



Le corpus des Ouvriers des deux mondes : des images et des URLs

Jean-Damien G      

► To cite this version:

Jean-Damien G      . Le corpus des Ouvriers des deux mondes : des images et des URLs. 2020. hal-03118736

HAL Id: hal-03118736

<https://hal.science/hal-03118736>

Submitted on 22 Jan 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destin     au d    t et    la diffusion de documents scientifiques de niveau recherche, publi     ou non,   manant des   tablissements d'enseignement et de recherche fran    ais ou   trangers, des laboratoires publics ou priv    s.



Distributed under a Creative Commons Attribution 4.0 International License

Le corpus des *Ouvriers des deux mondes* : des images et des URLs

Jean-Damien G  n  ro*

19 juin 2020

Si les documents d'archives ont une part pr  pond  rante dans le projet Time us, ils ne repr  sentent pas pour autant l'int  gralit   de sa documentation. Les imprim  s sont   galement pr  sents, sous la forme de trois importants dossiers : la collection de la presse ancienne lyonnaise, divers imprim  s portant sur le textile en France au XIX   si  cle, et le corpus des *Ouvriers des deux mondes*¹.

Les *Ouvriers des deux mondes* sont des enqu  tes sociologiques r  parties en 3 s  ries et 126 monographies². Initi  e par le sociologue Fr  d  ric Le Play (1806-1882), la publication est assur  e par la Soci  t   internationale des   tudes pratiques d'  conomie sociale de 1857    1928 et repr  sente un total de 13 volumes. Ceux-ci sont aujourd'hui int  gralement consultables sur le site *Internet Archive*³. Nous allons nous int  resser dans ce billet aux fichiers de transcription de ces volumes et au lien entre ceux-ci et les images num  ris  es d'origine.

Le script LSE OD2M,   crit par Alix Chagu  , avait automatiquement segment   et transcrit les images, puis encod   et structur   en XML-TEI les textes bruts ainsi obtenus⁴ ; la sortie avait r  sult   en 13 fichiers XML. Ces fichiers « sources » avaient ensuite   t   scind  s en 222 fichiers XML correspondant    autant de divisions logiques des volumes : les monographies bien s  r, mais   galement les introductions, tables des mati  res et autres   l  ments de paratexte. Des op  rations de v  rification ont permis de r  duire le nombre de fichiers    192.

*Ing  nieur de recherche et d  veloppement stagiaire, Inria,   quipe ALMA  CH ;   tudiant du Master "Technologies num  riques appliqu  es    l'histoire" de l'  cole nationale des chartes. Lien vers le billet original.

1. Aper  u des   tats sur le wiki Time Us.

2. Anthony, Lorry, « Les monographies des *Ouvriers europ  ens* (1855, 1877-1879) et des *Ouvriers des deux mondes* (1857-1930). Inventaire et classification », dans *Les   tudes sociales. Les monographies de famille de l'  cole de Le Play*, n  s 131-132, 2000, pp. 93-181, spec. p. 101.

3. Liens des diff  rents volumes sur le wiki Time Us.

4. Voir son billet    propos de la Constitution d'un corpus textuel sur les monographies de Le Play.

1 Les images : stockage local ou distant ?

Le schéma d'encodage retenu conserve le lien entre l'image et sa transcription. Celui-ci s'exprime sous la forme d'un élément `<facsimile>` englobant un ensemble de balises `<graphic>`, dont l'attribut `@url` indique la localisation de l'image.

```
<facsimile xml:id="facs_451">
  <surface lrx="2721" lry="4415">
    <figure />
    <zone rendition="printspace">
      <zone lrx="1357" lry="712" rendition="paragraph"
        ulx="352" uly="474" xml:id="facs_451_p_1"/>
      <zone lrx="1354" lry="864" rendition="paragraph"
        ulx="350" uly="750" xml:id="facs_451_p_2"/>
      <zone lrx="2430" lry="3952" rendition="paragraph"
        ulx="1419" uly="3772" xml:id="facs_451_p_24"/>
    </zone>
  </surface>
  <graphic url="../images/bin/
    lesouvriersdesde01sociuoft_0454.tif"/>
</facsimile>
```

Le script LSE OD2M a travaillé à partir d'images stockées localement après avoir été téléchargées depuis *Internet Archive* : chaque attribut `@url` contient donc un chemin local vers une image.

Ce stockage local répondait à un besoin spécifique lors de la phase de transcription. Il pose néanmoins problème pour la suite du traitement, dans la mesure où il ne garantit pas la portabilité du corpus.

Le site *Internet Archive* met à disposition du plus grand nombre des ressources digitalisées ou numériquement natives depuis son lancement en 1996. Son idée fondatrice est d'être un centre stable et durable d'archives digitales ; stocker en local les ressources qui en sont issues ne semble ainsi pas nécessaire. Il a donc très vite été question de substituer au chemin local l'url de l'image sur *Internet Archive*.

Pour ce faire, deux étapes étaient nécessaires :

- Rechercher les urls ;
- Écrire un script pour :
 1. Itérer sur l'ensemble des fichiers XML ;
 2. Comparer les chemins locaux aux urls ;
 3. Effectuer la substitution lorsque les deux correspondaient à la même image.

2 Recherche des urls

Effectuer la substitution sur la base d'une expression régulière n'était pas envisageable, car la dénomination des images dans les urls sur *Internet Archive* et dans les des fichiers images n'était pas similaire.

La piste de la librairie python `internetarchive` a été explorée, mais là encore sans succès. C'est finalement l'analyse du code source des pages d'*Internet Archive* qui a permis de remonter jusqu'à un fichier JSON contenant les urls des images.

Le format JSON présente une information structurée qui s'apparente à des dictionnaires et des listes pour le langage Python. Celui-ci intègre un module (`json`) permettant de lire les données des fichiers JSON. Dans le cas qui nous intéresse, les url se trouvaient à ce chemin :

```
|----['data ']  
    |----['brOptions ']  
        |----['data ']  
            |----[index de la double page]  
                |----[index de la page]  
                    |----['uri ']
```

Une fois cette information connue, il devenait possible de remplacer dans le fichier XML les chemins locaux par ces urls. Encore fallait-il s'assurer que les deux correspondaient à la même image.

3 Le script de substitution

Le script de comparaison des chemins et des urls (fig. 1) requiert trois arguments :

- Un fichier `.csv` contenant une liste des liens vers les fichiers JSON ;
- Un deuxième fichier `.csv` avec une liste de fichiers `.xml` ;
- Le chemin local menant au dossier de ces fichiers.

Dans un premier temps, le script itère sur les urls des JSON et sur les dénominations des fichiers XML. Dans un second temps et à condition que l'identifiant du JSON se trouve dans l'intitulé du fichier XML — et donc si les deux ont bien pour objet le même volume — une requête GET est effectuée vers le JSON via le module `requests`⁵ et le fichier XML est ouvert et parsé via la librairie `Beautiful soup`⁶.

5. Méthode du protocole HTTP permettant de demander une ressource stockée sur un serveur(documentation).

6. Parser un fichier consiste à le lire et à interpréter son contenu afin d'en extraire certains éléments.

Ensuite, une nouvelle itération est effectuée sur chaque double-page dans le JSON et sur chaque balise `<graphic>` dans le fichier XML. Les valeurs obtenues (url de la page de droite, url de la page de gauche, chemin local de l'attribut `@url` de `<graphic>`) sont stockées et comparées : lorsqu'une correspondance est trouvée, `@url` prend pour nouvelle valeur l'adresse de l'image sur *Internet Archive*.

La comparaison finale est effectuée grâce à une expression régulière fondée sur l'identifiant unique à quatre digits de l'image, précédé par un tiret bas (`_d{4}`), présent tant dans le chemin local que dans l'url de l'image.

Conclusion

Le script a effectué 6503 insertions dans 192 fichiers en un peu moins d'une minute.

Il a été écrit pour répondre à un besoin posé par les fichiers XML d'un corpus spécifique ; sa réutilisation est néanmoins envisageable si l'utilisateur peut satisfaire aux arguments du script. Sur les trois requis, les deux premiers sont essentiels pour son bon fonctionnement (deux tableaux listant les JSON et les fichiers XML qui seront comparés). Le troisième, le chemin absolu vers le dossier contenant les XML, peut en revanche être rendu optionnel et converti en une variable intégrée dans l'exécution du script.

En dernier lieu, il faut noter que les urls listées dans les JSON présentent l'avantage de pointer directement vers les images sources, et non vers l'interface de consultation d'Internet Archive. C'est une garantie de pérennité (aucune dépendance vis à vis des mises à jour de la visionneuse). Il devient possible d'utiliser ces images dans une édition en ligne des *Ouvriers des deux mondes*, tout en évitant le coût d'hébergement des images.

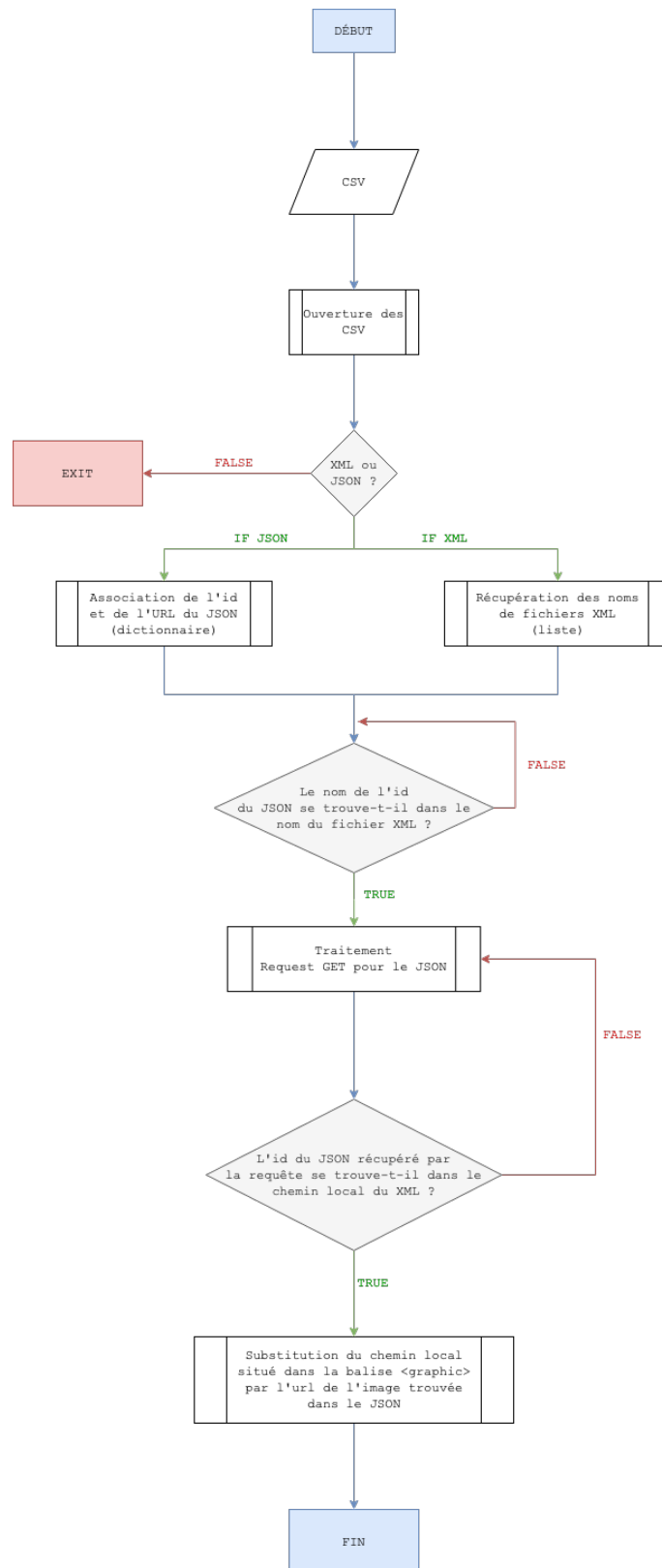


FIGURE 1 – Algorithme du script