



and the second



Est-ce que cette information est sérieuse?



Tracer l'origine et le résultat de cette information...





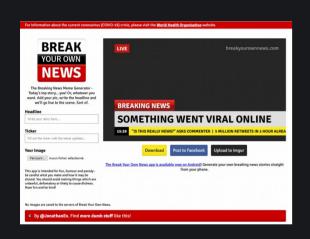


Boire du méthanol, de l'éthanol ou de l'eau de Javel NE PERMET PAS de prévenir ou de guérir la COVID-19 et peut être extrêmement dangereux

Le méthanol, l'éthanol et l'eau de Javel sont des poisons. Leur ingestion peut entraîner des lésions voire la mort. Le méthanol, l'éthanol et l'eau de Javel sont parfois utilisés dans des produits d'entretien qui servent à détruire le virus sur les surfaces, mais il ne faut jamais en boire. Ils ne détruisent pas le virus dans l'organisme et entraînent des lésions des organes internes.

Pour se protéger de la COVID-19, il faut désinfecter les objets et les surfaces, en particulier ceux que vous touchez régulièrement. Lavez-vous les mains souvent et soigneusement et évitez de vous toucher les yeux, la bouche et le nez.

... Accéder à de meilleurs points de vue...







... qualifier la source (un site parodique) ...

« L'alcool ne tue pas le coronavirus »





... c'est construire le savoir autour de l'information. c'est la valider ou la rejeter.

Et c'est automatisable!

Un déluge d'information...

Aujourd'hui, tout le monde peut publier du contenu en ligne et être lu/vu.

Mais l'économie de l'attention favorise les contenus émotionnels (au dépend des contenus d'experts).

Seuls 24 % des Français interrogés indiquent faire confiance aux médias (TV, presse papier et en ligne) – chiffre le plus bas en Europe.



... qui favorise la diffusion des fausses informations

Baisse de la confiance générale

Polarisation de la société

Défiance envers les marques, les institutions et les médias

Coût humain, politique et financier (\$78Milliards/an)



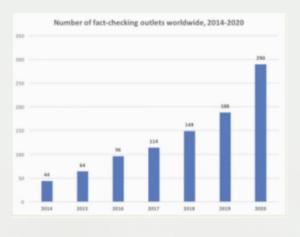
Qui lutte contre les fausses infos?

Des communautés dédiées s'organisent

Mais les fact checkers pros sont débordés

Difficultés pour le grand public qui n'a pas les méthodes

Pas d'outils dédiés, surtout pour le monde francophone





l'équipe

Jean-Marc Guerin, fondateur

13 ans d'expérience à concevoir des solutions d'Intelligence Artificielle dans des startups

Passionné par le fact checking et l'intelligence collaborative



Validalab

s'informer en confiance



L'app Validalab

C'est le moteur de recherche qui accélère le fact checking :



Qualifier les sources

Accéder à de meilleurs points de vue

Tracer les citations

Validalab trouve

- o ce qui fait qu'une information est fiable ou
- ⊗ s'il s'agit d'une fausse rumeur

Information fiable ou fausse rumeur?

Pour tout type de contenu

- Grand public
- Factcheckers
- Producteur de contenu

- ❷ Web et réseaux sociaux
- ❷ Vidéos, tweets, articles, et documents
- ❷ Spécialisés ou grand public

Pour tout profil d'utilisateur

Collaboratif et totalement scalable – ready for infowar



Le Validagraph : concept





L'information est fiable si toute la chaîne d'information est traçable et crédible

Le Validagraph : pratique





Recontextualisation

Validalab scanne et recontextualise les sites internets, blogs, twitter, youtube sous forme de graph de connaissance : le Validagraph.



Validation

L'objectif du Validagraph est de vous fournir simplement et rapidement toutes les informations dont vous avez besoin pour comprendre et valider ou non votre information.

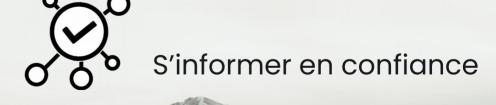


Participation

Vous pouvez participer à la construction du Validagraph en partageant vos résultats de façon anonyme ou en créant votre profil pour suivre votre crédibilité.

Validalab, c'est

- Ou un moteur de recherche basé sur la crédibilité
- O Des outils de contextualisation pour mieux comprendre
- O Des résultats adaptés à l'utilisateur
- Une plateforme de collaboration pour amateurs et professionnels





Nalidalab

s'informer en confiance

Jean-Marc Guerin, fondateur jm@validalab.fr





Présentation technique

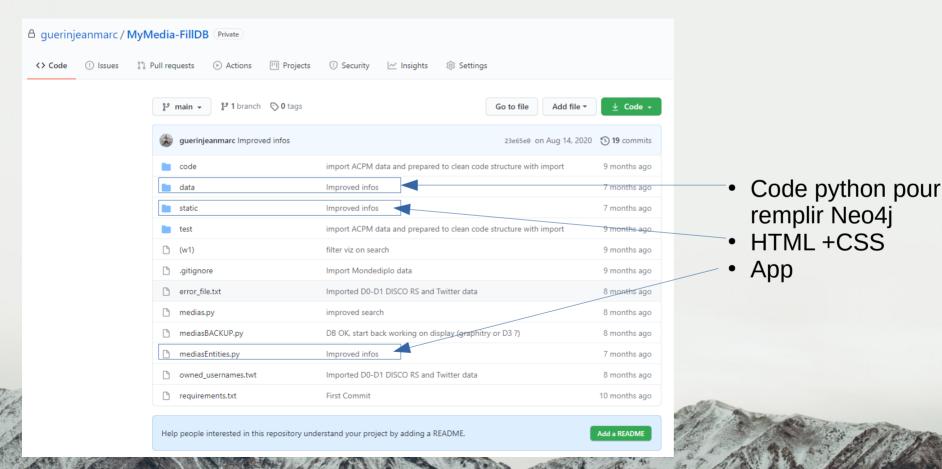
- 1.L'app Validalab, lancement et demo
- 2.La base Neo4j
- 3.Le code pour remplir la base
- 4.Le code de l'app

1. L'app Validalab

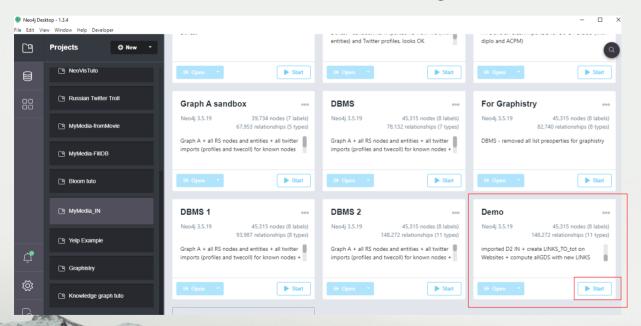
L'app Validalab

- Base Neo4j
- Code python + flask
- Code html + CSS

Repo Github



Lancer la base Neo4j



Activer l'environnement virtuel venv

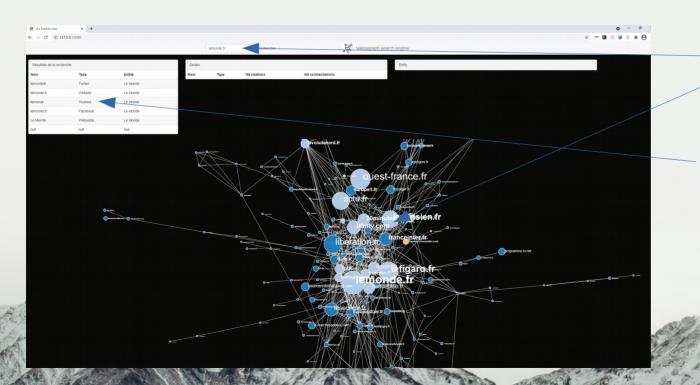
```
Directory of C:\Users\Jo\Documents\Tech\Atom prj\MyMedia-FillDB
21/08/2020 20:37
                    <DIR>
21/08/2020 20:37
16/06/2020 16:25
                                 0 (w1)
07/07/2020 18:38
                               162 .gitignore
08/07/2020 17:21
                    <DIR>
                                   code
10/12/2020 07:47
                    <DIR>
                                   data
                            37,096 error file.txt
23/07/2020 11:42
10/08/2020 18:40
                            5.541 medias.pv
                            3,901 mediasBACKUP.py
05/08/2020 17:35
08/10/2020 11:14
                            29,071 mediasEntities.py
23/07/2020 00:02
                            10,168 owned usernames.twt
18/05/2020 18:41
                               804 requirements.txt
25/09/2020 07:55
                    <DIR>
                                   static
10/07/2020 17:47
                    <DIR>
                                   test
19/08/2020 13:16
                             9,812 test.csv
06/07/2020 16:10
                    <DIR>
                                   venv
05/08/2020 18:20
                    <DIR>
                                   pycache
                                96,555 bytes
              9 File(s)
              8 Dir(s) 59,016,966,144 bytes free
C:\Users\Jo\Documents\Tech\Atom prj\MyMedia-FillDB>venv\Scripts\activate.bat
(venv) C:\Users\Jo\Documents\Tech\Atom prj\MyMedia-FillDB>_
```

Executer mediasEntities.py

```
(venv) C:\Users\Jo\Documents\Tech\Atom_prj\MyMedia-FillDB>python mediasEntities.py
* Serving Flask app "mediasEntities" (lazy loading)
* Environment: production
    WARNING: This is a development server. Do not use it in a production deployment.
    Use a production WSGI server instead.
* Debug mode: off
* Running on http://127.0.0.1:8080/ (Press CTRL+C to quit)
```

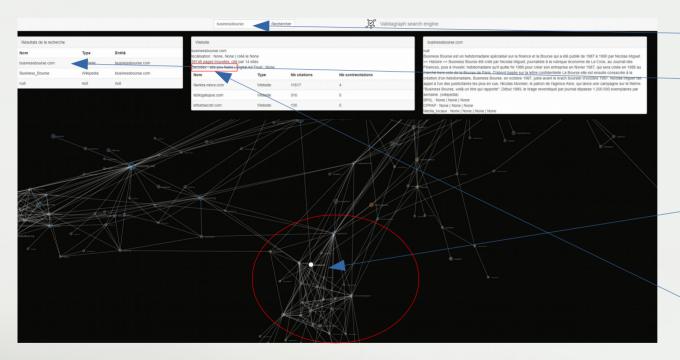


Dans chrome, aller sur http://127.0.0.1:8080/



- Barre de recherche
- Graph interactif (drag/ drop, click, zoom avec roulette...)
- Cliquer sur les résultats de la recherche permet d'afficher les infos sur les autres paneaux

Exemple



- Chercher « businessbourse »
- cliquer sur le résultat de recherche
 « businessbourse.com »
- =>
- positionne le nœud dans la branche liée à la complosphere
- affiche l'info suivantes : site peu fiable pour Decodex

Pour visualiser d'autres graphs

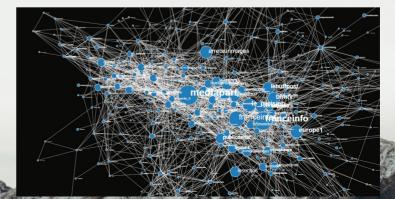
- Dans static/index_medias_entities.html :
- Chercher la ligne 240

```
var catColScale = d3.scaleOrdinal(d3.schemeCategory20);
var contColLinkScale = d3.scaleLinear()
    .domain([2, 100])
    .range(["white", "red"]);
var contLinkScale = d3.scaleLinear()
        .domain([2, 10000])
d3.json("/web_graph", function(error, graph){
var baseNodes = graph.nodes
var baseLinks = graph.links
console.log(graph.nodes)
console.log(graph.links)
var nodes = [...baseNodes]
var links = [...baseLinks]
```

Remplacer « /web_graph » par :

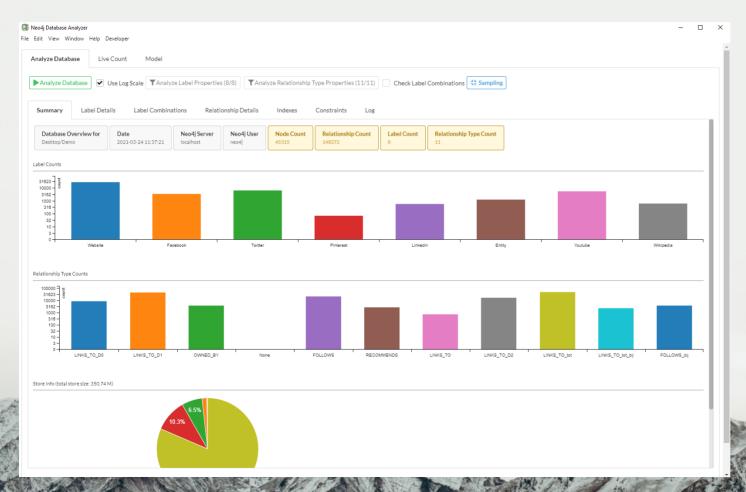
- « ent_graph » pour visualiser les entités
- « twit_graph » pour les comptes twitter
- « yt_graph » pour les comptes youtube

```
239 | 240 d3.json("/twit_graph", function(error, graph){
241
```



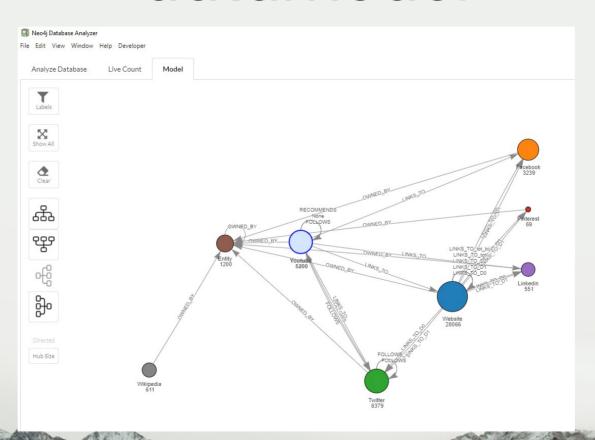
2. La base Neo4j

Statistiques globales (database analyzer)

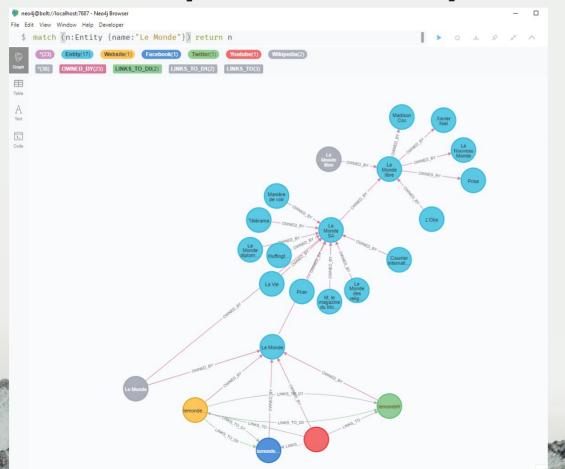


- 45 315 nœuds / 8 labels
- 148 272 relations / 11 types

datamodel



Exemple: Entity « Le Monde »



Dans la base, l'entité Le Monde possède

- Le site web lemonde.fr
- Le site web links les comptes :
 - Youtube : lemonde
 - Twitter: lemondefr
 - Facebook : lemonde.fr
- La page wikipedia Le Monde

L'entité Le Monde appartient à

- Le Monde SA
- Qui appartient à Le Monde libre
- Qui appartient à Xavier Nieal, Madison Cox, Prisa, Le Nouveau Monde...

Exemple: le website lemonde.fr

count(m) Links_to 304 nœuds match (n:Website {name:"lemonde.fr"})-[r:LINKS_TO_tot]→(m) return count(m) 304 Est linké par 270 nœuds count(m) match (n:Website {name:"lemonde.fr"}) ← [r:LINKS_TO_tot]-(m) return count(m) 270 Dont 44 de façon bijective count(m) match (n:Website {name:"lemonde.fr"}) ← [r:LINKS_TO_tot_bij]-(m) return count(m) liste 44 $\label{eq:match} \textbf{match (n:Website } \{name: \texttt{"lemonde.fr"}\}) \leftarrow \texttt{[r:LINKS_T0_tot_bij]-(m) } \ \textbf{return (m.name)}$ "leparisien.fr" "conspiracywatch.info" "lavie.fr" "lesinrocks com"

"huffingtonpost.fr"

"dreuz info"

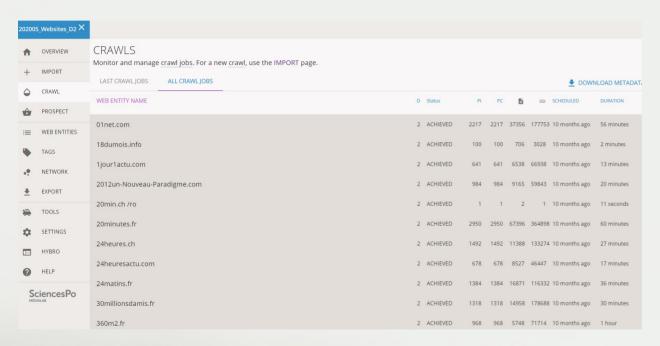
3. Le code pour remplir la base

Résumé : pour créer une database clean

For clean fresh new database : 1. run all clean Import DB. by (carefull, one procedure by one, connection bb with by 2neo.); tested with D0 and D1 DISCO, missing some websites compared to MyMedia-FillDB\data\202007WebsitesRS D0 2. run sandboximportRSformSiteList.pv, all cells starting with A. This imports file firstPageRS df_edit.csv 3 Twitter in twitter by o if not done, download twitter profiles with cell imports all Twitter profile files (with cell) o use twecoll for follow network (using python twecoll init -g NAME (where NAME is NAME.twt file) + python twecoll fetch NAME) o import twecoll dat file with import twecoll dat Youtube: in youtube.pv · download yt profiles and subscriptions from cell push profile info in db import youtube recommendations from profile o download profile / subscriptions for newly added youtube nodes (from previous step) o again: push profile in db + import youtube reco from profile · download links from youtube import links in db import subscriptions info wikipedia : wikipedia.py search for wikipedia page corresponding to Entities + create wiki search.csv Manual QC / Edit of the csv o import edited file as wikinodes and owned by connection Download wikipages as files o import info in db : infobox, categories, wikidata ref Import wikipedia summary 6 Whois o download whois files TODO: parse and import! CPPAP, Decodex, SPIIL, Mediaslocaux, DigitalAdTrust; classements.pv o follow script to import as Website properties

+ computeGDS.py pour des calculs de data science sur graph (pagerank, centrality, community detection, clustering...)

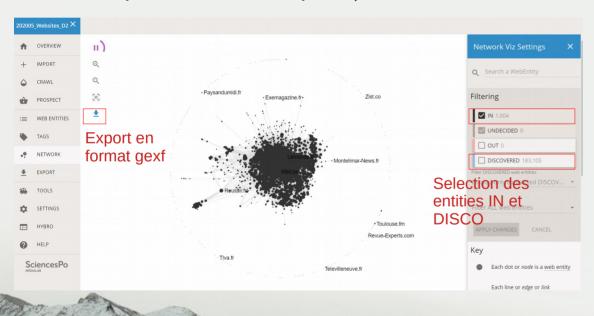
Le links crawler: hyphe

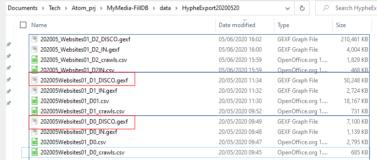


- De sciencepo medialab : site, github et demo
- Plateforme installé via docker en local, interface intuitive et robustesse des crawls
- Crawle une liste d'url et récupère tous les liens hypertextes.
 Possibilité de choisir la profondeur (0 pour les url, 1 pour les liens dans les url, 2 pour les liens des liens...)
- Itération pour une liste de 1000 sites de medias francophones

Le links crawler: hyphe

Exports de hyphe (dans network, selection des entities IN et DISCO pour projet crawl depth 0 et crawl depth 1)





Depth 2

Depth 1

Depth 0

Les fichiers importés dans la base : D0_DISCO.gexf et D1_DISCO.gexf.

D2_DISCO mériterait d'être importé mais le code non optimisé met trop de temps à l'import...

Import des données hyphe, monde diplo et ACPM

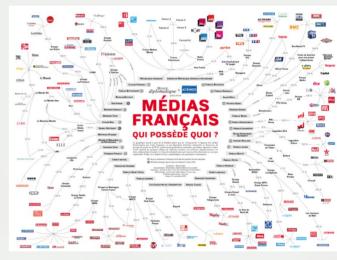
- 1. Le code utilisé est le début de data/cleanImportDB.py
- À partir des exports de Hyphe (D0_DISCO.gexf et D1_DISCO.gexf), création des nœuds (Website) et des liens LINKS_TO_D0 (depth 0) et LINKS_TO_D1 (depth D1)
- Import des données mondediplo disponible ici https://github.com/mdiplo/Medias_francais : création des entités et des relations OWNED_BY
- Import des données ACPM : statistiques des grands sites d'information française disponible ici :

https://www.acpm.fr/Les-chiffres/Frequentation-internet/Sites-Grand-Public/Classement-unifie

et ici:

https://www.acpm.fr/Les-chiffres/Frequentation-internet/Sites-Pro/Class ement-unifie

importés comme propriété dans les noeuds



Import des nœuds réseaux sociaux

site	fb	twit	yt_user	yt url	sc inst	ta dm site_name	id twit_hyphe	fb_hyphe	fb_final	twit_final	pi_hyphe
0 http://agri71.fr	agri71	agri71			70. 700	Agri71 fr	34 agri71	agri71	agri71	agri71	
1 http://axonais.fr	laxonaisofficiel					Axonais.fr	73	laxonaisofficiel	laxonaisofficiel		
2 http://bretons.bzh	bretons	magazinebretons				Bretons.bzh	98 magazinebretons	7000000000	bretons	magazinebretons	
3 http://c4magazine.org						C4magazine.org	111			7002000000	
4 http://cameditsport.com	cameditsport	lev egi1	cameditsport.com	https://www.youtube.com/channel/UCrCQO-y5BgQgH5xVqo47LxQ		Cameditsport.com	117		cameditsport	lev egi1	
5 http://chouard.org/blog/	700000000					Chouard.org /blog	137	juan.branco.98	700000000	700000	
6 http://courrier-français.com						Courrier-Français.com	165				
7 http://cgfg-journal.org	cgfd-mensuel-167147	cgfdjournal				Cgfd-Journal.org	177 cqfdjournal	cgfd-mensuel-1671478	cgfd-mensuel-16714783	26 cgfdjournal	
8 http://creuse-agricole.com	creuseagricole	creuseagricole				Creuse-Agricole com	186 creuseagricole	creuseagricole	creuseagricole	creuseagricole	
9 http://darons.net		parantsmagazine				Darons.net	188 parantsmagazine	/00000	7000000000	700000000000000000000000000000000000000	
10 http://diktacratie.com	diktacratiecom	diktacratie				Diktacratie.com	187		diktacratiecom	diktacratie	
11 http://entraid.com	entraid	journal entraid	entraid.com	https://www.youtube.com/channel/UCtk948ORqQ0oeeEWtqNEAuQ)	Entraid.com	219 journal entraid	entraid	entraid	journal entraid	
12 http://factuel.info	factuelinfo	factuelinfo				Factuel.info	237 factuelinfo	~~~	factuelinfo	factuelinfo	
13 http://haute-loire-paysanne.fr	,0000000	,000000				Haute-Loire-Paysanne.fr	305		,0000000	,0000000	
14 http://ilfattoquotidiano.fr						Ilfattoquotidiano.fr	338 estoyche			estoyche	
15 http://information-en-direct-france.com						Information-En-Direct-France.com	352	flashinfogiletsjaunes		~~~~	mickaelmick
16 http://kezako.unisciel.fr	unisciel	unisciel	unisciel	http://www.youtube.com/unisciel		kezako.Unisciel.fr	998 unisciel	unisciel	unisciel	unisciel	700000000
17 http://lagauchematuer.fr	la-gauche-ma-tuer	lagauchematuer	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,			Lagauchematuer.fr	416 lagauchematuer	700000	la-gauche-ma-tuer	lagauchematuer	
18 http://lalterego.fr				https://www.youtube.com/channel/UCUXg6UxOBQYeilcSqZvg4wQ		Lalterego.fr	430 lalteregofr	lalteregofr	lalteregofr	lalteregofr	
19 http://larenaissanceduloiretcher.fr	la-renaissance-du-loir-et-cher-111130649012373			700000700000000000000000000000000000000		Larenaissanceduloiretcher.fr	452	la-renaissance-du-loir-e	la-renaissance-du-loir-e	-cher-1111306490123	373
20 http://larevuedesmedias.ina.fr	larevuedesmedias	ina revuemedias		https://www.youtube.com/channel/UCD9xpWa_mjbbH0sskVjZxhg		larevuedesmedias.lna.fr	340 ina revuemedias	larev uedes medias	larev uedesmedias	ina revuemedias	
21 http://latelelibre.fr						Latelelibre.fr	470				
22 http://laterredecheznous.com						Laterredecheznous.com	471				
23 http://le-crestois.fr						Le-Crestois,fr	487	journal.lecrestois	journal.lecrestois		
24 http://leblogdenestor.com	leblogdenestor	marion_nestor		https://www.youtube.com/channel/UCj0HNkgz8eikcMM81NA4CEA		Leblogdenestor.com	497 leblogdenestor	leblogdenestor	leblogdenestor	leblogdenestor	
25 http://lechodelaboucle.fr	lécho-de-la-boucle	echodelaboucle				Lechodelaboucle.fr	503 echodelaboucle		lécho-de-la-boucle	echodelaboucle	
26 http://ledemocratevernonnais.fr	democratev ernonnais	Idmocrate				Ledemocratevernonnais.fr	514		democratev ernonnais	Idmocrate	
27 http://ledr.fr	ledr.fr	enfantsdurhone				Ledr.fr	518	525389480826266 339	∳ ledr.fr	enfantsdurhone	

2. code data/SanboximportRSfromSlteList.py

Récupère les noms de comptes réseaux sociaux crawlés par hyphe

Créé un fichier pour QC et edition

Importe les comptes dans Neo4j et créés liens

Import des données twitter

3. code data/twitter.py

Connection à l'API twitter (créer un compte dev et utiliser vos api keys pour le moment : remplacer les xxxxx en début de code)

Si non dispo, télécharge les infos de profile les noeuds twitter de la base comme des propriétés

Importe dans Neo4j

Puis utiliser twecol (code ici : https://github.com/jdevoo/twecoll) pour télécharger les follow network

Le code importe les dat files de twecol dans Neo4j (relationship FOLLOWS)

Import des données youtube

Code data/youtube.py

Connection à l'API youtube (se créer un compte dev sur google et utilisez vos propres clefs quand demandées)

Pour les comptes youtube, télécharge les infos de profil et subscriptions et les importe dans Neo4j

Import des données wikipedia

- Code data/wikipedia.py
- Cherche et importe les infos wikipedia sur les entités de la base Neo4j
- TODO: importer touts les infos wikidata utiles

Autres données

- Whois: download les infos whois dans des fichiers. Non importé dans Neo4j (TODO: a parser) → adresse des serveurs, personnes responsables...
- Classemen.py : import des groupes de crédibilité

