

# **Constructing the InterPARES thesaurus: Challenges and opportunities**

Jonathan Furner

# Four terminological tools

- Register
  - Dictionary
  - Glossary
  - Thesaurus
- 
- four views of one database

# Purpose

- to improve the effectiveness of IP-related communication and dissemination
- by helping authors, readers, and intermediaries to select and use terms **consistently**
- thus increasing the likelihood ...
  - that authors will use terms in the manner expected by their readers: effective **authorship**
  - that readers will find documents of the kind they want: effective **retrieval**
  - that readers will interpret documents in the manner intended by their authors: effective **understanding**

# Some general challenges

- potential **conflict** of purpose: is partly descriptive, partly prescriptive
  - founded on current usage ...
  - ... but specifying preferred future usage
- **specificity** of (often multiple) senses of technical terms
- **domain**-specific variation
- **region**-specific variation
- facet analysis

# Terms, qualifiers, definitions, concepts

- a qualified term is made up of ...
  - one **term**: e.g., “record”
  - three **qualifiers**:
    - a part of speech: e.g., “noun”
    - a facet: e.g., “object”
    - a context: e.g., “archival science”
- a qualified term + a **definition** = a **concept**
- a single term (homonym) may have multiple definitions
- different terms (synonyms) may have the same definition
- aim: disambiguate homonyms and collocate synonyms

# Dictionary records: Example

Ter m	Po S	Facet	C ontext	Defini t i o n	So u r c e
rec odr	n.	objec t	ar d. s c i .	A do c u m e n t c o n t a i n i n g i n f o r m a t i o n r e c e i v e d a n d m a i n t a i n e d b y a n a g e n c y , o r i g i n a t i n g o r i n d i v i d u a l i n p u r s u a n c e o f l e g a l o b l i g a t i o n s o r i n t e r t r a n s a c t i o n o f b u s i n e s s .	IP
rec odr	n.	objec t	ar d. s c i .	D o c u m e n t c o n t a i n e d b y a p h y s i c a l o r j u r i d i c a l p e r s o n n o t i n c o n f o r m a n c e w i t h a c t a c t i v i t y .	UBC
rec odr	n.	objec t	c o m p s c i .	A g r o u p i n g o f i n t e r r e l a t e d d a t a e l e m e n t s f o r m i n g t h e b a s i s u n i t o f a f i e l d .	SAA

# Dictionary

- vocabulary:
  - IP- “accepted” terms
  - one or more definitions per qualified term
  - one record per definition
- arrangement:
  - alphabetical
- data
  - definitions + sources

# Glossary records: Example

Ter m	Po S	Facet	C o t e x t	Defini t i o
rec odr	n.	objec t	ar d. s c i .	A do c u m e n t c o n t a i n i n g i n f o r m a t i o n r e c e i v e d a n d m a i n t a i n e d b y a n a g e n c y , c o m m u n i c a t i o n o r i n d i v i d u a l i n p u r s u a n c e o f l e g a l l i a b i l i t i e s o r i n t e r a c t i o n o f b u s i n e s s .
rec odr	n.	objec t	c o m p s c i .	A g r o u p i n g o f i n t e r r e l a t e d d a t a e l e m e n t s f o r m i n g t h e b a s i c u n i t o f a f i e l d



# Glossary

- vocabulary:
  - IP- “preferred” terms
  - one definition per qualified term
  - one record per qualified term
- arrangement:
  - alphabetical
- data:
  - IP- “preferred” definition

# Thesaurus records: Example

Qual i f d i e r m	L i n k t y p e	L i n k e d t e r m	S u b f a c e t
rec odr	BTG	do am et	
ac t i v e e c o r d	BTG	rec odr	level o f a c t i v i t
rel a b e r e o r d	BTG	rec odr	level o f r e l i a b i l i t

# Thesaurus

- vocabulary:
  - IP- “accepted” terms (cf. Dictionary)
- arrangement:
  - faceted and hierarchical
- data:
  - links among paradigmatically-related terms
    - hierarchical
      - genus-species
      - whole-part
      - class-instance
    - from IP-nonpreferred to IP-preferred

# Current facets

- agents
  - “juridical person”; “creator”
- objects
  - “fonds”; “active record”
- properties
  - “authenticity”; “date of receipt”
- actions
  - “authentication”; “emulation”
- disciplines
  - “archival science”; “records management”

# Register record: example

Ter m	rec odr
Par OfSpeech	n
Facet	objec t
C ontext	ar b. s c i .
Reg i ertC eat el:Date	20030901
Reg i ertC eat el:By	jpm
Reg i ertLast Modifi el:Date	20030901
Reg i ertLast Modifi el:By	jpm
D it inary:C eat el:Date	20031001
D it inary:C eat el:By	ym c
D it inary:Las Modifi el:Date	20031001
D it inary:Las Modifi el:By	ym c
D it inary:Deleted:Date	
D it inary:Deleted:By	
Gl o asyCr ated:Dat e	20031001
Gl o asyCr ated:By	ym c
Gl o asyLas Mo dfied:Da e	20031001
Gl o asyLas Mo dfied:By	ym c
Gl o asyApp rved:Date	20031101
Gl o asyApp rved:By	ld
Gl o asyDel eet:Dat e	
Gl o asyDel eet:By	
Thesaurus:C eat el:Date	20040101
Thesaurus:C eat el:By	nr
Thesaurus:Las Modi fed:Date	20040101
Thesaurus:Las Modi fed:By	nr
Thesaurus:Appr oed:Dat e	20040201
Thesaurus:Appr oed:By	ld
Thesaurus:Delet el:Date	
Thesaurus:Delet el:By	

# Register

- **vocabulary:** all terms considered for acceptance into Dictionary
  - terms selected manually from automatically-created concordance of IP documents (case-study reports, data models, etc.)
  - “entailed” terms extracted from Glossary definitions and suggested by Thesaurus structure
  - terms selected manually from non-IP term lists
  - terms suggested by IP researchers
  - homonyms disambiguated by part-of-speech, facet, and contextual qualifiers
  - new terms added continuously as set of source documents grows
- **arrangement:** alphabetical
- **data:** dates and agents of creation, modification, approval, deletion of entries in Dictionary, Glossary, Thesaurus

# Some specific challenges

- selection of additional vocabulary
  - which sources should be consulted?
  - what procedures and criteria should be used?
  - should composite terms be selected?
- establishing consistency among products of various IP groups
  - how?
- assignment of terms to facets
  - is the selection of facets useful and/or appropriate?
- categorization of terms within facets
  - within each facet, what set of categories would be useful and/or appropriate?

# Some specific challenges, cont'd

- reconciliation of conflicting recommendations, in the various standards, as to form of terms
- maintenance of integrity of terminological database
- evaluation
  - how are tools to be evaluated w.r.t.
    - complying with international standards?
    - meeting users' requirements?