

Knowledge graph construction for research literatures

Alisson Oldoni

A research project report submitted for
the degree of
Master of Computing and Information Technology



School of Computer Science and Engineering
The University of New South Wales

1 November 2016

Abstract

Contents

Contents	iv
1 Introduction	1
2 Information Extraction	5
2.1 Natural Language Processing	5
2.2 Information Extraction	12
2.3 Knowledge Graphs	18
3 Analysis and Related Work	21
3.1 Analysis of Academic Text	21
3.2 Peculiarities with Academic Text	24
3.3 Related Work	25
4 Developed Workflow	27
4.1 Tools	27
4.1.1 PDF extraction generates noisy output	27
4.1.2 Programming Languages	27
4.1.3 Python Modules	27
4.1.4 NLTK	28
4.1.5 Brat	28
4.1.6 Jupyter	28
4.1.7 Parsey	28
4.1.8 SpaCy	28
4.1.9 Stanford CoreNLP	28
4.1.10 Graphviz	28
4.2 Developed process and programs	28
5 Results	29
6 Conclusion and Future Work	31
Bibliography	33

Chapter 1

Introduction

Information extraction (IE) is the process of obtaining in an automatic fashion facts and information from unstructured text that can be read by a machine [18].

Historically, it mostly started with exercises on template filling based on raw natural text [25] as part of the Message Understanding Conferences (MUC) from the late 1980s and 1990s. As part of the MUC, competitions would take place in which a corpus would be made available of a specific domain, and different teams with different programs would try to extract the information from the natural text as to fill in the intended templates.

Note the following text from a news report regarding the result of a soccer match:

‘Though Brazilian star Diego Tardelli’s equaliser denied the Sky Blues victory at Jinan Olympic Sports Centre Stadium on Wednesday night, David Carney banked a precious away goal that will bode well for Graham Arnold’s side when they host Shandong in next week’s second round-of-16 leg. Sydney FC have taken a sizeable step towards a maiden Asian Champions League quarter-final berth after securing a 1-1 draw with Shandong Luneng in China.’

Team 1: _____
Team 2: _____
Winner: _____
Location: _____
Final Score: _____

Figure 1.1: An example of template to be filled in the sports domain.

An example of a task would be, to solely based on the above raw text, to fill in the template shown in Figure 1.1. The MUC competition would be

based in various corpus and tasks based on varieties of news reports, such as satellite launches, plane crashes, joint ventures and other different data in these specific domains.

On the above example, one can observe that the *Team 1* is *Sydney FC*, *Team 2* is *Shandong Luneng*, there was no *Winner*, and consequently the *Location* and the *Final Score*. It gets more interesting as you observe the same type of information being delivered by a different reported:

‘SYDNEY FC take the advantage of an away goal in China, leaving the second leg of their Asian Champions League Round of 16 tie with a 1-1 draw with Shandong Luneng.’

Although roughly similar in this case, the approach to retrieve the data from Natural Language text needs to be able to generalise to the various ways a reporter might write such information. This effort becomes more complex as one moves through different domains and audiences of a text, such as: technical manuals, academic papers from different areas, legal text, contracts, financial news, biomedical, among others.

More recently, the output of such Information Extraction systems are used as to build other systems, more prominently Knowledge Graphs. A Knowledge Graph (KG), also known as the knowledge base, is a collection of the machine-readable database that contains entities, the attributes of entities and the relationships between entities [14]. Information Extraction tools would harvest data from unstructured or semi-structured text and provide such databases.

Popular search engines such as Google [14] and Bing [4] leverage Knowledge Graphs as to provide entity summary information and the related entities based on the query that the user is searching for. It is an essential foundation for many applications that requires machine understanding.

The use of Knowledge Graphs then allow users to be able to see extra information in a summarised table-like form, as to resolve their query without having to navigate to other sites. Note in the example in Figure 1.2 how the right column represents a sequence of facts of the ‘*Jimi Hendrix*’ entity, in this case an entity of the class (or type) *PERSON*, such as: his official website; where and when he was born; where and when he died; and a list of movies where this person is the subject of.

Modern pipelines for building Knowledge Graphs from raw text would then encompass several Information Extraction techniques, such as the ones below:

1. Discover entities in the text;
2. Discover relationships between these entities;
3. Perform entity disambiguation

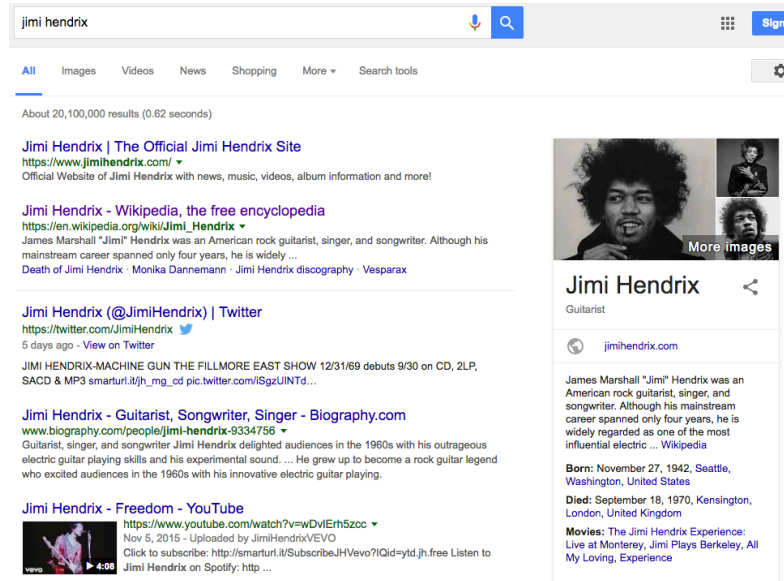


Figure 1.2: An example of knowledge graph application in the Google’s result page.

4. Link entities to a reference Knowledge Graph (e.g., Yago2 [32] or DB-Pedia [19]).
5. Improve the quality of the output via input data cleaning, robust extraction, and learning-based post-processing methods;
6. Reason about how accurate these facts are;
7. Finally presenting the facts in a graph (the Knowledge Graph).

Some of these techniques will be explained further as part of this document.

In this project, we focus on studying and presenting some Information Extraction techniques as to build a domain-specific verb-centric information extraction tool that extracts relations from academic papers. More specifically, we focus on papers from the topic of databases and attempt to extract information from these papers for posterior usage by other systems with applications such as:

- Allow for structured and fast search of techniques in the papers and the possible relations between them;
- Possibly group papers by their used techniques;
- Discovery of techniques to improve performance on a certain problem;

- Generate a hierarchy of concepts, and their use;
- Among others.

An existing service that organises data from academic papers is the Semantic Scholar [29] project, by Professor Oren Etzioni from Allen Institute for AI. However, Semantic Scholar only understands a limited number of relationships (such as 'cite', 'comment', 'use_data_set', and 'has_caption') which are also more closely related to the meta-data about the paper, but not from the knowledge that the paper itself presents. Other similar services are Microsoft Academic Graph [23], Google Scholar [15], and CiteSeerX [8].

In the next sections, this document will give some background information on the techniques needed to achieve the above (Chapter 2), and it will also define the problem more precisely (Chapter 3), building as to introduce the development of this research (Chapter 4). In Chapter 5 we will describe some of the results, followed by the some final remarks in Chapter 6.

Chapter 2

Information Extraction

Information Extraction, a term already defined in the introduction, is a hard problem which mostly relies in attempting to use a computer to understand information explicitly stated in the form of natural language.

It is interesting to observe that, before writing down thoughts in a paper, an academic form the ideas of what facts he/she wants to express in his/hers head, and then attempt a structure to most clearly state these in text form. These multiple facts, and the relations between them, are then stated in a sentence in what is assumed to be a somewhat logical format, following the semantics of the language, whether English or any other. Following this example, one must then also assume that the future reader of this paper will use the reverse process to decode this information into facts or ideas to be understood. In fact, this assumption is what justifies the attempt of Information Extraction.

Several initiatives in the Natural Language Processing area attempt to understand and map what these semantic rules are, and how one could use a computer for tackling natural language related tasks. These initiatives are fruitful and provide advanced tools and techniques in which some will be described in this chapter.

2.1 Natural Language Processing

Natural Language Processing (or *NLP*) can be a term used to discuss any kind of computer manipulation of natural language text, also called raw text. It can mean simple things such as counting words and obtain their frequency distribution to compare different writing styles, or in a more complex sense it would require the understanding of human writings, to the extent of being able to extract information and meaning from it, or also give useful responses to them [5]. On this more complex end of the spectrum, where one wants to understand raw text, language technology and existing tools rely on formal models, or representations, of knowledge of language at the levels of

morphology, syntax, semantics, among others linguistic concepts. A number of formal models including state machines, formal rule systems, logic, and probabilistic models are used to capture this knowledge from text and reason with it [18].

When information is laid out in natural language form, one start the analysis of the information presented by constructing a phrase or sentence based on smaller pieces of information such as verbs, nouns, and adjectives which are called the constituents. These then build up to form sequences of simple and complex sentences.

Observing more carefully a simple sentence, such as ‘*The police chased him.*’, it is possible to attempt to sample the different syntactic information presented in it. As a first step it is possible to dissect its constituent parts as per Figure 2.1.

The/DT police/NN chased/VBD him/PRP ./.

Figure 2.1: An example of tagged sentence.

The first word in the sentence, *The*, is a *DT* or determiner. Other possible determiners include ‘*my*’, ‘*your*’, ‘*his*’, ‘*her*’. The second word is ‘*police*’ is a *NN* which is the tag for a singular noun. With this information we can already tell that this sentence is speaking about something, and this something is the noun ‘*police*’. Subsequently the tag *VBD* is presented which specifically indicates the word ‘*chased*’ is a verb (an action) in the past tense. At this point one can observe that something or someone (in this case the ‘*police*’) did something in the past.

This is already great information to have about the sentence. These tags that were added to the text in Figure 2.1 are called Part-Of-Speech tags, or POS tags [18]. The standardization of these tags and work to develop and tag existing text with them is done by the Penn Treebank project [21].

Note how specific the tags are, dictating the type of the word, a verb for an example, and its variation either in quantity or tense. Some other examples of these tags are shown in Table 2.1. Although most of them are simple to understand, note that *ADP* are the adpositions, which encompasses prepositions and postpositions. Some systems, such as the spaCy Natural Language Processing parser [17, 30] also maps these more specific tags into more general ones, for an example, while three different words in a sentence are tagged independently as *VBD*, *VBG* and *VBZ*, they are also tagged with a *VERB* tag. This is useful, if the user is not too interested in the detail of which verb variation was used.

A Part-Of-Speech Tagger is then a system that, given a raw text as input, assigns parts of speech to each word (or token) and is able to produce as output the tagged version of this text. The text in Figure 2.1 was tagged

POS tag	Meaning	Sample
ADP	Adpositions	<i>at, or in</i>
CONJ	Coordinating conjunction	<i>and, or or</i>
DT	Determiner	The
JJ	Adjective	She is <i>tired</i>
JJR	Adjective, comparative	That one is <i>larger</i>
JJS	Adjective, superlative	That is the <i>largest</i>
NN	Noun, singular or mass	Car
NNS	Noun, plural	Cars
NNP	Proper noun, singular	Microsoft
NNPS	Proper noun, plural	The <i>Kennedys</i>
RB	Adverb	She said <i>firmly</i>
VB	Verb, base form	Attack
VBD	Verb, past tense	Attacked
VBG	Verb, gerund or present participle	Attacking
VCN	Verb, past participle	Broken
VBP	Verb, non-3rd person singular present	I <i>attack</i>
VBZ	Verb, 3rd person singular present	He <i>attacks</i>

Table 2.1: List of some of the possible Part-Of-Speech (POS) tags.

using the Stanford Log-linear Part-Of-Speech Tagger [35].

POS Tag	Word	Prev. Word	Prev. Tag
DT	The	<START>	<START>
NN	police	The	DT
VBD	chased	police	NN
PRP	him	chased	VBD
.	.	him	PRP

Table 2.2: Features for sequential POS tagging.

The task of assigning these tags starts by deciding what are the tokens in a raw text sentence, and what are its sentences. As an example, the tokenizer needs to decide if a period symbol near a word represents an abbreviation (e.g. *Dr.*) or a sentence boundary - in case of an abbreviation, this period is then considered simply a token within the sentence. Another common problem in this step is deciding if a single quote is part of a word, (e.g. *'It's'*), or is delimiting a quoted part of the sentence, thus potentially hinting other semantic meanings. The Stanford POS Tagger used in this example also contain a tokenizer, which is part of the Stanford CoreNLP [20], a set of natural language analysis tools.

Modern POS Taggers tackle this task using a technique called Sequence Classification. A machine learning classifier model is then trained with a cor-

pus of manually tagged text and has as an input certain features that might indicate which tag a token being currently analysed should be assigned with. Observing again the example from Figure 2.1, but now presented in table format in 2.2, it is easier to see how this learning algorithm would be trained to predict tags on unseen test. Note the second word ‘*police*’. For this word we are providing 4 features for the learning algorithm: the manually labelled POS Tag, the word itself, the previous word and the previous POS Tag. Suppose now that this sentence is included in a bigger corpus, and this pattern is a common one and the learning algorithm is provided with a substantial amount of labelled data in which this situation repeats itself: a *NN* is the second word in a sentence, with ‘*The*’ and *DT* being the previous word and tag respectively.

After the training, this model would behave in a similar fashion once presented with unseen data. Suppose now the first column on Table 2.2 is not presented. The model would pick the first word ‘*The*’ and observe the features: $\langle START \rangle$ and $\langle START \rangle$ respectively and given that in our previously described corpus this is a common occurrence, it would then label this word with *DT*. Now, for the second word ‘*police*’, the features would be ‘*The*’ for previous word, and *DT* for previous tag. Again, it is common for a noun to be placed after a determiner so the model assigns the label *NN* to the word ‘*police*’. The features used by the Sequence Classifier both while training and when using the model may vary and they impact the quality of the predictions it makes. The Stanford POS Tagger uses a broad use of lexical features, including jointly conditioning on multiple consecutive words [35].

A helpful concept at this stage is known as the lemma. A lemma is a canonical way of representing a word which strips out variations for quantity, tense, among others [18]. For an example, given the words ‘*running*’ or ‘*runs*’, NLTK [5] outputs *run* as their lemma. The lemma can, for an example, be used together or instead of the existing features for training models.

Jetstar/NNP Airways/NNPS ,/PUNC a/DET unit/NN of/ADP Qantas/NNP
Airways/NNP Limited/NNP

Figure 2.2: An example of another tagged sentence.

When reading a text, it is also important to understand the relation between the words or sub-sentences that are contained in the phrase, and this is specially useful for information extraction. As an example, suppose the sentence ‘*Jetstar Airways, a unit of Qantas Airways Limited*’. There are different ways of breaking down the relationship of the words in this sentence.

Starting again with the POS tagging discussed earlier, one can see in

2.2 what are the types of the words that constitutes this sentence. As a further step, it's possible to see what are the sub-sentences or parts that form the phrase structure, also called constituency parsing. The sub-sentences are connected upwards to one head (also called, *parent*), and downwards to one or more governors (also called *dependants*, or *children*), in a recursive structure. In Figure 2.3, the phrase '*of Qantas Airways Limited*', which is part of the bigger phrase we are using as an example, is a prepositional phrase. A prepositional phrase lacks either a verb or a subject, and serve to assert binary relations between their heads and the constituents to which they are attached, in this case '*a unit*' [18]. The sub-sentence '*a unit of Qantas Airways Limited*' is then a noun phrase, since it contains and talks about a noun '*unit*'.

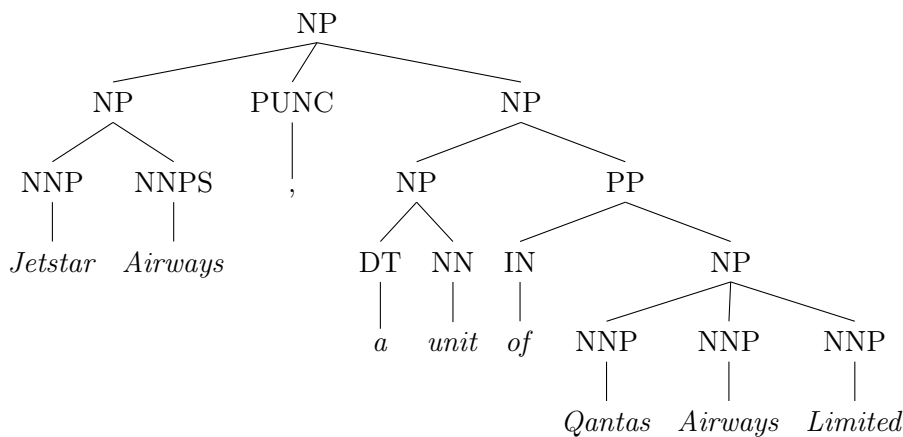


Figure 2.3: A sentence broken down to its Phrase Structure (sub-sentences), also known as constituency parsing.

After discovering the structure between the phrases and its sub-phrases, finding the syntactical dependency between words themselves is also a very interesting and useful task. Continuing on the same example, one can see how the sentence states the simple fact that one company (*Jetstar*) is a unit of another company (*Qantas*). The dependency tree in this case tells us that *Jetstar*/*NN* is the head of another noun *Airways* and that the relation between them is of the *compound* type. This relation is held between any noun that serves to modify the head noun and in this case indicate that this single entity is formed by two different words that are nouns. Note that *Jetstar* has no parents and thus is the root of the sentence. The dependencies can be fully viewed in Figure 2.4.

The next relation is the *appos* from '*Airways*' to '*unit*'. The appositional modifier relation indicates that the noun immediately to the right of the first noun that serves to define or modify the meaning of the former. One now already knows that '*Jetstar Airways*' is a '*unit*', and, in the busi-

ness context, it probably means that it is a company that belongs to a bigger company.

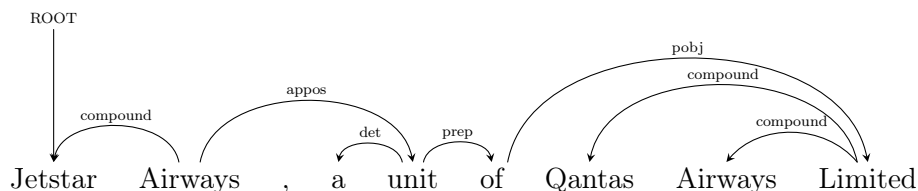


Figure 2.4: A sentence and the dependencies between the words.

Continuing the analysis, the next relation is of the *prep* type and it indicates a prepositional modifier of a verb, adjective, or noun (our case), and it serves to modify the meaning of the verb, adjective, noun, or even another preposition. Note that this relation simply indicates the word pointed to by the edge is the preposition (in this case ‘*of*’). The next relation, *pobj*, indicates the actual object of the preposition and the noun phrase following the preposition and what it related to. This tree was created by the spaCy dependency parser [17, 30]. This same tree can be visualised in a more traditional tree structure in Figure 2.5. Although not in this example, another very common relation is the relative clause modifier *recmod* (or *relcl* in some notations) relation. A relative clause modifier of an NP (*Noun Phrase*) is a relative clause modifying the NP. The relation then in the tree points from the head noun of the NP to the head of the relative clause, normally a verb. The explanations for the cited relations in this document were obtained from the Stanford typed dependencies manual [22], and the Universal Dependencies (UD) project [28], both which are reference for the possible relation and contain a full lists of their meanings.

The software that is able to output a syntactical dependency tree, given a sentence, is called a dependency parser. Several different methods can be used to achieve this. One way is by defining dependency grammars, and then parsing text using these grammar. The grammar would contain words and its possible heads, and it would be applied repeatedly into the text in a process called cascaded chunking [5]. More recent methods use a process called Shift-reduce, in which the sentence is kept in a queue with the left-most token in front of the queue. The model could then decide between applying 3 operations:

1. Shift: move one token from the queue to stack.
2. Reduce left: top word on stack is head of second word.
3. Reduce right: second word on stack is head of top word.

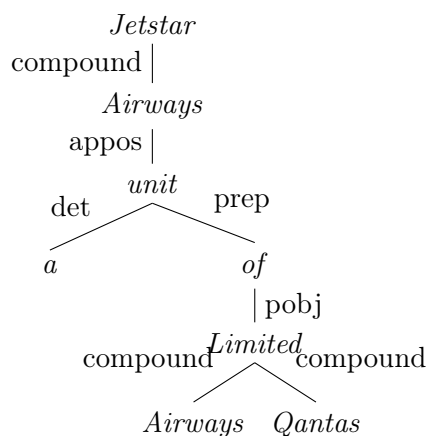


Figure 2.5: A sentence and the dependency tree, showing the syntactical relation between these words.

A model is then trained to predict, given a text that is added to the queue, what is the next move that it should take, and what is the sequence of moves that will result in the best possible final dependency tree. This is done in a monotonic manner, in the sense that once a decision is made by the model, it cannot change it. Full working examples of this method are described in [7]. Other methods also use these 3 possible decisions, but allow the parser to be non-monotonic and go back in the tree and change previous decisions given new evidence from the features, such as spaCy and its dependency parser [17, 30]. Another method also uses the same set of decisions, but does a beam-search observing multiple partial hypotheses and keeping them at each step, with hypotheses only being discarded when there are several other higher-ranked hypotheses under consideration, such as the Syntaxnet parser [34, 3]. All these cited models use neural networks as the method to form the model.

As a further example, besides providing dependency parsing tree, spaCy also provides an iterator so you can obtain what is called the noun chunks of the document. The noun chunks are smaller pieces of the sentences within the document that are base noun phrases, or a 'NP chunk' (as per Figure 2.3). They are noun phrases that do not permit other NPs to be nested within it so no NP-level coordination (e.g.: '*cat/NN and/CONJ dog/NN*'), no prepositional phrases, and no relative clauses [30].

Other problems tackled by the Natural Language Processing discipline and that are relevant for Information Extraction are *Coreference resolution* and *pronominal anaphora resolution*. Coreference resolution intends to define all possible entities that a text can reference to in some sort of definitive list, or more precisely a *discourse model*, find in the text all the chained references to these entities, and link them to the specific entities. While very

similar in nature, pronominal anaphora resolution is more simple as it is the problem of resolving in a given sentence to which previous NN (noun) or NNP (proper noun) a single PRP (pronoun) refers to [18].

Take for an example the sentence ‘*John is a quiet guy, but today he is furious.*’, the initial mention of the entity *John* appears in the first token. Token five talks about a *guy*, which although not a pronoun, is still a reference to the previous entity in the first word. The ninth token is a pronoun and again refers to the same *John*, so is part of the chain of mentions. The full resolution chain would be denoted in Figure 2.6.

Normally these systems work towards analysing pairs of tokens using a probabilistic model, and then decide how likely they are references for the same entity. More recent approaches also group possible tokens in a cluster and use cluster-level features to determine the chains of coreferences [9]. The tool used to extract Figure 2.6 is the Stanford Coreference Resolution annotator [9], which is part of the latter group of tools that use cluster level features. This annotator is also part of the CoreNLP toolset [20].

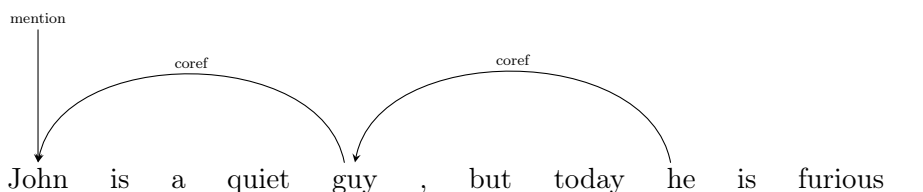


Figure 2.6: A sentence, the mentions of an entity, and the proposed coreference resolutions.

All linguistic data mentioned in this section is data that can be obtained from the raw text itself, and is then the base for several features in which methods for Information Extraction act on.

2.2 Information Extraction

The IE (Information Extraction) process is described by the following subtasks: Named Entity Recognition (NER), Coreference Resolution, Entity Disambiguation, Relation Extraction (RE), Event Detection, and Temporal Analysis [18]. The main subtasks relevant to this report will be described further in this section.

Once the information is extracted it is then used for tasks such as Template Filling [18], Question and Answering systems [25], or stored as a Knowledge Graph for downstream logical reasoning or for further queries.

Named Entity Recognition is the process of, given a sentence, identify and extract what are the entities that are part of it. Once the entity is detected, it needs to be classified within the classes of the given domain - in

[_{PER} James Cook] was born on 27 October 1728 in the village of [_{LOC} Mar-
ton] in [_{COUNTY} Yorkshire].

Table 2.3: An example of Named Entity Recognition (NER).

the spirit of the previous examples this would be e.g.: *CITY* or *PERSON*. Different types of entities are relevant to context of the data being worked with. A few approaches exist for the problem of NER, mostly related to Pattern Matching or Sequence Classification.

Pattern	Would yield ENTITY of type
[PERSON] was born	PERSON
in the village of [LOC]	LOCATION
in [LOC]	LOCATION

Table 2.4: Examples of Named Entity Recognition (NER) patterns, based on the sentence from Table 2.3.

Observe, for an example, the sentence in Table 2.3. Several articles regarding prominent figures, either historical or of our current society, can be of the format ‘*Jimi Hendrix was born*’. One approach might be Pattern Matching, which is to mine the input natural language text while looking for the pattern ‘[*ENTITY*] was born’, using Regular Expressions (Finite-State Automata) [18]. The entities found by this pattern would then also receive the *PERSON* class. This pattern would miss the sentence ‘*Jimi Hendrix, born in Seattle*’ since it does not fit the pattern and, because of this, one generally needs to build a list or database of patterns to work with in a corpus. An example of such database was generated by the PATTY system [27]. Table 2.4 depicts other possible similar patterns.

Another way to extract entities from text is to frame the NER problem as a Sequence Classification problem, similar to the POS tagging problem described earlier. It requires the training of a classifier in which, given the class of the previous word, and other surrounding features of the current word, will attempt to guess if the current word is an entity, and if it is, also guesses its class.

To achieve this, previously annotated data with existing sentences and its entities is needed. This can be obtained by manually labelling data, or by semi-automated methods, like the one proposed later by this document. The format in which this annotated data is provided varies, however the IOB format (Table 2.5) is more commonly used in several of the NER tools, including NLTK [5] and the popular Stanford Named Entity Recognizer (NER) [13], part of the Stanford CoreNLP [20]. Stanford CoreNLP provides a set of natural language analysis and information extraction tools.

The IOB format also helps remove ambiguity in case there are two con-

Word	Tag
James	B-PERSON
Cook	I-PERSON
was	O
born	O
on	O
27	B-DATE
October	I-DATE
1728	I-DATE
in	O
the	O
village	O
of	O
Marton	B-LOC
in	O
Yorkshire	B-LOC
.	O

Table 2.5: Example of IOB-formatted sentence used to train classifiers for the Named Entity Recognition (NER) task, based on the sentence from Table 2.3.

tiguous entities of same class without any word tagged as *O* in between. In practice these cases are somewhat rare in several domains, and even when trained with such tags classifiers struggle to accurately determine the boundaries of an entities, and thus a simplified version of this annotation without the *B*- and *I*- prefixes is more commonly used [33].

The Stanford Named Entity Recognizer (NER), also known as CRF-Classifier [13], provides a general implementation of (arbitrary order) linear chain Conditional Random Field (CRF) sequence models. A CRF is a conditional sequence model which represents the probability of a hidden state sequence given some observations.

Several relevant features can be used as an input during the training of a NER CRF classifier model. In Table 2.6, examples are presented. The Word Shape feature is an interesting addition from recent research, as it captures the notion that most entities are written in capital letters, or starting with capital letter, or containing numbers in the middle of the word, and other specific shapes.

In addition to the above methods another useful technique is the use of gazetteers. Gazetteers are common for geographical data, where government provided lists of names can contain millions of entries names for all manner of locations along with detailed geographical, geologic and political information [18].

Feature	Description
Word	The current word being classified.
N-grams	A feature from n-grams, i.e., sub-strings of the word.
Previous Class	The class of the immediate previous word.
Previous Word	The previous word.
Disjunctive	Disjunctions of words anywhere in the left or right.
Word Shape	The shape of the word being processed captured using. In general replaces numbers with <i>d</i> , <i>x</i> to lower-case letters, and <i>X</i> to upper-case letters.

Table 2.6: Examples of features used to train the CRFClassifier [13].

Relation Extraction (RE) is the ability to discern the relationships that exist among the entities detected in a text [18], and is naturally the next challenge after being able to detect entities. It is generally denoted as a triplet: two entities, and the one relation between them (Table 2.7). It can be done using Pattern Matching, Classifiers, or purely by exploiting linguistic data available from a sentence. The previously described Pattern Matching technique from NER can be improved upon in the Relation Extraction step, and involve more than one entity, yielding binary relations. This approach is used in tools such as PROSPERA [26] or those mined by PATTY [27]. More specifically, examples of patterns mined by PATTY for the *graduatedFrom* relation are seen in Figure 2.7.

Located_In(Kiel, Germany)

Table 2.7: An example of a triplet that represents a relation.

Pattern	Domain	Range	▲ Confidence	Support/Co-occurrence
actedIn	graduated [[con]] entered;	person university	1	4
created	completed [[prp]] university studies in;	person organization	1	2
dealsWith	attended before studying law at;	person organization	1	2
diedIn	sociology at;	person university	1	2
directed	speaking [[con]] representing;	person university	1	2
graduatedFrom	earned in economics from;	person organization	1	2
happenedIn	graduated from [[det]] department of;	person university	1	2
hasAcademicAdvisor	pursued [[det]] degree at;	person university	1	2
hasCapital	met [[prp]] [[adj]] wife [[det]];	person organization	1	2
hasChild	[[det]] member [[det]] governing body of;	person university	1	2
hasWonPrize	worked [[con]] received from;	person university	1	2
holdsPoliticalPosition	[[det]] degree in economics from;	person university	1	2
influences	entered;	person organization	0.975	24
isCitizenOf	obtained [[det]] doctorate at;	person university	0.974	2
isKnownFor	majoring at;	person university	0.966	7
isLeaderOf	[[det]] graduate student in;	person organization	0.965	4
isLocatedIn	received in mathematics from;	scientist university	0.963	3
isMarriedTo	graduated [[con]] with honors from;	person organization	0.947	3
isPoliticianOf	accepted [[det]] chair in;	person university	0.947	2
livesIn				
participatedIn				
playsFor				
produced				
wasBornIn				
worksAt				

Figure 2.7: An example of patterns extracted from PATTY for the *graduatedFrom* relation.

PROSPERA's main technique is that not only it obtain facts based on

a small set of initial seed patterns, but also obtain new candidate patterns that can be extrapolated from the corpus based on the mined known facts. Once the process of obtaining new candidate patterns finishes, these are evaluated and then added to the existing pattern repository for re-use. The whole process then iterates again finding even more facts from these new patterns, and new candidate patterns [26]. Moreover, another interesting characteristic of the PROSPERA’s approach is the care in Entity Disambiguation. For an example, given that in a text it finds a name, such as ‘*Captain James Cook*’. It then uses a knowledge base such as YAGO [32] to compare this with existing known entities, using techniques such as N-gram comparison [26]. With such effort, PROSPERA is able to know that *Captain James Cook* and *James Cook* are actually the same entity with a certain confidence, and thus don’t differentiate these and assigns ‘*Captain James Cook*’ as a representation of the canonical unambiguous entity ‘*James Cook*’. This helps in several ways, such as: it can then know other information about this entity, such as the fact that it is of the class *PERSON*; and it can also facilitate future queries in this knowledge base, centralizing the new information found about this existing entity.

Another tool in the Stanford CoreNLP package, the Relation Extractor [33] is a classifier to predict relations in sentences. This program has a model that extracts binary relations between entity mentions in the same sentence. The output is normally in the XML [12] format and denotes the tokens of each sentence, the possible relations, and the confidence level of these relations. The XML demonstrated in Figure 2.8 depicts a guess that 2 words in the sentence ‘..., including approaches that use parallel computation [1, 2, 6, 13, 24].’ have the *Uses* relation with a confidence above 70%. The classifier in this case also indicates that one of the entities is of the class *CONCEPT*.

As part of its training, annotated relation mentions are used as an input together with the text it belongs to, and the annotation becomes a positive example for the corresponding label, while all the other possible combinations between entity mentions in the same sentence become negative examples. The feature set models the relation between the arguments by using the distance between the relation arguments, the syntactic path between the two arguments, using both constituency and dependency representations.

Note how at this point it’s important to note the notion of a *pipeline* of natural language processing tasks. Stanford’s CoreNLP is able to, by itself, perform all the steps needed to produce such relation extraction output from the raw text input. The steps of this pipeline are executed in a certain order, as they depend on the previous step (e.g.: POS tagging is needed for Dependency Parsing). In this case, the process was:

1. Tokenize;
2. Sentence Splitter;

```

1 <relation id="RelationMention-6728">Uses
2   <arguments>
3     <entity id="EntityMention-1324">O
4       <span start="46" end="47"/>
5       <probabilities />
6     </entity>
7     <entity id="EntityMention-1319">CONCEPT
8       <span start="36" end="37"/>
9       <probabilities />
10    </entity>
11  </arguments>
12  <probabilities>
13    <probability>
14      <label>Uses</label>
15      <value>0.7379587661527921</value>
16    </probability>
17    <probability>
18      <label>Improves</label>
19      <value>0.06475441029357998</value>
20    </probability>
21  </probabilities>
22 </relation>

```

Figure 2.8: An example of relation extracted with the Stanford Relation Extractor that demonstrated the ‘Uses’ relation.

3. Part-of-Speech tagging;
4. Lemmatization;
5. Constituency parsing;
6. Dependency Parsing;
7. Named Entity Recognition;
8. and finally, Relation Extraction.

The Stanford Relation Extractor comes with a model that was trained to extract the following relations: *Live_In*, *Located_In*, *OrgBased_In*, *Work_For* - and the following classes: *PERSON*, *ORGANIZATION*, *LOCATION*. There are big challenges if one attempts to train the model for any relation outside of these, mainly in obtaining or generating annotated data to train the classifier as to generate a useful model. There are attempts in which relation extraction is not based on annotated data, but on linguistic characteristics of the text itself, such as its semantics. These tools are normally called Open Relation Extractors and will be further described in Section 3.3.

2.3 Knowledge Graphs

Knowledge Graphs contain a valuable of information in a structured format, traditionally originally mined from table-like structures form places like Wikipedia [38] tables [19], or from processes like Information Extraction as described in the previous section. It can be used for a diverse range of applications, such as helping other systems reason about quality of harvested facts [32], provide table-like facts about an entity [14], and question-answering systems [16]. Moreover, recent years have witnessed a surge in large scale knowledge graphs, such as DBpedia [19], Freebase [6], Googles Knowledge Graph [14], and YAGO [32].

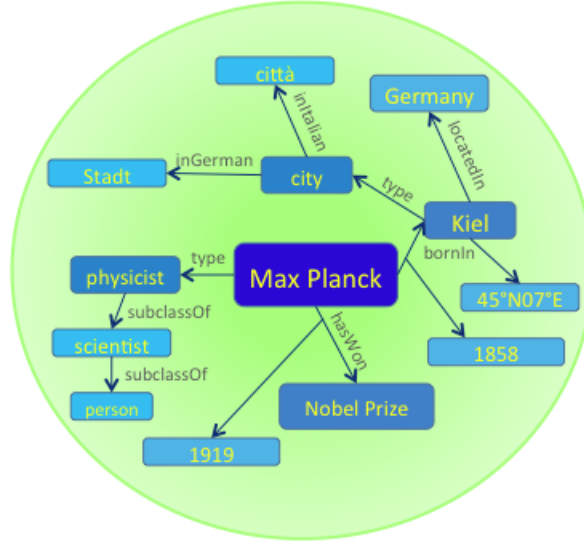


Figure 2.9: An example of knowledge graph from [32] plotted with vertices and edges.

The Knowledge Graph name follows from the data structure that is created from the facts in its final form, a graph with nodes representing entities and edges representing various relations between entities. In Figure 2.9, it is possible to observe an example plotted in this form. The list of possible entities classes, and allowable relations between entities is known as a schema. The schema represented in Figure 2.9 is detailed in Table 2.9; one can observe that, as an example, ‘*Max Planck*’ is an entity of the type *physicist*.

`type(A, D) :- type(A, B), subclassOf(B, C), subclassOf(C, D)`

Table 2.8: This entailment example allows one to assert that `type(Max Planck, person)` is also true, based on the fact tuples presented in Table 2.9.

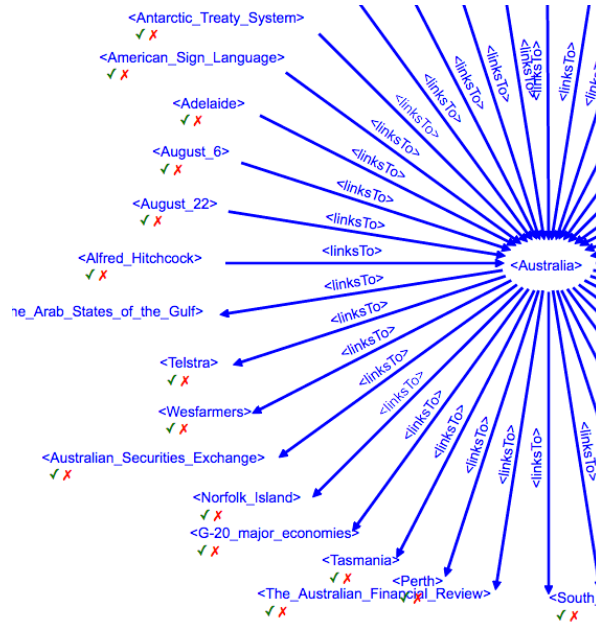


Figure 2.10: An example of patterns existants in YAGO.

Moreover, based on the facts presented, entailments can be made and one trivial example is denoted in Table 2.8. More complex examples of possible reasoning can be seen in [33]. This is equivalent to traversing the graph from a node that represents a more specific information, to a node that represents a more general information - e.g.: another possible child node of ‘*scientist*’ could be the type ‘*biologist*’.

```
type(Max Planck, physicist)
subclassOf(physicist, scientist)
subclassOf(scientist, person)
bornIn(Max Planck, Kiel, 1858)
type(Kiel, city)
locatedIn(Kiel, Germany)
hasWon(Max Planck, Nobel Prize, 1919)
```

Table 2.9: Some facts regarding Max Planck, also depicted in Figure 2.9.

This example denotes a classical domain, more precisely important persons, companies, locations, and the relations between them, in which Information Extraction (IE) tools have been very successful on.

As mentioned previously, YAGO [32] is a prominent Knowledge Graph database, and possesses several advanced characteristics. Every relation in its database is annotated with its confidence value. See the example of the

resulting graph in Figure 2.10. Moreover, YAGO combines the provided taxonomy with WordNet [24] and with the Wikipedia category system [38], assigning the entities to more than 350,000 classes. This allow for very powerful querying. Finally, it attaches a temporal and a spacial dimension to many of its facts and entities, being then capable to answer questions such as *when* and *where* such event took place.

WordNet is a semantically-oriented dictionary of English, similar to a traditional thesaurus but with a richer structure [5]. More specifically, it provides relations to synonyms, hypernyms and hyponyms, among others.

Chapter 3

Analysis and Related Work

This work intends to deliver a tool or a process in which one can extract information from academic text, more specifically Computer Science papers from the Database and Data Mining topics. The intention is to obtain entities, and relations between these entities. The motivation is that, with such tool, one could for an example:

- Historically research algorithms that were mostly used during a certain time period;
- Find which algorithms are used to resolve, or related to, a certain problem;
- Find techniques that improve a certain algorithm problem, among others.

3.1 Analysis of Academic Text

The corpus of text used was generated utilising papers published from the following conferences during various years: ACL [37], EMNLP [11], ICDE [10], SIGMOD [1], VLDB [36]. More specifically, the section of *Related Work* of these papers were the ones used to build the corpus. This was done due to the characteristics and patterns of this section compared to the rest of the paper. After careful reading, we observed that the *Related Work* section generally contains objective comparisons between other algorithms or softwares in contrast with more opaque or abstract explanations from other parts of the paper. This would mean that this section was a good candidate to start the analysis from. Note the following examples of sentences from the *Related Work* section of papers from the corpus:

1. *‘Bergsma et al (2013) show that large-scale clustering of user names improves gender, ethnicity and location classification on Twitter.’*

2. *‘N-Best ROVER (Stolcke et al, 2000) improves the original method by combining multiple alternatives from each combined system.’*
3. *‘By partitioning the velocity space, the Bdual -tree improves the query performance of the B x -tree.’*

Entities from academic text in this setting are not as straightforward to define as in, for an example, business news, or criminal news. Observe the following sentence:

- *‘Japan’s Toshiba Corp said it had nominated Satoshi Tsunakawa, a former head of its medical equipment division, to be its next chief executive officer.’*

Text	Entity Type
Japan	LOCATION
Toshiba Corp	ORGANIZATION
Satoshi Tsunakawa	PERSON

Table 3.1: Examples of Named Entity Recognition (NER) from the business news text example.

From the news text example above, Table 3.1 lists the entities that are clearly notes in the text. One can observe a very strong feature which is the common capitalization of the first letter of each of these entities. Another characteristic is how entities from this news text example are *global* or *unconditional*: ‘Japan’ is a location regardless of any condition or any context in this document. Another observation is that, referring to the Stanfords Relation Extraction default relations, ‘Toshiba Corp’ is an organisation *Located_In* ‘Japan’ regardless of other context in this document. This contrasts with concepts and their relations observed in academic papers, thus that while ‘*large-scale clustering*’ has the *Improves* with ‘*gender classification*’ in the context of the paper where this data is presented, it might not be true in all cases.

Text	Entity Type
Bergsma et al (2013)	AUTHOR
large-scale clustering	CONCEPT
gender classification	CONCEPT
ethnicity classification	CONCEPT
location classification	CONCEPT

Table 3.2: Examples of Named Entity Recognition (NER) from the academic text example.

Moreover, the entities in Table 3.2 are harder to classify in universally agreed classes. For an example, ‘*gender classification*’ can be considered an action, or a task, or an algorithm. More generally, one can simply classify these as concepts.

```
IsA(Concept,Concept)
SimilarTo(Concept,Concept)
Improves(Concept,Concept)
Employs(Concept,Concept)
Uses(Concept,Concept)
Supports(Concept,Concept)
Proposes(Author,ComplexConcept)
Introduces(Author,ComplexConcept)
```

Table 3.3: Some observed and possible relations between concepts.

Other possible relations from the Stanford Relation Extractor standard relations that are applicable to the above news text example are: *OrgBased_In* (again for ‘*Toshiba Corp*’ and ‘*Japan*’) and *Work_For* regarding its newly placed chief executive officer. Again, contrasting with the academic text, one might consider relations such as the one possible between concepts as denoted in Table 3.3. In fact, by analysing the corpus for the top 50 words in the singular third-person form, such as ‘*improves*’ or ‘*employs*’, one can have an idea of the possible relations that can be extracted. This process is illustrated in Figure 3.1. The same data is illustrated in Figure 3.2, but with the top two words removed (‘*is*’ and ‘*has*’).

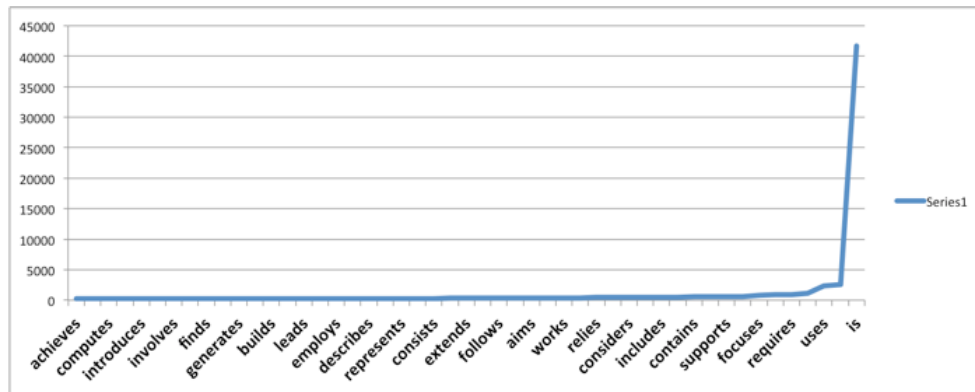


Figure 3.1: Samples of the most common words in the singular third-person form.

One of the initial attempts to explore how to extract information from the generated corpus was to use the Stanford Named Entity Recognizer (NER) to recognize the concepts discussed so far in the academic text. To

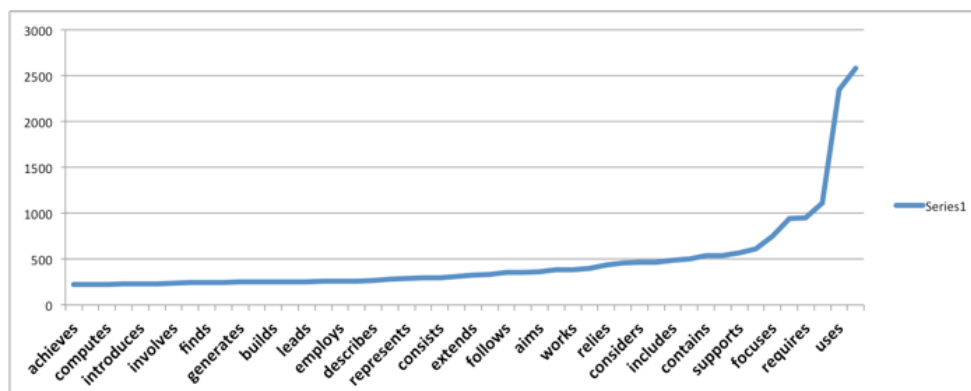


Figure 3.2: Samples of the most common words in the singular third-person form, after removing the top 2 words.

do so, a small set of around 20 papers' *Related Work* section was annotated for the concepts contained in them using Brat [31].

- Attempt of using Stanford NER

- Attempt of using Stanford REL

- Attempt of using Stanford OpenIE

To be completed: Querying a corpus? How to check for similar words statistically, etc

Please refer to Section 4.1 for more information on tools mentioned in this section.

3.2 Peculiarities with Academic Text

Problem accross domains, strict language vs. creative language.

Sarcasm.

We are extracting explicit information, and not trying to infer implicit information. That would be a different logical problem that would need to reason about using a database of information to compare what is being learned from the natural text with.

Difficult to manually tag.

Existing tools that are based in models don't come equipped to predict relation in this kind of text. E.g.:

Relations that could be extracted, e.g. denes, don't appear with relevant Entities.

A lot of coreference problems, e.g. "their work". Even when their work is simply a paper reference - with unclear. One could trivially parse the above reference with the actual Entity name of the system or algorithm or technique elaborated in the reference paper as to mine the relations between the proper entities.

Methods:

* Use semi-automatic methods to collect research literature and convert them into plain text with certain markups. * Use existing open source solutions to parse the store the input data. * Build a pipeline to extract entities and relationship from the input data. * Perform entity disambiguation and linking to a reference Knowledge Graph (e.g., Yago2 or DBPedia). * Design effective postprocessing methods to improve the quality of the extraction. * Evaluate the entire extraction system.

3.3 Related Work

- Open Information Extraction

Other methods for relation extraction such as the ones from papers:

ClausIE: Clause-Based Open Information Extraction

Extreme Extraction: Only One Hour per Relation

IKE - An Interactive Tool for Knowledge Extraction

Chapter 4

Developed Workflow

The first step to in achieving our goal is obtaining the raw natural language text from papers in the area. Most research products such as thesis, papers, or any other report, are mostly available primarily in the Portable Document Format (PDF) format [2] - which then needs to then be parsed into a raw text in an automated manner.

4.1 Tools

4.1.1 PDF extraction generates noisy output

4.1.2 Programming Languages

Java

Python

4.1.3 Python Modules

requests

BeautifulSoup4

standoff2conll

corpkit

Tregex

corenlp-xml

4.1.4 NLTK**4.1.5 Brat****4.1.6 Jupyter****4.1.7 Parsey****4.1.8 SpaCy****4.1.9 Stanford CoreNLP****4.1.10 Graphviz****4.2 Developed process and programs**

We leverage on existing pipeline for NER and Relation Extraction (Stanford) , instead of independent tools.

We use IO notation instead of IOB notation for entity mention labels (e.g., the labels for the tokens over the Seattle Seahawks on Sunday (from Figure 1) are encoded as O O NFLTEAM NFLTEAM O DATE). The IO notation facilitates faster inference than the IOB or IOB2 notations with minimal impact on performance, when there are fewer adjacent mentions with the same type. AS PER STANFORD RELATION EXTRACTOR PAPER.

TALK ABOUT FEATURES WE PICKED FOR NER.

TALK ABOUT FEATURES WE PICKED FOR REL.

Open information extraction (open IE) has been shown to be useful in a number of NLP tasks, such as question answering (Fader et al., 2014), relation extraction (Soderland et al., 2010), and information retrieval (Etzioni, 2011).

We did not implement coreference.

Expected Outcome:

1. A system that can build a knowledge graph from research literatures.
2. A written report about the detailed designs and implementations of this system.
3. A seminar to present the process and outcome of this project.

Chapter 5

Results

- Evaluation

Chapter 6

Conclusion and Future Work

- What drives the field forwardd

Furhter ideas: - Research relations through time. You could have a certain feature. - Reinforcements made to definitions in other papers. - Events, such as changes in conclusions - e.g.: this was the better technique, now this other technique is the best.

Bibliography

- [1] *About SIGMOD*. <https://sigmod.org/about-sigmod/>. [Online; accessed 09-October-2016]. 2016.
- [2] *Adobe: What is PDF?* <https://acrobat.adobe.com/us/en/why-adobe/about-adobe-pdf.html>. [Online; accessed 14-August-2016]. 2016.
- [3] Daniel Andor et al. ‘Globally Normalized Transition-Based Neural Networks’. In: *CoRR* abs/1603.06042 (2016).
- [4] *Bing Knowledge and Action Graph*. <https://www.bing.com/partners/knowledgegraph>. [Online; accessed 14-August-2016]. 2016.
- [5] Steven Bird, Ewan Klein and Edward Loper. *Natural Language Processing with Python*. O’Reilly Media, 2009.
- [6] Kurt Bollacker et al. ‘Freebase: A Collaboratively Created Graph Database for Structuring Human Knowledge’. In: *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*. SIGMOD ’08. New York, NY, USA: ACM, 2008, pp. 1247–1250.
- [7] Danqi Chen and Christopher Manning. ‘A Fast and Accurate Dependency Parser using Neural Networks’. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Doha, Qatar: Association for Computational Linguistics, 2014, pp. 740–750.
- [8] *CiteSeerX*. <http://csxstatic.ist.psu.edu/about>. [Online; accessed 27-September-2016]. 2016.
- [9] Kevin Clark and Christopher D. Manning. ‘Entity-Centric Coreference Resolution with Model Stacking’. In: *Association for Computational Linguistics (ACL)*. 2015.
- [10] *Data Engineering, International Conference on*. <http://ieeexplore.ieee.org/xpl/conhome.jsp?reload=true&punumber=1000178>. [Online; accessed 09-October-2016]. 2016.
- [11] *Empirical Methods in Natural Language Processing*. <https://www.aclweb.org/website/emnlp>. [Online; accessed 09-October-2016]. 2016.

- [12] *Extensible Markup Language (XML) 1.0 (Fifth Edition)*. <https://www.w3.org/TR/REC-xml/>. [Online; accessed 09-October-2016]. 2008.
- [13] Jenny Rose Finkel, Trond Grenager and Christopher Manning. ‘Incorporating Non-local Information into Information Extraction Systems by Gibbs Sampling’. In: *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*. ACL ’05. Stroudsburg, PA, USA: Association for Computational Linguistics, 2005, pp. 363–370.
- [14] *Google Knowledge Graph*. <https://www.google.com/insidesearch/features/search/knowledge.html>. [Online; accessed 14-August-2016]. 2016.
- [15] *Google Scholar*. <https://scholar.google.com/intl/en/scholar/about.html>. [Online; accessed 27-September-2016]. 2016.
- [16] Ben Hixon, Peter Clark and Hannaneh Hajishirzi. ‘Learning Knowledge Graphs for Question Answering through Conversational Dialog’. In: *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Denver, Colorado: Association for Computational Linguistics, 2015, pp. 851–861.
- [17] Matthew Honnibal and Mark Johnson. ‘An Improved Non-monotonic Transition System for Dependency Parsing’. In: *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Lisbon, Portugal: Association for Computational Linguistics, 2015, pp. 1373–1378.
- [18] Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. 1st. Upper Saddle River, NJ, USA: Prentice Hall PTR, 2000.
- [19] Jens Lehmann et al. ‘DBpedia - A Large-scale, Multilingual Knowledge Base Extracted from Wikipedia’. In: *Semantic Web Journal* 6.2 (2015), pp. 167–195.
- [20] Christopher D. Manning et al. ‘The Stanford CoreNLP Natural Language Processing Toolkit’. In: *Association for Computational Linguistics (ACL) System Demonstrations*. 2014, pp. 55–60.
- [21] Mitchell P. Marcus, Mary Ann Marcinkiewicz and Beatrice Santorini. ‘Building a Large Annotated Corpus of English: The Penn Treebank’. In: *Comput. Linguist.* 19.2 (June 1993), pp. 313–330.
- [22] Marie-Catherine De Marneffe and Christopher D. Manning. *Stanford typed dependencies manual*. 2008.
- [23] *Microsoft Academic Graph*. <https://www.microsoft.com/cognitive-services/en-us/academic-knowledge-api>. [Online; accessed 27-September-2016]. 2016.

- [24] George A. Miller. ‘WordNet: A Lexical Database for English’. In: *Commun. ACM* 38.11 (Nov. 1995), pp. 39–41. ISSN: 0001-0782.
- [25] Marie-Francine Moens. *Information Extraction: Algorithms and Prospects in a Retrieval Context (The Information Retrieval Series)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [26] Ndapandula Nakashole, Martin Theobald and Gerhard Weikum. ‘Scalable Knowledge Harvesting with High Precision and High Recall’. In: *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining*. WSDM ’11. New York, NY, USA: ACM, 2011, pp. 227–236.
- [27] Ndapandula Nakashole, Gerhard Weikum and Fabian Suchanek. ‘PATTY: A Taxonomy of Relational Patterns with Semantic Types’. In: *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*. EMNLP-CoNLL ’12. Stroudsburg, PA, USA: Association for Computational Linguistics, 2012, pp. 1135–1145.
- [28] Joakim Nivre et al. *Universal Dependencies 1.3*. <http://universaldependencies.org/>. LINDAT/CLARIN digital library at the Institute of Formal and Applied Linguistics, Charles University in Prague. 2016.
- [29] *Semantic Scholar*. <http://allenai.org/semantic-scholar/>. [Online; accessed 20-August-2016]. 2016.
- [30] *spaCy: Industrial-strength Natural Language Processing*. <https://spacy.io/>. [Online; accessed 27-September-2016]. 2016.
- [31] Pontus Stenetorp et al. ‘BRAT: A Web-based Tool for NLP-assisted Text Annotation’. In: *Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics*. EACL ’12. Stroudsburg, PA, USA: Association for Computational Linguistics, 2012, pp. 102–107.
- [32] Fabian M. Suchanek, Gjergji Kasneci and Gerhard Weikum. ‘Yago: A Core of Semantic Knowledge’. In: *Proceedings of the 16th International Conference on World Wide Web*. WWW ’07. New York, NY, USA: ACM, 2007, pp. 697–706.
- [33] Mihai Surdeanu et al. ‘Customizing an Information Extraction System to a New Domain’. In: *Proceedings of the ACL 2011 Workshop on Relational Models of Semantics*. RELMS ’11. Stroudsburg, PA, USA: Association for Computational Linguistics, 2011, pp. 2–10.
- [34] *SyntaxNet: Neural Models of Syntax*. <https://github.com/tensorflow/models/tree/master/syntaxnet>. [Online; accessed 01-October-2016]. 2016.

- [35] Kristina Toutanova et al. ‘Feature-rich Part-of-speech Tagging with a Cyclic Dependency Network’. In: *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1*. NAACL ’03. Stroudsburg, PA, USA: Association for Computational Linguistics, 2003, pp. 173–180.
- [36] *Very Large Data Base Endowment Inc. (VLDB Endowment)*. <http://www.vldb.org/>. [Online; accessed 09-October-2016]. 2016.
- [37] *What is the ACL and what is Computational Linguistics?* <https://www.aclweb.org/website/what-is-cl>. [Online; accessed 09-October-2016]. 2016.
- [38] *Wikipedia, The Free Encyclopedia*. <http://www.wikipedia.org/>. [Online; accessed 21-August-2016]. 2010.