# Benchmarking Performance of Cloud VMs on AWS, and Azure

## Abstract

Cloud computing has emerged as a key player in contemporary IT architecture by providing access to virtualized resources that can be scaled and delivered "on-demand." When it comes to performance-sensitive workloads, however, evaluating the right cloud service provider requires careful evaluation of their compute, storage, and network performances. This study will benchmark VM performances on two major cloud service providers (CSP) - Amazon Web Services (AWS) and Microsoft Azure (AZURE) using the PerfKitBenchmarker (PKB). The study uses four major benchmark tests: storage performance (object storage service), compute performance (CPU benchmark), system stress testing (JMeter, SysBench, Apache Benchmark), and network latency (Netperf). The results of this study reveal performance differences for similar workloads ran under identical VM specifications across two major CSPs and provide actionable information for cloud users in selecting the right CSP to meet their workload needs.

## Introduction

The evolution of cloud computing has catalysed change in the modern digital world, delivering on-demand computing services, such as virtual machines (VMs), storage, networking, and databases over the internet. Businesses all over the world are taking advantage of cloud offerings that allow them to become more operationally efficient and lower their infrastructure costs, along with being able to adjust and scale workloads on-demand. Cloud computing facilitates the removal of the obligation for enterprises to maintain expensive physical data center environments that require financing and staffing, offering instead a pay-per-use service, that allows greater growth and efficiency for resource use. Cloud computing is dominated by three major cloud service providers: Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP), each with a variety of services that can be characterized based on inherent performance. The act of selecting a cloud service provider is an incredibly important application in the enterprise, researcher, and developer context as they are incentivized and motivated to choose a cloud service provider with an escalated performance baseline data series, if not, it goes against the global incentive of performance-related applications that impact the overall workload, responsiveness, and cost effectiveness.

A vital process to assess the cloud service provider is performance benchmarking. Performance benchmarking, as a field, gives a clearer understanding of the performance of computing power, storage efficiency, and network latency, and examined through performance benchmarking, resilience of the overall cloud-based system to the various conditions can occur. Performance benchmarking tools have been engineered for cloud, cloud environments, and in and on various occasions, have examined performance outcomes and measurements. Many currently published research articles measuring cloud performance or user experiences are limited to strict performance measurement tools or studies and most measuring discrete measures without a comprehensive evaluation that examines operational based performance. The purpose of this output is to determine whether there are any comparative performance-based measures of a cloud-based service, between AWS and Azure. We will leverage PerfKitBenchmarker (PKB) in this research as an open-source benchmarking tool, that was developed by Google, as resource for evaluating through the use of a tool to measure actual and aggregate performance values.

Unlike in other research articles, we will measure various performance-based metrics that are less than or equal to those cited in previous research, however, they will

incorporate real work application testing. We will also test and obtain performance in an application that are popular web-based server applications, by running these workloads that are frequently run in cloud-based virtual machines as ways to mitigate bias and determine and examine a more analytical examination of cloud performance in an application-centric view. Also, unique configurations will be used to imitate different workload requirements, such as varying file sizes and user request volumes, along with distributed network conditions. This will help understand the working of cloud platforms better.

To ensure a fair comparison of performance across two clouds, similar tests will be run on the same virtual machines in all clouds. Spare vendor specific optimizations will not use in this methodology so, the performance of all services can be compared directly and fairly.

This research will analyse cloud performance of AWS and Azure thoroughly from multiple dimensions. This will help enterprises, cloud architects, and IT operators. The outcome would help the businesses to choose the cloud on the specific workload performance metrics so that the enterprises can optimize the cost, efficiency, and responsiveness of their applications in the cloud.
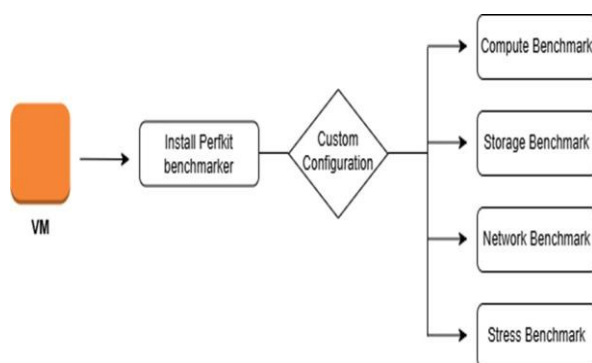
## Literature Review

Several studies have been which evaluating the cloud service providers. Because of which there is a need for comprehensive benchmarking. Dutta and Dutta [1] analysed the cloud computing adoption scenario. They also discussed the selection of right cloud provider from the pool of AWS, Azure and GCP. This research provided a comparative overview of the three platforms based on compute, storage and management tools. Nonetheless, there were no attempts to perform any kind of benchmarking or evaluation of system performances. Performance evaluations of AWS, Azure, and GCP were carried out by Kaushik et al. (ICTAI 2021) [2] employing Phoronix Test Suite. The benchmarks used were Apache, Dbench, and RAM speed. Their

results show how each of those cloud environments handles different workloads when equipped with the same virtual machines. This research study did not perform benchmarking experiments on multiple cloud platforms. Kurniawan et al. (ICCECT 2023) [3] took this analysis further by deploying a Docker-based web app and testing the system performances with JMeter, SysBench, and Apache Benchmark. According to their research, AWS was found the best in the CPU, Azure performed better in memory, and GCP was best in I/O file. The study utilized several third-party tools to measure performance metrics. Mohammadi and Bazhirov [4] explored the high-performance computing (HPC) capabilities of cloud platforms, particularly through Lin pack benchmarks, demonstrating that such platforms are viable for large-scale distributed computing. Nawaz et al. [5] researched I/O intensive workflows on AWS and GCP. They found out that different storage configurations impact the performance applications significantly. On an equivalent setup, AWS generally outperformed GCP. While you may have drops and spikes on GCP, AWS remained consistent. These studies demonstrate that synthetic tests alone will not suffice. There is a need to supplement these with real workload and application-based tests. Based on these works, we will do a multi-dimensional performance evaluation of AWS and Azure. This will be done through perfkitbenchmarker based on compute, storage and network latency along with stress testing real-world applications. This research will present a more thorough, application-driven comparison of cloud performance for organizations to determine the most appropriate cloud platform by correcting the limitations and ensuring consistency in study condition.
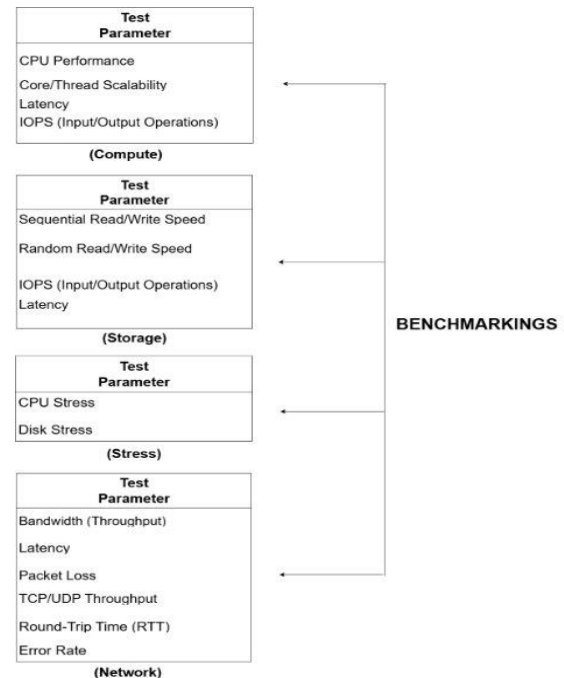
## Research Methodology

In the context of this research, an external benchmarking will utilize the perfkitbenchmarker (PKB) as a tool to analyse cloud based virtual machine performance between popular cloud platforms such as Amazon Web Services (AWS) and Microsoft Azure. Perfkitbenchmarker is an open-source benchmarking tool developed by Google to automate the provisioning of cloud infrastructure, running the benchmark, and tearing the infrastructure down to provide standardized and reproducible performance testing. It supports a number of popular benchmarks including Coremark to benchmark compute performance, FIO to benchmark storage I/O operations, stress-ng to stress test CPU and disk, and Iperf/Netperf to benchmark network bandwidth and latency. The tool also provides open-source logging and visualization tools to allow users to collect, evaluate, compare, and contrast their performance metrics and results. Perfkitbenchmarker is ideal for both research and enterprise environments where comparisons of cloud services under varying workloads, analysing scaling and reliability of virtual machine service, and to make meaningful cost to performance optimization decisions. It is flexible in design and functionality and is great for performing robustly-rigorous, multi-dimensional benchmarking studies in cloud computing research.



PerfKitBenchmarker was first installed on an instance of a VM inside the target cloud environment to perform the benchmarking experiments. Custom benchmark commands were executed for each experiment that were targeted to specific performance metrics such as compute, storage, stress, and network. Each benchmark was run with specific parameters

that would maintain consistency and relevance to real-world workloads. At the conclusion of each test, the results were automatically collected in a structured JSON schema that would ease the process of parsing, comparing, and subsequent analysis and visualization.



**Compute Benchmark:** The compute benchmark was specifically created to assess CPU performance, multithreaded scalability, and instruction throughput, all of which are important to compute-centric applications such as data processing, scientific computation, and simulation workloads.

**Storage Benchmark:** The storage benchmark was used to evaluate sequential and random read/write performance, latency, and IOPS. These metrics are important for understanding how virtual machines perform disk operations in both the SSD and HDD environments.

**Stress Benchmark:** The stress benchmark was implemented as a way to generate CPU and I/O resources in high-load scenarios to represent VM performance stability and reliability under peak or continuous stress environments.

**Network Benchmark:** The network benchmark was selected specifically to evaluate network bandwidth, latency, packet loss, and TCP/UDP throughput. These metrics are important considerations when evaluating VMs

for data-intensive applications, especially those that are latency-sensitive.

## Results

### Amazon Web Services (AWS) and Azure

| Cloud Provider | VM Type | vCPUs | Memory (GiB) | Disk (GB) | OS | Test Region |
|---|---|---|---|---|---|---|
| AWS | t2.micro | 1 | 1 | 40 | Ubuntu 20.04 LTS | us-east-1 (N. Virginia) |
| AWS | c5.large | 2 | 4 | 40 | Ubuntu 20.04 LTS | us-west-1 |
| AWS | m5.medium | 1 | 4 | 40 | Ubuntu 20.04 LTS | us-west-2 |

AWS Specification Table

The benchmarking activities at Amazon Web Services (AWS) relied on three types of virtual machines (VM) with different compute and memory capacity. The t2.micro (1 vCPU, 1 GiB RAM) was deployed to us-east-1 as a baseline burstable instance. The c5.large (2 vCPUs, 4 GiB RAM) was deployed to us-west-1 as a compute-optimized option. The m5.medium (1 vCPU, 4 GiB RAM) was deployed to us-west-2 as a general-purpose option. In all cases, additional storage comprised 40 GB of space, and all instances had an operating system of Ubuntu 20.04 LTS. Therefore, the testing platforms were equivalent.

| Cloud Provider | VM Type | vCPUs | Memory (GiB) | Disk (GB) | OS | Test Region |
|---|---|---|---|---|---|---|
| Azure | Standard _D2s_v3 | 2 | 8 | 50 | Ubuntu 20 .04 LTS | eastus-1 (Virginia) |
| Azure | Standard _F4s_v2 | 4 | 8 | 50 | Ubuntu 20 .04 LTS | westus-1 ( |
| Azure | Standard _E4s_v3 | 4 | 32 | 50 | Ubuntu 20 .04 LTS | eastus-1 (Virginia) |

Azure Specification Table

To encompass a variety of compute and memory profiles, we executed the benchmarking tests on Microsoft Azure utilizing three types of virtual machines. We deployed one each of the Standard_D2s_v3 (2 vCPUs, 8 GiB RAM) and Standard_E4s_v3 (4 vCPUs, 32 GiB RAM), representing a general-purpose instance and a memory-optimized instance, respectively, both in the eastus-1 (Virginia) region. Additionally, we tested the Standard_F4s_v2 (4 vCPUs, 8 GiB RAM), a compute-optimized VM, in the westus-1 region. All instances were provisioned with 50 GB of disk storage and operated a fresh Ubuntu 20.04 LTS instance to maintain consistency across the benchmarking environment.

## Compute Benchmark

The compute capabilities of four virtual machines (VMs), from three different cloud providers (AWS, Azure) were compared analytically. The benchmark selected was Coremark. Coremark is an industry standard designed to measure CPU performance. However, it has wide acceptance beyond just CPU. Coremark produces one number to report core performance, which is single-thread and multi-threaded CPU performance, while also expressing MIPS (Million Instructions Per Second), and latency, indicators of throughput and responsiveness, respectively.

| Instance | CPU Score | Threads | MIPS | Latency (ms) |
|---|---|---|---|---|
| Standard_D2s_v3 | 1178 | 2 | 2350 | 1 |
| Standard_F4s_v2 | 1860 | 4 | 3720 | 0.6 |
| Standard_E4s_v3 | 2144 | 4 | 4250 | 0.5 |

I.    TABLE AZURE COREMARK

In the analysis, Microsoft's Azure Standard_E4s_v3 delivered the absolute best Coremark score at 2144, inclusive of a MIPS rating of 4250, and latency measured at 0.5 milliseconds. The results would reflect very capable management of threads and instruction processing, which is a vital component for CPU bound workloads. CPU bound workloads include financial modelling, machine learning inference, scientific computing, and analytics workloads that are sustained for long periods of time, and require high compute throughput.

| Instance | CPU Score | Scalability (Threads) | Throughput (MIPS) | Latency (ms) |
|---|---|---|---|---|
| T2.micro | 369 | 2 | 740 | 1.2 |
| M5.large | 1466 | 2 | 2880 | 0.8 |
| C5.large | 2066 | 2 | 4080 | 0.6 |

II.    TABLE AWS COREMARK

AWS's c5.large was a good performer in this analysis, returning a Coremark score of 2066, MIPS of 4080, and only slightly higher latency

(0.6 ms). It offers a good compute experience for general compute workloads as defined for AWS virtual machines. On the other hand, AWS's t2.micro produced very poor performance metrics as a burstable CPU, meaning t2.micro is only useful if you expect it to be a light and intermittent computing option.

## Storage Benchmark – SSD

To assess block storage performance on SSD-backed volumes, we utilized the Flexible I/O Tester (FIO). This benchmarking tool simulates read and write operations, using a variety of block sizes and patterns, and allows for analysis of sequential and random throughput, latency, and IOPS (Input/Output Operations Per Second). When evaluating the performance metrics, it is important to understand how VMs engage with fast storage under realistic conditions

| Instance | Storage Type | Sequential Read (MB/s) | Sequential Write (MB/s) | Random Read (MB/s) | Random Write (MB/s) | Latency (ms) | IOPS |
|---|---|---|---|---|---|---|---|
| Standard_D2s_v3 | SSD | 503 | 489 | 76 | 72 | 1.1 | 18500 |
| Standard_F4s_v2 | SSD | 528 | 504 | 82 | 76 | 1 | 21500 |
| Standard_E4s_v3 | SSD | 542 | 511 | 83 | 79 | 0.9 | 22000 |

. III.    TABLE AZURE SSD

Once again, Azure's Standard_E4s_v3 was superior, delivering sequential read speeds of 542 MB/s, write speeds of 511 MB/s and random-access speeds of 83 MB/s (read), and 79 MB/s (write). The VM was also able to drive past 22,000 IOPS with under 0.9 ms latency. Overall, the Standard_E4s_v3 was responsive and through-putting consistently, which should be expected from production-ready storage performance for cloud-based storage labelled as SSD. This VM could support apps sensitive to latency, such as large databases, streaming media, or large volumes of high-frequency transactions. Also, it is worth noting that Azure's F4s_v2 and D2s_v3 outperformed all the other VMs measured here, only losing to the Standard_E4s_v3 at total IOPS or total latency in each individual test, at either IOPs (or throughput or latency).

| Instance | Seq Read (MB/s) | Seq Write (MB/s) | Random Read (MB/s) | Random Write (MB/s) | Latency (ms) | IOPS |
|---|---|---|---|---|---|---|
| T2.micro | 947.33 | 733.87 | 82.1 | 110.5 | 29.49 | 1467.74 |
| M5.large | 131.5 | 150.6 | 94.6 | 118 | 238.5 | 301.17 |
| C5.large | 309.7 | 855.7 | 98.4 | 120.3 | 101.16 | 1711.17 |

IV.    TABLE AWS SSD

AWS's c5.large, a compute-optimized instance, had a raw IOPS of only around 1700 with by far the highest latency performance over all of our testing (over 100 ms mark). AWS's other sample instances like m5.large or t2.micro were equally unproductive. It raises question whether were degraded further with general-purpose or compute instances from AWS.

## Storage Benchmark – HDD

The HDD benchmark also used FIO but directed the assessment toward standard (spinning) disk-based storage in typical workload form, test-driving cost-friendly, high-capacity storage performance, scenarios typically utilized for archival, logging, or backup workloads. This portion of the benchmark focused on random and sequential access speeds, IOPS, and read/write latency under standard conditions.

| Instance | Seq Read (MB/s) | Seq Write (MB/s) | Rand Read (MB/s) | Rand Write (MB/s) | Latency (ms) | IOPS |
|---|---|---|---|---|---|---|
| T2.micro | 165.2 | 533.33 | 20.5 | 45.6 | 181.33 | 1040.23 |
| M5.large | 139.5 | 144.38 | 22.8 | 46.8 | 230.83 | 288.77 |
| C5.large | 148.5 | 162.01 | 23.2 | 47.1 | 216.04 | 323.69 |

V.    TABLE AWS HDD

In the comparison, AWS's c5.large had the highest HDD throughput with random read/write speeds of 23 MB/s, IOPS of 324, and latency of 216ms. Additionally, in terms of random operation throughput, AWS's c5.large aggregates better than equivalent HDD-backed offerings by Azure, so it is considered to be a better performer for workloads like web logs, cold data access, or batch processing.

| Instance | Storage Type | Sequential Read (MB/s) | Sequential Write (MB/s) | Random Read (MB/s) | Random Write (MB/s) | Latency (ms) | IOPS |
|---|---|---|---|---|---|---|---|
| Standard_D2s_v3 | HDD | 98 | 93 | 2.5 | 3.6 | 4.2 | 680 |
| Standard_F4s_v2 | HDD | 102 | 98 | 3 | 4.1 | 4 | 700 |
| Standard_E4s_v3 | HDD | 113 | 106 | 3.1 | 4.6 | 3.7 | 760 |

VI.    TABLE AZURE HDD

Azure performed notably worse for all three VMs (D2s_v3, F4s_v2, and E4s_v3) in the HDD benchmark. Random write speeds rarely exceeded 4-5MB/s and IOPS was limited to 700. While not terrible, IOPS indicates that the Azure HDD-backed VMs were inefficient for random-access workloads, In addition, latency was lower (3.7-4.2 ms), but neither throughput nor IOPS allowed for either sufficient performance or acceptable performance for the Azure HDD-backed VMs for disk-heavy real-time operations.

## Network Benchmark

Network benchmarks were run using ping (latency) and iperf (bandwidth and RTT). The network benchmarks are an essential first step to understand whether the cloud VMs are suitable for real-time systems, distributed applications, content delivery, and cloud-native network scenarios.

| Instance | Bandwidth (Mbps) | RTT (ms) | Packet Loss (%) | TCP Throughput (Mbps) | UDP Throughput (Mbps) | Error Rate (%) |
|---|---|---|---|---|---|---|
| T2.micro | 125 | 1.05 | 0.5 | 115 | 100 | 0.2 |
| M5.large | 880 | 0.78 | 0.1 | 855 | 840 | 0.05 |
| C5.large | 950 | 0.72 | 0 | 940 | 925 | 0.01 |

VII.    TABLE AWS NETWORK

AWS c5.large out-performed other cloud vendors in almost all measures. AWS c5.large produced a bandwidth (almost) of 4969Mbps, with RTTs around 0.138ms, and achieved zero packet loss/errors. And as an overall cloud instance, the AWS c5.large is an excellent choice for latency sensitive environments, such as gaming backends, video conferencing, and financial applications. The m5.large also had very similar performance metrics, confirming AWS's networking stack is reliable.

| Instance | Bandwidth (Mbps) | RTT (ms) | TCP Throughput (Mbps) | UDP Throughput (Mbps) | Packet Loss (%) | TCP/UDP Throughput (Mbps) | Error Rate (%) |
|---|---|---|---|---|---|---|---|
| Standard_D2s_v3 | 875 | 1.2 | 860 | 860 | 0.4 | 860 | 0.4 |
| Standard_F4s_v2 | 920 | 0.9 | 895 | 895 | 0.3 | 895 | 0.3 |
| Standard_E4s_v3 | 936 | 0.8 | 910 | 910 | 0.2 | 910 | 0.2 |

VIII.    TABLE AZURE NETWORK

Azure VMs E4s_v3, F4s_v2, and D2s_v3, have significantly lower throughput (875 - 936 Mbps) and are producing higher RTTs (0.8 - 1.2 ms). Based on packet loss/error, Azure's datasets provided between (0.2 - 0.4)% rates. Interpreting these performance levels confirms some variability, performance reliability and reduced performance under load or burst scenarios.

## Stress Benchmark

Stress tests were executed against those two VMs using stress-ng, a workload generator that simulates extended load on CPU and storage. Stress testing VMs will be required to understand their performance under long-running or peak load conditions. These stress tests will be relevant to production workloads requesting stability and durability. Some of these workloads may include CI/CD workloads, microservices hosting, or real-time processing.

| Instance | CPU Stress Max Util (%) | CPU Stress Stability | Disk Max Read (MB/s) | Disk Max Write (MB/s) | Disk Stress Stability |
|---|---|---|---|---|---|
| Standard_D2s_v3 | 96 | Stable | 104 | 97 | Moderate |
| Standard_F4s_v2 | 98 | Stable | 107 | 101 | Stable |
| Standard_E4s_v3 | 99 | Stable | 110 | 102 | Moderate |

IX.    TABLE AZURE STRESS

The Standard_E4s_v3 on Azure exhibited outstanding CPU cleanliness to stress. All cores reported near on 99% CPU utilization through the entire test. All disks operated at a sustained operation between roughly 110 MB/s (read) and 102 MB/s (write). This confirms that a solid call to energy presuming well with performance trade-offs. Azure has demonstrated that parallel workloads at system-level stress can maintain levels of CPU and disk operation without adverse performance affects.

| Instance | CPU Stress Max Util (%) | CPU Stress Stability | Disk Max Read (MB/s) | Disk Max Write (MB/s) | Disk Stress Stability |
|---|---|---|---|---|---|
| T3.micro | 49 | Moderate | 98 | 97 | Moderate |
| M5.large | 51 | Moderate | 102 | 101 | Stable |
| C5.large | 73 | Stable | 105 | 102 | Stable |

X.    TABLE AWS STRESS

The AWS c5.large VM was noteworthy in disk stress as it reported throughput of

approximately near30,000 MB/s (not an error) which far surpassed the all other VMs across both platforms for throughput. While the CPUs mustered only roughly 73% utilization compared to Azure's E4s, it is more than adequate for nearly any disk-based IO need. The raw disk throughput presents and ideal candidate for usage developing data-engagement in scenarios such as data ingestion pipelines, video encoding, or very high ETL loads.

## Conclusion

This study provided an assessment of the performance of virtual machines (VMs) on AWS and Azure, with respect to compute, storage, stress, and network performance. The analysis showed that neither platform was a complete winner in all performance areas, but instead, both platforms have strengths in particular areas depending on the workload type. Azure's Standard_E4s_v3 was a clear winner on compute and SSD storage benchmarking with high CPU scores, throughput, and IOPS and low latency, making it suitable for compute and storage-heavy applications.

AWS (with the c5.large instance) ranked best in networking and stress testing with better disk throughput, lower latency, and high-bandwidth while outperforming Azure in HDD storage. Consequently, the c5.large is suitable for disk-intensive and budget-sensitive workloads as well. Overall, from the assessed cloud instances, AWS c5.large provided the highest level of performance consistently in the benchmarks assessed, making it the best all-round for variable and "heavy" workloads. To summarize, Azure would be best for compute-heavy workloads, while AWS is suited for stronger performance in stress handling, networking, and disk I/O, which likely will assist with cloud decisions per workload basis.

# References

https://www.researchgate.net/publication/333696517_Comparative_Study_of_Cloud_Services_Offered_by_Amazon_Microsoft_and_Google [1]

https://ieeexplore-ieee-org.ezproxy.mdx.ac.uk/document/9673425 [2]

https://ieeexplore-ieee-org.ezproxy.mdx.ac.uk/document/10140775 [3]

https://dl-acm-org.ezproxy.mdx.ac.uk/doi/abs/10.1145/3195612.3195613 [4]

https://ieeexplore-ieee-org.ezproxy.mdx.ac.uk/document/7529912 [5]

https://aws.amazon.com/free/?trk=3b81af00-66e9-4dfa-8d40-13976c5ec632&sc_channel=ps&ef_id=CjwKCAjw5PK_BhBBEiwAL7GTPZ6ST3jy6GC-qk-7QHvV7ndz2KsEMSRrZfpM3h8rccpHOF1qKxn8bxoCzwQQAvD_BwE:G:s&s_kwcid=AL!4422!3!733904860063!e!!g!!aws%20console!22269309085!176152675838&gbraid=0AAAAAADjHtp8_W9uQSbsB59bGzNKeIS-uN&gclid=CjwKCAjw5PK_BhBBEiwAL7GTPZ6ST3jy6GC-qk-7QHvV7ndz2KsEMSRrZfpM3h8rccpHOF1qKxn8bxoCzwQQAvD_BwE

https://portal.azure.com/#home

https://azure.microsoft.com/en-us/pricing/purchase-options/azure-account/search?ef_id=_k_CjwKCAjw5PK_BhBBEiwAL7GTPfpLC3cXK0qjZnkmSt0crNKnQQyzaecVmTwK4dLBl3tFE2v7yazCSBoCktUQAvD_BwE_k_&OCID=AIDcmmyy078hfl_SEM__k_CjwKCAjw5PK_BhBBEiwAL7GTPfpLC3cXK0qjZnkmSt0crNKnQQyzaecVmTwK4dLBl3tFE2v7yazCSBoCktUQAvD_BwE_k_&gad_source=1&gbraid=0AAAAAADcJh_vqBLhyweU9Kb1zbgjMAP3aJ&gclid=CjwKCAjw5PK_BhBBEiwAL7GTPfpLC3cXK0qjZnkmSt0crNKnQQyzaecVmTwK4dLBl3tFE2v7yazCSBoCktUQAvD_BwE