# GNARLY
# Data_Waves

PRESENTED BY 🐬dremio

EPISODE 19

# Data Mesh In Practice: Accelerating Cancer Research with Dremio's Data Lakehouse

**Arfath Pasha**
Sr. Software Engineer, MSK

**Tony Truong**
Sr. Product Marketing Manager, Dremio

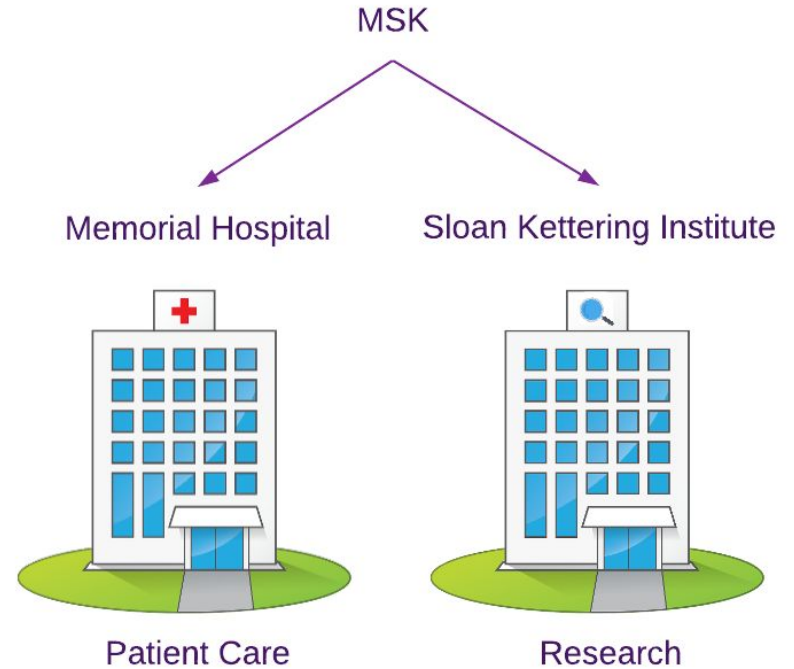📅 May 31 at 8AM PST | 11AM EST | 4PM GMT

🐬dremio | 🌲 Memorial Sloan Kettering Cancer Center
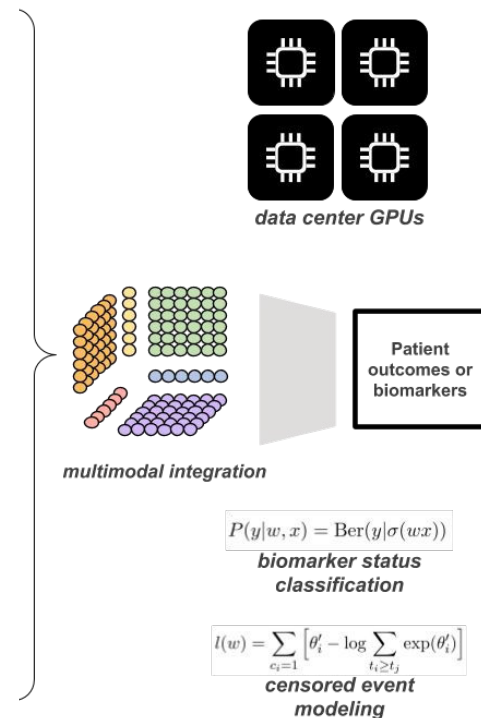
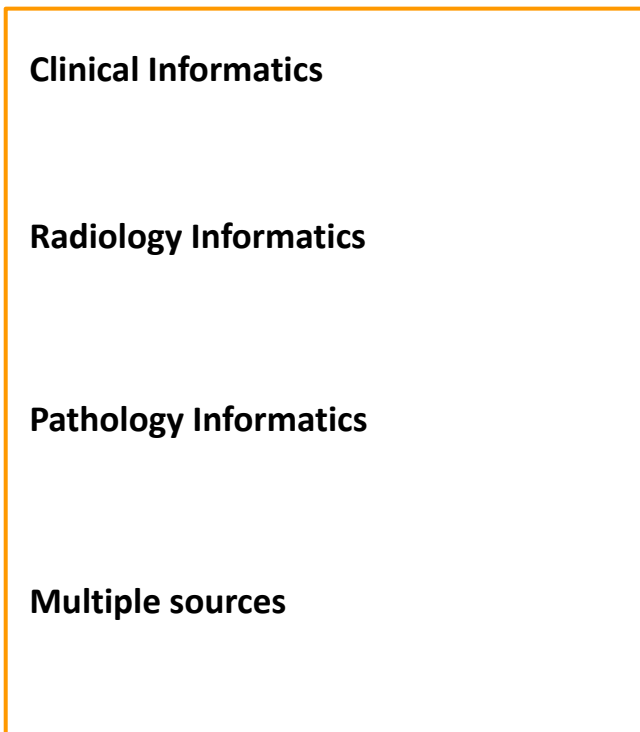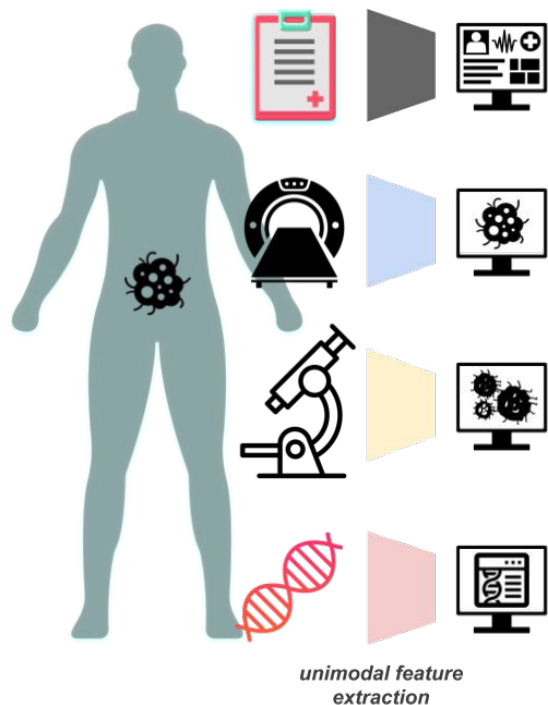# Memorial Sloan Kettering Cancer Center

Founded in 1884

Treats more than 400 cancer subtypes

- 20k inpatient and 700k outpatient visits
- > 1800 research protocols

# Research at MSK builds on top of existing resources in various departments at MSK



**Clinical Informatics**

**Radiology Informatics**

**Pathology Informatics**

**Multiple sources**

*unimodal feature extraction*

*data center GPUs*

*multimodal integration*

Patient outcomes or biomarkers

$$P(y|w,x) = \text{Ber}(y|\sigma(wx))$$

*biomarker status classification*

$$l(w) = \sum_{c_i=1} \left[ \theta_i' - \log \sum_{t_i \geq t_j} \exp(\theta_i') \right]$$

*censored event modeling*

**Challenges: regulatory, data governance, and technical barriers**
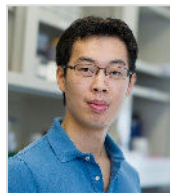
3

# Engineering for Research

Our goal:
- Build the right (scientific data management+compute) system

We need to be:
- An infrastructure team
- A data products team
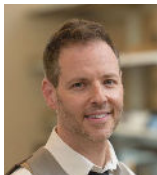- An analysis team

We are:
- Research Software Engineers

**Raymond Lim**
**Senior Engineer**

**Armaan Kohli**
**Engineer**

**Darin Moore**
**Engineer**

**Benjamin Gross**
**Lead Engineer**

**Anika Begum**
**Project Coordinator**

**Arfath Pasha**
**Senior Engineer**

# What is a Data Mesh?

A modern approach to data management that emphasizes distributed ownership and governance of data within domains, who then build, manage, and share data products across the organization.

"Dehghani, Z. (2022). Data Mesh: Delivering data-driven value at scale. O'Reilly Media, Incorporated.

# Data Challenges in Research

- High dimensional data (clinical, genomic, radiology, pathology, etc.)

- Teams with diverse skill-sets (engineers, scientists, pathologists, physicians, administrative staff, etc.)

- Highly iterative in nature (need for data versioning)

- Unstructured and structured data

  - Binary large objects and related tabular data
  - Long datasets (100-1 billion rows), wide datasets (10-1000 columns)
- Messy data (correctness, completeness issues)

- Siloed data (many data marts)

- Privacy

# Data Management Before Dremio



**Binary Large Objects**
- Pathology Slides
- Radiology Scans
- Genomics files ...

**Tables**
- Pathology metadata
- Genomics metadata
- Radiology metadata ...
- Data Warehouse

Intermediaries Request

Intermediaries Catalog

Intermediaries ETL

Intermediaries Deidentify

Intermediaries Gain Access

Intermediaries ETL

Intermediaries Gain Access

**Data Lake**
- Filestore
- Database

Data Managers

copies circulated via email

Physicians

Administrators

Data Scientists

Software Engineers

Researchers

Data Curators

**Time (data is made available to researchers in weeks to months per project)**

8

# Solution Considerations

### Architecture

- On-premises deployment
  - Data mesh
  - Query engine
- A no-copy data architecture

### Decentralized Data Management

- Eliminate siloed ETL pipelines, provide faster access to data
- Documentation support - data sheets for data sets
- Simple, mature governance model

### People

- Easy interface for all data consumers

# Why Dremio

Easy barrier to entry (satisfied our 1 hour rule for evaluating a new technology)

Supports on-prem deployment with path to cloud

Access control (unified semantic layer)

Data democratization (almost spreadsheet like interface plus connection to Tableau)
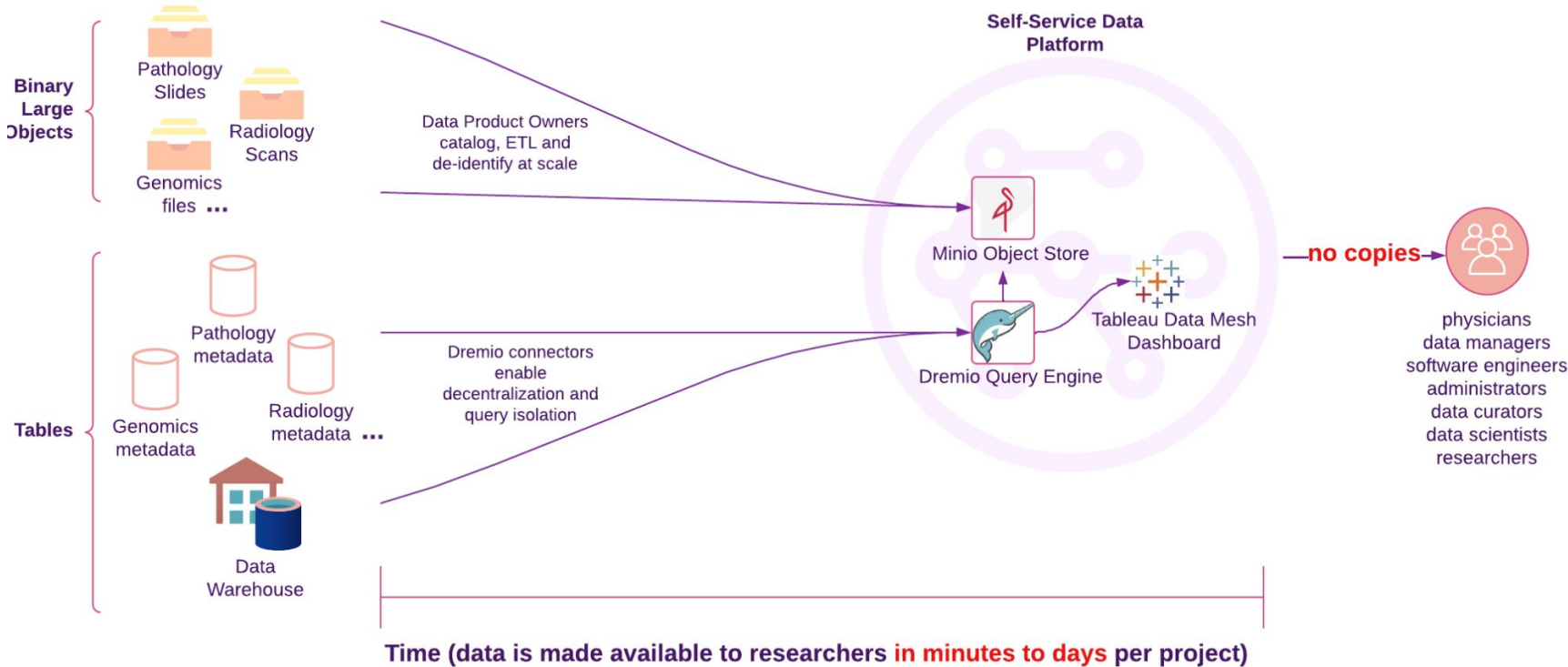
No copies (no more emailing copies)

Low code / no code (easy data inspection/curation/integration without pandas code)

Performance (horizontally scalable; Arrow Flight access)

Datasheets for datasets (support for documentation through catalog wiki)
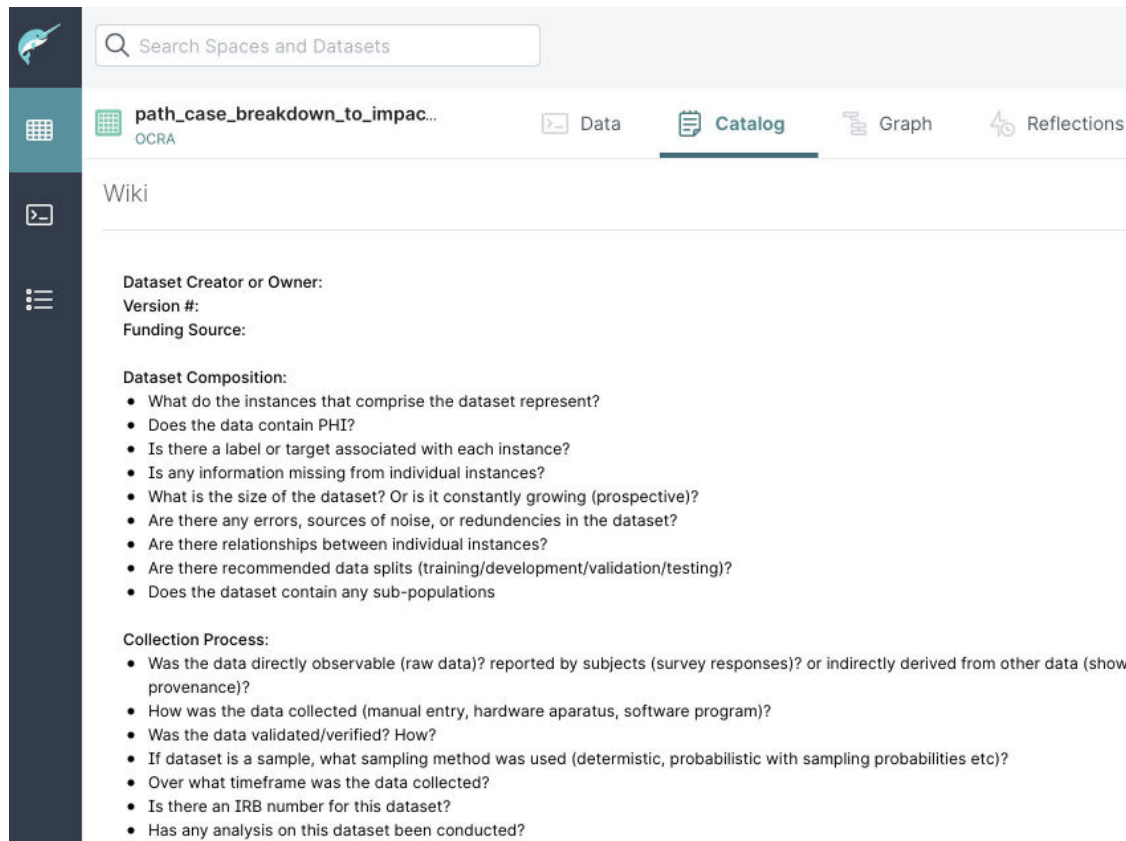
Data versioning (Iceberg and Nessie)

# How We Are Decentralizing Data Management With Dremio



**Binary Large Objects**
- Pathology Slides
- Radiology Scans
- Genomics files ...

**Tables**
- Pathology metadata
- Radiology metadata ...
- Genomics metadata
- Data Warehouse

Data Product Owners catalog, ETL and de-identify at scale

Dremio connectors enable decentralization and query isolation

**Self-Service Data Platform**

Minio Object Store

Dremio Query Engine

Tableau Data Mesh Dashboard

no copies →

physicians
data managers
software engineers
administrators
data curators
data scientists
researchers

**Time (data is made available to researchers in minutes to days per project)**

11

# Lessons Learned Along the Way

✓ Building trust between data product owners and consumers with data mesh

✓ Reduced data product delivery time by eliminating siloed ETL with Dremio

✓ User data copies eliminated, now easier to track and share across domains

12

# Datasheets for Datasets



Reference: https://arxiv.org/abs/1803.09010

# What's Next - Data Mesh Evangelism

Federated Computational Governance

- Standardizing data governance between domains and make decisions about how data is used and shared across the enterprise.
- Need for 'tenant' admins for multi-tenancy



14