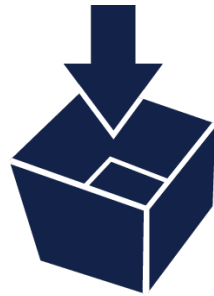# CyVerse Data Store

## Managing Your 'Big' Data

# Download Slides and Follow Along

# mcbios.readthedocs.org

# Welcome to the Data Store

Manage and share your data across all CyVerse platforms

# Working with Big Data

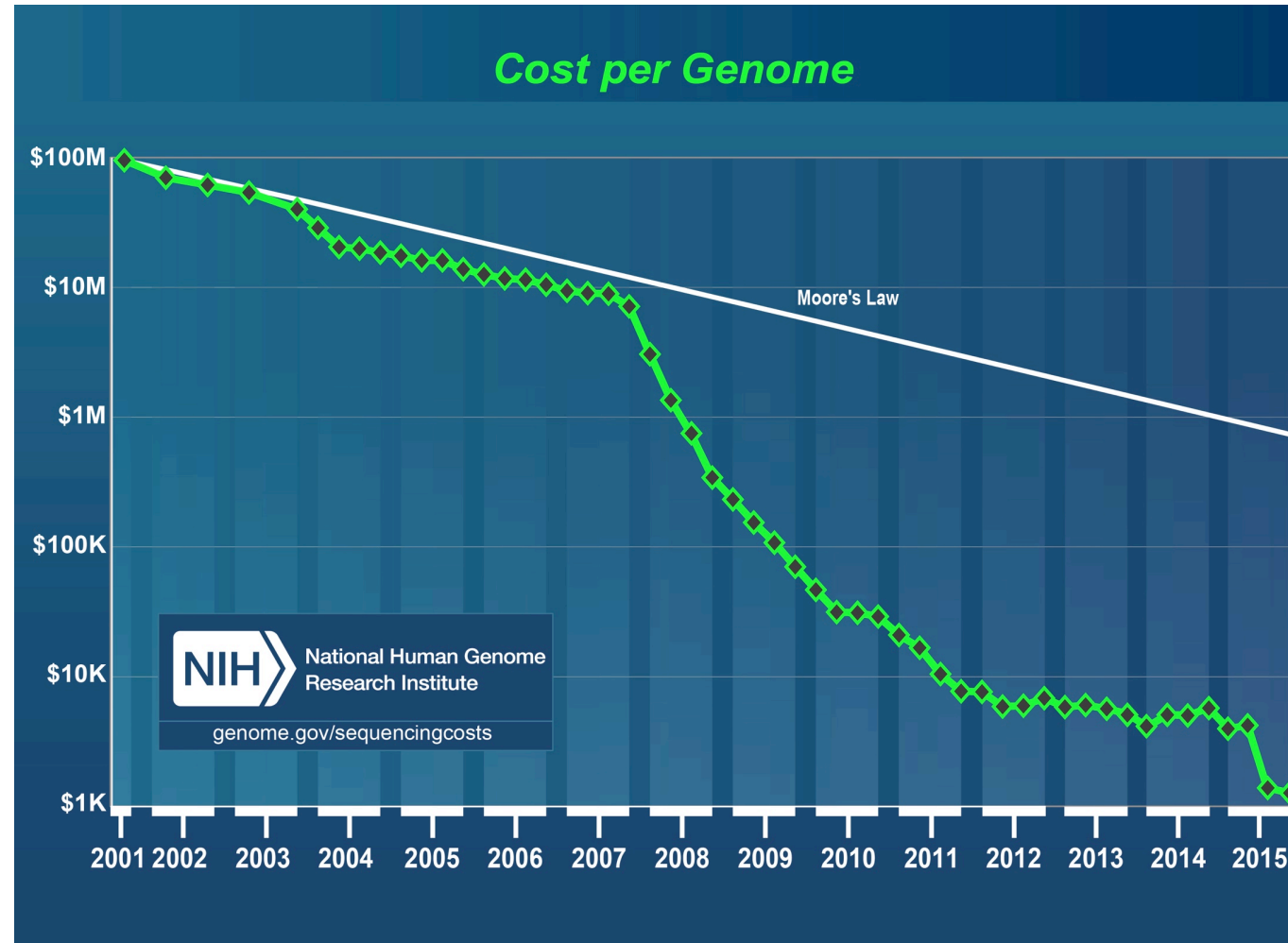Challenges: the scope and scale of life sciences data continue to grow

- Big data a term applied to data sets whose size is beyond the ability of commonly used software tools to capture, manage, and process the data within a tolerable elapsed time

- Big data sizes are a constantly moving target currently ranging from a few dozen terabytes (TB) to many petabytes of data in a single data set.

"'Big Data': Big gaps of knowledge in the field of Internet". *International Journal of Internet Science* 7: 1–5.

# Working with Big Data
## Challenges: data generation is cheaper and faster



**Cost per Genome**

Moore's Law

NIH — National Human Genome Research Institute

genome.gov/sequencingcosts
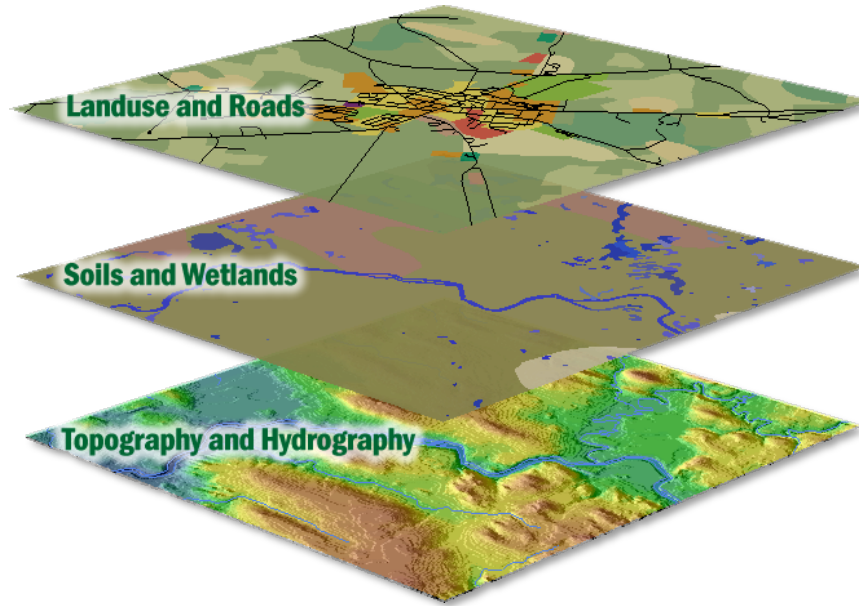
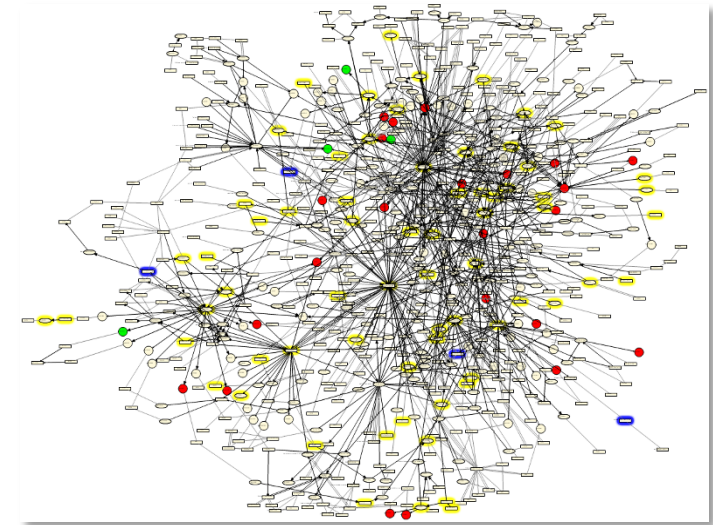http://www.genome.gov/sequencingcosts/

# Working with Big Data

Challenge: biology encompasses more than sequence data
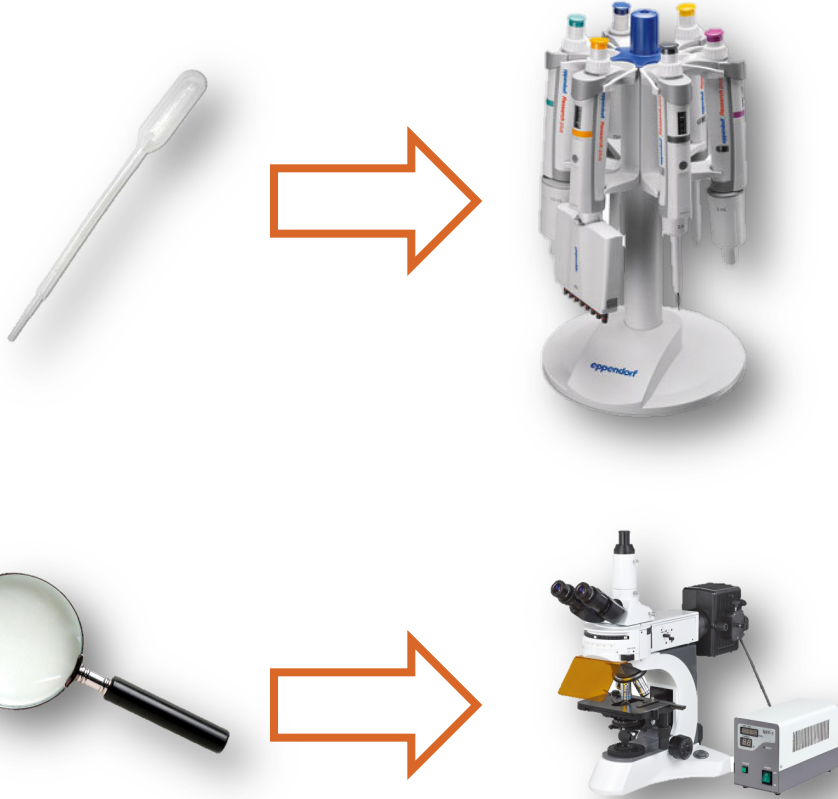


Advanced Imaging

Geospatial

Network

Biologists work with and require access to diverse data types

# Working with Big Data

Challenges: changes in data require changes in tools
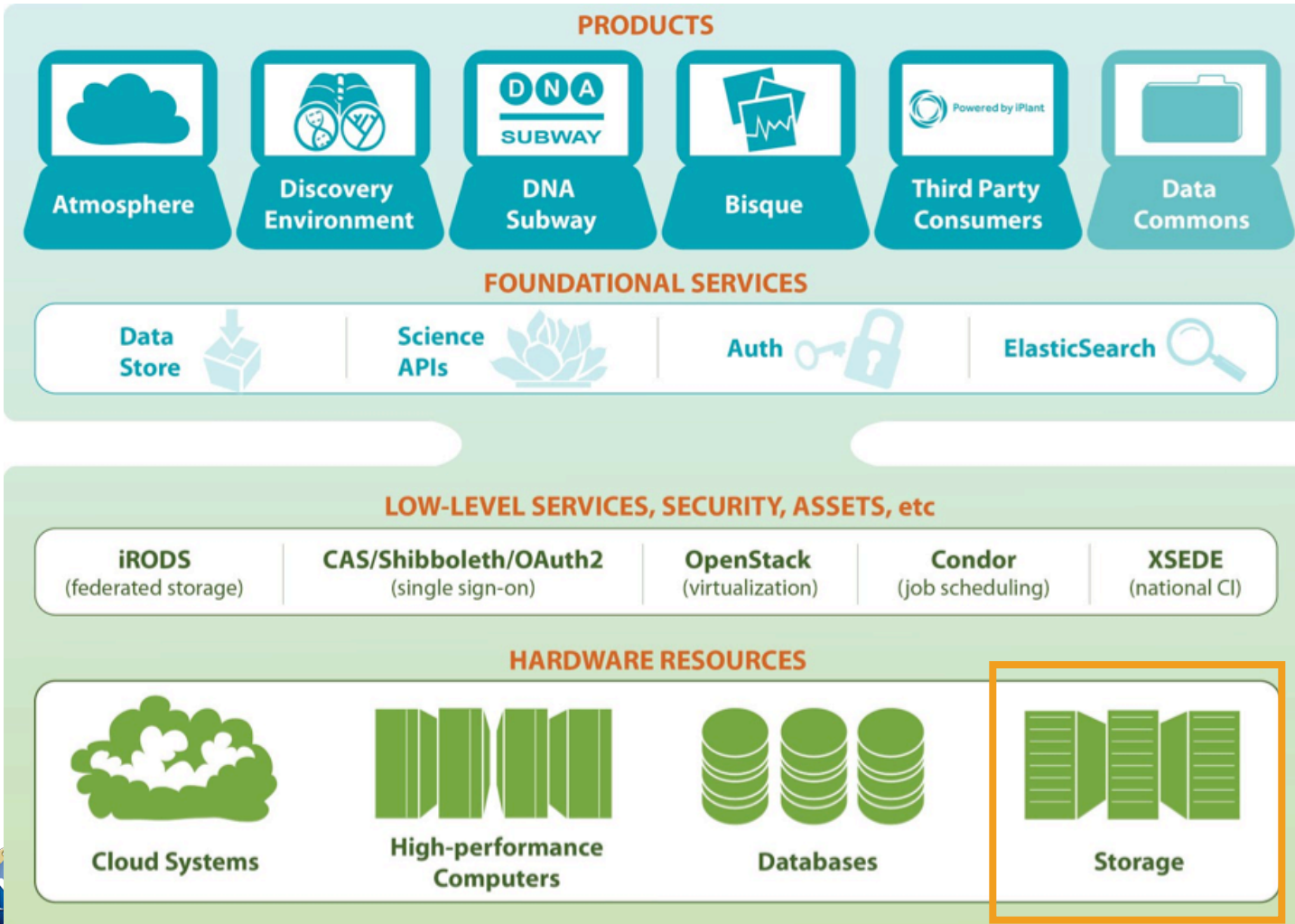
- Difficult / slow transfers

- Expense for storage / backup

- Difficult to share and publish

- Analysis

- Metadata (What Is metadata?)

Changes in scale introduce quantitative and qualitative complications

# Data Store Overview
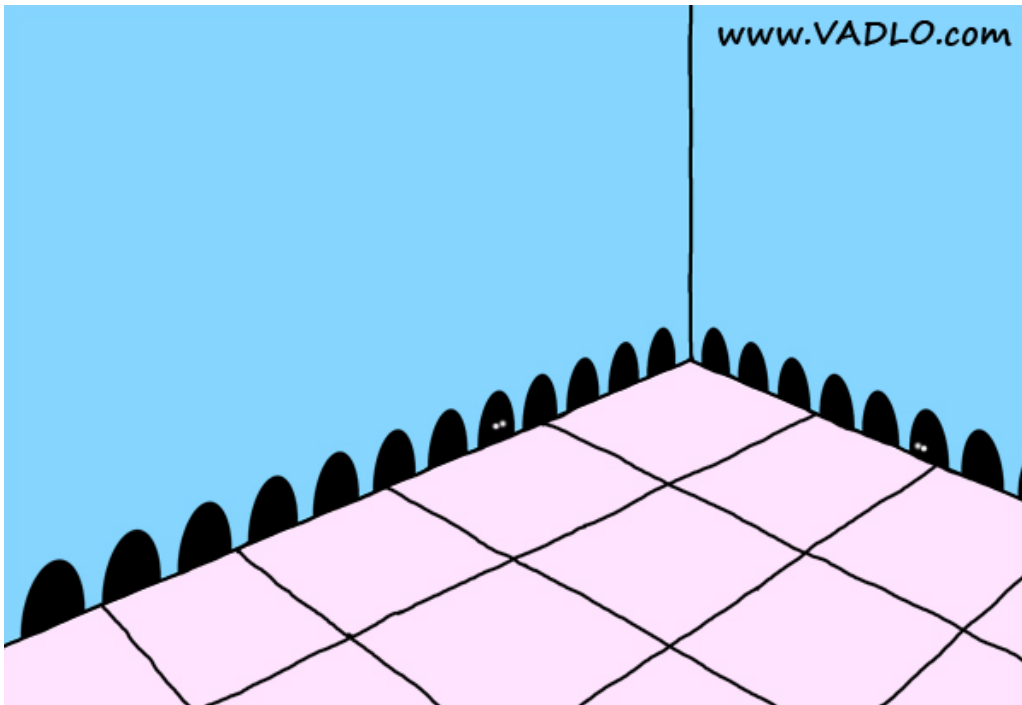
## The Data Store services all CyVerse platforms



- Access your data from multiple CyVerse services

- Automatic backup (redundant between University of Arizona and University of Texas

- Default 100 GB allocation, > 1 TB allocations available with justification

# Data Store Overview

Avoid reinventing the wheel



Mouse house that did not receive infrastructure funding.

- iRODS (integrated Rule-Oriented Data System) is an established, scalable, open-source data management sytem

- iRODS supports many data intensive projects

- iRODS abstracts data services from data storage to facilitate executing services across heterogeneous, distributed storage systems

# Data Store Overview
Benefits

**Get Science Done**

- Store any type of files related to your research

- An evolving "Data Commons" lets you access important datasets

**Reproducibility**

- Metadata captures information needed for reproducibility

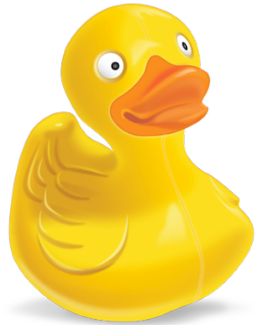- Automatic backup and accessibility support your data management plan

**Productivity**

- iRODS makes high-speed transfers possible (100 GB in ~30 min)

- Share data instantly with collaborators within CyVerse

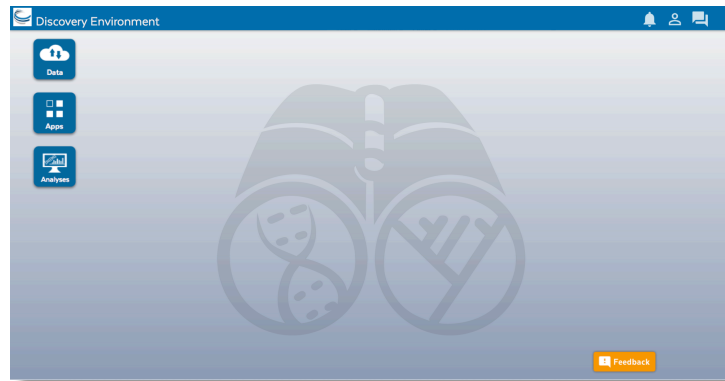# Data Store Overview
## Multiple ways to access

## Point-and-click

Cyberduck

Discovery Environment
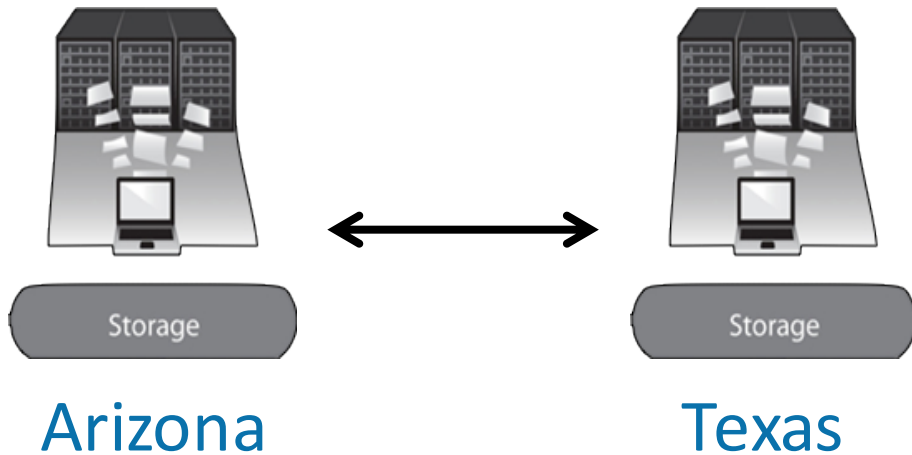
## Command line

iCommands

# Data Store Overview

Some important things we will not "see" in the demo

## Data Backups



Arizona ↔ Texas

Key component of your data management

Worry-free

## Data Transfer

| Source | Destination | Copy Method | Time (seconds) |
|---|---|---|---|
| CD | My Computer | cp | 320 |
| Berkeley Server | My Computer | scp | 150 |
| External Drive | My Computer | cp | 36 |
| USB 2.0 Flash | My Computer | cp | 30 |
| Data Store | My Computer | iget | 18 |
| My Computer | My Computer | cp | 15 |

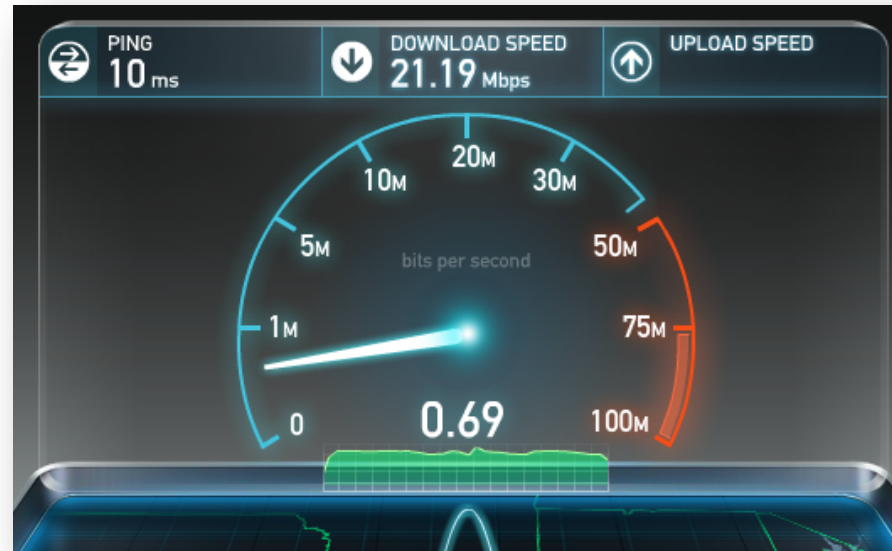Closer to optimum conditions: transfers between University of Arizona and UC Berkeley

100 GB:  26m15s, 1 GB 17.5s

# Data Store Overview

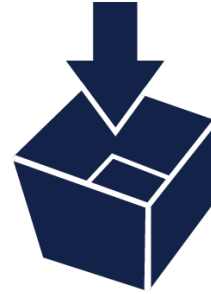Some important things we will not "see" in the demo



http://www.speedtest.net/

Local connections and institutional policies limit data transfer

# Data Store Overview



## Hands-on demo

# Data Store Overview

User perspectives and potential applications

**Bench Scientist**



- Uploads all of his .fastq files along with 50GB of root growth videos
- Shares all his analyses results with his thesis advisor

**Bioinformatician**



- Created a metadata template for assembled genomes her students and collaborators will place in a shared folder
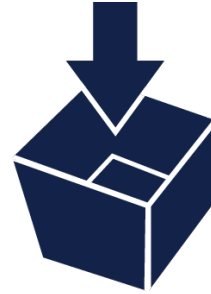- Uses public links in the supplemental materials of her publications

**Core Facilities**



- Developed a script to automate transfer of data to core users
- Uses a shared folder to make large datasets accessible

# Data Store Overview

## Time for Summaries and Tips

# Tips for any transfer method

## Spaces / Special Characters

- Many software packages are sensitive to spaces in files names and/or the special characters below

- Rename uploaded files before using them in an analysis

~ ` ! @ # $ % ^ & * ( ) + =

{ } [ ] | \ : ; " ' < > , ? /

# Tips

When sharing, use this chart to decide appropriate permissions

| Permission | Read | Download | Metadata | Rename | Move | Delete |
|---|---|---|---|---|---|---|
| Read | x | x | | | | |
| Write | x | x | x | | | |
| Own | x | x | x | x | x | x |

# Keep asking: ask.iplantcollaborative.org

## Detailed instructions with videos, manuals, documentation in Learning Center