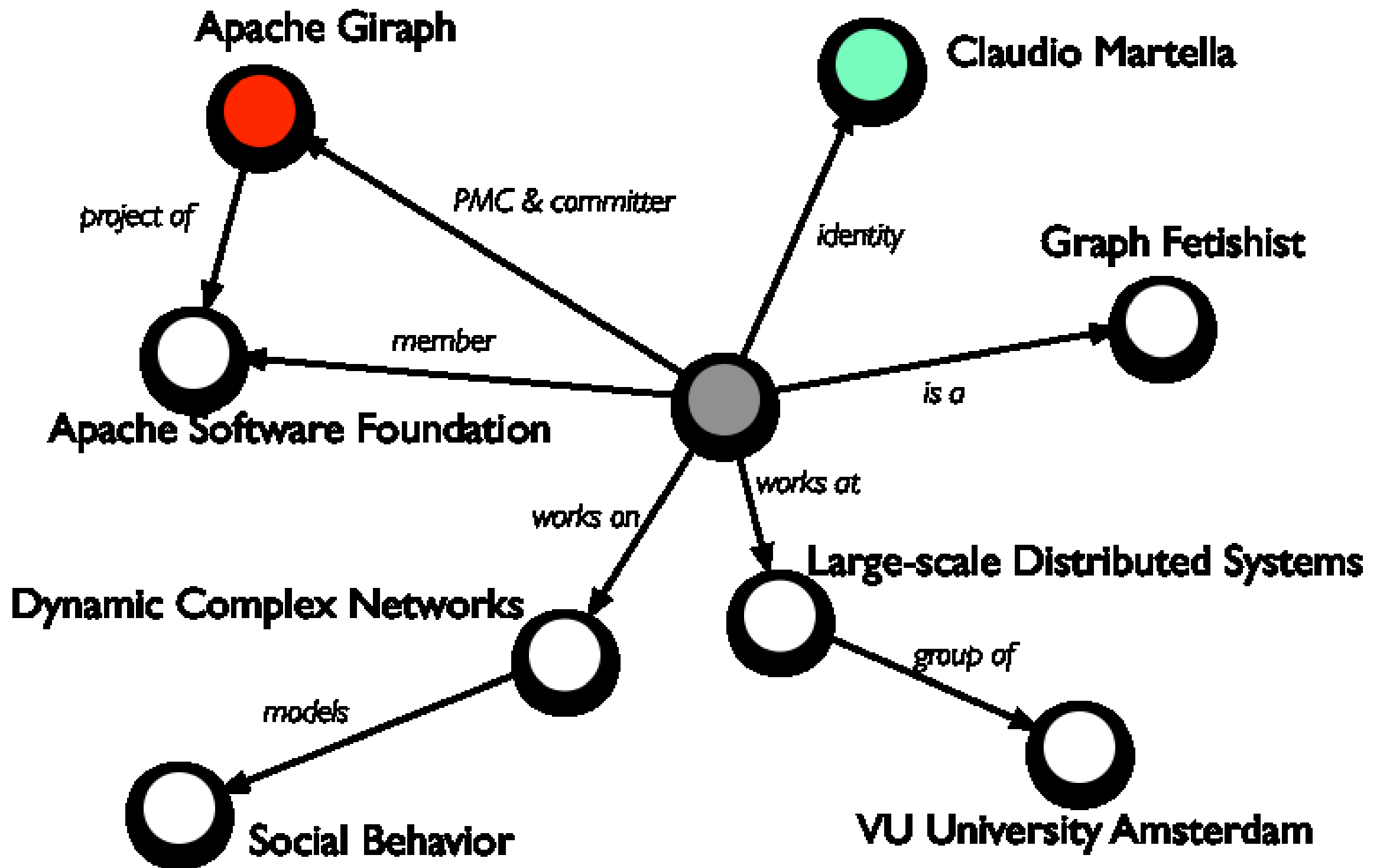


Apache Giraph

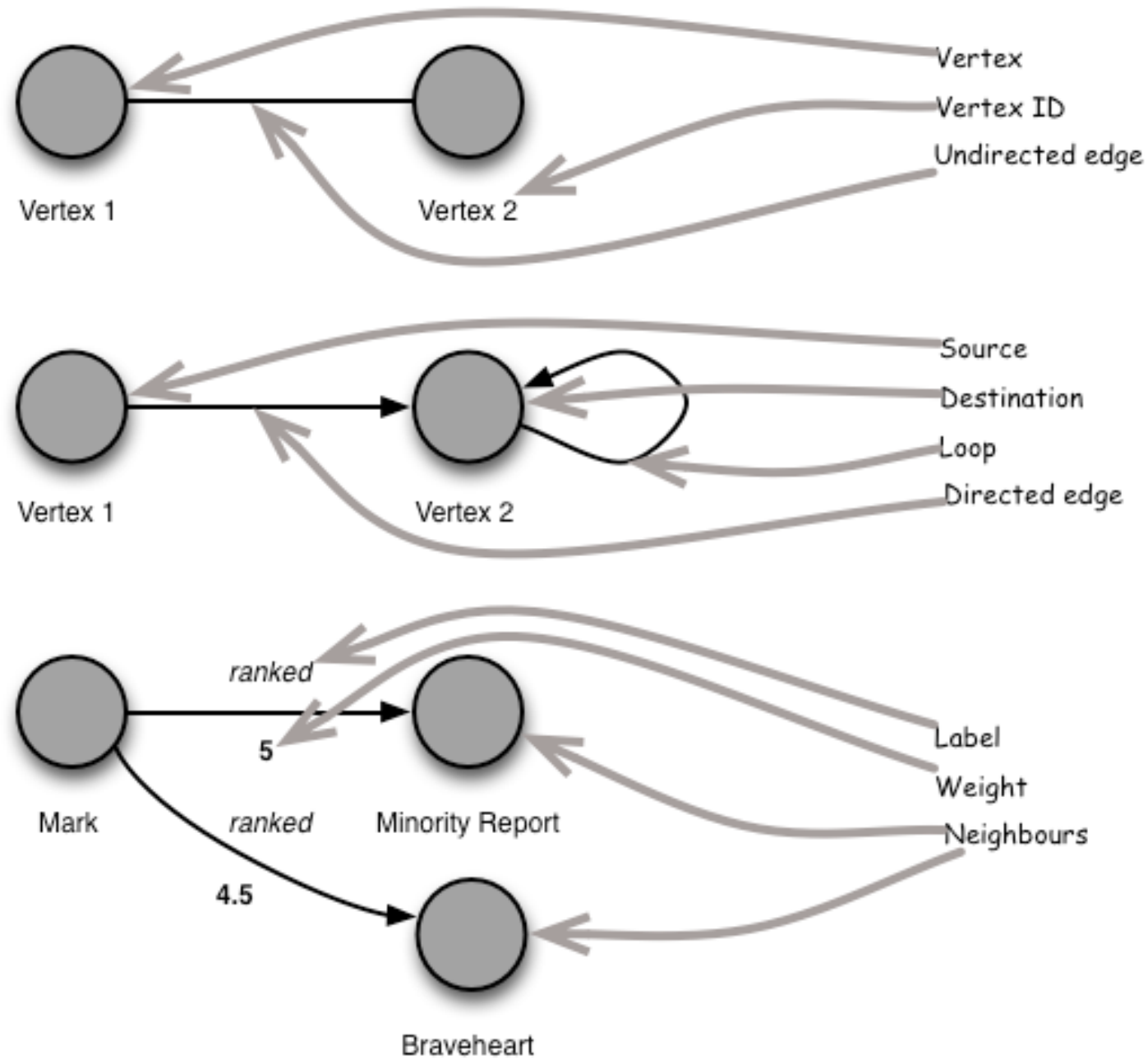
Large-scale **Graph** Processing on Hadoop

Claudio Martella <claudio@apache.org>
[@claudiomartella](#)

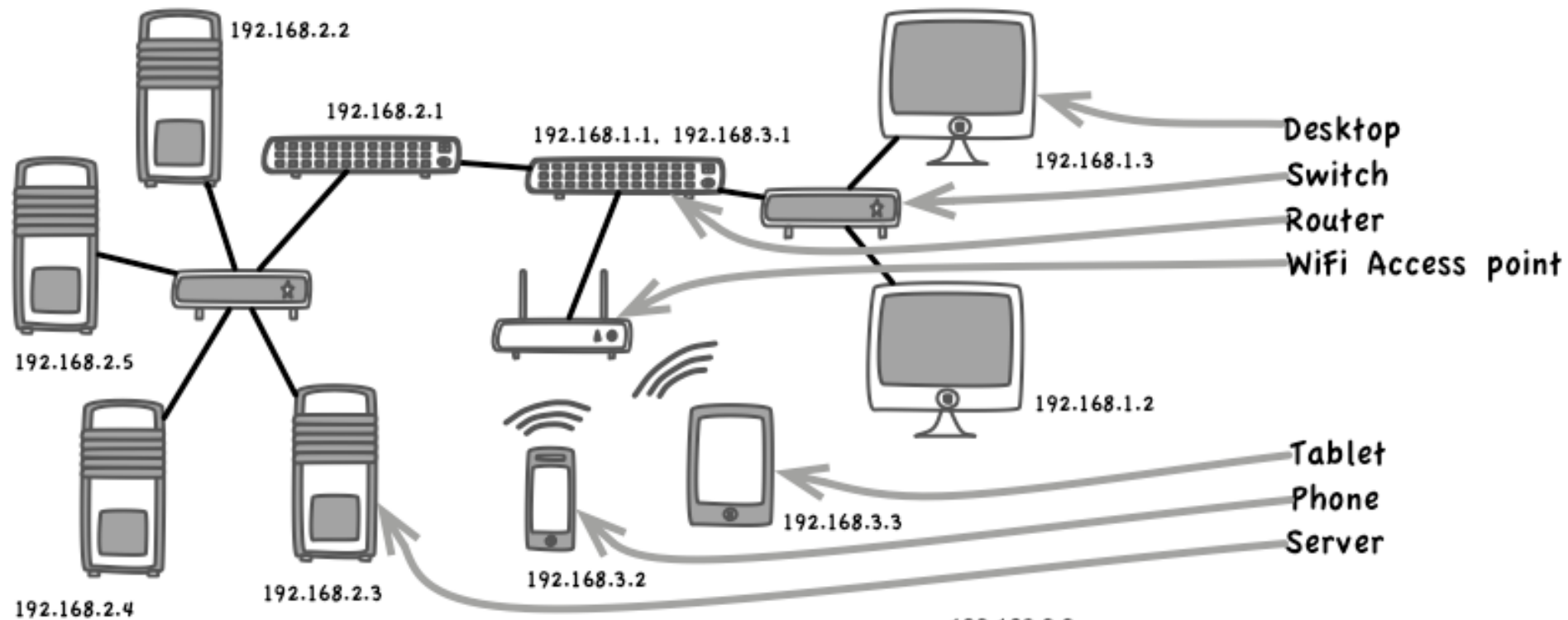
Hadoop Summit @ **Amsterdam** - 3 April 2014



Graphs are **simple**

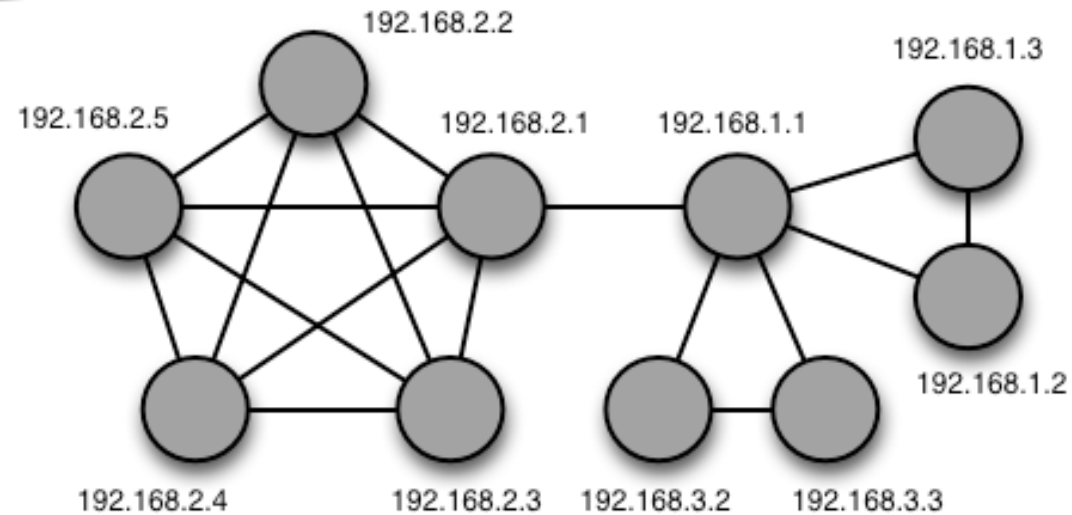


A computer network

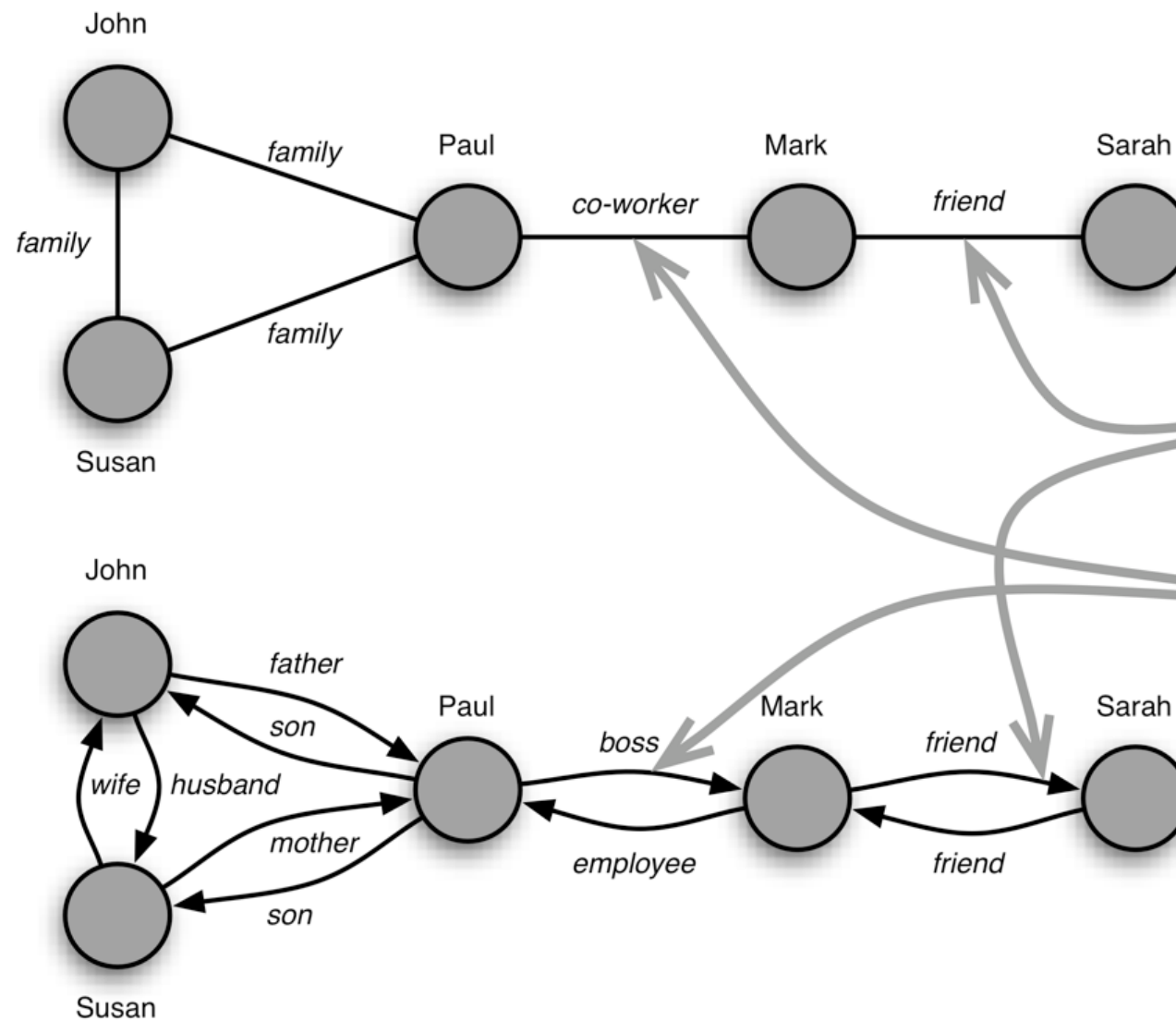


Note

1. There are three networks: servers, desktops and mobile.
2. They are connected through two routers/firewalls.
3. We ignored the switches and the access point in the graph.
4. Router 192.168.1.1 has two interfaces but we used one as vertex ID.



A social network

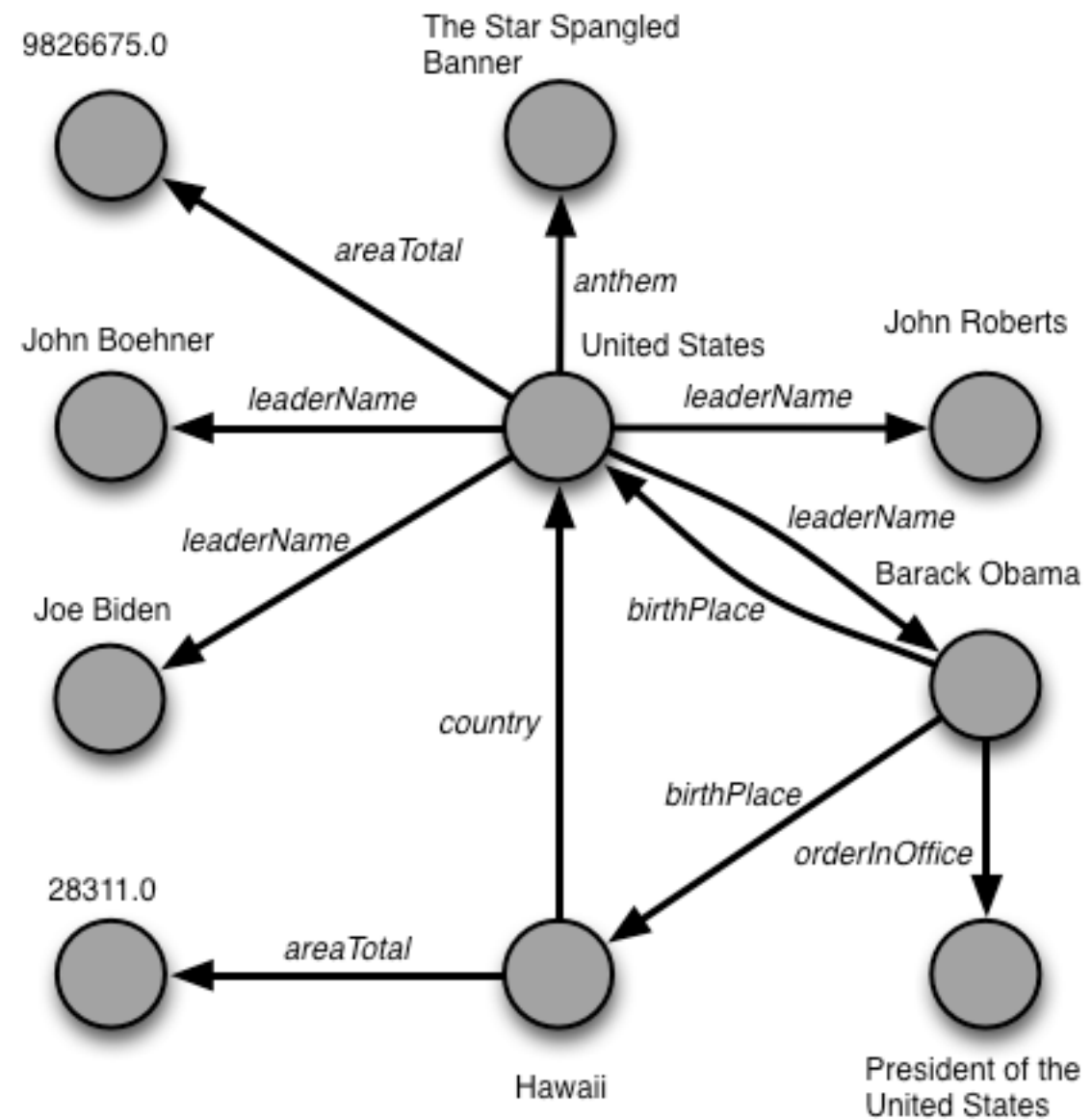


Note

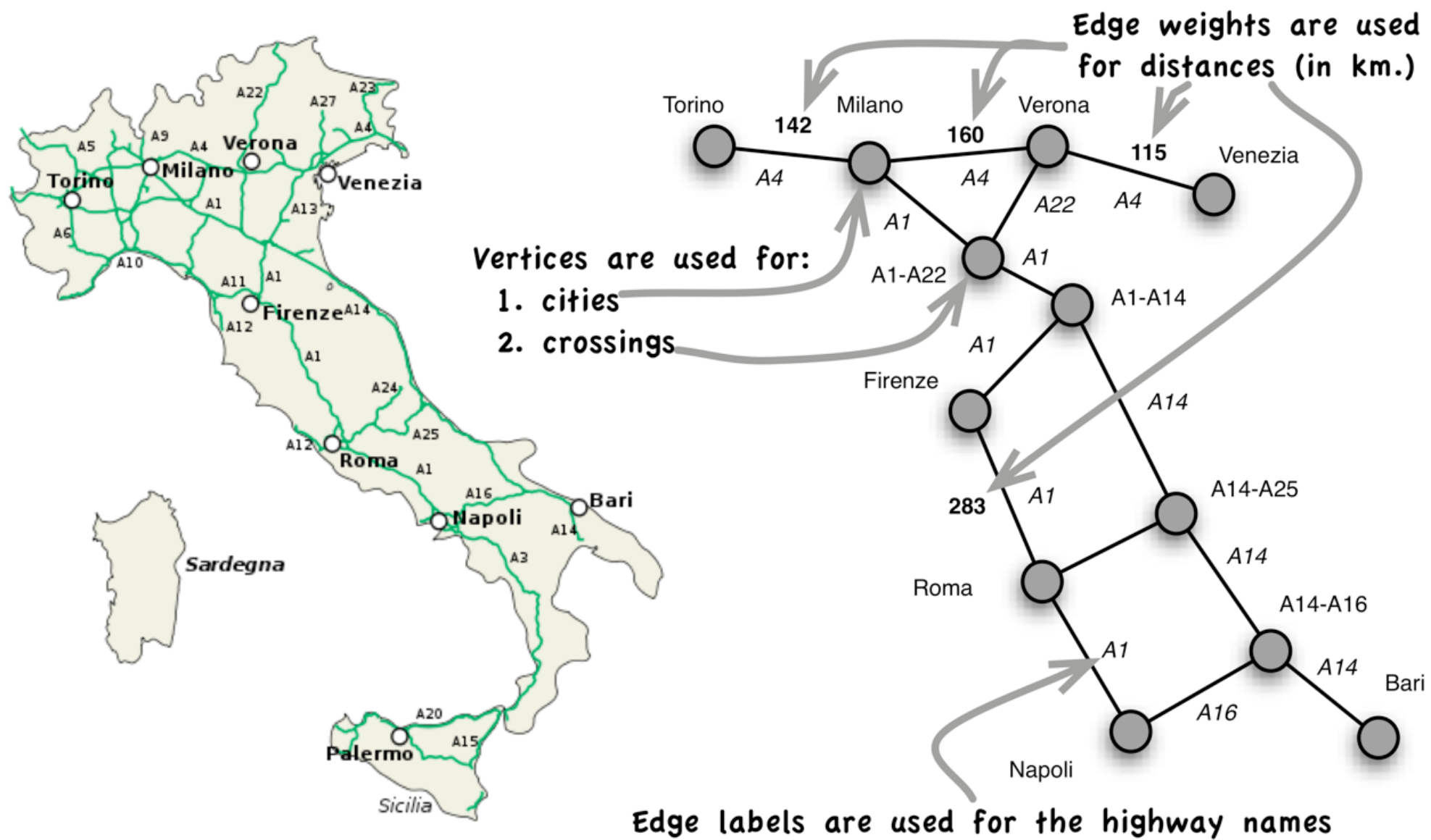
1. A symmetric relationship is substituted by two directed edges.
2. A relationship does not have to be substituted by two edges, but e.g. by a more specific one.

A semantic network

Subject	Predicate	Object
United States	areaTotal	9826675.0
United States	anthem	The Star Spangled Banner
United States	leaderName	Barack Obama
United States	leaderName	Joe Biden
United States	leaderName	John Boehner
United States	leaderName	John Roberts
Barack Obama	birthPlace	United States
Barack Obama	birthPlace	Hawaii
Barack Obama	orderInOffice	President of the United States
Hawaii	areaTotal	28311.0
Hawaii	country	United States



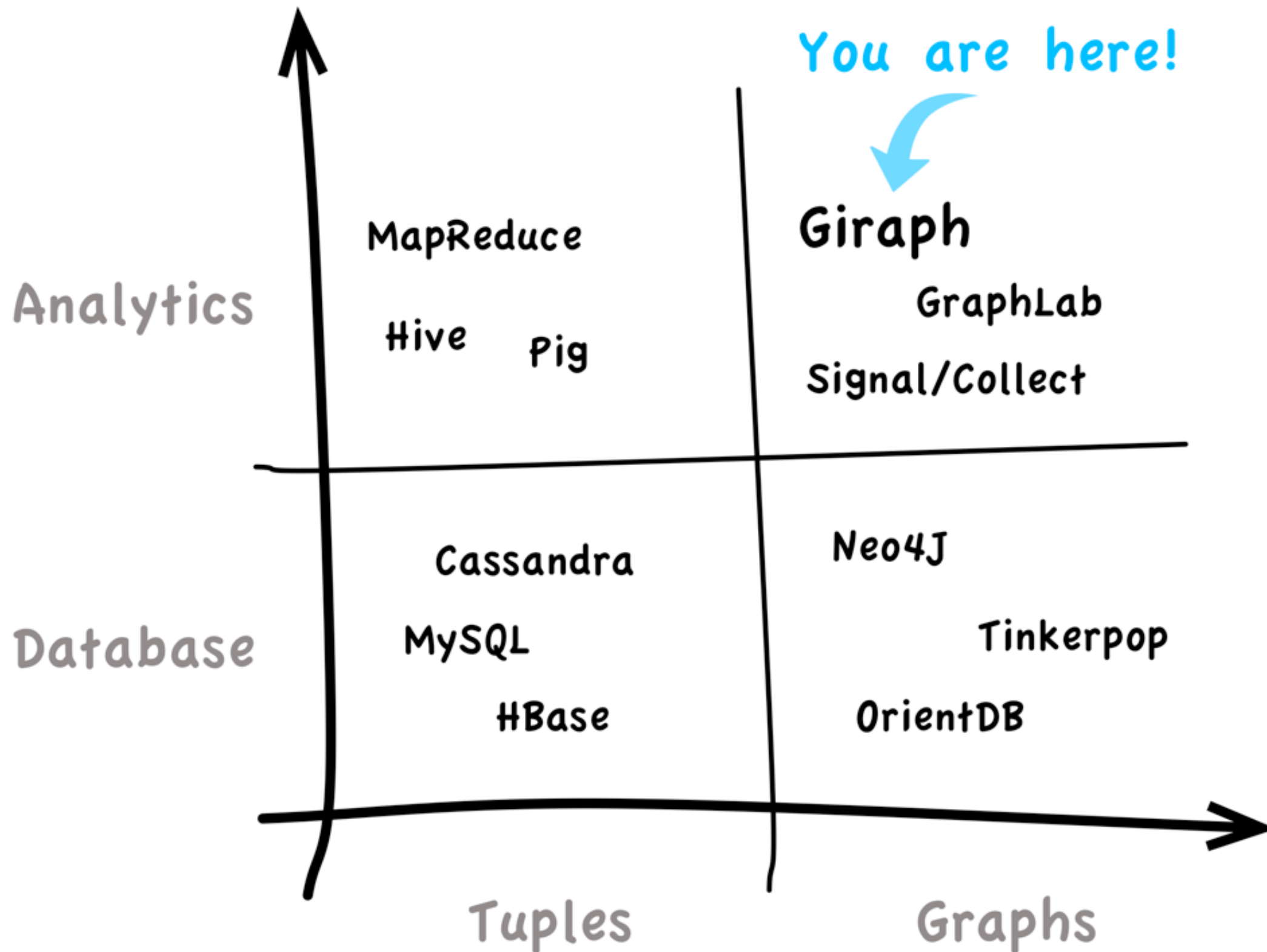
A map



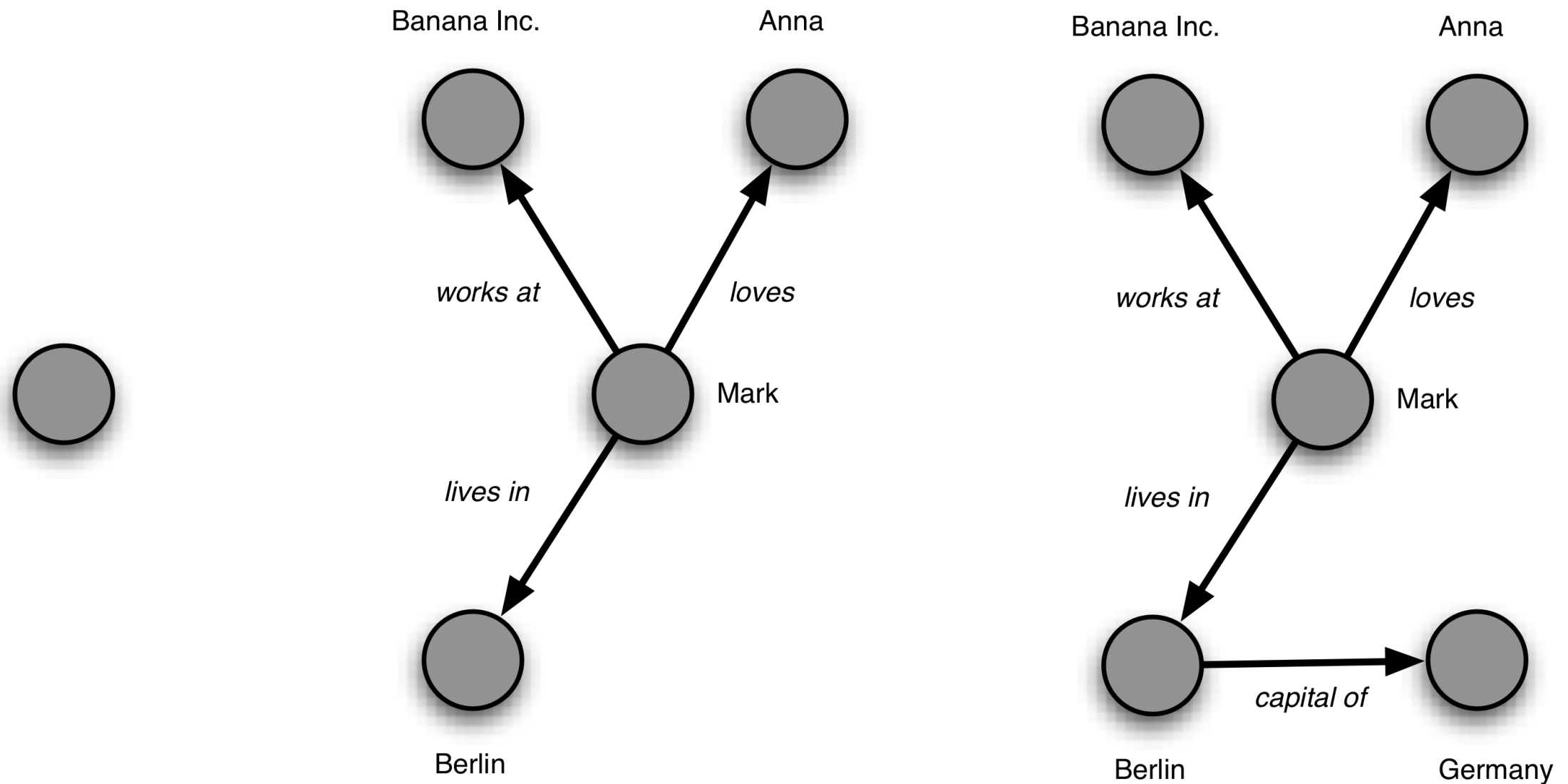
Graphs are huge

- Google's index contains 50B pages
- Facebook has around 1.1B users
- Google+ has around 570M users
- Twitter has around 530M users

VERY rough
estimates!



Graphs aren't *easy*



Graphs are *nasty*.

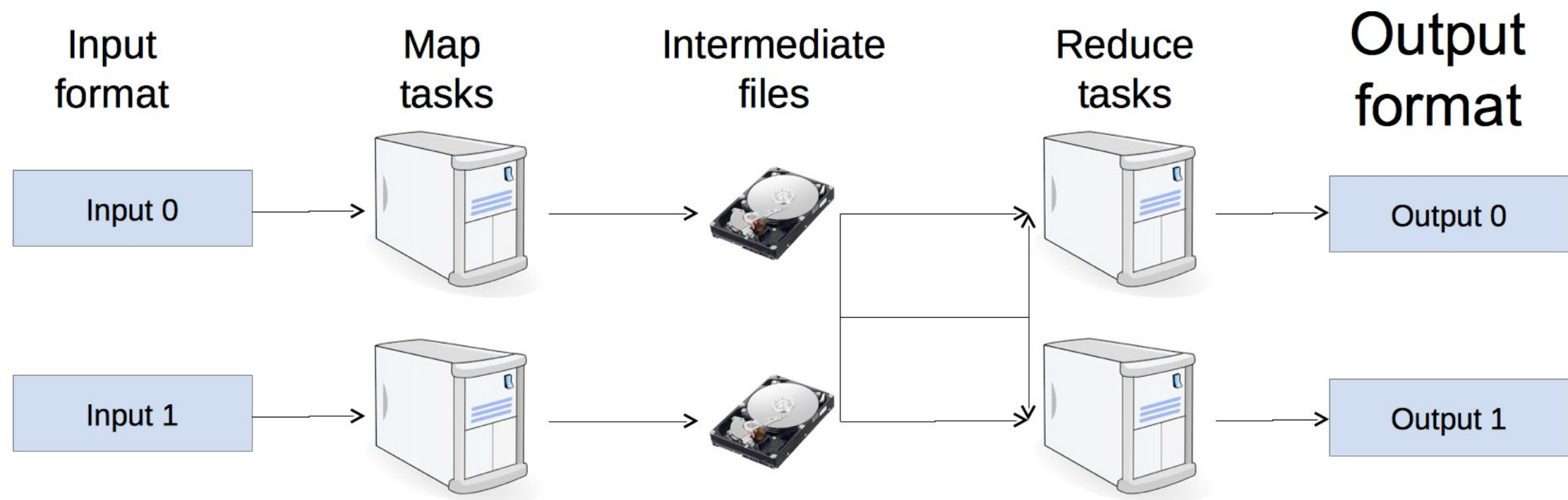
Each vertex depends
on its neighbours,
recursively.

Recursive problems
are nicely solved
iteratively.

PageRank in MapReduce

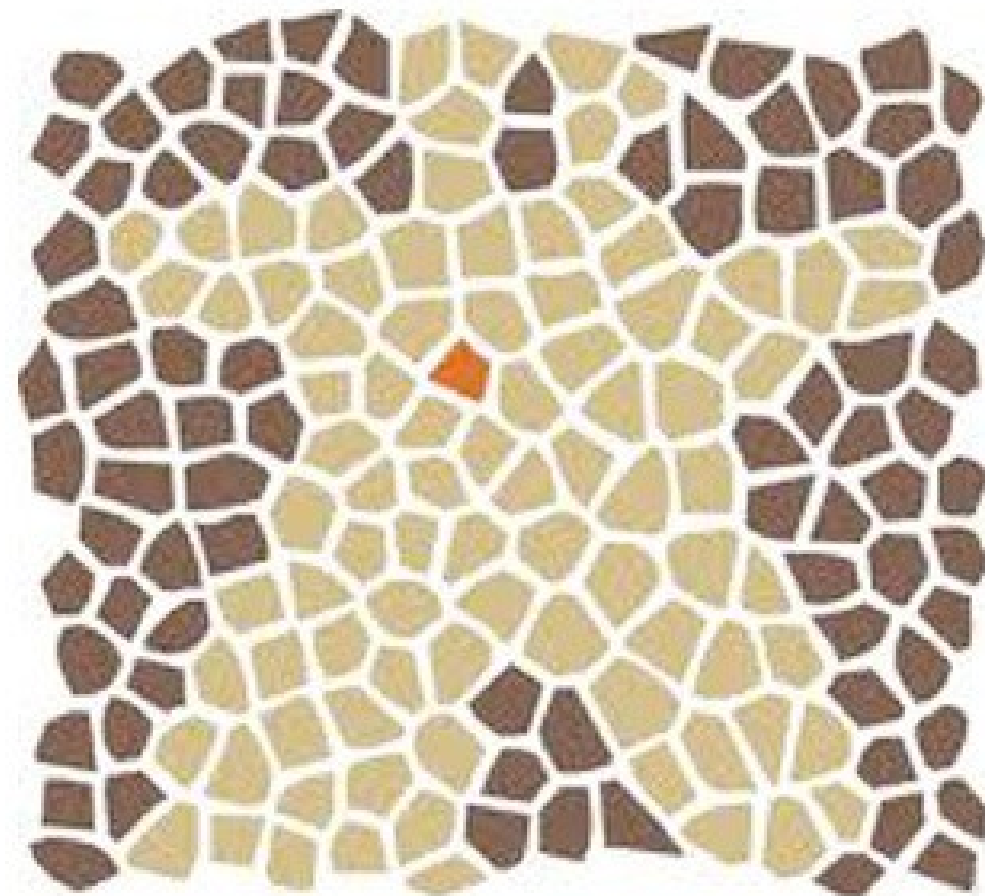
- **Record:** $\langle v_i, pr, [v_j, \dots, v_k] \rangle$
- **Mapper:** emits $\langle v_j, pr / \#neighbours \rangle$
- **Reducer:** sums the partial values

MapReduce dataflow



Drawbacks

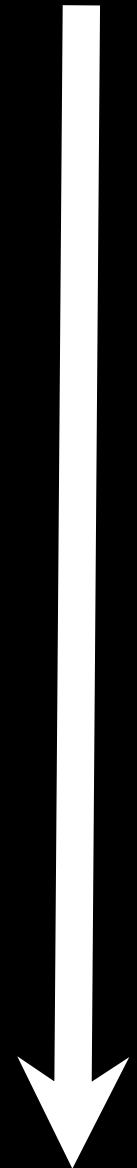
- Each job is executed **N** times
- Job **bootstrap**
- Mappers send PR values and **structure**
- Extensive **IO** at input, shuffle & sort, output



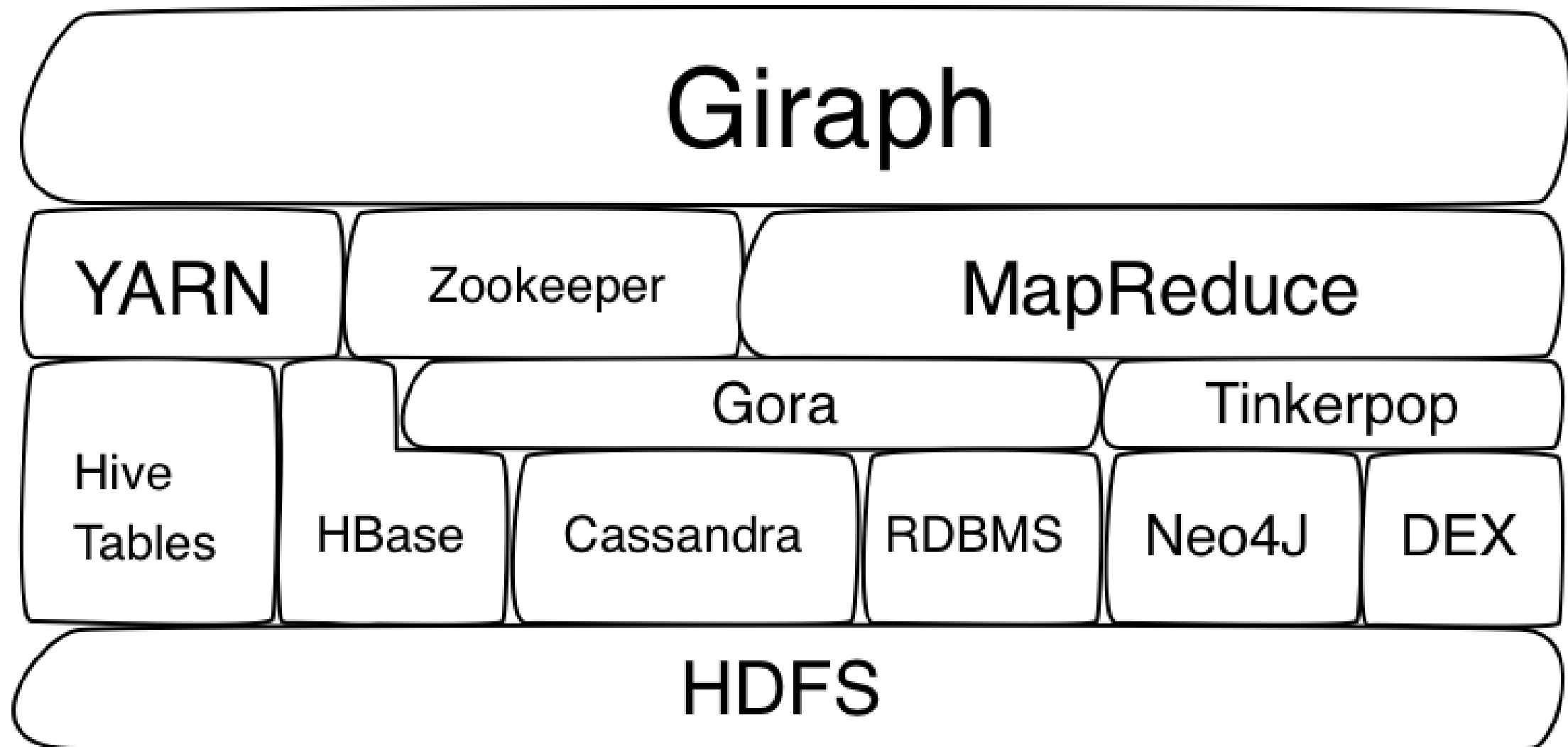
A P A C H E
G I R A P H

Timeline

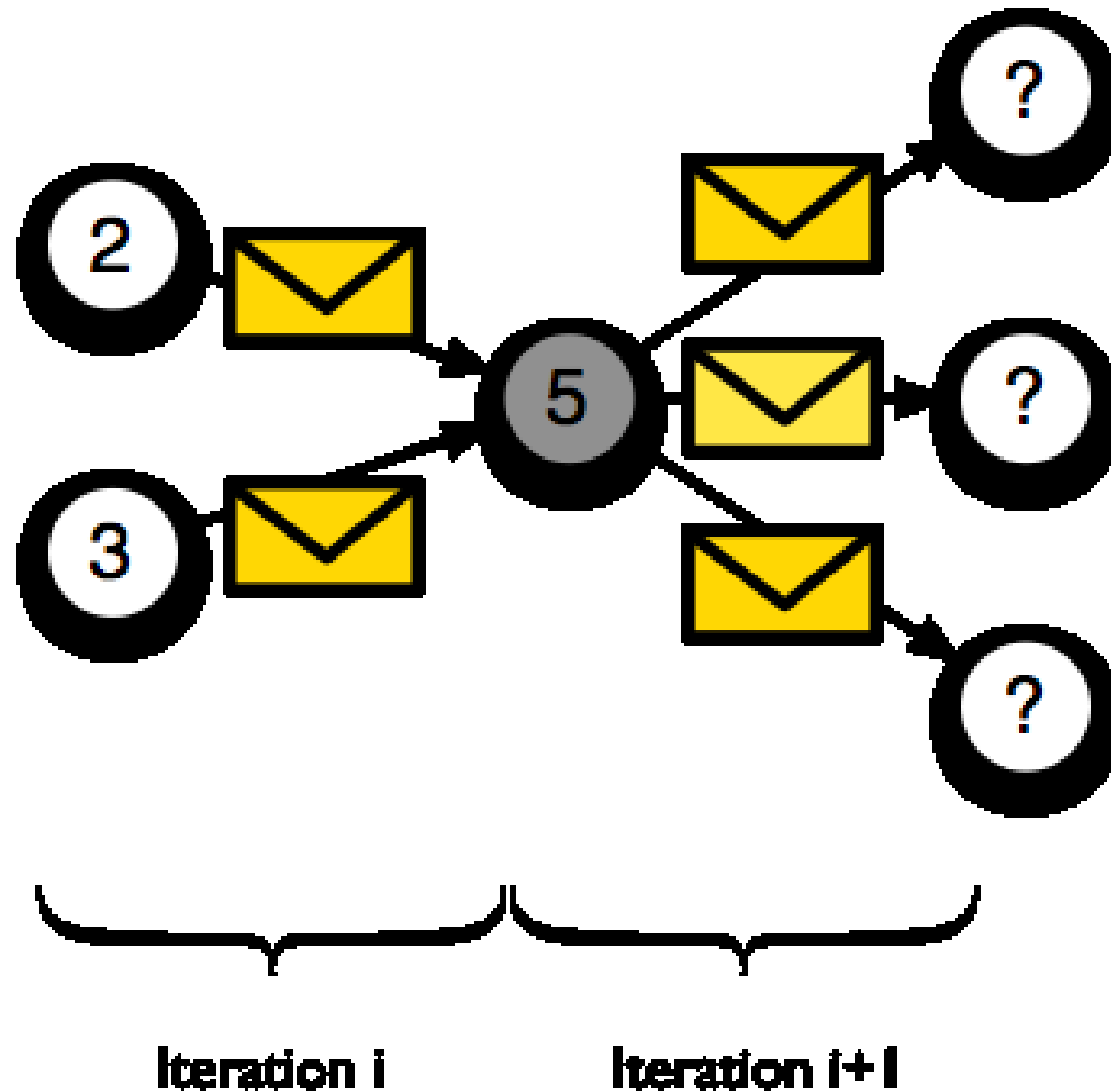
- Inspired by Google Pregel (2010)
- Donated to ASF by Yahoo! in 2011
- Top-level project in 2012
- 1.0 release in January 2013
- 1.1 release in days 2014



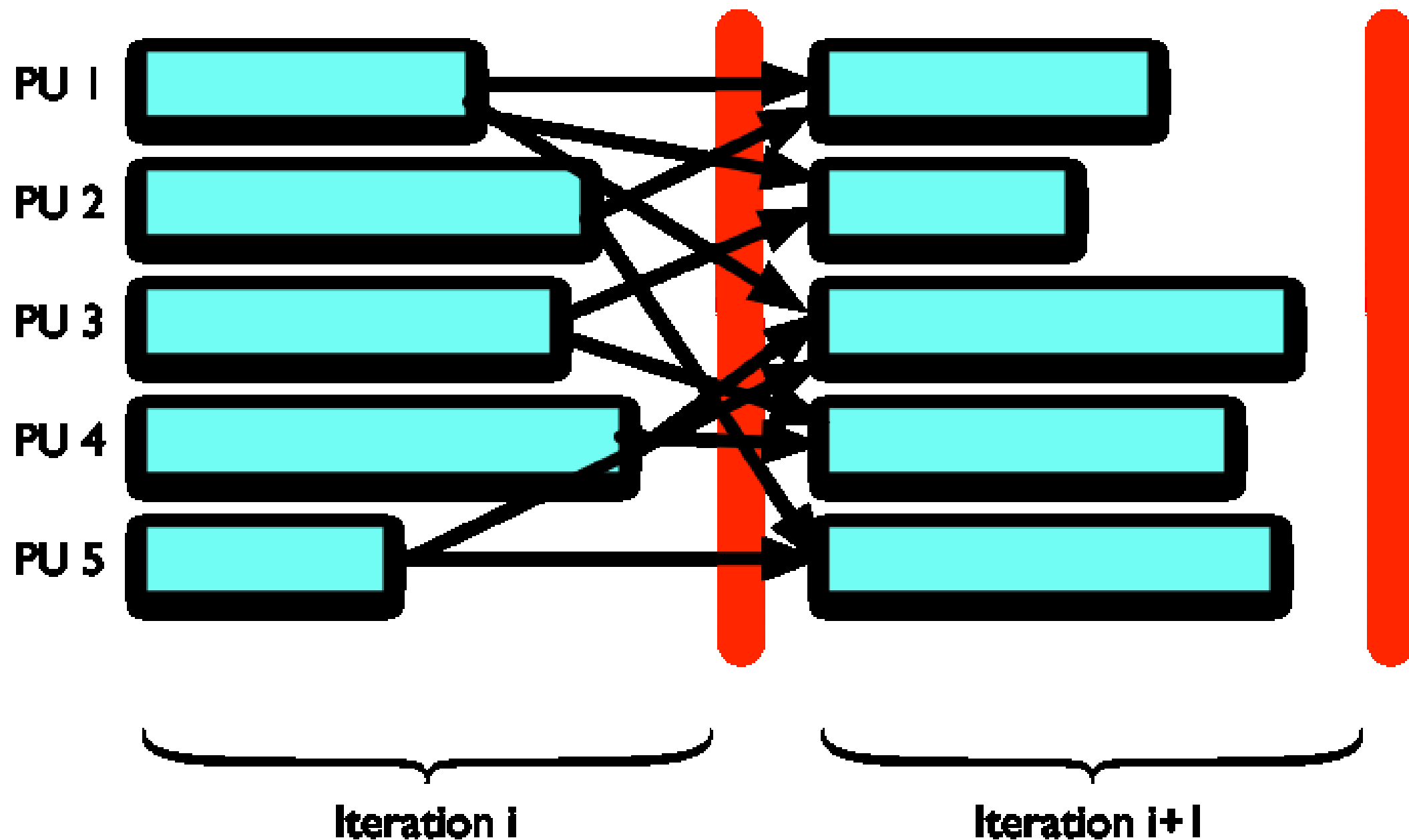
Plays well with Hadoop



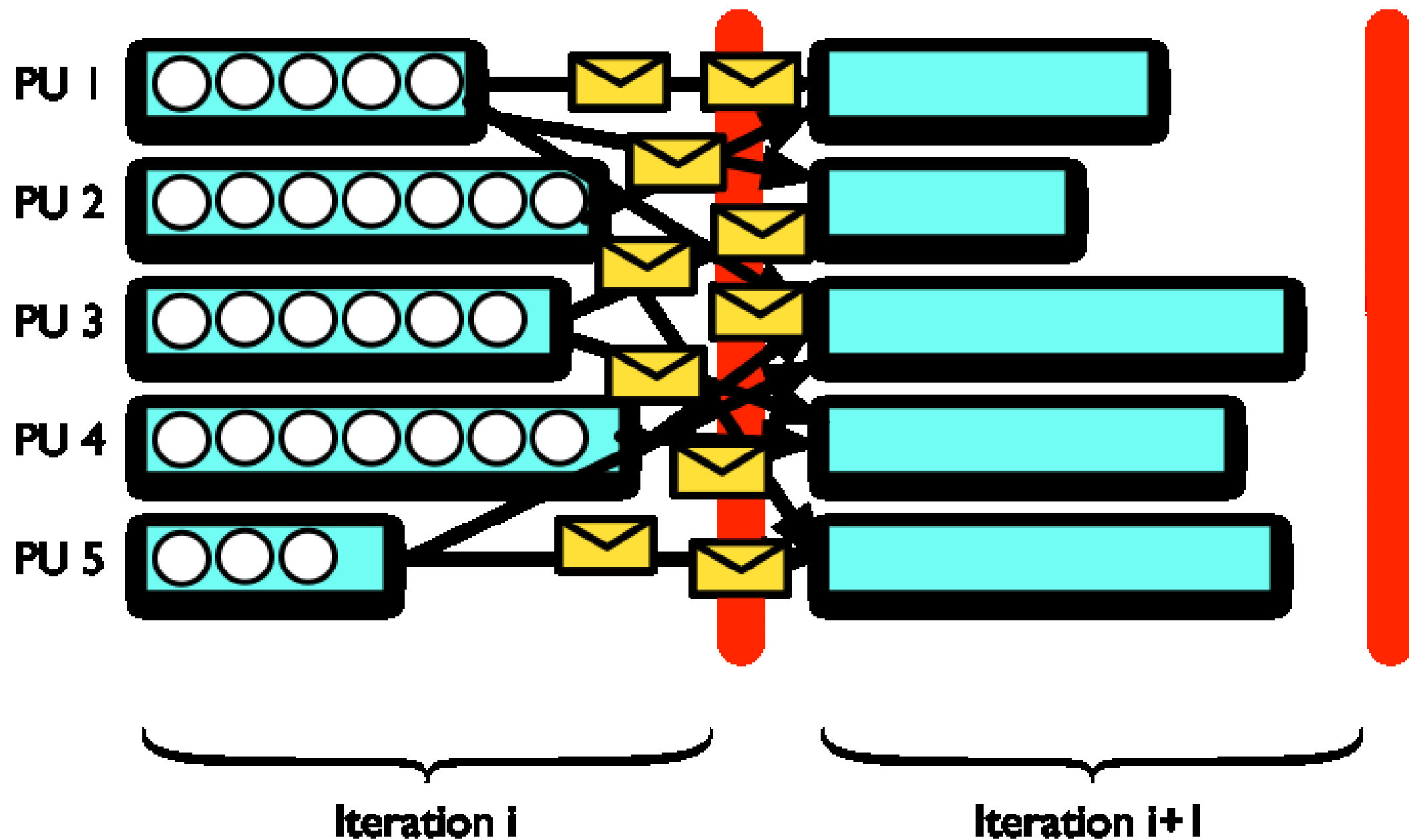
Vertex-centric API



BSP machine



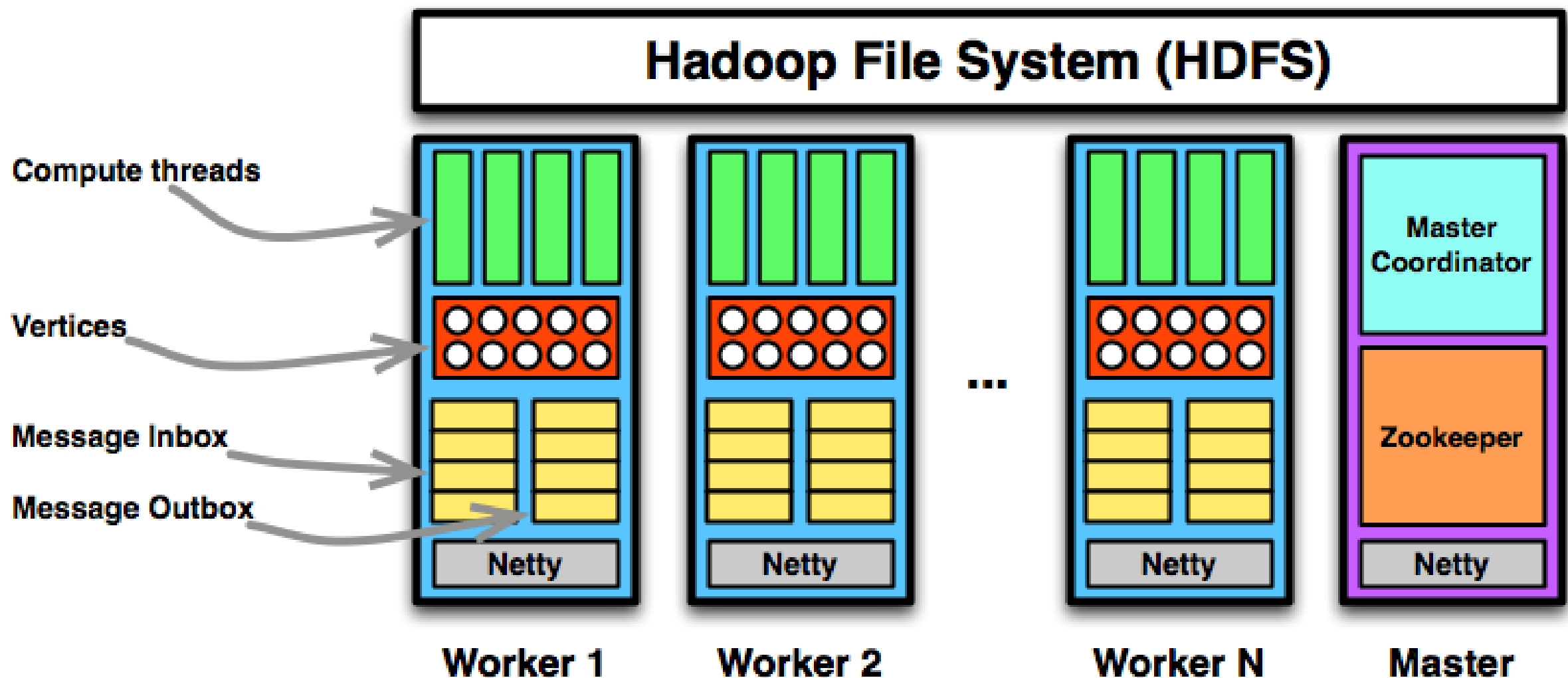
BSP & Giraph



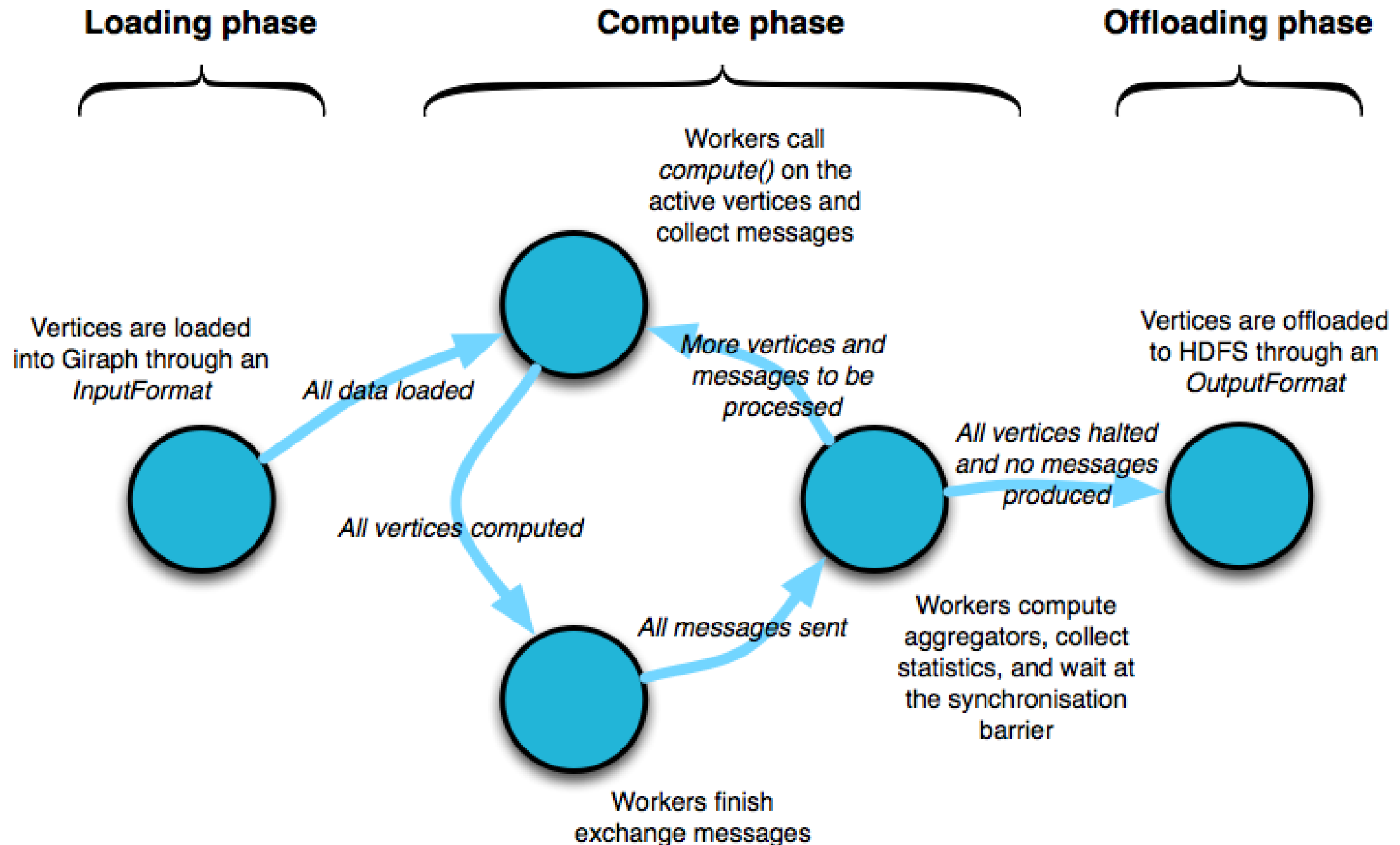
Advantages

- **No locks:** message-based communication
- **No semaphores:** global synchronization
- **Iteration isolation:** massively parallelizable

Architecture



Giraph job **lifetime**



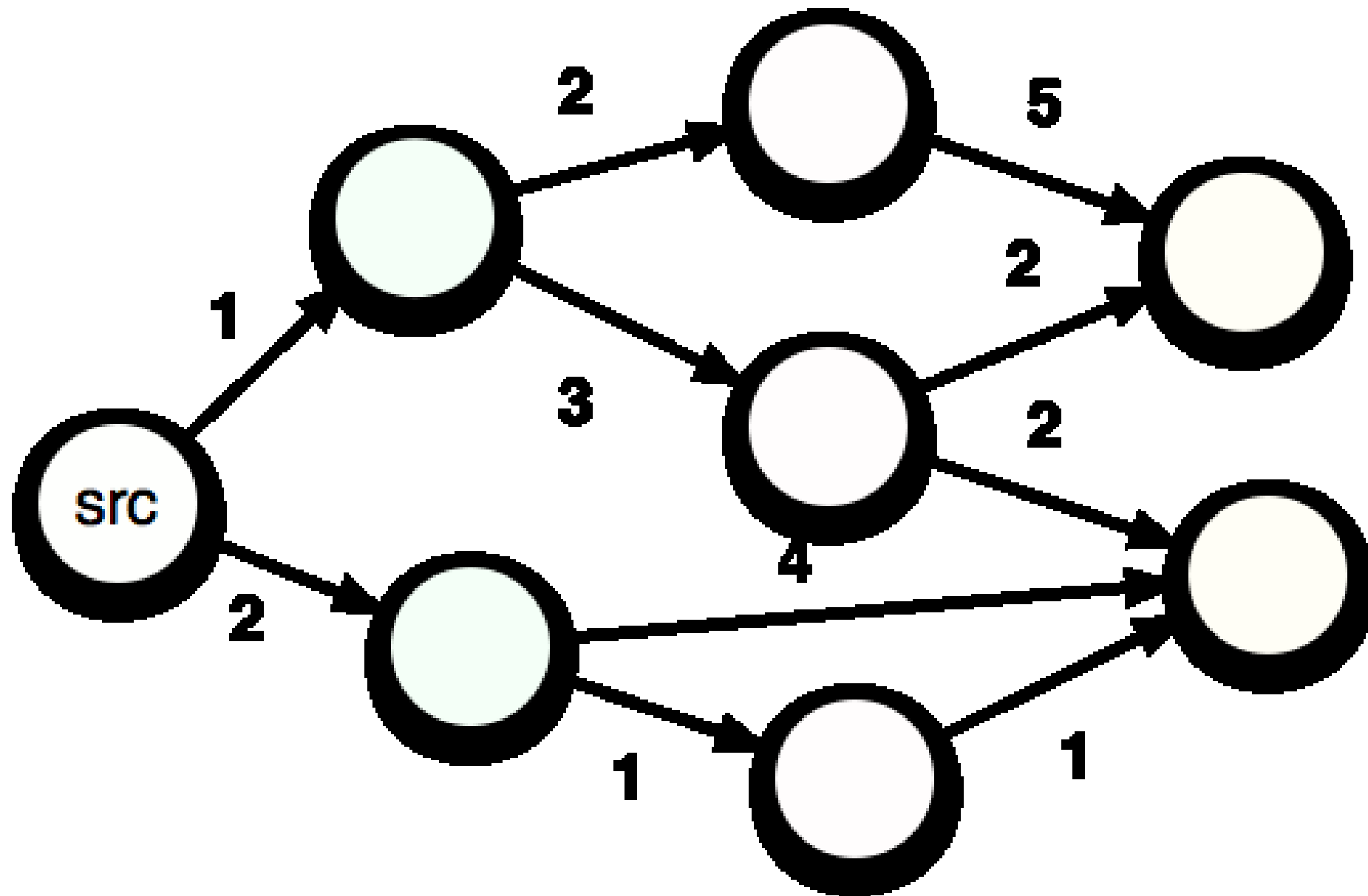
Designed for iterations

- **Stateful** (in-memory)
- **Only** intermediate values (messages) sent
- Hits the **disk** at input, output, checkpoint
- **Can** go out-of-core

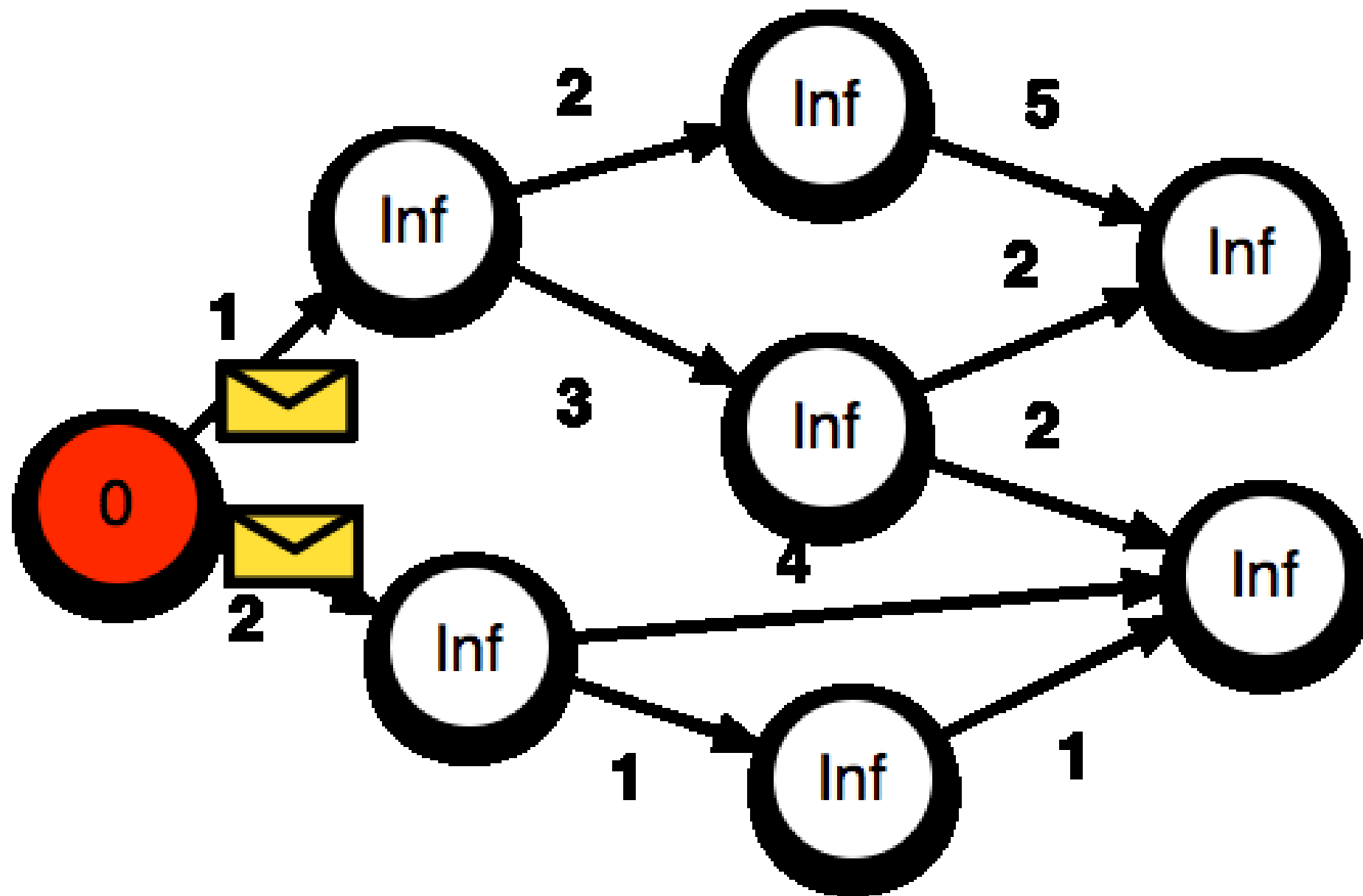
A bunch of other things

- **Combiners** (minimises messages)
- **Aggregators** (global aggregations)
- **MasterCompute** (executed on master)
- **WorkerContext** (executed per worker)
- **PartitionContext** (executed per partition)

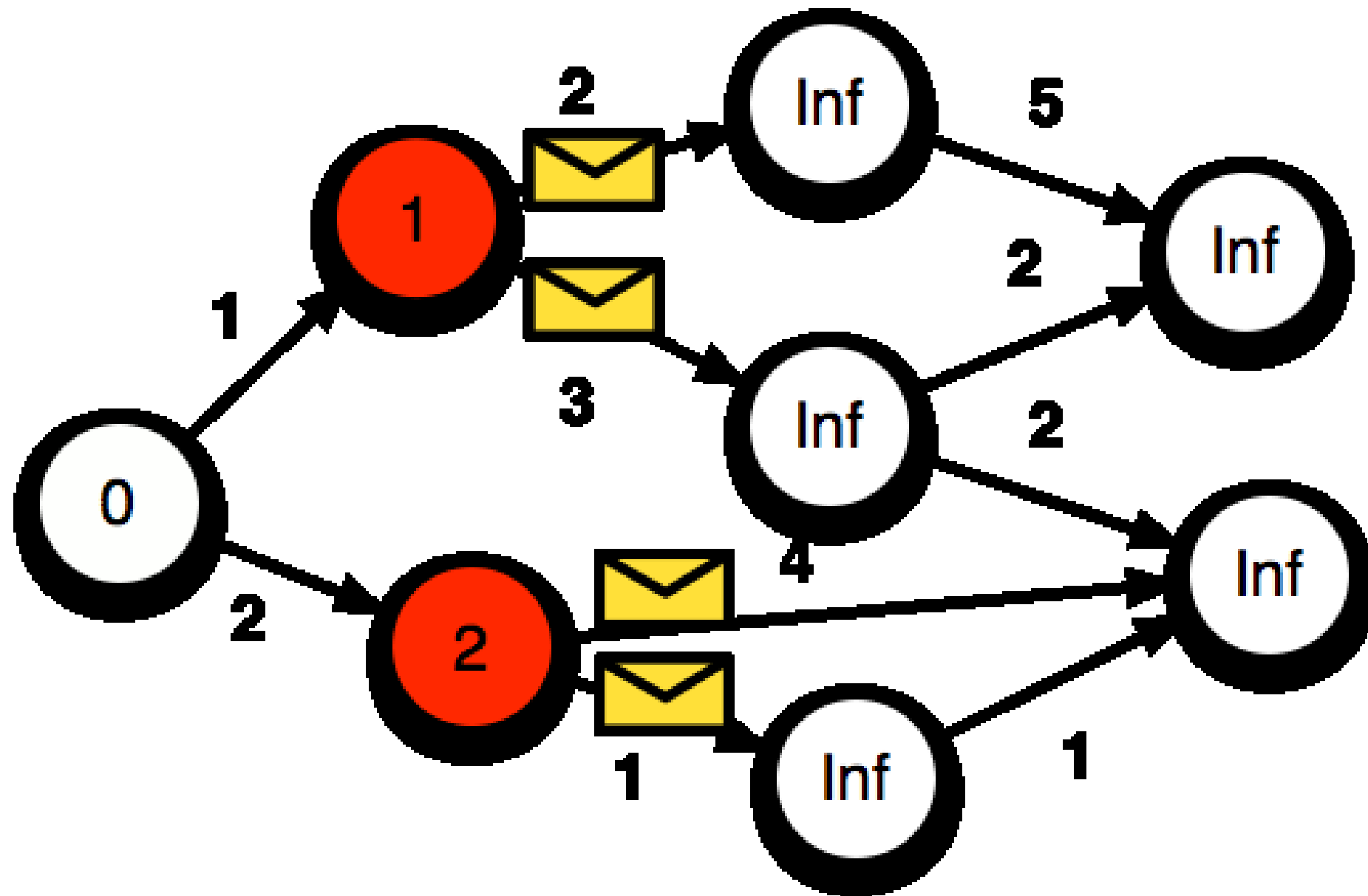
Shortest Paths



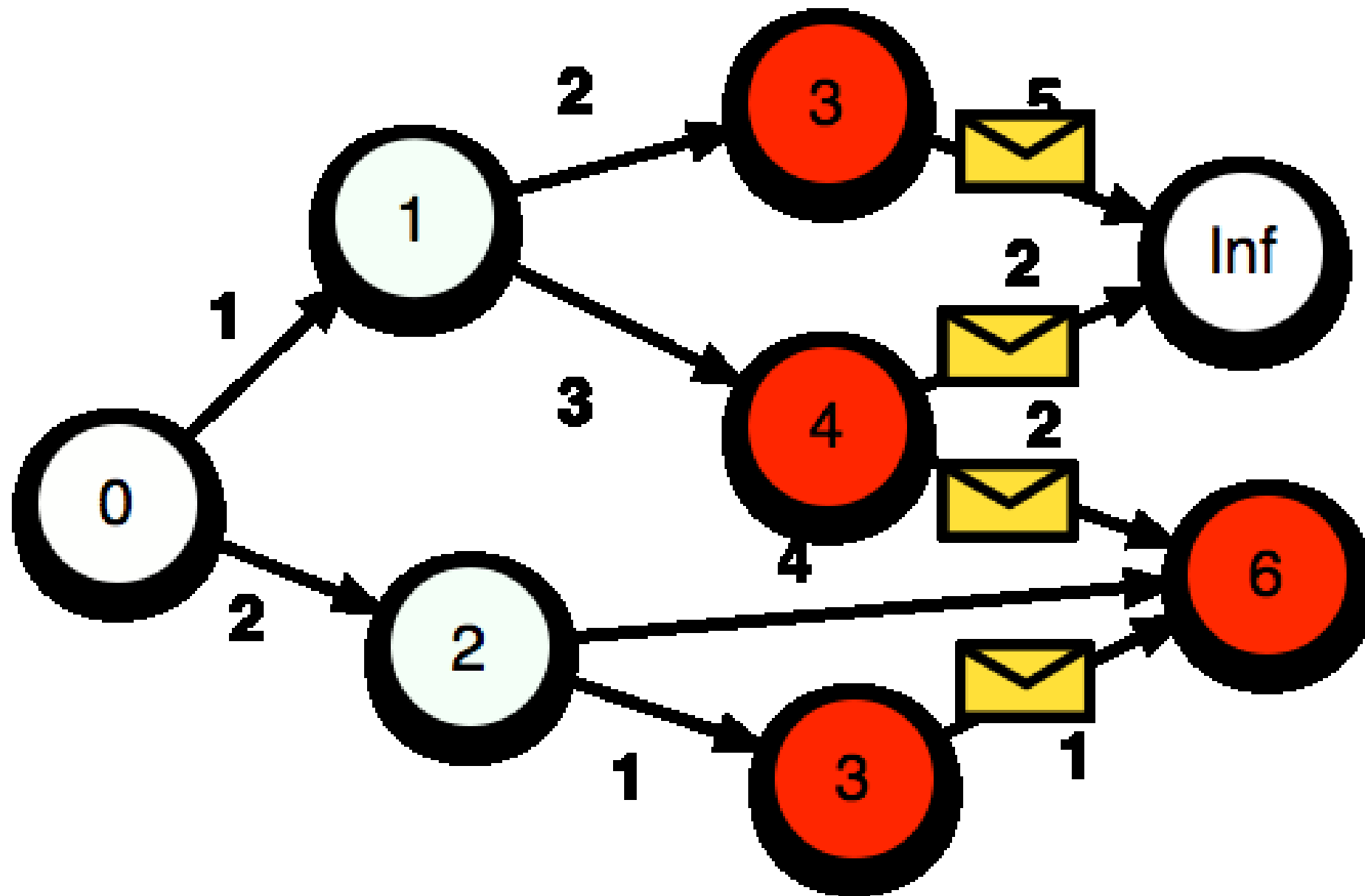
Shortest Paths



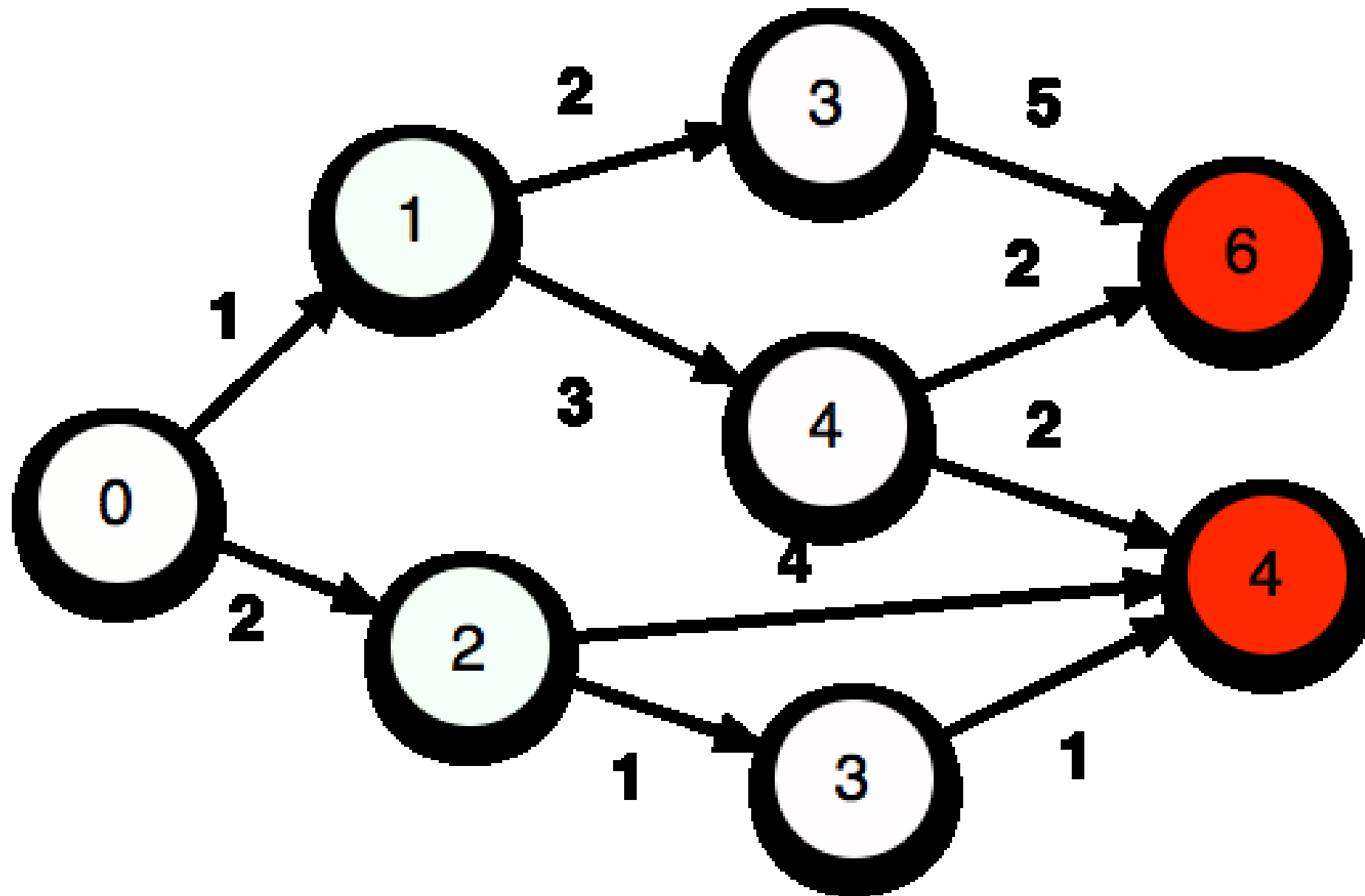
Shortest Paths



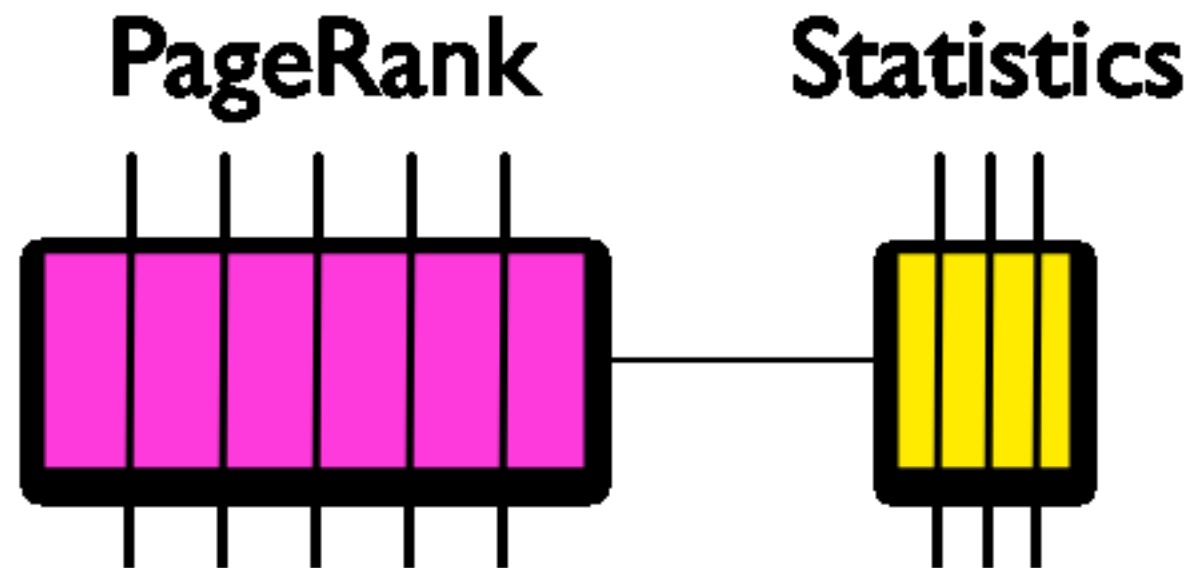
Shortest Paths



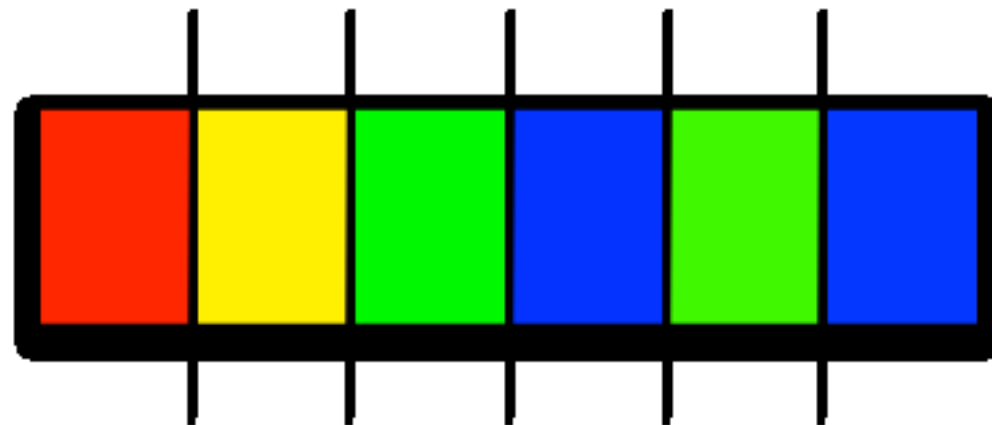
Shortest Paths



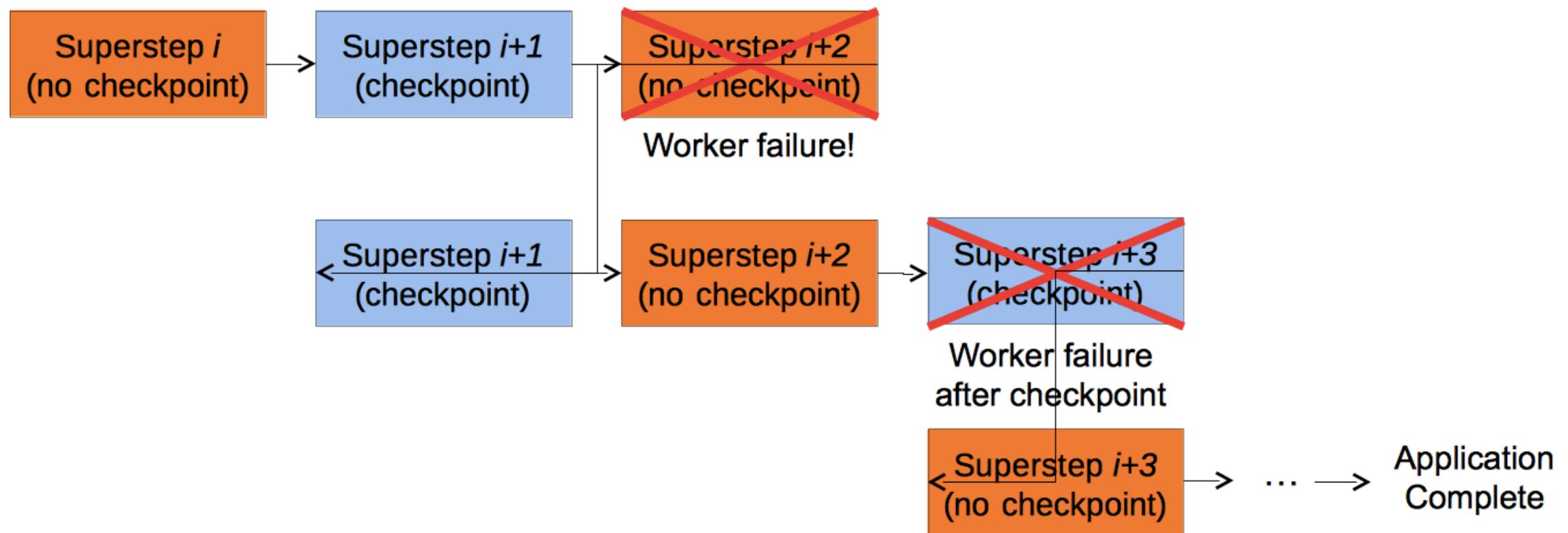
Composable API



Multi-phase algorithm

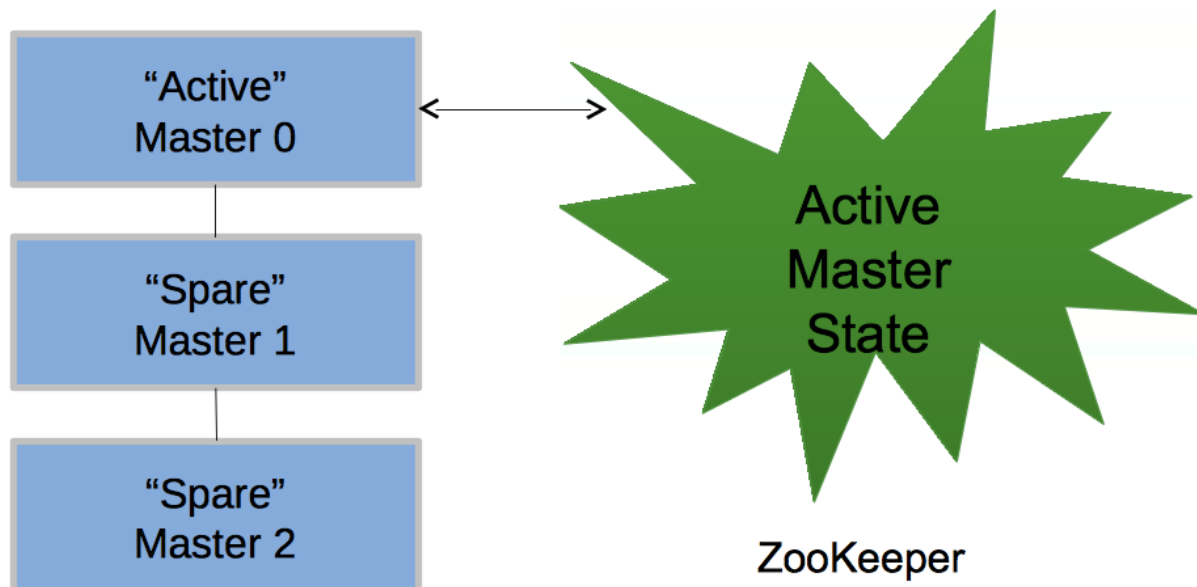


Checkpointing

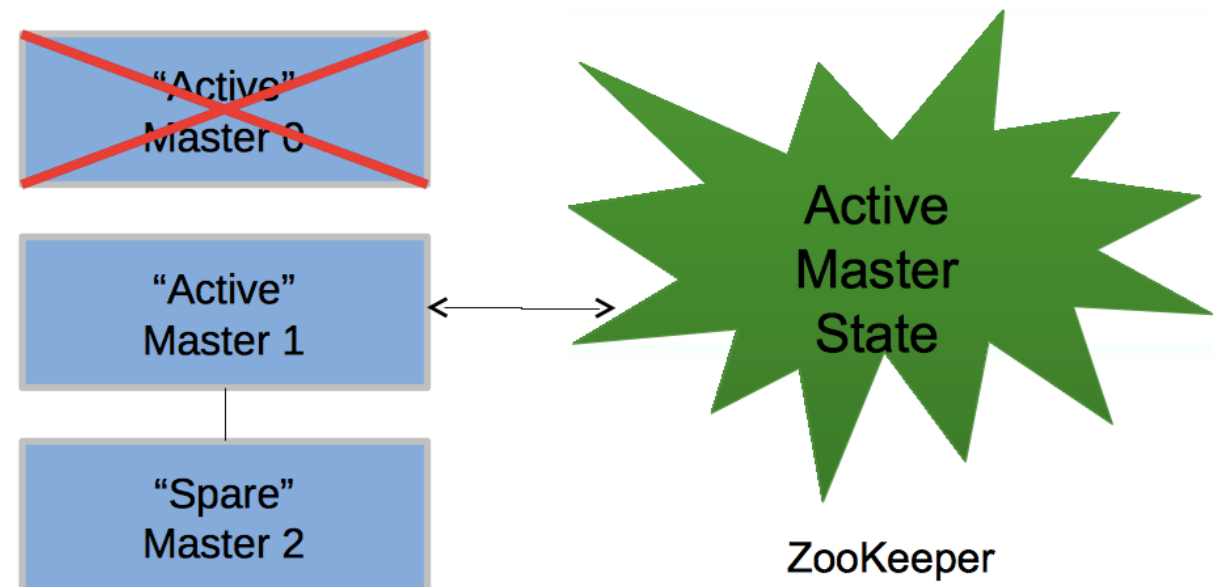


No SPoFs

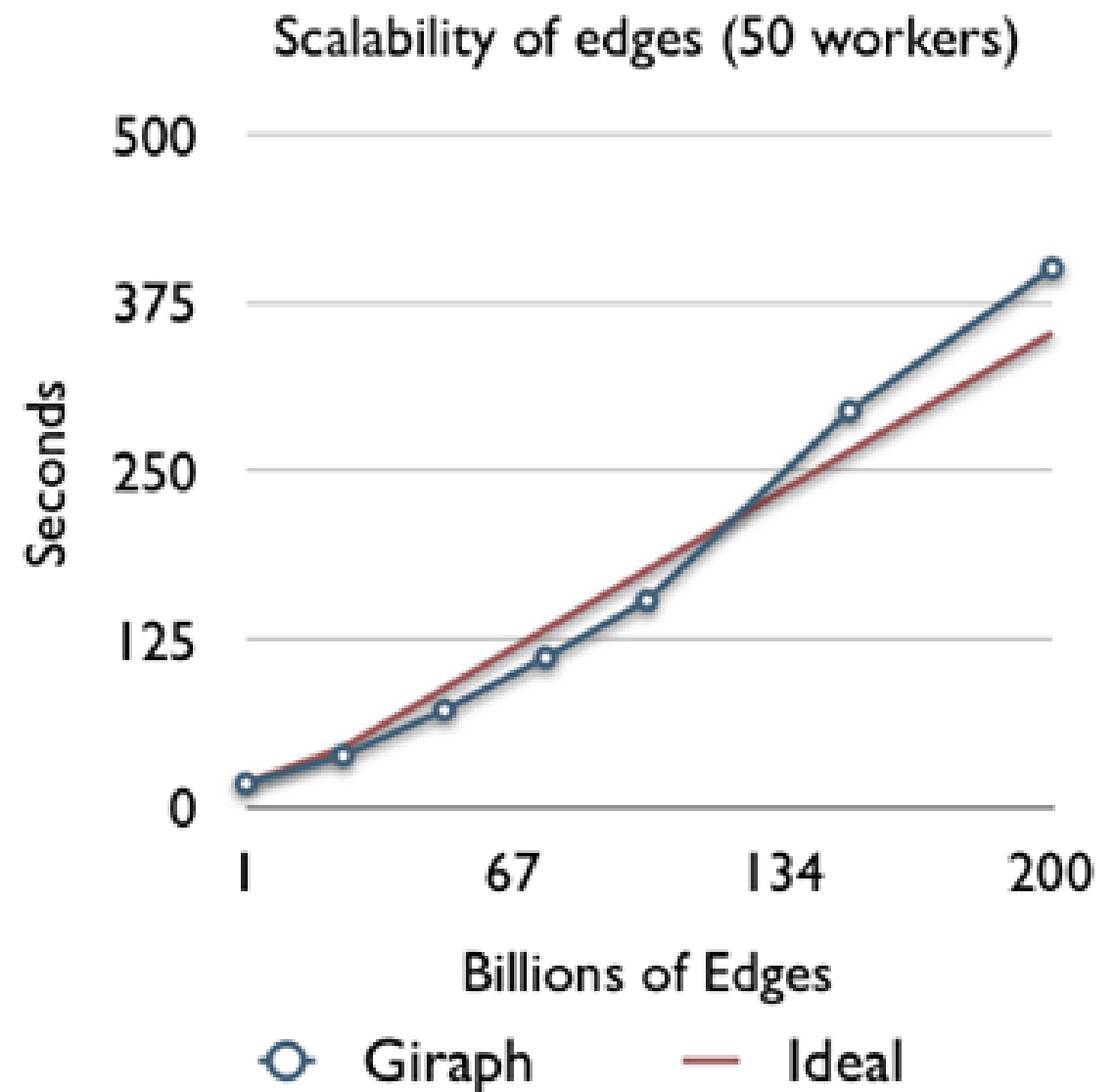
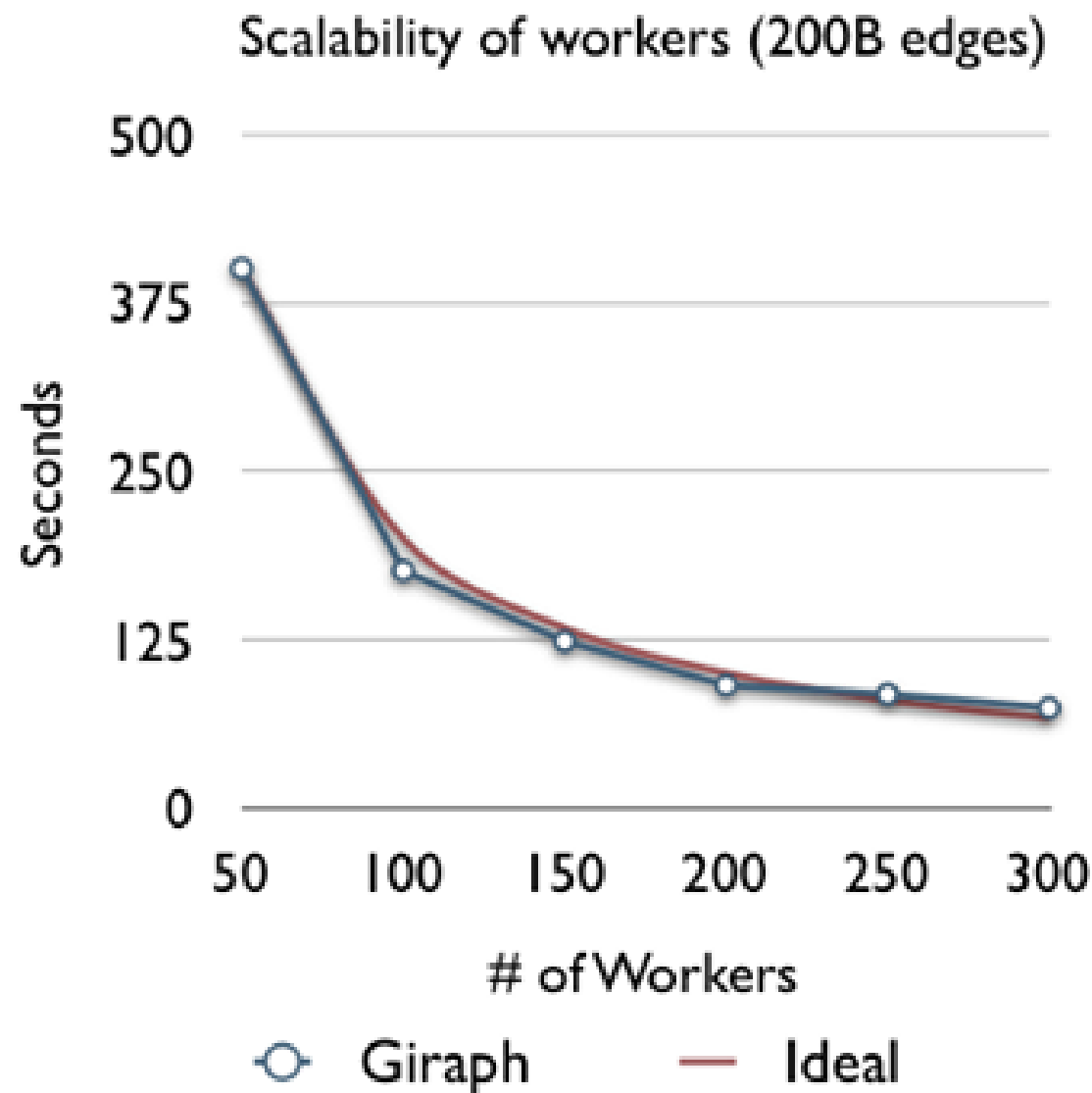
Before failure of active master 0



After failure of active master 0



Giraph **scales**



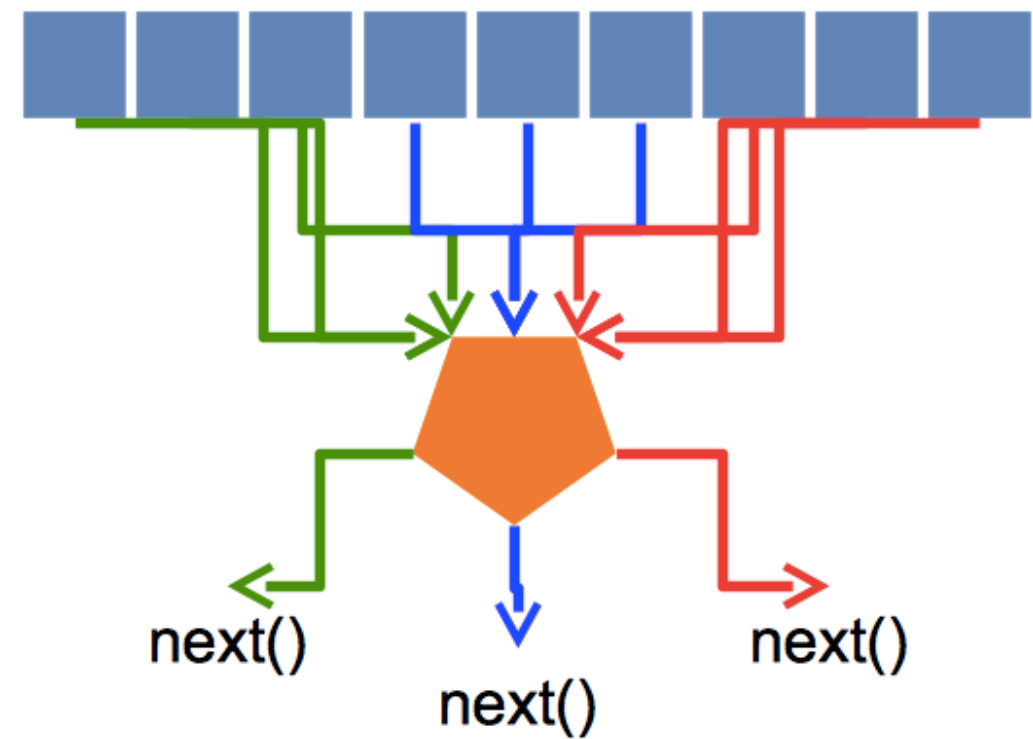
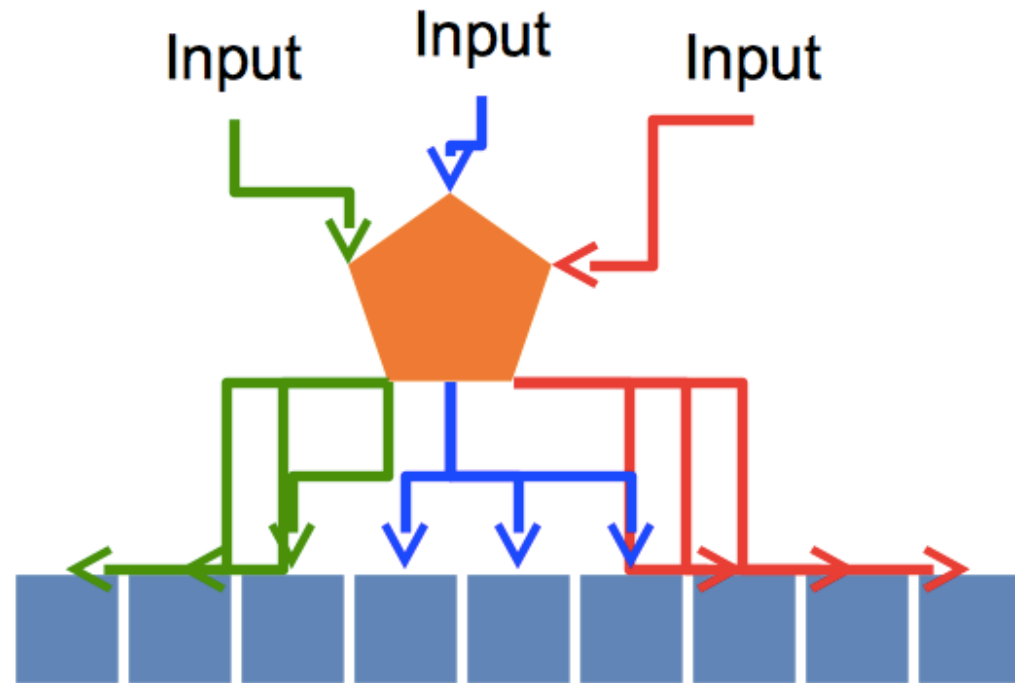
ref: <https://www.facebook.com/notes/facebook-engineering/scaling-apache-giraph-to-a-trillion-edges/10151617006153920>

Giraph is fast

- 100x over MR (Pr)
- jobs run within minutes
- given you have resources
;-)



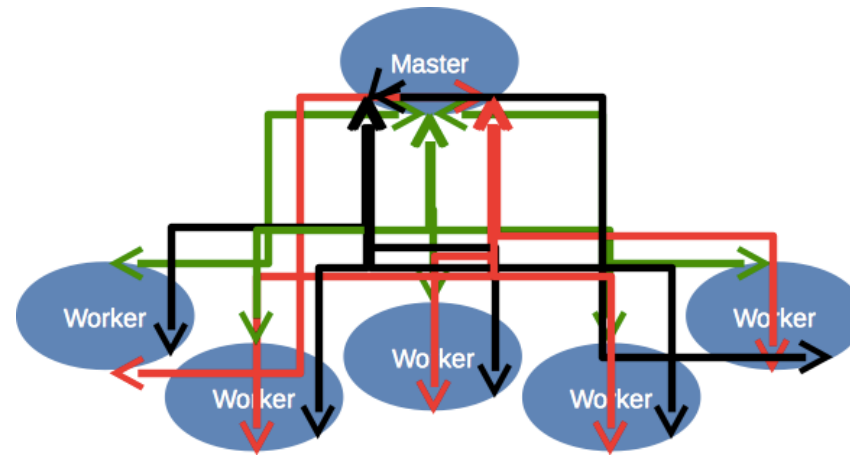
Serialised objects



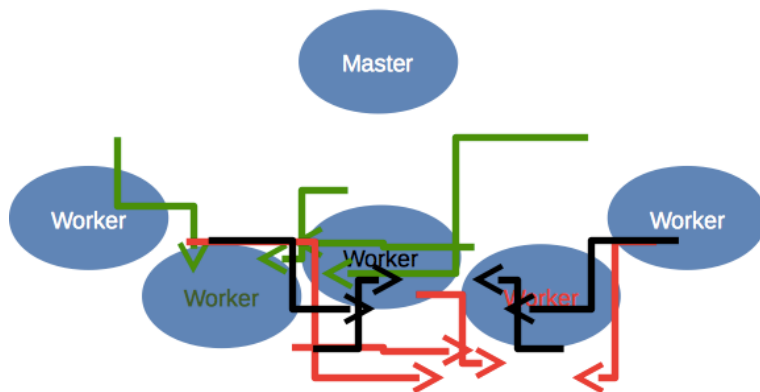
Primitive types

- Autoboxing is expensive
- Objects overhead (JVM)
- Use primitive types on your own
- Use primitive types-based libs (e.g. fastutils)

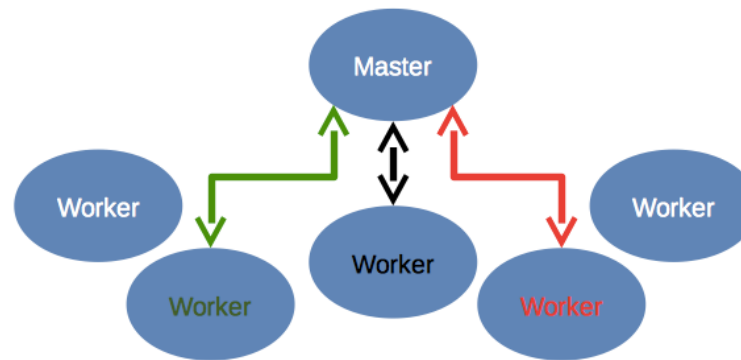
Sharded aggregators



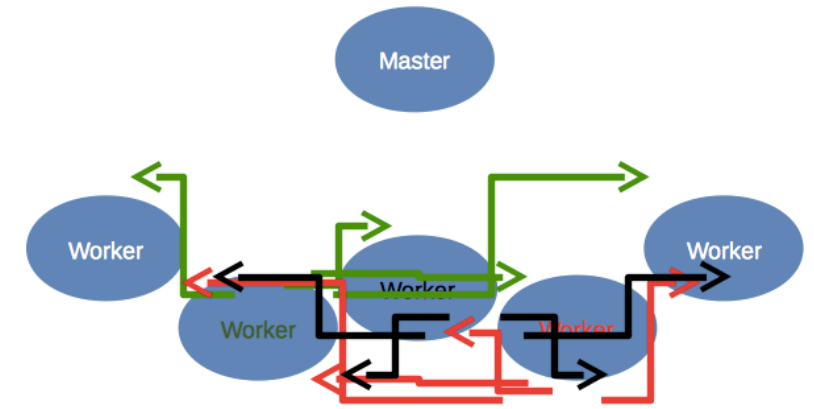
Workers own aggregators



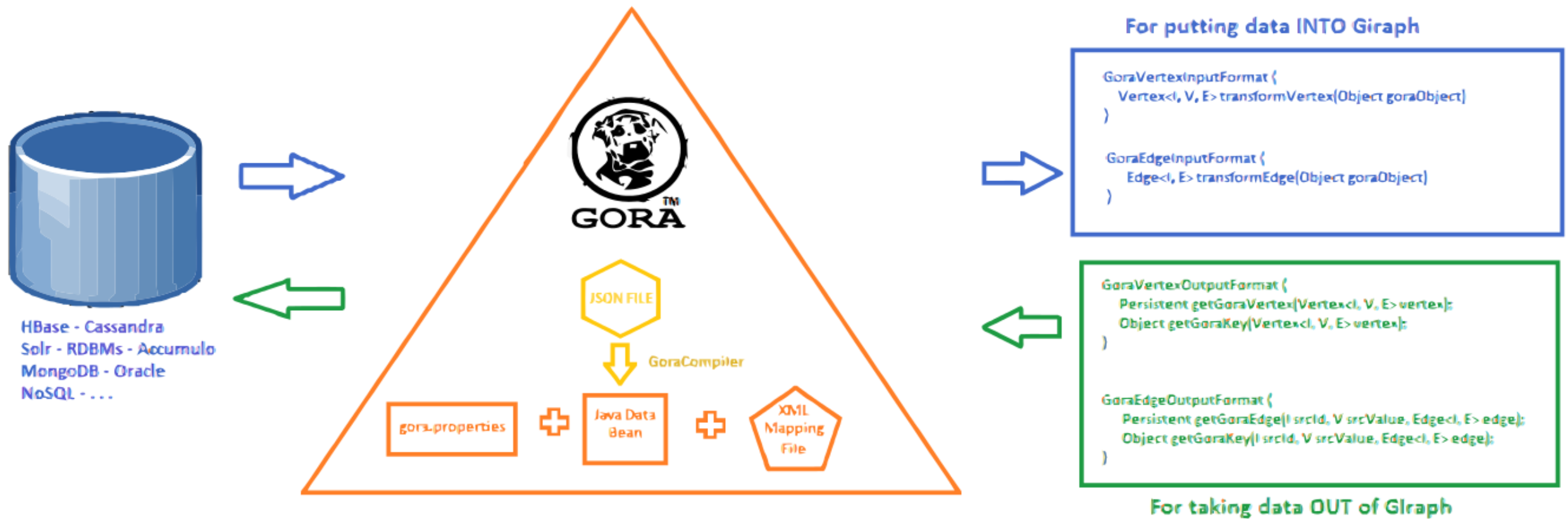
Aggregator owners communicate with Master



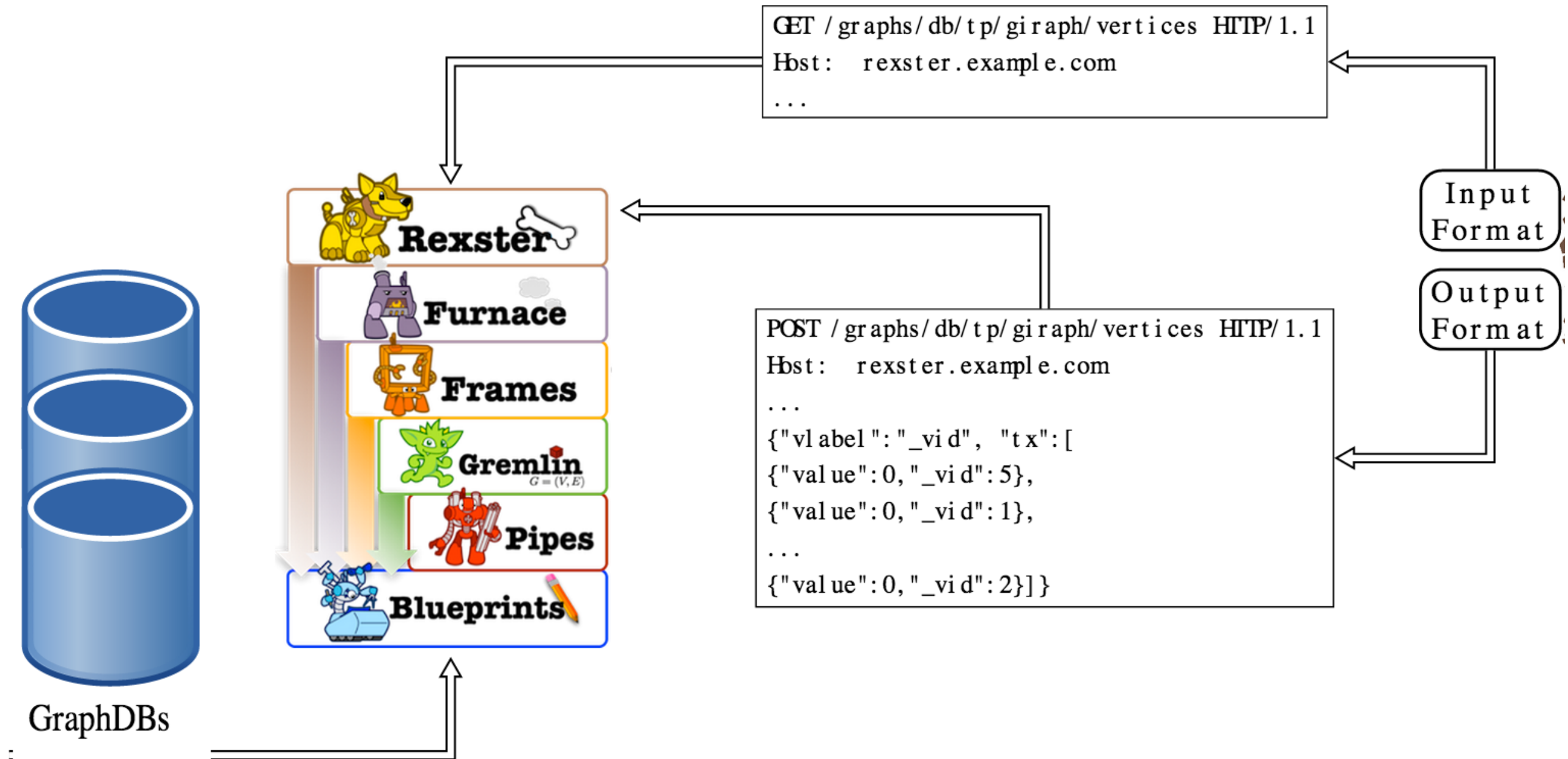
Aggregator owners distribute values



Many stores with Gora



And graph databases



Current and next steps

- Out-of-core graph and messages
- Jython interface
- Remove Writable from $\langle I V E M \rangle$
- Partitioned supernodes
- More documentation

Giraph in Action



- Published by **Manning**
- MEAP **now**
- Complete **Q3 2014** (well...)
- **Part 1**: Graphs and Algorithms
- **Part 2**: Giraph Basic Topics
- **Part 3**: Giraph Advanced Topics
- <http://www.manning.com/martella>

Okapi

- Apache **Mahout** for graphs
- Graph-based **recommenders**: ALS, SGD, SVD++, etc.
- Graph **analytics**: Graph partitioning, Community Detection, K-Core, etc.



Thank you

<http://giraph.apache.org>

Claudio Martella <claudio@apache.org>
@claudiomartella

Some figures gently borrowed from Nitay Joffe:
<http://www.slideshare.net/nitayj/20130910-giraph-at-london-hadoop-users-group>