# Latin America and the Caribbean
# Scientific Data Management Workshop

17-18 April 2018

Brazilian Academy of Sciences, Rio de Janeiro, Brazil

©Museum of Tomorrow

# PPBio's Metacat Data Repository

- Tim Vincent


- Instituto Nacional de Pesquisas da Amazonia
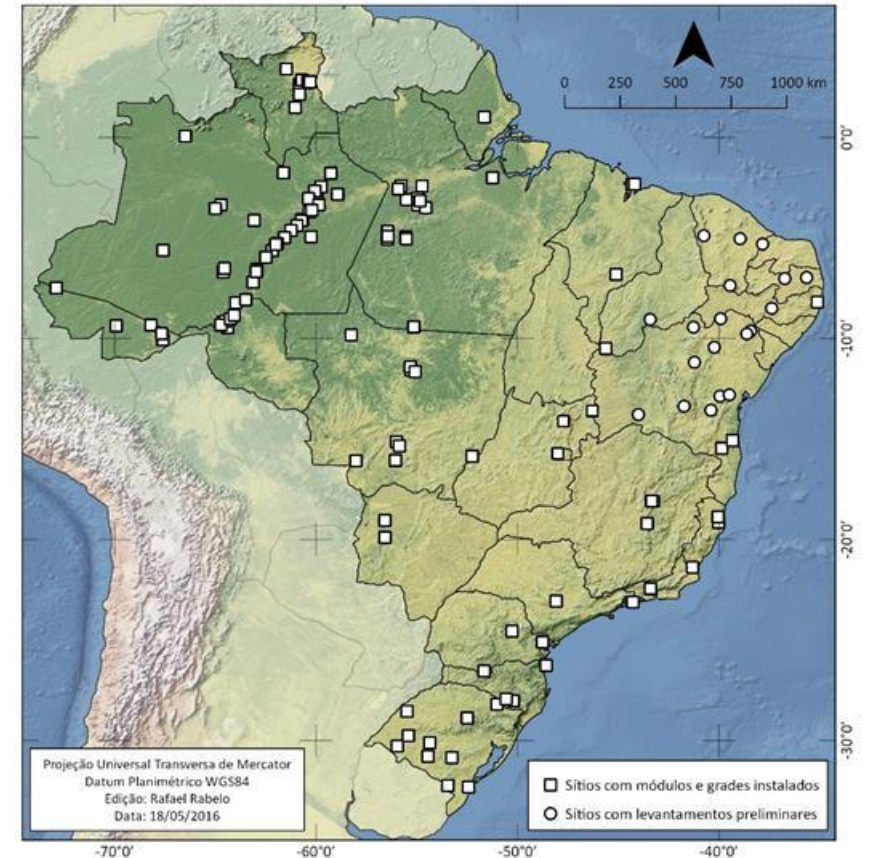

- CENBAM


- PPBio

# PPBio

- The Biodiversity Research Program was created in 2004 with the aims of furthering biodiversity studies in Brazil, decentralizing scientific production from already-developed academic centers, integrating research activities and disseminating results across a variety of purposes, including environmental management and education.

- The Executive Nucleus of the Western Amazon Biodiversity Research Program (PPBio-AmOc) maintains a data repository used by researchers involved with several regional program cores and some other executing nuclei.

# Why did PPBio implement a Metacat biodiversity data repository?

- Although this is explained in some detail in chapter 7 of the book, "Biodiversity and Integrated Environmental Monitoring" by William Magnusson et al. (2013), some of the important points are listed below:

- The most important thing to make available is the metadata.

- Cost

- Easy to upload and download. Data is in a common format.

- Ability to review and update the data package (metadata and data).

# SiBBr

- The repository's dataset are collected and also made available on the Brazilian Biodiversity Information System (SiBBr).

- SiBBr's Ecological data repository also uses Metacat.

Unlike worksheet style data repositories, Metacat allows a wide variety of information to be uploaded; maps, scans of raw data such as field notes, photographs and so on.

# Morpho

Morpho is the Java based application that we use for creating, uploading and editing Metacat datasets.

# Morpho

Morpho is the Java based application that we use for creating, uploading and editing Metacat datasets.

Morpho is used for organizing and preparing the metadata and for uploading the dataset to the Metacat server. It can also used for accessing the dataset and updating or adding information at a later date or for making controlled access data publicly available.

This standalone tool is downloaded to the researchers own computers where the metadata and data can be prepared without the need for an internet connection.

# Morpho

Morpho is the Java based application that we use for creating, uploading and editing Metacat datasets.

Researchers use Morpho to prepare a dataset which they then send to a data manager who is responsible to checking that the metadata has been entered correctly and that there are no issues with the format of the data before uploading to the server.

This human intervention in the process is very important for quality control. Ensuring that metadata is entered correctly is important as it can influence how easily the data-set can be discovered in the future.

# Morpho

Morpho is the Java based application that we use for creating, uploading and editing Metacat datasets.

Natural- language metadata may describe data using an ad-hoc set of descriptive terms but there may be subsequent issues with recall. Morpho enables the data to be described using terms that facilitate retrieval, but this must be correctly done by the person entering the data.

The PPBio website provides detailed instructions on how to use Morpho, but person-to- person training is very useful since there are some small hacks that are necessary in order to resolve issues with Morpho's user interface and the terminology used for defining and describing the data.

# When there are problems however, the help you get from the NCEAS/KNB team is very good.

# NCEAS

- National Center for Ecological Analysis and Synthesis
- Metacat and Morpho are products of NCEAS.

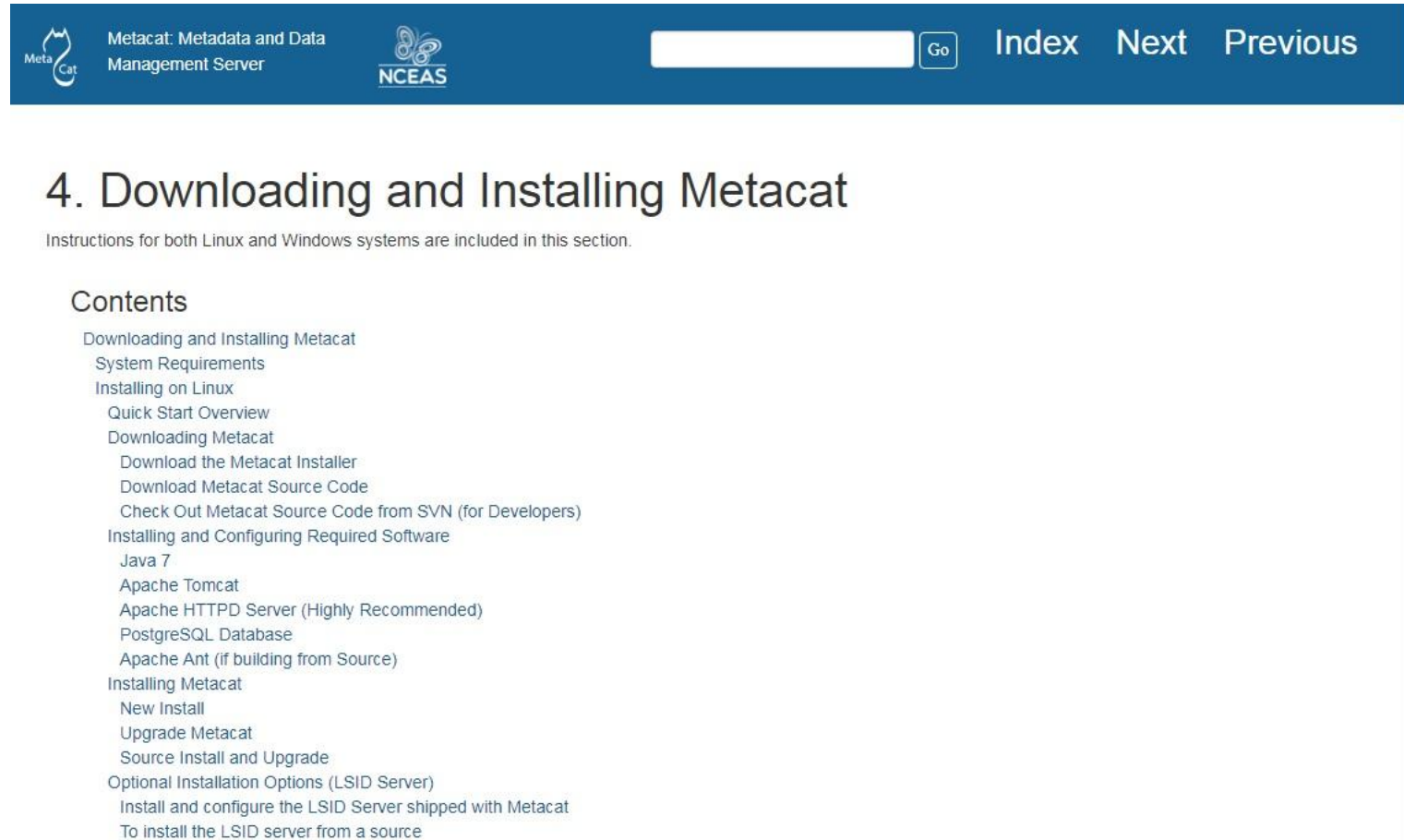

NCEAS - About

# Metacat

- "...flexible, open source metadata catalogue and data repository..."

- "...utilizes a relational database management system to store XML and associated meta-level information."

- Technical Expertise Required: Basic programming skills.

- Cost: Free

# Follow the instructions for installing Metacat...
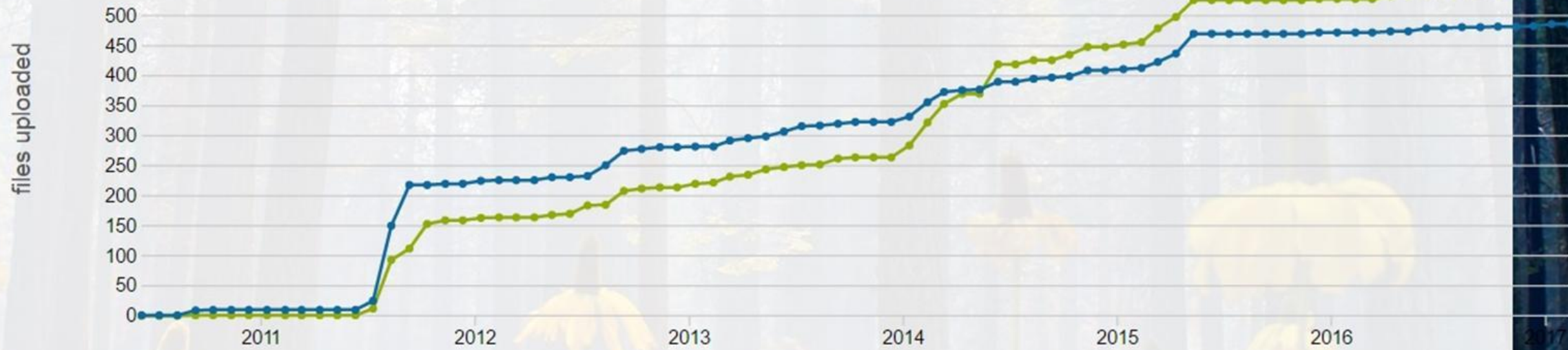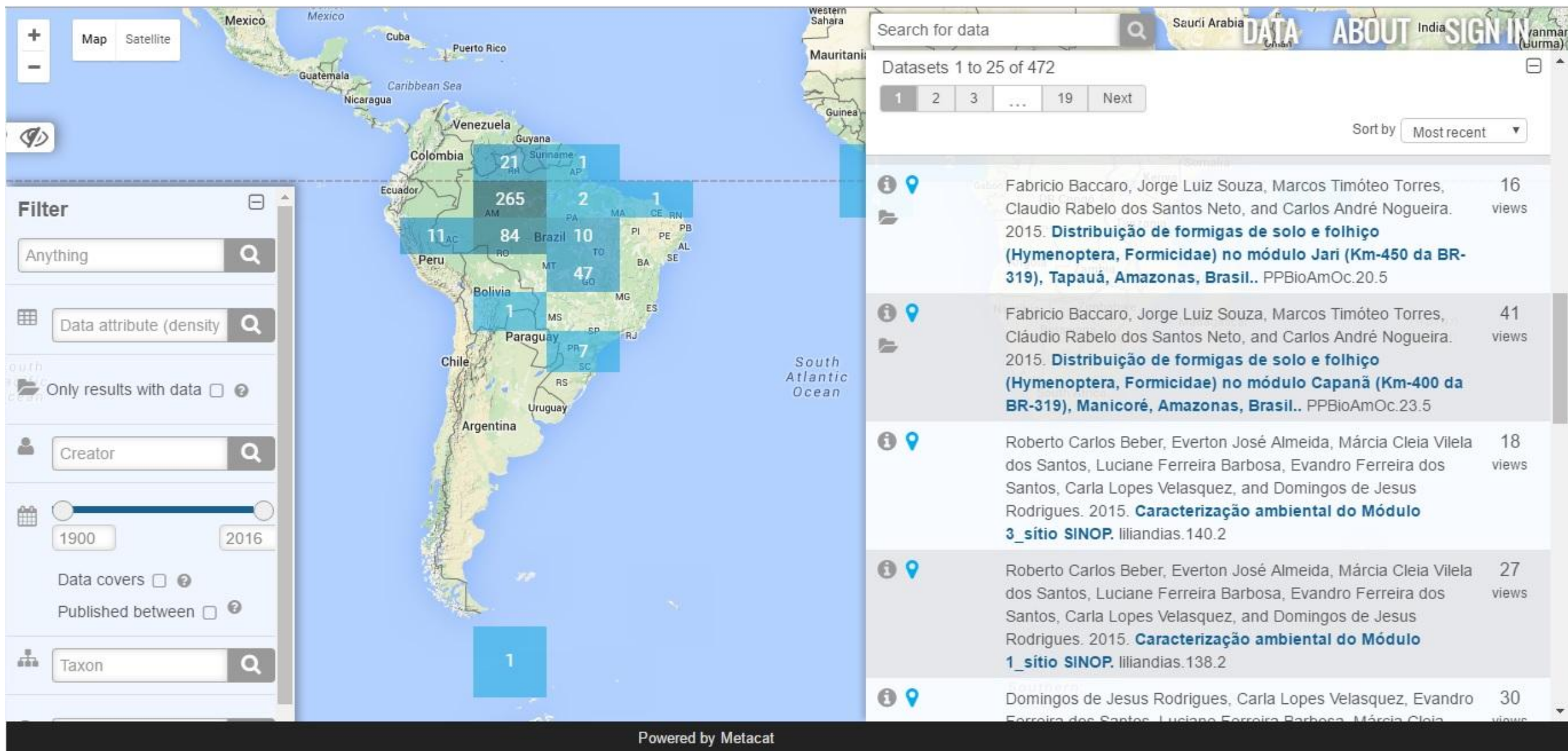
# Uploads

The number of individual metadata and data files uploaded over time. Only the first version of each file is counted.

**486** metadata

**548** data

+
−

Datasets 1 to 25 of 472

| 1 | 2 | 3 | … | 19 | Next |

Sort by Most recent ▼

21 RR
1 AP
265 AM
2 PA
1 MA
CE RN
11 AC
84 Brazil
10 TO
PI
PE PB
AL
RO
MT
BA SE
47 GO
MG
1
MS
ES
SP
RJ
PR
7 SC
RS

1

**Filter**

Anything 🔍

⊞ Data attribute (density 🔍

📁 Only results with data ☐ ❓

👤 Creator 🔍

📅
1900        2016

Data covers ☐ ❓
Published between ☐ ❓

⛓ Taxon 🔍

ℹ 📍 📁  Fabricio Baccaro, Jorge Luiz Souza, Marcos Timóteo Torres, Claudio Rabelo dos Santos Neto, and Carlos André Nogueira. 2015. **Distribuição de formigas de solo e folhiço (Hymenoptera, Formicidae) no módulo Jari (Km-450 da BR-319), Tapauá, Amazonas, Brasil..** PPBioAmOc.20.5   16 views

ℹ 📍 📁  Fabricio Baccaro, Jorge Luiz Souza, Marcos Timóteo Torres, Cláudio Rabelo dos Santos Neto, and Carlos André Nogueira. 2015. **Distribuição de formigas de solo e folhiço (Hymenoptera, Formicidae) no módulo Capanã (Km-400 da BR-319), Manicoré, Amazonas, Brasil..** PPBioAmOc.23.5   41 views

ℹ 📍  Roberto Carlos Beber, Everton José Almeida, Márcia Cleia Vilela dos Santos, Luciane Ferreira Barbosa, Evandro Ferreira dos Santos, Carla Lopes Velasquez, and Domingos de Jesus Rodrigues. 2015. **Caracterização ambiental do Módulo 3_sítio SINOP.** liliandias.140.2   18 views

ℹ 📍  Roberto Carlos Beber, Everton José Almeida, Márcia Cleia Vilela dos Santos, Luciane Ferreira Barbosa, Evandro Ferreira dos Santos, Carla Lopes Velasquez, and Domingos de Jesus Rodrigues. 2015. **Caracterização ambiental do Módulo 1_sítio SINOP.** liliandias.138.2   27 views

ℹ 📍  Domingos de Jesus Rodrigues, Carla Lopes Velasquez, Evandro Ferreira dos Santos, Luciane Ferreira Barbosa, Márcia Cleia   30 views
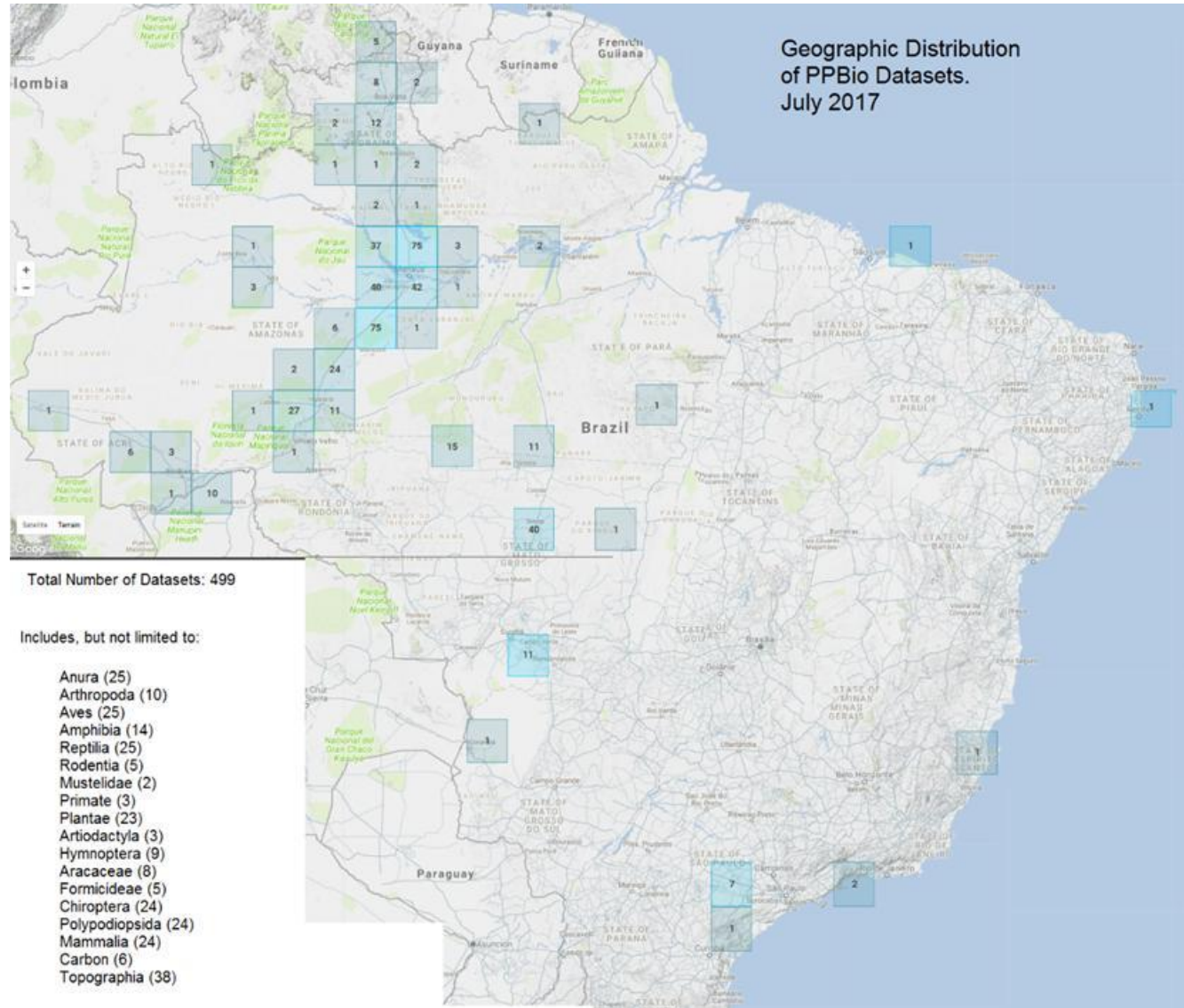
**Powered by Metacat**

# Example:

Claudia Keller, Francisco Villamarín, Rafael Bernhard, and Daniely Félix-Silva. 2016. Chelonian records from the upper Madeira River and the Madeira-Purus interfluvium, 2011- 2015.

urn:node:PPBIO. PPBioAmOc.50.30.

Geographic Distribution of PPBio Datasets.
July 2017

Total Number of Datasets: 499

Includes, but not limited to:

Anura (25)
Arthropoda (10)
Aves (25)
Amphibia (14)
Reptilia (25)
Rodentia (5)
Mustelidae (2)
Primate (3)
Plantae (23)
Artiodactyla (3)
Hymnoptera (9)
Aracaceae (8)
Formicideae (5)
Chiroptera (24)
Polypodiopsida (24)
Mammalia (24)
Carbon (6)
Topographia (38)

PPBio-AmOc is a node of the Earth Data Observation Network (DataONE).

Data Observation Network for Earth (DataONE) is the foundation of new innovative environmental science through a distributed framework and sustainable cyberinfrastructure that meets the needs of science and society for open, persistent, robust, and secure access to well-described and easily discovered Earth observational data.

DataONE will ensure the preservation, access, use and reuse of multi-scale, multi-discipline, and multi-national science data via three primary cyberinfrastucture elements and a broad education and outreach program.

DataONE comprises a distributed network of data centers, science networks or organizations. These organizations can expose their data within the DataONE network through the implementation of the DataONE Member Node service interface. In addition to scientific data, Member Nodes can provide computing resources, or services such as data replication, to the DataONE community.

Member Nodes

# DataOne GUI

# Deliverable 2.2 (D2.2) Data sharing tools

- Metacat, Morpho and DataOne were included in the February 2106 EUBON publication "Data sharing tools".

- The paper gives a comprehensive review of the existing tools for metadata, occurrence data, and ecological data.

- It includes a detailed description of the tools, their pros and cons, is followed by recommendations on their deployment and enhancement.

- Available online at

- http://www.eubon.eu/news/13351_D2.2 Data sharing tools

# Search Metacat using Solr

- Metacat uses the Solr open source search platform to organise and index the metadata and data contained in the database and to support search operations.

- This means that anyone can query a public data repository such as ours and get information according to their specific requirements.

- A Solr query begins with a base URL , in our case, this is "https://ppbiodata.inpa.gov.br/metacat/d1/mn/v2/query/solr/" which is followed by a specific set of questions, followed by the output that is required.

# Search Metacat using Solr

- So,the questions for example might ask for records that contain any kind of taxonomic information, for example ?q=class:* OR family:* OR genus:* OR kingdom:* OR order:* OR phylum:* OR scientificName:* OR species:*

- and the output takes the following format, which includes the ID number of the record, the geographic boundaries, taxomonic information related tto the record and the time that the information was collected:

- &fl=identifier,northBoundCoord,eastBoundCoord,southBoundCoord,westBoundCoord,class,family,genus,kingdom,order,phylum,scientificName,species,beginDate,

- &rows=1000 (the default value is 10) - You might need to increase this if your output file has exactly 1000 rows!

- &wt=csv indicates a request for the output to be in csv format.

# Search Metacat using Solr

- So, the complete URL for the query is:
- https:// ppbiodata.inpa.gov.br/metacat/d1/mn/v2/query/solr/?q=class:* OR family:* OR genus:* OR kingdom:* OR order:* OR phylum:* OR scientificName:* OR species:* &fl=identifier,northBoundCoord,eastBoundCoord,southBoundCoord, westBoundCoord,class,family,genus,kingdom,order,phylum,scientificN ame,species,beginDate,dateModified &rows=1000 &wt=csv
- Metacat and Tableau - Using the free application Tableau, data can be visualised in a variety of different ways.

# Tableau

# Tableau

# Certification of repositories