

How to deal with uncertainties and dynamicity ?

Frédéric Vivien

December 4, 2006

Outline

- 1 The problem
- 2 Analyzing sensitivity and robustness
- 3 Dynamic load-balancing
- 4 How useful is old information ?
- 5 Conclusion

Outline

- 1 The problem
- 2 Analyzing sensitivity and robustness
- 3 Dynamic load-balancing
- 4 How useful is old information ?
- 5 Conclusion

The problem : the world is not perfect !

▶ Uncertainties

- ▶ On the platforms' characteristics
(Processor power, link bandwidth, etc.)
- ▶ On the applications' characteristics
(Volume computation to be performed, volume of messages to be sent, etc.)

▶ Dynamicity

- ▶ Of network (interferences with other applications, etc.)
- ▶ Of processors (interferences with other users, other processors of the same node, other core of the same processor, etc.)
- ▶ Of applications (on which detail should the simulation focus?)

Solutions : to prevent or to cure ?

To prevent

- ▶ Algorithms tolerant to uncertainties and dynamicity.

To cure

- ▶ Algorithms auto-adapting to actual conditions.

Leitmotiv : the more the information, the more precise we can statically define the solutions, the better our chances to “succeed”

Outline

- 1 The problem
- 2 Analyzing sensitivity and robustness**
- 3 Dynamic load-balancing
- 4 How useful is old information ?
- 5 Conclusion

Analyzing the sensitivity

Question : we have defined a solution, how is it going to behave “in practice” ?

Possible approach

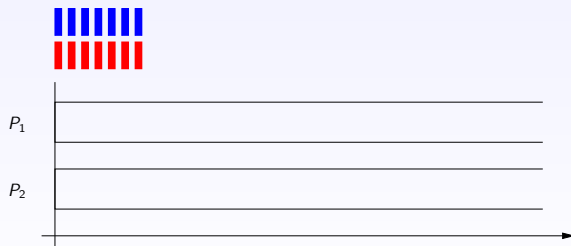
- 1 Definition of an algorithm \mathcal{A} .
- 2 Modeling the uncertainties and the dynamicity.
- 3 Analyzing the sensitivity of \mathcal{A} as follows :
 - ▶ For each theoretical instance of the problem
 - ▶ Evaluate the solution found by \mathcal{A}
 - ▶ For each “actual” instance corresponding to the given theoretical instance, find the optimal solution and the relative performance of the solution found by \mathcal{A} .

Sensitivity of \mathcal{A} : worst relative performance, or (weighted) average relative performance, etc.

Analyzing the sensitivity : an example

Problem

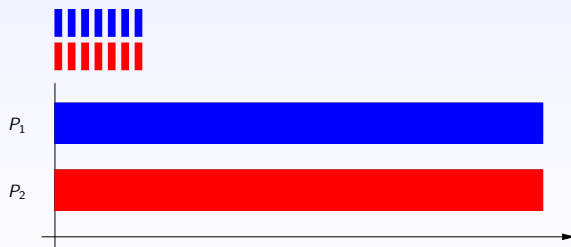
- ▶ Master-slave platform with two identical processors
- ▶ Flow of two types of identical tasks
- ▶ Objective function : maximum minimum throughput between the two applications (*max-min fairness*)



Analyzing the sensitivity : an example

Problem

- ▶ Master-slave platform with two identical processors
- ▶ Flow of two types of identical tasks
- ▶ Objective function : maximum minimum throughput between the two applications (*max-min fairness*)

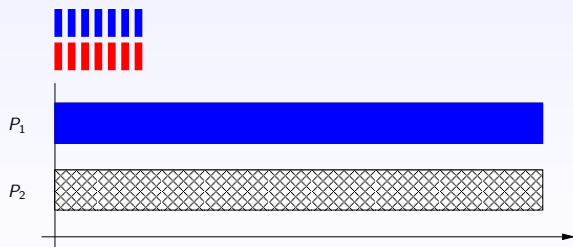


A possible solution...

Analyzing the sensitivity : an example

Problem

- ▶ Master-slave platform with two identical processors
- ▶ Flow of two types of identical tasks
- ▶ Objective function : maximum minimum throughput between the two applications (*max-min fairness*)



A possible solution... null if processor P_2 fails.

Analyzing the sensitivity : the case of Backfilling (1)

The EASY backfilling scheme

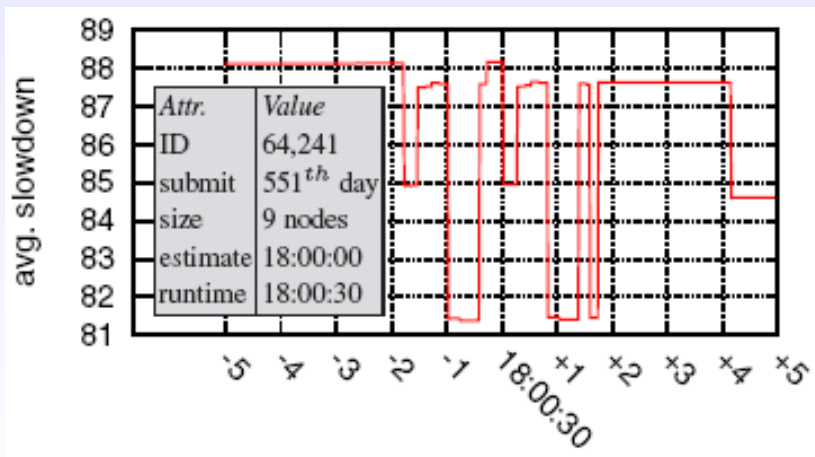
- ▶ The jobs are considered in First-come first-serve order
- ▶ Each time a job arrives or a job completes, a reservation is made for the first job that cannot be immediately started, later jobs that can be started immediately are started.
- ▶ In practice jobs are submitted with runtime estimates. A job exceeding its estimate is automatically killed.

Analyzing the sensitivity : the case of Backfilling (2)

The set-up

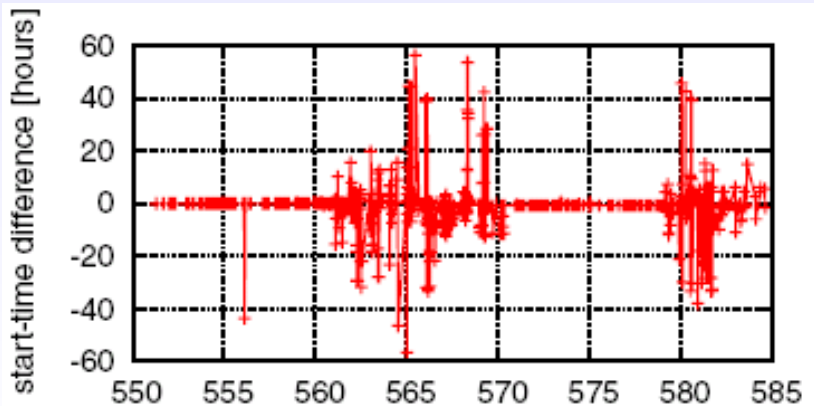
- ▶ 128-node IBM SP2 (San Diego Supercomputer Center)
- ▶ May 1998 to April 2000 log : 67,667 jobs
- ▶ Job runtime limit : 18 hours.
(Some dozens of seconds may be needed to kill a job.)
- ▶ Performance measure : average slowdown (=average stretch).

Analyzing the sensitivity : the case of Backfilling (3)

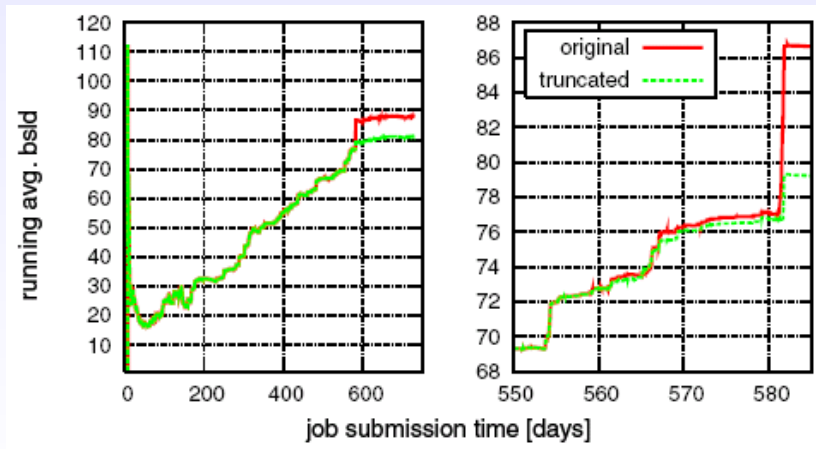


The length of a job running for 18 hours and 30 seconds is shortened by 30 seconds.

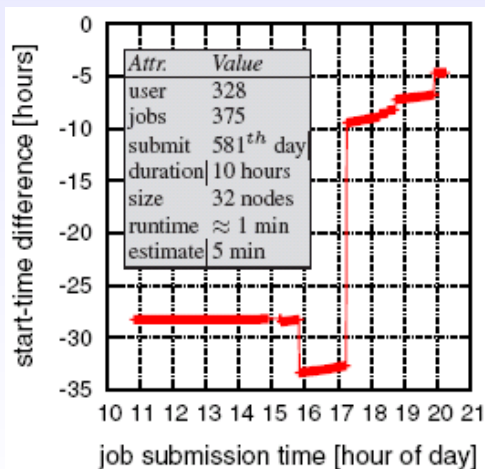
Analyzing the sensitivity : the case of Backfilling (3)



Analyzing the sensitivity : the case of Backfilling (3)



Analyzing the sensitivity : the case of Backfilling (3)

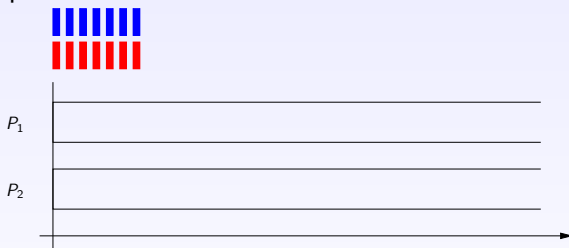


Robust solutions

An algorithm is said to be robust if its solutions stay close to the optimal when the actual parameters are slightly different from the theoretical parameters.

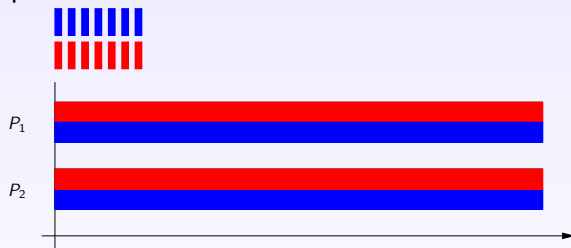
Robust slutions

An algorithm is said to be robust if its solutions stay close to the optimal when the actual parameters are slightly different from the theoretical parameters.



Robust slutions

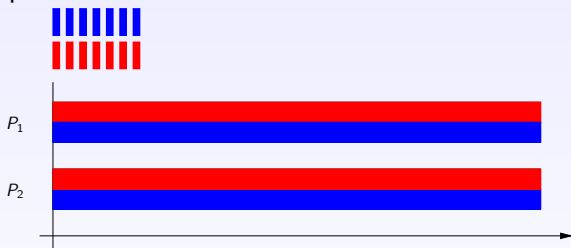
An algorithm is said to be robust if its solutions stay close to the optimal when the actual parameters are slightly different from the theoretical parameters.



This solution stays optimal whatever the variations in the processors' performance : it is not sensitive to this parameter !

Robust slutions

An algorithm is said to be robust if its solutions stay close to the optimal when the actual parameters are slightly different from the theoretical parameters.



This solution stays optimal whatever the variations in the processors' performance : it is not sensitive to this parameter !

It is at worst at two from the optimal when considering variations on the relative size of the two types of tasks.

An example of resolution : the problem

The problem :

- ▶ A master which has an output bandwidth B .

An example of resolution : the problem

The problem :

- ▶ A master which has an output bandwidth B .
- ▶ p slaves, slave P_i being link with a bandwidth b_i and having a computational speed of c_i .

An example of resolution : the problem

The problem :

- ▶ A master which has an output bandwidth B .
- ▶ p slaves, slave P_i being link with a bandwidth b_i and having a computational speed of c_i .
- ▶ n flows of tasks, each flow being a set of identical tasks.
A task from the flow k needs β_k units of communications and γ_k units of computation.

An example of resolution : the problem

The problem :

- ▶ A master which has an output bandwidth B .
- ▶ p slaves, slave P_i being link with a bandwidth b_i and having a computational speed of c_i .
- ▶ n flows of tasks, each flow being a set of identical tasks.
A task from the flow k needs β_k units of communications and γ_k units of computation.
- ▶ $\rho_i^{(k)}$: throughput of application k on processor i .
 $\rho^{(k)} = \sum_i \rho_i^{(k)}$ is the throughput of application k on the whole platform.

An example of resolution : the problem

The problem :

- ▶ A master which has an output bandwidth B .
- ▶ p slaves, slave P_i being link with a bandwidth b_i and having a computational speed of c_i .
- ▶ n flows of tasks, each flow being a set of identical tasks.
A task from the flow k needs β_k units of communications and γ_k units of computation.
- ▶ $\rho_i^{(k)}$: throughput of application k on processor i .
 $\rho^{(k)} = \sum_i \rho_i^{(k)}$ is the throughput of application k on the whole platform.
- ▶ Objective : maximize $\min_k \rho^{(k)}$.

An example of resolution : the problem

The problem :

- ▶ A master which has an output bandwidth B .
- ▶ p slaves, slave P_i being link with a bandwidth b_i and having a computational speed of c_i .
- ▶ n flows of tasks, each flow being a set of identical tasks.
A task from the flow k needs β_k units of communications and γ_k units of computation.
- ▶ $\rho_i^{(k)}$: throughput of application k on processor i .
 $\rho^{(k)} = \sum_i \rho_i^{(k)}$ is the throughput of application k on the whole platform.
- ▶ Objective : maximize $\min_k \rho^{(k)}$.

Dynamicity : a processor may fail (definitively).

An example of resolution : classical solution

- 1 Resource constraint : processors' capacities

$$\forall i, \quad \sum_k \rho_i^{(k)} \gamma_k \leq c_i$$

An example of resolution : classical solution

- 1 Resource constraint : processors' capacities

$$\forall i, \quad \sum_k \rho_i^{(k)} \gamma_k \leq c_i$$

- 2 Resource constraint : processors' bandwidths

$$\forall i, \quad \sum_k \rho_i^{(k)} \beta_k \leq b_i$$

An example of resolution : classical solution

- 1 Resource constraint : processors' capacities

$$\forall i, \quad \sum_k \rho_i^{(k)} \gamma_k \leq c_i$$

- 2 Resource constraint : processors' bandwidths

$$\forall i, \quad \sum_k \rho_i^{(k)} \beta_k \leq b_i$$

- 3 Resource constraint : the master output bandwidth

$$\sum_{i,k} \rho_i^{(k)} \beta_k \leq B$$

An example of resolution : classical solution

- 1 Resource constraint : processors' capacities

$$\forall i, \quad \sum_k \rho_i^{(k)} \gamma_k \leq c_i$$

- 2 Resource constraint : processors' bandwidths

$$\forall i, \quad \sum_k \rho_i^{(k)} \beta_k \leq b_i$$

- 3 Resource constraint : the master output bandwidth

$$\sum_{i,k} \rho_i^{(k)} \beta_k \leq B$$

- 4 Objective

$$\text{Maximize } \min_k \sum_i \rho_i^{(k)}$$

An example of resolution : robust solution

$$\textcircled{1} \quad \forall i, \quad \sum_k \rho_i^{(k)} \gamma_k \leq c_i$$

$$\textcircled{2} \quad \forall i, \quad \sum_k \rho_i^{(k)} \beta_k \leq b_i$$

$$\textcircled{3} \quad \sum_{i,k} \rho_i^{(k)} \beta_k \leq B$$

An example of resolution : robust solution

$$\textcircled{1} \quad \forall i, \quad \sum_k \rho_i^{(k)} \gamma_k \leq c_i$$

$$\textcircled{2} \quad \forall i, \quad \sum_k \rho_i^{(k)} \beta_k \leq b_i$$

$$\textcircled{3} \quad \sum_{i,k} \rho_i^{(k)} \beta_k \leq B$$

$\textcircled{4}$ Objective if processor P_p fails :

$$\rho_{\bar{p}} = \min_k \sum_{i \neq p} \rho_i^{(k)}$$

An example of resolution : robust solution

$$\textcircled{1} \quad \forall i, \quad \sum_k \rho_i^{(k)} \gamma_k \leq c_i$$

$$\textcircled{2} \quad \forall i, \quad \sum_k \rho_i^{(k)} \beta_k \leq b_i$$

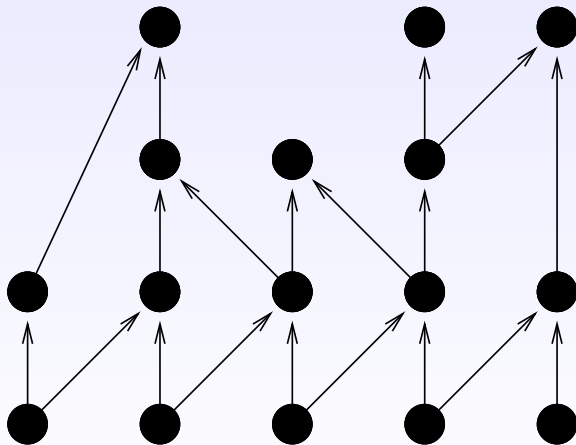
$$\textcircled{3} \quad \sum_{i,k} \rho_i^{(k)} \beta_k \leq B$$

$\textcircled{4}$ Objective if processor P_p fails :

$$\rho_{\bar{p}} = \min_k \sum_{i \neq p} \rho_i^{(k)}$$

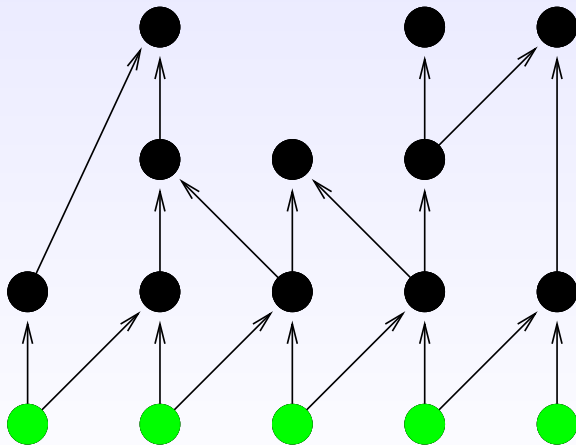
$$\textcircled{5} \quad \text{Objective : Maximize } \min \left\{ \min_p \frac{\rho_{\bar{p}}}{\rho_{\bar{p}}^{(\text{opt})}}, \min_k \sum_i \frac{\rho_i^{(k)}}{\rho^{(\text{opt})}} \right\}$$

Internet-Based Computing : problem motivation



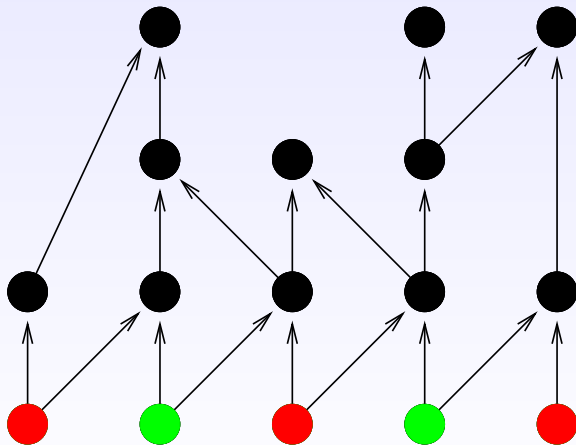
A possible schedule

Internet-Based Computing : problem motivation



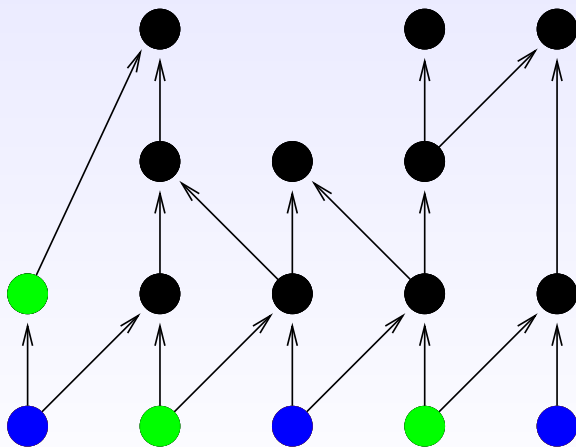
A possible schedule

Internet-Based Computing : problem motivation

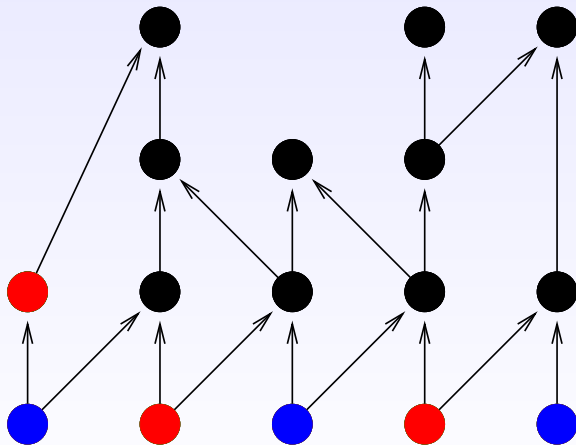


A possible schedule

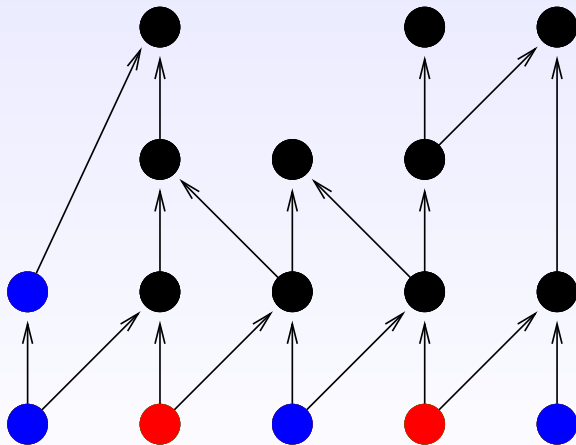
Internet-Based Computing : problem motivation



Internet-Based Computing : problem motivation

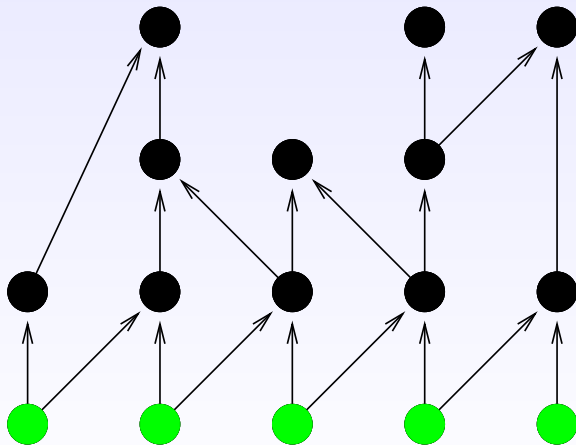


Internet-Based Computing : problem motivation



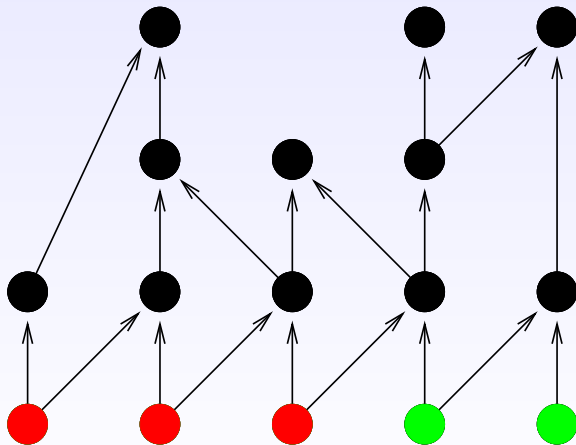
A possible schedule

Internet-Based Computing : problem motivation



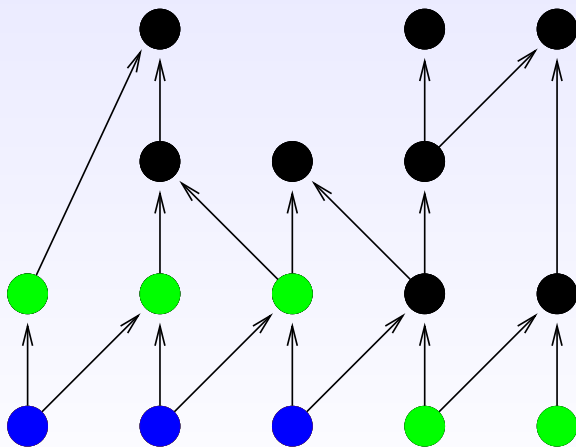
Another possible schedule

Internet-Based Computing : problem motivation



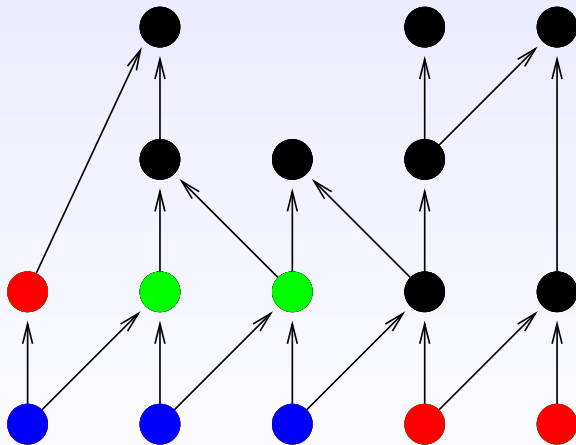
Another possible schedule

Internet-Based Computing : problem motivation



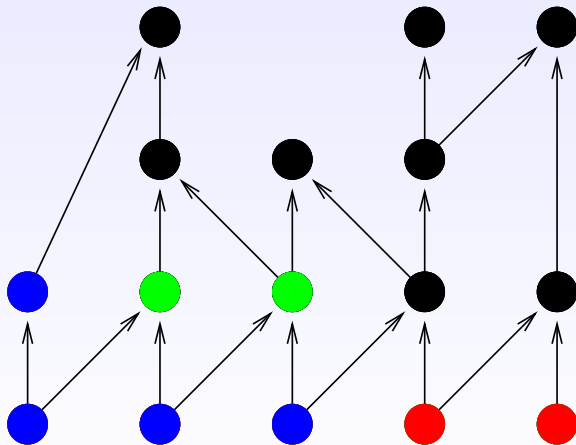
Another possible schedule

Internet-Based Computing : problem motivation



Another possible schedule

Internet-Based Computing : problem motivation



Another possible schedule

Internet-Based Computing : the principle

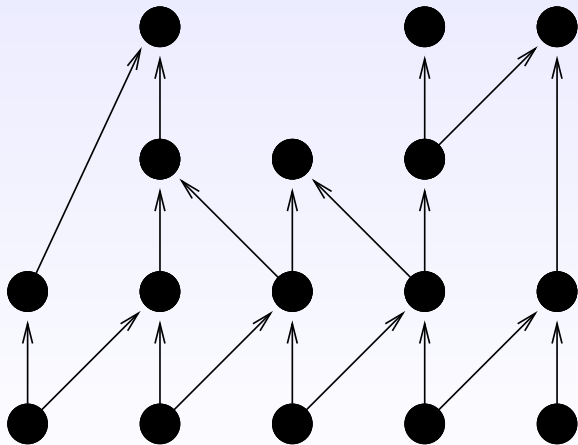
Framework

- ▶ Internet-Based computing : processor capacities are unknown, running times are thus unpredictable.
- ▶ Execution of a task graph (aka workflow) containing n tasks.

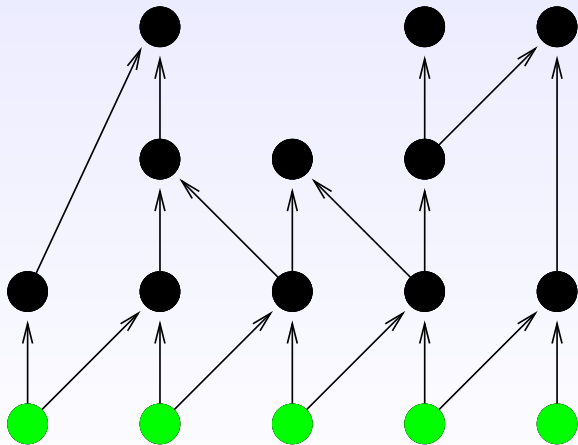
The principle

- ▶ Motivation : lessening the likelihood of the “gridlock” that can arise when a computation stalls pending computation of already allocated tasks.
- ▶ IC-optimal schedule : after t tasks have been executed, the number of eligible (=executable) tasks is maximal, for any $t \in [1, n]$

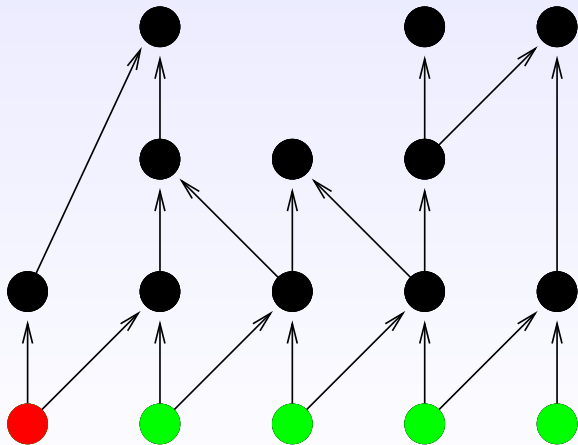
Internet-Based Computing : back to the example



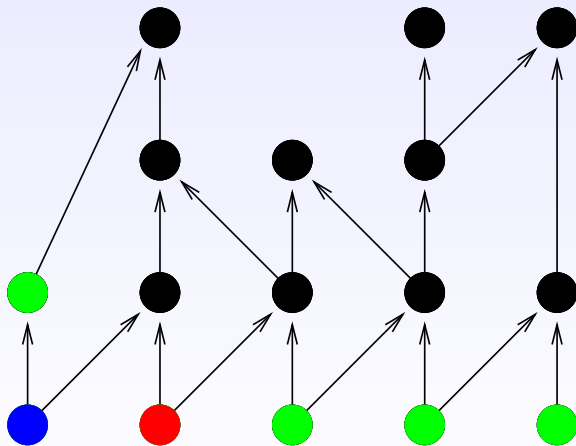
Internet-Based Computing : back to the example



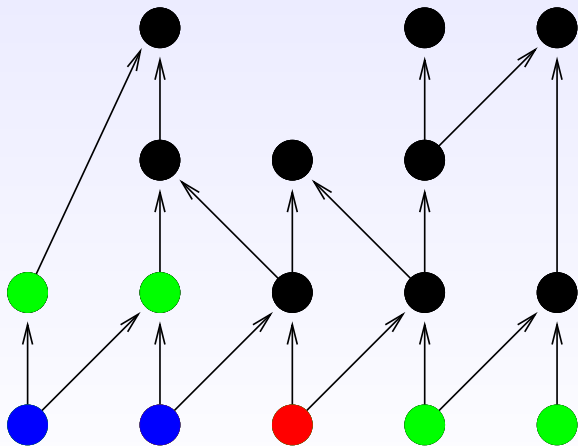
Internet-Based Computing : back to the example



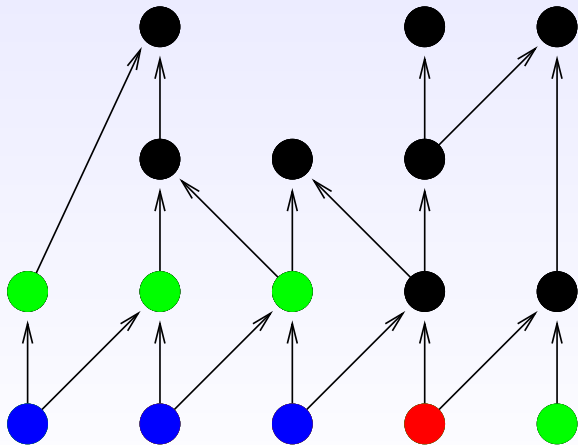
Internet-Based Computing : back to the example



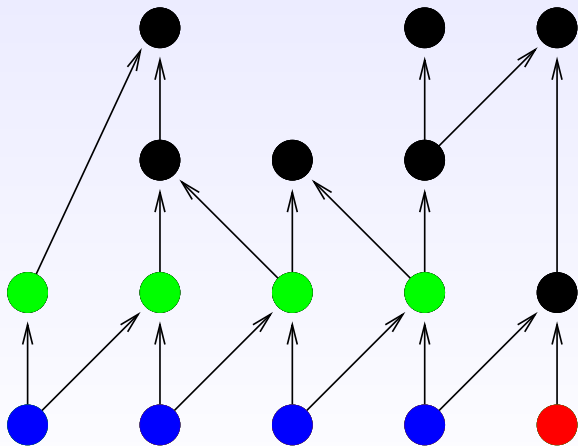
Internet-Based Computing : back to the example



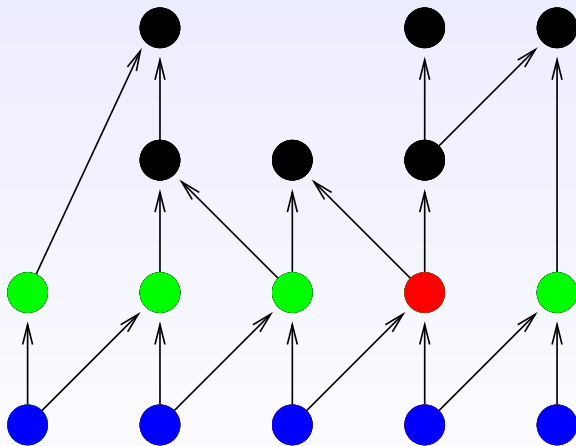
Internet-Based Computing : back to the example



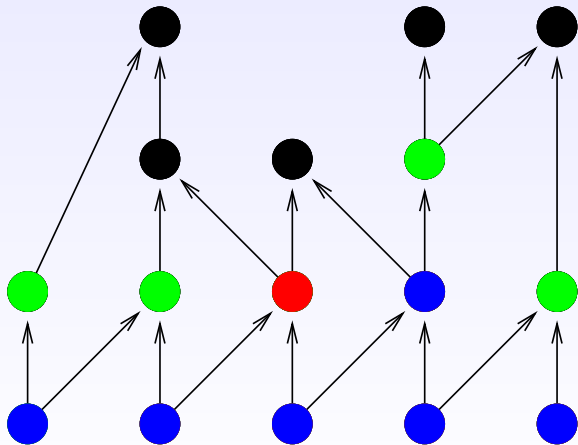
Internet-Based Computing : back to the example



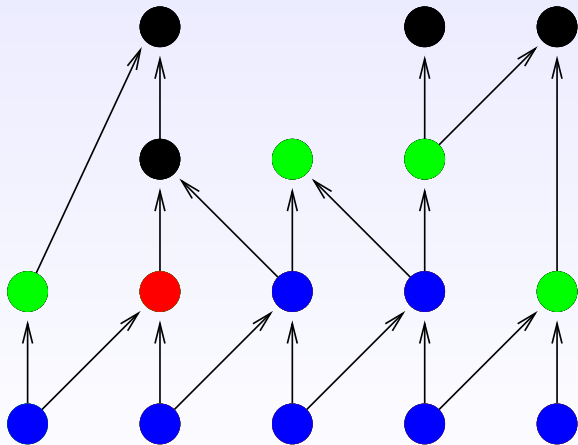
Internet-Based Computing : back to the example



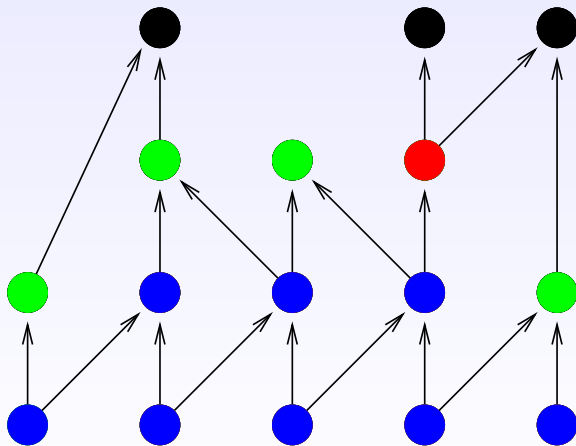
Internet-Based Computing : back to the example



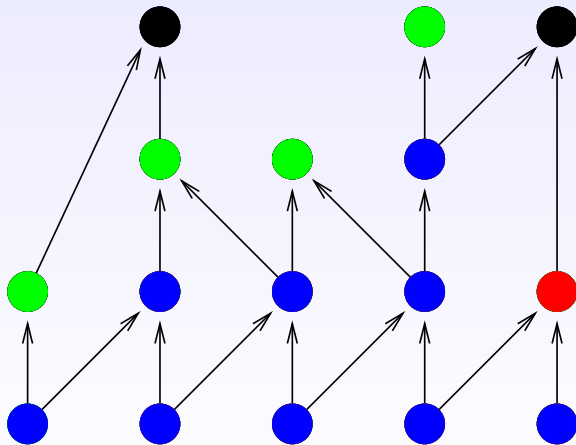
Internet-Based Computing : back to the example



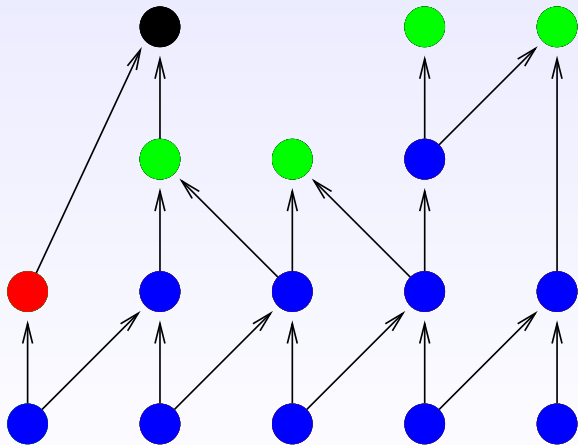
Internet-Based Computing : back to the example



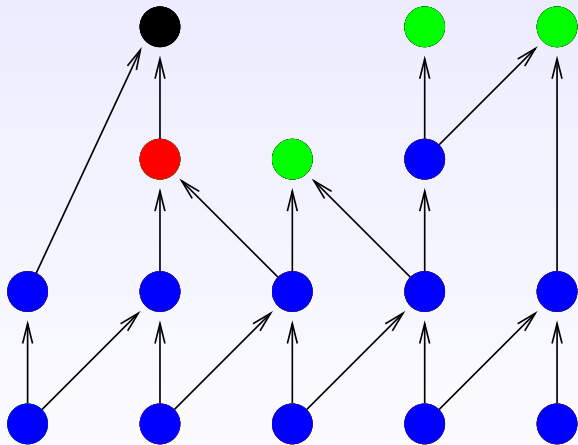
Internet-Based Computing : back to the example



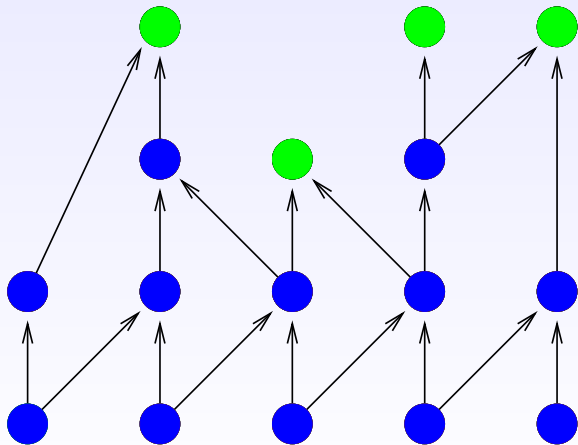
Internet-Based Computing : back to the example



Internet-Based Computing : back to the example



Internet-Based Computing : back to the example



Need of models

How do reality differ from the theoretical views?

- ▶ Model on uncertainties on the length of tasks?

Need of models

How do reality differ from the theoretical views?

- ▶ Model on uncertainties on the length of tasks?
- ▶ Model on the variations of the network capacities?

Need of models

How do reality differ from the theoretical views ?

- ▶ Model on uncertainties on the length of tasks ?
- ▶ Model on the variations of the network capacities ?
- ▶ Model on the variations of the processors capacities ?

Need of models

How do reality differ from the theoretical views?

- ▶ Model on uncertainties on the length of tasks?
- ▶ Model on the variations of the network capacities?
- ▶ Model on the variations of the processors capacities?

The model may have a big impact on the results...

What does the reality look like?

Modeling the availability of machines

Question :

- ▶ Being given a set of homogeneous machines for which we have a record of the availability ;
- ▶ Being given a percentage p and a confidence level c ;
- ▶ What is the longest time t during which we can assure with a confidence c that at least $p\%$ of the machines are available ?


Considered models

- ▶ Approximation by a Weibull function
density function $f_w(x) = \alpha\beta^{-\alpha}x^{\alpha-1}\epsilon^{\left(\frac{x}{\beta}\right)^\alpha}$
distribution $F_w(x) = 1 - \epsilon^{-\left(\frac{x}{\beta}\right)^\alpha}$
- ▶ Binomial law : let X_q be the q -th percentile of X distribution. Let (x_1, x_2, \dots, x_n) be an independent sample of X . The probability that at most k of the x_i are no larger than X_q is :

$$\sum_{j=0}^k \binom{n}{j} \cdot (1-q)^{n-j} \cdot q^j$$

- ▶ Re-sampling method : for a value q and a sample n , the asymptotic distribution of the sample of the q -th percentile is a normal law of mean value X_q and of variance :

$$\frac{1}{[f(X_q)]^2} \frac{q(1-q)}{n}$$

where $f(X_q)$ is the density function of the population. 

Results

Learning on 20 records (for each machine).

Data	# of machines	Weibull	Binomial	Re-sampling
CSIL	16	56.25	87,5	62.5
Condor	97	95.92	98.9	60.2
Long	83	57.95	94.3	53.4

TAB.: Percentage of machines' traces, for each set, for which the considered method gave an estimate between 95% and 100% of the experimental measure.

It had previously been shown that the third data set could not be modeled by a Poisson law or an exponential law.

Outline

- 1 The problem
- 2 Analyzing sensitivity and robustness
- 3 Dynamic load-balancing**
- 4 How useful is old information ?
- 5 Conclusion

Dynamic load-balancing

To cure (and no longer preventing) : the algorithm balance the load to take into account the des uncertainties and the dynamicity.

General scheme

- ▶ From time to time, do :

General scheme

- ▶ From time to time, do :
 - ▶ Compute a good solution using the observed parameters.

General scheme

- ▶ From time to time, do :
 - ▶ Compute a good solution using the observed parameters.
 - ▶ Evaluate the cost of balancing the load

General scheme

- ▶ From time to time, do :
 - ▶ Compute a good solution using the observed parameters.
 - ▶ Evaluate the cost of balancing the load
 - ▶ If the gain is larger than the cost : load-balance

General scheme

- ▶ From time to time, do :
 - Each invocation has a cost : the invocations should only take place at “useful” instants
 - ▶ Compute a good solution using the observed parameters.
 - ▶ Evaluate the cost of balancing the load
 - ▶ If the gain is larger than the cost : load-balance

General scheme

- ▶ From time to time, do :
 - Each invocation has a cost : the invocations should only take place at “useful” instants
 - ▶ Compute a good solution using the observed parameters.
How do we predict the future from the past ?
 - ▶ Evaluate the cost of balancing the load
 - ▶ If the gain is larger than the cost : load-balance

General scheme

- ▶ From time to time, do :
Each invocation has a cost : the invocations should only take place at “useful” instants
 - ▶ Compute a good solution using the observed parameters.
How do we predict the future from the past ?
 - ▶ Evaluate the cost of balancing the load
 - ▶ If the gain is larger than the cost : load-balance
If the objective is to minimize the running time, the comparison is obvious. How do we compare a time and some QoS ?

Network Weather Service

Distributed system which periodically monitors/records network and processor performance.

Also, allows to predict the future performance of the network and of the processors.

Network Weather Service

Distributed system which periodically monitors/records network and processor performance.

Also, allows to predict the future performance of the network and of the processors.

Does the past enable to predict the future ?

Outline

- 1 The problem
- 2 Analyzing sensitivity and robustness
- 3 Dynamic load-balancing
- 4 How useful is old information ?**
- 5 Conclusion

How useful is old information ?

The problem

- ▶ The values used when taking decisions have already “aged” .
- ▶ Is it a problem ? Should we take this ageing into account ?

Framework : the platform

- ▶ A set of n servers.
- ▶ Tasks arrive according to a Poisson law of throughput λn , $\lambda < 1$.
- ▶ Task execution time : exponential law of mean 1.
- ▶ Each server executes in FIFO order the tasks it receives.
- ▶ We look at the time each task spent in the system (=flow).

Framework : information

There is a *bulletin board* on which are displayed the loads of the different processors.

This information may be wrong or approximate.

We only deal with the case in which this information is *old*.

This is the only information available to the tasks : they cannot communicate between each other and have some coordinated behavior.

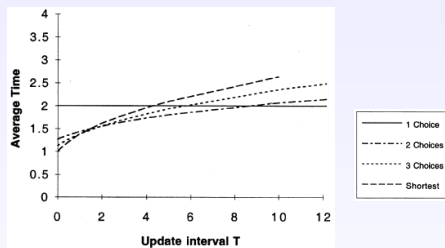
The obvious strategies

- ▶ Random and uniform choice of the server.
- ▶ Random and uniform choice of d servers, the task being sent on the least loaded of the d servers.
- ▶ Task sent on the least loaded server.

First model : periodic updates

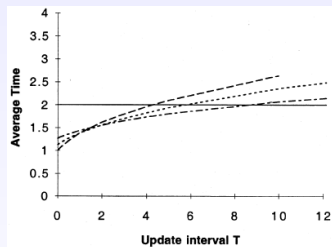
- ▶ Each T units of time the bulletin board is updated with correct information.

First model : results

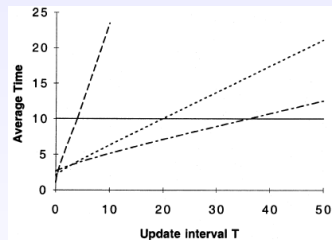


$n = 100$ and $\lambda = 0.5$

First model : results

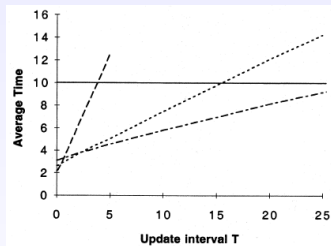


$n = 100$ and $\lambda = 0.5$

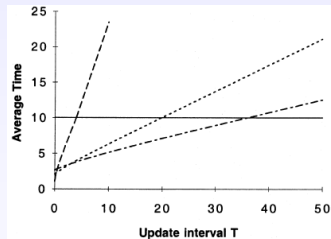


$n = 100$ and $\lambda = 0.9$

First model : results



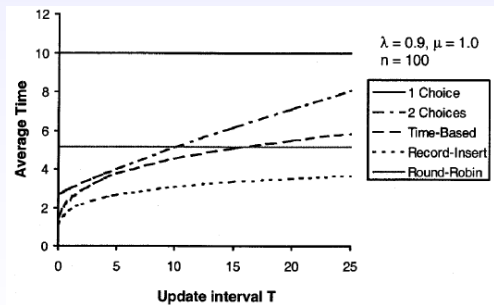
$n = 8$ and $\lambda = 0.9$



$n = 100$ and $\lambda = 0.9$

First model : more elaborated strategies

- ▶ *Time-based* : random choice among the servers which are supposed to be the least loaded.
- ▶ *Record-Insert* : centralized service in which each task updates the bulletin board by indicating on which server it is sent.

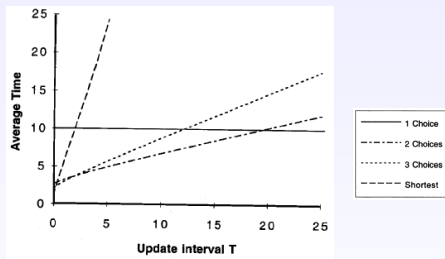


$n = 100$ and $\lambda = 0.9$

Second model : continuous updates

Model : continuous updates, but the information used is T units of time old.

$$n = 100 \text{ and } \lambda = 0.9$$



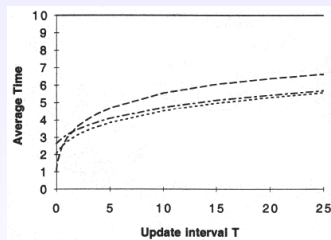
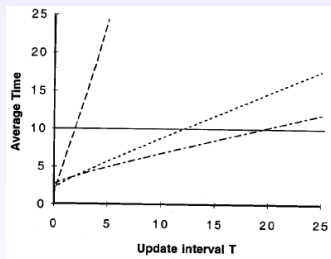
Age of information :

exactly T

Second model : continuous updates

Model : continuous updates, but the information used is T units of time old.

$$n = 100 \text{ and } \lambda = 0.9$$



Age of information :

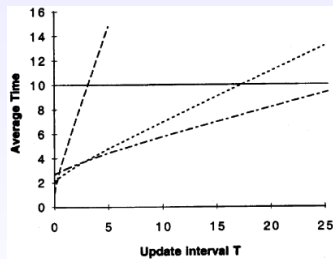
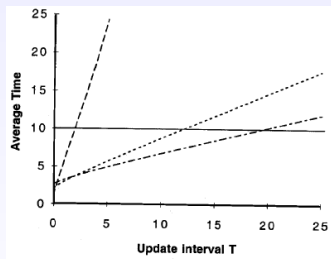
exactly T

exponential distribution, average T

Second model : continuous updates

Model : continuous updates, but the information used is T units of time old.

$$n = 100 \text{ and } \lambda = 0.9$$



Age of information :

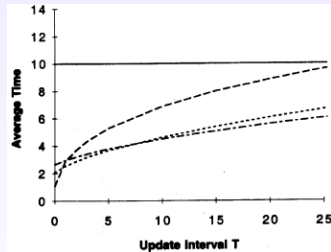
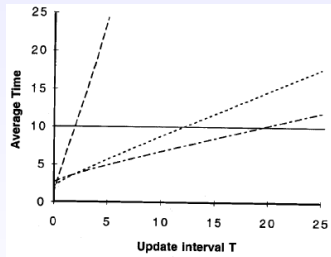
exactly T

uniform distribution on $[\frac{T}{2}; \frac{3T}{2}]$

Second model : continuous updates

Model : continuous updates, but the information used is T units of time old.

$$n = 100 \text{ and } \lambda = 0.9$$



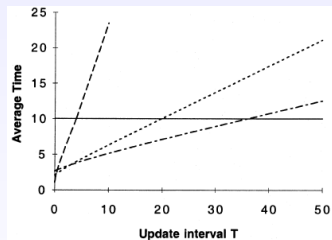
Age of information :

exactly T

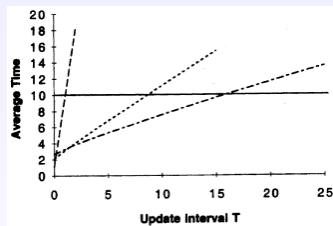
uniform distribution on $[0; 2T]$

Third model : de-synchronized updates

The different servers updates their information in a de-synchronized manner, each following an exponential law of average T .



Updates each T units of time.



De-synchronized updates.

$$n = 100 \text{ and } \lambda = 0.9$$

And if some were cheating?

With a probability p a task does not choose between two randomly determined servers, but takes the least loaded of all servers.

T	p	Avg. Time All Tasks	Avg. Time 2 Choices	Avg. Time Shortest	Variance All Tasks	Variance 2 Choices	Variance Shortest
1	0.00	3.23286					
1	0.01	3.21072	3.22093	2.19877	5.73117	5.74186	3.63718
1	0.05	3.17061	3.21389	2.34814	5.62948	5.67621	4.02956
1	0.10	3.14132	3.20978	2.52474	5.58554	5.65450	4.54205
1	0.25	3.20098	3.25693	3.03311	6.05849	5.94553	6.35980
5	0.00	4.94051					
5	0.01	4.95386	4.95677	4.66575	13.8029	13.8821	11.8131
5	0.05	5.05692	5.05668	5.06154	14.4591	14.5105	13.4837
5	0.10	5.21456	5.17956	5.52974	15.6083	15.6597	15.7552
5	0.25	6.06968	5.70758	7.15609	23.6380	22.0182	26.9240
10	0.00	6.74313					
10	0.01	6.80669	6.80588	6.88703	26.4946	26.5391	22.0827
10	0.05	7.00344	6.97692	7.50776	28.4836	28.6189	25.6448
10	0.10	7.36957	7.26152	8.34185	32.7326	32.7395	31.6201
10	0.25	8.91193	8.23577	10.9422	54.8097	52.0265	57.6721

Some memory always help

- ▶ Studied scenario : a task is allocated to the “best” of two randomly determined servers.
- ▶ New scenario : a task is allocated to the “best” of two servers, one being randomly chosen and the other one being the least-loaded one —after the previous task was allocated— of the two processors considered by the previous task.
- ▶ The problem : the memorization requires some communications and centralization.

Complete vs. incomplete information

Complete information

- ▶ Requires some centralization (or total replication);

Complete vs. incomplete information

Complete information

- ▶ Requires some centralization (or total replication);
- ▶ Communications of the most remote elements to the “center” ;

Complete vs. incomplete information

Complete information

- ▶ Requires some centralization (or total replication);
- ▶ Communications of the most remote elements to the “center” ;
- ▶ Obsolescence of the information.

Complete vs. incomplete information

Complete information

- ▶ Requires some centralization (or total replication) ;
- ▶ Communications of the most remote elements to the “center” ;
- ▶ Obsolescence of the information.

Decentralized schedulers

Complete vs. incomplete information

Complete information

- ▶ Requires some centralization (or total replication);
- ▶ Communications of the most remote elements to the “center” ;
- ▶ Obsolescence of the information.

Decentralized schedulers

- ▶ The local data are more up-to-date ;

Complete vs. incomplete information

Complete information

- ▶ Requires some centralization (or total replication) ;
- ▶ Communications of the most remote elements to the “center” ;
- ▶ Obsolescence of the information.

Decentralized schedulers

- ▶ The local data are more up-to-date ;
- ▶ A local optimization does not always lead to a global optimization...

Outline

- 1 The problem
- 2 Analyzing sensitivity and robustness
- 3 Dynamic load-balancing
- 4 How useful is old information ?
- 5 Conclusion**

Conclusion

- ▶ An obvious need to be able to cope with the dynamicity and the uncertainties.
- ▶ Crucial need to be able to model the dynamicity and the uncertainty.
- ▶ The static world is already complex enough!!!
- ▶ Where is the trade-off between the precision of the models and their usability?
- ▶ Trade-off between static and dynamic approaches?