
DMP du projet "E-NdP"

Plan de gestion de données créé à l'aide de DMP OPIDoR, basé sur le modèle "ANR - Modèle de PGD (français)" fourni par Agence nationale de la recherche (ANR).

Renseignements sur le plan

Titre du plan DMP du projet "E-NdP"

Langue fra

Date de création 2021-07-26

Date de dernière modification 2021-11-08

Identifiant

Renseignements sur le projet

Titre du projet E-NdP

Résumé Le projet e-NDP (Anr 2020) engage l'examen des registres capitulaires, de la documentation et des livres de Notre-Dame de Paris au Moyen Âge, pour étudier la société, l'économie, le bâti du cloître et de ses dépendances. Porté par le LaMOP, le projet a pour partenaires les Archives nationales, la Bibliothèque nationale de France (Département des Manuscrits, Bibliothèque de l'Arsenal), l'École nationale des Chartes et la Bibliothèque Mazarine.

Sources de financement

- Agence nationale de la recherche (ANR) : ANR-20-CE27-0012

Produits de recherche :

1. Numérisations des différents manuscrits (Image)
2. Transcriptions des différents manuscrits (Texte)
3. Modèles de transcription automatique de l'écriture manuscrite (Modèle IA)
4. Base de données comprenant les Entités Nommées (Jeu de données)
5. Données Cartographiques (Jeu de données)

Contributeurs

Nom	Affiliation	Rôles
Julie Claustre		<ul style="list-style-type: none">• Coordinateur du projet• Personne contact pour les données (Manuscrits, Modèles, Textes, Base de données, Cartes)
Pierre Brochard (https://orcid.org/0000-0003-1955-556X)		<ul style="list-style-type: none">• Responsable du plan

Droits d'auteur :

Le(s) créateur(s) de ce plan accepte(nt) que tout ou partie de texte de ce plan soit réutilisé et personnalisé si nécessaire pour un autre plan. Vous n'avez pas besoin de citer le(s) créateur(s) en tant que source. L'utilisation de toute partie de texte de ce plan n'implique pas que le(s) créateur(s) soutien(nen)t ou aient une quelconque relation avec votre projet ou votre soumission.

DMP du projet "E-NdP"

1. Description des données et collecte ou réutilisation de données existantes

Numérisations des différents manuscrits

Les campagnes de numérisation des manuscrits sont réalisées par les Archives Nationales, par la Bibliothèque Nationale de France et par la Bibliothèque Mazarine.

Les données produites seront

- des numérisations de 26 registres du Chapitre Notre-Dame de Paris sur la période de 1326 à 1504. Ces numérisations seront au format TIFF, au nombre de 9618 et d'une volumétrie totale de 400 Go. Le format TIFF choisi permet une numérisation de la meilleure qualité possible.
- des numérisations issues du fonds 'Notre-Dame', soit 30 manuscrits à la Bibliothèque nationale de France et 17 manuscrits à la Bibliothèque de l'Arsenal. Ces numérisations seront au format TIFF.

Transcriptions des différents manuscrits

Les textes produits sont des transcriptions des manuscrits à l'aide de deux plateformes : eScriptorium (<https://escriptorium.fr/>) projet de recherche soutenu par PSE et l'Inria et Transkribus (<https://transkribus.eu>), plateforme payante proposée par la coopérative READ-COOP.

La transcription avec ces deux outils s'effectuent en deux étapes : la première étape consiste à établir une transcription manuelle de plusieurs manuscrits. Dans un deuxième temps, à partir de ces vérités de terrain créées (données d'entraînement), les outils effectuent une transcription supervisée automatique.

L'utilisation de la plateforme eScriptorium sera gratuite pendant la durée de l'ANR.

Le coût d'une transcription effectuée par la plateforme Transkribus à la date du 18/10/2021 serait de 0.2 euros par page (<https://readcoop.eu/transkribus/credits/>)

Une transcription pour chaque page de manuscrit sera proposée dans plusieurs formats, un format Texte, un format TEI et un format PageXML.

Le format Texte permettra un usage simple de la transcription effectuée, le format TEI sera destiné à un usage plus élaboré avec l'inclusion de métadonnées décrivant le texte.

Enfin le format PageXML est destiné à un partage de données d'apprentissage pour les outils permettant la reconnaissance de l'écriture automatique. En effet, les données d'apprentissage ainsi créées pourront permettre l'enrichissement d'autres corpus d'apprentissage ou servir à la transcription automatique de textes à la graphie proche.

Modèles de transcription automatique de l'écriture manuscrite

Les modèles seront produits à partir des données d'apprentissage créées dans l'outil eScriptorium sur l'infrastructure CREMMA (Consortium pour la reconnaissance d'écriture manuscrite des matériaux anciens), infrastructure portée par l'École des Chartes et l'Inria.

Les modèles seront proposés dans un format propre au logiciel Kraken (<http://kraken.re/>). Ces modèles d'apprentissage pourront servir à la transcription automatique de textes à la graphie proche.

Base de données comprenant les Entités Nommées

Une liste d'Entités Nommées (Lieux et Personnes) sera extraite des transcriptions.

Cette base de données est construite en deux étapes : la première étape consiste à établir une détection manuelle de différentes entités nommées (Personnes, Lieux) afin de créer des données d'apprentissage.

Dans un deuxième temps, à partir de ces données d'entraînement, des bibliothèques de détection d'entités nommées tel la librairie Spacy (<https://spacy.io/>) effectuent une détection automatique.

Ces données seront proposées dans un format CSV.

Les Personnes et les Lieux seront alignés dans la mesure du possible avec le référentiel d'autorité Idref (<https://www.idref.fr/>) et avec création d'un identifiant pérenne de type IdRef pour les entités manquantes.

Données Cartographiques

En croisant les données du cadastre de Vasserot (1810-1836) avec les données des registres capitulaires et d'autres documents du chapitre, il sera possible de reconstituer les données spatiales concernant le bâti et le parcellaire du cloître Notre-Dame avec des degrés de précision variables selon les états chronologiques.

Les données produites seront de type Shapefile. Elles seront intégrées et distribuées par la plateforme d'information géohistorique sur Paris, Alpage (<https://alpage.huma-num.fr/>). La volumétrie des données produites sera précisé ultérieurement.

2. Documentation et qualité des données

Les manuscrits seront disponibles soit sur le site Internet des Archives Nationales, soit sur le site de la BNF, soit sur le site Internet de la Bibliothèque Mazarine. Les manuscrits seront décrits par des métadonnées au format EAD utilisées par ces institutions pour l'ensemble de leurs documents.

Les manuscrits, les textes, les modèles IA et les bases de données seront déposés sur l'entrepôt de données Nakala.

Cet entrepôt de données permet de décrire les données à l'aide du vocabulaire Dublin Core (<https://www.dublincore.org/>), soit le titre de la page, une description, des mots clefs associés, l'ensemble des contributeurs, la date de création et l'historique des dates de modifications du document, un identifiant DOI, etc.). Les liens entre les différents éléments (transcriptions, bases de données, numérisations) seront renseignés par différents champs (dc:source, dc:isreferencedby,): ces métadonnées permettront par exemple à un utilisateur d'effectuer la comparaison entre une transcription d'un manuscrit et ce manuscrit.

Le dépôt sur un entrepôt de données ainsi que l'utilisation d'un vocabulaire contrôlé permettront de répondre aux principes FAIR.

Le contrôle de la qualité des numérisations des manuscrits sera effectué par les services compétents des Archives Nationales, de la BNF et de la Bibliothèque Mazarine.

Le contrôle de la qualité de la transcription des registres capitulaires sera effectué par un comité éditorial constitué d'experts scientifiques des corpus étudiés.

Les bases de données seront accompagnées par un schéma de données (de type schema.json) décrivant les différents champs ainsi que leurs valeurs possibles. Ces schémas de données permettront de valider le contenu des bases de données. Un historique des modifications des données sera également mis en place.

Le 1er volume du catalogue de 1733 édité au format TEI sera proposé au projet THECAE (Corpus d'inventaires anciens de livres manuscrits et imprimés <https://www.unicaen.fr/services/puc/sources/thecae/accueil>), cette proposition sera examinée et soumise à l'approbation d'un comité éditorial.

3. Stockage et sauvegarde pendant le processus de recherche

Stockage et sauvegarde pendant le processus de recherche :

Les données de recherche (Transcriptions, Modèles, Base de données) produites pendant le processus de recherche seront stockées sur les différents outils de travail proposés par Huma-Num (Gitlab, Sharedocs, Huma-Num Box, Nakala).

Les manuscrits et les données (à l'exclusion de celles déposés sur le dispositif Huma-Num Box) sont sauvegardées par un automate une fois par nuit sur bande magnétique. Une copie de secours est effectuée dans un second bâtiment adjacent au premier.

Les caractéristiques de la sauvegarde journalière sur le dispositif de bandes magnétiques proposés par Huma-Num sont

- en mode incrémental
- par le logiciel TSM (Tivoli) d'IBM
- de tous les fichiers modifiés par rapport à la veille
- avec conservation de 6 versions du même fichier
- avec rétention durant 1 an d'un fichier totalement supprimé à la source
- sur deux robotiques de bandes LTO dans les deux bâtiments distincts

Une sauvegarde de ces données de recherche sera aussi effectué mensuellement sur un serveur de Paris 1 Panthéon-Sorbonne, cette redondance permettant de se prémunir d'un incident grave.

Les données présentes sur le dispositif Huma-Num Box seront copiés sur deux lieux de stockage différents sur disque ainsi que sur bande.

En cas d'incident, les données seront restaurées sur demande auprès des services d'Huma-Num ou depuis le serveur de l'université Paris 1 - Panthéon Sorbonne.

Les données ne comportent pas de données sensibles.

Seuls les membres participant au projet E-NdP pourront accéder en écriture aux données de recherche par l'intermédiaire d'une authentification de type login/motdepasse. L'authentification est géré dans la majorité des outils par l'outil de gestion de login Human-Id (<https://humanid.huma-num.fr/>)

4. Exigences légales et éthiques, codes de conduite

Question sans réponse.

Les données (manuscrits et données) sont diffusées sous Licence Creative Commons Attribution 4.0 (CC BY 4.0) : les données sont librement ré-utilisables, sous condition de citer leurs auteur-es et/ou leurs origines.

A l'exception des manuscrits de la Bibliothèque Mazarine qui sont diffusés sous licence Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International (CC BY-NC-ND 4.0).

L'ensemble des données n'implique pas de problèmes éthiques.

5. Partage des données et conservation à long terme

Les données produites seront partagées au fur et à mesure de leur finalisation afin de permettre leurs réutilisations et de favoriser leurs études au delà des personnes impliquées dans l'ANR E-NdP. Les numérisations des manuscrits seront partagées lors la première année du projet.

Les numérisations des manuscrits issues des Archives Nationales seront accessibles sur le site des Archives Nationales par navigation sur le site des Archives Nationales (<https://www.archives-nationales.culture.gouv.fr/>) et sur l'entrepôt de données d'Huma-Num Nakala par navigation sur le site Nakala (<https://www.nakala.fr/>) et par le protocole IIIF.

Les numérisations des manuscrits issues de la Bibliothèque Mazarine et de la BNF seront accessibles par le protocole IIIF et par navigation sur le site de la Bibliothèque Mazarine (<https://www.bibliotheque-mazarine.fr/fr/>) et de la BNF (<https://gallica.bnf.fr/>)

Les textes, base de données, modèles d'Intelligence Artificielle seront versés sur l'entrepôt de données Nakala.

Un identifiant pérenne de type DOI sera attribué à chaque numérisation, textes, base de données, modèles d'Intelligence Artificielle déposée sur l'entrepôt de données Nakala. Leurs métadonnées décrivant ces données seront moissonnées par les moteurs de recherche usuels (Google, Bing, etc..) mais aussi par le moteur de recherche d'Huma-Num Isidore (<https://isidore.science>) et par le consortium Datacite (<https://search.datacite.org>) afin de faciliter leur dissémination.

Cet ensemble de données sera téléchargeable par leur DOI, par l'API d'accès proposé par Nakala (<https://api.nakala.fr>) et par navigation sur le site de l'entrepôt de données Nakala.

Un identifiant pérenne de type ark sera attribué à chaque numérisation déposée sur le site de la Bibliothèque Mazarine.
Il n'est prévu aucune restriction sur le partage de données.

Les numérisations des manuscrits seront conservés à la fois par les Archives Nationales et sur l'entrepôt de données Nakala.
Les transcriptions finales des registres, les modèles d'IA ainsi que la base finale des entités nommées seront conservés sur l'entrepôt de données Nakala.

La sauvegarde des données des Archives Nationales et de la Bibliothèque Mazarine est effectuée par le centre national du microfilm et de la numérisation (CNMN) sur le site de St-Gilles-du-Gard.

L'ensemble des données sera accessible à tout public :

soit à l'aide d'un simple navigateur à travers les sites des Archives Nationales, de la Bibliothèque Mazarine et des services d'exposition des données proposés par Huma-Num (Nakala, Recherche-Isidore, Site Internet du Projet)

soit par l'API de Nakala (<https://api.nakala.fr/doc>), cet interface permettant aux métadonnées des données du projet d'être machine-actionable.

A chaque manuscrit sera attribué soit un permalien (Archives Nationales), soit un ark (Bibliothèque Mazarine) soit un DOI (Nakala).
A chaque texte, modèle d'IA et base de données sera attribué un DOI (Nakala).

6. Responsabilités et ressources en matière de gestion des données

Pierre Brochard (<https://orcid.org/0000-0003-1955-556X>) est le responsable des données issues du traitement des registres capitulaires.

Hélène Noizet (<https://orcid.org/0000-0003-2221-6104>) est la responsable des données spatiales du projet.

Patrick Latour est le responsable des données issues de la transcription du catalogue de la Bibliothèque de 1733.

Question sans réponse.