

Ben-Gurion University of the Negev

Faculty of Natural Science

Department of Computer Science

A STUDY IN HEBREW PARAPHRASE IDENTIFICATION

Thesis submitted as part of the requirements for the
M.Sc. degree of Ben-Gurion University of the Negev

by

Gabriel Stanovsky

May 2012

Subject: A Study in Hebrew Paraphrase Identification

Written By: Gabriel Stanovsky

Advisor: Professor Michael Elhadad

Department: Computer Science

Faculty: Natural Science

Ben-Gurion University of the Negev

Signatures:

Author: Gabriel Stanovsky

Date

Advisor: Professor Michael Elhadad

Date

Dept. Committee Chairman

Date

Contents

1	Introduction	1
1.1	Formal Definitions	2
1.1.1	Textual Entailment	2
1.1.2	Paraphrase	2
1.1.3	Related Shared Tasks	3
1.1.4	Paraphrasing in Hebrew	3
1.2	Motivation	3
2	Research Questions and Objectives	4
2.1	Research Questions	4
2.1.1	Hebrew Specific properties	4
2.1.2	Hebrew Datasets for Paraphrasing	5
2.1.3	Method Adaptation for Hebrew	5
2.1.4	Learning Hebrew Paraphrasing Rules	5
2.1.5	Producing Hebrew Resources for Other NLP Uses	5
2.2	Objectives	6
3	Previous Work	7
3.1	Paraphrase Identification	7
3.2	Paraphrase Generation	8
3.3	Paraphrases Extraction	8
3.4	Hebrew Paraphrasing	9
4	Linguistic Background	10
4.1	Paraphrasing Through the Prism of Different Linguistic Schools	10
4.1.1	Functionalism	10
4.1.1.1	Defining Paraphrases in Functionalist Terms	11
4.1.2	Generativism	11
4.1.2.1	Transformational Grammar	12
4.1.2.2	Defining Paraphrases in Generativist Terms	12
4.1.3	Cognitive	12

4.1.3.1	Natural Semantic Metalanguage (NSM)	13
4.2	Reasons for Paraphrasing	14
4.2.1	Principle of Conventionality	14
4.2.2	Principle of Contrast	14
4.2.3	Small Scale paraphrasing	14
4.2.4	Large Scale paraphrasing	15
4.3	Parse Trees of Paraphrases	16
4.3.1	Phrase Structure Grammar	16
4.3.2	Dependency Grammar	17
4.3.3	Examples of Paraphrases Pair in Phrase Structure Grammar vs Dependency Grammar	18
4.3.4	Discussion	20
5	Methodology	23
5.1	Plan of Work	23
5.2	Learning Methods	25
5.2.1	Neural Networks	25
5.2.1.1	Outline	25
5.2.1.2	Neuron	25
5.2.1.3	Weights	26
5.2.1.4	Network Computation	26
5.2.1.5	Learning	26
5.2.1.6	Backpropagation	26
5.2.1.7	Layered Networks	27
5.2.2	Deep Learning	27
5.2.2.1	Curriculum Learning	27
5.2.2.2	Multi Task Learning	28
5.2.2.3	Word Embeddings	28
5.2.3	Auto Encoders	28
5.2.3.1	Outline	29
5.2.3.2	Neural Network	29
5.2.3.3	Encoding	29
5.2.3.4	Decoding	30
5.2.3.5	Uses in NLP	30
5.2.4	Tree Matching	30
5.2.4.1	Definition	30
5.2.4.2	Tree Match Score	31
5.2.4.3	Tree Matching is NP-Complete	31
5.3	Hebrew Paraphrasing	31

5.3.1	Producing Deep Learning Embedding for Hebrew	31
5.3.1.1	Overview	31
5.3.1.2	Preprocessing and Hebrew Adaptation	32
5.3.1.3	Language Model	33
5.3.1.4	Curriculum Learning	33
5.3.1.5	Stochastic Sampling of the Error Function	33
5.3.1.6	Deep Learning Structure	34
5.3.2	Binarizing Parse Trees	35
5.3.2.1	Overview	35
5.3.2.2	Binarization of Constituency Parse Trees	36
5.3.2.3	Binarization of Dependency Parse Trees	36
5.3.3	Learning to Parse on Embeddings Using Autoencoders	37
5.3.3.1	Overview	37
5.3.3.2	Training Process	38
5.3.4	Tree Matching	38
5.3.4.1	Overview	38
5.3.4.2	Finding Tree Matching	39
5.4	Experiments	39
5.4.1	Testing Embeddings	40
5.4.1.1	Improving POS tagging	40
5.4.2	Testing Autoencoding of Parse Trees on Embeddings	40
5.4.3	Testing Tree Matching as Paraphrase Indicator	41
5.4.4	Testing Dependency vs Constituency Parse Trees	41
6	Experimental Setting and Results	42
6.1	Experimental Setting	42
6.1.1	Embeddings Generation Corpus	42
6.1.2	Autoencoding of Parse Trees Training Corpus	42
6.1.3	Labeled Paraphrase Corpus	43
6.2	Results	43
7	Conclusion	44
	Bibliography	46
A	Proof of Tree Matching NP Completeness	50
B	Paraphrase Tagging Guidelines	53
B.1	Goal	53
B.2	Rules for Tagging	53

B.2.1	Paraphrase	53
B.2.1.1	Restrictions	53
B.2.2	Partial Paraphrase	54
B.2.3	Negative	54
B.3	Examples	54

List of Tables

4.1	Proposed English Semantic Primes[19]	13
5.1	Baseline results for English paraphrase recognition on the MSRP corpus	24
6.1	Size of the Different Corpora used for Embeddings computation	42

List of Figures

4.1	Simple Transformational Grammar Example	12
4.2	Constituent Parse Tree for the Hebrew Sentence “no bears and no forest” . .	17
4.3	Dependency Parse Tree for the Hebrew Sentence “no bears and no forest” . .	18
4.4	Constituency Emphasis Paraphrasing Example	21
4.5	Dependency Emphasis Paraphrasing Example	21
4.6	Constituent Transliteration Paraphrasing Example	21
4.7	Dependency Transliteration Paraphrasing Example	22
4.8	Constituent Participle and Reported Speech Paraphrasing Example	22
4.9	Dependency Participle and Reported Speech Paraphrasing Example	22
5.1	Deep Network Structure	34
5.2	Constituency Binarization Paraphrasing Example	36
5.3	Dependency Binarization Paraphrasing Example	37
5.4	Encoding of a Binary Dependency Tree	39
5.5	Possible Match of a Paraphrase Pair	40

Chapter 1

Introduction

The topic of this research is the processing and generation of natural language paraphrases from an algorithmic point of view. Focusing mainly on the implication of adapting and applying such algorithms to the Hebrew language.

Paraphrasing is the action of restating sentences or paragraphs using different sentence structures or different choice of words. Paraphrases are used by an author of a single text in order to elaborate on a given subject, to provide examples, or to explain a topic by using a different approach while conveying the same message. When comparing text written by different authors describing similar events, one will also find occurrences of paraphrases, reflecting style differences between the authors.

Paraphrase identification is the task of determining whether two given text consist of a paraphrase.

In this research we deal with instances of both tasks for which it is trivial and natural for a native speaker to either either paraphrase a given sentence or to determine if a given pair of text fragments are a paraphrase.

For example, the following two news snippets (dated from 14.11.11) from two different News web sites are considered paraphrases:

- אזרח ומאבטח השתלטו על שודד בסניף בנק הדואר המרכזי בב"ש (מעריב)
- A citizen and a guard took over a robber at the general post office in the city of Beer-Sheva. (here appears in abbreviation)
- אזרח ומאבטח השתלטו על אדם שרצה לשדוד סניף דואר בבאר שבע (ידיעות אחרונות)
- A citizen and a guard took over a man who tried to rob the post office in the city of Beer-Sheva.

1.1 Formal Definitions

Recent research in the field of paraphrasing in English had yielded the formulation these definitions [2][4]:

1.1.1 Textual Entailment

Given two text fragments A, B , determine if one could be inferred from the other. A will entail B if a human being who trusts A , on all its parts, will therefore have to trust that B is also true.

For example:

- דורון קטש הגיע בשנת 2000 לקבוצת פנאתנייקוס וזכה יחד עימה באליפות אירופה
- During 2000 Doron Katash joined Panathinaikos team, and won the European championship with it
- קבוצת פנאתנייקוס זכתה באליפות אירופה בשנת 2000
- Panathinaikos team won the European championship in 2000.

1.1.2 Paraphrase

Commonly defined as bidirectional textual entailment. E.g., given two text fragments (A, B) , it could be said the A entails B and vice-versa. A few restrictions upon these texts have been defined:

- This bidirectional entailment can rely also on knowledge which is considered to be common:

- סין הגיעה להסכם סחר עם רוסיה
- China has reached a trading agreement with Russia
- המדינה המאוכלסת בעולם הגיעה להסכם סחר עם רוסיה
- The world most populated nation has reached a trading agreement with Russia
- Yet the entailment must not rely solely on prior knowledge:

- סין הגיעה להסכם סחר עם רוסיה
- China has reached a trading agreement with Russia
- סין היא המדינה המאוכלסת בעולם
- China is the most populated nation.

For the comprehensive definition of textual entailment, see [4].

1.1.3 Related Shared Tasks

Following these definitions several shared tasks had been targeted, amongst them:

- *Novelty detection* [4], Given a fraction T of text and a corpus C - determine if the T contains new information with respect to C .
- *Knowledge base population* [3], Given a list of entities and a corpus C extract values for a pre-defined set of attributes (a.k.a. “slots”) for target those entities.

1.1.4 Paraphrasing in Hebrew

No formal definitions for paraphrasing in Hebrew were given in previous research. Due to its agglutinated nature and single word abmigation, it seems reasonable to expect the task of paraphrasing to be an interesting and challenging one in the Hebrew language as well.

1.2 Motivation

The tasks of generation and identification of paraphrases help in various NLP fields, such as:

- Automatic text generation , data from a knowledge base can be expressed in a variety of forms using paraphrases. The capability of producing text variability is important for text generators to adapt the generated text to various target audience needs, by varying the style of the text, its complexity and its level of detail.
- Automatic summarization: while scanning through a document, paraphrases found in text body could be detected, and then omitted, in order to provide a shorter version of the document. This is particularly important in the context of multi-document summarization: in this task, several similar documents (describing the same events) are taken as input, and we expect them to contain many paraphrases. The summary of the cluster of documents should avoid repetitions and contain only one instance of each group of paraphrases from the source documents.
- Automatic construction of thesaurus: Identifying paraphrases from freely occurring text in conjunction with exploiting knowledge of the sentence structure can be used to yield a bank of Hebrew words which are, with high probability, synonyms. The result of such analysis would lead to automatic construction of a thesaurus in Hebrew, similar to the Wordnet [16] resource available in English.

Chapter 2

Research Questions and Objectives

2.1 Research Questions

As part of this work, we address the research questions related to analysis of paraphrasing in Hebrew:

2.1.1 Hebrew Specific properties

Are there specific properties of the Hebrew language that allow paraphrasing:

- **Morphological derivation** (e.g., using participle forms to operate as nouns, or verbs). For example, the following pair of sentences is a case of a participle taking form of a verb to create a paraphrase of a given sentence:
 - כשילדים מתעייפים הם נרדמים בקלות.
 - When children get tired, they will fall asleep easily.
 - ילדים עייפים נרדמים בקלות.
 - Tired children fall asleep easily.
- **Syntactic variations exploiting free-word order in Hebrew.** Sentences in Hebrew may be expressed in different word orderings, as a tool to emphasize different notions within the same occurrence, which is a common target of paraphrasing:
 - אני הכנתי את העוגה.
 - I made the cake.
 - את העוגה הכנתי.
 - The cake, I made it.
- **Lexical replacement.** Replacing a Hebrew word with another derived from another language with transliterations, with another part of speech.

- המכונית נהרסה לחלוטין בתאונה

- The car was completely ruined during the accident.

- המכונית עברה טוטאל-לוס בעקבות התאונה

- Beacuse of the accident, the car is a total loss.

2.1.2 Hebrew Datasets for Paraphrasing

Which datasets can be used to collect and identify a database of paraphrases in Hebrew?

2.1.3 Method Adaptation for Hebrew

Could approaches taken on other languages (especially English) for paraphrasing identification and generation be applied on Hebrew? What aspects of the methods must be adapted to account for specific properties of Hebrew?

- Word agglutination
- High morphological ambiguity
- Free-word order syntax

2.1.4 Learning Hebrew Paraphrasing Rules

How should paraphrasing rules be represented and encoded so that paraphrasing knowledge can be used to achieve other NLP tasks? We specifically investigate how paraphrasing rules can be used for multi-document summarization in Hebrew and quantitatively measure the impact of paraphrasing on the quality of produced summaries. In order to evaluate the summaries ROUGE (a package for automatic evaluation of summaries)[27] is used. This package automatically evaluates the quality of a computer generated summary in a recall-oriented evaluation: it counts the number of common units between the computer summarizer's output and a human generated summarization. It is now being commonly used to asses the quality of such algorithms.

2.1.5 Producing Hebrew Resources for Other NLP Uses

Can the process of learning paraphrases yield resources applicable for other NLP Hebrew tasks? Recent research in the field of NLP tries to reuse training information and share it across common NLP tasks (this is known as "Deep Learning"), can such information be obtain and encoded in such a way to aid other Hebrew NLP tasks?

2.2 Objectives

The objective of this work is to develop computational methods to identify paraphrases and generate paraphrases in Hebrew text. The database of learned paraphrasing rules is then applied to help achieve additional tasks in natural language processing (NLP) such as automatic acquisition of a Hebrew thesaurus and multi-document summarization in Hebrew. Training information gained is encoded in order to enhance the already existing algorithms for other NLP tasks.

Chapter 3

Previous Work

Relatively recent research has explored the possibility to learn paraphrases ([6, 33, 32, 26, 23]). These could be divided into three main categories [26]:

- Identification of a paraphrase.
- Generation of paraphrases
- Extraction of paraphrases from large texts.

We review work done in each of these interconnected fields in the rest of this chapter, given definitions formulated in previous chapters.

3.1 Paraphrase Identification

This task involves receiving a pair of text fragments and determining if the pair consists a paraphrase. [6] have developed an unsupervised paraphrase identification algorithm. Their dataset consists of multiple English translations of foreign books (e.g., books that were translated into English more than once). Their assumption was that different translators will introduce paraphrases when translating the same source text. First they used the alignment method of [17] to align the corresponding translations. They continued by applying an iterative model for extracting paraphrases features from the aligned sentences. The iteration starts with a simple rule, which extracts the surroundings of identical words (namely a number of tokens before and after an identical word) in both sentences. From these surroundings they code a "contextual rule", by exploiting the surrounding words POS. They continue by reapplying the new rules learned in the previous iteration upon their corpus once more, allowing for a predefined number of separating tokens to appear, thus allowing new rules to arise. This process ends when no more rules are discovered during a single iteration. During these iterations their classifier is trained from the examples extracted from the text.

[33] created a system for paraphrase identification composed of two parts: the first is an autoencoder [31], trained to compress sentences taken from the WSJ which were parsed by

the stanford parser. The leaves (representing original words in the sentence, retaining the original word order) were replaced by 100 dimensional number vector embeddings of the words (obtained by language model, [12, 36]). The autoencoder then traverses the tree in a bottom-up manner, at each step saving the encoding of two already computed nodes at thier parent node. Thus, the system outputs a tree in which each node contains an encoding of the subtree which it spans.

The second part of their system (after the autoencoder has been trained), recevies a sentence pair and computes two parsed auto-encoded trees, as explained above. The euclidean distance between each node i of the first sentence, and every node j in the second is then computed to form a similarity matrix. Since this matrix is not of a fixed size (because sentence pairs are not of fixed size), this matrix is then sampled by running a sliding window of fixed, predetermined size, over it and picking the minimum of every window to a “dynamic - pooling” matrix. This matrix, now being of fixed size, can be fed to another classifier trained to distinguish between paraphrases pair pooling matrix and non-paraphrase pair pooling matrix.

3.2 Paraphrase Generation

The task of generating a paraphrase of a given text fragment. The Microsoft NLP team [32] created a system to produce paraphrases of an input English sentence. Their system gathered a large automated training set from news sites, upon which they performed subsequently: sentence alignment, word alignment and phrasal replacement identification. They eventually learned a "phrasal translation database" from this dataset. To create this database they have used methods from the area of statistical machine translation (SMT). Given an input sentence, they create a "generation lattice" which describes all possible paraphrases for each possible phrasal alignment where each label is marked with a probability assigned by the system. The lattice is then exploited to generate traversal paths with ranked by probability.

3.3 Paraphrases Extraction

The task of extracting parphrase from a large given corpus. In the field of paraphrase extraction, [22] is a typical example: they developed a system for extraction of Japanese paraphrases from the web. They begin by scanning the web for what they call a "definition sentence", a sentence which describes a term. This is done by searching for a certain sequence of part of speech tags which their hypothesis claim that a definition sentence should adhere to. All the definitions of the same terms are considered candidate paraphrases. These possible paraphrases are parsed for dependency structures and classified by an SVM to decide which candidate should be declared as a paraphrase. Using this method they have achieved a large collection of 300K paraphrases with estimated precision of 94

3.4 Hebrew Paraphrasing

In comparison to the vast research efforts invested in English paraphrasing, very little work has been done in the field of Hebrew paraphrase. [29] have developed a medium scale Wordnet for Hebrew, consisting of 5300 groups of synonymous lexical items (synsets). The approach they have taken was to form the Wordnet by aligning English and Hebrew expressions, and infer relations from the English available Wordnet onto their created Hebrew Wordnet. They state that this method (called MultiWordNet) is preferable over building the Wordnet from scratch since the Hebrew language is poor on computational linguistic resources. The lack of monolingual dictionaries in Hebrew is given as an example of such resource.

Chapter 4

Linguistic Background

This chapter aims at giving the linguistic point of view of paraphrases.

It starts by surveying several main linguistic schools of thought, trying to explain the paraphrase phenomena through the school's concepts.

The chapter continues by attempting to give reasons for the existence of paraphrases, claiming that true equivalence between different phrases does not exist.

The chapter ends by reviewing two common parse tree structures, the dependency parse tree, and the constituent parse tree, displaying parse trees of paraphrases on both structures.

4.1 Paraphrasing Through the Prism of Different Linguistic Schools

We very briefly review several main linguistic schools of thought which provide different theories with respect to language acquisition, notation, and language formation among its native speakers.

We focus on trying to explain the paraphrasing phenomena under the theories' guidelines.

4.1.1 Functionalism

The functionalist school, whose primary formulation was done by Halliday, sees language as a tool which provides its speakers with a **system of choices** with which to transfer **meaning** between the participants of a conversation [21].

Both highlighted terms' description follows.

Meaning The meaning in the speaker's mind consists of *functional elements*. These items can be, for instance: *agent*, *object*, and *action*. In the simple sentence: "John ate an apple" it is clear that functional analysis will model that John is the *agent* while apple and eating are the *object* and *action*.

Lexical Choices These are categorized into three classes [21] :

Field Who does what to whom ?

Tenor What is the relationship between the speakers ?

Mode What theme and register is taken by the speaker ?

A meaning is thus transformed from *functional elements* in the speaker mind to natural language via deciding on each of the *lexical choices*. This analysis yields a view which sets its focus on the relation between the language and the meaning it set out to deliver. Analyzing text according to this school, is trying to parse from a sentence the *functional elements* it contains.

4.1.1.1 Defining Paraphrases in Functionalist Terms

Paraphrases can be seen as the same (or almost the same) meaning conveyed in two occurrences using different tools of language. Analyzing each of the sentences in a functionalist point of view is discerning the “input” to the linguistical functions. Looking at a paraphrase pair yields that both sentences initiated from the same input.

Pick [30] gives the following paraphrases as example:

- I “Julius Caesar was assassinated by Brutus”
II “Julius Caesar lay dead at the hands of Brutus”
- I “John lent Mary the book”
II “Mary had the book on loan from John”

According to Pick both sentences of each pair consist of the same *functional elements*.

4.1.2 Generativism

This is one of the current dominating linguistic school. Originally formulated by Chomsky ([10]), the theory had undergone many changes and evolutions through the years. At its core, the theory proposes that all human languages must comply to a certain set of rules, and that language can not be learned and spoken without adhering to these rules. These are called the grammar of a language, and are written using two disjoint sets of symbols. Constituent symbols (non-terminals), and the language vocabulary, as terminal symbols. In the spoken language we observe only the terminals. These were *generated* from non-terminal symbols in an hidden mechanism (which Chomsky referred to as an organ called “Language acquisition device”) in the mind of the native speaker, via a set of production rules, specific to a certain language. According to the Generativist approach, human beings are born with an innate notion of **Universal Grammar**, learning a language is thus instantiating this grammar with the language specific rules, and the learning the set of terminals, namely, the words of the language.

$$\begin{array}{lcl}
\mathbf{S} & \rightarrow & \mathbf{NP VP} \\
\mathbf{VP} & \rightarrow & \mathbf{Verb NP} \\
\mathbf{NP} & \rightarrow & \mathbf{ProperNoun}
\end{array}$$

Figure 4.1: Simple Transformational Grammar Example

4.1.2.1 Transformational Grammar

Transformational grammar is an inherent part of the generativist school, this grammar consists of rules to map between a set of symbols, where not all are terminals, onto a different series of symbols. By recursively applying this rule a valid sentence can be obtained (see below for a more formal definition). The input of the transformational grammar is thought of as higher abstract than language (at some time was referred to as *Deep Structure*), the target of these transformations is to bring this high abstract notion onto a spoken language (at some time was referred to as *Surface structure*). Through these rules a speaker of the language determines, though not necessarily consciously, if a given sentence is a valid sentence within that language.

A few examples of such rules for the English language appear at figure 4.1.

4.1.2.2 Defining Paraphrases in Generativist Terms

Paraphrases can be modeled in Generativist concepts by: Same *Deep Structure* reaching an individual at two points in time, he would like to convey this information by means of spoken language. Therefore he needs to transform the high level abstraction into a *Surface Structure*. Assume he faces more than one applicable production rule from the transformational grammar of the language. At one time he decides to use one applicable rule, and at a different time he decides on using another. Thus obtaining two different *Surface Structure* (phrases) conveying the same *Deep Structure* (meaning) - which is a paraphrase pair as we have defined in a previous chapter.

4.1.3 Cognitive

A broad field which emerged in recent time, tries to find relations between cognitive actions and language instances. As put by Evans[15] : "It is concerned with investigating the relationship between human language, the mind and socio-physical experience". We survey a particular field in this school - the Natural Semantic Metalanguage.

Table 4.1: Proposed English Semantic Primes[19]

Lexical Category	Words
<i>Substantives</i>	I, you, someone,
	something/thing, people,body
<i>Relational substantives</i>	kind, part
<i>Determiners</i>	this, the same, other /else
<i>Quantifiers</i>	one, two, much/many, some, all
<i>Evaluators</i>	good, bad
<i>Descriptors</i>	big, small
<i>Mental predicates</i>	think, know, want,
	feel, see, hear
<i>Speech</i>	say, words, true
<i>Action, events,movement, contact</i>	do, happen, move, touch
<i>Location, existence,possession, specification</i>	be (somewhere), there is/exist,
	have, be (someone/something)
<i>Life and death</i>	live, die
<i>Time</i>	when/time, now, before, after,
	a long time, a short time,
	for some time, moment
<i>Space</i>	where/place, here, above, below,
	far, near, side, inside
<i>Logical concepts</i>	not, maybe, can, because, if
<i>Augmentor, intensifier</i>	very, more
<i>Similarity</i>	like/as

4.1.3.1 Natural Semantic Metalanguage (NSM)

Originated from the work of Wierzbicka[37] and Goddard[18], this theory argues that all human languages lexicons can be categorized onto hierarchical levels. Each level can be explained using words in layers “beneath” it. At the bottom level of every language there is the layer of *semantic primitives*. This level consists of words which meaning is clear to every native speaker and cannot be explained using other words, thought of as the primary basic lexicon a speaker carries in his or her mind [38]. Research in the field led to the definition of about 60 such words across various languages, such as Chinese, English, and Polish. Table 4.1 shows a list of proposed English semantic primes formulated by Goddard[19].

4.1.3.1.1 Defining Paraphrases in NSM Terms A core term in the NSM field is the **reductive paraphrase**. These are different expressions which convey the same meaning.

From NSM point of view the reductive paraphrase is thought of as an “explanation”, or “definition”, in that the paraphrase must map a phrase composed of words of higher levels onto words from lower levels. Continuing to paraphrase is supposed to yield a sentence composed only of the semantic primes of the language.

4.2 Reasons for Paraphrasing

In this section we try and follow the reasons why paraphrases are so common when dealing with natural language. The starting point is that true equivalence between different phrases does not exist. This is formulated in Clark’s *Principle of Conventionality* and *Principle of Contrast* ([11][?]):

4.2.1 Principle of Conventionality

In every *language community* there are certain meanings (i.e. referants), for which there are conventional words to use when expressing these meanings.

4.2.2 Principle of Contrast

“Every two forms contrast in meaning”. I.e, when there is a deviation from the conventional word used to express a meaning, a member of the language community will assume the speaker is seeking to convey a (perhaps slightly) different notion than that of the conventional meaning.

Many experiments in the linguistic field were conducted to provide empirical evidence of the validity of these two principles. Among these were studies of the way children acquire language (and seemingly avoid attaching two different words to the same meaning), and studies across different language communities speaking the same language (showing different conventional names for the same meaning). Experiments supporting these principles were also conducted among Hebrew native speakers [14].

If every two forms of words indeed differ in meaning as research suggests, one can induce that two different phrases must also differ in meaning. Thus complete and true meaning preserving paraphrase does not exist.

4.2.3 Small Scale paraphrasing

Following Hirst’s division[23], small scale paraphrasing are two phrases which are identical except a few words exchanged for their synonym. Following the *principle of contrast* this divergence between the otherwise identical phrases serves the speaker intention to give the phrase a different meaning. For example:

I *Yossi bought a carpet at the market*

• יוסי קנה שטיח בשוק

• יוסי קנה מרבד בשוק

- The word “shatiach” in the first sentence is substituted with the word “marvad” in the second. The former being the Hebrew contemporary conventional word for carpet. The deviation in the latter may suggest that Yossi bought a rather exclusive artifact instead of a usual rug.

II *The destruction of the building occurred last year*

• הרס הבניין קרה בשנה שעברה

• חורבן הבניין קרה בשנה שעברה

- The word “herez” in the first sentence is substituted with the word “hurban” in the second. As before, The former being the Hebrew contemporary conventional word for destruction. The deviation in the latter may suggest that the speaker would like to convey a sense of grief related to the destruction, as opposed to the latter, which may sound as a more objective description.

III *The last few days indicate that the fighting is over*

• הימים האחרונים מעידים על סוף הקרבות

• הימים האחרונים מעידים על קץ הקרבות

- The word “sof” in the first sentence is substituted with the word “ketz” in the second. Again, the former being the Hebrew contemporary conventional word for ending. The deviation in the latter may suggest that the speaker believes that the break in the fighting is a more constant state, as oppose to the former, which may be understood as a less definite state.

4.2.4 Large Scale paraphrasing

Large scale paraphrasing, according to Hirst’s division, is the act of different phrasing of full sentences or paragraphs. This instance of paraphrasing can be found when an author, either consciously or unconsciously, inserts his own ideology and perception into the details of an occurrence. For example:

I Emphasis

• אני הכנתי את העוגה.

- I made the cake.

• את העוגה הכנתי.

- The cake, I made it.

The latter emphasizes the product of the baking (the cake), while the former emphasizes the maker (I made the cake).

II Ideology

- בייניש פורשת מנשיאות העליון. (הארץ, 29.2.2012)
- דורית בייניש הולכת הביתה. (כיכר השבת, 28.2.2012)
- Dorit Beinish retires from supreme court. (Haaret'z, 29.2.2012)
- Dorit Beinish goes home. (kikar hashabat, 28.2.2012)

Although the two news snippets above report the same occurrence - the retirement of Dorit Beinish from the supreme court, it is safe to say that the latter phrase takes a more supporting approach to her departure.

4.3 Parse Trees of Paraphrases

Parsing (also known as syntax analysis) is a common stage in the process of understanding a certain text. It build on the language grammar and aims at constructing higher level hierarchies present within an input text - normally, to form a tree representation of the data. This tree structure is then used by next stages of the process. Parsing is a common stage for compilation, which has a clear formal grammar, as well as natural language, for which the underlying grammar is often obscure and not fully defined. There are various definitions and conventions concerning the grammar definition which will affect the construction of the parse tree.

We survey two common approaches, phrase structure grammar (constituent grammar), and dependency grammar. We focus mainly on instances of paraphrases pairs, and review the resulting parse trees in each of the approaches. From here on, when speaking of constituents of parse trees we will refer to a part of sentence which can be seen as single unit in parse tree, i.e. a subtree of a parse tree.

4.3.1 Phrase Structure Grammar

Phrase structure grammar, also known as constituent grammar, is the grammar originally defined by Chomsky as part of the generative school. The formal definition of the phrase structure grammar follows:

A Phrase structure grammar is a 4-tuple $G = (N, T, S, P)$ where

- I $N \cap T = \emptyset$, N - The *non-terminal set*, T - the *terminal set*
- II $S \in N$, S being the *start symbol*
- III $P = \{(u, v) : u, v \in (N \cup T)^*\}$, P is finite and called the *production rules*

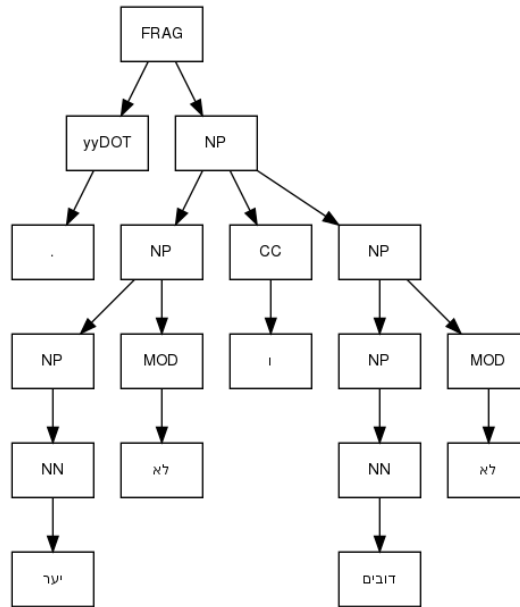


Figure 4.2: Constituent Parse Tree for the Hebrew Sentence “no bears and no forest”

According to these grammars, the leaves (*terminals*) of the parse tree are the words of the original sentence, appearing in the original sentence order. I.e, the leftmost word will be represented by the leftmost leaf, and so forth. The rules by which an input sentence is parsed onto a parse tree are defined via a transformational grammar like the one defined in figure 4.1.

Computational tools which parse input sentences onto a possible constituent parse tree were devised with good results for Hebrew and English languages ([9],[20], respectively). Figure 4.2 shows a constituent parse tree of the Hebrew sentence “No bears and no forest”.

4.3.2 Dependency Grammar

Dependency grammars in modern linguistics date to the work of Tesni‘re (1959) [28]. This approach towards parsing view the syntax analysis of a sentence as consisting of binary asymmetrical relations between words. In each of these relations there is a word which functions as a *head* (also known as governor, regent) and a second word which functions as a *dependent* (sometimes appear as modifier). The relation between the two words is called a *dependency*. According to this linguistic theory, the speaker of a language analyzes syntax by perceiving connections between words, the dependency relation aims at modeling this connection. This binary relation of head - dependent will appear in the parse tree as father - son relation.

This definition suggests that dependency parse trees will be formed of words in the sentence. As oppose to the constituency grammar trees, words in the dependency tree will appear also as the inner nodes, and not only as terminals. Because of this property, dependency trees will be of smaller size than that of the matching constituency tree, in which syntactic nodes

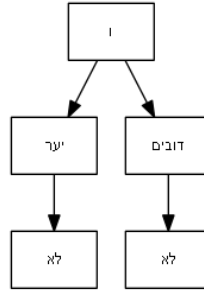


Figure 4.3: Dependency Parse Tree for the Hebrew Sentence “no bears and no forest”

are present in addition to the original words of the sentence. As an example of this, the dependency parse tree for the sentence “no bears and no forest”, which was shown previously for constituency tree, is shown in figure 4.3 (18 nodes in constituency tree versus 5 nodes in the dependency tree) .

The root of every subtree will represent the head (and will normally be a verb) of that subtree (by transitivity), where his direct dependents will be his sons.

Various definitions exist for defining when two words will appear in a dependency relation, yet, as is often the case with NLP, none of the rules covers all cases. Following are a few of these definitions [28] (where: **H** marks the head, **D** marks the dependent, and **C** marks their relation) :

- **H** determines the syntactic category of **C** and can often replace **C**.
- **H** determines the semantic category of **C**, **D** gives semantic specification.
- **H** is obligatory, **D** may be optional.
- The form of **D** depends on **H**.
- The linearize position of **D** is specified with reference to **H**.

These should be seen as guidelines for identifying words in dependency relation: Failure to satisfy one of these rules does not necessarily indicate the pair of words is not in a dependency relation. Accordingly, success to satisfy a rule does not necessarily deem the pair as being in a dependency relation.

As for Constituency parsing, tools to automatically parse a dependency structure of a sentence were devised ([25],[20]).

4.3.3 Examples of Paraphrases Pair in Phrase Structure Grammar vs Dependency Grammar

In the following section we review paraphrase pairs’ parse trees in both formats defined above (i.e, constituency and dependency).

1. A simple example of paraphrase pair was described before, as a pair which demonstrate the change of emphasis:

• אני הכנתי את העוגה.

- I made the cake.

• את העוגה הכנתי.

- The cake, I made it.

Figure 4.4 and 4.5 show the corresponding constituent and dependency parse trees of both sentences. In fig. 4.4 we can see the same constituents appearing as subtrees of both sentences, while in fig. 4.5 we can see that “the cake” in the latter sentence receives a higher node in the parse tree, when comparing to the former sentence. This may correspond to its increased emphasis in the second sentence.

2. Recall the example of Hebrew paraphrasing arising from replacing a term with its transliterated equivalent:

• המכונית נהרסה לחלוטין בתאונה

- The car was completely ruined during the accident.

• המכונית עברה טוטאל-לוס בעקבות התאונה

- Because of the accident, the car is a total loss.

Figure 4.6 and 4.7 show the corresponding constituent and dependency parse trees of both sentences.

In fig. 4.6 we can see that the general structure is similar within the trees (with respect to the inner nodes structure and type). There is a change in one constituent $S \Rightarrow PP \Rightarrow PP$ on the right is substituted with $S \Rightarrow ADVP \Rightarrow RB$ on the left. This can be seen as a substitution of the Hebrew “lahalutin” (completely) with the transliterated term “total loss”.

In fig. 4.7 The governor of the sentence is changed from “neherza” (wrecked), on the left, to “avra” (suffered, in this context) on the right. Apart from that we see, again, the subtree of “lahalutin” (completely) on the left, changed with “total loss” on the right.

3. In this example we examine a relatively complicated example of paraphrasing, which include participles changing part of speech between the two sentences and reported speech.

• שופטים קבעו היום כי החשוד אשם, למרות טענותיו כי הוא זכאי

- Judges have decided today that the defendant is guilty, although he claimed being not guilty.

- החשוד שהואשם היום אמר כי שופטים אותו למרות שהוא זכאי

- The defendant that was accused today said that he was not guilty.

Figure 4.8 and 4.9 show the corresponding constituent and dependency parse trees of both sentences.

In fig. 4.8 a few similarities can be found, although not as easy as in previous examples. The first level below the root is almost identical, and the prepositional phrase is also similar.

In fig. 4.9 it is hard to find any similarities between the parse trees, nor to map between matching constituents.

In both figures it can be seen that the participle “shoftim” (judge) indeed appear in different nodes in the parse tree.

4.3.4 Discussion

As a conclusion of this section, we notice a few relevant observations with regards to both formats of parse tree which were shown, while keeping in mind the task of paraphrase identification.

- Dependency parsing yields a tree with significantly less nodes when comparing to constituency parsing.
- A change in a dependency parse tree due to paraphrasing can inflict a change in higher nodes than that of a constituency parse tree.
- Dependency parse tree seem more sensitive to changes in the sentence structure.
- A note on binarization of parse trees: While the binarization of a constituency parse tree is rather straight forward, this is not the case when binarizing dependency parse trees. This topic is elaborated in 5.3.2.

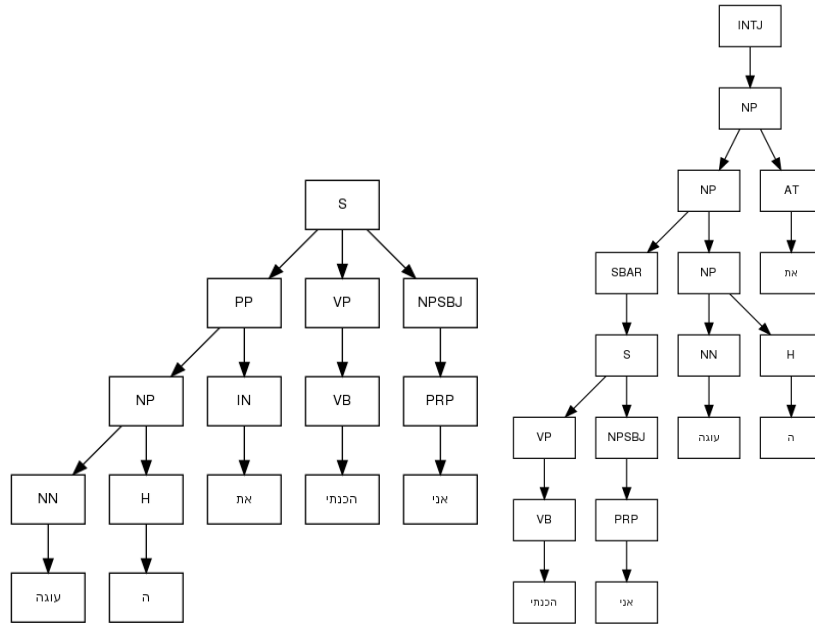


Figure 4.4: Constituency Emphasis Paraphrasing Example

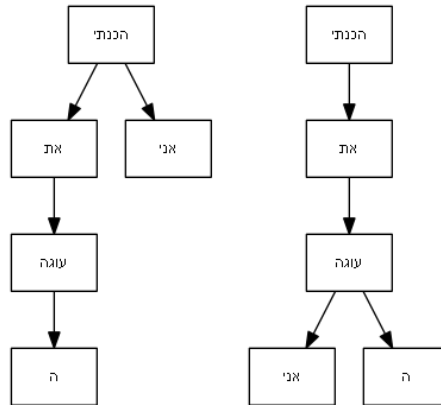


Figure 4.5: Dependency Emphasis Paraphrasing Example

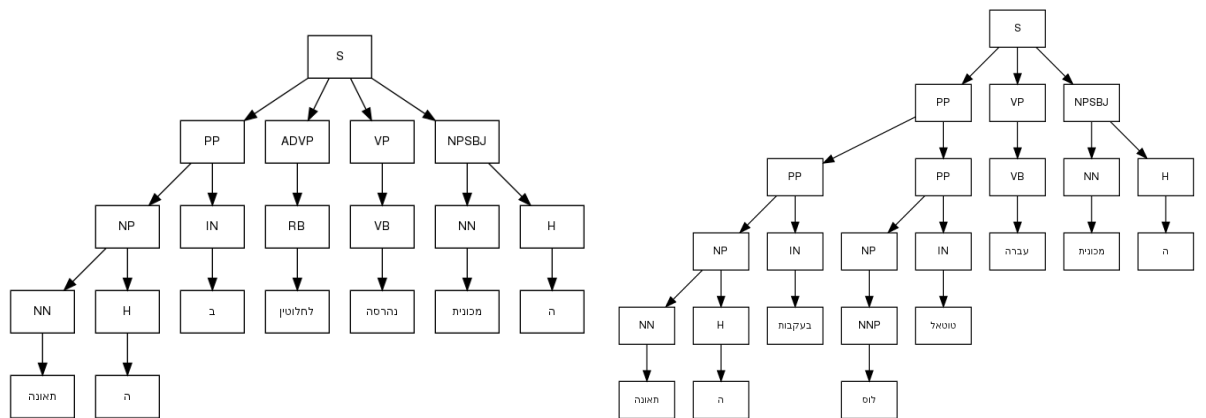


Figure 4.6: Constituent Transliteration Paraphrasing Example

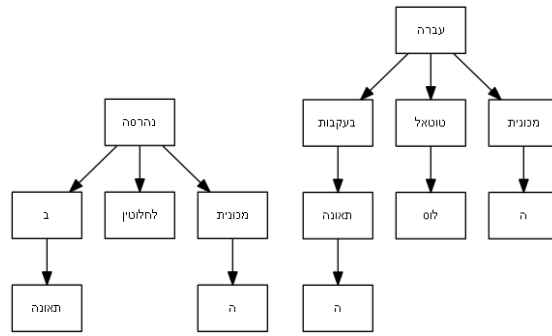


Figure 4.7: Dependency Transliteration Paraphrasing Example

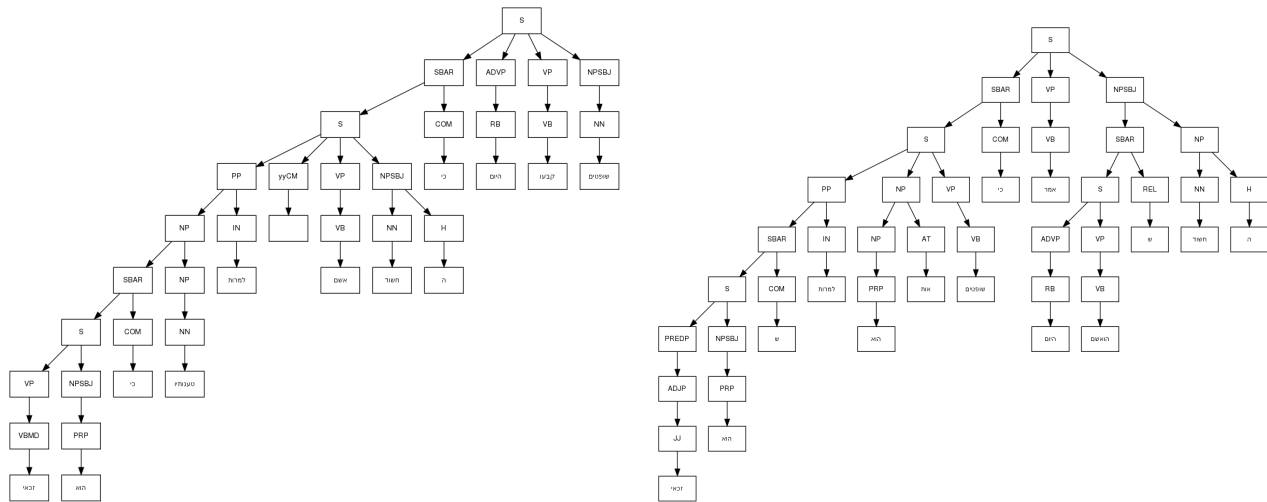


Figure 4.8: Constituent Participle and Reported Speech Paraphrasing Example

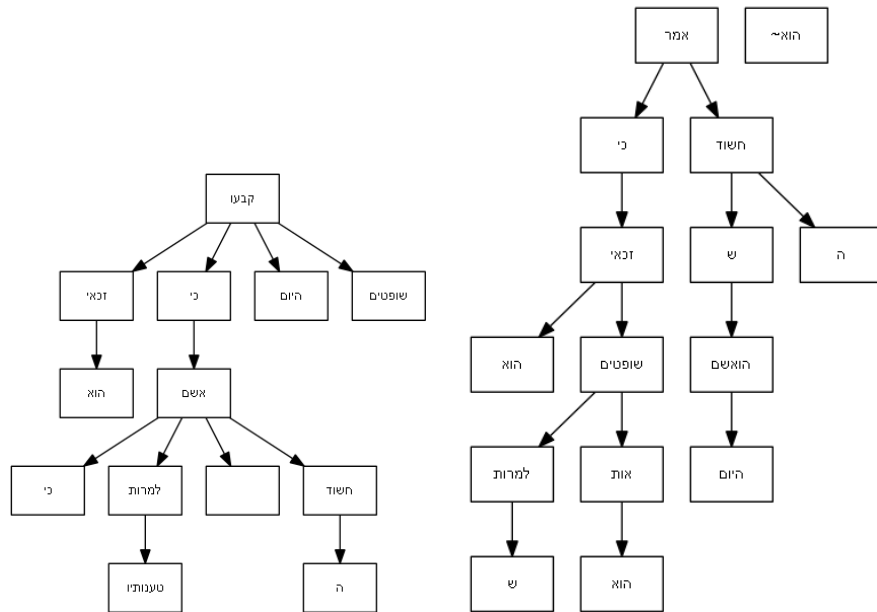


Figure 4.9: Dependency Participle and Reported Speech Paraphrasing Example

Chapter 5

Methodology

This chapter covers the main methods of learning which are applied in this work, description of the adaption technique and proposed hypotheses, and finally the experiments taken to prove those hypotheses.

The structure of this chapter is as follows. In section 5.1 we describe the method and outline of the work in general lines. Following this, section 5.2 reviews the methods of learning applied in NLP which are relevant to task of paraphrasing. Section 5.3 contains an in depth plan for adapting the algorithms discussed to answer the research questions posed in previous chapter er. Section 5.4 concludes this chapter by depicting the experiments by which we test correctness of the methods proposed.

5.1 Plan of Work

The target of this work is to face the task of paraphrase recognition in Hebrew. A task which was under great focus of research in recent years, and contributes to many NLP tasks, as was shown in previous chapters.

To attack this problem, we apply, adapt, and devise methods of machine learning (neural networks, autoencoders), and computational linguistic (language models, automatic parsing).

Since this task was yet to be faced in the Hebrew language, it is missing comparative benchmark results and adequate corpora. Thus, in order to properly attack the task in Hebrew we collect and publish such corpus. Datasets consist of a monolingual Hebrew parallel corpus of paraphrases. These are acquired in an unsupervised manner by obtaining hourly news flashes and corresponding news stories from leading Israeli News sites. Each news flash in the dataset is aligned against an assumed matching news flash from the other websites. This yields a method for obtaining an automatically growing dataset. These datasets will be hand tagged by human judges in order to obtain a baseline comparable corpus, guidelines for formally defining what is considered as an Hebrew paraphrase will also be published as part

Model	ACC	F1
Wan et al. (2006)	75.6	83.0
Das and Smith (2009)	76.1	82.7
Socher et al. (2011)	76.8	83.6

Table 5.1: Baseline results for English paraphrase recognition on the MSRP corpus

of this work (see appendix B). The methods and hypotheses raised in this work are trained and tested against this corpus. The desired size of the corpus is 6000 sentences, with 67% of the pairs being in a paraphrase relationship. These numbers are an Hebrew equivalent of the MicroSoft Research Paraphrase Corpus (MSRP corpus [1]). Recent research in the field of English paraphrase uses this corpus as baseline comparative parallel corpus.

In order to recognize Hebrew paraphrases, we go along the lines of [33] that used auto encoders on embedded parse trees of the input sentences. They have used word feature vector representation of the words, as obtained from a deep learning approach. In the course of this work we reproduce the settings described in that paper for the English language. We offer some deviation and improvement of the algorithms they have used, introducing a new search problem on the pair of trees. Following the reproduction of that experiment we implement it for the Hebrew language. This implementation requires generating resources not yet available for the Hebrew language, such as Hebrew word embeddings, binarization of dependency parse trees, and the parallel paraphrase corpus already mentioned above.

We test both Socher’s implementation’s equivalent , as well as the proposed improvement of it, on the Hebrew corpus. The main experiment is the task of identifying **paraphrase identification** as described in previous section. The baseline results of several English applications on the MSRP corpus appear in table 5.1.

5.2 Learning Methods

Learning methods applied in Artificial Intelligence (AI) in general are used to detect some recurring pattern in the input data. If the process of learning is **supervised** then the pattern is first *learned* from labeled training data (by adjusting the specific parameters of the learning model). Then the model is tested on an unlabeled, and previously unseen test data, to measure its ability to infer and deduce from the training phase. We hope that by properly modeling the input data and correctly adjusting the model parameters during training we will learn the real underlying pattern embedded in the data. Pattern which separates correctly between the samples.

In Natural Language Processing (NLP) in particular these methods are employed to natural language corpora. Training data may be labeled for the specific tasks - i.e pairs of sentences are marked a-priori as being paraphrase or not. Additional labels on input data may include semantic information, which is naturally not present in raw natural language. These may include Part of Speech (POS) tags, parsing of the sentences, anaphora resolution, to name a few. Over this data and potentially labels, learning methods are used to explore recurring patterns in the input text.

Next we review the learning methods applied in this work to identify recurring patterns in Hebrew paraphrases.

5.2.1 Neural Networks

5.2.1.1 Outline

The concept of neural networks is to try and imitate the function of neurons in an intelligent being's brain. Each neuron is modeled as a predetermined differentiable function which receives a predetermined number of inputs. Neurons are connected to other neurons, forming what can be seen as graph whose vertices are the neurons, and whose edges are the connections of those neurons. This graph is also determined in an a-priori manner. A neuron network is the set of networks and their connections. A special case of neurons are the *input neurons* - which receive their input from an input instance of predetermined format, and *output neurons* which output the result of the total network. Most of the information in this section is based on the book [34].

5.2.1.2 Neuron

The function of a specific neuron can be seen as composition of two functions:

- Integration function: $I : \mathbb{R}^n \rightarrow \mathbb{R}$ - receives as input the incoming edges of this neuron and outputs a single output.
- Activation function: $A : \mathbb{R} \rightarrow \mathbb{R}$ - receives the output of the integration function and

returns the final output of this neuron.

Thus, the function of this neuron can be seen as $A \circ I$.

5.2.1.3 Weights

As we've seen, the output of a neuron is in turn part of the input of another neuron. The two of the neurons are connected by an edge in the network graph on which the information “travels”. As part of this model, the graph is seen as a **weighted** graph. Each of the edges e has a weight $w_e \in \mathbb{R}$, the information x traveling an edge e is multiplied by w_e to reach the integration function as the input $x \cdot w_e$.

5.2.1.4 Network Computation

Neural network computation is the process of giving an input instance to the input neurons, continuing by computing at each step the neuron function, and transferring their output to their connected neighbors. The process is terminated when the output neurons output a value. This vector of values is the output of the network and the result of its computation on the input.

5.2.1.5 Learning

Most of the parameters of the network are predetermined, decided upon during design time, and remain constant throughout all iterations of network computations. These parameters are the neural network graph, including the number of neurons and their connections, and the neuron functions (both integration and activation). The learning is achieved by modifying the previously mentioned weights of the edges, these change during a number of iterations which is called “training” in which tagged input (i.e, that the desired output for this input is known), is fed into the system. Knowing the desired outcome can yield an error value for this specific computation, by means of some distance metric. Thus, after a number of these iterations, the value of an error function is known at specific points of input space. This error function can be seen as function of the weighted edges existing in the system, as these are the only values which can be modified according to the neural network model. Thus by summing the received errors of all inputs we can arrive at an error value for these weights. changing the network weights and recalculating the error yields an error function in weight space whose values are known at specific points. This error function can reach a minimum using numerical methods. Using for instance, the discrete gradient of the network function.

5.2.1.6 Backpropagation

We have seen in the previous section that learning in a neural network can be reduced to calculating partial derivatives of the error function over weight space, in order to perform gradient descent on it, and minimize the error on the training dataset. A common way to

calculate these partial derivatives is *Backpropagation*. This method exploits the chain rule that applies in calculating derivatives, and stores the derivative of its input during network computation. After the computation step begins a backpropagation step which starts from the output nodes, traveling the network backward, multiplying at each neuron the saved value of the derivative. Thus, the derivative of the network is obtained at the input neurons.

5.2.1.7 Layered Networks

A special case of neural networks, and one of common use in machine learning is the case of n layered networks. Several extra restrictions are applied on these networks:

- The Neuron vertices V of the graph can be seen as a union of disjoint sets: $V = V_1 \cup \dots \cup V_n$, $V_i \cap V_j = \emptyset$. In layered networks the inner layers: V_2, \dots, V_{n-1} are referred to as “hidden layers”, and their respective neurons are referred to as “hidden units”.
- Edges between neurons can be set only between nodes of adjacent groups, i.e, the edge set E of the graph can be seen as $E = E_1 \cup \dots \cup E_{n-1}$, $E_i \cap E_j = \emptyset$, such that $E_i = \{(v_i, v_{i+1}) | v_i \in V_i, v_{i+1} \in V_{i+1}\}$
- Normally all neurons in a layer are connected to all neurons of adjacent layers.

Special optimizations for calculating the backpropagation step on layered networks exists, and we use this types of networks throughout this work.

5.2.2 Deep Learning

Recent research in AI has yielded the approach of deep learning [7]. This approach aims at modeling the human perception of complicated notions in several levels of representation. This approach is implemented using several neural networks connected in such a way that one network’s outputs it transferred to another’s input. The backpropagation is carried across networks, starting from the topmost network to the bottom. This process can give flexibility in determining the networks parameters according to their location and task.

5.2.2.1 Curriculum Learning

As part of the model of deep learning, the concept of curriculum learning ([8]) was established. This school claims that by giving the deep learning networks inputs in an increasing order of complexity (as one teaches a child), the model will be able to obtain better results. This is due to establishing firmer representations on the “simple” samples, and improving that solid ground when facing higher degrees of complexity in the input.

5.2.2.2 Multi Task Learning

It is common in recent research to attack several tasks at once, when using deep learning. [12, 13] created an architecture where the lower levels of the deep learning extract features of the input text, and embeds them in a multidimensional vectors - these vectors are the output of the bottom networks. These bottom networks are shared across the tasks which they confronted (POS tagging, Semantic role labeling, Chunking, Name entity recognition, and language model), and requested the same features. The upper levels, which are task specific classifiers take as input these embeddings instead of the original text. Thus the lower levels change during training (which is done in an interleaved manner) of all these tasks. The rationale is that information learned for one task for representing feature embedding, may be useful for a second task which classifies according to the same feature - thus the tasks classifiers “profit” from the multi training.

5.2.2.3 Word Embeddings

As a by-product of this process - The embeddings ([12]) of words were published for English dictionaries [36] to use as a plug-in enhancer for use in all NLP tasks which treated words simply as index to a finite dictionary.

As part of this work we intend to use a deep learning architecture to learn an Hebrew language model from a large corpora of unlabeled text. This will yield a dictionary of Hebrew word embeddings which will aid the process of paraphrase identification. This dictionary can also be published to aid other research in the field of Hebrew NLP. The language model and embedding networks will be presented with dictionaries of growing size (words are ordered from the common to the less common), following the curriculum learning concept.

5.2.3 Auto Encoders

A major drawback on the model of neural networks as presented above is the constant size of the network, on all its parts. These must be defined during the design of the system. This constant size of input and inner representations (known as the *connectionist approach*) does not fit the nature of most AI tasks, for which the instances of a task are not of constant bounded size. NLP tasks possess this property as well. The instances of natural texts are of variable size and unbounded length - no known bounds exist for word length, nor for sentence length, nor paragraph length. In the field of this work, it is easy to observe from the examples given in previous chapters that paraphrases must not be of same length, although by definition they convey the same amount of information. This property seems to limit, or cancel the use of neural networks in the field of arbitrary length instances, such as NLP.

To attack this problem, several connectionist system were devised to cope with the unbounded input size. One of these is the *Recursive Auto-Associative Memory* (RAAM, due to [31]), also know and extended in later papers as *Auto Encoders*.

5.2.3.1 Outline

Auto Encoders aim at creating recursive fixed size representations of variable size input. The input is assumed to be a variable length list of fixed size elements (a character string for example), mark each element representation size as K . They do so by going over the input instance of size $K \cdot n$ and “encoding” every pair of consecutive input elements onto one compressed representation. Thus obtaining a second level of representation of size $K \cdot \frac{n}{2}$, this process is repeated until the last level contains one element of fixed size which represent the entire variable size input, in a fixed size representation.

5.2.3.2 Neural Network

The system Pollack devised was composed of a two concatenated neural networks:

- Encoder - a neural network composed of an input layer of $2K$ elements fully connected to an output layer of K elements.
- Decoder - a neural network composed of an input layer of K elements fully connected to an output layer of $2K$ elements.

These two networks are trained concurrently as a single network: a network of $2K$ input elements, K hidden units in one hidden layer, and $2K$ elements in the output layer. During training on an input element, the element itself is presented as the network desired outcome - thus training the network to output the input it received, after going through the hidden layer. The hidden layer (recall that this is actually the output layer of the encoder), now holds a compact representation of the input.

5.2.3.3 Encoding

Encoding of a variable size length input is giving to every two consecutive elements (concatenated to get $2K$ elements) to the encoder network, thus producing an element of size K , this is marked as the father of the two nodes which it encodes. This process repeats until we are left with one element of size K . This process yields a **binary** tree of the elements, in which the leaves are the original input and the inner nodes are representing some subtree they encode. If the encoding phase is flawless one can reconstruct the entire input sentence given only the root of this binary tree.

5.2.3.4 Decoding

The process of decoding is to take the root of the previously mentioned binary tree, and pass its K elements to the decoder network to receive $2K$ elements (the decoded representations of its sons). This process is repeated on the sons until the full tree is obtained again, and the input sentence is reconstructed.

5.2.3.5 Uses in NLP

Applications for such a system in the NLP field consist of taking a natural text structure, such as a parse tree, and encoding it. Later stages of the algorithms can refer to the entire structure, or parts of it, as a fixed size input.

[33, 35] have used autoencoders in two NLP applications:

- Predicting sentiment distributions ([35]) : Input sentences were encoded, and during training the inner representations of the tree were learned, as indicators for sentiment classification of the sentences.
- Identifying paraphrases ([33]) : Each sentence of the given input sentence were encoded. Afterwards, the inner representations were compared and euclidean distance was computed, to train as indicators of paraphrasing.

In this work we train an Auto Encoder for Hebrew parsed sentences, and try to compare inner representations as indicators for paraphrasing. The words in the leaves are replaced by their multi-dimensional vector embeddings.

5.2.4 Tree Matching

Tree matching is a new concept devised for better exploitation of the usage of autoencoder, compared with the approach taken in (Socher et al, 2011).

5.2.4.1 Definition

Given two binary trees - t_1, t_2 (corresponding to sentences s_1, s_2 , and obtained from them by auto encoding) in which the leaves contain word embeddings and the inner nodes contain encoding of the subtree they span. Define a “Tree Match” M , to be a set of tuples (n_1, n_2) , where n_1, n_2 are nodes of t_1, t_2 accordingly, s.t. for every word w in $s_1(s_2)$, M contains exactly one tuple which contains a node in the path from w to the root of t_1 (t_2). For instance the tuple which contains the root of both sentences is always a match.

This definition captures the idea that a paraphrase pair consist of sentences which parts are interchangeable, from the sentence level down to the word level (including word-reordering).

5.2.4.2 Tree Match Score

Following this definition, a score of a match can be defined:

$$S(M) = \sum_{(n_1, n_2) \in M} (||n_1, n_2|| \cdot (\text{number of spanned leaves by } n_1 \text{ and } n_2))$$

During training, a simple classifier can learn the threshold below which minimal matches represents pairs which are paraphrases. After training, the classifier would hopefully yield not only a binary value, but also significant matching between the pair, "explaining" why they are a paraphrase.

5.2.4.3 Tree Matching is NP-Complete

It can be shown that Tree Matching, as defined above, is an NP-Complete problem by showing:

$$\text{Positive SubsetSum} \leq_p \text{Tree Matching}$$

For proof, see Appendix A.

5.3 Hebrew Paraphrasing

This section describes the specific algorithms, modifications and adaptations, taken in the course of this work to answer the research questions. The general approach for paraphrase identification is as follows:

The rest of this section explains each of these in detail. The formation of the chapter is that every section's output is the next section's input. Every section start with a "black box" examination of the currently described logic, and its part in the overall process.

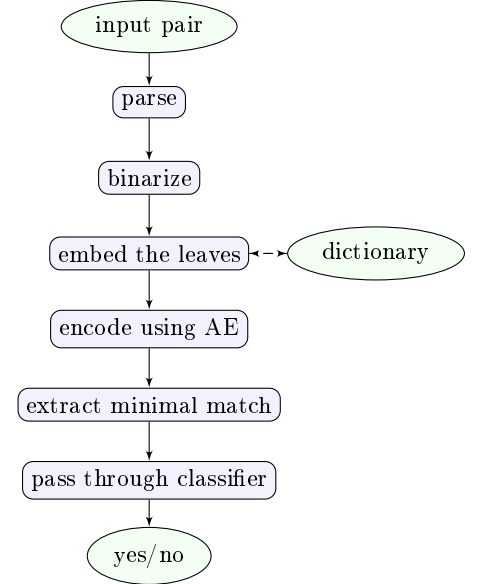
5.3.1 Producing Deep Learning Embedding for Hebrew

5.3.1.1 Overview

As described above, word embeddings are mapping from a finite word dictionary D onto a d - dimensional feature vector space. Each word $w \in D$ is mapped onto a vector $r_w \in \mathbb{R}^d$. This mapping is produced as a by product of deep learning architecture, and can be seen as a goal on its own when wanting to provide a plug-in enhancer for common NLP tasks. As a part of this work we intend to provide such a resource for the Hebrew language by learning a language model over a large Hebrew corpus.

The language model on its own won't be used for further tasks, but the embeddings extracted from the architecture will be published, and their ability to enhance NLP tasks will be tested on a common NLP task (POS tagging) in the course of this work.

- Given an input pair (s_1, s_2) of sentences in Hebrew.
- Parse them for dependency / constituency to receive two trees (t'_1, t'_2) .
- Binarize these trees to obtain two binary parse trees (t_1, t_2) . **(explained in section 5.3.2)**
- Change the leaves of the trees to contain Collobert & Weston language model embeddings. **(explained in section 5.3.1)**
- Use an auto encoder to receive compact representations at each of the inner nodes. **(explained in section 5.3.3)**
- Search for a minimal Tree Matching on those two trees, and use its value to decide if the pair is a paraphrase. **(explained in section 5.3.4)**



5.3.1.1.1 Input Dictionaries are pre constructed in a curriculum learning approach. These are used to filter an n -gram corpus, by selecting from the corpus only n -grams which consist of words from the dictionary. These filtered n - grams are the input for this stage.

5.3.1.1.2 Output For each word in the dictionary an embedding is produced onto a 100-dimensional hyperplane. This embedding encompass an encoding of richer features other than just index into a dictionary, as is usually the case in NLP applications. These embeddings will serve the next stages by replacing each word by its embedding representation.

5.3.1.2 Preprocessing and Hebrew Adaptation

The corpus is pretagged with [5] for POS, and tokenization. The preprocessing steps taken on this corpus are:

- We take into account only the segmentation of the words - the following Hebrew prefixes appear on their own as a different word:

ב,כ,ל,מ,ה,ו,ש

In, As, To, From, The, And, That

- Number tokens appear as NUMBER

After applying this segmentation on the corpus, it holds 131M tokens.

5.3.1.3 Language Model

Language models try to give a score to a phrase in a specific language. This score represents the model's estimation of the phrase being a valid phrase in this language. We build a probabilistic language model based on a large corpus of 131 million words obtained from news wire information. The language model is built using a neural network, which is trained to separate between positive and negative examples of n -grams. The correct n - grams are taken directly from the corpus. The negative examples (termed "corrupt" n - grams) are created artificially, by introducing noise into positive examples. The requirement of the language model is that it can separate between positive and negative examples by a difference score of at least 1.

Therefore, we would like to minimize the following expression:

$$\sum_{x \in S} \sum_{w \in D} \max(0, 1 - f(s) + f(s^w))$$

S being the pool of all available overlapping n - grams which are composed only of words in the current dictionary D . f is a network computation function, and s^w is the correct n - gram, where we corrupt the last word of it with the word w . Minimizing this aims at ensuring a difference of 1 between the score of a valid n - gram, and that of all "corrupt" ones.

5.3.1.4 Curriculum Learning

Following the curriculum learning doctrine ([8]), we want to present the language model learning environment with samples in an increasing order of complexity. The complexity of the samples in the language model learning context, is in terms of the probability of the appearance of a word in the corpus. E.g., a more common word is considered as a more "simple" example, whereas less common word is considered as "complicated" example.

Thus the learning is carried out in iterations:

- In the first iteration, start from random word embeddings in the range of $[-0.01, 0.01]$ and train only on n - grams which contain words from a dictionary of the 5000 most common words.
- In the next iterations, initialize the embeddings as those achieved in the previous step, and increase the dictionary size to (10K, 30K, 50K and finally 100K)

5.3.1.5 Stochastic Sampling of the Error Function

The calculation of an error function of the model, as described above, is a computationally heavy task, as it involves iterating over all the dictionary on every n - gram update. This fact will make training an unfeasible task when reaching large dictionary sizes. In order to make the task possible, even on large sizes of the dictionary, we instead sample the error

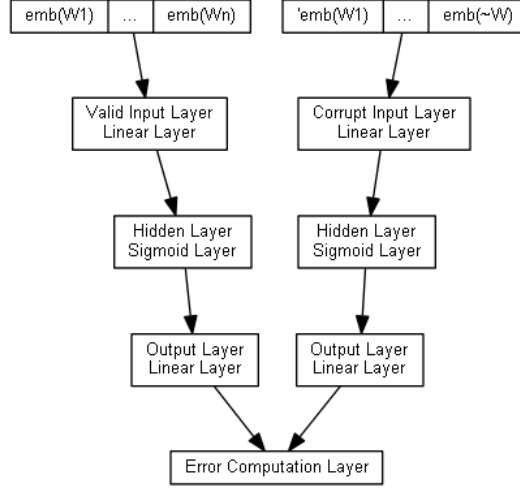


Figure 5.1: Deep Network Structure

function using only one corrupt n - gram at each update. Thus stochastically “peeking” into the model’s behavior instead of directly calculating it (this approach is taken by [12]). The error function for a single iteration on the n - gram s is:

$$\max(0, 1 - f(s) + f(s^w)) \quad (5.1)$$

Where w is a word chosen uniformly at random from the dictionary.

This term is then backpropogated through the network to alter its parameters accordingly.

5.3.1.6 Deep Learning Structure

This section describes in greater detail the neural network structure devised in order to achieve a model which is capable of computing the language model score according to the requirements described above. As a by product of the language model training, d - dimensional Hebrew word embeddings will be created in the deep structure hidden layers. Figure 4.1 graphically shows the structure described below.

- **Embedding networks** - These are $|D|$ distinct neural networks with input and output layers of size d , and no hidden layers. Each of which can be seen as a matrix $M \in \mathbb{R}^{d \times d}$ of the network’s weight.
- **Embedding Dictionary** - This will be the output of the embedding network computation - the mapping between words onto the feature vectors:
 - Its keys - are the words w of the current dictionary D
 - Its values - are $(v_w \in \mathbb{R}^d, M_w \in \mathbb{R}^{d \times d})$. These represents an initial word embedding v_w and a matrix of weights - this is adjusted during the embedding networks training phase.

To obtain the embedding of w , one can perform $v_w \cdot M_w$. We will mark the embedding of a word w as $emb(w)$ and the embedding of an n gram $G = (w_1, \dots, w_n)$ as $emb(G) = (emb(w_1), \dots, emb(w_n))$

- **Language model network** - A neural network with:
 - $d \cdot n$ input neurons (recall that we are looking at n -grams from the corpus, and would like to obtain an embedding of d features for each word in the dictionary). This will be used to input the network the current n -gram representation.
 - An hidden sigmoid layer of 100 hidden units.
 - An output layer of 1 neuron - this will be used to calculate the score for this n -gram.
- **Complete deep network** The complete deep learning network computation is as follows:
 - **Input:** $G = (w_1, \dots, w_n)$ - a valid n - gram from the corpus.
 - **Algorithm:**
 - * Obtain a random word \tilde{w} from the dictionary to create the corrupt n -gram $\tilde{G} = (w_1, \dots, w_{n-1}, \tilde{w})$
 - * Pass G through the **embedding dictionary** to receive $emb(G)$, concatenate the vectors to receive a $d \cdot n$ length vector - pass this vector through the **Language model network** to obtain a score S_G .
 - * Perform the same operations on \tilde{G} to obtain $S_{\tilde{G}}$
 - * Backpropagate the error in formula (4.1) through through the **Language model network** to obtain an error $e_{lm} \in \mathbb{R}^{d \cdot n}$ at its input nodes.
 - * split e_{lm} to n vectors: $(e_{lm_1}, \dots, e_{lm_n})$, where $e_{lm_i} \in \mathbb{R}^d$
 - * backpropagate e_{lm_i} through the appropriate embedding dictionary item (e.g, w_i) and update w_i 's embedding.

This algorithm is run for every n -gram in the corpus.

5.3.2 Binarizing Parse Trees

5.3.2.1 Overview

Neither constituency parse trees nor dependency parse trees are necessarily binary, as seen in the examples shown at chapter 4. This presents a problem if one wishes to perform autoencoding of those parse trees. This is due to the fact that the autoencoding of a tree structure representation of data requires a fixed size number of sons for each node. This limitation is a derivation of the fixed size limitation of neural networks, as stated in section

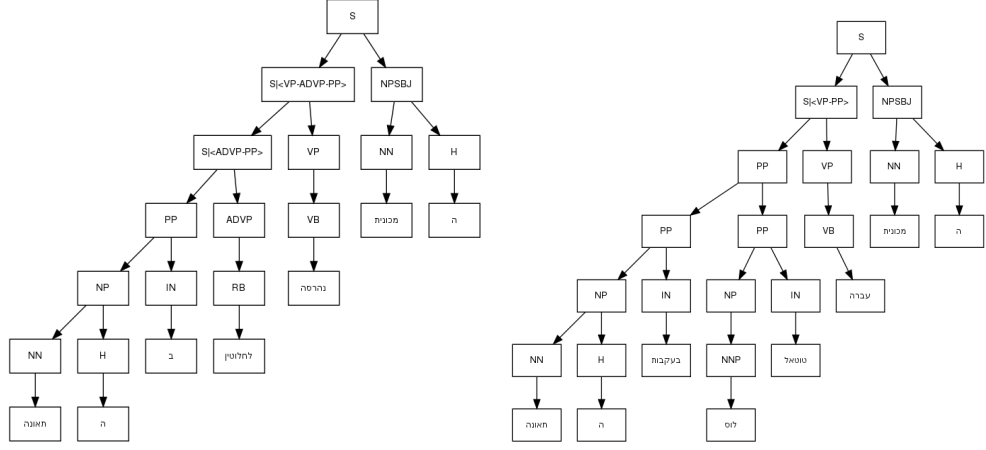


Figure 5.2: Constituency Binarization Paraphrasing Example

5.2.1.

Therefore, in order to be able to autoencode parse tree of both formats, some sort of “binarization” process of each format needs to be formulated.

5.3.2.1.1 Input The filtered n -grams mentioned in section 5.3.1 are pre parsed (using [20]), which yields a parse tree. This parse tree is further processed to replace the leaves (the original words) with their corresponding embedding. E.i, each leaf $w_i \in (w_1, \dots, w_n)$ is preprocessed to obtain $emb(w_1)$. Thus the input of this stage is a parse tree T which leaves were changed to the corresponding embedding.

5.3.2.1.2 Output A binary tree representing the same information as T is output.

5.3.2.2 Binarization of Constituency Parse Trees

The binarization of constituency parse trees is done in a rather simple approach, when compared with binarization of dependency parse trees. This is due to the fact that the words of the original sentence are already present in the tree leaves. Intuitively, binarization of such tree T over n nodes, is generating for each node which has $m > 2$ sons - An hierarchy of sons which introduce $O(\log m)$ new synthetic nodes. Each of which having only two sons. The obvious setback with this approach is the increase in the size of the parse tree by factor of $\log(n)$. Figure 5.2 revisits figure 4.6, this time the trees are binarized. Notice the number of synthetic nodes inserted into the trees - these are marked to point the pair of sons which were originally at that level of the tree.

5.3.2.3 Binarization of Dependency Parse Trees

The binarization of dependency parse trees is not an intuitive task, in contrast to binarization of constituency parse trees. No known method was found to propose an algorithm for such

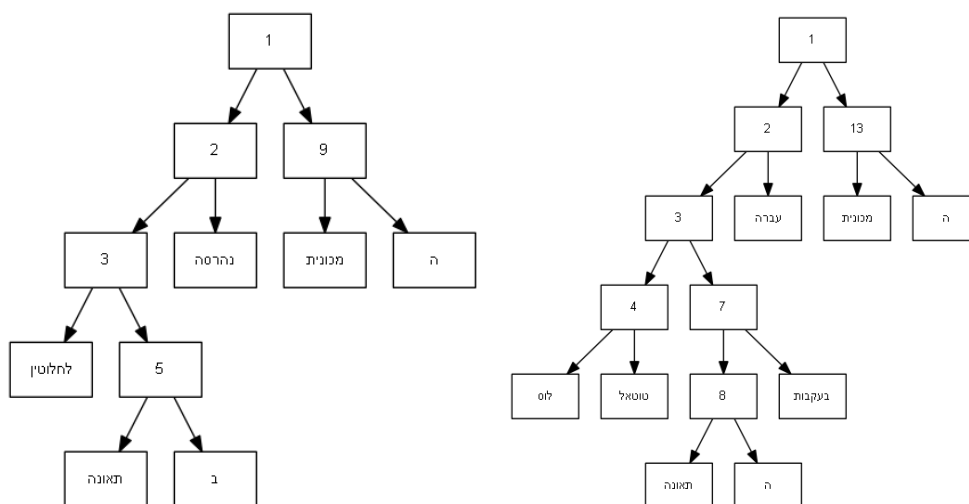


Figure 5.3: Dependency Binarization Paraphrasing Example

binarization. There are several special features that would be desired from the parse tree, output by a binarization algorithm for this task:

- The words of the sentence must appear only at the nodes of the tree, as opposed to the general definition of dependency parse trees (which was given at chapter 4). This is due to the fact that we employ autoencoding on this tree.
- Keep the head-dependency relation embedded in the format of the tree.
- Linear ordering of the leaves.
- Keep high nodes in the original dependency tree also high on the binary tree. This requirement is due to the fact that higher nodes in the parse tree are more likely to be encoded better. Since they are involved in less encodings, less noise is added along their way to the final encoding of the root of the tree. Naturally, it is desired that primary notions in the sentence (such as the head of relations), be encoded better.

Examples of the binary trees output by this process appear in 5.3 (again, a revision of figure 4.7). Notice the added syntactic leaves, marked with a running index (the index only for convenience of reference, and does not mark any other information encoded in that node). It can be seen that in this binary example of a paraphrase pair, transformed from dependency relation, we can still see the matching subtrees across the pair, as was discussed in section 4.3.2.

5.3.3 Learning to Parse on Embeddings Using Autoencoders

5.3.3.1 Overview

This section explains the process of training autoencoders on binary parse trees. The leaves of the trees are replaced with the embedding of an input sentence, as was described in the

previous section. The target of this, is to generate fixed size encodings of constituents of the parse trees. Later stages can use this encoding of a subtree as way to reference a subtree of variable size in a fixed size notation.

5.3.3.1.1 Input The binary parse tree which were the output of section 5.3.2 serve as input for this stage.

5.3.3.1.2 Output An autoencoder encoding from 200 to 100 features is trained on the set of all input binary n -grams parse trees. This autoencoder is later used to obtain fixed size representations of higher nodes in the parse tree, these will serve in giving a metric to compare inner nodes between a given test pair.

5.3.3.2 Training Process

For an input binary tree T , the process of training follows.

- Initialize an autoencoder network which encodes from $2d$ to d (and naturally decodes back from $2d$ to d), as was described in 5.2.3
- Encode each pair of sibling leaves of the parse tree, and put the result encoding in the father.
- Continue by encoding inner nodes of which both sons were encoded.
- At the end of this process, all nodes of the parse tree contain an encoding of its spanned subtree.

Along this process many encodings are performed, use these encodings to train the autoencoder, again as described in 5.2.3. Thus at the end of this process we obtain an autoencoder which is trained in efficiently encoding a parse tree in such a from to best reconstruct it during decoding. Since we parsed along a parse tree of linguistic meaning, the root of subtree contain encoding of the constituent they represent, hopefully encompassing interesting features of the constituent, in a fixed size notation.

Figure 5.4 shows an example of such an encoding. Dotted edges $(u_1, v), (u_2, v)$ mark that the encoding of $u_1 u_2$ will appear at the data of v .

5.3.4 Tree Matching

5.3.4.1 Overview

As described in section 5.2.4, tree matching is a new concept described in this work. It aims at formulating a paraphrase pair in terms of a matching function between subtrees of the parse trees, in which the words were changed for their embedding, binarized, and autencoded. Upon this matching a score is defined. This score tries to predict the probability of the given

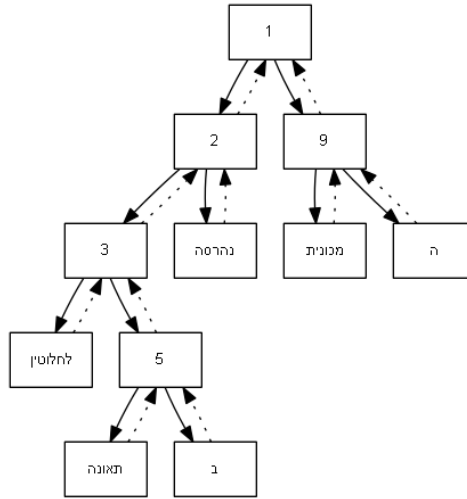


Figure 5.4: Encoding of a Binary Dependency Tree

pair being in paraphrase relation, based on the given matching. The hypothesis claims that the problem of paraphrase identification reduces to finding a best match between two given texts' parse trees, and deciding on a score threshold below which the texts are considered to be paraphrases.

5.3.4.1.1 Input A parsed, binarized, and autoencoded pair of texts is the input of this section.

5.3.4.1.2 Output The score of the **minimal** match is output. On a given training set, a classifier can then be trained to find a threshold in the data below which pairs are considered to be paraphrases. Thus yielding paraphrase identification.

5.3.4.2 Finding Tree Matching

Since tree matching is an NP-Complete problem, we use a very simple heuristic. Exhaustively searching the possible matches and recording the best one found. In order to reduce this exponential run time, we stop the search after a predefined number of iterations, and output the minimal match found thus far.

Figure 5.5 depicts a possible (and probably desired) match - nodes framed in box of the same color are matched against each other. It can be seen that this matching is according with the definition given for tree matching at section 5.2.4.

5.4 Experiments

This chapter describes the experiments which were conducted in order to prove the validity of the hypotheses that were introduced in the course of this work. These are targeted to test the separate contribution of each of the sections of the work.

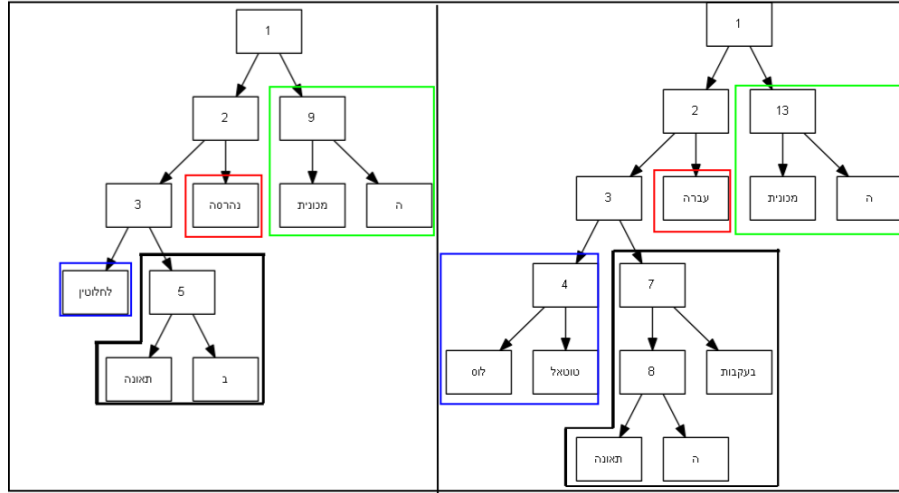


Figure 5.5: Possible Match of a Paraphrase Pair

5.4.1 Testing Embeddings

Collobert & Weston embeddings were shown to enhance the performance of many NLP tasks in English (C&W,2008,2011)(Bengio)(Turian). Although in the original design (C&W,2008) the embeddings were designed to change during the simultaneous supervised training process of the NLP tasks (thus exploiting the training process of one task to serve another's), it has been shown that these embeddings may as well serve as an off-the-shelf “plug-in” enhancer (Turian). This approach takes an already proven system for a common NLP task which treats words simply as index to a lexicon (thus using a very sparse representation), and exchanging its representation for words with the corresponding embedding (thus using a more condensed, linguistically based, encoding).

We test the “plug in” approach, testing the improvement over the task of POS tagging.

5.4.1.1 Improving POS tagging

We test the Conditional Random Field (CRF,(Lafferty et al., 2001) and (Sha and Pereira,2003)) model performance, when changing the words indexes for our Hebrew embeddings over labeled data. We compare against baseline of using the CRF model without the embeddings, and against state of the art Hebrew POS tagging ([5]).

5.4.2 Testing Autoencoding of Parse Trees on Embeddings

The goal of this section is to test the validity and contribution of the autoencoding stage in the algorithm, Following (Socher), we test the autoencoding in several criteria:

I Sentence level nearest neighbors:

- 1 Embed a corpus of sentences, and record a dictionary which contains also embeddings of sentences and all subtrees present in them, along with the original word

embeddings.

- 2 Given an input sentence, search for its nearest neighbor in this dictionary.
- 3 Compare the neighbor sentence with the input sentence.

II Ability to reconstruct from embedding:

- 1 Given a test parse tree of a sentence.
- 2 Encode it using the Hebrew autoencoder.
- 3 Decode the encoded parse tree, starting only from the root node.
- 4 For each of the leaves, which now contain some feature vector, search for their nearest neighbor (using euclidean metric) in the embedding dictionary - thus forming a “decoded sentence”.
- 5 Compare this decoded sentence against the original sentence.

5.4.3 Testing Tree Matching as Paraphrase Indicator

In order to test if tree matching reduction of paraphrasing is indeed a valid hypothesis, we test its validity using a common train-test partition of a paraphrase corpus.

5.4.4 Testing Dependency vs Constituency Parse Trees

As was discussed throughout the work, there are several possible definitions of parse trees. We have focused on dependency and constituency parse trees. In addition we have defined a binarization on both formats, which is an obligatory stage in order to later autoencode the trees. In order to compare both formats and their respective binarization. We run the previous tests both on binary dependency and binary constituency parse trees.

Chapter 6

Experimental Setting and Results

This chapter gives practical information regarding the corpora on which experiments are conducted, and the results obtained.

6.1 Experimental Setting

This section describes in detail the corpora used in this work, by means of its statistical distribution and its size. Each sub section will describe a single corpus which was used.

6.1.1 Embeddings Generation Corpus

The corpus used for the embeddings generation is a large corpus of the news domain. The details of this corpus are given in table 6.1. This corpus is divided into 5 sections by the complexity (how common the words are) of the words which appear in it. Since the segmented words are part of the most common words, the dictionary which divide by number of common words do not contain a value for “unsegmented” entry.

6.1.2 Autoencoding of Parse Trees Training Corpus

The same corpus described in section 6.1.1 is used for the training of parse trees. Since this corpus is formed from natural language sentences, we can parse it and use it as a training corpus for sentences, as opposed to the previous section which took it as a source for 5-grams.

Table 6.1: Size of the Different Corpora used for Embeddings computation

Dictionary	Unsegmented	Segmented
Full	115M	131M
5K	55.3M	63.6M

6.1.3 Labeled Paraphrase Corpus

This corpus is used for testing the proposed method for paraphrase identification in a supervised manner. It is composed of matched headlines from leading news sites. The size of the overall news wire corpus, before searching for paraphrases is of size of about 1.4 Million news headlines.

6.2 Results

This section describes the expected results for the aforementioned experiments in section 5.4. The experiments are still running and we expect to get results by the end of June.

Each of the following sections deals with the expected results of a single experiment:

- I **Embeddings:** We expect to see that commonly exchangeable words will receive a “close” embedding with respect to other words in the dictionary. For instance, days of the weeks should be “packed” rather close to each other. The POS tagging experiment should show significant improvement when exchanging words for their embeddings.
- II **Autoencoding:** The trained autoencoder should be able to reproduce correctly very short sentences (such as 4 word sentences), and to map longer sentences onto grammatically valid, related sentences. in the nearest neighbors experiment we should observe that the neighboring sentences of a given phrase should convey a meaning closer to its own meaning.
- III **Paraphrasing:** We expect to see results which come close to the state of the art results for paraphrase identification with respect to standard statistical metrics, such as accuracy and F1.
- IV **Comparing Parse Trees:** This experiment sets out to find which parse tree gives better results. Thus, any result shall be of interesting implication. We expect to see better performance when autoencoding constituency tree, since they seem to keep a more stabilized tree structure when looking at examples of paraphrase pair parse trees. A few examples of this were given in previous sections.

Chapter 7

Conclusion

The objective of this work is to research the field of paraphrasing in Hebrew. Providing first results in a field which English counterpart has gained a lot of interest and focus in recent years. The focus was set on the task of paraphrase identification.

We developed a method to identify if an input pair is in a paraphrase relationship, via adapting English state of the art algorithms, and providing modification of these.

The datasets we use are:

- A large news wire corpus, preprocessed with POS information, and segmentation.
- A large corpus of Hebrew paraphrases, formed in a parallel corpus of news wire headlines. This corpus was obtained in an unsupervised manner, by scraping major Hebrew news sites.

We run several experiments on these datasets, aiming at proving the hypotheses raised during the course of this work. The main experiment aims at assessing the validity of the proposed system as Hebrew paraphrase identifier. This consists of splitting the corpus of paraphrases in a train-test split and checking its performance via the standard statistical metrics.

The results indicate that Hebrew paraphrase identification is a feasible task, when comparing against result in the English equivalent field.

The main contributions of this work is:

- Showing that paraphrasing is research worthy field, also in Hebrew, while providing first results.
- Providing Hebrew resources for later paraphrase tasks (a paraphrase corpus), as well as for more general Hebrew NLP tasks (Collobert and Weston embeddings as a plug in enhancer).

As future work, we would like to further the research in the following optional directions:

- Explore the usage of the products of this work as tools for other paraphrase related tasks, such as paraphrase generation, and textual entailment.
- Test the paraphrasing identifier as a component in a more complex system of NLP or NLG, such as automatic multi-text text summarization, text generation system which is targeted for a specific audience.
- Use the word embeddings produced in this work in other NLP tasks.
- Formulate the Hebrew Paraphrase Corpus to have more features like those of the MSRP (in terms of average length of sentence, for instance).

Bibliography

- [1] Microsoft research paraphrase corpus, <http://research.microsoft.com/en-us/>. Technical report.
- [2] Textual entailment search pilot, guidelines. Technical report, 5th Textual Entailment Challenge, 2009.
- [3] Knowledge base population validation task, guidelines. Technical report, 7th Textual Entailment Challenge, 2011.
- [4] Main task and novelty detection subtask, guidelines. Technical report, 7th Textual Entailment Challenge, 2011.
- [5] Meni Adler. *Hebrew Morphological Disambiguation: An Unsupervised Stochastic Word-based Approach*. PhD thesis, Ben-Gurion University of the Negev, Beer-Sheva, Israel, 2007.
- [6] Regina Barzilay and Kathleen R. Mckeown. Extracting paraphrases from a parallel corpus. In *In Proc. of the ACL/EACL*, pages 50–57, 2001.
- [7] Yoshua Bengio. Learning deep architectures for ai. *Foundations and Trends in Machine Learning*, 2(1):1–127, 2009.
- [8] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *ICML*, page 6, 2009.
- [9] Eugene Charniak and Mark Johnson. Coarse-to-fine n-best parsing and maxent discriminative reranking. In *ACL*. The Association for Computer Linguistics, 2005.
- [10] Noam Chomsky. Three models for the description of language. *IRE Transactions on Information Theory*, 2:113–124, 1956.
- [11] Eve V. Clark. On the logic of contrast. *Journal of Child Language*, 15:317–335, 1988.
- [12] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: deep neural networks with multitask learning. In *Proceedings of the*

- 25th international conference on Machine learning*, ICML '08, pages 160–167, New York, NY, USA, 2008. ACM.
- [13] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu, and Pavel P. Kuksa. Natural language processing (almost) from scratch. *CoRR*, abs/1103.0398, 2011.
 - [14] Gil Diesendruck. The principles of conventionality and contrast in word learning: An empirical examination. *Developmental Psychology*, 41(3):451–463, 2005.
 - [15] V. Evans, B.K. Bergen, and J. Zinken. *The Cognitive Linguistics Reader*. Advances in Cognitive Linguistics. Equinox, 2007.
 - [16] Christiane Fellbaum, editor. *WordNet: An Electronic Lexical Database (Language, Speech, and Communication)*. The MIT Press, illustrated edition edition, May 1998.
 - [17] William A. Gale and Kenneth W. Church. A program for aligning sentences in bilingual corpora. *Computational Linguistics*, 1993.
 - [18] C. Goddard and A. Wierzbicka. *Meaning and Universal Grammar: Theory and Empirical Findings*. J. Benjamins Pub., 2002.
 - [19] Cliff Goddard. *Semantric Molecules*. Selected Papers of the 2006 Annual Meeting of the Australian Linguistic Society.y, 2006.
 - [20] Yoav Goldberg and Michael Elhadad. Easy-first dependency parsing of modern hebrew. In *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages*, pages 103–107, Los Angeles, CA, USA, June 2010. Association for Computational Linguistics.
 - [21] R. Harman and University of Massachusetts Amherst. Education. *Systemic Functional Linguistics and the Teaching of Literature in Urban School Classrooms*. University of Massachusetts Amherst, 2008.
 - [22] Chikara Hashimoto, Kentaro Torisawa, Stijn De Saeger, Jun’ichi Kazama, and Sadao Kurohashi. Extracting paraphrases from definition sentences on the web. In *ACL*, pages 1087–1097, 2011.
 - [23] Graeme Hirst. Paraphrasing paraphrased. Technical report, University of Toronto, 2003.
 - [24] R. Karp. Reducibility among combinatorial problems. In R. Miller and J. Thatcher, editors, *Complexity of Computer Computations*, pages 85–103. Plenum Press, 1972.

- [25] Dan Klein and Christopher D. Manning. Accurate unlexicalized parsing. In *IN PROCEEDINGS OF THE 41ST ANNUAL MEETING OF THE ASSOCIATION FOR COMPUTATIONAL LINGUISTICS*, pages 423–430, 2003.
- [26] Prodromos Malakasiotis and Ion Androutsopoulos. A generate and rank approach to sentence paraphrasing. In *EMNLP*, pages 96–106, 2011.
- [27] Marie-Francine Moens and Stan Szpakowicz, editors. *ROUGE: A Package for Automatic Evaluation of Summaries*, Barcelona, Spain, July 2004. Association for Computational Linguistics.
- [28] Joakim Nivre. Dependency grammar and dependency parsing. Technical report, Växjö University: School of Mathematics and Systems Engineering, 2005.
- [29] Noam Ordan, Bar Ilan, Noam Ordan, and Shuly Wintner. S.: Hebrew wordnet: a test case of aligning lexical databases across languages. *International Journal of Translation*, 19:39–58, 2007.
- [30] Anthony Pick. *Discourse and Function. A Framework of Sentence Structure*. <http://www.discourseandfunction.com>, 2009.
- [31] Jordan B. Pollack. Recursive distributed representations. *Artificial Intelligence*, 46:77–105, 1990.
- [32] Chris Quirk, Chris Brockett, and William Dolan. Monolingual machine translation for paraphrase generation. In *In Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 142–149, 2004.
- [33] Richard Socher and Eric H. Huang and Jeffrey Pennington and Andrew Y. Ng and Christopher D. Manning. Dynamic Pooling and Unfolding Recursive Autoencoders for Paraphrase Detection. In *Advances in Neural Information Processing Systems 24*. 2011.
- [34] Raul Rojas. *Neural Networks: A Systematic Introduction*. Springer, 1 edition, July 1996.
- [35] Richard Socher, Jeffrey Pennington, Eric H. Huang, Andrew Y. Ng, and Christopher D. Manning. Semi-Supervised Recursive Autoencoders for Predicting Sentiment Distributions. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2011.
- [36] Joseph Turian, Lev Ratinov, and Yoshua Bengio. Word representations: a simple and general method for semi-supervised learning. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, ACL ’10, pages 384–394, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics.

- [37] Anna Wierzbicka. *Semantic primitives*. (Frankfurt/M.) Athenaum-Verl., 1972.
- [38] Anna Wierzbicka. *Mental Lexicon*. Australian National University, 2009.

Appendix A

Proof of Tree Matching NP Completeness

This section provides a proof for the NP completeness of tree matching as described in section 5.2.4. We will do this by showing that

I Tree Matching \in NP

II Positive SubsetSum \leq_p Tree Matching

Since Positive SubsetSum (PS) is an NP-Complete problem (is a variation of one of Karp's 21 NP-Complete problems[24]), proving these yields that Tree Matching (TM) is NP-Complete. We continue by formally defining the PS decision problem, and the TM decision problem.

Definition Given a set of k positive integers $S = (n_1, \dots, n_k)$ and a target positive integer t , the *PS decision problem* is finding if there is a subset of S which sums exactly to t .

Formally:

$$PS = \{ (\{n_1, \dots, n_k\} \in \mathbb{N}^k, t \in \mathbb{N}) \mid \exists T \subseteq [1, k] : \sum_{j \in T} n_j = t \}$$

Definition Given a two binary trees T_1, T_2 which each nodes contain a 100 dimensional vector of numbers, and a target positive integer t , the *TM decision problem* is finding if there is a Tree Matching M of T_1, T_2 for which $S(M) = t$. (Recall section 5.2.4 for definitions)

Formally:

$$TM = \{ ((T_1, T_2), t \in \mathbb{N}) \mid \exists S \text{ a match of } T_1, T_2 : S(M) = t \}$$

Proof. Following these definitions we go on to prove that $PS \leq_p TM$.

Reduction: We will build a function f that receives ($S = \{n_1, \dots, n_k\} \in \mathbb{N}^k, t \in \mathbb{N}$) and returns ($(T_1, T_2), t' \in \mathbb{N}$) for which:

$$(S, t) \in PS \Leftrightarrow ((T_1, T_2), t') \in TM$$

f produces its output in the following manner:

1. $t' \leftarrow t$
2. T_1 is a binary tree which complies to:
 - * Contains k leaves, marked l_1, \dots, l_k .
 - * The fathers of all leaves have only one son. Mark them as p_1, \dots, p_k .
 - * All inner nodes contain $\vec{0}$ as their 100 dimensional vector.
 - * l_i contains $(\frac{n_i}{2}, 0, \dots, 0)$ as its 100 dimensional vector.
2. T_2 is a binary tree which complies to:
 - * Contains k leaves, marked r_1, \dots, r_k .
 - * All inner nodes contain $(0, t, 0, \dots, 0)$ as their 100 dimensional vector.
 - * r_i contains $\vec{0}$ as its 100 dimensional vector.

Correctness:

1. $(S, t) \in PS \Rightarrow ((T_1, T_2), t') \in TM$

Assume $(S = \{n_1, \dots, n_k\}, t) \in PS \Rightarrow \exists T \subseteq [1, k] : \sum_{i \in T} n_i = t$

We will define a matching M of T_1, T_2 :

- $\forall i \in T : (l_i, r_i) \in M$
- $\forall i \notin T : (p_i, r_i) \in M$

It is clear that M is indeed a Match, since for each leaf we either add it (first case), or its father (second case). Thus, we calculate the score of the match:

$$\begin{aligned}
 S(M) &= \sum_{i \in T} \|l_i, r_i\| \cdot 2 + \sum_{i \notin T} \|p_i, r_i\| \cdot 3 = \\
 &= \sum_{i \in T} \frac{n_i}{2} \cdot 2 + \sum_{i \notin T} 0 \cdot 3 = \\
 &= \sum_{i \in T} n_i = t = t'
 \end{aligned}$$

Therefore $((T_1, T_2), t') \in TM$

2. $((T_1, T_2), t') \in TM \Rightarrow (S, t) \in PS$

Assume $((T_1, T_2), t') \in TM \Rightarrow \exists M$ a Match, for which $S(M) = t' = t$.

We will use a lemma that will be proven later that M must be formed of exactly k pairs of the form $(*, r_i)$ (otherwise, by the construction of the trees, the score must be greater).

Following the lemma it is clear that $M = M_1 \cup M_2$, where M_1 contains only pairs of the form (l_i, r_j) and $M_2 = M \setminus M_1 = \{ (x_i, r_j) \in M : x_i \notin \{l_1, \dots, l_k\} \}$. Thus:

$$\begin{aligned}
 t = S(M) &= \sum_{(l_i, r_j) \in M_1} \|l_i, r_j\| \cdot 2 + \sum_{(x_i, r_j) \in M_2} \|x_i, r_j\| \cdot 3 = \\
 &= \sum_{(l_i, r_j) \in M_1} \|l_i, r_j\| \cdot 2 = \sum_{(l_i, r_j) \in M_1} n_i
 \end{aligned}$$

We see there is a subset of S which sums to $t \Rightarrow (S, t) \in PS$

□

Lemma A.1 $S(M) = t \Rightarrow M$ must be formed of exactly k pairs of the form $(*, r_i)$

Proof. Assume by negation that there exists a match M of T_1, T_2 , for which $S(M) = t$ and that there exists a pair $(x, y) \in M$ s.t

$$x \in T_1, y \in T_2, y \notin \{r_1, \dots, r_k\}$$

Mark as $c > 2$ - the number of spanned leaves by x and y .

Note, that by the formation of the trees, we get $\|x, y\| = \sqrt{(t-0)^2 + * } \geq t$

Thus:

$$S(M) \geq \|x, y\| \cdot c \geq \|x, y\| \cdot 2 \geq 2t > t$$

In contradiction with the assumption that $S(M) = t$.

□

Appendix B

Paraphrase Tagging Guidelines

B.1 Goal

The goal of tagging is to find sentences and parts of sentences which convey identical, or near-identical, information, articulated in different forms.

B.2 Rules for Tagging

Given two sentences S_1, S_2 , the annotator should decide between the following options:

B.2.1 Paraphrase

the pair would be considered a paraphrase in the following case:

1. A human who fully accepts S_1 , must therefore accept S_2 as a whole.
2. The same holds for the second direction as well. Namely, accepting S_2 , must yield acceptance of S_1 .

B.2.1.1 Restrictions

- The acceptance can derive also by former knowledge, yet it must not be based solely on it. (examples 1,2)
- The pair will be considered as a paraphrase, when there is sufficient confidence that they indeed convey the same information. (example 3)
- Given a pair for which their meaning is ambiguous (for example - there is not enough information to be certain regarding the identities which appear in them). One must assume that the ambiguity is solved in such a way that the pair relate to same incident. (example 4)

- The pair should be judged from a timeless point of view. Namely, if the pair relate the same information, yet treat it from a different point in time - it should be tagged as paraphrase.
- The pair should be judged from a location-less point of view. Namely, if the same incident is reported, yet it is regarded from a different relative point of view, it should be considered paraphrase.
- If one of the sentences, or both of them, contain additional information, or contradictory information, the pair should not be tagged as a paraphrase.

B.2.2 Partial Paraphrase

If S_1, S_2 are not paraphrases, by the definition given above, yet certain clauses can be removed to form a paraphrase, then S_1, S_2 should be tagged as partial paraphrases. The additional clauses should be annotated as such. (example 6)

B.2.3 Negative

A catch-all rule. If S_1, S_2 do not fall under any of the previous tags, then they should be tagged as Negative, meaning they do not convey any shared information.

B.3 Examples

1. S_1 - סין הגיעה להסכם סחר נפט עם רוסיה.

China has reached an oil trade contract with Russia

- S_2 - המדינה המאוכלסת בעולם חתמה על הסכם מסחר נפט עם רוסיה.

The world most populated nation has reached an oil trade contract with Russia

Tag: Paraphrase.

(It is of common knowledge that China is the world most populated nation)

2. S_1 - סין הגיעה להסכם סחר נפט עם רוסיה.

China has reached an oil trade contract with Russia

- S_2 - סין היא המדינה המאוכלסת בעולם.

China is the world most populated nation.

Tag: Negative.

(Although the second sentence is true, it does not follow from the first, nor does the second derives the first).

3. S_1 - ארה"ב סגרה השגרירות בסוריה.

USA has closed its embassy in Syria.

S_2 - ארה"ב החזירה את שגרירה מסוריה.

USA has called its ambassador back from Syria.

Tag: Paraphrase.

(It is very probable to assume that closing the embassy yields calling back the ambassador, and vice versa)

4. S_1 - הרשת סיקרה את הארועים האחרונים בסוריה.

The network has covered latest events in Syria.

S_2 - רשת אלג'זירה דיווחה על ההתרחשויות האחרונות בסוריה.

Al Jazira network has passed reports of the latest happenings in Syria.

Tag: Paraphrase.

(It should be assumed that the network mentioned in the second sentence is the same network mentioned in the first sentence)

5. S_1 - הרכבת תשבית את הקו לבאר שבע ביום ראשון.

Trains will not work on the line to Beer Sheva this Sunday, due to strike.

S_2 - הרכבת תפסיק מחר את הקו לבאר שבע.

Trains will not work on the line to Beer Sheva tomorrow, due to strike.

Tag: Paraphrase.

(Although the pair convey the same information from different points in time, it should be regarded as paraphrase)

6. S_1 - בנימין נתניהו נפגש עם עמיתו הרוסי, ולאחר מכן חזר לארץ.

Benjamin Netanyahu has met with his Russian counterpart, and later returned to Israel.

S_2 - בנימין נתניהו נועד עם מקבילו הרוסי, וציין בפניו את סוגיית העולים החדשים.

Benjamin Netanyahu has met with his Russian counterpart, and mentioned to him the issue of Israeli immigrants from Russia.

Tag: Partial Paraphrase.

(Without the marked parts, the pair consists a paraphrase)

7. S_1 - אהוד ברק נאם בפני המתגייסים הצעירים.

Ehud Barak spoke in front of new army recruits.

S_2 - אהוד ברק החליט להגיש מועמדותו לקדנציה הבאה.

Ehud Barak decided to run for the next elections.

Tag: Negative.

(Both sentences do not convey the same information, nor can any clauses be removed from them so that they do)