

Data-Aware Support for Hybrid HPC and Big Data Applications

Silvina Caíno-Lores, Florin Isaila, Jesús Carretero

{scaino, fisaila, jcarrete}@inf.uc3m.es

Computer Architecture and Technology Area

Computer Science and Engineering Department, University Carlos III of Madrid

Background and motivation

- Scientific computing is becoming data-intensive (data analysis, visualization...).
- Big Data applications require increasing performance.
- Traditional HPC and BD approaches do not suffice for hybrid applications.



Opportunity to exploit Big Data application models and infrastructures to increase HPC scalability

Objectives

1. Explore the effects of BD paradigms in current scientific applications.
2. Define a common interface for HPC applications and analysis jobs.
3. Build a middleware for performance, with a focus on I/O and locality.
4. Exploit the upcoming advances in supercomputing infrastructures.

Preliminary Requirement Analysis

- **Case study:** Monte Carlo hydrology workflow with legacy kernels
- **Computing models:**
 - Local cluster
 - Private cloud (OpenNebula)
- **Programming models:**
 - MPI
 - Apache Spark

BIG DATA PARADIGM (SPARK, HADOOP)	
Pros	Fault-tolerance by design
	Transparent data-locality Job and task scheduling at platform level
Cons	Low resource management control
	Significant memory overhead
	Poor integration with kernels
	Key-value only
	Deep software and communication stack
HPC PARADIGM (MPI)	
Pros	Low resource consumption
	Efficient communication
	Generalist and tailorable processes
Cons	Limited parallel abstractions
	No native provenance or replication

Research Roadmap

1. Analyze the data flow and I/O patterns of current scientific applications.
2. Define a common interface to support:
 - Traditional/legacy MPI applications.
 - Complementary data analysis and visualization.
 - Filtering stages.
3. Provide data awareness with **localization and multi-level caching**.
4. Harmonize **fault-tolerance** techniques found in HPC and BD.
5. Coordinate data management, fault-tolerance, and execution within the **task and I/O scheduler**.
6. Evaluate with artificial benchmarks and meaningful use cases.

Data-Aware Support for Hybrid HPC and Big Data Applications

Silvina Caíno-Lores, Florin Isaila, Jesús Carretero

{scaino, fisaila, jcarrete}@inf.uc3m.es

Computer Architecture and Technology Area

Computer Science and Engineering Department, University Carlos III of Madrid