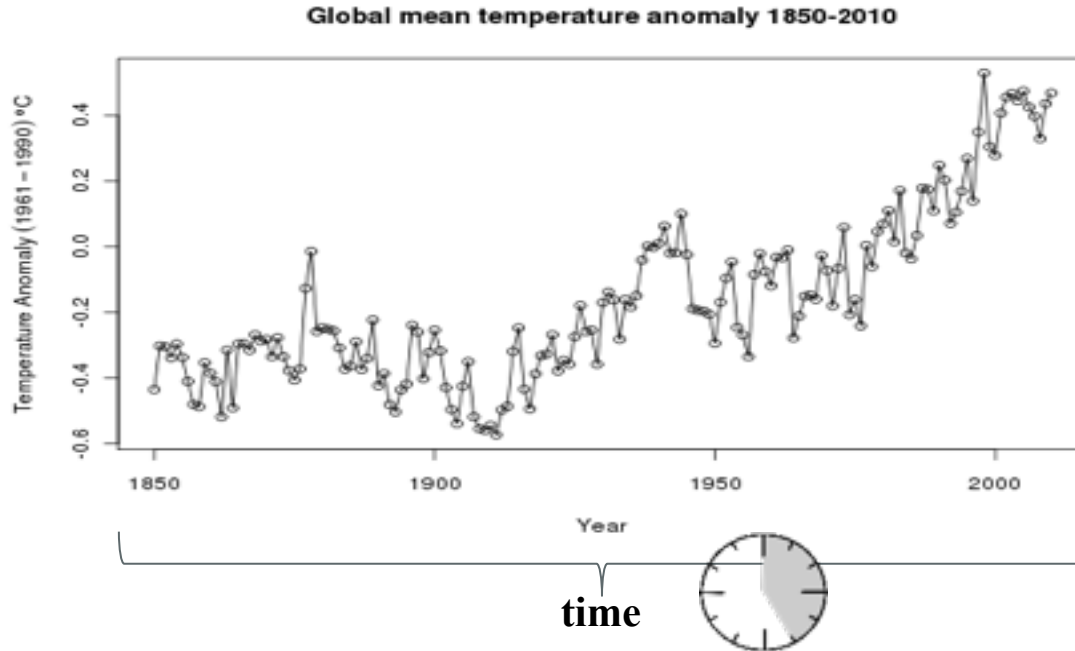


Interactive Visual Exploration of Large Data Series

Anna Gogolou

Time series

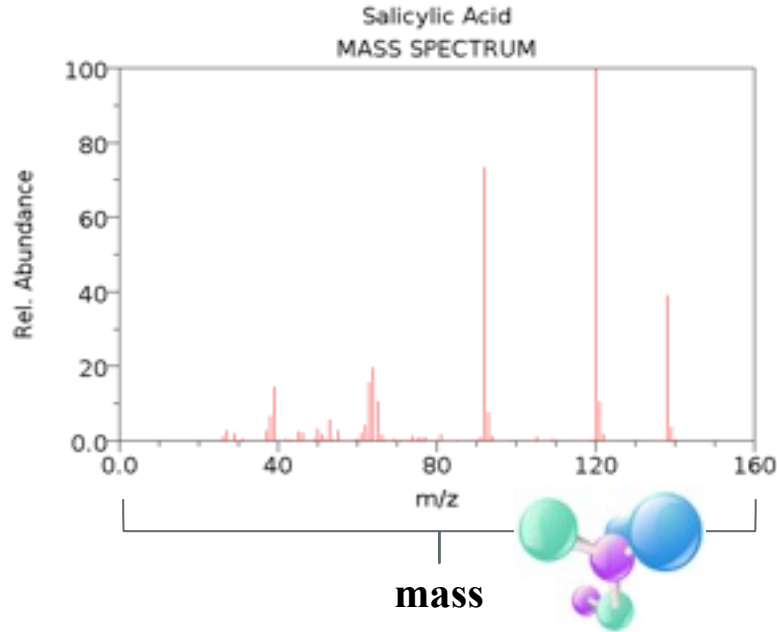
What are **time** series?



Data sequences
ordered along the
dimension of **time**

Data series

Observations, measurements ordered along a dimension (time, angle, mass, position, etc.)

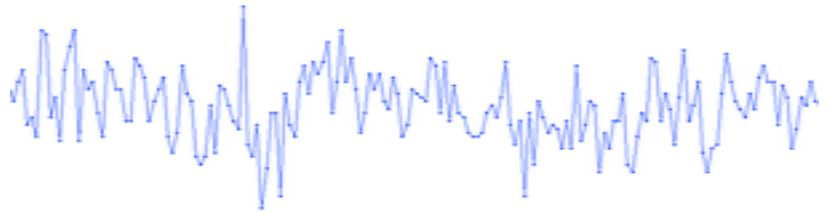


Why do we care about data series?

They occur in many scientific, medical, business, and social domains:

astronomy, biology, neuroscience, smart cities, nuclear power plants, finance, politics, ...

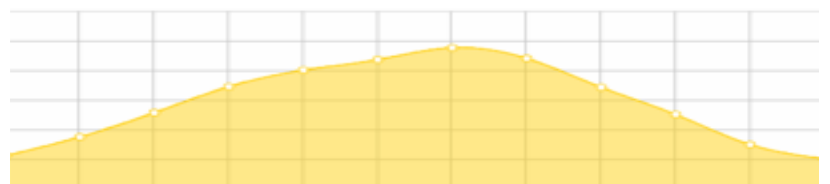
→ Our systolic/diastolic blood pressure



→ Donald Trump's popularity rating



→ The annual sunshine in Paris



Outline

Motivation

Our goal

Existing data series visualization tools

Scalability problems

State-of-the-art in data series management

Approximate vs exact search

Progressive visualizations

Motivation

How do we manage, explore and analyze large data series (order of terabytes)?

- More and more data series are produced every day:
 - 1 hour of ElectroCardioGraphy (ECG) data: 1 Gigabyte
 - Typical Weblog: 5 Gigabytes per week
 - EDF Database: ~100TB in total for all reactors (58), all sensors (10000 sensors per reactor), over all years and growing ...

Data series analysis

Analysis Tasks:

Pattern Matching

Frequent Pattern Mining

Clustering

Classification

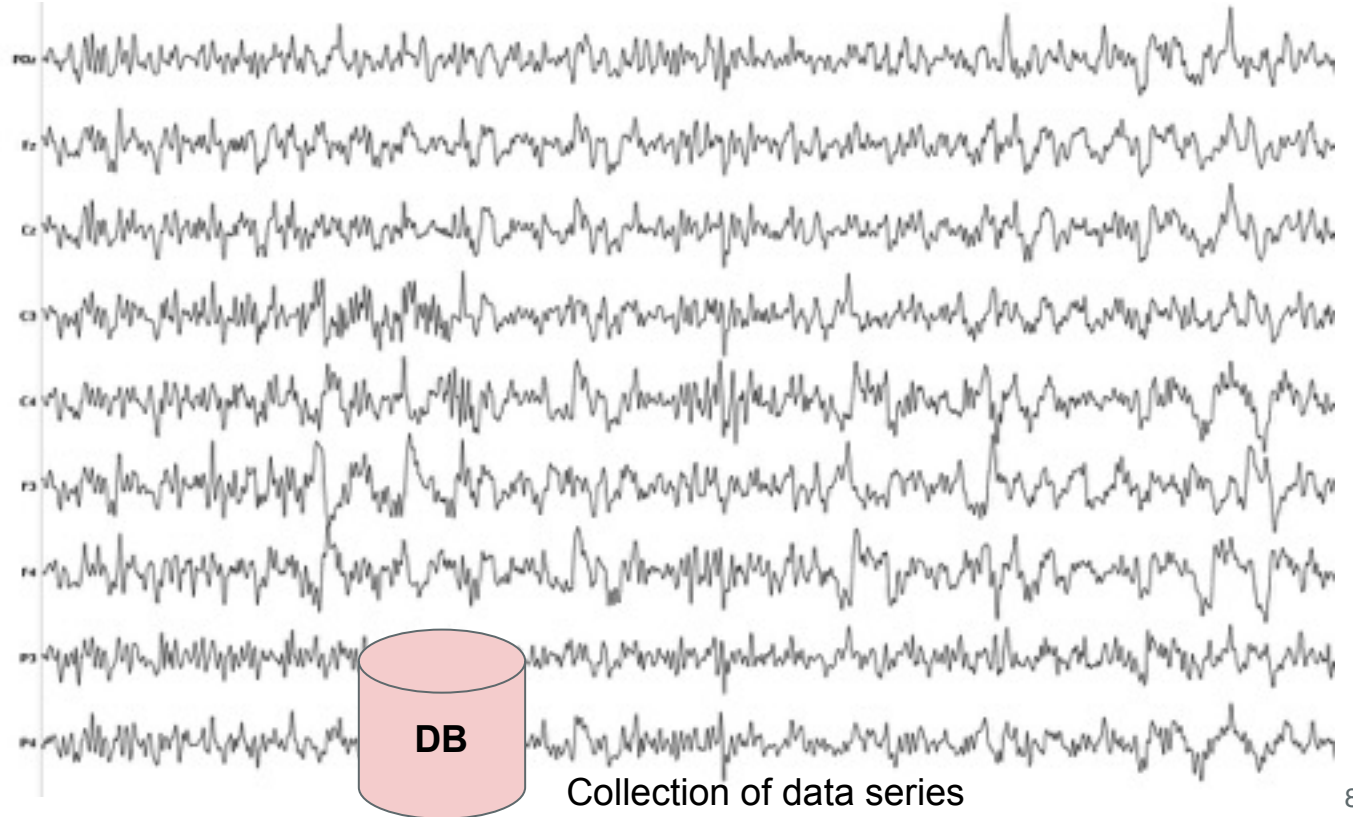
Outlier Detection

...

Pattern matching



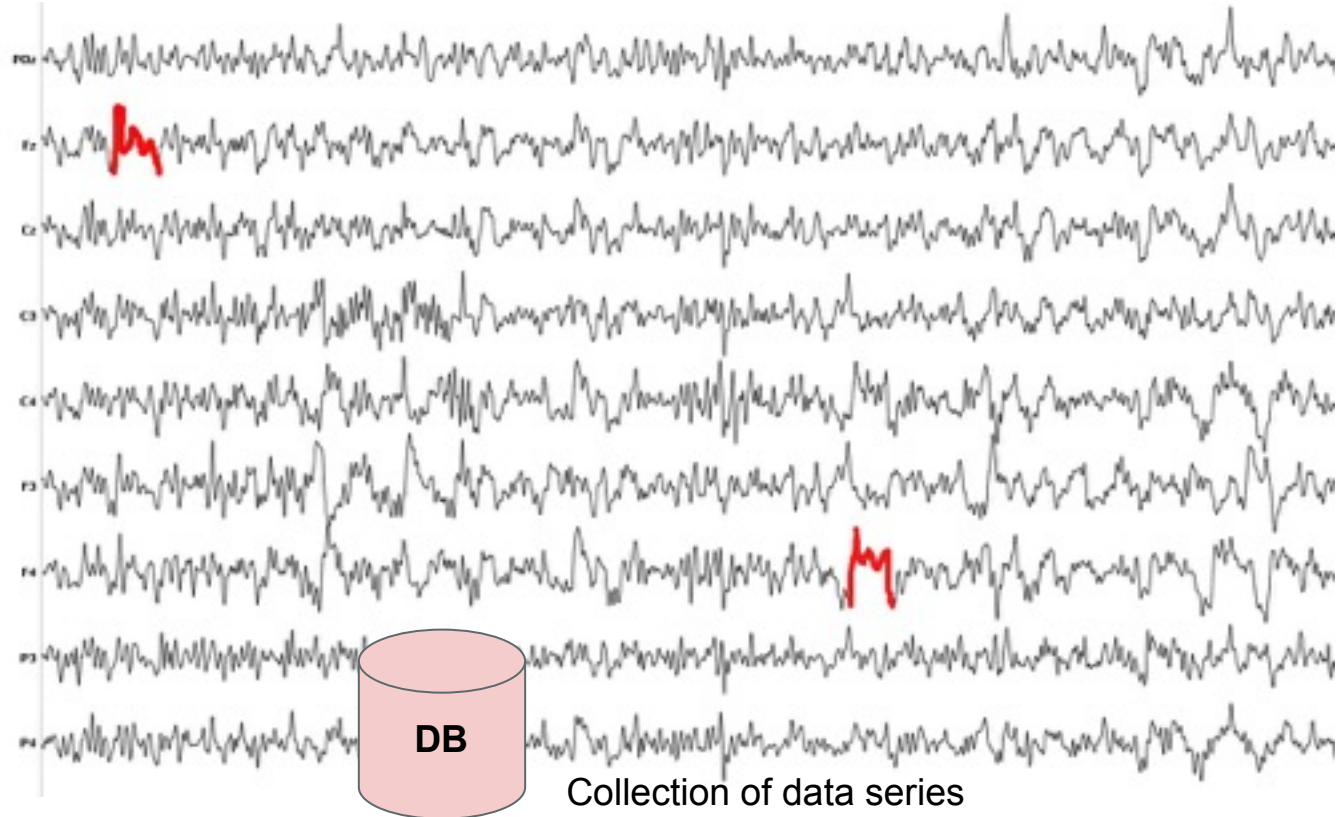
User-defined
pattern



Pattern matching



User-defined
pattern



Outline

Motivation

Our goal

Existing data series visualization tools

Scalability problems

State-of-the-art in data series management

Approximate vs exact search

Progressive visualizations

Our goal

Interactive visual exploration of large data series

Interdisciplinary topic in the areas of:

Human-Computer Interaction (HCI) & Information Visualization

Understand what analysts do with their data

Support interactive exploration of large datasets

Databases

Analyze efficiently large data series

Get answers quickly

Outline

Motivation

Our goal

Existing data series visualization tools

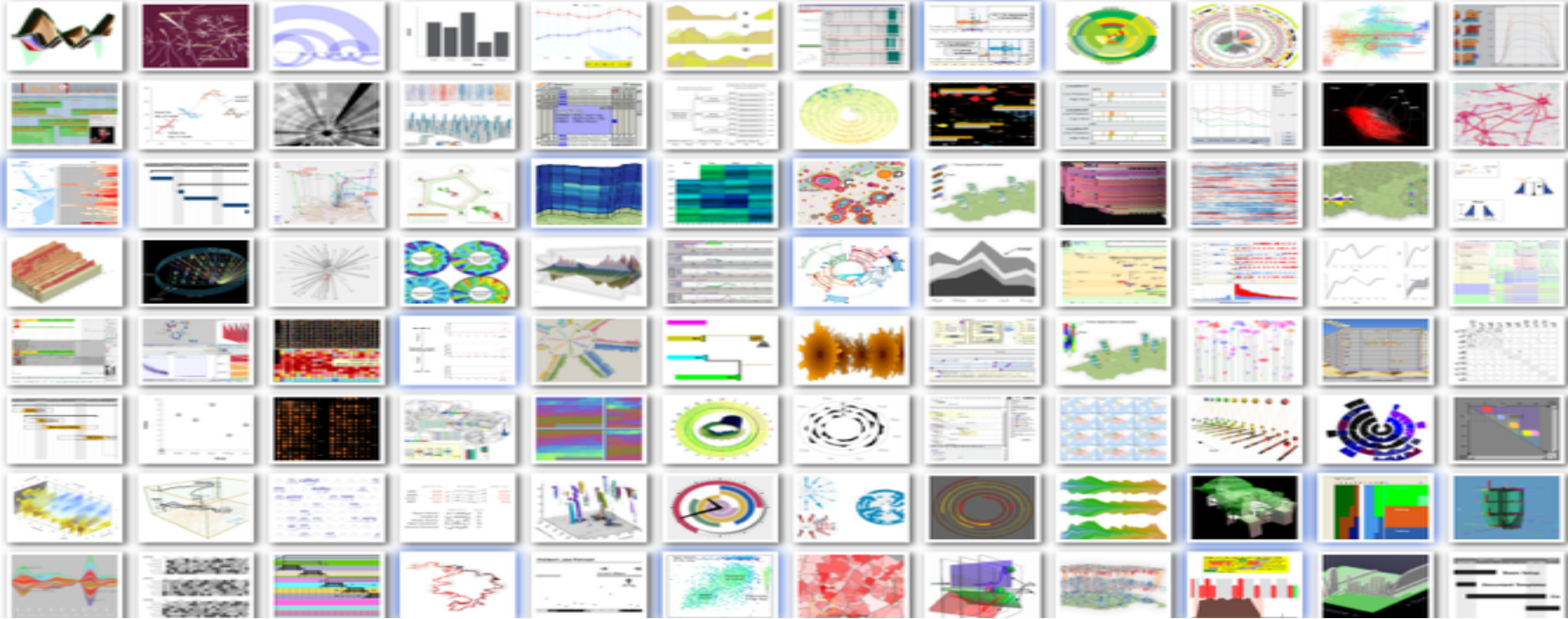
Scalability problems

State-of-the-art in data series management

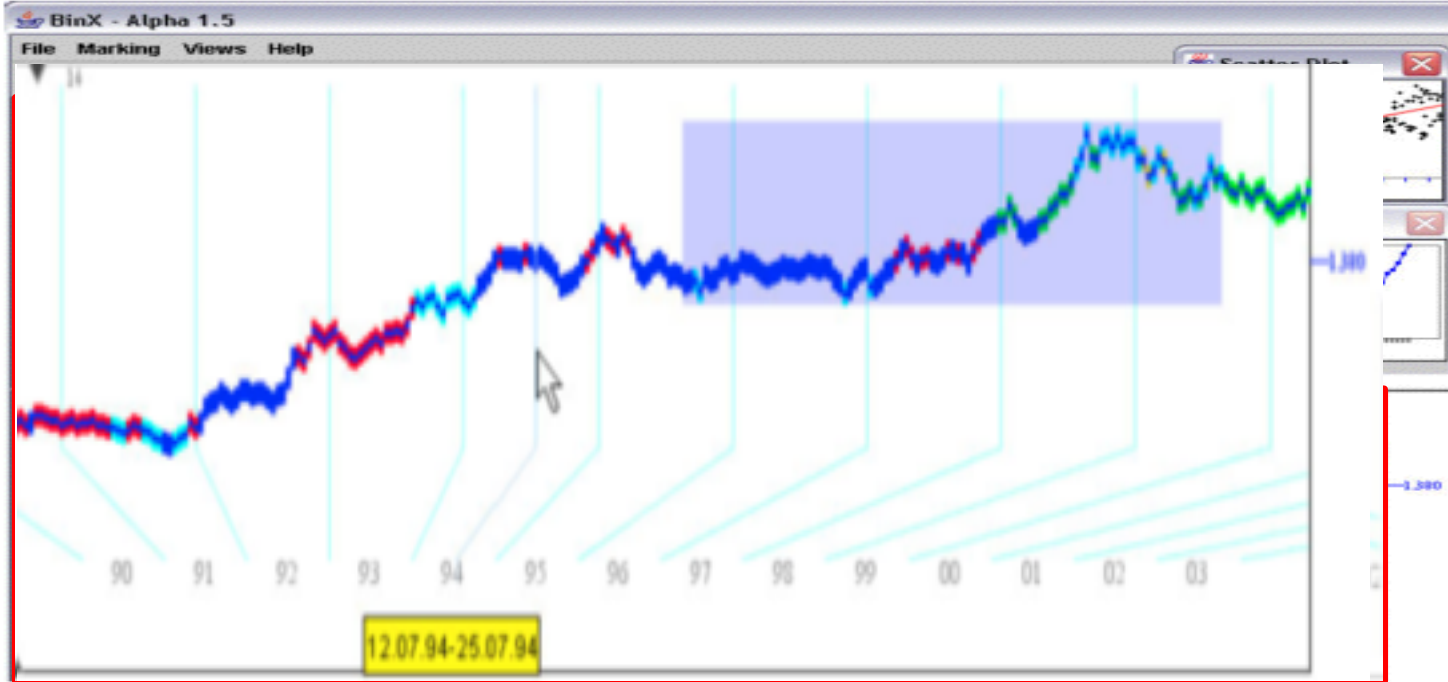
Approximate vs exact search

Progressive visualizations

Existing data series visualization tools

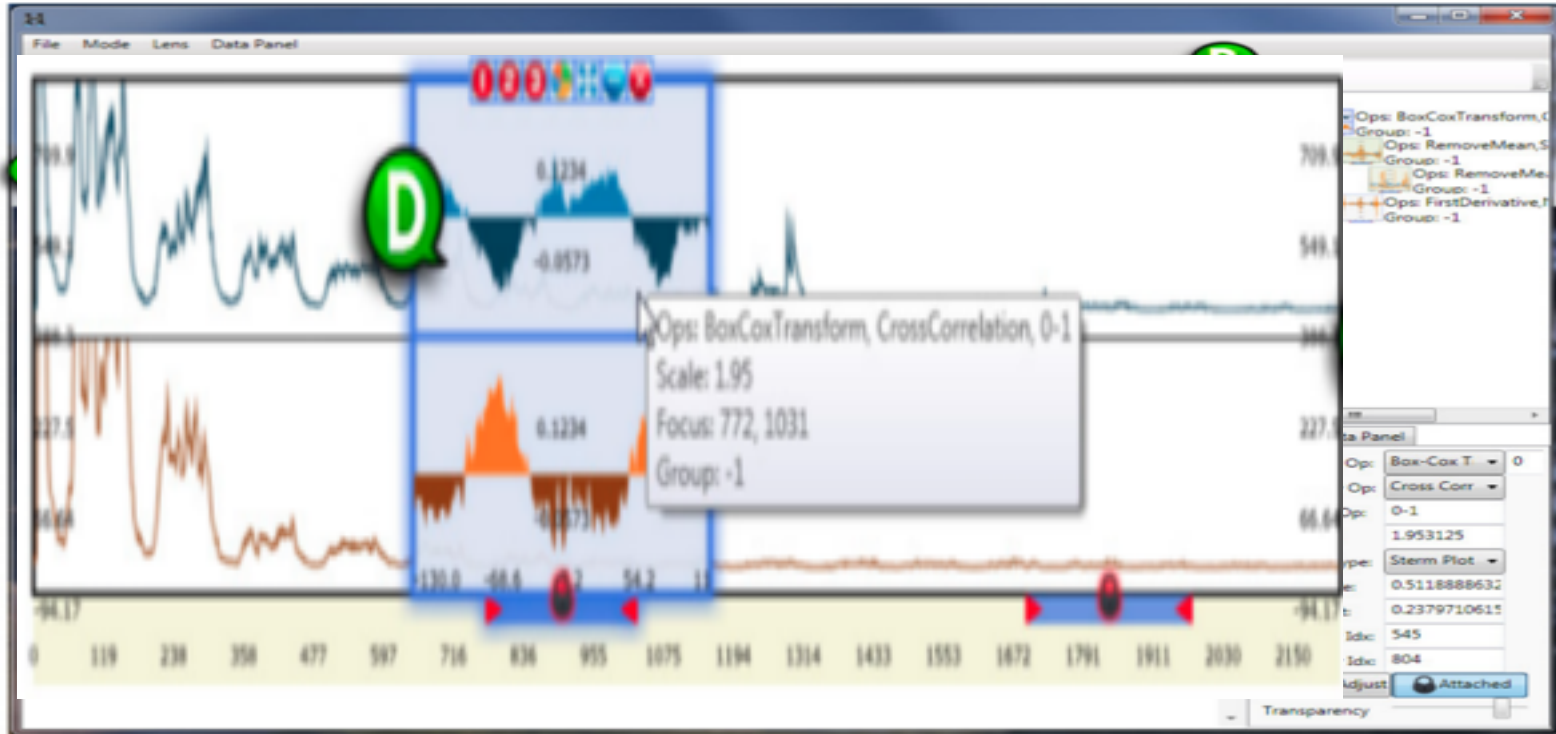


From <http://survey.timeviz.net/>



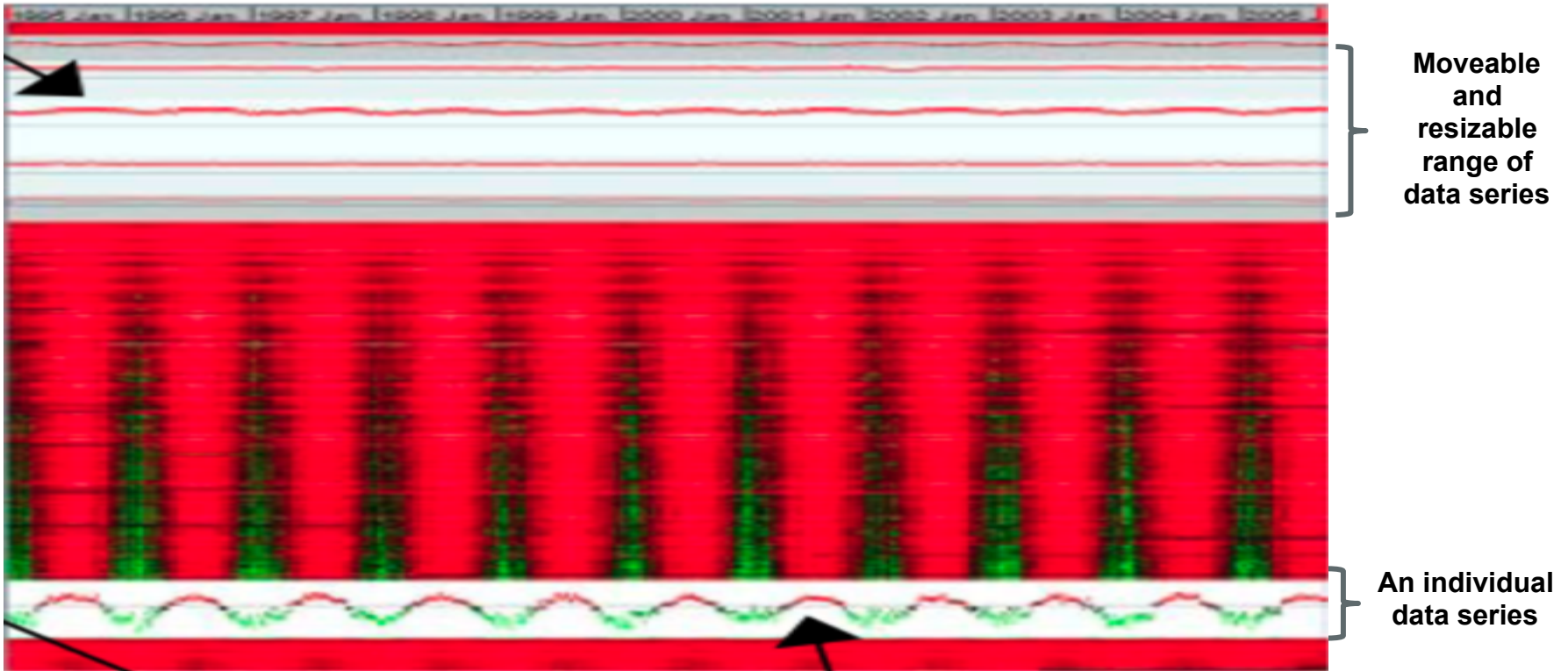
ChronoLenses

Data series
chart panel



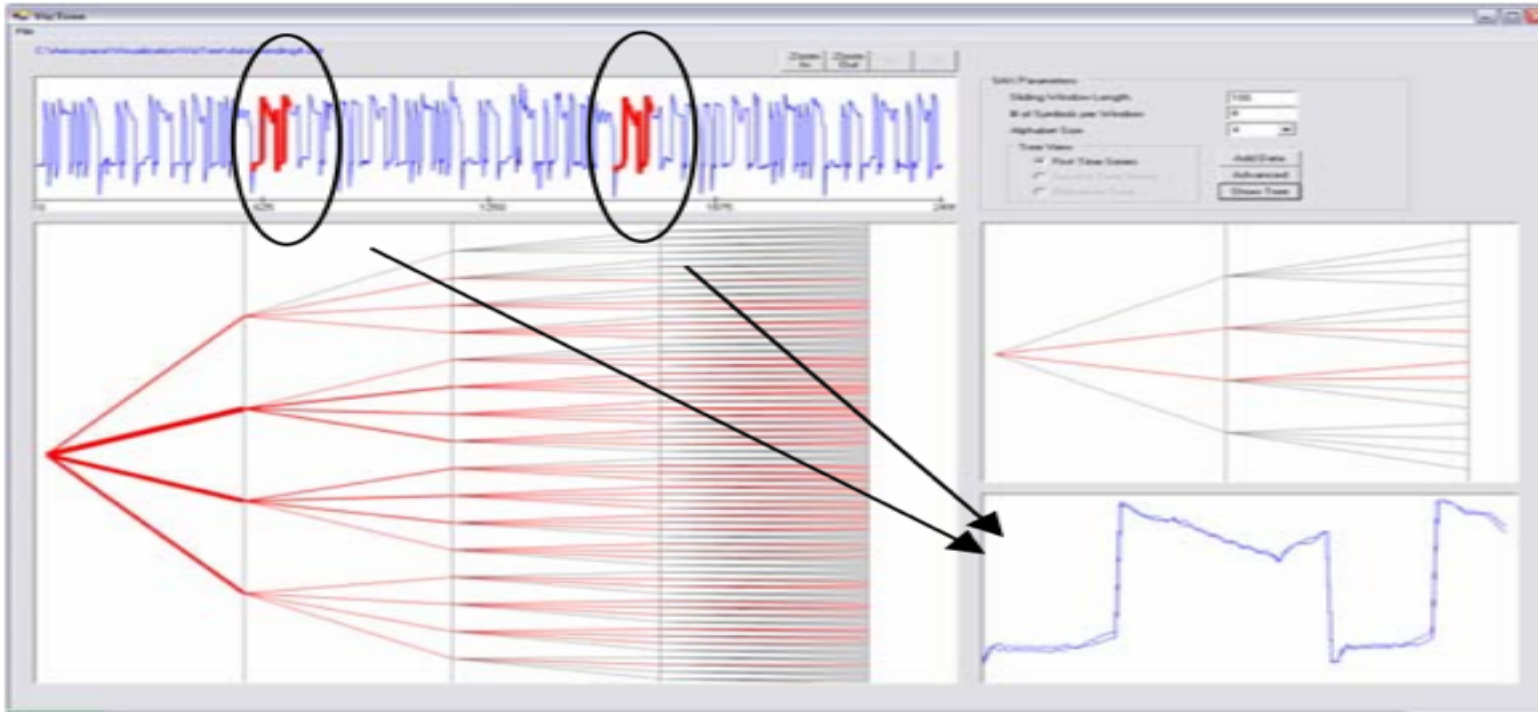
Line graph explorer

Line graph
panel

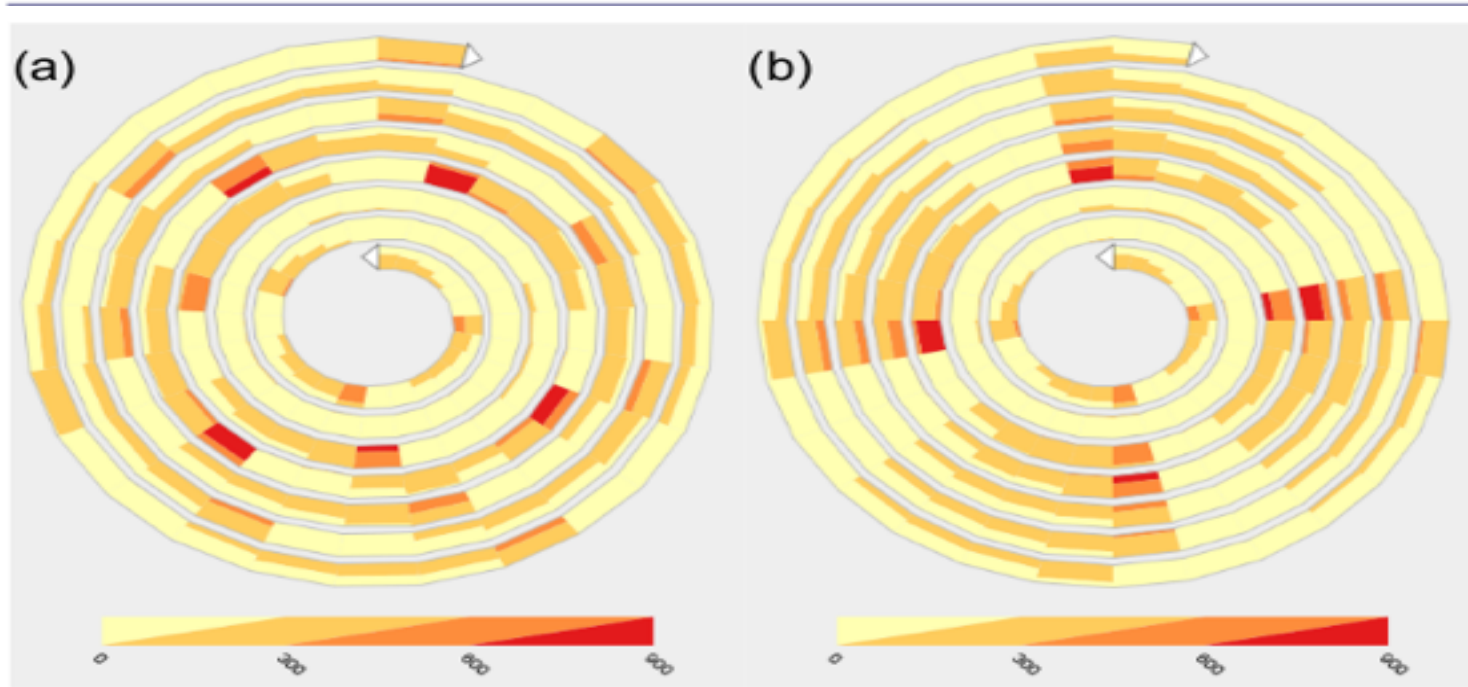


Robert Kincaid, Heidi Lam, Line graph explorer: scalable display of line graphs using Focus+Context, *Proceedings of the working conference on Advanced visual interfaces*, May 23-26, 2006, Venezia, Italy [doi>[10.1145/1133265.1133348](https://doi.org/10.1145/1133265.1133348)]

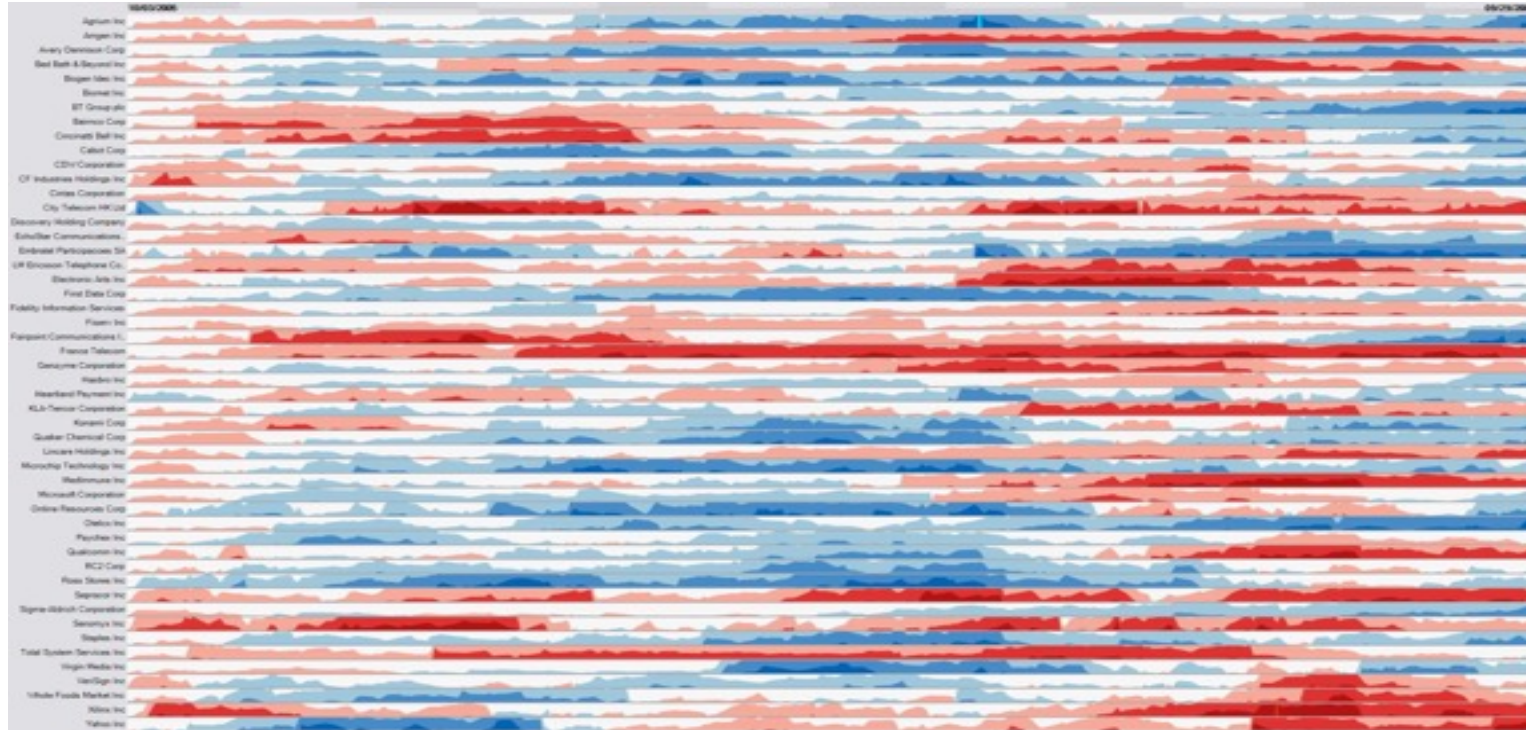
Viztree



Spirals



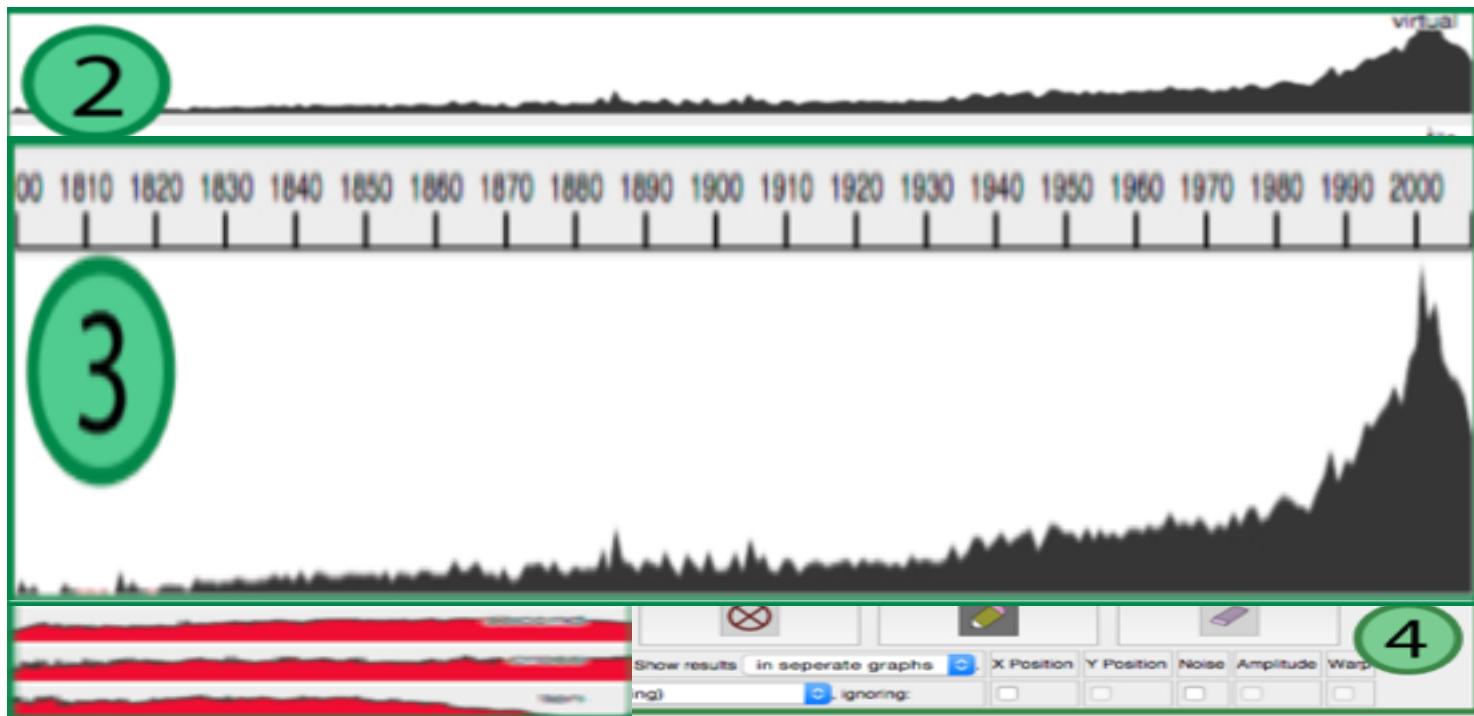
Horizon charts



Jeffrey Heer, Nicholas Kong, Maneesh Agrawala, Sizing the horizon: the effects of chart size and layering on the graphical perception of time series visualizations, *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, April 04-09, 2009, Boston, MA, USA [doi>[10.1145/1518701.1518897](https://doi.org/10.1145/1518701.1518897)]

The semantics of sketch

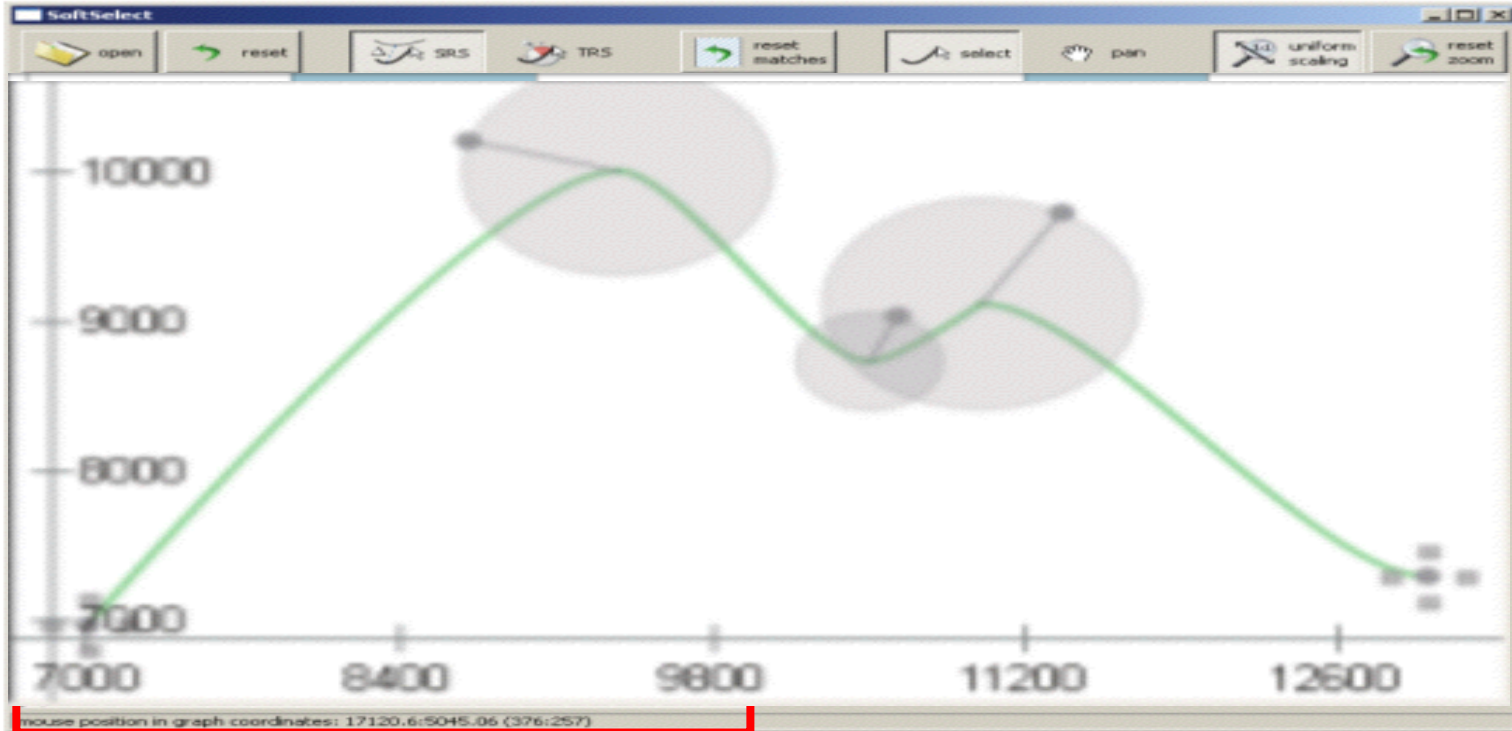
The matching panel



Relaxed queries

The graph-
display and
main-interaction
area

SoftSelect



Christian Holz, Steven Feiner, Relaxed selection techniques for querying time-series graphs, *Proceedings of the 22nd annual ACM symposium on User interface software and technology*, October 04-07, 2009, Victoria, BC, Canada [doi>[10.1145/1622176.1622217](https://doi.org/10.1145/1622176.1622217)]

Outline

Motivation

Our goal

Existing data series visualization tools

Scalability problems

State-of-the-art in data series management

Approximate vs exact search

Progressive visualizations

What about handling billions of data series?

None of these tools support scalability to terabytes of data:

How do we visualize them?

Limited number of pixels

Visual scalability

Limited human cognitive resources

Interactive response time scalability

How do we interact with them?

Days to answer a single query

Limit: < 100 ms

J. Nielsen, Response times: The 3 important limits, <https://www.nngroup.com/articles/response-times-3-important-limits/>,

January 1, 1993

Outline

Motivation

Our goal

Existing data series visualization tools

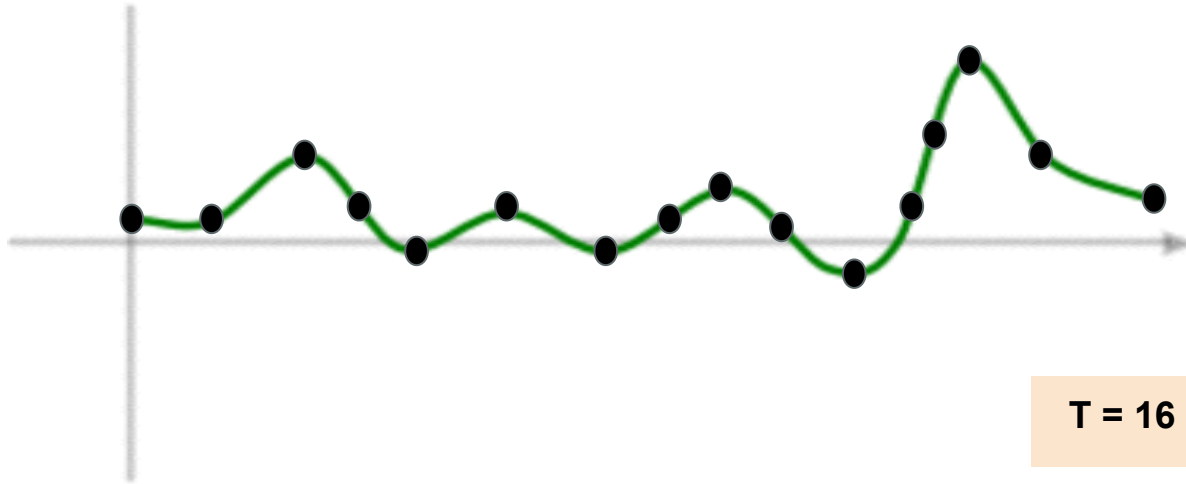
Scalability problems

State-of-the-art in data series management

Approximate vs exact search

Progressive visualizations

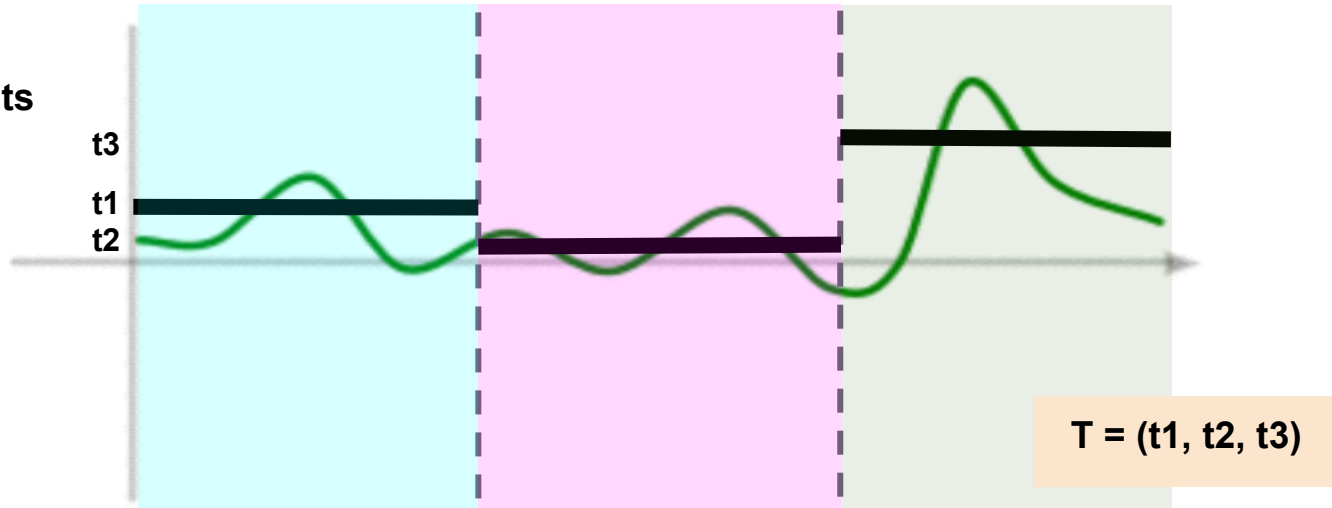
Data series summarization



A raw data series T

Data series summarization

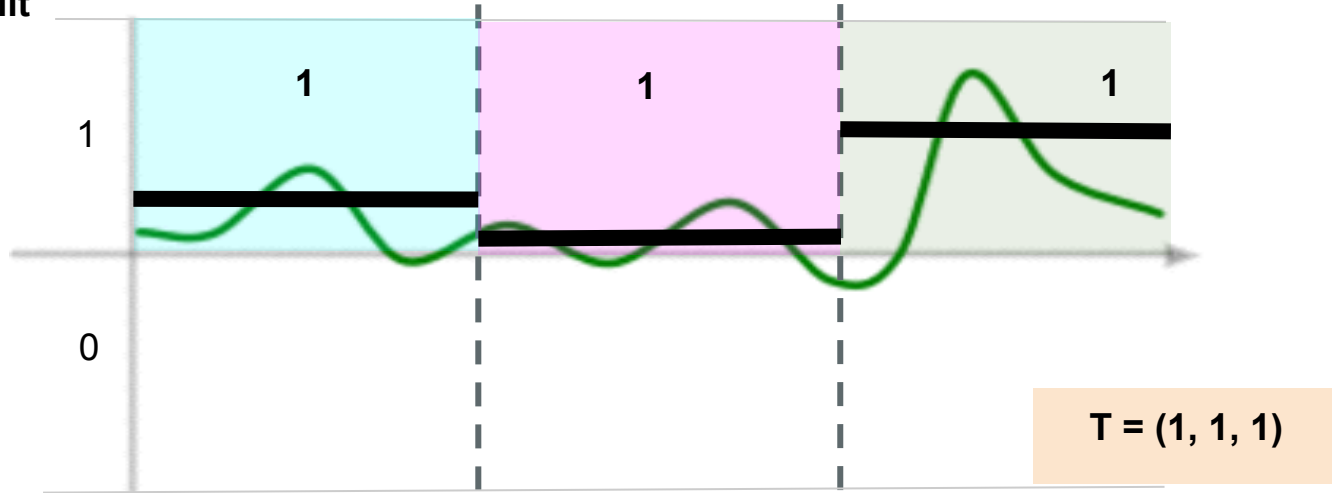
Let's split
X-axis into
3 equal segments



Piecewise Aggregate Approximation (PAA)

Data series summarization

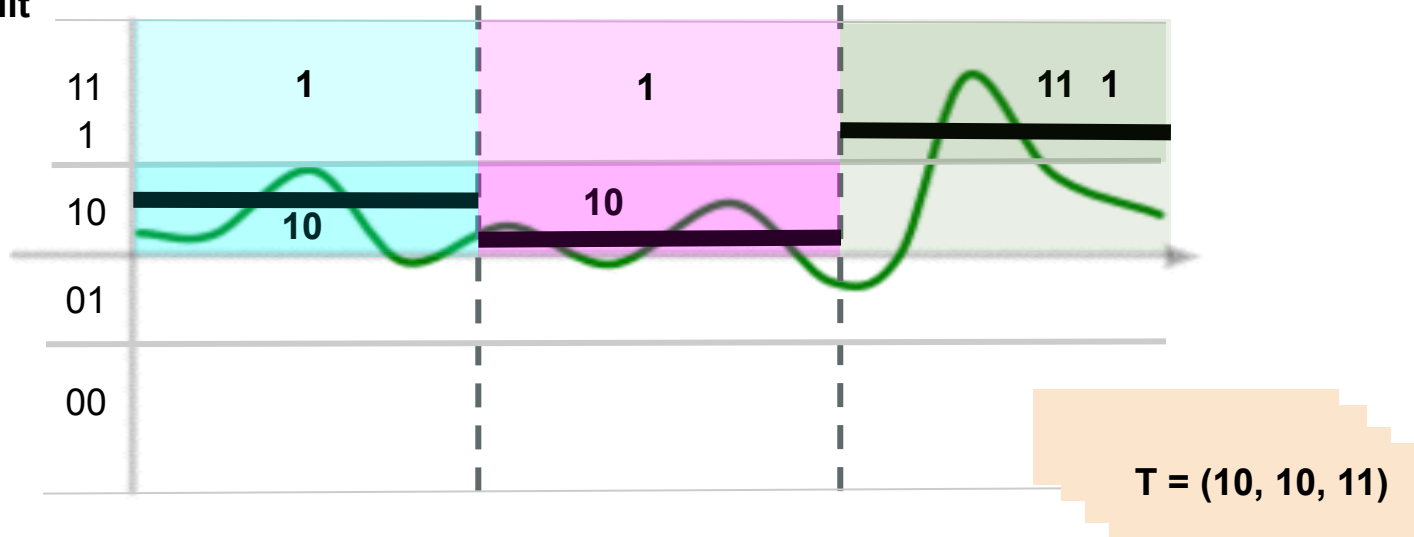
Let's further split
Y-axis



Symbolic Aggregate approxImation (SAX)

Data series summarization

Let's further split
Y-axis



Symbolic Aggregate approximation (SAX)

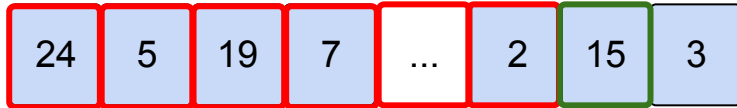
Indexing

How can queries be answered faster?

Using index structures (usually trees), which improve the speed of database operations on database tables (insert, select, delete, ...).

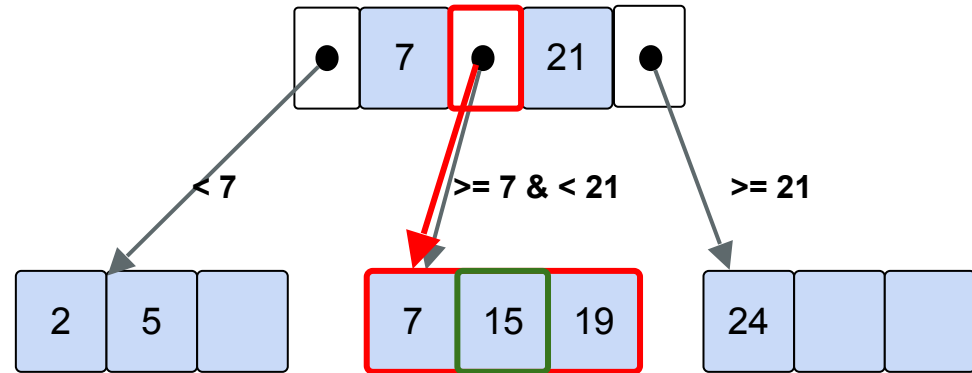
How do indexes work?

No structure



Serial scan

Index structure



Index search

Data series indexes

Data series indexes follow a similar tree structure.

They are built on top of data series summarizations.

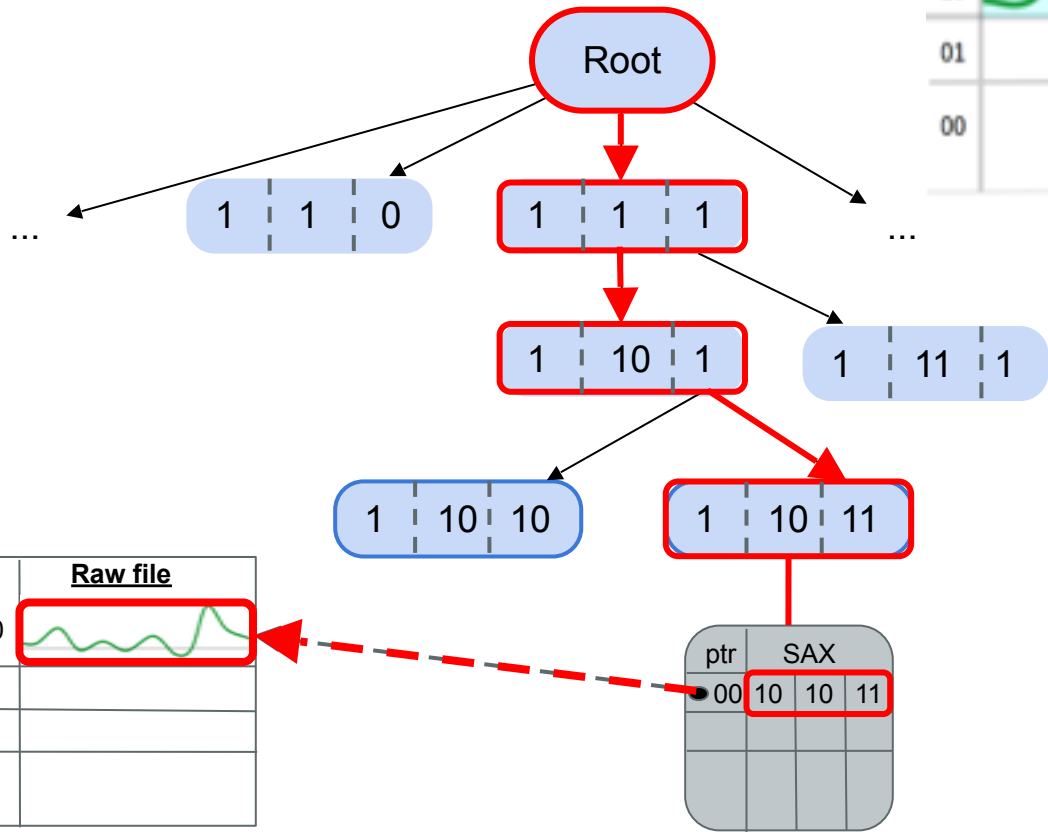
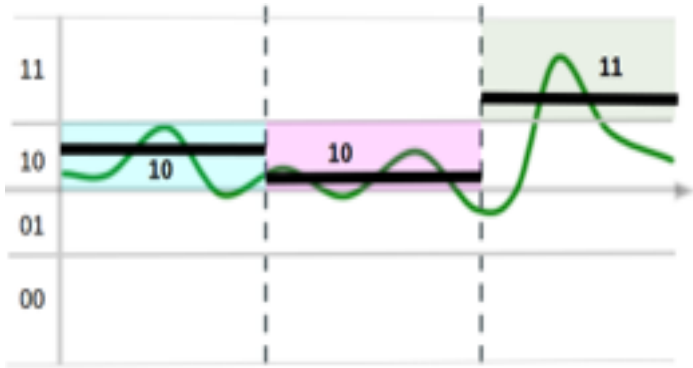
ADS+¹ is the state-of-the-art data series index. It is based on:

- SAX representation

- Euclidean Distance (ED) between SAX representations

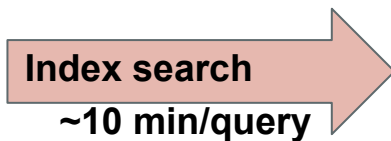
¹Zoumpatianos, K., Idreos, S., Palpanas, T., ADS: the adaptive data series index, *The VLDB Journal* (2016) 25: 843. [doi>[10.1007/s00778-016-0442-5](https://doi.org/10.1007/s00778-016-0442-5)]

ADS+ index



Response times

Finding the most similar pattern on 1 **billion** data series:



**Still not interactive response times
(not in the order of milliseconds)**

Outline

Motivation

Our goal

Existing data series visualization tools

Scalability problems

State-of-the-art in data series management

Approximate vs exact search

Progressive visualizations

Approximate vs exact search

Approximate search:

- Returns answer not guaranteed to be the best one
- Associated with an approximation error

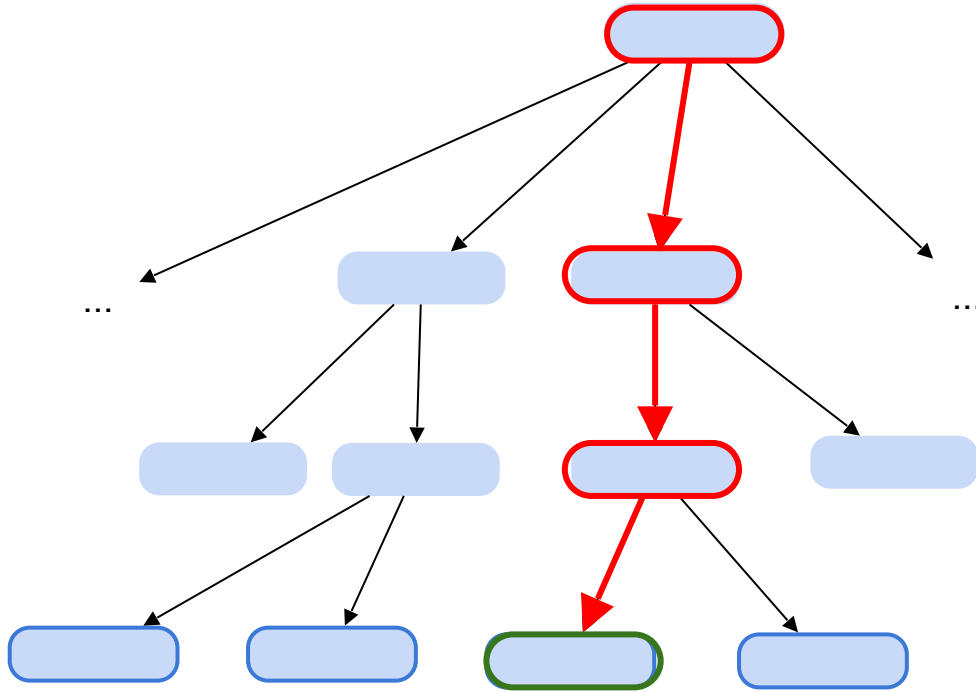
- Very fast (order of milliseconds)

Exact search:

- Always returns exact correct answer

- Slow (order of minutes)

Approximate search



Outline

Motivation

Our goal

Existing data series visualization tools

Scalability problems

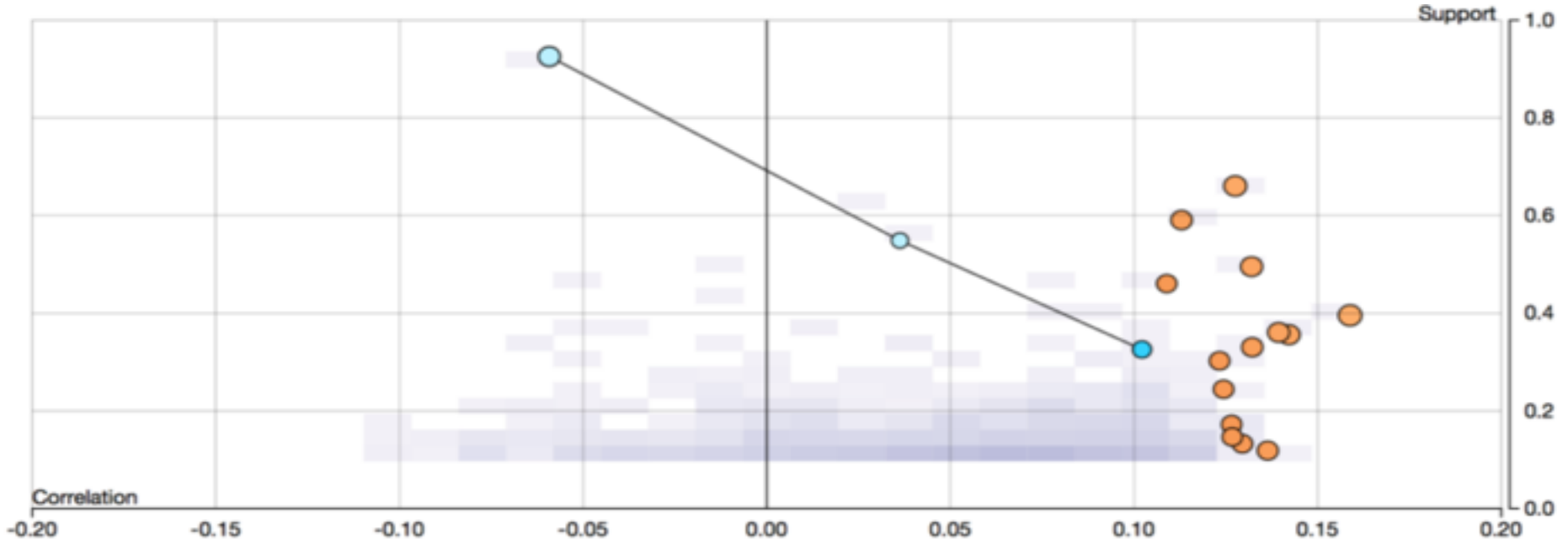
State-of-the-art in data series management

Approximate vs exact search

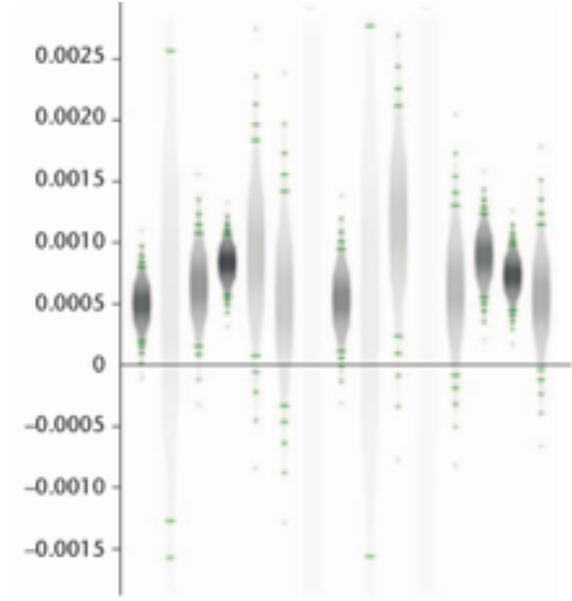
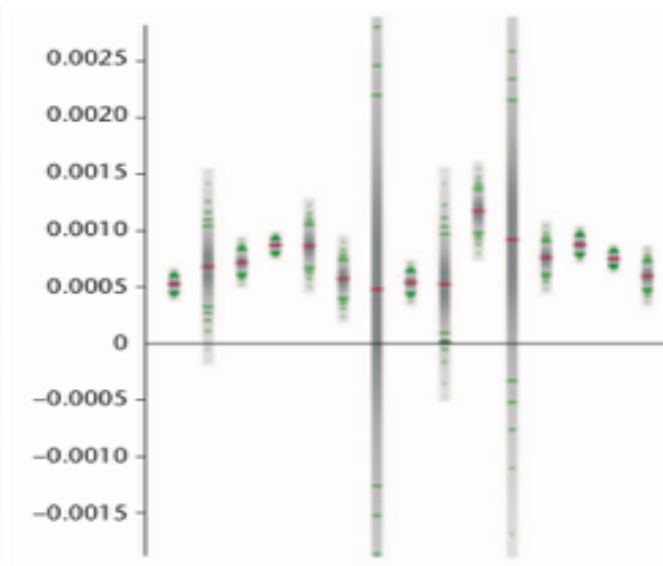
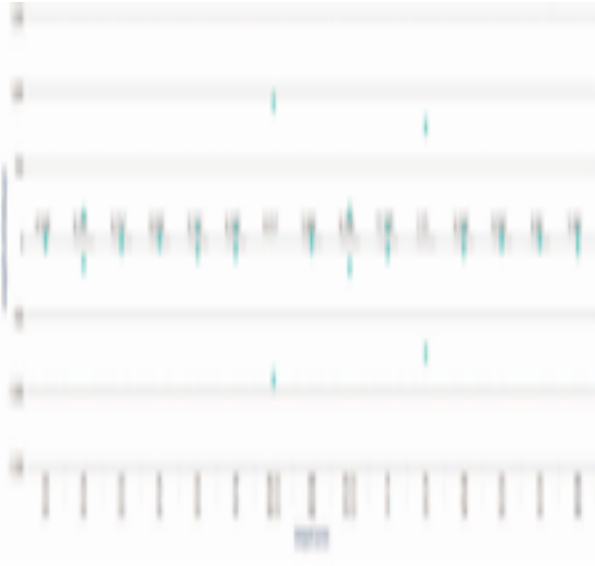
Progressive visualizations

Progressive visualizations

How can users progressively move from approximate to exact results?



Visualizing approximation error



Challenges

Provide quickly partial results

Visualize approximation error

Support iteration and refinement over approximate & progressive results

Easily focus on subspaces of interest

Conclusions

Open problems:

Scalability to terabytes of data

Current visualization tools cannot handle large data series collections

Need techniques that:

- effectively visualize large volumes of data

- have interactive response times ($< 100\text{ms}$)

Support for iteration and refinement of approximate & progressive results

Currently no such support for data series analytics

Need techniques that:

- visually inform users on progress of task

- convey information on accuracy of current results