



Provided by the author(s) and NUI Galway in accordance with publisher policies. Please cite the published version when available.

| | |
|-----------------------------|---|
| Title | Developing a Dataset for Technology Structure Mining |
| Author(s) | QasemiZadeh, Behrang; Buitelaar, Paul; Monaghan, Fergal |
| Publication Date | 2010 |
| Publication Information | Qasemizadeh, Behrang; Buitelaar, Paul; Monaghan, Fergal (2010) Developing a Dataset for Technology Structure Mining. Conference Paper |
| Publisher | IEEE |
| Link to publisher's version | http://dx.doi.org/10.1109/ICSC.2010.73 |
| Item record | http://www.deri.ie/sites/default/files/publications/bq_ieee_icsc_bare_conf.pdf ; http://hdl.handle.net/10379/4514 |
| DOI | http://dx.doi.org/10.1109/ICSC.2010.73 |

Downloaded 2022-05-21T07:05:54Z

Some rights reserved. For more information, please see the item record link above.



Developing a Dataset for Technology Structure Mining

Behrang QasemiZadeh, Paul Buitelaar, Fergal Monaghan
Unit for Natural Language Processing, DERI
National University of Ireland, Galway
Email: firstname.lastname@deri.org

Abstract—This paper describes steps that have been taken to construct a development dataset for the task of Technology Structure Mining. We have defined the proposed task as the process of mapping a scientific corpus into a labeled digraph named a Technology Structure Graph as described in the paper. The generated graph expresses the domain semantics in terms of interdependencies between pairs of technologies that are named (introduced) in the target scientific corpus. The dataset comprises a set of sentences extracted from the ACL Anthology Corpus. Each sentence is annotated with at least two technologies in the domain of Human Language Technology and the interdependence between them. The annotations - technology mark-up and their interdependencies - are expressed at two layers: lexical and termino-conceptual. Lexical representation of technologies comprises varying lexicalizations of a technology. However, at the termino-conceptual layer all these lexical variations refer to the same concept. We have adopted the same approach for representing Semantic Relations; at the lexical layer a semantic relation is a predicate i.e. defined based on the sentence surface structure; however at the termino-conceptual layer semantic relations are classified into conceptual relations either taxonomic or non-taxonomic. Moreover, the contexts that interdependencies are extracted from are classified into five groups based on the linguistic criteria and syntactic structure that are identified by the human annotators. The dataset initially comprises of 482 sentences. We hope this effort results in a benchmark that can be used for the technology structure mining task as defined in the paper.

Keywords-Technology Structure Mining, Text Mining, NLP

I. INTRODUCTION

Management of Technology (MoT) [1] is a strategic research topic dealing with innovation, efficiency and organization structure management in rapidly changing technology world. Started in the 60s, a long discussed topic in MoT is technology-structure relationships [2]. Among the category definitions for empirical technology-structure research is Technology Interdependence. Technology Interdependence potentially can be used for “minding the technology gap” as defined by Bailey et al [3]:

“We define a technology gap as the space in a work flow between two technologies wherein the output of the first technology is meant to be the input to the second one.”

The automatic extraction of such information faces several established research challenges in Information Extraction

and Natural Language Processing: Named Entity Recognition (NER) [4]; Semantic Role Identification [5]; and Relation Extraction (RE) [6], [7]. In a broader sense, solutions to these challenges feed into two emerging Natural Language Understanding and Semantic Computing research application areas: Open (Domain) Information Extraction (OIE) [8]; and Ontology Learning (OL) [9]. We classify the task of Technology Structure Mining as an activity situated between OIE and OL.

One of the main challenges to pursuing such tasks is the lack of linguistic resources for evaluation and development. While any task like the one we will introduce here tackles the problem of knowledge acquisition and tries to engineer solutions to the bottleneck of knowledge acquisition through automated methodologies and algorithms, the development and evaluation of such methods relies closely on the provided dataset for testing and training e.g. [10], [11]. In addition, understanding and evaluation of the outcome of an IE/OL task is subject to the understanding of domain experts and the sort of information they are looking for; generally speaking, these activities are more task-driven rather than fact-driven. In addition, research studies in these domains usually focus on evaluation of engaged activities such as NER or RE in isolation. There is no report on the impact of the quality of these activities in the overall quality of the task performance.

For the reasons mentioned above, we have developed a dataset that will ideally result in a benchmark to evaluate the proposed task in section 3. The dataset comprises of sentences in the domain of Human Language Technology from the ACL Anthology Reference Corpus (ACL ARC) [12]. The annotations are provided at two layers, lexical and termino-conceptual. At the lexical layer the representation of an identical technology may comprise of lexical variants e.g. Human Language Technology may be signaled by HLT, Human Language Technology, Natural Language Processing, and NLP. However, at the conceptual level all these lexical variations refer to the same concept i.e. *HLT*. We have adopted the same approach for representing Semantic Relations; at the lexical level a semantic relation is a predicate i.e. defined based on the sentence surface structure. However at the termino-

conceptual level, semantic relations are classified into conceptual relations, either taxonomic or non-taxonomic e.g. lexical relations such as *used_in*, *applied_in*, and *employed_by* are classified under a conceptual relation *DEPEND_ON*. This layered representation will assist us in modularizing the task of Technology Structure Mining into several sub-tasks, including detecting technologies at the lexical level, mapping the technology lexicalizations to concepts, relation extraction between pairs of technology concepts at the lexical level, and finally mapping the lexical relations to conceptual semantic relations.

The rest of the paper is organized as follows. Section II briefly introduces related work. Section III proposes a formal task definition. Section IV describes the methodology for generating the dataset out of the ACL ARC corpus. Finally, section V concludes and gives the direction of future work.

II. RELATED WORK

Besides existing research in information extraction from patents e.g. [13], there is not much research reported towards extracting information from scientific publications for mining technology interdependence. Considering technology as applied science then it is not far from reality to consider scientific publications as a primary source of information for the task of technology structure mining. Research in this area could result in methodologies for smoothing the process of domain-semantic modeling in terms of technologies that are involved in a scientific domain. This may result in a strategic tool for intelligent information retrieval as well as for assisting the process of technology management.

As stated in [6], the information science research community and the Natural Language Processing (NLP) community [14] have focused on concepts and terms, but “the focus is increasingly shifting to the identification, process and management of relations to achieve greater effectiveness”. However, none of the literature in these domains explicitly mentions the correlation between concepts and relations, particularly in their task formalization. They either have considered this as an obvious fact, or this has not been the focus of their theoretical foundation. What is required here is a model that can combine and express properties of semantic relations from both the lexical and logical perspectives at a scalable size. Our research strives towards this goal. The most prominent research in recent years has approached the problem from the ontology engineering and population point of view. The main power of this research resides in the use of ontologies as a foundation for expressing domain-semantics. However, until recently [15] this research lacked concern about the lexical properties of concepts.

In [16], Hobbs and Riloff provide an overview of research in the Information Extraction (IE) domain. With emphasis on diversity in IE tasks, they have identified *named entity recognition*, *relation extraction*, and the task of *event identification* under the IE research topic and provide a

classification over the existing approaches from various perspectives and a comparison between finite state based methods versus machine learning approaches. They have discussed the complexity of the tasks of detecting complex words, basic phrases, complex phrases, as well as event detection and assigning them a unique identifier and a semantic type. The importance of real-world knowledge and its encoding into such systems is also emphasized.

In [9], Cimiano et al. give a survey of current methods in ontology construction and discuss the relation between ontologies and lexica as well as ontology and natural language. They illustrate different engineering approaches to ontology design and enumerate their advantages and disadvantages. On the topic of ontology learning, the authors contemplate controversies in concept identification and relation extraction. They emphasize the difference between linguistic representation of concepts and the concepts themselves and make a distinction between concept hierarchy and relation extraction since they see these as the difference between paradigmatic versus syntagmatic relations. The importance of selectional restriction and choosing the right level of abstraction are mentioned as other challenges in this field.

Khoo and Na [6] provide a survey on semantic relations. Their survey describes the nature of semantic relations from the perspective of linguistics and psychology, in addition to a detailed discussion of types of semantic relations including lexical-semantic relations, case relations, and relations between larger text segments. They clarify the definition of semantic relations in knowledge structures such as thesauri and ontologies. Although some semantic relations can be extracted/inferred from syntactic structures, there are other semantic relations that require a multi-step sequence of reasoning. Their survey enumerates a number of approaches for automatic/semi-automatic extraction of relations and describes the application of semantic relations in applications such as question-answering, query-expansion, and text summarization.

Finally, we consider much of the work in BioNLP as the closest to the proposed task here. Bio texts are usually written to describe a specific phenomenon e.g. gene expression, protein pathways etc. in a very specific context. Extracting such information, e.g. extracting instances of specific relations or interactions between genes and proteins, from Bio-literature is similar to the task of technology structure mining. However, in contrast to the proposed application here, Bio-Text Mining is well supported by ontologies and language resources; the context and concepts are usually clearly defined and tools which are tuned for the domain are available. The availability of knowledge resources such as well-defined ontologies in this domain lets Bio-Text miners build new semantic layers on top of already existing semantic resources (ontologies).

III. TASK DEFINITION

We define the task of technology structure extraction as comprising four major processes: 1) Identification of technology terms at the lexical level; 2) Mapping the lexical representation of technologies into a termino-conceptual level; 3) Extracting relations between pairs of termino-conceptual technologies at the lexical level (i.e. at sentence surface structure); 4) Finally, mapping/grouping relations at the lexical level into canonical relation classes at the conceptual level. We name the result of the proposed processes the *Technology Structure Graph* (TSG). Therefore, we define the task of technology structure extraction as the process of mapping a scientific corpus into a *TSG* graph with the following definition:

Definition 1: A *Technology Structure Graph* (TSG) is a tuple

$G = \langle V, P, S, \Sigma, \alpha, \beta, \omega \rangle$ where:

- 1) V is a set of pairs $\langle W, T \rangle$ where $\langle W, T \rangle$ is a uniquely identifiable terminology from a set of identifiers N and T is the terminology semantic type, e.g., $\langle \text{NLP}, \text{TECHNOLOGY} \rangle$ or $\langle \text{Lexicon}, \text{RESOURCE} \rangle$ or $\langle \text{Quality}, \text{PROPERTY} \rangle$. To support different levels of granularity of information abstraction we also consider V can contain pairs $\langle G_i, \text{GRAPH} \rangle$ where G_i has the same definition as G above.
- 2) P is a set of technology terms at the lexical level, uniquely identifiable from a set of identifiers R , e.g., Natural Language Processing, NLP, Human Language Technology.
- 3) S is a set of lexical relations, uniquely identifiable from a set of identifiers Q , e.g., used by, applied for, is example of.
- 4) Σ is a set of relations, i.e., the canonical relations vocabulary, e.g., $\{\text{DEPEND_ON}, \text{KIND_OF}, \text{HAS_A}\}$.
- 5) α is a partial function that maps $\langle W, T \rangle$ to a label of Σ annotated by a symbol from a fixed set M , i.e., $\alpha : N \times N \rightarrow \Sigma \times M$. M can be, e.g., the symbols $\{\square, \diamond\}$ from modal logic.
- 6) β is a function that maps P to a tuple in V i.e., $\beta : R \rightarrow N$.
- 7) ω is a function that maps S to a term in Σ i.e., $\omega : S \rightarrow \Sigma$.

Consider the following example input sentence:

“There have been a few attempts to integrate a speech recognition device with a natural language understanding system.” [17]

With M defined as *possible* and *certain* modalities, i.e., $\{\square, \diamond\}$, then the expected output of analysis of this sentence will be as follows:

$V = \{\langle \text{NLU}, \text{TECHNOLOGY} \rangle, \langle \text{SR}, \text{TECHNOLOGY} \rangle\}$
 $P = \{\text{natural language understanding, speech recognition}\}$

$\Sigma = \{\text{MERGE}\}$
 $S = \{\text{integrate with}\}$
 $\beta = \text{natural language understanding} \mapsto \langle \text{NLU}, \text{TECHNOLOGY} \rangle, \text{speech recognition} \mapsto \langle \text{SR}, \text{TECHNOLOGY} \rangle$
 $\omega = \text{integrate with} \mapsto \text{MERGE}$
 $\alpha = \langle \langle \text{SR}, \text{Technology} \rangle, \langle \text{NLU}, \text{Technology} \rangle \rangle \mapsto \langle \text{MERGE}, \diamond \rangle$

The main goal of the introduced task is in giving unstructured data (i.e. natural language text) a machine tractable structure in a way that we can semantically interpret this input data. Any semantic interpretation by machines is limited to our definition of symbols and their interpretations. In fact, since our knowledge of (natural language) understanding is limited, we move towards human understanding of language through an engineering approach. The proposed definition above can provide us with a base-line to perform and evaluate this task.

As with previous research in this domain, our task definition deals with two major sub-tasks: concept and relation identification/definition. It considers concepts as the building blocks of knowledge and relations as the elements that are connecting these concepts into a structure. However, we emphasize the interaction between concept definition and relation definition. In addition, we make the boundaries in the process more visible so we can divide the task into sub-tasks in a more modular manner enabling us to study their interconnections in a more systematic way. We argue it is not possible to define what we call relations vocabulary Σ without considering the definition of V .

The task of semantic interpretation of a natural language text involves an eco-system that comprises concepts, relations and linking/connecting concepts to each other through these relations, in addition to the user’s understanding of the provided symbols in V , and Σ . The other research challenge resides in mapping lexically introduced “concepts and relations” to a canonical termino-conceptual format. As stated in the given definition, we only focus on binary relations; the proposed model only concentrates on the relation between two technologies and we are aware of the limitations of the proposed model e.g. in modeling and representing the following example sentence:

“This method eliminates possible errors at the interface between speech Recognition and machine translation(component technologies of an AUTOMATIC Telephone Interpretation system) and selects the most appropriate candidate from a lattice of typical phrases output by the speech Recognition system.” [18]

In the above sentence, the author(s) addresses the interaction between two technologies and provides information about an interdependence. Our definition does not support representation of such information.

As mentioned, *Definition 1* provides us with a base-line to approach the task of Technology Structure Mining. Our first attempt towards this goal starts with developing a dataset for further experiments as described in the next section.

IV. DATASET DEVELOPMENT

As mentioned above the dataset comprises of sentences with at least two technology terms and their interdependencies. The sentences are extracted from the ACL Anthology Reference Corpus (ACL ARC) i.e. a corpus of scholarly publications about Computational Linguistics consisting of 10,921 articles which can be downloaded from [19]. The ACL ARC is represented in three different formats: source PDF files of articles, plain text, and an XML version of the articles i.e. the OCR output of PDF files with additional information of visual features of the text e.g. font face, font size, the position of text etc. The corpus is further divided into different sections in directories labeled with a single letter, with 11 sections in total.

The dataset development essentially comprised 4 steps (Figure 2): 1) Text Processing; 2) Indexing and Storage; 3) Concept (technology) Identification; and 4) Compilation of dataset. Then we studied the selected sentences manually, verified the processes, and annotated the sentences with the lexical/semantic relations between pairs of technologies. In the remainder of this section we give a description of each step of the task with results on the corpus.

We followed an iterative process for the dataset development. In the first step, the main issue is to find the optimum boundary size of text for dataset development e.g. should we focus at paragraph level or sentence level. To answer this question, in the first step we chose 1,424 random papers from the corpus and performed the following analysis. The selected papers consist of 45,031 paragraphs, 168,028 sentences, 4,524,062 tokens, and 124,525 types¹. We studied the distribution of terms that can be considered as a representation of a technology in the domain. Our experiment showed that the co-occurrences of pairs of technologies tend to happen at sentence level (Figure 1). This means that if two technologies occur within a text segment then it is more likely that this happens within a sentence. In addition, studying the relations at a greater boundary such as paragraph level imposes computational costs that may not be desirable considering the size of the corpus, the cost of annotating a dataset, and the current state of technologies such as anaphora resolution. This has been also discussed from another perspective in [20]. In the remainder of this section we describe each step of the analysis in detail.

¹The numbers proposed here are subject to the errors that are imposed by text processing/extraction process and may not be identical using different approaches for text extraction

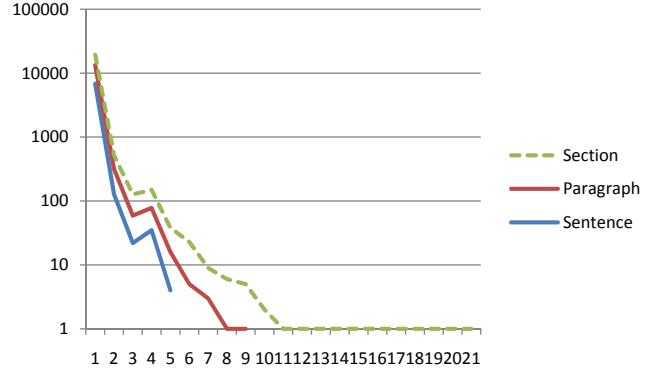


Figure 1. Distribution of co-occurrences of technology terms: The analysis shows that the co-occurrences of two technology terms tend to be at the boundary of sentences; The above diagram shows that if two technologies appeared together in a text boundary then it is most probable that these two terms are situated within a sentence. Here, the vertical axis shows the number of technology terms and the horizontal axis shows the number of terms (in logarithmic scale) in sentence, paragraph and sections segments e.g. the diagram shows that we have 10,000 sections, paragraphs, and sentences with one technology term while there are no paragraphs or sentences with more than 10 technology terms within their boundaries.

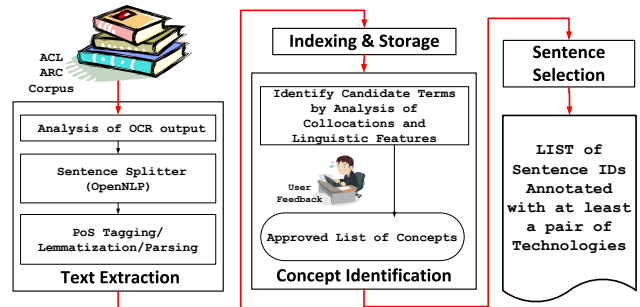


Figure 2. Dataset Development: Steps that have been taken for selecting sentences

A. Text Processing

The ACL ARC corpus does not provide text sections and segments. The first stage of our process therefore involved text sectioning, and structuring. The text sectioning step involved converting provided XML files in ACL ARC into a more structured XML document where different sections of a paper such as titles, abstract, references etc. were identified using a set of heuristics. The heuristic rules are based on provided visual information in the source XML files such as font face, font size, position of text segments, and their frequency distribution. As for any other text sectioning task, this step involves noise and error in the output. In the next step, we performed text segmentation including the detection of boundaries of paragraphs, sentences, and tokens. We have also performed part-of-speech tagging and lemmatization. For detecting paragraph boundaries we have used a set of heuristics. However sentence segmentation and tokenization has been carried out with OpenNLP [21]. Since OpenNLP

tools are trained on scientific publications, they tended to perform better when compared to other available tools. Then, We used the Stanford Part of Speech (POS) tagger [22] for tagging and lemmatization. The generated files are available for download¹. The indexed sentences were also processed with open source dependency parsers: Malt Parser [23], BioLG [24], and Stanford Dependency Parser [25].

B. Indexing and Storage

The next step of the process involved indexing and storage of the corpus. We have used a data model -available at the URL in the footnote- that lets us dynamically generate a lexicon out of the POS tagged and lemmatized tokens in the corpus, along with the frequency of words. This also enables us to keep track of the position of words, sentences, paragraphs, and sections within a document. For example, we can easily identify all the sentences, paragraphs, and sections that have the word *technology* with a specific linguistic annotation such as part of speech. We have used the model to retrieve data from the corpus with queries similar to the Corpus Query Language [26] but at uniquely indexed text segments. Improved performance, reduced processing time, ability for concurrent parsing of sentences, as well as flexibility in modification of metadata have been among the other reasons for using the proposed model in Figure ??.

C. Concept Identification

The concept identification (technology term recognition) process starts with selecting all the phrases in the corpus with the word “technology/ies”. In fact we queried the corpus for the chain of tokens/lexemes that end with a token that has “technology” as its lemma. In addition, we applied a set of filters which have been defined based on part of speech and the position of the tokens. For example, if we found a lexeme chain starting with a *verb in gerund or present participle form* (i.e. VBG part of speech in Penn Style Treebank [27]) then the chain would be accepted only if a determiner appeared before the token with VBG part of speech. In the next step, the extracted technology terms were manually refined. Among the 147 extracted lexeme chains, 31 terms were rejected manually (this includes meaningless terms in addition to very specific terms such as “Japaneses sentence parsing technology”). Then, we manually grouped the remaining terms into 43 different classes, each class refers to a specific technology in the domain of Human Language Technology e.g. finite-state, segmentation, parsing, entity-extraction, etc. As a matter of fact, this processing step comprises the evaluation of P , V , and the function β in Definition 1 in section III. As an example, at the end of this step, P includes these strings: *information retrieval technology, information retrieval technologies, information retrieval, IR technology, IR*, while V has

a member $\langle IR, TECHNOLOGY \rangle$ and function β maps all the given values above for P to $\langle IR, TECHNOLOGY \rangle$ in V . This processing step has been carried out on the sub-corpus of 1,424 random papers described above.

D. Sentence Selection

After choosing the technology classes and defining P , V and β for the corpus, we identified sentences that contain more than one string term from P . In this step, we extracted the sentences for each section of the ACL ARC; e.g. we were able to extract text from 2,435 papers out of section C (failing on 432 papers; either because of errors in the source XML files or deficiency in our heuristics for corpus processing). This step has been carried out on all sections of the corpus. Table I and Table II show summarized statistics of the performed processes. Table I shows the overall number of articles that have been extracted from the XML source files (*ARTICLES#*), the number of documents successfully segmented and indexed (*SUC-ARTICLE#*), and the number of documents that failed to segment and index (*UNSUC-ARTICLE#*). Table II² shows statistics for the successfully indexed documents. This includes the numbers of tokens, types, identical sentences (*SENT*), identical sentences with a minimum of 1 technology term (*SST1*) and identical sentences with more than one technology term (*SST2*) for each section of the corpus.

Table I
STATISTICS FOR TEXT PROCESSING STEP

| Section | ARTICLES# | SUC-ARTICLE# | UNSUC-ARTICLE# |
|--------------|---------------|--------------|----------------|
| A | 404 | 265 | 139 |
| C | 2,435 | 2,003 | 432 |
| E | 846 | 463 | 383 |
| H | 897 | 828 | 69 |
| I | 146 | 113 | 33 |
| J | 922 | 114 | 808 |
| M | 180 | 168 | 12 |
| N | 371 | 365 | 6 |
| P | 2028 | 1873 | 155 |
| T | 120 | 81 | 39 |
| W | 2281 | 2121 | 160 |
| Total | 10,630 | 8,394 | 2,236 |

E. Manual Verification of Analysis, Annotation and Grouping of Relations

In the final step of dataset development, we chose and annotated sentences from section C of the corpus. This section of the corpus comprises papers from different conferences from the years 1965 to 2004. Among the 230,936 sentences in this section of the corpus, only 2,012 sentences contain a technology term, and amongst these sentences only

²The total numbers of articles proposed here are not identical to the numbers proposed in [12] due to corruptions in the source XML files; we have excluded these files from the corpus

¹http://nlp.deri.ie/behrang/sepid_arc.html

Table II
STATISTICS FOR EXTRACTED TEXT FROM ACL-ARC SECTIONS

| Section | Token# | Type# | SENT# | SST1# | SST2# |
|---------|---------|--------|--------|-------|-------|
| A | 955761 | 40938 | 35439 | 2012 | 134 |
| C | 6168312 | 172077 | 230936 | 7514 | 482 |
| E | 1901481 | 61854 | 67588 | 1646 | 81 |
| H | 2107057 | 56470 | 78797 | 4777 | 330 |
| I | 358358 | 20299 | 14258 | 721 | 52 |
| J | 612692 | 23702 | 22061 | 496 | 25 |
| M | 400398 | 20807 | 14903 | 592 | 52 |
| N | 1164215 | 38772 | 44103 | 2349 | 180 |
| P | 7446189 | 152890 | 272706 | 8833 | 603 |
| T | 122969 | 10882 | 4693 | 65 | 1 |
| W | 8169591 | 167107 | 300612 | na | na |

482 have two or more lexical chains that signal appearance of technologies of different classes in the sentence. We manually read the extracted sentences and annotated them with the following information:

- 1) Whether the text processing step has been performed correctly: this comprised checking the sectioning/segmentation of the source XML files, sentence splitting and tokenization.
- 2) Technology Mark-up: whether the applied approach for detecting the technologies has been successful.
- 3) Type of Relation: whether the sentence implies/expresses a relation between marked-up technologies. Moreover, this gives the linguistic context for the relation as described below.
- 4) Lexical Relation: whether a sentence implies a relation and how that is expressed.
- 5) Grouping of Lexical Relations into Semantic Relations: classification of detected lexical relations into semantic relations.

As mentioned earlier, we have identified and classified 5 different types of contexts for relation extraction as follows:

- 1) *Noun-Compound*: This context refers to a relation that can be inferred from the combination of nouns in a compound, e.g.:

“Since a model of machine translation called translation by Analogy was first proposed in Nagao(1984), much work has been undertaken in *Example-Based NLP* (e.g. Sato and Nagao (1990) and Kurohashi and Nagao (1993)).” [28]

The above sentence suggests a relation as follows:
 <<(NLP, technology),has – sub – class,
 (EB-NLP, technology)>>
 Noun-Compound is the only context that provides termino-conceptual relations directly.
- 2) *Prepositional*: This class of relations can be inferred from prepositional attachment, e.g.:

“NLP components of a machine translation system are used to automatically generate semantic representations of text corpus that can be given

directly to an ILP system.” [29]

the above sentence suggests a relation as follows:

<<(MT, technology),has – component,(NLP, technology)>>

- 3) *Verb-based*: This refers to contexts where two technology terms are directly/indirectly related to each other by a verb, e.g.:

“lexical Knowledge acquisition *plays an important role* in Corpus-Based NLP.” [30]

However, extracting relations of this type may not be as straight-forward because other relations e.g. noun-compounds may occur at the same time. For example, relations in the above sentence are as follows:

<<(lexical-KA, technology),is – sub – class,
 (KA, technology)>>
 <<(CB-NLP, technology),is – sub – class,
 (NLP, technology)>>
 <<(lexical-KA, technology),plays – role – in,
 (CB-NLP, technology)>>

- 4) *Structural*: this context refers to relations that can be inferred based on the structure of a sentence, e.g.:

“Transformation-Based learning has been used to tackle *a wide range of* NLP problems, *ranging from* part-of speech tagging (Brill, 1995) to parsing (Brill, 1996) *to* segmentation and message understanding (Day et al., 1997).” [31]

This suggests the relation:
 <<(POS-tagging, technology),
 is – problem – example – of,(NLP, technology)>>

- 5) *Residuals*: this category refers to relations that do not fit into any of the first three above categories and/or are too complicated to be automatically inferred via structure, e.g.:

“finite-state rules are represented Using regular expressions and they are transformed into finite-state automata by a rule compiler.” [32]

This conveys a relation between *Finite Automata* and *Compiler*. Consider another example sentence:

“In translation memory or Example-Based machine translation systems, one of the decisive tasks is to retrieve from the database ,the example that best approaches the input sentence.” [33]

This expresses a relation between *Database Technology* and *Machine Translation Technology*. However, we believe that the expressed relations in these sentences are too complex: automatic extraction and expression of such relations by TSG may be far from reality. It is worthwhile to mention that we have identified some of the relations expressed by sentence

structure that are difficult to extract automatically. For example, the temporal relation between the time of introducing “translation by Analogy” and “Example-Based NLP” expressed in the above sentence, and the temporal relation conveyed by the sentence given previously as an example of a noun-compound relation. We have grouped these relations under the residuals category.

These different contexts have been studied in previous research e.g. [7], [34]–[36] and [37]. However, the authors are unaware of any reported research on the analysis of the distribution of these contexts, nor any corpus that provides linguistic context annotations for relation extraction.

Among the 482 annotated sentences, the text extraction process has been carried out correctly for 425 sentences, and it fails for 57 cases. This gives the precision of 89% for this process step. Unfortunately, our approach does not allow the measurement of the recall for text extraction at the sentence level. However, Table I may be used for measuring recall at the document level. The process of concept identification (technology recognition) has been done correctly for 385 sentences: this gives precision of 81% at the sentence level. However, among the total number of 982 instances of technologies, 78 cases were marked up incorrectly; this will give the precision of 92% for technology recognition ignoring the text segmentation error.³

Among the 482 sentences, 201 sentences are annotated with at least one relation context (summarized in table III): 37 *Noun-Compounds*, 26 *Prepositional*, 59 *Verb-based*, and 79 *Structural* relations. 55 sentences are annotated with relations of the type of *Residual*. Other sentences are not accompanied by a relation since they do not express any relation between the marked-up technologies, e.g.:

“the result could be helpful to solve the variant problems of information retrieval , information extraction , question answering , and so on.” [38]

Table III
FREQUENCY OF RELATION CONTEXTS IN THE DATASET OF 482 SENTENCES

| Context | Frequency |
|---------------|-----------|
| Noun-compound | 37 |
| Verb-based | 26 |
| Prepositional | 59 |
| Structural | 79 |
| Residual | 55 |

We finally mapped the lexical relations into the termino-conceptual relations manually (Defining $\omega : S \rightarrow \Sigma$ in Definition 1 in section III). For example, the lexical relations, S , such as `incorporate`, `is_combined_with`,

³We have defined precision as the number of correct annotations divided by the total number of annotations

and `integrate_with` are mapped into the termino-conceptual relation MERGE in Σ .

V. CONCLUSION AND FUTURE WORK

We introduce the task of *Technology Structure Mining* as an example of a broader task of extracting concepts and the relationships between them for a given text corpus. We propose a “Technology Structure Graph” for formalizing the task. The major challenge is the lack of a benchmark dataset for evaluation and development purposes. The paper reports steps taken for constructing such a dataset which comprises 482 sentences from section C of the ACL ARC corpus. Each sentence is annotated with at least two technology terms and their interdependencies. We have also annotated the sentences with a linguistic context category that relations may be inferred from. Moreover, sentences are accompanied by other miscellaneous annotations such as the modality of the relations, and the position of the sentence in the article.

Future work will include the manual correction/annotation of part of speech tags and dependency parses for the selected sentences. This will enable us to study the performance of generic parsers on our dataset. Since the proposed task consists of several steps including text sectioning and segmentation, part of speech tagging etc. and as each of these processes is subject to error, there may be the danger of accumulated errors. The annotated dataset will enable us to study this in details.

ACKNOWLEDGMENT

The authors would like to thank Dr. Antoine Zimmermann and Nuno Lopes for fruitful discussions. This research is funded by Science Foundation Ireland under grant number SFI/08/CE/I1380(Líon-2).

REFERENCES

- [1] A. M. Badawy, “Technology management simply defined: A tweet plus two characters,” *J. Eng. Technol. Manag.*, vol. 26, pp. 219–224, 2009.
- [2] L. W. Fry, “Technology-structure research: three critical issues,” *Academy of Management Journal*, vol. 25, pp. 532–52, 1982.
- [3] D. Bailey, P. M. Leonardi, and J. Chong, “Minding the gaps: Understanding technology interdependence and coordination in knowledge work,” *Forthcoming Organization Science*, 2009.
- [4] D. Nadeau and S. Sekine, “A survey of named entity recognition and classification,” *Linguisticae Investigationes*, vol. 30, pp. 3–26, 2007.
- [5] D. Gildea and D. Jurafsky, “Automatic labeling of semantic roles,” *Computational Linguistics*, vol. 28, no. 3, pp. 245–288, 2002.
- [6] C. S. G. Khoo and J.-C. Na, “Semantic relations in information science,” *Annual Review of Information Science and Technology*, vol. 40, no. 1, pp. 157–228, 2006.

- [7] D. Zelenko, C. Aone, and A. Richardella, "Kernel methods for relation extraction," *J. Mach. Learn. Res.*, vol. 3, pp. 1083–1106, 2003.
- [8] M. Banko, M. J. Cafarella, S. Soderland, M. Broadhead, and O. Etzioni, "Open information extraction from the web," *IJCAI*, pp. 2670–2676, 2007.
- [9] P. Cimiano, P. Buitelaar, and J. Völker, "Ontology construction," in *Handbook of Natural Language Processing, Second Edition*, N. Indurkha and F. J. Damerau, Eds., 2010, pp. 577–605.
- [10] R. Hwa, "Learning probabilistic lexicalized grammars for natural language processing," Ph.D. dissertation, Harvard University, Cambridge, MA, USA, 2001, adviser-Shieber, Stuart.
- [11] C. Zhang, "Extracting chinese-english bilingual core terminology from parallel classified corpora in special domain," in *WI-IAT '09*. Washington, DC, USA: IEEE Computer Society, 2009, pp. 271–274.
- [12] S. Bird, R. Dale, B. Dorr, B. Gibson, M. Joseph, M.-Y. Kan, D. Lee, B. Powley, D. Radev, and Y. F. Tan, "The acl anthology reference corpus: A reference dataset for bibliographic research in computational linguistics," in *LREC'08*, Marrakech, Morocco, May 2008.
- [13] Y.-H. Tseng, C.-J. Lin, and Y.-I. Lin, "Text mining techniques for patent analysis," *Information Processing & Management*, vol. 43, no. 5, pp. 1216 – 1247, 2007, patent Processing.
- [14] I. Hendrickx, S. N. Kim, Z. Kozareva, P. Nakov, D. O. Séaghdha, S. Padó, M. Pennacchiotti, L. Romano, and S. Szpakowicz, "Semeval-2010 task 8." *ACL*, 2009, pp. 94–99.
- [15] P. Buitelaar, P. Cimiano, P. Haase, and M. Sintek, "Towards linguistically grounded ontologies," in *ESWC*, June 2009, pp. 111–125.
- [16] J. R. Hobbs and E. Riloff, "Information extraction," in *Handbook of Natural Language Processing, Second Edition*, N. Indurkha and F. J. Damerau, Eds. CRC Press, 2010.
- [17] M. Tomita, M. Kee, H. Saito, T. Mitamura, and H. Tomabechi, "The universal parser compiler and its application to a speech translation system," in *Proc. of the 2nd Inter. Conf. on Theoretical and Methodological Issues in Machine Translation of Natural Languages*, 1988, pp. 94–114.
- [18] K. Kakigahara and T. Aizawa, "Completion of japanese sentences by inferring function words from content words," in *COLING*, 1988, pp. 291–296.
- [19] "Acl anthology reference corpus (acl arc)," <http://acl-arc.comp.nus.edu.sg/>.
- [20] T. M. Mitchell, J. Betteridge, A. Carlson, E. Hruschka, and R. Wang, "Populating the semantic web by macro-reading internet text," in *ISWC*, 2009, pp. 998–1002.
- [21] "The opennlp project," <http://opennlp.sourceforge.net/>.
- [22] "Stanford log-linear part-of-speech tagger," <http://nlp.stanford.edu/software/tagger.shtml/>.
- [23] J. Nivre, J. Hall, S. Kbler, and E. Marsi, "Maltparser: A language-independent system for data-driven dependency parsing," in *In Proc. of the Fourth Workshop on Treebanks and Linguistic Theories*, 2005, pp. 13–95.
- [24] S. Pyysalo, T. Salakoski, S. Aubin, and A. Nazarenko, "Lexical adaptation of link grammar to the biomedical sublanguage: a comparative evaluation of three approaches," *CoRR*, vol. abs/cs/0606119, 2006.
- [25] M.-C. de Marneffe, B. MacCartney, and C. D. Manning, "Generating Typed Dependency Parses from Phrase Structure Parses," in *IEEE / ACL 2006 Workshop on Spoken Language Technology*, 2006.
- [26] "Using corpus query language for complex searches," <http://www.fi.muni.cz/~thomas/corpora/CQL/>.
- [27] M. P. Marcus, B. Santorini, and M. A. Marcinkiewicz, "Building a large annotated corpus of english: The penn treebank," *Computational Linguistics*, vol. 19, 1994.
- [28] T. Utsuro, K. Uchimoto, M. Matsumoto, and M. Nagao, "Thesaurus-based efficient example retrieval by generating retrieval queries from similarities," 1994.
- [29] Y. Sasaki and Y. Matsuo, "Learning semantic-level information extraction rules by type-oriented ilp," in *COLING*, 2000, pp. 698–704.
- [30] A. Sarkar and W. Tripasai, "Learning verb argument structure from minimally annotated corpora," in *COLING*, 2002, pp. 1–7.
- [31] D. Wu, G. Ngai, and M. Carpuat, "Why nitpicking works: evidence for occam's razor in error correctors," in *COLING*, 2004, p. 404.
- [32] K. Koskenniemi, P. Tapanainen, and A. Voutilainen, "Compiling and using finite-state syntactic rules," in *COLING*, 1992, pp. 156–162.
- [33] E. P. Cyber and E. Planas, "Multi-level similar segment matching algorithm for translation memories and example-based machine translation," in *COLING*, 2000.
- [34] M. A. Hearst, "Automatic acquisition of hyponyms from large text corpora," in *COLING*, 1992, pp. 539–545.
- [35] D. I. Moldovan and R. Girju, "An interactive tool for the rapid development of knowledge bases," *International Journal on Artificial Intelligence Tools*, vol. 10, no. 1-2, pp. 65–86, 2001.
- [36] P. Sazedj and H. S. Pinto, "Mining the web through verbs: A case study," in *ESWC*, ser. Lecture Notes in Computer Science, E. Franconi, M. Kifer, and W. May, Eds., vol. 4519. Springer, 2007, pp. 488–502.
- [37] V. Srikumar, R. Reichart, M. Sammons, A. Rappoport, and D. Roth, "Extraction of entailed semantic relations through syntax-based comma resolution," in *ACL*, 2008, pp. 1030–1038.
- [38] T. Masuyama, S. Sekine, and H. Nakagawa, "Automatic construction of japanese katakana variant list from large corpus," in *COLING*, 2004, pp. 1214–1219.