

CSEG601 & CSE5601

Spatial Data Management & Application :

RNN Query Processing method using R-tree-

Sungwon Jung

Big Data Processing & Database Lab
Dept. of Computer Science and Engineering
Sogang University
Seoul, Korea
Tel: +82-2-705-8930
Email : jungsung@sogang.ac.kr

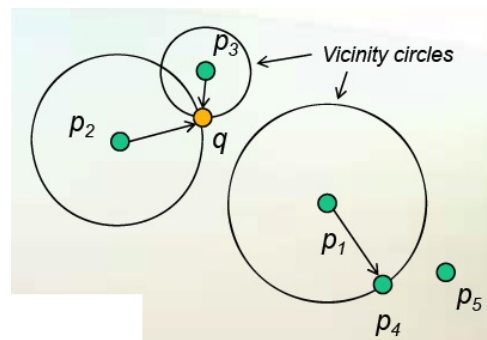
1

RNN Definition

- A data point p is the reverse nearest neighbor of query point q , if there is no point p' such that $\text{dist}(p', p) < \text{dist}(q, p)$, i.e. q is the NN of p .

$$\text{NN}(p_2) = \text{NN}(p_3) = q$$

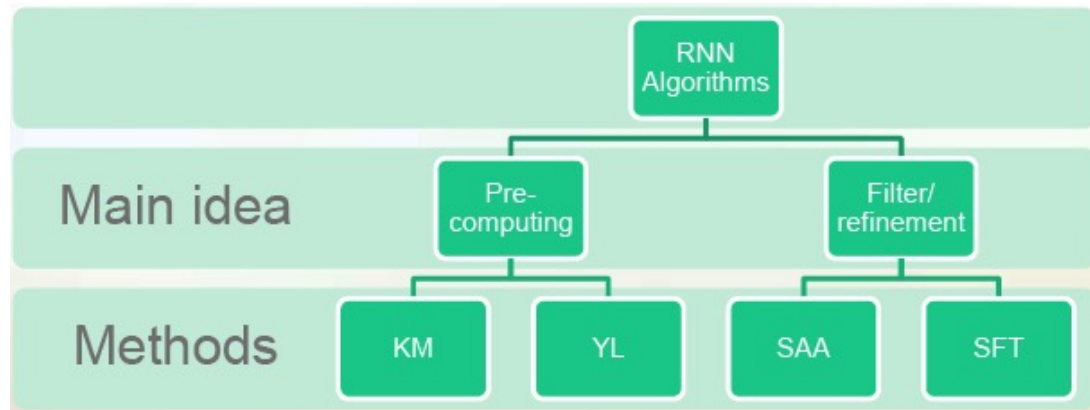
$$\text{RNN}(q) = \{p_2, p_3\}$$



- In our example, p_2, p_3 are the houses for which q is the nearest restaurant
- Is RNN a symmetric relation?

2

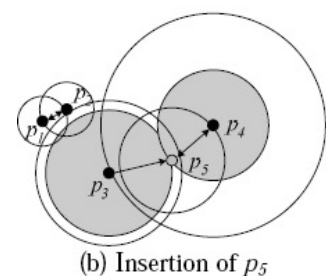
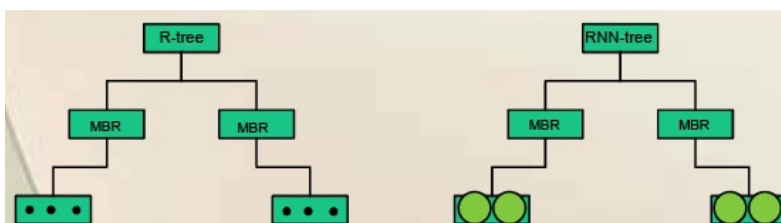
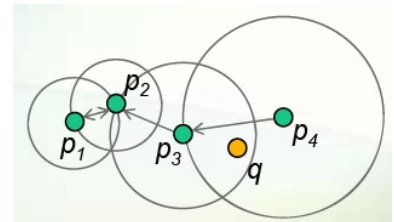
Related Works



3

Original RNN method [KM00]

- For all p :
 1. Pre-compute $NN(p)$
 2. Represent p as a vicinity circle
 3. Index the MBR of all circles by an R-tree(Named RNN-tree)
 4. $RNN(q)$ = all circles that contain q
- Needs two trees: RNN-tree & R-tree



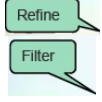
4

SAA

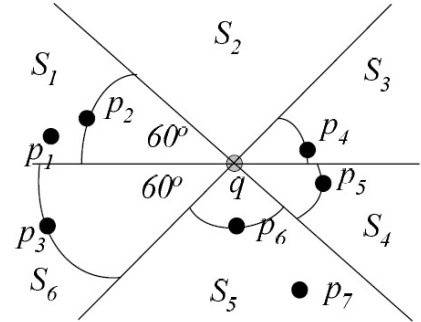
- Elimination of the need for pre-computing all NNs in filter/refinement methods

- SAA:

- Divide the space around query into six equal regions
- Find $NN(q)$ in all regions (candidate keys)
- Either (i) or (ii) holds for each candidate key p



- (i) p is in $RNN(q)$
- (ii) No $RNN(q)$ in S_i
- $RNN(q) = \{p_6\}$



- Any Drawbacks?

- The number of regions increases exponentially with the dimensionality

5

SFT



1. Find the k NNs of the query q (k candidates)



2. Eliminate the points that are closer to other candidates than q .

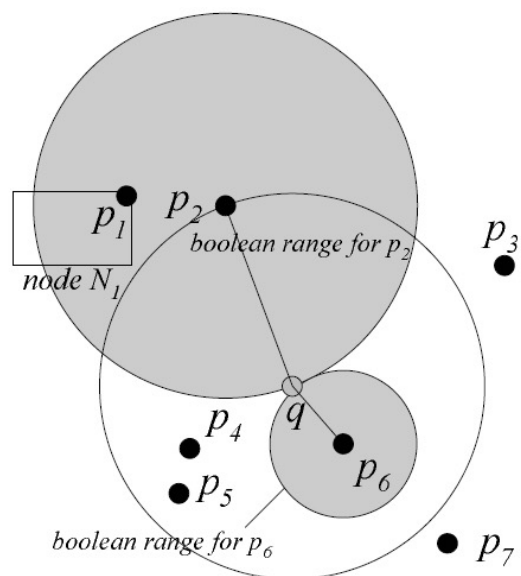


3. Apply *Boolean range queries* to determine the actual RNNs

- A Boolean range query terminates immediately when
 1. the first data point is found
 2. The entire side of a node MBR lies within a circle
 e.g., $\text{minmaxdist}(N_1, p_2) \leq \text{dist}(p_2, q)$

- Drawbacks?

- False misses
- Choosing a proper k

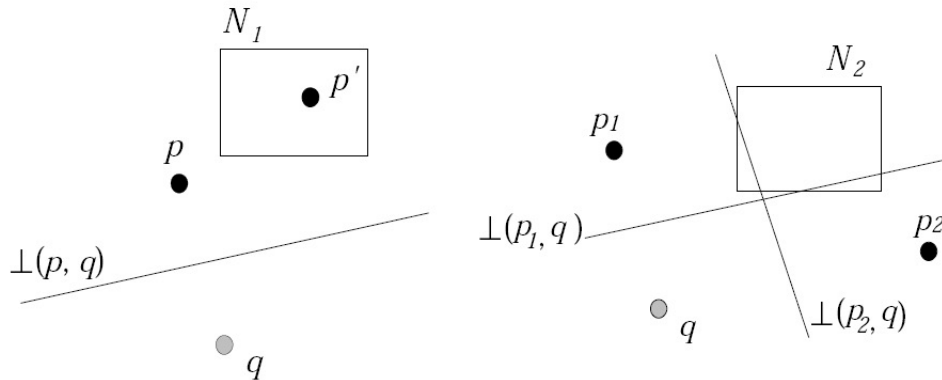


Step 1: Given $k=4$, Find 4NN Candidates Set = $\{p_6, p_5, p_4, p_2\}$
 Step 2: Discard p_4 and p_6 since they are closer to each other than q
 Step 3: Boolean Range Queries ($p_2, \text{dist}(p_2, q)$) and ($p_6, \text{dist}(p_6, q)$)
 Step 4: Discard p_2 since $\text{minmaxdist}(N_1, p_2) \leq \text{dist}(p_2, q)$
 False miss: p_3 due to $k=4$.

6

Half-plane pruning

- Can p' be closer to q than p ?

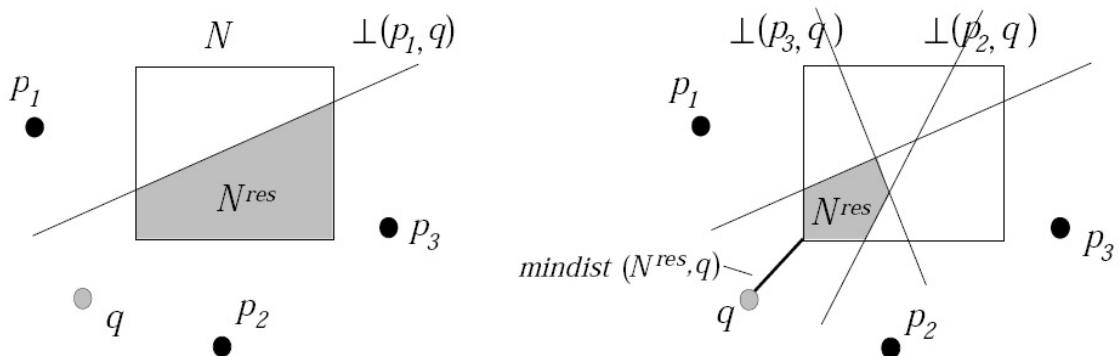


(a) Pruning with one point (b) Pruning with two points

- If p_1, p_2, \dots, p_n n data points, then any node whose MBR falls inside $\bigcup_{i=1, \dots, n} PL_{p_i}(p_i, q)$ cannot contain any RNN result.

7

Computing the residual region

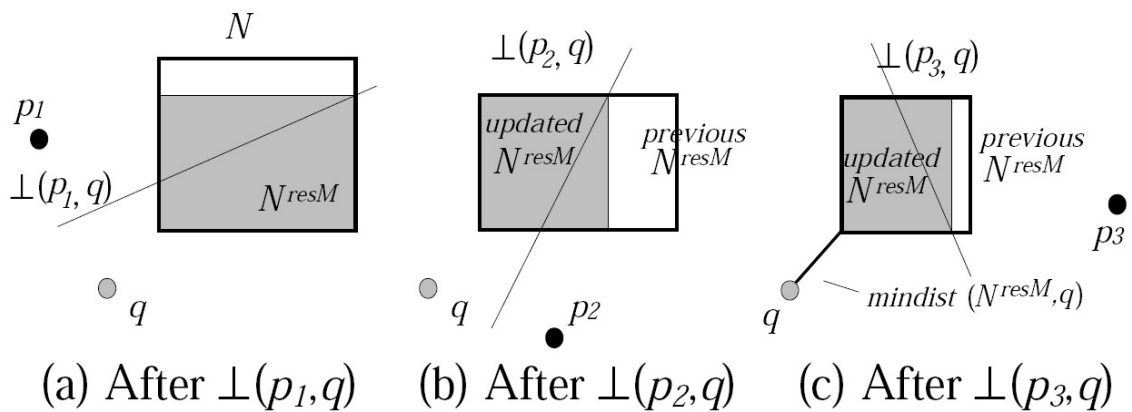


(a) After processing $\perp(q, p_1)$ (b) The final polygon

- $O(n^2)$ processing time in terms of bisector trimming for computing N^{res}
- Computation of intersections does not scale with dimensionality

8

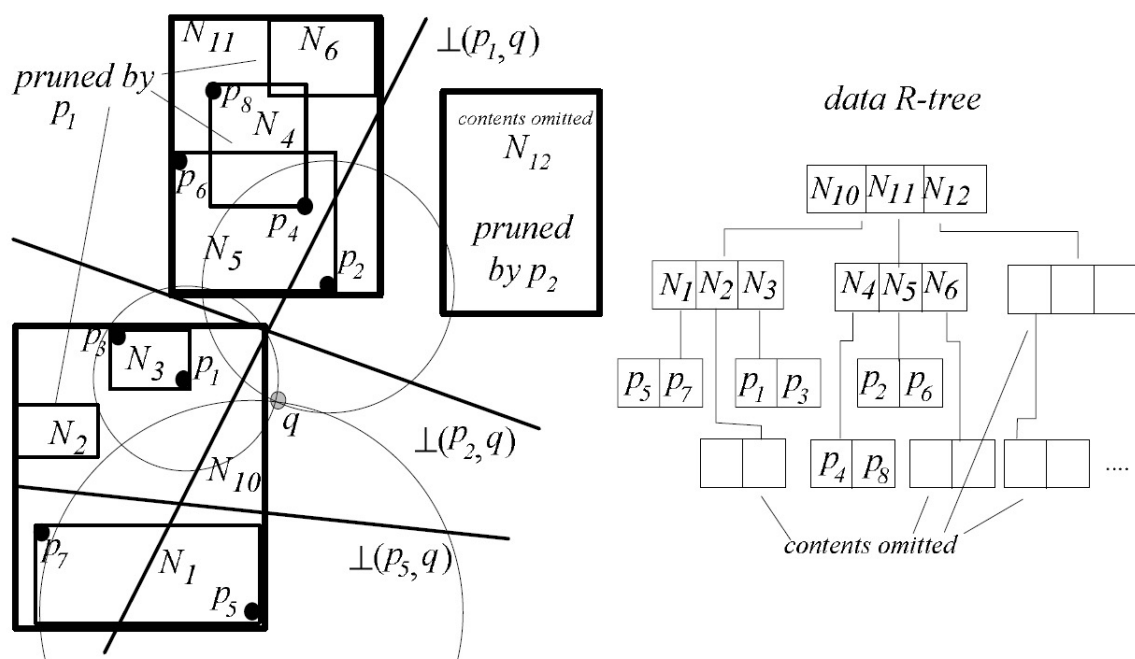
Approximating the residual MBR



- An MBR can be pruned if its residual region is empty
- The approximation is a superset of the real residual region
- We can prune an MBR if its approximate residual is empty
- Good news:
 - $O(n)$ processing time for computing N^{resM}
 - No more hyper-polyhedrons to make the intersection computation complex

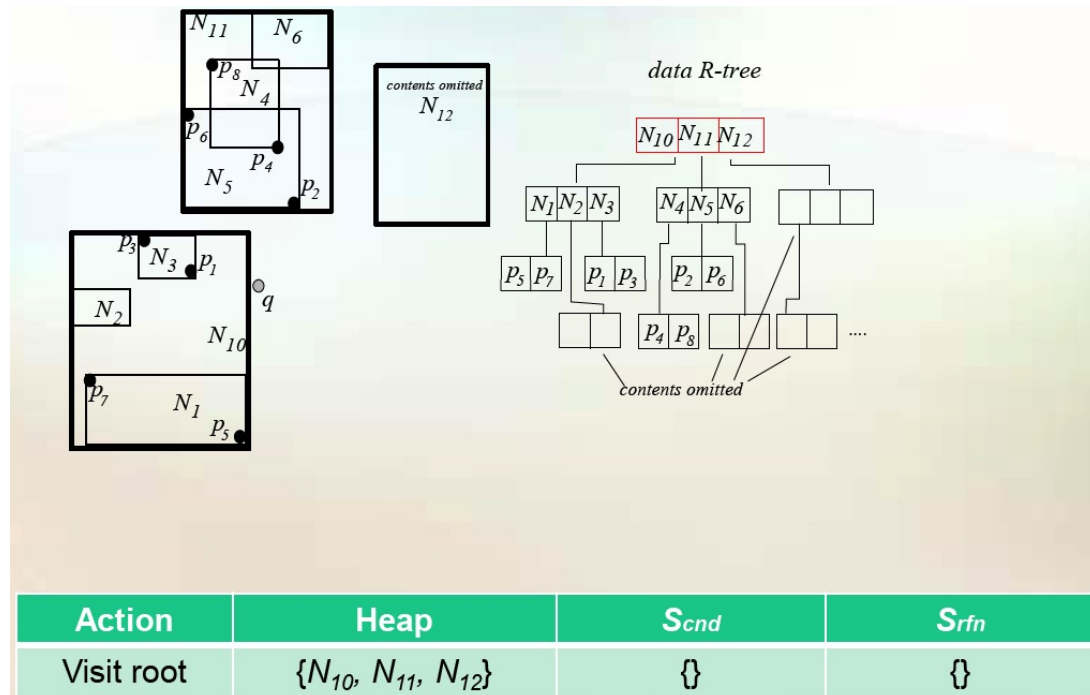
9

TPL algorithm for Single RNN – example(1)



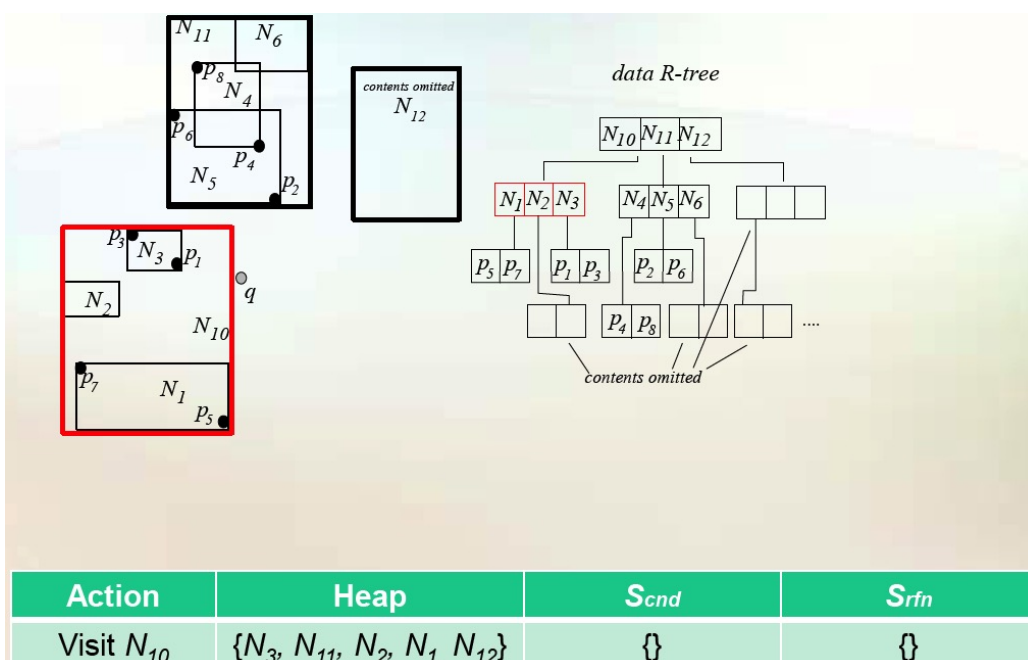
10

Flitering step - 1



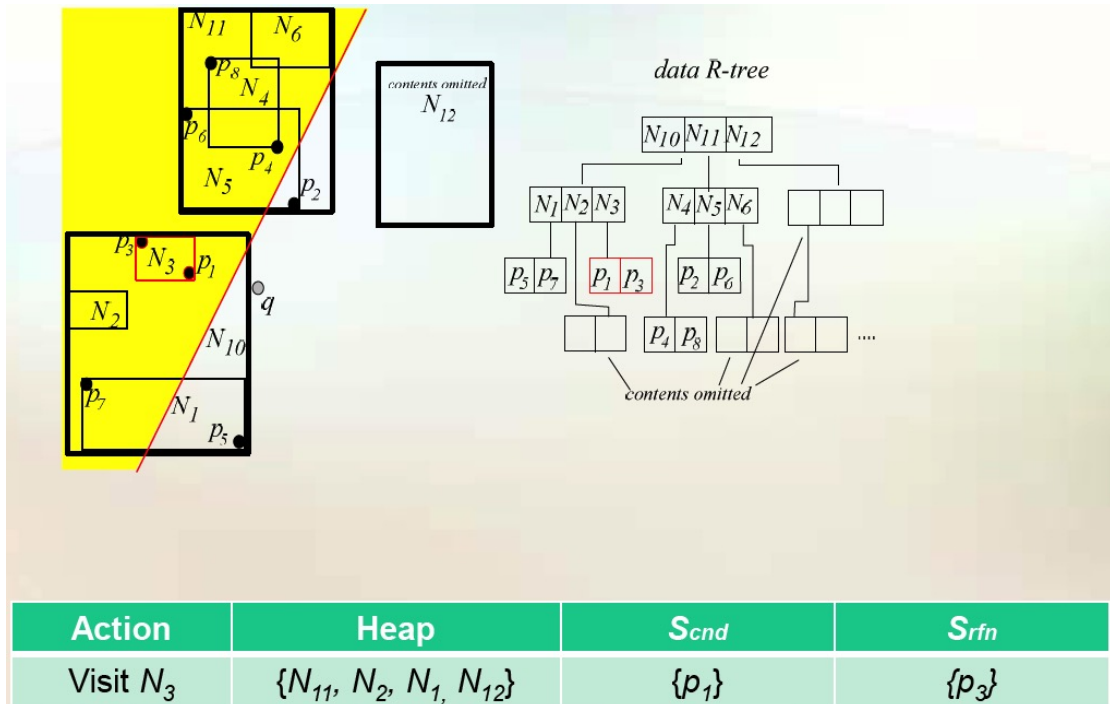
11

Flitering step - 2



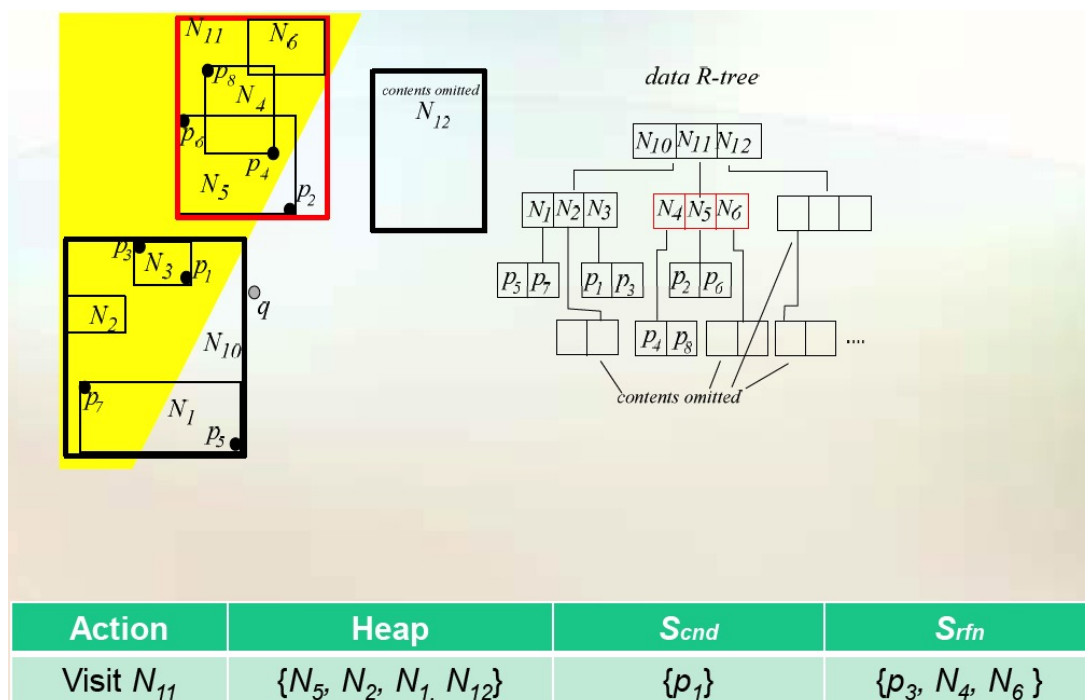
12

Flitering step - 3



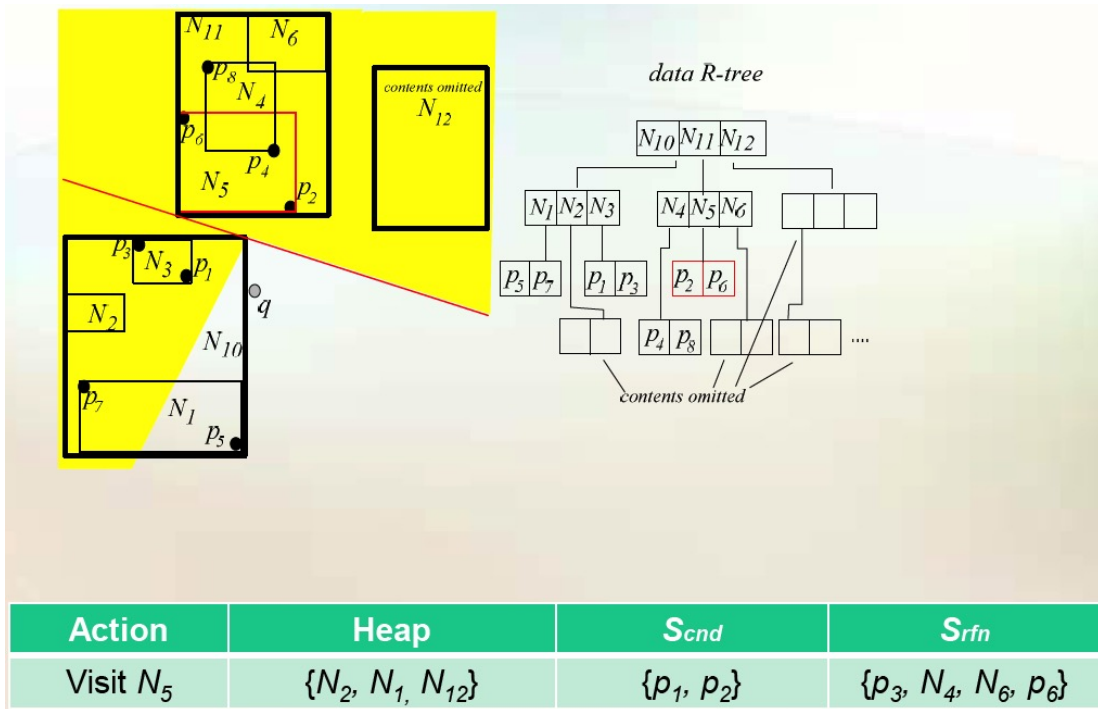
13

Flitering step - 4



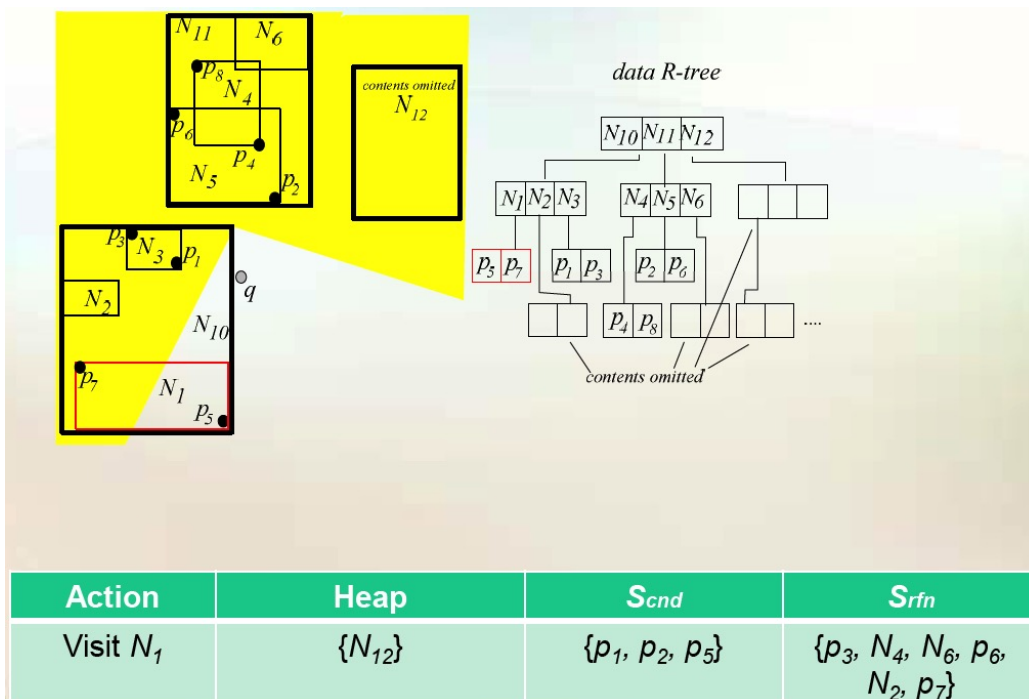
14

Flitering step - 5



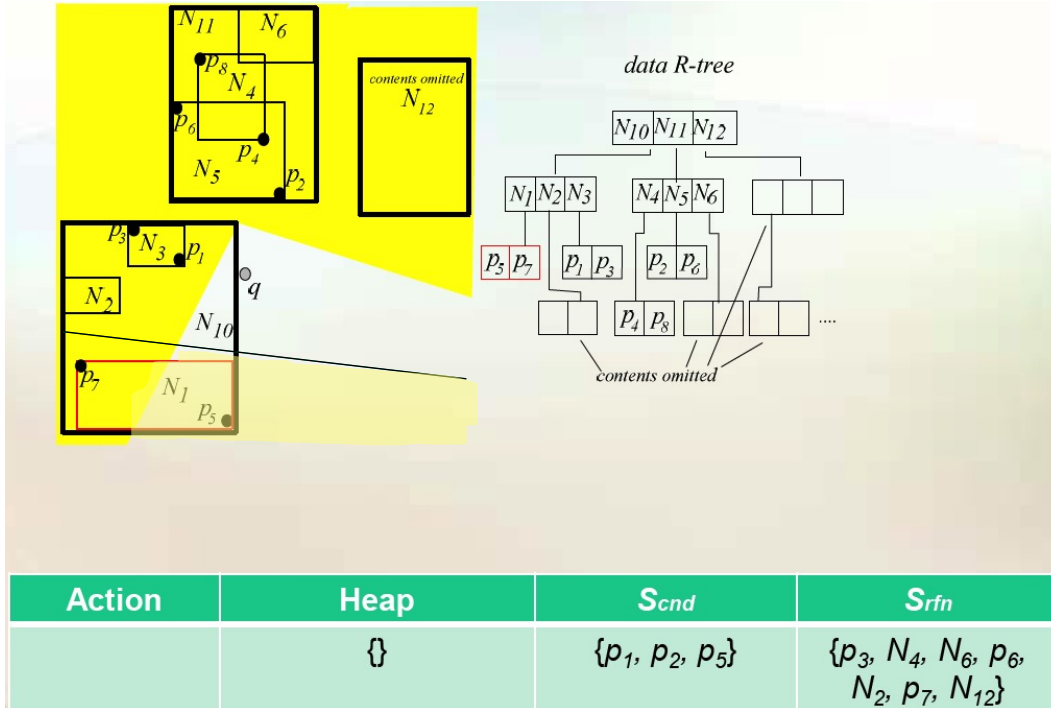
15

Flitering step - 6



16

Flitering step - 7



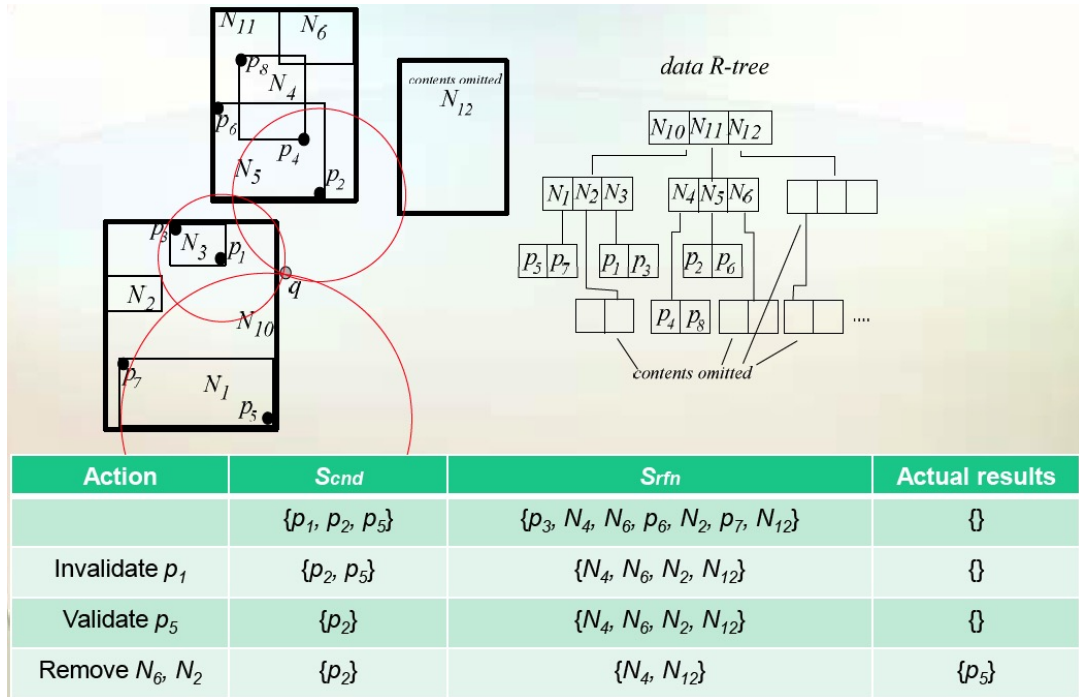
17

Refinement Heuristics

- Let P_{rfn} be the set of points and N_{rfn} be the set of nodes in S_{rfn}
- A point p from S_{cnd} can be discarded as a false hit if there is a point $p' \in P_{rfn}$ such that either of the following hold:
 - (i) $dist(p, p') < dist(p, q)$
 - (ii) There is a node MBR $N \in N_{rfn}$ such that $minmaxdist(p, N) < dist(p, q)$
- A candidate point can be eliminated if it is closer to another candidate point than to the query
- A point p from S_{cnd} can be reported as an actual result if the following two conditions hold:
 - (i) There is no point $p' \in P_{rfn}$ such that $dist(p, p') < dist(p, q)$
 - (ii) For every node $N \in N_{rfn}$: $mindist(p, N) \geq dist(p, q)$
- If none of the above works, visit all node MBRs $N \in N_{rfn}$ where $mindist(p, N) < dist(p, q)$ and use the mentioned heuristics considering the newly visited entries

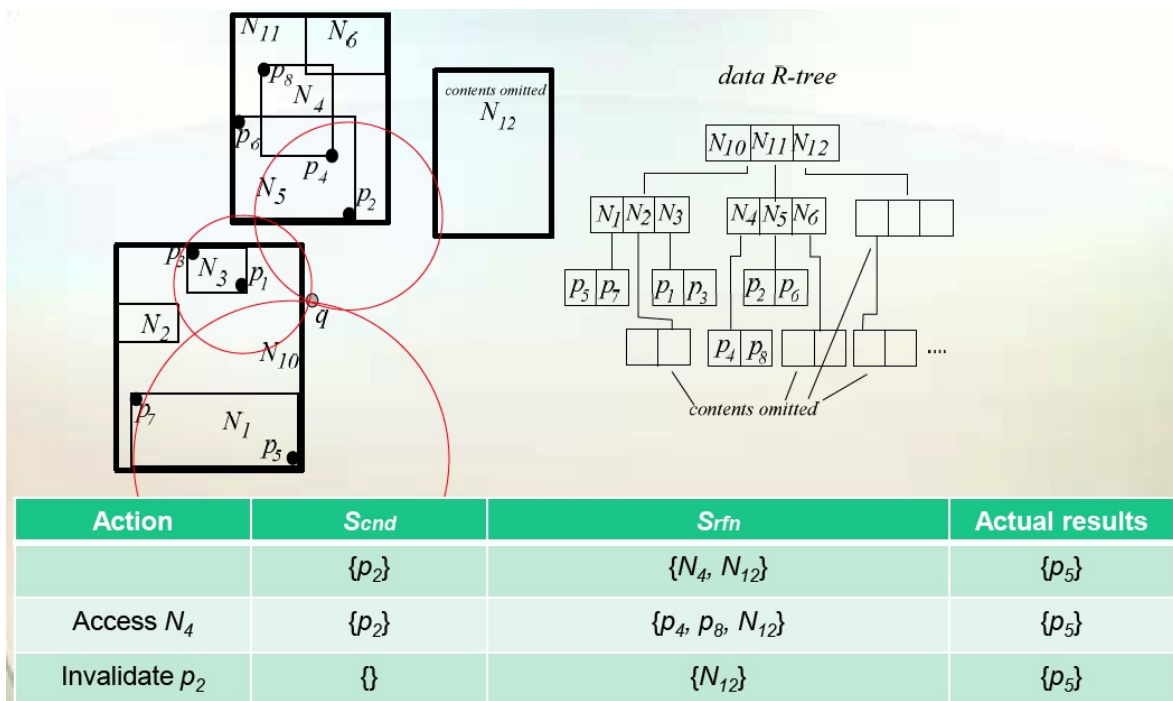
18

Refinement step - 1



19

Refinement step - 2



20