

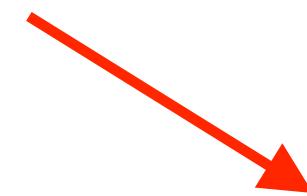
# Introduction to Time Series Mining

Slides from Keogh  
Eamonn's tutorial:

The slide is a presentation from VLDB2006, titled "Eamonn Keogh's VLDB06 Tutorial A Decade of Progress in Indexing and Mining Time Series Data". It features a red header with the VLDB logo and a blue footer. The main content area contains several sections: "Your CD-rom contains:" (listing DFT, DWT, SVD, APCA, PAA, HLP, CLIPPED, SAX, CIEB, PLA), "... excellent tutorial concerning Temporal mining..." (by Dr. Margaret Dunham, in her book, Data Mining Introductory and Advanced Topics), "Why is Time series Data?" (with diagrams of trees and animals), "Practical Examples of Time Series Data Mining", "Defining Distance Measures" (with a definition of distance and a person thinking), and a quote from Ben Shneiderman: "Awesome tutorial!! It's just wonderful... playful AND deep! I couldn't stop looking at it, even though I've got other things to do.... it was a well spent hour!"

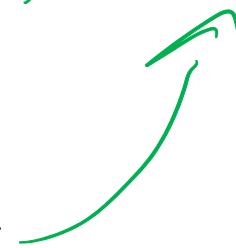
# Outline of Tutorial

- Introduction, Motivation
- The Utility of Similarity Measurements
  - Properties of distance measures
  - The Euclidean distance
  - Preprocessing the data
  - Dynamic Time Warping
  - Uniform Scaling



- Data Mining
  - Anomaly/Interestingness detection
  - Motif (repeated pattern) discovery

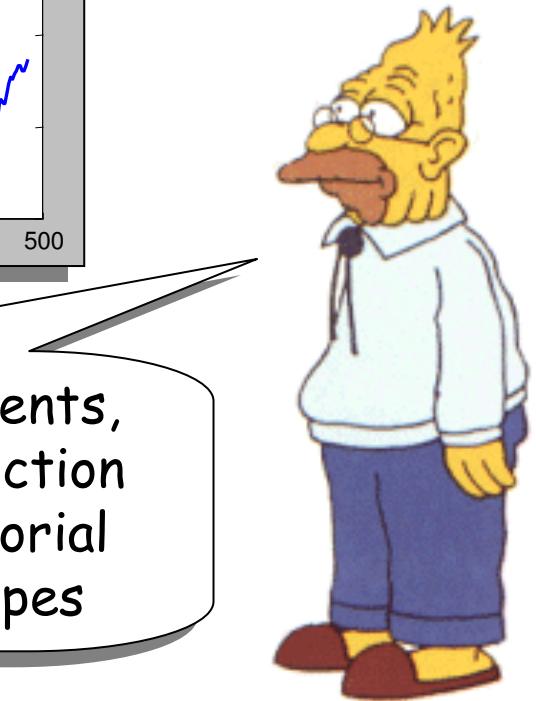
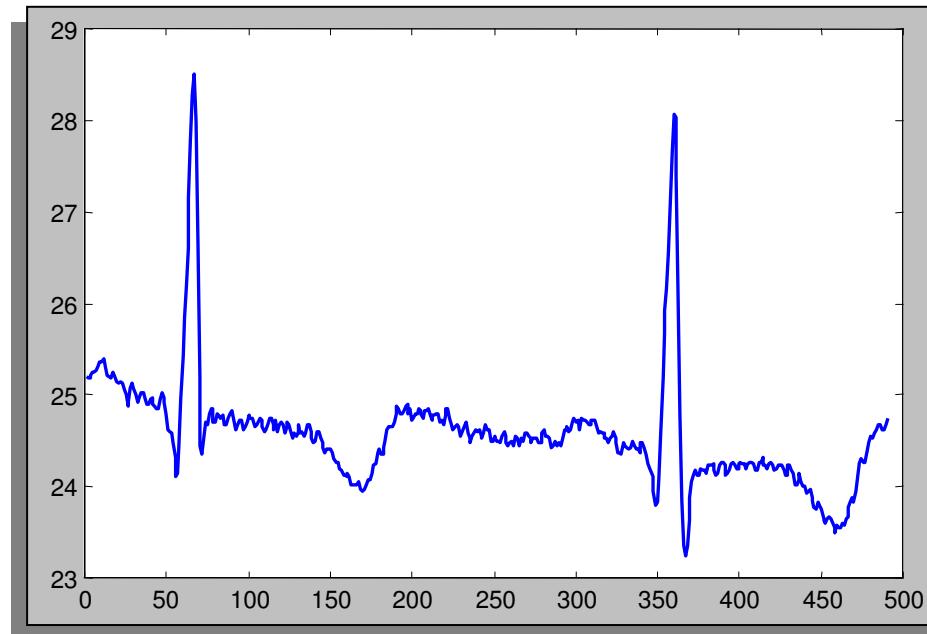
CLASSICAL  
APPROACHES



25.1750  
25.2250  
25.2500  
25.2500  
25.2750  
25.3250  
  
25.3500  
25.3500  
25.4000  
25.4000  
25.3250  
25.2250  
25.2000  
25.1750  
  
..  
..  
24.6250  
24.6750  
24.6750  
24.6250  
24.6250  
24.6250  
24.6750  
24.7500

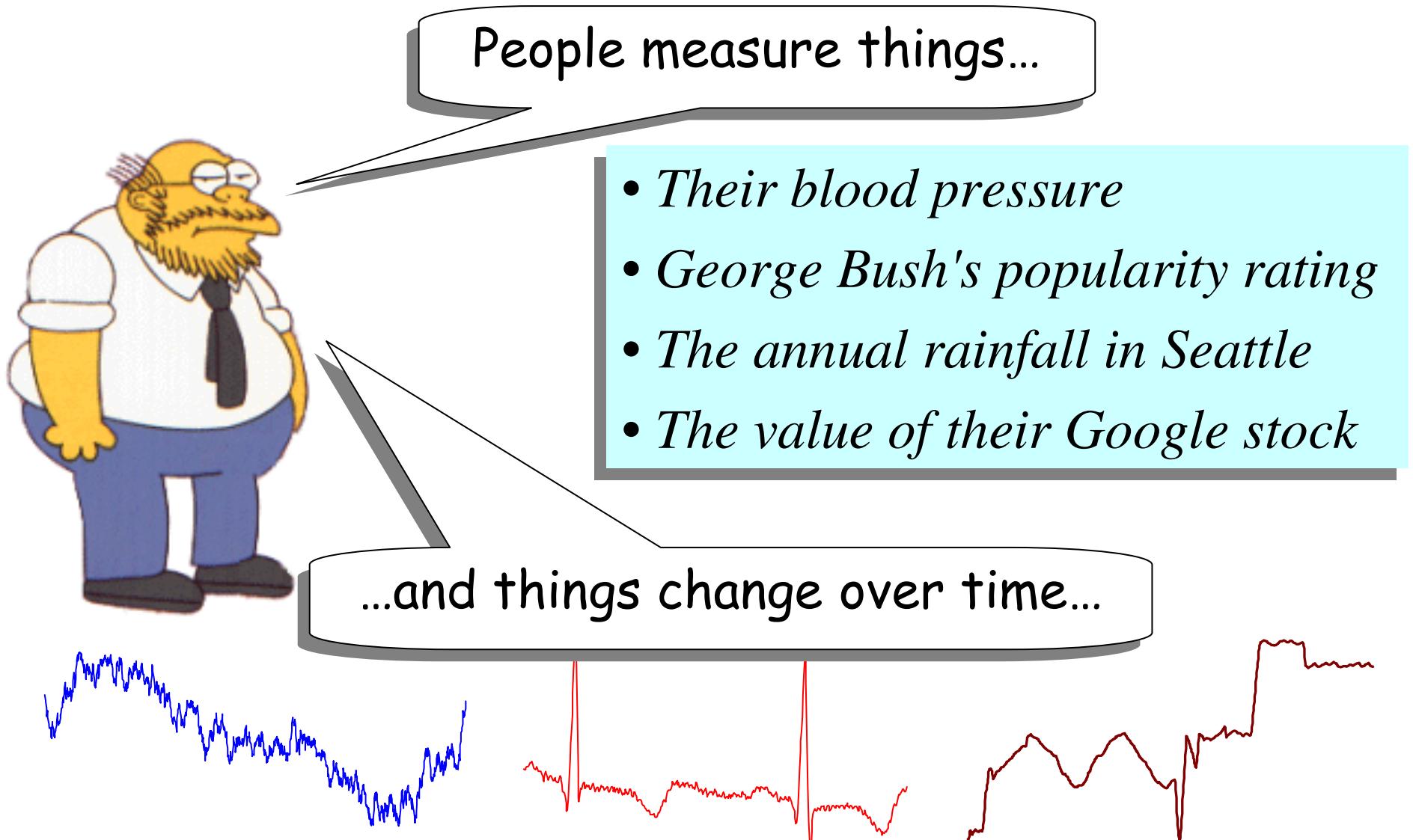
# What are Time Series?

A time series is a collection of observations made sequentially in time.



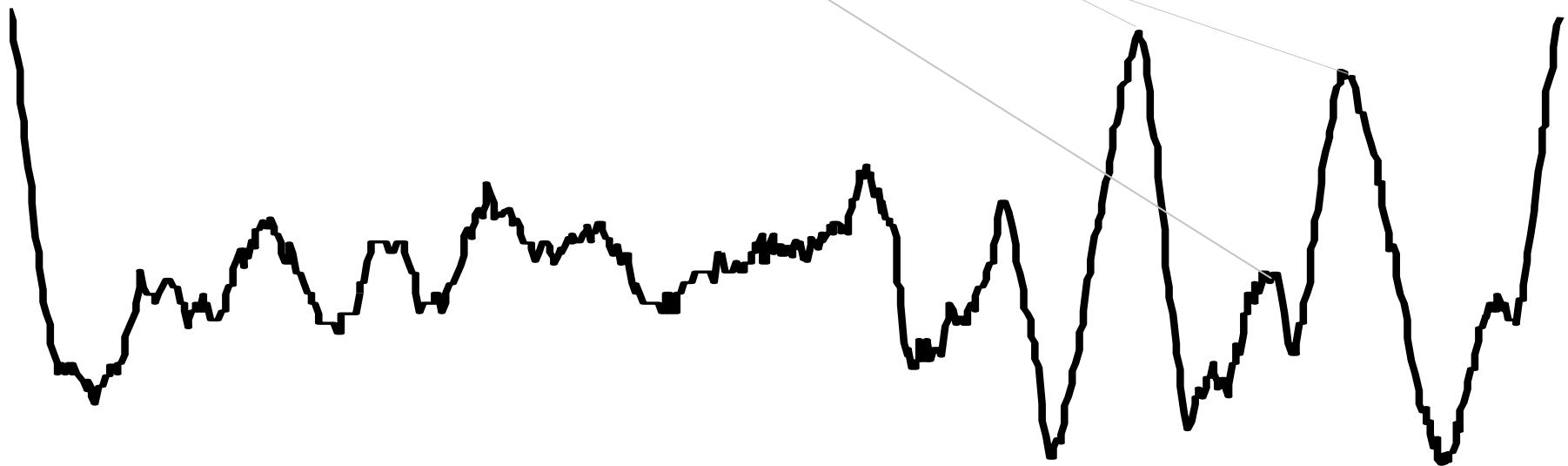
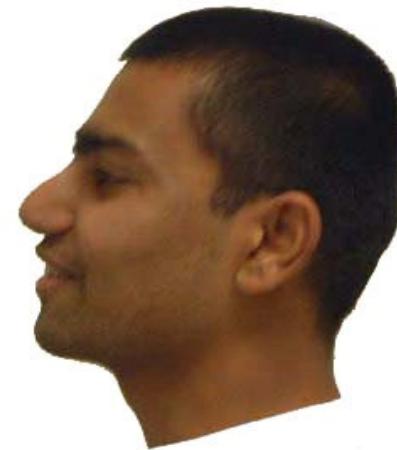
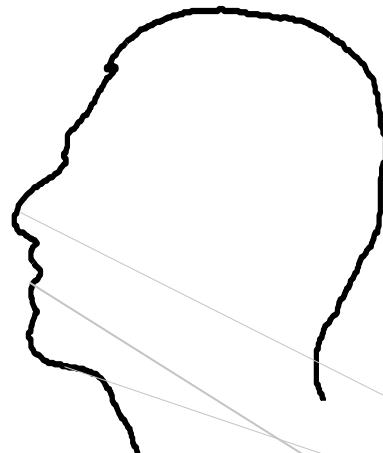
Virtually all similarity measurements,  
indexing and dimensionality reduction  
techniques discussed in this tutorial  
can be used with other data types

# Time Series are Ubiquitous! I



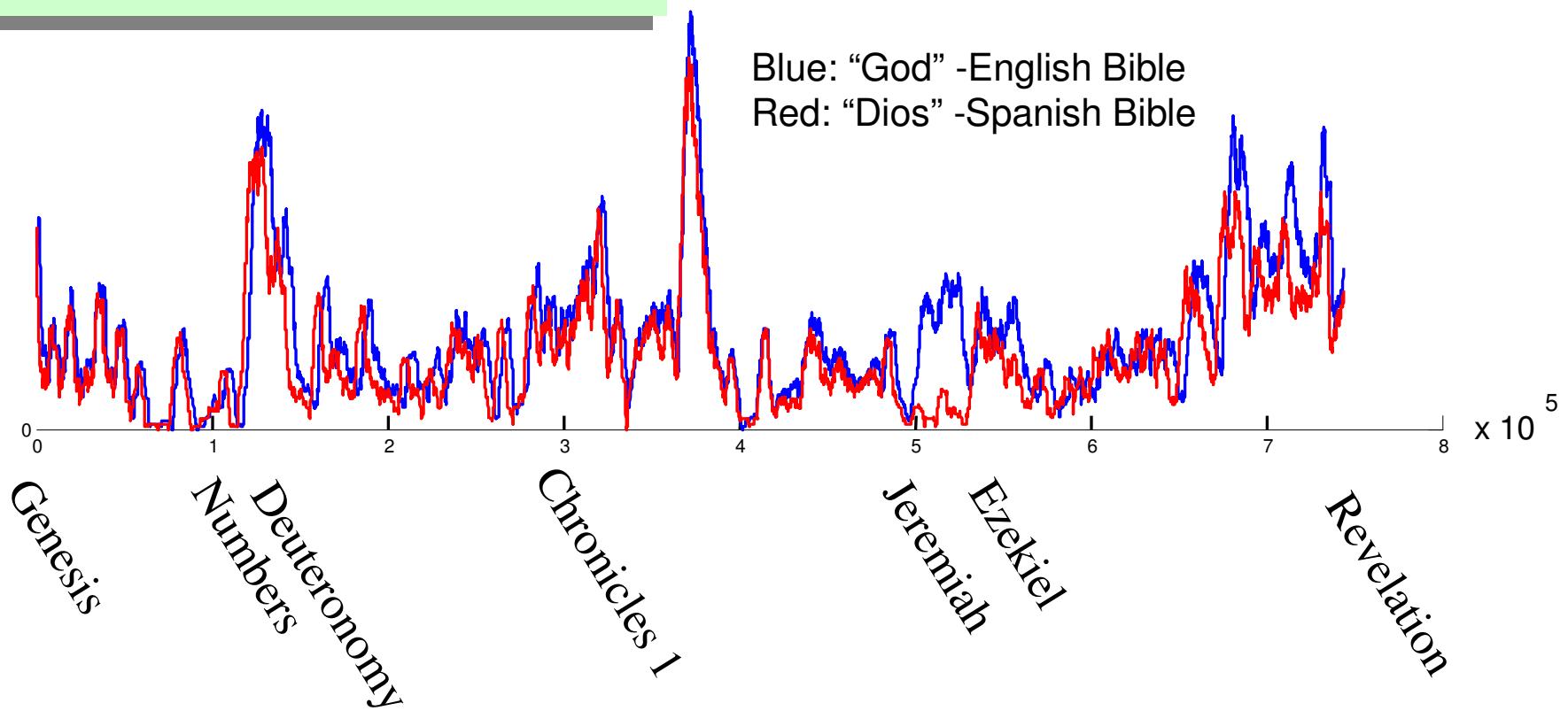
Thus time series occur in virtually every medical, scientific and businesses domain

Image data, may best be thought of as time series...



# Text data, may best be thought of as time series...

The local frequency  
of words in the Bible



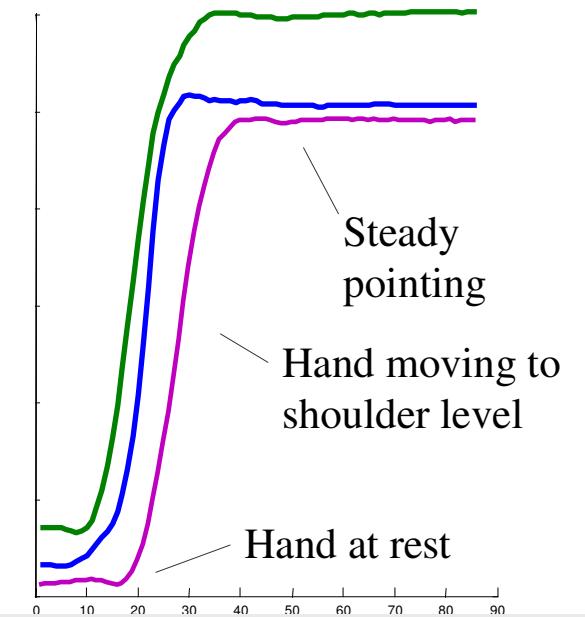
Gray: "El Señor" -Spanish Bible



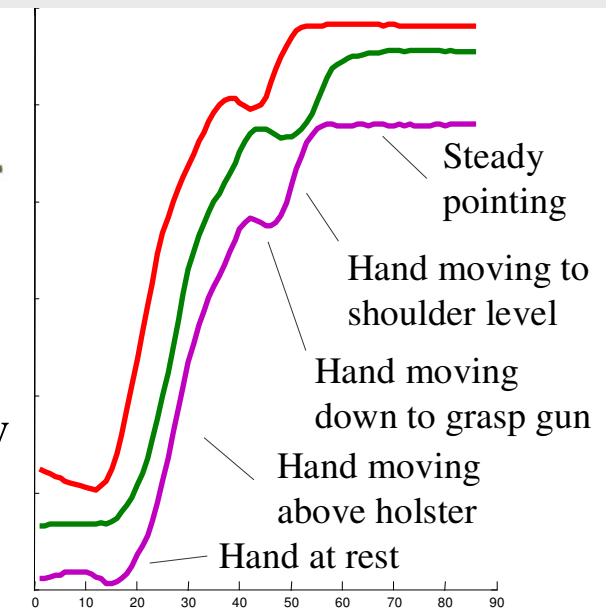
# Video data, may best be thought of as time series...



**Point**



**Gun-Draw**



# Why is Working With Time Series so Difficult? Part I

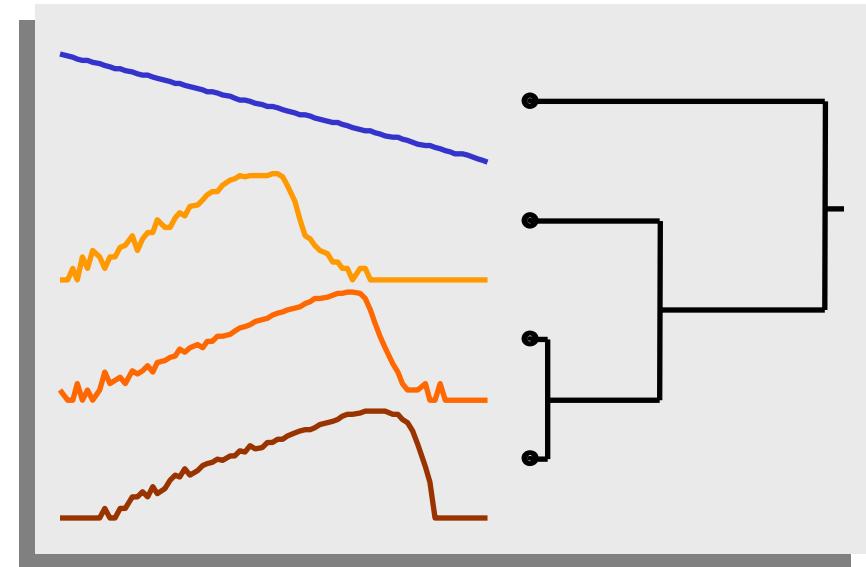
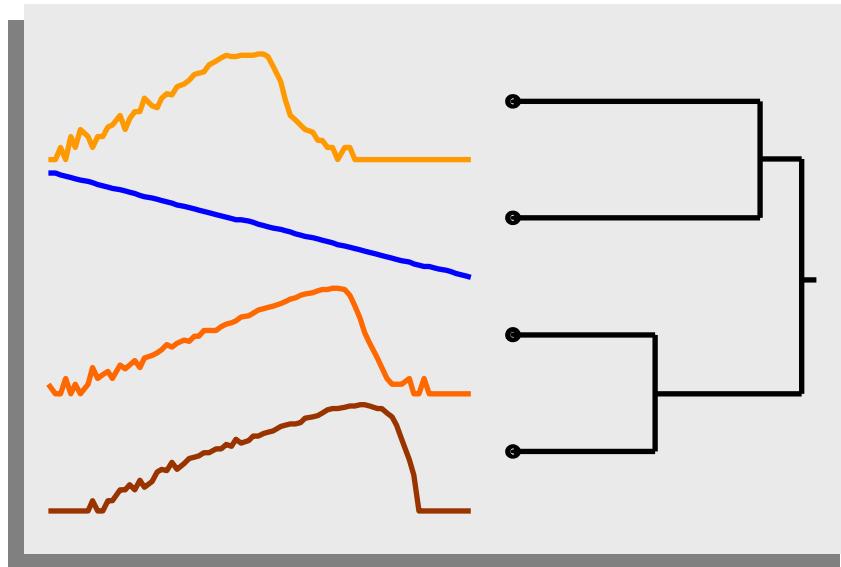
**Answer:** How do we work with very large databases?

- ◆ 1 Hour of EKG data: 1 Gigabyte.
- ◆ Typical Weblog: 5 Gigabytes per week.
- ◆ Space Shuttle Database: 200 Gigabytes and growing.
- ◆ Macho Database: 3 Terabytes, updated with 3 gigabytes a day.

Since most of the data lives on disk (or tape), we need a representation of the data we can efficiently manipulate.

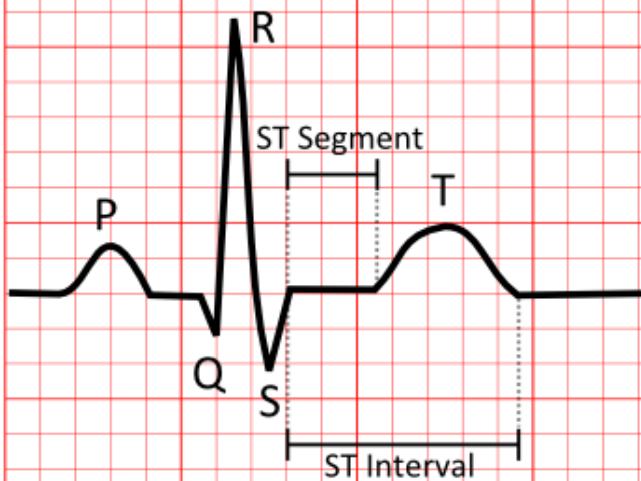
# Why is Working With Time Series so Difficult? Part II

**Answer:** We are dealing with subjectivity



The definition of similarity depends on the user, the domain and the task at hand. We need to be able to handle this subjectivity.

Normal



ST elevation



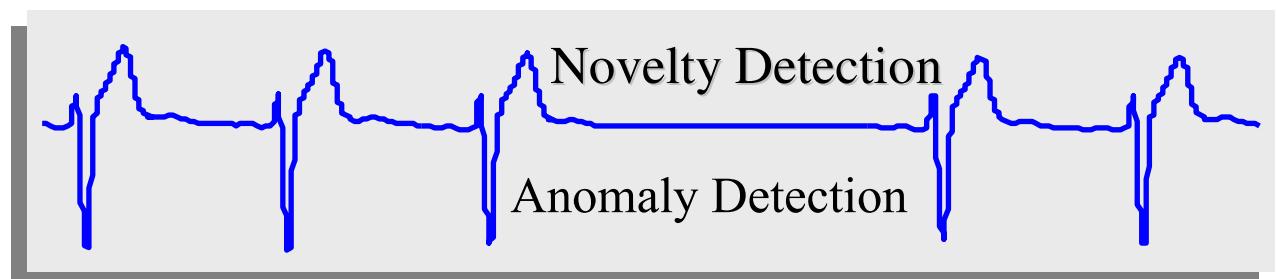
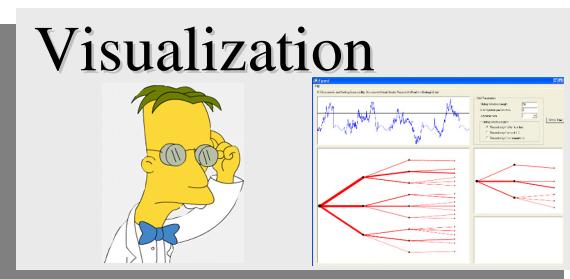
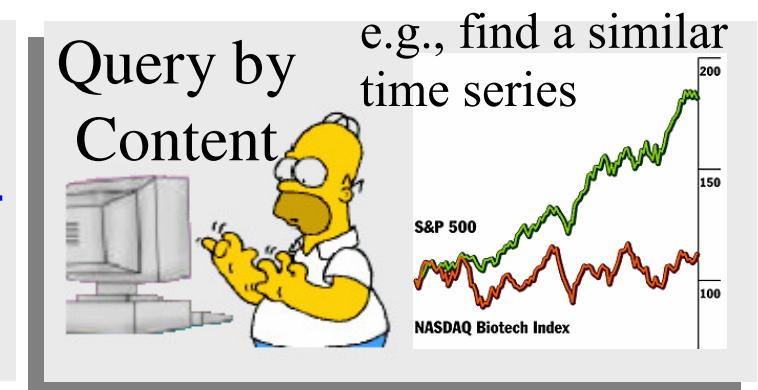
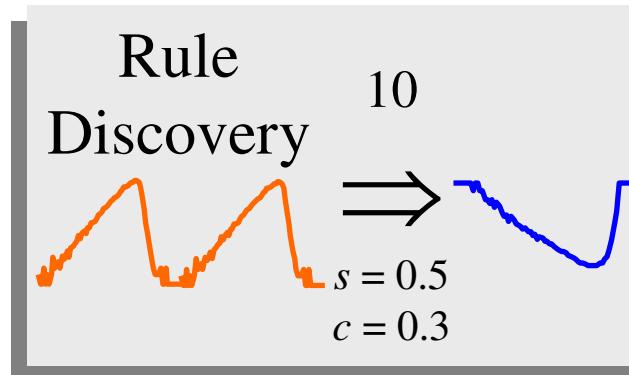
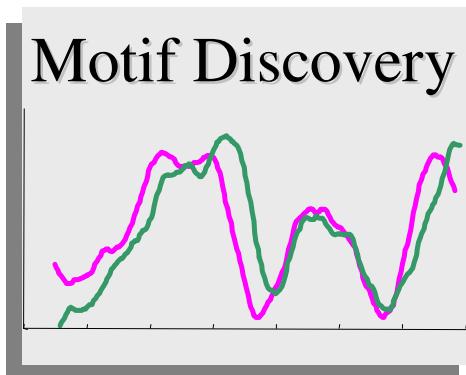
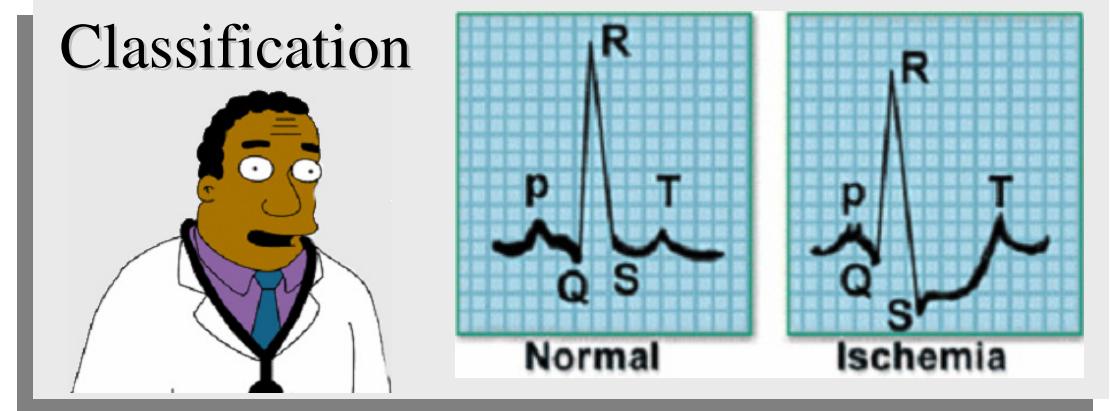
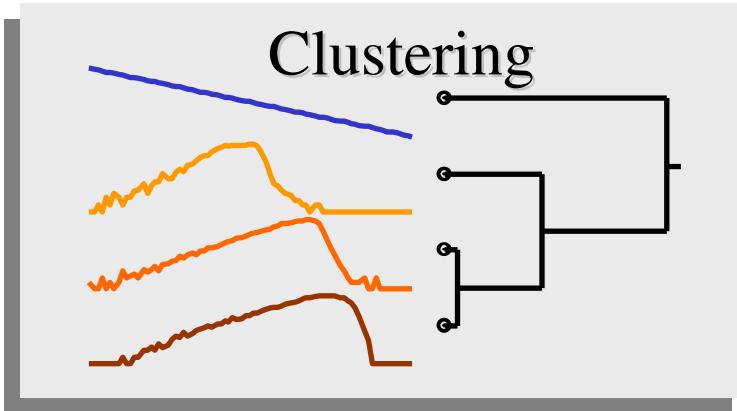
# Why is working with time series so difficult? Part III

**Answer:** Miscellaneous data handling problems.

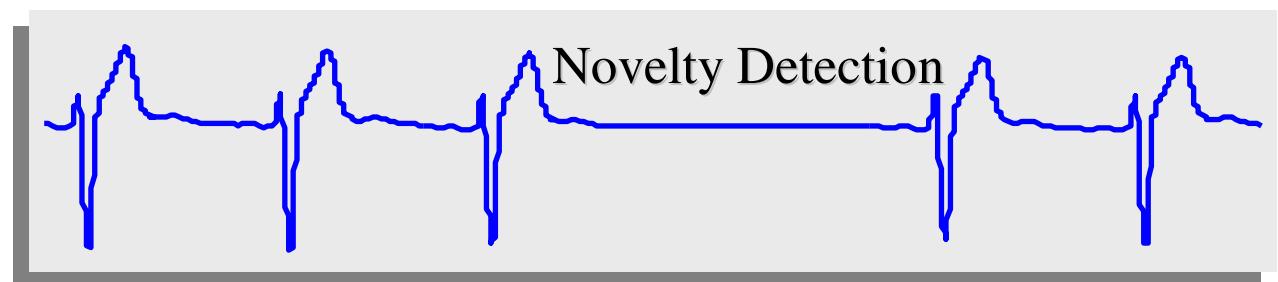
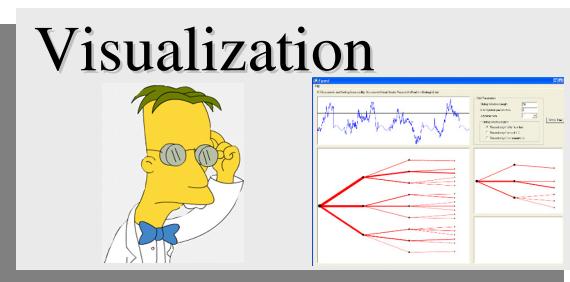
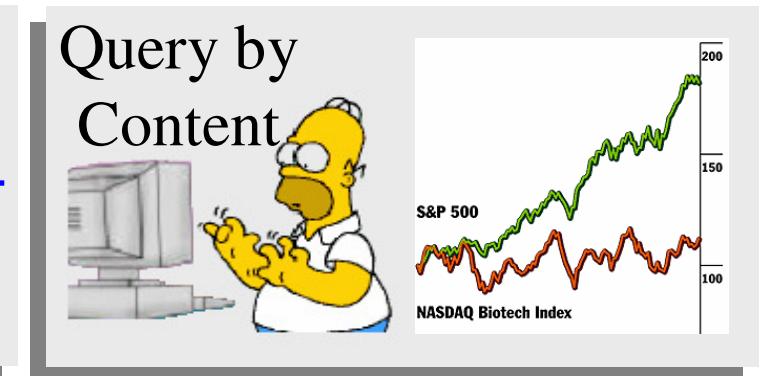
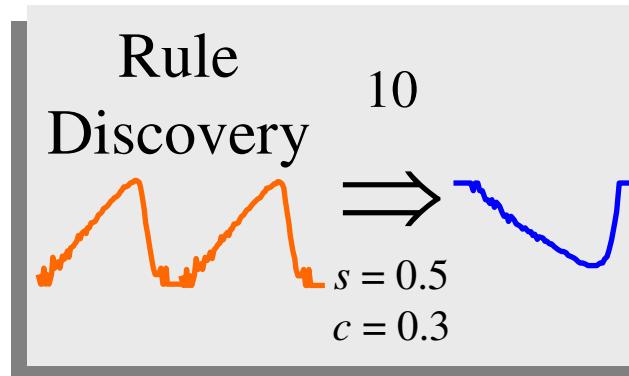
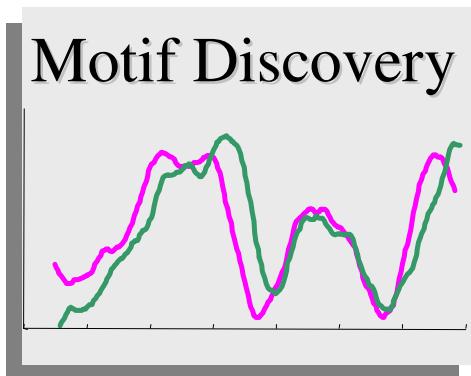
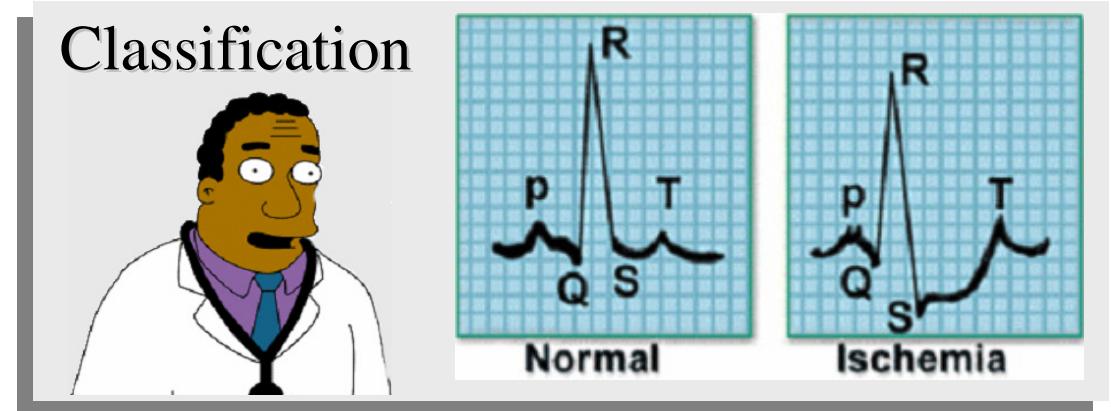
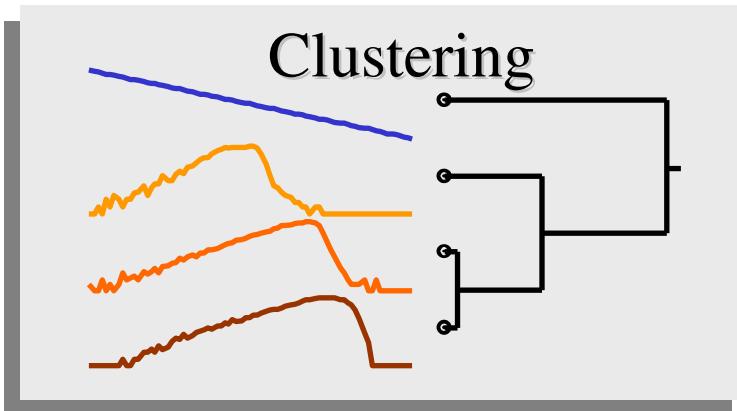
- Differing data formats.
- Differing sampling rates.
- Noise, missing values, etc.

We will not focus on these issues in this tutorial.

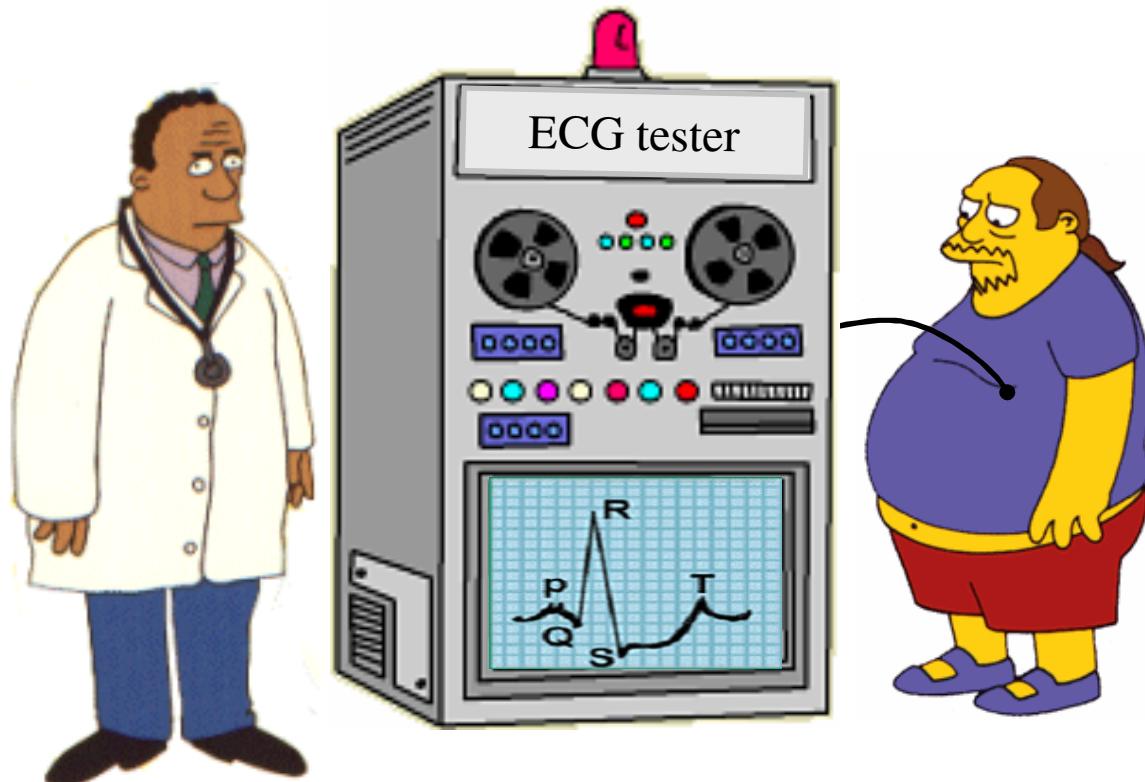
# What do we want to do with the time series data?



# All these problems require similarity matching



# Here is a simple motivation for the first part of the tutorial



You go to the doctor because of chest pains. Your ECG looks strange...

You doctor wants to search a database to find **similar** ECGs, in the hope that they will offer clues about your condition...

- **How do we define similar?**

# What is Similarity?

The quality or state of being similar; likeness; resemblance; as, a similarity of features. Webster's Dictionary



Similarity is hard to define, but...  
*“We know it when we see it”*

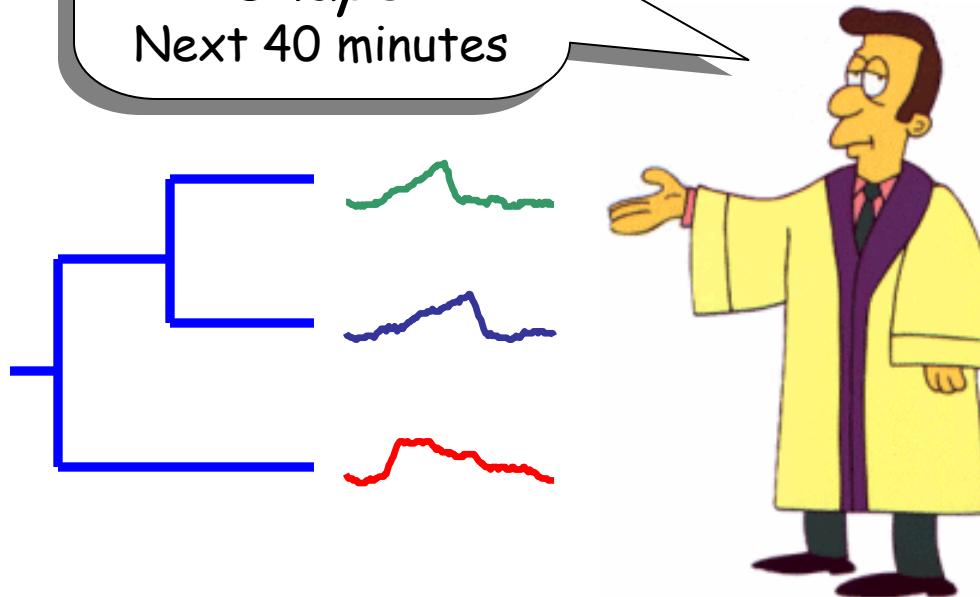
The real meaning of similarity is a philosophical question.

We will take a more pragmatic approach.

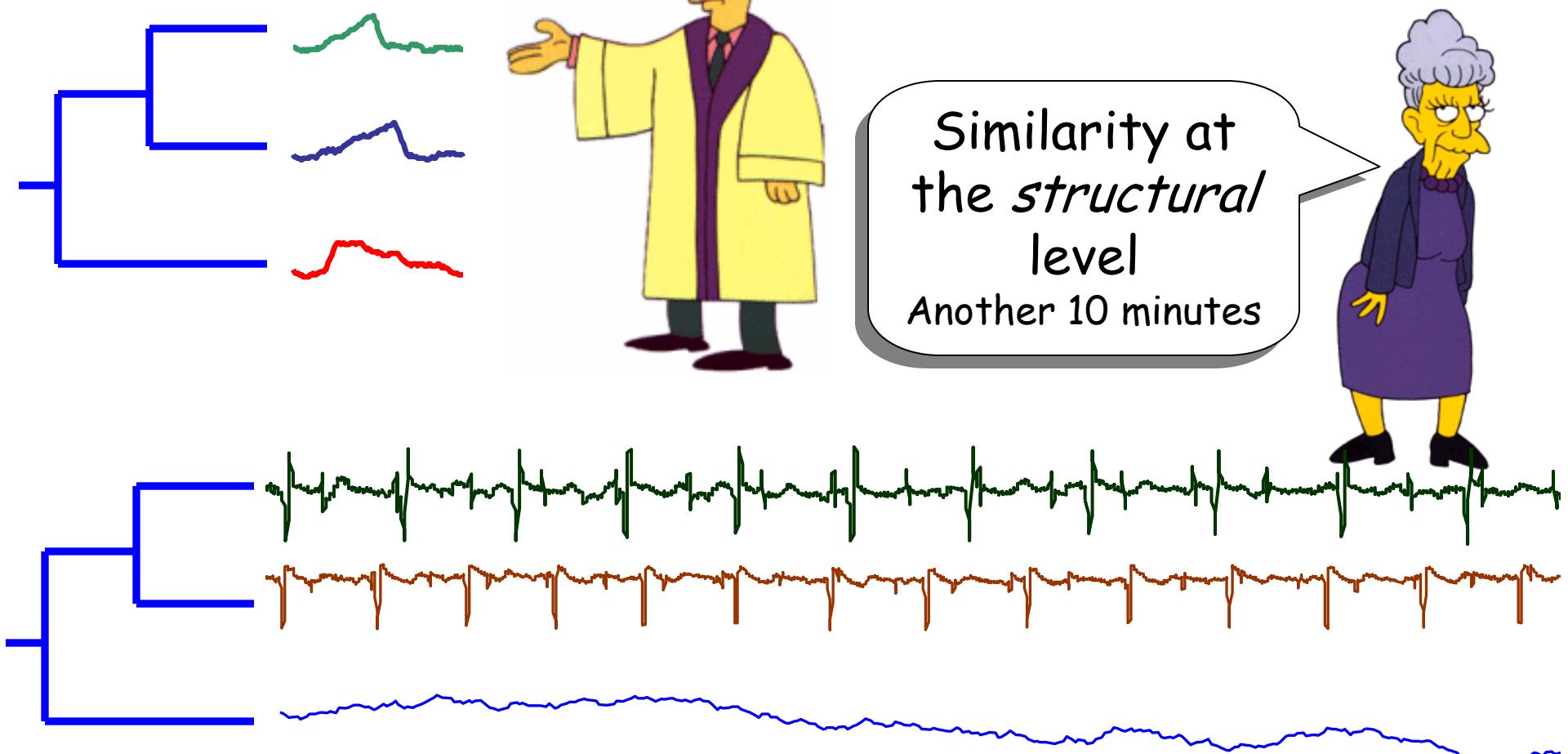
# Two Kinds of Similarity

time series

Similarity at  
the level of  
*shape*  
Next 40 minutes



Similarity at  
the *structural*  
level  
Another 10 minutes



# Euclidean Distance Metric

Given two time series:

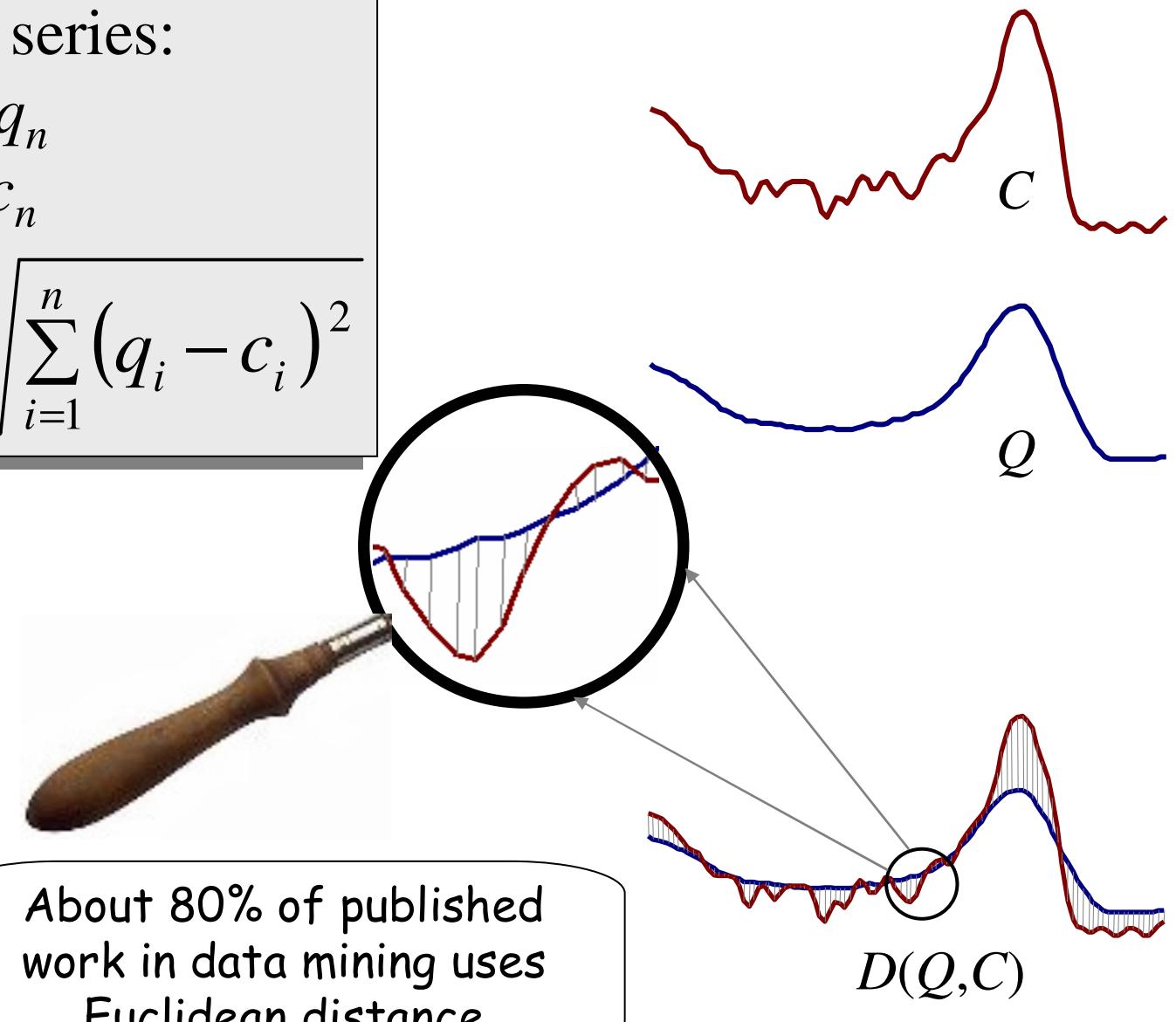
$$Q = q_1 \dots q_n$$

$$C = c_1 \dots c_n$$

$$D(Q, C) \equiv \sqrt{\sum_{i=1}^n (q_i - c_i)^2}$$



About 80% of published  
work in data mining uses  
Euclidean distance



# Preprocessing the data before distance calculations



If we naively try to measure the distance between two "raw" time series, we may get very unintuitive results



This is because Euclidean distance is very sensitive to some "distortions" in the data. For most problems these distortions are not meaningful, and thus we can and should remove them

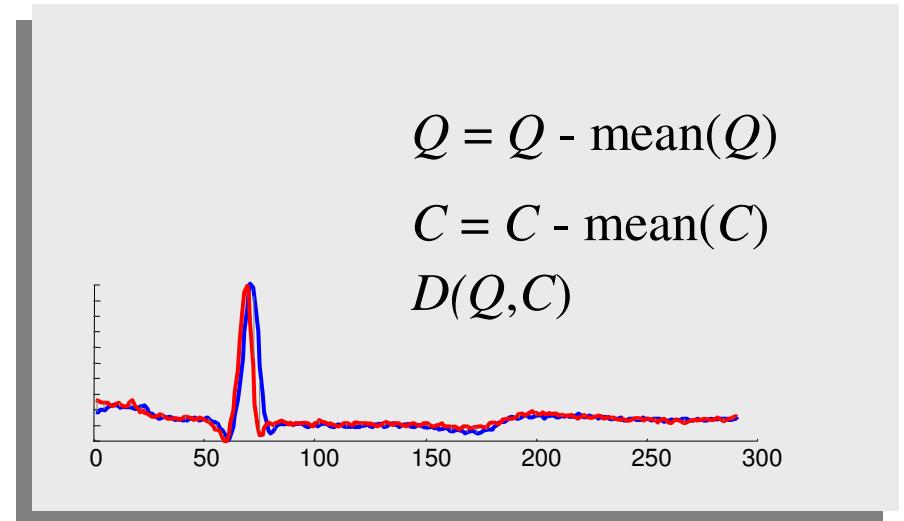
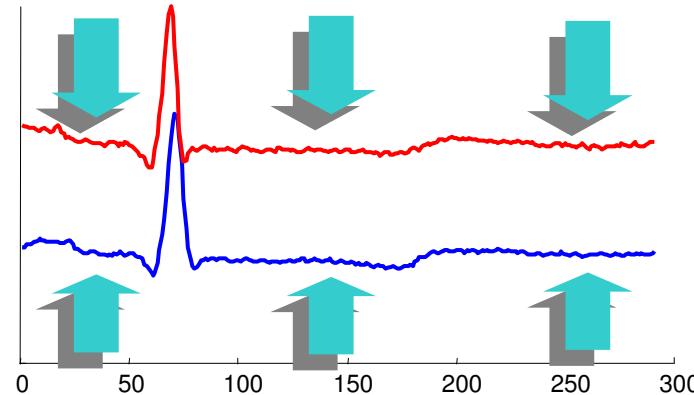
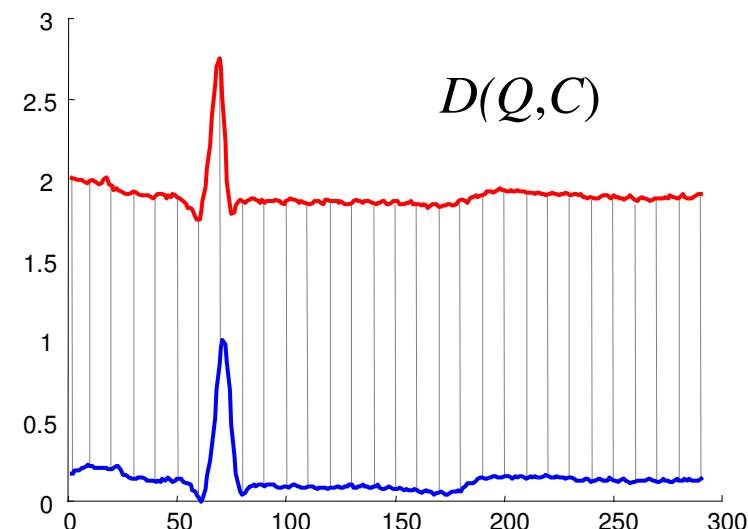
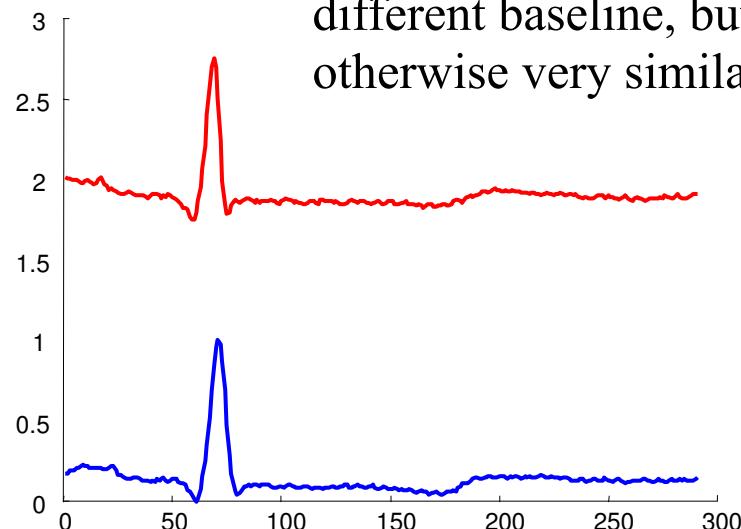


In the next few slides we will discuss the 4 most common distortions, and how to remove them

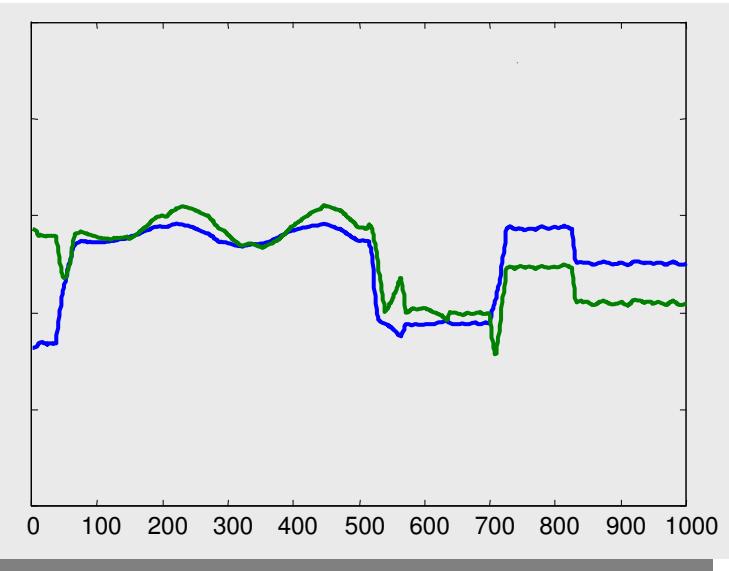
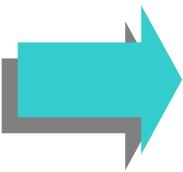
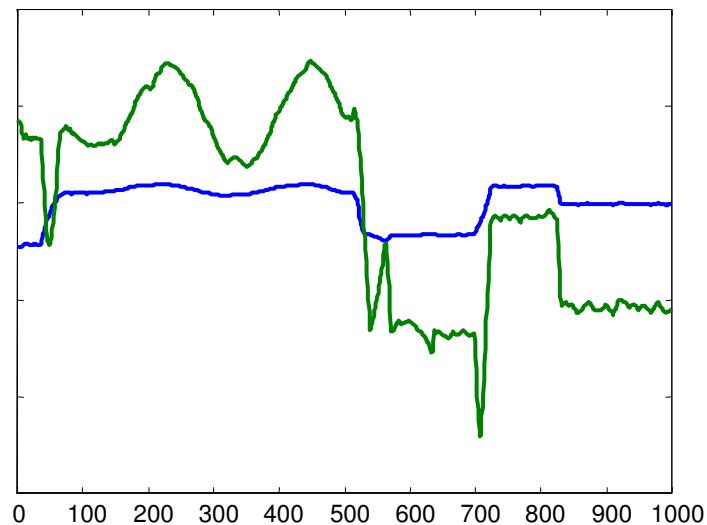
- Offset Translation
- Amplitude Scaling
- Linear Trend
- Noise

# Transformation I: Offset Translation

The two time series have a different baseline, but are otherwise very similar



# Transformation II: Amplitude Scaling



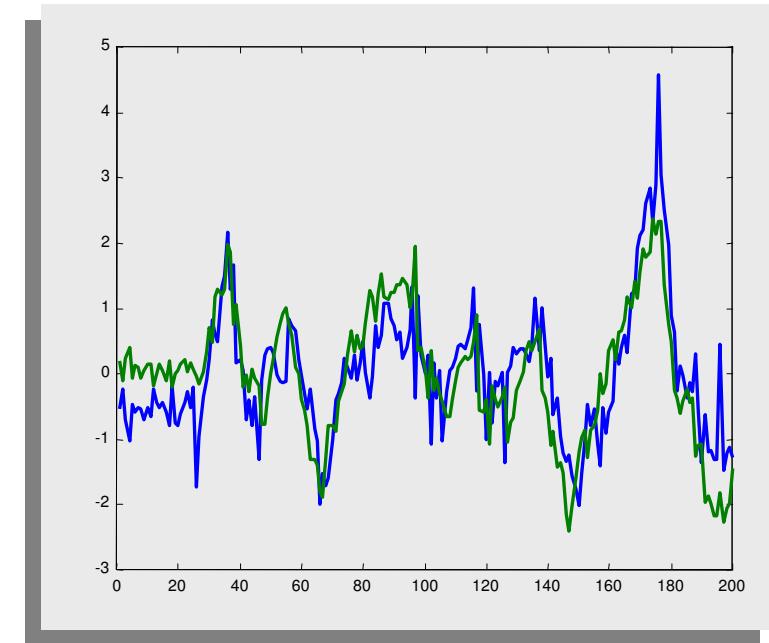
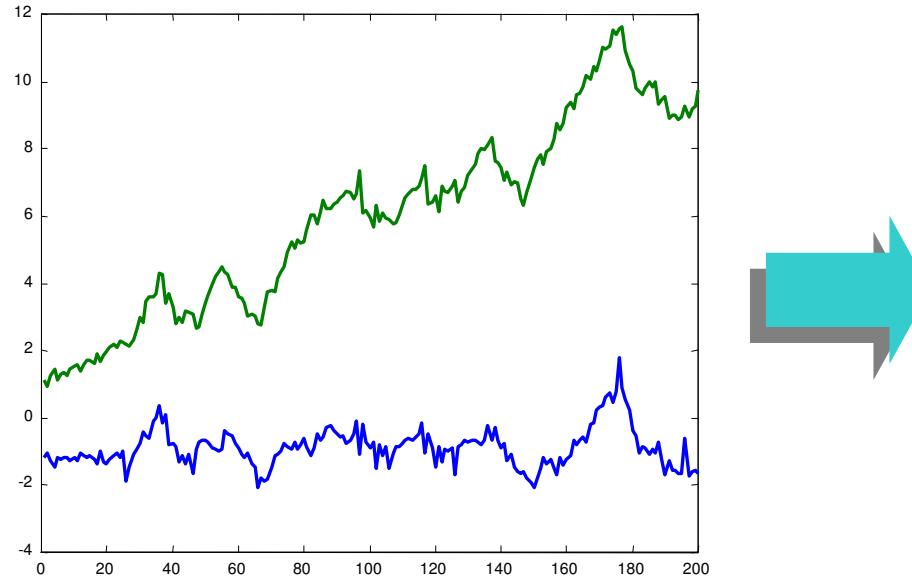
OF COURSE, WHETHER OR NOT YOU  
ARE DESTROYING INFORMATION BY  
MEANS OF THESE TRANSFORMATIONS  
REALLY DEPENDS ON YOUR SPECIFIC  
ANALYSIS TASK -

$$Q = (Q - \text{mean}(Q)) / \text{std}(Q)$$

$$C = (C - \text{mean}(C)) / \text{std}(C)$$

$$D(Q, C)$$

# Transformation III: Linear Trend



The intuition behind removing  
linear trend is...

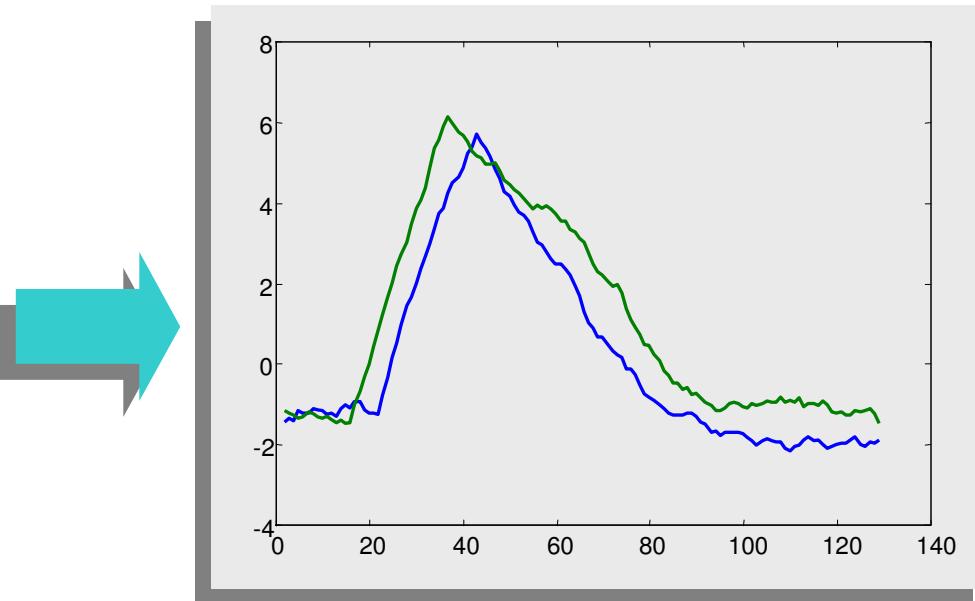
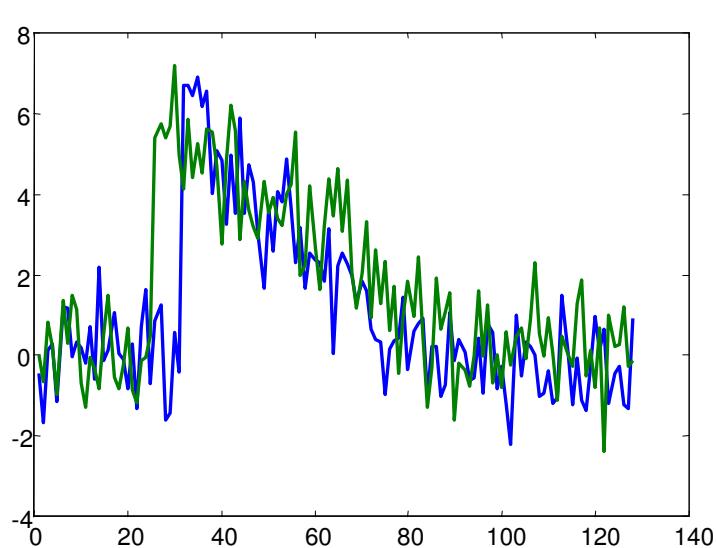
Fit the best fitting straight line to the  
time series, then subtract that line  
from the time series.

**Removed linear trend**

Removed offset translation

Removed amplitude scaling

# Transformation IIII: Noise



The intuition behind removing noise is...

Average each datapoint's value with its neighbors.

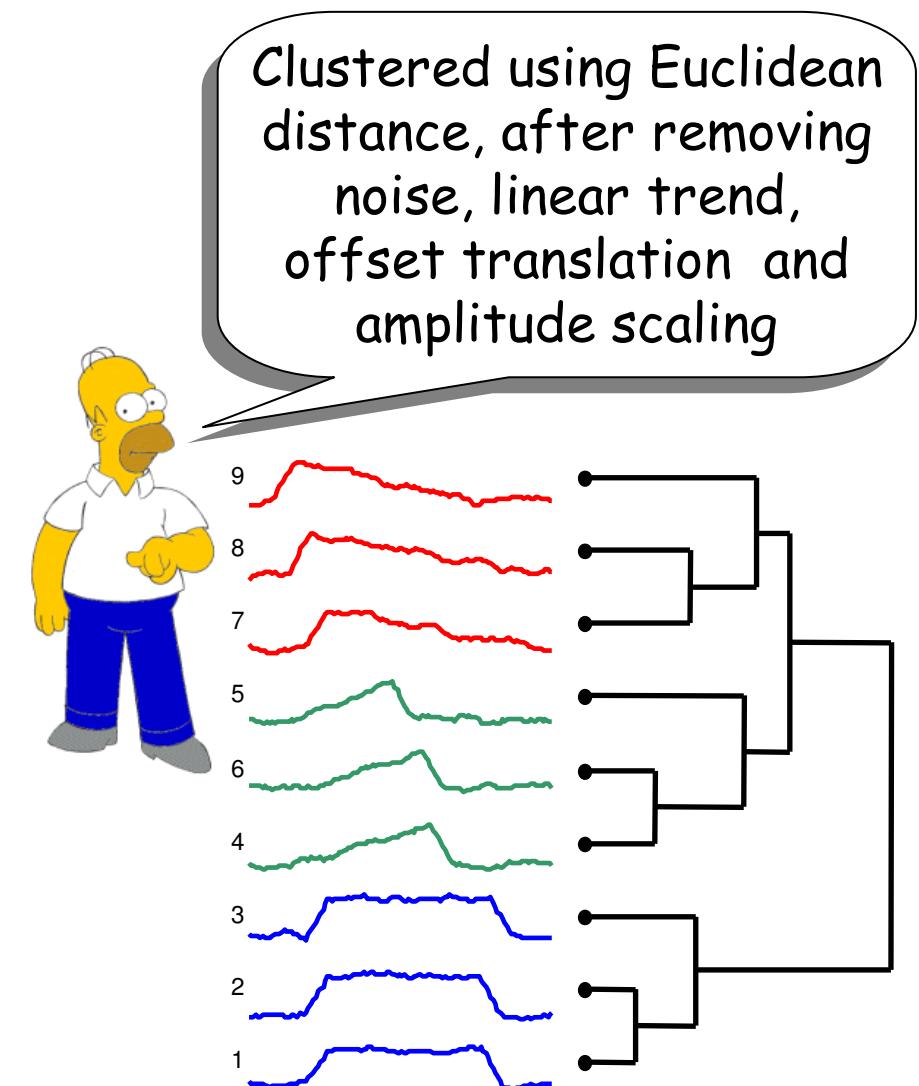
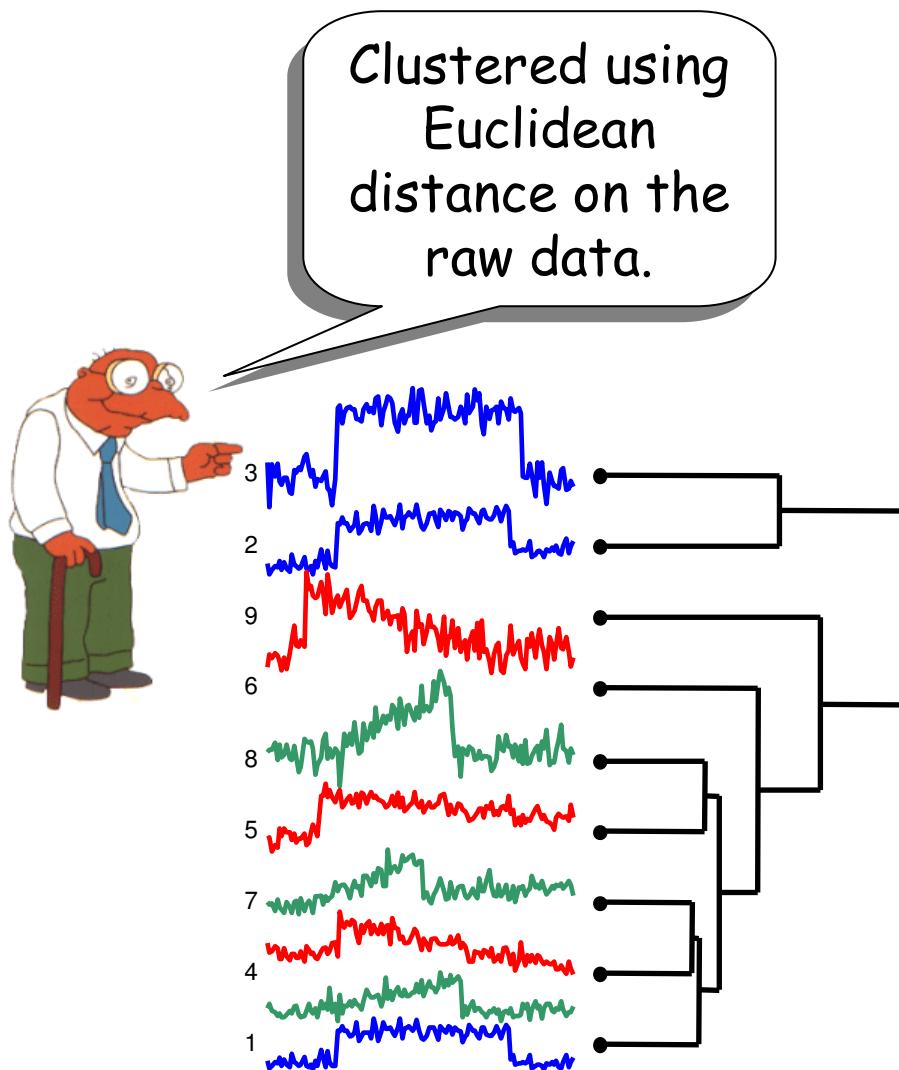
$$Q = \text{smooth}(Q)$$

$$C = \text{smooth}(C)$$

$$D(Q, C)$$

- C<sub>LINEAR</sub> SMOOTHING
- EXPONENTIAL SMOOTHING
  - HOLT
  - HOLT-WINTER
  - SAVGOL

# A Quick Experiment to Demonstrate the Utility of Preprocessing the Data



# Summary of Preprocessing

The “raw” time series may have distortions which we should remove before clustering, classification etc



Of course, sometimes the distortions are the most interesting thing about the data, the above is only a general rule

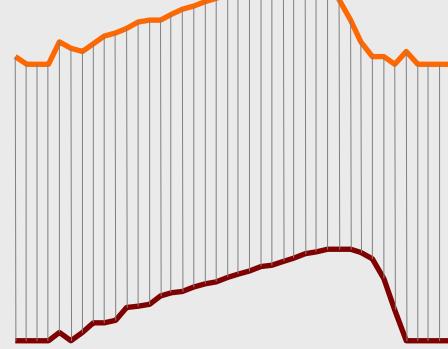
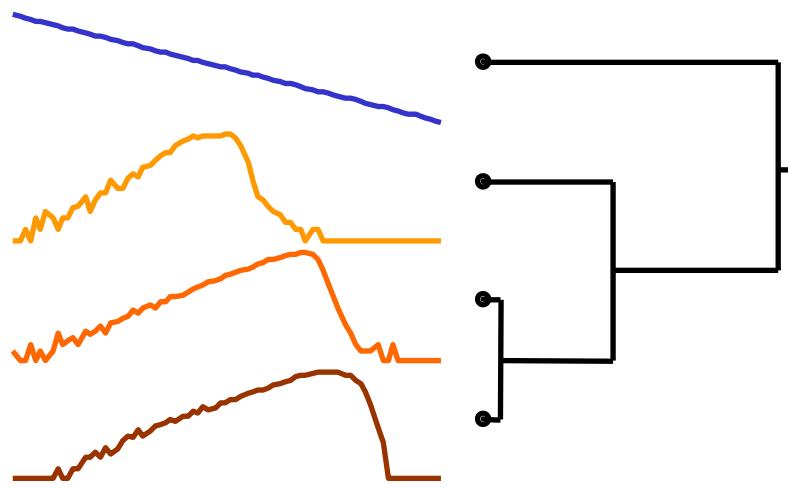
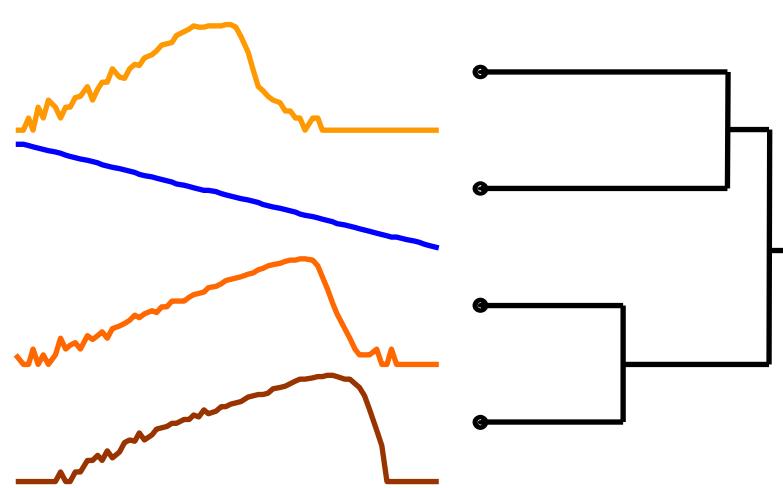


We should keep in mind these problems as we consider the high level representations of time series which we will encounter later (DFT, Wavelets etc). Since these representations often allow us to handle distortions in elegant ways

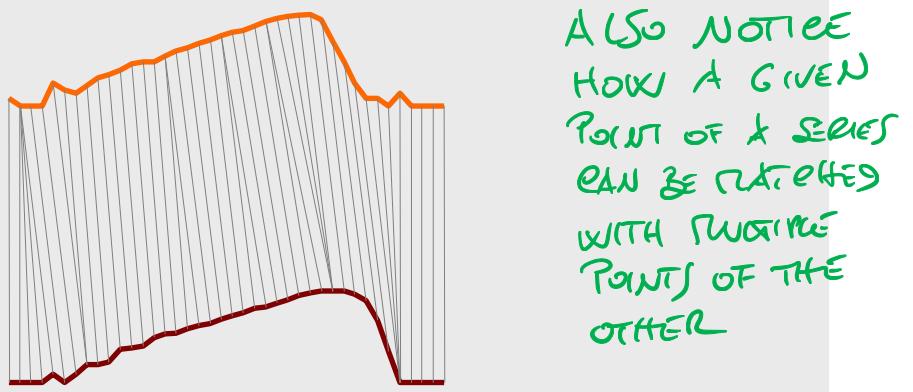


Sometimes, two time series can show a very similar behaviour, except for a slight temporal offset, or some temporal misalignments.

# Dynamic Time Warping



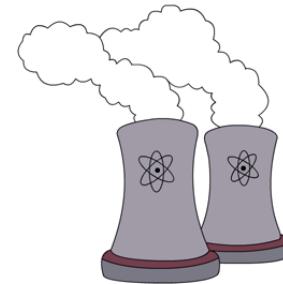
Fixed Time Axis  
*Sequences are aligned “one to one”.*



“Warped” Time Axis  
*Nonlinear alignments are possible.*

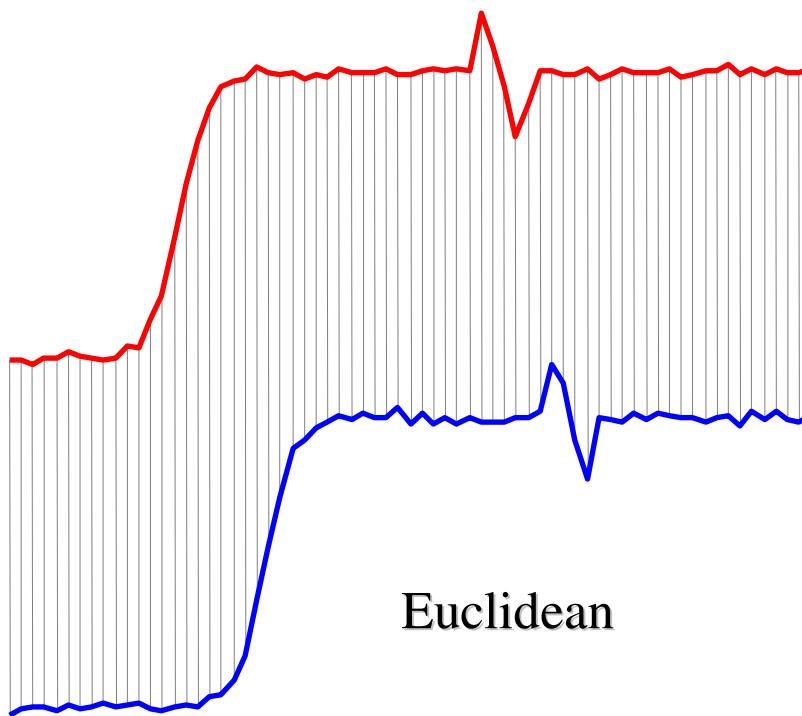
ALSO NOTICE  
HOW A GIVEN  
POINT OF A SERIES  
CAN BE RELATED  
WITH MULTIPLE  
POINTS OF THE  
OTHER

Note: We will first see the utility of DTW, then see how it is calculated.

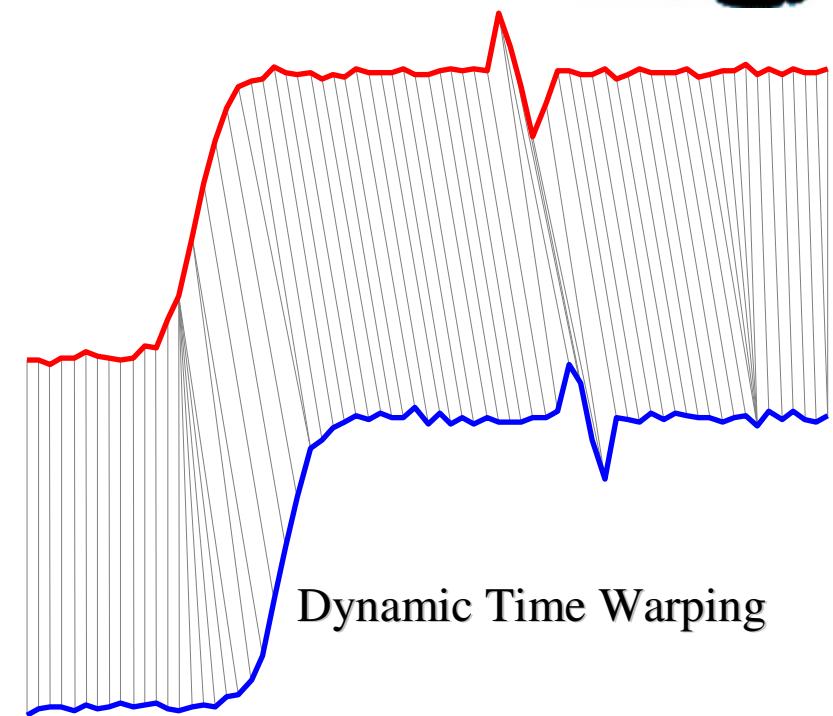


Here is another example on  
nuclear power plant trace  
data, to help you develop an  
intuition for DTW

Nuclear  
Power  
Excellent!

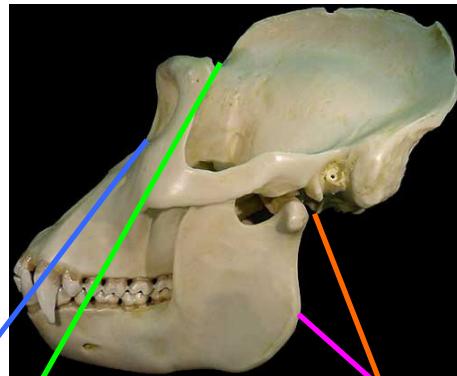


Euclidean

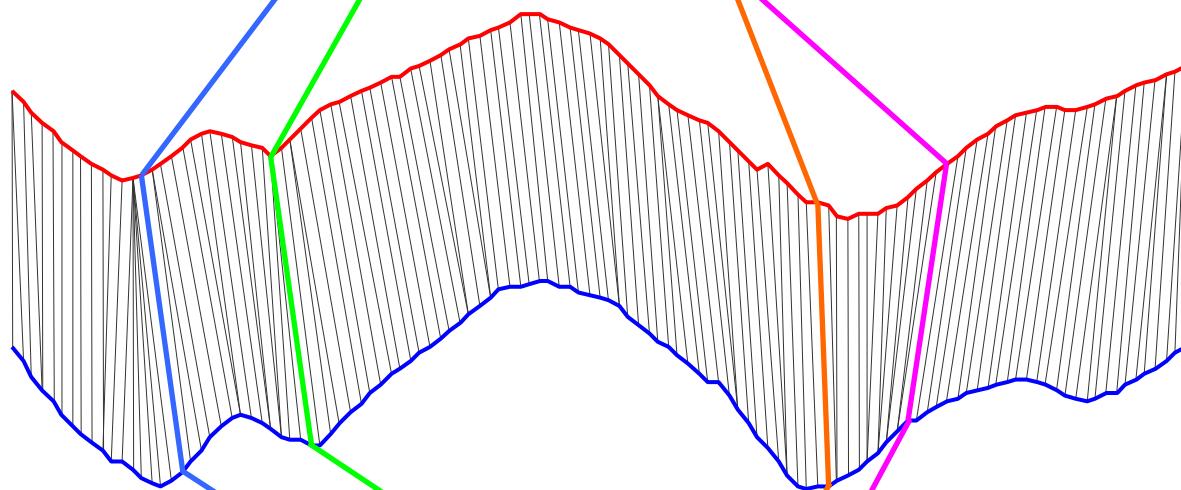


Dynamic Time Warping

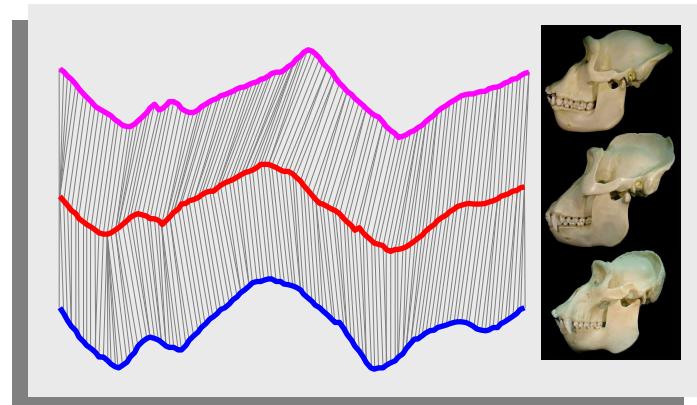
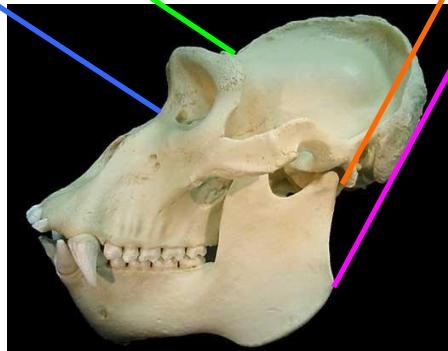
Lowland Gorilla  
*Gorilla gorilla graueri*



DTW is needed  
for most natural  
objects...

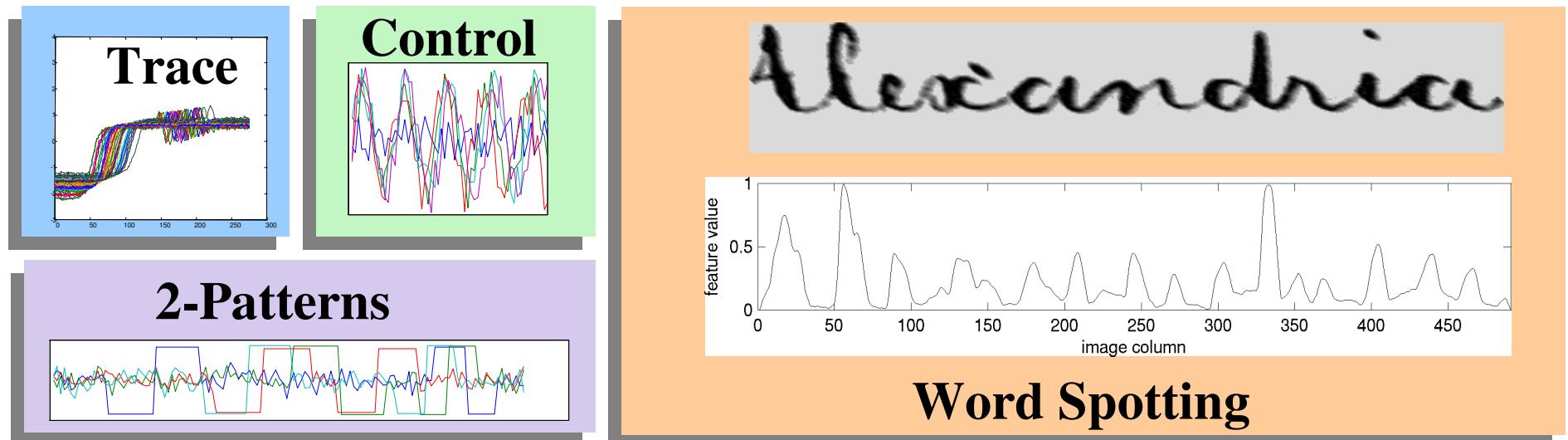
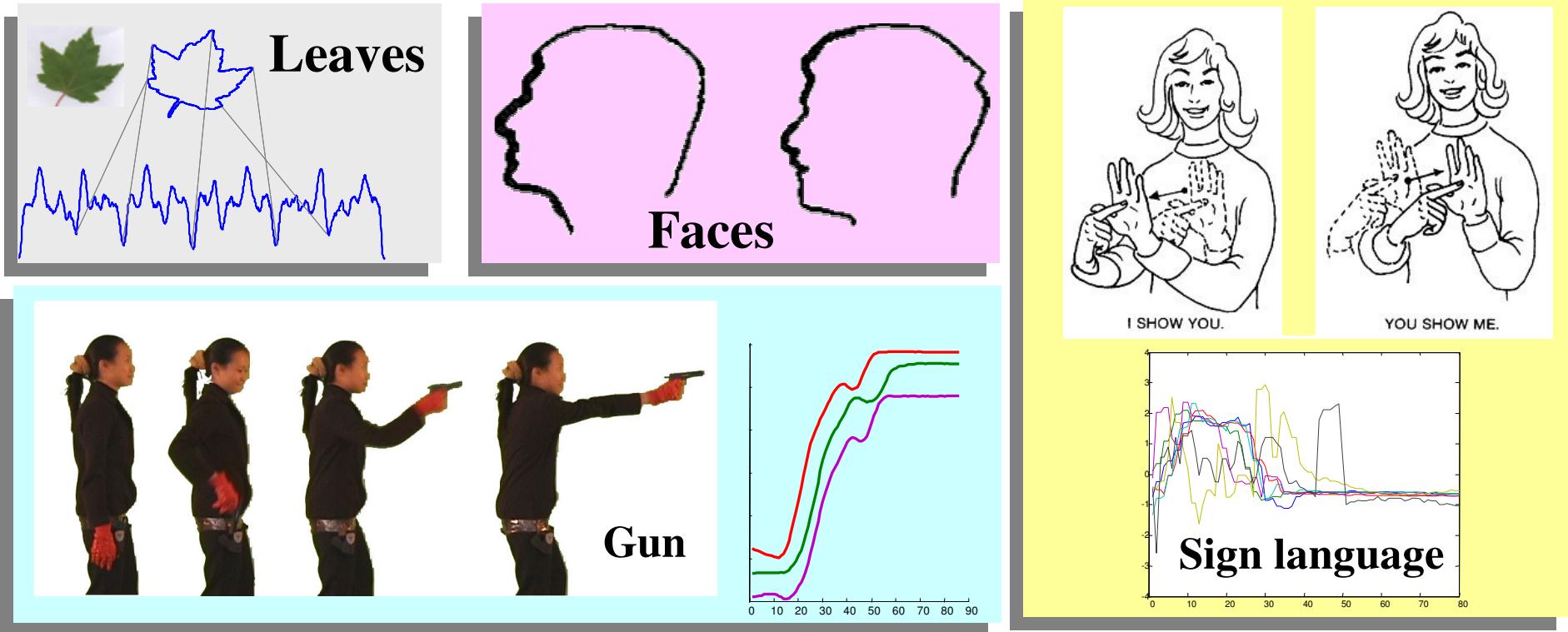


Mountain Gorilla  
*Gorilla gorilla beringei*



You can find all these datasets online on the UCR Time Series Classification Repository

Let us compare Euclidean Distance and DTW on some problems



# Results: Error Rate

Dataset	Euclidean	DTW
Word Spotting	4.78	1.10
Sign language	28.70	25.93
GUN	5.50	1.00
Nuclear Trace	11.00	0.00
Leaves <sup>#</sup>	33.26	4.07
(4) Faces	6.25	2.68
Control Chart*	7.5	0.33
2-Patterns	1.04	0.00

Using 1-nearest-neighbor,  
leaving-one-out  
evaluation!



# Results: Time (msec )

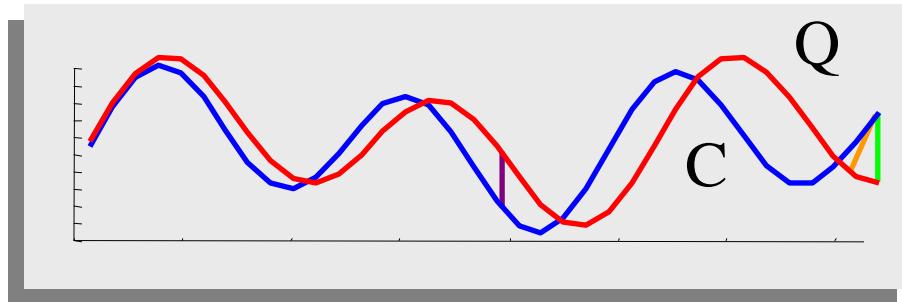
<b>Dataset</b>	<b>Euclidean</b>	<b>DTW</b>
Word Spotting	40	8,600
Sign language	10	1,110
GUN	60	11,820
Nuclear Trace	210	144,470
Leaves	150	51,830
(4) Faces	50	45,080
Control Chart	110	21,900
2-Patterns	16,890	545,123

DTW is  
two to  
three  
orders of  
magnitude  
slower  
than  
Euclidean  
distance

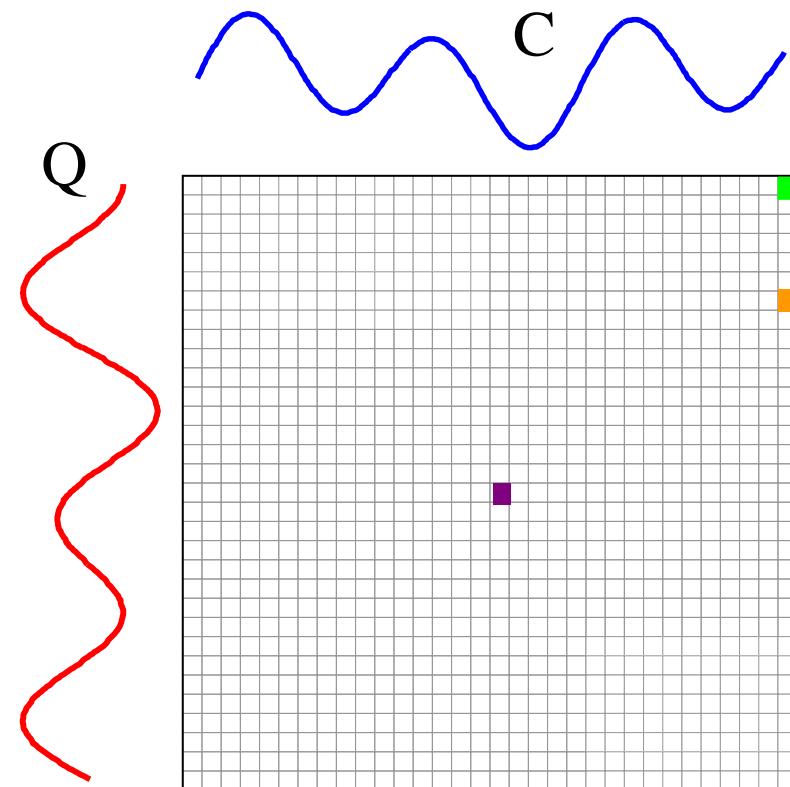


# How is DTW Calculated? I

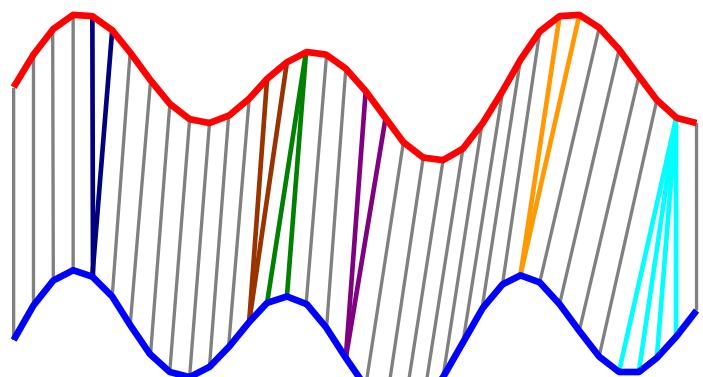
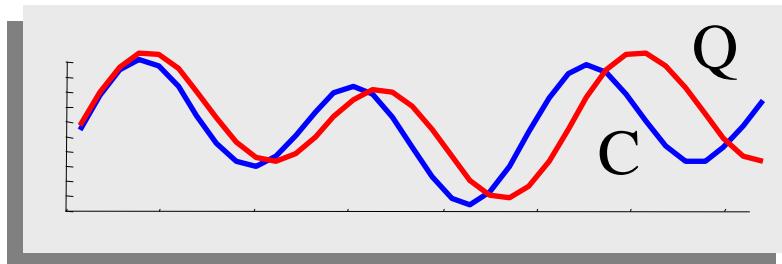
We create a matrix the size of  $|Q|$  by  $|C|$ , then fill it in with the distance between every pair of point in our two time series.



Thus, on the diagonal, you have the results of the classical Euclidean Distance metric applied on the two time series



# How is DTW Calculated? II



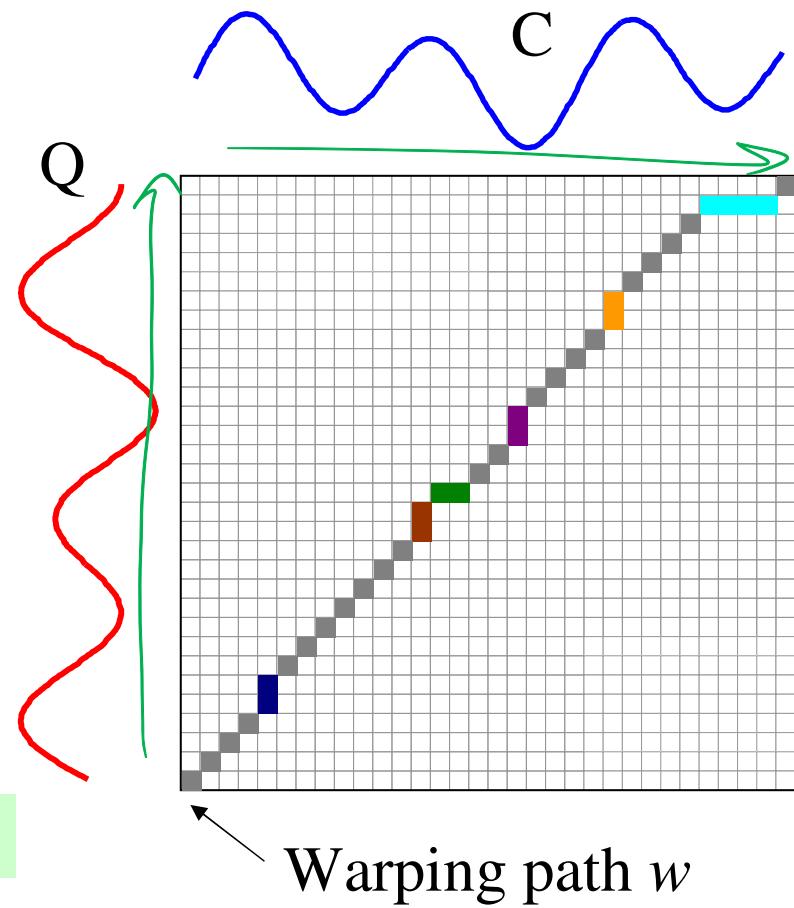
This recursive function gives us the minimum cost path

$$\gamma(i,j) = d(q_i, c_j) + \min\{ \gamma(i-1, j-1), \gamma(i-1, j), \gamma(i, j-1) \}$$

Every possible warping between two time series, is a path through the matrix. We want the best one...

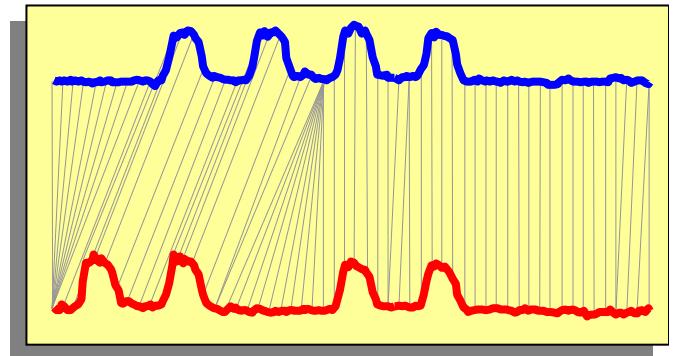
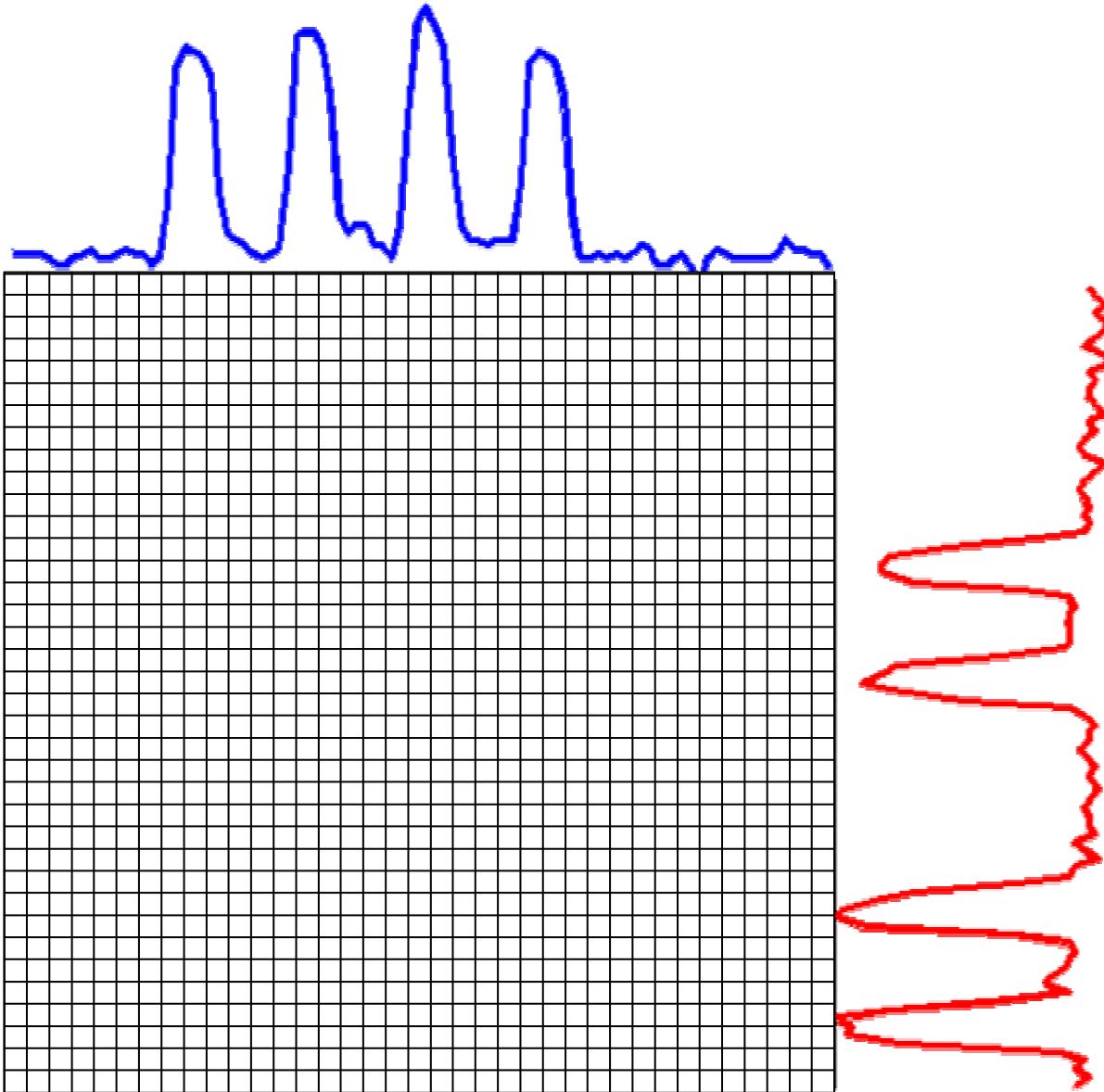
$$DTW(Q, C) = \min \left\{ \sqrt{\sum_{k=1}^K w_k} \right\} / K$$

where K is the warping path length



"THE DISTANCE AT THE  $(i,j)$  CELL, PLUS THE MINIMUM COST IN TERMS OF DISTANCES ALONG THE PATHS THAT IN THE MATRIX ARRIVE AT THAT CELL."

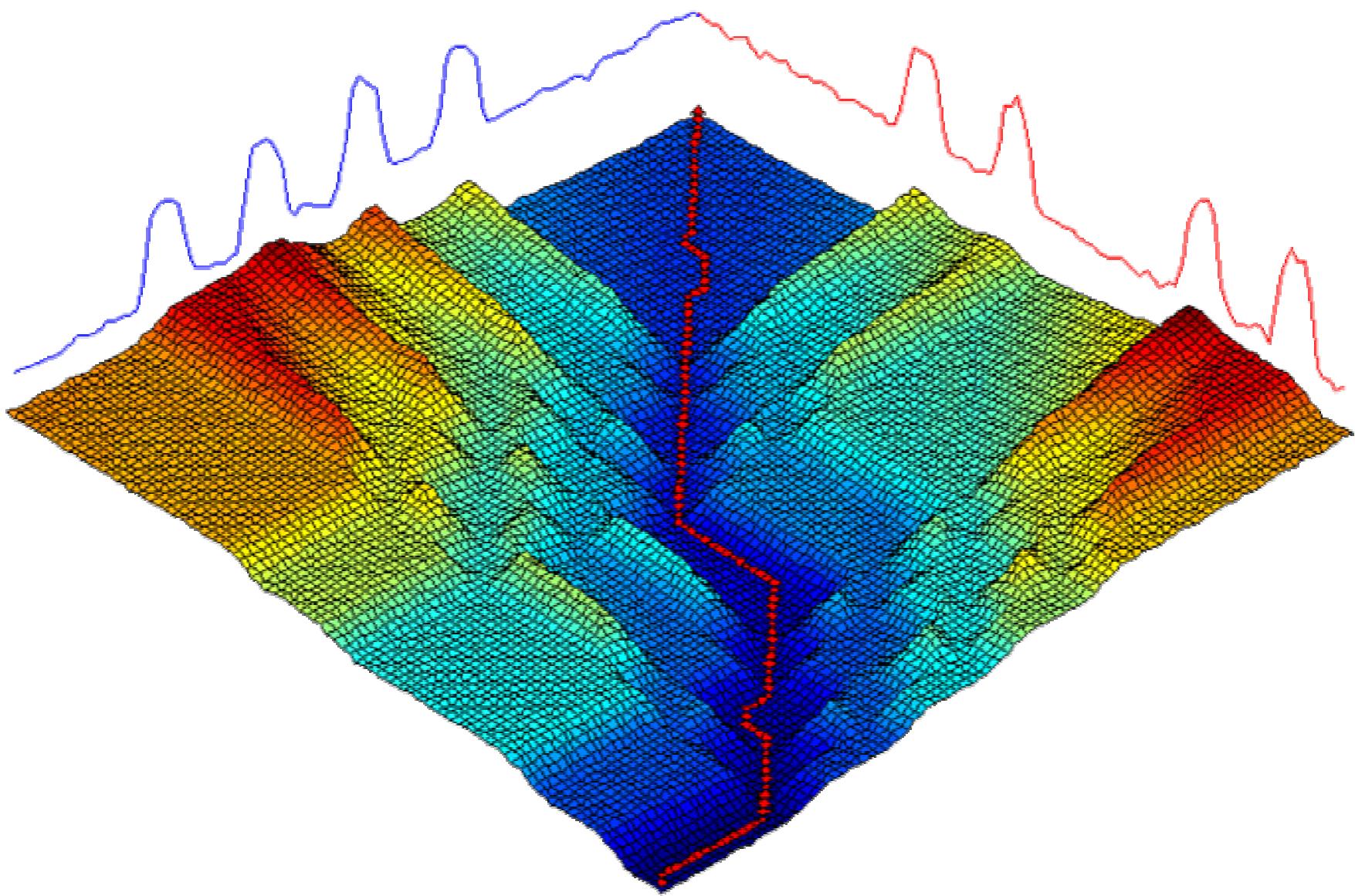
# Let us visualize the cumulative matrix on a real world problem I



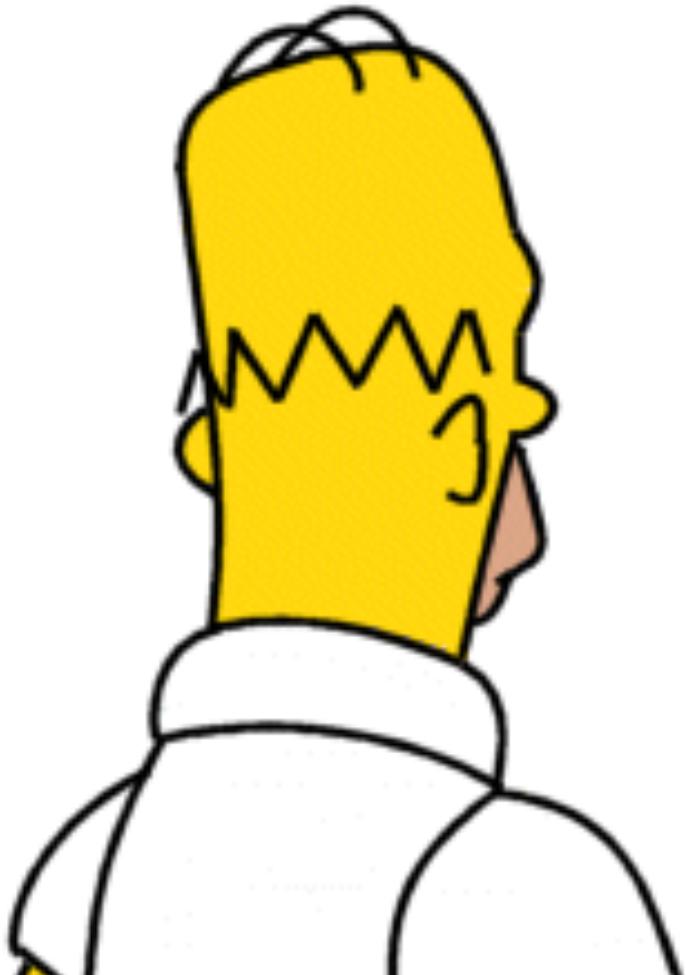
This example shows 2 one-week periods from the power demand time series.

Note that although they both describe 4-day work weeks, the blue sequence had Monday as a holiday, and the red sequence had Wednesday as a holiday.

Let us visualize the cumulative matrix on a real world problem II



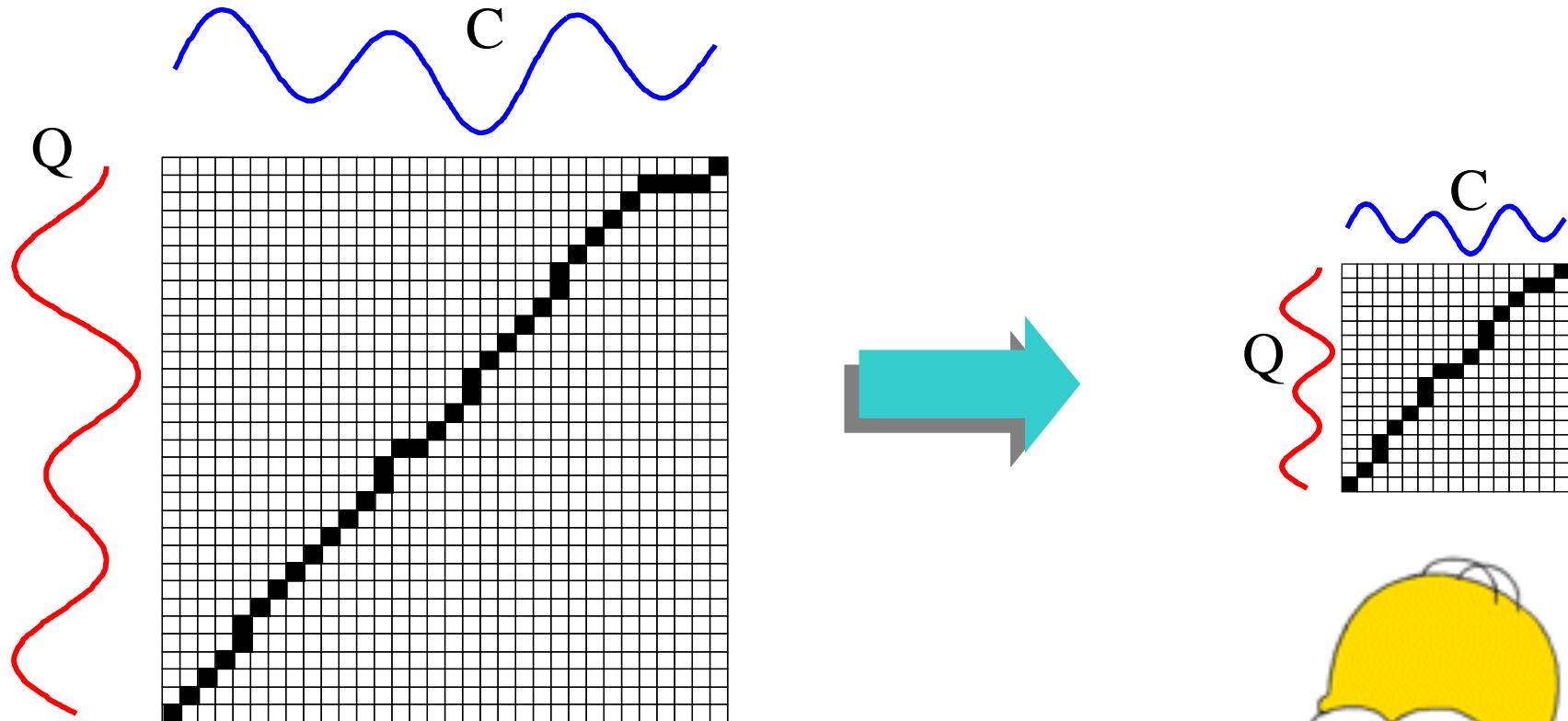
# What we have seen so far...



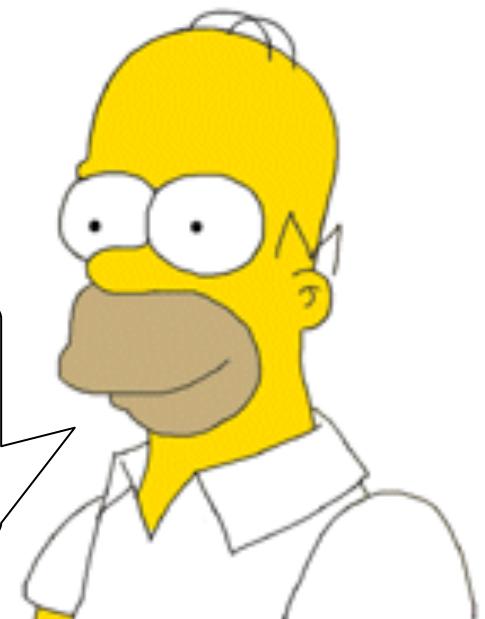
- Dynamic Time Warping gives **much better** results than Euclidean distance on virtually all problems.
- Dynamic Time Warping is very very slow to calculate!

Is there anything we can do to speed up similarity search under DTW?

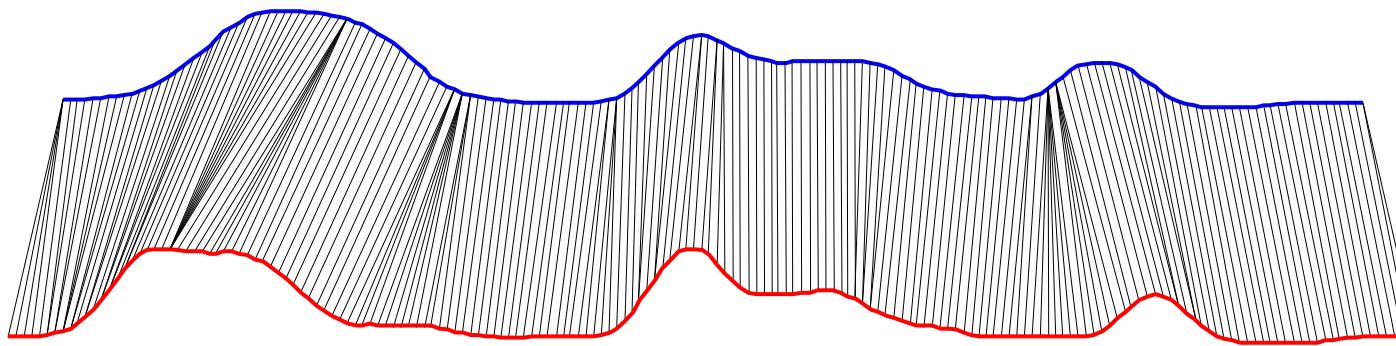
# Fast Approximations to Dynamic Time Warp Distance I



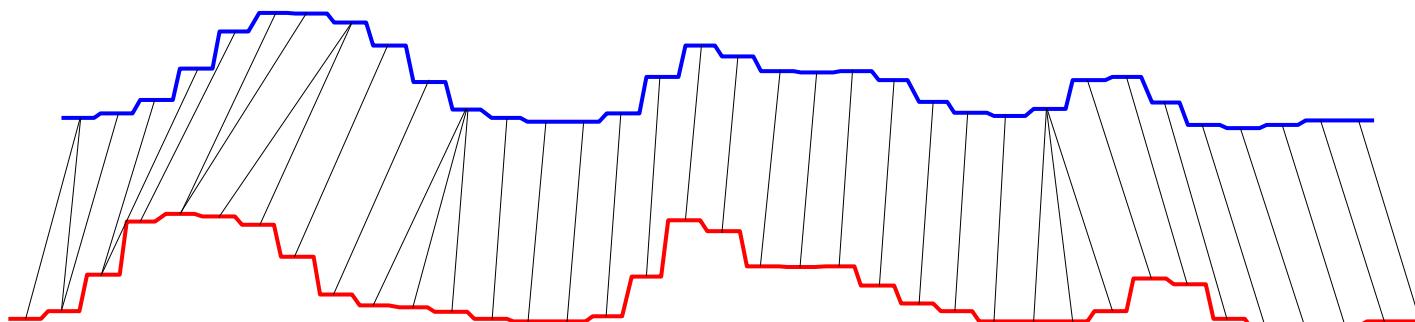
Simple Idea: Approximate the time series with some compressed or downsampled representation, and do DTW on the new representation. How well does this work...



## Fast Approximations to Dynamic Time Warp Distance II



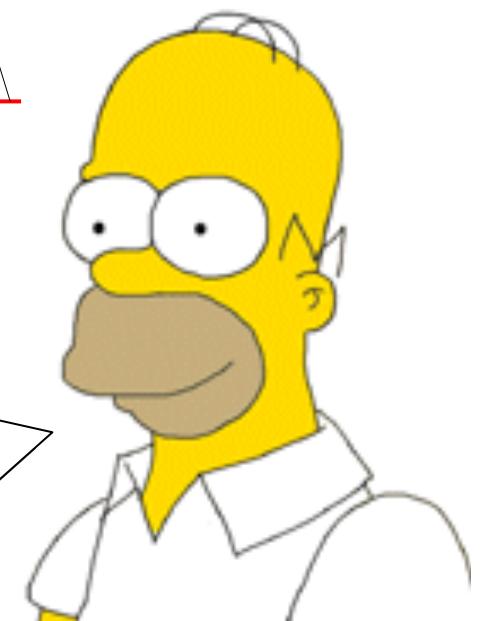
1.03 sec



0.07 sec

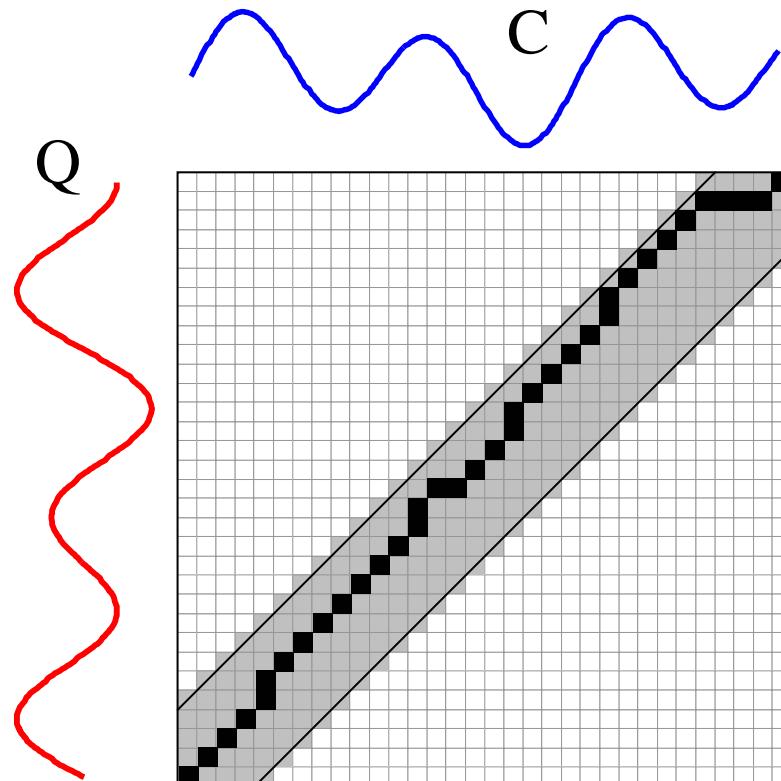
... there is strong visual evidence to suggests it works well

There is good experimental evidence for the utility of the approach on clustering, classification, etc

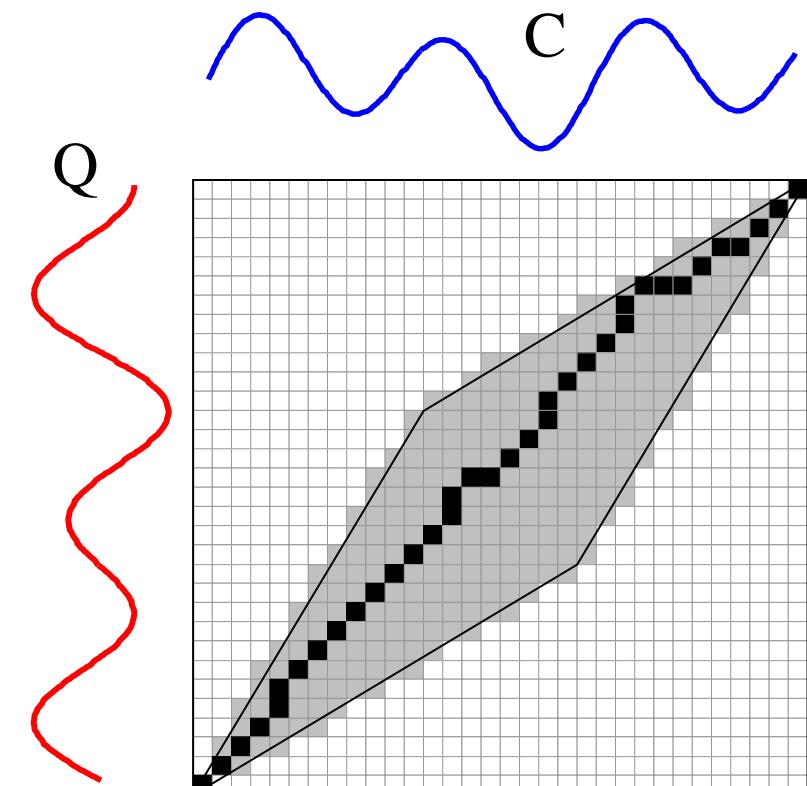


# Global Constraints

- Slightly speed up the calculations
- Prevent pathological warpings



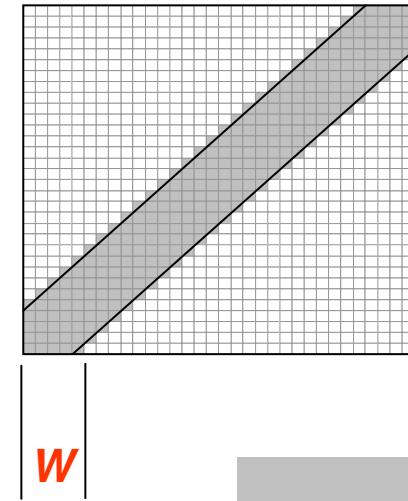
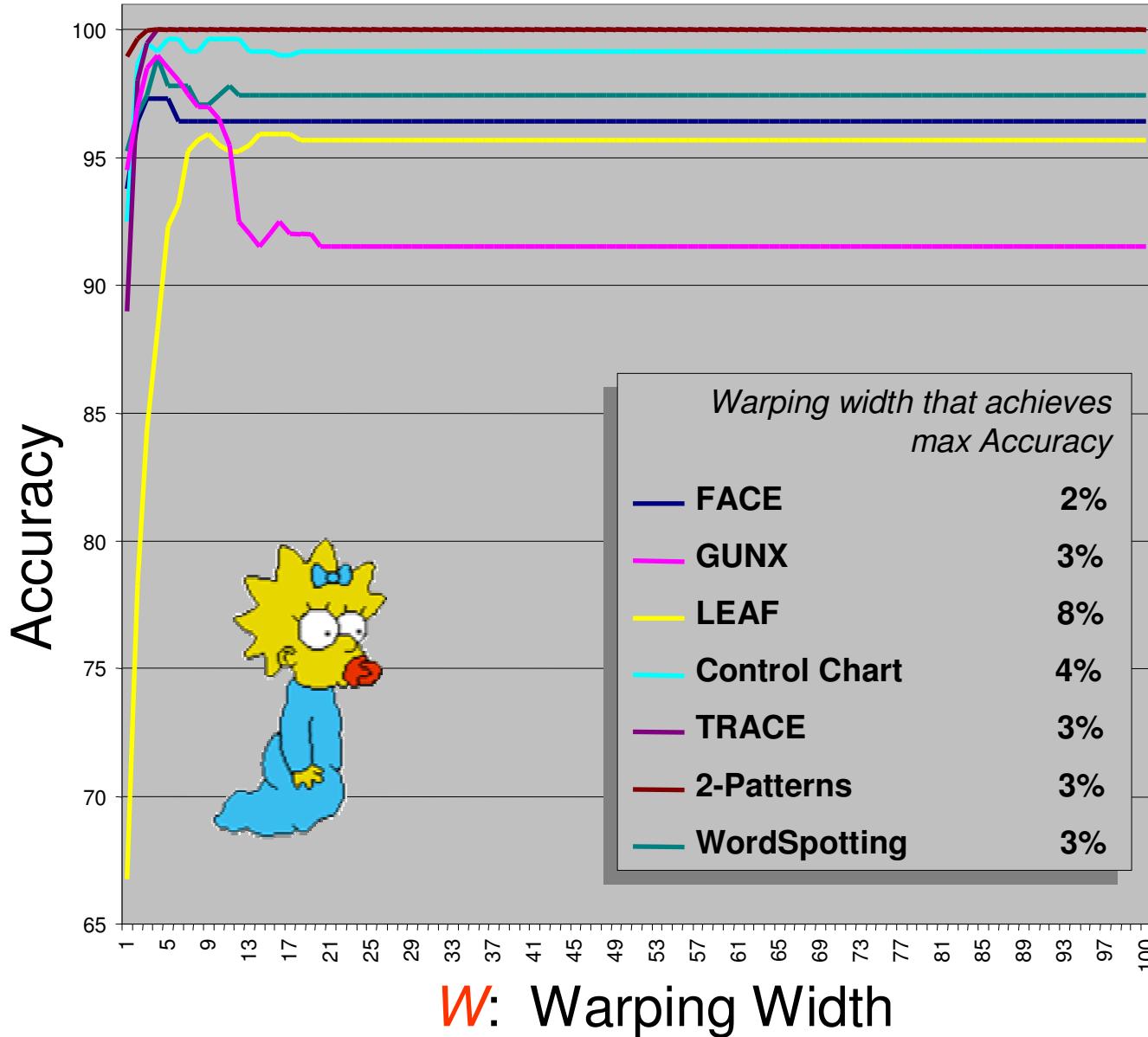
Sakoe-Chiba Band



Itakura Parallelogram

SOMETIMES, THE "WARPING" APPROXIMATION  
ACCURSES US EVEN TO IMPROVE THE PERFORMANCE!

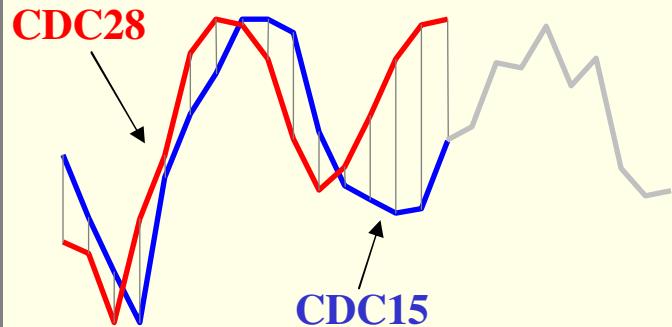
## Accuracy vs. Width of Warping Window



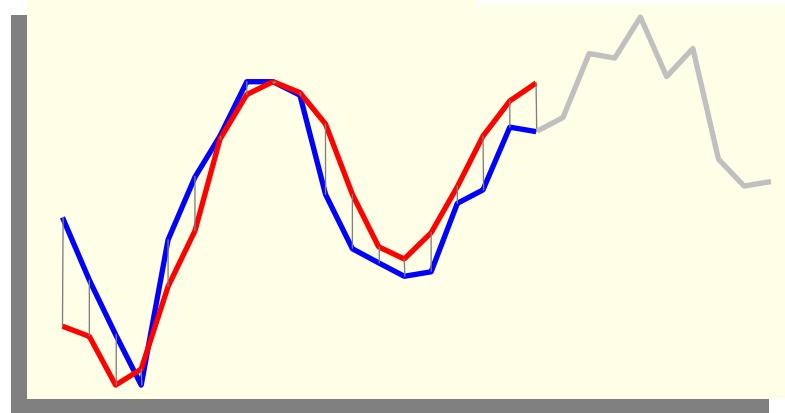
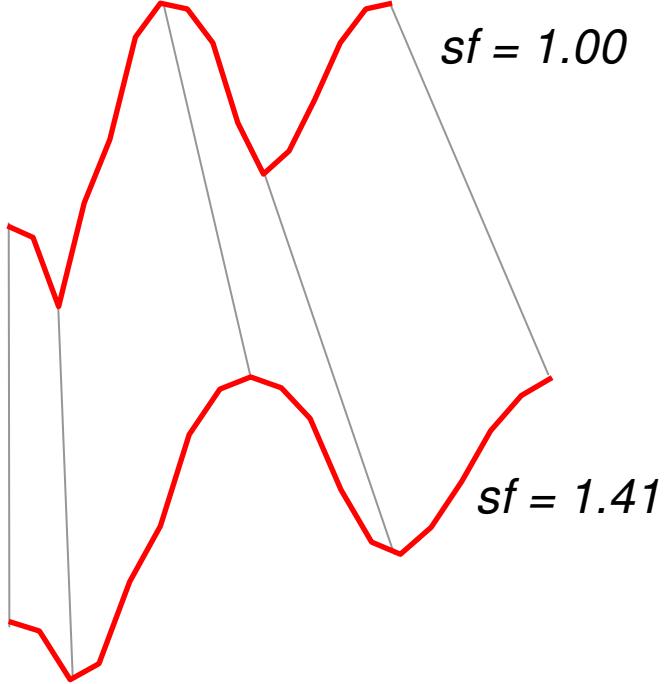
A global time scaling on a time series uniformly squeezes or stretches the time series on the time axis.

# Uniform Scaling I

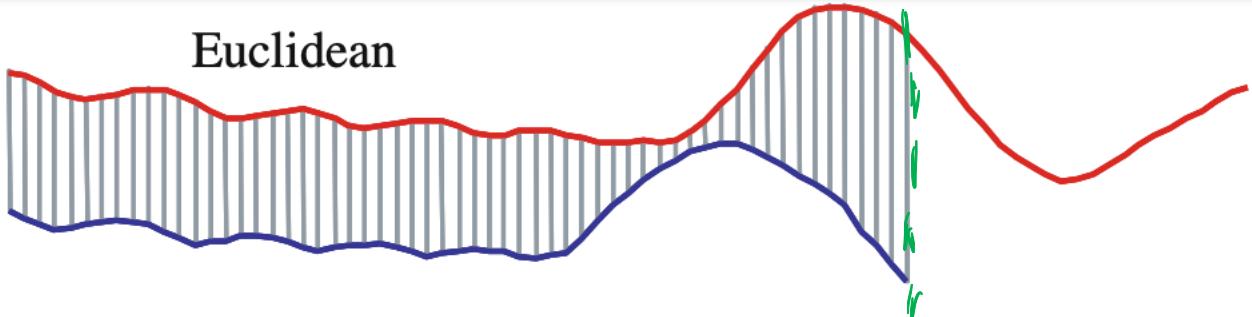
Two genes that are known to be functionally related...



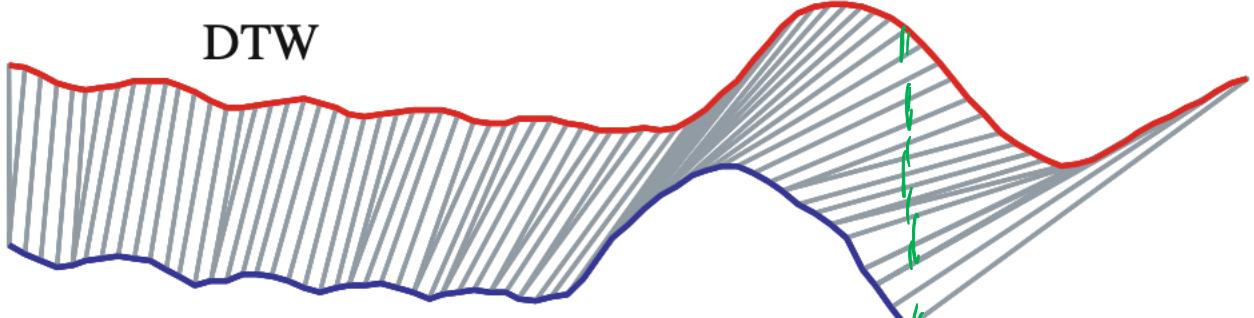
Sometimes  
global or  
*uniform scaling*  
is as important  
as DTW



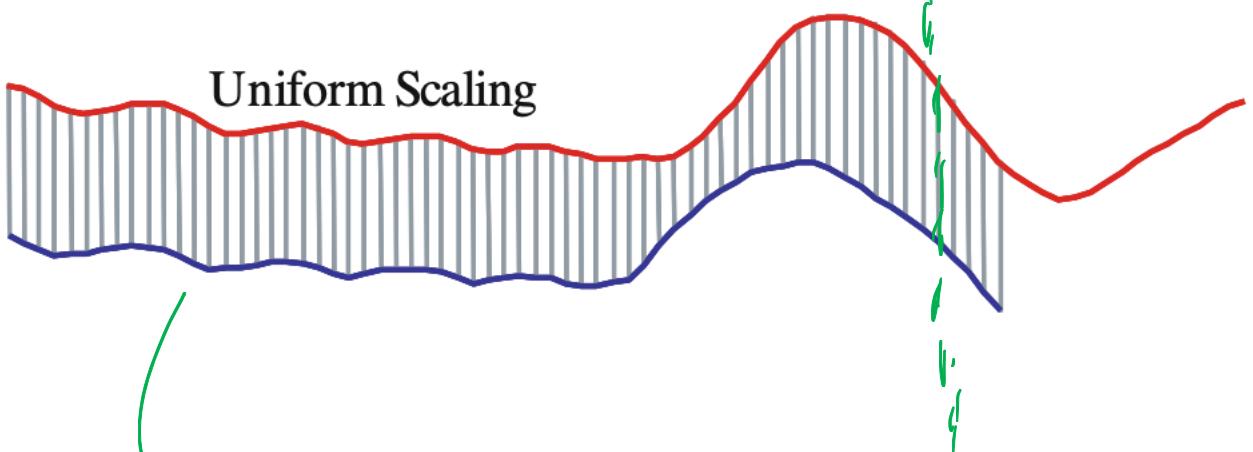
Euclidean



DTW



Uniform Scaling



→ HERE, THE BLUE SERIES  
HAS BEEN "STRETCHED"  
IN ORDER TO ALIGN  
IT WITH THE RED ONE -

# Only Euclidean and DTW Distance are Useful

**Stop!**

What about the dozens of other techniques for measuring time series shape similarity?

NO  
FREE  
LUNCH  
THEOREM



Unfortunately, none of them appear to be useful!



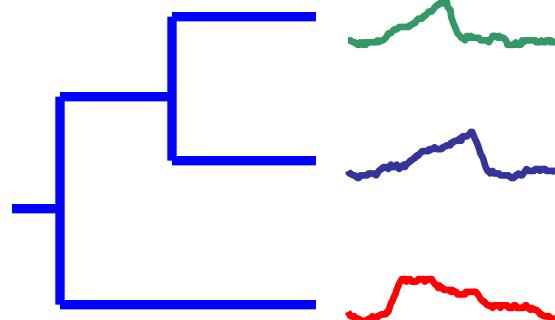
# Classification Error Rates on two publicly available datasets



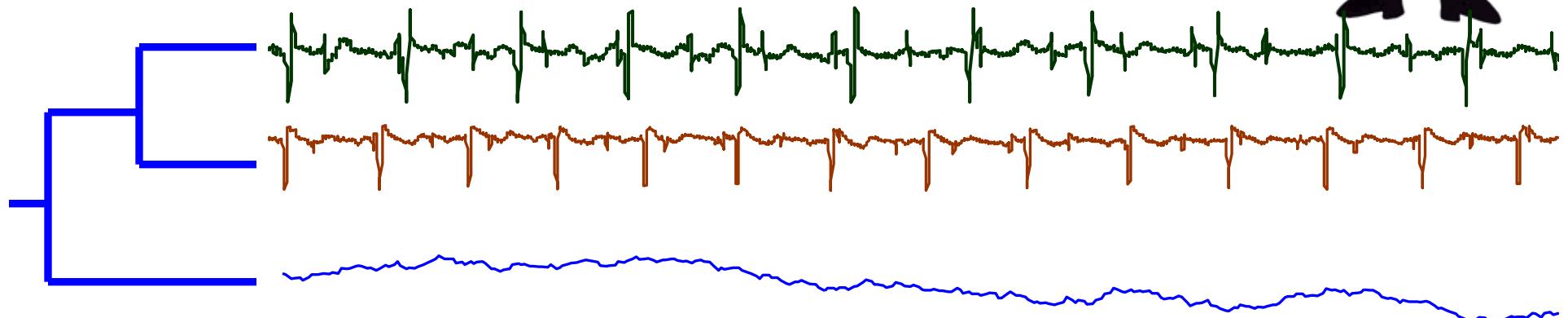
Approach	Cylinder-Bell-F'	Control-Chart
<i>Euclidean Distance</i>	<b>0.003</b>	<b>0.013</b>
Aligned Subsequence	<b>0.451</b>	<b>0.623</b>
Piecewise Normalization	<b>0.130</b>	<b>0.321</b>
Autocorrelation Functions	<b>0.380</b>	<b>0.116</b>
Cepstrum	<b>0.570</b>	<b>0.458</b>
String (Suffix Tree)	<b>0.206</b>	<b>0.578</b>
Important Points	<b>0.387</b>	<b>0.478</b>
Edit Distance	<b>0.603</b>	<b>0.622</b>
String Signature	<b>0.444</b>	<b>0.695</b>
Cosine Wavelets	<b>0.130</b>	<b>0.371</b>
Hölder	<b>0.331</b>	<b>0.593</b>
Piecewise Probabilistic	<b>0.202</b>	<b>0.321</b>

# Two Kinds of Similarity

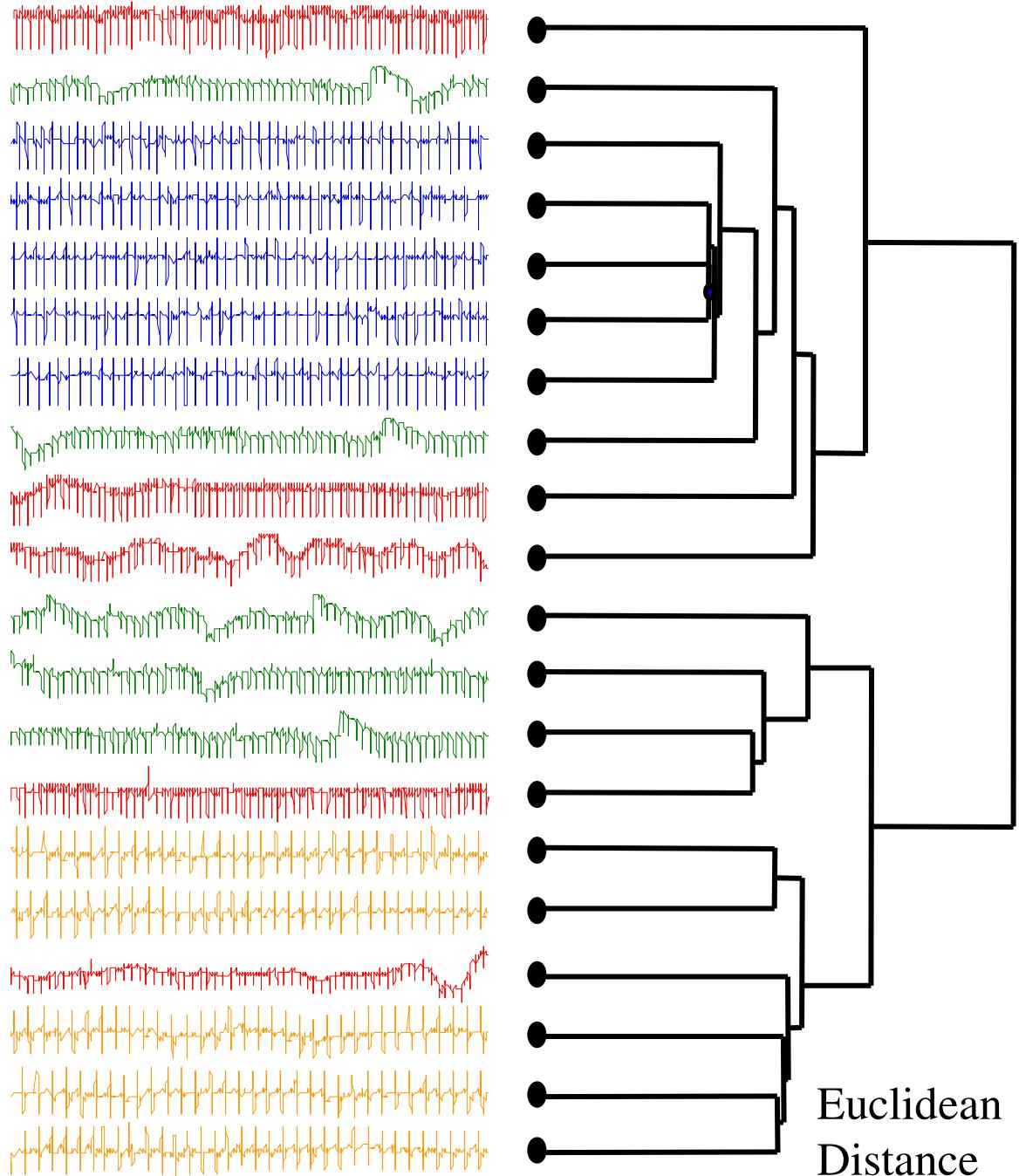
We are  
done with  
*shape*  
similarity



Let us consider  
similarity at  
the *structural*  
level for the  
next 10 minutes



For long time series, shape based similarity will give very poor results. We need to measure similarly based on high level structure

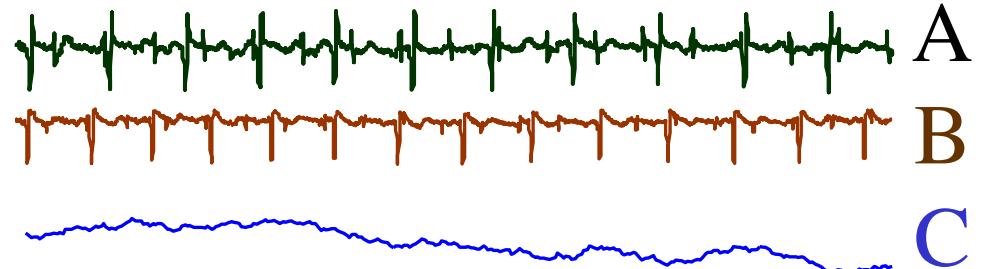


# Structure or Model Based Similarity

The basic idea is to extract *global* features from the time series, create a feature vector, and use these feature vectors to measure similarity and/or classify



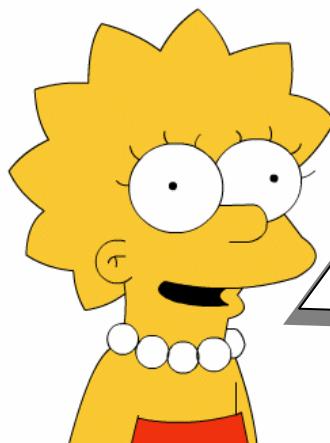
But which  
• **features**?  
• **distance measure/learning algorithm?**



Feature \ Time Series	A	B	C
Max Value	11	12	19
Autocorrelation	0.2	0.3	0.5
Zero Crossings	98	82	13
Average	0.3	0.4	0.1
...	...	...	...

# Feature-based Classification of Time-series Data

Nanopoulos, Alcock, and Manolopoulos



Makes sense, but when we looked at the *same* dataset, we found we could be better classification accuracy with Euclidean distance!

## Features

mean

variance

skewness

kurtosis

mean (1<sup>st</sup> derivative)

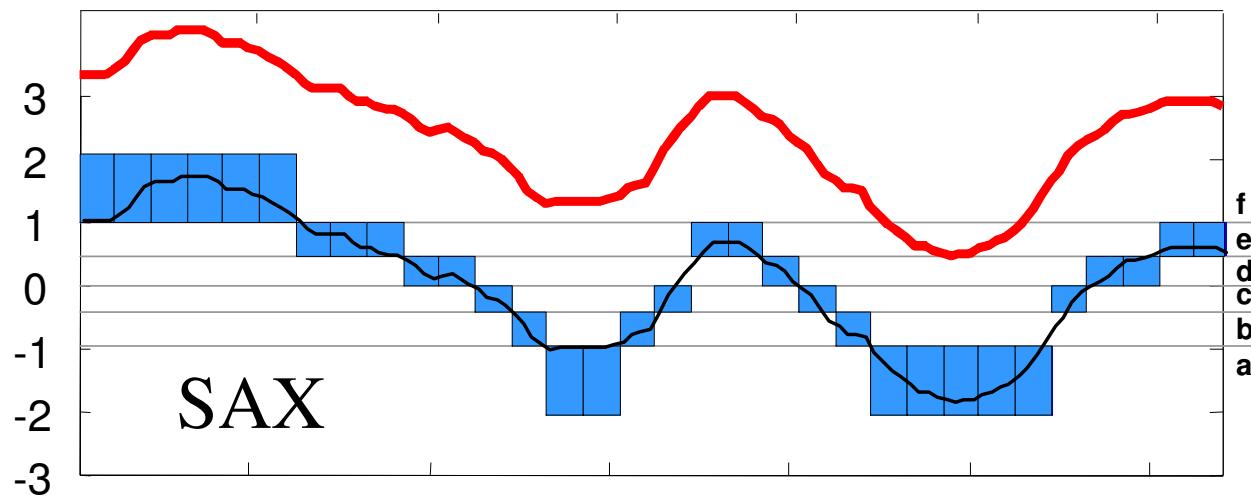
variance (1<sup>st</sup> derivative)

skewness (1<sup>st</sup> derivative)

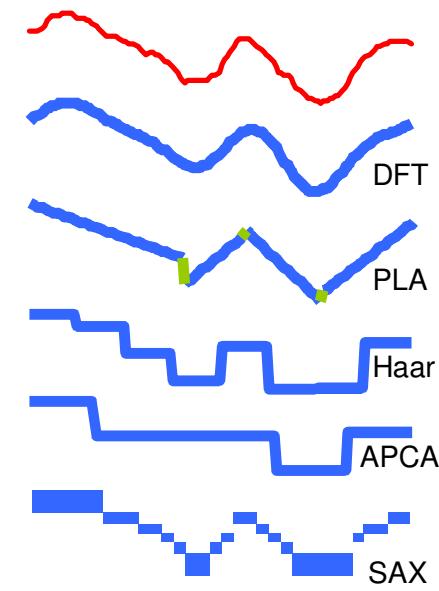
kurtosis (1<sup>st</sup> derivative)

# Exploiting Symbolic Representations of Time Series

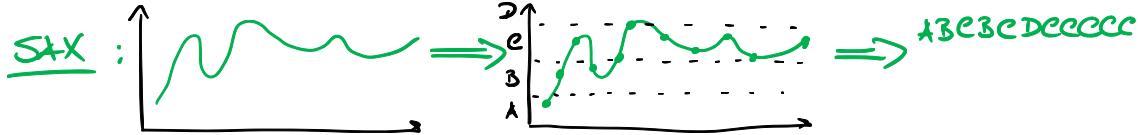
There is now a lower bounding dimensionality reducing time series representation! It is called SAX (Symbolic Aggregate ApproXimation)



**fffffffeeeddcbaabceeedcbaaaaacddee**



> Then you can rely on algorithms developed for text mining and bioinformatics!



# Compression Based Dissimilarity

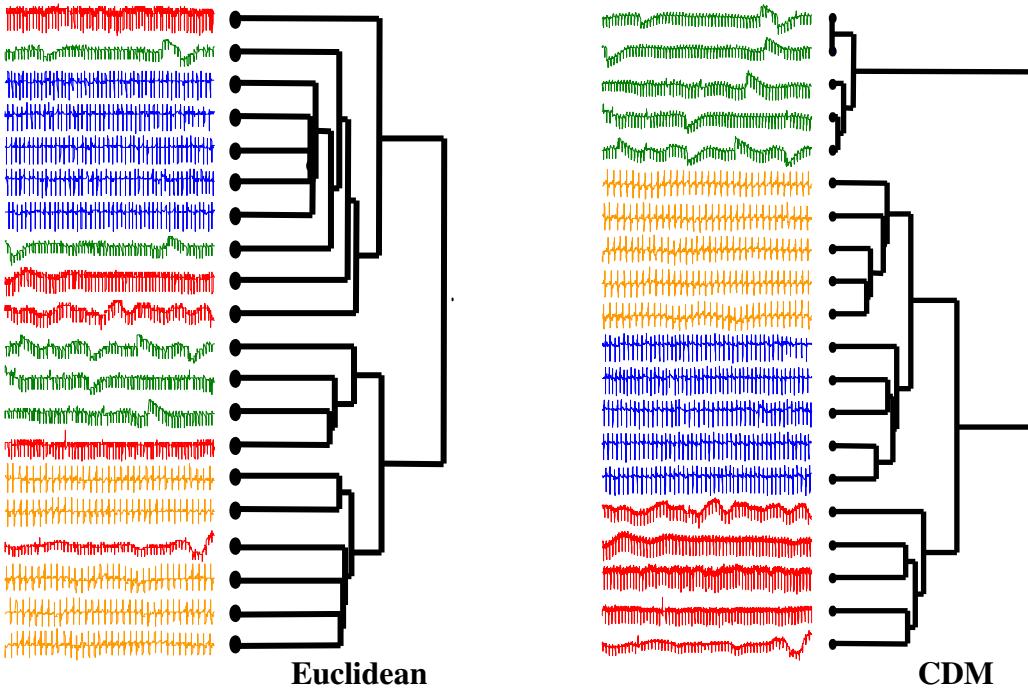
(In general) Li, Chen, Li, Ma, and Vitányi: (For time series) Keogh, Lonardi and Ratanamahatana

- features?
- distance measure/  
learning algorithm?

## Distance Measure

### Co-Compressibility

$$CDM(x, y) = \frac{C(xy)}{C(x) + C(y)}$$



## Features

Whatever structure  
the compression  
algorithm finds...

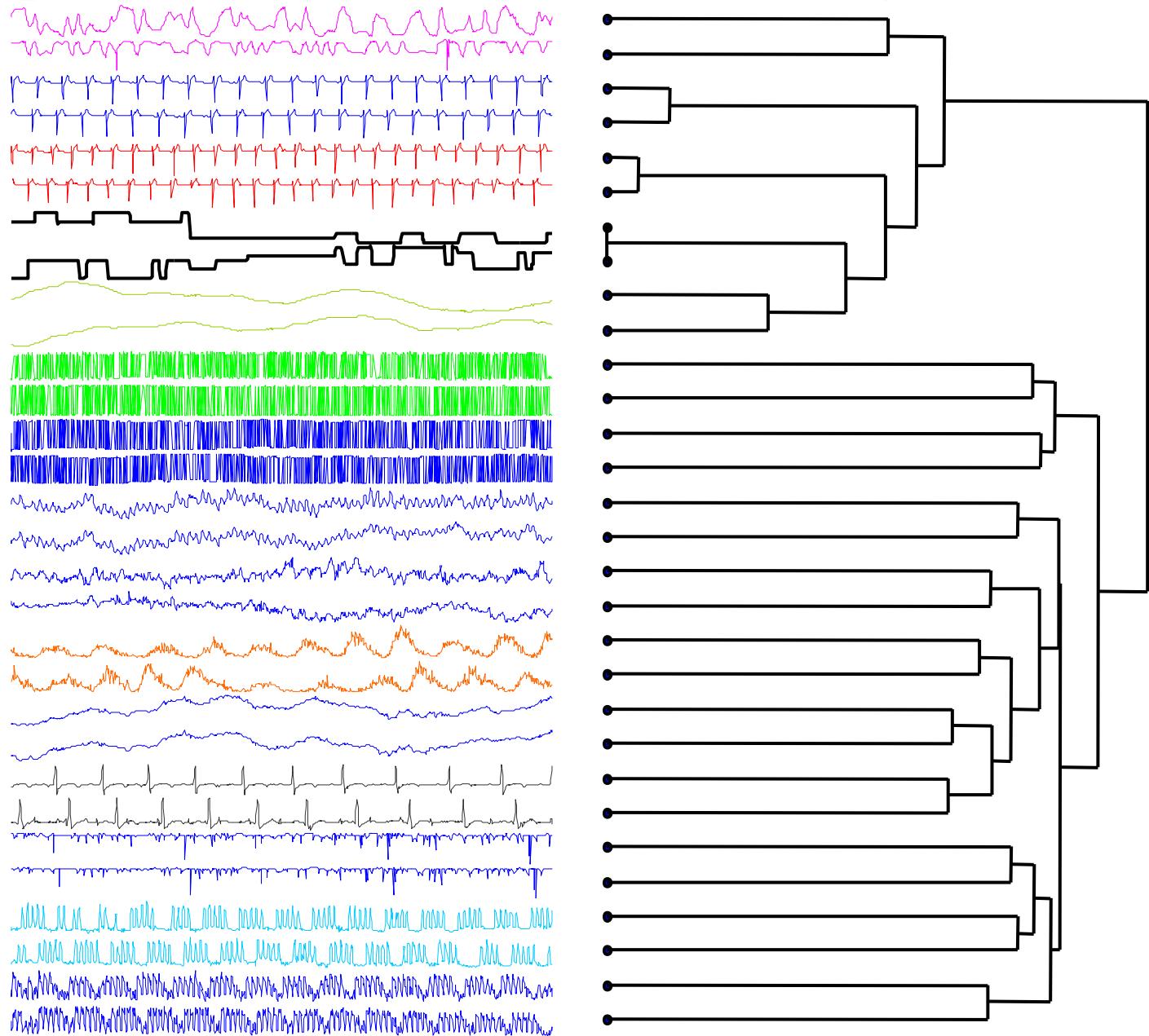
The time series is first converted  
to the SAX symbolic  
representation\*

SAX transforms a time-series  $X$  of length  $n$  into a string of arbitrary length  $\omega$ , where  $\omega \ll n$  typically, using an alphabet  $A$  of size  $a > 2$ .

Then, the principle of CBD is: the more patterns two strings share, the smaller is the compressed file size of their concatenated string.

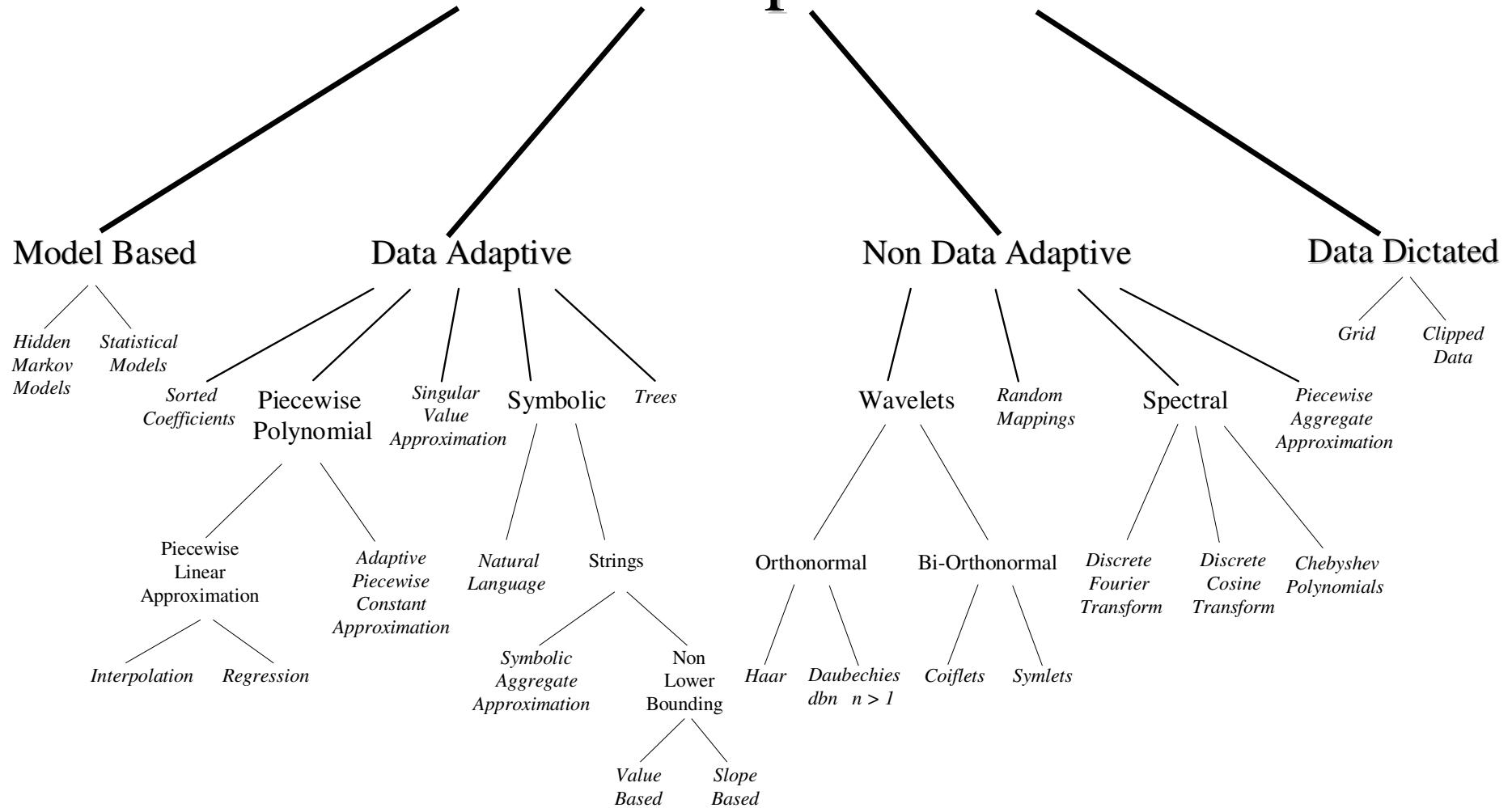
# Compression Based Dissimilarity

Reel 2: Tension  
Reel 2: Angular speed  
Koski ECG: Fast 2  
Koski ECG: Fast 1  
Koski ECG: Slow 2  
Koski ECG: Slow 1  
Dryer hot gas exhaust  
Dryer fuel flow rate  
Ocean 2  
Ocean 1  
Evaporator: vapor flow  
Evaporator: feed flow  
Furnace: cooling input  
Furnace: heating input  
Great Lakes (Ontario)  
Great Lakes (Erie)  
Buoy Sensor: East Salinity  
Buoy Sensor: North Salinity  
Sunspots: 1869 to 1990  
Sunspots: 1749 to 1869  
Exchange Rate: German Mark  
Exchange Rate: Swiss Franc  
Foetal ECG thoracic  
Foetal ECG abdominal  
Balloon2 (lagged)  
Balloon1  
Power : April-June (Dutch)  
Power : Jan-March (Dutch)  
Power : April-June (Italian)  
Power : Jan-March (Italian)



OF COURSE, SAX IS NOT THE ONLY APPROACH TO ENCODE THE CONTENT OF A TIME SERIES AND CREATE A MODEL OF IT...

# Time Series Representations



# Summary of Time Series Similarity

- If you have *short* time series, use DTW after searching over the warping window size
  - Also, consider Uniform Scaling, and preprocessing
- If you have *long* time series, and you know nothing about your data, try compression based dissimilarity. (after converting it with SAX)
- If you do know something about your data, try to leverage of this knowledge to extract features.

⇒ SOME OTHER TIME-SERIES RELATED TASKS

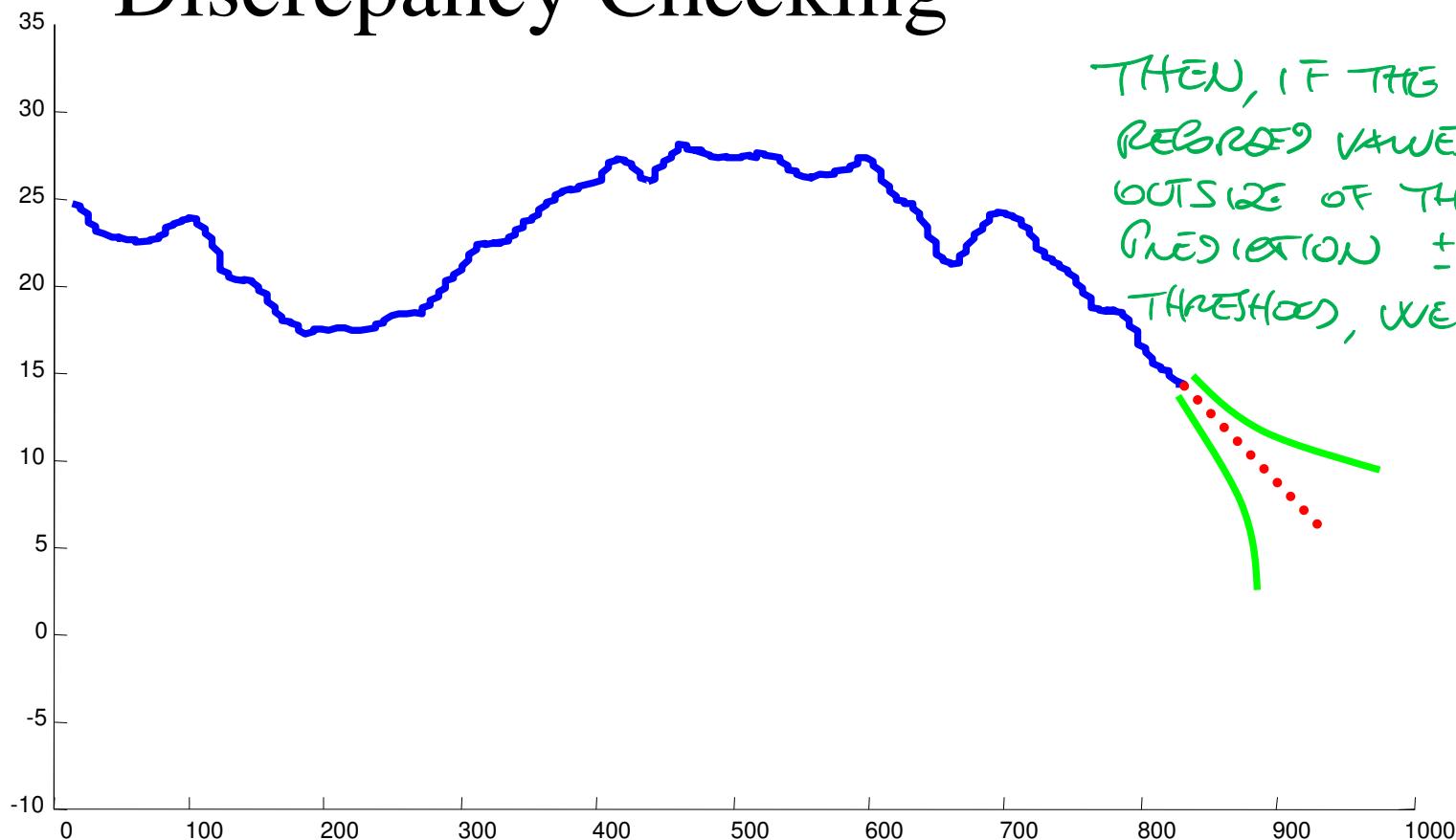
# Simple Approaches I

## Limit Checking



# Simple Approaches II

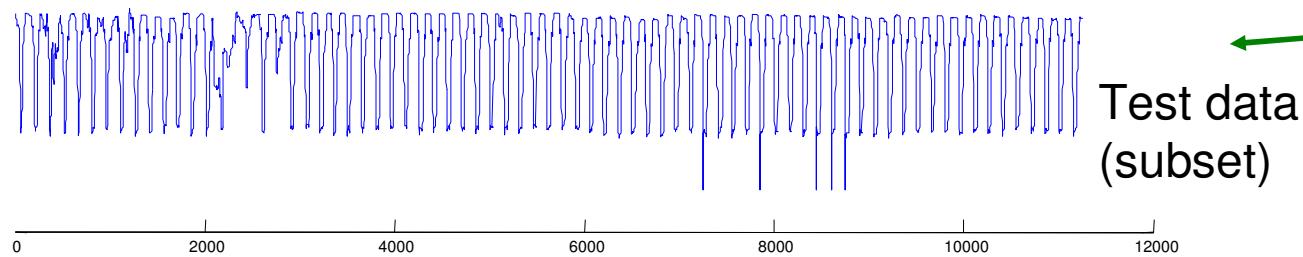
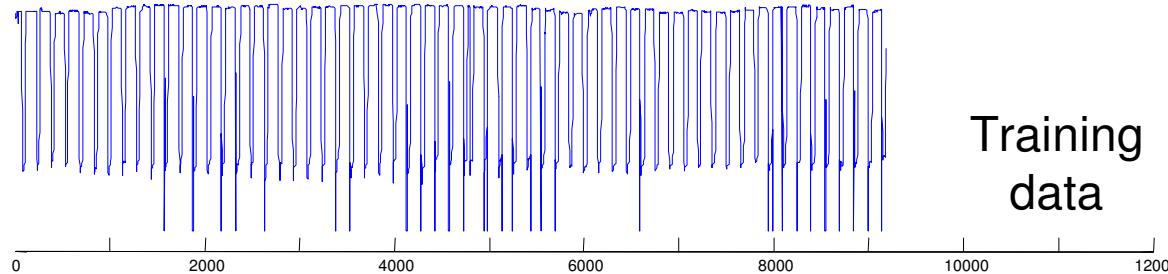
## Discrepancy Checking



WE BUILD A MODEL  
CAPABLE OF PREDICTING  
THE FUTURE VALUES OF  
A SERIES -  
THEN, IF THE ACTUAL  
RECORDED VALUES FALL  
OUTSIDE OF THE  
PREDICTION  $\pm$  A TOLERANCE  
THRESHOLD, WE DO SOMETHING

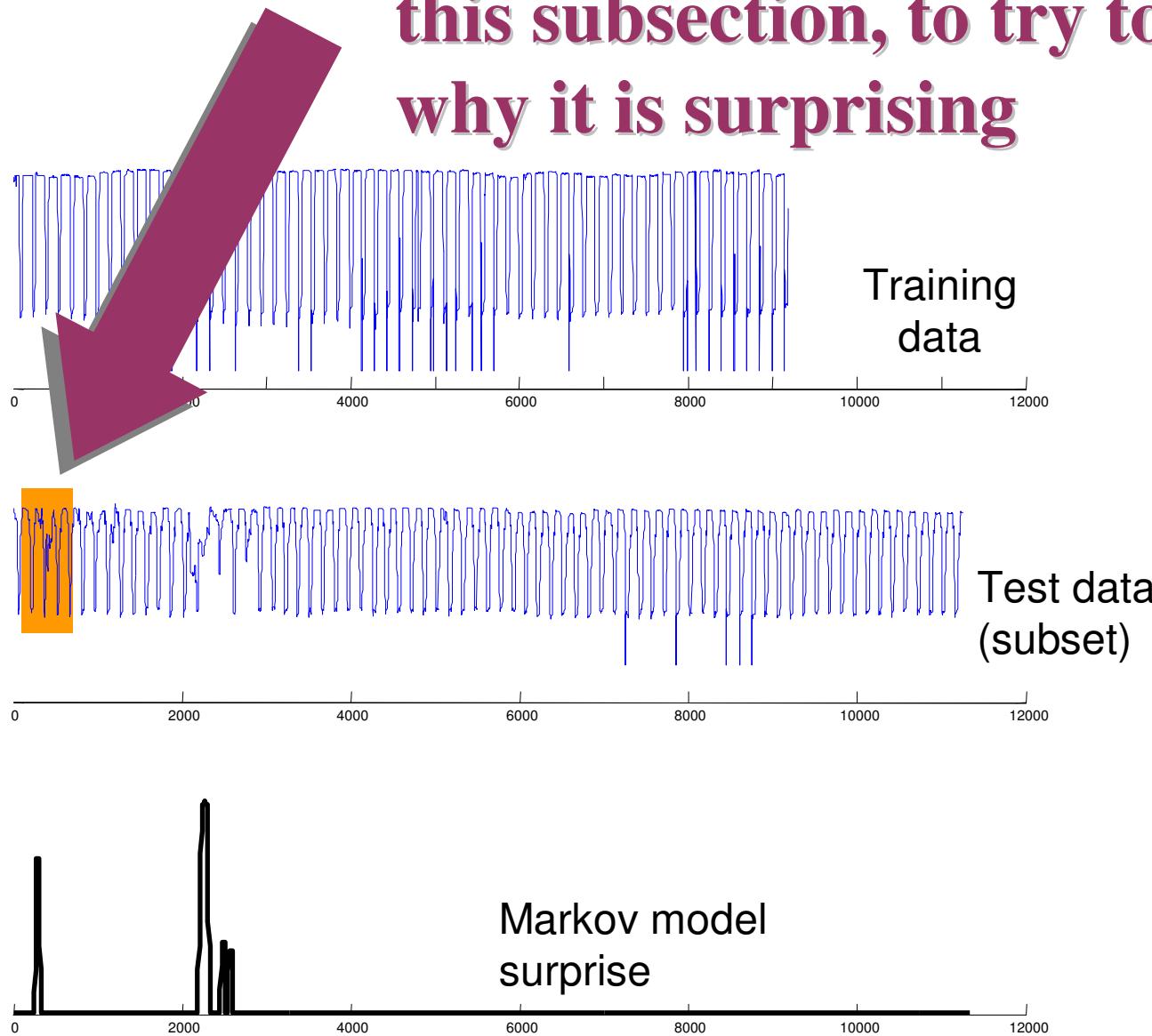
IN THIS WAY, WE CAN DETECT SOME SURPRISING,  
UNEXPECTED VALUES (E.G., ANTHRAX OUTBREAKS IN THE U.S.)

- Note that this problem has been solved for text strings
- You take a set of text which has been labeled “normal”, you learn a Markov model for it.
- Then, any future data that is not modeled well by the Markov model you annotate as *surprising*.
- Since we have just seen that we can convert time series to text (i.e SAX). Lets us quickly see if we can use Markov models to find surprises in time series...



These were converted to the symbolic representation.  
I am showing the original data for simplicity

In the next slide we will zoom in on this subsection, to try to understand why it is surprising



# Anomaly (interestingness) detection

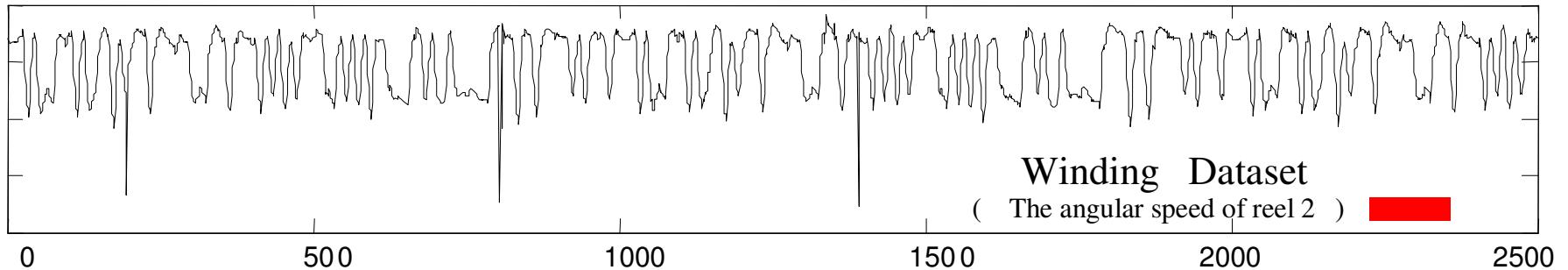
In spite of the nice example in the previous slide, the anomaly detection problem is wide open.

How can we find interesting patterns...

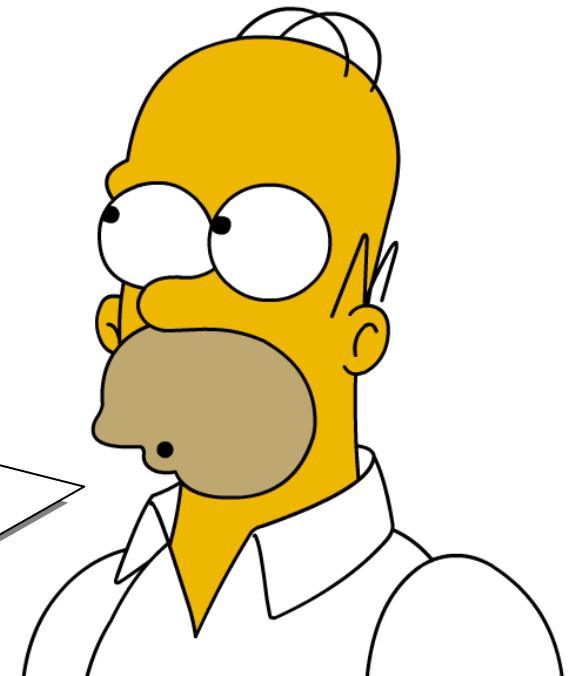
- Without (or with very few) false positives...
- In truly massive datasets...
- In the face of concept drift...
- With human input/feedback...
- With annotated data...

# Time Series Motif Discovery

(finding repeated patterns)

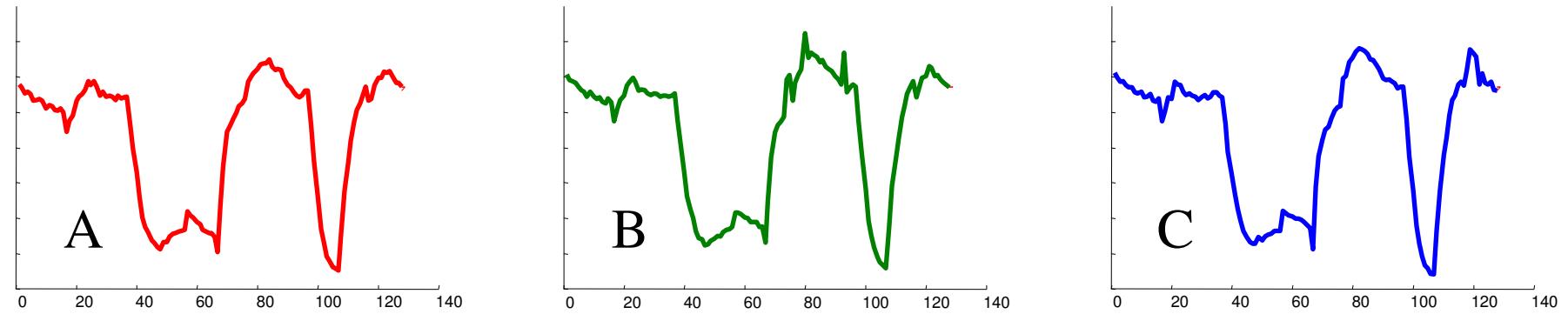
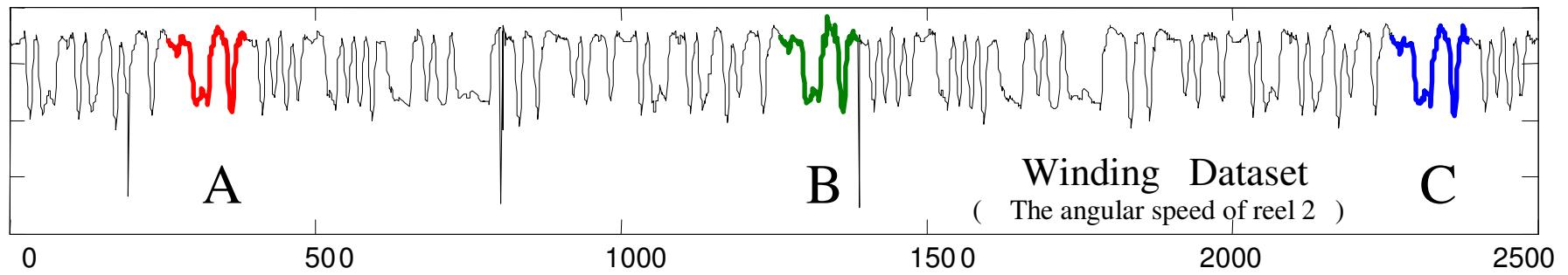


Are there any repeated patterns, of about this length — in the above time series?

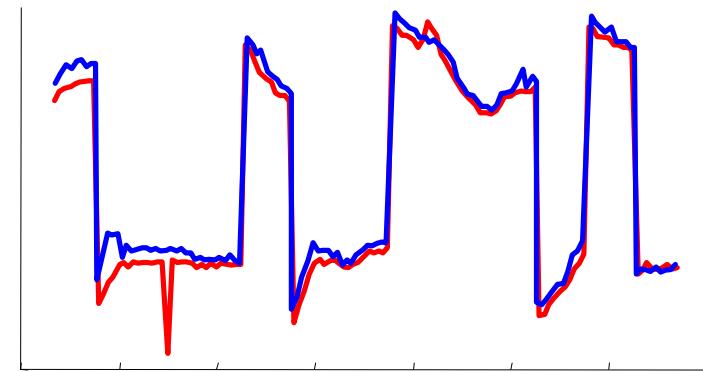
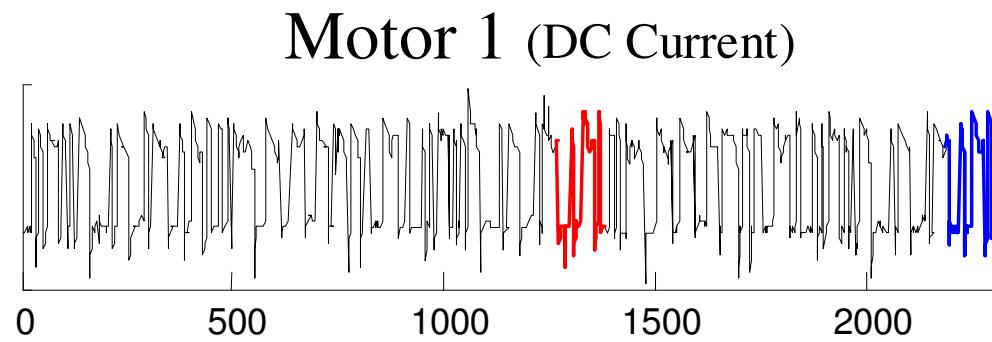


# Time Series Motif Discovery

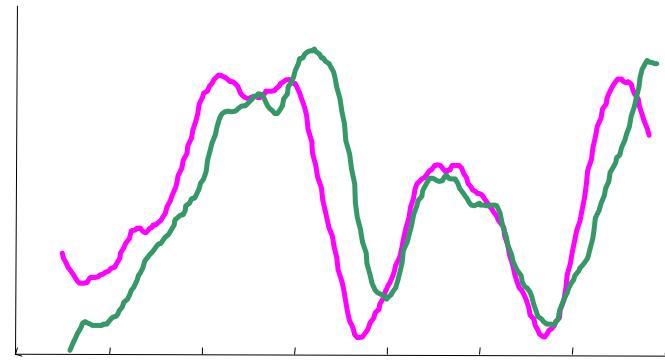
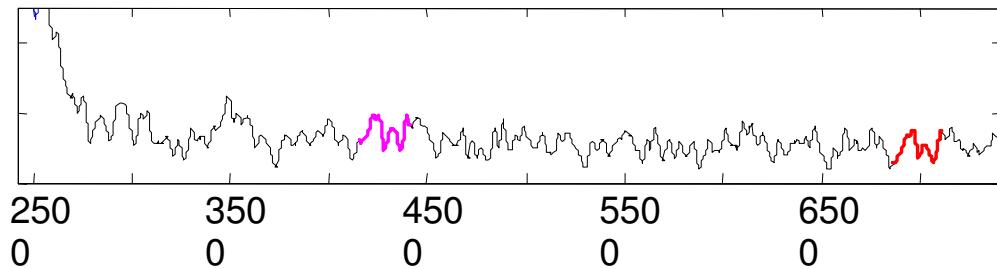
(finding repeated patterns)



# Some Examples of Real Motifs



Astrophysics (Photon Count)



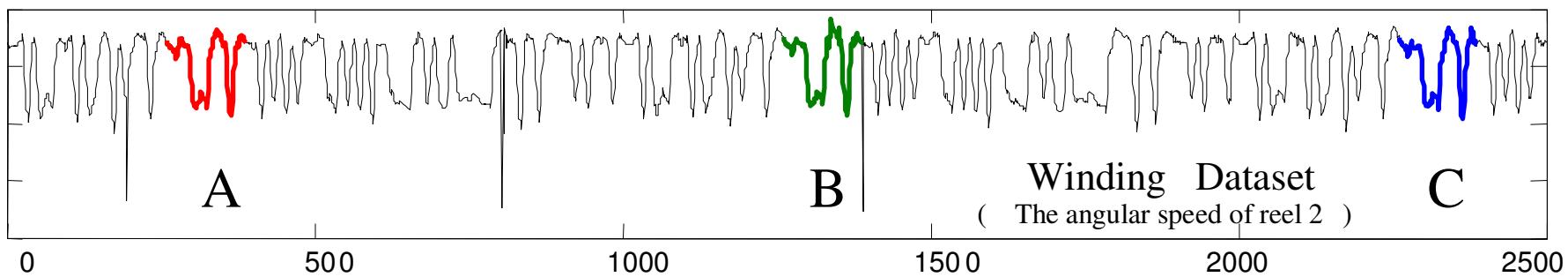
# Why Find Motifs?

- Mining **association rules** in time series requires the discovery of motifs. These are referred to as *primitive shapes* and *frequent patterns*.
- Several time series **classification algorithms** work by constructing typical prototypes of each class. These prototypes may be considered motifs.
- Many time series **anomaly/interestingness detection** algorithms essentially consist of modeling normal behavior with a set of typical shapes (which we see as motifs), and detecting future patterns that are dissimilar to all typical shapes.
- In **robotics**, Oates et al., have introduced a method to allow an autonomous agent to generalize from a set of qualitatively different *experiences* gleaned from sensors. We see these “*experiences*” as motifs.
- In **medical data mining**, Caraca-Valente and Lopez-Chavarrias have introduced a method for characterizing a physiotherapy patient’s recovery based on the discovery of *similar patterns*. Once again, we see these “*similar patterns*” as motifs.
- **Animation and video capture...** (Tanaka and Uehara, Zordan and Celly)

# Motifs Discovery Challenges

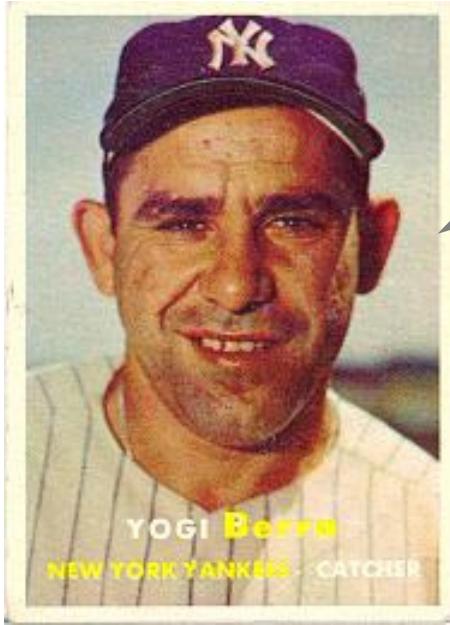
How can we find motifs...

- Without having to specify the length/other parameters
- In massive datasets
- While ignoring “background” motifs (ECG example)
- Under time warping, or uniform scaling
- While assessing their significance



Finding these 3 motifs requires about 6,250,000 calls to the Euclidean distance function

# Time Series Prediction



Yogi Berra  
1925 -

Prediction is hard, especially about the future

There are two kinds of time series prediction

- **Black Box:** Predict tomorrow's electricity demand, given *only* the last ten years electricity demand.
- **White Box (side information):** Predict tomorrow's electricity demand, given the last ten years electricity demand *and* the weather report, *and* the fact that the world cup final is on and...