

Automatic Parallelization of Graph Queries with MapReduce

Tao Zan¹ Yu Liu¹ Zhenjiang Hu²

¹The Graduate University for Advanced Studies ²National Institute of Informatics

Problem

◆ Graphs in real world can be quite huge with millions of nodes and edges.

Twitter Social Network, 20K nodes 250K edges

◆ How to parallelize graph queries efficiently is an important issue.

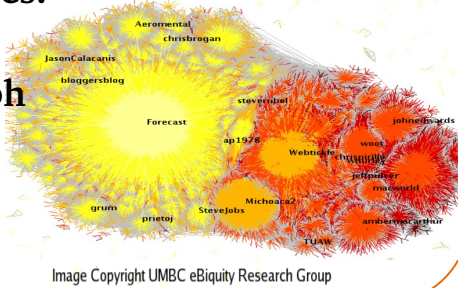


Image Copyright UMBC eBiquity Research Group

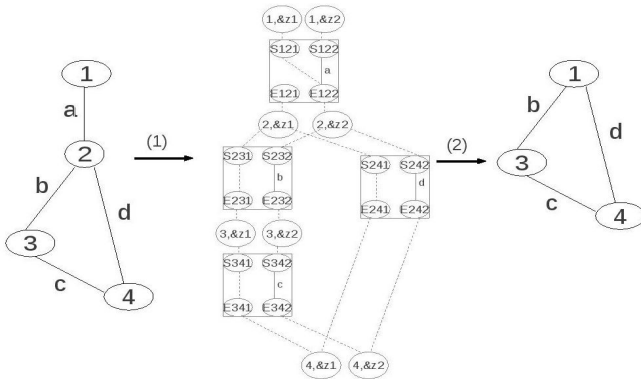
Goal

We aim at a new way of processing graph queries with MapReduce by translating UnQL query to structural recursion, then parallelizing structural recursion. When achieving this goal, thus can we use UnQL for large scale data processing by MapReduce.

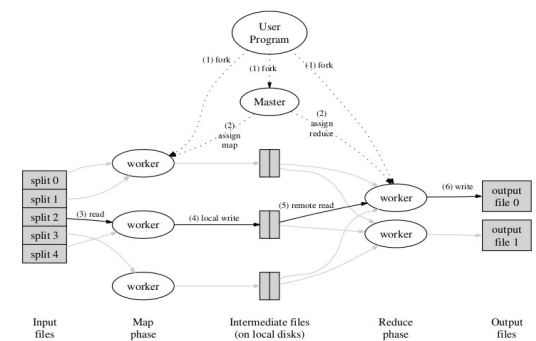
Background

◆ UnQL is a language that was designed for querying unstructured data.

An example:
Q1 = select t
where *.a.t in \$db



◆ MapReduce was introduced for processing huge datasets on clusters of computers.



Methods, Results and Future Work

◆ Methods

● Properties for Parallelization

All join-free queries can be highly parallelized in MapReduce based on three decomposable properties for structural recursion.

$$\begin{aligned} \text{rec}(e)(\{ \$l : \$g \}) &= e @ \text{rec}(e)(\$g) \\ \text{rec}(e)(\$g1 @ \$g2) &= \text{rec}(e)(\$g1) @ \text{rec}(e)(\$g2) \\ \text{rec}(e)(\text{cycle}(\$g)) &= \text{cycle}(\text{rec}(e)(\$g)) \end{aligned}$$

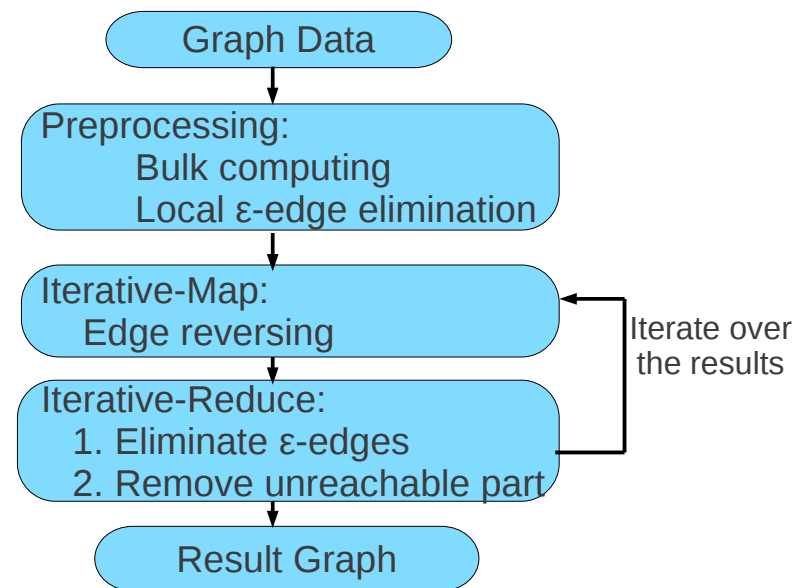
A graph can be expressed as:

$$g = \&x @ \text{cycle}(g1 \oplus g2 \oplus \dots \oplus gn)$$

By induction:

$$\begin{aligned} \text{rec}(e)(g) &= \text{rec}(e)(\&x) @ \text{rec}(e)(\text{cycle}(g1 \oplus g2 \oplus \dots \oplus gn)) \\ &= \text{rec}(e)(\&x) @ \text{cycle}(\text{rec}(e)(g1 \oplus g2 \oplus \dots \oplus gn)) \\ &= \text{rec}(e)(\&x) @ \text{cycle}(\text{rec}(e)(g1) \oplus \text{rec}(e)(g2) \oplus \dots \oplus \text{rec}(e)(gn)) \end{aligned}$$

● Parallel Computation with MapReduce



◆ Results

We evaluated in Hadoop cluster constructed with 8 virtual machines (VM). Each VM has 3 GB memory and one single-core 2.8 GHz CPU.

Q1 = select {a : t} where *.a.t in \$db

Q2 = select {a : {select b where *.b.t2 in \$t1}} where *.a.t1 in \$db

Q3 = select {a : {select {b : {select c where c.t3 in \$t2}} where b.t2 in \$t1}} where *.a.t1 in \$db

Graph Queries on Hadoop

Queries	1+3 million	10+30 million	20+60 million
Q ₁	104 sec.	545 sec.	1091 sec.
Q ₂	138 sec.	773 sec.	1429 sec.
Q ₃	251 sec.	1073 sec.	1683 sec.

(x+y)million : a graph with x million vertexes and y million edges

◆ Future work

- Parallelization of graph query on cyclic graphs.
- A combinational approach for parallelizing graph query UnQL.