

ÉCOLE NATIONALE DES CHARTES

Maxime Challon

licencié ès histoire

Les référentiels en institutions patrimoniales : évolution des pratiques et repositionnement

L'exemple des référentiels de l'Institut National de l'Audiovisuel

Mémoire pour le diplôme de master

« Technologies numériques appliquées à l'histoire »

2020

Résumé

Ce mémoire, réalisé pour l'obtention du diplôme de Master 2 « Technologies numériques appliquées à l'histoire » de l'École nationale des Chartes, retrace l'évolution des pratiques documentaires sur les référentiels en institution patrimoniale à travers l'étude des référentiels de l'Institut national de l'Audiovisuel (INA) et leurs alignements. Cette étude de l'évolution des formes et des structures des référentiels est liée à l'évolution de la place de ces référentiels au sein des systèmes documentaires, ainsi qu'aux besoins qui leur sont liés.

Mots-clés : institut national de l'audiovisuel ; référentiel ; thésaurus ; vocabulaire contrôlé ; vocabulaire hiérarchique ; ontologie ; web de données ; Wikidata ; liens ; alignement.

Informations bibliographiques : Maxime Challon, *Les référentiels en institutions patrimoniales : évolution des pratiques et repositionnement. L'exemple des référentiels de l'Institut National de l'Audiovisuel.*, mémoire de master « Technologies numériques appliquées à l'histoire », dir. Gautier Poupeau, École nationale des chartes, 2020.

Remerciements

MES remerciements vont tout d'abord à Gautier POUPEAU, mon maître de stage, qui m'a accueilli, guidé, conseillé et intégré à son équipe malgré le travail à distance imposé par le contexte actuel. Je souhaite également remercier Axel ROCHE-DIORÉ pour ses explications et son soutien dans la réalisation technique de mon stage.

J'adresse aussi mes remerciements aux membres du pôle « Ingénierie de la Donnée », Lauryne LEMOSQUET, Otmane ELABBOUBI et Akli ABDI, ainsi qu'à Florence BRÉANT, cheffe de projet pour le *Lac de données*, pour le temps qu'ils m'ont accordé.

Que soit également remercié l'ensemble du département « Architecture et Innovation » de l'INA pour l'accompagnement fourni tout au long de mon stage, notamment Stanislas DE MAIGRET et Matthieu BORICAUD pour le déploiement de l'application, et Olivio SÉGURA pour la présentation des archives de l'INA.

Liste des abréviations

DDCOL Direction déléguée aux collections

DJ Direction juridique

DSI Direction des systèmes d'information

ÉPIC Établissement public à caractère industriel et commercial

INA Institut national de l'Audiovisuel

ISAN *International Standard Audiovisual Number*

ORTF Office de la radio-télévision française

Introduction

« Toutefois pour ne laisser cette quantité infinie ne la définissant point, [et] aussi pour ne jeter les curieux hors d'espérance et pouvoir acco[m]plir [et] venir à bout de cette belle entreprise, il me semble qu'il est à propos de faire comme les Médecins, qui ordonnent la quantité des drogues suivant la qualité d'icelles, [et] de dire que l'on ne peut manquer de recueillir tous ceux qui auront les qualitez [et] conditions requises pour estre mis dans une Bibliotheque.¹ »

EN 1627, Gabriel NAUDÉ compare le médecin au bibliothécaire, semblables par leur nécessité d'ordonner pour sélectionner, de classer pour retrouver, au milieu d'une masse d'objets. Cet ordonnancement, ce classement, passent pas une hiérarchisation de leur connaissance ou de leurs outils, dans le but de faciliter la recherche d'un médicament ou d'un livre pour l'utilisateur final. Cependant, plusieurs siècles plus tard, la hiérarchisation de la connaissance, ayant pour but de référencer une instance de la vie réelle, ne fonctionne plus : l'utilisateur ne part plus que très rarement d'un terme de la hiérarchie pour trouver son document ; il utilise le plus souvent un mot ou un concept qui le renverront vers une liste de résultats correspondant à sa requête. Alors, la notion de graphe prend le dessus sur celle de hiérarchie.

La notion évoquée de « quantité infinie » est aujourd'hui d'autant plus valable avec le web et l'explosion des quantités de données produites et stockées : avec cette mort de la notion de ressource, et par conséquent de celle de référentiel, la donnée structurée est implantée, peut être exploitée à la fois par une machine et par une personne, et est divisible et modulable à l'infini.

Cette transition de la ressource à la donnée, des référentiels hiérarchiques aux référentiels en graphe, est observable à l'INA. Créé en 1975 suite au démantèlement en sept sociétés de l'Office de la radio-télévision française (ORTF) par la loi du 07 août 1974,

1. Gabriel Naudé, *Advis pour dresser une bibliotheque. Tome 1 / . Présenté à monseigneur le president de Mesme. Par G. Naudé P....* T. 1, Chez François Targa, au premier pillier de la grand'salle du Palais, devant les Consultations, Paris, 1627, URL : <https://gallica.bnf.fr/ark:/12148/bpt6k1041429f> (visité le 01/09/2020), p.41-42.

l'INA est désigné comme un Établissement public à caractère industriel et commercial (ÉPIC) et « chargé de la conservation des archives, des recherches de créations audiovisuelles et de la formation professionnelle »². À ces missions est ajouté à partir de 1992 le dépôt légal de la télévision, de la radio, de la télévision satellite, par câble et numérique. Cette massification continue de documents et de données nécessite un classement et un référencement efficace des collections, ce qui a conduit à la création de plusieurs référentiels dans l'Institut.

Face à la croissance de l'utilisation du numérique, à l'accroissement des collections et des données à l'INA depuis les numérisations des collections au début des années 2000, aux nouveaux besoins exprimés par les professionnels et le public, une refonte du système documentaire est mise en place à la Direction des systèmes d'information (DSI) au sein du département « Architecture et Innovation » : les données et leurs métadonnées sont extraites des anciens silos de conservation, puis transformées et migrées dans un nouveau système d'information centralisé. Ainsi, les référentiels, descripteurs de chaque document, identificateurs de personnes ou d'instances des collections, subissent également ce traitement pour les uniformiser et permettre une homogénéisation et une meilleure valorisation des données de l'INA.

Cette migration massive permet d'observer l'évolution des pratiques documentaires de référencement et de description de ces dernières décennies, suivant la même évolution que l'ensemble du milieu bibliothéconomique en France, ainsi que les changements de structure des référentiels utilisés. La diversité de formes et de structures des référentiels montre que ces derniers sont considérés seulement comme des outils à disposition du documentaliste pour décrire ses fonds ; périphériques et éclatés, ils ne permettent pas une centralisation uniforme des données de l'INA.

Le projet du *Lac de données*, débuté en 2014, a pour but de centraliser l'ensemble des données de l'INA, les référentiels prenant alors une place centrale dans le nouveau système d'information. Ce projet s'inscrit dans l'évolution des besoins, tant chez les documentalistes que chez les utilisateurs, avec une utilisation désormais massive du web par tous les publics - chercheurs, professionnels des médias, jeunesse, ...- pour la recherche et la consultation de contenus. Cette éditorialisation croissante et indispensable nécessite de nombreuses données de référence, par lesquelles les contenus sont recherchables et trouvables.

2. Loi n°74-696 du 7 août 1974 RELATIVE A LA RADIODIFFUSION ET TELEVISION, 1974, URL : https://www.legifrance.gouv.fr/jo_pdf.do?id=JORFTEXT000000333539&pageCourante=08355 (visité le 01/09/2020), art.3.

Ce mémoire offre une réflexion sur ces évolutions des pratiques et des usages des référentiels à l'INA, et plus généralement dans une institution patrimoniale. Au-delà de ces évolutions sensibles, c'est le positionnement du référentiel au sein des systèmes documentaires qu'il est nécessaire d'interroger, de manière à faire face aux nouveaux enjeux et aux nouveaux besoins exprimés ces dernières années : d'un rôle périphérique, pensé comme un outil, le référentiel devient désormais un pivot autour duquel les données documentaires se raccrochent.

Mon stage, débuté en mai 2020 et terminé fin août 2020, à la DSI de l'INA, m'a permis d'intégrer le département « Architecture et Innovation » de Gautier POUPEAU, et plus particulièrement le pôle « Ingénierie de la Donnée » dirigé par Axel ROCHE-DIORÉ, afin d'effectuer une réflexion sur les méthodes d'alignement de plusieurs référentiels, et de mettre en œuvre ces méthodes. Les échanges avec mes collègues du pôle « Ingénierie de la Donnée » et les professionnels de la documentation de la Direction déléguée aux collections (DDCOL) et de la Direction juridique (DJ) m'ont permis de naviguer dans les référentiels, d'observer leurs différences, leurs structures, de comprendre les besoins qui leurs étaient associés ainsi que les difficultés impliquées par chaque référentiel dans l'opération d'alignement en vue de leur migration vers le *Lac de données*. Plusieurs missions m'ont ainsi été confiées :

- Extraire les fonctions et les occupations de personnes physiques depuis les notes qualité en texte libre du référentiel des personnes physiques et morales de la DDCOL, puis aligner ces fonctions extraites avec un thésaurus de noms communs propre à la DDCOL
- Aligner les personnes physiques de la DDCOL avec les entités correspondantes de Wikidata
- Aligner les fictions et les séries conservées à l'INA avec Wikidata de manière à récupérer également l'identifiant *International Standard Audiovisual Number* (ISAN)
- Aligner les référentiels de personnes physiques de la DJ et de la DDCOL, puis développer une interface de vérification et de complétion des alignements réalisés automatiquement

Ce mémoire retrace l'évolution des usages et des pratiques documentaires concernant les référentiels dans les institutions patrimoniales, en s'appuyant sur l'exemple des référentiels de l'INA. Dans un premier temps, dans une période allant jusqu'au début des années 2000, les référentiels sont uniquement considérés comme des fournisseurs de clés entre les données de manière à les contrôler plus facilement. Puis, jusqu'au milieu des années 2010, le web et le web de données permettent une mise en commun des référentiels qui se retrouvent alors liés entre eux. Enfin, depuis le milieu des années 2010, les

référentiels sont placés au centre des systèmes d'information : ils sont devenus les pivots des systèmes documentaires.

Première partie

**CONTRÔLER. A la recherche de
clés (années 1960 – fin des années
1990)**

Chapitre 1

Le référentiel comme clé

- 1.1 Du langage libre au langage contrôlé : vers l'indexation
- 1.2 Une clé entre les données : une terminologie maîtrisée, objectif des vocabulaires contrôlés
- 1.3 Une clé entre les jeux de données : l'interopérabilité par les fichiers d'autorité et les portails

Chapitre 2

L'arbre, un vocabulaire contrôlé hiérarchique

Chapitre 3

Les référentiels à l'INA

Deuxième partie

**RELIER. Vers le partage de
référentiels communs (début des
années 2000 – milieu des années
2010)**

Chapitre 4

Le web de données : une exposition commune des référentiels

Chapitre 5

Partager des structurations
similaires de jeux de données par les
classes et les propriétés : les
ontologies, grammaires communes
mais spécifiques

Chapitre 6

Relier ses données à Wikidata

Troisième partie

CENTRALISER. Le référentiel, clé de voûte et pivot (depuis le milieu des années 2010)

Chapitre 7

Les labyrinthes comme réseaux de données et de liens

Chapitre 8

Le Lac de données de l'INA : le référentiel au centre du modèle

Chapitre 9

Centraliser les référentiels de l'INA
dans le Lac de données : l'exemple
de l'alignement de deux référentiels
de personnes physiques

Conclusion

Annexes

Table des figures

Table des matières

Résumé	iii
Remerciements	v
Liste des abréviations	vii
Introduction	ix
 I CONTRÔLER. A la recherche de clés (années 1960 – fin des années 1990)	 1
1 Le référentiel comme clé	3
1.1 Du langage libre au langage contrôlé : vers l’indexation	3
1.2 Une clé entre les données : une terminologie maîtrisée, objectif des vocabulaires contrôlés	3
1.3 Une clé entre les jeux de données : l’interopérabilité par les fichiers d’autorité et les portails	3
2 L’arbre, un vocabulaire contrôlé hiérarchique	5
3 Les référentiels à l’INA	7
 II RELIER. Vers le partage de référentiels communs (début des années 2000 – milieu des années 2010)	 9
4 Le web de données : une exposition commune des référentiels	11
5 Partager des structurations similaires de jeux de données par les classes et les propriétés : les ontologies, grammaires communes mais spécifiques	13
6 Relier ses données à Wikidata	15

III CENTRALISER. Le référentiel, clé de voûte et pivot (depuis le milieu des années 2010)	17
7 Les labyrinthes comme réseaux de données et de liens	19
8 Le Lac de données de l'INA : le référentiel au centre du modèle	21
9 Centraliser les référentiels de l'INA dans le Lac de données : l'exemple de l'alignement de deux référentiels de personnes physiques	23
Conclusion	25
Table des figures	29
Table des matières	31