# Slides and Tutorials

www.cyverse.org/sol2016

# CyVerse Evolution

**CyVerse 2016**
Transforming Science
Through Data-Driven
Discovery

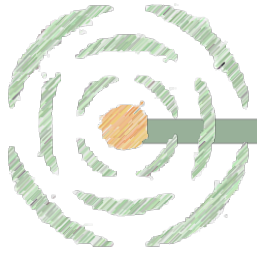Vision:

Transforming science through data-driven discovery

Mission:

Design, develop, deploy, and expand a national cyberinfrastructure for life science research, and train scientists in its use

More than 30K users, PB of data, and hundreds of publications, courses, and discoveries

# CyVerse Evolution



**iPlant 2008**
Empowering a New
Plant Biology

**iPlant 2013**
Cyberinfrastructure for
Life Science

**CyVerse 2016**
Transforming Science
Through Data-Driven
Discovery

# CyVerse is Built for Data

**Plant / Microbial**  **Animal**  **Biomedical**  **Ecological/Climate**

CyVerse supports all domains of life science

# CyVerse Evolution

## We are funded by the National Science Foundation

- We are your colleagues and collaborators!
- $100 Million in investment
- Freely available to the community
- Spur national/international collaboration
- Cite CyVerse:
    CyVerse.org/acknowledge-cite-cyverse

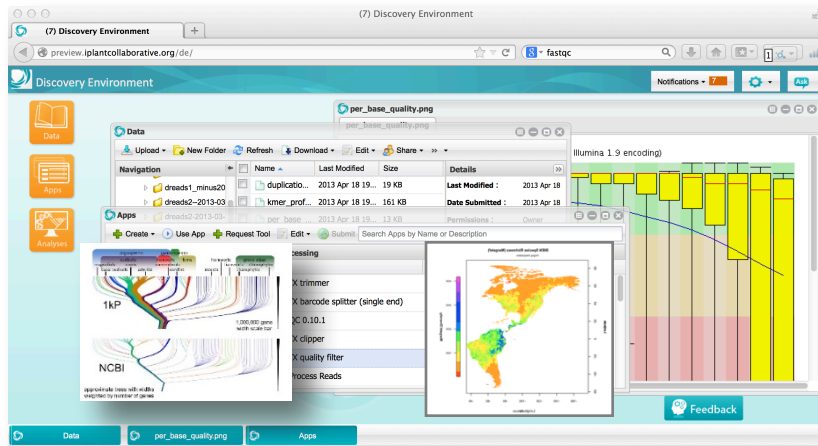# What is Cyberinfrastructure?

# Cyberinfrastructure for Data-intensive Biology

- Data storage
- Software
- High-performance computing
- People

organized to solve problems of size and scope not otherwise solvable.

# Cyberinfrastructure for Data-intensive Biology



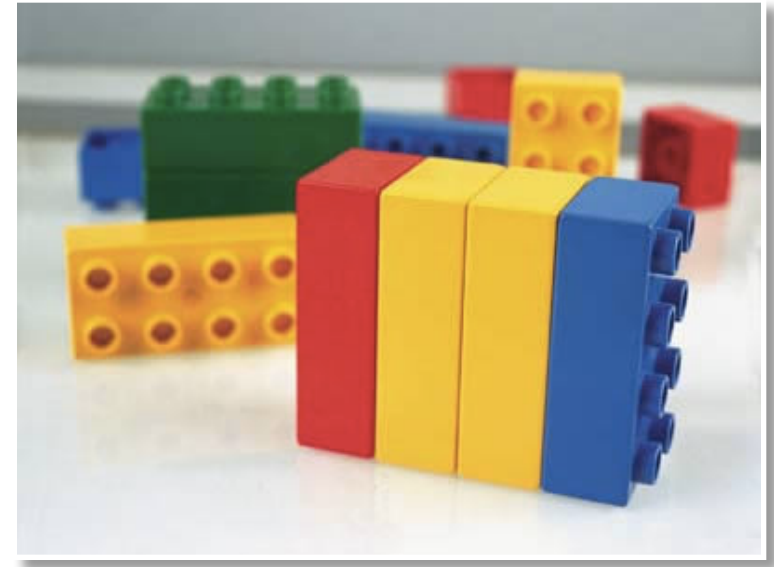Platforms, tools, datasets



Storage and compute



Training and support

# What is available in CyVerse?

# CyVerse Products

- We strive to be the **CI Lego blocks**
- Danish 'leg godt' - '**play well**'
- Also translates as '**I put together**' in Latin
- If a solution is not available you can craft your own using CyVerse CI components

# CyVerse Features

## Get Science Done

- Ability to access and manage data
- Software to analyze data
- Computing resources
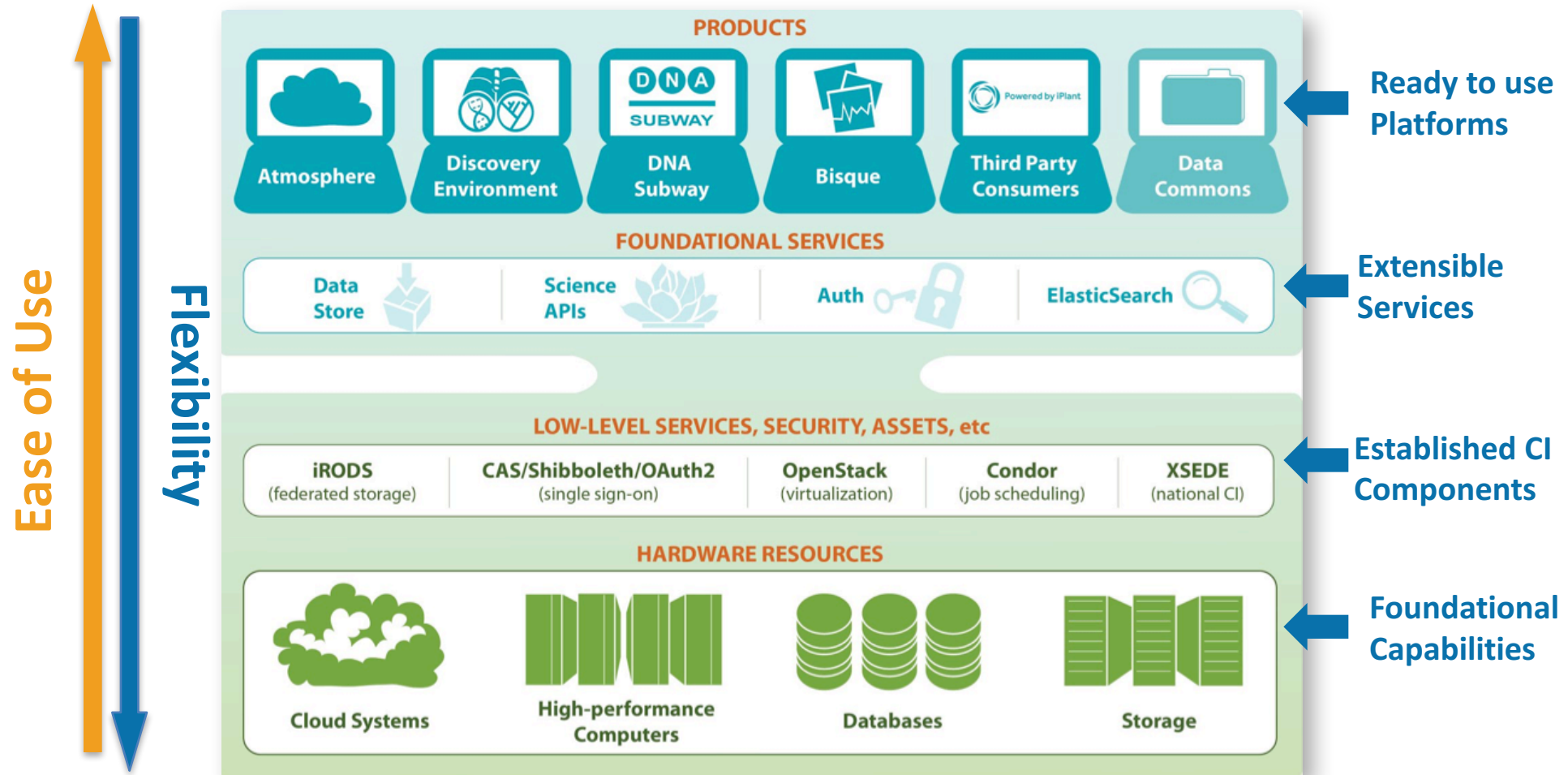- Skills and help to use software and interpret results

## Ensure Reproducibility

- Metadata management
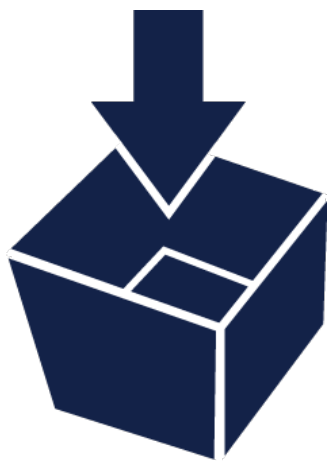- Ability to share data and workflows
- Open source sustainable tools

## Increase Productivity

- High-performance and scalable computing
- Ability automate and collaborate
- Funding spent on science, not software or hardware

# CyVerse product stack

# Data Store

The resources you need to share and manage data with your lab, colleagues and community

- ✓ Initial 100 GB allocation – TB allocations available

- ✓ Automatic data backup

- ✓ Easy upload /download and sharing

# Discovery Environment

Hundreds of bioinformatics Apps in an easy-to-use interface

- ✓ A <u>platform</u> that can run almost any bioinformatics application

- ✓ Seamlessly integrated with data and high performance computing

- ✓ User extensible – add your own applications

# Atmosphere

Cloud computing for the life sciences

- ✓ Simple: quick access to hundreds of virtual machine images

- ✓ Flexible: fully customize your software setup

- ✓ Powerful: integrated with CyVerse computing and data resources

# Science APIs

Fully customize CyVerse resources

- ✓ Science-as-a-service platform

- ✓ Define your own compute, and storage resources (local and *CyVerse*)

- ✓ Build your own app store of scientific codes and workflows

# DNA Subway

**Educational workflows for Genomes, DNA Barcoding, RNA-Seq**

- ✓ Commonly used bioinformatics tools in streamlined workflows

- ✓ Teach important concepts in biology and bioinformatics

- ✓ Inquiry-based experiments for novel discovery and publication of data

# Bisque

Image analysis, management, and metadata

✓ Secure image storage, analysis, and data management

✓ Integrate existing applications or create new ones

✓ Custom visualization and image handling routines and APIs

# How can you do science in CyVerse?

# CyVerse Data Commons

Making data discoverable and reusable

Data Store architecture is built for metadata management

- Data within CyVerse maintains provenance and can be labeled using a variety of metadata templates or using custom schemas

- Direct submission of sequencing data to NCBI Sequence Read Archive

- Request DOIs and ARKs (Archival Resource Key) for permanent data identification

- Support for elastic search, data sharing, and anonymous access

# CyVerse Data Commons

## Supporting the lifecycle of data

**Simple upload**

**markup, search, share**

**publish**

# Discovery Environment
## User friendly, developer friendly

The Agave Platform API

CyVerse Platforms
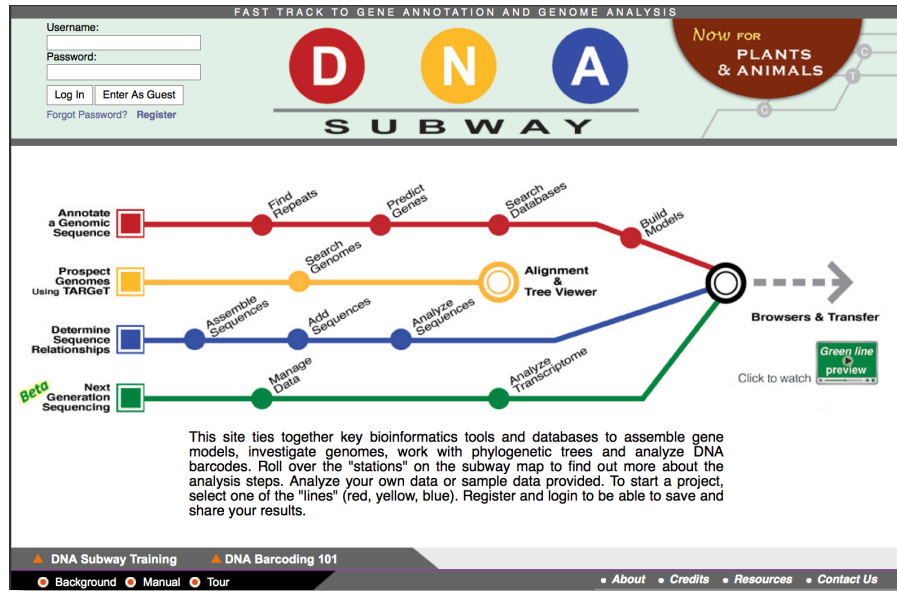
XSEDE Supercomputing

External cloud systems

# DNA Subway
## Classroom applications



**RNA-Seq** for the
**Next Gen**eration

- Classroom-friendly RNA-Seq workflows
- More than 900Gbs Student-faculty projects analyzed in 36 genomes
- Online educational materials



- DNA Barcoding workflow
- ~10,000 student-analyzed barcodes and 125 novel NCBI submissions
- Online educational materials

# How can I use CyVerse?

# Hands-on training
Workshops that reach users where they are

## Tools & Services



- Two days; targeted to researchers
- Hands-on learning modules tailor to interests
- Individual consultations
- 1026 participants since 2011

## Genomics in Education



- Two days; targeted to educators
- Pair bioinformatics with classroom labs
- Help for generating lesson plans
- 748 participants since 2010

# Online learning materials
## Documentation, tutorials, and help

### Learning Center



- cyverse.org/learning-center
- Tutorials and quick-starts

### Wiki



- wiki.cyverse.org
- Documentation

### Forum



- ask.iplantcollaborative.org
- Community help/Q&A

# CyVerse Executive Team

**Parker Antin**
**Nirav Merchant**
**Eric Lyons**

**Matt Vaughn**

**Doreen Ware**
**Dave Micklos**

# Getting Data into CyVerse

# CyVerse Data Store

- Store any type of file related to your research

- Move files seamlessly between CyVerse platforms

- Automate file transfers

- Share files with lab members, collaborators, and communities
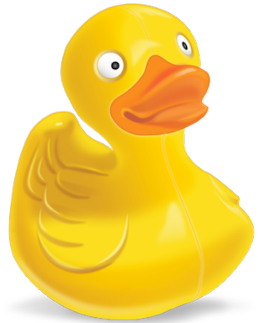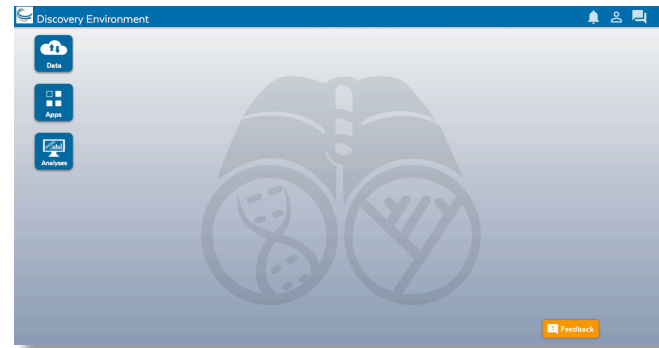
# CyVerse Data Store

Multiple ways to access

## Point-and-click



Cyberduck



Discovery Environment

## Command line



iCommands

# Discovery Environment

- Simple upload/download for small files

- Bulk upload files and folders (<10GB)

- Import from URL (no size limit)



**Advantage +**

Covers most upload/download sharing needs

**Disadvantage -**

Some size/speed limitations

# Cyberduck

- Drag and drop files and folders

- No size limit, file editing/previews

- Easy Desktop functionality



**Advantage +**

More like desktop file systems

**Disadvantage -**

No permissions/metadata control

# iCommands

- Full flexibility

- Ability to script and automate

- Access from terminal/server


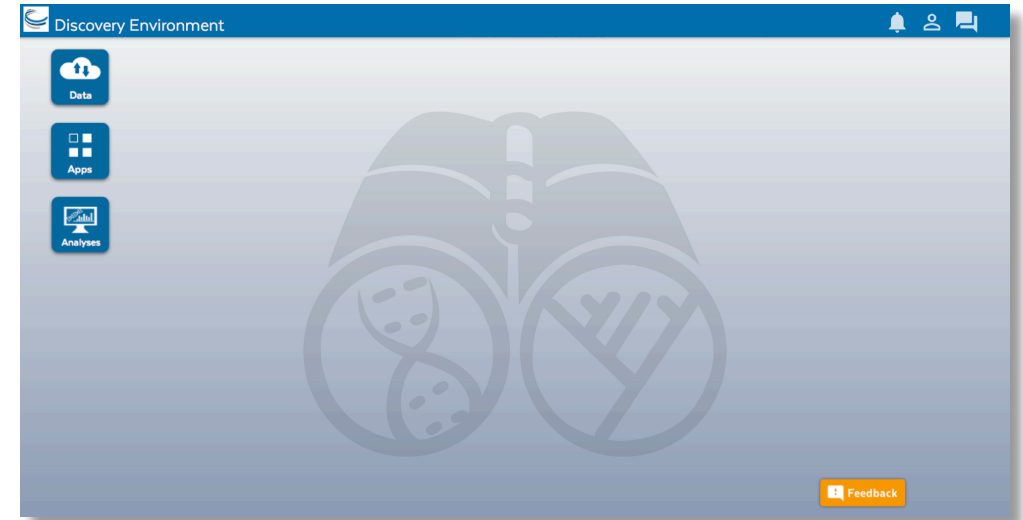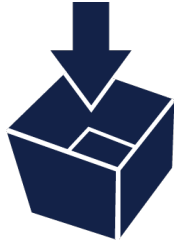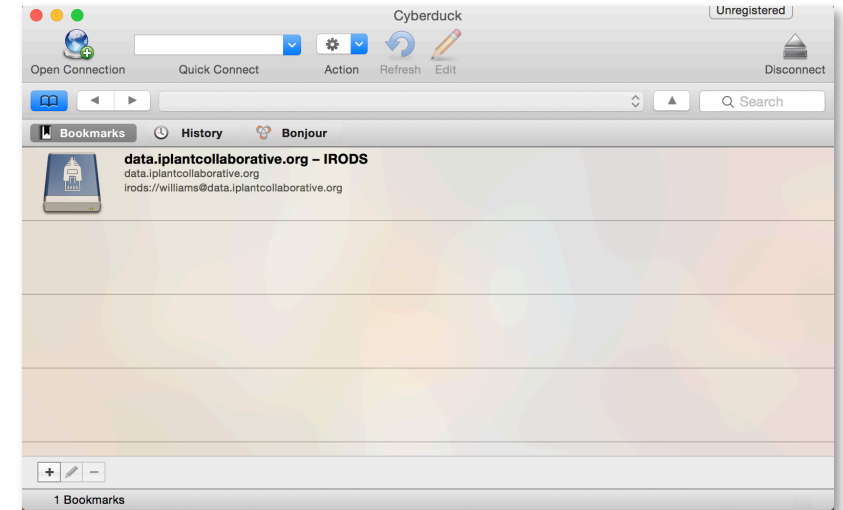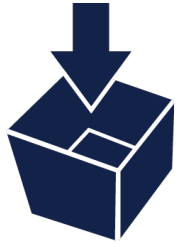
```
jasonwilliams — bash — 44×11
   C- /iplant/home/williams/analyses/Soapdeno
vo_2.04b_analysis1-2013-09-18-13-35-36.219
   C- /iplant/home/williams/analyses/Soapdeno
vo_2.04b_analysis1_47-2013-09-18-22-50-52.01
6
   C- /iplant/home/williams/analyses/TASSEL_4
.3.0__MLM__analysis1-2013-09-11-20-17-30.232
   C- /iplant/home/williams/analyses/TASSEL_4
.3.0__MLM__analysis1-2013-09-12-14-52-35.844
   C- /iplant/home/williams/analyses/Test_of_
New_App_analysis1-2013-10-25-14-40-49.857
```

**Advantage +**

Customizability

**Disadvantage -**

Requires some command line expertise

# Cyberduck and iCommands Demo

# Discovery Environment

# Discovery Environment
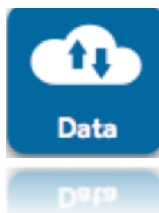
✓ A <u>platform</u> that can run almost any bioinformatics application

✓ Seamlessly integrated with data and high performance computing
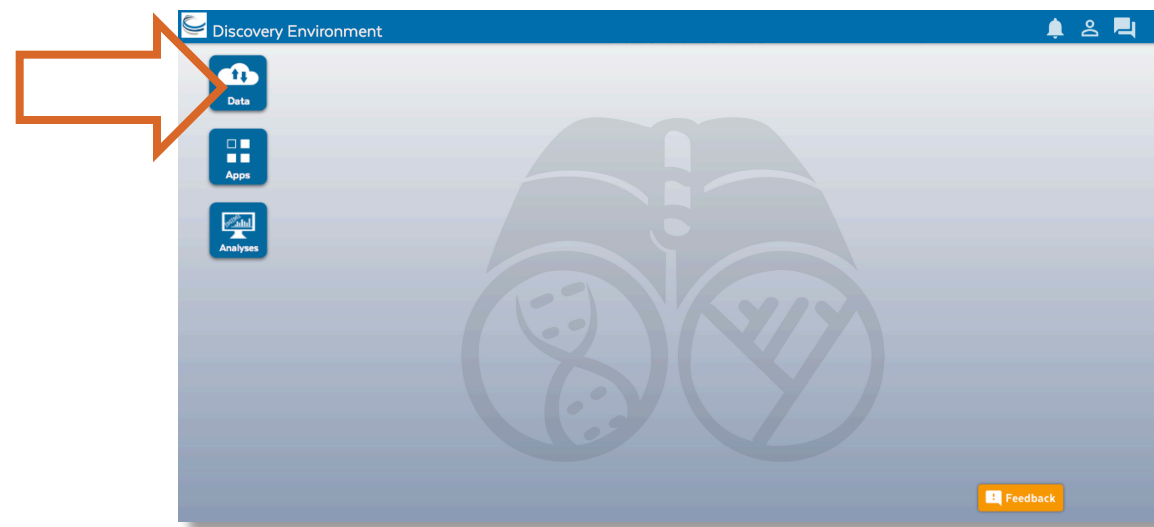
✓ User extensible – add your own applications

# Discovery Environment Overview
## Manage data

**Data**

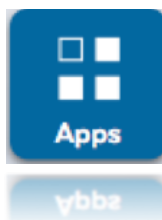- Upload / Download files and folders

- Share files via URL (Public Links)

- Share files/folders with other users

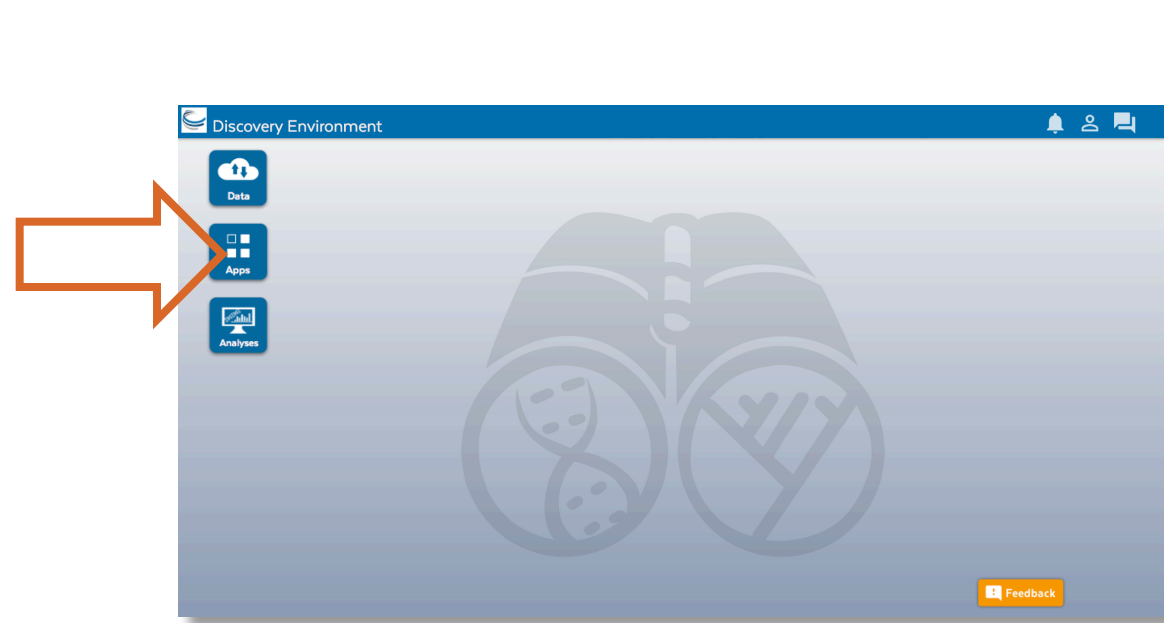# Discovery Environment Overview
## Analyze data and customize Applications

**Apps**

- Run hundreds of bioinformatics Apps

- Build automated workflows

- Modify Apps or integrate new ones
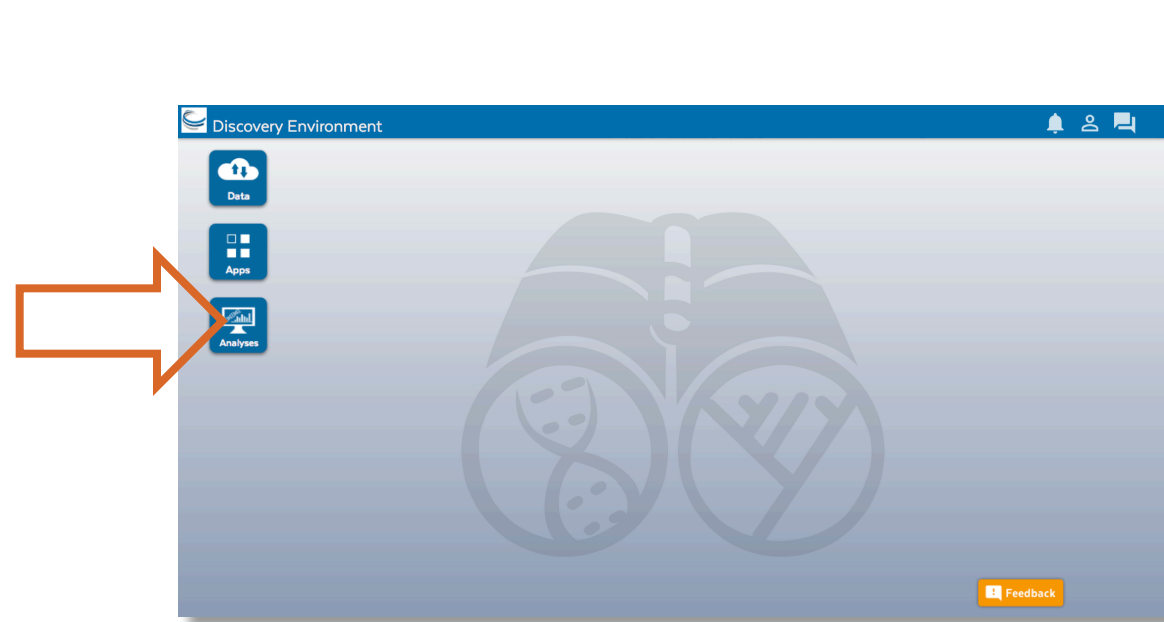
# Discovery Environment Overview
## View history, find results, reproduce analyses, optimize parameters

## **Analyses**

- Monitor job status and find results

- Cancel jobs or re-launch jobs

- Detailed job history

# Discovery Environment Demo

# Discovery Environment

Demo analysis – sequence alignment using MUSCLE

**Task:** Take unaligned DNA sequences in FASTA format and create a multiple alignment

- ✓ View sample data in iPlant Data Store

- ✓ Launch a job using the MUSCLE sequence alignment app

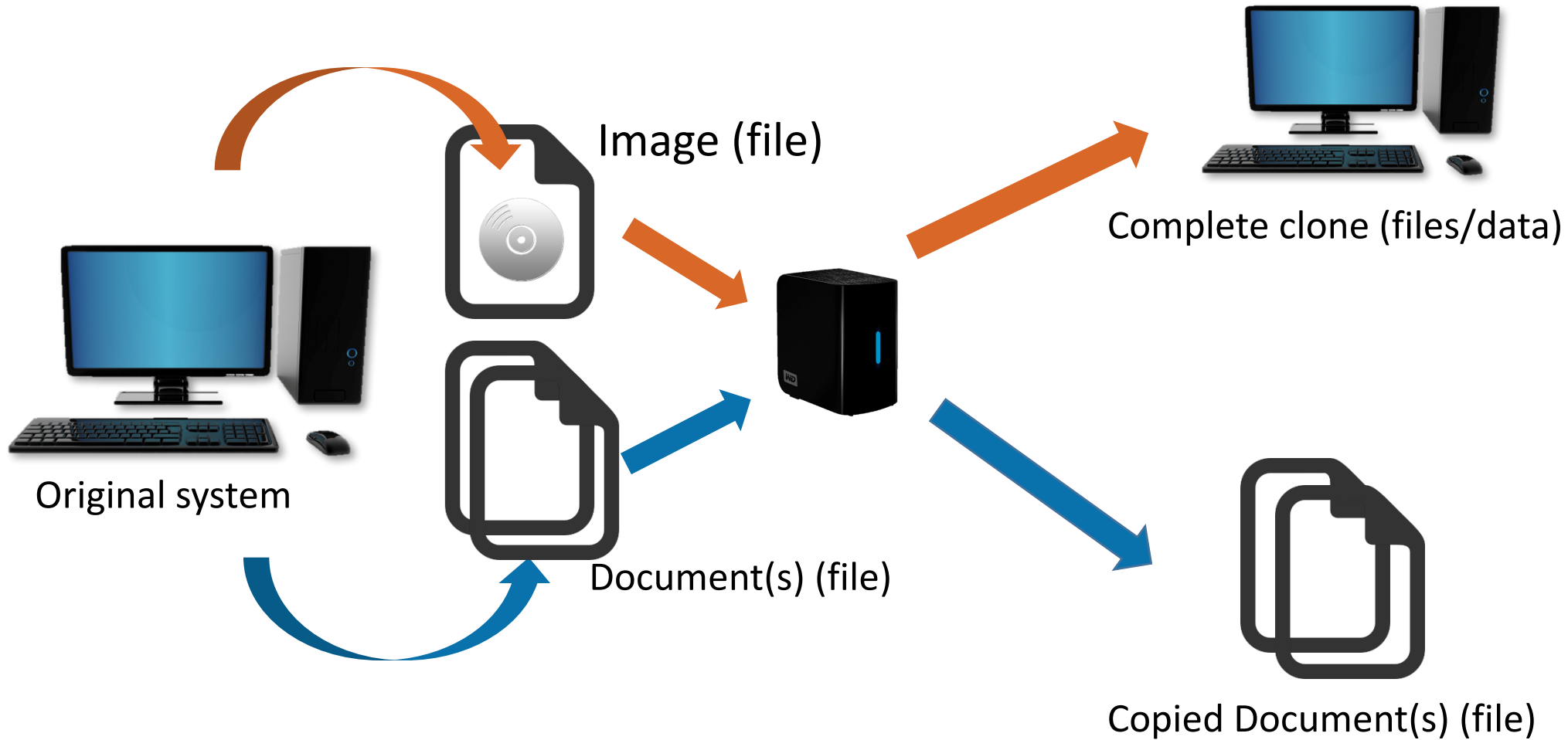- ✓ Monitor the job progress and view results

# Atmosphere

# Atmosphere

- ✓ Simple: 1-click access to hundreds of virtual machine images

- ✓ Flexible: Fully customize your software setup

- ✓ Powerful: Integrated with CyVerse computing and data resources
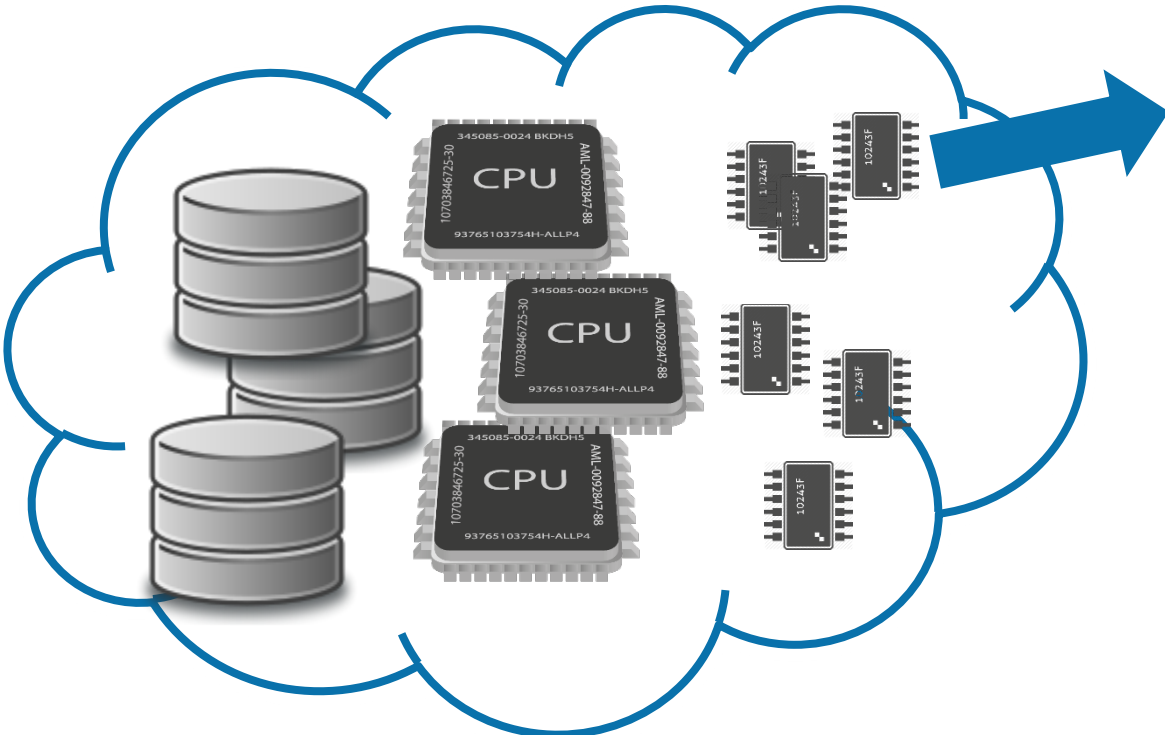
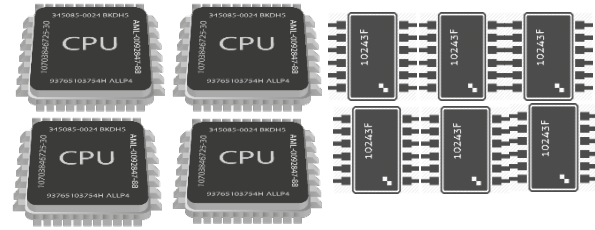# What is Cloud Computing?

Important concepts: Image



Image (file)

Original system

Document(s) (file)

Complete clone (files/data)

Copied Document(s) (file)

# What is Cloud Computing?

Important concepts: Instance



**(Disk + CPU + Memory) + (Image)**

**128.196.34.158**

**CyVerse Cloud**

**Atmosphere Instance
(virtual machine)**

# Atmosphere Overview

Largest, easiest to use cloud for Life Sciences



- Choose an existing image or customize

- Instances up to 16-Core / 128 GB RAM

- Access via shell or VNC

- Share you image with selected users, or make them public

# Atmosphere

Cloud computing for life sciences: sample use cases

- Run the software and data that are monopolizing your laptop/desktop

- Use desktop enabled images to run visually oriented programs (GUI)

- SUDO access – manage complex dependencies

- Uniform computing setups for your lab, collaborators, and students

- Make your own software available to a larger user community

# Atmosphere Demo

# Atmosphere
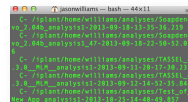Cloud computing for life sciences
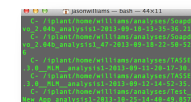
## Windows    Mac    Linux

128.196.65.175

VNC Viewer    VNC Viewer    VNC Viewer

PuTTY    Shell/terminal    Shell/terminal

VNC Viewer:    www.realvnc.com/download/viewer
PuTTy:    www.putty.org

# Where to go from here:



## Learning Center

- Get Started Guide
- Tutorials and Videos
- Documentation

## Upcoming Events

- Workshops
- Webinars

# Slides and Tutorials

www.cyverse.org/sol2016