



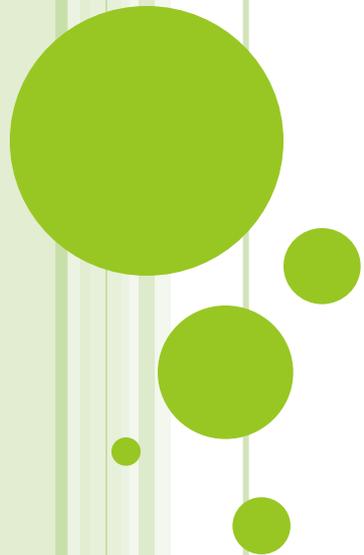
L'INTÉGRATION DES ARK DANS ISTEX EN QUELQUES QUESTIONS

Nicolas Thouvenin

 [thouveni](#)

CNRS/INIST

QU'EST-CE QUE ISTEEX ?



ISTEX

L'excellence documentaire pour tous

“ Construire le socle
de la bibliothèque scientifique
numérique nationale. ”



ANR-10-IDEX-0004-02

WWW.ISTEX.FR

www.licencesnationales.fr

“ ISTEEX plus de **18,5 millions de documents** provenant de **19 éditeurs**
plus de **7 500 titres** et **13 000 ebooks** entre 1406 et 2015. ”



abes
agence bibliographique de l'enseignement supérieur

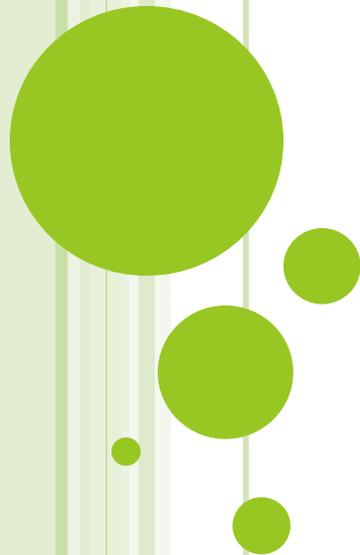
CU CONFÉRENCE
DES PRÉSIDENTS
D'UNIVERSITÉ

UNIVERSITÉ
DE LORRAINE

couperin.org
Consortium universitaire
de publications numériques



COMMENT IDENTIFIER UN “OBJET DOCUMENTAIRE” ISTEEX ?



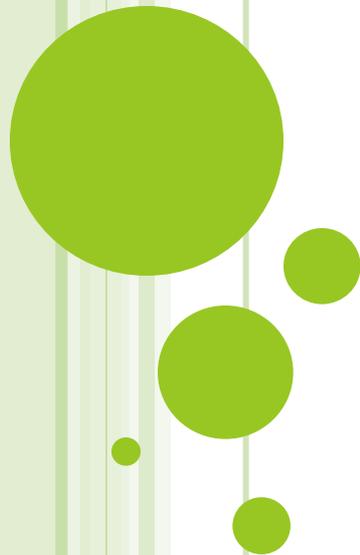
AVEC UN SYSTÈME AUTOMATIQUE !

- 7AF88EC8E547846C736021BF26ED71EFABB36991
- B14CC18B27A2A5A13D87E5DE4052EE9E340C02CE
- 4471A7AC4A966A679545A477299512E5C0C38230
- 7E195F8F9F9A75BBB8AFDDC56A5A6B7F30E6F01B
- A0DEFD7E09821E0A5F50D31FA70B86B7F960B493
- DF0E8799A961756E13EF36BDB170C2A30E2C3D44
- DB7FB95E6ABD3065F754578ACF4BE00257286346
- C0CE73FE3B5B1D1EBDC9A5B3CFC079391F36CCF5
- 5E780B4F23B9A4E198EC8742064ED77015B5D872
- 9E66D0A322F35469F4BC12AC03C9961FFF2800DE

OU BIEN AVEC UN SYSTÈME NORMALISÉ

- DOI,
- HANDLE
- URI
- ARK
- ...

POURQUOI CHOISIR LE SYSTÈME ARK?



ANALYSE DE NOS CONTRAINTES

- Présence d'un DOI sur 96% des documents ISTEEX
 - DOI
- Une architecture informatique existante (non java)
 - HDL
- Un système d'identifiant ad-hoc existant faiblement normalisé
 - URI

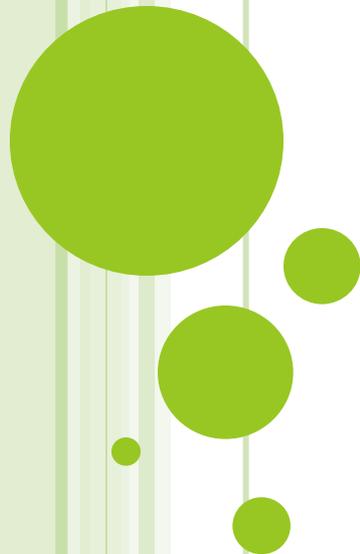
ANALYSE DE NOS BESOINS

- Pouvoir s'intégrer facilement dans l'existant
 - HDL
- Éviter une dépendance avec un système externe
 - HDL, DOI
- Identifier avec le même système des données différentes
 - ADHOC

ANALYSE DES AVANTAGES DU SYSTÈME ARK

- système décentralisé
- sans contrainte technique
- dissocier d'un type de données

QUELLE AUTORITÉ NOMMANTE CHOISIR ?



CHOIX DU NAAN

- ISTEEX est un **projet** et non une **institution**
- L'Inist-CNRS une **unité** du CNRS (ups076)
- Le CNRS est l'un des **acteurs** du projet avec l'ABES, l'université, Couperin, le ministère.

Qui a autorité sur les identifiants ISTEEX ?

Pour ISTEK

l'Inist-CNRS est chargé :

- de mise en oeuvre de la plateforme
- de réceptionner les documents ISTEK
- ...

⇒ **il est l'autorité nommante**

```
naa:  
who: Institut de l'information scientifique et technique (=)  
Institute for scientific and technical information (=) INIST  
what: 67375  
when: 2016.03.30  
where: http://www.inist.fr  
how: NP | (:unkn) unknown | 2016 |
```

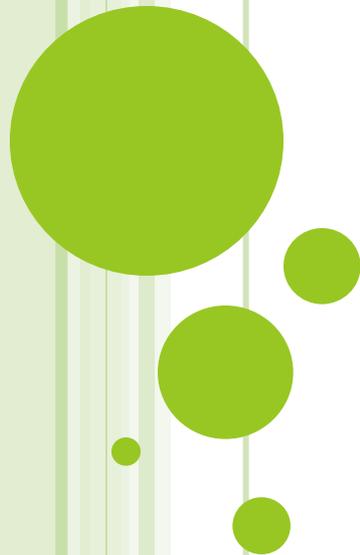
Pour le CNRS

l'Inist-CNRS est chargé :

- de mise en oeuvre d'autres plateformes
- de traiter des documents non ISTEEX
- de traiter des données différentes (terminologie, etc.)
- ...

Le NAAN Inist ne peut pas être associé uniquement à ISTEEX

COMMENT GÉRER UNE SOUS AUTORITÉ NOMMANTE ?



REGISTRE

Un **registre unique** à l'Inist
partagé par tous les projets et services
pour **enregistrer & conserver** les préfixes.

> un fichier Excel partagé !

ou via une application similaire...

ARK Registry

Registre officiel des préfixes ARK pour l'Inist-CNRS

Cette interface centralise les subpublishers ARK (prefix) du NAAN de l'Inist-CNRS (67375). Ce registre permet de générer, recenser et garantir l'unicité de ces subpublishers. Chaque équipe ou projet de l'INIST peut ainsi venir déclarer un subpublisher pour ensuite venir associer des ARK aux données gérées par une plateforme donnée (ex : plateforme ISTEEX ou plateforme Loterre).

NAAN

67 375

Nombre de subpublishers

74

Liste des subpublishers déclarés

LODEX
property-dataset

ISTEX
Corpus IOP

ISTEX
Corpus Springer

Loterre, Transfert de chaleur
jeu de données

Vocabulaire des sciences de
l'éducation
terminologie

Loterre
jeu de données, terminologie

ARK Registry

GRAPH LIST

Found 74 on 74

enter your search
term here

Resources

URI

PL

09J

LODE

0T8

ISTE

1BB

ISTE

1WB

Loter

216

Vocat

26L

Loter

27X

Loter

2BK

DPI

3JP

Loter

jeu de données, terminologie

3VV

Vocabulaire d'écotoxicologie

jeu de données

/ Equipe

du
er ARK

Add a new resource to the dataset

uri

Leave empty to autogenerate uri

Plateforme / Equipe

Description du subpublisher ARK

Commentaire

URL d'accès aux ressources de ce subpublisher

SAVE

CANCEL

UN PREFIX DE 3 caractères :

1BB

ezark -

<https://github.com/Inist-CNRS/ezark>

*Une configuration particulière de l'outil
Lodex*



COMMENT GÉNÉRER DES IDENTIFIANTS ?

UNE NOMENCLATURE

ark:/67375/**XXX**-**YYYYYYYYYYY**-**Z**

- **Un préfixe**
- **Un identifiant sur 10 caractères**
- **un caractère de contrôle - NCDA**
NOID CHECK DIGIT ALGORITHM



QUELS QUALIFICATIFS UTILISÉS ?

UN OBJET DOCUMENTAIRE

peut-être composé de :

- ❖ Métadonnées
 - XML, MODS
- ❖ Texte intégral
 - PDF, TEI, TXT, OCR, ZIP, TIFF
- ❖ Annexes
 - PDF, TXT, DOC, JPEG, QT, MPEG, MP4, PPT, XLS, XLSX, AVI, XML, RTF, GIF, WMV
- ❖ Couvertures
 - PDF, GIF, JPEG, TIFF, HTML
- ❖ Enrichissements
 - Entités nommées, Références bibliographiques, etc.

ÉTUDES DES URLs dynamiques pré-existantes

TYPLOGIE	FORMAT	dans l'url Istex	
record	MODS	metadata	mods
record	XML	metadata	xml
record	JSON		
annexes	PDF	annexes	pdf
annexes	JPEG	annexes	jpeg
annexes	TXT	annexes	txt
annexes	XLS	annexes	xls
annexes	XLSX	annexes	xlsx
annexes	GIF	annexes	gif
annexes	ZIP	annexes	zip
annexes	PPT	annexes	ppt
annexes	QT	annexes	qt
annexes	DOC	annexes	doc
annexes	MP4	annexes	mp4
annexes	MPEG	annexes	mpeg
annexes	AVI	annexes	avi
annexes	WMV	annexes	wmv
annexes	RTF	annexes	rtf
covers	TIFF	covers	tiff
covers	HTML	covers	html
covers	GIF	covers	gif
covers	PPT	covers	ppt
fulltext	PDF	fulltext	pdf
bundle	ZIP	fulltext	zip
fulltext	TEI	fulltext	tei
fulltext	TXT	fulltext	txt
fulltext	OCR	fulltext	ocr
entities	TEI	enrichments	unitex
refbibs	TEI	enrichments	refbibs
keywords-nb	TEI	enrichments	nb
keywords-multicat	TEI	enrichments	multicat
keywords-teeft	TEI	enrichments	teeft
keywords-abes	TEI	enrichments	abesSubjects
authors	TEI	enrichments	abesAuthors

TYPLOGIE	FORMAT	dans l'url Istex	
fulltext-original	PDF	fulltext	original
record-original	XML	metadata	original
covers-original	TIFF	cover	original
annexes-original	ZIP	annexes	original
annexes-original	PDF	annexes	original
annexes-original	JPEG	annexes	original
annexes-original	TXT	annexes	original
annexes-original	XLS	annexes	original
annexes-original	XLSX	annexes	original
annexes-original	GIF	annexes	original
annexes-original	ZIP	annexes	original
annexes-original	PPT	annexes	original
annexes-original	QT	annexes	original
annexes-original	DOC	annexes	original
annexes-original	MP4	annexes	original
annexes-original	MPEG	annexes	original
annexes-original	AVI	annexes	original
annexes-original	WMV	annexes	original
annexes-original	RTF	annexes	original
covers-original	TIFF	covers	original
covers-original	HTML	covers	original
covers-original	GIF	covers	original
covers-original	PPT	covers	original

*Beaucoup de
combinaisons possibles*

UNE NOMENCLATURE SIMPLIFIÉE

ark:/67375/XXX-YYYYYYYYYYY-Z/**TYPOLOGIE . FORMAT**

- **3 Typologies :**
fulltext, metadata, bundle
- **7 formats :**
zip, pdf, xml, txt, json, tei, mods



QUELS FORMATS UTILISER PAR DÉFAUT

PROBLÈMES DES QUALIFICATIFS POUR CHAQUE DOCUMENT

- non obligatoire
- pas toujours disponible

⇒ Ils varient en fonction

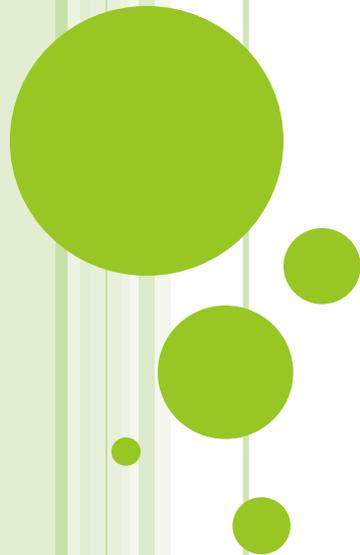
- des données,
- des outils,
- de la plateforme

L'API ISTEEX RENVOIE POUR

un ARK :

- ❖ soit le document demandé dans le format sélectionné
 - <ark:/67375/VQC-ZJHCRXRV-B/fulltext.pdf>
- ❖ soit la liste des formats disponibles pour la typologie sélectionnée (en JSON)
 - <ark:/67375/VQC-ZJHCRXRV-B/fulltext>
- ❖ soit la liste des typologies et formats disponibles pour le document sélectionné (en JSON)
 - <ark:/67375/VQC-ZJHCRXRV-B>

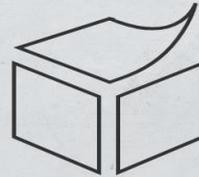
QUELLES “RESSOURCES” IDENTIFIER AVEC DES ARK ?



ALLER AU DELÀ DU PDF

Les objets documentaires ISTEEX, mais aussi

- les référentiels documentaires
- les données (concept) générées automatiquement
- **presque** tout ce que l'on diffuse au travers du portal <https://data.istex.fr>

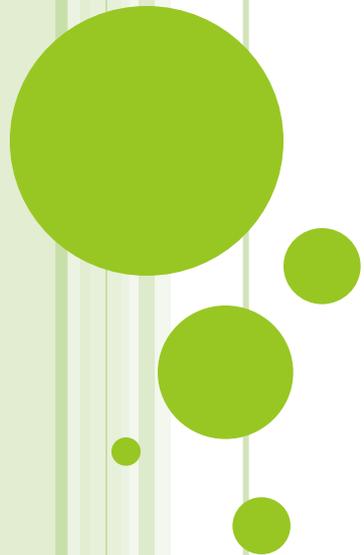


data.istex.fr

LISTE DES JEUX DE DONNÉES

Catégories Science Metrix + EN SAVOIR PLUS	Entité PlaceName + EN SAVOIR PLUS	Editeurs Scientifiques + EN SAVOIR PLUS
Catégories Web Of Science + EN SAVOIR PLUS	Catégories Scopus + EN SAVOIR PLUS	Catégories Inist + EN SAVOIR PLUS
Référentiel Des Corpus Chargés Dans ISTEX + EN SAVOIR PLUS	Tutoriels ISTEX + EN SAVOIR PLUS	Entités Nommées + EN SAVOIR PLUS
Publications Remarquables Dans ISTEX + EN SAVOIR PLUS	Types De Contenu + EN SAVOIR PLUS	Domaines Scientifiques + EN SAVOIR PLUS
Types De Publication + EN SAVOIR PLUS	Enrichissements ISTEX + EN SAVOIR PLUS	Langues De Publication + EN SAVOIR PLUS
Ayants Droit À L'usage D'ISTEX + EN SAVOIR PLUS		Entité GeogName + EN SAVOIR PLUS

QUE RESTE-T-IL À FAIRE ?



BEAUCOUP !

- résolution ARK
- diffusion /utilisation
- formation
- “respect” des identifiants
- politique de conservation
- (...)



MERCI

