

Rapid re-purposing: How the foresight to replicate their HPC cluster in the cloud fueled the University of Liverpool's COVID-19 research

1st Yalin Zheng

University of Liverpool

Department for Eye and Vision Science

Liverpool, UK

Yalin.Zheng@liverpool.ac.uk

2nd Joshua Bridge

University of Liverpool

Department for Eye and Vision Science

Liverpool, UK

joshua.bridge@liverpool.ac.uk

3rd Yanda Meng

University of Liverpool

Department for Eye and Vision Science

Liverpool, UK

Yanda.Meng@liverpool.ac.uk

4th Cliff Addison

University of Liverpool

Advanced Research Computing

Liverpool, UK

C.Addison@liverpool.ac.uk

5th Manhui Wang

University of Liverpool

Advanced Research Computing

Liverpool, UK

Manhui.Wang@liverpool.ac.uk

6th Cristin Merritt

Alces Flight Limited

Bicester, UK

cristin.merritt@alces-flight.com

7th Stu Franks

Alces Flight Limited

Bicester, UK

stu.franks@alces-flight.com

8th Maria Mackey

Amazon Web Services

London, UK

mmacke@amazon.com

9th Steve Messenger

Amazon Web Services

London, UK

messteph@amazon.com

Abstract—The Advanced Research Computing team at the University of Liverpool created a mission-critical cloud replica of their on-premises HPC cluster, named Cloud Barkla, to function as a rapid-deployment service in October, 2019. Originally deployed to assist the relocation of their on-premises cluster, the team had the foresight to build a cloud cluster that could be repurposed at speed - and that need came sooner than expected. Dr. Yalin Zheng and his team utilised this capability to launch their AI-based medical image analysis research for accurate diagnosis of COVID-19 on the Cloud Barkla system in record time to allow lead researchers Joshua Bridge and Yanda Meng the ability to run scaling tests on their developed toolset in order to bring their concepts to peer review. This talk covers the rapid deployment of an HPC cluster fit for Dr. Zheng's AI requirements - demonstrating how cloud offers flexibility within a complete HPC services offering.

Index Terms—HPC Project Management, HPC in the Cloud, Artificial Intelligence

I. INTRODUCTION

The University of Liverpool made the decision in 2017 to build cloud computing into their overall HPC service, something which became fully realised in 2019 as the University's Computing Services Department began the process of completing their 2-year expansion on their HPC data center. The Advanced Research Computing (ARC) team's core HPC Cluster, Barkla, received a mission-critical cloud replica: the aptly named Cloud Barkla. Through a test cycle and complete production run the team - alongside managed service providers

Alces Flight - created and documented optimal settings for Cloud Barkla as well as potential candidates for future use. Their efforts paid off during the COVID-19 outbreak, which drove up usage requests for HPC Services and allowed the knowledge gained on cloud to be leveraged to meet the AI requirements of Dr. Yalin Zheng and his team at the Department for Eye and Vision Science.

Through a thorough exploration of GPU scaling capabilities using on-premise facilities to laying foundational rules on data requirements lead researchers Joshua Bridge and Yanda Meng were able to run a series of scaling tests on curated patient data sets in order to hone tools intended for public distribution in the rapid diagnosis of COVID-19. Along the way the team made insights into data organisation and cloud optimisation, collaborating with ARC, Alces Flight and cloud platform provider Amazon Web Services to reach their goal of achieving next-stage peer review.

II. THE INITIAL BUILD OF CLOUD BARKLA

After a two-year, £2 million refurbishment of data centre facilities at the University of Liverpool their core HPC cluster, Barkla, was set to move into its new location. In order to ensure that no core functionality was lost, Dr. Cliff Addison and Dr. Manhui Wang of the ARC team worked alongside their managed services and integration partner, Alces Flight, to replicate the physical HPC cluster using public cloud. This cloud cluster, designed to work on both Amazon Web Services (AWS) and Microsoft Azure platforms went through

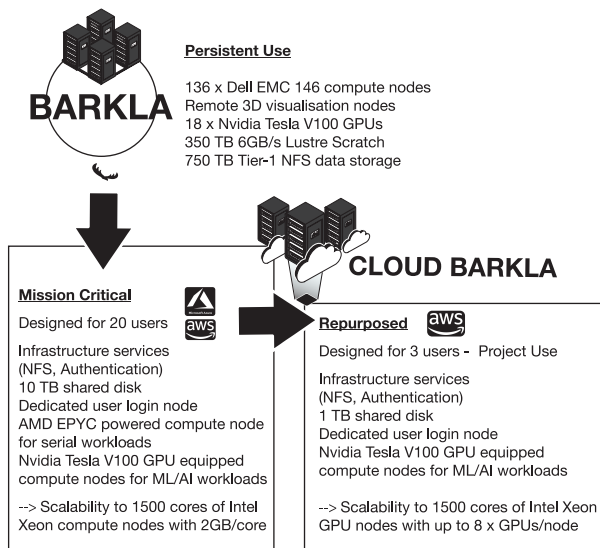


Fig. 1. The Barkla HPC Cluster and Cloud Barkla Variation

an 8-day testing and optimisation cycle prior to temporarily replacing the Barkla on-premises solution for two weeks as the physical system was moved into its new home. The transition was hugely successful, ensuring no research time was lost during the move - but also resulting in a template for future project work; something which came into play far sooner than anticipated following the outbreak of COVID-19¹.

III. RE-PURPOSING CLOUD BARKLA FOR COVID-19

Following the outbreak of COVID-19, demand quickly increased for the majority of Barkla’s resources and the capabilities of the Cloud Barkla solution were brought into play to assess incoming projects focused on research, treatment, and diagnosis. Of these, Dr. Zheng’s AI techniques for fast and accurate diagnosis of COVID-19 using automated interpretation of CT and X-ray images were identified as ideal for Cloud Barkla. Through the award of the AWS COVID-19 Global Disaster Fund for Researchers, lead researchers Joshua Bridge and Yanda Meng received USD 50,000 worth of AWS resources towards GPU-enabled research in toolsets they wished to test at scale in order to move their initial research towards full peer review. The ARC team at University of Liverpool then engaged Alces Flight to repurpose Cloud Barkla to meet the requirements set out by Dr. Zheng and his team.

IV. WORKING WITH SCALING DIAGNOSTIC MODELS

AI, especially deep learning, has shown great promise in the automatic diagnosis of a wide range of diseases. The main benefit of AI diagnosis lies in its ability to attain human-level performance while reducing strain on clinicians. The most

significant drawbacks of deep learning are the computational complexity of modern networks and the amount of data required, which requires the use of HPC facilities for training. With a time-critical problem such as COVID-19, GPUs are vital. One Nvidia V100 GPU can speed up training throughput 32 times over a CPU, meaning that models are trained in days rather than months. Training algorithms on single images, such as X-ray images, is not hugely computationally intensive; however, when volumes such as CT scans are used, the amount of computing power required increases considerably, and multiple GPUs are often needed for computation to complete promptly. GPU computing is expensive, and powering down nodes when not in use was critical in reducing the amount spent on training.

Deep learning is infamously data-hungry with hundreds or even thousands of cases needed to train. A significant amount of project time was spent ensuring that the data, from multiple sources, was in the appropriate format for training. Once the data was ready, it could be uploaded through the initial Barkla cloud setup to the AWS instance; this provided a secure and easy way to upload data, either remotely or on-site. Initially, a small dataset was trained on a single GPU, and training was then scaled up to use all the available data on 4 GPUs. Multiple GPUs allowed larger batch sizes to be used, which decreased computation time, ensuring that the project timely and useful conclusions.

Replicating the Barkla cloud meant that the computing resources of AWS were available in a format familiar to the researchers; time was not wasted learning a new system. This type of rapid repurposing substantially decreased the amount of time taken for the project to reach meaningful conclusions.

V. COST MONITORING OF GPUS

A big concern for most public cloud projects is cost of use. Through Alces Flight Center the teams rigorously tracked the rise and fall of GPU usage. By having this data to hand the team noticed points in the project in which the utilisation of Nvidia v100 quad GPUs would work best to speed time to science. This close observation allowed for the team to extend their research time by swapping the quad GPU instance type for a single GPU instance type in the initial modeling phase. They then revisited the quad GPU instance type in the final model in order to ensure that the research deadline was achieved. (See graph 2 page 3)

The utilisation of quad GPU types (which are noted in green on Fig.2 on page 3) covered the initial period of testing with an interim single GPU instance used (which is noted in blue on Fig.2 on page 3) as time to results wasn’t as imperative and allowed for a final model run utilising the quad GPUs due to cost savings.

VI. LESSONS LEARNED

This project brought into focus four key lessons learned within AI research and cloud:

¹For original Cloud Barkla build see: <https://sc19.supercomputing.org/presentation/?id=imp103sess=sess406>

Time & Cost Associated with Single and Quad GPU Usage

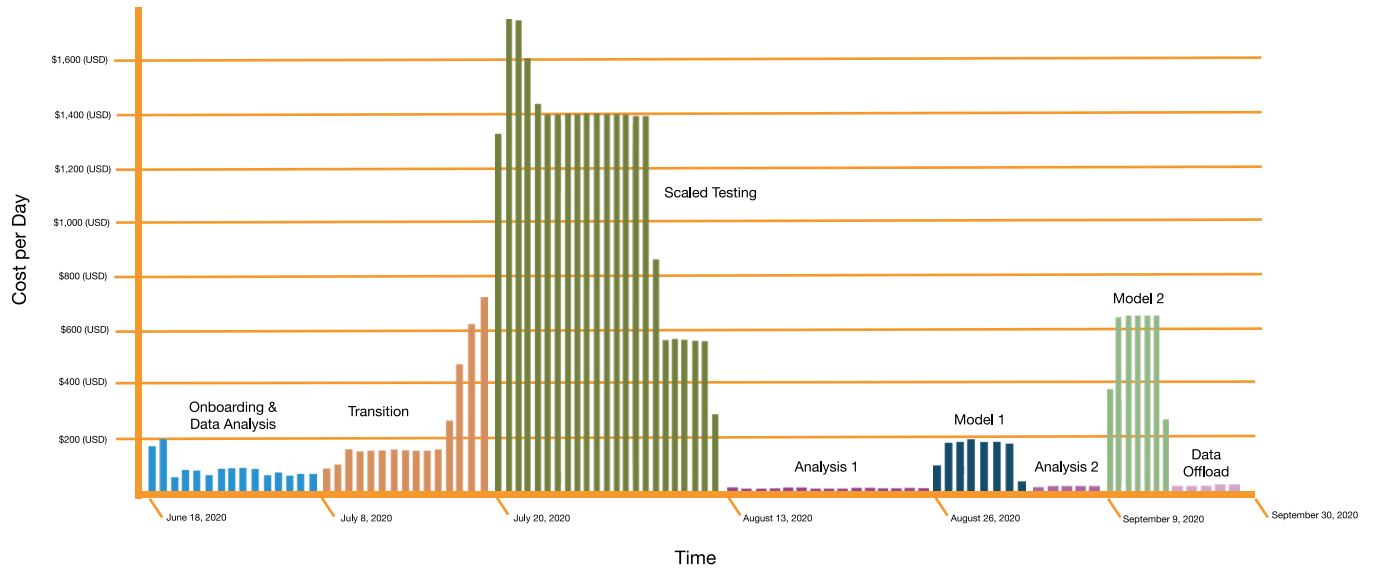


Fig. 2. Cloud Barkla GPU cost per day over time.

- Data cleansing and optimisation is key prior to commencement of a project - especially when the data is combined from multiple sources.
- Strictly monitoring GPU optimisation kept costs associated with runtimes low and offered more testing opportunity.
- Having quad GPU access for testing has led to both speed and improvement in prediction models and will allow progression of this research to continue.
- The models can be easily replicated through the the Alces Flight Center system that is managed by the ARC team to ensure proper storage and cloud scaling can be actioned quickly and efficiently.

The group was grateful for the foresight of the initial Cloud Barkla creation as it saved build time as well as additional cost and resource optimisation steps. Being able to store knowledge with Alces Flight will further allow repurposing requests for the team to revisit their work as the project reaches peer review stage.

VII. TESTING RESULTS IN BRIEF

This project is still undergoing analysis with initial reviews showing improvements in providing accurate diagnosis. The University of Liverpool team is currently working to calculate these improvements via the following equation:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

Precision is the number of COVID-19 predictions that are correct. Recall is the actual proportion of COVID-19 cases identified. As of this publication the following results from testing can be noted in brief:

$$Precision - 0.897(0.840, 0.935)$$

$$Recall - 0.873(0.812, 0.916)$$

VIII. CONCLUSION

Through foresight and appropriately aligned partners and collaborators the Cloud Barkla capability was shifted to meet the requirements of Dr. Zheng and his team as they explored AI in COVID-19 diagnosis. Initial publication for their concepts in toolsets is scheduled for end of 2020 with further peer review based on the work performed here due in early 2021.