

Extraction and Representation of Financial Entities from Text

Tim Repke and Ralf Krestel

Abstract In our modern society, almost all events, processes, and decisions in a corporation are documented by internal written communication, legal filings, or business and financial news. The valuable knowledge in such collections is not directly accessible by computers as they mostly consist of unstructured text. This chapter provides an overview of corpora commonly used in research, highlights related work, and state-of-the-art approaches to extract and represent financial entities and relations.

The second part of this chapter considers applications based on knowledge graphs of automatically extracted facts. Traditional information retrieval systems typically require the user to have prior knowledge of the data. Suitable visualization techniques can overcome this requirement and enable users to explore large sets of documents. Furthermore, data mining techniques can be used to enrich or filter knowledge graphs. This information can augment source documents and guide exploration processes. Systems for document exploration are tailored to specific tasks, such as investigative work in audits or legal discovery, monitoring compliance, or providing information in a retrieval system to support decisions.

1 Introduction

Data is frequently called the oil of the 21st century.¹ Substantial amounts are produced by our modern society each day and stored in big data centers. However, the actual value is only generated through statistical analyses and data mining.

Tim Repke
Hasso Plattner Institute, University of Potsdam, Germany, e-mail: tim.repke@hpi.de

Ralf Krestel
Hasso Plattner Institute, University of Potsdam, Germany, e-mail: ralf.krestel@hpi.de

¹ E.g., <https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data>

Computer algorithms require numerical and structured data, such as in relational databases. Texts and other unstructured data contain a lot of information that is not readily accessible in a machine-readable way. With the help of *text mining*, computers can process large corpora of text. Modern *natural language processing* (NLP) methods can be used to extract structured data from text, such as mentions of companies and their relationships. This chapter outlines the fundamental steps necessary to construct a *knowledge graph* (KG) with all the extracted information. Furthermore, we will highlight specific state-of-the-art techniques to further enrich and utilize such a knowledge graph. We will also present text mining techniques that provide numerical representations of text for structured semantic analysis.

Many applications greatly benefit from an integrated resource for information in exploratory use cases and analytical tasks. For example, journalists investigating the Panama papers needed to untangle and sort through vast amounts of data, search entities, and visualize found patterns hidden in the large and very heterogeneous leaked set of documents and files [10]. Similar datasets are of interest for data-journalists in general or in the context of computational forensics [19, 13]. Auditing firms and law enforcement need to sift through massive amounts of data to gather evidence of criminal activity, often involving communication networks and documents [28]. Current computer-aided exploration tools,² offer a wide range of features from data ingestion, exploration, analysis, to visualization. This way, users can quickly navigate the underlying data based on extracted attributes, which would otherwise be infeasible due to the often large amount of heterogeneous data.

There are many ways to represent unstructured text in a machine-readable format. In general, the goal is to reduce the amount of information to provide humans an overview and enable the generation of new insights. One such representation are *knowledge graphs*. They encode facts and information by having nodes and edges connecting these nodes forming a graph.³ In our context, we will consider nodes in the graph as named entities, such as people or companies, and edges as their relationships. This representation allows humans to explore and query the data on an abstracted level and run complex analyses. In economics and finance, this offers access to additional data sources. Whereas internally stored transactions or balance sheets at a bank only provide a limited view of the market, information hidden in news, reports, or other textual data may offer a more global perspective.

For example, the context in which data was extracted can be a valuable additional source of information that can be stored alongside the data in the knowledge graph. Topic models [8] can be applied to identify distinct groups of words that best describe the key topics in a corpus. In recent years, embeddings significantly gained popularity for a wide range of applications [64]. Embeddings represent a piece of text as a high dimensional vector. The distance between vectors in such a vector space can be interpreted as semantic distance and reveals interesting relationships.

This chapter focuses on the construction and application of knowledge graphs, particularly on company networks. In the first part, we describe an NLP pipeline's

² e.g., extraction and indexing engine (<https://www.nuix.com/>), network analysis and visualization (<https://linkurio.us/>), or patent landscapes (<https://clarivate.com/derwent/>)

³ Knowledge graphs are knowledge bases whose knowledge is organized as a graph.

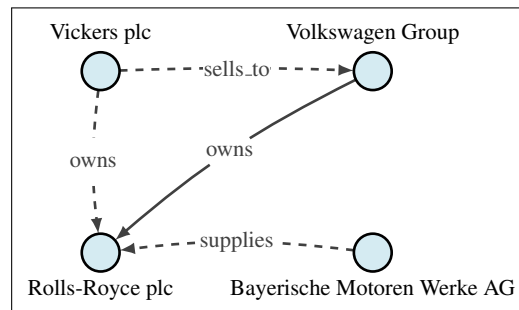
key steps to construct (see Sect. 2) and refine (see Sect. 3) such a knowledge graph. In the second part, we focus on applications based on knowledge graphs. We differentiate them into syntactic and semantic exploration. The syntactic exploration (see Sect. 5) considers applications that directly operate on the knowledge graph’s structure and meta-data. Typical use cases assume some prior knowledge of the data and support the user by retrieving and arranging the relevant extracted information. In Sect. 6 we extend this further to the analogy of semantic maps for interactive visual exploration. Whereas syntactic applications follow a localized bottom-up approach for the interactions, semantic exploration usually enables a top-down exploration, starting from a condensed global overview of all the data.

2 Extracting Knowledge Graphs from Text

Many business insights are hidden in unstructured text. Modern NLP methods can be used to extract that information as structured data. In this section, we mainly focus on named entities and their relationships. These could be mentions of companies in news articles, credit reports, emails, or official filings. The extracted entities can be categorized and linked to a knowledge graph. Several of those are publicly accessible and cover a significant amount of relations, namely Wikidata [77], the successor of DBpedia [34], and Freebase [9], as well as YAGO [76]. However, they are far from complete and usually general-purpose, so that specific domains or details might not be covered. Thus, it is essential to extend them automatically using company-internal documents or domain-specific texts.

The extraction of named entities is called *names entity recognition* (NER) [23] and comprises two steps: First, detecting the boundaries of the mention within the string of characters and second, classifying it into types such as ORGANIZATION, PERSON, or LOCATION. Through *named entity linking* (NEL) [70],⁶ a mention is matched to its corresponding entry in the knowledge graph (if already known). An

Fig. 1 Network of information extracted from the excerpt: *VW purchased Rolls-Royce & Bentley from Vickers on 28 July 1998. From July 1998 until December 2002, BMW continued to supply engines for the Rolls-Royce Silver Seraph.*⁵



⁵ Excerpt from https://en.wikipedia.org/wiki/Volkswagen_Group. Accessed on 22.02.2020

⁶ Also called *entity resolution*, *entity disambiguation*, *entity matching*, or *record linkage*.

unambiguous assignment is crucial for integrating newly found information into a knowledge graph. For the scope of this chapter, we consider a fact to be a relation between entities. The most naïve approach is to use entity co-occurrence in text. *Relationship extraction* (RELEX) identifies actual connections stated in the text, either with an *open* or *closed* approach. In a closed approach, the relationships are restricted to a pre-defined set of relations, whereas the goal with an open approach is to extract all connections without restrictions.

Figure 1 shows a simplified example of a company network extracted from a small text excerpt. Instead of using the official legal names, quite different colloquial names, acronyms, or aliases are typically used when reporting about companies. There are three main challenges in entity linking: 1) *name variations* as shown in the example with “VW” and “Volkswagen”; 2) *entity ambiguity*, where a mention can also refer to multiple different knowledge graph entries; and 3) *unlinkable entities* in the case, that there is no corresponding entry in the knowledge graph yet. The resulting graph in Fig. 1 depicts a sample knowledge graph generated from facts extracted from the given text excerpt. Besides the explicitly mentioned entities and relations, the excerpt also contains many implied relationships; for example, a sold company is owned by someone else after the sell. Further, relationships can change over time, leading to edges that are only valid for a particular time. This information can be stored in the knowledge graph and, e.g., represented through different types of edges in the graph. Through *knowledge graph completion*, it is possible to estimate the probability whether a specific relationship between entities exists [74].

In the remainder of this section, we provide a survey of techniques and tools for each of the three steps mentioned above: NER (Subsect. 2.2), NEL (Subsect. 2.2), and RELEX (Subsect. 2.3).

2.1 Named Entity Recognition (NER)

The first step of the pipeline for knowledge graph construction from text is to identify mentions of named entities. Named entity recognition includes several subtasks, namely identifying proper nouns, the boundaries of named entities, and classifying the entity type. The first work in this area was published in 1991 and proposed an algorithm to automatically extract company names from financial news to build a database for querying [54, 46]. The task gained interest with MUC-6, a shared task to distinguish not only types, such as person, location, organization, but also numerical mentions, such as time, currency, and percentages [23]. Traditionally, research in this area is founded in computational linguistics, where the goal is to parse and describe the natural language with statistical rule-based methods. The foundation for that is to correctly tokenize the unstructured text, assign part-of-speech tags (also known as POS-tagging), and create a parse tree that describes the sentence’s dependencies and overall structure. Using this information, linguists defined rules that describe typical patterns for named entities.

Handcrafted rules were soon replaced by machine learning approaches that use tags mentioned above and so-called surface features. These surface features describe syntactic characteristics, such as the number of characters, capitalization, and other derived information. The most popular supervised learning methods for the task are hidden Markov models and conditional random fields due to their ability to derive probabilistic rules from sequence data [6, 42]. However, supervised learning requires large amounts of annotated training data. Bootstrapping methods can automatically label text data using a set of entity names as a seed. These semi-supervised methods do so by marking occurrences of these seed entities in the text and using contextual information to annotate more data automatically. For an overview of related work in that area, we refer to Nadeau et al. [47]. In recent years, deep learning approaches gained popularity. They have the advantage that they do not require sophisticated pre-processing, feature engineering, or potentially error-prone POS-tagging and dependency parsing. Especially recurrent neural networks are well suited since they take entire sequences of tokens or characters into account. For an overview of currently researched deep learning models for NER, we refer readers to the extensive survey by Yadav and Bethard [80].

Although the task of automatically identifying company names in text, there is still a lot of research dedicated to named entity recognition. Due to their structural heterogeneity, recognizing company names is particularly difficult compared to person or location names. Examples of actual German company names show the complexity of the task. Not only are some of the names very long (“Simon Kucher & Partner Strategy & Marketing Consultants GmbH”), they interleave abstract names with common nouns, person names, locations, and legal forms, for example: “Loni GmbH”, “Klaus Traeger”, “Clean-Star GmbH & Co Autowaschanlage Leipzig KG”. Whereas in English almost all capitalized proper nouns refer to named entities, it is significantly harder to find entity mentions in other languages, for example, in German, where all common nouns are capitalized [17]. Loster et al. [65, 36] dedicate a series of papers to the recognition of financial entities in text. In particular, they focus on correctly determining the full extent of a mention by using tries, which are tree structures, to improve dictionary-based approaches [39].

The wealth of publications and the availability of open-source libraries reflect the importance and popularity of NER. The following overview shows the most successful projects used in research and industry alike.

GATE ANNIE The General Architecture for Text Engineering (GATE),⁷ first released in 1995, is an extensive and mature open-source Java toolkit for many aspects of natural language processing and information extraction tasks. ANNIE, A Nearly-New Information Extraction system, is the component for named entity extraction implementing a more traditional recognition model [15]. GATE provides all necessary tools to build a complete system for knowledge graph construction in combination with other components.

⁷ <https://gate.ac.uk/>

NLTK The Natural Language Toolkit (NLTK),⁸ first released in 2001, is one of the most popular Python libraries for natural language processing [7]. It provides a wealth of easy to use API for all traditional text processing tasks and named entity recognition capabilities.

OpenNLP The Apache OpenNLP project⁹ is a Java library for the most common text processing tasks and was first released in 2004 [27]. It provides implementations of a wide selection of machine learning-based NLP research designed to extend data pipelines built with Apache Flink or Spark.

CoreNLP The StanfordNLP research group released their first version of CoreNLP¹⁰ in 2006 as a Java library with Python bindings [52]. It is actively developed and provides NLP tools ranging from traditional rule-based approaches to models from recently published state-of-the-art deep learning research.

spaCy Only recently in 2015, spaCy¹¹ was released as a Python/Cython library, focusing on providing high-performance in terms of processing speed. It also includes language models for over fifty languages and a growing community of extensions. After an initial focus on processing speed, the library now includes high-quality language models based on recent advances in deep learning.

For a detailed comparison of the frameworks mentioned above, we refer readers to the recently published study by Schmitt et al. [68].

2.2 *Named Entity Linking (NEL)*

The problem of linking named entities is rooted in a wide range of research areas. Through named entity linking, the strings discovered by NER are matched to entities in an existing knowledge graph or extend it. Wikidata is a prevalent knowledge graph for many use cases. Typically, there is no identical string match for an entity mention discovered in the text and the knowledge graph. Organizations are rarely referred to by their full legal name, but rather an acronym or colloquial variation of the full name. For example, VW could refer to Vorwerk, a manufacturer for household appliances, or Volkswagen, which is also known as Volkswagen Group or Volkswagen AG. At the time of writing, there are close to 80 entries in Wikidata¹² when searching for “Volkswagen”, excluding translations, car models, and other non-organization entries. Entity linking approaches use various features to match the correct real-world entity. These features are typically based on the entity mention itself or information about the context in which it appeared. Thereby, they face similar challenges and use comparable approaches as research in record linkage and duplicate detection. Shen et al. [70] provide a comprehensive overview of applications, challenges, and

⁸ <http://nltk.org/>

⁹ <https://opennlp.apache.org>

¹⁰ <https://nlp.stanford.edu/software/>

¹¹ <https://spacy.io/>

¹² <https://www.wikidata.org/w/index.php?search=volkswagen>

a survey of the main approaches. As mentioned earlier, there are three main challenges when linking named entities, namely name variations, entity ambiguity, and unlinkable entities. In this subsection, we discuss these challenges using examples to illustrate them better. We also present common solutions to resolve them and close with an overview of entity linking systems.

Name Variations. A real-world entity is referred to in many different ways, such as the full official name, abbreviations, colloquial names, various known aliases, or simply with typos. These variations increase the complexity of finding the correct match in the knowledge base. For example, Dr. Ing. h.c. F. Porsche GmbH, Ferdinand Porsche AG, Porsche A.G. are some name variations for the German car manufacturer Porsche commonly found in business news. Entity linking approaches traditionally take two main steps [70]. The first step selects candidate entries for the currently processed mention from the knowledge base. The second step performs the actual linking by choosing the correct candidate. The candidate generation reduces the number of possible matches, as the disambiguation can become computationally expensive. The most common approach is to use fuzzy string comparisons, such as an edit-distance like the Levenshtein distance or the Jaccard index for overlapping tokens. Additionally, a few rules for name expansion can generate possible abbreviations or extract potential acronyms from names. These rules should use domain-specific characteristics, for example, common legal forms (Ltd. \mapsto Limited) as well as names (International Business Machines \mapsto IBM). If an existing knowledge base is available, a dictionary of known aliases can be derived.

Entity Ambiguity. A mentioned entity could refer to multiple entries in the knowledge graph. For example, Volkswagen could not only refer to the group of car manufacturers, but also the financial services, international branches, or to the local car dealership. Only the context, the company mention appears in, may help identify the correct entry, by taking keywords within the sentence (local context) or the document (global context) into account. The entity disambiguation, also called *entity ranking*, selects the correct entry among the previously generated set of candidates

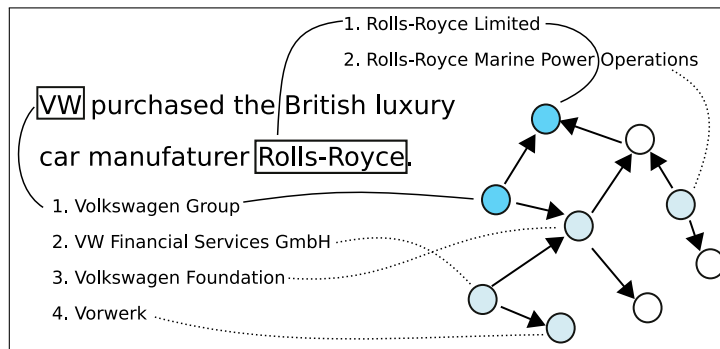


Fig. 2 Example for ranking and linking company mentions to the correct entity in a set of candidates from the knowledge graph.

of possible matches from the knowledge base. This second linking step aims to estimate the likelihood of a knowledge base entry being the correct disambiguation a given mention. These scores create a ranking of candidates. Typically, the one with the highest score is usually chosen to be the correct match. Generally, ranking models follow either a supervised or unsupervised approach. Supervised methods use annotated data mentions are explicitly linked to entries in the knowledge base to train classifiers, ranking models, probabilistic models, or graph-based methods. When there is no annotated corpus available, data-driven unsupervised learning or information retrieval methods can be used. Shen et al. [70] further categorize both approaches into three paradigms. *Independent ranking methods* consider entity mentions individually without leveraging relations between other mentions in the same document and only focusing on the text directly surrounding it. On the other hand, *collective ranking methods* assume topical coherence for all entity mentions in one document and link all of them collectively. Lastly, *collaborative ranking methods* leverage the textual context of similar entity mentions across multiple documents to extend the available context information.

Unlinkable Entities. Novel entities have no corresponding entries in the knowledge graph yet. It is important to note that NEL approaches should identify such cases and not just pick the best possible match. Unlinkable entities may be added as new entries to the knowledge graph. However, this depends on the context and its purpose. Suppose HBO_2 was found in a sentence and is supposed to be linked to a knowledge base of financial entities. If the sentence is about inorganic materials, this mention most likely refers to metaboric acid and should be dismissed. Whereas, in a pharmaceutical context, it might refer to the medical information systems firm HBO & Company. In that case, it should be added as a new entity and not linked to the already existing television network HBO. Entity linking systems deal with this in different ways. They commonly introduce a NIL entity, which represents a universal unlinkable entity, into the candidate set or a threshold for the likelihood score.

Other challenges. Growing size and heterogeneity of KGs are further challenges. Scalability and speed is a fundamental issue for almost all entity ranking systems. A key part to solve this challenge is a fast comparison function to generate candidates with a high recall to reduce the number of computations of similarity scores. State-of-the-art approaches that use vector representations have the advantage that nearest neighborhood searches within a vector space are almost constant [41]. However, training them requires large amounts of data, which might not be available in specific applications. Furthermore, targeted adaptations are not as trivial as with rule-based or feature-based systems. Another challenge for entity ranking systems are heterogeneous sources. Whereas multi-language requirements can be accounted for by separate models, evolving information over time impose other difficulties. Business news or other sources continuously generate new facts that could enrich the knowledge graph further. However, with a growing knowledge graph, the characteristics of the data change. Models tuned on specific characteristics or trained on a previous state of the graph may need regular updates.

Approaches. There are numerous approaches for named entity linking. Traditional approaches use textual fragments surrounding the entity mention to improve the linking quality over just using a fuzzy string match. Complex joint reasoning and ranking methods negatively influence the disambiguation performance in cases with large candidate sets. Zou et al. [83] use multiple bagged ranking classifiers to calculate a consensus decision. This way, they can operate on subsets of large candidate sets and exploit previous disambiguation decisions whenever possible. As mentioned before, not every entity mention can be linked to an entry in the knowledge graph. On the other hand, including the right entities in the candidate set is challenging due to name variations and ambiguities. Typically, there is a trade-off between the precision (also called linking correctness rate) of a system and its recall (also called linking coverage rate). For example, simply linking mentions of VW in news articles to the most popular entry in the knowledge graph is probably correct. All common aliases are well known and other companies with similar acronyms appear less frequently in the news, which leads to high precision and recall. In particular applications, this is more challenging. Financial filings often contain references to numerous subsidiaries with very similar names that need to be accurately linked. CoHEEL is an efficient method that uses random walks to combine a precision-oriented and a recall-oriented classifier [25]. They achieve wide coverage while maintaining a high precision, which is of high importance for business analytics.

The research on entity linking shifted towards deep learning and embedding based approaches in recent years. Generally, they learn high-dimensional vector representations of tokens in the text and knowledge graph entries. Zwicklbauer et al. [85] use such embeddings to calculate the similarity between an entity mention and its respective candidates from the knowledge graph. Given a set of training data in which the correct links are annotated in the text, they learn a robust similarity measure. Others use the annotated mentions in the training data as special tokens in the vocabulary and project words and entities into a common vector space [81, 21]. The core idea behind DeepType [53] is to support the linking process by providing type information about the entities from an existing knowledge graph to the disambiguation process, which they train in an end-to-end fashion. Such approaches require existing knowledge graphs and large sets of training data. Although this can be generated semi-automatically with open information extraction methods, maintaining a high quality can be challenging. Labeling high-quality training data manually is infeasible while maintaining high coverage. Active learning methods can significantly reduce the required amount of annotated data. DeepMatcher offers a ready-to-use implementation of a neural network that makes use of fully automatically learned attribute and word embeddings to train an entity similarity function with targeted human annotation [45].

2.3 Relationship Extraction (RELEX)

Relationship extraction identifies triples of two entities and their relation that appear in a text. Approaches follow one of two strategies: mining of *open-domain* triples or *fixed-domain* triples. In an open-domain setting, possible relations are not specified in advance and typically just use a keyword between two entities. Stanford’s OpenIE [3] is a state-of-the-art information extraction system that splits sentences into sets of clauses. These are then shortened and segmented into triples. Figure 3 shows the relations extracted by OpenIE from the example used in Fig. 1. One such extracted triple would be (BMW, supply, Rolls-Royce).

Such a strategy is useful in cases where no training data or no ontology is available. An ontology is a schema (for a knowledge graph) that defines the types of possible relations and entities. In the following section, we provide more details on standardized ontologies and refinement. One disadvantage of open-domain extraction is that synonymous relationships lead to multiple edges in the knowledge graph. Algorithms can disambiguate the freely extracted relations after enriching the knowledge graph with data from all available text sources. In a fixed-domain setting, all possible relation types are known ahead of time. Defining a schema has the advantage that downstream applications can refer to pre-defined relation types. For example, in Fig. 1 we consider relations such as ORG owns ORG, which is implicitly matched by “*VW purchased Rolls-Royce*”.

The naïve way to map relations mentioned in the text to a schema is to provide a dictionary for each relation type. An algorithm can automatically extend a dictionary from a few manually annotated sentences with relation triples or a seed dictionary. Agichtein and Gravano published the very famous *Snowball* algorithm, which follows this approach [1]. In multiple iterations, the algorithm grows the dictionary based on an initially small set of examples. This basic concept is applied in semi-supervised training to improve more advanced extraction models. The collection of seed examples can be expanded after every training iteration. This process is also called distant supervision. However, it can only detect relationship types already contained in the knowledge graph and can not discover new relationship types. A comprehensive discussion of distant supervision techniques for relation extraction is

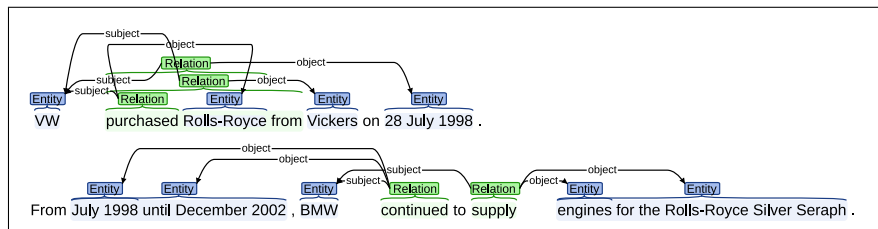


Fig. 3 Relations recognized by OpenIE in text from Fig. 1, output is visualized by CoreNLP¹³

¹³ An online demo of CoreNLP is available at <https://corenlp.run/>

provided by Smirnova [71]. Zuo et al. demonstrated the domain-specific challenges of extracting company relationships from text [84].

Recent approaches mostly focus on deep learning architectures to identify relations in a sequence of words. Wang et al. [78] use convolutional layers and attention mechanisms to identify the most relevant syntactic patterns for relation extraction. Others employ recurrent models to focus on text elements in sequences of variable length [33]. Early approaches commonly used conditional random fields (CRF) on parse trees, representing the grammatical structure and dependencies in a sentence. Nguyen et al. [48] combine modern neural BiLSTM architectures with CRFs for an end-to-end trained model to improve performance. Based on the assumption that if two entities are mentioned in the same text segment, Soares et al. [73] use BERT [16] to learn relationship embeddings. These embeddings are similar to dictionaries with the advantage, that embedding vectors can be used to easily identify the matching relation type for ambiguous phrases in the text.

3 Refining the Knowledge Graph

In the previous section, we described the key steps in constructing a knowledge graph, namely named entity extraction, entity linking, and relationship extraction. This process produces a set of triples from a given text corpus that forms a knowledge graph's nodes and edges. As we have shown in the previous section, compiling a duplicate-free knowledge graph is a complex and error-prone task. Thus, these triples need refinement and post-processing to ensure a high-quality knowledge graph. Any analysis based on this graph requires the contained information to be as accurate and complete as possible.

Manual refinement and standards are inevitable for high-quality results. For better interoperability, the Object Management Group, the standards consortium that defined UML and BPMN, among other things, specified the Financial Industry Business Ontology (FIBO).¹⁴ This ontology contains standard identifiers for relationships and business entities. The Global Legal Identifier Foundation (GLEIF)¹⁵ is an open resource that assigns unique identifiers to legal entities and contains statistics for around 1.5 million entries at the time of writing.

Using existing knowledge graphs as a reference together with standardized ontologies is a good foundation for the manual refinement process. However, the sheer size of these datasets requires support by automated mechanisms in an otherwise unattainable task. With CurEx, Loster et al. [37] demonstrate the entire pipeline of curating company networks extracted from text. They discuss the challenges of this system in the context of its application in a large financial institution [38]. Knowledge graphs about company relations are also handy beyond large scale analyses of the general market situation. For example, changes in the network, as reported in

¹⁴ <https://www.omg.org/spec/EDMC-FIBO/BE/>

¹⁵ <https://search.gleif.org/>

SEC filings,¹⁶ are of particular interest to analysts. Sorting through all mentioned relations is typically impractical. Thus, automatically identifying the most relevant reported business relationships in newly released filings can significantly support professionals in their work. Repke et al. [60] use the surrounding text, where a mentioned business relation appears, to create a ranking to enrich dynamic knowledge graphs. There are also other ways to supplement the available information about relations. For example, a company network with weighted edges can be constructed from stock market data [29]. The authors compare the correlation of normalized stock prices with relations extracted from business news in the same time-frame and found, that frequently co-mentioned companies oftentimes share similar patterns in the movements of their stock prices.

Another valuable resource for extending knowledge graphs are internal documents, as they contain specialized and proprietary domain knowledge. For example, the graph can also be extended beyond just company relations and include key personnel and semantic information. In the context of knowledge graph refinement, it is essential to provide high quality and clean input data to the information extraction pipeline. The Enron corpus [30], for example, has been the basis for a lot of research in many fields. This corpus contains over 500,000 emails from more than 150 Enron employees. The text's structure and characteristics in emails are typically significantly different from that of news, legal documents, or other reports. With Quagga,¹⁷ we published a deep learning-based system to pre-process email text [55]. It identifies the parts of an email text that contains the actual content. It disregards additional elements, such as greetings, closing words, signatures, or automatically inserted meta-data when forwarding or replying to emails. This meta-data could extend the knowledge graph with information about who is talking to whom about what, which is relevant for internal investigations.

4 Analyzing the Knowledge Graph

Knowledge about the structure of the market is a highly valuable asset. This section focuses on specific applications in the domain of business intelligence for economics and finance. Especially financial institutions have to have a detailed overview of the entire financial market, particularly the network of organizations in which they invest. Therefore, Ronnqvist et al. [63] extracted bank networks from text to quantify interrelations, centrality, and determinants.

In Europe, banks are required by law to estimate their systemic risk. The network structure of the knowledge graph allows the investigation of many financial scenarios, such as the impact of corporate bankruptcy on other market participants within the network. In this particular scenario, the links between the individual market participants can be used to determine which companies are affected by bankruptcy and

¹⁶ <https://www.sec.gov/edgar.shtml>

¹⁷ <https://github.com/HPI-Information-Systems/QuaggaLib>

to what extent. Balance sheets and transactions alone would not suffice to calculate that risk globally, as it only provides an ego-network and thus a limited view of the market. Thus, financial institutions have to integrate their expertise in measuring the economic performance of their assets and a network of companies to simulate how the potential risk can propagate. Constantin et al. [14] use data from the financial network and market data covering daily stock prices of 171 listed European banks to predict bank distress.

News articles are a popular source of information for analyses of company relationships. Zheng and Schwenkler demonstrate that company networks extracted from news can be used to measure financial uncertainty and credit risk spreading from a distressed firm [82]. Others also found that the return of stocks reflects economic linkages derived from text [67]. We have shown that findings like this are controversial [29]. Due to the connectedness within industry sectors and the entire market, stock price correlation patterns are very common. Large companies and industry leaders heavily influence the market and appear more frequently in business news than their smaller competitors. Additionally, news are typically slower than movements on the stock market, as insiders receive information earlier through different channels. Thus, observation windows have to be in sync with the news cycle for analyses in this domain.

News and stock market data can then be used to show, for example, how the equity market volatility is influenced by newspapers [4]. Chahrour et al. [11] make similar observations and construct a model to show the relation between media coverage and demand-like fluctuations orthogonal to productivity within a sector. For models like this to work, company names have to be detected in the underlying texts and linked to the correct entity in a knowledge graph. Hoberg and Phillips extract an information network from product descriptions in 10-K statements filed with the SEC [26]. With this network, they examine how industry market structure and competitiveness change over time.

These examples show that knowledge graphs extracted from text can model existing hypotheses in economics. A well-curated knowledge graph that aggregates large amounts of data from a diverse set of sources would allow advanced analyses and market simulations.

5 Exploring the Knowledge Graph

Knowledge graphs of financial entities enable numerous downstream tasks. These include automated enterprise valuation, identifying the sentiment towards a particular company, or discovering political and company networks from textual data. However, knowledge graphs can also support the work of accountants, analysts, and investigators. They can query and explore relationships and structured knowledge quickly to gather the information they need. For example, visual analytics helps to monitor the financial stability in company networks [18]. Typically, such applications display small sub-graphs of the entire company network as a so-called

node-link diagram. Circles or icons depict companies connected by straight lines. The most popular open-source tools for visualizing graphs are Cytoscape [49] and Gephi [5]. They mostly focus on visualization rather than capabilities to interactively explore the data. Commercial platforms, on the other hand, such as the NUIX Engine¹⁸ or Linkurious,¹⁹ offer more advanced capabilities. These include data pre-processing and analytics frequently used in forensic investigations, e.g., by journalists researching the Panama papers leak.

There are different ways to visualize a network. Most commonly, by node-link diagrams as described above. However, already with a small number of nodes and edges, the readability is hard to maintain [51]. Edge bundling provides for better clarity of salient high-level structures [35]. The downside is that individual edges can become impossible to follow. Other methods for network visualization focus on the adjacency matrix of the graph. Each row and column corresponds to a node in the graph and cells are colored according to the edge weight or remain empty. The greatest challenge is to arrange the rows and columns in such a way, that salient structures become visible. Sankey diagrams are useful to visualize hierarchically clustered networks to show the general flow of information or connections [12]. For more details have a look at the excellent survey of network visualization methods by Gibson et al. [22]. There is no one best type of visualization. It depends on the specific application to identify the ideal representation to explore the data.

Repke et al. developed “Beacon in the Dark” [59], a platform incorporating various modes of exploration. It includes a full system pipeline to process and integrate structured and unstructured data from email and document corpora. It also consists of an interface with coordinated multiple views to explore the data in a topic sphere (semantics), by tags that are automatically assigned, and the communication and entity network derived from meta-data and text. The system goes beyond traditional approaches by combining communication meta-data and integrating additional information using advanced text mining methods and social network analysis. The objectives are to provide a data-driven *overview* of the dataset to determine initial leads without knowing anything about the data. The system also offers extensive filters and rankings of available information to focus on relevant aspects and finding necessary data. With each interaction, the interface components update to provide the appropriate context in which a particular information snippet appears.

6 Semantic Exploration using Visualisations

Traditional interfaces for knowledge graphs typically only support node-to-node exploration with basic search and filtering capabilities. In the previous section, we have shown that some of them also integrate the underlying source data. However, they only utilize the meta-data, not the semantic context provided by the text in which

¹⁸ NUIX Analytics extracts and indexes knowledge from unstructured data (<https://www.nuix.com>)

¹⁹ Linkurious Enterprise is a graph visualization and analysis platform (<https://linkurio.us/>)

entities appear. Furthermore, prior knowledge about the data is required to formulate the right queries as a starting point for exploration. In this section, we focus on methods to introduce semantics into the visualization of knowledge graphs. This integration enables users to explore the data more intuitively and provides better explanatory navigation.

Semantic information can be added based on the context in which entity mentions appear. The semantics could simply be represented by keywords from the text surrounding an entity. The benefit is that users not only interact with raw network data but with descriptive information. Text mining can be used to automatically enrich the knowledge about companies, for example, by assigning the respective industry or sentiment towards products or an organization itself. Such an approach is even possible without annotated data. Topic models assign distributions of topics to documents and learn which words belong to which topics. Early topic models, such as latent semantic indexing does so by correlating semantically related terms from a collection of text documents [20]. These models iteratively update the distributions, which is computationally expensive for large sets of documents and long dictionaries. Latent semantic analysis, the foundation for most current topic models, uses co-occurrence of words [31]. Latent Dirichlet allocation jointly models topics as distributions over words and documents as distributions over topics [8]. This allows to summarize large document collections by means of topics.

Recently, text mining research has shifted towards embedding methods to project words or text segments into a high-dimensional vector space. The semantic similarity between words can be calculated based on distance measures between their vectors. Initial work in that area includes word2vec [44] and doc2vec [32]. More recent popular approaches such as BERT [16] better capture the essential parts of a text. Similar approaches can also embed graph structures. RDF2vec is an approach that generates sequences of connections in the graph leveraging local information about its sub-structures [62]. They show some applications that allow the calculation of similarities between nodes in the graph. There are also specific models to incorporate text and graph data to optimize the embedding space. Entity-centric models directly assign vectors to entities instead of just tokens. Traditionally, an embedding model assumes a fixed dictionary of known words or character n-grams that form these words. By annotating the text with named entity information before training

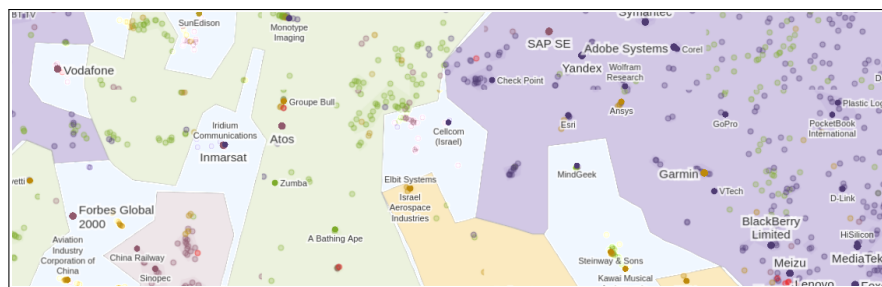


Fig. 4 Screenshot of part of the Cartograph map of organizations and their sectors

the model, unique multi-word entries in the dictionary directly relate to known entities. Almasian et al. propose such a model for entity-annotated texts [2]. Other interesting approaches build networks of co-occurring words and entities. TopExNet uses temporal filtering to produce entity-centric networks for topic exploration in news streams [75]. For a survey of approaches and applications of knowledge graph embeddings, we refer readers to [79].

Topic models, document embeddings, and entity embeddings are useful tools for systematic data analysis. However, on their own, they are not directly useable. In the context of book recommendations, embeddings have been used to find similar books using combinations of embeddings for time and place of the plot [61]. Similar approaches could be applied in the domain of financial entities, for example, to discover corresponding companies in a different country. In use-cases without prior knowledge, it might be particularly helpful to get an overview of all the data. Also, for monitoring purposes, a bird's-eye view of the entire dataset can be beneficial. The most intuitive way is to organize the information in the form of an interactive map. Sarlin et al. [66] used self-organizing maps to arrange economic sectors and countries to create maps. Coloring the maps enables them to visually compare different financial stability metrics across multiple time-frames around periods with high inflation rates or an economic crisis.

The idea of semantic landscapes is also popular in the area of patent research. The commercial software Themescape by Derwent²⁰ produces landscapes of patents that users can navigate similar to a geographical map. Along with other tools, they enable experts to find related patents or identify new opportunities quickly. Smith et al. built a system to transform token co-occurrence information in texts to semantic patterns. Using statistical algorithms, they generate maps of words that can be used for content analysis in knowledge discovery tasks [72]. Inspired by that, the New York public library made a map of parts of their catalog.²¹ Therefore they use a force-based network layout algorithm to position the information. It uses the analogy of forces that attract nodes to one another when connected through an edge or otherwise repel them. The network they use is derived from co-occurring subject

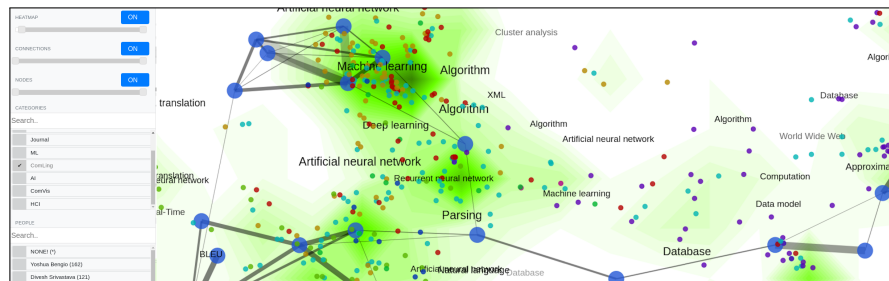


Fig. 5 Screenshot of the *MODiR* interface prototype showing an excerpt of a citation network.

²⁰ <https://clarivate.com/derwent>

²¹ <https://www.nypl.org/blog/2014/07/31/networked-catalog>

headings and terms, which were manually assigned tags to organize their catalog. Sen et al. created a map of parts of the Wikipedia in their Cartograph project [69]. This map, as shown in Fig. 4, uses embedded pages about companies and dimensionality reduction to project the information on a two-dimensional canvas [40, 50]. Structured meta-data about pages is used to compute borders between “countries” representing different industry sectors. Maps like this provide an intuitive alternative interface for users to discover related companies. Most recently, the OpenSyllabus Project²² released their interactive explorer. Like Cartograph, this enables users to navigate through parts of the six million syllabi collected by the project. To do so, they first create a citation network of all publications contained in the visualization. Using this network, they learn a node embedding [24] and reduce the number of dimensions for rendering [43].

The approaches presented above offer promising applications in business analytics and exploring semantically infused company networks. However, even though the algorithms use networks to some extent, they effectively only visualize text and rely on manually tagged data. Wikipedia, library catalogs and the syllabi corpus are datasets that are developed over many years by many contributors who organize the information into structured ontologies. In business applications, data might not always have this additional information available, and it is too labor-intensive to curate the data manually. Furthermore, when it comes to analyzing company networks extracted from text, the data is comprised of both the company network and data provenance information. The methods presented above only visualize either the content data or the graph structure. In data exploration scenarios, the goal of getting a full overview of the dataset at hand is insurmountable with current tools. We provide a solution that incorporates both, the text sources *and* the entity network, into exploratory landscapes [56]. We first embed the text data and then use multiple objectives to optimize for a good network layout and semantically correct layout of source documents during the dimensionality reduction [58]. Figure 5 shows a small demonstration of the resulting semantic-infused network layout [57]. Users exploring such data, e.g., journalists investigating leaked data or young scientists starting research in an unfamiliar field, need to be able to interact with the visualization. Our prototype allows users to explore the generated landscape as a digital map with zooming and panning. The user can select from categories or entities to shift the focus, highlight characterizing keywords, and adjust a heatmap based on the density of points to only consider related documents. We extract region-specific keywords and place them on top of the landscape. This way, the meaning of an area becomes clear and supports fast navigation.

²² Open Syllabus Explorer visualization shows the 164,720 texts (<http://galaxy.opensyllabus.org/>)

7 Conclusion

In this chapter, we provided an overview of methods to automatically construct a knowledge graph from text, particularly a network of financial entities. We described the pipeline starting from named entity recognition, over linking and matching those entities to real-world entities, to extracting the relationships between them from text. We emphasized the need to curate the extracted information, which typically contains errors that could negatively impact its usability in subsequent applications. There are numerous use-cases that require knowledge graphs connecting economic, financial, and business related information. We have shown how these knowledge graphs are constructed from heterogeneous textual documents and how they can be explored and visualized to support investigations, analyses, and decision making.

References

- [1] Agichtein, E., Gravano, L.: *Snowball*: Extracting relations from large plain-text collections. In: Proceedings of the Joint Conference on Digital Libraries (JCDL), pp. 85–94. ACM Press (2000)
- [2] Almasian, S., Spitz, A., Gertz, M.: Word embeddings for entity-annotated texts. In: Proceedings of the European Conference on Information Retrieval (ECIR), *Lecture Notes in Computer Science*, vol. 11437, pp. 307–322. Springer-Verlag (2019)
- [3] Angeli, G., Premkumar, M.J.J., Manning, C.D.: Leveraging linguistic structure for open domain information extraction. In: Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL), pp. 344–354. Association for Computational Linguistics (2015)
- [4] Baker, S.R., Bloom, N., Davis, S.J., Kost, K.J.: Policy news and stock market volatility. Working Paper 25720, National Bureau of Economic Research (2019)
- [5] Bastian, M., Heymann, S., Jacomy, M.: Gephi: An open source software for exploring and manipulating networks. In: Proceedings of the International Semantic Web Conference (ISWC). The AAAI Press (2009)
- [6] Bikel, D.M., Miller, S., Schwartz, R.M., Weischedel, R.M.: Nymble: A high-performance learning name-finder. In: Applied Natural Language Processing Conference (ANLP), pp. 194–201. Association for Computational Linguistics (1997)
- [7] Bird, S., Klein, E., Loper, E.: *Natural Language Processing with Python*. O’Reilly (2009)
- [8] Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent Dirichlet allocation. *Journal of machine Learning research* **3**(Jan), 993–1022 (2003)
- [9] Bollacker, K.D., Evans, C., Paritosh, P., Sturge, T., Taylor, J.: Freebase: A collaboratively created graph database for structuring human knowledge. In:

- Proceedings of the ACM Conference on Management of Data (SIGMOD), pp. 1247–1250 (2008)
- [10] Chabin, M.A.: Panama papers: A case study for records management? *Brazilian Journal of Information Science: Research Trends* **11**(4), 10–13 (2017)
 - [11] Chahrour, R., Nimark, K., Pitschner, S.: Sectoral media focus and aggregate fluctuations. *Swedish House of Finance Research Paper Series 19-12*, SSRN (2019)
 - [12] Chang, C., Bach, B., Dwyer, T., Marriott, K.: Evaluating perceptually complementary views for network exploration tasks. In: *Proceedings of the SIGCHI conference on Human Factors in Computing Systems (CHI)*, pp. 1397–1407. ACM Press (2017)
 - [13] Coddington, M.: Clarifying journalism’s quantitative turn: A typology for evaluating data journalism, computational journalism, and computer-assisted reporting. *Digital Journalism* **3**(3), 331–348 (2015)
 - [14] Constantin, A., Peltonen, T.A., Sarlin, P.: Network linkages to predict bank distress. *Journal of Financial Stability* **35**, 226–241 (2018)
 - [15] Cunningham, H., Maynard, D., Bontcheva, K., Tablan, V.: GATE: A framework and graphical development environment for robust NLP tools and applications. In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 168–175. Association for Computational Linguistics (2002)
 - [16] Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT)*, pp. 4171–4186. Association for Computational Linguistics (2019)
 - [17] Faruqui, M., Padó, S.: Training and evaluating a german named entity recognizer with semantic generalization. In: *Proceedings of the Conference on Natural Language Processing (KONVENS)*, pp. 129–133 (2010)
 - [18] Flood, M.D., Lemieux, V.L., Varga, M., Wong, B.W.: The application of visual analytics to financial stability monitoring. *Journal of Financial Stability* **27**, 180–197 (2016)
 - [19] Franke, K., Srihari, S.N.: Computational forensics: Towards hybrid-intelligent crime investigation. In: *Proceedings of the International Symposium on Information Assurance and Security (IAS)*, pp. 383–386. IEEE, New York City, USA (2007)
 - [20] Furnas, G.W., Deerwester, S.C., Dumais, S.T., Landauer, T.K., Harshman, R.A., Streeter, L.A., Lochbaum, K.E.: Information retrieval using a singular value decomposition model of latent semantic structure. In: *Proceedings of the ACM Conference on Information Retrieval (SIGIR)*, pp. 465–480. ACM Press (1988)
 - [21] Ganea, O., Hofmann, T.: Deep joint entity disambiguation with local neural attention. In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 2619–2629. Association for Computational Linguistics (2017)

- [22] Gibson, H., Faith, J., Vickers, P.: A survey of two-dimensional graph layout techniques for information visualisation. *Information Visualization* **12**(3-4), 324–357 (2013)
- [23] Grishman, R., Sundheim, B.: Message understanding conference- 6: A brief history. In: *Proceedings of the International Conference on Computational Linguistics (COLING)*, pp. 466–471 (1996)
- [24] Grover, A., Leskovec, J.: node2vec: Scalable feature learning for networks. In: *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pp. 855–864. ACM Press (2016)
- [25] Grütze, T., Kasneci, G., Zuo, Z., Naumann, F.: CohEEL: Coherent and efficient named entity linking through random walks. *Journal of Web Semantics* **37-38**, 75–89 (2016)
- [26] Hoberg, G., Phillips, G.: Text-based network industries and endogenous product differentiation. *Journal of Political Economy* **124**(5), 1423–1465 (2016)
- [27] Ingersoll, G., Morton, T., Farris, A.: *Taming text*. Manning Publications (2012)
- [28] Karthik, M., Marikkannan, M., Kannan, A.: An intelligent system for semantic information retrieval information from textual web documents. In: *International Workshop on Computational Forensics (IWCF)*, pp. 135–146. Springer-Verlag, Heidelberg, Germany (2008)
- [29] Kellermeier, T., Repke, T., Krestel, R.: Mining business relationships from stocks and news. In: V. Bitetta, I. Bordino, A. Ferretti, F. Gullo, S. Pascolutti, G. Ponti (eds.) *Proceedings of the European Conference on Machine Learning (ECML), Lecture Notes in Computer Science*, vol. 11985, pp. 70–84. Springer-Verlag, Heidelberg, Germany (2019)
- [30] Klimt, B., Yang, Y.: The Enron corpus: A new dataset for email classification research. In: *Proceedings of the European Conference on Machine Learning (ECML)*, pp. 217–226. Springer-Verlag, Heidelberg, Germany (2004)
- [31] Landauer, T.K., Foltz, P.W., Laham, D.: An introduction to latent semantic analysis. *Discourse processes* **25**(2-3), 259–284 (1998)
- [32] Le, Q., Mikolov, T.: Distributed representations of sentences and documents. In: *Proceedings of the International Conference on Machine Learning (ICML)*, pp. 1188–1196. JMLR Inc. and Microtome Publishing, Brookline, USA (2014)
- [33] Lee, J., Seo, S., Choi, Y.S.: Semantic relation classification via bidirectional LSTM networks with entity-aware attention using latent entity typing. *Symmetry* **11**(6), 785 (2019)
- [34] Lehmann, J., Isele, R., Jakob, M., Jentzsch, A., Kontokostas, D., Mendes, P.N., Hellmann, S., Morsey, M., van Kleef, P., Auer, S., Bizer, C.: DBpedia - A large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web* **6**(2), 167–195 (2015)
- [35] Lhuillier, A., Hurter, C., Telea, A.: State of the art in edge and trail bundling techniques. *Computer Graphics Forum* **36**(3), 619–645 (2017)
- [36] Loster, M., Hegner, M., Naumann, F., Leser, U.: Dissecting company names using sequence labeling. In: *Proceedings of the Conference "Lernen, Wis-*

- sen, Daten, Analysen” (LWDA), *CEUR Workshop Proceedings*, vol. 2191, pp. 227–238. CEUR-WS.org (2018)
- [37] Loster, M., Naumann, F., Ehmueller, J., Feldmann, B.: Curex: A system for extracting, curating, and exploring domain-specific knowledge graphs from text. In: *Proceedings of the International Conference on Information and Knowledge Management (CIKM)*, pp. 1883–1886. ACM Press (2018)
- [38] Loster, M., Repke, T., Krestel, R., Naumann, F., Ehmueller, J., Feldmann, B., Maspfuhl, O.: The challenges of creating, maintaining and exploring graphs of financial entities. In: *Proceedings of the International Workshop on Data Science for Macro-Modeling with Financial and Economic Datasets (DSMM@SIGMOD)*, pp. 6:1–6:2. ACM Press (2018)
- [39] Loster, M., Zuo, Z., Naumann, F., Maspfuhl, O., Thomas, D.: Improving company recognition from unstructured text by using dictionaries. In: *Proceedings of the International Conference on Extending Database Technology (EDBT)*, pp. 610–619. OpenProceedings.org (2017)
- [40] Maaten, L.v.d., Hinton, G.: Visualizing data using t-SNE. *Journal of Machine Learning Research (JMLR)* **9**, 2579–2605 (2008)
- [41] Malkov, Y.A., Yashunin, D.A.: Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *CoRR abs/1603.09320* (2016)
- [42] McCallum, A., Li, W.: Early results for named entity recognition with conditional random fields, feature induction and web-enhanced lexicons. In: *Proceedings of the Conference on Computational Natural Language Learning (CoNLL)*, pp. 188–191. Association for Computational Linguistics (2003)
- [43] McInnes, L., Healy, J.: UMAP: Uniform manifold approximation and projection for dimension reduction. *CoRR abs/1802.03426* (2018)
- [44] Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: *Proceedings of the Conference on Neural Information Processing Systems (NIPS)*, pp. 3111–3119. NIPS Foundation, Inc., San Diego, USA (2013)
- [45] Mudgal, S., Li, H., Rekatsinas, T., Doan, A., Park, Y., Krishnan, G., Deep, R., Arcaute, E., Raghavendra, V.: Deep learning for entity matching: A design space exploration. In: *Proceedings of the ACM Conference on Management of Data (SIGMOD)*, pp. 19–34. ACM Press (2018)
- [46] Nadeau, D., Sekine, S.: A survey of named entity recognition and classification. *Linguisticae Investigationes* **30**(1), 3–26 (2007)
- [47] Nadeau, D., Turney, P.D., Matwin, S.: Unsupervised named-entity recognition: Generating gazetteers and resolving ambiguity. In: L. Lamontagne, M. Marchand (eds.) *Proceedings of the Conference of the Canadian Society for Computational Studies of Intelligence, Lecture Notes in Computer Science*, vol. 4013, pp. 266–277. Springer-Verlag (2006)
- [48] Nguyen, D.Q., Verspoor, K.: End-to-end neural relation extraction using deep biaffine attention. In: *Proceedings of the European Conference on Information Retrieval (ECIR), Lecture Notes in Computer Science*, vol. 11437, pp. 729–738. Springer-Verlag (2019)

- [49] Otasek, D., Morris, J.H., Bouças, J., Pico, A.R., Demchak, B.: Cytoscape automation: empowering workflow-based network analysis. *Genome biology* **20**(1), 1–15 (2019)
- [50] Pezzotti, N., Lelieveldt, B.P., van der Maaten, L., Höllt, T., Eisemann, E., Vilanova, A.: Approximated and user steerable t-SNE for progressive visual analytics. *IEEE Transactions on Visualization and Computer Graphics (TVCG)* **23**(7), 1739–1752 (2017)
- [51] Pohl, M., Schmitt, M., Diehl, S.: Comparing the readability of graph layouts using eyetracking and task-oriented analysis. In: *Computational Aesthetics 2009: Eurographics Workshop on Computational Aesthetics*, Victoria, British Columbia, Canada, 2009, pp. 49–56 (2009)
- [52] Qi, P., Dozat, T., Zhang, Y., Manning, C.D.: Universal dependency parsing from scratch. In: *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pp. 160–170. Association for Computational Linguistics (2018)
- [53] Raiman, J., Raiman, O.: DeepType: Multilingual entity linking by neural type system evolution. In: *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, pp. 5406–5413. AAAI Press (2018)
- [54] Rau, L.F.: Extracting company names from text. In: *Proceedings of the IEEE Conference on Artificial Intelligence Application*, vol. 1, pp. 29–32. IEEE (1991)
- [55] Repke, T., Krestel, R.: Bringing back structure to free text email conversations with recurrent neural networks. In: *Proceedings of the European Conference on Information Retrieval (ECIR)*, pp. 114–126. Springer-Verlag, Heidelberg, Germany (2018)
- [56] Repke, T., Krestel, R.: Topic-aware network visualisation to explore large email corpora. In: *International Workshop on Big Data Visual Exploration and Analytics (BigVis), Proceedings of the International Conference on Extending Database Technology (EDBT)*, pp. 104–107. CEUR-WS.org (2018)
- [57] Repke, T., Krestel, R.: Exploration interface for jointly visualised text and graph data. In: *Proceedings of the International Conference on Intelligent User Interfaces (IUI)*, pp. 73–74. ACM Press, Geneva, Switzerland (2020)
- [58] Repke, T., Krestel, R.: Visualising large document collections by jointly modeling text and network structure. In: *Proceedings of the Joint Conference on Digital Libraries (JCDL)*, pp. 279–288. ACM Press, Geneva, Switzerland (2020)
- [59] Repke, T., Krestel, R., Edding, J., Hartmann, M., Hering, J., Kipping, D., Schmidt, H., Scordialo, N., Zenner, A.: Beacon in the dark: A system for interactive exploration of large email corpora. In: *Proceedings of the International Conference on Information and Knowledge Management (CIKM)*, pp. 1871–1874. ACM Press (2018)
- [60] Repke, T., Loster, M., Krestel, R.: Comparing features for ranking relationships between financial entities based on text. In: *Proceedings of the International Workshop on Data Science for Macro-Modeling with Financial and Economic Datasets (DSMM@SIGMOD)*, pp. 12:1–12:2. ACM Press (2017)

- [61] Risch, J., Garda, S., Krestel, R.: Book recommendation beyond the usual suspects - embedding book plots together with place and time information. In: Proceedings of the International Conference on Asia-Pacific Digital Libraries (ICADL), *Lecture Notes in Computer Science*, vol. 11279, pp. 227–239. Springer (2018)
- [62] Ristoski, P., Rosati, J., Noia, T.D., Leone, R.D., Paulheim, H.: RDF2Vec: RDF graph embeddings and their applications. *Semantic Web* **10**(4), 721–752 (2019)
- [63] Rönqvist, S., Sarlin, P.: Bank networks from text: interrelations, centrality and determinants. *Quantitative Finance* **15**(10), 1619–1635 (2015)
- [64] Ruder, S., Vulic, I., Søgaard, A.: A survey of cross-lingual word embedding models. *Journal of Artificial Intelligence Research (JAIR)* **65**, 569–631 (2019)
- [65] Samiei, A., Koumarelas, I., Loster, M., Naumann, F.: Combination of rule-based and textual similarity approaches to match financial entities. In: Proceedings of the International Workshop on Data Science for Macro-Modeling, (DSMM@SIGMOD), pp. 4:1–4:2. ACM Press (2016)
- [66] Sarlin, P.: Exploiting the self-organizing financial stability map. *Engineering Applications of Artificial Intelligence* **26**(5-6), 1532–1539 (2013)
- [67] Scherbina, A., Schlusche, B.: Economic linkages inferred from news stories and the predictability of stock returns (2015). SSRN
- [68] Schmitt, X., Kubler, S., Robert, J., Papadakis, M., Traon, Y.L.: A replicable comparison study of NER software: StanfordNLP, NLTK, OpenNLP, SpaCy, Gate. In: International Conference on Social Networks Analysis, Management and Security (SNAMS), pp. 338–343. IEEE (2019)
- [69] Sen, S., Swoap, A.B., Li, Q., Boatman, B., Dippenaar, I., Gold, R., Ngo, M., Pujol, S., Jackson, B., Hecht, B.: Cartograph: Unlocking spatial visualization through semantic enhancement. In: Proceedings of the International Conference on Intelligent User Interfaces (IUI), pp. 179–190. ACM Press, Geneva, Switzerland (2017)
- [70] Shen, W., Wang, J., Han, J.: Entity linking with a knowledge base: Issues, techniques, and solutions. *IEEE Transactions on Knowledge and Data Engineering (TKDE)* **27**(2), 443–460 (2015)
- [71] Smirnova, A., Cudré-Mauroux, P.: Relation extraction using distant supervision: A survey. *ACM Computing Surveys* **51**(5), 106:1–106:35 (2019)
- [72] Smith, A.E., Humphreys, M.S.: Evaluation of unsupervised semantic mapping of natural language with leximancer concept mapping. *Behavior research methods* **38**(2), 262–279 (2006)
- [73] Soares, L.B., FitzGerald, N., Ling, J., Kwiatkowski, T.: Matching the blanks: Distributional similarity for relation learning. In: Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL), pp. 2895–2905. Association for Computational Linguistics (2019)
- [74] Socher, R., Chen, D., Manning, C.D., Ng, A.Y.: Reasoning with neural tensor networks for knowledge base completion. In: Proceedings of the Conference on Neural Information Processing Systems (NIPS), pp. 926–934 (2013)

- [75] Spitz, A., Almasian, S., Gertz, M.: TopExNet: Entity-centric network topic exploration in news streams. In: Proceedings of the International Conference on Web Search and Data Mining (WSDM), pp. 798–801. ACM Press (2019)
- [76] Suchanek, F.M., Kasneci, G., Weikum, G.: YAGO: a core of semantic knowledge. In: Proceedings of the International World Wide Web Conference (WWW), pp. 697–706 (2007)
- [77] Vrandečić, D., Krötzsch, M.: Wikidata: A free collaborative knowledgebase. *Communications of the ACM* **57**(10), 78–85 (2014)
- [78] Wang, L., Cao, Z., de Melo, G., Liu, Z.: Relation classification via multi-level attention CNNs. In: Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL). ACM Press (2016)
- [79] Wang, Q., Mao, Z., Wang, B., Guo, L.: Knowledge graph embedding: A survey of approaches and applications. *IEEE Transactions on Knowledge and Data Engineering (TKDE)* **29**(12), 2724–2743 (2017)
- [80] Yadav, V., Bethard, S.: A survey on recent advances in named entity recognition from deep learning models. In: Proceedings of the International Conference on Computational Linguistics (COLING), pp. 2145–2158. Association for Computational Linguistics (2018)
- [81] Yamada, I., Shindo, H., Takeda, H., Takefuji, Y.: Joint learning of the embedding of words and entities for named entity disambiguation. In: Proceedings of the Conference on Computational Natural Language Learning (CoNLL), pp. 250–259. Association for Computational Linguistics (2016)
- [82] Zheng, H., Schwenkler, G.: The network of firms implied by the news. ESRB Working Paper Series 108, European Systemic Risk Board (2020)
- [83] Zuo, Z., Kasneci, G., Grütze, T., Naumann, F.: BEL: Bagging for entity linking. In: J. Hajic, J. Tsujii (eds.) Proceedings of the International Conference on Computational Linguistics (COLING), pp. 2075–2086. Association for Computational Linguistics (2014)
- [84] Zuo, Z., Loster, M., Krestel, R., Naumann, F.: Uncovering business relationships: Context-sensitive relationship extraction for difficult relationship types. In: Lernen, Wissen, Daten, Analysen (LWDA) Conference Proceedings, *CEUR Workshop Proceedings*, vol. 1917, p. 271. CEUR-WS.org (2017)
- [85] Zwicklbauer, S., Seifert, C., Granitzer, M.: Robust and collective entity disambiguation through semantic embeddings. In: Proceedings of the ACM Conference on Information Retrieval (SIGIR), pp. 425–434. ACM Press (2016)