## My historical perspective

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

## My historical perspective

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

## My historical perspective

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

## My historical perspective

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

## My historical perspective

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

## My historical perspective

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

## My historical perspective

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

## My historical perspective

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

## My historical perspective

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

- I made it to the 8 CCGSC workshops!
- I talked about a nice little scheduling problem in 1992
- I talked about a nice little scheduling problem in 1994
- I talked about a nice little scheduling problem in 1996
- I talked about a nice little scheduling problem in 1998
- I talked about a nice little scheduling problem in 2000
- I talked about a nice little scheduling problem in 2002
- I talked about a nice little scheduling problem in 2004
- I wondered what I should do this year?
- Maybe I should find a nice little scheduling problem! ☺

- Parallel algorithm design and scheduling were already difficult tasks with homogeneous machines

- **On heterogeneous platforms, it gets worse**

- Patrick Geoffray went from kindergarten to Myricom but he's still a kid! ☺

- He says that only embarrassingly parallel applications can be deployed on the grid

- Clearly, he is over optimistic! ☹

## Did anything change in the meantime?

- Parallel algorithm design and scheduling were already difficult tasks with homogeneous machines

- **On heterogeneous platforms, it gets worse**

- Patrick Geoffray went from kindergarten to Myricom but he's still a kid! 🙂

- He says that only embarrassingly parallel applications can be deployed on the grid

- Clearly, he is over optimistic! ☹

## Did anything change in the meantime?

- Parallel algorithm design and scheduling were already difficult tasks with homogeneous machines
- **On heterogeneous platforms, it gets worse**
- Patrick Geoffray went from kindergarten to Myricom but he's still a kid! ☺
- He says that only embarrassingly parallel applications can be deployed on the grid
- Clearly, he is over optimistic! ☹

## Did anything change in the meantime?

- Parallel algorithm design and scheduling were already difficult tasks with homogeneous machines
- **On heterogeneous platforms, it gets worse**
- Patrick Geoffray went from kindergarten to Myricom but he's still a kid! 😊
- He says that only embarrassingly parallel applications can be deployed on the grid
- Clearly, he is over optimistic! ☹

## Did anything change in the meantime?

- Parallel algorithm design and scheduling were already difficult tasks with homogeneous machines
- **On heterogeneous platforms, it gets worse**
- Patrick Geoffray went from kindergarten to Myricom but he's still a kid! ☺
- He says that only embarrassingly parallel applications can be deployed on the grid
- Clearly, he is over optimistic! ☹

## A nice little embarrassingly parallel application



- One (divisible load) application running on each cluster
  ⇒ Which fraction of the job to delegate to other clusters?
- Different communication-to-computation ratios
  ⇒ How to ensure fair scheduling and good resource utilization?

# Revisiting matrix product on heterogeneous platforms

Jack Dongarra, Zhiao Shi, UT Knwoville

Jean-François Pineau, Yves Robert, Frédéric Vivien, ENS Lyon

The great talk you've been expecting     4/35

# Revisiting matrix product on heterogeneous platforms

**Eh wait!**

Experiments are not ready?!

Jean-François Pineau, Yves Robert, Frédéric Vivien, ENS Lyon

# Scheduling and Data Redistribution Strategies on Star Platforms

Loris Marchal, Veronika Rehn,
Yves Robert and Frédéric Vivien

GRAAL team, LIP
École Normale Supérieure de Lyon

September 2006

## Outline

1. Target problem
   - Fully homogeneous platforms
   - Bus platforms
   - General platforms

2. Simulations

3. Divisible Loads Using the Multiport Switch-Model

4. Conclusion

## Outline

# Example

## Architecture



$P_1$

$P_2$

$P_3$

$P_4$

$T = 0$

# Example

## Architecture

# Example

# Example

## Architecture

# Example

## Architecture

# Example

## Architecture

## Architecture

# Framework

- Master-slave platforms
- New: Distributed loads

## Problem

Redistribution of data

Goal: Minimize overall processing time

## Data models

- Independent tasks
- Divisible loads

# Framework

- Master-slave platforms
- New: Distributed loads

## Problem

Redistribution of data
Goal: Minimize overall processing time

## Data models

- Independent tasks
- Divisible loads

# Framework

- Master-slave platforms
- New: Distributed loads

### Problem

Redistribution of data
Goal: Minimize overall processing time

### Data models

- Independent tasks
- Divisible loads

## Related Work

### Known results



### Our approach



### Independent tasks

- Application: BOINC (e.g. Einstein@home)
- NP-completeness for different task sizes

### Divisible load theory

- Perfect parallel jobs
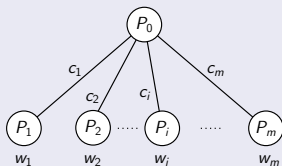- Optimal algorithms for video processing (Altilar, Paker)

### Redistribution algorithms

- NP-completeness (Kremer)
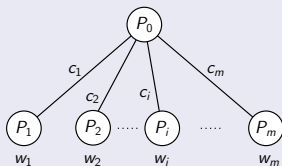- Optimality for particular cases: homogeneous ring topologies

# Related Work

### Known results



### Our approach



### Independent tasks

- Application: BOINC (e.g. Einstein@home)
- NP-completeness for different task sizes

### Divisible load theory

- Perfect parallel jobs
- Optimal algorithms for video processing (Altilar, Paker)

### Redistribution algorithms

- NP-completeness (Kremer)
- Optimality for particular cases: homogeneous ring topologies

## Related Work

### Known results



### Our approach



#### Independent tasks

- Application: BOINC (e.g. Einstein@home)
- NP-completeness for different task sizes

#### Divisible load theory

- Perfect parallel jobs
- Optimal algorithms for video processing (Altilar, Paker)

#### Redistribution algorithms

- NP-completeness (Kremer)
- Optimality for particular cases: homogeneous ring topologies

## Model 11/35

### Star network



- Star network $S = P_0, P_1, \ldots, P_m$
- Communication cost $c_i$
- Computing power $w_i$
- Initial data $L_i$
- Independent and identical tasks
- Linear cost model
- Bidirectional one-port model
- Objective function:
  Minimize makespan

# Model

### Star network



- Star network $S = P_0, P_1, \ldots, P_m$
- Communication cost $c_i$
- Computing power $w_i$
- Initial data $L_i$
- Independent and identical tasks
- Linear cost model
- Bidirectional one-port model
- Objective function:
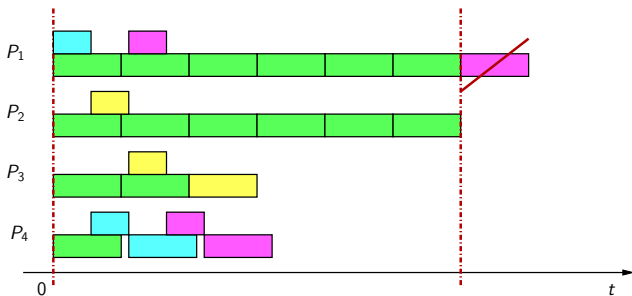  Minimize makespan
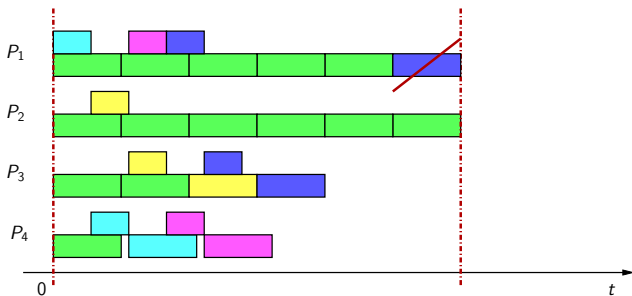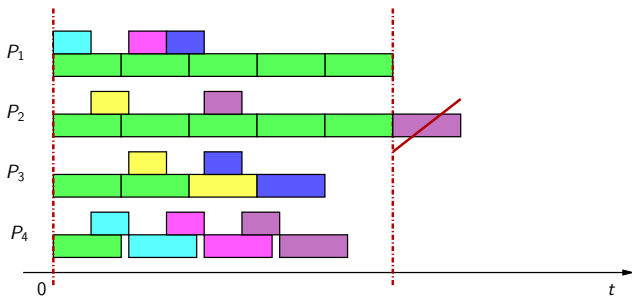
# Best-Balance Algorithm - BBA

- Homogeneous communication links
- Homogeneous workers



Principle: Local optimization of current makespan

# Best-Balance Algorithm - BBA

- Homogeneous communication links
- Homogeneous workers



**Principle:** Local optimization of current makespan

# Best-Balance Algorithm - BBA

- Homogeneous communication links
- Homogeneous workers



Principle: Local optimization of current makespan

# Best-Balance Algorithm - BBA

- Homogeneous communication links
- Homogeneous workers



Principle: Local optimization of current makespan

# Best-Balance Algorithm - BBA

- Homogeneous communication links
- Homogeneous workers



**Principle:** Local optimization of current makespan

# Best-Balance Algorithm - BBA

- Homogeneous communication links
- Homogeneous workers



Principle: Local optimization of current makespan

# Best-Balance Algorithm - BBA

- Homogeneous communication links
- Homogeneous workers



Principle: Local optimization of current makespan

# BBA- Optimality

### Theorem

BEST-BALANCE ALGORITHM *calculates an optimal schedule S on a fully homogeneous star network.*

# Moore-Based Binary-Search Algorithm - MBBSA

- Homogeneous communication links
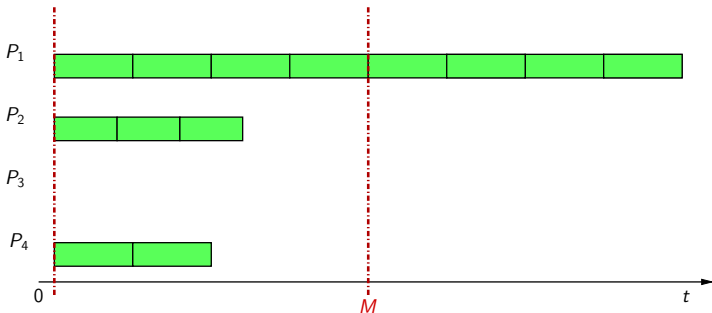- Heterogeneous workers
- Makespan $M$

Principle:

- Moore's algorithm
- Schedule within $M$
- Binary search

### Moore's algorithm

1: Order the jobs by non-decreasing deadlines: $d_1 \leq d_2 \leq \cdots \leq d_d$
2: $\sigma \leftarrow \emptyset;\ t \leftarrow 0$
3: **for** $i := 1$ to $n$ **do**
4: $\quad \sigma \leftarrow \sigma \cup \{i\}$
5: $\quad t \leftarrow t + w_i$
6: $\quad$ **if** $t > d_i$ **then**
7: $\quad\quad$ Find job $j$ in $\sigma$ with largest $w_j$ value
8: $\quad\quad \sigma \leftarrow \sigma \backslash \{j\}$
9: $\quad\quad t \leftarrow t - w_j$
10: $\quad$ **end if**
11: **end for**

# Moore-Based Binary-Search Algorithm - MBBSA

- Homogeneous communication links
- Heterogeneous workers
- Makespan $M$

Principle:

- Moore's algorithm
- Schedule within $M$
- Binary search

### Moore's algorithm

1: Order the jobs by non-decreasing deadlines:
$d_1 \leq d_2 \leq \cdots \leq d_d$
2: $\sigma \leftarrow \emptyset; \ t \leftarrow 0$
3: **for** $i := 1$ to $n$ **do**
4: $\quad \sigma \leftarrow \sigma \cup \{i\}$
5: $\quad t \leftarrow t + w_i$
6: $\quad$ **if** $t > d_i$ **then**
7: $\qquad$ Find job $j$ in $\sigma$ with largest $w_j$ value
8: $\qquad \sigma \leftarrow \sigma \backslash \{j\}$
9: $\qquad t \leftarrow t - w_j$
10: $\quad$ **end if**
11: **end for**

Determination of senders and receivers

Determination of senders and receivers

Computation of deadlines

Bus platforms

# MBBSA- Phase 4

Scheduling step

Scheduling step

Scheduling step

Scheduling step

Scheduling step

Target problem · · · · · · · · · · · · · · Simulations · · · · · · · · · · · · · · Divisible Loads · · · · · · · · · · · · · · Conclusion

Bus platforms

# MBBSA- Phase 4 · · · · · · · · · · · · · · · · · · · · · · · · · · 17/35

Scheduling step

Scheduling step

# MBBSA- Optimality

### Theorem

(i) *MBBSA succeeds to build a schedule $\sigma$ for a given makespan M, if and only if there exists one.*

(ii) *Binary search algorithm returns in polynomial time an optimal schedule $\sigma$ for bus platforms (homogeneous communication links and heterogeneous workers).*

# Dealing with fully heterogeneous platforms

Difficulty: Who is sender, who is receiver?

M = 12

| Worker | c | w | load |
|--------|---|---|------|
| $P_1$ | 1 | 1 | 13 |
| $P_2$ | 8 | 1 | 13 |
| $P_3$ | 1 | 9 | 0 |
| $P_4$ | 1 | 10 | 0 |

**Target problem**
Simulations
Divisible Loads
Conclusion

General platforms

NP-completeness
20/35

Scheduling Problem for Master-Slave Tasks on a Star of
Heterogeneous Processors

### Definition (SPMSTSHP)

Let $N$ be a star-network. Let $T$ be a deadline.
"Is it possible to redistribute tasks and process them in time $T$?".
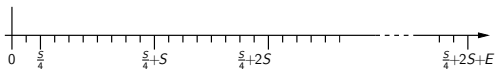
### Theorem

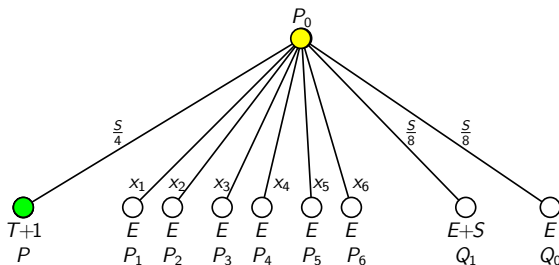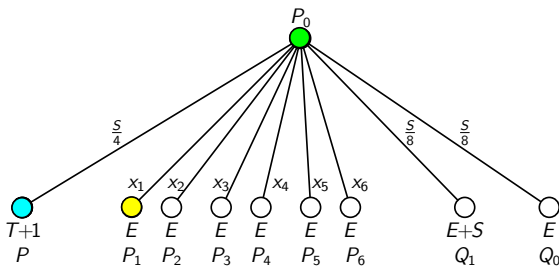*NP-complete in the strong sense.*

## Proof: Reduction to 3-partition

Proof: Reduction to 3-partition problem

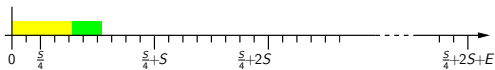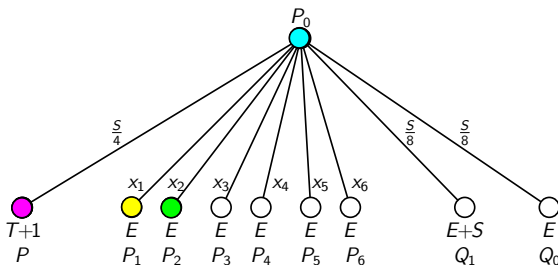# Proof: Reduction to 3-partition

Proof: Reduction to 3-partition problem

Proof: Reduction to 3-partition problem

Proof: Reduction to 3-partition problem

Proof: Reduction to 3-partition problem

# Proof: Reduction to 3-partition

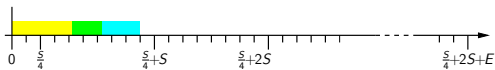Proof: Reduction to 3-partition problem

Proof: Reduction to 3-partition problem

# Proof: Reduction to 3-partition

Proof: Reduction to 3-partition problem

**Target problem**
○○○○○○○○○○●○

Simulations

Divisible Loads

Conclusion

General platforms

## Proof: Reduction to 3-partition

21/35

Proof: Reduction to 3-partition problem

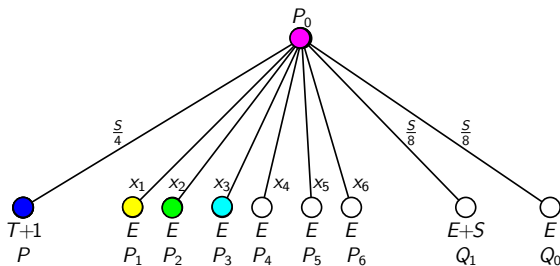Proof: Reduction to 3-partition problem

# Proof: Reduction to 3-partition

Proof: Reduction to 3-partition problem

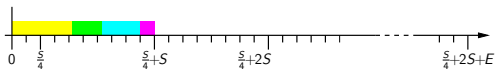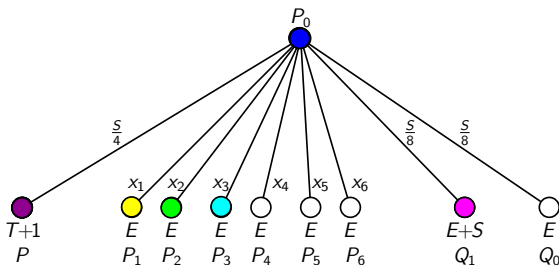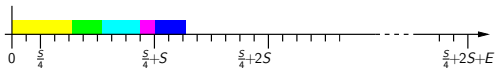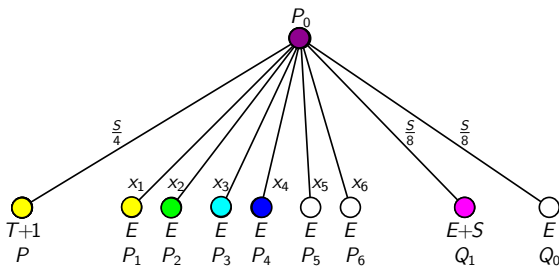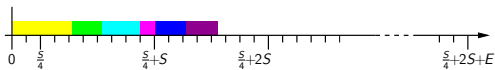# Proof: Reduction to 3-partition

Proof: Reduction to 3-partition problem

# Impact of Heterogeneity

| Platform type | | Difficulty |
| Comm. | Comp. | |
| --- | --- | --- |
| Hom. | Hom. | simple greedy algorithm |
| Hom. | Het. | complicated algorithm |
| Het. | Hom. | ? |
| Het. | Het. | NP-strong |

## Heuristics

- BBA

- MBBSA

- R-BSA: Reversed-Binary Search Algorithm
  Combination of greedy algorithm and binary search

# Impact of Heterogeneity

| Platform type | | Difficulty |
| Comm. | Comp. | |
|---|---|---|
| Hom. | Hom. | simple greedy algorithm |
| Hom. | Het. | complicated algorithm |
| Het. | Hom. | ? |
| Het. | Het. | NP-strong |

## Heuristics

- BBA

- MBBSA

- R-BSA: Reversed-Binary Search Algorithm
  Combination of greedy algorithm and binary search

# Impact of Heterogeneity

| Platform type | | |
|---|---|---|
| Comm. | Comp. | Difficulty |
| Hom. | Hom. | simple greedy algorithm |
| Hom. | Het. | complicated algorithm |
| Het. | Hom. | ? |
| Het. | Het. | NP-strong |

## Heuristics

- BBA

- MBBSA

- R-BSA: Reversed-Binary Search Algorithm
  Combination of greedy algorithm and binary search

## Impact of Heterogeneity

| Platform type | | Difficulty |
|---|---|---|
| Comm. | Comp. | |
| Hom. | Hom. | simple greedy algorithm |
| Hom. | Het. | complicated algorithm |
| Het. | Hom. | ? |
| Het. | Het. | NP-strong |

### Heuristics

- BBA

- MBBSA

- R-BSA: Reversed-Binary Search Algorithm
  Combination of greedy algorithm and binary search

# Impact of Heterogeneity

| Platform type | | Difficulty |
| Comm. | Comp. | |
| --- | --- | --- |
| Hom. | Hom. | simple greedy algorithm |
| Hom. | Het. | complicated algorithm |
| Het. | Hom. | ? |
| Het. | Het. | NP-strong |

### Heuristics

- BBA

- MBBSA

- R-BSA: Reversed-Binary Search Algorithm
  Combination of greedy algorithm and binary search

## Outline

## Simulations

### SimGrid:

Simulator for distributed applications

- 4 platform types
- 1000 instances
- 10 workers
- Random variables

- $c_i$: 1..100
- $w_i$: 1..100
- $L_i$: 0..50

## Simulations

### SimGrid:

Simulator for distributed applications

- 4 platform types
- 1000 instances
- 10 workers
- Random variables

- $c_i$: 1..100
- $w_i$: 1..100
- $L_i$: 0..50

# Trace Tests

## BBA



## MBBSA

# Distance from the Best Heuristic

Heterogeneous platform

# Distance from the Best Heuristic

Heterogeneous platform

## Heterogeneous platform

## Distance from the Best Heuristic

Heterogeneous platform

# Distance from the Best Heuristic

Heterogeneous platform

# Standard Deviation

| Platform type | | Standard deviation | | |
|---|---|---|---|---|
| Comm. | Comp. | BBA | MBBSA | R-BSA |
| Hom | Hom | 0 | 0 | 0.0107 |
| Hom | Het | 0.0006 | 0 | 0.0181 |
| Het | Hom | 0.4007 | 0.0208 | 0.0173 |
| Het | Het | 0.3516 | 0.0327 | 0.0284 |

# Standard Deviation

| Platform type | | Standard deviation | | |
|---|---|---|---|---|
| Comm. | Comp. | BBA | MBBSA | R-BSA |
| Hom | Hom | 0 | 0 | 0.0107 |
| Hom | Het | 0.0006 | 0 | 0.0181 |
| Het | Hom | 0.4007 | 0.0208 | 0.0173 |
| Het | Het | 0.3516 | 0.0327 | 0.0284 |

# Standard Deviation

| Platform type | | Standard deviation | | |
|---|---|---|---|---|
| Comm. | Comp. | BBA | MBBSA | R-BSA |
| Hom | Hom | 0 | 0 | 0.0107 |
| Hom | Het | 0.0006 | 0 | 0.0181 |
| Het | Hom | 0.4007 | 0.0208 | 0.0173 |
| Het | Het | 0.3516 | 0.0327 | 0.0284 |

# Standard Deviation

| Platform type | | Standard deviation | | |
| --- | --- | --- | --- | --- |
| Comm. | Comp. | BBA | MBBSA | R-BSA |
| Hom | Hom | 0 | 0 | 0.0107 |
| Hom | Het | 0.0006 | 0 | 0.0181 |
| Het | Hom | 0.4007 | 0.0208 | 0.0173 |
| Het | Het | 0.3516 | 0.0327 | 0.0284 |

## Outline

# Framework

### Star network

- Switch as master
- $m$ workers
- Computation speed $s_i$
- Bandwidth $b_i$
- Divisible load $\alpha_i$
- Linear cost model
- Overlapped unbounded switch model

# Redistribution Strategy

Goal: Every worker finishes at the same time

# Redistribution Strategy

Goal: Every worker finishes at the same time

## Solution for Divisible Loads          31/35

Imbalance of a worker $\delta_i$

### Linear program

MINIMIZE $T$,

UNDER THE CONSTRAINTS

$$\begin{cases} (1a) & |\delta_i| \leq T \times b_i \\ (1b) & \delta_i \geq \alpha_i - T \times s_i \\ (1c) & \sum_i \delta_i = 0 \end{cases}$$

$$(1)$$

Fraction of load $f_{i,j}$

$$f_{i,j} = \delta_i \times \frac{\delta_j}{\sum_{k \in R} \delta_k} = \delta_i \times \frac{\delta_j}{-L}$$

Communication rate $\lambda_{i,j}$

$$\lambda_{i,j} = \frac{f_{i,j}}{T_0}$$

Computation rate $\gamma_{i,j}$

$$\gamma_{i,j} = \frac{f_{i,j}}{T_0}$$

## Solution for Divisible Loads

Imbalance of a worker $\delta_i$

### Linear program

$\text{Minimize } T,$
$\text{under the constraints}$
$$\begin{cases} (1a) & |\delta_i| \leq T \times b_i \\ (1b) & \delta_i \geq \alpha_i - T \times s_i \\ (1c) & \sum_i \delta_i = 0 \end{cases}$$

$$(1)$$

Fraction of load $f_{i,j}$

$$f_{i,j} = \delta_i \times \frac{\delta_j}{\sum_{k \in R} \delta_k} = \delta_i \times \frac{\delta_j}{-L}$$

Communication rate $\lambda_{i,j}$

$$\lambda_{i,j} = \frac{f_{i,j}}{T_0}$$

Computation rate $\gamma_{i,j}$

$$\gamma_{i,j} = \frac{f_{i,j}}{T_0}$$

# Outline

# Conclusion

## Complete study of a difficult load-balancing problem

Scheduling and redistributing data on master-slave platforms

Independent tasks:

- General case: Proof of NP-completeness in the strong sense
- Special platforms: Optimal algorithms
- Simulations: Verification of theoretical results

Divisible loads:

- Solution for general case: LP + analytical formulas

## Perspectives

Beyond the NP-completeness: Search for approximation algorithms

Extension to dynamic master-slave platforms

Extension to more general interconnection networks

# Conclusion

## Complete study of a difficult load-balancing problem

Scheduling and redistributing data on master-slave platforms

Independent tasks:

- General case: Proof of NP-completeness in the strong sense
- Special platforms: Optimal algorithms
- Simulations: Verification of theoretical results

Divisible loads:

- Solution for general case: LP + analytical formulas

## Perspectives

Beyond the NP-completeness: Search for approximation algorithms

Extension to dynamic master-slave platforms

Extension to more general interconnection networks

## Recent work at our place

#### On-line scheduling heuristics for master-slave platforms

- Competitive ratios and inapproximability results
- Communication-aware heuristics

#### Collective communications

- Broadcast, multicast on heterogeneous clusters
- Resource selection for future MPI2 routines

#### Load-balancing

- Optimize BOINC-like applications
- Data redistribution strategies

#### Steady-state scheduling

- Multiple applications competing for resources
- Centralized vs fully distributed heuristics

## Scheduling for large-scale platforms

#### Assess the impact of new architectural characteristics

- Heterogeneity
- Irregular network topologies
- Hierarchy
- Variability (volability)

#### Inject **static** knowledge in a (mostly) dynamic environment

- Divisible loads vs bag of tasks
- Steady-state scheduling
- Resource selection

#### Evaluation

- Evaluate strategies through simulation
- SimGrid software co-developed with UCSD
- Large-scale experiments with Grid'5000