

# Référentiels et vocabulaires

Pistes méthodologiques pour  
la construction  
d'un vocabulaire contrôlé



Blandine NOUVEL, Centre Camille Jullian (Aix-Marseille Univ, CNRS) & Frantiq  
7 novembre 2023

Méthode

Concepts &  
termes

Relations  
sémantiques

Outils

Normalisation  
& FAIR







# Rappel sur ce qu'est un référentiel

- Types de Référentiels et vocabulaires contrôlés
  - Du dictionnaire au **thésaurus**
  - Organisation +/- structurée
- Objectifs
  - Normaliser une terminologie pour éviter les ambiguïtés du langage naturel (synonymie, polysémie)
  - Partager un langage validé par une communauté
- Usages
  - Enrichir les métadonnées d'une ressource (indexation) en qualifiant son contenu (avec des mots-clés)
  - Faciliter la recherche documentaire en évitant le bruit et le silence (index et filtres)

# Une affaire de longue haleine

## PACTOLS 1987-2019

59727 concepts














- >  Anthroponymes
- >  Chronologie
- >  Lieux
- >  Oeuvres
- >  Peuples
- >  Sujets



## Pactols 2022







### Sujets

12037 concepts

- >  activités
- >  entités immatérielles
- >  entités matérielles
- >  entités nommées
- >  entités sociales collectives
- >  entités spatiales
- >  entités temporelles
- >  matériaux
- >  organismes vivants
- >  processus naturels
- >  rôles
- >  unités géopolitiques
- >  ~[termes dépréciés]

### Lieux

50305 concepts

- >  entités actuelles
- >  entités du passé
- >  espaces aquatiques
- >  espaces terrestres
- >  sites archéologiques
- >  ~[termes dépréciés]

# Par quel bout commencer ?

- Etat des lieux : qu'existe-il comme outil sur mon sujet ?

- Les dictionnaires et lexiques spécialisés...
- Les sites de référencement



- Les vocabulaires numériques disponibles sont-ils normalisés, évolutifs, ouverts, réutilisables (FAIR) ?
- Les adopter, y participer ou les utiliser comme source d'inspiration ?

- Le cas Pactols :

- Recréer un nouveau à partir d'un lexique (mettre à jour et réorganiser)
- Conserver le lien avec le Catalogue Collectif Indexé
- Exploiter les fonctions de Opentheso

# Définir le·s public·s, le·s usage·s

La transformation de Pactols d'un thésaurus documentaire à un référentiel du web pour l'archéologie a eu des incidences sur son contenu et son organisation

## Évolution des ressources à indexer



## De nouveaux utilisateurs



# Constituer une équipe

- Des compétences
  - Compétences sur le web sémantique et la structuration de terminologies
  - Compétences techniques et scientifiques sur le domaine
  - Utilisateurs finaux (métiers IST et Données, archéologues)
- Un calendrier
  - Objectifs tenables
  - Séances de travail régulières
  - Processus itératif de discussion/validation
    - Lexique (ateliers de spécialistes)
    - Structuration (formalisation, définition, relations sémantiques)

# Constituer son réservoir de mots-clés

- Extraction (automatique) de mots-clés dans un corpus (publications scientifiques) ou un ou des vocabulaires existants
- Viser l'exhaustivité terminologique du domaine?
  - Lien avec un fonds documentaire = pas de concept s'il n'y a pas au moins une ressource correspondante dans le catalogue
- Combien de termes ? de 100 à +10 000 aines

« La collecte des termes ne peut être efficace que si elle prend aussi en compte le vocabulaire des utilisateurs. » (M. Hudon, 2013)



# Ré-organiser un vocabulaire

- Top-down:
  - Définir les branches principales, puis les sous-branches
  - Y positionner les concepts
- Down-top :
  - Rassembler les concepts par « famille »
  - Identifier des thématiques, de la plus précise à la plus large
- Un mélange des 2...

# Les concepts & les termes

- Un terme représente un concept
  - réception de l'Antiquité ; vente d'esclave ; nucléus
- Se méfier des termes apparemment sans difficultés
  - Préférer des concepts plébiscités par les spécialistes
    - crémation plutôt qu'incinération
    - inhumation ≠ sépulture à inhumation
  - Mais quand la terminologie n'est pas [encore] stabilisée ou que le terme a des acceptions différentes selon les époques...
    - sarcophage vs cercueil vs coffre funéraire
    - charnier médiéval ≠ charnier moderne
- Exploiter les variantes et les notes d'application

# L'écriture du concept

- Un substantif
  - simple ou composé
  - au singulier, sauf cas particulier (catacombes)
  - ni adjectif : pas **gravé**, mais galet gravé, plaquette gravée, signe gravé
  - ni verbe : pas **fouiller**, mais fouille, fouille subaquatique, fouille programmée
- Justesse des libellés et de leurs variantes
  - skos:prefLabel = la forme choisie/préférée
  - skos:altLabel = formes rejetées (synonymes, équivalent (ou antonymes) dans le thésaurus, formes obsolètes)
  - skos:hiddenAltLabel = formes erronées

bouton de préhension // bouton (récipient) vs bouton de costume // bouton (vêtement)

instrument de musique // musique (instrument)

oppidum plutôt qu'habitat fortifié de hauteur protohistorique

archéologie préventive // archéologie de sauvetage

# Définition et notes d'application

- Norme ISO 1087: 2019

## Terminology work and terminology science — Vocabulary

Travaux terminologiques et science de la terminologie — Vocabulaire

### \* Les critères de qualité

Une définition...

- **Ne décrit qu'un seul concept**
- **Concise** : tout ce qui n'est pas indispensable à la compréhension est à supprimer
- Est de forme affirmative
- Mentionne le concept supérieur
- Réutilise des termes définis par ailleurs, si possible dans un dictionnaire de langue générale ou dans le même référentiel
- Tient en une seule phrase

INRAE DipS  
Rédiger une définition  
2021 / Vocabulaires Ouverts @INRAE



### \* Bonnes pratiques

- La définition d'un nom débute par un nom, celle d'un verbe par un verbe
- La définition ne doit commencer ni par un article, ni par un adjectif démonstratif, ni par un pronom démonstratif.
- La définition ne doit pas commencer par "espèce de", "type de" ... (à quelques exceptions près)
- Une définition ne doit pas être introduite par le terme à définir ni comprendre ce terme (ou un de ses synonymes) ou un terme de la même famille.

INRAE DipS  
Rédiger une définition  
2021 / Vocabulaires Ouverts @INRAE



p 5

## Péloponnèse dans Pactols

### Péloponnèse (périphérie)

- **Définition** : Région administrative de Grèce ou périphérie comprenant cinq des sept districts régionaux de la péninsule du Péloponnèse.
- **Note d'application** : À différencier du concept "péninsule du Péloponnèse" qui est un espace géographique et qui ne couvre pas la même étendue que la périphérie du Péloponnèse.

### Péninsule du Péloponnèse

- **Définition** : Espace géographique désignant la péninsule du Péloponnèse, bordée par la mer Ionienne et la mer Égée et reliée à la Grèce continentale par l'isthme de Corinthe.
- **Note d'application** : À différencier du concept "Péloponnèse (périphérie)" qui est une région administrative de Grèce et qui ne couvre pas la même étendue géographique que la péninsule du Péloponnèse.

# Structurer le thésaurus

- Norme ISO25964 -1 et 2, *Information and documentation – Thesauri and interoperability with other vocabularies*
- Langage SKOS, *Simple Knowledge Organisation System* (Modèle FAIR développé par le W3C) [Version FR](#)
  - des objets : thésaurus, concepts, groupes
  - des propriétés : les relations sémantiques
- Bonnes pratiques pour structurer un thésaurus (<https://opentheso.hypotheses.org/67>)

# Branches thématiques vs organisation ontologique

2 ARCHIVES DE FRANCE. THESAURUS POUR L'INDEXATION DES ARCHIVES LOCALES

**SOMMAIRE**

1. Administration	7. Extérieur
1.1 Droit public	7.1 Défense du territoire
1.2 Administration générale	7.2 Guerre
1.3 Finances publiques	7.3 Relations internationales
1.4 Fiscalité	
1.5 Police	
1.6 Protection civile	
1.7 Régime seigneurial	
2. Agriculture	8. Justice
2.1 Economie rurale	8.1 Condition pénitentiaire
2.2 Forêt	8.2 Justice civile
2.3 Production agricole	8.3 Justice pénale
	8.4 Décision de justice
	8.5 Organisation judiciaire
3. Communications	9. Opinion
3.1 Messagerie	9.1 Mouvement d'idées
3.2 Transport	9.2 Election
	9.3 Vie politique
	9.4 Vie publique
	9.5 Vie religieuse
	9.6 Croyances et sciences parallèles
4. Economie	10. Société
4.1 Action économique	10.1 Condition des personnes et des biens
4.2 Commerce	10.2 Emploi
4.3 Entreprise	10.3 Population
4.4 Industrie	10.4 Protection sociale
4.5 Energie	10.5 Santé
	10.6 Travail
5. Education	11. Temps libre et sociabilité
5.1 Enseignement	11.1 Culture
5.2 Organisation scolaire	11.2 Loisir
5.3 Recherche scientifique	11.3 Tourisme
5.4 Vie scolaire	11.4 Vie quotidienne
6. Equipement	
6.1 Environnement	
6.2 Immobilier	
6.3 Urbanisme	
6.4 Voie de communication	

[https://francearchives.gouv.fr/file/efb9242a3d6d4e5557940d65a95d4f830dc1a30b/static\\_5372.pdf](https://francearchives.gouv.fr/file/efb9242a3d6d4e5557940d65a95d4f830dc1a30b/static_5372.pdf)

Backbone Thesaurus		
Alphabetical	Hierarchy	Groups
+ 000001 activities		
- 000010 disciplines		
- 000013 functions		
- 000011 human interactions		
- 000012 intentional destructions		
+ 000006 conceptual objects		
- 000024 concepts		
- 000023 methods		
- 000022 propositional objects		
- 000021 symbolic objects		
+ 0000049 geometric extents		
- 0000052 3d-volumes		
- 0000053 linear extents		
- 0000050 points		
- 0000051 surface areas		
- 000009 geopolitical units		
- 000007 groups and collectivities		
+ 000004 material things		
- 000018 built environment		
- 000017 mobile objects		
- 000019 physical features		
- 000020 structural parts of material objects		
- 000003 materials		
+ 000002 natural processes		
- 000016 geneeses		
- 000015 natural disasters		
+ 000008 roles		
- 000025 offices		
- 000026 roles of interpersonal relations		
- 000005 types of epochs		

<https://vocabs.dariah.eu/bbt/en/?clang=fr>

# La subordination

- Caractéristique du thésaurus, représenté en arbre hiérarchique  
terme générique (TG) / terme spécifique (TS & TSI)
- Subordination en ISO :
  - Générique = IS A  
plan d'eau *est un* lac / lac *est un* plan d'eau
  - Instantielle = IS A (déclinaison)  
lac Érié *est une instance de* (iso-thes:narrowerInstantial) lac
  - Partitive = IS PART OF  
lac Érié *est une partie de* (iso-thes:narrowerPartitive) Grands lacs
- SKOS définit la direction de la subordination:
  - skos:broader / skos:narrower  
plan d'eau <skos:narrower = lac > / lac <skos:broader = plan d'eau >



[https://upload.wikimedia.org/wikipedia/commons/c/c3/Canada\\_relief\\_map\\_2.svg](https://upload.wikimedia.org/wikipedia/commons/c/c3/Canada_relief_map_2.svg)

# Les niveaux de profondeur

- La norme :
  - 4 niveaux ou plus
- L'usage et la visibilité,
  - Un compromis entre le nombre de niveaux et le nombre de branches
- La nécessité :
  - Respecter les besoins du public (spécialiste)

- > 📁 activités
- > 📁 entités immatérielles
- ✓ 📁 entités matérielles
  - > 📁 caractéristiques physiques
  - > 📁 environnements bâtis
  - > 📁 objets mobiles
  - ✓ 📁 parties structurales d'entités matérielles
    - ✓ 📁 élément de préhension
      - ✓ 📁 anse
        - 📄 anse en flûte de pan
        - 📄 anse surbaissée
      - ✓ 📁 bouton (récipient)
        - 📄 bouton perforé
        - 📄 languette perforée
      - ✓ 📁 manche
        - 📄 hampe
        - 📄 préhampe
      - 📄 oreille (préhension)
      - 📄 poignée
      - 📄 queue (récipient)



# La polyhiérarchie

- Un concept a plusieurs génériques

France

Bourgogne-Franche-Comté

Département de la Nièvre

Glux-en-Glenne

**Bibracte**

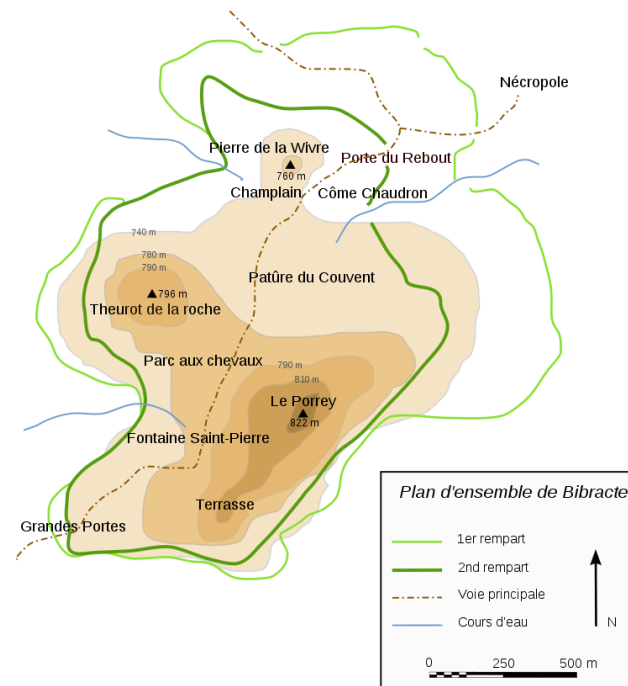
Bibracte EPCC

Larochemillay

**Bibracte**

Saint-Léger-sous-Beuvray

**Bibracte**



La grande *domus* de Bibracte

Photo et plan du site <https://www.bibracte.fr>

# Les relations d'association

Équivalent à un « Voir aussi »

- Lier des termes de niveau équivalent

entités matérielles > objets mobiles > **matériel de broyage** > meule > meule va-et-vient

TA

entités immatérielles > méthodes > technique de fabrication > **broyage (technique)**

- Une solution alternative à la polyhiérarchie ou à des relations génériques transitives fautives

entités actuelles > ... > France > ... > Département de la Nièvre > Glux-en-Glenne > **Bibracte**

entités actuelles > ... > France > ... > Département de la Nièvre > Larochemillay > **Bibracte**

entités actuelles > ... > France > ... > Département de la Nièvre > Saint-Léger-sous-Beuvray > **Bibracte**

TA

espaces terrestres > ... > Massif central > Morvan > **mont Beuvray**

**ET PAS (dans l'organisation Pactols)**

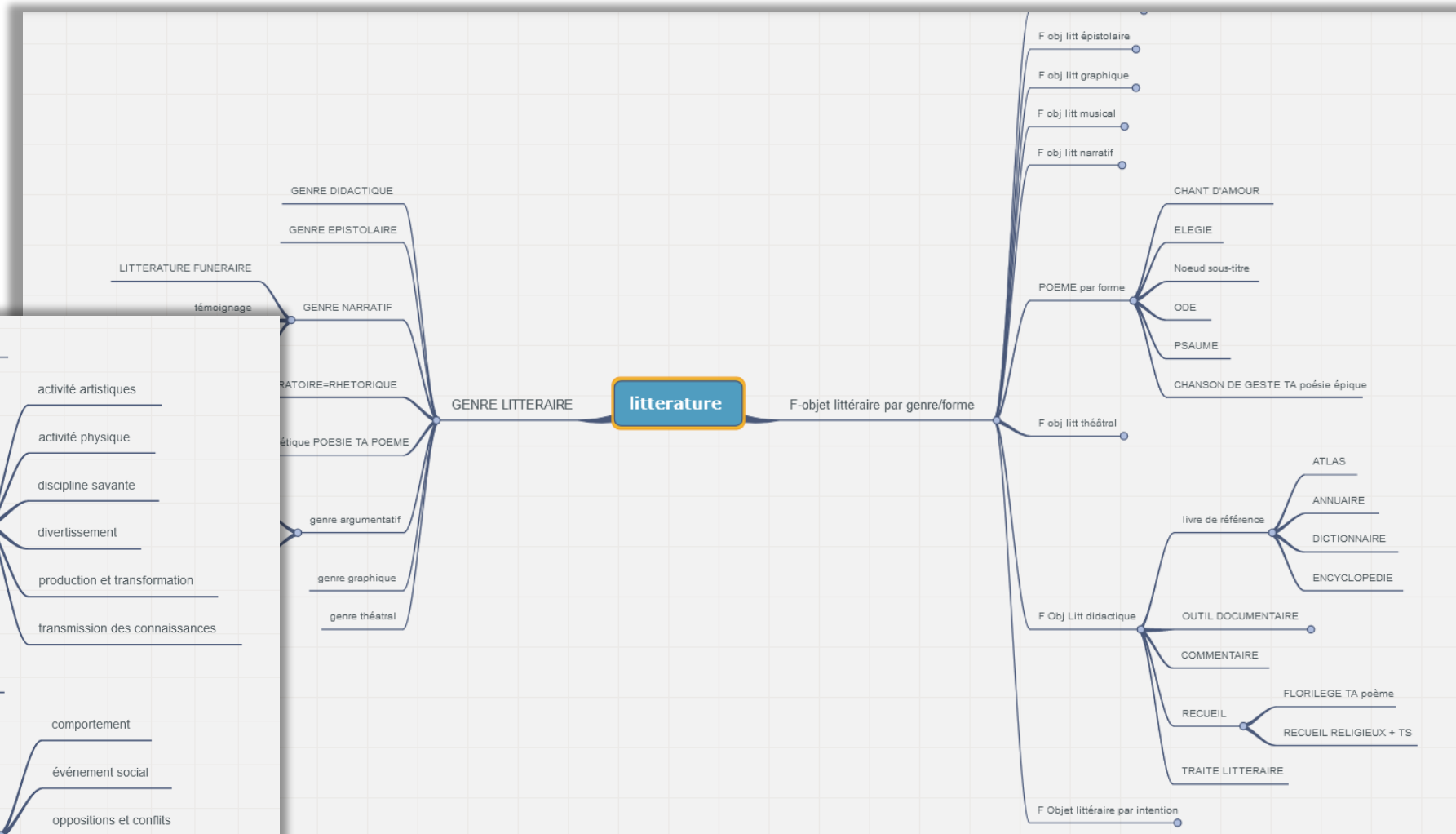
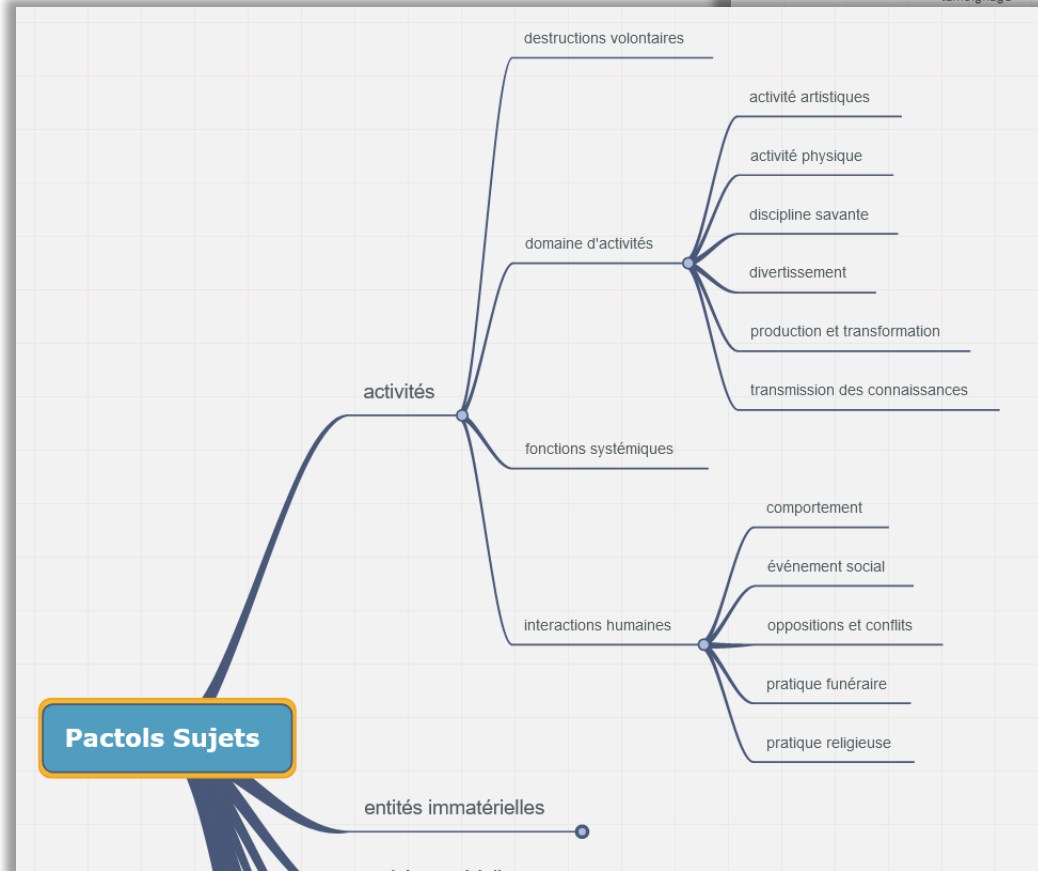
entités actuelles > ... > France > ... > Département de la Nièvre > Glux-en-Glenne > **Bibracte**

entités actuelles > ... > France > ... > Département de la Nièvre > Larochemillay > **Bibracte**

entités actuelles > ... > France > ... > Département de la Nièvre > Saint-Léger-sous-Beuvray > **Bibracte**

espaces terrestres > ... > Massif central > Morvan > **mont Beuvray**

# Préparer-1



# Préparer-2

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	Code couleur	bbs : noir gras ; candidat ventilé : vert ; nouveau terme : rouge ; pactols : normal	à débattre : fête/manifestation ?																
2	facettes BBT	tt BBT										TA	EM	autre TG	Définition	commentaire			
3	F - activités																		
4	F - activités	TT - activités																	
5	F - activités		TT - disciplines																
6				construction de phénomènes															
7				histoire du sport									sport (histoire)			plutôt compréhension de phénomènes non ? tout ce qui est spécifique			
8				compréhension de phénomènes															
9				fourniture de savoir et/ou d'expertise															

9																	
10																	
11																	
12																	
13																	
14																	
15																	
16																	
17																	
18																	
19																	
20																	
21																	
22																	
23																	
24																	
25																	
26																	
27																	
28																	
29																	
30																	
31	F - activités		TT - interactions humaines														
32			TS événements sociaux														
33			fête														
34			commémoration														

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	
1	MOT CLE / CANDIDAT	STATUT	ID INTERNE	POIDS	FRAN	ACTION	EN FAVEUR	DEFINITION	TRADUCTION	SYN	TERME ASSOC.	POSITION (BBT 1)	POSITION (BBT 2)	POSITION (BBT 3)	TG	COMMENTAIRE	PROPOSITION
2	agriculture	concept	13130									Activités	Disciplines	construction d'objets matériels et installations			
3	exploitation agricole	concept	14676									entités matérielles	complexes bâtis				
4	plantation	concept	178384					Vaste domaine fondé, en général, par des colons, autour d'une exploitation où l'on pratique la monoculture à usage alimentaire ou industriel.			esclavagisme	entités matérielles	complexes bâtis		exploitation agricole		
5	fertilité	concept	18262									processus naturels	génèses				
6	produit agricole	concept	16416									entités matérielles	objets mobiles	objets mobiles par matériaux ?		en attendant la révision de la branche Objets	
7	saison agricole	concept	13131									objets conceptuels	méthodes		unité de temps		
8	labour	concept	15381									objets conceptuels	méthodes	techniques	technique agricole		
9	récolte	concept	16740									objets conceptuels	méthodes	techniques	technique agricole		
10	battage	concept	13557									objets conceptuels	méthodes	techniques	technique agricole		
11	moisson	concept	15699									objets conceptuels	méthodes	techniques	technique agricole		
12	vannage	concept	177758									objets conceptuels	méthodes	techniques	technique agricole		
13	vendange	concept	17485									objets conceptuels	méthodes	techniques	technique agricole		
14	structure agraire	concept	17056					Aménagement du terroir caractérisé par la forme, les dimensions, la disposition des parcelles constituant l'espace exploité par un groupe d'agriculteurs ; au sens large, ensemble des relations entre l'homme et le terroir, qui se rapportent à l'exploitation du terroir, à la structure foncière, à la taille des exploitations agricoles, et qui se traduisent par des paysages ruraux caractéristiques. L'ARROUSSE.			terroir, agriculture	Activités	Fonctions				
15	colonat	concept	13966									Activités	Fonctions		structure agraire		
16	fermage	concept	14728									Activités	Fonctions		structure agraire		
17	technique agricole	concept	13132									objets conceptuels	méthodes	techniques			
18	amendement	concept	13168									objets conceptuels	méthodes	techniques	Technique agricole		
19	assolement	concept	13434									objets conceptuels	méthodes	techniques	Technique agricole		
20	culture en terrasse	concept	14232									objets conceptuels	méthodes	techniques	Technique agricole		
21	culture extensive	concept	14233									objets conceptuels	méthodes	techniques	Technique agricole		

+

≡

BBT

définitions ajoutées

dépréciés sans renvoi

propositions déplacements

économie (faits)

12 Savoir

Re-Indexation

agriculture élevage éco de subsistance

transport

Vie administrative

# De CSV à SKOS

Skos Play! &gt; « convertir une feuille de calcul Excel en fichier SKOS »

SKOS Play !

Accueil

▶ Play!

📄 Convertir

🔍 Tester

💡 A propos

🗨 Forum

fr ▼

Où se trouvent les données excel à convertir ?

☒

Dans un des fichiers d'exemple fourni

Example 1 (simple exemple, in english)▼

Télécharger l'exemple : Example 1 (simple exemple, in english)

☐

Dans un fichier sur mon ordinateur

Sélectionner un fichier

(Extensions supportées: .xls ou .xlsx. Les fichiers OpenOffice ne fonctionnent pas !)

☐

Sur le web

http://...

## OpenRefine + extension RDF-openRefine

**OpenRefine** Jeu de donnés d'exemples pour Outils de FAIRisation.xlsx    Permalink    Open... Export Help

---

Facet / Filter    Undo / Redo 12 / 12
Extensions: RDF Wikidata

Refresh    Reset All    Remove All
Show as: rows records    Show: 5 10 25 50 rows
« first » previous 1 - 12 next » last »

**X Libellé FR**

12 choices Sort by: name count Cluster

- Activité 1
- Andainage 1
- Andaineuse 1
- Aplatisseur de grains 1
- Broyeur 1
- Conditionneur de fourrage 1
- Faneuse 1
- Matériel 1
- Matériel Agricole 1
- Matériel d'élevage 1
- Matériel de récolte 1
- Tritrateur 1

All	ID	Libellé FR	Synonyme FR	Libellé EN	Définition
☆	1. c_2631	Matériel	équipement, outillage	equipment	Les objets nécessaires à une exploitation.
☆	2. c_25753	Matériel Agricole		farm equipment	
☆	3. c_25889	Broyeur		grinders	Qui permet de broyer.
☆	4. c_37683	Tritrateur		crushers	Appareil servant à la trituration des substances.
☆	5. c_25746	Matériel d'élevage		animal husbandry equipment	
☆	6. c_3499	Matériel de récolte		harvesters	Machine agricole qui rassemble une culture vivrière des champs
☆	7. c_8400	Andaineuse		windrowers	
☆	8. c_3036	Conditionneur de fourrage		forage conditioners	Un outil agricole qui utilise des fourches mobiles pour aérer ou "ébouffier" l'andainage.
☆	9. c_7649	Faneuse		tedders	
☆	10. c_330634	Activité		activities	
☆	11. c_25675	Aplatisseur de grains		grain crushing	
☆	12. c_25674	Andainage		windrowing	

## VocBench + Sheet2RDF

About VocBench

VocBench

Projects

Data

Meladata

SPARQL

Tools

Spreadsheet file:

Browse

Jeu de données d'exemples pour Outils de FAIRisation.xlsx

Spreadsheet preview (Rows: 20 of 26)

Subject mapping

Pearl

ID	Libellé FR	Synonyme FR	Libellé EN	Définition	Hérarchie (parents)
c_2631	Matériel	équipement, outi...	equipment	Les objets néc...	
c_25753	Matériel Agricole		farm equipment		c_2631
c_25889	Broyeur		grinders	Qui permet de ...	c_25753
c_37683	Triturateur		crushers	Appareil serva...	c_25753

Legend

Generated triples preview

Subject

Predicate

## Opentheso

✓ **Correction 1 - Restructuration**

**Permet de corriger les incohérences dans le thésaurus en cours :**

- 1- détecter les concepts TT erronés : si le concept n'a pas de BT, alors, il est forcément TopTerm,
- 2- compléter le thésaurus par les relations qui manquent NT ou BT,
- 3- supprimer les relations en boucle (100 -> BT -> 100) ou (100 -> NT -> 100) ou (100 -> RT -> 100).

**!!! Ne pas oublier de recharger le thésaurus à la fin du traitement !!!**

[Lancer la correction ✓](#)

<https://vocabulaires-ouverts.inrae.fr/2022/10/13/3-outils-transformer-tabule-skos/>

# FAIRiser un thésaurus

Un vocabulaire FAIR doit être

- **Facile à trouver :**
  - être enregistré (indexé, répertorié) dans un entrepôt dédié aux vocabulaires ou généraliste
  - avoir des identifiants uniques et pérennes (PID)
- **Accessible :**
  - être accessible sur le web (https)
  - téléchargeable dans des formats standardisés (csv, rdf, Json)
  - requêtable via des APIs (open API)
- **Interopérable :**
  - encodé dans une représentation standard, telle que le langage d'ontologie Web (OWL) ou SKOS
  - et les extensions spécifiques aux domaines
- **Réutilisable :**
  - avec une licence explicite
  - une documentation
  - et une organisation pour sa maintenance

# Outils spécialisés open-source

## Construire et gérer « AtoZ »

- Intelligent Taxonomy Manager - <https://mondeca.com/logiciels/>
- TemaTres - <https://sourceforge.net/projects/tematres/>
- Protégé - <https://protege.stanford.edu/>
- Ginco, Gestion informatisée de nomenclatures collaboratives ouvertes - <https://github.com/culturecommunication/ginco>
- Opentheso - <https://github.com/miledrousset/Opentheso2>

## Contrôler

- SKOSPlay - <https://skos-play.sparna.fr/play/>

## Publier et consulter

- Skosmos - <https://www.skosmos.org/>
- VocBench - <https://bitbucket.org/art-uniroma2/vocbench3/downloads/>

# Opentheso & Pactols

- Depuis 2005 pour Frantiq
- Outil complet
  - Création
  - Gestion et administration
  - Collaboration (4 niveaux d'autorisation)
  - Publication
- Outil normalisé
  - ISO 25964
  - Multilingue
  - SKOS
  - API
- Diffusé
  - Github
  - IR\* Huma-Num
  - Projets nationaux & internationaux
- Documenté  
<https://opentheso.hypotheses.org/>

The screenshot displays the Opentheso web interface. On the left, a green sidebar contains a tree view of concepts, with 'objet en cuir' selected. The main content area shows the details for 'objet en cuir (fr)'. It includes fields for 'Libellé', 'Variante du libellé', 'Collection' (P2-ENTITÉS MATÉRIELLES, P1-SUJETS), 'Facettes' (objets mobiles par matériau), 'Total de la branche' (Catalogue Frantiq (16 notices)), 'Corpus lié', 'Concept générique' (objets mobiles), 'Concept spécifique', 'Concept associé' (caveçon, cuir, tannerie, travail des peaux), 'Traduction' (leather object (en), Lederobjekt (de), objeto de cuero (es), oggetto di cuoio (it)), 'Notes', 'Alignement' (exactMatch: https://www.wikidata.org/wiki/Q99689314 (Wikidata)), 'Notation' (ID Interne: 178439, Uri: https://ark.frantiq.fr/ark:/26678/ctZFzfpJ0lgB, ID Ark: 26678/ctZFzfpJ0lgB), and 'Exporter le concept en SKOS' (RDF/XML, json, jsonLd, Turtle). A QR code is also present. At the bottom, it shows 'Créé le: 2018-10-26' and 'Dernière modification le: 2023-08-02'. A small orange button 'Proposer une amélioration' is at the bottom right. The footer indicates 'Copyright ©CNRS Opentheso V23.09.03'.



# Faire vivre le thésaurus

## Mise à jour le lexique

- Gérer les candidats
- Administrer les propositions d'améliorations

## Révision des branches

## Mise à niveau des traductions

## Suivi des alignements

- Wikidata
- IdRef
- Geonames, Pleiades
- PeriodO

## Communication

- En interne GT Pactols
- [utilisateurs-pactols@listes.services.cnrs.fr](mailto:utilisateurs-pactols@listes.services.cnrs.fr)
- Web et RS

## Formation

- Les catalogueurs Frantiq
- Les autres utilisateurs (MASA)

# Expériences publiées

- MENG Fan, ZHOU Kaile, BU Yi, HUANG Win-Bin, ZHANG Pengyi, LONG Fei et Li Yan, « Keywords Extraction and Thesaurus Construction for Domain News », *Procedia Computer Science*, vol. 214, 2022, p. 837-844.  
<https://www.sciencedirect.com/science/article/pii/S1877050922019597>
- Building the DEFC thesaurus / Digitizing Early Farming Cultures project  
<https://defc.acdh.oeaw.ac.at/blog/post03/>
- KOMBOLO Moise, YON Jérémy, LANDRIEU François, RICHON Brigitte, AUBIN Sophie et HOCQUETTE Jean-François, « Le Thésaurus de la viande: un nouvel outil accessible à tous Une nouvelle ressource sémantique répondant aux principes de la science ouverte: le thésaurus de la viande comme outil informatique de dialogue entre les acteurs de la filière », *Viandes & produits carnés*, 4 avril 2022, p. 14.  
<https://hal.inrae.fr/hal-03670992/document>
- Thesaurus sur la diaspora canadienne sud-asiatique :  
<https://doi.org/10.18357/kula.223>

# Quelques sites de référencement

- AgroPortal / Univ de Montpellier, CNRS, Inrae <https://agroportal.lirmm.fr/ontologies>
  - + 160 référentiels et ontologies sur l'agronomie
- BARTOC, Basic Register of Thesauri, Ontologies & Classifications <https://bartoc.org/>
  - Recense + 3400 vocabulaires
- FAIRsharing <https://fairsharing.org/FAIRsharing>
  - Plusieurs milliers de standards, terminologies, formats... tous domaines
- Glossary links / Term Coord (European Parliament, DG TRAD, Terminology coordination) <https://termcoord.eu/glossarylinks/>
  - Recense + 2000 glossaires sur tous les sujets traités par l'UE, évolutif
- Loterre / CNRS-Inist <https://www.loterre.fr/>
- TemaTres / Diego Ferreyra <https://vocabularyserver.com/web/items/browse?collection=1>
  - Recense les utilisateurs de l'application (690)
- Data Geo-science Vocabularies / BRGM <https://data.geoscience.fr/ncl/>