TECHILA TECHNOLOGIES LTD

BENCHMARK REPORT

DEPLOYING 40,000 SPOT vCPUs

IN GOOGLE CLOUD

## Background and Motivation

In September 2022, Google Cloud published a joint customer study with the University of Pittsburgh (UPitt) and Techila Technologies. The case study is about the willingness of UPitt's central IT department to offer computing services to researchers outside those used by the traditional high-performance computing sciences. The shared high-performance computing resources of UPitt continuously experience very high utilization, so it was an excellent opportunity to demonstrate the economics and capacity of cloud computing for research projects that can't wait. UPitt IT also wanted to encourage research that transcends disciplines and impacts the whole community. As stated in the story, it's all about "speed to science" -- demonstrating the possibilities of cloud computing to provide immediate access to large-scale compute resources that massively reduce the time to results.

As a result, it was demonstrated that using 40,000 CPUs on Google Cloud instead of the local 673 CPU on-premises cluster, Techila Distributed Computing Engine (DCE) enabled a UPitt researcher to reduce the execution time of a seemingly impossible MATLAB® simulation **from 6 months to just 48 hours**.

Furthermore, being able to use preemptible compute instances was also far more cost-effective: UPitt IT estimated that the project cost less than a third of what it would have cost on-premises and around 80% cheaper than using on-demand cloud computing instances. More details can be found in:

[https://edu.google.com/intl/ALL_us/why-google/customer-stories/university-of-pittsburgh/](https://edu.google.com/intl/ALL_us/why-google/customer-stories/university-of-pittsburgh/)

## Introduction

As UPitt IT stated in the above customer study, for its' researchers, it's all about the time-to-science, just as in industry, it's all about the time-to-market. Researchers are often under extreme pressure to complete their research within a given period because of the deadline for publishing their results. The need to deliver high-quality research and the tough competition in science easily lead researchers to desire access to 100s or even 1000s of times more computing power than existing available resources. It shows a need for an economical and just-in-time solution that provides another option for researchers and scientists.

To support the notion of having 40,000 CPUs instantly available, Techila Technologies engineers ran several large-scale benchmarks in Google Cloud to determine how quickly the compute capacity can be deployed and made available to run workloads. This test report aims to provide performance statistics to demonstrate the practical aspects of deploying and using a large cloud capacity. Tests were done using a similar 40,000 vCPU Techila DCE environment, as mentioned in the UPitt customer study.

## The Test Plan Setup

To evaluate the economics and potential of using the Google Cloud Platform (GCP) for high-throughput computing, it was essential to design the tests accordingly. Like UPitt, Techila chose to perform tests using preemptible Spot instances available at a fraction of the cost of the regular high Service Level Agreement (SLA) compute instances.

---

**What are Spot instances?**
Cloud providers must have spare capacity available for any surge in customer demand. Spot instances represent the cloud providers' excess capacity. To offset the loss of idle infrastructure, cloud providers offer this excess capacity at a massive discount to drive usage. This discounted Spot instance pricing comes with a caveat. Cloud providers can "pull the plug" and terminate Spot instances with a 30-second warning. These interruptions occur when cloud providers need to draw from the excess capacity to service customers who purchased on-demand instances. Spot instances are not covered by any SLA like on-demand instances.

---

More in-depth information about Google Spot instances can be found at: https://cloud.google.com/spot-vms

Techila DCE includes an automatic and efficient ability to recover from compute infrastructure failures in microseconds. This capability allows customers to benefit from these highly cost-effective Spot instances.

There are two operating systems available in GCP: Linux and Windows. Using the Windows operating system in the cloud causes additional license costs. Considering this, it was economical to choose a free Linux operating system for the tests, as UPitt did.

Economics is one side of the story; another is the just-in-time aspect. Nobody likes waiting, especially if there is a rush due to schedule pressures or, for example, a great need for quick prototyping. Waiting for your turn in an on-premises cluster batch queue does not support the intuitive way of working that people have used to when working with their laptops. Likewise, waiting for cloud computing capacity to be ready for computing does not support the target of immediate access to computing power which is critical enablement for time-to-science.

**Techila Technologies engineers performed several large-scale benchmarks with Techila Distributed Computing Engine in Google Cloud to find answers the following questions:**
- What time is required to deploy and make 40,000 vCPUs (virtual CPU) available?
- How much time does it take to shut down the deployed capacity?
- How high is the utilization rate that Techila DCE can provide when processing High Throughput Computing (HTC) workloads using the 40,000 vCPU capacity?
- What is the price of computing, on-demand vs. Spot?

## How quickly can 40,000 vCPUs be deployed and shut down?

Techila Technologies performed several tests where a capacity of 40,000 Spot vCPUs was deployed to answer this question. Tests were performed using either *n1-standard-32* or *n1-standard-64* instances in Google Cloud *europe-west1* and *us-central1* regions. In the tests, only one region and one zone were used at the time. All tests were performed using Techila DCE Advanced Edition, publicly available as a pay-as-you-go solution in Google Marketplace.

The results exceeded expectations. First, it should have been asked if it is even realistic to get 40,000 Spot instances up and running. Second, the speed with which this capacity was brought into computational condition was a surprise. On average, the entire Techila Worker capacity of 40,000 vCPUs was online and executing workloads in **under 100 seconds** from hitting the start button. This included the time required to provision the instances and spin up the Techila Worker processes to establish a working cluster. Third, when the test workload was complete, shutting down the total worker capacity was completed in just **under 30 seconds**.
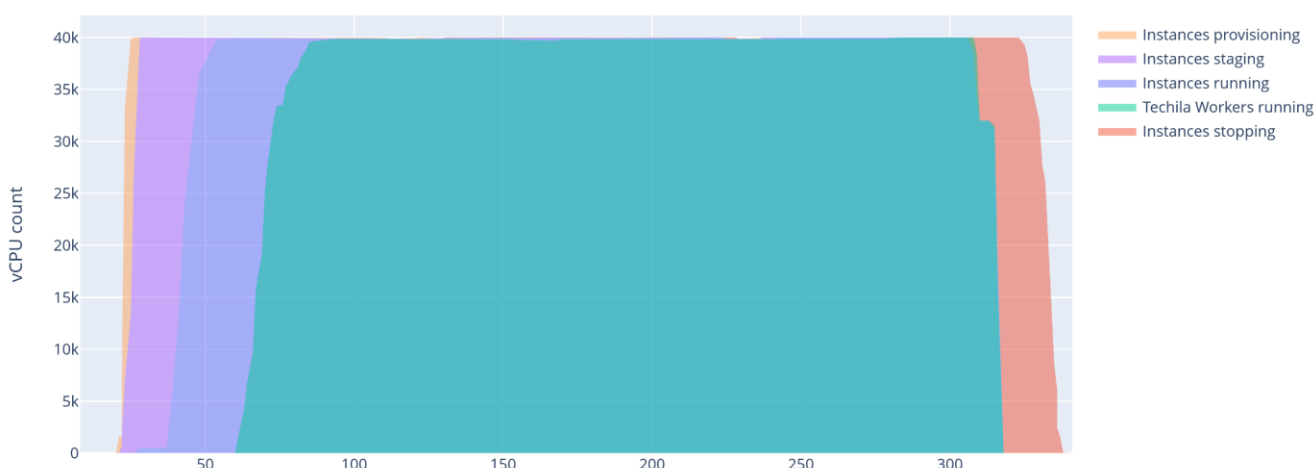


*Figure 1: Deploying a 40,000 vCPU Techila Distributed Computing Environment in Google Cloud. Test performed in us-central1 using n1-standard-64 instances on September 1, 2022, at 12:00 CEST.*

| Instances stage | Explanation of the stage |
|---|---|
| Instances provisioning | Resources are allocated for the virtual machine (VM). The VM is not running yet. |
| Instances staging | Resources are acquired, and the VM is preparing for first boot |
| Instances running | The VM is booting up or running. |
| Techila Workers running | Techila Worker is ready for computing |
| Instances stopping | The VM is being stopped after which the VM enters the terminated status. |

More information about the instance lifecycles can be found in the Google Cloud documentation at:

https://cloud.google.com/compute/docs/instances/instance-life-cycle

## How efficiently can Techila DCE utilize the 40,000 Spot vCPU capacity?

The proprietary scheduler included in Techila Distributed Computing Engine ensures that the maximum amount of available CPU power can be utilized with minimal scheduler overhead. Previously released benchmarks in **2021** and **2022** show that the system utilization rate provided by the Techila DCE scheduler can be several times greater than traditional HPC solutions. This high system utilization rate reduces overall runtime compared to open source and batch-type schedulers, which in turn means **lower infrastructure costs** when using pay-as-you-go cloud capacity.

The figure below illustrates the system utilization rate of 40,000 Spot vCPUs when processing very high throughput workloads consisting of 400,000 short 20-40 second sleep tasks. The main driver for testing with relatively short sleep tasks was focusing on the scheduler's performance and verifying that the capacity remains fully utilized even with the most demanding workloads. It is also notable that with Spot instances, it is wise to use as short tasks as possible because of the possible sudden terminations of Spot instances.
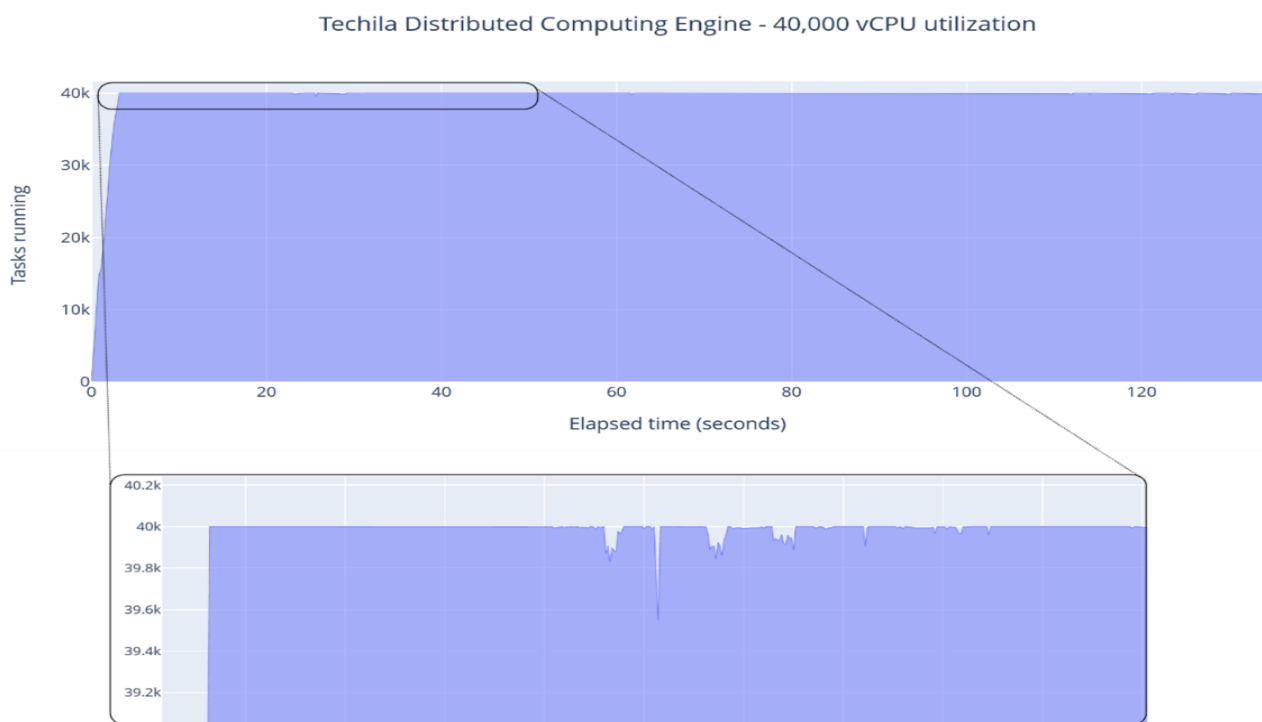


*Figure 2: Number of tasks running in 40,000 vCPU Techila DCE environment. The system consisted of 1250 n1-standard-32 instances running in Google Cloud. Zoomed area shows minimal fluctuations in the running tasks count.*

When starting computations, Techila DCE could allocate a task for each of the 40,000 vCPUs in just 3.20 seconds. While processing the tasks, Techila DCE was able to maintain the average **system utilization rate at 99.977%.** The system was processing an average of 39,990 computational tasks at any given time, and Techila DCE scheduling overhead was just 0.023%.

## How much does an hour of 40,000 Spot vCPUs cost?

The general perception of Spot instances is that they are not ideal for high-performance computing. But the built-in and automatic fault-tolerant features of Techila DCE allow customers to fully utilize also Spot instances. The maximum utilization enabled by Techila DCE combined with the significantly lower cost of Spot instances means that even massive computing power won't break the bank. In our tests, using Spot instances **reduced the cloud infrastructure costs, in this case, 78.95%.**

|  | On-Demand | Spot |
|---|---|---|
| Instance type | n1-standard-32 | n1-standard-32 |
| Price / instance / h | $1.519992 | $0.320000 |
| Price / vCPU /h | $0.047410 | $0.010000 |
| Instance count | 1250 | 1250 |
| Total vCPU count | 40,000 | 40,000 |
| **Total cost / h** | **$1899.90** | **$400.00** |

*Table 1: Example calculation for cloud instance costs. Costs were calculated using us-central1 prices on August 28, 2022.*

## Conclusions

As outlined in the introduction, these tests were performed to obtain quantitative results about how quickly Techila DCE can be deployed at scale and how efficiently performs when running large-scale computations in Google Cloud.

Taking into consideration the temporal nature of cloud computing and the differences in data center loads at different times, it is impossible for cloud providers to guarantee Spot capacity availability at any given time. However, as illustrated by the deployment test results shown in Figure 1, Techila DCE was able to provision 40,000 Spot vCPUs reliably and incredibly fast. Furthermore, the scheduler performance statistics illustrated in Figure 2 show that the proprietary scheduler of Techila DCE can efficiently utilize this 40,000 vCPU capacity even in demanding HTC workloads.

All tests were performed using Techila Distributed Computing Engine Advanced Edition, publicly available as a pay-as-you-go solution in Google Marketplace. Please contact your local Techila representative to find out more or arrange a demonstration.

If you can't wait you can try out the Techila DCE now by visiting the link below:
**https://console.cloud.google.com/marketplace/product/techila-public/techila**

## Appendix

More information about Techila DCE can be found in the Introduction to Techila Distributed Computing Engine document. Additional details about the scheduler functionality can be found in the Techila DCE Scheduler document. Below are links to short demo videos explaining how Techila DCE integrates with different programming languages.

**Techila DCE with Python :**

https://www.youtube.com/watch?v=5zALrtFTeso

**Techila DCE with Jupyter Notebook :**

https://www.youtube.com/watch?v=WsCkjmTIfko

**Techila DCE with R :**

https://www.youtube.com/watch?v=D1xk4_VUCtI

**Techila DCE with MATLAB® :**

https://www.youtube.com/watch?v=LLMXV3o2FT0

**Techila DCE with Command Line API :**

https://www.youtube.com/watch?v=X42jVetFY8g

## Disclaimer and parting words

This benchmark was done independently by Techila Technologies. Google Cloud did not sponsor or contribute financially to this benchmark report.

Benchmarking studies conducted by companies about their product are understandably difficult to take at face value, and their scientific value may well be disputed. If you are looking to improve the utilization and throughput of your computing environment but are hesitant about the results or methods presented here, **please contact sales@techilatechnologies.com to schedule your own tests.**