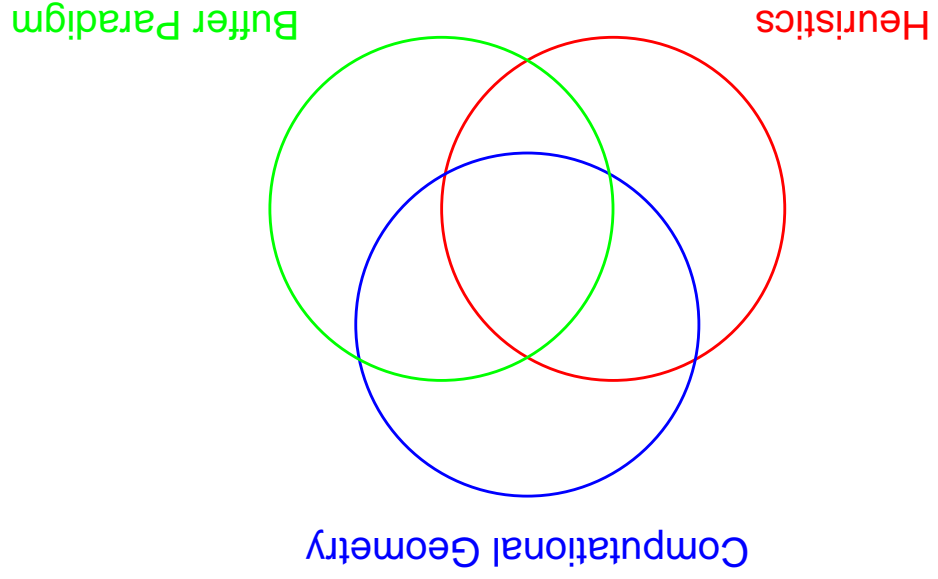


# R-trees



# Spatial Data

**Spatial data:** points, lines, polygons/polyhedra in  $\mathbb{R}^d$ .



For simplicity: assume  $d = 2$ .

**Typical queries:**

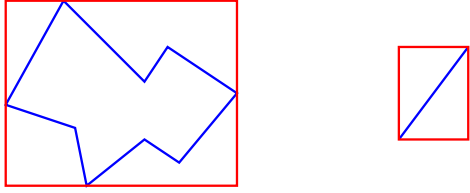
- Point:** All objects containing query point
- Intersection:** All objects intersecting query object
- Enclosure:** All objects containing query object
- Spatial Join:** All intersections between stored objects and query set of objects

**Example:** windowing in GIS.

# R-trees

Approximate objects by **Bounding Box**:

[Gutman, 84]

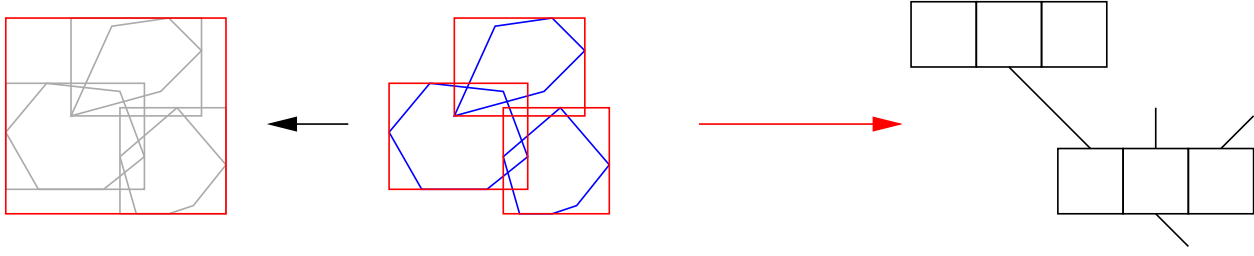


**R-tree** = B-tree where

- Leaves store objects and their BBs.

- Internal nodes store for each child a BB containing all objects in

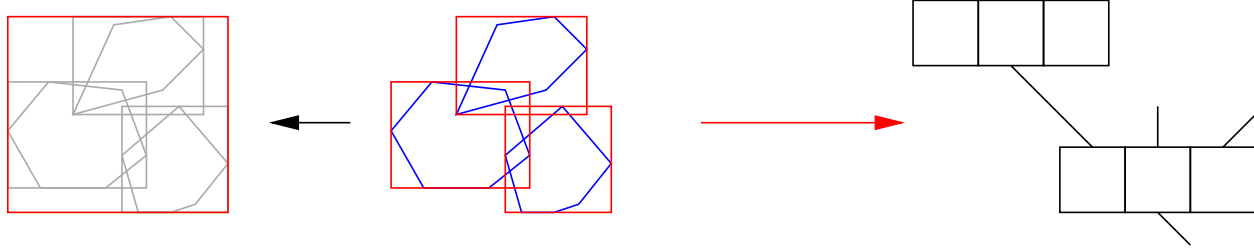
that subtree.



# Searching in R-trees

Example: **Point Query**

Search recursively in all subtrees of node where BB contains the query point.



Observe:

- Worst case query time:  $\Theta(N)^i$
- Average case performance depends heavily on distribution of objects in leaves.

Heuristic data structure.

# Static R-Trees

**Build** ("bulk load" in database community):

1. Distribute objects in leaves ( $\Leftrightarrow$  assign linear order).
2. Build R-tree bottom up.

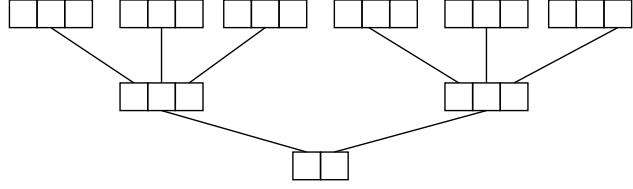
Number of I/O's:

1.  $O(\text{Sort}(N))$
2.  $O(\text{Scan}(N))$

**Query efficiency**: depends heavily on the linear order.

Example: random distribution in leaves  $\Rightarrow$  most leaves has BB close to BB of entire set of objects  $\Rightarrow$  queries inefficient.

**Goal**: linear order where objects close in order  $\Leftrightarrow$  close in Euclidean distance.



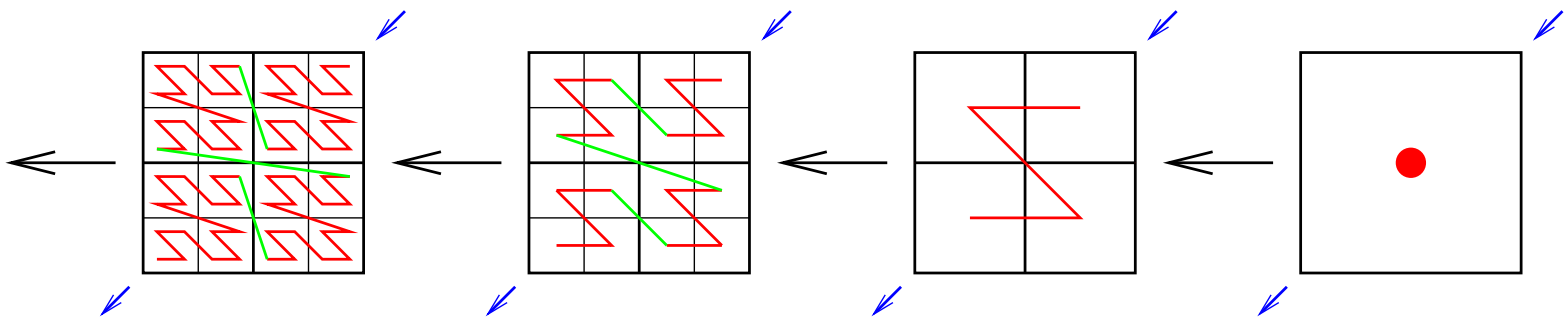
“Bit interleaving”

Compare  $y_1x_1y_2x_2 \dots x_k$  and  $y'_1x'_1y'_2x'_2 \dots x'_k$



Compare  $y_1$  and  $y'_1$ , then  $x_1$  and  $x'_1$ , then  $y_2$  and  $y'_2$ , ...

Comparing  $(y, x)$  and  $(y', x')$  having binary expansions  $(y_1x_2 \dots x_k, y_1y_2 \dots y_k)$  and  $(y'_1x'_2 \dots x'_k, y'_1y'_2 \dots y'_k)$ :

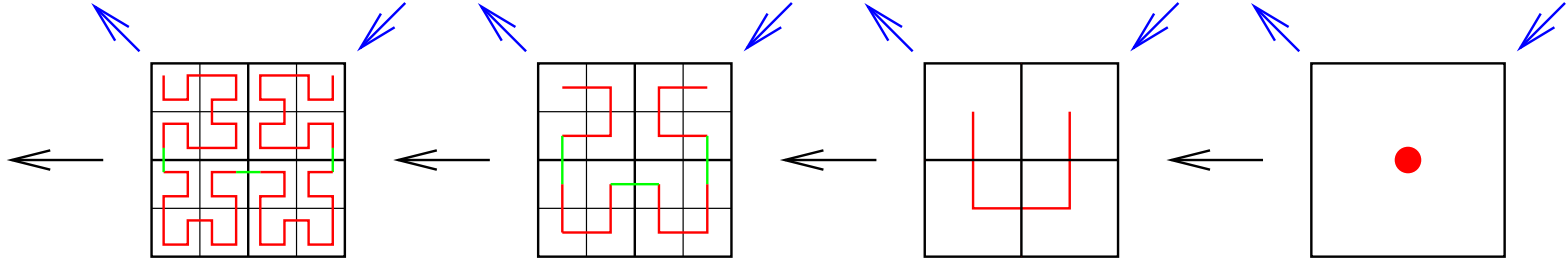


The z-curve:

# Space Filling Curves

# Space Filling Curves

The Hilbert-curve:



Note: subpattern = pattern flipped along a diagonal. Generates four possibilities (as  $\text{flip}_1^2 = \text{flip}_2^2 = \text{ID}$  and as  $\text{flip}_1 \text{flip}_2 = \text{flip}_2 \text{flip}_1$ ).

Comparing  $(x, y)$  and  $(x', y')$  having binary expansions  $(x_1 x_2 \dots x_k, y_1 y_2 \dots y_k)$  and  $(x'_1 x'_2 \dots x'_k, y'_1 y'_2 \dots y'_k)$ :

Top level: Compare  $x_1$  and  $x'_1$ , then  $y_1$  and  $y'_1$ . Advance to next bit, choosing proper variant of this scheme.

# Space Filling Curves and R-Trees

Choose ordering in leaves by sorting **midpoints** of bounding boxes of objects according to position on some space filling curve.

**Empirically:**

A static R-tree build this way has good query time in practice.



# Dynamic R-trees

Insert:

We need two subroutines:

**Route()** Given a node  $v$  and a new object  $x$ , decide which subtree of  $v$  that  $x$  should belong to.

**Split()** Given a node  $v$  that is overflowing, decide how to split the subtrees in two groups.

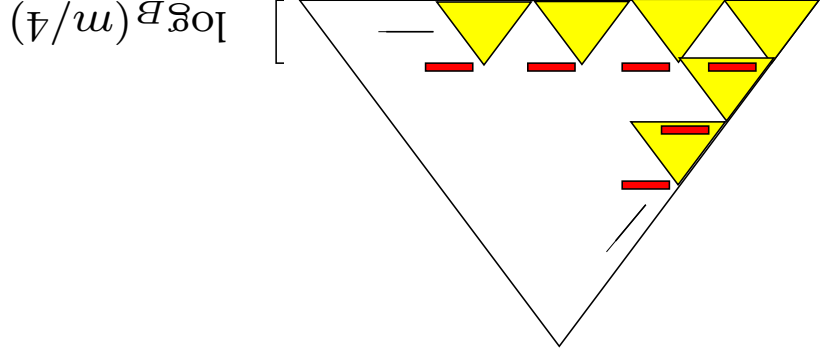
Various **heuristics** for these have been proposed, based on minimizing various properties of changed BBS (total area, area not used by objects, circumference, overlap with siblings  $BB_1, \dots$ ).

**Bottom line:**  $\exists$  fast heuristics giving good query times in practice.

# Buffered R-Trees

[AHV, 99]

Add buffers of size  $\Theta(M)$  to layers  $i \cdot \log_B(m/4)$ , for  $i = 1, 2, \dots$



Principle and analysis: exactly as for buffer trees (using **Route** and **Split**)

**Advantage:** Efficient bulk updates and queries. For  $\Theta(N)$  size bulks: Number of I/O's is  $O(\frac{B}{1} \log^m(N))$  per insert instead of  $\Theta(\log_B(N))$ .

## Buffered R-Trees

Similar statements for queries and deletions (see paper).

Note that spatial join is a bulk query.

**Empirically:** Faster bulk insertions than for previous proposals. Query time in resulting structure competitive with repeated insertions.