# Data Sharing and Analysis Enabling Data Driven Agricultural Innovation While Respecting IP

Kevin Silverstein, PhD, GEMS Operations Manager

GEMS led by: Philip Pardey, Jim Wilgenbusch and Kevin Silverstein
College of Food Agricultural and Natural Resource Science, CFANS
Minnesota Supercomputing Institute, MSI
University of Minnesota

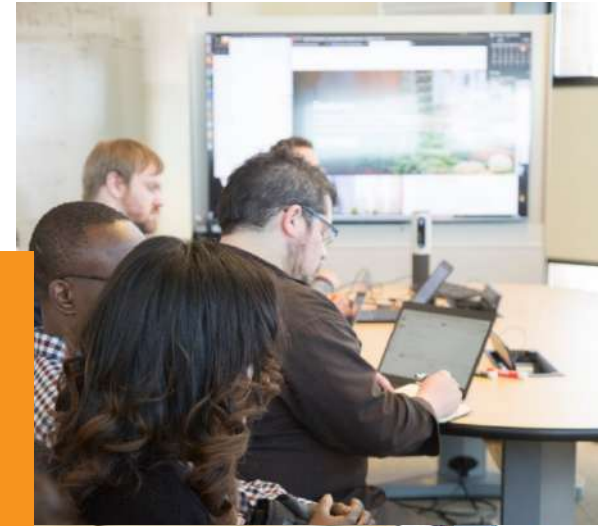PAG XXVII Workshop Connecting Crop Phenotype and Genotype Data
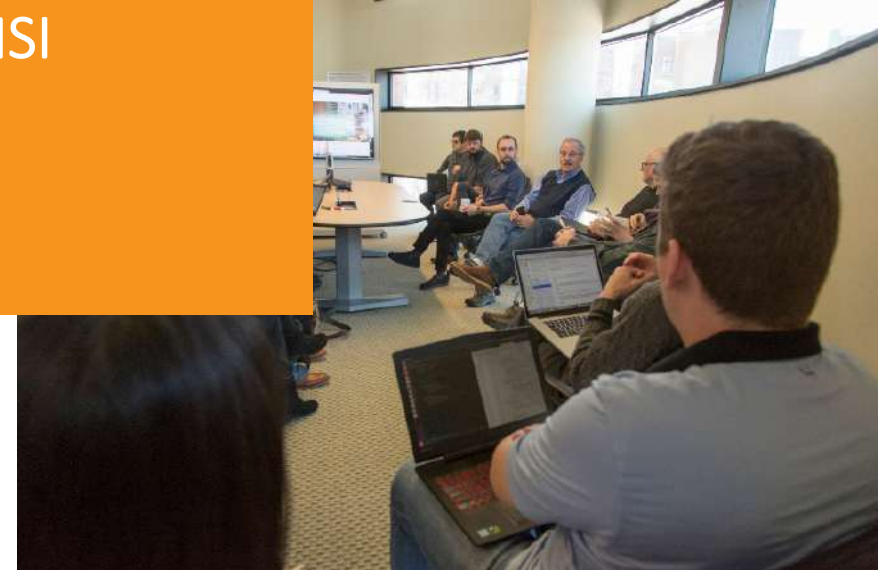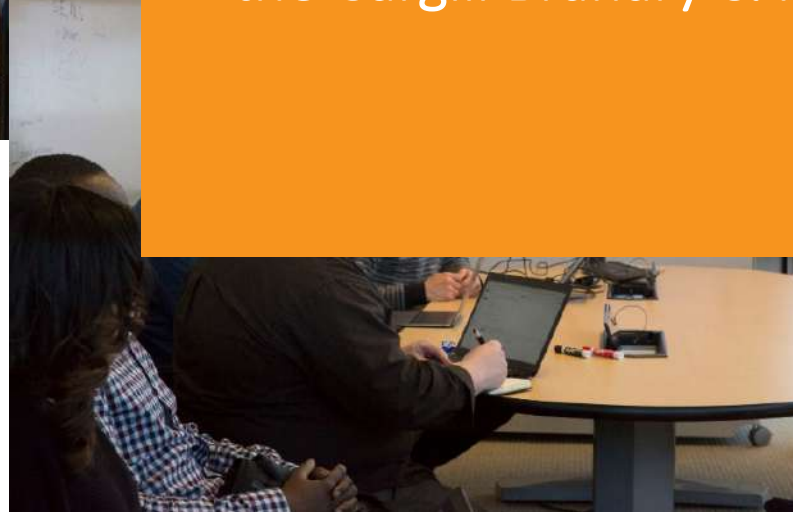
January 16, 2019

# What is G.E.M.S?

A novel data sharing and big data analytical platform that enables public-private research collaborations for innovation in food and agricultural production, and other domain areas



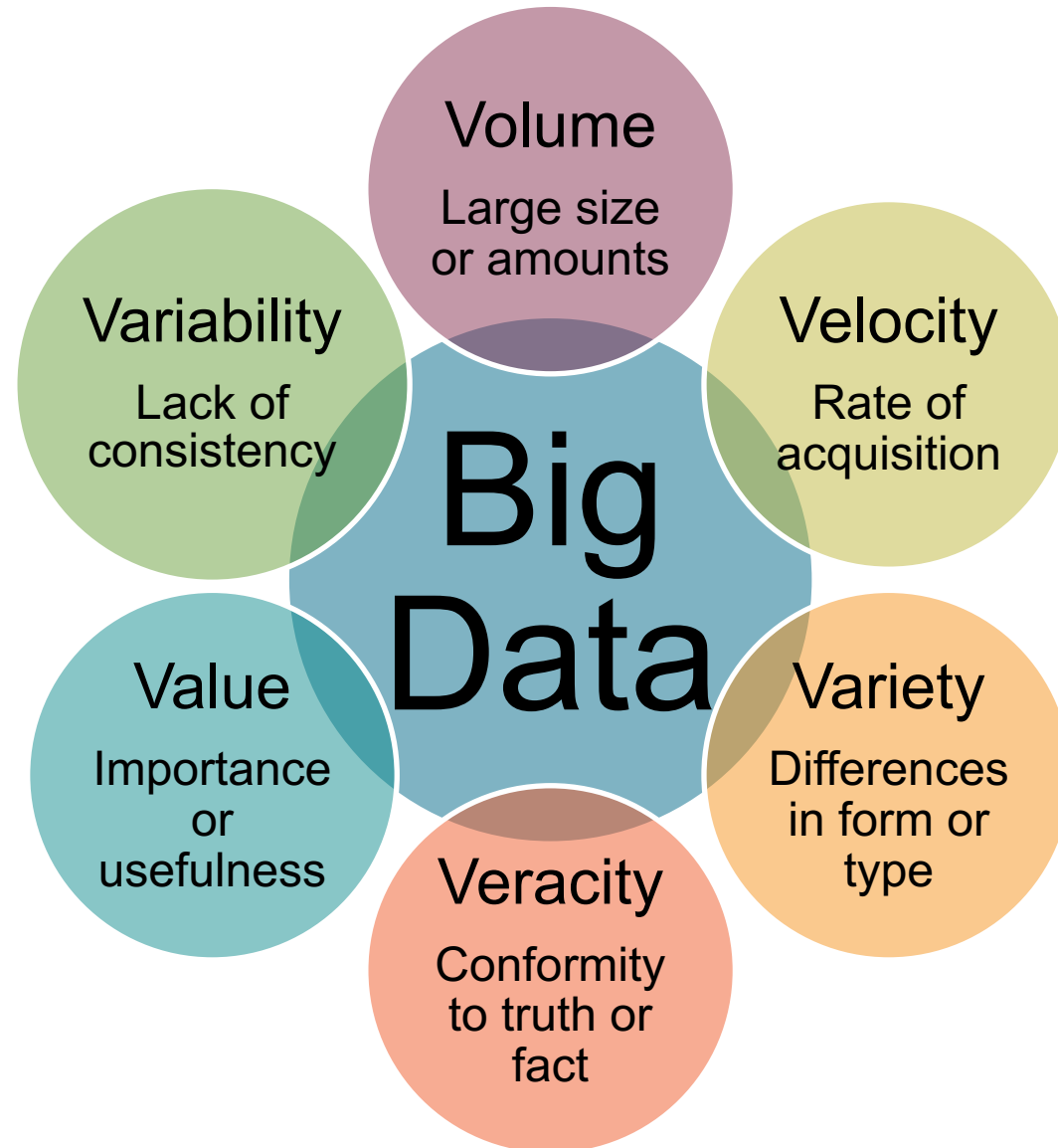Genomics     Environment     Management     Socio-Economics

Time ✖ Space

# The G.E.M.S Team (more than 20 brains strong!)

- Bi-Weekly build meetings
- Weekly technical meetings
- Numerous ad hoc consultations in the Cargill Branary & MSI

**GEMS**

# Big Data Challenges in Food & Agriculture

# Realizing the Big Data Revolution

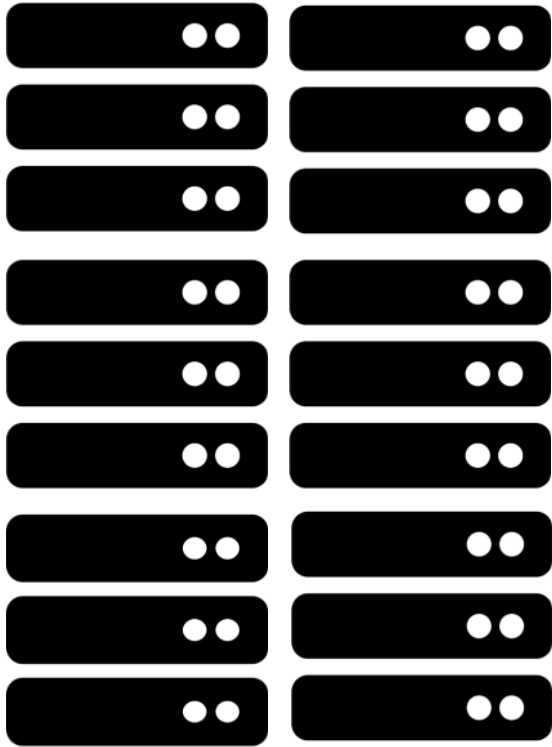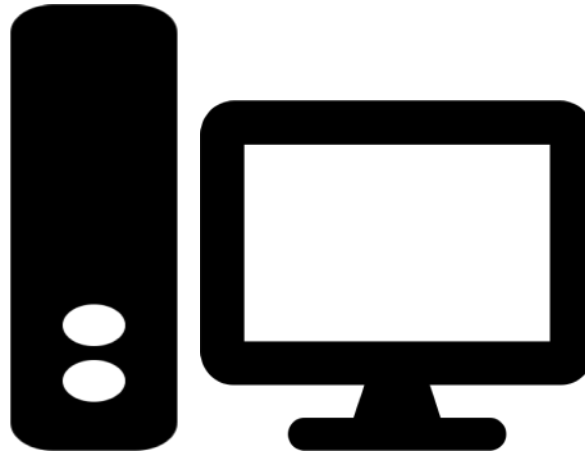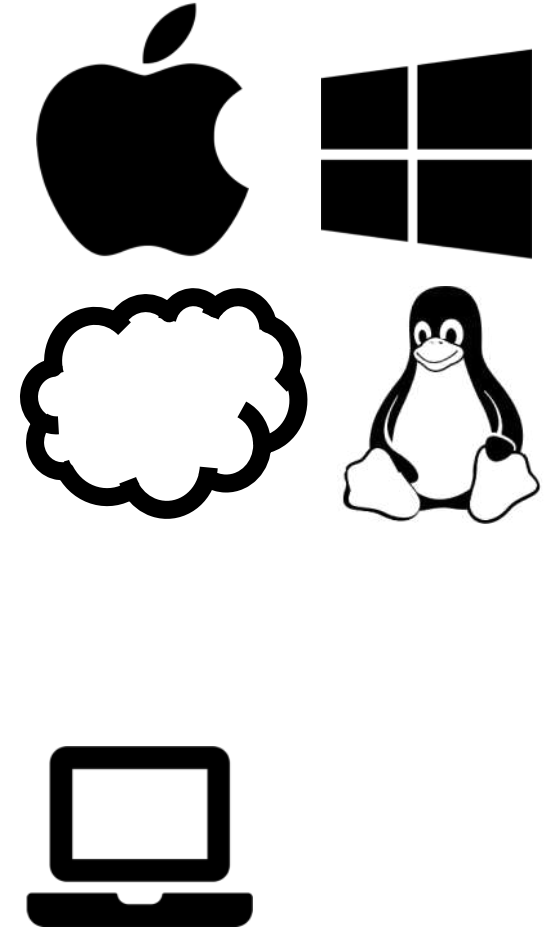| Get the data to the tool or get the tool to the data | Reconcile file formats, units, vocabularies, languages, and ontologies | Access to complex software and ability to replicate analyses | Facilitate complex partnerships and respecting data ownership and privacy |
|---|---|---|---|
| Data Transfer | Data Interoperability | Data Analysis | Data Sharing |

GEMS

# Data Transfer – Platform Portability

**GEMs on Clusters**
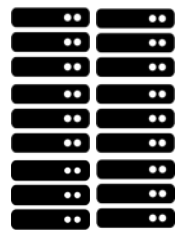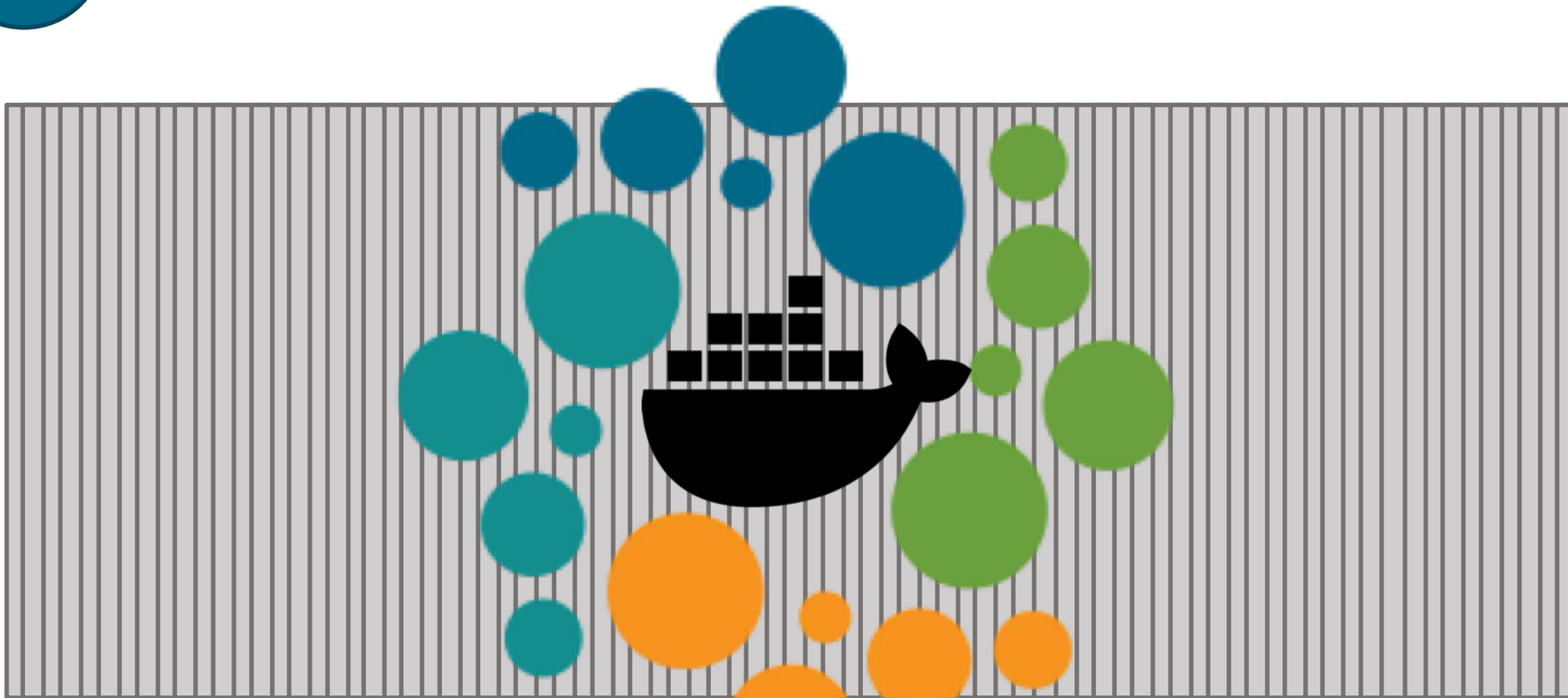
**GEMs on Workstations**

**GEMs on Laptops**

# Data Transfer – Containerize the Platform

# Data interoperability: column metadata

# Data interoperability: spelling correction

# Data interoperability - ontology matching

# Data Interoperability via Automated (units & terms) Standardizations

**Total Phosphorus**

207 lb/A

46lb per acre

46 pound/A

22 lbs

68 kg/ha

54lbs/acre

55.5 lbs P per Acre

80lb/acre

None

40 pounds

none applied

192 lbs;17-Apr-14

 ...

| | | | | TABLE 7 Conversion Table |
|---|---|---|---|---|
| Length | 1 yd | = | 0.9144 m | |
| | 12 in. | = | 1 ft | |
| | 5280 ft | = | 1 mile | |
| | 1 m | = | 3.281 ft | |
| | 1 in. | = | 0.0254 m | |
| Time | 60 sec | = | 1 min | |
| | 3600 sec | = | 1 hr | |
| Mass | 1 lbm | = | 0.4535 kg | |
| | 2.205 lbm | = | 1 kg | |
| | 1 kg | = | 1000 g | |
| Area | 1 ft$^2$ | = | 144 in.$^2$ | |
| | 10.764 ft$^2$ | = | 1 m$^2$ | |
| | 1 yd$^2$ | = | 9 ft$^2$ | |
| | 1 mile$^2$ = | | 3.098 X 10$^6$ yd$^2$ | |
| Volume | 7.48 gal | = | 1 ft$^3$ | |
| | 1 gal | = | 3.785 l (liter) | |
| | 1 l | = | 1000 cm$^3$ | |

GEMS

# Data analysis – ad hoc investigation

# Data and Tools Sharing

**GEMShare**™

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

- Supports open, private and pooled data

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

**GEMSTools**™ is an ever-expanding suite of 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**

# Your Data, Your Tools, Your Choice!

## Technical Security

- Staff trained to handle sensitive data
- Ability to move the platform to the data
- Analyses run in isolated containers
- Servers hosted in a robust and secure data center
- Data encrypted at rest and in flight
- Systems constantly monitored

## Legislative/Legal Privacy



2018 Minnesota Session Laws

BE IT ENACTED BY THE LEGISLATURE OF THE STATE OF MINNESOTA:

Section 1. Minnesota Statutes 2016, section 13.643, subdivision 7, is amended to read:

Subd. 7. **Research, monitoring, or assessment data.** (a) Except as provided in paragraph (b), the following data created, collected, and or maintained by the Department of Agriculture or the University of Minnesota during research, monitoring, or the assessment of farm practices and related to natural resources, the environment, agricultural facilities, or agricultural practices are classified as private or nonpublic:

(1) names, addresses, telephone numbers, and e-mail addresses of study participants or cooperators; and

(2) location of research, study site, and global positioning system data; and

(3) data created, collected, or maintained by the University of Minnesota for inclusion on an agricultural data analysis platform maintained and hosted by the University of Minnesota that identify or could identify an individual or business.
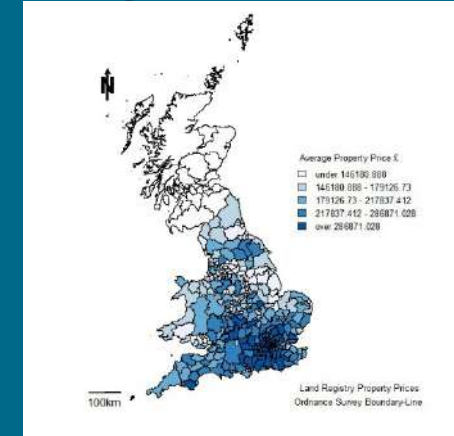
Law came into effect
August 1, 2018

# Data Sharing: Data Fuzzing



De-identifying

Aggregating

GEMS

# Use Cases
## Data to Actionable Information



**The GEMS Cluster**
Now up and running in UMN's supercomputer facility

Jim Wilgenbusch , Director Minnesota Supercomputer Institute (left), chatting with Jan Greyling University of Stellenbosch, South Africa, local GEMS coordinator

# Agroinformatics Support for G2F


the Genomes to Fields initiative



**2017 Academic & Federal Institutions**

Arkansas State University (2016–2017)
Clemson University (2016–2017)
Colorado State University (2017)
Cornell University (2014–2017)
Iowa State University (2014–2017)
Kansas State University (2015–2016)
Michigan State University (2016–2017)
Mississippi State University (2017)

North Carolina State University (2014–2017)
Ohio State University (2015–2017)
Pennsylvania State University (2015–2017)
Purdue University (2014–2017)
South Dakota State University (2015)
Texas A&M University (2014–2017)
University of Arizona (2015 & 2017)
University of Delaware (2014–2017)

University of Georgia (2014–2017)
University of Guelph (2014–2017)
University of Illinois (2014–2017)
University of Minnesota (2014–2017)
University of Missouri (2014–2017)
University of Nebraska (2014–2017)
University of Wisconsin (2014–2017)
USDA-ARS (2014–2017)

- Standardizing nomenclature, units etc

- Outlier detection

- Data interoperability (weather, soil, management, phenotypic measurement)

- Pilot linkage of field measurements from tablet to GEMSTools

- Manage data distribution among G2F partners and to the world

- Data mining and other (predictive) analytics

23 States, 37 experimental sites

**IOWA CORN**

**GEMS**

# International AgroInformatics Alliance



IAA 2.0 March 20-21, 2017, St. Paul MN



IAA 3.0 May 2-3, 2018, St. Paul MN

# GEMS Web site: now online!

Thanks

G.E.M.S:
https://agroinformatics.org