

Google Cloud computing
at scale.

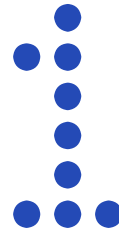
*How Punch scaled test
simulations to hundreds of
machines.*

numin

punch case study

Contents

2	Contents
3	About Numin
5	A solution in need of scale
7	The fix for Numin
18	Results



About Numin

NUMIN IS A LEADING QUANT FUND DEDICATED TO ENGINEERING AND TRADING THE LATEST STATISTICAL AND ARTIFICIAL INTELLIGENCE DRIVEN STRATEGIES IN PUBLIC FINANCIAL MARKETS.

FINANCE MEETS CLOUD

Numin's sophisticated algorithmic backtesting systems required a level of scale and distributed cloud computing that wasn't achievable through a manual testing process.

Punch worked as Numin's platform engineering solutions provider to help Numin dramatically expand its cloud infrastructure on demand. We focused on performance, cost, and ease of use in architecting a robust Google Cloud Platform solution.

The result is a fault tolerant backtesting system that is incredibly low cost, highly sophisticated, and easy to use.

PUNCH SERVICES PROVIDED

Punch provided expertise in engineering and platform engineering to help Numin meet deadlines and goals for rapid development.

Google Cloud Platform
engineering,
Platform engineering,
Developer Operations,
Site Reliability, QA,
Project Management



A solution in need of scale

NUMIN WAS FACING A PROBLEM. THEIR DATA SCIENTISTS AND ENGINEERS HAD DEVELOPED A SOPHISTICATED BACKTESTING AND ARTIFICIAL INTELLIGENCE TRAINING SYSTEM: THE PROBLEM WAS, HOW TO SCALE?

The number of trials required to train a population of agents and achieve model convergence was well beyond what any local hardware setup could achieve. And the data would be silo'd.

- 1 Data Silos.** Each test was silo'd to each engineers machine. While the backtesting code was in the cloud, the implementation and results of the tests driven by that code were not.
- 2 Insufficient hardware.** The scale of machinery required to run hundreds of parallel trials was too much for the team's computers to achieve on any reasonable time horizon. The team looked into building their own local data center, but that carried with it a whole new set of challenges, such as further scaling.
- 3 Human error.** Translation of each test trial to the master record was tedious and laborious. Mistranslation was a common problem that could result in faulty decisions by management.

800

NUMBER OF VIRTUALIZED CPUS

The system needed to accommodate hundreds of CPUs to perform all the mathematics required for accurate backtests.



The fix for Numin

PUNCH WORKED WITH NUMIN TO SOLVE BACKTESTING ISSUES THROUGH A DYNAMIC RESOURCE ALLOCATION MODEL.

Punch worked closely with Numin and their team to break-apart the problem into several steps, resulting in a streamlined, first-principles approach to the problem.

PROBLEM PROCESS

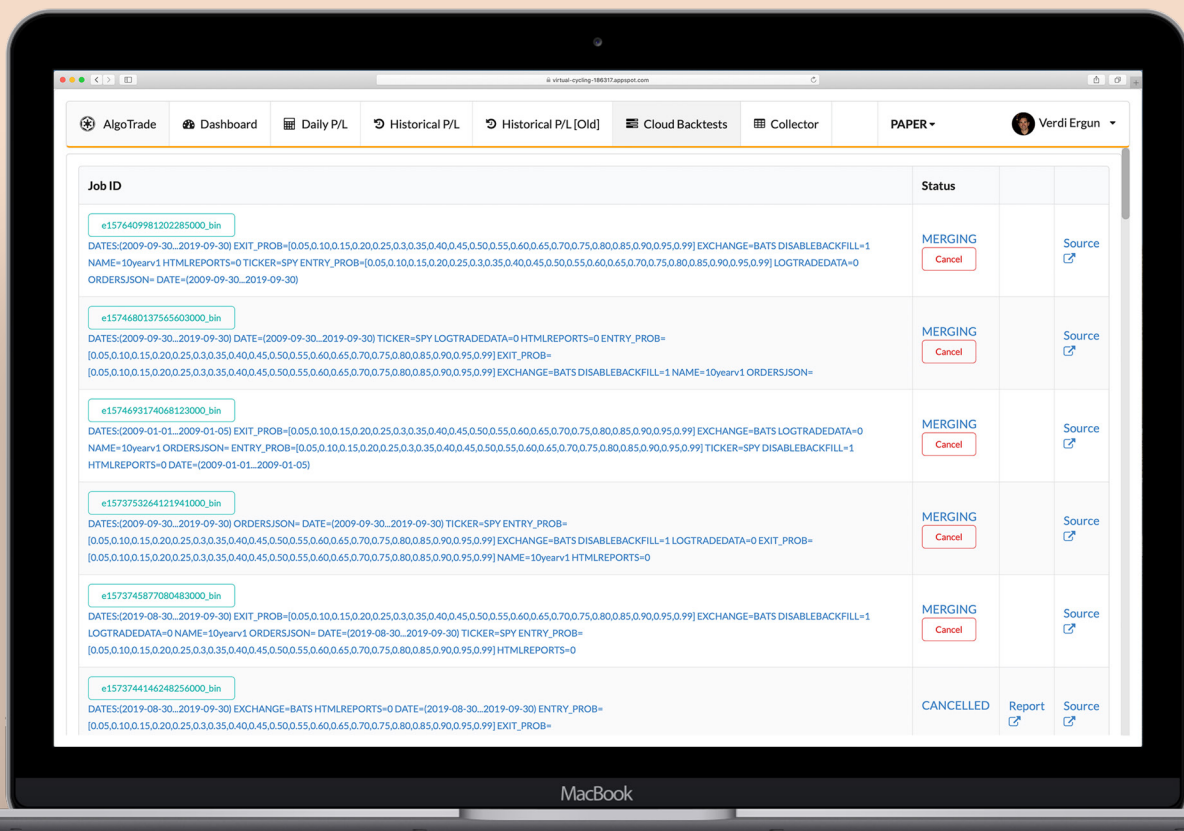
Demand	Engineers from Numin would queue a backtest job to the cloud.
Supply	A series of scripts would begin allocating resources to the job.
Scaling	The CPU utilization of each resource was monitored, as the CPU utilization grew beyond certain thresholds, more virtualized machines were dynamically spawned. Each virtualized machine would pick up jobs from the stack when its job was complete. Machines were fault-tolerant.
Descaling	As the tests finished, machines would be automatically despawned to minimize cost.
Collection	Results from each machine would be reported back in real-time to a centralized machine that would combine the results into a single results graph.

1

DEMAND

SPAWN A CLOUD JOB FROM ANYWHERE IN THE WORLD.

Numin engineers could spawn a cloud job from their local terminal using Google Auth and Google Cloud SDK. A single cloud machine would begin the test and self-monitor for CPU utilization. A separate reporting dashboard would show the progress of each individualized backtest over time.

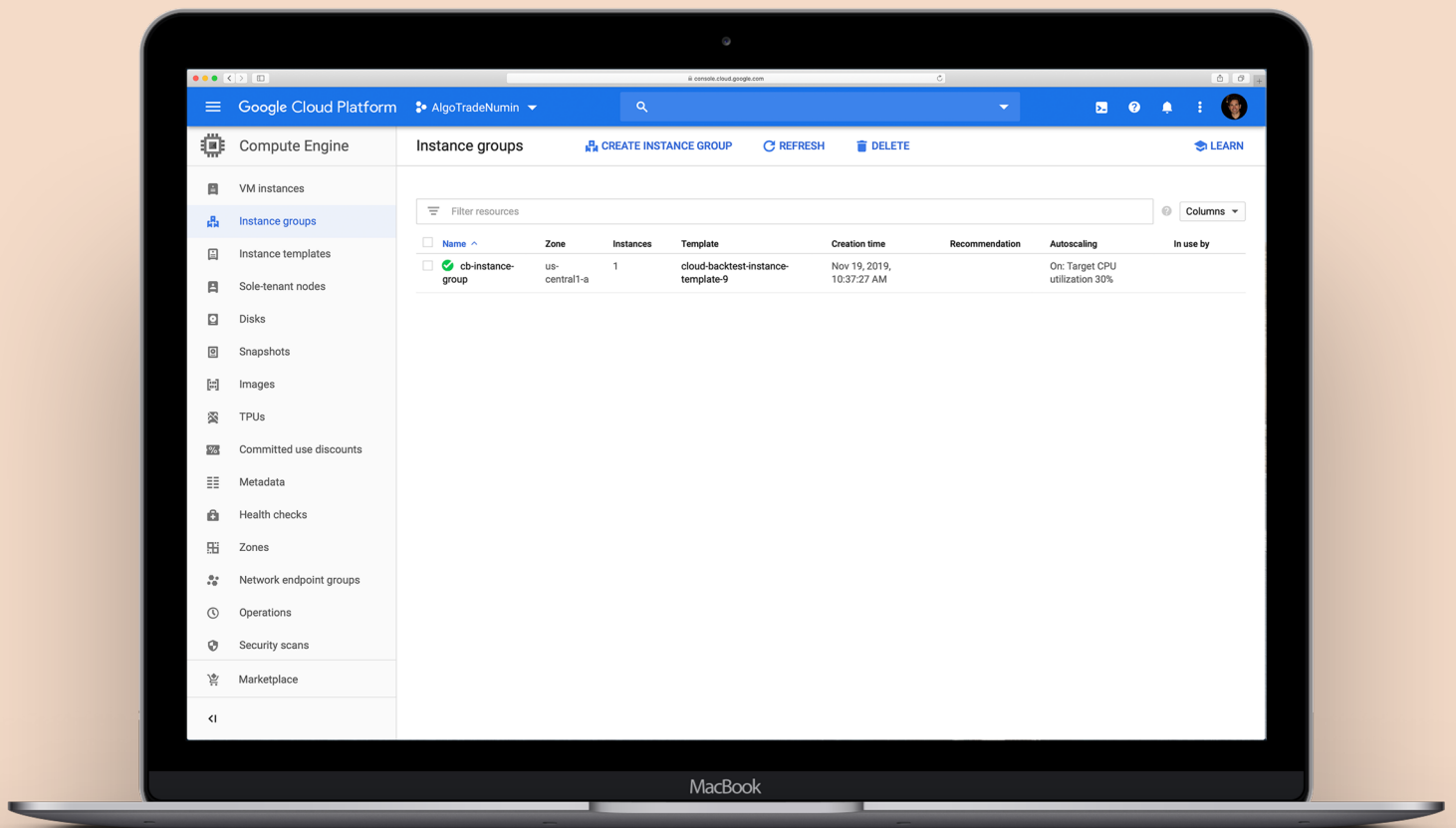


2

SUPPLY

HAVE THE JOB EXECUTED VIA SCRIPTS TO GOOGLE CLOUD.

The backtesting job would hit Google's Cloud servers. A job queue would allocate resources to the machine which would begin the job.

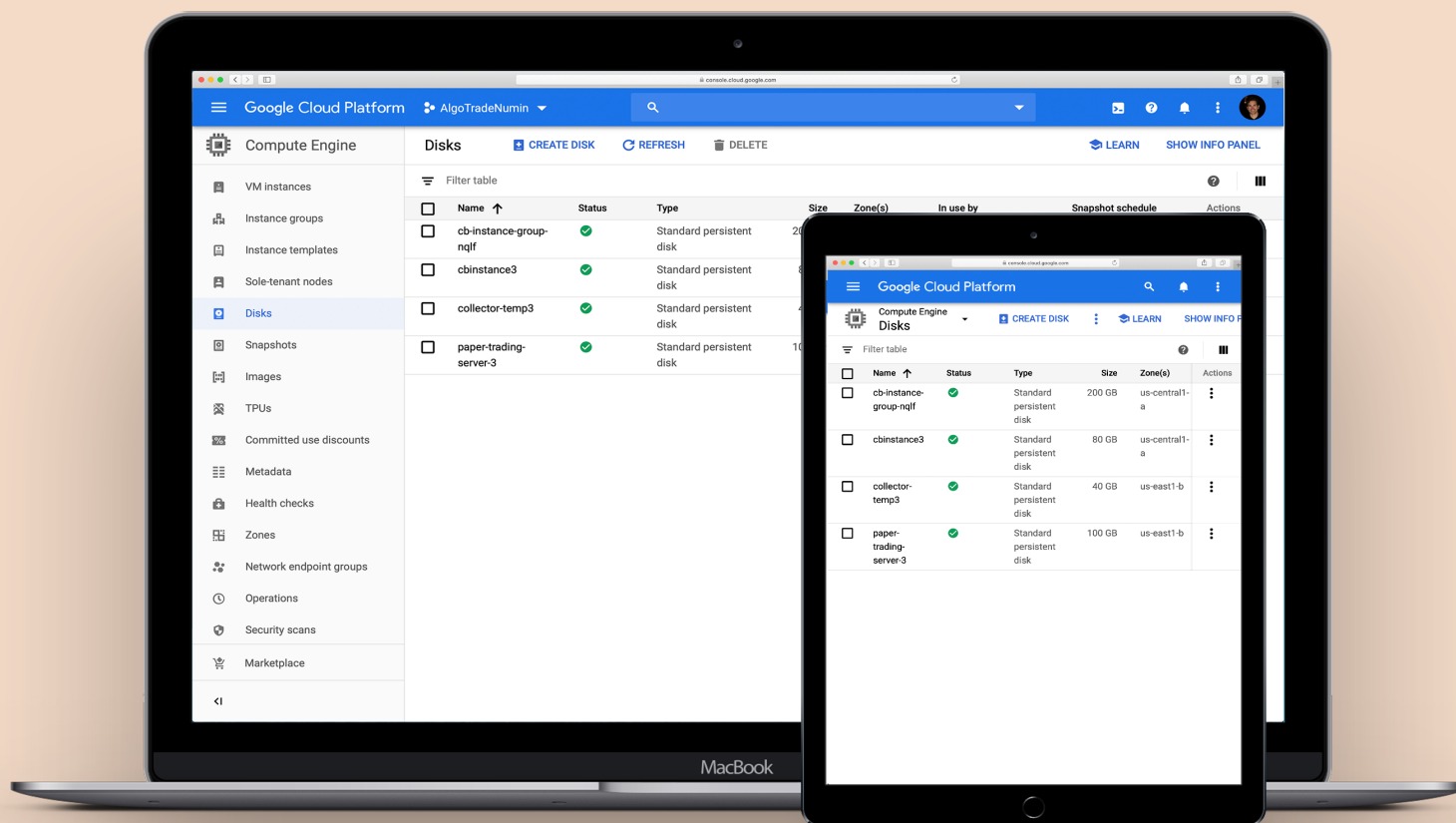


3

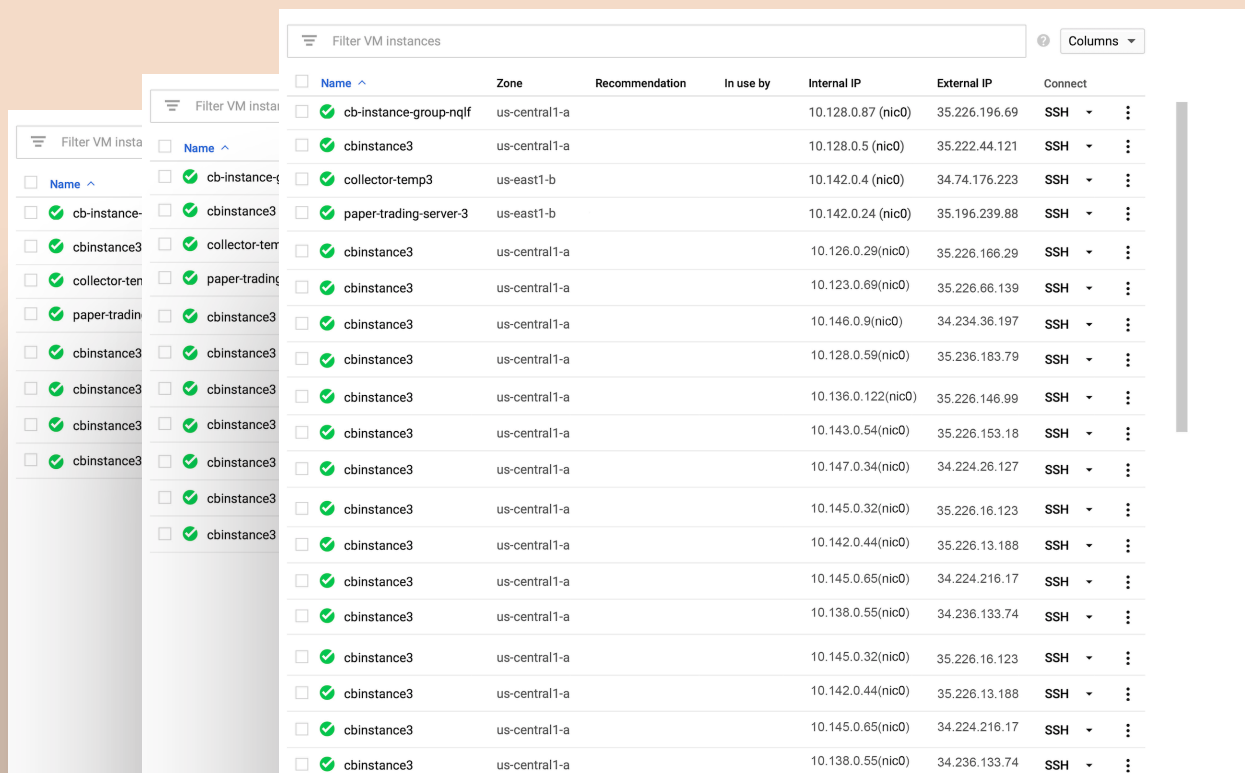
SCALING

FROM 1 TO 800.

Each instance, once its CPU utilization reached beyond a certain threshold, would self-spawn a copy spot instance which would then start pulling job queues off a job stack. As spot instances can be dequeued by Google at any moment, if a machine was dequeued a new machine would be spawned in its place. In cases where a new spot instance was not available, the stack would wait and attempt to respawn at key time intervals.



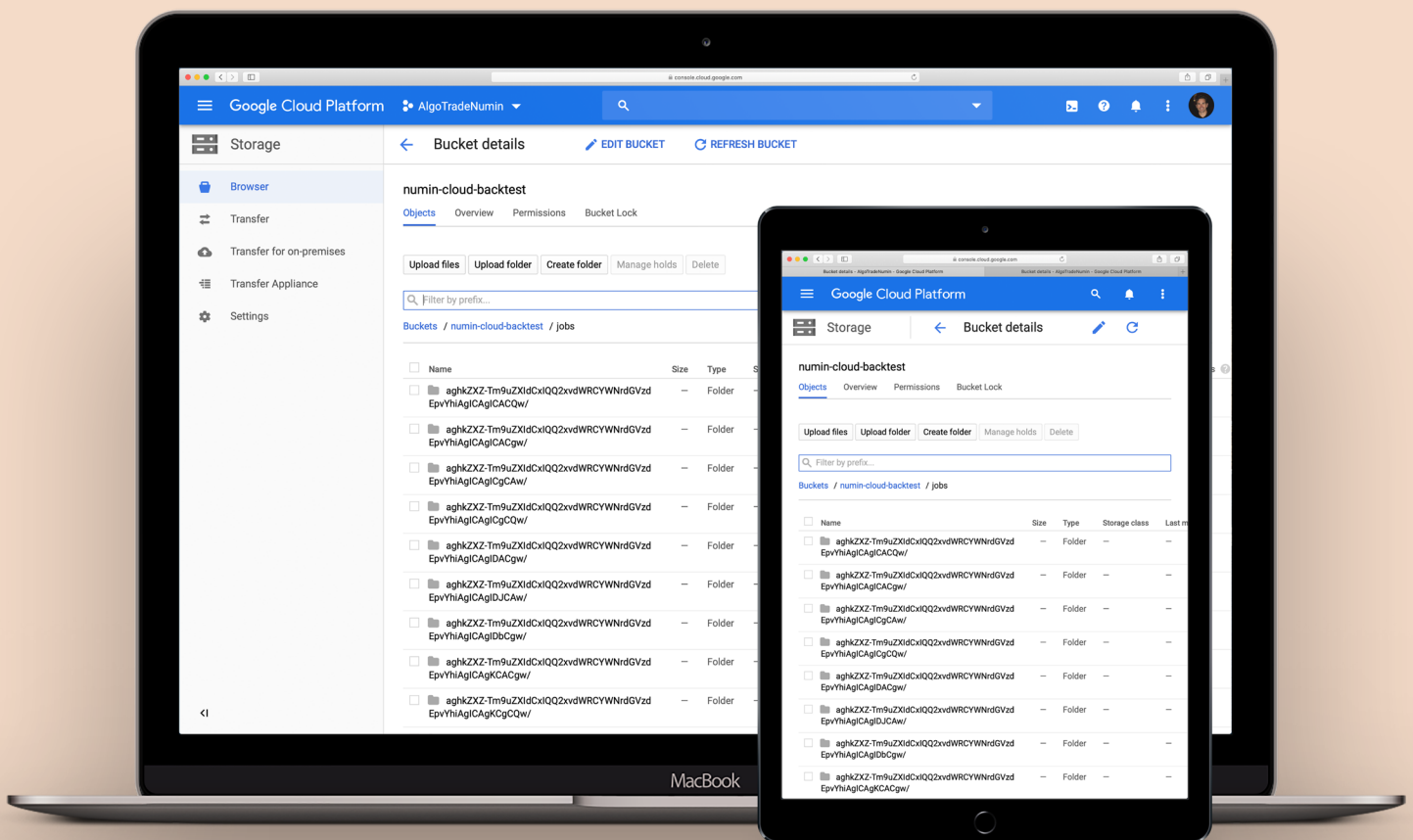
Hundreds of spawned machines would coordinate together to complete the task.



The screenshot displays the Google Cloud VM instances management interface. It features a search bar at the top labeled "Filter VM instances" and a "Columns" dropdown menu. Below the search bar is a table listing VM instances. The table has columns for Name, Zone, Recommendation, In use by, Internal IP, External IP, and Connect. The instances are listed in a table with 7 columns: Name, Zone, Recommendation, In use by, Internal IP, External IP, and Connect. The instances are sorted by Name. The first instance is "cb-instance-group-nqlf" in the "us-central1-a" zone. The second instance is "cbinstance3" in the "us-central1-a" zone. The third instance is "collector-temp3" in the "us-east1-b" zone. The fourth instance is "paper-trading-server-3" in the "us-east1-b" zone. The fifth instance is "cbinstance3" in the "us-central1-a" zone. The sixth instance is "cbinstance3" in the "us-central1-a" zone. The seventh instance is "cbinstance3" in the "us-central1-a" zone. The eighth instance is "cbinstance3" in the "us-central1-a" zone. The ninth instance is "cbinstance3" in the "us-central1-a" zone. The tenth instance is "cbinstance3" in the "us-central1-a" zone. The eleventh instance is "cbinstance3" in the "us-central1-a" zone. The twelfth instance is "cbinstance3" in the "us-central1-a" zone. The thirteenth instance is "cbinstance3" in the "us-central1-a" zone. The fourteenth instance is "cbinstance3" in the "us-central1-a" zone. The fifteenth instance is "cbinstance3" in the "us-central1-a" zone. The sixteenth instance is "cbinstance3" in the "us-central1-a" zone. The seventeenth instance is "cbinstance3" in the "us-central1-a" zone. The eighteenth instance is "cbinstance3" in the "us-central1-a" zone. The nineteenth instance is "cbinstance3" in the "us-central1-a" zone. The twentieth instance is "cbinstance3" in the "us-central1-a" zone. The twenty-first instance is "cbinstance3" in the "us-central1-a" zone. The twenty-second instance is "cbinstance3" in the "us-central1-a" zone. The twenty-third instance is "cbinstance3" in the "us-central1-a" zone. The twenty-fourth instance is "cbinstance3" in the "us-central1-a" zone. The twenty-fifth instance is "cbinstance3" in the "us-central1-a" zone. The twenty-sixth instance is "cbinstance3" in the "us-central1-a" zone. The twenty-seventh instance is "cbinstance3" in the "us-central1-a" zone. The twenty-eighth instance is "cbinstance3" in the "us-central1-a" zone. The twenty-ninth instance is "cbinstance3" in the "us-central1-a" zone. The thirtieth instance is "cbinstance3" in the "us-central1-a" zone. The thirty-first instance is "cbinstance3" in the "us-central1-a" zone. The thirty-second instance is "cbinstance3" in the "us-central1-a" zone. The thirty-third instance is "cbinstance3" in the "us-central1-a" zone. The thirty-fourth instance is "cbinstance3" in the "us-central1-a" zone. The thirty-fifth instance is "cbinstance3" in the "us-central1-a" zone. The thirty-sixth instance is "cbinstance3" in the "us-central1-a" zone. The thirty-seventh instance is "cbinstance3" in the "us-central1-a" zone. The thirty-eighth instance is "cbinstance3" in the "us-central1-a" zone. The thirty-ninth instance is "cbinstance3" in the "us-central1-a" zone. The fortieth instance is "cbinstance3" in the "us-central1-a" zone. The forty-first instance is "cbinstance3" in the "us-central1-a" zone. The forty-second instance is "cbinstance3" in the "us-central1-a" zone. The forty-third instance is "cbinstance3" in the "us-central1-a" zone. The forty-fourth instance is "cbinstance3" in the "us-central1-a" zone. The forty-fifth instance is "cbinstance3" in the "us-central1-a" zone. The forty-sixth instance is "cbinstance3" in the "us-central1-a" zone. The forty-seventh instance is "cbinstance3" in the "us-central1-a" zone. The forty-eighth instance is "cbinstance3" in the "us-central1-a" zone. The forty-ninth instance is "cbinstance3" in the "us-central1-a" zone. The fiftieth instance is "cbinstance3" in the "us-central1-a" zone. The fifty-first instance is "cbinstance3" in the "us-central1-a" zone. The fifty-second instance is "cbinstance3" in the "us-central1-a" zone. The fifty-third instance is "cbinstance3" in the "us-central1-a" zone. The fifty-fourth instance is "cbinstance3" in the "us-central1-a" zone. The fifty-fifth instance is "cbinstance3" in the "us-central1-a" zone. The fifty-sixth instance is "cbinstance3" in the "us-central1-a" zone. The fifty-seventh instance is "cbinstance3" in the "us-central1-a" zone. The fifty-eighth instance is "cbinstance3" in the "us-central1-a" zone. The fifty-ninth instance is "cbinstance3" in the "us-central1-a" zone. The sixtieth instance is "cbinstance3" in the "us-central1-a" zone. The sixty-first instance is "cbinstance3" in the "us-central1-a" zone. The sixty-second instance is "cbinstance3" in the "us-central1-a" zone. The sixty-third instance is "cbinstance3" in the "us-central1-a" zone. The sixty-fourth instance is "cbinstance3" in the "us-central1-a" zone. The sixty-fifth instance is "cbinstance3" in the "us-central1-a" zone. The sixty-sixth instance is "cbinstance3" in the "us-central1-a" zone. The sixty-seventh instance is "cbinstance3" in the "us-central1-a" zone. The sixty-eighth instance is "cbinstance3" in the "us-central1-a" zone. The sixty-ninth instance is "cbinstance3" in the "us-central1-a" zone. The seventieth instance is "cbinstance3" in the "us-central1-a" zone. The seventy-first instance is "cbinstance3" in the "us-central1-a" zone. The seventy-second instance is "cbinstance3" in the "us-central1-a" zone. The seventy-third instance is "cbinstance3" in the "us-central1-a" zone. The seventy-fourth instance is "cbinstance3" in the "us-central1-a" zone. The seventy-fifth instance is "cbinstance3" in the "us-central1-a" zone. The seventy-sixth instance is "cbinstance3" in the "us-central1-a" zone. The seventy-seventh instance is "cbinstance3" in the "us-central1-a" zone. The seventy-eighth instance is "cbinstance3" in the "us-central1-a" zone. The seventy-ninth instance is "cbinstance3" in the "us-central1-a" zone. The eightieth instance is "cbinstance3" in the "us-central1-a" zone. The eighty-first instance is "cbinstance3" in the "us-central1-a" zone. The eighty-second instance is "cbinstance3" in the "us-central1-a" zone. The eighty-third instance is "cbinstance3" in the "us-central1-a" zone. The eighty-fourth instance is "cbinstance3" in the "us-central1-a" zone. The eighty-fifth instance is "cbinstance3" in the "us-central1-a" zone. The eighty-sixth instance is "cbinstance3" in the "us-central1-a" zone. The eighty-seventh instance is "cbinstance3" in the "us-central1-a" zone. The eighty-eighth instance is "cbinstance3" in the "us-central1-a" zone. The eighty-ninth instance is "cbinstance3" in the "us-central1-a" zone. The ninetieth instance is "cbinstance3" in the "us-central1-a" zone. The ninety-first instance is "cbinstance3" in the "us-central1-a" zone. The ninety-second instance is "cbinstance3" in the "us-central1-a" zone. The ninety-third instance is "cbinstance3" in the "us-central1-a" zone. The ninety-fourth instance is "cbinstance3" in the "us-central1-a" zone. The ninety-fifth instance is "cbinstance3" in the "us-central1-a" zone. The ninety-sixth instance is "cbinstance3" in the "us-central1-a" zone. The ninety-seventh instance is "cbinstance3" in the "us-central1-a" zone. The ninety-eighth instance is "cbinstance3" in the "us-central1-a" zone. The ninety-ninth instance is "cbinstance3" in the "us-central1-a" zone. The hundredth instance is "cbinstance3" in the "us-central1-a" zone.

Name	Zone	Recommendation	In use by	Internal IP	External IP	Connect
cb-instance-group-nqlf	us-central1-a			10.128.0.87 (nic0)	35.226.196.69	SSH
cbinstance3	us-central1-a			10.128.0.5 (nic0)	35.222.44.121	SSH
collector-temp3	us-east1-b			10.142.0.4 (nic0)	34.74.176.223	SSH
paper-trading-server-3	us-east1-b			10.142.0.24 (nic0)	35.196.239.88	SSH
cbinstance3	us-central1-a			10.126.0.29(nic0)	35.226.166.29	SSH
cbinstance3	us-central1-a			10.123.0.69(nic0)	35.226.66.139	SSH
cbinstance3	us-central1-a			10.146.0.9(nic0)	34.234.36.197	SSH
cbinstance3	us-central1-a			10.128.0.59(nic0)	35.236.183.79	SSH
cbinstance3	us-central1-a			10.136.0.122(nic0)	35.226.146.99	SSH
cbinstance3	us-central1-a			10.143.0.54(nic0)	35.226.153.18	SSH
cbinstance3	us-central1-a			10.147.0.34(nic0)	34.224.26.127	SSH
cbinstance3	us-central1-a			10.145.0.32(nic0)	35.226.16.123	SSH
cbinstance3	us-central1-a			10.142.0.44(nic0)	35.226.13.188	SSH
cbinstance3	us-central1-a			10.145.0.65(nic0)	34.224.216.17	SSH
cbinstance3	us-central1-a			10.138.0.55(nic0)	34.236.133.74	SSH
cbinstance3	us-central1-a			10.145.0.32(nic0)	35.226.16.123	SSH
cbinstance3	us-central1-a			10.142.0.44(nic0)	35.226.13.188	SSH
cbinstance3	us-central1-a			10.145.0.65(nic0)	34.224.216.17	SSH
cbinstance3	us-central1-a			10.138.0.55(nic0)	34.236.133.74	SSH

Results would be collated and pushed to a unique Cloud Storage bucket specific to that backtest.



“*The compartment alization of the quants on both teams have **owned their domain**, creating a **research factory** of **continuous improvement** that would have been difficult to achieve otherwise.*”

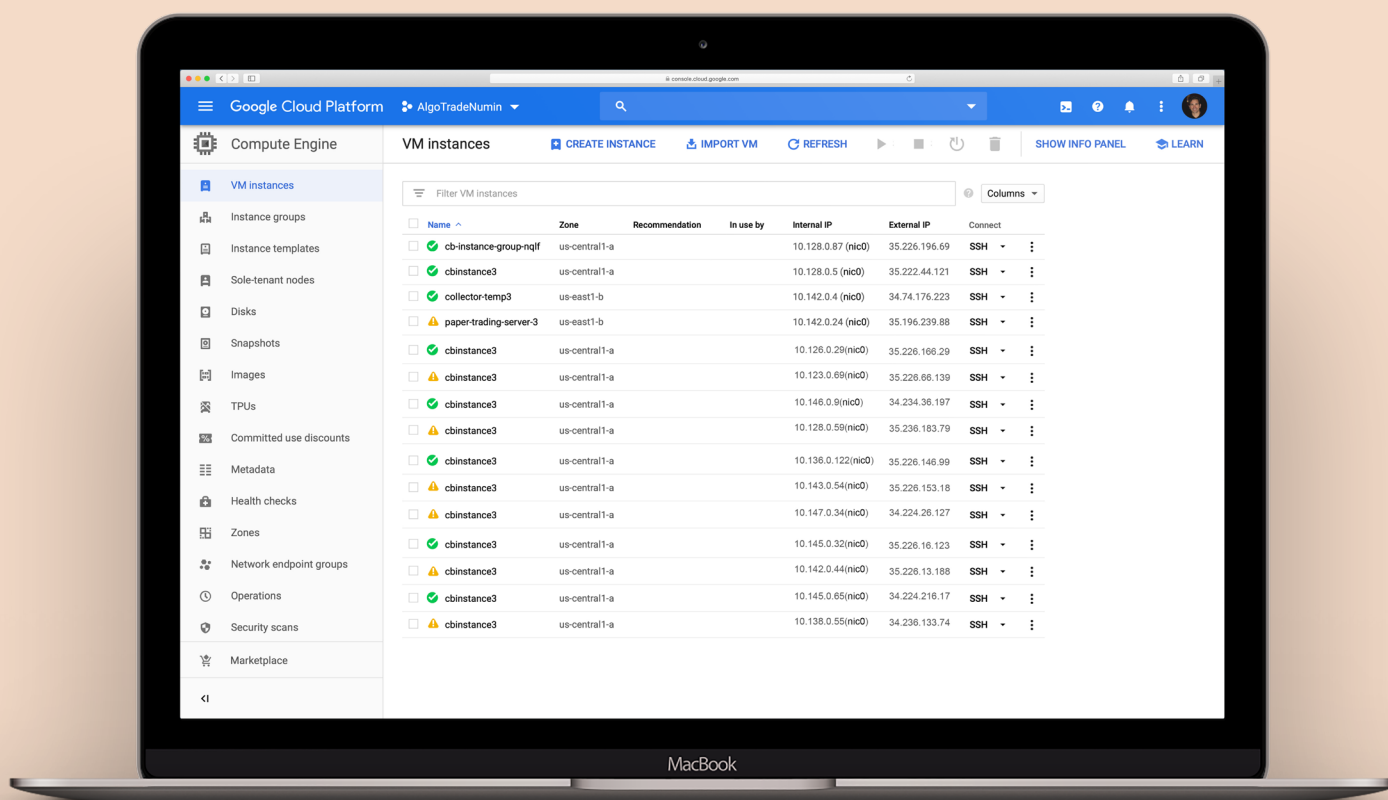
ANUP

*Director of Engineering
Numin*

4

DESCALING

As the job neared completion virtualized machines would dequeue at the rate of spawning to keep costs to an absolute minimum.

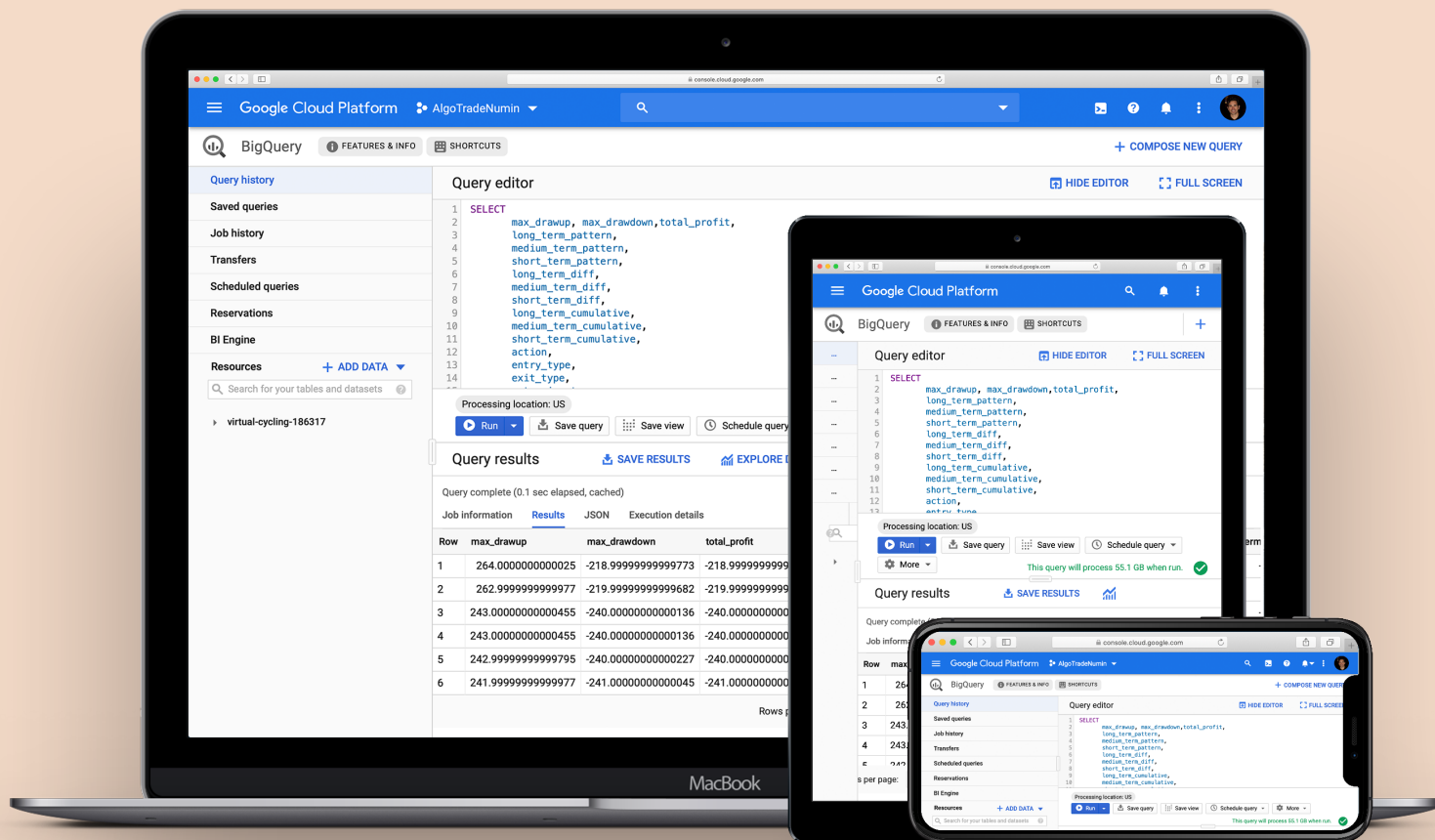


5

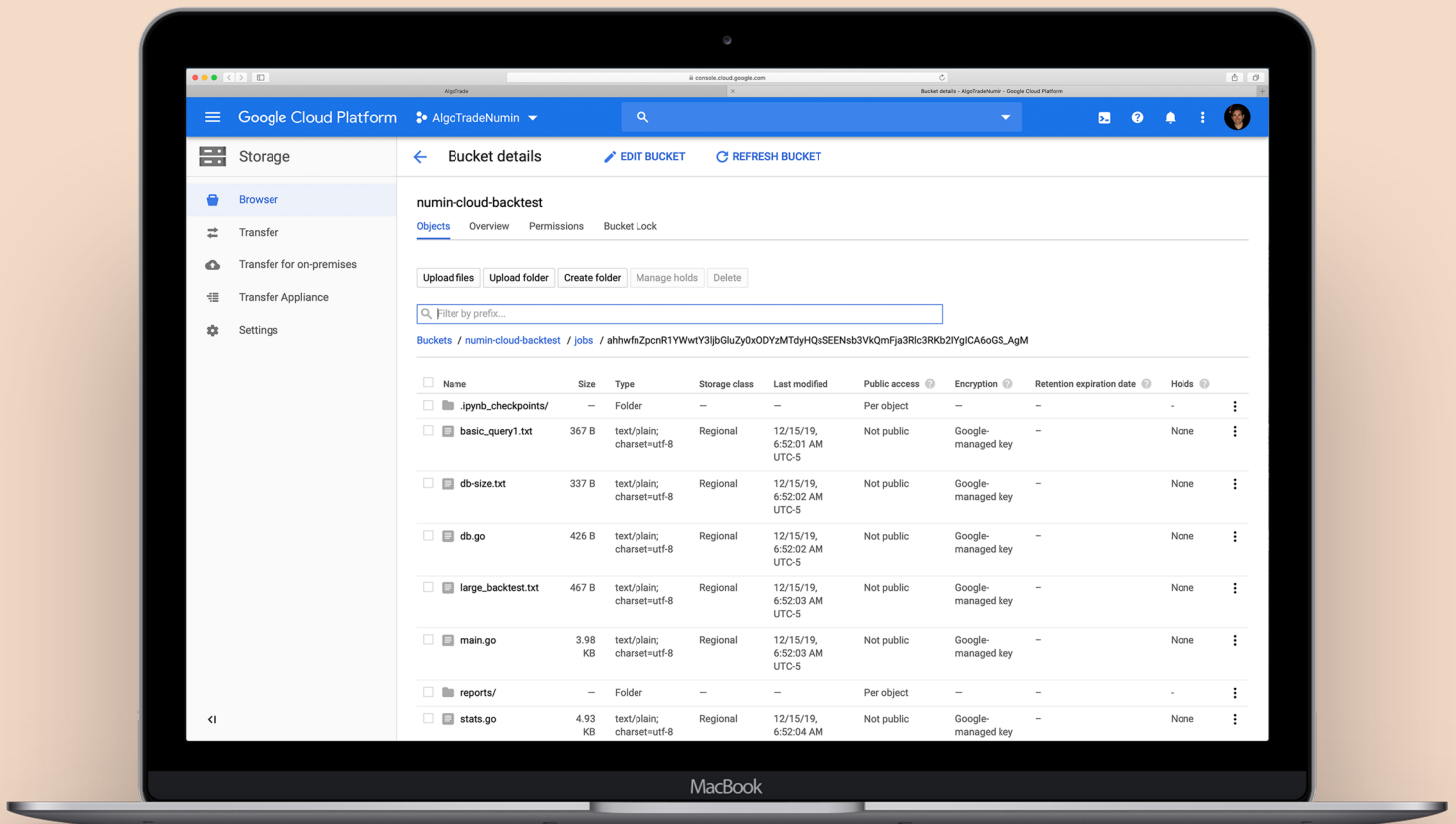
COLLECTION

RETRIEVING THE DATA IN THE RIGHT FORMAT WAS AS IMPORTANT AS CREATING IT.

Total backtest results would be automatically broadcast to Big Query for storage, assigned a trial ID with metadata about the trial itself.



Metadata included how many virtualized CPUs were used, the test parameters, how long the test took, number of machines that were terminated during the test (through error or by Google), and how long dequeuing took as a percentage of queueing. Test metadata was replicated back into Google Cloud Storage unique to the backtest job.





Results

PUNCH PROVIDED A LARGE-SCALE ENTERPRISE BACKTESTING AND CLOUD PROVISIONING AND DEPROVISIONING SOLUTION TO HELP POWER CUTTING EDGE FINANCIAL TRADING ALGORITHMS.

APP STATISTICS

	Items
Spawned Virtualized CPUs	800
Length of Project	4 months
Tech stack	Google Cloud Platform, Big Query, Google Cloud SK
Teams involved	San Francisco & Lahore

Legal

Prepared on February 14, 2020. This document is confidential. It contains material intended solely for the original recipient. Ideas presented here are the property of Punch and are copyrighted.

© 2020 Punch. All rights reserved. Confidential.

punch