

# Computação em Nuvem e Virtualização

## MEIC-T

### Grupo 14

Rafael Soares  
87675

Miguel Barros  
87691

Gonçalo Pires  
94112

## Architecture

In this phase, we have the Solver classes instrumented and the WebServer code up and running on the EC2 AWS with the default Load-Balancer and Auto-Scaler.

In the next section we will explain what classes we instrumented, why we instrumented those classes, which metrics we chose and why and how we modified the WebServer to handle multi-thread. Finally, we explain our plans for future implementation.

### Instrumentation

On a first approach, we instrumented the whole package solver for all metrics available in order to study the behaviour of each type of solving technique.

We tested these solvers by running three puzzles of each size with each of the available solving techniques (except for the 25x25 puzzle size, which only has 2 puzzles available) and an ample sample size of unsigned entries.

We did this to understand how the solver performance is affected by different puzzle sizes and number of unassigned entries. (this is not the only thing that affects performance but it is the only one we can control)

We save these results in JSON files and generate graphics from these results.

By analysing these graphics and calculating the correlation between Instruction Count (which approximately gives us the execution cost) and all other metrics, we came to the conclusion that most metrics show a correlation of >95%. With that in mind, we chose the metric with least amount of overhead to be the one instrumented during runtime, giving us a close approximation of the cost that puzzle will have on the system while giving us the least amount of overhead from instrumentation, giving a significant performance boost. With that said, we have picked the following metrics for each solving strategy:

BFS - Method Count

CP - New Count

DLX - New Array Count

These metrics were the ones with the least number of occurrences, hence will give us the

least amount of overhead, while still having a high correlation to the instruction count

We have decided to only instrument the three solver classes: SudokuSolverDLX, SudokuSolverCP and SudokuSolverBFS.

This is due to these classes having the majority of the changes between solver strategies, given that the distribution of the metrics count is the same for the whole package and for the metrics of only the solver classes, as seen by the comparison between the graphs of Annex A and the ones of Annex B. This allows us to reduce the overhead further.

### WebServer

The WebServer keeps an HashMap which associates each Thread solving a puzzle with the parameters of that puzzle.

In the instrumentation classes we also keep an HashMap that associates each Thread with the metrics for it's request.

When the thread finishes solving the puzzle, it saves the parameters of the puzzle with the associated metrics in a JSON file for future statistics by obtaining them from the HashMaps with their correspondent Thread ID. Finally the entries of the current thread in the two HashMaps are removed. This way, the server does not mix metrics from different requests while being multi-threaded (since the server uses a cached thread pool which can reuse threads).

## Future Implementation

### Metrics Storage System

The collected metrics must be permanently stored so that the load balancer can access them and use them to decide where to forward each request. We plan to store the metrics obtained from each request in the dynamoDB. To reduce the overhead, we plan to send the metrics in chunks, i.e send all metrics every X requests or only send when the server is not massively overloaded.

Alternatively, we could have instances dedicated to collecting metrics and sending them to dynamoDB, which would synchronize

among themselves to ensure that the dynamoDB stays always consistent. Each register in the database will have the parameters of the puzzle, the values of the metrics obtained and the cost associated with that puzzle. When the Load Balancer receives a query with the same parameters, the LB can estimate more precisely the cost of the query.

### Request Cost

All correlations we have seen between metrics are linear. As such, we can assume the cost of each request to be the number of the instrumented metric for its respective strategy. The instrumented metric must be converted into raw instruction cost, which can be done by a simple multiplication for each algorithm. The instruction cost for each strategy was obtained using the following formula:

$$InstructionCost = \#ChosenMetric * \frac{\#Instructions}{\#Metric}$$

Since the correlation is linear, the ration  $\frac{\#Instructions}{\#Metrics}$  will be constant.

With that in mind, these were the obtained Instruction Costs for each strategy:

BFS:

$$InstructionCost = \#Methods * 2.5 * 10^2$$

CP:

$$InstructionCost = \#NewCount * 3 * 10^2$$

DLX:

$$InstructionCost = \#NewCount * 2 * 10^6$$

The  $\frac{\#Instructions}{\#Metrics}$  ratio for each strategy was obtained from the ratio between the instruction count and the metric chosen, as obtained from the graphs in Annex A.

### Load Balancer

When the load balancer receives a request:

If it is the first request of its kind (i.e a unique configuration of board and unassigned entries) the load balancer must assume the worst case, by choosing the cost of a request with similar parameters (we bootstrap this process by providing the experimental results already obtained).

Otherwise, since the cost of each request is constant for the same request, we can know for sure what it will cost (the cost will be stored in dynamoDB).

The load balancer can then keep track of the load of each server and choose the least busy server for each request it receives. To

decrease the accesses of the load balancer to the dynamoDB, we can cache responses in a LRU manner and use those to infer the cost of new requests.

Since requests can be long running, we must periodically update the cost of each server to account for requests that started a long time ago. This can be done by periodically subtracting the cost of each long running request of each server by a fixed amount (since it is expected to be a direct correlation between the instrumented metrics and the cost in CPU and time of each request).

### Auto Scaling

Since the load balancer can only distribute work between the running instances, it is the job of the auto scaler to detect when the computing nodes are overloaded and launch new instances to handle the demanding workload. Conversely, it must also shut down instances when demand is small to spare resources. Launching new instances is simple: as every instance in AWS has the same computational power, we just detect when the load on the system as a whole is higher than the capacity of the current running instances and launch new instances until we can cope with the current running cost. To scale down, we must take into account several factors to decide which instances to terminate:

→ Longer running instances are more preferable than new instances, since they are more likely to have reduced performance (the program might open files it never closes or allocate memory that is never collected) while also checking for the last payment date, as we do not want to shut down instances that we have already paid for. Conversely, if we plan on shutting down an instance we can give it work to do for it's current time slot, relaxing the load balancing for that instance but still ensuring that all requests will be replied before it's time slot expires and it is terminated.

→ Instances with lower load are preferable, since we must for the instance to not be working before shutting it down or it's state will be lost.

## Annex A - Solver Class Metrics

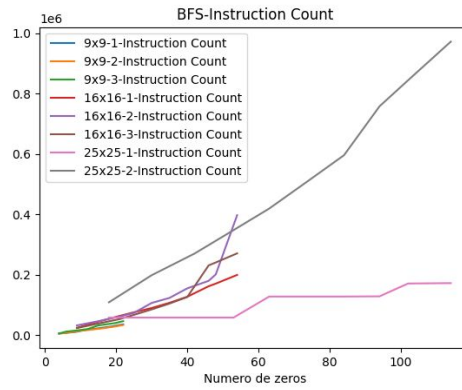


Figure 1 - Variation of Instruction Count with the number of unassigned entries for the BFS Strategy

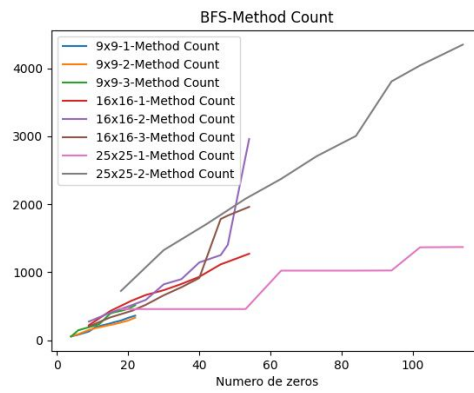


Figure 2 - Variation of Method Count with the number of unassigned entries for the BFS Strategy

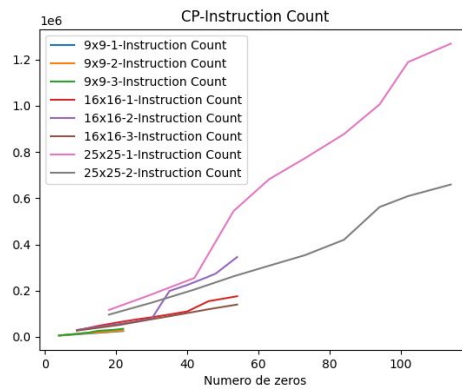


Figure 3 - Variation of Instruction Count with the number of unassigned entries for the CP Strategy

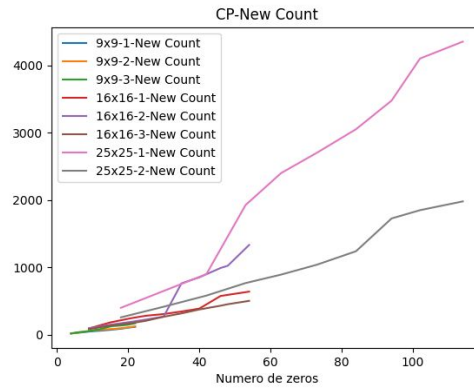


Figure 4 - Variation of New Count with the number of unassigned entries for the CP Strategy

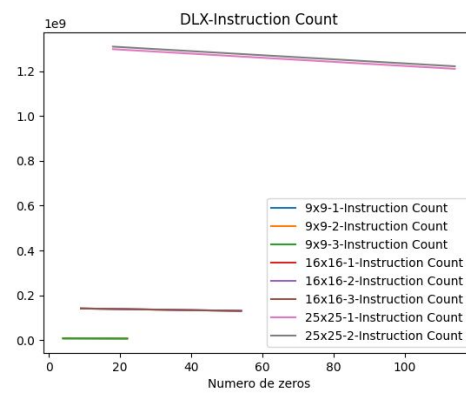


Figure 5 - Variation of Instruction Count with the number of unassigned entries for the DLX Strategy

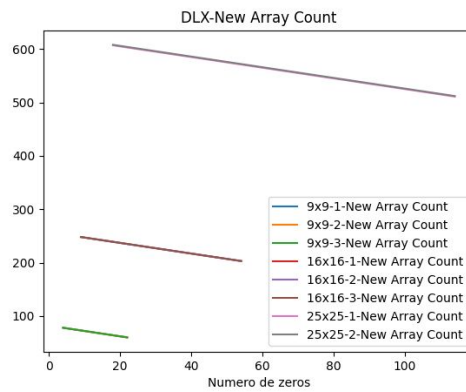


Figure 6 - Variation of New Array Count with the number of unassigned entries for the DLX Strategy

## Annex B - Solver Package Metrics

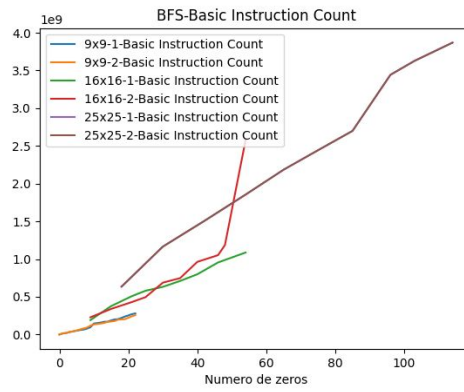


Figure 7 - Variation of Instruction Count with the number of unassigned entries for the BFS Strategy

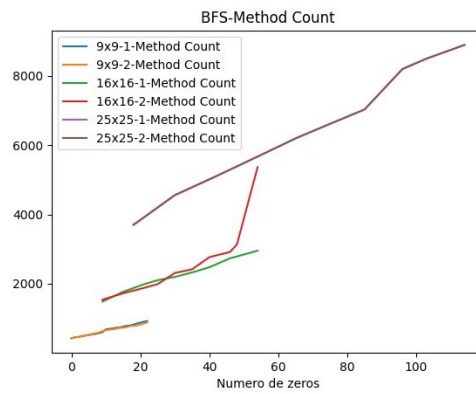


Figure 8 - Variation of Method Count with the number of unassigned entries for the BFS Strategy

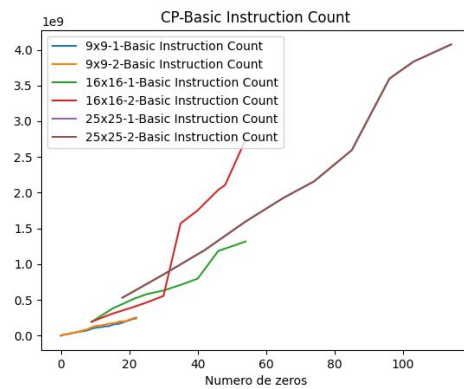


Figure 9 - Variation of Instruction Count with the number of unassigned entries for the CP Strategy

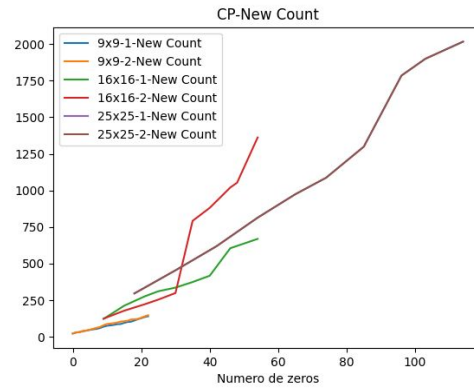


Figure 10 - Variation of New Count with the number of unassigned entries for the CP Strategy

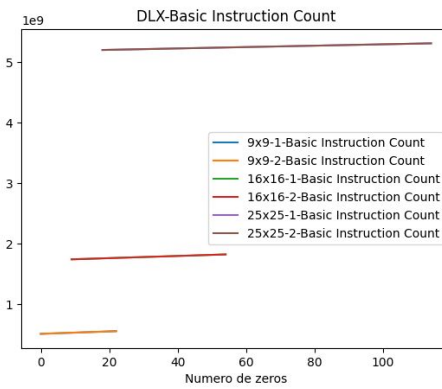


Figure 11 - Variation of Instruction Count with the number of unassigned entries for the DLX Strategy

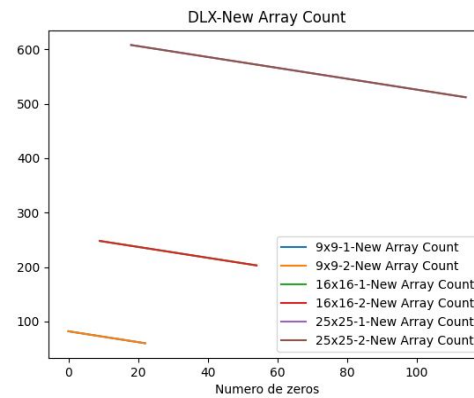


Figure 12 - Variation of New Array Count with the number of unassigned entries for the DLX Strategy

## Annex C - Correlations of Instrumented Metrics and Instruction Count

	BFS	CP	DLX
Store Count	0.9765	0.9885	0.9999
Load Count	0.9999	0.9999	0.9999
Basic Block Count	0.9999	0.9999	0.9999
Branch Count	0.9999	0.9999	0.9999
Method Count	0.9800	0.9951	0.9602
Field Load Count	0.9998	0.9999	0.9981
Stack Depth	0.7319	0.6222	0.3056
New Count	0	0.9971	0.9638
Field Store Count	0	0	0.9160
New Array Count	0	0	0.9655

Table 1 - Correlation of each instrumented metric for each strategy with the Instruction Count of each strategy



## **Annex D - AWS Settings**

### Load Balancer:

- Ping Protocol - HTTP
- Ping Port - 8000
- Ping Path - /lb
- Response Timeout - 5 seconds
- Interval - 30 seconds
- Unhealthy threshold - 3
- Healthy threshold - 10

### Scaling Policy:

- Execute policy : CPUUtilization  $\geq 30$  for 180 seconds
- Take the action: Remove 1 capacity units when  $30 \leq \text{CPUUtilization} < +\infty$

### Increase Group Size:

- CPUUtilization  $\geq 60$  for 180 seconds
- Add 1 capacity unit when  $60 \leq \text{CPUUtilization} \leq +\infty$
- Instances need 30 seconds to warm up after each step
- Grace Period: 30 seconds
- Minimum instances: 1
- Maximum instances: 5

## **Annex E - Explanation of each metric**

New Count - Number of objects allocated

New Array Count - Number of arrays allocated

Field Load Count - Number of loads of object fields

Field Store Count - Number of stores of object fields

Store Count - Number of stores

Load Count - Number of loads

Basic Block Count - Number of basic blocks executed

Branch Count - Number of branches executed

Method Count - Number of methods executed

Stack Depth - Maximum possible stack depth achieved by the program (this metric was decided to not be very useful)