



# Analysis of Researchers' Online Behavior

Semester project in Data Science

Author: Mohamed Beyrem Kaabachi

Supervisor: Dr. Simon Dumas Primbault

# TABLE OF CONTENTS

01

INTRODUCTION

02

LITERATURE REVIEW

03

DATA ACQUISITION  
& DESCRIPTION

04

METHODOLOGY

05

RESULTS &  
DISCUSSION

06

QUALITATIVE RESEARCH

07

FUTURE WORK

# INTRODUCTION

The image features a dark, moody background with a stack of books on the right side. The books are slightly out of focus, showing their spines and edges. The word 'INTRODUCTION' is prominently displayed in the center-left in a white, classic serif typeface. The overall aesthetic is academic and sophisticated.

# INTRODUCTION

- Studying libraries and library users is an important subject in the field of ethnography. [3] The goal of **ethnography** is to study the culture and social organization of a particular group or community.
- While observing the behavior of researchers in libraries and laboratories is commonplace in the humanities, the **computational** study of digital research practices is only in its early days.
- This project aims at documenting the behavior of scientists on online platforms by making sense of the **digital traces** they generate while navigating.

# INTRODUCTION

- In the context of this project, work has been done on the **browsing logs** of **Gallica** which is the digital library for online users of the Bibliothèque nationale de France and its partners.
- The documents present in that website form an encyclopaedic and comprehensive library, representative of major French authors and of the various trends of reflection and research over the centuries.
- 4.9+ million documents.

# LITERATURE REVIEW

The background of the image shows a stack of several books resting on a dark wooden surface. The books are slightly out of focus, with their spines and edges visible. The lighting is soft, creating a warm, scholarly atmosphere. The text 'LITERATURE REVIEW' is prominently displayed in the foreground in a white, classic serif typeface, arranged in two lines.

# LITERATURE REVIEW

- The study of user's behavior in conventional libraries has already been done in previous research. According to a survey [7], there has been a total of 81 studies that researched library users' behaviors as of 2012
- The researcher could either observe or watch subjects in their natural settings [5], conduct interviews with library users [6] or conduct fieldwork in order to contextualize research results. During this project, we are trying to implement a combination of these methods while applying them to a digital resource.

# LITERATURE REVIEW

- Interest in digital resources is rising as the platforms become more and more adopted by its users.
- Few notable resources that went over similar subjects:
  - OpenEdition platform - Romain Deveaud
  - Wikipedia - Robert West
  - Gallica - Nouvellet et al.
- Although our research works on similar data, our goal during this project is to interpret users' behaviors through paths that specifically rely on document themes and how one could go from a theme to another.



# DATA MANIPULATION AND ENRICHMENT

# DATA DESCRIPTION

- Browsing logs over the month of **April 2016**.

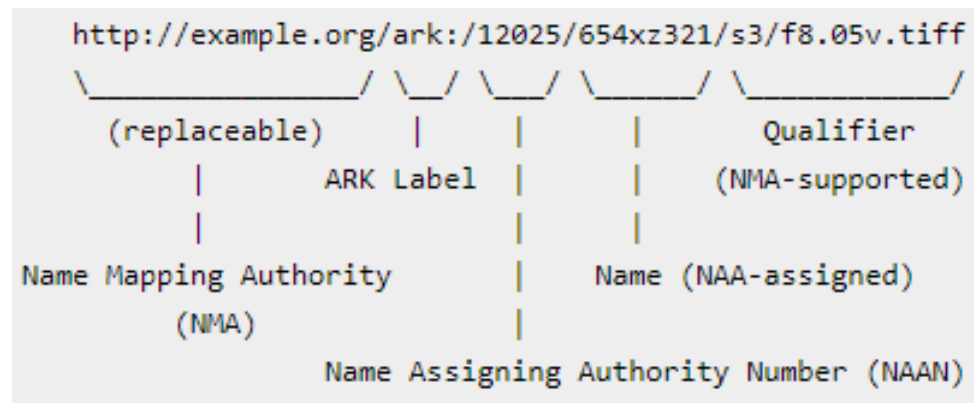
Information available	Description
Hashed IP	Anonymized IP addresses
Country	Country of requester
City	City of requester
Complete date	Date of request
Request	HTTP requests / ARK requests etc..
Protocol	Communication Protocol
Answer Code	Response Status Codes
Length	Length of request
Referring Website	Website from which the user comes

# DATA DESCRIPTION

- Request example

```
##6958a5de61066cceb1831af6e2f0fc76##United States##Madison##- -  
[04/Mar/2017:00:31:03 +0100] "GET /ark:/12148/bpt6k9708547.lowres HTTP/1.1" 200 78109  
"http://gallica.bnf.fr/" "Mozilla/5.0 (Windows NT 6.1; Trident/7.0; rv:11.0) like Gecko"
```

- ARKs the way documents can be requested through Gallica's website. Each ARK is an identifier for a document and does not contain any semantically recognizable information about the document itself. Although it is possible to collect metadata about a document from an ARK by querying the website.



# DATA ENRICHMENT

- Query Gallica website using the OAI record: Service for retrieving bibliographic information from a document.

Source: [API Document de Gallica | Api \(bnf.fr\)](#)

- XML answer, several interesting information. Some of them are optional.
- Year
- Language
- Theme – Dewey Classification
- Title

Exemple : <https://gallica.bnf.fr/services/OAIRecord?ark=bpt6k5738219s>

```
<?xml version="1.0" encoding="UTF-8" ?>
<results countResults="1" resultType="LuceneOAIRecordSearch" searchTime="0:00:00.001">
  <notice>
    <record xmlns="http://www.openarchives.org/OAI/2.0/">
      <header>
        <identifier>oai:bnf.fr:gallica/ark:/12148/bpt6k5738219s</identifier>
        <timestamp>2012-01-27</timestamp>
        <setSpec>gallica:theme:8:84</setSpec>
        <setSpec>gallica:typedoc:monographies</setSpec>
      </header>
      <metadata>
        <oai_dc:dc xmlns:dc="http://purl.org/dc/elements/1.1/" xmlns:oai_dc="http://www.openarchives.org/OAI/2.0/oai_dc/">
          <dc:identifier>https://gallica.bnf.fr/ark:/12148/bpt6k5738219s</dc:identifier>
          <dc:title>La plage d'Etretat par l'auteur de "Monsieur X et Mme ***"</dc:title>
          <dc:publisher>Michel Levy (Paris)</dc:publisher>
          <dc:date>1868</dc:date>
          <dc:format>In-18</dc:format>
          <dc:language>fre</dc:language>
          <dc:relation>Notice du catalogue : http://catalogue.bnf.fr/ark:/12148/cb33539190h</dc:relation>
          <dc:type xml:lang="eng">text</dc:type>
          <dc:type xml:lang="fre">monographie imprimée</dc:type>
          <dc:type xml:lang="eng">printed monograph</dc:type>
          <dc:format>application/pdf</dc:format>
          <dc:source>Bibliothèque nationale de France, département Littérature et art, Y2-59413</dc:source>
          <dc:rights xml:lang="fre">domaine public</dc:rights>
          <dc:rights xml:lang="eng">public domain</dc:rights>
        </oai_dc:dc>
      </metadata>
    </record>
  </notice>
  <mode_indexation>text</mode_indexation>
  <nqamoyen>092.57</nqamoyen>
  <provenance>bnf.fr</provenance>
  <source>Bibliothèque nationale de France, département Littérature et art, Y2-59413</source>
  <typedoc>monographies</typedoc>
  <date>1868</date>
  <title>La plage d'Etretat par l'auteur de "Monsieur X et Mme ***"</title>
  <sdewey>84</sdewey>
</results>
```

# TOOLS USED

- Large data provided in logs, reliance on EPFL IC Cluster to do our calculations.



- BeautifulSoup for XML Parsing / Data Manipulation using Pandas and Numpy

# METHODOLOGY

The image features a dark, moody background with a stack of books on the right side. The books are stacked horizontally, with their spines and edges visible. The top book has a dark cover, while the ones below it appear lighter. The entire scene is set on a wooden surface, with the wood grain visible in the lower right corner. Overlaid on the left side of the image is the word 'METHODOLOGY' in a large, white, serif font. The text is centered vertically and horizontally relative to the left half of the image.

# METHODOLOGY

- Since we are interested in understanding researchers' behaviors, we first start our data manipulation by **creating sessions**. The goal here is to assign a session to each request, where a session is a group of requests emanating from the same user in a constrained period of time.
- A **session** is a series of request that a user made in Gallica website without hitting the inactivity threshold. ( Threshold = 60 minutes. )
- Filter only ARK requests

# METHODOLOGY

Session				
Session_1	ARK1	ARK2	ARK4	ARK5...
Session_2	ARK2	ARK3	ARK1	ARK6...



# METHODOLOGY

Session				
Session_1	ARK1 - Theme1 - Lang1	ARK2 – Theme2 – Lang2	...	...
Session_2	ARK2 - Theme2 - Lang2	ARK3 – Theme3 – Lang3	...	...

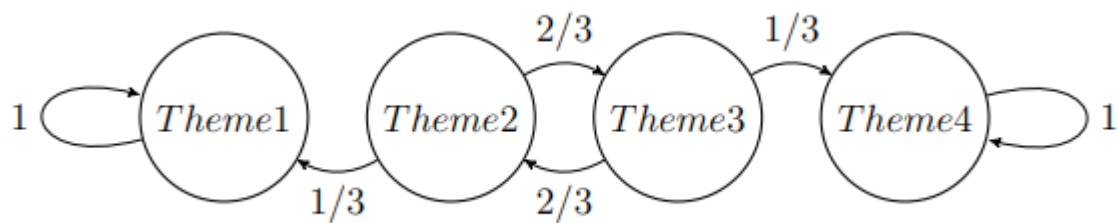
# RESULTS & DISCUSSION



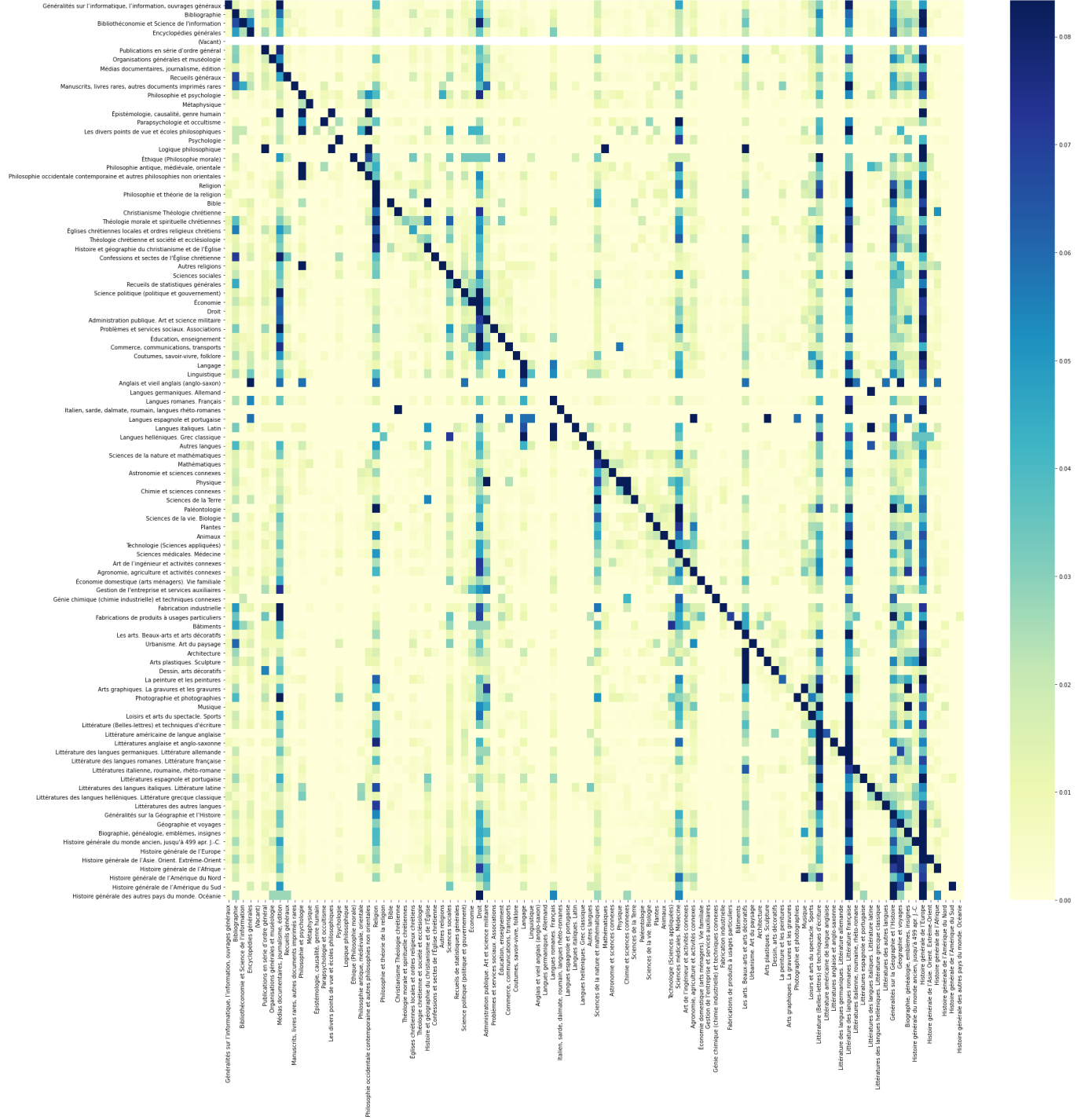
# USER PATHS AS MARKOV CHAINS

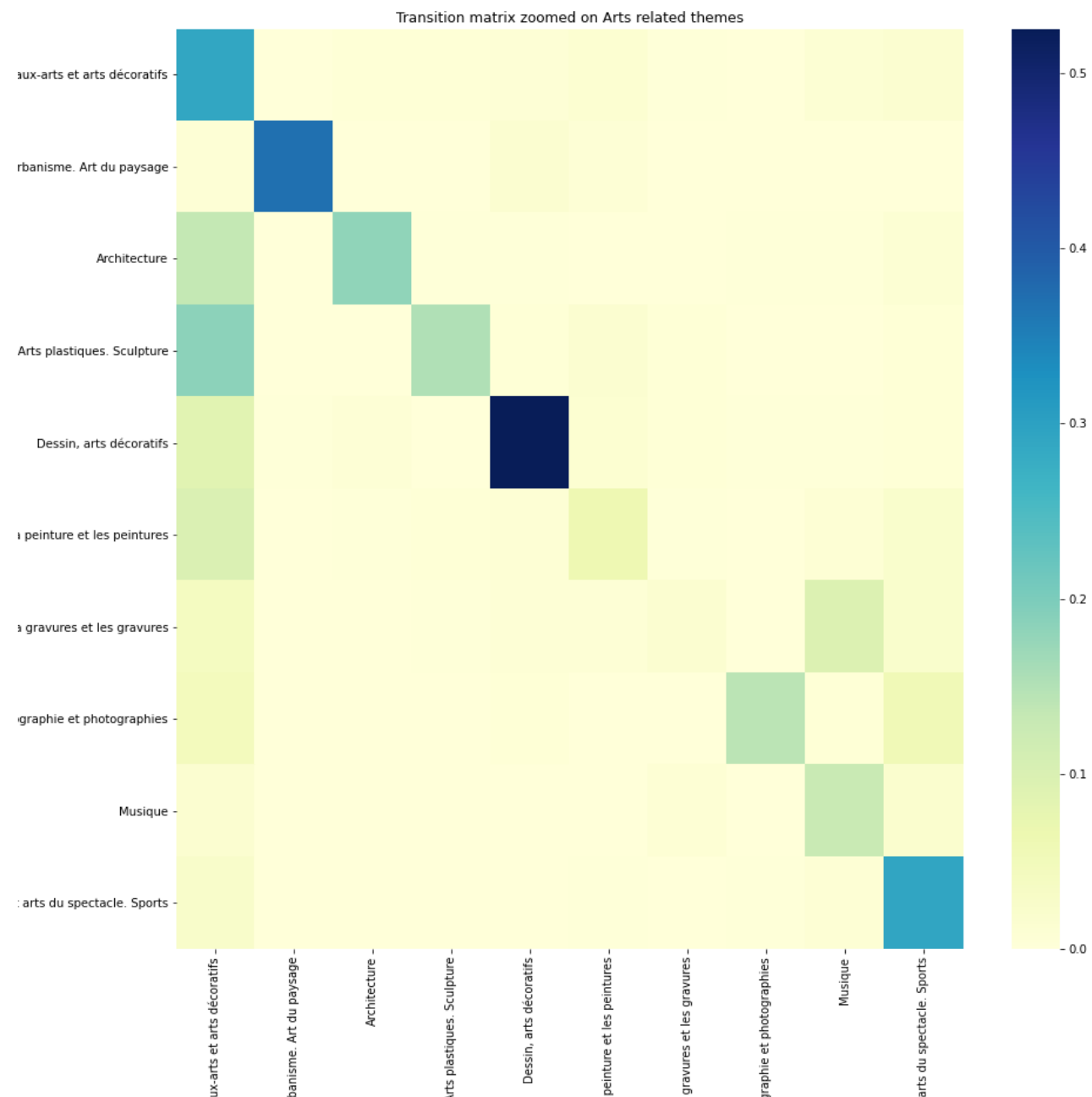
$$P_{(Theme_i, Theme_j)} = P(X_{n+1} = Theme_j | X_n = Theme_i)$$

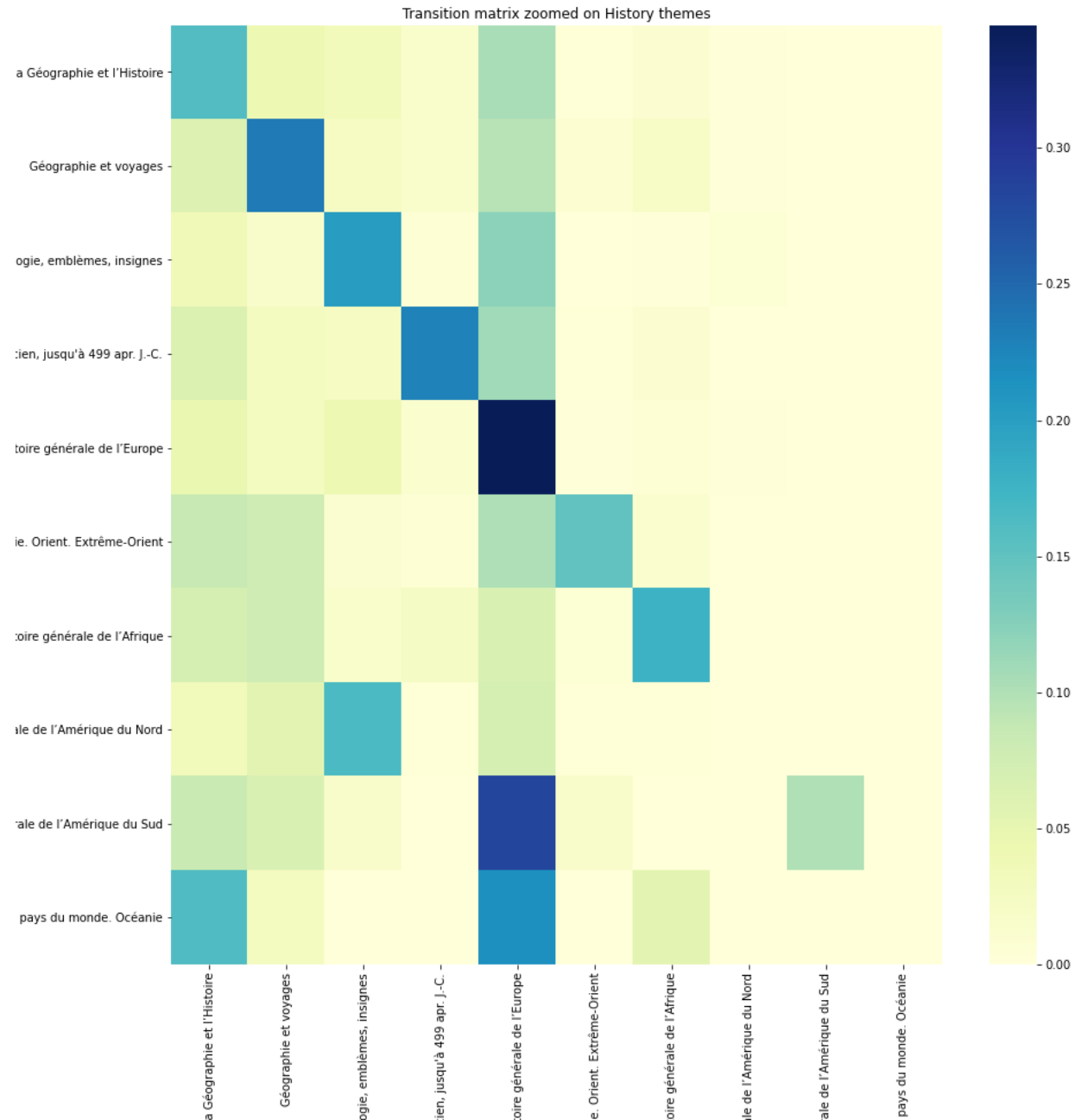
$$\hat{P}_{Theme_i, Theme_j} = \frac{n_{Theme_i, Theme_j}}{\sum_{k=1}^m n_{Theme_i, Theme_k}}$$



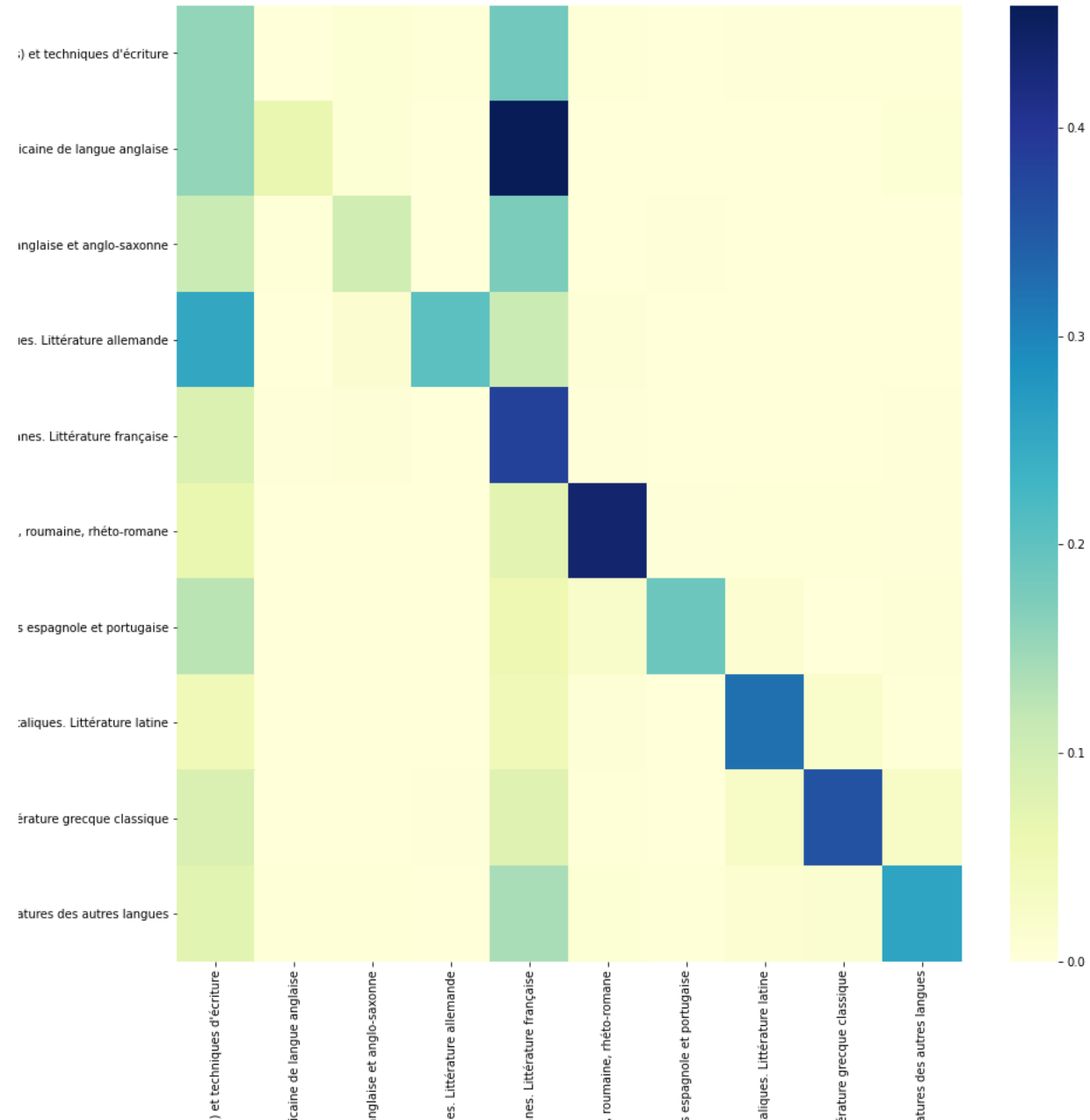
Transition matrix by dewey themes

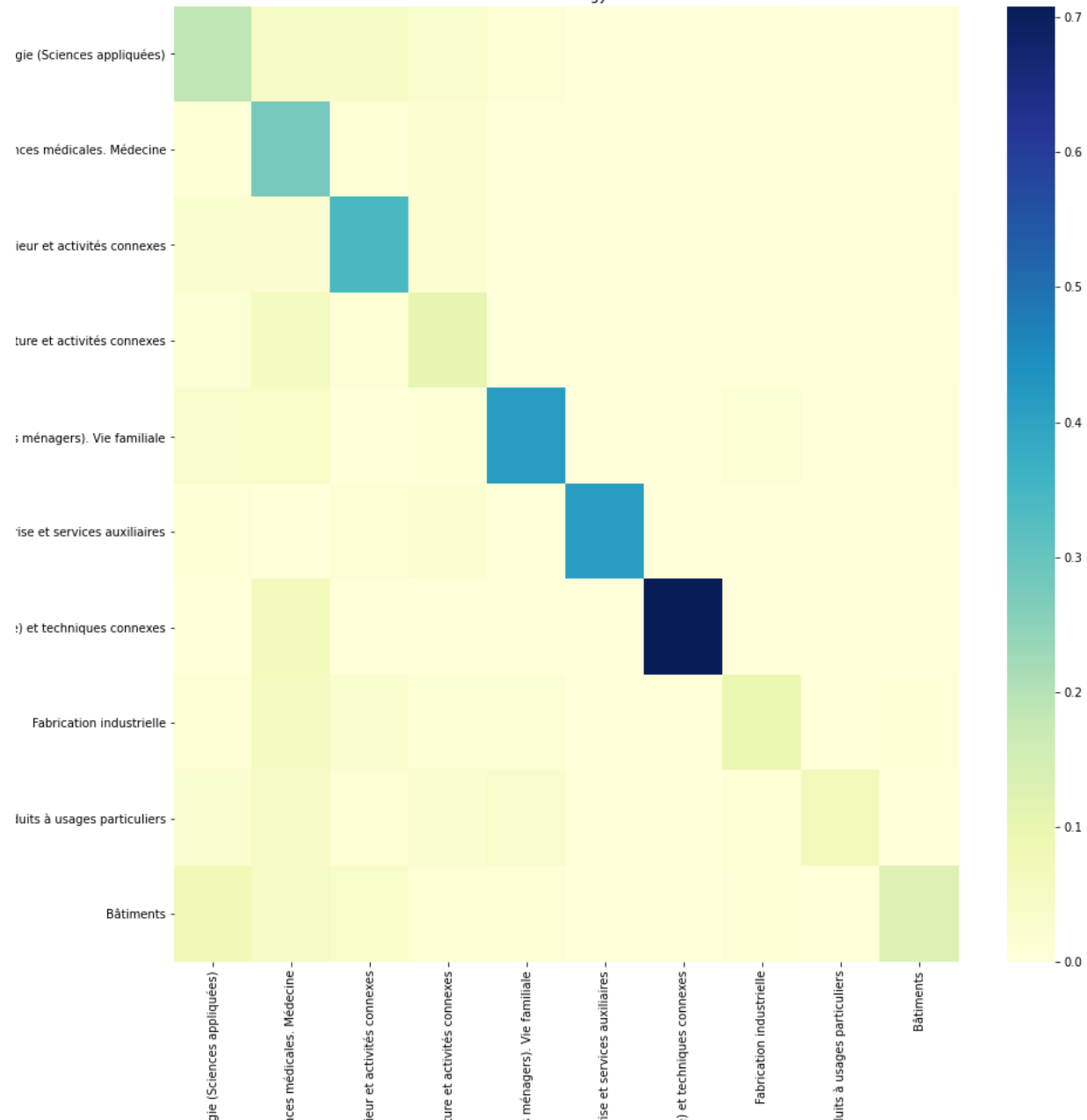






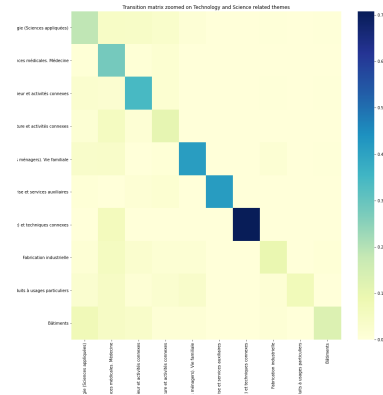
Transition matrix zoomed on Literature themes



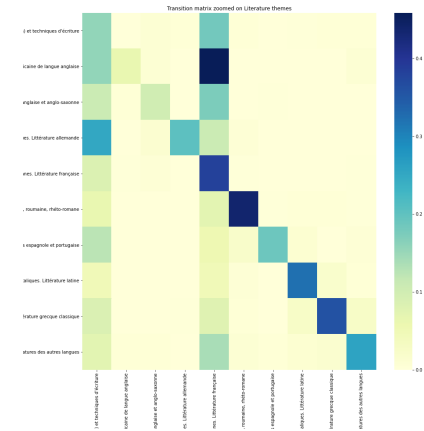




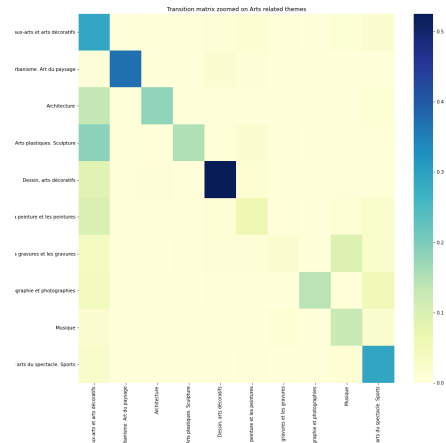
## Technology and Sciences



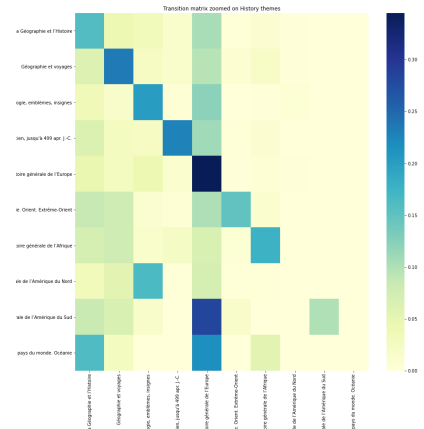
## Literature



## Arts



## History



# USER PATHS AS MARKOV CHAINS

$$PopularTheme = \operatorname{argmax}_j \left( \sum_{k=1}^m P_{Theme_k, Theme_j} \right)$$

- Recueils généraux: 0.01% probability of transition on average

$$PopularOutwardsTheme = \operatorname{argmax}_j \left( \sum_{k=1}^m P_{Theme_j, Theme_k} \right)$$

- Histoire Générale de l'Europe: 0.064% probability of transition on average

# REPRESENTATION OF PATHS IN A WORD EMBEDDING FORM

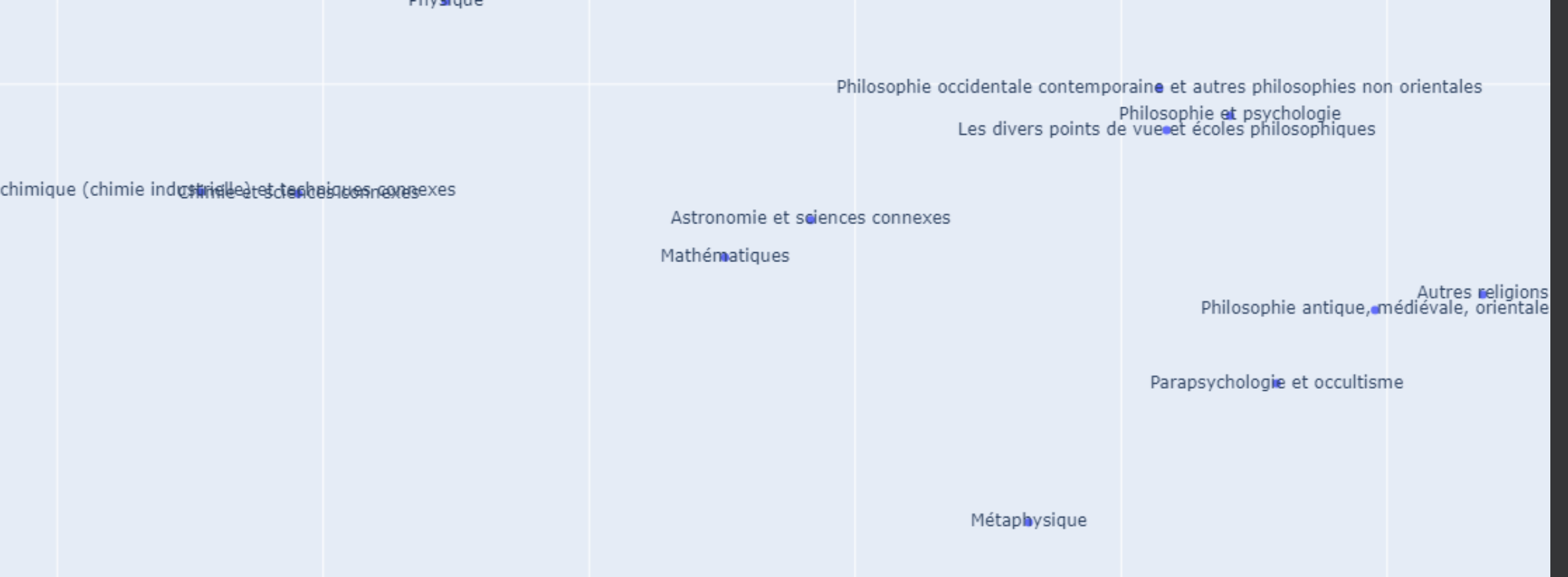
- Creating a corpus

Corpus	
Sentence 1	Theme1, Theme2, Theme3, Theme1, Theme4...
Sentence 2	Theme5, Theme2, Theme7, Theme8, Theme4...
Sentence 3	Theme6, Theme9, Theme7, Theme2, Theme1...

- Word2vec using Skip-gram
- Each Theme is represented as a vector of 200 dimensions.
- Visualization can be obtained with t-SNE.



# REPRESENTATION OF PATHS IN A WORD EMBEDDING FORM

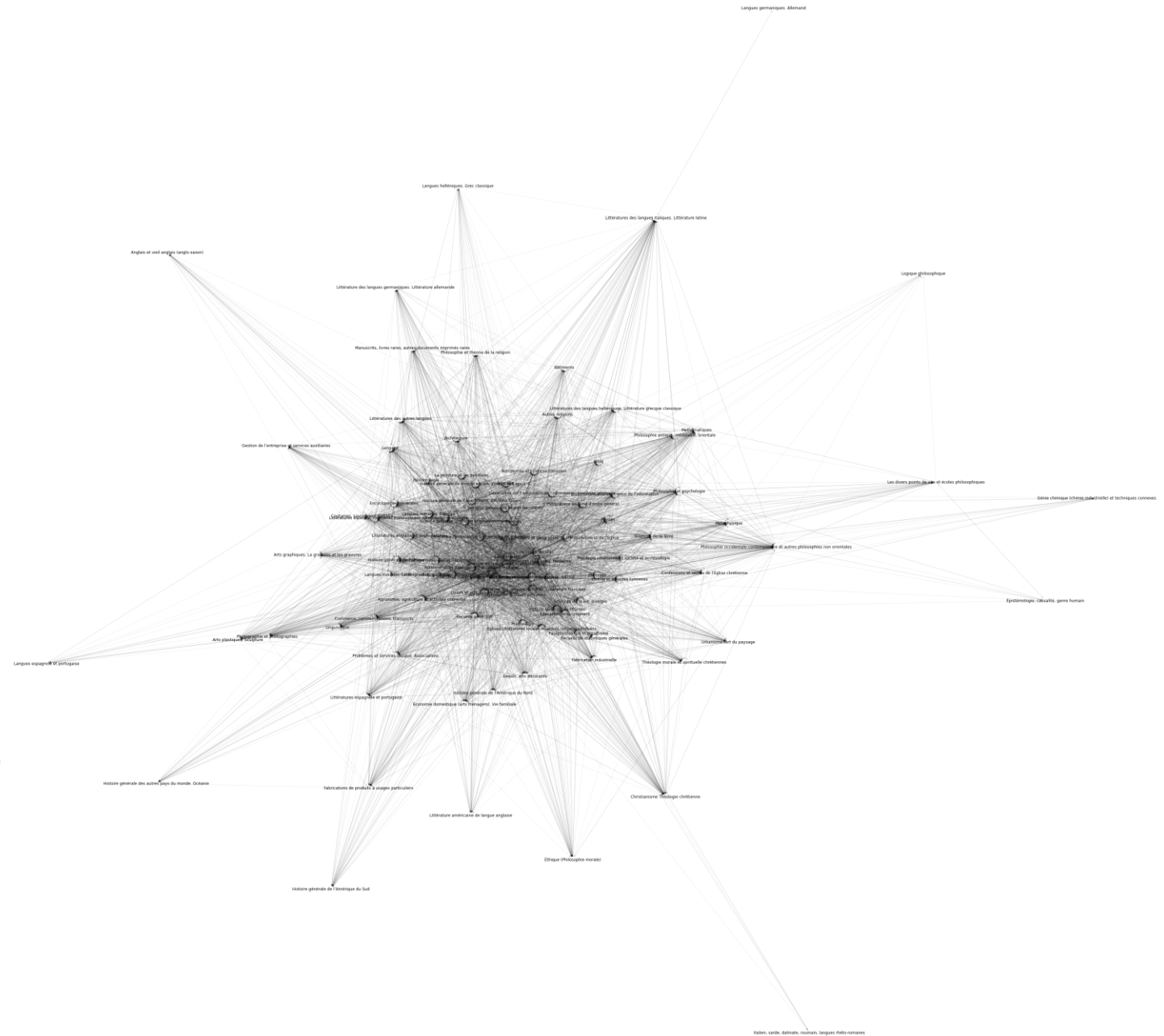


## REPRESENTATION OF PATHS IN A WORD EMBEDDING FORM

“Mathematics is, as it were, a sensuous logic, and relates to philosophy as do the arts, music, and plastic art to poetry.” — K. Shegel

# REPRESENTING THEMES AS A NETWORK

- All themes in the graph are added as vertices.
- Number of transitions from a theme to another as edges.
- Nose positions done using Fruchterman-Reingold[1] force-directed algorithm.
- Make use of graph theory features such as betweenness centrality.



# REPRESENTING THEMES AS A NETWORK

The betweenness centrality of a node  $v$  is given by the expression:

$$g(v) = \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

where  $\sigma_{st}$  is the total number of shortest paths from nodes  $s$  to node  $t$  and  $\sigma_{st}(v)$  is the number of those paths that pass through  $v$ .

A node strength is then given by the sum of the weights of its adjacent edges. With  $a_{ij}$  and  $w_{ij}$  being adjacency and weight matrices between nodes  $i$  and  $j$

$$s_i = \sum_{j=1}^N a_{ij} w_{ij}$$

# REPRESENTING THEMES AS A NETWORK

Themes	Betweenness centrality
Histoire générale de l'Europe	0.027
Littérature des langues romanes. Littérature française	0.024
Littératures des langues italiques. Littérature latine	0.0216
Les arts. Beaux-arts et arts décoratifs	0.0214
Médias documentaires, journalisme, édition	0.0174



# QUALITATIVE RESEARCH

A stack of several books is visible on the right side of the image, resting on a wooden surface. The books are slightly out of focus, with their spines and edges visible. The overall lighting is soft and warm, creating a scholarly or academic atmosphere.

# QUALITATIVE RESEARCH

- Qualitative research: collecting and analyzing non-numerical data (e.g., text, video, or audio) to understand concepts, opinions, or experiences.
- Interview with a Gallica user, key points:
  - Non-linearity of consultations
- Participant observation

Le Félibre Félix Gras : poètes administratifs des comités des départements, en temps de paix et en temps de guerre / par le Dr Duchausse  
 La Fleur mauve / par Véra Dérelle',  
 Principj elementari di musica ... compilati da B. Asioli',  
 L'informateur des aliénistes et des neurologistes",  
 Une brigade allemande d'infanterie au combat : Borny, Noisseville, Villers-Bretonneux, Saint-Quentin / par le capitaine Grange,,  
 Recherches sur les générations spontanées et sur la matière, ses propriétés et ses lois / par Michel-Hyacinthe Deschamps,...',  
 Honneur à Faïdherbe',  
 1792 à propos de 1892 : les martyrs de septembre / Le Père J. Delbrel,... ; avec une lettre-préface de Mgr d'Hulst,...",  
 Style de Pascal / par M. Nault,...',  
 Itinéraire de Paris à Jérusalem et de Jérusalem à Paris, en allant par la Grèce et revenant par l'Égypte, la Barbarie et l'Espagne  
 Du culte de la Sainte Vierge dans l'Église catholique / par le cardinal J. H. Newman, traduction revue et corrigée par un bénédictin  
 Les Trois mousquetaires. Tome 5-6 / par Alexandre Dumas',  
 Le Concordat, sa négociation, ses dix-sept articles, son histoire de 1801 à 1903 , par Auguste Body',  
 L'ami des hommes, ou Exposé simple des moyens de conserver la santé et de prolonger autant que possible la durée de la vie. Traité  
 La maison au Moyen âge dans le Midi de la France : actes du colloque de Cahors des 6, 7 et 8 juillet 2006 / [organisées par la] Académie  
 Adèle / par l'auteur de Jean Sbogard",  
 Histoire et description des voies de communication aux États-Unis... par Michel Chevalier. Table analytique et alphabétique des  
 Épisodes de la guerre de Trente ans. Le maréchal de Guébriant (1602 à 1643) / par le vicomte de Noailles',  
 Sommation du château de Montsaugéon par Guébriant (19 avril 1639) . (Signé : Emile Longin.)',  
 Les papiers de Pierre Rotrou de Saudreville, secrétaire du maréchal de Guébriant : introduction / publiés par Léonce Person,...'  
 Histoire généalogique de la maison de Lantivy, de ses alliances et des seigneuries qu'elle a possédées, Bretagne, Maine, Anjou et  
 son (Écosse et France)... par Théodore Courtaux et le Cte de Lantivy de Trédion",  
 Revue de l'art français ancien et moderne",  
 Bulletin... / Société de l'histoire du protestantisme français : études, documents, chronique littéraire",  
 Le Cabinet historique : revue... contenant, avec un texte et des pièces inédites, intéressantes ou peu connues, le catalogue géométrique  
 toire de l'ancienne France et de ses diverses localités, avec les indications de sources, et des notices sur les bibliothèques et  
 Les papiers de Pierre Rotrou de Saudreville, secrétaire du maréchal de Guébriant : introduction / publiés par Léonce Person,...'  
 Épisodes de la guerre de Trente ans. Le maréchal de Guébriant (1602 à 1643) / par le vicomte de Noailles',  
 Sommation du château de Montsaugéon par Guébriant (19 avril 1639) . (Signé : Emile Longin.)',  
 Revue d'histoire moderne et contemporaine / Société d'histoire moderne",  
 Mémoires de Messire Philippe de Comines, seigneur d'Argenton. Tome 1 / , où l'on trouve l'histoire des rois de France Louis XI et  
 titres, contrats et instructions... par messieurs Godefroy, augmentée par M. l'abbé Lenglet Du Fresnoy",  
 Catalogue de la bibliothèque de la ville de Pau. HISTOIRE Partie 2 / par L. Soulice,... [et Gabriel Loirette]',  
 Histoire de dix ans de la Franche-Comté de Bourgogne, 1632-1642 / par Girardot de Nozeroy ["sic"], seigneur de Beauchemin... ; [et]  
 Revue de Gascogne : bulletin mensuel du Comité d'histoire et d'archéologie de la province ecclésiastique d'Auch",  
 Revue d'histoire littéraire de la France",  
 Mémoires du maréchal de Turenne. Tome I. 1643-1653 / publiés, pour la Société de l'histoire de France, d'après le manuscrit autographe  
 Mélanges historiques : choix de documents. 1, Tome premier. Tome 1',  
 Revue de l'art français ancien et moderne",  
 Histoire généalogique de la maison de Lantivy, de ses alliances et des seigneuries qu'elle a possédées, Bretagne, Maine, Anjou et  
 son (Écosse et France)... par Théodore Courtaux et le Cte de Lantivy de Trédion",  
 Le Cabinet historique : revue... contenant, avec un texte et des pièces inédites, intéressantes ou peu connues, le catalogue géométrique  
 toire de l'ancienne France et de ses diverses localités, avec les indications de sources, et des notices sur les bibliothèques et  
 Catalogue de la bibliothèque de la ville de Pau. HISTOIRE Partie 2 / par L. Soulice,... [et Gabriel Loirette]',  
 Histoire de dix ans de la Franche-Comté de Bourgogne, 1632-1642 / par Girardot de Nozeroy ["sic"], seigneur de Beauchemin... ; [et]  
 Mémoires de Messire Philippe de Comines, seigneur d'Argenton. Tome 1 / , où l'on trouve l'histoire des rois de France Louis XI et  
 titres, contrats et instructions... par messieurs Godefroy, augmentée par M. l'abbé Lenglet Du Fresnoy",  
 Revue d'histoire moderne et contemporaine / Société d'histoire moderne",  
 Revue de Gascogne : bulletin mensuel du Comité d'histoire et d'archéologie de la province ecclésiastique d'Auch",  
 Mélanges historiques : choix de documents. 1, Tome premier. Tome 1',  
 Revue d'histoire littéraire de la France",  
 Mémoires du maréchal de Turenne. Tome I. 1643-1653 / publiés, pour la Société de l'histoire de France, d'après le manuscrit autographe

# QUALITATIVE RESEARCH

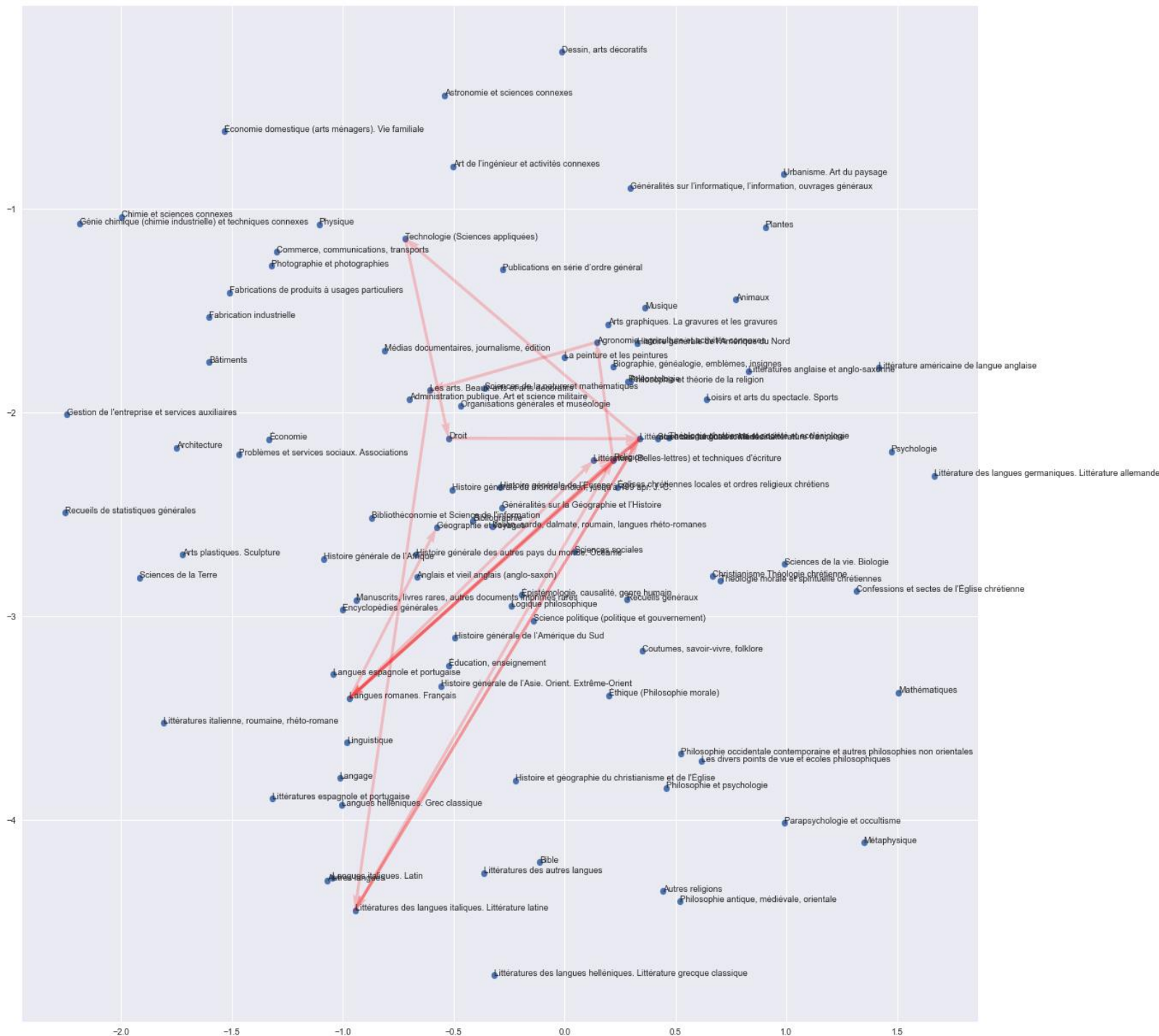
- Path presented to user
- Maréchal de guébriant
- Maréchaux de France

- Path presented to user

- Path presented to user

- Paths where it is more difficult to understand what the user is looking for.

- Paths where it is more difficult to understand what the user is looking for.



# Key points

- Mixed methods
- Understanding user movements, exploratory research/targeted research.  
Making the data “talk” and validation from a Gallica user.
- Topological Data Analysis

A stack of several books is visible on the right side of the image, resting on a wooden surface. The books are slightly out of focus, with the top book's cover and the edges of the pages visible. The background is a dark, textured surface, possibly a desk or table, with some faint lines and shadows. The text 'FUTURE WORK' is overlaid in a large, white, serif font, centered horizontally and positioned in the lower half of the image.

# FUTURE WORK

# Future work

- Distances between sessions using Word Mover's distance.
- Clustering on sessions and identifying trends
- Topological Data Analysis



A stack of several books is visible on the right side of the image, resting on a wooden surface. The books are slightly out of focus. A dark, semi-transparent overlay covers the entire image, with the text 'Thank you.' written in white serif font on the left side.

Thank you.

# REFERENCES

- Interest in digital resources is rising as the platforms become more and more adopted by its users.
- Few notable resources that went over similar subjects:
  - OpenEdition platform - Romain Deveaud
  - Wikipedia - Robert West
  - Gallica - Nouvellet et al.
- Although our research works on similar data, our goal during this project is to interpret users' behaviors through paths that specifically rely on document themes and how one could go from a theme to another.

## References

- [1] Thomas M. J. Fruchterman and Edward M. Reingold. “Graph drawing by force-directed placement”. In: *Software: Practice and Experience* 21.11 (1991), pp. 1129–1164. DOI: <https://doi.org/10.1002/spe.4380211102>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/spe.4380211102>. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/spe.4380211102>.
- [2] Andrew Y. Ng, Michael I. Jordan, and Yair Weiss. “On Spectral Clustering: Analysis and an Algorithm”. In: *Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic*. NIPS’01. Vancouver, British Columbia, Canada: MIT Press, 2001, pp. 849–856.
- [3] Craig Calhoun. *Dictionary of the social sciences*. Oxford University Press, 2002.
- [4] Laurens van der Maaten and Geoffrey Hinton. “Visualizing data using t-SNE”. In: *Journal of Machine Learning Research* 9 (Nov. 2008), pp. 2579–2605.
- [5] Lauren H. Mandel. “Toward an understanding of library patron wayfinding: Observing patrons’ entry routes in a public library”. In: *Library Information Science Research* 32.2 (2010), pp. 116–130. ISSN: 0740-8188. DOI: <https://doi.org/10.1016/j.lisr.2009.12.004>. URL: <https://www.sciencedirect.com/science/article/pii/S0740818810000034>.
- [6] Kylie Bailin. “Changes in Academic Library Space: A Case Study at The University of New South Wales”. In: *Australian Academic & Research Libraries* 42.4 (2011), pp. 342–359. DOI: 10.1080/00048623.2011.10722245. eprint: <https://doi.org/10.1080/00048623.2011.10722245>. URL: <https://doi.org/10.1080/00048623.2011.10722245>.
- [7] Michael Khoo, Lily Rozaklis, and Catherine Hall. “A survey of the use of ethnographic methods in the study of libraries and library users”. In: *Library Information Science Research* 34 (Apr. 2012), pp. 82–91. DOI: 10.1016/j.lisr.2011.07.010.
- [8] Tomas Mikolov, Quoc V. Le, and Ilya Sutskever. *Exploiting Similarities among Languages for Machine Translation*. 2013. arXiv: 1309.4168 [cs.CL].
- [9] Matt J. Kusner et al. “From Word Embeddings to Document Distances”. In: *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*. ICML’15. Lille, France: JMLR.org, 2015, pp. 957–966.
- [10] Adrien Nouvellet et al. “Analyse des traces d’usage de Gallica : Une étude à partir des logs de connexions au site Gallica.” In: (2017).
- [11] *Archival Resource Key (ARK) Identifiers*. URL: [https://n2t.net/e/ark\\_ids.html](https://n2t.net/e/ark_ids.html).
- [12] Romain Deveaud. “Modalités d’accès au savoir ouvert sur les plateformes d’OpenEdition”. In: ().
- [13] *Gallica – The BnF digital library*. URL: <http://www.bnf.fr/en/gallica-bnf-digital-library>.
- [14] Stuart Hannabuss. *How to... Use ethnographic methods and participant observation*. URL: <https://www.emeraldgrouppublishing.com/how-to/observation/use-ethnographic-methods-participant-observation>.
- [15] Robert West - HUMAN NAVIGATION OF INFORMATION NETWORKS, 2016