

Wiktionnaire, RDF, Linked Data: une colonne vertébrale pour le lexique ?

GETALP-LIG, Université Joseph Fourier, Grenoble, France

A word cloud featuring the word "Language" in large blue letters at the top center. Below it, the word "parole" appears in green. Other words include "lingua" in red, "bahasa" in purple, "język" in yellow, "lóngos" in orange, "palabra" in light blue, "sprak" in dark blue, "الكلام" in brown, "اللغة" in grey, "ภาษา" in pink, "word" in teal, "speech" in light green, "γλώσσα" in dark green, "言語" in black, and "שפה" in white. The background is a solid light blue.



- 1 Introduction
- 2 Apparté: le projet Papillon
- 3 État actuel de mes réflexions
- 4 Dbnary
- 5 RDF, Linked data, késako ?
- 6 Conclusion

Propos liminaires

Cette présentation à été préparée avec une méthode ancestrale...

C'est la méthode **ALARACHE** !!

Cela signifie que je compte sur vous! Rendez la vivante et complète en m'interrompant!

Ou nous nous ferons tous Magn(ennuyer)



Propos liminaires

Cette présentation à été préparée avec une méthode ancestrale...

C'est la méthode **ALARACHE** !!

Cela signifie que je compte sur vous! Rendez la vivante et complète en m'interrompant!

Ou nous nous ferons tous ch... comme la mort (bored to death)



Qui suis je ?

Je travaille sur la structuration/gestion/... des bases lexicales multilingues (MLDB) depuis 1990.

Thèse: définition d'une approche pivot pour les MLDB et d'un système générique de gestion des MLDB (1994)

Responsable projet Papillon depuis 2000.



Mes contributions principales ?

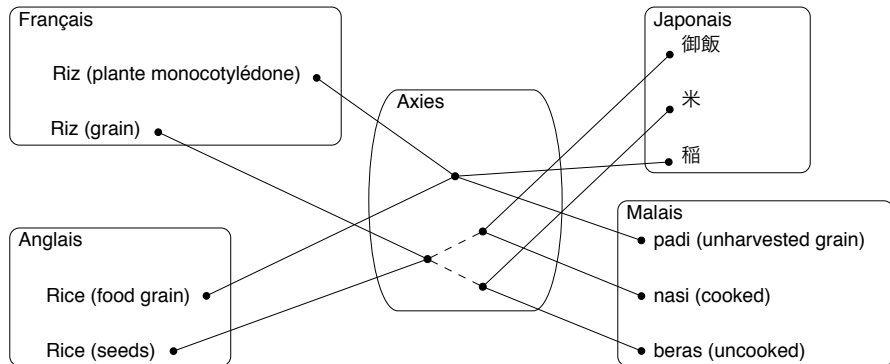
- Approche par acceptions interlingues (Axes)
 - ▶ évite les biais de l'utilisation d'une langue naturelle comme pivot
 - ▶ Reprise dans le standard Lexicon Markup Framework
- Outils génériques pour les MLDB
 - ▶ plateforme JIBIKI
 - ▶ utilisée pour le projet Papillon
 - ▶ réutilisée (ainsi que l'approche par acception) pour le projet LexALP
 - ▶ ... GDEF, DILAF (Mathieu Mangeot)
- Réflexion sur les lexiques, leur acquisition
 - ▶ Modèle de graphes: Lexical systems,
 - ▶ corpus vu comme un graphe
 - ▶ extraction d'informations lexicales vu comme des opérations sur des graphes
 - ▶ ... Thèse de Vincent Archer



- 1 Introduction
- 2 Apparté: le projet Papillon
- 3 État actuel de mes réflexions
- 4 Dbnary
- 5 RDF, Linked data, késako ?
- 6 Conclusion

Architecture of the Data

Macrostructure: An acception based multilingual lexical database



- 1 Introduction
- 2 Apparté: le projet Papillon
- 3 État actuel de mes réflexions
- 4 Dbnary
- 5 RDF, Linked data, késako ?
- 6 Conclusion



Mes réflexions actuelles

- Interopérabilité "syntaxique"

- ▶ Biais d'XML dans jibiki
- ▶ ... Passage en RDF...
- ▶ ... Lexical Linked Data.

- Interopérabilité "sémantique"

- ▶ Relier/aligner les sens de différentes ressources
- ▶ techniques proches du WSD
- ▶ Besoin de sens décrits, **avec des définitions**
- ▶ ... difficulté d'en construire.

- Multilinguisme

- ▶ modèle par acceptions toujours d'actualité
- ▶ besoin de bases lexicales multilingues encore d'actualité (même si on est passé d'une disette à une sur-abondance d'information lexicale)
- ▶ importance des informations adjectivales, verbales, adverbiales
- ▶ focalisation du domaine sur le nominal (wikipedia)



- 1 Introduction
- 2 Apparté: le projet Papillon
- 3 État actuel de mes réflexions
- 4 Dbnary**
- 5 RDF, Linked data, késako ?
- 6 Conclusion

Dbnary: Wiktionary comme graphe lexical

English [\[edit\]](#)

Pronunciation [\[edit\]](#)

- (*UK*) IPA: /kæt/, [kʰæt]

- (*US*) IPA: /kæt/, [kʰæʔ(̩)], [kʰeəʔ]

- Audio (UK)  0:00 [MENU](#)

- Audio (US)  0:00 [MENU](#)

- Audio (US-Inland North)  0:00 [MENU](#)

- Rhymes: –æʔ

Etymology 1 [\[edit\]](#)

From Middle English *cat*, *catte*, from Old English *catt* ("male cat") and *catte* ("female cat"), from Late Latin *cattus* ("domestic cat"), from Latin *catta* (c.75 B.C., Martial)^[1], from Afro-Asiatic (compare Nubian *kadis*, Berber *kaddiska* "wildcat"), from Late Egyptian *ḥaute*,^[2] feminine of *ḥaus* "jungle cat, African wildcat", from earlier Egyptian *teṣau* "female cat". Cognate with Scots *cat* ("cat"), Welsh *cath* ("cat"), West Frisian *kat* ("cat"), North Frisian *kāt* ("cat"), Dutch *kat* ("cat"), Low German *Katt*, *Katte* ("cat"), German *Katze* ("cat"), Danish *kat* ("cat"), Swedish *katt* ("cat"), Icelandic *köttur* ("cat"), Armenian *կատու* (*katu*, "cat").

Noun [\[edit\]](#)

cat (*plural* **cats**)

1. A domesticated subspecies, *Felis silvestris catus*, of feline animal, commonly kept as a house [pet](#). [from 8th c.]
2. Any similar animal of the family *Felidae*, which includes [lions](#), [tigers](#), bobcats, etc.
3. A catfish. [\[quotations ▼\]](#)
4. (*offensive*) A spiteful or angry [woman](#). [from earlier 13th c.]
5. An enthusiast or player of [jazz](#).
6. (*slang*) A person (usually male).
7. (*nautical*) A strong tackle used to hoist an anchor to the [cathead](#) of a ship.
8. (*chiefly nautical*) Short form of [cat-o'-nine-tails](#). [\[quotations ▼\]](#)
9. (*slang*) Any of a variety of earth-moving [machines](#). (from their manufacturer [Caterpillar Inc.](#))
10. (*archaic*) A sturdy merchant sailing vessel (now only in "catboat").
11. (*archaic, uncountable*) The game of "[trap and ball](#)" (also called "cat and dog").
12. (*archaic, uncountable*) The trap of the game of "trap and ball".
13. (*slang*) Prostitute. [from at least early 15th c.]
14. (*slang, vulgar, African American Vernacular*) A [vagina](#); female external genitalia [\[quotations ▼\]](#)
15. A double tripod (for holding a [plate](#), etc.) with six feet, of which three rest on the ground, in whatever position it is placed.



Wikipedia has an article on:

Cat



A domestic cat (1)



cat
penn
C18
jazz
Ingua
arole
men
Abyoc
jazz
Ingua
arole

Dbnary: Wiktionary comme graphe lexical

Synonyms [\[edit\]](#)

- (any member of the *suborder* (sometimes *superfamily*) *Feliformia* or *Feloidea*): **feliform** ("cat-like" *carnivoran*), **feloid** (compare *Caniformia*, *Canoidea*)
- (any member of the *family* *Felidae*): **felid**
- (any member of the *subfamily* *Felinae*, *genera* *Puma*, *Acinonyx*, *Lynx*, *Leopardus*, and *Felis*): **feline cat**, a **feline**
- (any member of the *subfamily* *Pantherinae*, *genera* *Panthera*, *Uncia* and *Neofelis*): **pantherine cat**, a **pantherine**
- (technically, all members of the *genus* *Panthera*): **panther** (i.e. *tiger*, *lion*, *jaguar*, *leopard*), (*narrow sense*) **panther** (i.e. **black panther**)
- (any member of the *extinct* *subfamily* *Machairodontinae*, *genera* *Smilodon*, *Homotherium*, *Miomachairodus*, *etc.*): *Smilodontini*, *Machairodontini* (=Homotherini), *Metailurini*, "**saber-toothed cat**" (**saber-tooth**)
- (*domestic species*): **housecat**, **puss**, **pussy**, **malkin**, **kitten**, **kitty**, **pussy-cat**, **mouser**, **tomcat**, **grimalkin**
- (*man*): **bloke** (*UK*), **chap** (*British*), **cove** (*UK*), **dude**, **fellow**, **fella**, **guy**
- (*spiteful woman*): **bitch**
- See also **Wikisaurus:cat**
- See also **Wikisaurus:man**

Derived terms [[edit](#)]

Terms derived from *cat* in the above senses

[show ▼]



Dbnary: Wiktionary comme graphe lexical

Translations [\[edit\]](#)

domestic species	[show ▼]
member of the suborder (or superfamily) Feliformia (Feloidea), "cat-like" carnivorans	[show ▼]
member of the family Felidae	[show ▼]
member of the subfamily Felinae	[show ▼]
member of the subfamily Pantherinae	[show ▼]
member of the extinct subfamily Machairodontinae	[hide ▲]
Select targeted languages	
<ul style="list-style-type: none"> Esperanto: maĥairodono ^(eo) French: machairodontiné ^(fr) <i>m</i>, machairodontinés ^(fr) <i>pl</i> 	<ul style="list-style-type: none"> German: Säbelzahnkatze ^(de) <i>f</i>, Machairodontine ^(de) <i>m</i>, Machairodontine ^(de) <i>f</i>, Machairodontinen ^(de) <i>pl</i>, Machairodontinae ^(de) <i>pl</i> Add translation <input type="text"/> : <input type="text"/> Preview translation More Script template: <input type="text"/> (e.g. Cyril for Cyrillic, Latn for Latin)
type of fish — <i>see</i> catfish	
spiteful woman — <i>see</i> bitch	
jazz enthusiast	[show ▼]
guy, fellow	[show ▼]
strong tackle used to hoist an anchor to the cathead of a ship	[show ▼]
cat-o'-nine-tails — <i>see</i> cat-o'-nine-tails	
type of boat — <i>see</i> catboat	
game of "trap and ball" (or "cat and dog")	[show ▼]
the trap in the game of "trap and ball"	[show ▼]



Dbnary: Wiktionary comme graphe lexical

==English==

[[Category:English three-letter words]]{{rfc-auto}}

{{wikipedia}}

[[Image:Cat03.jpg|thumb|A domestic cat (1)]]

===Pronunciation===

* {{a|UK}} {{IPA|/kæt|/[kʰæt]}}

* {{a|US}} {{IPA|/kæt|/[kʰæɾ(ɾ)]/[kʰeɾt]}}

* {{audio|En-uk-a cat.ogg|Audio (UK)}}

* {{audio|En-us-cat.ogg|Audio (US)}}

* {{audio|En-us-inlandnorth-cat.ogg|Audio (US-Inland North)}}

* {{rhymes|æt}}

===Etymology 1===

From {{etyl|enm|en}} {{term|cat|lang=enm}}, {{term|catte|lang=enm}}, from {{etyl|ang|en}} {{term|catt|male cat|lang=ang}} and

====Noun====

{{en-noun}}

A domesticated [[subspecies]], "[[Felis silvestris catus]]", of [[feline]] animal, commonly kept as a house [[pet]]. {{defdate|fro

Any similar animal of the family [[Felidae]], which includes [[lion]]s, [[tiger]]s, bobcats, etc.

A [[catfish]].

"'1913'", [[w:Willa Cather|Willa Cather]], "[[s:O Pioneers!|O Pioneers!]]", [[s:O Pioneers!/The Wild Land, II|chapter 2]]:

#*: She missed the fish diet of her own country, and twice every summer she sent the boys to the river, twenty miles to the sou

{{context|offensive|lang=en}} A spiteful or angry [[woman]]. {{defdate|from earlier 13th c.}}

An enthusiast or player of [[jazz]].

Dbnary: Wiktionary comme graphe lexical

=====Synonyms=====

- * {{sense|any member of the [[suborder]] (sometimes [[superfamily]]) [[Feliformia]] or {{taxlink|Feloidea|suborder}}}} [[feliform]]
- * {{sense|any member of the [[family]] [[Felidae]]}} [[felid]]
- * {{sense|any member of the [[subfamily]] [[Felinae]], genera "[[Puma]]", "[[Acinonyx]]", "[[Lynx]]", "[[Leopardus]]", and "[[Felis]]"}}
- * {{sense|any member of the subfamily [[Pantherinae]], genera "[[Panthera]]", [[Uncia]]" and "[[Neofelis]]"}} [[pantherine cat]],
- * {{sense|technically, all members of the genus "Panthera"}} [[panther]] (i.e. [[tiger]], [[lion]], [[jaguar]], [[leopard]]), {{qualifier|scientific name|}} [[scientific name|pantherine cat]]
- * {{sense|any member of the [[extinct]] subfamily "{{taxlink|Machairodontinae|subfamily}}", genera {{taxlink|Smilodon|genus|}} and {{taxlink|Homotherium|genus|}}}}
- * {{sense|domestic species}} [[housecat]], [[puss]], [[pussy]], [[malkin]], [[kitten]], [[kitty]], [[pussy-cat]], [[mouser]], [[tomcat]],
- * {{sense|man}} [[bloke]] {{qualifier|UK}}, [[chap]] {{qualifier|British}}, [[cove]] {{qualifier|UK}}, [[dude]], [[fellow]], [[fella]], [[guy]],
- * {{sense|spiteful woman}} [[bitch]]
- * See also [[Wikisaurus:cat]]
- * See also [[Wikisaurus:man]]



Dbnary: Wiktionary comme graphe lexical

====Translations====

{{trans-top|domestic species}}

* {{trreq|ab}}

* Acehnese: {{tø|ace|mië}}

* Adyghe: {{tø|ady|КІЭТЫ|sc=Cyrl}}

* Afrikaans: {{t+|af|kat}}

* Ainu: {{tø|ain|チヤペ|tr=cape}}

* Akan: [[agyanamo]] {{n}}

* Albanian: {{t+|sq|mace|f}}

* Alemannic German: {{tø|gsw|Chätz}}

* Amharic: {{t-|am|ደመት|tr=dəmət|sc=Ethi}}

* Apache:

*: Western Apache: {{tø|apw|gídí}}

* Arabic: {{t+|ar|قط|m|tr=qīṭṭ|sc=Arab}}, {{t+|ar|قطه|f|tr=qīṭṭa|sc=Arab}}

* Egyptian: {{tø|arz|قط|m|tr='uṭṭ}}, {{tø|arz|قطه|f|tr='uṭṭa}}

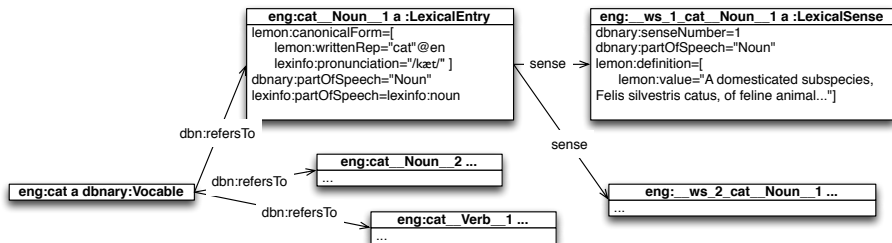
*: Libyan: {{t-|ar|قطوس|m|tr=gatṭūs|sc=Arab}}, {{t-|ar|قطوسة|f|tr=gatṭūsa|sc=Arab}}

* Moroccan Arabic: {{tø|ary|مش|tr=mešš}}, {{tø|ary|مشة|f|tr=mešša}}

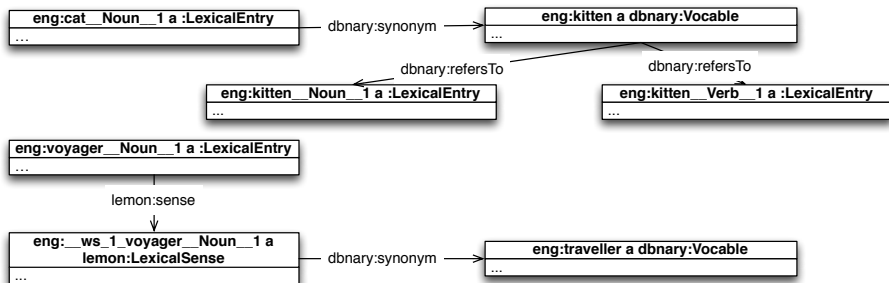
* Aramaic:

*: Syriac: [[ܫܢܪܐ]] (šūnārā') {{m}}, [[ܫܢܪܬܐ]] (šūnārtā') {{f}}

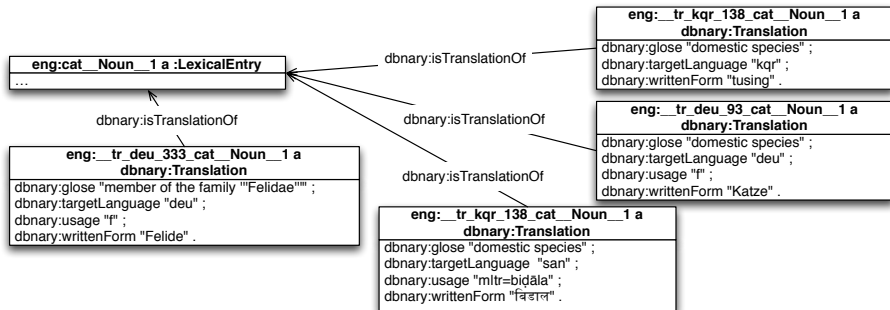
Dbnary: Wiktionary comme graphe lexical



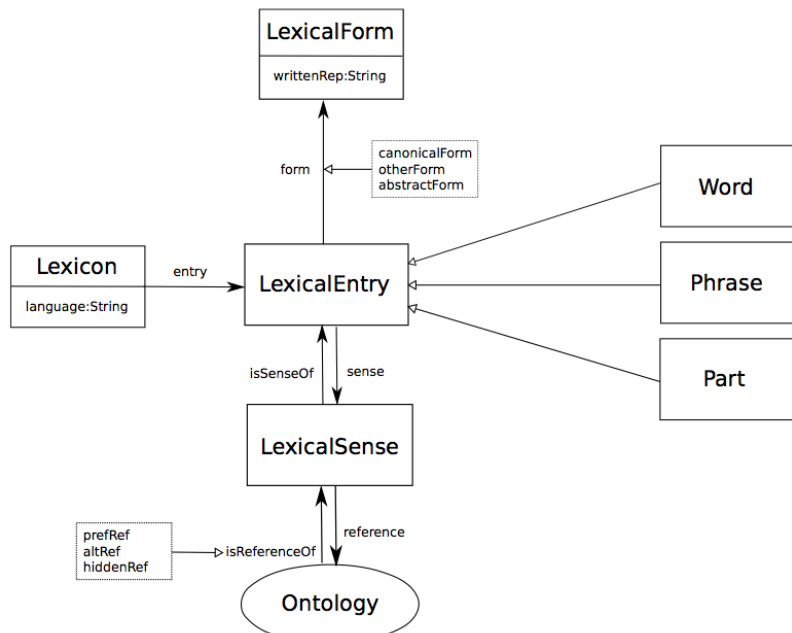
Dbnary: Wiktionary comme graphe lexical



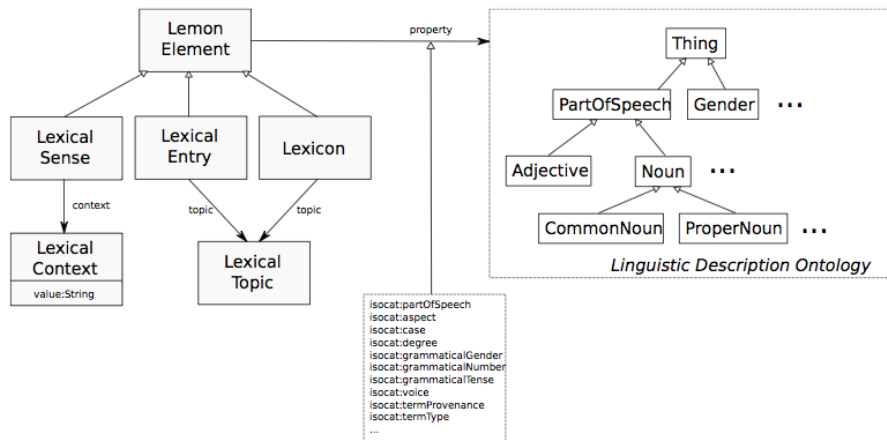
Dbnary: Wiktionary comme graphe lexical

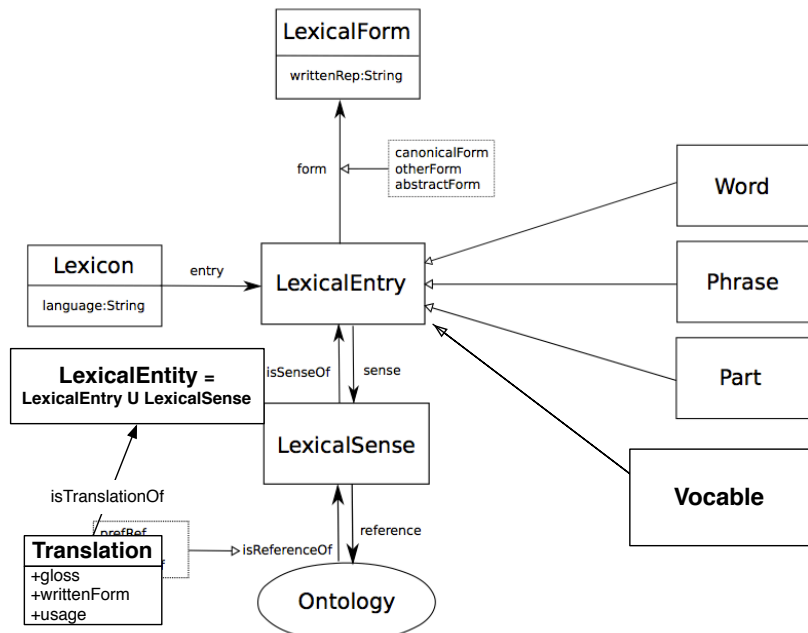


LEMON



LEMON





Taille des données (aujourd'hui)

	syn	ant	hyper	hypo	mero	holo
eng	31461	6877	959	1103	114	0
fra	30088	6735	8215	3557	943	1847
deu	27516	14315	30202	9509	0	0
rus	22631	9204	21028	4756	0	0
ell	3975	1116	0	0	0	0
fin	2255	0	0	0	0	0
por	3527	575	6	3	0	0
ita	7091	2337	0	0	0	0

Table : Number of lexicon-semantic relations. Languages are sorted according to their number of lexical entries.



Taille des données (aujourd'hui)

Source/Target	deu	ell	eng	fin	fra	ita	por
eng	62501	23794	1	74938	57959	37467	30256
fra	34608	7063	74687	7589	12	18806	17784
deu	0	2675	81015	4947	67143	41485	8872
rus	23056	3295	48559	3966	14776	12643	5567
ell	2242	2	10090	1056	8436	1470	1149
fin	8046	918	30103	0	6700	3856	2196
por	7000	2816	11284	4607	8720	7096	4
ita	4619	506	17539	925	4461	75	1219

Table : Number of translations from/to the 8 currently extracted languages. Source languages are sorted according to their number of lexical entries. Target languages are sorted by their ISO 639-3 language code. The number of different target languages is also given.



Taille des données (aujourd'hui)

Source/Target	por	rus	others	Total	# of lang
eng	30256	74837	764710	1126463	1143
fra	17784	7783	296624	464956	952
deu	8872	17354	248401	471892	355
rus	5567	0	206709	318571	490
ell	1149	1315	29892	55652	246
fin	2196	7997	58912	118728	329
por	4	4396	179142	225065	695
ita	1219	938	27514	57796	315

Table : Number of translations from/to the 8 currently extracted languages. Source languages are sorted according to their number of lexical entries. Target languages are sorted by their ISO 639-3 language code. The number of different target languages is also given.



Qualité des données

- Évaluer la qualité des données est difficile
- Qualité de l'extraction \neq qualité des données
- Pas d'application particulière actuellement
- Pas encore d'alignement avec d'autres ressources (wordnet, jeuxdemots, ...)
- Mais quelques indices...



Qualité des données (hier)

language	# of transl.
eng	5110 (99.1 %)
fra	5799 (107.0 %)
deu	10287 (99.2 %)
rus	8436 (24811.7 %)
ell	2598 (64.3 %)
fin	7245 (28980 %)
por	17720 (93.2 %)
ita	7855 (3167.3 %)

Table : Extracted translations vs interwiki links, on a random sample of 1000 entries.

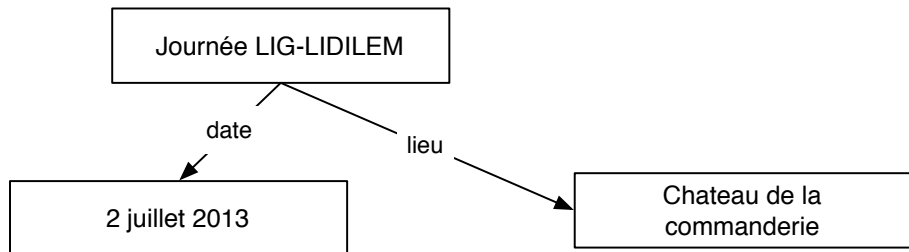


- 1 Introduction
- 2 Apparté: le projet Papillon
- 3 État actuel de mes réflexions
- 4 Dbnary
- 5 **RDF, Linked data, késako ?**
- 6 Conclusion

RDF: Resource Description Framework

Un graphe = un ensemble de **Statements**

un statement = un triplet: (S, P, O)

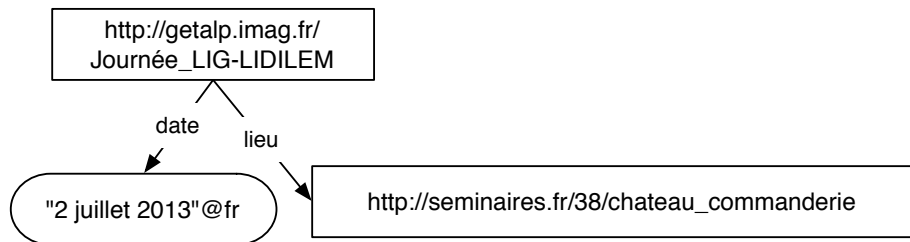


RDF: Resource Description Framework

Un graphe = un ensemble de **Statements**

un statement = un triplet: (S, P, O)

S: forcément une URI; P: une URI; O: une URI ou un littéral



RDF: Resource Description Framework

Un graphe = un ensemble de **Statements**

un statement = un triplet: (S, P, O)

Attention: une URI n'est pas forcément une document URL

`http://getalp.imag.fr/Journée_LIG-LIDILEM` : EROR 404: Not found !



RDF: Resource Description Framework

Exemple en turtle...

```
@prefix getalp: <http://getalp.imag.fr/>
```

```
@prefix sem: <http://seminaires.fr/38/>
```

```
@prefix dcterm: <http://purl.com/dcterm>
```

getalp:Journée_LIG-LIDILEM

```
dcterm:date "2 juillet 2013"@fr
```

```
dcterm:lieu sem:chateau_commanderie .
```



RDF: Resource Description Framework

Quoi d'autre ?

On peut aussi décrire les meta-donnée en RDF (OWL)

par exemple:

l'objet d'une relation "lieu" est forcément un lieu géographique.
un "dcterm:event" a forcément une date, un lieu et au moins un participant...

Un raisonnement est possible:

getalp:Journée_LIG-LIDILEM a dcterm:event
implique qu'il existe un participant...



Et le "Linked Data" ? Koikesse ?

En gros (très gros), les URIs deviennent des URLs

On peut donc "résoudre" (déréférencer) les noeuds de l'arc

Tout se passe comme si chaque noeud était à la fois
un objet du modèle... (comme en RDF)

une page web décrivant l'objet pour un utilisateur...

un document RDF décrivant l'objet pour un outil du web sémantique

C'est le serveur web qui "négocie" et redirige vers l'info adéquate.



- 1 Introduction
- 2 Apparté: le projet Papillon
- 3 État actuel de mes réflexions
- 4 Dbnary
- 5 RDF, Linked data, késako ?
- 6 Conclusion



