# Data packages
# and
# EML Data Manager

Ben Leinfelder[1]

Jing Tao[1], Duane Costa[2], Matthew B. Jones[1], Mark Servilla[2], Margaret O'Brien[3], Chad Burt[3]

[1] *National Center for Ecological Analysis and Synthesis, University of California Santa Barbara*
[2] *Long Term Ecological Research Network, University of New Mexico*
[3] *Santa Barbara Coastal LTER, University of California Santa Barbara*

Originally @ EIM 2008
September 10th, 2008

UC Santa Barbara

LTER

Wednesday, August 27, 14

# DATA PACKAGES

# Data packages

- # Data Package (eml-dataset)
  - collection of data Entities

- # Entity (eml-dataTable)
  - tabular data
  - other

- # Attribute
  - data column

- Resource map ***describes*** an aggregation
  - aggregation ***aggregates*** objects
- metadata **<--?-->** data



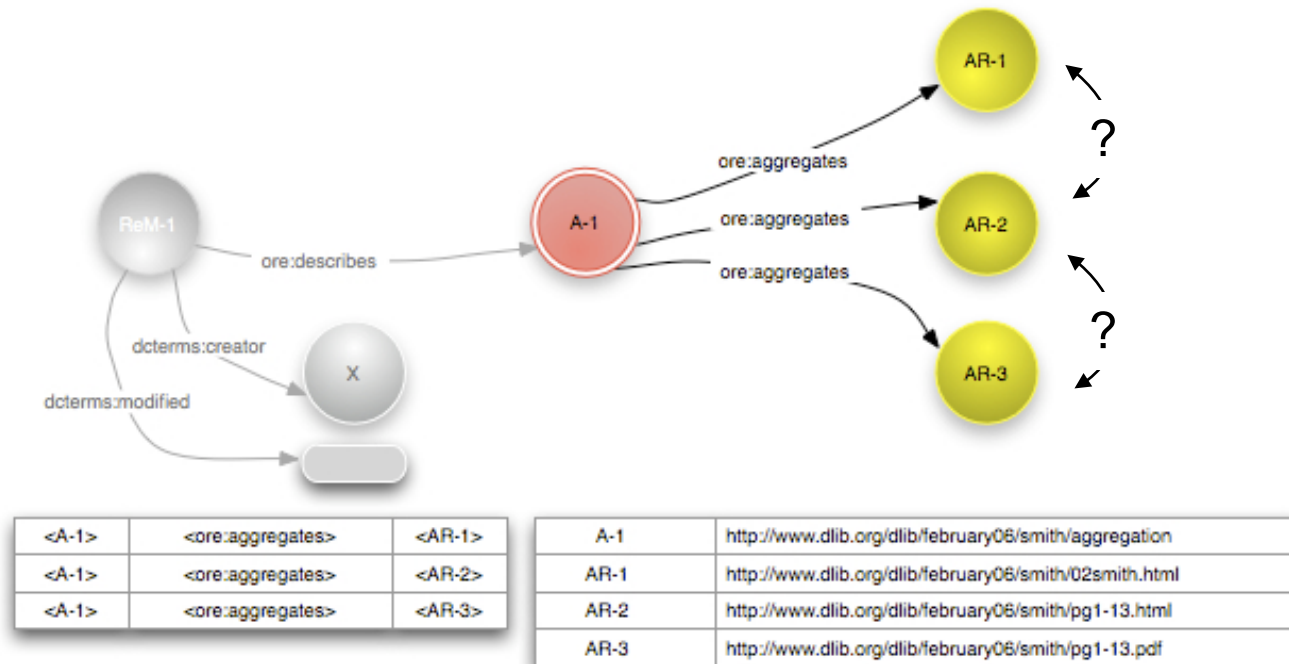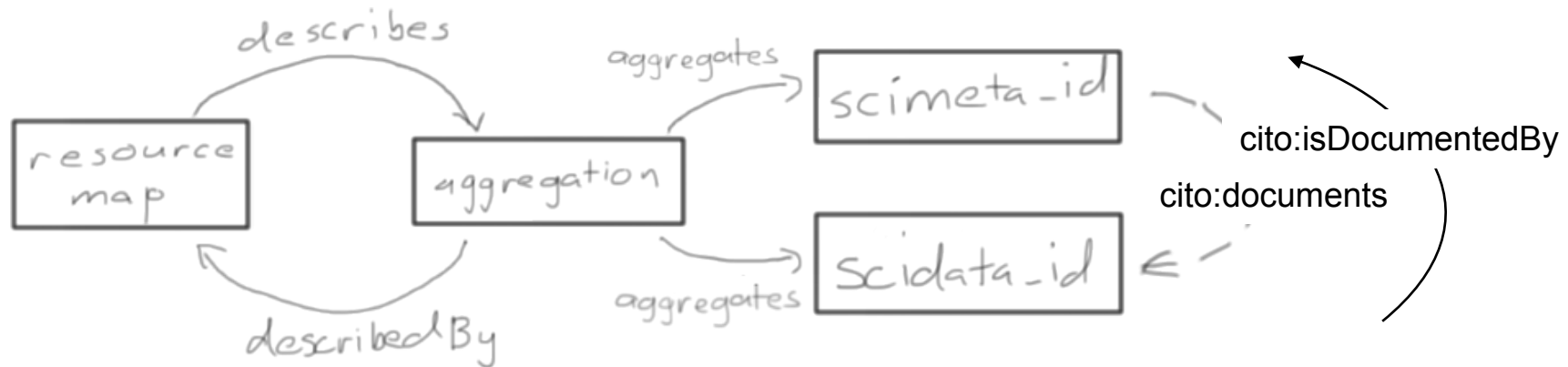| | | |
|---|---|---|
| \<A-1\> | \<ore:aggregates\> | \<AR-1\> |
| \<A-1\> | \<ore:aggregates\> | \<AR-2\> |
| \<A-1\> | \<ore:aggregates\> | \<AR-3\> |

| | |
|---|---|
| A-1 | http://www.dlib.org/dlib/february06/smith/aggregation |
| AR-1 | http://www.dlib.org/dlib/february06/smith/02smith.html |
| AR-2 | http://www.dlib.org/dlib/february06/smith/pg1-13.html |
| AR-3 | http://www.dlib.org/dlib/february06/smith/pg1-13.pdf |

- **Add relationships**
  - metadata **documents** data
  - data **isDocumentedBy** metadata



cito:isDocumentedBy

cito:documents

# EML DATA MANAGER

"Provide metadata-based query access to data with the ability to filter, join, and concatenate across data sets."
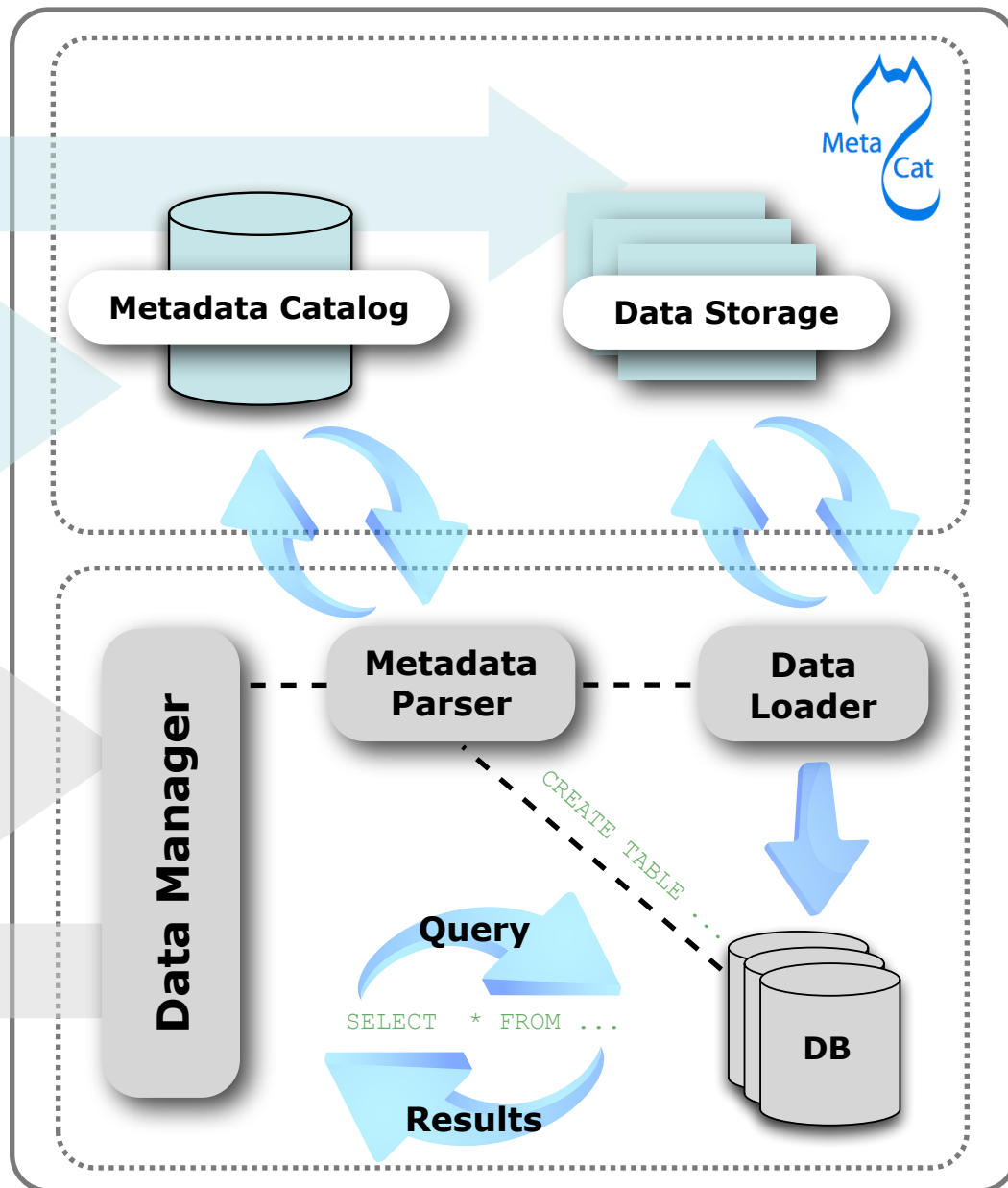
# Dynamic Data Retrieval



Store Data

Store Metadata

Metadata Catalog

Data Storage

Meta Cat

User Client

Data Query

Results

att1 | attr2 | attr3
.... | .... | ......
.... | .... | ......
.... | .... | ......
.... | .... | ......

Data Manager

Metadata Parser

Data Loader

CREATE TABLE ....

Query

SELECT * FROM ...

Results

DB

<eml>
  <dataset>
  .........
  </dataset>
</eml>

attr1 | attr2
attr1 | attr2
attr1 | attr2
.... | ....
.... | ....
.. | ....
.... | ....
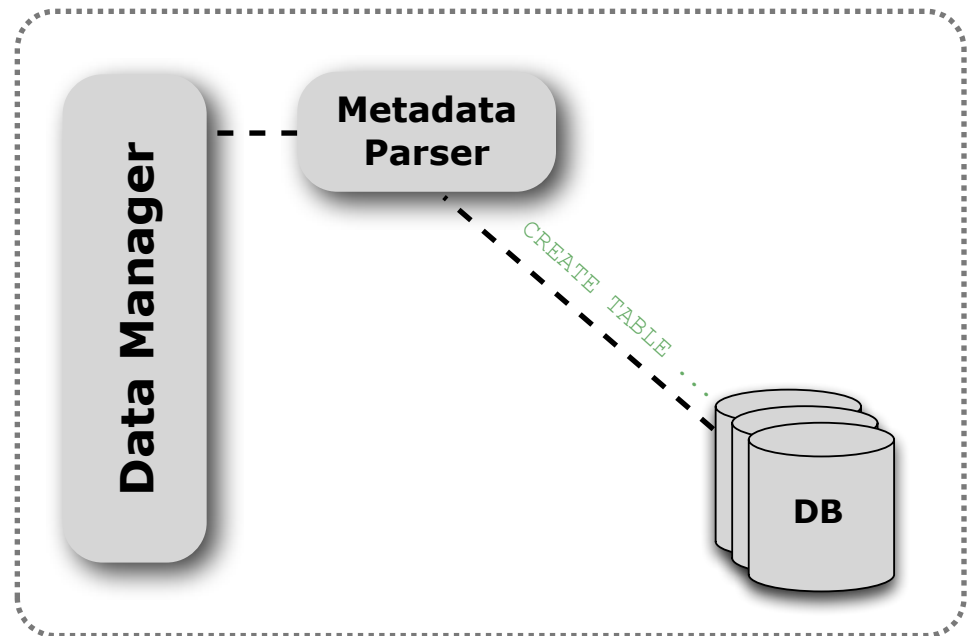.... | ....

- Uses EML metadata to:
  - Automatically create database tables
  - Download data from remote sites
  - Load data into the database
  - Manage table space via caching

- Client applications can:
  - Inspect table structures
  - Pose SQL-like queries
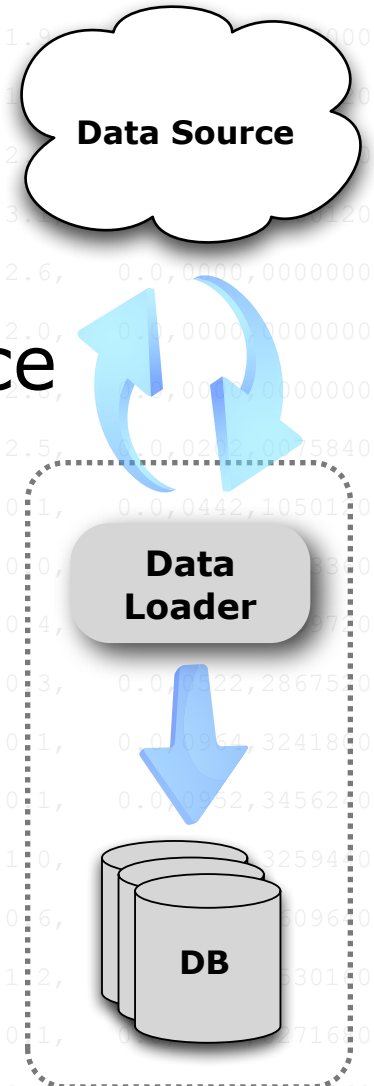  - Join and concatenate data packages

- Create underlying database table
  - Schema derived from EML

- Well-described attributes
  - type
  - precision
  - range
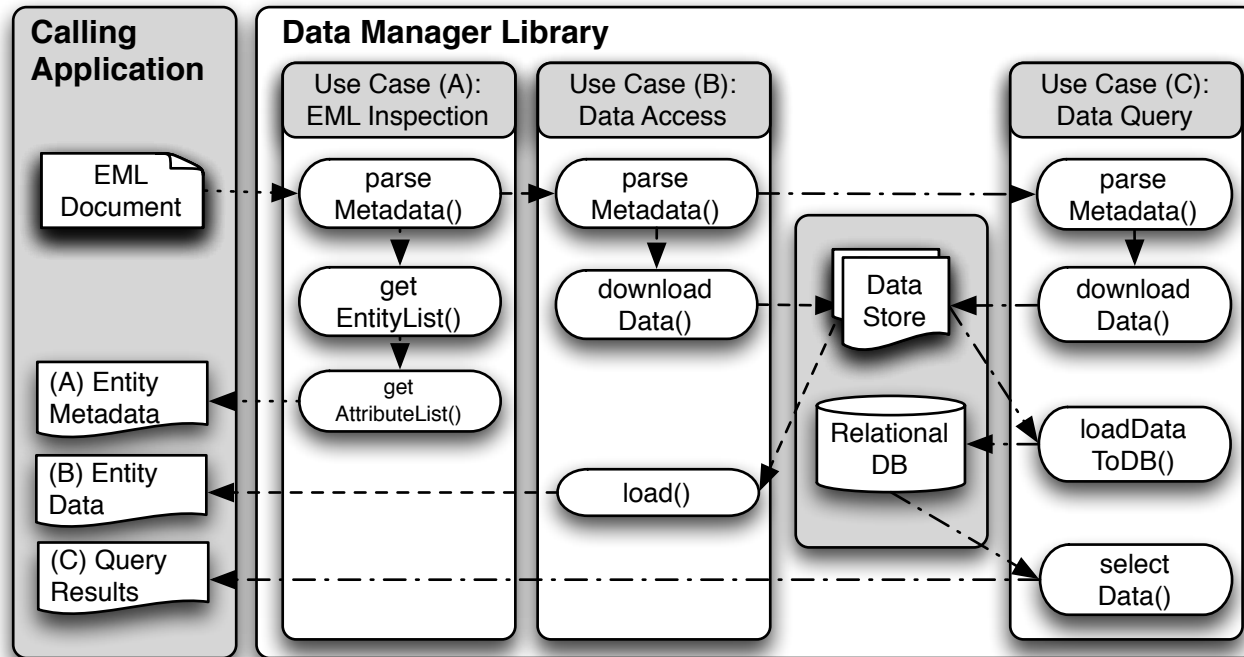  - format

# Download and Stage Data

- Retrieve data from source
  - Metacat with web services
  - Other external [accessible] source

- Insert data rows
  - Date formatting

```
Date: 09/10/2008
September 10th?
October 9th?
```

(A) Inspect metadata

(B) Download data files

(C) Query and join tabular data files

# Join Query

# Data Query Specification

- # XML syntax for describing queries

```xml
<?xml version="1.0" encoding="UTF-8"?>
<dataquery>
  <query>
    <selection>
      <datapackage id="tao.1.1">
        <entity index="0">
          <attribute index="0"/>
          <attribute index="1"/>
        </entity>
      </datapackage>
    </selection>
    <where>
      <condition type="condition">
        <left>
          <datapackage id="tao.1.1">
            <entity index="0">
              <attribute index="0"/>
            </entity>
          </datapackage>
        </left>
        <operator>=</operator>
        <right>
          <value>11/12/2008</value>
        </right>
      </condition>
    </where>
  </query>
</dataquery>
```
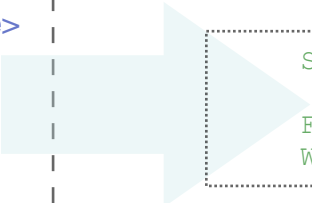
- ## Alternative to Java API
  - ✓ Attribute selection
  - ✓ Joins
  - ✓ Conditions
  - ✓ Subqueries
  - ✓ Metadata promotion

```sql
SELECT  Datos_Meteorologicos.DATE,
        Datos_Meteorologicos.TIME
FROM Datos_Meteorologicos
WHERE Datos_Meteorologicos.DATE = '11/12/2008';
```

# Data Query Specification

**Attribute(s) Info:**

| Name | Site | Year | Month | Day | Transect | Species_Code | Count |
|---|---|---|---|---|---|---|---|
| **Column Label** | | | | | | | |
| **Definition** | GCE Sampling Site | Calendar year of the observation | Calendar month of the observation | Calendar day of the observation | Transect number (randomly placed) | Coded species name | Number of individuals observed |
| **Type of Value** | integer | integer | integer | integer | integer | string | integer |
| **Measurement Type** | ordinal | datetime | datetime | datetime | nominal | nominal | ratio |
| **Measurement Domain** | Domain Info | Format YYYY / Precision 1 | Format MM / Precision 1 | Format DD / Precision 1 | Def Transect number (randomly placed) | Domain Info | Unit number / Precision 1 / Type whole / Min 0 / Max |
| **Missing Value Code** | Code NaN / Expl value not recorded or invalid | Code NaN / Expl value not recorded or invalid | Code NaN / Expl value not recorded or invalid | Code NaN / Expl value not recorded or invalid | Code NaN / Expl value not recorded or invalid | Code NaN / Expl value not recorded or invalid | Code NaN / Expl value not recorded or invalid |

```xml
<selection>
    <datapackage id="knb-lter-gce.1.9">
        <entity name="INS-GCEM-0011_1_3.TXT">
            <attribute index="0"/>
            <attribute index="1"/>
            <attribute index="5"/>
            <attribute index="6"/>
        </entity>
    </datapackage>
</selection>
```
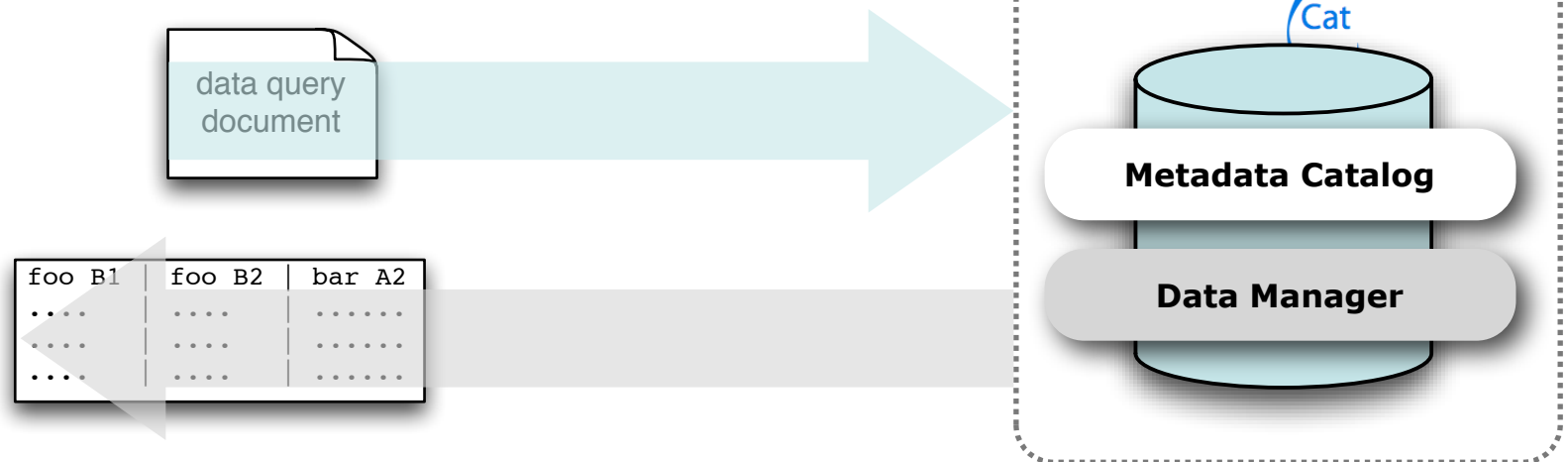
| Site | Year | Species_Code | Count |
|---|---|---|---|
| 6 | 2000-01-01 | G1 | 5 |
| 6 | 2000-01-01 | G1 | 8 |
| 6 | 2000-01-01 | G1 | 5 |
| 6 | 2000-01-01 | G1 | 8 |
| 6 | 2000-01-01 | G1 | 5 |
| 6 | 2000-01-01 | G1 | 9 |
| 6 | 2000-01-01 | G1 | 7 |

# Metacat Integration

- # Embed Data Manager
  - ## Eliminate individual client deployments

- # Expose data query capabilities
  - ## via web
  - ## additional clients

# Santa Barbara Coastal LTER



Temporal and location queries

Wednesday, August 27, 14

# LTER "PASTA"

# QUALITY CHECKS

- Verify metadata fields
  - Title (length)
  - Publication date
  - Keywords present
  - Coverage elements present

- Data files
  - Accessible/downloadable
  - Attribute metadata match data (types)

NCEAS

## • Check for optional fields in metadata

```xml
<qualityCheck qualityType="metadata" system="lter" statusType="warn" >
    <identifier>entityDescriptionPresent</identifier>
    <name>An entity description is present</name>
    <description>Check for presence of an entity description.</description>
    <expected>Field should have a data file description</expected>
    ...
  </qualityCheck>
```

```
Package ID: knb-lter-xyz.10013.1  Entity: NoneSuchBugCount
  Entity: NoneSuchBugCount  Quality Check:     entityNameLength  Status: valid
  Entity: NoneSuchBugCount  Quality Check: entityDescriptionPresent  Status: valid
  Entity: NoneSuchBugCount  Quality Check:   numHeaderLinesPresent  Status:  info
  Entity: NoneSuchBugCount  Quality Check:   numFooterLinesPresent  Status:  info
  Entity: NoneSuchBugCount  Quality Check:     fieldDelimiterValid  Status: valid
  Entity: NoneSuchBugCount  Quality Check:  recordDelimiterPresent  Status: valid
  Entity: NoneSuchBugCount  Quality Check:   attributeNamesUnique  Status: valid
```

- Verify data matches metadata

Online Distribution Info **doi:10.5063/AA/tao.2.1**

Size 188860 bytes

```
Data identifier is: tao.2.1

Entity: NoneSuchBugCount  Quality Check:    displayDownloadData  Status:  info
Entity: NoneSuchBugCount  Quality Check:        urlReturnsData  Status: valid
Entity: NoneSuchBugCount  Quality Check:            onlineURLs  Status: valid
Finished testDownloadData(), success = true

Entity: NoneSuchBugCount  Quality Check:   databaseTableCreated  Status: valid
Entity: NoneSuchBugCount  Quality Check:            onlineURLs  Status: valid
Entity: NoneSuchBugCount  Quality Check: examineRecordDelimiter  Status: valid
Entity: NoneSuchBugCount  Quality Check:   displayFirstInsertRow  Status:  info
Entity: NoneSuchBugCount  Quality Check:          tooFewFields  Status: valid
Entity: NoneSuchBugCount  Quality Check:         tooManyFields  Status: valid
Entity: NoneSuchBugCount  Quality Check:        dataLoadStatus  Status: valid
Entity: NoneSuchBugCount  Quality Check:       numberOfRecords  Status: valid
Finished testLoadDataToDB(), success = true
```

Number Of Records 100

# **INTERNATIONALIZATION**

- Important for
  - Discovery
  - Interpretation
  - Attribution

Title

Keyword

Contact information

   (names, organizations, addresses)

# EML 2.1.1 documents

- # Mark entire metadata record
  - ## – all elements inherit Portuguese
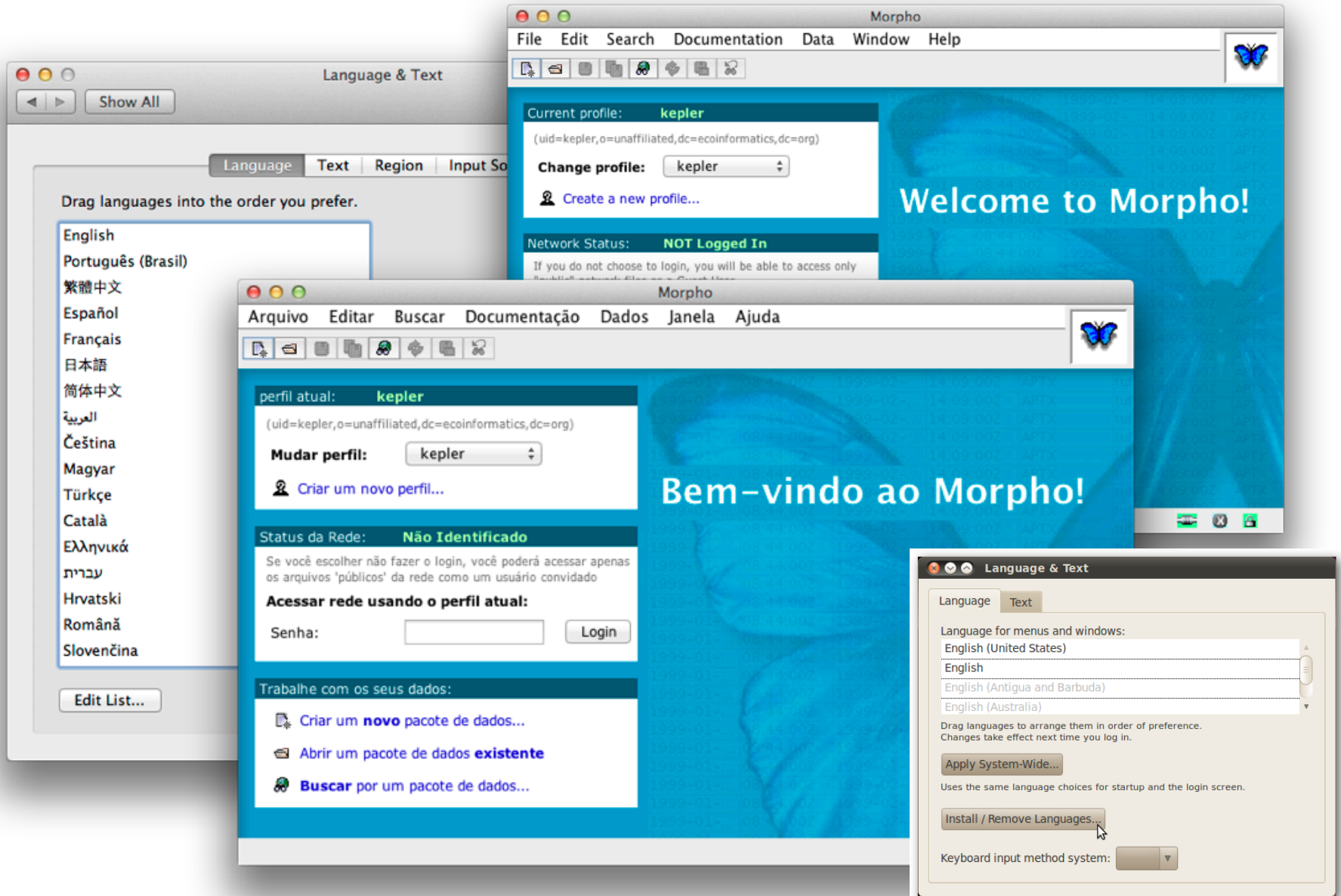
```xml
<?xml version="1.0"?>
<eml:eml
    packageId="eml.1.1" system="knb"
    xml:lang="pt_BR"
    xmlns:eml="eml://ecoinformatics.org/eml-2.1.1"
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    xsi:schemaLocation="eml://ecoinformatics.org/eml-2.1.1 eml.xsd">
```

## • Provide translations with original

```
<!-- English title with Portuguese translation -->
<title xml:lang=""en_US">
    Sample Dataset Description
    <value xml:lang="pt_BR">Exemplo Descrição Dataset</value>
</title>
```

```
<!-- Portuguese abstract with English translation -->
    <abstract>
    <para>
        Neste exemplo, a tradução em Inglês é secundário
        <value xml:lang="en_US">
            In this example, the English translation is secondary
         </value>
    <para>
    </abstract>
```

# Tool Support

# Tool Support

# Acknowledgements

- Resources
  - http://www.nceas.ucsb.edu/ecoinfo/
  - https://knb.ecoinformatics.org/
  - http://lno.lternet.edu/projects/pasta
  - http://sbc.lternet.edu/

Wednesday, August 27, 14