

## Introduction

This project/Experimentation involved simulating and generating synthetic genomic data (Using Varsim) in a highly scalable fashion to test Genomic pipeline tools used in Next generation Sequencing (NGS). It helps newly establishing projects in this space to test their NGS platforms with synthetic genomic data. The project would also discuss the architecture in achieving the big data scale of operations and the tools used and experimented (Docker, Kubernetes, Container engine & Registry, Cloud functions, storage, Google cloud) with to achieve the end results.

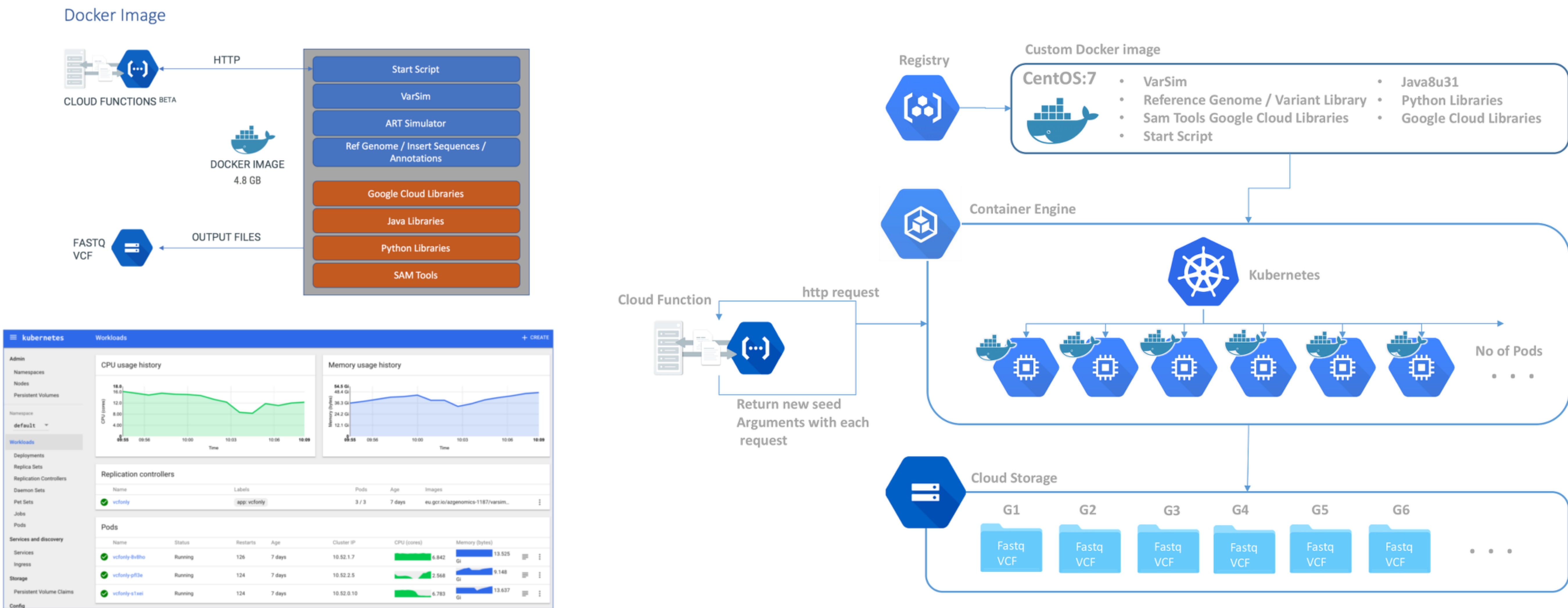


### Purpose

- Ability to simulate Genome Sequencing
- Generated Genomic Data can be used to test pipelines and infrastructure
- If run in a distributed mode can generate sufficient data to create near production line scenarios
- Test various analytics tools
- Synthetic so no issues with privacy, patient de-identification or transfer across regions

### Methods

- VarSim is a Python/Java based tool that would simulate one genome per run.
- We looked into ways how we could parallelize the execution to generate more genomes.
- Build Docker container/Image for the tool.
- Execution on Google Cloud – Container engine, Kubernetes.
- Parameters like Coverage, Unique ID & Seed value were externalized in a Lambda function that the Docker images could talk to and receive arguments before execution.
- Output FASTQ files and VCF's would then be stored in Cloud storage.
- Ability to choose to generate FASTQ & VCF or just VCF.



## Results

We have around 1000 unique VCF files generated so far on the trial run. We have the ability to generate the required number of Genome FASTQ and VCF's The Architecture can be re-used for other distributed compute applications.

Buckets / genedata / cluster2		
G360.vcf	199.1 MB	
G361.vcf	195.48 MB	
G362.vcf	197.79 MB	
G363.vcf	196.16 MB	
G364.vcf	194.39 MB	
G365.vcf	195.31 MB	
G366.vcf	199.79 MB	

Buckets / genedata / cluster / G7		
Name	Size	Type
lane0.read1.fq.gz	86.72 GB	application/octet-stream
lane0.read2.fq.gz	88.71 GB	application/octet-stream
simu.truth.vcf	197.07 MB	text/vcf+xml

Varsim was pickedup as the tool for Genome Simulation, as it provided ability to insert variations into the Genome Simulation and output FASTQ and VCF files