

Constitution de corpus - Défis méthodologiques et solutions apportées par l'infrastructure Istex

Mathilde Huguin

▶ To cite this version:

Mathilde Huguin. Constitution de corpus - Défis méthodologiques et solutions apportées par l'infrastructure Istex. École thématique. Modèles de langue pour le traitement sémantique et l'intégration de connaissances et données en agriculture, alimentation et environnement, Montpellier, France. 2025. hal-05332761

HAL Id: hal-05332761 https://cnrs.hal.science/hal-05332761v1

Submitted on 27 Oct 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Constitution de corpus

Défis méthodologiques et solutions apportées par l'infrastructure Istex

Mathilde Huguin

IR PhD Inist - CNRS mathilde.huguin@inist.fr

2025 | École thématique Modèles de langue pour le traitement sémantique et l'intégration de connaissances et données en agriculture, alimentation et environnement















Inist- CNRS @ (1) (6)

Plan

PARTIE 1: ÉLÉMENTS CONTEXTUELS

- 1. Les corpus
- 2. Défis méthodologiques
- 3. Istex

PARTIE 2 : CONSTITUTION D'UN CORPUS SPÉCIALISÉ

- 4. Présentation du cas d'usage
- 5. Sélection des documents
- 6. Exploration, visualisation & affinage
- 7. Premières annotations : extractions d'EN
- 8. Conclusion
- 9. Annexes et références





Les corpus



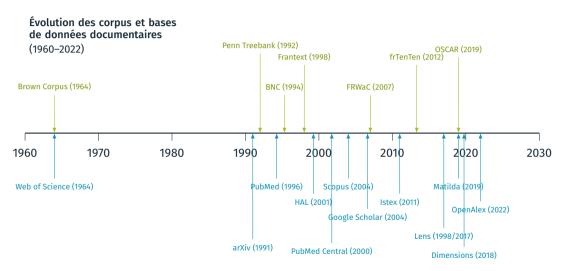
Les corpus

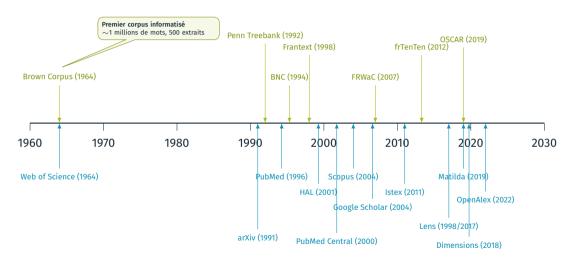
Qu'est ce qu'un corpus?

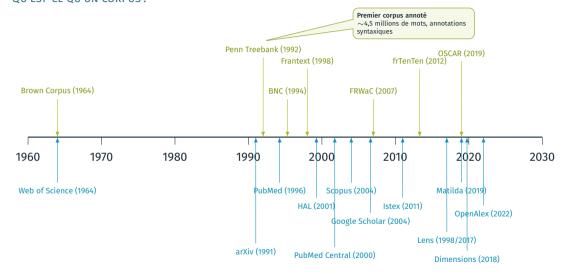
- Employé depuis l'antiquité dans la langue savante pour désigner des collections de textes (juridiques, épigraphiques, ...) : "Recueil formé d'un ensemble de données sélectionnées et rassemblées pour intéresser une même discipline" (Mellet, 2002)
- Apparition du terme chez les linguistes dans les années 60 alors que des collectes existent depuis longtemps (premiers corpus informatisés; cf. Léon, 2008; Benzitoun et Cappeau, 2025)
- Définition très discutée en linguistique et en linguistique de corpus (taille, représentativité, etc. cf. Sinclair, 1996; Habert, 2000); distinction entre réservoir (ou archive) et corpus (Rastier, 2004)
- Depuis 2000, évolution exponentielle des réservoirs et des corpus (documentaires ou non)
- Matière première des LLM : ils structurent, apportent savoir, diversité linguistique et contextuelle, conditionnent l'efficacité des modèles
- Utilisation des LLM sur les corpus (augmentation de 47% des publications en 6 ans, entre 12 et 82% des publications jamais lues, cf. Larivière *et al.*, 2009)

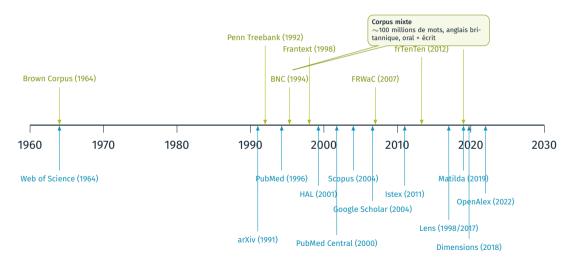
Partie 1 : éléments contextuels | Les corpus

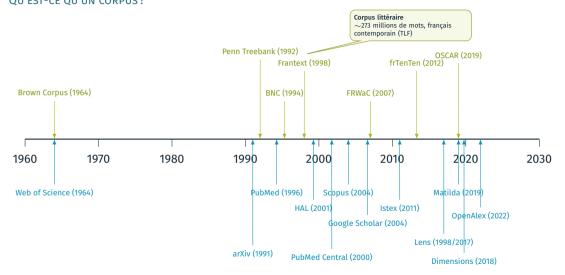
Qu'est-ce qu'un corpus?

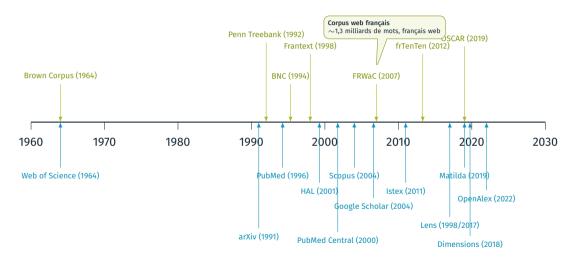


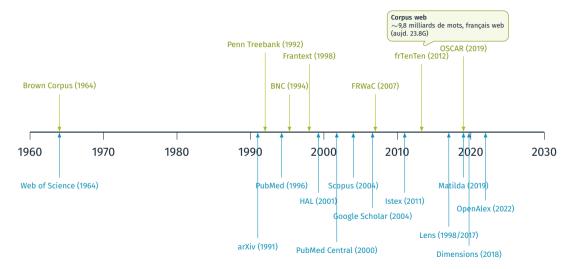


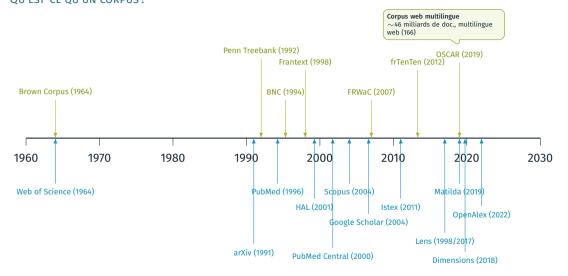














Les corpus

Corpus & LLM

Partie 1: éléments contextuels | Les corpus CORPUS & LLM

Type de corpus	Fonction	Exemple
Pre-training Cor- pora	Apprentissage initial du modèle sur un vaste corpus généraliste multi-domaines	Donner au modèle une connaissance gé- nérale du monde → Common Crawl pour ChatGPT (données massives issues du web)
Instruction Fine- tuning Datasets	Adaptation du modèle via des exemples su- pervisés ou consignes explicites	Répondre à des instructions claires → Alpaca (52 000 instructions-réponses générées par GPT-3.5, utilisé pour fine-tuner LLaMA)
Preference Data- sets	Apprentissage des préférences utilisateur, renforcement par feedback humain	Privilégier les réponses utiles → Stiennon et al. (2022) (réponse A plus pertinente que B)
Evaluation Data- sets	Évaluation de la qualité, précision et robus- tesse du modèle à l'aide de benchmarks stan- dardisés	rester la compréhension multi-domaines → MMLU (questions multi-choix couvrant 57 domaines, cf. Hendrycks <i>et al.</i> , 2021)

Usages majeurs selon Liu et al. (2024)



Les corpus

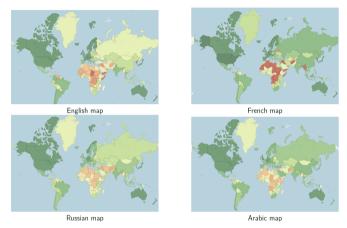
Biais

Partie 1 : éléments contextuels | Les corpus

Exemples de biais induits par les datasets

- IA du "pôle emploi" autrichien: "Le Berufsinfomat a par exemple proposé à des candidats masculins de postuler dans le domaine de l'informatique, tandis qu'ils conseillaient à des candidates ayant un CV identique de se tourner vers les études de genre ou la restauration" (cf. Courrier international 2024)
- Certains outils médicaux basés sur l'IA (OpenAl's GPT-4, Meta's Llama 3 and Palmyra-Med) ont tendance à minimiser les symptômes des femmes ou des minorités ethniques (cf. Financial Times 2025)

Partie 1 : éléments contextuels | Les corpus



« Mapping the answers to the prompt "Could you assign a positivity score to these words" for every country in the world, in 4 languages » (Boussidan et al., 2024)



Défis méthodologiques



Questions incontournables

Partie 1 : éléments contextuels | Défis méthodologiques OUESTIONS INCONTOURNABLES

La constitution d'un corpus peut être plus longue que son analyse

- Quels genres de texte doit-on inclure dans le corpus (scientifique, journalistique, conversationnel, littéraire, etc.)? Comment gérer la qualité linguistique des données (erreurs, bruit, doublons, textes tronqués, OCR de mauvaise qualité)?
- Quels sont les droits et licences associés aux données, et peut-on les utiliser légalement pour l'entraînement ou la recherche?
- Sous quels formats les données sont-elles disponibles (texte brut, PDF, HTML, XML, JSON, etc.), et sont-ils adaptés au traitement automatique?
- Quelle méthode de récupération des documents est employée (web scraping, API, dépôts institutionnels, partenariats, crowdsourcing), et quelles en sont les limites?

Partie 1 : éléments contextuels | Défis méthodologiques QUESTIONS INCONTOURNABLES

Genre	Atouts	Risques et limites
Juridique	Précision, structure, faible ambiguïté	Faible variété, manque d'expression naturelle et d'exemples pratiques
Web	Grande diversité thématique, registres variés	Qualité inégale, bruit, duplication, informations erronées
Journalistique	Clarté, structuration éditoriale, capacité argumentative	Biais idéologique/politique, diversité de formats, subjectivité
Scientifique	Rigueur factuelle, données structurées, terminologie ex- perte	Jargon/terminologie, spécialisation, généralisa- tion limitée
Brevet	Descriptions techniques dé- taillées, structure normée	Complexité technique, difficulté d'interprétation

Influence du genre de texte sur la qualité du corpus (e.g. Labadie et Prince (2008) sur la segmentation)

Partie 1: éléments contextuels | Défis méthodologiques QUESTIONS INCONTOURNABLES

Droits et licences associés aux données

- Textes littéraires : copyright souvent strict (ex. éditeurs modernes), Open Access rare (Project Gutenberg pour œuvres tombées dans le domaine public)
- Données du web : protégées par copyright, conditions d'utilisation souvent restrictives (sauf API ouvertes)
- Journaux quotidiens : protégés par copyright, reproduction interdite
- Publications scientifiques: copyright par l'éditeur, TDM possible pour la recherche (directive UE 2019; transposition de la directive européenne en droit français 2020/2021)

Formats

- CSV adapté à des contenus courts ou tabulaires mais ne gère pas la hiérarchie du texte et peut être limitant pour des documents complexes
- PDF, DOCX extractions sont souvent imparfaites : risques liés au formatage, aux éléments non textuels et à la segmentation non standardisée
 - TXT universel et simple à manipuler mais sans structure balisée, perte des tableaux ou figures
 - XML/TEI facilite l'extraction structurée mais nécessite souvent des parsers complexes, et la gestion de schémas différents selon les sources peut s'avérer fastidieuse
 - JSON permet de structurer les textes et les métadonnées mais requiert un nettoyage intensif si la structure varie

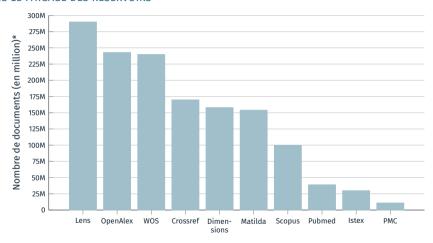
Partie 1 : éléments contextuels | Défis méthodologiques QUESTIONS INCONTOURNABLES

Moyen	Description / Avantages	Limites	Exemple
API	Accès structuré aux données, rapide et fiable	Limité en volume et quotas	HAL, arXiv API, Twitter API
Web scraping	Extraction flexible depuis des pages web	robots.txt bloquant, nettoyage nécessaire	Articles de presse, blogs, forums
Interface web	Téléchargement direct de lots de documents	Conversion parfois nécessaire (PDF → texte), métadonnées li- mitées	Project Gutenberg, PubMed Central
Crowdsour- cing	Création de corpus via contribu- tions humaines	Temps et coût humain, contrôle de qualité	

Moyens techniques pour récupérer des données



Istex dans le paysage des réservoirs



Comparaison quantitative des réservoirs (les chiffres peuvent varier en fonction des mises à jour et des méthodologies de comptabilisation)

Istex	Négociations, Achats, Plan de soutien, GIS CollEX-Persée, OA	Istex
PMC	Négociation, OA	National Library of Medicine
Pubmed	Références sélectionnées (MEDLINE)	National Library of Medicine
Scopus	Revues sélectionnées	Elsevier
WOS	Revues sélectionnées	Clarivate
Crossref	Crossref	Crossref
манца	Crossref, PubMed Central, arXiv et RePEc, BASE + Unpaywall, OR- CID	13 & Huma-Num
Matilda		I3 & Huma-Num
Dimensions	ROR, Internet archive, crawls web, arXiv, Zenodo Crossref. PubMed. DataCite. arXiv + 130 éditeurs	Digital Science
OpenAlex	ORCID Microsoft Academic, Crossref, PubMed, Unpaywall, DOAJ, ORCID,	OurResearch
Lens	Microsoft Academic, PubMed, Crossref, OpenAlex, Unpaywall,	Cambia
Réservoirs	Origine des données	Fournisseurs

Origine des données (Bouchard, 2024, Gusenbauer, 2022)

ISTEX DANS LE PAYSAGE DES RÉSERVOIRS

	Lens	OpenAlex	Dimensions	Matilda	Pubmed	PMC	Istex	
Filtres	+	+	+	+/-	+	+	+	

Propriétés des réservoirs : des usages différents

ISTEX DANS LE PAYSAGE DES RÉSERVOIRS

	Lens	OpenAlex	Dimensions	Matilda	Pubmed	PMC	Istex	
Filtres	+	+	+	+/-	+	+	+	
Tris	+	+	+	+	+	+	+	

Propriétés des réservoirs : des usages différents

ISTEX DANS LE PAYSAGE DES RÉSERVOIRS

	Lens	OpenAlex	Dimensions	Matilda	Pubmed	PMC	Istex
Filtres	+	+	+	+/-	+	+	+
Tris	+	+	+	+	+	+	+
Enrichissements	+	+	+	+	+	+	++

Propriétés des réservoirs : des usages différents

ISTEX DANS LE PAYSAGE DES RÉSERVOIRS

	Lens	OpenAlex	Dimensions	Matilda	Pubmed	PMC	Istex	
Filtres	+	+	+	+/-	+	+	+	
Tris	+	+	+	+	+	+	+	
Enrichissements	+	+	+	+	+	+	++	
Recherche assistée	++	+/-	_	_	+	+	++	

Propriétés des réservoirs : des usages différents

Formats biblio. (e.g. BIB, MODS)

 Filtres
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 +
 <t

Propriétés des réservoirs : des usages différents

	Lens	OpenAlex	Dimensions	Matilda	Pubmed	PMC	Istex
Filtres	+	+	+	+/-	+	+	+
Tris	+	+	+	+	+	+	+
Enrichissements	+	+	+	+	+	+	++
Recherche assistée	++	+/-	_	_	+	+	++
Formats biblio. (e.g. BIB, MODS)	+	+	+	+	+	+	+
Indicateurs biblio.	+	+	+	+	_	_	_

Propriétés des réservoirs : des usages différents

	Lens	OpenAlex	Dimensions	Matilda	Pubmed	PMC	Istex
Filtres	+	+	+	+/-	+	+	+
Tris	+	+	+	+	+	+	+
Enrichissements	+	+	+	+	+	+	++
Recherche assistée	++	+/-	_	_	+	+	++
Formats biblio. (e.g. BIB, MODS)	+	+	+	+	+	+	+
Indicateurs biblio.	+	+	+	+	_	_	_
Indicatours TDM			_		_		

Propriétés des réservoirs : des usages différents

	Lens	OpenAlex	Dimensions	Matilda	Pubmed	PMC	Istex
Filtres	+	+	+	+/-	+	+	+
Tris	+	+	+	+	+	+	+
Enrichissements	+	+	+	+	+	+	++
Recherche assistée	++	+/-	_	_	+	+	++
Formats biblio. (e.g. BIB, MODS)	+	+	+	+	+	+	+
Indicateurs biblio.	+	+	+	+	_	_	_
Indicateurs TDM	_	_	_	_	_	_	+
Texte intégral	_	_	_	_	_	+	+

Propriétés des réservoirs : des usages différents

	Lens	OpenAlex	Dimensions	Matilda	Pubmed	PMC	Istex
Filtres	+	+	+	+/-	+	+	+
Tris	+	+	+	+	+	+	+
Enrichissements	+	+	+	+	+	+	++
Recherche assistée	++	+/-	_	_	+	+	++
Formats biblio. (e.g. BIB, MODS)	+	+	+	+	+	+	+
Indicateurs biblio.	+	+	+	+	_	_	_
Indicateurs TDM	_	_	_	_	_	_	+
Texte intégral	_	_	_	_	_	+	+
Formats TDM standardisés (e.g. TEI)	_	_	_	_	_	_	+

Propriétés des réservoirs : des usages différents

	Lens	OpenAlex	Dimensions	Matilda	Pubmed	PMC	Istex
Filtres	+	+	+	+/-	+	+	+
Tris	+	+	+	+	+	+	+
Enrichissements	+	+	+	+	+	+	++
Recherche assistée	++	+/-	_	_	+	+	++
Formats biblio. (e.g. BIB, MODS) Indicateurs biblio.	+	+	+	+	+	+	+
Indicateurs biblio.	+	+	+	+	_	_	_
Indicateurs TDM	_	_	_	_	_	_	+
Texte intégral	_	_	_	_	_	+	+
Formats TDM standardisés (e.g. TEI)	_	_	_	_	_	_	+
Passerelles	_	_	_	_	_	_	+

Propriétés des réservoirs : des usages différents

Istex : infrastructure dédiée à la constitution et à l'analyse de corpus

⊕ Des données vérifiées accessibles un seul lieu pour de nombreuses sources

Istex : infrastructure dédiée à la constitution et à l'analyse de corpus

- ⊕ Des données vérifiées accessibles un seul lieu pour de nombreuses sources
- ⊕ Des données interopérables, des formats homogénéisés et des données corrigées (moins de prétraitements)

Istex : infrastructure dédiée à la constitution et à l'analyse de corpus

- ⊕ Des données vérifiées accessibles un seul lieu pour de nombreuses sources
- ⊕ Des données interopérables, des formats homogénéisés et des données corrigées (moins de prétraitements)
- ⊕ Des données enrichies (réocérisation, structuration de texte, métadonnées) : des documents retrouvés et analysés plus facilement

Istex : infrastructure dédiée à la constitution et à l'analyse de corpus

- ⊕ Des données vérifiées accessibles un seul lieu pour de nombreuses sources
- Des données interopérables, des formats homogénéisés et des données corrigées (moins de prétraitements)
- ⊕ Des données enrichies (réocérisation, structuration de texte, métadonnées) : des documents retrouvés et analysés plus facilement
- ⊕ Des millions de textes et de métadonnées téléchargeables facilement

Istex : infrastructure dédiée à la constitution et à l'analyse de corpus

- ⊕ Des données vérifiées accessibles un seul lieu pour de nombreuses sources
- Des données interopérables, des formats homogénéisés et des données corrigées (moins de prétraitements)
- ⊕ Des données enrichies (réocérisation, structuration de texte, métadonnées) : des documents retrouvés et analysés plus facilement
- ⊕ Des millions de textes et de métadonnées téléchargeables facilement
- ⊕ Des connexions vers des outils / plateformes du monde académique

Istex : infrastructure dédiée à la constitution et à l'analyse de corpus

- ⊕ Des données vérifiées accessibles un seul lieu pour de nombreuses sources
- ⊕ Des données interopérables, des formats homogénéisés et des données corrigées (moins de prétraitements)
- ⊕ Des données enrichies (réocérisation, structuration de texte, métadonnées) : des documents retrouvés et analysés plus facilement
- ⊕ Des millions de textes et de métadonnées téléchargeables facilement
- ⊕ Des connexions vers des outils / plateformes du monde académique
- ⊕ Un cadre juridique sécurisé par une licence appropriée et déjà négociée



Istex



Istex

Qu'est-ce qu'Istex?

Partie 1 : éléments contextuels | Istex Ou'est-ce ou'istex?

Istex est une infrastructure dédiée à la constitution et à l'analyse de corpus (ou fouille de textes, TDM) qui donne accès :

- à des ressources textuelles et lexicales enrichies (publications et terminologies)
- à des outils pour sélectionner et analyser ces ressources

Ou'est-ce ou'Istex? **Constituer un corpus** 30M de publications Recherche et télécharaement facilités API Istex Search 39 corpus prêts à l'emploi Annoter, enrichir et analyser Corpus Istex Istex Réservoir de publications Istex Loterre multilingues et multidisciplinaires Data Istex 73 terminologies Publication des résultats Istex TDM Lodex 44 web services de TDM

Se documenter

Datalake

Constitution de corpus Mathilde Huguin 28 / 90

Visualiser



Istex

Pour quels usages?

Partie 1: éléments contextuels | Istex Pour quels usages?

Une ressource documentaire

- Trouver des articles scientifiques
- Alimenter des bases documentaires, des agrégateurs (Google Scholar, Click & Read, ENT)

Un matériau de recherche

- Constituer des corpus
- Faire de la fouille de textes ou du TDM, du TAL (entraîner des LLM, faire de la RI, etc. e.g. Bénard et al., 2023)

Projet ANR MaTOS



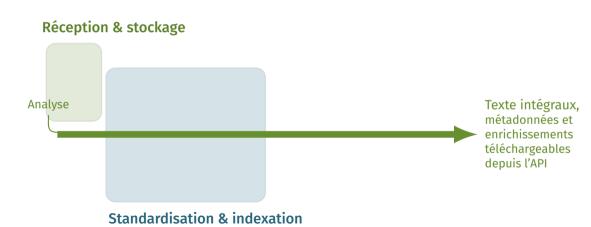
Istex

Chaîne de traitements

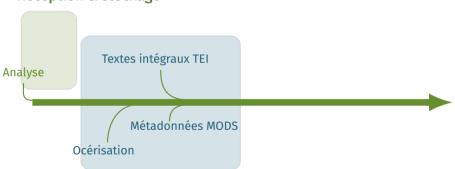
Réception & stockage



Texte intégraux, métadonnées et enrichissements téléchargeables depuis l'API

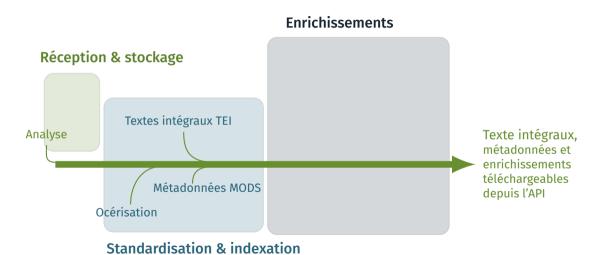


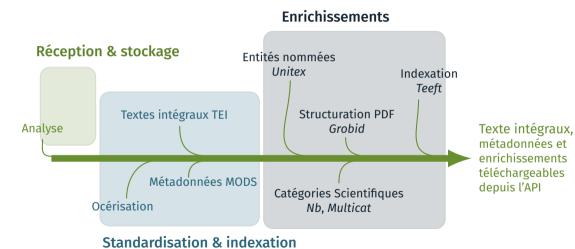
Réception & stockage



Texte intégraux, métadonnées et enrichissements téléchargeables depuis l'API

Standardisation & indexation





PARTIE 2 : CONSTITUTION D'UN CORPUS SPÉCIALISÉ



Partie 2 : constitution d'un corpus spécialisé | Présentation du cas d'usage

Objectif Cette recherche en sociologie de l'alimentation vise à comprendre comment les motivations des personnes véganes et végétariennes évoluent au fil du temps (justifications éthiques, sanitaires ou écologiques). L'étude s'intéresse spécifiquement à la manière dont ces transformations sont traitées dans la littérature scientifique internationale.

Moyen Afin d'explorer ces évolutions, un corpus de publications académiques internationales sera constitué et analysé.

Partie 2 : constitution d'un corpus spécialisé | Présentation du cas d'usage

Q	Requêtage et première exploration des résultats dans Istex Search
*	Téléchargement du corpus au format approprié
•	Visualisation et exploration des résultats dans Lodex
•	Affinage du corpus
	Premières annotations : extractions des EN

Sélection des documents

Exploration, visualisation & affinage

Annotations



Sélection des documents

Présentation d'Istex Search

Partie 2 : constitution d'un corpus spécialisé | Sélection des documents Présentation d'Istex Search

Application en ligne qui permet :

- D'interroger l'API grâce à trois modes de requêtage
- De consulter un document dans une notice détaillée
- De naviguer dans les résultats avec différents modes de tri et différents filtres
- De vérifier la compatibilité TDM grâce à des indicateurs techniques
- De sélectionner / exclure des documents pour éliminer simplement le bruit
- De télécharger massivement des documents au format approprié

https://search.istex.fr



Requêtage

C

Requêtage et première exploration des résultats dans Istex Search

Requête 1

végan vegan véganisme veganism végétarien vegetarian végétarisme vegetarism végétalien vegetalian végétalisme vegetalism végétarianisme vegetarianism Description Liste de mots-clés anglais et français, coordonnés par l'opérateur booléen OR (= ou inclusif), insensibilité à la casse mais sensibilité aux diacritiques (chat = CHAT: éclair ≠ eclair)

Résultats 41 739 documents

Problèmes Il y a du silence; Variantes morphosyntaxiques ou dérivationnelles non recherchées (tokenisation lacto-végétarien mais absence de lemmatisation)

Requête 2

végan* vegan* végétar* vegetar* végétalien* vegetalian* végétalisme* vegetalism*

Description Recherche des formes morphologiquement apparentées grâce à la troncature symbolisée par *

Résultats 53 594 documents

Problèmes Il y a du bruit (documents récents)

Original Article

Effect of slope on the forest structure of the Atlantic Forest domain in southern Brazil

Efeito da inclinação na estrutura de uma floresta do domínio Mata Atlântica no su Brasil

D. C. Souza^a O, L. R. Souza^b O, E. V. Couto^a O, M. G. Caxambú^a O and A. P. Peron^{ac} O

⁴Universidade Tecnológica Federal do Paraná — UTFPR, Departamento de Biodiversidade e Conservação da Natureza — DABIC, PR, Brasil

¹Universidade Tecnológica Federal do Paraná – UTFPR, Curso de Engenharia Ambiental, Campo Mourão, PR, Brasil
¹Universidade Tecnológica Federal do Paraná – UTFPR, Programa de Pós-graduação em Inovações Tecnológicas – PPGIT, Campo N

Effet de la pente sur la structure d'une forêt brésilienne

Original Article

Effect of slope on the forest structure of the Atlantic Forest domain in southern Brazil

Efeito da inclinação na estrutura de uma floresta do domínio Mata Atlântica no su Brasil

D. C. Souza** @, L. R. Souzab @, E. V. Coutoa @, M. G. Caxambúa @ and A. P. Peronac @

'Universidade Tecnológica Federal do Paraná — UTFPR, Departamento de Biodiversidade e Conservação da Natureza — DABIC, PR. Brasil

¹Universidade Tecnológica Federal do Paraná – UTFPR, Curso de Engenharia Ambiental, Campo Mourão, PR, Brasil

'Universidade Tecnológica Federal do Paraná – UTFPR, Programa de Pós-graduação em Inovações Tecnológicas – PPGT. Campo N

Effet de la pente sur la structure d'une forêt brésilienne

the temporary crop begins, being these collections every 20 meters along the slope. In these samples, macro- and micro-nutrients were analyzed.

The physical analysis of the remaining soil was carried out through three trenches called Tt, T2 and T3; T1 located at high slope, and T2 and T3 at medium slope, soil particle size analyses were made for each trench, for soil classification according to the Brazilian Soil Classification System (Embrapa, 2006), in the trenches, color and horizon were determined in situ. As for the color, the Munsell Color Chart was used as a reference (USDA, 1988), The slope of the area was determined along a transact that covered all the plots, using GNSS spectra promark 220, accurate to 0.06m. The same equipment was used for the follow the same plots can be used.

2.3. Phytosociology

The forest structure was sampled in 12 continuous rectangular plots, with 30 X 20 m, totaling 0.72 ha. The

2.4 Data analysis

Species were analyzed for clusters according to their distribution in the area by means of a Detrended Correspondence Analysis (DCA). These distribution data were used together with soil chemical properties to set some pattern between the plots. These were analyzed by Canonical Correspondence Analysis (CCA), which allows associating species distribution with abiotic variables. Pedological variables with the greatest pressure on phytosociological characteristics were tested for patterns in the environmental gradient by Regression Analysis to identify the influence of the variable and creating prediction. models. All analyses were run in the R program. To test normality and perform regression analysis, the shapiro test function and the Im function of the R base were used (R Core Team, 2020). For multivariate analysis, the decorana and cca functions of the vegan package were used (Oksanen et al., 2020).

3. Results and Discussion

Vegan est le nom d'un package R cité dans le corps du texte

Research paper

Preventing intense false positive and negative reactions attributed to the principle of ELISA to re-investigate antibody studies in autoimmune diseases

Kuniaki Terato ^{a,*}, Christopher T. Do ^a, Dawn Cutler ^a, Takaki Waritani ^a, Hiroshi

^a Chondrex Inc., 2607 151st Place NE, Redmond, WA 98052, United States

Publication sur les maladies auto-immunes

Research paper

Preventing intense false positive and negative reattributed to the principle of ELISA to re-investigatudies in autoimmune diseases

Kuniaki Terato ^{a,*}, Christopher T. Do ^a, Dawn Cutler ^a, Takaki W ^a Chondrex Inc. 2607 151st Place NE. Redmond. WA 98052. United States

Publication sur les maladies auto-immunes

Kijima, Y., Iwatsuki, S., Akamatsu, H., Terato, K., Kuwabara, Y., Ueda, S., Shionoya, H., 2009. Natural antibodies to pathogenic bacteria and their toxins in whey protein concentrate. Food Sci. Technol. Res. 56, 475.

Kjeldsen-Kragh, J., 1999. Rheumatoid arthritis treated with vegetarian diets. Am. I. Clin. Nutr. 70, 594S.

Komiya, İ., Arimura, Y., Nakabayashi, K., Yamada, A., Osaki, T., Yamaguchi, H., Kamiya, S., 2011. Increased concentrations of antibody against heat shock protein in patients with myeloperoxidase anti-neutrophil cytoplasmic autoantibody positive microscopic polyangiitis. Microbiol. Immunol. 55, 531.

Ligier, S., Fortin, P., Newkirk, M., 1998. A new antibody in rheumatoid arthritis targeting glycated IgG: IgM anti-IgG-AGE. Br. J. Rheumatol. 37, 1307

Lucke, K., Miehlke, S., Jacobs, E., Schuppler, M., 2006. Prevalence of Bacteroides and Prevotella spp. in ulcerative colitis. J. Med. Microbiol. 55, 617.

Maes, M., Twisk, F., Kubera, M., Ringel, K., Leunis, J., Geffard, M., 2011. Increased IgA responses to the LPS of commensal bacteria is associated with inflammation and activation of cell-mediated immunity in chronic fatigue syndrome. I. Affect. Disord. 136. 909.

Terme présent dans les références bibliographiques

Requête 3

```
title:(végan* vegan* végétar* vegetar*
végétalien* vegetalian* végétalisme*
vegetalism*) abstract:(végan* vegan*
végétar* vegetar* végétalien*
vegetalian* végétalisme* vegetalism*)
subject.value:(végan* vegan* végétar*
vegetar* végétalien* vegetalian*
végétalisme* vegetalism*)
keywords.teeft:(végan* vegan*
végétar* vegetar* végétalien*
vegetalian* végétalien*
```

Description Cibler les champs les plus à même de renvoyer des documents pertinents: title, abstract, subject.value, keywords.teeft

Résultats 8 143 documents

Problèmes Il y a du bruit (documents anciens)

An epystell of ye famous doctor Erasmus of Roterdam: vnto the reuerende father & excellent prince, Christofer bysshop of Basyle, co[n]cernyng the forbedynge of eatynge of flesshe, and lyke constitutions of men. &c

An epystell of ye famous doctor Erasmus of Roterdam : vnto the reuerende father & excellent prince, Christofer bysshop of Basyle, co[n]cernyng the forbedynge of eatynge of flesshe, and lyke constitutyons of men. &c

Erasmus, Desiderius d. 1536

★ REVENIR AUX RÉSULTATS

PARTAGER LE DOCUMENT

Consulter ce document
Tecte Intégral

France Consulter ce document
Tecte Intégral

France Consulter ce document
Tecte Intégral

France Consulter ce document
Tecte Intégral

France Consulter ce document

France Consul

Informations sur ce document

Publication de 1534

An epystell of ye famous doctor Erasmus of Roterdam: vnto the reuerende father & excellent prince, Christofer bysshop of Basyle, co[n]cernyng the forbedynge of eatynge of flesshe, and lyke constitutions of men. &c.

An epystell of ye famous doctor Erasmus of Roterdam: vnto the reuerende father & excellent prince, Christofer bysshop of Basyle, co[n]cernyng the forbedynge of eatynge of flesshe, and lyke constitutyons of men. &

Erasmus, Desiderius d. 1536

→ PAST DE L'ESULTATS

✓ PARTAGER LE DOCUMENT

✓ PARTAGER LE DOCUMENT

Publication de 1534







Document sans résumé, non océrisé



Indicateur des langues de publication



Indicateur des langues de publication

404 Жива Бенчич

В анализе короткого рассказа Евгения Полова "Как съели петуха" мы могли бы прибегнутъ к интересному положению Клода Леви-Стросса, согласно которому кухия любого общества представляет собой особый язык, на котором оно бессознательно раскрывает свою структуру (Lévistruss 1983: 373). Дело в том, что, по утверждению французского ученого, кулинарные приемы, как и язык, можно анализировать в конститутивных элементах, так называемых тустемах, организованных по определенным структурам оппозиции и коррелярии. Так, например, английскую кухню можно отличать от французской посредством трех оппозиций:

»идоленное/эколенное (то есть национальное или заграничное сырье); µентральное/периферцийное (главнюе блюдо и сопутетвующие блюда); выраженное/невъргаженное (то есть вкусное или безвкусное) [...]. Другими словами, в английской кухие на основные, безвкусное ризтовленные блюда, днут отечественные продукты; а на сопутствующие – импортная продукция, в которой все дифференциальные качества ярко выраженные чай, фруктовый пирог, анслыснювый марменад, портвейи). И наоборот, во французской кухне оппозиция эндоленное/эколенное значительно ослаблена или отсутствукта, а даннаково выраженные густемы сочетаются друг с другом, как в центральной, так и в перименейской скарами.

Document russe

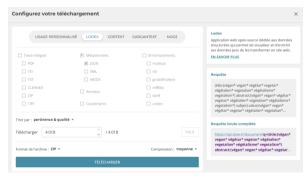
Requête 4

```
(title:(végan* vegan* végétar* vegetar* végétalien* vegetalian* végétalisme* vegetalism*) abstract:(végan* vegan* végétar* végétar* végétalien* vegetalian* végétalien* vegetalian* végétalisme* vegetalism*) subject.value:(végan* vegan* végétar* vegetar* végétalien* vegetalian* végétalisme* vegetalian* végétalisme* vegetalism*) keywords.teeft:(végan* vegan* végétar* végétar* végétalien* vegetalian* végétalien* végétalien* végétalien* végétalisme* vegetalian*) AND publicationDate:[1950 TO 2025] AND abstract:* AND language:(eng fre)
```

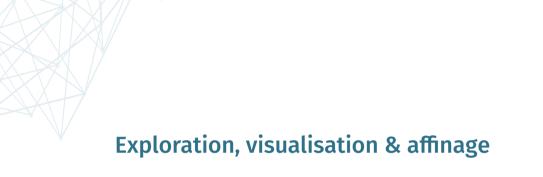
Description Contraindre la présence d'un résumé, contraindre les langues, cibler les dates de publications grâce à l'intervalle [1950 TO *] le tout grâce au parenthésage et à l'opérateur AND

Résultats 4 018 documents





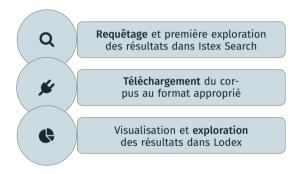
Modale de téléchargement des données dans Istex Search



Exploration, visualisation & affinage

Présentation de Lodex

PRÉSENTATION DE LODEX



PRÉSENTATION DE LODEX

« Linked Open Data EXperiment »

Application web open-source dédiée à la **visualisation** de données structurées (Gregorio *et al.*, 2019)

https://www.lodex.fr

PRÉSENTATION DE LODEX



Transformer ses données en site web À partir de différents formats (CSV, JSON, TXT,...)

PRÉSENTATION DE LODEX



Transformer ses données en site web

À partir de différents formats (CSV, JSON, TXT,...)



Explorer ses données

À l'aide de différents graphiques, filtres et au travers de données complémentaires

PRÉSENTATION DE LODEX



Transformer ses données en site web

À partir de différents formats (CSV, JSON, TXT,...)



Explorer ses données

À l'aide de différents graphiques, filtres et au travers de données complémentaires



Enrichir ses données

Grâce à un catalogue de web services

PRÉSENTATION DE LODEX



Transformer ses données en site web

À partir de différents formats (CSV, JSON, TXT,...)



Explorer ses données

À l'aide de différents graphiques, filtres et au travers de données complémentaires



Enrichir ses données

Grâce à un catalogue de web services



Aligner ses données

Avec des données similaires ou connexes (requête SPARQL)

PRÉSENTATION DE LODEX

A	$ \leftarrow $
ч	0
V	LV

Transformer ses données en site web À partir de différents formats (CSV, JSON, TXT,...)



Explorer ses donnéesÀ l'aide de différents graphiques, filtres et au travers de données complémentaires



Enrichir ses données Grâce à un catalogue de web services



Aligner ses données



Exporter ses données

Avec des données similaires ou connexes (re-

quête SPARQL)

Totalement ou partiellement (JSONL, CSV...)

PRÉSENTATION DE LODEX



Transformer ses données en site web

À partir de différents formats (CSV, JSON, TXT,...)



Explorer ses données

À l'aide de différents graphiques, filtres et au travers de données complémentaires



Enrichir ses données

Grâce à un catalogue de web services



Aligner ses données

Avec des données similaires ou connexes (requête SPARQL)



Exporter ses données

Totalement ou partiellement (JSONL, CSV...)

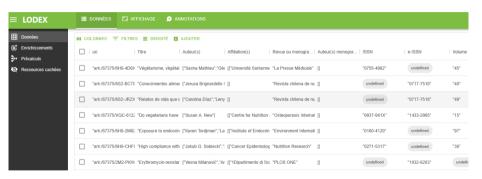
Exemple

PRÉSENTATION DE LODEX



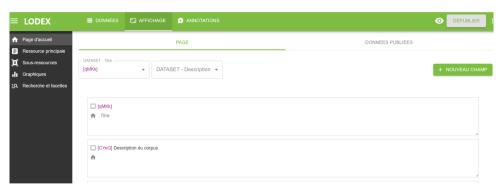
Interface d'import des données dans Lodex

PRÉSENTATION DE LODEX



Transformation des données par le loader

PRÉSENTATION DE LODEX



Import du modèle

Exploration, visualisation & affinage

Exploration



Corpus Véganisme et Végétarisme

Description du corpus

Le présent travail s'appuie sur le **corpus « Véganisme et végétarisme »**, conçu comme ressource documentaire destinée à l'analyse des transformations sociales et culturelles lièes aux régimes aillementaires alternatifs. Ce corpus s'inscrit dans une perspective de sociologie de l'alimentation, en mettant l'accent sur l'évolution des reordsentations, des pratiques et des aspirations associales au véganisme et ut végétafsime.

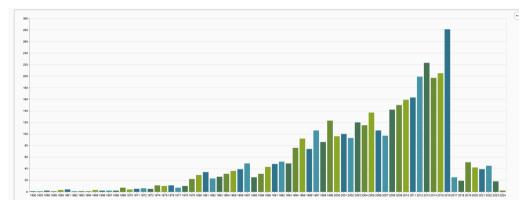
Constitué à partir des publications indexées dans la base **Istex Search**, il bénéficie d'un enrichissement automatisé grâce à plusieurs services de **text and data mining (TDM)**, notamment :

- Teeft, qui assure l'extraction de mots-clés et facilite l'identification des thématiques dominantes;
 entityTag qui procède à la reconnaissance des entités nommées (individus, lieux, organismes.).
- Ce dispositif permet non seulement de cartographier la production scientifique consacrée à ces pratiques alimentaires, mais également d'envisager une analyse dischronique des discours académiques et de leurs évolutions. L'objectif est ainsi de fournir un matériau robust en tinterdisciplinaire, apte à nourrir un effektion sur les dynamiques sociales, culturelles et orbitiques eu l'aversent les débats contemporaire, autour du véantiment du véaletiment.

CONSULTER LA REQUÊTE INITIALE

Site web du corpus créé avec Lodex

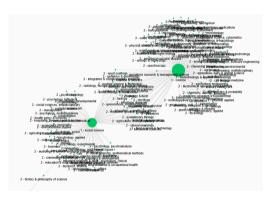
Partie 2 : constitution d'un corpus spécialisé | Exploration, visualisation & affinage EXPLORATION



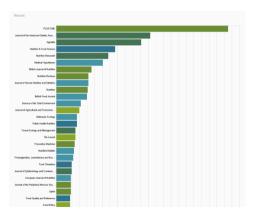
Années de publication



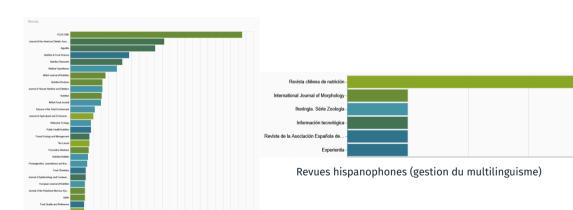
Mots-clés d'auteurs



Catégories WoS



Revues du corpus



Revues du corpus

Rev Chil Nutr 2020: 47(5): 782-791.

http://dx.doi.org/10.4067/\$0717-75182020000500782

Artículo Original / Original Articl e

Estado nutricional, hábitos de alimentación y de estilo de vida en vegetarianos de Asunción y Gran Asunción, Paraguay

Nutritional status, eating habits and lifestyle factors among vegetarians from Asunción and Great Asunción, Paraguay

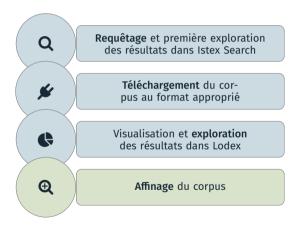
RESUMEN

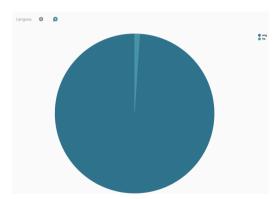
Una deta vegetarian se define por la exclusión total o parcial de allemonte de origen animal. Cuando es palnificada cal de allemonte de origen animal. Cuando es palnificada adecuadamente puede prevenir y tratar enfermedades. Sin embargo, una didar astrictiva en circos allemonto delire entrago, una parámetros antroponetricos, hábitos de allemontación se describir de vida en ovo-lásten-vegetarianos y veganos y explora diferencias en la negesta de nutrientes criticos entre anogrupos. Estado descriptivo de conte transversal, nalizado en condicionamento de sucreta de su consecuencia de sucreta de condicionamento de sucreta de sucreta de sucreta de sucreta de condicionamento de sucreta de sucreta de sucreta de sucreta de condicionamento de sucreta de sucreta de sucreta de sucreta de condicionamento de sucreta de sucreta de sucreta de sucreta de condicionamento de sucreta de sucreta de sucreta de sucreta de condicionamento de sucreta de sucreta de sucreta de sucreta de sucreta de condicionamento de sucreta de sucreta de sucreta de sucreta de sucreta de condicionamento de sucreta de sucreta de sucreta de sucreta de sucreta de condicionamento de sucreta Meliessa Penner Teichgräf!, Natalia Elizabeth González Cañete!".

Carrera de Nutrición, Facultad de Ciencias Médicas,
 Universidad del Pacífico, Asunción, Paraguay.

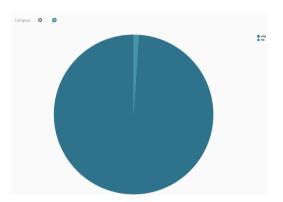
Exemple d'article issu de la revue Revista chilena de nutrición

AFFINAGE

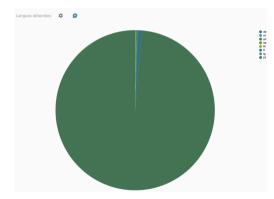




Langues: français (42), anglais (3 976)



Langues: français (42), anglais (3 976)



Langues issues du WS (8 langues; 3 955 en/4 018)



Suppression des documents dans le dataset



Partie 2 : constitution d'un corpus spécialisé | Premières annotations : extractions d'entités nommées

PRÉSENTATION D'ISTEX TDM



Premières annotations: extractions d'EN

Présentation d'Istex TDM

Les services de fouille de textes

- Istex met a disposition 44 web services dans Istex TDM (Cuxac, 2024)
 - Prétraitement (textExtract)
 - Classification (aiAbstractCheck)
 - Validation (bibCheck)
 - Extraction (diseaseTag)
 - Alignement (loterre-resolvers)
- Un web service est un programme mono-tâche, frugal avec un paramétrage minimal
- Les web services peuvent être utilisés dans Lodex, via TDM Factory, ou en ligne de commande

Premières annotations: extractions d'EN

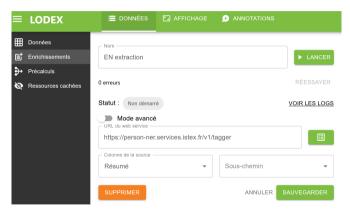
Extraction d'information

EXTRACTION D'INFORMATION



Istex contient des annotations mais elles sont parfois incomplètes (Unitex)

EXTRACTION D'INFORMATION



Extraction des entités nommées de lieux, organismes et personnes grâce au web service multilingue entitytag

EXTRACTION D'INFORMATION



Résultats de la détection des entités nommées

EXTRACTION D'INFORMATION

Grand-Duché de Luxembourg

cartographyCode: LUX
about: http://data.loterre.fr

about: http://data.loterre.fr/ark:/67375/9SD-HXXBRCFQ-F

prefLabel@fr: Luxembourg
prefLabel@en: Luxembourg

id: grandduchedeluxembourg

wikidataURI: https://www.wikidata.org/wiki/Q1842

geonameURI: https://www.geonames.org/2960313
countryCode: LU

latitude: 49.765705074151 longitude: 5.965223432344 localization@en: Western Europe

localization@fr: Europe occidentale

Traitement des entités nommées de lieux : utilisation du web service Pays et subdivision

EXTRACTION D'INFORMATION



Résultats de l'alignement avec Loterre

EXTRACTION D'INFORMATION



EXTRACTION D'INFORMATION







EXTRACTION D'INFORMATION

Rapport entre les mots-clés et pays cités

'I don't eat meat': Discourse on food among transnational Hindus

Jennifer B. Saunders

Hindu transmigrants use discourse on diet as a way to maintain connections with India, as well as to construct Indian, Hindu and caste identities. In this article, I argue that such discourse on food is a meta-discourse that reframes the symbolic meaning of food in the transnational context. This article examines a transnational Hindu community's discourse on food, and pairs R.S. Khare's arguments about the communicative function of food in a South Asian context with transnational and performance theories, as well as with Arjun Appadurai's argument about the significance of imagination in creating lived realities. Through their narratives involving food, this community is actively engaged in shifting the meanings of what it eats to emphasise their connections with each other, and with India. Thus, a vegetarian diet and the use of 'authentic' Indian ingredients become the symbols of Indian identity through discourse, which is then solidified through the acts of cooking and active Tries is board on feldovor conducted with an extended transmissional Hindu

Exemple d'article en sociologie



Une méthodologie itérative : examen des données ↔ révision de la requête

• La pertinence scientifique : circonscrire au maximum le sujet quitte à créer du silence?

Une méthodologie itérative : examen des données ↔ révision de la requête

- La pertinence scientifique : circonscrire au maximum le sujet quitte à créer du silence?
- La pertinence technique : se limiter aux documents possédant un résumé, un PDF texte, un texte nettoyé, un nombre de mots / pages raisonnable?

Une méthodologie itérative : examen des données ↔ révision de la requête

- La pertinence scientifique : circonscrire au maximum le sujet quitte à créer du silence?
- La pertinence technique : se limiter aux documents possédant un résumé, un PDF texte, un texte nettoyé, un nombre de mots / pages raisonnable?

Des aides précieuses développées dans Istex

• L'interface Istex Search

Une méthodologie itérative : examen des données ↔ révision de la requête

- La pertinence scientifique : circonscrire au maximum le sujet quitte à créer du silence?
- La pertinence technique : se limiter aux documents possédant un résumé, un PDF texte, un texte nettoyé, un nombre de mots / pages raisonnable?

Des aides précieuses développées dans Istex

- L'interface Istex Search
- L'appui de Lodex et des web services



Annexes et références

Les outils Istex

ISTEX	Istex	Site général, vitrine du possible	www.istex.fr
~	Data	Informations sur les données Istex (collections, enrichissements) & accès aux corpus (Cuxac <i>et al.</i> , 2019)	data.istex.fr
	Search	Application pour interroger l'API et télécharger des corpus	search.istex.fr
	Lodex	Outil de création de sites web (visualisation / enrichissement de données) (Gregorio <i>et al.</i> , 2019)	www.lodex.fr
	TDM	Bibliothèque de web services (Cuxac, 2024)	services.istex.fr
	Lo- terre	Bibliothèque de terminologies multilingues (Khayari <i>et al.</i> , 2021)	www.loterre.fr

Documentation & tutoriels

Se documenter

- Documentation Istex: https://doc.istex.fr/tdm
- Documentation API Istex: https://doc.istex.fr/api
- Documentation Lodex: https://www.lodex.fr/docs/documentation/

S'authentifier

• Vérifier ses droits d'accès: https://api.istex.fr/auth

Informations & contacts

Se tenir informé

• Article d'actualité: https://www.istex.fr/category/actualites/

Chercher de l'aide / Contribuer à l'amélioration

- Via le formulaire: https://www.istex.fr/contact/
- Via les listes: contact@listes.istex.fr, https://groupes.renater.fr/sympa/subscribe/lodex?previous_action=info

- BENZITOUN, C. et CAPPEAU, P. (2025). Les corpus et leur exploitation. *Encyclopédie Grammaticale du Français*, En ligne. Artwork Size: 995084 bytes Medium: application/pdf Publisher: NAKALA https://nakala.fr (Huma-Num CNRS) Version Number: 1.
- BOUCHARD, A. (2024). Sortir de Google Scholar, Scopus ou Web of Science : que valent Lens, Dimensions, OpenAlex et Matilda pour la recherche bibliographique? [à distance] | URFIST de Paris.
- BOUSSIDAN, A., DUCEL, F., NÉVÉOL, A. et FORT, K. (2024). What ChatGPT tells us about ourselves.
- BÉNARD, M., MESTIVIER, A., KUBLER, N., ZHU, L., BAWDEN, R., DE LA CLERGERIE, E., ROMARY, L., HUGUIN, M., NOMINÉ, J.-F., PENG, Z. et YVON, F. (2023). MaTOS: Traduction automatique pour la science ouverte. In BOUDIN, F., DAILLE, B., DUFOUR, R., EL, O., HOUBRE, M., JOURDAN, L. et KOOLI, N., éditeurs: Actes de CORIA-TALN 2023. Actes de l'atelier "Analyse et Recherche de Textes Scientifiques" (ARTS)@TALN 2023, pages 8–15, Paris, France. ATALA.

CUXAC, P. (2024). La fouille de textes en IST : les outils Istex-TDM. Nancy, France.

- CUXAC, P., COLLIGNON, A., GREGORIO, S. et PARMENTIER, F. (2019). Des bases de données massives au Web de données : désambiguïsation et alignement d'entités géographiques dans les textes scientifiques. In 12ème Colloque international d'ISKO-France : Données et mégadonnées ouvertes en SHS : de nouveaux enjeux pour l'état et l'organisation des connaissances?
- GREGORIO, S., COLLIGNON, A., PARMENTIER, F. et THOUVENIN, N. (2019). LODEX : des données structurées au web sémantique. In Atelier Web des Données de la 19ème Conférence sur l'Extraction et la Gestion des Connaissances (EGC 2019), Metz, France.
- GUSENBAUER, M. (2022). Search where you will find most: Comparing the disciplinary coverage of 56 bibliographic databases. *Scientometrics*, 127(5):2683–2745.
- HABERT, B. (2000). Des corpus représentatifs : de quoi, pour quoi, comment? *In* BILGER, M., éditeur : Linguistique sur corpus. Etudes et réflexions - Mireille Bilger, pages 11–58. Presses Universitaires de Perpignan, Perpignan.
- HENDRYCKS, D., BURNS, C., BASART, S., ZOU, A., MAZEIKA, M., SONG, D. et STEINHARDT, J. (2021). Measuring massive multitask language understanding.

- KHAYARI, M., RESZETKO, V., VACHEZ, D., VEDOVOTTO, N., YON, J. et AUBIN, S. (2021). De TermSciences à Loterre : comment l'Inist-CNRS a rendu les terminologies ouvertes plus conformes aux principes FAIR.
- LABADIE, A. et PRINCE, V. (2008). The impact of corpus quality and type on topic based text segmentation evaluation. *In 2008 International Multiconference on Computer Science and Information Technology*, pages 313–319.
- LARIVIÈRE, V., GINGRAS, Y. et ARCHAMBAULT, E. (2009). The decline in the concentration of citations, 1900–2007. *Journal of the American Society for Information Science and Technology*, 60(4):858–862.
- LIU, Y., CAO, J., LIU, C., DING, K. et JIN, L. (2024). Datasets for Large Language Models: A Comprehensive Survey. arXiv:2402.18041 [cs].
- LÉON, J. (2008). Aux sources de la « Corpus Linguistics » : Firth et la London School. *Langages*, 171(3) :12–33. Place : Paris Publisher : Armand Colin.
- MELLET, S. (2002). Corpus et recherches linguistiques: Introduction. Corpus. (1).

- RASTIER, F. (2004). Enjeux épistémologiques de la linguistique de corpus. *In* WILLIAM, G., éditeur : *La linguistique de corpus.*, pages 31–46. Presses Universitaires de Rennes, Rennes.
- SINCLAIR, J. (1996). Preliminary recommendations on Corpus Typology. Rapport préliminaire, Expert Advisory Group on Language Engineering Standards.
- STIENNON, N., OUYANG, L., WU, J., ZIEGLER, D. M., LOWE, R., VOSS, C., RADFORD, A., AMODEI, D. et CHRISTIANO, P. (2022). Learning to summarize from human feedback.