



Scalable 3D Spatial Queries for Analytical Pathology Imaging with MapReduce

Yanhui Liang¹, Hoang Vo¹, Ablimit Aji², Jun Kong³, Fusheng Wang¹

¹Department of Computer Science, Stony Brook University, ²Hewlett Packard Labs, ³Department of Biomedical Informatics, Emory University

Introduction

Background

□ Massive amount of 3D spatial objects derived from 3D analytical pathology imaging GIS" generic queries for pathologically meaningful spatial analysis on 3D big data e.g., for each 3D cell, return the nearest 3D blood vessel and the distance

Challenges

• Explosion of 3D data: hundreds of millions of 3D objects in the scale of terabytes • Complex structures and representations: 3D mesh model with numerous structural details □ High computation complexity: computationally intensive 3D geometric operations



Objectives

To provide a scalable and efficient 3D spatial query system for querying massive 3D spatial data based on MapReduce



Scalable 3D Spatial Queries with MapReduce

System Overview

□ Spatial data partitioning: partition 3D input into cuboids to increase the level of parallelism □ 3D on-demand spatial query engine: multi-level indexing and spatial query processors

3D Spatial Partitioning

• Each partitioned cuboids as a processing unit for querying tasks on MapReduce • Cuboids are processed in parallel without data dependence or communication requirements

Multi-level Indexing

Global storage indexing: HDFS-level data filtering (point query) Cuboid indexing: MapReduce-level computation filtering (containment query) • On-demand object indexing: to build indexes on-the-fly for objects within a cuboid □ Structural indexing: for geometric computation on 3D objects with complex structures

3D Spatial Queries



□ Spatial join or spatial cross-matching

□ Nearest neighbor: 3D R*-tree (with AABB tree) and Voronoi diagram (with skeleton)





AABB Tree Illustration with 2D Data



Hadoop/MapReduce

System Architecture of Hadoop-GIS 3D





Structure Index of 3D Blood Vessel with Skeleton

Experimental Results

Performance of 3D Spatial Queries

- □ The cluster has five nodes with 124 cores in total; Each node has 5TB hard drive and 128GB memory
- □ Five datasets with millions of 3D cells, and several thousand of 3D blood vessels
- □ Three 3D datasets for Level of Details (LoD) testing





Fig. 1: System Performance of Spatial Join



Acknowledgements

This research is supported in part by grants from National Science Foundation ACI 1443054 and IIS 1350885, National Institute of Health K25CA181503, the Emory University Research Committee, Pitney Bowes, Amazon and Google.













UNIVERSITY