

Analyse Visuelle de Réseaux Sociaux Historiques: Traçabilité, Exploration et Analyse

*Visual Analytics for Historical Social Networks:
Traceability, Exploration, and Analysis*

Thèse de doctorat de l'université Paris-Saclay et de Telecom Paris

École doctorale n°580 : Sciences et technologies de l'information et de la communication (STIC)
Spécialité de doctorat: Informatique
Graduate School : Informatique et Sciences du Numérique
Référent : Faculté des sciences d'Orsay

Thèse préparée au Laboratoire interdisciplinaire des sciences du numérique (Université Paris-Saclay, CNRS, Inria), et à Telecom Paris, sous la direction de Jean-Daniel FEKETE, Directeur de recherche et la co-direction de Christophe Prieur, Professeur des universités.

Thèse soutenue à Paris-Saclay, le JJ décembre 2022, par

Alexis PISTER

Composition du jury

Ulrik Brandes	Rapporteur & Examinateur
Professeur, ETH Zürich	
Guy Melançon	Rapporteur & Examinateur
Professeur, Université de Bordeaux	
Wendy Mackay	Examinateuse
Directrice de recherche, Univ. Paris-Saclay, CNRS, Inria, LISN	
Uta Hinrichs	Examinateuse
Reader, University of Edinburgh	
Laurent Beauguitte	Examinateur
Chargé de recherche, CNRS	
Jean-Daniel Fekete	Directeur de thèse
Directeur de recherche, Univ. Paris-Saclay, CNRS, Inria, LISN	
Christophe Prieur	Directeur de thèse
Professeur, Université Gustave Eiffel	

Titre: Analyse Visuelle pour l'Analyse de Réseaux Sociaux Historiques

Mots clés: 3 à 6 mots clefs (version en français)

Résumé:

Title: Visual Analytics for Historical Social Networks: Traceability, Exploration, and Analysis

Keywords: visual analytics, social network analysis, social network visualization, social history.

Abstract:

Historical Social Network Analysis is a method followed by social historians to model relational phenomena of the past such as kinship, political power, migrations, or business affiliations with networks using the content of historical documents. Through visualization and analytical methods, social historians are able to describe the global structure of such phenomena and explain individual behaviors through their network position. However, the inspection, encoding, and modeling process of the historical documents leading to a finalized network is complicated and often results in inconsistencies, errors, distortions, and traceability issues. For these reasons and usability issues, social historians are often not able to make thorough historical conclusions with current visualization tools. In this thesis, I aim to identify how visual analytics—the combination of data mining capabilities integrated into visual interfaces—can support social historians in their process, from the collection of their data to the answer to high-level historical questions. Towards this goal, I first formalize the workflow of historical network analysis in collaboration with social historians, from the acquisition of their sources to their final visual analysis, and propose to model historical sources into bipartite multivariate dynamic social networks with roles to satisfy traceability, simplicity, and document reality properties. This modeling allows a concrete representation of historical documents, hence letting users encode, correct, and analyze their data with the same abstraction and tools. I, therefore, propose two interactive visual interfaces to manipulate, explore, and analyze this type of data with a focus on usability for social historians. First, I present ComBiNet, which allows an interactive exploration leveraging the structure, time, localization, and attributes of the data model with the help of coordinated views, a visual query system, and comparison mechanisms. Finding specific patterns easily, social historians are able to find inconsistencies in their data and answer their questions. The second system, PK-Clustering, is a concrete proposition to increase the usability and effectiveness of clustering mechanisms in social network visual analytics systems. It consists in a mixed-initiative clustering interface that let social scientists create meaningful clusters with the help of their prior knowledge, algorithmic consensus, and exploration of the network. Both systems have been designed with continuous feedback from social historians, and aim to increase the traceability, simplicity, and document reality of the historical social network analysis process. I conclude with discussions on the potential merging of both systems and more globally on research directions towards better integration of visual analytics systems on the whole workflow of social historians. Such systems with a focus on usability can lower the requirements for the use of quantitative methods for historians and social scientists, which has always been a controversial discussion among practitioners.

Contents

1	Introduction	9
1.1	Social History and Historical Social Network Analysis	10
1.2	Visualization and Visual Analytics	12
1.3	Visual Analytics Supported Historical Network Research	15
1.4	Contributions and Research Statement	17
2	Related Work	23
2.1	Visualization	23
2.1.1	Information Visualization	24
2.1.2	Visual Analytics	26
2.2	Quantitative Social History	28
2.2.1	History, Social History, and Methodology	28
2.2.2	Quantitative History	29
2.2.3	Digital Humanities	30
2.3	Historical Social Network Analysis	32
2.3.1	Sociometry to SNA	33
2.3.2	Methods and Measures	34
2.3.3	Historical Social Network Analysis	36
2.3.4	Network Modeling	38
2.4	Social Network Visualization	39
2.4.1	Graph Drawing	39
2.4.2	Social Network Visual Analytics	41
3	HSNA Process and Network Modeling	45
3.1	Context	46
3.2	Related Work	47
3.2.1	History Methodology	47
3.2.2	Historian Workflows	48
3.3	Historical Social Network Analysis Workflow	48
3.3.1	Textual Sources Acquisition	49
3.3.2	Digitization	49
3.3.3	Annotation	50
3.3.4	Network Creation	50
3.3.5	Network Analysis and Visualization	51
3.4	Network modeling and analysis	51
3.4.1	Network Models	52
3.4.2	Bipartite Multivariate Dynamic Social Network	54
3.4.3	Examples	55
3.5	Applications	58

3.6	Discussion	58
3.7	Conclusion	59
4	ComBiNet: Visual Query and Comparison of Bipartite Dynamic Multivariate Networks with Roles	61
4.1	Context	62
4.2	Related Work	64
4.2.1	Graphlet Analysis	64
4.2.2	Visual Graph Querying	64
4.2.3	Visual Graph Comparison	65
4.2.4	Provenance	65
4.3	Task Analysis and Design Process	65
4.3.1	Use Cases	66
4.3.2	Tasks Analysis	68
4.4	The ComBiNet System	69
4.4.1	Visualizations	70
4.4.2	Query Panel	72
4.4.3	Comparison	78
4.4.4	Implementation	81
4.5	Use Cases	81
4.5.1	Construction sites in Piedmont (#1)	81
4.5.2	French Genealogy (#2)	83
4.5.3	Marriage acts in Buenos Aires (#3)	85
4.5.4	Sociology thesis in France	86
4.6	Formative Usability Study	87
4.6.1	Feedback	88
4.7	Discussion	89
4.8	Conclusion and Future Work	89
5	PK-Clustering	91
5.1	Context	91
5.2	Related Work	94
5.2.1	Graph Clustering	94
5.2.2	Semi-supervised Clustering	95
5.2.3	Mixed-Initiative Systems and Interactive Clustering	95
5.2.4	Groups in Network Visualization	96
5.2.5	Ensemble Clustering	96
5.2.6	Summary	97
5.3	PK-clustering	97
5.3.1	Overview	97
5.3.2	Specification of Prior Knowledge	99
5.3.3	Running the Clustering Algorithms	99
5.3.4	Matching Clustering Results and Prior Knowledge	100

5.3.5	Ranking the Algorithms	101
5.3.6	Reviewing the Ranked List of Algorithms	101
5.3.7	Reviewing and Consolidating Final Results	103
5.3.8	Wrapping up and Reporting Results	107
5.4	Case studies	108
5.4.1	Marie Boucher Social Network	108
5.4.2	Lineages at VAST	109
5.4.3	Feedback from practitioners	110
5.5	Discussion	112
5.5.1	Limitations	113
5.5.2	Performance	113
5.6	Conclusion	114
6	Conclusion	115
6.1	Summary	115
6.2	Discussion	115
6.3	Perspectives	117
6.4	Conclusion	119

List of Figures

1.1	Business contract originated from Nantes (France) during the 17th century. See [34] for more detail of the historian process to analyze her sources.	12
1.2	Marriage, partnership, trading, banking, and real estate networks of the powerful families of Florence from [105]. We can see the central position in the network of the Medici Family.	13
1.3	Abstraction of the VA process. It is characterized by continuous interactions between the data, visualizations, models, and knowledge. Image from [71].	14
1.4	Node-link diagram of a medieval social network of peasants, produced with a force-directed layout, commonly used in SNA softwares. Image from [?].	16
2.1	Categorization of visual variables which can be used to represent network data, resulting in many different network representations. Image from [8].	25
2.2	Anscombe quartet. The four datasets have the same descriptive statistics (average, variance, correlation coefficient) but very different structures. Image from [4].	26
2.3	TULIP software designed for application-independant network visual analytics [?]. The view shows a dataset among multiple interactive coordinated views. Users can also apply data mining algorithms on the data to extract interesting patterns.	27
2.4	Correspondence letters of Benjamin Franklin and his close relationships, using a map and an histogram, accessible online on the republic of letter website [?].	31
2.5	Moreno's original sociogram of a class of first grades from [94] (left). The diagram shows 21 boys (triangles) and 14 girls (circles). The same sociogram using modern practices generated from Gephi from [53] (right). The color encodes the number of connections incoming.	33
2.6	All possible graphlets of size 2 to 5 for undirected graphs	35
2.7	Cicero personal communication network represented with a node-link diagram. Image from [?]	37
2.8	Different criteria are proposed to enhance node-link diagram readability. Image from [75]	40
2.9	NodeTrix system showing a scientific collaboration social network with clusters. Each cluster is represented as a matrix, Image from [64].	41
2.10	Vistorian interface [128] used to explore a historical social network of business trades in the 17th century, with a coordinated node-link diagram and a matrice view.	42
3.1	HSNA workflow is split into five steps: textual sources acquisition, digitization, annotation, network creation, and network visualization/analysis. Practitioners typically have to do back and forth during the process. I list potential pitfalls for each step.	49

3.2 bipartite multivariate dynamic network modeling for two cases of marriage acts of example #3. Some marriage acts mention the parents of the spouses, which is a relationship different than the marriage in itself. This case can be modeled using a document model (a) or an event model (c) by splitting the document into several different event nodes. The other case refers to documents that do not mention the parents (b) and in that case, the network represents both the documents and the events with the same model. M: Marriage, H: Husband, W: Wife, T: Witness, (H/W)(M/F): Husband/Wife Mother/Father. Yellow links refer to parenting mentions/relationships.	59
4.1 The ComBiNet system used to compare two subgroups of a social network of contracts from [25], extracted with dynamic visual queries. (A) and (B) show the two visual queries created by the user in the query panel using an interactive node-link diagram editor (V6), dynamic query widgets (V7), and the equivalent Cypher script (V8). The right part shows ComBiNet's global interface in <i>comparison</i> mode: (V1) Network visualization panel, (V2) Map of the geolocalized nodes, (V3) Table of persons, (V4) Graph measures comparison, (V5) Attribute distribution plots, and (V9) Provenance tree. The two visual queries on the left, translated into Cypher queries below, select the "Menafoglio" family on the left, and the "Zo" family on the right, along with their construction contracts and close collaborators.	70
4.2 ComBiNet interface with the dataset of collaboration #1. The user selected the <i>_year</i> attribute, showing the distribution of document years with a histogram (bottom right), and coloring the documents node on the bipartite view (left) and map view simultaneously (right).	72
4.3 All link creation possibilities: Any link type (left), one selected link type, here guarantor (middle left), the union of several link types (middle right), several links with different types (right)	74
4.4 Visual queries created to answer questions 2 and 6 of our collaboration #1. (a) The visual query retrieves individuals who are mutually guarantor to each other in separate construction contracts. (b) The two visual queries retrieve the documents—along with the signatories—of Torino (<i>Turin</i> in french) (left) and of Torino surroundings (<i>Turin Territoire</i> and <i>Piemont</i>) (right)	74
4.5 Widget designs for the different attribute types: checkboxes for categorical attributes (top), text input for nominal attributes (middle), and a double slider for numerical attributes (bottom). The categorical attribute example shows the inputs letting users create new constraints for other attributes and other nodes.	75
4.6 Results of question 2 of collaboration #1: (a) shows a subset of the table view with every occurrence of the pattern found. (b) shows the summary panel, with the graph measures and the attributes view with the <i>origin</i> attribute selected and the Sankey option checked. It allows us to see the attribute distribution of the persons included in the pattern and see if there is a relationship between persons who are mutually guarantors and their origin.	77

4.7	Two ways of showing the distribution of “type chantier” (type of works), a categorical attribute with three possible values “ <i>religious</i> ”, “ <i>military</i> ”, and “ <i>civilian</i> ”. (a) A query matching the contracts made by the same person (<i>per1</i>) as an “approbator” (green link to <i>doc2</i>) after being a “guarantor” (blue link to <i>doc1</i>) using the constraint (<i>doc2._year > doc1._year</i>). (b) Stacked bar chart for the matches, the earlier contract (<i>doc1</i>), the older contract (<i>doc2</i>), and (c) Sankey diagram with the early values on the left and the last on the right. The Sankey diagram reveals the value changes between the two documents: the guarantor who worked initially on religious work switched to military work.	78
4.8	Provenance tree to answer question 2 of collaboration #1: left branch leads to Torino documents (the node is labeled as A) while right branch leads to surrounding documents (the node is labeled as B). The user hovers over one node, revealing a tooltip that shows the visualization of the node’s query..	79
4.9	Comparison table of the networks measures for Torino subgraph (A) and Torino surroundings subgraph (B).	80
4.10	Distribution of the type of constructions, the years, and the betweenness centrality for the documents and signatories of Torino (A), Torino surroundings (B), and the whole graph (top).	81
4.11	Menafoglio (a) and Zo (b) families retrieved with queries and highlighted in the bipartite node-link and map views.	82
4.12	Attributes distributions plots between the whole graph, the <i>Menafoglio</i> family (A), the <i>Zo</i> family (B), and $A \cap B$, for the <i>region</i> , <i>type_chantier</i> , <i>material type</i>	83
4.13	Migrations across departments over three generations	84
4.14	Sankey diagrams showing the migration of people in the 18th and 19th centuries, extracted from their birth and death places.	85
4.15	ComBiNet used to request persons appearing as husband, wife, or witness in two marriages which occurred at 70 years apart or more.	86
4.16	ComBiNet used for exploring theses of sociology defended in France between 2016 and 2021. The bipartite and map views show an overview of two visions of the network. The user selects the <i>region</i> attribute, showing the geographical distribution of the defended thesis.	87
4.17	Sociology thesis dataset explored with ComBiNet. The user constructed a visual query to see if there are symmetrical relationships between thesis directors and referees (or jury directors). The <i>region</i> attribute is selected with the Sankey option, letting the user see if there are correlations between the regions of the thesis found in this pattern.	88
5.1	Traditional Clustering. The output is a clustering, usually from a randomly chosen algorithm.	92
5.2	PK-clustering. The output is a clustering supported by algorithms and validated (fully or partially) according to the user’s Prior Knowledge.	93
5.3	Prior Knowledge specification, the user defined two groups composed of two members.	98

5.4	Red edges represent the prior knowledge matching	100
5.5	Two different modalities for the ranked list of algorithms. Top: persons are shown as circles. Bottom: aggregated view. Colors indicate the matching group. Gray indicates no match. White indicates extra nodes or clusters.	101
5.6	Reviewing and comparing results of multiple algorithms. One algorithm is selected to order the names and group them, but icons show how other algorithms cluster the nodes differently, summarized in the consensus bar on the left.	104
5.7	The user quickly drags on consecutive icons (in yellow) representing the suggestions made by one algorithm to validate node clustering. Once the cursor is released the validated nodes appear as squares icons in the Consolidated Knowledge column.	105
5.8	Suggestion of extra clusters. The two PK-groups (red and blue) are validated (nodes in the consensus column are all squared). One extra clusters is proposed by the Louvain algorithm, labeled as 2. Hovering over the cluster 2, the consensus is displayed by the green diamonds. This feedback is also visible in the graph.	106
5.9	The dataset has been fully consolidated. The persons are grouped and colored by the consolidated knowledge. The user decided to assign Claude, Guillaume, Madeleine and Renexent to cluster C, by taking into account the graph and the consensus of the algorithms.	107
5.10	Computing the Lineages of VAST authors: Prior Knowledge from Alice and results of the clusterings matching it.	109
5.11	Four consolidated groups in the VAST dataset: C North, RVAC, Andrienko and London	111

List of Tables

2.1	Comparison table of most widely used visualization and analytical tool for HSNA. Visualizations: number of different visualization techniques, layout, and interactions. SNA and Models: Number of proposed SNA measures and algorithms. Clustering: Number of proposed clustering algorithms. Filtering: Possibilities of filtering according various criteria. Interaction/Direct Manipulation: Number of possible interactions mechanisms directly applicable on the visualizations.	42
3.1	Resulting networks using different models produced by one document of the examples detailed in §3.4.3: co-occurrence, unipartite and bipartite models. The first column shows the partial transcription of real documents. Colors represent annotations concerning the persons mentioned, their roles, and attributes. Underline refer to information related to the events and which can be encoded as document/event attributes. H: Husband, W: wife, T: Witness, M: Marriage, A_N : Associate, G: Guarantor, Ap: Approbator, C: Construction, F: Father, M: Mother, C: Child.	57
4.1	Tasks to support during exploration, according to our expert collaborators, split into 3 main high level tasks.	69
4.2	Comparison of the data model of several VA systems aimed at exploring bipartite social networks.	69

4 ComBiNet: Visual Query and Comparison of Bipartite Dynamic Multivariate Networks with Roles

Contents

3.1 Context	46
3.2 Related Work	47
3.2.1 History Methodology	47
3.2.2 Historian Workflows	48
3.3 Historical Social Network Analysis Workflow	48
3.3.1 Textual Sources Acquisition	49
3.3.2 Digitization	49
3.3.3 Annotation	50
3.3.4 Network Creation	50
3.3.5 Network Analysis and Visualization	51
3.4 Network modeling and analysis	51
3.4.1 Network Models	52
3.4.2 Bipartite Multivariate Dynamic Social Network	54
3.4.3 Examples	55
3.5 Applications	58
3.6 Discussion	58
3.7 Conclusion	59

In the previous chapter chapter 3, I showed that bipartite multivariate dynamic networks allow to model well historical documents, with *simplicity*, *reality* and *traceability* properties. However, no visual tools currently exist to specifically explore and manipulate this type of data. In this chapter, I propose a VA interface aimed at exploring historical documents modeled as bipartite multivariate dynamic networks, for historians to be able to reflect on their data and encoding, and to follow in depth-analysis. I answer Q2 by analyzing tasks and questions historians have on their data and providing interactions mechanisms which would allow them to answer their historical questions while finding errors in their data.

This chapter is an updated version of an article currently submitted to the journal Computer Graphics Forum, and a poster presented at the conference EuroVis 2022 [?]. It was a joint work with my advisors Christophe Prieur and Jean-Daniel Fekete. I did the development of the interface, and led in the discussions, evaluation, and writing of the paper.

4.1 Context

Social scientists such as historians aim at characterising the structure and dynamics of social groups of interest, on a region and period of time they focus on [?]. Their work essentially relies on documents—such as marriage acts, census records, surveys, and business contracts—to gather information about the life of important actors that they explore in-depth, or to draw conclusions on social aspects of groups in the society of that period and place. Instead of drawing conclusions from their gathered knowledge and interpretations of the documents, a more systematic approach consists in constructing a social network from the documents and following a network analysis approach [149]. For this, they need to encode their documents to extract the persons and any other useful information in the text and transfer it into a structured file or a database. Social scientists can then explore, validate, or refute their hypotheses by visualizing and analyzing the network structure and the connectivity patterns between the entities of the resulting network.

Currently, social scientists often model their datasets as simple networks where the nodes are the persons mentioned in the documents (see chapter 3). Usually, Two persons are then connected together in the network when they appear in shared documents. This representation is easy to visualize and analyze but simplifies and distorts the information by hiding the documents that witness the relationships between the persons. Thus, another approach consists in modeling the data as bipartite networks, where both the documents and the persons are represented as nodes and are connected together when a document mentions a given person [54, 120, 130].

In addition, historical documents include time and geospatial information corresponding to the date and location of the events they refer to, and potential additional information on the mentioned individuals such as their gender, profession, and date of birth. These are often essential to understanding underlying social phenomena, as time, space, and classical social categories play an important role in sociological structures and dynamics [80]. For these reasons, as we discussed in chapter 3, historical sources and the underlying social events they refer to can be modeled well by *bipartite* with *roles*, *multivariate* *dynamic* networks. *Bipartite* means that both persons and documents (or events, that are often witnessed by physical documents) are modeled as typed nodes. *Multivariate* means that the nodes and links can carry additional attributes. *Dynamic* means that time is a mandatory attribute of documents. Furthermore, a link created between a person's node and a document's node (when the person is mentioned in the document), has an associated link type that models the *role* of the person in the document/event. Additionally, documents can optionally carry a geographical location. This model unifies several social network models and allows to represent historical documents with simplicity, traceability, and document reality, i.e., the relationships appear as they are mentioned in the documents without distortions implied by projections [23].

More complicated models exist such as knowledge graphs [?], which are very expressive but hard to manipulate, especially for social scientists. In contrast, most visual and analytics tools widely used by social scientists such as Gephi [5] and NodeXL [133] enforce too simplistic network models and only provide limited interactions for exploring the network data, even if they provide the computation of several network measure. This results in many social historians ending their analysis by plotting the network using a node-link diagram—which is hard to read

with dense networks—, and identifying the most important actors with the help of centrality measures [81]. Lemercier & al describe this phenomena the following: “Network graphs of the “spaghetti monster” variety are a case in point. Often, in historical papers, they are used to show that a network is dense (and that the author has mastered the new technology). The narrative then comments on the individuals identified as central in the network. This approach can indeed be quite interesting, but it is hardly the only possible use for network analysis. [81]”. VA tools guiding social scientists towards more complex exploration with the help of interaction mechanisms are therefore needed.

In this chapter, I present a VA system to explore and analyze Bipartite Multivariate Dynamic Social Networks, with the aim of supporting more complex historical analysis based on easy-to-use interactions, but also potential data correction. I elaborated the tool based on four collaborations with social scientist colleagues. I first collected important questions they each had on their data and transcribed them from a network analysis perspective. The majority of the questions raised consisted in either finding specific patterns in the network—corresponding to specific groups or individual exhibiting intriguing behaviours—or in comparing several subsets of the network, in terms of network measures, attribute distributions, and their overlaps.

I hence focus on three high-level tasks: exploration, queries, and comparison of this type of network. Users can explore the data using two layouts: a node-link bipartite view showing the sociological structure of the network, and a map layout based on the geolocation of documents. I designed and implemented a new visual graph query system that allows us to build both topological and attribute constraints, based respectively on a node-link interactive representation, and dynamic widgets. By easy-to-create queries, social historians are able to 1) detect erroneous patterns and reflect on their encoding process and 2) find relevant patterns which can answer their historical questions. For this, I rely on the Neo4j graph database [100] and its query language *Cypher*. Most visualization systems offer dynamic queries to hide the complexity of query languages. However, using a rich data model, some queries are easier to refine using scripting than dynamic queries. I implemented dynamic queries that also show the translated Cypher queries, and inversely, can translate textual queries into visual queries. With that interface, social scientists can start building their queries with simple widgets and, if needed, complement them by editing the query, alone or with the help of power users. Furthermore, they can export their query, the associated results, and its history at any point to share it with someone else or to start an analysis session from a previous result. ComBiNet also implements subgraph comparison techniques, allowing the comparison of networks, network-related measures, and attribute distributions between the entities returned by the queries. I validate the query and comparison system with a formative usability study and I demonstrate ComBiNet can be used to answer sociological questions by describing in depth several real-world use cases.

After the related work section, I describe our design process in §4.3, present the system ComBiNet in §4.4 with the design of the visual query and comparison features, and present four use cases demonstrating the utility of the system in §4.5 showing it can be used to explore complex historical data and allowing historians to answer several of their questions using queries and comparisons. I finish with the results of the formative usability study in §4.6.

4.2 Related Work

As I already discussed the related work on network modeling and social network visualization in chapter 2 I only discuss in this section the related work on graphlet analysis, visual graph querying, visual graph comparison, and provenance.

4.2.1 Graphlet Analysis

One of the inspirations for this project came after participating in the 2020 VAST challenge¹ where our team used graphlets to measure similarity between several networks [139].

Graphlets are small connected induced, non-isomorphic subgraphs composing any network (see §2.3.2 for more details). They were first introduced by Milo et al. [91] to explore the structural differences between biological networks, but they are now used in several disciplines involving networks such as sociology [?].

One of the aims of the VAST 2020 challenge was to compare several multivariate networks. However, by using graphlets we realized that 1) it was not very efficient to compare several networks in contrast to other measures and 2) the interpretation of all graphlets patterns that are found in a network is complicated given the fact that one specific pattern can have various interpretations given the nodes involved and their positions in the network [68]. This is especially true that the number of potential graphlets grows exponentially as we increase the number of nodes considered (there are 6 graphlets of size 4 and 21 graphlets of size 5 for example) and if we add complexity to the network model, for example by using directed links or node and link types [118].

Instead of counting every graphlet occurrence and interpreting them with sociological meanings, it appeared more efficient to let social scientists finding specific patterns to answer questions they already ask themselves on the data.

4.2.2 Visual Graph Querying

Graph pattern matching is an important task in SNA, which consists in finding a subgraph of interest in a larger graph [?].

Several scripting languages, such as R [114] and Python [143], have been extended to support the exploration of social networks using specialized libraries such as igraph [27] and NetworkX [57], and provide functions for graph pattern matching. However, social scientists are often challenged to use scripting languages and programming, as they often do not have formal training in such technologies.

Graph databases allow to store and manipulate network data with the use of query languages, such as the Cypher language for Neo4j [100]. To lower the complexity barrier of their usage, several visual graph query systems have been developed to allow analysts to rapidly build and refine their queries visually. Some systems hide the scripting language such as

¹This is a challenge organized in the context of the IEEE Visual Analytics Science and Technology (VAST) conference. The challenge consisted of a series of analytical questions united under an overarching cyber threat scenario. We participated in the Mini-Challenge 1 which asked participants to identify a group of people that accidentally caused an internet outage. To identify this group, we were given a network profile and a large multi-variate social network to search in.

GRAPHITE [20] and VERTIGO [28] that allow specifying a graph query as a node-link diagram that the user creates interactively. Shadoan and Weaver [129] use a similar concept with hypergraphs to filter multidimensional data. Other systems, such as VIGOR [111] only visualize the query after it has been written by the user with a scripting language and do not allow direct interaction of the visual representation of the query. All these visual query systems are limited to topological queries with constraints on the vertex and edge types and do not allow to make constraints on other dimensions such as attributes and time associated with vertices and edges.

4.2.3 Visual Graph Comparison

Gleicher et al. [51] propose a taxonomy of visual comparison designs for complex objects. They claim any visual comparison system can be classified into one (or a mix) of the three following categories: juxtaposition, superposition, or explicit design. Yet, few visual systems support comparison tasks for social networks.

Andrews et al. [3] describe a technique to compare several networks, using a combination of juxtaposition and superposition techniques. The two candidate networks are shown side by side, along with a third view composed of a fusion network highlighting both the shared nodes along with the non-shared nodes with different colors. Freire et al [44] describe the ManyNets system to compare many networks by using a table where each describes one graph and each column shows graph measures in terms of small visualizations, from simple bars to distributions, allowing the comparison of a large number of graphs. However, ManyNets does not visualize the networks per se (no layout shown), and do not take into account attributes, node types, or time. Hascoët and Dragicevic [60] describe a system to match and compare graphs using superposition, focusing on the topology, not taking into account attributes or time. Tovanich et al. [139] propose a visual analytics tool to compare multivariate, sometimes bipartite, dynamic networks and find common structures. However, the tool does not handle roles and is designed for the specific task of matching a subgraph into a larger network.

4.2.4 Provenance

Provenance in the context of Visual Analytics consists in the logging of the sequence of actions of users on an interactive visualization system during analysis sessions. Collecting provenance information has proven to benefit users by providing them action recovery (undo) plus collaborative and reproducibility capabilities [115]. For example, the VisTrails system allows users to reproduce their visual analyses by providing an executable history graph of their actions, [17] while GraphTrail provides provenance tools to ease collaborative analysis [36]. Provenance can also be beneficial for visualization designers and researchers, as it gives them a tool to understand users' behaviors [7, 12] and evaluate/improve visualization systems [117]. All the reasons and concrete implementations of provenance are discussed in depth in Xu's survey [150].

4.3 Task Analysis and Design Process

I designed the ComBiNet tool in collaboration with social historians who wanted to follow a network analysis on their historical semi-structured documents, that are well modeled by

bipartite multivariate dynamic networks. I first collected questions they had about their data and what they wanted to see in a visual interface. By analyzing the questions we leveraged tasks and requirements that I used to design and implement the interface, with continuous feedback from our collaborators.

4.3.1 Use Cases

We elaborated this interface from the collaborations with historians we described in §3.4.3. These collaborations involved regular meetings and multiple discussions over two years. All these datasets are textual corpora constituted of historical documents mentioning people with complex relationships, which are well modeled with bipartite multivariate dynamic networks. We give more details about the datasets of these collaborations in this section, along our collaborators' main questions with the associated network queries to answer them. The full answers involve visualizations of the query results and attribute summaries that I describe in the next section. We categorized the questions according to four dimensions: global (G)/local (L) (do they want to categorize a group of nodes or retrieve specific persons/documents), if the question can be answered using the topology (T), and/or the attributes (A), and finally if a comparison (C) using several filters is needed or not (N).

1. Analysis of the social dynamics from **construction contracts in Italy in the 18th century (141 documents, 272 persons)** [25]. The corpus is made of contracts (manuscript documents) for different types of constructions in the Piedmont area in Italy. People are mentioned in three different roles: *Associates*, who participate in the construction; *Guarantors*, who bring financial guarantees; and *Approvers*, who vouch for the guarantors. Along with the time and location of the construction site, documents have a construction type (military, religious, and civil), work type (big work, small work, reparation, transportation, etc.), and material (wood, stone, metal). People also have an origin attribute (the place they come from), manually extracted from the original documents.

Question 1 Do approbators act as bridges compared to associates and guarantors? (G, T, C)

Query 1.1 Request all approbators occurrences

Query 1.2 Request all associates and guarantors occurrences

Question 2 Are there people mutually guarantors to each other in different contracts? (G, AT, N)

Query 2.1 Select pairs of people connected each to the two same document, with a guarantor role and any other role

Question 3 Who are the persons of the extended Zo family (G, AT, N)

Query 3.1 Request all the persons of the Zo family and their N+2 ego network

Question 4 Compare the Menafoglio and Zo families in terms of contracts and activities (G, AT, C)

Query 4.1 Request all the persons of the Menafoglio family and the documents that mention them

Query 4.2 Request all the persons of the Zo family and the documents that mention them

Question 5 Who are the persons having the 3 roles? (G, AT, N)

Query 5.1 Select persons with associate, guarantor, and approbator roles in 3 different documents

Question 6 What are the differences between Torino and Torino surroundings, concerning the types of constructions and actors involved? (G, AT, C)

Query 6.1 Request all documents located in Torino, with the persons mentioned

Query 6.2 Request all documents located in the Torino area, with the persons mentioned

2. Analysis of migrations from the genealogy of a french family between the 17th–20th centuries (2053 events, 957 persons from a private source). The corpus is made of family trees referring to several document/event types: birth and death certificates, marriage acts, military mobilization, and census report. The roles are different for each event type and consist of *children*, *father*, *mother* for the birth events, *deceased* for the death event, *spouse* and *witnesses* for the marriages, and *family members* for the census events.

Question 7 What is the trajectory of life for a given specific individual (birth, living, marriage, death) (L, A, N)

Query 7.1 Select one person and all her/his documents (to use the mentioned places)

Question 8 What is the trajectory of live for a family (L, A, N)

Query 8.1 Select birth certificates with the child, parents, and birthplace

Question 9 Where are located the main migrations, and at which time do they occur? (G, A, N)

Query 9.1 Select persons with a geolocated birth and death certificate

Question 10 Are there differences in volume and location between migrations in the 18th and 19th centuries? (G, A, C)

Query 10.1 Select persons with a geolocated birth and death certificate from the 18th century

Query 10.2 Select persons with a geolocated birth and death certificate from the 19th century

Question 11 In the Haute-Vienne and Cote d'Armor administrative areas, are there cycles in living places every 10/20 years? (G, A, N)

Query 11.1 Select persons with their census reports located in Cote d'Armor and Haute-Vienne

Question 12 In the 19th century, was there an overall decrease in the social status and professions of persons in the dataset? (G, A, C)

Query 12.1 Select persons in the first half of the 19th century with a profession mentioned

Query 12.2 Select persons in the second half of the 19th century with a profession mentioned

3. Analysis of migrations from Spain to Argentina through the marriage acts at Buenos Aires in the 17–19th centuries (1396 acts, 6731 persons) [96]. The corpus is made of acts that mention the spouses and the witnesses of the wedding, which are the roles modeled by the links. The origin, date of birth, and parents' names are specified for both spouses.

Question 13 How are spouses and witnesses linked in their family network? (G, T, N)

Query 13.1 Select marriages with spouses and witnesses, where the spouse and witnesses have the same parents

Query 13.2 Select marriages with spouses and witnesses, where the spouse and witnesses have the same grandparents

Question 14 Who are the persons with 2 marriages with a long delay? (L, A, N)

Query 14.1 Select persons in 2 marriages as husband or wife. Put a constraint on the difference of time in the marriages

Question 15 Where are the persons marrying in Buenos Aires coming from? (G, A, N)

Query 15.1 Select persons with a birth certificate not located in Buenos Aires

4. Socio-political analysis of **migration of ethnic Germans from communist Romania to West Germany in the 20th century (ongoing work)** [32]. The corpus is made of administrative forms that mention persons requesting to migrate, along with the persons they want to join, and the administrative persons of the ministry in charge of the forms (3 roles). The family members of the aspiring migrant are also mentioned in the forms, with their respective dates of birth.

Question 16 What member of their family do emigrants usually join? (G, AT, N)

Query 16.1 Select all migration documents with the emigrant and the person they are joining

Question 17 What price does the emigrant have to pay, given their socio-economic profiles? (G, A, C)

Query 17.1 Select people who are mentioned in a budget and a migration document

4.3.2 Tasks Analysis

Most of the questions we collected from our collaborators could be answered by isolating a subgroup of entities and analyzing them in the context of the whole network, or by comparing two subgraphs, in terms of their entities, structure, and attribute distributions. From discussions with our collaborators and the analysis of their questions on their data, we elaborated a list of requirements for the visual interface, split into three main parts: 1) Exploration of the data, 2) Queries, and 3) Comparisons. The elaboration of the tasks was an iterative process, as we showed the interface to our collaborators several times in the development phase to get feedback. The tasks are described here and summarized in Table 4.1:

1. **Exploration of bipartite multivariate dynamic network.** The visual interface must allow exploration of this specific type of network, using every aspect of the data, i.e. its topology (T1.1), node attributes (T1.2), roles (T1.3), geolocation of the documents/events (T1.4) and time (T1.5). Common interactions such as selection and zooming are also needed for the exploration.
2. **Applying filters.** To answer their questions, users need to be able to apply filters to the data, to isolate specific groups of entities having specific behaviors or characteristics. To answer the diversity of questions, they should be able to put constraints on every aspect of the data, i.e. the topology, the roles (T2.1), and the attributes (including time and geolocation) (T2.2). Access to provenance information can also help them in their query construction, by going to previous states and exploring different paths more easily (T2.3). Once they are satisfied with their query, they want to explore the results, usually in the context of the whole network (T2.4).

Main Tasks	Subtasks	Views	Constraints
Bipartite Graph Exploration	T1.1 Overview of the network	V1	A node-link representation is expected. The geolocation of events has to be done according to the historical period.
	T1.2 Overview of nodes attribute values and distributions	V1,V2,V4	
	T1.3 Show the persons' roles in the documents they appear in	V1	
	T1.4 Show the location of the different documents	V2	
	T1.5 Show the time of the documents	V1,V2,V4	
Apply filters to isolate subgraphs	T2.1 Filter on topological patterns	V6,V8	Constraints must be easy to set and visual.
	T2.2 Filter on attribute values	V7,V8	
	T2.3 Show the provenance of filters	V9	
	T2.4 Show the subgroups alone or in network's context	V1,V2	
Compare several subgroups	T3.1 Show the shared and exclusive entities	V1/V2	
	T3.2 Compare the node attribute distributions	V4	
	T3.3 Compare the subgraph measures	V3	

Table 4.1 – Tasks to support during exploration, according to our expert collaborators, split into 3 main high level tasks.

	Bipartite	Node Attributes	Links Attributes	Dynamic	Geolocated
Jigsaw	✓	Only some	✗	✓	✓
Puck	✓	✗	✗	✓	✗
ComBiNet	✓	✓	Encode roles	✓	✓

Table 4.2 – Comparison of the data model of several VA systems aimed at exploring bipartite social networks.

3. **Comparison of several subgraphs.** Users should be able to compare several subgraphs isolated after applying filters, to see the similarities and differences between groups of entities of interest. The system should be able to easily see the common and shared entities of the two subgraphs (T3.1), their respective place in the network, their structural differences (T3.2), and their different attribute distributions (T3.3).

4.4 The ComBiNet System

ComBiNet is designed to visualize, explore, and analyze social networks encoded as bipartite multivariate dynamic network. Some other systems exist to explore bipartite social networks such as Jigsaw and Puck, but do not encode every aspects of historical documents historians are interested in. Table 4.2 shows a comparison of their data model compared to ComBiNet.

When started, ComBiNet dynamically collects the node types, roles, sub-types, and attributes when reading the network from the database. The interface is constituted of four main panels, split into different views as shown in Figure 4.1: the query and comparison panel (V6,

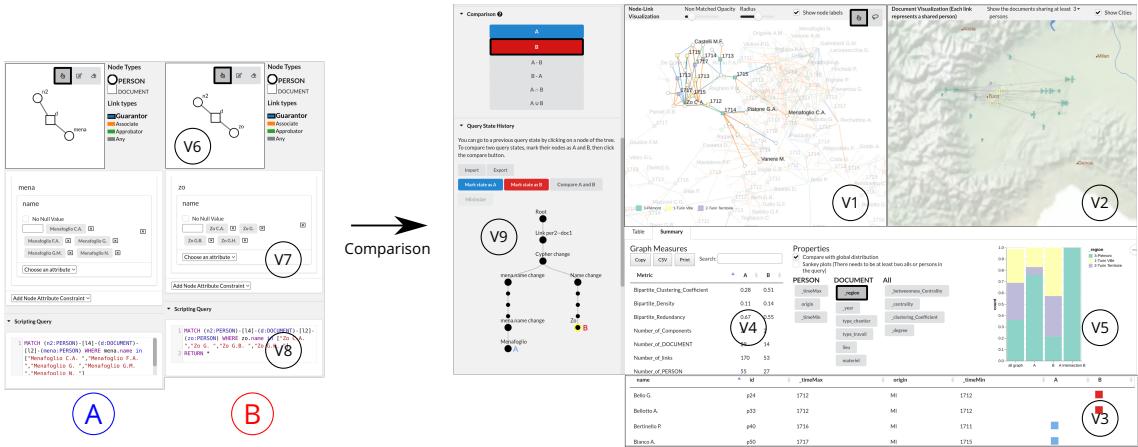


Figure 4.1 – The ComBiNet system used to compare two subgroups of a social network of contracts from [25], extracted with dynamic visual queries. (A) and (B) show the two visual queries created by the user in the query panel using an interactive node-link diagram editor (V6), dynamic query widgets (V7), and the equivalent Cypher script (V8). The right part shows ComBiNet’s global interface in *comparison* mode: (V1) Network visualization panel, (V2) Map of the geolocalized nodes, (V3) Table of persons, (V4) Graph measures comparison, (V5) Attribute distribution plots, and (V9) Provenance tree. The two visual queries on the left, translated into Cypher queries below, select the “Menafooglio” family on the left, and the “Zo” family on the right, along with their construction contracts and close collaborators.

V7, V8, and V9), the bipartite visualization panel (V1), the map visualization panel (V2), and the query results panel (V3, V4, V5). I present in the following the different views, according their visualization or query functions. Comparison features are incorporated in the same views with different comparison mechanisms.

4.4.1 Visualizations

ComBiNet presents a social network with multiple visualizations and views highlighting different aspects of the data. The views are linked when it makes sense so that interactions such as selection done on one propagate to other views.

V1: Bipartite Node-Link View The bipartite node-link visualization panel shows the network using the DrL force layout from igraph [27] with overlap removal using D3 [13]. Node-link representations are very common in social sciences [5, 24, 98, 133] and were a specific request from our collaborators. In the context of our bipartite model, the persons are represented as circles and the documents/events as squares, while the roles are encoded as link colors. A link models the mention of a person in a document. This view provides an overview of the data by showing the structure of the network (T1.1) and the roles of the persons in their different documents (T1.2). The view also provides pan & zoom and selection interactions for effective navigation. Nodes’ labels are displayed (name of person and ids of documents by default) on

the canvas with an occlusion-free mechanism which hide nodes with low degree when two or more nodes labels overlap.

V2: Map View The map visualization panel on the right shows an event-centric view, displaying only the geolocalized event nodes on a map. By default, only event nodes are shown, but users can select a threshold to show links between nodes when they share at least a given number of persons in their mentions. Persons are not directly shown in this view as they do not have a unique location. This map view presents a transformation of the bipartite network, focused on the geospatial information that is very important to social scientists (T1.3).

As we collaborate with historians who study different periods, we cannot use modern map backgrounds such as the default one provided by OpenStreetMap or Google Maps since many features are anachronistic (e.g., roads, administrative areas, borders). We, therefore, provide a map background with only these non-administrative features: elevation, lakes, rivers, and types of environment. We also show the most important cities as most of them existed in the past and provide landmarks. The map uses Natural Earth tiles and vector data [99].

The map has the same interactions mechanisms than the bipartite node-link view. The two views are also coordinated: selecting/hovering an event node in the graph view highlights it on the map and vice versa, while hovering a person node highlights all its corresponding documents on the map, rapidly showing the person's events' locations.

V3: Entities Tables All the persons and the documents of the loaded dataset are listed in two separate tables, showing the attributes of the entities (person or document). This way, users can order the entities according to any attribute they want (T1.2). The tables are linked to the visualizations, meaning that selecting a row highlights the respective entity in the visualizations and vice-versa. Selecting a node hence highlights the corresponding row and pushes it at the top of the table. Tables in social network visualization systems have been proven to be efficient and useful for social scientists when exploring their data [9] and is a feature that have been asked by our collaborators. It allows them to link the visualization to the network entities more easily, and dive deeper into one entity's attribute values after selecting it in the network. For example, if the visualization reveal an intriguing person connected to two distant components through two documents, the user can rapidly see the information of this persons and documents on the tables, to see if this could be an error from the annotations or an interesting person he or she could investigate more in depth in the original sources. It also makes ranking entities according to various criteria easier and more straightforward. Finally, the tables are exportable in csv, pdf, or directly in the clipboard, which was a request to our collaborators.

V4: Graph Measures The Graph Measures view shows measures related to the network and gives insights into its structure to users (T1.1). We report simple measures like the number of persons, documents, links, and components, and more sophisticated bipartite network measures asked by our users, that they can report for their analysis: the bipartite density, bipartite average clustering coefficient, and bipartite average redundancy [?]. These measures are updated in real-time when filters and comparisons are applied.

V5: Attributes View All the attributes in the network are shown as buttons in the bottom right of the interface, sorted by their associated node type (person, document, and "All" for both types). They can be quickly visualized by hovering over the button, producing two effects:

it colors all the nodes on the two views according to their attribute values, and it shows a plot of the distribution of the selected attribute. Figure 4.2 shows the construction dataset of collaboration #1 where the user selected the `_year` attribute, coloring the documents nodes with their year in the node-link diagram (left) and the map view (right), revealing for example that most construction occurring in 1714 occurred in Torino and Torino close surrounding. By clicking on the button, the visual encoding and the distribution plot remain selected. This interaction is inspired by the x-ray technique of the Vizster system [61]. Users can follow a first exploration of their data by visually detecting correlations between attribute values and some groups of persons or between attribute values and some specific areas in the map view (T1.2, T1.4, T1.5).

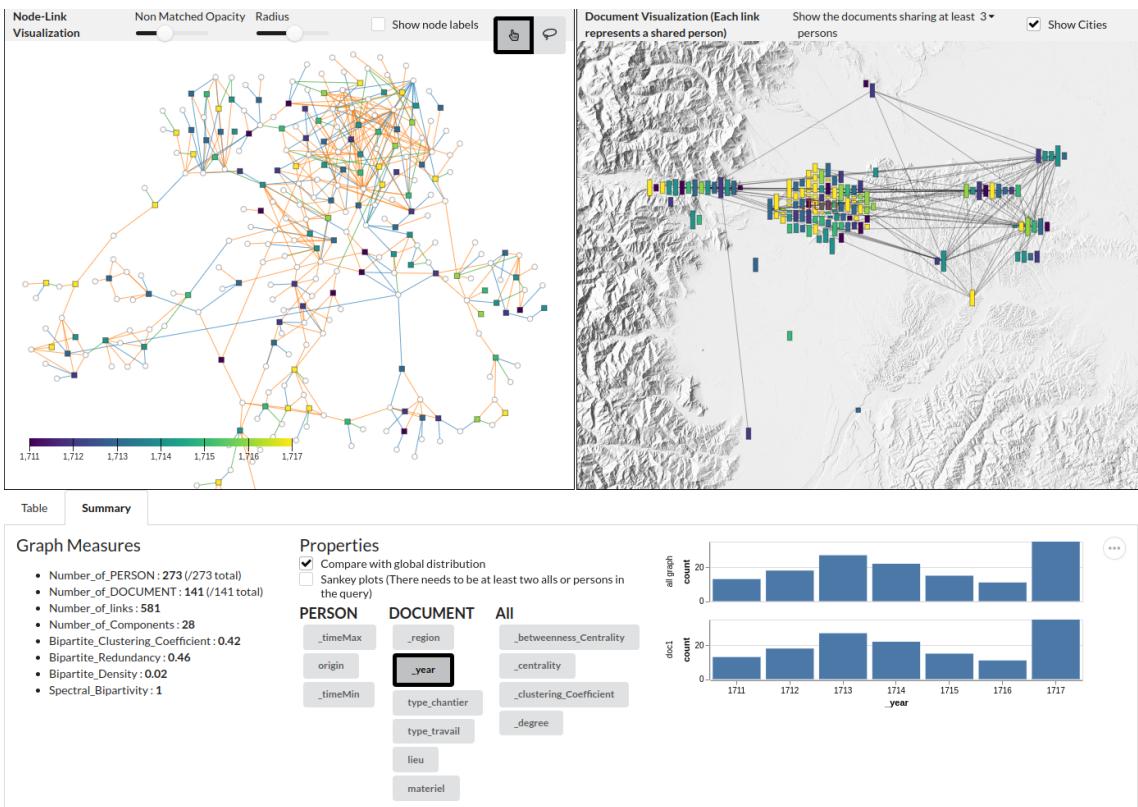


Figure 4.2 – ComBiNet interface with the dataset of collaboration #1. The user selected the `_year` attribute, showing the distribution of document years with a histogram (bottom right), and coloring the documents node on the bipartite view (left) and map view simultaneously (right).

4.4.2 Query Panel

The query panel allows users to rapidly build queries visually, with topological and attribute constraints. The visualization of the query is synchronized with the Cypher query sent to the database. Modifying one representation update the other, allowing users to build a query visually and refine it in Cypher when appropriate. Experts users who know the Cypher language can

also start to construct their query textually and modify it visually later on. In this subsection, I describe all the features and interactions allowing ComBiNet to build a query and illustrate them with questions 2 and 6 of our collaboration #1. Our collaborator wants to *find the persons who are mutually Guarantor to each other in separate contracts* (2) and to know *how Torino and Torino's surroundings differ according to their contracts?* (6). Figure 4.4 (left) shows the final queries, but first, I explain how to create them.

V6: Node-Link Dynamic Query

The interactive node-link diagram allows the construction of a subgraph query graphically, that represents a topological constraint (T2.1). The query subgraph is built and edited interactively. At each modification, the visual query is converted into a Cypher query and run in the database which return the results. All the matches are displayed in the entities tables (V3) and highlighted in the main visualizations views (V1, V2). Three modes of interaction are available through the top-right menu: *selection*, *addition*, and *deletion*. The *selection* mode allows to drag the nodes in the panel, while the *addition* and *deletion* modes allow the following actions:

Node Creation: In *addition* mode, clicking on an empty area creates a new node. The node will be of the selected type from the legend on the right (Person or Document).

Node Deletion: In *deletion* mode, clicking on a node deletes it and its links.

Change Node type: In *selection* mode, clicking on a node opens a menu allowing to change its type.

Link Creation: In *addition* mode, clicking on a node and dragging the mouse to another node will connect the two with a link. Its type (color) will be the link type selected on the legend.

Link Deletion: In *deletion* mode, clicking on a link deletes it.

Change link type: In *selection* mode, clicking on a link opens a menu to change its type.

Users build concrete subgraphs with the same representation as in the bipartite network view: a visual query is a network template with additional attribute constraints. Each role (link type) is rendered using a color (Figure 4.3 left). Users can also create untyped links using the *Any* value, which will match all the existing link types (Figure 4.3 left). Created links can be matched by different selected link types, by checking several possible types for one link. These links are represented by a dashed line with the colors of the possible types (Figure 4.3 middle right). Several links with different types can also be created among two nodes to query a person with more than one role in the same event (Figure 4.3 right). When a node or link is created in the query, it is given an identifier starting with *per* for a person, *doc* for a document, *link* for a link, followed by a number. These identifiers are used in the attribute constraint panels and the textual query and can be changed through their textual representations.

To find persons who are mutually guarantors in our collaboration #1, we first create one person and two documents using the addition mode and by clicking on the canvas. We then link the person node to the first document with a link that is not typed (Figure 4.3 left),

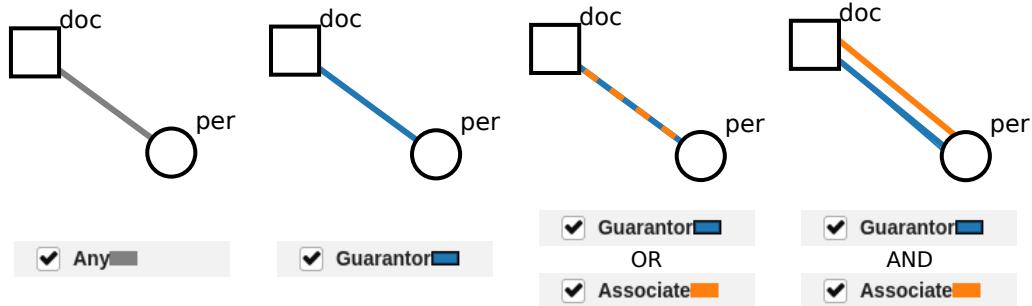


Figure 4.3 – All link creation possibilities: Any link type (left), one selected link type, here guarantor (middle left), the union of several link types (middle right), several links with different types (right)

and link it to the second document with a *guarantor* link (Figure 4.3 middle left). We then create a second person node and link it to the two documents with opposite link types. The resulting visual query is presented in Figure 4.4 (a). To answer the question 6 (comparing Torino and Torino surroundings), we start to request all the links in the graph, no matter the type, as shown in Figure 4.4 (b). The database then returns all the links in the graph with their attached nodes. Putting attribute constraints on the location of the contracts will then let us answer the question.

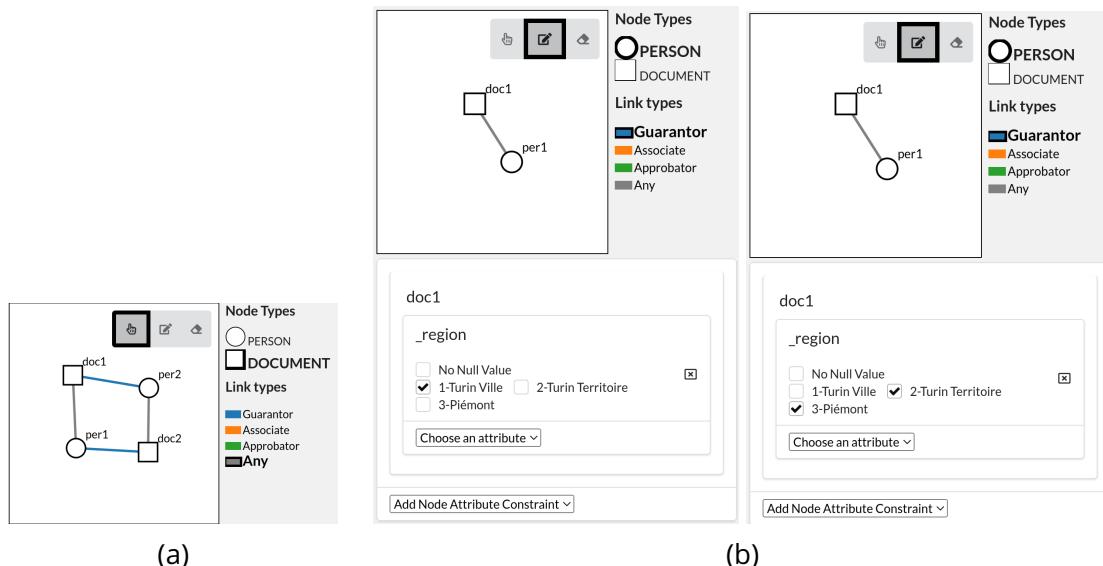


Figure 4.4 – Visual queries created to answer questions 2 and 6 of our collaboration #1. (a) The visual query retrieves individuals who are mutually guarantor to each other in separate construction contracts. (b) The two visual queries retrieve the documents—along with the signatories—of Torino (*Turin* in french) (left) and of Torino surroundings (*Turin Territoire* and *Piemont*) (right)

V7: Attribute Constraint Widgets Users can also add attribute constraints (T2.2) on the created nodes with the help of interactive widgets. An input button is created for each node and link identifier from the node-link query panel. It allows to create a dynamic query widget for any of its attributes. The widget design vary according to the three possible attribute types: numeric, categorical, and nominal, as in the original dynamic queries formalization by Shneiderman [131]:

1. **Numeric constraints** are modeled as range sliders, allowing the selection of lower and upper bounds to the filter.
2. **Categorical constraints** are modeled as a set of checkboxes. Each possible value has a corresponding checkbox.
3. **Nominal constraints** are modeled as text input, where the user can write any desired value. All the possible values are shown at the same time and filtered as the user writes.

For the categorical and nominal widgets, selecting several values corresponds to the union of the filters. The three widget types are shown in Figure 4.5.

The figure displays three attribute constraint widgets:

- region**: A categorical filter with checkboxes. It includes options for "No Null Value", "1-Turin Ville" (which is checked), "2-Turin Territoire", and "3-Piémont".
- Entity**: A nominal filter with a text input field containing "No Null Value".
- id**: A numeric filter with a double slider set between 1712 and 1714.

Figure 4.5 – Widget designs for the different attribute types: checkboxes for categorical attributes (top), text input for nominal attributes (middle), and a double slider for numerical attributes (bottom). The categorical attribute example shows the inputs letting users create new constraints for other attributes and other nodes.

To answer our collaborator's second question (*how do Torino and Torino's surroundings differ according to their contracts?*), we first filter the documents which are located in Torino

(*Turin*). For this, we select the whole dataset by linking a person and document node with an untyped link. Then, we select the id *doc1* of the document of our visual node-link query, and the *region* attribute. It initializes a categorical widget including all the values found in the dataset for this attribute with associated checkboxes. We check the region of interest “1-*Turin Ville*” to select all the documents from this region. The first widget of Figure 4.5 illustrates the created constraint. To select the documents of Torino’s surroundings, we can simply uncheck the “1-*Turin Ville*” value for the *region* attribute and check the two other values “2-*Turin Territoire*” and “3-*Piemont*” which are areas corresponding to the surroundings of Torino. Both queries are represented in Figure 4.4 (b).

V8: Cypher Editor Users can build or modify a query using the Cypher query language, with the Cypher text editor. This allows users to start creating a query visually and refining it by text for complex constraints which can not be represented by a visual form easily. The parts of the Cypher query which are not visually expressible appear in Cypher widgets next to the other widgets. The editor supports autocomplete to e.g., help to discover and spell the attribute names. The visual and textual representations are synchronized, meaning that modifying one update the other and return the results in the visualizations, tables, and attribute distributions.

Query Results

Each modification of the query, whether from the node-link dynamic query, the widgets, or the Cypher text box, update the two visualization panels (V1, V2), the entities tables (V3), the network measures view (V5), and the attribute plots (V6). The nodes and links that do not match (are not retrieved by the query) are grayed out in V1 and V2 and are removed from the persons and documents tables (V3). A third table shows every occurrence found of the created pattern that we call the occurrence table. The occurrence table for question 1 of collaboration #1 is shown in Figure 4.6 (a). It tells us that the occurrence has been found 72 times, meaning that the pattern exists 36 times in the network by taking into account the symmetry of the subgraph query. Users can switch between the three tables in the table view using the tabs. The network measures are computed on the new subgraph formed by the union of all patterns found and updated on the network measures view (V5). Figure 4.6 (b) (left) shows to the user the different graph measures of the subgraph induced by the patterns found. Since some measures can be long to compute, the values are computed iteratively in the backend and shown progressively [41] to avoid blocking the interface. The distribution plots in the attributes view (V6) are updated, showing the values of the entities of the latest constructed query, next to the global distributions.

Attributes Visualization When users select an attribute in the attributes view (V5), its distribution is visualized for the queried entities and the whole network with an histogram. However, these plots show the aggregated values and we lose the potential value transitions between the query nodes. For example, Figure 4.7 shows a query to list the persons with the role of “approbator” (green) in a contract after being a “guarantor” (blue) in another contract (using a time constraint). We may want to see if the locations or types of the two contracts are the same or if they change, case by case. Unfortunately, we lose this information with the aggregated plots. By checking the “Sankey” option on top of the distribution visualization, the plots are transformed into Sankey diagrams, giving information on how the attribute values

Persons Documents Occurrences

CSV Copy Unselect Resize Search:

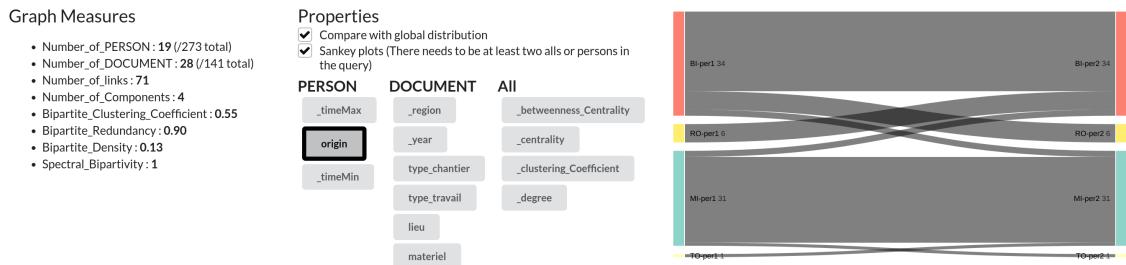
Option	doc1	doc2	I	I3	I5	I6	per1	per2
[+]	DOCUMENT (a25)	DOCUMENT (a4)	Guarantor	Associate	Guarantor	Guarantor	Bello P. (p27)	Gillonio G.P. (p175)
[+]	DOCUMENT (a5)	DOCUMENT (a4)	Guarantor	Guarantor	Guarantor	Guarantor	Bello P. (p27)	Gillonio G.P. (p175)
[+]	DOCUMENT (a5)	DOCUMENT (a4)	Guarantor	Guarantor	Guarantor	Guarantor	Gillonio G.P. (p175)	Bello P. (p27)
[+]	DOCUMENT (a25)	DOCUMENT (a5)	Guarantor	Associate	Guarantor	Guarantor	Bello P. (p27)	Gillonio G.P. (p175)
[+]	DOCUMENT (a4)	DOCUMENT (a5)	Guarantor	Guarantor	Guarantor	Guarantor	Bello P. (p27)	Gillonio G.P. (p175)
[+]	DOCUMENT (a4)	DOCUMENT (a5)	Guarantor	Guarantor	Guarantor	Guarantor	Gillonio G.P. (p175)	Bello P. (p27)

Showing 1 to 10 of 72 entries

Order by ▾

Previous 1 2 3 4 5 ... 8 Next

(a)



(b)

Figure 4.6 – Results of question 2 of collaboration #1: (a) shows a subset of the table view with every occurrence of the pattern found. (b) shows the summary panel, with the graph measures and the attributes view with the *origin* attribute selected and the Sankey option checked. It allows us to see the attribute distribution of the persons included in the pattern and see if there is a relationship between persons who are mutually guarantors and their origin.

relate between the nodes (person or event) of the same query. A Sankey diagram showing the attribute distributions is particularly useful for queries where nodes have a relationship in regard to time, such as birth certificates, marriage, or death certificates where we know the order in which these events occurred. It is also useful for queries with user-defined time order constraints as in Figure 4.7.

The attributes plots are exportable in SVG, while the tables are exportable in csv, pdf, or directly in the clipboard. This was a demand of our collaborators, who can share the results of their queries easily or use them in another software or tool, according the *traceability* principle.

The graph measures and attribute views for the results of question 2 of collaboration #1 are shown in Figure 4.6. The Sankey view of the *origin* attribute shows that mutual guarantors come from 4 regions only and that usually, people have mutual guarantor relationships only with persons of the same origin. This is especially true for persons from *Milano*, and with some reciprocal links between persons from *Bioglio* and the *Comune di Ro*.

V9: Provenance Tree Each change in the query panel is saved with the computed results so that the history of the query construction can be shown in the form of a provenance tree

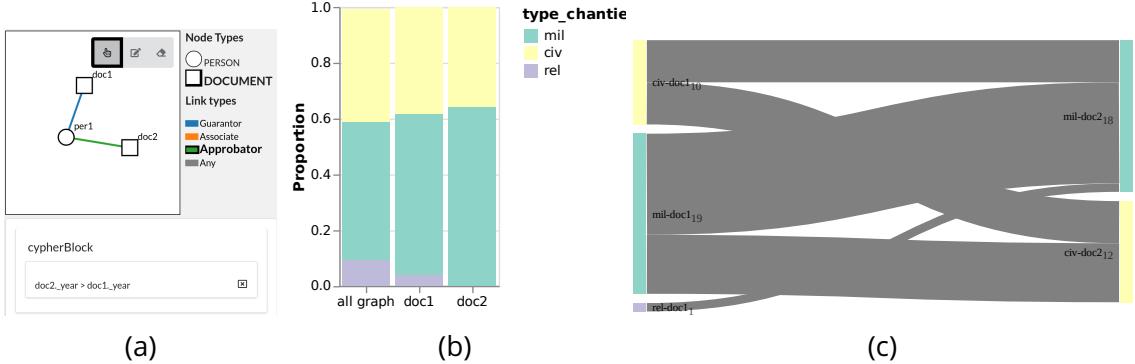


Figure 4.7 – Two ways of showing the distribution of “type chantier” (type of works), a categorical attribute with three possible values “religious”, “military”, and “civilian”. (a) A query matching the contracts made by the same person (*per1*) as an “approbator” (green link to *doc2*) after being a “guarantor” (blue link to *doc1*) using the constraint (*doc2._year > doc1._year*). (b) Stacked bar chart for the matches, the earlier contract (*doc1*), the older contract (*doc2*), and (c) Sankey diagram with the early values on the left and the last on the right. The Sankey diagram reveals the value changes between the two documents: the guarantor who worked initially on religious work switched to military work.

(T2.4), managed with the Ttrack library [29]. Each node of the tree represents a query change, with a description label such as “New Link”. It allows to rapidly visualize the succession of filters applied with their refinements. At any moment, users can rename a tree node or click on it to go back to the previous state; allowing them to explore different query possibilities easily and iteratively. Hovering over a node shows a tooltip with the query panel associated with the selected query state. It let users rapidly see what query is associated with each node of the tree. If a new change is made on the query from a previous state, a new branch is created on the tree, allowing to revisit and refine explorations. Figure 4.8 shows the provenance tree made to answer question 2, split into 2 branches, with the tooltip showing one of the node query state. The whole provenance tree is exportable and importable in json format, allowing to 1) start a session from a previous exploration and not from scratch, 2) share exploration sessions and results with others, and 3) provide a trace of the exploration leading to a potential interesting result, hence providing *traceability* in the results.

4.4.3 Comparison

In addition to comparing the results of a query to the whole graph, ComBiNet allows comparing the results of two queries. Users can select two query states in the provenance tree and mark them either as “A” or “B”. Clicking on the button “Compare State A and B” compares them. The interface changes to *comparison mode*. Several buttons appear on top of the provenance tree: *A*, *B*, *A – B*, *B – A*, *A ∩ B*, and *A ∪ B* for exploring the combinations of the two results of A and B in the two visualizations panels.

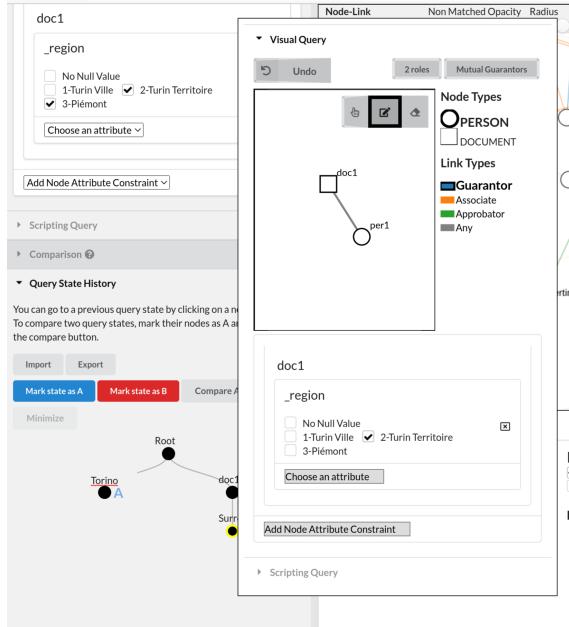


Figure 4.8 – Provenance tree to answer question 2 of collaboration #1: left branch leads to Torino documents (the node is labeled as A) while right branch leads to surrounding documents (the node is labeled as B). The user hovers over one node, revealing a tooltip that shows the visualization of the node's query..

To answer several of the questions raised by our collaborators, we need to compare two subsets of the network.

In the question 6 of collaboration #1, we want to compare the constructions in Torino with the ones in Torino surrounding. Since we previously constructed the query returning all the contracts from Torino (*Turin*) with the mentioned people, we can return to this point in the provenance tree, and change the constraint of the *region* attribute from *Turin* to *Turin Territoire* and *Piemont* using the checkboxes to get the documents of Torino surroundings in a second query. Both queries are shown in Figure 4.4 (b). The user can then rename the provenance tree nodes with explicit names such as “Torino” and “Surroundings”, and mark them as A and B using the appropriate buttons. Clicking on the “Compare State A and B” will make the interface compare the two query results.

Topological Comparison In comparison mode, users can rapidly switch between the visual filters of (A) and (B) by hovering over their respective buttons on the comparison menu and thus compare the structure of the two resulting subgraphs (T3.1). Similarly, different boolean comparison operations are available by hovering their respective buttons (Figure 4.1-C), such as the intersection, union, and differences between the two filters. Moreover, the summary tab allows comparing the different graph measures of the two subgraphs by showing them side by side (T3.3). Figure 4.9 shows the comparison table for the queries returning the subgraph of Torino (A) and Torino surroundings (B). Comparing these measures, such as the number of matched documents or the densities, is crucial for SNA. For example, the table indicates that

the density is two times higher for Torino, suggesting that less persons participate in the same construction compared to Torino surroundings.

Metric	A	B
Bipartite_Clustering_Coefficient	0.52	0.42
Bipartite_Density	0.04	0.02
Bipartite_Redundancy	0.45	0.46
Number_of_Components	13	22
Number_of_DOCUMENT	42	97
Number_of_links	153	419
Number_of_PERSON	99	214
Spectral_Bipartivity	1	1

Showing 1 to 8 of 8 entries Previous 1 Next

Figure 4.9 – Comparison table of the networks measures for Torino subgraph (A) and Torino surroundings subgraph (B).

Attribute-Based Comparison The comparison of one or several attribute distributions between (A) and (B) is also useful for answering the historical questions of our users. In the attribute view (V5) of the results panel, hovering or clicking on an attribute name will show the distribution of this attribute in four contexts: the nodes of the whole graph, the queries (A), (B), and the currently selected Boolean operator (e.g., intersection or union) if one is selected. This allows users to compare attribute distributions between several subsets of interest (T3.2). For example, we can compare the attributes between the contracts of Torino and the ones of its surroundings. We can also compare the persons who worked in Torino, in Torino's close territory, and in both areas, by selecting the intersection operator. Figure 4.10 illustrates the comparison plots for different attributes. The first plot indicates that the types of construction sites differ between the two regions: the city of Torino clearly has a lot of military sites compared to the surroundings of Torino, which has almost none. This is the opposite for the number of religious sites, which are almost all localized in the surroundings of Torino. If we now look at the year distribution of the contracts, we can see a difference in the distributions. The years of Torino's construction contracts were steady between 1711 and 1717 with a little spike in 1713, while the constructions were more scarce in the surroundings before 1716. We can see a big spike in construction in 1717. This is interesting to our users, as it shows the dynamic of the construction in the area: the center of the city started to be constructed before other constructions arose in the surroundings.

We can also compare the profiles of persons who collaborated at Torino and Torino surroundings by selecting the intersection of those two queries. One of the questions the historian had (question 2 of Table 4.1) was to know if those persons were a group with specific attributes and characteristics, or were inseparable from other persons working in the two areas. If we look at the betweenness centrality, on average, the values are higher for this group of people, meaning that the persons who work on the construction site at Torino and Torino's

territory are clearly two distinct groups, and the persons collaborating in the two areas act as bridges between these groups. This visual demonstration was convincing and revealing for our users.

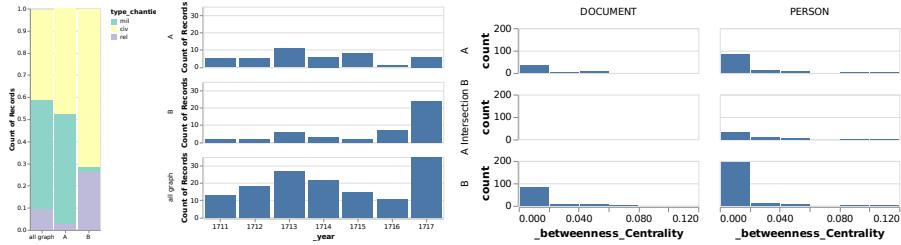


Figure 4.10 – Distribution of the type of constructions, the years, and the betweenness centrality for the documents and signatories of Torino (A), Torino surroundings (B), and the whole graph (top).

4.4.4 Implementation

ComBiNet is made of three components: a web visual interface, a python server, and a Neo4j [100] graph database instance. The client interface is written in JavaScript using D3 [13], Vega [125], and the Ttrack library [29]. The python server is written in Flask and interacts with the Neo4j instance for query processing before sending the results to the frontend. We implemented the Cypher parser with the ANTLR parser generator [107]. The abstract syntax tree of the Cypher query is used as a representation of the query. Modifying the query visually update the tree, which is translated into Cypher in the textual query panel. Similarly, a manual change in the Cypher query update the abstract syntax tree which is translated into a visual query.

4.5 Use Cases

In this section, I describe how our system has been able to specifically answer questions from three of our collaborators and one other use case. The tool was mostly operated by the developers working side by side with the collaborators to test the expressiveness of the queries and the value of the results visualizations. The tool was refined as needed along the way.

4.5.1 Construction sites in Piedmont (#1)

One of the main questions of our collaborator was to compare two families which he knew played a big role in the structure of the network: the *Menafoglio* and *Zo* families (question 4 in Table 4.1). Specifically, he was interested in knowing if there were differences in specialization in type of contracts and area of work for the core members of these families, and to what extent the two families were collaborating. Moreover, he was very interested in characterizing the group of people collaborating with both families.

To answer those questions, we first selected the core members of the *Menafoglio family*, by checking the people known by the historian, and their close neighbors. Looking at the bipartite

view (see Figure 4.11 (a), we can see that the group is pretty dense with people collaborating a lot between them. Looking at the map, we can clearly see that the family has been mostly active in Piedmont outside of Torino and Torino's close territory. We also have a first view of the attribute distribution of the persons in the group and their contracts.

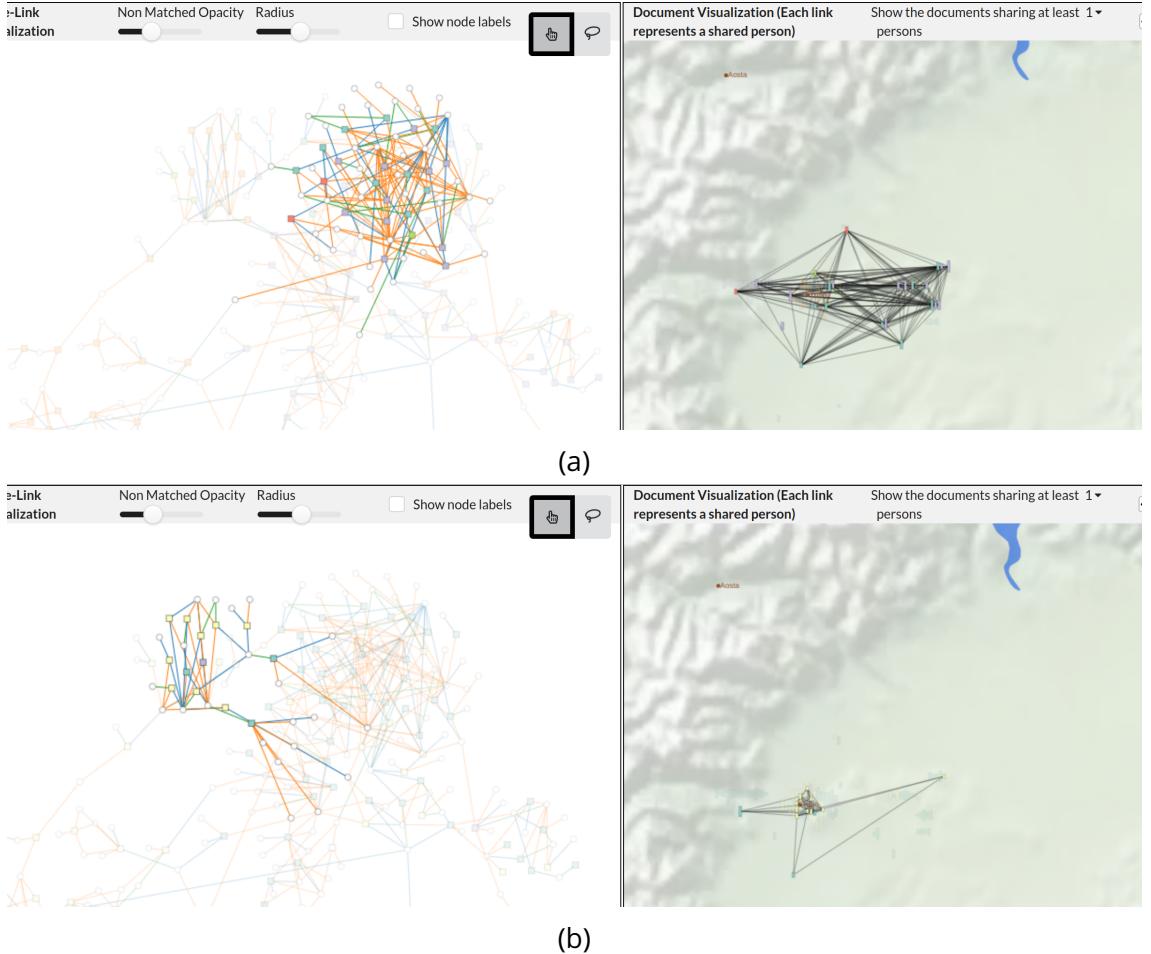


Figure 4.11 – Menafoglio (a) and Zo (b) families retrieved with queries and highlighted in the bipartite node-link and map views.

We then do the same query for the Zo family. We keep the same topological filter and replace the name filters with the core members of the Zo family known by the historian. We see on the graph view (Figure 2 of the supplementary material) that the group is smaller and is in a different area in the graph. The map enriched with a selection of the *region* attribute shows that, contrary to the Menafoglio, the Zo family has been more active in Torino and Torino territory (very close area of the city).

The two groups can be compared using the *comparison mode* by selecting the two queries in the provenance tree. This opens the comparison menu to quickly navigate between the visual selection of (A), (B), and the set $A \cap B$ that interests our collaborator. The table showing the

graph measures of the two subsets confirms what is shown visually: the Menafoglio group is more populated but less dense than the Zo family.

Our user is then interested in comparing the distribution of several attributes between the two groups. We can clearly see in Figure 4.12 (top right) that the Menafoglio family is more specialized in military (*mil*) sites, while the Zo family is doing more civil (*civ*) constructions. This is confirmed by the *material* distribution that shows that the contracts of the Menafoglio are often using stones, whereas it is never the case for Zo contracts. Finally, the persons collaborating in the two groups have a betweenness centrality higher on average (bottom right, middle chart). This makes sense as they act as bridges linking the two families.

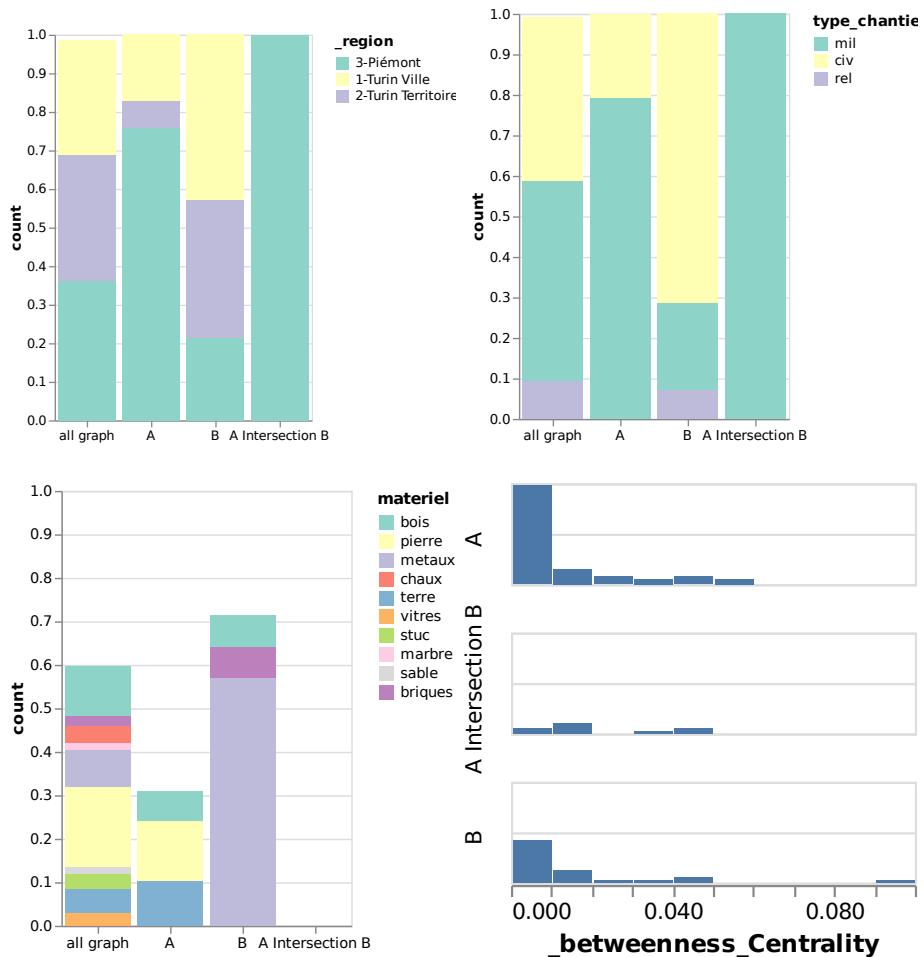


Figure 4.12 – Attributes distributions plots between the whole graph, the Menafoglio family (A), the Zo family (B), and $A \cap B$, for the *region*, *type_chantier*, *material* type.

4.5.2 French Genealogy (#2)

We describe how ComBiNet allowed us to answer an important question of the use case #2: to detect the largest migrations across several generations, in which areas, and at what

time they occurred (question 7 in Table 4.1). The map view shows at a glance (Figure 3 in the supplementary material) that the majority of events have taken place in three specific regions west, mid-north, and mid-south.

To find patterns of migrations within families, we first make a query representing a simple family by linking a person node to a birth event, connected to the parents using a link of *father* or *mother* type. We repeat the process to the new parent node to add another generation. Finally, we connect the latest generation child with a death event, to have another date and location to compare to (see Figure 4.13a). This query returns every person with their parents and grandparents, along with their respective birth and death data for the latest person. We also create a constraint on the *department* attribute on the documents to only retrieve the events that have a non-null associated location. This request returns a subgraph of 64 persons and 88 documents. The user can now select the *department* attribute to create a Sankey diagram that shows the change of departments across the different generations of the families. Figure 4.13b shows that the majority of families are from *Haute-Vienne* (which can easily be confirmed by checking the map), and do not move much across generations. Our collaborator however detected interesting patterns of people moving from the department *Creuse* to *Haute-Vienne* across two generations. She refined the query by adding an attribute filter on this specific department using a widget. The table view then showed her who these migrants were and when it occurred. The bipartite visualization panel allowed exploring more in-depth this specific group of people.

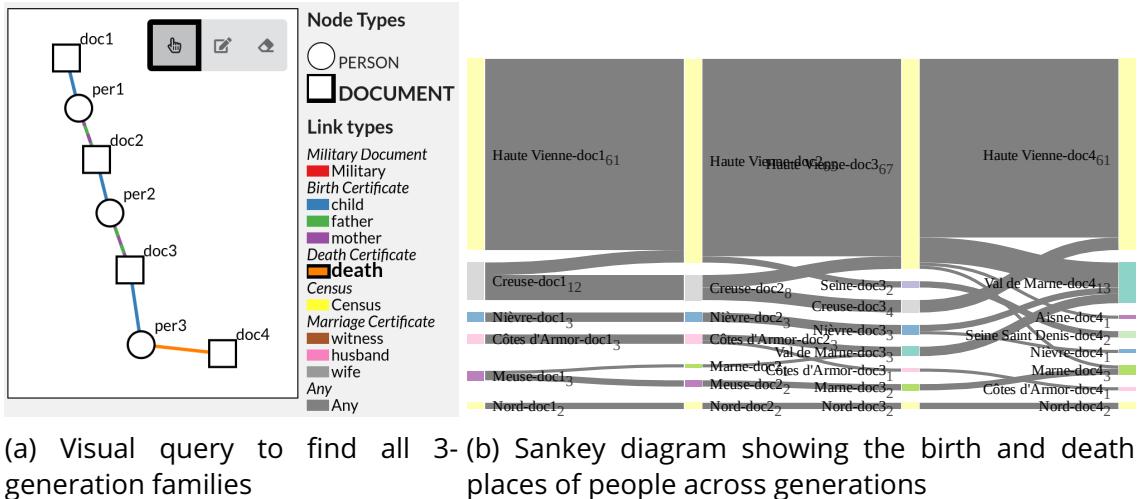


Figure 4.13 – Migrations across departments over three generations

Afterward, we answered question 8 (Table 4.1), to compare the migrations between the 18th and 19th centuries. She thought people started moving in the 19th century and wanted to confirm it. To answer this, we first created a query to retrieve the people with birth and death certificates from a specified department. We then applied a time filter on the death certificate node, first for the 18th century and then the 19th century, compared the two query results using the comparison mode, and looked side by side the Sankey graphs related to *departments*

(Figure 4.14). We can clearly see that people do not move at all in the 18th century, while in the 19th century even if the majority of people stay in the same place from their birth to their death, more than half moved.

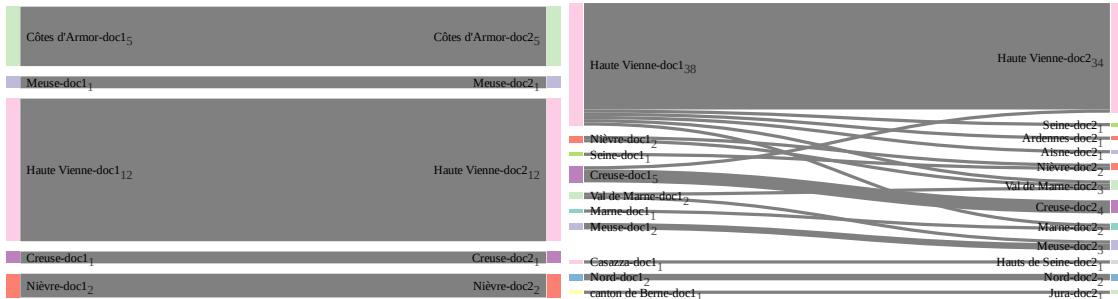


Figure 4.14 – Sankey diagrams showing the migration of people in the 18th and 19th centuries, extracted from their birth and death places.

4.5.3 Marriage acts in Buenos Aires (#3)

I present in this subsection how ComBiNet has been used to detect erroneous encoding during the annotation process of the marriage acts. The documents mention 6659 individuals, whom can have the same name, especially between fathers and sons in this period and region as specified by our collaborator. During the annotation process, the historian and his collaborators gave identifiers to the persons mentioned in the documents—which is the common annotation procedure. However, for this case of homonyms, it can be hard to know if some mentions between different documents refer to the same or different persons. Historians cross the information contained in the different documents to disambiguate the persons [?], but errors can easily be made, i.e., giving the same identifier to different persons or giving different identifiers to the same person. I used ComBiNet in collaboration with researchers of this project to detect erroneous patterns and help them refine their encoded data. For this, we can find the persons mentioned in two acts either as *husband*, *wife*, or *witness* with a time difference of 70 years or more. Such person nodes in the network are with almost full certainty representing two different persons who lived in different generations. We constructed a request to find this pattern with the visual query view and added the time constraint between the two marriage acts with the Cypher textual input. Figure 4.15 shows the visual query constructed (left), the bipartite view with the persons and documents matching the query highlighted (middle), and the table listing all the documents having mentions of people with erroneous identifiers (bottom right). The table permit to browse through all persons nodes (29 have been found) who correspond in fact to more than one person and to the documents which contain the wrongly given identifiers. Using the exporting capabilities, our collaborator have exported the occurrence table indicating him with their identifiers the pairs of documents mentioning two different persons who have been given the same identifier. Using this table, he has been able to rapidly correct those errors in his annotation framework.

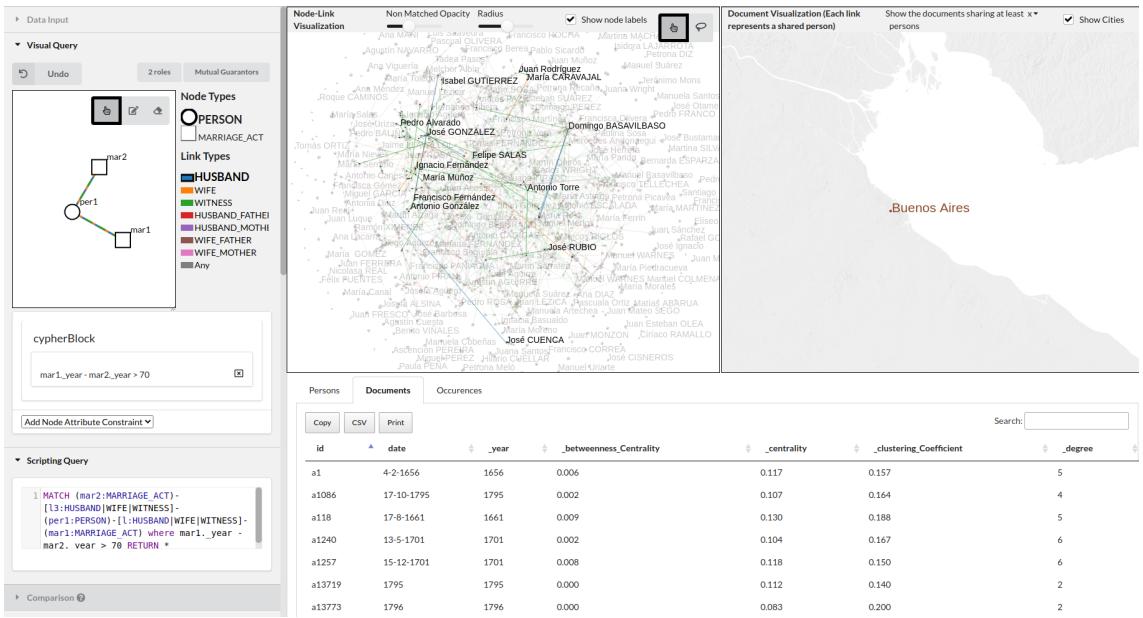


Figure 4.15 – ComBiNet used to request persons appearing as husband, wife, or witness in two marriages which occurred at 70 years apart or more.

4.5.4 Sociology thesis in France

I describe in this third use case how ComBiNet can be used to answer questions about thesis in France between 2016 and 2022. Indeed, some sociological datasets made of documents can also be well modeled as bipartite multivariate dynamic networks like for example thesis dissertations: a thesis is a document with specific attributes such as the subject, the doctoral school, the domain, the university, and the date of defense, and mention several peoples who are socially connected through the thesis defense with different roles: author (*auteur* in french), director(s) (*directeur*), referees (*rapporteur*), and jury president (*président de jury*). We present here an exploration of the data by ourselves using ComBiNet. A first look at the graph measures tells us that 896 theses have been defended in sociology in France between 2016 and 2021 in France, with 2453 persons included in the defenses (see Figure 4.16 bottom). The bipartite node-link view shows us an overview of the network but is hard to parse due to the network's size. Zoom actions though allow centering the view for specific parts of the network. The map view allows us to see that theses has been defended all around France, even though the majority are defended in Paris. This is confirmed by a look at the distribution of the cities (Figure 4.16 bottom right): around half of the defenses are in Paris, compared to the rest of the country which is more or less homogeneous. By setting the threshold to link creation to one (meaning that a link is created between two documents if they mention at least one common person), a lot of links are created as seen in Figure 4.16 (right). It means that a lot of thesis defenses include referees and juries from different cities.

Let's now try to answer an interesting question: "Do referees and jury presidents often ask thesis directors to be referees and jury presidents in their turn of another thesis where they are

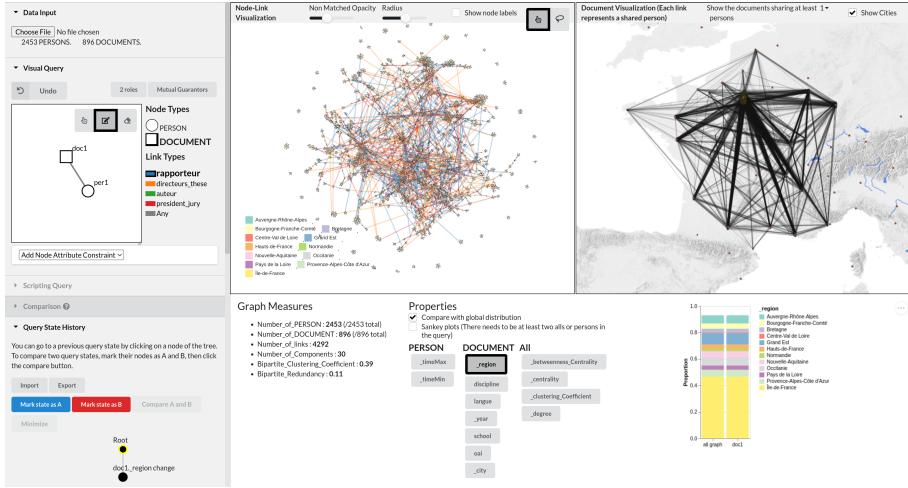


Figure 4.16 – ComBiNet used for exploring theses of sociology defended in France between 2016 and 2021. The bipartite and map views show an overview of two visions of the network. The user selects the *region* attribute, showing the geographical distribution of the defended thesis.

directors ?". For this, we can construct a visual query representing this pattern by creating two person nodes and two document nodes, and by connecting them with two president links and two referee or jury director links in a symmetrical way, as shown in Figure 4.17 (right). The occurrence table tells us that this pattern has been found 76 times in the network, meaning that this is a recurrent behavior. We are now interested in characterizing the thesis occurring in this pattern, by their regions. We can look at the *city* attribute distribution for this thesis by selecting it in the attribute view as shown in Figure 4.17 (bottom right). We can first see on the map that this pattern occurs mainly in the biggest cities of the country. By selecting the Sankey view option, we can investigate if this pattern occurs between thesis defended in different regions or if it occurs mainly in the same ones. We learn that it depends mainly on the regions: in Bourgogne-Fanche-Comté 26 out of 29 theses are connected with the thesis of another region. On contrary, in *Occitanie* it is the case for only 4 out of 17. On average, we can see that this pattern occurs a lot for theses of the same region. In *Île-de-France*, it is the case for around half of the thesis (28/50). This exploratory analysis shows that ComBiNet can be used to explore and gain insight into such datasets.

4.6 Formative Usability Study

I performed a formative usability study with two historians and one expert in visualization. I had 3 meetings with each and gave them control of the tool to see if they could use it to explore their data—the visualization expert used the interface with the dataset of construction contracts #1—and perform queries and comparisons. At the first meeting, I explained to them the panels of the system and each features. During each session, they were free to explore the data as they wanted. If they were stuck using one feature, I helped them by explaining

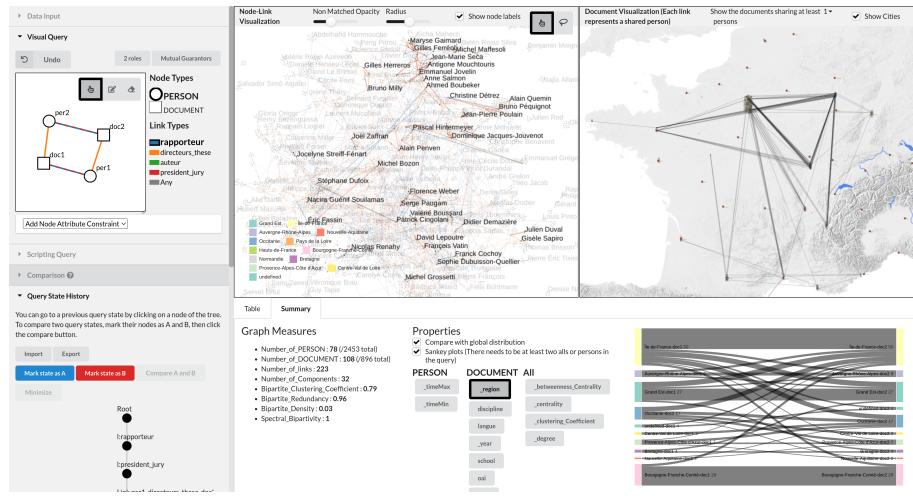


Figure 4.17 – Sociology thesis dataset explored with ComBiNet. The user constructed a visual query to see if there are symmetrical relationships between thesis directors and referees (or jury directors). The *region* attribute is selected with the Sankey option, letting the user see if there are correlations between the regions of the thesis found in this pattern.

how to use it. If they seemed to not know what to do next, I asked them to answer a specific question on the data (for example “can you find the persons who collaborated in more than two contracts between 1711 and 1714 in Torino”). At each meeting, I asked them to speak aloud, commenting on their aims and actions. At the end of each session, they reported their general feedback, what they did not like or understand, and what other features they would like to have to explore their data.

I improved the system and made the changes asked by the users before setting up new appointments. This usability study led to the redesign of some core features, like the activation of the comparison mode which is now started by first marking the state nodes in the provenance tree. It also led to the implementation of new features, such as the person and document tables (which are updated after each query), the persistent selection of nodes across the two views and the tables, and the undo feature for visual queries. At the final meetings, the three users were able to perform exploration, queries, and comparisons to answer socio-historical questions by themselves.

4.6.1 Feedback

All three users liked the table views and were exploiting them to study in depth who were the person and documents found in their specific queries. Both historians liked the Sankey diagram of the attributes, allowing them to see the evolution of distributions and answering several of their questions. Our collaborator of the use case #2 was making sense of it by linking the migration patterns she was seeing in the Sankey diagram with specific persons of the dataset she knew in depth. She was also curious about other migration patterns she was not aware of and wanted to know who these persons were, the system allowing her to select

them and follow a deeper exploration. One other historian outside of our direct collaborators liked the overlay of node attributes on the bipartite and map views, and the distribution plots. She said: “With this data model, even if historians realize the structure of their network do not allow to answer their research questions, they can still visualize and compare attributes of documents and persons with this interface, which is always useful in quantitative history”.

4.7 Discussion

I discuss in this section several points of potential limitation for the system.

Query Expressiveness. The visual query system currently allows finding occurrences of attributed subgraphs, with potential union operations on constraints (links and node attribute values can be set at one value or as a set of values). Being able to express attribute constraints (other than for labels and ids) and unions is new compared to other visual graph query systems. More complex constraints are then expressible using the Cypher editor, such as dependent constraints, e.g., if one node attribute value has to be greater or lower than another attribute value. The visual query system could be extended by introducing more complex time constraints capabilities, such as in [93].

Scalability. We assess the scalability in network size (number of nodes and links) concerning the cluttering and readability of the network visualizations. Our biggest dataset from #3 comprises 7212 nodes (4886 persons and 2326 events) and 7790 links, after splitting the documents into birth and marriage event nodes. The system allows the exploration of networks of this size with a decent frame rate. ComBiNet allows navigating relatively large sparse graphs (thousands of nodes) with the node-link visualization using zoom & pan and filtering with the query system. It lets users focus on subsets of the data, one or two at a time.

Generalizability. The system has been designed specifically for bipartite multivariate dynamic networks, which models well a diversity of historical sources we encountered via our collaborations: marriage acts, birth/death certificates, construction/work contracts, census, and migrations forms. Moreover, bipartite multivariate dynamic network can also be used to model other similar data types, such as scientific publications or thesis data. However, other kinds of historical textual data exist where documents can mention each other, such as in private letters for example. The model and interface would need to be slightly modified to take into account document-to-document links for these datasets. Bipartite networks are also used in various other disciplines, such as biology [73] and chemistry [74]. ComBiNet could be extended to these other application domains, in particular by modifying the map view to show other location data related to the entities of the network, or removing it altogether if it makes no sense for a particular domain.

4.8 Conclusion and Future Work

I presented in this chapter ComBiNet, a VA system for exploring social networks modeled from historical textual sources, primarily aimed at social historians. It relies on modeling documents as bipartite, multivariate, dynamic, geolocated social networks where persons are linked

to documents or events using typed links that express roles. With this data model, ComBiNet let historians explore a concrete representation of their annotated documents (i.e., the model express the *reality* of the documents, without the use of projections or distortions) with *traceability* to the original sources and *simplicity*. Historians can hence reflect on their encoding process and answer their socio-historical questions using 1) dynamic queries on the network structure and attributes to highlight groups of interests or erroneous patterns and 2) visual comparisons to contrast selected groups according to their structure, location, time, or any other attribute. The results can be visualized as a bipartite node-link diagram, a geographical map, graph measures, and distributions of values for the attributes. I have shown that complex explorations and analyses were easy to perform, and validated our approach by first describing four use cases among several other projects we are collaborating with and by performing a formative usability study showing that after many improvements the system is usable by social scientists.

By specifying a unifying data model and novel high-level visual and interactive mechanisms for comparing topology, attributes, and time, social scientists were able to correct their data more easily with exploration and querying errors-induced patterns. Thanks to the document-centered model, it was easy for them to trace back the errors and inconsistencies to the sources for corrections. With the same representation, they were able to operate explorations and analyses using easy-to-use interactions implemented in ComBiNet such as coordinated views, visual querying, and comparison. Using these mechanisms, social scientists performed visual exploratory analyses of their network based on topological and attribute descriptions, and comparisons of subgroups of interests—between them or the overall network. This methodology allows them to either ground or refute their hypotheses in network measures and attribute distributions, or to generate new ones from new insight revealed thanks to the exploratory and interaction mechanisms.

We believe ComBiNet leads the way toward a new generation of highly interactive exploration tools applicable to wrangle and analyze a wide variety of real social networks modeled from textual sources, with a focus on the *reality* of the documents, *traceability* of the network and results, and *simplicity* of use, which are essential for historical work.

For future work, ComBiNet could be extended to support more SNA measures and computations such as clustering; it would create a new attribute containing a cluster identifier. The interface currently proposes two layouts based on the topology and the geolocations of the entities. Providing more layout options could be interesting, especially one to highlight better the time, similar to the PAOHvis technique [142]. Finally, the interface could in the future make suggestions on the query construction process based on frequent subgraphs similar to VERTIGo [28], and within a mixed-initiative perspective [?].

In the next chapter, I present a visual interface following this mixed-initiative framework for network clustering, to answer Q3 and showing that such approaches can help social scientists use data mining tools while controlling their biases—with the condition of an explicit results' traceability.

Bibliography

- [1] NodeXL: Simple network analysis for social media.
- [2] Mashael Alkadi, Vanessa Serrano, James Scott-Brown, Catherine Plaisant, Jean-Daniel Fekete, Uta Hinrichs, and Benjamin Bach. Understanding barriers to network exploration with visualization: A report from the trenches. *IEEE Transactions on Visualization and Computer Graphics*, 27(2), 2022. 9, 15, 17, 23, 38, 43, 45, 46, 48, 51, 59
- [3] Keith Andrews, Martin Wohlfahrt, and Gerhard Wurzinger. Visual Graph Comparison. In *2009 13th International Conference Information Visualisation*, pages 62–67, July 2009. 65
- [4] F. J. Anscombe. Graphs in Statistical Analysis. *The American Statistician*, 27(1):17–21, February 1973. 7, 24, 26
- [5] Mathieu Bastian, Sébastien Heymann, and Mathieu Jacomy. Gephi: An open source software for exploring and manipulating networks. In Eytan Adar, Matthew Hurst, Tim Finin, Natalie S. Glance, Nicolas Nicolov, and Belle L. Tseng, editors, *ICWSM*. The AAAI Press, 2009. 14, 41, 62, 70, 91
- [6] Sugato Basu, Ian Davidson, and Kiri Wagstaff. *Constrained Clustering: Advances in Algorithms, Theory, and Applications*. Chapman & Hall/CRC, first edition, 2008. 95
- [7] Leilani Battle and Jeffrey Heer. Characterizing Exploratory Visual Analysis: A Literature Review and Evaluation of Analytic Provenance in Tableau. *Computer Graphics Forum*, 38(3):145–159, 2019. 65
- [8] Jacques Bertin. *Sémiologie graphique: les diagrammes, les réseaux, les cartes*. Paris: Gauthier-Villars, 1967. 7, 24, 25
- [9] A. Bezerianos, F. Chevalier, P. Dragicevic, N. Elmquist, and J.d. Fekete. GraphDice: A System for Exploring Multivariate Social Networks. *Computer Graphics Forum*, 29(3):863–872, 2010. 71
- [10] Marc Bloch. *Apologie Pour l'histoire*. A. Colin, 1949. 10
- [11] Christian Böhm and Claudia Plant. HISSCLU: A hierarchical density-based method for semi-supervised clustering. In *Proceedings of the 11th International Conference on Extending Database Technology: Advances in Database Technology*, EDBT '08, pages 440–451, New York, NY, USA, 2008. ACM. 95
- [12] Christian Bors, John Wenskovitch, Michelle Dowling, Simon Attfield, Leilani Battle, Alex Endert, Olga Kulyk, and Robert S. Laramee. A Provenance Task Abstraction Framework. *IEEE Computer Graphics and Applications*, 39(6):46–60, November 2019. 65

- [13] Michael Bostock, Vadim Ogievetsky, and Jeffrey Heer. D³ Data-Driven Documents. *IEEE Transactions on Visualization and Computer Graphics*, 17(12):2301–2309, December 2011. 70, 81
- [14] Pierre Bourdieu. Sur les rapports entre la sociologie et l'histoire en Allemagne et en France. *Actes de la Recherche en Sciences Sociales*, 106(1):108–122, 1995. 29
- [15] Ulrik Brandes, Daniel Delling, Marco Gaertler, Robert Gorke, Martin Hoefer, Zoran Nikoloski, and Dorothea Wagner. On Modularity Clustering. *IEEE Transactions on Knowledge and Data Engineering*, 20(2):172–188, February 2008. 51
- [16] Peter Burke. *History and Social Theory*. Polity, 2005. 29
- [17] Steven P. Callahan, Juliana Freire, Emanuele Santos, Carlos E. Scheidegger, Cláudio T. Silva, and Huy T. Vo. VisTrails: Visualization meets data management. In *Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data - SIGMOD '06*, page 745, Chicago, IL, USA, 2006. ACM Press. 65
- [18] Charles-Olivier Carbonell. *L'Historiographie*. FeniXX, January 1981. 28
- [19] Stuart-K. Card, Jock-D. Mackinlay, and Ben Shneiderman. *Readings in Information Visualization: Using Vision to Think*. Morgan Kaufmann Publishers Inc, San Francisco, Calif, February 1999. 12, 23, 24
- [20] Duen Horng Chau, Christos Faloutsos, Hanghang Tong, Jason I. Hong, Brian Gallagher, and Tina Eliassi-Rad. GRAPHITE: A Visual Query System for Large Graphs. In *2008 IEEE International Conference on Data Mining Workshops*, pages 963–966, December 2008. 65
- [21] J. S. Coleman. Introduction to mathematical sociology. *Introduction to mathematical sociology.*, 1964. 34
- [22] TEI Consortium. TEI P5: Guidelines for electronic text encoding and interchange, February 2021. 50
- [23] Pascal Cristofoli. Aux sources des grands réseaux d'interactions. *Reseaux*, 152(6):21–58, 2008. 9, 15, 36, 46, 48, 51, 56, 62
- [24] Pascal Cristofoli. Principes et usages des dessins de réseaux en SHS. *La visualisation des données en histoire*, page 35, 2015. 10, 40, 70
- [25] Pascal Cristofoli and Nicoletta Rolla. Temporalités à l'œuvre dans les chantiers du bâtiment. *Temporalités. Revue de sciences sociales et humaines*, (27), June 2018. 8, 53, 55, 66, 70
- [26] Tarik Crnovrsanin, Chris W. Muelder, Robert Faris, Diane Felmlee, and Kwan-Liu Ma. Visualization techniques for categorical analysis of social networks with multiple edge sets. *Social Networks*, 37:56–64, 2014. 53

- [27] Gabor Csardi and Tamas Nepusz. The igraph software package for complex network research. *InterJournal, Complex Systems*:1695, 2006. 64, 70
- [28] Erick Cuenca, Arnaud Sallaberry, Dino Ienco, and Pascal Poncelet. VERTIGO: A Visual Platform for Querying and Exploring Large Multilayer Networks. *IEEE Transactions on Visualization and Computer Graphics*, pages 1–1, 2021. 65, 90
- [29] Zach Cutler, Kiran Gadhave, and Alexander Lex. Trrack: A Library for Provenance-Tracking in Web-Based Visualizations. In *2020 IEEE Visualization Conference (VIS)*, pages 116–120, October 2020. 78, 81
- [30] Mandeep K. Dhami, Ian K. Belton, and David R. Mandel. The “analysis of competing hypotheses” in intelligence analysis. *Applied Cognitive Psychology*, 33(6):1080–1090, 2019. 11
- [31] Jana Diesner, Craig Evans, and Jinseok Kim. Impact of Entity Disambiguation Errors on Social Network Properties. *Proceedings of the International AAAI Conference on Web and Social Media*, 9(1):81–90, 2015. 15, 43, 46, 50
- [32] Dana Diminescu. The migration of ethnic germans from romania to west germany: Insights from the archives of the former communist regime. In *CERS, Public Lecture, UCLA*, Los Angeles, United States, March 2020. 56, 68
- [33] Nicole Dufournaud. La recherche empirique en histoire à l’ère numérique. *Gazette des archives*, 240(4):397–407, 2015. 9
- [34] Nicole Dufournaud. Comment rendre visible le rôle économique des femmes sous l’Ancien Régime ? Étude méthodologique sur les marchandes à Nantes aux XVI^e et XVII^e siècles. In Bernard Michon and Nicole Dufournaud, editors, *Femmes et Négoce Dans Les Ports Européens (Fin Du Moyen Age - XIXe Siècle)*, pages 65–84. Peter Lang, 2018. 7, 12, 15, 48, 51
- [35] Nicole Dufournaud and Jean-Daniel Fekete. Comparaison d’outils pour la visualisation de sources historiques codées en XML/TEI. *Document numérique*, 9(2):37–56, April 2006. 50
- [36] Cody Dunne, Nathalie Henry Riche, Bongshin Lee, Ronald Metoyer, and George Robertson. GraphTrail: Analyzing large multivariate, heterogeneous networks while supporting exploration history. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI ’12, pages 1663–1672, New York, NY, USA, May 2012. Association for Computing Machinery. 65
- [37] P. Erdős and A. Rényi. On the evolution of random graphs. In *On the Evolution of Random Graphs*, pages 38–82. Princeton University Press, October 2011. 34
- [38] Emily Erikson and Peter Bearman. Malfeasance and the Foundations for Global Trade: The Structure of English Trade in the East Indies, 1601–1833. *American Journal of Sociology*, 112(1):195–230, July 2006. 52

- [39] Michael Eve. Deux traditions d'analyse des réseaux sociaux. *Réseaux*, 115(5):183–212, 2002. 17, 35, 36
- [40] Lucien Febvre. VERS UNE AUTRE HISTOIRE. *Revue de Métaphysique et de Morale*, 54(3/4):225–247, 1949. 28
- [41] Jean-Daniel Fekete, Danyel Fisher, Arnab Nandi, and Michael Sedlmair. *Progressive Data Analysis and Visualization*. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, April 2019. 76
- [42] L. Freeman. Visualizing Social Networks. *J. Soc. Struct.*, 2000. 10, 12, 13, 39
- [43] L.C. Freeman. *The Development of Social Network Analysis: A Study in the Sociology of Science*. Empirical Press, 2004. 9, 11, 23, 32, 33, 34, 35, 51
- [44] Manuel Freire, Catherine Plaisant, Ben Shneiderman, and Jen Golbeck. ManyNets: An interface for multiple network analysis and visualization. In *CHI '10*, CHI '10, pages 213–222, New York, NY, USA, 2010. ACM. 65
- [45] Michael Friendly. Visions and Re-Visions of Charles Joseph Minard. *Journal of Educational and Behavioral Statistics*, 27(1):31–51, March 2002. 24
- [46] Michael Friendly. A Brief History of Data Visualization. In Chun-houh Chen, Wolfgang Härdle, and Antony Unwin, editors, *Handbook of Data Visualization*, Springer Handbooks Comp.Statistics, pages 15–56. Springer, Berlin, Heidelberg, 2008. 24
- [47] GEDCOM: The genealogy data standard. 39
- [48] Mohammad Ghoniem, J.-D. Fekete, and Philippe Castagliola. A comparison of the readability of graphs using node-link and matrix-based representations. In *IEEE Symposium on Information Visualization*, pages 17–24. Ieee, 2004. 40
- [49] Carlo Ginzburg and Carlo Poni. La micro-histoire. *Le Débat*, 17(10):133, 1981. 11, 36, 47
- [50] Barney G. Glaser and Anselm L. Strauss. *The Discovery of Grounded Theory: Strategies for Qualitative Research*. Aldine Transaction, New Brunswick, 5. paperback print edition, 2010. 11
- [51] Michael Gleicher. Considerations for visualizing comparison. *IEEE Trans. Vis. Comput. Graphics*, 24(1):413–423, 2018. 65
- [52] Claudia Goldin. Cliometrics and the Nobel. *Journal of Economic Perspectives*, 9(2):191–208, June 1995. 10
- [53] Martin Grandjean. Social network analysis and visualization: Moreno's Sociograms revisited, 2015. 7, 33

- [54] Martin Grandjean. Analisi e visualizzazioni delle reti in storia. L'esempio della cooperazione intellettuale della Società delle Nazioni. *ME*, (2/2017), 2017. 62
- [55] Maurizio Gribaudo and Alain Blum. Des catégories aux liens individuels : l'analyse statistique de l'espace social. *Annales*, 45(6):1365–1402, 1990. 9
- [56] Jo Guldi and David Armitage. *The History Manifesto*. Cambridge University Press, October 2014. 9
- [57] Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics, and function using NetworkX. In Gaël Varoquaux, Travis Vaught, and Jarrod Millman, editors, *Proceedings of the 7th Python in Science Conference*, pages 11–15, Pasadena, CA USA, 2008. 64, 91
- [58] Klaus Hamberger, Cyril Grange, Michael Houseman, and Christian Momon. Scanning for patterns of relationship: Analyzing kinship and marriage networks with Puck 2.0. *The History of the Family*, 19(4):564–596, October 2014. 37, 39, 52
- [59] Klaus Hamberger, Michael Houseman, and R. White, Douglas. Kinship network analysis. In John Scott & Peter J. Carrington, editor, *The Sage Handbook of Social Network Analysis*, pages 533–549. Sage Publications, 2011. 9, 39
- [60] Mountaz Hascoët and Pierre Dragicevic. Interactive graph matching and visual comparison of graphs and clustered graphs. In Genny Tortora, Stefano Levialdi, and Maurizio Tucci, editors, *AVI '12*, pages 522–529. ACM, 2012. 65
- [61] J. Heer and D. Boyd. Vizster: Visualizing online social networks. In *IEEE Symposium on Information Visualization, 2005. INFOVIS 2005.*, pages 32–39, October 2005. 72
- [62] Jeffrey Heer. Agency plus automation: Designing artificial intelligence into interactive systems. *Proceedings of the National Academy of Sciences*, 116(6):1844–1850, 2019. 118
- [63] Louis Henry and Michel Fleury. Des registres paroissiaux à l'histoire de la population: Manuel de dépouillement et d'exploitation de l'état civil ancien. *Population (French Edition)*, 11(1):142–144, 1956. 10
- [64] Nathalie Henry, Jean-Daniel Fekete, and Michael J. McGuffin. NodeTrix: A Hybrid Visualization of Social Networks. *IEEE Transactions on Visualization and Computer Graphics*, 13(6):1302–1309, November 2007. 7, 40, 41
- [65] Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia D'amato, Gerard De Melo, Claudio Gutierrez, and Sabrina Kirrane et al. Knowledge graphs. *ACM Comput. Surv.*, 54(4), July 2021. 53
- [66] Pat Hudson and Mina Ishizu. *History by Numbers: An Introduction to Quantitative Approaches*. Bloomsbury Publishing, November 2016.

- [67] Infovis SC policies FAQ. 109
- [68] Piers J Ingram, Michael PH Stumpf, and Jaroslav Stark. Network motifs: Structure does not determine function. *BMC Genomics*, 7:108, May 2006. 64
- [69] Frédéric Kaplan. The Venice Time Machine. In *Proceedings of the 2015 ACM Symposium on Document Engineering*, DocEng '15, page 73, New York, NY, USA, September 2015. Association for Computing Machinery. 14, 32
- [70] Karine Karila-Cohen, Claire Lemercier, Isabelle Rosé, and Claire Zalc. Nouvelles cuisines de l'histoire quantitative. *Annales. Histoire, Sciences Sociales*, 73(4):773–783, December 2018. 9, 47, 48, 52
- [71] Daniel Keim, Gennady Andrienko, Jean-Daniel Fekete, Carsten Görg, Jörn Kohlhammer, and Guy Melançon. Visual Analytics: Definition, Process, and Challenges. In Andreas Kerren, John T. Stasko, Jean-Daniel Fekete, and Chris North, editors, *Information Visualization: Human-Centered Issues and Perspectives*, Lecture Notes in Computer Science, pages 154–175. Springer, Berlin, Heidelberg, 2008. 7, 14, 27
- [72] Florian Kerschbaumer, Linda von Keyserlingk-Rehbein, Martin Stark, and Marten Düring. *The Power of Networks. Prospects of Historical Network Research*. Routledge, December 2021. 9, 10, 23, 38
- [73] Steffen Klamt, Utz-Uwe Haus, and Fabian Theis. Hypergraphs and cellular networks. *PLoS computational biology*, 5(5):e1000385, 2009. 89
- [74] Elena V. Konstantinova and Vladimir A. Skorobogatov. Application of hypergraph theory in chemistry. *Discrete Mathematics*, 235(1-3):365–383, May 2001. 89
- [75] C. Kosak, J. Marks, and S. Shieber. Automating the layout of network diagrams with specified visual organization. *IEEE Transactions on Systems, Man, and Cybernetics*, 24(3):440–454, March 1994. 7, 40
- [76] Ernest Labrousse. *La Crise de l'économie Française à La Fin de l'Ancien Régime et Au Début de La Révolution*, volume 1. Presses Universitaires de France-PUF, 1990. 29
- [77] Charles-Victor Langlois and Charles Seignobos. *Introduction aux études historiques*. ENS Éditions, February 2014.
- [78] Katherine A. Larson. Thomas F. Tartaron, Maritime Networks in the Mycenaean World. New York: Cambridge University Press, 2013. *Comparative Studies in Society and History*, 56(4):1064–1065, October 2014. 9
- [79] Emmanuel Lazega. *Réseaux sociaux et structures relationnelles*. Presses universitaires de France, Paris, 1998. 13, 36

- [80] Claire Lemercier. 12. Formal network methods in history: Why and how? In Georg Fertig, editor, *Social Networks, Political Institutions, and Rural Societies*, volume 11, pages 281–310. Brepols Publishers, Turnhout, January 2015. 9, 15, 23, 35, 37, 38, 39, 43, 45, 53, 62
- [81] Claire Lemercier and Claire Zalc. *Quantitative Methods in the Humanities: An Introduction*. University of Virginia Press, March 2019. 10, 11, 15, 30, 32, 39, 46, 47, 48, 63, 116
- [82] Claire Lemercier and Claire Zalc. Back to the Sources: Practicing and Teaching Quantitative History in the 2020s. *Capitalism*, 2(2):473–508, 2021. 15, 29, 45, 46, 47, 48, 52
- [83] Bernard Lepetit. L'histoire quantitative : deux ou trois choses que je sais d'elle. *Histoire & Mesure*, 4(3):191–199, 1989. 9, 47
- [84] Carola Lipp. Kinship Networks, Local Government, and Elections in a Town in Southwest Germany, 1800–1850. *Journal of Family History*, 30(4):347–365, October 2005. 37
- [85] Gribaudi Maurizio. *Espaces, Temporalités, Stratifications : Exercices Méthodologiques Sur Les Réseaux Sociaux*. Editions de l'Ecole des Hautes Etudes en Sciences Sociales, Paris, January 2000. 35
- [86] Philip Mayer. Migrancy and the Study of Africans in Towns. *American Anthropologist*, 64(3):576–592, 1962. 36
- [87] Fintan McGee, Benjamin Renoust, Daniel Archambault, Mohammad Ghoniem, Andreas Kerren, and Bruno Pinaud et al. *Visual Analysis of Multilayer Networks*. Synthesis Lectures on Visualization. Morgan & Claypool Publishers, 2021. 53
- [88] Michael J. McGuffin. Simple algorithms for network visualization: A tutorial. *Tsinghua Science and Technology*, 17(4):383–398, August 2012. 40
- [89] Pierre Mercklé and Claire Zalc. Peut-on modéliser la persécution ?: Apports et limites des approches quantifiées sur le terrain de la Shoah. *Annales. Histoire, Sciences Sociales*, 73(4):923–957, December 2018.
- [90] R. Michalski, P. Kazienko, and D. Krol. Predicting Social Network Measures Using Machine Learning Approach. In *2012 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pages 1056–1059, Istanbul, August 2012. IEEE. 118
- [91] R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon. Network Motifs: Simple Building Blocks of Complex Networks. *Science*, 298(5594):824–827, October 2002. 35, 64
- [92] Christoph Molnar. *Interpretable Machine Learning. A Guide for Making Black Box Models Explainable*. Lulu.com, 2019. 94

- [93] Megan Monroe, Krist Wongsuphasawat, Catherine Plaisant, Ben Shneiderman, Jeff Millstein, and Sigfried Gold. Exploring point and interval event patterns: Display methods and interactive visual query. *University of Maryland Technical Report*, 2012. 89
- [94] J. L. Moreno. *Who Shall Survive?: A New Approach to the Problem of Human Interrelations*. Who Shall Survive?: A New Approach to the Problem of Human Interrelations. Nervous and Mental Disease Publishing Co, Washington, DC, US, 1934. 7, 33, 39
- [95] J. L. Moreno. Foundations of Sociometry: An Introduction. *Sociometry*, 4(1):15, February 1941. 33, 52
- [96] Zacarias Moutoukias. Buenos Aires, port between two oceans: Mobilities, networks, stratifications (2nd half of the 18th century). *E-SPANIA-REVUE ELECTRONIQUE D'ETUDES HISPANIQUES MEDIEVALES*, 25, 2016. 56, 67
- [97] Zacharias Moutoukias. Réseaux personnels et autorité coloniale : Les négociants de Buenos Aires au XVIIIe siècle. *Annales. Histoire, Sciences Sociales*, 47(4-5):889–915, October 1992. 37
- [98] Andrej Mrvar and Vladimir Batagelj. Analysis and visualization of large networks with program package Pajek. *Complex Adaptive Systems Modeling*, 4(1), April 2016. 41, 70
- [99] Natural earth. 71
- [100] Neo4j graph data platform. 63, 64, 81
- [101] Rolla Nicoletta. Mobilité et conflits. Travailler sur les chantiers de construction piémontais dans la première moitié du XVIIIe siècle. In Andrea Caracausi and Marco Schnyder, editors, *Travail et Mobilité En Europe (XVle-XIXe Siècles)*, Coll. Histoire et Civilisations. Presses universitaires du Septentrion, Villeneuve d'Ascq, 2018. 55
- [102] Carolina Nobre, Marc Streit, and Alexander Lex. Juniper: A Tree+Table Approach to Multivariate Graph Visualization. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):544–554, January 2019. 40
- [103] Gérard Noiriel. Naissance du métier d'historien. *Genèses. Sciences sociales et histoire*, 1(1):58–85, 1990.
- [104] Maryjane Osa. *Solidarity And Contention: Networks Of Polish Opposition*. Univ Of Minnesota Press, Minneapolis, first edition edition, July 2003. 9, 51
- [105] John F. Padgett and Christopher K. Ansell. Robust Action and the Rise of the Medici, 1400-1434. *American Journal of Sociology*, 98(6):1259–1319, May 1993. 7, 9, 13, 37, 38
- [106] Pajek — Analysis and visualization of very large networks. 14, 16, 91

- [107] Terence J. Parr and Russell W. Quong. ANTLR: A predicated-LL (k) parser generator. *Software: Practice and Experience*, 25(7):789–810, 1995. 81
- [108] Vanessa Peña-Araya, Tong Xue, Emmanuel Pietriga, Laurent Amsaleg, and Anastasia Bezerianos. HyperStorylines: Interactively untangling dynamic hypergraphs. *Information Visualization*, 21(1):38–62, January 2022. 58
- [109] Cindarella Sarah Maria Petz. *On Combining Network Research and Computational Methods on Historical Research Questions and Its Implications for the Digital Humanities*. PhD thesis, Technische Universität München, 2022. 10, 38, 47
- [110] James P. Philips and Nasseh Tabrizi. Historical Document Processing: Historical Document Processing: A Survey of Techniques, Tools, and Trends, September 2020. 118
- [111] Robert Pienta, Fred Hohman, Alex Endert, Acar Tamersoy, Kevin Roundy, Chris Gates, Shamkant Navathe, and Duen Horng Chau. VIGOR: Interactive Visual Exploration of Graph Query Results. *IEEE Transactions on Visualization and Computer Graphics*, 24(1):215–225, January 2018. 65
- [112] Alexis Pister, Nicole Dufournaud, Pascal Cristofoli, Christophe Prieur, and Jean-Daniel Fekete. From Historical Documents To Social Network Visualization: Potential Pitfalls and Network Modeling. In *7th Workshop on Visualization for the Digital Humanities (VIS4DH)*, 2022. 45
- [113] Antoine Prost. *Douze Leçons sur l'histoire*. Média Diffusion, April 2014. 10, 23, 28
- [114] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2014. 64, 91
- [115] Eric Ragan, Endert Alex, Jibonananda Sanyal, and Jian Chen. Characterizing Provenance in Visualization and Data Analysis: An Organizational Framework of Provenance Types and Purposes. *IEEE Transactions on Visualization and Computer Graphics*, 22(1), January 2016. 65
- [116] Ramana Rao and Stuart K. Card. The table lens: Merging graphical and symbolic representations in an interactive focus + context visualization for tabular information. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '94*, pages 318–322, New York, NY, USA, 1994. Association for Computing Machinery. 114
- [117] Donghao Ren, Bongshin Lee, and Matthew Brehmer. Charticulator: Interactive Construction of Bespoke Chart Layouts. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):789–799, January 2019. 65
- [118] Pedro Ribeiro and Fernando Silva. Discovering Colored Network Motifs. In Pierluigi Contucci, Ronaldo Menezes, Andrea Omicini, and Julia Poncela-Casasnovas, editors, *Complex Networks V, Studies in Computational Intelligence*, pages 107–118, Cham, 2014. Springer International Publishing. 64

- [119] Giulio Rossetti and Rémy Cazabet. Community discovery in dynamic networks: A survey. *ACM Comput. Surv.*, 51(2), February 2018. 95
- [120] Fabrice Rossi, Nathalie Vialaneix, and Florent Hautefeuille. Exploration of a large database of French notarial acts with social network methods. *Digital Medievalist*, 9:2013, July 2014. 62
- [121] Juan A. Rubio-Mondejar and Josean Garrues-Irurzun. Women entrepreneurs and family networks in Andalusia (Spain) during the second industrial revolution. *Business History*, pages 1–22, May 2022. 9
- [122] C.J. Rueda and Catedral de Buenos Aires. *Matrimonios de La Catedral de Buenos Aires, 1747-1823*. Number v. 2 in Fuentes Históricas y Genealógicas Argentinas. Fuentes Históricas y Genealógicas Argentinas, 1989. 56
- [123] Anni Sairio. Methodological and practical aspects of historical network analysis: A case study of the Bluestocking letters. In Arja Nurmi, Minna Nevala, and Minna Palander-Collin, editors, *Pragmatics & Beyond New Series*, volume 183, pages 107–135. John Benjamins Publishing Company, Amsterdam, 2009. 52
- [124] Bahador Saket, Paolo Simonetto, and Stephen Kobourov. Group-level graph visualization taxonomy. In N. Elmquist, M. Hlawitschka, and J. Kennedy, editors, *EuroVis - Short Papers*. The Eurographics Association, 2014. 96
- [125] Arvind Satyanarayan, Dominik Moritz, Kanit Wongsuphasawat, and Jeffrey Heer. Vega-lite: A grammar of interactive graphics. *IEEE Trans. Vis. Comput. Graphics*, 23(1):341–350, 2016. 26, 81
- [126] Shruti S. Sawant and Manoharan Prabukumar. A review on graph-based semi-supervised learning methods for hyperspectral image classification. *The Egyptian Journal of Remote Sensing and Space Science*, 2018. 95
- [127] John Scott. Social Network Analysis. *Sociology*, 22(1):109–127, February 1988. 23, 33, 34, 41, 51
- [128] Vanessa Serrano Molinero, Benjamin Bach, Catherine Plaisant, Nicole Dufournaud, and Jean-Daniel Fekete. Understanding the use of the vistorian: Complementing logs with context mini-questionnaires. In *Visualization for the Digital Humanities Workshop*, Phoenix, United States, October 2017. 7, 40, 41, 42, 43, 50, 58
- [129] Rachel Shadoan and Chris Weaver. Visual Analysis of Higher-Order Conjunctive Relationships in Multidimensional Data Using a Hypergraph Query System. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2070–2079, December 2013. 65
- [130] Termeh Shafie, David Schoch, Jimmy Mans, Corinne Hofman, and Ulrik Brandes. Hypergraph Representations: A Study of Carib Attacks on Colonial Forces, 1509-1700. *Journal of Historical Network Research*, pages 52–70 Pages, October 2017. 62

- [131] Ben Shneiderman. Dynamic queries for visual information seeking. *IEEE Softw.*, 11(6):70–77, November 1994. 75
- [132] Georg Simmel. *Soziologie: Untersuchungen über die Formen der Vergesellschaftung*. Duncker & Humblot, Berlin, 7. aufl edition, 2013. 35
- [133] Marc A. Smith, Ben Shneiderman, Natasa Milic-Frayling, Eduarda Mendes Rodrigues, Vladimir Barash, Cody Dunne, Tony Capone, Adam Perer, and Eric Gleave. Analyzing (social media) networks with NodeXL. In John M. Carroll, editor, *Proceedings of the Fourth International Conference on Communities and Technologies, C&T 2009, University Park, PA, USA, June 25-27, 2009*, pages 255–264. ACM, 2009. 14, 41, 62, 70
- [134] SNA — Tools for social network analysis.
- [135] John Snow. On the Mode of Communication of Cholera. *Edinb Med J*, 1(7):668–670, January 1856. 24
- [136] John T. Stasko, Carsten Görg, and Zhicheng Liu. Jigsaw: Supporting investigative analysis through interactive visualization. *Inf. Vis.*, 7(2):118–132, 2008. 43, 54, 58, 117
- [137] Alexander Strehl and Joydeep Ghosh. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *Journal of machine learning research*, 3(Dec):583–617, 2002. 96
- [138] Shazia Tabassum, Fabiola S. F. Pereira, Sofia Fernandes, and João Gama. Social network analysis: An overview. *WIREs Data Mining and Knowledge Discovery*, 8(5):e1256, 2018. 9, 13, 34
- [139] Natkamon Tovanich, Alexis Pister, Gaelle Richer, Paola Valdivia, Christophe Prieur, Jean-Daniel Fekete, and Petra Isenberg. VAST 2020 Contest Challenge: GraphMatchMaker: Visual Analytics for Graph Comparison and Matching. *IEEE Computer Graphics and Applications*, pages 1–1, 2021. 64, 65
- [140] John W. Tukey. The Future of Data Analysis. *The Annals of Mathematical Statistics*, 33(1):1–67, 1962. 24
- [141] John W. Tukey. *Exploratory Data Analysis*. Pearson, Reading, Mass, 1er édition edition, January 1977. 13, 26
- [142] Paola Valdivia, Paolo Buono, Catherine Plaisant, Nicole Dufournaud, and Jean-Daniel Fekete. Analyzing Dynamic Hypergraphs with Parallel Aggregated Ordered Hypergraph Visualization. *IEEE Trans. Visual. Comput. Graphics*, 27(1):1–13, January 2021. 40, 58, 90, 117
- [143] Guido van Rossum. Python tutorial. Technical Report CS-R9526, Centrum voor Wiskunde en Informatica (CWI), Amsterdam, May 1995. 64

- [144] Ingeborg van Vugt. Using multi-layered networks to disclose books in the republic of letters. *Journal of Historical Network Research*, 1(1):25–51, October 2017. 53
- [145] Corinna Vehlow, Fabian Beck, and Daniel Weiskopf. The state of the art in visualizing group structures in graphs. In R. Borgo, F. Ganovelli, and I. Viola, editors, *Eurographics Conference on Visualization (EuroVis) - STARs*. The Eurographics Association, 2015. 96
- [146] VisMaster: Visual analytics — Mastering the information age. 110
- [147] Kiri Wagstaff, Claire Cardie, Seth Rogers, Stefan Schrödl, et al. Constrained k-means clustering with background knowledge. In *Icmi*, volume 1, pages 577–584, 2001. 95
- [148] Stanley Wasserman and Katherine Faust. *Social Network Analysis: Methods and Applications*. Cambridge University Press, November 1994. 13, 35, 41
- [149] Charles Wetherell. Historical Social Network Analysis. *Int Rev of Soc His*, 43(S6):125–144, December 1998. 9, 11, 23, 36, 37, 46, 51, 52, 62
- [150] Kai Xu, Alvitta Ottley, Conny Walchshofer, Marc Streit, Remco Chang, and John Wenskovitch. Survey on the Analysis of User Interactions and Visualization Provenance. *Computer Graphics Forum*, 39(3):757–783, June 2020. 65
- [151] Michelle X. Zhou. “Big picture”: Mixed-initiative visual analytics of big data. In *Proceedings of the 6th International Symposium on Visual Information Communication and Interaction*, VINCI ’13, page 120, New York, NY, USA, 2013. Association for Computing Machinery. 95