



Spatial and Geographic Databases

- Spatial databases store information related to spatial locations, and support efficient storage, indexing and querying of spatial data.
- Special purpose index structures are important for accessing spatial data, and for processing spatial join queries.
- **Computer Aided Design (CAD)** databases store design information about how objects are constructed E.g.: designs of buildings, aircraft, layouts of integrated-circuits
- Geographic databases store geographic information (e.g., maps): often called **geographic information systems** or **GIS**.

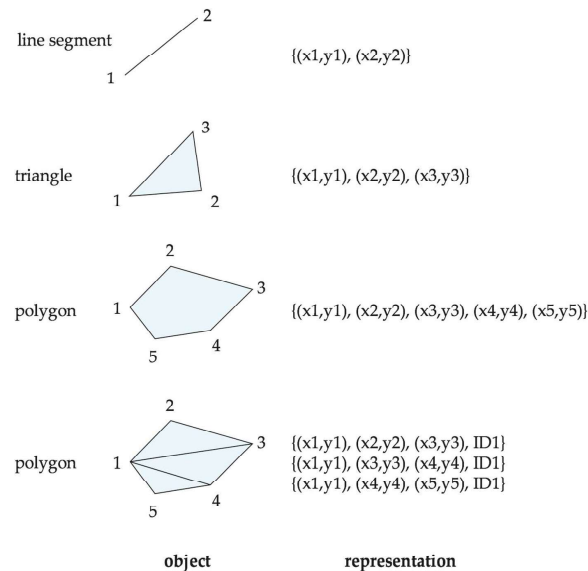


Represented of Geometric Information

- Various geometric constructs can be represented in a database in a normalized fashion.
- Represent a line segment by the coordinates of its endpoints.
- Approximate a curve by partitioning it into a sequence of segments
 - Create a list of vertices in order, or
 - Represent each segment as a separate tuple that also carries with it the identifier of the curve (2D features such as roads).
- Closed polygons
 - List of vertices in order, starting vertex is the same as the ending vertex, or
 - Represent boundary edges as separate tuples, with each containing identifier of the polygon, or
 - Use **triangulation** — divide polygon into triangles
 - ▶ Note the polygon identifier with each of its triangles.



Representation of Geometric Constructs



Representation of Geometric Information (Cont.)

- Representation of points and line segment in 3-D similar to 2-D, except that points have an extra z component
- Represent arbitrary polyhedra by dividing them into tetrahedrons, like triangulating polygons.
- Alternative: List their faces, each of which is a polygon, along with an indication of which side of the face is inside the polyhedron.

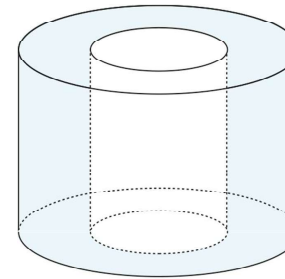


Design Databases

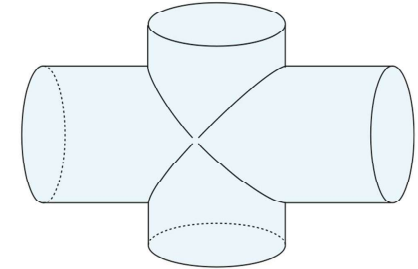
- Represent design components as objects (generally geometric objects); the connections between the objects indicate how the design is structured.
- Simple two-dimensional objects: points, lines, triangles, rectangles, polygons.
- Complex two-dimensional objects: formed from simple objects via union, intersection, and difference operations.
- Complex three-dimensional objects: formed from simpler objects such as spheres, cylinders, and cuboids, by union, intersection, and difference operations.
- Wireframe models represent three-dimensional surfaces as a set of simpler objects.



Representation of Geometric Constructs



(a) Difference of cylinders



(b) Union of cylinders

- Design databases also store non-spatial information about objects (e.g., construction material, color, etc.)
- Spatial integrity constraints are important.
 - E.g., pipes should not intersect, wires should not be too close to each other, etc.



Geographic Data

- **Raster data** consist of bit maps or pixel maps, in two or more dimensions.
 - Example 2-D raster image: satellite image of cloud cover, where each pixel stores the cloud visibility in a particular area.
 - Additional dimensions might include the temperature at different altitudes at different regions, or measurements taken at different points in time.
- Design databases generally do not store raster data.



Geographic Data (Cont.)

- **Vector data** are constructed from basic geometric objects: points, line segments, triangles, and other polygons in two dimensions, and cylinders, spheres, cuboids, and other polyhedrons in three dimensions.
- Vector format often used to represent map data.
 - Roads can be considered as two-dimensional and represented by lines and curves.
 - Some features, such as rivers, may be represented either as complex curves or as complex polygons, depending on whether their width is relevant.
 - Features such as regions and lakes can be depicted as polygons.



Applications of Geographic Data

- Examples of geographic data
 - map data for vehicle navigation
 - distribution network information for power, telephones, water supply, and sewage
- Vehicle navigation systems store information about roads and services for the use of drivers:
 - **Spatial data**: e.g., road/restaurant/gas-station coordinates
 - **Non-spatial data**: e.g., one-way streets, speed limits, traffic congestion
- **Global Positioning System (GPS)** unit - utilizes information broadcast from GPS satellites to find the current location of user with an accuracy of tens of meters.
 - increasingly used in vehicle navigation systems as well as utility maintenance applications.



Spatial Queries

- **Nearness queries** request objects that lie near a specified location.
- **Nearest neighbor queries**, given a point or an object, find the nearest object that satisfies given conditions.
- **Region queries** deal with spatial regions. e.g., ask for objects that lie partially or fully inside a specified region.
- Queries that compute intersections or **unions** of regions.
- **Spatial join** of two spatial relations with the location playing the role of join attribute.



Spatial Queries (Cont.)

- Spatial data is typically queried using a graphical query language; results are also displayed in a graphical manner.
- Graphical interface constitutes the front-end
- Extensions of SQL with abstract data types, such as lines, polygons and bit maps, have been proposed to interface with back-end.
 - allows relational databases to store and retrieve spatial information
 - Queries can use spatial conditions (e.g., contains or overlaps).
 - queries can mix spatial and nonspatial conditions

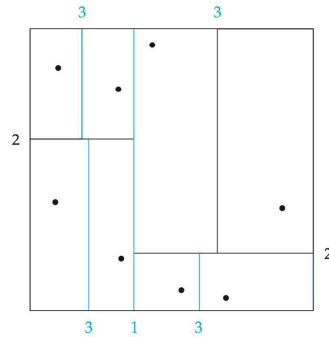


Indexing of Spatial Data

- **k-d tree** - early structure used for indexing in multiple dimensions.
- Each level of a *k-d* tree partitions the space into two.
 - choose one dimension for partitioning at the root level of the tree.
 - choose another dimensions for partitioning in nodes at the next level and so on, cycling through the dimensions.
- In each node, approximately half of the points stored in the sub-tree fall on one side and half on the other.
- Partitioning stops when a node has less than a given maximum number of points.
- The **k-d-B tree** extends the *k-d* tree to allow multiple child nodes for each internal node; well-suited for secondary storage.



Division of Space by a k-d Tree



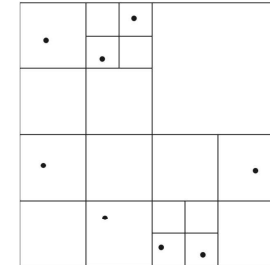
- Each line in the figure (other than the outside box) corresponds to a node in the *k-d* tree.
 - The maximum number of points in a leaf node has been set to 1.
- The numbering of the lines in the figure indicates the level of the tree at which the corresponding node appears.



Division of Space by Quadrees

Quadrees

- Each node of a quadtree is associated with a rectangular region of space; the top node is associated with the entire target space.
- Each non-leaf node divides its region into four equal sized quadrants
 - Correspondingly each such node has four child nodes corresponding to the four quadrants and so on
- Leaf nodes have between zero and some fixed maximum number of points (set to 1 in example).



Quadrees (Cont.)

- **PR quadtree**: stores points; space is divided based on regions, rather than on the actual set of points stored.
- **Region quadrees** store array (raster) information.
 - A node is a leaf node if all the array values in the region that it covers are the same. Otherwise, it is subdivided further into four children of equal area, and is therefore an internal node.
 - Each node corresponds to a sub-array of values.
 - The sub-arrays corresponding to leaves either contain just a single array element, or have multiple array elements, all of which have the same value.
- Extensions of *k-d* trees and PR quadrees have been proposed to index line segments and polygons
 - Require splitting segments/polygons into pieces at partitioning boundaries
 - Same segment/polygon may be represented at several leaf nodes



R-Trees

- **R-trees** are a N-dimensional extension of B⁺-trees, useful for indexing sets of rectangles and other polygons.
- Supported in many modern database systems, along with variants like R⁺-trees and R*-trees.
- Basic idea: generalize the notion of a one-dimensional interval associated with each B⁺-tree node to an N-dimensional interval, that is, an N-dimensional rectangle.
- Will consider only the two-dimensional case ($N = 2$)
 - generalization for $N > 2$ is straightforward, although R-trees work well only for relatively small N



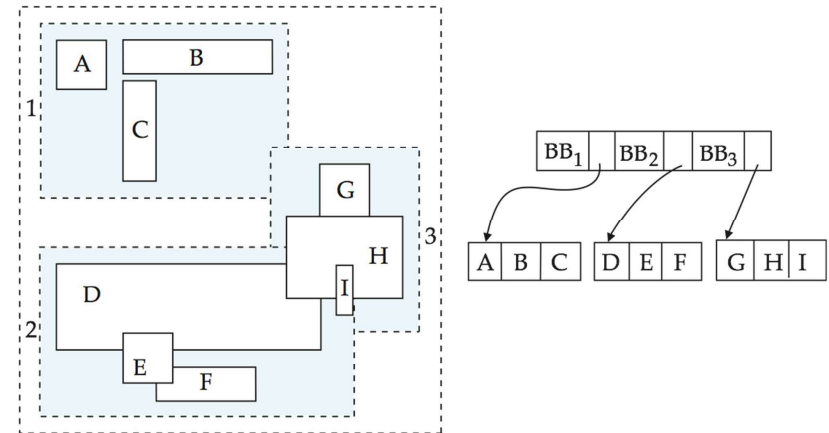
R Trees (Cont.)

- A rectangular **bounding box** is associated with each tree node.
 - Bounding box of a leaf node is a minimum sized rectangle that contains all the rectangles/polygons associated with the leaf node.
 - The bounding box associated with a non-leaf node contains the bounding box associated with all its children.
 - Bounding box of a node serves as its key in its parent node (if any)
 - *Bounding boxes of children of a node are allowed to overlap*
- A polygon is stored only in one node, and the bounding box of the node must contain the polygon.
 - The storage efficiency of R-trees is better than that of k-d trees or quadtrees since a polygon is stored only once.



Example R-Tree

- A set of rectangles (solid line) and the bounding boxes (dashed line) of the nodes of an R-tree for the rectangles. The R-tree is shown on the right.



Search in R-Trees

- To find data items (rectangles/polygons) intersecting (overlaps) a given query point/region, do the following, starting from the root node:
 - If the node is a leaf node, output the data items whose keys intersect the given query point/region.
 - Else, for each child of the current node whose bounding box overlaps the query point/region, recursively search the child
- Can be very inefficient in worst case since multiple paths may need to be searched
 - but works acceptably in practice.
- Simple extensions of search procedure to handle predicates *contained-in* and *contains*



Insertion in R-Trees

- To insert a data item:
 - Find a leaf to store it, and add it to the leaf
 - ▶ To find leaf, follow a child (if any) whose bounding box contains bounding box of data item, else child whose overlap with data item bounding box is maximum
 - Handle overflows by splits (as in B+-trees)
 - ▶ Split procedure is different though (see below)
 - Adjust bounding boxes starting from the leaf upwards
- Split procedure:
 - Goal: divide entries of an overfull node into two sets such that the bounding boxes have minimum total area
 - ▶ This is a heuristic. Alternatives like minimum overlap are possible
 - Finding the "best" split is expensive, use heuristics instead
 - ▶ See next slide



Splitting an R-Tree Node

- **Quadratic split** divides the entries in a node into two new nodes as follows
 1. Find pair of entries with “maximum separation”
 - ▶ that is, the pair such that the bounding box of the two would have the maximum wasted space (area of bounding box – sum of areas of two entries)
 2. Place these entries in two new nodes
 3. Repeatedly find the entry with “maximum preference” for one of the two new nodes, and assign the entry to that node
 - Preference of an entry to a node is the increase in area of bounding box if the entry is added to the *other* node
 4. Stop when half the entries have been added to one node
 - Then assign remaining entries to the other node
- Cheaper **linear split** heuristic works in time linear in number of entries,
 - Cheaper but generates slightly worse splits.



Deleting in R-Trees

- Deletion of an entry in an R-tree done much like a B⁺-tree deletion.
 - In case of underfull node, borrow entries from a sibling if possible, else merging sibling nodes
 - Alternative approach removes all entries from the underfull node, deletes the node, then reinserts all entries