# The iPlant Collaborative: A model for collaborative science cyberinfrastructure

### (at 30 seconds a slide)

*Dan Stanzione*

*Co-Director, The iPlant Collaborative*

*Deputy Director*

*Texas Advanced Computing Center, UT-Austin*

American Geophysical Union 2012 Meeting

# What is iPlant

- The iPlant Cyberinfrastructure Collaborative is building a comprehensive informatics *infrastructure* for plant biology.

- (and lately, some animals as well).

- This rapidly evolving infrastructure is sometimes very visible in your work, and sometimes hides in the background.

# iPlant – EarthCube

- iPlant began with many goals and a community building process, somewhat similar to EarthCube.

- The starting point in bio is wayyy behind geosciences.

  - Biology wasn't really computational before the late '90s; pre-sequencer not much digital data, very little modeling and simulation, etc.

- Presuming we have done anything right, there might be some lessons for Earthcube (or perhaps just what not to do!).

Data-intensive biology will mean getting biologists comfortable with new technology…

One key goal in our infrastructure, training and outreach is to minimize the emphasis on technology and return the focus to the biology.



**1958**
**Matt Meselson &**
**Ultracentrifuge, $500,000**



**1973**
**Sharp, Sambrook, Sugden**
**Gel Electrophoresis Chamber,**
**$250**
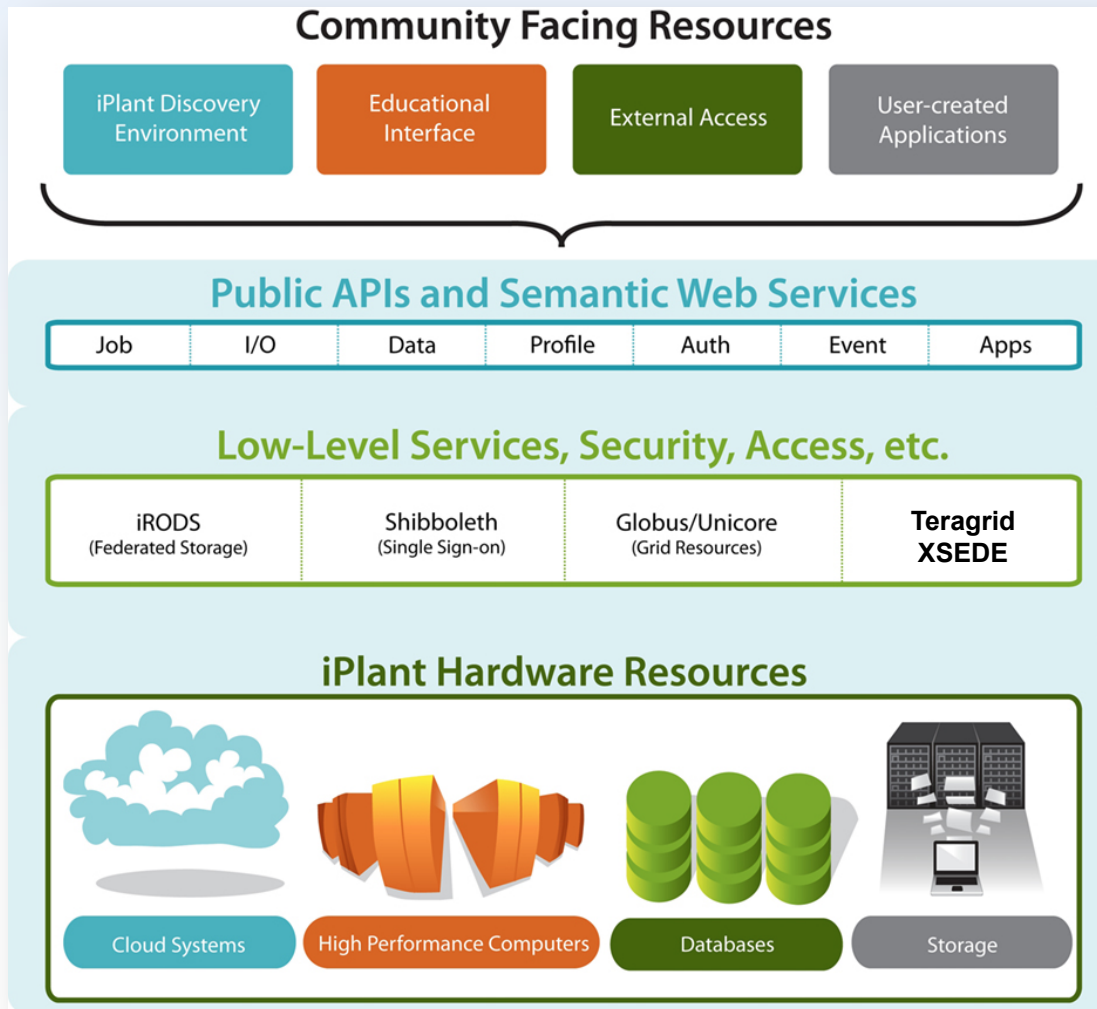
# What does iPlant Provide

- DATA
  - iPlant Data Storage: All data large and small
- COMPUTING:
  - *Large Scale:* Up to hundreds of thousands of processors
  - *Virtual:* "Cloud Style" server hosting
- A Programmer's Interface
  - Easily embed iPlant resources in your applications
- User Interfaces
  - The Discovery Environment: Integrated Web apps.
  - More than 200 bioinformatics applications
  - MyPlant, DNASubway, TNRS, TreeViewer, PhytoBisque, etc

# The iPlant Cyberinfrastructure



**Community Facing Resources**

| iPlant Discovery Environment | Educational Interface | External Access | User-created Applications |

**End Users**

**Public APIs and Semantic Web Services**

| Job | I/O | Data | Profile | Auth | Event | Apps |

**Low-Level Services, Security, Access, etc.**

| iRODS (Federated Storage) | Shibboleth (Single Sign-on) | Globus/Unicore (Grid Resources) | Teragrid XSEDE |

**iPlant Hardware Resources**

| Cloud Systems | High Performance Computers | Databases | Storage |

**Computational Users**

# Ways to Access iPlant

- **Atmosphere:** a free cloud computing platform

- **Data Store:** secure, cloud-based data storage

- **Discovery Environment:** a web portal to many integrated applications

- **DNA Subway:** genome annotation, DNA bar-coding (and more) for science educators

- **The API:** For programmers embedding iPlant infrastructure capabilities

- **Command line:** for expert access (thru TeraGrid/XSEDE)

# The iPlant Discovery Environment

- A rich web client
  - Consistent interface to bioinformatics tools
  - Portal for users who won't want to interact with lower level infrastructure

- An integrated, extensible system of applications and services
  - Additional intelligence above low level APIs – Provenance, Collaboration, etc.

# Workflows within the DE; Phylogenetics
## Trees also present computational challenges



It can take weeks or months to analyze data sets with

> 100, 000 species. Example of iPlant contribution:

NINJA/WINDJAMMER                              (Neighbor-Joining)

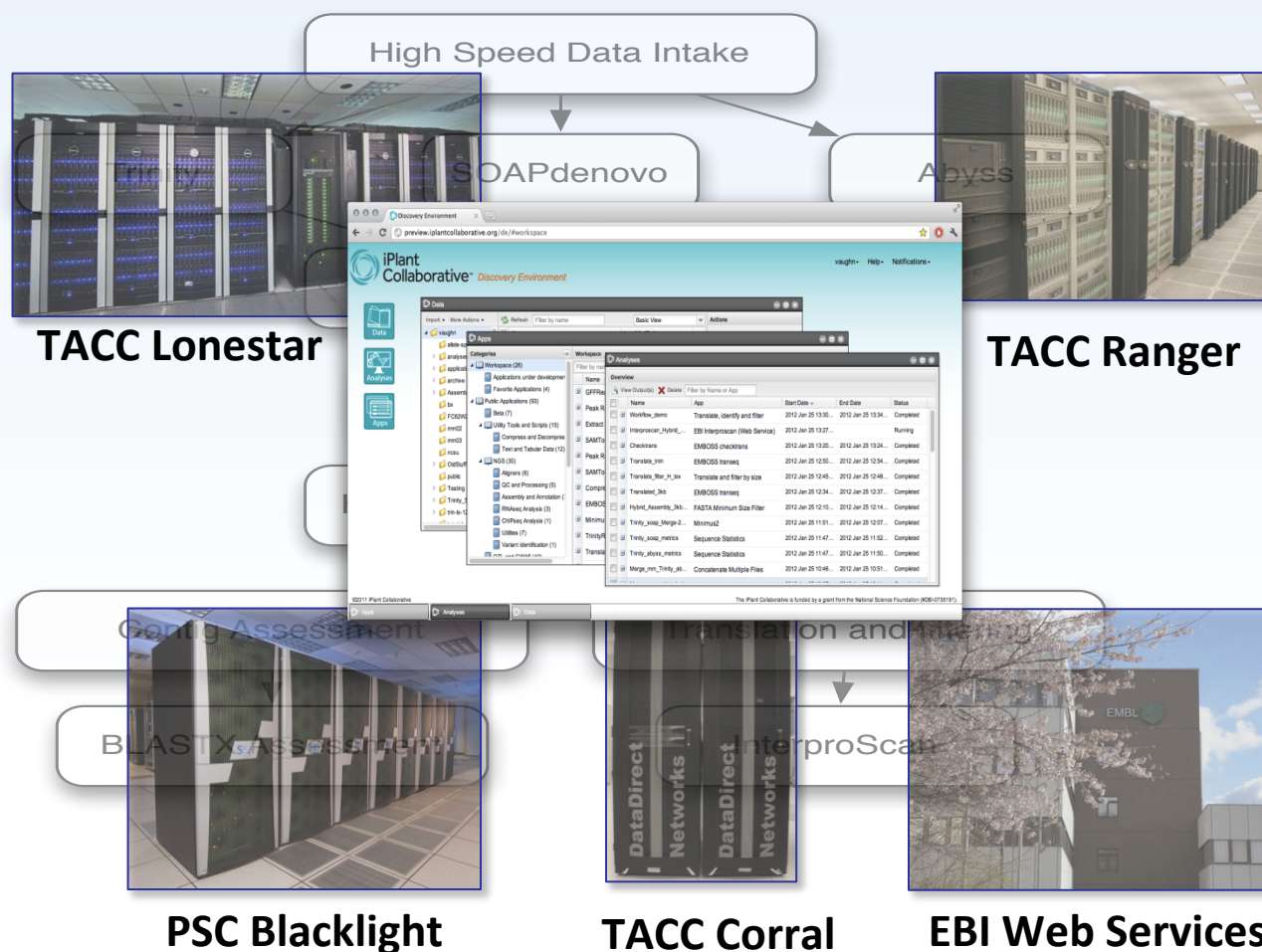-- NINJA                                               216K species, ~8 days

-- WINDJAMMER                                216K species, ~4 hours
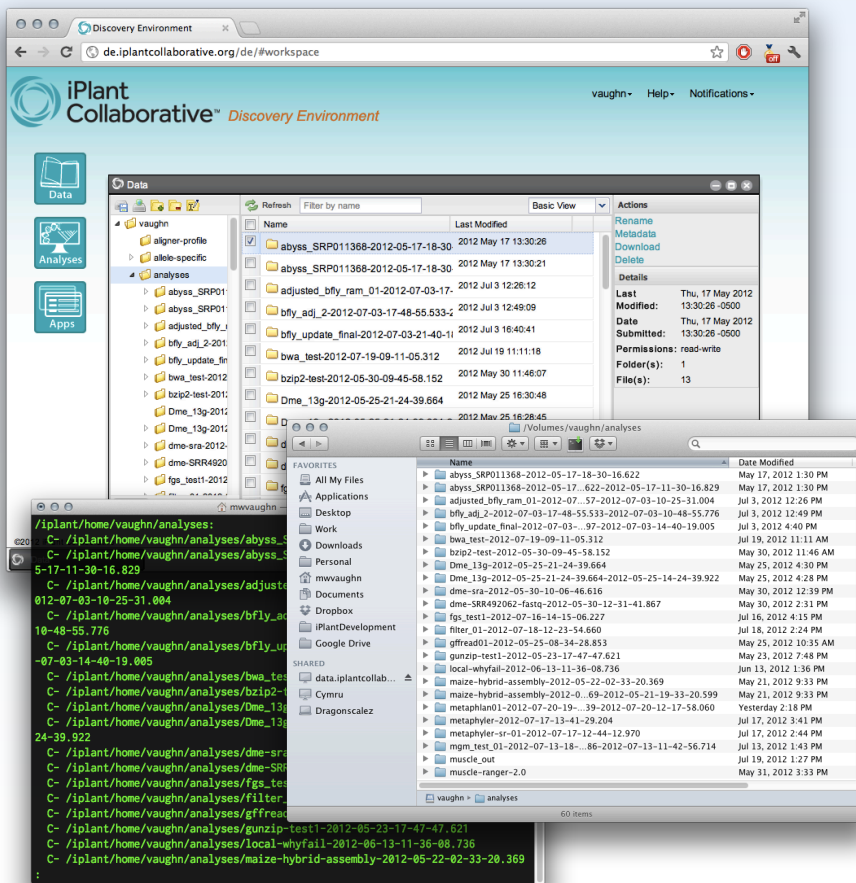
# Scalable Computation for High-Throughput Inquiry

- 90,000 Compute Cores

- Up to 1TB shared memory

- Growing to ~500,000 cores by end of 2012



**TACC Lonestar**

**TACC Ranger**

**PSC Blacklight**

**TACC Corral**

**EBI Web Services**

# The iPlant Data Store



**Fast data transfers via parallel, non-TCP file transfer**

- Move large (>2 GB) files with ease

**Multiple, consistent access modes**

- iPlant API
- iPlant web apps
- Desktop mount (FUSE/DAV)
- Java applet (iDrop)
- Command line

**Fine-grained ACL permissions**

- Sharing made simple

## Access and a storage allocation is automatic with your iPlant account

# Powered by iPlant

- The iPlant CI is designed as infrastructure. This means it is a platform upon which other projects can build.

- Use of the iPlant infrastructure can take one of several forms:
  - Storage
  - Computation
  - Hosting
  - Web Services
  - Scalability

# Powered by iPlant

- Other major projects are beginning to adopt the iPlant CI as their underlying infrastructure (some completely, some in limited ways):

  – BioExtract *(web service platform)*

  – CiPRES *(computation)*

  – Gates Integrated Breeding Platform *(hosting, development)*

  – Galaxy (*storage, for now)*

  – CoGE *(authentication, hosting*o

  – TAIR

# Biological Range Maps

***Objective: Compute range Maps for >120k species in Botanical Information and Ecology Network (BIEN) database***

- This is **Big Data**
  - Over 120k species
  - Multiple approaches required: Maximum entropy, convex hull, one and two-point algorithms
  - 11 map products (convex hull, latitude extent, etc)
  - Estimated >400 days of desktop computation
- Using TACC Longhorn: 6 hours for ~72k species test
  - Scale-up and increase in efficiency underway

# The iPlant Collaborative
## Project Atmosphere™: Custom Cloud Computing

- API-compatible implementation of Amazon EC2/S3 interfaces

- Virtualize the execution environment for applications and services

- Up to 12 core / 48 GB instances

- Access to Cloud Storage + EBS

- Big data and the desktop are co-local again
  - Bring your data to Atmosphere VM for interactive access and analysis
  - Send it back to the DE for transactional analysis



>60 hosted applications in Atmosphere today, including users from USDA, Forest Service, database providers, etc.

(30 more for postdocs and grad students for training classes)

# Atmosphere: Motivations

- Standalone GUI-based applications are frequently required for analysis
- GUI apps not easily to transform into web apps
- Need to handle complex software dependencies (e.g specific bioperl version and R modules)
- Users needing full control of their software stack (occasional sudo access)
- Need to share desktop/applications for collaborative analysis (remote collaborators)

# Atmosphere: What is it?

- Self-service cloud infrastructure
- Designed to make underlying cloud infrastructure easy to use by novice user
- Built on open source Eucalyptus
- Fully integrated into iPlant authentication and storage and HPC capabilities
- Enables users to build custom images/ **appliances** and share with community
- Cross-platform desktop access to GUI applications in the cloud (using VNC)
- Provide easy web based access to resources

# Atmosphere: Launch a new VM

# Atmosphere: Access a running VM

# Atmosphere: Log in via shell

# Collaborations within iPlant

- At the conception of the project, there were, ahem, "wide ranging" notions of what CI collaboration meant.

- Videoconferencing, chat, social networking of all shapes and sizes, wikis, interactive visualization, matchmaking services, virtual conferences, etc.

- Remember that "behind" starting point

# Data sharing is the key collaboration mechanism

- Lots of these things were essentially solved problems.

- Re-implementing them was usually a bad idea.

- The ability to have a shared dataspace mattered more than the communication mechanisms.

  – A working dataspace, not just an archive.

  – Ability to share the workflows used to analyze the data, not just the data.

# A 15 second case study

- A large plant transcriptome project chose to work with iPlant (who will remain nameless to protect the guilty).

- They wanted to use iPlant storage, and wanted help building ingest pipelines.

- They were not interested in changing their workflow, so they didn't use things like metadata support.

- Now with 5,000+ directories, no search capability beyond "grep", no unique identifier scheme, and no notion of how to pull data to release from the overall corpus, they are seeing the pain of not using proper data management, metadata, etc.

# Additional Lessons

- Consensus is not going to merge on the collaboration issue (or most other ones).

- Listen carefully to the community, but don't be paralyzed.

  - Build some things. Convince people to use them. Give up on the ones that don't work, and focus on the ones that do.

  - Then try some more.

  - Find users in the community who want to *invest* their time in doing things differently, and make them successful.

- There are other projects in the world; build with federation in mind from day 1 (how does your data repository talk to other ones?).

# Additional Lessons(2)

- There are other projects in the world; build with federation in mind from day 1

  – How does your data repository talk to other ones?.

  – Can you use a standard interface? Or build on top of one? (i.e. REST).

- (Strongly) encouraging best practices is really important… but requiring compliance can be self-defeating

  – Semantic web technologies is a good example)

  – Data formats, Ontologies, Metadata/tagging, versioning, unique IDS, retention policies, transfer mechanisms, etc.

# Does it work?

- 7,000+ users have DE accounts.

- 200+ Atmosphere images exist.

- 300+ TB is stored in the Data store.

- iPlant users have run thousands of jobs using millions of hours on supercomputers through iPlant in the last year.

- 100+ NSF awards mention iPlant

- It's hard to authoritatively say if this means transformational science... but something is happening.

# Questions?

dan@tacc.utexas.edu

dan@iplantcollaborative.org