

Scheduling the I/O of HPC applications under congestion

Ana Gainaru, Guillaume Aupy, Anne Benoit,
Yves Robert, Franck Cappello & Marc Snir

IPDPS - May 2015

Motivation

Model

Platform Applications Objectives

Algorithms

Simulations

Applications

Assessment of heuristics

Experiments

Conclusion

1 Motivation

2 Model

Platform
Applications
Objectives

3 Algorithms

4 Simulations

- Applications
- Assessment of heuristics

5 Experiments

6 Conclusion

G. Aupy

Motivation

Model

- Platform
- Applications
- Objectives

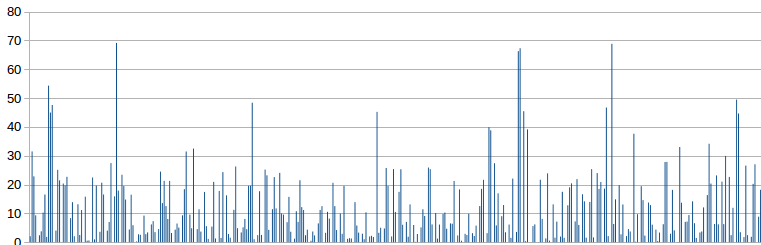
Algorithms

Simulations

- Applications
- Assessment of heuristics

Experiments

Conclusion



Analysis of the Intrepid system @Argonne: I/O throughput decrease (percentage per application, over 400 applications).

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

1 Motivation

2 Model
Platform
Applications
Objectives

3 Algorithms

4 Simulations
Applications
Assessment of heuristics

5 Experiments

6 Conclusion

G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

Simulations

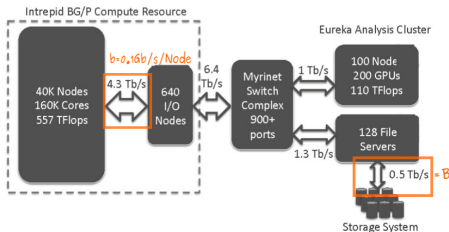
Applications

Assessment of
heuristics

Experiments

Conclusion

- N unit-speed processors, equipped with an I/O card of bandwidth b
- Centralized I/O system with total bandwidth B



Model instantiation for the Intrepid platform.

G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion

K applications competing for I/O. For application $\text{App}^{(k)}$:

- Released at time r_k ;
- Executed on $\beta^{(k)}$ procs;
- $n_{\text{tot}}^{(k)}$ instances: $\mathcal{I}_i^{(k)}$ consists of $w^{(k,i)}$ units of computation followed by the transfer of a volume $\text{vol}_{\text{io}}^{(k,i)}$;
- The minimum time to execute $\text{vol}_{\text{io}}^{(k,i)}$ is:

$$\text{time}_{\text{io}}^{(k,i)} = \frac{\text{vol}_{\text{io}}^{(k,i)}}{\min(\beta^{(k)} b, B)};$$

- Last instance finishes at time d_k .

G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

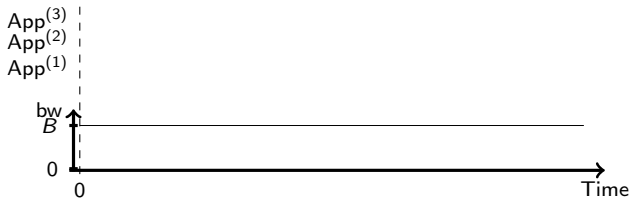
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

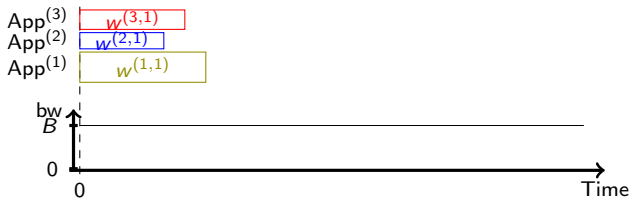
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

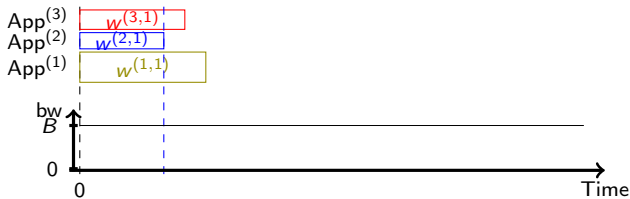
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

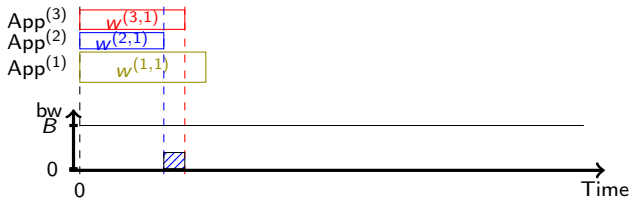
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

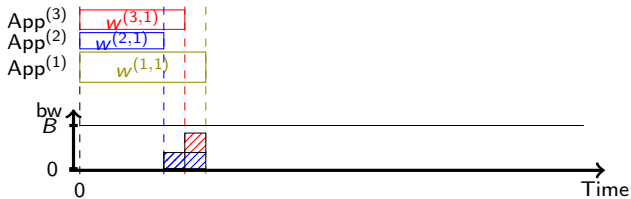
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

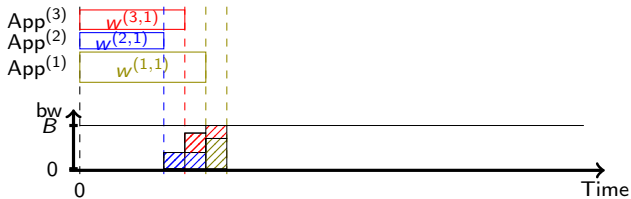
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

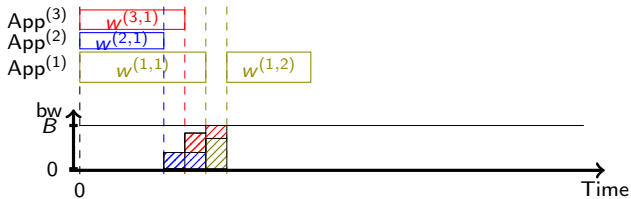
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

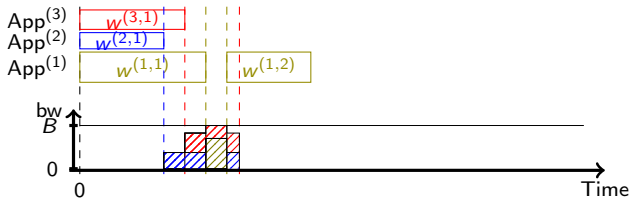
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

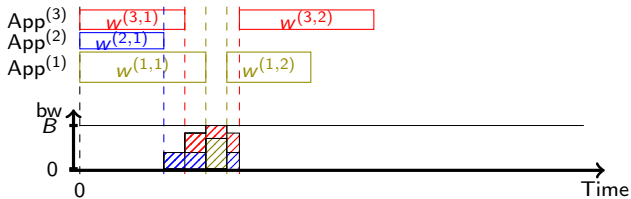
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

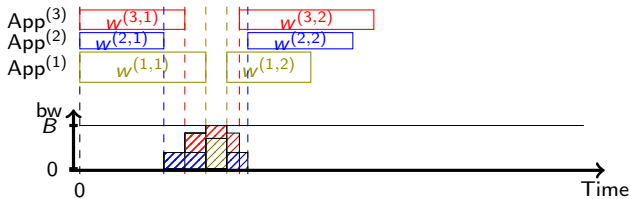
Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform

Applications

Objectives

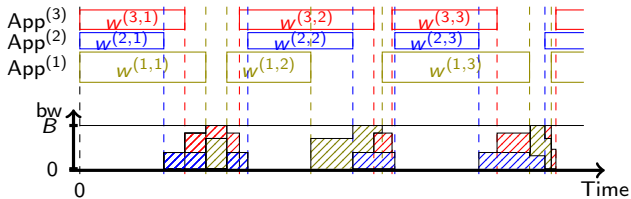
Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion



G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

Definition (Application efficiency)

$$\tilde{\rho}^{(k)}(t) = \frac{\sum_{i \leq n^{(k)}(t)} w^{(k,i)}}{t - r_k},$$

where $n^{(k)}(t)$ is the number of instances of $\text{App}^{(k)}$ executed at time t .

G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

Definition (Application efficiency)

$$\tilde{\rho}^{(k)}(t) = \frac{\sum_{i \leq n^{(k)}(t)} w^{(k,i)}}{t - r_k},$$

where $n^{(k)}(t)$ is the number of instances of $\text{App}^{(k)}$ executed at time t .

Obviously: $t - r_k \geq \sum_{i \leq n^{(k)}(t)} \left(w^{(k,i)} + \text{time}_{\text{io}}^{(k,i)} \right)$.

Hence:

$$\tilde{\rho}^{(k)}(t) \leq \rho^{(k)}(t) = \frac{\sum_{i \leq n^{(k)}(t)} w^{(k,i)}}{\sum_{i \leq n^{(k)}(t)} \left(w^{(k,i)} + \text{time}_{\text{io}}^{(k,i)} \right)}.$$

G. Aupy

Motivation

Model

Platform

Applications

Objectives

Algorithms

Simulations

Applications

Assessment of
heuristics

Experiments

Conclusion

- SYSEFFICIENCY:

$$\text{maximize } \frac{1}{N} \sum_{k=1}^K \beta^{(k)} \tilde{\rho}^{(k)}(d_k).$$

- DILATION:

$$\text{minimize } \max_{k=1..K} \frac{\rho^{(k)}(d_k)}{\tilde{\rho}^{(k)}(d_k)}.$$

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

1 Motivation

2 Model

Platform
Applications
Objectives

3 Algorithms

4 Simulations

Applications
Assessment of heuristics

5 Experiments

6 Conclusion

G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

The scheduler monitors the stream of I/O calls; decides on the fly which applications can perform I/O.

- At each time step, it has access to the state of the system (each application efficiency, $\tilde{\rho}^{(k)}$).
- Based on a given strategy, chooses a subset of applications that are allowed to perform I/O.

G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

The scheduler monitors the stream of I/O calls; decides on the fly which applications can perform I/O.

- At each time step, it has access to the state of the system (each application efficiency, $\tilde{\rho}^{(k)}$).
- Based on a given strategy, chooses a subset of applications that are allowed to perform I/O.

When a strategy *favors* $\text{App}^{(k)}$, it means that $\text{App}^{(k)}$ is executed as fast as possible ($\min(b\beta^{(k)}, \text{bw}_{\text{avail}})$).

G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

- **ROUNDROBIN**: Similar to the current scheduler in HPC systems. Applications are served following the “First-Come, First Served” principle.

G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

- **ROUNDROBIN**: Similar to the current scheduler in HPC systems. Applications are served following the “First-Come, First Served” principle.
- **MINDILATION**: favors applications with high values of $\frac{\rho^{(k)}(t)}{\tilde{\rho}^{(k)}(t)}$.

G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

- **ROUNDROBIN**: Similar to the current scheduler in HPC systems. Applications are served following the “First-Come, First Served” principle.
- **MINDILATION**: favors applications with high values of $\frac{\rho^{(k)}(t)}{\tilde{\rho}^{(k)}(t)}$.
- **MAXSYSEFF**: favors applications with low values of $\beta^{(k)}\tilde{\rho}^{(k)}(t)$.

G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

- **ROUNDROBIN**: Similar to the current scheduler in HPC systems. Applications are served following the “First-Come, First Served” principle.
- **MINDILATION**: favors applications with high values of $\frac{\rho^{(k)}(t)}{\tilde{\rho}^{(k)}(t)}$.
- **MAXSYSEFF**: favors applications with low values of $\beta^{(k)} \tilde{\rho}^{(k)}(t)$.
- **MINMAX- γ** : same as **MAXSYSEFF**, unless there exists an applications with $\frac{\tilde{\rho}^{(k)}(t)}{\rho^{(k)}(t)}$ below a threshold γ . In that case, switches to **MINDILATION**.

G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

- **ROUNDROBIN**: Similar to the current scheduler in HPC systems. Applications are served following the “First-Come, First Served” principle.
- **MINDILATION**: favors applications with high values of $\frac{\rho^{(k)}(t)}{\tilde{\rho}^{(k)}(t)}$.
- **MAXSYSEFF**: favors applications with low values of $\beta^{(k)}\tilde{\rho}^{(k)}(t)$.
- **MINMAX- γ** : same as **MAXSYSEFF**, unless there exists an applications with $\frac{\tilde{\rho}^{(k)}(t)}{\rho^{(k)}(t)}$ below a threshold γ . In that case, switches to **MINDILATION**.

PRIORITY variant: if an application has started to do some I/O, then it is prioritized.

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

1 Motivation

2 Model

Platform
Applications
Objectives

3 Algorithms

4 Simulations
Applications
Assessment of heuristics

5 Experiments

6 Conclusion

G. Aupy

We use Darshan to capture the behavior of applications that ran on Intrepid (2013).

Motivation

Model

Platform
Applications
Objectives

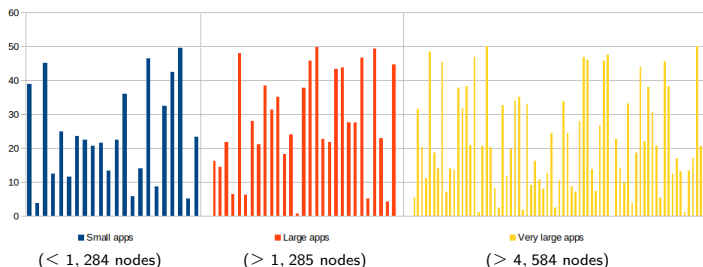
Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion



Percentage time spent doing I/O per application type.

G. Aupy

We use Darshan to capture the behavior of applications that ran on Intrepid (2013).

Motivation

Model

Platform
Applications
Objectives

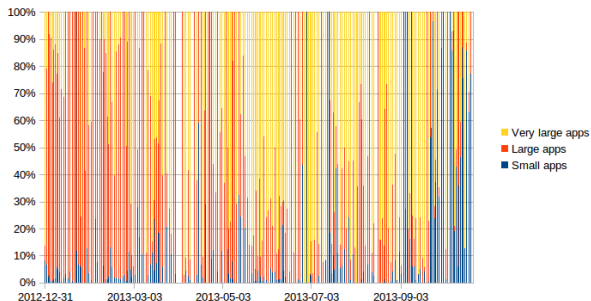
Algorithms

Simulations

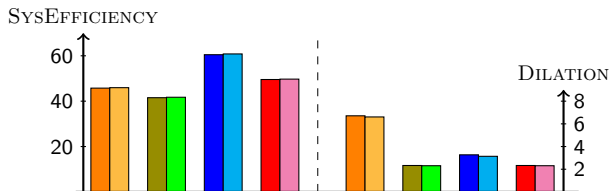
Applications
Assessment of
heuristics

Experiments

Conclusion

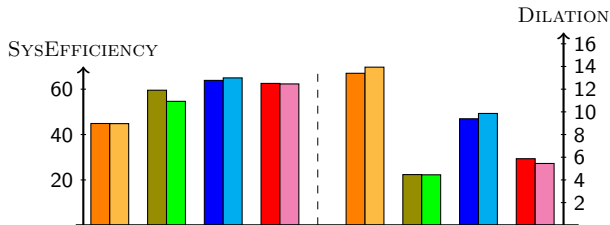


System usage per day for each application type



(a) 10 large applications, ratio of 20%

Objectives for different mixes of applications and I/O
computation ratios.

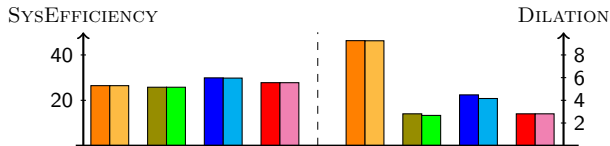


(b) 50 small and 5 large applications, ratio of 20%

Objectives for different mixes of applications and I/O computation ratios.

■ ROUNDROBIN
■ PRIORITY-ROUNDROBIN
■ MAXSYSEFF
■ PRIORITY-MAXSYSEFF

■ MINDILATION
■ PRIORITY-MINDILATION
■ MINMAX- γ
■ PRIORITY-MINMAX- γ



(c) 50 small and 5 large applications, ratio of 35% Objectives for different mixes of applications and I/O computation ratios.

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

We then compared our results with the Intrepid and Mira scheduler when congestion occurs.

Note that Intrepid and Mira use an architectural enhancement to improve the behavior of applications with large bursts of I/O: *Burst Buffers*.

Comparison of the heuristics on current platforms

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

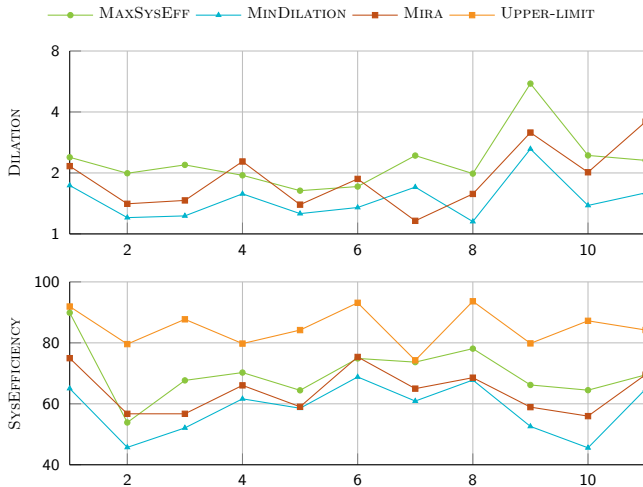
Applications
Assessment of
heuristics

Experiments

Conclusion

	DILATION (minimize)	SYS EFFICIENCY (maximize)
MAXSYSEFF	2.46	85.35
PRIORITY variant	3.13	82.98
MINMAX-0.25	2.33	83.08
PRIORITY variant	2.93	80.31
MINMAX-0.5	1.99	77.2
PRIORITY variant	2.43	72.96
MINMAX-0.75	1.69	71.66
PRIORITY variant	2.03	66.94
MINDILATION	1.63	70.45
PRIORITY variant	1.96	65.64
INTREPID	2.55	71.12
UPPER-LIMIT	-	91.59

Table: Averages over 56 different congested moments on Intrepid.



Comparison of the PRIORITY heuristics over the current
DILATION and SYSEFFICIENCY of Mira.

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

1 Motivation

2 Model

Platform
Applications
Objectives

3 Algorithms

4 Simulations

Applications
Assessment of heuristics

5 Experiments

6 Conclusion

Model

Platform Applications Objectives

Algorithms

Simulations

Applications

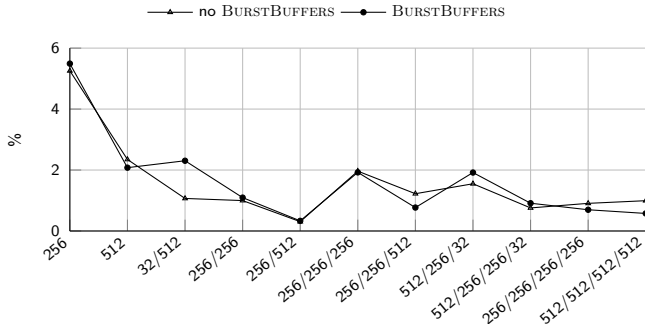
Assessment of heuristics

Experiments

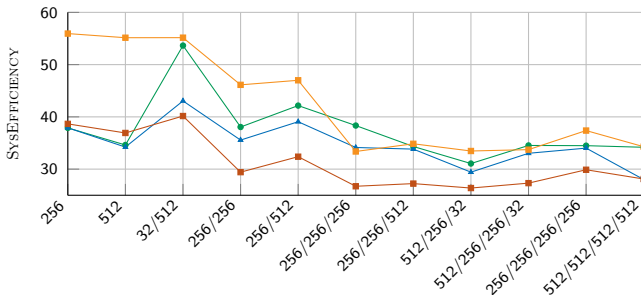
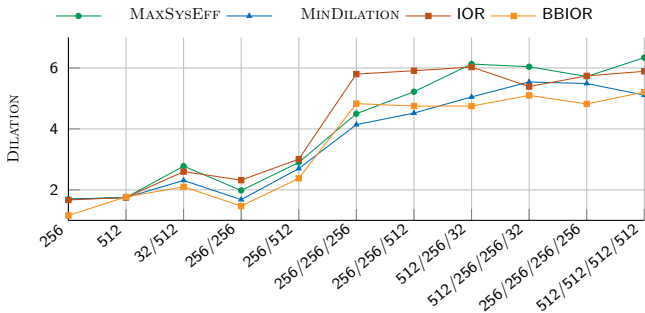
Conclusion

- Experiments on Vesta (development platform for Mira)
- Vesta is using hard disks and is affected by locality: we only used the PRIORITY variant of heuristics
- We implemented the heuristics as an additional layer on top of Vesta I/O scheduler

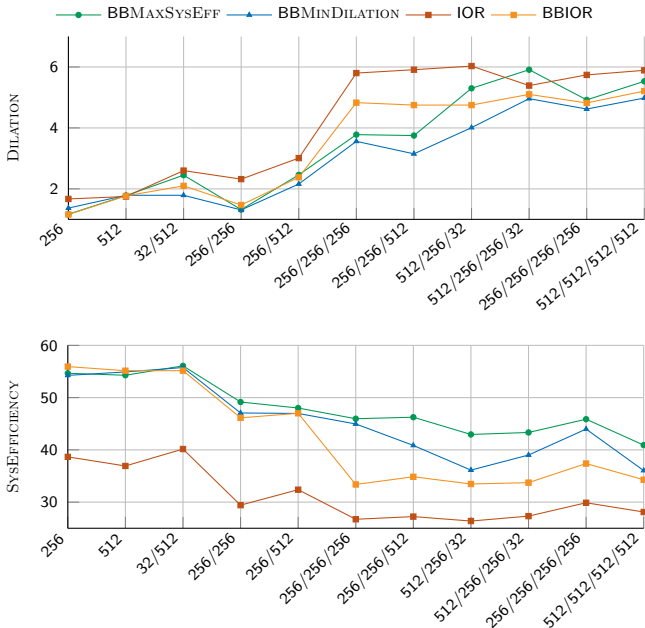
- Experiments on Vesta (development platform for Mira)
- Vesta is using hard disks and is affected by locality: we only used the PRIORITY variant of heuristics
- We implemented the heuristics as an additional layer on top of Vesta I/O scheduler



Execution time overhead of our implementation of the IOR benchmark.



System efficiency and dilation for different scenarios on Vesta.



System efficiency and dilation for different scenarios on Vesta.

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

- 1 Motivation
- 2 Model
 - Platform
 - Applications
 - Objectives
- 3 Algorithms
- 4 Simulations
 - Applications
 - Assessment of heuristics
- 5 Experiments
- 6 Conclusion

G. Aupy

Motivation

Model

Platform
Applications
Objectives

Algorithms

Simulations

Applications
Assessment of
heuristics

Experiments

Conclusion

- New I/O scheduler taking global view of system into account
- Outperforms current scheduler
- More experiments needed on larger application sets
- Window-based schedules for periodic applications?