# Ontology and Taxonomy: Strange Bedfellows

## by Michael Uschold
## Semantic Arts
www.semanticarts.com

# The Situation

- Knowledge assets in large enterprises are very complex
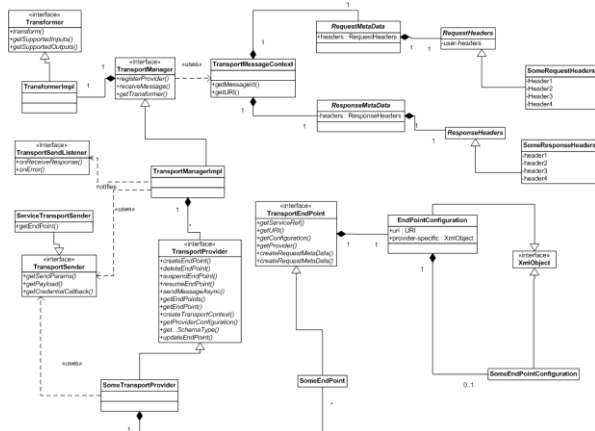- It got that way for many reasons

Mega Corp

# Roots of Complexity

- Ambiguity is pervasive

- Systems are developed independently

- One database for each application

- Lots of metadata but:
  - No reuse of data models
  - Heterogeneity reigns supreme

- Lets look at Mega Corp

# Themes at Mega Corp

- People would lament the growing complexity of their information systems

- But their focus was on short term results

- They realized they needed some good models, ideally *one model to rule them all*…

- But they kept acquiring more companies

- "Let's not re-invent the wheel" led to more models (and more wheels)

# Many Modeling Structures

*Glossaries / Controlled Vocabularies*

*Data and Document Metamodels*

ad hoc
Hierarchies
(DMOZ)

structured
Glossaries

XML
Schema

Restricted
Logics
(OWL, Flogic)

formal
Taxonomies

Terms

Thesauri

XML DTDs

'ordinary'
Glossaries

Principled,
informal
taxonomies

Data Models
(UML, STEP)

Data
Dictionaries

Frames

General
Logic

DB
Schema

*Informal Taxonomies and Thesauri*

Formal Knowledge Bases & Inference

# Many Modeling Structures

*Glossaries / Controlled Vocabularies*

*Data and Document Metamodels*

ad hoc
Hierarchies
(DMOZ)

structured
Glossaries

XML
Schema

Restricted
Logics
(OWL, Flogic)

formal
Taxonomies

Terms

Thesauri

XML DTDs

'ordinary'
Glossaries

Principled,
informal
taxonomies

Data Models
(UML, STEP)

Data
Dictionaries

DB
Schema

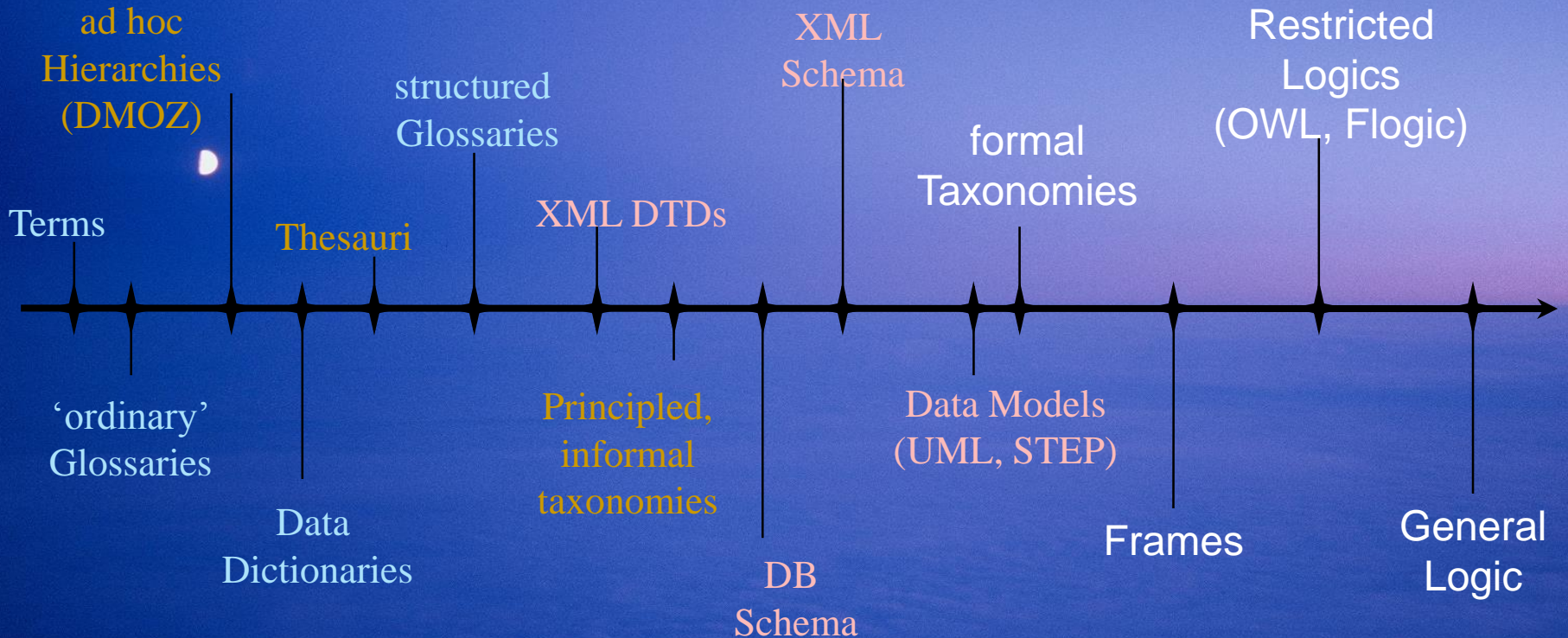Frames

General
Logic

*Informal Taxonomies and Thesauri*

Formal Knowledge Bases & Inference

Many Tools and Approaches: Informal to Formal

# Tools and Approaches

- **There are many tools**
  - Spreadsheets, Spreadsheets & more Spreadsheets
  - Vocabulary managers
  - Indexing and search
  - XML editors
  - ER modeling
  - Taxonomy and Ontology tools
- **Reuse and sharing is next to impossible**
- **These tools and approaches usually mix like oil & water**

# Oil & Water

- Different reasons for organizing knowledge
- Different cultures: both technological & social
- Different levels of formality (neats/scruffies)
- Conceptual vs. Design vs. Implementation
- Governance: who gets to control what?
- The menu vs. the meal

# The Menu vs. the Meal

## Taxonomy and Thesauri:

- focus is on words not concepts (the menu)
- relationships are between terms: synonym, hyponym, broader/narrower term
- each term <u>should</u> refer to just one concept

Don't eat the menu...

## Ontology:

- focus is on concepts (the meal)
- relationships are between concepts
- formal definitions
- automated inference

Eat the meal

# Holy Grail: Bring It All Together

- Understand where each approach adds the most value

- Find the touch points and link them all up

- Can everyone and every tool can live in harmony?

- An impossible dream?

- We are pushing hard on this
  It's getting a lot less impossible

- Lets look at a Case Study at Mega Corp

# Case Study

- A certain kind of thing needs to be managed

- There are millions of them, and 1000s of new ones arrive every day

- They wanted to track these items and see which ones were having what impact where (purpose A).

- So they created many hundreds of "buckets" which they would use to classify the items.

- It must be possible to classify every item into exactly one bucket.

# Case Study

- The set of buckets were defined and enshrined in a spreadsheet where each row represented a bucket and there were two main columns:
  - Name of the bucket
  - Text description of the bucket
- And they saw that it was good (for their purpose)

# Others Noticed

- There were some other groups that managed similar items.

- They went ahead and tried to use those same buckets for their different purposes.

Purpose A

Purpose B

Purpose C

Single-purpose
Knowledge Asset

# Reuse not so Easy: Why?

- Despite strong similarity in the underlying items for all groups, there were large differences in how they managed the items

- This was maddening:
  Similarity so near, yet Reuse so far away

- Much head-scratching ensued

- How to get to the bottom of this?

# Lets Walk before we Run

- While it is important that the asset become reusable across different groups.

- We wanted to first look carefully at that first asset purely in the context of its original purpose.

- Then we will get back to reuse.

# Gather Background Information

- ## Subject and Scope of what is being organized

  - what is known to be out of scope?

  - criteria for deciding in/out of scope?

- ## Intended audience, purpose & current uses

- ## Notation: syntax and semantics

- ## Provenance

# Evaluate along Specific Criteria

Critera:

- Clarity & Focus

- Scope coverage

- Right level of detail

- Categorization rigor

- Consistency and Uniformity

- Rate on a scale: 1-5

- Review with client

# Some Problems with Initial Asset

- It turned out that the original set of buckets was hard to use and manage…

- .. even for the original intended purpose

- If a change was needed, it was hard to see which of the hundreds of buckets would be affected.

- Impact analysis was next to impossible

- But the core asset was very important

# What's Going On?

- The buckets were mainly text descriptions.

- No structure means no automation,
  every bucket had to be examined manually.

- Yet, there *was* evidence of much structure
  lurking behind the text.

- It was plain to see, if you were paying attention
  and knew what to look for

# Finding Structure

- Look closely at the text descriptions and look for patterns.

- There were some recurring themes.

- Frequent mention of Goal, Region and Product, but these ideas are not captured or used in a uniform manner.

# Making Implicit Structure Explicit

- Try to reword descriptions for the buckets in a uniform way.  For example:

  – Work on iPhones in Africa to reduce service call wait times.

  – Work on product P in region R to achieve goal G

- The manual rewording was not always so easy.

- Not all ways to capture structure are equal

- Let's consider an old and familiar structure…

# Dewey Decimal: Geography

**G**

| | |
|---|---|
| Graphic novel genres | 741.53 |
| Graphic novelists | |
|     biography | 741.593–.599 |
|     *see Manual at 741.593–.599* | |
| Graphic novels | 741.5 |
|     Argentina | 741.598 2 |
|     Belgium | 741.594 93 |
|     England | 741.594 2 |
|     France | 741.594 4 |
|     geographic treatment | 741.593–.599 |
|       *see Manual at 741.593–.599* | |

Geography will turn up in many different places.

# Dewey Decimal: Geography

- Use in many places

- Manage in one place

| GeographicRegion | Argentina | Belgium | France |

| | Novel |
| --- | --- |
| | HistoricalNovel |
| | Argentina |
| | Belgium |
| | Mystery Novel |
| | Graphic Novel |
| | ScienceFiction |
| | Argentina |
| | France |

| | Novel |
| --- | --- |
| | Historical Novel |
| | Mystery Novel |
| | Graphic Novel |
| | Science Fiction |

- These repeating ideas are called "facets"

- E.g. faceted search

# Faceted Search: www.newegg.com

# A Faceted Taxonomy for Laptops

The item being classified

Sample category or "bucket":
*Under 5 pounds, Battery life more than 11 hours, 1600 x 900 resolution, i7 CPU*

Each bucket is highly structured

Facet

Possible values

Laptop

hasFacet → CPU → Core i3 / Core i5 / **Core i7**

hasFacet → Screen Resolution → 1376 x 768 / **1600 x 900** / 1910 x 1080

hasFacet → Weight → < 3 lb / < 4 lb / **< 5 lb**

hasFacet → Battery Life → > 5 hours / > 7 hours / > 9 hours / **> 11hours**

# Back to our Structured English

- We have Goal, Region and Product

- They are candidate "facets" for characterizing the items in question.

- But there were about a dozen other potential facets that we saw in the text descriptions

- Which ones really mattered?

- Which ones are just incidental?

- Can facets really help?

# Facet Math

Without facets: there are exponentially many buckets:
3 x 3 x 3 x 4 = 108

108 things to learn and remember is a lot.

The faceted approach means there are:
- 4 facets + 13 values
- = 17 things to learn and remember

# This is a Big Deal

Exponentially reducing the number of things to learn to classify things has numerous benefits

- Faster to train people

- More accurate classification

- Easier to evolve and maintain moving forward.

- The more facets & values, the greater the savings

- But how *do* we know we have the right facets?

# Many Meetings with Stakeholders

- Get experts about the items in question.

- Ask them to identify the ways that items are different from one another

- Brainstorm to identify candidate facets

- Then evaluate them

- Example Criteria:  Ideally each facet value should be unique for a given facet.

# Uniqueness: An Example

Suppose we are classifying quality control activities. One facet is the goal. Values might be:

- timeliness

- completeness

- accuracy

- timeliness

- completeness only

- accuracy and possibly completeness

What happens if some control actions are for both completeness and accuracy?

Then it is hard to uniquely classify the item.

# A Faceted Taxonomy for Control Activity



The item being classified

Facet

Possible values

| Amazon |
| Best Buy |
| NewEgg |
| ... |
| ... |

Customer

| Africa | → North Africa, East Africa, West Africa |
| Europe |
| Asia |
| Americas | → North America, Central America |
| AustrailAsia |

Region

Control Activity

hasFacet

| Compliance |
| Timeliness |
| Completeness Only |
| Accuracy and perhaps Completness |

Goal

| Smart Phone |
| Tablet |
| Laptop |

Product

# Space of Possible Values for a Facet

The values for the facets/properties may be:

- <u>Flat</u>: a flat list of a handful of possible values (*e.g. Amazon, Best Buy, New Egg*)

- <u>Hierarchical</u>: a simple taxonomy (*e.g. geographic regions*)

- Can anyone think of another possibility?

- What about Laptops?

# Space of Possible Values for a Facet

The values for the facets/properties may be:

- <u>Flat</u>: a flat list of a handful of possible values (*e.g. Amazon, Best Buy, New Egg*)

- <u>Hierarchical</u>: a simple taxonomy (*e.g. geographic regions*)

- <u>Faceted</u>: another faceted taxonomy embedded in the prior faceted taxonomy (*e.g. products*)

# A Faceted Taxonomy for Control Activity

# A Faceted Taxonomy for Controls

# A Faceted Taxonomy for Controls



- *We decomposed the original asset*
- *Next: re-compose it from the pieces*

# Re-Characterizing the Items

- For each of hundreds or thousands of item descriptions, re-characterize them using the facets.

- Many ways to do this:
  - Manually reword them one by one

  - Use a spreadsheet to create a form
    - One field in the form for each facet
    - Values may be selected from a dropdown or entered into a text field

  - Build a simple app that automates the form
    - Create a simple ontology
    - Use it to drive the form
    - The taxonomy becomes a set of triples that can be queried

# Decomposing & Re-characterizing

Purpose A

Re-Composed for Purpose A

Single-purpose Knowledge Asset

- Exposed the hidden structure.
- Much easier to use and evolve.
- BUT: still a single purpose asset

Re-characterize

Decompose

Reusable Parts (facets)

# What about other uses?

# What about other uses?

Purpose A

Purpose B

Purpose C

Re-Composed for Purpose A
(still single purpose)

Composed for Purpose B

Composed for Purpose C

Reusable set of building blocks

Reusable Parts (facets)

Compose

Compose

# More than Just a Story

- In our early work at Mega, this was just a story, a nice idea we hoped would come true.

- Several months later, we were back at Mega and asked them how things were going.

- They are doing just what the picture depicts

- Taking the facets and applying them to classify the items for their own purposes

- But wait, there's more!

- What about an ontology?

# Linking Taxonomies to an Ontology

- Normally, a taxonomy of terms, or a faceted taxonomy would live independently from an ontology.

- Our vision is to have every thing connected.
  - spreadsheets with a semantic underpinning
  - multiple applications & databases
  - data models and messages

- Opens up vast possibilities for querying and analyzing data across an enterprise

# Enterprise-Wide Ontology

- We were also building an enterprise ontology for a major part of their business.

- They are now linking the facets to the ontology so that faceted taxonomies are living in harmony with formal ontologies.

- All the way from text definitions to informal taxonomies to faceted taxonomies to ontologies, everything linked together.

# Triple as Common Denominator



XML

Databases

Spreadsheets

Free text

:predicate

:subject

:object

# Linking Taxonomies to an Ontology

# An Ontology Perspective

- We have been talking from a taxonomy perspective, and then linking to an ontology.

- The reverse is when we are building an enterprise ontology and we want to identify where there are potential taxonomies lurking.

# Example: *Corporations and Charities*



WA SOS Corporations and Charities Semantic Model

# Example: Corporations and Charities



WA SOS Corporations and Charities Semantic Model Summary
November 3, 2013

# Governance

- Taxonomies can be independently governed
- When changes occur, the touch points are limited so there is minimal disruption

# One More Example: Codes

- A key application at Mega has over 24,000 codes grouped into over 700 code categories.

- There are only two or three people in the company who understand them

- Very time-consuming to learn,
  a risk if people are no longer around

- Impossible to do any serious analytics

# One More Example: Codes

- We uploaded the codes into a triple store and explored using SPARQL queries .

- This dramatically reduces learning curve, eases risk and burden on the few experts.

- We also found some errors

- Lucky that no one used those fields (or maybe they did, and no one noticed!)

# The Learning

- A little bit of semantics can make a big difference in a surprising way

- Codes are notoriously difficult to understand (hence the name?)

- But they really *do* mean something, and we are starting the process of giving them meaning by linking them to the enterprise ontology and the taxonomies.

- In the long term, *every* one of those codes could turn into a facet value in a taxonomy.

- The vision for ontology and taxonomy to live in harmony is unfolding in another division in Mega

# Summary Themes

- Knowledge assets are often a mess

- Hard to use, reuse, maintain and evolve

- Decompose into the essential components and using facet analysis

- Re-characterize the original asset so it is easier to use, maintain and evolve.

- Use the essential components as common building blocks to purpose-build other assets for other uses.

# Summary Themes

- The common building blocks are linked to and/or become part of an enterprise wide ontology.

- The Enterprise Ontology has many uses:
  - reaching a common understanding
  - basis for semantic integration of heterogeneous knowledge and data assets (including countless spreadsheets)
  - supports automated inference for consistency, completeness and enhanced analytics

# Summary Themes

- Most taxonomy work is about search and navigation

- We broadened it to help manage knowledge assets more generally, whatever their purpose.

- Improve understandability, use and reuse

# Ontology vs. Taxonomy

- Most ontologists are not very interested in taxonomy

- Many traditional taxonomists don't understand ontology

- We are applying ontological analysis to design better taxonomies

- We find that both are critically important in the modern enterprise.

- Thus we have Strange Bedfellows…

# One Happy Family of Models

- Ontology: best for core classes and properties
- Taxonomies, often faceted
  - for fine scale distinctions on the edges
  - to be governed by separate and sometimes external parties
- Data models and messages derived from the ontology, using fine grained distinctions from the taxonomies as needed.
- When Mega started doing this in their Enterprise…

# *Something Magic Happened*

- Rather than the 1,000,000 concepts Mega had baked into all the schema of all their current systems

- Or the 100,000 elements they had captured in a metadata repository (so far)

- Or the 200,000 taxonomic distinctions they had either collected or subscribed to

- Or the 50,000 attributes they had in their fully attributed Entity Data Model

- Or the 20,000 elements they had in the sum total of all the messages in their SOA

# It turned out…

- There were less that 1000 concepts that they ran their whole business on

- And of these 1000 there were 70 classes and 30 properties that shaped all other information

- Anyone who was a bit motivated could find they concepts they needed in this new simplified knowledge-scape

# Semantic Computing Writ Large

Our focus today: taxonomy, ontology, the semantic web.

Complementary technology:

- Machine learning
- NLP
- Big Data
- The Cloud

Early in 2014, Gartner called out the following as major trends: Cloud, Big Data and Semantics

Big Semantic Cloud

# Thank You

Our website: www.semanticarts.com

- We do consulting and training, specializing in helping large companies find their information core

- Leave me a card for a copy of this presentation

- Questions?