# Pilot-Streaming: Design Considerations for a Stream Processing Framework for High-Performance Computing

**Andre Luckow, Peter M. Kasson, Shantenu Jha**
**STREAMING 2016, 03/23/2016**
**RADICAL, Rutgers, http://radical.rutgers.edu**

# Motivation

There is a need to couple data sources, HPC, analytics! 20+ applications identified at STREAM16
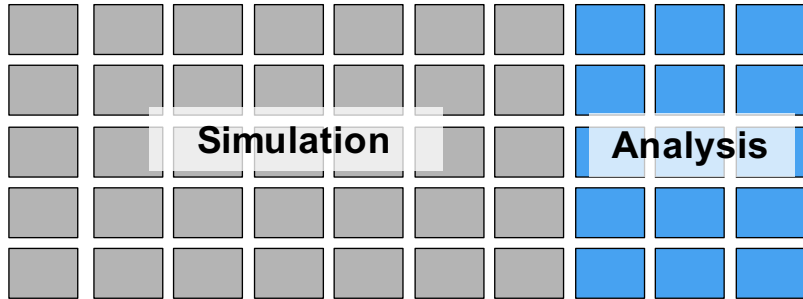
**Challenges:**
- Data applications and pipelines are **complex**
- **Scalability and Elasticity:** dynamic changes in resource demands
- **Scheduling and provisioning of resources:** right amount of resources at right time
- **Programming models:** HPC (MPI, OpenMP, GPU) vs. Big Data (Java, Python, R)
- **Interoperability:** Data sources sinks often in different environments (IoT, cloud, HPC, HPDC) than compute
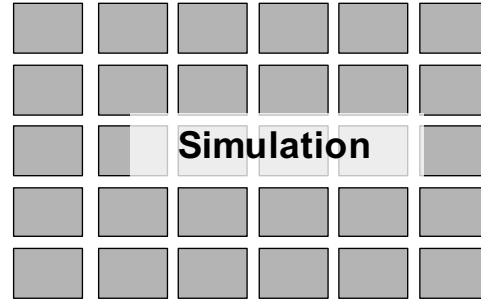
**Current State:**
- Streaming (in sciences) often implemented on application-level (w/ limited re-use)
- Manifold landscape of streaming tools (Apache Open Source Tools, Cloud Tools)
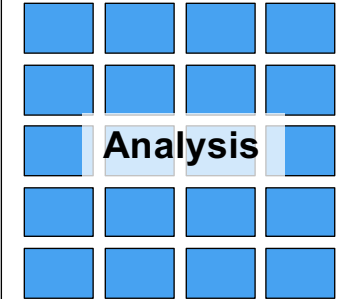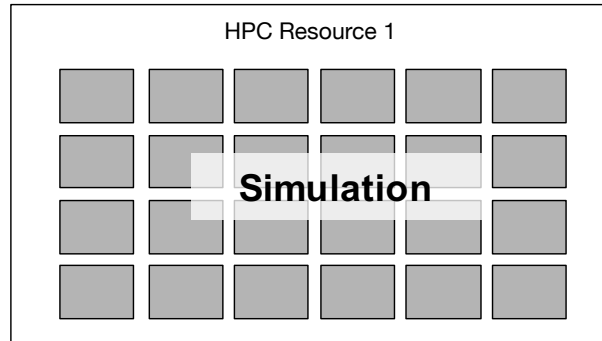
# Workload Characteristics

# Workload Characteristics
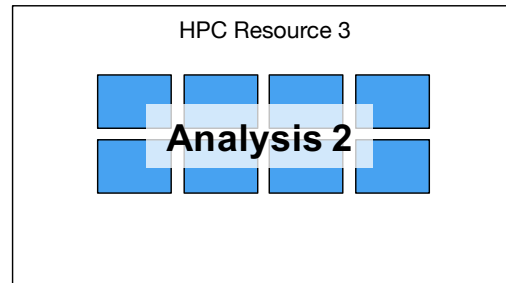
# Introduction Pilot Abstraction

# The Convergence of HPC and "Data Intensive" Computing



| Applications | | | | Applications | | | | |
|---|---|---|---|---|---|---|---|---|
| Orchestration (Pegasus, Taverna, Dryad, Swift) | | | | Orchestration (Oozie, Pig) | | | | |

High-Performance Computing | Apache Hadoop Big Data

A Tale of Two Data-Intensive Paradigms: Data Intensive Applications, Abstractions and Architectures In collaboration with Geoffrey Fox (Indiana), http://arxiv.org/abs/1403.1528

# Pilot-Abstraction for HPC and Hadoop Interoperability

| Map Reduce | Spark-App | Other YARN App | Hadoop/Spark App | HPC App (e.g. MPI) | Application |

| YARN | Spark | Hadoop Application Scheduler (e.g. Spark, Tez, LLama) | Pilot-Job | Application-level Scheduling |
| Pilot-Job | | | |

| HPC Scheduler (Slurm, Torque, SGE) | YARN/HDFS | System-level Scheduling |

**Mode I: Hadoop on HPC**  **Mode II: HPC on Hadoop**

http://arxiv.org/abs/1602.00345

# Streaming and Batch Computing

Data



**Questions:**
- How to manage batch and streaming frameworks side-by-side?
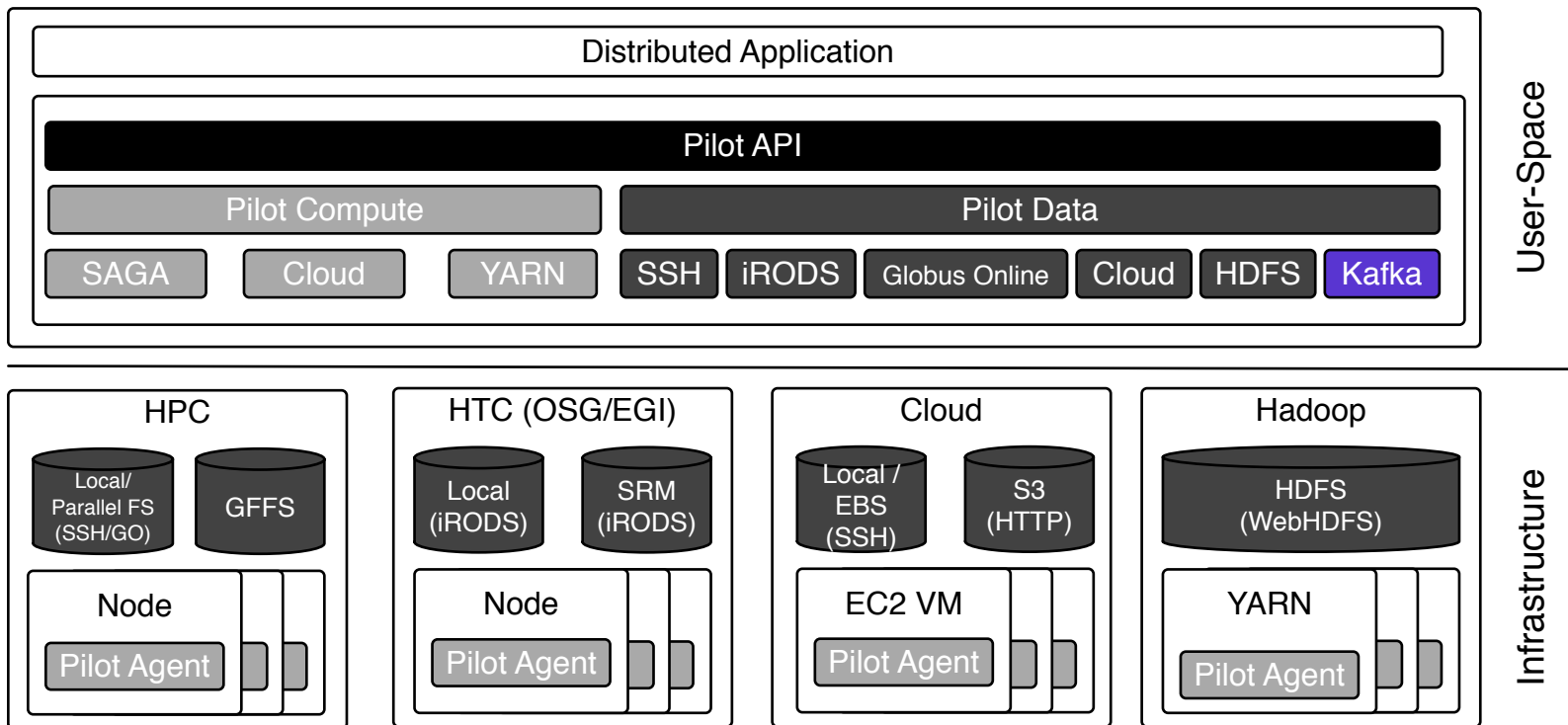- How to enable interoperability between different programming system/models/middleware/schedulers?
- How to enable elasticity?

# Pilot-Streaming

Distributed Application

User-Space

Pilot API

Pilot Compute

Pilot Data

SAGA | Cloud | YARN | SSH | iRODS | Globus Online | Cloud | HDFS | Kafka

Infrastructure

### HPC

Local/ Parallel FS (SSH/GO)

GFFS

**Node**

Pilot Agent

### HTC (OSG/EGI)

Local (iRODS)

SRM (iRODS)

**Node**

Pilot Agent

### Cloud

Local / EBS (SSH)

S3 (HTTP)

**EC2 VM**

Pilot Agent

### Hadoop

HDFS (WebHDFS)

**YARN**

Pilot Agent

# Conclusion

1. Pilot-Jobs enable the co-location of HPC/Simulations and Big Data Tools (Hadoop, Spark, higher-level tools)

2. Pilot-Streaming will support message-broker as data source/sink that enables the de-coupling of applications

3. Dynamic resource management provided by the Pilot-Abstraction is critical for stream environments

# Thank you!