



# HiBench: A Representative and Comprehensive Hadoop Benchmark Suite

**Bhaskar D. Gowda**

*Enterprise Architect.*

*DCSG/IASI*



# The HiBench benchmark suite

1

## **Micro Benchmarks**

- Sort
- WordCount
- TeraSort
- Enhanced DFSIO
- TextSort\*

2

## **Web Search**

- Nutch Indexing
- Page Rank

\* HiBench 2.0



HiBench

3

## **Machine Learning**

- Bayesian Classification
- K-Means Clustering





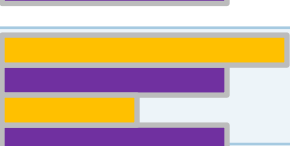


4

## **Analytical Query\***

- (Hive) Join\*
- (Hive) Aggregation\*

HiBench 1.0 published ("The HiBench Suite: Characterization of the MapReduce-Based Data Analysis") and presented in **ICDE'10** workshops

# Workload traits drive optimization approach

Workload	System Resource Utilization	Data Access Patterns	Map/Reduce Stage Time Ratio
Sort	I/O bound	➤ M ➤ R ➤	
WordCount	CPU bound	➤ M ➤ R ➤	
TeraSort	Map stage : CPU-bound; Red stage : I/O-bound	➤ M ➤ R ➤	
Nutch Indexing	I/O bound, high CPU utilization in map stage	➤ M ➤ R ➤	
Page Rank (1 <sup>st</sup> & 2 <sup>nd</sup> job)	CPU-bound in all jobs	➤ M ➤ R ➤ ➤ M ➤ R ➤	
Bayesian Classification (1 <sup>st</sup> & 2 <sup>nd</sup> job)	I/O bound, with high CPU utilization in map stage in the 1 <sup>st</sup> job	➤ M ➤ R ➤ ➤ M ➤ R ➤	
K-means Clustering	CPU bound in iteration; I/O bound in clustering	➤ M ➤ R ➤ ➤ M ➤ R ➤	
Enhanced DFSIO	I/O-bound	trivial	trivial

➤ data    ➤ less data    ➤ even less data    ➤ compressed

no reducer

See our whitepaper "Optimizing Hadoop Deployments" (<http://communities.intel.com/docs/DOC-5645>)

# Tradeoffs and Discussions

	<b>GridMix3</b>	<b>SWIM</b>	<b>HiBench</b>
Synthetic workloads	Yes	Yes	No
Ease of scaling	No	Yes	Partially
Mix of jobs	Yes	Yes	Hard
Platform implications	I/O subsystem	I/O subsystem	Full



Amazing things happen with Intel inside®

