



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

ALOJA: Cost-effectiveness study of Hadoop deployments

Nicolas Poggi, Senior Researcher



EXCELENCIA
SEVERO
OCHOA

BSC~Microsoft Research
Centre

December 2014



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

INSTITUTIONAL ABOUT THE PROJECT

Barcelona Supercomputing Center (BSC)

- ❧ 22 year history in Computer Architecture research
 - European Center for Parallelism of Barcelona (CEPBA)
 - Based at the Technical University of Catalonia (UPC) in 1991
 - Long track record with IBM in chip Architecture & Parallelism
- ❧ Led by Mateo Valero
 - ACM fellow, Eckert-Mauchly award 2007, Goode award 2009
 - Active research staff with more than 1000 publications
 - Large ongoing life science computational projects
 - Computational Genomics
 - Molecular modeling & Bioinformatics
 - Protein Interactions & Docking
 - In place computational capabilities
 - Mare Nostrum Super Computer



- ❧ Prominent body of research activity around Hadoop since 2008
 - Previous to ALOJA
 - SLA-driven scheduling (Adaptive Scheduler), in memory caching, etc.
 - Group page: <http://www.bsc.es/autonomic>

BSC-MSRS Centre and ALOJA



- ⌘ Long-term relationship between
 - BSC and Microsoft Research and Microsoft product teams
- ⌘ Previous research at the intersection of computer architecture, language implementation, and systems software, and performance profiling
- ⌘ Open model:
 - **No patents, public IP, publications and open source main focus**
 - 87 publications, 4 Best paper awards

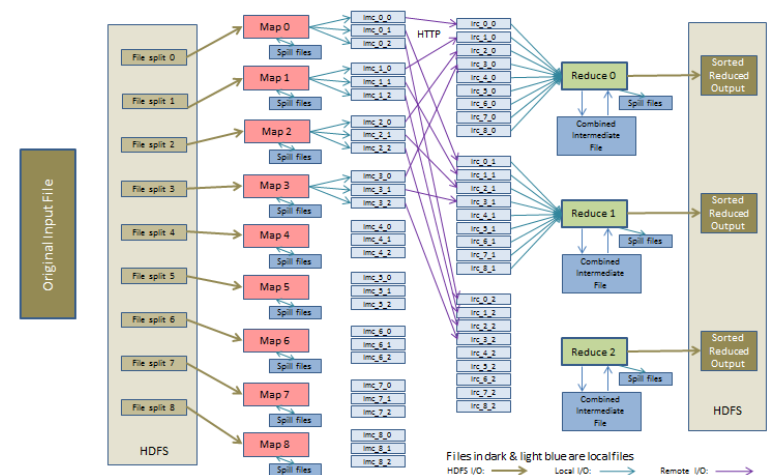
- ⌘ **ALOJA** is the latest phase of the engagement
- ⌘ With intent to explore:
 - upcoming hardware architectures
 - and building automated mechanism
- ⌘ for deploying cost-effective Hadoop clusters.



Motivation

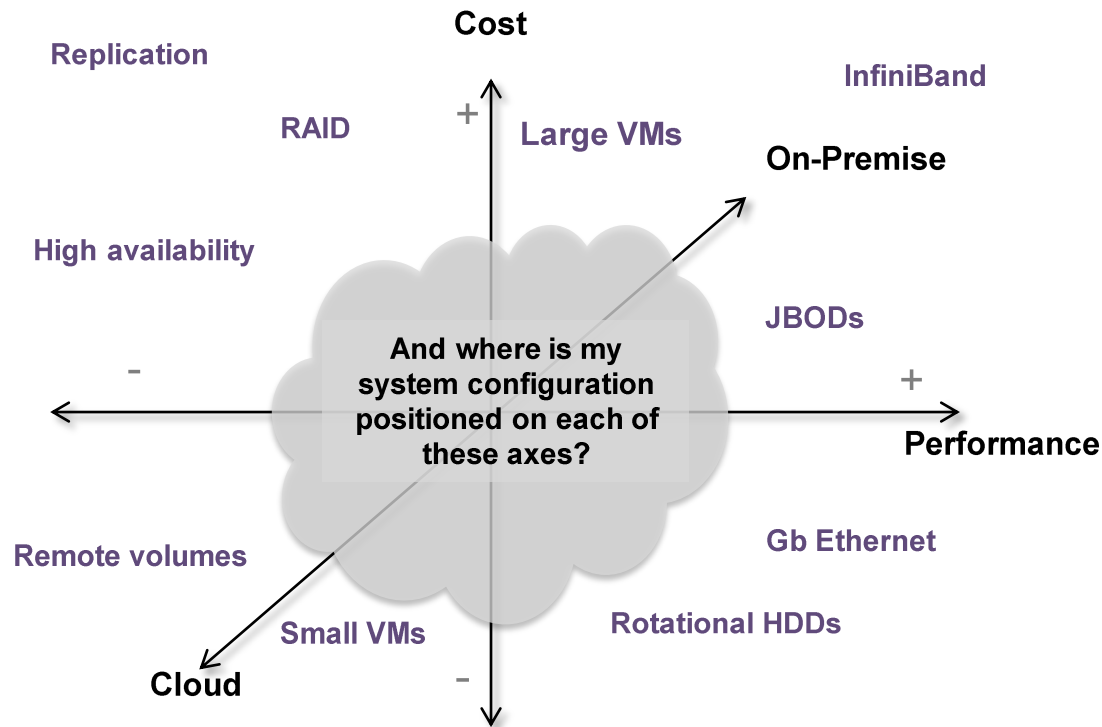
- ❧ The Hadoop framework implements a complex distributed execution model
 - Over 100 interrelated config parameters
 - Requires manual iterative benchmarking and tuning
- ❧ Early results show that Hadoop's price/performance
 - are affected by relatively simple SW $> 3x$
 - and HW configuration choices $> 3x$
- ❧ Commodity HW no longer low-end
 - new affordable hardware from original design (ie., SSDs)
 - Hadoop performs poorly on scale-up
 - or low power HW
- ❧ New Cloud services for Hadoop
 - IaaS and PaaS
 - Direct vs. remote attached volumes
- ❧ Spread Hadoop ecosystem
 - Dominated by vendors
 - Lack of verifiable benchmarks

Hadoop – Map / Reduce – Overview of I/O flows

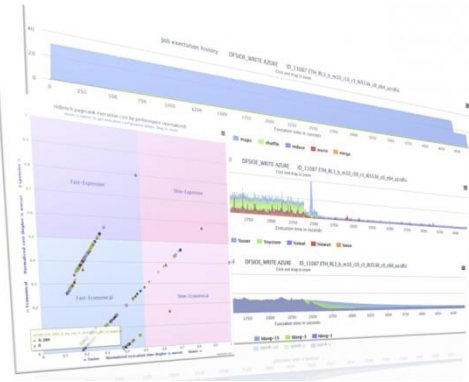


Current scenario

« What is the most cost-effective configuration for my needs?



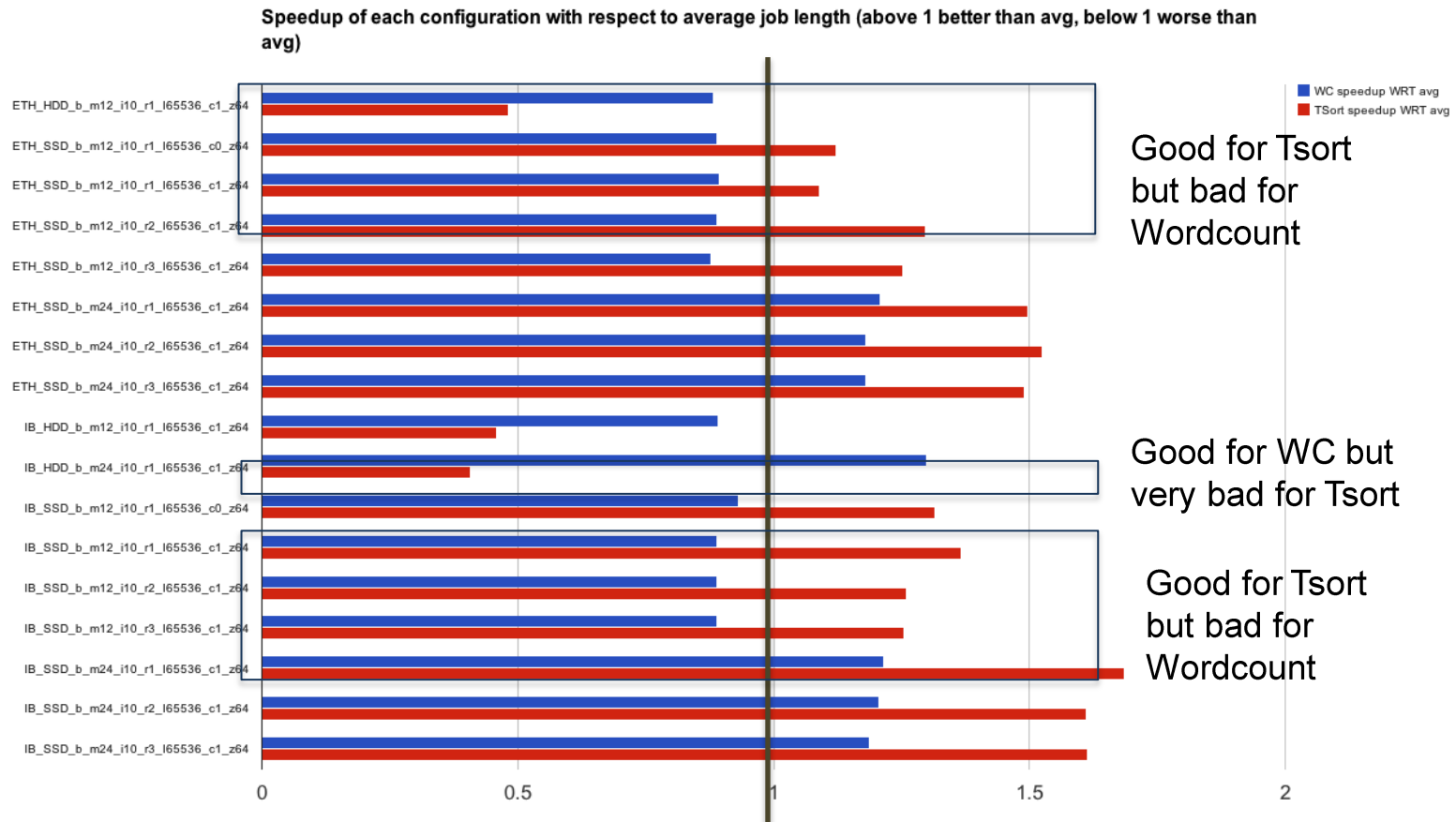
Project ALOJA



- ⌘ Initiative to produce mechanisms for an **automated characterization of cost-effectiveness** of Big Data deployments with focus on Hadoop
- ⌘ Results from of a growing need of the community to understand job execution
- ⌘ Explore different configuration deployment options and their tradeoffs
 - Both Software and Hardware
 - Cloud and on-premise
 - High-end, mid range, low-end
- ⌘ Seeks to provide both **knowledge**, tools, and an online service
 - to with which users make better informed decisions
 - and reduce the TCO for their Big Data infrastructures

Early exploration: Job resource configs

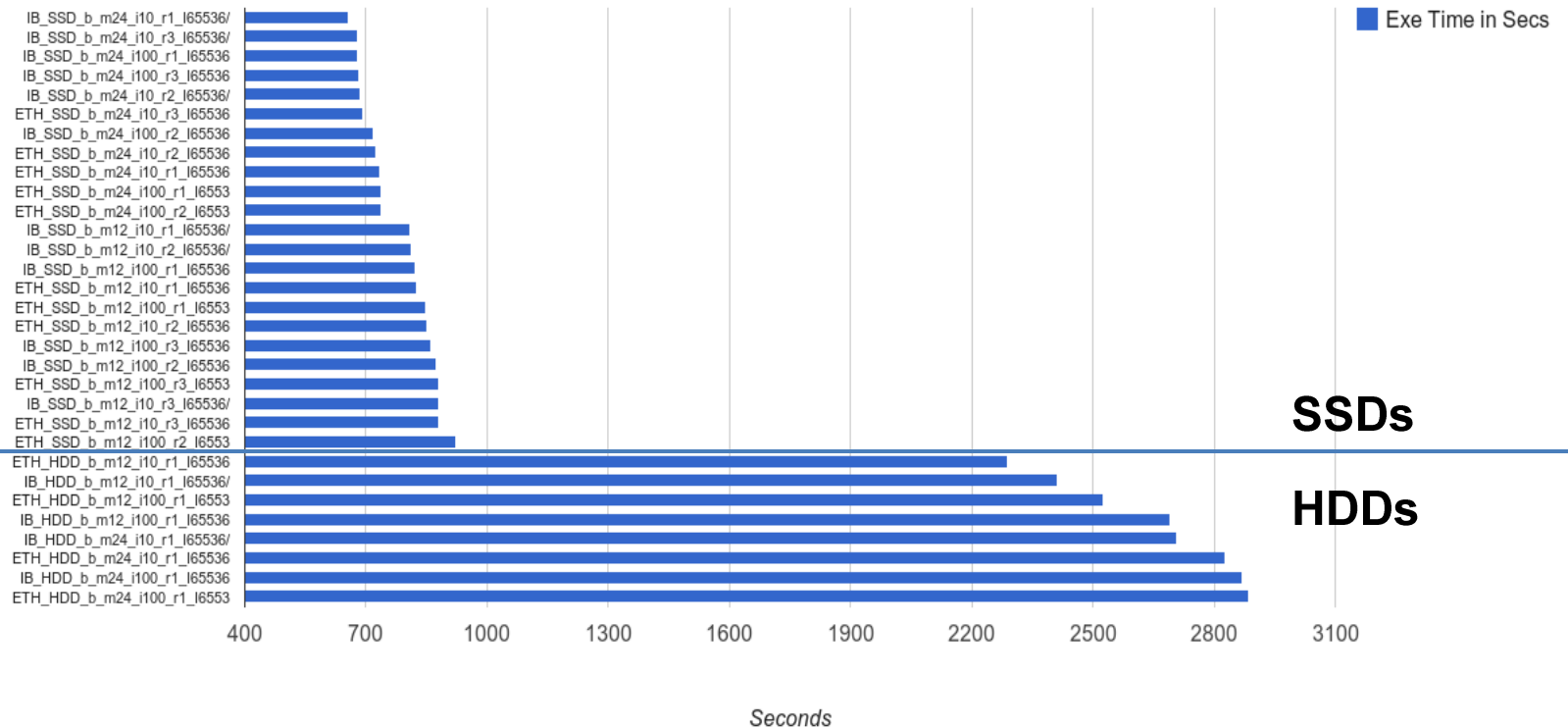
Is there one software configuration iteration that fits everybody?



Vertical line: Average performance for this workload across configurations
Values to the right: above average
Values to the left: below average

Early exploration: HW technology impact

Terasort runs





**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

PLATFORM

ALOJA Platform

⌘ Benchmarking, Repository, and Analytics tools for Big Data



⌘ Composed of open-source

- benchmarking and configuration management tools,
- high-level system performance metric collection,
- low-level Hadoop instrumentation based on **BSC Tools**
- and Web based data analytics tools

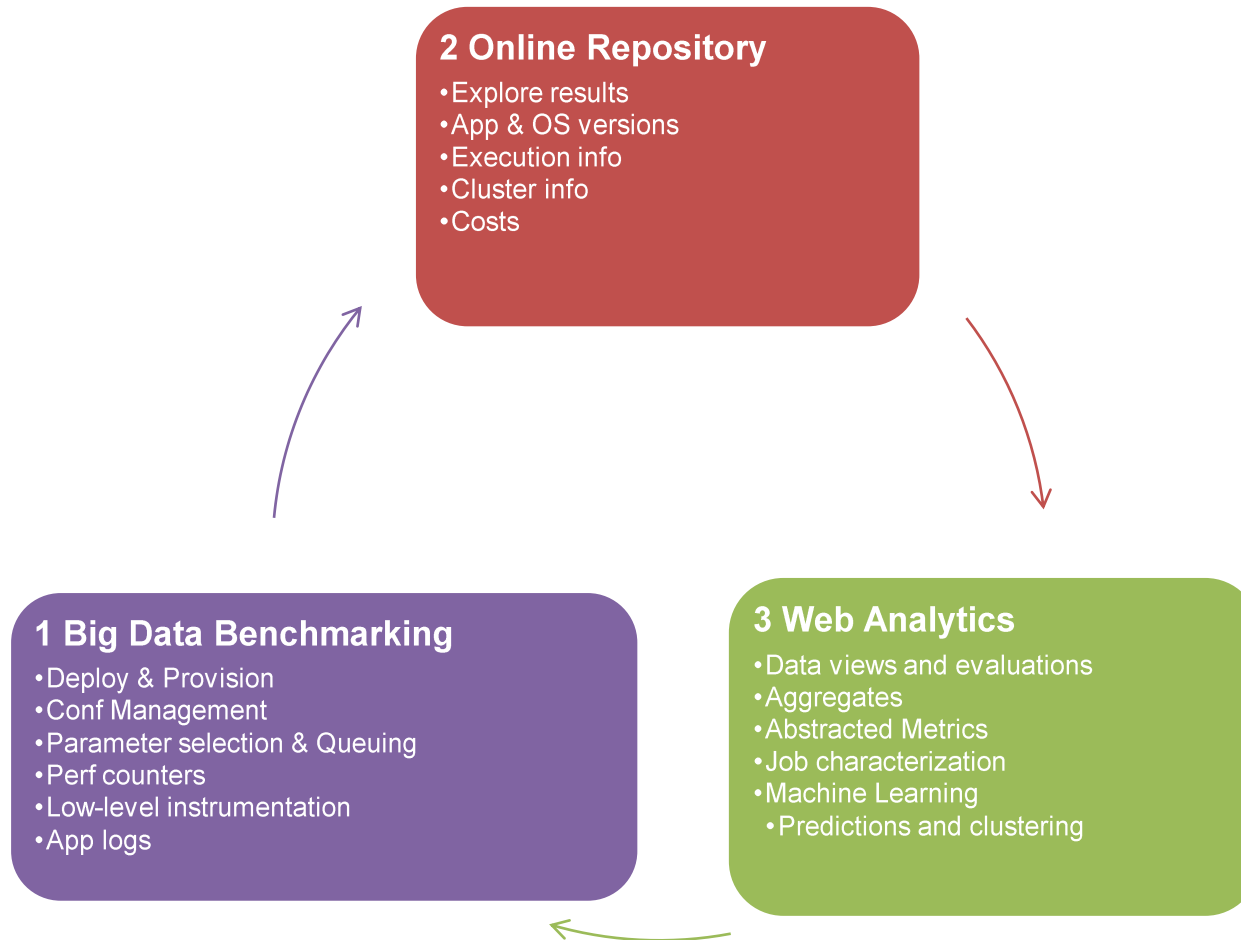
⌘ Online Big Data Benchmark repository of:

- + 5000 runs (from HiBench)
- Sharable, comparable, repeatable, verifiable executions

⌘ Not reinventing the wheel, but

- most current tools designed for production, not for benchmarking
- leverages current compatible tools and projects

ALOJA main components



ALOJA-WEB

⌋ <http://hadoop.bsc.es>

⌋ Entry point for explore the results collected from the executions, providing insights on the obtained results through continuously evolving data views.

HiBench Executions on Hadoop BSC - Microsoft Research Centre

Navigation: **HiBench Runs Details** | Hadoop Job Counters | Cost Evaluation | Performance Charts

Click on a **benchmark name** to see execution details.
Select different rows and click **compare**, to compare charts.
Search to filter results. Shift+Click to order by multiple columns

Show / hide columns

Show **10** entries

Filter all columns:

ID	Benchmark	Exe Time	Running Cost \$	Net	Disk	Maps	IO SFac	Rep	IO FBuf	Comp	Blk size	Cluster	Files	PARAVER	
<input type="checkbox"/>	20372	dfsio_read	2990	5.81	ETH	RL1	8	10	2	65536	3	32	Azure L	files	PRV .ZIP
<input checked="" type="checkbox"/>	20371	pagerank	2809	5.46	ETH	RL1	8	10	2	65536	3	32	Azure L	files	PRV .ZIP
<input type="checkbox"/>	20370	sort	657	1.28	ETH	RL1	8	10	2	65536	3	32	Azure L	files	PRV .ZIP
<input type="checkbox"/>	20369	wordcount	1336	2.60	ETH	RL1	8	10	2	65536	3	32	Azure L	files	PRV .ZIP
<input type="checkbox"/>	20368	kmeans	3002	5.84	ETH	RL1	8	10	2	65536	3	32	Azure L	files	PRV .ZIP
<input type="checkbox"/>	20367	dfsio_write	2640	5.13	ETH	RL1	8	10	2	65536	3	64	Azure L	files	PRV .ZIP
<input type="checkbox"/>	20366	dfsio_read	3139	6.10	ETH	RL1	8	10	2	65536	3	64	Azure L	files	PRV .ZIP
<input checked="" type="checkbox"/>	20365	pagerank	2612	5.08	ETH	RL1	8	10	2	65536	3	64	Azure L	files	PRV .ZIP
<input type="checkbox"/>	20364	sort	627	1.22	ETH	RL1	8	10	2	65536	3	64	Azure L	files	PRV .ZIP
<input type="checkbox"/>	20363	wordcount	1279	2.49	ETH	RL1	8	10	2	65536	3	64	Azure L	files	PRV .ZIP

Showing 1 to 10 of 3,696 entries (filtered from 4,019 total entries)

Compare executions:

Select rows by clicking on checkboxes and click: **Compare Executions**

Job execution history: PAGERANK AZURE

map shuffle reduce waste merge



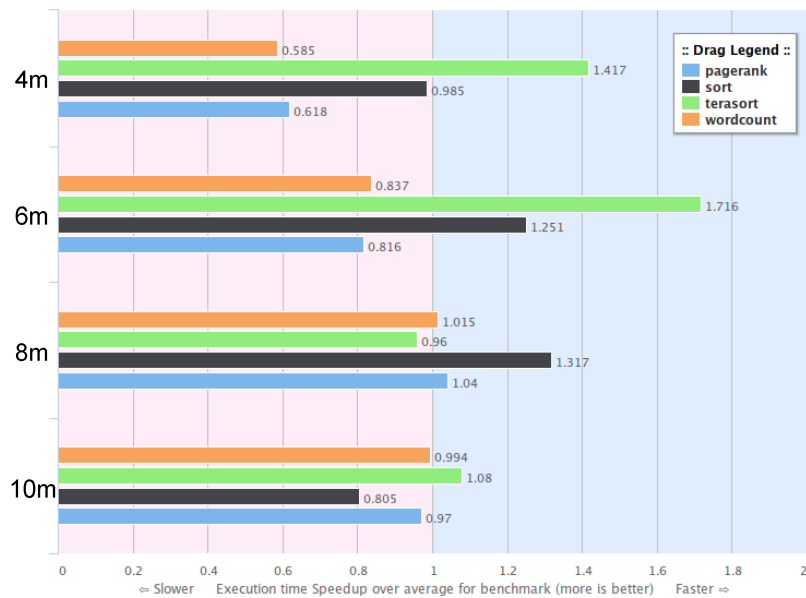
**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

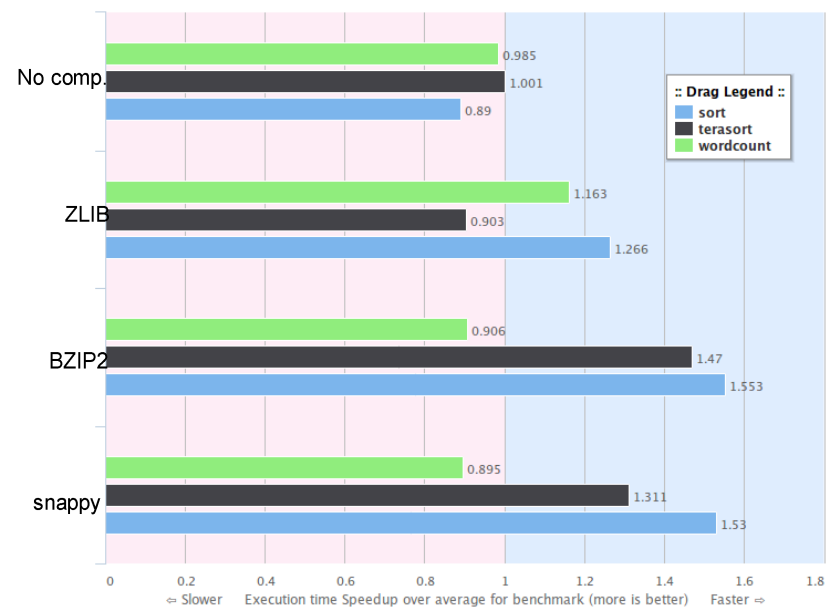
ONLINE DEMO

Impact of SW configurations in Speedup

Number of mappers

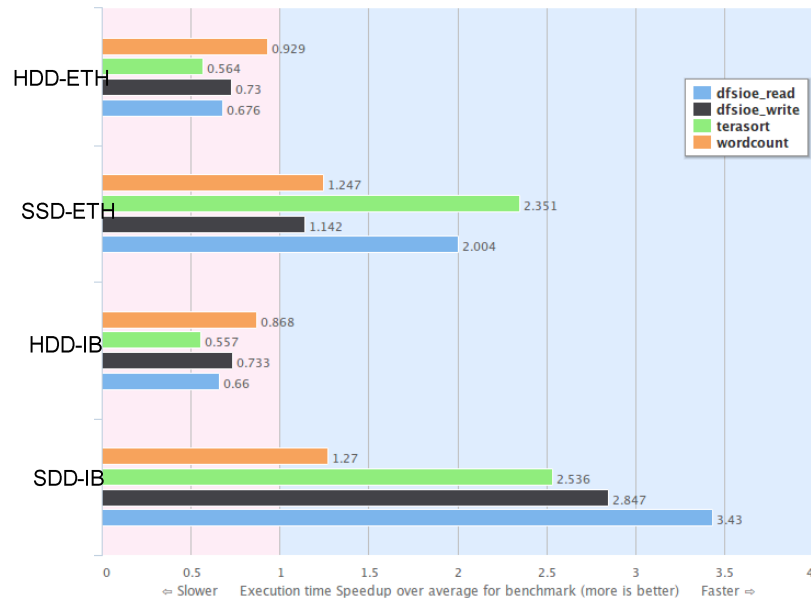


Compression algo.

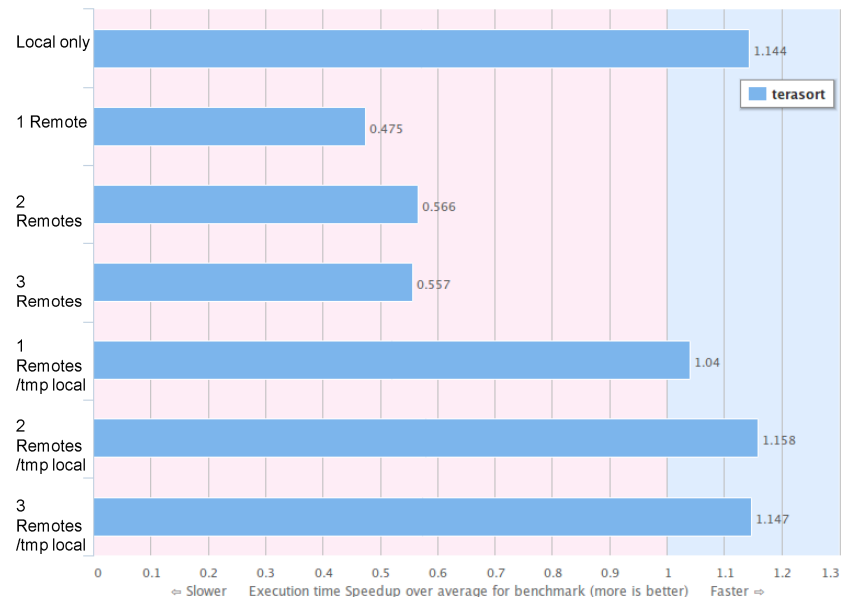


Impact of HW configurations in Speedup

Disks and Network



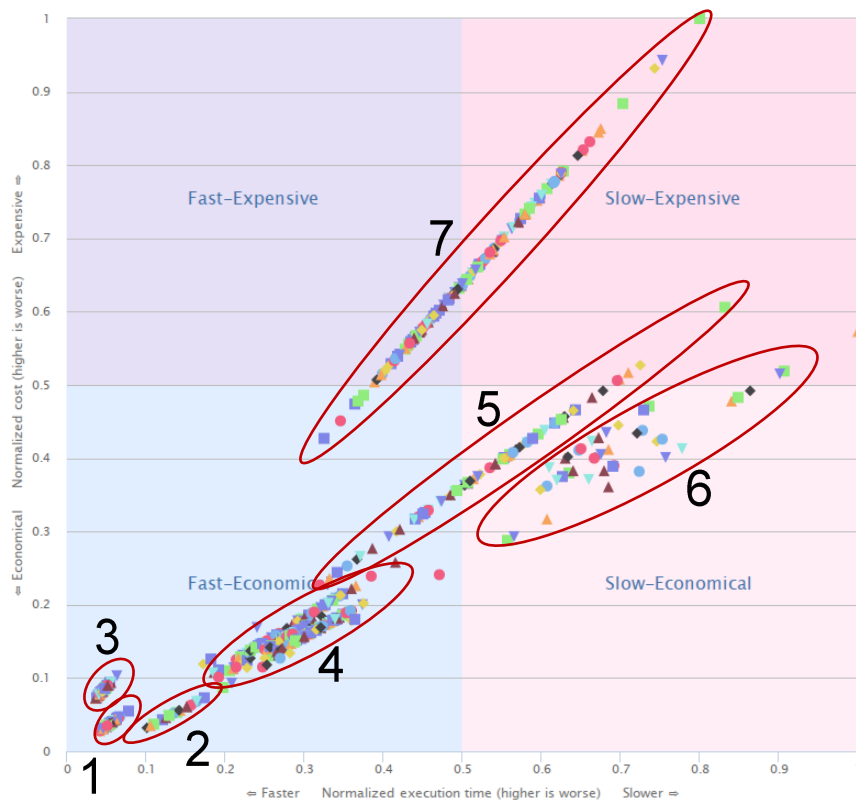
Cloud remote volumes



Cost-effectiveness of SW and HW

Point (0,0) represents most cost-effective execution

Terasort



- 1) On-premise cluster: SSD disks + GbE.
- 2) Azure IaaS: Only local disk, virtualized SSD and GbE (baseline).
- 3) On-premise cluster: SSD disks + InfiniBand.

Wordcount



- 4) Azure IaaS: 1-3 remote vol. and Hadoop /tmp to local disk (SSD) and GbE
- 5) On-premise cluster: 1 SATA disk + GbE.
- 6) Azure IaaS: 1-3 remote volumes (Blob storage).
- 7) On-premise cluster: 1 SATA disk + InfiniBand.

Initial testing infrastructure

High-End Cluster:

- 4 nodes, 12 real cores, 64GB RAM, 6x SSDs , 56Gb InfiniBand, 4Gb GbE (bonding)

Mid-end Cluster *:

- 18 nodes, 8 real cores, 24GB RAM, 1x SSD, 2x HDDs, 1Gb GbE

Cloud IaaS (Azure)

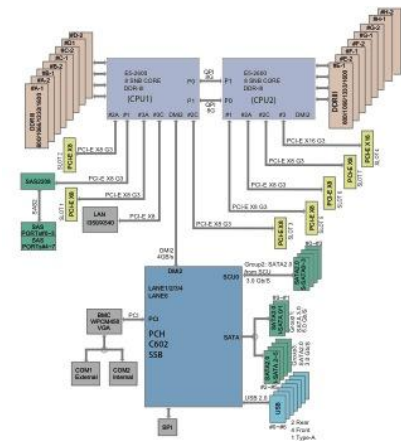
- 4 node A7s, 9 node A7s*, 9 node A3s*

Cloud PaaS (HDInsight)

- 4, 8, 16, 32 data nodes*

Low-powered cluster:

- 10-node ARM based cluster*





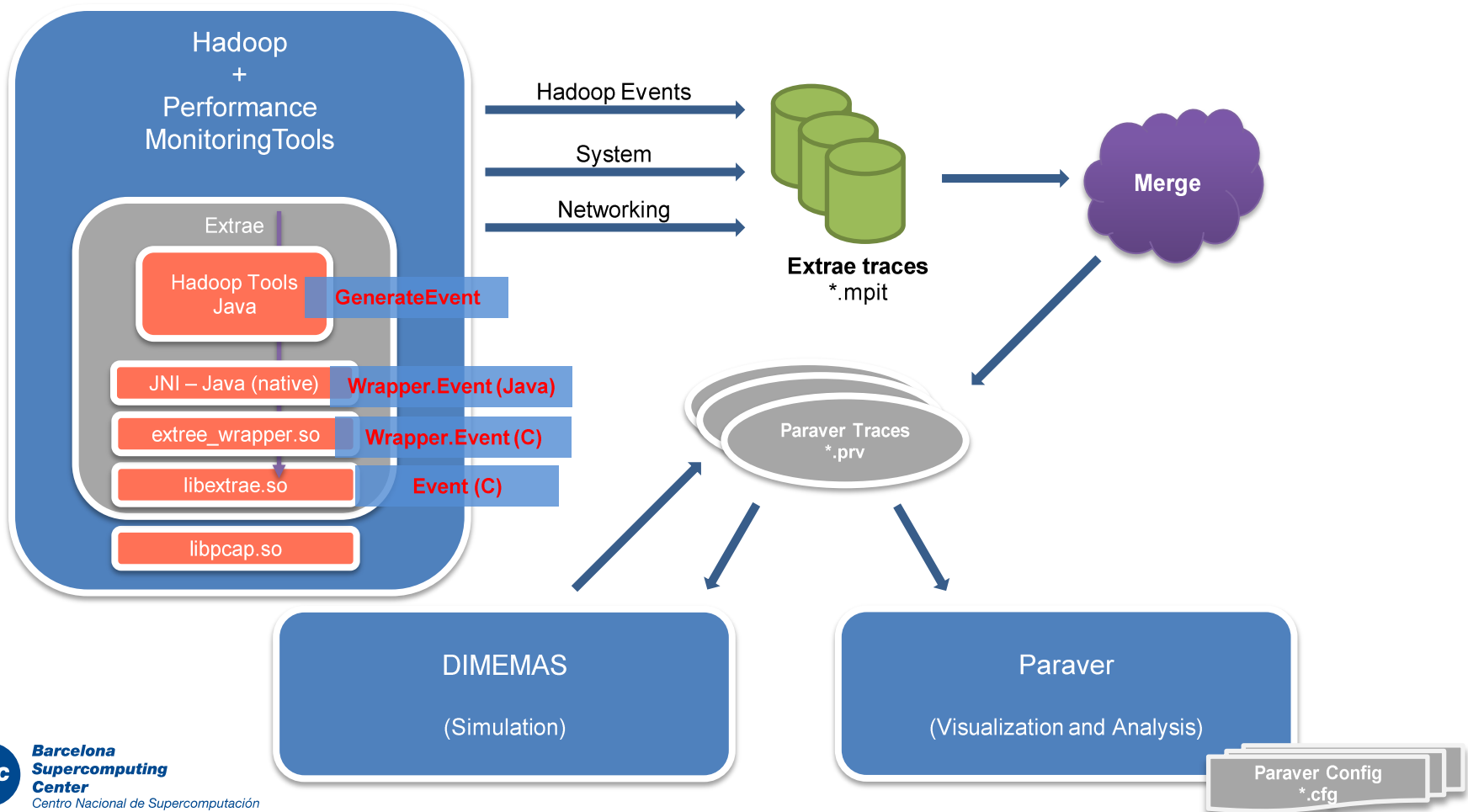
**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

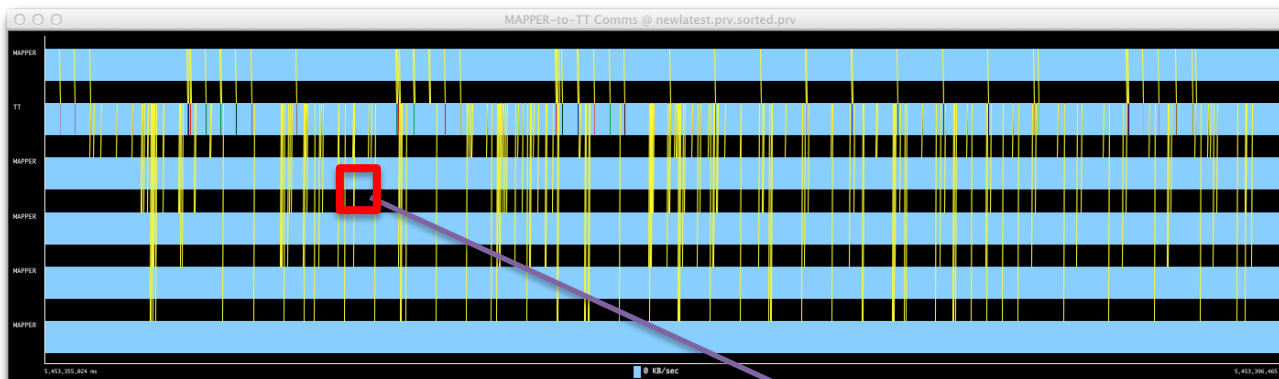
DEEP INSTRUMENTATION

Overview

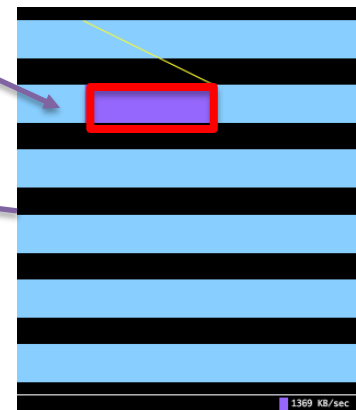
« Hadoop Analysis Toolkit and BSC tools



Example: Packet level communications

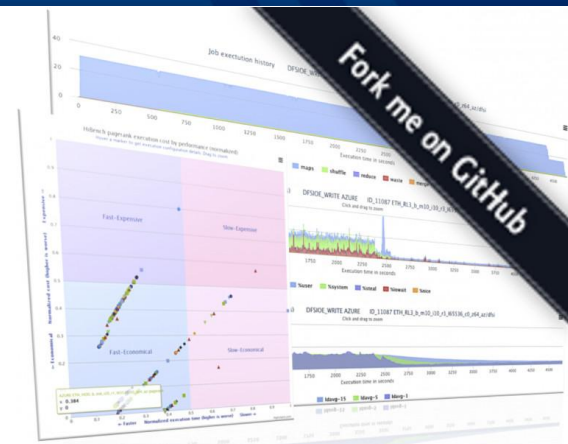


	MAPPER	MAPPER	MAPPER	MAPPER	MAPPER	MAPPER
JT	--	--	--	--	--	--
NN	--	--	--	--	--	--
SNN	--	--	--	--	--	--
MAPPER	--	--	--	--	--	--
REDUCER	--	--	--	--	--	--
DN	--	--	--	--	--	--
MAPPER	--	--	--	--	--	--
TT	792,866.50	453,508.20	450,403.59	528,267.08	457,474.00	484,062.35
MAPPER	--	--	--	--	--	--
MAPPER	--	--	--	--	--	--
MAPPER	--	--	--	--	--	--
MAPPER	--	--	--	--	--	--
REDUCER	--	--	--	--	--	--
Total	792,866.50	453,508.20	450,403.59	528,267.08	457,474.00	484,062.35



Extending ALOJA and Collaborating

1. Install prerequisites
2. git clone <https://github.com/Aloja/aloja.git>
3. cd aloja/vagrant
4. vagrant up
5. Go to: <http://localhost:8080>



Concluding remarks

- « The early findings of the project already show significant value in understanding Hadoop's runtime
 - for optimizing executions times
 - understanding the cost-effectiveness of different configuration and deployment options

- « Our intent is that researchers and organizations evaluating or deploying the Hadoop stack will benefit
 - from this growing database of performance results
 - and configuration guidance

www.bsc.es



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

Thanks

Q&A

Contact: hadoop@bsc.es