

CILLEX

Exploiter les réseaux documentaires
pour mieux comprendre les résultats d'une recherche

Emmanuel Navarro, Bruno Gaume

CLLE, CNRS & Université de Toulouse

Séminaires ISTEK

25 avril 2016

1 Objectifs

2 Etat Courant

3 Résultats exploitables pour les « Chantiers d'usage »

- API Cillex
- Annotations et système d'annotation
- Graphes globaux
- Bibliothèque Python, client API ISTEEX (et plus)

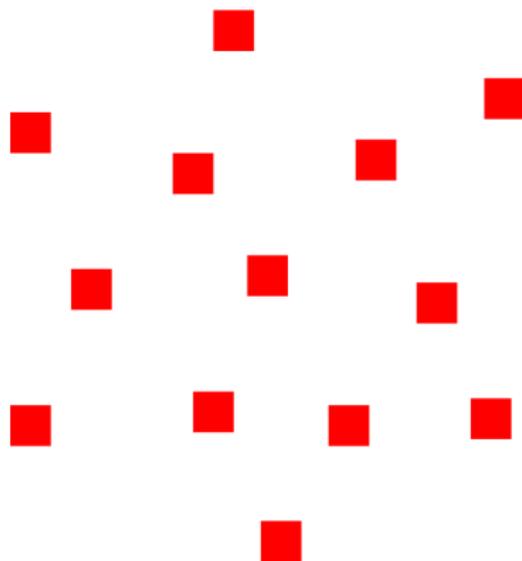
1 Objectifs

2 Etat Courant

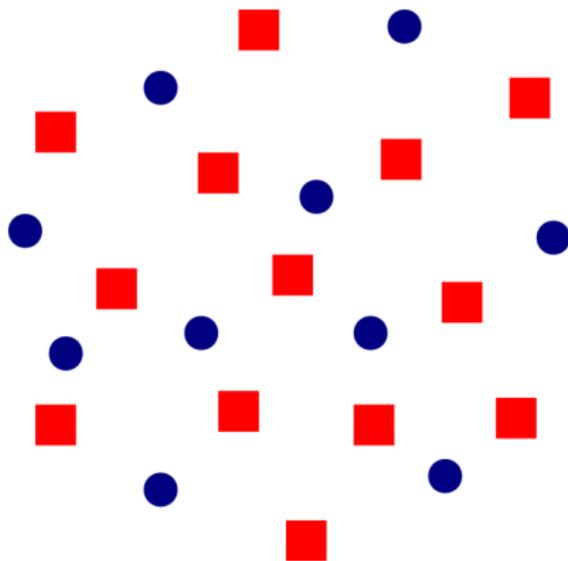
3 Résultats exploitables pour les « Chantiers d'usage »

- API Cillex
- Annotations et système d'annotation
- Graphes globaux
- Bibliothèque Python, client API ISTEEX (et plus)

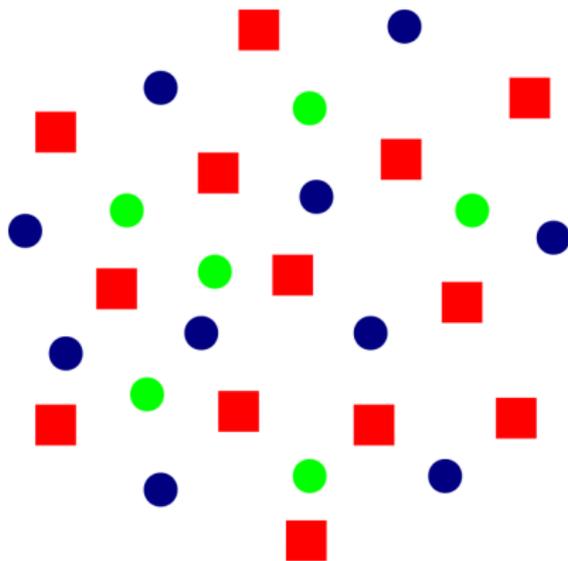
Des Documents



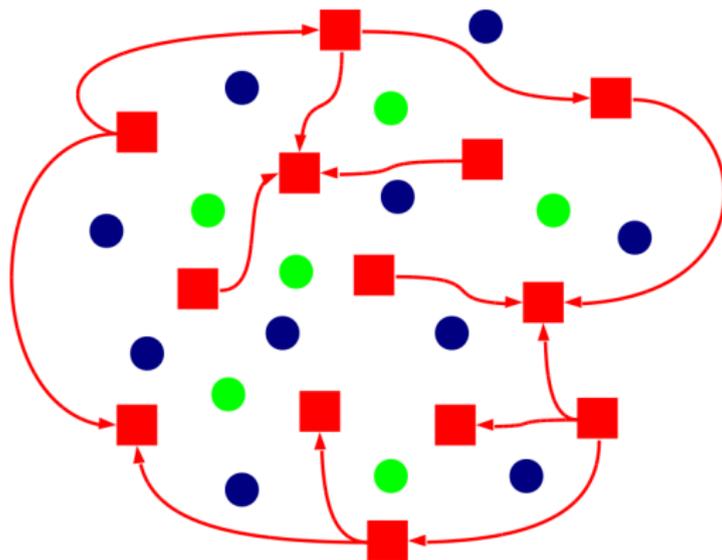
Des Documents, Des Mots



Des Documents, Des Mots, Des Auteurs

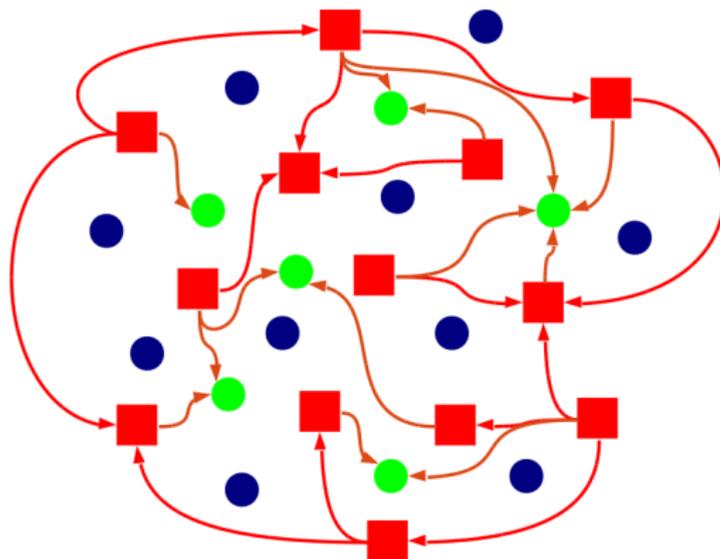


Des Documents, Des Mots, Des Auteurs, ... Des Liens

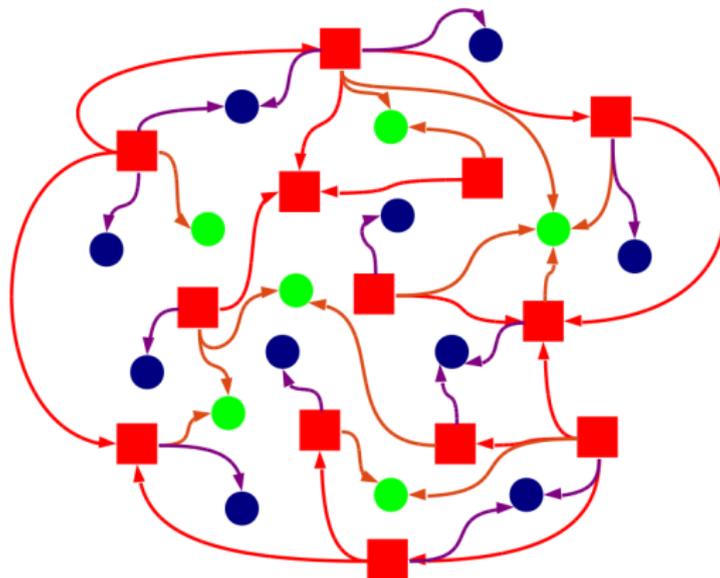


Objectifs

Des Documents, Des Mots, Des Auteurs, ... Des Liens, Des Liens

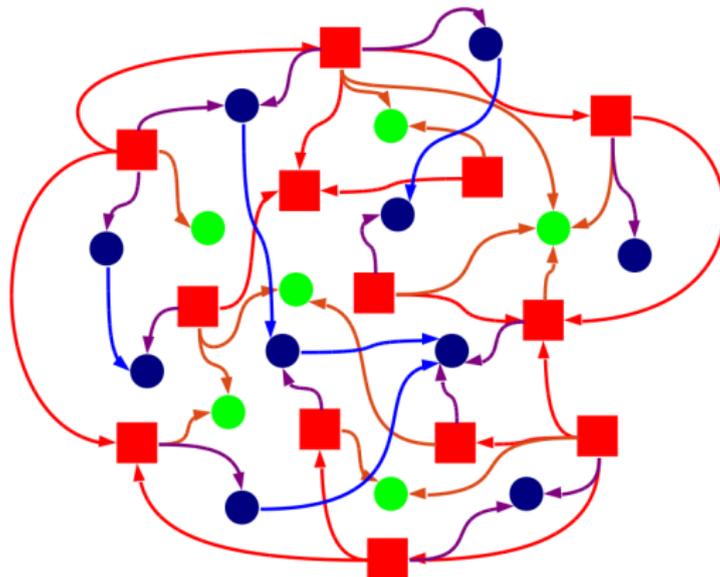


Des Documents, Des Mots, Des Auteurs, ... Des Liens, Des Liens

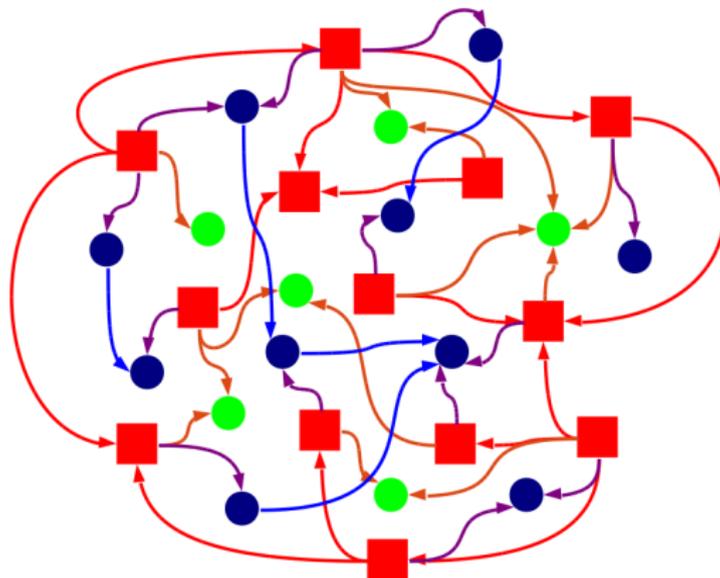


Objectifs

Des Documents, Des Mots, Des Auteurs, ... Des Liens, Des Liens



Graphe de terrain



Graphe de terrain

- ▶ Les graphes d'acointance d'un groupe d'humains
- ▶ Le graphe du World Wide Web
- ▶ Le graphe de *Caenorhabditis elegans*
- ▶ Les graphes Lexicaux
- ▶ Les graphes extraits des bases documentaires

Graphe de terrain

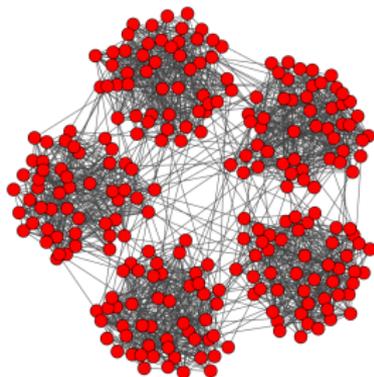
Quatre propriétés fondamentales

- ▶ Faible densité
- ▶ Chemins courts
- ▶ Distribution des degrés à queue lourde (loi de puissance)
- ▶ Fort coefficient de clustering : zones denses en arêtes

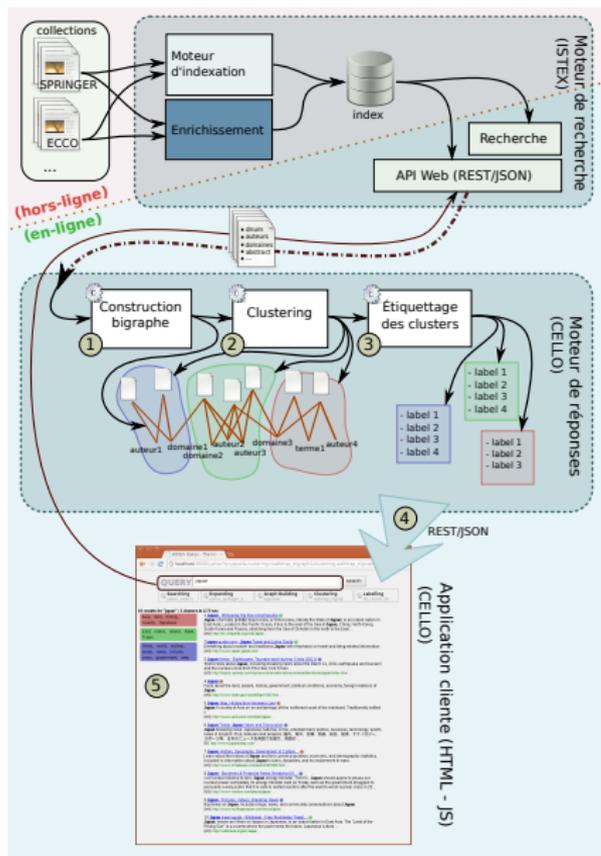
Graphe de terrain

Quatre propriétés fondamentales

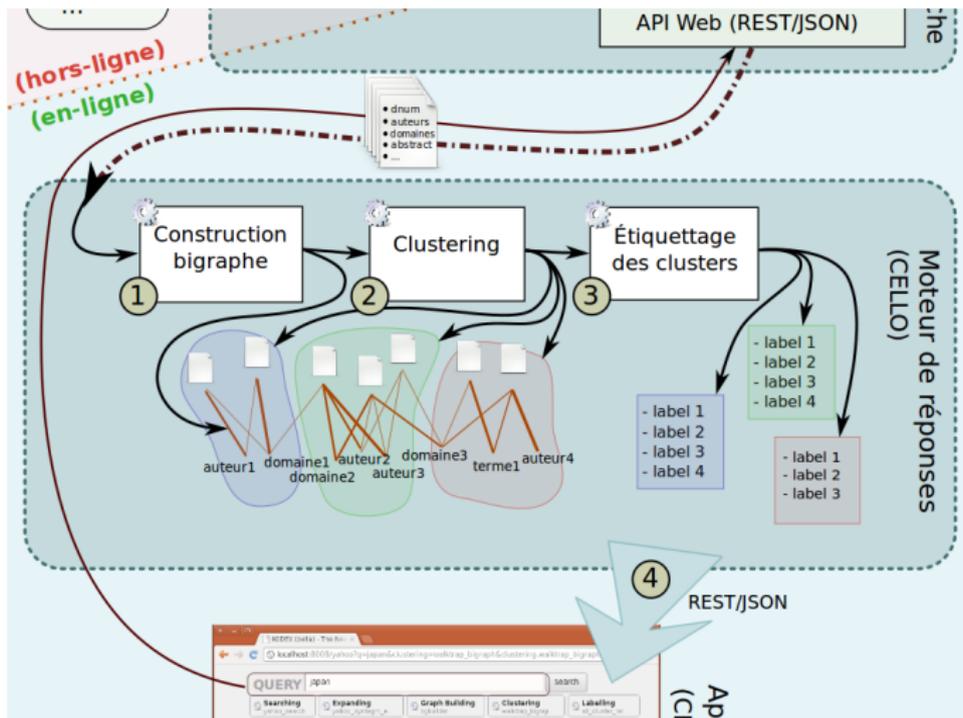
- ▶ Faible densité
- ▶ Chemins courts
- ▶ Distribution des degrés à queue lourde (loi de puissance)
- ▶ **Fort coefficient de clustering : zones denses en arêtes**



Schema général



Schema général



1 Objectifs

2 Etat Courant

3 Résultats exploitables pour les « Chantiers d'usage »

- API Cillex
- Annotations et système d'annotation
- Graphes globaux
- Bibliothèque Python, client API ISTEEX (et plus)

<http://cillex.kodexlab.com/>

Un point sur les données disponibles/exploitable

- ▶ **Cat. WOS**
 - Résultats intéressants,
 - Couverture partielle,
 - PB : liées aux « host » non aux articles.
- ▶ **Mots des titres et abstracts**
 - Couverture de 100% (ou presque...),
 - Mais, sélection/extraction nécessaire... (travail en cours)
- ▶ **« Subjects »**
 - Couverture partielle,
 - PB : peu de consistance, données éparses,
 - Utilisation du graphe global (travail en cours).
- ▶ **Ref. bib.**
 - Couverture partielle,
 - PB de match/normalisation,
 - On a uniquement les citations sortantes...
- ▶ **Auteurs, affiliations,**
 - PB : données éparses, normalisation.
- ▶ **« Host », revue, série,**
- ▶ **Enrichissement par les autres projets, à venir...**

Un point sur les données disponibles/exploitable

- ▶ **Cat. WOS**
 - Résultats intéressants,
 - Couverture partielle,
 - PB : liées aux « host » non aux articles.
- ▶ **Mots des titres et abstracts**
 - Couverture de 100% (ou presque...),
 - Mais, sélection/extraction nécessaire... (travail en cours)
- ▶ **« Subjects »**
 - Couverture partielle,
 - PB : peu de consistance, données éparses,
 - Utilisation du graphe global (travail en cours).
- ▶ **Ref. bib.**
 - Couverture partielle,
 - PB de match/normalisation,
 - On a uniquement les citations sortantes...
- ▶ **Auteurs, affiliations,**
 - PB : données éparses, normalisation.
- ▶ **« Host », revue, série,**
- ▶ **Enrichissement par les autres projets, à venir...**

Un point sur les données disponibles/exploitable

- ▶ Cat. WOS
 - Résultats intéressants,
 - Couverture partielle,
 - PB : liées aux « host » non aux articles.
- ▶ **Mots des titres et abstracts**
 - Couverture de 100% (ou presque...),
 - Mais, sélection/extraction nécessaire... (travail en cours)
- ▶ « Subjects »
 - Couverture partielle,
 - PB : peu de consistance, données éparses,
 - Utilisation du graphe global (travail en cours).
- ▶ Ref. bib.
 - Couverture partielle,
 - PB de match/normalisation,
 - On a uniquement les citations sortantes...
- ▶ Auteurs, affiliations,
 - PB : données éparses, normalisation.
- ▶ « Host », revue, série,
- ▶ Enrichissement par les autres projets, à venir...

Un point sur les données disponibles/exploitable

- ▶ Cat. WOS
 - Résultats intéressants,
 - Couverture partielle,
 - PB : liées aux « host » non aux articles.
- ▶ Mots des titres et abstracts
 - Couverture de 100% (ou presque...),
 - Mais, sélection/extraction nécessaire... (travail en cours)
- ▶ « **Subjects** »
 - Couverture partielle,
 - PB : peu de consistance, données éparses,
 - Utilisation du graphe global (travail en cours).
- ▶ Ref. bib.
 - Couverture partielle,
 - PB de match/normalisation,
 - On a uniquement les citations sortantes...
- ▶ Auteurs, affiliations,
 - PB : données éparses, normalisation.
- ▶ « Host », revue, série,
- ▶ Enrichissement par les autres projets, à venir...

Un point sur les données disponibles/exploitable

- ▶ Cat. WOS
 - Résultats intéressants,
 - Couverture partielle,
 - PB : liées aux « host » non aux articles.
- ▶ Mots des titres et abstracts
 - Couverture de 100% (ou presque...),
 - Mais, sélection/extraction nécessaire... (travail en cours)
- ▶ « Subjects »
 - Couverture partielle,
 - PB : peu de consistance, données éparses,
 - Utilisation du graphe global (travail en cours).
- ▶ Ref. bib.
 - Couverture partielle,
 - PB de match/normalisation,
 - On a uniquement les citations sortantes...
- ▶ Auteurs, affiliations,
 - PB : données éparses, normalisation.
- ▶ « Host », revue, série,
- ▶ Enrichissement par les autres projets, à venir...

Un point sur les données disponibles/exploitable

- ▶ Cat. WOS
 - Résultats intéressants,
 - Couverture partielle,
 - PB : liées aux « host » non aux articles.
- ▶ Mots des titres et abstracts
 - Couverture de 100% (ou presque...),
 - Mais, sélection/extraction nécessaire... (travail en cours)
- ▶ « Subjects »
 - Couverture partielle,
 - PB : peu de consistance, données éparses,
 - Utilisation du graphe global (travail en cours).
- ▶ Ref. bib.
 - Couverture partielle,
 - PB de match/normalisation,
 - On a uniquement les citations sortantes...
- ▶ **Auteurs, affiliations,**
 - PB : données éparses, normalisation.
- ▶ « Host », revue, série,
- ▶ Enrichissement par les autres projets, à venir...

Un point sur les données disponibles/exploitable

- ▶ Cat. WOS
 - Résultats intéressants,
 - Couverture partielle,
 - PB : liées aux « host » non aux articles.
- ▶ Mots des titres et abstracts
 - Couverture de 100% (ou presque...),
 - Mais, sélection/extraction nécessaire... (travail en cours)
- ▶ « Subjects »
 - Couverture partielle,
 - PB : peu de consistance, données éparses,
 - Utilisation du graphe global (travail en cours).
- ▶ Ref. bib.
 - Couverture partielle,
 - PB de match/normalisation,
 - On a uniquement les citations sortantes...
- ▶ Auteurs, affiliations,
 - PB : données éparses, normalisation.
- ▶ « Host », revue, série,
- ▶ **Enrichissement par les autres projets, à venir...**

1 Objectifs

2 Etat Courant

3 Résultats exploitables pour les « Chantiers d'usage »

- API Cillex
- Annotations et système d'annotation
- Graphes globaux
- Bibliothèque Python, client API ISTEEX (et plus)

- ▶ API HTTP/JSON, avec « en plus » de l'API ISTEEX :
 - Clustering
 - Labelling
 - Graphe biparti (absent par défaut)

```
1 $ curl -GET http://cillex.kodexlab.com/api/search/brain
2 {
3   "meta": {...},
4   "results": {
5     "clusters": {
6       "clusters": [...],
7       "labels": [...],
8     },
9     "docs": [...],
10    "meta": {...},
11    "query": "br ain"
12  }
13 }
```

Listing 1 – Exemple simple avec requete GET

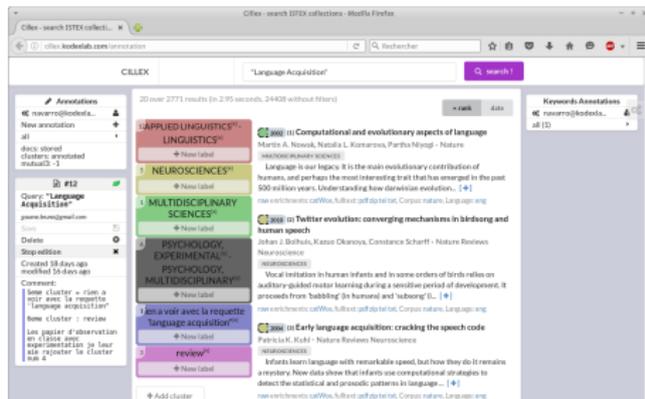
Déjà en place, mais :

- ▶ Peu documentée,
- ▶ Pas encore de version stabilisée.

Tout ça se corrige si il y en a la besoin...

Annotations et système d'annotation

- ▶ Approche hybride entre non-regression & campagne d'éval
- ▶ Sys. en trois parties
 - Serveur d'annotation (API JSON)
 - Client annotation des clusters
 - Client annotation des titres & abstracts (mots-clés)



Potentiellement réutilisable :

- ▶ Réutilisation du système (serveur) pour construire d'autres annotations,
- ▶ Jeu d'annotation de clusters (~ 15 requêtes pour le moment)
- ▶ Jeu d'annotation des titres & abstracts (à venir...)

- ▶ Possible avec l'API ISTEEX de construire des graphes globaux :
 - Subject-Subject,
 - Auteurs-revues,
 - Cat. WOS-subjects,
 - Cat. WOS-revues,
 - *Ref. Bib. (pb de match)*,
 - ...
- (lien si présent ensemble sur -au moins- un document)
- ▶ En fct des données demandées, de quelques heures à quelques jours.
 - ▶ Pour nous : expérimentation (en cours) utilisation du graphe global de *Subjects* pour répondre au problème de données clairsemées.

Envisageable de diffuser :

- ▶ Des graphes (format ? licence ?),
- ▶ Les scripts de constructions.

Reliex : client Python pour l'API ISTEEX (et plus)

- ▶ Basé sur *Reliure*¹, framework minimal pour construire des chaînes de traitement de données,
- ▶ Jeu de composants pour :
 - Construire des requêtes (filtres),
 - Intérogger l'API,
 - Post-Traitements minimaux.

Fonctionnalités principales :

- ▶ Système modulaire, enchaînement de composants (opérateur « | »)
- ▶ Déclaration/découverte d'options :
 - construction auto. d'API HTTP/JSON,
 - construction auto. de parser de param. en ligne de cmd (basé sur *argparse*)
- ▶ Système de cache HTTP local très simple avec *requests-cache*.

1. <http://reliure.readthedocs.org/>

```
1 from reliex.search import ISTEKScanQuery
2 from reliex.query import LanguageFilter, HasSubject
3
4 # Composant d'enrichissement de la requete
5 query_builder = LanguageFilter() | HasSubject(default=True)
6 # Composant complet de recherche
7 search = query_builder | ISTEKScanQuery(outputs="subject")
8
9 if __name__ == '__main__':
10     from argparse import ArgumentParser
11     from reliure.utils.cli import arguments_from_optionable,
12         get_config_for
13     parser = ArgumentParser()
14     arguments_from_optionable(parser, search, prefix="search-")
15     parser.add_argument('QUERY', type=str, help='The search query')
16
17     args = parser.parse_args()
18     search_config = get_config_for(args, search, prefix="search-")
19     query = args.QUERY
20
21     # Run it
22     for doc in search(query, **search_config):
23         pprint(doc)
```

Listing 2 – simple_example.py

```
$ python simple_example.py -h
2 usage: simple_example.py [-h]
                             [--search-language
                             {ALL,eng,fre,deu,lat,spa,ita,dut,rus,wel}]
4                             [--search-not-has_subject]
                             [--search-has_host_subject]
6                             [--search-operator {AND,OR}]
                             QUERY
8
positional arguments:
10  QUERY                    The search query

optional arguments:
12  -h, --help                show this help message and exit
14  --search-language {ALL,eng,fre,deu,lat,spa,ita,dut,rus,wel}
                             Documents language
16  --search-not-has_subject
                             'subject' needed
18  --search-has_host_subject
                             'host_subject' needed
20  --search-operator {AND,OR}
                             Default search operator
```

Listing 3 – Usage

```
$ python simple_example.py --search-language fre brain | head -20
2 {'id': 'DFBCC1D6D7A029F5B26CC3BEE331559893293220 ',
   'rank': 0,
4  'subject': [{'lang': ['fre'], 'value': 'Anura'},
               {'lang': ['fre'], 'value': 'ontogeny'},
6              {'lang': ['fre'], 'value': 'heterochrony'},
               {'lang': ['fre'], 'value': 'variation'},
8              {'lang': ['fre'], 'value': 'Ceratophryinae'}]}
{'id': 'CA69329F917E7FE9694CB50C047AA711ACAB3781 ',
10  'rank': 1,
   'subject': [{'lang': ['fre'], 'value': 'phylogeny'},
               {'lang': ['fre'], 'value': 'Leptodactylus fuscus
12                  group'},
               {'lang': ['fre'], 'value': 'osteology'},
14              {'lang': ['fre'], 'value': 'morphology'}]}
{'id': '4335E3679C9B3B080D158F3DC2EEA07656A32ED8 ',
16  'rank': 2,
   'subject': [{'lang': ['fre'], 'value': 'AUFSÄTZE'}]}
18 {'id': 'E60B5E6F0BBFDE96BB245A09D7C160FB6F0C6072 ',
   'rank': 3,
20  'subject': [{'lang': ['fre'], 'value': 'Escherichia coli'},
               {'lang': ['fre'], 'value': 'facteurs de virulence'},
```

Listing 4 – Exemple

Merci

Questions ?