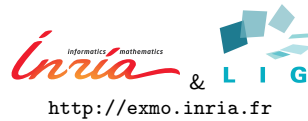
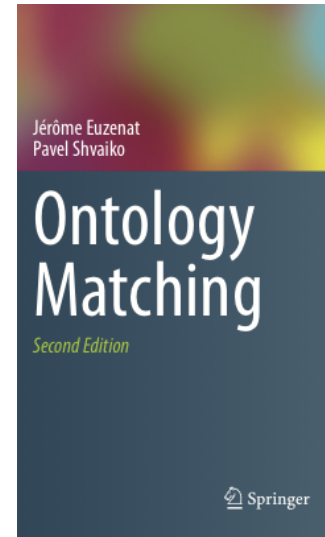


Aligner les ontologies pour communiquer

Jérôme Euzenat



May 16, 2014



- ▶ Échange de connaissance structurées médiatisées par ordinateur;
- ▶ Équipe INRIA associé au LIG (Grenoble);
- ▶ **Web sémantique – Données liées**
- ▶ Spécialiste de l'alignement d'ontologies (Ontology matching)
- ▶ Théorie – Logiciels – Applications.
- ▶ <http://exmo.inria.fr>

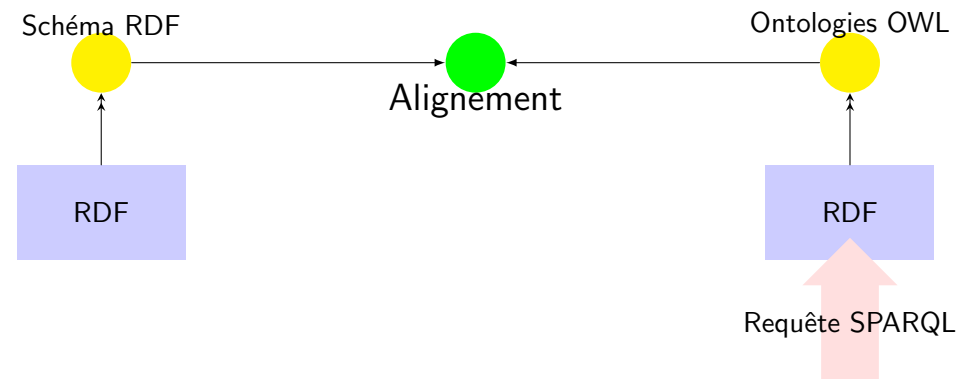
Problème

Faire coopérer des systèmes d'information indépendants

- ▶ Créer un nouveau système depuis zéro
- ▶ Échanger les données entre systèmes (une fois/continuellement)
- ▶ Créer une interface commune à ces systèmes

En général, il n'y a pas de raison de ne pas laisser ces systèmes indépendants.

Technologies du web sémantique?



Recommandations W3C

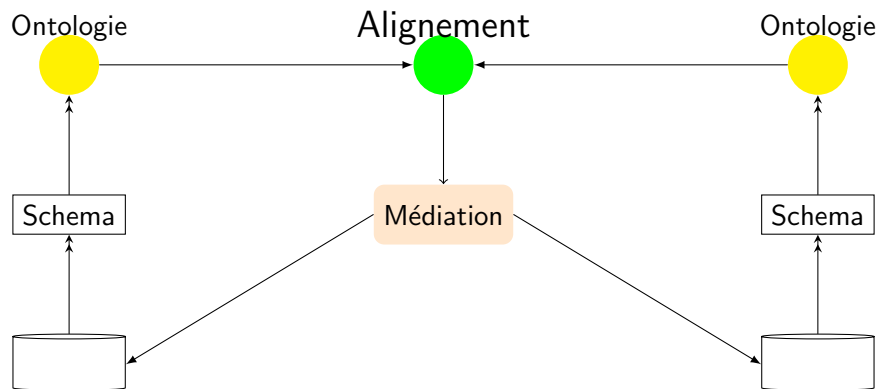
- ▶ Des technologies ouvertes (partout où le web passe);
- ▶ Des langages standardisés au niveau mondial;
- ▶ Elles sont suffisamment mûres et définies;
- ▶ Elles sont utilisées;
- ▶ Des solutions logicielles permettent de les mettre en œuvre.

Les technologies sémantiques sont là pour rester!

Les technologies du web sémantiques
(Représenter déclarativement les modèles par des "ontologies"
Lier ces ontologies par des alignements)

sont des outils appropriés

pour approcher l'interopérabilité.



Alignement d'ontologies

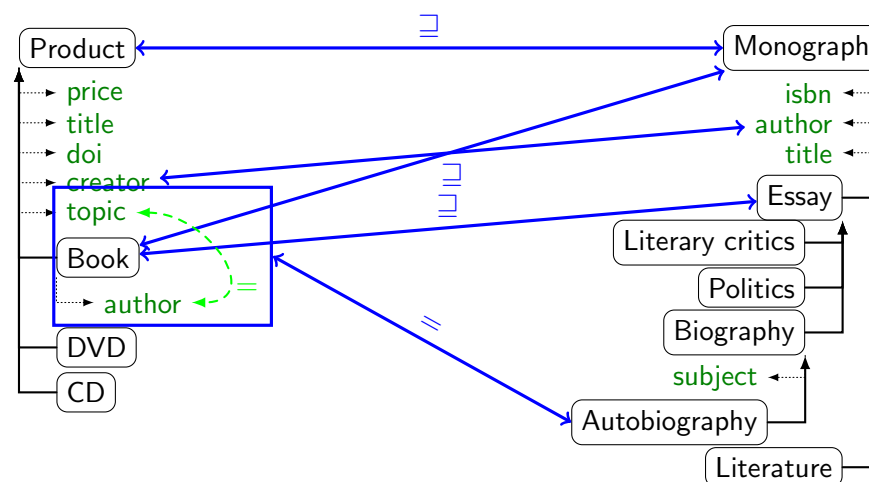
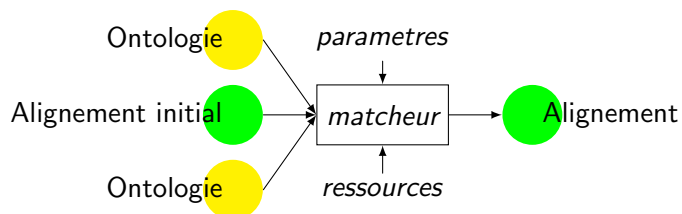
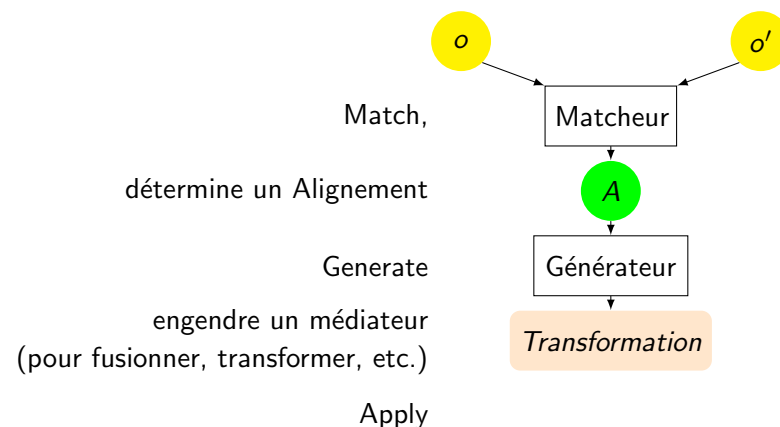
Exemple: liage de données

Conclusions

Les ressources exprimées de manières différentes doivent être réconciliées avant d'être utilisées.

L'hétérogénéité peut apparaître entre:

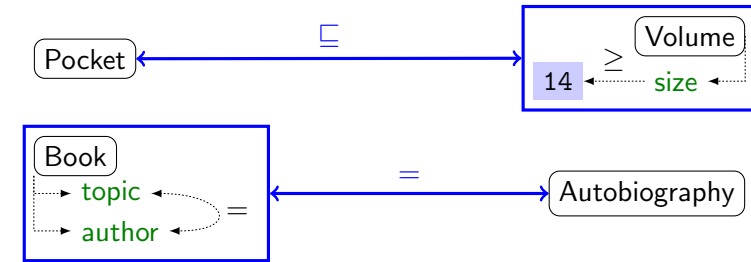
- ▶ différent langages de représentation de connaissance (OWL vs. RDFS);
- ▶ **différentes terminologies:**
 - ▶ Anglais vs. Chinois;
 - ▶ Livre vs. Monographie.
- ▶ **différents modèles peuvent être utilisées:**
 - ▶ différentes classes: Autobiographie vs. Livre de poche;
 - ▶ classes/propriétés: Essai vs. genrelittéraire = essai;
 - ▶ classe/nstances: Un livre physique ou une œuvre.
- ▶ **différents domaine ou granularité.**
 - ▶ Livres vs. artefact culturels vs. produits;
 - ▶ Livres détaillés jusqu'à la production et la traduction ou livres comme œuvres.



- ▶ ensemble de correspondances entre entités:

$$\langle e, \sqsubseteq, e' \rangle$$

- ▶ expriment déclarativement les relations entre entités;
- ▶ peuvent être obtenus par une multitude de techniques;
- ▶ utilisables pour transformer données, requêtes ou ontologies;
- ▶ pas de format normalisé.



$$\forall x, Pocket(x) \Leftarrow Volume(x) \wedge size(x, y) \wedge y \leq 14$$

$$\forall x, Book(x) \wedge author(x, y) \wedge topic(x, y) \equiv Autobiography(x)$$

```
SELECT ?d
WHERE {?x rdf:type o:Book .
      ?x o:creator ?y .
      ?x o:topic ?y .
      ?y o:name "Bertrand Russell" .
      ?x o:doi ?d .}
```

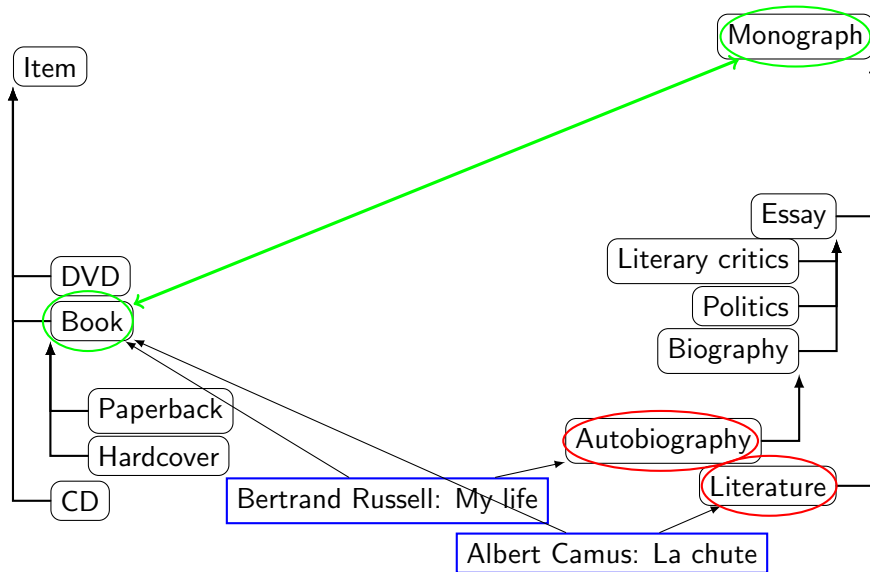
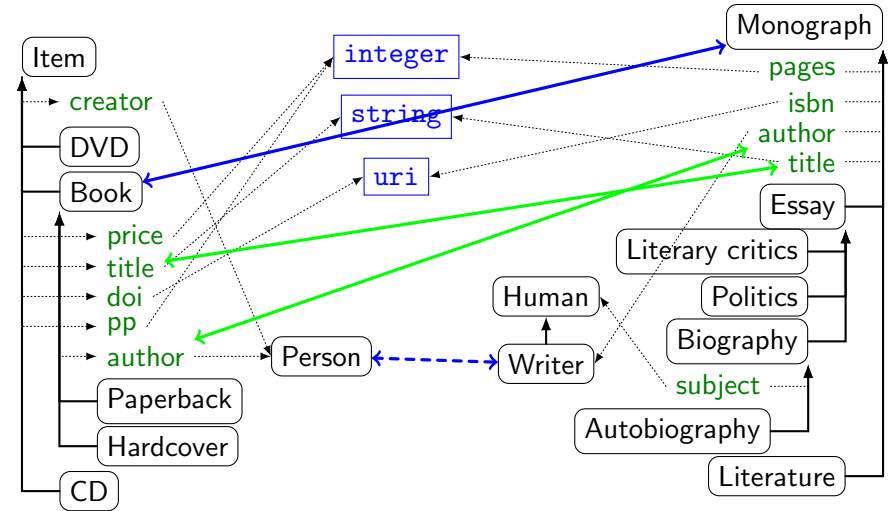
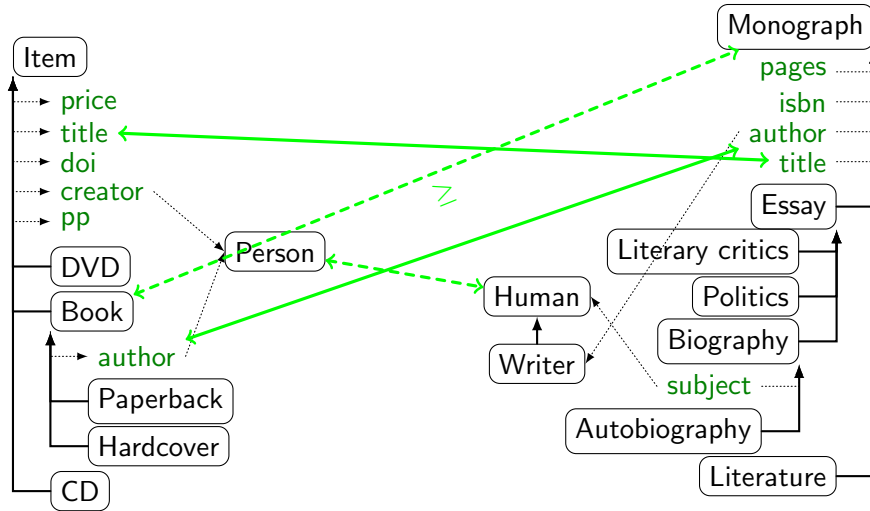
```
SELECT ?i
WHERE { ?x rdf:type o':Autobiography .
      ?x o':author/o':name "Bertrand Russell" .
      ?x o':isbn ?i .}
```

mediateur

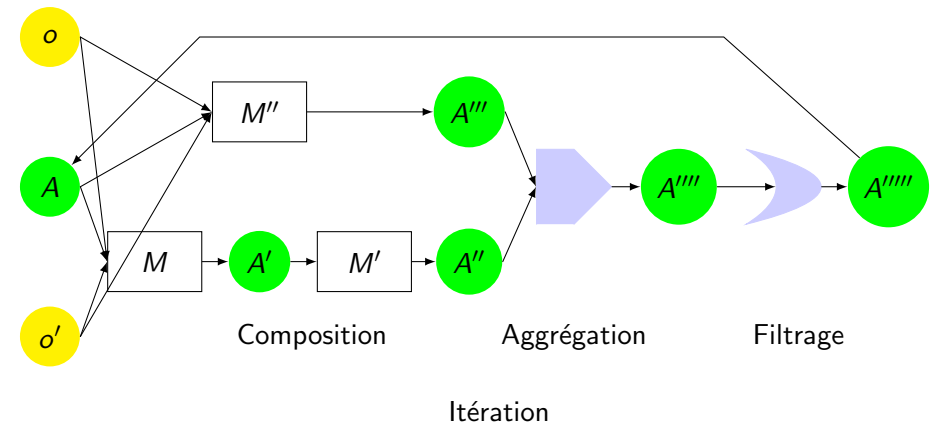
x.doi=http://dx.doi.org/10.1080/041522862X

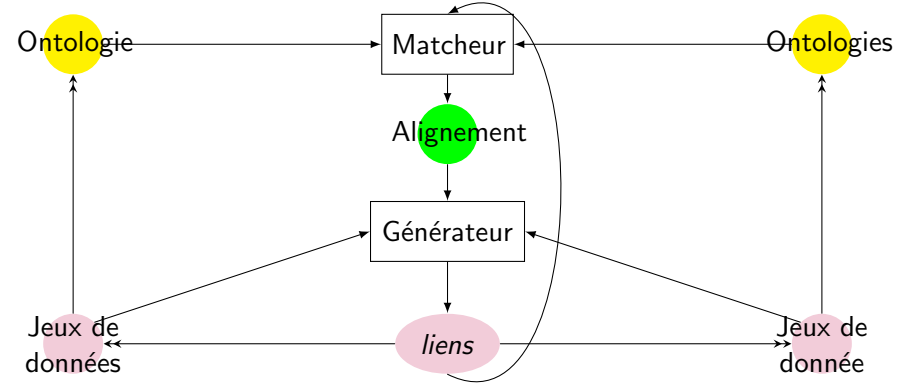
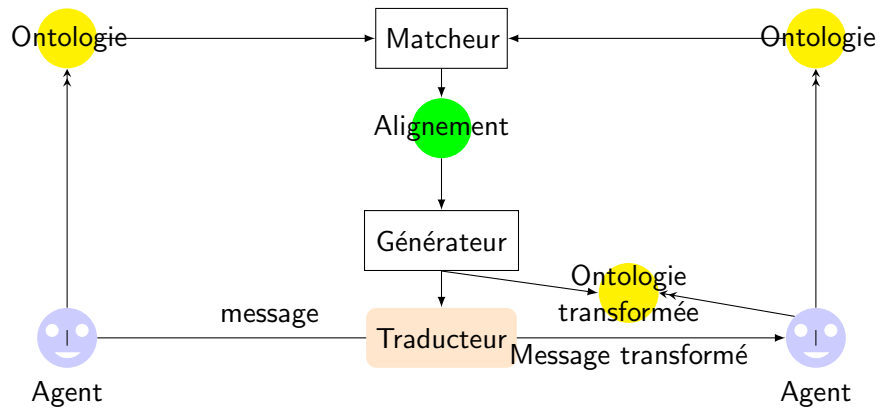
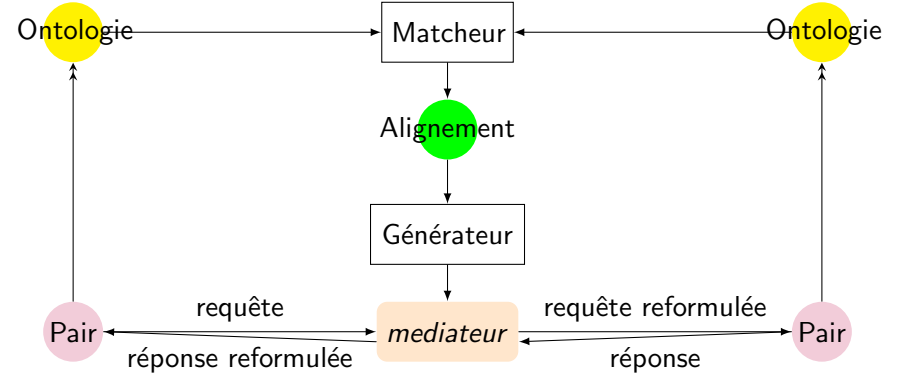
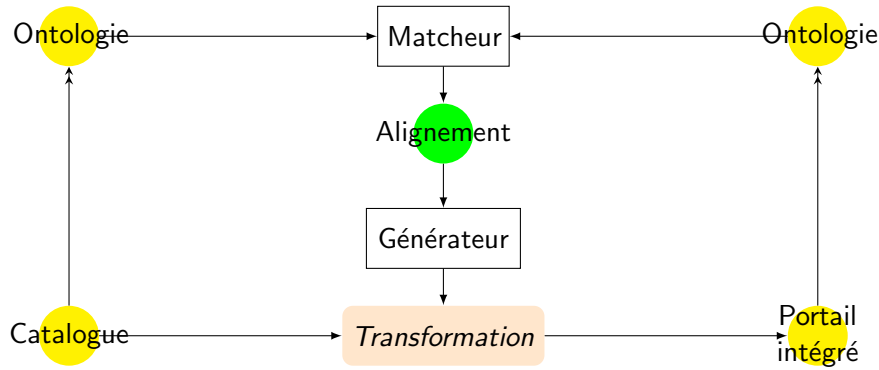
x.isbn=041522862X

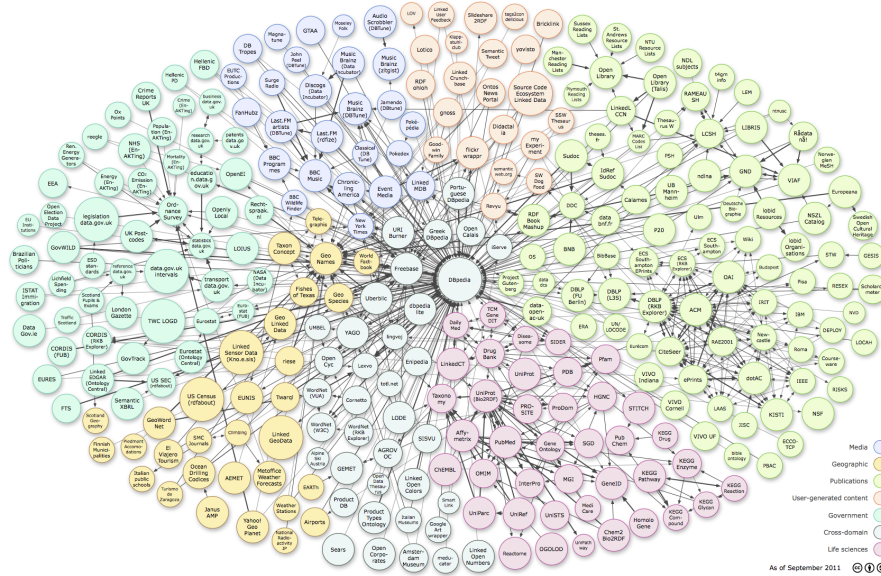
- ▶ Contenu: ce qui est exprimé dans les ontologies
 - ▶ **Noms**, commentaires, noms alternatifs, noms des entités connexes: NLP, RI, etc.
 - ▶ **Structure interne** (contraintes entre relations, typage)
 - ▶ **Structure externe** (relations entre entités): Fouille de données, mathématiques discrètes
 - ▶ **Extension**: Statistiques, analyse de données, fouille de données, apprentissage
 - ▶ **Sémantique** (models): Raisonnement automatique
- ▶ Contexte: les relations entre les ontologies et leur environnement:
 - ▶ Le **web**
 - ▶ **Ontologies** externes: LOV, dbpedia, etc.
 - ▶ **ressources** externes: WordNet, etc.
 - ▶ **Usage**: requêtes, services, choix



Les matcheurs élémentaires produisent des correspondances candidates. La plupart des systèmes en utilisent plusieurs et combinent leurs résultats.







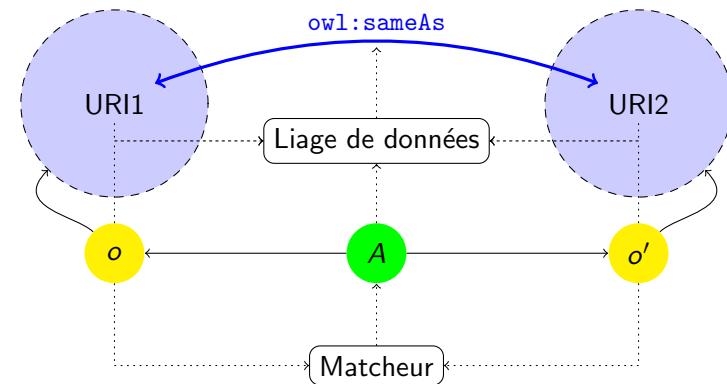
- ▶ Bibliothèques (comp. sci.: DBLP, ACM + nationales: Europeana, BNF, etc.);
- ▶ Musées (British museum, Rijksmuseum)
- ▶ Encyclopédique (BDPedia, Freebase, UMBEL, Cyc)
- ▶ Biologie (principalement moléculaire: GO, bioportal)
- ▶ Médecine (Pubmed, UMLS)
- ▶ Média (NYT, Reuters, BBC)
- ▶ Géographiques (geonames, Ordnance Survey, IGN)
- ▶ Administration (data.gov, data.gov.uk, data.fr)
- ▶ Statistiques (insee.fr, Eurostat)
- ▶ Agriculture (FAO)
- ▶ Produits (Schema.org, GoodRelations)

Problème: Soient deux sources de données RDF, trouver les paires d'entités de chaque source qui dénotent les mêmes ressources.

Trois directions:

- ▶ Extraction de clés ou de clés de liage;
- ▶ Similarités entre entités;
- ▶ Lier à partir d'alignements expressifs.

Projets: Datalift, Lindicle.



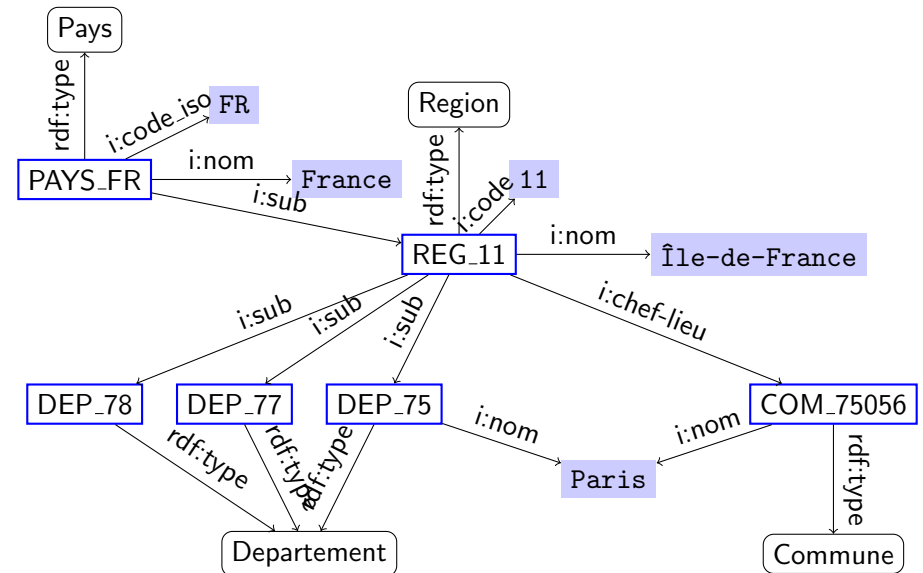
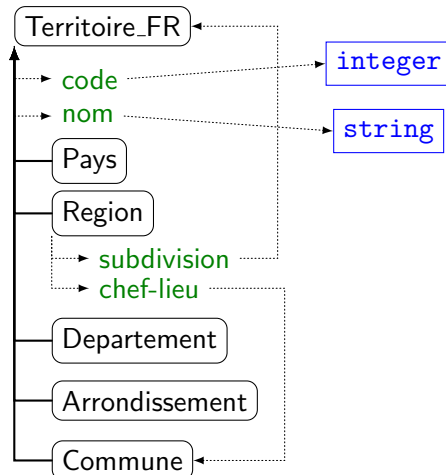
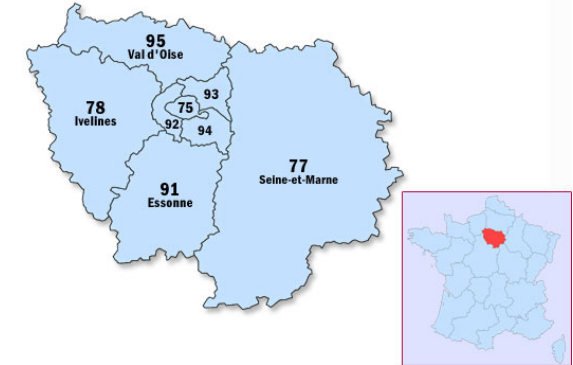
- ▶ Trouver des concepts équivalents [concept matching];
- ▶ Pour chacun, déterminer les propriétés équivalentes (basé sur les similarités entre leurs valeurs dans les jeux de données) [property matching];
- ▶ Trouver les combinaison de propriétés qui identifient des entités correspondantes [key extraction];
- ▶ Lier les entités qui correspondent [link generation].

Région table:

code	nom	chef-lieu
11	Île-de-France	75056
21	Champagne-Ardenne	51108
22	Picardie	80021

Sous-région table:

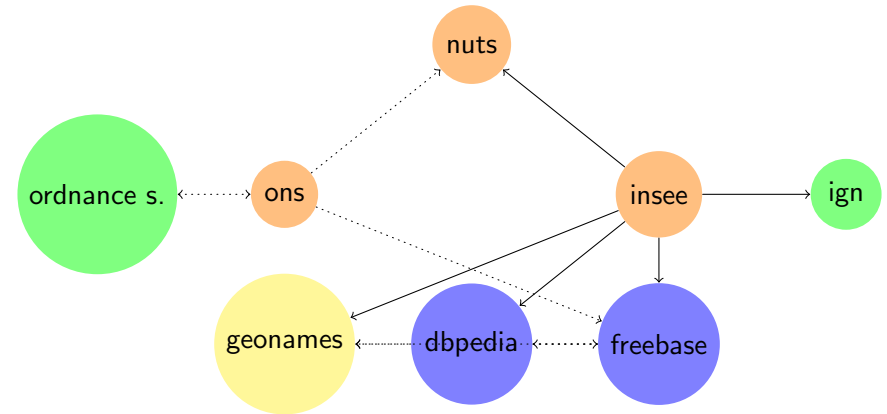
région	département
11	75
11	77
11	78
11	91
11	92
11	93



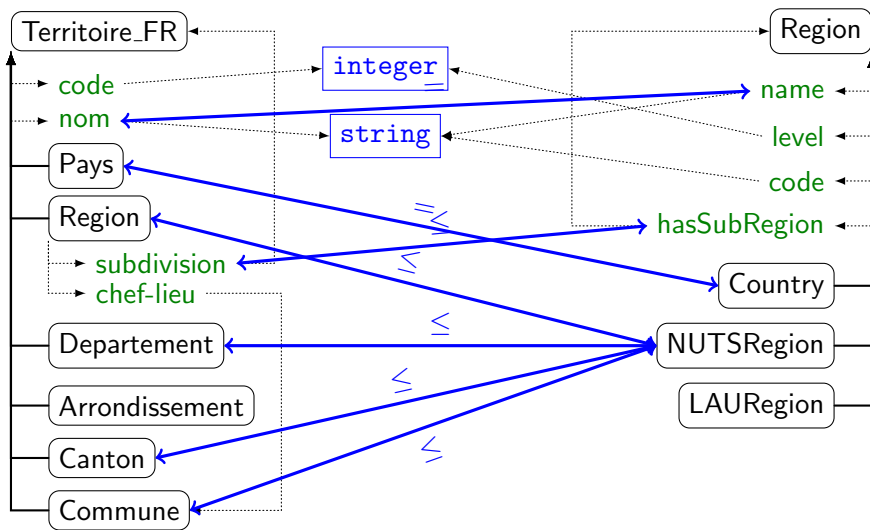
Jeux de données spécifiques contenant des liens entre entités.

```
<http://www.example.org/linkset/INSEE-NUTS>
  a void:Linkset ;
  void:target <http://rdf.insee.fr/geo/regions-2011.rdf>;
  void:target <http://nuts.psi.enakting.org/id/>;
```

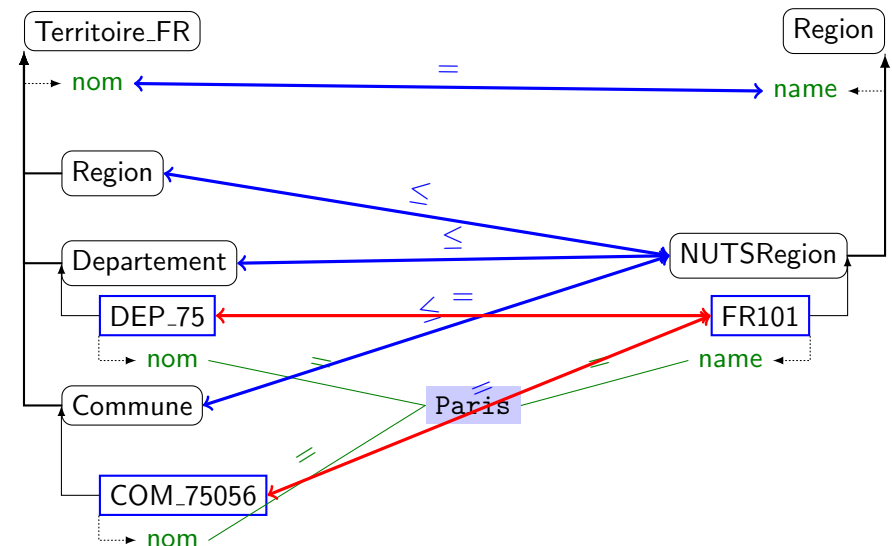
```
insee:PAYS_FR owl:sameAs nuts:FR
insee:REG_11 owl:sameAs nuts:FR10
insee:DEP_75 owl:sameAs nuts:FR101
insee:DEP_77 owl:sameAs nuts:FR102
insee:DEP_78 owl:sameAs nuts:FR103
```



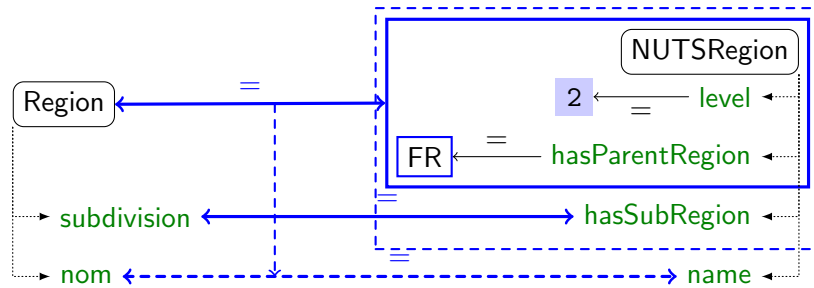
Alignement entre les ontologies INSEE et NUTS



Les alignements simples ne sont pas suffisants



Des alignements expressifs sont nécessaires (EDOAL)



Génération de requêtes (SPARQL)

```
PREFIX insee: <http://rdf.insee.fr/ontologie-geo-2006.rdf#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
SELECT ?r
FROM <http://rdf.insee.fr/geo/regions-2011.rdf>
WHERE {
  ?r rdf:type insee:Region .
}
```

```
PREFIX nuts: <http://ec.europa.eu/eurostat/ramon/ontologies/geographic.rdf#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
SELECT ?n
FROM <http://ec.europa.eu/eurostat/ramon/rdfdata/nuts2008/>
WHERE {
  ?n rdf:type nuts:NUTSRegion .
  ?n nuts:level 2^^xsd:int .
  ?n nuts:hasParentRegion nuts:FR1 .
}
```

Transformation de données

```
PREFIX insee: <http://rdf.insee.fr/ontologie-geo-2006.rdf#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX insee: <http://rdf.insee.fr/ontologie-geo-2006.rdf#>
CONSTRUCT {
  ?r rdf:type nuts:NUTSRegion .
  ?r nuts:level 2^^xsd:int .
  ?r nuts:hasParentRegion nuts:FR1 .
}
FROM <http://rdf.insee.fr/geo/regions-2011.rdf>
WHERE {
  ?r rdf:type insee:Region .
}
```

Génération de lien sameAs

```
CONSTRUCT { ?r owl:sameAs ?n . }
PREFIX insee: <http://rdf.insee.fr/ontologie-geo-2006.rdf#>
PREFIX nuts: <http://ec.europa.eu/eurostat/ramon/ontologies/geographic.rdf#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
FROM <http://rdf.insee.fr/geo/regions-2011.rdf>
FROM <http://ec.europa.eu/eurostat/ramon/rdfdata/nuts2008/>
WHERE {
  ?r rdf:type insee:Region .
  ?r insee:nom ?l .
  ?n rdf:type nuts:NUTSRegion .
  ?n nuts:name ?l .
  ?n nuts:level 2^^xsd:int .
  ?n nuts:hasParentRegion nuts:FR1 .
}
```

- ▶ Les alignements d'ontologies sont des expressions de correspondance au niveau des schémas;
- ▶ Ils sont utiles pour focaliser la recherche d'entités;
- ▶ Des alignements expressifs sont nécessaires;
- ▶ Ils peuvent être transformés en générateurs de liens en SPARQL.

mais il est aussi nécessaire d'exprimer des contraintes s'appliquant aux instances:

- ▶ pour convertir les données, par exemple de mph en m/s;
- ▶ pour exprimer les contraintes d'alignement sur les données, par exemple, par similarité.

- ▶ SILK: outil de liage de données partant de requête SPARQL et de mesures de similarité;
- ▶ RDB2RDF: standard de lien d'un SGBD relationnel et d'un schéma RDF (permet d'exporter les données ou de transformer les requêtes).

Les technologies du web sémantiques
(Représenter déclarativement les modèles par des "ontologies"
Lier ces ontologies par des alignements)

sont des outils appropriés

pour approcher l'interopérabilité.

- Jena API RDF + OWL API + SPARQL (Apache)
- Corese RDF+RDFS+ moteur SPARQL (INRIA Sophia)
- OWL API (U. Manchester)
- Sesame Triple store + moteur SPARQL (Aduna)
- Virtuoso Triple/Rel store + moteur SPARQL (OpenLink)
- Pellet Raisonneur OWL (Clark & Parsia)
- HermiT Raisonneur OWL (Oxford U.)
- Alignment API Gestion d'alignements (INRIA Grenoble)
- Protégé Éditeur RDF+OWL (Stanford)
- Datalift Plateforme de publication de données (public)

Jerome.Euzenat@inria.fr

<http://exmo.inria.fr>