

“happy” - “unhappy”, “agree”- “disagree”) and knowledge bases such as WordNet [Miller, 1995], WikiData [Vrandečić, 2012], DBpedia [Auer et al., 2008] and ConceptNet [Speer and Havasi, 2012]. The tasks and resources for lexical meaning relations are also relevant for the research on textual meaning relations.

Textual Meaning Relations also have an impact on **Other Areas of CL and NLP** Systems that can successfully process meaning relations can also be used in other tasks in CL and NLP, such as text summarization [Lloret et al., 2008, Harabagiu and Lacatusu, 2010], text simplification [Yimam and Biemann, 2018], plagiarism detection [Barrón-Cedeño et al., 2013], question answering [Harabagiu and Hickl, 2006], and machine translation evaluation [Padó et al., 2009], among others.

1.2 Motivation and Objectives of the Thesis

This thesis arose from an interest in applying linguistic knowledge to the empirical studies of textual meaning relations. My research was motivated by two gaps in the research field:

- a lack of large-scale corpora for research on decomposing textual meaning relations and, as a consequence, a lack of machine learning experiments.
- insufficient resources (annotation guidelines and corpora) and a lack of empirical studies on multiple textual meaning relations.

I address both these gaps in turn by combining theoretical knowledge and empirical, data-driven approaches (human judgments, statistical analysis, and machine learning experiments). First, I bring together paraphrase typology and the task of Paraphrase Identification. Second, I present a joint study on the textual meaning relations of Paraphrasing, Textual Entailment, and Semantic Similarity. In the rest of this section I present in more detail the objectives behind each of the two research directions of my thesis.

1.2.1 Paraphrase Typology and Paraphrase Identification

The work on Paraphrase Typology (PT) uses knowledge from theoretical linguistics to understand the Paraphrasing phenomenon. Paraphrase Identification (PI) is an empirical task that aims to produce systems capable of recognizing paraphrasing in an automatic manner. However, at the time of beginning this dissertation, there had been almost no interaction or intersection between these two areas of

Paraphrasing research. PT research, prior to this dissertation, was mostly theoretical, with very limited practical implications and applications. PI research in the era of deep learning is radically empirical, focused on quantitative performance, with little to no interpretability and theoretical justification. My intuition was that these two research areas are not mutually exclusive, however there was no previous work trying to combine them. My **objectives** in combining PT and PI were twofold:

- Obj1 To use linguistic knowledge and paraphrase typology in order to improve the evaluation and interpretation of automated PI systems.
- Obj2 To empirically validate and quantify the difference between the various linguistic and reason-based phenomena involved in paraphrasing.

1.2.2 Joint Study on Meaning Relations

Meaning relations, such as Paraphrasing, Textual Entailment, and Semantic Similarity, have attracted a lot of attention from the researchers in Computational Linguistics (CL) and Natural Language Processing (NLP). There is a substantial amount of theoretical and empirical research on these meaning relations and many resources, datasets, and automated systems. Traditionally, these relations have been studied in isolation and the transfer of knowledge and resources between them has been very limited. My intuition was that these textual meaning relations can be brought together in a single corpus and compared empirically. My **objectives** in this part were twofold:

- Obj3 To empirically determine the interactions between Paraphrasing, Textual Entailment, Contradiction, and Semantic Similarity in a corpus of multiple textual meaning relations.
- Obj4 To propose and evaluate a novel shared typology of meaning relations. The shared typology would then be used as a conceptual framework for joint research on meaning relations.

1.3 Thesis Development

My research has three separate phases, described in parts I, II, and III of this thesis. First, I explore the basic concepts of Distributional Semantics and the notion of Semantic Similarity at the level of words and short phrases in Part I. Second, I present my empirical research on bringing together Paraphrase Typology and

Paraphrase Identification in Part II. Finally, I describe the setup and results of my joint study on multiple textual meaning relations in Part III.

The order of the chapters follows the chronological order in which the articles were written. At the same time, the order of the chapters follows the logical progression of my dissertation. Each of the articles is self sufficient: it poses its own research questions, presents related work, proposes a methodology, and describes the experimental results. However, there is also a clear thread that connects all the articles. When brought together, the articles tell a coherent story about the linguistic phenomena involved in textual meaning relations and how these phenomena can be used to improve the evaluation and interpretation of automated systems and bring together multiple textual meaning relations.

In the rest of this section I briefly present the main motivation, research questions and findings for each of the three parts and the logical progression of the thesis. I also discuss how each article fits within the more general objectives and how the different articles interact with each other.

Part I: Lexical Relations and Distributional Semantics

The two articles presented in this part of the thesis serve as an introduction to the research on meaning relations and aim to familiarize the reader with the core concepts and theories used in the whole thesis.

In the article “*Comparing Distributional Semantics Models for identifying groups of semantically related words*” (Chapter 2), I explore the theoretical concepts and the empirical tools within the framework of Distributional Semantics (DS). DS is the most popular framework in contemporary Natural Language Processing (NLP) and Computational Linguistics (CL) and in the research on lexical and textual meaning relations. I experiment with different methodologies for representing the meaning of individual words, different ways to quantitatively compare meaning representations, and different approaches to measuring semantic similarity at the level of words. Lexical similarity is the most “atomic” form of semantic similarity. Many aspects of lexical similarity are also important for semantic similarity at the level of longer pieces of texts.

In the article “*DISCover: DIStributional approach based on syntactic dependencies for discovering COnstructions*” (Chapter 3), I present a successful data-driven methodology that can compose individual words into short phrases. The methodology is based on Distributional Semantics, lexical semantic similarity, and syntactic similarity between words. Many of the resulting short phrases are novel and have never been observed in the training data, indicating that the system is composing as opposed to memorizing. This article demonstrates the importance of lexical similarity in the context of complex language expressions