# A Scalable Algorithm for Maximizing Range Sum in Spatial Databases
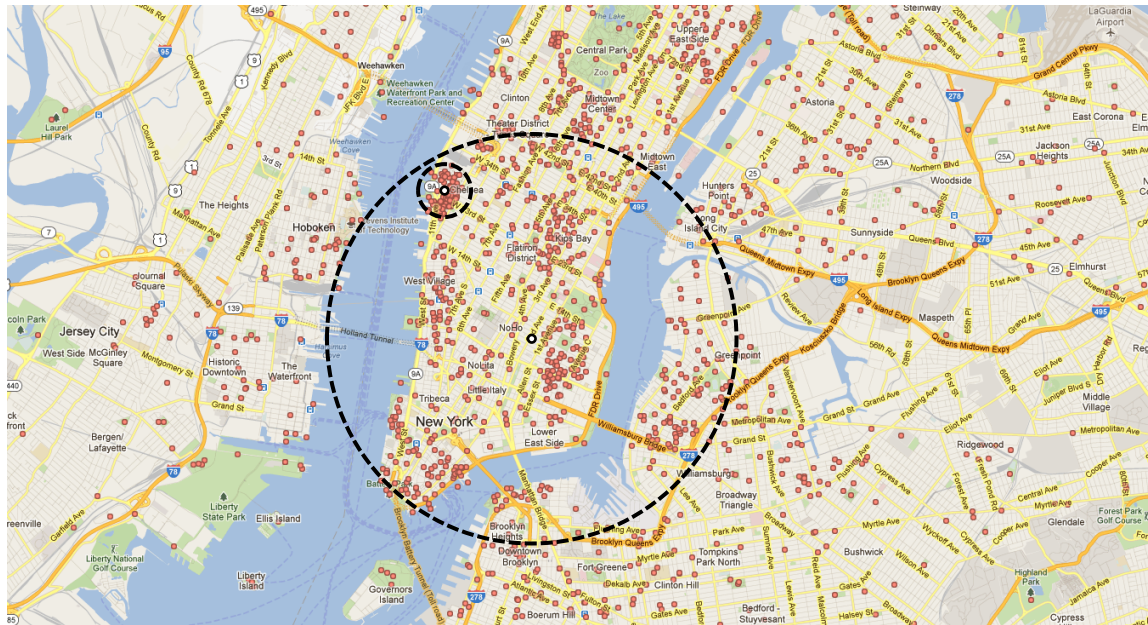
Dong-Wan Choi, KAIST, Korea

Chin-Wan Chung, KAIST, Korea

Yufei Tao, KAIST & CUHK, Hong Kong
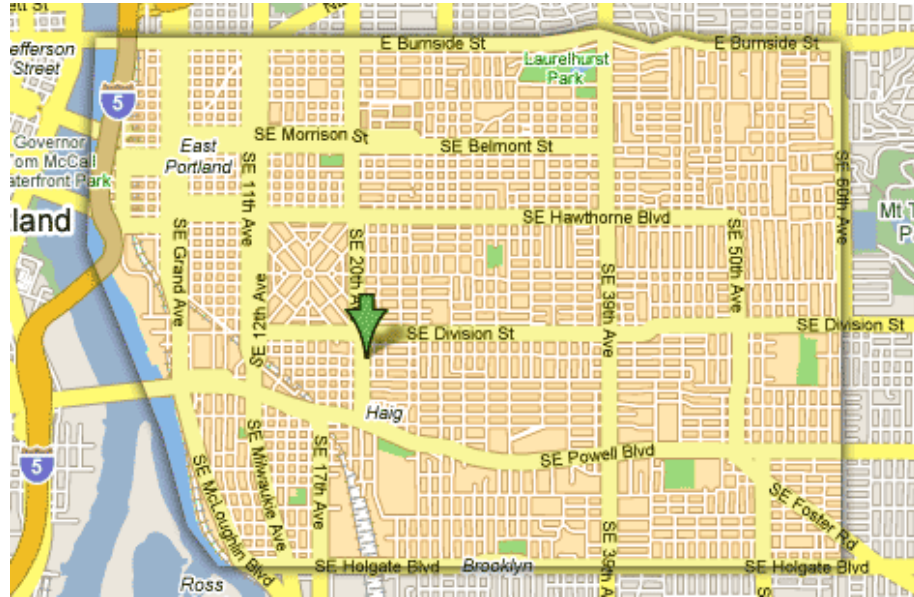
# MOTIVATION

‣ Where is <span style="color:red">the most representative spot</span> in a city for a tourist?



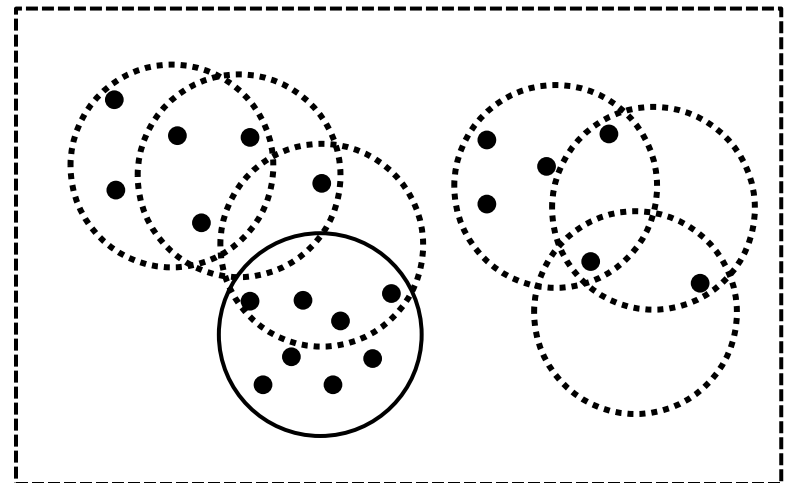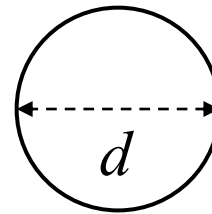‣ Maximize the number of tourist attractions within a movable range

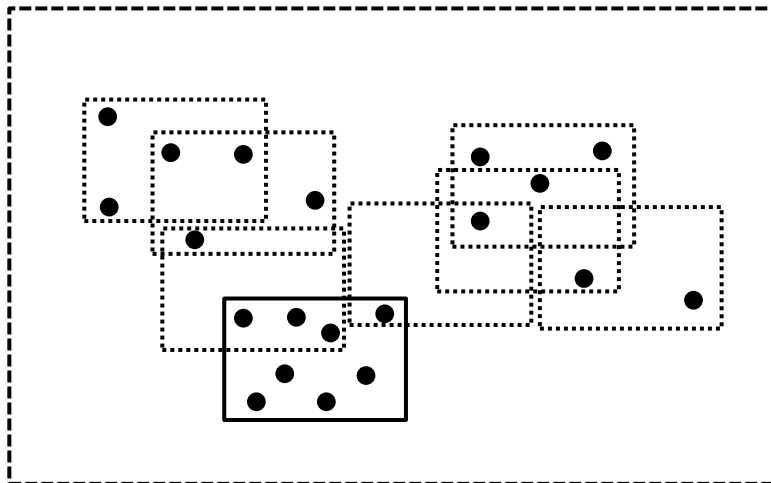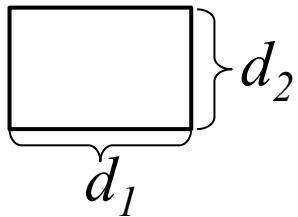# MOTIVATION



▸ Where is the most profitable place to set up a new pizza store?

  ▸ Maximize the number of residents within a delivery range

# MOTIVATION

▸ Where is the best location that maximizes the number of objects covered by a given range?

# FORMAL DEFINITION

▸ MaxRS (Maximizing Range Sum) problem
  ▸ Given a set $O$ of weighted objects and a rectangle $r$ of a given size,
  ▸ Find a location $p$ of $r$ that maximizes:
$$\sum_{o \in O_{r(p)}} w(o), \text{ where}$$

$r(p)$ is the rectangle centered at a location $p$, $O_{r(p)}$ is the set of objects covered by $r(p)$, and $w(o)$ is the weight of $o \in O$

▸ MaxCRS (Maximizing Circular Range Sum) problem
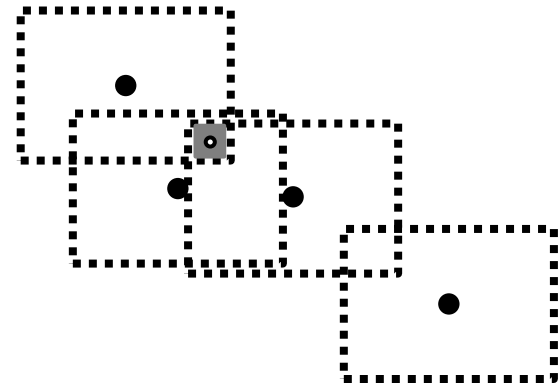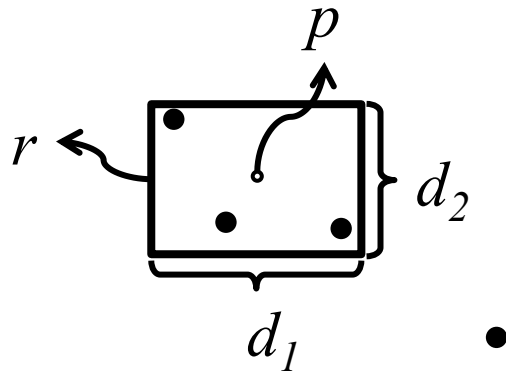  ▸ The circle version of the MaxRS problem

# CONTENTS

❖ ~~Motivation~~

❖ ~~Formal Definition~~

❖ Preliminary

❖ Our Algorithms

  ❖ ExactMaxRS

  ❖ ApproxMaxCRS

❖ Theoretical Results

❖ Experimental Results

❖ Conclusion and Future works

# PRELIMINARY

▶ A naïve solution
  ▶ Issuing range aggregate queries for every location
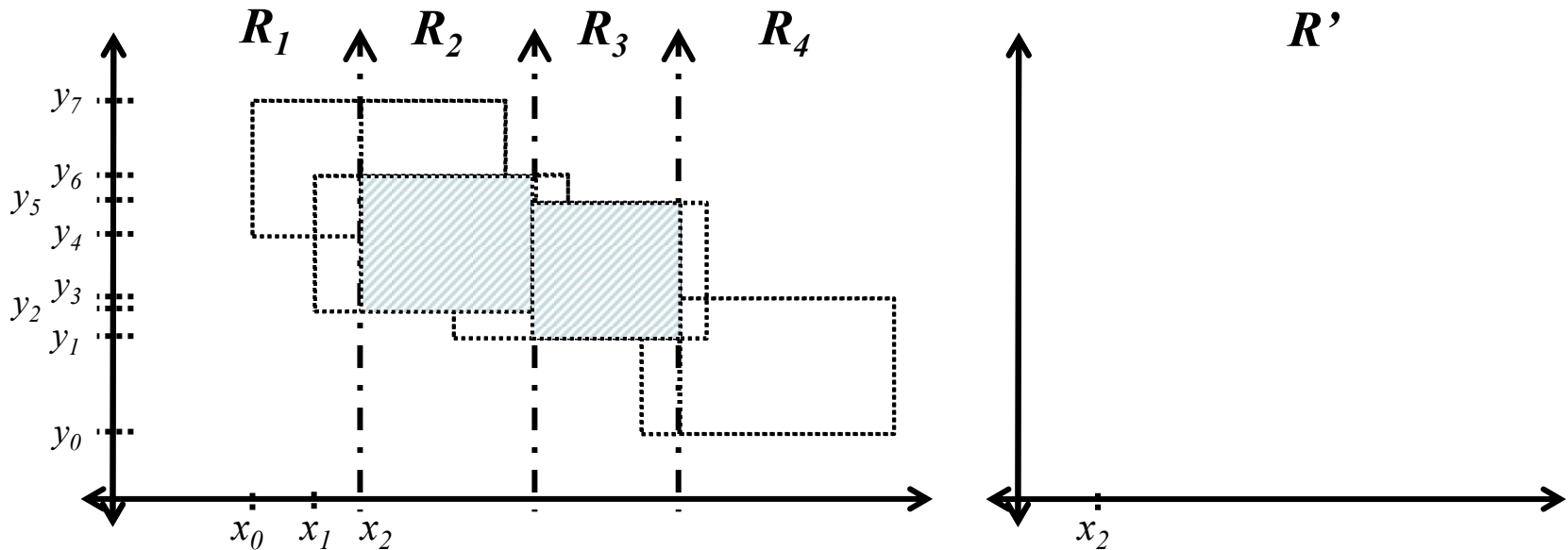  ▶ Problem: Infinite # of locations!

▶ Problem transformation



  ▶ Given a set of rectangles, find the most dense region where the most rectangles intersect

# ExactMaxRS

▸ Exact algorithm that solves MaxRS

▸ External-memory algorithm
  ▸ Scalable for a massive dataset

▸ Follows the divide-and-conquer strategy:
  ▸ Recursively divide the entire dataset into smaller subsets
  ▸ Compute a local solution for each subset
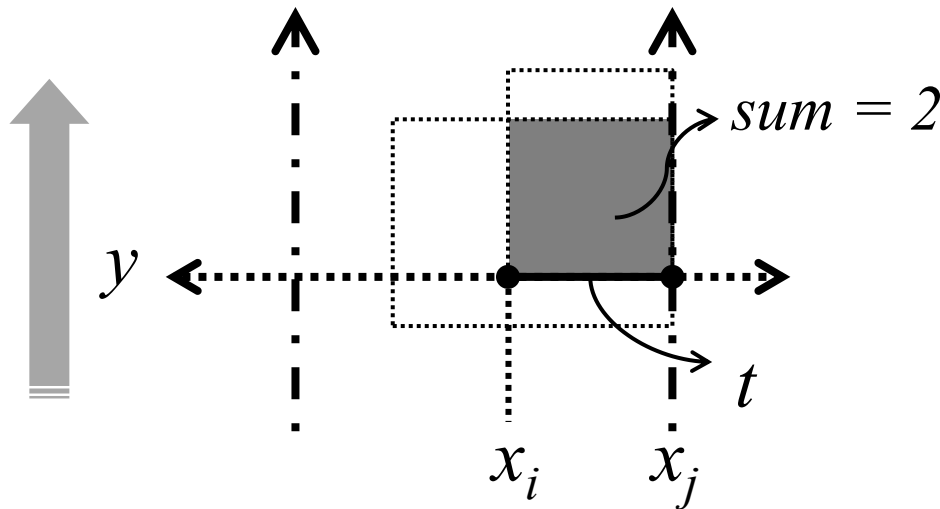  ▸ Merge local solutions of subsets

# DIVISION PHASE



- Divide the space vertically into $m$ sub-spaces, called *slab*s, <u>each of which has roughly the same # of rectangles</u>
  - Until the # of rectangles can fit in the main memory
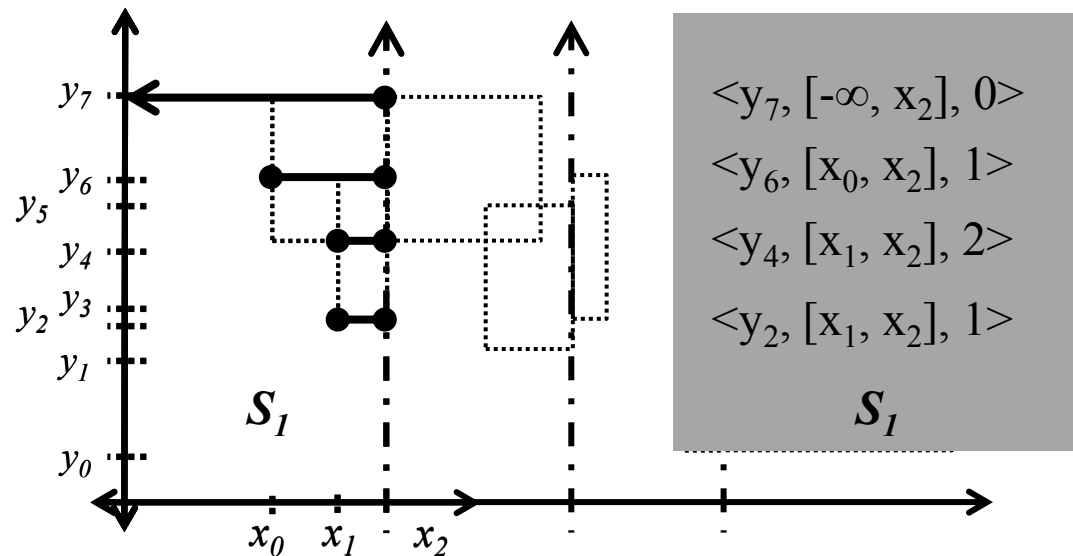- Do not pass spanning rectangles to the next level of recursion

# SLAB-FILES

▶ The structure to be returned <span style="color:red">after conquering the sub-problem</span> w.r.t. a slab
  ▶ The set of tuples, each of which is $t = <y, [x_i, x_j], sum>$
    ▶ <span style="color:red">In the upward direction, after $y$, the most dense region (whose total weights is $sum$) is in the x-range $[x_i, x_j]$.</span>
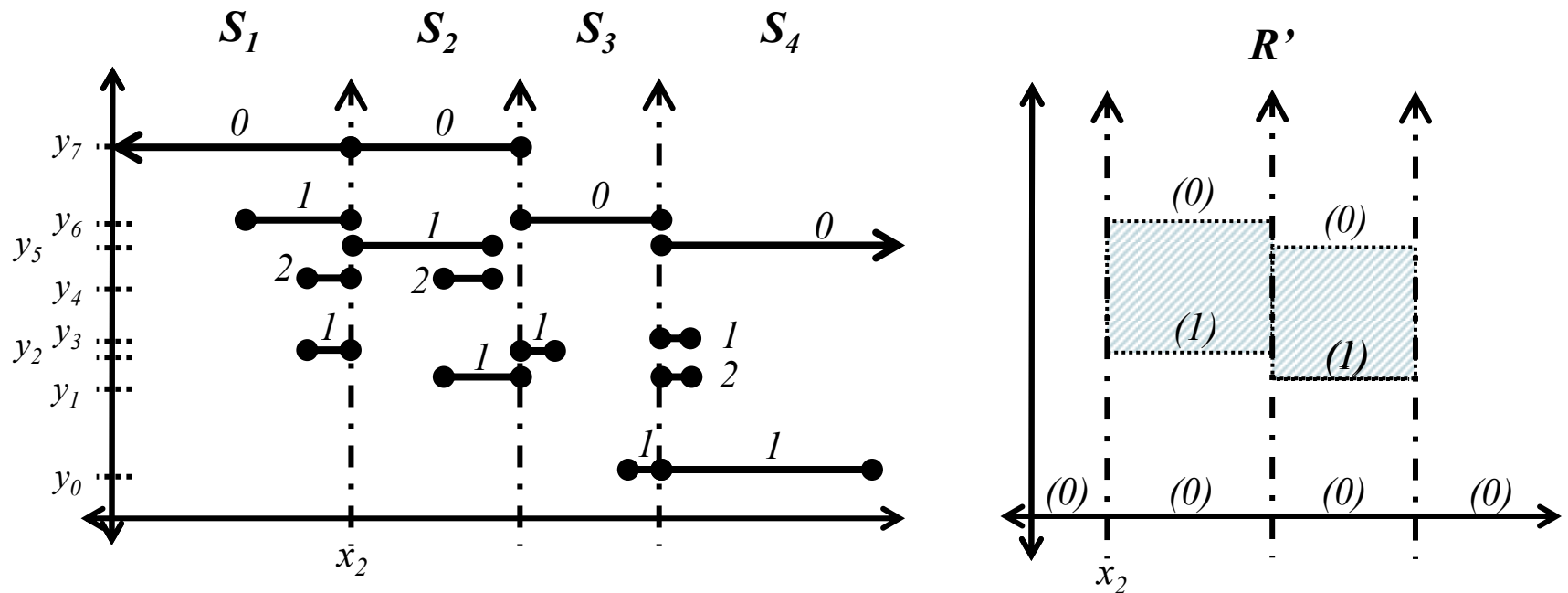


$sum = 2$

$y$

$t$

$x_i$    $x_j$

# SLAB-FILES

▸ An example of a slab-file



$<y_7, [-\infty, x_2], 0>$

$<y_6, [x_0, x_2], 1>$

$<y_4, [x_1, x_2], 2>$
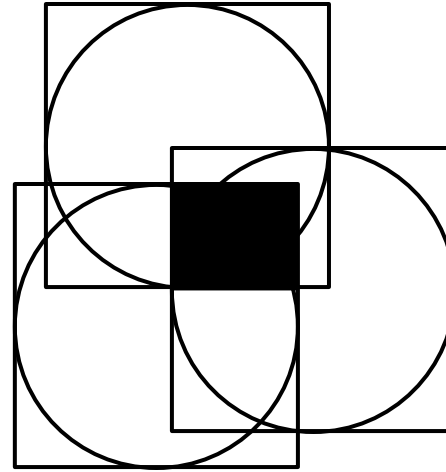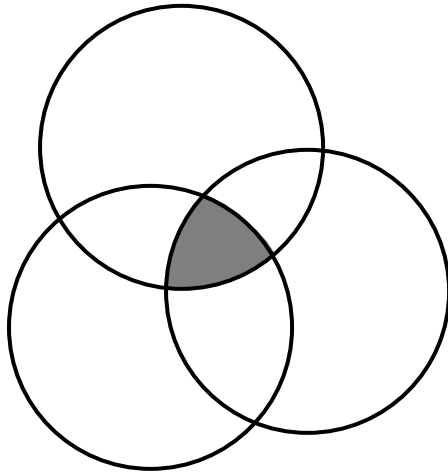
$<y_2, [x_1, x_2], 1>$

$S_1$

# MERGING PHASE



▸ Sweep a horizontal line across the slab-files ($S_1, ..., S_4$) and the spanning rectangle file ($R'$)

▸ When encountering several tuples at a horizontal line, choose a tuple with a maximum sum

# APPROXMAXCRS

▸ Approximation algorithm for MaxCRS

▸ Uses the ExactMaxRS algorithm as a tool

▸ Overall Flow

  1. Transform MaxCRS into MaxRS with MBRs

  2. Do ExactMaxRS on the transformed dataset

  3. Generate candidate points based on the result from ExactMaxRS

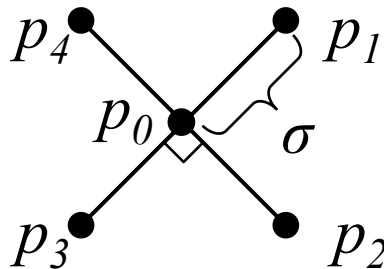  4. Choose the best point among the candidate points

# TRANSFORMATION

▸ MaxCRS ➔ MaxRS with MBRs
   ▸ Construct the MBR for each circle



▸ Find the most dense region of MaxRS with MBRs

# CANDIDATE POINTS

‣ Generate five candidate points based on the result from ExactMaxRS

- ‣ $p_0$ : the center point of the most dense region returned from ExactMaxRS

- ‣ $p_1, p_2, p_3, p_4$ : four shifted points from $p_0$



$$\sigma \in \left( (\sqrt{2} - 1)\frac{d}{2}, \frac{d}{2} \right)$$

, where $d$ is the diameter of circles

‣ Return the best point $p_i$ among $p_0, ..., p_4$

- ‣ such that the total weight of the circles covering $p_i$ is maximized.

# THEORETICAL RESULTS

▸ ExactMaxRS

  ▸ Optimal in terms of the I/O complexity

    ▸ $O((N/B)\log_{M/B}(N/B))$ I/O's, where N is the # of objects, M is the memory size, and B is the block size

      ☐ The counterpart of <u>O(nlogn)</u> in the main memory
      
      → optimal time complexity
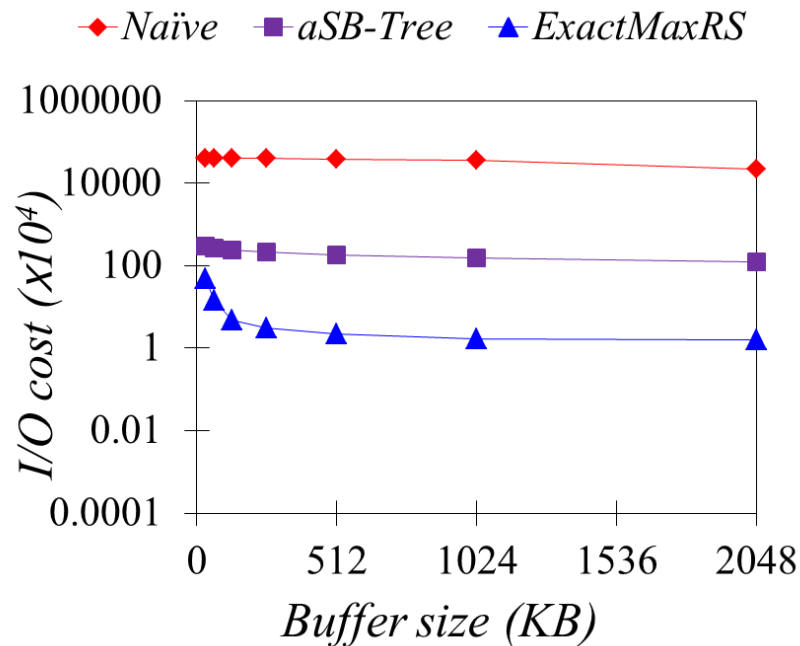
▸ ApproxMaxCRS

  ▸ (1/4)-approximation bound
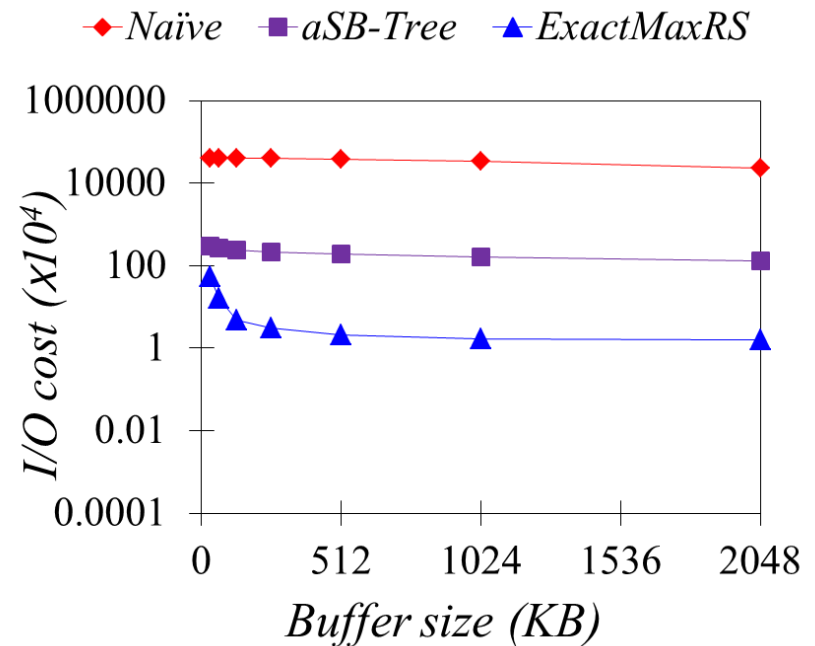
    ▸ Much better in practice

# EXPERIMENTAL RESULTS

▸ Environment Setting
  ▸ Synthetic datasets
    ▸ Uniform distribution and Gaussian distribution
    ▸ Cardinality: 100,000~500,000
  ▸ Real datasets
    ▸ NE and UX datasets from the R-tree Portal site
    ▸ Cardinality: 123,593 for NE, 19,499 for UX
  ▸ Compared approaches
    ▸ Naïve plane sweep, $O(n^2)$
    ▸ Plane sweep using aSB-tree, $O(n\log n)$
      ☐ <u>Adaptation of the optimal in-memory algorithm</u> for the MaxRS problem

# EXPERIMENTAL RESULTS

▸ I/O cost with varying the buffer size
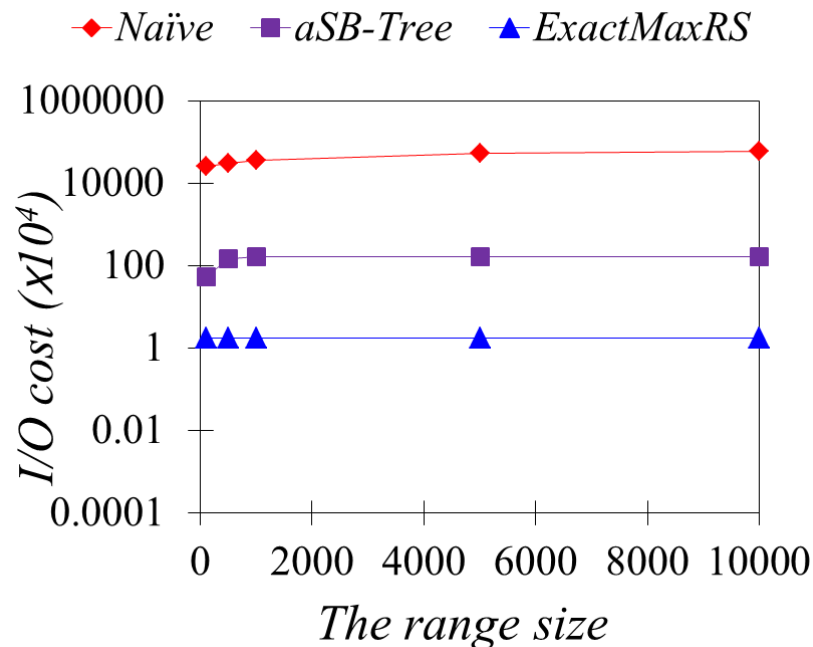


<Gaussian distribution>

<Uniform distribution>

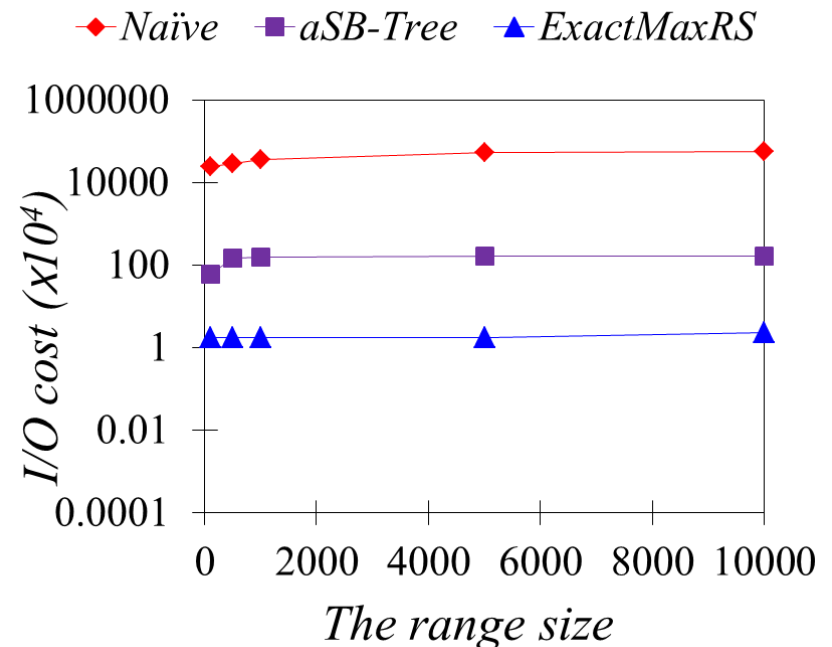# EXPERIMENTAL RESULTS

▸ I/O cost with varying the range size
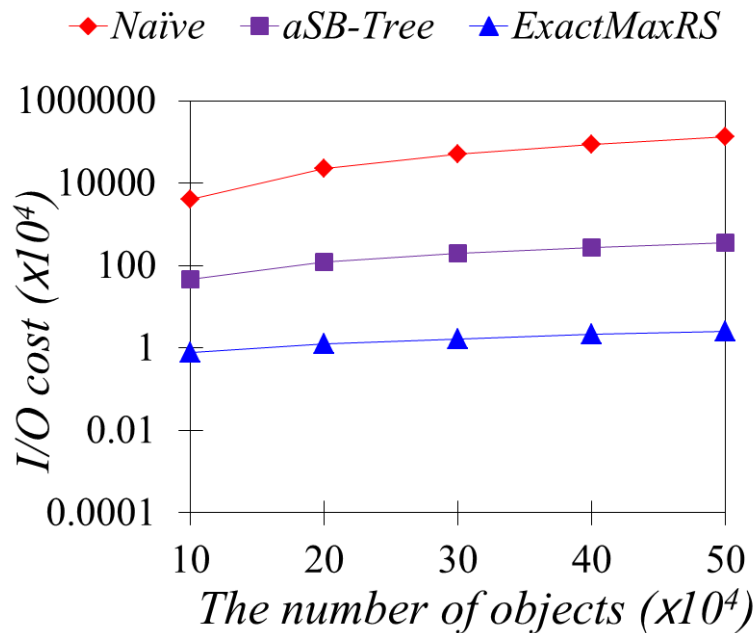


<Gaussian distribution>        <Uniform distribution>

# Experimental Results

‣ I/O cost with varying the cardinality



&lt;Gaussian distribution&gt;



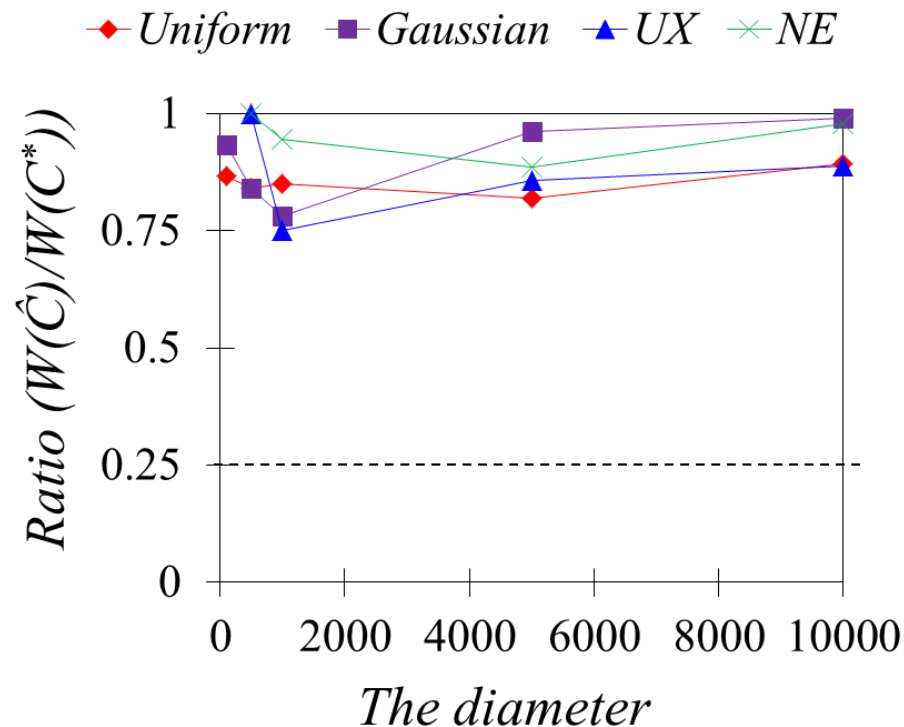&lt;Uniform distribution&gt;

# Experimental Results

▶ Error ratio of ApproxMaxCRS with varying the diameter

# CONCLUSION & FUTURE WORKS

▸ Contributions
  ▸ Proposed <span style="color:red">the first optimal external-memory algorithm</span> for the MaxRS problem
  ▸ Proposed <span style="color:red">the (1/4)-approximation algorithm</span> for the MaxCRS problem
  ▸ Proved theoretically the <span style="color:red">correctness</span>, <span style="color:red">optimality</span>, <span style="color:red">approximation bound</span>, and <span style="color:red">tightness</span> of the bound
  ▸ Evaluated experimentally the efficiency and accuracy
▸ Future works
  ▸ Max*k*RS problem
  ▸ MinRS problem

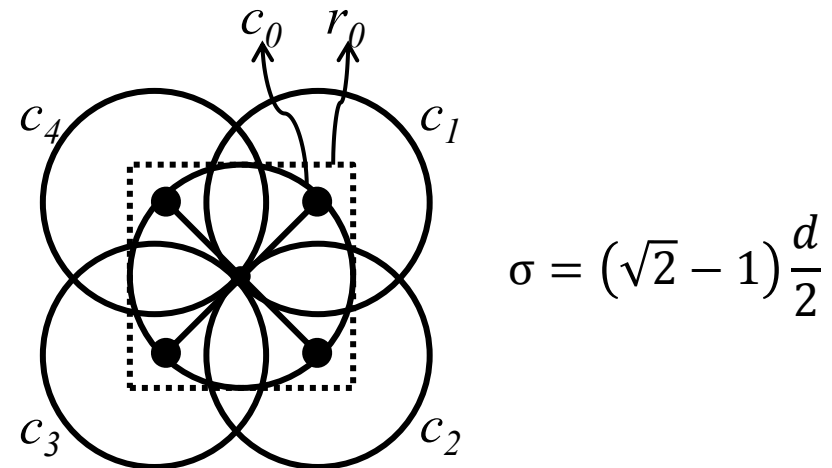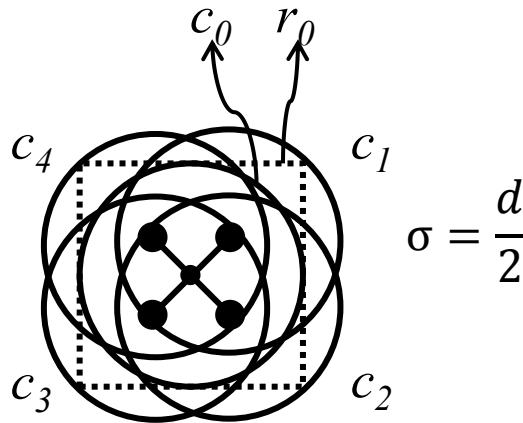# Thank You

# APPENDIX

- The proof sketch of the correctness of ExactMaxRS
  - The x-range of a tuple in a slab-file is called "max-interval".
  - Let $I^*$ be the max-interval w.r.t. entire space and $I^\wedge$ be a piece of I for a recursion level. Then $I^\wedge$ is also the max-interval w.r.t. the sub-space corresponding to the recursion level.
    - There cannot be an interval $I'$ whose sum is larger than $I^\wedge$, since spanning rectangles can affect all or none of intervals in a slab.
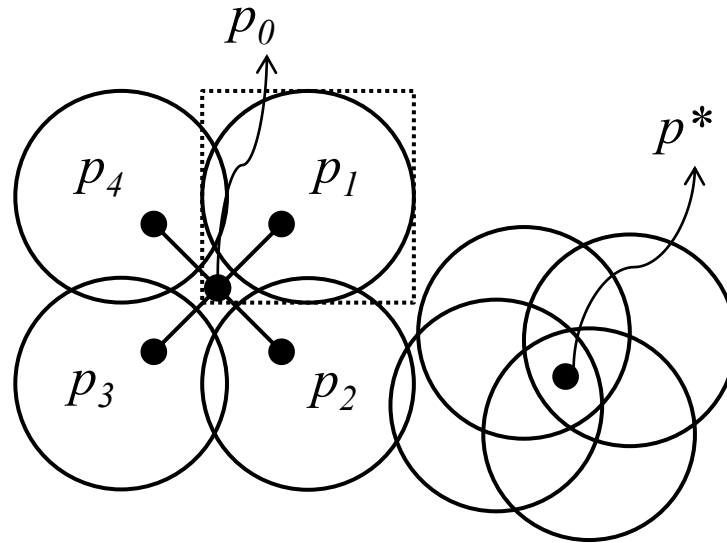
# APPENDIX

‣ The proof sketch of (1/4)-approximation bound

  ‣ The four circles together cover the region for MaxRS, which covers *k* points. Hence, at least one of those circles covers *k/4* points.



$$\sigma = \frac{d}{2}$$

$$\sigma = (\sqrt{2} - 1)\frac{d}{2}$$

# APPENDIX

▸ The proof sketch of the tightness of the approximation bound

  ▸ $4p_i \leq p^*$

# APPENDIX

▸ Results on real datasets



<UX>

<NE>

<UX>

<NE>