



First experimental assessments of the adaptivity of the scheduling with AS4DR

Daniel Millot, Christian Parrot

► To cite this version:

Daniel Millot, Christian Parrot. First experimental assessments of the adaptivity of the scheduling with AS4DR. PDCAT '12: The Thirteenth International Conference on Parallel and Distributed Computing, Applications and Technologies, Dec 2012, Beijing, China. pp.487-492, 10.1109/PD-CAT.2012.75 . hal-00863179

HAL Id: hal-00863179

<https://hal.science/hal-00863179>

Submitted on 18 Sep 2013

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

First experimental assessments of the adaptivity of the scheduling with AS4DR

Daniel Millot and Christian Parrot
Telecom sudParis
Institut Mines-Telecom
France
{Daniel.Millot, Christian.Parrot}@mines-telecom.fr

Abstract—The AS4DR (Adaptive Scheduling for Distributed Resources) scheduling method experimented in this paper aims at maximizing the CPU use efficiency when executing divisible load applications on heterogeneous distributed memory platforms. AS4DR adapts the scheduling to: the unawareness of the total workload, both the unspecification and the variation over time of the execution parameters (available communication speed, available computing speed, etc.). This paper presents the first experimental assessments of the adaptivity of the scheduling with this method.

Keywords—parallel application; multi-round divisible load scheduling; heterogeneous platform; adaptive scheduling; unspecified distributed memory platform;

I. INTRODUCTION

This paper addresses the problem of maximizing the CPU utilization with useful work when scheduling a divisible load over a set of heterogeneous distributed resources, according to a master-worker model.

On the one hand, we suppose that the master receives a continuous input stream of data to be processed by the workers (such as some video stream, for instance). So the size of the total workload happens to be known only when the last item is acquired by the master, and as it is unknown when scheduling starts, the master must proceed to an iterated distribution as the workload flows in. Hence the resulting algorithm is necessarily multi-round, where we call "round" a sequence of consecutive actions leading the master to feed all the workers with chunks once and collect the corresponding results from the workers.

On the other hand, we suppose that the execution platform is a heterogeneous distributed memory platform whose communication and computation resources have inaccurately specified characteristics: available communication speeds, available computation speeds, latencies, etc., liable to vary over time and called execution parameters in the sequel.

In this paper, we assume that all the available workers are used, in a predefined order. From now on, for a given scheduling, we call CPU-efficiency the ratio of the time spent in useful computation over the corresponding elapsed time. Our goal is to maximize the CPU-efficiency. We

suppose that the duration of communicating or processing a chunk is affine according to the chunk size, and we assume that computation can overlap communication. Besides, we consider a 1-port bi-directional communication model, which allows a communication from master to worker to overlap a communication from worker to master; a risk of contention may then appear when workers compete to access the master.

The next section is a reminder of the operating principle of the AS4DR scheduler. Section III gives hints on related works whereas section IV presents experimental results. Finally we conclude.

II. PRESENTATION OF THE AS4DR METHOD

A. Operating principle

The ultimate goal of the AS4DR method is to automatically adapt the scheduling of a divisible load application to: the heterogeneity of the workers, the unawareness of the total workload, both the unspecification of the execution parameters and their evolution over time. Let $\alpha_{w,i}$ be the size of the chunk sent to a worker w for round i . Let τ and $\sigma_{w,i}$ be respectively the wanted and the estimated time durations between the start of the sending of a chunk of size $\alpha_{w,i}$ and the end of the reception of the corresponding result by the master. The basic idea of the AS4DR multi-round method is to adapt $\alpha_{w,i}$ according to:

$$\alpha_{w,i} := \alpha_{w,i-1} \frac{\tau}{\sigma_{w,i-1}} \quad \text{for } i > 1. \quad (1)$$

A special feature of AS4DR is that it splits each chunk it has to deliver to a worker for a round into two subchunks that it delivers in a row to the worker. So, sending subchunks of arbitrarily chosen sizes $\dot{\alpha}_{w,1}$ and $\ddot{\alpha}_{w,1}$ to each worker w for the first round, the AS4DR scheduler then sends to worker w , for each round i , two subchunks \dot{s} and \ddot{s} of respective sizes $\dot{\alpha}_{w,i}$ and $\ddot{\alpha}_{w,i}$, such that

$$\dot{\alpha}_{w,i} + \ddot{\alpha}_{w,i} = \alpha_{w,i}. \quad (2)$$

Dividing the chunks in two parts allows the computation to overlap the communications between a worker and the master as can be seen in Figure 3. Let us suppose that the

ratio between $\dot{\alpha}_{w,i}$ and $\alpha_{w,i}$ is constant, and let us denote θ_w this ratio:

$$\theta_w \equiv \frac{\dot{\alpha}_{w,i}}{\alpha_{w,i}}. \quad (3)$$

Figures 1 and 2 give the AS4DR scheduling algorithm for a M workers platform, and Figure 3 shows how computation overlaps communication with this method. As can be seen

- With CIP set τ and $(\alpha_{w,1}, \theta_{w,1}) \dots \dots \dots$ (Figure 6)
 $0 \leq w \leq M-1$
- Set $(\dot{\alpha}_{w,1}, \ddot{\alpha}_{w,1})_{0 \leq w \leq M-1} \dots \dots \dots$ (3), (2)
- After an appropriate delay, post \dot{s} and \ddot{s} data to each worker

while (the last data item has not been acquired) **do**

- Get a \dot{s} result from some worker w
- Compute the size of the next \dot{s} and \ddot{s} for worker $w \dots \dots \dots$ (4), (1), (3), (2)
- Post \dot{s} and \ddot{s} data to worker w
- Get previous \ddot{s} result from worker w

end while

Figure 1. AS4DR scheduling: master

while (the last subchunk has not been posted) **do**

- Get a subchunk data from master
- Process the subchunk data

if (\dot{s}) **then**

- Post \dot{s} and previous \ddot{s} results to master

end if

end while

Figure 2. AS4DR scheduling: worker

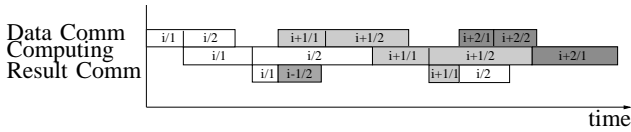


Figure 3. Overlapping between communication and computation

in Figure 3, worker round i for worker w is composed of three phases:

- transmission of the data from master to worker, lasting $\dot{D}_{w,i}$ and $\ddot{D}_{w,i}$ for subchunks \dot{s} and \ddot{s} respectively,
- worker computation on the received data, lasting $\dot{C}_{w,i}$ and $\ddot{C}_{w,i}$ for subchunks \dot{s} and \ddot{s} respectively,
- transmission of the computation result from worker to master, lasting $\dot{R}_{w,i}$ and $\ddot{R}_{w,i}$ for subchunks \dot{s} and \ddot{s} respectively.

It is worth noticing in Figure 3 that the result corresponding to the \ddot{s} subchunk of some round is returned to the master just after the result corresponding to the \dot{s} subchunk of the

next round has itself been returned. As the scheduler does not make use of the return of \ddot{s} results in any way, AS4DR delays this return in order to make all actions have the same period; this helps contentions avoidance. Let us denote F_w the available computation speed (relative to the processing of one workload unit) of worker w . Likewise, B_w^D (resp. B_w^R) is the available communication speed (relative to one workload unit) of the link from the master to worker w (resp. from worker w to the master). Finally b_w^D , b_w^R and f_w are the respective latencies for a transfer of data from the master to worker w , for a transfer of result from worker w to the master and for a computation on worker w . We define

$$\sigma_{w,i} \equiv \frac{\dot{C}_{w,i} - f_w}{\theta_w} + 2f_w. \quad (4)$$

When the communications between master and workers are contention-free,

$$\lim_{i \rightarrow +\infty} \sigma_{w,i} = \tau.$$

B. Prevention of idleness

The AS4DR method could experience either of the workers idlenesses illustrated in figures 4: inter-round idleness and intra-round idleness. Let us respectively define:

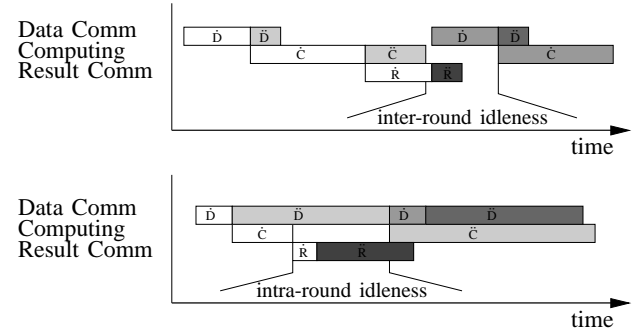


Figure 4. Example of inter-round and intra-round idlenesses

$$\theta_{w,i+1}^{\min} \equiv \frac{D_{w,i+1} - (b_w^D + f_w)}{D_{w,i+1} + C_{w,i+1} - 2(b_w^D + f_w)}, \quad (5)$$

$$\theta_{w,i+1}^{\max} \equiv \frac{C_{w,i} - (b_w^D + f_w + b_w^R)}{D_{w,i+1} + C_{w,i} + R_{w,i} - 2(b_w^D + f_w + b_w^R)}. \quad (6)$$

Let us suppose that the communications between master and workers are contention-free. Then the AS4DR method prevents idleness, if and only if, for each worker w

$$\theta_{w,i+1}^{\max} \geq \theta_w \geq \theta_{w,i+1}^{\min}, \quad \forall i. \quad (7)$$

C. Prevention of contentions

In order to take advantage of this result, we need to address the problem of contention avoidance. To make the instant each worker accesses to the master far enough from the instants the others access too, time delays d_w are

introduced before posting the very first subchunk to each worker w . Figure 5 illustrates the model of asymptotic τ -

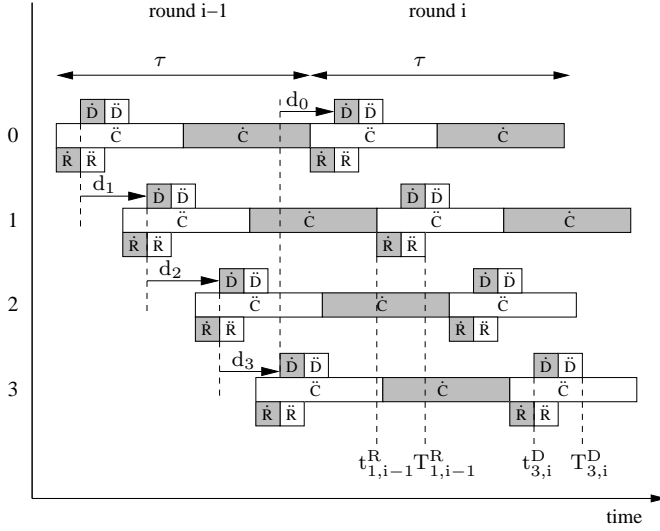


Figure 5. Contention-free asymptotic schedule

periodic (thus round-robin) scheduling we are looking for, in the case of a four workers platform. The delay d_w , to be computed by the CIP algorithm is defined (modulo M) as follows:

$$d_w \equiv (1 + \lambda_w) \max \left(\overline{D_{w-1}} + \overline{D_{w-1}}, \overline{R_{w-1}} + \overline{R_w} \right); \quad (8)$$

where: λ_w stands for a positive constant factor, $\overline{B_w^D}$ (resp. $\overline{F_w}$, $\overline{B_w^R}$, $\overline{b_w^D}$, $\overline{f_w}$ and $\overline{b_w^R}$) denotes an estimate of B_w^D (resp. F_w , B_w^R , b_w^D , f_w and b_w^R) and $\overline{D_w}$ (resp. $\overline{D_w}$, $\overline{R_w}$ and $\overline{R_w}$) is an estimate of $\dot{D}_{w,1}$ (resp. $\dot{D}_{w,1}$, $\dot{R}_{w,1}$ and $\dot{R}_{w,1}$).

$$\overline{D_w} := \theta_w \alpha_w \frac{1}{\overline{B_w^D}} + \overline{b_w^D}, \quad \overline{D_w} := (1 - \theta_w) \alpha_w \frac{1}{\overline{B_w^D}} + \overline{b_w^D}, \quad (9)$$

$$\overline{R_w} := \theta_w \alpha_w \frac{1}{\overline{B_w^R}} + \overline{b_w^R}, \quad \overline{R_w} := (1 - \theta_w) \alpha_w \frac{1}{\overline{B_w^R}} + \overline{b_w^R}. \quad (10)$$

$$\alpha_w := \frac{\tau - 2\overline{f_w}}{\frac{1}{\overline{F_w}}} \quad (11)$$

$$\theta_w := \phi_w \theta_w^{\max} + (1 - \phi_w) \theta_w^{\min}, \quad (12)$$

$$\theta_w^{\min} := \frac{\frac{\alpha_w}{\overline{B_w^D}} + \overline{b_w^D} - \overline{f_w}}{\alpha_w \left(\frac{1}{\overline{F_w}} + \frac{1}{\overline{B_w^D}} \right)}, \quad \theta_w^{\max} := \frac{\frac{\alpha_w}{\overline{B_w^R}} + \overline{b_w^R} - \overline{f_w}}{\alpha_w \left(\frac{1}{\overline{F_w}} + \frac{1}{\overline{B_w^D}} + \frac{1}{\overline{B_w^R}} \right)}.$$

$$\phi_w = 0.5. \quad (13)$$

So, before launching AS4DR, a preliminary step: CIP (for Contentions and Idleness Prevention), determines τ , $(\alpha_w)_{0 \leq w \leq M-1}$, $(\theta_w)_{0 \leq w \leq M-1}$ and $(d_w)_{0 \leq w \leq M-1}$ as shown by Figure 6.

• Estimation of $(B_w^D, F_w, B_w^R, b_w^D, f_w, b_w^R)_{0 \leq w \leq M-1}$
 • Set $(\phi_w)_{0 \leq w \leq M-1}$ (13)
 • Set $(\lambda_w)_{0 \leq w \leq M-1}$ (15)
 • $\tau := 2 \max_{0 \leq s \leq M-1} f_w$
repeat
 • $\tau := \tau + \hat{\tau}$
 • Set $(\alpha_w)_{0 \leq w \leq M-1}$ (11)
 • Set $(\theta_w)_{0 \leq w \leq M-1}$ (12)
 • Set $(\overline{D_w}, \overline{D_w}, \overline{R_w}, \overline{R_w})_{0 \leq w \leq M-1}$ (9),(10)
 • Set $(d_w)_{0 \leq w \leq M-1}$ (8)
until $(\tau \geq \sum_{w=0}^{M-1} d_w)$

Figure 6. CIP algorithm

Besides, the time intervals d_w should allow all the workers to be served during the first round, i.e. within a τ period. Thus τ must verify:

$$\tau \geq \sum_{w=0}^{M-1} d_w. \quad (14)$$

So, starting from an initial value of τ , the CIP algorithm enters an iterative process which increments τ with an arbitrarily fixed value $\hat{\tau}$, then computes $(\alpha_w)_{0 \leq w \leq M-1}$ and $(d_w)_{0 \leq w \leq M-1}$ successively and loops until (14) holds.

Once CIP is processed, proper AS4DR scheduling starts with a first round which sets the initial time-lags between successive round beginnings, according to the values $(d_w)_{0 \leq w \leq M-1}$ previously computed by CIP.

Taking into account the previous results about contentions and idleness avoidance, the CIP algorithm can avoid contentions and idleness during the AS4DR scheduling first round. For the next rounds, thanks to assignment (1), the AS4DR method helps maintain the duration of each round close to the reference value τ .

Let us define the value Λ_w for each worker w as

$$\Lambda_w \equiv \frac{1}{M} \frac{1}{\overline{F_w} \max \left(K_w + \frac{1}{\overline{B_w^D}}, \frac{1}{\overline{B_w^R}} \right)} - 1;$$

$$\text{where } K_w \equiv \frac{1}{\overline{B_w^R}} \left(\phi_w \frac{1}{1 + \overline{F_w} \left(\frac{1}{\overline{B_w^D}} + \frac{1}{\overline{B_w^R}} \right)} + (1 - \phi_w) \frac{1}{1 + \frac{\overline{b_w^D}}{\overline{F_w}}} \right).$$

Let us assume that $(\lambda_w)_{0 \leq w \leq M-1}$ verifies:

$$0 \leq \lambda_w \leq \min(\Lambda_w, \Lambda_{w-1}) \quad (\text{modulo } M), \quad (15)$$

then the CIP preliminary step provides a τ value and sets of values $(\theta_w)_{0 \leq w \leq M-1}$, $(\alpha_w)_{0 \leq w \leq M-1}$ and $(d_w)_{0 \leq w \leq M-1}$

that allow the AS4DR scheduling to start with neither contention nor idleness. Let us recall that the value of λ_w characterizes the importance of both the unspecification of the execution parameters and the variation of their magnitude over time, for any worker w . For a given set of execution parameters values, increasing the number M of workers will ultimately make the relation (15) false.

III. RELATED WORKS

Details about related works can be found in [1] and the proof of the results presented in the previous section has been established in [2]. To the best of our knowledge, no scheduler aims at maximizing the efficiency of the use of the CPUs in a similar context.

IV. EXPERIMENTAL ASSESSMENT OF THE ADAPTIVITY

In order to experimentally assess the ability of the AS4DR method to cope with inaccurately estimated execution parameters and with variation of their value over time, simulations have been conducted, with the SimGrid framework [3]. Let us consider a star-shaped platform (Figure 7 with $M=1000$) and 10 sets $(S_k)_{0 \leq k \leq 9}$ of values for the execution parameters of the worker nodes, given in Table I. Each of these sets is randomly allocated to 100 workers among the 1000 workers which constitute the platform. In order

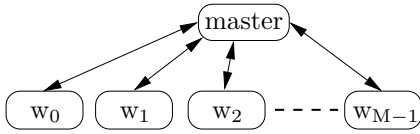


Figure 7. Star-shaped platform

	computation		communication master $\rightarrow w_i$		communication master $\leftarrow w_i$		number of workers
	speed	latency	speed	latency	speed	latency	
S_0	1.0e+4	1.0e-4	1.0e+8	1.0e-4	1.0e+8	1.0e-4	100
S_1	0.9e+4	1.5e-4	0.9e+8	1.0e-4	0.9e+8	1.0e-4	100
S_2	0.8e+4	1.4e-4	0.8e+8	1.0e-4	0.8e+8	1.0e-4	100
S_3	0.7e+4	1.3e-4	0.7e+8	1.0e-4	0.7e+8	1.0e-4	100
S_4	0.6e+4	1.2e-4	0.6e+8	1.0e-4	0.6e+8	1.0e-4	100
S_5	0.5e+4	0.8e-4	0.5e+8	1.0e-4	0.5e+8	1.0e-4	100
S_6	0.4e+4	0.9e-4	0.4e+8	1.0e-4	0.4e+8	1.0e-4	100
S_7	0.3e+4	0.8e-4	0.3e+8	1.0e-4	0.3e+8	1.0e-4	100
S_8	0.2e+4	0.4e-4	0.2e+8	1.0e-4	0.2e+8	1.0e-4	100
S_9	0.1e+4	0.5e-4	0.1e+8	1.0e-4	0.1e+8	1.0e-4	100

Table I
REFERENCE VALUES OF THE EXECUTION PARAMETERS

to assess the relevance of AS4DR adapting the workload (of each worker at each round), we compare this method to a scheduler called “Baseline”. The Baseline method is identical to AS4DR except that, on the one hand, the wanted period τ is not computed but just set, and that, on the other hand, the workload is not adapted at each round; in

other words, neither CIP algorithm nor assignment (1) are run. With τ and estimates of the execution parameters, it is possible, due to (11), to compute α_w an initial workload for each worker w . For our simulations, we set the estimate of each execution parameter for each worker as the time-average of the value of this execution parameter for this worker, during the simulation. In order to be able to compare the two methods efficiently, the value of τ for the Baseline scheduler is set to the one obtained by the CIP method when running AS4DR.

For the comparison of the simulation results, let us define CPU_{eff} the CPU-efficiency:

$$\text{CPU}_{\text{eff}} \equiv 1 - \frac{\text{CPU idleness} + \text{CPU latencies}}{\text{elapsed time}}.$$

When the execution parameters are exactly known and steady, both schedulers make the workers process data without idleness; except time spent in latencies. So, in the sequel, we will successively compare the schedulers when the parameters are either poorly estimated or time-varying.

A. Poor estimates in a steady context

In order to assess the effect of estimates’ inaccuracy of the execution parameters on both schedulings, the effective workload initially allocated to each worker w is set by penalizing the value $(\alpha_w)_{\text{ref}}$ computed using assignment (11) and the reference values contained in Table I:

$$\alpha_{w,1} := (1 \pm \iota_k) (\alpha_w)_{\text{ref}}; \quad (16)$$

where ι_k is a strictly positive real number which values are given in Table II and where the operator \pm means that the operation performed is randomly chosen between either addition or subtraction. The values of ι_k characterize esti-

ι_0	ι_1	ι_2	ι_3	ι_4	ι_5
0.0	0.18	0.36	0.54	0.72	0.9

Table II
ESTIMATES’ INACCURACY PARAMETER

mates’ inaccuracy of the execution parameters; inaccuracy is minimum with ι_0 , whereas it is maximum with ι_5 . It is supposed to be the same for all the workers.

Figure 8 shows the measured CPU-efficiency of the whole platform for each scheduler, as a function of estimates’ inaccuracy. For this simulation, which lasted 2000 seconds, the value of τ computed by CIP equals 3 seconds. When the execution parameters are exactly known (ι_0), both schedulers offered the same CPU-efficiency: 99.99%. The computation latencies are responsible for the gap with the theoretical maximum CPU-efficiency: 100%. As expected, CPU-efficiency decreases for both methods when estimates’ inaccuracy increases. This decrease is considerably faster for the Baseline scheduler than for AS4DR. For instance, when the inaccuracy is maximum (ι_5), the CPU-efficiency

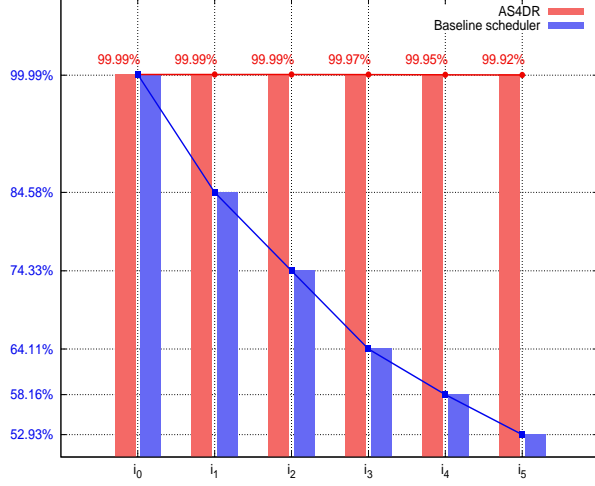


Figure 8. CPU-efficiency of the platform as a function of estimates' inaccuracy

with AS4DR is 2.13 times higher than with the Baseline scheduler.

B. Time-varying context

This subsection assesses to what extent one can put this ability to good use in adapting to the variations of the values of execution parameters over time, given that it has already been proved that the outbreak of these variations will not make the AS4DR scheduling unstable.

In this subsection, the reference values for the execution parameters, which are contained in Table I, are still randomly allocated to the workers at the initial instant. But, from this initial instant these values are likely to vary over time, according to the 10 profiles \$(P_k)_{0 \leq k \leq 9}\$ of variation shown by Figure 9. Whatever the profile \$P_k\$ being allocated to a worker \$w\$ and whatever an execution parameter, the higher value of \$P_k\$ equals the reference value (in Table I) allocated to the worker \$w\$, for the execution parameter under consideration; the lower value of \$P_k\$ is computed as a perturbation of the reference value:

$$B_w^D := (1 - \delta_k) (B_w^D)_{\text{ref}}, B_w^R := (1 - \delta_k) (B_w^R)_{\text{ref}}, \\ F_w := (1 - \delta_k) (F_w)_{\text{ref}};$$

where \$\delta_k\$ is a strictly positive number, the values of which are given in Table III. In this context, the coefficient \$\delta_k\$ characterizes the variations of the execution parameters and is called “dynamicity” in the sequel; with \$\delta_4\$ the amplitude of the variation of the execution parameters is maximum, whereas it is minimum with \$\delta_0\$. For each simulation the dynamicity is the same for all the workers. The same profile is used to simultaneously perturb the speeds of computation and communication. The profiles \$(P_k)_{0 \leq k \leq 9}\$ are randomly allocated to the 1000 workers as well. Figure 10 shows the

δ_0	δ_1	δ_2	δ_3	δ_4
0.0	0.2	0.4	0.6	0.8

Table III
DYNAMICITY PARAMETER

measured CPU-efficiency of the whole platform for each scheduler, as a function of dynamicity. For this simulation,

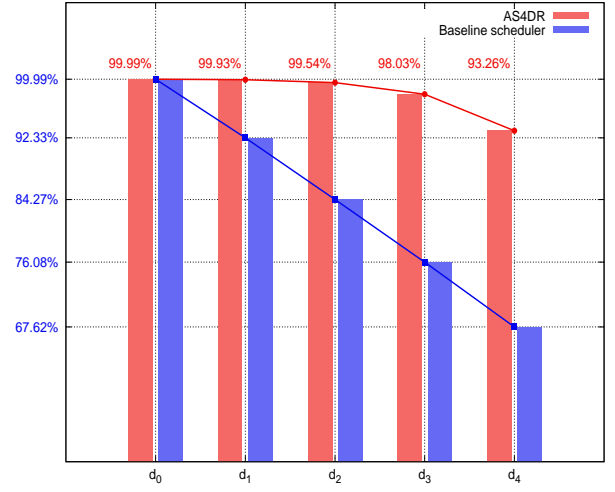


Figure 10. CPU-efficiency of the platform as a function of dynamicity

which lasted 2000 seconds, the value of \$\tau\$ computed by CIP equals 3 seconds. When the execution parameters are steady (\$\delta_0\$), both schedulers offered the same CPU-efficiency: 99.99%. As expected, the decrease of the CPU-efficiency, when the dynamicity of the execution parameters increases, is significantly faster for the scheduler Baseline than for AS4DR. For instance, when the dynamicity equals \$\delta_4\$, the CPU-efficiency with AS4DR is 1.38 times higher than with the Baseline scheduler.

V. CONCLUSION

The AS4DR method experimentally assessed in this paper succeeds in maximizing the CPU-efficiency when scheduling a divisible load of unknown total size on distributed resources with inaccurately specified or time-varying characteristics. Despite the fact that a bidirectionnal 1-port communication model is prone to contention, AS4DR can avoid the idleness of the CPU due to contentions, thanks both to the asymptotic periodicity it installs (for both data and results) and to its preliminary step CIP. AS4DR can only be efficient when the unspecification of the execution parameters and their variation over time are small enough; especially AS4DR cannot manage the complete stop of a worker. The pseudo-periodicity of the rounds limits the workload discrepancy when the last data is acquired by the master; when the cleanup phase begins. Thus AS4DR limits

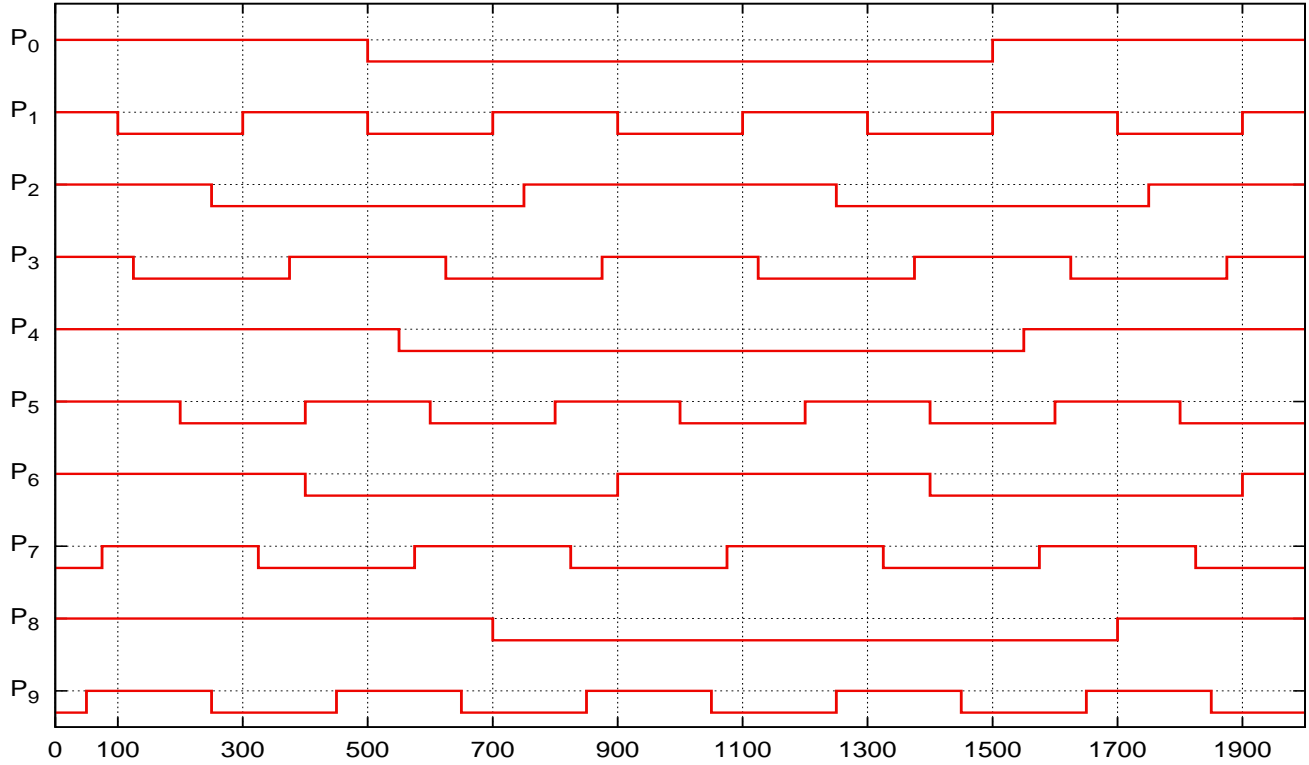


Figure 9. Perturbed execution parameter as a function of time (in seconds)

the cost of the communications for reallocating the extra load of the overloaded workers among the others ones; to make workers ending their processing simultaneously. The experimental results of this paper confirm that the adaptation to either poor estimates of the characteristics of the platform, or to their time-variation, are similar problems. Compared to using pure hardware performance figures, such as measured bandwidth or CPU frequency to adapt the workload at each round, AS4DR has the extra advantage of taking into account characteristics of the software such as algorithmic complexity.

Up to now we have considered that the whole set of resources is used. Of course, for a given set of execution parameters values, using all the available resources will ultimately become impossible with an evermore increasing amount of resources. The CIP algorithm could help to select

relevant subsets of resources. It is the aim of a future work.

REFERENCES

- [1] D. Millot and C. Parrot, "Scheduling on unspecified heterogeneous distributed resources," in *Proceeding of the 25th International Symposium on Parallel and Distributed Processing Workshops (IPDPSW'11)*, vol. 1, no. 1, IEEE Computing Society Press, May 2011, pp. 45–56.
- [2] D. Millot and C. Parrot, "Fundamental results on the AS4DR scheduler," TELECOM sudParis, Évry(France), Tech. Rep. RR-11005-INF, December 2011.
- [3] H. Casanova, A. Legrand, and L. Marchal, "Scheduling distributed applications: the simgrid simulation framework," in *Proceedings of the 3th International Symposium on Cluster Computing and the Grid (CCGrid03)*, IEEE Computing Society Press, 2003.