



HAL
open science

Some tests of adaptivity for the AS4DR scheduler

Daniel Millot, Christian Parrot

► **To cite this version:**

Daniel Millot, Christian Parrot. Some tests of adaptivity for the AS4DR scheduler. ICPPW '12: The 41st International Conference on Parallel Processing Pittsburg, PA 10-13 september 2012, Sep 2012, Pittsburgh, États-Unis. pp.323-331, 10.1109/ICPPW.2012.48 . hal-00822659

HAL Id: hal-00822659

<https://hal.science/hal-00822659v1>

Submitted on 15 May 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.

Some tests of adaptivity for the AS4DR scheduler

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

Abstract—This paper presents some tests of adaptivity for the AS4DR (Adaptive Scheduling for Distributed Resources) scheduler. The objective of AS4DR is to maximize the CPU use efficiency when executing divisible load applications on heterogeneous distributed memory platforms. Furthermore, this scheduler can operate when the total workload is unknown and when the execution parameters (available communication speed, available computing speed, etc.) are unspecified or may vary through time.

The paper analyzes results obtained when simulating AS4DR scheduling on platforms characterized by such type of execution parameters.

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

I. INTRODUCTION

Let us consider the problem of maximizing the CPU utilization with useful work when scheduling a divisible load over a set of distributed resources, according to a master-worker model. The master receives a continuous input stream of data to be processed by the workers. 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. The considered execution platform is a 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. 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. As numerous schedulers do, AS4DR attempts to install periodic accesses to the master in order to ease the contention avoidance inherent to the master-workers model. Finally, we assume that master/workers have an unlimited buffering capability.

The rest of the paper is organized as follows. The next section deals with related works. Section III reminds the fundamental results on the AS4DR scheduler whereas section IV presents experimental results on a platform with 1000 workers; they assess the adaptivity of the scheduling obtained with this method. Finally, we conclude.

II. RELATED WORKS

A review of related works can be found in [1]. To the best of our knowledge, no scheduler aims at maximizing the efficiency of the use of the CPUs in a similar context.

III. THE AS4DR METHOD

The proof of the results presented in this section has been established in [2].

A. AS4DR method principle

Let's first recall briefly the context in which the AS4DR method has been proposed and studied. Its ultimate goal is to automatically adapt the scheduling of a divisible load application to the evolution of the execution parameters of the platform 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 periodicity and the estimated time duration 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 (1):

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

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

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

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 7)
- Set $(\hat{\alpha}_{w,1}, \hat{\alpha}_{w,1})_{0 \leq w \leq M-1} \dots \dots \dots$ (3), (2)
- After an appropriate delay, post \hat{s} and \hat{s} data to each worker

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

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

- Post \hat{s} and previous \hat{s} results to master

end if

end while

Figure 2. AS4DR scheduling : worker

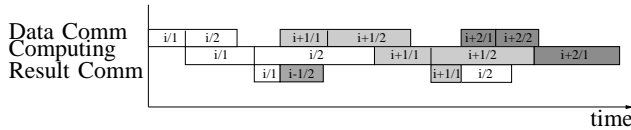


Figure 3. Overlapping between communication and computation

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

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

It is worth noticing in Figure 3 that the result corresponding to the \hat{s} subchunk of some round is returned to the master just after the result corresponding to the \hat{s} subchunk of the next round has itself been returned. As the scheduler does not make use of the return of \hat{s} results in any way, AS4DR delays this return in order to make all actions have the same period; this helps contentions avoidance.

As we suppose that the communication and computation costs are both roundwise affine in the size of the corresponding chunk, we have

$$\begin{aligned} \hat{D}_{w,i} &= \theta_w \frac{\alpha_{w,i}}{B_w^D} + b_w^D, & \hat{D}_{w,i} &= (1 - \theta_w) \frac{\alpha_{w,i}}{B_w^D} + b_w^D, \\ \hat{C}_{w,i} &= \theta_w \frac{\alpha_{w,i}}{F_w} + f_w, & \hat{C}_{w,i} &= (1 - \theta_w) \frac{\alpha_{w,i}}{F_w} + f_w, \\ \hat{R}_{w,i} &= \theta_w \frac{\alpha_{w,i}}{B_w^R} + b_w^R, & \hat{R}_{w,i} &= (1 - \theta_w) \frac{\alpha_{w,i}}{B_w^R} + b_w^R. \end{aligned}$$

Where F_w is the available computation speed (relative to the processing of one workload unit) of worker w . It is worth noticing that this computation speed depends both on the platform and the application, by means of the clock frequency of worker w and the algorithmic complexity of the application respectively. 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

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

The value of $\sigma_{w,i}$ is an extrapolation of the measured value of $\hat{C}_{w,i}$, in order to estimate the computation cost $C_{w,i}$ for the whole chunk. So,

$$\sigma_{w,i} = \alpha_{w,i} \frac{1}{F_w} + 2f_w. \quad (5)$$

The AS4DR method computes the next chunk size using equation (1). From (5) and (1), we get

$$\sigma_{w,i} = 2f_w + \tau - \frac{2f_w \tau}{\sigma_{w,i-1}}. \quad (6)$$

This sequence can be defined if and only if $\sigma_{w,i}$ takes non zero values for all i . So, in the sequel, we will suppose that $\alpha_{w,i}$ and f_w do not simultaneously equal zero.

Let us suppose next that the communications between master and workers are contention-free. The sequence $(\sigma_{w,i})_i$, built according to the AS4DR method, linearly converges. More, the scheduling built according to the AS4DR method is asymptotically stable.

$$\begin{aligned} \text{if } \tau &\geq 2f_w, \text{ then } & \lim_{i \rightarrow +\infty} \sigma_{w,i} &= \tau, \\ \text{if } \tau &< 2f_w, \text{ then } & \lim_{i \rightarrow +\infty} \sigma_{w,i} &= 2f_w. \end{aligned}$$

Practically, the duration τ , typically targeted for a round, is much larger than that of the computation latency f_w . Therefore, only the very first case: $\tau \geq 2f_w$, will be considered in the sequel. The dates when the master posts data, as well as those when it gets results, are asymptotically periodic, with the same period τ for all the workers.

B. Prevention of idleness and contentions

The AS4DR method could experience either of the workers idlenesses illustrated in Figures 4 and 5 : inter-round idleness and intra-round idleness. When the communications

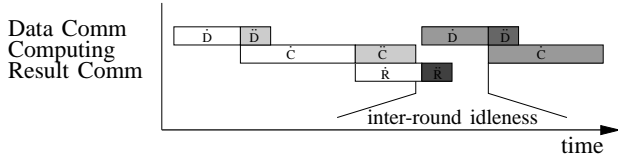


Figure 4. Example of inter-round idleness

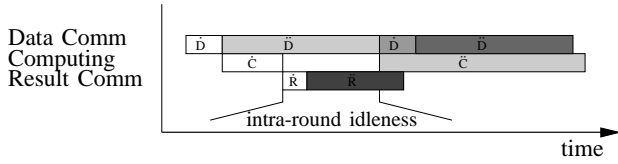


Figure 5. Example of intra-round idleness

between master and workers are contention-free, the value of the ratio θ_w can be set (e.g. according to (14)) such that the AS4DR method prevents idleness. So, in order to take advantage of both this result and the one of the previous section, the problem of contentions avoidance needs to be adressed.

To make the instant each worker accesses 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 . Introducing such laxity in the scheduling makes the AS4DR method more compliant with the errors on the estimate of the execution parameters. Figure 6 illustrates the model of asymptotic τ -periodic (thus round-robin) scheduling we are looking for, in the case of a four workers platform.

The AS4DR method prevents contentions, if and only if, $\forall i$ and for each worker w (modulo M),

$$d_w \geq \max(\dot{D}_{w-1,i} + \ddot{D}_{w-1,i}, \ddot{R}_{w-1,i-1} + \dot{R}_{w,i}). \quad (7)$$

The larger the d_w time intervals, the lower the risk of contentions in case of inaccurate execution parameters estimation, or in case of variation of these parameters over time. Before the launch of the AS4DR scheduler (see Figure 1) a step, called CIP (for Contentions and Idleness Prevention), determines the value for τ and the three sets of

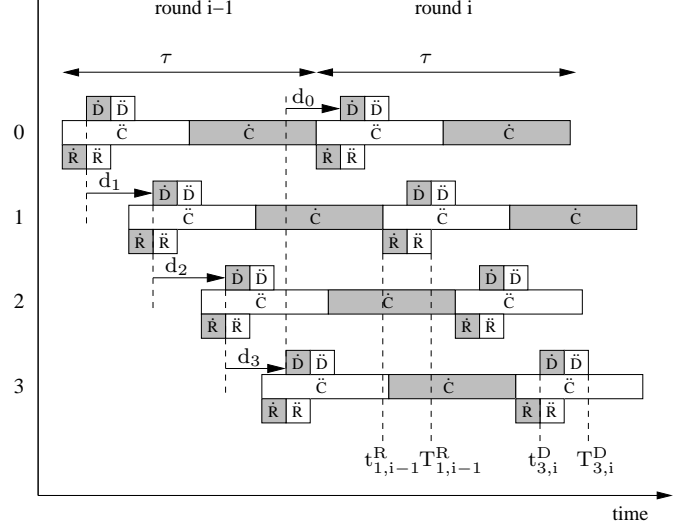


Figure 6. Contention-free asymptotic scheduling

values $(\alpha_{w,1})_{0 \leq w \leq M-1}$, $(\theta_w)_{0 \leq w \leq M-1}$ and $(d_w)_{0 \leq w \leq M-1}$, as shown by Figure 7. Let us suppose that we have at our disposal an estimate of B_w^D (resp. F_w , B_w^R , b_w^D , f_w and b_w^R) which is denoted \overline{B}_w^D (resp. \overline{F}_w , \overline{B}_w^R , \overline{b}_w^D , \overline{f}_w and \overline{b}_w^R). The

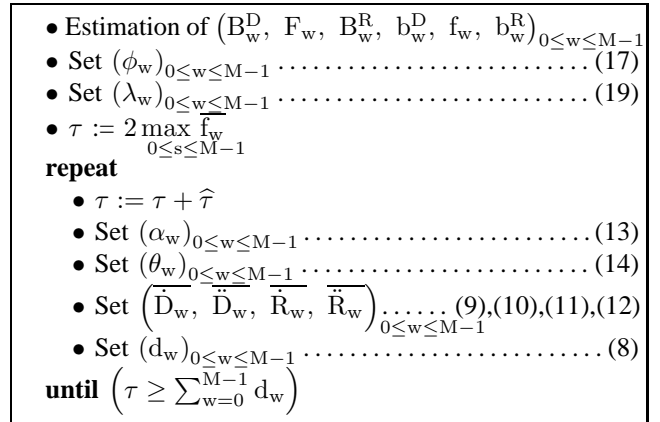


Figure 7. CIP algorithm

delay d_w , to be computed by the CIP algorithm is defined as follows:

$$d_w \equiv (1 + \lambda_w) \max \left(\overline{\dot{D}}_{w-1} + \overline{\dot{D}}_{w-1}, \overline{\dot{R}}_{w-1} + \overline{\dot{R}}_w \right) \text{ (modulo } M); \quad (8)$$

where λ_w is a positive constant factor and $\overline{\dot{D}}_w$ (resp. $\overline{\dot{D}}_w$, \overline{R}_w and $\overline{\dot{R}}_w$) are estimates of $\dot{D}_{w,1}$ (resp. $\dot{D}_{w,1}$, $\dot{R}_{w,1}$ and

$\ddot{R}_{w,1}$) obtained thanks to

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

$$\overline{\overline{D}}_w := (1 - \theta_w) \alpha_w \frac{1}{B_w^D} + \overline{b}_w^D, \quad (10)$$

$$\overline{R}_w := \theta_w \alpha_w \frac{1}{B_w^R} + \overline{b}_w^R, \quad (11)$$

$$\overline{\overline{R}}_w := (1 - \theta_w) \alpha_w \frac{1}{B_w^R} + \overline{b}_w^R. \quad (12)$$

For relation (7) to hold, the greater the inaccuracy on the estimate of the execution parameters, the greater λ_w should be chosen.

The computation of \overline{D}_w , $\overline{\overline{D}}_w$, \overline{R}_w and $\overline{\overline{R}}_w$ requires the values of α_w and θ_w :

$$\alpha_w := \frac{\tau - 2\overline{f}_w}{\frac{1}{F_w}} \quad (13)$$

$$\theta_w := \phi_w \theta_w^{\max} + (1 - \phi_w) \theta_w^{\min}; \quad (14)$$

where θ_w^{\min} and θ_w^{\max} are respectively deduced from (15) and (16) as follows

$$\theta_w^{\min} := \frac{\frac{\alpha_w}{B_w^D} + \overline{b}_w^D - \overline{f}_w}{\alpha_w \left(\frac{1}{F_w} + \frac{1}{B_w^D} \right)}, \quad (15)$$

$$\theta_w^{\max} := \frac{\frac{\alpha_w}{F_w} + \overline{f}_w - \overline{b}_w^R - \overline{b}_w^D}{\alpha_w \left(\frac{1}{F_w} + \frac{1}{B_w^D} + \frac{1}{B_w^R} \right)}. \quad (16)$$

When

$$\phi_w = 0.5, \quad (17)$$

the risks of intra-round idleness and inter-round idleness are well balanced. But the value of θ_w could be obtained by averaging θ_w^{\min} and θ_w^{\max} with other weights than 0.5. 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:

$$\tau \geq \sum_{w=0}^{M-1} d_w. \quad (18)$$

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, thanks to (13) respectively, and loops until (18) holds. Once CIP is processed, proper AS4DR scheduling starts with a first round which sets the initial time-lags between successive worker 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

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

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

Let us assume that for each worker w

$$0 \leq \lambda_w \leq \min(\Lambda_w, \Lambda_{w-1}) \quad (\text{modulo } M). \quad (19)$$

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. Inequality (19) cannot hold when the right hand side member is negative. Roughly, this can take place when the number of workers M is too big compared to the computation over communication ratio.

IV. TESTS OF ADAPTIVITY FOR AS4DR

Simulations have been conducted in order to assess experimentally the AS4DR adaptivity. All the results presented in this section have been obtained with the SimGrid framework [3]. We consider a star-shaped platform (Figure 8 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 set is randomly allocated to 100 workers among the 1000 workers which constitute the platform.

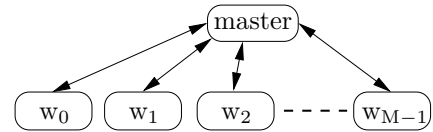


Figure 8. Star-shaped platform

In order 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

	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

is possible, due to (13), 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 the inaccuracy of the execution parameters estimates 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 (13) and the reference values contained in Table I:

$$\alpha_{w,1} := (1 \pm \iota_k) (\alpha_w)_{\text{ref}}; \quad (20)$$

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.

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
VALUES OF ι_k

The values of ι_k characterize the inaccuracy of the execution parameters estimates; 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 9 shows the measured CPU-efficiency of the whole platform for each scheduler, as a function of estimates' inaccuracy. For this simulation, which lasted 2000 sec-

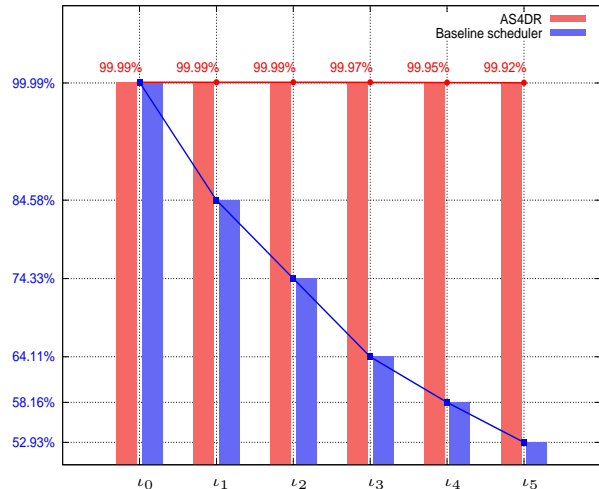


Figure 9. CPU-efficiency of the platform as a function of inaccuracy

onds, the value of τ computed by CIP equals 3 seconds. When the execution parameters are exactly known (ι_0), both schedulers offer 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 with AS4DR is 2.13 times higher than with the Baseline scheduler.

Figures 10 and 11 show, for both schedulers, the variation (during the very first rounds) of the CPU-efficiency of a worker with maximum estimates' inaccuracy (ι_5). In abscissa are reported the dates of CPU-efficiency measurement; after the master received the result for the subchunk 1 of the first round and after it received the result for the subchunk 2 of the other rounds. Beyond 26 seconds, the CPU-efficiency of this worker becomes steady, for both schedulers. With AS4DR (Figure 11), the existence of an asymptotic limit to $\sigma_{w,i}$ explains the asymptotical character of the evolution of the CPU-efficiency. For both schedulers this asymptotic CPU-efficiency is numerically reached from the very first rounds. With AS4DR, the smallness of the computation latency of the worker under consideration: 10^{-4} seconds, in front of τ : 3 seconds, explains the magnitude of the speed of convergence; after just a few rounds, the CPU-efficiency with AS4DR for this worker is and stays 17.86 times higher than the one obtained with the Baseline

scheduler. On Figures 10 and 11 it can be noticed, for

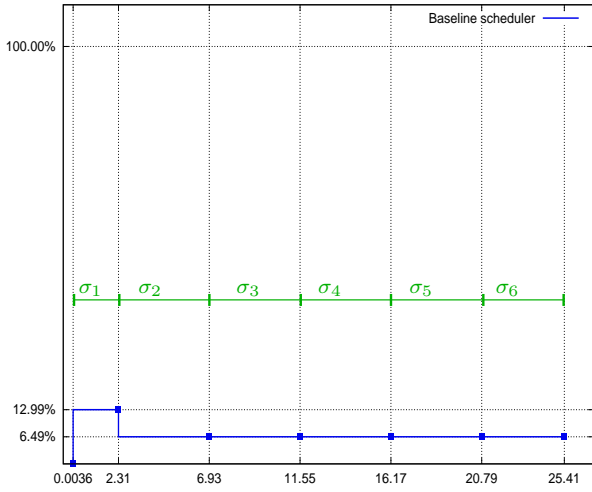


Figure 10. Baseline: CPU-efficiency of a worker as a function time (in seconds)

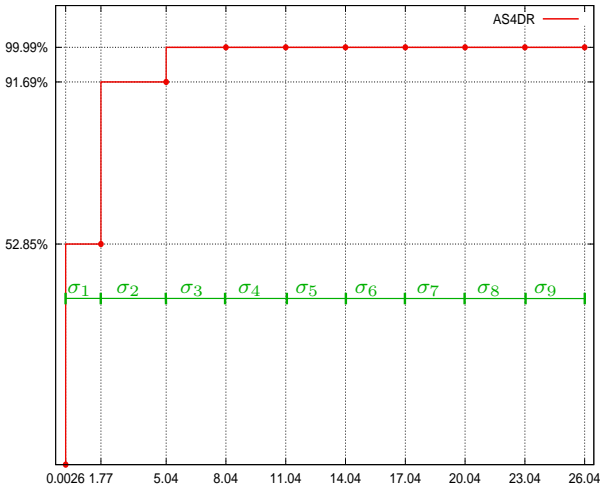


Figure 11. AS4DR: CPU-efficiency of a worker as a function time (in seconds)

the worker under consideration, that the dates of end of rounds (in abscissa) become asymptotically periodic. For instance, with AS4DR, the values $(\sigma_i)_{3 \leq i \leq 9}$ are identical and equal τ (computed with CIP): 3, while the values $(\sigma_i)_{2 \leq i \leq 6}$ are equal to 4.62, with the Baseline scheduler. The wanted asymptotical period τ : 3, is not reached for this worker with the Baseline scheduler.

In order to estimate the duration $(\sigma_{w,i})_{i \geq 1}$ of a round for all the workers of the platform, Figure 12 shows its average and its standard deviation over all the rounds and for all the workers, for different values of estimates' inaccuracy. For both methods, Figure 12 shows, on the one hand, that the

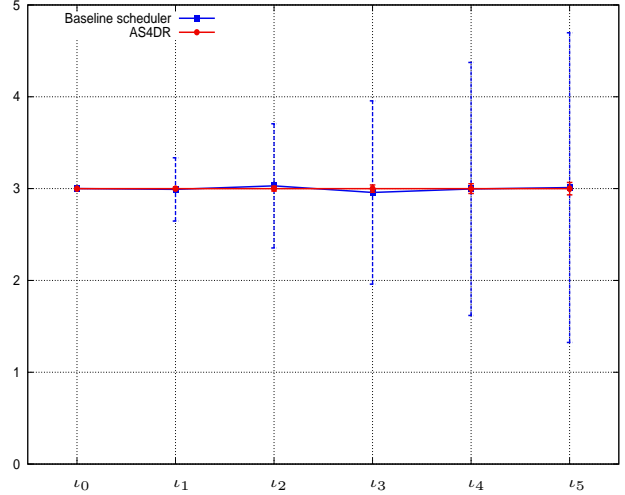


Figure 12. Average and standard deviation of σ as a function of inaccuracy

average of the values of $(\sigma_{w,i})_{i \geq 1}$ is close to the wanted value τ and, on the other hand, that the standard deviation of these values increases when estimates' inaccuracy of the execution parameters increases. It can be noticed that the standard deviation increases much faster with the Baseline scheduler, in relation to estimates' inaccuracy, than with AS4DR. For instance, the standard deviation with the Baseline scheduler is 24.98 times higher than with AS4DR for l_5 . The ratio: 17.86, observed between the CPU-efficiency of both schedulers for l_5 on Figure 9, is a consequence of this disparity between the standard deviations of the values of $(\sigma_{w,i})_{i \geq 1}$. By adapting at each round the workload for each worker AS4DR reduces, on the one hand, the risk of contentions inherent to the 1-port communication model and, on the other hand, the risk (in Round-Robin mode) of waiting for the result of a previous worker.

B. Time-varying context

The previous subsection has shown the ability of the AS4DR method to adapt to a priori poor estimates of execution parameters in steady execution contexts. 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 13. 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:

$$\begin{aligned} B_w^D &:= (1 - \delta_k) (B_w^D)_{\text{ref}}, B_w^R := (1 - \delta_k) (B_w^R)_{\text{ref}} \quad (21) \\ F_w &:= (1 - \delta_k) (F_w)_{\text{ref}}; \quad (22) \end{aligned}$$

where δ_k is a strictly positive number, the values of which are given in Table III. In this context, the coefficient δ_k characterizes the variations of the execution parameters and is called ‘‘dynamicity’’ in the sequel; with δ_4 the amplitude of the variation of the execution parameters is maximum, whereas it is minimum with δ_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.

δ_0	δ_1	δ_2	δ_3	δ_4
0.0	0.2	0.4	0.6	0.8

Table III
VALUES OF δ_k

Figure 14 shows the measured CPU-efficiency of the whole platform for each scheduler, as a function of the dynamicity. For this simulation, which lasted 2000 seconds, the value of

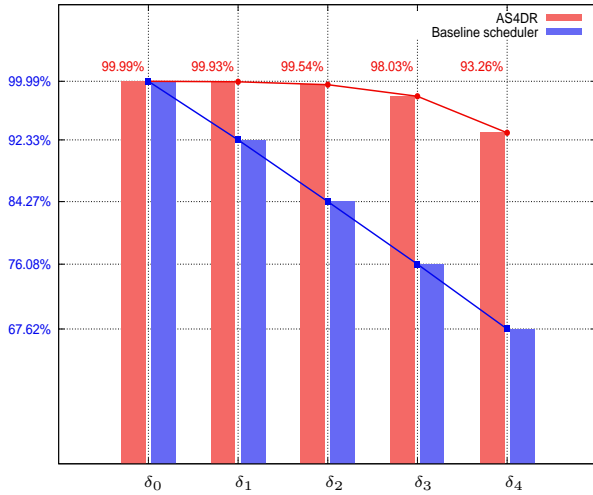


Figure 14. CPU-efficiency of the platform as a function of dynamicity

τ computed by CIP equals 3 seconds. When the execution parameters are steady (δ_0), both schedulers offered the same CPU-efficiency: 99.99%. As expected, the CPU-efficiency for both methods decreases when the dynamicity of the execution parameters increases. This decrease is significantly faster for the scheduler Baseline than for AS4DR. For instance, when the dynamicity equals δ_4 , the CPU-efficiency with AS4DR is 1.38 times higher than with the Baseline scheduler.

Figures 15 and 16 show, for both schedulers, the variation of the CPU-efficiency of a worker with profile P_9 during the first reduction of the value of the execution parameters, when the dynamicity is maximum: δ_4 . When time becomes equal

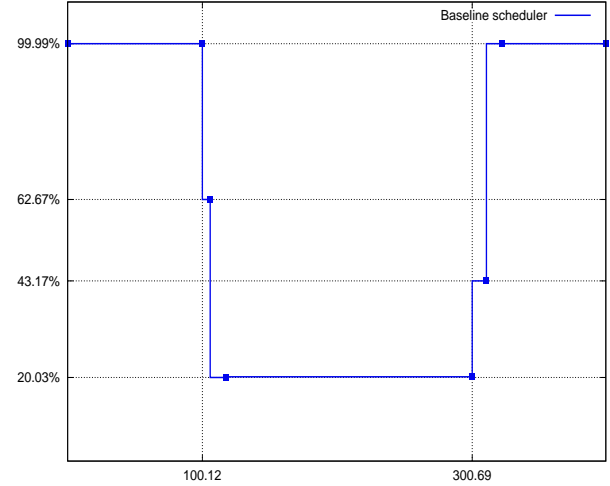


Figure 15. Baseline: CPU-efficiency of a worker as a function of time (in seconds)

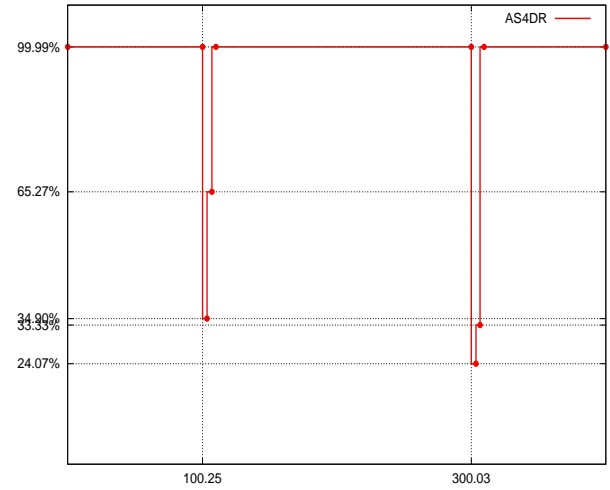


Figure 16. AS4DR: CPU-efficiency of a worker as a function of time (in seconds)

to 100 seconds each execution parameter, for the worker under consideration, is reduced by 80%. Then, when time becomes equal to 300 seconds each execution parameter recovers its reference value. With AS4DR (Figure 16), the smallness of the computation latency of the worker under consideration in front of τ , explains the quickness of the adaptation (3 rounds) to the change of context; on the falling edge as well as the rising one. In the absence of adaptation to the change of the characteristics of the platform, the CPU-

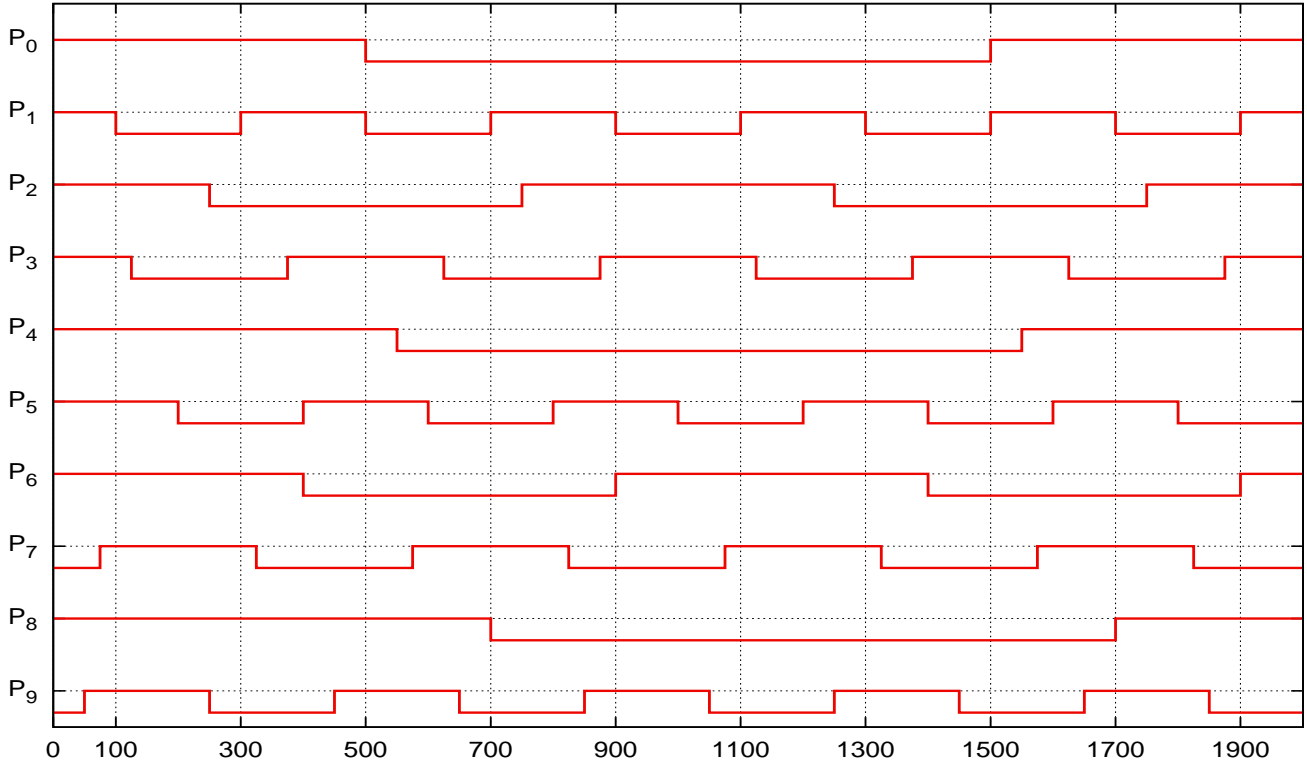


Figure 13. Perturbed execution parameter as a function of time (in seconds)

efficiency observed (Figure 15) with the Baseline scheduler is reduced during a long while (68 rounds).

Figure 17 shows the average and the standard deviation of $(\sigma_{w,i})_{i \geq 1}$ over all the rounds and for all the workers, for different values of the dynamcity. It can be noticed

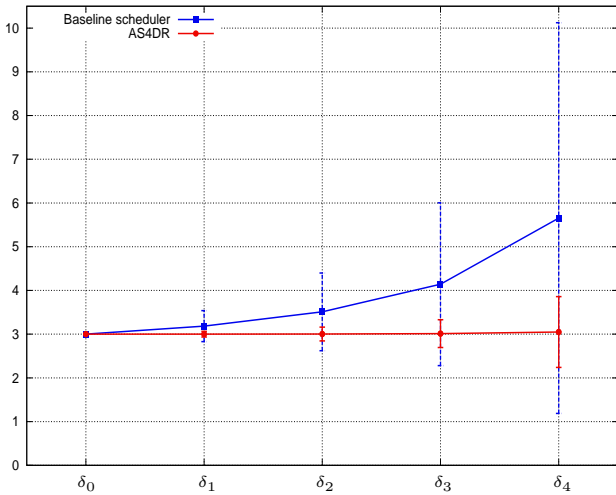


Figure 17. Average and standard deviation of σ as a function of dynamcity

that, with the Baseline scheduler, the average of the values $(\sigma_{w,i})_{i \geq 1}$ is an increasing function of the dynamcity. The reason for this is that the perturbation of the context, due to (21) and (22), is a reduction of the reference value of the execution parameters, given in Table I; which systematically increases the duration of the rounds. On the contrary, with AS4DR, the average of the values $(\sigma_{w,i})_{i \geq 1}$ stays close to the wanted value τ . For both methods, Figure 17 shows that the standard deviation of these values is an increasing function of the dynamcity. Once again, the better CPU-efficiency observed with AS4DR in relation to the one observed with the Baseline scheduler, for instance for δ_4 on Figure 14, is a consequence of the disparity between the standard deviations of the values of $(\sigma_{w,i})_{i \geq 1}$, shown by Figure 17.

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 bidirectional 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. The experimental

results presented in 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. The optimality of such a selection with respect to the global throughput will be 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.