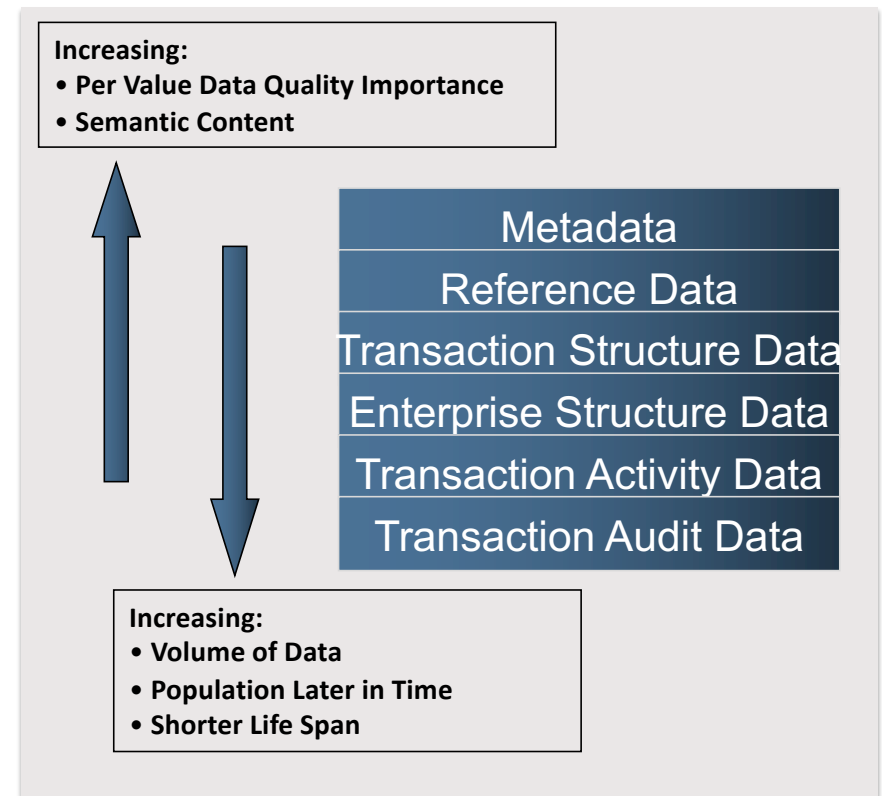# How Reference Data is Used

- Reference data is found in practically every enterprise application

- Mission-critical applications use reference data to classify and categorize master and operational data

- Data warehouses use reference data for aggregation, reporting and analytics

# Why Connect Reference Datasets

- Legacy code tables are defined on an application-by-application basis

- How do you then align, aggregate, integrate data?

- Thus, the need to connect or map disparate reference data

**Increasing:**
- Per Value Data Quality Importance
- Semantic Content

Metadata
Reference Data
Transaction Structure Data
Enterprise Structure Data
Transaction Activity Data
Transaction Audit Data

**Increasing:**
- Volume of Data
- Population Later in Time
- Shorter Life Span

*Reproduced with permission from Malcolm Chisholm*

# TopBraid EDG Crosswalks

Enterprise Data Governance (**EDG**) addresses data governance needs using heterogeneous data stores, data processing, and applications
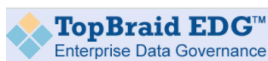
**Reference Datasets** in EDG are graphs

**Crosswalks** are special graphs that connect terms from two different vocabularies/reference datasets

# From Spreadsheet to RDF

| | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Currency | alphabetic code | numeric code | minor unit exponent | country using currency | issuing country | | | | |
| 2 | Brunei Dollar | BND | 096 | 2 | BN | BN | | | | |
| 3 | Codes specifically reserve | XTS | 963 | | | | | | | |
| 4 | Bond Markets Unit Europ | XBB | 956 | | | | | | | |
| 5 | Croatian Kuna | HRK | 191 | 2 | HR | HR | | | | |
| 6 | Taka | BDT | 050 | 2 | BD | BD | | | | |
| 7 | Yuan Renminbi | CNY | 156 | 2 | CN | CN | | | | |
| 8 | Rufiyaa | MVR | 462 | 2 | MV | MV | | | | |
| 9 | Philippine Peso | PHP | 608 | 2 | PH | PH | | | | |
| 10 | CFA Franc BCEAO | XOF | 952 | 0 | ML, GW, SN, CI, NE, BJ, TG, BF | SN | | | | |
| 11 | Bahamian Dollar | BSD | 044 | 2 | BS | BS | | | | |
| 12 | Balboa | PAB | 590 | 2 | PA | PA | | | | |
| 13 | North Korean Won | KPW | 408 | 2 | KP | KP | | | | |
| 14 | Boliviano | BOB | 068 | 2 | BO | BO | | | | |
| 15 | Euro | EUR | 978 | 2 | ES, LU, CY, TF, GF, YT, MT, AD, PT, VC, GP, LV, IT, | ES, LU, LV, IT, GR, FR, CY, SI, MC, SK, ME, EE, IE, MT, DE, FI, PT | | | | |

**TopBraid EDG™** Enterprise Data Governance ✚ ☰ ☆ ▦ **Currency Codes**   Global Lookup   Hello, 👤 **Administrator**

Codes  Dashboard  Settings  Users  **Import**  Transform  Export  Reports  Workflows  Tasks  Comments  Manage

**Import RDF File**
Adds RDF triples from a Turtle, JSON-LD or RDF/XML file.

**Import Spreadsheet using Template**
Loads a given spreadsheet file and converts its content based on a pre-defined mapping template.

**Import Spreadsheet using Pattern**
Takes a spreadsheet and converts its rows based on one out of several common spreadsheet patterns, including hierarchical patterns. Lets you map columns and save mappings as a re-usable template.

# Reference Data as a Graph

**US Dollar**
ID currency:CurrencyCode-USD

## ▼ Currency Code Info

| | |
|---|---|
| **currency name:** | US Dollar |
| **alphabetic code:** | USD |
| **numeric code:** | 840 |
| **country using currency:** | American Samoa ⌄ |
| | British Indian Ocean Territory ⌄ |
| | British Virgin Islands ⌄ |
| | Ecuador ⌄ |
| | El Salvador ⌄ |
| | Guam ⌄ |
| | Haiti ⌄ |
| | Marshall Islands (the) ⌄ |
| | Micronesia (Federated States of) ⌄ |
| | Northern Mariana Islands ⌄ |
| | Palau ⌄ |
| | Panama ⌄ |
| | Puerto Rico ⌄ |
| | Timor-Leste ⌄ |
| | Turks and Caicos Islands ⌄ |
| | United States of America ⌄ |
| | United States Virgin Islands ⌄ |
| | US Minor Is. ⌄ |
| **minor unit exponent:** | 2 |

## ▼ Issuing Authority

| | |
|---|---|
| **issuing country:** | United States of America ⌄ |

# Creating Crosswalks



Create New Crosswalk

*Crosswalks store connections between items into two different vocabularies/assets. Connections always use a predefined property: crosswalk:closeMatch.*

This creates a new Crosswalk with yourself as the manager.

**Label:** FIPS Country Codes to ISO Country Codes

**Description:**

**From Graph:** FIPS Country Codes

**From Entity Type:** FIPS Country

**To Graph:** Country Codes

**To Entity Type:** ISO Country

☑ Also generate property shapes for the match predicate

- Select source and target data
- Specify source and target classes.  Only instances of these classes will be candidates for matching

# Generating matches

# Accepting Matches



- Can be automatic or Can require user input
  - Filter by confidence / "Accept all" option
  - Matches can be entered / changed manually

# Many to Many Mappings

# Determining Candidate Matches

- Based on the label

  - rdfs:label, skos:prefLabel, skos:altLabel

  - Users can specify custom labels

  - Alternative labels (skos:altLabel) are given lower confidence

- Uses a Lucene index

  - Create an index holding labels of all instances of target class

- Query the index for matches for each instance of the source class

  - By default we do a fuzzy query on the labels

- Customers can customise/extend matching behaviour

# Accessing Crosswalk Information in a Reference Dataset

# How Well Does Automated Mapping Perform in Real Life Scenarios?

**TopQuadrant™**

- For a pharma customer, we analyzed performance of automated mapping from a customer specific terminology used for clinical trials to SNOMED and MESH. Calculating:
  - % found mappings from the local terminology to SNOMED and MESH i.e., recall
  - % of correct matches vs Expected results i.e., precision
  - Mappings for 2 columns

## Correct match % Calculation

- We considered a correct match one that presents expected mapping, regardless of the confidence.
  - For example, if we returned 3 mappings from the same input value to different concepts (at different confidence levels), as long as one of them matches the expected term, the match is considered *correct*.

# How Well Does Automated Mapping Perform in Real Life Scenarios?

## Matches Found

| Local Terminology | Standard Vocabulary | Matches found (%) |
|---|---|---|
| **Col1 - 105 unique cell values** | MESH | 102 out of 105 (97%) |
| **Col2 - 150 unique cell values** | SNOMED | 150 out of 150 (100%) |
| **Col2 - 150 unique cell values** | MESH | 150 out of 150 (100%) |

## Correct matches (expected results)

| Local Terminology | Standard Vocabulary | Correct Matches found (%) |
|---|---|---|
| **Col1** | MESH - 132 expected matches | 116 out of 132 (88%) |
| **Col2** | SNOMED - 69 expected matches | 62 out of 69 (90%) |
| **Col2** | MESH - 147 expected matches | 116 out of 147 (79%) |

# Accessing Crosswalk Information in a Reference Dataset

# Querying and Exporting via SPARQL

# Querying and Exporting via GraphQL

# Crosswalk Services

```
1 ▾ query getSearchResults {
2 ▾   results: fipscountries(where: {closeMatch: {exists: {iso3166_3AlphaCode: {hasValue: "AUS"}}}}) {
3        uri
4        label
5        fipsCode
6      }
7    }
8
```

```json
{
  "data": {
    "results": [
      {
        "uri": "http://topbraid.org/data/country/FIPSCode-AS",
        "label": "Australia",
        "fipsCode": "AS"
      },
      {
        "uri": "http://topbraid.org/data/country/FIPSCode-CR",
        "label": "Coral Sea Islands",
        "fipsCode": "CR"
      }
    ]
  }
}
```

**What is mapped to a reference data item with the ISO 3-character code = AUS?**

# Thank You!

- Irene Polikoff

  irene at topquadrant.com

- TopQuadrant

  https://www.topquadrant.com/

  **Request demo:** sales at topquadrant.com

  **Request evaluation account:**
  https://www.topquadrant.com/products/topbraid-enterprise-data-governance/request-edg-evaluation-account/

- Enterprise Data Governance (EDG)

  https://www.topquadrant.com/products/topbraid-enterprise-data-governance/