

ABSTRACT

The recent evolution of ubiquitous positioning technologies has led to the rapid growth of the usage of spatial data by various systems, such as geographic information systems (GIS) and location based services (LBS). The spatial queries of these systems are complex and computationally intensive, this caused the need of handling them using high performance platforms.

MapReduce and its open source implementation Hadoop, have been extensively used by large enterprises in processing large-scale data sets in parallel and therefore, they are suitable for performing the intensive computations of spatial data. However, Hadoop is not aware with the spatial properties of spatial data, which is why computation benefits of considering the spatial properties of the data when processing, may be wasted by this limitation. This introduces the need of extending Hadoop to be aware of spatial features of the data in order to obtain better performance.

Several research work have proposed approaches to extend Hadoop to support spatial data. For example, SpatialHadoop extends Hadoop to add support to spatial data and spatial queries. This is achieved through the following: (1) using spatial indexes to efficiently access the spatial data stored on the Hadoop distributed file system (HDFS) and (2) adding an operation layer that supports spatial operation.

This thesis introduces Co-SpatialHadoop; it is a Hadoop-based spatial approach and an extension for SpatialHadoop platform. It enhances the network usage of spatial complex queries such as “spatial join” by colocating related spatial files together on HDFS, which may be joined later. It enhances the saving operation of HDFS by introducing a new spatial placement policy which colocates data using spatial attributes of spatial records.

Co-SpatialHadoop also builds inverted indexes using non-spatial attributes of data records to enhance the response time of non-spatial queries. It adds non-spatial operations to operation layer of SpatialHadoop. These modules are implemented using the functionality of MapReduce.

Co-SpatialHadoop uses real spatial data obtained from OpenStreetMap (OSM) in experiments. The experiments demonstrate good results; colocation algorithm enhances network overhead of spatial join queries by 60%. Also, inverted indexes enhance response time and network overhead of non-spatial queries.

List of Figures

Figure 1.1: Technologies and services that use spatial data	2
Figure 2.1: GIS system	7
Figure 2.2: OSM cities layer for USA	8
Figure 2.3: Spatial data model	9
Figure 2.4: Hadoop Ecosystem [13]	12
Figure 2.5: Hadoop physical architecture	13
Figure 2.6: A client reading data from HDFS.....	14
Figure 2.7: A client writing data to HDFS.....	15
Figure 2.8: Spatial file Partitioning.....	20
Figure 2.9: Architecture of [24] approach	22
Figure 2.10 : SpatialHadoop [1] Architecture.....	24
Figure 3.1: Co-Spatial Hadoop Architecture	30
Figure 3.2: Indexed Files "A" and "B"	31
Figure 3.3: Spatial block distribution among nodes without colocation.....	32
Figure 3.4: Spatial block distribution among nodes with colocation.....	33
Figure 3.5: An example demonstrating the overlapping between data blocks and locator items.....	37
Figure 3.6: Locator Table is partitioned to five partitions	40
Figure 3.7: Locator Table is partitioned to six partitions.....	41
Figure 3.8: Use vertical partitions only	41
Figure 3.9: Physical architecture for new placement of a new block	44
Figure 4.1: Some records of Parks file	47
Figure 5.1: The Ganglia monitoring system	51
Figure 5.2: Network overhead measured by Ganglia for a query	52
Figure 5.3: Network overhead chart	54
Figure 5.4: Response time chart of the experiment	56
Figure 5.5: Network overhead comparison using two different approaches	57
Figure 5.6: Network overhead comparison using different number of locator items	58
Figure 5.7: Comparison between SpatialHadoop and Co-SpatialHadoop indexing time.....	59
Figure 5.8: Response time comparison with/without inverted index.....	60
Figure 5.9: Network overhead comparison with/without inverted index	60

List of Tables

Table 2.1: Comparison between Parallel SDBMS & Hadoop	11
Table 2.2: Spatial approaches challenges	25
Table 3.1: use case Locator Table	35
Table 5.1: Experiments Queries.....	53
Table 5.2: some statistics for the experiment.....	55
Table 5.3: blocks distribution over Locator Table items	55

List of Algorithms

Algorithm 1: default placement for new spatial block.....	38
Algorithm 2: Build Locator Table	39
Algorithm 3: Assign locator item to new block and choose suitable DataNodes.....	42
Algorithm 4: Spatial placement for new spatial block.....	43

LIST OF ACRONYMS

GFS	Google File System
GIS	Geographic Information System
GPS	Global Positioning System
HDFS	Hadoop Distributed File System
KNN	K Nearest Neighbors
LBS	Location Based Services
MBR	Minimum Bounding Rectangle
ORDBMS	Object Relational Database Management Systems
SDBMS	Spatial Database Management Systems
WKB	Well Known Binary Language
WKT	Well Known Text Language