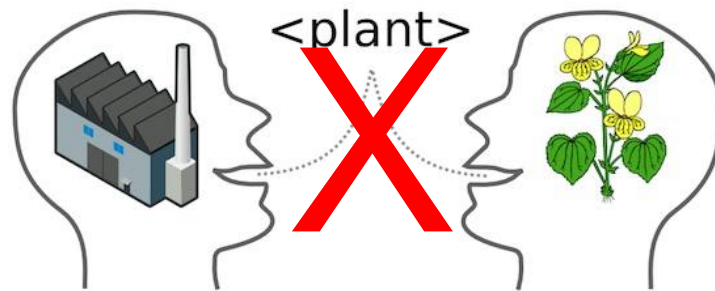


# Metadata Interoperability



# Agenda

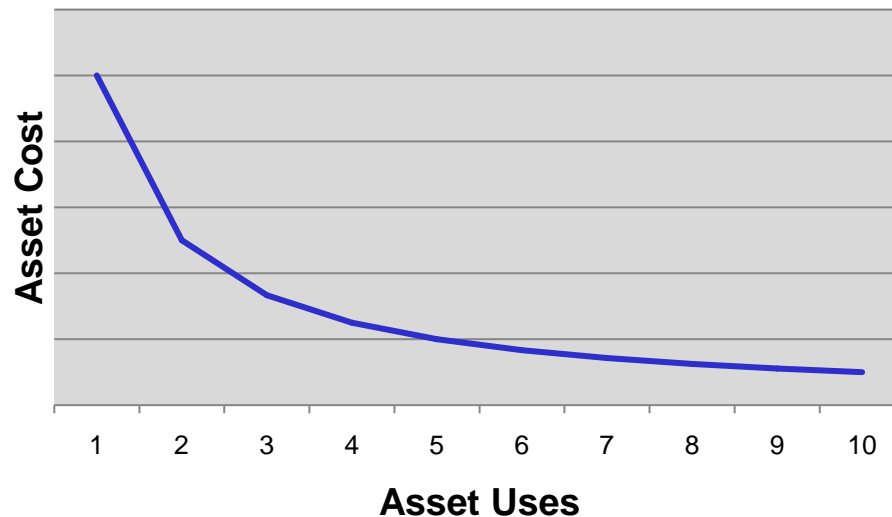
- ❖ Interoperability
- ❖ Dublin Core
- ❖ Dates, roles and topics
- ❖ Tagging assets
- ❖ Tools for automating tagging

# Interoperability

- ❖ **The ability of diverse systems and organizations to work together by exchanging information.**
- ❖ Semantic interoperability is the ability to automatically interpret the information exchanged meaningfully and accurately.

# Interoperability ROI

- ❖ Content assets are expensive to create so it's critical that they can be found, so they can be used and re-used.
- ❖ Every re-use decreases the content asset creation cost and increases the asset value.



## Interoperability (2)

- ❖ If content assets are so important, why can't they be found?
  - They contain no searchable text, e.g., images.
  - They exist in different applications, file shares and/or desktops.
  - ... Other reasons?
- ❖ When they are found why can't content assets be reused?
  - When there are multiple versions, it's difficult to choose which one to use.
  - The usage rights may not be clear.
  - ... Other reasons?

# Interoperability (3)

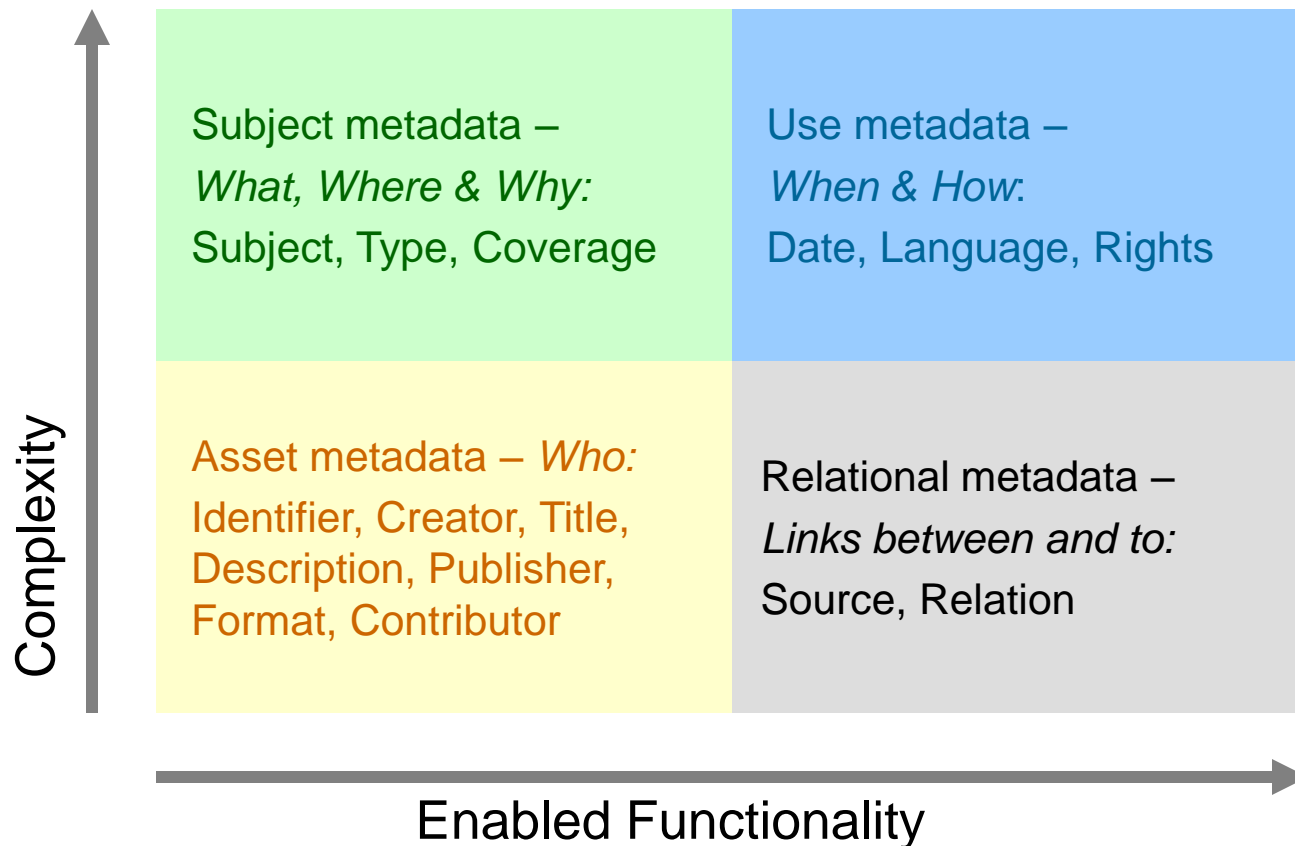
- ❖ Digital content assets are sourced from multiple applications and locations
  - Desktop applications
  - Enterprise content management applications
  - Product lifecycle management (PLM) applications
  - Product information management (PIM) applications
  - Third party contractors' systems
  - Other divisions and departments in your organization
  - Marketing and Communications servers
  - ...Other applications and locations?

# Interoperability vision

- ❖ I want to easily find any content assets in a particular format that can be used for a specific purpose regardless of where they are located.
- ❖ Challenges:
  - How to align different metadata properties
    - E.g., Author and Creator; Title and Caption; Location and Setting; etc.
  - How to align different vocabularies
    - E.g., CA and California; RiM and Research in Motion; etc.

# What is metadata?

- ❖ Metadata provides enough information for any user, tool, or program to find and use any piece of content.

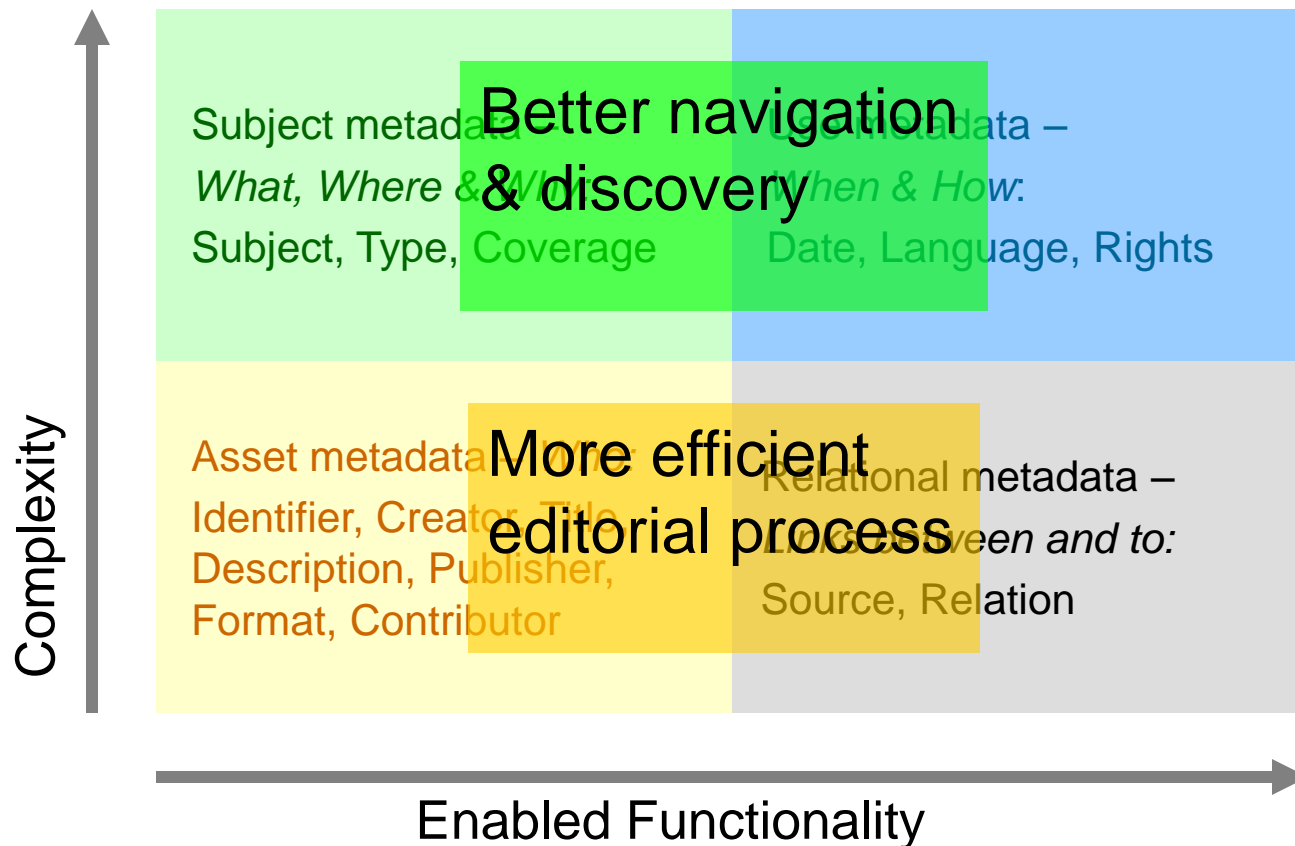


<http://dublincore.org/documents/dces/>



# What is metadata

- ❖ Metadata provides enough information for any user, tool, or program to find and use any piece of content.

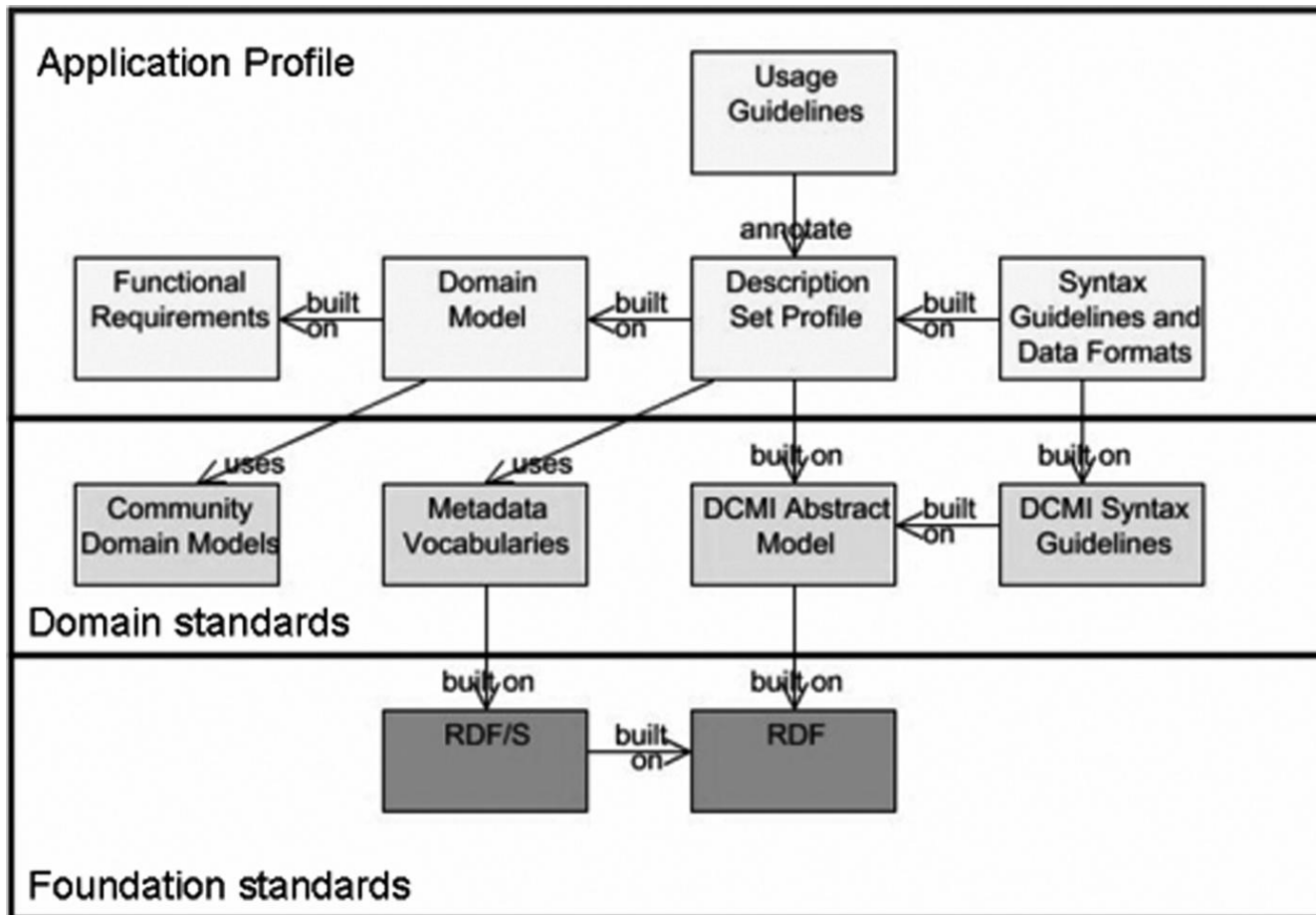


<http://dublincore.org/documents/dces/>

# But Dublin Core is a little more complicated

Elements	Refinements		Encodings	Types
1. Identifier	Abstract	Is referenced by	Box	Collection
2. Title	Access rights	Is replaced by	DCMIType	Dataset
3. Creator	Alternative	Is required by	DDC	Event
4. Contributor	<i>Audience</i>	Issued	IMT	Image
5. Publisher	Available	Is version of	ISO3166	Interactive
6. Subject	Bibliographic citation	License	ISO639-2	Resource
7. Description	Conforms to	Mediator	LCC	Moving Image
8. Coverage	Created	Medium	LCSH	Physical Object
9. Format	Date accepted	Modified	MESH	Service
10. Type	Date copyrighted	<i>Provenance</i>	Period	Software
11. Date	Date submitted	References	Point	Sound
12. Relation	Education level	Replaces	RFC1766	Still Image
13. Source	Extent	Requires	RFC3066	Text
14. Rights	Has format	<i>Rights holder</i>	TGN	
15. Language	Has part	Spatial	UDC	
	Has version	Table of contents	URI	
	Is format of	Temporal	W3CTDF	
	Is part of	Valid		

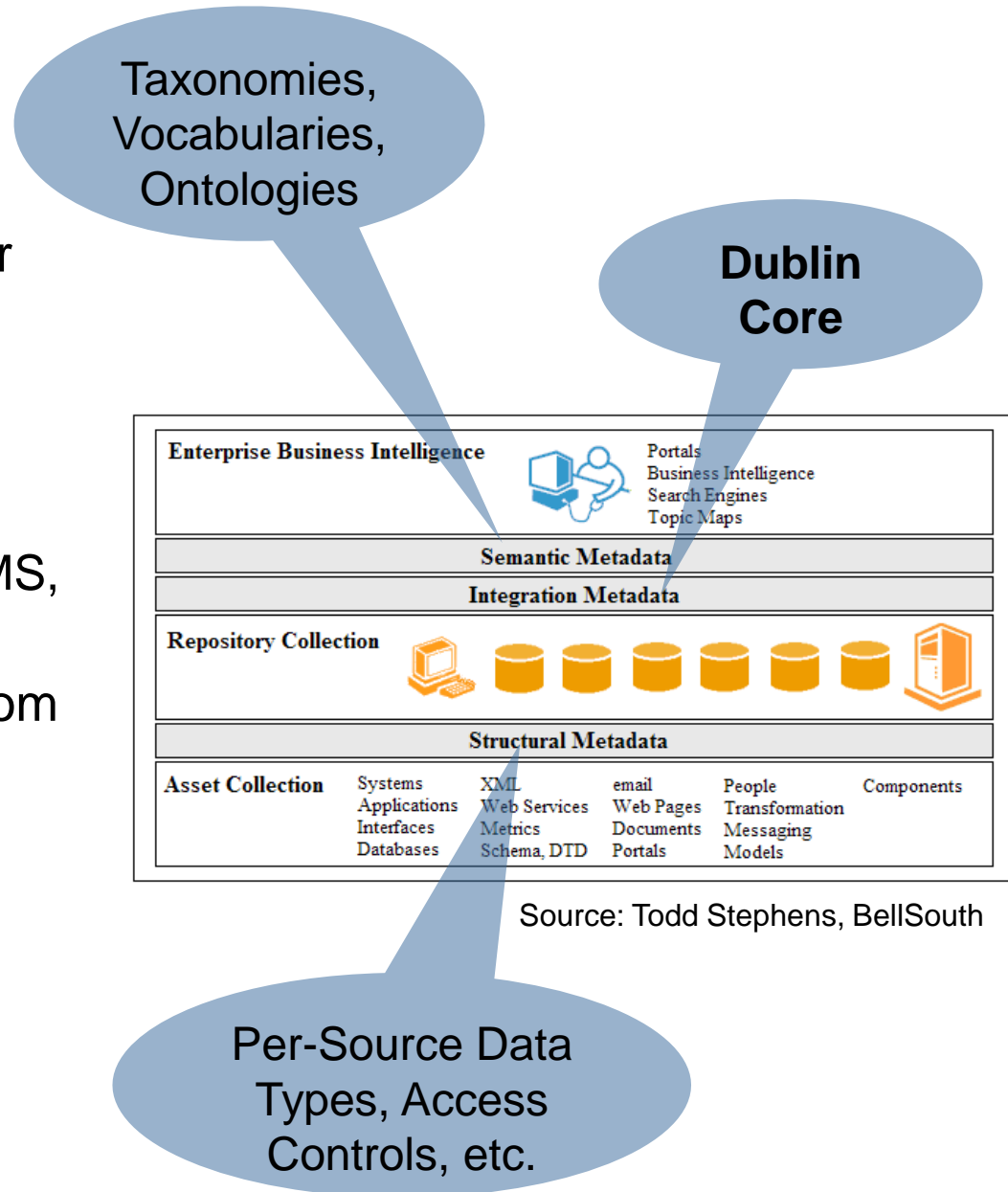
# ... and Dublin Core has gotten more abstract



DCAM (Dublin Core Abstract Model) Singapore Framework

# So why Dublin Core?

- ❖ Dublin Core is a de-facto standard across many other systems and standards
  - RSS (1.0), OAI (Open Archives Initiative)
  - Inside organizations – SharePoint, ECMS, DAMS, etc.
- ❖ Mapping to DC elements from most existing schemes is simple
- ❖ Metadata already exists in enterprise applications
  - MS Office, SharePoint, MarkLogic, OpenText, Documentum, SAP, etc.



# Dublin Core is the top vocabulary in the linked data cloud



Vocabulary prefix	Vocabulary link	Number of usages in data sets	Data sets that use the vocabulary
dc	<a href="http://purl.org/dc/elements/1.1/">http://purl.org/dc/elements/1.1/</a>	92 (31.19 %)	Data sets that use dc
foaf	<a href="http://xmlns.com/foaf/0.1/">http://xmlns.com/foaf/0.1/</a>	81 (27.46 %)	Data sets that use foaf
skos	<a href="http://www.w3.org/2004/02/skos/core#">http://www.w3.org/2004/02/skos/core#</a>	58 (19.66 %)	Data sets that use skos
geo	<a href="http://www.w3.org/2003/01/geo/wgs84_pos#">http://www.w3.org/2003/01/geo/wgs84_pos#</a>	25 (8.47 %)	Data sets that use geo
xhtml	<a href="http://www.w3.org/1999/xhtml/vocab#">http://www.w3.org/1999/xhtml/vocab#</a>	19 (6.44 %)	Data sets that use xhtml
akt	<a href="http://www.aktors.org/ontology/portal#">http://www.aktors.org/ontology/portal#</a>	17 (5.76 %)	Data sets that use akt
bibo	<a href="http://purl.org/ontology/bibo/">http://purl.org/ontology/bibo/</a>	14 (4.75 %)	Data sets that use bibo
mo	<a href="http://purl.org/ontology/mo/">http://purl.org/ontology/mo/</a>	13 (4.41 %)	Data sets that use mo
vcard	<a href="http://www.w3.org/2006/vcard/ns#">http://www.w3.org/2006/vcard/ns#</a>	10 (3.39 %)	Data sets that use vcard
sioc	<a href="http://rdfs.org/sioc/ns#">http://rdfs.org/sioc/ns#</a>	10 (3.39 %)	Data sets that use sioc
cc	<a href="http://creativecommons.org/ns#">http://creativecommons.org/ns#</a>	8 (2.71 %)	Data sets that use cc
geonames	<a href="http://www.geonames.org/ontology#">http://www.geonames.org/ontology#</a>	6 (2.03 %)	Data sets that use geonames
frbr	<a href="http://purl.org/vocab/frbr/core#">http://purl.org/vocab/frbr/core#</a>	6 (2.03 %)	Data sets that use frbr
xsd	<a href="http://www.w3.org/2001/XMLSchema#">http://www.w3.org/2001/XMLSchema#</a>	6 (2.03 %)	Data sets that use xsd
time	<a href="http://www.w3.org/2006/time#">http://www.w3.org/2006/time#</a>	5 (1.69 %)	Data sets that use time
event	<a href="http://purl.org/NET/c4dm/event.owl#">http://purl.org/NET/c4dm/event.owl#</a>	5 (1.69 %)	Data sets that use event
dbpedia	<a href="http://dbpedia.org/resource/">http://dbpedia.org/resource/</a>	5 (1.69 %)	Data sets that use dbpedia
gr	<a href="http://purl.org/goodrelations/v1#">http://purl.org/goodrelations/v1#</a>	4 (1.36 %)	Data sets that use gr

<http://www4.wiwiss.fu-berlin.de/lodcloud/state/#structure>

## Dates, roles and topics

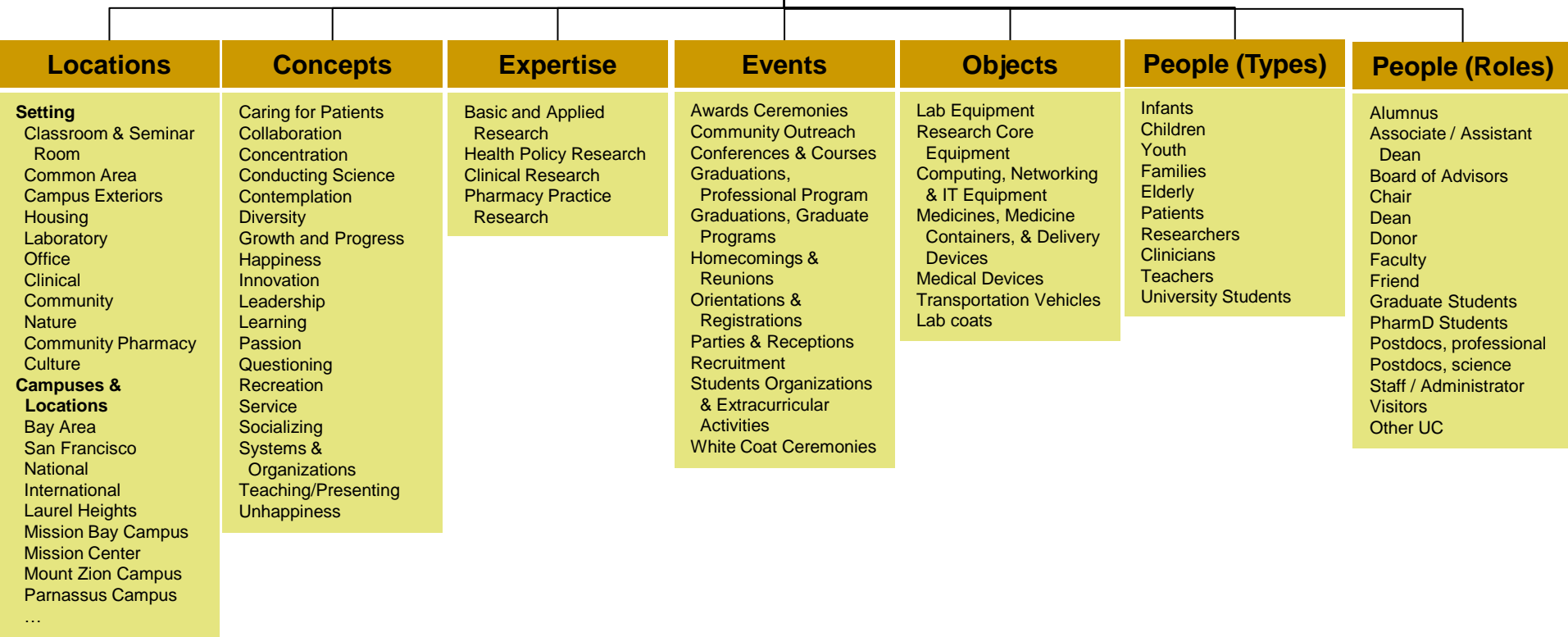
Explicit dates associated with an event in the life cycle of a resource.  
Best practice: YYYY-MM-DD (*W3C DTF Profile of ISO 8601*)

Property	Description	Set By
date.added	Date the asset was first added to the DAM.	DAM
date.lastModified	Date the asset was last reviewed for accuracy and relevance. Used for provenance and to validate content or rights.	DAM
date.reviewed	Date the content was last reviewed for accuracy and relevance. Used for provenance, and to compute a future date to recheck the content.	DAM
date.nextReviewed	Date of next scheduled review for accuracy and relevance.	Rule
date.embargoed	Date and time that content is scheduled to become available on the site. Content can be prepared in advance and system will push it out once the embargo date is reached.	Manual
date.subject	Date of the event, data, or other information depicted in the asset. Used for search and recall purposes. (This is not the date the asset was uploaded or last updated).	Manual

# Dates, roles and topics

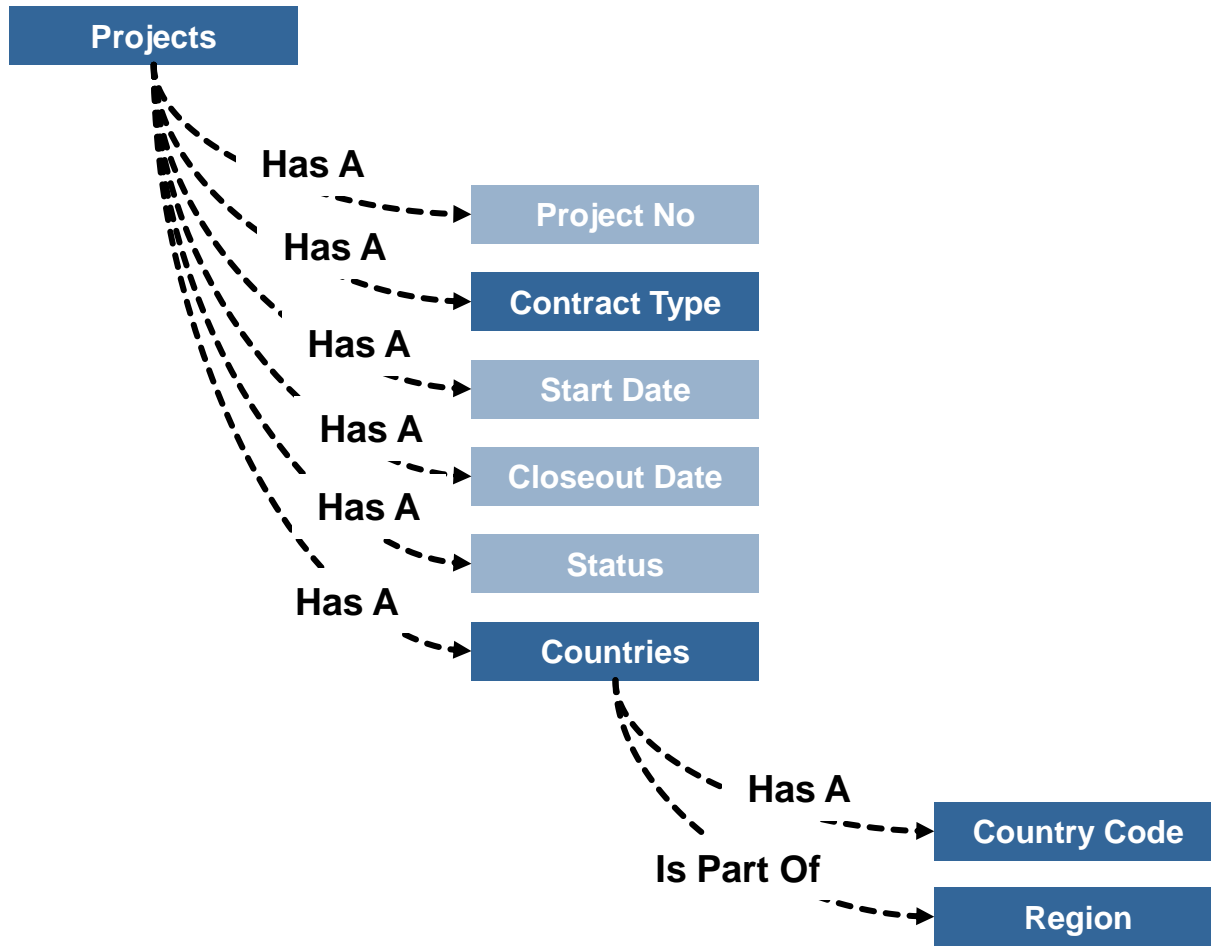
Role	Description	Admin	Add	Edit	Delete	Approve	Review
Administrator	Technical administration of the DAM. Generally allowed to do anything, to keep the system running and up-to-date.	Y	Y	Y	Y	Y	Y
Approver	Senior DAM staff with the authority to approve assets for publication. In small shops Contributors may also be Approvers. Larger shops, and those using outsider contractors will have many Contributors but just a few Approvers.	N	Y	Y	Y	Y	Y
Contributor	Editorial staff with authority to contribute new assets to the DAM. Their work must be approved by an Approver before it can be published. Administrators have the authority to approve content for publication, but only as an exception not the rule.	N	Y	Y	N	N	Y

# Dates, roles and topics

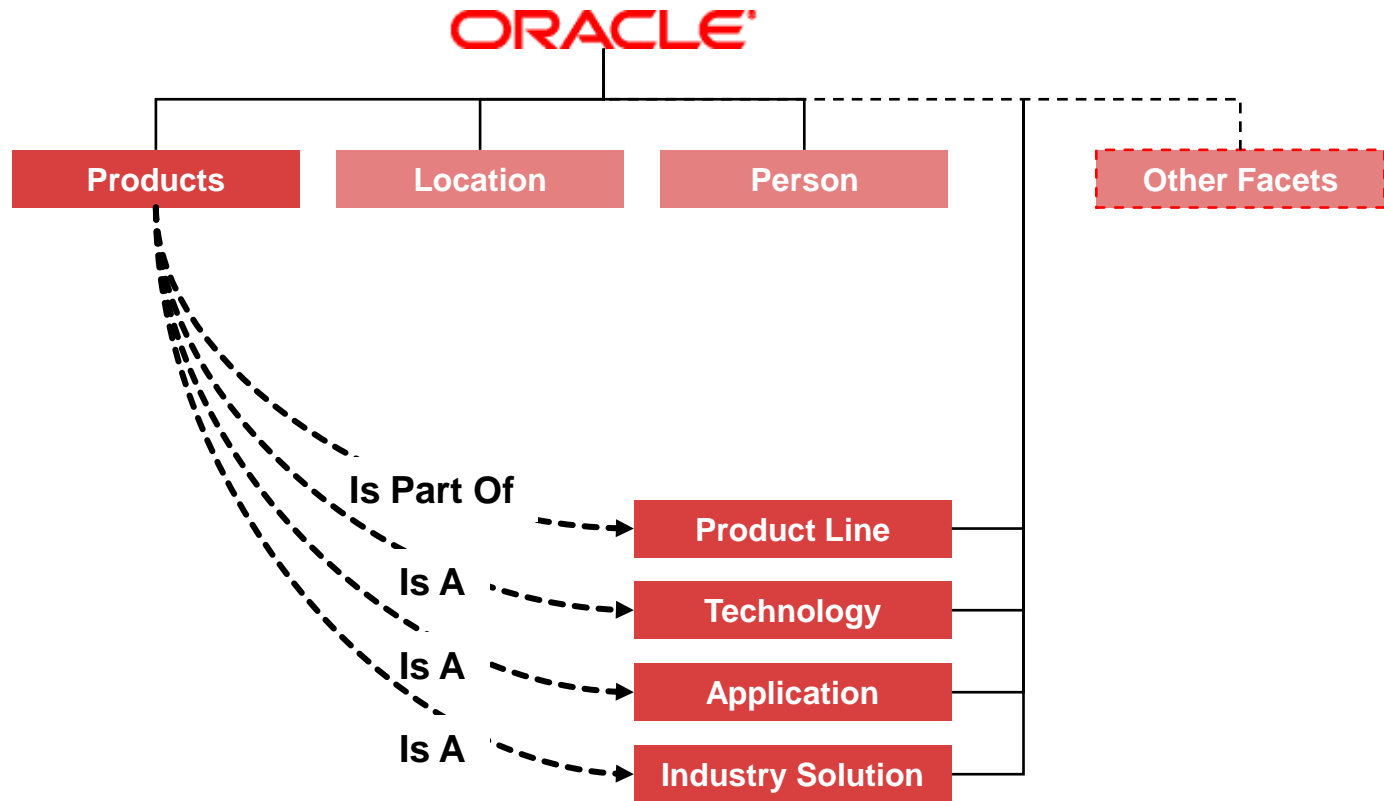




# Semantic relationships among topics



# Semantic relationships among topics



- Oracle events <http://events.oracle.com/>
- Oracle press releases <http://pressroom.oracle.com/>

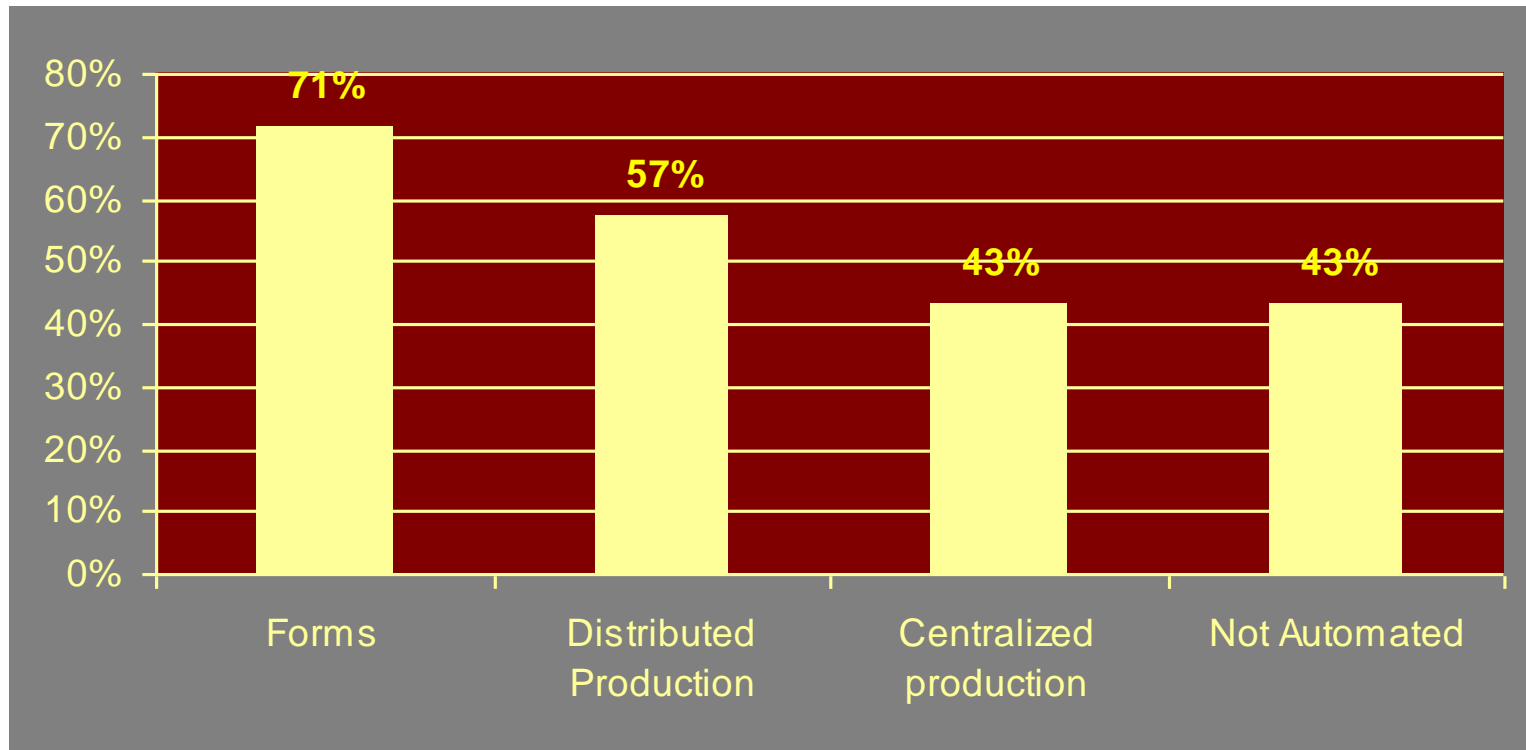
# The Tagging Problem

- ❖ How are we going to populate metadata elements with complete and consistent values?
- ❖ What can we expect to get from automatic classifiers?

# Cheap and Easy Metadata

- ❖ Some fields will be constant across a group or collection of content items
  - E.g., project, project no, contract type, start date, closeout date, status, country and region
  - E.g., product, product line, technology, application and industry solution
- ❖ In the context of a single group or collection those kinds of elements may add little value, but they add tremendous value when many groups or collections are brought together into one place, and they are cheap to create and validate.

# Methods used to create & maintain metadata



- ❖ Web-based forms are most widely used:
  - Distributed resource origination metadata tagging
  - Centralized clean-up and metadata entry.

**Source:** CEN/ISSS Workshop on Dublin Core.

# Tagging considerations

- ❖ Who should tag assets? Producers or editors?
- ❖ Taxonomy is often highly granular to meet task and re-use needs, but with detailed taxonomy it's difficult to get complete and consistent tags.
- ❖ The more tags there are (and the more values for each tag), the more hooks to the content, but the more difficult it is to get completeness and consistency.
- ❖ If there are too many tags or tags are too detailed, producers will resist and use “general” tags (if available)
- ❖ Vocabulary is often dependent on originating department, but the lingo may not be readily understood by people outside the department (who are often the users).

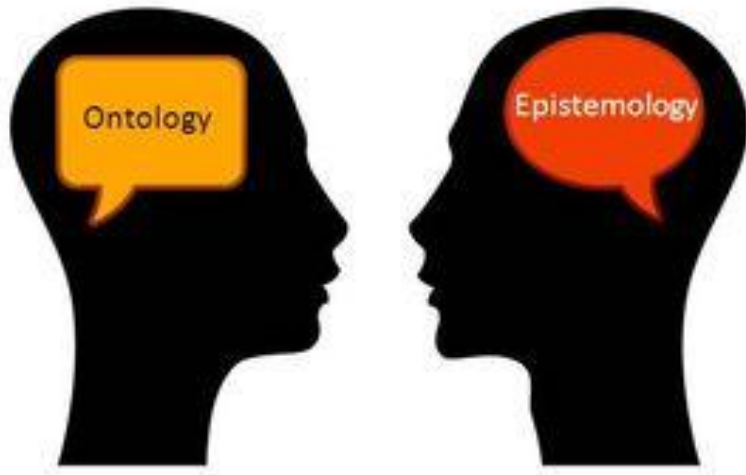
## Tagging considerations (2)

- ❖ Automatic classification tools exist, and are valuable, but results are not as good as people can do.
  - “Semi-automated” is best.
  - Degree of human involvement is a cost/benefit tradeoff.

# Tools for tagging

Vendor	Taxonomy Editing Tools	URL
	Autonomy Collaborative Classifier	<a href="http://www.autonomy.com/content/Functionality/idol-functionality-categorization/index.en.html">www.autonomy.com/content/Functionality/idol-functionality-categorization/index.en.html</a>
	ConceptSearching	<a href="http://www.conceptsearching.com">www.conceptsearching.com</a>
	Data Harmony M.A.I. <sup>TM</sup> (Machine Aided Indexing)	<a href="http://www.dataharmony.com/products/mai.html">www.dataharmony.com/products/mai.html</a>
	Microsoft Office Properties	<a href="http://office.microsoft.com/en-us/access-help/view-or-change-the-properties-for-an-office-file-HA010354245.aspx?CTT=1">office.microsoft.com/en-us/access-help/view-or-change-the-properties-for-an-office-file-HA010354245.aspx?CTT=1</a>
	Intelligent Topic Manager	<a href="http://www.mondeca.com/Products/ITM">www.mondeca.com/Products/ITM</a>
	nStein TME (Text Mining Engine)	<a href="http://www.nstein.com/en/products-and-technologies/text-mining-engine/">www.nstein.com/en/products-and-technologies/text-mining-engine/</a>
	PoolParty Extractor	<a href="http://poolparty.biz/products/poolparty-extractor/">poolparty.biz/products/poolparty-extractor/</a>
	Semaphore Classification and Text Mining Server	<a href="http://www.smartlogic.com/home/products/semaphore-modules/classification-and-text-mining-server/overview">www.smartlogic.com/home/products/semaphore-modules/classification-and-text-mining-server/overview</a>
	Temis Luxid <sup>®</sup> for Content Enrichment	<a href="http://www.temis.com/?id=201&amp;selt=1">www.temis.com/?id=201&amp;selt=1</a>





Joseph Busch

[jbusch@taxonomystrategies.com](mailto:jbusch@taxonomystrategies.com)

[\(415\) 377-7912](tel:(415)377-7912)

[twitter.com/joebusch](https://twitter.com/joebusch)

Thank You

**QUESTIONS**