

Graphalytics = From Benchmarking to Performance Engineering, leading to Massivizing Graph-Processing Systems



COMMIT/



Tim Hegeman, Wing-Lung Ngai, and Stijn Heldens.



Presentation developed jointly with Ana Lucia Varbanescu.

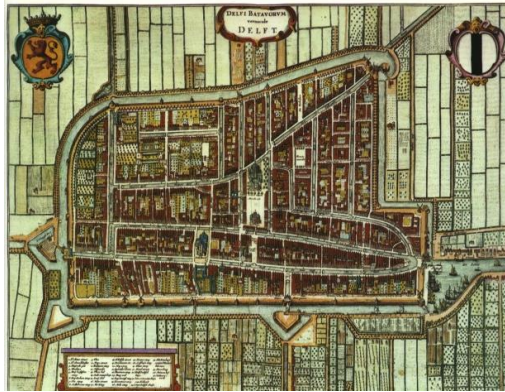


Several slides developed jointly with Yong Guo.



dr. ir. Alexandru Iosup
Distributed Systems Group

(TU) Delft – the Netherlands – Europe



founded 13th century
pop: 100,000



founded 1842
pop: 15,000



pop: 16.5 M



Barcelona

Graphs Are at the Core of Our Society: The LinkedIn Example

The State of LinkedIn



A very good resource for matchmaking workforce and prospective employers

Vital for your company's life,
as your Head of HR would tell you

Vital for the prospective employees

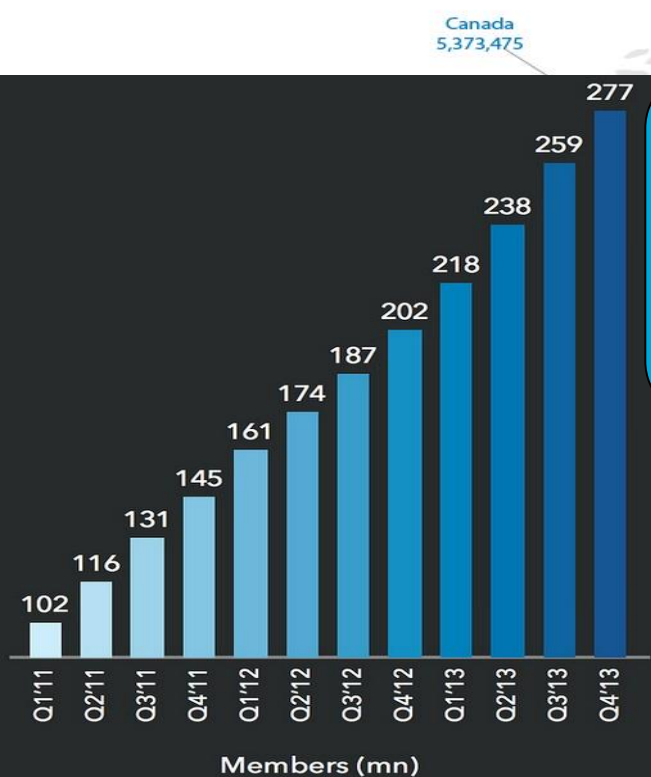
Tens of “specialized LinkedIns”: medical, mil, edu, gov, ...

LinkedIn's Service/Ops Analytics

The State of LinkedIn

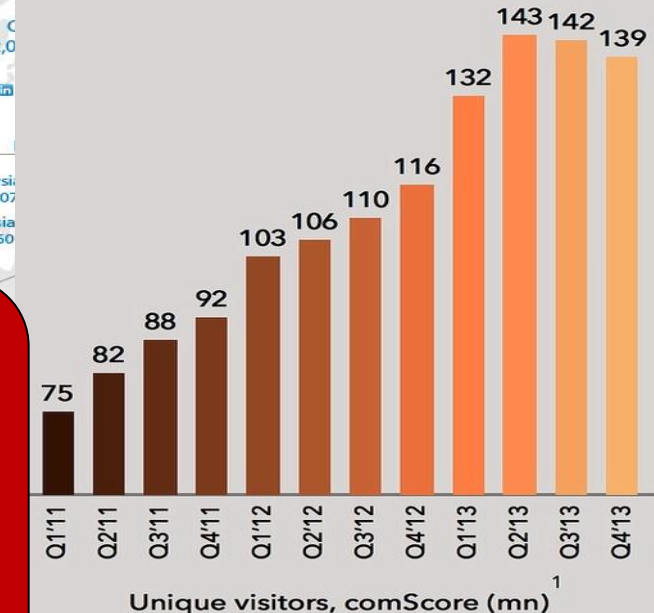
3-4 new users every second

but fewer visitors (and page views)



By processing the graph:
opinion mining, hub detection,
etc. Always new questions
about whole dataset.

100+ million questions of
customer retention,
of (lost) customer influence, of ...
Plus new hypotheses to test on
whole dataset.



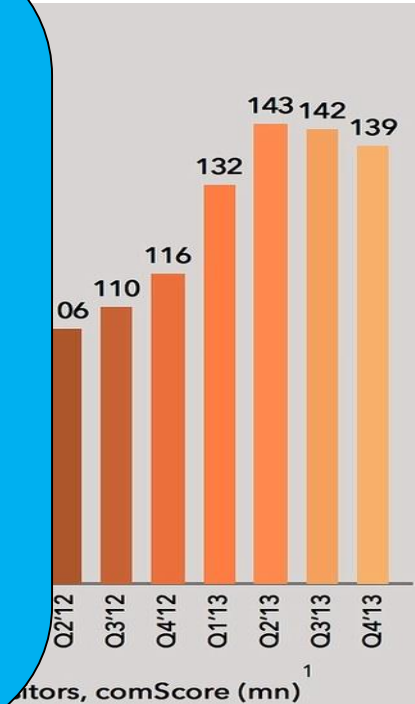
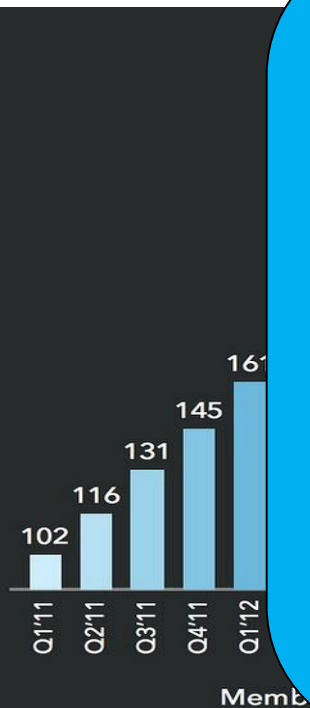
Why Analytics?

The State of LinkedIn

3-4 new users every second

but fewer visitors (and page views)

Periodic and/or continuous
full-graph **analytics**



How to do Analytics? Graph Processing @large

Linked 



A Graph Processing Platform



friendster 

 XFIRE™

6

Interactive processing not considered in this presentation.
Streaming not considered in this presentation.

Graph-processing is at the core of our society

The Data Deluge vs. Analytics

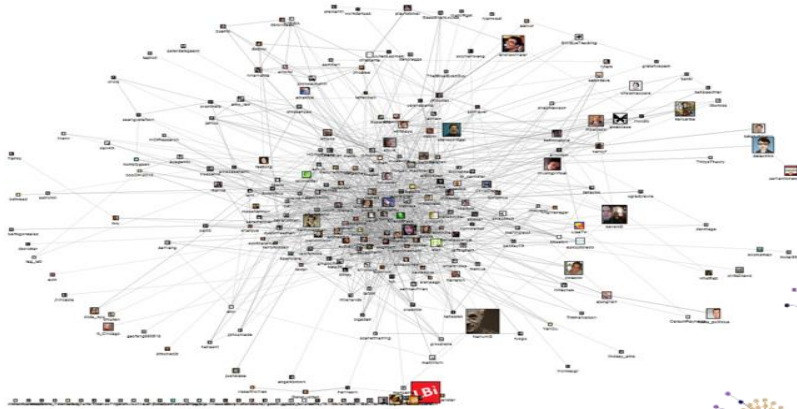
Graph processing @large

Which to select? What to tune? What to re-design?

A performance comparison of graph-processing systems

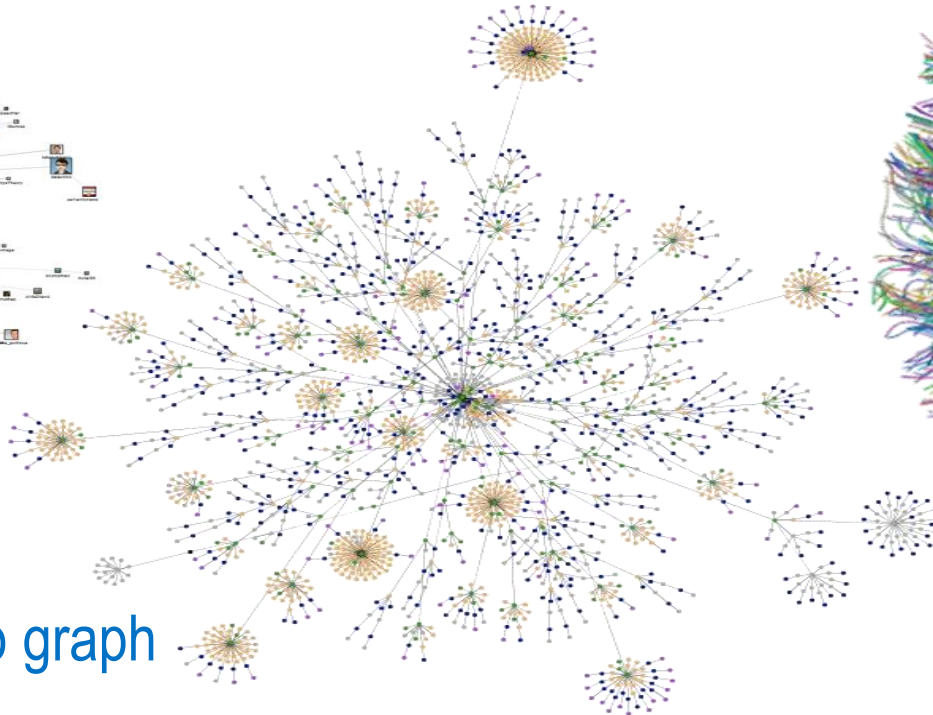
Take-home message

The data deluge: large-scale graphs



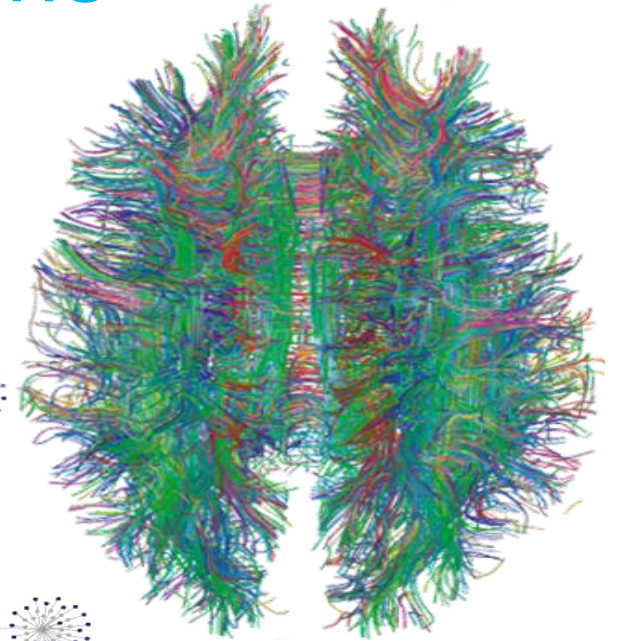
Social network

~1 billion vertices
~100 billion connections



Web graph

~50 billion pages
~1 trillion hyperlinks



Brain network

~100 billion neurons
~100 trillion connections

The data deluge: graphs everywhere!

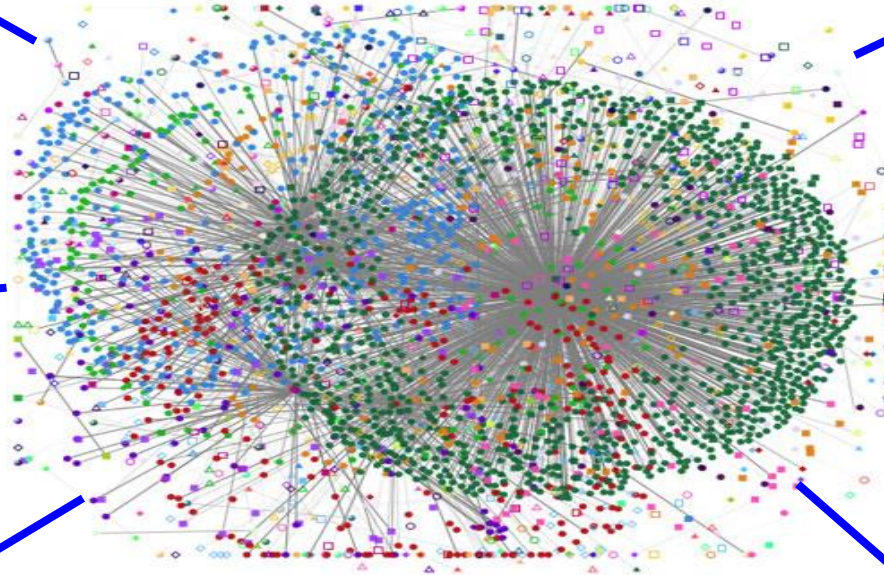
Linked in

400M users

??? edges

YAHOO!

friendster



270M MAU

200+ avg followers

>54B edges



1.2B MAU 0.8B DAU

200+ avg followers

>240B edges



9

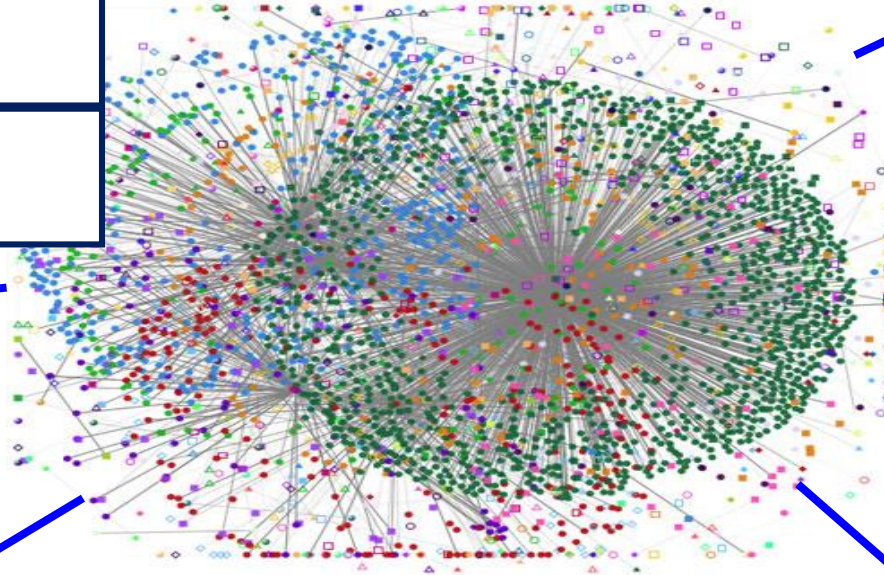
The data deluge: graphs everywhere!

LinkedIn

ORACLE Oracle 1.2M followers,
132k employees
company/day:
40-60 posts, 500-700 comments

YAHOO!

friendster



270M MAU
200+ avg followers
>54B edges



1.2B MAU 0.8B DAU
200+ avg followers
>240B edges



10

The data deluge vs. Analytics

LinkedIn

Oracle 1.2M followers



270M MAU

Data-intensive workload
10x graph size → 100x—1,000x slower

YAHOO!



1.2B MAU 0.8B DAU
200+ avg followers

>240B edges

friendster



The data deluge vs. Analytics

Linked in



Oracle 1.2M followers

270M MAIL

Data-intensive workload
10x graph size → 100x—1,000x slower

Compute-intensive workload
more complex analysis → ?x slower

friendster

>240B edges



12

The data deluge vs. Analytics

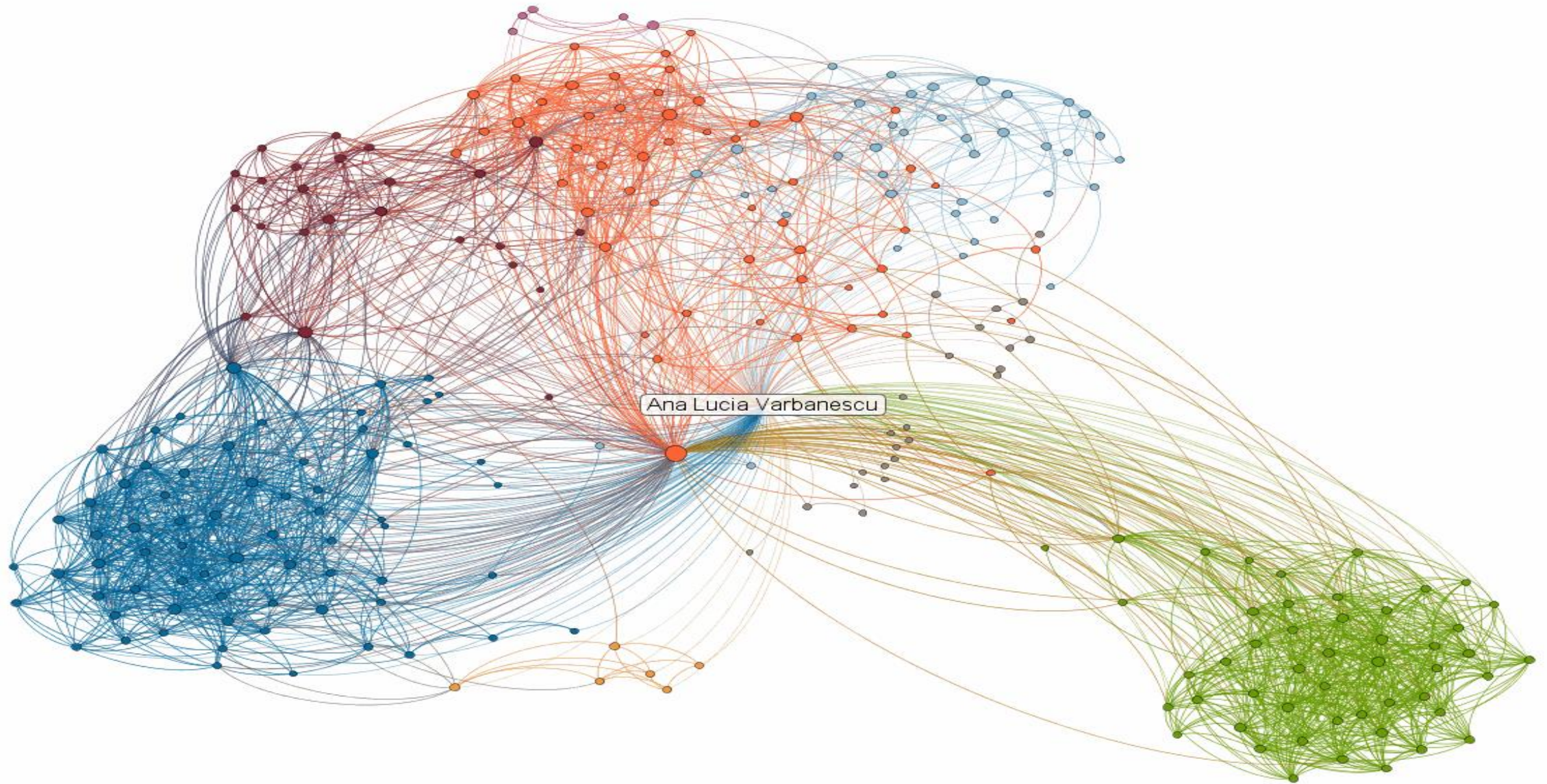
Linked in



Data-intensive workload
10x graph size → 100x—1,000x slower

Compute-intensive workload
more complex analysis → ?x slower

Dataset-dependent workload
unfriendly graphs → ??x slower



Your network is so large...

Sorry, but your network is too large to be computed, we are working to increase the limit, stay tuned!

The “sorry, but...” moment

Supporting multiple users

10x number of users → ???x slower

What would **you** do to solve this?

Data-intensive workload

10x graph size → 100x—1,000x slower

Compute-intensive workload

more complex analysis → ?x slower

Dataset-dependent workload

unfriendly graphs → ??x slower

Supporting multiple users

10x number of users → ???x slower

Graph-processing is at the core of our society
The data deluge vs. Analytics

Graph Processing @Large

Which to select? What to tune? What to re-design?
A performance comparison of graph-processing systems
Take-home message

Graph Processing @large

Linked **in**



A Graph Processing Platform

**Ideally,
N cores/disks \rightarrow Nx
faster**

(partitioning, compression,
replication, caching)

Distribution
to processing
platform

**Ideally, Parallel/
Distributed/
Heterogeneous**

**Ideally,
N cores/disks \rightarrow Nx
faster**

friendster 



19

Interactive processing not considered in this presentation.
Streaming not considered in this presentation.

Graph processing systems

Performance ↑

- Specify application
- Choose the hardware
- Implement & optimize
- Think Graph500 performers

Dedicated
Systems

Custom
Systems

Generic
Systems

- Systems for graph processing
- Separate users from backends
- Think Giraph

- Use existing distributed systems
- Mapping is difficult
- Parallelism is “free”
- Think Hadoop/Spark

Development Effort →

System diversity (CPU-based)

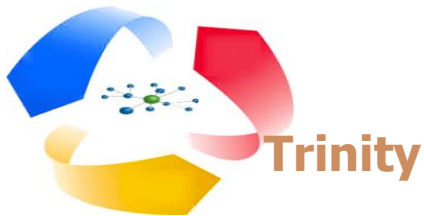
Dedicated Systems

Generic

ORACLE PGX

Intel Graphmat

Neo4j
the graph database



System diversity (GPU-enabled)

Dedicated Systems

Generic



medusa-gpu

Medusa: Simplified Graph Processing on GPUs



NetSysLab

mapgraph ^{Beta}

Massively Parallel Graph processing on GPUs

TOTEM

Gunrock




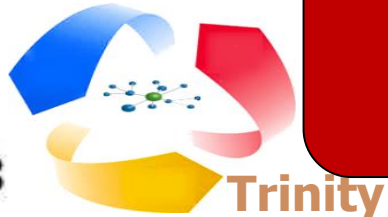








High-performance Graph Primitives on GPU



VertexAPI2

Graph-Processing Platforms

Platform: the combined hardware, software, and programming system that is being used to complete a graph processing task



Which to choose?

What to tune?

What to re-design?

23

Graph-processing is at the core of our society
The data deluge vs. Analytics
Graph processing @large

Graphalytics:

Which system to select?

What to tune? What to re-design?

A performance comparison of graph-processing systems
Take-home message

Benchmarking, but ... What Is a Benchmark?

Benchmark definition must include:

Data schema: data representation

Workloads: formalize datasets + algorithms

Performance metrics: from performance to non-traditional to cost-related

Execution rules: how to run the benchmark tests, parameter values, etc.

Desirable support for stakeholders:

Live addition of results

Curation of added results

Auditing results

What is the performance of graph-processing platforms?

Metrics
Diversity

Graph
Diversity

Algorithm
Diversity

- Graph500
 - Single application (BFS), Single class of synthetic datasets. @ISC16: future diversification.
- Few existing platform-centric comparative studies
 - Prove the superiority of a given system, limited set of metrics
- GreenGraph500, GraphBench, XGDBench
 - Issues with representativeness, systems covered, metrics, ...

What is the performance of graph-processing platforms?

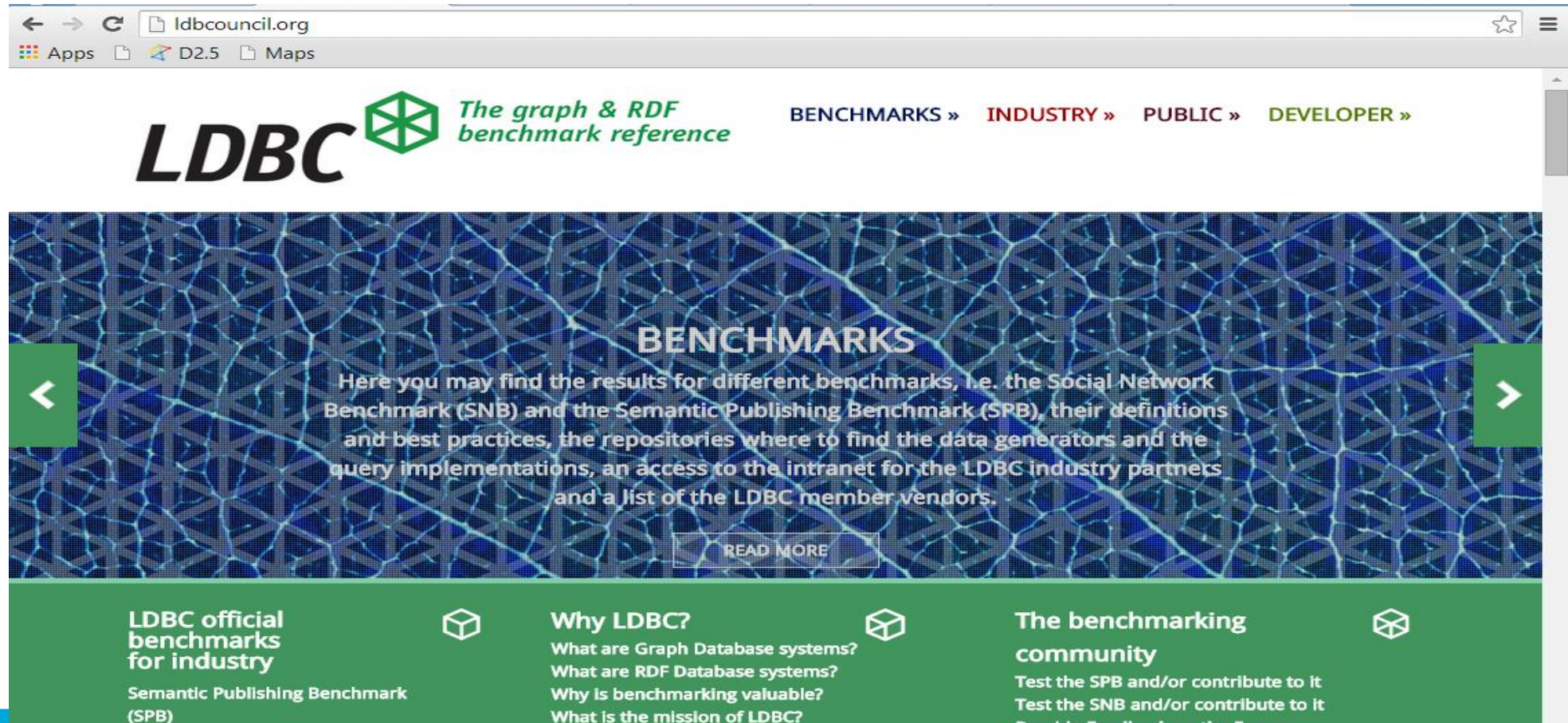
Metrics
Diversity

Graph
Diversity

Algorithm
Diversity

Graphalytics = comprehensive benchmarking suite for graph processing across many platforms

ldbcouncil.org



The screenshot shows the LDBC website homepage. At the top, the browser address bar displays 'ldbcouncil.org'. The LDBC logo, consisting of the text 'LDBC' and a green hexagonal icon, is followed by the tagline 'The graph & RDF benchmark reference'. Navigation links for 'BENCHMARKS', 'INDUSTRY', 'PUBLIC', and 'DEVELOPER' are visible. The main content area features a large banner with a blue and green geometric pattern. The banner title is 'BENCHMARKS', and the text describes the availability of benchmark results, definitions, best practices, data generators, query implementations, intranet access for industry partners, and a list of member vendors. A 'READ MORE' button is located at the bottom of the banner. Below the banner, there are three green boxes with white text and icons. The first box is titled 'LDBC official benchmarks for industry' and mentions the 'Semantic Publishing Benchmark (SPB)'. The second box is titled 'Why LDBC?' and lists questions about Graph Database systems, RDF Database systems, the value of benchmarking, and the mission of LDBC. The third box is titled 'The benchmarking community' and mentions testing the SPB and SNB, and contributing to the community.

LDBC The graph & RDF benchmark reference

BENCHMARKS » INDUSTRY » PUBLIC » DEVELOPER »

BENCHMARKS

Here you may find the results for different benchmarks, i.e. the Social Network Benchmark (SNB) and the Semantic Publishing Benchmark (SPB), their definitions and best practices, the repositories where to find the data generators and the query implementations, an access to the intranet for the LDBC industry partners and a list of the LDBC member vendors.

[READ MORE](#)

LDBC official benchmarks for industry

Semantic Publishing Benchmark (SPB)

Why LDBC?

What are Graph Database systems?
What are RDF Database systems?
Why is benchmarking valuable?
What is the mission of LDBC?

The benchmarking community

Test the SPB and/or contribute to it
Test the SNB and/or contribute to it
Provide Feedback on the Forum



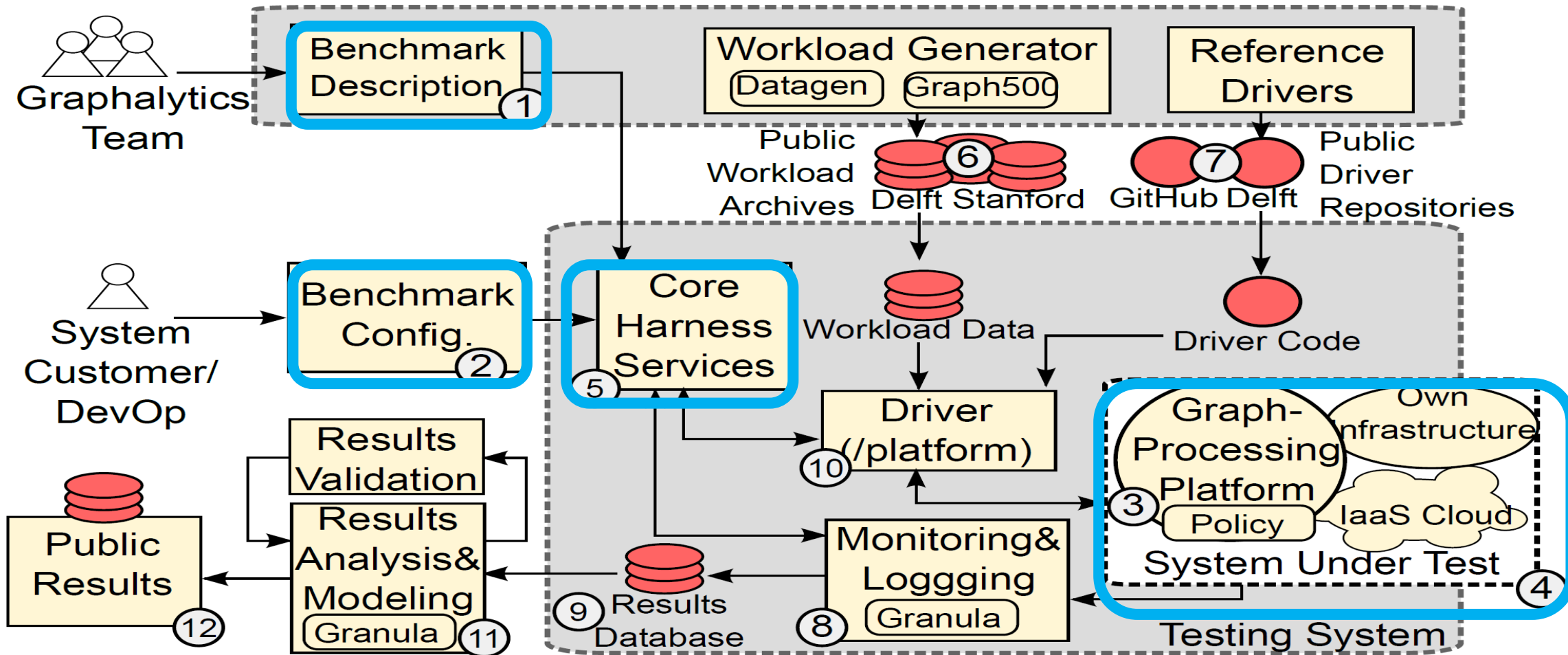
Graphalytics, in a nutshell

- An LDBC benchmark*
- Advanced benchmarking harness
- Many classes of algorithms used in practice
- Diverse real and synthetic datasets
- Diverse set of experiments representative for practice
- Granula for manual choke-point analysis
- Modern software engineering practices
- Supports many platforms
- Enables comparison of community-driven and industrial systems



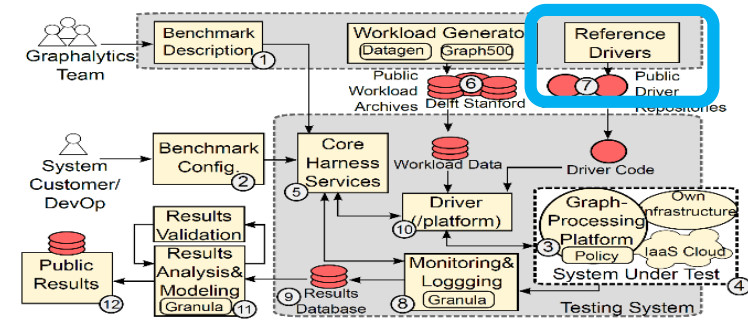
<http://graphalytics.ewi.tudelft.nl>
<https://github.com/tudelft-atlarge/graphalytics/>

Benchmarking Harness



Iosup et al. LDBC Graphalytics: A Benchmark for Large Scale Graph Analysis on Parallel and Distributed Platform, VLDB'16.

Graphalytics = Representative Classes of Algorithms and Datasets



- 2-stage selection process of **algorithms** and datasets

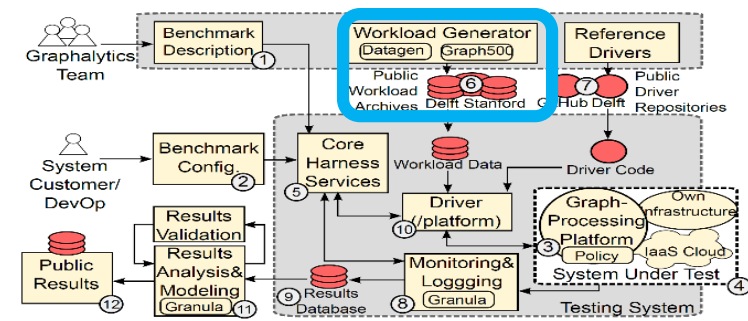
Class	Examples	%
Graph Statistics	Diameter, Local Clust. Coeff, PageRank	20
Graph Traversal	BFS, SSSP, DFS	50
Connected Comp.	Reachability, BiCC, Weakly CC	10
Community Detection	Clustering, Nearest Neighbor, Community Detection w Label Propagation	5
Other	Sampling, Partitioning	<15

+ weighted graphs: Single-Source Shortest Paths (~35%)

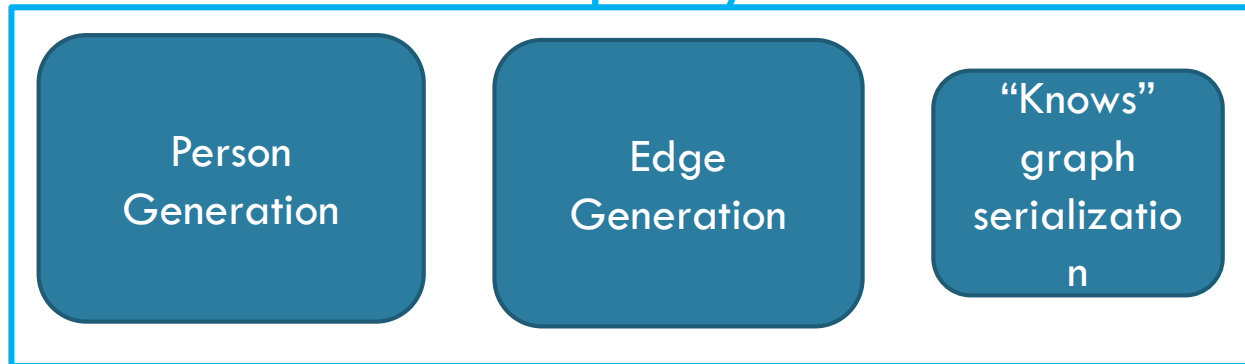
Graphalytics = Distributed Graph Generation w DATAGEN



- Rich set of configurations
- More diverse degree distribution than Graph500
- Realistic properties, e.g., clustering coefficient and assortativity



Graphalytics

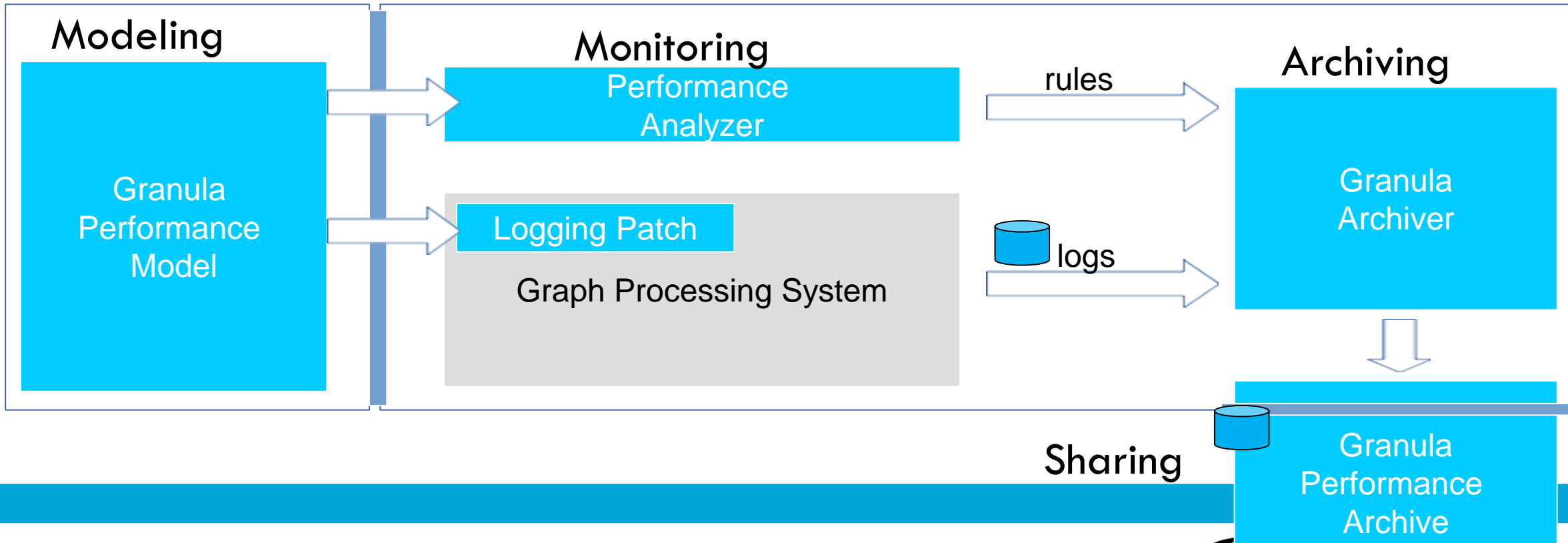
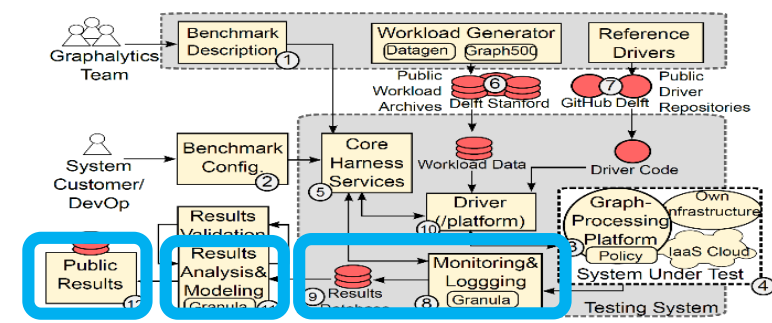


Level of
Detail

Graphalytics = Diverse Automated Experiments

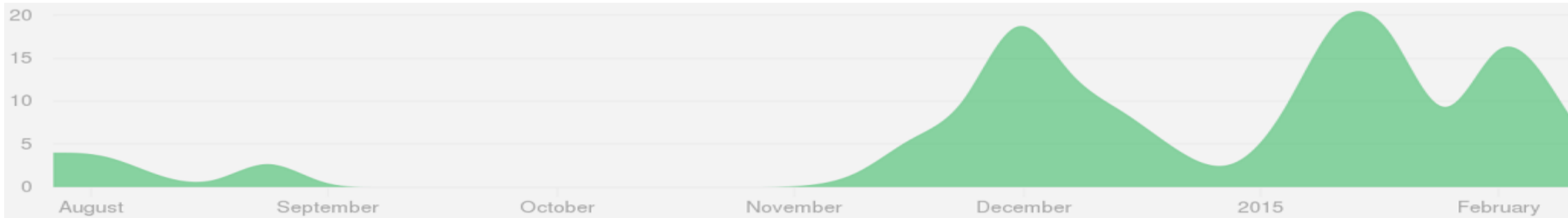
Category	Experiment	Algo.	Data	Nodes/ Threads	Metrics
Baseline	Dataset variety	BFS,PR	All	1	Run, norm.
	Algorithm variety	All	R4(S), D300(L)	1	Runtime
Scalability	Vertical vs. horiz.	BFS, PR	D300(L), D1000(XL)	1—16/1—32	Runtime, S
	Weak vs. strong	BFS, PR	G22(S)— G26(XL)	1—16/1—32	Runtime, S
Robustness	Stress test	BFS	All	1	SLA met
	Variability	BFS	D300(L), D1000(L)	1/16	CV
Self-Test	Time to run/part	--	Datagen	1—16	Runtime

Graphalytics = Portable Performance Analysis w Granula



Minimal code invasion + automated data collection at runtime
+ portable archive (+ web UI) → portable bottleneck analysis

Graphalytics = Modern Software Engineering Process



Graphalytics code reviews

Internal release to LDBC partners (first, Feb 2015; last, Feb 2016)

Public release, announced first through LDBC (Apr 2015)

First full benchmark specification, LDBC criteria (Q1 2016)

Jenkins continuous integration server

SonarQube software quality analyzer



Implementation status

G=validated, on GitHub
V=validation stage

	MR 2	Gi-raph	Graph X	Power Graph	Graph Lab	Neo4j	PGX. D	Graph Mat	Open G	TOTEM	Map Graph	Me du sa
LCC	G	G	G	G	G	G	--	G	G	--	--	--
BFS	G	G	G	G	G	G	G	G	G	V	V	V
WCC	G	G	G	G	G	G	G	G	G	V	V	V
CDLP	G	G	G	G	G	G	G	G	G	--	--	--
P'Ran k	--	G	G	G	V	--	G	G	G	V	V	V
SSSP	--	G	G	G	--	--	G	G	G	--	--	--

<https://github.com/tudelft-atlarge/graphalytics/>

Implementation status

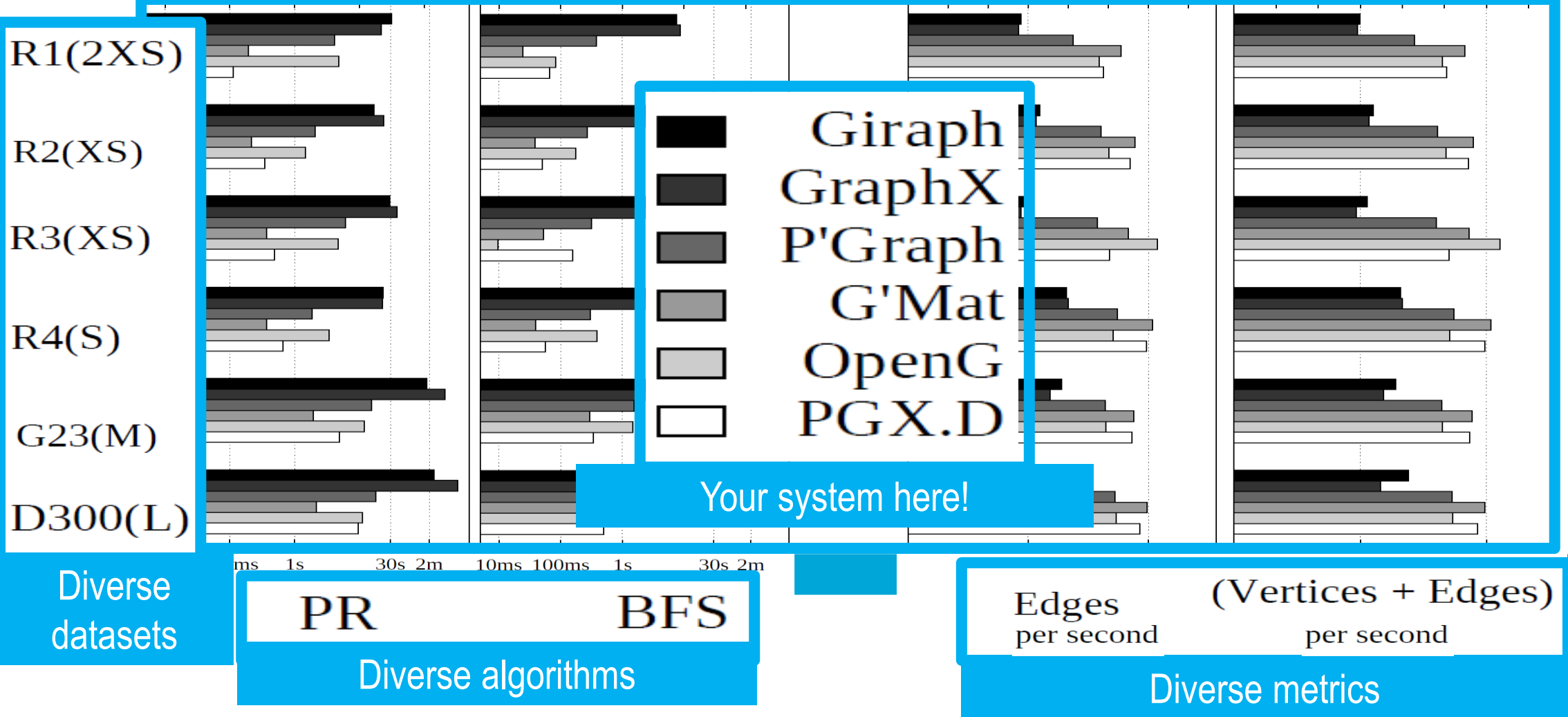
G=validated, on GitHub
V=validation stage

	MR 2	Gi-raph	Graph X	Power Graph	Graph Lab	Neo4j	PGX. D	Graph Mat	Open G	TOTEM	Map Graph	Me du sa
LCC	G	G	G	G	G	G	--	G	G	--	--	--
BFS	G	G	G	G	G	G	G	G	G	V	V	V
WCC	G	G	G	G	G	G	G	G	G	V	V	V
CDLP	G	G	G	G	G	G	G	G	G	--	--	--
P'Ran k	--	G	G	G	V	--	G	G	G	V	V	V
SSSP	--	G	G	G	--	--	G	G	G	--	--	--

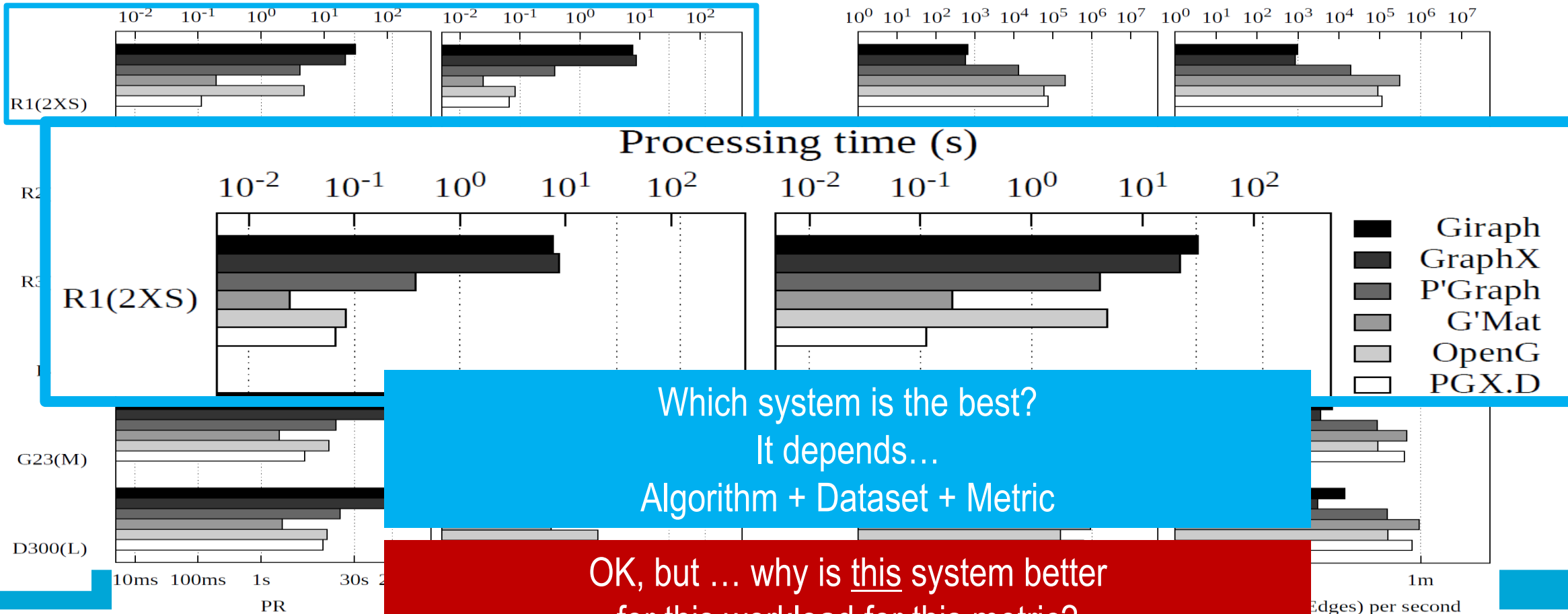
Benchmarking and tuning performed by vendors

Graphalytics Capabilities: An Example

Graphalytics enables deep comparison of many systems at once, through diverse experiments and metrics



Processing time (s) + Edges[+Vertices]/s



Graph-processing is at the core of our society
The data deluge vs. Analytics
Graph processing @large

Graphalytics:

Which system to select?

What to tune? What to re-design?

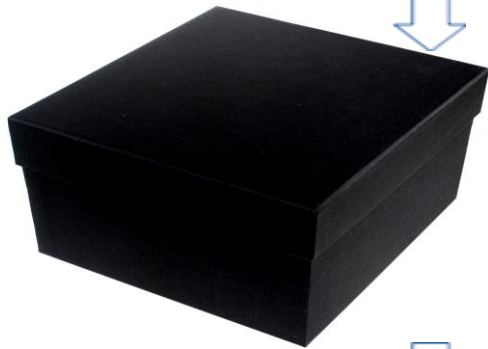
A performance comparison of graph-processing systems
Take-home message

Coarse-grained vs Fine-grained Evaluation (1)

Coarse-grained Method

system viewed as a black-box

Algorithms, Datasets, Resources



Graph
processing
system

Coarse-grained metrics

(Overall Execution Time)

Fine-grained Method

system viewed as a white-box

Algorithms, Datasets, Resources

IO operations

Processing

Overheads



Fine-grained metrics

(Stage 3 time, straggler tasks)

Fine-grained evaluation method is more comprehensive

Coarse-grained vs Fine-grained Evaluation (2)

Abstract

Coarse-grained Method
knowledge at conceptual level

Graph
Processing
Systems

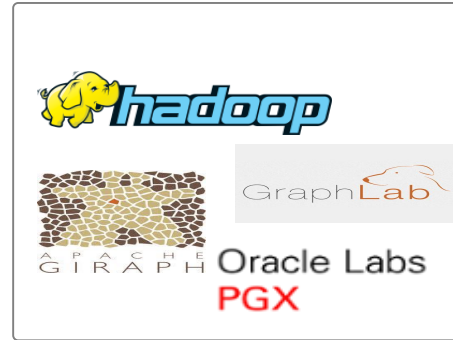
Distributed
Infrastructure



Few, coarse-grained results

Granular

Fine-grained Method
knowledge at technical level



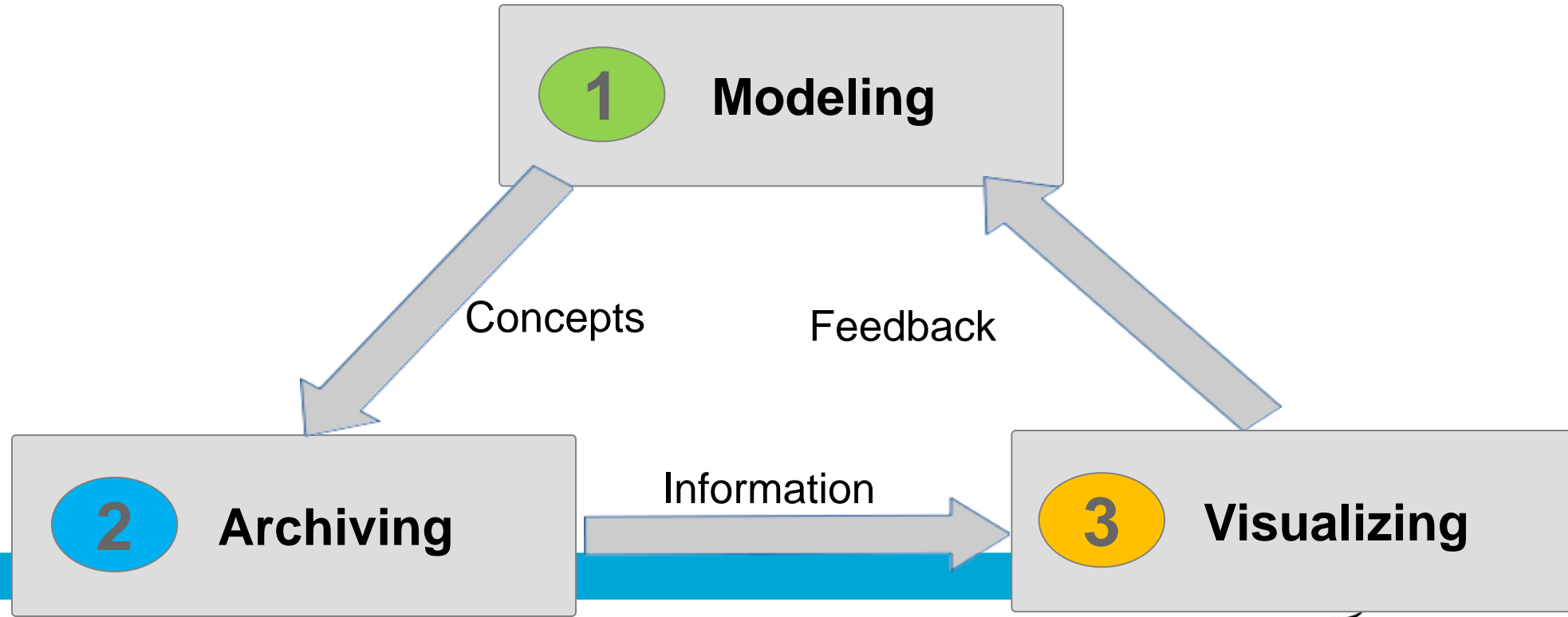
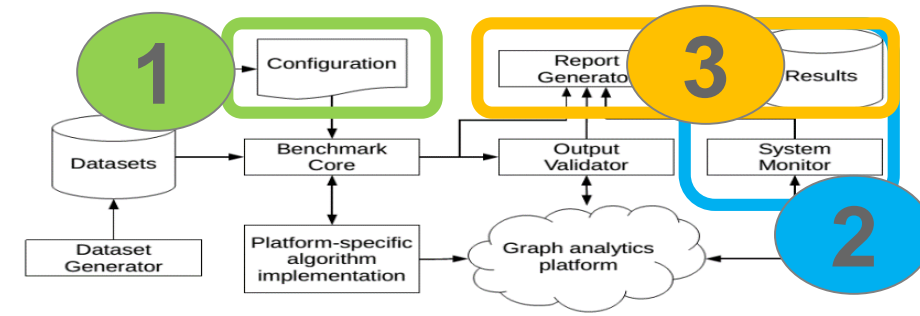
Many, fine-grained results

Fine-grained evaluation method is more comprehensive
... but more time-consuming, esp. to implement

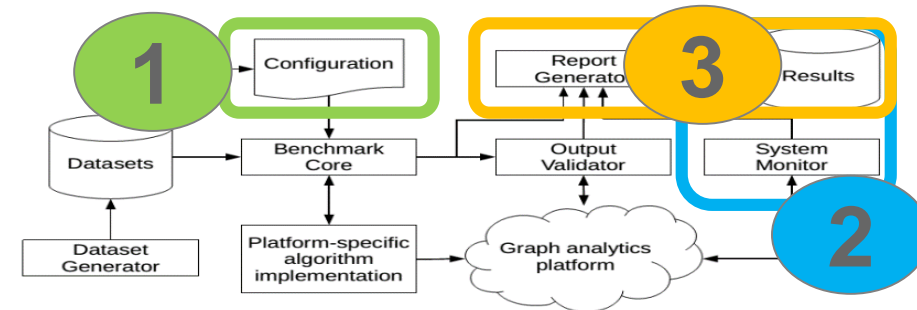
Graphalytics: Granula Overview

Granular

Fine-grained Method



1 Granula Modeller



Job

Operation

Operation [Actor @ Mission]

Info [StartTime]

Info [EndTime]

Info [.....]

Visual

Visual

Visual

Operation

Operation

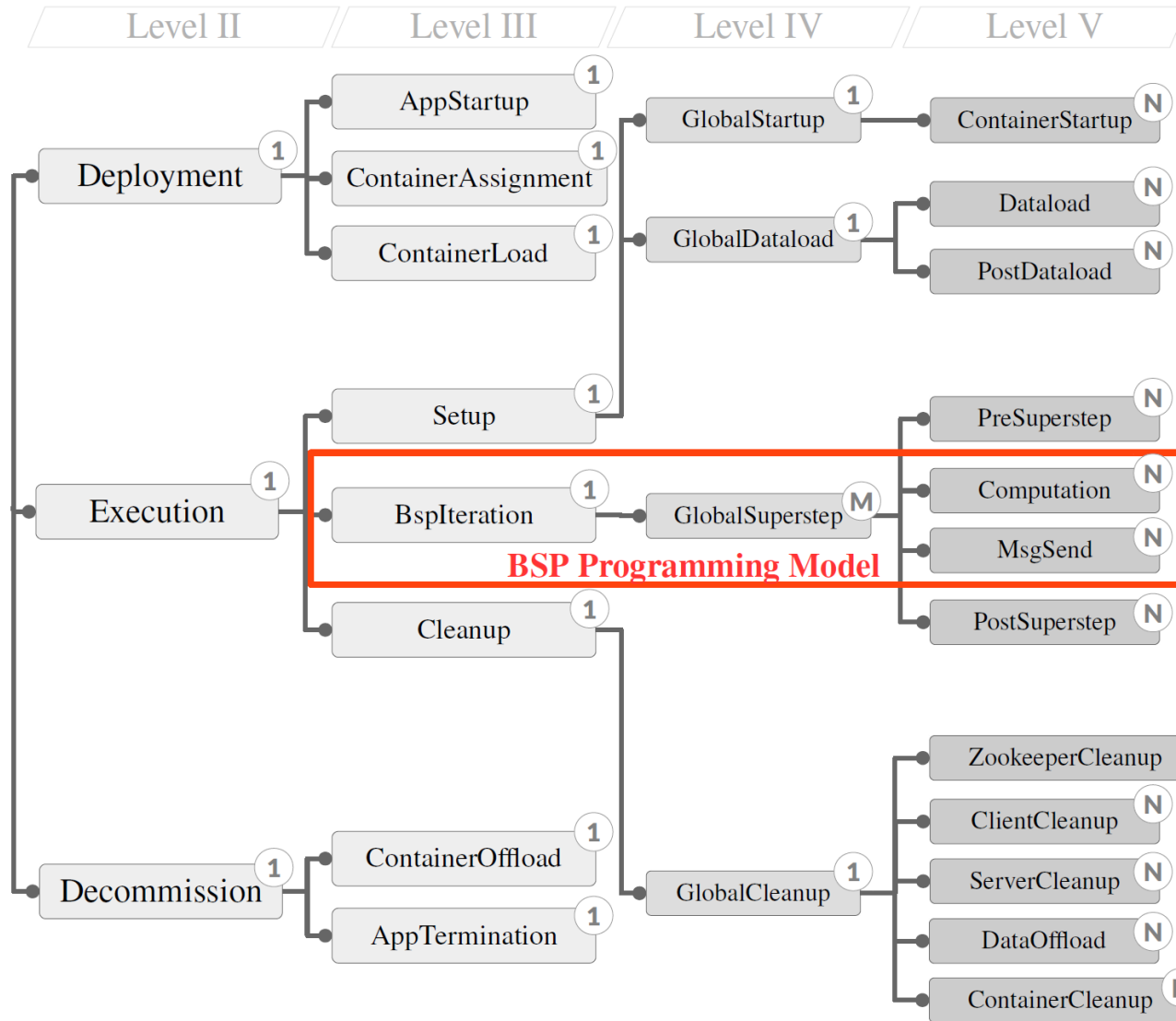
Operation

Time-consuming, expert-only,
Done only once per platform, but incrementally

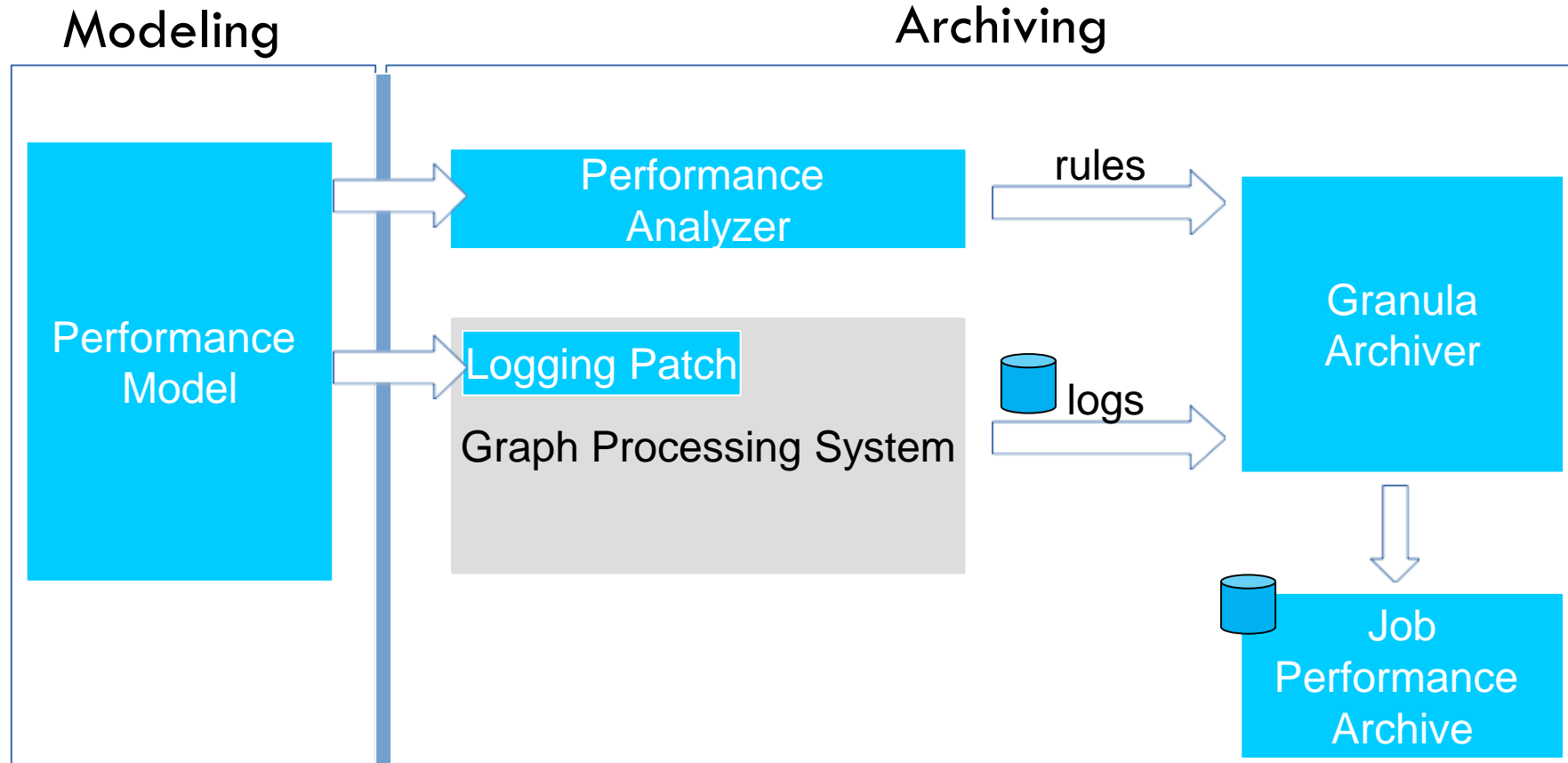
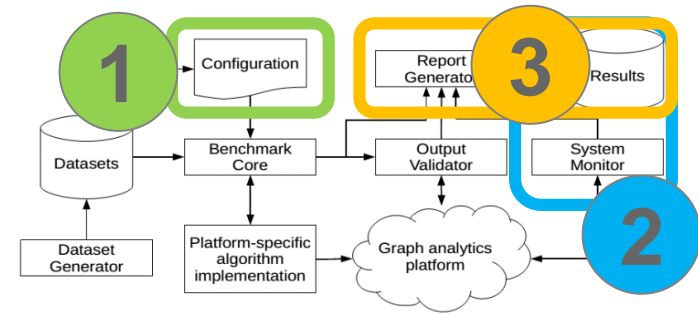
1

Granula Modeller

Incremental Model of Graph-Processing in Giraph



2 Granula Archiver

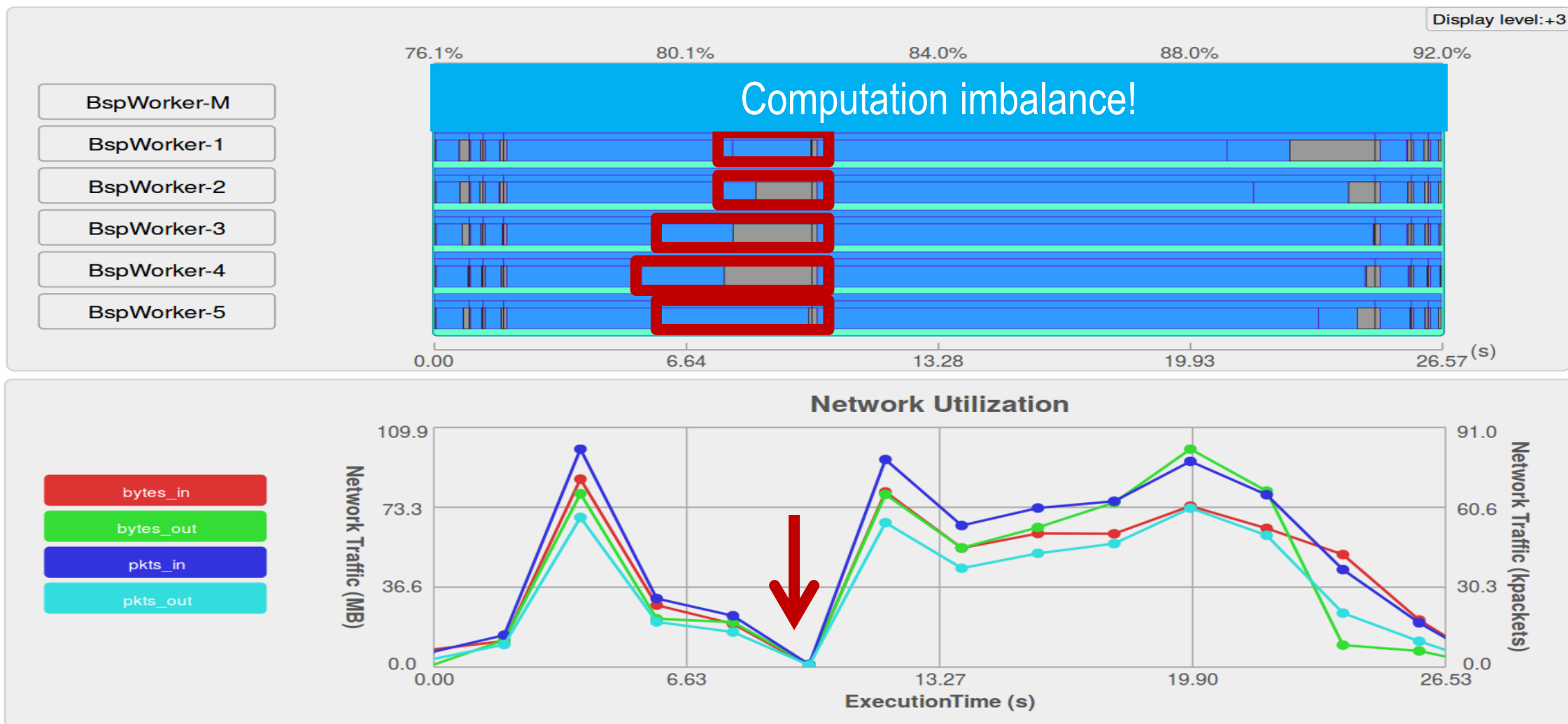


Time-consuming, minimal code invasion,
automated data collection at runtime, portable archive

3

Granula Visualizer

Portable choke-point analysis for everyone!



Graph-processing is at the core of our society

The data deluge vs. Analytics

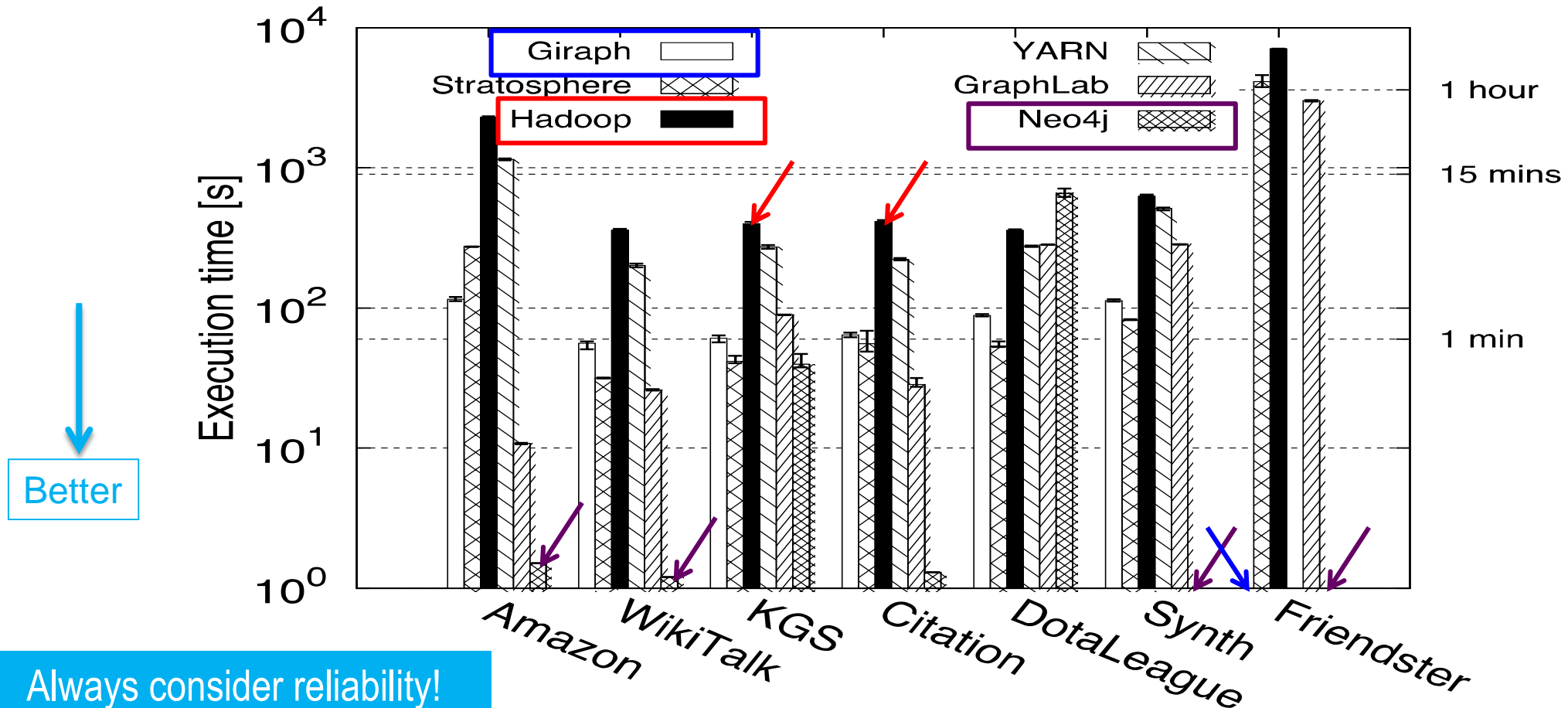
Graph processing @large

Graphalytics: Which system to select? What to tune? What to re-design?

A Performance Comparison of Graph-Processing Systems

Take-home message

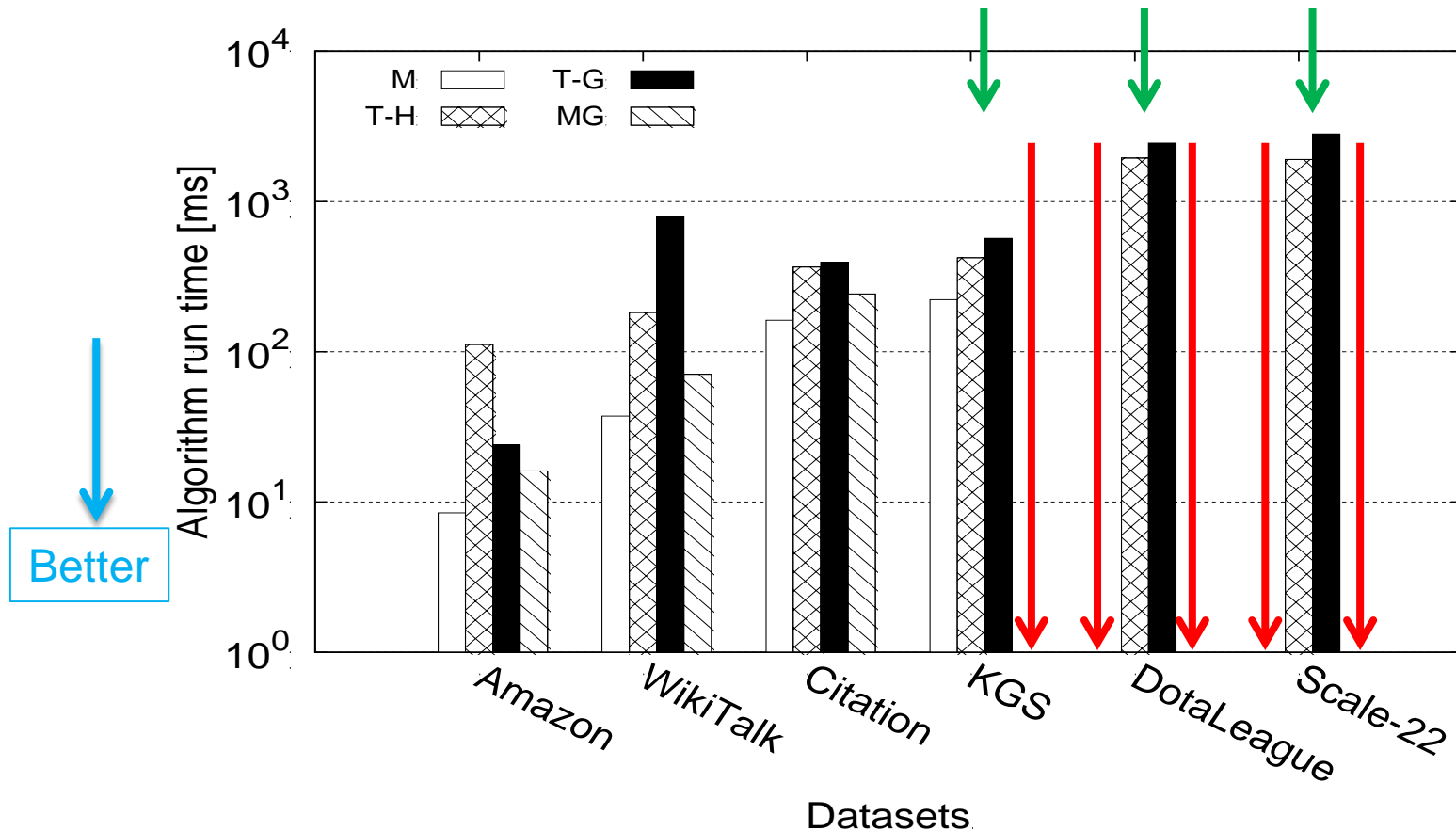
BFS Algorithm: All CPU Platforms, All Datasets



Always consider reliability!

No platform runs fastest for all graphs, but **Hadoop is the worst performer**.
Not all platforms can process all graphs, but **Hadoop processes everything**.

PageRank Algorithm: All GPU Platforms, Datasets



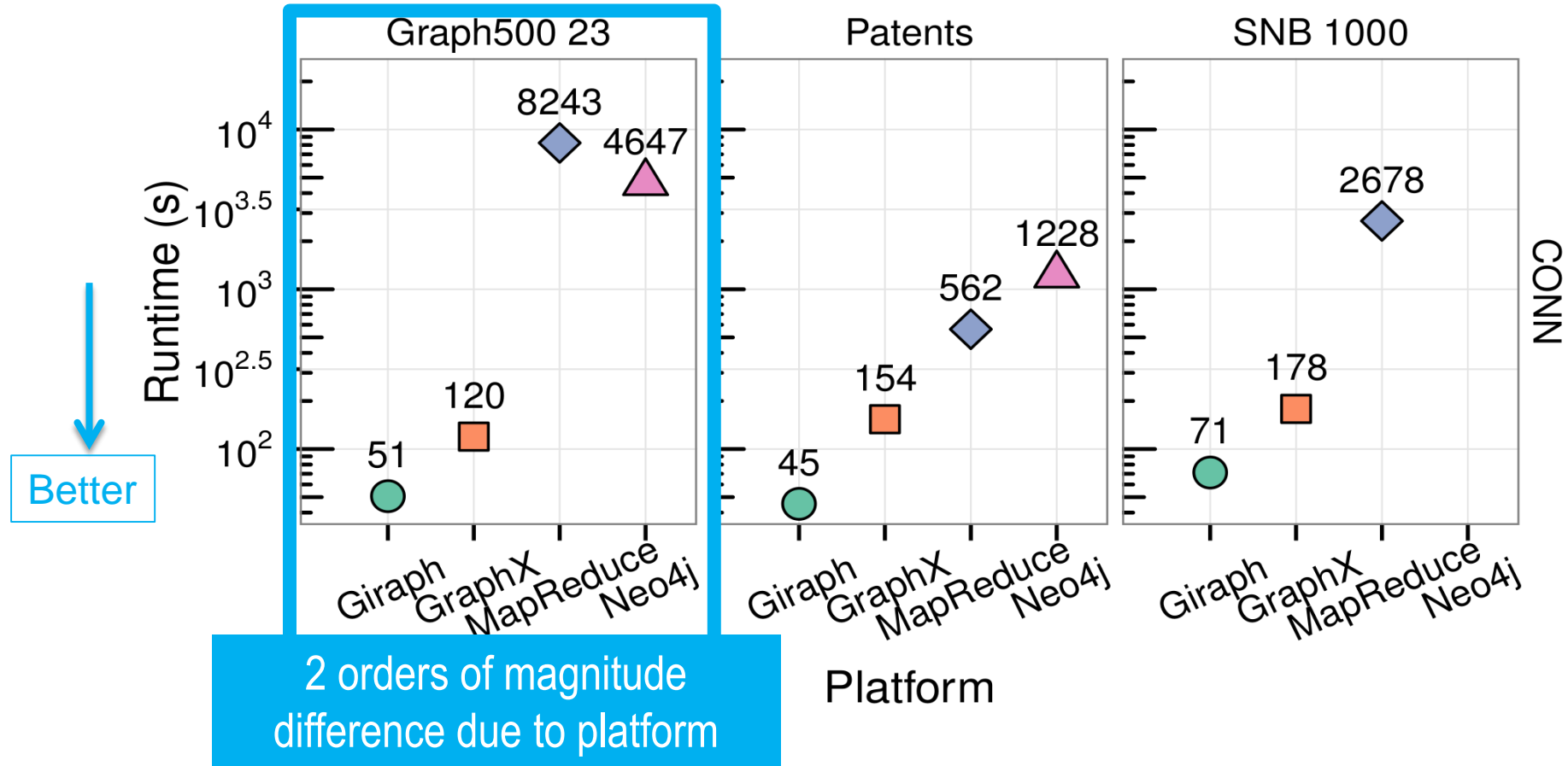
- **Medusa, MapGraph** fail on larger datasets, with MG failing earlier

- **Medusa** better for small datasets

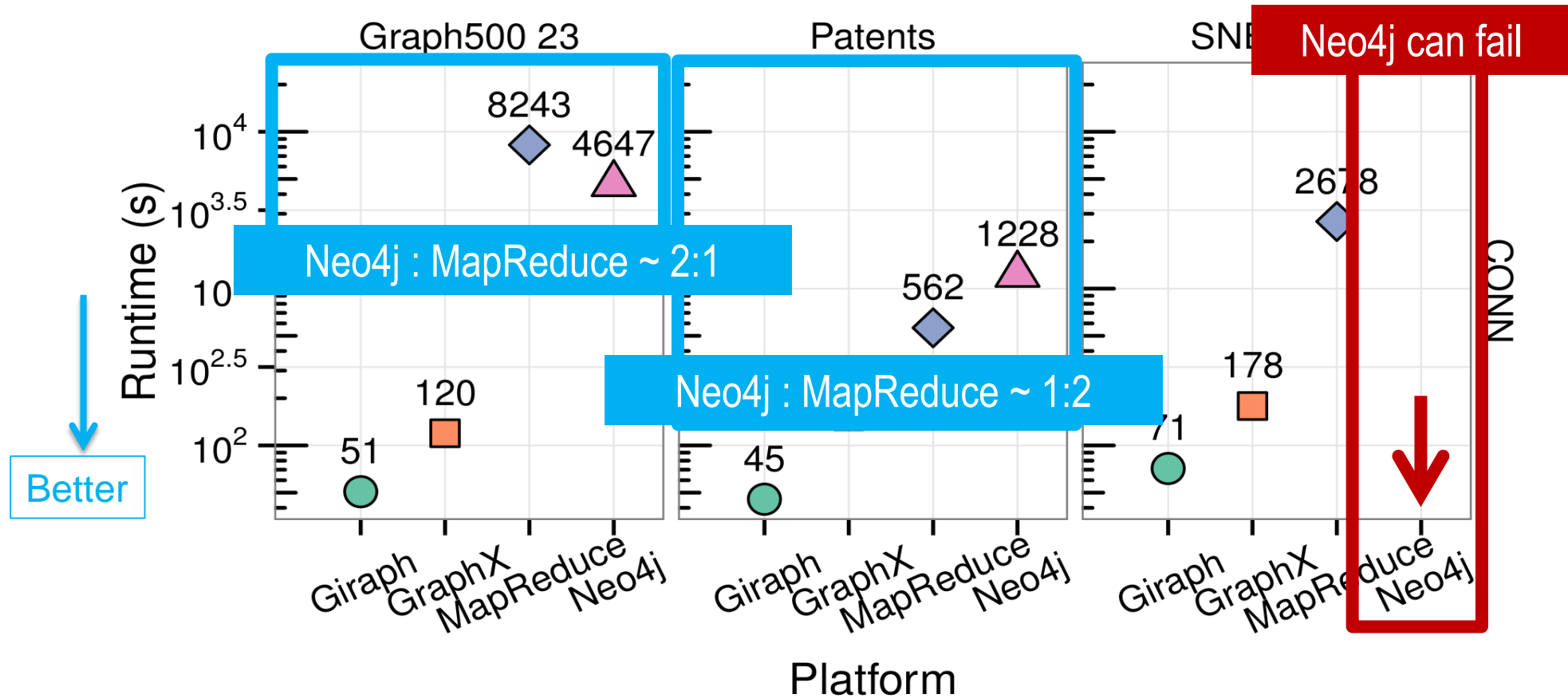
- **Totem** the only system able to process all tried datasets

Always consider reliability!

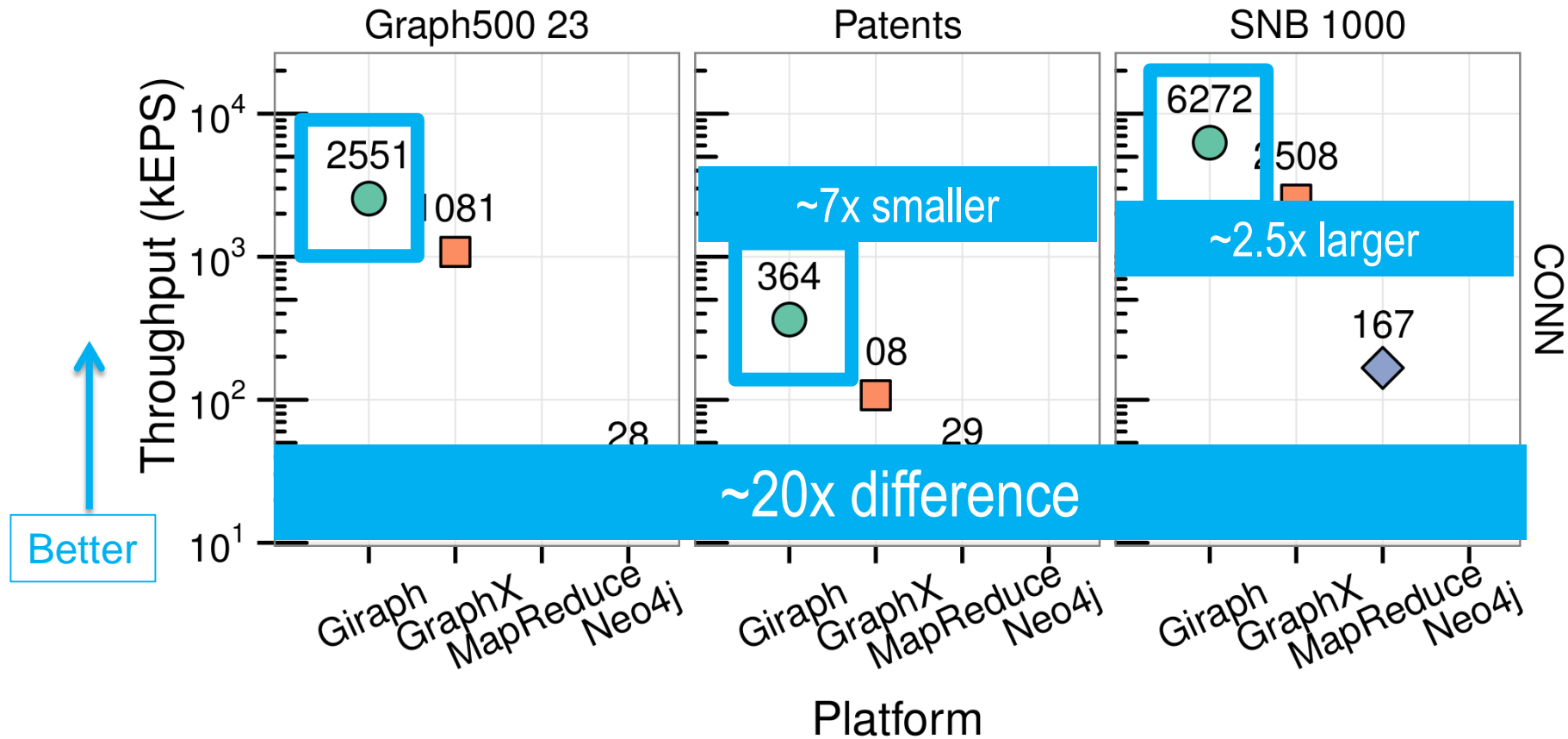
Runtime: The Platform Has Large Impact



Runtime: The Dataset Has Large Impact



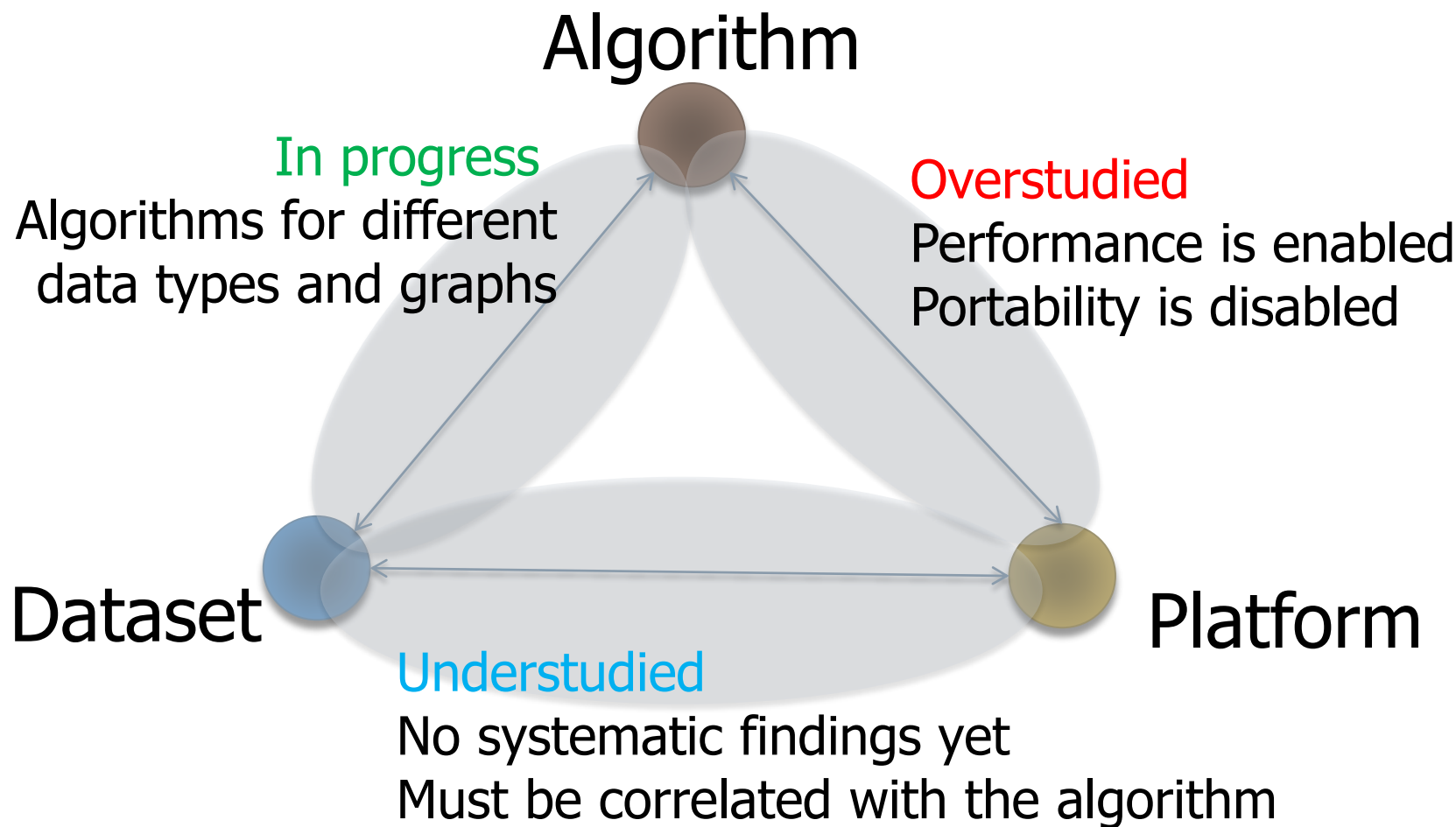
Throughput: The Dataset Has Large Impact





Introduced
by Ana Lucia
Varbanescu.

The Platform-Algorithm-Dataset (PAD) Triangle for Performance Engineering of Graph-Processing Systems



Lessons learned*

Performance of graph processing is function of
(Dataset, Algorithm, Platform, Deployment)

All current platforms have important drawbacks

(crashes, long execution time, tuning, etc.)

Best-performing is not only low response time

Scalability with cluster size/number of cores varies per system

Ease-of-use of a platform is very important

Graph-processing is at the core of our society

The data deluge vs. Analytics

Graph processing @large

Graphalytics: Which system to select? What to tune? What to re-design?

A performance comparison of graph-processing systems

Take-Home Message

Take-Home Message

Performance ↑

Graphalytics: unified view and benchmarking of tens of systems, with little effort

Granula: iterative, fine-grained, shareable performance evaluation to enable performance engineering

The P-A-D Triangle:

Performance = f (Algorithm, Dataset, Platform, ...)

Towards addressing the graph data deluge with high performance, low development effort

Development Effort →

Reading List

- Alexandru Iosup, Tim Hegeman, Wing Lung Ngai, Stijn Heldens, Arnau Prat Perez, Thomas Manhardt, Hassan Chaffi, Mihai Capotă, Narayanan Sundaram, Michael Anderson, Ilie Gabriel Tanase, Yinglong Xia, Lifeng Nai, Peter Boncz, *LDBC Graphalytics: A Benchmark for Large-Scale Graph Analysis on Parallel and Distributed Platforms*, VLDB'16.
- Mihai Capotă, Tim Hegeman, Alexandru Iosup, Arnau Prat-Pérez, Orri Erling, and Peter Boncz, *Graphalytics: A Big Data Benchmark for Graph-Processing Platforms*, International Workshop on Graph Data Management Experiences and Systems (GRADES), 2015.
- Guo et al., *An Empirical Performance Evaluation of GPU-Enabled Graph-Processing Systems*. CCGRID'15.
- A. Iosup, A. L. Varbanescu, M. Capotă, T. Hegeman, Y. Guo, W. L. Ngai, and M. Verstraaten, *Towards Benchmarking IaaS and PaaS Clouds for Graph Analytics*, Workshop on Big Data Benchmarking (WBDB), 2014.
- Y. Guo, M. Biczak, A. L. Varbanescu, A. Iosup, C. Martella, and T. L. Willke, *How Well Do Graph-Processing Platforms Perform? An Empirical Performance Evaluation and Analysis*, IEEE International Parallel and Distributed Processing Symposium (IPDPS), 2014, pp. 395–404.
- Y. Guo, A. L. Varbanescu, A. Iosup, C. Martella, and T. L. Willke, *Benchmarking graph-processing platforms*, ACM/SPEC International Conference on Performance Engineering (ICPE), 2014, pp. 289–292.

System diversity (GPU-enabled)

Dedicated Systems

Generic



medusa-gpu

Medusa: Simplified Graph Processing on GPUs



NetSysLab

mapgraph ^{Beta}

Massively Parallel Graph processing on GPUs

TOTEM

Gunrock

High-performance Graph Primitives on GPU



VertexAPI2

Medusa

- Enables the use of GPUs for graph processing
 - Single-node, multiple GPUs
 - In-memory processing
- Simple API that hides GPU programming
 - Edge- / vertex-granularity that enables fine-grained parallelism.
 - API calls are grouped in kernels
 - Kernels are scheduled on one or multiple GPUs
- Run-time for communicating with the GPU

Totem

- Enables use of GPUs (T-G)
- Enables *single-node* heterogeneous (T-H) computing on graphs
- Programming requires expert knowledge of all types of systems
 - C+CUDA+API for specifying applications
 - Based on BSP
- Partitions the data (edge-based) between CPUs and GPUs
 - Based on processing capacity
 - Minimizing the overhead of communication
 - Buffer schemes, aggregation, smart partitioning

MapGraph

- Target at high performance graph analytics on GPUs.
- Single GPU available and Multi-GPU ready
 - Also available in a CPU-only version
- API based on the Gather-Apply-Scatter (GAS) model as used in GraphLab.
 - Productivity-oriented API

Lessons learned from GPU-based systems

Brave attempts to enable the use of GPUs *inside* graph processing *systems*

Every system has its own quirks

- Lower level programming allows more optimizations, better performance

- Higher level APIs allow more productivity

No clear winner, performance-wise

Challenge:

- Distributed accelerated graph-processing

System diversity (CPU-based)

Dedicated Systems

Generic

ORACLE PGX



Neo4j
the graph database



GraphLab



StratoSphere
Above the Clouds



Hadoop (Generic)

The most popular MapReduce implementation

Generic system for large-scale computation

Pros:

Easy to understand model

Multitude of tools and storage systems

Cons:

Express the graph application in MapReduce

Costly disk and network operations

No specific graph processing optimizations



Hadoop2 with YARN (Generic)

Next generation of Hadoop

- Supports old MapReduce jobs

- Designed to facilitate multiple programming models (frameworks, e.g., Spark)

- Separates resource management (YARN) and job management

- MapReduce uses resources provided by YARN



Stratosphere (Generic)

Now Apache Flink (now 6M\$ investment from Intel)

Nephele resource manager

- Scalable parallel engine

- Jobs are represented as DAGs

- Supports data flow in-memory, via network, or via files

PACT job model

- 5 second-order functions (MapReduce has 2):

- Map, Reduce, Match, Cross, and CogGroup

- Code annotations for compile-time plans

- Compiled as DAGs for Nephele



Pregel: dedicated graph-processing + Apache Giraph (Dedicated)

Proposed a **vertex-centric** model for graph processing

Graph-to-graph transformations

Front-end:

Write the computation that runs on all vertices

Each vertex can vote to halt

All vertexes halt => terminate

Can add/remove edges and vertices

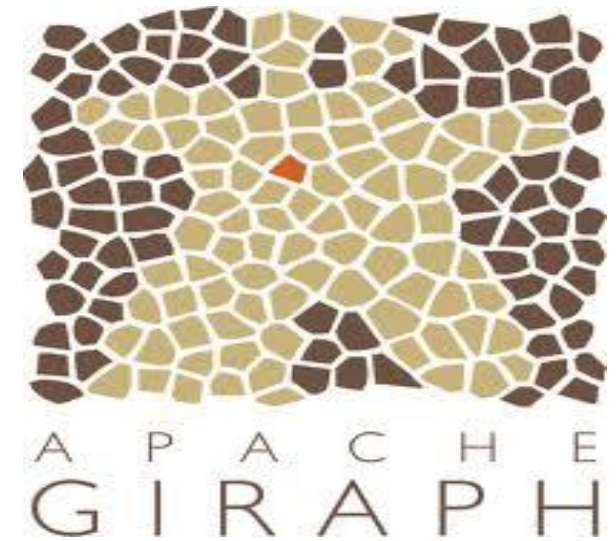
Back-end:

Uses the BSP model

Message passing between nodes

Combiners, aggregators

Checkpointing for fault-tolerance



GraphLab (Dedicated)

Distributed programming
model for machine learning

Provides an API for graph processing, C++ based
(now Python)

All in-memory

Supports asynchronous processing

GraphChi is its single-node version,
Dato as GraphLab company



Neo4J (Dedicated)

Very popular graph **database**

Graphs are represented as relationships and annotated vertices

Single-node system

Uses parallel processing

Additional caching and query optimizations

All in-memory

The most widely used solutions for medium-scale problems

Cluster version in development



PGX.D (Dedicated)

Designed for beefy clusters

- Fully exploits the underlying resources of modern beefy cluster machines

Low-overhead communication mechanism

- Lightweight cooperative context switching mechanism

Support for data-pulling (also data-pushing)

- Intuitive transformation of classical graph algorithms

Reducing traffic and balancing workloads

- Several advanced techniques: Selective Ghostnodes, edge based partitioning, edge chunking

Attend presentation of SC15 article!



PGX.D: Programming Model

High level programming model for Neighborhood Iteration Tasks

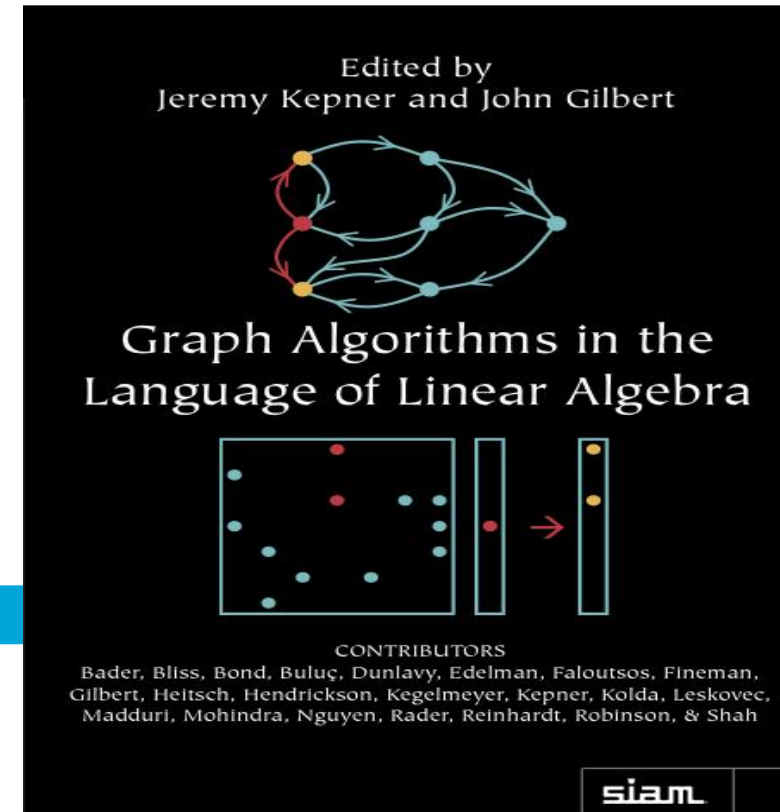
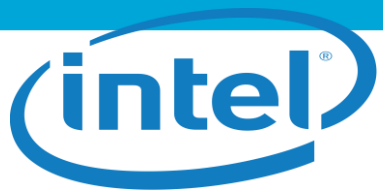
```
foreach(n: G.nodes)
  foreach(t: n.Nbrs)
    n.foo += t.bar
```



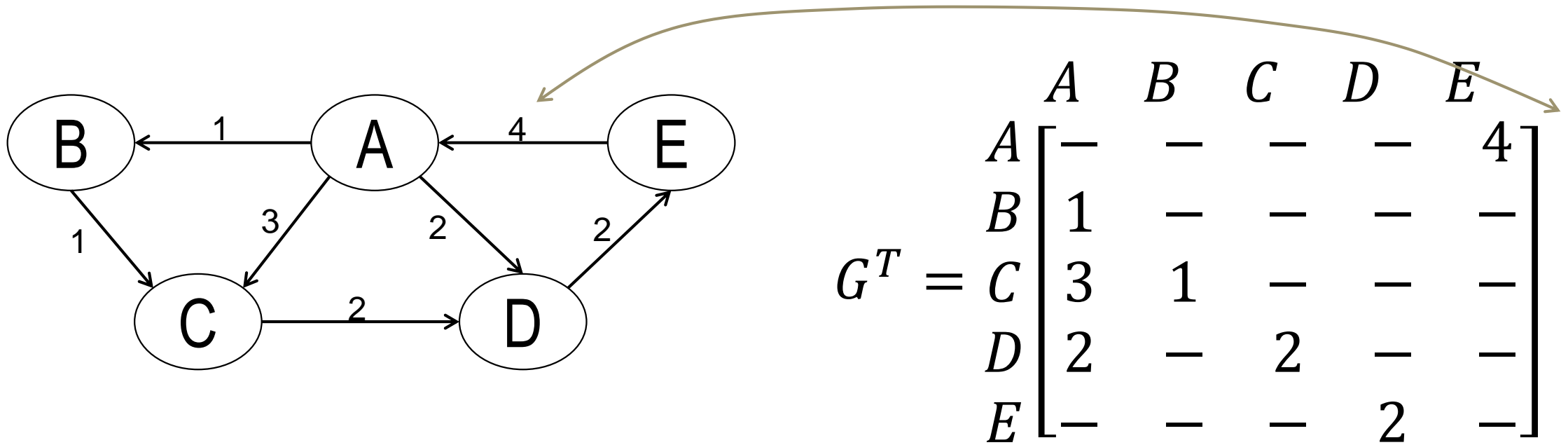
```
class my_task_pull : public innbr_iter_task {
  void run(..) {
    read_remote(get_nbr_id(), bar);
  }
  void read_done(void* buffer,..) {
    int foo_v = get_local<int>(node_id, foo);
    int bar_v = get_data<int>(buffer);
    set_local(node_id, foo_v + bar_v, foo);
  }
}
```

GraphMat (Dedicated)

- Vertex programming as front-end and sparse matrix operations as back-end
 - ▣ “Matrix level performance with vertex program productivity”
 - ▣ Unifying vertex programming w linear algebra is new



Example



A Vertex Program (Single Source Shortest Path) ~ Giraph

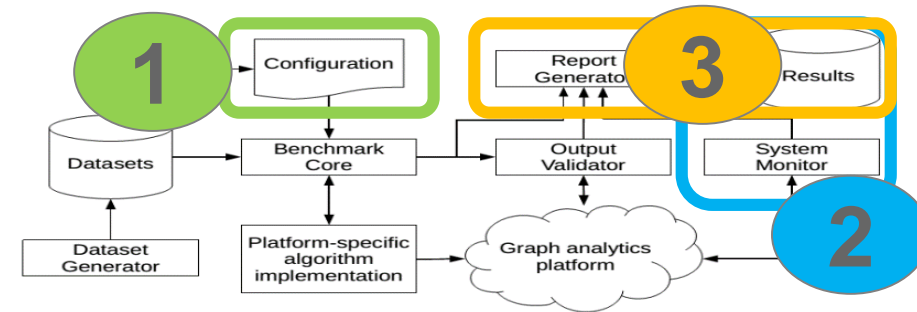
SEND_MESSAGE : message := vertex_distance

PROCESS_MESSAGE : result := message + edge_value

REDUCE : result := min(result, operand)

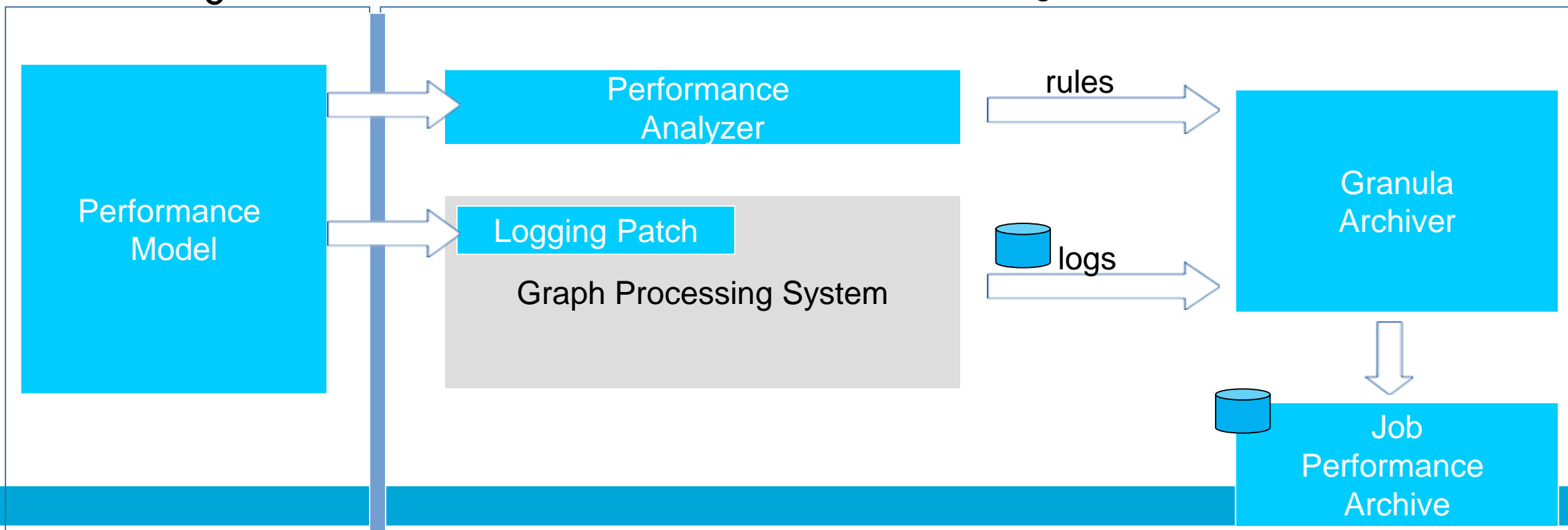
APPLY : vertex_distance = min(result, vertex_distance)

2 Granula Archiver



Modeling

Archiving

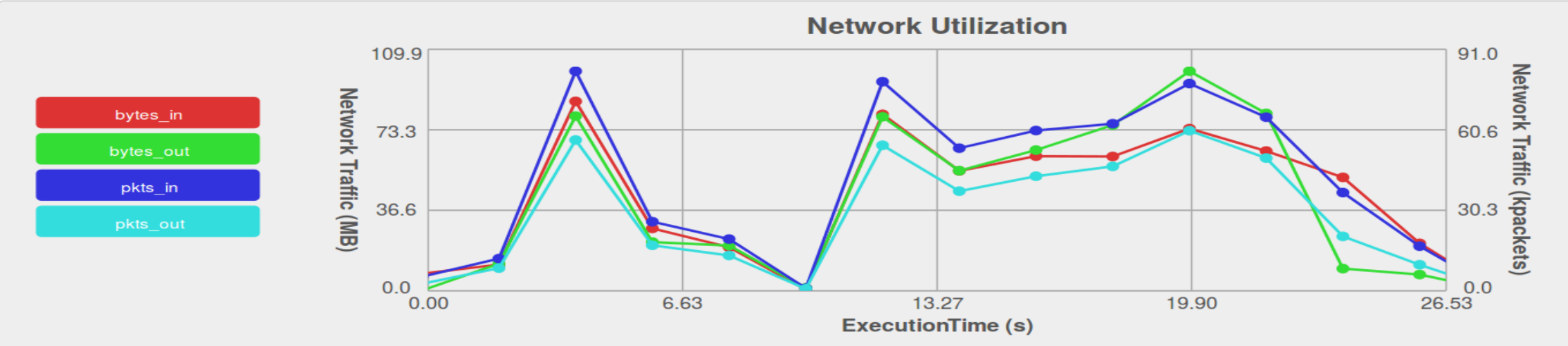
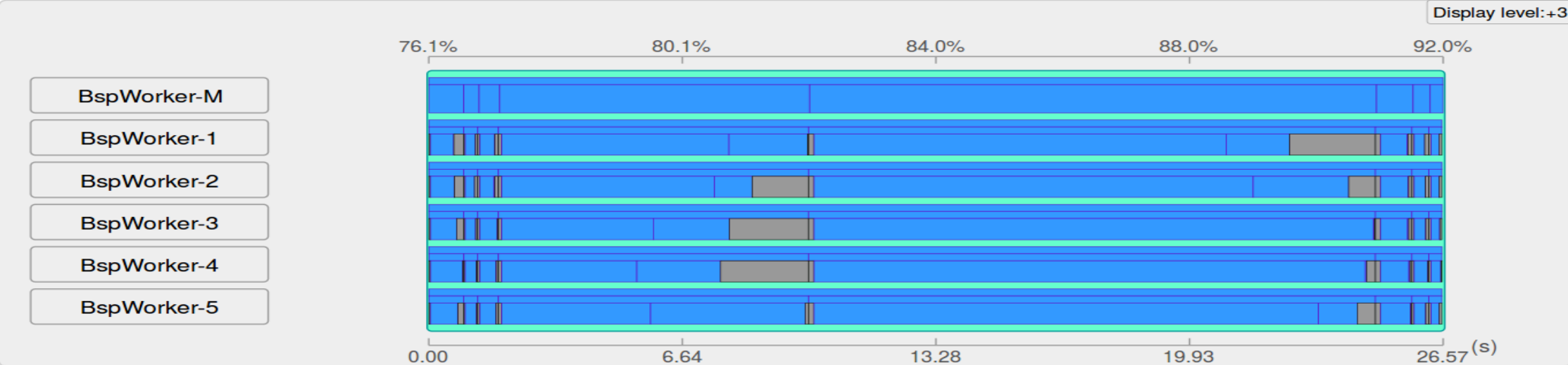
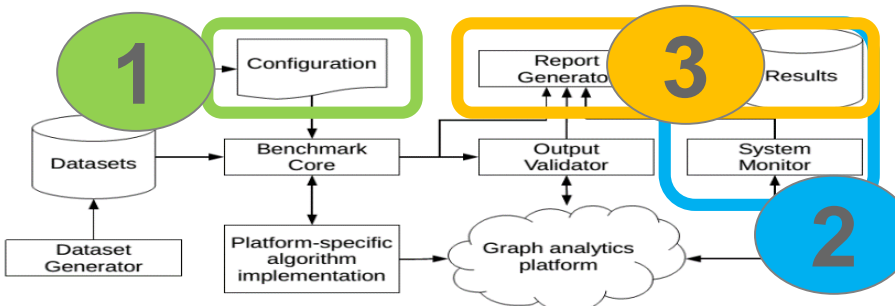


Time-consuming, minimal code invasion,
automated data collection at runtime, portable archive

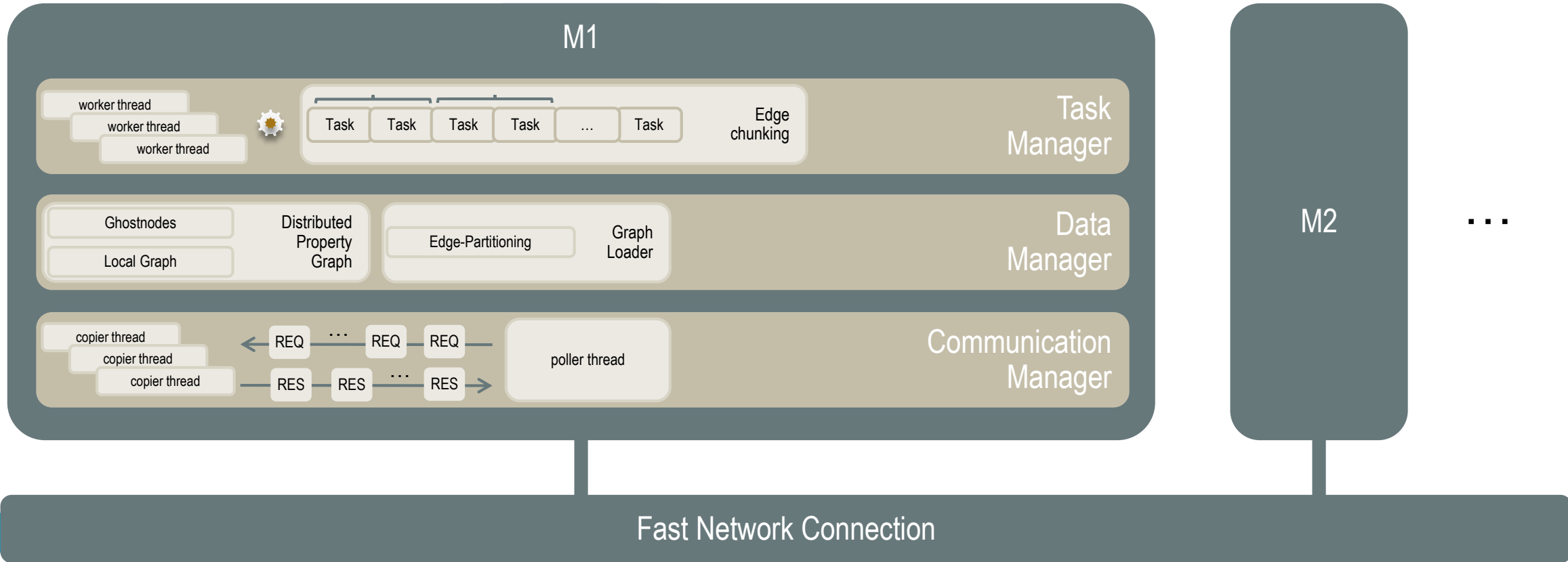
3

Granula Visualizer

Portable choke-point analysis for everyone!



PGX.D: System Design Overview



Icons

