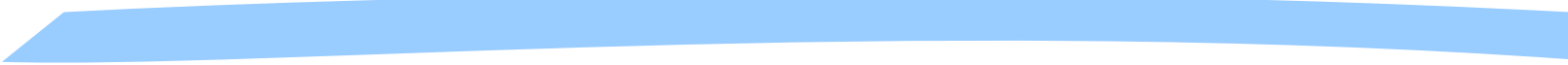


Construction d'une ontologie à partir d'un thésaurus de l'astronomie: vers une représentation plus sémantique des données textuelles



Nathalie Hernandez, Josiane Mothe

IRIT, 118 route de Narbonne, 31040 Toulouse cedex

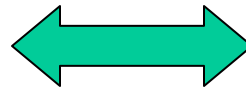
{hernande,mothe}@irit.fr

Contexte général



Systeme d'information

- Systeme de recherche d'information
- Systeme d exploration



- Contenu des documents
- Meta-données

→ Représentation des données textuelles

Problématiques



- Vocabulaire (besoin utilisateur, corpus)
 - Approches de la littérature : mots clés de la requête recherchés dans les documents [Salton 1971] [Spärck Jones 2003]
- perte du contexte et du vocabulaire
- ⇒ Maintenir ces liens pour s'adapter au vocabulaire de l'utilisateur, de la tâche, de la collection

Problématiques



- Indexation

- Indexation intégrant la variabilité des mots utilisés dans les collections vs indexation par mots clés
- ⇒ Meilleure représentation pour un meilleur appariement

- Accès

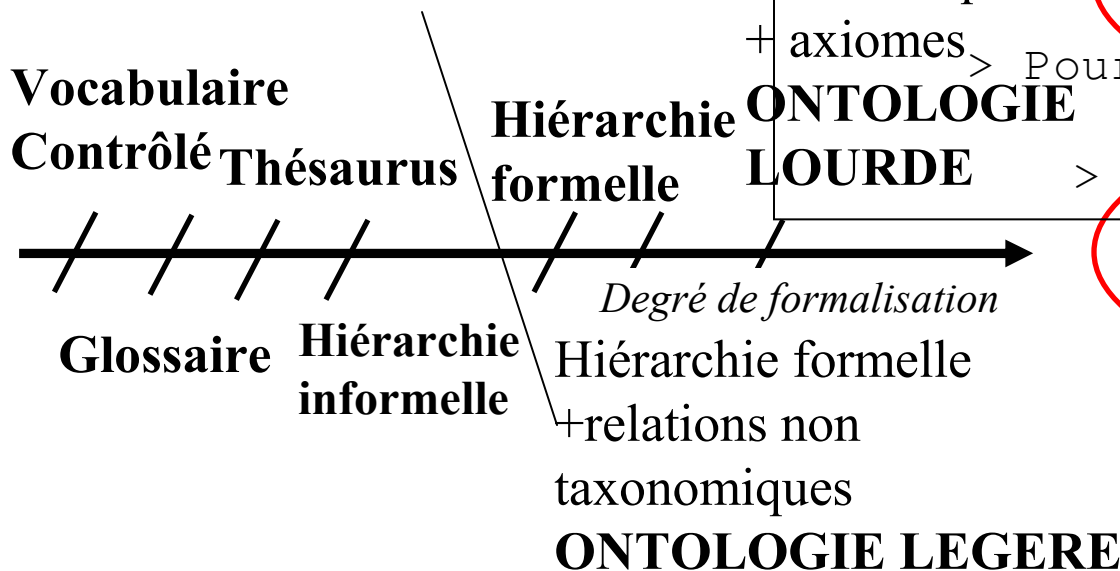
- Vues globales et spécifiques sur la connaissance modélisée et sur le contenu de la collection
- ⇒ Exploration du corpus à partir du contexte

⇒ Connaissances à représenter et intégrer dans le système

Différentes représentations de la connaissance

- Différents niveaux de formalisation

[Lassila 2001]



	X RAYS	
UF	Rayon X un photon à haute énergie	t
BT	ELECTROMAGNETIC WAVES	
NT	SPECTRE X	
RT	Accessoires	
	X RAY BACKGROUND	
	X RAY PHENOMENON	
	X RAY SOURCES	
	X RAY SPECTRA	
	Légende	
UF	terme utilisé pour	
BT	terme plus générique	
NT	terme plus spécifique	
RT	terme lié	

Plan



- Thesaurus vs ontologie
- Présentation générale de la méthode
- Description des 3 étapes principales
- Utilisation d'ontologies en Recherche d'Information
- Conclusion

Thesaurus

- Thesauri = ressources lexicales
 - Collection de termes organisés hiérarchiquement
 - Relations entre les termes
- Thesaurus IAU utilisés par des documentalistes pour
 - Indexer manuellement des documents
 - Formuler manuellement des requêtes

Principaux inconvénients

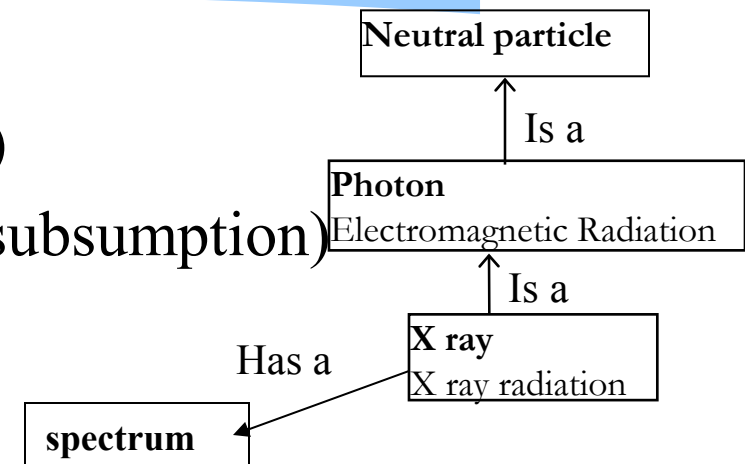
- Créés dans les années 1990
 - Ne contiennent pas la connaissance actuelle
- Normes sur leur contenu (ISO 2788 - ANSI Z39), MAIS aucun format uniforme (ascii, html, bases de données)
 - Outils pouvant les utiliser limités (visualisation, annotation, ...)
 - Utilisation limitée par les systèmes d'information (phase d'adaptation)

Principaux inconvénients

- Faible degré de formalisation dans la représentation de la connaissance
 - Pas de niveau d’abstraction (concepts, concepts génériques)
 - Pas de distinction entre un concept et sa lexicalisation
 - Relations entre termes ambigus (“est lié à”)
 - ⇒ représentation de domaines de connaissance en terme de terminologie et de catégories d’indexation et non pas en terme de sens
 - difficile à utiliser par des applications automatiques (eg indexation)

Potentiel des ontologies légères

- Ontologies légères =
 - concepts définis par des labels (termes)
 - concepts organisés hiérarchiquement (subsumption)
 - Relations associatives entre termes



- Référence pour la communication entre machines et humains
- Indexation sémantique de données hétérogènes

Elaboration d'ontologies

- Approches existantes
 - A partir de “rien” [Uschold 1996] [Guarino1998a] [Fernandez 1997]
 - A partir de textes [Maedche 2000] [Velardi 2001]
 - A partir de thésaurus (mais sans mise à jour de la connaissance) [Soergel 2004] [SKOS schéma w3c] [Hahn 2004]
- Notre méthode :
 - Tirer le bénéfice des termes du thésaurus
 - Extraire la connaissance implicite des thésaurus
 - Mettre à jour cette connaissance à partir de documents textuels

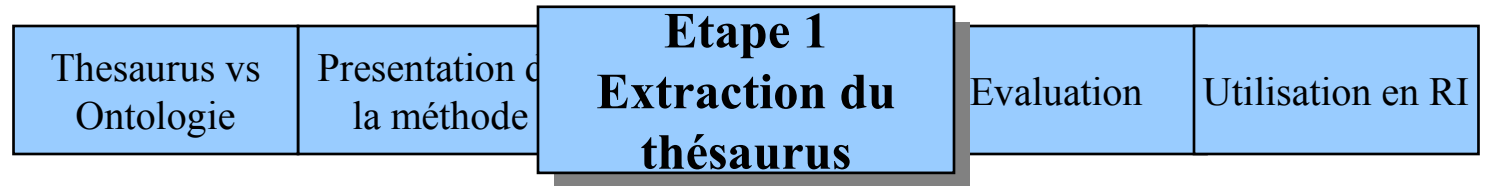
Méthode pour transformer un thesaurus

- Principales étapes à partir de la méthodologie Terminae [Aussenac 2000]
 - Spécification des besoin: indexation de documents (termes du domaine, concepts, relations entre concepts)
 - Choix d'un corpus de référence du domaine: A&A 1995, 2002
 - Analyse linguistique du domaine: Analyse syntaxique du corpus (Syntex) + termes et relations extraits du thesaurus
 - Normalisation (concepts et relations)
 - Formalisation : OWL-Lite [w3c]

Thésaurus v Ontologie	Presentation de la méthode	Étapes	Evaluations	Utilisation en RI
--------------------------	---------------------------------------	--------	-------------	-------------------

3 étapes semi-automatiques

- Extraction des concepts de l'ontologie et de sa structure (relations entre concepts) à partir du thésaurus
- Capture de nouvelles relations entre concepts non présentes dans le thésaurus
- Mise à jour de l'ontologie à partir de nouveaux termes et relations

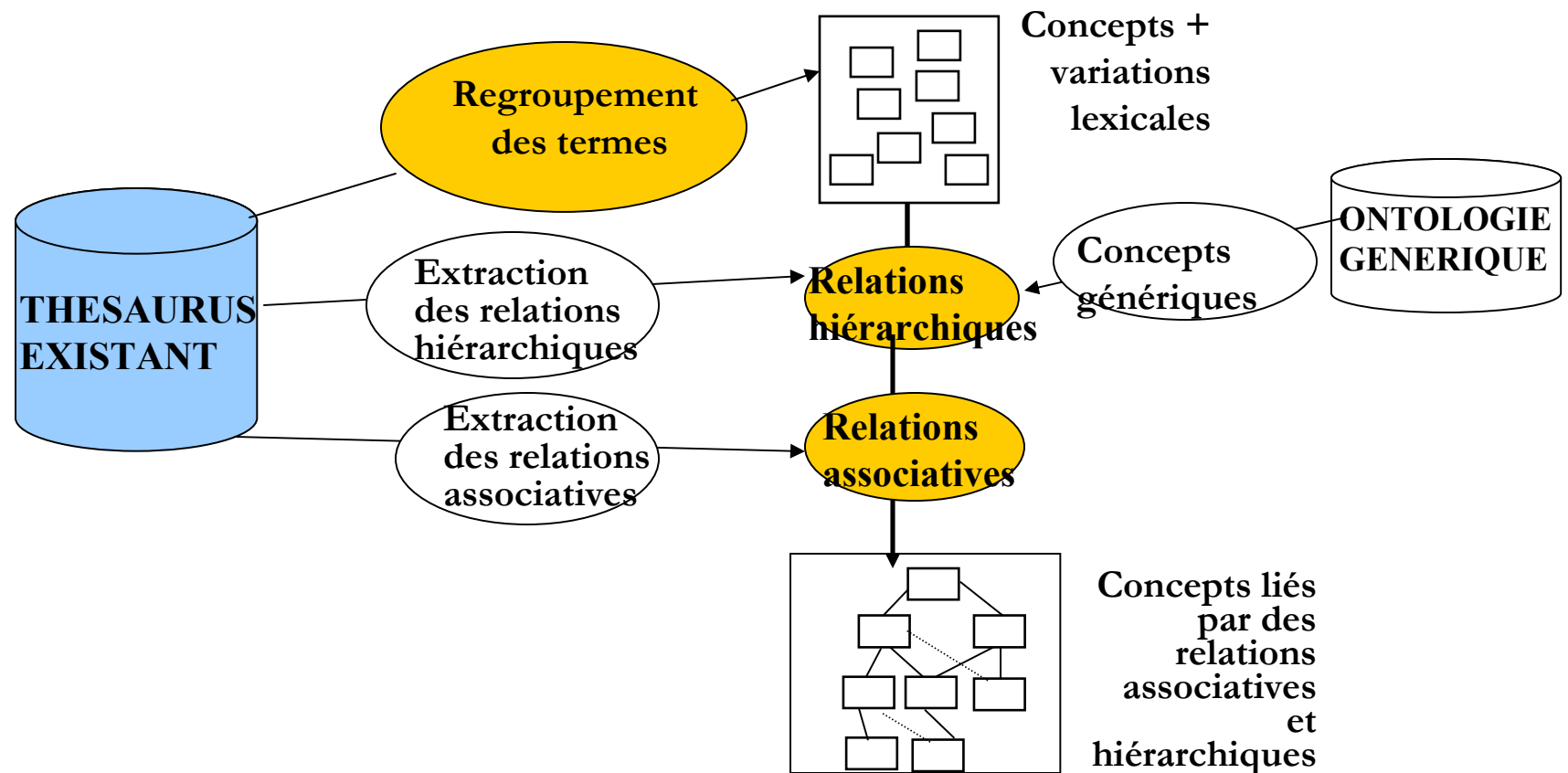


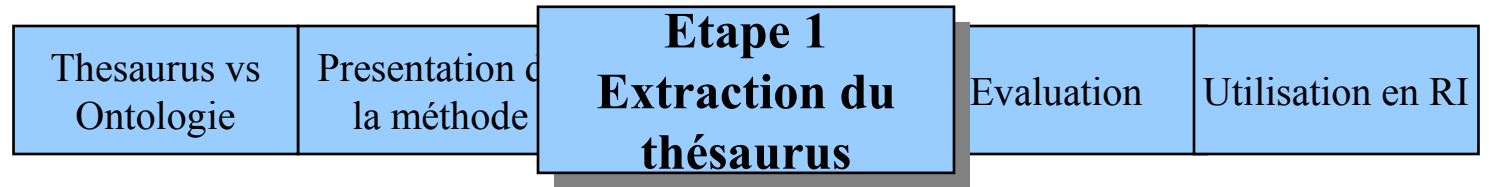
3 étapes semi-automatiques

- Extraction des concepts de l'ontologie et de sa structure (relations entre concepts) à partir du thésaurus
- Capture de nouvelles relations entre concepts non présentes dans le thésaurus
- Mise à jour de l'ontologie à partir de nouveaux termes et relations

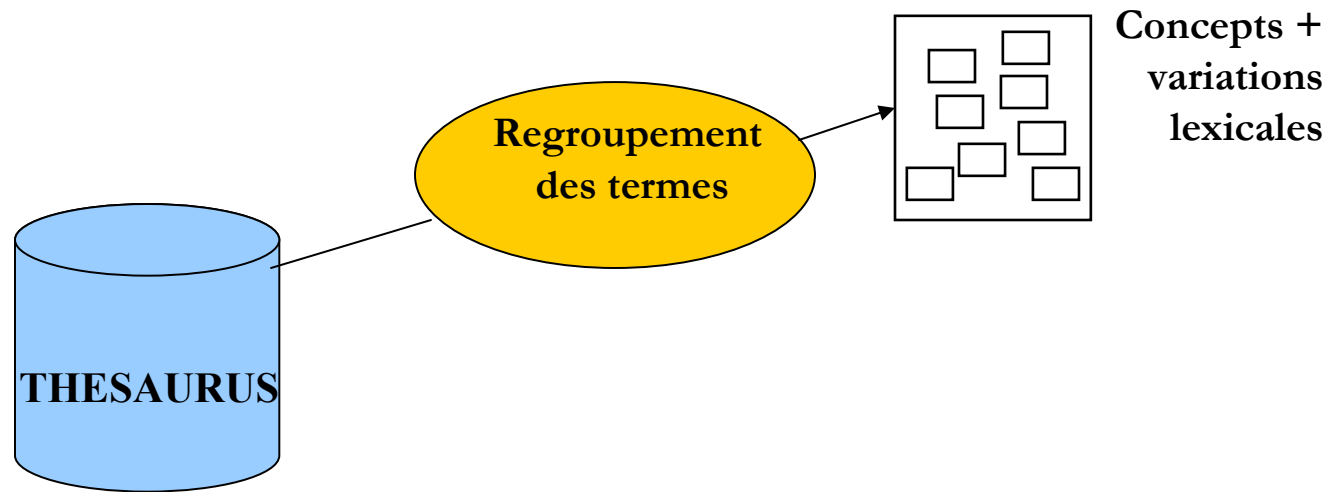
Thesaurus vs Ontologie	Presentation de la méthode	Etape 1 Extraction du thésaurus	Evaluation	Utilisation en RI
---------------------------	-------------------------------	--	------------	-------------------

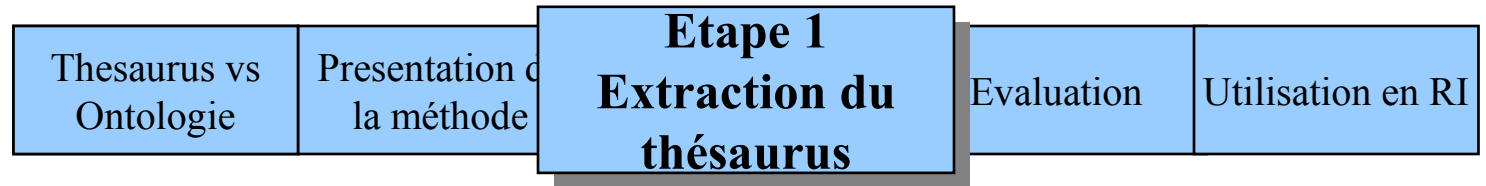
Etape 1: extraction des concepts et relations du thésaurus





Etape 1: Extraction des concepts



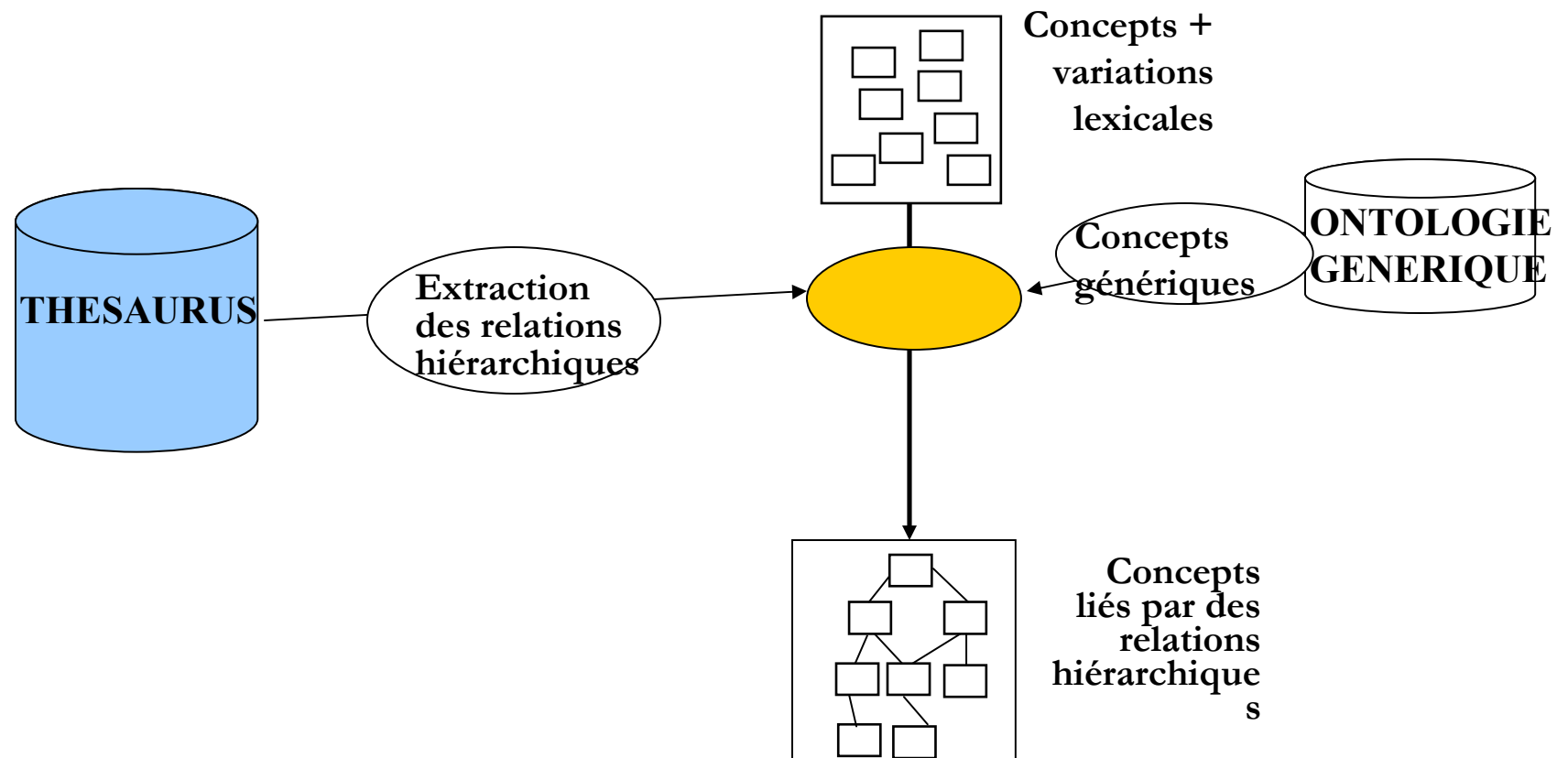


Extraction des concepts

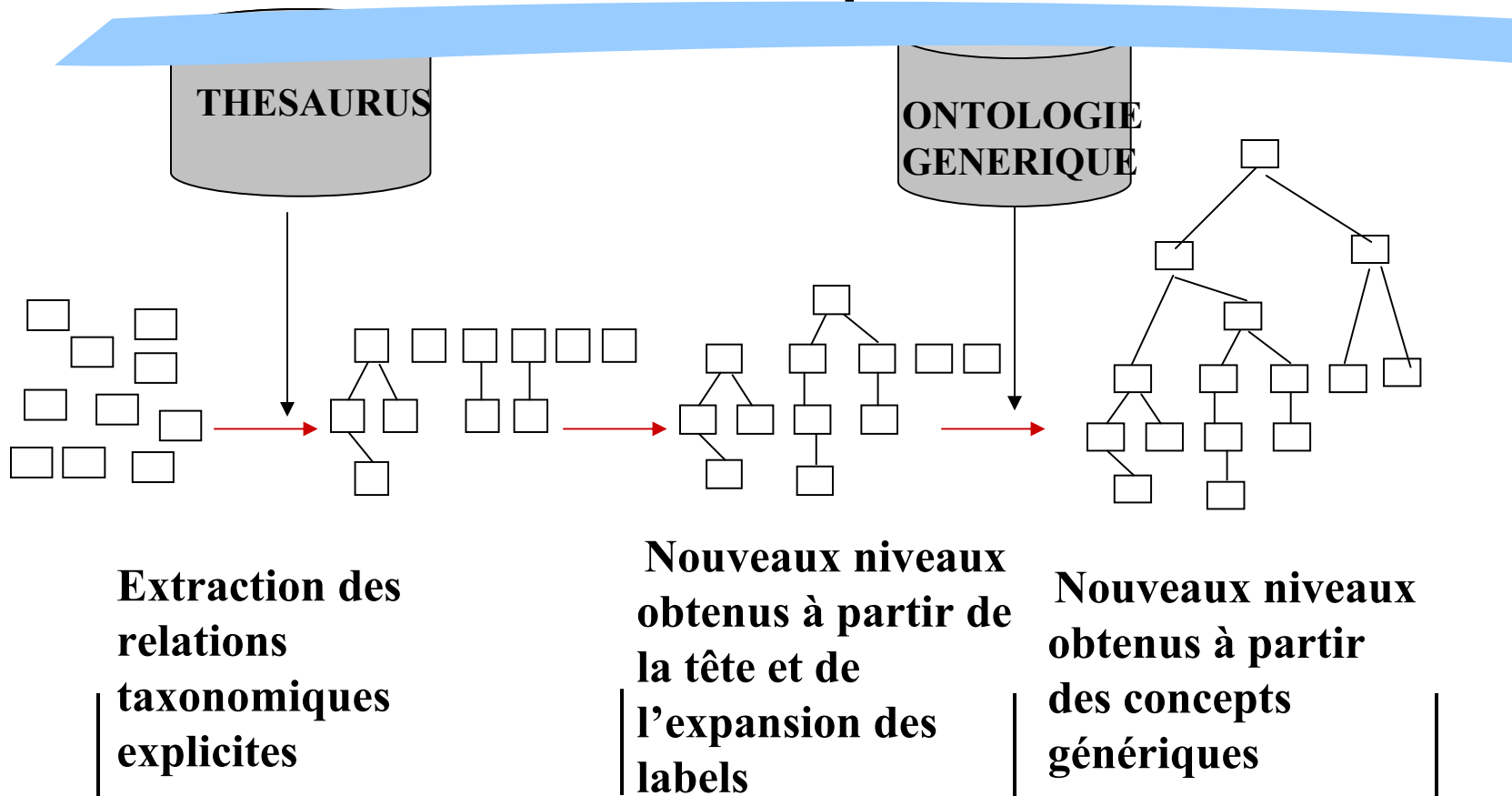
- Conceptualisation du lexique
 - Relations du thesaurus
 - Term1 **USE** term2
 - Term3 **USED FOR** term2
 - Regroupements à partir de la fermeture transitive de ces relations
 - *Exemple :*
 - ELLIPSOIDAL VARIABLE STARS **USE** photometric binary stars
 - ellipsoidal binary stars **USED FOR** ELLIPSOIDAL VARIABLE STARS
 - ⇒ Concept : ELLIPSOIDAL VARIABLE STARS
 - labels : photometric binary stars, ellipsoidal binary stars, ellipsoidal variable stars

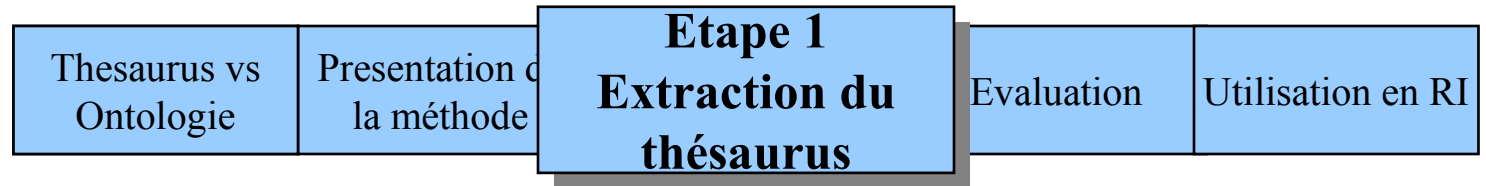
Thesaurus vs Ontologie	Presentation de la méthode	Etape 1 Extraction du thésaurus	Evaluation	Utilisation en RI
---------------------------	-------------------------------	--	------------	-------------------

Extraction des relations hiérarchiques



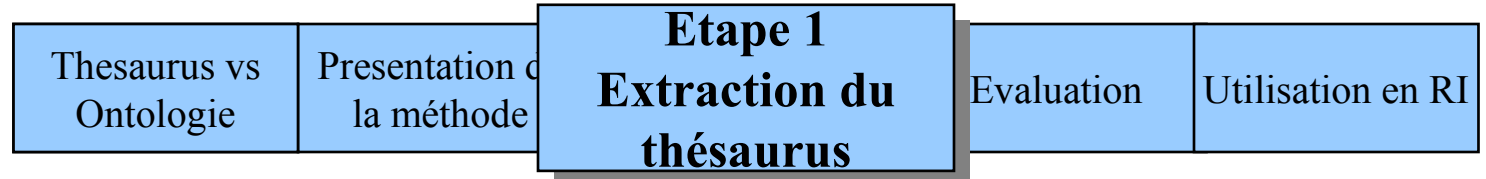
Extraction des relations hiérarchiques





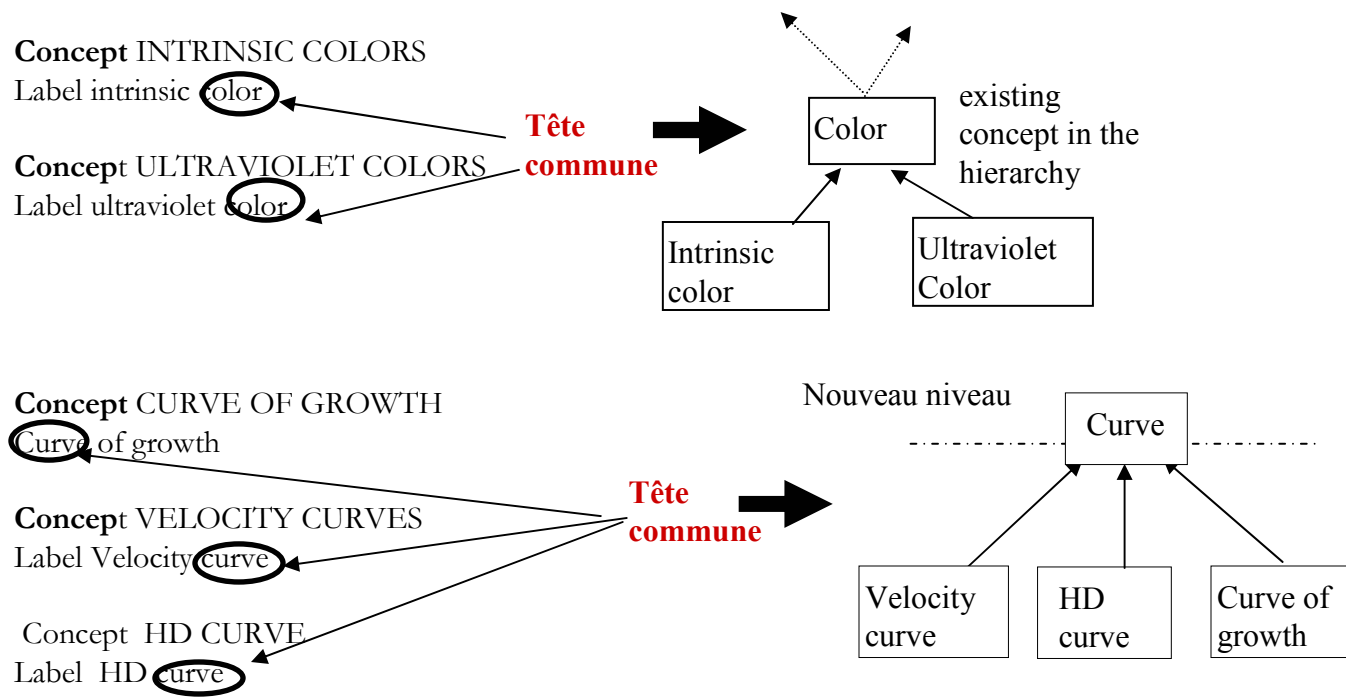
Extraction des relations explicitées dans le thesaurus

- Relations du thesaurus
 - Term1 **Broader Term** Term2
 - Term3 **Narrower Term** Term4
 - Création de liens hiérarchiques entre concepts dont les labels sont liés par ces relations



Nouveaux niveaux hiérarchiques

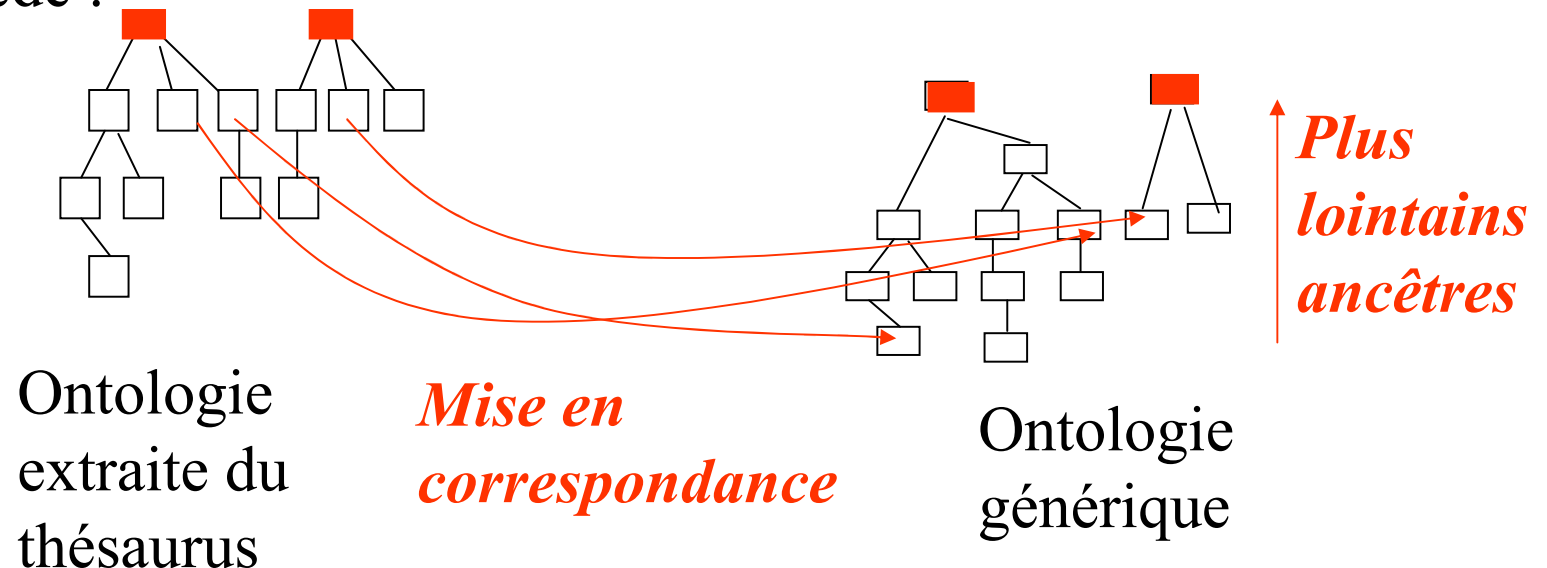
- A partir des labels des concepts



Thesaurus vs Ontologie	Presentation de la méthode	Etape 1 Extraction du thésaurus	Evaluation	Utilisation en RI
---------------------------	-------------------------------	--	------------	-------------------

Extraction des concepts génériques

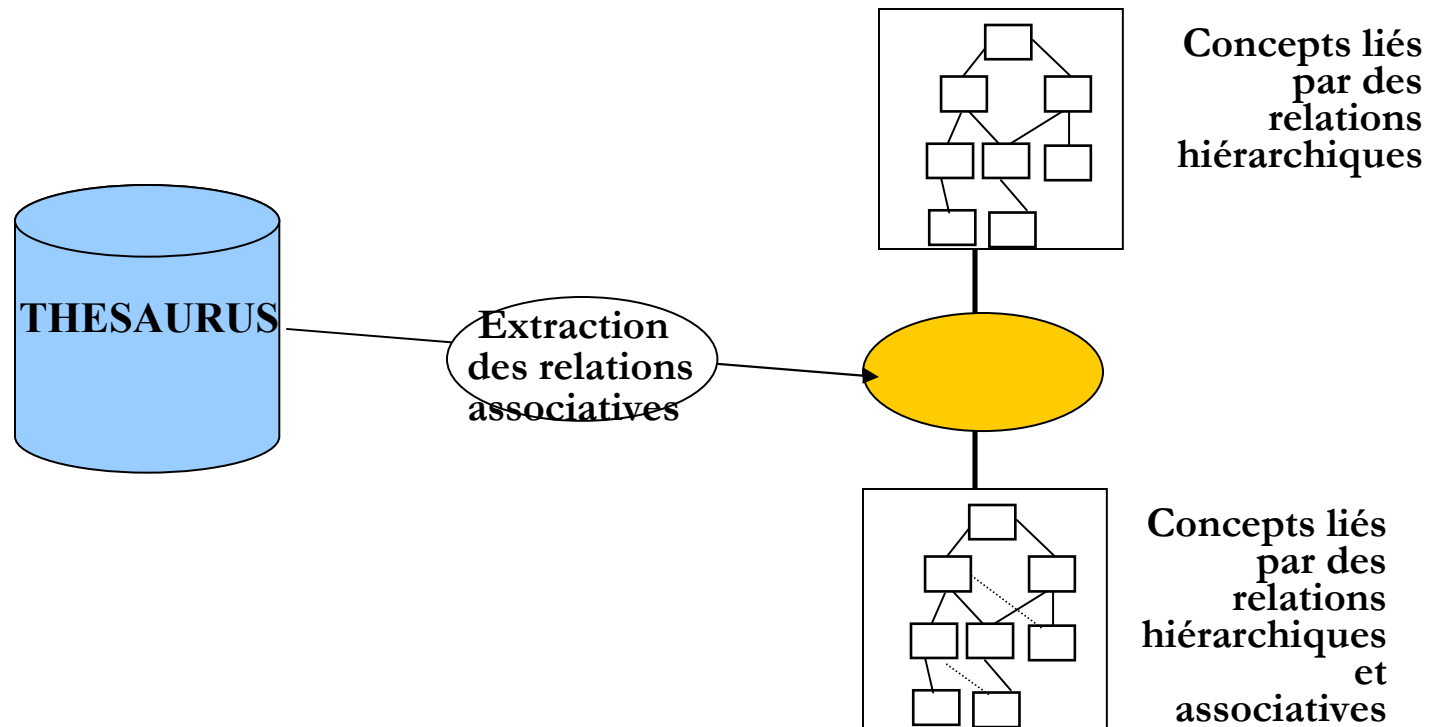
- Concepts génériques = structuration de l'ontologie
- Extraction des concepts génériques à partir d'ontologies génériques (WordNet)
- Procédé :

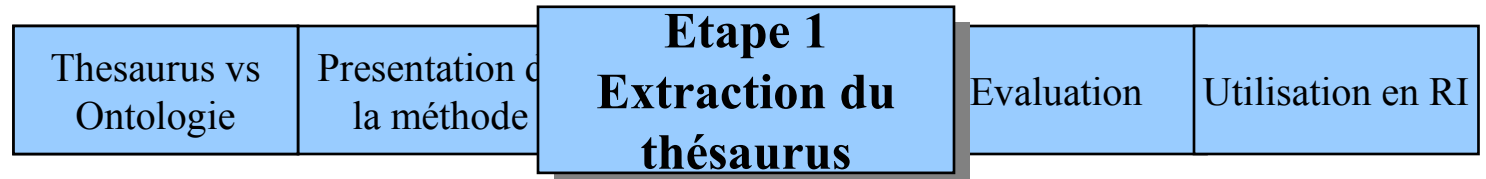


Exemple : *phenomenon, natural object*

Thesaurus vs Ontologie	Presentation de la méthode	Etape 1 Extraction du thésaurus	Evaluation	Utilisation en RI
---------------------------	-------------------------------	--	------------	-------------------

Etape 1: extraction des relations associatives

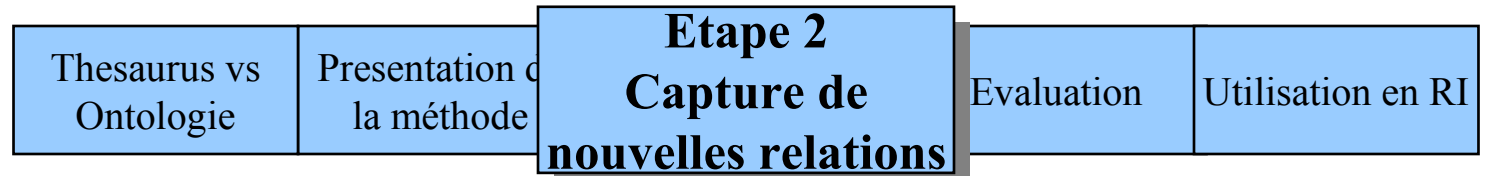




Extraction des relations associatives

- Relations explicitées dans le thésaurus
 - Term1 **RELATED TO** term2
 - Extraction de ces relations entre concepts dont les labels sont liés dans le thésaurus
 - relations vagues et ambigus
- Désambiguïisation de ces relations à partir des concepts génériques

*Exemple : coronagraph **RT** solar corona →
coronagraph **OBSERVES** solar corona*

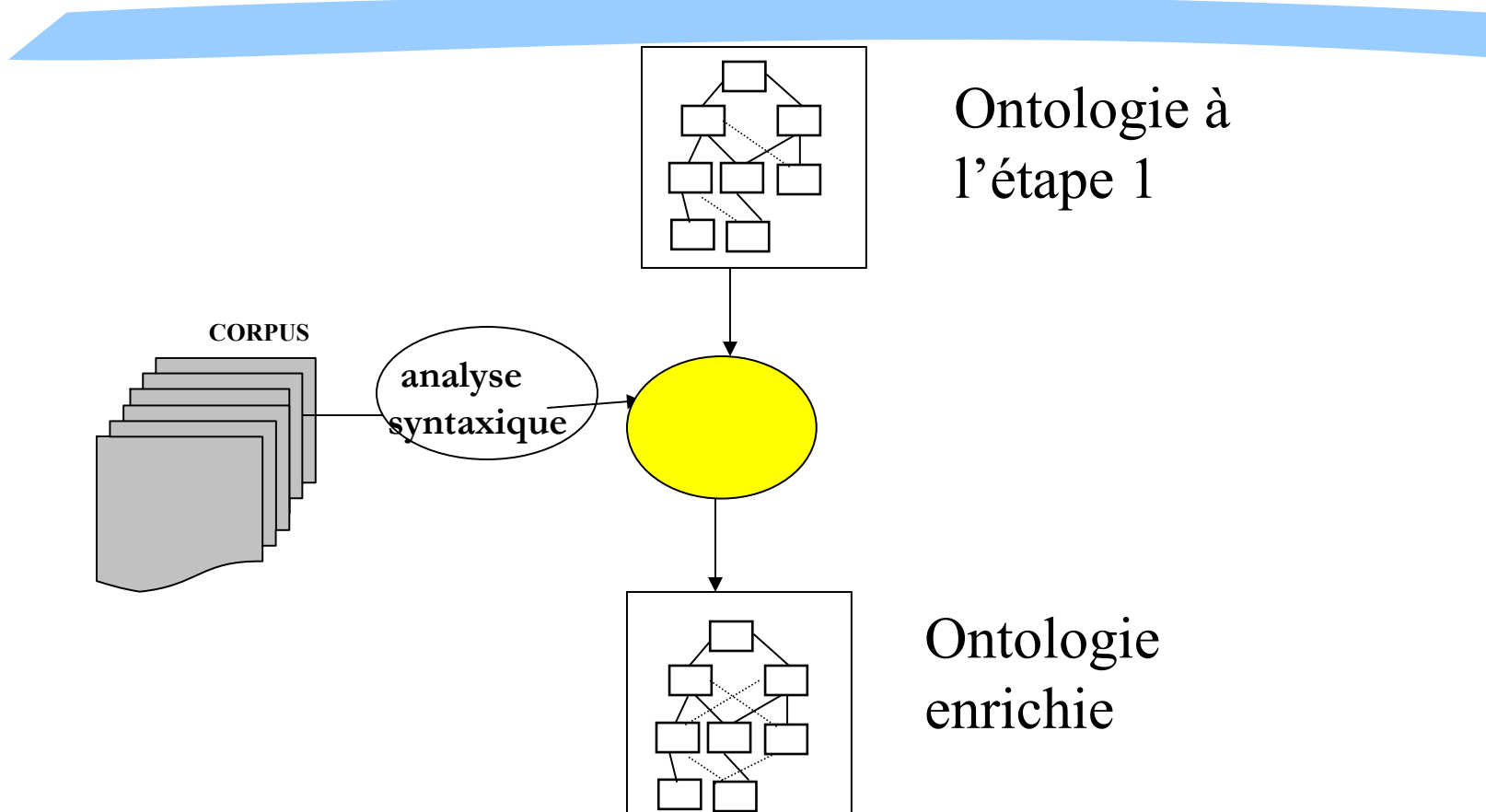


3 étapes semi-automatiques

- Extraction des concepts de l'ontologie et de sa structure (relations entre concepts) à partir du thésaurus
- **Capture de nouvelles relations entre concepts non présentes dans le thésaurus**
- Mise à jour de l'ontologie à partir de nouveaux termes et relations

Thesaurus vs Ontologie	Presentation de la méthode	Etape 2 Capture de nouvelles relations	Evaluation	Utilisation en RI
---------------------------	-------------------------------	---	------------	-------------------

Etape 2 : Capture de nouvelles relations



Thesaurus vs Ontologie	Presentation de la méthode	Etape 2 Capture de nouvelles relations	Evaluation	Utilisation en RI
---------------------------	-------------------------------	---	------------	-------------------

Capture des relations

- Analyse syntaxique du contexte des labels de concepts dans le corpus de référence
- Si un label apparaît dans le contexte d'un label d'un autre concept

⇒ Création d'une nouvelle relation entre les deux concepts

Exemple : « **intensity** » trouvé dans le contexte de « **radial velocity** » (the intensity of radial velocity)

⇒ Création de la relation « **is a property of** » entre les concepts « **radial velocity** » et « **intensity** »

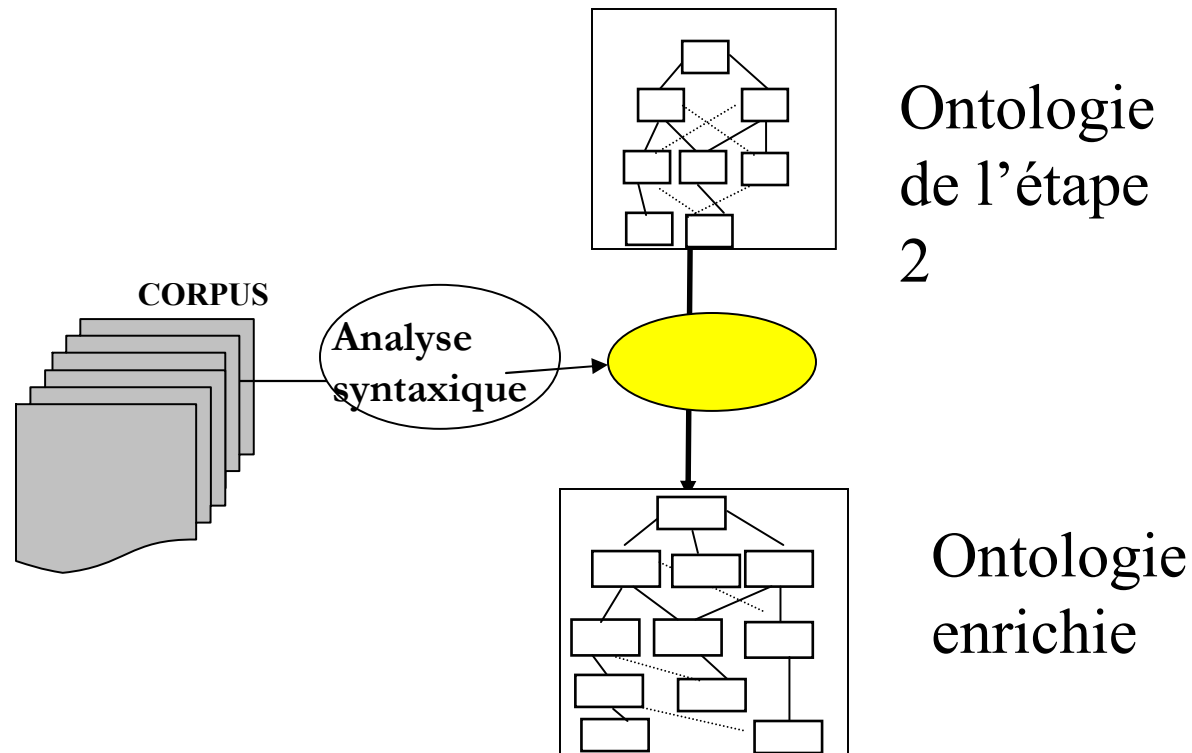


3 étapes semi-automatiques

- Extraction des concepts de l'ontologie et de sa structure (relations entre concepts) à partir du thésaurus
- Capture de nouvelles relations entre concepts non présentes dans le thésaurus
- Mise à jour de l'ontologie à partir de nouveaux termes et relations

Thesaurus vs Ontologie	Presentation de la méthode	Etape 3 Mise à jour	Evaluation	Utilisation en RI
---------------------------	-------------------------------	--------------------------------	------------	-------------------

Etape 3 : mise à jour de l'ontologie



Thesaurus vs
Ontologie

Presentation de
la méthode

Etape 3
Mise à jour

Evaluation

Utilisation en RI

Extraction de nouveaux termes

- Deux fonctions d'extraction

Termes généraux

Termes spécifiques

high resolution
globular cluster
binary system
soft X ray
orbital period
stellar population
power law
absorptance line
line emission
active region

Yarkovsky force
Relativistic gravity
Suprathermal
electron
Halpna knot
Penumbral wave
Mean free path
Integral magnitude
Mixing layer
stellar population



Intégration des termes dans l'ontologie

- 2 approches :
 - Nouveaux sous-concepts de concepts existants
 - Nouvelles relations entre concepts

Thesaurus vs
Ontologie

Presentation de
la méthode

Etape 3
Mise à jour

Evaluation

Utilisation en RI

Intégration des termes dans l'ontologie

- Nouveaux concepts sous-concepts de concepts existants
→ tête d'un terme = label d'un concept existant

Exemple : nouveau terme “soft **X Ray**”

concept existant “**X Ray**”

⇒ création du concept “soft X Ray” sous-
concept de “X Ray”

Thesaurus vs
Ontologie

Presentation de
la méthode

Etape 3
Mise à jour

Evaluation

Utilisation en RI

Intégration des termes dans l'ontologie

- Nouvelles relations associatives entre concepts existants
 - tête et expansion d'un terme = labels de concepts existants

Exemple :

Nouveau terme : **star mass**

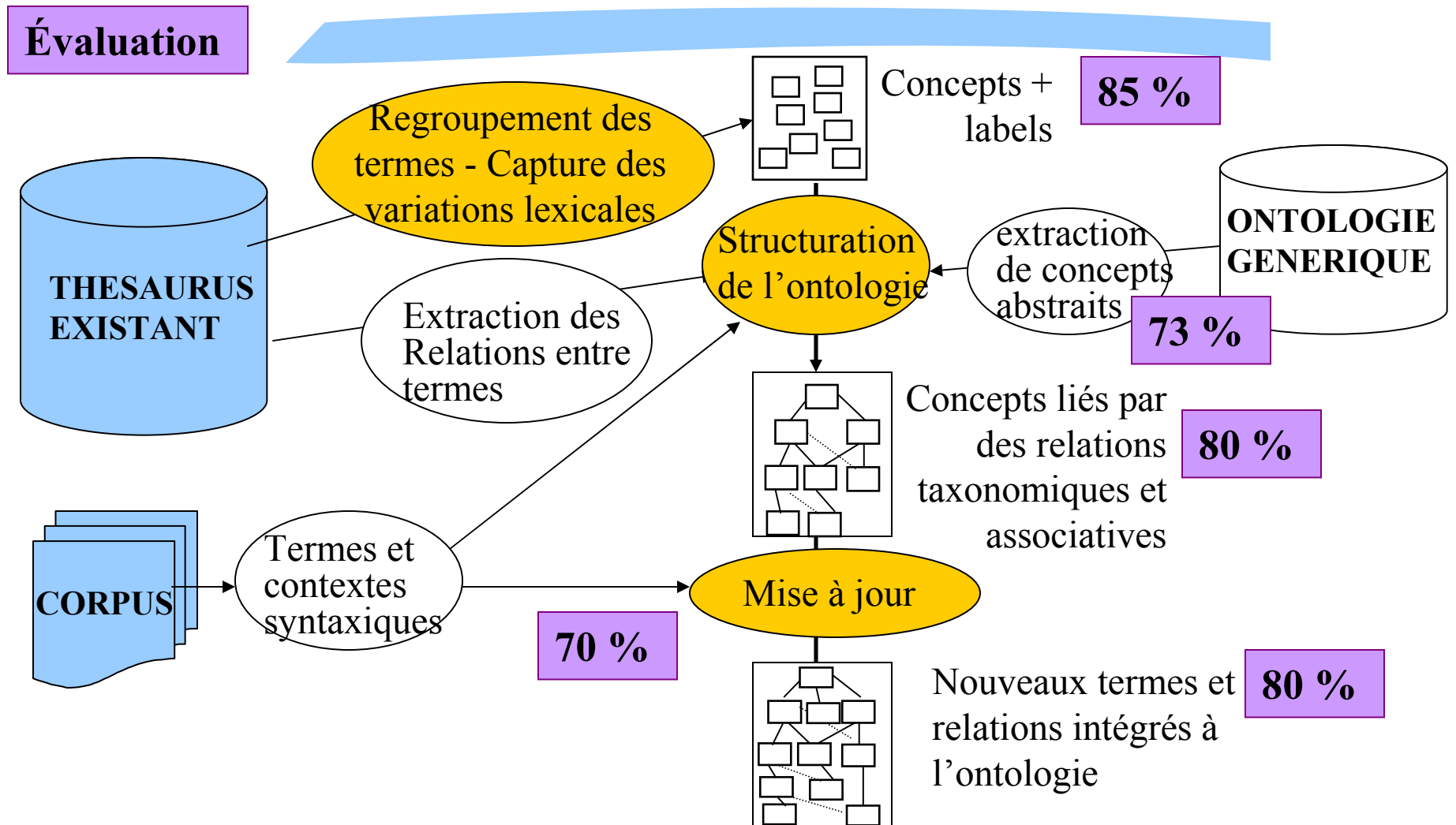
Concepts existants :

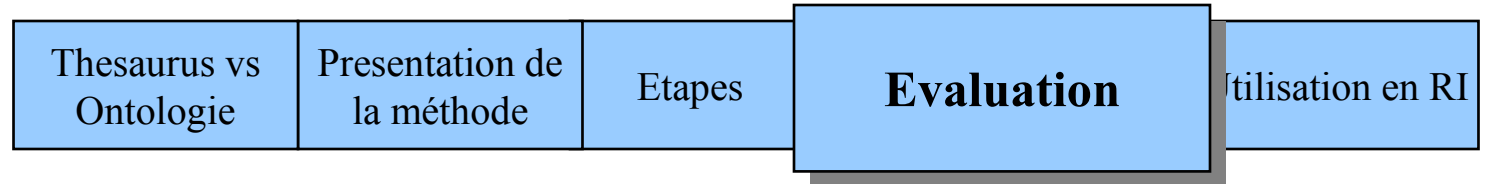
- **star** (natural_object)
- **mass** (property)

⇒ creation de la relation « **has property** » entre les concepts star et mass

Thesaurus vs Ontologie	Presentation de la méthode	Etapes	Evaluation	Utilisation en RI
------------------------	----------------------------	--------	-------------------	-------------------

Evaluation





- Résultats :
 - Ontologie de l’astronomie évaluée à partir de 10% du thésaurus
 - Listes des nouveaux termes
 - Interface permettant de naviguer dans l’ontologie en OWL

Thesaurus vs
Ontologie

Presentation de
la méthode

Etapas

Evaluation

Utilisation en RI

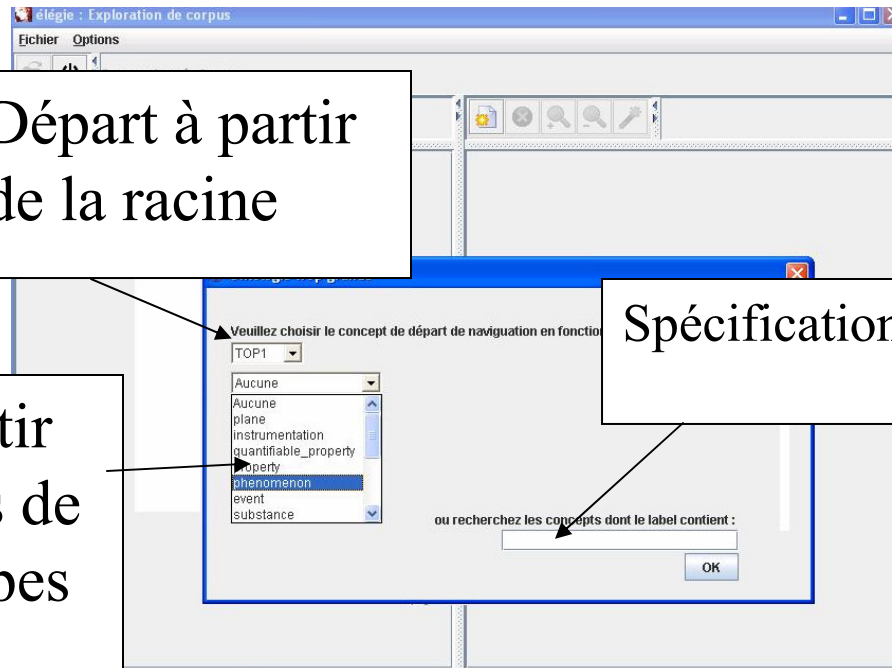
Exploration d'ontologies volumineuses

**Départ de
la
navigation**

Départ à partir
de la racine

Départ à partir
des concepts de
niveau 1 (types
abstrais)

Spécification d'un label



Thesaurus vs
Ontologie

Presentation de
la méthode

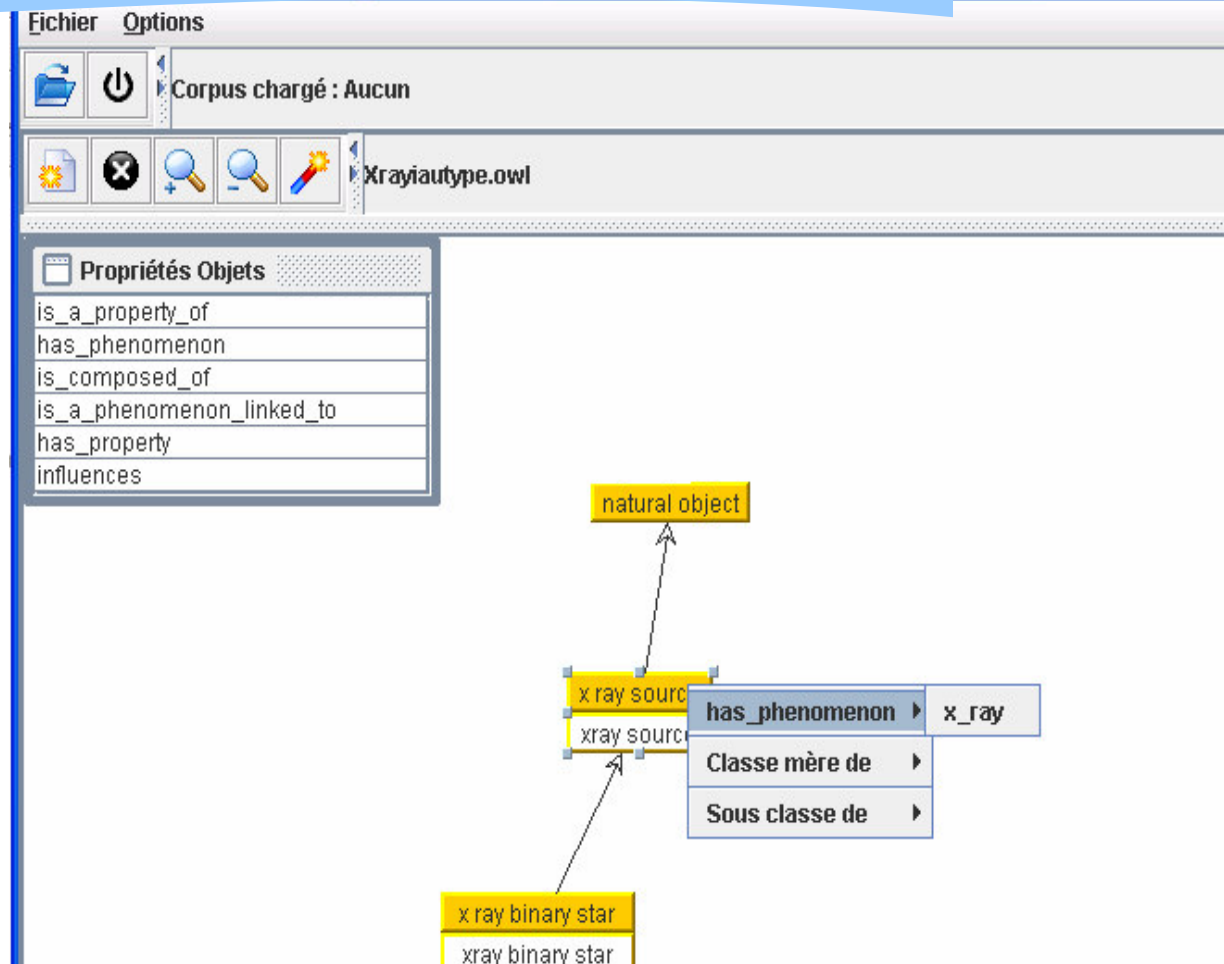
Etapas

Evaluation

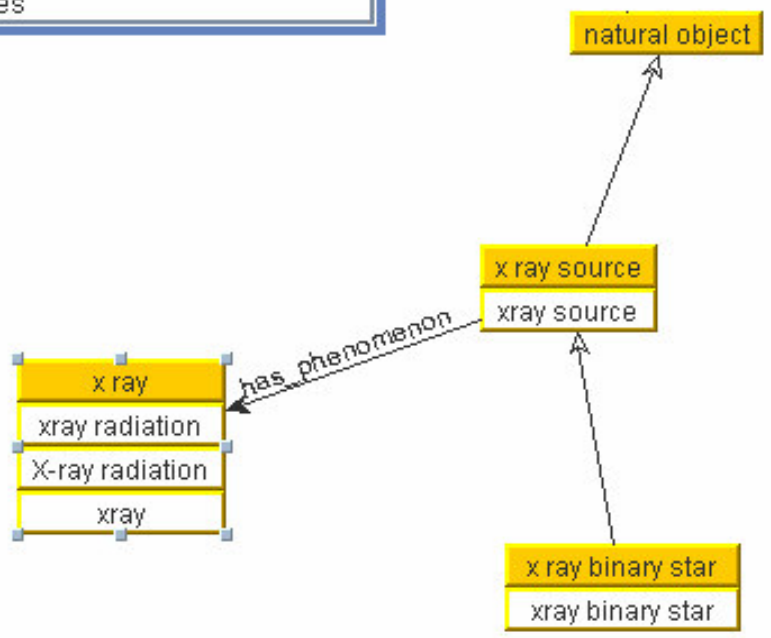
Utilisation en RI

Exploration d'ontologies volumineuses

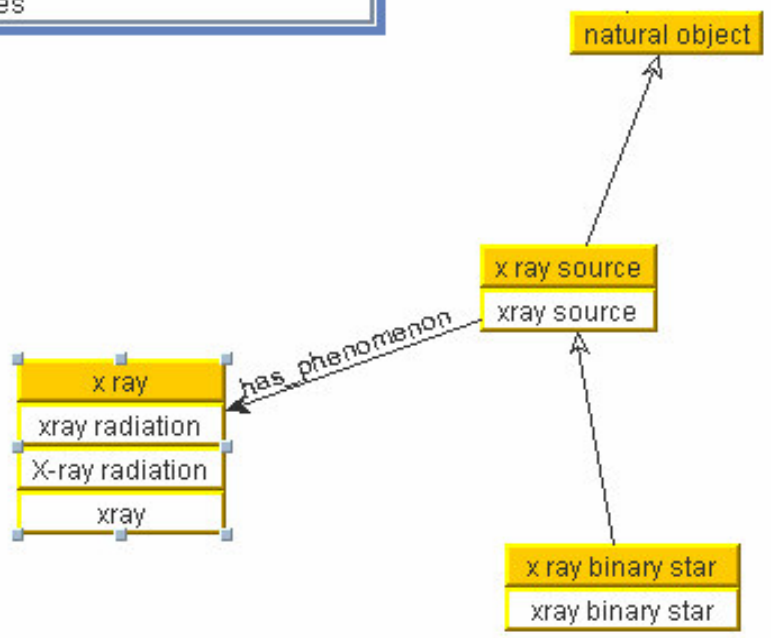
- Visualisation
du contexte
d'un concept



- Propriétés Objets**
- is_a_property_of
 - has_phenomenon
 - is_composed_of
 - is_a_phenomenon_linked_to
 - has_property
 - influences



- Propriétés Objets
- is_a_property_of
 - has_phenomenon
 - is_composed_of
 - is_a_phenomenon_linked_to
 - has_property
 - influences



Thesaurus vs
Ontologie

Presentation de
la méthode

Etapes

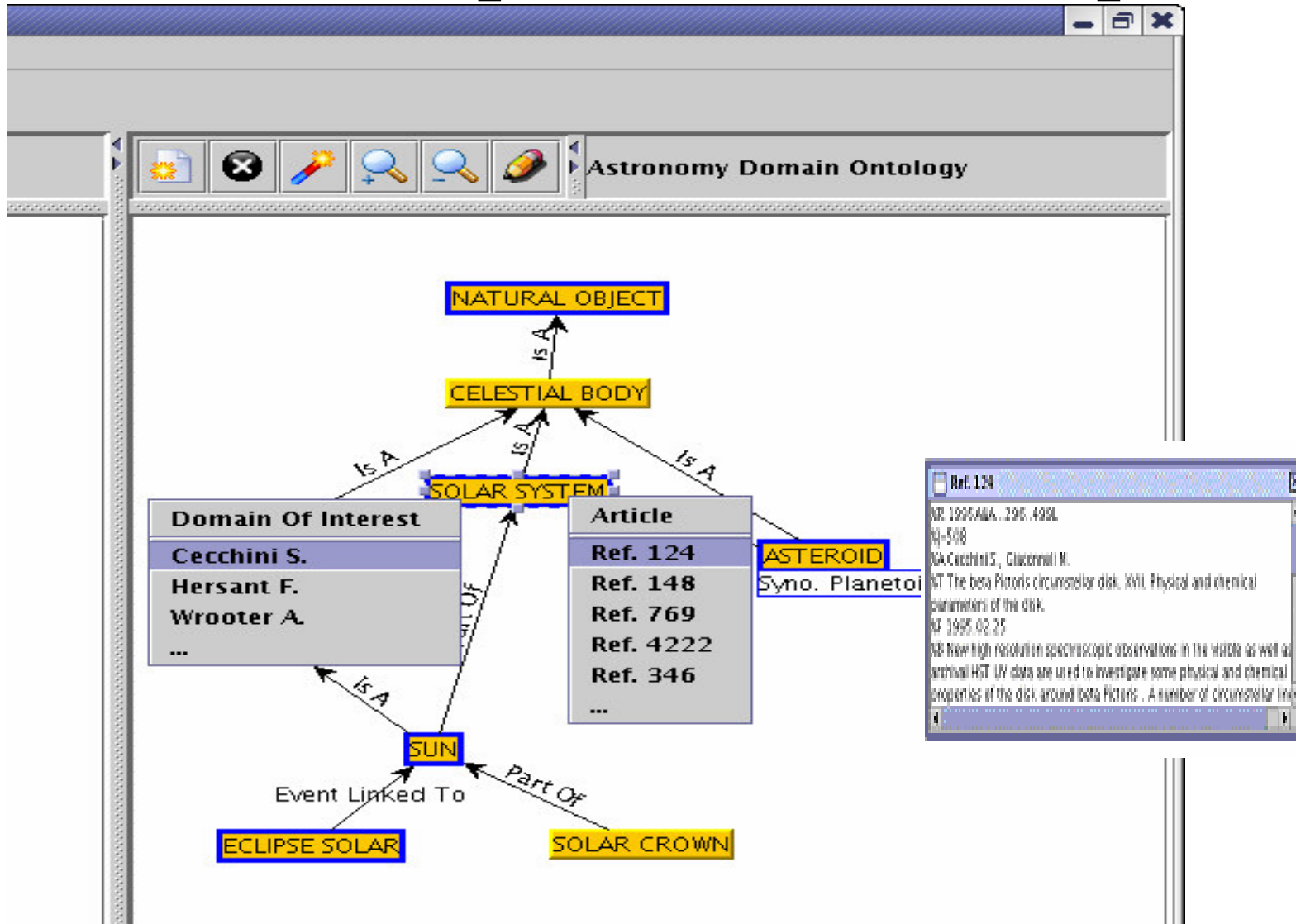
Evaluation

**Utilisation en
RI**

Indexation sémantique à partir de l'ontologie de l'astronomie

- Indexation sémantique [Haav 2001]
 - Non plus uniquement à partir de considérations statistiques
 - A partir des objets du monde référencés
- Descripteurs des granules documentaires choisis à partir des concepts de l'ontologie de l'astronomie
 - Détection des concepts dans les granules
 - Identification des labels référençant des concepts
 - Désambiguïsation des labels
 - Pondération des concepts à partir de leur représentativité des granules
 - Poids statistique
 - Poids sémantique : liens avec les autres concepts

Exploration du corpus



Thesaurus vs Ontologie	Presentation de la méthode	Etapes	Evaluation	Utilisation en RI
---------------------------	-------------------------------	--------	------------	----------------------

Conclusion

- Méthode pour le transformation d'un thésaurus en une ontologie légère
- Évaluation dans le cadre de l'astronomie résultat efficace (+80% des propositions validées)
- Utilisation de l'ontologie créée pour l'indexation et l'exploration de corpus
- Application aux types d'objets

Désambiguïisation des labels retrouvés dans les documents

« *Polarization* varies noticeably with emergent *photon* energy below 40keV, being up to 30% and down to -10% for different angles of view; these variations cover the range of observed *magnitudes*.»

Label ambigu :

polarization concerning wave

polarization concerning charge separation

Labels non
ambigus

→ Nombre d'arcs dans l'ontologie séparant les deux concepts candidats aux concepts identifiés dans la phrase

→ Choix du concept ayant le nombre d'arcs minimum

Désambiguïisation des relations RT

