

# **Genomics Deployments: Enabling Precision Medicine with IBM Spectrum Scale**

**Sandeep Patil**  
**IBM STSM, Master Inventor**



# Acknowledgement

## IBM

Frank N Lee, Ulf Troppens, Yael Shani, Carl Zetie, Kevin Gildea, Piyush Chaudhary, Theodore Hoover Jr, Kumaran Rajaram, Joanna Wong, Monica Lemay, Luis Bolinches

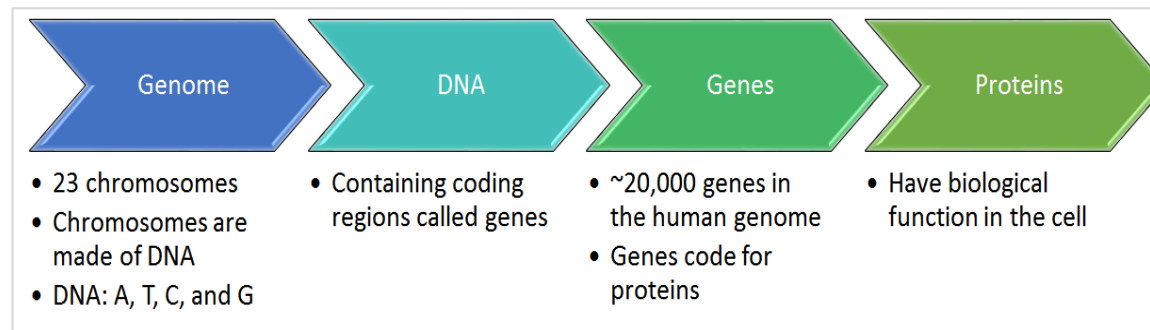
# Agenda

- ❑ Genomic Introduction
  - ❑ What is Genomics
  - ❑ Genomics – An Emerging Market
- ❑ Understanding Genomic Sequencing Workloads
- ❑ Requirements on Infrastructure
- ❑ Solution Approach
- ❑ Solution Architecture
- ❑ Performance for GATK based on proposed solution



# Genomics - Introduction

- Genomics is a branch of biotechnology focusing on genomes.
- Genomics involves applying the techniques of genetics and molecular biology to sequence, analyze or modify the DNA of an organism.
- It finds its use in a number of fields, such as, diagnostics, personalized healthcare, agricultural innovation, forensic science and others.

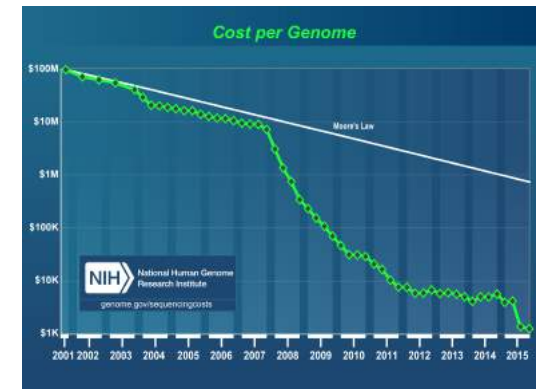


Frost & Sullivan: Global Precision Medicine Growth Opportunities, Forecast to 2025

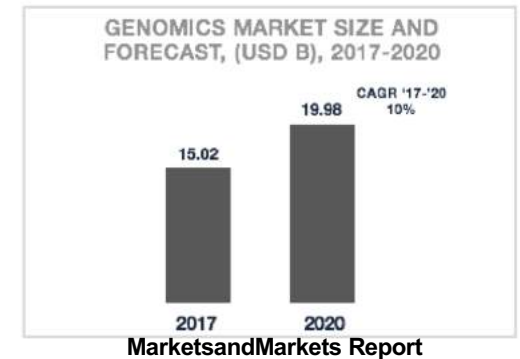
# Genomics – Path Towards Precision Medicine

## Why Genomics is becoming more relevant ?

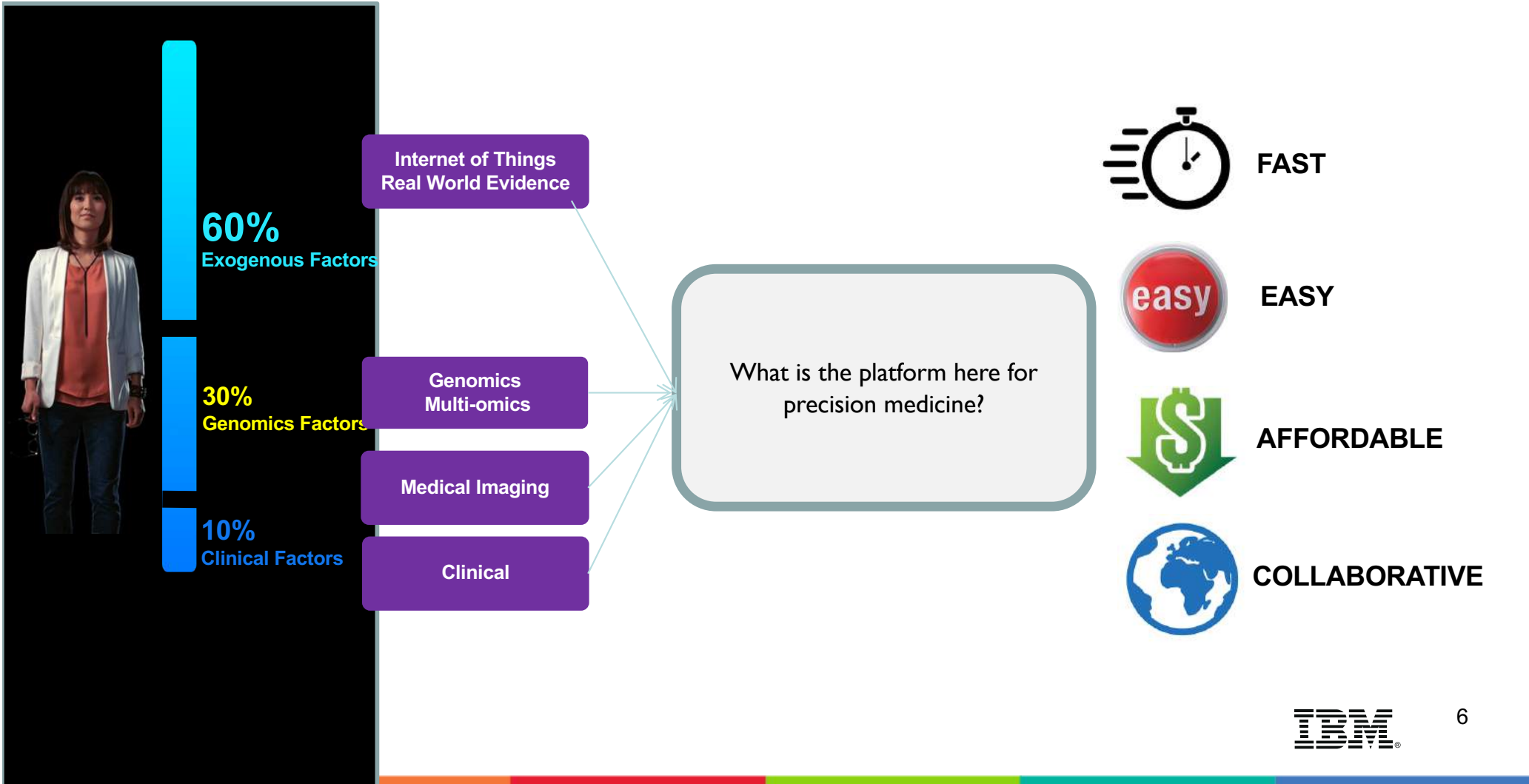
- **Feasibility:** Decreased cost of sequencing.
  - First sequencing of the whole human genome in 2003 cost roughly \$2.7 billion
  - Today Genome sequence can cost around 1000 to 1500 USD
  - DNA sequencing players target to get it down to 100 USD
- **Value:**
  - Genomics is bringing in an era of proactive and personalized medicine (among other fields) – Potential of disruption.
- **Investment:**
  - The market for genomic products and services is growing at 10% and is predicted to become a \$20 billion opportunity.
  - Growth in the Genomics market is majorly attributed to increasing government initiatives and increasing research.



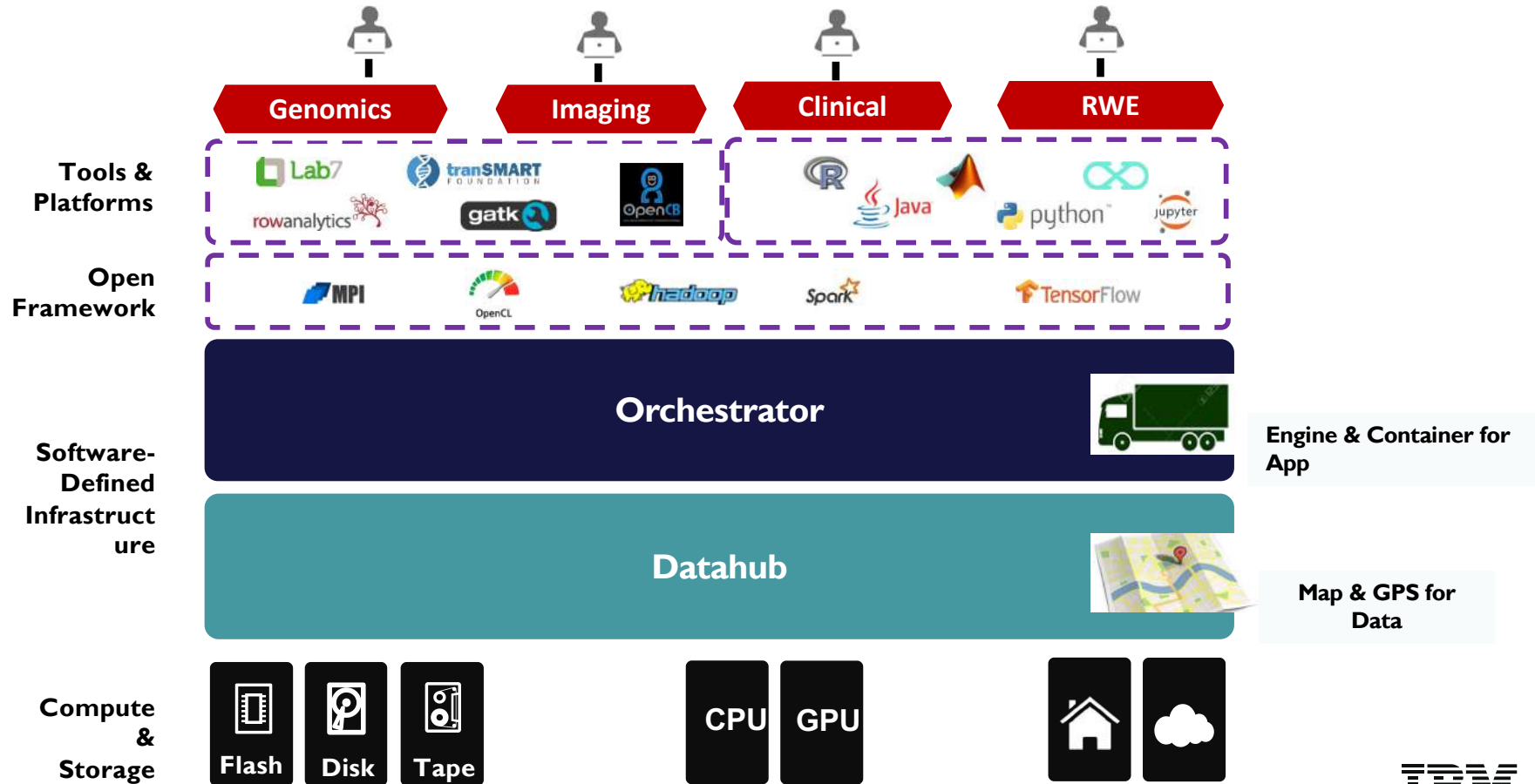
<https://www.genome.gov/27565109/the-cost-of-sequencing-a-human-genome/>



# Moon Shot: Building Cognitive Platform for HC/LS



# Reference Architecture for the Foundation



# For Data Hub Part: Questions for Storage Community

## Questions that come in mind....

- ❑ **Customer:** Is there a reference architecture or approach.
- ❑ **Solution Architect :** How do I solutionize storage for Genomics ... What is the workload requirements ?
- ❑ **Storage Developer:** what I developed meets the genomics requirement... What is the workload looks like ?
- ❑ **Storage Tester:** did my testing cover the requirements for genomic workload...What is the workload, what are the tools,





# Questions for Storage Community (Cont.)

Answer to the Questions in mind :

Need to understand the Genomics Sequencing  
Workload from Storage Perspective !



# Genomic Sequencing Workload– High Level

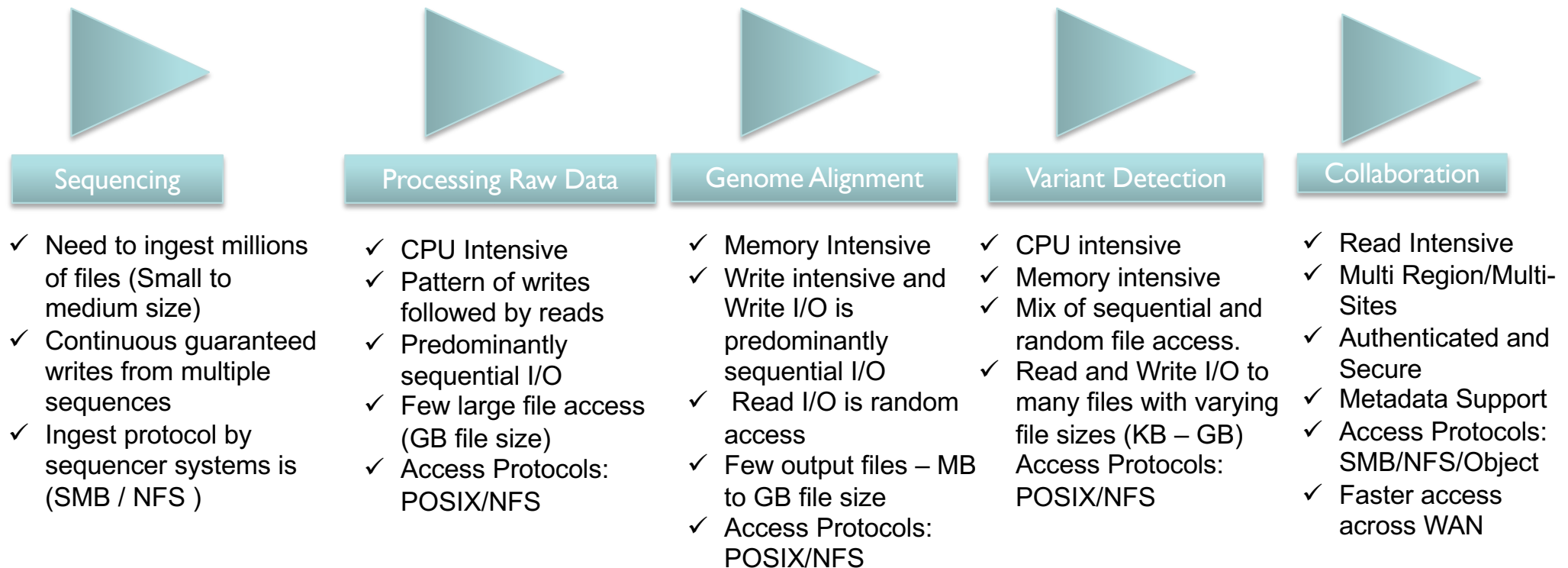
## Key Characteristics of Genomic Workload

- Genomics requires a significant focus on big data management as the sequencing of the genome results in the production of a large amount of data.
- Genomic [data analysis](#) requires 3 process steps:
  1. Sequencers convert the physical sample to raw data. ‘
  2. Raw data is put in a sequence corresponding to the genome.
  3. Analytics (example: matching mutations with certain diseases), is then performed.

Requires easy to use and scalable IT Infrastructure for:  
1) Owning, managing and accessing PBs of file storage  
2) High throughput batch processing to analyze data.

# Genomic Sequencing Workload to Storage Requirement Mapping

(based on GATK3 pipeline reference)



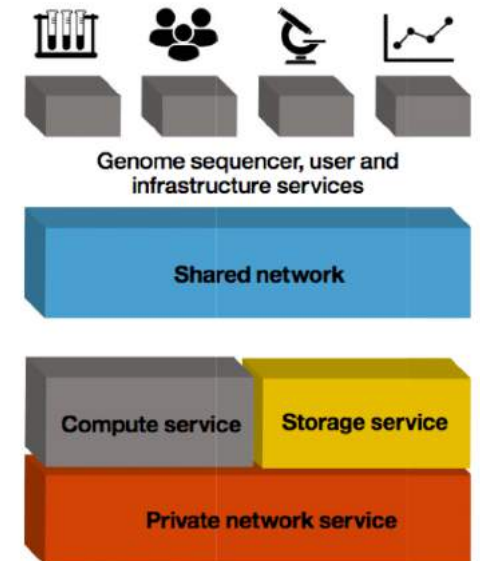
## Need for Optimal Solution.

- Need to think end to end which included Compute, Network and Storage as the key building blocks.
- The infrastructure (Compute, Network & Storage) should allow elasticity to scale-in / scale-out of the building blocks similar to “Lego” blocks.
- For Storage Building Block : Need for a high performance file storage with multiple access interfaces/protocols – Not a typical Network Attach Storage (NAS) as genomic sequencing workload is not a NAS workload but a **Technical Computing** workload.

# Need for Elasticity like 'Lego' Blocks

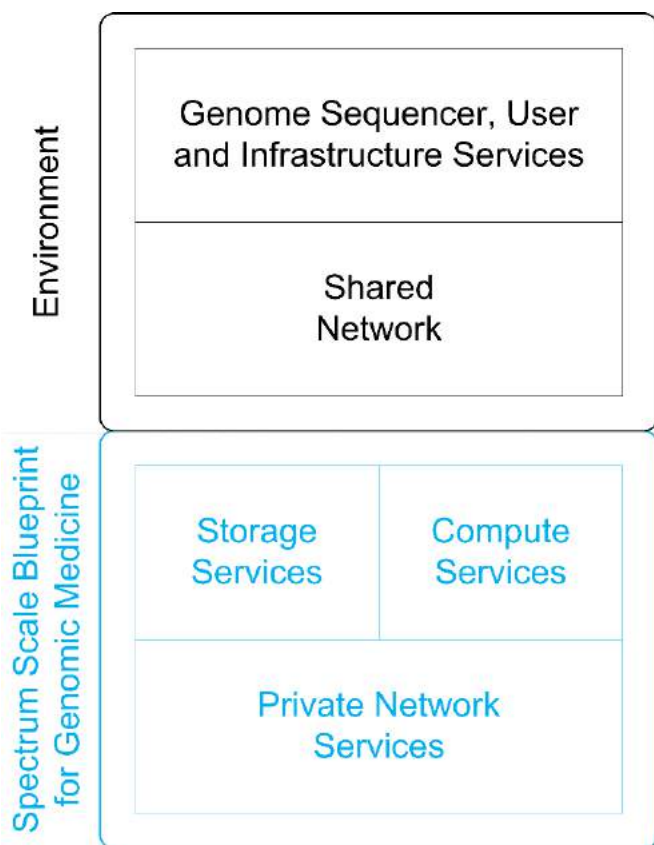
## ... Choosing the Composable Infrastructure Principal

- Composable solutions are built in a way that disaggregates the underlying building blocks viz. compute, storage, and network services.
- These disaggregated services provide the required granularity allowing the infrastructure that can be sliced, diced, expanded and contracted at will and based on the actual need.
- It facilitates ease in deployment with well defined configuration and tuning templates per building block.
- Genomic workloads benefit from composable principals as one can grow and shrink the building blocks based on the needs.



Composable building block for genomics.

# For Genomics – A Composable Building Block Approach



## Shared Network

- **High-speed NFS , SMB , Object Data Access**, connected to shared campus network.
- **User Login** to submit and manage batch jobs and to access interactive applications.

## Compute Services

- Scale-able **Compute Cluster** to analyze genomics data.

## Storage Services

- Scale-able **Storage Cluster** to store, manage and access genomic data.

## Private Network Services

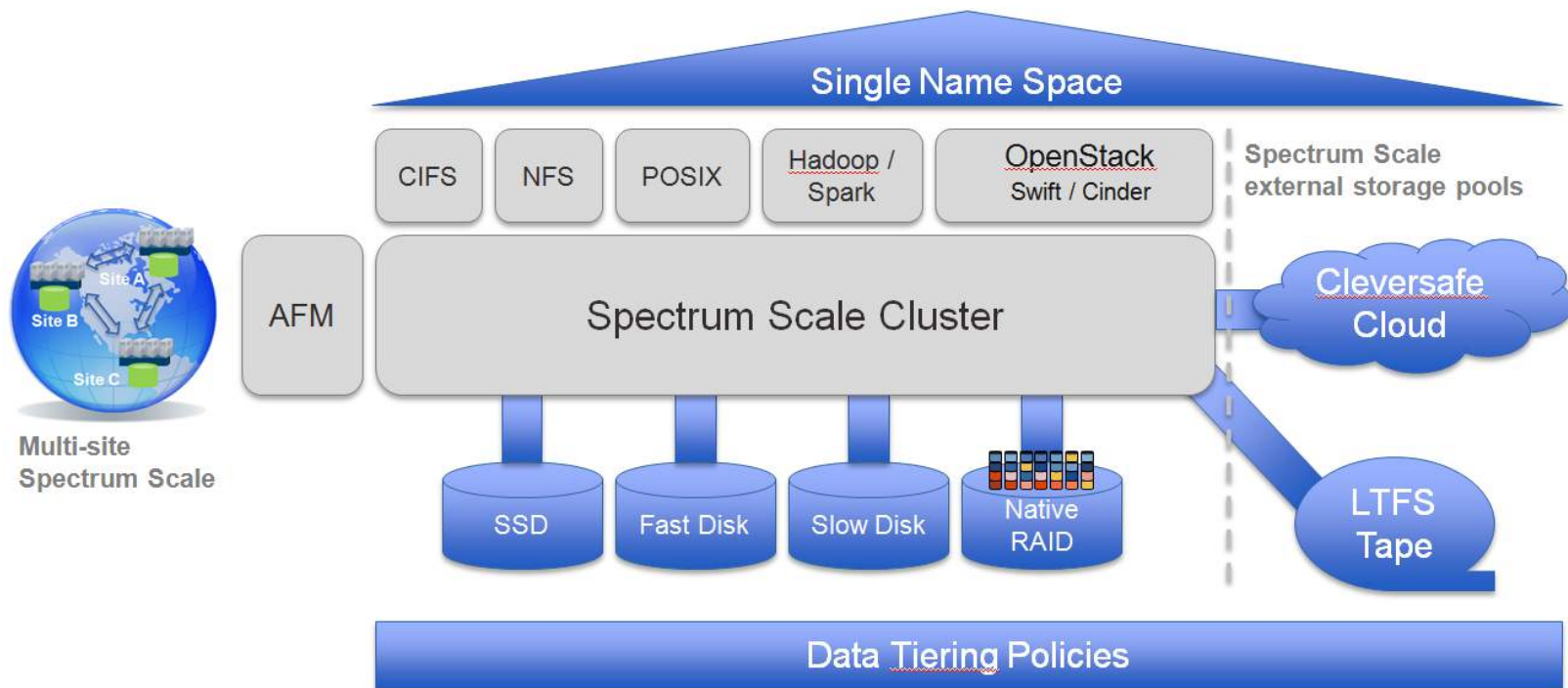
- **High-speed Data Network**, not connected to data center network.
- **Provisioning Network** and **Service Network** for administrative login and hardware services, optionally connected to shared campus network.

## Storage Service: Need for High Performance File Storage aligning to Composable Infrastructure Principals

- **Key requirements per genomic sequencing workload**
  - High Performance & high throughput is key – Technical Computing workload , HPC-like , not a typical NAS workload.
  - Should support scale-in and scale-out to adhere to composable infra principals.
  - Ability to support different type of storage backend (need to be software defined)
  - Support global namespace across different stages of sequencing.
  - Multiprotocol support like NFS,SMB,HDFS,POSIX,Object for data ingestion, collaboration and computing the sequencing.
  - Easy ability for archive to low cost tier (like Tape), cloud integration.

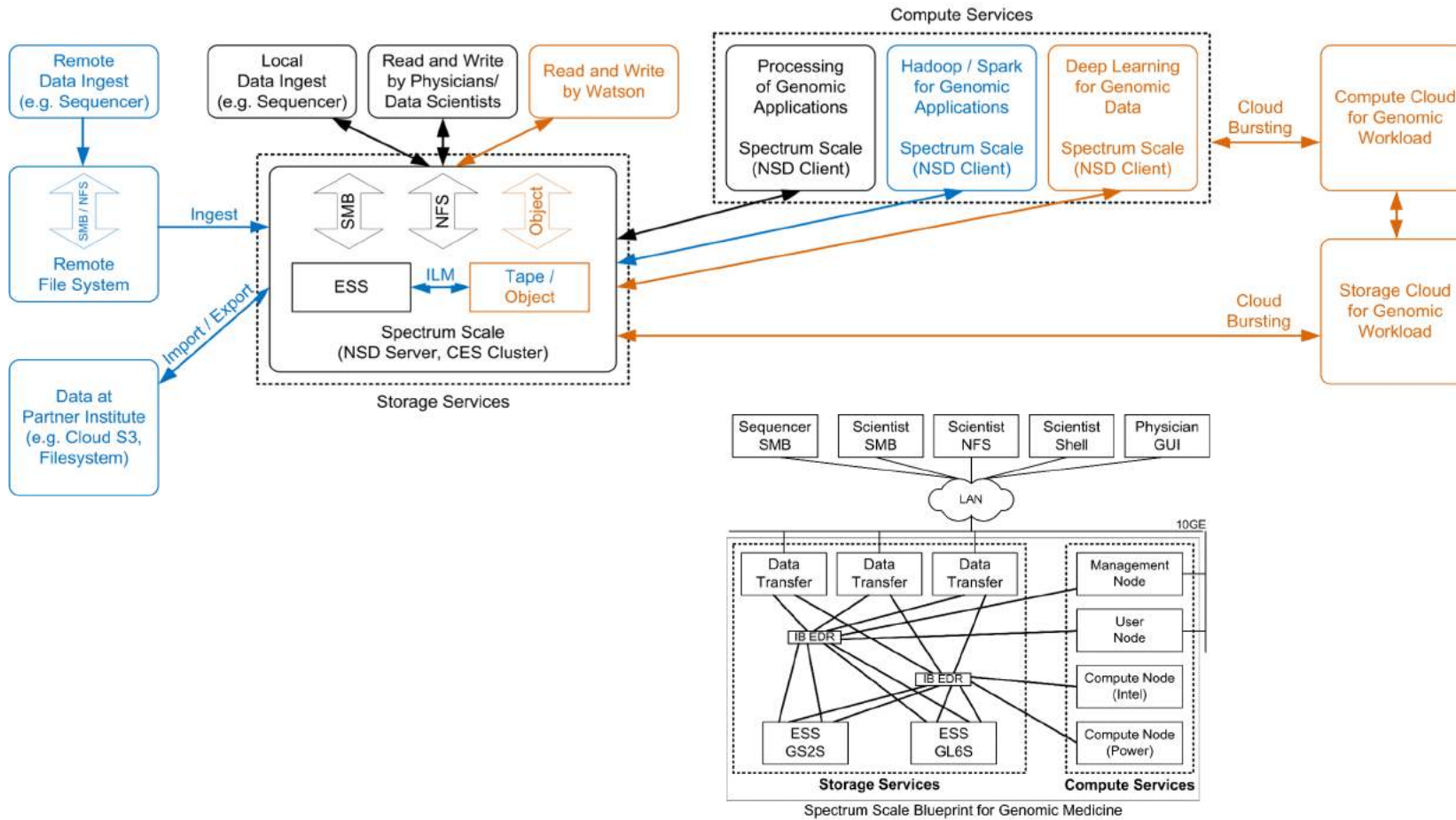
***Storage Solution: Taking the “Software Defined” approach and choosing a clustered filesystem that meets the above requirements***

# Choosing a High Performance Clustered Filesystem for Storage – IBM Spectrum Scale By Design is The Ideal Candidate !



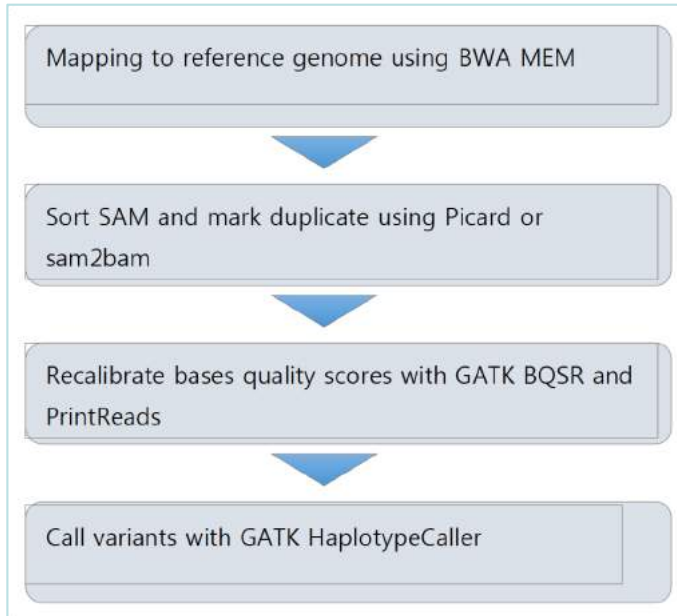


# Solution Architecture: Putting it all together



## Accelerated Performance for Genomics Sequencing

### GATK Workflow – Execution Time on Profiling Environment using the Proposed Solution Architecture for single sample



#### Profiling environment:

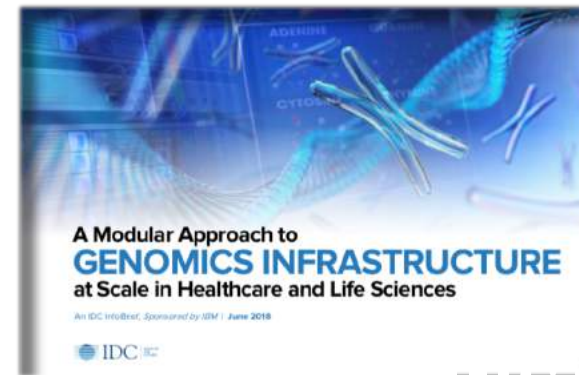
- 1x Power8 Node (IBM 8247-22L with SMT=8) with 256GB memory to execute whole workflow.
- 1x IBM ESS GS4 storage based on SSD ( $\geq 23$  GB/s write bandwidth and  $\geq 30$  GB/s read bandwidth)
- Dual rail FDR InfiniBand aggregating to  $\sim 13$  GB/s

	Solexa WGS Broad dataset with b37 reference
BWA-Mem	303 min 47 sec
sam2bam (storage mode)	35 min 53 sec
GATK BaseRecalibrator (java setting -Xmn10g -Xms10g -Xmx10g)	87 min 21 sec
GATK PrintReads (java setting -Xmn10g -Xms10g -Xmx10g)	97 min 1 sec
GATK HaplotypeCaller (java setting -Xmn10g -Xms10g -Xmx10g)	261 min 37 sec
GATK mergeVCF (java setting -Xmn10g -Xms10g -Xmx10g)	0 min 51 sec

**Note:** Execution time was measured on the testbed configuration (detailed in profiling environment). The actual Genomics application performance will depend on testbed configuration, tunings, and other factors.

# Spectrum Scale for Genomics - Collaterals

- **Solution Brief: Deeper, Faster insights with compostable building blocks based in IBM Spectrum Scale**
  - Gives a quick overview of the solution its advantages and references.
  - Download from:  
<https://public.dhe.ibm.com/common/ssi/ecm/ts/en/tss03239usen/systems-hardware-system-storage-ts-solution-brief-tss03239usen-20180105.pdf>
- **An IDC InfoBrief: A Modular Approach to GENOMICS INFRASTRUCTURE at Scale in Healthcare and Life Sciences**
  - Download from:  
<https://www-01.ibm.com/common/ssi/cgi-bin/ssialias?htmlfid=37016937USEN>



# Blueprint : Best Practices Guide

- **Spectrum Scale Best Practices Guide for Genomics Medicine Workload**

- 1) Solution Overview
- 2) Best Practices for Compute Services
- 3) Best Practices for Storage Services
- 4) Best Practices for Private Network Services



## References

- Genome Analysis Toolkit Variant Discovery in High-Throughput Sequencing Data.  
<https://software.broadinstitute.org/gatk/>
- IBM Redpaper: IBM Spectrum Scale Best Practices for Genomics Medicine Workloads:  
<http://www.redbooks.ibm.com/abstracts/redp5479.html>
- Performance optimization of Broad Institute GATK Best Practices on IBM reference architecture for healthcare and life sciences: <https://www-01.ibm.com/common/ssi/cgi-bin/ssialias?htmlfid=TSW03540USEN>
- IBM Reference Architecture for Genomics: Speed, Scale, Smarts:  
<http://www.redbooks.ibm.com/abstracts/redp5210.html?Open>

**Thank You!**



## Workload profile for each GATK processing step for one sample

	BWA-Mem	sam2bam (storage mode)	GATK BaseRecalibrator	GATK PrintReads	GATK HaplotypeCaller	GATK mergeVCF
CPU	Intensive. Close to 100% CPU utilization	~93% (initial phase) and ~40% in later phases	~70% CPU utilization	~70% CPU utilization	~40% CPU utilization	Less than 1% CPU utilization
Memory	Low memory consumption	Higher memory consumption with ~223 GB consumed	Total of 18 x Java threads with each thread customized with 10 GB → 180 GB	Total of 18 x Java threads with each thread customized with 10 GB → 180 GB	Not memory intensive	Not memory intensive
File data I/O access pattern	Pattern of writes followed by reads. Predominantly sequential I/O.	Write I/O predominantly sequential I/O. Read I/O is random access in units of 512 KiB	Predominantly read intensive. Read is mix of sequential and random I/O	Mix of read and write. Write I/O is mostly 512 KiB with mix of sequential and random. Read is mostly sequential	Mix of read and write. Write I/O is mix of sequential and random. Read is mostly sequential	Mix of read and write. Read and write I/O is predominantly sequential I/O.
File I/O bandwidth	<= 200 MB/s (read and write)	Write < 2.5 GB/s. Sustained read < 300 MB/s. High degree of pagepool cache hits during reads (< 36 GB/s).	<= 100 MB/s (read and write)	Write < 150 MB/s and read < 75 MB/s.	Write < 100 MB/s and read < 100 MB/s.	Write < 1.5 GB/s and read < 2 GB/s.
File Metadata	<=2 inode updates	Initial phase <= 60 inode updates. Later phase, <=2 inode updates.	~24 file open and ~24 file closes.	~24 file open and ~24 file closes.	~20 file open and ~20 file closes.	~2 file open and ~2 file closes.
Output file(s)	Single output file (*.sam) <= 380 GB file size	Two output files. ~77 GB (.bam) and ~9 MB (.bam.bai).	Total of 52 files. 26 x *.table.log-4* files (<200 KB) and 26 x *.table* files (< 300 KB)	Total of 78 files. 26 x *.recal_reads*.bam* files (< 15 GB), 26 x *.bai* files (< 750 KB), and 26 x *.recal_reads*.bam.log* files (< 200 KB)	Total of 78 files. 26 x *.raw_variants*.vcf* files (< 6 GB), 26 x *.raw_variants*.vcf*.log* files (< 400 KB), and 26 x *.raw_variants*.vcf*.idx* files (< 20 KB)	Single output file (*.raw_variants.vcf) with ~66 GiB file size

Source: IBM Redpaper: IBM Spectrum Scale Best Practices for Genomics Medicine Workloads:

<http://www.redbooks.ibm.com/abstracts/redp5479.html>

