

# From Textual to Historical Networks: Social Relations in the Biographical Dictionary of Republican China

Cécile Armand, Christian Henriot

# ▶ To cite this version:

Cécile Armand, Christian Henriot. From Textual to Historical Networks: Social Relations in the Biographical Dictionary of Republican China. Journal of Historical Network Research, 2021, Beyond Guanxi: Chinese Historical Networks, 5 (1), pp.114-153. 10.25517/jhnr.v5i1.117. halshs-03213995

# HAL Id: halshs-03213995 https://shs.hal.science/halshs-03213995

Submitted on 8 Jun2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# Armand, Cécile Henriot, Christian

From Textual to Historical Networks: Social Relations in the Biographical Dictionary of Republican China







## Abstract

In this paper, we combine natural language processing (NLP) techniques and network analysis to do a systematic mapping of the individuals mentioned in the *Biographical Dictionary of Republican China,* in order to make its underlying structure explicit. We depart from previous studies in the distinction we make between the subject of a biography (bionode) and the individuals mentioned in a biography (object-node). We examine whether the bionodes form sociocentric networks based on shared attributes (provincial origin, education, etc.). Our major contribution consists in annotating the links between individuals in order to (1) question the assumption that word cooccurrences equate with actual relations; (2) define a more accurate classification of relationships among elites in republican China. We demonstrate that political and professional relations in this population outweigh the types of social ties opmonly accepted in the scholarship on modern China. We eventually levelop a method that can be applied to similar corpora in a critical and professional relations.

## 1 Introduction\*

A biographical dictionary is by definit n individuals whose inct biographic lives provide the back narhtives. The amount and -10-0. rt of the breadth and depth of full scope of information SI. ll sh biographical wor ich ι n individual is minutely described and in. e of social, political, economic events of the usually closely nten éff ntention to offer a macro-reading of times. Even with the bo **m** and as an explicit goal of the editors of the Biographical historical event (D) — the format of more or less short Dictionary of Re ub ...co *A*h unailed this ambition.<sup>1</sup> This holds especially biographical note ity bv true for the social real tion. and corracts that an individual had in the course of

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

<sup>\*</sup> Acknowledgements: We would like to acknowledge the contribution of Pierre Magistry (Aix-Marseille University) and Baptiste Blouin (Aix-Marseille University) who prepared the data used in this paper, advised on the methods, and organized the annotation workflow. This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No 788476).

Corresponding author: Cécile Armand, Aix-Marseille University, IrAsia, cecile.armand@gmail.com

<sup>1</sup> Howard L. Boorman and Richard C Howard, *Biographical Dictionary of Republican China* (New York: Columbia University Press, 1967), I, vii.

his/her life. In the condensed biographical notes that make up a dictionary, all the related historical actors are reduced to brief and often unique mentions in the body of the text. Moreover, due to the involvement of many contributors vs. a single author in a biographical work — such mentions are unsystematic with no apparent rationale as to the selection of the people included beyond the subjects of the biography.

In this paper, we propose a systematic mapping of all the individuals whose names appear in the biographical notes in order to make the networks underlying the BDRC explicit. We argue that the links between individuals in the biographical texts create an interlinked reference network of the biographical texts.<sup>2</sup> This network can in turn be used to examine to what degree the cooccurrence of names is constitutive of relations between individuals and whether these relations can be further qualified. We follow in the steps of previous experiments on the relevance of pervork analysis in exploring a world of word cooccurrences (named indivi uals) in biographical texts and n these named establishing the existence of actual soci etworks based entities.<sup>3</sup> Our approach, however, in he distinction fron that we make between the two d of 🔊 the dictionary: ferent nds latid those who were the subject of a biography and those were just mentioned in a biography. There is a considerable in alth of information on each group. The atter is real ed simply to name, except when they belonged to the grou b10 rgue that this distinction als. W raphed i diviď is necessary whe llysis. letworl

This paper to structured as follows. Section 1 describes how we build a reference network from projectly recognised person-to-person cooccurrences at the highest possible accuracy. In purceise, we built the network as a directed

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

<sup>2</sup> Christopher N. Warren et al., "Six Degrees of Francis Bacon: A Statistical Method for Reconstructing Large Historical Social Networks," *Digital Humanities Quarterly* 010, no. 3 (July 12, 2016).

Matje van de Camp and Antal van den Bosch, "The Socialist Network," Decision Support Systems 53, no. 4 (November 2012): 761–69; Matje van de Camp and Antal van den Bosch, "A Link to the Past: Constructing Historical Social Networks," in Proceedings of the 2Nd Workshop on Computational Approaches to Subjectivity and Sentiment Analysis, WASSA '11 (Stroudsburg, PA, USA: Association for Computational Linguistics, 2011), 61–69; Minna Tamper, Eero Hyvönen, and Petri Leskinen, "Visualizing and Analyzing Networks of Named Entities in Biographical Dictionaries for Digital Humanities Research," EasyChair Preprints, EasyChair Preprints (EasyChair, April 8, 2019); Pablo Aragon et al., "Biographical Social Networks on Wikipedia: A Cross-Cultural Study of Links That Made History," in Proceedings of the Eighth Annual International Symposium on Wikis and Open Collaboration - WikiSym '12 (Eighth Annual International Symposium, Linz, Austria: ACM Press, 2012), 1.

network in which the nodes are individuals, and when individual B is mentioned in the biography of A, we added a directed edge from A to B. We examine this network of cooccurrences (textual links) in the first section of the paper to study the underlying structure of the BDRC and propose an alternative reading of the dictionary and its population at a global scale. In section 2, we explore whether sociocentric subnetworks form on the basis of the specific attributes (provincial origins, education, etc.) that we extracted from the biographies. In section 3, we shift the point of observation from the study of networks of cooccurrences to that of social networks. To this end, we enriched the network of cooccurrences with annotations that qualified the nature of relations, in order to (1) distinguish mere cooccurrences from actual social relationships, and (2) build subnetworks based on the nature of relations, which we can compare with the corresponding subnetworks of attributes.

# 2 The BDRC as a network of forcurrences

In this section, we proceed in two steps. First, we describe the workflow for extracting named entities from the BDRC and building there erence network of cooccurrences. Second, we experiment with various methods (global and local metrics, pruning tables, clustering) in order to analyze its structure.

The BDRC consists of four you mes published etween 1967 and 1971<sup>4</sup> and has served genera It N as produced under the Cl ha l storià editorship of Ho ith contributions from about 100 different authors. Even if the bagraphi rent ally went through the hands of a small ur digital forensics has revealed is group of edite tı inconsistencies ab ulary sed f b describe individuals, positions, and co. ta. 589 individual biographies of unequal institutions. The four volumes length — from 57, to 2,000 tok ns - that feature "eminent Chinese" of the Republican period (1, 12-, 249) This constitutes a very small sample of the Chinese Republican elites, by any standard, and the criteria for selecting this group of historical figures have proved debatable. Yet, a great number of people (3,178) are mentioned in these 589 biographies, which come under three main categories: family members, authors, and other individuals.

As a rule, the biographers provided information on the family of the biographed individuals, usually starting the biographical notes with the genitors or those who raised them, if known. The latter may not be the same as

<sup>4</sup> Howard L Boorman and Richard C Howard, *Biographical Dictionary of Republican China* (New York: Columbia University Press, 1967-1971).

the genitors due to death or adoption. Each biography thus starts with birth and childhood, with a discussion of the family background. In most cases, the biographers cited only the name of the father, almost never that of the mother, even in the case of prominent families. The father and mother of people of humble origin were simply not named. Generally, at the end of each biography, there is also often, but not always, a list of the biographed character's direct family members, namely wife/wives and children. Most of the time, the information is sketchy, especially for the wives, except when they were themselves prominent figures, socially, intellectually, or politically (the Song sisters, Ding Ling, etc.).

Under the category of authors, we grouped all the individuals who wrote about the biographed character and whose works are cited in the biography. Some of them are people who were effectively in contact with the biographed person in the course of his/her life. This who case of former students who compiled and edited the writings of their sometimes next-ofrmer mentor, or kin (son-in-law, nephew, etc.). The majo of such w however, were ork produced ex post facto by individuarsw rson. Finally, nrel the group of authors also include iographers who historia and prof ssior es in he BDRC. wrote extensive monographs or papers on the major h

lentify all the n med entities in the text, we To index the cont and to processed the 589 þ. Stat ford CoreNLP.<sup>5</sup> Data extraction gra hie Wh t listed all the biographed persons and produced a raw th biographies. The high number of all the individuals entione τı worlth of data available in the BDRC beyond cooccurrences i line cat something that one may perceive through the 589 biograph sons, ρe at win fan to morace to its full extent. The number of conventional reading bic graphy varies greatly, from 124 for Jiang Jieshi individuals mentioned ea (蔣介石) to just one for Livizhi and Ma Buging. Based on this list, we built a directed network linking each biographed person (thereafter "bionode") to the individuals mentioned in his/her biography (thereafter "object-nodes").<sup>6</sup> The direction of arrows indicates whether an individual mentions (outgoing edges)

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

<sup>5 &</sup>lt;u>https://stanfordnlp.github.io/CoreNLP</u>

<sup>6</sup> We borrowed this distinction between bionodes and nodes from Henrike Rudolph, "Structures of Empowerment: A Network Exploration of the Collective Biographies of Women Activists in Twentieth- Century China," *Elites, Knowledge, and Power in Modern China: The Formation and Transformation of Elites in Modern China* (Aix en Provence, 2019), 25.

or is mentioned by (incoming edges) another individual.<sup>7</sup> We counted each pair of individuals only once, even if an individual is mentioned several times in a biography. Given the high number of cooccurrences and the nature of the BDRC — a collection of individual biographies — what is the relational structure of the dictionary? Is it merely an aggregate of multiple ego-networks or does it form a interconnected global network?

The network of cooccurrences generated from the extracted data comprises 3,254 nodes and 9,524 edges. It is made up of a total of 11 components, with one giant component (3,177 nodes with 9,377 edges) and ten disconnected components. All of the latter, except one, are in fact isolated ego-networks built around one single bionode. The exception is a small component consisting of two small ego-networks that centered on — Kang Cheng (Ida Kahn) and Shi Meiyu (Mary Stone) respectively. In fact, these two figures were the first Chinese women physicians trained in the United Traces at the end of the 19<sup>th</sup> century. Both received the help and support of the same woman missionary (Gertrude Howe), through whom their ego-network an interconnected and form a small component.<sup>8</sup> The other ego-networks Iduals with very that unfolded mos the Republican specific profiles: some had career v bè nization (Wei Zhuomin, era (Ye Changchi, Wang Ganchang) crin religious org Zheng Hefu), others were scientists who ht à ime abroad or even made most of their areer of tside of China (ph) sicist. Wu Jianxiong, Qian Xuesen) or whose profile fully caverg minstream population" in d fror the " the BDRC (Li Yiz di, Pel We z op **.**. cannot be aid that these individuals were poorly connected and what ye tchildre are just mentions of names in their said is that weither they nor the individuals named in biographies. W nt ca re related to any of the nodes in the main component. Why their biographie were these individual set stead if a ev emed quite off the mark? The probable answer lies between the dites' archion to have "representatives" of different hild inty of source materials. sectors of society and up a

Within the main component, how and to what extent are the biographies interconnected? How far do they rely on object-nodes to be interconnected? Do the latter contribute to the connectedness of the global network? Previous studies of biographical dictionaries have generally focused only on the

<sup>7</sup> The extraction process was done in R-studio. The data and the script are available on GitLab (https://gitlab.com/enpchina/brelations).

<sup>8</sup> Connie Anne Shemo, *The Chinese Medical Ministries of Kang Cheng and Shi Meiyu, 1872-1937: On a Cross-Cultural Frontier of Gender, Race, and Nation* (Bethlehem: Lehigh University Press, 2011).

biographed individuals and their relations.<sup>9</sup> In this paper, we move a step further and compare the entire network of cooccurrences with the network consisting only of bionodes. Once the object-nodes are excluded, the number of disconnected components increases to 16, with 15 isolated individuals. The network of bionodes, however, remains highly connected. If we compare the global metrics of the two networks, the density increases tenfold (entire network = 0.002 / bionode network = 0.028) and the clustering coefficient multiplies by a factor of four (entire network = 0.083 / bionode network = 0.337).<sup>10</sup> In both networks, there is still a high degree of connectedness given the considerable number of nodes in each.

Beyond global metrics, we use various centrality measures in order to examine the relative position of nodes and bionodes: edge count (number of neighbors), indegree (number of incoming edges), outdegree (number of outgoing edges), and betweenness centrality. The edge count displays a long-tailed distribution, with a minimum of one tie and a maximum of 367 (Jiang Jieshi).<sup>11</sup> We defined six thresholds as shown in Table 1:

Table 1 Edge tount

It can be observed from Table 1 that re mentioned only livid once in the BDRC. These r ach lii ked on to a single biography and de are y character. Although the vast majority are strictly related to the life of mis ve the BDC network by virtue of their becomes part of the m. n лe association with a massone bi presence makes sense only in relation ST. th e **I** dividuals all came from the object-node to this particul .. Th inc 882-1938), who happened to have the category, except fo (en vin er, lowest number of edges an ong he hio, odes. The group of 496 individuals with bie notes and 37 bionodes. The latter represents a 2 to 5 ties includes 159 er historical importance, more peripheral group of individuals individuals (scientists), or individuals with a short lifespan. Within the network

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

<sup>9</sup> Aragon et al., "Biographical Social Networks on Wikipedia"; van de Camp and van den Bosch, "The Socialist Network"; Tamper, Hyvönen, and Leskinen, "Visualizing and Analyzing Networks of Named Entities in Biographical Dictionaries for Digital Humanities Research."

<sup>10</sup> Technically, network density shows how densely the network is populated with edges. It is a value between 0 and 1. A network which contains no edges and solely isolated nodes has a density of 0. In contrast, the density of a clique is 1. The clustering coefficient measures the ratio of the number of edges between neighbors and the maximum number of edges that could possibly exist in the network.

<sup>11</sup> We do not include isolated nodes (13) in this table.

of only bionodes, a significant percentage (19 percent) of bionodes also have only this range of ties.

In the next group of 396 individuals with between 6 and 24 ties, we find mostly bionodes (345) who directly connected to one another. Yet, it also includes 51 object-nodes who appear in a good number of biographies. This category comprises a wide range of profiles and includes both Chinese and foreign elites.<sup>12</sup> Among the foreigners, we can distinguish two groups of people that are highly connected, although they are from different social circles -American philosophers or military advisers and Soviet/Comintern agents. In both cases, this is about foreign experts involved in Chinese politics. The Chinese in this group include political figures from the late imperial period (emperor Guangxu, Li Hongzhang, Zeng Guofan) and intellectual figures from the Republican period (Mao Dun, Zhang Junmai, Yan Huiqing, Ding Wenjiang), some with political connections. Zhang Zhang appears in the biographies of individuals of a similar generation with w on he had direct (Sheng Xuanhuai) or indirect contact (Shen Jiaben, Yan Hu), Nut lso he was a fig re of inspiration to younger people (Guo Bingwen

Outside bionodes, the two ind duals vith the hishest mber of edges ue to their direct or include Joseph Stalin (30) and Michael Bo odin indirect role in Chines politice 5. mostly in the biographies of lin is mentioned members of the Charles CP) (1) but also in relation with C. nm nist arty of the Guimingang (Jung Jieshi, Jiang Jingguo, Song major political f opelin, s Song Qingling, Liao Chengzhi, Deng Ziwen, etc.) and left-ying per-Yanda, etc.). Bo och in opt as is hard, mentioned in the biographies of CCP close interaction with communist leaders in China as the figures despite h letwork includes most of the major main Cominteri tal ve. Hi rep Thoughan to Jiang Jieshi, which reflects his Guomindang figu es, í m ın lati ns with these individuals in the mid-1920s in activity and direct Guangzhou.

Finally, the top three categories with more than 50 edges include exclusively bionodes. This group consists of the main figures of the Guomindang (Sun Zhongshan, Jiang Jieshi, Wang Jingwei, Hu Hanmin), two communist leaders (Mao Zedong, Zhou Enlai), all the main warlords (Duan Qirui, Zhang Zuolin,

<sup>12</sup> By order of importance: Zhang Zhidong (23), Li Hongzhang (18), general George Marshall (17), Zhang Junmai (Carson Chang), Mao Dun (16 each), John Dewey and Yan Huiqing (W.W. Yen) (15 each), Wu Chaoshu (C. C. Wu) and Zeng Guofan (14), Emperor Guangxu (13), Komintern agent Gregory Voitinsky, warlord Lu Yongxiang, and the Indian writer Rabindranath Tagore (12 each).

Feng Yuxiang, Wu Peifu, Zhang Xueliang), the transitional figure of Yuan Shikai, and two intellectuals (Liang Qichao, Hu Shi).

In order to better understand the underlying structure of the network of bionodes, we apply Marilyn Levine's method of dissecting the "hairball" by using pruning tables.<sup>13</sup> We take the edge count as the criterion for pruning the network. As shown on the pruning tables and the pruned graphs below, no significant change occurs until we remove the bionodes with 25 ties or more. From this point, the number of ties and nodes in the network is more than halved at each step. In the final step (Graph F), there only remains the four pivotal figures in the BDRC — Yuan Shikai, Sun Zhongshan, Jiang Jieshi, and Mao Zedong — i.e. the four major leaders that shaped the conventional narrative of the Republican period. This is just a preliminary exploration. More systematic utilization of the pruning method will help penetrate more deeply the structure of such complex networks.

**Table 2.** Pruning table with the range of edge counts, he ratio between remaining ties and nodes, and the remaining biorodes in he BDAC network at easy threshold.

Figure 1. Prunings of network graph. A. 574 50 bionodes with >=15 ties. C. 205 bionodes with> 118 bionod 54 bionodes with >= 35ong [M ties. F. 4 bionodes with 200 Tse-tu Yuan Shikai [Yuan Shih-k'ai], NO Ze Sun Zhongshan [Su **Jiang Jiesh** Chian ai-she . Color and size of node proan th ir edge

ole one count, however, do not take into The hierar hies bas d nature of the network. In order to refine our analysis, we account the direct we en norm grand outgoing edges. To this end, we need to disting ish b selected the individual with an our legree above 20 and calculated the ratio between indegree and edge court. Table 3 presents the results for the top 25 bionodes ranked by edge count (i.e., bionodes with more than 70 mentions in the BDRC). The ratio between indegree and edge count (last column) serves as a general indicator of how often an individual was mentioned in other biographies. It reinforces the impression that individuals with a very high rate of indegree are those who play an important role in the biographies of a large number of other individuals. Based on this ratio, we identified three major profiles: (1) individuals with a relative balance between incoming and outgoing edges (ratio  $\approx 50\%$ ); (2) "source" figures with a greater number of outgoing

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

<sup>13</sup> Marilyn Levine, "Post WWI Chinese Revolutionary Leaders in Europe," *Journal of Historical Network Research* xx, no. xx (n.d.): xx–xx.

edges (ratio <50%), which include mostly political or military leaders who controlled the chain of command at the top of institutions and were often the source of action or decision; (3) "referential" figures with a greater number of incoming edges (ratio >50%).

The referential figures are those who are mentioned far more frequently in other people's biographies than other people get mentioned in their biography. For instance, Duan Qirui (400%), Wang Jingwei (371%), Yan Xishan (321%), Li Yuanhong (300%), and Liu Bocheng (300%) were each mentioned three to four times more than they mentioned other individuals. In absolute number, the most frequently mentioned individuals include, in descending order, Jiang Jieshi, Sun Zhongshan, Yuan Shikai, Wang Jingwei, Feng Yuxiang, Mao Zedong, and Duan Qirui. As we elaborate later, these individuals were either solicited for advice or mentioned as contextual references (that is, they are not mentioned as part of an actual relationship, but as an area pent of historical context). This is a central question that we discuss in the third section. It points to a major shortcoming of previous studies that ofte sume that tooc arrences are the expression of social relationshi hallenge this the assumption through a close analy os, based on the s of the natu**r**e c atio computer-assisted annotations of biograp

#### **Fable 3.** The 25 b mode, with an edge count above 70

Note: Bionodes are one of a reds, count in a scending order with their respective indegree (number of inconting edges), where we can be used ing edges) and ratio indegree/edge count.

Edge count ey the importance of certain individuals, CO does not rely solely on the number of whose significa the n. twor ofi alternative way of measuring the ties. Betweenne are ty importance of obje es, 1 t just bonodes, who hold a central position in the t-no network.<sup>14</sup> Although ave relatively low number of neighbors, certain W individuals play a structuring role as mediators between different parts of the global network. For example, Paul Pelliot, the French archaeologist, with only 4 ties, holds a connecting position that joins two peripheral branches to the main component. If we remove him, the main component loses 95 nodes and 1,488 links. Similarly, the Qing official Zeng Guofan presents another intriguing case of a object-node with just 14 edges, who nevertheless connects 544 nodes and

<sup>14</sup> Betweenness centrality scores are particularly informative because they highlight the individuals who served as essential bridges ("brokers") between individuals and communities. Technically, betweenness centrality measures the number of shortest paths that travel through a node.

7,795 links. This is due to his connection to major bionodes such as Jiang Jieshi and Mao Zedong.

The list and range of bionodes with the highest scores of betweenness centrality remain very much aligned with the hierarchy defined by edge count. This group comprises 19 political, military, and intellectual figures who appear highly connected among themselves (123 edges). When their direct neighbors are included, these 19 individuals form a network of 870 nodes (including both the bionodes and the object-nodes) (26.7 percent of all nodes) and 6,226 edges (65.4% of all edges). The distinctive feature in the ranking by betweenness centrality is the emergence of a few prominent intellectuals such as Hu Shi, Liang Qichao, and Guo Moruo, who are placed 4<sup>th</sup>, 6<sup>th</sup>, and 8<sup>th</sup> respectively, well ahead of the major political and military figures who, by edge count, rank far above other nodes. It can be argued that the more versatile profiles of these three intellectuals who had a foot in various circles explain their position as eminent "brokers" in the BDRC.

The last method we apply to lustering. We ioni nlv seek to identify sub-communities of me bionodes. The algorithm (GLav) detected 23 comr in their size.<sup>15</sup> nitie vith gr The largest cluster (cluster 4) comprise 279 ties, but the 14 8 smallest "clusters" each const tor nt largest clusters are The st one blonode listed in Table 4. Fre t the number of nodes and ties it ch lus repo contains, list the n hd give it a label that best ost representative individu ls. or the names of its members. Some of them are describes its composition based to calfigures or clearly identified social groups. clearly centered on nave b nce, revolves a oundriamous military and political leaders Cluster 4, for it st (Jiang Jieshe, lot gsh n, Yi Shikai). Cluster 1 includes major un j intellectuals such as Hershi, Cai Yumpei, and Liang Qichao. Cluster 3 groups together major Communit leaders (Mao Zedong, Zhou Enlai). Cluster 6 represents the business circle, whereas cluster 9 connects several prominent scientists (physician Wu Liande, biochemist Wu Xian). Other clusters, however, exhibit more complex patterns with less obvious rationale for their grouping (such as clusters 7 and 16). Little can be said about these communities relying

<sup>15</sup> We used Cytoscape Glay algorithm (Fast-greedy), which relies on the greedy optimization of modularity score, with different corrections on edge density and cluster size. Previous studies have demonstrated its dramatic performance advantage in handling large networks. Gang Su et al., "GLay: Community Structure Analysis of Biological Networks," *Bioinformatics* 26, no. 24 (December 15, 2010): 3135–37.

From Textual to Historical Networks

solely upon the names of their members. In the next section, we address the question of whether they coalesce on the basis of specific attributes.

**Table 4.** The eight largest clusters of bionodes detected by the GLay algorithm.

 Note: For each cluster, the table reports the number of nodes and ties, the most representative individuals, and the label that best describes its composition, based on the names of its members.

**Figure 2.** Clustered network of bionodes Note: Size of node is proportionate to degree centrality.

To conclude the first section, our analysis points to the dual structure of the BDRC as a network of cooccurrences, featuring, on the one hand, several small ego-networks isolated from the main component, and on the other hand, a polycentric network (the main component) that present a relatively wellconnected group of biographies polarized by a limited number of prominent figures. In brief, because of the consid number of individuals mentioned in the BDRC, the cooccurrences of nar relatively wellally cons es ei connected network, with a giant omport ent of interconne d individuals. The pruning method based on the edge cour t and a f betweenness alv centrality have helped define distinct who are important ps a roi in varying degrees maintaining the to interconnectedness of the g beneath this massive "hair stru ure. I oreover ba. ball," clustering subgroups of more densely al veale connected indiv dua bgroups that contribute to the global DK entric networks based on specific structure of tl CO .es SOC attributes? This the ore question we appress in the next section.

# 3 Networks of an ribute.

After we identified the person on the BDRC, we used NLP techniques to retrieve a wide range of information related to these persons, such as institutions, positions, locations, events, etc. We propose to use this data as attributes to enrich the network of cooccurrences. Our analysis focuses on bionodes only, because the BDRC provides information on the provincial origin, education, and other details of all the biographed individuals, but such information is entirely missing for those mentioned in their biographies. Retrieving such information at this stage would require a huge amount of time, especially because for many Chinese people mentioned in the biographies of others, the BDRC gives only the initials of their given names.

In this section, we examine whether individuals who share common attributes tend to group together so as to form "sociocentric" networks in the BDRC.<sup>16</sup> We focus on five major attributes: provincial origin (well-studied by historians and often recognized as central in Chinese society<sup>17</sup>), military background (well represented in the BDRC<sup>18</sup>), education abroad (also a frequent feature in the BDRC population<sup>19</sup>), CCP affiliation (a self-contained, easily identifiable, and well-studied group<sup>20</sup>), and gender (women<sup>21</sup>). For each attribute, we built the network in two steps. First, we identified all the bionodes with the selected attribute, and second, we built an extended network that includes these bionodes and their first neighbors. It is these extended networks that we study below.

The hypothesis that provincial origin could provide the basis for specific networks did not pan out in general. For example, natives of Zhejiang are well represented in the BDRC, with 79 biographed individuals (13.9 percent of all the bionodes) who are connected to 434 other individuals in the extended network. Two isolated ego-networks around two contists (Qian Xuesen and Wang Ganchang) from Zhejiang have no direct indirect connection with the other Zhejiang natives. Furthermore, almost all t ovinces of Chin, are represented in the main component (1,040 no he ang network, 6 with Jiangsu (57) and Hunan (56) inces, followed as the i st repi

- 16 A "sociocentrice, two k is a network based on shared so ial attributes. This notion is borrowed from Tumpe. Hyvoren, and heski en, "Visualizing and Analyzing Networks of Named Entities in Biographical Disconaries or Digital Humanities Research."
- 17 William T Room, Hacker : Commerciana accuration in a Chinese City, 1796-1889 (Stanford, Calif.: Stanford University Press (1994); Bryn: Good man, Native Place, City, and Nation: Regional Networks and Identities in marchin, 1833-1957, Berkeley: University of California Press, 1995); Richard Belsky, pocalition at the Center, Native Place, Space, and Power in Late Imperial Beijing (Cambridge, Mass.: Harvied University Asia Center: Distributed by Harvard University Press, 2005).
- 18 Diana Lary, Region and Nation: The Kwangsi Clique in Chinese Politics, 1925-1937 (London; New York: Cambridge University Press, 1974); Jerome Ch'ên, The Military-Gentry Coalition: China under the Warlords (Toronto: University of Toronto-York University Joint Centre on Modern East Asia, 1979); Edward Allen McCord, The Power of the Gun: The Emergence of Modern Chinese Warlordism (Taipei: SMC Pub., 1997).
- 19 Y.C. Wang, *Chinese Intellectuals and the West*, *1872-1949* (Chapel Hill, University of North Carolina Press, 1966).
- 20 Marilyn Levine, *The Found Generation: Chinese Communists in Europe during the Twenties.* (Seattle, Wash.: University of Washington, 1993); Steve Smith and Taylor & Francis, *A Road Is Made: Communism in Shanghai* 1920-1927, 2018.
- 21 Gail Hershatter, Berkeley University of California, and Area Global and International Archive, *Women in China's Long Twentieth Century* (Berkeley: Global, Area, and International Archive : University of California Press, 2007); Barbara Mittler, Michael Hockx, and Joan Judge, eds., *A Space of Their Own: Women and the Periodical Press in China's Long Twentieth Century* (Cambridge: Cambridge University Press, forthcoming).

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

by Guangdong (48). Due to their importance in the Guomindang elites, Guangdong natives apparently offered the prospect of a tighter-knit group. Their network includes 376 bionodes and 4,050 edges, with 67 Cantonese nodes that possess a total of 239 edges. As in the case of the Zhejiang natives, however, a wide range of 23 provinces are represented, and the Cantonese account for only 18 percent of the total. We also observed that the network was actually made up of seven separate components, with six ego-networks built around specific personalities who had no link with each other.<sup>22</sup> In other words, there is little evidence of homophily by provincial origin among the Zhejiang and Guangdong natives in the BDRC.

Only one group based on the same provincial origin stands out in the BDRC: the natives of Hunan province form a large network of six components with 722 nodes and 4,908 edges. The presence of such prominent figures as Mao Zedong may have introduced a bias in terms of provincial origin and political affiliation (CCP). Yet, even after removing lao Zedong, the main component of the Hunan network still reveals a group of de sely connected Munanese whose pillars are also CCP members 1 Sh Pehuai, Peng esen Shuzhi, etc.). In this network, 72 the bid who have links odes an unans to Zhejiang (47), Jiangsu (40), and Cartonese (37 Altogether 21 tives provinces are represented. Yet the main the derable number of CCP members in this etwork 112. ionodes), of wi ich the Hunan natives claim a substantial share 7. One cal argu that e Hunanese clearly form a more homophilic network that is te ith politica liation. sec

The BDRC accelest high retime roumilitary figures. Individuals with any kind of military be itions in their careers form a population of 466 bionodes with 4,035 edges. The size of chils ne werk is an indication of the high level of cooccurrences in the bigraphie of these individuals. Of these individuals, however, only a much smaller purpler (43 bionodes with 94 edges) received a military education or held military positions continuously. Taking only these forty-three individuals into account, this smaller network displays a high density (0.322 vs. 0.074 in the larger network), although five individuals constitute isolated components with few links each (Xie Bingying, Huang Kecheng, Sun Lanfeng, Sun Liren, He Zhonghan). The major broker in the main

<sup>22</sup> The six individuals were Wei Zhuomin (Religious leader), Hu Die (actress), Xu Guangping (women writer), Li Fanggui (linguist), Dai Ailian (woman writer), and Luo Dengxian (labor activist).

<sup>23</sup> On the place of Hunanese in CCP networks, see Levine, "Post WWI Chinese Revolutionary Leaders in Europe."

component is Li Zongren, who has the highest degree and betweenness centrality. At the next level and almost at par, three lesser-known figures emerge: Liu Zhi, Xue Yue, and Yang Sen, who each form their own clusters. Only Xue Yue and Yang Sen are directly connected. In terms of indegree, Li Zongren still ranks first, followed by Duan Qirui and Bai Chongxi.

There is no clear evidence that graduation from the same academic institution was a strong connecting factor, except for the military officers who graduated from Japanese military academies. Graduates from Shinbun Gakko and Shikan Gakko, for instance, show a strong propensity to group together, but it needs to be established whether this is the sign of actual social relationships or an artifact of contextual mentions. The generational factor may reinforce the impact of cooccurrences. There is a clear overrepresentation of individuals born during the decade 1886-1896 who graduated mostly from the new military schools and academies established by the one or from the Japanese military academies between 1906 and 1920. generational group includes 25 Thbionodes (58%) and 60 edges (67%) and ad in brokers in the ounts for all the ma military-only network.

The Communists form a ve Identifiable tained set of self individuals who are included in the x because of their DR affiliation with the party (CCP). The network of members includes 865 nodes and 5,391 edges. It is e network that reflects the share of CCP ٦ Vt y lai members in the InRC (121 bring tes) and the high number of individuals tied to them (760 edges). The network or CCP members alone is made up of nine components, with eight solved individuals, and it has only four women in it. The outdegree a stribut on high ignts 20 individuals with more than 30 neighbors. We can define the burge ups. The first group with an outdegree above 35 includes vix a time includes (Mao Zedong, Zhou Enlai, Lin Biao, Zhu De, He Long, L. Daz, ao). These six individuals actually connect almost every CCP member (93/99). At the next level (outdegree between 25 and 34), we find a second group of seven tightly knit individuals (Liu Shaoqi, Ye Ting, Li Lisan, Zhang Guotao, Li Jishen, Chen Duxiu, Guo Moruo). Taken together, these thirteen individuals (Guo Moruo as an outlier) form the backbone of the CCP network in the BDRC. Betweenness centrality reveals a limited number of mediators (10), yet with large discrepancies between them. The whole network is clearly centered around Mao Zedong, who serves as the main broker (0.3), followed by Zhou Enlai (0.12), Zhou Yang (0.06), Ye Ting (0.05), and Lin Biao (0.04). Within the CCP network, as discussed above, the Hunanese lead the pack with 65 individuals, followed by natives of Zhejiang (52), Guangdong (44), Jiangsu (44), Hebei (27), Hubei (21), and Jiangxi (18). The CCP network includes

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

13 military officers, but these CCP military figures do not form a specific community.

Most of the biographed characters in the BDRC received a high level of education. 519 (88 percent) received a college degree. Among them, many had the opportunity to study abroad, which place them in the particular category of "returned students." These returned students are commonly grouped according to where they studied (country, university). In the BDRC, they form a population of 200 individuals who attended and graduated from a total of 343 different academic programs. Since the returned students established alumni associations or held events in China that brought together those who had studied in the same country or region, one can hypothesize that networks may have been built on this basis.<sup>24</sup>

The United States ranks first in the PORC with 70 returned students, followed closely by Japan with 67 individuals. Europe received the next largest batch, but the returned students from were distributed across several - TNC countries: United Kingdom (24), viet Union (8), and a host of other countries (6) We bu base n the country of nè study to examine to what extent the en country had m a` turn a propensity to connect with each oth k includes bionodes who were not educ or example, for the a 11 untry of refe ence American-trained Chin neighbors were not trained in the of the Зp cen ted their propensity to mingle United States. P. ork h nonstr de rates our observation.25 We found the with diverse co ba imuh orr same ratio am ng is only among the Japan-returned Eur The ower, indicating possibly a greater students that htly sli homophily. Yet, er factors may explain this higher level of homogeneity.

The American-trained students form one of the most interesting networks. In fact, it can be read as a miniature of the BDRC global network. On the one hand, a fair number of individuals (14) are not mentioned in any of the

<sup>24</sup> Stacey Bieler, "Patriots" or "Traitors"?: A History of American-Educated Chinese Students (New York: Routledge, 2003); Liu Xiaoqin, "Minguo Liumei Shetuan Yu Liumei Sheng de Shehui Wangluo -- Yi Chengzhihui Zhang Boling Fenxi Wei Zhongxin (Social networks and student associations in the United States in the republican era: A study of Zhang Boling and the Chengzhihui)," Huaqiao Huaren Lishi Yanjiu, no. 4 (2019): 88-95.

<sup>25</sup> Cécile Armand, "Foreign Clubs with Chinese Flavor: The Rotary Club of Shanghai and the Politics of Language," in *Knowledge*, *Power*, *and Networks: Elites in Transition in Modern China*, ed. Cécile Armand, Christian Henriot, and Huei-min Sun (Leiden: Brill, 2021).

biographies of their peers. Their networks do not intersect with the individuals in the main component, nor with the dyads formed by another five individuals. The 51 bionodes in the main component, however, tend to exhibit a higher level of connectedness. Hu Shi, one of China's leading intellectuals, serves as the main broker in this network (betweenness centrality = 0.15). He is connected to 21 peers through a total of 37 edges. His peers include mostly intellectuals but very few political figures. He is not linked to the second most important broker, Kong Xiangxi, an eminent figure in the dual world of business and politics. Kong is connected essentially to political figures, including his family relations (Song Ziwen, the Song sisters, and their father). The only intellectuals in his network are scholars with a foot in administration (Guo Bingwen, Ma Yinchu). Two other figures also play an important role in the network of American returnees, each in a different register: Song Ziwen, with a profile quite similar to his brother-inlaw, Kong Xiangxi, and Jiang Menglin, a multi-face d intellectual bridging the worlds of education, culture, and politics. While an exclusive pattern of homophily based on the country of edu cannot be established among the American-trained students, their netw ared a common ork ests that cultural, educational and linguistic ba have served for ckgrð d that establishing professional and politi e of their life. cal relat onships

### **Figure 3A.** Network of American-returned students (main component) Notes 5. 2 of des a proportionate a betweenness centrality.

de ts presents a very different structure. The networ ot ot e trality of Jiang Jieshi (by degree and A striking but ous featu ected al host all the individuals in the Japan betweenness m con. sur ly in his network, those who received military network, and more my or training. If we remove Jia g from the network, two other figures emerge: Li Liejun and Wang Jin wei, each at the center of a substantially different network. pontical figures, with very few military leaders. Li, Wang is connected mostly by contrast, reaches out to all the military leaders. If we enlarge the scope of observation to include all the bionodes connected to the Japan-returned students, we find a network made up largely of military figures trained in China or elsewhere. In other words, the single most important factor in the Japantrained individuals is less the country where they studied than the military education that they received there, which put them on a career path that connected them to a wider circle of military figures. If we compare it with the American returnees, one could say that the degree of heterophily based on the country of education is higher among the Japan returnees than their American counterparts.

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

#### Figure 3B. Network of Japan-returned students (main component) Note: Size of nodes is proportionate to betweenness centrality.

Very few women were selected for a biography in the BDRC. Altogether, they account for only 25 of the 589 biographies. What compelled the editors to select these women? Was it for their own profile and intrinsic importance, or because they were related to prominent men?<sup>26</sup> How do these women fit into the network structure of the BDRC? How do they contribute to male-dominated specific communities? In other words, can we identify a women's network in the BDRC?

The network of women includes 189 nodes and 990 edges, with four components: one main component and three ego or bi-ego networks. The bi-ego networks are those identified previously in the global analysis of the BDRC network of cooccurrences, namely the first two wonth physicians trained in the United States who graduated in 1896. Their nework is composed exclusively of foreigners — which reflects the fact that Borman focus an their period of education and training before the e other small retu to component revolves around Wu ianxiolog, a who made most hanphy of her career outside China, also arge nu eigners in her with a r of network.

Their limited number Women as such to not form an coherive net wrk. may be part of the explanation for the tack of more obvious networking. Five of them stand along and four wo district pairs. The main component mr . f m is made up of a long of very sha hunder of highly connected women — Song ing r), and Sing Meiling. Except for Ding Ling, who Qingling, Ding .ri+ stands apart, Song Dinging and Song Neiling - two sisters from an influential family and among the host prominent and politically active women of are note, only cted at the same level. Their marriage with Republican China <sub>5</sub>shan and Jiang Jieshi) placed them within a men of prominence (Sur larger network that included many main figures of the Republican period. Yet, even after removing the three main male figures (Mao, Jiang, and Sun), the centrality of the three women remains the same. On the other hand, they fail to connect directly with any significant number of women, and they even do not connect with each other. In the case of Song Meiling, it is through He Xiangning, the wife/widow of Liao Zhongkai (d. 1925), that she makes the connection with

<sup>26</sup> Henrike Rudolph, "Structures of Empowerment: A Network Exploration of Women Activists' Collective Biographies in Twentieth-Century China," in *Knowledge, Power, and Networks: Elites in Transition in Modern China*, ed. Cécile Armand, Christian Henriot, and Huei-min Sun (Leiden: Brill, 2021).

Chen Bijun (wife of Wang Jingwei) and Deng Yingchao (wife of Zhou Enlai). The remaining group of five women are even more tenuously connected.

Ding Ling's network branches out in two main directions: CCP members, including the men of letters in the party (Zhou Yang, Hu Feng) and literary figures (two other women writers, [Xie Wanying and Su Xuelin] and three male writers [Ye Shengcao, Lu Xun, Cao Yu]). Song Meiling's network is made up of powerful men that include all her direct and indirect next-of-kin (father, husband, brother-in-law, etc.), as well as military and political figures who served her husband or her more directly (Yu Hongjun, Wu Guozhen). There is no CCP figure in her network. Song Qingling, the older sister, presents a similar profile in terms of next-of-kin relations, but her network branches out to both Guomindang and CCP figures, which quite accurately reflects her positioning in Republican politics and in the People's Republic of China when she became the willing pawn of the CCP. In brief, the aranysis of the main component reveals four main profiles of women that may have erved as a guide for including them in the BDRC: scientists, artists, writers, and olitical women. Political women appear in the BDRC for their oversa narriage with e bu io di important political figures in Republican llectuals such as lhina. w as 🖪 Ding Ling appear only due to their own m

coon, he sudy of a tribute based networks reveals that To conclude this neither provincial wight r, et uca or any single attribute alone is sufficient on, i to constitute sin ific no su communities n the BDRC. It is only the a can account for the most densely multiple attri **utes** combination of the BDRC. For instance, this paper has connected clusters of b gr ph as h reasserted, in the win viewicus scholarship, the effect of Hunanese origin combined with CCP affiliation, that of study-abroad experience in Japan combined with military training and the conjunction of marriage, political influence, but also professional kills in the case of women. The patterns delineated in the study of attribute-based networks tend to support the hypothesis that there is more to the cooccurrence of names than the aggregated mentions of individuals in the various biographies. The repeated mentions suggest the existence of actual relationships. The network of cooccurrences, however, does not permit us to fully ascertain this. A more in-depth examination of the nature of the relationships is needed. This is the purpose of the next section.

## 4 Cooccurrences or relations?

Our analysis in the previous sections suggests the existence of two major categories of mentions associated with two dominant types of links:

textual/contextual references and actual social contacts. But it has proved difficult, if not impossible, to firmly establish the difference as long as we remain focused on the network of cooccurrences. In the first section, we made the hypothesis that some individuals may be mentioned as elements of historical context or were a source of inspiration for the biographed characters. We also pointed to historical characters from the past, such as Adam Smith or Zhuangzi, who could not actually have met with the individuals biographed in the BDRC. In the particular case of "referential" figures with high indegree centrality, we highlighted that they were most often sought for advice or served as contextual references. The close reading of their biographies further reinforces this impression. These observations led us to question the more general assumption that the cooccurrence of names in a biography could be systematically considered as the expression of a genuine social relationship.

nd explore in greater depth the In this section, we move a step further nature of the links between individuals ir selected sample of 36 biographies that we annotated manually (see Append Our approach is to examine the actual ground for the mention of man and qualify the e in` zen 1 relation between the named indiv duals. build networks ur ultim zoaľ from these annotations that better inflect historical in ionships, and not just textual cooccurrences. This section for ste st, we present the method for annotating the reactions in the biograph ies. Second, we analyze the results statistically. third, we build an as networks based on the anal ze vario (most significant extra tec a no ati which ompare with the attributensbased networks discussed in sicti

The first chillinge was to constitute a "representative" sample to annotate. We relied on two criteria crites we sale ted the individuals on the basis of the edge count. A high edge count was a sign that these individuals were involved in the widest and richost ange a possible relations. Second, we refined the sample to include the greatest possible variety of individual profiles in order to correct the bias produced by relying solely on the edge count. The selected biographies represent only 6 percent of the total number of biographies, but 18 percent of the total number of words in the BDRC.

In each biography, we focused only on the relations involving the biographed individual. For instance, in the biography of A, we annotated the relations between A and B and between A and C, but we discarded any potential relation between B and C. The selection of qualifying terms was based on the terms identified through close reading. We ensured that the manual annotations reflect only the language and the terms used in the text without adding any layer of interpretation or external knowledge, because the model for automatic annotations would ultimately rely only on the text itself and the particular

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

combinations of words through which different relations are expressed. As shown in Table 5, we classified the relationships into twelve categories: acquaintance, protégé, friendship, liaison, kinship, *tongxiang* (same native place), education (master/disciple, co-disciple), professional, political, military (military conflict), indirect (no direct relation: context, third-party reference, etc.), and neutral (uncategorized mentions). We were aware that our categories represented a wide range of choices, but previous experiments had convinced us that a narrow range of terms might produce results too broad for analysis.

For instance, a previous study based on a similar corpus — a biographical dictionary of Dutch socialists - had also attempted to qualify the relations between individuals through manual and automatic annotations. All the cases of cooccurrences, however, were considered as meaningful relations. Relying on sentiment analysis, the authors chose to classify the relations into three main types: positive, antagonistic, and neutral. The eventually found that the trilogy was too reductive, especially the neutral of hat regrouped too many cases to make it a significant marker.<sup>27</sup> This c teg tion may have seen relevant for relationships within a coherent and the BDRC presented mind us with a wider array of very dist ct profi ere interested in s. In out qualifying the relationships along a Cher of terms ed on the words used did not really pan in the text itself to describe the relations ca liaison, protégé, tongxiang). out due to the limited number of su h relationships Yet the process provide a a in the form of pre-determined rel...inai sche categories to be odel fð bmatic annotations.

#### 

For the annotation workflow, verselied on InCeption, a machine-assisted interactive annotation protocol.<sup>28</sup> and biography was annotated manually by a pair of annotators who work an independently, and then we curated together their respective biographies.<sup>29</sup> There were significant variations in the

<sup>27</sup> Matje van de Camp and Antal van den Bosch, "A Link to the Past: Constructing Historical Social Networks," in *Proceedings of the 2nd Workshop on Computational Approaches to Subjectivity and Sentiment Analysis*, WASSA '11 (Stroudsburg, PA, USA: Association for Computational Linguistics, 2011), 61–69.

<sup>28</sup> Jan-Christoph Klie et al., "The INCEpTION Platform: Machine-Assisted and Knowledge-Oriented Interactive Annotation," n.d., 5. See also: <u>https://inception-project.github.io/</u>

<sup>29</sup> The annotators included Cécile Armand (ENP-China, Aix-Marseille University), Guo Weiting (ENP-China, Aix-Marseille University), Christian Henriot (Aix-Marseille University), Jiang Jie (Shanghai Normal University), David Serfass (Inalco), Sun Huei-min (IMH, Academia Sinica).

annotations, mostly due to the inevitable propensity to "interpret" based on prior knowledge of the individuals mentioned in the biography and the difficulty to disentangle multifaceted relationships. Eventually, all the biographies went through the hands of a single curator who homogenized the annotations. These manual annotations were used to train a model for expanding the annotations automatically to the whole corpus in the future.

The annotation workflow produced a total of 3,227 annotated relations. As shown in table 6, the three most frequent types of relations represent 77 percent of the total. They include political relations (34 percent), indirect relations (26 percent), and professional relations (17 percent). Although military relations and kinship relations garner a good number of annotations, they represent only 7 percent and 6 percent of the total, respectively. All the other categories, including that of *tongxiang* (same native place), failed to produce a significant number of annotations. They were not sed to train our model. The high percentage of the top three types of relation indicates that the contributors who wrote the biographies centered primarily the work and sublic life of the individuals, especially their profession w, although a act significant share (one quarter) of l cooccurrences annota ons boit té (indirect relations), they do not preval in BDRO emaining 75% are the expressions of actual social relationship ctors.

# Table 6. Distribution of an otate brelations of a sample of 3 biographies from the BDRC.

An Irsis (PCA) on the data extracted from We ran a P Comp 'n in pa. This ugle PCA, we sought to construct relational the annotated quant correlations of specific relationships.<sup>30</sup> For profiles based on the m this PCA, we retained as active variables only the most frequent categories of relations: political (1994, 4%, indirect (823, 26%), professional (563, 17%), and military (239, 7%) and indered the minor relations as supplementary variables. The PCA graph plots all 35 biographed individuals and their relations in a two-dimensional space.<sup>31</sup> We selected the first two dimensions, as they best explain the variance among individuals and capture more than 70% of the information (46% and 24% on each dimension). The graph clearly contrasts individuals with many indirect and political relations on the right and those with few such relations on the left. The second dimension, which is primarily

<sup>30</sup> For conducting our PCA, we relied on the package "Factominer" in R Studio: <u>http://fac-tominer.free.fr/factomethods/hierarchical-clustering-on-principal-components.html</u>

<sup>31</sup> We chose to remove Guo Moruo whose profile was too specific, too different from the rest of the sample.

determined by military and professional relations, separates individuals with strong military and professional relations above the x axis from those with few such relations below it.

# **Figure 4.** PCA analysis of the 36 annotated biographies: A. Graph of variables. *Note: Black arrows: active variables; blue arrows: supplementary variables.* B. Hierarchical clustering of individuals.

Based on the PCA, we performed hierarchical clustering on all four dimensions. We observed that professional relations contribute the most to the partition at large (0.693), followed by indirect (0.664), military (0.633), and political (0.506) relations. The algorithm grouped individuals into four clusters based on their relational characteristics. In figure 3 below, each individual is color-coded by cluster. Cluster 4 on the right of the graph isolates two "big names" — Jiang Jieshi and Yuan Shikai — from all the others. What sets them apart is the wealth of professional relations (4.52), the weight of indirect relations (2.88), and finally, the rap political relations (2.41). What ofestional relations dominates cluster 2 is the conjunction strong rough which they (individuals who held a variety of positions he (rm) came into contact) and military confrontations, a enemies. This ies on a oth nationalists and cluster brings together almost all th mi itary lea communists, except for the Vu Peifu, and Zhang rlords ( elia Zuolin). Individuals in v few military and aracterned vre cl usi professional relation, and a stronger v eight of political and indirect relations. This cluster in luces 401 Iu Hannan, Zhou Enlai, and Sun Ma Zhongshan. Wi Perfu and Zhang Zuoh two major warlords of the post-Yuan be ig io th dest r, which suggests their respective Shikai era, al biographies define the n ch mure b their political relations than by their professional or military on

In the last step, we we are reased of networks based on annotated relations. We constructed one network for each of the most significant types of relationships. We expected annotated relations to determine with greater confidence the reality of a relationship and to delineate more precisely the nature of the links between the biographies and between the individuals. The annotated relations may not alter the overall structure of the co-occurrence networks fundamentally, nor the dominant position of high-profile individuals such as Jiang Jieshi or Yuan Shikai — we could only compare them with the co-occurrence network of 36 annotated biographies — but they substantiated the hypothesis that a relationship could, in fact, reveal different configurations of proximities and interactions.

The category of "indirect relations" lumps together different types of relations (contextual mention, source of inspiration, connection through a third party, etc.). In future analyses, we will refine this category to separate the purely contextual relations from the other types of relations, and we will adjust the annotation workflow accordingly. Still, the network of indirect relations confirms the weight of Sun Zhongshan, Mao Zedong, and Jiang Jieshi (in descending order of betweenness centrality) in purely contextual mentions. In this function, they have no direct relationship with the person in whose biography they appear. Most such mentions come under formulas such as "When A came to power," "At the time of A's death," "Under A's regime," etc. Moreover, the clustering of this network produced very interesting subgroups centered on specific individuals who shared common characteristics. One of these clusters is dominated by Jiang Jieshi and Sun Zhongshan, a second by the northern military leaders *cum* warlords, another or by CCP leaders, whereas two less densely connected groups revolved abound intellectual figures (Hu Shi, Li Shizeng, Li Dazhao) and late Qing re ionary activist Huang Xing, Liang lu Qichao, Zhang Binlin). These indirect relation not denote an even if t actual social relationship, provide a sort of vance in terms of index" of contextual mentions.

Political relations are the most prominent ire **N** onnection between individuals. They delineate newly two separate wo ds: CP leaders on one side and all the other main istorica. figures on he other side. They highlight the centrality of Jiarg noi gshan T SU .7 petweenness centrality) in the es hile Mao Zenory and Thou Enlai are significant only within whole network on of cempan out very much into non-communist the CCP sub-ne VOI. circles, except h Epilar Vet the is nainly due to Zhou's later career as premier of the People's Pep bl. h. t considerably extended his contacts internationally. In ontrict, Darhao, more of a secondary figure who was executed in 1927 at age deptoad network of relations across political lines although he was executed in 1927 at age 38. Yuan Shikai also built an extensive and very diverse network of political relations with political and mostly nonmilitary figures, including a good number of Qing officials, but also major opponents such as Huang Xing, Song Jiaoren, or allies-turned opponents like Liang Qichao and, of course, Sun Zhongshan.

The exploration of professional relations redraws the previous configurations and leads to a very different network structure. Jiang Jieshi and Yuan Shikai are the two most central figures in this network. By placing emphasis on their professional relations, the resulting network gives more weight to their careers in government and the army. Their respective networks, however, diverge significantly. Yuan Shikai is connected to all the military leaders of the early republican period, except Wu Peifu. Many were his protégés. His professional relations include only a few political or intellectual figures such as Sun Zhongshan, Hu Hanmin, Wang Jingwei, or Cai Yuanpei. Yuan shares these figures with Jiang Jieshi's professional network. They are the connecting points, though indirectly, between Jiang and Yuan. Jiang's network includes two major military figures on the nationalist side — He Yingqin and Zhang Fakui — and a host of less central individuals. What distinguishes the network based on professional relations is the greater diversity one can see in the multiple sub-networks built around individuals (Hu Shi, Wu Peifu, Zhang Xueliang) who are only remotely connected to Jiang and Yuan. This observation also holds true for CCP leaders. Mao Zedong's professional relations place him as a secondary figure in the network, and he connects mostly to the triad formed by Zhou Enlai, Liu Shaoqi, and Zhu De. Yet again, this reflects the nature of post-1949 relations.

Military relations provide a good to compare the networks of cooccurrences and annotations. We comp e the network based on annotated military relations with the corresponding work of conccurrences (based on military attributes) that we buil down to the 36 ctio annotated biographies. In the ne the number of work D sed otati on ( nodes and edges decreases greatly, Nom 2 nodes a 11 edges to 101 nodes and 225 edges (Table 7. This is mostly a have focused only **r**to ct 1 on the relations involving the vione canth, the lower number of les. More signi ustering coel edges and the low that military operations uggest cient played but a light d aphies of military elites. The BDRC a range of relations instead. The power emphasizes the vo. vement in wic. and prestige re m t. eir i litery deeds — victories or defeats on the 1tec of influence (political negotiations, battlefield her suurce cts, lecolomend. professional cont

#### Table 7. Comparative analysis, military networks (global metrics).

Who are the most important players in these two networks? In the network of cooccurrences based on military attributes, betweenness centrality places Jiang Jieshi as the domineering broker, connecting a host of second-rank military and non-military actors (including Sun Zhongshan). In the network of annotations, Jiang remains central, but at the same level with other military leaders. Moreover, non-military actors are relegated to a secondary position. Quite interestingly, Yuan Shikai becomes a minor and more marginal broker with a very limited range of military relations.

**Figure 5.** Military networks: A. Based on attributes; B. Based on annotated relations (sample of 36 biographies). Note: Size of nodes is proportionate to betweenness centrality.

To conclude this section, the analysis based on annotated relations in biographies reveals that individual mentions in the BDRC are not just cooccurrences – i.e., names connected through textual links – but also refer to historical actors who actually came into contact in the course of their life. There is, however, a substantial number of cases of indirect relations (25% of all annotations) and even of dropped names that argue for the need to exercise greater caution about considering any cooccurrence as the expression of an actual historical relation. On the other hand, the BDRC features a great variety of actual links within and between the biographies. The relational patterns we delineated through PCA and SNA also demonstrated the intermingling of multiple relationships among individuals, which complicates the previous typologies based on sentiment analysis. Clearly, annotations provide a necessary and efficient way to add historical substance to the analysis of networks simply based on cooccurrences.

## 5 Concluding Remarks

Reading the BDRC through the ns of s vial netv may seem like ana putting old wine in a new bottle. The major biase this work have been established in previous academic revie ont n is not to reassert these biases in terms a context (pr blems of same ing and representativeness) but rather to uncover the underlying book, namely the hidden e of the ructu relations between h dh idenlato 1a hier and b tween the individuals therein. This allows us t (1) extract a " of ctiv portrait" of the entire population based on their indivi ributes); (2) assess whether and how far al .ist. (a constructive of petworks; and (3) propose an approach that cooccurrences v en more accurately. defines and qualifies relationships nu

We demonstrated that in the DRC, political and professional relations far outweigh the "three sames" (native place, education, trade/business) commonly accepted in the historical literature. This goes against the grain, but we argue that this is not due solely to the nature of the elites selected in the

<sup>32</sup> J. K. Fairbank, "Biographical Dictionary of Republican China. Volume I, Howard L. Boorman, Editor. Richard C. Howard, Associate Editor. (New York: Columbia University Press. 1967)," *The American Historical Review* 73, no. 2 (December 1, 1967): 565–66 ; Lucien Bianco, "Howard L. Boorman, editor, Richard C. Howard, associate editor, Biographical Dictionary of Republican China, vol. 1.," *Annales* 23, no. 5 (1968): 1133–35; D. C. Twitchett, review of *Review of Biographical Dictionary of Republican China*, by Howard L. Boorman and Richard C. Howard, *Political Science Quarterly* 84, no. 4 (1969): 650–52.

BDRC.<sup>33</sup> The persistence of these types of relations, especially native place ties, was undeniable in Republican China. Chinese society, however, departed increasingly from the patterns studied for late imperial China. There was a flurry of new types of social organization that offered the possibility for individuals to get involved in multiple groups and networks. Political parties are a prime example of a completely novel type of organization, but professional or cultural associations also provided numerous arenas based on non-partisan grounds. While it is difficult to escape the conventional categorization (politician, merchant, military, etc.) that historians use to define elites — some individuals do fit in such categories — the relational profiles we have revealed challenge the relevance of such narrow categorizations for the Republican elites. The complex web of relations in which the individuals in the BDRC were enmeshed cut across such categories and their careers often followed more than one path, sometimes in parallel.

From a methodological perspective have also demonstrated that much original knowledge can be gained from bi phies through this approach. On the one hand, network analysis reg ind self-contained ts the biographies and open pathways through It falls short of the BD mos the creating a global narrative, but K revisits and function of a n biographical dictionary. The relations v nas /eil open a new way of e BDRC in a digital version, navigating through t such as we are planning to release. On the other hand, allowed us to put the BDRC to a etwork a alvsi truth test. It re ea his wo etween a largely densely 510 Ing a group of leading elites in the on ponent' connected "main fea. f .d political and m itar nd ffe ent subsets of individuals, and even ind vi Juals various disconn

The BDRC contains a page but finite volume of biographical data. We processed only what was to the toot, with no addition of external information. Yet the implementation of data mining and annotation methods based on NLP, followed by exploration with network analysis and PCA, allowed us to identify and trace patterns of relationships, to question some assumptions about the types of relations among this composite elite population, and to breathe life into the stock of knowledge contained in the BDRC. The set of manual annotations of just a small sample of the biographies proved highly instructive, as a learning experience for historians to "define" the nature of relations in a text. People are multifaceted and it proved very challenging to reduce the nature of relations to a single word. This also demonstrates the need for human intervention at every

33 Boorman and Howard, Biographical Dictionary of Republican China, viii.

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

From Textual to Historical Networks

step, from close reading to defining terms and expressions, to annotating the text.

This is an on-going experiment, but we believe that the models that we trained on our set of manual annotations offer a enormous potential for moving toward automatic annotations of large biographical corpora. It paves the way for the exploration of similar corpora such as the main English language dictionaries and the numerous Chinese language works published both before and after 1949.

# 6 Appendix: List of 36 Biographies

Name	Chinese	Wade-Giles
Zhang Fakui	張發奎	Chang Fa <b>t</b> uei
Zhang Xueliang	張學良	Chang Hsueh-liang
Zhang Xun	張勳	Chang Hsün
Zhang Binglin	章炳麟	Chang Ping In
Zhang Zuolin	張作霖	Chang Tro-in
Chen Duxiu	陳獨秀	Ch'en Tu-hsiu
Jiang Jieshi	蒋介石	Chiong Kai-shek
Zhou Enlai	尼恩外	chou En-lei
Zhou Shuren	月樹人	Chou Shujjen
Zhu De	朱蕊	Chu Ten
Feng Yuxiang	馬云谷	Feng Yü-hsiang
He Long	賀龍	Ho Lung
He Yingqin	仙德女	Ho Ying-ch'in
Xu Shichang	余世号	Hsü Shih-ch'ang
Hu Hanmin	助展民	Hu Han-min
Hu Shi	胡適	Hu Shih
Huang Xing	黃興	Huang Hsing
Kong Xiangxi	孔祥熙	H. H. K'ung
Guo Moruo	郭沫若	Кио Мо-јо
Li Liejun	李烈鈞	Li Lieh-chün
Li Shizeng	李石曾	Li Shih-tseng
Li Dazhao	李大釗	Li Ta-chao
Li Zongren	李宗仁	Li Tsung-jen
Li Yuanhong	黎元洪	Li Yuan-hung
Liang Qichao	梁啓超	Liang Ch'i-ch'ao

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx.

Liu Shaoqi	劉少奇	Liu Shao-ch'i
Mao Zedong	毛澤東	Mao Tse-tung
Song Ziwen	宋子文	T. V. Soong
Sun Zhongshan	孫中山	Sun Yat-sen
Cai Yuanpei	蔡元培	Ts'ai Yuan-p'ei
Duan Qirui	段祺瑞	Tuan Ch'i-jui
Wang Jingwei	汪精衛	Wang Ching-wei
Wu Peifu	吳 <b>佩孚</b>	Wu P'ei-fu
Ye Ting	葉挺	Yeh T'ing
Yan Xishan	閻錫山	Yen Hsi-shan
Yuan Shikai	袁世凱	Yuan Shih-k'ai

## References

Aragon, Pablo, David Laniado, Arabea, Kalunterunnar, an Uran, Volkovich. "Biographical Social Networks on Wikipedia: A Cross Cultural Study of Links That Made History." In *Proceedings of the Eighth Annual International Symposium on Wikis and Open Collaboration - WikiSys: 12*, 12, 12, 10, 20, 2000, 2012.

Armand, Cécile, a progra Clubs with Cainese Navo: The Rotary Club of Shanghai and the Polyics of Language "In *Knowledge, Power, and Networks: Elites in Transit on I. Molern chica*, edited by Cécile Armand, Christian Henriot, and Huei-Line Suc. Luden: Bill, 2021.

Belsky, Richard. Localities in the Senar: Mative Place, Space, and Power in Late Imperial Beijing. Cambridge, Mass.: Howard University Asia Center, 2005.

Bianco, Lucien. "Howard L. Boorman, editor, Richard C. Howard, associate editor, Biographical Dictionary of Republican China, vol. 1." *Annales* 23, no. 5 (1968): 1133–35.

Bieler, Stacey. "Patriots" or "Traitors"?: A History of American-Educated Chinese Students. New York: Routledge, 2003.

Blouin, Baptiste, Pierre Magistry, and Nora Van den Bosch. "Creating Biographical Networks from Chinese and English Wikipedia." *JHNR*, n.d.

Boorman, Howard L, and Richard C Howard. *Biographical Dictionary of Republican China*. New York: Columbia University Press, 1967.

eISSN: 2535-8863 DOI: 10.25517/jhnr.xxxxx.xx. Journal of Historical Network Research No. x • 202x • xx-xx

7

#### From Textual to Historical Networks

Camp, Matje van de, and Antal van den Bosch. "A Link to the Past: Constructing Historical Social Networks." In *Proceedings of the 2Nd Workshop on Computational Approaches to Subjectivity and Sentiment Analysis*, 61–69. WASSA '11. Stroudsburg, PA, USA: Association for Computational Linguistics, 2011.

———. "The Socialist Network." *Decision Support Systems* 53, no. 4 (November 2012): 761–69.

Ch'ên, Jerome. *The Military-Gentry Coalition: China under the Warlords*. Toronto: University of Toronto-York University Joint Centre on Modern East Asia, 1979.

Fairbank, J. K. "Biographical Dictionary of Republican China. Volume I, Ai-Ch'Ü. Howard L. Boorman, Editor. Richard C. Howard, Associate Editor. (New York: Columbia University Press. 1967. Pp. Xv, 483. \$20.00)." *The American Historical Review* 73, no. 2 (December 1, <u>19</u>67): 56–66.

Goodman, Bryna. *Native Place, City, and Nation : Regional Networks and Identities in Shanghai, 1853-1937.* Berkeley: University of California Press 1995.

Hershatter, Gail, Berkeley University of California, and Area Clobal and International Archive. *Women in China's Long Twentieth Century*. Berkeley: Global, Area, and International Archive : University of California Press, 2007.

Klie, Jan-Christoph, Muhae Bugert, Leto Bordlosa, Richard Eckart de Castilho, and Iryna Gurev (C.) The INCEPTION Platform: Machine-Assisted and Knowledge-Ori integrative Annatation," n.d., 5.

Lary, Diana. *Rept. 2 and Neuron*: The Kounnesi Clique in Chinese Politics, 1925-1937. London; New York. Camp and the timers ty Press, 1974.

Levine, Marilyn. "I st w VI chin, e Revolutionary Leaders in Europe." *Jour*nal of Historical Network Revealed xx, no. xx (n.d.): xx–xx.

Levine, Marilyn Avra. *The Found Generation: Chinese Communists in Europe during the Twenties.* Seattle, Wash.: University of Washington, 1993.

Liu Xiaoqin. "Minguo Liumei Shetuan Yu Liumei Sheng de Shehui Wangluo --Yi Chengzhihui Zhang Boling Fenxi Wei Zhongxin (Les Réseaux Sociaux et Les Clus d'étudiants Chinois Aux Etats-Unis Sous La République : Une Étude Centrée Sur La Chengzhihui et Zhang Boling)." *Huaqiao Huaren Lishi Yanjiu*, no. 4 (2019): 88-95.

McCord, Edward Allen. *The Power of the Gun: The Emergence of Modern Chinese Warlordism*. Taipei: SMC Pub., 1997.

Mittler, Barbara, Michael Hockx, and Joan Judge, eds. *A Space of Their Own: Women and the Periodical Press in China's Long Twentieth Century*. Cambridge: Cambridge University Press, forthcoming.

Rowe, William T. *Hankow: Commerce and Society in a Chinese City*, 1796-1889. Stanford, Calif.: Stanford University Press, 1984.

Rudolf, Henrike. "Structures of Empowerment: A Network Exploration of the Collective Biographies of Women Activists in Twentieth- Century China," 25. Aix en Provence, 2019.

———. "Structures of Empowerment: A Network Exploration of Women Activists' Collective Biographies in Twentieth-Century China." In *Knowledge*, *Power*, and Networks: Elites in Transition in Modern China, edited by Cécile Armand, Christian Henriot, and Huei-min Sun. Leider Brill, 2021.

Shemo, Connie Anne. *The Chinese Medical Ministries of Kang Cheng and Shi Meiyu, 1872-1937: On a Cross-Cultural Frontier of Gender, Rice, and Nation.* Bethlehem: Lehigh University Press, 2011.

Smith, Steve and Taylor & Francis. *XRoadels Made: Communisment Shanghai* 1920-1927, Honolulu, University of Hawar Press, 2018

Su, Gang, Allan Kuchusky, John H. Morris, David J. States, and Fan Meng. "GLay: Communa, Surgure Analysis of Biological Networks." *Bioinformatics* 26, no. 24 (December 15, 2017): 7135–57

Tamper, MinnanSeroLyyvinen, and OttoLeskinen. "Visualizing and Analyzing Networks of Numer Exclusion Liographical Dictionaries for Digital Humanities Research" EasyConir Liepting. EasyChair Preprints. EasyChair, April 8, 2019.

Twitchett, D. C. Review of *Review of Biographical Dictionary of Republican China*, by Howard L. Boorman and Richard C. Howard. *Political Science Quarterly* 84, no. 4 (1969): 650–52.

Wang, Y.C. *Chinese Intellectuals and the West*, 1872-1949. Chapel Hill, University of North Carolina Press, 1966.

Warren, Christopher N., Daniel Shore, Jessica Otis, Lawrence Wang, Mike Finegold, and Cosma Shalizi. "Six Degrees of Francis Bacon: A Statistical Method for Reconstructing Large Historical Social Networks." *Digital Humanities Quarterly* 010, no. 3 (July 12, 2016).

31