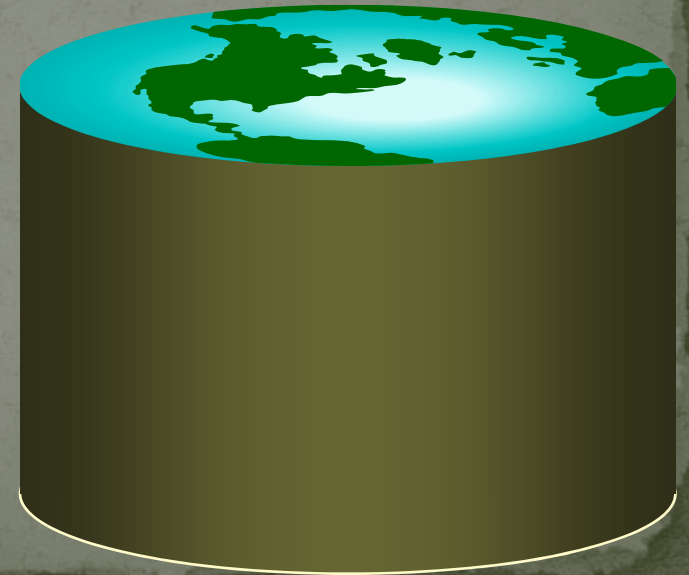


CS 505: Intermediate Topics to Database Systems

Instructor: Jinze Liu

Fall 2008



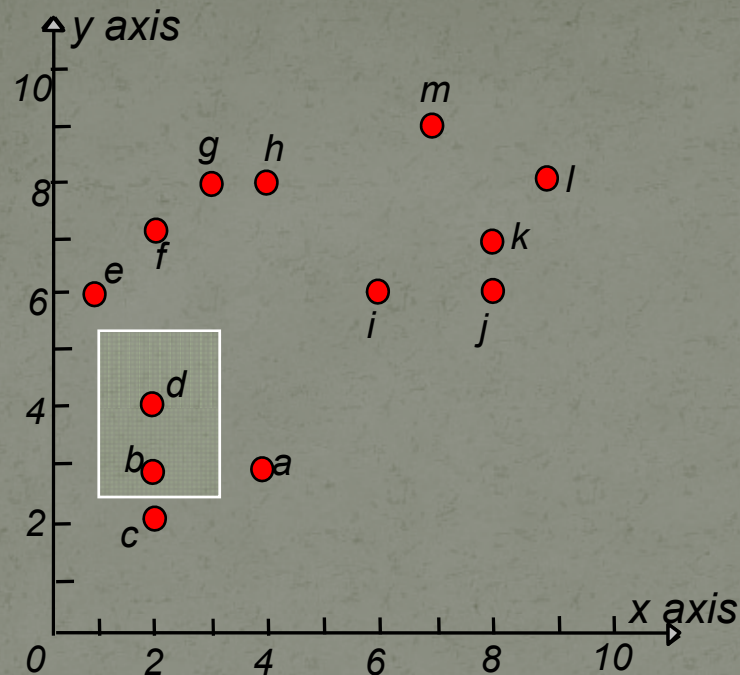
Beyond ISAM, B-, and B⁺-trees

- Other tree-based indexes: R-trees and variants, GiST, etc.
- Hashing-based indexes: extensible hashing, linear hashing, etc.
- Text indexes: inverted-list index, suffix arrays, etc.
- Other tricks: bitmap index, bit-sliced index, etc.
 - How about indexing subgraph search?

R-Tree

- The R-tree
 - Range Query
 - Aggregation Query
- NN Query
- RNN Query
- Closest Pair Query
- Close Pair Query
- Skyline Query

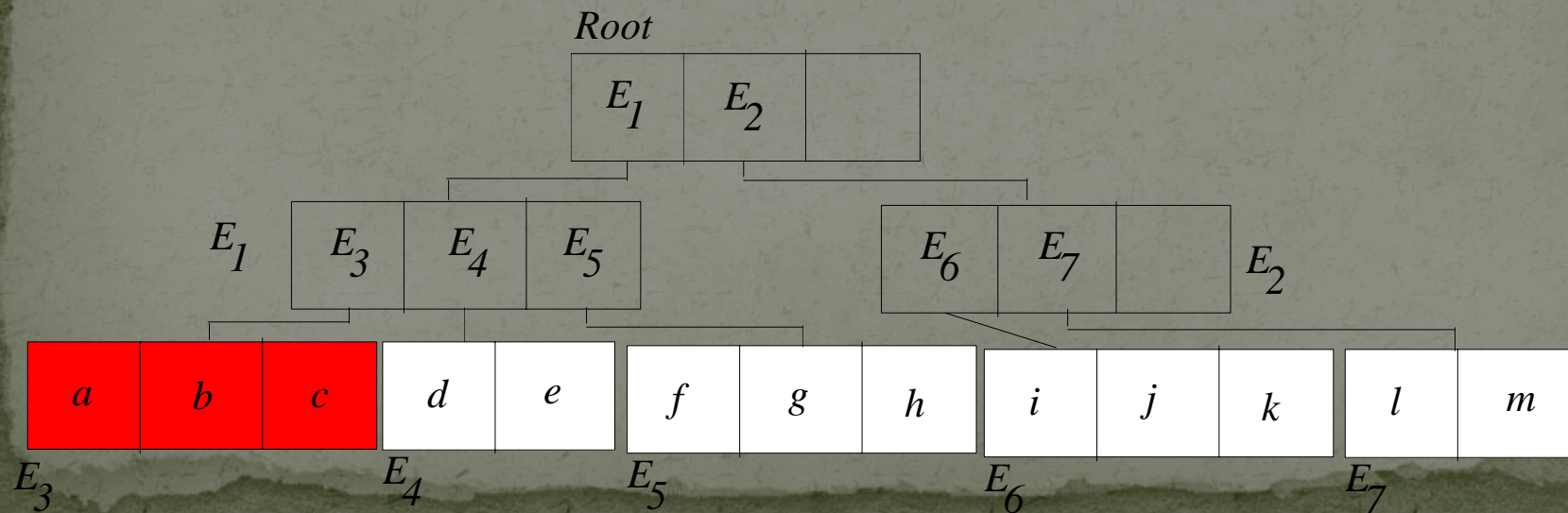
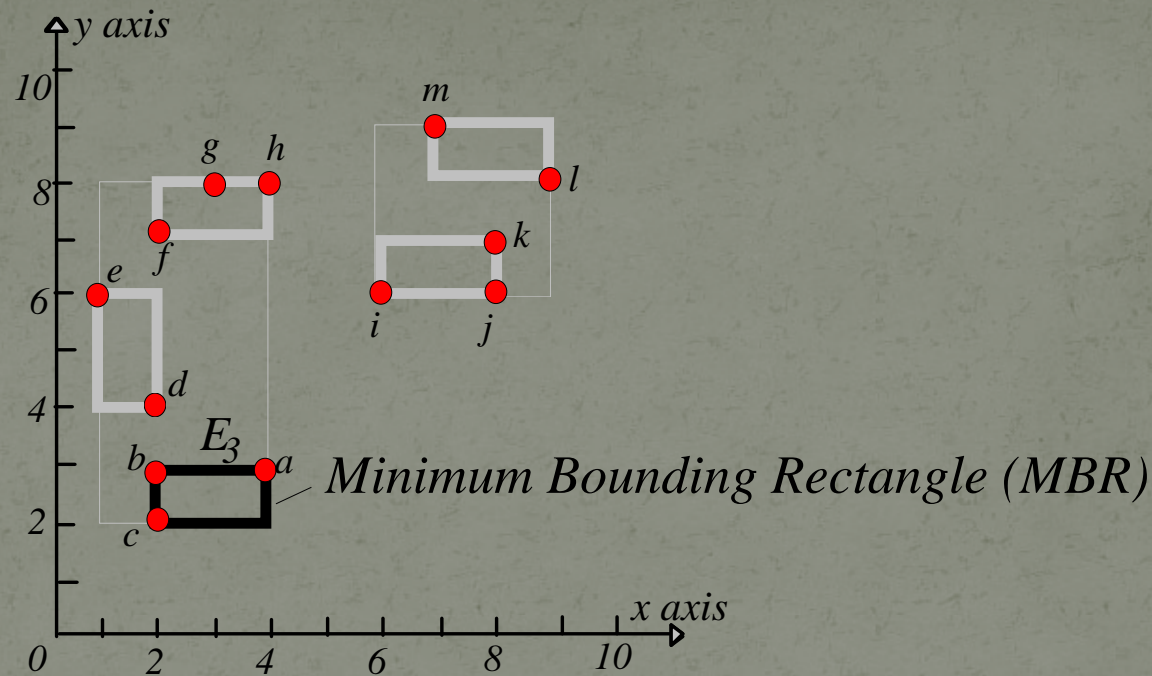
R-Tree Motivation



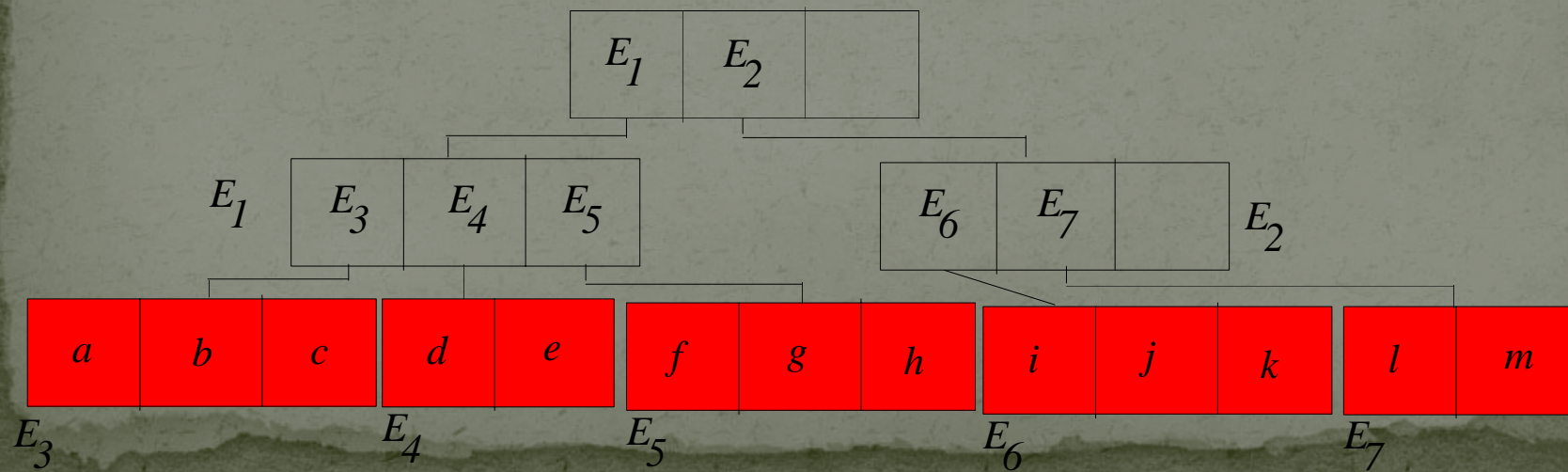
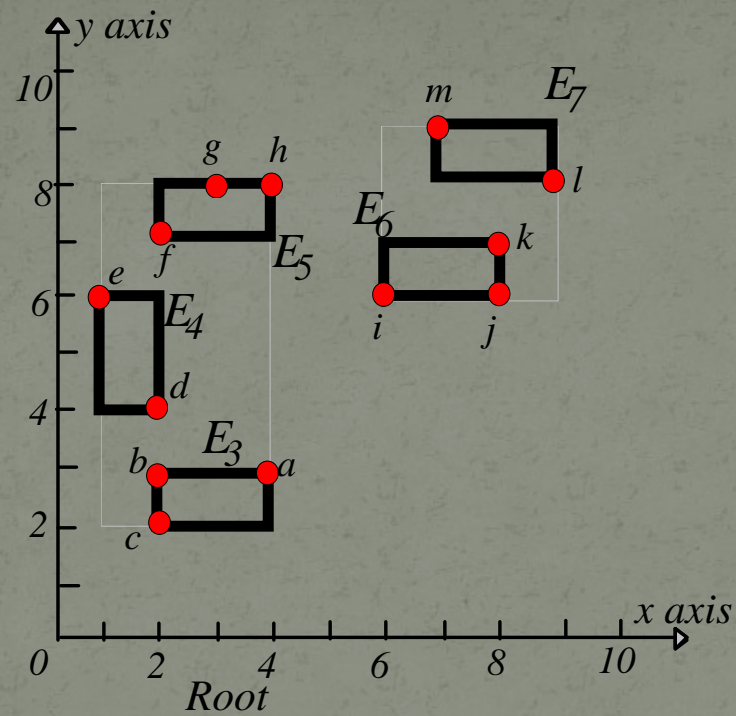
Range query: find the objects in a given range.
E.g. find all hotels in Boston.

No index: scan through all objects. NOT EFFICIENT!

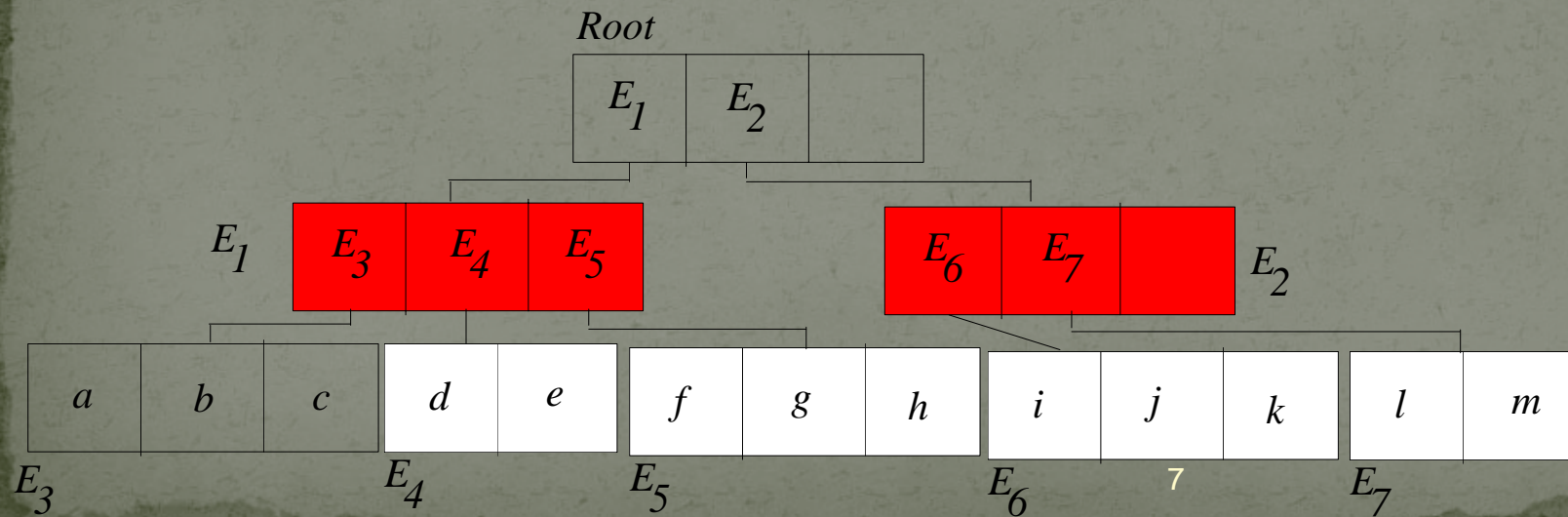
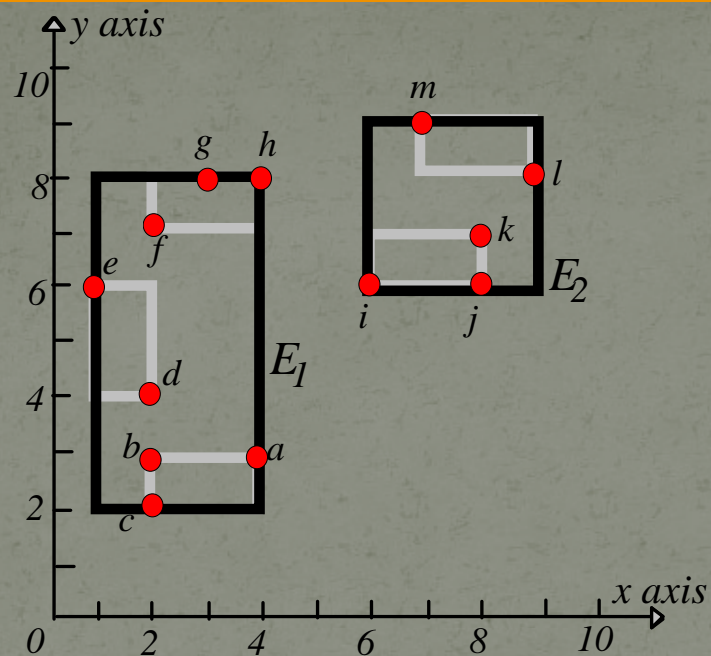
R-Tree: Clustering by Proximity



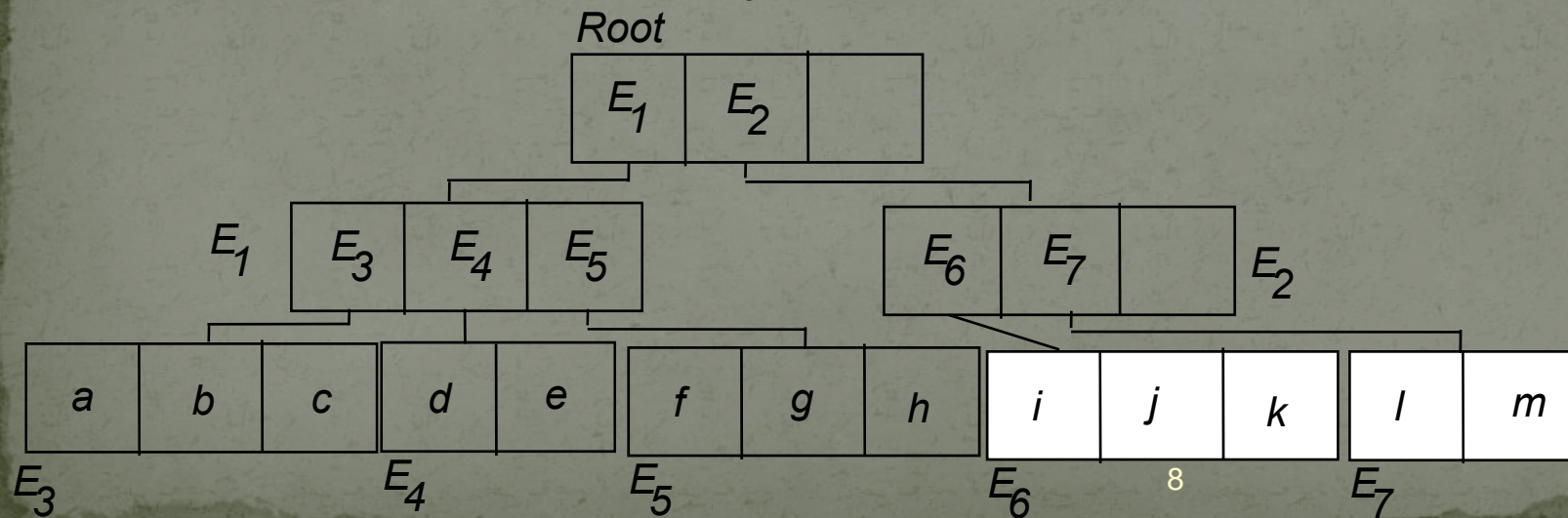
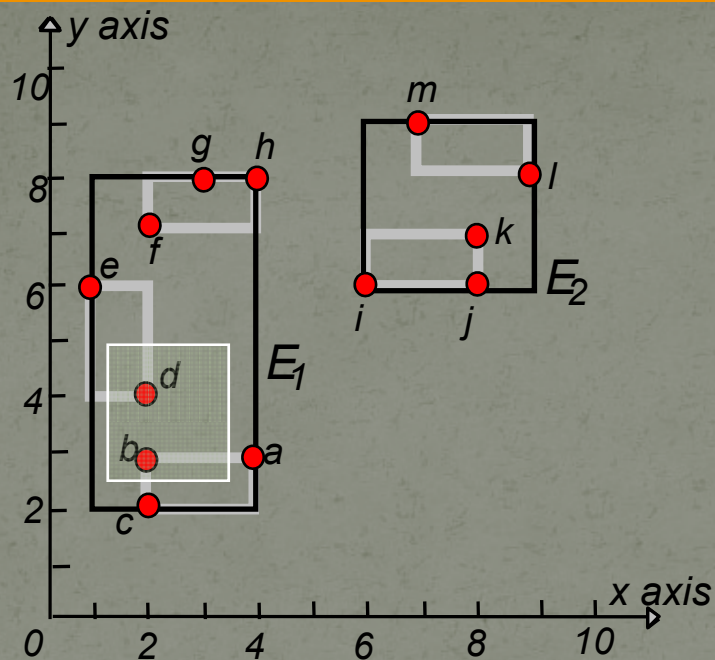
R-Tree



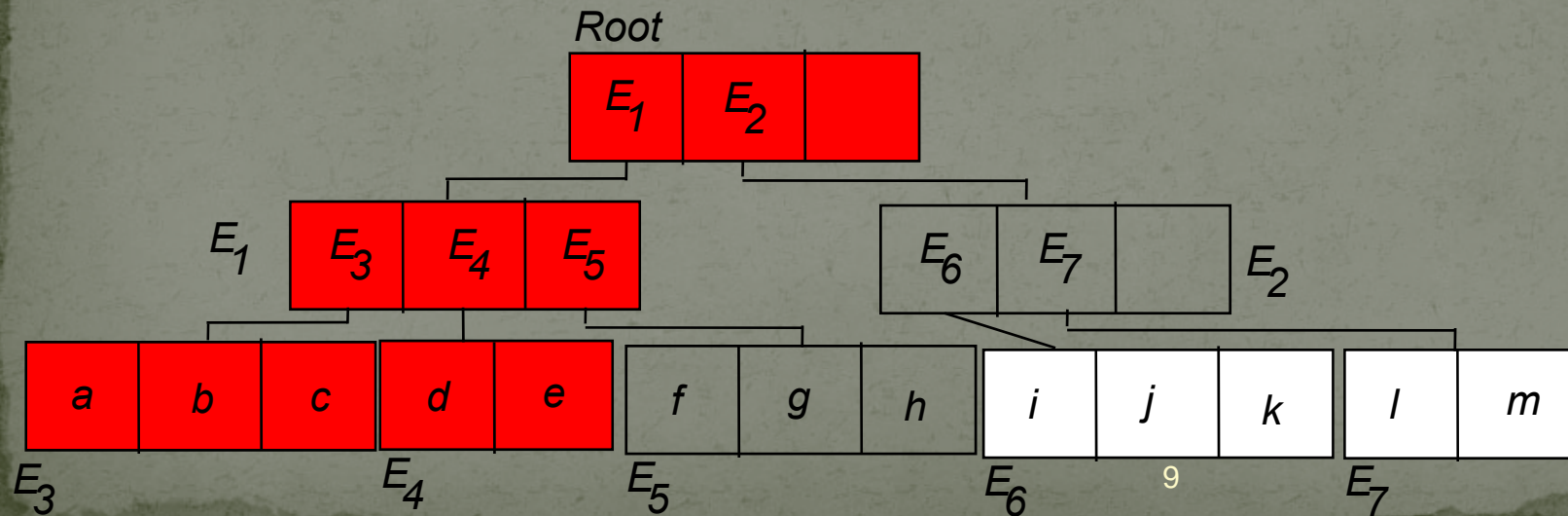
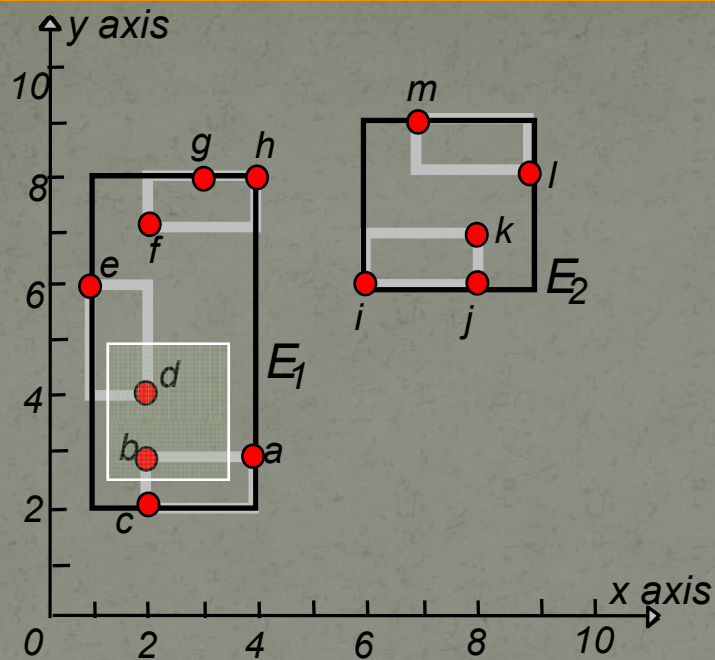
R-Tree



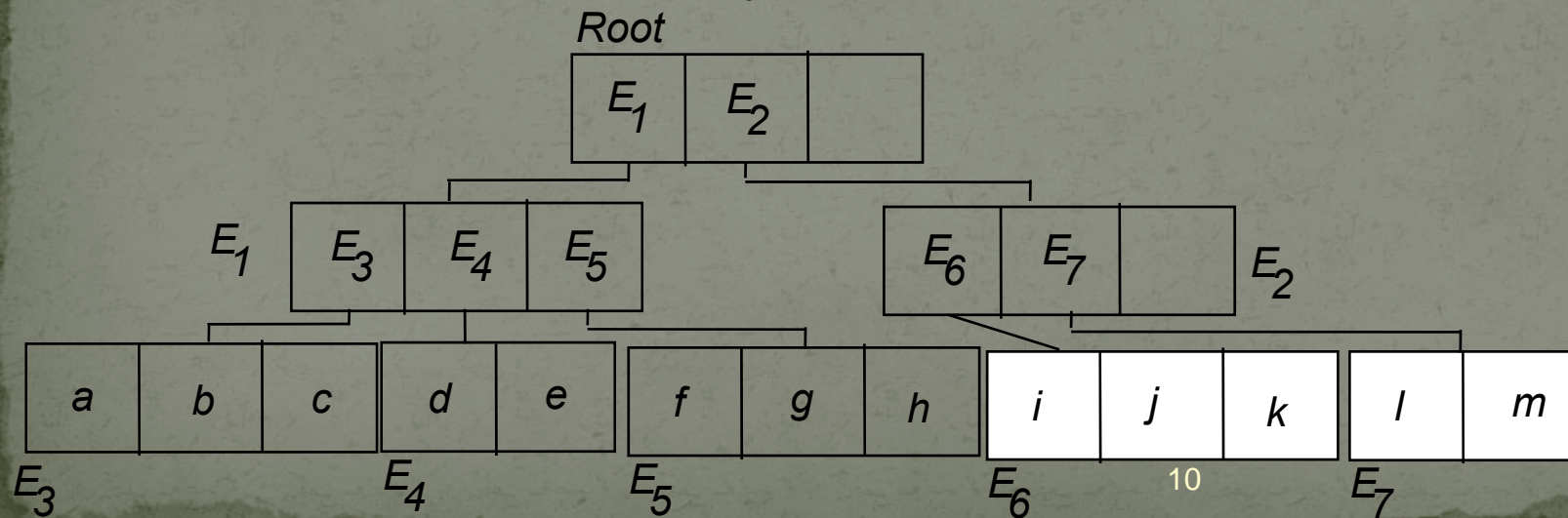
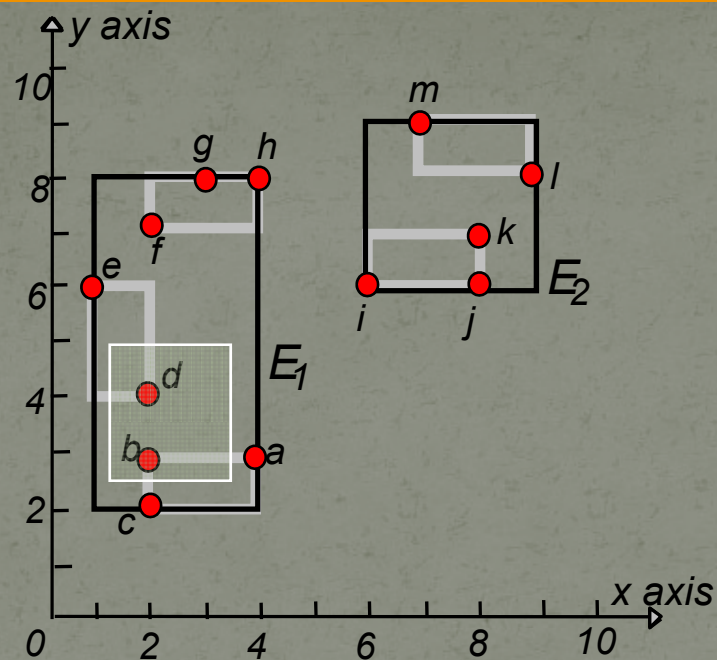
Range Query



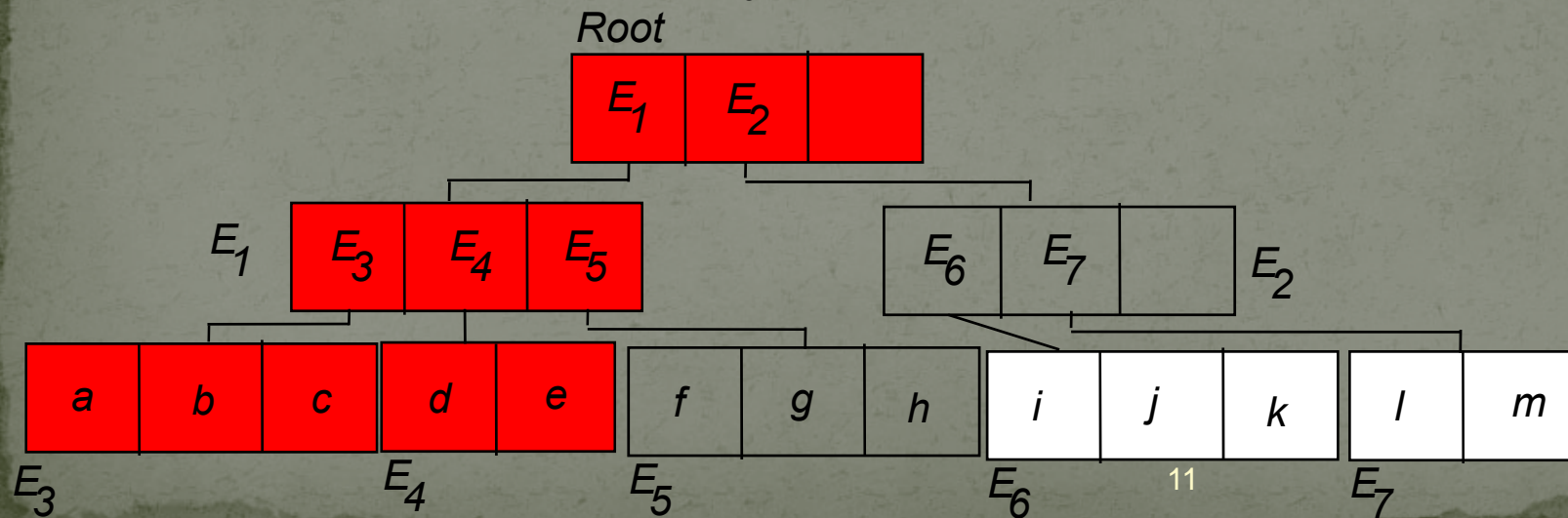
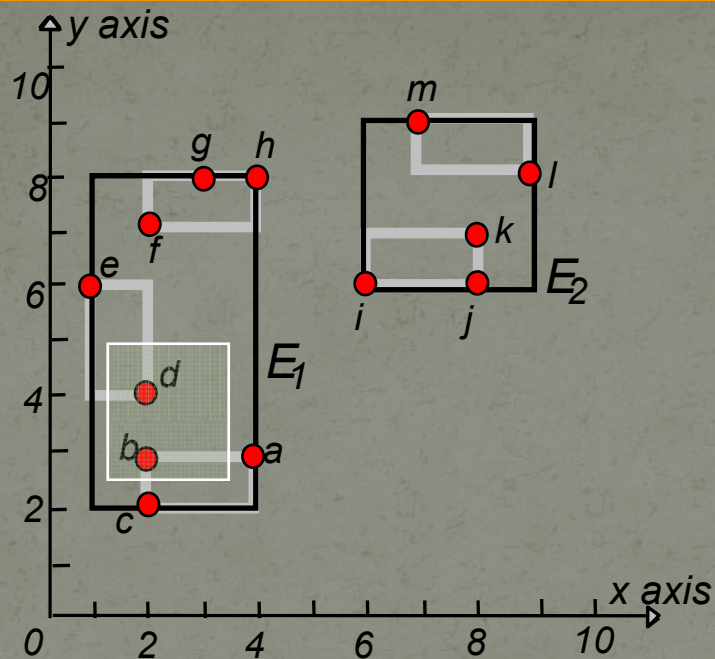
Range Query



Range Query



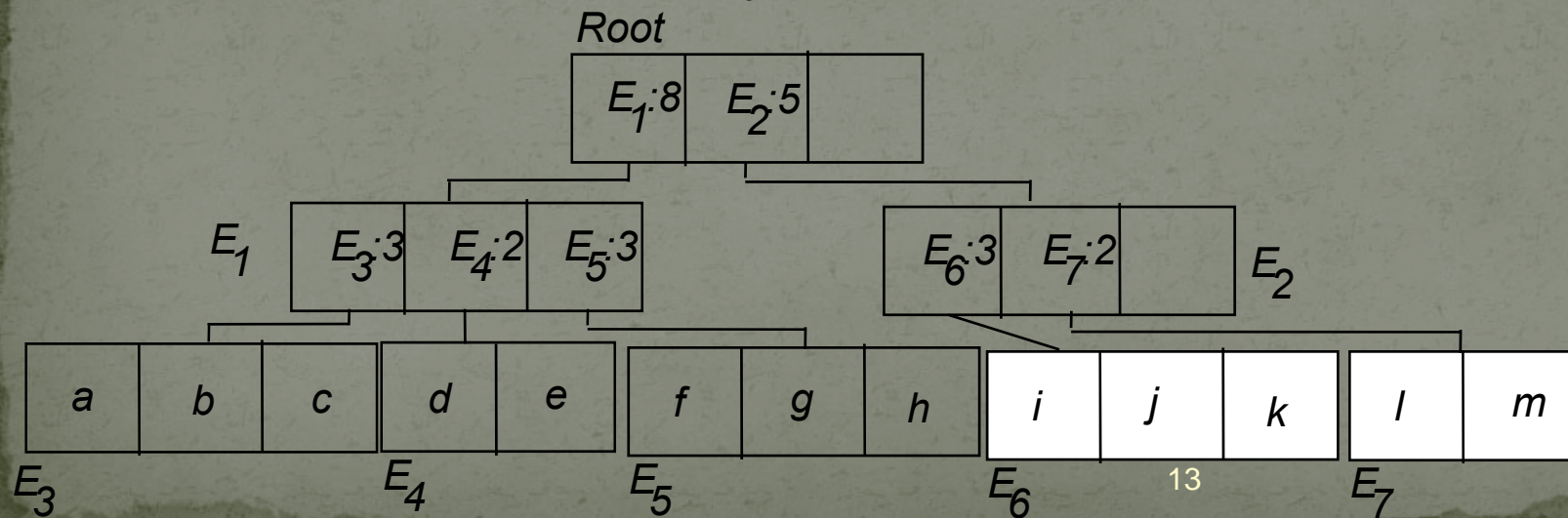
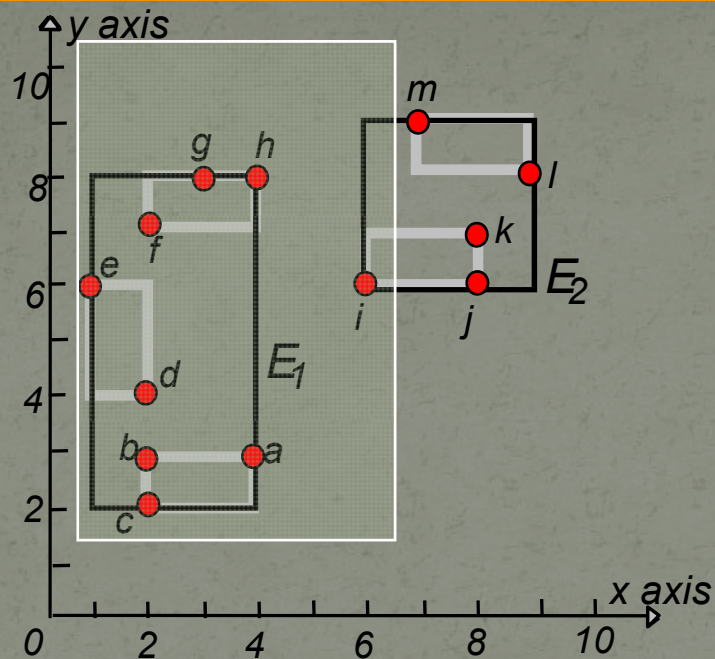
Range Query



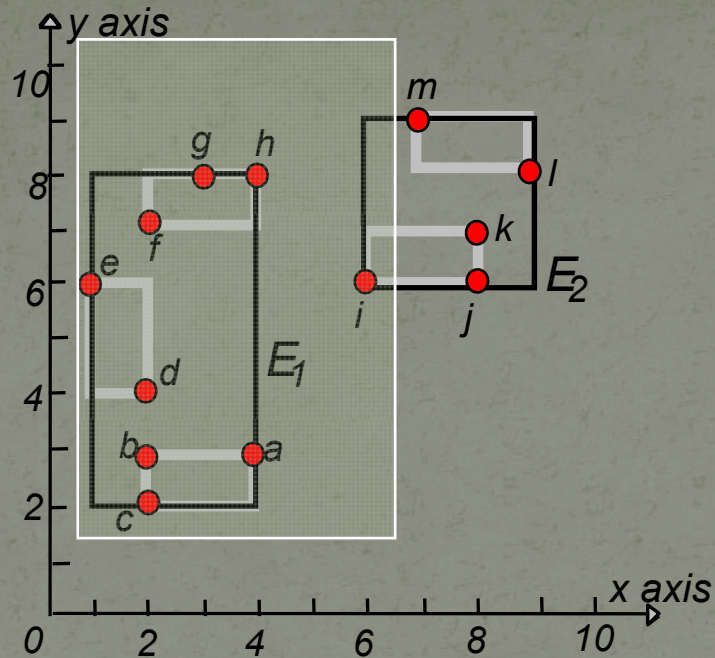
Aggregation Query

- *Given a range, find some aggregate value of objects in this range.*
 - COUNT, SUM, AVG, MIN, MAX
 - E.g. find the total number of hotels in Massachusetts.
-
- Straightforward approach: reduce to a range query.
 - Better approach: along with each index entry, store aggregate of the sub-tree.

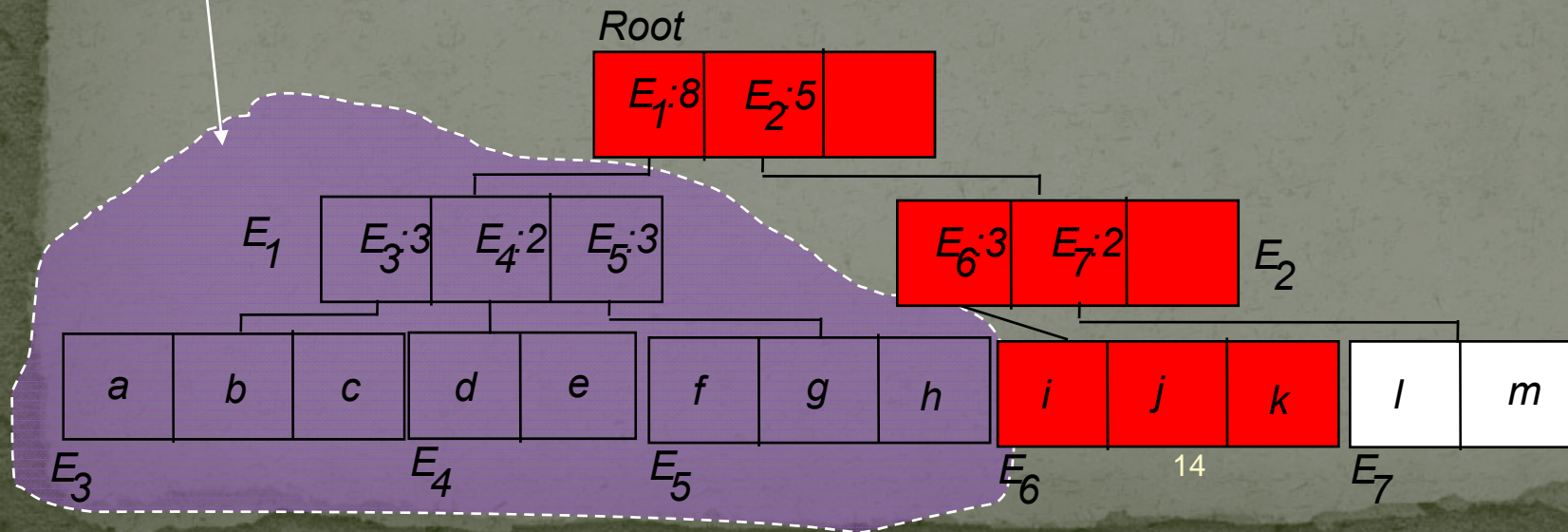
Aggregation Query



Aggregation Query



Subtree pruned!



Content

- The R-tree
 - Range Query
 - Aggregation Query
- **NN Query**
- RNN Query
- Closest Pair Query
- Close Pair Query
- Skyline Query

Nearest Neighbor (NN) Query

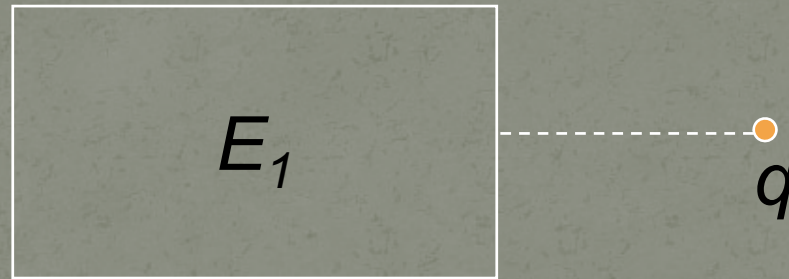
- Given a query location q , find the nearest object.



- E.g.: given a hotel, find its nearest bar.

A Useful Metric: MINDIST

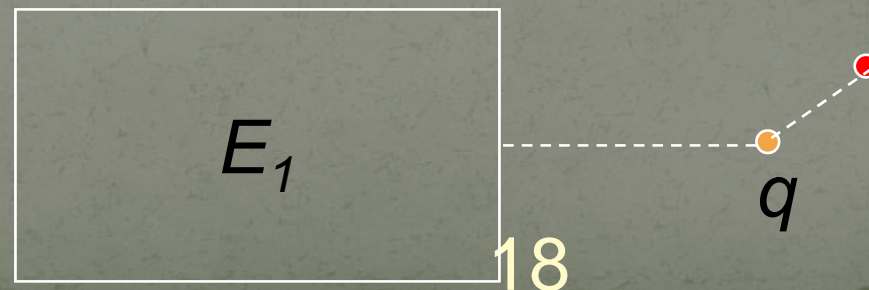
- Minimum distance between q and an MBR.



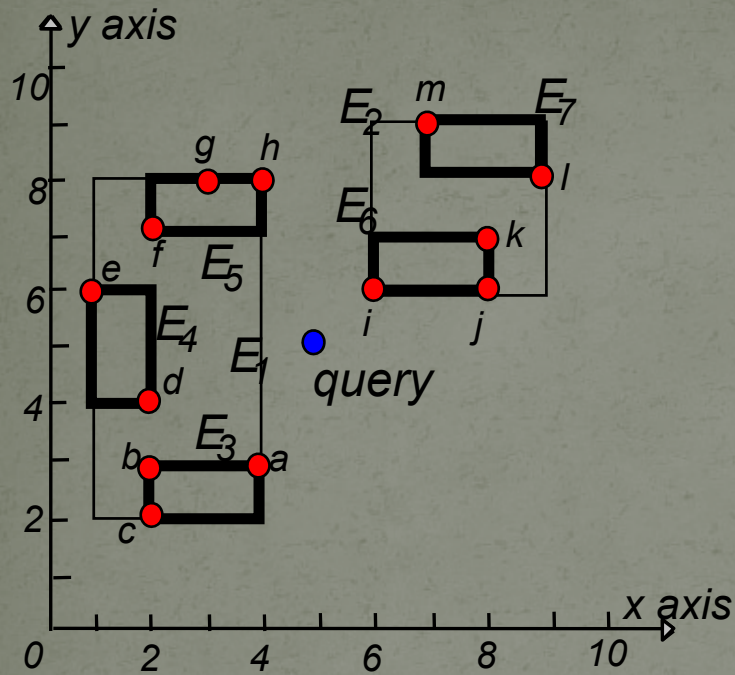
- It is an lower bound of $d(o, q)$ for every object o in E_1 .
- $\text{MINDIST}(o, q) = d(o, q)$.

NN Basic Algorithm

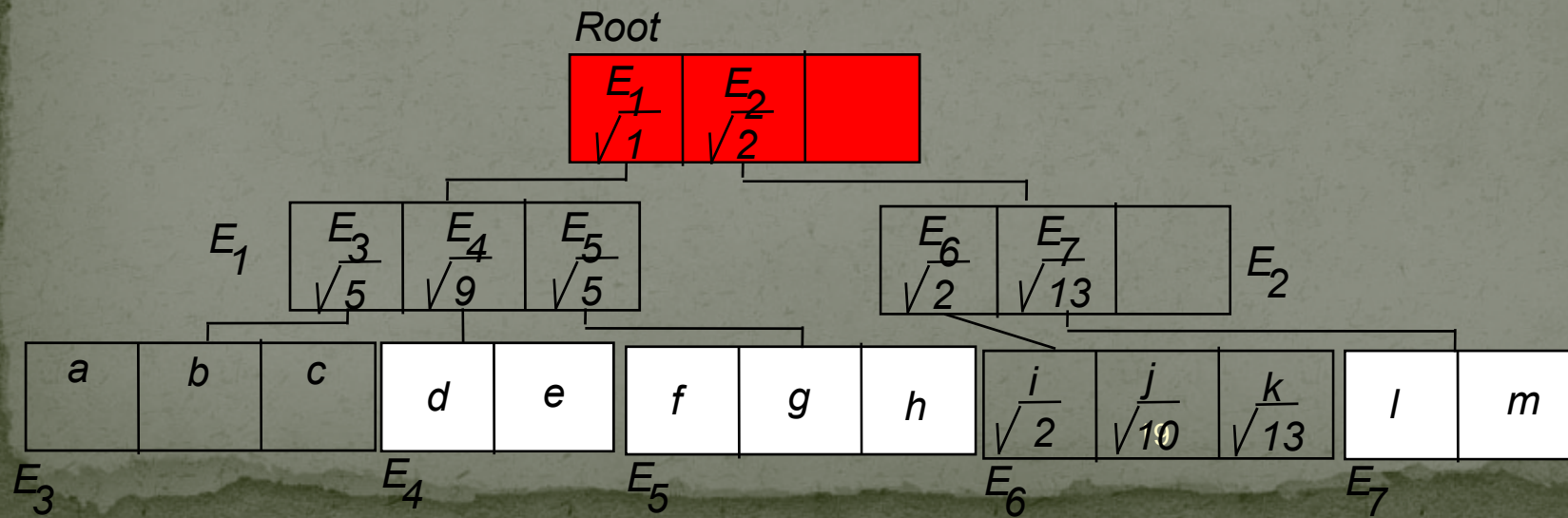
- Keep a heap H of index entries and objects, ordered by MINDIST.
- Initially, H contains the root.
- While $H \neq \phi$
 - Extract the element with minimum MINDIST
 - If it is an index entry, insert its children into H .
 - If it is an object, return it as NN.
- End while



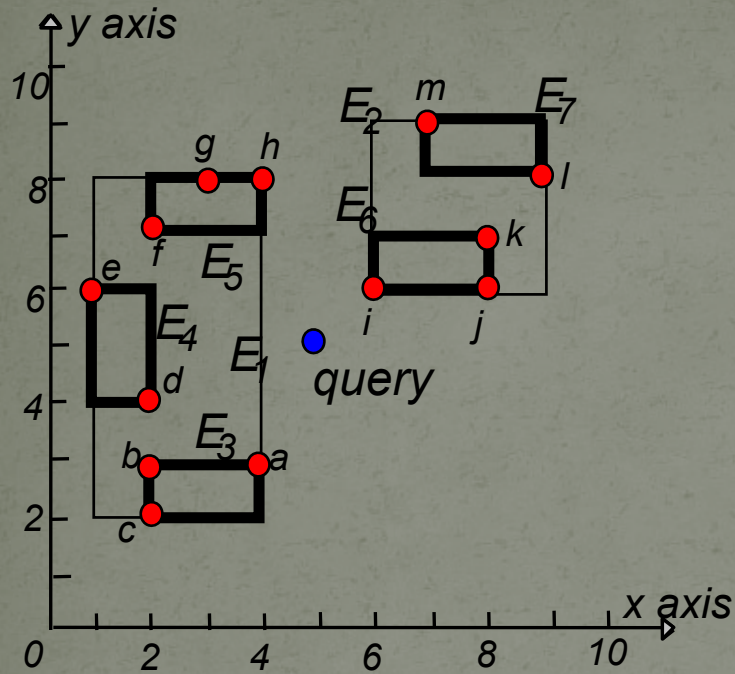
NN Query Example



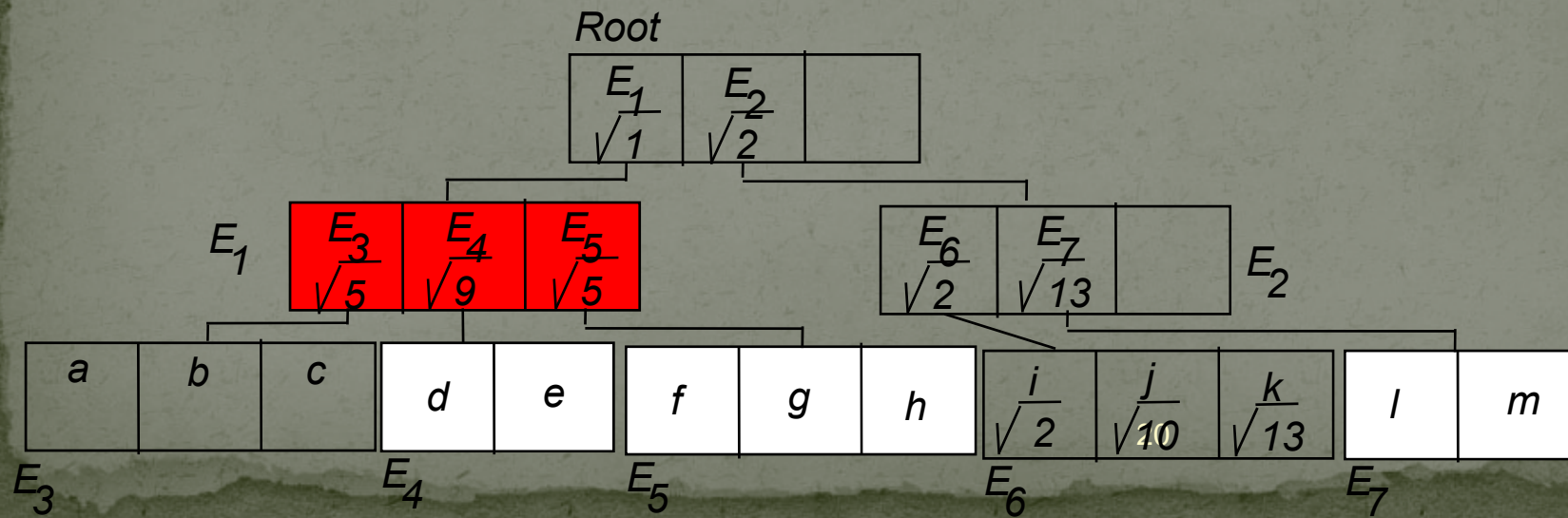
Action	Heap
Visit Root	$E_1\sqrt{1}$ $E_2\sqrt{2}$



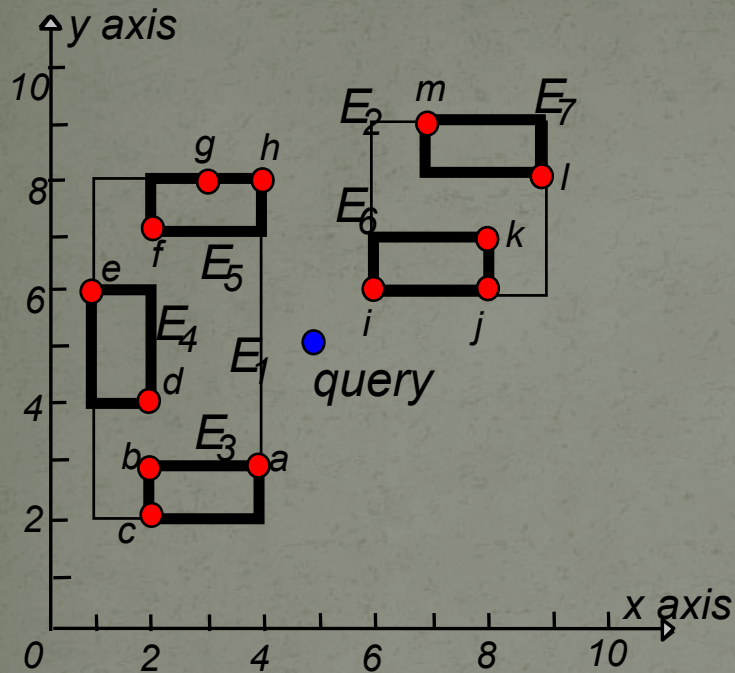
NN Query Example



Action	Heap						
Visit Root	$E_1\sqrt{1}$	$E_2\sqrt{2}$					
follow E_1	$E_2\sqrt{2}$	$E_3\sqrt{5}$	$E_5\sqrt{5}$	$E_4\sqrt{9}$			



NN Query Example



Action

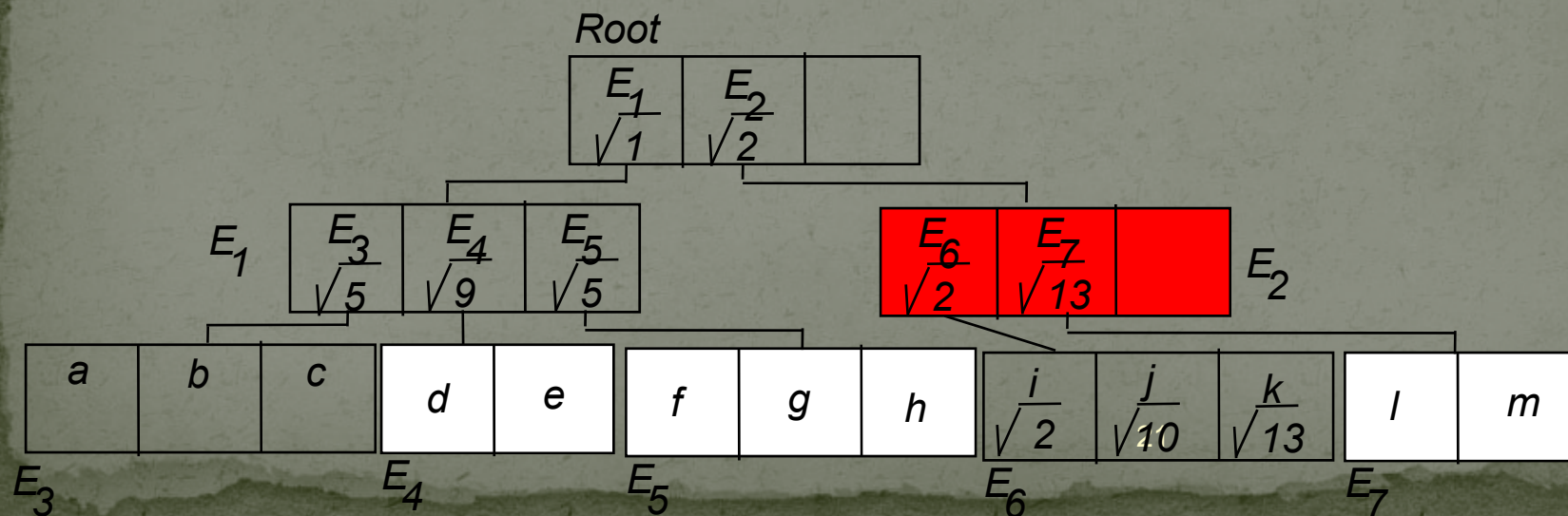
Visit Root

follow E_1

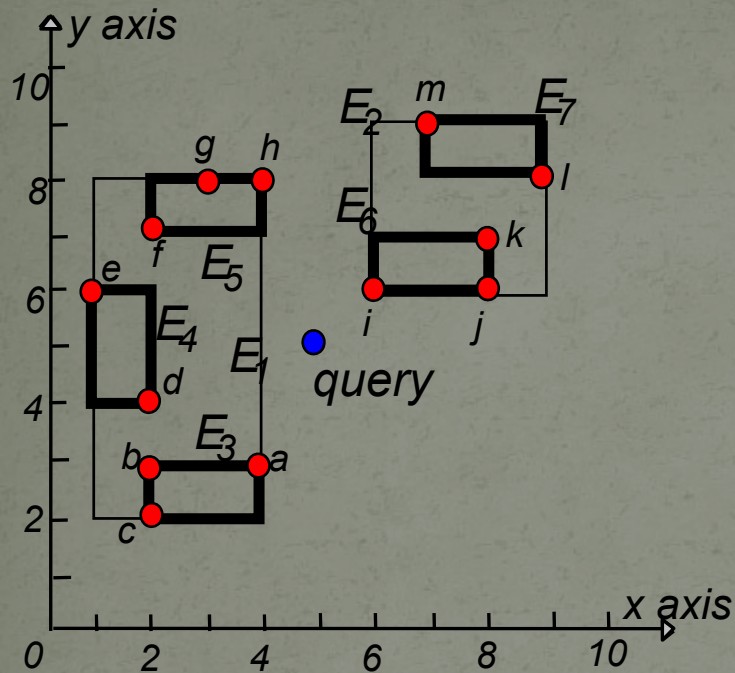
follow E_2

Heap

$E_1\sqrt{1}$	$E_2\sqrt{2}$					
$E_2\sqrt{2}$	$E_3\sqrt{5}$	$E_5\sqrt{5}$	$E_4\sqrt{9}$			
$E_6\sqrt{2}$	$E_3\sqrt{5}$	$E_5\sqrt{5}$	$E_4\sqrt{9}$	$E_7\sqrt{13}$		



NN Query Example



Action

Visit Root

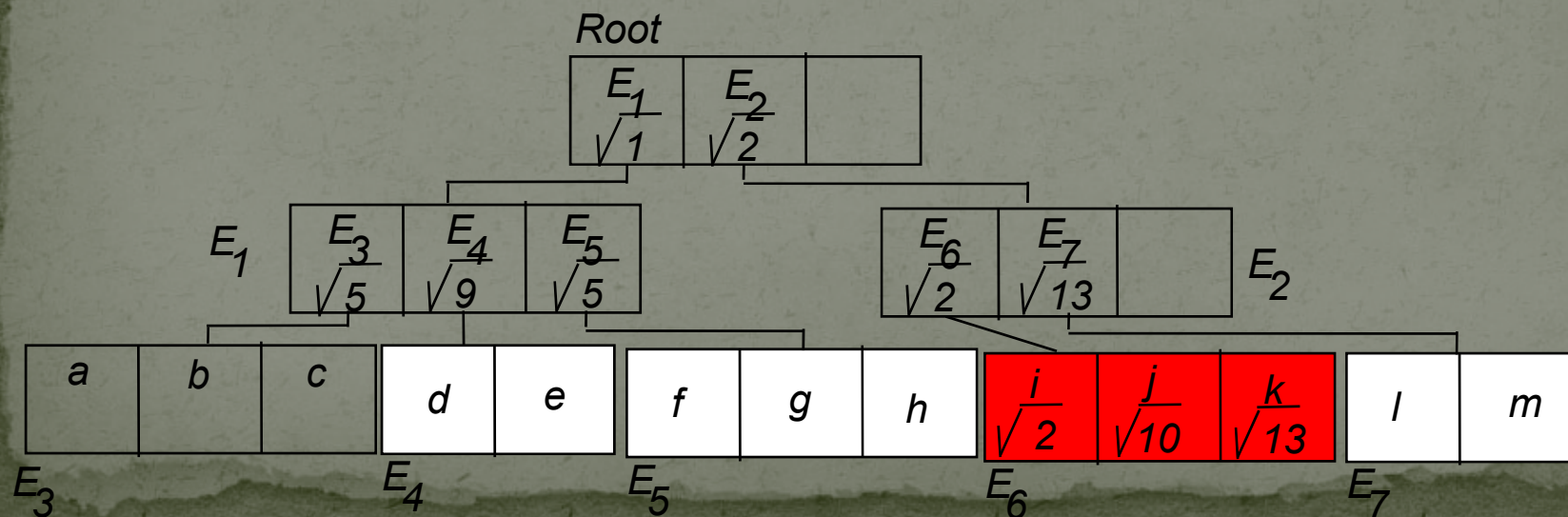
follow E_1

follow E_2

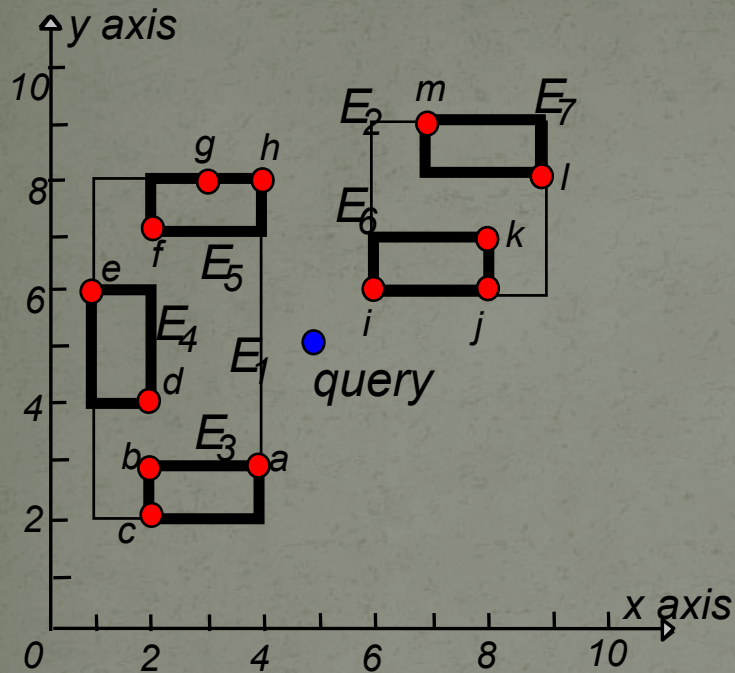
follow E_6

Heap

$E_1\sqrt{1}$	$E_2\sqrt{2}$					
$E_2\sqrt{2}$	$E_3\sqrt{5}$	$E_5\sqrt{5}$	$E_4\sqrt{9}$			
$E_6\sqrt{2}$	$E_3\sqrt{5}$	$E_5\sqrt{5}$	$E_4\sqrt{9}$	$E_7\sqrt{13}$		
$i\sqrt{2}$	$E_3\sqrt{5}$	$E_5\sqrt{5}$	$E_4\sqrt{9}$	$j\sqrt{10}$	$E_7\sqrt{13}$	$k\sqrt{13}$

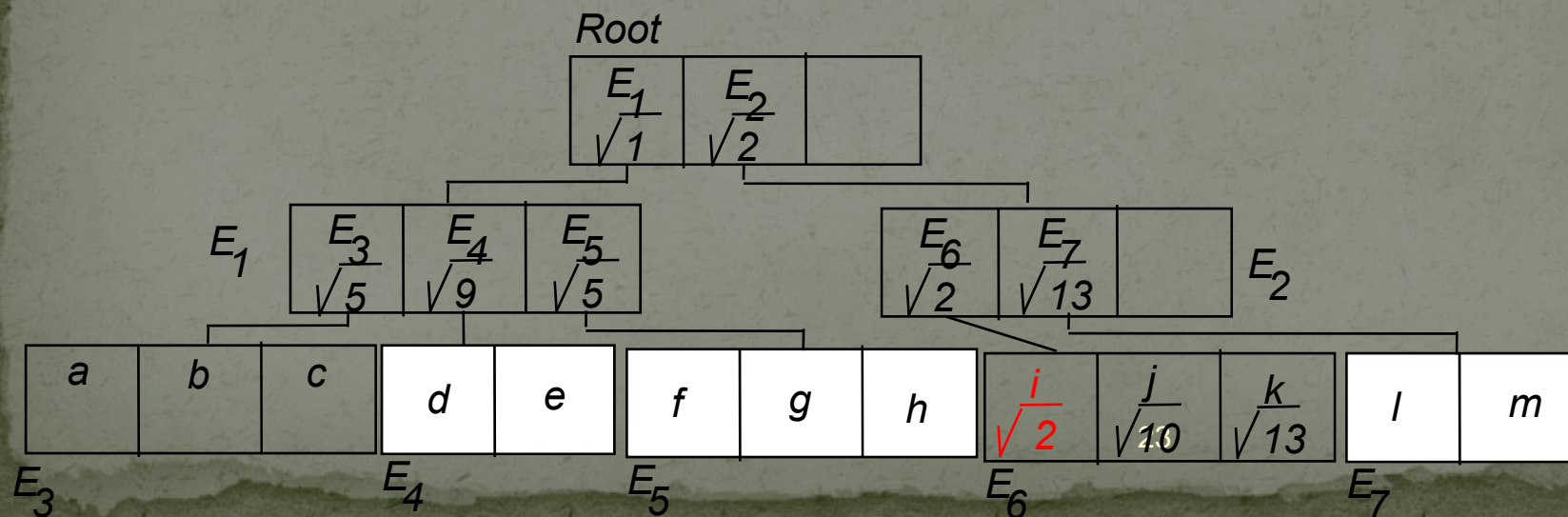


NN Query Example



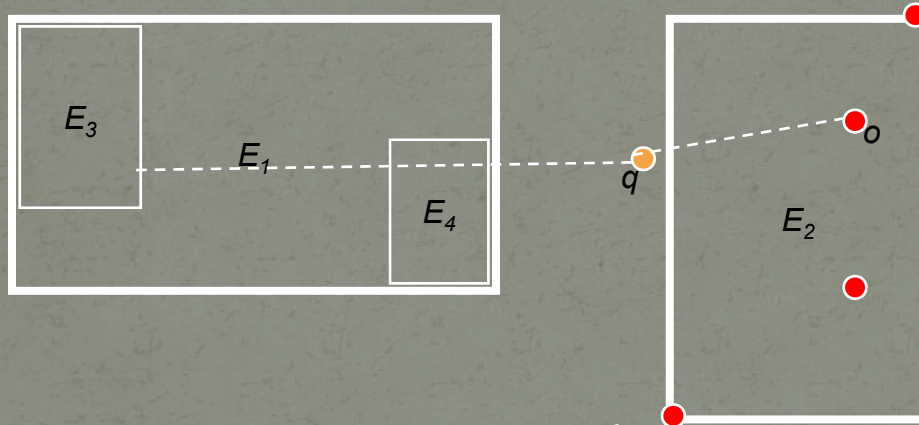
Action	Heap						
Visit Root	$E_1\sqrt{1}$	$E_2\sqrt{2}$					
follow E_1	$E_2\sqrt{2}$	$E_3\sqrt{5}$	$E_5\sqrt{5}$	$E_4\sqrt{9}$			
follow E_2	$E_6\sqrt{2}$	$E_3\sqrt{5}$	$E_5\sqrt{5}$	$E_4\sqrt{9}$	$E_7\sqrt{13}$		
follow E_6	$i\sqrt{2}$	$E_3\sqrt{5}$	$E_5\sqrt{5}$	$E_4\sqrt{9}$	$j\sqrt{10}$	$E_7\sqrt{13}$	$k\sqrt{13}$

Report i and terminate



Pruning 1 in NN Query

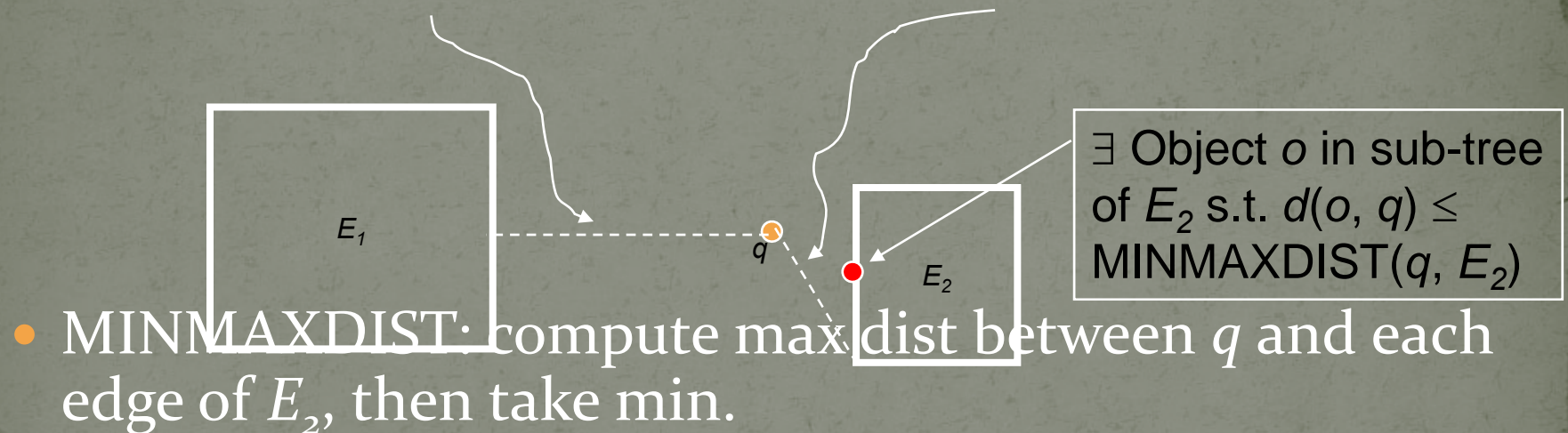
- If we see an object o , prune every MBR whose $\text{MINDIST} > d(o, q)$.



- Side notice: at most one object in H !

Pruning 2 using MINMAXDIST

- Prune even before we see an object!
- Prune E_1 if exists E_2 s.t.
 $\text{MINDIST}(q, E_1) > \text{MINMAXDIST}(q, E_2)$.



NN Full-Blown Algorithm

- Keep a heap H of index entries and objects, ordered by MINDIST.
- Initially, H contains the root.
- Set $\delta = +\infty$.
- While $H \neq \emptyset$
 - Extract the element e with minimum MINDIST.
 - If it is an object, return it as NN.
 - For every entry se in $\text{PAGE}(e)$ whose $\text{MINDIST} \leq \delta$
 - Insert se into H .
 - Decrease δ to $\text{MINMAXDIST}(q, se)$ if possible.
- End while

Content

- The R-tree
 - Range Query
 - Aggregation Query
- NN Query
- **RNN Query**
- Closest Pair Query
- Close Pair Query
- Skyline Query

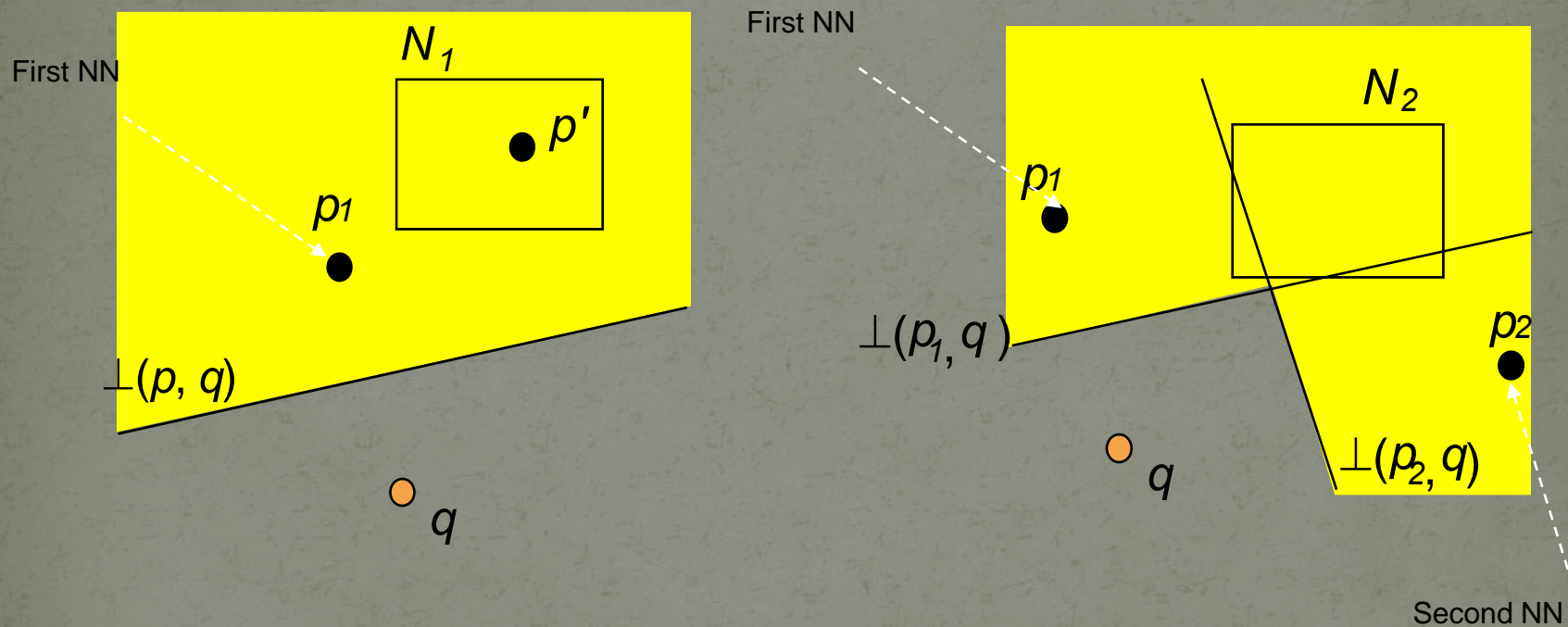
Reverse NN: Definition

- Given a set of points, and a query location q .
- Find the points whose NN is q .



- $RNN(q) = \{p_1, p_2\}$, $NN(q) = p_3$.

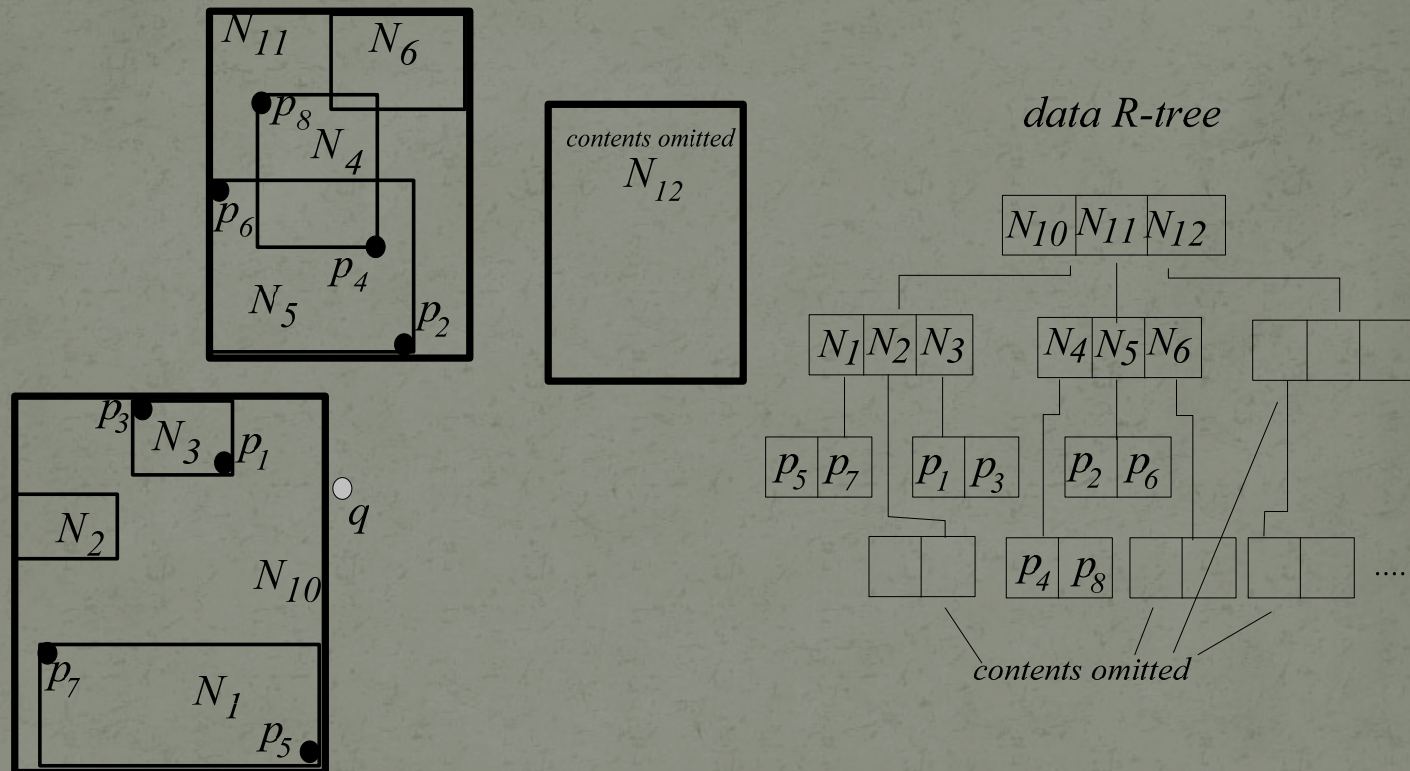
Half-plane pruning



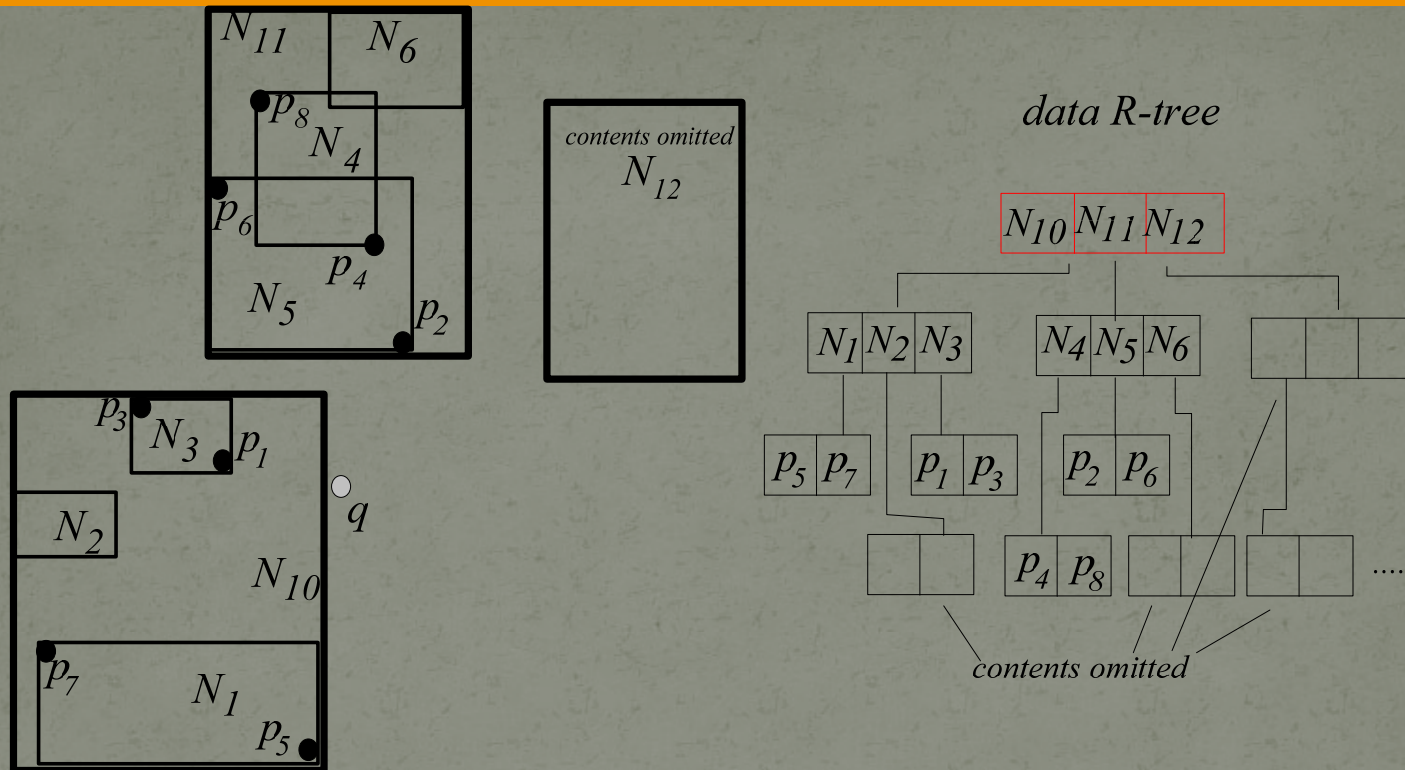
(TPL) Algorithm

- Two logical steps:
 - **Filter step: Find the set S_{cnd} of candidate points**
 - Find NN;
 - Prune space;
 - Find NN in unpruned space;
 - ...
 - Till no more object left.
 - **Refinement step: eliminate false positives**
 - For each point p in S_{cnd} , check whether its NN is not q .
- The two steps are combined in a single tree traversal, by keeping all pruned MBRs/objects in S_{rfn} .

Example



Example



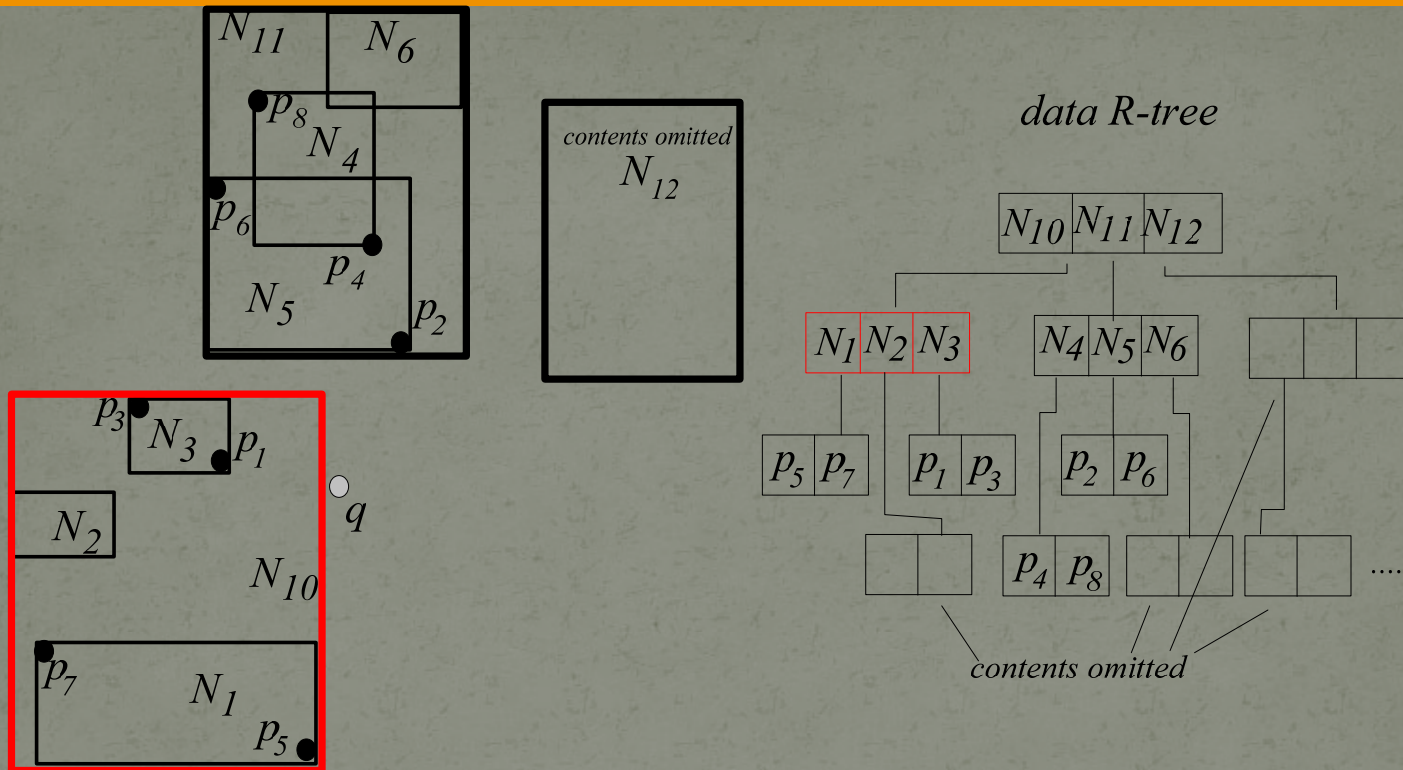
Action
visit root

Heap
 $\{N_{10}, N_{11}, N_{12}\}$

S_{cnd}
 \emptyset

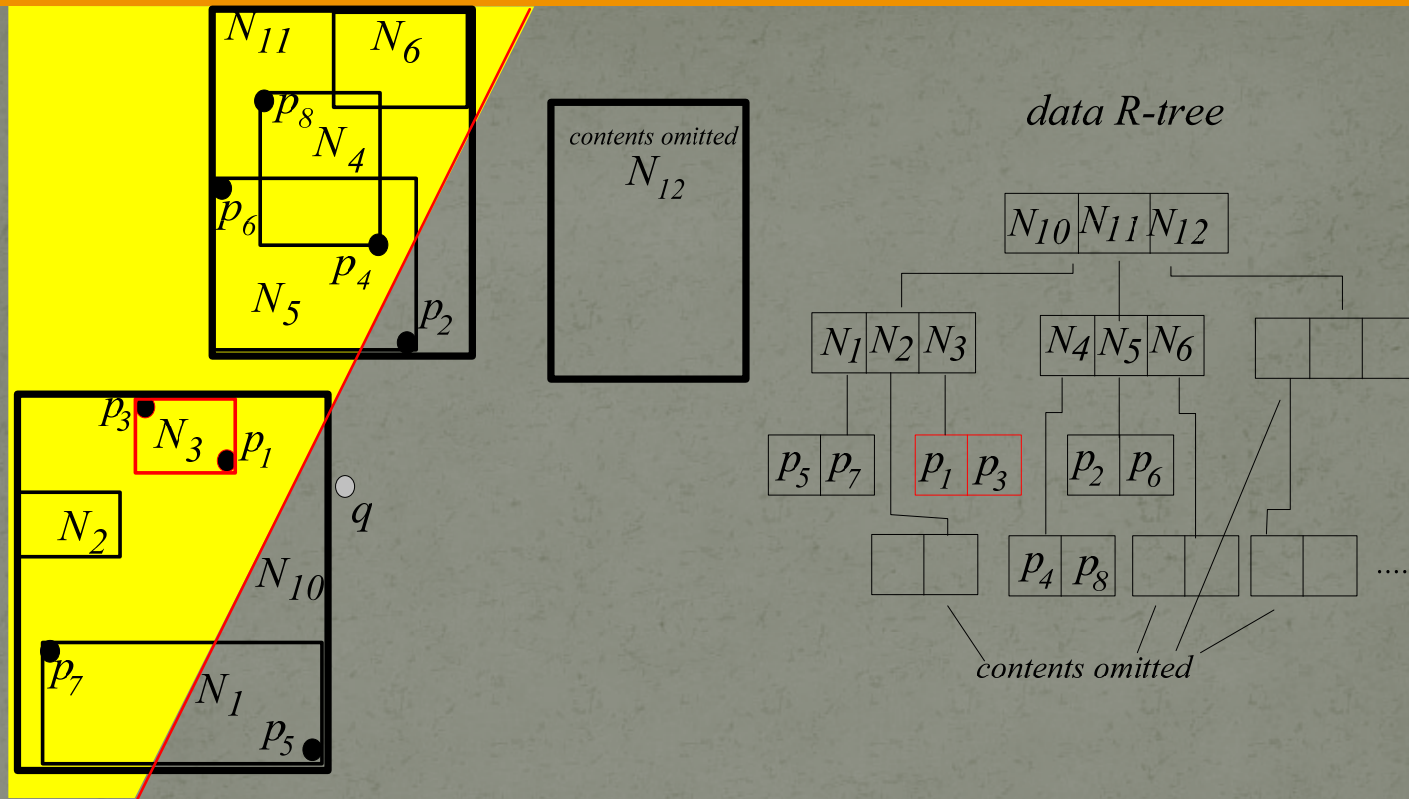
S_{rfn}
 \emptyset

Example



Action	Heap	S_{cnd}	S_{rfn}
visit N_{10}	$\{N_3, N_{11}, N_2, N_1, N_{12}\}$	\emptyset	\emptyset

Example



Action

visit N_3

Heap

$\{N_{11}, N_2, N_1, N_{12}\}$

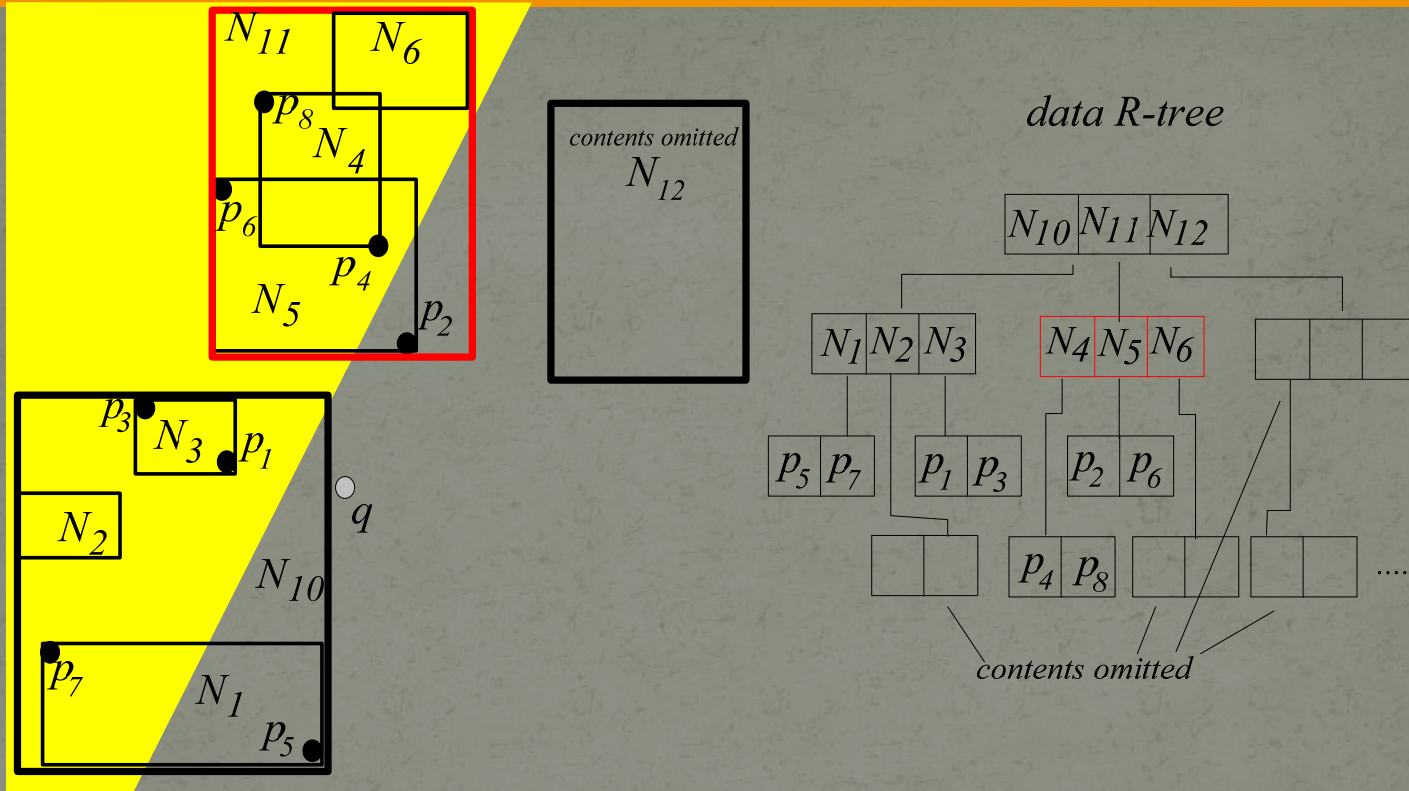
S_{cnd}

$\{p_1\}$

S_{rfn}

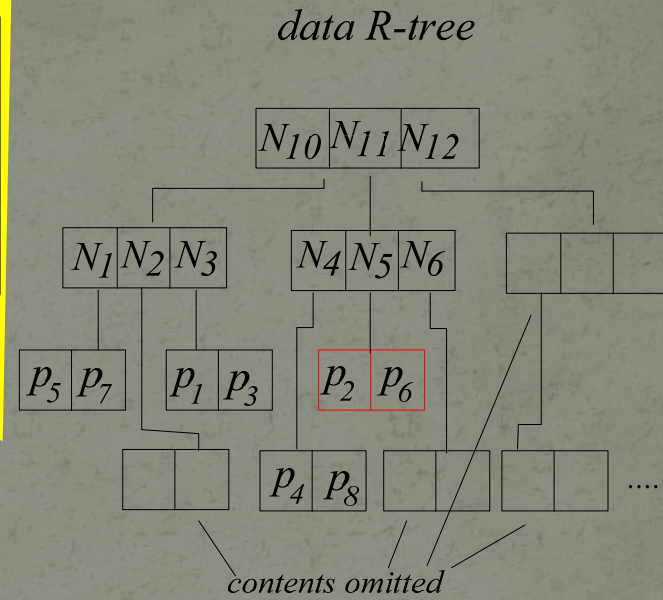
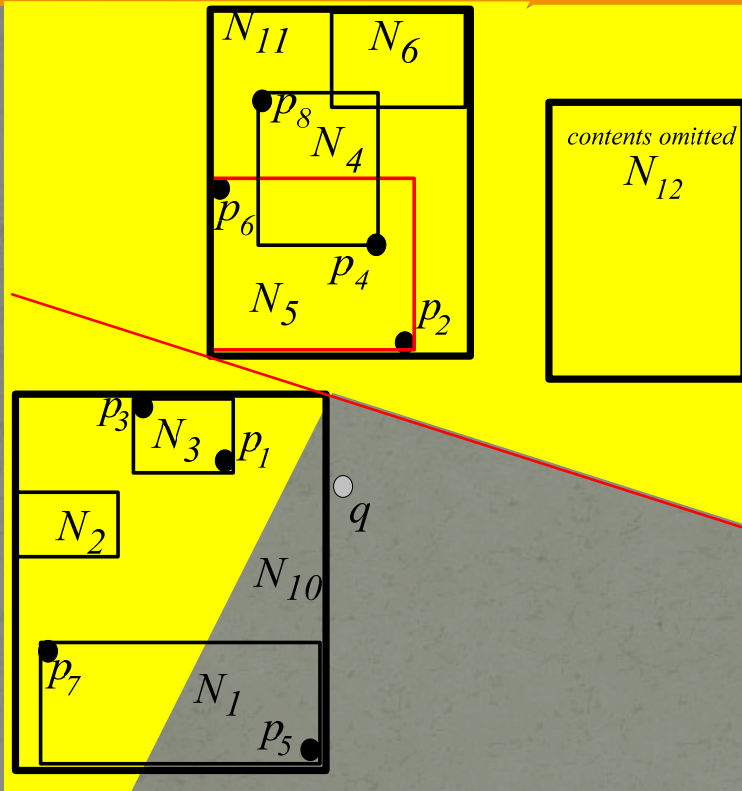
$\{p_3\}$

Example



Action	Heap	S_{cnd}	S_{rfn}
visit N_{11}	$\{N_5, N_2, N_1, N_{12}\}$	$\{p_1\}$	$\{p_3, N_4, N_6\}$

Example



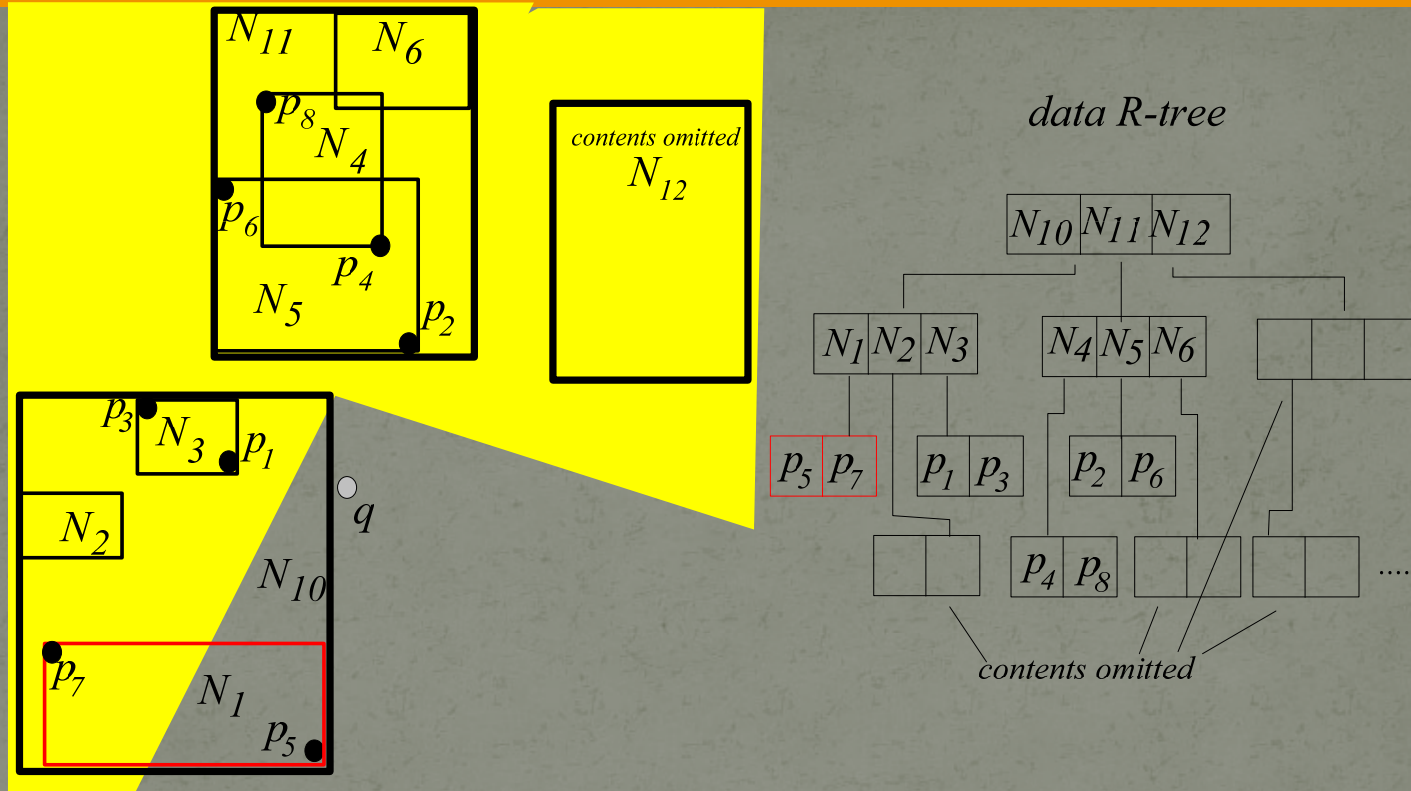
Action

visit N_5

Heap

$$\{N_2, N_1, N_{12}\}$$
 S_{cnd}
$$\{p_1, p_2\}$$
 S_{rfn}
$$\{p_3, N_4, N_6, p_6\}$$

Example (end of filter step)



Action

visit N_1

Heap

$\{N_{12}\}$

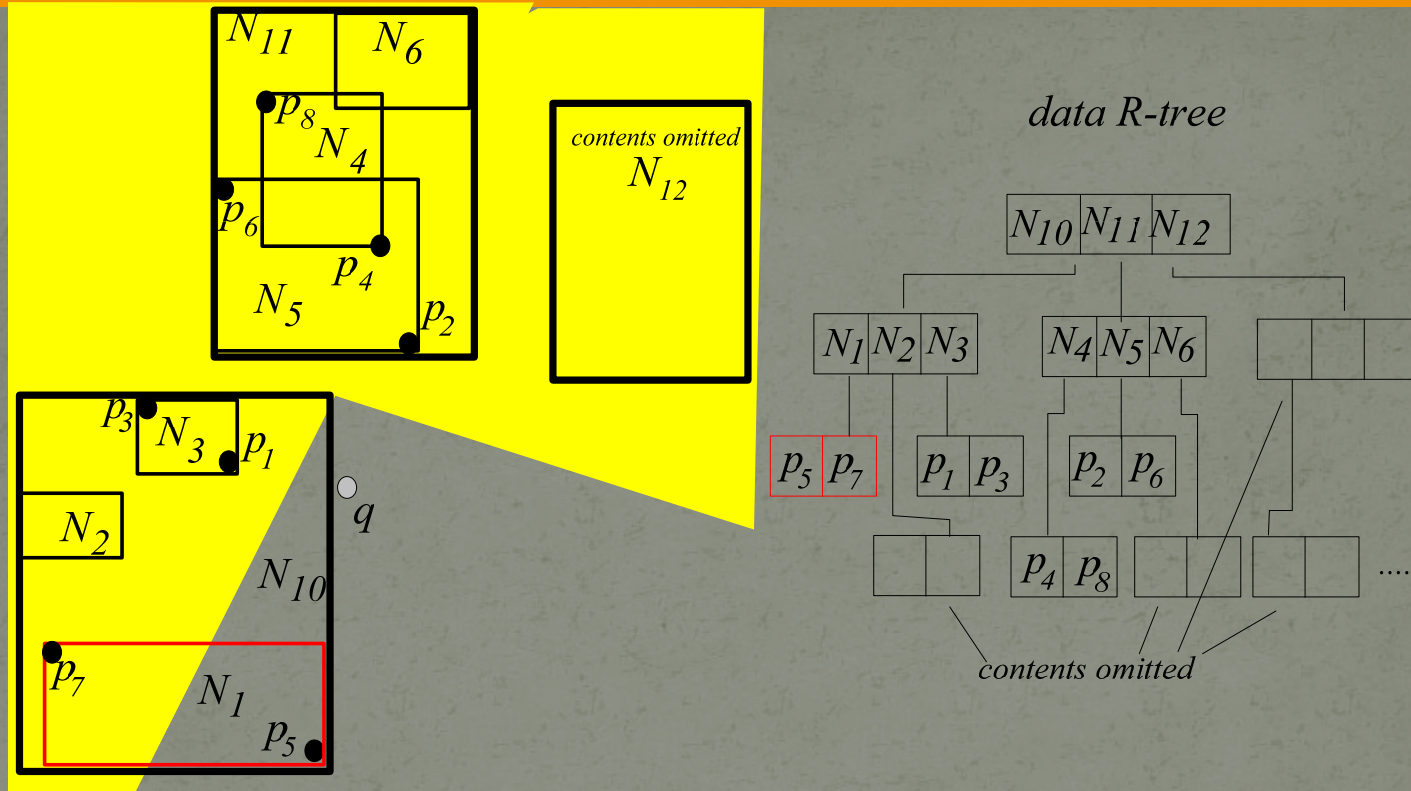
S_{cnd}

$\{p_1, p_2, p_5\}$

S_{rfn}

$\{p_3, N_4, N_6, p_6, N_2, p_7\}$

Example (end of filter step)



Action

Heap

\emptyset

S_{cnd}

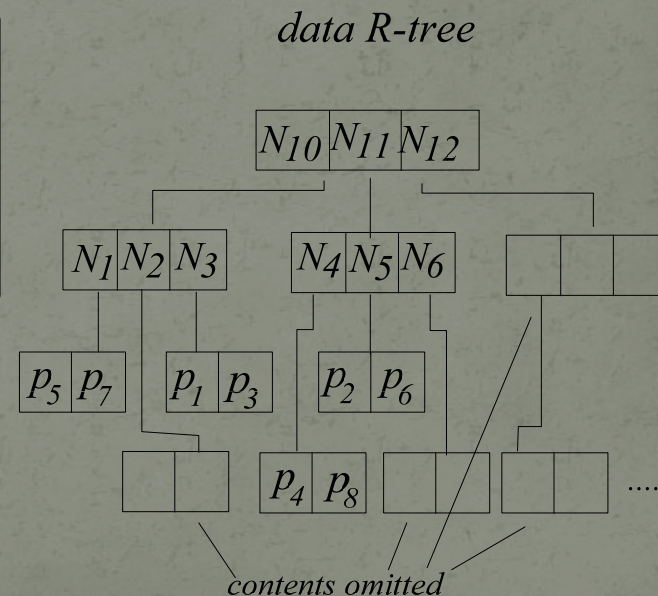
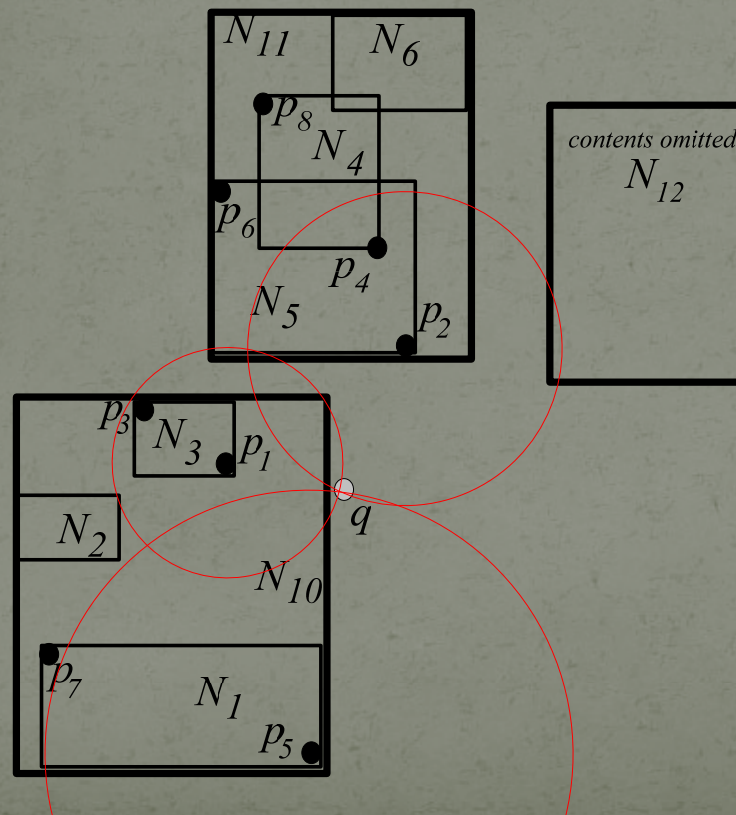
$\{p_1, p_2, p_5\}$

S_{rfn}

$\{p_3, N_4, N_6, p_6, N_2, p_7, N_{12}\}$

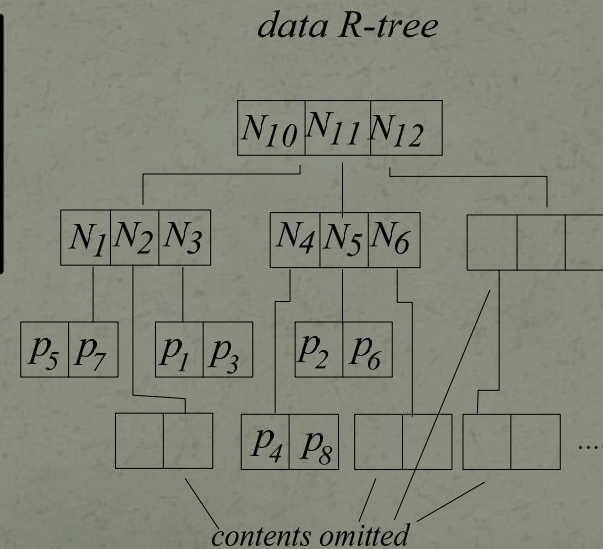
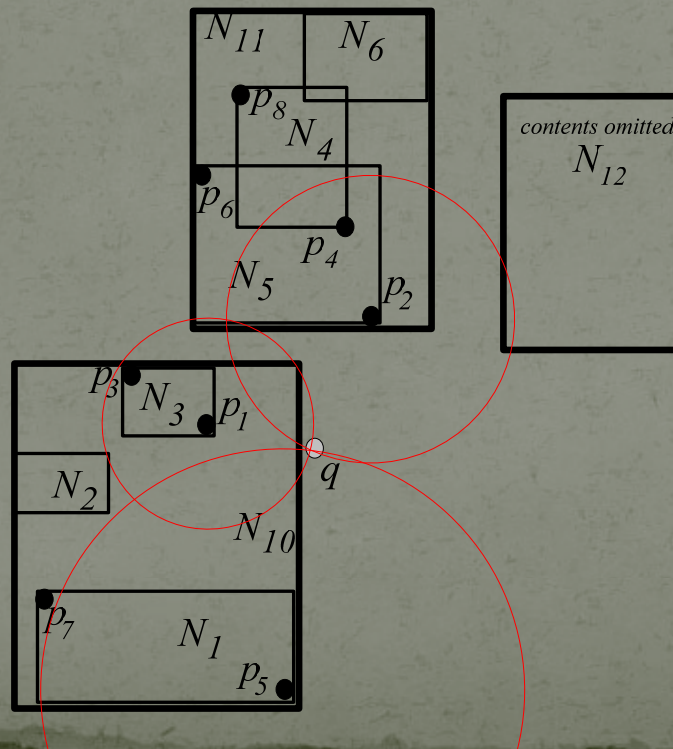
Example (refinement)

Action	S_{cnd}	S_{rfn}
	$\{p_1, p_2, p_5\}$	$\{p_3, p_6, p_7, N_4, N_6, N_2, N_{12}\}$
invalidate p_1	$\{\cancel{p_1}, p_2, p_5\}$	$\{N_4, N_6, N_2, N_{12}\}$
remove N_6, N_2	$\{\cancel{p_1}, p_2, p_5\}$	$\{N_4, N_{12}\}$



Example (refinement)

Action	S_{cnd}	S_{rfn}
	$2 \{ \cancel{p_1}, p_2, p_5 \}$	$\{ N_4, N_{12} \}$
access N_4	$\{ \cancel{p_1}, p_2, p_5 \}$	$\{ p_4, p_8, N_{12} \}$
invalidate p_2	$\{ \cancel{p_1}, \cancel{p_2}, p_5 \}$	$\{ N_{12} \}$
remove N_{12}	$\{ \cancel{p_1}, \cancel{p_2}, p_5 \}$	\emptyset

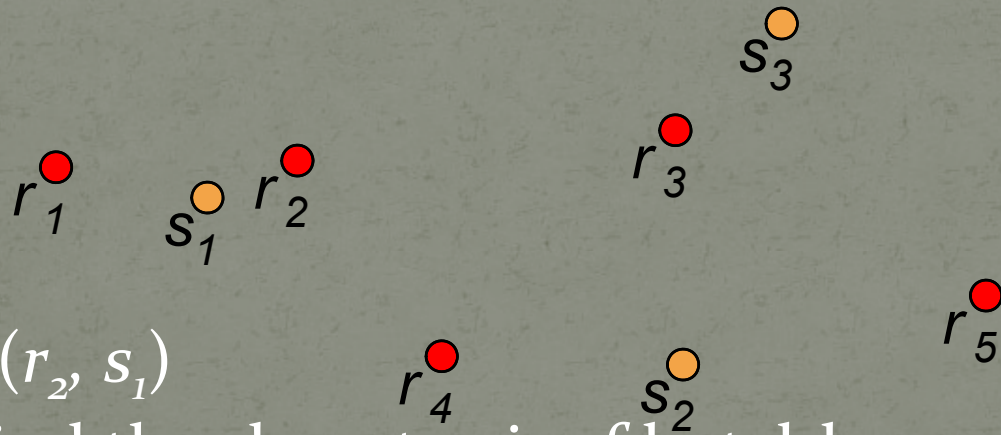


Content

- The R-tree
 - Range Query
 - Aggregation Query
- NN Query
- RNN Query
- **Closest Pair Query**
- Close Pair Query
- Skyline Query

Closest Pair (CP) Query

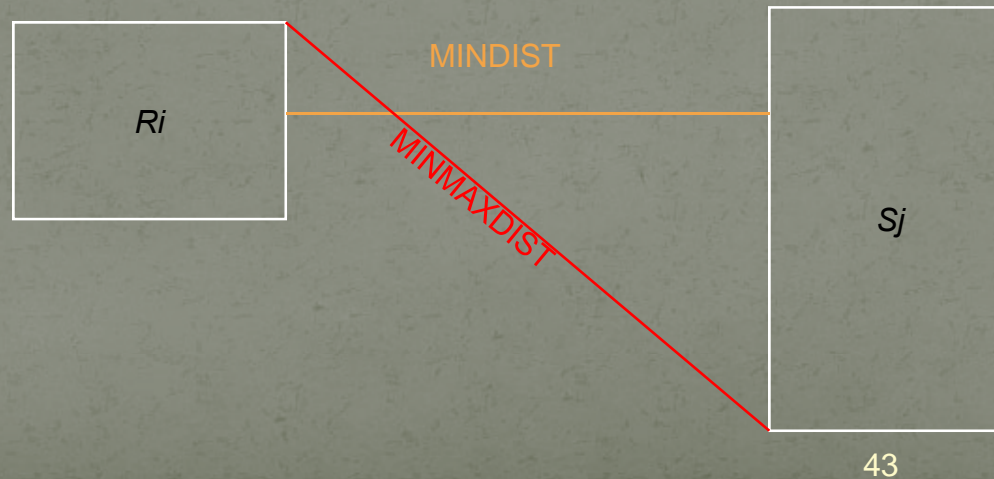
- Given two sets of objects R and S ,
- Find the pair of objects ($r \in R, s \in S$) with minimum distance.



- CP = (r_2, s_1)
- E.g. find the closest pair of hotel-bar.

CP Solution Idea

- Assume R and S are indexed by R-trees with same height.
- Similar to the NN query algorithm.
- MINDIST, MINMAXDIST for a pair of MBRs:



CP Basic Algorithm

- Keep a heap H of pairs of index entries and pairs of objects, ordered by MINDIST.
- Initially, H contains the pair of roots.
- While $H \neq \emptyset$
 - Extract the pair (e_R, e_S) with minimum MINDIST.
 - If it is a pair of objects, return it as CP.
 - For every entry se_R in $\text{PAGE}(e_R)$ and every entry se_S in $\text{PAGE}(e_S)$
 - Insert (e_R, e_S) into H .
- End while

CP Full-Blown Algorithm

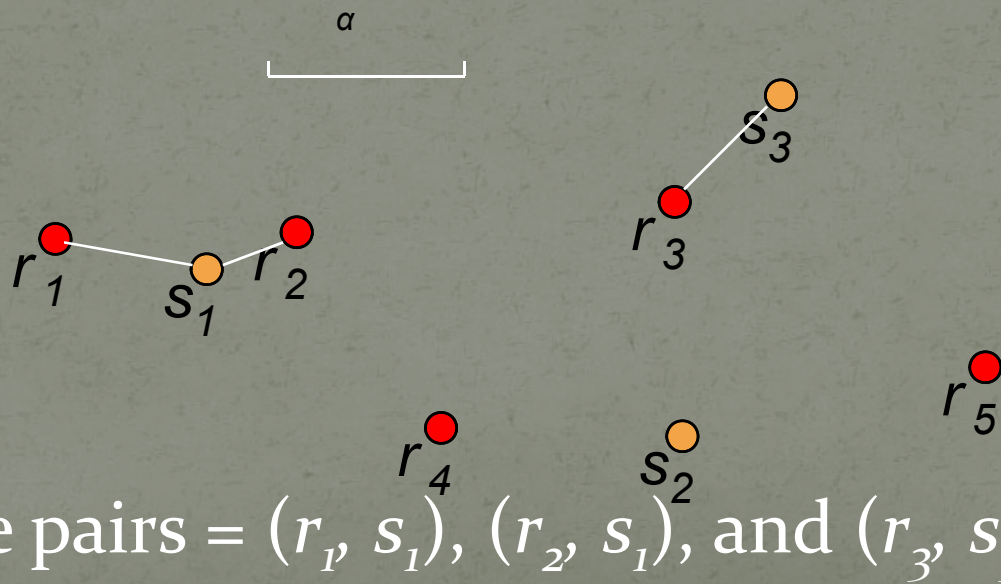
- Keep a priority queue H of pairs of index entries and pairs of objects, ordered by MINDIST.
- Initially, H contains the pair of roots.
- Set $\delta = +\infty$.
- While $H \neq \emptyset$
 - Extract the pair (e_R, e_S) with minimum MINDIST.
 - If it is a pair of objects, return it as CP.
 - For every entry se_R in $\text{PAGE}(e_R)$ and every entry se_S in $\text{PAGE}(e_S)$ whose $\text{MINDIST} \leq \delta$
 - Insert (se_R, se_S) into H .
 - Decrease δ to $\text{MINMAXDIST}(se_R, se_S)$ if possible.
- End while

Content

- The R-tree
 - Range Query
 - Aggregation Query
- NN Query
- RNN Query
- Closest Pair Query
- **Close Pair Query**
- Skyline Query

Close Pair Query

- Given two sets of objects R and S , plus a threshold α ,
- Find every pair of objects $(r \in R, s \in S)$ with distance $< \alpha$.



- Close pairs = (r_1, s_1) , (r_2, s_1) , and (r_3, s_3) .

Close Pair Solution Idea

- Observation: if $d(r, s) < \alpha$, $\forall mbr_R, mbr_S$ that contain r and s , respectively, we have: $MINDIST(mbr_R, mbr_S) < \alpha$.
- Solution idea:
 - start with the pair of root nodes,
 - Join pairs of index entries whose $MINDIST < \alpha$,
 - Till we reach leaf level.

Close Pair Algorithm

- Push the pair of root nodes into *stack*.
- While $stack \neq \emptyset$
 - Pop a pair (e_R, e_S) from *stack*.
 - For every entry se_R in $PAGE(e_R)$ and se_S in $PAGE(e_S)$ where $MINDIST(se_R, se_S) < \alpha$
 - Push (se_R, se_S) into *stack* if se_R is an index entry;
 - Otherwise report (se_R, se_S) as one close pair.
- End while

Content

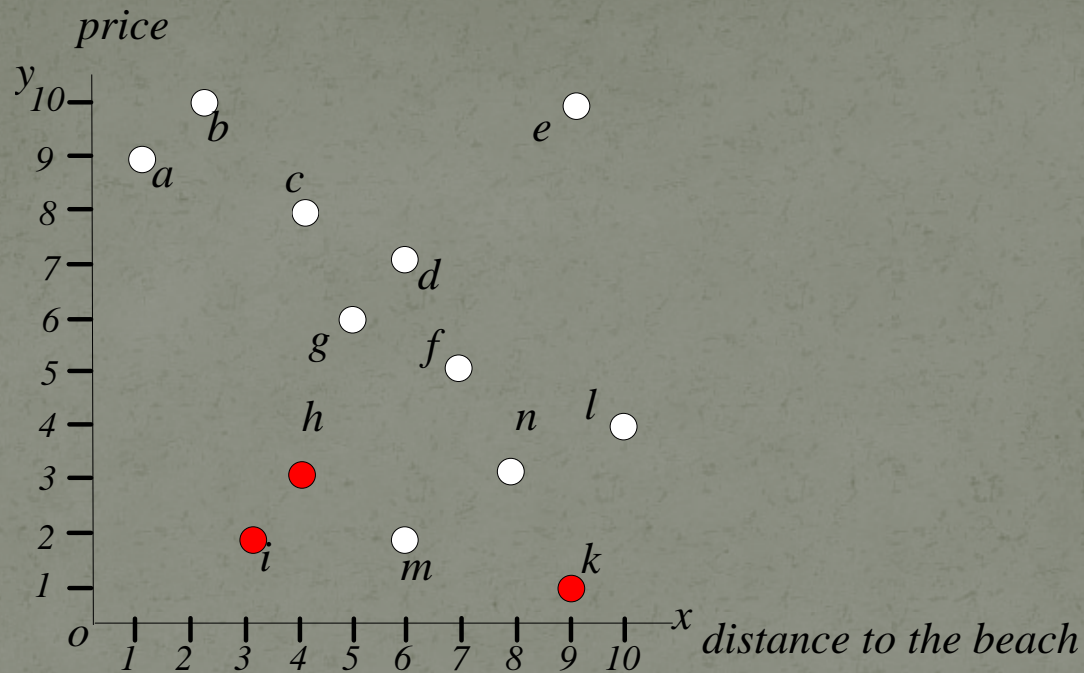
- The R-tree
 - Range Query
 - Aggregation Query
- NN Query
- RNN Query
- Closest Pair Query
- Close Pair Query
- Skyline Query

Skyline of Manhattan



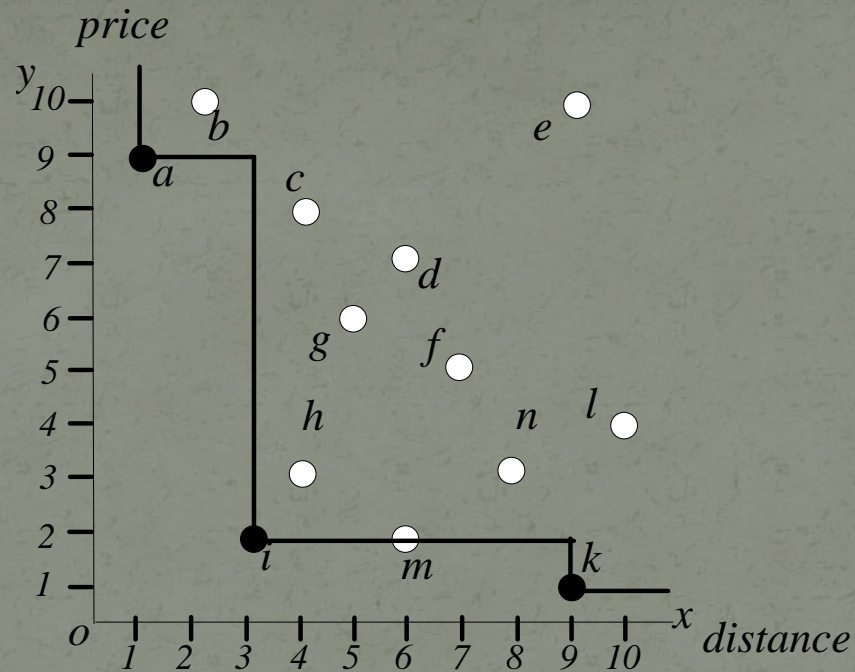
- Which buildings can we see?
 - *Higher or nearer*

Finding A Hotel Close to the Beach



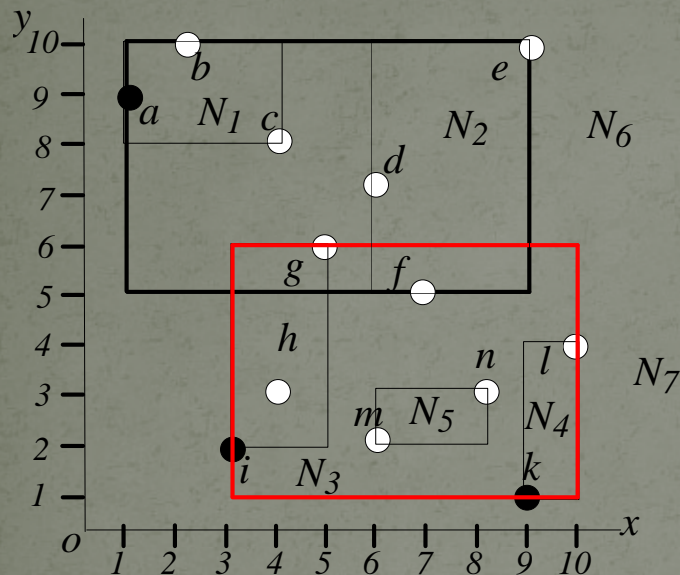
- Which one is better?
 - *i* or *h*? (*i*, because its price and distance **dominate** those of *h*)
 - *i* or *k*?

Skyline Queries

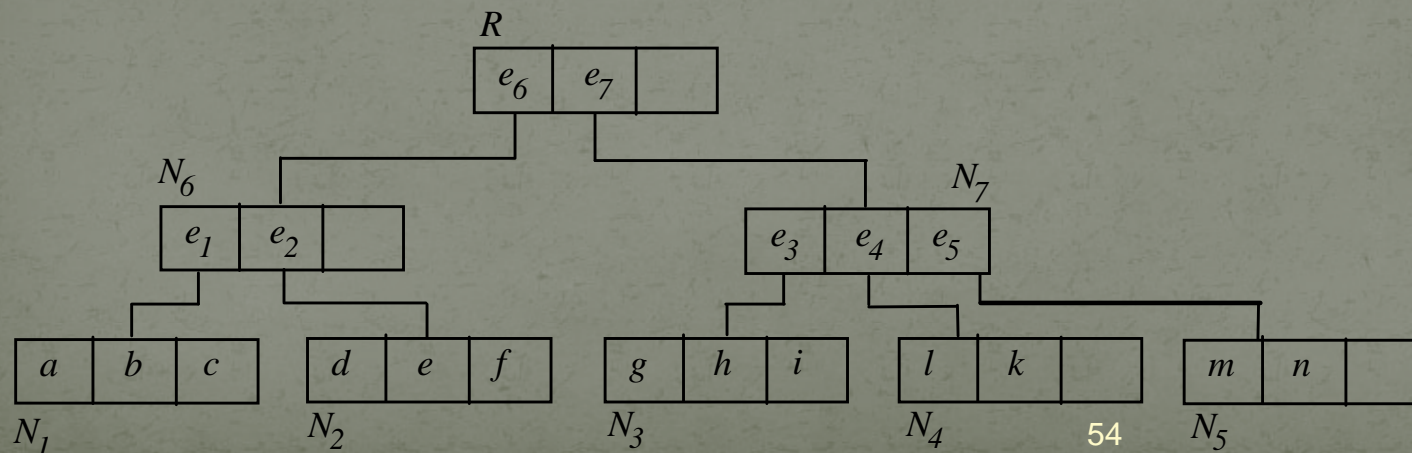


- Retrieve points **not dominated** by any other point.

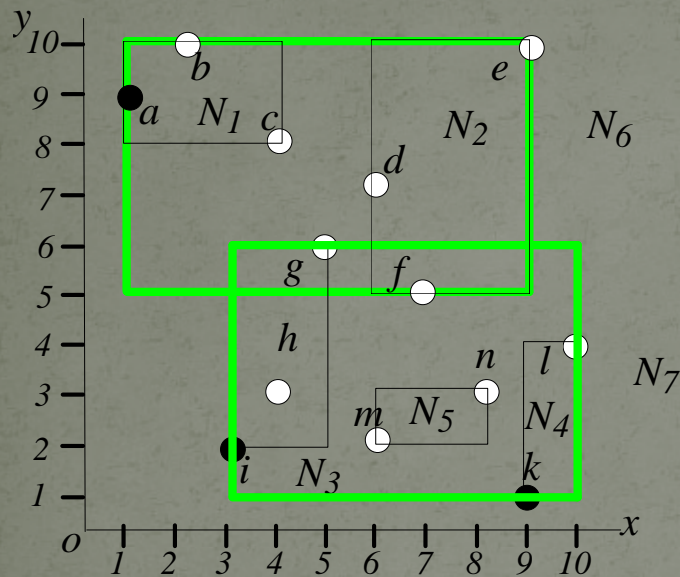
Branched and Bound Skyline (BBS)



- Assume all points are indexed in an R-tree.
- $\text{mindist}(\text{MBR}) = \text{the } L_1 \text{ distance between its lower-left corner and the origin.}$



Branched and Bound Skyline (BBS)

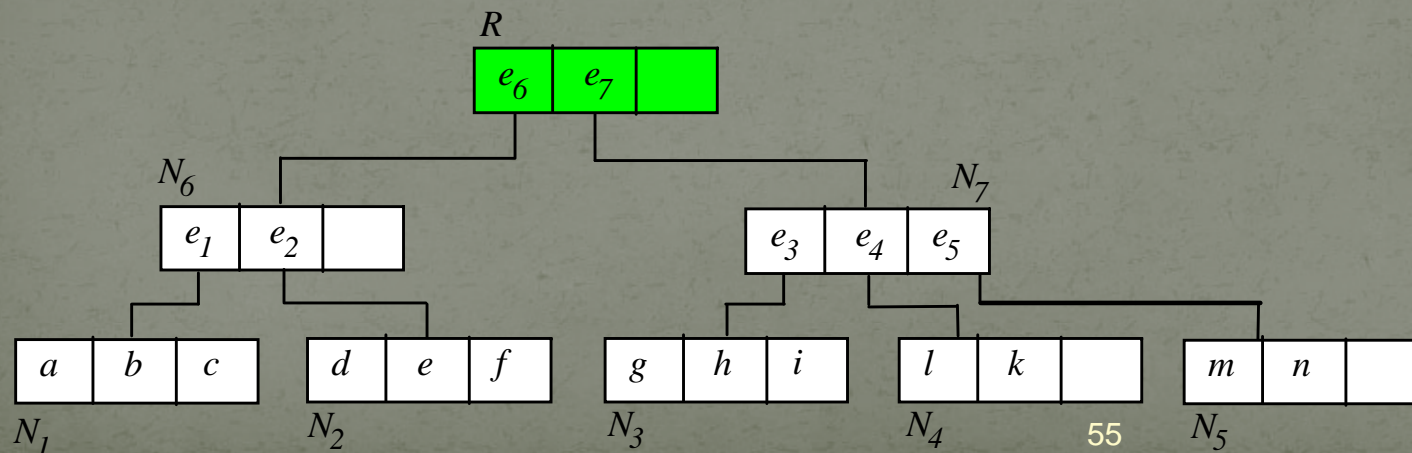


- Each heap entry keeps the mindist of the MBR.

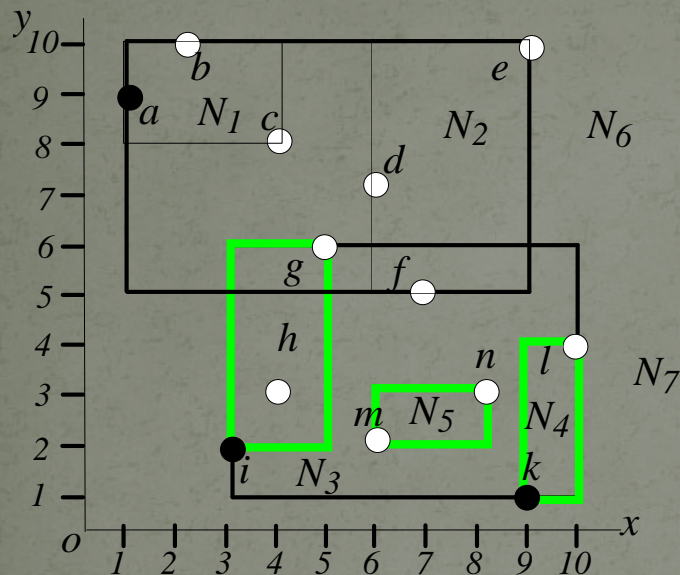
action
access root

heap contents
<e₇,4><e₆,6>

S
∅



Example of BBS



- Process entries in ascending order of their mindists.

action
access root
expand e_7

heap contents

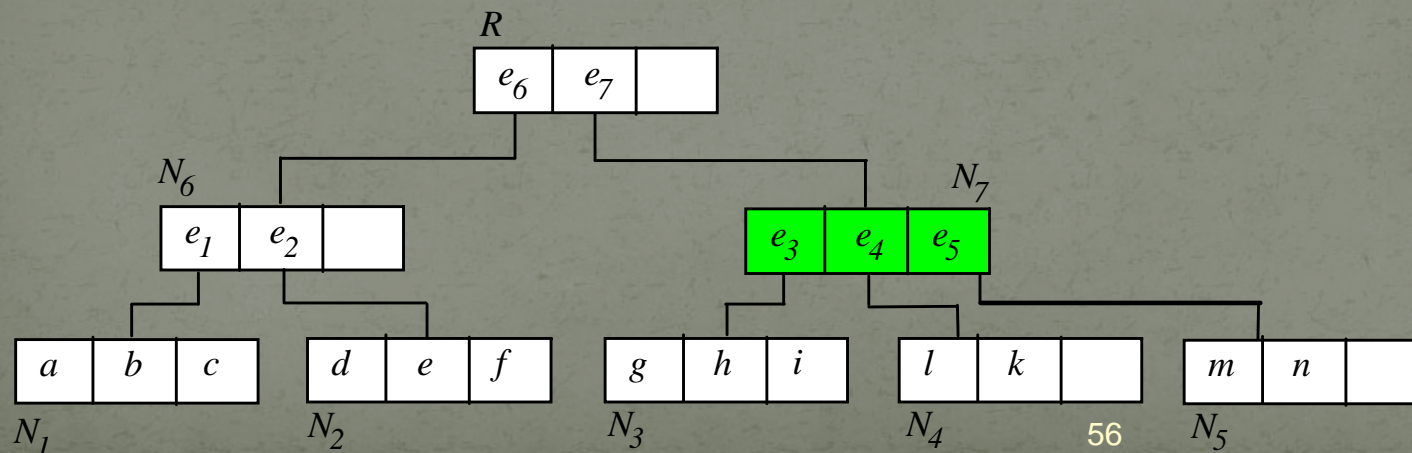
$\langle e_7, 4 \rangle \langle e_6, 6 \rangle$

$\langle e_3, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$

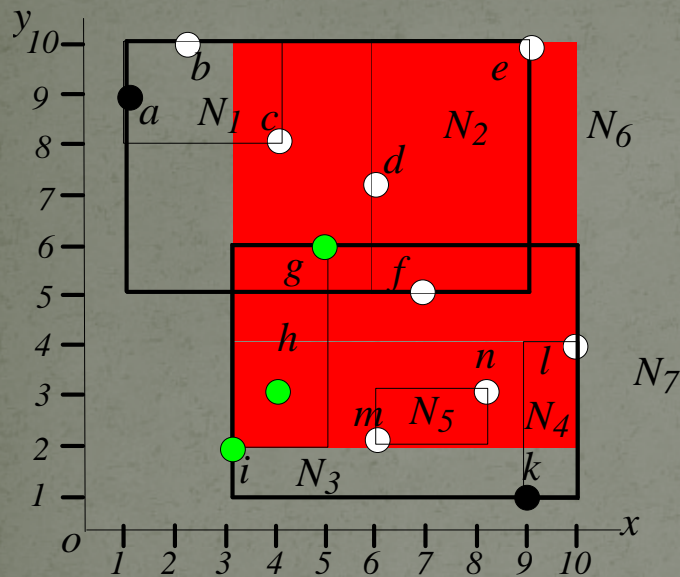
S

\emptyset

\emptyset



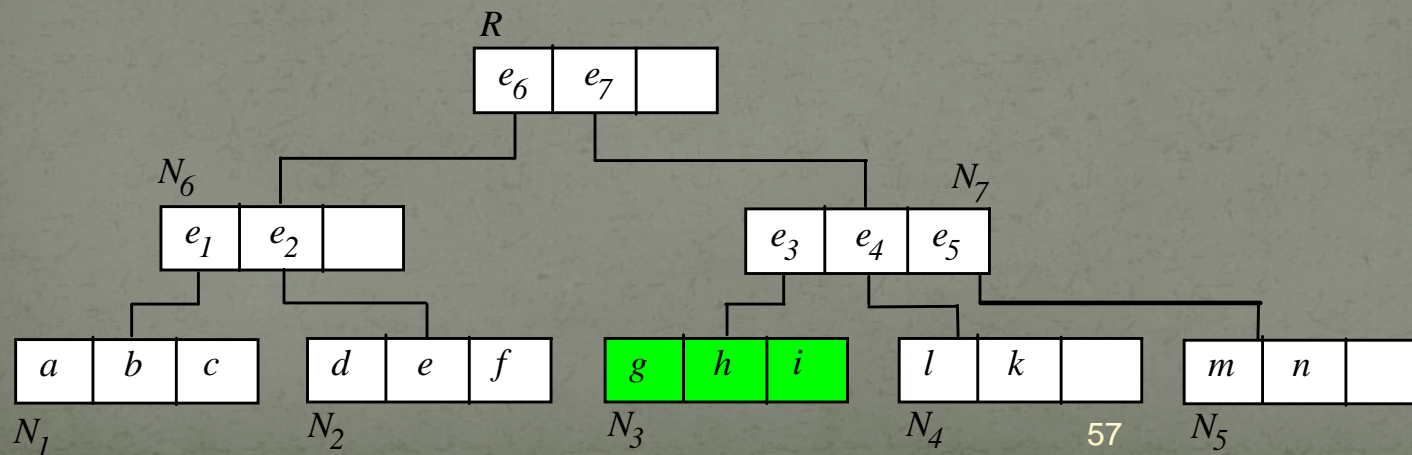
Example of BBS



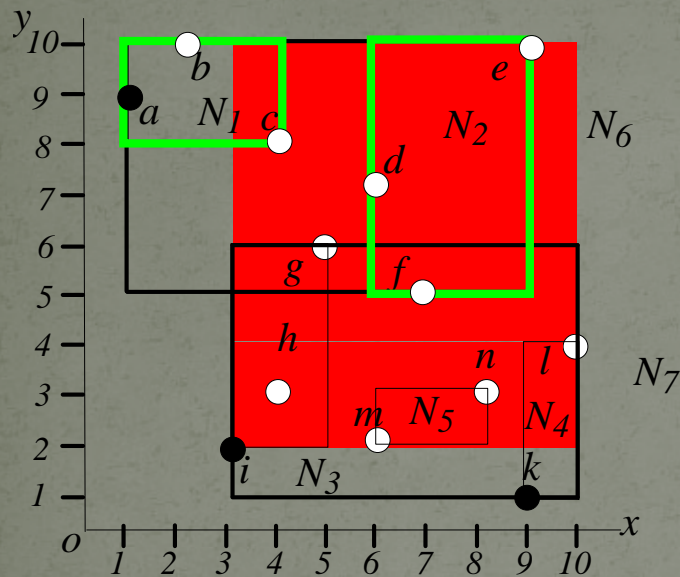
action
 access root
 expand e_7
 expand e_3

heap contents
 $\langle e_7, 4 \rangle \langle e_6, 6 \rangle$
 $\langle e_3, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$
 $\langle i, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$

S
 \emptyset
 \emptyset
 $\{i\}$



Example of BBS

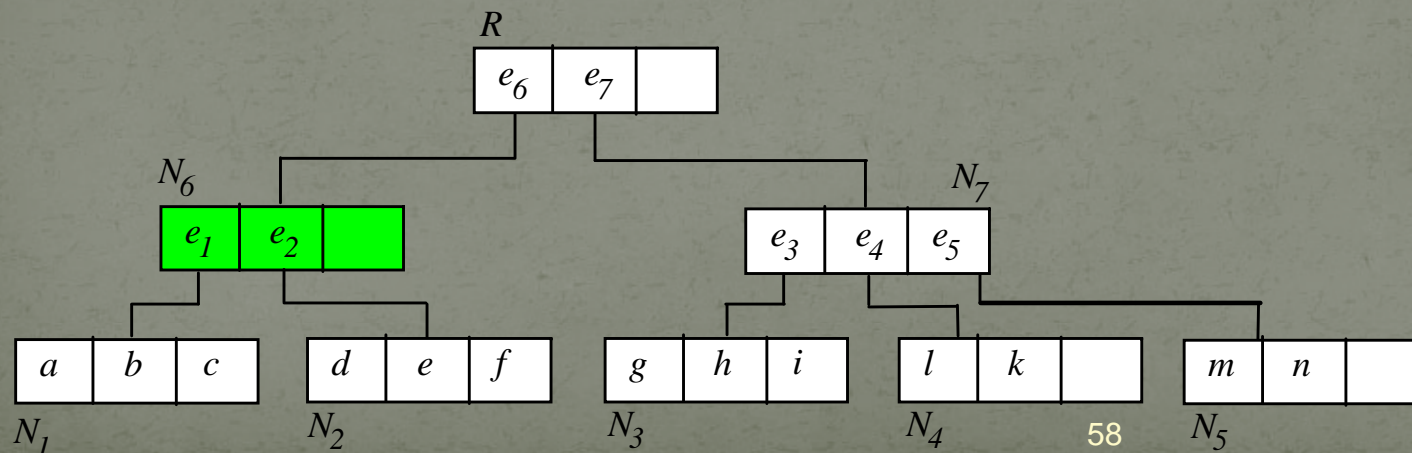


action
 access root
 expand e_7
 expand e_3
 expand e_6

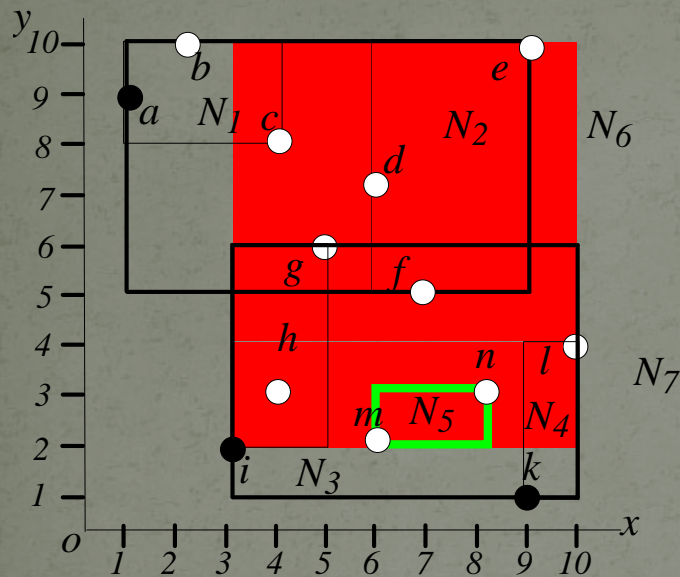
heap contents

$\langle e_7, 4 \rangle \langle e_6, 6 \rangle$
 $\langle e_3, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$
 $\langle i, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$
 $\langle e_5, 8 \rangle \langle e_1, 9 \rangle \langle e_4, 10 \rangle$

S
 \emptyset
 \emptyset
 $\{i\}$
 $\{i\}$



Example of BBS

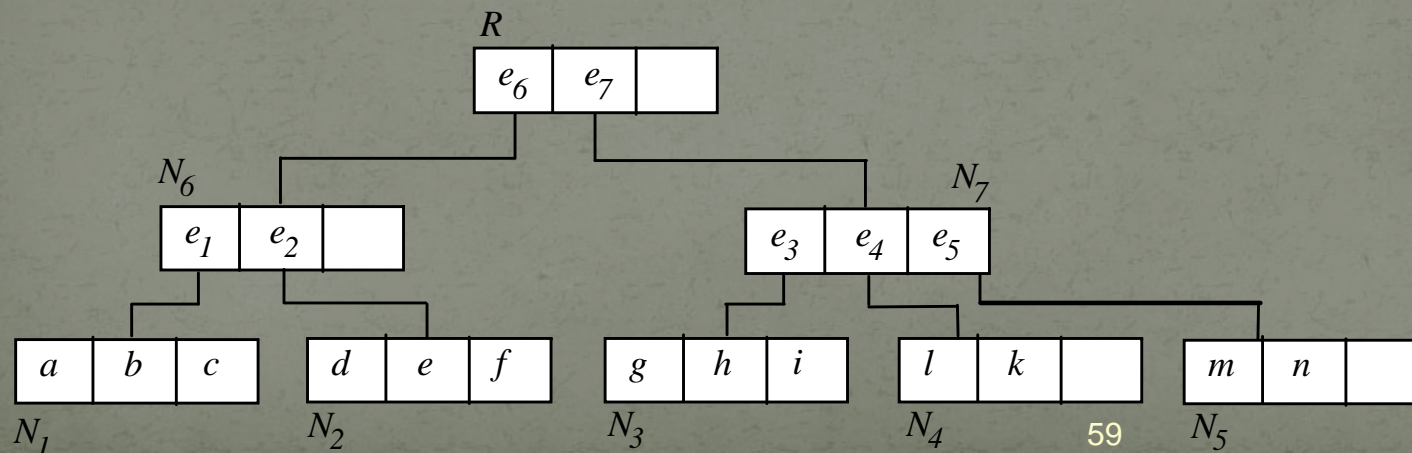


action
 access root
 expand e_7
 expand e_3
 expand e_6
 remove e_5

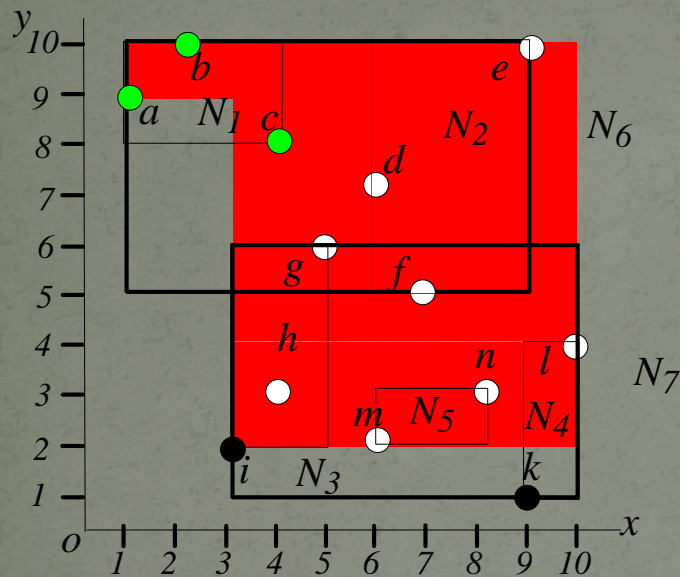
heap contents

$\langle e_7, 4 \rangle \langle e_6, 6 \rangle$
 $\langle e_3, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$
 $\langle i, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$
 $\langle e_5, 8 \rangle \langle e_1, 9 \rangle \langle e_4, 10 \rangle$
 $\langle e_1, 9 \rangle \langle e_4, 10 \rangle$

S
 \emptyset
 \emptyset
 $\{i\}$
 $\{i\}$
 $\{i\}$



Example of BBS



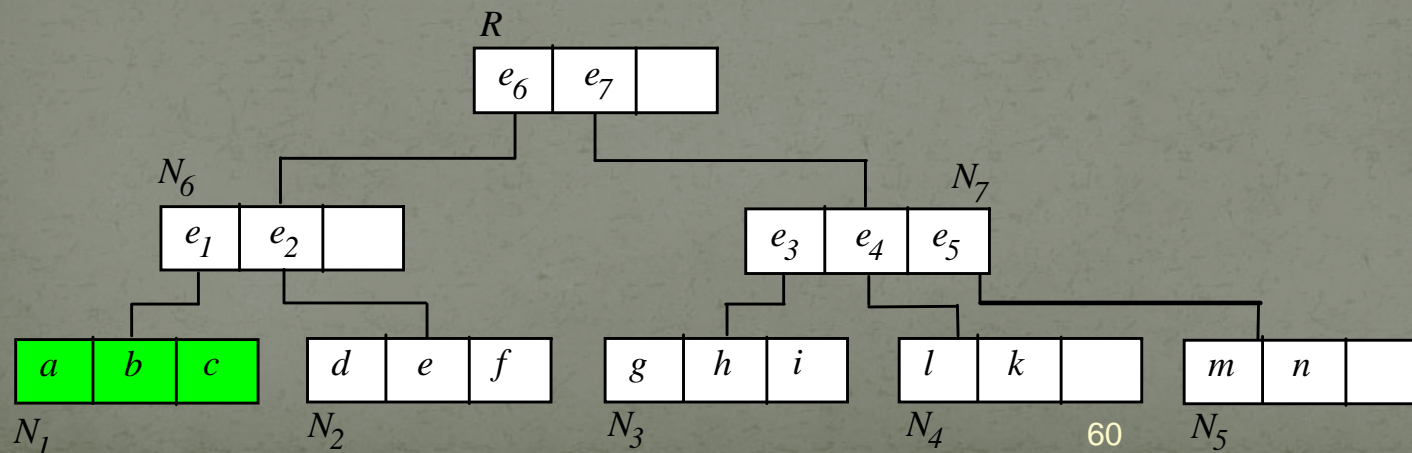
action
 access root
 expand e_7
 expand e_3
 expand e_6
 remove e_5
 expand e_1

heap contents

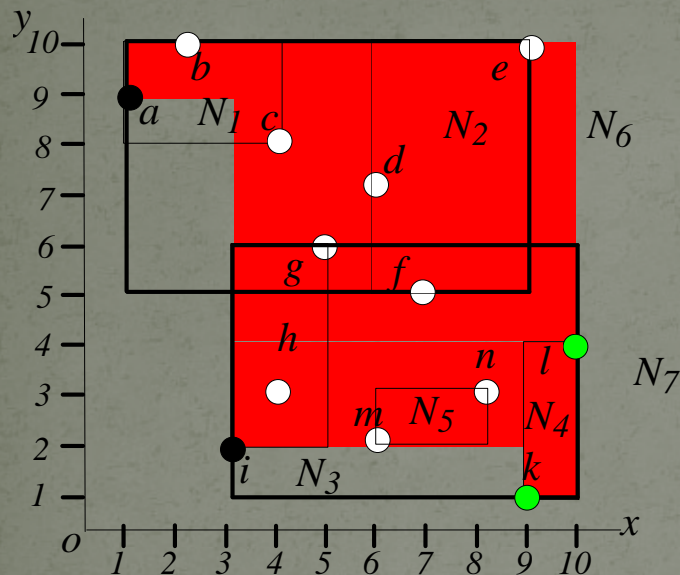
$\langle e_7, 4 \rangle \langle e_6, 6 \rangle$
 $\langle e_3, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$
 $\langle i, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$
 $\langle e_5, 8 \rangle \langle e_1, 9 \rangle \langle e_4, 10 \rangle$
 $\langle e_1, 9 \rangle \langle e_4, 10 \rangle$
 $\langle a, 10 \rangle \langle e_4, 10 \rangle$

S

\emptyset
 \emptyset
 $\{i\}$
 $\{i\}$
 $\{i\}$
 $\{i, a\}$



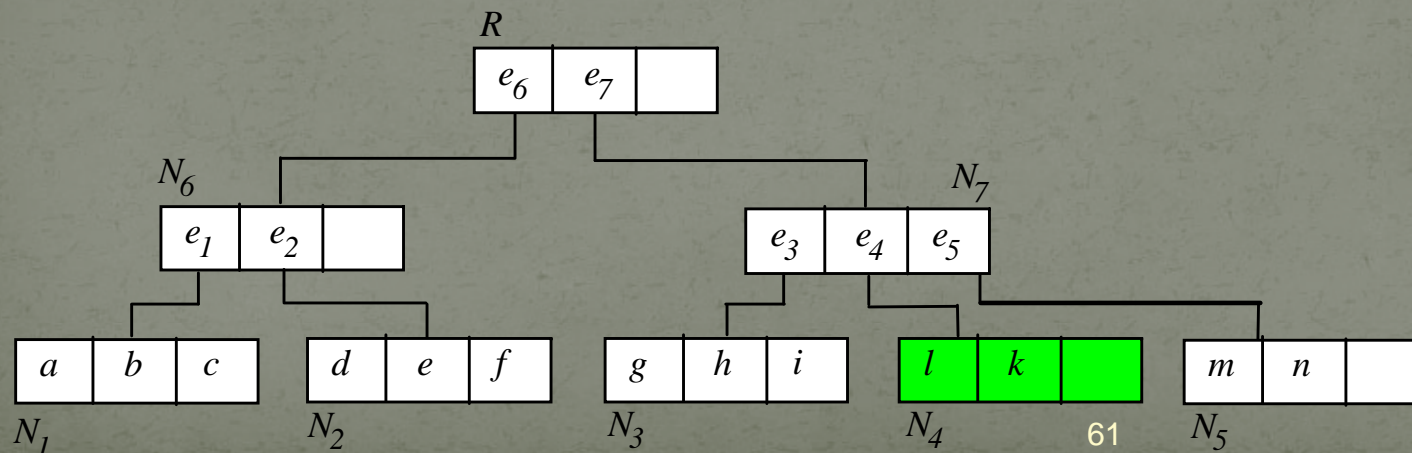
Example of BBS



action
 access root
 expand e_7
 expand e_3
 expand e_6
 remove e_5
 expand e_1
 expand e_4

heap contents
 $\langle e_7, 4 \rangle \langle e_6, 6 \rangle$
 $\langle e_3, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$
 $\langle i, 5 \rangle \langle e_6, 6 \rangle \langle e_5, 8 \rangle \langle e_4, 10 \rangle$
 $\langle e_5, 8 \rangle \langle e_1, 9 \rangle \langle e_4, 10 \rangle$
 $\langle e_1, 9 \rangle \langle e_4, 10 \rangle$
 $\langle a, 10 \rangle \langle e_4, 10 \rangle$
 $\langle k, 10 \rangle$

S
 \emptyset
 \emptyset
 $\{i\}$
 $\{i\}$
 $\{i\}$
 $\{i, a\}$
 $\{i, a, k\}$



Summary & Reference

- NN Query: N. Roussopoulos, S. Kelley and F. Vincent, “Nearest Neighbor Queries”, SIGMOD’95. (It didn’t talk about best-first search).
- RNN Query: Y. Tao, D. Papadias and X. Lian, “Reverse kNN Search in Arbitrary Dimensionality”, VLDB’04.
- Skyline Query: D. Papadias, Y. Tao, G. Fu, and B. Seeger, “An Optimal and Progressive Algorithm for Skyline Queries”, SIGMOD’03.
- Also talked about Closest/Close Pair Queries.