

GLOBAL  
BIODIVERSITY  
INFORMATION  
FACILITY

Information Technologies  
and GBIF

*Francisco Pando*

*COSTECH Training Event*

*Dar Es Salaam (Tanzania)*

*25-29 February 2008*

[WWW.GBIF.ES](http://WWW.GBIF.ES)



# Summary

---

- GBIF objectives and vision
- Data network
- Data schemas
- Points and names
- Protocols
- Database registration
- Portals
- Strategies, options and implementations when sharing data

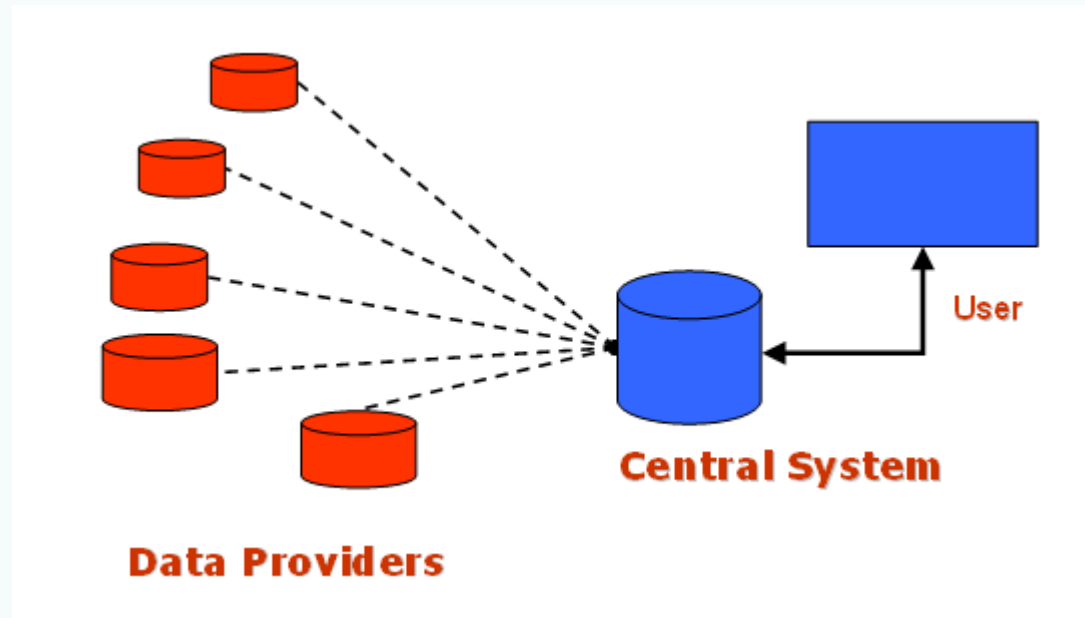
# GBIF Objectives

---

GBIF has the objective of making all information about all known living organisms available on the Internet at a global level.

In other words, GBIF tries to change this old idea of “My data is mine, look at my results” (in science and management) into “Everybody’s data for everybody”

# Data model: centralised networks



# Distributed data network

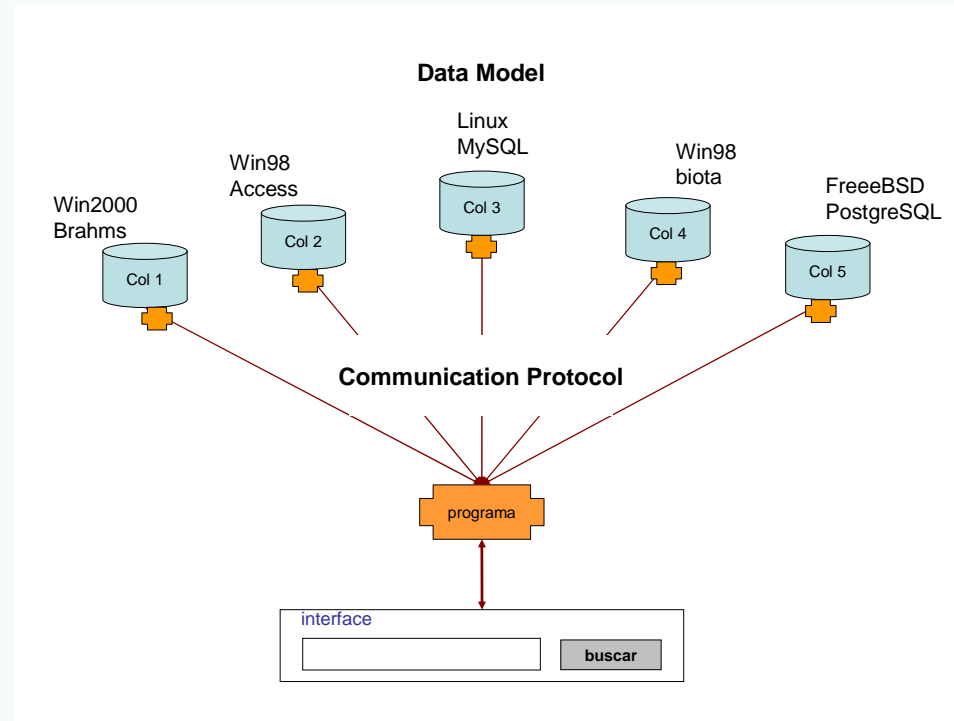
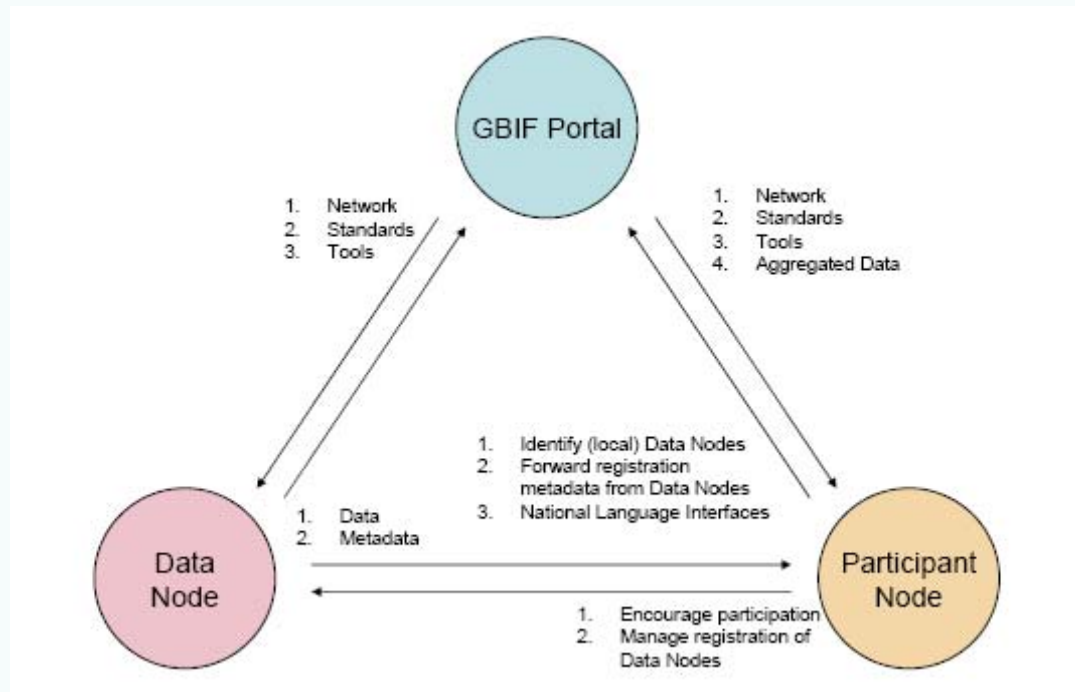


Figure 4. Diagram showing the complexity of integrating data from biological collections

# Network elements



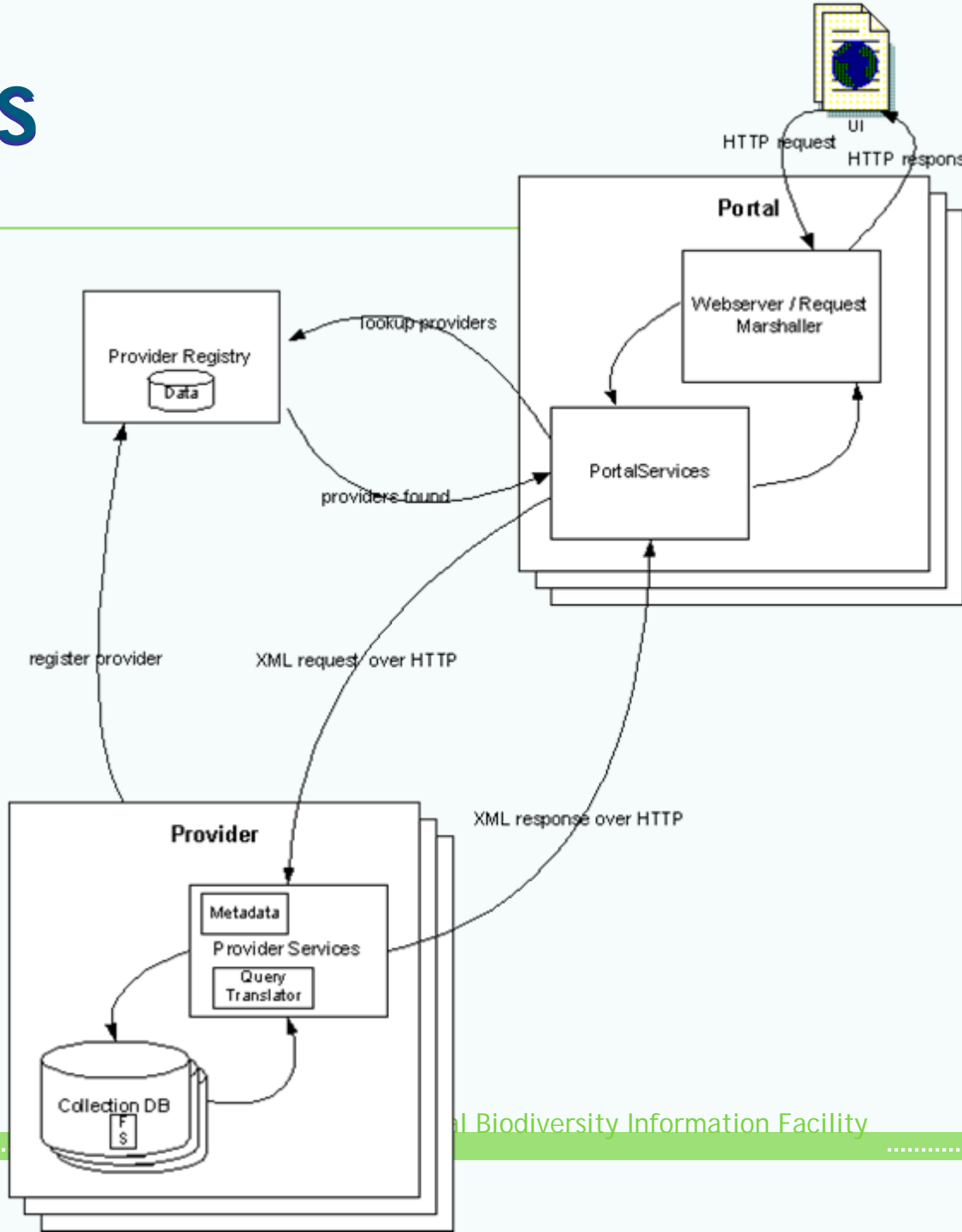
**Figure 5. GBIF Network: major classes of nodes**

GBIF is responsible for running the network, establishing standards, and developing tools. The portal is the hub for the development of any service that must be centralized such as the registry of metadata and for serving data from the biodiversity data index to the end user.

# Further details



- Protocol
- Provider
- Portal
- Registry



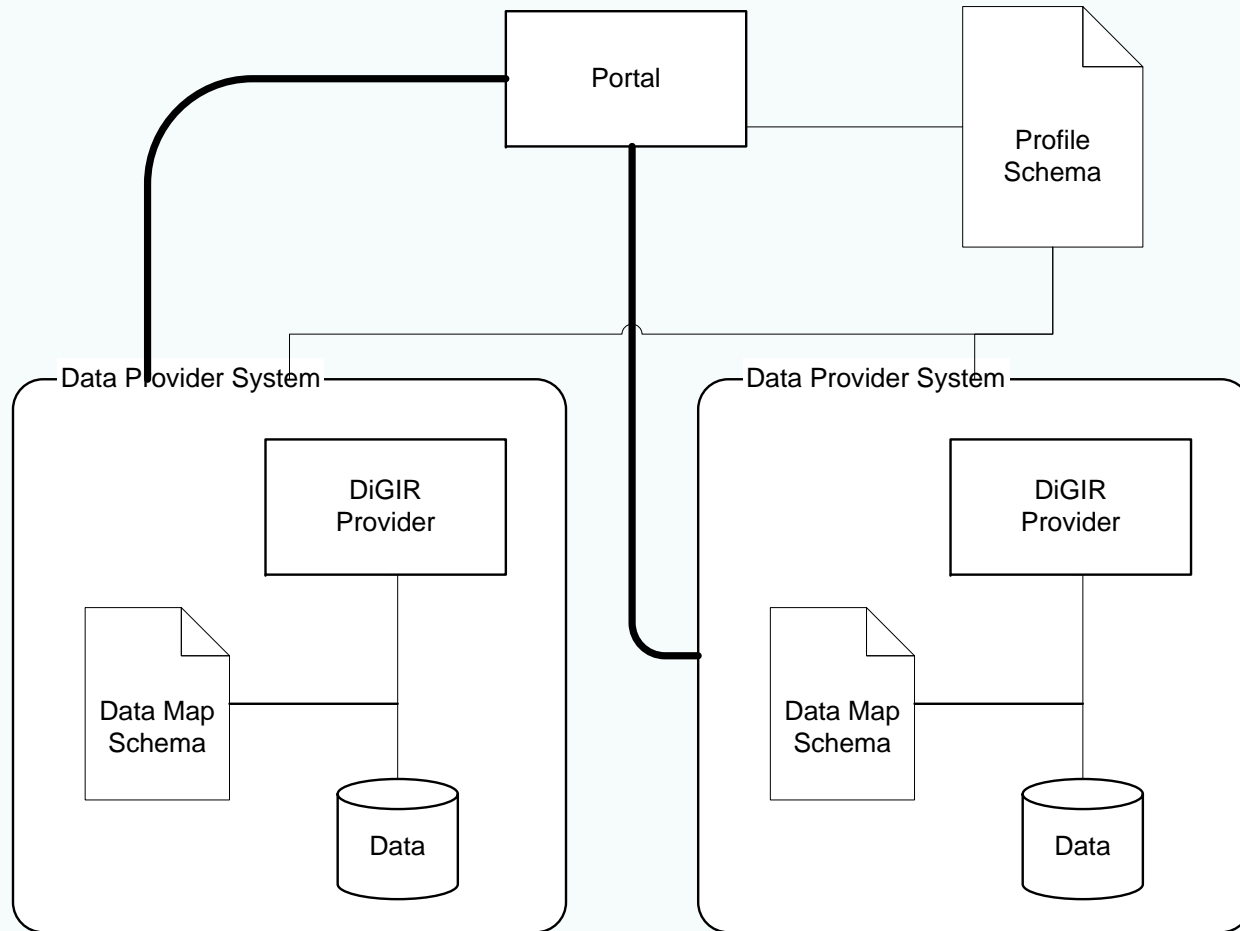
# Technologies used

---

- XML data exchange based on Providers, Services and Biodiversity Data Records
- UDDI registry for technical (access) metadata
- Descriptive metadata retrieved through service interfaces
- Specimen/observation exchange using DiGIR-Darwin Core or BioCASE-ABCD
- Taxonomic name data from Catalogue of Life (annual checklist for first release, moving to service-based approach as possible)
- Java (and JSP) components being developed centrally for GBIF Portal
- Current portal development using Tomcat, Xerces, Log4J, MySQL
- Components to be packaged for reuse as appropriate



# Data mapping



# Data schemas

---

- Darwin Core
  - Simple (50 elements)
  - The basic unit is the record
  - 500 databases in data.gbif.org
- ABCD
  - Elaborated and detailed content (+500 elements)
  - The basic unit is the file
  - 170 databases in data.gbif.org

# The base for unified access:

- Common profile:  
Every database is translated into a “common profile”, which is a table with a standardised field list that can be queried in a unified manner
- Standards:
  - “Darwin Core”
  - ABCD Schema
  - [www.tdwg.org](http://www.tdwg.org)

The Species Analyst - Database Scan Results - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Address <http://habanero.nhm.ukans.edu/TSA/scanDB.asp?tar> Go

| Attr#              | Name            |   |
|--------------------|-----------------|---|
| <a href="#">1</a>  | ScientificName  | Genus [+ " " + species [+ " " + subspecies]]. This                              |
| <a href="#">2</a>  | Kingdom         | The kingdom to which the organism belongs.                                      |
| <a href="#">3</a>  | Phylum          | The phylum (or division) to which the organism                                  |
| <a href="#">4</a>  | Class           | The class name of the organism.   |
| <a href="#">5</a>  | Order           | The order name of the organism.   |
| <a href="#">6</a>  | Family          | The family name of the organism.  |
| <a href="#">7</a>  | Genus           | The genus name of the organism.   |
| <a href="#">8</a>  | Species         | The species name of the organism.   |
| <a href="#">9</a>  | Subspecies      | The subspecies name of the organism.  |
| <a href="#">10</a> | InstitutionCode | A unique identifier for you institution.  |
| <a href="#">11</a> | CollectionCode  | Unique identifier for the collection within the ins                             |
| <a href="#">12</a> | CatalogNumber   | Unique identifier for the specimen record within                                |
| <a href="#">13</a> | Collector       | The name of the collector or collectors that were (observation) from the field. |
| <a href="#">14</a> | Year            | The year (four digit) in which the specimen was                                 |
| <a href="#">15</a> | Month           | The month of the year (1..12) in which the speci                                |
| <a href="#">16</a> | Day             | The day of month that the specimen was collec                                   |
| <a href="#">17</a> | Country         | The country or major political unit (ocean) from                                |
| <a href="#">18</a> | StateProvince   | The state, province or region (i.e. next political r collected.                 |
| <a href="#">19</a> | County          | The county (or shire, or next political region sma                              |
| <a href="#">20</a> | Locality        | The locality description (place name plus option specimen was collected.        |
| <a href="#">21</a> | Longitude       | The longitude in decimal degrees of the locatio                                 |
| <a href="#">22</a> | Latitude        | The latitude in decimal degree of the location fr                               |

# Darwin Core, current developments

- Digital Image support

| References Elements       |   |
|---------------------------|---|
| <u>ImageURL</u>           | A Universal Resource Locator reference to digital image associated with the specimen or observation.                                  |
| <u>RelatedInformation</u> | Free text references to information not delivered by the Darwin Core conceptual schema, including URLs to specimen publications, etc. |

- Extensions

# Darwin Core, extensions

## Geospatial Extension

Schema document and element definition table for the geospatial extension to the Darwin Core 2.

## Curatorial Extension

Schema document and element definition table for the curatorial extension to the Darwin Core 2.

## Paleontology Extension

Schema document and element definition table for the paleontological extension to the Darwin Core 2.

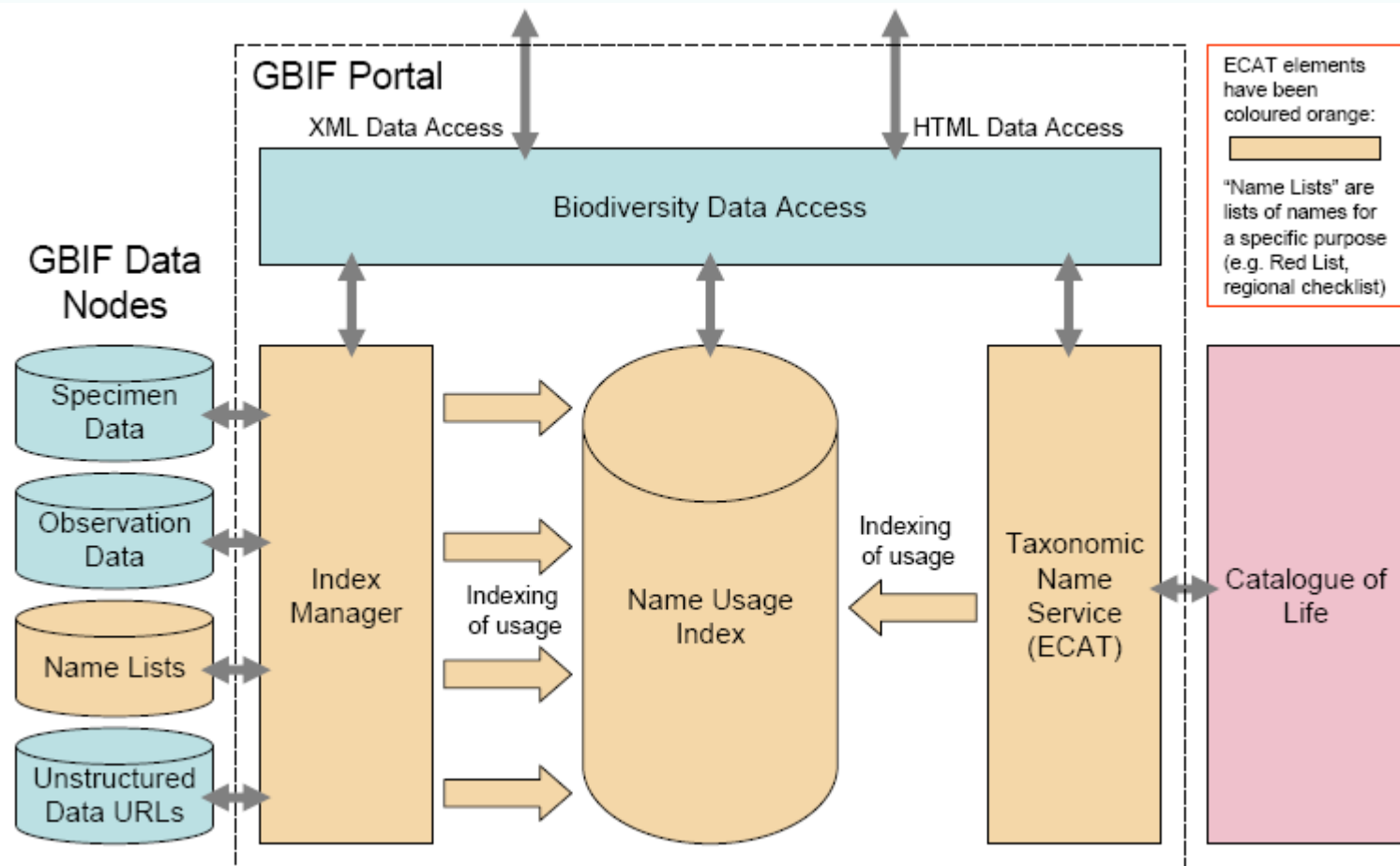
## Microbial Extension

An extension to the Darwin Core 2 for microbiology culture collections. Posted by Renato de Giovanni. Integrated with the DiGIR protocol.

## Observation/Monitoring Extension

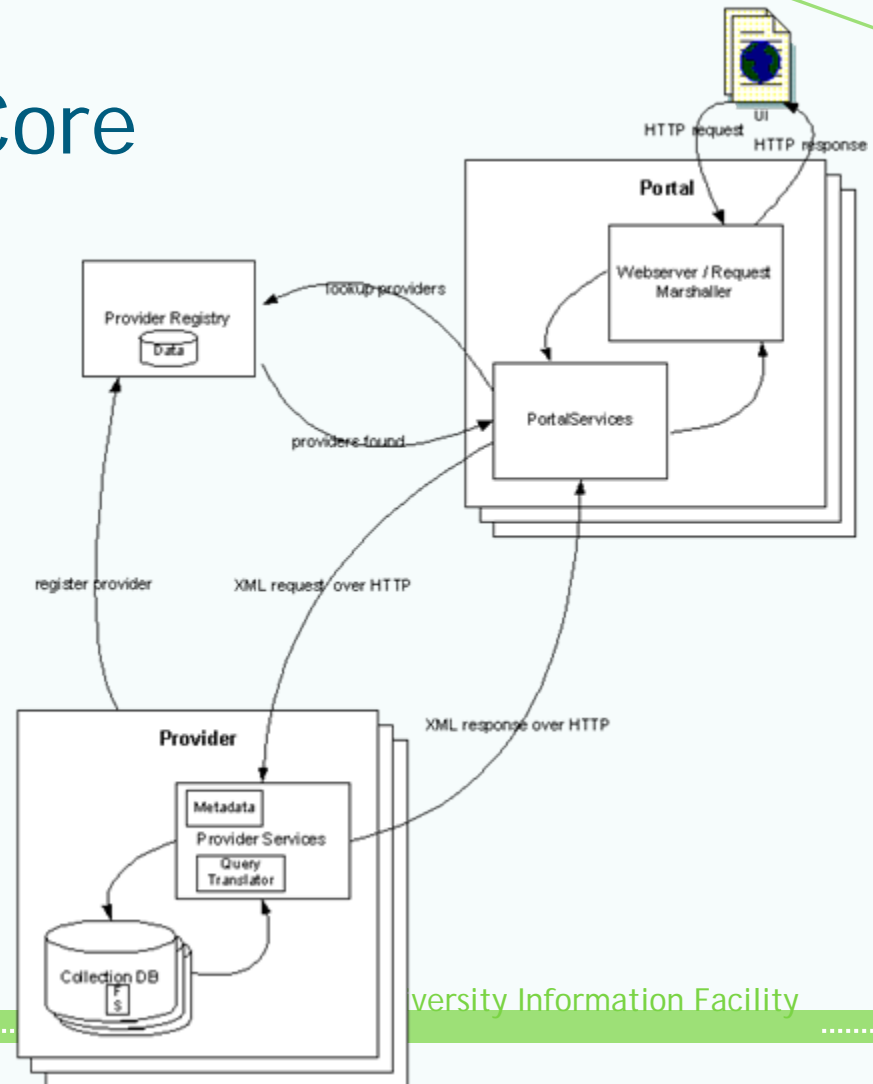
Link to the Avian Knowledge Network BMDE schema page which poses an extension to the Darwin Core 2 for observations and monitoring.

# Names and specimens integration



# Protocols

- DIGIR for Darwin Core
- Biocase for ABCD



# Life after DIGIR

---

- TAPIR
  - <http://www.gbif.org/News/NEWS1129877273>
  - <http://ww3.bgbm.org/protocolwiki/>
  - Unify protocols in GBIF data network
  - Registry expansion (UDDI), thematic networks support, national portals, Darwin Core extensions
  - Toolkit for data portals being developed



# Registry

To highlight:

- A characteristic name
- A user-oriented description
- Additional data use constraints
- How to cite this resource

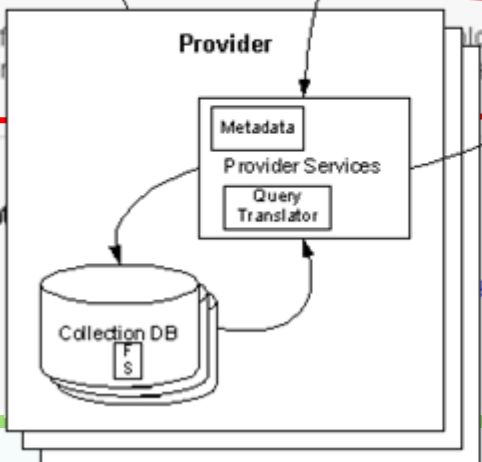
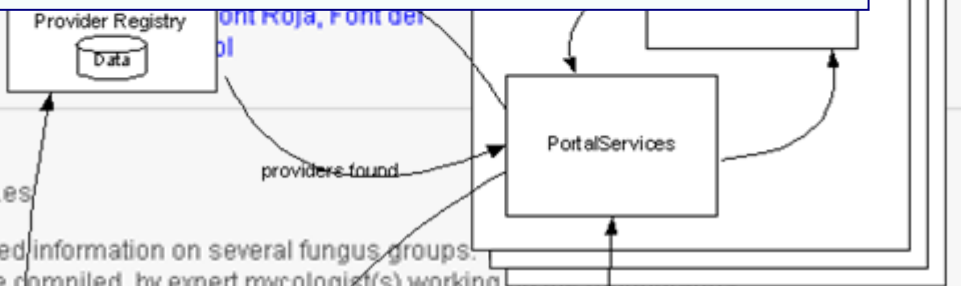
|         |                                       |              |    |
|---------|---------------------------------------|--------------|----|
| 15865-2 | Echinostelium brooksii<br>K.D.Whitney | Nov 2, 1984  |    |
| 16059-1 | Echinostelium brooksii<br>K.D.Whitney | Feb 1, 1984  |    |
| 38579-1 | Echinostelium brooksii<br>K.D.Whitney | Jul 26, 1997 | ES |

**Service** GBIF-Spain (taray.csic.es) Datasources provided by GBIF.es

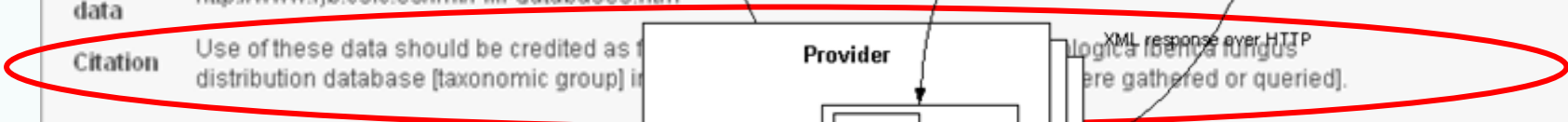
**Resource** Flora Mycologica Iberica  
The databases contains detailed information on several fungus groups. been revised, and in many time compiled, by expert mycologist(s) working for the Flora Mycologica Iberica project (FMI). Thus, the quality and currentness of the information here presented is high. Records were gathered or queried from the following sources: <http://www.rjb.csic.es/fmi/FMI-databases.htm>

**Use of data** <http://www.rjb.csic.es/fmi/FMI-databases.htm>

**Citation** Use of these data should be credited as follows: <http://www.rjb.csic.es/fmi/FMI-databases.htm> distribution database [taxonomic group] in the GBIF registry



| Latitude | Longitude | User feedback |
|----------|-----------|---------------|
| 41.64    | 2.52      | ✉             |
| 41.82    | 2.4       | ✉             |



# GBIF Data Portals

---

- New Data Portal - [data.gbif.org](http://data.gbif.org)
- Nodes Portal
- Old prototype - [www.gbif.net](http://www.gbif.net)  
(no longer available)

# New Data Portal

---

- Similar to [www.biologybrowser.com](http://www.biologybrowser.com)
- Web services
- API interface
- Added indexing and validation services

# New Data Portal

119,662,684 occurrence records from 235 data providers - GBIF Portal - Windows Internet Explorer

http://data.gbif.org/welcome.htm;jsessionid=A476A1E5A3E414ACDE76A0E71D0A2899

GLOBAL BIODIVERSITY INFORMATION FACILITY

SPECIES COUNTRIES DATASETS OCCURRENCES SETTINGS ABOUT

... free and open access to biodiversity data

Search  
species/country/dataset

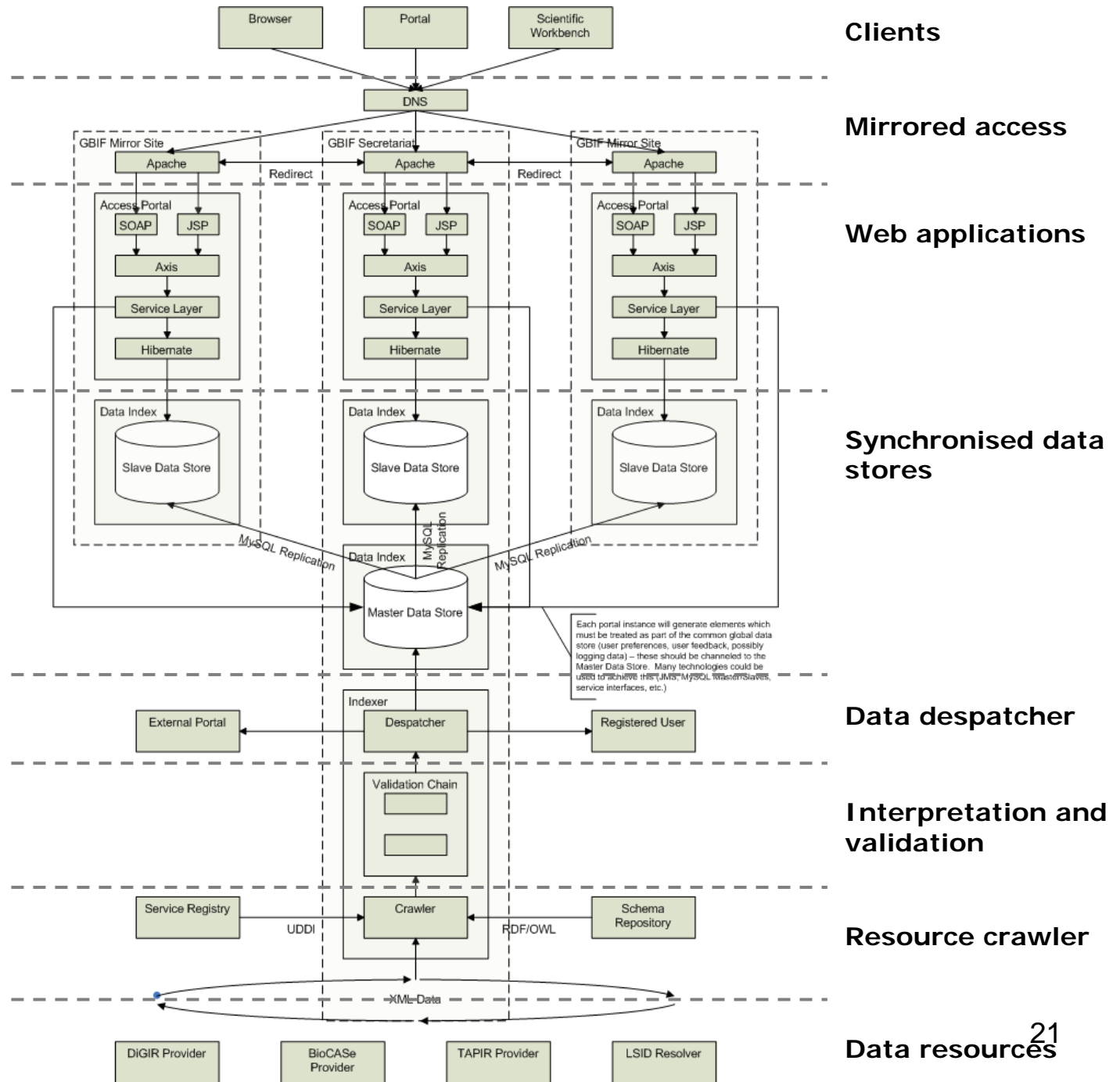
Welcome to the GBIF Data Portal  
Access millions of data records shared via the GBIF network.  
To learn how to use this site, please see [About](#).  
To tune this site for smaller displays, see [Settings](#).

**Explore Species**  
Find data for a species or other group of organisms.  
**Species**  
Information on species and other groups of plants, animals, fungi and micro-organisms, including species occurrence records, as well as classifications and scientific and common names.  
**Example species:**  
*Puma concolor* (Linnaeus, 1771)

**Explore Countries**  
Find data on the species recorded in a particular country.  
**Countries**  
Information on the species recorded in each country, including records shared by providers from throughout the GBIF network.  
**See data for:**  
France

**Explore Datasets**  
Find data from a data provider, dataset or data network.  
**Datasets**  
Information on the data providers, datasets and data networks that share data through GBIF, including summary information on 1693 datasets from 235 data providers.  
**Latest dataset added:**  
Lichenes North-Eastern Poland

# Portal architecture (new version D. Hobern)



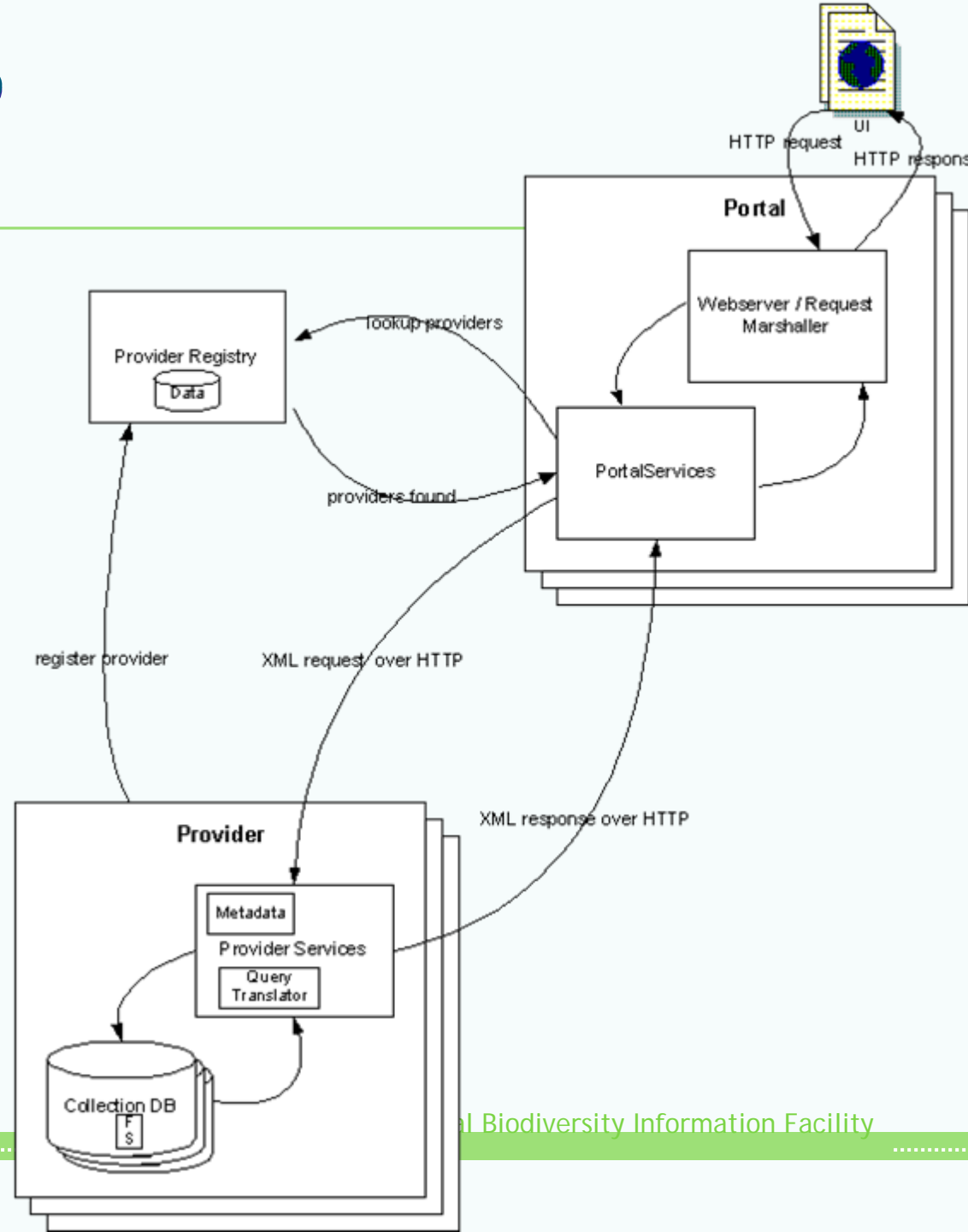
# Nodes Portal

---

- Under development
- Such as the nodes can provide:
  - Data from the collections in their area (country)
  - Data relevant under their scope
- Adaptable to their needs (language, common names,...)
- Help nodes to serve their communities


# Clear enough?

- Protocol
- Provider
- Portal
- Registry



# Old prototype

HOME | GBIF | BROWSE TAXONOMY | SEARCH | DATA PROVIDERS | COUNTRIES | DATA USE



## Prototype data portal Global Biodiversity Information Facility

---

**BROWSE TAXONOMY**

- Kingdom: [Animalia](#)
- Kingdom: [Archaea](#)
- Kingdom: [Bacteria](#)
- Kingdom: [Chromista](#)
- Kingdom: [Fungi](#)
- Kingdom: [Plantae](#)
- Kingdom: [Protozoa](#)
- Kingdom: [Viruses](#)

**DATA PROVIDERS**

- [Taxonomic name](#)
- [Specimen/observation](#)

**COUNTRIES**

- [Data by country](#)

**GBIF REGISTRY**

- Data providers 167
- Records 89656226
- [Become a data provider](#)
- [Advanced search](#)

**ABOUT GBIF**

- [Press](#)
- [What is GBIF?](#)
- [Demonstration Projects](#)
- [GBIF's Current Members](#)
- [How to Join GBIF](#)
- [GBIF Directory](#)
- [Documents & History](#)
- [Ebbe Nielsen Prize](#)
- [Symposia & Workshops](#)
- [Work Programme](#)

**ABOUT GBIF**

The Global Biodiversity Information Facility (GBIF) is an international non-profit organisation to provide free and universal access to data regarding the world's biodiversity. To learn more about GBIF itself, visit the GBIF Communications Portal: [www.gbif.org](http://www.gbif.org)




A wide range of countries and organisations participate in GBIF and have made their data available here. These bodies maintain ownership for all of the data they share. Any feedback provided through this web site will be passed back to the data provider concerned.

**ABOUT THIS PORTAL**

This portal is a prototype service providing access to the two types of data which are already being shared through the GBIF Network:

- Taxonomic names.** GBIF developing an 'Electronic Catalogue of Taxonomic Names'. This will provide access to authoritative information about both scientific and common names for all organisms, and will integrate data from a wide range of different organisations. The portal already includes data for over 983,000 scientific names and 253,000 common names from the Catalogue of Life Partnership Annual Checklist. Some names are listed with the words 'Tentative position in taxonomy'. This indicates that the name is only known to the portal from specimen/observation records and should not be treated as authoritative simply on the basis of being listed here.
- Specimens and Observations.** The GBIF Network already provides access to over 40 million records of occurrences of different organisms. Many of these relate to specimens in natural history museums and herbaria around the world, or to living cultures of micro-organisms, but at least a third come from observations of wild organisms. Wherever possible these records include information about the locality where the organisms were found and are used to generate maps of the distribution of these occurrences. Counts of occurrence records are listed against the organism names to which they apply. For genera and taxa of higher ranks, these counts include only those occurrences which have been identified to the taxon concerned. For species these counts include all occurrences for the species and also for any included infraspecific taxa, as well as for any known synonyms.

This is a work in progress. Please explore what is already present and visit us again in the coming months as we integrate more data and provide more flexible interfaces to search and browse the data. The following icons are used on many pages in the portal:

-  Send feedback on a data item to the original data provider
-  Get further details
-  Download data

**USE OF DATA**

GBIF Participants have made their data available for use according to the terms of the [GBIF Data Use Agreement](#). Please understand these terms before using GBIF data: [GBIF Data Use Agreement](#)

**SEARCH**

Search for name:   Country/Territory:

Selecting a country limits results to those taxa which specimen/observation records show have been identified from the country concerned.

- Search by scientific name (any rank)
- Search by common name in any language
- Search by English name
- Find names starting with search string (minimum 3 characters)
- Find names containing search string (minimum 3 characters)
- Find exact matches

**Google Earth**

You can search and visualise distributions using [Google Earth link](#) to GBIF data

**Search only by scientific or vernacular names (optional country filter)**



# Old prototype

## Additional names from specimen/observation data (unreviewed)

**Rank** **Name**  
 Subspecies *Balaenoptera musculus pribilofensis*

## Common names

**Language** **Name**  
 English **Blue whale**  
 French **Rorqual bleu**

**Countries from which species is recorded**

**Authority**  
 Catalogue of Life: Integrated Taxonomic Information System

## Specimens/observations

Including records from: [Antarctica](#); [Australia](#); [Canada](#); [France](#); [Mexico](#); [Norway](#); [Portugal](#); [Spain](#); [United Kingdom](#); [United States](#)

## Service

- ABRS DiGIR Provider ([www.deh.gov.au](http://www.deh.gov.au))
- Australian Antarctic Data Centre ([aadc-maps.aad.gov.au](http://aadc-maps.aad.gov.au))
- Australian Antarctic Data Centre ([aadc-maps.aad.gov.au](http://aadc-maps.aad.gov.au))
- Avian Knowledge Network ([akn.ornith.cornell.edu](http://akn.ornith.cornell.edu))
- California Academy of Sciences (CAS) ([www.calacademy.org](http://www.calacademy.org))
- EUNIS 2 DiGIR Provider ([woodpecker.eea.eu.int](http://woodpecker.eea.eu.int))
- Los Angeles County Museum of Natural History (LACM)
- MCZ-Harvard University Provider ([digir.mcz.harvard.edu](http://digir.mcz.harvard.edu))
- Museum of Natural Science - Louisiana State University Mammal Collection (LSUMZ)
- Museum of Vertebrate Zoology (MVZ) (128.32.146.144)
- National Chemical Laboratory ([digir.indobis.org](http://digir.indobis.org))
- NatureServe ([services.natureserve.org](http://services.natureserve.org))
- NLBIF (145.18.162.60)
- OBIS/DIGIR Data Provider Server ([www.iobis.org](http://www.iobis.org))
- OBIS/DIGIR Data Provider Server ([www.iobis.org](http://www.iobis.org))
- OBIS/DIGIR Data Provider Server ([www.iobis.org](http://www.iobis.org))
- OBIS/DIGIR Data Provider Server ([www.iobis.org](http://www.iobis.org))
- OZCAM Provider ([digir.austmus.gov.au](http://digir.austmus.gov.au))
- Royal Ontario Museum ([digir.rom.on.ca](http://digir.rom.on.ca))
- University of Washington Burke Museum (UWBM) ([biology.burke.washington.edu](http://biology.burke.washington.edu))

**Total**

**Georeferenced records**

**Occurrence Detail**

**User feedback**

**User feedback**

| AM          | Lat         | Long        |
|-------------|-------------|-------------|
| 3           | 0           | 0           |
| 4           | 4           | 4           |
| 1968        | 1968        | 1968        |
| 34          | 12          | 12          |
| 8           | 0           | 0           |
| 21          | 0           | 0           |
| 3           | 2           | 2           |
| 5           | 0           | 0           |
| 1           | 0           | 0           |
| 1           | 0           | 0           |
| 1           | 1           | 1           |
| 1           | 0           | 0           |
| 2           | 0           | 0           |
| 1137        | 1051        | 1051        |
| 4           | 4           | 4           |
| 1           | 1           | 1           |
| 171         | 171         | 171         |
| 9           | 9           | 9           |
| 4           | 0           | 0           |
| 7           | 1           | 1           |
| <b>3385</b> | <b>3224</b> | <b>3224</b> |

[Contact info](#) | [Webmaster](#)

**Data resources with specimens or observations**

Global Biodiversity Information Facility  
**Download records**

# Standards for data integration and interoperability



**International Working Group on Taxonomic  
Databases**

**International Union of Biological Sciences Taxonomic  
Database Working Group**

Select a topic:

[Home Page](#)  
[What's New](#)  
[Listservs](#)  
[Standards](#)  
[Subgroups](#)  
[Secretariat](#)  
[2005 Officers](#)  
[Constitution](#)  
[Membership](#)  
[Related Links](#)  
[Comments](#)  
[History](#)  
[Archives](#)  
  
[2004 Meeting](#)  
[2005 Meeting](#)

[TDWG STANDARDS](#)

<http://www.tdwg.org/standrds.html>

[Authors of plant names](#)

[Botanico-periodicum-huntianum](#)

[Botanico-periodicum-huntianum/supplementum](#)

[Economic botany data collection standard](#)

[Floristic regions of the world](#)

[Herbarium information standards and protocols for interchange of data](#)

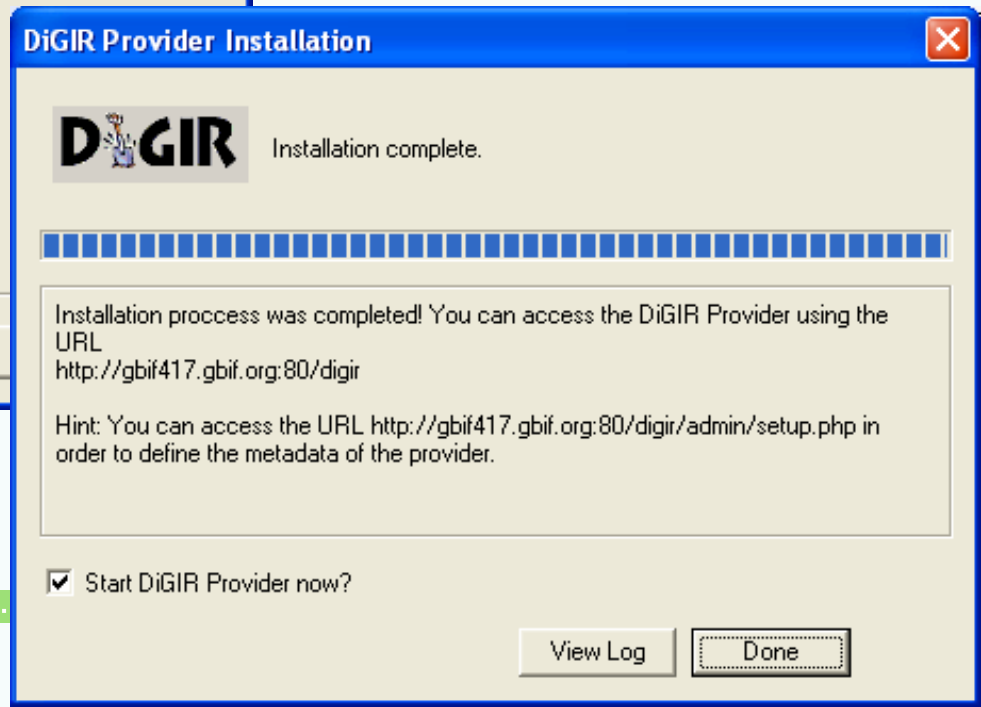
[Index Herbarium, Part 1: The herbaria of the world](#)

[International transfer format for botanic garden plant records](#)

[Plant names in botanical databases](#)



- "Turn-key package"
- Based on PHP and code from DiGIR project
- For Linux and Windows
- Register with GBIF UDDI



- Technical support: helpdesk@gbif.org

# About access to databases via Internet

---

- Whose data is this?
- All, Everything?
- How clean is it?
- How?

# Whose data is this?

Address  http://www.secretariat.gbif.net/portal/ecat\_search.jsp?nextTask=ecat\_search.jsp

## Global Biodiversity Information Facility (GBIF)

### Data Use Agreement

#### Background

The goals and principles of making biodiversity data openly and universally available have been defined in the Memorandum of Understanding on GBIF (MoU; see the relevant excerpts in [Annex](#)).

The Participants who have signed the MoU have expressed their willingness to make biodiversity data available through their nodes to foster scientific research development internationally and to support the public use of these data.

GBIF data sharing should take place within a framework of due attribution.

Therefore, using data available through the GBIF network requires agreeing with the following:

#### 1. Data Use Agreements

1. The quality and completeness of data cannot be guaranteed. Users employ these data at their own risk.
2. Users shall respect restrictions of access to sensitive data.
3. In order to make attribution of use for owners of the data possible, the [identifier](#) of ownership of data must be retained with every data record.
4. Users must publicly acknowledge, in conjunction with the use of the data, the data providers whose biodiversity data they have used. Data providers may require additional attribution of specific collections within their institution.
5. Users must comply with additional terms and conditions of use set by the data provider. Where these exist they will be available through the metadata associated with the data.

#### 2. Definitions

- GBIF Participant: Signatory of the GBIF-establishing Memorandum of Understanding (MoU).
- GBIF Secretariat: Legal entity empowered by the GBIF Participants to enter into contracts, execute the Work Programme, and maintain the central services for the GBIF network.
- GBIF network: The infrastructure consisting of the central services of the GBIF Secretariat, Participant Nodes and data providers. Making data available through GBIF network means registering and advertising the pertinent services via the GBIF central services..

# All, everything?

---

- Data provider keeps control over what is made available:
  - Tests
  - Decides to dilute the record precision for endangered or economically relevant species
  - Does not publish data from current research

# How clean is it?

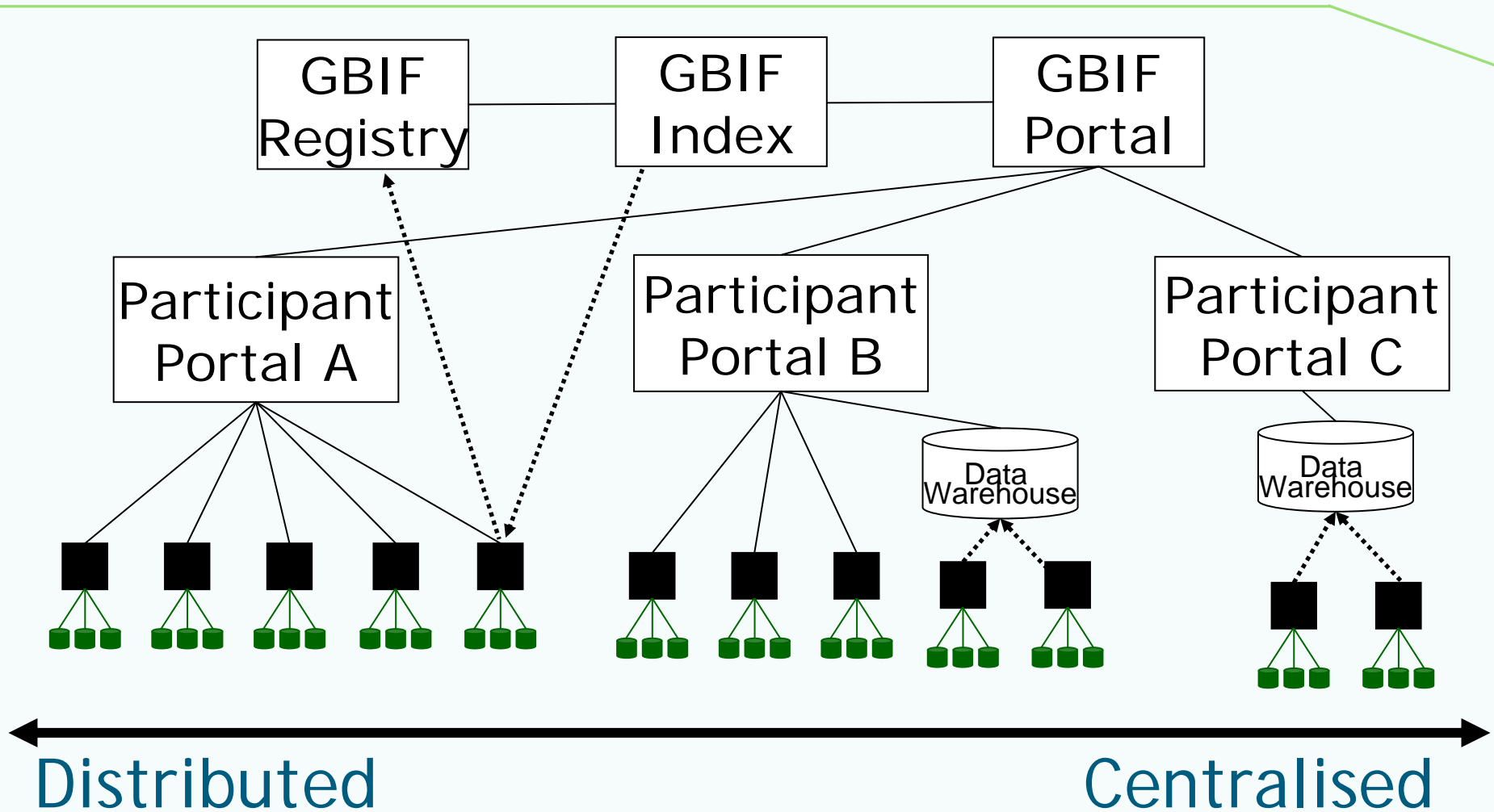
---

- Perfection does not exist
- Neither bad data exist; a record is not good or bad intrinsically, its validity depends on its use
- Making data publicly available helps to improve them
- There are some tools to improve data and GBIF is working on this:

<http://www.secretariat.gbif.net/datatester/index.jsp>

[http://www.gbif.org/prog/digit/data\\_quality](http://www.gbif.org/prog/digit/data_quality)

# How?





# At your disposal:

---

Francisco Pando  
Spanish GBIF Node  
Royal Botanical Garden (CSIC)  
Plaza de Murillo, 2  
E-28014 Madrid, Spain

[pando@gbif.es](mailto:pando@gbif.es)

Tél.: + 34 91 420 30 17

