

# Néoveille - documentation technique, structure et fonctionnement du site web

Emmanuel Cartier, emmanuel.cartier@lipn.univ-paris13.fr

14 octobre 2018

## Table des matières

<b>1</b>	<b>Présentation générale des programmes</b>	<b>2</b>
<b>2</b>	<b>Installation de l'interface web de Néoveille</b>	<b>4</b>
2.1	Prérequis . . . . .	4
2.2	Procédure d'installation . . . . .	4
<b>3</b>	<b>Installation des programmes Néoveille en backend</b>	<b>4</b>
3.1	Pré-requis . . . . .	4
3.2	Procédure d'installation . . . . .	5
3.3	Récupération dynamique des fichiers sur le web : fichier <b>corpus_fr.py</b> . . . . .	5
3.3.1	Paramètres à régler dans le fichier . . . . .	5
3.3.2	Installation du fichier corpus_fr.py en tâche récurrente . . . . .	5
3.3.3	Structure du fichier corpus_fr.py . . . . .	5
3.4	Détection des néologismes, stockage des néologismes candidats et mise à jour des articles dans Apache Solr : fichier <b>detect_neologisms_all.py</b> . . . . .	6
3.4.1	Paramètres à régler dans le fichier . . . . .	6
3.4.2	Installation du fichier detect_neologisms_all.py en tâche récurrente . . . . .	6
3.4.3	Structure du fichier detect_neologisms_all.py . . . . .	6
<b>4</b>	<b>Détails des fichiers : interface web</b>	<b>8</b>
4.1	Organisation des dossiers et fichiers principaux . . . . .	8
4.2	Organisation page d'accueil . . . . .	8
4.2.1	Feuilles de style . . . . .	8
4.2.2	Librairies javascript . . . . .	9
4.3	Bandeau supérieur (<nav class="navbar...">) . . . . .	9
4.4	Menu gauche (<div class="side-menu">) . . . . .	10
4.4.1	Action Javascript associée aux éléments de menu . . . . .	10
4.5	Page de contenu (<div class="container-fuild">) . . . . .	10
4.6	Gestionnaire de corpus . . . . .	11
4.7	Gestionnaire des néologismes candidats . . . . .	13
4.8	Gestionnaire des néologismes sémantiques . . . . .	13
4.9	Gestionnaire des néologismes . . . . .	13
4.10	Moteur de recherche . . . . .	13
<b>5</b>	<b>Détails des bases de données</b>	<b>13</b>
5.1	Détails de la base de données rssdata (sources des données) . . . . .	13
5.2	Détails de la base de données datatables (néologismes candidats) . . . . .	14
5.3	Détails de la base de données neo3 (base des néologismes validés) . . . . .	16
<b>6</b>	<b>Détails des collections Apache Solr</b>	<b>16</b>

# 1 Présentation générale des programmes

La figure 1 présente l'architecture générale de Néoveille. Elle comprend une série de programmes en back-end (tout ce qui se trouve au-dessus de la ligne horizontale pointillée) et une série de programmes en front-end (interface web, partie en-dessous de la ligne horizontale pointillée).

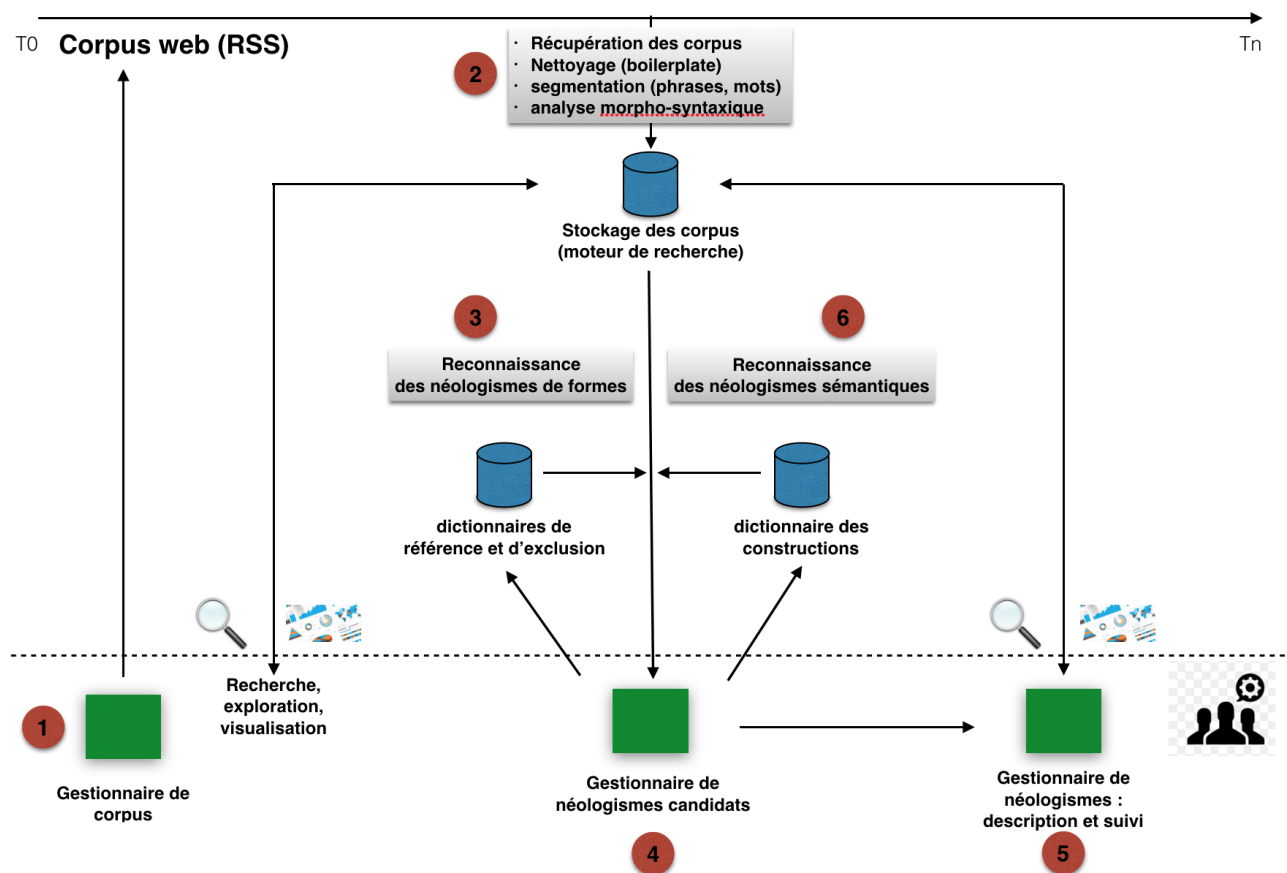


FIGURE 1 – Architecture de Néoveille

Bien qu'inter-agissant sur les mêmes données, les deux applications sont distinctes (backend et frontend).

Dans le schéma, les cylindres bleus représentent des ressources de deux types : corpus (stockés dans un moteur de recherche Apache Solr) et des ressources linguistiques de type dictionnaires. Les ressources linguistiques sont stockées dans des bases de données Mysql. On détaille l'architecture fonctionnelle dans la figure 2.

Nous détaillons ci-après comment installer l'interface web (frontend) puis les programmes en backend..

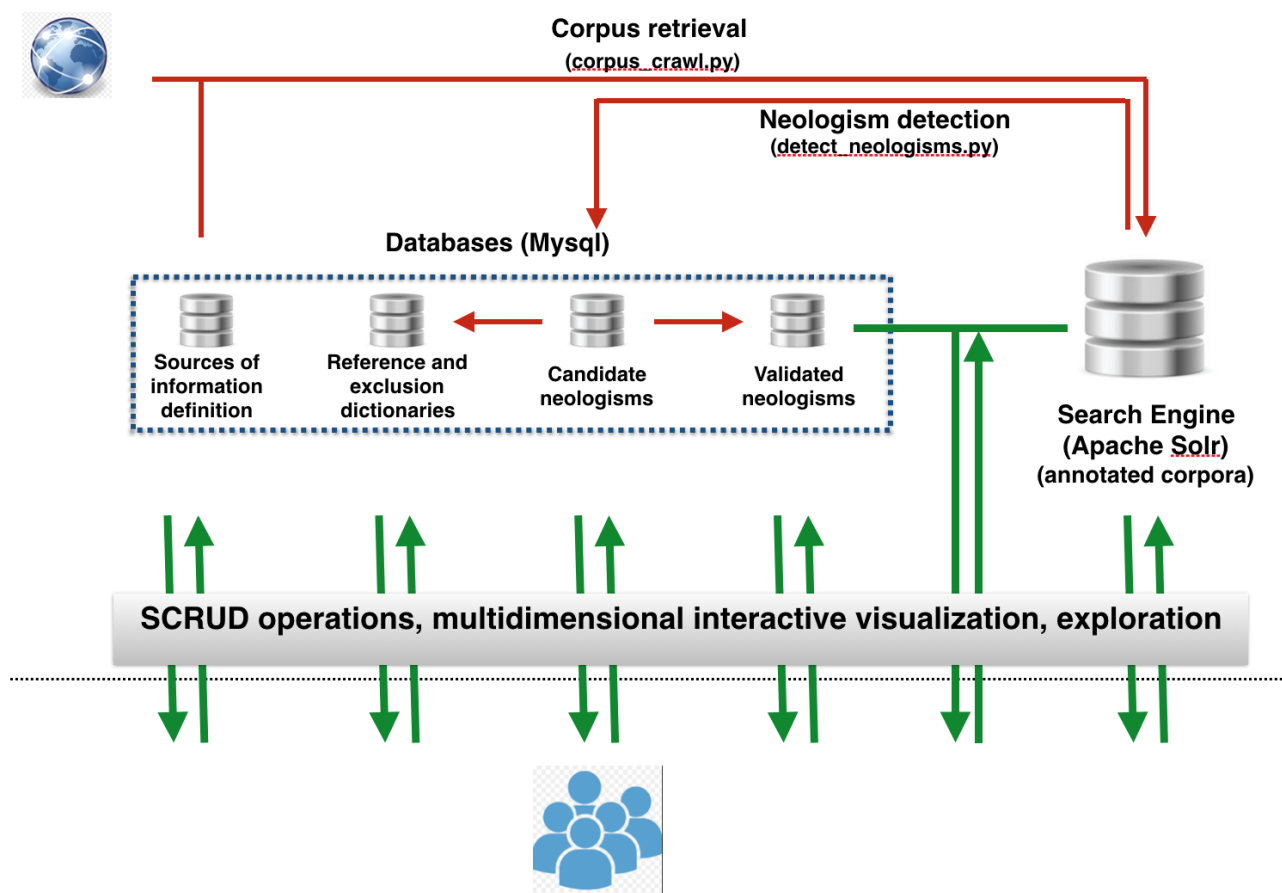


FIGURE 2 – Architecture fonctionnelle de Néoveille

## 2 Installation de l'interface web de Néoveille

Pour avoir tous les fichiers en local, faites un clone du dépôt github :

```
git clone https://github.com/ecartierlipn/neoveille2016.git
```

Les fichiers pour l'interface web se trouvent dans le sous-dossier "frontend".

### 2.1 Prérequis

L'installation a été testée sur une machine Linux Ubuntu 16.04 code Xenial.

**Serveur Apache** : un serveur web Apache 2 doit être installé avec une gestion de PHP (>version 5) et Mysql (> version 5.7).

On peut soit installer séparément ces différents éléments, soit les installer ensemble : voir procédure par exemple ici : <https://www.linode.com/docs/web-servers/lamp/install-lamp-stack-on-ubuntu-16-04/>

**Apache Solr** : le moteur de recherche Apache Solr doit également être installé (> 5.3 (version 5.5 actuellement utilisée), non testé sur les versions 7+) : <http://lucene.apache.org/solr/>

### 2.2 Procédure d'installation

Pour l'installation, il faut effectuer la suite d'opérations suivantes :

1. placer le contenu du dossier frontend dans un répertoire web du serveur Apache (par défaut : /var/www/).
2. modifier les paramètres d'accès au serveur Mysql dans le fichier credentials.php (qui se trouve, en prenant /var/www comme base d'installation des programmes du site ./html/html/credentials.php), à savoir la variable `$usermysql`, et la variable `$password`.
3. création de trois bases de données dans Mysql : rssdata (contiendra les informations sur les corpus à récupérer), datatables (contiendra les néologismes candidats et les dictionnaires d'exclusion), neo3 (contiendra les néologismes et l'ensemble des descriptions de ceux-ci).
4. lancer les trois scripts sql (qui se trouvent à `resources/mysql`) à partir de la racine du clone, pour créer les tables dans les trois bases de données<sup>1</sup>.

Ensuite, si le serveur Apache est lancé, il faut se rendre à : <http://localhost/html/html/index.php>. Vous n'avez plus qu'à naviguer sur le site (il y a un id/mot de passe par défaut : admin/neoveille2016).

## 3 Installation des programmes Néoveille en backend

Si ce n'est déjà fait, faites un clone du dépôt github :

```
git clone https://github.com/ecartierlipn/neoveille2016.git
```

Les programmes du backend se trouvent dans le sous-répertoire backend.

### 3.1 Pré-requis

Pour que les programmes de récupération et d'analyse des néologismes fonctionnent, il faut préalablement avoir installé :

- Apache Solr : une fois le programme installé, il faut créer une collection Solr (par défaut une par langue). Créer un répertoire `rss_french` dans `<répertoire d'installation solr>/server/solr`. Allez dans ce répertoire et dézipper le fichier de configuration exemple (`resources/apachesolr/solr_configs.zip`). Cela va créer un sous-répertoire `conf` avec les fichiers de configuration (le plus important est `schema.xml` qui contient les champs d'index). Il faut relancer le serveur Solr et vérifier que la collection a bien été chargée (elle sera vide).
- MySQL avec les bases de données listées ci-dessus
- TreeTagger pour l'analyse morphosyntaxique du français : voir <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>.
- Hunspell avec les fichiers de la langue traitée. Par défaut installé sur linux. Les fichiers dicos pour le français sont disponibles dans `resources/hunspell`.
- librairies pour Python : **à vérifier en lançant le programme.**

Les programmes Python ont été confectionnés en python 2.7. le portage est possible en python 3.4+ mais nécessite de changer plusieurs librairies, notamment pour la récupération web.

1. Pour lancer les scripts, saisir (exemple de neo3) : `mysql -u <id> -p <password> neo3 < neo3.complete.truct.sql`

## 3.2 Détails des fichiers du backend

Le répertoire backend contient tous les fichiers du backend. Il contient un sous-répertoire pour crawl des corpus (`corpus_crawler`) et un sous-répertoire pour la détection des néologismes formels (`formal_neology`).

Le répertoire `corpus_crawler` contient :

- **corpus\_crawler** : contient les fichiers principaux (fichiers `.py`) pour récupérer les corpus par langue (`corpus_<lang>.py`), et les `.sh` pour faciliter l'installation d'une tâche cron (`corpus_crawl_<lang>.sh`). le fichier `URLutils.py` est une librairie additionnelle utilisée par les fichiers `corpus_<lang>.py`.
- **corpus\_crawler/log** : contient les fichiers de log qui sont créés chaque fois que les programmes python sont lancés. Voir le nom de ces fichiers de log déclarés en fin des fichiers `corpus_<lang>.py` ;
- **corpus\_crawler/tobeindexed** : contient les fichiers à indexer pour Apache Solr. Répertoire à nettoyer de temps à autres.
- **corpus\_crawler/webcrawl** : travail en cours pour récupérer des données sur le web sans le support des fils RSS..

Le répertoire `formal_neology` contient :

- **fichiers Python** : contient les fichiers principaux (fichiers `.py`) pour détecter les néologismes par langue (`detect_neologisms_<lang>.py`, all pour le français), et les `.sh` pour faciliter l'installation d'une tâche cron (`detect_neologisms_<lang>.sh`) ;
- **formal\_neology/log** : contient les fichiers de log qui sont créés chaque fois que les programmes python sont lancés. Voir le nom de ces fichiers en fin des `detect_neologisms_<lang>.py` ;
- **formal\_neology/tobeindexed** : contient les fichiers à indexer pour Apache Solr. Répertoire à nettoyer de temps à autres.

## 3.3 Récupération dynamique des fichiers sur le web : fichier `corpus_fr.py`

Le programme `corpus_fr.py` effectue les opérations suivantes : il récupère, pour les sources définies dans la table `rssdata/RSS_INFO` (dans cette table il récupère toutes les métainfos + le champ `NAME_RSS`, ie url du fil RSS)) les fichiers qui sont listés dans les fils rss : il récupère également les méta-infos de l'item du fil rss (title, author, description, url, keywords, etc.) puis va chercher l'article complet via le champ url. Il effectue un nettoyage de la page (boilerplate removal, avec `justext`) pour ne conserver que les zones textes utiles. Ensuite, le programme stocke le texte et toutes les métainformations dans la table `rssdata/rss_data2`, en donnant au champ `IS_INDEXED` la valeur 0, indiquant que le fichier doit encore être traité pour repérer les néologismes. Il stocke les mêmes informations dans le moteur Apache Solr. NB - le stockage dans la table `rssdata_2` est temporaire (il faudrait prévoir un nettoyage régulier de cette table, actuellement aucun process ne s'en occupe).

Le programme génère un fichier de log qui doit permettre d'identifier les erreurs. Il se trouve dans le sous-répertoire `log`.

### 3.3.1 Paramètres à régler dans le fichier

:

- **accès à la base de données MySQL** : à modifier dans trois méthodes : `get_corpus_list_fromDB`, `get_last_indexed_fromDB` (dans la classe `corpus`) et dans `save_rsscorpus_to_DB` (classe `rssfeeds`).
- **accès à la collection Apache Solr** : à modifier dans la variable `lang_solr2`. (par défaut, pour le français, la collection s'appelle `rss_french`).

### 3.3.2 Installation du fichier `corpus_fr.py` en tâche récurrente

: Le plus simple sous linux est de créer une tâche CRON :

```
1 crontab -e
```

et (par exemple, en utilisant le script shell fourni, pour lancer le programme chaque nuit à 0 :14) :

```
1 0 14 * * * <chemin vers le fichier>corpus_crawl_fr.sh 2> <chemin vers le fichier>/  
errors_corpus_crawl_all_crontab.txt
```

### 3.3.3 Structure du fichier `corpus_fr.py`

Le fichier se compose de variables globales, de deux classes (`corpus` et `rssfeeds`) et d'une fonction `main()` détaillée ci-dessous.

```

1  for lang in ['Rép. Tchèque', 'pologne', 'france']:
2      c = corpus(lang, 'db', 'rss') # init. corpus avec info langue, type stockage info
3          et type fichiers
4      c.get_corpus_list_fromDB() # recup des fils rss à recuperer dans bd
5      if str(len(c.rssfeeds))>0:
6          c.get_last_indexed_fromDB() # recup des fils rss recuperes (>un jour avant)
7          c.retrieve_corpus()# recup des nouveaux articles et stockage dans bd et
8              Apache Solr
9      else :
10         log.info("No RSS feeds for this language. Check your configuration" )
11         log.info("All is done! Quitting program.")

```

### 3.4 Détection des néologismes, stockage des néologismes candidats et mise à jour des articles dans Apache Solr : fichier detect\_neologisms\_all.py

Le programme `detect_neologisms_all.py` effectue les opérations suivantes (pour le français) : il récupère dans la table `datatables/rssdata_2` les articles à analyser, il les analyse morphosyntactiquement (avec Treetagger), il récupère l'analyse morphosyntaxique du texte qu'il uniformise, y récupère les mots inconnus, fait un passage par Hunspell pour enlever des candidats néologismes les coquilles, ensuite filtre les mots qui se trouvent dans le dictionnaire d'exclusion, enfin stocke les candidats néologismes dans la table `datatables/neologismes_fr` et met à jour l'article dans Apache Solr avec l'analyse morphosyntaxique.

#### 3.4.1 Paramètres à régler dans le fichier

- :
  - **accès à la base de données MySQL** : à modifier dans deux fonctions : `get_corpus_list_fromDB`, et dans `add_neologisms_to_db`.
  - **accès à la collection Apache Solr** : à modifier dans la variable `lang_solr2`. (par défaut, pour le français, la collection s'appelle `rss_french`).
  - **Accès à l'exécutable paramétré de Treetagger** : il s'agit des shell script du Treetagger lançant le programme avec une configuration de langue spécifique. Par exemple pour le français, il s'agit du fichier : '`<chemin d'installation Treetagger>/cmd/tree-tagger-french-utf8`'. On explicite le chemin vers ce fichier dans la variable : `pos_ana['france']`;
  - **Accès aux dictionnaires Hunspell** : indiquer le chemin pour accéder aux dictionnaires Hunspell. Par exemple pour le français : `hunspell['france']='<chemin>/opt/nlp_tools/dictionaries/fr_FR/fr_FR'`.

#### 3.4.2 Installation du fichier detect\_neologisms\_all.py en tâche récurrente

: Voir programme précédent.

#### 3.4.3 Structure du fichier detect\_neologisms\_all.py

Le fichier se compose de variables globales, d'une classe (document), d'une série de fonctions utilitaires et d'une fonction `main()` détaillée ci-dessous.

```

1  conn=''
2  langs = ['pologne', 'france', 'brésil', 'Rép. Tchèque']
3  for lang in langs:
4      res = get_corpus_data_fromDB(lang)
5      if res:
6          load_exclusion_dico(iso[lang])
7          i=0
8          ### preparation of solrxmlfile
9          now = datetime.now().strftime("%d-%m-%y_%H")
10         filenameXML = './tobeindexed/update_rss-data.' + lang_solr[lang] + "." + now
11             + '.xml'
12         fout = codecs.open(filenameXML, mode="w+", encoding="utf-8")
13         fout.write("<add>")
14         for data in res: # on parcourt les fichiers à analyser
15             i+=1
16             metas={}
17             source_link = data[0]
18             metas['title'] = data[1]
19             metas['subject'] = data[2]
20             metas['category'] = data[3]
21             #metas['description'] = data[4]
22             cts = data[1] + "\n" + data[5]

```

```

22 contents1 = re.sub(r'[ńž]', r' ', unicode(cts), flags=re.UNICODE)
23 contents = re.sub(r'\s+', r' ', unicode(contents1), flags=re.
    UNICODE)
24 metas['ID_RSS'] = data[6]
25 d = document(source_link, lang, metas, contents)
26 langd = d.detect_language(contents)
27 if (langd not in lang_detect[lang]):
28     log.info("Problem with language detection for contents : "
        + unicode(contents) + "\nAutomatic detection says : [" +
        langd + "] whereas expected language is " + str(
        lang_detect[lang]) + "\nSkipping analysis for this
        document and deleting it from database.")
29     continue
30 d.ling_analyze('pos_tagging') ## to be done : specific return value
    if fails
31 #print str(d)
32 ### retrieve info so as to create a dict for saving into solr xml
    format
33 doc = {}
34 doc['link']=source_link
35 if d.postagger==False:
36     log.info("Pos tagger has returned false for content")
37     continue # go to next document
38 else:
39     doc['lemmes']= d.lemmas
40     doc['noms_propres']= d.np
41     doc['pos-text']= d.postaggedText # too heavy for apache
        solr 15-12-2016
42 # neo exclusion list
43 if len(d.unk)>0: # if some unknown words
44     log.info(d.unk)
45     unk2=d.filter_unknown(d.unk) # hunspell and
        exclusion dictionaries filtering
46     if len(unk2)>0:
47         add_neologisms_to_db(unk2, iso[lang])
48         # update command for apache solr
49         updatecmd={'lemmes': 'add', 'noms_propres': '
            add', 'pos-text': 'set'}
50         res = write_xml_chunk(fout, doc, lang_solr[
            lang], updatecmd)
51         if res:
52             # if solr indexing ok, mark text as
                1 in db table
53             db_tag_doc_as_processed(source_link
                ,1)
54             else: # else mark as 2
55             db_tag_doc_as_processed(source_link
                ,2)
56             # no neologismes : index-update anyway with other
                linguistic information (lemmes, noms-propres)
57         else:
58             log.info("Pas de mot inconnu dans : " + source_link
                + ". indexing anyway.")
59             updatecmd={'lemmes': 'add', 'noms_propres': 'add', 'pos
                -text': 'set'}
60             res = write_xml_chunk(fout, doc, lang_solr[lang],
                updatecmd)
61             if res:
62                 db_tag_doc_as_processed(source_link,1)
63             else:
64                 db_tag_doc_as_processed(source_link,2)
65             # no unknown words but change IS_INDEXED to value 1
66         else:
67             log.info("Pas de mot inconnu dans : " + source_link + ".
                indexing anyway.")
68             updatecmd={'lemmes': 'add', 'noms_propres': 'add', 'pos-text': '
                set'}
69             res = write_xml_chunk(fout, doc, lang_solr[lang], updatecmd)
70             #
                continue
71             if res:
72                 db_tag_doc_as_processed(source_link,2)
73             else:
74                 db_tag_doc_as_processed(source_link,2)
75 fout.write("</add>") # level of for data in res
76 fout.close()
77 # now indexing for lang

```

```

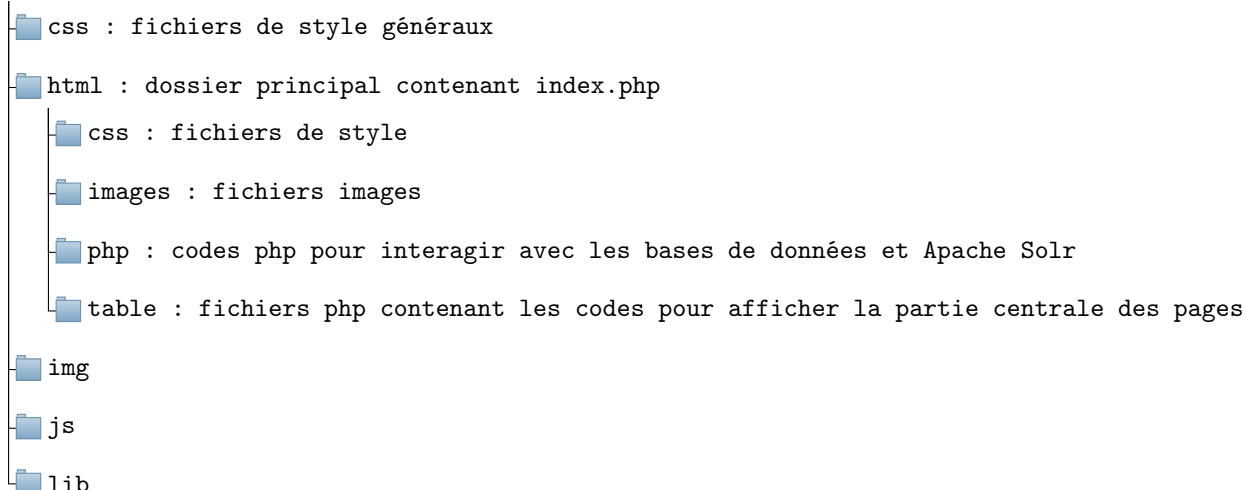
78         res = index_data_solr(filenameXML, lang_solr[lang], "xml")
79     if res :
80         log.info("Indexation succeeded for file : " + filenameXML)
81     else :
82         log.error("Indexation failed for file : " + filenameXML)
83
84 else :
85     log.info("No corpus data to process for " + lang)
86     log.info("All is done. Exiting.")

```

## 4 Détails des fichiers : interface web

### 4.1 Organisation des dossiers et fichiers principaux

Les fichiers sont, à partir de la racine d'installation, répartis dans des sous-dossiers :  
(racine)



Nous présentons ci-après l'organisation des pages web. Nous présentons tout d'abord la structure de la page d'accueil et le mécanisme général d'organisation des fichiers.

### 4.2 Organisation page d'accueil

La page d'accueil (index.php quand on est connecté, login.php dans les autres cas) se présente sous forme de quatre zones (figure 3). Le bandeau supérieur est couvert par l'élément `<nav class="navbar navbar-default...">`, le menu gauche par l'élément `<div class="side-menu...">`, le contenu principal par l'élément `<div class="container-fluid...">`, et le pied de page par `<footer class="app-footer">`.

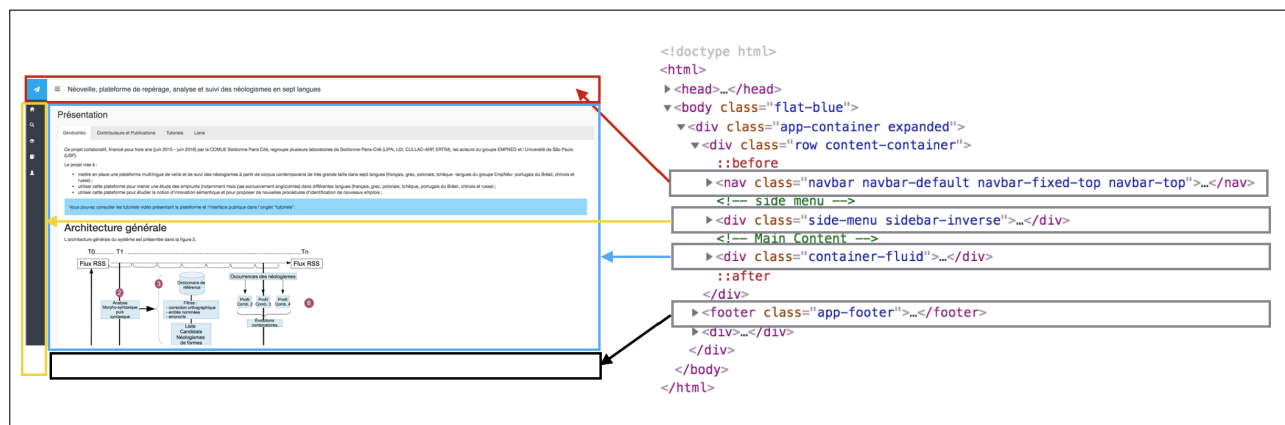


FIGURE 3 – Organisation page d'accueil (login.php et index.php)

#### 4.2.1 Feuilles de style

Les feuilles de style associées sont (à vérifier...) déclarées dans le header, il s'agit d'un mixte entre bootstrap.css et des thèmes définis dans `../css/style.css` et `../css/themes/flat-blue.css`.



Il serait nécessaire d'uniformiser les styles en ne prenant qu'une librairie basée sur bootstrap.

#### 4.2.2 Bibliothèques javascript

Les bibliothèques javascript chargées sont les suivantes (figure 4).

```
<!-- Javascript Libs -->
<!-- jquery -->
<script type="text/javascript" src="../../lib/js/jquery.min.js"></script>
<!-- bootstrap -->
<script type="text/javascript" src="../../lib/js/bootstrap.min.js"></script>
<!-- datatables -->
<script type="text/javascript" charset="utf-8" src="js/dataTables.min.js">
<!-- datatables style |bootstrap -->
<script type="text/javascript" src="../../lib/js/dataTables.bootstrap.min.js"></script>
<!-- datatables editor -->
<script type="text/javascript" charset="utf-8" src="js/dataTables.editor.min.js"></script>

<!-- bibliothèques spécifiques datatables -->
<script type="text/javascript" charset="utf-8" src="js/editor.title.js"></script>
<script type="text/javascript" src="js/jquery.dataTables.columnFilter.js"></script>
<script type="text/javascript"
src="https://cdn.datatables.net/buttons/1.3.1/js/dataTables.buttons.min.js"></script>
<script type="text/javascript"
src="https://cdn.datatables.net/buttons/1.3.1/js/buttons.flash.min.js"></script>
<script type="text/javascript"
src="https://cdn.datatables.net/buttons/1.3.1/js/buttons.print.min.js"></script>
<script type="text/javascript"
src="https://cdn.datatables.net/buttons/1.3.1/js/buttons.html5.min.js"></script>
<script type="text/javascript">
```

FIGURE 4 – Liste bibliothèques javascript chargées dans login.php et index.php

#### 4.3 Bandeau supérieur (<nav class="navbar...">)

Le lien entre les composants de l'élément navbar et les zones affichées est donné dans la figure 5.

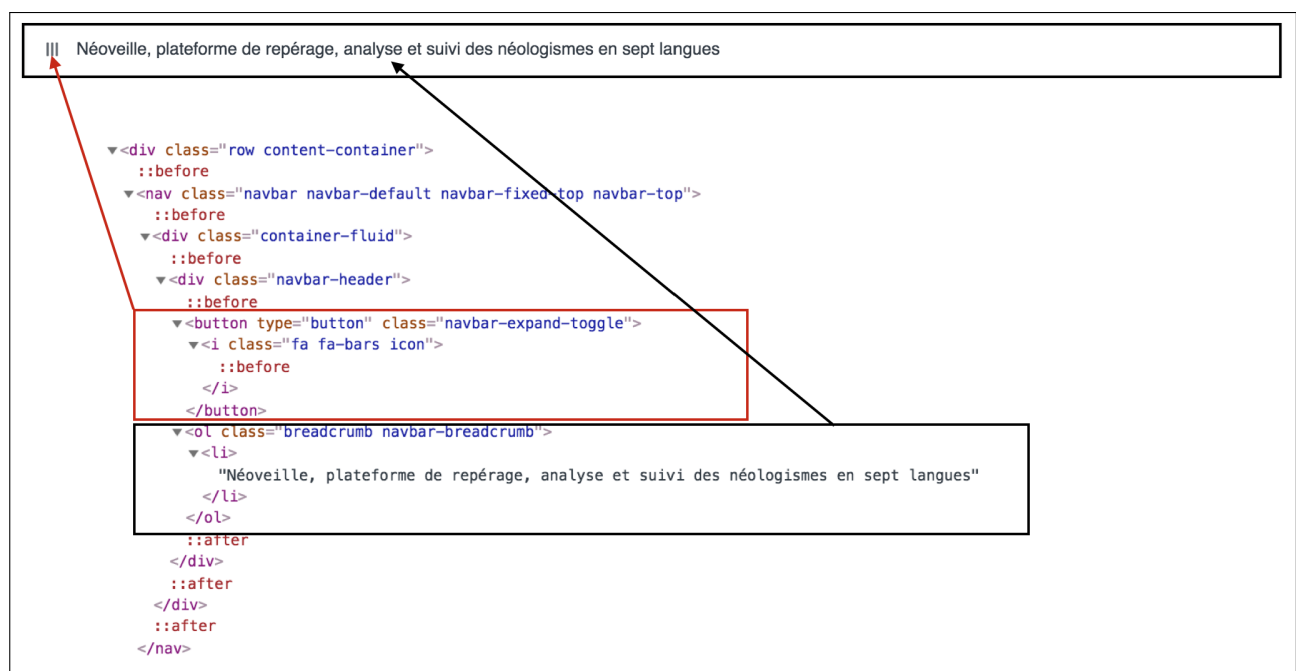


FIGURE 5 – Organisation Barre de navigation (login.php et index.php)

≡ Néoveille, plateforme de repérage, analyse et suivi des néologismes en sept langues



#### 4.4 Menu gauche (<div class="side-menu">)

#### 4.4.1 Action Javascript associée aux éléments de menu

Deux autres fonctionnements plus simples consistent :

- #### 4.5 Page de contenu (<div class="container-fuild">)

10

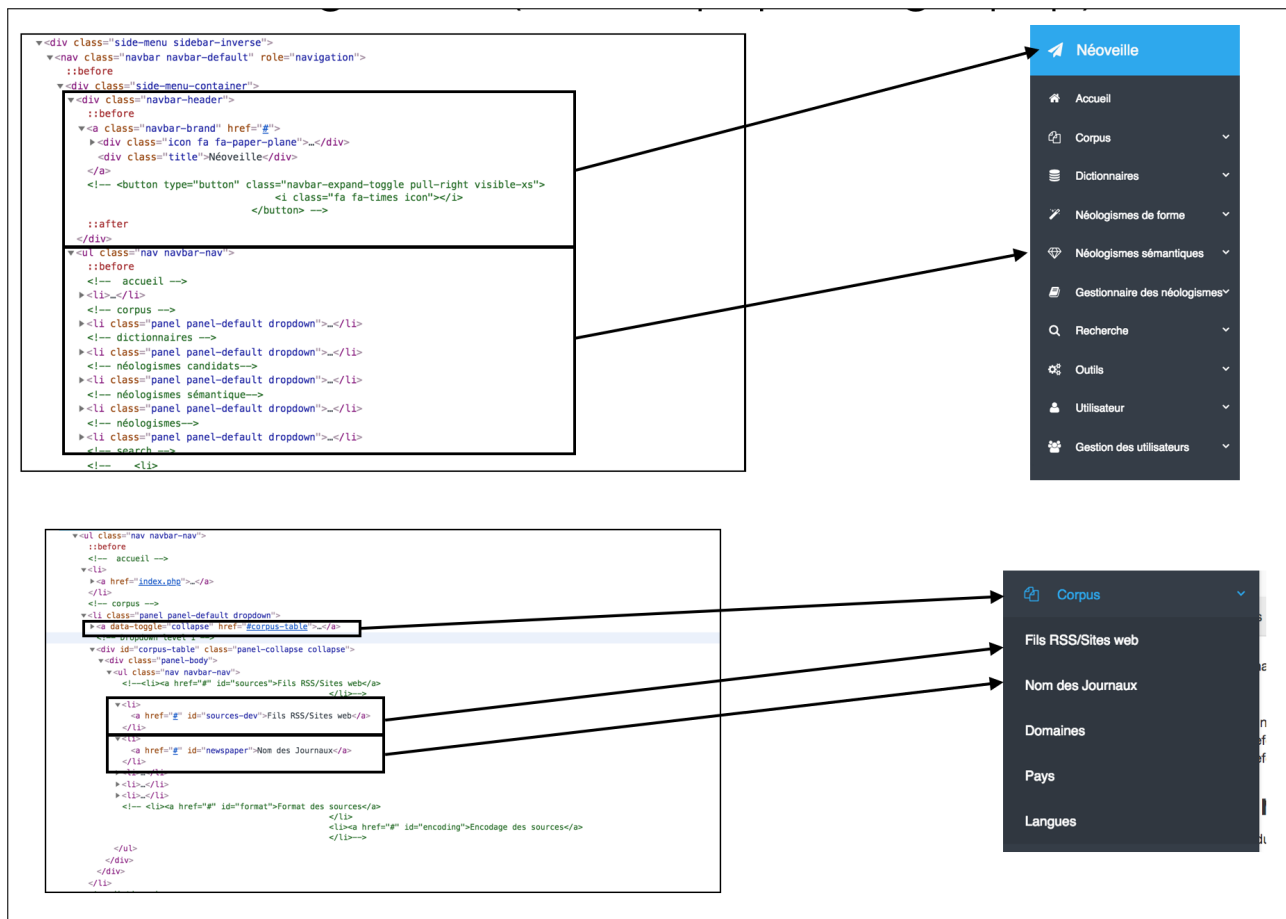


FIGURE 7 – Organisation Menu gauche (login.php et index.php)

bootstrap. La figure 11 montre la structure de base de cet élément (`<div class="side-body">`). dans le `<div class="page-title">` on inclut le titre, et dans le `<div class="row">` on inclut le contenu lui-même.

Lorsque le contenu est dans un fichier php (qui se trouve toujours dans le sous-répertoire "table") il est généralement lui-même lié à deux autres fichiers : un fichier javascript qui va prendre en charge les différentes actions associées et la création des tables éditables (via les bibliothèques `datatables.js` et `datatables.editor.js`), cette bibliothèque s'occupant de charger les données à afficher/éditer, et de mettre en place les éléments d'interaction (tri des tables, recherche, navigation, édition).

**Remarque :** on peut retrouver les fichiers liés à chaque des éléments de menus, avec les outils de développements web sous Google Chrome, dans l'onglet Network. Pour plus d'infos sur l'utilisation de ces outils : <https://www.malekal.com/chrome-firefox-outils-de-developpement/>.

## 4.6 Gestionnaire de corpus

Le menu Corpus/Fils RSS/Sites web permet d'accéder à l'interface présentée dans la figure ?? . Cette page est chargée en chargeant le fichier `table/datatables-corpus-dev.php`, qui lui-même charge le fichier javascript `js/table.RSS_INFO-dev.js`, lui-même chargeant le fichier `php/table.RSS_INFO.php`. Le rôle respectif de ces trois fichiers est détaillé ci-dessous ;

- **table/datatables-corpus-dev.php** : fichier chargeant les éléments html à placer dans la zone principale d'affichage.
- **js/table.RSS\_INFO-dev.js** : fichier chargé via le précédent et contenant la définition initiale de la grille à afficher (`datatables.editor`), et les actions/fonctions associées aux éléments.
- **php/table.RSS\_INFO.php** : fichier qui assure l'interaction entre `datatables.editor` et les données elles-mêmes, stockées dans une base de données (en l'occurrence la table `rssdata/RSS_INFO` et les tables liées). Le fichier php fait partie du framework `datatables.editor`.

Les bibliothèques javascript pour la visualisation interactive des résultats sont également chargées par le fichier `table/datatables-corpus-dev.php`.

- **../js/d3.js** : bibliothèque pour créer des graphes en svg dans des pages html ;
- **../js/dc.js** : bibliothèque utilisant `d3.js` pour générer plusieurs graphes en interaction ;

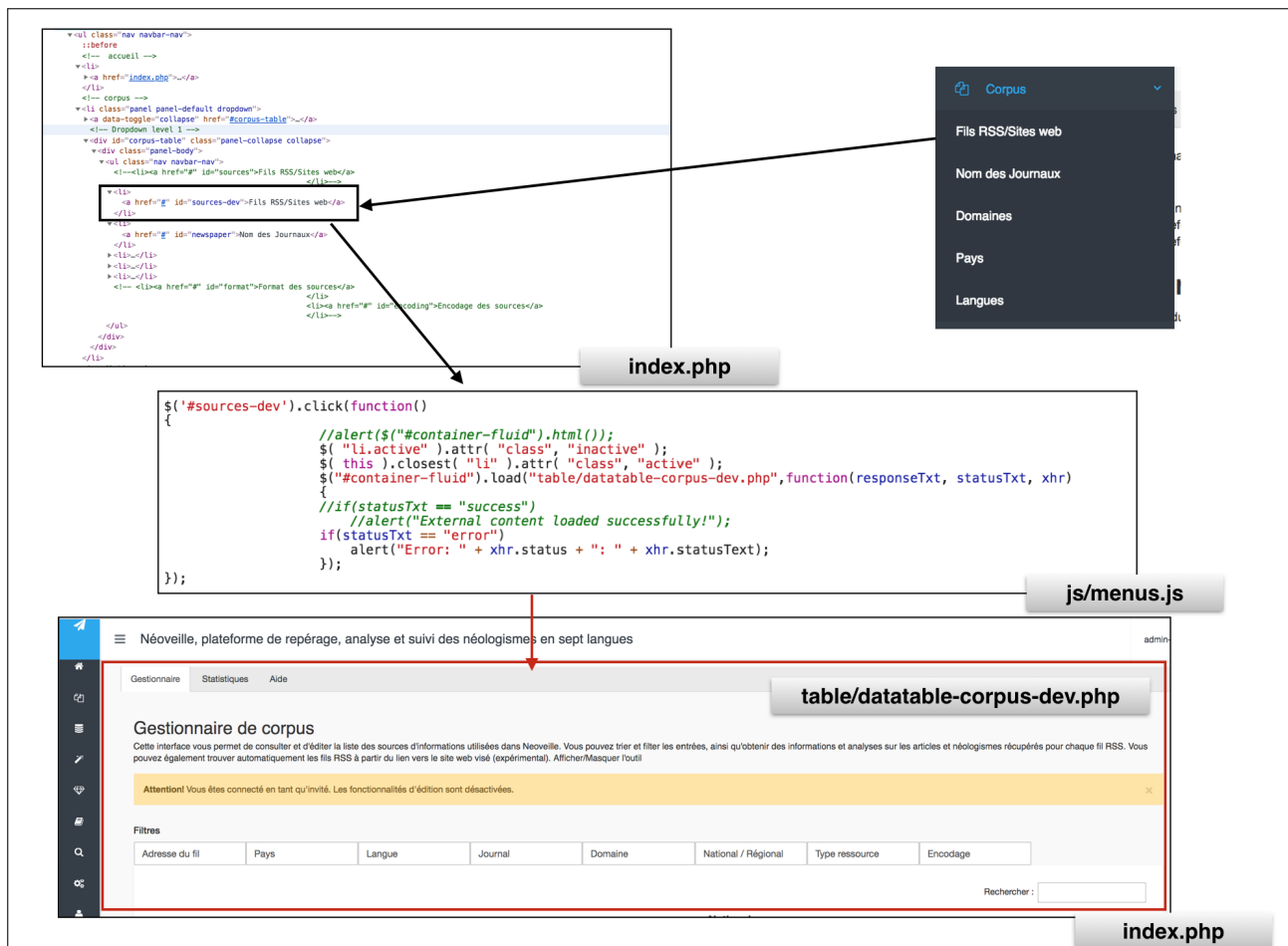


FIGURE 8 – Détail de l'action de menu (exemple corpus)

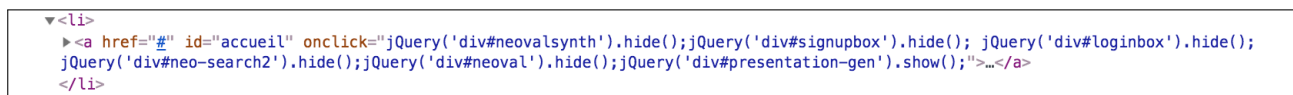


FIGURE 9 – Action de menu directement dans l'élément (exemple accueil)



FIGURE 10 – Action javascript de menu intégré au fichier login.php (exemple search)

- `../js/crossfilter.js` : librairie pour agréger les données brutes provenant de d3.js et permettant l'interaction entre les graphes résultants;

```

<!-- Main Content -->
▼<div class="container-fluid" id="container-fluid">
  ::before
  ▼<div class="side-body"> == $0
    ►<div class="page-title">...</div>
    ►<div class="row">...</div>
  </div>
  ::after
</div>
::after
</div>

```

FIGURE 11 – Structure globale de l'élément `<div class="container-fluid">`

Gestionnaire des néologismes candidats

Cette interface vous permet de consulter et d'éditer la liste des sources d'informations utilisées dans Neoville. Vous pouvez trier et filtrer les entrées, ainsi qu'obtenir des informations et analyses sur les articles et néologismes récupérés pour chaque fil RSS. Vous pouvez également trouver automatiquement les fils RSS à partir du lien vers le site web visé (expérimental). Afficher/Masquer l'outil

Filtres

Adresse du fil Pays Langue Journal Domaine National / Régional Type ressource Encodage

Nouveau Modifier Afficher 10 éléments Rechercher :

Adresse du fil	Pays	Langue	Journal	Domaine	Fréquence	National / Régional	Type ressource	Encodage		
http://rss.usinenouv...	France	Français	L'Usine Nouvelle	Industrie	hebdomadaire	National	rss	utf-8	○	✎
http://ticetsociete...	France	Français	TIC&Société	Informatique	hebdomadaire	National	rss	utf-8	○	✎
http://www.inserm.fr...	France	Français	Science et Santé (Inserm)	Recherche	hebdomadaire	National	rss	utf-8	○	✎
http://www.inserm.fr...	France	Français	Science et Santé (Inserm)	Société	hebdomadaire	National	rss	utf-8	○	✎
http://www.lcp.fr/rs...	France	Français	LCP	Politique	hebdomadaire	National	rss	utf-8	○	✎
http://www.lemondein...	France	Français	Le Monde Informatique	Informatique	hebdomadaire	National	rss	utf-8	○	✎
feed//lematin.ma/co...	Maroc	Français	Le Matin	Général	quotidien	National	rss	utf-8	○	✎
http://actusen.com...	Sénégal	Français	Actusen	Général	quotidien	National	web	utf-8	○	✎
http://algeriasong.o...	Algérie	Français	Tadakt - le Kabyle Magazine	Général	quotidien	National	rss	utf-8	○	✎
http://animeand.com...	France	Français	AnimeLand	Langue des jeunes	quotidien	National	web	utf-8	○	✎

Affichage de l'élément 1 à 10 sur 474 éléments

Précédent 1 2 3 4 5 ... 48 Suivant

FIGURE 12 – Exemple de contenu dans le fichier `table/datatables-corpus-dev.php`

## 4.7 Gestionnaire des néologismes candidats

## 4.8 Gestionnaire des néologismes sémantiques

## 4.9 Gestionnaire des néologismes

## 4.10 Moteur de recherche

# 5 Détails des bases de données

Trois bases de données permettent de stocker les différentes informations :

- **Les sources d'informations** : base de donnée rssdata.
- **Les néologismes candidats et les dictionnaires de référence "utilisateurs" et d'exclusion** : base de données datatables.
- **La base de description des néologismes** : base de données neo3.

## 5.1 Détails de la base de données rssdata (sources des données)

Cette base de données est dédiée à la représentation des sources de données (corpus). Elle est organisée autour de la table `RSS_INFO`, avec des clés étrangères liées aux différentes informations (voir figure 13). Ce modèle est extensible, si l'on souhaite ajouter de nouvelles métadonnées aux sources d'information.

Actuellement, le fichier `RSS_INFO` comprend les sources d'informations pour toutes les langues (voir champ langue).

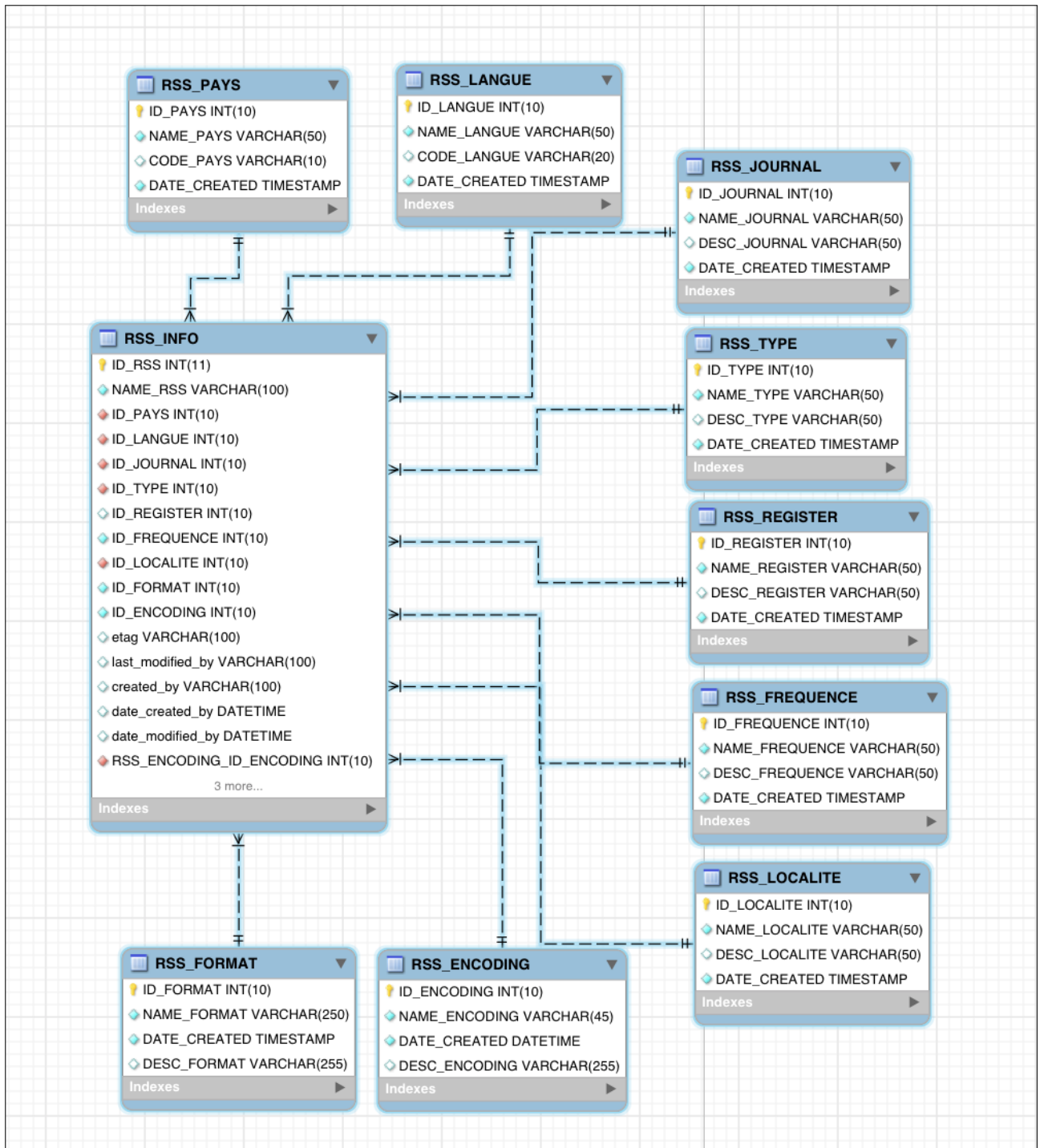


FIGURE 13 – Modèle de données MySQL pour rssdata

## 5.2 Détails de la base de données datatables (néologismes candidats)

Cette base de données est dédiée à la représentation des néologismes candidats repérés par le programme backend (detect\_neologisms.py), d'une part (table neologismes) et à la représentation des dictionnaires de référence "utilisateurs" (dico simple, dico composé, dico termino et dico prefixes et suffixes, ces deux derniers non-utilisés actuellement mais qui pourraient l'être pour détecter les affixations) et dictionnaires d'exclusion (excluded\_fr). La même structure est suivie pour chaque langue (avec le suffixe \_<lang>).

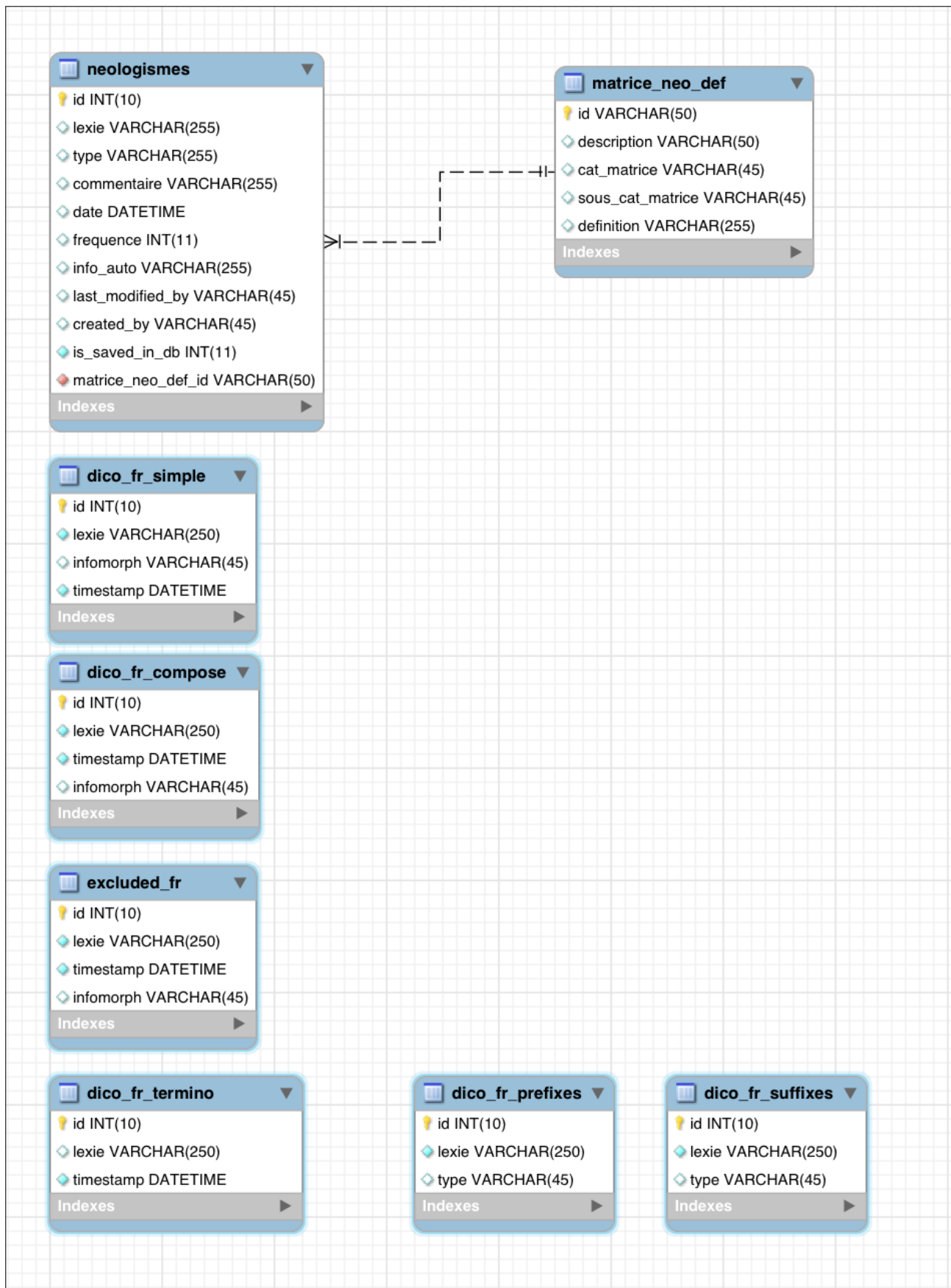


FIGURE 14 – Modèle de données MySQL pour datatables



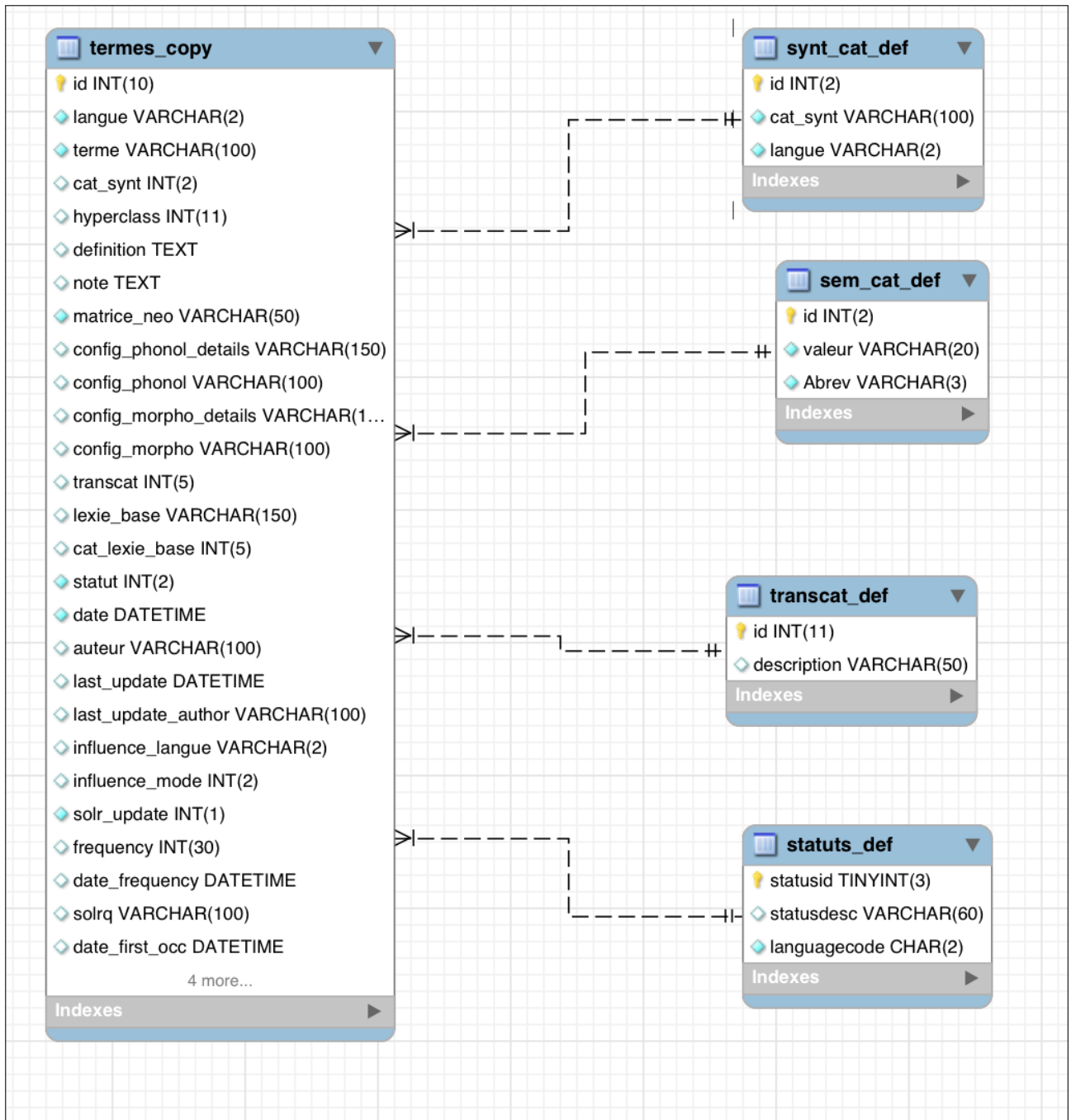


FIGURE 15 – Modèle de données MySQL pour neo3

### 5.3 Détails de la base de données neo3 (base des néologismes validés)

## 6 Détails des collections Apache Solr

Les corpus récupérés sont stockées dans une collection dans Apache Solr. nous présentons ci-dessous l'architecture générale de ce système et le format des collections (susceptible d'évoluer). La version d'Apache Solr utilisée est actuellement la 5.3.