# Beyond Today's Internet
# Experiencing a Smart Future

Prototype SDX Bioinformatics Exchange: Demonstrating an Essential Use-Case for Personalized Medicine

Robert Grossman – University of Chicago
Joe Mambretti – Northwestern University
Piers Nash – University of Chicago
Jim Chen – Northwestern University
Allison Heath – University of Chicago
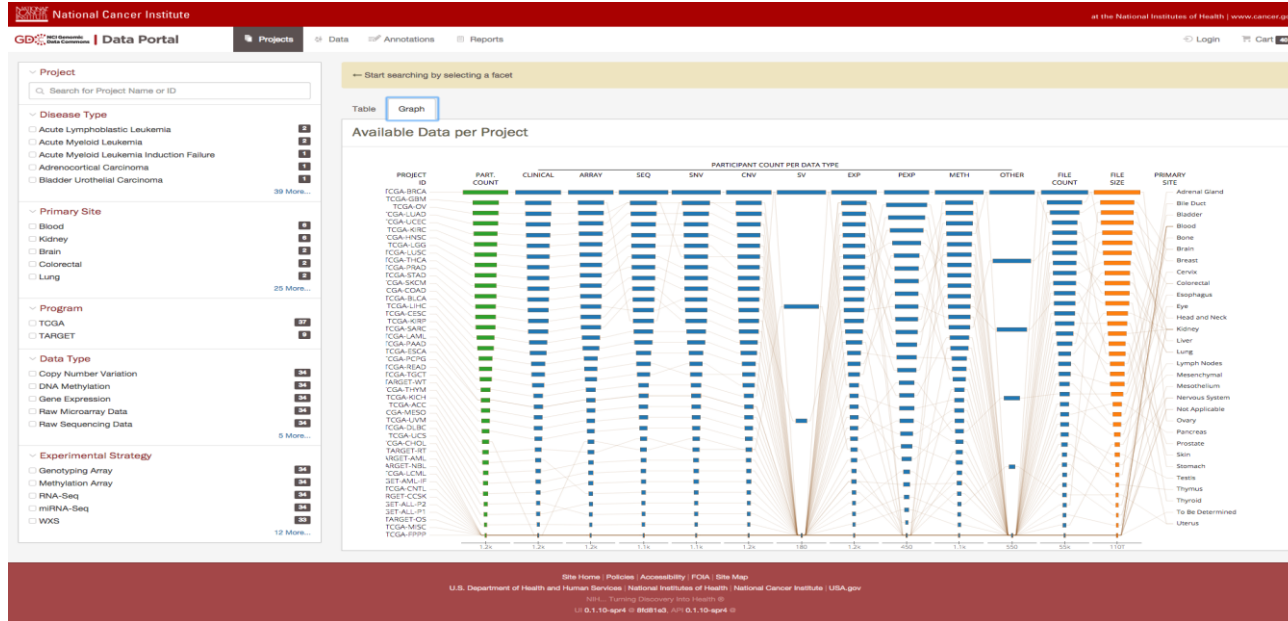
# Precision Medicine

- Precisely match treatments to patients and their specific disease
- Genomic data promises optimal matching.
- 1.7 million cancer cases diagnosed in America each year.
- A single RNA-seq file is 10-20 GB, Whole genome raw data files are > 100 GB.
- Analysis has become the bottleneck and data size is an issue.
  - 2,000,000 genomes ≈ 1 Exabyte (1,000,000,000,000 MB)
  - Cost to sequence 1 genome less than $5,000 and falling fast.
  - Cost to analyze 1 genome is approx. $100,000 and rising.
- A key step towards Algorithm-assisted Personalized medicine is building Data Commons/Cloud analytics and the *Programmable* Networks & Communication Exchanges (SDXs) for high performance, flexible data transport.

usignite  geni

Infrastructure for Precision Medicine

# NCI Genomic Data Commons



- Harmonization and storage for the Nations Cancer Genomic Data GDC 1.6PB of cancer genomic data and associated clinical data.
- **Precision Medicine Enabled By Precision Networking**

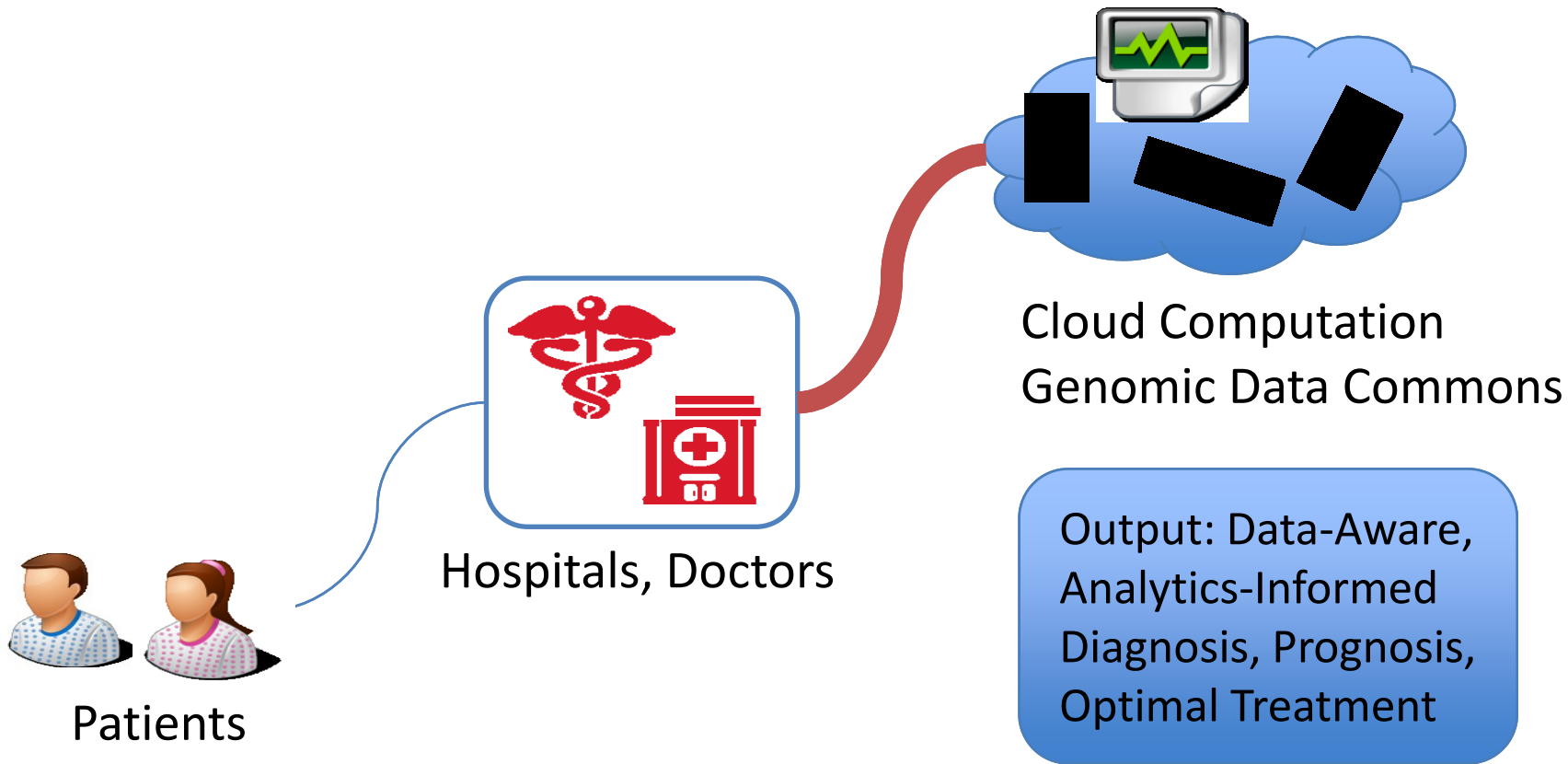# Bionimbus Protected Data Cloud



- Petabyte-scale, secure compliant biomedical cloud that interoperates with dbGaP controlled access data at NIH.

Infrastructure for Precision Medicine

# Future Vision: A Nationwide Virtual Comprehensive Cancer Center



Cloud Computation
Genomic Data Commons

Hospitals, Doctors

Patients

Output: Data-Aware,
Analytics-Informed
Diagnosis, Prognosis,
Optimal Treatment

# Opportunity: Close Integration of Research Workflows and Foundation Networks

- Opportunity: Using GENI To Develop Innovative Techniques for Extremely Close Integration of Research WorkFlows and Dynamic Programmable Network Resources, Enabling Precision Networking

- Network Foundation Architecture: GENI + Innovative Customized Software Defined Networking Exchange (SDX)

- For This Demonstration: Specifically To Meet The Requirements of Bioinformatic Workflows

# GENI Network <u>Programmability</u> Is Key

## GENI Programmability

- GENI Provides A Platform for Building the Required Precision Communication Services, Networks and Exchanges (SDXs)

- GENI OpenFlow Network

  - National Overlay Infrastructure Comprised of Shared VLANs Interconnected With OpenFlow Switches

  - FOAM/FlowVisor Enabling Sliced OF Switches (e.g., via Subnet, VLAN, Tunnel, etc)

- Discoverable, Integratable, Configurable, Programmable, Virtual Devices: Click Routers, OVS Switches, Mobile Devices, Instrumentation, and Other Resources

- Dynamic Edge Process Topology Design and Implementation

Precision Networks for Precision Medicine

# Biomedical Data Commons

Data Repository A (West Coast)

Data Repository B (South)

*Required Resources (Data & Tools) Are Highly Distributed*
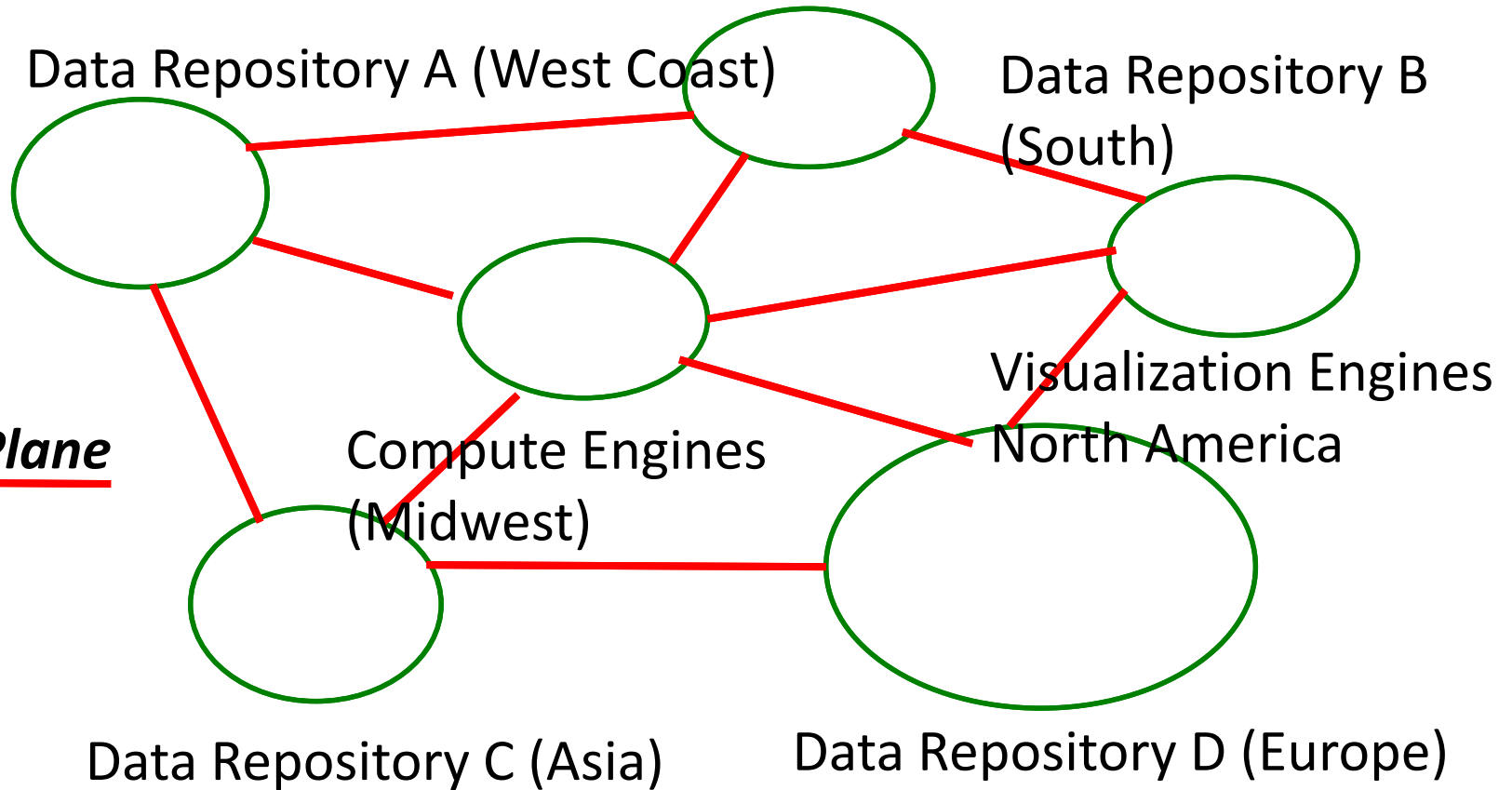
Compute Engines (Midwest)

Visualization Engines North America
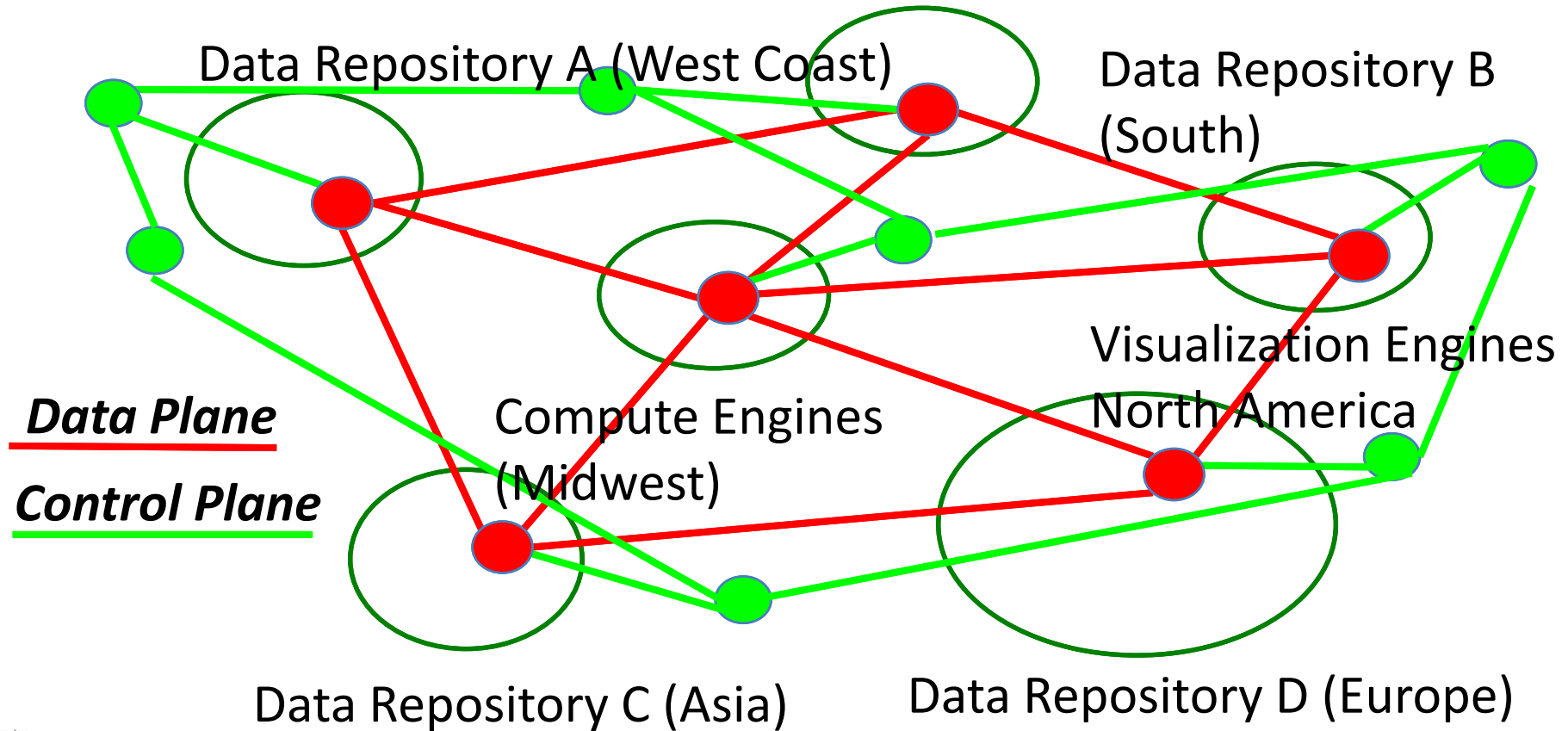
Data Repository C (Asia)

Data Repository D (Europe)
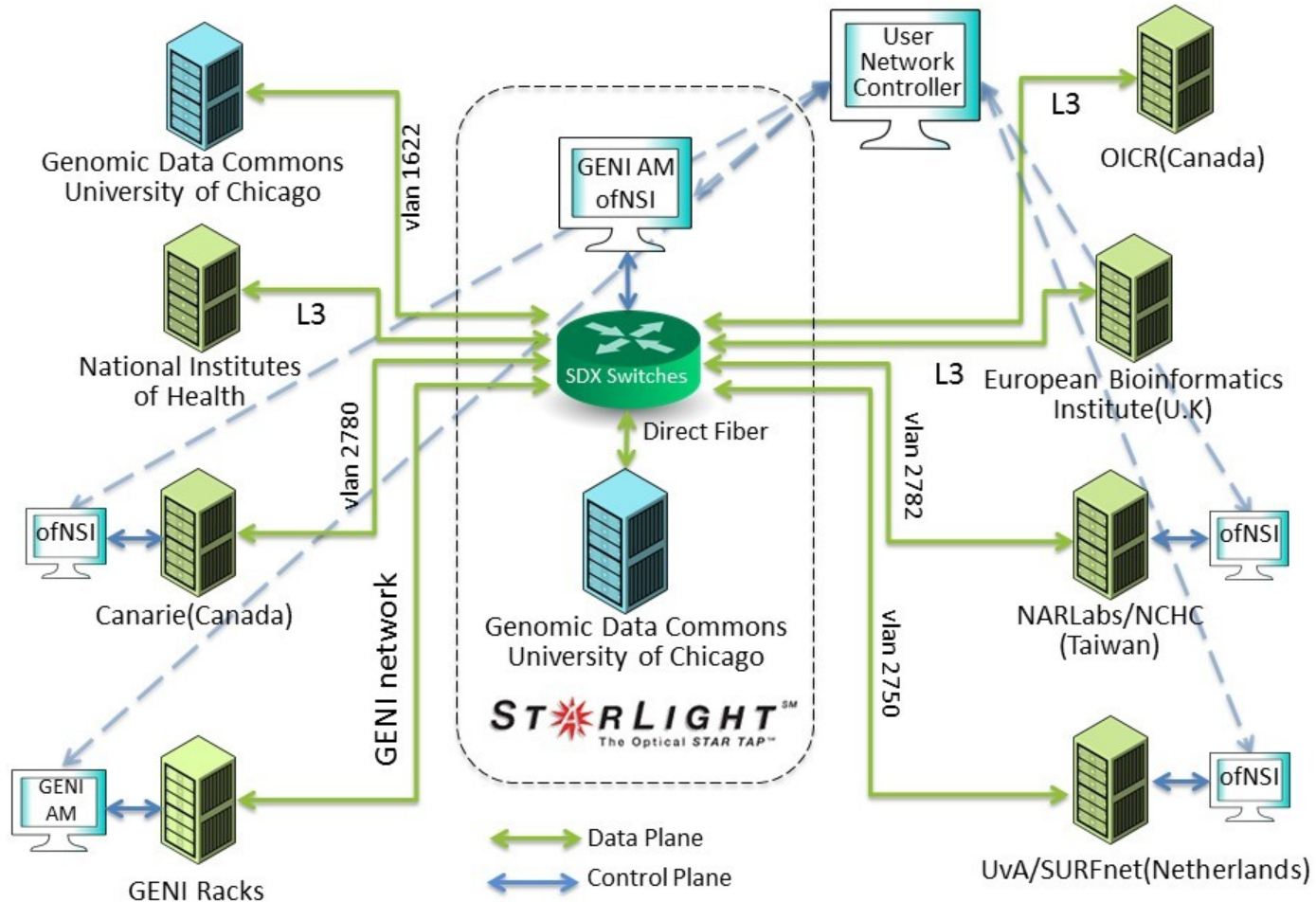
Biomedical Data Commons:
Flow Orchestration: Data Plane

Data Repository A (West Coast)

Data Repository B (South)

Data Plane

Visualization Engines North America

Compute Engines (Midwest)

Data Repository C (Asia)

Data Repository D (Europe)

# Biomedical Data Commons:
## Flow Orchestration: Control Plane + Data Plane

Data Repository A (West Coast)

Data Repository B (South)

Visualization Engines North America

*Data Plane*

*Control Plane*

Compute Engines (Midwest)

Data Repository C (Asia)

Data Repository D (Europe)

# GEC22 Bioinformatics SDXs Demo Network

# Today's Demonstration

- A) Dynamically Moving Core Data Files Among Multiple Sites Around the World Via StarLight SDX

- B) Moving RNA-seq Data Files From NCI (Bethesda, MD) and EBI (Hinxton, UK) Through SDX Switch/Routers to The University of Chicago.

  - Analysis By Comparison To Known Data Correlated To Drug Response.

  - Determine Possible Actionable Therapeutic Options.

  - Return Viable Treatment Options To the Originating Site.
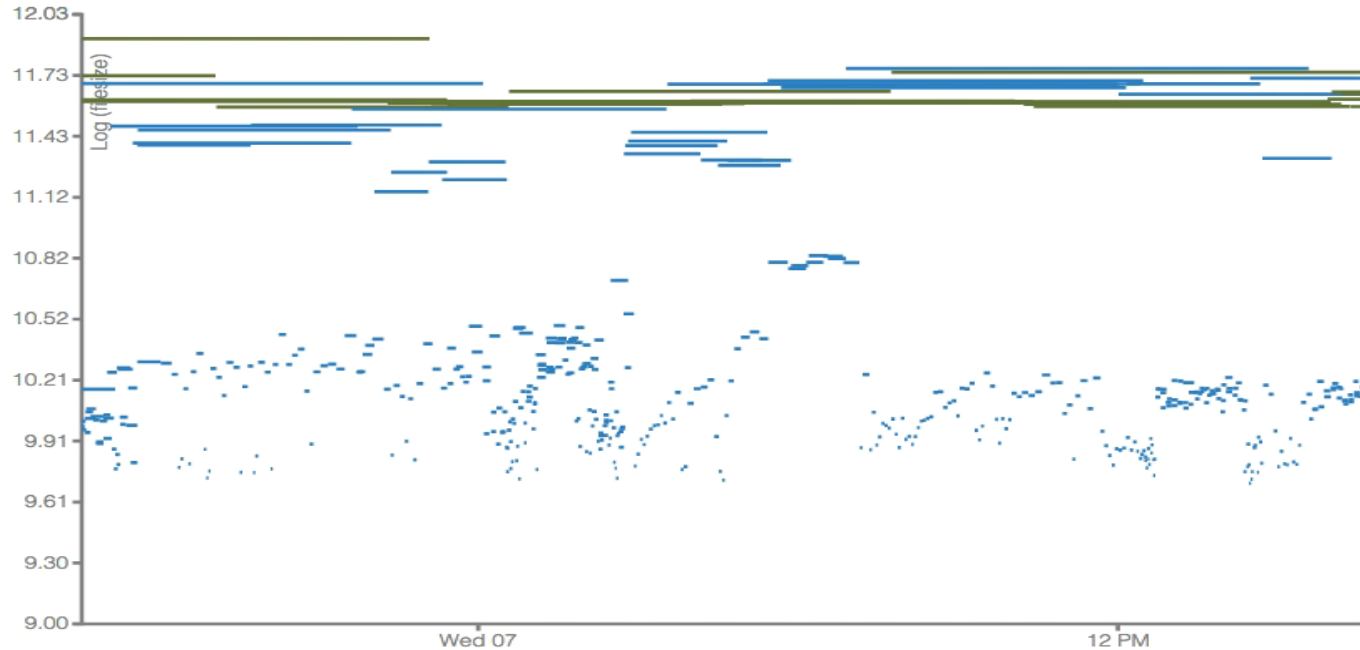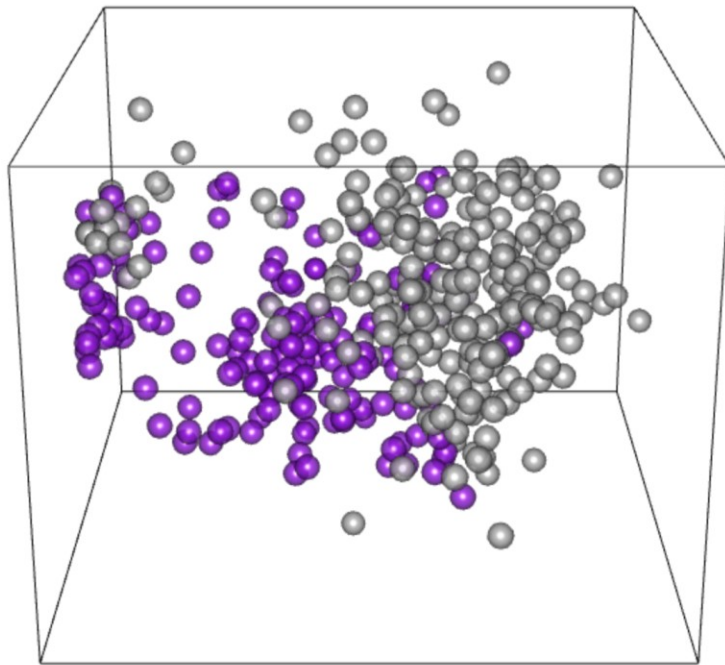
# Genomic Data Commons Data Transfer

# Gene Expression Clustering of Lung Cancers



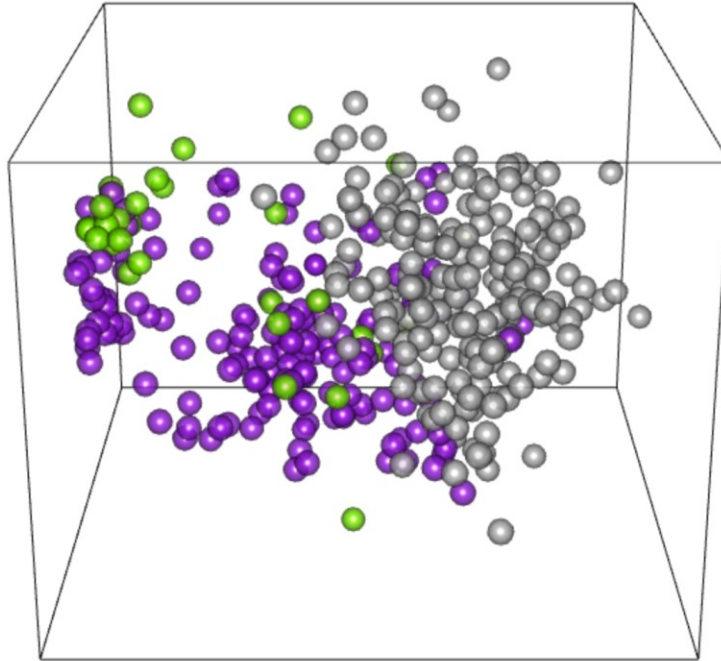Color By:

Original Diagnosis

⬜ Lung squamous cell carcinoma (LUSC)
🟪 Lung adenocarcinoma (LUAD)
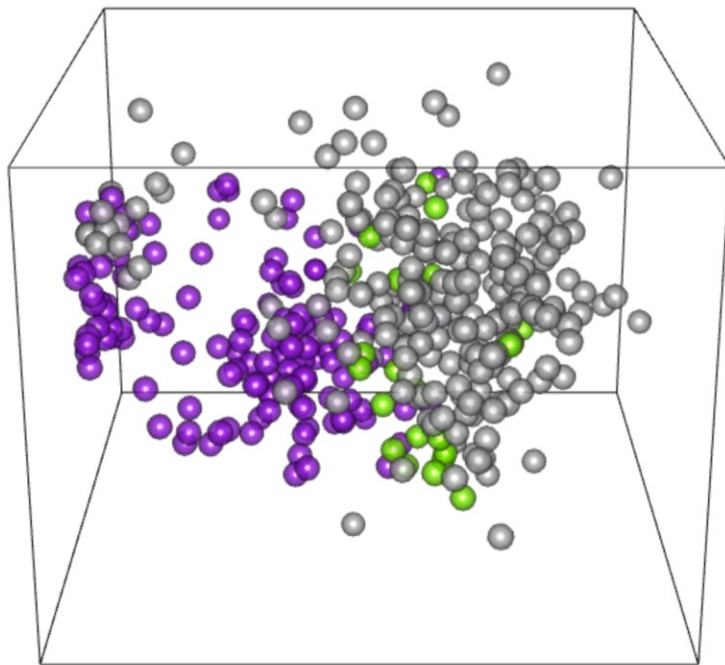
# Gene Expression Clustering of Lung Cancers

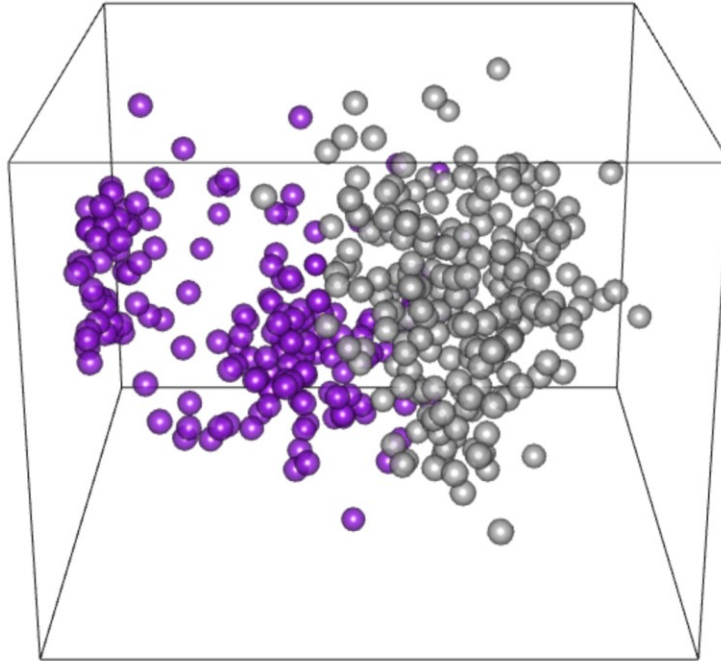# Gene Expression Clustering of Lung Cancers



**Color By:**

Potential Misdiagnosis of LUAD ▼

☐ Lung squamous cell carcinoma (LUSC)
■ Lung adenocarcinoma (LUAD)

# Gene Expression Clustering of Lung Cancers



**Color By:**

By Expression

☐ Lung squamous cell carcinoma (LUSC)
■ Lung adenocarcinoma (LUAD)

# Results

- Precision medicine requires data commons that scale to hundreds of petabytes scale, with programmable networks and data peering to support data sharing.

- Speed discovery and support analytics-driven healthcare to recommend treatment.

- Large Scale Data Analysis and Dynamic Pipelines For Workflows Are Essential For Determining Optimal Results.

# Summary and Future

- <u>What you saw</u>: An innovative approach to advanced knowledge discovery and medical treatment*: **Precision medicine being supported by precision networking**

- <u>Why GENI/US Ignite is important</u>: Precision mapping of communication services to BI workflow requirements across the world using advanced analytics, the Genomics Data Commons & a programmable dynamic SDX

- What happens looking forward, for the application and its integration with GENI:

  - A) Further development/refinement of basic capabilities

  - B) Transition to *actual production services*

  - **C) The Genomics Data Commons and Bionimbus Protected Data Cloud is Being Developed As a Key Production Knowledge Discovery/Transformational Medical Treatment Facility**

# Using GENI To Invent the Future...

## Thank You!