

A Moldable Online Scheduling Algorithm and Its Application to Parallel Short Sequence Mapping

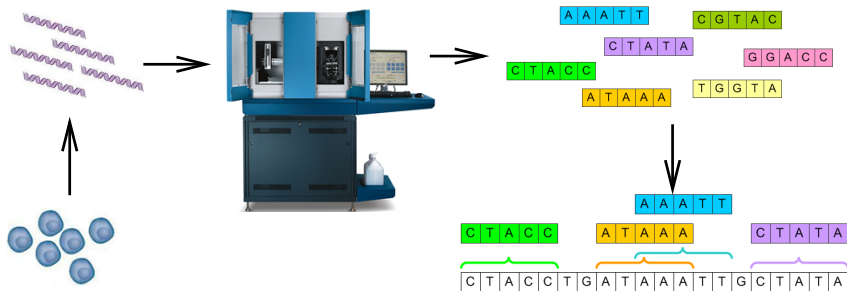
Erik Saule, Doruk Bozdağ, Umit V. Catalyurek

Department of Biomedical Informatics, The Ohio State University
{esaule,bozdagd,umit}@bmi.osu.edu

Scheduling for Large Scale Systems, May 2009

Supported by the U.S. DOE SciDAC Institute, the U.S. National Science Foundation and the Ohio Supercomputing Center

Motivation



Sequencing

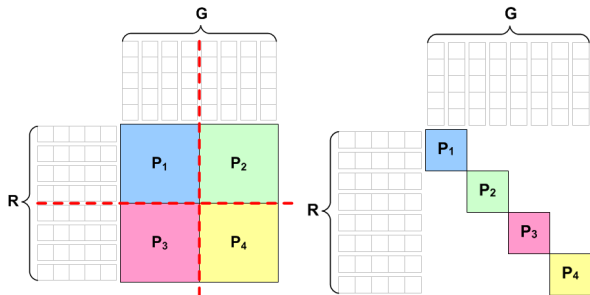
- Next generation sequencing instruments (SOLiD, Solexa, 454) can sequence up to 1 billion bases a day
 - Hundreds of millions of 35-50 base reads

Mapping

- Map reads to a reference genome efficiently (Human genome: 3Gb)
- Sequential mapping takes about a day
- Need fast, parallel algorithms that can handle mismatches

Parallel Short Sequence Mapping [Bozdag *et al.*, IPDPS 09]

Three partitioning dimensions:



$$P(m_g, m_r, m_s) = c_{gs} \frac{G}{m_g} + c_g \frac{G}{m_g m_s} + c_{rs} \frac{R}{m_r} + \left(c_r + c_c \frac{G}{m_g m_s} \right) \frac{R}{m_r m_s}$$

Partitioning on m processors is finding minimum $P(m_g, m_r, m_s)$ such that $m_g m_r m_s \leq m$

A cost efficient approach

To reduce cost, Ohio SuperComputing Center is building a bioscience dedicated cluster. It will host a Short Sequence Mapping service.

- Laboratories submits mapping request over the network.
- The service computes the mapping using the parallel algorithm.
- And sends the result back.

This talk

How to schedule the mapping request ?

This talk

A cost efficient approach

To reduce cost, Ohio SuperComputing Center is building a bioscience dedicated cluster. It will host a Short Sequence Mapping service.

- Laboratories submits mapping request over the network.
- The service computes the mapping using the parallel algorithm.
- And sends the result back.

This talk

How to schedule the mapping request ?

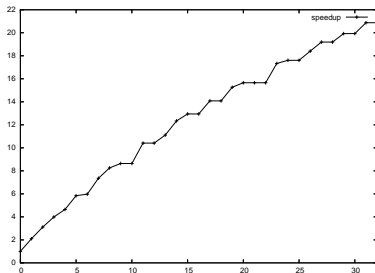
Outline of the Talk

- 1 Introduction
- 2 A Moldable Scheduling Problem
- 3 Deadline Based Online Scheduler (DBOS)
- 4 Experiments
- 5 Conclusion

Parallel Short Sequence Mapping

The important facts:

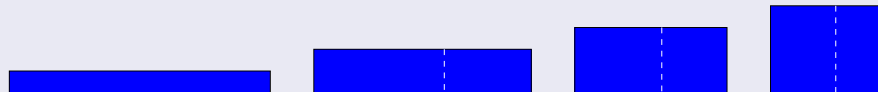
- can adapt to different number of processor
- good runtime prediction function
- no super linear speed up
- non convex speedup function (steps)
- no preemption



Moldable Scheduling

Instance

- m processors
- n tasks
- Task i arrives at r_i
- The execution of i on j processors takes $p_{i,j}$ time units



Solution

- Task i is executed on π_i processors
- Task i starts at σ_i
- Task i finishes at $C_i = \sigma_i + p_{i,\pi_i}$

Objective Function

Flow time

The flow time is the time spent in the system per a task $F_i = C_i - r_i$.

- Does not take task size into account.
- Optimizing the maximum flow time is unfair to small tasks.
- Optimizing the average flow time should starve large tasks.

Stretch [Bender *et al.* SoDA 98]

The stretch is the flow time normalized by the processing time of the task

$$s_i = \frac{C_i - r_i}{p_{i,1}}$$

- It provides a better fairness between tasks.
- Optimizing maximum stretch avoids starvation.

Objective Function

Flow time

The flow time is the time spent in the system per a task $F_i = C_i - r_i$.

- Does not take task size into account.
- Optimizing the maximum flow time is unfair to small tasks.
- Optimizing the average flow time should starve large tasks.

Stretch [Bender *et al.* SoDA 98]

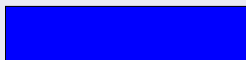
The stretch is the flow time normalized by the processing time of the task

$$s_i = \frac{C_i - r_i}{p_{i,1}}.$$

- It provides a better fairness between tasks.
- Optimizing maximum stretch avoids starvation.

Online maximum stretch can not be approximated

Adversary technique on one processor



A large task enters in the system

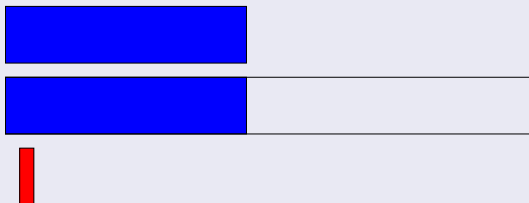
On several processors

There are similar techniques on several processors but there are more complicated and thus less prone to appear in practice.

The key point: if all processors are busy, a small task entering the system will have a large stretch.

Online maximum stretch can not be approximated

Adversary technique on one processor



If it is scheduled immediately, a small task is sent

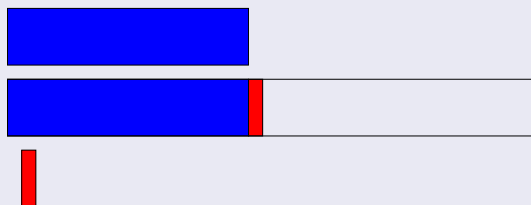
On several processors

There are similar techniques on several processors but there are more complicated and thus less prone to appear in practice.

The key point: if all processors are busy, a small task entering the system will have a large stretch.

Online maximum stretch can not be approximated

Adversary technique on one processor



It suffers a large delay (and an unbounded stretch)

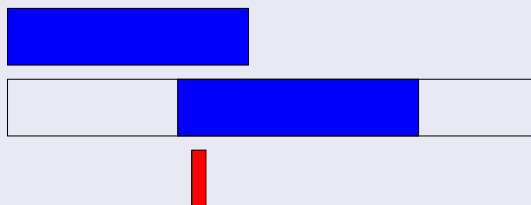
On several processors

There are similar techniques on several processors but there are more complicated and thus less prone to appear in practice.

The key point: if all processors are busy, a small task entering the system will have a large stretch.

Online maximum stretch can not be approximated

Adversary technique on one processor



If the large task is scheduled later, a small task is sent accordingly

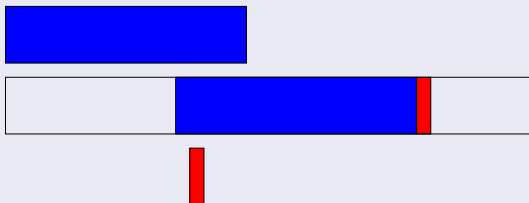
On several processors

There are similar techniques on several processors but there are more complicated and thus less prone to appear in practice.

The key point: if all processors are busy, a small task entering the system will have a large stretch.

Online maximum stretch can not be approximated

Adversary technique on one processor



It suffers a large delay (and an unbounded stretch)

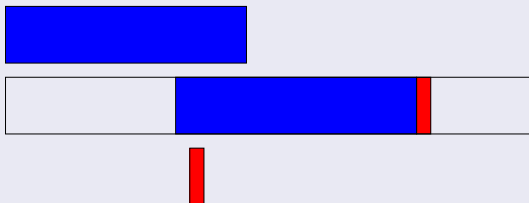
On several processors

There are similar techniques on several processors but there are more complicated and thus less prone to appear in practice.

The key point: if all processors are busy, a small task entering the system will have a large stretch.

Online maximum stretch can not be approximated

Adversary technique on one processor



It suffers a large delay (and an unbounded stretch)

On several processors

There are similar techniques on several processors but there are more complicated and thus less prone to appear in practice.

The key point: if all processors are busy, a small task entering the system will have a large stretch.

Outline of the Talk

- 1 Introduction
- 2 A Moldable Scheduling Problem
- 3 Deadline Based Online Scheduler (DBOS)**
- 4 Experiments
- 5 Conclusion

Principle of the Deadline Based Online Scheduler (DBOS)

- All tasks running concurrently should get the same stretch to maximize efficiency
- Using the optimal maximum stretch as an instant measure of the load
- Aim at a more efficient schedule than the optimal instant maximum stretch one to deal with still-to-arrive tasks

The DBOS Algorithm

Targeting a maximum stretch S

Task i must complete before the deadline $D_i = r_i + \rho_{i,1}S$.

Moldable Earliest Deadline First (MEDF)

- Considers task in deadline order.
- Allocates the minimum number of processors to each task to completes before the deadline.
- Schedules the task as soon as possible without moving any other task.

DBOS(ρ)

- Estimate the optimal maximum stretch S^* using a binary search.
- The deadline problem is solved by MEDF.
- Build a schedule of good efficiency of stretch ρS^* .
 - ρ is the online parameter

The DBOS Algorithm

Targeting a maximum stretch S

Task i must complete before the deadline $D_i = r_i + \rho_{i,1}S$.

Moldable Earliest Deadline First (MEDF)

- Considers task in deadline order.
- Allocates the minimum number of processors to each task to completes before the deadline.
- Schedules the task as soon as possible without moving any other task.

DBOS(ρ)

- Estimate the optimal maximum stretch S^* using a binary search.
- The deadline problem is solved by MEDF.
- Build a schedule of good efficiency of stretch ρS^* .
 - ρ is the online parameter

The DBOS Algorithm

Targeting a maximum stretch S

Task i must complete before the deadline $D_i = r_i + p_{i,1}S$.

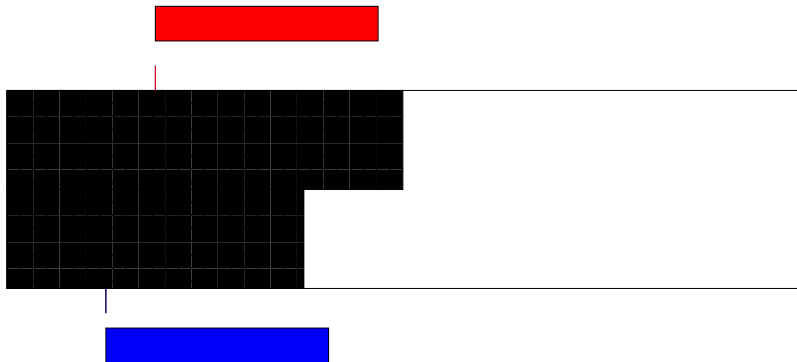
Moldable Earliest Deadline First (MEDF)

- Considers task in deadline order.
- Allocates the minimum number of processors to each task to completes before the deadline.
- Schedules the task as soon as possible without moving any other task.

DBOS(ρ)

- Estimate the optimal maximum stretch S^* using a binary search.
- The deadline problem is solved by MEDF.
- Build a schedule of good efficiency of stretch ρS^* .
 - ρ is the online parameter

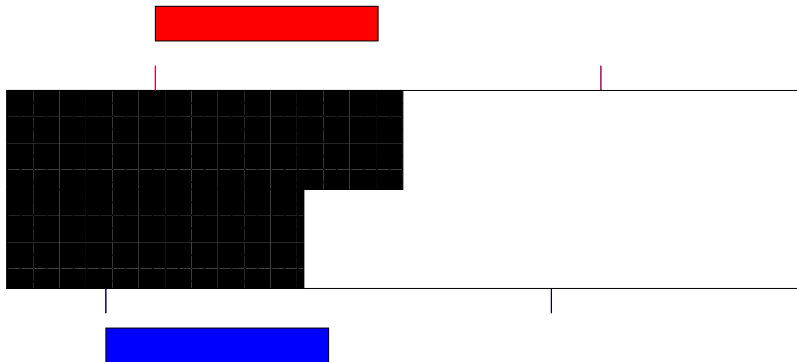
An example



A system with two pending tasks

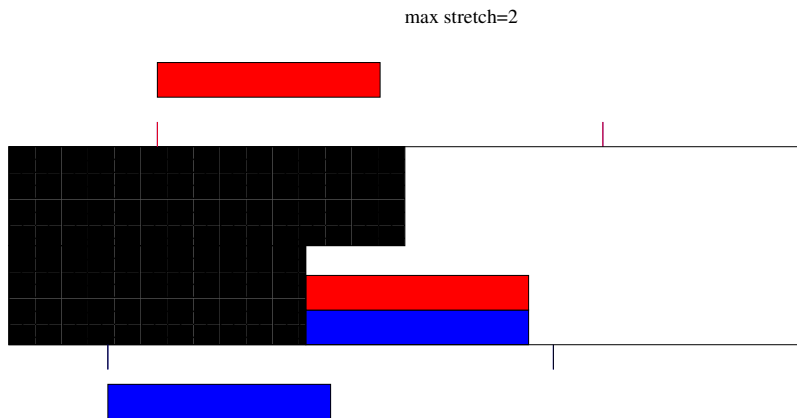
An example

max stretch=2



Deadlines induced by a stretch of 2

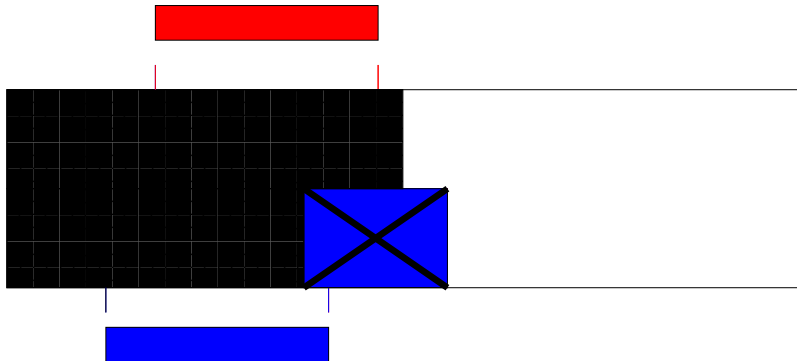
An example



A maximum stretch of 2 is reachable

An example

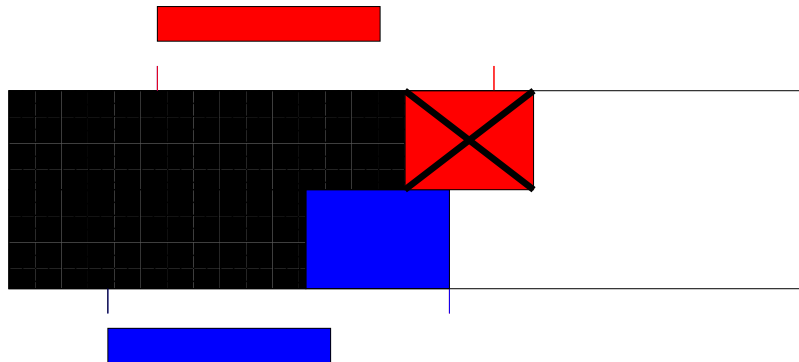
max stretch=1



But 1 is not

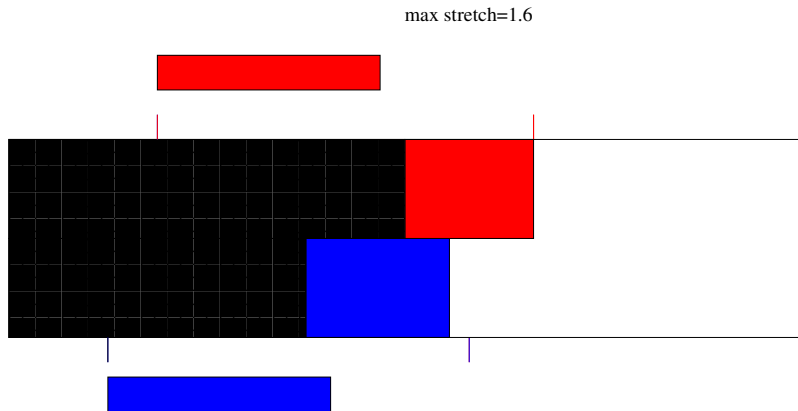
An example

max stretch=1.5



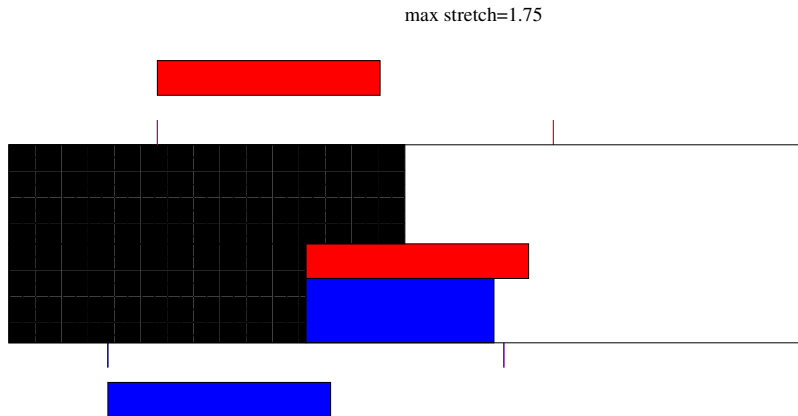
Neither 1.5

An example



The optimal stretch is 1.6

An example



The online parameter $\rho = 1.1$ leaves much more space (thanks to MEDF).

Outline of the Talk

- 1 Introduction
- 2 A Moldable Scheduling Problem
- 3 Deadline Based Online Scheduler (DBOS)
- 4 Experiments**
- 5 Conclusion

The algorithm

- $\forall i, \pi_i \leftarrow 1, \text{mark}[i] \leftarrow \text{false}$
- $\sigma \leftarrow \text{schedule}(\pi)$
- while $\exists i \mid \text{mark}[i] = \text{false}$
 - Get unmarked i such that $p_{i,\pi_i} - p_{i,\pi_i+1}$ is maximal and positive
 - $\pi_i \leftarrow \pi_i + 1$
 - $\sigma' \leftarrow \text{schedule}(\pi)$
 - if $\text{avgflow}(\sigma') < \text{avgflow}(\sigma)$
 - $\sigma \leftarrow \sigma'$
 - else
 - $\pi_i \leftarrow \pi_i + 1; \text{mark}[i] \leftarrow \text{true}$

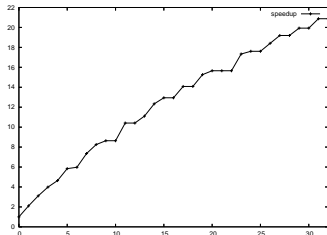
schedule

schedule is a conservative backfilling algorithm. Unspecified, we used FCFS.

An Iterative Process [Sabin et al, JSSPP 06]

Properties

- Optimizing flow time
- Claimed to outperform fair share
- Parameter-less



Improvement

If the speedup function is non convex or has steps. The algorithm gets stuck. (It was originally tested with a model where the speedup is convex)

Modification:

- Get unmarked i and k such that $(p_{i,\pi_i} - p_{i,\pi_i+k})/k$ is maximal and positive
- $\pi_i \leftarrow \pi_i + k$

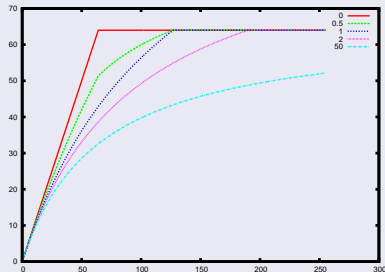
First Experimental Setting

Goal: assess performance on a well known setting

Downey model

Two parameters:

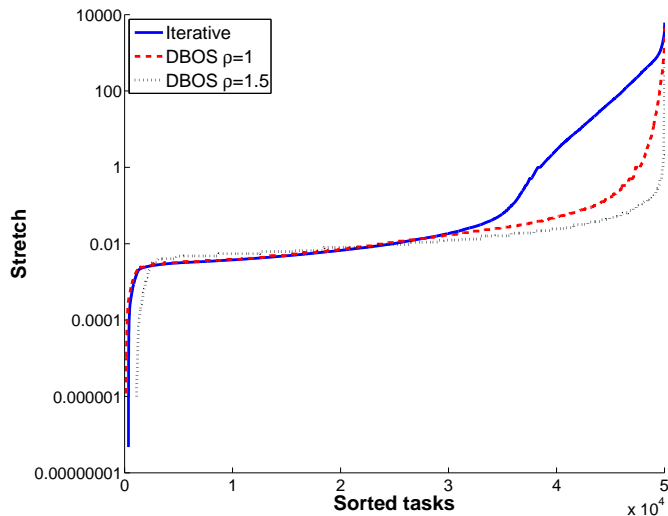
- Average parallelism (64)
- Distance to linear speedup



Generation

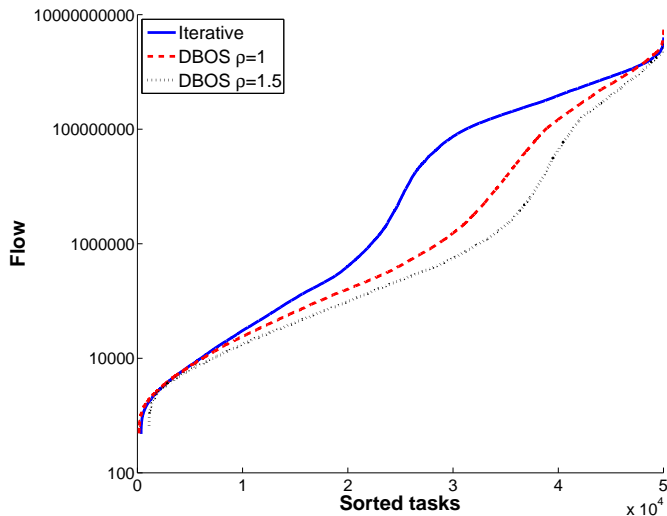
- 512 processors
- First 5000 tasks of SDSC Par 96 (From the Feitelson archive)
- Sequential time : total execution time
- Average parallelism : between number of used processor and 512
- Distance to linear speedup : between 0 and 2

Downey model results



DBOS generates less tasks with high stretch.

Downey model results



DBOS leads to better flow time. Iterative could be improved.

Second Experimental Setting

Goal: test case reflecting the cluster usage

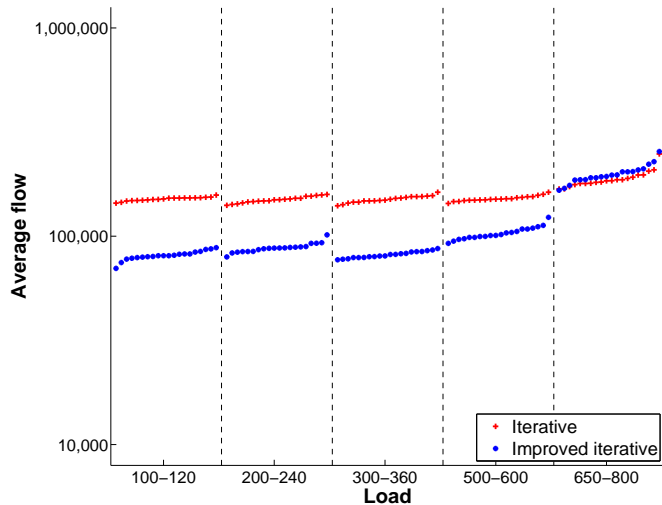
Generation

- 512 processors
- Each task corresponds to one lab studying one genome
- Speedup according to the runtime prediction function
- 2000 tasks are uniformly distributed in an time interval
- Changing the span of the interval to control the load

Real data

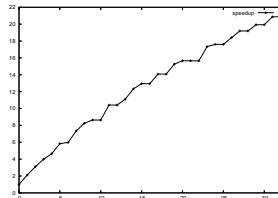
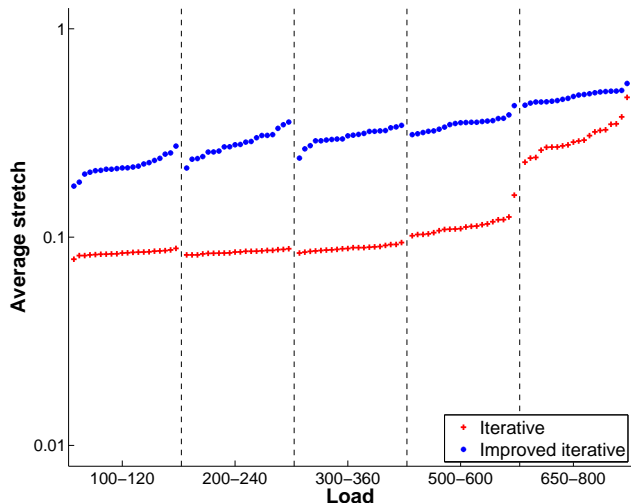
Sequencing machine	Reads	Genome	Size
454 GS FLX Genome Analyzer	1 million	E. Coli	4.6 million
Solexa IG sequencer	200 million	Yeast	15 million
SOLiD system	400 million	A. Thaliana	100 million
		Mosquito	280 million
		Rice	465 million
		Chicken	1.2 billion
		Human	3.4 billion

Mapping : Improvement on Iterative (flow)



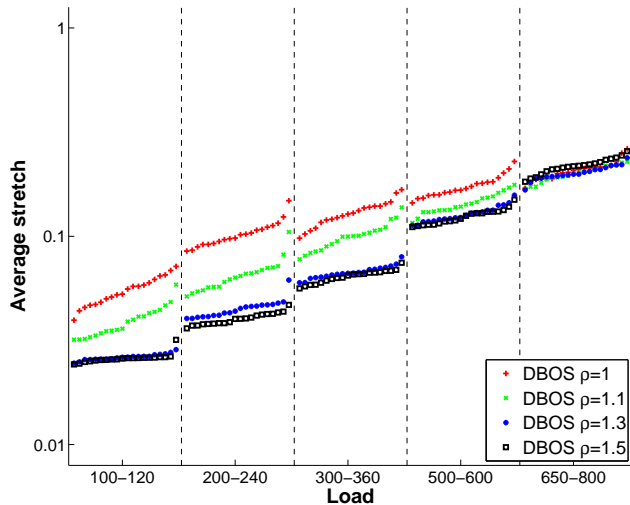
The improvement really improves. The iterative got stuck.

Mapping: Improvement on Iterative (stretch)



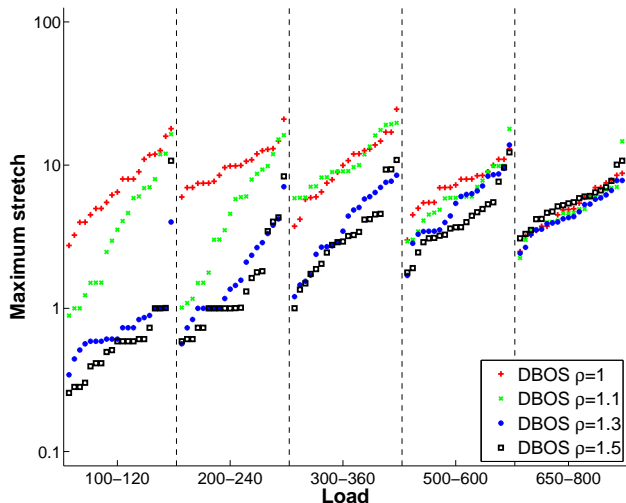
Getting stuck is good for stretch since it avoids interrupting tasks. They are just lucky.

Mapping : the online parameter (average stretch)



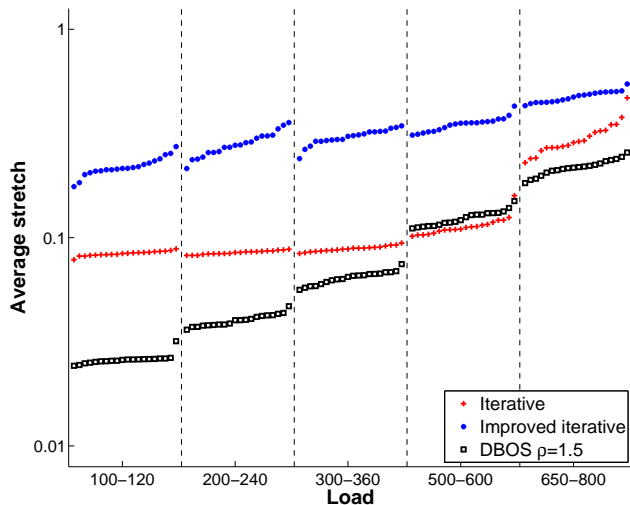
Quickly drops with ρ . Step at $\rho = 1.3$.

Mapping : the online parameter (maximum stretch)



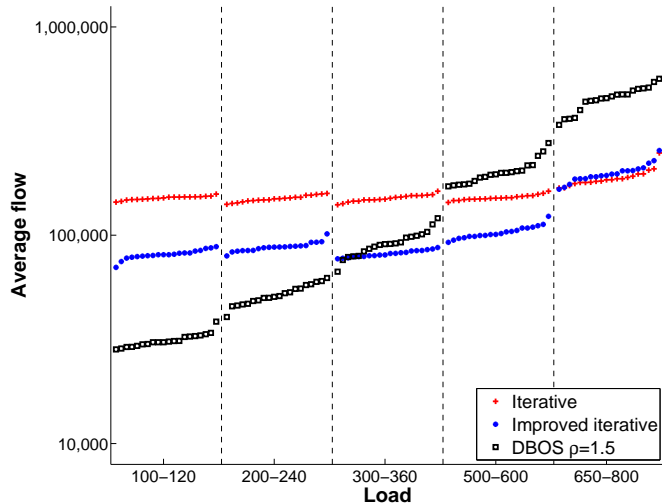
Max stretch is kept at a reasonable level. The online parameter ρ is very helpful here.

Mapping : DBOS vs Iterative (average stretch)



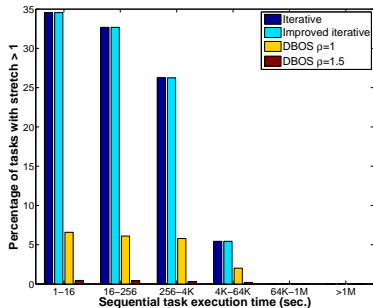
DBOS leads to much better stretch (even when iterative got stuck).

Mapping : DBOS vs Iterative (average flow)

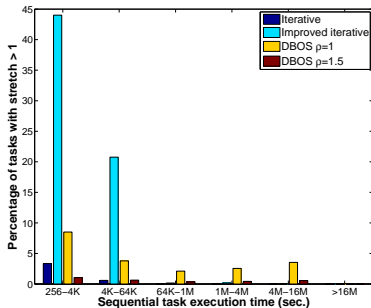


Confirm there is room for improvement for Iterative. DBOS is not bad.

Mapping : DBOS vs Iterative (Fairness Issues)



Downey model



Short Sequence Mapping

The Iterative algorithm leads to high stretch for a lot of the smaller tasks. DBOS has better performance and less fairness issues thanks to stretch optimization.

Outline of the Talk

- 1 Introduction
- 2 A Moldable Scheduling Problem
- 3 Deadline Based Online Scheduler (DBOS)
- 4 Experiments
- 5 Conclusion**

Conclusion

- A cluster dedicated to bioscience will be built.
- To provide fairness stretch should be considered instead of flow time.
- An scheduling algorithm is proposed to optimize stretch and avoid worst case online scenario.
- Which performs well on Short Sequence Mapping application.

Perspective

- Investigate other way to avoid worst case scenarios.
- Study more simple algorithm to get reference points.
- Build the service !

Conclusion

- A cluster dedicated to bioscience will be built.
- To provide fairness stretch should be considered instead of flow time.
- An scheduling algorithm is proposed to optimize stretch and avoid worst case online scenario.
- Which performs well on Short Sequence Mapping application.

Perspective

- Investigate other way to avoid worst case scenarios.
- Study more simple algorithm to get reference points.
- Build the service !