# Scalable Vector Analytics
## A Story of Twists and Turns

Themis Palpanas

*Université Paris Cité*
*French University Institute*
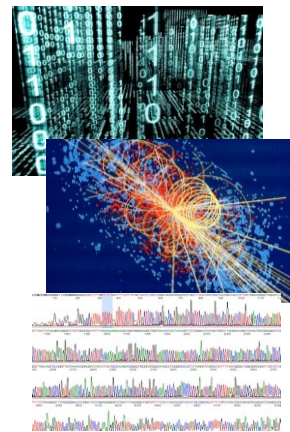
LIPADE
Laboratoire d'Informatique PAris DEscartes

Extraction et Gestion des Connaissances (EGC) – Strasbourg (France), January 2025

diNo
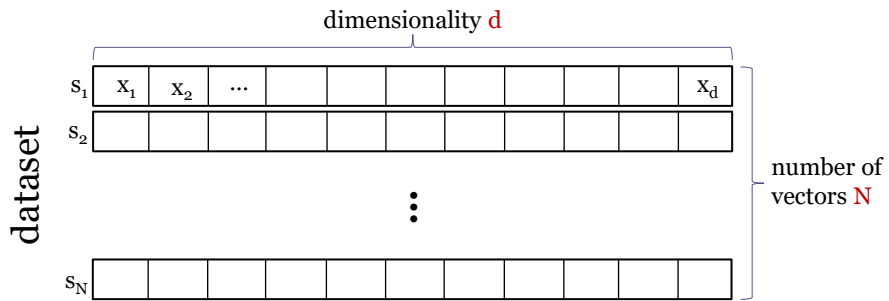
---

diNo  2

## In a Nutshell



- data collected at unprecedented rates

- they enable data-driven scientific discovery

- lots of these data are high-d vectors
  - takes **days-weeks** to analyze big high-d vector collections

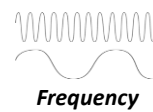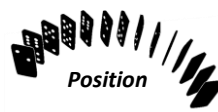goal: analyze big high-d vectors in **seconds**

# Vector Collections

- represented as *N d*-dimensional vectors

dimensionality d

dataset

$s_1$ | $x_1$ | $x_2$ | ... | | | | | | | | $x_d$

$s_2$

⋮

$s_N$

number of vectors N

# Data Series

- Sequence of points ordered along some dimension

value

$x_1$

$x_2$

$x_n$

sequence dimension

Time    Angle    Position    Mass    Frequency

# Data Centers

- cloud utilization/operation/health monitoring
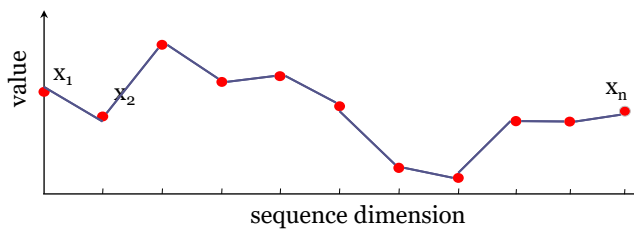


Themis Palpanas - EGC - Jan 2025

*Time*

# Neuroscience

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



*Time*

s - EGC - Jan 2025

# Remote Sensing

- Earth monitoring



*Time*

# Astrophysics



*Angle*

Schinnerer et al.

Themis Palpanas - EGC - Jan 2025

# Medicine



**Mass**

**Frequency**

Themis Palpanas - EGC - Jan 2025

# 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

Themis Palpanas - EGC - Jan 2025

5

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

- a solved(?) problem
  - your favorite DBMS
  - …
  - InfluxData
  - kx
  - Riak TS
  - OpenTSDB
  - Gorilla/Beringei
  - TimescaleDB
  - KairosDB
  - Druid
  - …

Themis Palpanas - EGC - Jan 2025

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



Clustering

Outlier Detection

Classification

Similarity Search

Frequent Pattern Mining

Themis Palpanas - EGC - Jan 2025

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

query

similar sequences

sequence
collection

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

Euclidean

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

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

Dynamic Time
Warping (DTW)

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

09-Feb-25

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

Euclidean

Dynamic Time
Warping (DTW)

Themis Palpanas - EGC - Jan 2025

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

Euclidean

equivalence to
- Pearson's Correlation
- Cosine Similarity
- Maximum Inner Product Similarity

Themis Palpanas - EGC - Jan 2025

09-Feb-25

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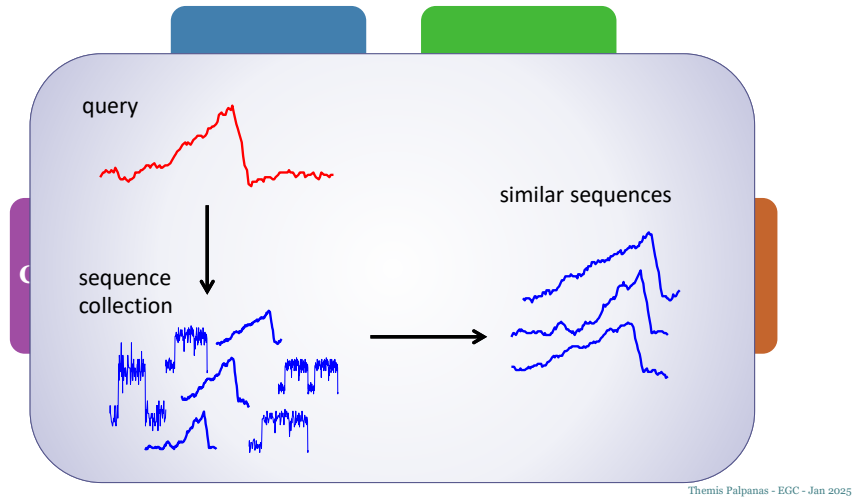## - complex analytics



Themis Palpanas - EGC - Jan 2025

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



Themis Palpanas - EGC - Jan 2025

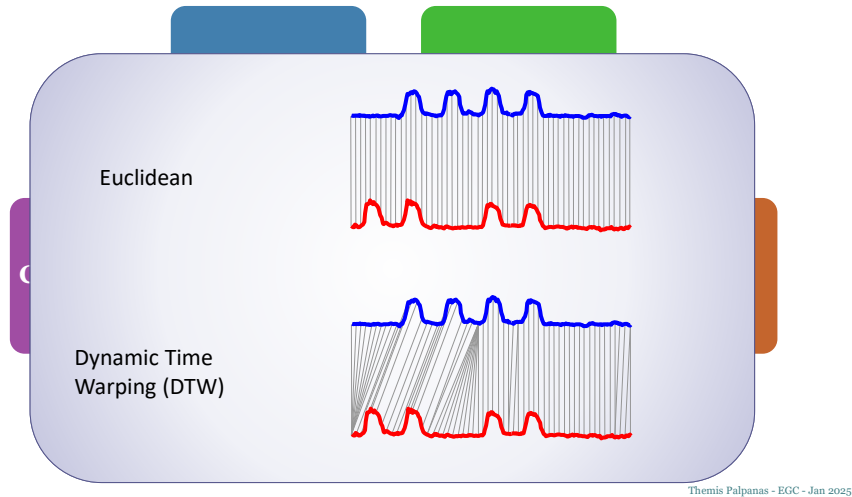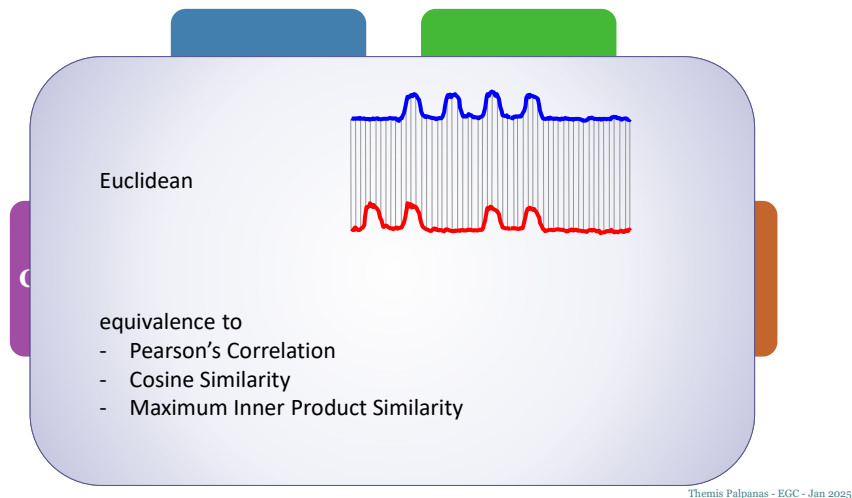# What do we want to do with them?
# - complex analytics

**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 - EGC - Jan 2025

# Real Use-Cases

astrophysics: gravitational waves, **TB/hour**
partner: **European Gravitational Observatory (EGO)**
*Pisa, Italy*

seismology: seismic sequences, **100s of TB**
partner: **Atomic Energy Commission (CEA)**
*Paris, France*

neuroscience: intracranial EEG sequences, **TB/patient**
partner: **Paris Brain Institute (ICM)**
*Paris, France*

engineering: operation monitoring, **TB-PB**
partners: **Airbus / Électricité de France (EDF)**
*Toulouse / Paris, France*

Themis Palpanas - EGC - Jan 2025

# Nearest Neighbor (NN) Queries...

Publications

PVLDB'19

**Prob($d_ε <= d_x (1+ε)$) = 1**

result within **(1+ ε) of exact NN**
with **probability 1**

ε-approximate neighbors

$O_ε$

$d_x (1+ε)$  $d_ε$  $d_x$  exact NN  $O_x$

$O_Q$

ng-approximate neighbors

**Prob($d_{ng}$ <>= ?) = ?**

result within **? of exact NN**

$d_{ng}$  $O_{ng}$

$d_{δε}$  $(1+ε) r_δ(O_Q)$

$O_{δε}$
**δ-ε-approximate neighbor**

**Prob( $d_x = min\{d_i\}$ ) = 1**

result is exact NN

**Prob($d_ε <= d_x (1+ε)$) >= δ**

result within **(1+ ε) of exact NN**
with **probability at least δ**

Themis Palpanas - EGC - Jan 2025

21

# Similarity Search via
# Serial Scan



Themis Palpanas - EGC - Jan 2025

22

## Similarity Search via
## **Serial Scan**

## Similarity Search via
## **Serial Scan**

## Similarity Search via
# Indexing

## Similarity Search via
# Indexing

## Similarity Search via
# Indexing

## Similarity Search via
# Indexing

# Query answering process



**data-to-query** time        **query answering** time

*these times are big!*

Themis Palpanas - EGC - Jan 2025

29

---

diNo 30

# High-dimensional Indexes

**KD-Tree**      **R-Tree**      **TV-Tree**



Yu Kai Him Otto, https://medium.com/

Publications

Bentley-CACM'75

Publications

Lin et al-VLDBJ'94

Publications

Guttman-SIGMOD'84

Themis Palpanas - EGC - Jan 2025

# High-dimensional Indexes

**X-Tree**



Normal Directory Nodes  Supernodes  ○ Data Nodes

Publications

Berchtold et al - VLDB'96

**M-Tree**



Publications

Ciaccia et al- VLDB'97

**Pyramid technique**



Publications

Berchtold et al - SIGMOD'98

Themis Palpanas - EGC - Jan 2025

---

# High-dimensional Indexes in a new Era

- their world
  - focused on exact query answering
  - used relatively small dataset sizes (hundreds of thousand) and dimensionalities (few dozen)
  - tested for curse of dimensionality on uniform datasets(!)

- new world
  - looking for sublinear scalability performance on 1000x larger datasets with 100x more dimensions
  - some of these indexes (R-Trees, M-Trees) used for data series with less than impressive results
  - time series community proposed new indexes

Themis Palpanas - EGC - Jan 2025

# 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!*

Themis Palpanas - EGC - Jan 2025

33

---

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

  - ▫ **(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 - EGC - Jan 2025

# *i*SAX Index

- non-balanced tree-based index with non-overlapping regions, and controlled fan-out rate
  - base cardinality $b$ (optional), segments $w$, threshold $th$
  - hierarchically subdivides SAX space until num. entries ≤ $th$

Themis Palpanas - EGC - Jan 2025

---

# *i*SAX Index

- non-balanced tree-based index with non-overlapping regions, and controlled fan-out rate
  - base cardinality $b$ (optional), segments $w$, threshold $th$
  - hierarchically subdivides SAX space until num. entries ≤ $th$

- Approximate Search
  - Match *i*SAX representation at each level

- Exact Search
  - Leverage approximate search
  - Prune search space
    - Lower bounding distance

Themis Palpanas - EGC - Jan 2025

# iSAX Index Family Lineage Tree

**Publications**

**Palpanas-ISIP'19**

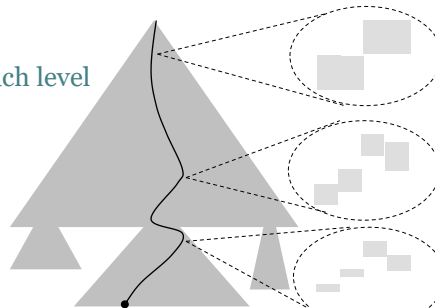| timeline | 2008 | 2010 | 2014 | 2015 | 2017 | 2018 | 2019 | 2020 | 2022 | |
|---|---|---|---|---|---|---|---|---|---|---|
| **basic index** | iSAX | | | | | | | | | C# |
| **+ Bulk Loading** | | iSAX 2.0 | iSAX2+ | | | | iSAX2+* | | | C#, C |
| **+ Adaptive** | | | ADS / ADS+ | ADSFull | | | ADS+* | Dumpy | | C |
| **+ Distributed** | | | | DPiSAX *(Spark)* | | | *(MPI)* Odyssey | | | Java C |
| **+ Multi-Core, Multi-Socket, SIMD** | | | | ParIS | ParIS+ | MESSI | Hercules | | | C |
| **+ Graphics Processing Units (GPUs)** | | | | | | SING | | | | C |
| **+ Sortable Summarizations, Streaming Data Series** | | | | Coconut-Trie / Coconut-Tree | Coconut-LSM | | | | | C |
| **+ Variable-Length Queries** | | | | ULISSE | | | | | | C |
| **+ Graph-based Similarity Search** | | | | | | | Elpis | | | C |

Themis Palpanas - EGC - Jan 2025

# iSAX Index Family Lineage Tree

**Publications**

**Palpanas-ISIP'19**

| timeline | 2008 | 2010 | 2014 | 2015 | 2017 | 2018 | 2019 | 2020 | 2022 | |
|---|---|---|---|---|---|---|---|---|---|---|
| **basic index** | iSAX | **600 sec** | | | | | | | | C# |
| **+ Bulk Loading** | | | iSAX2+ | | | | iSAX2+* | | | C#, C |
| **+ Adaptive** | | | ADS / ADS | **50 sec** Full | | | ADS+* | Dumpy | | C |
| **+ Distributed** | | | | DPiSAX *(Spark)* | **15 sec** | | *(MPI)* Odyssey | | | Java C |
| **+ Multi-Core, Multi-Socket, SIMD** | | | | | S+ | MESSI | Hercules | | | C |
| **+ Graphics Processing Units (GPUs)** | | | | | | SING | | | | C |
| **+ Sortable Summarizations, Streaming Data Series** | | | | Coconut-Trie / Coconut-Tree | Coconut-LSM | | | | | C |
| **+ Variable-Length Queries** | | | | ULISSE | | | | | | C |
| **+ Graph-based Similarity Search** | | | | | | | Elpis | | | C |

execution time for **1 similarity search query on a 100GB dataset _on disk_**

Themis Palpanas - EGC - Jan 2025

09-Feb-25

**Slide 40**

# iSAX Index Family Lineage Tree

Publications: Palpanas-ISIP'19



timeline: 2008 2010 2014 2015 2017 2018 2019 2020 2022

- basic index — iSAX — **600 sec** — C#
- + Bulk Loading — iSAX2+ — iSAX2+* — C#, C
- + Adaptive — ADS / ADS...Full — **50 sec** — ADS+* — Dumpy — C
- + Distributed — DPiSAX (Spark) — **15 sec** — (MPI) Odyssey — Java, C
- + Multi-Core, Multi-Socket, SIMD — ...S+ — M... les — **0.035 sec** — C
- + Graphics Processing Units (GPUs) — C
- + Sortable Summarizations, Streaming Data Series — Coconut-Trie / Coconut-Tree — Coconut-LSM — C
- + Variable-Length Queries — ULISSE — C
- + Graph-based Similarity Search — Elpis

execution time for **1 similarity search query on a 100GB dataset _in memory_**

Themis Palpanas - EGC - Jan 2025

---

**Slide 41**

# Further Advances

Publications: Gogolou - BigVis'19, Gogolou - SIGMOD'20, Gogolou - SIGMOD'20



Average Times of 100 queries (in sec) — 1-NN Time, Total Time; seismic (100M, 256 p.), SALD (200M, 128 p.), deep1B (267M, 96 p.)

- how do we further reduce the wasted (gray) effort?
  - **progressive query answering**
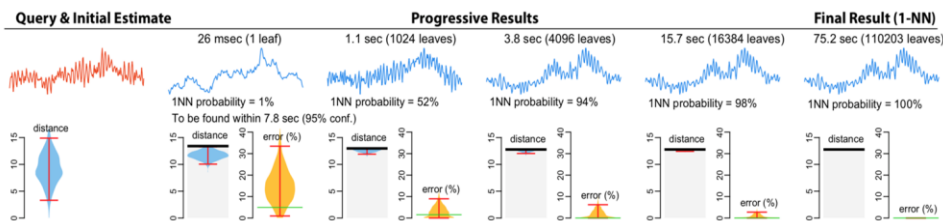    - produce intermediate answers with (probabilistic) quality guarantees

Themis Palpanas - EGC - Jan 2025

20

# Further Advances:
# Progressive Query Answering

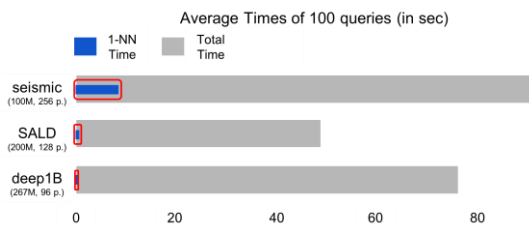- interaction with users offers new opportunities
  - ▫ progressive answers
    - · produce intermediate results
      - · iteratively converge to final, correct solution
    - · provide bounds on the errors (of the intermediate results) along the way



Themis Palpanas - EGC - Jan 2025

---

# Further Advances

- how do we further reduce the wasted (gray) effort?
  - ▫ **progressive query answering**
    - · produce intermediate answers with (probabilistic) quality guarantees
  - ▫ **learned summarizations + index structures**
    - · adapt to data characteristics
    - · build more efficient indexes
    - · perform more effective pruning
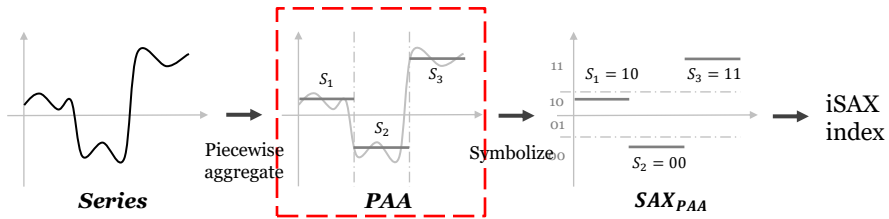
Themis Palpanas - EGC - Jan 2025

# Further Advances: Learning

- Series Approximation Network (SEAnet)
  - ▫ novel autoencoder architecture
  - ▫ learns deep embedding approximations
  - ▫ uses those for similarity search

# Further Advances: Learning

- Series Approximation Network (SEAnet)
  - ▫ novel autoencoder architecture
  - ▫ learns deep embedding approximations
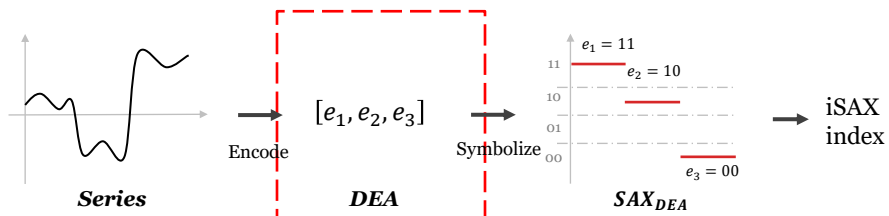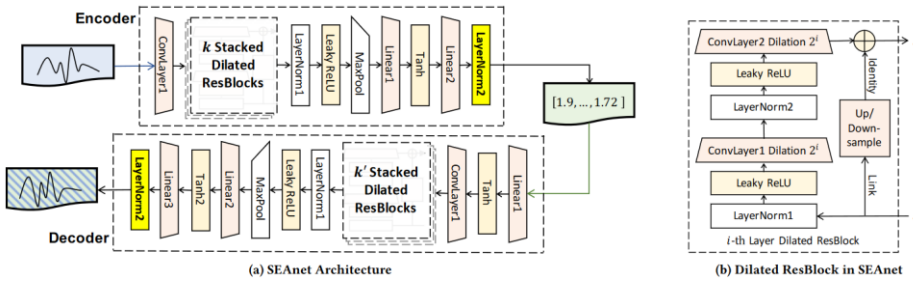  - ▫ uses those for similarity search

# Further Advances: Learning

- Series Approximation Network (SEAnet)
  - ▫ is an exponentially dilated ResNet architecture + Sum of Squares regularization
  - ▫ minimizes
    - · reconstruction error
    - · difference between distance of two vectors in embedded space and distance in original space



(a) SEAnet Architecture

(b) Dilated ResBlock in SEAnet

Themis Palpanas - EGC - Jan 2025

# Further Advances: Learning

- Series Approximation Network (SEAnet)
  - ▫ is an exponentially dilated ResNet architecture + Sum of Squares regularization
  - ▫ minimizes
    - · reconstruction error
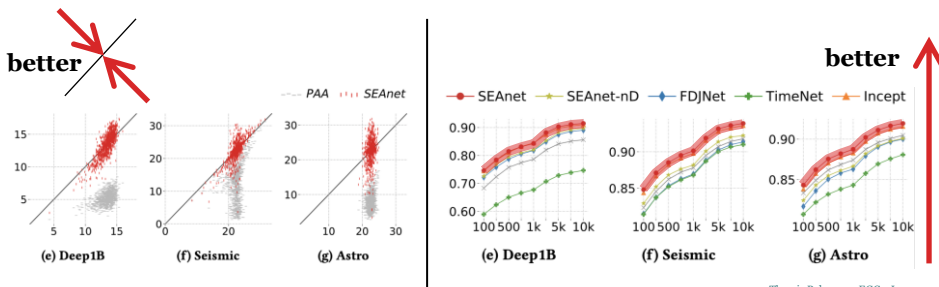    - · difference between distance of two vectors in embedded space and distance in original space



Themis Palpanas - EGC - Jan 2025

# Further Advances: Learning

- Learned Filters (LeaFi)
  - machine learning models that make pruning decisions
  - applied when pruning based on lower bounding is not possible



up to 20x more pruning
up to 32x faster query answering

# High-d Vectors
# Beyond Data Series

- two sides of the same(?) coin
  - data series as multidimensional points
  - for a specific ordering of the dimensions
- everything we discussed applicable to high-d vectors, too!

# High-d Vectors
# Beyond Data Series

- two sides of the same(?) coin
  - data series as multidimensional points
  - for a specific ordering of the dimensions
- everything we discussed applicable to high-d vectors, too!

**sequences**
**text**
**images**
**video**
**graphs**
**molecules**
**...**

Themis Palpanas - EGC - Jan 2025
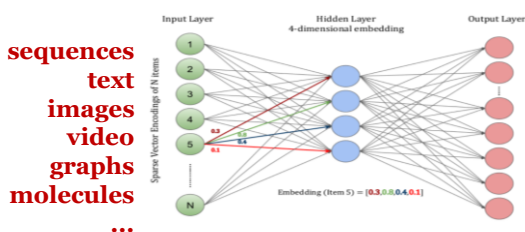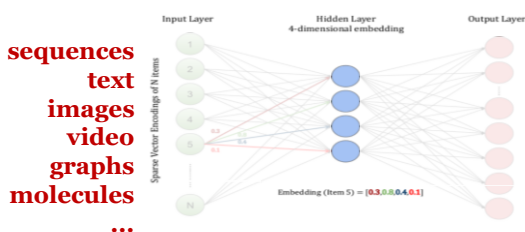
---

# High-d Vectors
# Beyond Data Series

- two sides of the same(?) coin
  - data series as multidimensional points
  - for a specific ordering of the dimensions
- everything we discussed applicable to high-d vectors, too!

**sequences**
**text**
**images**
**video**
**graphs**
**molecules**
**...**

**deep embeddings**
high-d vectors learned using a DNN

Themis Palpanas - EGC - Jan 2025

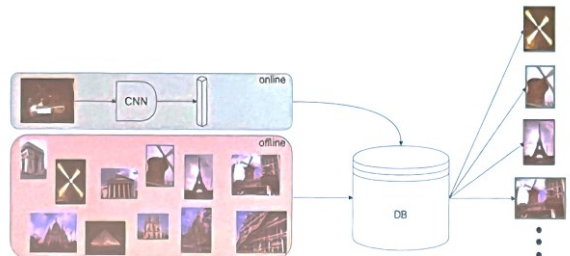# Deep Embeddings
# Similarity Search Applications

- image retrieval

Image Retrieval: the task

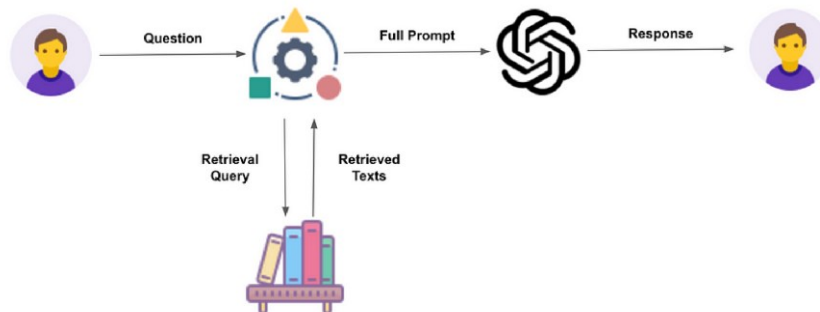Given a query image, rank images of a database from most to least similar.

Ramzi et al., Cap&RFIAP'22

Themis Palpanas - EGC - Jan 2025

# Deep Embeddings
# Similarity Search Applications

- image retrieval
- retrieval augmented generation (RAG)

https://pvhil.com/

Themis Palpanas - EGC - Jan 2025

54

# Deep Embeddings
# Similarity Search Applications

- image retrieval
- retrieval augmented generation (RAG)
- recommendations
- entity matching
- fraud detection
- drug discovery
- ...

Ramzi et al., Cap&RFIAP'22

Themis Palpanas - EGC - Jan 2025

---

55

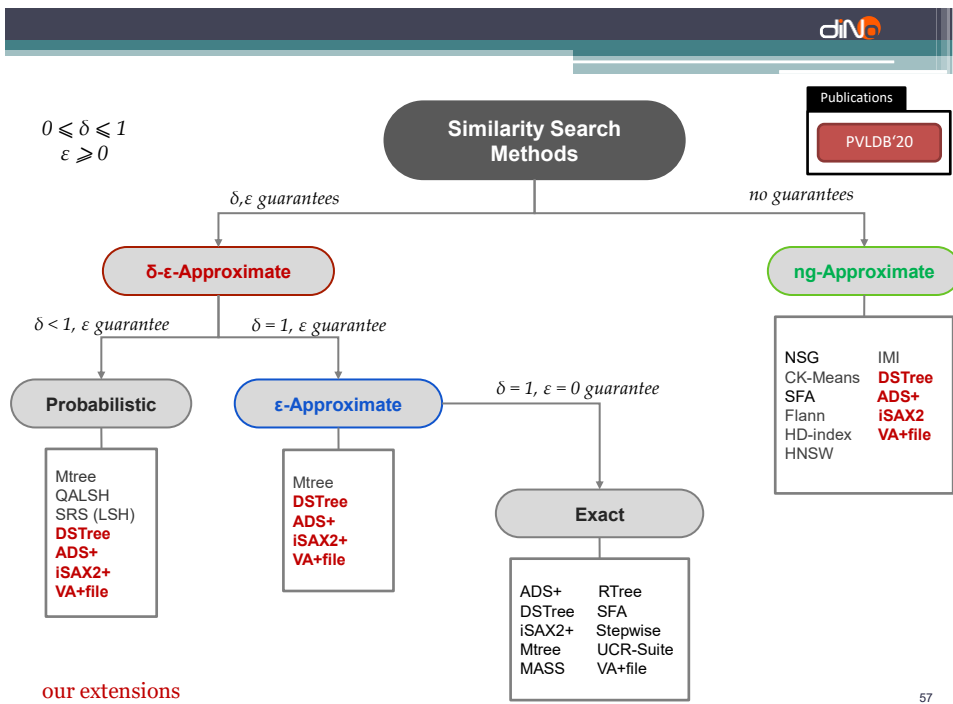# Data Series Indexes in a new Era

- their world
  - focused on **exact** query answering
  - centered discussion around data series **shapes/patterns**

- new world
  - looking for **ultra-fast** performance for applications that tolerate **approximate** answers
  - machine learning and related communities **proposed new** indexes

Themis Palpanas - EGC - Jan 2025

# High-d Vectors Indexes

- techniques for approximate similarity search in high-d vectors
  - [LSH (SRS)]
  - space quantization and inverted files (IMI)
  - k-NN graphs (HNSW)

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

Themis Palpanas - EGC - Jan 2025

---

$0 \leqslant \delta \leqslant 1$
$\varepsilon \geqslant 0$

**Similarity Search Methods**

*δ,ε guarantees* · · · · · · · · · · · *no guarantees*

**δ-ε-Approximate** · · · · · · · · · · · **ng-Approximate**

*δ < 1, ε guarantee* · · · *δ = 1, ε guarantee*

**Probabilistic** · · · · · · **ε-Approximate** · · · *δ = 1, ε = 0 guarantee*

| | |
|---|---|
| NSG | IMI |
| CK-Means | **DSTree** |
| SFA | **ADS+** |
| Flann | **iSAX2** |
| HD-index | **VA+file** |
| HNSW | |

| |
|---|
| Mtree |
| QALSH |
| SRS (LSH) |
| **DSTree** |
| **ADS+** |
| **iSAX2+** |
| **VA+file** |

| |
|---|
| Mtree |
| **DSTree** |
| **ADS+** |
| **iSAX2+** |
| **VA+file** |

**Exact**

| | |
|---|---|
| ADS+ | RTree |
| DSTree | SFA |
| iSAX2+ | Stepwise |
| Mtree | UCR-Suite |
| MASS | VA+file |

our extensions

Themis Palpanas - EGC - Jan 2025

57

# Data Series vs. high-d Vectors

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

---

# Data Series vs. high-d Vectors

- **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(δε)**

DSTree — HNSW — IMI — iSAX2+ — SRS — VA+file

# Data Series vs. high-d Vectors

Publications

PVLDB'20

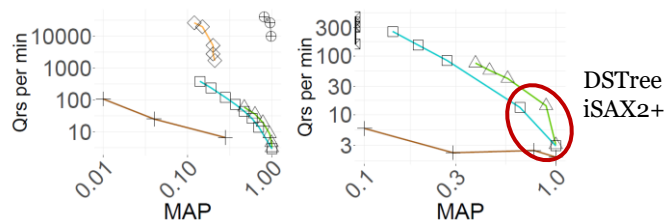- **data series techniques** are the **overall winners**, even on **general high-d vector** data
  - ▫ perform the best for approximate queries with probabilistic guarantees (δ-ε-approximate search), in-memory and on-disk
  - ▫ perform the best for long vectors, in-memory and on-disk



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

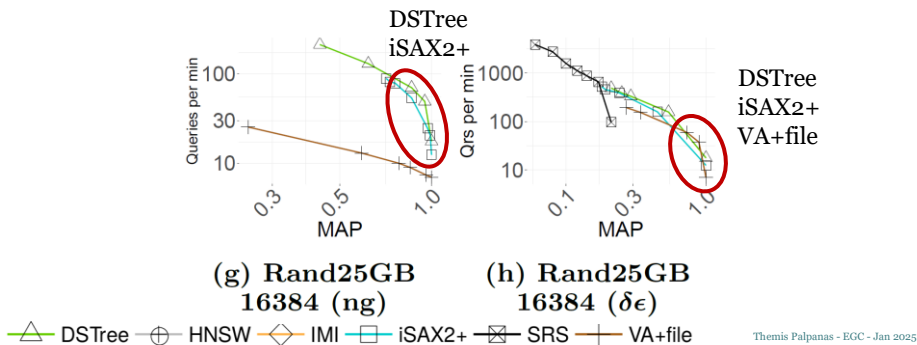—△— DSTree —⊕— HNSW —◇— IMI —⊟— iSAX2+ —⊠— SRS —+— VA+file
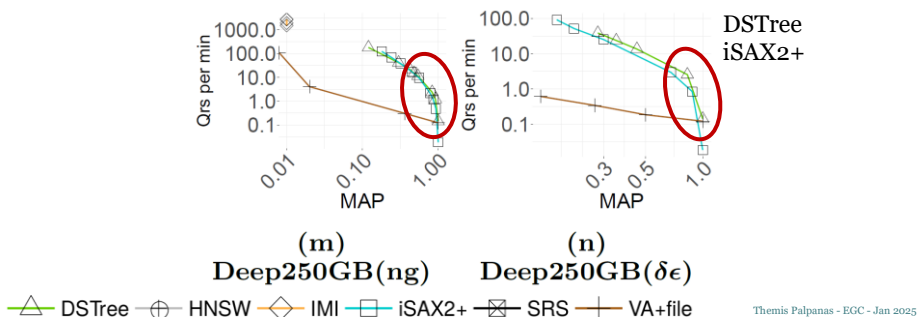
Themis Palpanas - EGC - Jan 2025

# 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
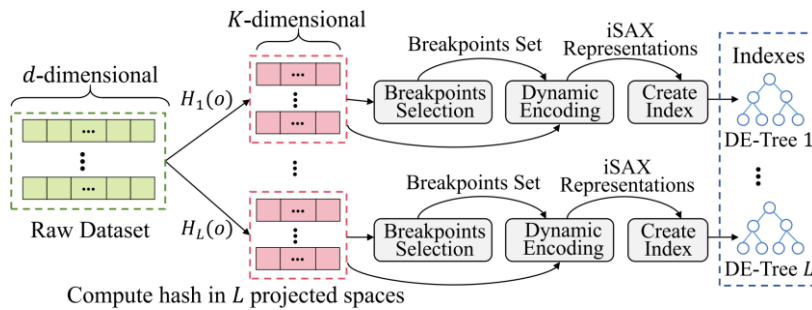  - ▫ perform the best for disk-resident vectors



(m) Deep250GB(ng)   (n) Deep250GB(δε)

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

Themis Palpanas - EGC - Jan 2025

# Hybrid (iSAX + LSH): DET-LSH

- **DET-LSH** combines tree and LSH for efficient indexing and approximate search with probabilistic guarantees

# Hybrid (iSAX + LSH): DET-LSH

- **DET-LSH** combines tree and LSH for efficient indexing and approximate search with probabilistic guarantees

# Hybrid (iSAX + LSH): DET-LSH

- **DET-LSH** combines tree and LSH for efficient indexing and approximate search with probabilistic guarantees



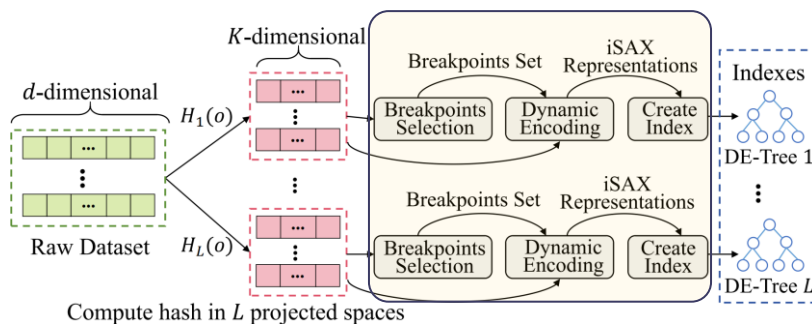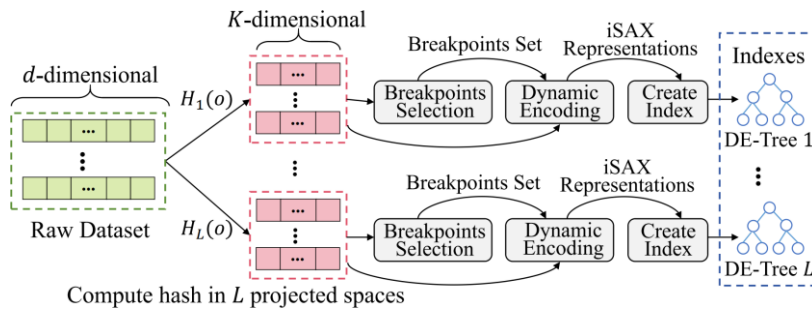up to <span style="color:red">6x faster indexing</span> and <span style="color:red">2x faster query answering</span> (than standard LSH methods)
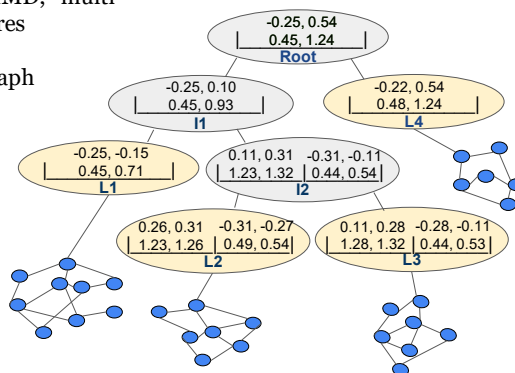
Themis Palpanas - EGC - Jan 2025

---

# Hybrid (DSTree + HNSW): ELPIS
## Parallel, In-Memory Indexing of Sequences

- ❑ In-memory solution for SIMD, multi-core, multi-socket architectures

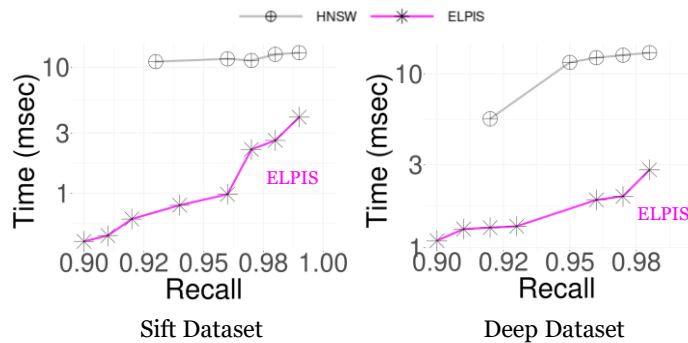- ❑ **ELPIS** combines tree and graph structures for efficient in-memory ng-approximate vector similarity search.



Themis Palpanas - EGC - Jan 2025

# Hybrid (DSTree + HNSW): ELPIS

Publications

Azizi-PVLDB'23

## Parallel, In-Memory Indexing of Sequences

○ Query Performance on 1B vectors datasets (Sift, Deep)



Sift Dataset                    Deep Dataset

ELPIS answers **10-NN queries** in **~3 msec** for a dataset of **1 billion vectors** with **recall 0.99**

Themis Palpanas - EGC - Jan 2025

# Conclusions

- high-d vectors is a very common data type
  - across several different domains and applications
- complex high-d vector analytics are challenging
  - have very high inherent complexity
- data series management/indexing techniques provide state-of-the-art performance
  - work for data series and general high-d vectors (and embeddings)
  - lead to fast complex analytics and machine learning

- several exciting research opportunities
  - distributed solutions
  - progressive analytics
  - learned (data-adaptive) summarizations/data structures

Themis Palpanas - EGC - Jan 2025

## Some more open questions…

- what are the theoretical properties of existing solutions?
  - best/expected/worst time performance
  - best/expected/worst accuracy (for approximate query answering)

- can we predict performance based on data characteristics?
  - analytical results (eg, based on different distributions)
  - learned

- how do we integrate vector similarity search in databases?
  - combine similarity search with other predicates
  - optimization

- what are the right benchmarks to evaluate high-d vector indices?
  - data and query workloads (currently: most queries are easy)
  - evaluation measures (currently: time and recall)

Themis Palpanas - EGC - Jan 2025

## Going forward

- high-d vector similarity search relevant to many communities
  - data management
  - time series
  - information retrieval
  - text search
  - machine learning
  - deep learning
  - parallel and distributed computing

Themis Palpanas - EGC - Jan 2025

# Going forward

- high-d vector similarity search <span style="color:red">relevant to many</span> communities
  - data management
  - time series
  - information retrieval
  - text search
  - machine learning
  - deep learning
  - parallel and distributed computing

- research on this problem <span style="color:red">fragmented</span> across communities
  - open communication channels among these communities
  - initiate discussions

- start <span style="color:red">collaborations</span>!

## Data-Intensive and Knowledge-Oriented systems

# thank you!

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

work supported by:

72