

# Usage des ressources ISTEX Ateliers pratiques





#### Programme

1. Présentation générale d'ISTEX

- 2. Constitution d'un corpus par interrogation d'ISTEX
- 3. Téléchargement du corpus avec ISTEX-DL
- 4. Exploration et valorisation du corpus avec l'outil LODEX

--- Une pause ! -----

- 5. Liens utiles
- 6. Annexes



#### **Documents utiles**

Télécharger ce support et les fichiers utiles aux TP

https://formation.lodex.fr



. Enregistrer le zip contenant les documents de travail : clic droit, puis" enregistrer la cible du lien sous" → Bureau ou dossier créé

2. extraire les 5 fichiers : clic droit, puis "extraire tout" → Bureau ou dossier créé

١.

Index of /				
/ Documents-atelier.zip		04-Mar-2020 08:	36	57956 <mark>0</mark> 9:
	0			



# **Introduction à ISTEX**





### ISTEX...

Initiative d'excellence en Information Scientifique et Technique

pour

Construire le socle de la bibliothèque scientifique numérique nationale





Construire le socle de la bibliothèque scientifique numérique nationale.

est né d'une impulsion nationale, dans le cadre du programme « Investissements d'avenir » du Ministère de l'Enseignement Supérieur, de la Recherche et de l'Innovation (MESRI), et financé par l'Agence Nationale de la Recherche (ANR)

est porté par quatre partenaires : CNRS, ABES, Couperin.org, l'Université de Lorraine pour la Conférence des Présidents d'Université (CPU). <u>En savoir plus</u>



SUPÉRIEUR ET DI

est lancé le 30

avril 2012

abes e



est réservé aux abonnés de

l'Enseignement supérieur et de la

recherche française (ESR).









6



de la bibliothèque scientifique numérique nationale.



#### **ISTEX : quels objectifs ?**

- Acquisition massive et centralisée d'archives scientifiques
  - collections rétrospectives de la littérature scientifique
  - toutes disciplines
- Construction d'une plateforme nationale d'hébergement et de mise à disposition des données (Inist)

https://www.istex.fr/



**Construire le socle** de la bibliothèque scientifique numérique nationale.



#### La plateforme ISTEX : quels services ?





de la bibliothèque scientifique numérique nationale.



#### La plateforme ISTEX : quels usages ?

Deux usages : l'un classique, l'autre plus avancé !

- Usage documentaire
- Usage en fouille de textes (text mining)



La fouille de textes désigne toute technique de <mark>traitement automatisé</mark> visant à <mark>extraire des connaissances</mark> d'un ensemble de textes sous forme numérique, produits par des humains pour des humains. Elle permet d'analyser parallèlement de <mark>vastes quantités de données</mark> selon un critère de nouveauté ou de similarité, et ainsi de dégager des informations de haute qualité difficiles à appréhender par la simple lecture cursive, telles que des constantes, des tendances et des corrélations.

(Sources : https://cnnumerique.fr/files/2017-10/CNNum\_Fiche\_TDM.pdf ;

https://fr.wikipedia.org/wiki/Fouille\_de\_textes;

https://adbu.fr/competplug/uploads/2016/12/v9-Designed-Exec-Summary-ADBU-8pp\_fra\_final.pdf)



# Quel contenu ? Quelques chiffres...



**ISTEX** aujourd'hui, c'est :

# **23 207 399**

documents





#### Le fonds ISTEX > répartition par provenance

*NB : Documents provenant de 27 corpus éditeurs mais chiffres des graphiques datant du 12/02/2020* 





#### Le fonds ISTEX > répartition par types de publication



Prépondérance des documents issus de revues

Définitions et chiffres sur data.istex.fr : Types de publication



#### Le fonds ISTEX > répartition par types de contenu

Types de documents Plus de 66% des documents sont des brief-communication articles conference 3,3% abstract 2,1% 3.6% review-article 4.0% research-article book-reviews 44,2% 5,3% Définitions et chiffres sur other data.istex.fr : 12,8% Types de contenu article 22,3%



#### Le fonds ISTEX > répartition par revues

#### Chiffres du 12-02-2020 424895 Dans le fonds de The Lancet 394747 Notes and Oueries plus de 9 000 revues ChemInform 394471 365272 Nature présentes dans **British Medical Journal** 287487 164333 Journal of the American Chemical Society ISTEX : 103845 **Analytical Chemistry** 94155 **Angewandte Chemie** 74468 The Journal of Organic Chemistry liste des **20** revues **Tetrahedron Letters** 67414 63089 The American Historical Review les plus importantes Annals of the New York Academy of Sciences 63316 Monthly Notices of the Royal Astronomical Society 59186 en **nombre de Biochemistry** 57959 documents Fresenius' Zeitschrift für analytische Chemie 54554 **Biochemical and Biophysical Research** 52662 51980 International Journal of Rock Mechanics and 51903 The Journal of Physical Chemistry 51668 Journal of the Franklin Institute Journal of Applied Polymer Science 47916 100000 200000 300000 400000 500000

#### Vingt premières revues



#### Le fonds ISTEX > répartition par domaines scientifiques - 1

**Scopus** est la classification la plus représentative des documents lstex, sachant que plus de 20% des documents ne comportent aucune catégorie scientifique



Nombre de documents Istex (en millions) et pourcentages



#### Le fonds ISTEX > répartition par domaines scientifiques - 2

54% des documents Istex sont classés avec une catégorie Scopus relative aux sciences physiques ou de la santé





#### Le fonds ISTEX > répartition par dates - 1

95% du fonds a été publié aux 20e et 21e siècles



#### 700 ans de publications



#### Le fonds ISTEX > répartition par dates - 2

95% des documents publiés entre 1900 et aujourd'hui (2017)

88% des documents publiés depuis 1950

53% des documents publiés sur les 30 dernières années

#### Dates de publication : par décennie (20&21e siècles)

Chiffres du 12-02-2020





#### Le fonds ISTEX > répartition par langues - 1

53 langues !

Anglais majoritaire

Information non renseignée par les éditeurs pour près de 2,2 millions de documents !





#### Le fonds ISTEX > répartition par langues - 2

**Autres langues** Chiffres du 12-02-2020 italien latin Les 49 autres langues espagnol 10,8% concernent 0,3% des 30,2% néerlandais documents restants (soit près de 57 000 russe documents) gallois 23,5% romanes (langues) portugais 26,0% 44 autres



#### Le fonds ISTEX > répartition par éléments descriptifs du contenu





#### Le fonds ISTEX > répartition par indicateurs - 1

Plus la version de PDF est élevée, meilleure devrait être la qualité du document

Près de 88% des documents existent dans une version de PDF comprise entre 1.2 et 1.4





#### Le fonds ISTEX > répartition par indicateurs - 2

Un nombre de mots réduit peut indiquer un PDF image ou bien un court résumé

La majorité des documents PDF (54%) ont un nombre de mots compris entre 1 000 et 5 000



#### Nombre de mots par PDF



#### Le fonds ISTEX > répartition par indicateurs - 3

Le score de qualité est un nombre décimal compris entre 0 et 10.

Son mode de calcul tient compte de la version du PDF, du nombre de mots du résumé, ainsi que du nombre de mots du texte intégral (PDF).

Voir la formule dans : <u>Calcul du</u> <u>score de qualité</u>





54 % des documents (12 475 415) ont un score supérieur à 5



#### La plateforme ISTEX > la chaîne de traitements



#### La plateforme ISTEX : les enrichissements

#### Junitex

Détection de 9 types d'entités nommées (lieux géographiques, organismes financeurs, etc.)

enrichment-process.data.istex.fr/unitex

⇒ <mark>68,6%</mark> des docs



#### Bayésien naïf (nb)

Catégorisation de l'article par apprentissage avec Pascal/Francis

enrichment-process.data.istex.fr/nb

#### $\Rightarrow$ **43,4%** des docs

Chiffres du 04-02-2020



#### Grobid

Structuration des références bibliographiques

enrichment-process.data.istex.fr/grobid

⇒ **58,1%** des docs



#### Teeft

Indexation des articles par extraction de termes à partir du texte intégral en anglais enrichment-process.data.istex.fr/teeft

#### ⇒ <mark>72,6%</mark> des docs



#### Multicat

Catégorisation de la revue par appariement (WoS / Science Metrix / Scopus)

enrichment-process.data.istex.fr/wos

enrichment-process.data.istex.fr/science metrix

enrichment-process.data.istex.fr/scopus

⇒ <mark>75,2%</mark> des docs

28

# STATEST SEALER

#### La plateforme ISTEX : deux usages



# Exemples d'usage documentaire



- Accéder à un article
- Rechercher les articles d'un auteur
- Constituer une bibliographie sur un domaine
  - Produire une synthèse documentaire

# STISSEARCH BURNES

#### La plateforme ISTEX : deux usages



## Exemples d'usage 🔅

- Aider au pilotage scientifique & au travail du chercheur
- Répondre à une question scientifique
- Usages hors Recherche

#### La plateforme ISTEX : exemples d'usage TDM\*

#### Pilotage & Recherche

- Détection de sujets de recherche émergents
- Détection de plagiat
- Recommandation d'articles pour des néophytes ou non

\* Ces cas d'usage sont documentés dans le fichier "Exemples-usage-TDM.pdf" présent dans le zip "Documents-atelier"

#### **Question scientifique**

- Applications inattendues pour un ancien médicament
- Cartographie du génome humain
- Nouveaux matériaux thermoélectriques
- Traitement médiatique des "Gilets Jaunes"
- Portrait-robot d'espèces de poissons sauvages domesticables



#### Autres usages

- Marketing ciblé : détection d'intention d'achats et personnalisation des réponses
- Aide au diagnostic médical
- Prévention de la cybercriminalité
- Gestion des risques de fuite dans les réseaux d'assainissement
- Traduction automatique
- Assistant personnel



#### La plateforme ISTEX : exemples d'usage TDM à l'Inist

#### Réaliser des analyses scientifiques

- Analyses syntaxiques et terminologiques
- Evolution temporelle de thématiques

#### Mettre au point ou tester des outils

- Entity-Fishing/IRC3 : détection de noms d'espèces animales et végétales
- Unitex : détection de définitions de concepts scientifiques
- TermSuite : détection de termes et de leurs variantes
- TermSuite/Teeft + Neurodoc : constitution et analyse de corpus
- Topic Modeling, Iramuteq : exploration thématique

#### Évaluer la performance d'outils

**Unitex** : Corpus gold annoté manuellement en entités nommées pour évaluer les résultats de détection de l'outil



Accessibles sur data.istex.fr

#### La plateforme ISTEX : exemples d'usage TDM à l'Inist

#### Corpus réalisés à l'Inist :

- Sur des **domaines** et des **thématiques** variés :
  - Astrophysique, orthophonie, botanique, zoologie, géosciences, sciences et techniques alimentaires
  - Zone polaire arctique, vieillissement
- De petits et de gros corpus :
  - De **34** à **52 238** documents



# **Interrogation d'ISTEX**

#### Interroger ISTEX > Comprendre la structuration du réservoir



### Le réservoir ISTEX contient des **objets documentaires**

Ces objets documentaires sont :

- Soit des **articles** (de revue)
- Soit des chapitres (de monographie, d'e-book, de monographies en série)

#### Différents niveaux de granularité

L'interrogation peut se faire sur des informations concernant :

- Le niveau article ou chapitre
- Le niveau monographie
- Le niveau série de monographies



\* Ces informations seront identiques pour tous les articles d'une même revue tous les chapitres d'une même monographie



36

#### Où interroger ISTEX ?

Le démonstrateur est une interface à **vocation pédagogique** branchée sur l'API ISTEX qui permet de :

- construire sa requête (en mode simple ou en mode avancé)
- visualiser et filtrer les résultats

https://demo.istex.fr






### Facettes pré-définies dans l'interface



→ donne une vision synthétique du corpus
 → permet de filtrer les résultats de la requête
 → mais possibilités limitées - exploratoire











## **Interroger ISTEX > la recherche basique**

#### Petit exercice pour démarrer :

• On souhaite rechercher les documents possédant le terme "polar"

Bienvenue sur le démonstrate En savoir plus polar Recherche avancée		Recherche sur : - Les métadonnées - Le texte intégral - Les références bibliographiques - Les enrichissements
Résultats	: + de 1 000 000 docs !!!	



# **Cas d'usage**

"Zoologiste spécialisé en milieu marin au département "Adaptations du vivant" du MNHN, je m'intéresse à l'adaptation au changement global en zone polaire"

**Objectif pédagogique :** 

Écrire une équation, testée pas à pas, utilisant un certain nombre d'opérateurs, d'astuces et de syntaxes, associés à une recherche sur un certain nombre de champs utiles pour délimiter un corpus pertinent et de taille raisonnable



## Interroger ISTEX > Cas d'usage

## **Application TDM**

- Outil de reconnaissance d'**entités** nommées scientifiques
- Détecter les espèces animales
- Dans le texte intégral

## Thématique

- zoologie : **espèces animales** de type **poisson ou mollusque**
- zone géographique : zone polaire ou subpolaire arctique

## Contraintes liées à l'application

1

- L'outil traite les documents en anglais
- Les documents devront être en format PDF
- Attention : l'outil ne traite pas les PDF image

Rédiger sa requête dans un bloc-note type "NotePad". Ou copier successivement les exemples de "**TP.txt**", et coller les requêtes dans le démonstrateur

## 1. je recherche les termes suivants

- fish, mollusc
- arctic, subarctic, Svalbard, Barents sea

## 2. dans les champs

- titre
- résumé
- mots-clés d'auteurs
- 3. je limite mon corpus aux documents
  - langue de publication = anglais



## Notions abordées **Opérateurs booléens** Parenthésage **Recherche sur champs** Recherche insensible à la casse Expressions multitermes / mots-composés



## Equation étape 1

Recherche sur des termes de zoologie

fish **OR** mollusc

#### Explications :

- OR cumule les résultats associés à "fish" et "mollusc"
- Les mots sont recherchés sur tout le document
- Si pas d'opérateur utilisé, l'opérateur par défaut OR s'applique

Les opérateurs doivent s'écrire **en MAJUSCULES** 



## Equation étape 1

Recherche sur des termes de zoologie et de géographie

(fish OR mollusc ) AND

(arctic OR subarctic )

- OR cumule les résultats associés à "fish" et "mollusc"
- Les mots sont recherchés sur tout le document
- Si pas d'opérateur utilisé, l'opérateur par défaut OR s'applique
- AND + parenthèses permet d'associer 2 critères, ici termes zoologiques et termes géographiques



## Equation étape 1

Recherche sur des termes de zoologie et de géographie contenant **majuscule** 

(fish OR mollusc ) AND

(arctic OR subarctic OR svalbard )

- OR cumule les résultats associés à "fish" et "mollusc"
- Les mots sont recherchés sur tout le document
- Si pas d'opérateur utilisé, l'opérateur par défaut OR s'applique
- AND + parenthèses permet d'associer 2 critères, ici termes zoologiques et termes géographiques
- insensibilité à la casse
  - "svalbard" cherche aussi "Svalbard"



## Equation étape 1

Recherche sur des termes de zoologie et de géographie contenant **majuscule, expression multiterme** 

```
(fish OR mollusc ) AND
```

```
(arctic OR subarctic OR svalbard OR "barents sea")
```

- OR cumule les résultats associés à "fish" et "mollusc"
- Les mots sont recherchés sur tout le document
- Si pas d'opérateur utilisé, l'opérateur par défaut OR s'applique
- AND + parenthèses permet d'associer 2 critères, ici termes zoologiques et termes géographiques
- insensibilité à la casse
  - "svalbard" cherche aussi "Svalbard"
- utilisation d'une expression multiterme entre guillemets
  - **barents sea** = 1 346 725 docs
  - "barents sea" = 12 752 docs



## Equation étape 1

Recherche sur des termes de zoologie et de géographie contenant **majuscule, expression multiterme, mot-composé** 

```
( fish OR mollusc ) AND
( arctic OR subarctic OR "sub-arctic"
  OR svalbard OR "barents sea")
```

#### Résultats (10 mars 2020) : 52 885 docs

- OR cumule les résultats associés à "fish" et "mollusc"
- Les mots sont recherchés sur tout le document
- Si pas d'opérateur utilisé, l'opérateur par défaut OR s'applique
- AND + parenthèses permet d'associer 2 critères, ici termes zoologiques et termes géographiques
- insensibilité à la casse
  - "svalbard" cherche aussi "Svalbard"
- utilisation d'une expression multiterme entre guillemets
  - barents sea = 1 347 835 docs
  - "barents sea" = 12 778 docs
- utilisation d'un **mot-composé** entre guillemets
  - **sub-arctic** = 2 116 617 docs
  - "sub-arctic" = 7 808 docs



## Equation étape 1

Recherche sur des termes de zoologie et de géographie contenant **majuscule, expression multiterme, mot-composé** 

#### *Résultats (10 mars 2020) : 52 885 docs*

#### Explications :

- OR cumule les résultats associés à "fish" et "mollusc"
- Les mots sont recherchés sur tout le document
- Si pas d'opérateur utilisé, l'opérateur par défaut OR s'applique
- AND + parenthèses permet d'associer 2 critères, ici termes zoologiques et termes géographiques
- insensibilité à la casse
  - "svalbard" cherche aussi "Svalbard"
- utilisation d'une expression multiterme entre guillemets
  - $\circ$  barents sea = 1 347 835 docs
  - "barents sea" = 12 778 docs
- utilisation d'un mot-composé entre guillemets
  - **sub-arctic** = 2 116 617 docs
  - "sub-arctic" = 7 808 docs

#### Plus de détails : <u>opérateurs</u> / <u>astuces</u>



## Equation étape 2

#### Recherche restreinte sur certains champs

```
title:((fish OR mollusc) AND (arctic OR
subarctic OR "sub-arctic" OR svalbard OR
"barents sea"))
```

```
OR abstract<mark>:</mark>(idem)
```

```
OR subject.value: (idem)
```

Résultats (10 mars 2020) : 52 885 docs > 1 243 docs

#### Explications :

- Recherche restreinte sur les champs :
  - titre de l'article : "title"
  - résumé : "abstract"
  - mots-clés d'auteur : "subject.value"
- Les noms de champs sont introduits par :
- Ajout de **parenthèses** pour regrouper les valeurs recherchées dans un même champ
- Pas de factorisation des noms de champs. Il faut répéter les termes pour chaque champ interrogé

Plus de détails : <u>exemples de contenus</u> / <u>recherche</u> <u>sur champs</u>



## Equation étape 3

Restriction aux documents publiés en anglais

(idem) AND language:eng

#### Explications :

- Ajout de parenthèses à l'ensemble de l'équation car on ajoute un critère qui porte sur toute l'équation précédente
- Recherche sur le champ : language

Liste des <u>codes langues</u>

Résultats (10 mars 2020) : 1 243 docs >> 1 133 docs



## Équation complète à cette étape

(title:((fish OR mollusc) AND (arctic OR subarctic OR svalbard OR "sub-arctic" OR "barents sea"))
OR abstract:((fish OR mollusc) AND (arctic OR subarctic OR svalbard OR "sub-arctic" OR "barents sea"))
OR subject.value:((fish OR mollusc) AND (arctic OR subarctic OR svalbard OR "sub-arctic" OR "barents sea")))

AND language:eng

Résultats (10 mars 2020) : 1 133 docs



- 4. je limite mon corpus aux documents - gui ne sont pas des PDF image
- 5. je cherche à ajouter un maximum de variantes pertinentes aux termes suivants
  - fish, mollusc
  - svalbard

Les variantes peuvent être de plusieurs atures :

- formes au **pluriel**
- formes non accentuées
- variantes **orthographiques**
- formes **composées**

#### Notions abordées

Interrogation sur des intervalles de nombres

**Utilisation de troncatures** 

Utilisation d'expressions régulières



## Equation étape 4

Un petit test pour commencer : Exploration de la facette "Qualité" dans le démonstrateur

ualité <del>-</del>	Status of Zoobenthos and Fish Populations in Subarctic Rivers of the Northernmost Finland:.     Abstract In 1990, a monitoring programme was initiated to survey the status of benthic invertebrate     communities and fish populations over a wide range of subarctic rivers in northernmost Finnish Lapland     (68°15′-70°N) A special emphasis was placed on detecting possible effects of aridification through suffur     research-article
Score	Fulltext Metadata Enrichments
Entre: 2.434 à 9.934	POT ZP XML MCCS TTE 2001
lombre de mots	
Entre : 0 à 174600	3 - Ouverture du PDF 2 - Vérification du et sélection du texte nombre de mots
	1 - Réglage du curseur à 500



## Equation étape 4

Recherche restreinte sur certains champs avec des **intervalles** de valeurs

qualityIndicators.pdfWordCount:[500 TO 10000]

qualityIndicators.pdfVersion:[1.2 TO \*]

**Explications** :

Recherche sur des intervalles de valeurs à l'aide de crochets []

[500 TO 10000] ← valeurs limites inférieure et supérieure

 [1.2 TO \*]
 ←

 valeur limite supérieure infinie

TO s'écrit obligatoirement en MAJUSCULES

Plus de détails : <u>intervalles</u>



## Equation étape 5

Recherche sur des formes au pluriel et des formes non accentuées à l'aide de **troncatures** 

fish\* OR mollus\*

#### Explications :

- ? remplace **1** caractère
- \* remplace **0 à n** caractère(s)

**mollus\*** = mollusque(s) mollusca mollusc(s) mollusk(s) molluscan...

**fish\*** = fish fishes fishing fishery fisheries fisherman fishermen... et Tully-Fisher ! **1** variante hors domaine

Plus de détails : <u>troncatures</u>



## Equation étape 5 (suite)

Recherche plus ciblée sur des formes au pluriel à l'aide d'**expressions régulières** 

#### /fish(es)?/



#### **Explications**:

La troncature par \* ne convient pas ici car ramène des variantes hors domaine

Les expressions régulières permettent de rendre compte de toutes les possibilités d'écriture d'un terme en ciblant des variantes spécifiques





## Equation étape 5 (suite)

On peut aller beaucoup plus loin avec la recherche sur des termes avec variantes orthographiques

> Exemple : Spitzberg, île principale de l'archipel du Svalbard, s'écrit avec 4 orthographes possibles

svalbard OR /spit[sz]berg(en)?/





## **Interroger ISTEX > Niveau Intermédiaire - affiner les valeurs** Équation complète à cette étape

```
(title:((/fish(es)?/ OR mollus*) AND (arctic OR subarctic OR "sub-arctic" OR svalbard OR
/spit[sz]berg(en)?/ OR "barents sea"))
```

```
OR abstract:((/fish(es)?/ OR mollus*) AND (arctic OR subarctic OR "sub-arctic" OR svalbard OR
/spit[sz]berg(en)?/ OR "barents sea"))
```

```
OR subject.value:((/fish(es)?/ OR mollus*) AND (arctic OR subarctic OR "sub-arctic" OR svalbard OR
/spit[sz]berg(en)?/ OR "barents sea")))
```

```
AND language:(eng)
```

AND qualityIndicators.pdfWordCount:[500 TO 10000]

Résultats (10 mars 2020) : 1 133 > 1 192 docs

- je cherche à ajouter des variantes supplémentaires aux termes suivants :
  - subarctic, barents sea



#### Notions abordées

**Recherche floue** 

Recherche de proximité



## Equation étape 6

Recherche sur des variantes d'écriture du terme "subarctic" à l'aide de la **recherche floue** :

"sub-arctic" OR subarctic~1

- permet de **récupérer** des articles **pertinents** bien que contenant des erreurs (numérisation)

- à utiliser **avec prudence** après **vérification** des résultats

#### Explications :

en fonction du nombre entier indiqué après le tilde, compris entre **0 et 2**, on recherche un terme ayant **au maximum 2 caractères de différence** (en plus, en moins ou dissemblables).

- subarctic~O récupère : subarctic
- subarctic~1 récupère : subarctic + subacrtic, subartic, subaratic, subarctica
- subarctic~2 récupère : subarctic, subacrtic, subartic, subaratic, subarctica + subaratica mais aussi subarcsec, subaortic

NB : subarctic~ (sans chiffre précisé) équivaut à subarctic~2

Plus de détails : <u>recherche floue</u>



## Equation étape 6 (suite)

Recherche sur des synonymes de "barents sea" à l'aide d'**opérateurs de proximité** 

"barents sea" OR "barents seas"~2

#### **Explications :**

en fonction du nombre entier indiqué après le tilde, compris entre **0 et n**, on recherche une expression comportant deux termes **plus ou moins distants** (distance de **n termes**).

"barents seas"~2 récupère :

- Barents-Kara Seas
- Barents and Kara seas
- Barents, and Kara seas



Le moteur de recherche ne tient pas compte de l'ordre des termes et recherche les termes dans les 2 sens

NB : pas de troncature possible

Plus de détails : recherche de proximité



## Équation complète

(title:((/fish(es)?/ OR mollus\*) AND (arctic OR "sub-arctic" OR subarctic~1 OR svalbard OR
/spit[sz]berg(en)?/ OR "barents sea" OR "barents seas"~2))

OR abstract:((/fish(es)?/ OR mollus\*) AND (arctic OR "sub-arctic" OR subarctic~1 OR svalbard OR
/spit[sz]berg(en)?/ OR "barents sea" OR "barents seas"~2))

OR subject.value:((/fish(es)?/ OR mollus\*) AND (arctic OR "sub-arctic" OR subarctic~1 OR svalbard OR
/spit[sz]berg(en)?/ OR "barents sea" OR "barents seas"~2)))

AND language:eng

AND qualityIndicators.pdfWordCount:[500 TO 10000]

Résultats (10 mars 2020) : 1 192 > 1 197 docs



## **Interroger ISTEX > Et après ?**

# Téléchargement du corpus avec ISTEX-DL

Exploration de l'interface

Réutilisation de la requête constituée

Sélection du format de fichier en fonction du cas d'usage

## Exploration du corpus dans LODEX S

Import du fichier dans Lodex Modélisation à l'aide du modèle prédéfini Exploration des différents graphiques



## **Téléchargement de corpus** ISTEX-DL & Autres outils



## ISTEX-DL ou ISTEX-DownLoad : "télécharger un corpus ISTEX en quelques clics"

"





## **ISTEX-DL** : Objectif

Interface web conviviale de type formulaire pour **extraire un corpus ISTEX** en **3 étapes** et récupérer son corpus sur son poste de travail dans un **format compressé (zip)**...avec un **minimum** de connaissance informatique !



## TROIS ÉTAPES,...



A

A

## construction d'une équation

## sélection des fichiers

### téléchargement du corpus

**ISTEX-DL** 

Requête Explicitez ci-dessous l'équation ou la liste d'identifiants qui décrit le corpus souhaité : Recherche classique () Recherche par ARK () ark:/67375/0T8-IMF4G14B-2 ark:/67375/0T8-RNCBH0V7-8 Nombre de caractères restants 🚯 : 67 000 🖌 Choisir le nombre de documents souhaités 🜖 : 50000 \$ Choisir les documents classés 🜖 Par pertinence 
Aléatoirement 2. Formats et types de fichiers Créez votre sélection en cochant ou décochant les cases ci-dessous : Texte intégral Métadonnées 🚯 Enrichissements PDF JSON multicat TEI XML nb TXT MODS refBibs ZIP teeft Annexes ① TIFE unitex Couvertures 🚺 3. Télécharger 1



## **Deux modes** de recherche

de recherche	Explicitez ci-dessous l'équation ou la liste d'identifiants qui décrit le	e corpus souhaité :
	R Recherche classique Recherche par ARK	0
	ark:/67375/0T8-JMF4G14B-2 ark:/67375/0T8-RNCBH0VZ-8	
	Nombre de caractères restants 🚯 : 67 000 🖌	
	Choisir le nombre de documents souhaités 🟮 : 50000 🗘	50000 0 100000
	Choisir les documents classés 🟮 :	

lacksquare



## Recherche par identifiants des documents

1. Requête		ARK attrib docu	(équivalent du DOI) oué à chaque ment ISTEX	
Explicitez ci-dessous l'équation o	ou la liste d'identifiants qui décrit le	corpus souhaité :		
Recherche classique <b>1</b>	Recherche par ARK <b>1</b>			
ark:/67375/0T8-JMF4G14B-2 ark:/67375/0T8-RNCBH0VZ-8				
Nombre de caractères restants	③:67 000 ✓			
Choisir le nombre de document	s souhaités 🚺 : 50000 🗘	50000	100000	
Choisir les documents classés				
citoton los accantentes classes				

## **ISTEX-DL** : étape 1
# AVENIE

6

## Recherche par identifiants des documents



**ISTEX-DL** : étape 1

## 1. Requête

Explicitez ci-dessous l'équation ou la liste d'identifiants qui décrit le corpus souhaité :

ark:/67375/018-JMF4G14B-2 ark:/67375/0T8-RNCBH0VZ-8		Réinitialiser
lombre de caractères restants	1 : 67 000 🖌	
Choisir le nombre de document	ts souhaités 🕄 : 50000 📮 🚽	no c Récupérer
Choisir les documents classés ( Par pertinence ) Aléatoir	<b>្ងំ</b> : rement	Partager
		D Historique

## **Recherche par** identifiants des documents





## **Recherche par** identifiants des documents



## **Recherche par** identifiants des documents

Explicitez ci-dessous l'équation	ou la liste d'identifiants qui décrit le corpu	s souhaité :		
Recherche classique <b>1</b>	Recherche par ARK <b>1</b>			
brain AND language:fre				
Nombre de caractères restants	(1) : 67 000 ✓			
Choisir le nombre de documen	ts souhaités 🜖 : 50000 🗘	50000	100000	
Choisir les documents classés	0 :			

## **ISTEX-DL** : étape 1



## Longueur équation : limitation

(nombre & couleur)





## **Recherche par** équation booléenne





## 1. Requête

Explicitez ci-dessous l'équation ou la liste d'identifiants qui décrit le corpus souhaité :

# Recherche classique Nonbre de caractères restants Recherche par ARK <

Syntaxe de la requête incorrecte sur : undefined . Si vous utilisez un caractère & dans votre requête, veuillez le remplacer par son équivalent ASCII : %26. Allez sur https://api.istex.fr/documentation/300-search.html pour plus d'informations sur la syntaxe des requêtes.

## Recherche par équation booléenne



A

## **Recherche par** équation booléenne



Requête



A

## Nombre de documents : limitation



Explicitez ci-dessous l'équation ou la liste d'identifiants qui décrit le corpus souhaité :



## **ISTEX-DL** : étape 1



## Nombre de documents : limitation



## **ISTEX-DL** : étape 1



0

## Nombre de documents : sélection

## **ISTEX-DL** : étape 1

## . Requête

• Par pertinence O Aléatoirement

Explicitez ci-dessous l'équation ou la liste d'identifiants qui décrit le corpus souhaité :

Recherche classique <b>0</b>	Recherche par ARK ()	
arctic AND gualityIndicators.p	dfVersion:{1.3 TO *]	no
Nombre de caractères restants	€ : 66 950 🖌	du
L'équation saisie correspond à	70 981 document(s)	Ca
Choisir le nombre de document	s souhaités 🚺 : 50000	50000 ST
Choisir les documents classés	Ð :	

Avec le sélecteur ou la réglette, adapter le choix du nombre de documents à la durée d'extraction souhaitée et à la capacité de stockage



A

## Deux modes de tri



Explicitez ci-dessous l'équation ou la liste d'identifiants qui décrit le corpus souhaité :

	Recherche classique ①     Recherche par ARK ①
mode par	brain AND language:fre
défaut qui trie selon un score	Nombre de caractères restants 🜖 : 67 000 🖌
de pertinence des documents en réponse à l'	Choisir les documents classés ():
STEX-DL : étage 1	mode à choisir pour constituer un échantillon représentatif (si réduction du nombre de documents)



85

## Recherche par équation booléenne



Explicitez ci-dessous l'équation ou la liste d'identifiants qui décrit le corpus souhaité :

Requête



A

## Sélection

**ISTEX-DL**:

## 2. Formats et types de fichiers

Créez votre sélection en cochant ou décochant les cases ci-dessous :

Texte intégral	Métadonnées 🚺	Enrichissements (1)
PDF	D JSON	multicat
TEI	XML	nb
TXT	MODS	refBibs
ZIP	Annexes (1	teeft
TIFF		unitex
	Couvertures 🚺	



#### Eight glacial cycles from an Antarctic ice core **Un exemple** The Antarctic Vostok ice core provided compelling evidence of the nature of climate, and of climate feedbacks, research-article nature de document over the past 420,000 years. Marine records suggest that the amplitude of climate variability was smaller before ark:/67375/GT4-NBDQR6K2-9 that time, but such records are often poorly resolved. Moreover, it is not possible to infer the abundance of ... Mots : 6131 Covers Metadata Enrichments ulltext Annexes formats & Publication : 2004 <> <> <> XML MODS JPG nultica, nb 8 ZIP fichiers TEI TEI PDF PDF disponibles efBibs TEI XLS DOC {} JSON TEI TEI TXT 4 formats Les 5 types GIF TEI (les plus d'enrichissement courants) présents dans lstex (au format TEI) 3 formats 5 annexes 1 couverture toujours disponibles (format image) présents (avec des ISTEX-DL : étape 2 disponible formats différents)







# Quels fichiers choisir ?

## Aide des infobulles



## Taper pour rechercher Documentation ISTEX Usage documentaire d'ISTEX Usage TDM d'ISTEX Construction d'une requête Extraction d'un corpus Vérification et mise en forme des r... Annexes Liste des codes langues dans IS... Liste des types de publication e... Liste des catégories scientifique... Liste des formats et types de ... API ISTEX Intégrer ISTEX dans un ENT FAO Published with GitBook

#### ISTEX

#### Texte intégral

#### o PDF (Portable Document Format \*Format de document portable)

Il s'agit d'un format de description de pages pouvant contenir du texte, des dessins, des images et photographies (noir et blanc, couleur, 3D). C'est un format ouvert, évolutif et multiplateforme, issu de l'imprimerie, qui conserve la mise en page du document original. Il offre une sécurité permettant à l'auteur d'un document d'empêcher sa modification par des utilisateurs. Il a été créé par Adobe Systems, Inc. Le logiciel Adobe® Acrobat® Reader est nécessaire pour lire et imprimer un fichier PDF.

Les fichiers en format PDF dans ISTEX sont des fichiers originaux fournis par l'éditeur.

#### o ZIP

Il s'agit d'un format permettant l'archivage et la compression de fichiers. L'archivage est l'utilisation d'un seul fichier pour stocker plusieurs fichiers. La compression des fichiers permet de réduire leur taille. Compresser les fichiers permet de gagner du temps dans le chargement des données et de la place dans le stockage de celles-ci. Le logiciel de compression analyse le fichier et compresse les parties qui se répètent. Lors de la décompression, la forme originale du fichier est restaurée. On peut le comparer à la

90



A

**Quels fichiers** choisir?

Aide des infobulles











**Fermeture** inopinée?

Bouton "Récupérer"









Copier

Annuler

Formulaire complété... 2 boutons s'activent

**Bouton** "Partager"









## S'authentifier si besoin

# mode normal & mode nomade

Authentification dans votre établissement :

• automatique par adresse IP connue

## **ISTEX-DL** : étape 3



#### m Sélectionnez votre établissement Sélection de votre établissemen X < → C ♠ ① A https://discovery.renater.fr/renater/?entityIE Pour accéder au service API ISTEX sélectionnez ou cherchez l'établissement auquel vous appartenez. C Les plus visités Débuter avec Firefox Veuillez entrer le nom de votre établissement. Fédération Éduc Veuillez entrer le nom de votre établissement... ABES - Agence Bibliographique de l'Enseignement Supérieur AUF Agence universitaire de la Francophonie (AUF) AGROCAMPUS OUEST Sélectionnez votre établissem AgroParisTech Institut des sciences et des industries du vivant et de l'environnement Agropolis International Authentification en dehors Agrosup DIJON INSSAAE AMUE - Agence de Mutualisation des Universités et Etablissements de votre établissement : Assistance Publique - Hôpitaux de Marseille bou Bibliothèque Nationale et Universitaire - Strasbourg via la **fédération** Serdeaux INP Bordeaux Sciences Agro d'identité lors du BRGM - Bureau de Recherches Géologiques et Minières CAMPUS CONDORCET premier CEA-Extra téléchargement Centrale Nantes se reconnecter CentraleSupelec CEREQ - Centre d'Etudes et de Recherches sur les Qualifications ensuite à l'application CHU de Lille ISTFX-DI CIHEAM / IAMM CINES CINES CIRAD **ISTEX-DL** : étape 3 CNES CNOUS

S'authentifier si besoin

mode normal & mode nomade







## Fichier zip du corpus extrait



Formation Urfist > Extraction Corpus > istex-subset-2019-05-15.zip



# Fichier zip du corpus extrait

## **ISTEX-DL** : étape 3

"arkIstex": "ark:/67375/VQC-271B662S-H", "title": "Feeding ecology of dominant groundfish in the northern Bering Sea". "abstract": "Abstract: We investigated current diets of the six most abundant benthic fish in the northern Bering Sea. Our objective was to explore feeding strategies and potential competition with other top predators as ecosystem changes occur in the northern Bering Sea ecosystem. Our approach used stomach content data collected from field sampling during spring 2006 and 2007. Calanoid copepods and ampeliscid amphipods were important prev of Arctic cod (Boreogadus saida) but in different proportions depending upon fish size, feeding location, and local environmental conditions. Snailfish (Liparidae) occupied a broad niche and fed on a variety of benthic amphipods. Arctic alligatorfish (Ulcina olrikii) and Arctic staghorn sculpin (Gymnocanthus tricuspis) consumed ampeliscid amphipods predominantly. Shorthorn sculpin (Myoxocephalus scorpius) had a less-diverse diet, with snow crab (Chionoecetes opilio) most important by weight. Finally, all Bering flounder (Hippoglossoides robustus) sampled had empty stomachs. Our results indicate that ampeliscid amphipods, which have high biomass in the central region of the northern Bering Sea, are the most important prev for the dominant groundfish in the Chirikov Basin. Generally, all dominant benthic fish in the northern Bering Sea had narrow feeding niches, except snailfish. High diet overlap was found among many of the fish species, including Arctic cod and snailfish, snailfish and Arctic alligatorfish, and Arctic alligatorfish and Arctic staghorn sculpin. These findings are consistent with a relatively short food chain for benthic fish that are for the most part specialized feeders with narrow preferences for food and may be affected by changes in benthic prey distributions.", "language": [

"eng"
],
"publicationDate": "2012",
"copyrightDate": "2012",
"doi": [
 "10.1007/s00300-012-1180-9"
],
"ark": [
 "ark:/67375/VQC-271B662S-H"
],

"articleId": [ "1180", "s00300-012-1180-9"

], "genre": [ "research-article"

"originalGenre": [ "OriginalPaper" ],

"author": [

**Conserver une** trace & Rejouer

Bouton "Historique"







## **Documentation :**

- Application ISTEX-DL
  - adresse: <u>https://doc.istex.fr/tdm/extraction/istex-dl.html</u>
  - directement via le pied de page

NB : Application en évolution constante...

Vos retours sont bienvenus !



## **Autres Outils : Quels Avantages ?**

## **1 - Fonction "extract"** (API ISTEX) :

à ajouter et à paramétrer directement dans l'URL d'interrogation

## https://api.istex.fr/document/?q=...&extract=...

- option rankBy = qualityOverRelevance
- option rankBy = random avec randomSeed (nombre unique à 13 chiffres)

NB : ISTEX-DL exploite la fonction "extract"



## Documentation

Fonction "extract" :

https://doc.istex.fr/tdm/ext raction/extract-feature.html

## **Autres Outils : Quels Avantages ?**

2 - **Moissonneurs** : utilitaires à installer et à lancer en ligne de commande avec la liste des options choisies :

### "istex-api-harvester" :

- extraction au-delà de 100 000 documents (scroll à paramétrer)
- tous modes d'option rankBy



## Documentation

Tous les détails pour installer et utiliser :

moissonneur
 istex-api-harvester

https://doc.istex.fr/tdm/ext raction/istex-api-harvester.

## **Autres Outils : Quels Avantages ?**

**2 - Moissonneurs** : utilitaires à installer et à lancer en ligne de commande avec la liste des options choisies :

## "harvestCorpus" :

- extraction au-delà de 100 000 documents
- fichier ".corpus" (liste des identifiants des documents)
- statsCorpus : outil de calcul de statistiques associé à cet utilitaire
- tirage aléatoire (avec "graine aléatoire" possible)
- possibilité de renommer les fichiers extraits





## Télécharger un corpus ISTEX > Et après ?

# Exploration du corpus dans LODEX 🔮

Import du fichier dans Lodex Modélisation à l'aide du modèle prédéfini Exploration des différents graphiques

## **Objectifs**:

- Explorer le contenu du corpus pour affiner à nouveau la requête
- Exposer le corpus pour le partager avec la communauté scientifique



## **Exploration et valorisation du corpus** LODEX


#### https://lodex.inist.fr/

#### https://github.com/Inist-CNRS/lodex





111

#### LODEX > on passe à l'action !

- 1. Accéder à une instance de l'outil
  - Adresse : http://formation-urfistn.lodex-dev.inist.fr/
  - Nom d'utilisateur : admin
  - Mot de passe : **secret**

instances numérotées de **2 à 20** 





#### Le jeu de données importé

Loc	<b>Iex</b> - Modélisez et publie	ez vos données vers le web	sémantique !	IMPORTER UN JEU DE DONNÉE	IMPORTER LE MODÈLE	EFFACER Y DÉCONNEXION			
IMPORTÉ	Affiliation(s)	ARK	Auteur(s)	Auteur(s) monographie	Catégories INIST	Catégories Science-Metrix	Catégories Scopus	Catégories WoS	Chapitre
	[["Department of Natural S	ark:/67375/WNG-CHKK1Z0N-H	["Jesper Hansen", "Nils-Mar		[{"Nom":"3 - sciences biol	[{"Nom":"3 - ecology","Cla	[{"Nom":"3 - Ecology","Cla	[{"Nom":"2 - geography, ph	
	[["Institute of Marine Res	ark:/67375/VQC-74LKF0CS-R	["Elena Eriksen","Bjarte B	۵	۵	П	0	۵	
	[["Institut für Polarökolo	ark:/67375/1BB-5CTZBF3P-Z	["C. F. von Dorrien"]		[{"Nom":"3 - sciences biol	. [{"Nom":"3 - marine biolog	[{"Nom":"3 - General Agric	[{"Nom":"2 - ecology","Cla	
	[["VNIRO, 17 V. Krasnosels	ark:/67375/WNG-L6WPQ8B3-3	["V. M. Borisov", "A. A. El	۵	[{"Nom":"3 - sciences biol		[{"Nom":"3 - Aquatic Scien		
	0	ark:/67375/HXZ-1MPJ41L6-4	["T.G. Sazykina"]		[{"Nom":"3 - sciences biol	. [{"Nom":"3 - nuclear medic	[{"Nom":"3 - Public Health	[{"Nom":"2 - public, envir	
	<		1197 ressou	56 colonnes	s chargées	nnes ajoutées			-+
									113



#### 1.1. Création du modèle

- Ajouter une colonne depuis le jeu de données original
- Ajouter une nouvelle colonne 🔍

			1197 ressources chargées 56 colonnes chargées 62 colonnes ajoutées						
uri	Nom du corpus	Nombre de documents	Titre de l'article	Lien vers le PDF	Type de corpus	Domaine(s)	Auteur(s)	Affiliation(s)	Langu
uid:/F8Z20TW8	Les poissons et les mollus	/api/run/count-all	Late Pleistocene and Holoc	https://api.istex.fr/ark:/	Corpus thématique	["Zoologie", "Géographie"]	["Jesper Hansen","Nils-Mar	[["Department of Natural S	anglais
uid:/TZDP19FM	Les poissons et les mollus	/api/run/count-all	Ecological significance of	https://api.istex.fr/ark:/	Corpus thématique	["Zoologie", "Géographie"]	["Elena Eriksen", "Bjarte B	[["Institute of Marine Res	anglais
uid:/CQSKD1FG	Les poissons et les mollus	/api/run/count-all	Reproduction and larval ec	https://api.istex.fr/ark:/	Corpus thématique	["Zoologie", "Géographie"]	["C. F. von Dorrien"]	[["Institut für Polarökolo	anglais
uid:/D297MCJP	Les poissons et les mollus	/api/run/count-all	Long-term variations of bi	https://api.istex.fr/ark:/	Corpus thématique	["Zoologie", "Géographie"]	["V. M. Borisov", "A. A. El	[["VNIRO, 17 V. Krasnosels	anglais
uid:/MD69C5G4	Les poissons et les mollus	/api/run/count-all	Long-Distance Radionuclide	https://api.istex.fr/ark:/	Corpus thématique	["Zoologie", "Géographie"]	["T.G. Sazykina"]	Ω	anglais
uid:/T0T5PT3G	Les poissons et les mollus	/api/run/count-all	Heavy metals in fish from	https://api.istex.fr/ark:/	Corpus thématique	["Zoologie", "Géographie"]	["GP Zauke", "V.M Savinov	[["Carl von Ossietzky Univ	anglais



# LODEX > Le modèle

Permet de configurer toutes les informations à afficher sur un corpus, à partir :

- Du jeu de données
- De nouvelles colonnes





Recherche



#### 1.2. Import d'un modèle existant





#### Le modèle importé





Les données publiées

Lodex - Modélisez et publiez vos données vers le web sémantique ! AJOUTER DES DONNÉES MODÉRATION RESSOURCES RETIRÉES VOIR LE MODÈLE EFFACER v DÉCONNEXION
Vos données ont été publiées en ligne. Vous pouvez configurer davantage l'affichage de ces dernières ou les mettre à jour en utilisant l'interface publique





# **LODEX > comment naviguer dans le corpus**

Accès à la page d'accueil







#### Page d'accueil

- Présentation du corpus
  - Identification du corpus
  - Rappel du cas d'usage
  - Possibilité de rejouer la requête
  - Droits d'utilisation des documents
  - Citation du corpus

Exposition du corpus



121

# **LODEX > modification des données publiées**

- 2 icônes disponibles pour chaque champ
  - Pour modifier la valeur du champ
  - Pour modifier le paramétrage du champ (format, facette, etc.)

Permet d'adapter le modèle d'un corpus à l'autre

#### • Exercices

10

- 1. Modifier le titre du corpus
- 2. Modifier le format du titre du corpus





#### Les graphiques

- Visualisation de la répartition des documents selon différents angles de vue
  - Dates de publication
  - Revues
  - Type de publication, type de document
  - Éditeurs
  - Mots-clés d'auteur
  - Catégories scientifiques (Science-Metrix, Scopus, WoS, Inist)
  - Indicateurs (version de PDF, score de qualité)

Analyse du corpus



Les graphiques

- Exercices
  - 1. Explorer le graphique "Revues"

□ utiliser la facette "Titre de revue" pour accéder aux documents d'une revue

□ utiliser la facette "Date de publication" pour cibler les revues récentes

□ changer le nombre de revues affichées



Les graphiques

- Exercices
  - 2. Explorer le graphique "Type de documents"

□ utiliser la facette "Type de documents" pour identifier les documents qui ne sont **pas** de type : article, article de recherche, article de synthèse

3. Explorer le graphique "Catégories Science-Metrix"

naviguer dans les différents niveaux de hiérarchie pour valider la qualité du corpus



# LODEX > de l'exploration à la ré-interrogation d'ISTEX

#### Équation affinée

(title:((/fish(es)?/ OR mollus\*) AND (arctic OR "sub-arctic" OR subarctic~1 OR svalbard OR
/spit[sz]berg(en)?/ OR "barents sea" OR "barents seas"~2 ))

OR abstract:((/fish(es)?/ OR mollus\*) AND (arctic OR "sub-arctic" OR subarctic~1 OR svalbard OR
/spit[sz]berg(en)?/ OR "barents sea" OR "barents seas"~2 ))

OR subject.value:((/fish(es)?/ OR mollus\*) AND (arctic OR "sub-arctic" OR subarctic~1 OR svalbard OR
/spit[sz]berg(en)?/ OR "barents sea" OR "barents seas"~2 )))

AND language:eng

AND qualityIndicators.pdfWordCount:[500 TO 10000]

NOT genre:abstract NOT categories.scienceMetrix:pediatrics

*Résultats (03 mars 2020) : 1 197 > 1 194 docs* 



## LODEX > exemples de corpus Inist

- Corpus avec enrichissement classification issue de Topic Modeling
  - Vieillissement
- Corpus avec annotation manuelle disponible
   <u>Unitex</u>
- Corpus avec enrichissement systématique animale
  - Échantillon Systématique Animale (version en développement)
- Corpus avec évolution diachronique des thématiques
  - Géosciences (version en développement)



# Liens utiles

Adresses & Co Documentation & Informations Contact & Discussion



#### Adresses & Co



#### Se connecter :

- ISTEX : <u>http://www.istex.fr</u>
- Démonstrateur Istex : <u>http://demo.istex.fr/</u>
- Interface Istex-DL : <u>https://dl.istex.fr/</u>
- Infos Lodex : <u>https://lodex.inist.fr/</u>
- Données Istex : <u>https://data.istex.fr/</u>

#### S'authentifier :

- Vérifier ses droits d'accès : <u>https://api.istex.fr/auth</u>
- Vérifier son accès par fédération d'identité : <u>https://api.istex.fr/auth?auth=fede</u>



#### **Documentation & Informations**



**B** 

#### Se documenter :

- Documentation "Usage TDM d'Istex" : <u>https://doc.istex.fr/tdm/</u>
- Documentation "API ISTEX" : <u>https://doc.istex.fr/api/</u>
- Documentation "Lodex" : <u>https://user-doc.lodex.inist.fr/</u>

#### Se tenir informé :

- Blog ISTEX : <u>https://blog.istex.fr/</u>
- Plateforme Twitter : <u>@ISTEX\_Platform</u>



#### **Contact & Discussion**



#### Chercher de l'aide / Contribuer à l'amélioration :

- Contact :
  - Via le formulaire : <u>https://www.istex.fr/contact/</u>
  - Via la liste : <u>contact@listes.istex.fr</u>

Liste de discussion : <u>users@listes.istex.fr</u> (liste publique)





# Annexes

S'aider...

# Interrogation > quelques champs utiles

Penser à utiliser les champs **"host.title" et "host.issn"** pour rechercher une revue (changement d'ISSN ou de titre au cours du temps, erreurs orthographiques, etc.) Titre d'article Auteur Résumé Mots-clés d'auteur Date de publication Langue Titre de revue ISSN revue

title author.name abstract subject.value publicationDate language host.title host.issn Les champs **.raw** permettent d'interroger sur la valeur exacte du champ.

Exemple :

q=host.title.raw:"Science" ne ramènera que les documents de la revue Science et non ceux de "The Science of the Total Environment", etc.



S'aider...

# Interrogation > quelques champs utiles

Éditeur Type de publication Type de document corpusName host.genre genre

Enrichissements

**Disciplines scientifiques** 

Entités nommées

categories.scopus, categories.inist, categories.sciencesMetrix, categories.wos

namedEntities.unitex, etc.



S'aider...

#### Interrogation > quelques champs utiles

Indicateurs (TDM) Score de qualité Version de PDF Nombre de mots par PDF

qualityIndicators.score qualityIndicators.pdfVersion qualityIndicators.pdfWordCount

À venir

Pdf image Nombre de mots par page qualityIndicators.pdfText qualityIndicators.pdfWordsPerPage

135

#### Sauter le pas...

# Interrogation > intérêts de l'API ISTEX

- Exploration du contenu du fonds ISTEX et répartition des documents en fonction d'une facette Exemple : connaître toutes les langues des documents présents dans ISTEX
- Interrogation selon une ou plusieurs facettes imbriquées

Exemple : Pour chaque langue présente dans ISTEX, connaître le nombre de documents par siècle

 Interrogation de toutes les facettes disponibles dans l'API (nombre réduit sur le démonstrateur)



#### **Tutos**

Pour comprendre ce qu'est l'API ISTEX et apprendre à l'utiliser, reportez-vous aux tutoriels :

> https://istex-tutorial.data.is tex.fr/

# Interrogation > Exemples de requêtes sur l'API

Éditeurs

https://api.istex.fr/document/?q=\*&facet=corpusName[\*]&size=0

Types de publication

https://api.istex.fr/document/?q=\*&facet=host.genre[\*]&size=0

Types de document

https://api.istex.fr/document/?q=\*&facet=genre[\*]&size=0





Exemple 1 : lister tous les

connaître le nombre de

(emploi de **size=0**)

éditeurs présents dans ISTEX et

documents pour chacun d'eux sans afficher les documents

# Interrogation > Exemples de requêtes sur l'API



Exemple 1 : explorer la répartition du nombre de documents Istex par période de 100 ans entre les 15e et 21e siècles

Dates de publication

par siècle (entre 1401 et 2100) https://api.istex.fr/document/?q=\*&size=0&facet=publicationD ate[1401-2100:100]

par période de 10 ans (20& et début 21 siècles) https://api.istex.fr/document/?q=\*&size=0&facet=publicationD ate[1901-2020:10]

Langues & dates les dates pour chaque langue **(facettes imbriquées)** 

https://api.istex.fr/document/?q=\*&facet=language[\*]>publica tionDate[1401-2100:100]&size=0

# Interrogation > Exemples de requêtes sur l'API



Exemples **1 & 2** : savoir si les documents présents dans ISTEX **comportent ou non un résumé et afficher** ces documents

Présence de résumé avec résumé sans résumé https://api.istex.fr/document/?q=(abstract:\*)
https://api.istex.fr/document/?q=(NOT (abstract:\*))

Présence deavec mots-clésmots-clés d'auteursans mots-clés

https://api.istex.fr/document/?q=(**subject.value.raw:\***)&size=0 https://api.istex.fr/document/?q=(**NOT** (subject.value.raw:\*))&size=0

# Interrogation > Exemples de requêtes sur l'API



Exemple 1 : quantifier les documents appartenant à chaque catégorie Scopus d'un corpus de documents sur l'Arctique

#### Domaines scientifiques

catégories Scopus d'un corpus particulier

https://api.istex.fr/document/?q=**arctic**&facet=**categories.scopus[\*]**& <u>size=0</u>

catégories Science-Metrix de l'archive ISTEX

https://api.istex.fr/document/?q=\*&facet=**categories.scienceMetrix[**\* ]&size=0

# Interrogation > Exemples de requêtes sur l'API



41

Exemple 1 : obtenir la distribution du nombre de documents ISTEX en fonction de leur **score de qualité**, et ce par **pas de 1** 

#### Indicateurs

score de qualité

https://api.istex.fr/document/?q=\*&size=0&facet=qualityIndica tors.**score[1-10:1**]

version de PDF

https://api.istex.fr/document/?q=\*&size=0&facet=qualityIndica
tors.pdfVersion[\*]

nombre de mots par PDF

https://api.istex.fr/document/?q=\*&size=0&facet=qualityIndica tors.pdfWordCount[1-10000:1000]



# **Merci** Avez-vous des questions ?