

Spatial Hadoop

A MapReduce Framework for Spatial Data

http://spatialhadoop.cs.umn.edu/

Mohamed F. Mokbel Ahmed Eldawy University of Minnesota shadoop@cs.umn.edu

Big Spatial Data





Elegant and Powerful

 $x < x^2$ AND $x > x^1$ AND Takes 193 y < y2 AND y > y1;seconds

Results



Storage Layer (Indexing)

Hadoop is designed to work with heap non-indexed files Traditional spatial indexes are designed for:

- Procedural programming while Hadoop uses MapReduce programming

- Traditional file system while Hadoop uses Hadoop Distributed File System (HDFS)

- A two layered approach (Global and local indexes) is used to build grid file, R-tree and R+-tree indexes in SpatialHadoop



Index Building





Grid Partitioning



R-tree Partitioning



MapReduce Layer

MapReduce programs in Hadoop can only deal with non-indexed files SpatialHadoop adds two new components to allow MapReduce programs to utilize spatial indexes

- 1. SpatialFileSplitter: Utilizes the global index by pruning file blocks that do not contribute to answer
- 2. SpatialRecordReader: Utilizes local indexes by efficiently selecting

Operations Layer

Range query

SpatialFileSplitter prunes blocks outside query range SpatialRecordReader passes local indexes to the map function Map function selects records in



Spatial Join

SpatialFileSplitter finds overlapping block pairs





records that need to be processed

Map phase in SpatialHadoop

Indexed Input







Map function joins each pair of overlapping blocks



and result is tested for correctness **×** Initial answer is incorrect **Second iteration processes** other blocks that might

contain an answer ✓ Final answer is correct

First iteration runs as before

kNN

SpatialFileSplitter selects the block that contains the query point Map function performs kNN in the selected block Answer is tested for correctness ✓ Initial answer is correct





University of Minnesota

This work is supported in part by the National Science Foundation under Grants IIS-0952977 and IIS-1218168