

Complicating the Social Networks for Better Storytelling: An Empirical Study of Chinese Historical Text and Novel

Chenhan Zhang^{ID}, *Student Member, IEEE*, Qingpeng Zhang^{ID}, *Senior Member, IEEE*,
Shui Yu^{ID}, *Senior Member, IEEE*, James J. Q. Yu^{ID}, *Senior Member, IEEE*, and Xiaozhuang Song

Abstract—Digital humanities is an important subject because it enables developments in history, literature, and films. In this article, we perform an empirical study of a Chinese historical text, *Records of the Three Kingdoms (Records)*, and a historical novel of the same story, *Romance of the Three Kingdoms (Romance)*. We employ deep-learning-based natural language processing (NLP) techniques to extract characters and their relationships. The adopted NLP approach can extract 93% and 91% characters that appeared in the two books, respectively. Then, we characterize the social networks and sentiments of the main characters in the historical text and the historical novel. We find that the social network in *Romance* is more complex and dynamic than that of *Records*, and the influence of the main characters differs. These findings shed light on the different styles of storytelling in the two literary genres and how the historical novel complicates the social networks of characters to enrich the literariness of the story.

Index Terms—Complex networks, social network analysis (SNA), text mining.

I. INTRODUCTION

DIGITAL humanities is a transdisciplinary subject between information technologies and humanities, such as literary classics. For instance, Google makes a contribution to digital humanities by promoting the “Google Books Library Project,” which includes millions of paper books scanned into electronic text [1]. Digital text is easier for researchers to explore than printed books since the development of information technology has provided numerous effective tools [2]. In the past decade, overwhelming data science techniques have advanced the research on digital humanities; thus, components can be extracted and analyzed from the literature.

Manuscript received September 7, 2020; revised December 28, 2020 and February 6, 2021; accepted February 20, 2021. (*Corresponding author: James J. Q. Yu.*)

Chenhan Zhang is with the Guangdong Provincial Key Laboratory of Brain-Inspired Intelligent Computation, Department of Computer Science and Engineering, Southern University of Science and Technology, Shenzhen 518055, China, and also with the University of Technology Sydney, Ultimo, NSW 2007, Australia (e-mail: zhangch@mail.sustech.edu.cn).

Qingpeng Zhang is with the School of Data Science, City University of Hong Kong, Kowloon, Hong Kong (e-mail: qingpeng.zhang@cityu.edu.hk).

Shui Yu is with the School of Computer Science, University of Technology Sydney, Ultimo, NSW 2007, Australia (e-mail: shui.yu@uts.edu.au).

James J. Q. Yu and Xiaozhuang Song are with the Department of Computer Science and Technology, Southern University of Science and Technology, Shenzhen 518055, China (e-mail: yujq3@sustech.edu.cn).

Digital Object Identifier 10.1109/TCSS.2021.3061702

A review of previous research reveals that some areas in digital humanities remain unexplored. First, mainstream studies are limited to the humanities works on the background of the Western world [3]. It is both interesting and constructive to investigate humanities works with oriental backgrounds. Second, only a few comparative studies on the literature with different styles of the same story are conducted [4]. In particular, previous studies focused more on longitudinal studies, in which researchers usually adopt a story series, such as *Harry Potter Books 1–7*, as the object of study [5]. A potential research interest about the same story that discovers varied features (narrative levels, characters, and events) or sentiments can arise from different literature, which may be driven by literary genres or authors’ opinion, among others. Third, network study is essential for the social network of a story and any network that possesses a topological structure, which can help gain an insight into the story’s characters based on its narration [6].

To fill the gap, this article introduces a social network and sentimental analysis work on two different texts of one of the most famous Chinese story, the Three Kingdoms. In particular, we leverage the state-of-the-art natural language processing (NLP)-based model to extract the social networks in the narratives of two books. In particular, the adopted NLP approach can extract the majority of characters (more than 90%) in both two books. Thereafter, a series of descriptive statistical analyses on the extracted networks is conducted, and we discover the homogeneity and heterogeneity in terms of topological features in these networks. In addition, we adopt the sentimental analysis to compare the evaluations on some of the main characters. The results reveal that the social network is more complicated in the narrative of the novel (*Romance*) than that of the historical text (*Records*). Consequently, it can be concluded that the literariness of stories has a tight relationship with the complexity of the social networks they entail.

The main contribution of this article is as follows.

- 1) We integrate the latest NLP and network science techniques to extract and analyze the social networks of historical text and novel.
- 2) We depict the difference in the dynamic social networks of the *Records* and the *Romance*, the classic historical text, and novel of the same story.

TABLE I
METRICS OF SNA

Connections	Homophily [19], Multiplexity [20], Reciprocity [21], Network Closure [22], Propinquity
Distributions	Bridge, Centrality [23], Density, Distance [24], Structural holes [25], Tie Strength
Segmentation	Cliques, Clustering Coefficient, Cohesion

- 3) A series of comprehensive case studies is performed, and we find that the historical novel complicates the social networks of characters to enrich the literariness of the story.

The remainder of this article is organized as follows. In Section II, the backgrounds of text mining and social network analysis (SNA) studies are presented. Section III elaborates on the network extraction approach. We perform a series of empirical studies in Section IV to demonstrate the thesis of this work. Finally, this article is concluded in Section V with a summary of potential future studies.

II. LITERATURE REVIEW

A. Social Network Analysis

Previous studies have demonstrated the importance of network analysis in different domains, such as complex network in supply chains [7] and risk identification in electric industries [8]. For networks that possess a social structure, SNA can be used to study social structures by analyzing the relationships, communities, and activities through topology graph theory [9], [10]. Initially, the study of SNA focuses on the network that actually exists, such as mobile social networks [11] (Table I categorizes the metrics of SNA according to various features of social network). The development of NLP enables the extraction of the latent social network in narratives, such as literary text and news text (narrative network analysis). Recently, studies focus on narrative networks in literary works such as novels. For example, studies on Harry Potter find salient, small-world, and scale-free features in its social network [12], [13], and these features reveal that the story is penetrated by compact character relationships. Gessey-Jones *et al.* [14] investigated the *A Song of Ice and Fire*, one of the most popular epic fantasy novel series, and they studied the distribution of time intervals between significant deaths of characters measured regarding the in-story timeline, which is consistent with power-law distributions. In the context of the spreading of COVID-19, Wang *et al.* [15] applied SNA to the information extracted from COVID-19 patients and explored the epidemic stages of emerging infectious disease.

Furthermore, there are many community detection algorithms have been proposed for topology analysis. For example, Seo *et al.* [16] proposed a general model that can exploit the deep structure of fiction novels based on graph topology. Bu *et al.* [17] proposed a community detection solution by using the graph k -means technique, and the locally Pareto-optimal community structure in social media networks can be detected effectively. Cao *et al.* [18] proposed a dynamic game model for prosumer–community groups’ detection in smart grids.

TABLE II
POS TAGS [29]

Tag	Meaning
CC	Coordinating conjunction
CD	Cardinal number
DT	Determiner
EX	Existential “there”
FW	Foreign word
IN	Preposition of subordinating conjunction
JJ	Adjective
JJR	Adjective, comparative
JJS	Adjective, superlative
LS	List item marker
MD	Modal
NN	Noun, singular or mass
NNS	Noun, plural
NNP	Proper noun, singular
NNPS	Proper noun, plural
TO	Infinitive marker “to”
UH	Interjection
VB	Verb, base form
VBD	Verb, past tense
VBG	Verb, gerund or present participle
VBN	Verb, past participle
VBP	Verb, non-third person singular present
VBZ	Verb, third person singular present
WDT	Wh-determiner
WP	Wh-pronoun
WP\$	Possessive wh-pronoun
WRB	Wh-adverb
XNOT	Not and n’t

B. Text Mining and NLP

1) *Named Entity Recognition*: Named entity recognition (NER) is among the core tasks in NLP. In story-oriented text mining, the NER task requires that the characters and sentiment representatives are treated as entities and can be identified in the texts [26]. A bulk of computational linguistic-based NER methods are developed, which plays vital roles in NER tasks, especially the token-level tasks [27], [28].

2) *Part-of-Speech Tagging*: Part-of-speech (POS) tagging is the process of tagging a token (a word) for a particular part of speech according to its context [34]. Table II shows each type of tag with its corresponding meaning. POS tagging helps from related grammatical rules for different language patterns.

C. Deep Learning-Based NLP Models

To extract the social network of a story, characters and their connections among one another must be identified. The distribution of characters in a story is scattered and sometimes connotative. NLP technologies automate the identification of this specific information in texts, which can be a useful weapon [35].

The popularity of deep learning has facilitated designing a number of related models to handle the subtasks of NLP, such as NER. Google proposed BERT [36], which substantially overcomes the limitations of existing models. BERT is based on Open AI GPT and performs attention mechanism on its model [33], [37]. It can predict the correct textual ID according to its entire context without a single directional limitation. In actual cases, BERT distinctly outperforms existing models in various metrics. For reference of the readers, Table III compares the capacities of the most widely adopted NLP models.

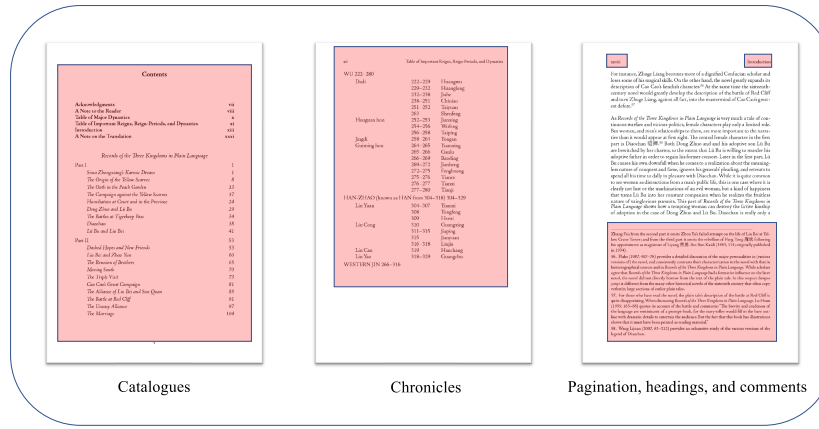


Fig. 2. Noisy content in text mining.

In this article, data from literary texts are limited. Therefore, we use a BERT + SQuAD method that can substantially address the problem because it can considerably improve prediction accuracy despite limited data [36]. Furthermore, some traditional methods are still adopted in such situations.

C. Text-Mining Algorithm

In this work, we propose a text-mining algorithm to extract the social networks in the narratives. We first preprocess the raw text to clean out the noise in the text and extract the accessible text from the narrative as the corpus. Then, we identify the characters from the corpus. Meanwhile, we also achieve sentiment extraction. Finally, the extracted characters are utilized to construct the social network. A schematic of our text-mining algorithm is shown in Fig. 1.

1) *Preprocessing*: Raw text is required to be cleaned and further normalized to a specific format (i.e., corpus) for the processing of the algorithm. Preprocessing work is relatively simplified in this work since the adoption of the deep learning-based tool enables that the related feature of the extracting objective can be learned by the model automatically. It is not required to further formulate each item of the corpus into a more easily parsing form (see examples in [12] and [42]) but keep the original text.

a) *Regular expression in data cleaning*: Noisy contents are expected to be adjusted or eliminated because they are mixed with useful data, which may mislead results. The most typical noisy contents in the text of a book include tables of contents, titles, headers, and so on. Fortunately, most of these noises usually follow regular formats. For example, as a translation of historical records, the *Records of the Three Kingdoms in Plain Language* includes a multitude of notes (see Fig. 2), where they follow the same format that starts with a serial number that leads the content. A similar phenomenon is also observed in broken words, which are all split by a hyphen or space. Regular expressions can be used to effectively eliminate this type of noisy text by developing corresponding rules [43]. Specifically, we analyze the specific pattern of different noises and apply regular expressions search algorithm to match the

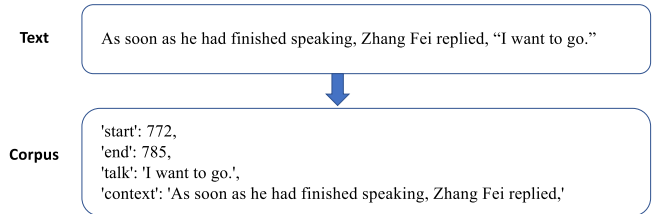


Fig. 3. Mapping from the text to the corpus.

contents matching these patterns. Finally, we can obtain clean text by eliminating these noises.

b) *Corpus extraction*: Independently building a character-oriented corpus instead of basing on the existing corpus is essential for the character extraction task in this work. We assume that in a narrative, characters usually perform in conversations; hence, their identification is focused on such conversations. Each conversation consists of several dialogs, each of which usually follows a specific double quotation mark format, that is, one paragraph starts with the double opening quote (“) and ends with the double closing quote (”). Following this rule allows dialogs to be extracted from where conversations are located. The *context* is the description that author tells audience when, where, and how a conversation occurs. The characters are usually contained in the contexts of conversations. Since the context usually leads, intersperses, or follows behind the dialog, they can be identified easily. Finally, we can construct the corpus, each of which consists of two parts, “context” and “talk” (i.e., the dialog), which map to their corresponding content in the text (see Fig. 3).

2) Speakers Extraction:

a) *Labeling the speakers*: Conversations’ portrayal varies in the storytelling. Similarly, the location of the speaker in a dialogical context differs considerably, thereby making it difficult to identify in an automated way. Therefore, a manual labeling process is required to locate the speaker in each context accurately. Given that this process is time-consuming, a GUI-labeling tool based on Jupyter Notebook is developed,

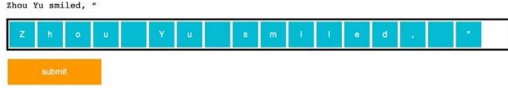


Fig. 4. Visual resolution for data labeling.

TABLE IV
DATA AUGMENTATION

	S	M	\mathcal{D}_A
<i>Romance</i>	664	1702	1130128 (approx.)
<i>Records</i>	806	1248	1005888 (approx.)

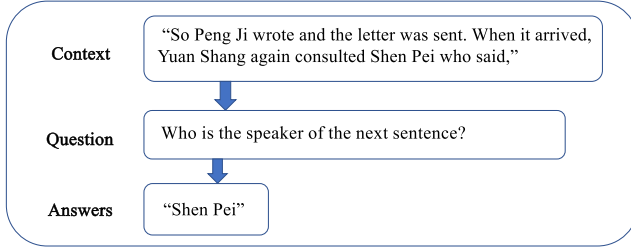


Fig. 5. Example of SQuAD.

and the visual operation substantially facilitates manual work (see Fig. 4). In this work, a total of 1702 items from *Romance* can be labeled within just 3 h.

b) Data augmentation: The size of data extracted from books is usually insufficient to reach a promising number of training samples, and it may result that the deep learning-based models cannot achieve a satisfactory prediction result [44]. In this work, the speaker corpus of *Records* collected only 1248 items, and a measly portion of 806 samples (64.5%) are labeled after transcribing the entire text. To address this issue, a data augmentation approach is introduced to generate a sufficient number of new annotated data.

How to generate new data and how much data should be generated are essential questions to answer. All speakers are assumed to be included in all the contexts. Supposing a total of S labeled speakers and M contexts, we can generate $\mathcal{D}_A = S * M$ new data samples. In this work, we use this data augmentation method and obtain over a million new data samples, as shown in Table IV.

c) Speakers identification: A BERT + SQuAD algorithm is used to build a speaker prediction model in this work. SQuAD provides a structure to answer the question (prediction) by comprehending the context. Referring to the structure of SQuAD, we structured a ternary data set (i.e., context, answer, and question), as shown in Fig. 5.

BERT provides a contextual prediction algorithm, and we use this model to predict the speakers from the text. Specifically, we employ Google’s **BERT–Base–Multilingual–Cased** model as the pretrained model, which incorporates 12-layer, 768-hidden, 12-heads, 110 million parameters; the pretrained model is then used to fine-tune our data set. Note that we omit the training procedure of BERT since it is not among the main focuses, interested readers can refer to our codes¹. No related baseline evaluates the training effort since this study is one



Fig. 6. Sentence tokenization with POS tags.

of the only existed projects relating to topic of the Three Kingdoms. The experiments are conducted with statistical significance (T-test, $\alpha = 5\%$). The predicted results after manual proofreading cover the vast majority of characters in the books (approximately 93% on *Romance* and 91% on *Records*). In this way, the appearing characters can be obtained from the prediction result.

d) Aliases association: More than a few characters in the two books possess one or multiple aliases. For example, “Xuande,” “Lord Liu,” and “The First Ruler” all refer to the character “Liu Bei.” To overcome this problem, an alias-matching mechanism is established to map aliases of the characters. Specifically, to guarantee a promising result, we manually find the aliases for corresponding characters. Then, a many-for-one mapping table is built to correlate these aliases to the character. Nonetheless, a flaw of this mechanism in practical use is that the shared family name or title may be mapped to multiple characters. For example, “Sima” can be mapped to “Sima Yi” and “Sima Zhongxiang.” We develop two solutions that can solve this problem. First, the aliases mapping is classified according to the chapters of the story. For instance, “Sima Zhongxiang” is a character who simply appears at the beginning of *Records*; hence, the mapping “Sima” to “Sima Zhongxiang” should solely be applied at the first few chapters. Second, the context is considered when mapping aliases. For example, when “Liu Bei” appears, the closest “Lord” should be “Lord Liu” (i.e., “Liu Bei”) with a high possibility. In addition, we notice that some uncommon (also known as rarely used) aliases are neglected when using this mechanism, which remains an interesting research question to be solved in the future.

3) Sentiment Extraction: While the extraction and analysis with respect to sentiment is not the main focus of this article, we still conduct related simple studies on some key characters to make the audience gain a deeper understanding of the story. Sentiments toward a character can be differently described. In this work, our sentiment analysis focuses on evaluative words. Other characters who comment about a certain character are a good entry to extract evaluations. Fig. 6 shows one of “Chen Gong”’s evaluations on “Cao Cao” in *Romance*.

Therefore, the extraction of such evaluative words is applied. First, all conversations involving a specific character are extracted and tokenized, and each token is tagged with the corresponding part of speech (i.e., POS tag). Subsequently, following the example shown in Fig. 6, words that possess an adjective POS (tagged with “JJ”) and collected since we consider them as “evaluative words” to characters, which can be utilized in sentiment analysis. Furthermore, to obtain more accurate semantic results, the evaluative words are further

¹Available at <https://github.com/GreatWizard9519/Social-network-extraction-and-analysis-of-Three-kingdoms>

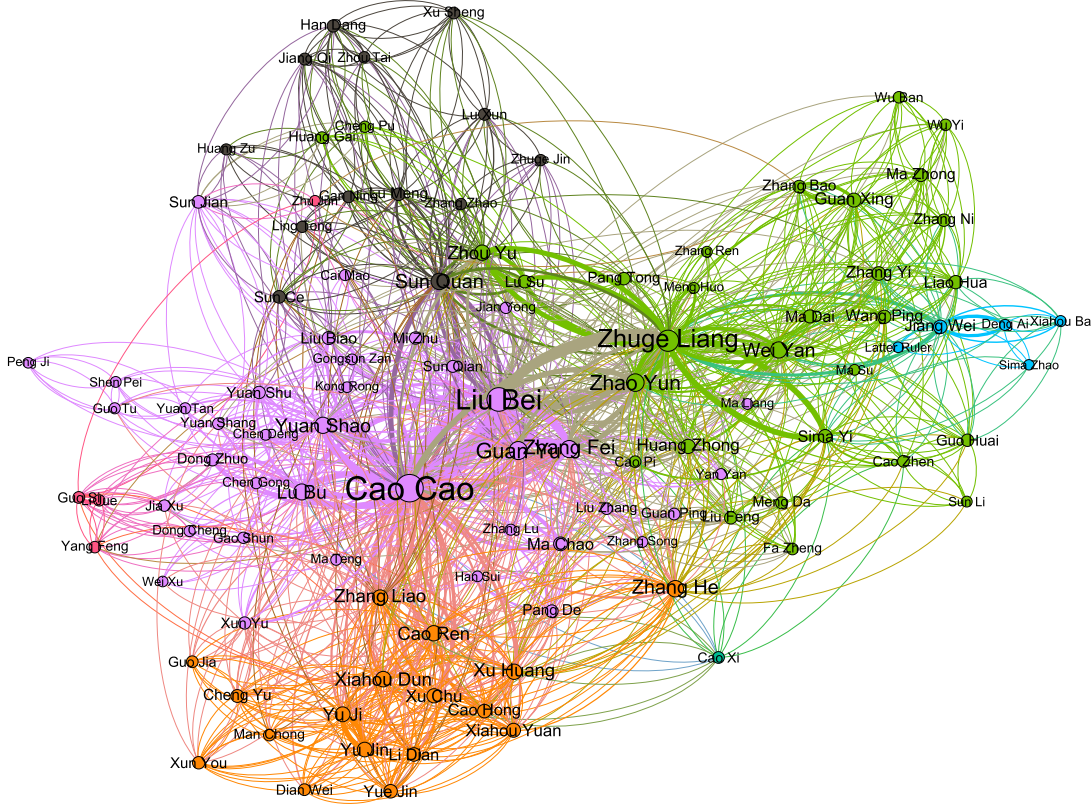


Fig. 7. Network extracted from *Romance* (only show nodes whose degree is larger than 6).

processed by using fasttext² that the semantic relatedness values of the evaluative words are computed. We give a threshold to only preserve the words with strong semantic relatedness to emotion and evaluation.

D. Network Building

1) *Representations*: Upon the collection of characters and the interaction that represents the nodes and edges, we can construct social networks. The essential representations for our extracted social networks are defined as follows.

- 1) *Nodes*: For each character coming on stage, a node is built. As aforementioned, all characters are from the identified speakers; hence, the social network merely describes the relationship between characters who have monologs or dialogs. It is worth noticing two phenomena when using this node representation. First, the number of nodes is less than the actual number of characters that appear in the books. Second, there appear some characters that are isolated without any interactions with other characters (i.e., nodes whose degree is 0) in our networks. To ensure the completeness of the social network, we manually append some of the missing characters and meanwhile include the isolated nodes when constructing the networks.
- 2) *Edge*: To correlating the nodes, namely, construct edges, we establish an assumption that the adjacent appearance

of characters will serve as the basis for creating interactions. Such an assumption is seemingly a coarse-grained solution. However, the outcome will have a high degree to match the actual situation when the size of the involved data is large enough. Based on this assumption, an algorithm established that an interaction (edge) is built when two adjacent characters are detected in the same context. Furthermore, on the account that the representation of edge describes a reciprocal relationship, the network is thereby considered as a bidirectional graph, in which the values of in-degree and out-degree of every single node are equivalent.

2) *Dynamic Network*: Unlike others, the social network extract from narrative will grow as the story carries on. Investigation of network dynamics can help us gain a better insight into the story. To this end, the texts of the two books are chronologically split into five stages, and their corresponding networks are extracted through the same method introduced above. Some key events are set aside as separate markers to normalize the distribution of each stage due to the difference in the chapter settings of the two books, for instance, *The Death of Dong Zhuo* and *The Death of Liu Bei*. Moreover, these five stages represent the five most prominent periods in the story of the Three Kingdoms. Joining the five separate networks, a dynamic network with evolving growth across the five stages is obtained.

3) *Network Visualization*: In this work, we use Gephi [45] to visualize the extracted social networks. To present a clear

²<https://fasttext.cc/docs/en/pretrained-vectors.html>

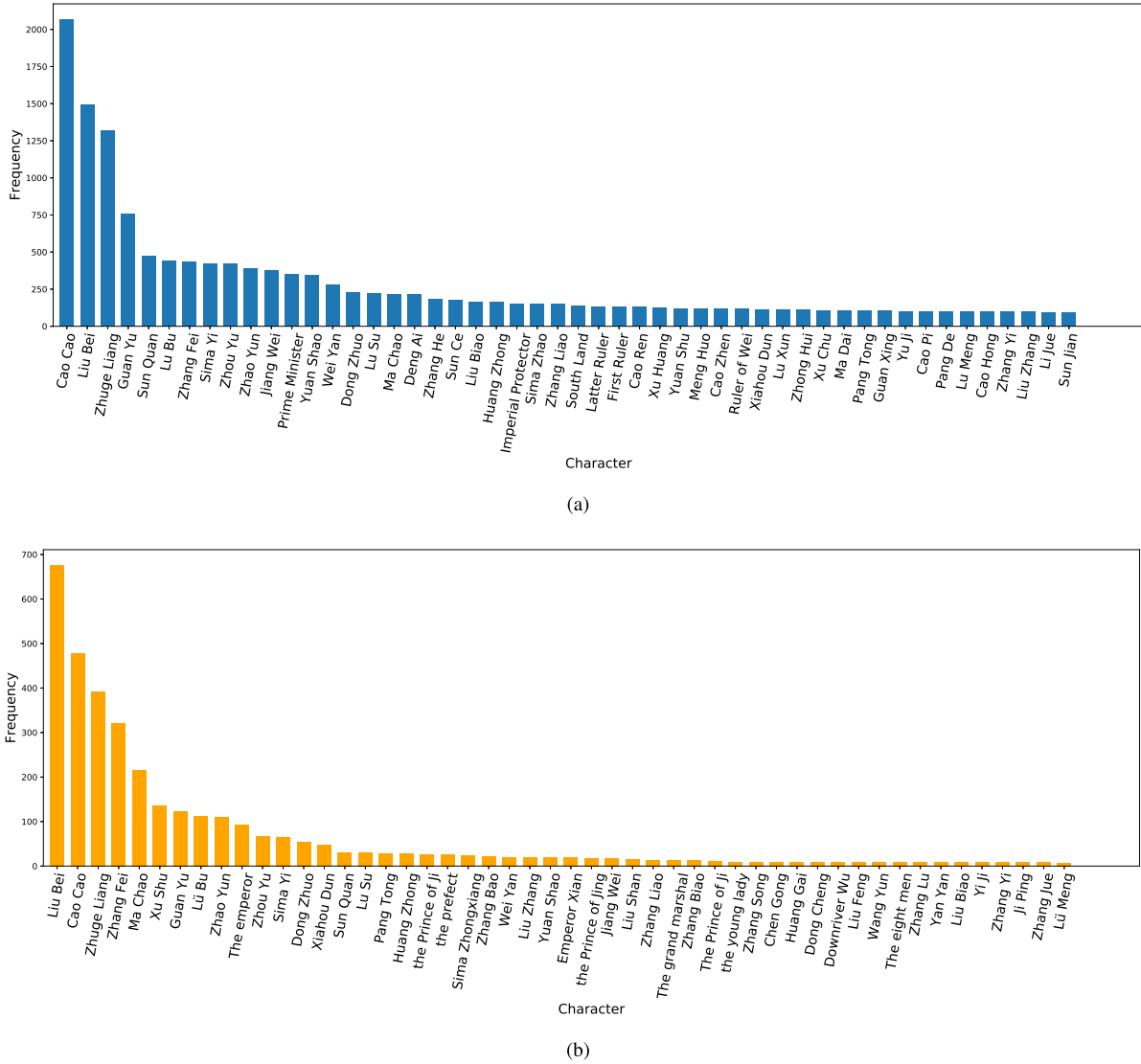


Fig. 9. Frequency of characters occurrences. (a) Top 50 in *Romance*. (b) Top 50 in *Records*.

larger than that of *Records*, the former covers approximately 71.4% characters of the latter. It indicates that the similarity of casts between the two books is very high. In addition, we rank the top 50 most frequently appearing characters in the two books (see Fig. 9) and find a 50% coincidence referring to Fig. 10. The protagonists (i.e., the top 20 characters) are notably similar; the top three most frequently appearing characters in both books are Liu Bei, Cao Cao, and Zhuge Liang.

3) *Complex Network Features of the Story*: Previous studies indicate that novelistic literature usually involves a social network that is more complex, and thus, the literariness and the dramaticism can be greatly enriched [12], [46]. In this article, two main topological features of the complex networks are considered, and related investigations are conducted on the two extracted social networks.

a) *Small-world*: It is a complex network feature that describes a random network with a highly clustered structure. In a small-world network, most nodes are not neighbors of each other, yet the neighbors of some random nodes are

probably going to be neighbors of one another, and most nodes can be reached from each other node by few jumps or steps. We can find out more homogeneity in the social structure when its social network possesses such a small-world feature. To measure the small-world feature in our extracted social networks, we introduce the two key metrics, namely, average clustering coefficient and average path length. Small-world networks are usually recognized as having large average path value length and low average clustering coefficient value. Moreover, an advanced metric, small-world index (SWI) [47], is introduced. SWI is capable of quantifying the small-world feature, which can provide a more straightforward recognition. The calculation of SWI is

$$SWI = \frac{(L - L_l)(C - C_r)}{(L_r - L_l)(C_l - C_r)} \quad (1)$$

where C and L are the clustering coefficient and average path length, respectively, which are derived from the observed network (note that we compute them by Gephi in this work); C_l and L_l refer to the clustering coefficient and mean path

"Romance"	"Records"	"Romance"	"Records"
1 Cao Cao	Liu Bei	26 Xu Huang	Liu Shan
2 Liu Bei	Cao Cao	27 Yuan Shu	Zhang Liao
3 Zhuge Liang	Zhugge Liang	28 Meng Huo	Zhang Biao
4 Guan Yu	Zhang Fei	29 Cao Zhen	The Prince of Ji
5 Sun Quan	Guan Yu	30 Xiahou Dun	the young lady
6 Lu Bu	Lu Bu	31 Lu Xun	Zhang Song
7 Zhang Fei	Zhao Yun	32 Zhong Hui	Chen Gong
8 Sima Yi	Zhou Yu	33 Xu Chu	Huang Gai
9 Zhou Yu	Sima Yi	34 Ma Dai	Dong Cheng
10 Zhao Yun	Dong Zhuo	35 Guan Xing	Xu Shu
11 Jiang Wei	Xiahou Dun	36 Pang Tong	Downriver Wu
12 Yuan Shao	Ma Chao	37 Yu Ji	Liu Feng
13 Wei Yan	Sun Quan	38 Cao Pi	Wang Yun
14 Dong Zhuo	Lu Su	39 Lu Meng	The eight men
15 Lu Su	Pang Tong	40 Pang De	Zhang Lu
16 Ma Chao	Huang Zhong	41 Cao Hong	Yan Yan
17 Deng Ai	the Prince of Ji	42 Zhang Yi	Liu Biao
18 Zhang He	Sima Zhongxiang	43 Liu Zhang	Yi Ji
19 Sun Ce	Zhang Bao	44 Li Jue	Zhang Yi
20 Liu Biao	Wei Yan	45 Sun Jian	Ji Ping
21 Huang Zhong	Liu Zhang	46 Gan Ning	Zhang Jue
22 Sima Zhao	Yuan Shao	47 Zhang Bao	Lu Meng
23 Zhang Liao	Emperor Xian	48 Xiaohou Yuan	Liu Cong
24 Liu Shan	the Prince of Jing	49 Wang Ping	Meng Huo
25 Cao Ren	Jiang Wei	50 Guo Huai	Kuai Yue

Fig. 10. Similar characters that appear in both books' top 50 ranks are highlighted in red.

TABLE VI

AVERAGE CLUSTERING COEFFICIENT, AVERAGE PATH LENGTH, AND SWI OF THE TWO NETWORKS

	Avg. Clustering coefficient	Avg. Path length	SWI
Romance	0.337	3.012	0.8713
Records	0.661	2.095	1.5679

length in a lattice reference network characterized by a high C and L , respectively, and C_r and L_r refer to the clustering coefficient and mean path length in a random reference graph characterized by a low C and L , respectively.

From the results shown in Table VI, we can observe that the *Records* has a significantly lower average path value and higher average clustering coefficient compared to those of *Romance*. Especially, the results of the calculated SWI indicate that the SWI of *Records* (1.5679) is higher than that of *Romance* (0.8713), thereby quantifiably confirming our assumption. Literature that focuses on a single character or a group of characters presents a higher SWI than those focused on a mass of characters. *Romance* involves plenty of protagonists to enrich its storytelling, features a much lower SWI than the *Records*, where the story follows only a few protagonists. It implies that *Romance* focuses more on storytelling by introducing a plethora of characters rather than the biographical narrative around several characters as presented by *Records*.

b) Scale-free: It describes a network whose degree distribution follows a power law. It reveals the Pareto principle that 20% of individuals commonly hold 80% of the total resources in a society, also known as, "the rich get richer" [48]. To investigate the scale-free feature of the two networks, we demonstrate the degree distribution of nodes in them, as shown in Fig. 11, where x is the degree of a node and y is the number of nodes that possess this degree.

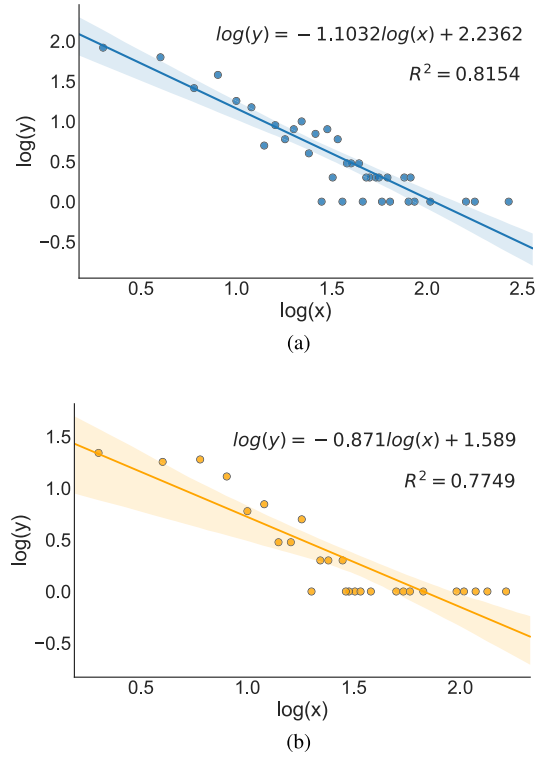


Fig. 11. Degree distribution. (a) Degree distribution of the network in *Romance*. (b) Degree distribution of the network in *Records*.

From the results, we can observe a salient power-law distribution in the diagrams of both networks. The satisfaction of the power law indicates that networks of the two books are both scale-free. However, the distribution of *Romance* has a more significant coefficient of determination (R^2 : 0.8154 > 0.7749) than that of *Records*, which means that *Romance* is relatively more in line with this law.

c) Rich-club coefficient: A number of scale-free networks exhibit a "rich-club" feature, indicating that a small number of nodes possessing a large number of edges also connect well to one another [49], [50]. The rich-club coefficient is used to measure this feature, which can be computed by

$$\phi(k) = \frac{2\mathcal{E}_{>k}}{\mathcal{N}_{>k}(\mathcal{N}_{>k} - 1)} \quad (2)$$

where $\mathcal{N}_{>k}$ is the number of nodes whose degree is not less than k , $\mathcal{E}_{>k}$ is the actual number of edges among the nodes whose degree are not less than k , and $\phi(k)$ is the ratio between the number of edges that exist among the nodes that have a degree larger than k and the total possible number among them. Considering the different sizes of the two networks, we compare the ratio of nodes that can form a fully connected network ($\phi(k) = 100\%$) deduced by the cutoff degree observed, which can be calculated by

$$r_{fc} = \frac{k_{\phi(k)=100\%}}{\mathcal{N}} \quad (3)$$

where $k_{\phi(k)=100\%}$ is the minimum k , which enables $\phi(k) = 100\%$, and \mathcal{N} is the number of nodes in the network.

The calculated Rich-club coefficients of the two networks are shown in Fig. 12. The calculated r_{fc} are 5.01% (*Romance*)

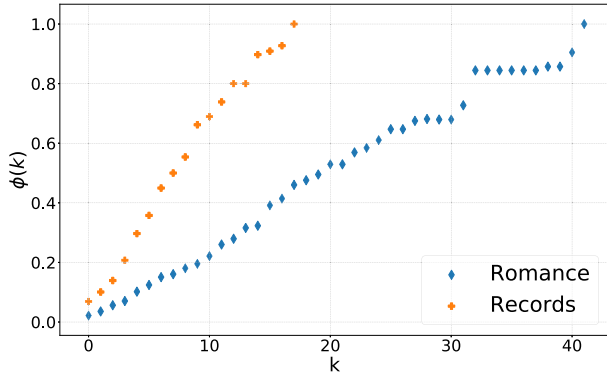


Fig. 12. Rich-club coefficients of the two networks.

and 20.30% (*Records*). It reveals that the top 5% rich nodes in *Romance* can approximately form a fully connected network, whereas the number has to be approximately the top 20% in *Records*. It can be concluded that both networks have a rich-club feature, which is more significant in *Romance*. These results reveal that despite more characters appearing in *Romance* than in *Records*, the story always revolves around a few protagonists in *Romance*.

4) *History or Romance: Which Is More Dramatic?*: Dramatic changes make stories splendid. The rise and fall of warlords constantly change the social structure of the story of the Three Kingdoms across all stages. To study the growth of the social structure, we investigate the social network according to the idea introduced in Section III-D2. As shown in Fig. 13, five metrics are adopted to observe the dynamic change of the networks.

Interesting phenomena are found in the results as follows. The average node growth rates in *Romance* and *Records* are 147% and 63%, respectively. This suggests that *Romance* has dramatic changes in terms of the number of characters, and the appearance of characters on each stage is overwhelming. The density comparison indicates that *Romance* has a larger network size through all the stages, yet its density is lower, where the gap in the last three stages is especially notable. The change of the average degree of two networks follows a similar pattern, demonstrating that they both increase at the beginning and then reach a plateau. *Records* has a considerably shorter average path length and higher average clustering coefficient in the majority of five stages except in the first two, which match its more significant performance in small-worldliness. Overall, we can observe that the growth of the social network in *Romance* is more rapid. Comparatively, *Records* has a tight and gradually clustered network.

B. Network Feature on Specific Characters

In this section, we assess the network feature on specific characters. While the story of the Three Kingdoms involves numerous forces, the main focus is the three force blocs, i.e., Wei, Shu, and Wu. Therefore, their respective sovereigns, namely, Cao Cao (Wei), Liu Bei (Shu), and Sun Quan (Wu), are chosen as the targets for our character-centric investigation.

1) *Who Is the Most Influential?*: In the story of the Three Kingdoms, the personal influence of each sovereign

considerably represents the influence of the forces they possess. Given this kind of influence in a social network, the sovereigns' interactions with other characters reflect their influence. Three related metrics are introduced to compare their influence.

- 1) *Degree*: Degree or degree centrality is a basic measure that counts the number of neighbors that a node character has. The weighted degree is additionally considered, which is calculated by considering the number of interactions that occur between two characters.
- 2) *Closeness Centrality*: Closeness centrality measures the extent of closeness of a node to a network. It is calculated as the reciprocal of the sum of the length of the shortest paths between the node and all other nodes in the graph. Its formula is expressed as

$$C(i) = \frac{1}{\sum_j d(i, j)} \quad (4)$$

where $C(i)$ is the closeness centrality of node i and $d(i, j)$ denotes the distance between nodes i and j .

- 3) *Betweenness Centrality*: For each pair of nodes in a network, at least one shortest path exists between nodes, in which either the number of edges that the path passes through (for unweighted networks) or the sum of the weights of the edges (for weighted networks) is minimized. Betweenness centrality is a measure of the number of the shortest path that passes through a node. Denoted by $g(v)$, the betweenness centrality of node v can be calculated by

$$g(v) = \sum_{i \neq v \neq j} \frac{\sigma_{ij}(v)}{\sigma_{ij}} \quad (5)$$

where σ_{ij} is the total number of shortest paths from node i to node j and $\sigma_{ij}(v)$ is the number of those paths that pass through node v . A character's property of "bridge" can be measured by betweenness centrality.

Table VII shows the results of *Romance* and Table VIII shows the results of *Records*. In *Romance*, Cao Cao exhibits the highest measures of the four metrics, followed by Liu Bei, whereas Sun Quan has the lowest measures. In *Records*, Liu Bei leads the performance instead of Cao Cao, and Sun Quan is far behind them. This can support us to conclude that Cao Cao is the most influential of the three sovereigns in *Romance*, Liu Bei is the one in *Records*, and the influence of Sun Quan is lower than the other two lords in both books.

C. Sentiment Analysis on Characters

1) *"Taste" of the Character Sentiment*: We commence by collecting and ranking the evaluative words to Liu Bei and Cao Cao. Specifically, we present the results by adopting the word cloud, as shown in Fig. 14. As the word cloud shown in Fig. 14, a sketchy sentimental opinion on Cao Cao and Liu Bei can be obtained. In particular, we obtain the following subjective observations. In *Romance*, evaluative words, such as "great" and "able," are mentioned for both two lords. However, we in addition find words such as "crafty" and "evil" on Cao Cao and "humble" on Liu Bei, which reveals

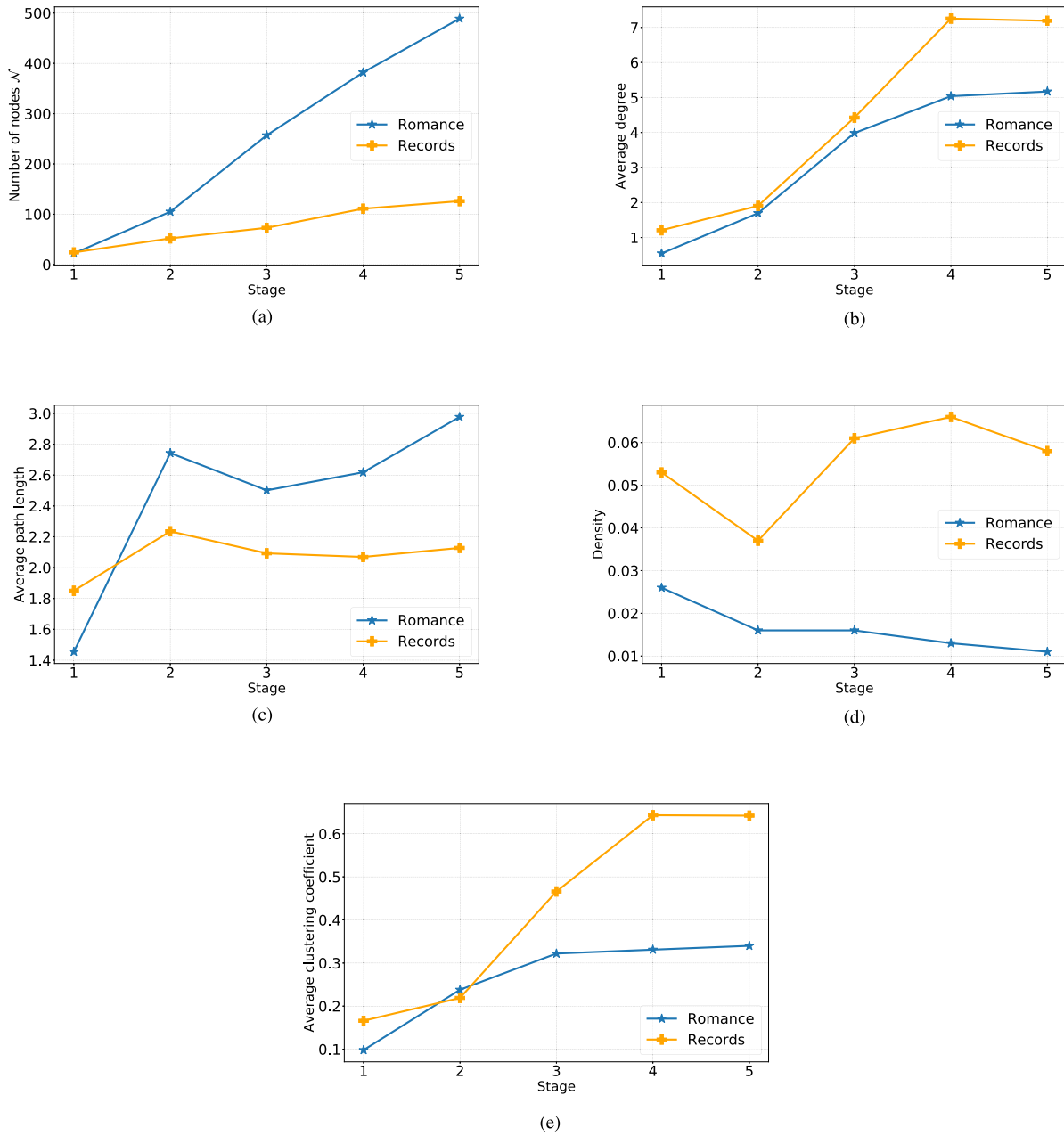


Fig. 13. Dynamic growth of the two networks (five stages). (a) Number of nodes. (b) Average degree. (c) Average path length. (d) Density. (e) Average clustering coefficient.

TABLE VII
DEGREE AND CENTRALITY OF CAO CAO, LIU BEI, AND SUN QUAN IN THE NETWORK OF *Romance*

	Degree	Weighted degree	Closeness centrality	Betweenness centrality
Cao Cao	268	3442	0.565371	18765.21
Liu Bei	182	2301	0.511279	7913.80
Sun Quan	80	908	0.448808	1688.14

the difference. Moreover, more negative words are obviously found about Cao Cao in *Romance* than in *Records*. While this observation cannot bring us to the conclusion that the authors of the two books have an evident preference to a character, we can at least find that there exist differences regarding the depiction of the same character in the two books.

Generally, historical records tend to lean toward objectivity, whereas fictional novels contain subjective emotions. The

creation of *Romance* began approximately toward the end of the Yuan Dynasty, which was a dark era for common people. The dissatisfaction with the ruling class can be reflected by the impressionable attitude of people to some forces (e.g., Shu) in the story of the Three Kingdoms. In this context, the author of *Romance* emotionally depicted a series of characters who are different in actual history. This phenomenon substantially occurs to Liu Bei and Cao Cao, which are lords of Shu

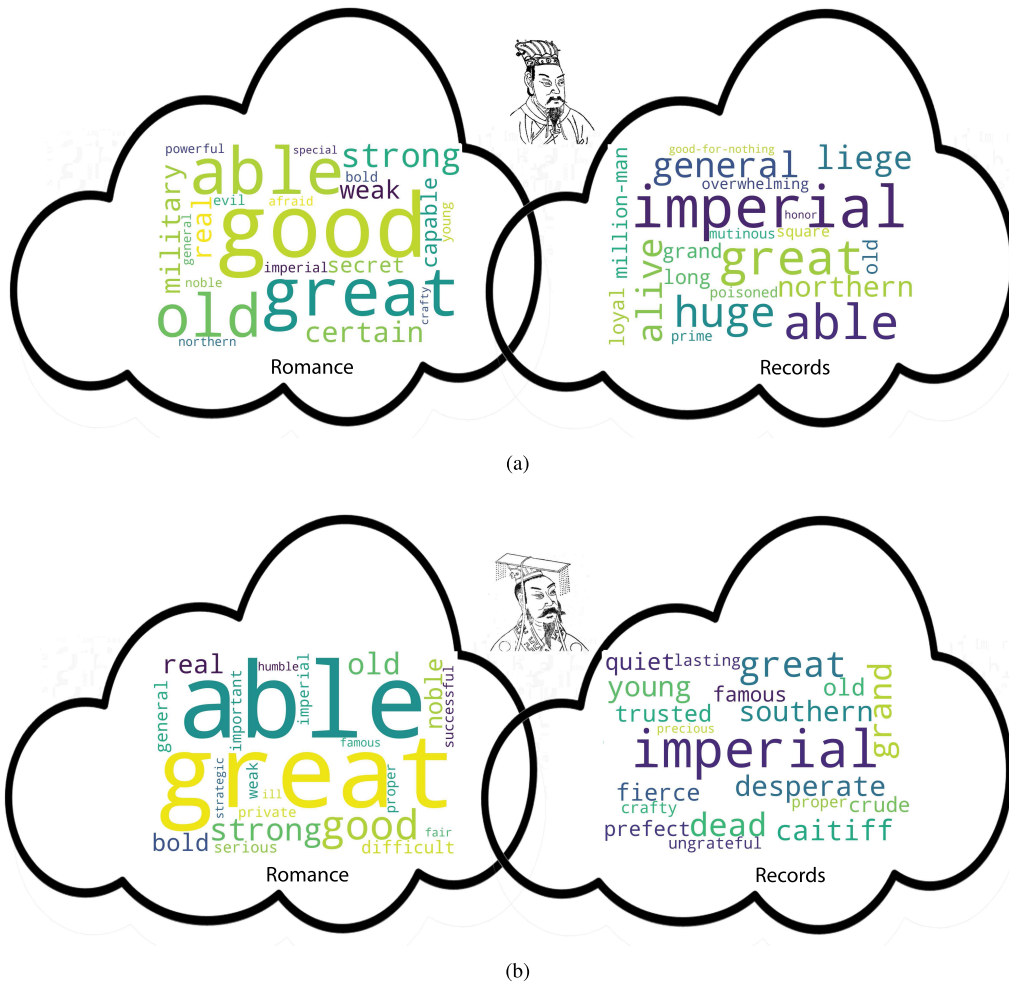


Fig. 14. Word clouds of the evaluative words on (a) Cao Cao and (b) Liu Bei.

TABLE VIII
DEGREE AND CENTRALITY OF CAO CAO, LIU BEI, AND SUN QUAN IN THE NETWORK OF *Records*

	Degree	Weighted degree	Closeness centrality	Betweenness centrality
Liu Bei	169	2762	0.795316	4779.68
Cao Cao	134	1920	0.707317	2764.84
Sun Quan	24	326	0.527273	20.007937

and Wei, respectively. Our investigation focuses on these two characters from the point of their evaluating words.

2) *Sentimental Quantification*: For a better understanding of the evaluative words to, we conduct a quantitative comparison. In particular, we introduce a sentimental scoring metric, SentiWordNet [51]. SentiWordNet score can be calculated by subtracting both polarities (positive and negative) of each token and subsequently calculating them

$$\text{score} = \frac{\sum_{i=1}^n (\text{posScore}_i - \text{negScore}_i)}{n} \quad (6)$$

where n denotes the number of involved evaluative words and posScore_i and negScore_i are the positive and negative scores of word i provided by SentiWordNet. The criterion of SentiWordNet gives Negative (i.e., -1), Neutral (i.e., 0), and Positive (i.e., 1) for users to classify the word.

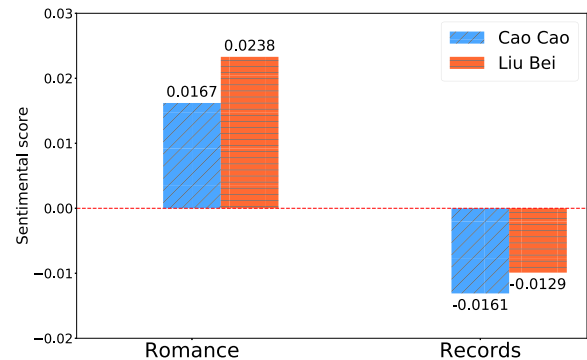


Fig. 15. SentiWordNet scores of Cao Cao and Liu Bei.

Fig. 15 implies that Cao Cao's score is lower than that of Liu Bei in *Romance*. Nonetheless, the score of Cao Cao is

higher in *Records*. This finding is consistent with the subjective perception obtained in Section IV-C1. In addition, the scores of the two characters are both higher in *Romance* than in *Records*. This possibly reveals the different sentimental tones of the authors' wording in the narrative.

V. CONCLUSION

Surrounding on the story of the Three Kingdoms, this article revives the research on digital humanities, which seeks to digitize working procedures of sociologists and historians in the field of humanities by using state-of-the-art data science technologies.

An algorithm is developed to extract social networks of stories narrated in two books with respect to the Three Kingdoms. In particular, the advanced NLP model BERT is employed in our character identification work, and a satisfying outcome is obtained. Subsequently, we conduct a series of topological analyses to quantify and characterize the extracted social networks, where we additionally present a quantitative comparison between the two books. Specifically, network topological features, such as small-world, scale-free, and centrality of specific characters, are measured. The results reveal that the social network is more entangled in the narrative of the *Romance* than that of the *Records*, especially, more protagonist-oriented. Moreover, this provides a quantitative reference for the macro (e.g., structural features of a story) and micro levels (e.g., the influence or sentiment of a specific character), and the extent of the grandness vividness of a story can be expressed scientifically.

This work can help both researchers and nonexpert readers gain an insight into the story of the Three Kingdoms and the procedure of its digital analysis. Moreover, numerous involved subworks can be refined in the future. First, the definition of interactions between characters is coarse-grained. Second, a mere five-slice dynamic network is built in this project, and hopefully, a large-scale dynamic network, which can incorporate hundreds even thousands of slices, can be obtained if the story is subdivided in fine granularity, for instance, year-to-year or day-to-day. In addition, we would like to involve more social relationships (e.g., forces, military conflicts, and marriages) in our future exploration.

REFERENCES

- [1] E. Rosati, "Google Books' library project is fair use," *J. Intellectual Property Law Pract.*, vol. 9, no. 2, pp. 104–106, Feb. 2014.
- [2] J. Strehovec, *Text as Ride: Electronic Literature and New Media Art*. Morgantown, WV, USA: Center for Literary Computing, 2016.
- [3] M. G. Kirschenbaum, "What is digital humanities and what's it doing in English departments," in *Defining Digital Humanities*. Evanston, IL, USA: Routledge, 2016, pp. 211–220.
- [4] S. E. Worth, "Storytelling and narrative knowing: An examination of the epistemic benefits of well-told stories," *J. Aesthetic Educ.*, vol. 42, no. 3, pp. 42–56, 2008.
- [5] T. Chowdhury, S. Muhuri, S. Chakraborty, and S. N. Chakraborty, "Analysis of adapted films and stories based on social network," *IEEE Trans. Comput. Social Syst.*, vol. 6, no. 5, pp. 858–869, Oct. 2019.
- [6] M. Polo, U. D. Iacono, G. Fiorentino, and A. Pierri, "A social network analysis approach to a digital interactive storytelling in mathematics," *J. e-Learn. Knowl. Soc.*, vol. 15, no. 3, pp. 239–250, 2019.
- [7] M. A. Bellamy and R. C. Basole, "Network analysis of supply chain systems: A systematic review and future research," *Syst. Eng.*, vol. 16, no. 2, pp. 235–249, Jun. 2013.
- [8] R. C. Basole and M. A. Bellamy, "Visual analysis of supply network risks: Insights from the electronics industry," *Decis. Support Syst.*, vol. 67, pp. 109–120, Nov. 2014.
- [9] E. Otte and R. Rousseau, "Social network analysis: A powerful strategy, also for the information sciences," *J. Inf. Sci.*, vol. 28, no. 6, pp. 441–453, Dec. 2002.
- [10] R. Franzosi, *Quantitative Narrative Analysis*, no. 162. Newbury Park, CA, USA: Sage, 2010.
- [11] N. Kayastha, D. Niyato, P. Wang, and E. Hossain, "Applications, architectures, and protocol design issues for mobile social networks: A survey," *Proc. IEEE*, vol. 99, no. 12, pp. 2130–2158, Dec. 2011.
- [12] M. C. Waumans, T. Nicodème, and H. Bersini, "Topology analysis of social networks extracted from literature," *PLoS ONE*, vol. 10, no. 6, Jun. 2015, Art. no. e0126470.
- [13] J. Zhang, H. Zhao, J.-Q. Xu, and J.-F. Wang, "Small-world and scale-free features in Harry Potter," *TELKOMNIKA Indonesian J. Elect. Eng.*, vol. 12, no. 8, pp. 6411–6416, Aug. 2014.
- [14] T. Gessey-Jones *et al.*, "Narrative structure of a song of ice and fire creates a fictional world with realistic measures of social complexity," *Proc. Nat. Acad. Sci. USA*, vol. 117, no. 46, pp. 28582–28588, 2020.
- [15] P. Wang, J.-A. Lu, Y. Jin, M. Zhu, L. Wang, and S. Chen, "Statistical and network analysis of 1212 COVID-19 patients in Henan, China," *Int. J. Infectious Diseases*, vol. 95, pp. 391–398, Jun. 2020.
- [16] J. Seo, G.-M. Park, S.-H. Kim, and H.-G. Cho, "Characteristic analysis of social network constructed from literary fiction," in *Proc. Int. Conf. Cyberworlds*, Oct. 2013, pp. 147–150.
- [17] Z. Bu, H.-J. Li, C. Zhang, J. Cao, A. Li, and Y. Shi, "Graph K-means based on leader identification, dynamic game, and opinion dynamics," *IEEE Trans. Knowl. Data Eng.*, vol. 32, no. 7, pp. 1348–1361, Jul. 2020.
- [18] J. Cao, Z. Bu, Y. Wang, H. Yang, J. Jiang, and H.-J. Li, "Detecting prosumer-community groups in smart grids from the multiagent perspective," *IEEE Trans. Syst., Man, Cybern. Syst.*, vol. 49, no. 8, pp. 1652–1664, Aug. 2019.
- [19] M. McPherson, L. Smith-Lovin, and J. M. Cook, "Birds of a feather: Homophily in social networks," *Annu. Rev. Sociol.*, vol. 27, no. 1, pp. 415–444, Aug. 2001.
- [20] A. Cardillo *et al.*, "Emergence of network features from multiplexity," *Sci. Rep.*, vol. 3, no. 1, pp. 1–6, Dec. 2013.
- [21] S. R. Dobrow, D. E. Chandler, W. M. Murphy, and K. E. Kram, "A review of developmental networks: Incorporating a mutuality perspective," *J. Manage.*, vol. 38, no. 1, pp. 210–242, Jan. 2012.
- [22] F. J. Flynn, R. E. Reagans, and L. Guillory, "Do you two know each other? Transitivity, homophily, and the need for (network) closure," *J. Personality Social Psychol.*, vol. 99, no. 5, p. 855, 2010.
- [23] F. Bloch, M. O. Jackson, and P. Tebaldi, "Centrality measures in networks," SSRN, Tech. Rep. 2749124, 2019.
- [24] I. M. Pepperberg, "Rethinking syntax: A commentary on E. Kako's 'Elements of syntax in the systems of three language-trained animals,'" *Animal Learn. Behav.*, vol. 27, no. 1, pp. 15–17, Mar. 1999.
- [25] R. S. Burt, *Structural Holes: The Social Structure of Competition*. Cambridge, MA, USA: Harvard Univ. Press, 2009.
- [26] O. Borrega, M. Taulé, and M. A. Martí, "What do we mean when we speak about named entities," in *Proc. Corpus Linguistics*, 2007.
- [27] E. F. Sang and F. De Meulder, "Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition," 2003, *arXiv:cs/0306050*. [Online]. Available: <https://arxiv.org/abs/cs/0306050>
- [28] M. Marrero, J. Urbano, S. Sánchez-Cuadrado, J. Morato, and J. M. Gómez-Berbís, "Named entity recognition: Fallacies, challenges and opportunities," *Comput. Standards Interfaces*, vol. 35, no. 5, pp. 482–489, Sep. 2013.
- [29] P. Pakray, A. Pal, G. Majumder, and A. Gelbukh, "Resource building and parts-of-speech (POS) tagging for the Mizo language," in *Proc. 14th Mex. Int. Conf. Artif. Intell. (MICAI)*, Oct. 2015, pp. 3–7.
- [30] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," 2013, *arXiv:1301.3781*. [Online]. Available: <http://arxiv.org/abs/1301.3781>
- [31] H. Zen and H. Sak, "Unidirectional long short-term memory recurrent neural network with recurrent output layer for low-latency speech synthesis," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Apr. 2015, pp. 4470–4474.
- [32] M. E. Peters *et al.*, "Deep contextualized word representations," 2018, *arXiv:1802.05365*. [Online]. Available: <http://arxiv.org/abs/1802.05365>
- [33] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, "Improving language understanding with unsupervised learning," OpenAI, San Francisco, CA, USA, Tech. Rep., 2018.

- [34] L. Màrquez and H. Rodríguez, “Part-of-speech tagging using decision trees,” in *Proc. Eur. Conf. Mach. Learn.* Berlin, Germany: Springer, 1998, pp. 25–36.
- [35] M. Mitri, “Story analysis using natural language processing and interactive dashboards,” *J. Comput. Inf. Syst.*, pp. 1–11, Jul. 2020.
- [36] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” 2018, *arXiv:1810.04805*. [Online]. Available: <http://arxiv.org/abs/1810.04805>
- [37] A. Vaswani *et al.*, “Attention is all you need,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2017, pp. 5998–6008.
- [38] S. Huang and J. Wu, “A pragmatic approach for classical Chinese word segmentation,” in *Proc. 11th Int. Conf. Lang. Resour. Eval. (LREC)*, 2018, pp. 1–8.
- [39] C. H. Brewitt-Taylor and T. Richard, *Romance of the Three Kingdoms*. North Clarendon, VT, USA: Tuttle, 1931.
- [40] W. L. Idema, S. H. West, *Records of the Three Kingdoms in Plain Language*. Indianapolis, IN, USA: Hackett, 2016.
- [41] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, “SQuAD: 100,000+ questions for machine comprehension of text,” 2016, *arXiv:1606.05250*. [Online]. Available: <http://arxiv.org/abs/1606.05250>
- [42] A. Doitch, R. Yazdi, T. Hazan, and R. Reichart, “Perturbation based learning for structured NLP tasks with application to dependency parsing,” *Trans. Assoc. Comput. Linguistics*, vol. 7, pp. 643–659, Nov. 2019.
- [43] Y. Li, R. Krishnamurthy, S. Raghavan, S. Vaithyanathan, and H. V. Jagadish, “Regular expression learning for information extraction,” in *Proc. Conf. Empirical Methods Natural Lang. Process. (EMNLP)*, 2008, pp. 21–30.
- [44] H.-L. Trieu, D.-V. Tran, A. Ittoo, and L.-M. Nguyen, “Leveraging additional resources for improving statistical machine translation on asian low-resource languages,” *ACM Trans. Asian Low-Resource Lang. Inf. Process.*, vol. 18, no. 3, pp. 1–22, Jul. 2019.
- [45] M. Bastian *et al.*, “Gephi: An open source software for exploring and manipulating networks,” in *Proc. ICWSM*, vol. 8, 2009, pp. 361–362.
- [46] D. Elson, N. Dames, and K. McKeown, “Extracting social networks from literary fiction,” in *Proc. 48th Annu. Meeting Assoc. Comput. Linguistics*, 2010, pp. 138–147.
- [47] Z. P. Neal, “How small is it? Comparing indices of small worldliness,” *Netw. Sci.*, vol. 5, no. 1, pp. 30–44, Mar. 2017.
- [48] A. Chakraborti and M. Patriarca, “Variational principle for the Pareto power law,” *Phys. Rev. Lett.*, vol. 103, no. 22, Nov. 2009, Art. no. 228701.
- [49] S. Zhou and R. J. Mondragon, “The rich-club phenomenon in the Internet topology,” *IEEE Commun. Lett.*, vol. 8, no. 3, pp. 180–182, Mar. 2004.
- [50] J. Alstott, P. Panzarasa, M. Rubinov, E. T. Bullmore, and P. E. Vértés, “A unifying framework for measuring weighted rich clubs,” *Sci. Rep.*, vol. 4, no. 1, p. 7258, May 2015.
- [51] S. Baccianella, A. Esuli, and F. Sebastiani, “SentiWordNet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining,” in *Proc. LREC*, vol. 10, 2010, pp. 2200–2204.