# Scalable Machine Learning on Large Sequence Collections

Themis Palpanas

*University of Paris*
*French University Institute*

LIPADE
Laboratoire d'Informatique PAris DEscartes

University of Tokyo – Tokyo (Japan), January 2020

diNo

1

---

diNo  2

# References 1

- papers
  - **iSAX 2.0: Indexing and Mining One Billion Time Series**. ICDM 2010
    - http://www.mi.parisdescartes.fr/~themisp/publications/icdm10-billiontimeseries.pdf
  - **Beyond One Billion Time Series: Indexing and Mining Very Large Time Series Collections with iSAX2+**. KAIS 2014
    - http://www.mi.parisdescartes.fr/~themisp/publications/kais14-isax2plus.pdf
  - **Indexing for Interactive Exploration of Big Data Series**. SIGMOD 2014
    - http://www.mi.parisdescartes.fr/~themisp/publications/sigmod14-ads.pdf
  - **RINSE: Interactive Data Series Exploration**. PVLDB 2015
    - http://www.mi.parisdescartes.fr/~themisp/publications/vldb15-rinse.pdf
  - **Query Workloads for Data-Series Indexes**. KDD 2015
    - http://www.mi.parisdescartes.fr/~themisp/publications/kdd15-bends.pdf
  - **Big Sequence Management: A Glimpse on the Past, the Present, and the Future. LNCS, 2016**
    - http://www.mi.parisdescartes.fr/~themisp/publications/sofsem16-bisem.pdf
  - **ADS: The Adaptive Data Series Index**. PVLDBJ 2016
    - http://www.mi.parisdescartes.fr/~themisp/publications/vldbj16-ads.pdf
  - **DPiSAX: Massively Distributed Partitioned iSAX**. ICDM 2017
    - http://www.mi.parisdescartes.fr/~themisp/publications/icdm17-dpisax.pdf
  - **Generating Data Series Query Workloads**. VLDBJ 2018
    - http://www.mi.parisdescartes.fr/~themisp/publications/vldbj18-benchmark.pdf
  - **Coconut: A Scalable Bottom-Up Approach for Building Data Series Indexes.** PVLDB 2018
    - http://www.mi.parisdescartes.fr/~themisp/publications/vldb18-coconut.pdf
  - **Comparing Similarity Perception in Time Series Visualizations.** VIS 2018
    - http://www.mi.parisdescartes.fr/~themisp/publications/vis2018.pdf
  - **Data Series Management: Fulfilling the Need for Big Sequence Analytics**. ICDM 2018
    - http://www.mi.parisdescartes.fr/~themisp/publications/icde18-sms.pdf
  - **ULISSE: ULtra compact Index for Variable-Length Similarity SEarch in Data Series**. ICDE 2018
    - http://www.mi.parisdescartes.fr/~themisp/publications/icde18-ulisse.pdf

Themis Palpanas - University of Tokyo - Jan 2020

2

---

# References 2

- papers
  - **Massively Distributed Time Series Indexing and Querying.** TKDE 2018
    - http://helios.mi.parisdescartes.fr/~themisp/publications/tkde18-dpisax.pdf
  - **ParIS: The Next Destination for Fast Data Series Indexing and Query Answering.** IEEE BigData 2018
    - http://helios.mi.parisdescartes.fr/~themisp/publications/bigdata18.pdf
  - **Progressive Similarity Search on Time Series Data.** BigVis 2019
    - http://helios.mi.parisdescartes.fr/~themisp/publications/bigvis19.pdf
  - **Scalable, Variable-Length Similarity Search in Data Series: The ULISSE Approach.** PVLDB 2019
    - http://www.mi.parisdescartes.fr/~themisp/publications/vldb19-ulisse.pdf
  - **The Lernaean Hydra of Data Series Similarity Search: An Experimental Evaluation of the State of the Art.** PVLDB 2019
    - http://www.mi.parisdescartes.fr/~themisp/publications/vldb19-lernaeanhydra.pdf
  - **Coconut Palm: Static and Streaming Data Series Exploration Now in your Palm.** SIGMOD 2019
    - http://www.mi.parisdescartes.fr/~themisp/publications/vldb19-lernaeanhydra.pdf
  - **Report on the First and Second Interdisciplinary Time Series Analysis Workshop (ITISA).** Sigmod Record 2019
    - http://helios.mi.parisdescartes.fr/~themisp/publications/sigrec19-itisa.pdf
  - **Distributed Algorithms to Find Similar Time Series.** PKDD 2019
    - http://helios.mi.parisdescartes.fr/~themisp/publications/ecmlpkdd19.pdf
  - **Evolution of a Data Series Index.** ISIP 2019
    - http://helios.mi.parisdescartes.fr/~themisp/publications/isip19.pdf
  - **Return of the Lernaean Hydra: Experimental Evaluation of Data Series Approximate Similarity Search.** PVLDB 2020
    - http://www.mi.parisdescartes.fr/~themisp/publications/vldb20-lernaeanhydra2.pdf
  - **MESSI: In-Memory Data Series Indexing.** ICDE 2020
    - http://helios.mi.parisdescartes.fr/~themisp/publications/icde20-messi.pdf
  - **Coconut: Sortable Summarizations for Scalable Indexes over Static and Streaming Data Series.** VLDBJ 2020
    - http://helios.mi.parisdescartes.fr/~themisp/publications/vldbj19-coconut.pdf

3

---

# References 3

- code and datasets
  - **MESSI:** http://helios.mi.parisdescartes.fr/~themisp/messi/
  - **ParIS:** http://helios.mi.parisdescartes.fr/~themisp/paris/
  - **Coconut:** https://github.com/kon0925/coconut
  - **ULISSE:** http://helios.mi.parisdescartes.fr/~mlinardi/ULISSE.html
  - **DPiSAX:** https://github.com/DPiSAX/DPiSAX.github.io
  - **ADS:** https://github.com/zoumpatianos/ADS
  - **iSAX2+:** http://www.mi.parisdescartes.fr/~themisp/isax2plus/

- data series toolbox
  - **DSStat:** https://github.com/zoumpatianos/DSStat

- demos
  - **Coconut Palm:** http://users.ics.forth.gr/~kondylak/Coconut_Palm/Coconut_Palm.html
  - **DPiSAX:** http://imitates.gforge.inria.fr/
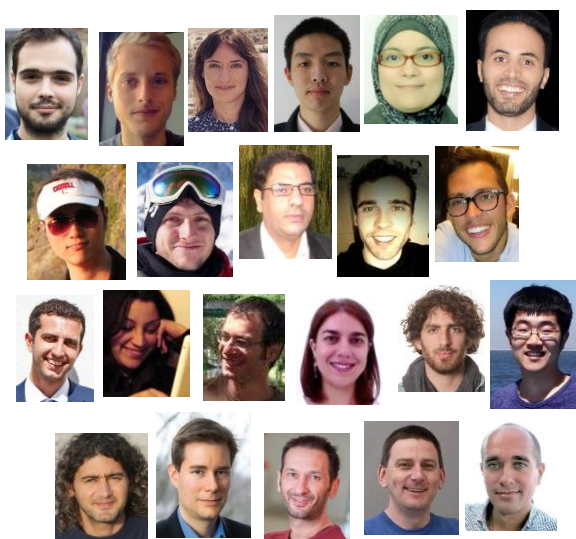  - **RINSE:** http://daslab.seas.harvard.edu/rinse/

- **nestor** project
  - http://nestordb.com

4

## Acknowledgements

**University of Paris**
- Michele Linardi
- Anna Gogolou
- Botao Peng
- Karima Echihabi
- Paul Boniol
- Federico Roncallo

**Harvard University**
- Kostas Zoumpatianos
- Stratos Idreos
- Niv Dayan

**Cornell University/Microsoft**
- Yin Lou
- Johannes Gehrke

**Inria**
- Anastasia Bezerianos
- Fanis Tsandilas
- Djamel-Edine Yagoubi
- Reza Akbarinia
- Florent Masseglia

**University of California at Riverside**
- Jin Shieh
- Eammon Keogh

**University of Crete/FORTH**
- Haridimos Kondylakis
- Panagiota Fatourou

**University of Trento**
- Alessandro Camerra

Themis Palpanas - University of Tokyo - Jan 2020

5

## Executive Summary

- data collected at unprecedented rates

- they enable data-driven scientific discovery

- lots of these data are sequences
  - takes **days-weeks** to analyze big sequence collections

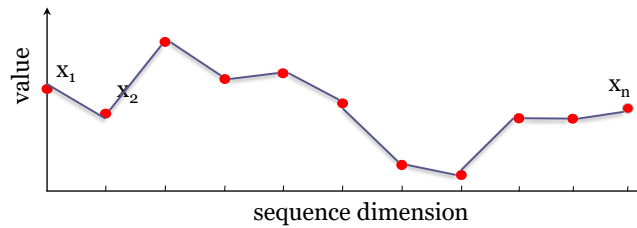goal: analyze big sequences in **minutes/seconds**

Themis Palpanas - University of Tokyo - Jan 2020

6

# Data series

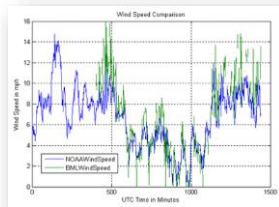- Sequence of points ordered along some dimension



Themis Palpanas - University of Tokyo - Jan 2020

7

# Scientific Monitoring

- meteorology, oceanography, astronomy, finance, sociology, …



Wind speed
From ocean observing node project
http://bml.ucdavis.edu/boon/wind.html

Historical stock quotes
http://money.cnn.com/2012/04/23/markets/walmart_stock/index.htm
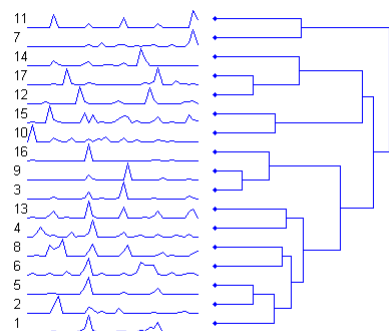
Themis Palpanas - University of Tokyo - Jan 2020

8

# Home Networks

- temporal usage behavior analysis of home networks
  - Portugal Telecom



(previously unknown) frequent behavior pattern
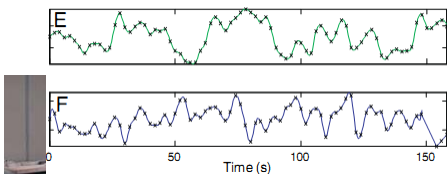
clustering based on user activity patterns
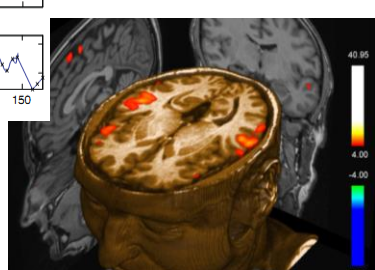
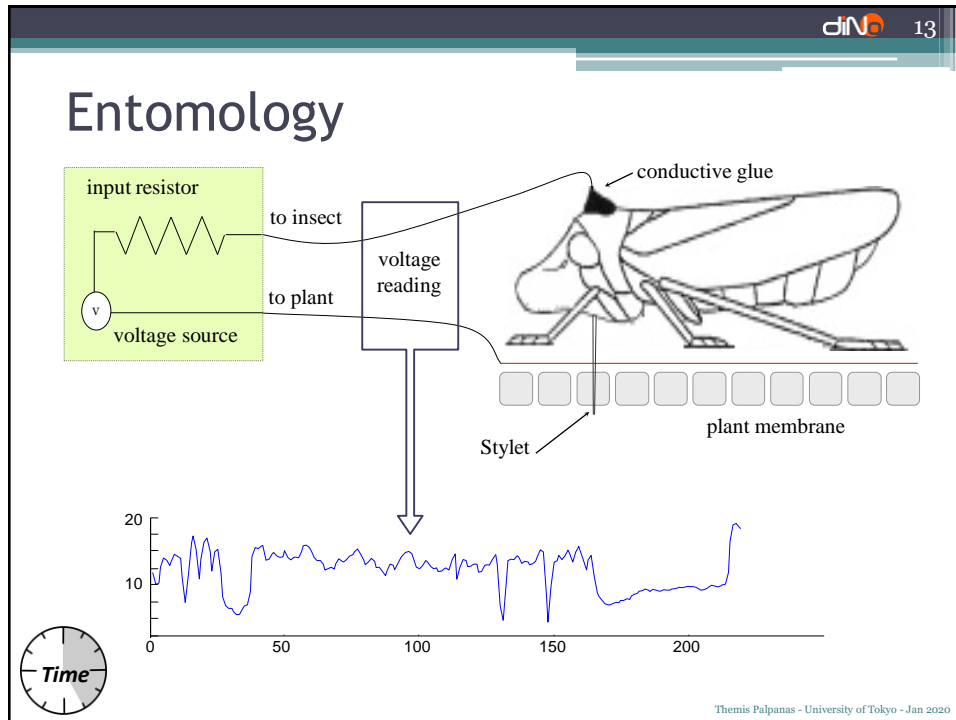Themis Palpanas - University of Tokyo - Jan 2020

10

# Neuroscience

- functional Resonance Magnetic Imaging (fMRI) data
  - primary experimental tool of neuroscientists
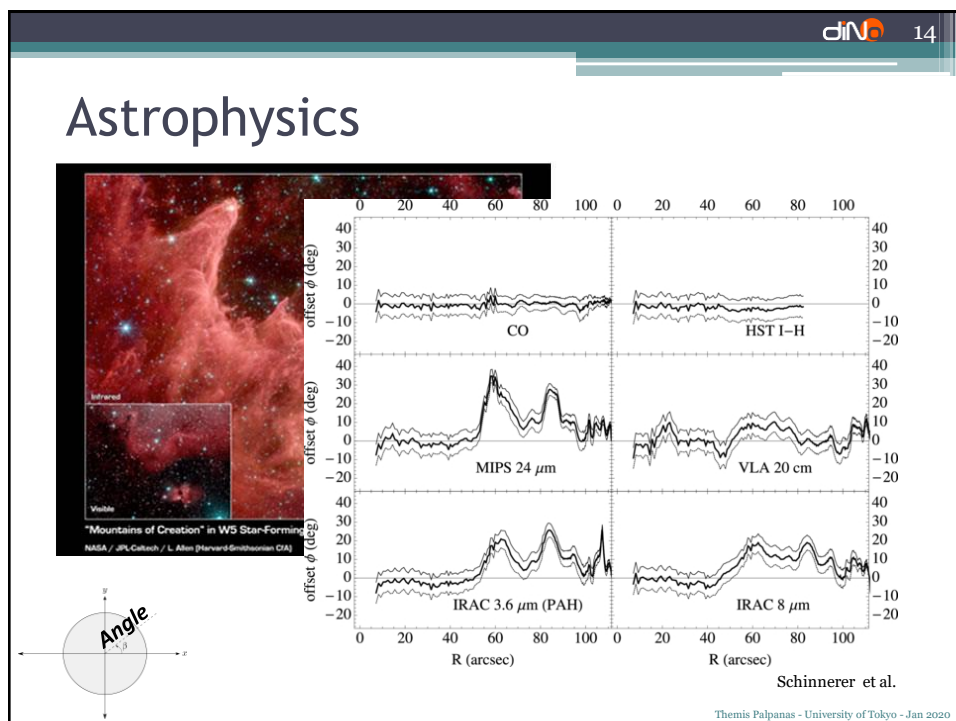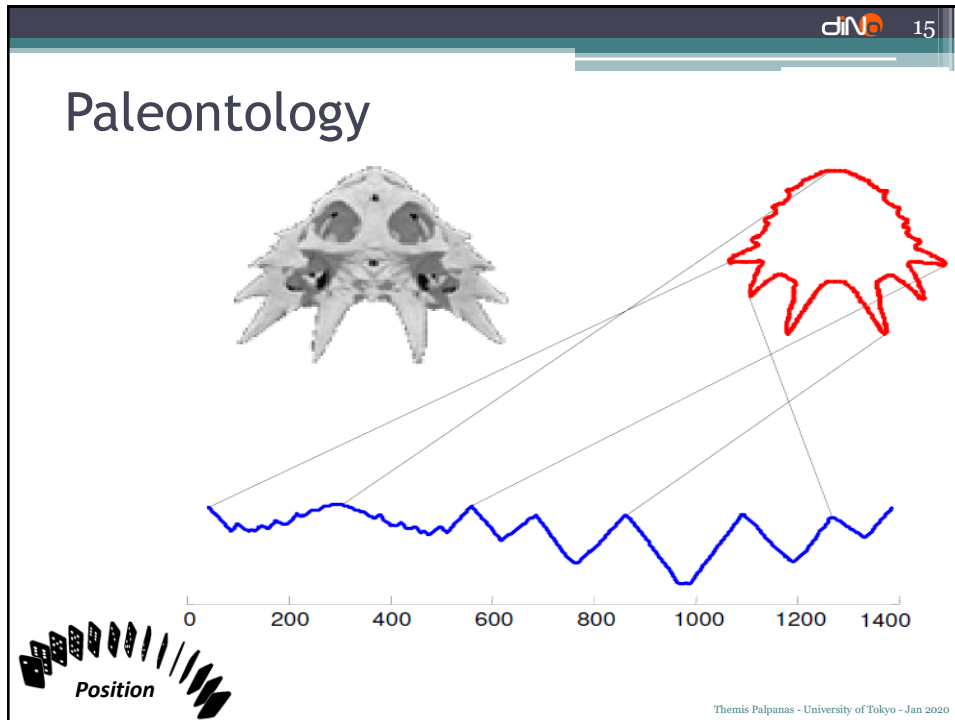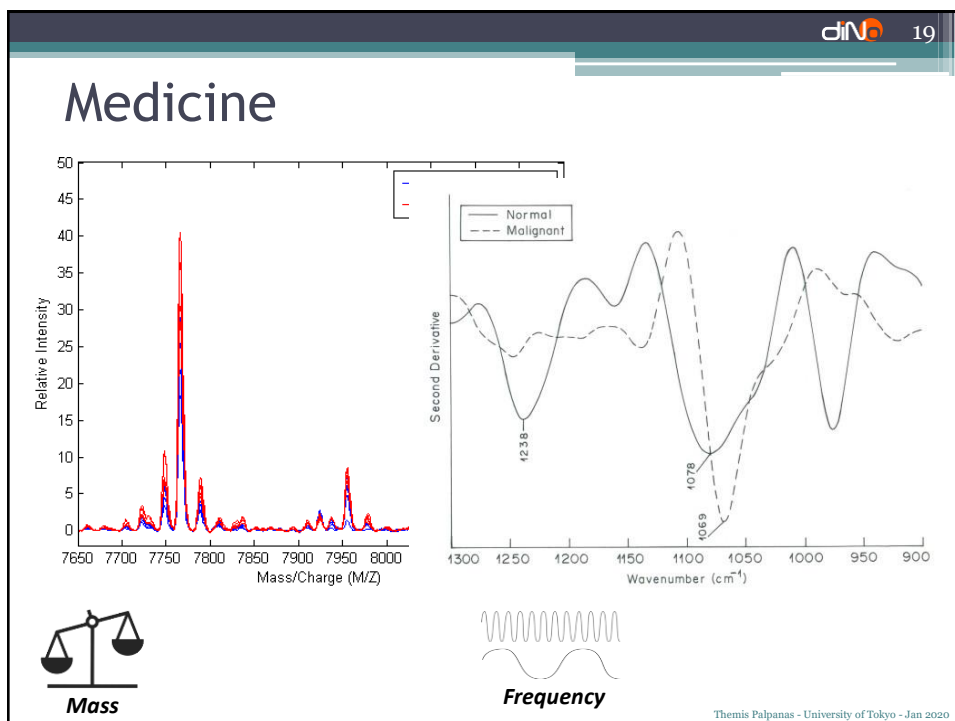  - reveal how different parts of brain respond to stimuli



12

13



14

## Paleontology



*Position*

Themis Palpanas - University of Tokyo - Jan 2020

15

## Medicine



*Mass*

*Frequency*

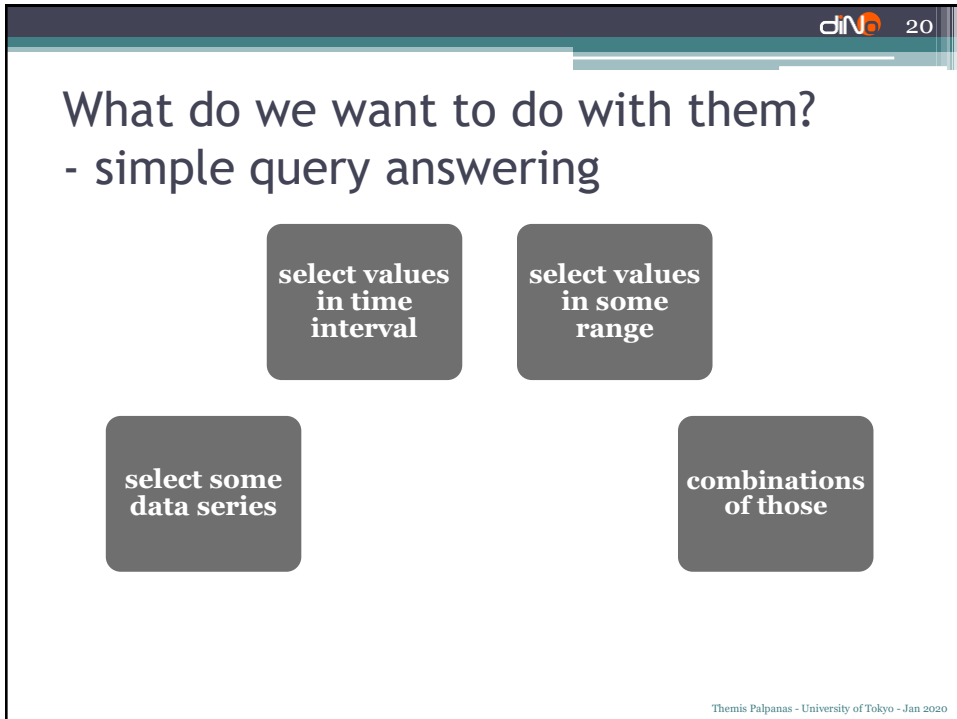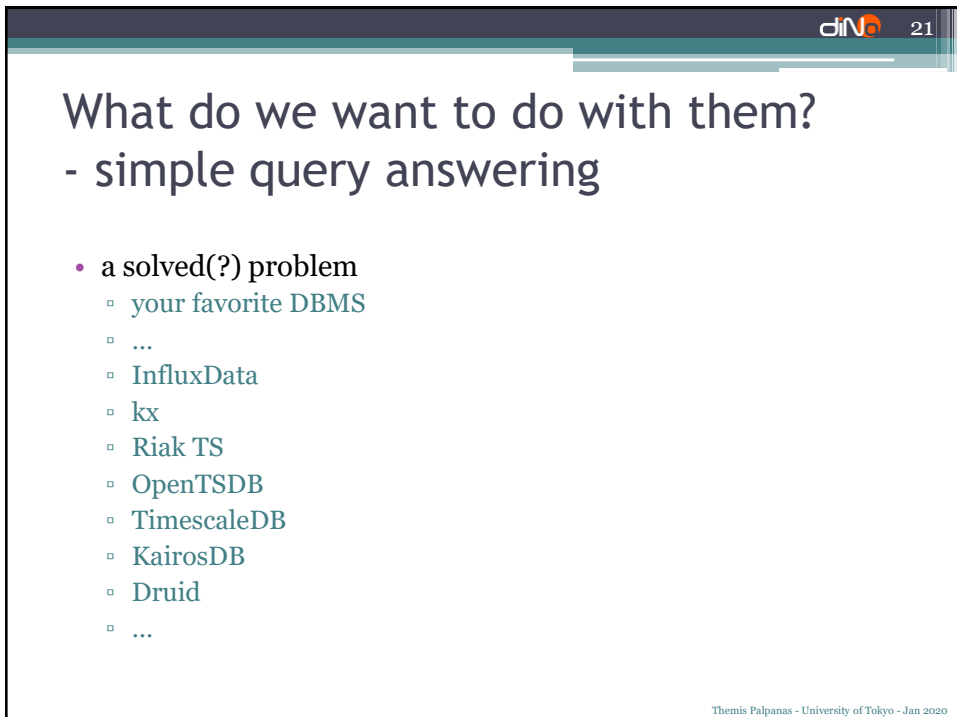Themis Palpanas - University of Tokyo - Jan 2020

19

# What do we want to do with them?
# - simple query answering

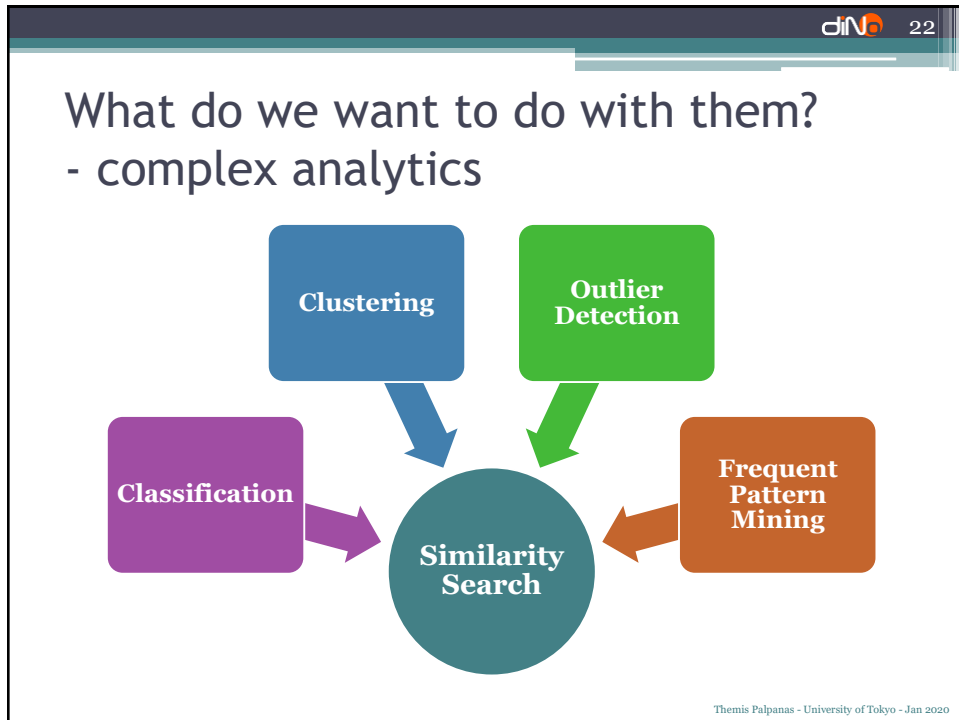| select values in time interval |
| select values in some range |
| select some data series |
| combinations of those |

20

# What do we want to do with them?
# - simple query answering
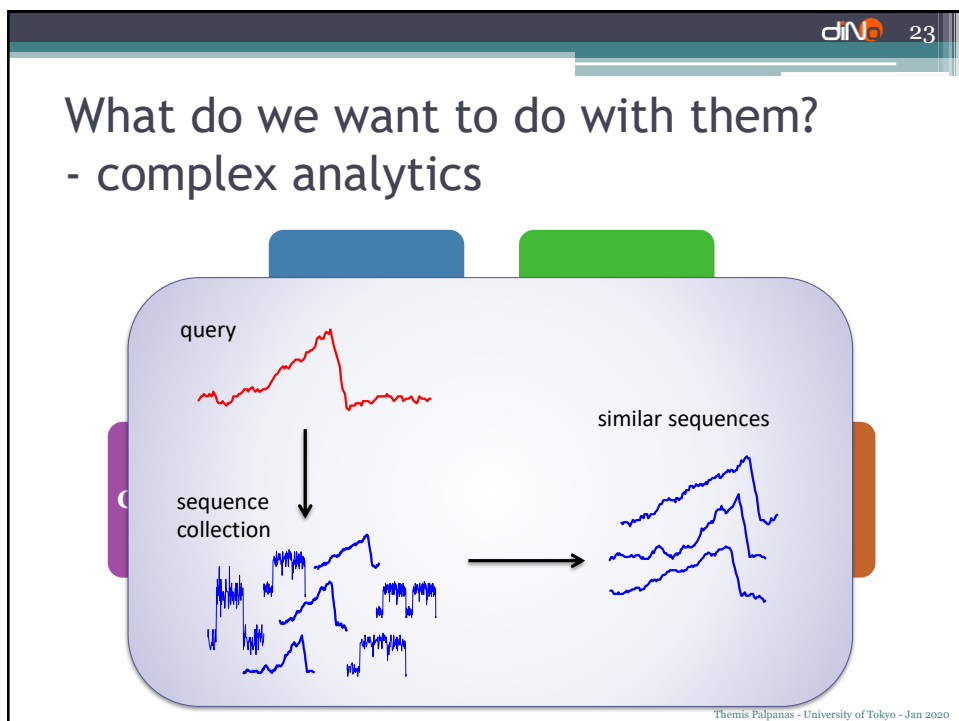
- a solved(?) problem
  - your favorite DBMS
  - ...
  - InfluxData
  - kx
  - Riak TS
  - OpenTSDB
  - TimescaleDB
  - KairosDB
  - Druid
  - ...

21

Slide 24: What do we want to do with them? - complex analytics

Euclidean

$$D(X,Y) \equiv \sqrt{\sum_{i=1}^{n}(x_i - y_i)^2}$$

Dynamic Time Warping (DTW)

$$D_{dtw}(X,Y) = f(n,m)$$

$$f(i,j) = \|x_i - y_j\| + \min \begin{cases} f(i, j-1) \\ f(i-1, j) \\ f(i-1, j-1) \end{cases}$$

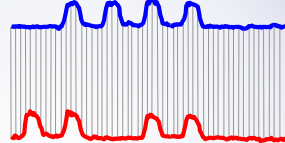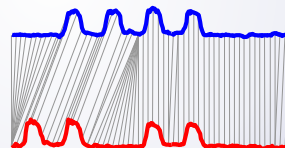Themis Palpanas - University of Tokyo - Jan 2020

24



Slide 25: What do we want to do with them? - complex analytics
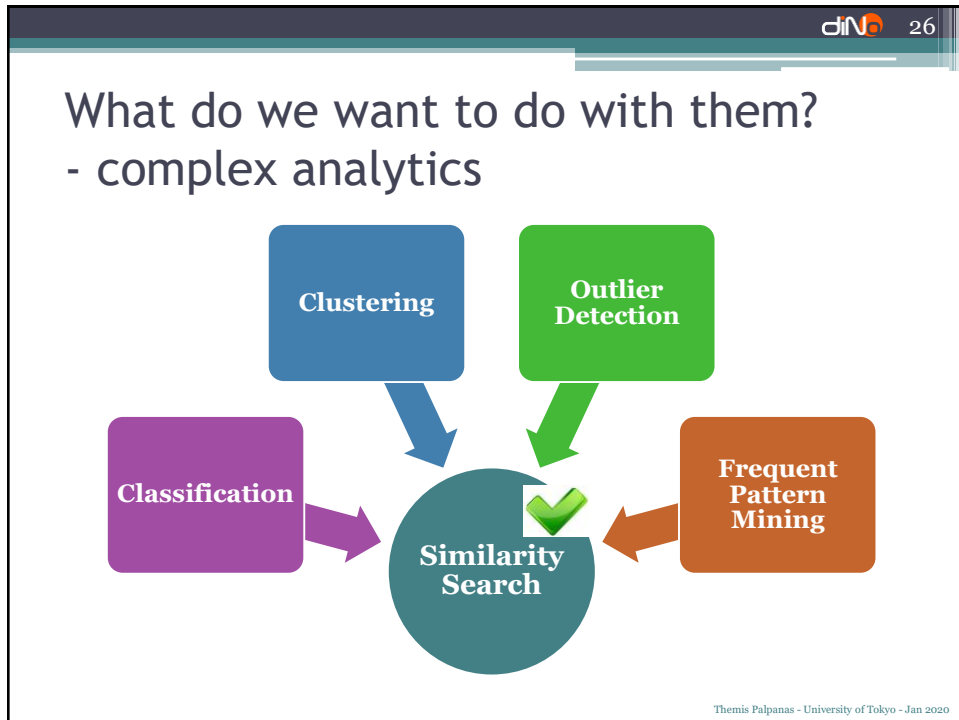
Euclidean

Dynamic Time Warping (DTW)

Themis Palpanas - University of Tokyo - Jan 2020

25

26



27

# What do we want to do with them? - complex analytics

**Clustering**

**Outlier Detection**

**HARD, because of very high dimensionality: each data series has 100s-1000s of points!**

Mining

**even HARDER, because of very large size: millions to billions of data series (multi-TBs)!**

Themis Palpanas - University of Tokyo - Jan 2020

28

28

# Query answering process

*Data Loading Procedure*

*Query Answering Procedure*

*Raw data*

*Data*

*Data Series Database/ Indexing*

*Queries*

*Answers*

**data-to-query** time

**query answering** time

*these times are big!*

Themis Palpanas - University of Tokyo - Jan 2020

29

29

# Similarity Search via
# **Serial Scan**

30

30

# Similarity Search via
# **Indexing**

31

31

13

# Traditional Approaches

answer **nearest neighbor queries** on a **1TB** dataset:

**serial scan** takes **45 minutes/query**

*Query Answering*

but building the index takes **too long!**

**indexing** a 1TB dataset takes **days**

a **data series index** can **reduce** querying time

*Query Answering*

**complex analytics in days…**

32

# Query answering process

*Data Loading Procedure*

*Query Answering Procedure*

*Raw data*

*Data*

*Data Series Database/ Indexing*

*Queries*

*Answers*

**data-to-query** time

**query answering** time

*we have proposed the state-of-the-art solutions for both problems!*

33

# State of the Art Approach: ADS+

Publications

SIGMOD'14

PVLDB'15

VLDBJ'16



**complex analytics in hours!**

Themis Palpanas - University of Tokyo - Jan 2020

34

34

---

diN◯  36

## SAX Representation

- **S**ymbolic **A**ggregate appro**X**imation (SAX)
  - **(1)** Represent data series $T$ of length $n$ with $w$ segments using Piecewise Aggregate Approximation (PAA)
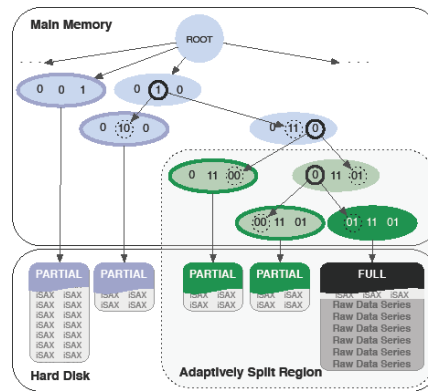    - $T$ typically normalized to $\mu = 0$, $\sigma = 1$

    - PAA$(T,w) = \overline{T} = \bar{t}_1, \ldots, \bar{t}_w$

      where $\bar{t}_i = \frac{w}{n} \sum_{j=\frac{n}{w}(i-1)+1}^{\frac{n}{w}i} T_j$

  - **(2)** Discretize into a vector of symbols
    - Breakpoints map to small alphabet $a$ of symbols



A data series $T$

PAA$(T,4)$

iSAX$(T,4,4)$

00
01
10
11

Themis Palpanas - University of Tokyo - Jan 2020

36

15

09-Jan-20



37



38

16

# Drawback of iSAX2+

- cannot start answering queries until *entire* index is built!

Themis Palpanas - University of Tokyo - Jan 2020

39

# Adaptive Data Series Index: ADS+

- novel paradigm for building a data series index
  - do not build entire index and then answer queries
  - start answering queries by building the part of the index needed by those queries

- still guarantee correct answers

Themis Palpanas - University of Tokyo - Jan 2020

40

# Adaptive Data Series Index: ADS+

Publications

SIGMOD'14

PVLDB'15

VLDBJ'16

- intuition for proposed solution

  - build the iSAX index using the iSAX representations
    - just like iSAX2+
  - but start with a large leaf size
    - minimize initial cost

  - postpone leaf materialization to query time
    - only materialize (at query time) leaves needed by queries
  - parts that are queried more are refined more
    - use smaller leaf sizes (reduced leaf materialization and query answering costs)

Themis Palpanas - University of Tokyo - Jan 2020

41

Start building an index with only the iSAX representations

ROOT

FBL

I1

I2

LBL

*RAM*

*DISK*

Raw data

Themis Palpanas - University of Tokyo - Jan 2020

43

43

Read the data-series one by one from the raw file

FBL

LBL

ROOT

I1

I2

RAM

DISK

Raw data

Themis Palpanas - University of Tokyo - Jan 2020

44

44



Convert them to iSAX

FBL

LBL

ROOT

I1

I2

RAM

DISK

Raw data

Themis Palpanas - University of Tokyo - Jan 2020

45

45

Store only iSAX in memory (64 times smaller) ~1%

FBL

ROOT

I1

I2

LBL

*RAM*

*DISK*

Raw data

46

46

Discard raw data and keep pointer to raw file

FBL

ROOT

I1

I2

LBL

*RAM*

*DISK*

Raw data

47

47

Continue loading data until we run out of memory

ROOT

FBL

I1  I2

LBL

*RAM*

*DISK*

Raw data

Themis Palpanas - University of Tokyo - Jan 2020

48

48



Expand each sub-tree and move data to LBL

ROOT

FBL

I1  I2

LBL

L1  L2  L3  L4

*RAM*

*DISK*

Raw data

Themis Palpanas - University of Tokyo - Jan 2020

49

49

We flush the data to the disk to free up memory

50



59

09-Jan-20

60

Themis Palpanas - University of Tokyo - Jan 2016

61

23

Query #1 — ROOT, FBL, I1, I2, LBL, L1, L2, L4, L5, RAM, DISK, Raw data, TOO BIG!, PARTIAL, PARTIAL, PARTIAL, PARTIAL

Themis Palpanas - University of Tokyo - Jan 2020

62

62



Query #1 — **Create** a **smaller** leaf — ROOT, FBL, I1, I2, *Adaptive split*, LBL, I3, L4, L5, L2, L4, L5, RAM, DISK, Raw data, PARTIAL, PARTIAL, PARTIAL, PARTIAL, PARTIAL

Themis Palpanas - University of Tokyo - Jan 2020

63

63

64



65

25

## Experimental Evaluation



- iSAX 2.0 needs more than 35 hours to answer 100K approximate queries
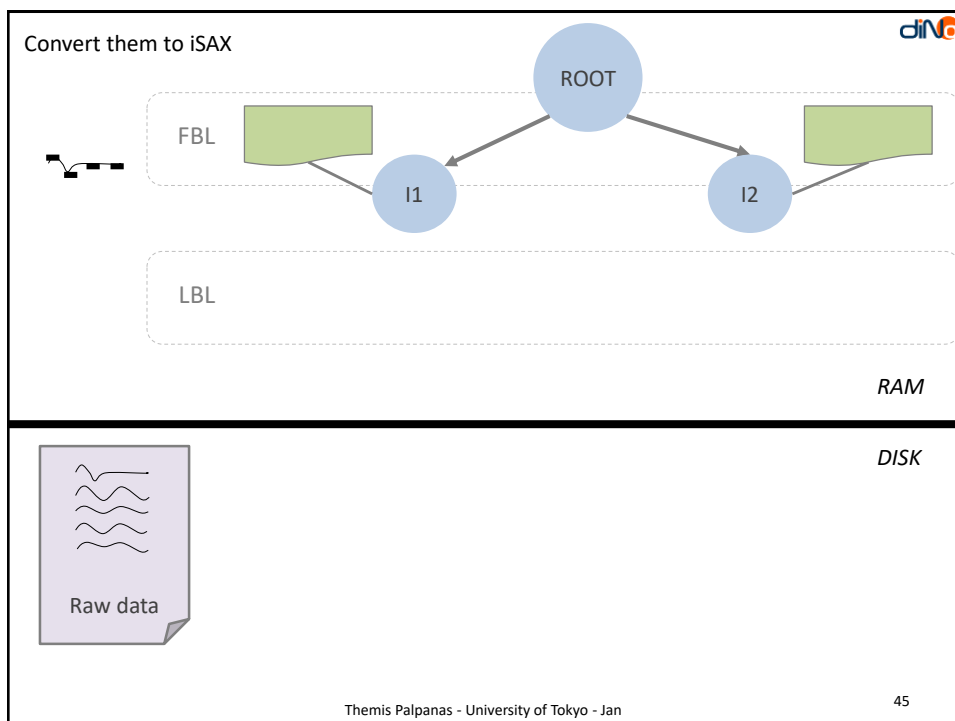- ADS+ answers 100K approximate queries in less than 5 hours

Themis Palpanas - University of Tokyo - Jan 2020

66

---

# Comparison to
# *multi-dimensional indices*

measure data-to-query time
(just index 1 **billion** data-series)



**1-3 orders of magnitude faster** than multi-dimensional indexing methods

Themis Palpanas - University of Tokyo -

67

09-Jan-20

diNo  69

Publications

PVLDB'18

SIGMOD'19

VLDBJ'20

# Extensions...

- Coconut: current solution for limited memory devices
  and streaming time series
  - bottom-up, succinct index construction based on sortable
    summarizations

Themis Palpanas - University of Tokyo - Jan 2020

69

---

diNo  70

Publications

PVLDB'18

SIGMOD'19

VLDBJ'20

# Extensions...

- Co



Themis Palpanas - University of Tokyo - Jan 2020

70

27

diN♥  71

**Extensions…**

Publications

PVLDB'18

SIGMOD'19

VLDBJ'20

- Coconut: current solution for limited memory devices
  and streaming time series
  - bottom-up, succinct index construction based on sortable summarizations
  - outperforms state-of-the-art in terms of index space, index construction time, and query answering time

Themis Palpanas - University of Tokyo - Jan 2020

71

diN♥  72

**Extensions…**

Publications

PVLDB'18

SIGMOD'19

VLDBJ'20

ICDE'18

PVLDB'19

- Coconut: current solution for limited memory devices
  and streaming time series
  - bottom-up, succinct index construction based on sortable summarizations
  - outperforms state-of-the-art in terms of index space, index construction time, and query answering time

- ULISSE: current solution for variable-length queries
  - single-index support of queries of variable lengths

Themis Palpanas - University of Tokyo - Jan 2020

72

## Slide 73

# Extensions…

**Publications**

- PVLDB'18
- SIGMOD'19
- VLDBJ'20
- ICDE'18
- PVLDB'19

- **Coconut**: current solut...
  and streamir...
  - bottom-up, succinct in...
    summarizations
  - outperforms state-of-t... ...ex
    construction time, and...

- **ULISSE**: current solut...
  - single-index support o...

## Slide 74

# Extensions…

**Publications**

- PVLDB'18
- SIGMOD'19
- VLDBJ'20
- ICDE'18
- PVLDB'19

- **Coconut**: current solution for limited memory devices
  and streaming time series
  - bottom-up, succinct index construction based on sortable
    summarizations
  - outperforms state-of-the-art in terms of index space, index
    construction time, and query answering time

- **ULISSE**: current solution for variable-length queries
  - single-index support of queries of variable lengths
  - orders of magnitude faster than competing approaches

## Slide 75

# Extensions…

Publications

PVLDB'18

SIGMOD'19

VLDBJ'20

ICDE'18

PVLDB'19

- Coconut: current solution for limited memory devices and streaming time series
  - bottom-up, succinct index construction based on sortable summarizations
  - outperforms state-of-the-art in terms of index space, index construction time, and query answering time

- ULISSE: current solution for variable-length queries
  - single-index s
  - orders of mag

**Michele Linardi:**
**BDA Best PhD Thesis Award (2019)**

Themis Palpanas - University of Tokyo - Jan 2020

75

## Slide 76

problems solved

declare success!

Themis Palpanas - University of Tokyo - Jan 2020

76

problems solved

declare success!

# well, not so fast...

Themis Palpanas - University of Tokyo - Jan 2020

77

# Massive Data Series Collections

- functional Resonance Magnetic Imaging (fMRI) data
  - primary experimental tool of neuroscientists
  - reveal how different parts of brain respond to stimuli
  - single experiment (1 subject, 1 test) produces
    - 60,000 data series of length 3,000: 12 GB

Tokyo - Jan 2020

78

## Slide 79

diN☺ 79

# Massive Data Series Collections

- ADHD-200 Global Competition
  - classification task: detect Attention Deficit Hyperactivity Disorder
    - 776 subjects: 9 TB
      - equivalent to: 4.5 billion non-overlapping data series of size 256
      - equivalent to: 1100 billion overlapping data series of size 256



Tokyo - Jan 2020

79

## Slide 80

# Massive Data Series Collections

Publications

SIGREC'19

**NASA's Solar Observatory**
**1.5 TB per day**

**Large Synoptic Survey Telescope (2019)**
**~30 TB per night**

**Human Genome project**
**130 TB**

**passenger aircrafts**
**20 TB per hour**

**data center and services monitoring**
**2B data series**
**4M points/sec**

facebook

Themis Palpanas - University of Tokyo - Jan 2020

80

80

## Slide 81

# The Road Ahead

Publications

ICDE'18

HPCS'17

SIGREC'15

*"enable practitioners and non-expert users to easily and efficiently manage and analyze massive data series collections"*

Themis Palpanas - University of Tokyo - Jan 2020

81

## Slide 82

# The Road Ahead

Publications

ICDE'18

HPCS'17

SIGREC'15

- Big Sequence Management System
  - general purpose data series management system

Data Model
Summarizations | Query Language

Data Structures
Access Methods

Holistic Optimization | Varying Length Queries | Uncertain Sequences

Distributed Processing

data sequences

Themis Palpanas - University of Tokyo - Jan 2020

82

## Slide 83

Publications

ICDE'18

HPCS'17

SIGREC'15

# The Road Ahead

• Big Sequence Management System

| Data Model | | |
|---|---|---|
| Summarizations | Query Language | |

**Holistic Optimization**

**Data Structures**

Access Methods

| Varying Length Queries | Uncertain Sequences |
|---|---|

**Distributed Processing**

Spark / Flink / (HDFS)

Themis Palpanas - University of Tokyo - Jan 2020

83

## Slide 84

Publications

ICDE'18

HPCS'17

SIGREC'15

PVLDB'19

# The Road Ahead

• Big Sequence Management System

Data Model

Summarization

**Holistic Optimization**

D

Varying Quer

Distr

Spark /

|  | Dataset | Idx | Scenarios | | | | |
|---|---|---|---|---|---|---|---|
|  |  |  | Exact 100 | Idx+ Exact 100 | Idx+ Exact 10K | Exact Easy-20 | Exact Hard-20 |
| HDD | Small | A | D | S | D | D | D |
|  | Large | A | D | S | D | D | D |
|  | Astro | A | U | U | V | V | U |
|  | Deep1B | A | U | U | U | D | U |
|  | SALD | A | D | I | D | D | D |
|  | Seismic | A | D | S | D | D | U |
| SSD | Small | S | D | I | D | I | D |
|  | Large | S | D | I | D | I | D |
|  | Astro | I | V | V | V | V | V |
|  | Deep1B | S | I | I | V | I | U |
|  | SALD | S | I | I | I | I | V |
|  | Seismic | A | V | V | V | D | V |

A: ADS  D: DSTree,  I: iSAX2+
S: SFA  U: UCR-Suite,  V: VA+file

84

# Parallelization/Distribution?

- discussion so far assumed serial execution in a single core
  - focus on efficient resource utilization
  - squeeze the most out of a single core
  - produce scalable solutions at lowest possible cost
    - also suitable for analysts with no access to/expertise for clusters

diNo 87

Themis Palpanas - University of Tokyo - Jan 2020

87

# Need for Parallelization/Distribution

diNo 88

Publications

HPCS'17

- take advantage of modern hardware!
  - Single Instruction Multiple Data (SIMD)
    - natural for data series operations
  - multi-tier CPU caches
    - design data structures aligned to cache lines
  - multi-core and multi-socket architectures
    - use parallelism inside each computation server
  - Graphics Processing Units (GPUs)
    - propose massively parallel techniques for GPUs
  - new storage solutions: SSDs, NVRAM
    - develop algorithms that take these new characteristics/tradeoffs into account
  - compute clusters
    - distribute operation over many machines

Themis Palpanas - University of Tokyo - Jan 2020

88

## Slide 89

Publications

HPCS'17

# Need for Parallelization/Distribution

- further scale-up and scale-out possible!
  - ▫ techniques inherently parallelizable
    - · across cores, across machines



compute nodes

compute node number

parallelized data series index

ends computation early, based on information from other nodes

data series collection

subset of collection that contains the answer

Themis Palpanas - University of Tokyo - Jan 2020

89

## Slide 90

Publications

HPCS'17

# Need for Parallelization/Distribution

- further scale-up and scale-out possible!
  - ▫ techniques inherently parallelizable
    - · across cores, across machines

- more involved solutions required when optimizing for energy
  - ▫ reducing execution time is relatively easy
  - ▫ minimizing total work (energy) is more challenging

Themis Palpanas - University of Tokyo - Jan 2020
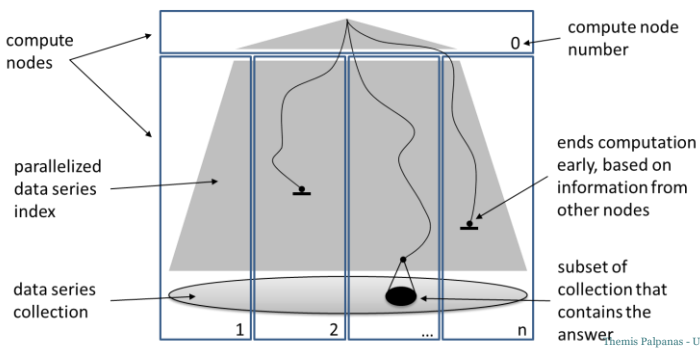
90

# Need for Parallelization/Distribution

Publications

ICDM'17

TKDE'18

PKDD'19

- DPiSAX: current solution for distributed processing (Spark)
  - balances work of different worker nodes

Themis Palpanas - University of Tokyo - Jan 2020

91

# Need for Parallelization/Distribution

Publications

ICDM'17

TKDE'18

PKDD'19

- DPiSAX: current solution for distributed processing (Spark)
  - balances work of different worker nodes



Themis Palpanas - University of Tokyo - Jan 2020

92

37

---

**diNo** 93

Publications

ICDM'17

TKDE'18

PKDD'19

BigData'18

# Need for Parallelization/Distribution

- **DPiSAX**: current solution for distributed processing (Spark)
  - ▫ balances work of different worker nodes
  - ▫ performs 2 orders of magnitude faster than centralized solution

- **ParIS**: current solution for modern hardware
  - ▫ completely masks out the CPU cost

Themis Palpanas - University of Tokyo - Jan 2020

93

---

**diNo** 95

Publications

ICDM'17

TKDE'18

PKDD'19

BigData'18

# Need for Parallelization/Distribution

- **DPiSAX**: current solution for distributed processing (Spark)
  - ▫ balances work of different worker nodes



Themis Palpanas - University of Tokyo - Jan 2020
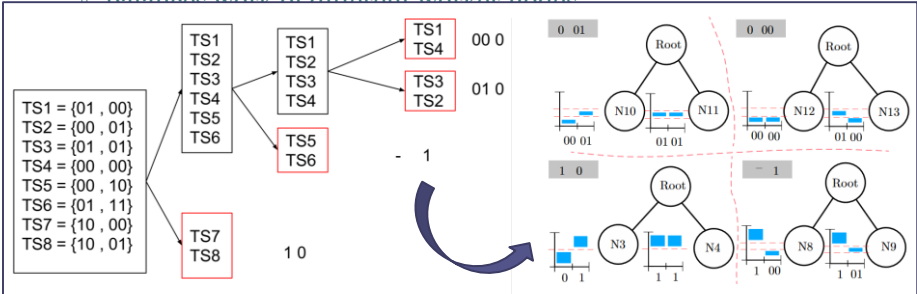
95

## Slide 96
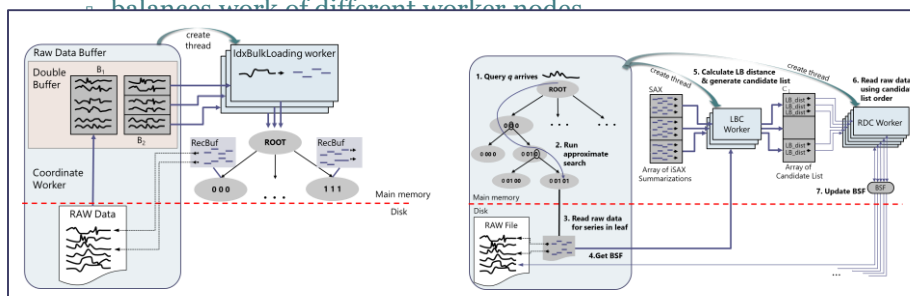
# Need for Parallelization/Distribution

Publications

ICDM'17

TKDE'18

PKDD'19

BigData'18

- DPiSAX: current solution for distributed processing (Spark)
  - balances work of different worker nodes
  - performs 2 orders of magnitude faster than centralized solution

- ParIS: current solution for modern hardware
  - masks out the CPU cost
  - answers exact queries in the order of a few secs
    - >1 order of magnitude faster then single-core solutions

Themis Palpanas - University of Tokyo - Jan 2020

96

## Slide 97

# Need for Parallelization/Distribution

Publications

ICDM'17

TKDE'18

PKDD'19

BigData'18

- DPiSAX: current solution for distributed processing (Spark)
  - balar
  - perfo

- ParIS:
  - mask
  - answ
    - >1



k-NN Classification

18x faster

Themis Palpanas - University of Tokyo - Jan 2020

97

## Slide 98

Publications

ICDM'17
TKDE'18
PKDD'19
BigData'18

# Need for Parallelization/Distribution

- DPiSAX: current solution for distributed processing (Spark)
  - balar ...

**k-NN Classification**

**classifying 100K objects using a 100GB dataset goes down from several days to few hours!**

- answ
  - >1

18x faster

1−NN    5−NN    10−NN    50−NN
Number of nearest neighbors

Themis Palpanas - University of Tokyo - Jan 2020

98

## Slide 99

Publications

ICDM'17
TKDE'18
PKDD'19
BigData'18
ICDE'20

# Need for Parallelization/Distribution

- DPiSAX: current solution for distributed processing (S
  - balances work of different worker nodes
  - performs 2 orders of magnitude faster than centralized solution

- ParIS: current solution for modern hardware
  - masks out the CPU cost
  - answers exact queries in the order of a few secs
    - >1 order of magnitude faster then single-core solutions

- MESSI: current solution for modern hardware + in-memory data
  - answers exact queries in the order of ms

Themis Palpanas - University of Tokyo - Jan 2020

99

# Interactive Analytics?

- data series analytics is computationally expensive
  - very high inherent complexity

- may not always be possible to remove delays
  - but could try to hide them!

100

# Need for
# Interactive Analytics

Publications

BigVis'19

- interaction with users offers new opportunities
  - progressive answers
    - produce intermediate results
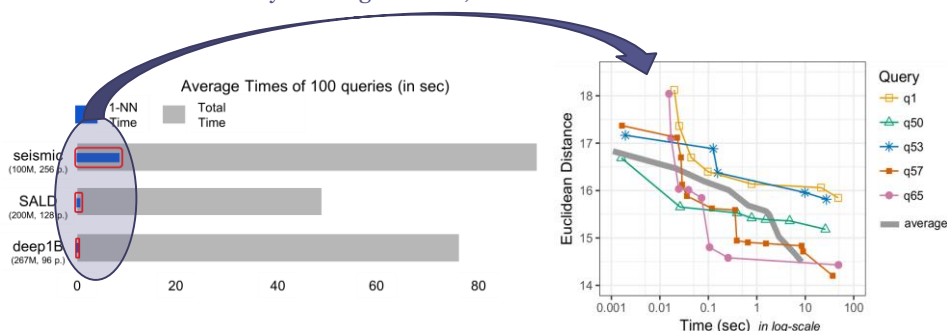      - iteratively converge to final, correct solution

101

# Need for Interactive Analytics

Publications

BigVis'19

VIS'18

- interaction with users offers new opportunities
  - progressive answers
    - produce intermediate results
      - iteratively converge to final, correct solution
    - provide bounds on the errors (of the intermediate results) along the way
  - imprecise queries
    - enable user to specify varying accuracy requirements for different parts of the same query

- several exciting research problems in intersection of visualization and data management
  - *frontend*: HCI/visualizations for querying/results display
  - *backend*: efficiently supporting these operations

Themis Palpanas - University of Tokyo - Jan 2020

102

# Data Series vs. high-d Vectors

- two sides of the same(?) coin
  - data series as multidimensional points
  - for a specific ordering of the dimensions

- several techniques for similarity search in high-d vectors
  - using LSH (SRS), space quantization (IMI), k-NN graphs (HNSW)

- how do these high-d vector techniques compare to data series techniques?
  - currently conducting extensive experimental comparison

Themis Palpanas - University of Tokyo - Jan 2020

103

09-Jan-20

# Data Series vs. high-d Vectors

Publications

PVLDB'20

- **data series techniques** are the **overall winners**, even on **general high-d vector** data

Themis Palpanas - University of Tokyo - Jan 2020

104

# Data Series vs. high-d Vectors

Publications

PVLDB'20

- **data series techniques** are the **overall winners**, even on **general high-d vector** data
  - perform the best for approximate queries with probabilistic guarantees (δ-ε-approximate search), in-memory and on-disk



(s) Deep25GB(ng)  (t) Deep25GB($\delta\epsilon$)

△ DSTree  ⊕ HNSW  ◇ IMI  ⊟ iSAX2+  ⊠ SRS  + VA+file

Themis Palpanas - University of Tokyo - Jan 2020

105

43

# Data Series vs. high-d Vectors

Publications

PVLDB'20

- **data series techniques** are the **overall winners**, even on **general high-d vector** data
  - ▫ perform the best for approximate queries with probabilistic guarantees (δ-ε-approximate search), in-memory and on-disk

HNSW: only for in-memory data, with no guarantees for the answers

(s) **Deep25GB(ng)** (t) **Deep25GB(δε)**

△— DSTree  ⊕— HNSW  ◇— IMI  □— iSAX2+  ⊠— SRS  +— VA+file  *mis Palpanas - University of Tokyo - Jan 2020*

---

# Data Series vs. high-d Vectors
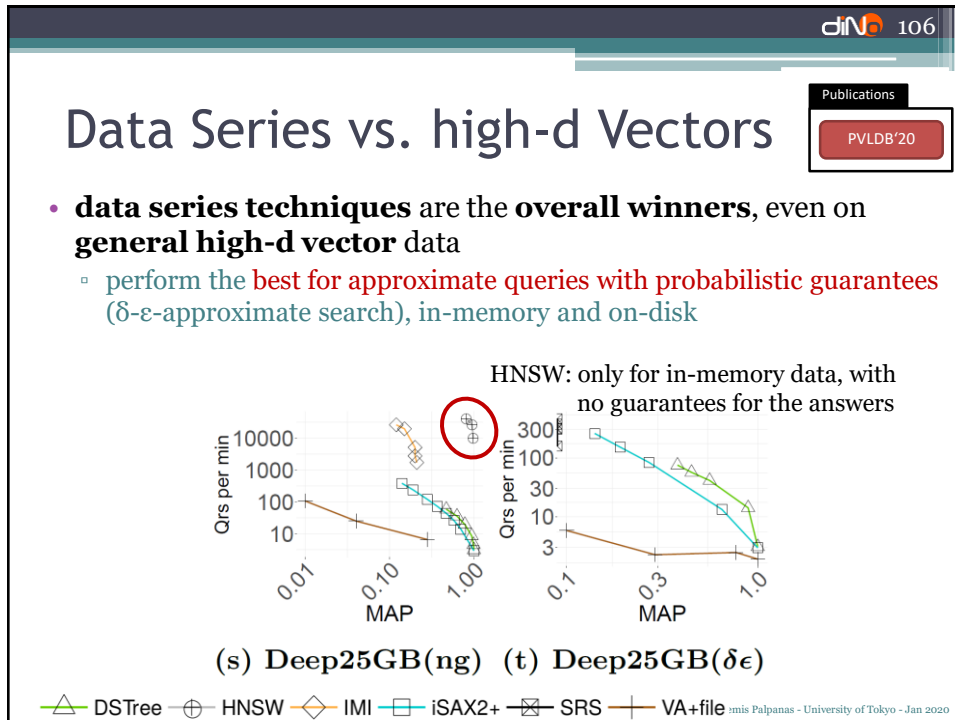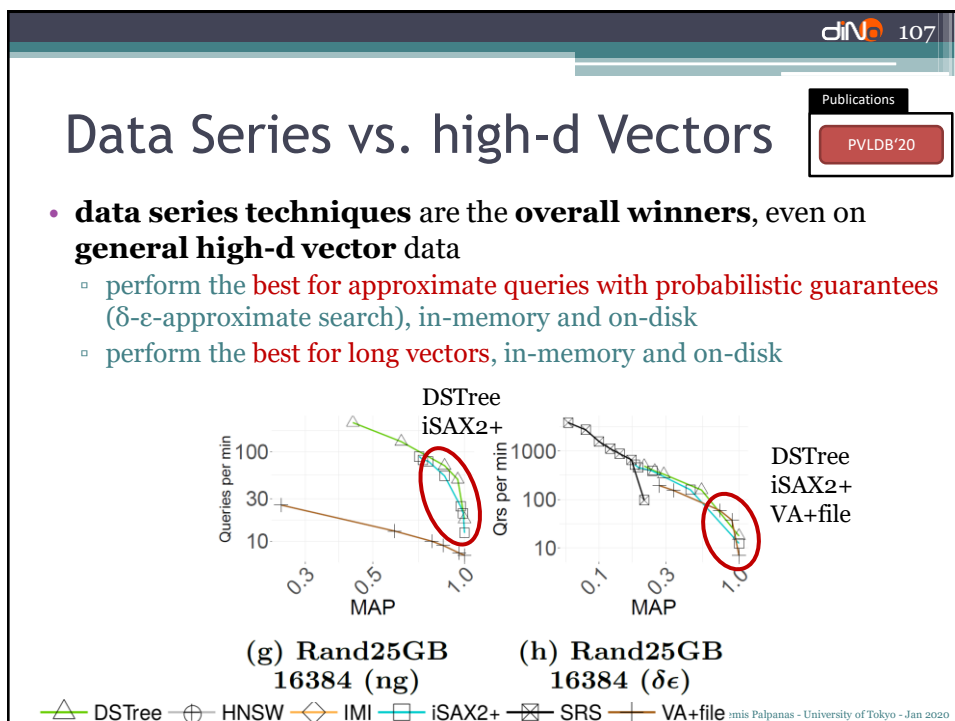
Publications

PVLDB'20

- **data series techniques** are the **overall winners**, even on **general high-d vector** data
  - ▫ perform the best for approximate queries with probabilistic guarantees (δ-ε-approximate search), in-memory and on-disk
  - ▫ perform the best for long vectors, in-memory and on-disk

DSTree
iSAX2+

DSTree
iSAX2+
VA+file

(g) **Rand25GB**
**16384 (ng)** (h) **Rand25GB**
**16384 (δε)**

△— DSTree  ⊕— HNSW  ◇— IMI  □— iSAX2+  +— SRS  +— VA+file  *mis Palpanas - University of Tokyo - Jan 2020*

# Data Series vs. high-d Vectors

Publications

PVLDB'20

- **data series techniques** are the **overall winners**, even on **general high-d vector** data
  - ▫ perform the best for approximate queries with probabilistic guarantees (δ-ε-approximate search), in-memory and on-disk
  - ▫ perform the best for long vectors, in-memory and on-disk
  - ▫ perform the best for disk-resident vectors (the only viable solution)



DSTree
iSAX2+

(m)
Deep250GB(ng)

(n)
Deep250GB($\delta\epsilon$)

△ DSTree  ⊕ HNSW  ◇ IMI  ☐ iSAX2+  ⊠ SRS  + VA+file   Themis Palpanas - University of Tokyo - Jan 2020

108

# iSAX Index Family Lineage Tree

Publications

ISIP'19



Lineage of the iSAX family of indexes. Timeline is depicted on the top; implementation languages are marked on the right. Solid arrows denote inheritance of the index design; dashed arrows denote inheritance of some of the design features; the two new versions of iSAX2+ and ADS+ marked with an asterisk support approximate similarity search with deterministic and probabilistic quality guarantees.

Themis Palpanas - University of Tokyo - Jan 2020

109

# Benchmarking Data Series Indexes?

110

# Previous Studies

evaluate **performance** of **indexing methods** using **random queries**

- chosen from the data (with/without noise)

111

111

# Previous Studies

**With or without noise**
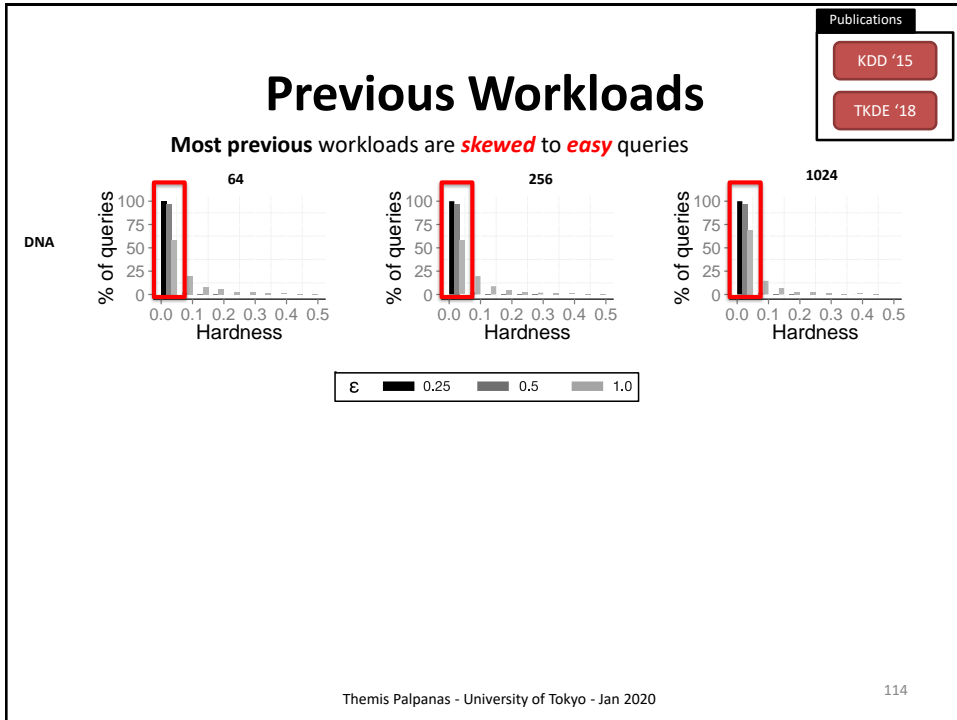
112

# Problem with
# Random Queries

*No control* on their *characteristics*

We **cannot properly evaluate** summarizations and indexes

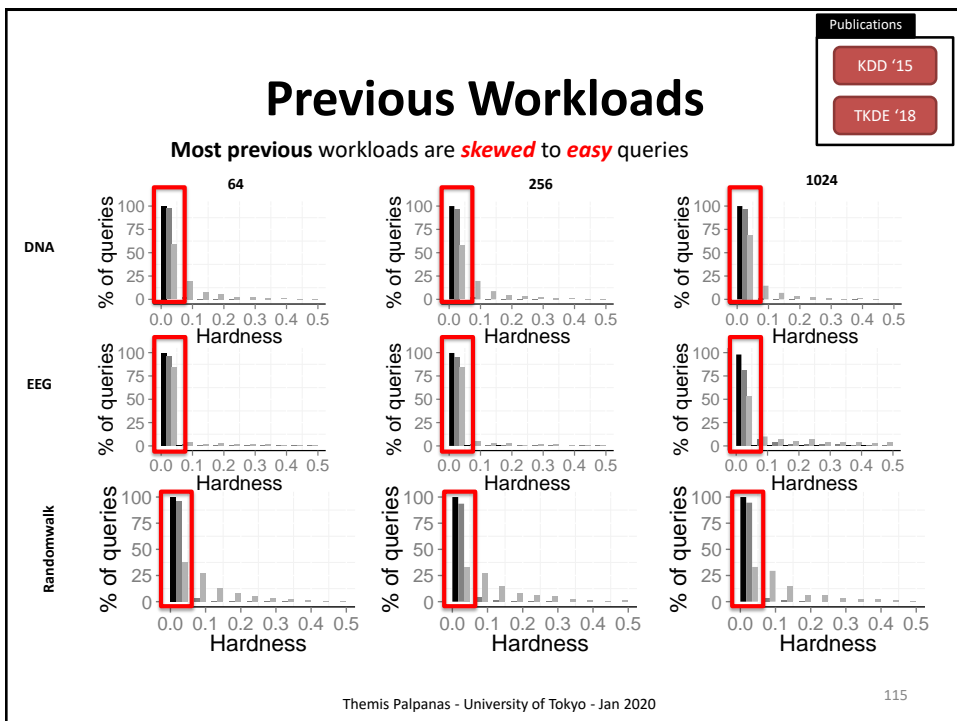**We need queries that cover the entire range
from easy to hard**

113

114

115

# Benchmark Workloads

If all queries are **easy**
all indexes look **good**

If all queries are **hard**
all indexes look **bad**

need **methods** for **generating** queries of **varying hardness**

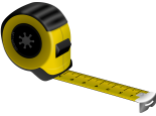Themis Palpanas - University of Tokyo - Jan 2020

116

116

# Contributions

**Theoretical background**
Methodology for characterizing
NN queries for data series indexes

**Nearest neighbor query workload generator**
Designed to stress-test data series indexes
at varying levels of difficulty

Themis Palpanas - University of Tokyo - Jan 2020

117

117

# Subsequence Anomaly Detection

144

## Data Series Anomalies Problem

● develop anomaly detection techniques based on sequences (data series), not on individual values

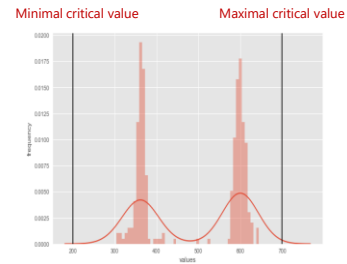○ individual values can be normal, but their sequence can be abnormal!

145

145

## Data Series Anomalies Problem

150 points in a sequence S

- develop anomaly detection techniques based on sequences (data series), not on individual values
  - ○ individual values can be normal, but their sequence can be abnormal!

Minimal critical value    Maximal critical value



values are not outside critical thresholds
values are normal

146

146

## Data Series Anomalies Problem

Sequence S

- develop anomaly detection techniques based on sequences (data series), not on individual values
  - ○ individual values can be normal, but their sequence can be abnormal!



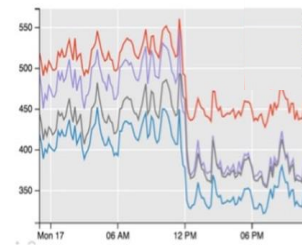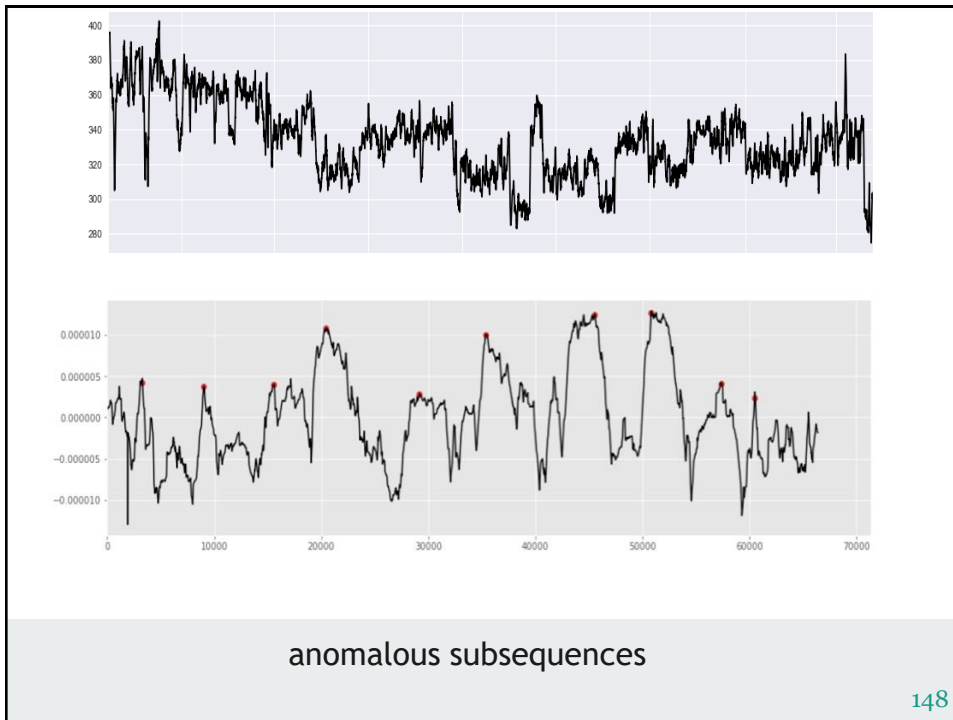values are not outside critical thresholds
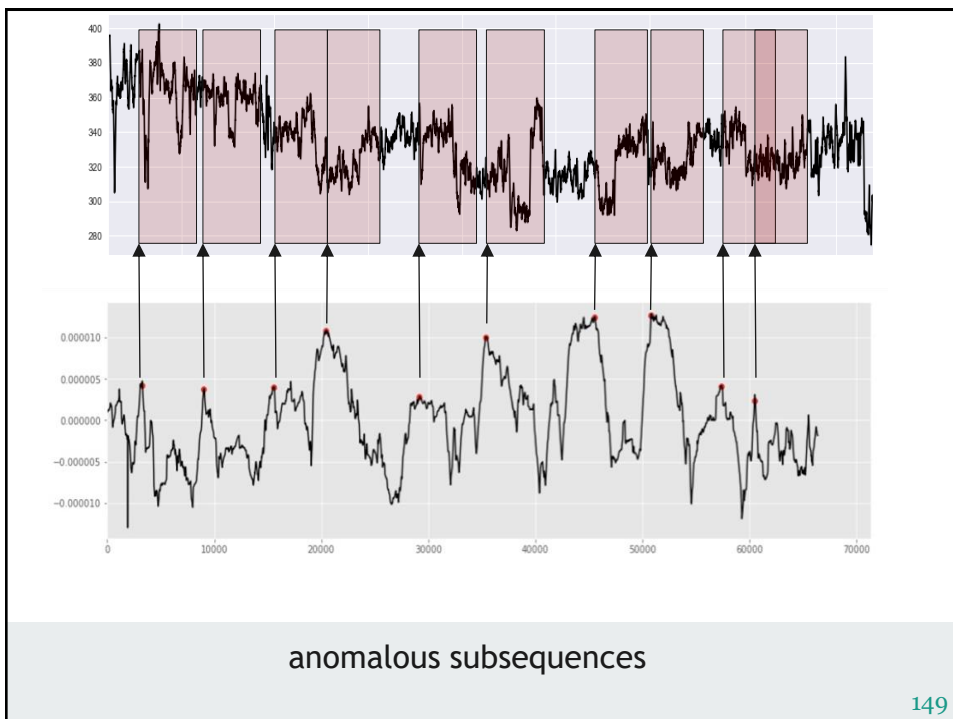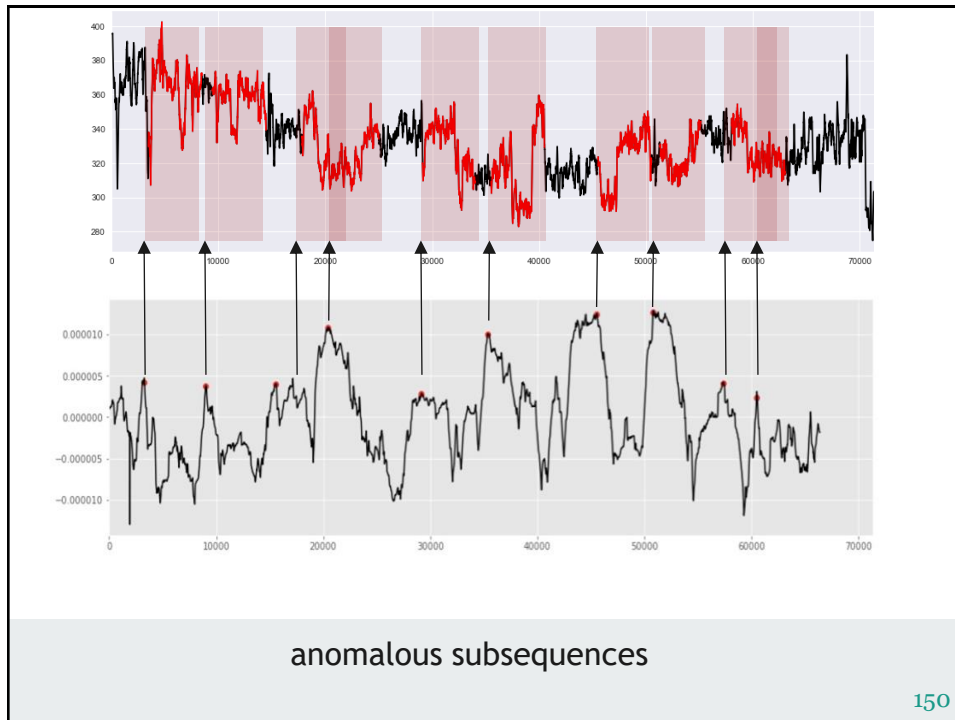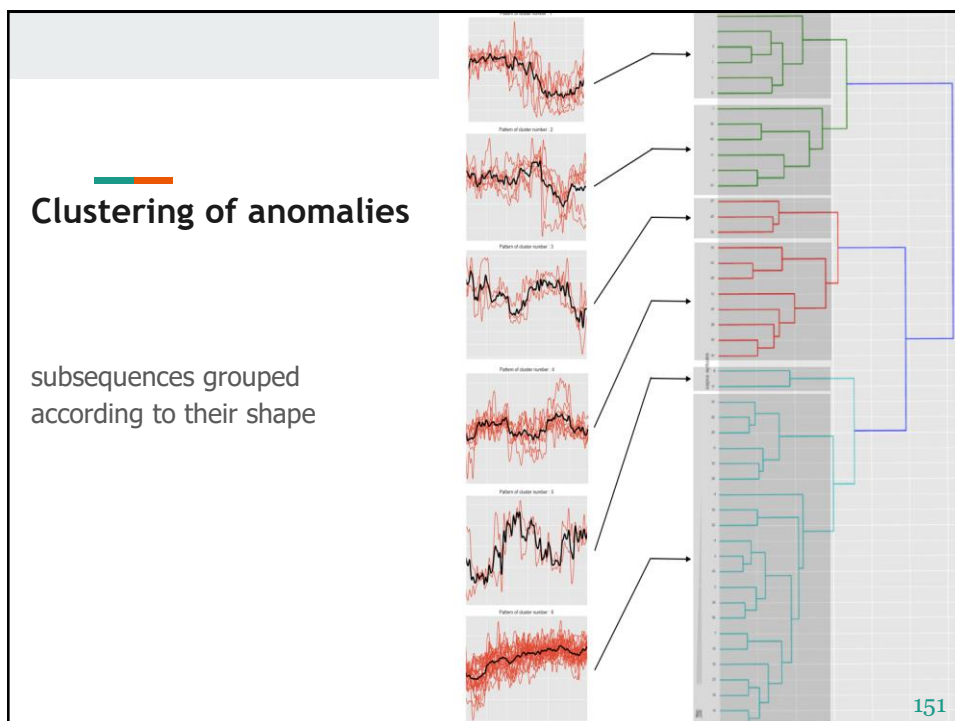values are normal
sequences are abnormal

147

147

anomalous subsequences

148

148



anomalous subsequences

149

149

anomalous subsequences

150

150



## Clustering of anomalies

subsequences grouped
according to their shape

151

151

## Conclusions

- data series is a very common data type
  - across several different domains and applications
- complex data series analytics are challenging
  - have very high complexity
  - efficiency comes from data series management/indexing techniques
- current approaches used in practice are ad-hoc
  - waste of time and effort
  - suboptimal solutions
- need for Sequence Management System
  - optimize operations based on data/hardware characteristics
  - transparent to user
- several exciting research opportunities

Themis Palpanas - University of Tokyo - Jan 2020

152

# collaborations!

Themis Palpanas - University of Tokyo - Jan 2020

153

## Data-Intensive and Knowledge-Oriented systems



155

thank you!

google: **Themis Palpanas**
visit: http://**nestordb.com**

160

160