

# Cloud Computing Assignment 2022-2023

## Assignment Details

The assignment requires you to design and implement an application for processing large data sets in parallel on a distributed Cloud environment. The solution to the assignment should meet the following specific requirements:

### **1. Implement an application that can decompose a data set and process it in parallel on a distributed Cloud environment.**

The data processing requirements for the assignment are simple matrix operations, implemented in a distributed Cloud environment:

- (i) Matrix addition
- (ii) Matrix multiplication

The matrices should be large and dense. You will need to decide how big the data sets should be according to how you want to explore the performance of your system. As a guide, a 1000x1000 matrix should be sufficient to start the investigation but you can increase this to simulate heavier computational loads.

The output of the system should be a single matrix resulting from the operations on two input matrices.

To get the best performance, it will be necessary to decompose the data into work packages and distribute them among several virtual machine instances. The decomposed data set should be sent to a maximum of 8 distributed processing nodes (virtual machine instances) in the AWS Cloud.

Your implementation should be able to queue multiple work packages and send them to each processing node once that node has finished its previous work. You can approach the solution from whatever direction you see best, including the implementation of your own queues or use of AWS Simple Queue Service (SQS).

The final step in your application will combine the outputs from each virtual machine instance to produce the completed result.

You should explain how you account for all the work units being completed, and how your system recovers from node crashes.

### **2. Optimise your solution using a range of different sized work packages to find the best balance of computation, communication and resource costs.**

How can you take advantage of the scalability of the Cloud platform to process the data in the shortest amount of time? For example, is it better to have smaller blocks of data and more processing nodes? To explore this, you should decompose your input data into a range of consistently sized blocks and test across different configurations (ie different numbers of available VMs, queue sizes, etc). Examine the bottlenecks in your design and explore how they can be mitigated.

In each case, the performance of the solution should be measured and reported. You will need to adopt a scientific approach to measuring the performance of the different configurations.

### **3. Ensure that your application that can monitor the progress of the calculation.**

The application should supply the following functionality:

- Set up the Cloud environment according to how many virtual machines the user wants to use. A more advanced approach could be to automatically spin up new VM instances according to the number of work packages currently in the queue.
- Submit the job to the distributed system
- Monitor the current progress of the calculation
- Report back the final result of the calculation
- Report the time taken to perform the calculation

Any tools you produce should be easy to use and provide clear outputs for the user.

### **4. Check the validity of your results.**

Ensure that the results from your calculation are correct by comparing them with a single-node based implementation. Use this implementation as a baseline for comparing the performance and correctness of your Cloud implementation.

#### **Submission:**

Include the following components in your submission:

- Produce a report that includes the following:
  1. Design documentation detailing the architecture of your solution. You should discuss how you came to your eventual design, highlighting how the Cloud platform influenced your decisions. Be specific about how you had to approach the solution for the matrix addition and matrix multiplication operations.
  2. Evaluation of the associated costs that your solution would incur due to use of computational, storage and communication resources - use the AWS costs website as a guide: <https://aws.amazon.com/pricing/>

You should use the costs of the platform to inform your design decisions and include this analysis in your design document. It is important that you can demonstrate that your solution is as cost effective as possible.

3. A validation plan explaining how you ensure that the results of your solution are correct, as well as examples of correct output. Include evidence that testing was carried out.
  4. Some measure of the computational performance of your solution with evidence to support your results
  5. Analysis and discussion of your solution, addressing the key advantages and disadvantages of the implementation. You should critically discuss design, performance and cost, supporting your argument with results from experiment. You may want to compare alternative designs and/or the optimisation steps you undertook.
- Any source code or script files you have generated, with full documentation and details of how they fit into your solution architecture.
  - The VM and applications available for testing.
    - You should ensure that the VM is not terminated (but stopped) after the submission deadline and ready to start for the lecturer to test when required; you should provide the necessary information for the lecturer to get access to VM and run the application for testing, including creation of an account where appropriate.
    - You shall turn off (i.e. stop) the AMI instance hosting the application at other times whenever you are not working on it. You will be informed after testing is done, at which point you will be able to terminate the AMI instance.

**Important: you have a limited budget for AWS cloud, so use it wisely.**

Note: There are no restrictions to the programming language or framework you may wish to use for the implementation of your assignment. However, consider the reporting requirements when choosing a suitable method.

### General Marking Scheme

- 40% software
- 40% report document
- 20% test design and results

Note: A more detailed marking scheme with marks broken down can be found in marking-breakdown.pdf on Canvas.

You must submit your assignment by:

- **9:30am, 9<sup>th</sup> January 2023 (full-time students)**
- **9:30am, 23<sup>rd</sup> January 2023 (part-time students)**