



Spatial Databases I

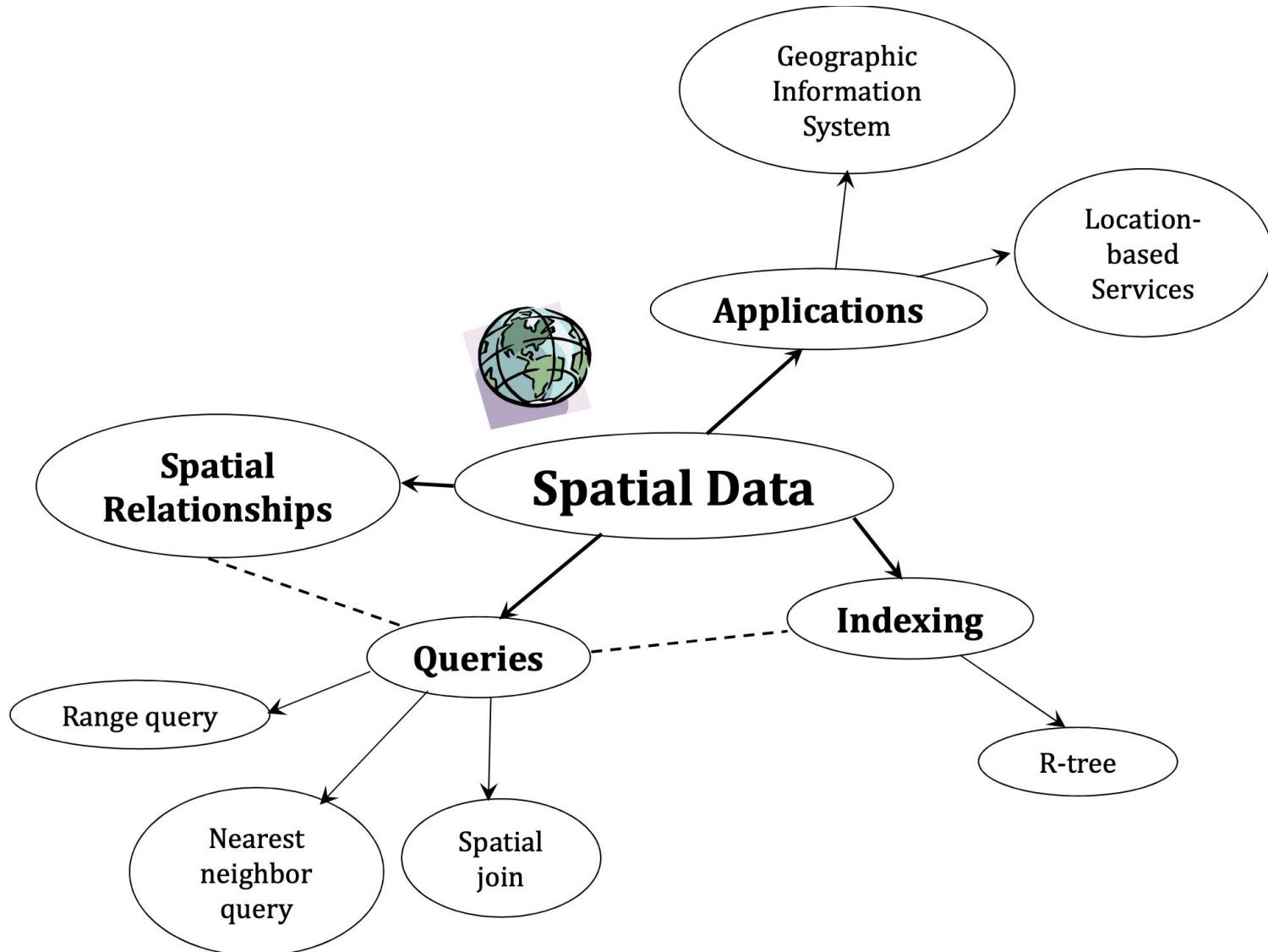
南方科技大学

唐 博

tangb3@sustech.edu.cn



The Concepts





Spatial Data Management



- ❖ Spatial Data, Relationships, and Queries
- ❖ Challenges of Indexing Spatial Data
- ❖ Filter-Refinement Processing
- ❖ The R-tree



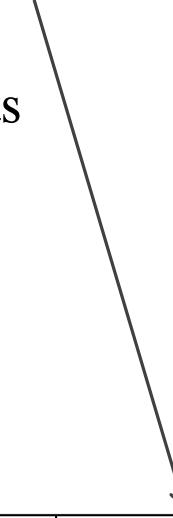
Spatial Data Management



- ❖ A *spatial object* has some spatial attribute that describes its location and geometry
 - ❖ Typically 2-dimensional or 3-dimensional data
- ❖ A *spatial relation* (dataset) is a collection of spatial objects of the same entity (e.g., rivers, cities, roads)
- ❖ What applications use spatial data?
 - ❖ Geographic information system (GIS), location-based services (LBS), computer-aided design (CAD)



Road segments of a region



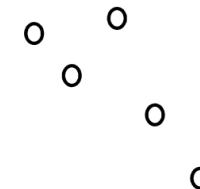
ID	Name	Type	Polyline
1	Boulevard	avenue	(10023,1094), (9034,1567), (9020,1610)
2	Leeds	highway	(4240,5910), (4129,6012), (3813,6129), (3602,6129)
...

A spatial relation

Types of Spatial Data

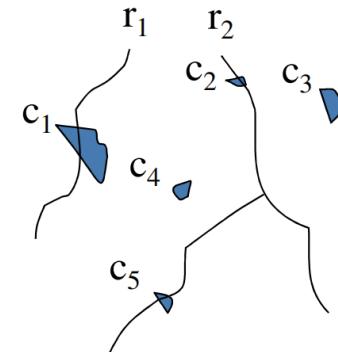
❖ Points

- ❖ Cities in a large-scale map

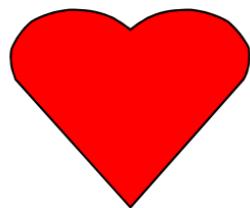


❖ Objects with extent

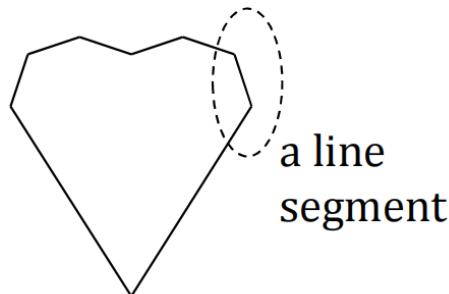
- ❖ include rivers, forests, districts, buildings in a city, etc



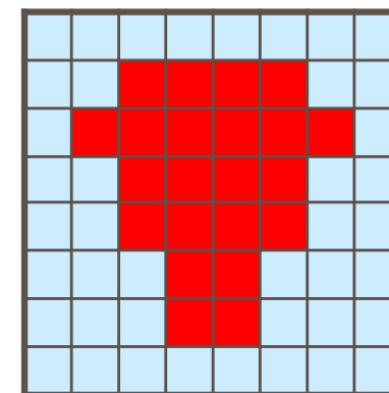
- ❖ **Vector representation:** approximation by geometric objects, e.g., polygons, poly-lines, rectangles, etc.
- ❖ **Raster representation:** divide the space by a grid and approximate the objects by a set of pixels in this grid



actual object



vector approximation



raster approximation

Vector Representation

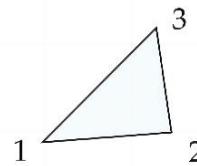
Line Segments and Polygons

line segment



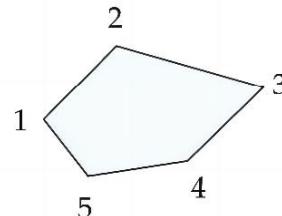
$\{(x_1, y_1), (x_2, y_2)\}$

triangle



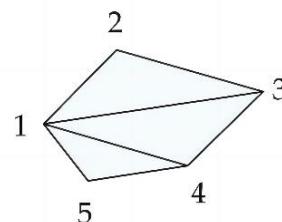
$\{(x_1, y_1), (x_2, y_2), (x_3, y_3)\}$

polygon



$\{(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4), (x_5, y_5)\}$

polygon



$\{(x_1, y_1), (x_2, y_2), (x_3, y_3), \text{ID1}\}$
 $\{(x_1, y_1), (x_3, y_3), (x_4, y_4), \text{ID1}\}$
 $\{(x_1, y_1), (x_4, y_4), (x_5, y_5), \text{ID1}\}$

object

representation



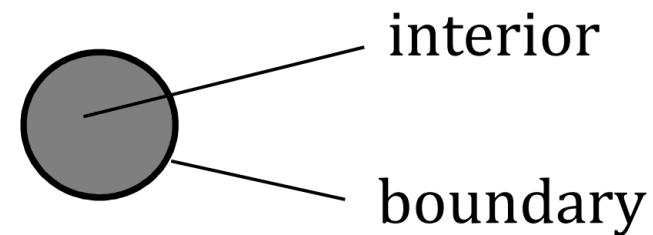
Spatial Relationships



- ❖ A **spatial relationship** links two objects based on their relative location and extent in space
 - ❖ Example: “My house is **close to** Central Park”
 - ❖ Different from the term “spatial relation”, which means a dataset of spatial objects

- ❖ Types of spatial relationships
 - ❖ topological relationships
 - ❖ distance relationships
 - ❖ directional relationships

- ❖ Each object is characterized by the space it occupies
 - ❖ a set of pixels (finite or infinite)
- ❖ Each object has a boundary and an interior
 - ❖ **boundary**: the set of object pixels that are adjacent to at least one pixel not occupied by the object
 - ❖ **interior**: the set of object pixels which are not boundary
- ❖ Some objects may not have interior
 - ❖ e.g., points, line segments





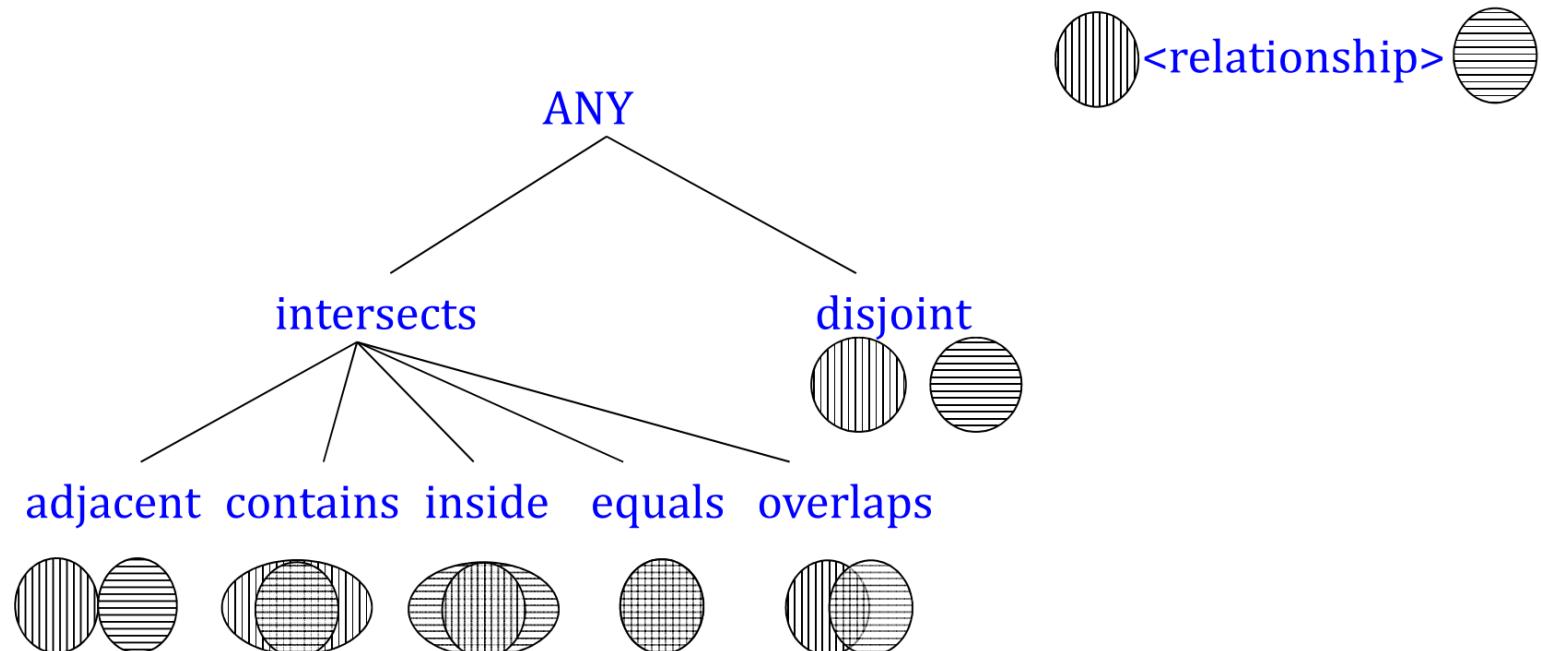
Topological Relationships



- ❖ A topological relationship between two objects is defined by a set-based relationship between their boundaries and interiors
 - ❖ E.g., o_1 is inside o_2 if $\text{interior}(o_1) \subset \text{interior}(o_2)$
- ❖ intersects includes: equals, inside, contains, adjacent, overlaps
- ❖ intersects $\Leftrightarrow \neg$ disjoint

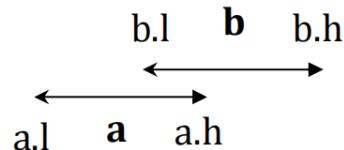
Topological relationship	equivalent boundary/interior relationships
$\text{disjoint}(o_1, o_2)$	$(\text{interior}(o_1) \cap \text{interior}(o_2) = \emptyset) \wedge (\text{boundary}(o_1) \cap \text{boundary}(o_2) = \emptyset)$
$\text{intersects}(o_1, o_2)$	$(\text{interior}(o_1) \cap \text{interior}(o_2) \neq \emptyset) \vee (\text{boundary}(o_1) \cap \text{boundary}(o_2) \neq \emptyset)$
$\text>equals}(o_1, o_2)$	$(\text{interior}(o_1) = \text{interior}(o_2)) \wedge (\text{boundary}(o_1) = \text{boundary}(o_2))$
$\text{inside}(o_1, o_2)$	$\text{interior}(o_1) \subset \text{interior}(o_2)$
$\text{contains}(o_1, o_2)$	$\text{interior}(o_2) \subset \text{interior}(o_1)$
$\text{adjacent}(o_1, o_2)$ (or $\text{meets}(o_1, o_2)$)	$(\text{interior}(o_1) \cap \text{interior}(o_2) = \emptyset) \wedge (\text{boundary}(o_1) \cap \text{boundary}(o_2) \neq \emptyset)$
$\text{overlaps}(o_1, o_2)$	$(\text{interior}(o_1) \cap \text{interior}(o_2) \neq \emptyset) \wedge (\exists p \in o_1 : p \notin \text{interior}(o_2) \wedge p \notin \text{boundary}(o_2))$ $\wedge (\exists p \in o_2 : p \notin \text{interior}(o_1) \wedge p \notin \text{boundary}(o_1))$

Hierarchy of topological relationships

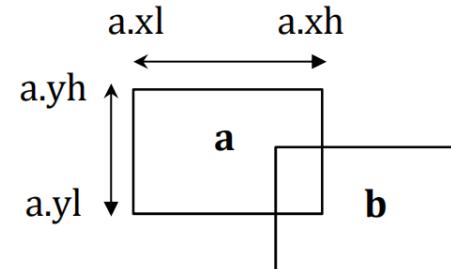


Intersection of Rectangles

- ❖ **intersects:** frequently-used topological relationship
- ❖ **rectangles:** used to represent objects/queries
- ❖ *Question:* how to check the “intersects” predicate without enumerating pixels?



Intersection condition for 1-dimensional intervals:
 $(b.l \leq a.h) \text{ and } (a.l \leq b.h)$



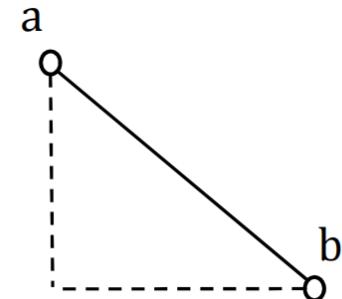
Intersection condition for 2-dimensional rectangles:
 $(b.xl \leq a.xh) \text{ and } (a.xl \leq b.xh) \text{ and }$
 $(b.yl \leq a.yh) \text{ and } (a.yl \leq b.yh)$

Distance Relationships

- ❖ Distance relationships associate two objects based on how far they are

- ❖ We focus on the **Euclidean distance**

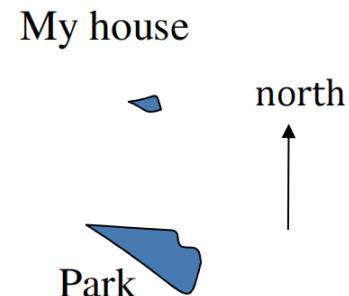
$$\diamond \quad dist(a, b) = \sqrt{(a.x - b.x)^2 + (a.y - b.y)^2}$$



- ❖ Another distance measure: the shortest path distance on a road network

- ❖ Directional relationships associate two objects based on their relative **orientation** according to a global reference system
 - ❖ My house is **north of** the park

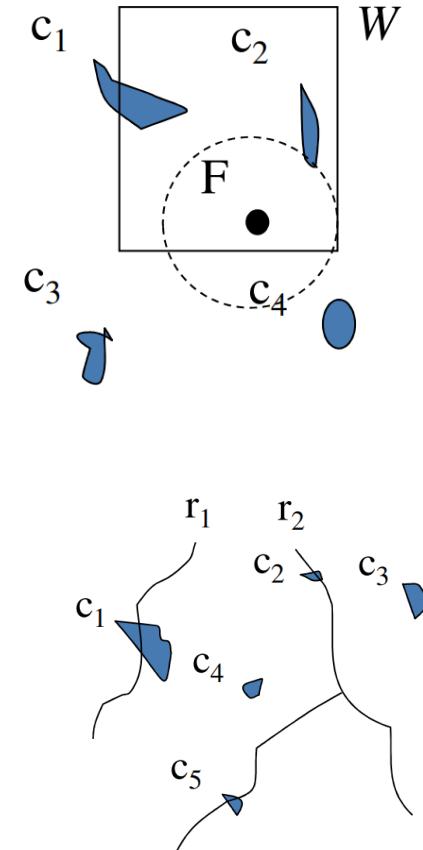
- ❖ Examples of directional relationships:
 - ❖ north, south, east, west, etc.
 - ❖ left, right, above, below, front, behind, etc.



- ❖ Combinations of topological, distance, and directional relationships can be used to describe a pair of spatial objects:
 - ❖ My house is **disjoint** with the park, **100 meters north of** it

Spatial Queries

- ❖ Input: spatial relation(s)
- ❖ Output: objects (or combinations of objects) satisfying some spatial relationships
 - ❖ with a reference query object, or
 - ❖ between two spatial datasets
- ❖ **Range query** (spatial selection, window query)
 - ❖ e.g., find all cities that intersect window W
 - ❖ *Answer set*: {c1, c2}
- ❖ **Nearest neighbor query**
 - ❖ e.g., find the city closest to the forest F
 - ❖ *Answer*: c2
- ❖ **Spatial join**
 - ❖ e.g., find all pairs of cities and rivers that intersect
 - ❖ *Answer set*: {(r1,c1), (r2,c2), (r2,c5)}





History

- ❖ Spatial index developed for range query
 - ❖ 1974-1975: point access method (quad-tree, kd-tree)
 - ❖ 1981: Space ordering method (B+-tree over the Z-order)
 - ❖ 1984: R-tree
 - ❖ 1990: R*-tree
- ❖ Query types on R-trees
 - ❖ 1984: Range query
 - ❖ 1993: Spatial join
 - ❖ 1995: Nearest neighbor query
 - ❖ 1998: Closest pair



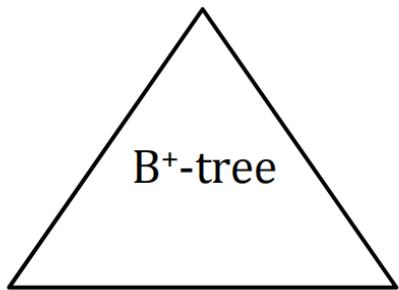
History



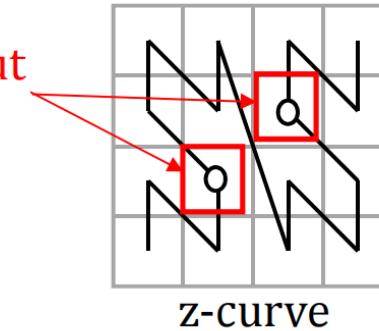
Commercial products

- ❖ Oracle Spatial (and Graph)
 - ❖ 1995: Point data only
 - ❖ 1996-2000: Lines, polygons,
 - ❖ 2002: Quadtree, R-tree index
 - ❖ https://docs.oracle.com/cd/B10500_01/appdev.920/a96630/sdo_index_query.htm
- ❖ SQL Server, Spatial Support
 - ❖ 2008: spatial data types (for flat/round Earth models)
 - ❖ 2012: geometry hierarchy, complicated curves
 - ❖ Index: B+-tree over grid index
 - ❖ <https://msdn.microsoft.com/en-us/library/bb895265.aspx>

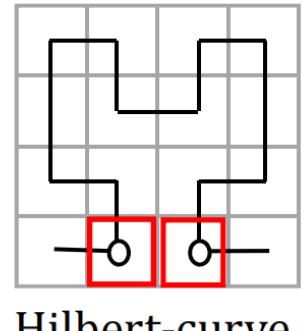
Space Ordering Method



Close in space, but
far away in curve



z-curve



Hilbert-curve

- ❖ Space ordering method
 - ❖ Use a **space filling curve** to convert each point to a value
 - ❖ Index the values of these points by a B+-tree
 - ❖ Range query: Decompose the query into multiple interval queries
- ❖ Problems
 - ❖ Dimensionality: No space ordering preserves spatial proximity
 - ❖ Complex spatial extent: Redundancy for assigning objects into disk pages

- ❖ Point access methods divide the space into disjoint regions and group the points by their regions
 - ❖ Grid-file
 - ❖ kd-tree
 - ❖ Quad-tree

Example of a grid-file:

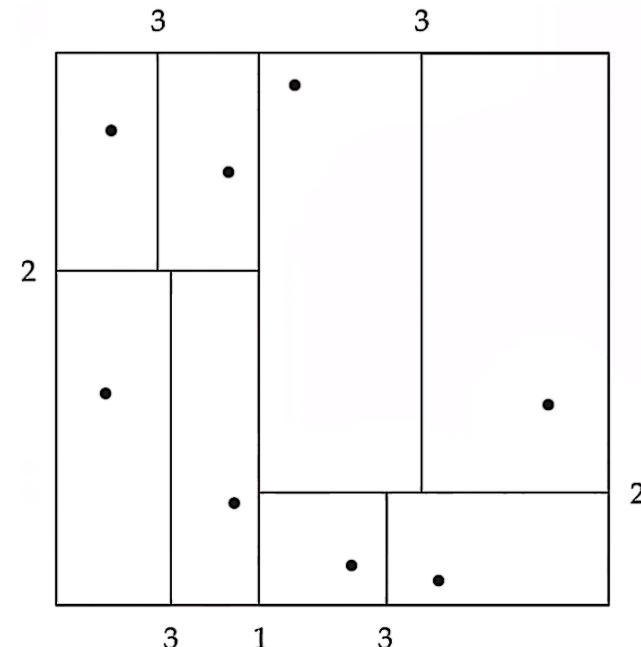
o	o	o	
oo		o	
	o	o	o



disk block 1
disk block 2
disk block 3

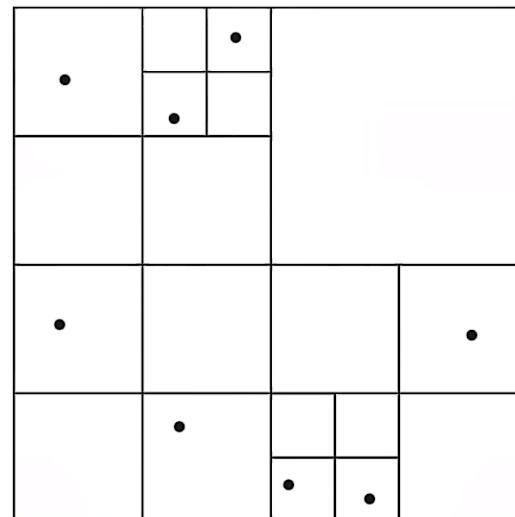
k-d tree

- ❖ Each tree level partitions the space into two regions with
 - ❖ (almost) equal number of points
 - ❖ choose one dimension for partitioning at the root level
 - ❖ choose another dimension for partitioning at the next level and so on
 - ❖ Partitioning stops when a node has less than a fixed maximum number of points
-
- ❖ Example figure of k-d tree
 - ❖ Each line (except the outside box) corresponds to a tree node
 - ❖ Each leaf node contains 1 point
 - ❖ Tree level: the numbering of the lines

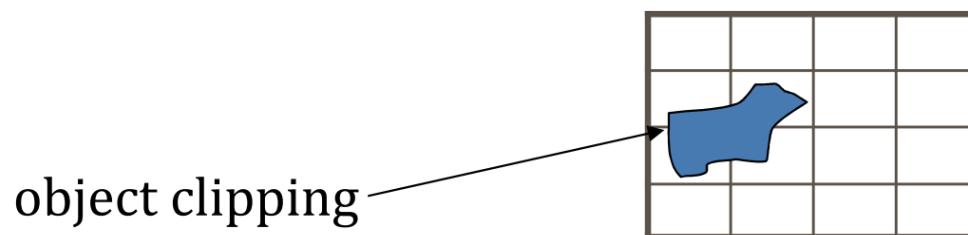


Quadtrees

- ❖ The top node is associated with the entire space
- ❖ Each non-leaf node divides its region into four rectangles with equal area
 - ❖ Each such node has four corresponding child
- ❖ Partitioning stops when a node has less than a fixed maximum number of points



- ❖ Extensions of $k-d$ trees and quadtrees can be used to index line segments and polygons
 - ❖ Require splitting segments/polygons into pieces at partitioning boundaries
 - ❖ The same segment/polygon may be stored at several leaf nodes
- ❖ Drawback: point access methods are not effective for extended objects
 - ❖ Objects may be *clipped* into several parts → data redundancy and affects performance badly





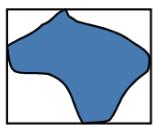
Spatial Access Methods



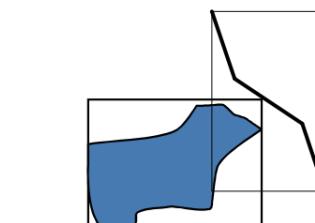
- ❖ Large spatial dataset, stored on disk
 - ❖ Query processing: I/O time is much larger than CPU time
- ❖ Spatial Access Methods (SAM): Indexing spatial data
 - ❖ aim at minimizing the expected I/O cost
- ❖ R-tree is the most popular spatial access method
 - ❖ Use the filter-refinement approach for processing queries

- ❖ It is slow to evaluate spatial relationships on geometric data; thus, we approximate a *spatial object* by its **minimum bounding rectangle** (MBR)
 - ❖ The spatial query is processed in two steps:
 1. Filter step: test the MBR against the query predicate
 2. Refinement step: [if passed the filter step]
then test the exact geometry of objects
- Fast!
- ↓ Slow!

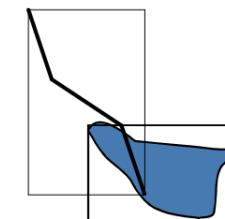
Example for spatial **intersection** joins:



filtered pair

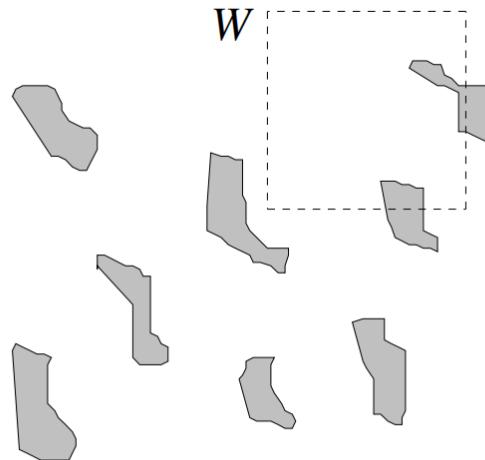


non-qualifying pair that
passes the filter step
(false hit)

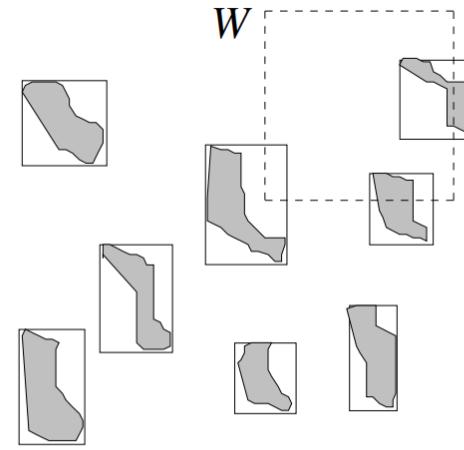


qualifying pair

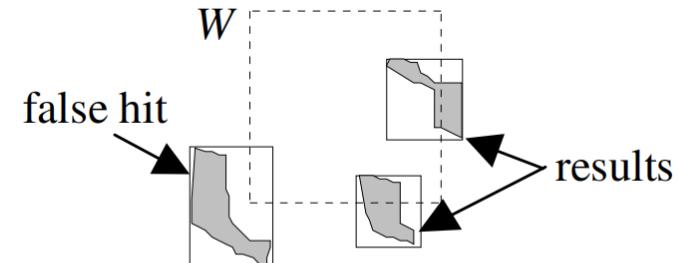
- ❖ Problem: Find all objects that intersect the query window W
- ❖ Comparison cost \approx number of examined vertices
- ❖ Suppose that each polygon has 100 vertices
 - ❖ Comparisons in Fig. (a) = $(100+4)*8 = 832$
 - ❖ Comparisons in Fig. (b) = $(4+4)*8 + (100+4)*3 = 376$



(a) objects and a query



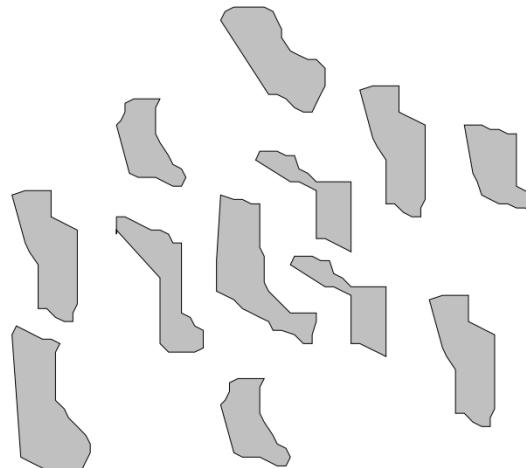
(b) object MBRs



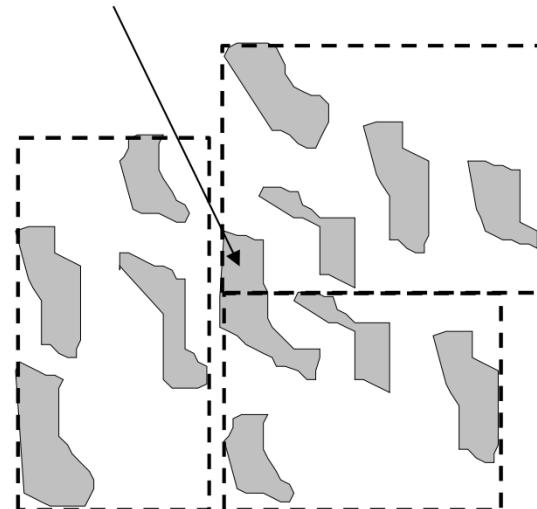
(c) candidates and results

❖ Problem:

- ❖ Group the objects below into 3 groups of 4 objects each such that the MBRs of the groups do not overlap
- ❖ not feasible in this example
- ❖ Avoid object clipping by allowing regions of object groups to overlap

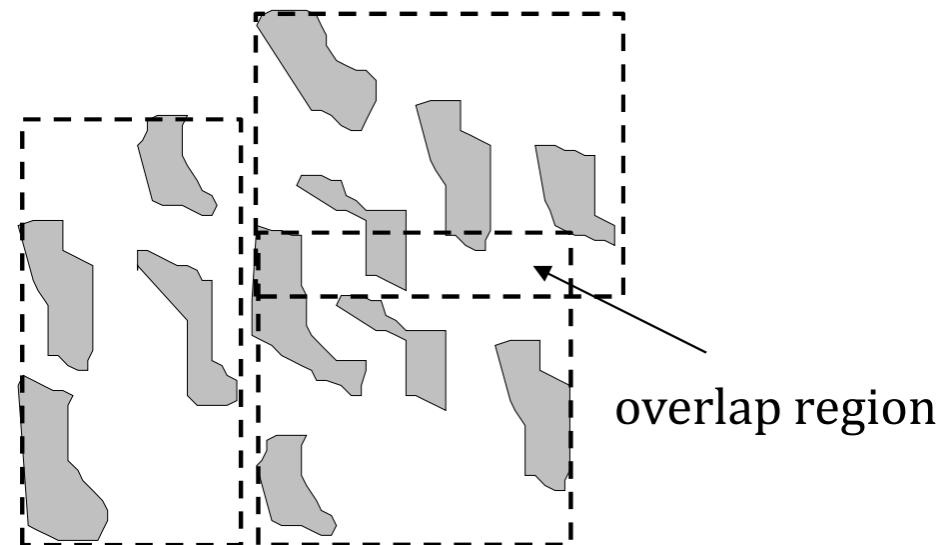
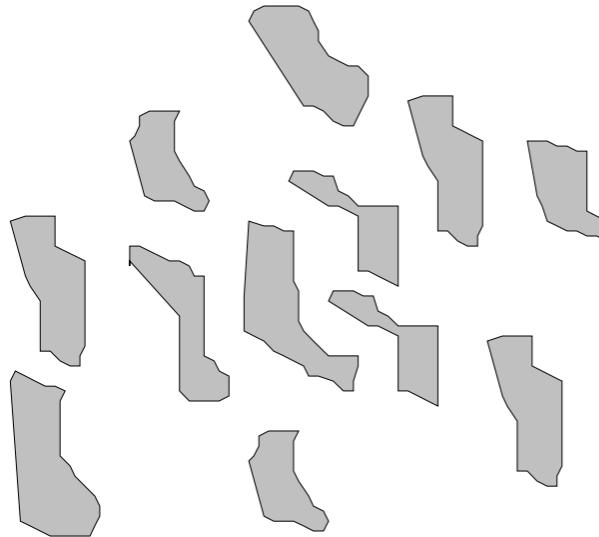


clipped object



❖ Revised problem:

- ❖ Group the objects below into 3 groups of 4 objects each such that the MBRs of the groups have the **minimum overlap**
- ❖ Hard optimization problem
- ❖ Use a fast, suboptimal solution instead (discussed later)





Spatial Indexing Methods



<i>Indexing method</i>	Space ordering method (B ⁺ -tree over space filling curve)	Point access methods (k-D tree, quad tree)	R-tree
<i>Point object</i>	yes	yes	yes
<i>Object with extent</i>	no	Use clipping (but it leads to redundant data and high query cost)	yes
<i>Point query</i>	yes	yes	yes
<i>Range query</i>	High cost	reasonable	Low cost



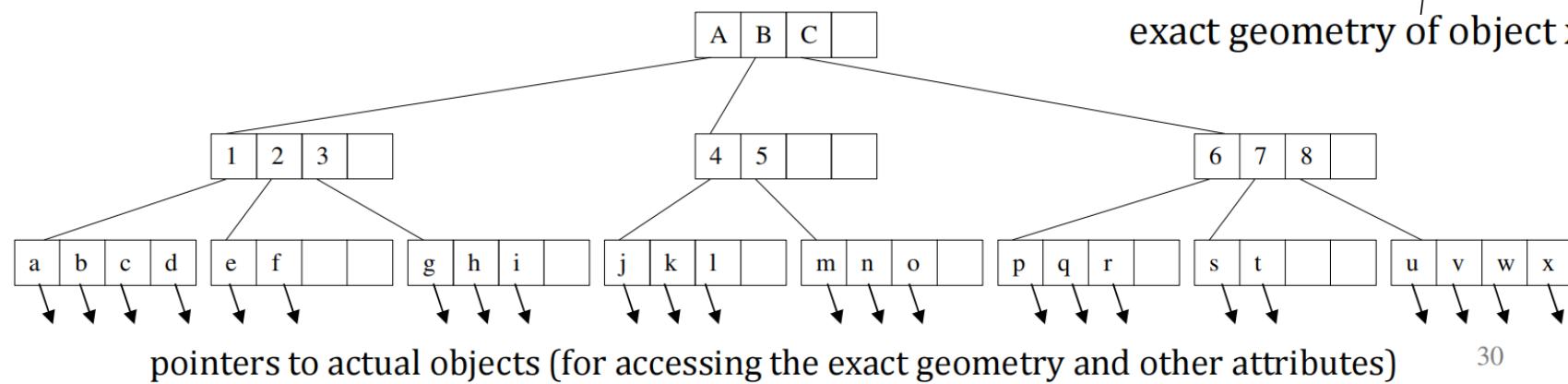
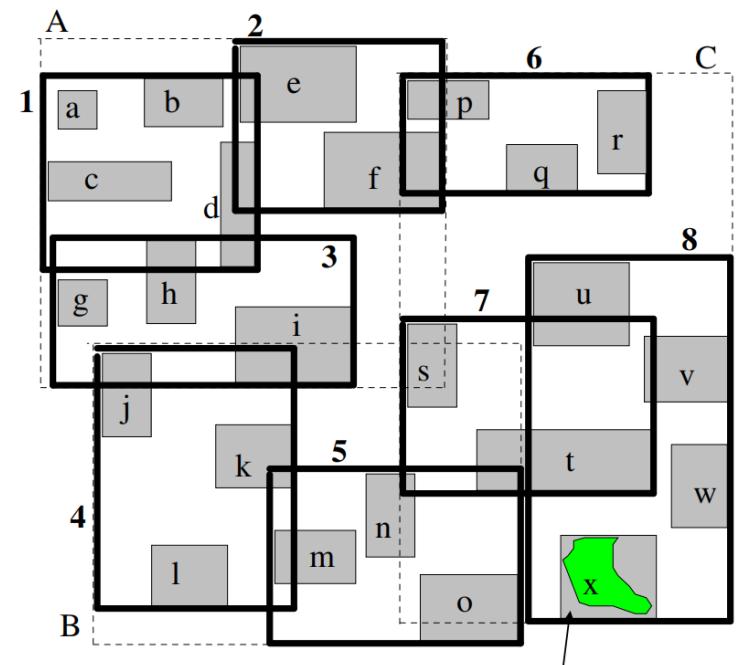
The R-tree



- ❖ Groups **object MBRs** to disk blocks
 - ❖ Each group of objects (a disk block) is a leaf of the tree
 - ❖ Objects in the same leaf node should be close together
- ❖ The MBRs of the leaf nodes are grouped to form nodes at the next level
 - ❖ MBRs located close together are grouped
- ❖ Grouping is recursively applied at each level until a single group (the root) is formed

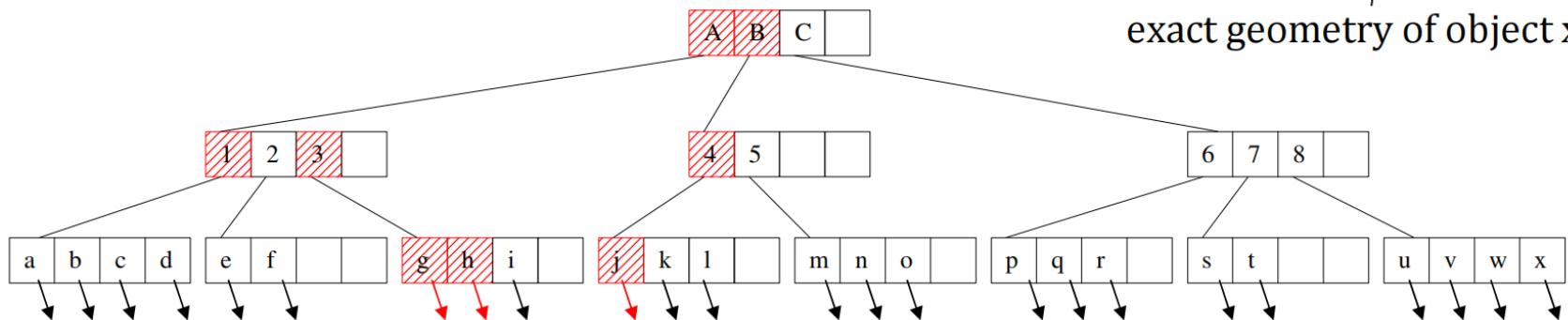
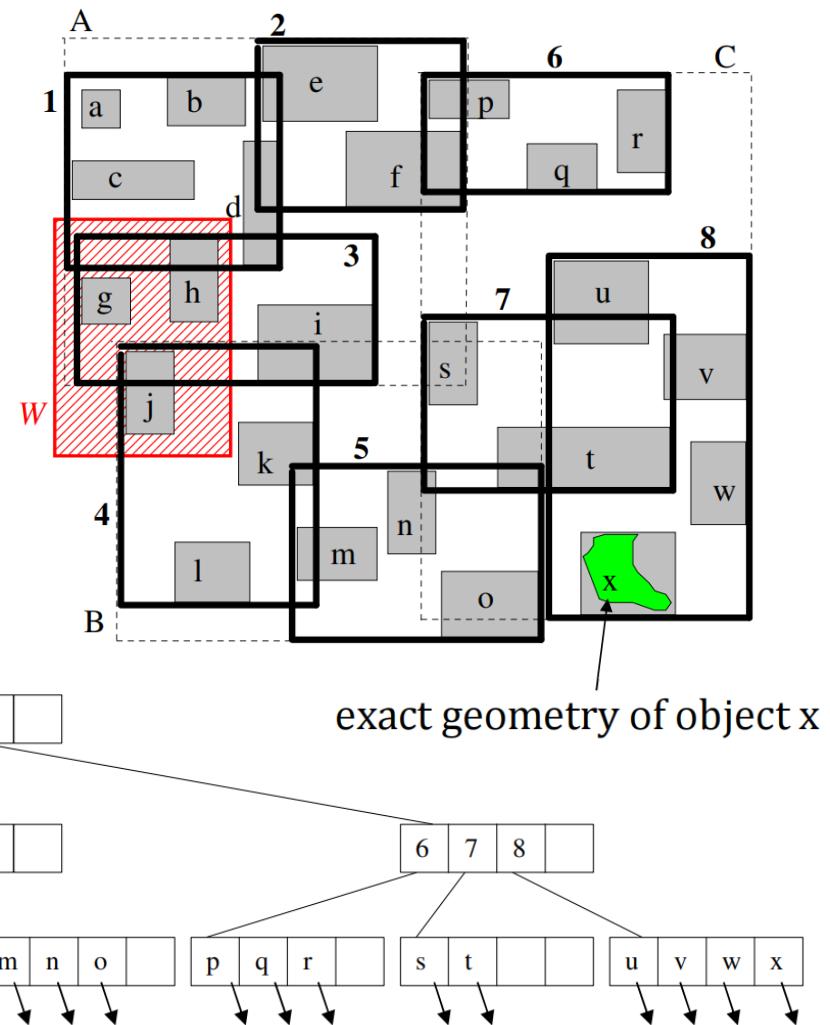
The R-tree

- ❖ Leaf node entries
 - ❖ <MBR, object-id>
- ❖ Non-leaf node entries
 - ❖ <MBR, ptr>
 - ❖ Its MBR = MBR of all entries in its child node
- ❖ Properties like B+-tree
 - ❖ balanced tree
 - ❖ 1 node → 1 disk block
 - ❖ Minimum fanout: 0.4 of capacity



Searching

- ❖ W is a range query
(window intersection query)
- ❖ Which nodes are visited during search?





Range searching on an R-tree



❖ Range_Query(query W , R-tree node n):

If n is not a leaf node

For each index entry e in n such that $e.\text{MBR}$ intersects W

visit node n' pointed by $e.ptr$;

Range_Query(W, n') ;

If n is a leaf node

For each index entry e in n such that $e.\text{MBR}$ intersects W

visit object o pointed by $e.object-id$;

test range query against exact geometry of o ;

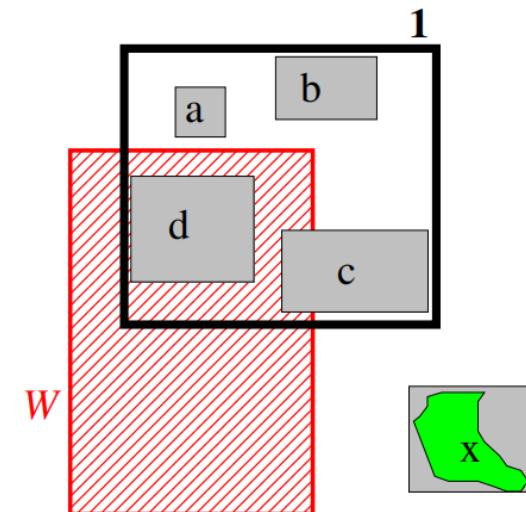
if o intersects W , then report o ;

❖ The search starts from the root node

❖ May follow *multiple paths* during search

Range search with other predicates

- ❖ The search algorithm considers only the intersection predicate
 - ❖ I.e., find all objects intersecting the query window W
- ❖ How to modify the algorithm for other search predicates?
- ❖ E.g., find all objects that are **contained** in W
 - ❖ Search predicate for an object o
 - ❖ $\text{contains}(W, o)$
 - ❖ Search predicate for an object's MBR R_o
 - ❖ $\text{contains}(W, R_o)$. Why?
 - ❖ Search predicate for a non-leaf MBR R_e
 - ❖ $\text{intersects}(W, R_e)$. Why?



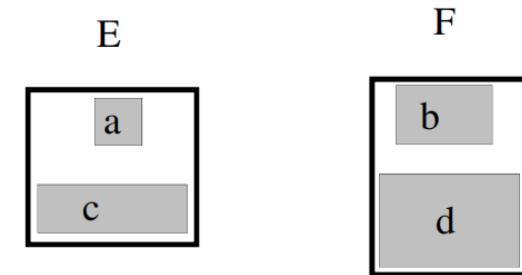
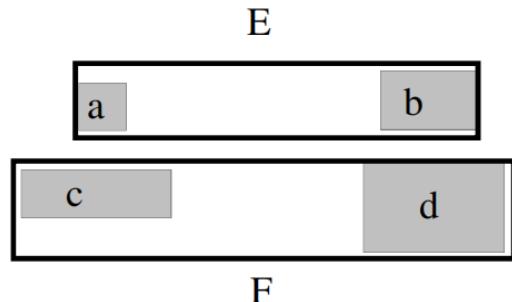


Construction of the R-tree

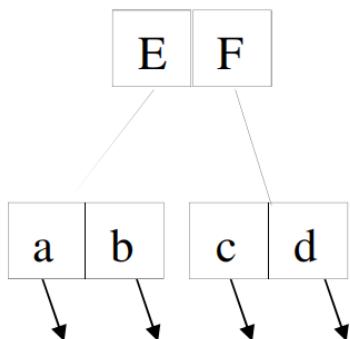


- ❖ Insertion similar to B+-tree
 - ❖ In addition, special optimization algorithms are designed for
 - ❖ choosing the path where a new MBR is inserted
 - ❖ splitting overflowed nodes
- ❖ Underflows in deletions are handled by
 - ❖ deleting the underflowed leaf node
 - ❖ re-inserting the remaining entries

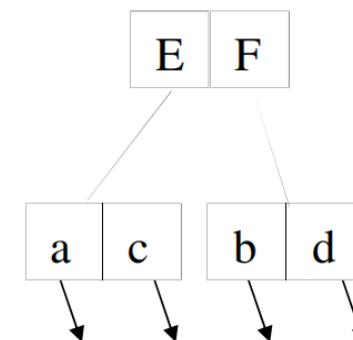
Which R-tree is better? Why?



Tree #1



Tree #2





R*-tree

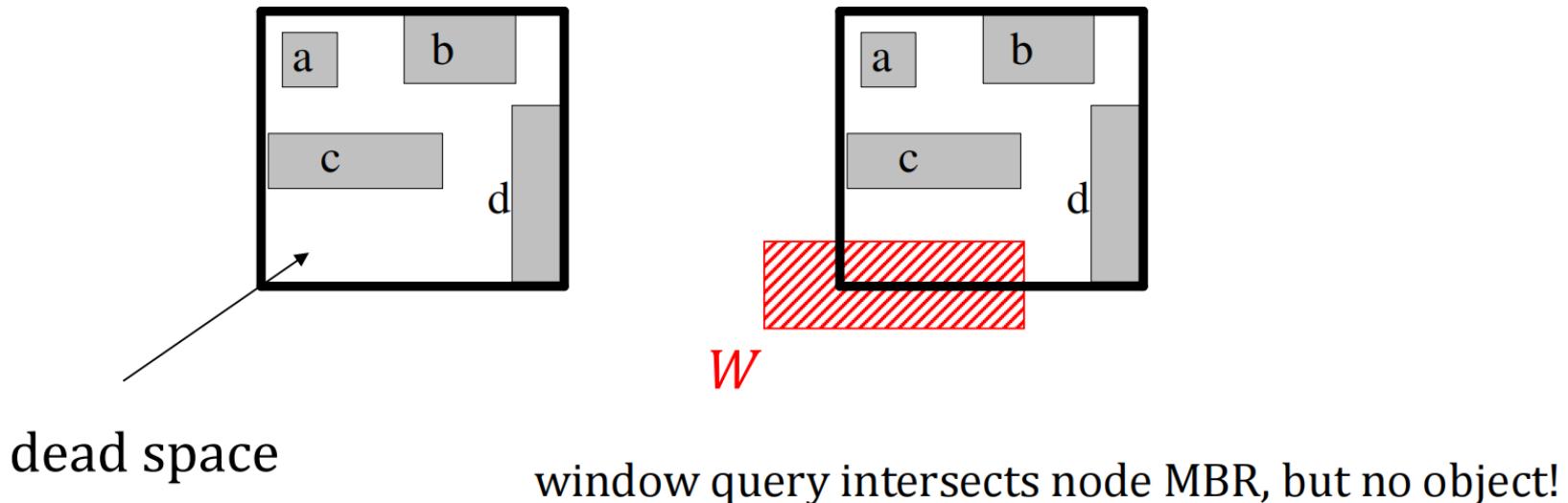


An optimized version of R-tree

- ❖ The R-tree and the R*-tree differ only in the insertion algorithm
- ❖ The improved insertion algorithm aims at constructing a tree of high quality
- ❖ A “good” tree should have
 - ❖ nodes with small MBRs
 - ❖ nodes with small overlap
 - ❖ nodes that look like squares
 - ❖ nodes as full as possible

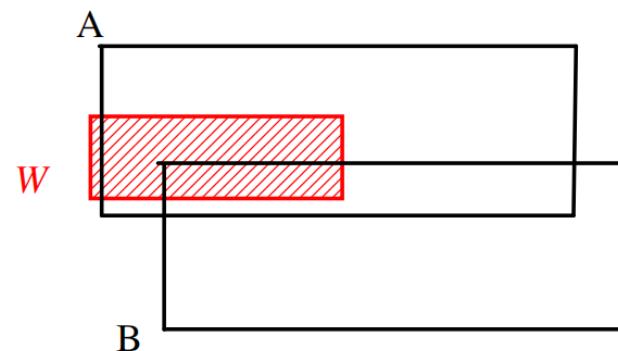
Optimization Criteria: Small Area

- ❖ Minimize the area covered by an index rectangle
 - ❖ small area means small dead space



Optimization Criteria: Small Overlap

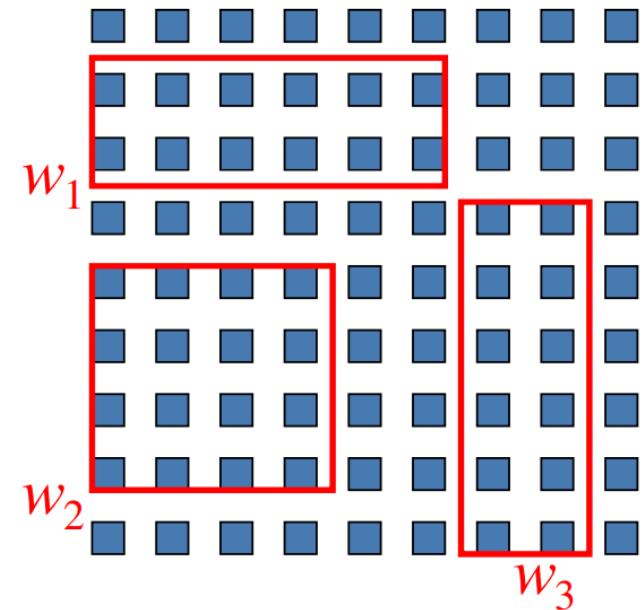
- ❖ Minimize overlap between node MBRs
 - ❖ Minimizes the number of traversed paths



Both nodes intersect query!

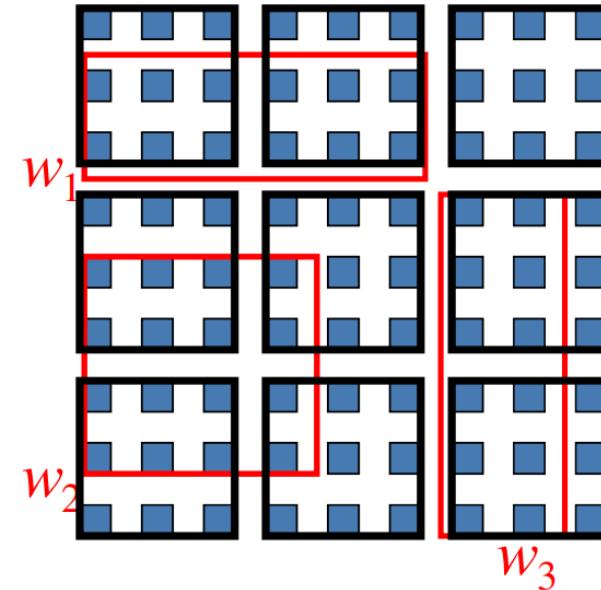
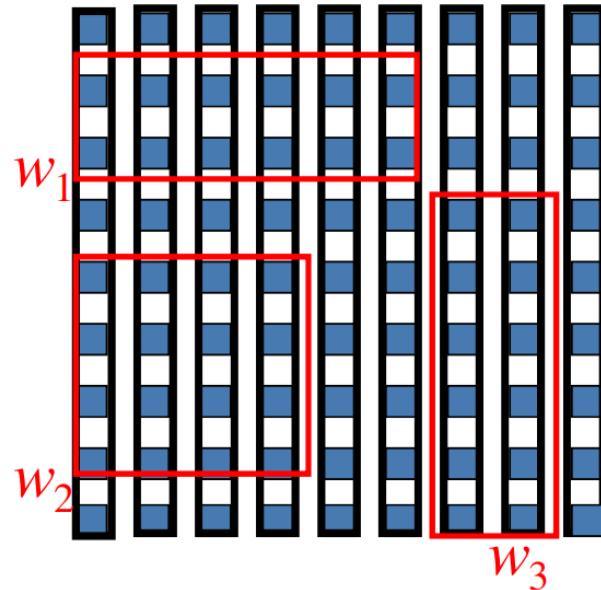
- ❖ Minimize the margins of node MBRs
 - ❖ Square-like nodes, smaller number of intersections for a random query, better structure

- ❖ Example: assign the rectangles (objects) into nodes (groups) of 9 rectangles,
 - ❖ in a way such that the expected number of group MBRs intersected by a (random) query is minimized



Effect of Margins Minimization

- ❖ Grouping 1 (minimal areas of group MBRs)
 - ❖ bad grouping for queries w_1, w_2
- ❖ Grouping 2 (minimization of margins)
 - ❖ minimizes the expected number of groups touching a random query





Optimization Criteria

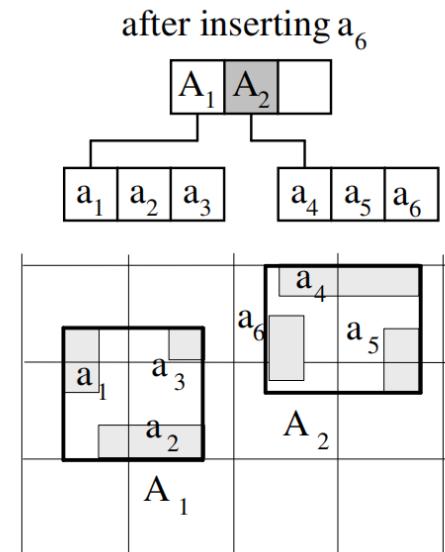
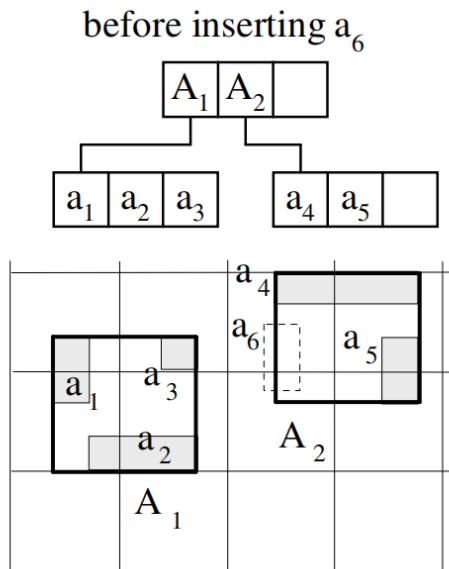


High Storage Utilization

- ❖ Optimize the storage utilization
 - ❖ Nodes in tree should be filled as much as possible
 - ❖ Minimizes tree height and potentially decreases dead space
 - ❖ This criteria may be in conflict with the first optimization criteria

Insertions in an R-tree

- ❖ What happens during an insertion?
 - ❖ Find an appropriate leaf node to store the new entry <MBR,object-id>
 - ❖ For all nodes along the path from the root to the leaf node, adjust (if necessary) their MBRs to fit the new MBR





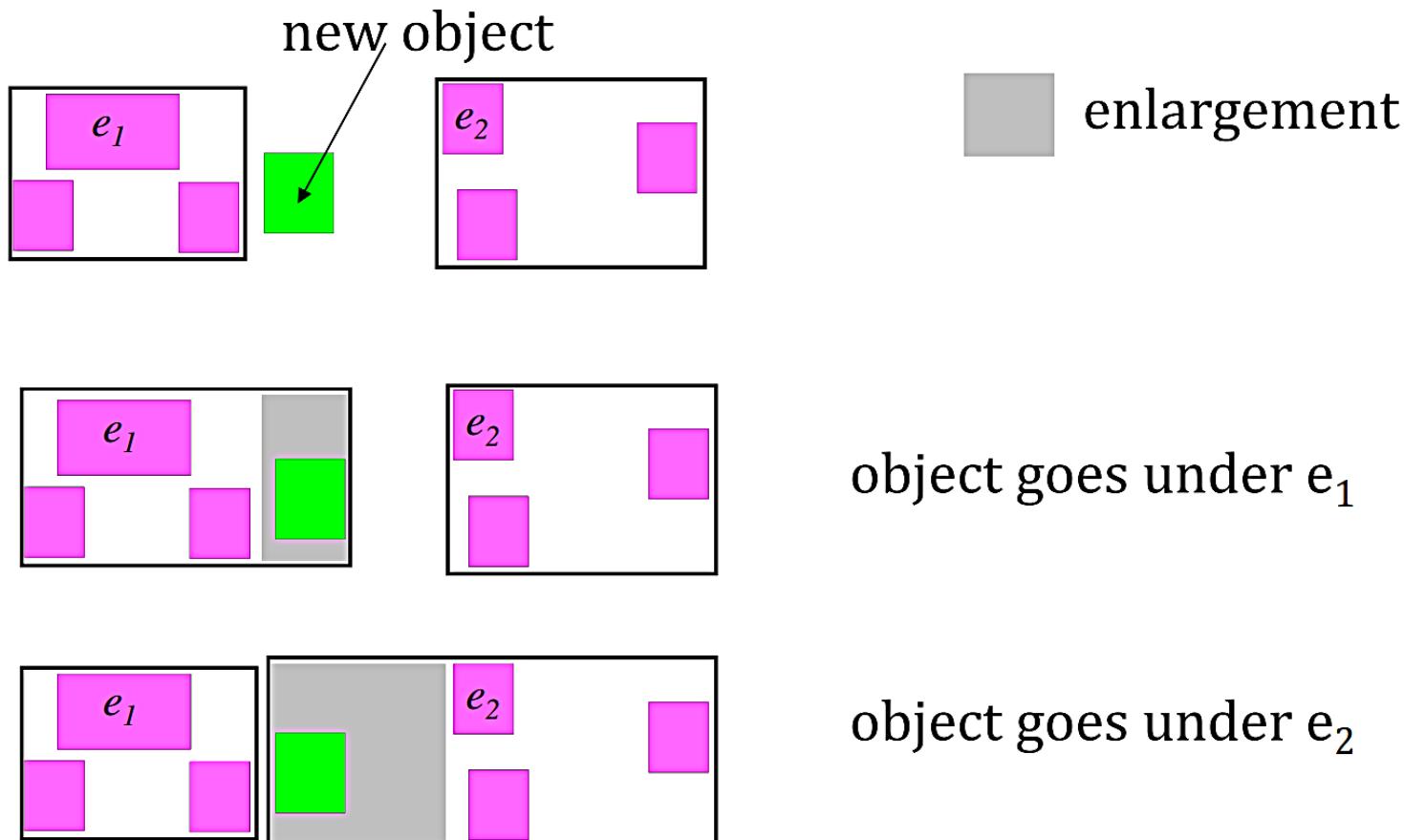
Insertion heuristics



- ❖ When inserting a new entry e into the tree, we follow one path from the root to a leaf
 - ❖ The entry is then inserted to the leaf
- ❖ Issue: How to choose the path?
- ❖ Let N be the current node, for inserting the new entry e :
 - ❖ Follow subtree pointed by entry e' of N :
 - ❖ whose MBR is **enlarged** the least after insertion
 - ❖ whose MBR enlargement will cause the **minimum overlap** with other entries of the same node
 - ❖ break any ties by choosing MBR with the **minimum area**

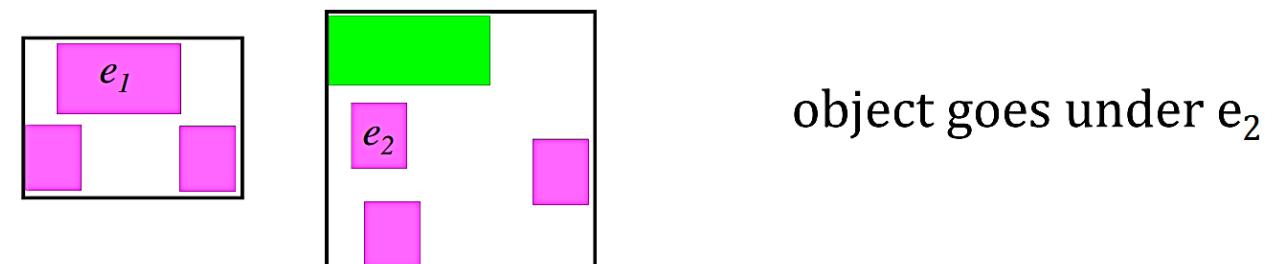
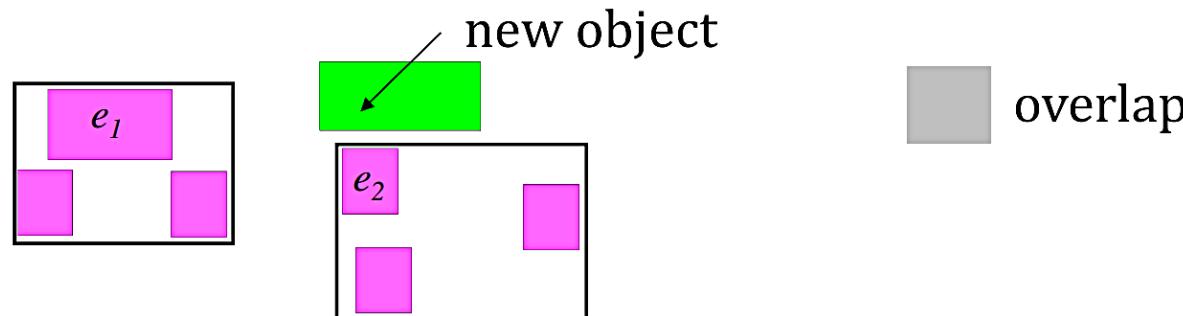
Insertion heuristics

MBR Enlargement



Insertion heuristics

MBR Overlap





Summary



- ❖ Spatial Data, Relationships, and Queries
 - ❖ Topological, distance, directional
- ❖ Challenges of Indexing Spatial Data
- ❖ Filter-Refinement Processing
- ❖ The R-tree
 - ❖ Tree structure
 - ❖ Searching, insertion



谢谢！

DBGroup @ SUSTech

Dr. Bo Tang (唐博)

tangb3@sustech.edu.cn

