

Apache Hama

Big Data 2nd Generation,
Big Compute and Big Insight!

2014 Samsung Open Source Conference

Edward J. Yoon @ *DataSayer*

 eddieyoon

 edwardyoon@apache.org

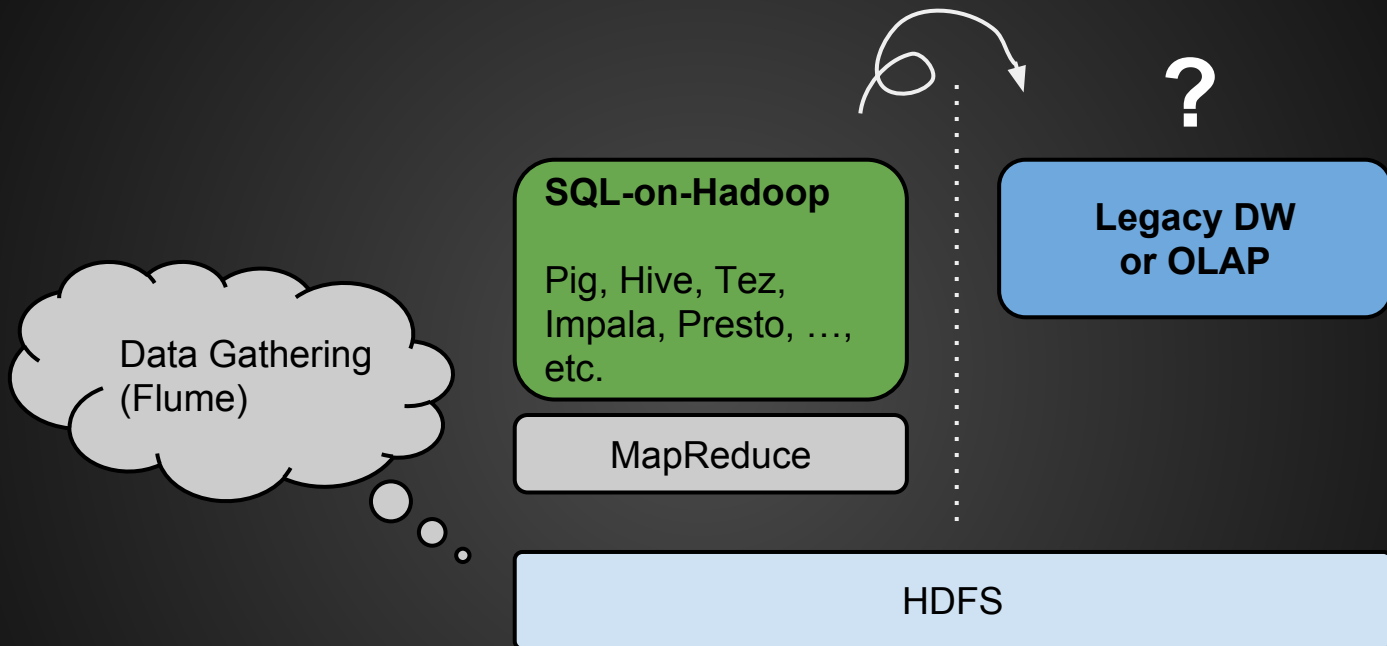
1. Big Compute Platform based on Apache Hama.
2. Big Data Crowdsourcing Service: datacrowds.com

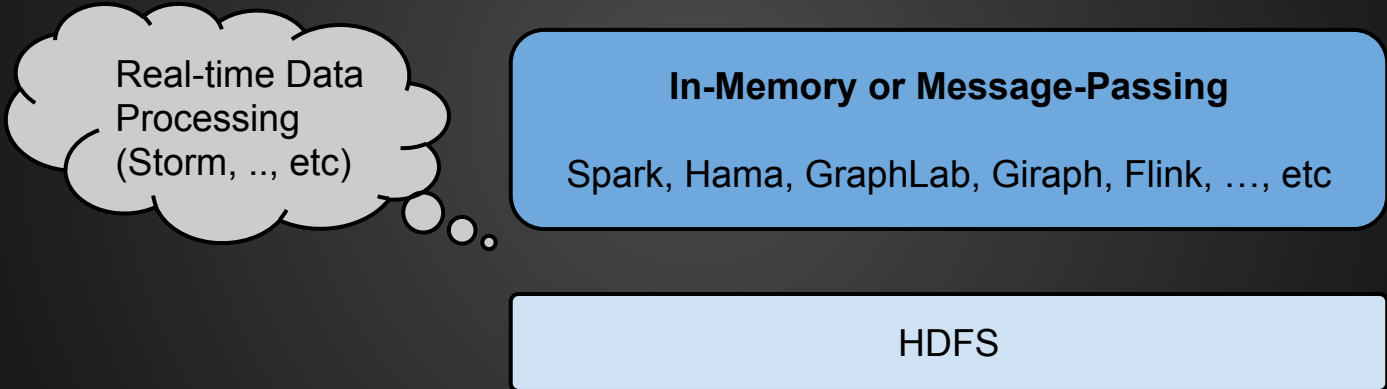
- 1. Big Data Trends**
- 2. What's Hama?**
- 3. Future Architecture of Hama**

Big Data Analytics

- Large-scale unstructured data processing
- Statistics and Data mining

To mine the valuable insights





Real-time Data
Processing
(Storm, ..., etc)

In-Memory or Message-Passing

Spark, Hama, GraphLab, Giraph, Flink, ..., etc

HDFS

2. In-Memory and Message-Passing

Spark,
Hama,
Giraph,
Storm

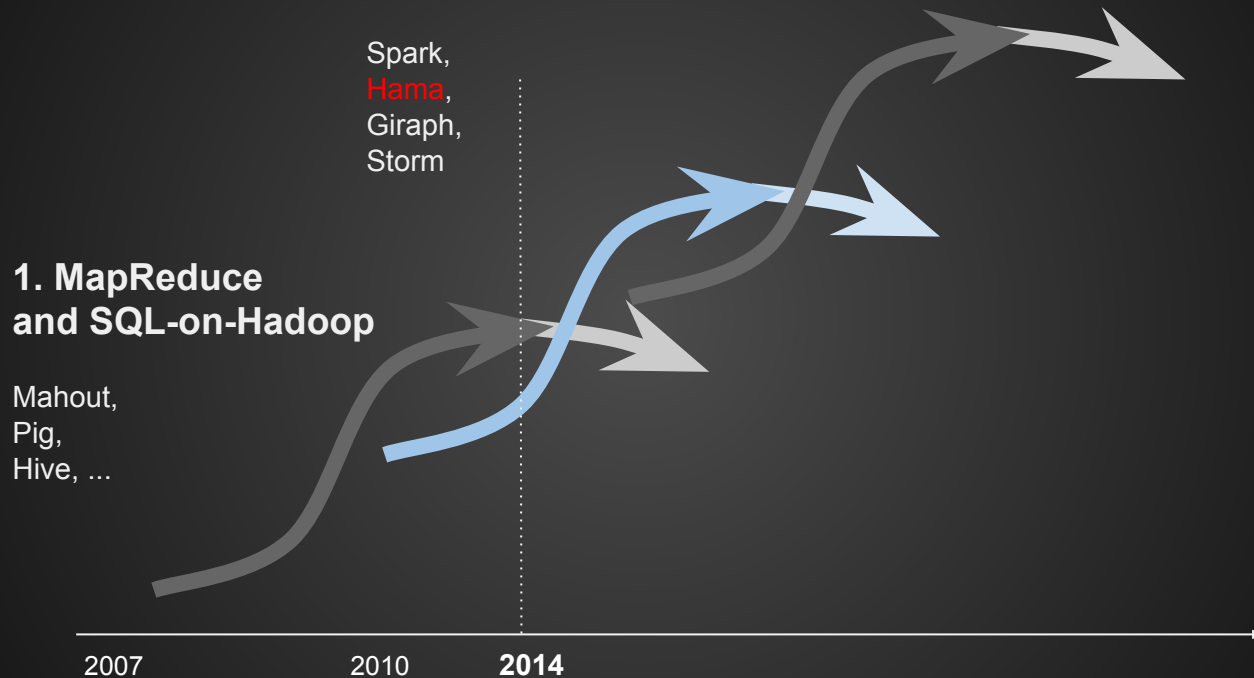
1. MapReduce and SQL-on-Hadoop

Mahout,
Pig,
Hive, ...

2007

2010

2014



- 1990 ~ : Web Documents

Web 2.0

Blog, Open API

Smartphone

Social Network

- ~ 2014 : Responsive Apps for multi-devices

- 1990 ~ : Server/Web Hosting

Google Apps

Cloud Computing

IaaS, PaaS, SaaS

- ~ 2014 : Cloud/App Hosting

- 1990 ~ : Text processing and mining

MapReduce

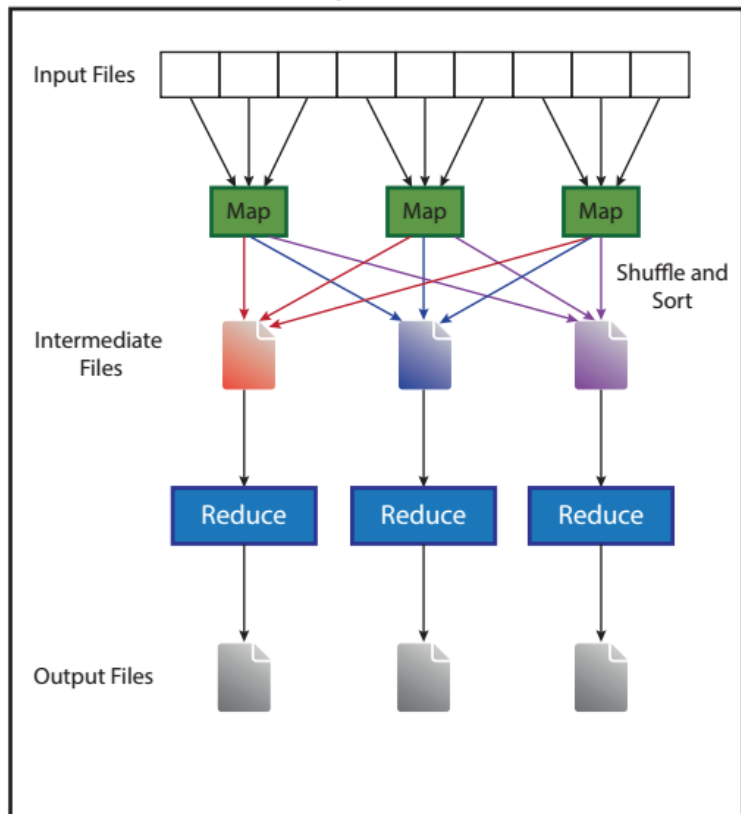
SQL-on-Hadoop

In-memory or Message-Passing

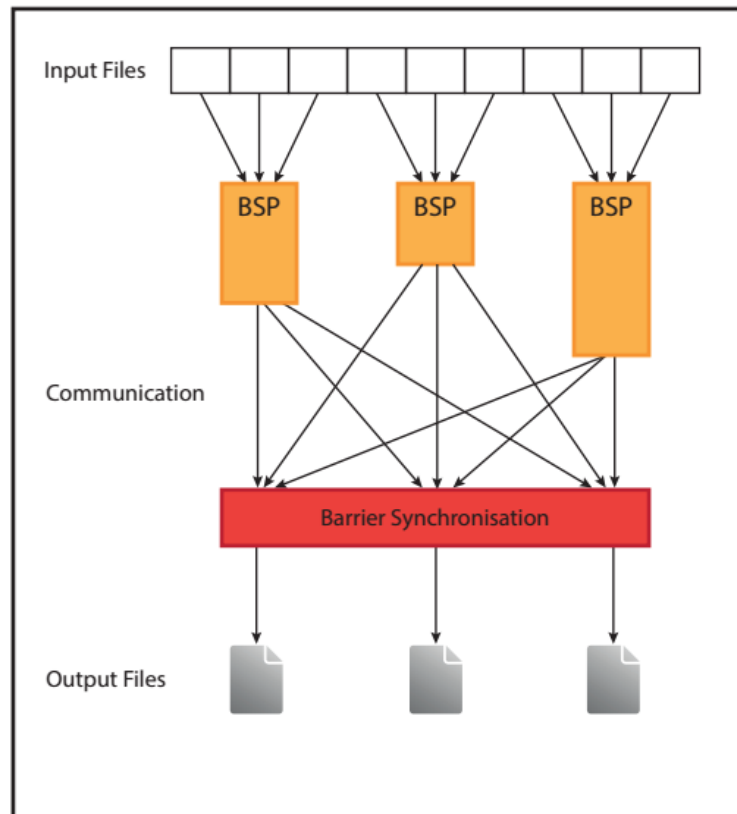
- ~ 2014 : Matrix, Mining Networks and Graphs

Hama^[hó:ma] is a general-purpose
BSP computing engine on Top of Hadoop

MapReduce



BSP



1 / 154



Apache Hama is listed on the Best Open Source Big Data tools, Bossie Awards 2013

Bossie Awards 2013: The best open source big data tools

InfoWorld's top picks in the expanding Hadoop ecosystem, the NoSQL universe, and beyond

By InfoWorld staff, [InfoWorld](#), September 17, 2013

Subscribe to slideshows: [RSS](#)

Slideshow

1 Comment



Share

92



Like

79



Previous

6 of 16

Next

Apache Hama

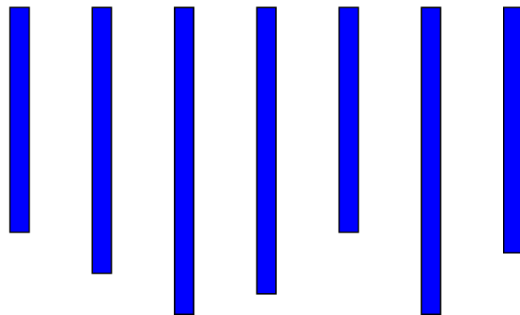
Like Graph, [Apache Hama](#) brings Bulk Synchronous Parallel processing to the Hadoop ecosystem and runs on top of the Hadoop Distributed File System. However, whereas Giraph focuses exclusively on graph processing, Hama is a more generalized framework for performing massive matrix and graph computations. It combines the advantages of Hadoop compatibility with a more flexible programming model for tackling data-intensive scientific applications.

— *Indika Kotakadeniya*

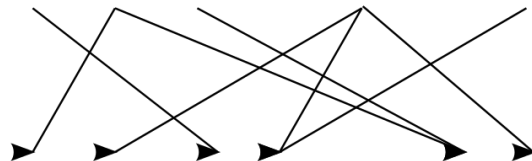


Processors

Local
Computation



Communication

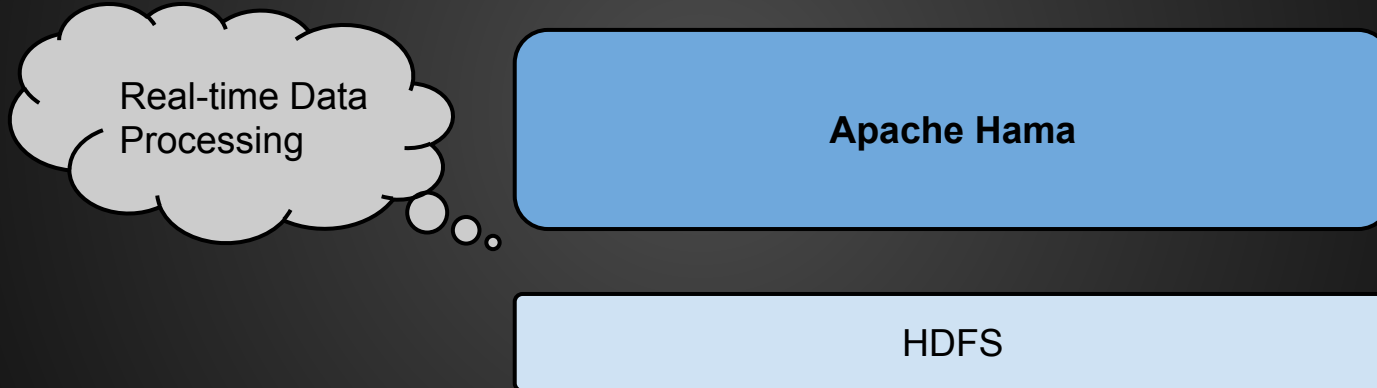


Barrier
Synchronisation

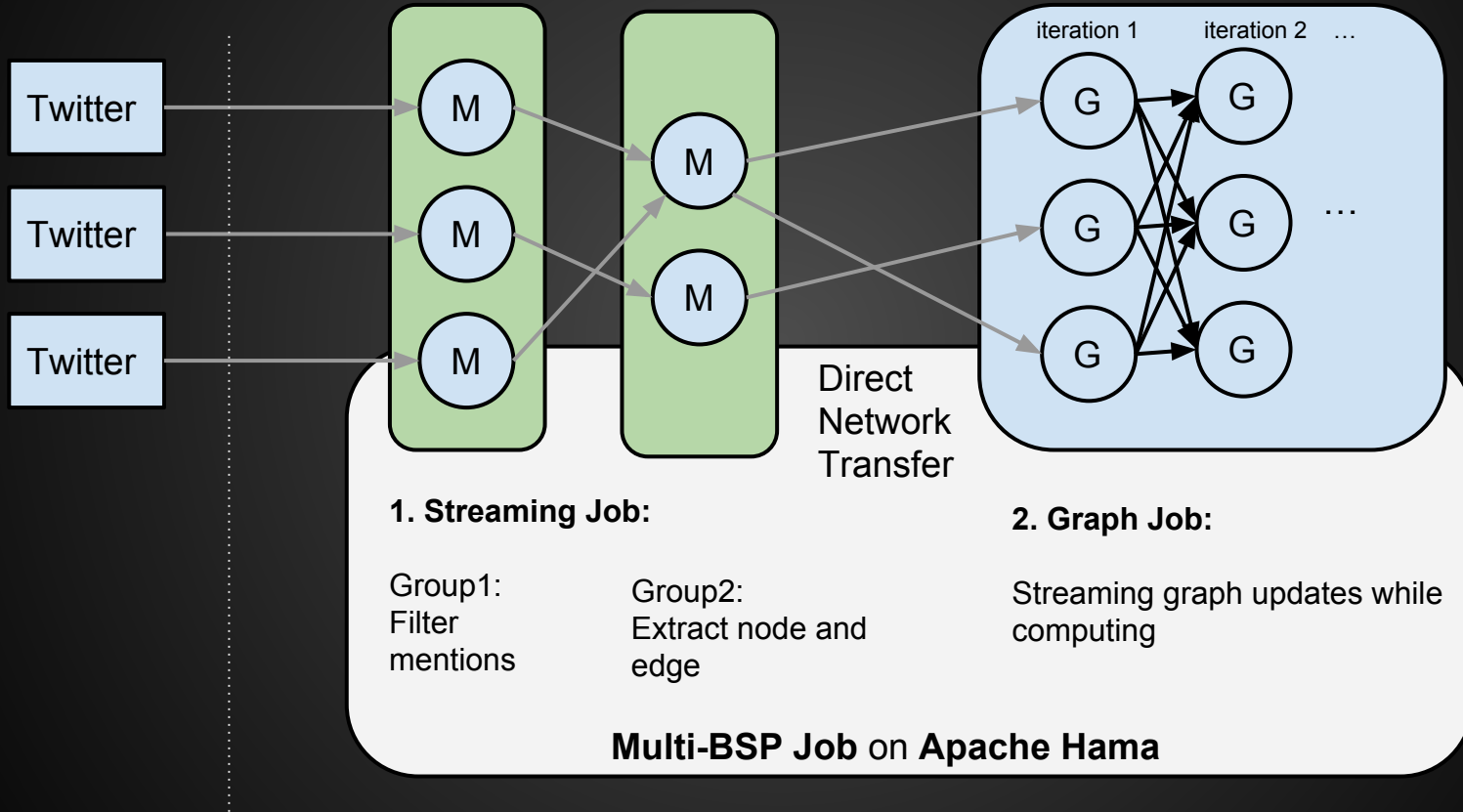


	Streaming	Graph	Machine Learning	Incremental Learning
Hama (General-purpose BSP)	○	○	○	○
Spark (Databricks) (In-memory MapReduce)	○	○ (GraphX)	○	✗
GraphLab (Asynchronous graph computing)	✗	○	○	○ (Limited)
Giraph (BSP-based graph computing)	✗	○	✗	✗

Think about **the Spam Filtering** of Google's Gmail!



Another Example



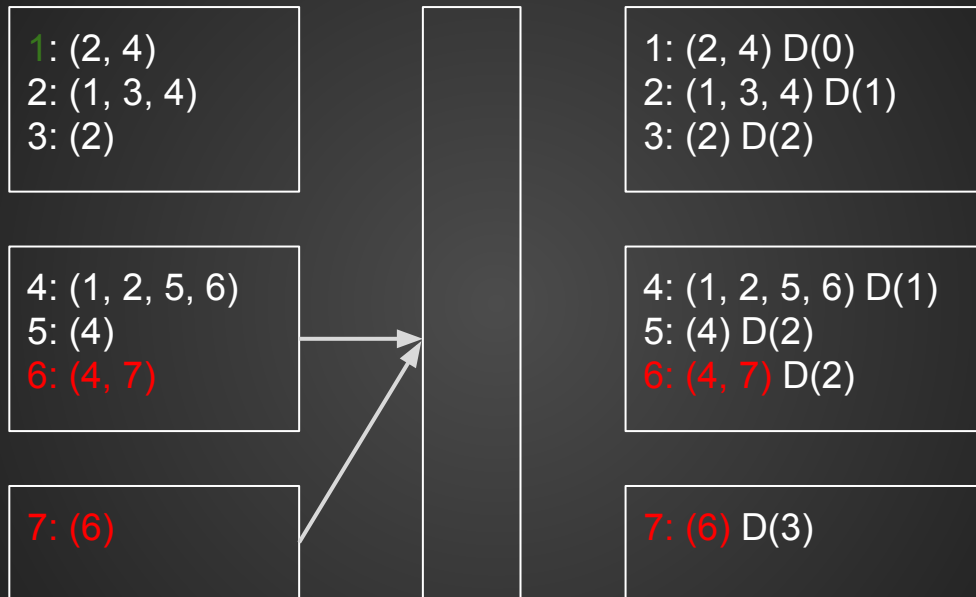
Appendix.

Why all platforms uses **BSP-style** for graph-parallel?

MR version: Shortest Path

- A map task receives a node n as a key, and $(D, \text{points-to})$ as its value
- D is the distance to the node from the start
- points-to is a list of nodes reachable from n
 - $\forall p \in \text{points-to}, \text{emit}(p, D+1)$
- Reduce task gathers possible distances to a given p and selects the minimum one

BSP version: Shortest Path



Why **Google's** Pregel (graph)
and DistBelief (deep learning) **uses BSP-style?**



Thomas Jungblut

@tjungblut



Follow

@Seoul_Tech @eddieyoon spark is garbage for that computation, because all your weight updates are in the RDD and that's a network b'neck



Reply



Retweet



Favorite



More

9:37 PM - 16 Jul 2014

Apache Hama's **Advanced Analytics** Examples:

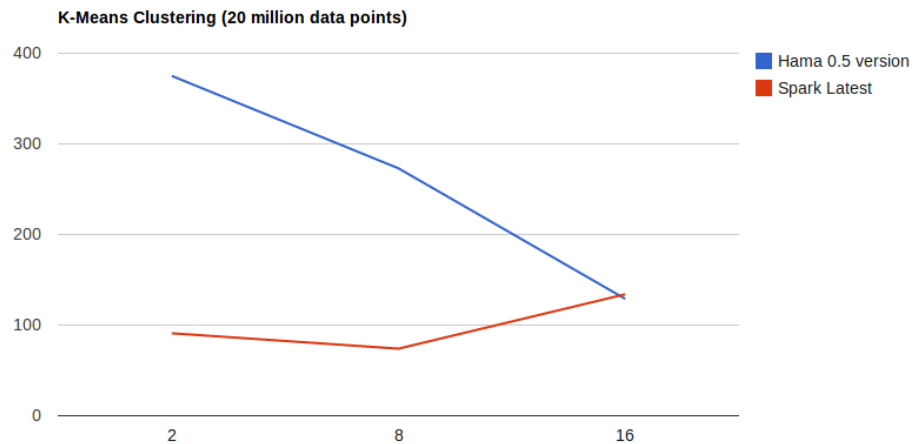
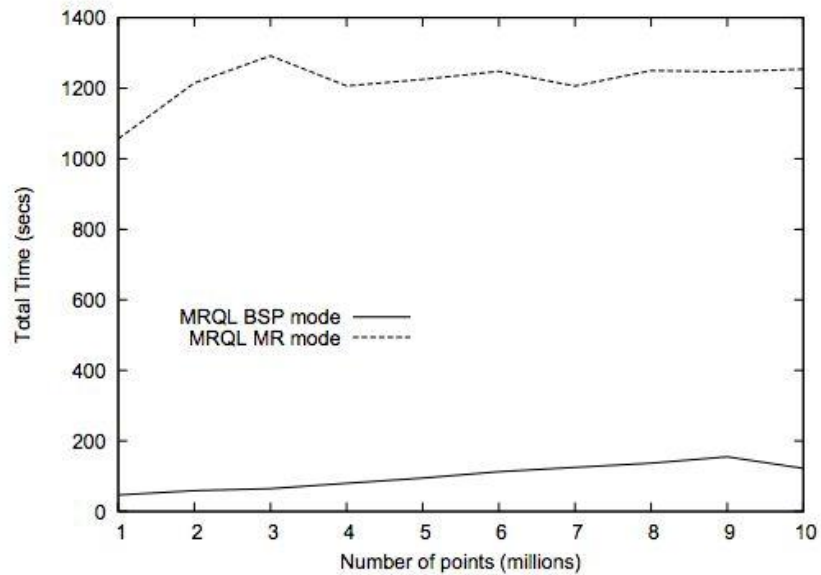
- Sparse Matrix-Vector Multiplication
- Semi-Clustering
- **K-means Clustering**
- Neural Networks
- **Gradient Descent**
- **PageRank**
- Single Source Shortest Path
- Bipartite Matching

Hama Supports:

Hadoop 1.0

Hama on Hadoop 2.0 YARN

Hama on Apache Mesos



Apache Hama

at Sogou





Sogou.com runs **7,200 cores Hama cluster** for SiteRank.

- SiteRank is the ranking generated by applying the classical PageRank algorithm to the graph of Web sites.
- Dataset is about 400GB, contains 6 Billion edges.

L3.1 The symbolic solution of $y' = -2xy$, $y(0) = 2$ is $y = 2e^{-x^2}$. Derive it and display a full answer check

Solution

$$y' + 2xy = 0$$

$$\frac{(wy)'}{w} = 0$$

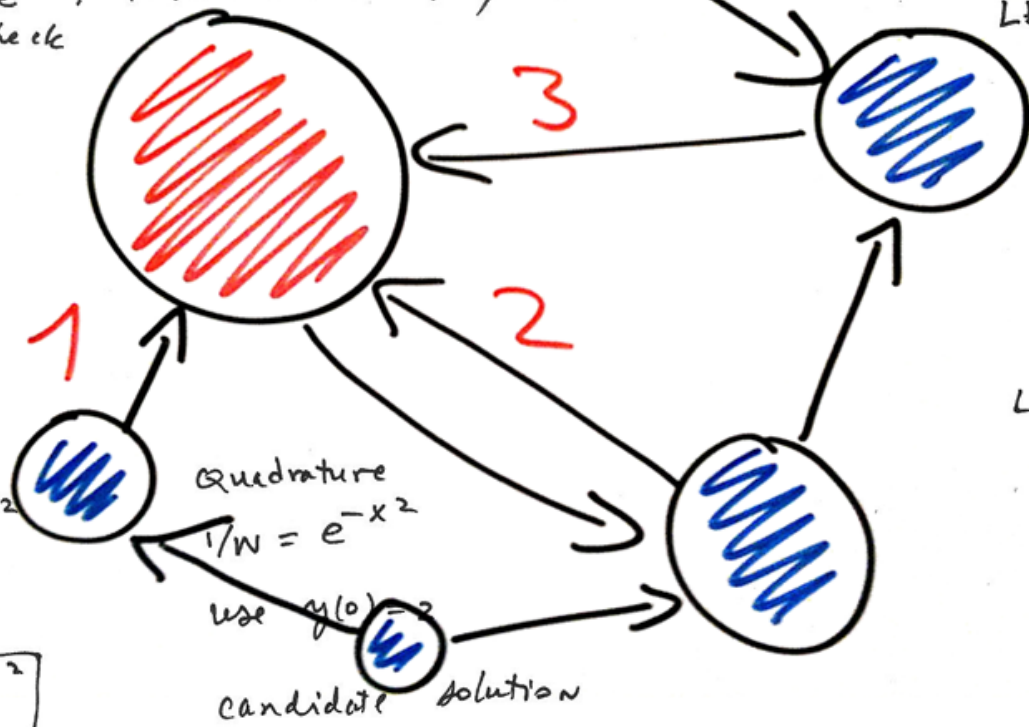
$$(wy)' = 0$$

$$wy = c$$

$$y = c e^{-x^2}$$

$$2 = c e^0$$

$$y = 2e^{-x^2}$$

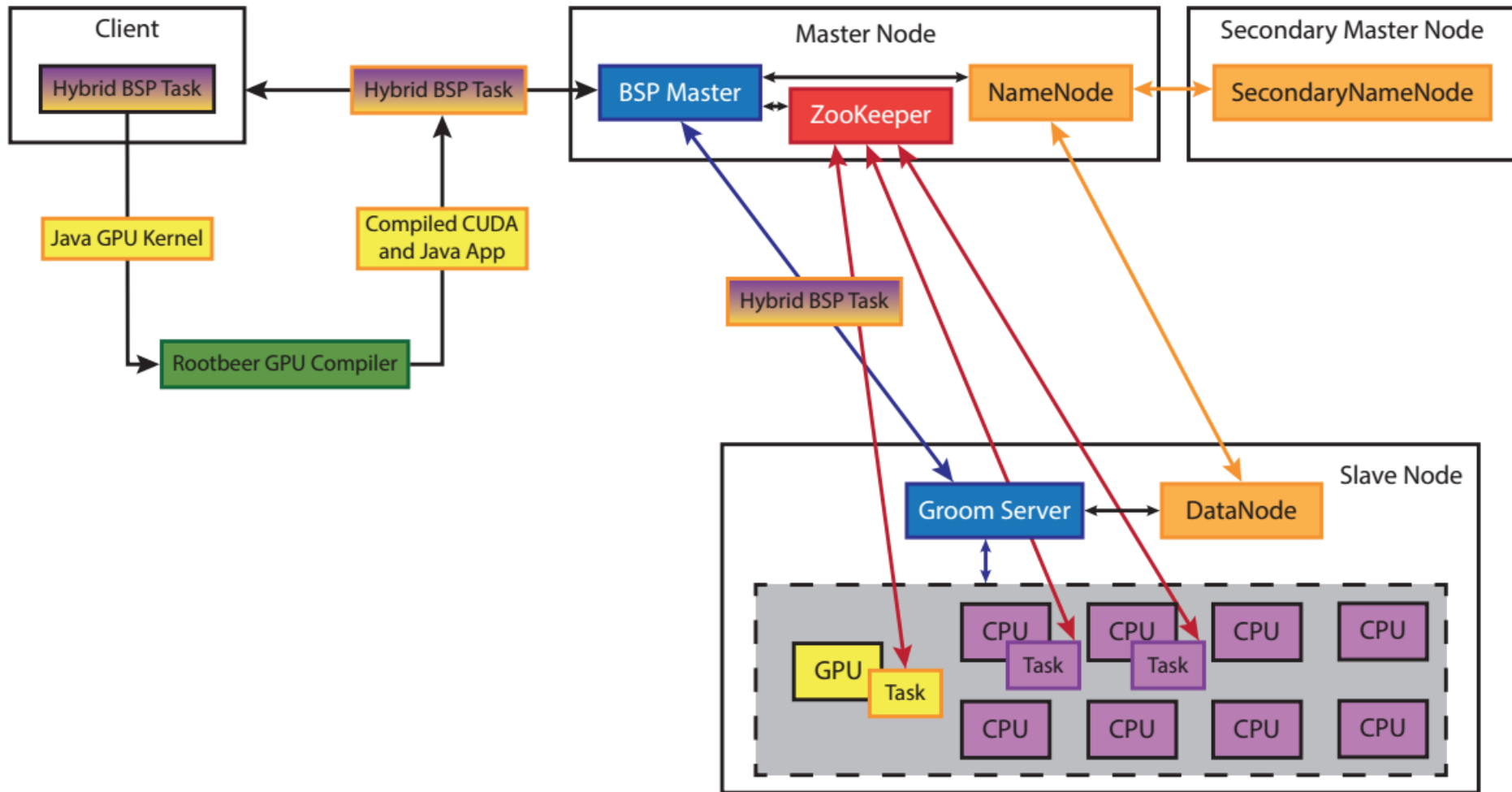


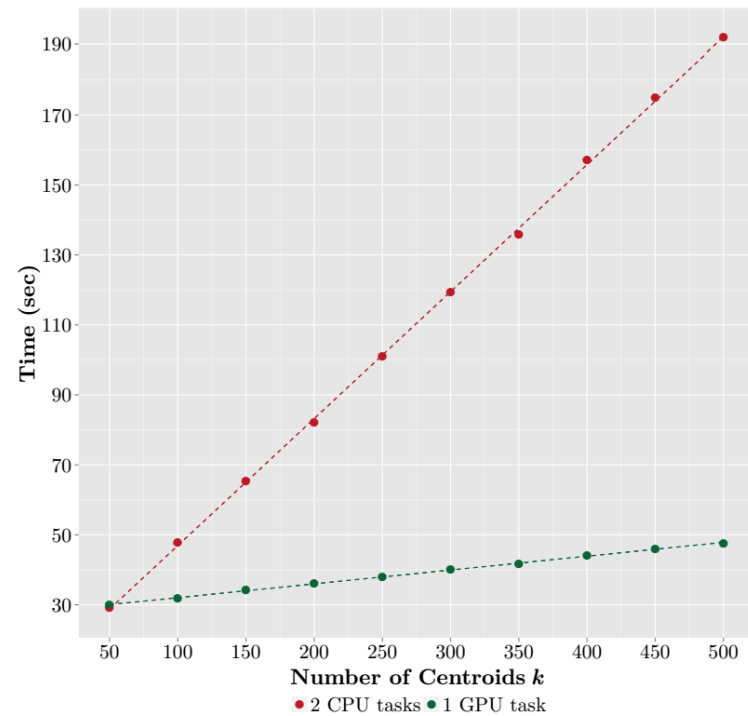
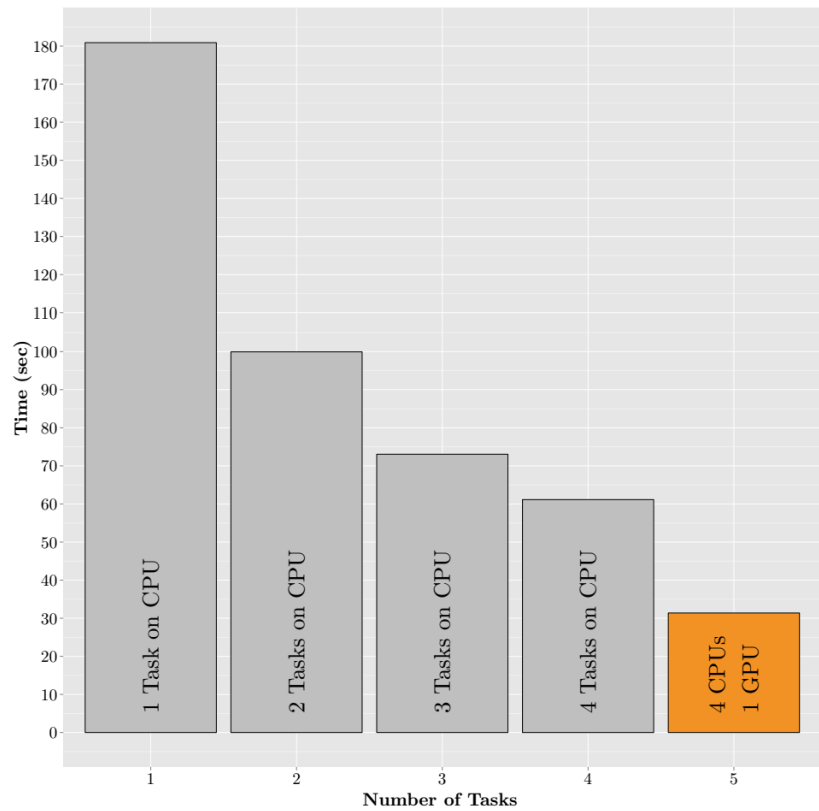
Answer check

$$\begin{aligned} \text{LHS} &= y' \\ &= (2e^{-x^2})' \\ &= 2(-2x)e^{-x^2} \\ &= (-2x)2e^{-x^2} \\ &= (-2x)y \\ &= \text{RHS} \end{aligned}$$

$$\begin{aligned} \text{LHS} &= y(0) \\ &= 2e^{-x^2} \Big|_{x=0} \\ &= 2e^0 \\ &= 2 \\ &= \text{RHS} \end{aligned}$$

The future features of Apache Hama:
Kryo serialization, Rootbeer GPU acceleration (**Martin
Illecker, University of Innsbruck**), ..., etc.





References

- Hama Website <http://hama.apache.org/>
- Scientific Computing in the Cloud with Apache Hadoop and Apache Hama on GPU by Martin Illecker

If you want to be one of us, **be one of us.**

Thanks!