


# High-Dimensional Index Structures: Database Support for Next Decade's Applications



**Stefan Berchtold** *stb software technologie  
beratung gmbh*

Stefan.Berchtold@stb-gmbh.de

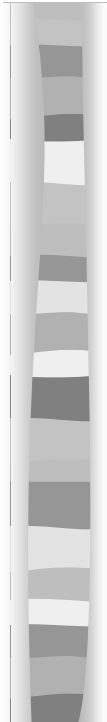
**Daniel A. Keim** *University of Halle-Wittenberg*

keim@informatik.uni-halle.de



## Modern Database Applications

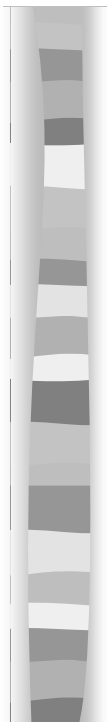
- |                         |                         |
|-------------------------|-------------------------|
| ■ Multimedia Databases  | ■ Data Warehouses       |
| – large data set        | – large data set        |
| – content-based search  | – data mining           |
| – feature-vectors       | – many attributes       |
| – high-dimensional data | – high-dimensional data |



# Overview

1. Modern Database Applications
2. Effects in High-Dimensional Space
3. Models for High-Dimensional Query Processing
4. Indexing High-Dimensional Space
  - 4.1 kd-Tree-based Techniques
  - 4.2 R-Tree-based Techniques
  - 4.3 Other Techniques
  - 4.4 Optimization and Parallelization
5. Open Research Topics
6. Summary and Conclusions

3



# Effects in High-Dimensional Spaces

- Exponential dependency of measures on the dimension
- Boundary effects
- No geometric imagination
  - ⇒ Intuition fails

*The Curse of Dimensionality*

4



## Notations and Assumptions

- $N$  data items
- $d$  dimensions
- data space normalized to  $[0, 1]^d$
- query types: range, partial range, NN
- for analysis: uniform data
- but not:  $N$  exponentially depends on  $d$

5



## Exponential Growth of Volume

- Hyper-cube

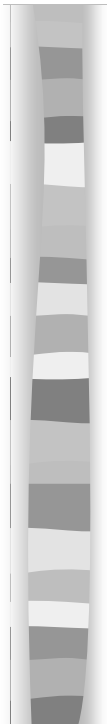
$$Volume_{cube}(edge, d) = edge^d$$

$$Diagonal_{cube}(edge, d) = edge \cdot \sqrt{d}$$

- Hyper-sphere

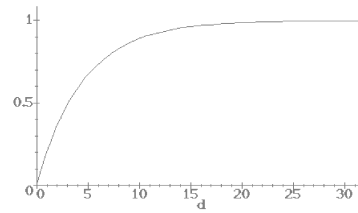
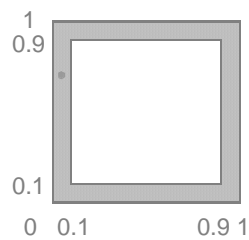
$$Volume_{sphere}(radius, d) = radius^d \cdot \frac{\sqrt{\pi^d}}{\Gamma(d/2 + 1)}$$

6

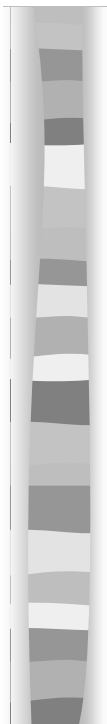


## The Surface is Everything

- Probability that a point is closer than 0.1 to a  $(d-1)$ -dimensional surface



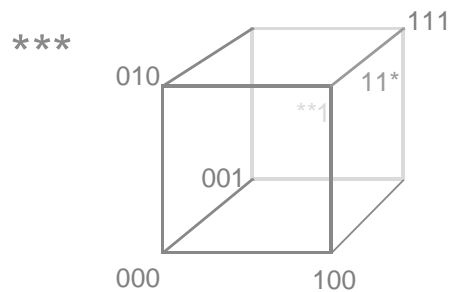
7



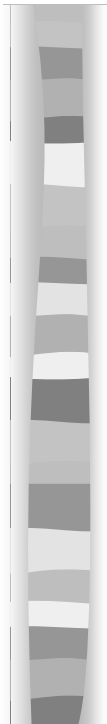
## Number of Surfaces

- How much  $k$ -dimensional surfaces has a  $d$ -dimensional hypercube  $[0..1]^d$  ?

$$\binom{d}{k} \cdot 2^{(d-k)}$$

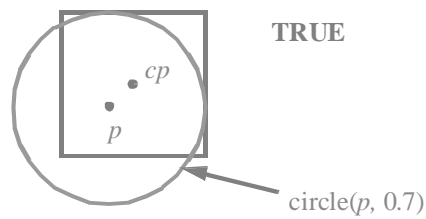


8

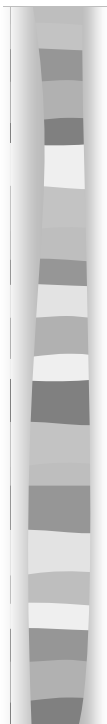


## “Each Circle Touching All Boundaries Includes the Center Point”

- $d$ -dimensional cube  $[0, 1]^d$
- $cp = (0.5, 0.5, \dots, 0.5)$
- $p = (0.3, 0.3, \dots, 0.3)$
- 16- $d$ : circle  $(p, 0.7)$ , distance  $(p, cp)=0.8$



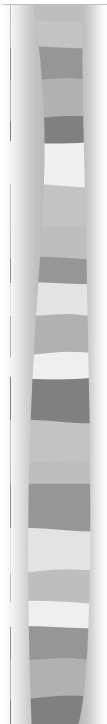
9



## Database-Specific Effects

- Selectivity of queries
- Shape of data pages
- Location of data pages

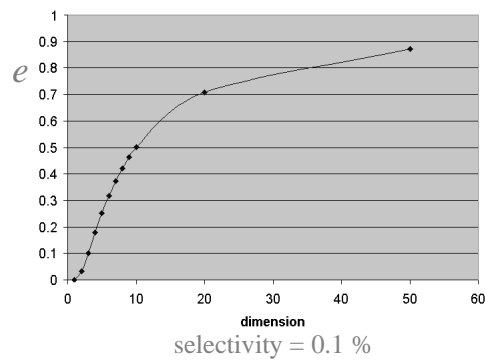
10



## Selectivity of Range Queries

- The selectivity depends on the volume of the query

$$e = d\sqrt{Vol_{cube}}$$

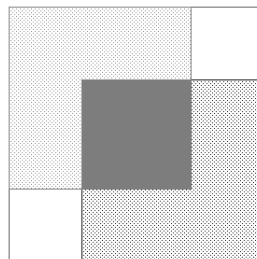


11

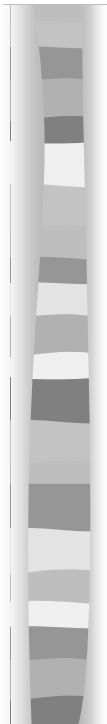


## Selectivity of Range Queries

- In high-dimensional data spaces, there exists a region in the data space which is affected by ANY range query (assuming uniformity)



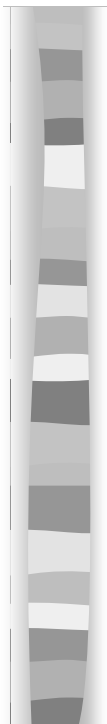
12



## Shape of Data Pages

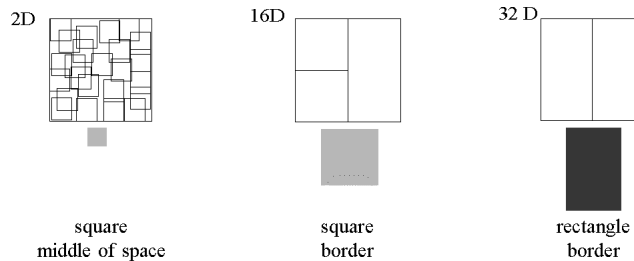
- uniformly distributed data  
⇒ each data page has the same volume
- split strategy: split always at the 50%-quantile
- number of split dimensions:  
$$d' = \log_2\left(\frac{N}{C_{eff}(d)}\right)$$
- extension of a “typical” data page: 0.5 in  $d'$  dimensions, 1.0 in  $(d-d')$  dimensions

13



## Location and Shape of Data Pages

- Data pages have large extensions
- Most data pages touch the surface of the data space on most sides



14



## Overview

1. Modern Database Applications
2. Effects in High-Dimensional Space
- 3. Models for High-Dimensional Query Processing**
4. Indexing High-Dimensional Space
  - 4.1 kd-Tree-based Techniques
  - 4.2 R-Tree-based Techniques
  - 4.3 Other Techniques
  - 4.4 Optimization and Parallelization
5. Open Research Topics
6. Summary and Conclusions

15



## Models for High-Dimensional Query Processing

- Traditional NN-Model [FBF 77]
- Exact NN-Model [BBKK 97]
- Analytical NN-Model [BBKK 00]
- Modeling the NN-Problem [BGRS 99]
- Modeling Range Queries [BBK 98]

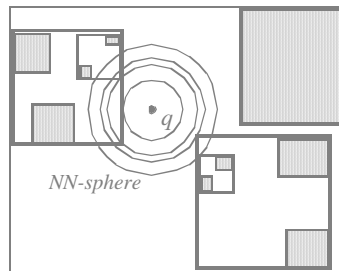
16





## Nearest-Neighbor Algorithms

- Algorithm by Hjaltason et Samet [HS 95]
  - loads only pages intersecting the NN-sphere
  - optimal algorithm



17



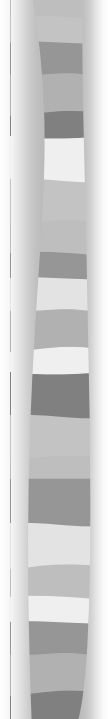
## Traditional NN-Model

- Friedman, Finkel, Bentley-Model [FBF 77]

### Assumptions:

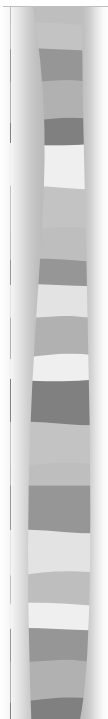
- number of data points  $N$  goes towards infinity  
( $\Rightarrow$  unrealistic for real data sets)
- no boundary effects  
( $\Rightarrow$  large errors for high-dim. data)

18



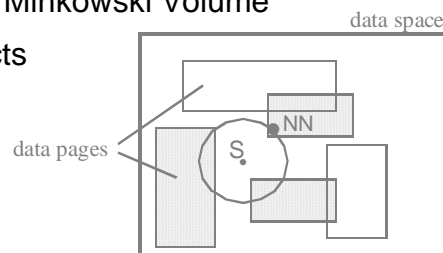
# Exact NN-Model [BBKK 97]

- Goal: Determination of the number of data pages which have to be accessed on the average
- Three Steps:
  1. Distance to the Nearest Neighbor
  2. Mapping to the Minkowski Volume
  3. Boundary Effects



# Exact NN-Model

1. **Distance to the Nearest Neighbor**
2. Mapping to the Minkowski Volume
3. Boundary Effects



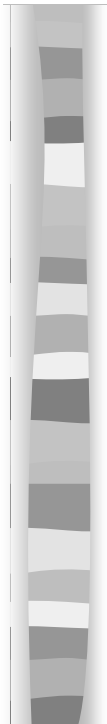
*Distribution function*

$$P(NN - dist = r) = 1 - P(\text{None of the } N \text{ points intersects NN - sphere})$$

$$= (1 - (1 - Vol_{avg}^d(r))^N)$$

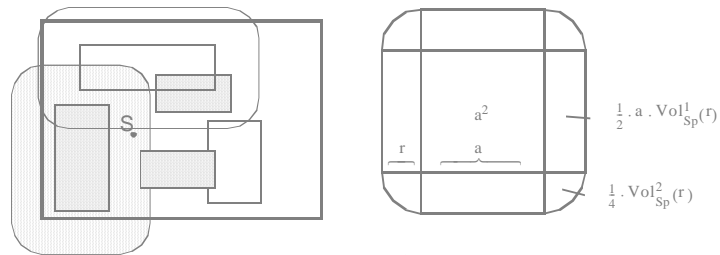
*Density function*

$$\frac{d}{dr} P(NN - dist = r) = \frac{d}{dr} Vol_{avg}^d(r) \cdot N \cdot (1 - Vol_{avg}^d(r))^{N-1}$$



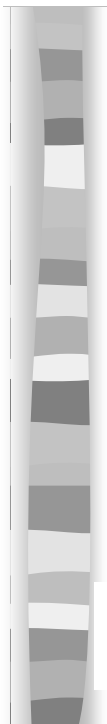
## Exact NN-Model

1. Distance to the Nearest Neighbor
2. **Mapping to the Minkowski Volume**
3. Boundary Effects



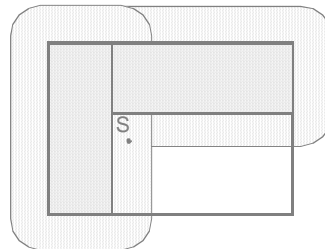
*Minkowski Volume:* 
$$Vol_{\text{Mink}}^d(r) = \sum_{i=0}^d \binom{d}{i} \cdot a^{d-i} \cdot Vol_{\text{Sp}}^i(r)$$

21



## Exact NN-Model

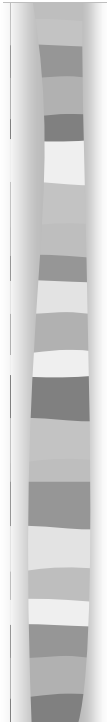
1. Distance to the Nearest Neighbor
2. Mapping to the Minkowski Volume
3. **Boundary Effects**



*Generalized Minkowski Volume with boundary effects:*

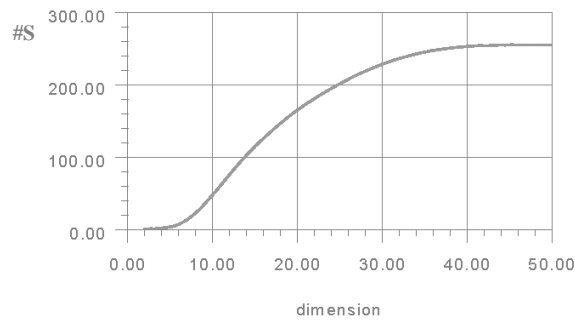
$$\#S(r) = \sum_{k=0}^{d'} \sum_{\{i_1, \dots, i_k\} \in \mathcal{P}(\{1, \dots, d'\})} Vol(SP^k([a_{i_1}, \dots, a_{i_k}], r) \cap DS) \quad \text{where } d' = \left\lceil \log_2 \left( \frac{N}{C_{\text{eff}}} \right) \right\rceil$$

22

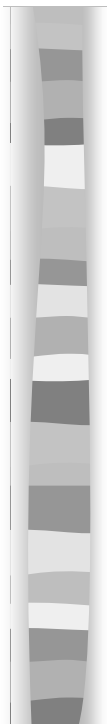


## Exact NN-Model

$$\begin{aligned}
 E(\#S) &= \int_0^{\infty} \#Pages(r) \cdot p(r) \, dr \\
 &= N \cdot \int_0^{\infty} Vol_{avg}^d(r) \cdot (1 - Vol_{avg}^d(r))^{N-1} \cdot \sum_{k=0}^d \sum_{\{i_1, \dots, i_k\} \in \mathcal{P}(\{1, \dots, d\})} Vol(SP^k(\{a_{i_1}, \dots, a_{i_k}\}, r) \cap DS) \, dr
 \end{aligned}$$



23



## Approximate NN-Model [BBKK 00]

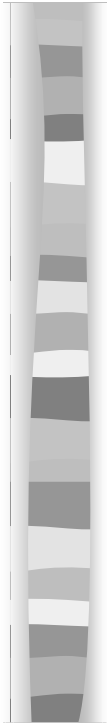
### 1. Distance to the Nearest-Neighbor

Idea:

Nearest-neighbor Sphere contains  $1/N$  of the volume of the data space

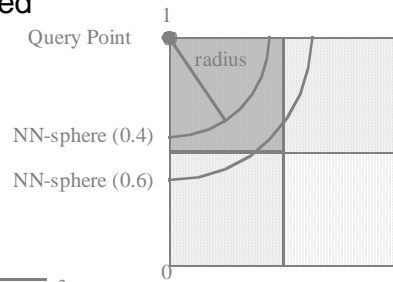
$$Vol_{Sp}^d(NN-dist) = \frac{1}{N} \Rightarrow NN-dist(N, d) = \frac{1}{\sqrt{\pi}} \cdot d \sqrt{\frac{\Gamma(d/2 + 1)}{N}}$$

25



## Approximate NN-Model

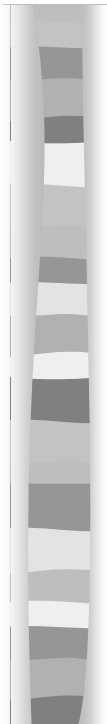
2. Distance threshold which requires more data pages to be considered



$$NN-dist(N, d) = 0.5 \cdot \sqrt{i}$$

$$\Leftrightarrow i = \left( \frac{1}{\sqrt{\pi}} \cdot \frac{d \sqrt{\Gamma(d/2 + 1)}}{N} \right)^2 \Rightarrow i \approx \frac{2 \cdot d}{e \cdot \pi} \cdot d \sqrt{\frac{\pi \cdot d^3}{4 \cdot N^2}}$$

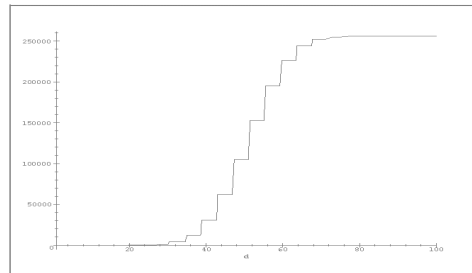
26



## Approximate NN-Model

3. Number of pages

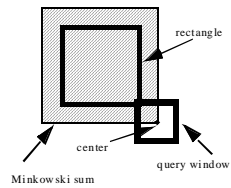
$$\#S(d) = \sum_{k=0}^{d'} \binom{d'}{k} = \sum_{k=0} \left( \left\lceil \log_2 \left( \frac{N}{C_{\text{eff}}} \right) \right\rceil \right)_k$$



27

## Modeling Range-Queries [BBK 98]

- Idea: Use Minkowski-sum to determine the probability that a data page (URC, LLC) is loaded



$$P_{\text{no\_bound\_eff}}(q) = \sum \prod_{i,j=0}^{d-1} (\text{URC}_{i,j} - \text{LLC}_{i,j} + q)$$

$$P_{\text{bound\_eff}}(q) = \sum \prod_{i,j=0}^{d-1} \frac{\min(\text{URC}_{i,j}, 1-q) - \max(\text{LLC}_{i,j} - q, 0)}{1-q}$$

30

## The Problem of Searching the Nearest Neighbor [BGRS 99]

- Observations:

- When increasing the dimensionality, the nearest-neighbor distance grows.
- When increasing the dimensionality, the farthest-neighbor distance grows.
- The nearest-neighbor distance grows FASTER than the farthest-neighbor distance.
- For  $d \rightarrow \infty$ , the nearest-neighbor distance equals to the farthest-neighbor distance.

31



## When Is Nearest Neighbor meaningful?

- Statistical Model:

- For the  $d$ -dimensional distribution holds:

$$\lim_{d \rightarrow \infty} (\text{var}(D_d^p) / E(D_d^p)^2) = 0$$

where  $D$  is the distribution of the distance of the query point and a data point and we consider a  $L_p$  metric.

- This is true for synthetic distributions such as normal, uniform, zipfian, etc.
- This is NOT true for clustered data.

32



## Overview

1. Modern Database Applications
2. Effects in High-Dimensional Space
3. Models for High-Dimensional Query Processing
- 4. Indexing High-Dimensional Space**
  - 4.1 kd-Tree-based Techniques
  - 4.2 R-Tree-based Techniques
  - 4.3 Other Techniques
  - 4.4 Optimization and Parallelization
5. Open Research Topics
6. Summary and Conclusions

33



## Indexing High-Dimensional Space

- **Criteria**
- **kd-Tree-based Index Structures**
- **R-Tree-based Index Structures**
- **Other Techniques**
- **Optimization and Parallelization**

34



## Criteria [GG 98]

- **Structure of the Directory**
- **Overlapping vs. Non-overlapping Directory**
- **Type of MBR used**
- **Static vs. Dynamic**
- **Exact vs. Approximate**

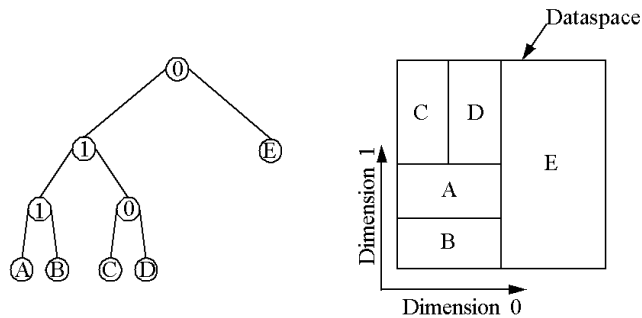
35



## The kd-Tree [Ben 75]

### ■ Idea:

Select a dimension, split according to this dimension and do the same recursively with the two new sub-partitions



36

## The kd-Tree

### ■ Plus:

- fanout constant for arbitrary dimension
- fast insertion
- no overlap

### ■ Minus:

- depends on the order of insertion  
(e.g., not robust for sorted data)
- dead space covered
- not appropriate for secondary storage

37

## The kdB-Tree [Rob 81]

### ■ Idea:

- Aggregate kd-Tree nodes into disk pages
- Split data pages in case of overflow (B-Tree-like)

### ■ Problem:

- splits are not local
- forced splits

38

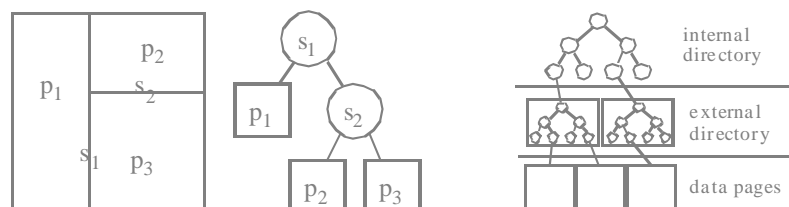
## The LSD<sup>h</sup>-Tree [Hen 98]

### ■ Two-level directory:

first level in main memory

### ■ To avoid dead space:

only actual data regions are coded



39



## The LSD<sup>h</sup>-Tree

- Fast insertion
- Search performance (NN) competitive to X-Tree
- Still sensitive to pre-sorted data
- Technique of CADR (Coded Actual Data Regions) is applicable to many index structures

40



## The VAMSplit Tree [JW 96]

- Idea:
  - Split at the point where maximum variance occurs (rather than in the middle)
- sort data in main memory
- determine split position and recurse
- Problems:
  - data must fit in main memory
  - benefit of variance-based split is not clear

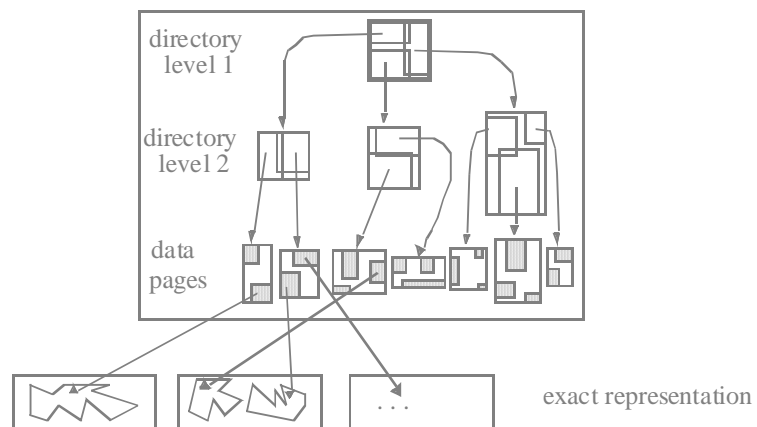
41

# Overview


1. Modern Database Applications
2. Effects in High-Dimensional Space
3. Models for High-Dimensional Query Processing
4. Indexing High-Dimensional Space
  - 4.1 kd-Tree-based Techniques
  - 4.2 R-Tree-based Techniques**
  - 4.3 Other Techniques
  - 4.4 Optimization and Parallelization
5. Open Research Topics
6. Summary and Conclusions

42

## R-Tree: [Gut 84] The Concept of Overlapping Regions



43



## Variants of the R-Tree

### Low-dimensional

- R<sup>+</sup>-Tree [SRF 87]
- R<sup>\*</sup>-Tree [BKSS 90]
- Hilbert R-Tree [KF94]

### High-dimensional

- TV-Tree [LJF 94]
- X-Tree [BKK 96]
- SS-Tree [WJ 96]
- SR-Tree [KS 97]

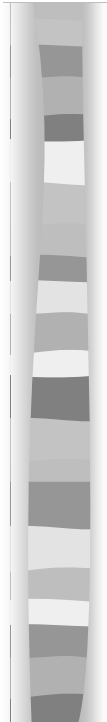
44



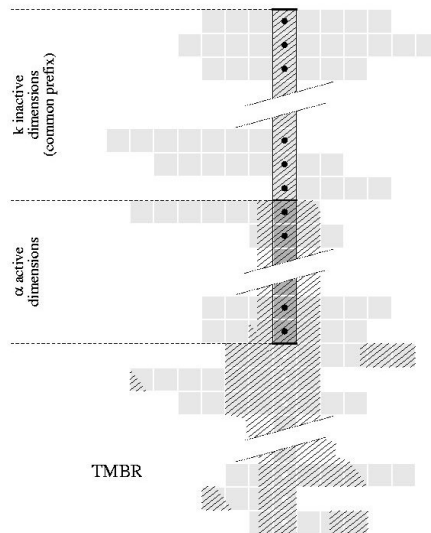
## The TV-Tree [LJF 94] (Telescope-Vector Tree)

- Basic Idea: Not all attributes/dimensions are of the same importance for the search process.
- Divide the dimensions into three classes
  - attributes which are shared by a set of data items
  - attributes which can be used to distinguish data items
  - attributes to ignore

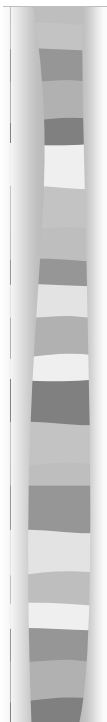
45



## Telescope Vectors



46



## The TV-Tree

- Split algorithm:  
either increase dimensionality of TV  
or split in the given dimensions
- Insert algorithm: similar to R-Tree
- Problems:
  - how to choose the right metric
  - high overlap in case of most metrics
  - complex implementation

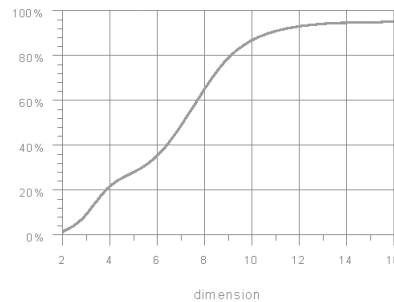
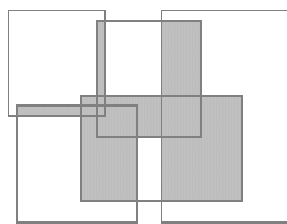
47

## The X-Tree [BKK 96] (eXtended-Node Tree)

### ■ Motivation:

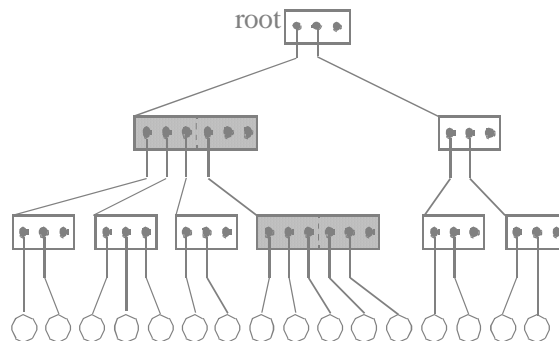
Performance of the R-Tree degenerates in high dimensions

### ■ Reason: overlap in the directory



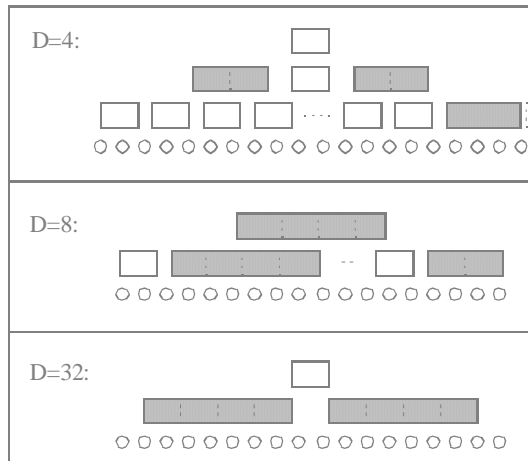
## The X-Tree

- X-tree avoids overlap in the directory by using
  - an overlap-free split
  - the concept of supernodes



# The X-Tree

## Examples for X-Trees with different dimensionality



51

# The X-Tree

## Overlap-Free Split

### □ Definition (Split):

The split of a node  $S = \{mbr_1, \dots, mbr_n\}$  into two subnodes  $S_1$  and  $S_2$  ( $S_1 \neq \emptyset$  and  $S_2 \neq \emptyset$ ) is defined as

$$Split(S) = \{(S_1, S_2) \mid S = S_1 \cup S_2 \wedge S_1 \cap S_2 = \emptyset\}.$$

The split is called

- (1) overlap-minimal iff  $\|MBR(S_1) \cap MBR(S_2)\|$  is minimal
- (2) overlap-free iff  $\|MUR(S_1) \cap MUR(S_2)\| = 0$
- (3) balanced iff  $-\varepsilon \leq |S_1| - |S_2| \leq \varepsilon$  (for small  $\varepsilon$ ).

52



# The X-Tree

- **Lemma 1** (for uniformly distributed point data)

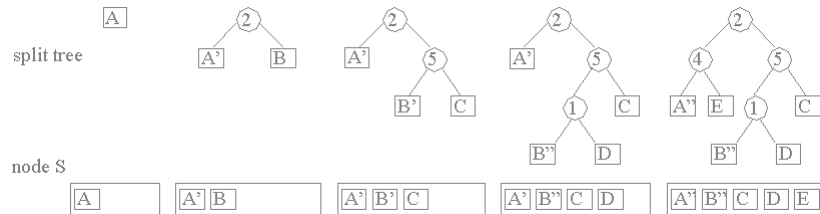
*Split(S) is overlap-free*  $\Leftrightarrow$

$\exists d \in \{1, \dots, D\} \forall mbr \in S: mbr \text{ has been split according to } d$

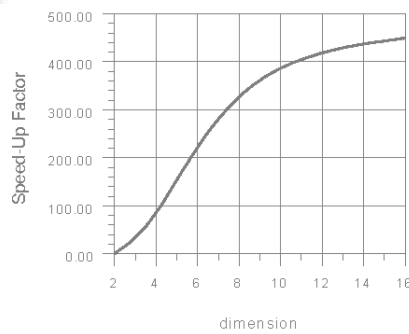
- **Lemma 2**

*For point data, an overlap-free split always exists.*

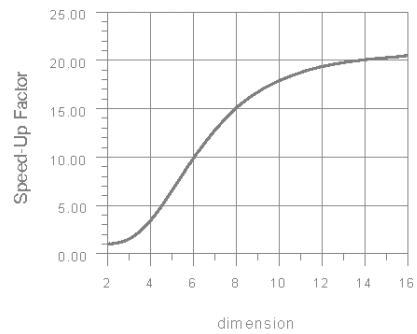
Example split history:



# Speed-Up of X-Tree over the R\*-Tree



*Point Query*



*10 NN Query*

## Bulk-Load of X-Trees [BBK 98a]

- **Observation:**

In order to split a data set, we do not have to sort it

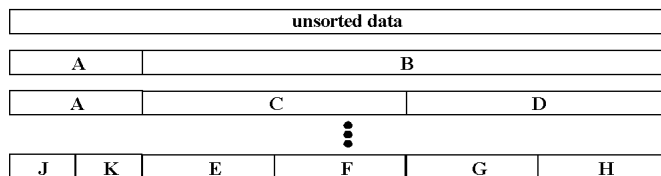
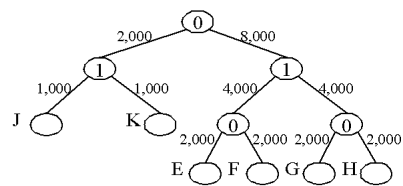
- Recursive top-down partitioning of the data set
- Quicksort-like algorithm
- Improved data space partitioning

56

## Example

external bisection split strategy

dimension 1	J 1,000	E 2,000	F 2,000
	K 1,000	G 2,000	H 2,000
dimension 0			



57

## Unbalanced Split

- Probability that a data page is loaded when processing a range query of edge length 0.6 (for three different split strategies)

	5/6	5/6
2/3	1	1
1/3	5/6	5/6
	1/2	

$$P_1(0.6) = 5.33$$

	5/12	
5/6	15/16	
5/8	1/2	1
5/16	1	25/32
	1/5	7/15

$$P_2(0.6) = 4.64$$

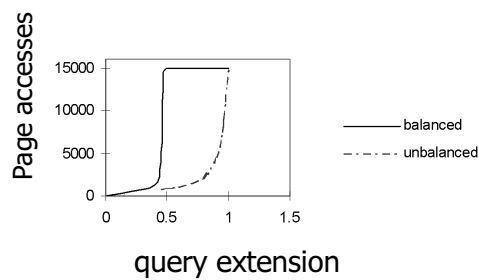
	5/12	
5/6	5/8	5/8
1/2	1	1
1/6	5/12	
	1/4	3/4

$$P_3(0.6) = 4.08$$

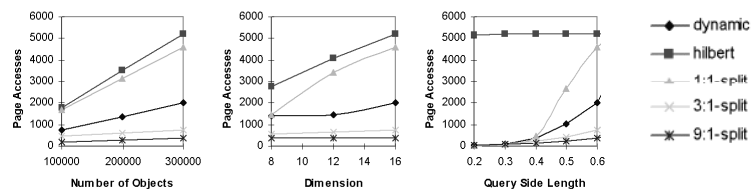
58

## Effect of Unbalanced Split

*In Theory:*



*In Practice:*

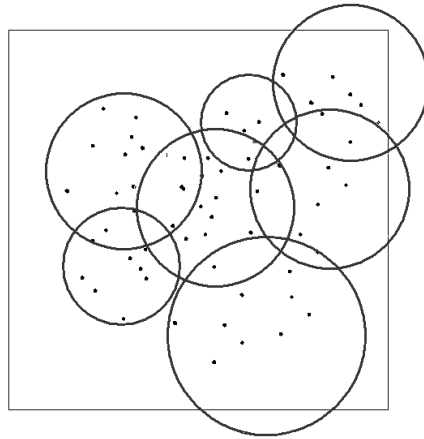


59



## The SS-Tree [WJ 96] (Similarity-Search Tree)

- Idea:  
Split data space into spherical regions
- small MINDIST
- high fanout
- Problem: overlap

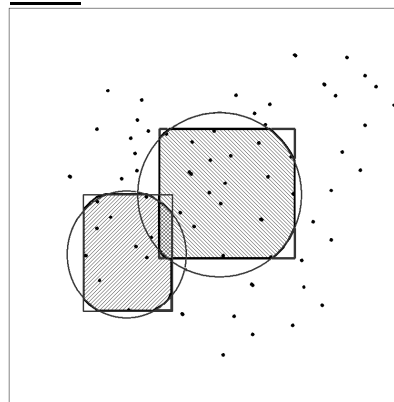


60



## The SR-Tree [KS 97] (Similarity-Search R-Tree)

- Similar to SS-Tree, but:
- Partitions are intersections of spheres and hyper-rectangles
- Low overlap



61



## Overview

1. Modern Database Applications
2. Effects in High-Dimensional Space
3. Models for High-Dimensional Query Processing
4. Indexing High-Dimensional Space
  - 4.1 kd-Tree-based Techniques
  - 4.2 R-Tree-based Techniques
  - 4.3 Other Techniques**
  - 4.4 Optimization and Parallelization
5. Open Research Topics
6. Summary and Conclusions

62



## Other Techniques

- Pyramid-Tree [BBK 98]
- VA-File [WSB 98]
- Voroni-based Indexing [BEK+ 98]

63

## The Pyramid-Tree [BBK 98]

### ■ Motivation:

Index-structures such as the X-Tree have several drawbacks

- the split strategy is sub-optimal
- all page accesses result in random I/O
- high transaction times (insert, delete, update)

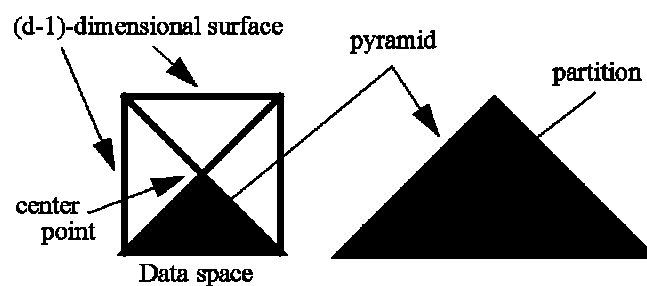
### ■ Idea:

Provide a data space partitioning which can be seen as a mapping from a  $d$ -dim. space to a 1-dim. space and make use of B<sup>+</sup>-Trees

64

## The Pyramid-Mapping

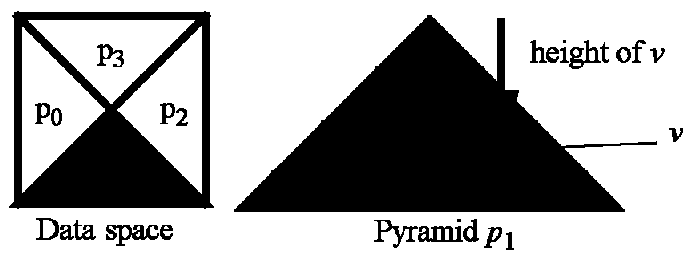
- Divide the space into 2d pyramids
- Divide each pyramid into partitions
- Each partition corresponds to a B<sup>+</sup>-Tree page



65

## The Pyramid-Mapping

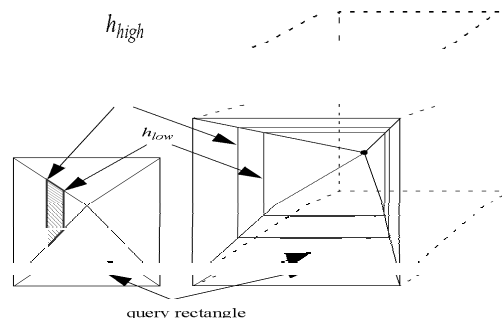
- A point in a high-dimensional space can be addressed by the number of the pyramid and the height within the pyramid.



66

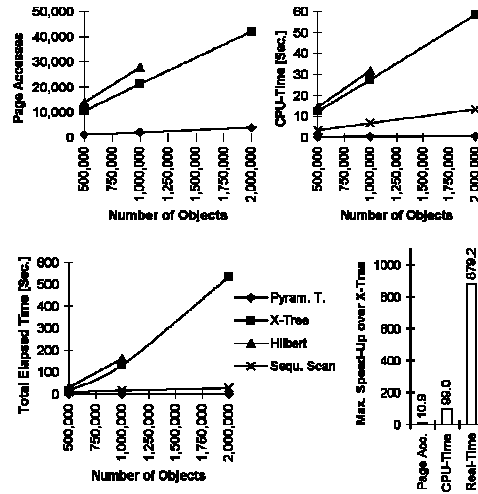
## Query Processing using a Pyramid-Tree

- Problem:  
Determine the pyramids intersected by the query rectangle and the interval  $[h_{high}, h_{low}]$  within the pyramids.



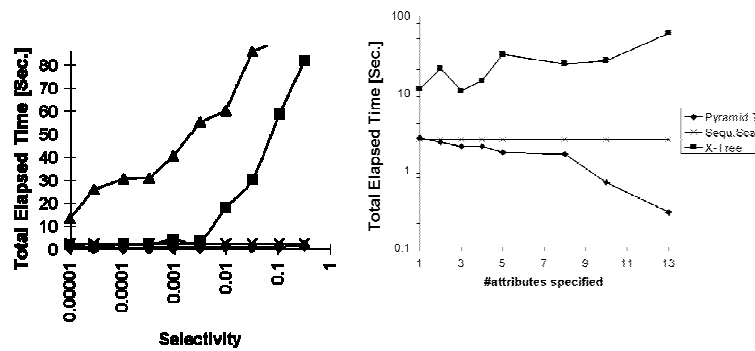
67

## Experiments (uniform data)



68

## Experiments (data from data warehouse)



69





## The VA-File [WSB 98] (Vector Approximation File)

- Idea:  
If NN-Search is an inherently linear problem, we should aim for speeding up the sequential scan.
- Use a coarse representation of the data points as an approximate representation (only  $i$  bits per dimension -  $i$  might be 2)
- Thus, the reduced data set has only the  $(i/32)$ -th part of the original data set

71



## The VA-File

- Determine  $(1/2^i)$ -quantiles of each dimension as partition boundaries
- Sequentially scan the coarse representation and maintain the actual NN-distance
- If a partition cannot be pruned according to its coarse representation, a look-up is made in the original data set

72

## The IQ-Tree [BBJ+ 00] (Independent Quantization)

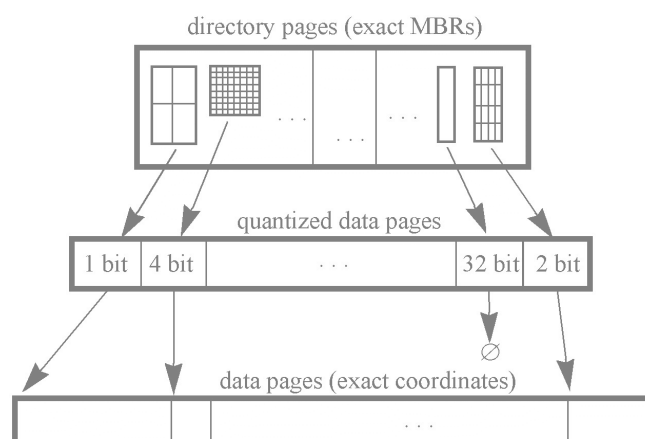
- Idea:

If the VA-file does a good job for uniform data and partitioning techniques do so for correlated data, let's find the optimum in between.

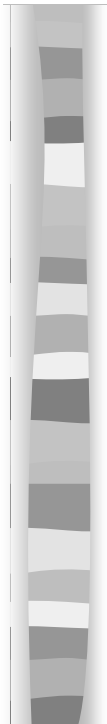
- Hybrid index / file structure
- 2-level directory: first level is a hierarchical directory, second level is an adaptive VA-file
- adapts the level of partitioning to the actual data

75

## The IQ-Tree - Structure



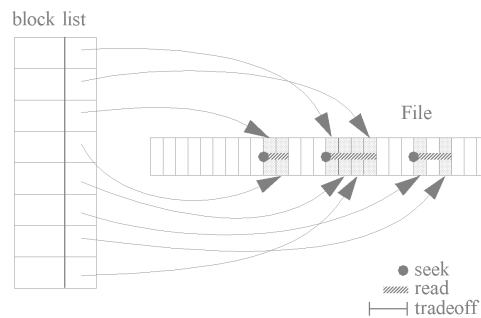
76



## New NN-Algorithm

### ■ Idea:

Overread pages if the (probabilistic) cost for overreading are smaller than the seek cost.



77

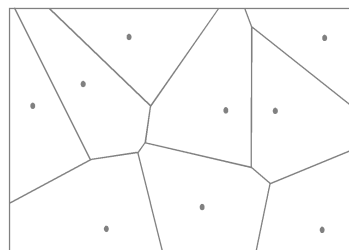


## Voronoi-based Indexing [BEK+ 98]

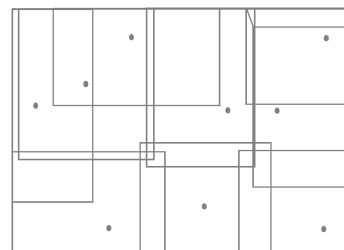
### ■ Idea:

Precalculation and indexing of the result space

⇒ Point query instead of NN-query



*Voroni-Cells*



*Approximated Voroni-Cells*

78



## Overview

1. Modern Database Applications
2. Effects in High-Dimensional Space
3. Models for High-Dimensional Query Processing
4. Indexing High-Dimensional Space
  - 4.1 kd-Tree-based Techniques
  - 4.2 R-Tree-based Techniques
  - 4.3 Other Techniques
  - 4.4 Optimization and Parallelization**
5. Open Research Topics
6. Summary and Conclusions

81



## Optimization and Parallelization

- Tree Striping [BBK+ 00]
- Parallel Declustering [BBB+ 97]
- Approximate Nearest Neighbor Search [GIM 99]

82

## Tree Striping [BBK+ 00]

- Motivation:

The two solutions to multidimensional indexing - inverted lists and multidimensional indexes - are both inefficient.

- Explanation:

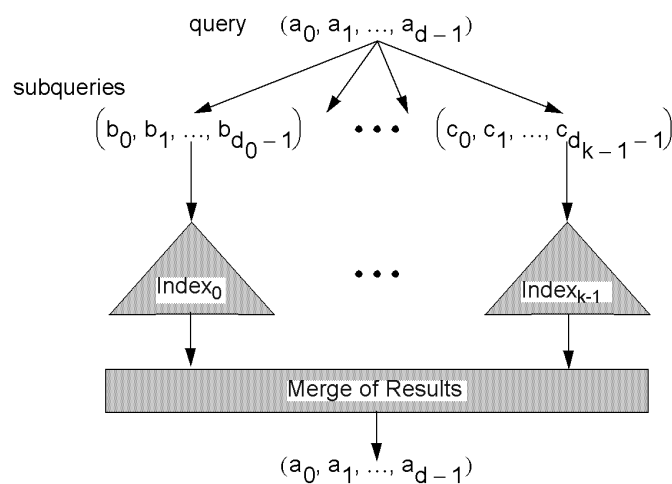
High dimensionality deteriorates the performance of indexes and increases the sort costs of inverted lists.

- Idea:

There must be an optimum in between high-dimensional indexing and inverted lists.

83

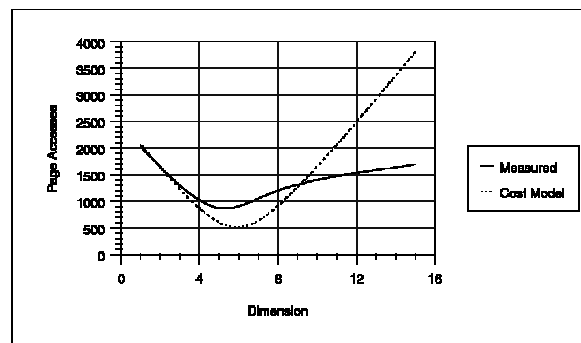
## Tree Striping - Example



84

## Experiments

- Real data, range queries,  $d$ -dimensional indexes



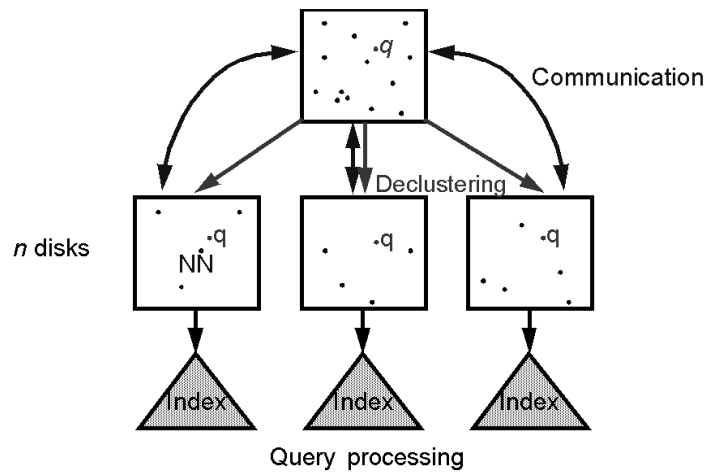
87

## Parallel Declustering [BBB+ 97]

- Idea:  
If NN-Search is an inherently linear problem, it is perfectly suited for parallelization.
- Problem:  
How to decluster high-dimensional data?

88

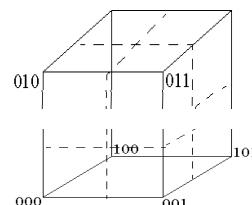
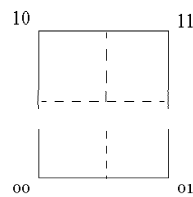
## Parallel Declustering



89

## Near-Optimal Declustering

- Each partition is connected with one corner of the data space  
Identify the partitions by their **canonical corner numbers**  
= bitstrings saying left = 0 and right = 1 for each dimension



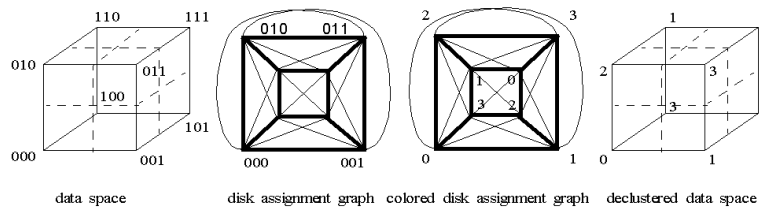
- Different degrees of neighborhood relationships:
  - Partitions are **direct** neighbors if they differ in exactly 1 dimension
  - Partitions are **indirect** neighbors if they differ in exactly 2 dimension

90

# Parallel Declustering

## Mapping of the Problem to a Graph:

partitions  $\Rightarrow$  vertices  
 neighborhood-relations  $\Rightarrow$  edges  
 disks  $\Rightarrow$  colors



91

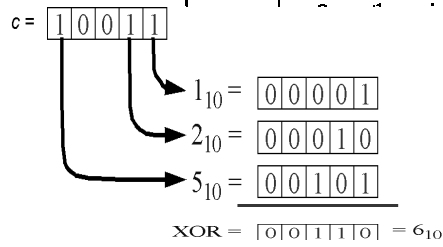
# Parallel Declustering

- Given: vertex number = corner number in binary representation

$$c = (c_{d-1}, \dots, c_0)$$

- Compute: vertex color  $\text{col}(c)$  as

$$\text{col}(c) = \left( \text{XOR}_{i=0}^{d-1} \left( \begin{matrix} i+1 & \text{if } c_i = 1 \\ 0 & \text{if } c_i = 0 \end{matrix} \right) \right)_{10}$$

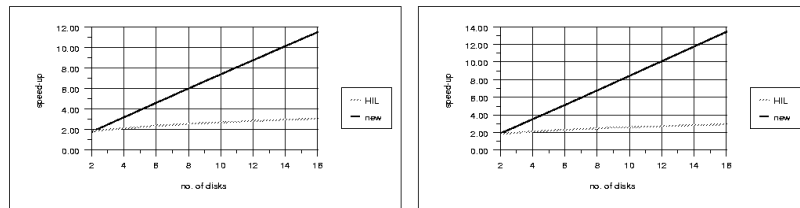


92



## Experiments

- Real data, comparison with Hilbert-declustering, # of disks vs. speed-up



93

## Approximate NN-Search

(Locality-Sensitive Hashing) [GIM 99]

- Idea:  
If it is sufficient to only select an approximate nearest-neighbor, we can do this much faster.
- Approximate Nearest-Neighbor: A point in distance  $(1 + \epsilon) \cdot NN_{dist}$  from the query point.

94



# Locality-Sensitive Hashing

## ■ Algorithm:

- Map each data point into a higher-dimensional binary space
- Randomly determine  $k$  projections of the binary space
- For each of the  $k$  projections determine the points having the same binary representations as the query point
- Determine the nearest-neighbors of all these points

## ■ Problems:

- How to optimize  $k$ ?
- What is the expected  $\epsilon$ ? (average and worst case)
- What is an approximate nearest-neighbor “worth”?

95



# Overview

1. Modern Database Applications
2. Effects in High-Dimensional Space
3. Models for High-Dimensional Query Processing
4. Indexing High-Dimensional Space
  - 4.1 kd-Tree-based Techniques
  - 4.2 R-Tree-based Techniques
  - 4.3 Other Techniques
  - 4.4 Optimization and Parallelization
- 5. Open Research Topics**
6. Summary and Conclusions

96



## Open Research Topics

- Partitioning strategies
- Parallel query processing
- Data reduction
- Approximate query processing
- High-dim. data mining & visualization
- The ultimate cost model

97



## Partitioning Strategies

- What is the optimal data space partitioning schema for nearest-neighbor search in high-dimensional spaces?
- Balanced or unbalanced?
- Pyramid-like or bounding boxes?
- How does the optimum changes when the data set grows in size or dimensionality?

98



## Parallel Query Processing

- Is it possible to develop parallel versions of the proposed sequential techniques?  
If yes, how can this be done?
- Which declustering strategies should be used?
- How can the parallel query processing be optimized?

99



## Data Reduction

- How can we reduce a large data warehouse in size such that we get approximate answers from the reduced data base?
- Tape-based data warehouses  
⇒ disk based
- Disk-based data warehouses  
⇒ main memory
- Tradeoff: accuracy vs. reduction factor

100



## Approximate Query Processing

- Observation:

Most similarity search applications do not require 100% correctness.

- Problem:

- What is a good definition for approximate nearest- neighbor search?
- How to exploit that fuzziness for efficiency?

101



## High-dimensional Data Mining & Data Visualization

- How can the proposed techniques be used for data mining?
- How can high-dimensional data sets and effects in high-dimensional spaces be visualized?

102



## Summary

### ■ Major research progress in

- understanding the nature of high-dim. spaces
- modeling the cost of queries in high-dim. spaces
- index structures supporting nearest-neighbor search and range queries

103



## Conclusions

### ■ Work to be done

- leave the clean environment
  - uniformity
  - uniform query mix
  - number of data items is exponential in  $d$
- address other relevant problems
  - partial range queries
  - approximate nearest neighbor queries

104



## Literature

- [AMN 95] Arya S., Mount D. M., Narayan O.: 'Accounting for Boundary Effects in Nearest Neighbor Searching', Proc. 11th Annual Symp. on Computational Geometry, Vancouver, Canada, pp. 336-344, 1995.
- [Ary 95] Arya S.: 'Nearest Neighbor Searching and Applications', Ph.D. Thesis, University of Maryland, College Park, MD, 1995.
- [BBB+ 97] Berchtold S., Böhm C., Braunmueller B., Keim D. A., Kriegel H.-P.: 'Fast Similarity Search in Multimedia Databases', Proc. ACM SIGMOD Int. Conf. on Management of Data, Tucson, Arizona, 1997.
- [BBK 98] Berchtold S., Böhm C., Kriegel H.-P.: 'The Pyramid-Tree: Indexing Beyond the Curse of Dimensionality', Proc. ACM SIGMOD Int. Conf. on Management of Data, Seattle, 1998.
- [BBK 98a] Berchtold S., Böhm C., Kriegel H.-P.: 'Improving the Query Performance of High-Dimensional Index Structures by Bulk Load Operations', 6th Int. Conf. On Extending Database Technology, in LNCS 1377, Valencia, Spain, pp. 216-230, 1998.
- [BBK+ 00] Berchtold S., Böhm C., Keim D., Kriegel H.-P., Xu X.: 'Optimal Multidimensional Query Processing Using Tree Striping', submitted for publication.

105



## Literature

- [BBKK 97] Berchtold S., Böhm C., Keim D., Kriegel H.-P.: 'A Cost Model For Nearest Neighbor Search in High-Dimensional Data Space', ACM PODS Symposium on Principles of Database Systems, Tucson, Arizona, 1997.
- [BBKK 00] Berchtold S., Böhm C., Keim D., Kriegel H.-P.: 'Optimized Processing of Nearest Neighbor Queries in High-Dimensional Spaces', submitted for publication.
- [BEK+ 98] Berchtold S., Ertl B., Keim D., Kriegel H.-P., Seidl T.: 'Fast Nearest Neighbor Search in High-Dimensional Spaces', Proc. 14th Int. Conf. on Data Engineering, Orlando, 1998.
- [BBJ+ 00] Berchtold S., Böhm C., Jagadish H.V., Kriegel H.-P., Sander J.: 'Independent Quantization: An Index Compression Technique for High-Dimensional Data Spaces: ', Int. Conf. on Data Engineering, San Diego, 2000.
- [BBKK 97] Berchtold S., Böhm C., Keim D., Kriegel H.-P.: 'A Cost Model For Nearest Neighbor Search in High-Dimensional Data Space', ACM PODS Symposium on Principles of Database Systems, Tucson, Arizona, 1997.
- [Ben 75] Bentley J. L.: 'Multidimensional Search Trees Used for Associative Searching', Comm. of the ACM, Vol. 18, No. 9, pp. 509-517, 1975.
- [BGRS 99] Beyer K., Goldstein J., Ramakrishnan R., Shaft U.: 'When Is "Nearest Neighbor" Meaningful?', Proc. Int. Conf. on Database Theory (ICDT), 1999, pp. 217-235.

106



## Literature

- [BK 97] Berchtold S., Kriegel H.-P.: '*S3: Similarity Search in CAD Database Systems*', Proc. ACM SIGMOD Int. Conf. on Management of Data, Tucson, Arizona, 1997.
- [BKK 96] Berchtold S., Keim D., Kriegel H.-P.: '*The X-tree: An Index Structure for High-Dimensional Data*', 22nd Conf. on Very Large Databases, Bombay, India, pp. 28-39, 1996.
- [BKSS 97] Berchtold S., Keim D., Kriegel H.-P.: '*Using Extended Feature Objects for Partial Similarity Retrieval*', VLDB Journal, Vol.4, 1997.
- [BKSS 90] Beckmann N., Kriegel H.-P., Schneider R., Seeger B.: '*The R\*-tree: An Efficient and Robust Access Method for Points and Rectangles*', Proc. ACM SIGMOD Int. Conf. on Management of Data, Atlantic City, NJ, pp. 322-331, 1990.
- [CD 97] Chaudhuri S., Dayal U.: '*Data Warehousing and OLAP for Decision Support*', Tutorial, Proc. ACM SIGMOD Int. Conf. on Management of Data, Tucson, Arizona, 1997.
- [Cle 79] Cleary J. G.: '*Analysis of an Algorithm for Finding Nearest Neighbors in Euclidean Space*', ACM Trans. on Mathematical Software, Vol. 5, No. 2, pp.183-192, 1979.

107



## Literature

- [FBF 77] Friedman J. H., Bentley J. L., Finkel R. A.: '*An Algorithm for Finding Best Matches in Logarithmic Expected Time*', ACM Transactions on Mathematical Software, Vol. 3, No. 3, pp. 209-226, 1977.
- [GG 98] Gaede V., Günther O.: '*Multidimensional Access Methods*', ACM Computing Surveys, Vol. 30, No. 2, 1998, pp. 170-231.
- [GIM 99] Gionis A., Indyk P., Motwani R.: '*Similarity Search in High Dimensions via Hashing*', Proc. 25th Int. Conf. on Very Large Data Bases, Edinburgh, GB, pp. 518-529, 1999.
- [Gut 84] Guttman A.: '*R-trees: A Dynamic Index Structure for Spatial Searching*', Proc. ACM SIGMOD Int. Conf. on Management of Data, Boston, MA, pp. 47-57, 1984.
- [Hen 94] Henrich, A.: '*A distance-scan algorithm for spatial access structures*', Proceedings of the 2nd ACM Workshop on Advances in Geographic Information Systems, ACM Press, Gaithersburg, Maryland, pp. 136-143, 1994.
- [Hen 98] Henrich, A.: '*The LSD<sup>h</sup>-tree: An Access Structure for Feature Vectors*', Proc. 14th Int. Conf. on Data Engineering, Orlando, 1998.

108





## Literature

- [HS 95] Hjaltason G. R., Samet H.: '*Ranking in Spatial Databases*', Proc. 4th Int. Symp. on Large Spatial Databases, Portland, ME, pp. 83-95, 1995.
- [HSW 89] Henrich A., Six H.-W., Widmayer P.: '*The LSD-Tree: Spatial Access to Multidimensional Point and Non-Point Objects*', Proc. 15th Conf. on Very Large Data Bases, Amsterdam, The Netherlands, pp. 45-53, 1989.
- [Jag 91] Jagadish H. V.: '*A Retrieval Technique for Similar Shapes*', Proc. ACM SIGMOD Int. Conf. on Management of Data, pp. 208-217, 1991.
- [JW 96] Jain R, White D.A.: '*Similarity Indexing: Algorithms and Performance*', Proc. SPIE Storage and Retrieval for Image and Video Databases IV, Vol. 2670, San Jose, CA, pp. 62-75, 1996.
- [KF 94] Kamel I, Faloutsos C.: '*Hilbert R-tree: An Improved R-tree using Fractals*'. Proc. 20th Int. Conf. on Very Large Databases, 1994, pp. 500-509.
- [KS 97] Katayama N., Satoh S.: '*The SR-tree: An Index Structure for High-Dimensional Nearest Neighbor Queries*', Proc. ACM SIGMOD Int. Conf. on Management of Data, pp. 369-380, 1997.
- [KSF+ 96] Korn F., Sidiropoulos N., Faloutsos C., Siegel E., Protopapas Z.: '*Fast Nearest Neighbor Search in Medical Image Databases*', Proc. 22nd Int. Conf. on Very Large Data Bases, Mumbai, India, pp. 215-226, 1996.

109



## Literature

- [LJF 94] Lin K., Jagadish H. V., Faloutsos C.: '*The TV-tree: An Index Structure for High-Dimensional Data*', VLDB Journal, Vol. 3, pp. 517-542, 1995.
- [MG 93] Mehrotra R., Gary J.: '*Feature-Based Retrieval of Similar Shapes*', Proc. 9th Int. Conf. on Data Engineering, 1993.
- [Ore 82] Orenstein J. A.: '*Multidimensional tries used for associative searching*', Inf. Proc. Letters, Vol. 14, No. 4, pp. 150-157, 1982.
- [PM 97] Papadopoulos A., Manolopoulos Y.: '*Performance of Nearest Neighbor Queries in R-Trees*', Proc. 6th Int. Conf. on Database Theory, Delphi, Greece, in: Lecture Notes in Computer Science, Vol. 1186, Springer, pp. 394-408, 1997.
- [RKV 95] Roussopoulos N., Kelley S., Vincent F.: '*Nearest Neighbor Queries*', Proc. ACM SIGMOD Int. Conf. on Management of Data, San Jose, CA, pp. 71-79, 1995.
- [Rob 81] Robinson J. T.: '*The K-D-B-tree: A Search Structure for Large Multidimensional Dynamic Indexes*', Proc. ACM SIGMOD Int. Conf. on Management of Data, pp. 10-18, 1981.
- [RP 92] Ramasubramanian V., Paliwal K. K.: '*Fast k-Dimensional Tree Algorithms for Nearest Neighbor Search with Application to Vector Quantization Encoding*', IEEE Transactions on Signal Processing, Vol. 40, No. 3, pp. 518-531, 1992.

110



## Literature

- [See 91] Seeger B.: 'Multidimensional Access Methods and their Applications', Tutorial, 1991.
- [SK 97] Seidl T., Kriegel H.-P.: 'Efficient User-Adaptable Similarity Search in Large Multimedia Databases', Proc. 23rd Int. Conf. on Very Large Databases (VLDB'97), Athens, Greece, 1997.
- [Spr 91] Sproull R.F.: 'Refinements to Nearest Neighbor Searching in  $k$ -Dimensional Trees', Algorithmica, pp. 579-589, 1991.
- [SRF 87] Sellis T., Roussopoulos N., Faloutsos C.: 'The  $R^+$  Tree: A Dynamic Index for Multi-Dimensional Objects', Proc. 13th Int. Conf. on Very Large Databases, Brighton, England, pp 507-518, 1987.
- [WSB 98] Weber R., Schek H.-J., Blott S.: 'A Quantitative Analysis and Performance Study for Similarity-Search Methods in High-Dimensional Spaces', Proc. Int. Conf. on Very Large Databases, New York, 1998.
- [WJ 96] White D.A., Jain R.: 'Similarity indexing with the  $SS$ -tree', Proc. 12th Int. Conf on Data Engineering, New Orleans, LA, 1996.
- [YY 85] Yao A. C., Yao F. F.: 'A General Approach to  $D$ -Dimensional Geometric Queries', Proc. ACM Symp. on Theory of Computing, 1985.

111



## Acknowledgements

We thank **Stephen Blott** and **Hans-J. Scheck** for the very interesting and helpful discussions about the VA-file.

We thank **Raghu Ramakrishnan** and **Jonathan Goldstein** for the very interesting and helpful comments on their work on "When Is Nearest-Neighbor Meaningful".

Furthermore, we thank **Andreas Henrich** for introducing us into the secrets of LSD and KDB trees.

Finally, we thank **Marco Poetke** for providing the nice figure explaining telescope vectors.

Last but not least, we thank **H.V. Jagadish** for encouraging us to put this tutorial together.

112