

Evolution of similarity search in databases

... from trivial applications to real problems

Tomáš Skopal

Siret Research Group, KSI MFF UK

3/2012



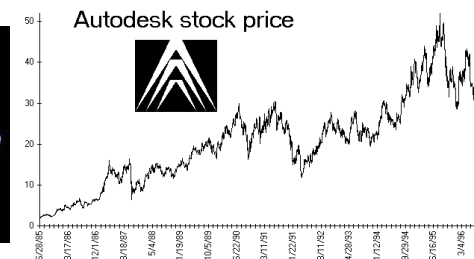
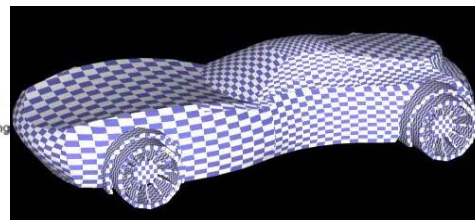
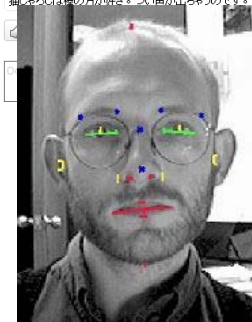
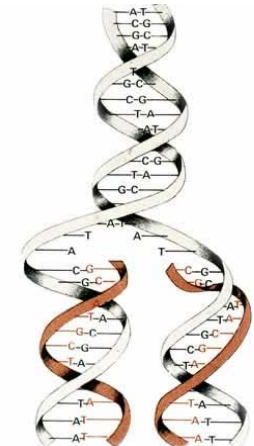
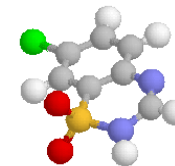
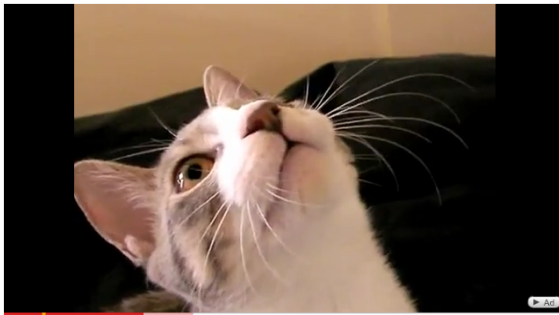
Talk outline

- the problem of similarity search
- similarity models
 - early models (vector space)
 - current models (metric space)
 - future models (non-metric models)
- applications
 - data-centric (multimedia search engines)
 - service-centric (“enterprise” applications)

Similarity search

- task
 - how to search non-structured data based on content?
 - we cannot use traditional approaches (e.g., linear order, SQL)

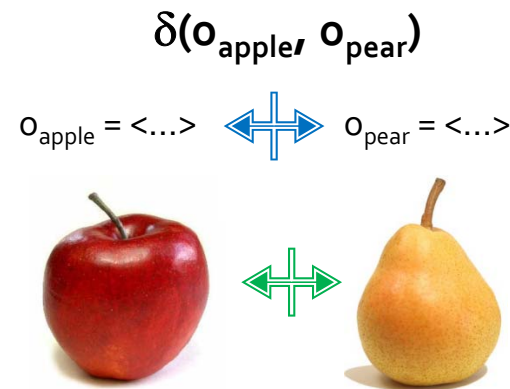
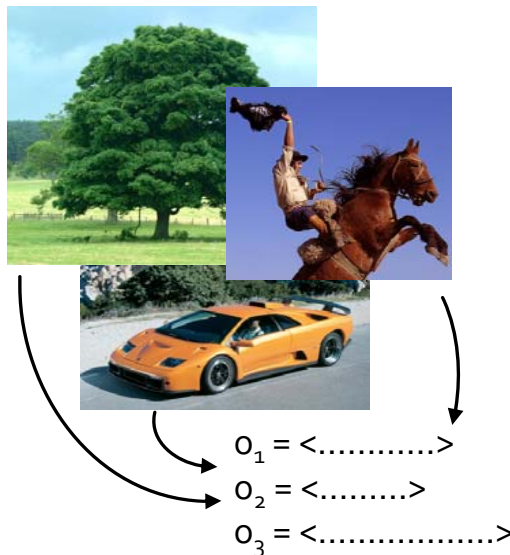
YouTube [Vyhledat](#) [Procházet](#) [Nahrát](#)
おしゃべりキャット - Talking Cat -
lowdope Počet videí: 175 [Odběr](#)



T.Skopal, Evolution of similarity search in databases

Similarity search

- task
 - how to search non-structured data based on content?
 - we cannot use traditional approaches (e.g., linear order, SQL)
 - similarity model
 - descriptor universe (feature extraction) $\mathbf{o}_i \in \mathbf{U}$
 - pair-wise similarity/distance function $\delta(\mathbf{x}, \mathbf{y})$



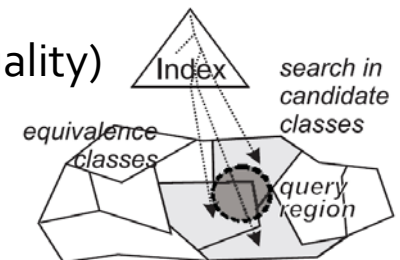
Similarity search

- task
 - how to search non-structured data based on content?
 - we cannot use traditional approaches (e.g., linear order, SQL)
 - similarity model
 - descriptor universe (feature extraction) $\mathbf{o}_i \in \mathbf{U}$
 - pair-wise similarity/distance function $\delta(\mathbf{x}, \mathbf{y})$
 - retrieval
 - querying (range query, kNN query)
 - exploration (and other modern means of retrieval)



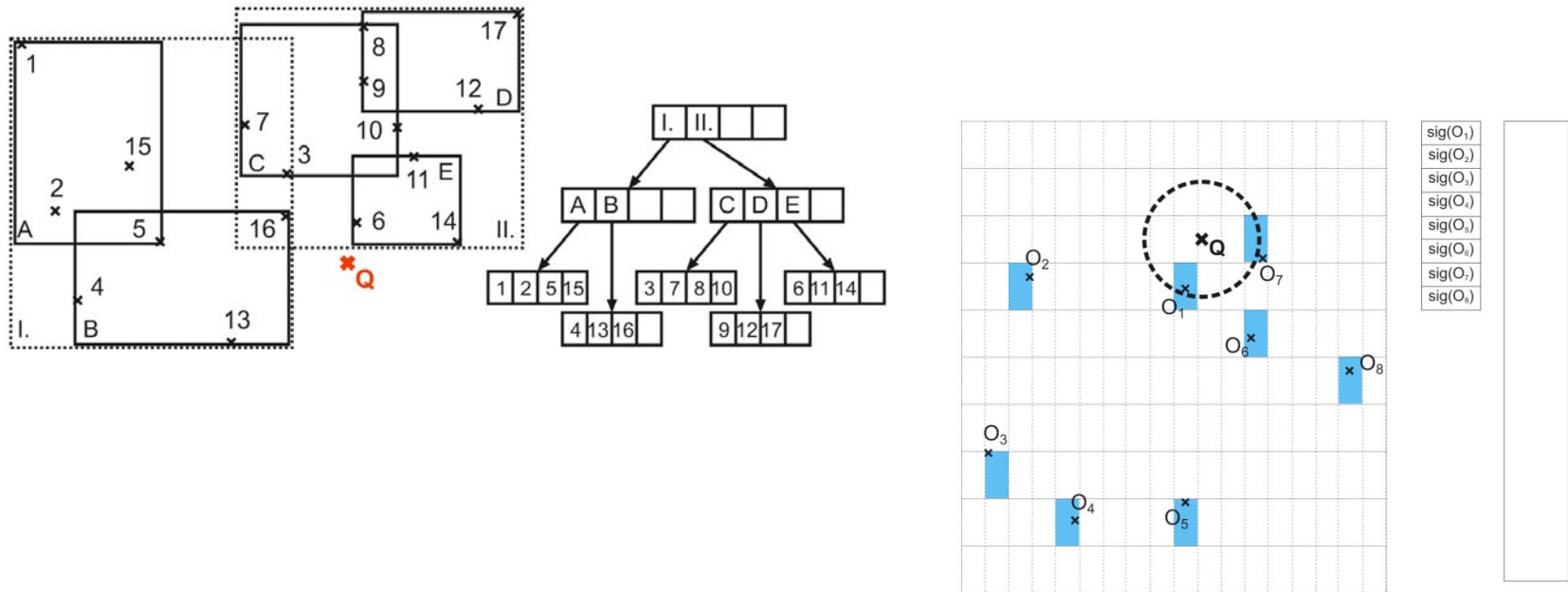
Similarity search

- task
 - how to search non-structured data based on content?
 - we cannot use traditional approaches (e.g., linear order, SQL)
 - similarity model
 - descriptor universe (feature extraction) $\mathbf{o}_i \in \mathbf{U}$
 - pair-wise similarity/distance function $\delta(\mathbf{x}, \mathbf{y})$
 - retrieval
 - querying (range query, kNN query)
 - exploration (and other modern means of retrieval)
- database solution
 - efficiency vs. effectiveness (performance vs. quality)
 - assumption: sequential scan is too slow I/O cost but also (mainly) CPU cost of $\delta(\mathbf{x}, \mathbf{y})$



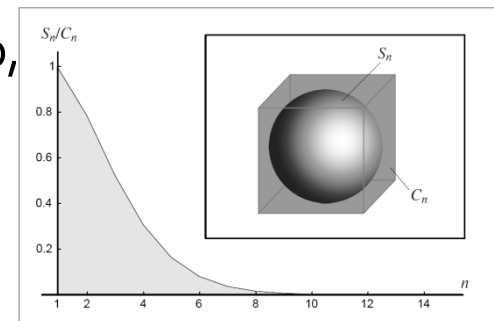
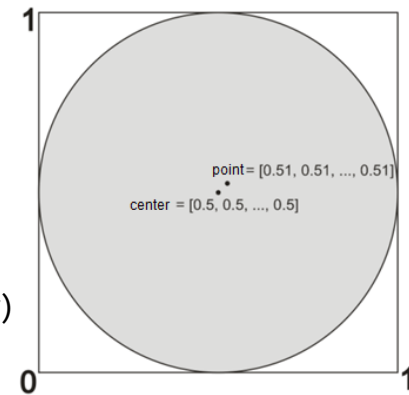
Early models – vector models

- similarity modeled in vector space
 - reuse of spatial indexes designed for GIS, CAD/CAM, etc.
 - R-tree, X-tree, SS-tree, Grid file, VA-file, inverted index, etc.



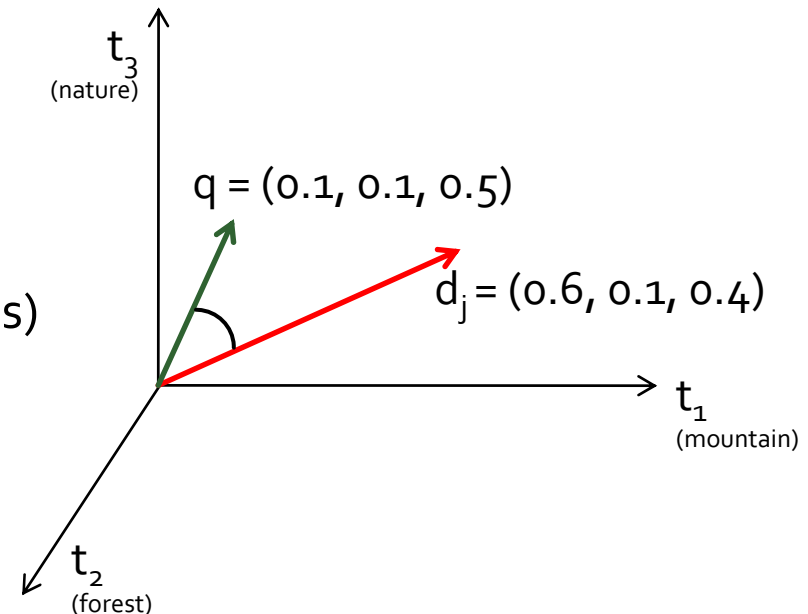
Early models – vector models

- similarity modeled in vector space
 - reuse of spatial indexes designed for GIS, CAD/CAM, etc.
 - R-tree, X-tree, SS-tree, Grid file, VA-file, inverted index, etc.
- cons
 - similarity modeling limited to L_p spaces
 - descriptors = numeric vectors, similarity = L_p metric
 - independent dimensions („smart“ descriptors, „stupid“ similarity)
 - curse of dimensionality
 - inefficient for high-dimensional data, say $\text{dim} > 10$, (not the case of original use in GIS/CAD, i.e., 2D, 3D data!)



Early models – vector models

- suitable application in similarity search
 - vector query (originated in Information retrieval)
 - inverted index
 - specific requirements
 - cosine measure/distance
 - sparse vectors (query and data)
 - applications
 - text retrieval (dimensions are terms)
 - **bag-of-words models** in multimedia retrieval



Early models – vector models

■ bag-of-words in image retrieval

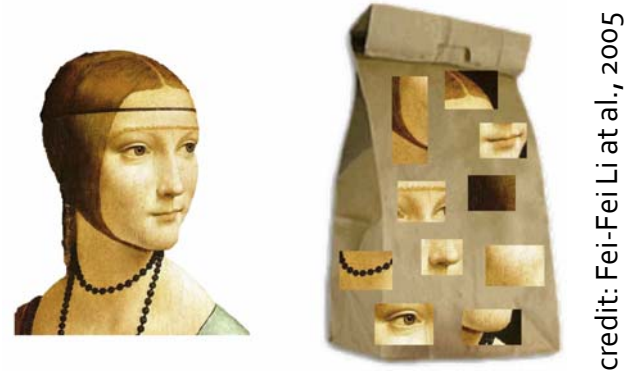
1. extract features
(from many images)



2. quantize/cluster features



T.Skopal, Evolution of similarity search in databases

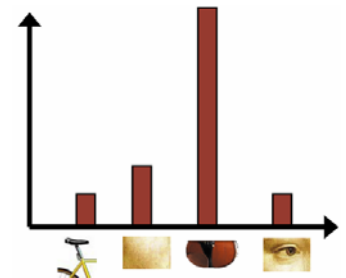


credit: Fei-Fei Li et al., 2005

3. learn visual vocabulary



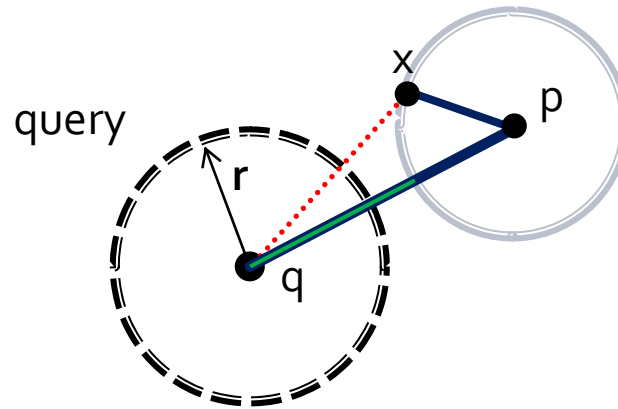
4. represent
images by
histograms
(vectors)



Metric space model

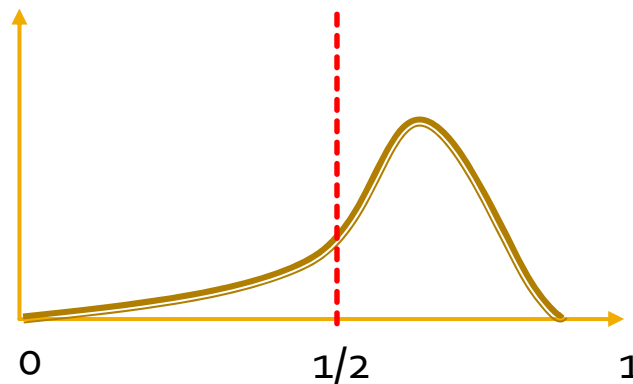
- more general than vector models
 - any metric space
 - black-box descriptor and black-box distance function δ
 - **the distance must be a metric (axioms)**
 - “lowerbounding” using pivots (triangle inequality)

$$\delta(x, y) + \delta(y, z) \geq \delta(x, z),$$



Metric space model

- more general than vector models
 - any metric space
 - black-box descriptor and black-box distance function δ
 - **the distance must be a metric (axioms)**
 - “lowerbounding” using pivots (triangle inequality)
 - curse of dimensionality reduced to some extent
 - generalized problem of **intrinsic dimensionality** (dist. distrib. matters)



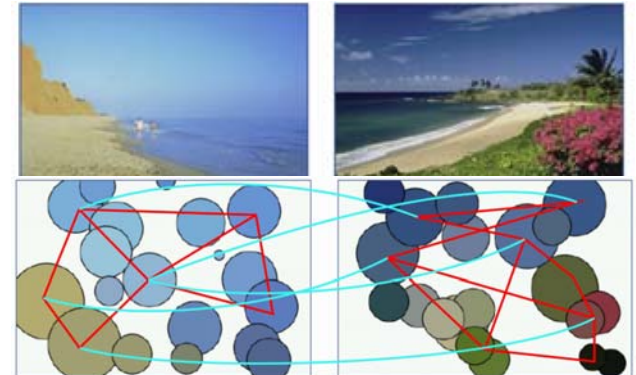
Metric space model

- more general than vector models
 - any metric space
 - black-box descriptor and black-box distance function δ
 - **the distance must be a metric (axioms)**
 - “lowerbounding” using pivots (triangle inequality)
 - curse of dimensionality reduced to some extent
 - generalized problem of **intrinsic dimensionality** (dist. distrib. matters)
 - many metric indexes proposed
 - main-memory/persistent, serial/parallel/distributed, static/dynamic, etc.
 - attempts extending SQL by similarity predicates + implementations
- for real problems often still not sufficient
 - metric axioms limit the modeling capabilities

Metric space model

- suitable application
 - image feature signatures + SQFD

$$\text{SQFD}_{f_s}(S^q, S^p) = \sqrt{(w_q \mid -w_p) \cdot A_{f_s} \cdot (w_q \mid -w_p)^T}$$



M. Kruliš, T. Skopal, J. Lokoč, Ch. Beecks. Combining CPU and GPU Architectures for Fast Similarity Search, **Distributed and Parallel Databases**, Springer, 2012

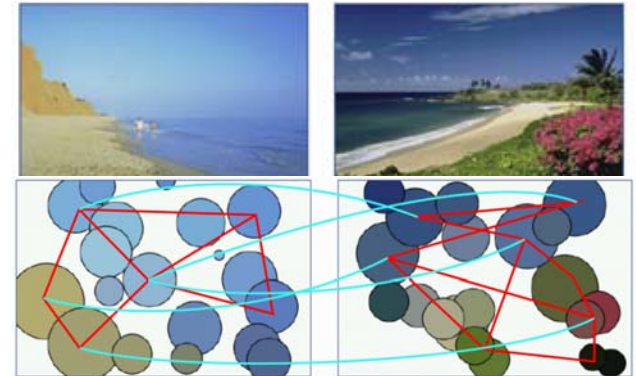
M. Kruliš, J. Lokoč, Ch. Beecks, T. Skopal, T. Seidl. Processing Signature Quadratic Form Distance on Many-Core GPU Architecture, **ACM CIKM 2011**, Glasgow, UK

Ch. Beecks, J. Lokoč, T. Seidl, T. Skopal. Indexing the Signature Quadratic Form Distance for Efficient Content-Based Multimedia Retrieval, **ACM ICMR 2011**, Trento, Italy

Metric space model

- suitable application
 - image feature signatures + SQFD

$$\text{SQFD}_{f_s}(S^q, S^p) = \sqrt{(w_q \mid -w_p) \cdot A_{f_s} \cdot (w_q \mid -w_p)^T}$$



- „mistyping vocabulary” using edit distance

$$\delta_{\text{edit}}(\text{'bank'}, \text{'bar'}) = 2 \quad \delta_{\text{edit}}(\text{'tank'}, \text{'bank'}) = 1$$

- bad application

- indexing of protein sequences using edit distance
 - edit distance too simple for biologic applications
 - extension needed like scoring matrices, gap penalties, local alignment
 - e.g., Smith-Waterman

N	P	H	G	I	I	M	G	L	A	E	
		+8	+6	+2							→ 16
		H	G	L							

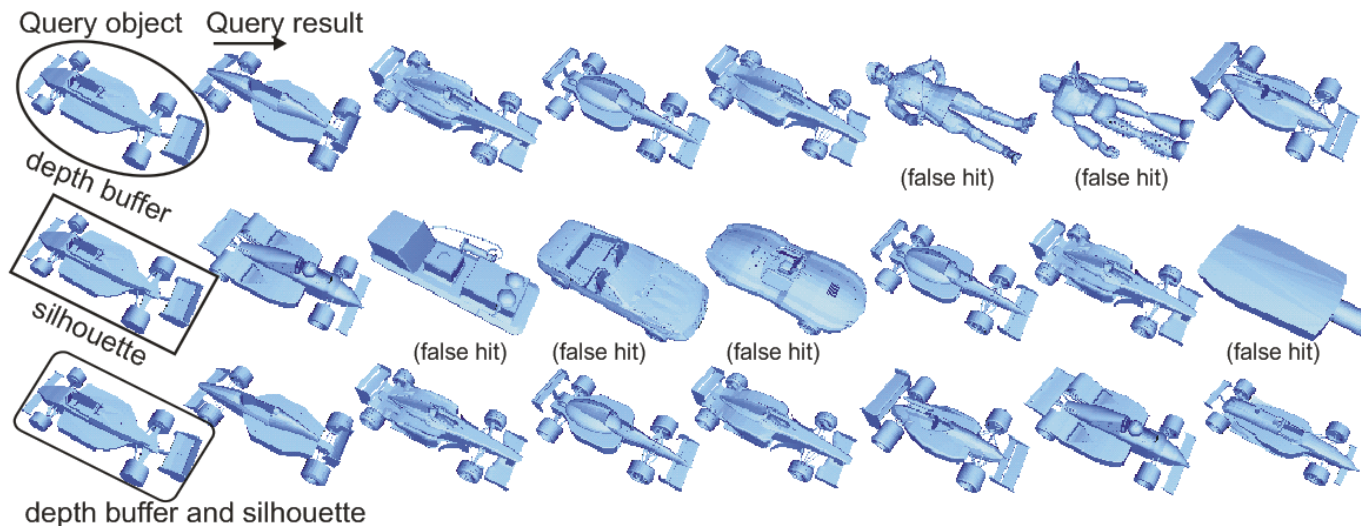
Non-metric similarity models

- absence of some metric axiom = non-metric function
- similarity parameterization
 - inclusion of user preferences/context at query time
- non-metric similarities “by design”
 - robust behavior (ignoring „noisy“ parts)
 - local behavior (favoring similarity)
 - comfort of black-box approach (domain experts outside CS)

T. Skopal, B. Bustos. On Nonmetric Similarity Search Problems in Complex Domains, **ACM Computing Surveys**, 43(4):34:1-34:50, 2011
B. Bustos, T. Skopal. Nonmetric Similarity Search Problems in Very Large Collections, **ICDE 2011**, Hannover, Germany, IEEE

Non-metric similarity models

- need to allow explicit preferences/user profile at query time
 - dynamic function $\delta(\mathbf{x}, \mathbf{y}, \text{other parameters})$
 - i.e., the same descriptors \mathbf{x}, \mathbf{y} lead to different similarity values (so it is not even a function, not yet metric distance)



Non-metric similarity models

- possible solution
 - multi-metric approach, linear combination of metrics
 - query time weighting

$$\Delta_{\mathbb{W}}(O_1, O_2) = \sum_{i=1}^m w_i \cdot \delta_i(O_1, O_2).$$

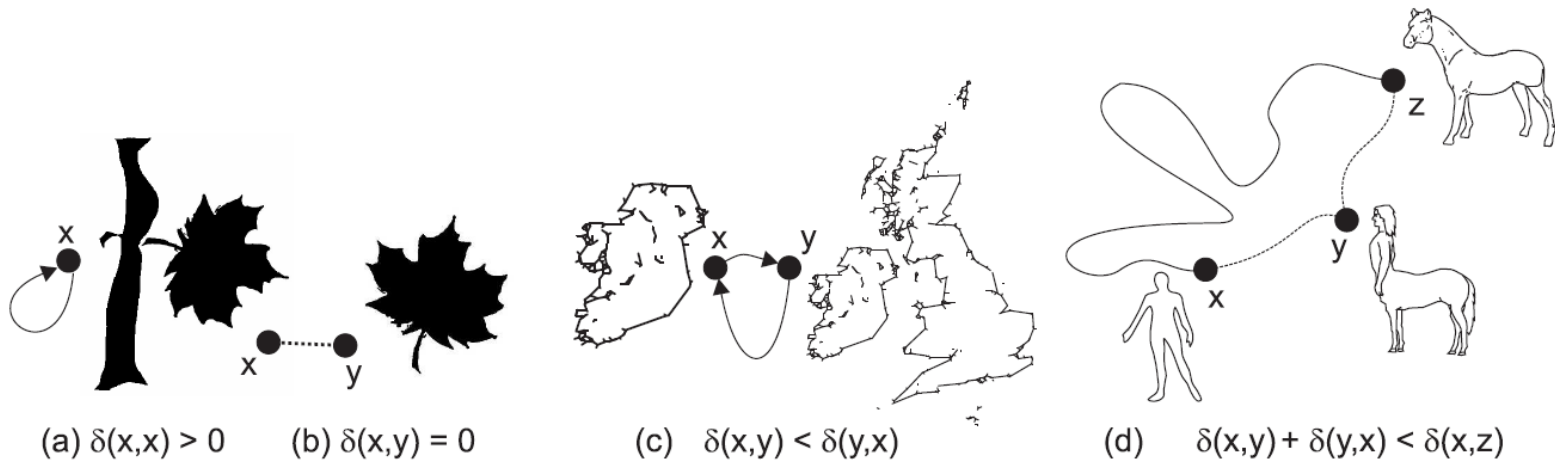
- multi-metric indexes (adaptation of metric ones, e.g., M³-tree)
 - storage of distance components,
i.e., final aggregation with weights at query time

B. Bustos, S. Kreft, T. Skopal. Adapting Metric Indexes for Searching in Multi-Metric Spaces, **Multimedia Tools and Applications**, Springer, 2012

B. Bustos, T. Skopal. Dynamic Similarity Search in Multi-Metric Spaces, **ACM MIR 2006**, Santa Barbara, CA, USA

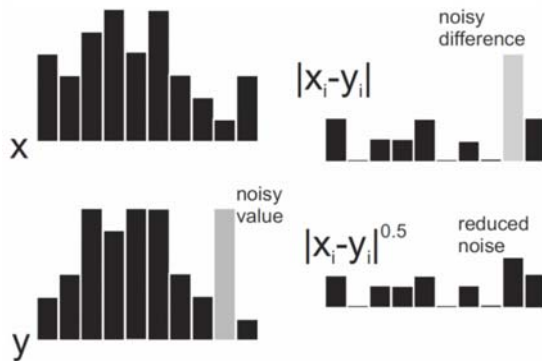
Non-metric similarity models

- arguments against metric axioms in theory (psychology, cognitive sciences)



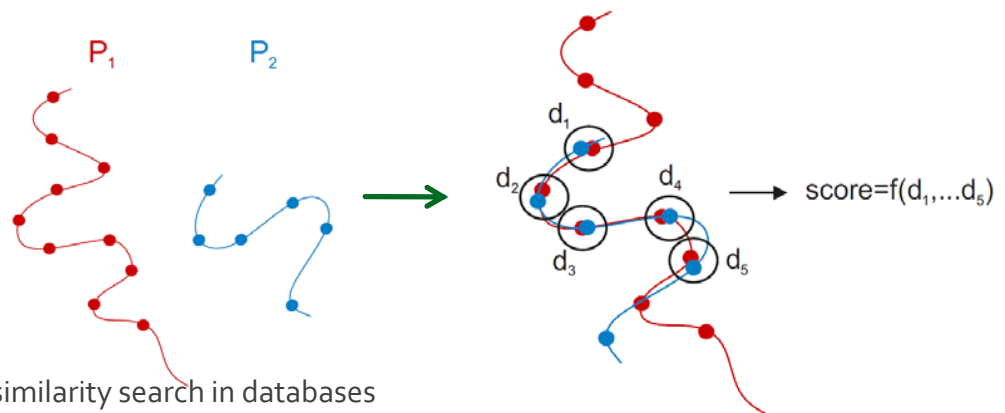
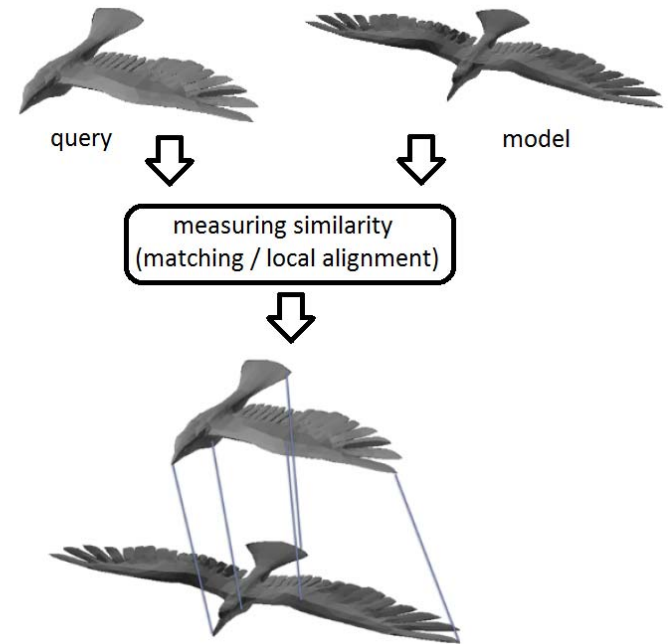
Non-metric similarity models

- and practically...



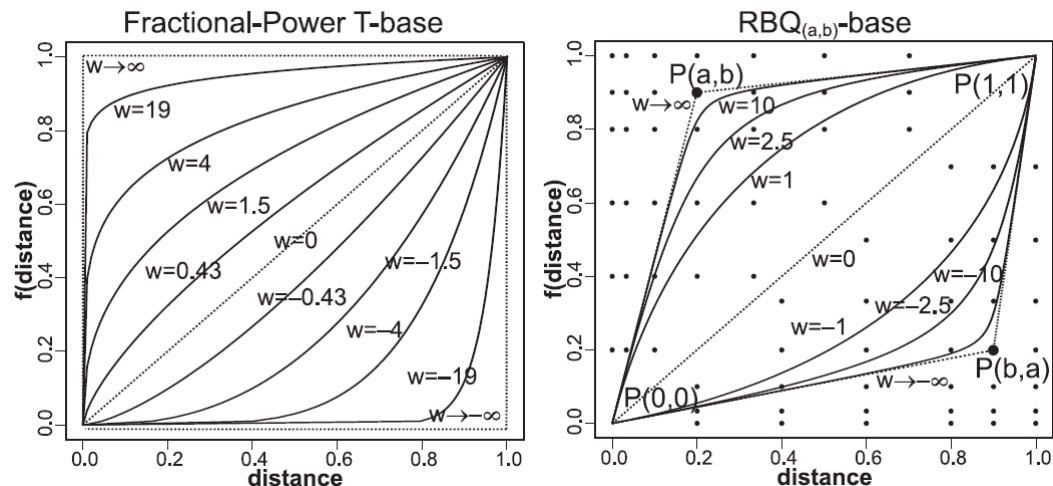
```

if (condition) {
    statements;
}
else if (condition) {
    statements;
}
else {
    statements;
}
}
if (x == 3) // curly braces not needed
    flag = 1; // when if statement is
else // followed by only one
    flag = 0; // statement
Repetition (while)
while (expression) { // loop until
    statements; // expression is false
}
Repetition (do-while)
do { // perform the statements
    statements; // as long as condition
} while (condition); // is true
Repetition (for)
init - initial value for loop control variable
condition - stay in the loop as long as condition
is true
increment - change the loop control variable
for (init; condition; increment) {
    statements;
}
Bifurcation (break, continue, goto, exit)
break; // ends a loop
continue; // stops executing statements
// in current iteration of
    
```



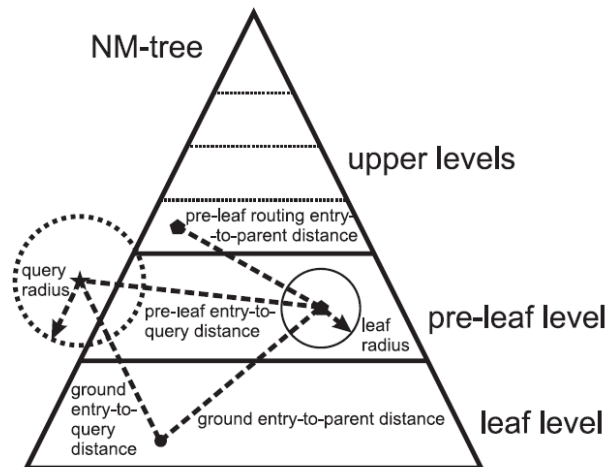
Non-metric similarity models

- general indexing of non-metric distances
 - TriGen algorithm
 - modification of semimetric distance δ to (partially) fulfill triangle inequality (T-error)
 - the lower T-error, the more precise but slower search, and vice versa



Non-metric similarity models

- general indexing of non-metric distances
 - TriGen algorithm
 - modification of semimetric distance δ to (partially) fulfill triangle inequality (T-error)
 - the lower T-error, the more precise but slower search, and vice versa
 - NM-tree (index+TriGen, multiple precisions at query time)



Non-metric similarity models

- general indexing of non-metric distances
 - TriGen algorithm
 - modification of semimetric distance δ to (partially) fulfill triangle inequality (T-error)
 - the lower T-error, the more precise but slower search, and vice versa
 - NM-tree (index+TriGen, multiple precisions at query time)
- pros – universal non-metric index
- cons – might lead to large intrinsic dimensionality

T. Skopal, J. Lokoč. NM-tree: Flexible Approximate Similarity Search in Metric and Non-metric Spaces, **DEXA 2008**, Turin, Italy, LNCS 5181, Springer

T. Skopal. Unified Framework for Fast Exact and Approximate Search in Dissimilarity Spaces, **ACM Transactions on Database Systems**, 32(4):1-47, 2007

T. Skopal. On Fast Non-Metric Similarity Search by Metric Access Methods, **EDBT 2006**, Munich, Germany, LNCS 3896, Springer

Non-metric similarity models

- applications for non-metric similarity search in bioinformatics
 - mass spectra interpretation (protein identification)
 - protein identification based on (3D) structure
 - RNA identification based on (3D) structure

J. Novák, T. Skopal, D. Hoksza, J. Lokoč. Non-metric Similarity Search of Tandem Mass Spectra Including Posttranslational Modifications, **Journal of Discrete Algorithms**, Elsevier, 2012

J. Galgonek, D. Hoksza, T. Skopal. SProt: sphere-based protein structure similarity algorithm, **Proteome Science**, 9(Suppl 1):S20, 2011

D. Hoksza, D. Svozil. SETTER - RNA SEcondary sTructure-based TERtiary Structure Similarity Algorithm, **ISBRA 2011**, Changsha, China, Springer

Data-centric applications

- similarity search as “exposed technology”
 - querying mechanism is the main functionality
 - the user must cooperate
(e.g., formulate a query, or browse/explore)
 - today mostly desktop applications
- data is the final product
 - is this the “killer application” for similarity search?
...depends

Data-centric applications

FindSounds

Search the Web for Sounds

Search for [Help](#)

See examples in [English](#)
[Deutsch](#) [Español](#) [Français](#) [Portugues](#)
[Chinese 中文](#) [Japanese 日本語](#) [Russian Русский язык](#)

File Formats	Number of Channels	Minimum Resolution	Minimum Sample Rate	Maximum File Size
<input checked="" type="checkbox"/> AIFF	<input checked="" type="checkbox"/> mono	8-bit	8000 Hz	2 MB
<input checked="" type="checkbox"/> AU	<input checked="" type="checkbox"/> stereo			
<input checked="" type="checkbox"/> MP3				
<input checked="" type="checkbox"/> WAVE				

Sounds 1-10 of 200 similar to this sound

-  http://sep800_mine.nu/files/sounds/carhornshort.wav
car horn
5k, mono, 8-bit, 8000 Hz, 0.6 seconds [show page](#) | [e-mail this sound](#) | [tweet this sound](#)
100%
-  <http://amazingsounds.iespana.es/carhornshort.wav>
car horn
5k, mono, 8-bit, 8000 Hz, 0.6 seconds [show page](#) | [e-mail this sound](#) | [tweet this sound](#)
93%
-  <http://amazingsounds.iespana.es/carhornrtice.wav>
car horn
11k, mono, 8-bit, 8000 Hz, 1.4 seconds [show page](#) | [e-mail this sound](#) | [tweet this sound](#)
93%
-  http://sep800_mine.nu/files/sounds/carhornrtice.wav
car horn
11k, mono, 8-bit, 8000 Hz, 1.4 seconds [show page](#) | [e-mail this sound](#) | [tweet this sound](#)
64%
-  http://cd.textfiles.com/10000gp2/500SNDS/BONG_3.WAV
11k, mono, 8-bit, 11025 Hz, 1.0 seconds [show page](#) | [e-mail this sound](#) | [tweet this sound](#)
63%

GazoPaβ

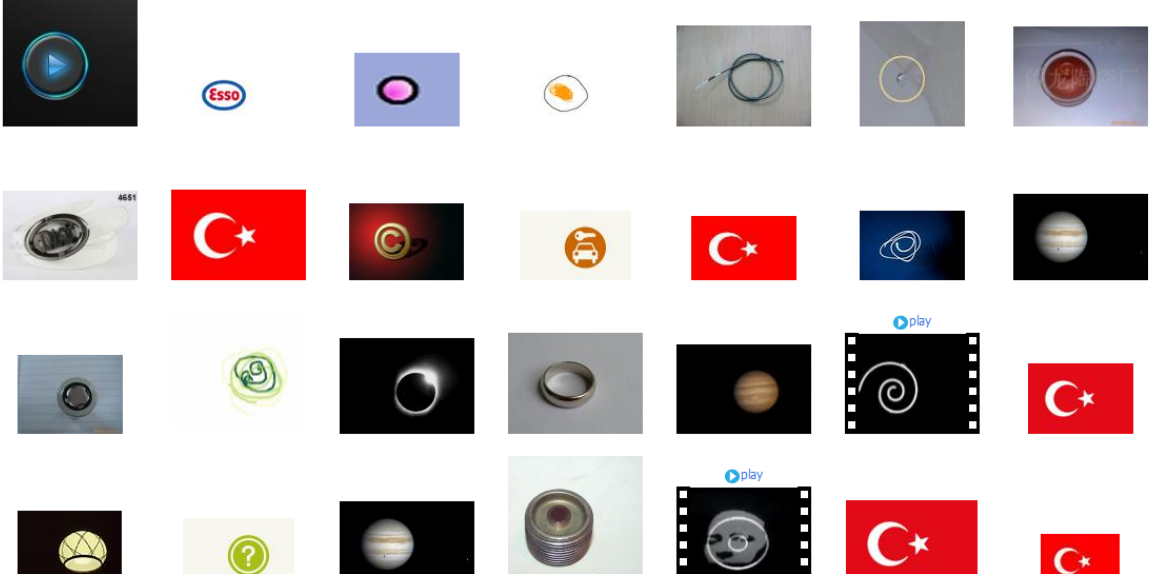
similar image search

Created: less than a minute ago
Size: 400x300

Options: **Shape** Any size Any time Gray scale only Omit same [Reset parameters](#) Safe search is on

[All](#) [Video](#) [News](#) [Sports](#) [Twitter](#) [Funny](#) [Flickr](#)

Results 1 - 30 of 1000 for **key image**



Data-centric applications

P3S: protein structure similarity search (ver.: 1.0.1)

► Query

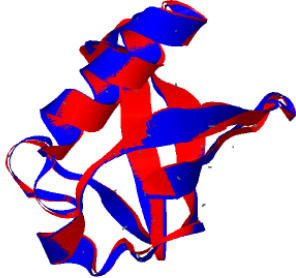
▼ Result

Result List

Name	Score	Coverage
d1xd3b_	0.95	98.7%
d1ogwa_	0.94	94.7%
d1zgub1	0.94	98.7%
d3by4b1	0.94	96.1%
d2c7nb1	0.93	96.1%
d2o6vd1	0.93	97.4%
d2d3ga1	0.93	94.7%
d2o6vh1	0.93	98.7%
d1ndda_	0.92	97.4%
d1bt0a_	0.92	96.1%
d1nndb_	0.91	97.4%
d1nndc_	0.91	96.1%
d1sifa_	0.91	93.4%
d1c3ta_	0.91	97.4%
d3cmmb1	0.91	98.7%
d1ud7a_	0.90	97.4%
d2faza1	0.89	97.4%
d1yqba1	0.89	98.7%
d2bweu1	0.88	94.7%
d2zccc1	0.88	93.4%

[Download results](#)

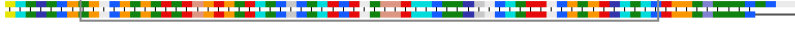
Superposition



[Open in external application](#)

Jmol [Settings](#)

Alignment



Query: T G K T T I E V E P S D T I E N V K A K I Q D K E G P P D Q Q R L F A G K Q L E D G R T I S D Y N I Q I
Result: T G K T T I E V E P S D T I E N V K A K I Q D K E G P P D Q Q R L F A G K Q L E D G R T I S D Y N I Q I

[Open in external application](#)

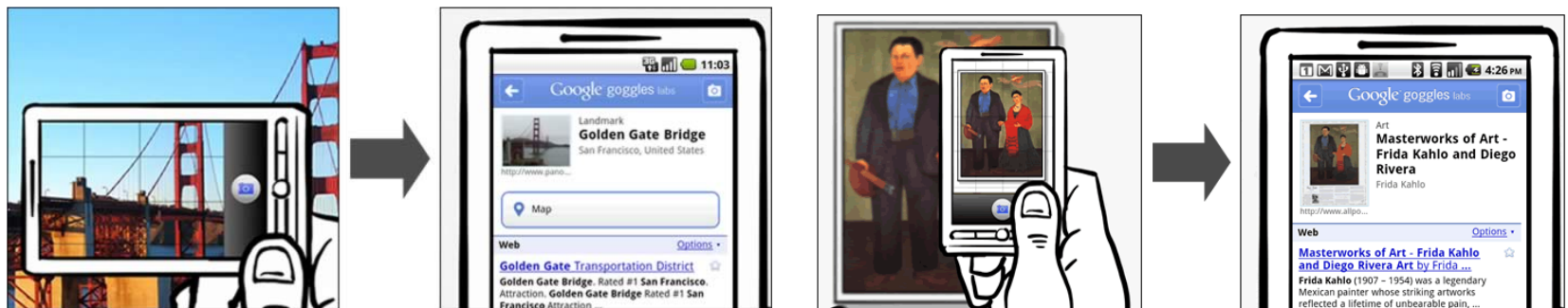
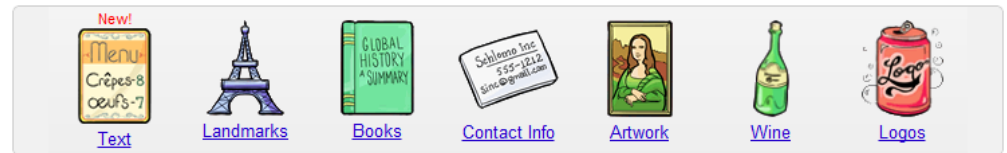
[Settings](#)

Service-centric applications

- similarity search as “hidden technology”
 - querying mechanism and even the data type hidden to the user
 - data is just means of search
- service is the final product
 - much wider applicability of similarity search?
 - e.g., e-commerce, “audiovisual wikipedia”, etc.

Service-centric applications

- Google Goggles – mobile version of Google
 - combination of OCR and pattern matching (similarity-based) techniques and content-based image retrieval
 - Android, iOS (iPhone, iPad)
- nowadays provides identification of:



Service-centric applications

- augmented reality
 - very near future

