

Parallel computation of the sequence of iterates of a function

Jean-Michel Fourneau Maël Guiraud Yann Strozecki

January 19, 2018

Abstract

1 Model

Let \mathcal{S} be a finite set of states and let f be a computable function from \mathcal{S} to \mathcal{S} . The aim of this short article is to propose practical methods to compute the sequence $(s, f(s), \dots, f^n(s))$ in parallel. When it is possible to compute efficiently $f^i(s)$ using only s and i , it is easy to compute the sequence in parallel by assigning one $f^i(s)$ to each processor. However, for general f we need $f^i(s)$ to compute $f^{i+1}(s)$ and there is no simple scheme to parallelize the computation of the sequence. Hence, we consider additional structure on \mathcal{S} to make the problem tractable. The set \mathcal{S} is equipped with a partial order \leq and it has a smallest element \perp and a greatest element \top . We assume that f is monotone, that is if s_1 and s_2 are elements of \mathcal{S} such that $s_1 \leq s_2$ then $f(s_1) \leq f(s_2)$. We call the problem of computing the sequence $(s, f(s), \dots, f^n(s))$ SIMULATIONPO when \mathcal{S} is a partially ordered set.

This models a concrete problem: the simulation of a random process. In that setting, the function f which describe the dynamic of the process also depends on the value of some random variable. In a simulation we use a random seed and then a random generator to generate the sequence of values of the random variable from the random seed. Say that the seed is a m bit integer in $[2^m]$, we have a pseudo random generator R which maps $[2^m]$ to $[2^m]$. A state of the system is $(s, y) \in \mathcal{S} \times [2^m]$ where s describes the state of the random process and y is the current random number. A transition of the random process is given by a function f from $\mathcal{S} \times [2^m]$ to \mathcal{S} thus a pair (s, x) is mapped to $(f(s, x), R(x))$ and we want to compute the iterates by this function.

To write parallel algorithm, we choose a simple PRAM model –exclusive read and exclusive write– with shared memory to neglect synchronization and communication problems. This simplification is reasonable in a multicore machine, since in our algorithms we use very few concurrent accesses to only limited informations, which can be dealt with atomic operations. However for the network of Raspberry Pi we use in experiments, we must also investigate the cost of communication.

2 Parallel computation

The main idea of the method is to divide the $n + 1$ values to compute into intervals of size t . If we know the initial state of each interval then we can compute the sequences of iterates independently. We denote by I_j the interval $\{tj, \dots, tj + t - 1\}$. Each processor will be assigned the task to compute the sequence on some interval I_j that is $(f^i(s))_{i \in I_j}$. We assume that we have some central memory in which the outputed sequence is stored and to which each processor can write. In the PRAM model, with an unbounded number of processors and zero cost of communication, it may be optimal to have t small and independent from n , but in our practical experiments we will set t to a much larger value.

2.1 Two bounds

We describe an algorithm which solves SIMULATIONPO in parallel, that is given s , n and an algorithm which computes f as inputs, it produces the sequence $(f^i(s))_{i \in [n+1]}$. For each interval I_j , we store the states s_j^{min} and s_j^{max} and one of the values ToCompute, INPROGRESS or DONE, which represents the status of the interval at some point of the algorithm.

The algorithm is the following: at the beginning, all intervals are marked ToCompute, for all $j > 0$, $s_j^{min} = \perp$, $s_j^{max} = \top$ and $s_0^{min} = s_0^{max} = s$. While there is a free processor P and an interval I_j in state ToCompute, select them. The status of I_j is set to INPROGRESS and the processor P computes the two sequences $(f^i(s_j^{min}))_{i \in [t]}$ and $(f^i(s_j^{max}))_{i \in [t]}$ iteratively. When the sequences are computed, we have access to $f^t(s_j^{min})$ and $f^t(s_j^{max})$ with one more application of f . If $f^t(s_j^{min}) > s_{j+1}^{min}$ or $f^t(s_j^{max}) < s_{j+1}^{max}$ then better bounds have been found and P sets $s_{j+1}^{min} = f^t(s_j^{min})$, $s_{j+1}^{max} = f^t(s_j^{max})$ and the state of I_{j+1} to ToCompute. Finally, if s_j^{min} is equal to s_j^{max} , the result of the simulation is stored as the solution on the interval I_j and I_j is set to DONE.

TODO: Ecrire l'algo dans un environnement algo

Algorithm 1 SMALLEST PARALLEL SANDWICH

Input: Size of the intervals t

```
// Initialisation
for  $i < \min(\text{nb\_machines}, \text{nb\_inter}-1)$  do
    Send  $I_{j+1}$  to the server  $i$ 
end for
// Main loop
while All the intervals are not DONE do
    Wait for a server to answer the results of  $\text{current\_interval}$ 
    if The server was computing a trajectory then
        set  $I_{\text{current\_interval}}$  to DONE
    end if
    if  $f^t(s_{\text{current\_interval}}^{\min}) > s_{\text{current\_interval}+1}^{\min}$  or  $f^t(s_{\text{current\_interval}}^{\max}) < s_{\text{current\_interval}+1}^{\max}$  // Better bounds have been found then
         $s_{\text{current\_interval}+1}^{\min} = f^t(s_{\text{current\_interval}}^{\min})$ 
         $s_{\text{current\_interval}+1}^{\max} = f^t(s_{\text{current\_interval}}^{\max})$ 
        set  $I_{\text{current\_interval}}$  to TOCOMPUTE
    end if
     $\text{next\_interval} \leftarrow$  search the first interval which is to TOCOMPUTE
    if  $s_{\text{next\_interval}+1}^{\min} = s_{\text{next\_interval}+1}^{\max}$  // The bounds are coupled then
        Wait a trajectory for  $I_{\text{next\_interval}}$ 
    else
        Wait some bounds for  $I_{\text{next\_interval}}$ 
    end if
    Send  $s_{\text{next\_interval}+1}^{\min}$  and  $s_{\text{next\_interval}+1}^{\max}$  to the current server
end while
```

Remark that the choice of an interval with state TOCOMPUTE by a free processor is not specified when there are several available. In practice we propose two ways to select it. The first is to chose the interval of smallest index which is available. We call this algorithm SMALLEST PARALLEL SANDWICH, smallest for the selection rule and parallel sandwich because it works by reducing the gap between the lower and the upper bound until a single correct state is found. The second is to cut $[n+1]$ into l consecutive meta-intervals (containing several I_j), where l is the number of processors used for the computation. Each processor is assigned to a meta-interval. A processor can be selected only to deal with an interval inside its meta-interval, and the smallest such interval is selected if there are several. We call this variant BALANCED PARALLEL SANDWICH, since we try to assign the same computing power to all parts of the sequence we compute.

Theorem 1. *Algorithms SMALLEST PARALLEL SANDWICH solve the problem SIMULATIONPO in time bounded by $O(cn)$ where c is a bound on the time to compute f .*

Proof. To prove that the algorithm terminates in time $O(cn)$, we prove that, at any point in time, there is a processor which is computing a new part of the sequence we must produce. We do the proof by induction on the computation time of the algorithm. Assume a processor P finishes to compute the sequence on an interval. Then consider the smallest interval I_j which is not marked DONE. If it is in state INPROGRESS, then a processor is computing the solution

on I_j . If it is in state `TOCOMPUTE`, it will be selected by the free processor P and since I_{j-1} is in state `DONE`, $x_j^{min} = x_j^{max}$ and P will compute the solution on I_j .

We must also prove that when an interval is set to `DONE`, the right sequences of states has been computed. It follows from the fact that $x_j^{min} \leq f^{jt}(s)$ and $x_j^{max} \geq f^{jt}(s)$. It can be proved by induction on the number of times the variables x_j^{min} and x_j^{max} are updated, using the monotonicity of f and the initialization of these variables to \top and \perp . Hence when $x_j^{min} = x_j^{max}$ then it is also equal to $f^{jt}(s)$ and this value is enough to compute the sequence on I_j . \square

Remark first that the time bound $O(cn)$ is the same as for the *sequential* algorithm which just applies f repeatedly. It is not worst, and it can be better, if at some point we have $x_j^{min} = x_j^{max}$, we say that the two sequences are coupling. If this phenomena happens frequently, the speed-up can be important. For instance, in the best case, we have a coupling on each interval the first time it is simulated, then the sequential time is bounded by $O(ct)$.

Finally, we can prove a theorem for `BALANCED PARALLEL SANDWICH` similar to Theorem 1, by using the fact that the steps of all processors are synchronized, therefore when one processor is free, they all are, which is enough to prove that one processor will always be computing a new part of the solution.

TODO: Faire un petit modèle probabiliste qui pour une proba de coupler donnée combien on va faire de calculs. Ça serait bien de se servir de ça pour couper des intervalles de la bonne taille. Peut être faire un calcul simple pour un tradoff entre nombre de processeurs et temps de couplage.

Remark that all processors do the simulation using only their private memory and the two states at the beginning of the interval. This scheme is thus reasonably adapted to a distributed computing environment where the cost of transmitting information between processors is high.

In our application to the simulation of a random process, the states of the process are often equipped with a partial order, but the random integers are not. However, we have the following property in our random system, if $s_1 \leq s_2$ then $f(s_1, x) \leq f(s_2, x)$. In fact, the random value alone is used to select an action to apply to the system and the action does not depends on the state of the system. Rather than simulating a given sequence $(s, f(s), \dots, f^n(s))$, we want to compute a random sequence that is to obtain a sequence with the right probability, which is a problem slightly different of `SIMULATIONPO`. We use the following trick to transform our problem: instead of using a single seed to generate all pseudo random values, one seed for each interval I_j is used. We have a collection of seeds $x_0, \dots, x_{n/t}$ and the i -th random number will be equal to $R^{i \bmod t}(x_{i/t})$ instead of $R^i(x)$. Since a random integer in the sequence depends only on the previous one, the random process is the same, that is all realizations of the random process still occur with the same probability as when only one seed is used. However, we can now abstract away the random numbers from our sequence, even when we divide the sequence into intervals, and use the algorithms we have described which solve `SIMULATIONPO`.

2.2 Algos avec intervalles de taille variables

We can slightly improve the algorithm `BALANCED PARALLEL SANDWICH`. First we make it simpler by cutting $[n+1]$ into l intervals if we have l processors, that

is at first intervals and meta-intervals are the same. Then when a processor a processor has correctly simulated its interval, we cut the remaining interval to simulate into intervals of similar sizes and affect one processor to each one.

2.3 One bound

3 Experiments

3.1 Random process

The random process we use in our experiments is the following one. We have a system composed of m finite queues of capacity `BUFF_MAX` in tandem. A queue is characterized by its number of client C_i . For each queue i , three events can occur:

- An arrival; C_i is increased by one, if it is not already to `BUFF_MAX`.
- A service; if there is at least one client in i , a client leave the queue i and goes in the queue $i + 1$. C_i is decreased by one and C_{i+1} is increased by one. If C_{i+1} is equal to `BUFF_MAX`, the client is lost.
- A departure: the client leave the system. The number of client of the queue is decreased by one if $C_i > 0$.

For the queue i , every event have a probability denoted by respectively a_i , s_i and d_i . The last queue has no service, thus $s_m = 0$. There is thus a total of $3m - 1$ different events that can change the system.

The sequence that we want to parallelize is a succession of drawing of one of those $3m - 1$ events. Every event is drawn randomly, following a pseudo random generator.

In the following experiments, the parameters a_i , s_i and d_i of the queues are set by the following rules:

- $P = 0.75$
- $\mu = 300$
- $\text{load} = 1.1$
- $a_0 = \mu \times \text{load}$
- $a_i = (1 - P) \times a_0, i \in [1 : m]$
- $s_i = P \times \mu, i \in [0 : m]$
- $d_i = (1 - P) \times \mu, i \in [0 : m]$

All those probabilities are then normalized such that $\sum_{i=0}^m a_i + s_i + d_i = 1$.

With those parameters, every queue have an average load of 1.1. This load has been experimentally calculated as the load for which the coupling time of the two bounds in the largest.

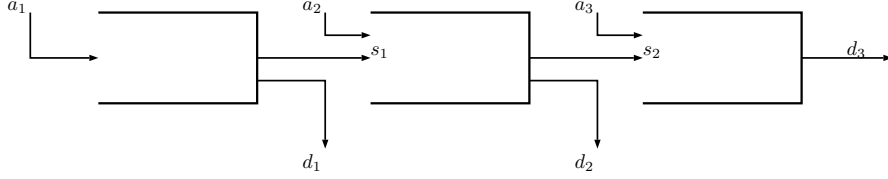


Figure 1: A system with 3 queues

3.2 Hardware and software used

3.3 Experimental settings

Our experimentations are made with the following set-up. The main program run on a MacBook Air, with a 2.2 GHz Intel Core i7, and 8 Go of RAM DDR3 at 1600 MHz. The operating system of the machine is macOS High Sierra v10.13.1. The source code is compiled with gcc version 7.1.0 (Homebrew GCC 7.1.0 – without-multilib).

The simulations of each interval are dispatched on up to 7 Raspberry Pi 3 , Model B, with 1GB of RAM. Their operating system is Raspbian GNU/Linux 8.0, installed on a micro SD card element14 with a size of 8GB. The source code is compiled with gcc version 4.9.2 (Raspbian 4.9.2 – 10). All the machines are connected to a local network through an HP 14-10 8G switch.

iiiiii HEAD All the machines are connected to a local network through an HP 14-10 8G switch.

During our experiments, we use the socket with the protocol TCP/IP to allow the servers and the master to communicate. All the machines were wired to a local network without access to any external networks. The source code of the simulation is written in C and available on the webpage of the authors.

TODO: lien vers site yann avec code

3.4 Network impact

Since our network uses a local network with TCP/IP, one must investigate the cost of the network in our simulations. In our source code, we fixed the size of a standard master/server message to 24 integers and the size of the server/master messages depends of the number of queues in the simulation. One must focus on the cost of sending a trajectory, which is by far the largest message, indeed, the size of a trajectory is equal to $t \times m$ (t is the size of the interval, and m is the number of queues). On the following experiment, we tried to compare the evolution of the cost of the computing time and the network cost for different number of queues going from 1 to 10. We arbitrary set the size of the interval t to 10,000, which is a good order of magnitude considering the following experiments. The points on fig— ?? are some average times computed on 100 different simulations.

As we can see, both the communication an computing times of a trajectory increases linearly with the number of queues, but the communication time is increasing faster. A linear regression allowed us to calculate the slope of those two point clouds; the slope of the communication time is 3.390585 while the slope of the computing time is only of 0.885585. Both of those two correlations

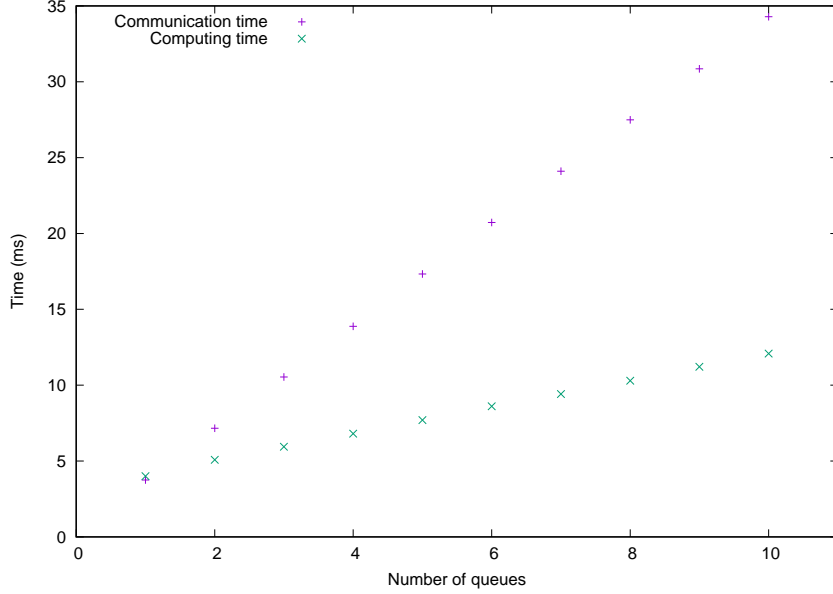


Figure 2: Evolution of the computing and communication time of a trajectory of 10,000 events

coefficients are greater 0.999. Thus, to avoid an over domination of the network time, which would hide every interesting results, in the following experiments, we will consider a little number of queues in tandem.

3.5 Algorithm performance evaluation

We now want to measure the performance of the SMALLEST PARALLEL SANDWICH and the BALANCED PARALLEL SANDWICH, compared to the One bound algorithm (cf sec 2.3). In the following experiment we thus take a random process with 5 queues, which allows us to have a reasonable amount of computation in the servers, and a not too high $\frac{\text{Network time}}{\text{Computing time}}$ ratio. In order to choose the size n of the simulation, we made an upstream experiment to determine the average time between the coupling of the two bounds, for two queues, with the parameters we give in sec 3.1. Thus, this average coupling time is to 28,049, calculated over 100 simulations.

3.5.1 Long simulation

We first want to look at the behavior of our algorithms in a simulation in which there is an high probability to have a coupling during the computing of a single interval. We then set the size of the simulation to $n = 210,000$. Fig. 3.5.1 and 3.5.1 show the results of respectively, the average sum of the lengths of

the intervals computed, and the average total time if simulation in regard of the number of intervals, over 100 instances. By the way, the size of the interval directly depends of the number of intervals. The total size of the simulation is always the same, that is, 210,000. We used 7 servers in those experiments.

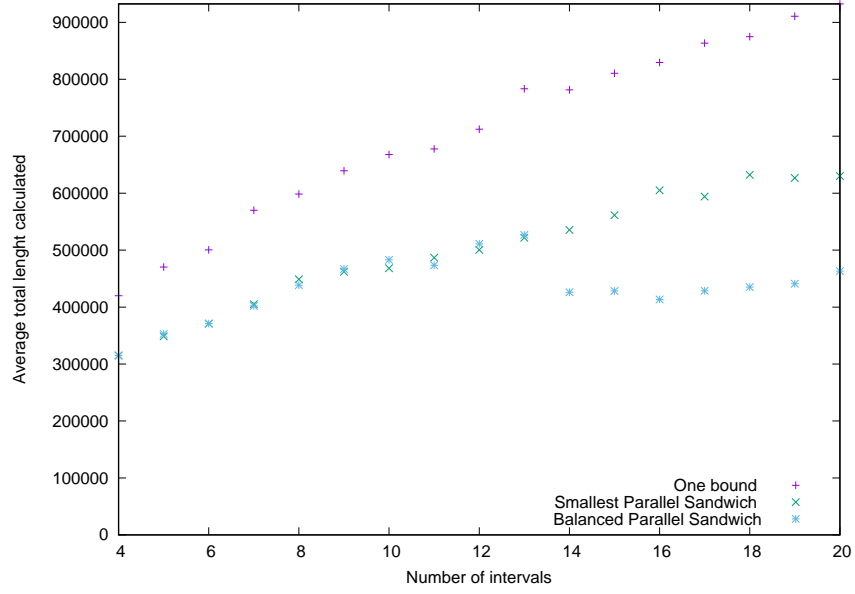


Figure 3: Evolution of the average of the sum of the lengths of the intervals computed for 100 simulations of 210,000 events.

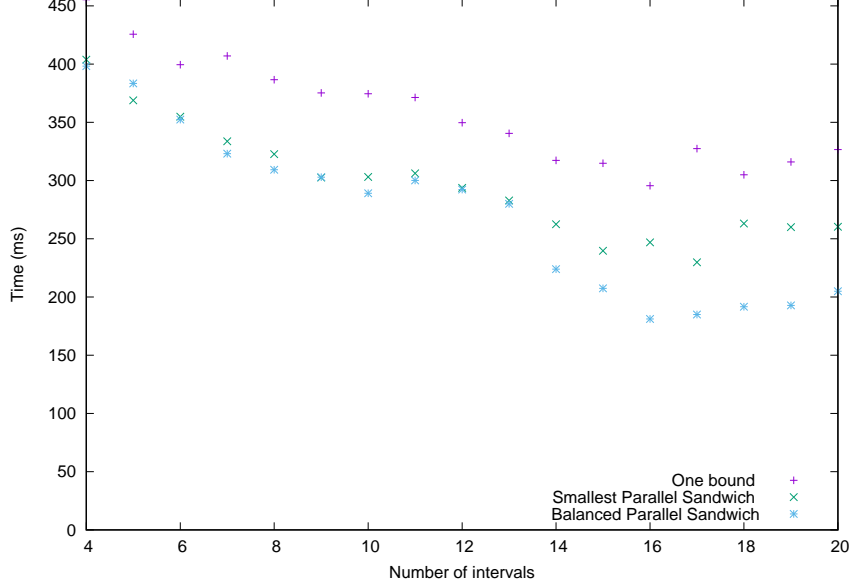


Figure 4: Evolution of average for 100 simulations of 210,000 events.

First, we remark that the SMALLEST PARALLEL SANDWICH and the BALANCED PARALLEL SANDWICH are always better than the One Bound algorithm. The weak increasing trend for One Bound and SMALLEST PARALLEL SANDWICH on fig. 3.5.1 might be due to the decreasing probability of coupling in one interval when the number of intervals increases. Nevertheless the BALANCED PARALLEL SANDWICH seems to have a different behavior after an undetermined threshold. We will study this phenomenon later.

TODO: pas vraiment en fait, j'utilise le fait que il y ait ce point du rupture pour comparer les performances en fonction du nombre de servers mais je ne sais pas du tout pourquoi ca fait ca??

In contrast to those results, as we can see on fig. 3.5.1, it looks like having few long intervals is more expensive in time than having more shorter intervals. This might be due to the network. Indeed, as we observed in the results of sec. 3.4, the more the interval is long, the more network increases. Moreover, the experiment in sec. 3.4 were made on a single server communicating with the master. Here, we have 7 servers which probably create some contention, and thus, some additional latency. On a single processor, the average time on 100 simulations for 210,000 events is to 162.13 ms.

To make sure that the network is dominating the simulation time, we investigate the activity of the servers during the experiments. Fig. 3.5.1 shows the average repartition of the computation, network and waiting time of the servers, during the previous experiments, for the computation of the BALANCED PARALLEL SANDWICH simulations.

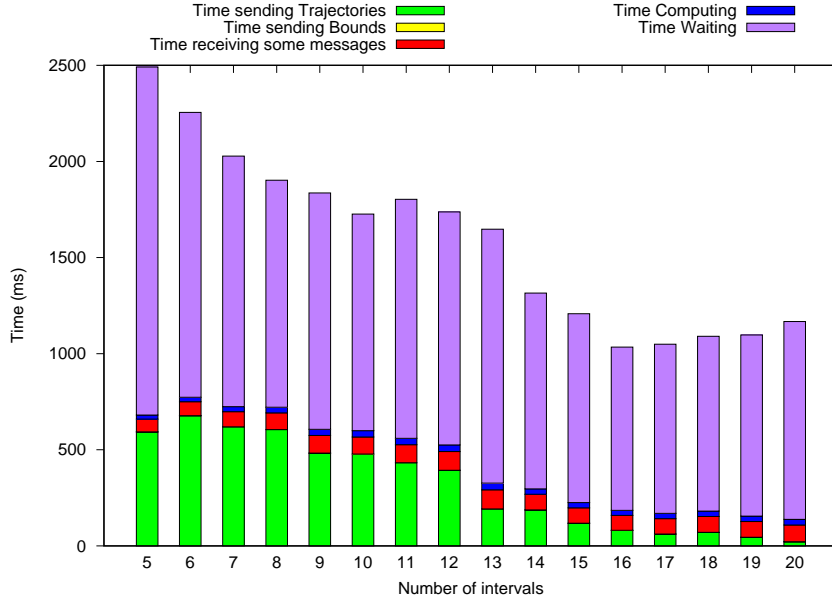


Figure 5: Repartition of the activities of the processors during the simulation.

First, we can see that, the most of the time, the servers are waiting for something to do. Since the master algorithm is not very complicated, we suppose that the network highly slows the communication between the master and the servers. The rest of the time, the servers are mainly computing some network operations (sending, receiving).

3.5.2 Short simulation

We then tried to look if there is a difference when the simulation is short. The following experiment are made with exactly the same parameters than upward, except the simulation length, which is now of 30,000. In this situation, the probability of coupling is very low.

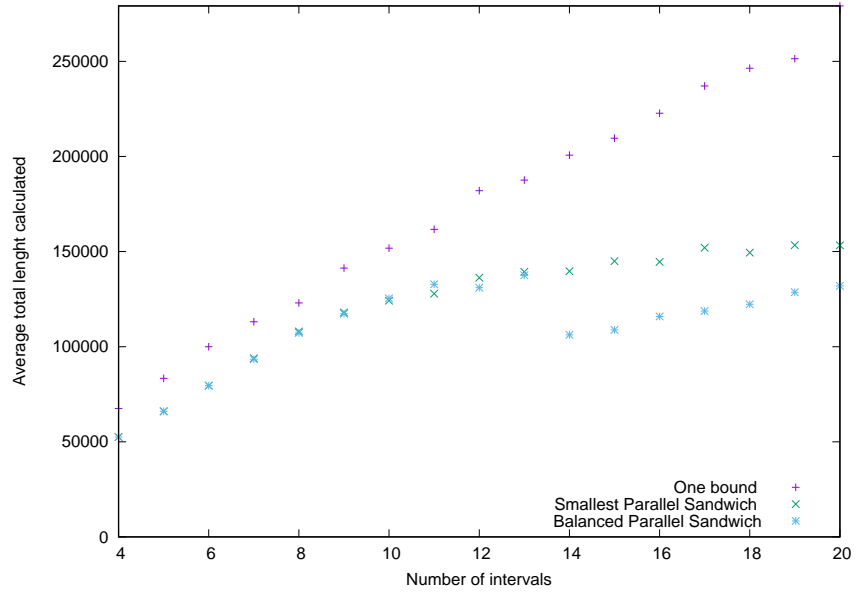


Figure 6: Evolution of the average of the sum of the lengths of the intervals computed for 100 simulations of 30,000 events.

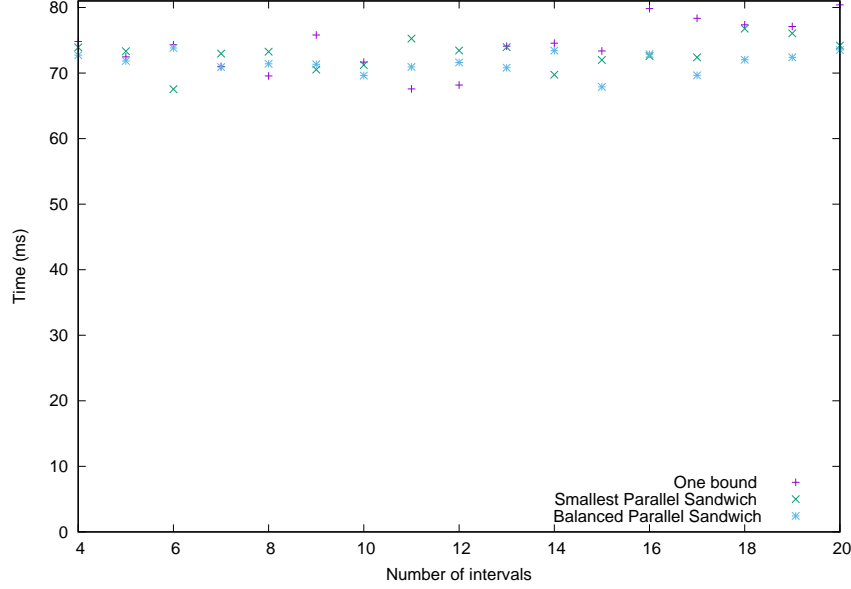


Figure 7: Evolution of average for 100 simulations of 30,000 events.

The trend of the curves on fig. 3.5.2 are notably similar to the curves of fig. 3.5.1. On the other hand, the three algorithm seems to run in a similar time on fig. 3.5.2. Once again, this is due to the domination of the network time, which is the same for all algorithms, since the messages are all short. On a single processor, the average time on 100 simulations for 30,000 events is to 23.06 ms. **TODO: expliquer un peu plus, de toute facon on calcule tout les intervals en attendant le gars du debut vu que ca a tres peu de chance de coupler**

=====

TODO: Dire ce qui est utilisé comme techno entre les raspberry et donc la latence induite (à mesurer) `43f1c719ccab17cdd0288dfaf16a9cd2cea56835`

3.6 Number of server used

We now want to study the impact of the number of server used. As we noticed in fig. 3.5.1, there is an interessant point when the number of intervals is twice the number of servers (this threshold was observed for any number of servers between 2 and 7). We then look at the execution time of 100 simulation from 2 to 7 servers, whith the BALANCED PARALLEL SANDWICH algorithm. The number of event n of the simulation is set to 210,000, and the number of intervals is always set to $2 \times$ number of servers, since it is the first number at which the execution time of the algorithm starts to stagnate to it's lower value, for any number of processors.

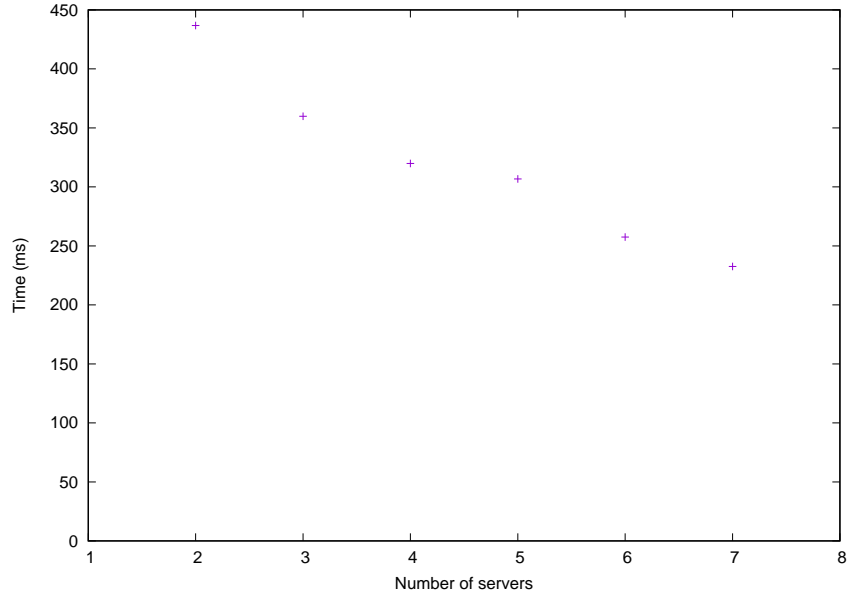


Figure 8: Impact of the number of servers on the execution time.

As we can see on fig.3.6, the more we have servers, the faster the algorithm runs. **TODO: bof cette conclusion, vu que c'est le reseau qui masque tout**

Practical problems: cost of the network transmission, especially for transmitting long sequences. -; measure the time of a two way trip for a small message and the time of sending an interval. To say that in practice we will not compute the whole sequence but statistic on it which could help reduce the use of the network.