Outline
Hadoop
BlobSeer as storage for Hadoop
Introducing support for append in Hadoop
Conclusions

# BlobSeer in the context of MapReduce applications

Diana Moise

KerData team, INRIA

**Outline**
Hadoop
BlobSeer as storage for Hadoop
Introducing support for append in Hadoop
Conclusions

Outline
**Hadoop**
BlobSeer as storage for Hadoop
Introducing support for append in Hadoop
Conclusions

Hadoop Core

# Hadoop

- $\star$ Yahoo!'s implementation of MapReduce
- $\star$ Open-source Java project
- $\star$ Large scale computation and data processing
- $\star$ Works on comodity hardware

Outline
**Hadoop**
BlobSeer as storage for Hadoop
Introducing support for append in Hadoop
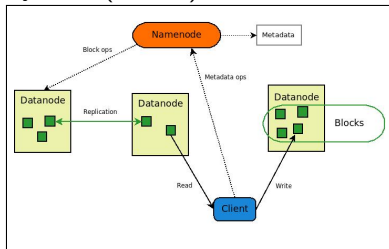Conclusions

Hadoop Core

# Hadoop

- $\star$ Yahoo!'s implementation of MapReduce
- $\star$ Open-source Java project
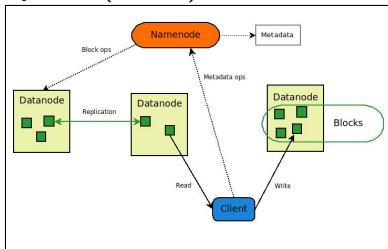- $\star$ Large scale computation and data processing
- $\star$ Works on comodity hardware

Outline
**Hadoop**
BlobSeer as storage for Hadoop
Introducing support for append in Hadoop
Conclusions

Hadoop Core

# Hadoop Core

- Hadoop Distributed File System (HDFS)

Outline
**Hadoop**
BlobSeer as storage for Hadoop
Introducing support for append in Hadoop
Conclusions

Hadoop Core

# Hadoop Core

- Hadoop Distributed File System (HDFS)



- Limitations
  1. one writer at a time
  2. no overwrites
  3. no appends

Outline
**Hadoop**
BlobSeer as storage for Hadoop
Introducing support for append in Hadoop
Conclusions

Hadoop Core

# Hadoop Core

- Hadoop Distributed File System (HDFS)



- Hadoop MR framework



- Limitations
  1. one writer at a time
  2. no overwrites
  3. no appends

Outline
**Hadoop**
BlobSeer as storage for Hadoop
Introducing support for append in Hadoop
Conclusions

Hadoop Core

# In-production use at...



Source: http://wiki.apache.org/hadoop/PoweredBy

## Integrating BlobSeer with Hadoop

- Implementing the HDFS API for BlobSeer
  - ⋆ Exposes basic file system operations: create, read, write...
  - ⋆ Introduces support for concurrent append operations

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

Integrating BlobSeer with Hadoop
Experimental evaluation

# Integrating BlobSeer with Hadoop

- Implementing the HDFS API for BlobSeer
  - ⋆ Exposes basic file system operations: create, read, write...
  - ⋆ Introduces support for concurrent append operations
- BlobSeer File System (BSFS)

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

**Integrating BlobSeer with Hadoop**
Experimental evaluation

## Integrating BlobSeer with Hadoop

- Implementing the HDFS API for BlobSeer
  - ⋆ Exposes basic file system operations: create, read, write...
  - ⋆ Introduces support for concurrent append operations
- BlobSeer File System (BSFS)
  - ✓ File system namespace - keeps metadata, maps files to BLOBs

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

**Integrating BlobSeer with Hadoop**
Experimental evaluation

## Integrating BlobSeer with Hadoop

- Implementing the HDFS API for BlobSeer
  - ⋆ Exposes basic file system operations: create, read, write...
  - ⋆ Introduces support for concurrent append operations
- BlobSeer File System (BSFS)
  - ✓ File system namespace - keeps metadata, maps files to BLOBs
  - ✓ Client-side buffering: data prefetching, write aggregation

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

Integrating BlobSeer with Hadoop
Experimental evaluation

# Integrating BlobSeer with Hadoop

- Implementing the HDFS API for BlobSeer
    - ⋆ Exposes basic file system operations: create, read, write...
    - ⋆ Introduces support for concurrent append operations
- BlobSeer File System (BSFS)
    - ✓ File system namespace - keeps metadata, maps files to BLOBs
    - ✓ Client-side buffering: data prefetching, write aggregation
    - ✓ Exposes data layout to Hadoop, just like HDFS

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

Integrating BlobSeer with Hadoop
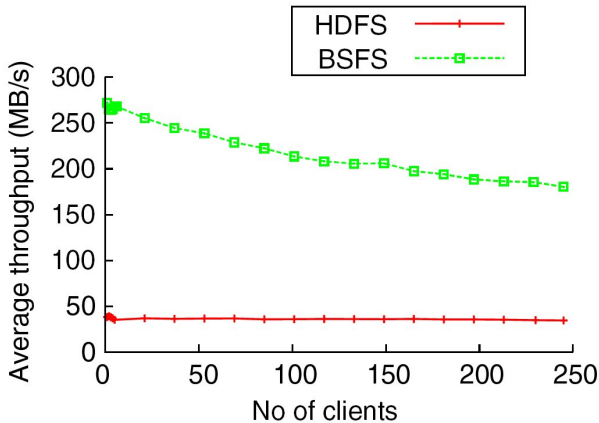**Experimental evaluation**

# Testing and evaluation - overview and goals

- Goal
  - Measure the throughput of HDFS and BSFS
  - Evaluate the impact of replacing HDFS with BSFS
- Test scenarios
  - Microbenchmarks
    - Direct access to the file system
    - Common access patterns in Map/Reduce applications
  - Real Map/Reduce Applications
    - Distributed sort

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

Integrating BlobSeer with Hadoop
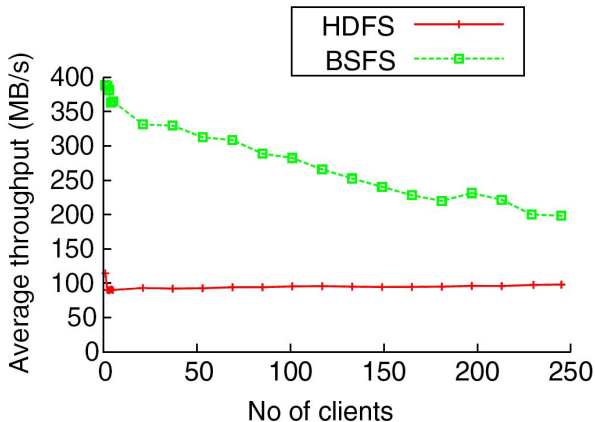**Experimental evaluation**

## Setup

- 270 nodes from the same cluster on Grid'5000
- HDFS:
  - one namenode on a dedicated machine
  - one datanode on each cluster node
- BSFS:
  - one vmanager, one pmanager, one namespace manager
  - 20 metadata providers
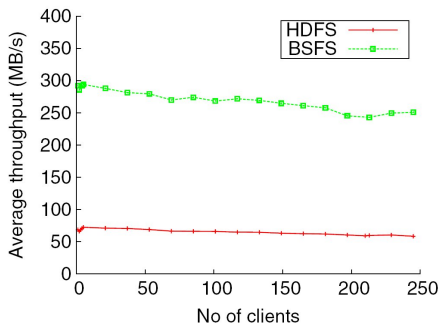  - providers on the rest of the nodes

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

Integrating BlobSeer with Hadoop
**Experimental evaluation**

# Concurrent clients writing to different files

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

Integrating BlobSeer with Hadoop
**Experimental evaluation**

# Concurrent clients reading from different files

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

Integrating BlobSeer with Hadoop
**Experimental evaluation**

# Concurrent clients reading parts from the same file

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

Integrating BlobSeer with Hadoop
**Experimental evaluation**

# Concurrent clients appending data to the same file

Outline
Hadoop
**BlobSeer as storage for Hadoop**
Introducing support for append in Hadoop
Conclusions

Integrating BlobSeer with Hadoop
**Experimental evaluation**

# Distributed sort

- Sorts key-value pairs
- Both read and write intensive

Outline
Hadoop
BlobSeer as storage for Hadoop
**Introducing support for append in Hadoop**
Conclusions

Application case

## Modifying Hadoop to use appends

- Append implemented at
  the file system level

Outline
Hadoop
BlobSeer as storage for Hadoop
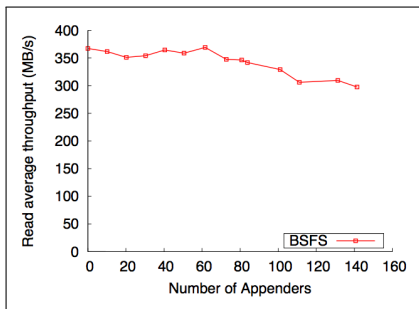**Introducing support for append in Hadoop**
Conclusions

Application case

# Modifying Hadoop to use appends

- Append implemented at the file system level
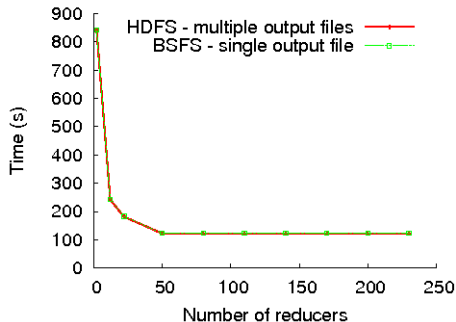- Modify reducer code in Hadoop to append the output to a single file

Outline
Hadoop
BlobSeer as storage for Hadoop
**Introducing support for append in Hadoop**
Conclusions

Application case

# Concurrent reads and appends to the same file

Outline
Hadoop
BlobSeer as storage for Hadoop
**Introducing support for append in Hadoop**
Conclusions

Application case

# Data join - Results

- Similar to outer join from the database context
- Merge two input files based on common keys
- 6.3 GB of output

Outline
Hadoop
BlobSeer as storage for Hadoop
Introducing support for append in Hadoop
**Conclusions**

## Conclusions

- BSFS improves Hadoop's throughput
- Support for append
- Work in progress
  - ⋆ Intermediate data management

    Store map output to BSFS

    Resume computation in case of failures
  - ⋆ Pipeline MapReduce applications

    Schedule mappers as soon as splits are produced

    Application study: Pig

# Thank you!