# Counting Your Assets
## Don't Forget About Data!

James (Jim) Wilgenbusch

Director of Research Computing:
- Minnesota Supercomputing Institute
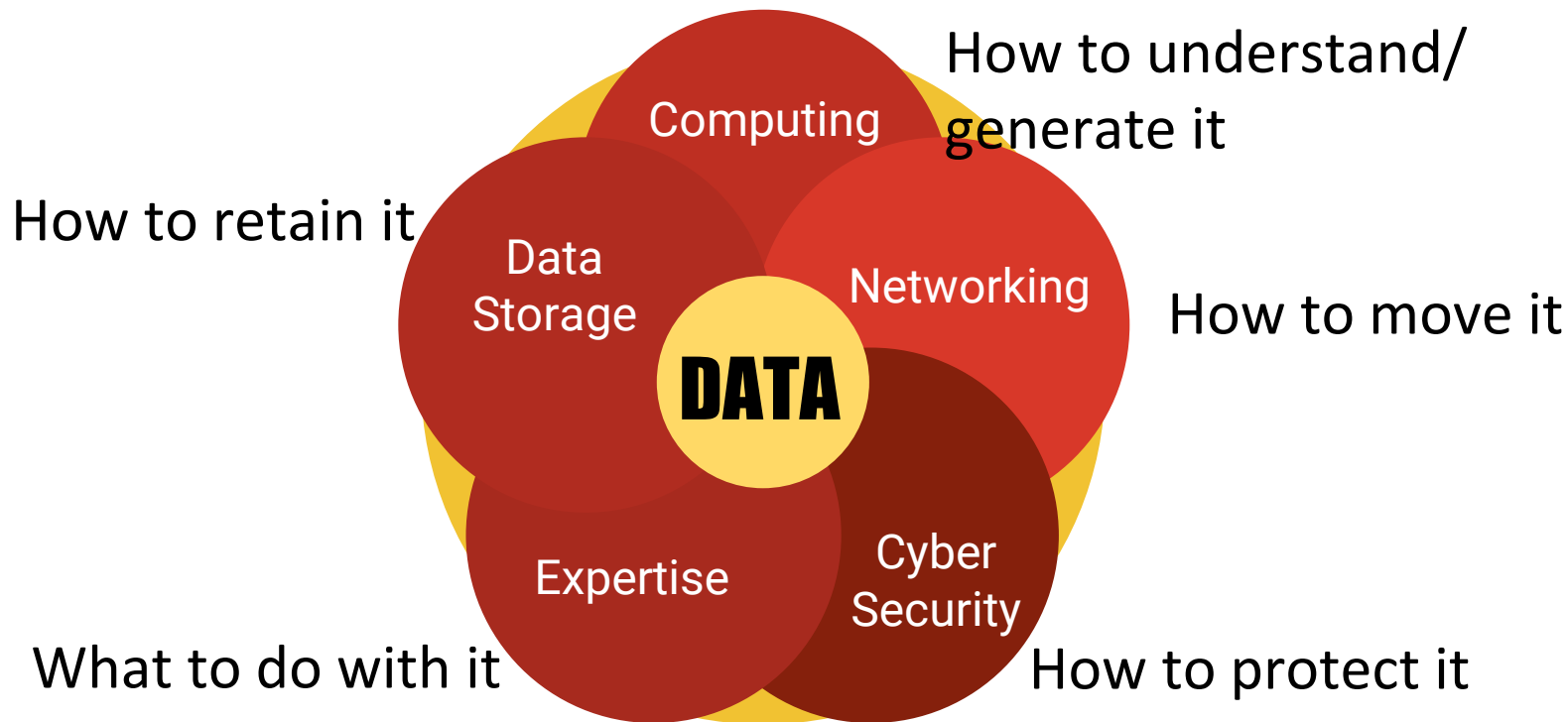- Minnesota Informatics Institute
- U-Spatial

*Office of the Vice President for Research*

UNIVERSITY OF MINNESOTA

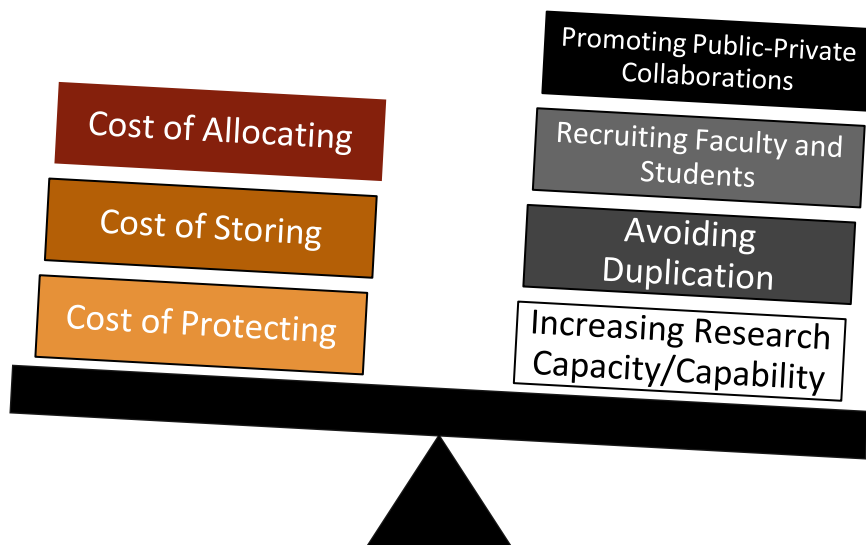**Driven to Discover**®

# Key Components of Research CI



How to understand/ generate it

How to retain it

How to move it

What to do with it

How to protect it

Computing

Data Storage

Networking

DATA

Expertise

Cyber Security

UNIVERSITY OF MINNESOTA

# Paradigm Shift
## Data as a Liability → Data as an Asset

**Liabilities**

**Assets**

Cost of Allocating

Cost of Storing

Cost of Protecting

Promoting Public-Private Collaborations

Recruiting Faculty and Students

Avoiding Duplication

Increasing Research Capacity/Capability

# Challenges are Real

- Breath of research data is enormous

  - Socioeconomic to stellar evolution

- Size of some data sets are daunting
  - 100 TB to Multi-PB data sets are common

- Data use agreements are onerous
  - Lots of time to setup and penalties for violations

Minnesota Supercomputing Institute

UNIVERSITY OF MINNESOTA

# Data as an Asset:
## Key Ingredients

- Sustained Research Infrastructure
- Well Trained Cyberinfrastructure Professionals Cyberpractitioners

# MSI Computing and Data Storage Assets

## Batch High Performance Computing

- Two Supercomputers
- 25,000 CPU Cores
- 230,400 GPU CUDA Cores
- 100 TB Memory
- Infiniband Network

## Big Data Storage & Analysis

- 6 PB Primary High Performance
- 3 PB Second Tier
- 30 PB Archive Tape Library

## Interactive & Cloud Computing

- Citrix VDI for Windows
- DCS Nice for Linux Desktops
- OpenStack for Secure Cloud
- 100 Gbps Campus Research Network
- Regional & National Optical Networks

## Web Portals & Databases

- Galaxy for Multi-omics
- Jupyter Hub
- Custom Interfaces & Applications

# Office of the Vice President for Research

## Research Computing

### Minnesota Supercomputing Institute

### Univ. of Minnesota Informatics Institute

### U-Spatial

**Scientific Computing Solutions**
6-FTEs

- Code Optimization
- Workflow & Platform Dev
- Project Leadership
- Dedicated Grant Support
- In Depth User Support, Consulting, and Troubleshooting

**Research Informatics Solutions**
12-FTEs

- Life Sciences Computing
- Workflow & Platform Dev
- Informatics Education
- Informatics Research
- Project Leadership
- Dedicated Grant Support
- In Depth User Support, Consulting, and Troubleshooting

**Application Development Solutions**
6-FTEs

- Web development
- User Dashboard Development
- Systems Programming
- Custom App Dev
- Project Leadership
- Dedicated Grant Support

**User Gateway Group**
5-FTEs

- Helpdesk Lead
- Onboarding and User Training
- Communications
- Outreach
- Administrative functions

**Advanced Systems Operations**
11-FTEs

- Systems Support
- Hosted Services
- Benchmarking
- Project Leadership
- Limited Dedicated Grant Support
- In Depth User Support and Troubleshooting

**Dedicated Solutions Groups**

**Core Operations Support**

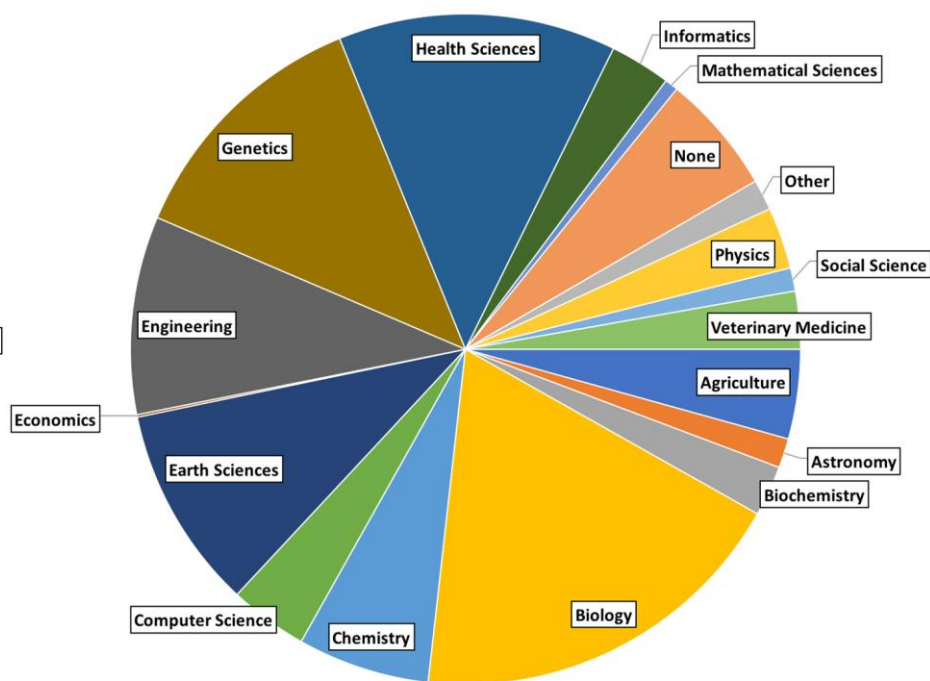# In 2018: 888 User Groups, 4,555 Active users



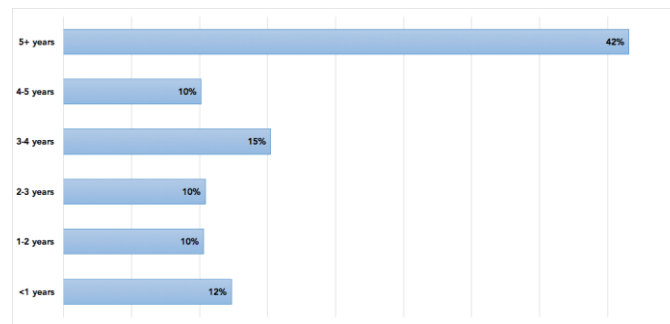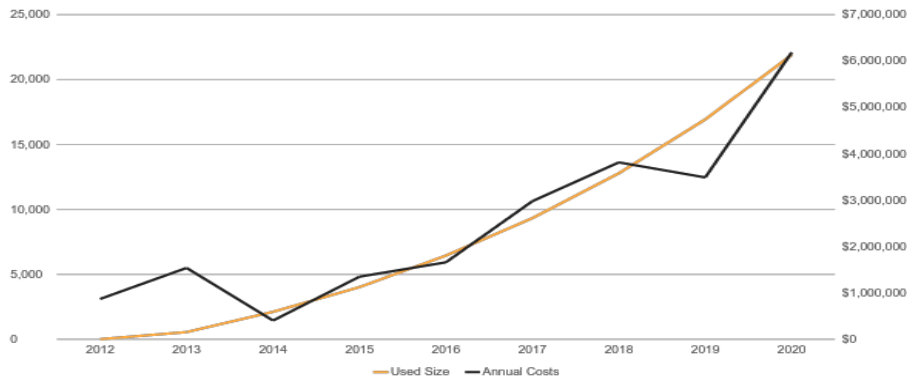**Biggest increasing in Life Sciences**

# Resource Utilization by Group



**CPU Hours
150 Million Total**
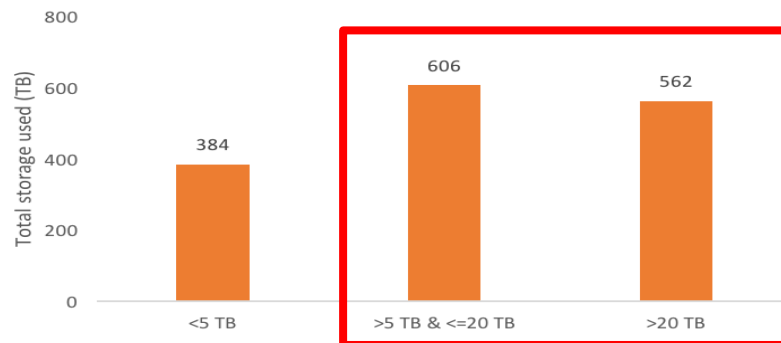
**HPC Storage
1.5 PetaBytes**

Minnesota Supercomputing Institute

UNIVERSITY OF MINNESOTA

# What's the Problem?



Number of Groups by Allocation Tier

Total Storage Used by Allocation Tier

Minnesota Supercomputing Institute

UNIVERSITY OF MINNESOTA

# Data as an Asset:
## Key Ingredients

- Sustained Research Infrastructure
- Well Trained Cyberinfrastructure Professionals Cyberpractitioners
- **Good Storage Governance that Spans Institutional Reporting Lines**

# It Doesn't Have to Be Scary

**Assemble Stakeholders**

AHC

MSI

OIT

Centers

Others

CSE

Storage Redesign and Restructure Committee (SRRC)

**Define areas that need attention**

Service Analysis

Campus Needs

Standards and Operating Procedures

Education and Storage Champions

# Outcomes

## Establish a University-wide Storage Council



### Council is Charged to:

- Develop a Storage Champion Program
- Enhancing Website Infrastructure
- Enhancing User Training and Onboarding
- Collaborating and Sharing Internal Knowledge Articles
- Promoting Marketing and Communication

# Outcomes



Storage Champions Network
Kickoff meeting
Septemeber 19, 2019
z.umn.edu/scn



Storage Selection Tool
z.umn.edu/storage-selection-tool

Minnesota Supercomputing Institute
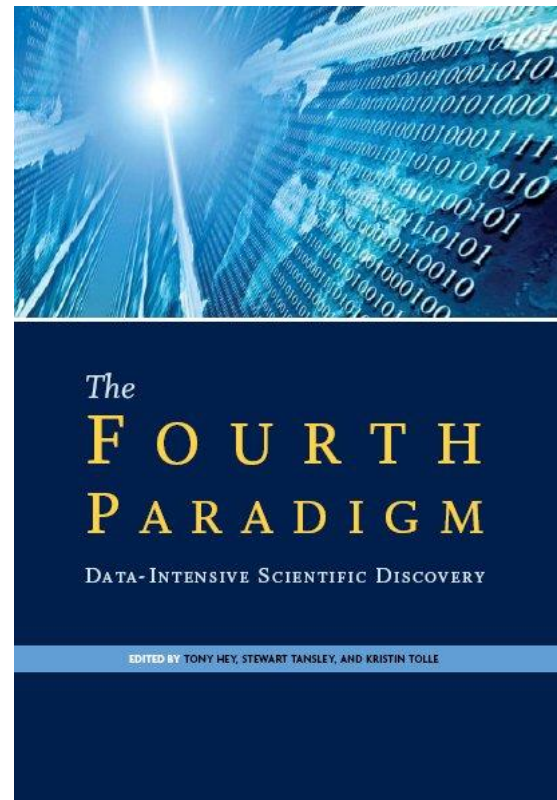
UNIVERSITY OF MINNESOTA

# Data as an Asset:
## Key Ingredients

- Sustained Research Infrastructure
- Well Trained Cyberinfrastructure Professionals -- AKA Cyberpractitioners
- Good Storage Governance that Spans Institutional Reporting Lines
- **Tools to Make Data Interoperable and to Facilitate Analyses and Sharing**

UNIVERSITY OF MINNESOTA

"Today, the tools for capturing data both at the mega-scale and at the milli-scale are just dreadful."

Jim Gray, 2007



The FOURTH PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

UNIVERSITY OF MINNESOTA

# The Challenge

UNIVERSITY OF MINNESOTA

GEMS

**Data-Driven Agricultural Innovation**
GENETICS · ENVIRONMENT · MANAGEMENT · SOCIOECONOMICS

A novel data sharing and analysis platform to enable public-private research collaborations for innovation in agricultural production and other domain areas.

**G** ✖ **E** ✖ **M** ✖ **S**

Genomics  Environment  Management  Socio-Economics

Time ✖ Space

# Our Specific Contributions

**GEMShare** ™

- Smart sharing -- Enables data providers to control who sees what, and when

- Data Versioning – Ensures reproducibility and ability to roll back from changes

- Supports -- open, private, and pooled data

- Beyond data -- Enables sharing of tools and workflows too

**GEMSTools**™ is an ever-expanding suite of web-based and command-line analytical tools designed to:

- Cleanup messy (meta-)data

- Intelligently impute missing data

- Enable data interoperability

- Apply advanced analytic methods to genomic, environmental, management and socio-economic data

# GEMS in Action







https://agroinformatics.org/

# GEMS™: Enabling External Partnerships Through the International AgroInformatics Alliance (IAA)

- CGIAR Big Data Platform
- Embrapa, Brazil
- Pepsico
- G2F (Genomes to Fields)
- Diversity Arrays Technology (DArT/KDDART)
- CIAT (cassava, edible beans, forages, rice, )
- University of Adelaide
- Oat Global
- Stellenbosch University, South Africa
- CIMMYT (corn, wheat, socio-economics, genetic resources, IT )

---

- GRDC (Grains Research Development Corporation), Australia
- MN Department of Agriculture
- PPIRC, Phenotyping and Imaging Center, Canada
- CIP (potatoes, sweet potatoes)
- ICRISAT (sorghum, millet, chickpeas, groundnut)

# Thank You

# Questions?

**jwilgenb@umn.edu**

Minnesota Supercomputing Institute