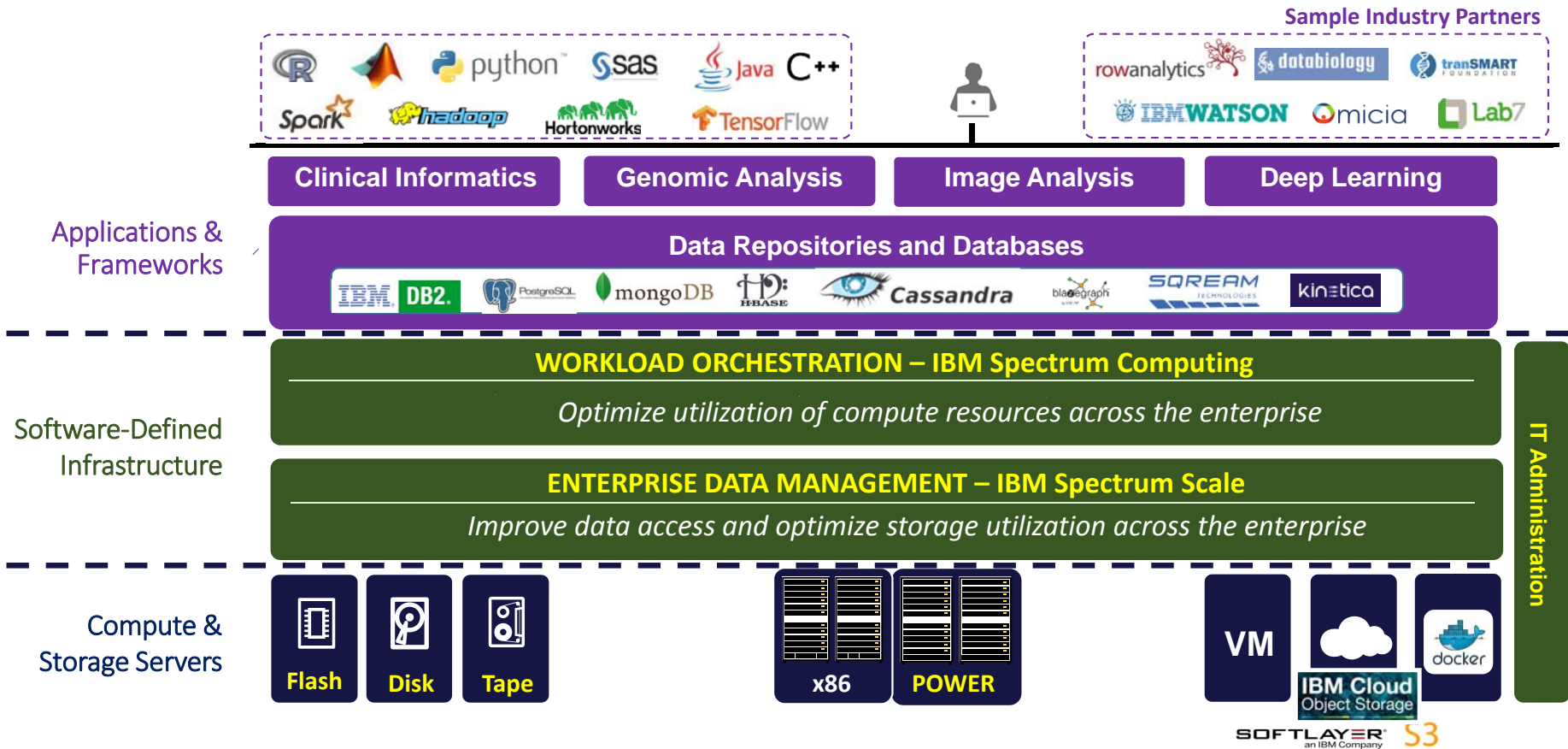# HPC 2.0 for Genomics
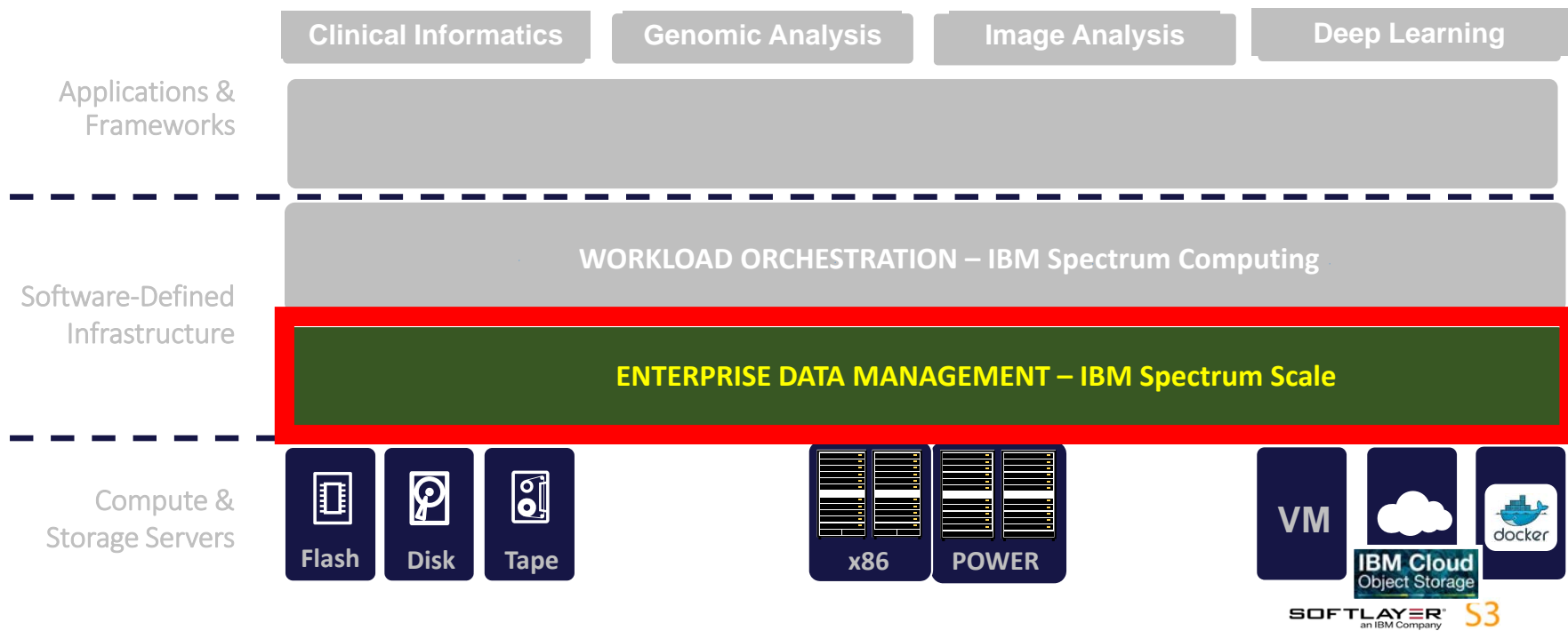
An Introduction to IBM HPDA Framework & Reference Architecture

Frank Lee, PhD
IBM Systems

# IBM Systems Builds the Foundation for the Cognitive Era

**IBM**

**Sample Industry Partners**

**Applications & Frameworks**

| Clinical Informatics | Genomic Analysis | Image Analysis | Deep Learning |

**Data Repositories and Databases**

IBM DB2. · PostgreSQL · mongoDB · H-BASE · Cassandra · bladegraph · SQREAM TECHNOLOGIES · kinetica

**Software-Defined Infrastructure**

**WORKLOAD ORCHESTRATION – IBM Spectrum Computing**

*Optimize utilization of compute resources across the enterprise*

**ENTERPRISE DATA MANAGEMENT – IBM Spectrum Scale**

*Improve data access and optimize storage utilization across the enterprise*

**IT Administration**

**Compute & Storage Servers**

| Flash | Disk | Tape | x86 | POWER | VM | IBM Cloud Object Storage | docker |

SOFTLAYER an IBM Company · S3

# Foundation for Data

**IBM**

| Clinical Informatics | Genomic Analysis | Image Analysis | Deep Learning |

**Applications & Frameworks**

**Software-Defined Infrastructure**

WORKLOAD ORCHESTRATION – IBM Spectrum Computing

ENTERPRISE DATA MANAGEMENT – IBM Spectrum Scale

**Compute & Storage Servers**

Flash | Disk | Tape

x86 | POWER

VM | IBM Cloud Object Storage | docker

SOFTLAYER an IBM Company | S3

# Challenge 1: High Speed for Big Data

**IBM**

**Next Generation Sequencing (NGS)**

- Raw Data: Up to 200 GB/file (compressed)
- Processed "Variant" Data: Up to 500 MB/file

**Biomedical Imaging**

- Medical Imaging: MRI, CT, Ultrasound, ....
- Microscopy

**Time-Varying Sensors**

- Medical Monitors
- Personal Sensors

**Curated Scientific Literature**

- Text files: CSV, TXT
- Online Web Crawls

**Super-speed**

**Super-capacity**

**policy engine for ingesting**

| Before | After |
|---|---|
| 50 | 5 |
| hours using 1 Node ~24cores, 1 QDR link, 256GB RAM | hours using 1 Node ~12cores, 1 FDR Link, 64GB RAM |

Mount Sinai



Big Omics Data Experience

# Data Machine for HPC 2.0

**IBM**

**Capability**

2014

8-10X

2017

20U
0.5PB
4GB/s

34U
5PB
34GB/s

**Fault Tolerance**

2014

>1000X

2017

HW RAID
Rebuild: weeks
MTTDL: 1 week (4+P)

SW RAID
Rebuild: minutes
MTTDL: 200M Y (8+3P)

*MTTDL: for 50,000 disk*

# Challenge 2: Data Sharing for Global Collaboration

**IBM**



**Policy Engine For Sharing**

A data scientist anywhere in the world can access the most nearby resource for data & compute

| C: | /data | HDFS | HTTP | S3 |

**Share Data Anywhere**

**AFM**

# Challenge 3: Cost Control

**policy engine: tiering, compression**

A Decade of Growth for GenBank and NCBI

Genomic data doubling
6-12 months

**But budgets are not!**

**Store Data Anywhere**

Local Tier

Cloud Tier

cleversafe

2014 → 2017

14PB    170PB    unlimited

# Foundation for Workload

**IBM**

| Clinical Informatics | Genomic Analysis | Image Analysis | Deep Learning |

**Applications & Frameworks**

**Software-Defined Infrastructure**

**WORKLOAD ORCHESTRATION – IBM Spectrum Computing**

**ENTERPRISE DATA MANAGEMENT – IBM Spectrum Scale**

**Compute & Storage Servers**

Flash    Disk    Tape

x86    POWER

VM    IBM Cloud Object Storage    docker

SOFTLAYER an IBM Company   S3

**Example #1:** Resource Utilization for Workflow



*time* →

Reduce processing times

*time* →

x86   POWER

VM

docker

Data-aware scheduling with API

IO-aware scheduling with real-time data

IO-aware scheduling with some math

$$\frac{\partial f}{\partial t} = \lim_{h \to 0} \frac{f(t+h, \vec{x}) - f(t, \vec{x})}{h}$$ *github/stjude*

# Challenge 5 Workflow Automation

**IBM**

**Example #1:** Genomic Analysis Pipelines

Whole Human Genome @30x coverage

Processing time per genome

**1 to 100 hours**\*

on 1 compute node

| | | |
|---|---|---|
| High-Throughput Sequencing | FastQ ~ 150 GB (compressed) | |
| Assembly & Alignment | SAM / BAM ~ 100 GB | |
| Variant Calling | VCF 100 to 200 MB | |
| Variant Annotations | Annotated VCF 500 MB | |



Sub-flow module

Job arrays

error-handling

Provides logic syntax

# Client Reference: Workflow Automation

# Foundation for Metadata & Provenance

**IBM**

**Clinical Informatics**  **Genomic Analysis**  **Image Analysis**  **Deep Learning**

Applications &
Frameworks

Software-Defined
Infrastructure

**WORKLOAD ORCHESTRATION – IBM Spectrum Computing**

**ENTERPRISE DATA MANAGEMENT – IBM Spectrum Scale**

Compute &
Storage Servers

**Flash**  **Disk**  **Tape**

**x86**  **POWER**

**VM**

IBM Cloud
Object Storage

docker

SOFTLAYER
an IBM Company  S3

**Filename = MDA1.vcf**
**Fileid = 100034589**
**ProjectID = AXVCBS**
**FlowID = 20170411-2**
**RefGenome = V38**
**UserV1 = EXP1**
**ActionTag = NeverMove**

**What**

**When**

- File name
- File size
- Filer owner
- File path
- File set name
- File permission
- File ctime
- File atime
- File mtime

- Cluster name
- Global file ID
- Job submission user
- Job ID
- Job name
- Flow ID (Parsed from job name)
- Job status
- Job start time
- Job finish time
- Job submission cmd
- Job working directory
- Input files
- User variables

- Cluster name
- Server name and port
- Global file ID
- Flow owner
- Flow ID
- Flow definition name
- Flow definition version
- Flow status
- Flow start time
- Flow finish time
- Flow working directory
- Input files
- User variables

# Beyond Provenance

METADATA-driven Cognitive Engine

What

When

How

Why

# Foundation for Cloud

**IBM**

| | Clinical Informatics | Genomic Analysis | Image Analysis | Deep Learning |
|---|---|---|---|---|

**Applications & Frameworks**

**Software-Defined Infrastructure**

WORKLOAD ORCHESTRATION – IBM Spectrum Computing

ENTERPRISE DATA MANAGEMENT – IBM Spectrum Scale

**Compute & Storage Servers**

Flash  Disk  Tape

x86  POWER

VM  IBM Cloud Object Storage  docker

SOFTLAYER an IBM Company  S3

# A Hybrid Cloud Architecture

**On-premise infrastructure**

**Cloud infrastructure**

Workloads

**On-Premise Cluster**

**Cloud Resident Cluster**

**Spectrum LSF**

Secure VPN tunnel

*IBM Aspera FASP*

**Spectrum Scale (GPFS)**

**Spectrum Scale (GPFS)**

**Spectrum Scale**

**Transparent Cloud Tiering**

**IBM Elastic Storage**

**IBM Elastic Storage**

SOFTLAYER®
an IBM Company

**IBM Cloud** Object Storage

†AFM = Active File Management

# Application-level Optimization



**Applications & Frameworks**

| Clinical Informatics | Genomic Analysis | Image Analysis | Deep Learning |

Data Repositories and Databases

**Software-Defined Infrastructure**

WORKLOAD ORCHESTRATION – IBM Spectrum Computing

ENTERPRISE DATA MANAGEMENT – IBM Spectrum Scale

IT Administration

**Compute & Storage Servers**

Flash | Disk | Tape

x86 | POWER

VM | IBM Cloud Object Storage

**IBM** **IBM POWER continues to develop technologies that accelerate compute for the next generation of analytics, including the latest deep learning & machine learning algorithms**

Sample POWER-based Technology:  NVLink

## OpenPOWER™

**IBM POWER S822LC**

### Basic Advantages over Intel Haswell

| Feature | Intel Haswell | IBM POWER8 |
|---|---|---|
| SMT / Core | 2 Threads | 8 Threads |
| L1d Cache / Core | 32 KB | 64 KB |
| L2 Cache / Core | 256 KB | 512 KB |
| L3 Cache / Processor | 16 to 45 MB | 80 to 96 MB |
| L4 Cache / System | - | 64 MB to 2 GB |
| Maximum Sustained Memory Bandwidth | 53 GB/s | 224 GB/s |

**Nvidia GPU with NVLink**

GPU Memory

NVIDIA    NVIDIA

GPU Memory

720 GB/s

**NVLink 80 GB/s**

720 GB/s

System Memory

115 GB/s

Power S822LC CPU

- Up to 4 integrated Nvidia P100 "Pascal" GPUs
- Delivering > 2.5X the bandwidth to GPUs
- Reduces common CPU-GPU bottlenecks

# Open Source Genomics Applications– Optimized with POWER8

- ALLPATHS-LG
- BarraCUDA
- bamtools
- Bedtools
- Bfast
- Bioconductor
- BioPerl
- BioPython
- BLAST (NCBI)
- Bowtie
- Bowtie2
- BreakDancer
- BWA
- Chimerascan
- Conda
- ClustalW
- Cufflinks
- DELLY2
- EMBOSS

- FASTA
- FastQC
- FASTX-Toolkit
- FreeBayes
- GenomicConsensus
- GenomeFisher
- GraphViz
- HMMER
- HTSeq
- Htslib
- IGV
- InterProScan
- ISAAC3
- iRODS
- Mothur
- MrBayes
- MrBayes5d
- MUSCLE
- Numpy

- Pandas
- PHYLIP
- PICARD
- Pindel
- PLINK
- PRADA
- Pysam
- Python
- R
- RNAStar/STAR
- RSEM
- SAMTools
- Sailfish
- Scalpel
- SHRiMP
- SIFT
- Snpeff
- SOAP3-DP
- SOAPaligner

- SOAPdenovo
- SoapFuse
- SQLite
- Sratoolkit
- STAR-fusion
- Tabix
- Tablet
- Tassel
- T-Coffee
- TMAP
- TopHat
- TranSMART
- Trinity
- Variant_tools
- Varscan
- Velvet/Oases

- bamkit
- bedops
- cutadapt
- diamond
- kraken
- lumpy
- parallel
- PLINK2
- primer3
- QIIME
- R cowplot
- R tidyverse

- Salmon
- Samblaster
- Scikit-bio
- Seqtk
- Spades
- Trimmonmatic
- Vcftools

Bio Builds™    http://biobuilds.org/

- **Turn-key:** Pre-built binaries and complete build scripts
- **Optimized:** POWER8 binaries
- **Long Term Support:** Community sponsorship and support contracts ensure ongoing support for tools

# Qatar Genome Project

**IBM**

**Industry Partners**

R Spark • python • Java • TensorFlow

NEXTCODE HEALTH

**Applications & Frameworks**

Clinical Informatics • Genomic Analysis • Image Analysis • Deep Learning

**Data Repositories and Databases**

mongoDB • neo4j

**Software-Defined Infrastructure**

**WORKLOAD ORCHESTRATION – IBM Spectrum Computing**

**ENTERPRISE DATA MANAGEMENT – IBM Spectrum Scale**

IT Administration

2013 | 2014 | 2014 | 2015 | 2016

قطر جينوم
QATAR GENOME