



## Introduction

- Big Data
- Existing Solutions
- Outline

## Motivation

- Distributed and Out-of-core Computing
- Insights

## BDMPI

- Overview
- Usage
- Implementation

## Results

- Experiments
- Single Node
- Cluster
- Scaling

## Conclusion

# BDMPI: Conquering Big Data with Small Clusters using MPI

Dominique LaSalle and George Karypis  
University of Minnesota, Minneapolis, MN, USA

November 18, 2013



# Big Data

## Introduction

### Big Data

Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## What is Big Data?

- Depends on your compute system:
  - Laptop/PC
  - Server
  - Cluster
  - Data Center
- Data > DRAM



# Existing Solutions

## Introduction

Big Data  
**Existing  
Solutions**  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## Big Data Solutions

- MapReduce/Hadoop
- GraphChi
- Giraph
- Hama
- Custom Solution



# Outline

## Introduction

Big Data  
Existing  
Solutions  
**Outline**

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

- 1 **Introduction**
- 2 Motivation
- 3 BDMPI
  - 1 Overview
  - 2 Usage
  - 3 Implementation
- 4 Results
- 5 Conclusion



# Distributed and Out-of-core Computing

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## Distributed Algorithms

- Minimize communication between processes.
- Extract independent tasks to perform in parallel.
- Organized into a series of compute and collective/point-to-point communication steps.

## Out-of-Core Algorithms

- Minimize reads and writes to disk.
- Extract independent tasks to perform serially.
- Organized into a series of compute and disk read/write steps.



## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing

### Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## The Graph Ordering Problem

- How can a graph be efficiently re-order in an out-of-core fashion?
- How can a graph be efficiently re-order in a distributed fashion?

## General Applications

- How can we treat a remote process as a disk?
  - Already supported by MPI's one sided communication (exchange `fread/fwrite` for `MPI_get/MPI_put`).
- Can we treat the disk as a remote process?
  - Need to handle remote computations/data movement.



# How it Works

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## BDMPI

- Transparent layer between an MPI program and an MPI runtime.
- For a problem of size  $n$  and a compute cluster with  $p$  processing nodes each with  $m$  memory:
  - 1 Divide the data into  $t = n/m$  blocks.
  - 2 Spawn a *master* process on each compute node.
  - 3 Spawn  $t/p$  *slave* processes on each compute node.
- Allow only one slave process to run at a time on each compute node.
  - That process will run until it blocks on a communication operation.



# Why it Works

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## Node-Level Cooperative Multi-Tasking

- Processes run until blocking for a collective communication or receive operation.
- Cost of loading data from disk is amortized over large blocks of computation.
- Since only one process runs at a time, the thrashing associated with multiple processes attempting to gain residency is avoided.





# BDMPI Usage

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview

**Usage**  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## Usage

- `bdmpiexec`

```
mpiexec -np 80  
    program [arg1] [arg2] ...
```

```
bdmpiexec -np 4 [-nr 2] -ns 20  
    program [arg1] [arg2] ...
```

- Executes `mpi` program on a cluster with four nodes as if it were on a cluster of 80 compute nodes.
- `libbdmpi`
  - Provides `MPI_X` functions.
- Replace `#include <mpi.h>` with `#include <bdmpi.h>`.



# BDMPI API

## Introduction

- Big Data
- Existing Solutions
- Outline

## Motivation

- Distributed and Out-of-core Computing
- Insights

## BDMPI

- Overview
- Usage**
- Implementation

## Results

- Experiments
- Single Node
- Cluster
- Scaling

## Conclusion

## MPI Subset Implemented by BDMPI

BDMPI\_Init , BDMPI\_Finalize

BDMPI\_Comm\_size , BDMPI\_Comm\_rank , BDMPI\_Comm\_dup ,  
BDMPI\_Comm\_free , BDMPI\_Comm\_split

BDMPI\_Send , BDMPI\_Isend , BDMPI\_Recv , BDMPI\_Irecv ,  
BDMPI\_Sendrecv

BDMPI\_Probe , BDMPI\_Iprobe , BDMPI\_Test , BDMPI\_Wait ,  
BDMPI\_Get\_count

BDMPI\_Barrier

BDMPI\_Bcast , BDMPI\_Reduce , BDMPI\_Allreduce ,  
BDMPI\_Scan , BDMPI\_Gather[v] , BDMPI\_Scatter[v] ,  
BDMPI\_Allgather[v] , BDMPI\_Alltoall[v]



# Implementation

## Introduction

- Big Data
- Existing Solutions
- Outline

## Motivation

- Distributed and Out-of-core Computing
- Insights

## BDMPI

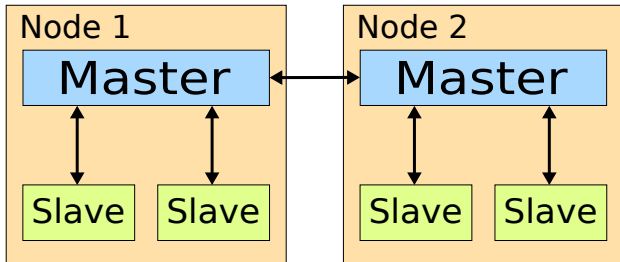
- Overview
- Usage
- Implementation**

## Results

- Experiments
- Single Node
- Cluster
- Scaling

## Conclusion

## Communication Model





# Implementation Cont.

## Introduction

- Big Data
- Existing Solutions
- Outline

## Motivation

- Distributed and Out-of-core Computing
- Insights

## BDMPI

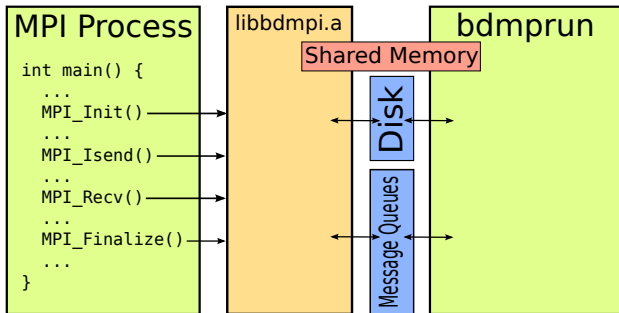
- Overview
- Usage
- Implementation**

## Results

- Experiments
- Single Node
- Cluster
- Scaling

## Conclusion

## Master-Slave Communication





# Point-to-point Communication

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage

## Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## Message Buffering

- Small messages buffered in memory.
- Large messages buffered on disk.

## Send and ISend

- Message buffering allows sending process to continue executing without blocking.

## Recv and IRecv

- If the master has already buffered the message, no blocking occurs.
- Otherwise the process becomes blocked, and another process is allowed to run.



# Benchmarks

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## PageRank

- Memory heavy operation.
- Multiplying a sparse matrix by a vector.

## KMeans Clustering

- Multiplying a sparse matrix by a dense matrix (100 clusters).

## SGD

- Matrix factorization  $A = UV$  (20 factors).
- Element-wise random traversal.
- SGD-row
  - Row-wise traversal.
  - Better locality than regular SGD.



# Test Codes

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

**Experiments**  
Single Node  
Cluster  
Scaling

## Conclusion

- **Serial-OOC** - Custom out-of-core solutions.
- **MPI** - MPI codes ran using MPICH.
- **GraphChi** - Kyrola et. al. 2012.
- **Hadoop**
  - Mahout for KMeans.
  - Pegasus for PageRank - Kang et. al. 2009.
- **BDMPI**
  - **BDMPI** - MPI codes ran using the BDMPI runtime.
  - **BDMPI-mlock** - MPI codes + `munlock()/mlock()`.
  - **BDMPI-OOC** - MPI codes + `fread()/fwrite()`.



# Experiment Setup

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

**Experiments**  
Single Node  
Cluster  
Scaling

## Conclusion

## Our Cluster

- Four machine cluster:
  - Intel i7 @ 3.4 GHz
  - 4 GB of DRAM
  - Seagate Barracuda 7200 RPM 1.0 TB (300GB swap and /scratch partitions)

## Our Datasets

- PageRank - 6.6B edges, ordered randomly (50GB CSR).
- KMeans - 30M  $\times$  83K with 7.3B non-zeros (56GB CSR).
- SGD - 3.8M  $\times$  284K with 12.8B non-zeros (50GB CSR).





# Single Node Results

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

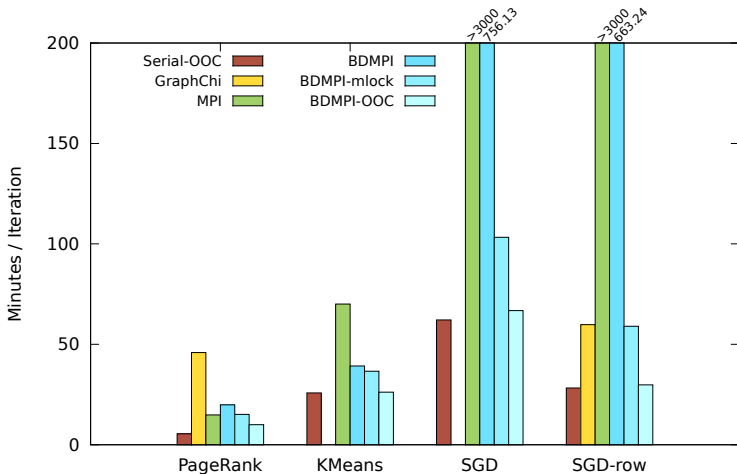
## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
**Single Node**  
Cluster  
Scaling

## Conclusion





# Cluster Results

## Introduction

- Big Data
- Existing Solutions
- Outline

## Motivation

- Distributed and Out-of-core Computing
- Insights

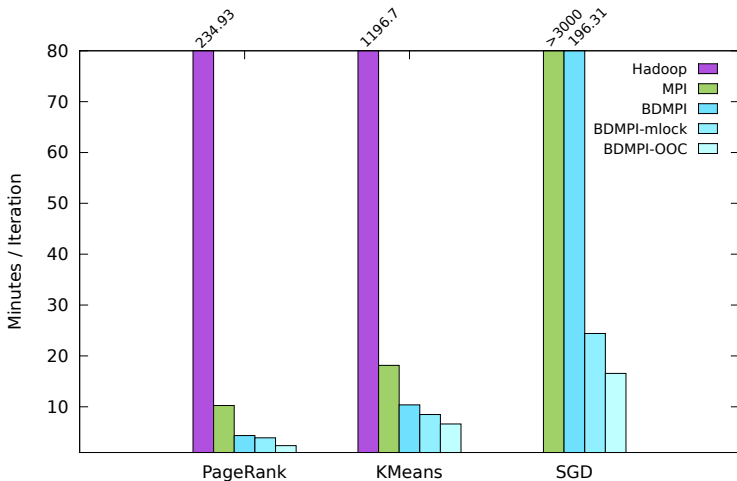
## BDMPI

- Overview
- Usage
- Implementation

## Results

- Experiments
- Single Node
- Cluster**
- Scaling

## Conclusion





# Scaling Results

## Introduction

- Big Data
- Existing Solutions
- Outline

## Motivation

- Distributed and Out-of-core Computing
- Insights

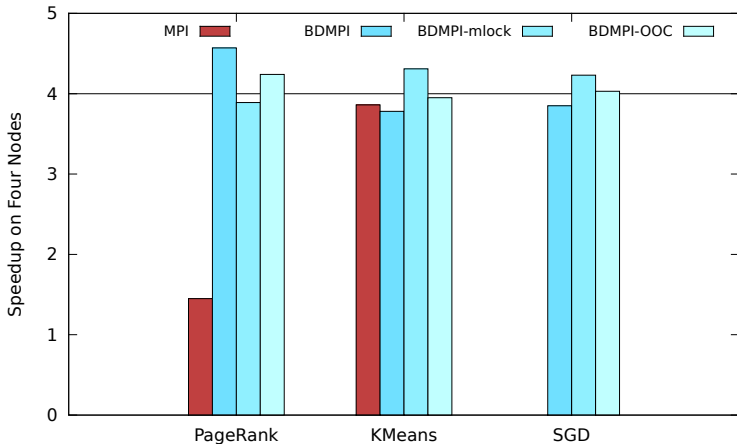
## BDMPI

- Overview
- Usage
- Implementation

## Results

- Experiments
- Single Node
- Cluster
- Scaling

## Conclusion





# Conclusion

## Introduction

Big Data  
Existing  
Solutions  
Outline

## Motivation

Distributed and  
Out-of-core  
Computing  
Insights

## BDMPI

Overview  
Usage  
Implementation

## Results

Experiments  
Single Node  
Cluster  
Scaling

## Conclusion

## BDMPI

- Utilizes existing MPI interface.
  - Turns existing MPI applications into distributed out-of-core applications.
  - Leverages 20 years worth of experience.
- Achieves speeds comparable to custom out-of-core solutions.
- Scales well across multiple machines.