

## Scalable Subgraph Enumeration in MapReduce



CSE, The University of New South Wales, Australia
QCIS, University of Technology, Sydney



L: Een Goed begin is het halve werk

- Problem Statement
- Methodology
- Evaluation



Given a pattern graph *P*, we aim at finding all matches of P in a large data graph *G* using MapReduce

Both P and G are **unlabeled**, **simple** graphs

*G* is arranged as (*u*; *N*(*u*)) for each data nodes, and stored in the HDFS







Given a pattern graph *P*, we aim at finding all matches of *P* in a large data graph *G* using MapReduce



G

Ρ





Given a pattern graph *P*, we aim at finding all matches of *P* in a large data graph *G* using MapReduce



G

Ρ

We can find three matches

 $(u_1, u_2, u_3, u_4)$ , the peripheral square





## **Subgraph Enumeration**

Given a pattern graph *P*, we aim at finding all matches of *P* in a large data graph *G* using MapReduce



We can find three matches

( $u_1$ ,  $u_2$ ,  $u_3$ ,  $u_4$ ), the peripheral square

 $(u_1, u_3, u_2, u_4)$ , the shadowed part





## **Subgraph Enumeration**

Given a pattern graph *P*, we aim at finding all matches of *P* in a large data graph *G* using MapReduce



We can find three matches

 $(u_1, u_2, u_3, u_4)$ , the peripheral square

 $(u_1, u_3, u_2, u_4)$ , the shadowed part

 $(u_1, u_2, u_4, u_3)$ , the white part





## Why Subgraph Enumeration



#### Social Networks

#### **Community Detection**

Nodes in the same communities tend to form specific subgraphs (clique, e.g.)

Network Similarity Study

The frequencies of specific pattern graphs can be good features to study network similarity

**Cluster Coefficient Calculation** 

# Triangles / (# neighbors X (#neighbors - 1))



## Why Subgraph Enumeration



#### **BioInformatics**

**Motif Discovery** 

Motifs are some very frequent patterns

Complex Network Design

Specific pattern graphs can facilitate complex network design

#### **Abnormity Detection**

Abnormity represents as unexpected statistical information of specific pattern graphs



## **Most Related Work**

#### Edge Join [T. Plantenga, JPDC 2012]

- Join edges one by one to get the final results
- Finish in m rounds, not scalable to handle complicated pattern graph
- Z Too many intermediate results

#### Multiway join [F.Afrati, D.Fotakis, J.D.Ullman, ICDE 2013]

- One-round computation by duplicating edges
- Do not scale well due to memory issues







L: Een Goed begin is het halve werk

- Problem Statement
- Methodology
- Evaluation



#### **Graph Decomposition**

Pattern Graph is decomposed into a sequence of **stars** 







## **Graph Decomposition**

Pattern Graph is decomposed into a sequence of **stars** 



#### Why Star?

When we have (u, N(u)) as input, we can search the matches of star rooted on u by enumerating the node combinations in N(u).



After decomposition, we will process (t-1)-round leftdeep join operations using MapReduce





After decomposition, we will process (t-1)-round leftdeep join operations using MapReduce





After decomposition, we will process (t-1)-round leftdeep join operations using MapReduce







## TwinTwig

Matching a star can still produce large number of intermediate results







Matching a star can still produce large number of intermediate results

Common for a node to contain many neighbors in a large social networks or web graphs

<u>A node with 100,000 neighbors  $=> 10^{15}$  matches of a 3-star</u>







Matching a star can still produce large number of intermediate results

Common for a node to contain many neighbors in a large social networks or web graphs



#### A star of **at most** two edges







#### **TwinTwig Decomposition**







#### **TwinTwig Decomposition**



<u>Star Join:</u> The join process w.r.t. star decomposition <u>TwinTwig Join:</u> The join process w.r.t. TwinTwig decomposition

























E









In order to estimate the cost, we need to compute the number of matches of any pattern graphs in the data graph <u>before running the algorithms</u>







G: Random Data graphM: Number of edges in Gm: Number of edges in P

P: Pattern graph

N: Number of nodes in G

*n*: Number of nodes in *P* 





#### **Instance Optimality**

Given a *star join*, *star*, there always exists a *TwinTwig join*, *tt*, s.t.

 $\Theta(cost(tt)) \le \Theta(cost(star))$ 





#### **Instance Optimality**

# We also show that the instance optimality holds in a power-law random graph model







L: Een Goed begin is het halve werk

- Problem Statement
- Methodology
- **Evaluation**



Progress: 75%



#### **Experimental Settings**

Algorithms	Desc.	Dataset G	V(G)	E(G)
Edge	Edge Join	sk	1.6M	11.1M
Mul	Multiway Join Star Join TwinTwig join	yt	3.2M	12.2M
Star TT		lj	4.8M	42.9M
		orkut	3.1M	261.8M
		uk	18.5M	261.8M
		fs	65.6M	1806.1M

Algorithms

Datasets





#### **Experimental Settings**



Queries





#### Vary Datasets



#### **Vary Queries**



#### **Vary Queries**





L: Een Goed begin is het halve werk

- Problem Statement
- Methodology
- Evaluation



#### Conclusions

A scalable MapReduce solution for subgraph enumeration

A left-deep-join framework to process star join

Leveraging the (power-law) random graph model, we show the instance optimality of TwinTwig join

Extensive performance studies to confirm the effectiveness of our approach









#### Thank you & Questions