# iDigBio Technology, Cloud and Appliances

Jose Fortes
(on behalf of the iDigBio IT team)

# iDigBio (idigbio.org)

- **Goal**: making data and images for millions of biological specimens available in electronic format for the biological research community, agencies, students, educators, and public

- **Mission:** leadership, coordination, and outreach in digitization of collections by implementing resources for communication, use of technology, access to data, research and education.
  - The "Hub" part of the NSF ADBC program aggregating TCNs and PENs
- A **resource**: permanent cloud computing infrastructure
  - to link biological data from collections across the USA
  - to use search and analytics tools to mine and reference data

# Research Questions

- How are species distributed in geographical and ecological space?
- What is the history of life on Earth?
- What factors lead to speciation, dispersal, and extinction?
- What are the impacts of climate change likely to be?
- What information is needed for effective conservation strategies?

Slide provided by Pam Soltis

# iDigBio IT Vision

- Cyberinfrastructure to enable
  - the collaborative creation, integration and management of digitized biocollections,
  - their use in scientific research, education and outreach
- Visible as a collection of persistent Internet-accessible services, data and resources
  - For biocollection "producers"
  - For biocollection "consumers"
  - For biocollection service providers
  - For cyberinfrastructure providers
  - For national/global data aggregators

# CI Stakeholders



Domain Data Producers
- Museums
- TCNs
- Collectors

Infrastructure Providers
- Amazon Turk
- Amazon WS
- Google
- DataONE
- TCNs
- Microsoft Azure
- Data Conservancy

National/Global Data Aggregators
- GBIF
- iPlant
- ALA
- EOL
- BISON

iDigBio

Domain Service Providers
- Georeferencing
- Imaging services
- Data quality
- NESCent
- OCR
- Translation
- iPlant
- Mapping
- TCNs

Domain Data Consumers
- Researchers
- Teachers
- Citizens
- TCNs
- Government

# Stakeholders APIs



**Domain Data Producers**

TCNs
Collectors
Museums

Amazon Turk
Amazon WS
Google
DataONE
TCNs
Microsoft Azure
Data Conservancy

GBIF
ALA
EOL
BISON

**National/Global Data Aggregators**

Domain-level data

Updates Notification Usage track

**Infrastructure Providers**

iDigBio

Domain data

BLOBs Appliances

Updates Notification

Query results

Customer Requests

Processed data

Researchers
Teachers
Citizens
TCNs

**Domain Data Consumers**

**Domain Service Providers**

Government

Mapping    TCNs

Georeferencing
Imaging services
Data quality    NESCent
OCR    Translation
iPlant

# Interface Model for iDigBio and TCNs

TCNs



iDigBio + Resources

| | | | | | |
|---|---|---|---|---|---|
| Archiving | Data Collections | Wiki | Workshop Resources | Workflow Engines | Taxonomic Validation |
| Learning Modules | Structured Data Services | Storage | Non-structured Data Services | Geographical Mapping | Data Conversion |
| Virtual Appliances | Machines | | | Networking | Collaboration Tools |

TCP    OCCIWG    HTTP    RDF    JPEG2000    X.509    OpenID    XMPP    ODBC

| | | | | | | |
|---|---|---|---|---|---|---|
| National History Museums | Federal Collections | Applied Innovations | Microsoft Azure | Amazon EC2/S3 | Google Apps | Microsoft Live | Google App Engine |

iPlant   NCBI   EOL   LifeMapper   ALA   XSEDE   TCNs   DataONE   Academic Clouds   NESCent

Infrastructure Providers, National/Global Data Aggregators, Domain Service Providers, Domain Data Consumers
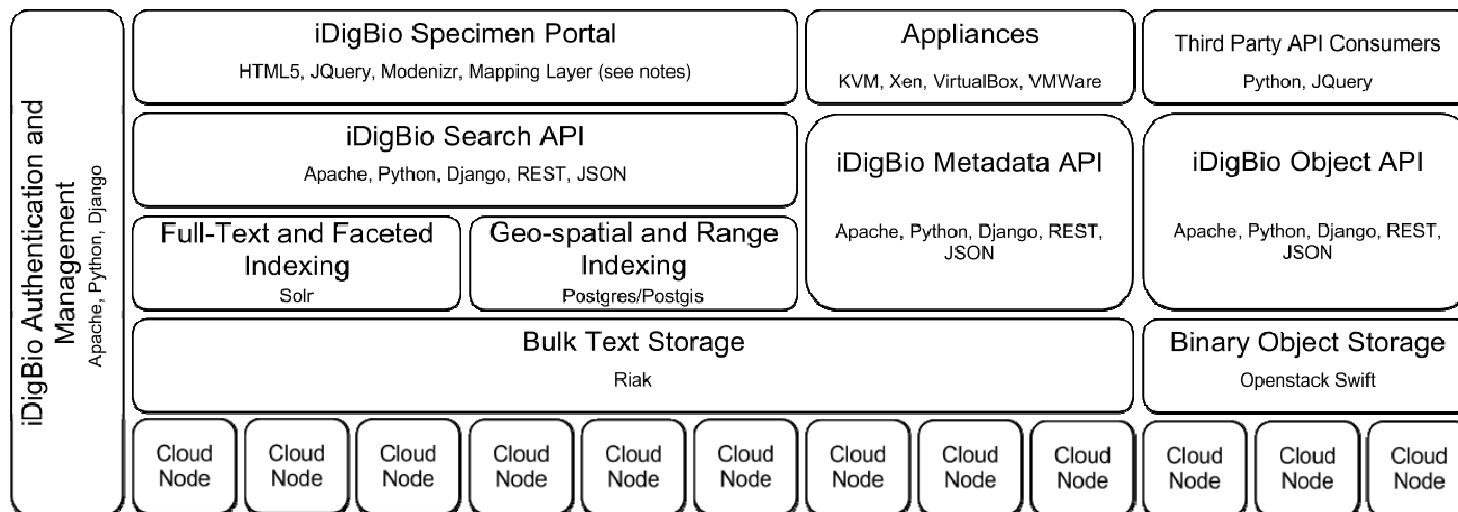
# Building the iDigBio Cloud

- Cloud-based strategy
  - Providing useful services/APIs (programmatic and web-based)
  - Federated scalable object storage and information processing
  - Digitization-oriented virtual appliances
  - Reliance on standards, proven solutions and sustainable software
- Continuous consultation with stakeholders
  - Surveys, workgroups, summit/workshops, person-to-person …

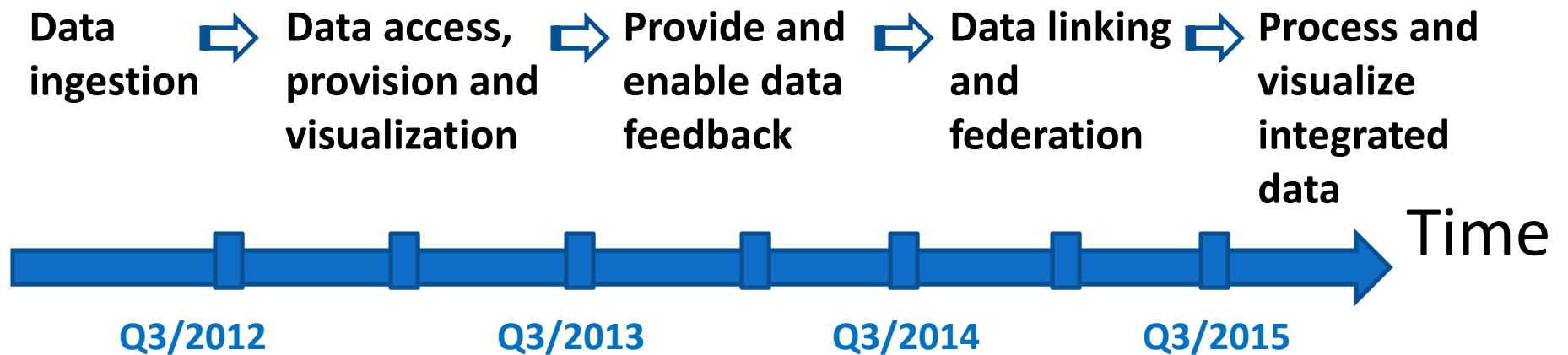| iDigBio Authentication and Management<br>Apache, Python, Django | iDigBio Specimen Portal<br>HTML5, JQuery, Modenizr, Mapping Layer (see notes) | | Appliances<br>KVM, Xen, VirtualBox, VMWare | Third Party API Consumers<br>Python, JQuery |
|---|---|---|---|---|
| | iDigBio Search API<br>Apache, Python, Django, REST, JSON | | iDigBio Metadata API<br><br>Apache, Python, Django, REST, JSON | iDigBio Object API<br><br>Apache, Python, Django, REST, JSON |
| | Full-Text and Faceted Indexing<br>Solr | Geo-spatial and Range Indexing<br>Postgres/Postgis | | |
| | Bulk Text Storage<br>Riak | | | Binary Object Storage<br>Openstack Swift |
| | Cloud Node · Cloud Node · Cloud Node · Cloud Node · Cloud Node · Cloud Node | | Cloud Node · Cloud Node · Cloud Node | Cloud Node · Cloud Node · Cloud Node |

# Keeping our eyes on the ball

Common/frequent needs: archival storage, server hosting, feedback on the data, data intensive transformations …

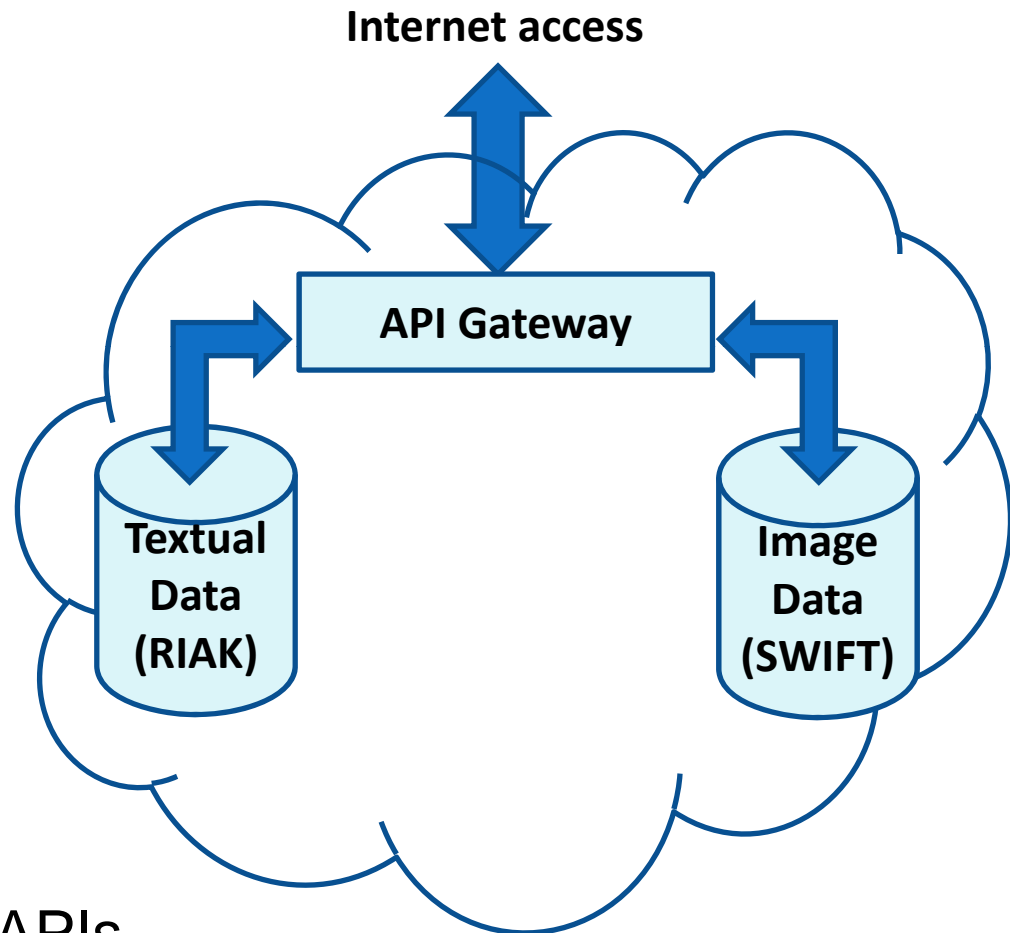10-year tsunami of requirements: from being on Facebook to multilingual search-and-compute across multiple data sets…

# Evolution of iDigBio capabilities

**Data ingestion** → **Data access, provision and visualization** → **Provide and enable data feedback** → **Data linking and federation** → **Process and visualize integrated data**

Time

Q3/2012          Q3/2013          Q3/2014          Q3/2015

Increasing storage and server  hosting in support of the above
Increasing number of appliances in support of the above
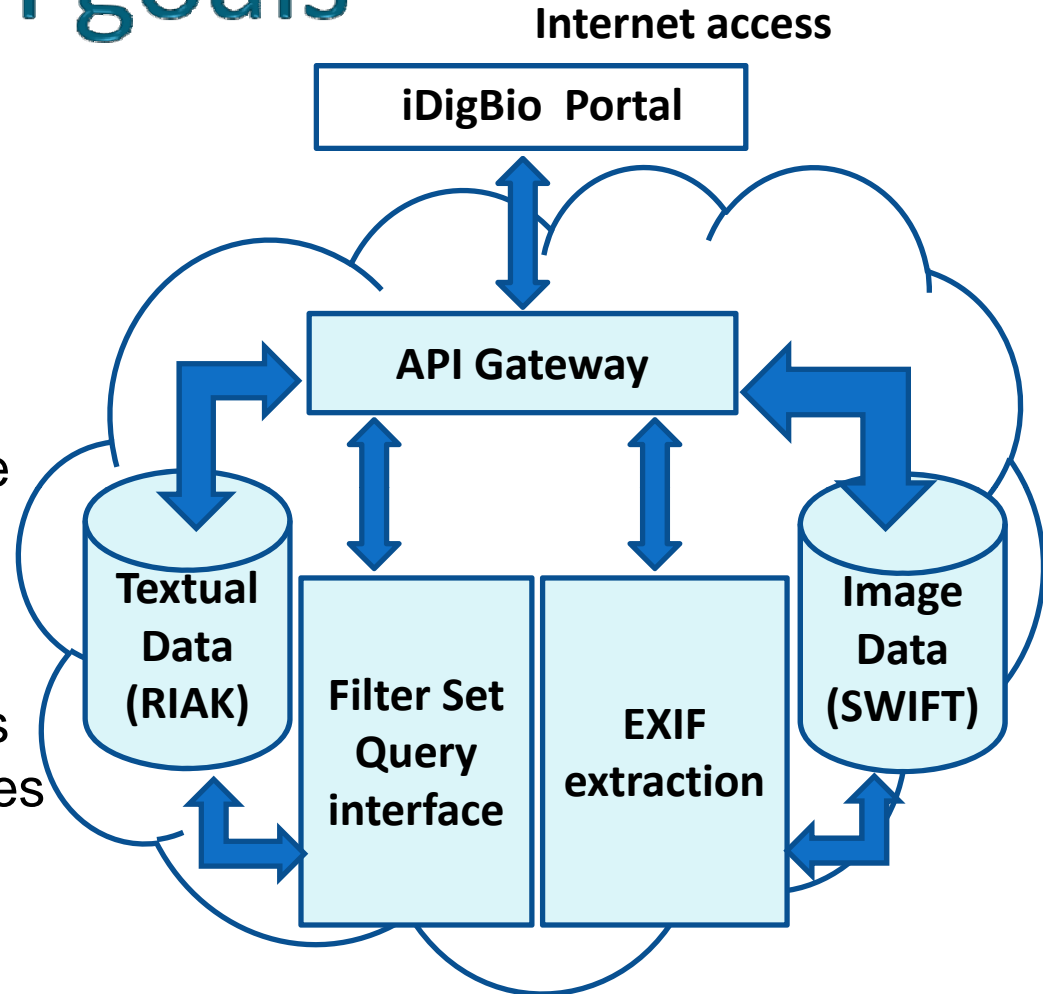Web site for interaction with public, community, education and above

# Near-term goals: ingest data

- Textual data
  - JSON document database
  - Data ingestion via DwC-a files
  - Get / Set API

- Image Data
  - Internet-accessible object storage
  - Upload appliance
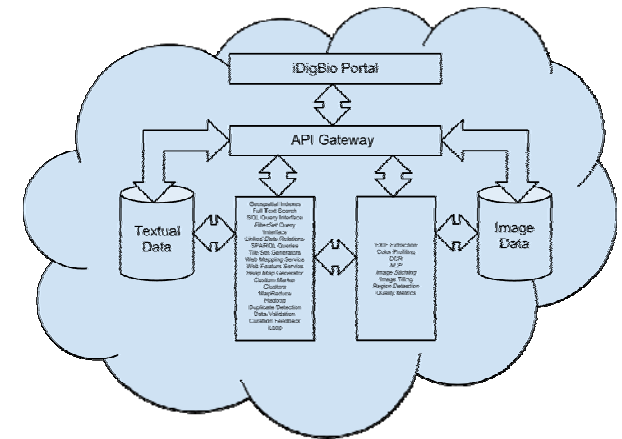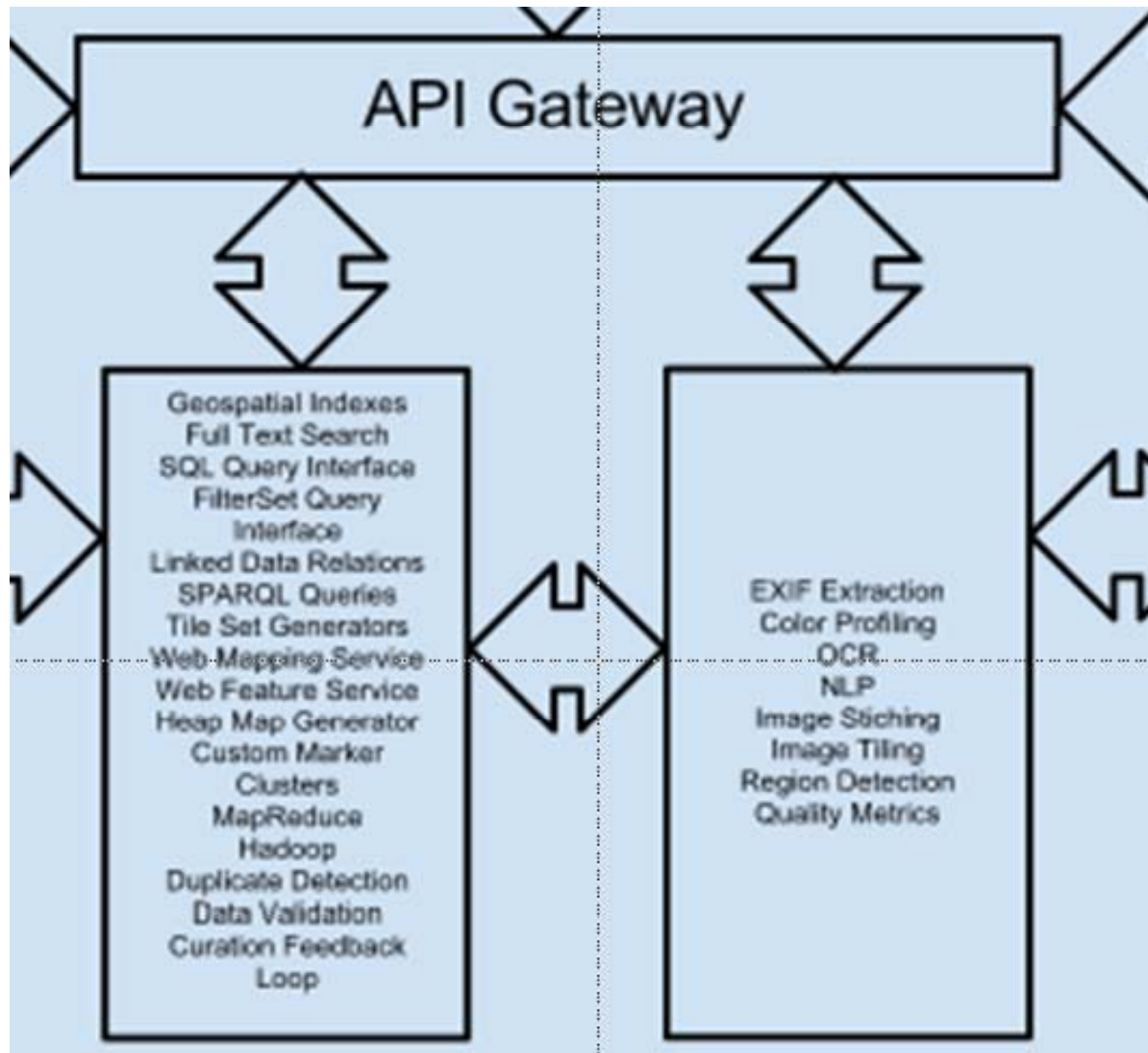  - Limited access to low-level APIs

**Internet access**

**API Gateway**

**Textual Data (RIAK)**

**Image Data (SWIFT)**

# Medium-term goals

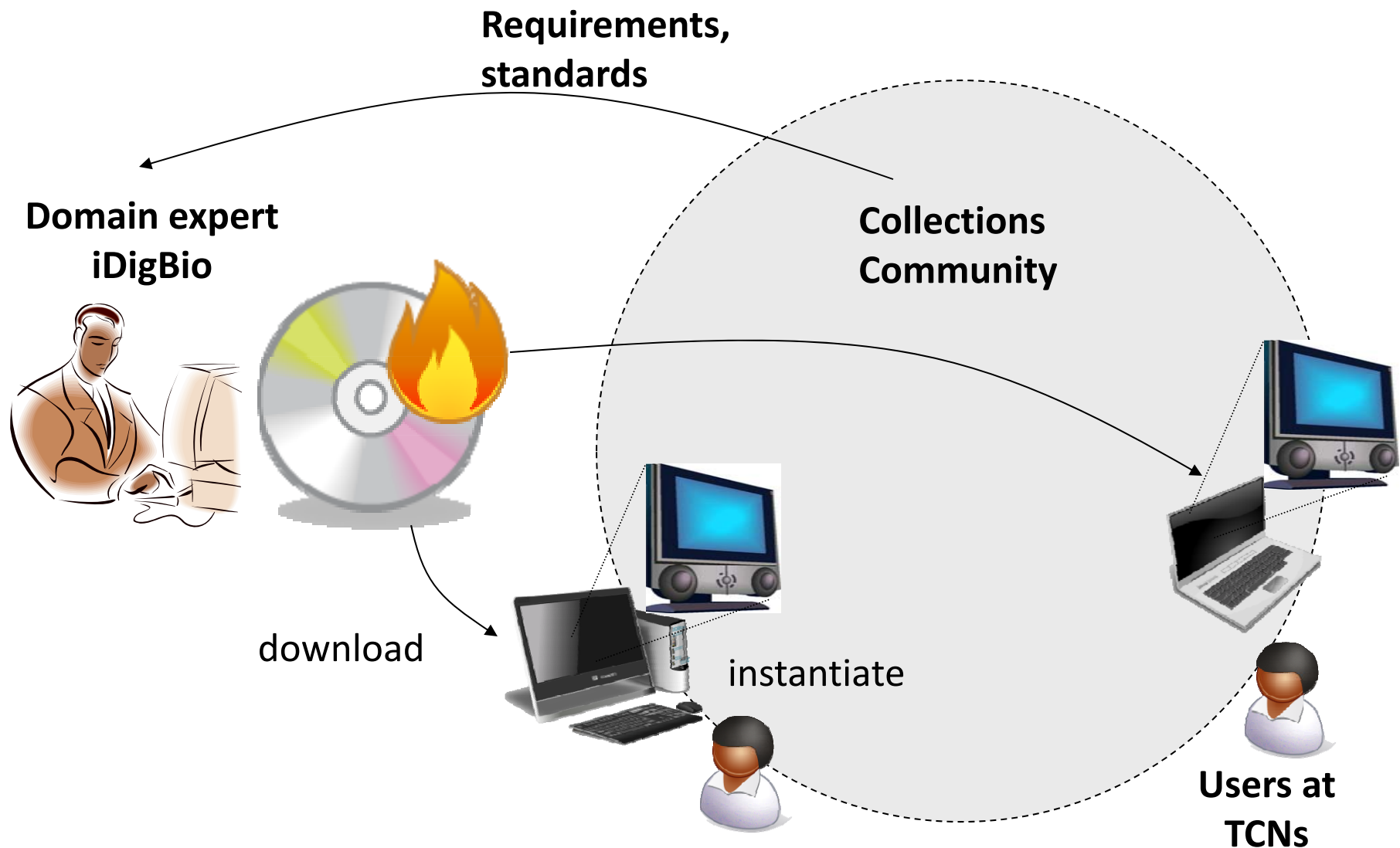**Internet access**

- Textual Data
  - JSON document database
  - Data Ingestion via DwC-a files
  - Rich RESTful API
- Image Data
  - Web-accessible object storage
  - Upload appliance
  - Fully abstracted storage
- Indexing and Search
  - Extract EXIF data from images
  - Limited but useful set of indexes
  - Intuitive search UI
  - Search available via API
- Portal
  - Consumes and interfaces text, image and search APIs (minimal server side code)
  - Web-based mapping - client side javascript limits useable record count to about 50k records at a time.
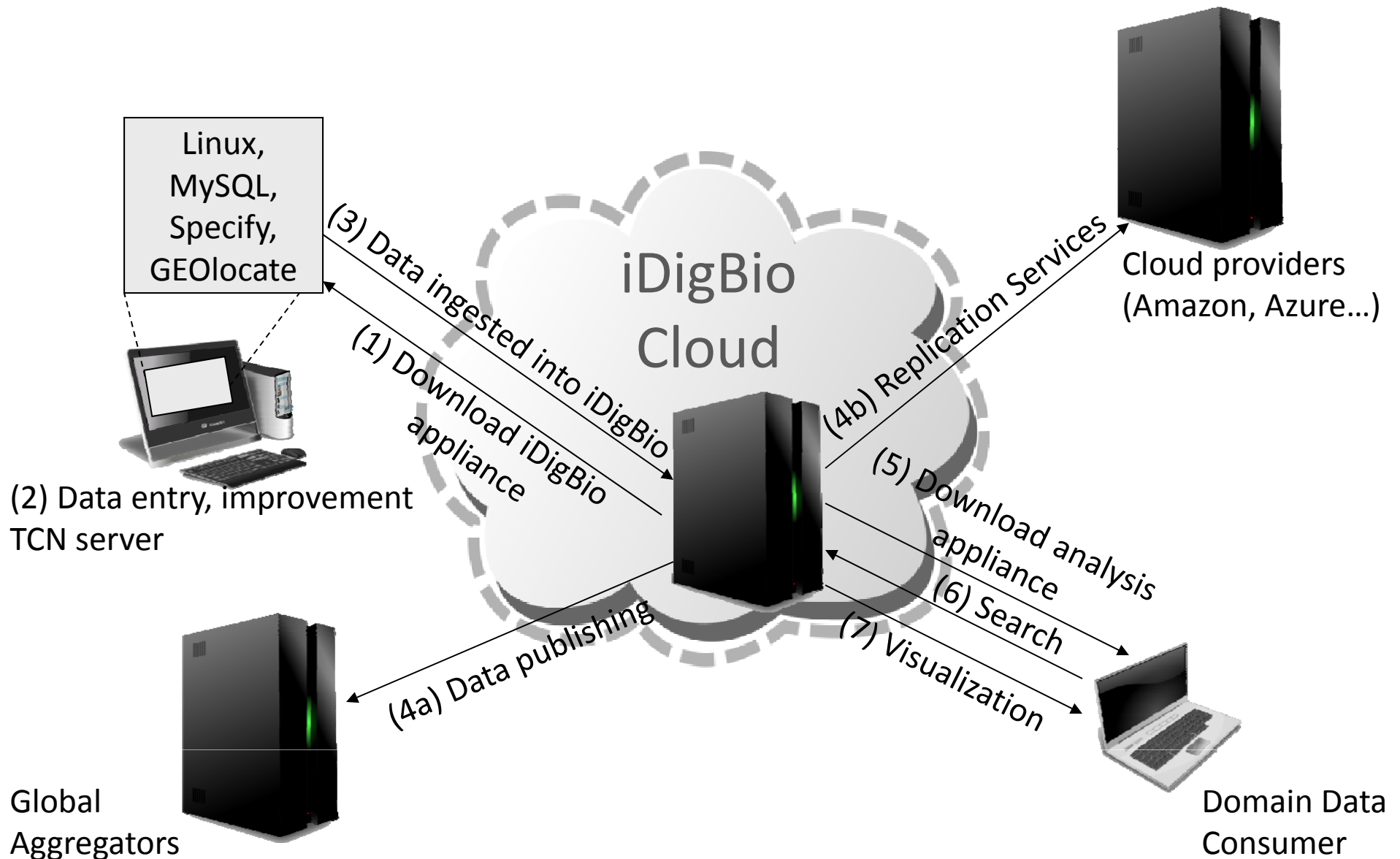
**iDigBio Portal**

**API Gateway**

**Textual Data (RIAK)**

**Filter Set Query interface**

**EXIF extraction**

**Image Data (SWIFT)**

# (Very) Long-term Goals

# Virtual appliance cycle



Requirements, standards

Domain expert
iDigBio

Collections
Community

download

instantiate

Users at
TCNs

# Toolbox Workflow Example



Linux, MySQL, Specify, GEOlocate

(2) Data entry, improvement TCN server

(3) Data ingested into iDigBio

(1) Download iDigBio appliance

iDigBio Cloud

(4b) Replication Services

Cloud providers (Amazon, Azure…)

(5) Download analysis appliance

(6) Search
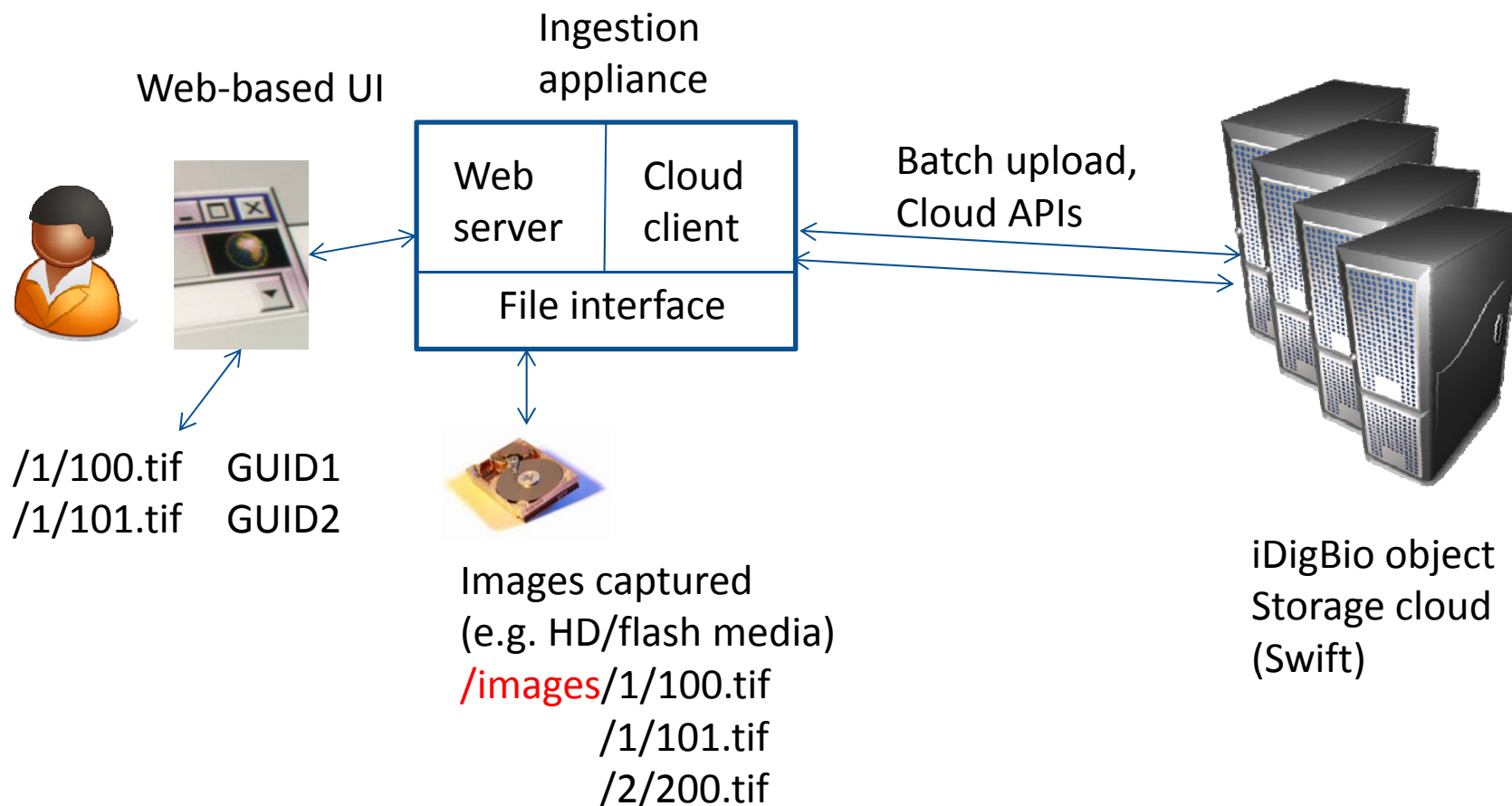
(7) Visualization
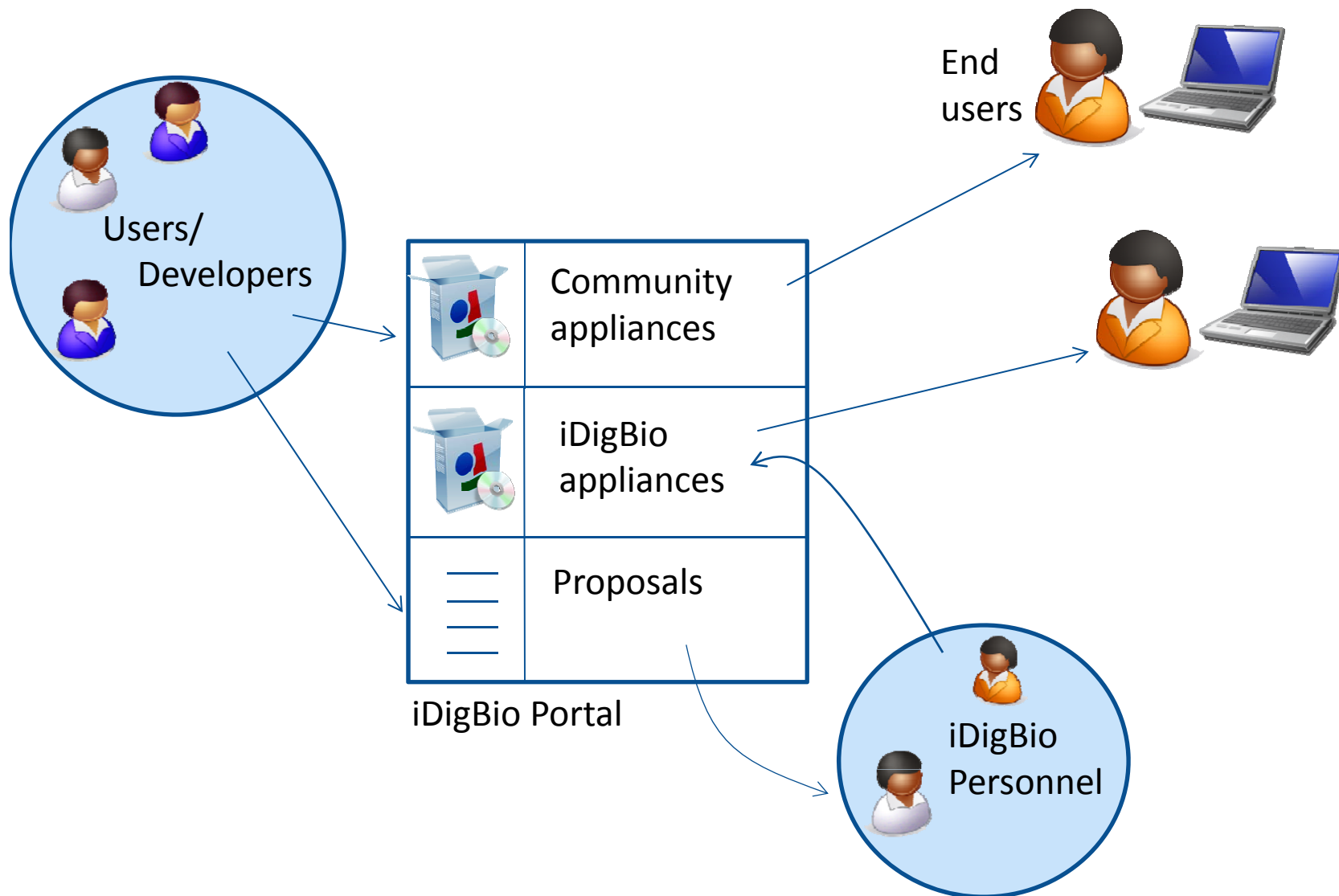
(4a) Data publishing

Global Aggregators

Domain Data Consumer

# Short term

- Facilitate data ingestion, interface with iDigBio
- Tools identified by community in workshops/groups

Web-based UI

Ingestion appliance

| Web server | Cloud client |
| --- | --- |
| File interface | |

Batch upload, Cloud APIs

/1/100.tif    GUID1
/1/101.tif    GUID2

Images captured
(e.g. HD/flash media)
/images/1/100.tif
        /1/101.tif
        /2/200.tif

iDigBio object
Storage cloud
(Swift)

# Medium-term – "Marketplace"



End users

Users/ Developers

Community appliances

iDigBio appliances

Proposals

iDigBio Portal

iDigBio Personnel

# Long-term – information processing



Users/
Developers

End users

Download

Workflows
Map/Reduce

Community appliances

Deploy

iDigBio Portal

iDigBio Personnel

Specimen Database

Advanced Computing and Information Systems laboratory

UF UNIVERSITY of FLORIDA

# Summary

- iDigBio cloud
  - Service-oriented <u>standards</u>-based cyberinfrastructure focused on the ADBC community needs
  - Scalable data management and information processing using <u>standard interfaces, data formats, protocols, tools</u>
- Toolboxes as appliances
  - Evolving collection of <u>community-selected</u> tools
  - Built-in <u>interfaces</u> for effortless iDigBio integration
  - Embedded <u>best practices and standards </u>in biocollections work
- Software re-use when open-source, well maintained, manageable, sustainable and efficient to re-purpose
- Feedback and suggestions welcome
  - <u>fortes@ufl.edu</u> and "Contacts" at idigbio.org

# Acknowledgments

- National Science Foundation
  - Judith Skog and Anne Maglia

- IDigBio team at University of Florida and Florida State University