# Graphalytics:
# A Big Data Benchmark for Graph-Processing Platforms

<u>Mihai Capotă</u>, Tim Hegeman, Alexandru Iosup,
Arnau Prat-Pérez, Orri Erling, Peter Boncz

Delft University of Technology
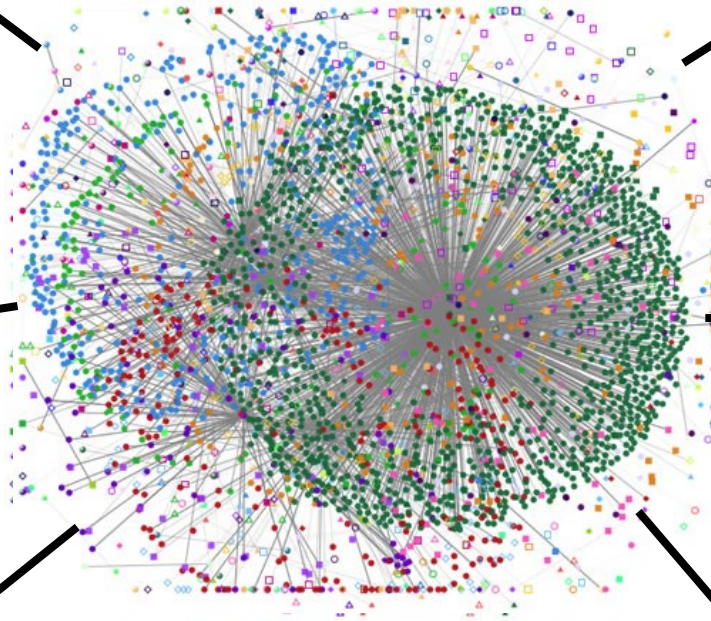Universitat Politècnica de Catalunya
OpenLink Software
CWI

TUDelft

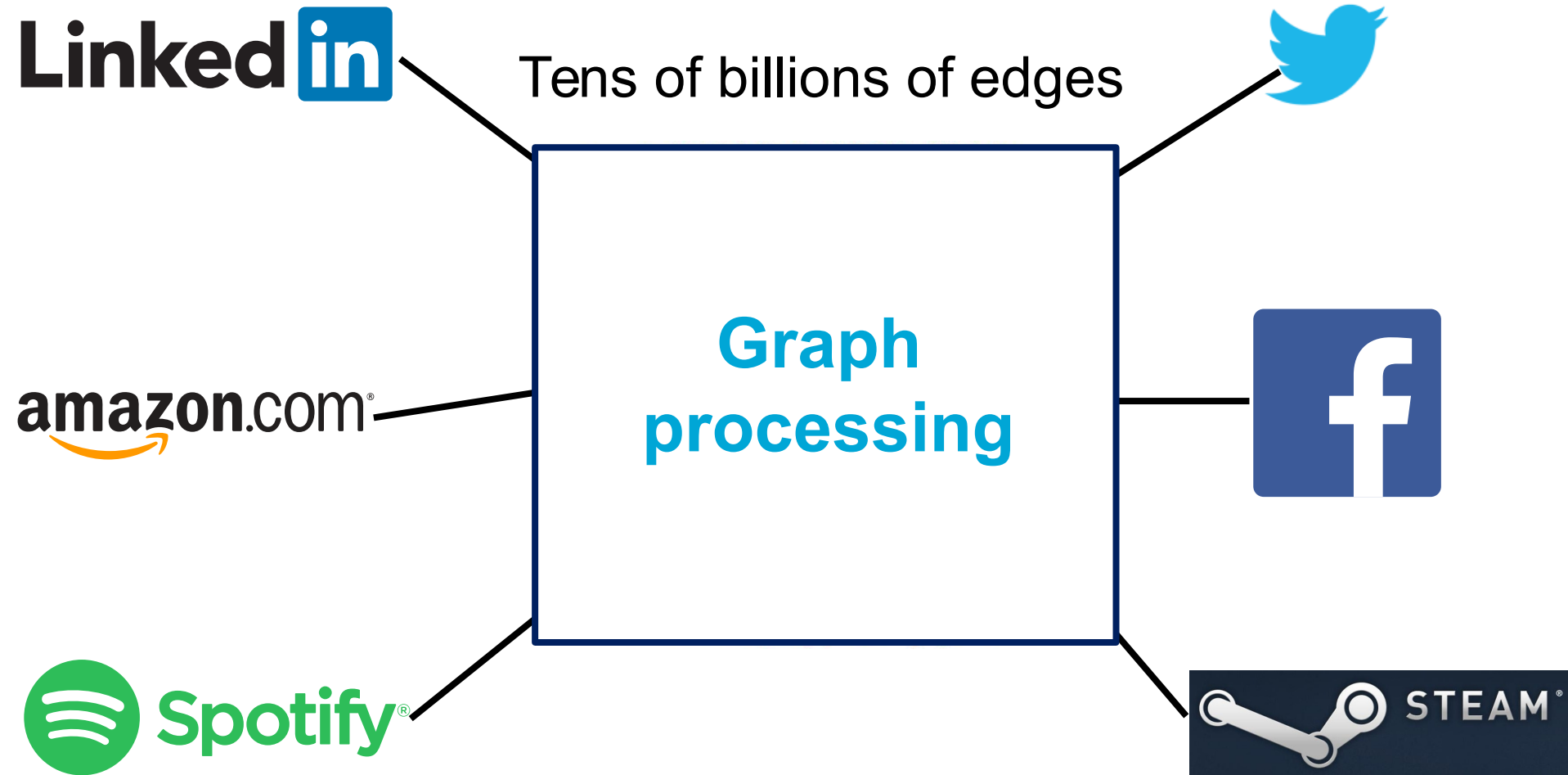# The data deluge: large-scale graphs



Tens of billions of edges

# The data deluge: large-scale graphs

Tens of billions of edges

**Graph processing**

# Platform diversity
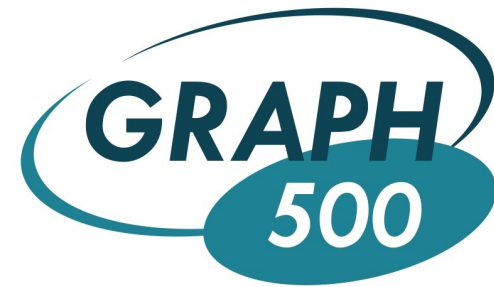
Oracle Labs
PGX

GraphLab

GraphX

# Yet another benchmark?

# Yet another benchmark?

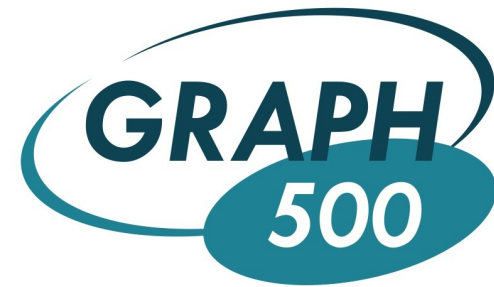- Lack of benchmarks for generic graph processing platforms

# Yet another benchmark?

- Lack of benchmarks for generic graph processing platforms

- Graph500

  - BFS

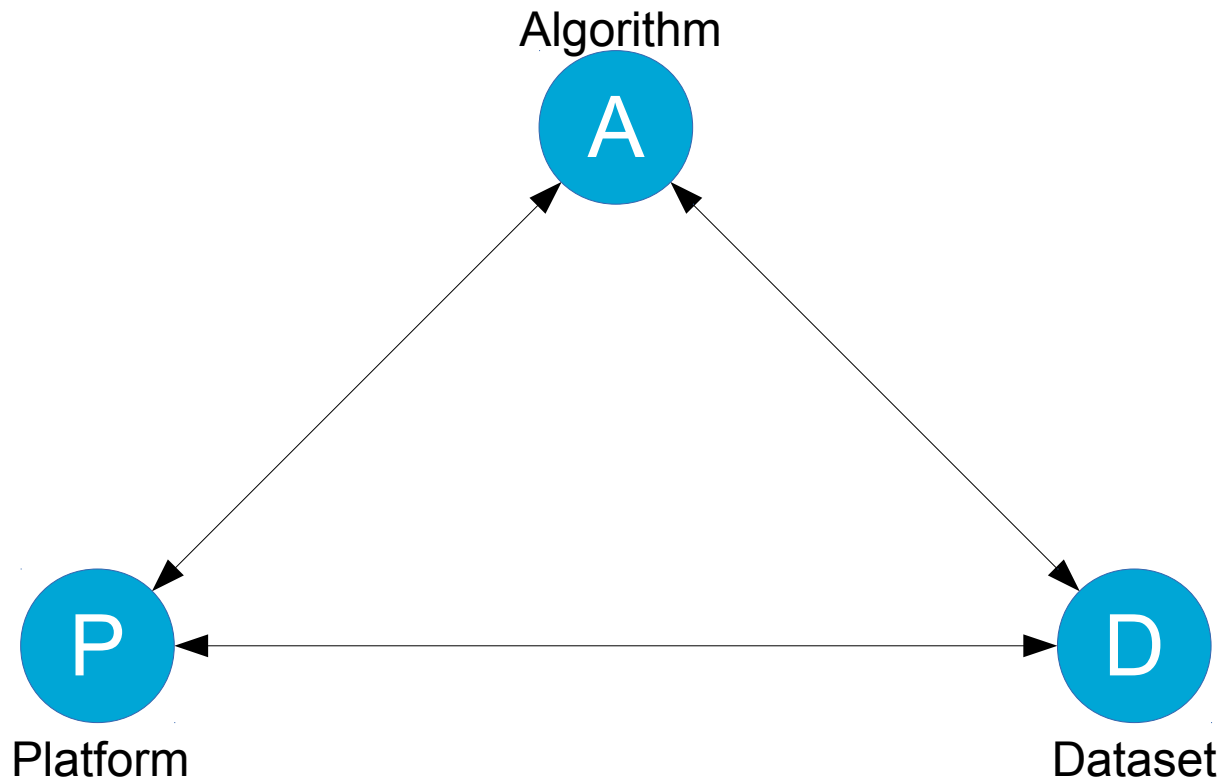  - Kroneker graph



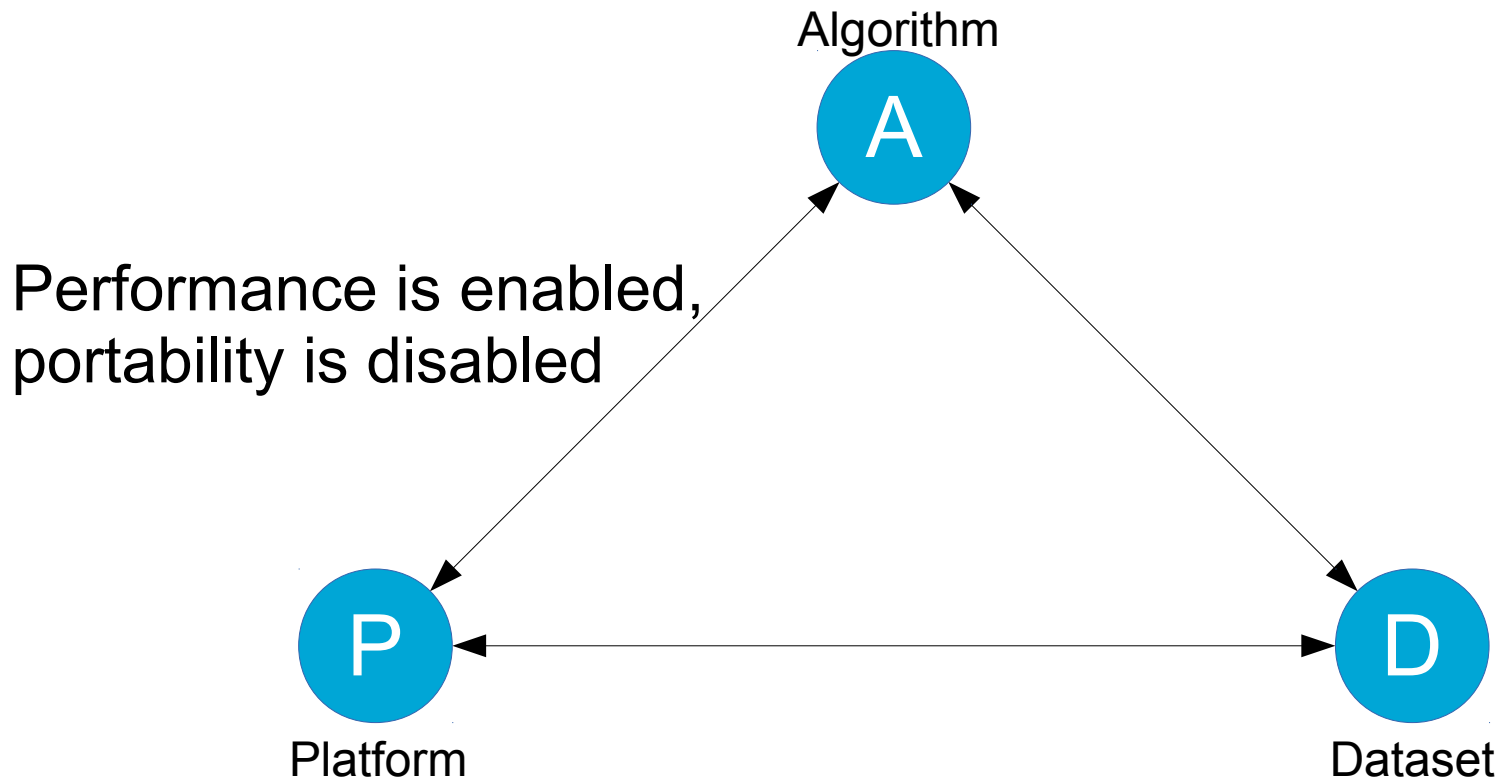**TU**Delft

# Yet another benchmark?

- Lack of benchmarks for generic graph processing platforms

- Graph500
  - BFS
  - Kroneker graph

- Several academic studies

  - Specific to graph or RDF databases
  - Ad hoc setup, difficult to extend

# P-A-D triangle

# P-A-D triangle
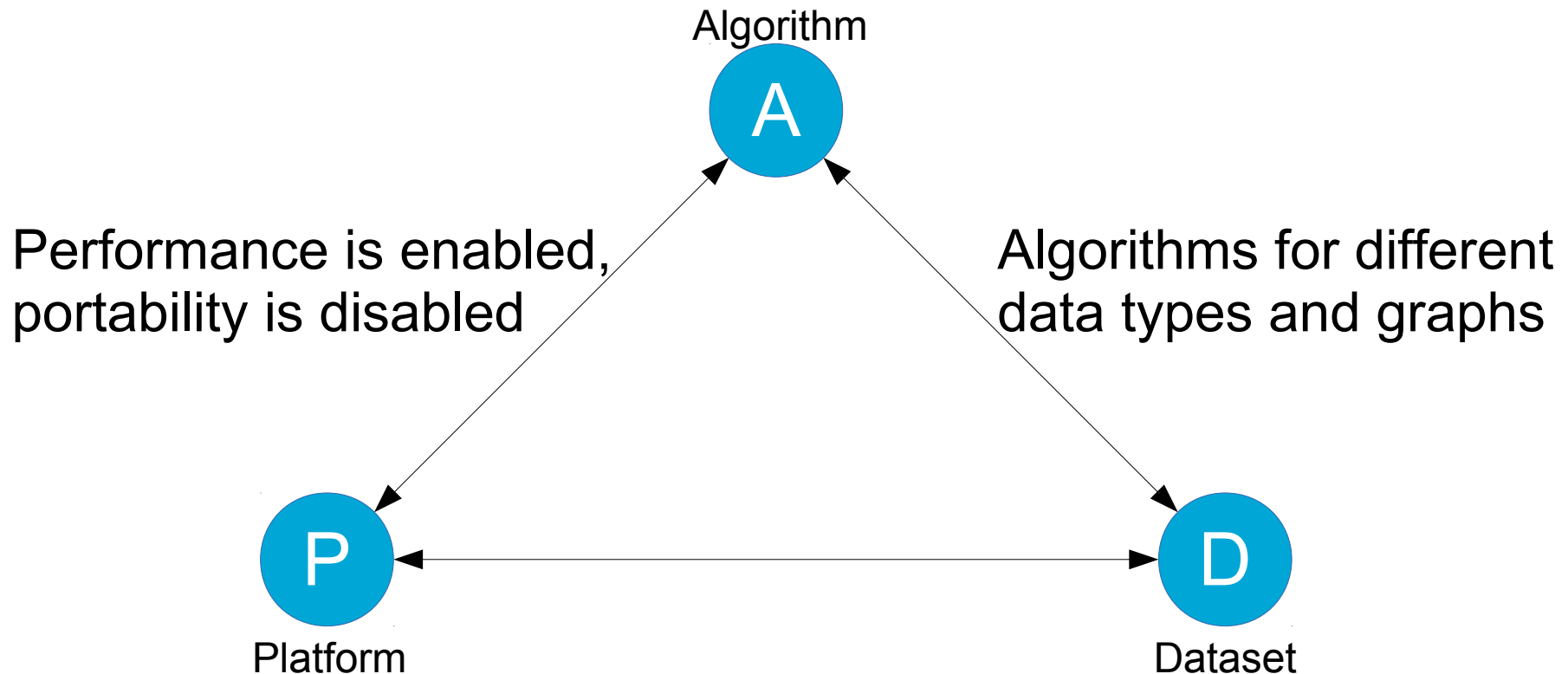
Algorithm

A

Performance is enabled,
portability is disabled

P

Platform

D

Dataset

# P-A-D triangle



Algorithm

A

Performance is enabled, portability is disabled

Algorithms for different data types and graphs

P

Platform

D

Dataset

# P-A-D triangle

Algorithm

**A**

Performance is enabled, portability is disabled

Algorithms for different data types and graphs

**P**

**D**

No systematic findings yet

Platform

Dataset

# P-A-D triangle

Algorithm

A

Performance is enabled, portability is disabled

Deployment

Algorithms for different data types and graphs

P

No systematic findings yet

Platform

D

Dataset

# Graphalytics

# Graphalytics

- The first comprehensive benchmark for big data graph-processing platforms

**TU**Delft

# Graphalytics

- The first comprehensive benchmark for big data graph-processing platforms

    - Advanced benchmark harness

**TU**Delft

# Graphalytics

- The first comprehensive benchmark for big data graph-processing platforms
  - Advanced benchmark harness
  - Choke-point analysis

**TU**Delft

# Graphalytics

- The first comprehensive benchmark for big data graph-processing platforms

  – Advanced benchmark harness

  – Choke-point analysis

  – Realistic graph generator

**TU**Delft

# Graphalytics

- The first comprehensive benchmark for big data graph-processing platforms
  - Advanced benchmark harness
  - Choke-point analysis
  - Realistic graph generator

# Graphalytics

- The first comprehensive benchmark for big data graph-processing platforms

  - Advanced benchmark harness

  - Choke-point analysis

  - Realistic graph generator
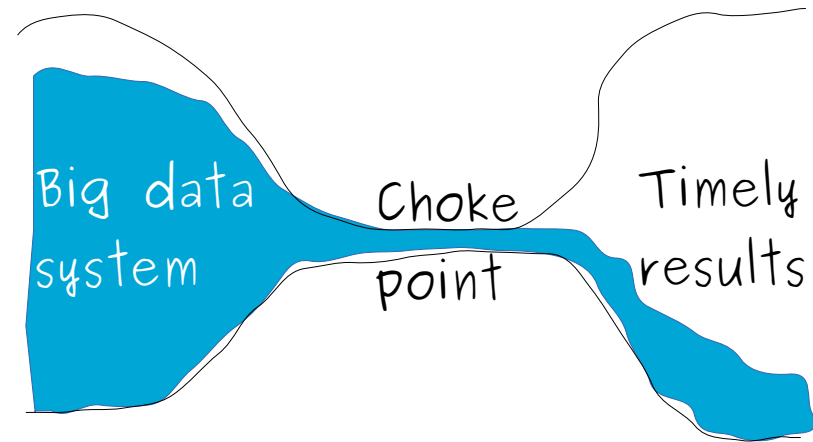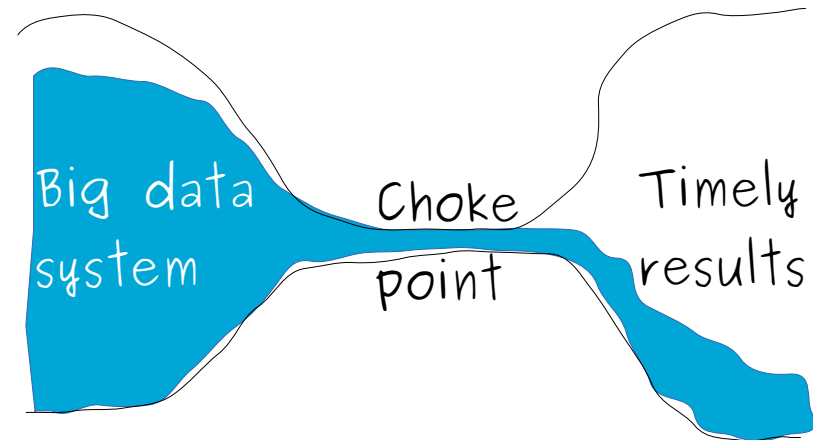
- Co-sponsored by Oracle

# Choke points

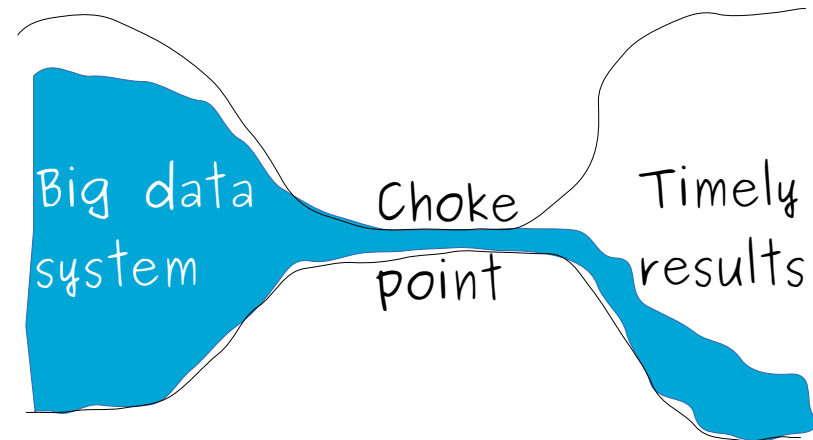# Choke points

- Technological challenges for platforms

Big data system

Choke point

Timely results

# Choke points

- Technological challenges for platforms
- Identified by experts
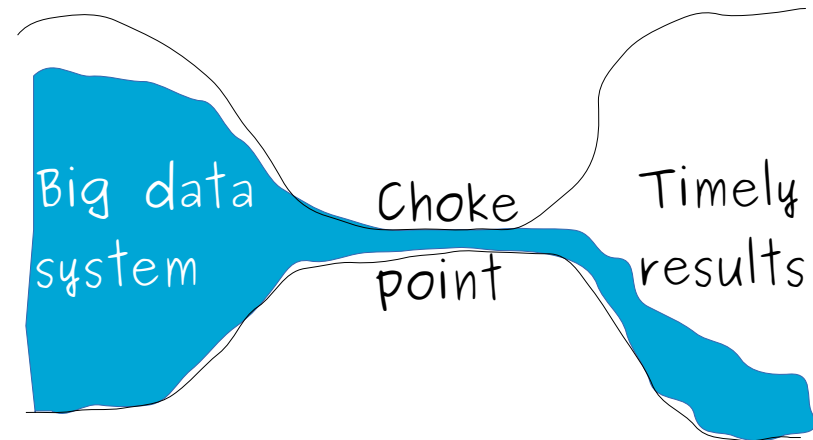
Big data system · Choke point · Timely results

# Choke points

- Technological challenges for platforms

- Identified by experts

- Real-world scenarios are enhanced to stress choke points

  - Prevent tunnel vision

  - Advance state of the art

Big data system

Choke point

Timely results

# Choke points

- Technological challenges for platforms

- Identified by experts

- Real-world scenarios are enhanced to stress choke points

  - Prevent tunnel vision

  - Advance state of the art

Big data system    Choke point    Timely results

Erling et al., The LDBC Social Network Benchmark: Interactive Workload , SIGMOD 2015

# Examples of choke points

# Examples of choke points

- Network utilization

**TU**Delft

# Examples of choke points

- Network utilization

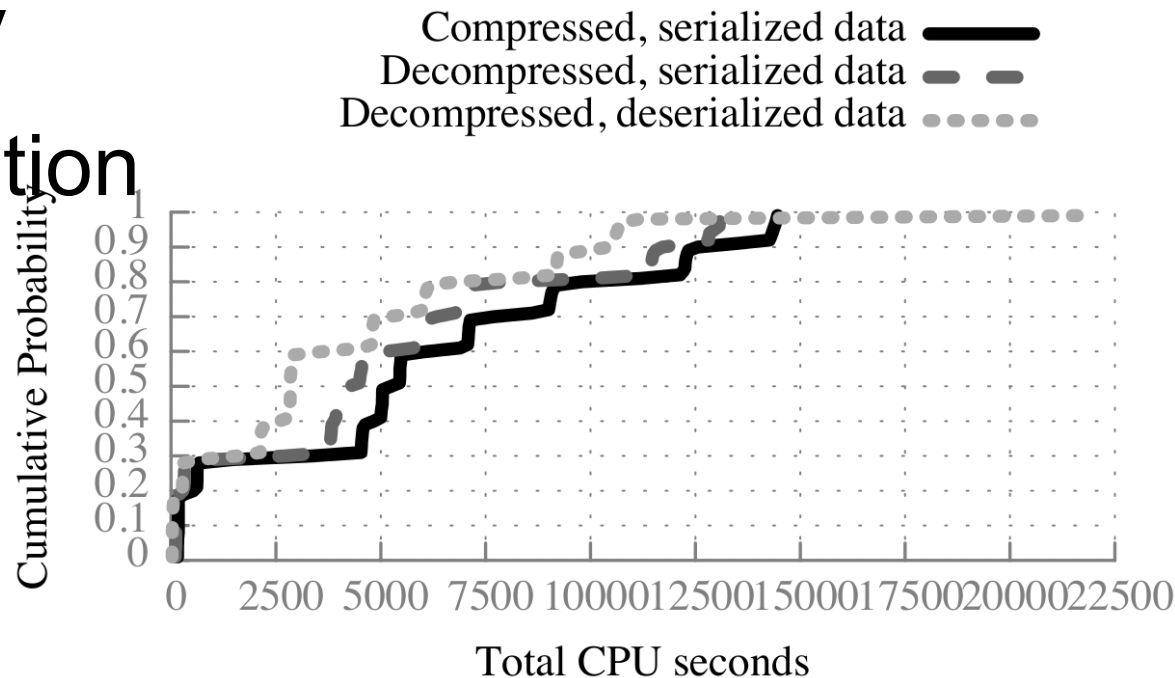- Memory footprint

**TU**Delft

# Examples of choke points

- Network utilization

- Memory footprint

- Access locality

**TU**Delft

# Examples of choke points

- Network utilization

- Memory footprint

- Access locality

- Skewed execution

# Examples of choke points

- Network utilization

- Memory footprint

- Access locality
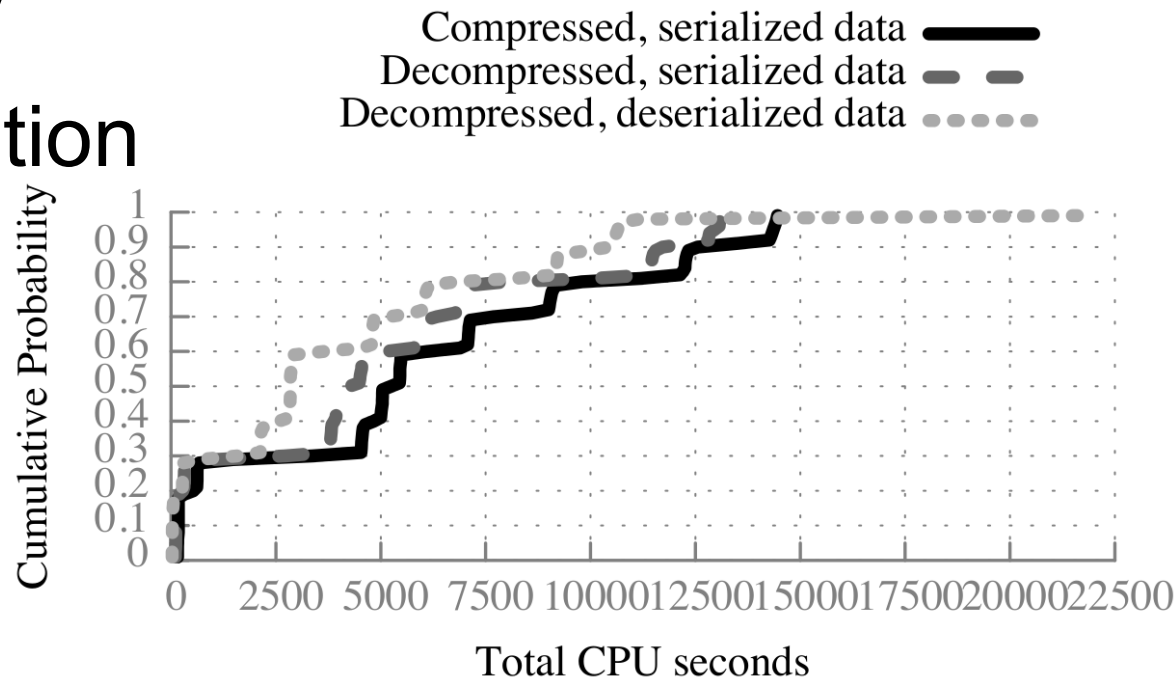
- Skewed execution

- CPU?

**TU**Delft

# Examples of choke points

- Network utilization

- Memory footprint

- Access locality

- Skewed execution

- CPU?



Ousterhout et al., "Making Sense of
Performance in Data Analytics Frameworks", NSDI 2015

# Examples of choke points

- Network utilization

- Memory footprint

- Access locality

- Skewed execution

- CPU?

- Others?



Ousterhout et al., "Making Sense of Performance in Data Analytics Frameworks", NSDI 2015

**TU**Delft

33

# Realistic graph generator

# Realistic graph generator

- LDBC Datagen

  - Synthetic social network similar to Facebook

# Realistic graph generator
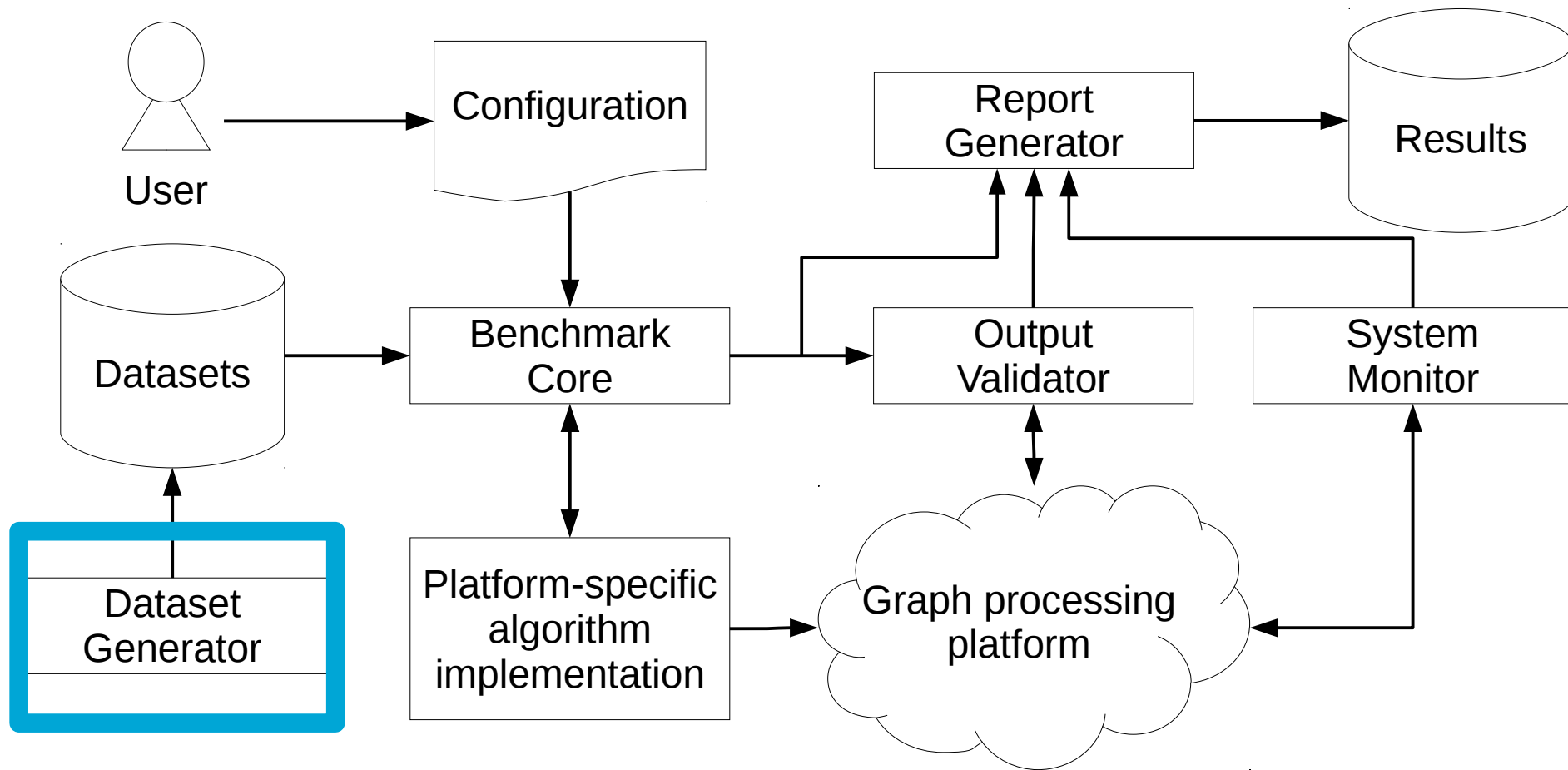
- LDBC Datagen

    - Synthetic social network similar to Facebook

- Graphalytics enhancements

    - Multiple degree distributions

        - Zeta and geometric implemented

# Realistic graph generator

- LDBC Datagen

  - Synthetic social network similar to Facebook

- Graphalytics enhancements

  - Multiple degree distributions

    - Zeta and geometric implemented

  - Other structural characteristics

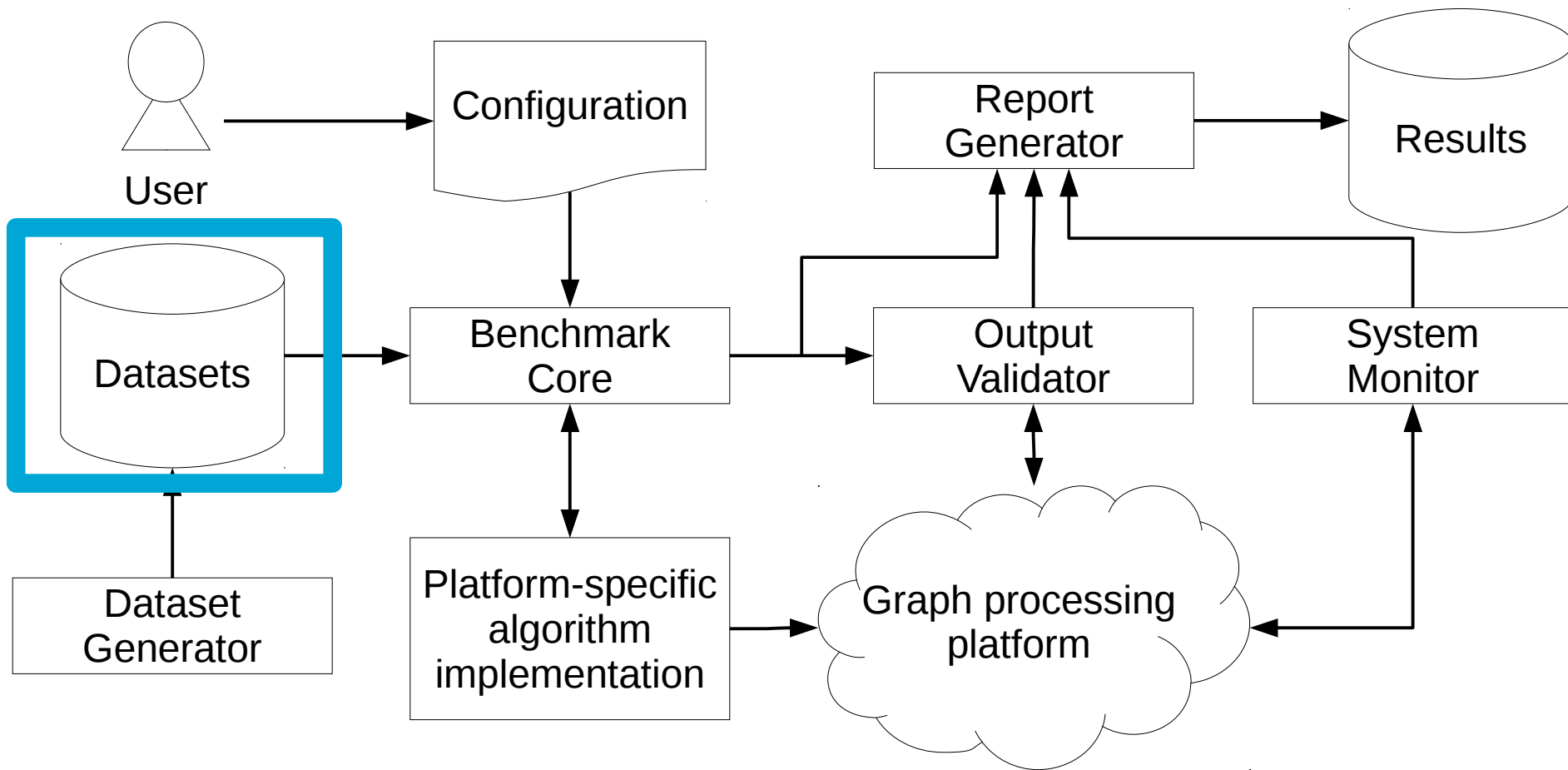    - Clustering coefficient
    - Assortativity

**TU**Delft
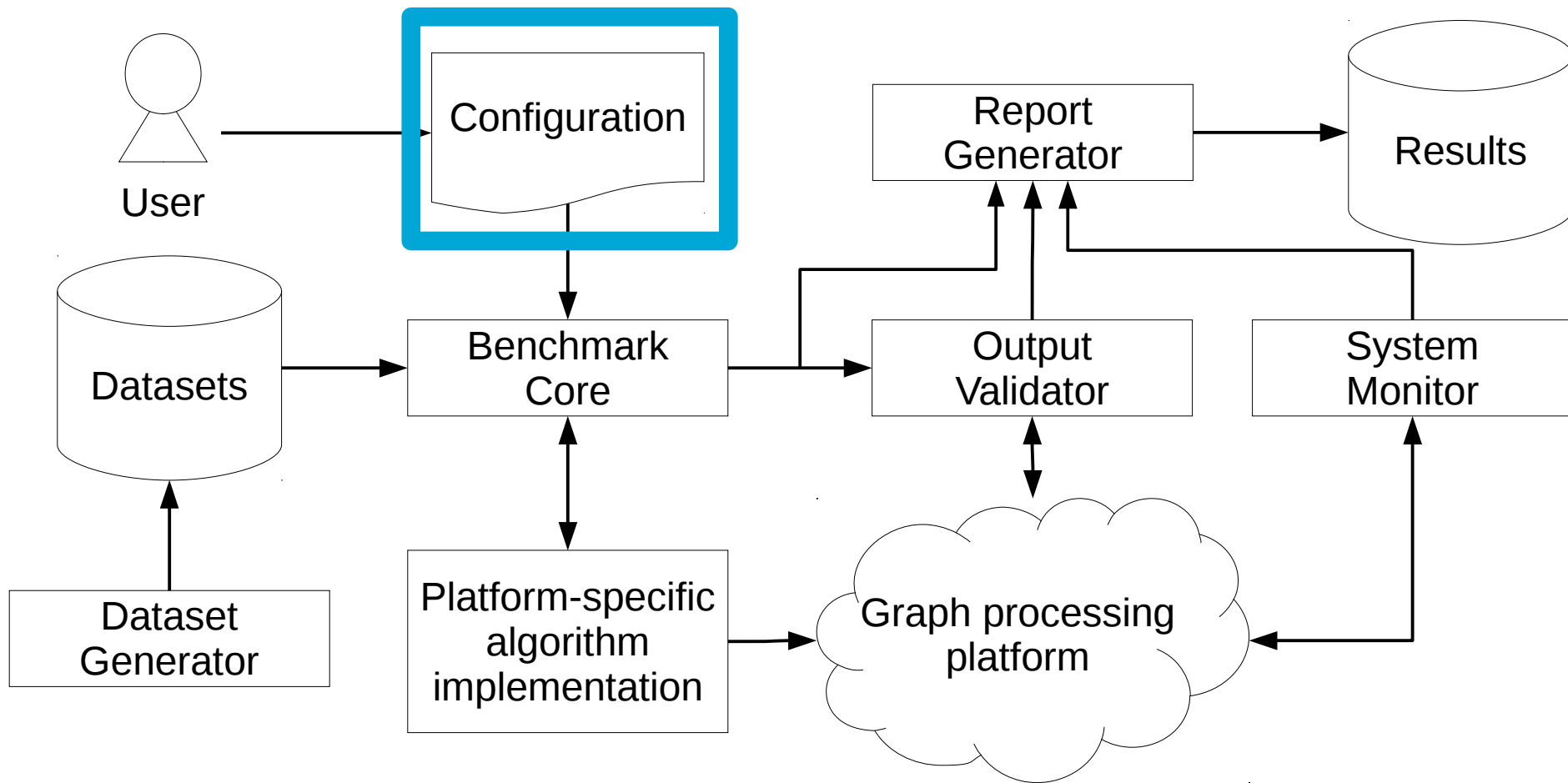
# Realistic graph generator

- ## LDBC Datagen

  - Synthetic social network similar to Facebook

- ## Graphalytics enhancements

  - Multiple degree distributions

    - Zeta and geometric implemented

  - Other structural characteristics

    - Clustering coefficient
    - Assortativity

Erling et al., The LDBC Social Network Benchmark:
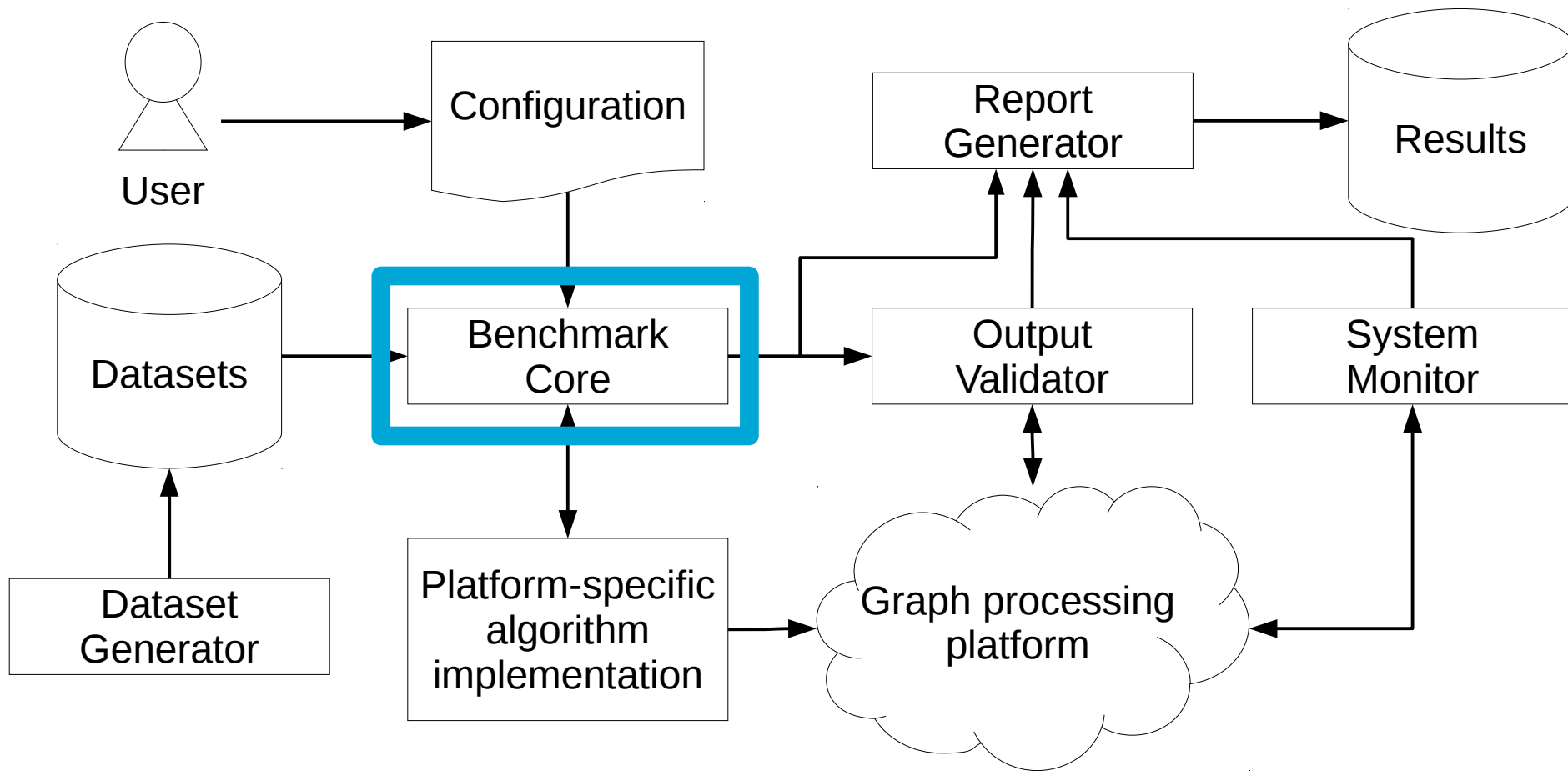Interactive Workload , SIGMOD 2015

# Advanced benchmark harness
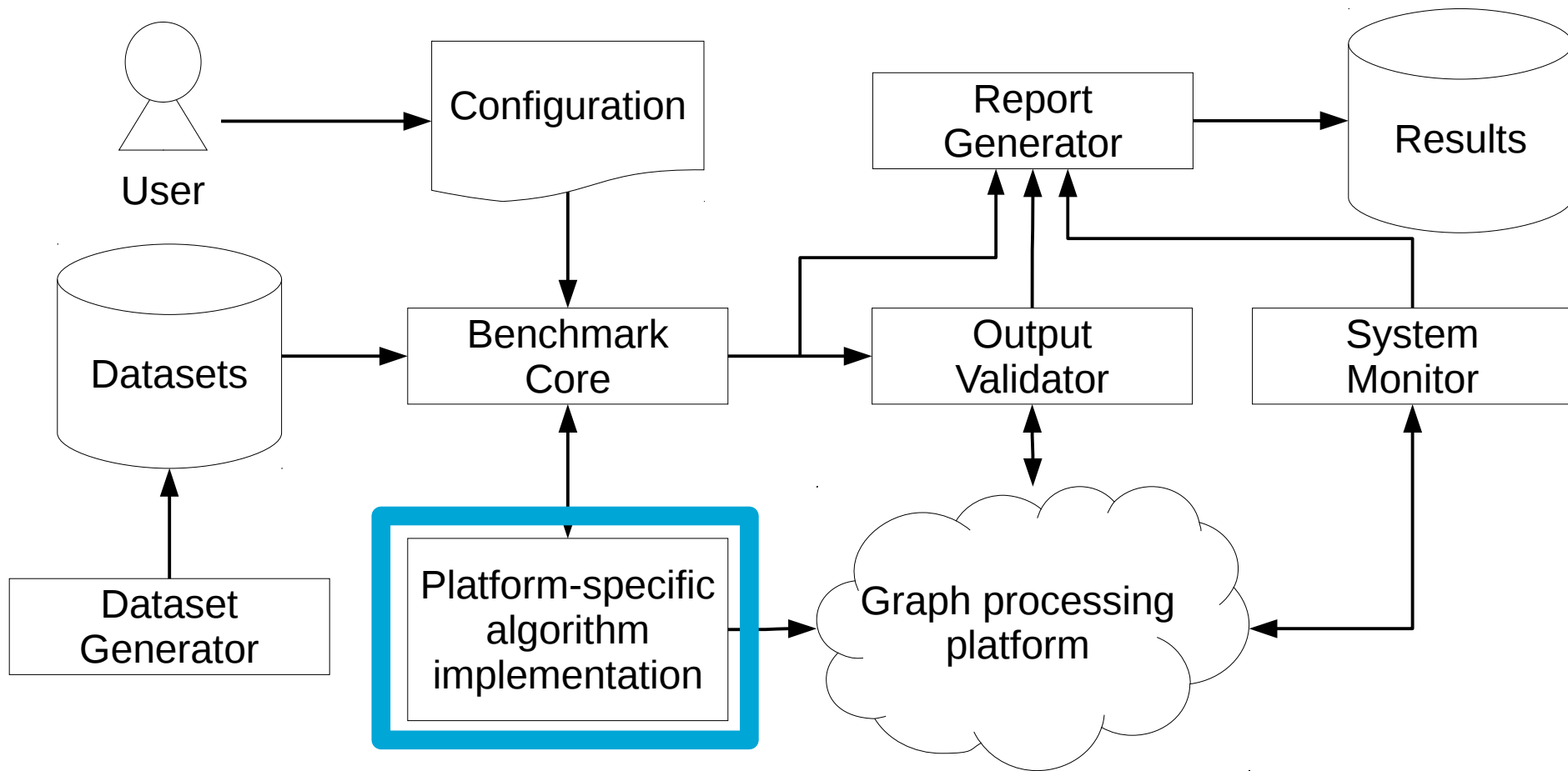
# Advanced benchmark harness

# Advanced benchmark harness



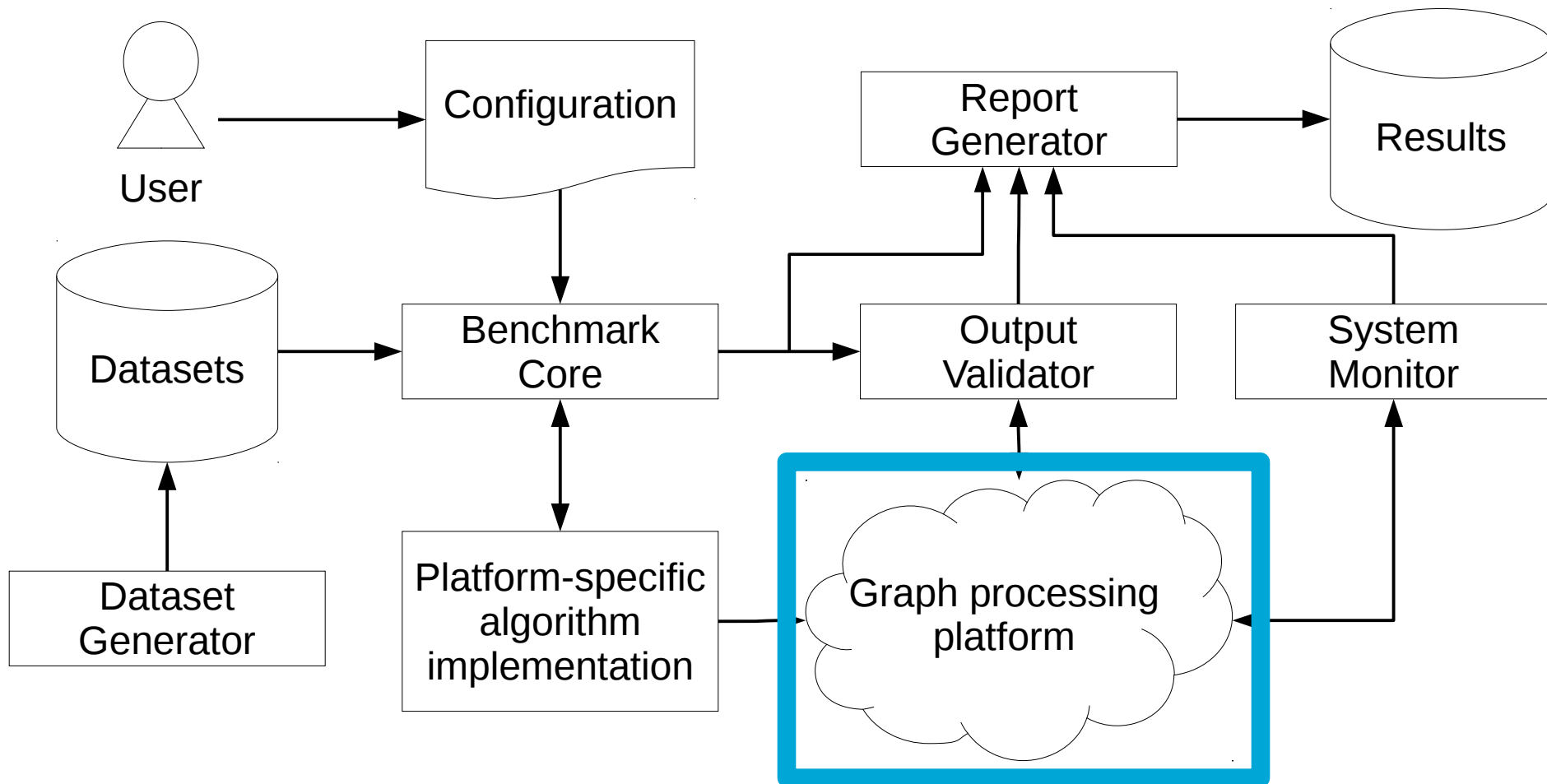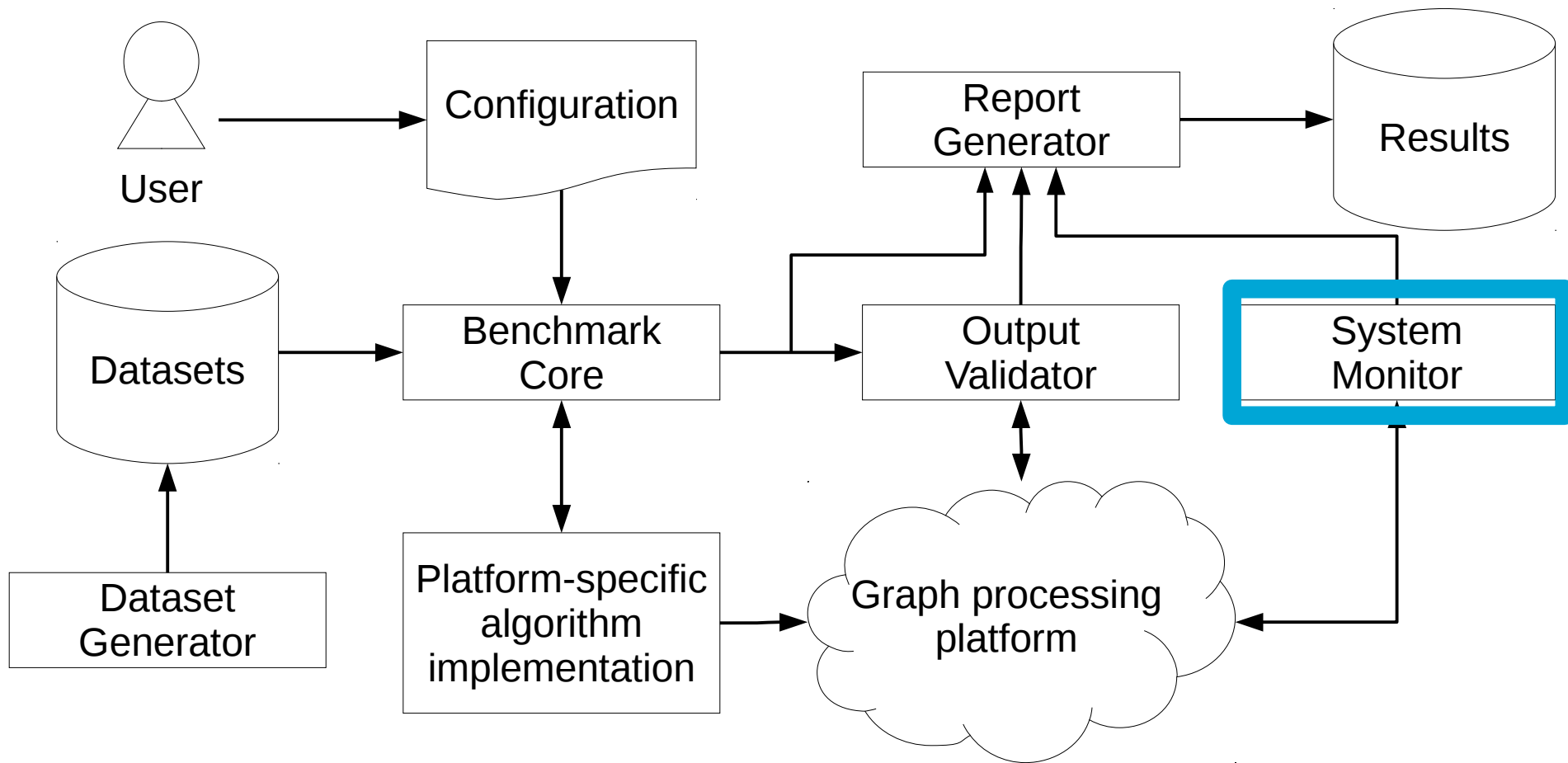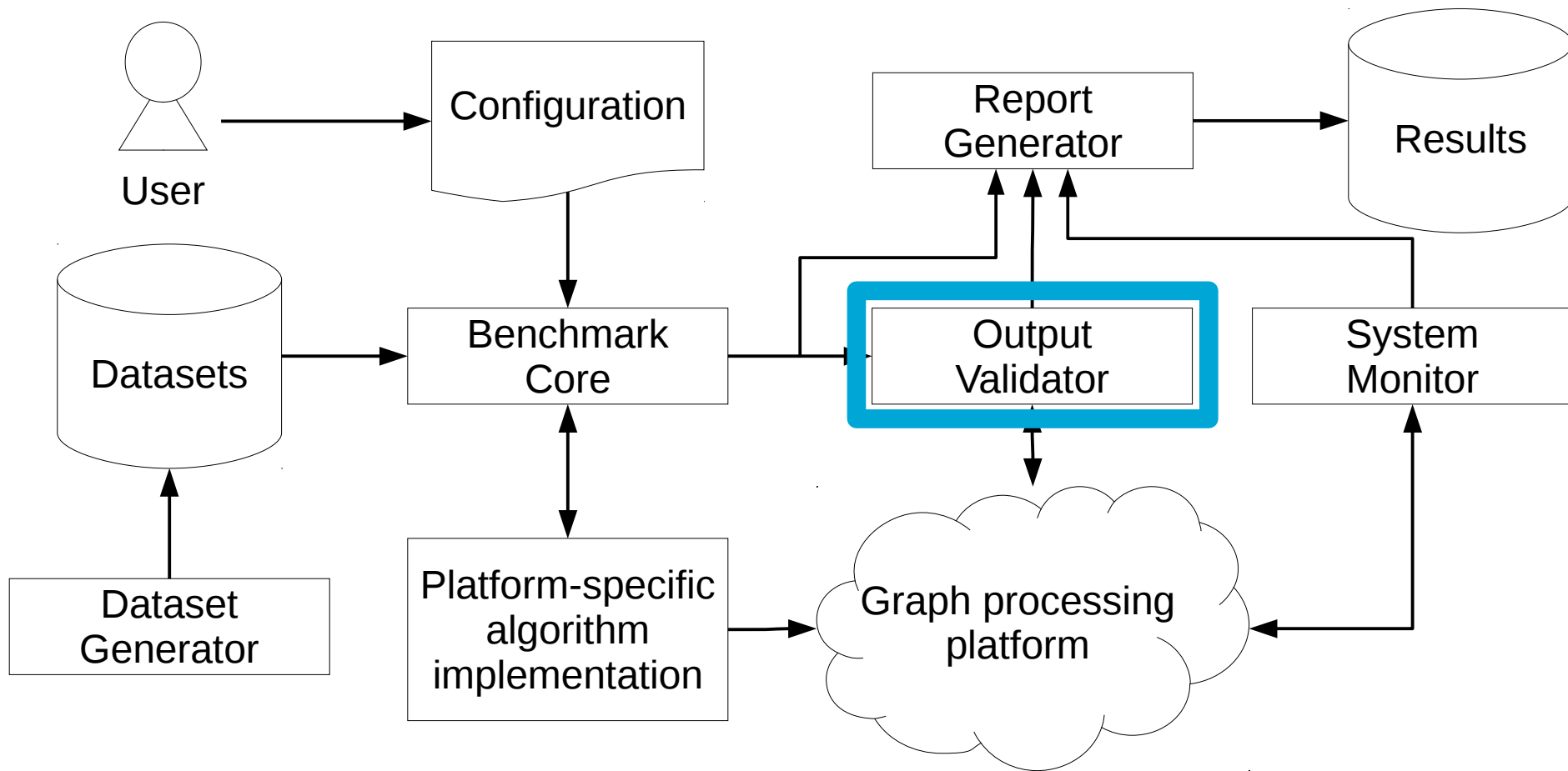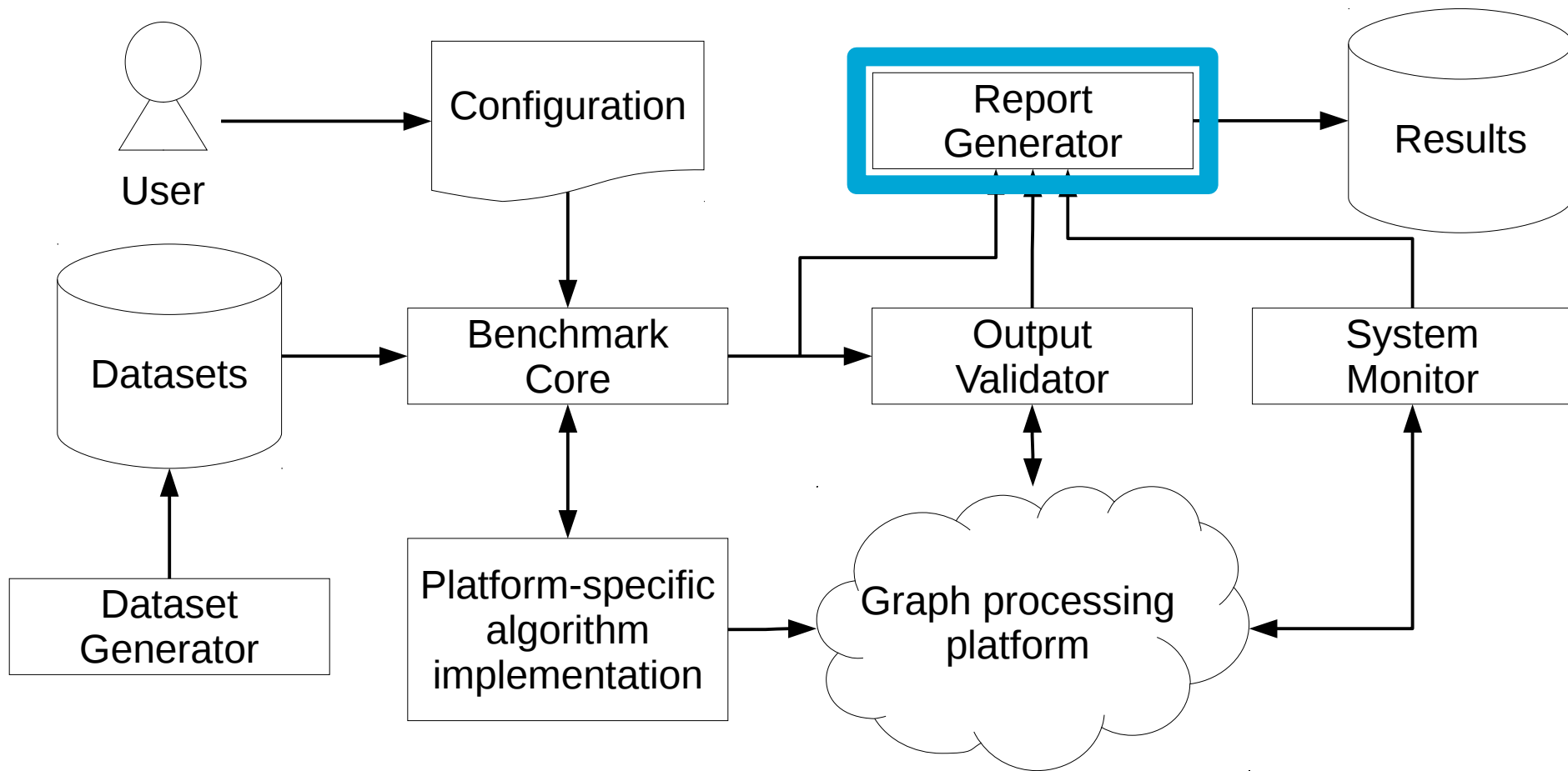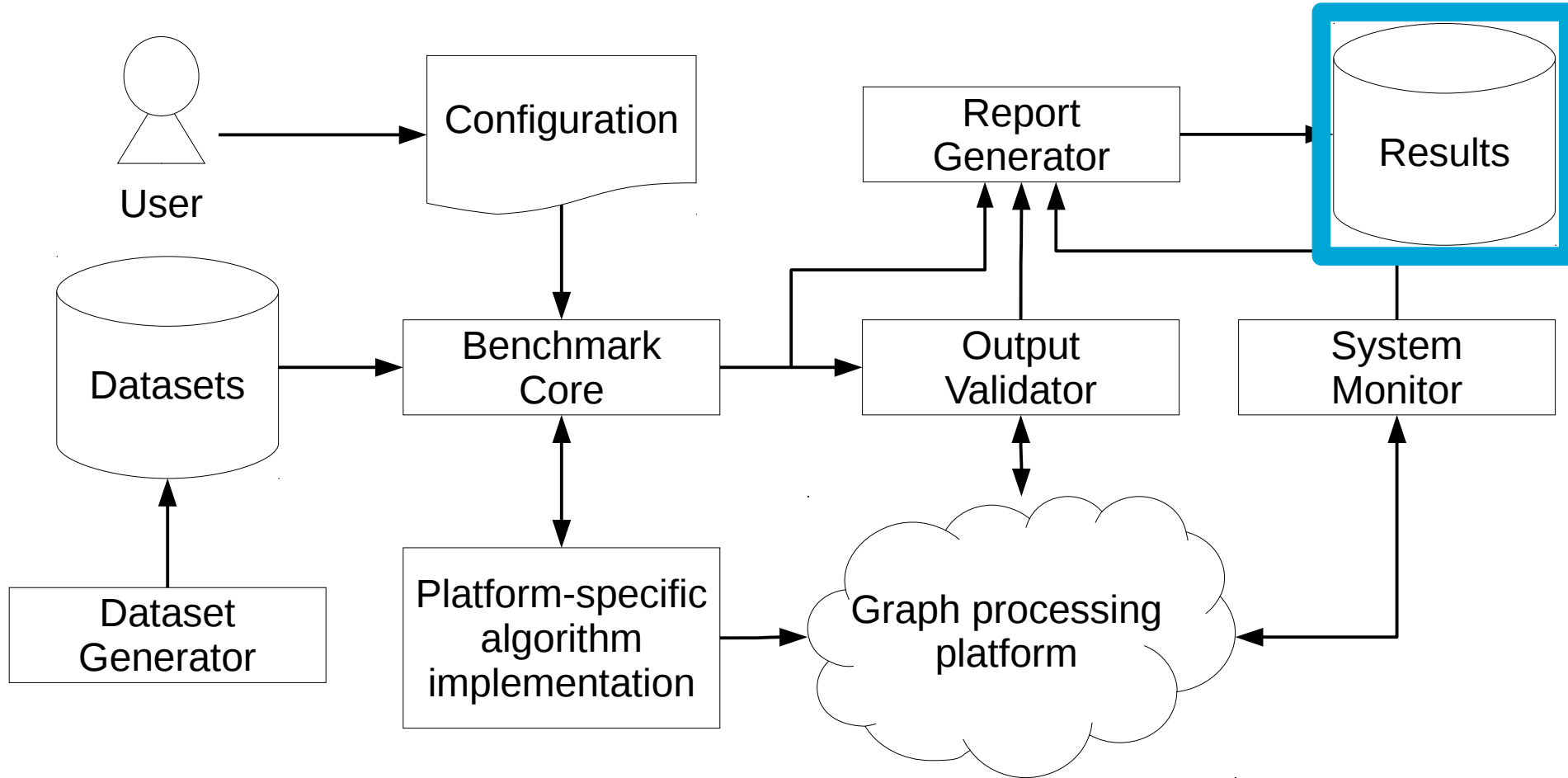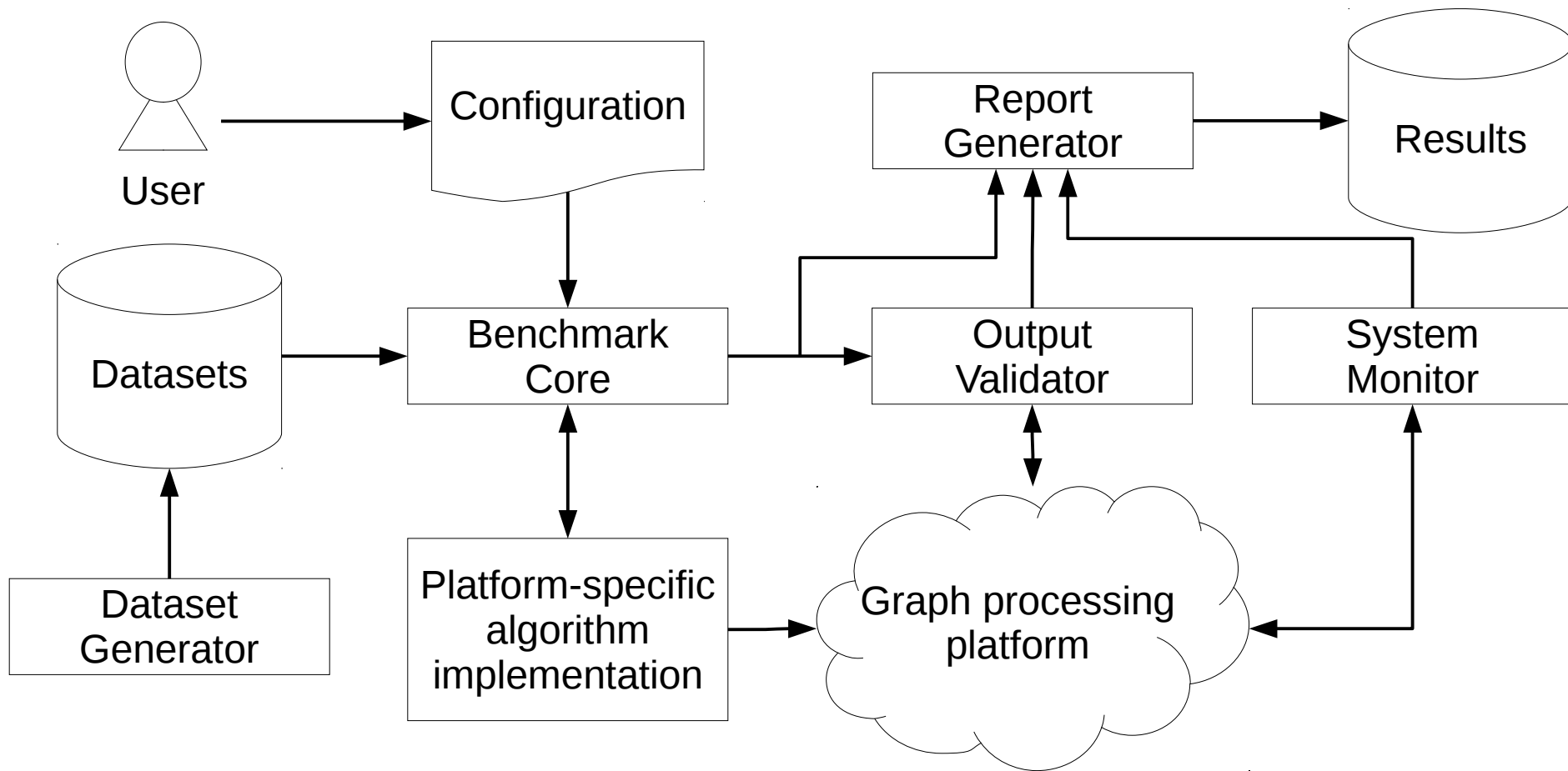User

Datasets

Dataset Generator

Configuration

Benchmark Core

Platform-specific algorithm implementation

Graph processing platform

Report Generator

Output Validator

System Monitor

Results

# Advanced benchmark harness

# Advanced benchmark harness



User

Configuration

Report Generator

Results

Datasets

Benchmark Core

Output Validator

System Monitor

Dataset Generator

Platform-specific algorithm implementation

Graph processing platform

**TU**Delft

# Advanced benchmark harness



User

Configuration

Report Generator

Results

Datasets

Benchmark Core

Output Validator

System Monitor

Dataset Generator

Platform-specific algorithm implementation

Graph processing platform

**TU**Delft

# Advanced benchmark harness



User

Datasets

Dataset Generator

Configuration

Benchmark Core

Platform-specific algorithm implementation

Report Generator

Output Validator

Graph processing platform

Results

System Monitor

**TU**Delft

# Advanced benchmark harness



User

Configuration

Report Generator → Results

Datasets → Benchmark Core → Output Validator

System Monitor

Dataset Generator → Datasets

Platform-specific algorithm implementation → Graph processing platform

**TU**Delft

# Advanced benchmark harness

# Advanced benchmark harness

# Advanced benchmark harness



User

Configuration

Report Generator

Results

Datasets

Benchmark Core

Output Validator

System Monitor

Dataset Generator

Platform-specific algorithm implementation

Graph processing platform

**TU**Delft

# Advanced benchmark harness



User

Configuration

Report Generator

Results

Datasets

Benchmark Core

Output Validator

System Monitor

Dataset Generator

Platform-specific algorithm implementation

Graph processing platform

**TU**Delft

# Supported algorithms

| ID | Algorithm | Class | Use (%) |
|---|---|---|---|
| BFS | Breadth-first search | Traversal | 46 |
| STATS | Local clustering coefficient | Statistics | 16 |
| CONN | Weakly connected components | Connected components | 13 |
| CD | Label propagation | Community detection | 5 |
| EVO | Forest fire evolution | Evolution | 4 |

**TU**Delft

Guo et al., How Well Do Graph-Processing Platforms Perform?
An Empirical Performance Evaluation and Analysis, IPDPS 2014

# Experimental setup

- DAS-4 cluster
  - Typical big data setup
  - 11 nodes, 24 GiB RAM, 2 x 8-core Xeon E5620
  - 1 Gbit/s Ethernet
- Single machine
  - HPC-like setup
  - 192 GiB RAM, 2 x 8-core Xeon E5-2450 v2

**T**UDelft

# Runtime

# Runtime

# Runtime

# Runtime

# Runtime

# Runtime

# Runtime

# Runtime
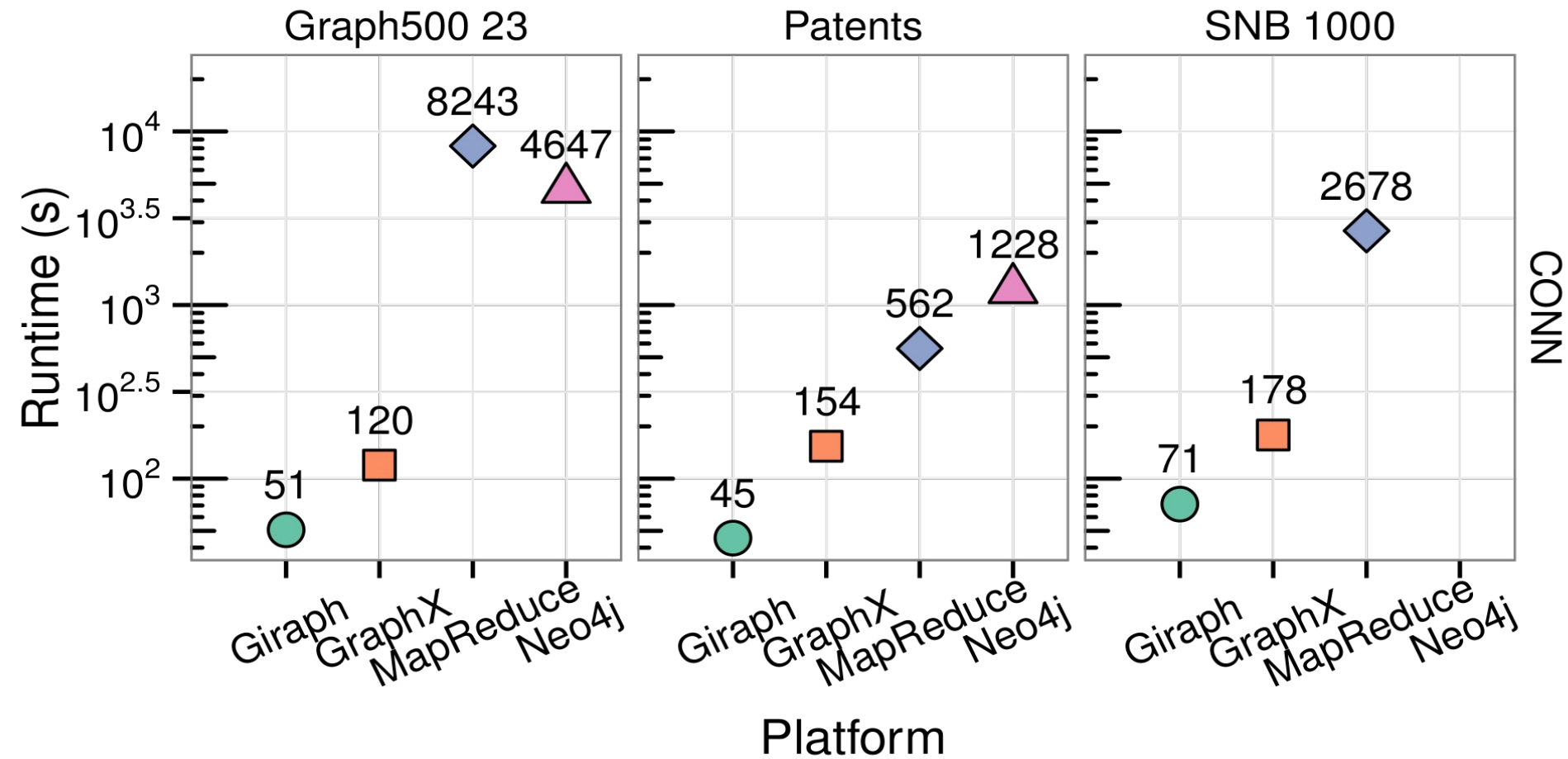


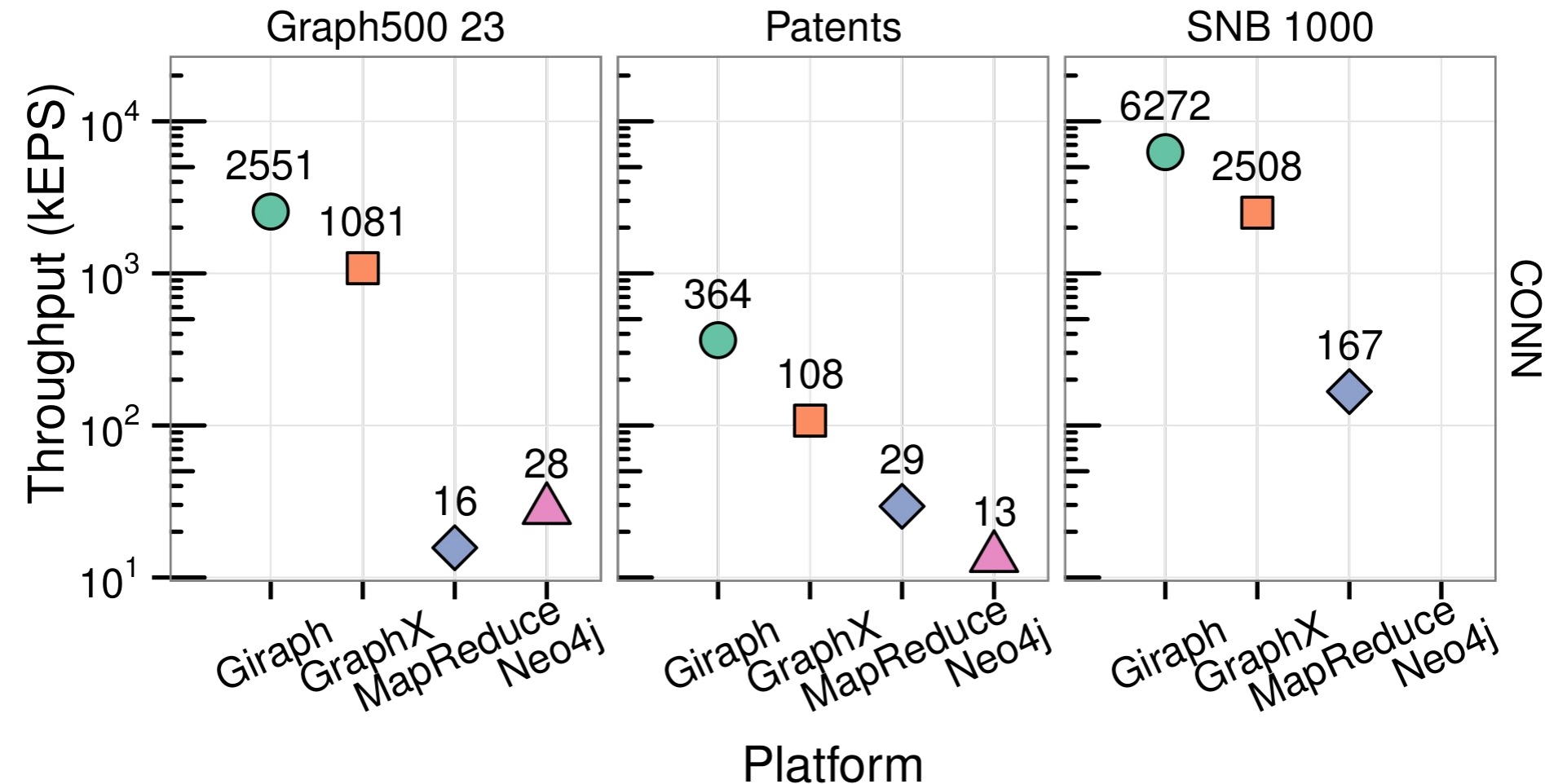Graph500 23     Patents     SNB 1000

**MR 2x faster than Neo4j**

TUDelft

61

# Runtime

# Runtime

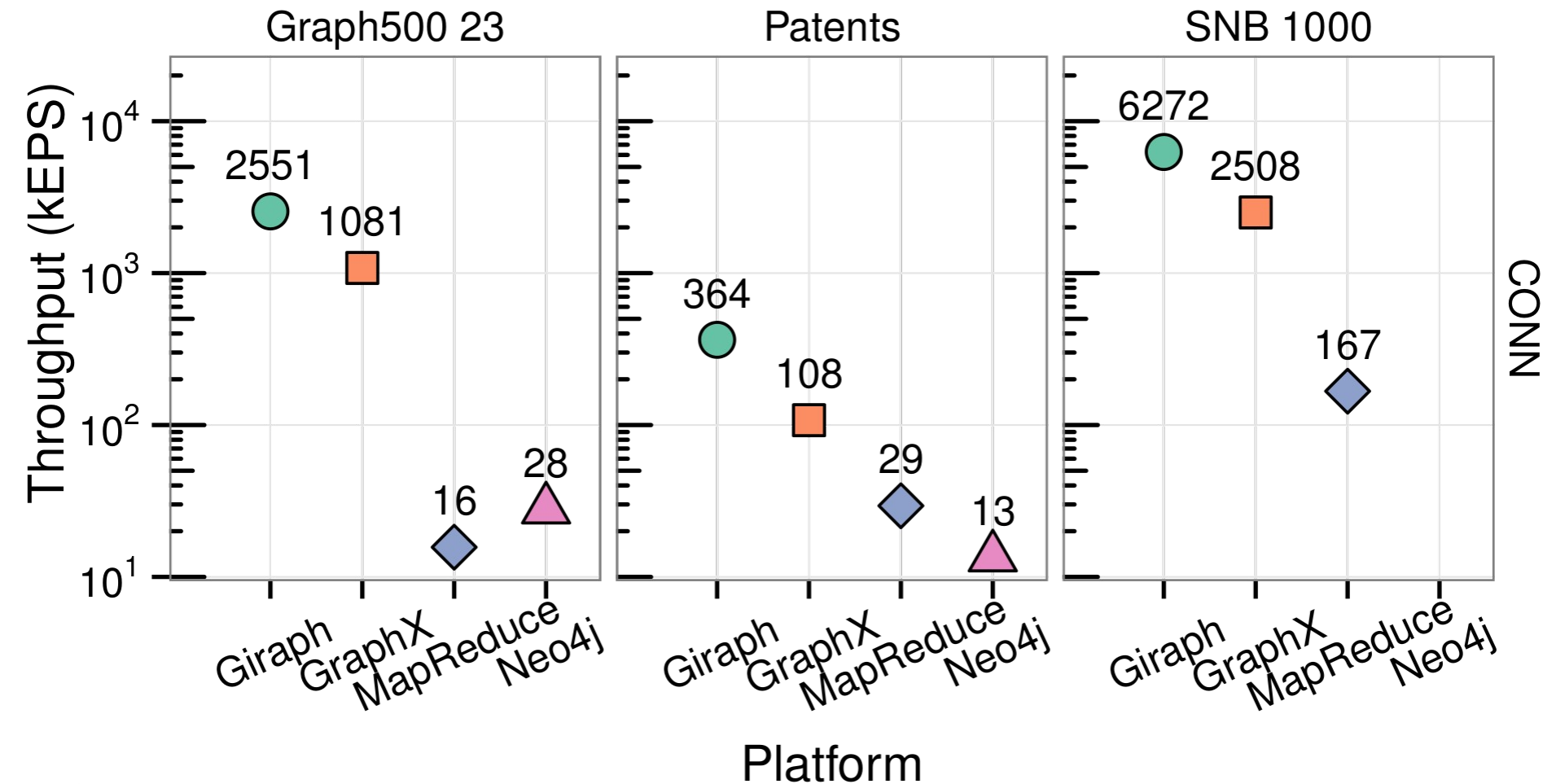# Edge-normalized performance

# Edge-normalized performance



Graph500 23 — Patents — SNB 1000

Throughput (kEPS) vs Platform (Giraph, GraphX, MapReduce, Neo4j)

Graph500 23: 2551, 1081, 16, 28
Patents: 29, 13
SNB 1000: 6272, 2508, 167

Number of edges in graph / runtime

**TU**Delft

# Edge-normalized performance



Graph500 23 | Patents | SNB 1000

Throughput (kEPS)

CONN

2551    081

6272   2508

364    108

167

20 x difference
because of dataset structure

Giraph  GraphX  MapReduce  Neo4j

Platform

*TU*Delft

66

# Edge-normalized performance

# Edge-normalized performance

See paper for more results

# Datagen scalability

# Conclusion

**TU**Delft

# Conclusion

- Use Graphalytics to:
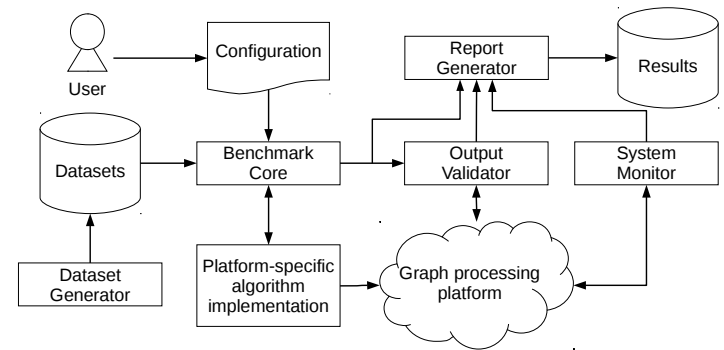  - Compare
  - Tune
  - Re-design

# Conclusion

- Use Graphalytics to:
  - Compare
  - Tune
  - Re-design



- Open source (Apache License 2.0)
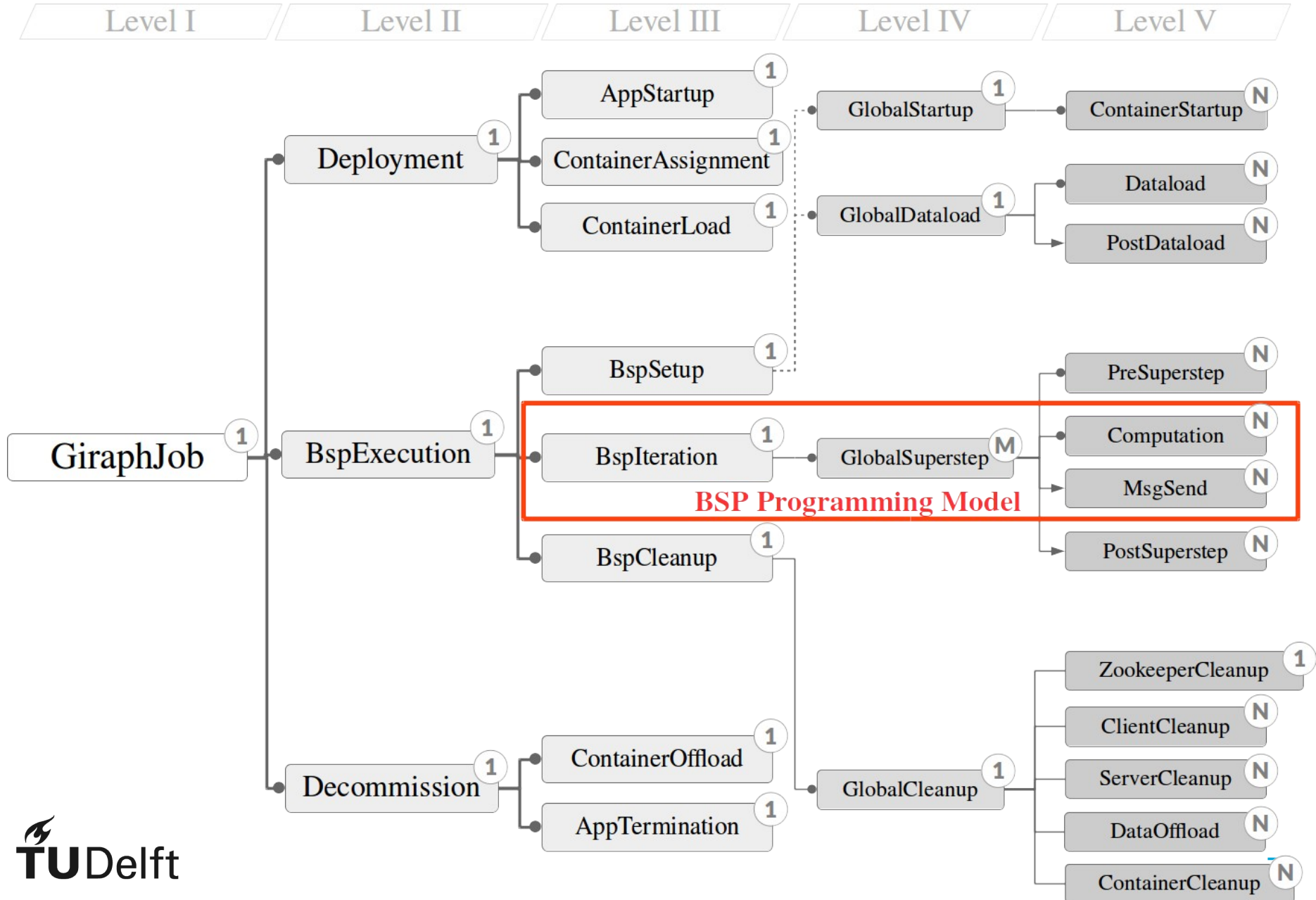  - Contribute an implementation for your platform

**T̃U**Delft

# Conclusion

- Use Graphalytics to:
  - Compare
  - Tune
  - Re-design

- Open source (Apache License 2.0)
  - Contribute an implementation for your platform
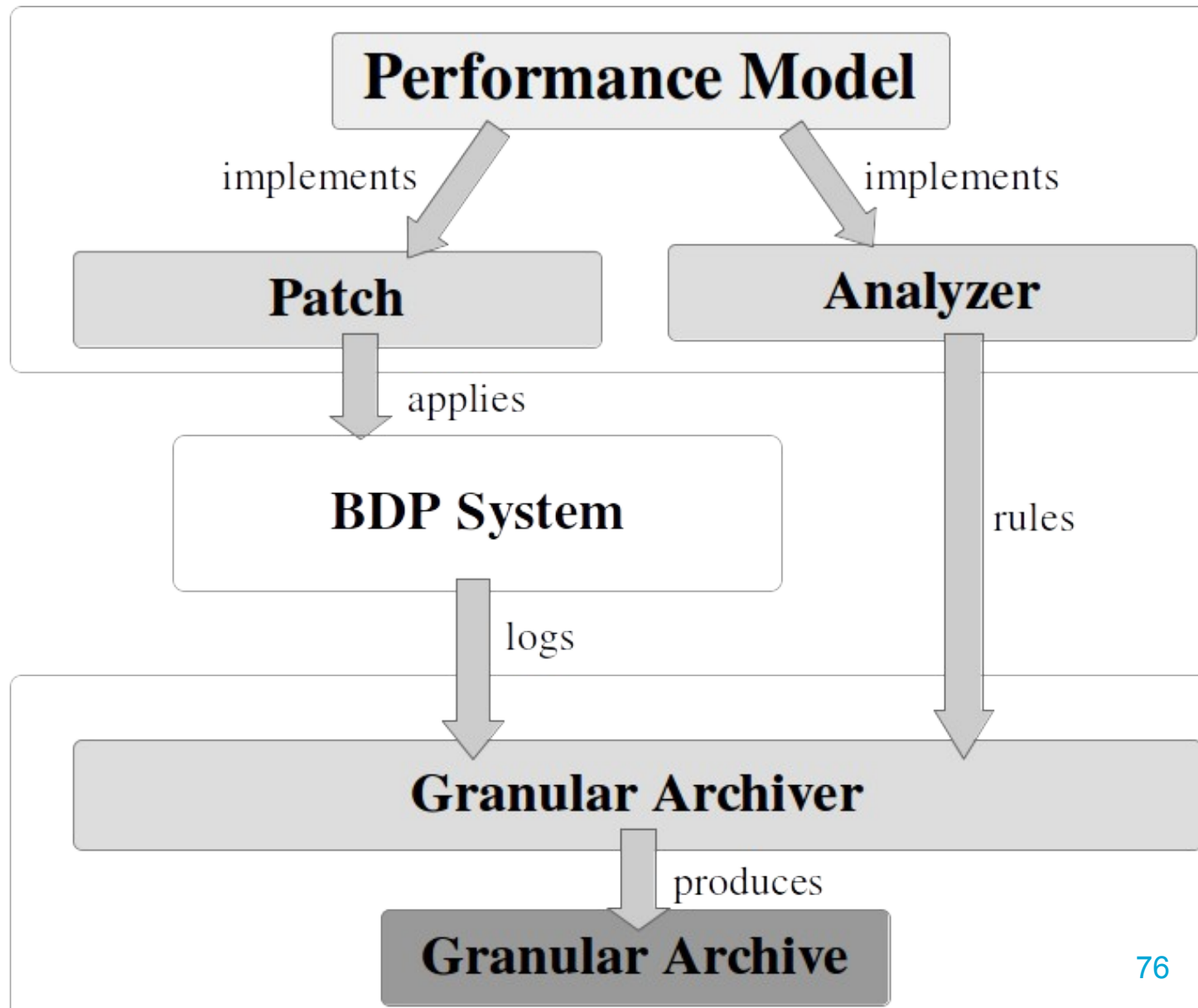
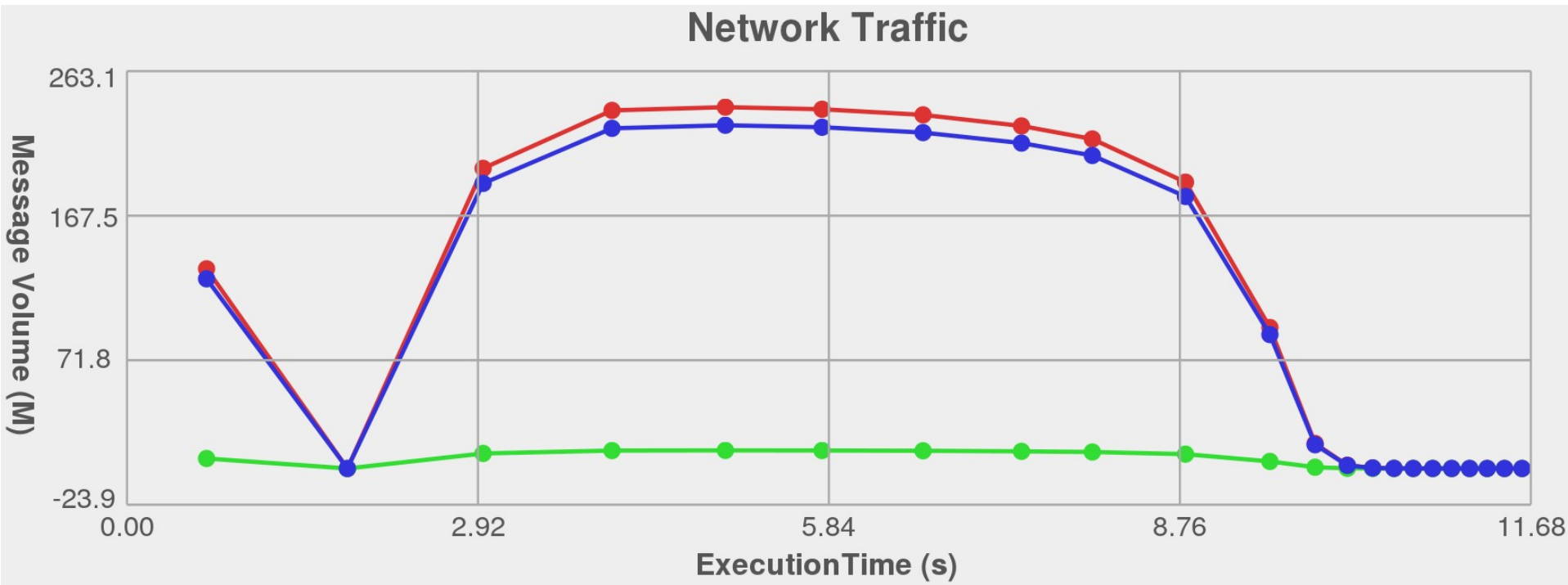- SPEC RG standardization coming

# Performance evaluation – Granular

# Deploying Granular

# Granular results

# Thank you!

mihai@mihaic.ro

graphalytics.ewi.tudelft.nl

# GitHub

github.com/tudelft-atlarge/graphalytics

TU Delft