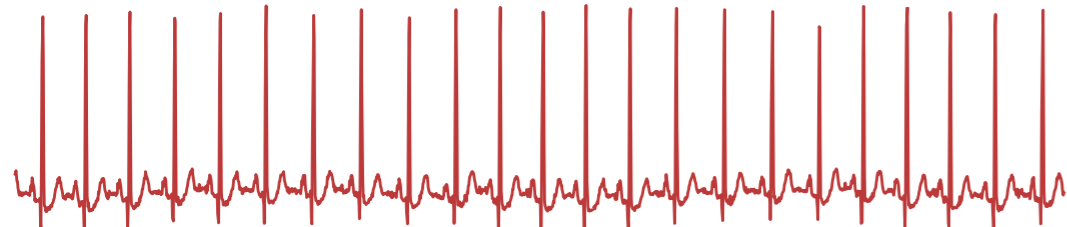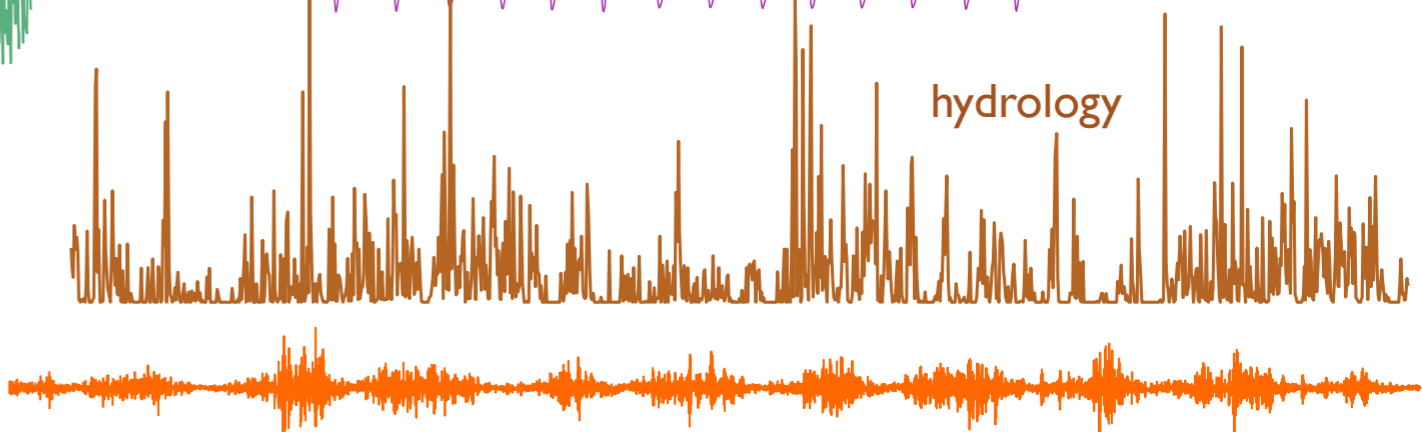- Constrained time-series datasets

# What time series?

> 30 000

medical CO₂ fluctuations

dynamical systems

AR processes

hydrology

medical: normal sinus rhythm
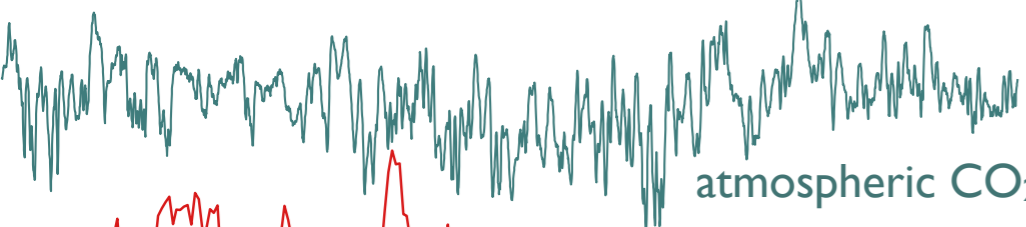
audio: brushing teeth

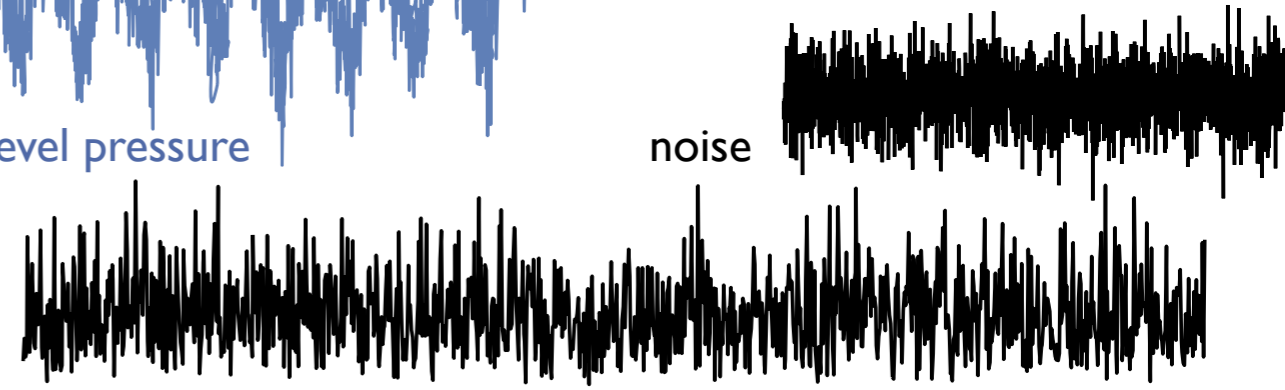finance: oil prices

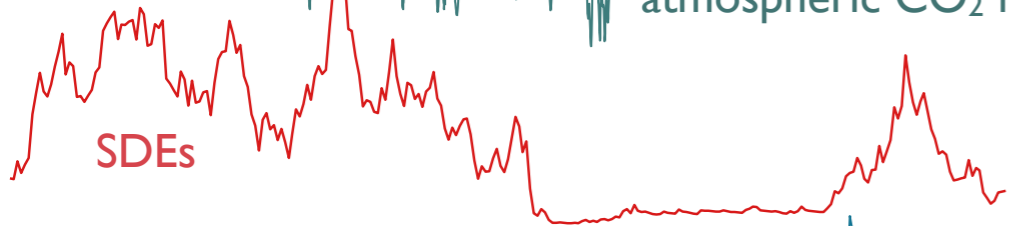text: sentence lengths

satellite position

climatology: sea level pressure

noise

atmospheric CO₂ fluctuations

SDEs

zooplankton growth

# What operations?

> 9 000

## Basic statistics

trimmed means     zero crossings

standard deviation

outliers     local extrema

## Stationarity

StatAv     sliding windows

bootstraps

distribution comparisons

## Static distribution

quantiles     moments

fits to standard distributions

hypothesis tests

## Basis Functions

wavelet transform

power spectrum peaks

spectral measures

low frequency power

## Correlation

linear autocorrelations     decay properties

automutual information

dependence on additive noise

nonlinear autocorrelations

time reversal asymmetry

generalized self-correlation function

recurrence structure

autocorrelation robustness

fluctuation analysis: scaling

randomization robustness

recurrence plots

seasonality testing

## Model fits

primitive forecasting

Fourier fits     GARCH modeling

step-ahead dependence

exponential smoothing     AR models

state space models

hidden Markov models

piecewise splines     'walker' statistics

ARMA modeling     Gaussian Processes

## Nonlinear

2D embedding structure     TSTOOL

TISEAN     fractal dimension

correlation dimension     Taken's estimator

Poincaré sections     surrogate data

nonlinear prediction error

Lyapunov exponent
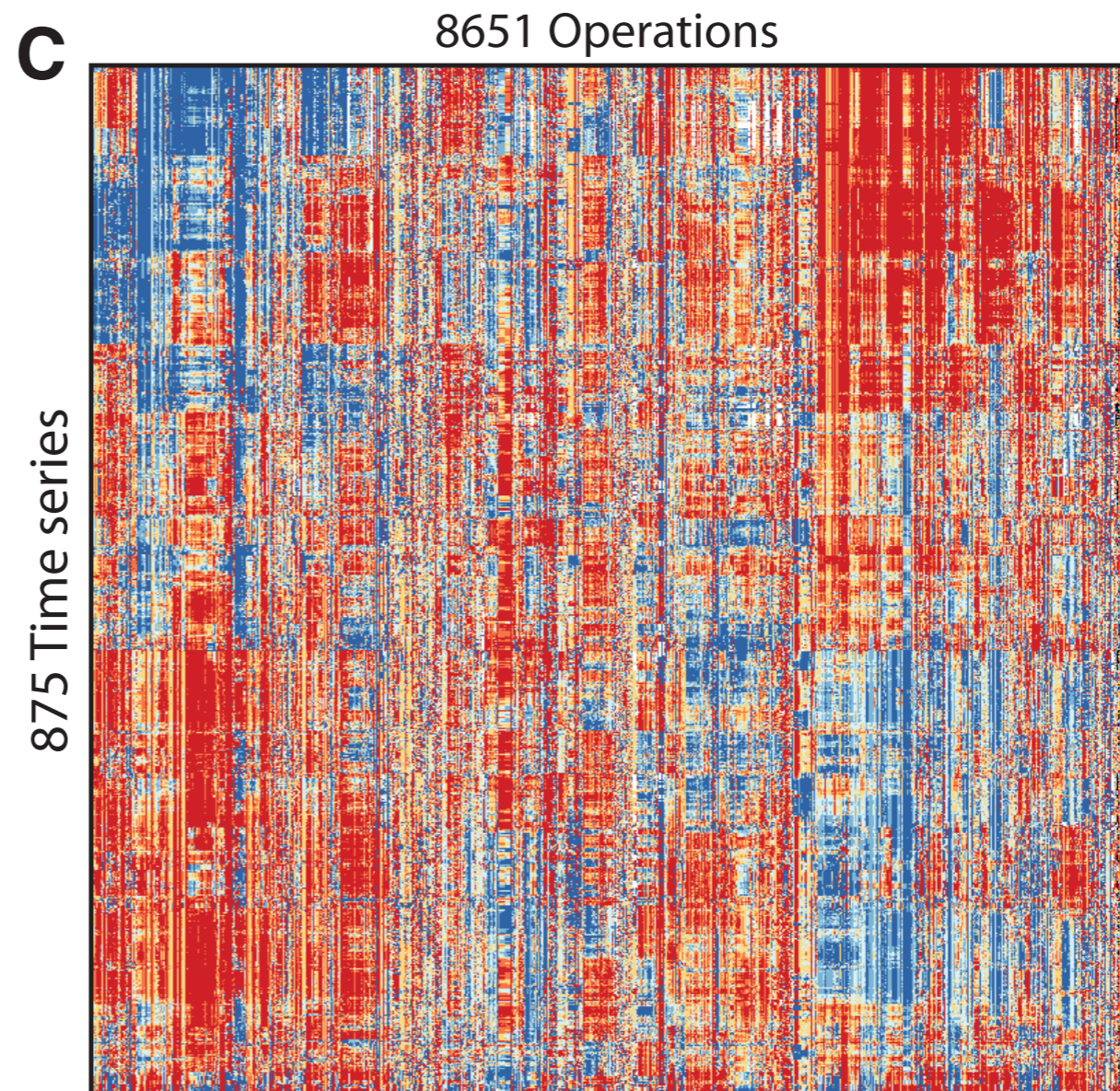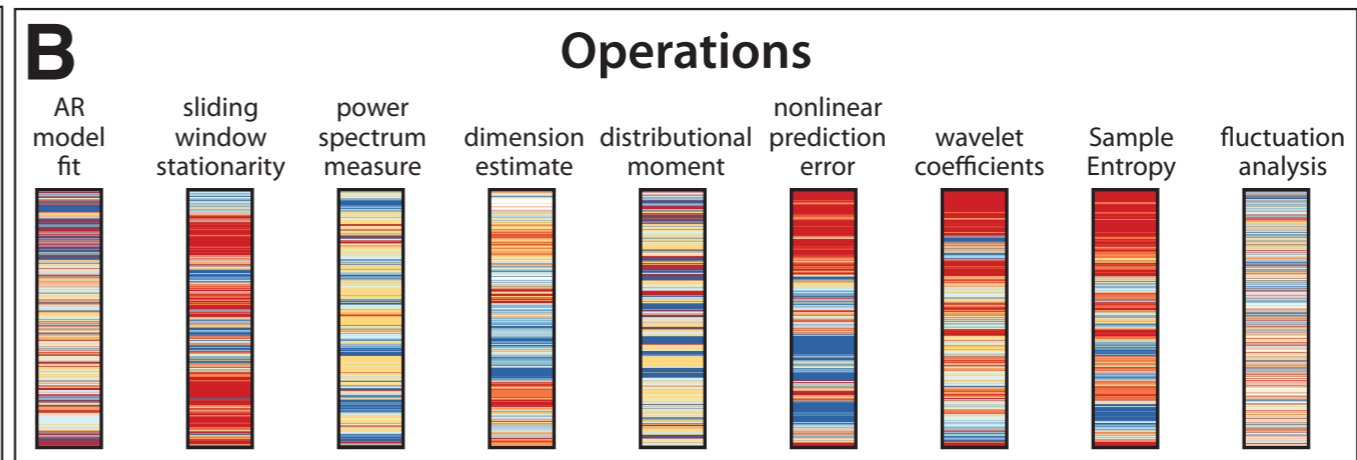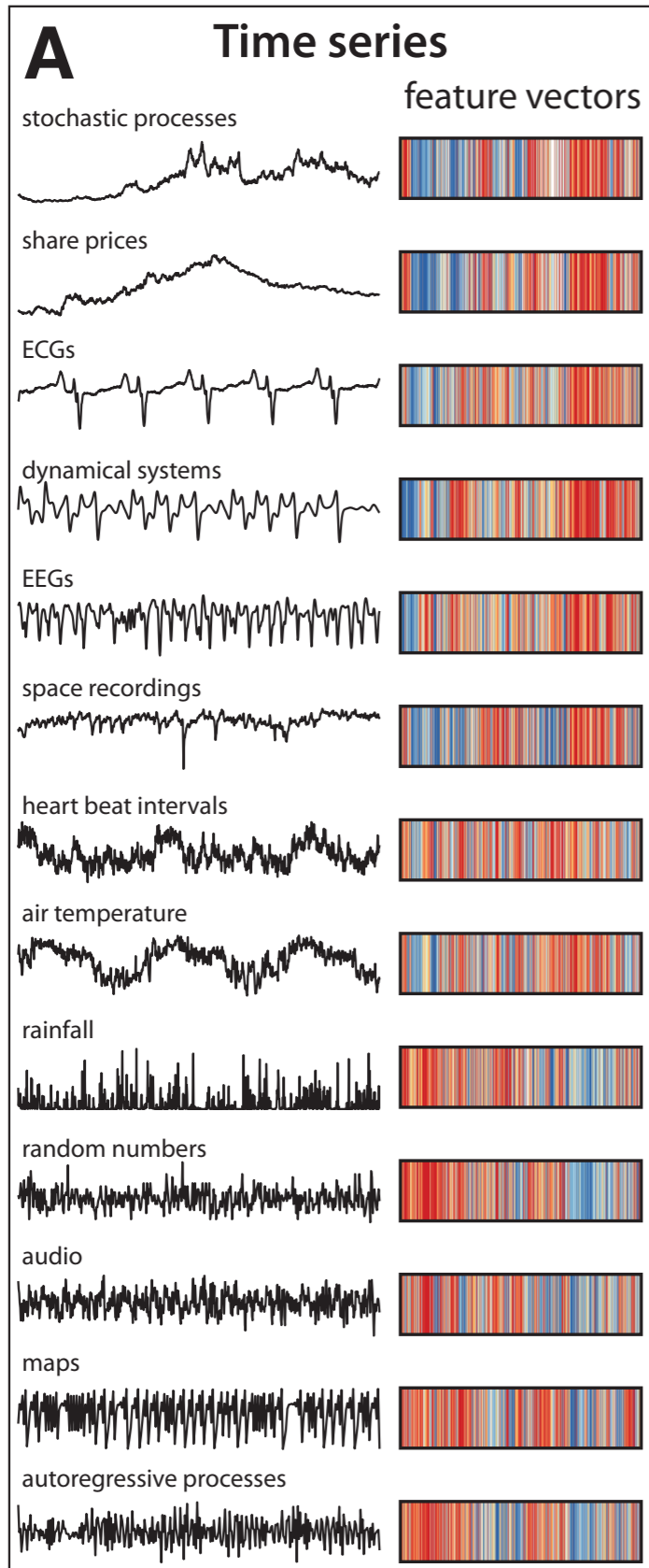
false nearest neighbours

## Others

course-grained transition matrices

motif distributions

couple to dynamical systems

visibility graph     stick angle distribution

step detection algorithms

extreme events     drifting mean tests

PCA of embedded signal

domain-specific standard metrics

## Information Theory

SampEn     distributional entropies

conditional entropies     binned entropies

kernel smoothed entropies

Tsallis entropies     ApEn

# Design Matrix



**A** Time series

stochastic processes

share prices

ECGs

dynamical systems

EEGs

space recordings

heart beat intervals

air temperature

rainfall

random numbers

audio

maps

autoregressive processes

feature vectors

**B** Operations

AR model fit

sliding window stationarity

power spectrum measure

dimension estimate

distributional moment

nonlinear prediction error

wavelet coefficients

Sample Entropy

fluctuation analysis

**C** 8651 Operations

875 Time series

# Organizing Our Methods

- We organize operations using their outputs on a diverse range of empirical time series.

- Clustering allows us to form reduced sets of operations that capture the dominant types of behavior in our database.

- Gives structure to an interdisciplinary field.

long-range scaling

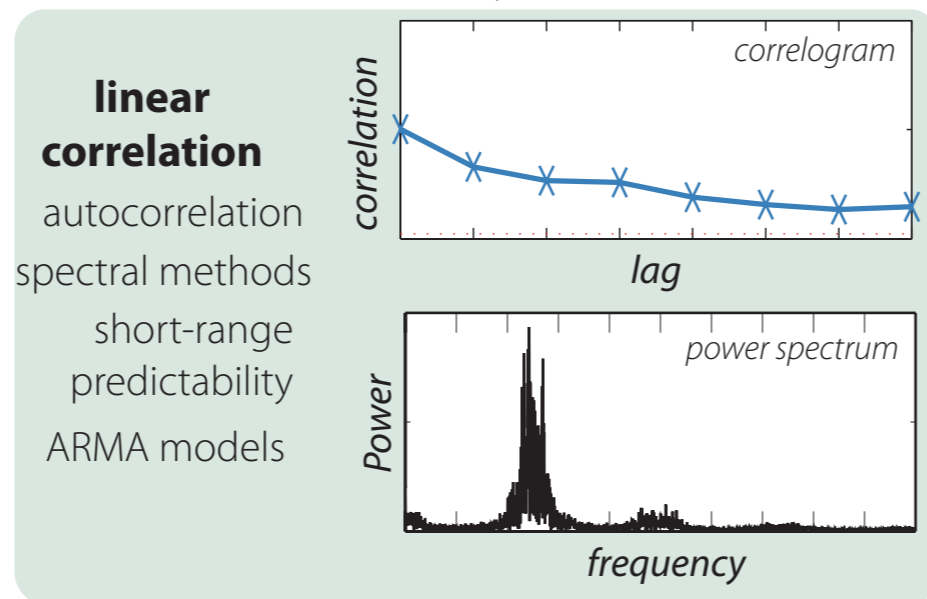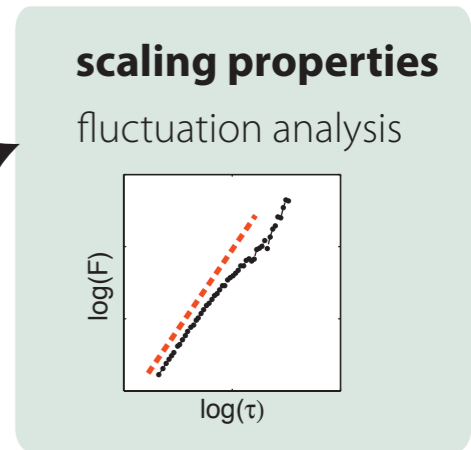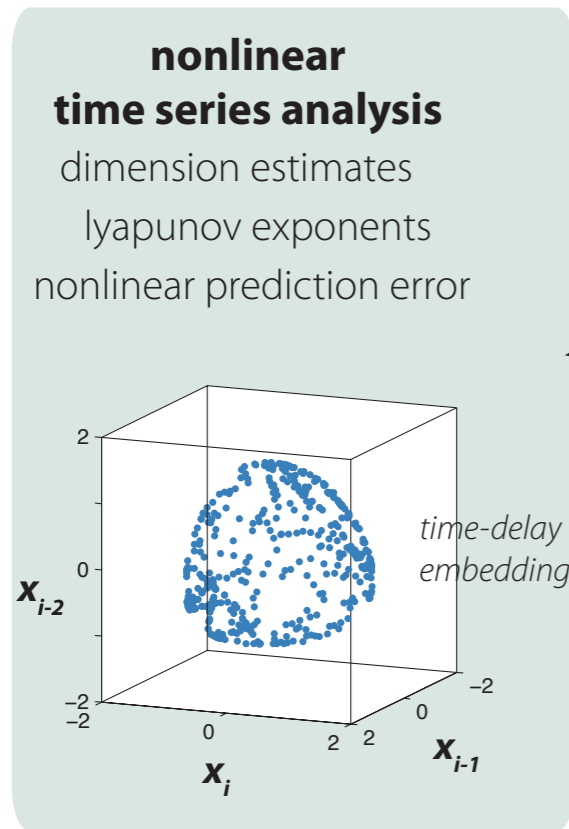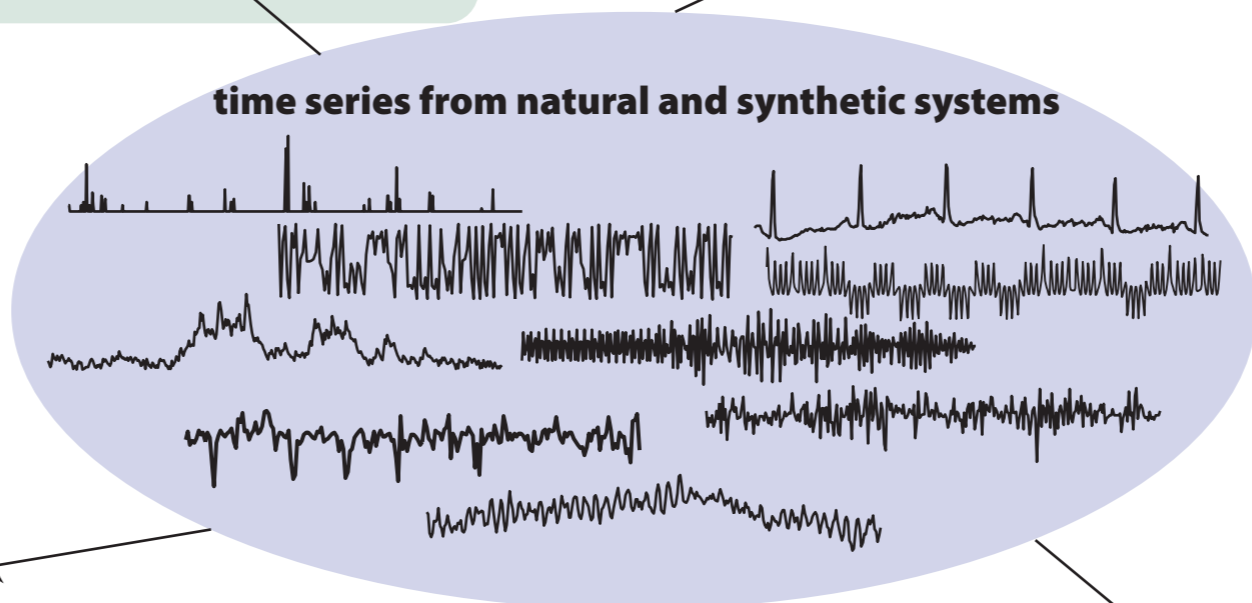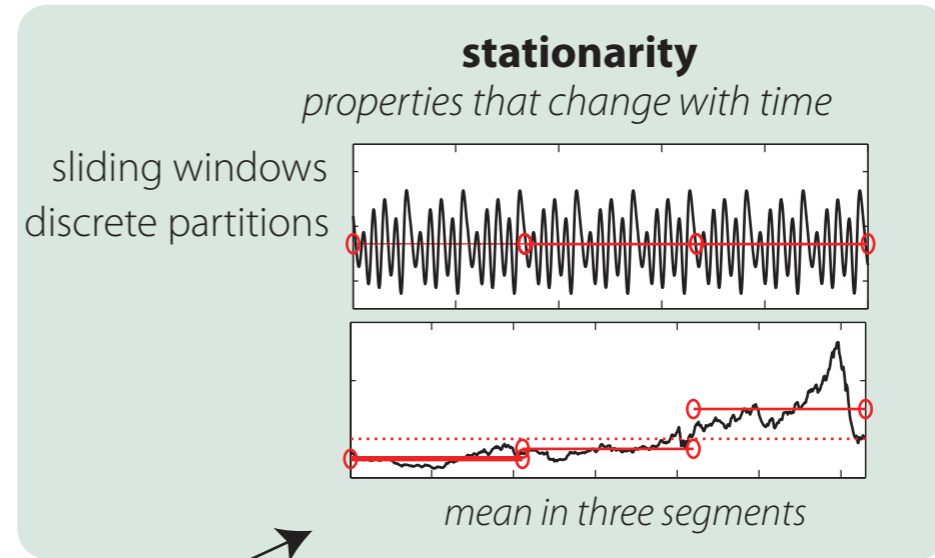power spectral density

linear models

stationarity

variance

entropy

correlation dimension

complexity

information theory

BIG PICTURE

**distribution**
*time-independent structure*

*kernel-smoothed distribution*

*histogram*

spread
location
outliers
distribution entropy
higher order moments

**stationarity**
*properties that change with time*

sliding windows
discrete partitions

*mean in three segments*

**time series from natural and synthetic systems**

**scaling properties**

fluctuation analysis

**nonlinear time series analysis**

dimension estimates
lyapunov exponents
nonlinear prediction error

*time-delay embedding*

**linear correlation**

autocorrelation
spectral methods
short-range predictability
ARMA models

*correlogram*

*power spectrum*

**Information theoretic**

auto mutual information
entropy

*symbolic entropy*

B A C B C A A A C B C A B B C C C A A B

$p(A|A), p(A|B), p(A|C)$

*How many operations are needed to efficiently summarize the structure we observe in empirical signals?*



*200 operations provide an efficient and interpretable summary*

# Local Neighborhoods



Automutual Information

Shannon Entropy

Approximate Entropy

Lempel-Ziv Complexity

*ApEn(2,0.2)*

Randomized Sample Entropy

Sample Entropy

*Organize our methods for time-series analysis*

# Local Neighborhoods

# Local Neighborhoods



*Organize our methods for time-series analysis*

Automutual Information Measures

# Visualize behavior



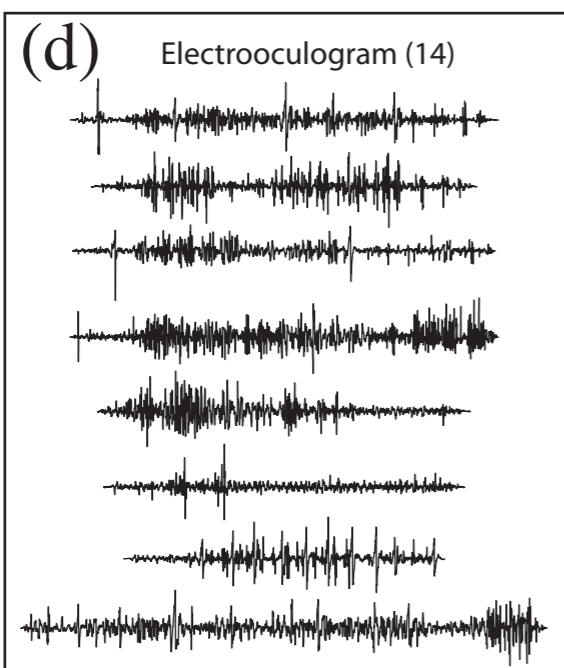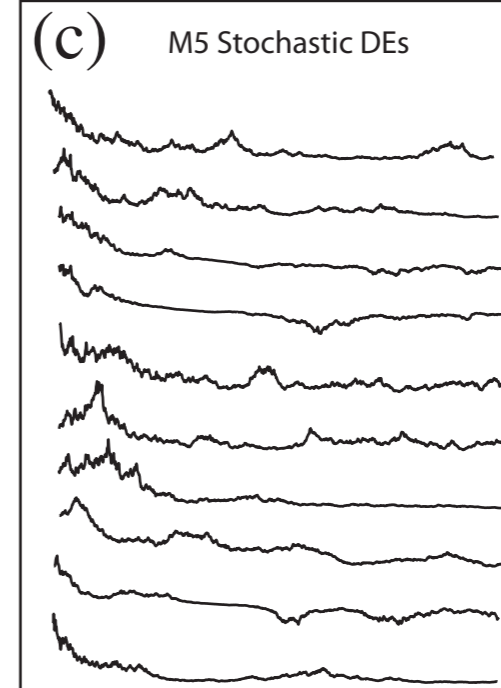Probability Density vs. ApEn(2,0.1) Outputs
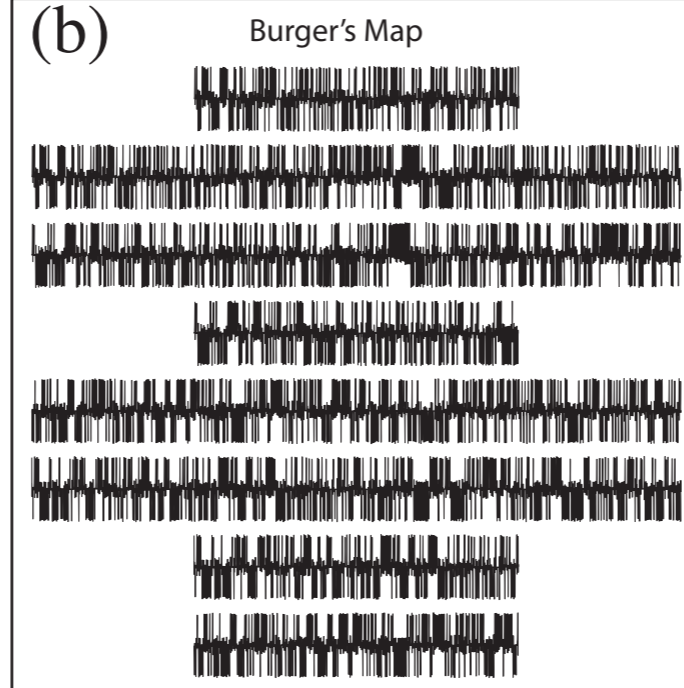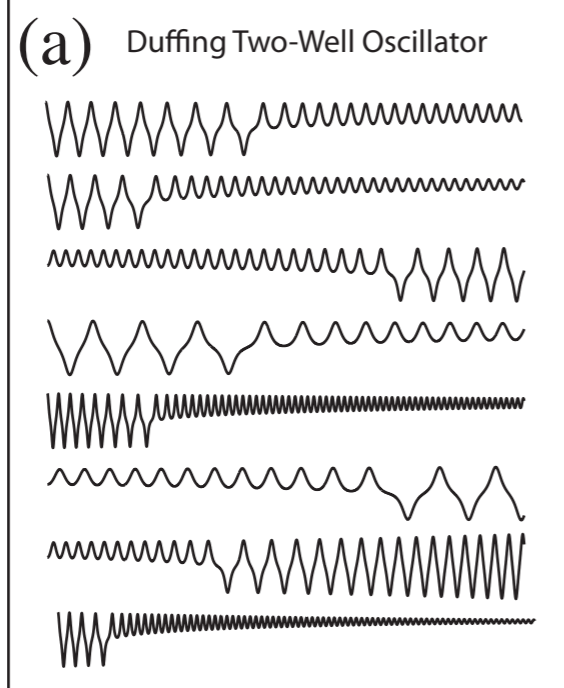
# Organizing Our Data

- Our reduced interdisciplinary set of operations is a powerful summary of the structure in empirical time series

- Links between real-world and synthetically-generated time series encourage a unified, collaborative framework for understanding the dynamics in time series
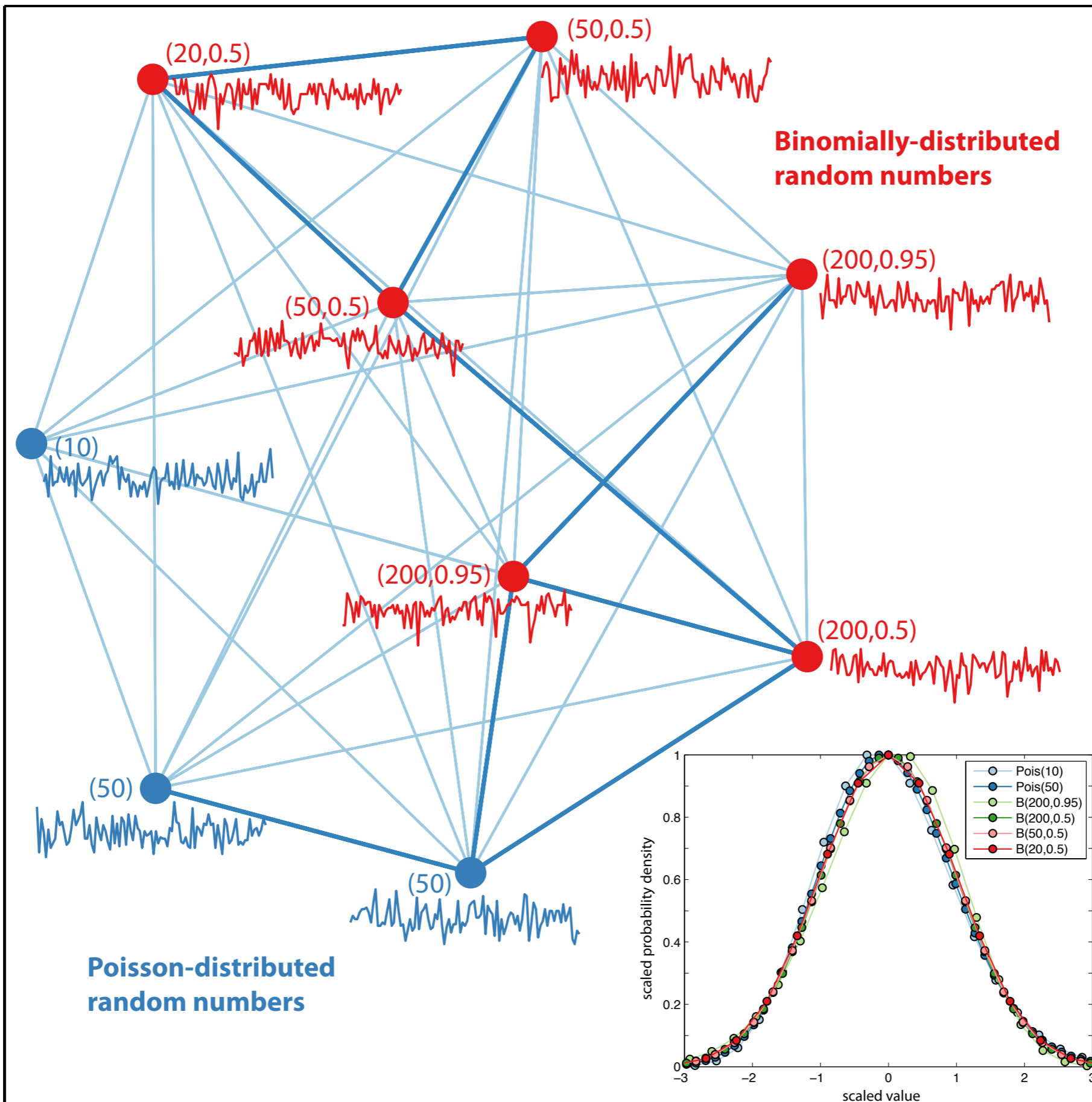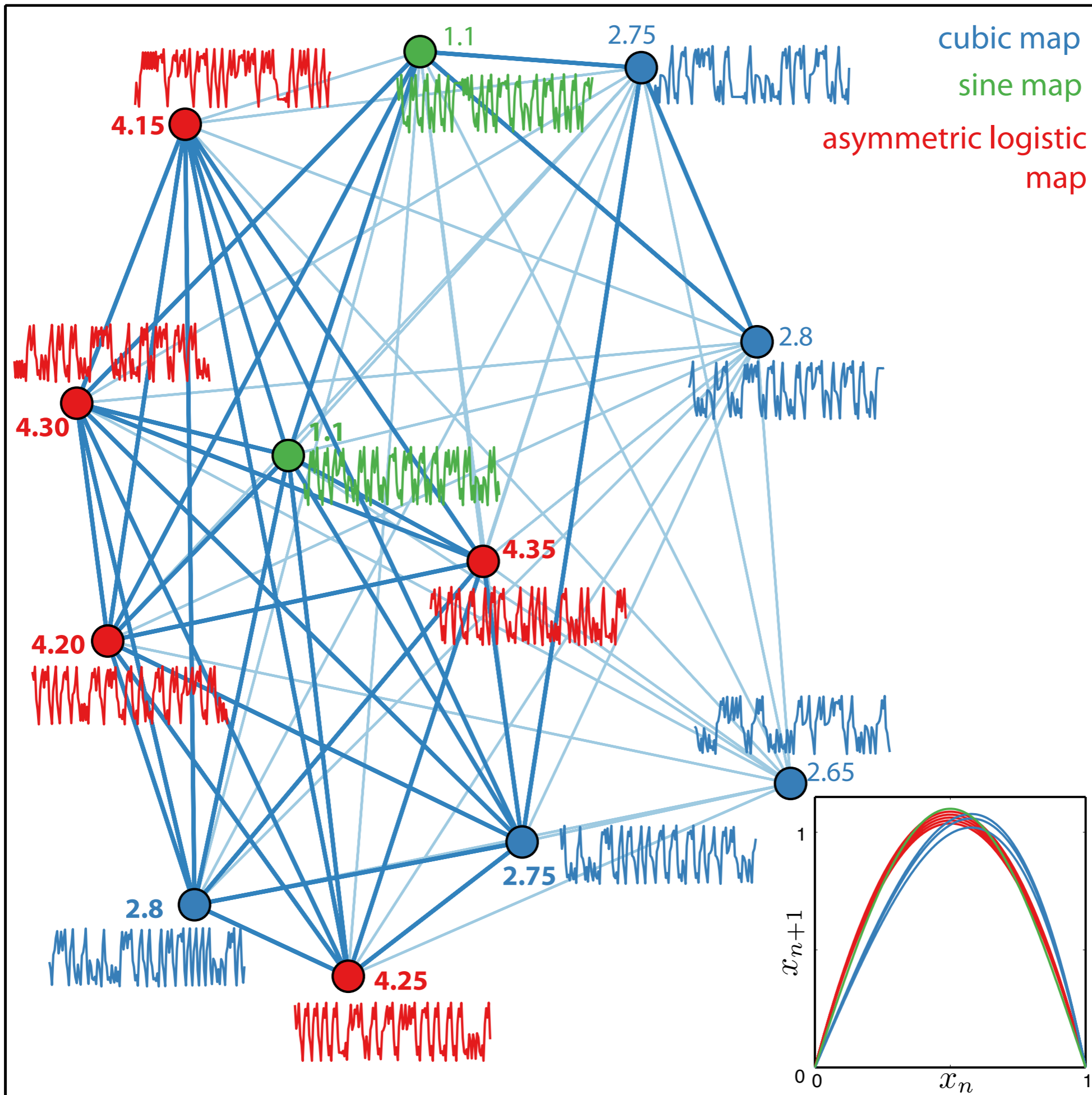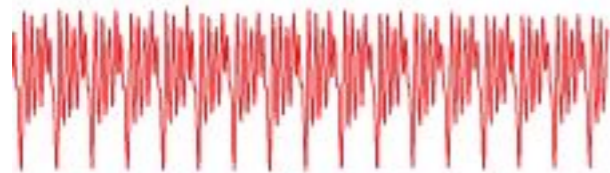
broad classes of dynamics are distinguished

194 operations

birdsong recordings

human speech

output from SDEs

heart beat intervals

air temperature

EEGs

normal sinus rhythm ECGs

rainfall

uncorrelated noise

log returns of financial time series

10 PCs (67%)
9 PCs (65.6%)
8 PCs (63.9%)
7 PCs (61.9%)
6 PCs (59.6%)
5 PCs (56.7%)
4 PCs (51.3%)
3 PCs (45%)
2 PCs (37.5%)
1 PC (25%)

random numbers
stochastic DEs
EEGs
ECGs
heart beat intervals
rainfall
airtemp
financial log returns
speech
birdsong

(a) Duffing Two-Well Oscillator

(b) Burger's Map

(c) M5 Stochastic DEs

(d) Electrooculogram (14)

(e) Speech (19)

(f) Airway $CO_2$

(g) Space: Power Index (15)

(h) Congestive Heart Failure ECGs (25)

(i) Music

(20,0.5)

(50,0.5)

**Binomially-distributed random numbers**

(200,0.95)

(50,0.5)

(10)

(200,0.95)

(200,0.5)

(50)

(50)

**Poisson-distributed random numbers**

scaled probability density

1

0.8

0.6

0.4

0.2

0

-3  -2  -1  0  1  2  3

scaled value

- Pois(10)
- Pois(50)
- B(200,0.95)
- B(200,0.5)
- B(50,0.5)
- B(20,0.5)

# Time-series models?

- Real data are recordings of dynamics

- Time-series models generate dynamics with a known mechanism



speech recording

$$dx/dt = y$$
$$dy/dt = z$$
$$dz/dt = -z - (T - R + Rx^2)y - Tx$$

*"he did it"*

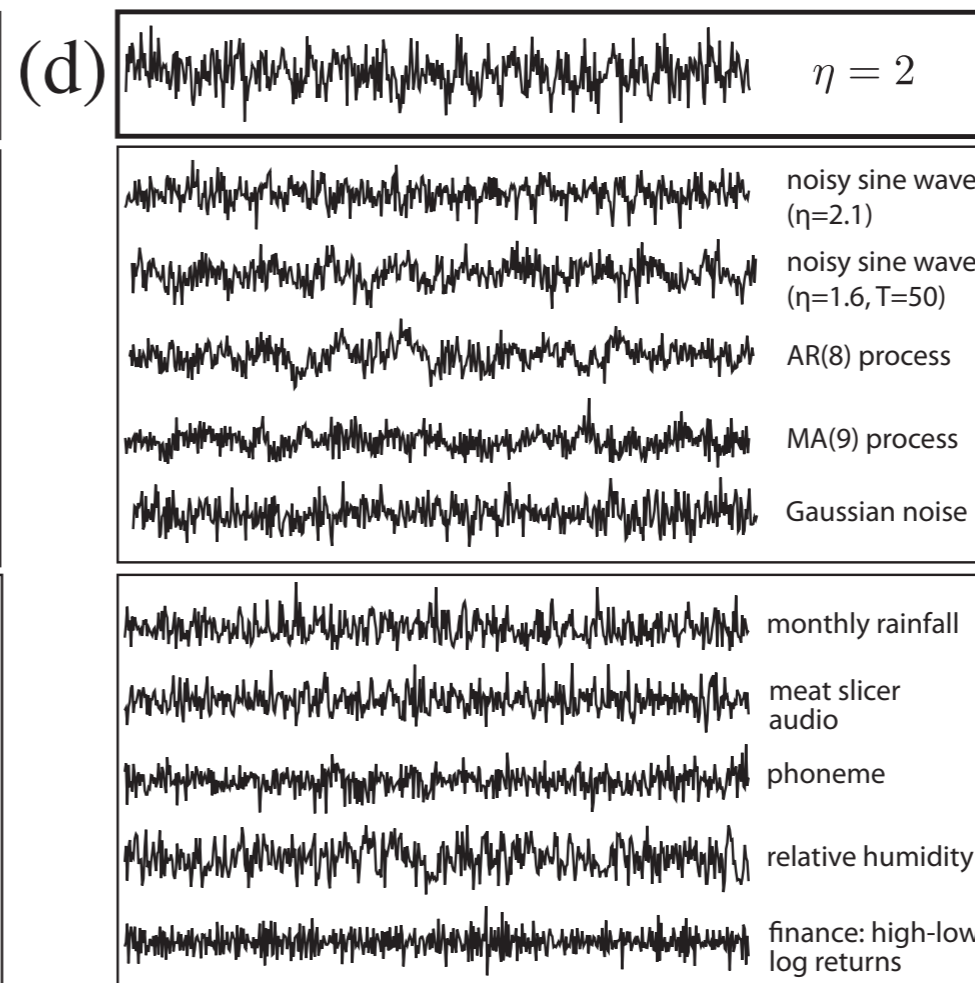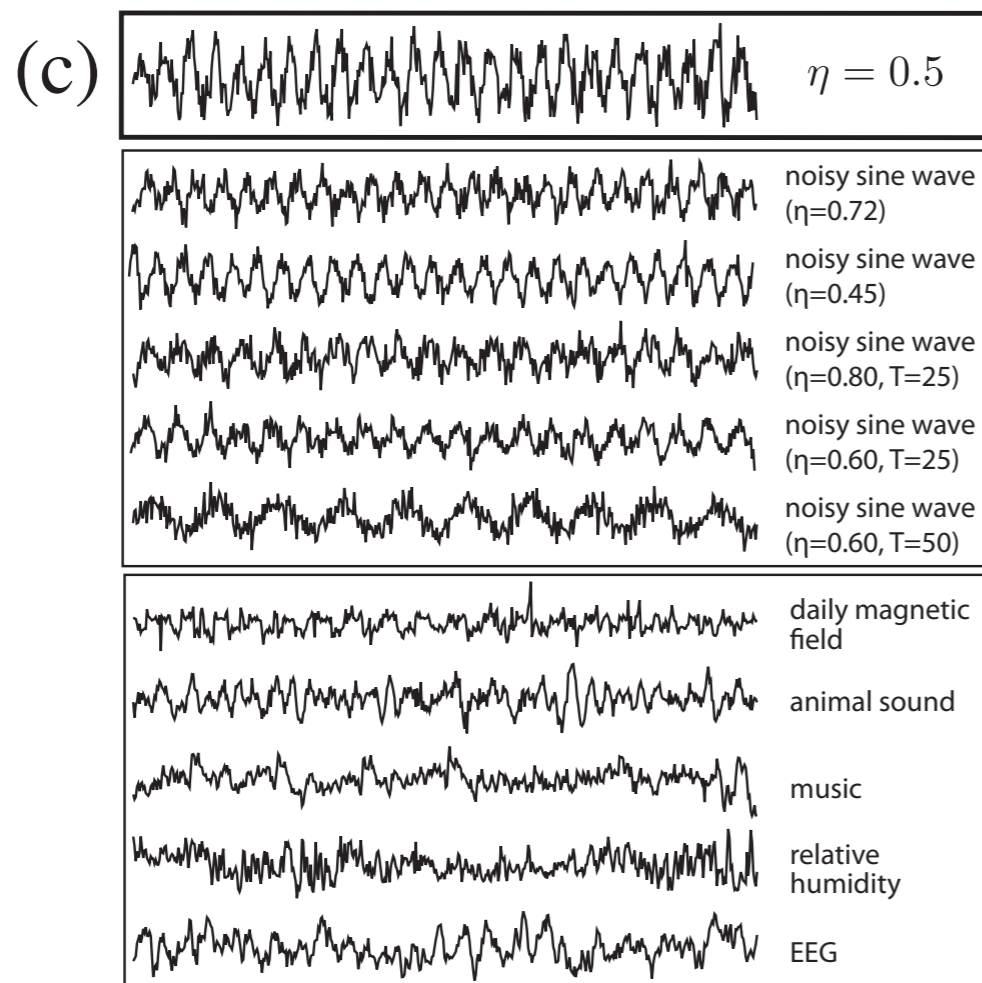*"he did it"*

*Pointing the finger*

# Similarity Search

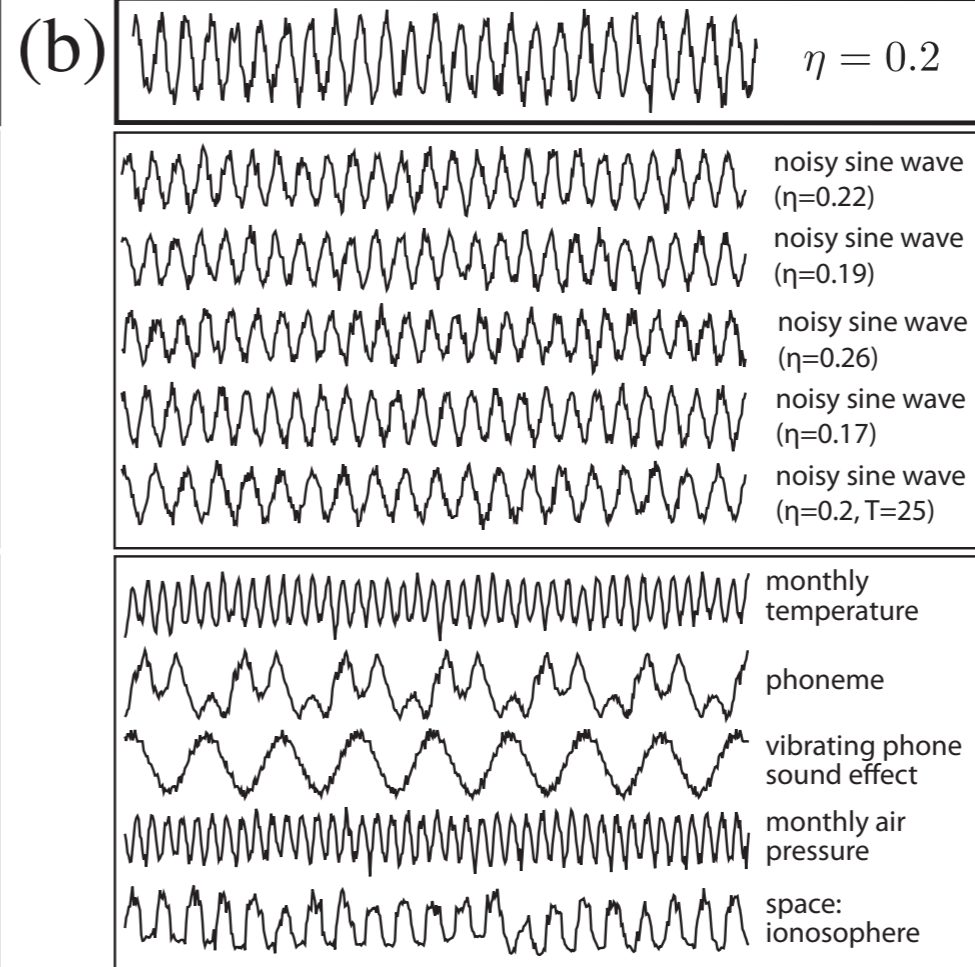# Similarity Search



*suggest models
for our data*

(a)

(b) **stochastic sine map target**

*synthetic matches:*     *real-world matches:*
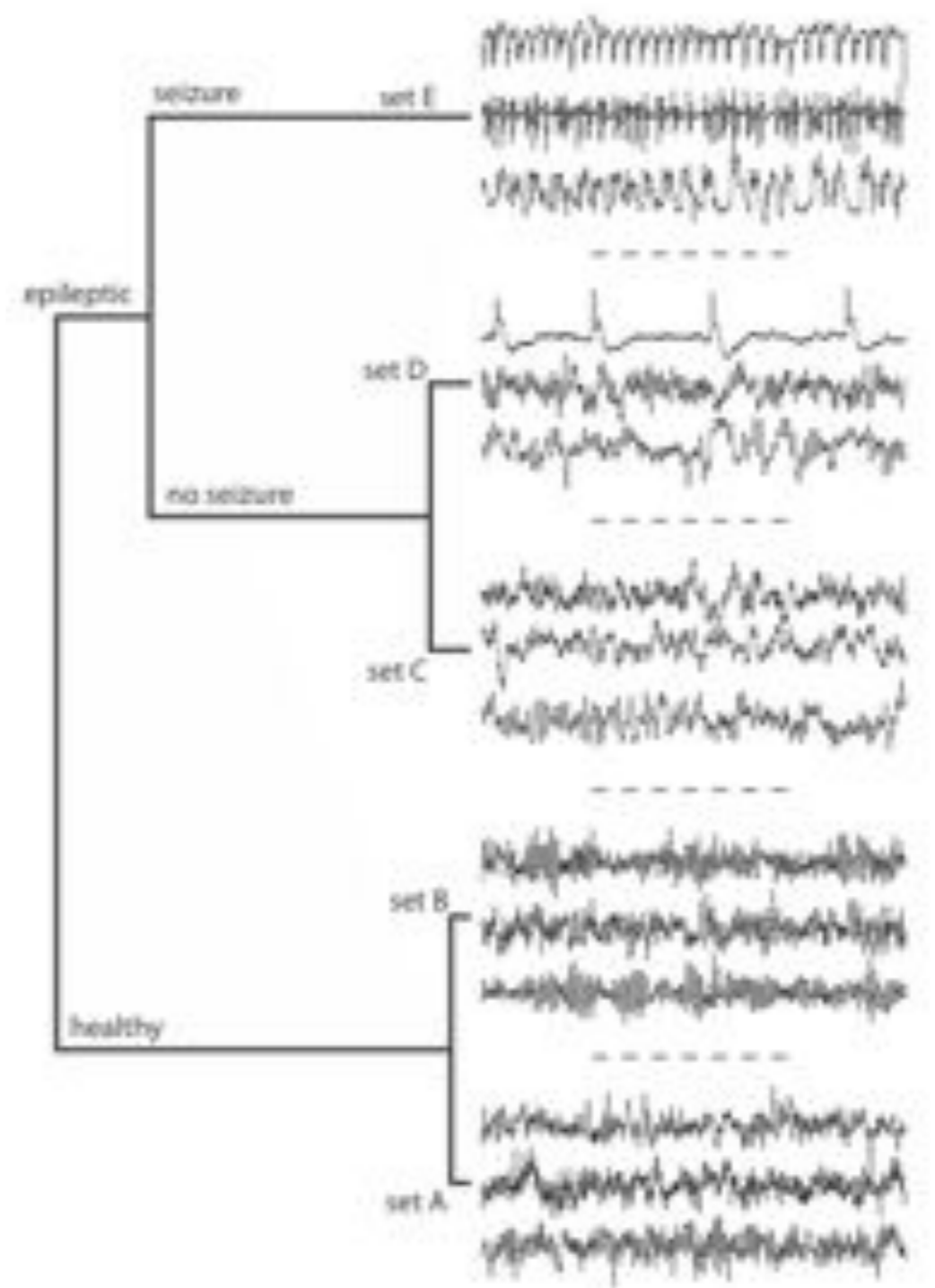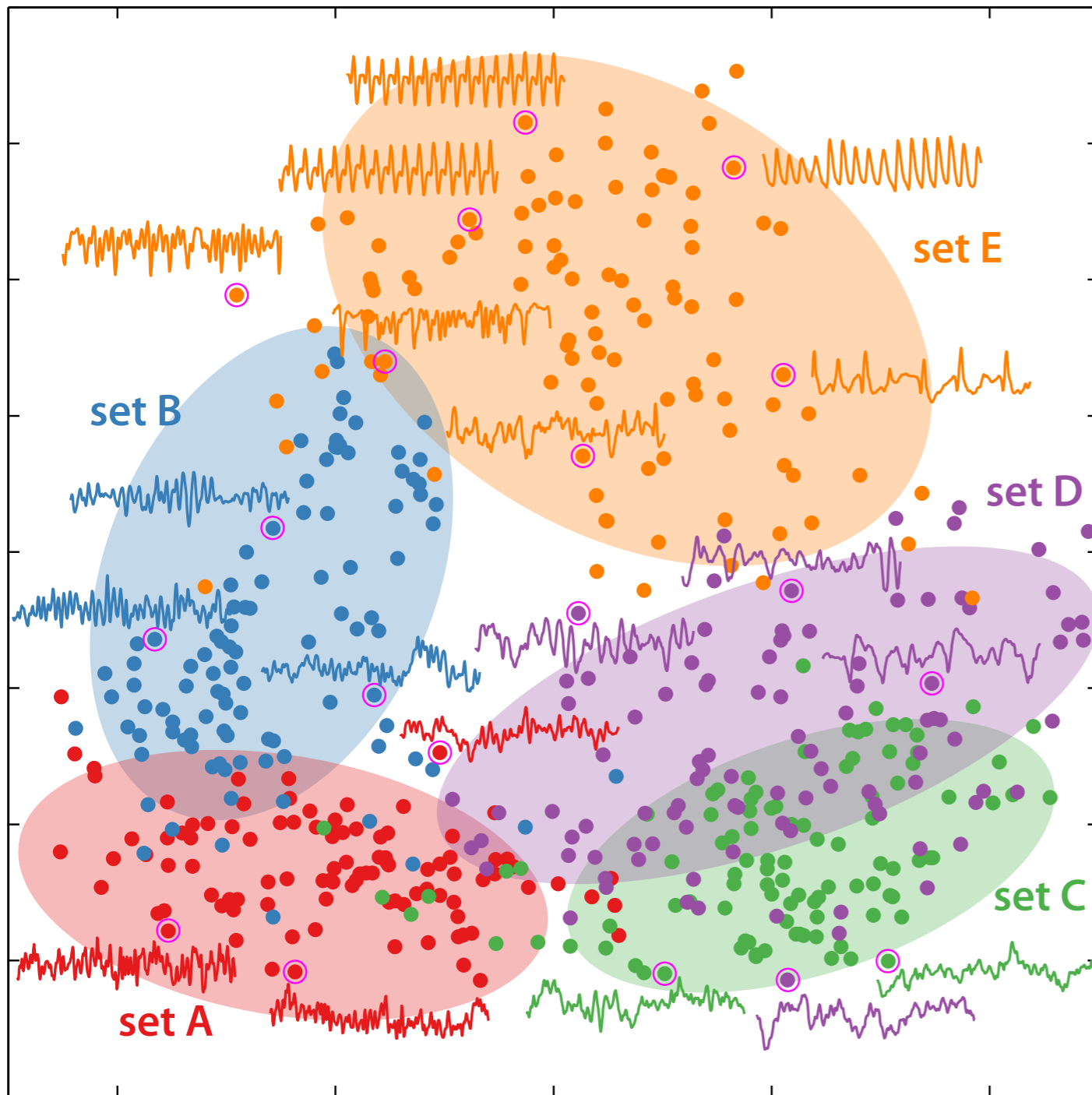
cloud amount

weather types

wind direction

low cloud cover

wind speed

*suggest data for our models*

**(a)** $\eta = 0$

sine wave (longer period)
Jerk system
Hadley flow
sine wave (longer period)
Rössler flow

speech phoneme
satellite position
speech phoneme
EEG: seizure
EEG: seizure

**(b)** $\eta = 0.2$

noisy sine wave ($\eta = 0.22$)
noisy sine wave ($\eta = 0.19$)
noisy sine wave ($\eta = 0.26$)
noisy sine wave ($\eta = 0.17$)
noisy sine wave ($\eta = 0.2$, $T = 25$)

monthly temperature
phoneme
vibrating phone sound effect
monthly air pressure
space: ionosophere

**(c)** $\eta = 0.5$

noisy sine wave ($\eta = 0.72$)
noisy sine wave ($\eta = 0.45$)
noisy sine wave ($\eta = 0.80$, $T = 25$)
noisy sine wave ($\eta = 0.60$, $T = 25$)
noisy sine wave ($\eta = 0.60$, $T = 50$)

daily magnetic field
animal sound
music
relative humidity
EEG

**(d)** $\eta = 2$

noisy sine wave ($\eta = 2.1$)
noisy sine wave ($\eta = 1.6$, $T = 50$)
AR(8) process
MA(9) process
Gaussian noise

monthly rainfall
meat slicer audio
phoneme
relative humidity
finance: high-low log returns

# Applications

- Drawing on a rich, interdisciplinary database of methods for time-series analysis allows datasets to be analyzed in new ways

- Reveal structure using PCA

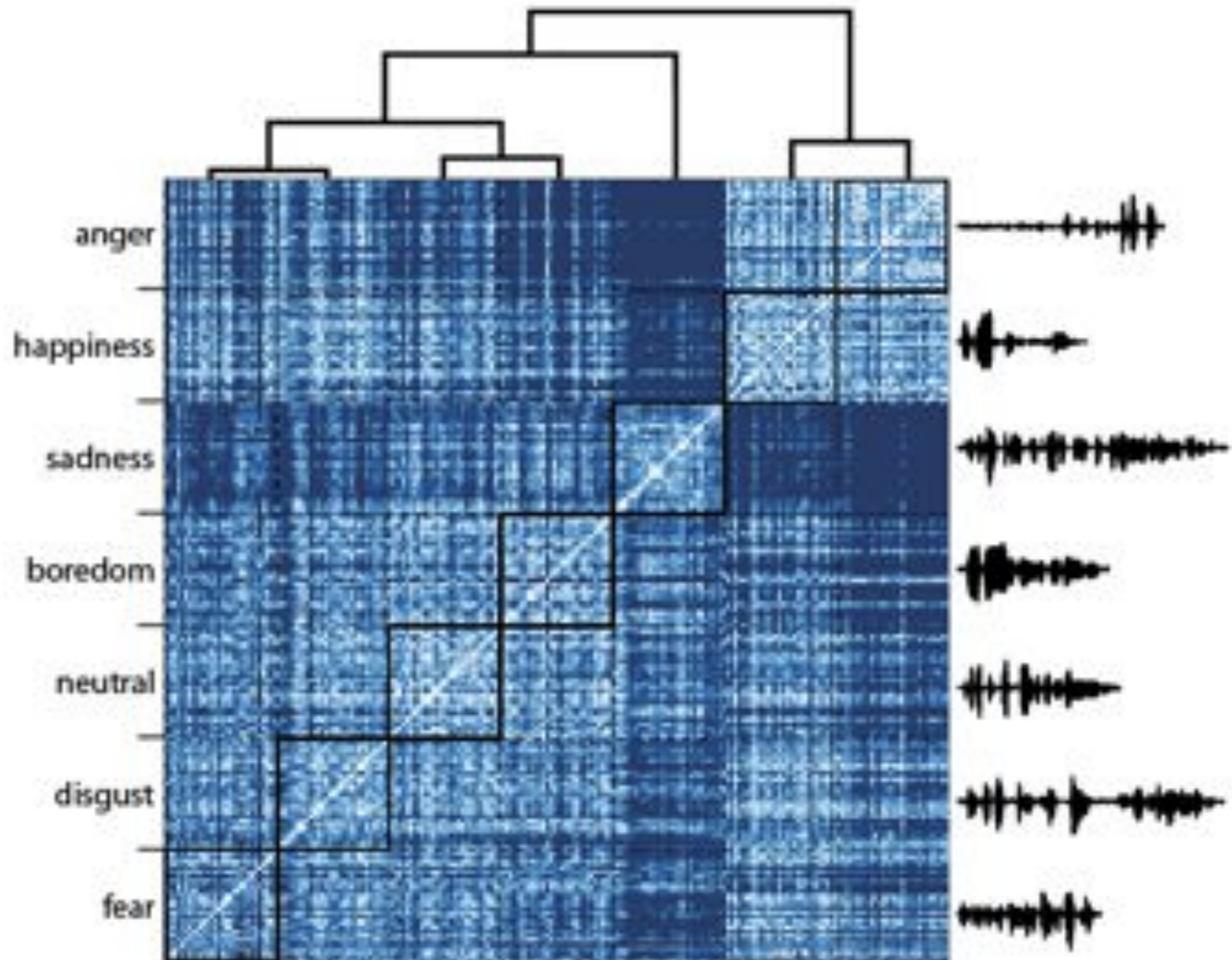- Select interpretable families of useful methods for a given classification/regression task
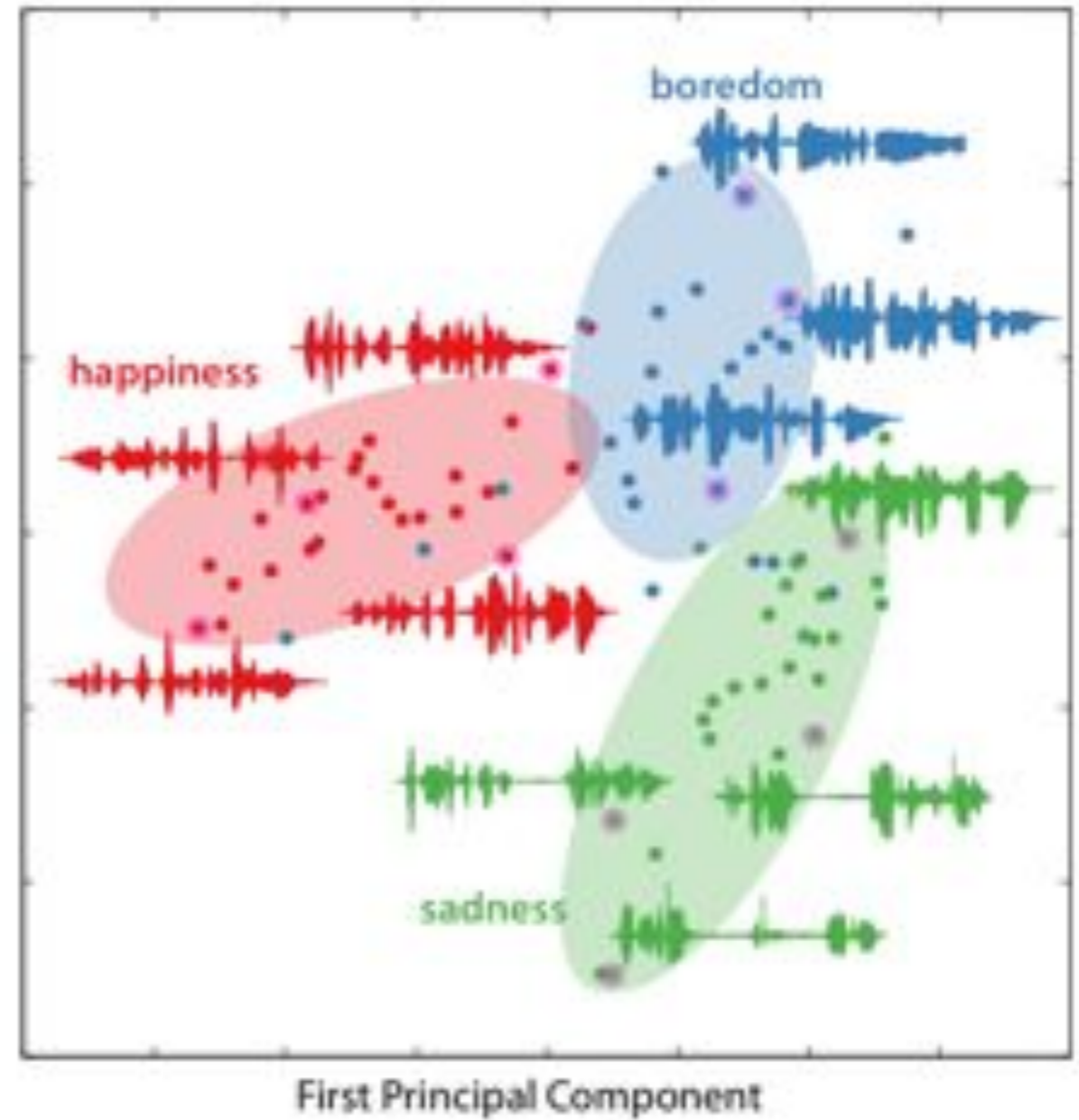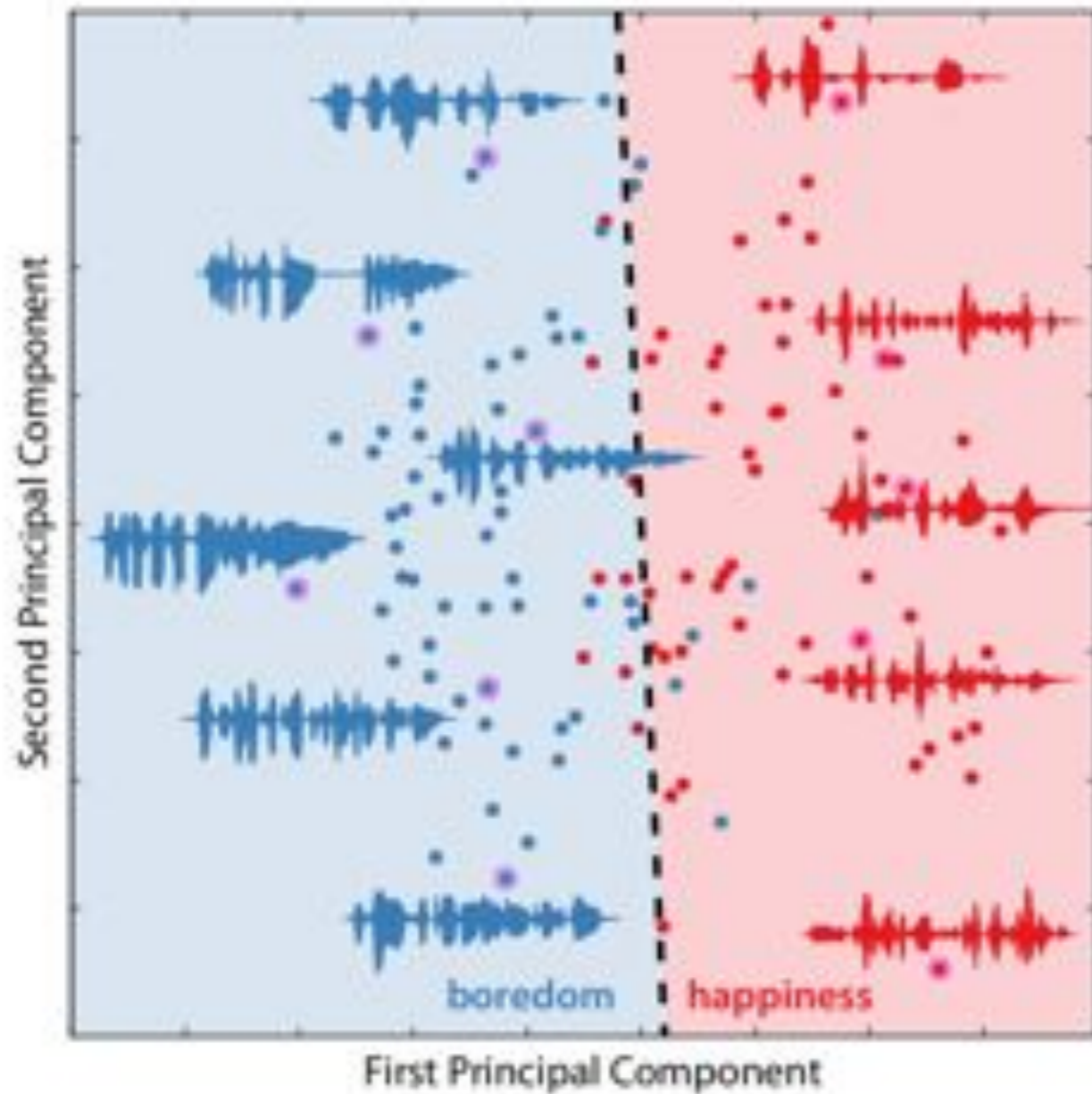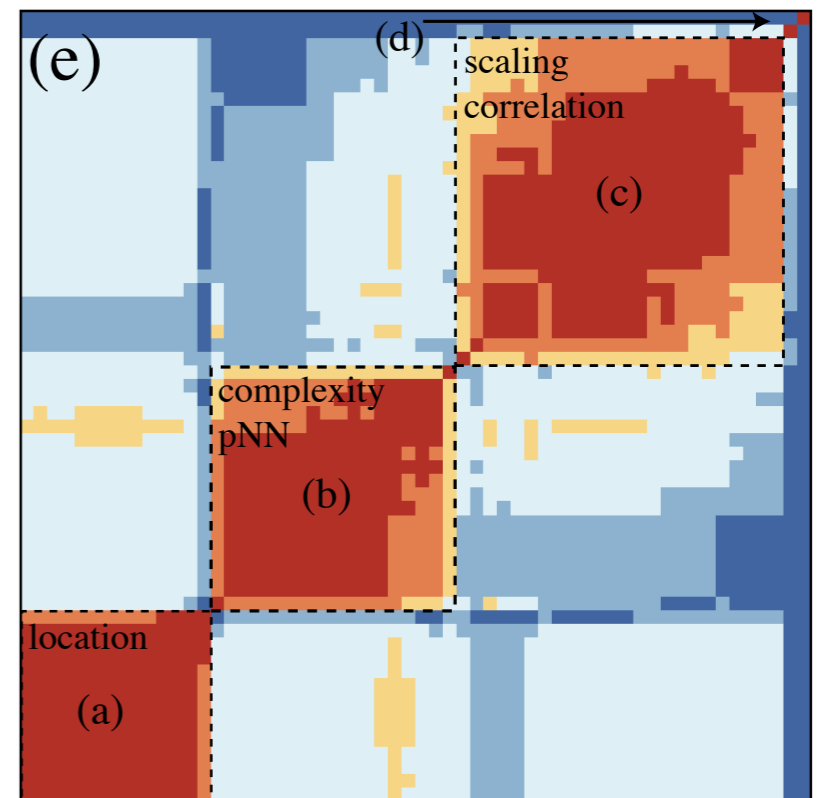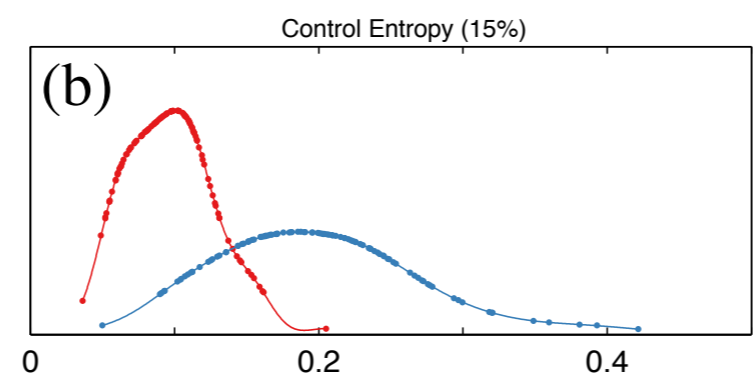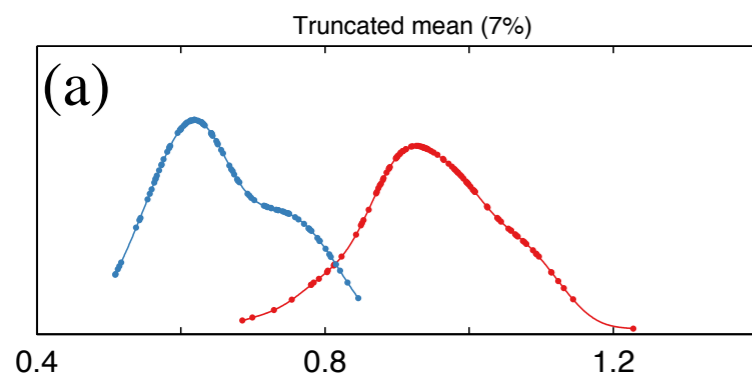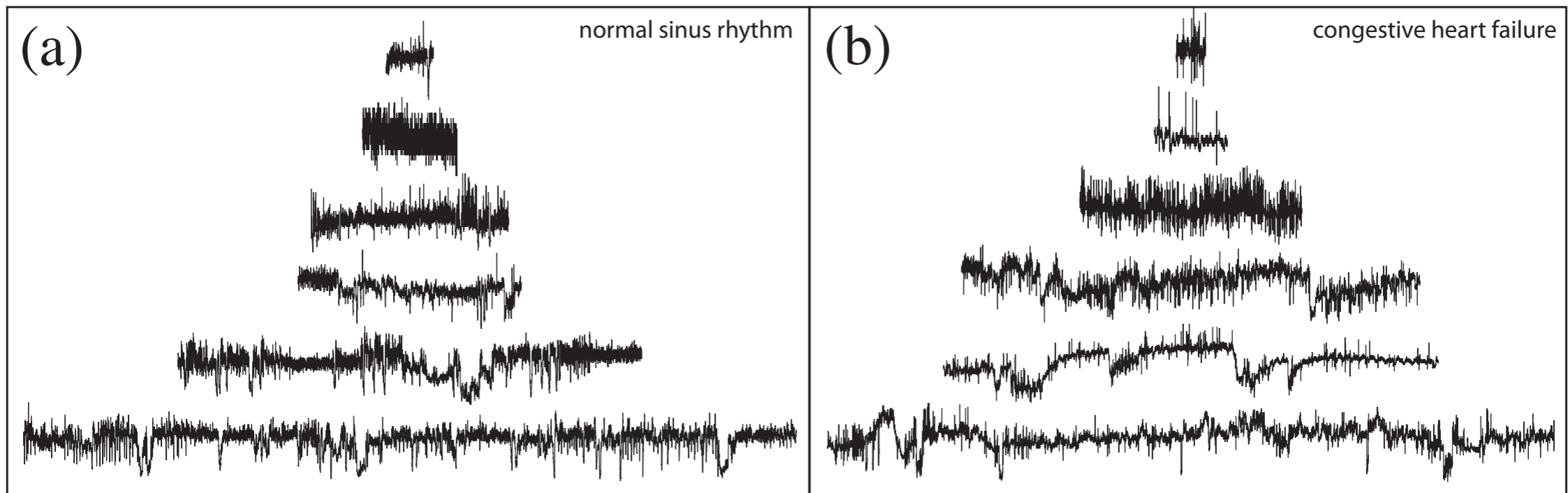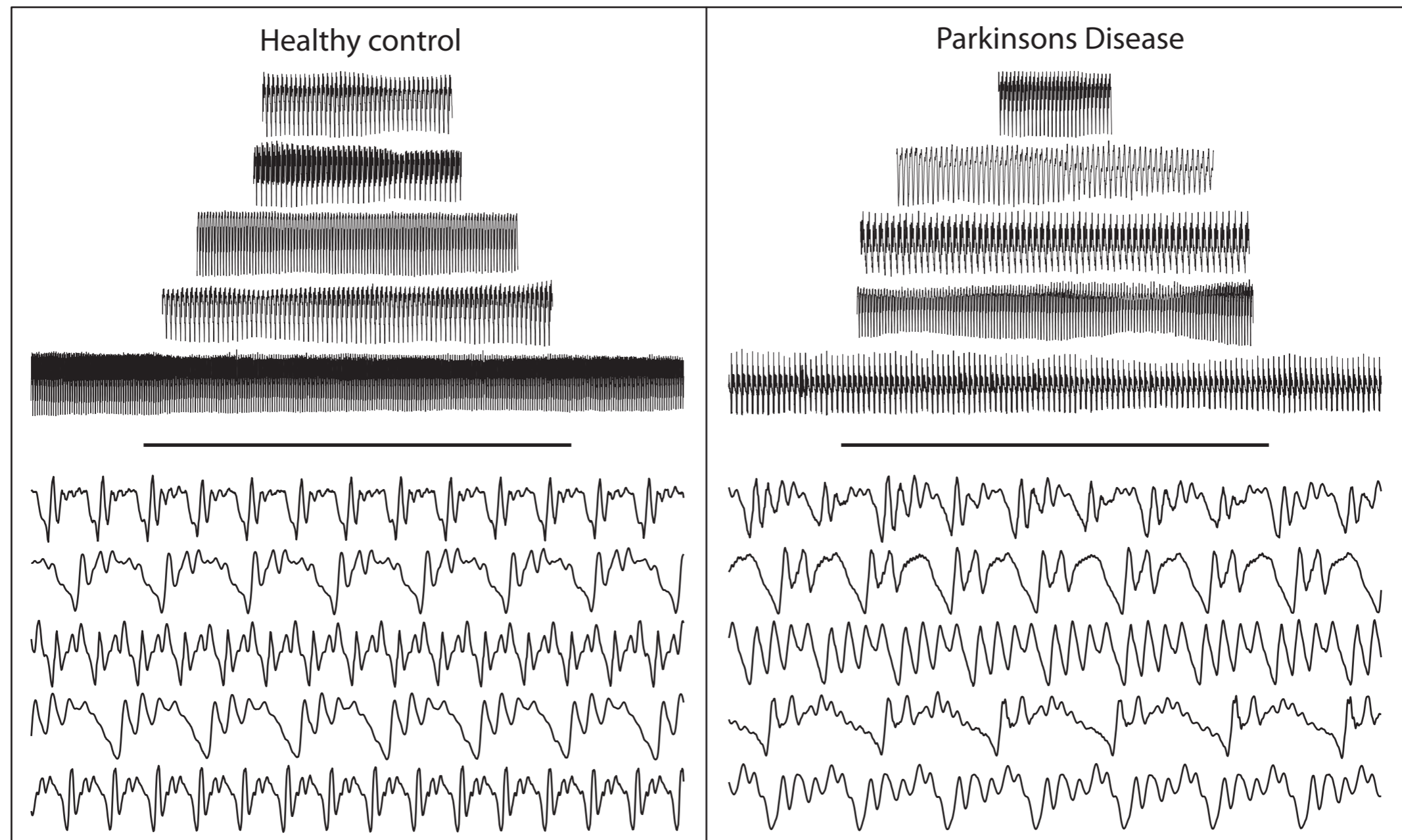
# EEGs

# Distinguishing seizures
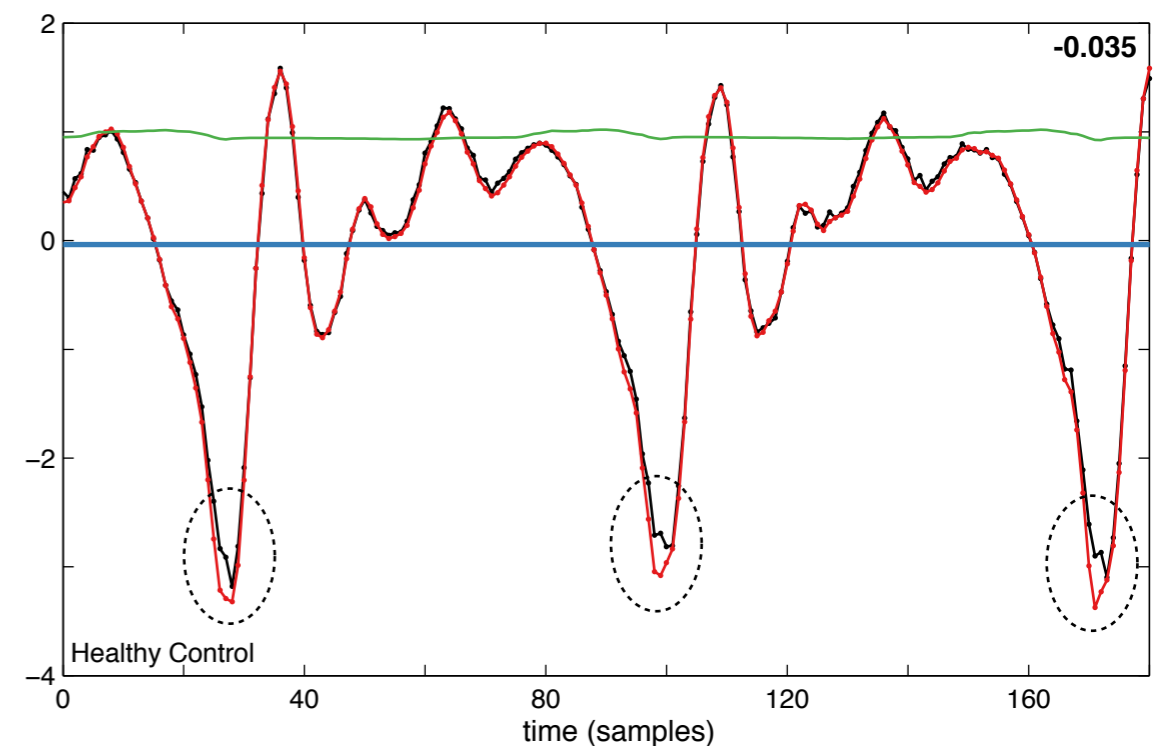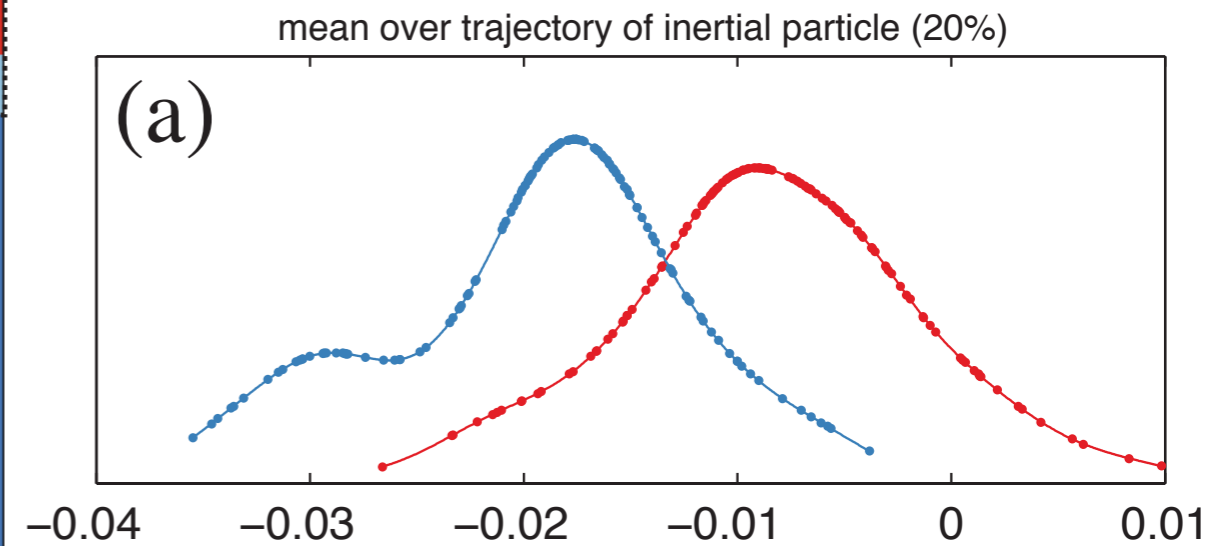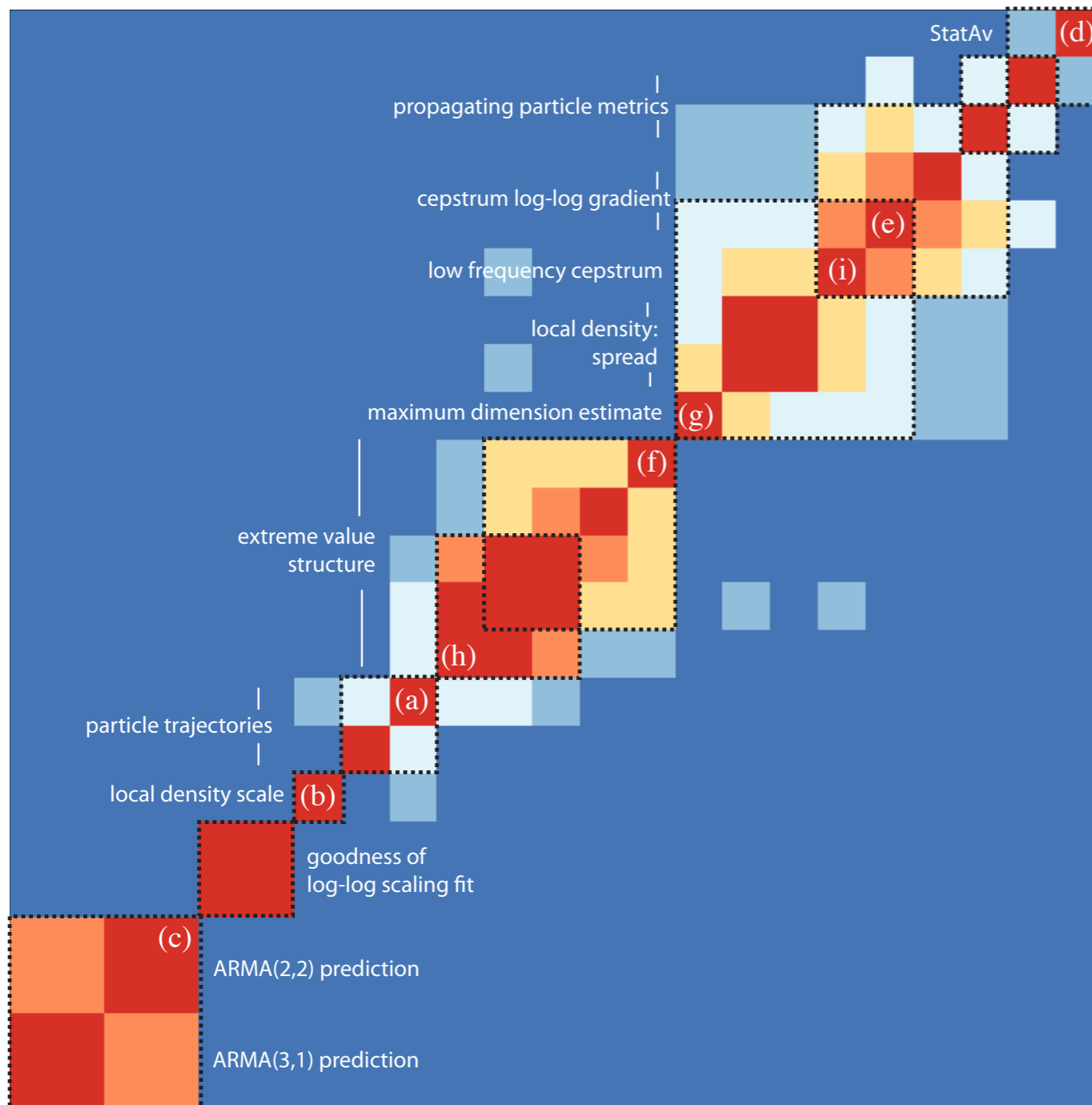
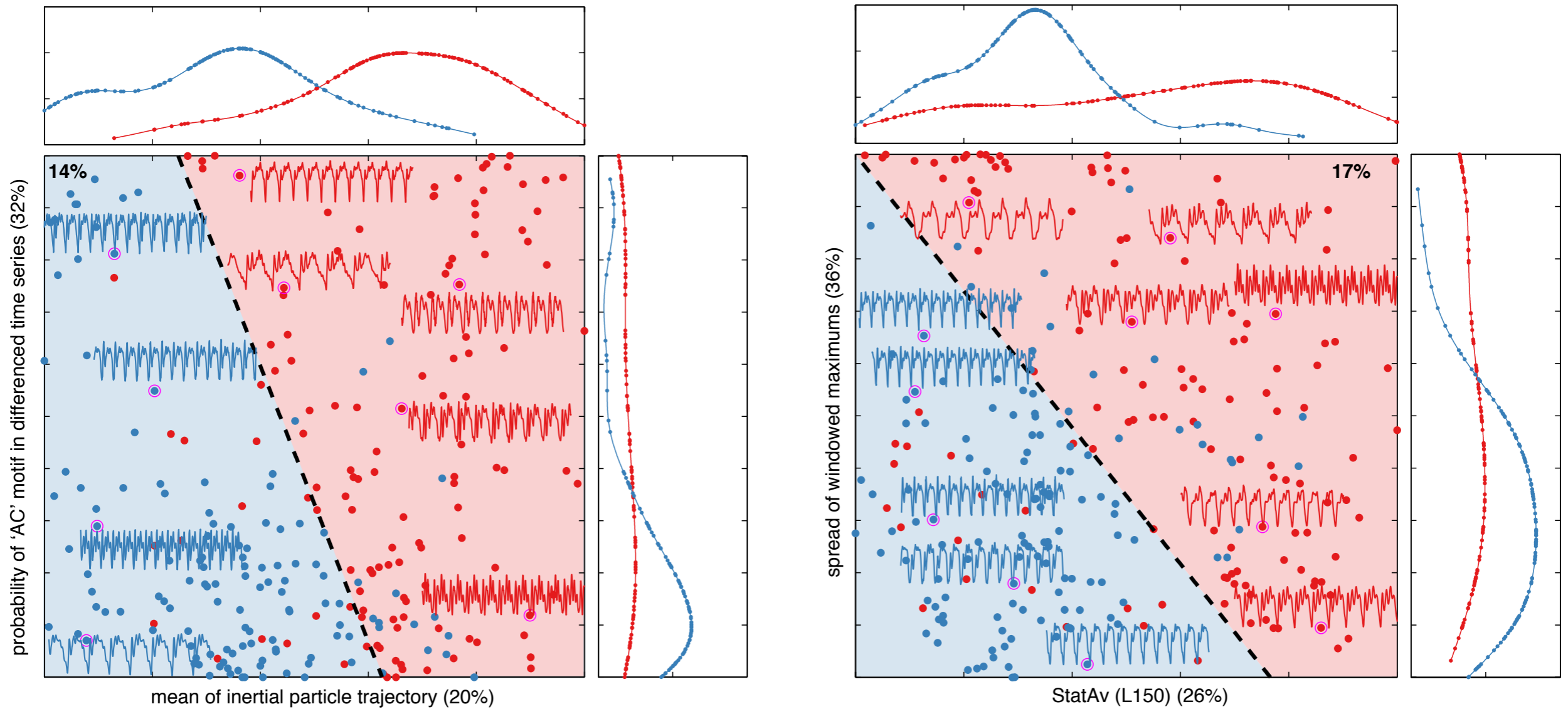# Emotional Speech

# Emotional Speech

# Heart Rate Variability
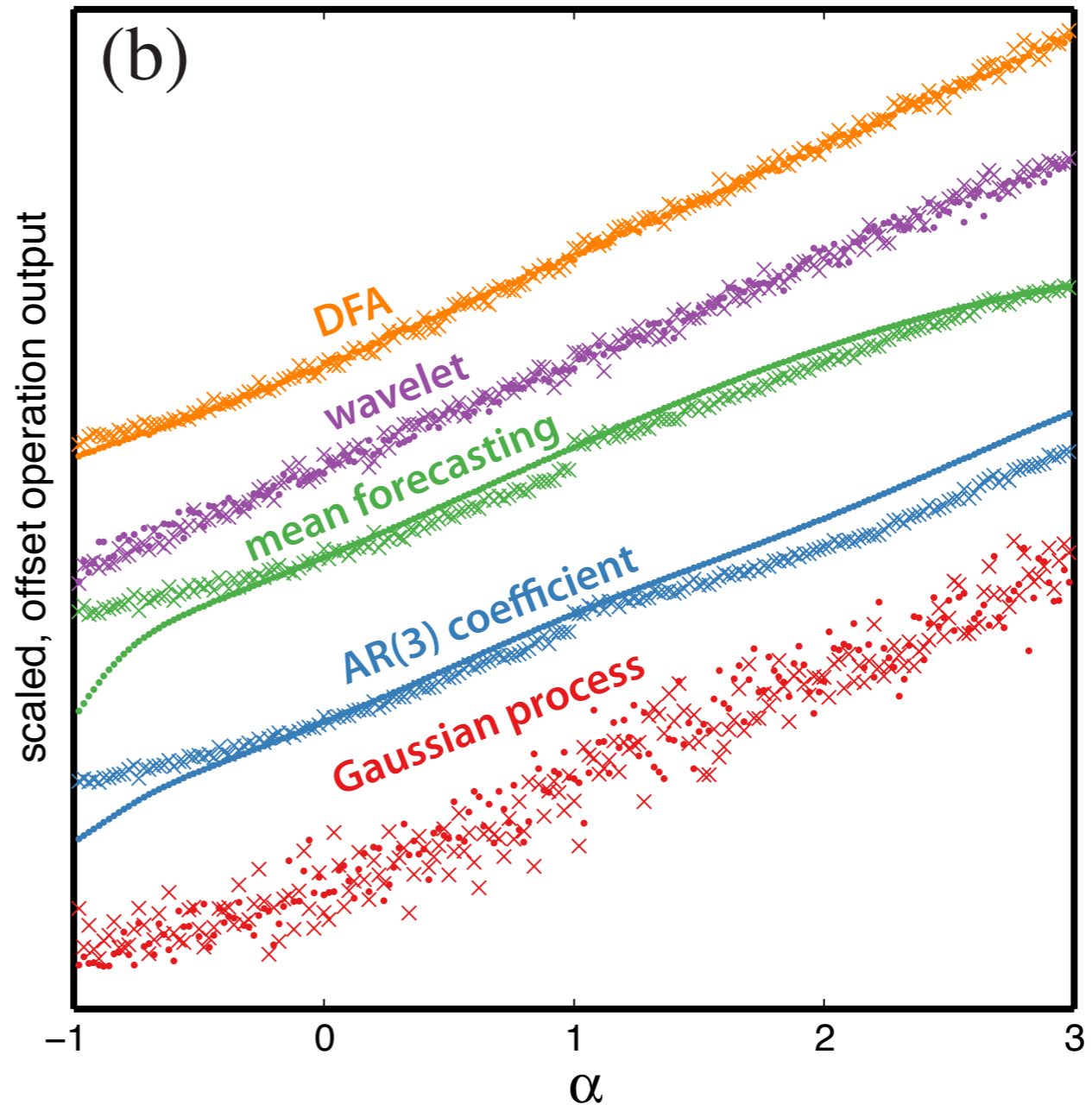
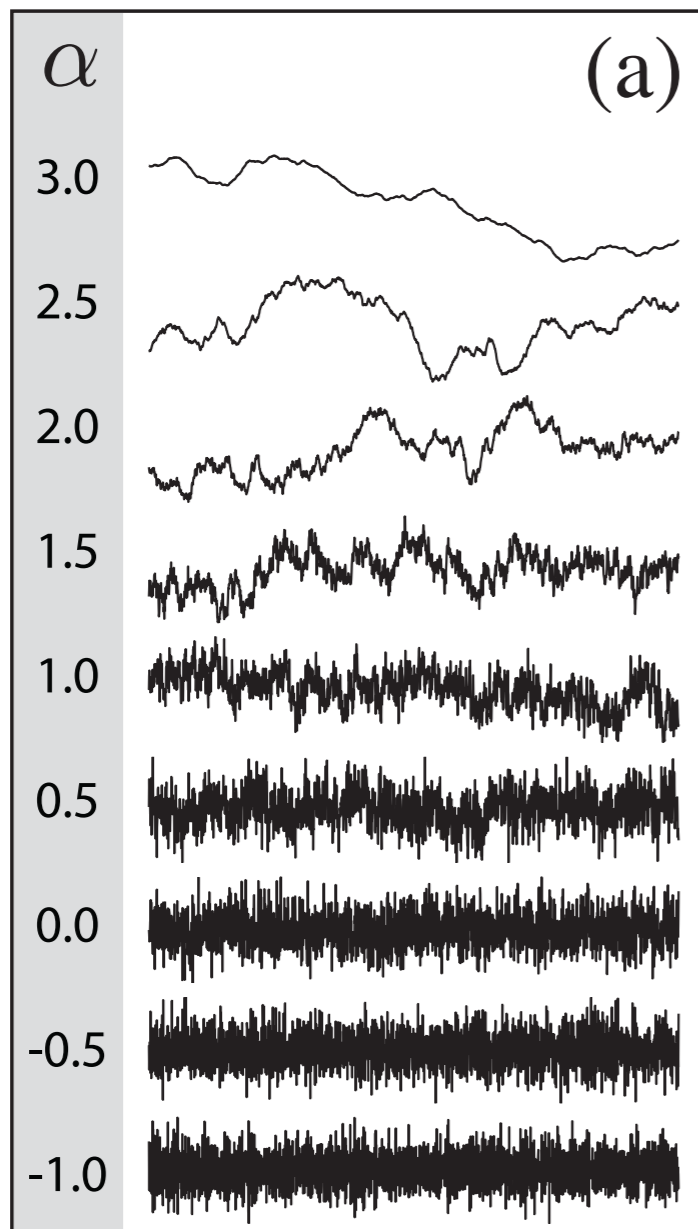# Parkinson's Disease Speech
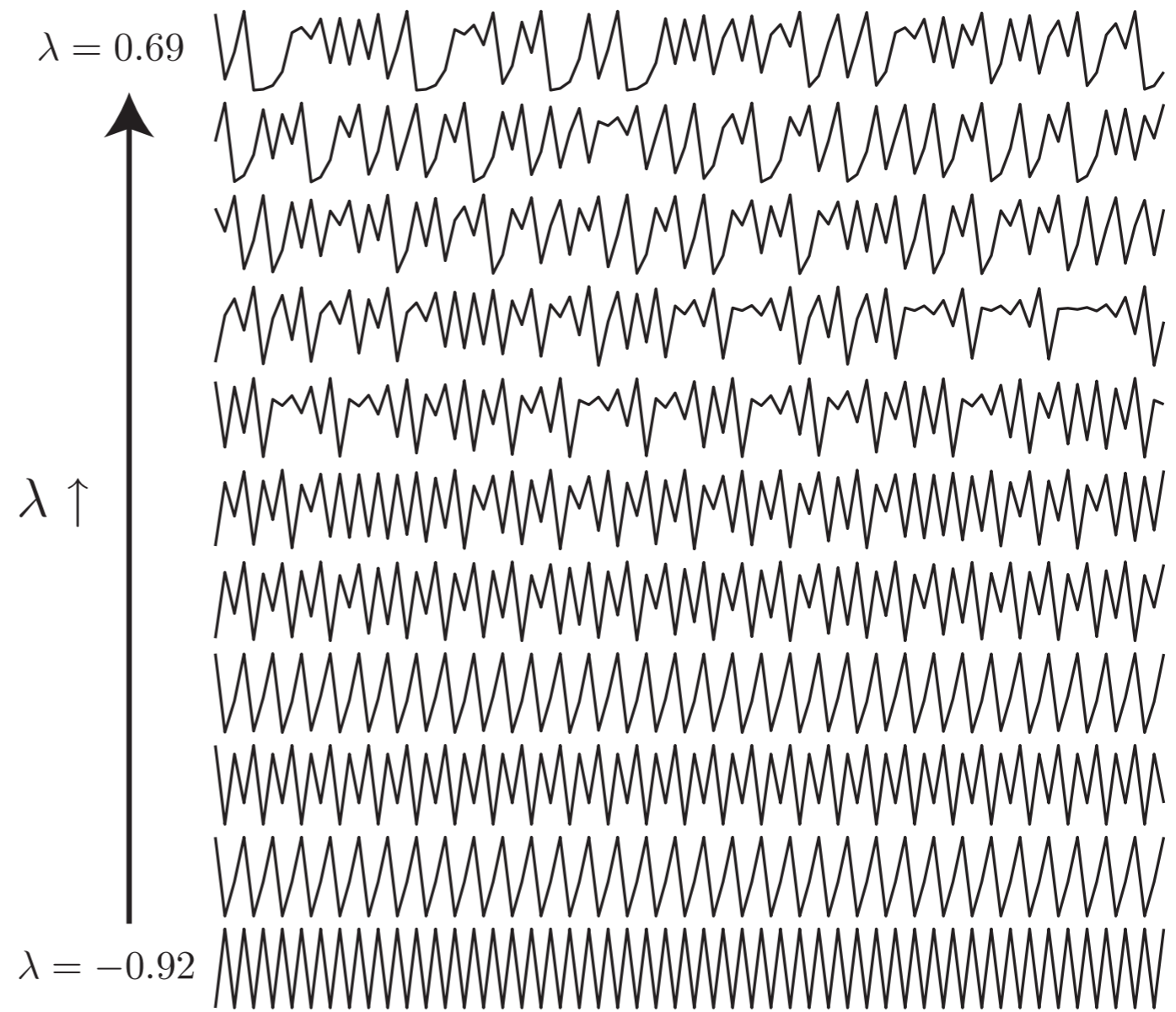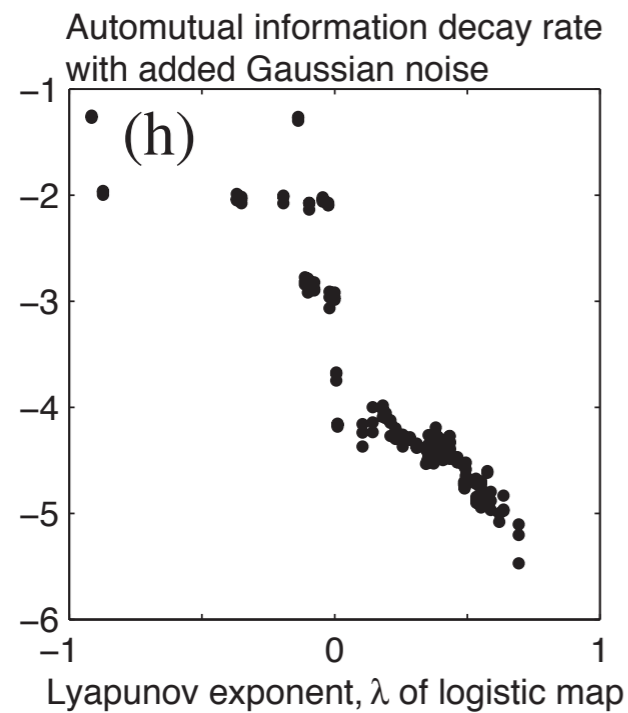
# Parkinsonian Speech

# Parkinsonian Speech



*classifiers mix methods developed in different disciplines*

# Self-Affine Time Series

# Logistic Map Regression



Largest Lyapunov exponent estimator

(a)

unscaled operation output

Automutual information decay rate with added Gaussian noise

(h)

Lyapunov exponent, $\lambda$ of logistic map

$\lambda = 0.69$

$\lambda \uparrow$

$\lambda = -0.92$

logistic maps: $x_{n+1} = A x_n (1 - x_n)$

(a)

Lyapunov exponent

(b)

45 operations

(c)

159 Logistic map time series
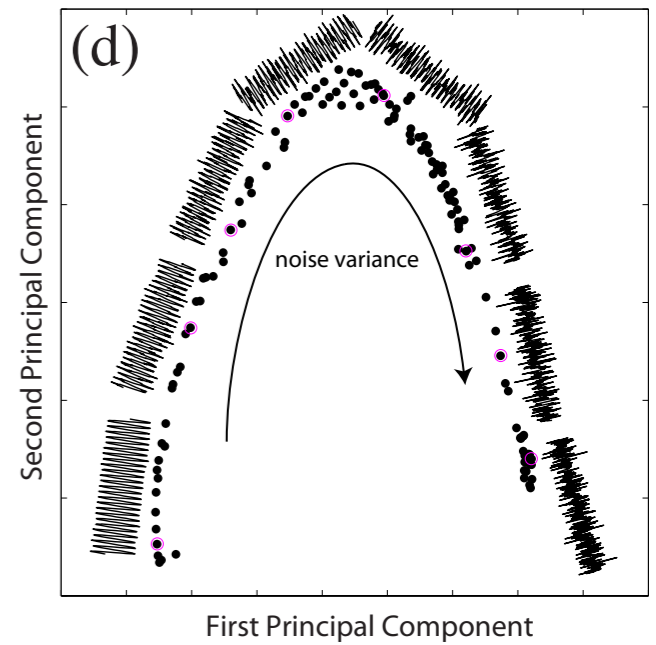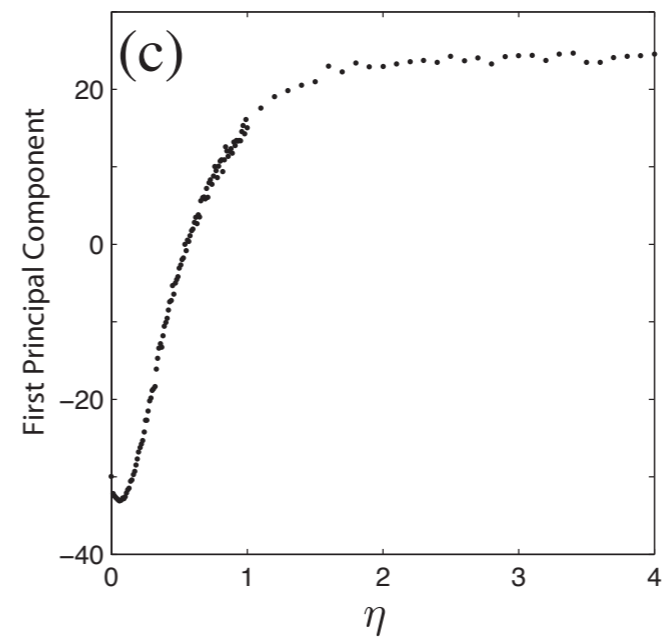
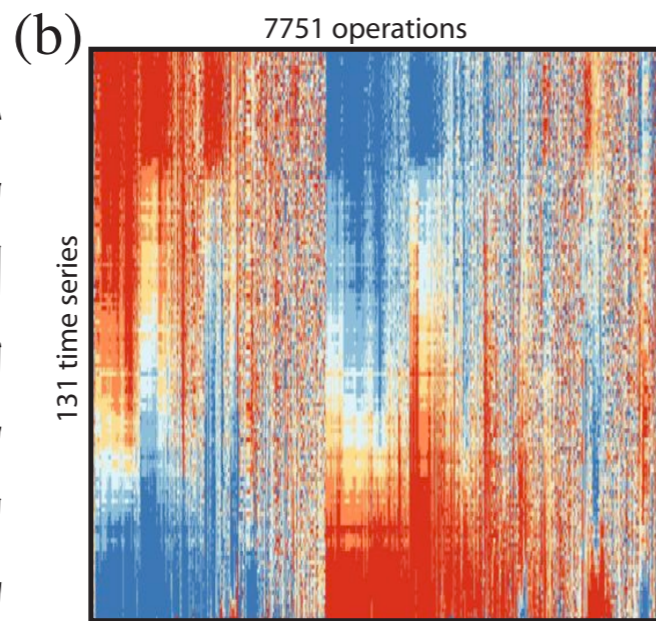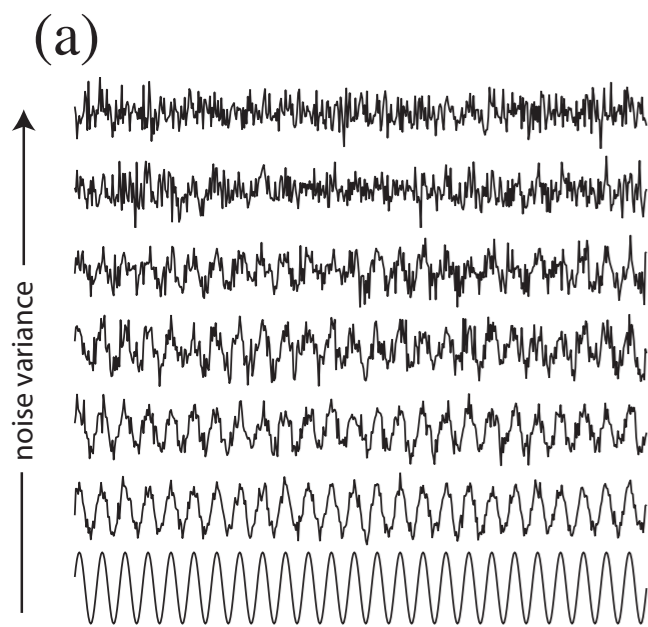increasing *A* ⟶

Logistic Map

# Constrained systems

- We've seen redundancy in set of methods for natural signals.

- What about systems that can be fully described by a small number of parameters?

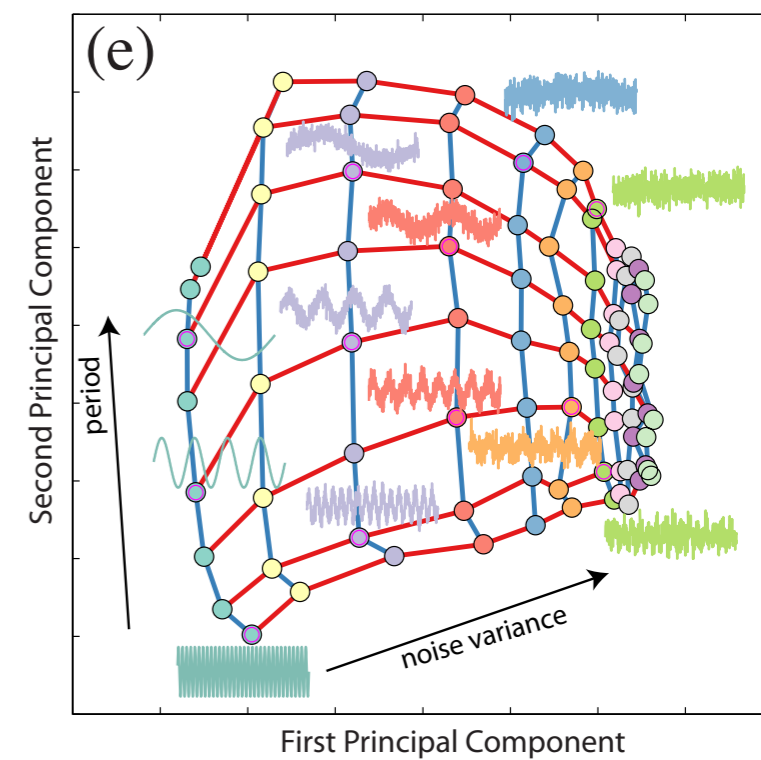- The structure of our database can hint at this.

# Many parameters



(a)

(b) 6237 operations

875 time series

(c) First Principal Component

Second Principal Component

# One parameter



(a) noise variance

(b) 7751 operations / 131 time series

(c) First Principal Component vs $\eta$

(d) Second Principal Component vs First Principal Component / noise variance

# Two parameters



(a) period / noise variance

(b) 7611 operations / 88 time series

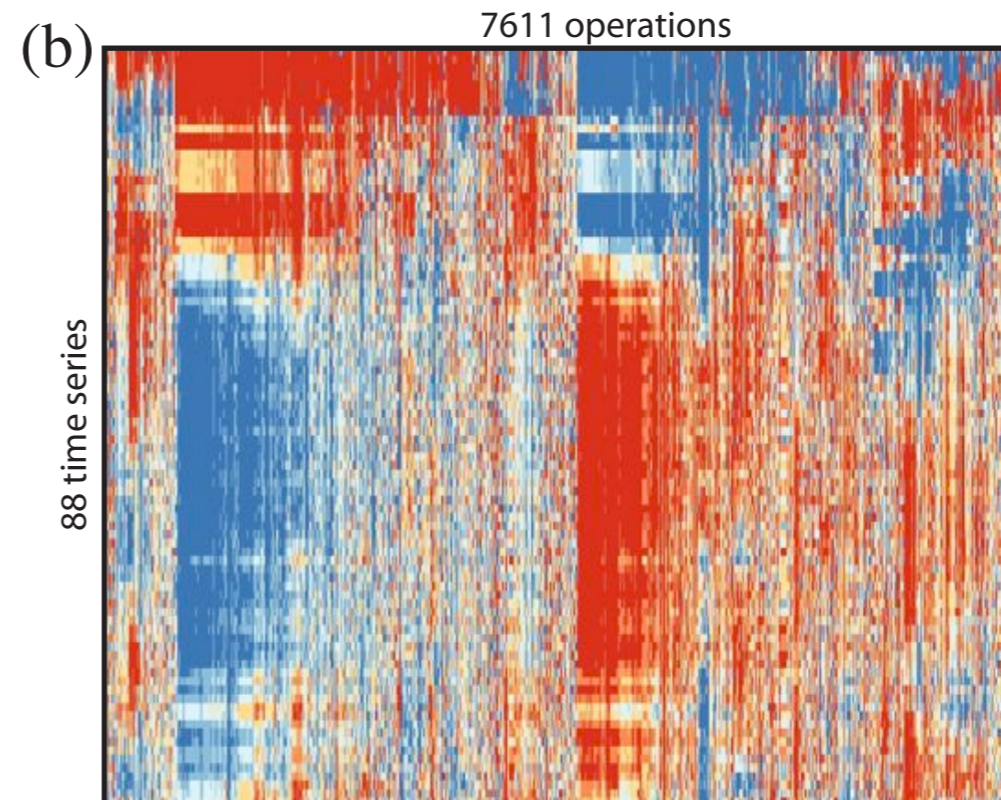(c) First Principal Component vs $\eta$

(d) Second Principal Component vs period (samples)

(e) Second Principal Component vs First Principal Component; period / noise variance
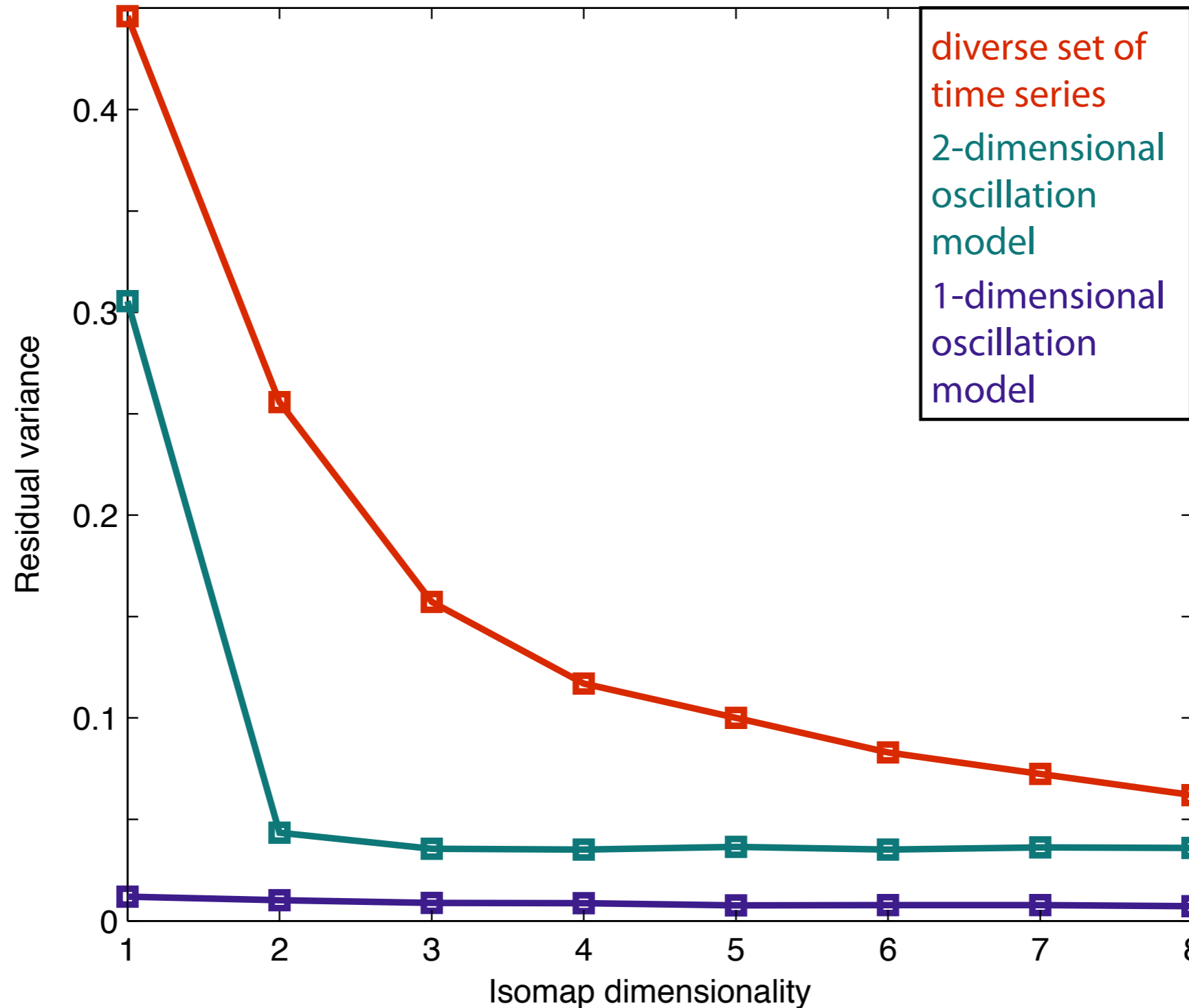
# *Isomap* can quantify this

# Conclusions

- Empirical organization of the methods we use in science

- Empirical organization of the time series and models we study in science

- Automatic classification and regression with the ability to give insights into underlying dynamics