



cikisi
Search, Watch, Explore

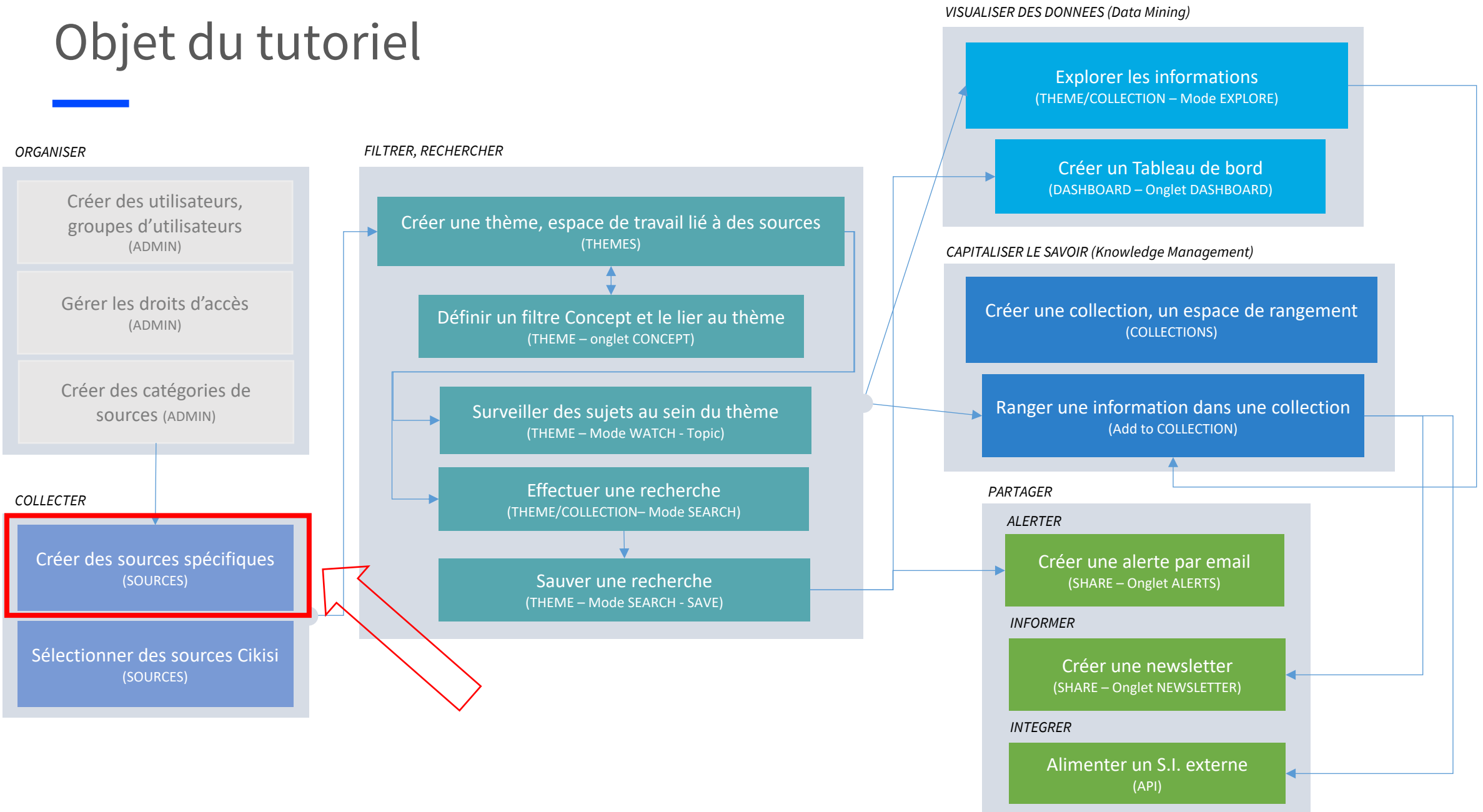


Sources – Robot de collecte partant d'un flux (geo)RSS

V1.0

March 2021

Objet du tutoriel



Plan

- **Définition** d'un flux (geo)RSS
- **Cas d'emploi**
- Vérification de la **validité**
- **Avantages**
- **Limitations**
- **Principes techniques**
- **Fréquence** de collecte
- **Trouver** le(s) flux RSS d'un site web

Définition d'un flux (geo)RSS

- Un flux RSS est un **fichier XML**, composé de balises encadrant du contenu.
- Les **balises** indiquent le type de contenu :
 - Titre
 - Description
 - Auteur
 - Date de publication
 - Lien
 - Lien de l'image associée
- Toutes les balises ne sont pas toujours présentes
- Les flux **geoRSS** sont ceux possédant également une balise reprenant des **coordonnées géographiques** (lat/long)
- Les symboles du flux RSS les plus fréquents sont les suivants :

XML = contenu structuré

```
<item>
  <title>
  ...
</title>
  <description>
  ...
</description>
  <guid isPermaLink="true">0200700c103de4ab1832bc960e5c6e879</guid>
  <pubDate>Thu, 11 Mar 2021 00:00:00 +0000</pubDate>
  <link>
  ...
</link>
</item>
```

balise

contenu



Cas d'emploi

- Le type de source « Flux geo(RSS) » **est utilisé lorsque**:
 - Un ou des **flux RSS sont disponibles** sur le site à mettre en surveillance
 - Les **flux RSS sont bien alimentés**, de façon régulière et exhaustive par le site
 - Les articles listés au sein du flux RSS couvrent une période au minimum égale à la fréquence de collecte
 - Ne pas utiliser un robot de collecte partant d'un flux RSS avec une fréquence de 12 heures si les articles listés au sein du flux RSS sont relatifs aux 2 dernières heures. Sans quoi, la collecte sera partielle (1/6^{ème})
- Le type de source « Flux geo(RSS) » **ne doit pas être utilisé lorsque** le flux RSS provient d'une alerte **Google**, d'un site **Reddit** ou **Blogspot** (utiliser le type de source Google Alert, Reddit ou Blogspot)

Vérification de la validité d'un flux (geo)RSS

Un flux RSS est valide si :

- Il contient plusieurs articles
- La date de publication des articles est récente ou correspond à la dernière publication du site

Flux RSS pas valide

```
<?xml version="1.0" encoding="utf-8"?>
<rss version="2.0" xml:base="https://www.navy.gov.au"
xmlns:dc="http://purl.org/dc/elements/1.1/">
<channel>
<title>Royal Australian Navy</title>
<link>https://www.navy.gov.au</link>
<description></description>
<language>en</language>
</channel>
</rss>
```

```
<?xml version="1.0" encoding="utf-8"?>
<rss version="2.0" xml:base="https://www.navy.gov.au"
xmlns:dc="http://purl.org/dc/elements/1.1/">
<channel>
<title>Royal Australian Navy</title>
<link>https://www.navy.gov.au</link>
<description></description>
<language>en</language>
</channel>
</rss>
```

Avantages

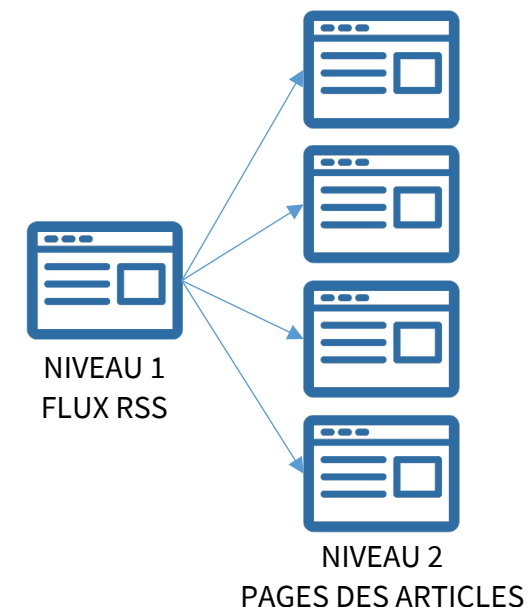
- **Configuration extrêmement simple**
 - Seule l'adresse du flux RSS est nécessaire
- En cas de changement de structure de site, **pas besoin de maintenance** (contrairement aux « scraping bots ») tant que le flux RSS existe toujours
- Permet de **collecter l'entièreté de l'article** et pas uniquement la description disponible dans le flux RSS
- **Fréquence de collecte élevée**
 - Jusqu'à 60 min
 - = en moyenne un article est collecté 30 min après sa publication

Limitations

- Parfois **le(s) flux RSS du site ne liste(nt) pas toutes les nouvelles pages créées**
- Si le **site web n'utilise pas des standards HTML** (de plus en plus rare), risque de ne pas reconnaître le contenu principal de la page et de limiter la collecte à la description
- Dans de rares cas, la **date de publication** de l'article n'est pas présente ou pas déclarée correctement dans le flux RSS et sera donc manquante.
- L'**adresse du flux RSS peut changer** au bout de quelques années il faut donc la mettre à jour (suivre le statut des sources)

Principes techniques

- Le robot Cikisi **extraie du contenu sur 2 niveaux**:
 - Le contenu présent au sein du **fichier XML** (flux RSS)
 - Contenu structuré
 - Le contenu présent sur **les liens (URLs) listés** au sein du flux RSS
 - Extraction automatique du contenu en interprétant les balises HTML
- **Pas d'écrasement**
 - Cikisi conserve la version la plus ancienne des pages (URL)
 - Une page déjà visitée ayant mené à la création d'un article n'est pas mise à jour
 - A noter qu'à l'exception des modifications mineures du contenu d'un article (orthographe, faute de frappe), les autres mises à jour font souvent l'objet d'une modification de l'URL et seront donc collectées par Cikisi.



Fréquence de collecte

- Choisissez la fréquence de collecte en adéquation avec vos besoins en veille mais aussi en fonction de la quantité d'articles publiés par heure/jour par le site. En effet, le fichier XML contient uniquement les 10, 20 ou 50 (rarement) derniers articles publiés par le site. Visiter ce fichier pas assez fréquemment pourrait mener à rater des nouveaux articles.
- La fréquence de 60 minutes est idéale

Frequency:

6 Hours

12 Hours

24 Hours

7Days



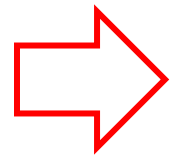
Trouver le flux RSS d'un site

■ Lorsque le site déclare son flux RSS

- Chercher « RSS » au sein de la page à l'aide de CTRL-F
- Identifier visuellement l'icône RSS quelque part sur la page
- Regarder dans le pied de page sous le nom « RSS », « ATOM » ou « abonnez-vous »
- Rendez-vous dans la section « news » / « actualités » où se situe fréquemment le flux

■ Lorsque le site possède un flux RSS mais ne le déclare pas

- Ajouter après le nom de domaine, un des chemins suivants :



- | | | | |
|-----------------------|------------------------|------------------|--------------------------|
| ■ /feed | ■ /xml/syndication.rss | ■ en/rss.xml | ■ /rss-feed |
| ■ /rss | ■ /feed/rss | ■ /index.rss | ■ /rss/all |
| ■ /rss.xml | ■ fr/rss.xml | ■ /rss/feed | ■ /news/rss.xml |
| ■ /blog/feed | ■ /rss.jsp | ■ /en/rss | ■ /rss/articles |
| ■ /fr/feed | ■ /rss/news | ■ /feeds/news | ■ /feeds/all |
| ■ news/rss | ■ /en/feed | ■ /rss/news.xml | ■ /feeds/rss |
| ■ news/feed | ■ /rss/home | ■ /blog/rss | ■ /rss.business.xml |
| ■ rss.aspx | ■ /fr/rss | ■ /rss.html | ■ /rss/news-releases.XML |
| ■ feeds/posts/default | ■ /rss/2.0 | ■ blog/rss.xml | ■ /feeds/rss.php |
| ■ rss.php | ■ /category/news/feed | ■ /comments/feed | ■ ... |

Questions?

Utilisez notre formulaire de support!

<https://wmt.cikisi.com/support>

Account Settings

Support form

Maintenance

Log Out

+ 🌙 🌞

Company name

Prenom / Firstname

Nom / Name

Email*

Sujet du Ticket / Ticket name*

Description du ticket / Ticket description*

Comment reproduire le probleme? / How to reproduce the issue*

Veuillez copier l'URL correspondant au probleme rencontre et donner les details essentiels afin de le reproduire

Attacher un fichier ou Capture d'ecran / Attached File or Screenshot
 Aucun fichier choisi



ATTENTION

Pour une compatibilité parfaite et une meilleure expérience utilisateur, Cikisi recommande l'emploi d'une version récente du navigateur Google Chrome