# Spatial Online Sampling and Aggregation

Lu Wang<sup>1</sup>, Robert Christensen<sup>2</sup>, Feifei Li<sup>2</sup>, Ke Yi<sup>1</sup>

<sup>1</sup>Hong Kong University of Science and Technology {luwang, yike}@cse.ust.hk

<sup>2</sup>University of Utah {robertc, lifeifei}@cs.utah.edu



# Motivation

- ► Geo-Spatial data is being collected on a massive scale.
- Approximate analysis is *fast* and often effective for this data.



#### Notation

**P** The raw data set in  $\mathbb{R}^d$ . k The number of samples to report. N|P|, the size of the raw data set. Q A range query in  $\mathbb{R}^d$ .  $P_Q | P \cap Q$ , elements in the query range.  $q \mid P_Q$ , the number of elements in the query range.  $u, v, \cdots$  Tree nodes. T(u) The subtree rooted at node u. P(u) The set of all data points covered by T(u).



### Hybrid R-Tree



- Top levels are *memory only nodes* which are persisted in main memory for fast access.
- Lower levels are disk nodes which are persisted on disk and loaded into main memory when needed.

#### **Baseline Algorithms**

Query First — Calculate  $S = P \cap Q$ . Repeatably extract a sample from S upon request.

r(N) The size of a canonical set in a R-tree of size N. *B* The size of a disk block. s The sample buffer size in RS-tree.

## **Comparison of various sampling algorithms**

Algorithm	Query Time	Update Time
QueryFirst	r(N) + q	log N
SampleFirst	kN/q	log N
RandomPath	k log N	log N
RandomShuffle	kN/q	log N
LS-tree	$\sum_{i=\log q/k}^{\ell} r\left(\frac{N}{2^{j}}\right) + k$	log N
RS-tree	r(kN/q) + k	log N

 $\blacktriangleright$  Full analysis of algorithms, including I/O cost can be found in the paper.

#### **Building Index Experiments**



Sample First — Upon Request, pick a point p where  $p \in P$ . If  $p \in Q$  report *p*, otherwise repeat.

#### **Data Structures**

R-tree — A Hybrid R-tree implementation based off of Hilbert R-tree.

RS-tree — A single R-tree T with a sample buffer s attached to each internal node. For each internal node  $u \in T$  with sample buffer s, we sample uniformly from the children of u,  $p(x = y | x \in s \land p \in P(u)) = \frac{1}{|P(u)|}$ 

LS-tree — a collection of R-trees where each R-tree indexes a set of samples from the original data set. The sample rates for these sets of samples form a geometric series.

#### Individual Nodes of RS-tree

#### from parent node

### Sample Query Experiments



**RS**-tree

LS-tree

15

10

(b)

 $q (\times 10^6)$ 

20

25



children within MBB are added to *frontier* 



Figure: vary q, the number of samples in the query region.