



NUI Galway
OÉ Gaillimh



Doctoral Thesis

Monolingual Alignment of Word Senses and Definitions in Lexicographical Resources

SINA AHMADI

M.Sc., Paris Descartes University, 2017

M.A., Sorbonne Nouvelle University, 2016

B.Eng., University of Kurdistan, 2014

Supervisor

Dr. John P. McCrae

External Examiner

Prof. Marie-Claude L'Homme

Internal Examiner

Dr. James McDermott

*Thesis Submitted in Partial Fulfillment of the Requirements for the Degree of
Doctor of Philosophy*

in the
School of Computer Science
College of Science and Engineering
National University of Ireland Galway

Spring 2022

ABSTRACT

Dictionaries are fundamental resources for people to learn and document languages as well as for computers to process natural languages. A dictionary provides a fine-grained structure and description of the vocabulary of a language. With decades of advances in electronic lexicography, a significant amount of lexicographical resources are currently available. Such resources are the fruits of elaborate and strenuous efforts of lexicographers and oftentimes, are costly projects to initiate and maintain. Moreover, given the increasing number of lexical semantic resources thanks to community-driven initiatives such as Wiktionary, the alignment of such resources is of importance to promote interoperability and increase their exploitation more effectively. On the other hand, the significant progress in the field of computer science, artificial intelligence and the semantic web has been tremendously beneficial to various scientific fields, particularly language technology. Therefore, there is a necessity to leverage the current techniques and resources to facilitate the automatic alignment, integration and enrichment of lexicographical data.

The focus of this thesis is broadly on the alignment of lexicographical data, particularly dictionaries. In order to tackle some of the challenges in this field, two main tasks of word sense alignment and translation inference are addressed. The first task aims to find an optimal alignment given the sense definitions of a headword in two different monolingual dictionaries. This is a challenging task, especially due to differences in sense granularity, coverage and description in two resources. After describing the characteristics of various lexical semantic resources, we introduce a benchmark containing 17 datasets of 15 languages where monolingual word senses and definitions are manually annotated across different resources by experts. In the creation of the benchmark, lexicographers' knowledge is incorporated through the annotations where a semantic relation, namely exact, narrower, broader, related or none, is selected for each sense pair. This benchmark can be used for evaluation purposes of word-sense alignment systems. The performance of a few alignment techniques based on textual and non-textual semantic similarity detection and semantic relation induction is evaluated using the benchmark. Finally, we extend this work to translation inference where translation pairs are induced to generate bilingual lexicons in an unsupervised way using various approaches based on graph analysis. This task is of particular interest for the creation of lexicographical resources for less-resourced and under-represented languages and also, assists in increasing coverage of the existing resources. From a practical point of view, the techniques and methods that are developed in this thesis are implemented within a tool that can facilitate the alignment task.

CONTENTS

Declaration	vii
Acknowledgements	ix
Acronyms	xi
List of Figures	xv
List of Tables	xviii
1 Introduction	1
1.1 The ELEXIS project	3
1.2 Motivation	5
1.3 Research Questions	6
1.4 Thesis Structure	7
1.5 Publications	8
2 Background	11
2.1 Introduction	11
2.1.1 Lexical Semantic Resources	13
2.1.2 Polysemy	14
2.2 Dictionaries	16
2.2.1 Lexicon vs. Dictionary	17
2.2.2 Content	18
Headwords	19
Senses	20
Definitions	21
2.2.3 Electronic Dictionaries	23
2.3 Network-based Resources	24
2.3.1 WordNet	24
2.3.2 FrameNet	26
2.3.3 JeuxDeMots	27
2.4 Explanatory Combinatorial Dictionary	27
2.5 Generative Lexicon	30
2.6 Natural Semantic Metalanguage	32
2.7 Terminologies	34
2.8 Ontological Resources	35
2.8.1 Linguistic Linked Data	36
2.9 Knowledge Graphs	38
2.10 Language Models	41
2.11 Conclusion	44
3 Systematic Literature Review	47

3.1	Introduction	47
3.2	Creation and Modeling	48
3.3	Enrichment	51
3.3.1	Semantic Similarity Detection	52
	Corpus-based approaches	53
	LSR-based approaches	53
	Embeddings-based approaches	53
	Datasets	54
3.3.2	Language Resource Alignment	55
3.3.3	Translation Inference	56
3.4	Publication and Storage	59
3.5	Dictionaries in NLP applications	60
3.5.1	Word Sense Disambiguation	60
3.5.2	Semantic Role Labeling	60
3.5.3	Reverse Dictionary	61
3.6	What is missing?	62
4	Leveraging the graph structure of lexicographical resources	63
4.1	Introduction	63
4.2	Related Work	64
4.3	Lexicographical network	65
4.3.1	Analysis of lexicographical networks	65
4.3.2	Experiments	66
4.4	Weighted bipartite b-matching	68
4.4.1	String-based Methods	69
4.4.2	The WBbM algorithm	71
4.4.3	Experiments	72
4.5	Translation Inference	74
4.5.1	Cycle-based approach	74
4.5.2	Path-based approach	75
4.5.3	Experiments	76
4.6	Conclusion and Contributions	77
5	A Benchmark for Monolingual Word Sense Alignment	81
5.1	Introduction	81
5.2	Related Work	84
5.3	Methodology	86
5.3.1	Semantic Relationships	87
5.3.2	Data Selection	89
5.3.3	Dictionaries used in the creation of the dataset	90
5.3.4	Dataset Structure	92
5.4	Case Studies	92
5.4.1	Danish	92

	Sense structure	93
	Definition content	93
	Data structure	95
	Manual annotation	96
5.4.2	Italian	98
5.4.3	Portuguese	99
	Sense structure	100
	Definition content	100
	Manual annotation	101
5.5	Evaluation	103
5.5.1	Sense Granularity	104
5.5.2	Sense Alignments	104
5.5.3	Inter-annotator Agreement	107
5.6	Conclusion and contributions	111
6	Monolingual Word Sense Alignment	113
6.1	Introduction	113
6.2	Related Work	117
6.3	Naisc architecture	119
6.4	Textual Similarity Methods	122
6.4.1	String-based Methods	122
6.4.2	Beyond String Similarity	122
6.4.3	Word Alignment	126
6.4.4	Monolingual Alignment	127
6.5	Non-Textual Similarity Methods	128
6.6	Linking Constraints	129
6.6.1	Bijection	130
6.6.2	Taxonomic	131
6.6.3	Axiomatic	132
6.7	Semantic Relation Induction	133
6.7.1	Feature Extraction	134
6.7.2	Feature Learning	136
6.8	Experiments	139
6.8.1	Baseline system	139
6.8.2	System 1: Classification and Feature Learning	140
6.8.3	System 2: Length-limited Alignment	146
6.9	ELEXIS Monolingual Word Sense Alignment Task	147
6.10	Conclusion and contributions	150
7	Conclusions	153
7.1	Research Contributions	153
7.1.1	Benchmarking Word Sense Alignment	153
7.1.2	Alignment of dictionaries at the sense level	156

7.1.3	Alignment of dictionaries at the entry level	157
7.1.4	NAISC	158
7.2	Revisiting the Research Questions	158
7.3	Limitations and Future directions	160
	Bibliography	163

DECLARATION

I declare that this thesis, titled “*Monolingual Alignment of Word Senses and Definitions in Lexicographical Resources*”, is composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification.

Galway, March 7, 2022

Sina Ahmadi

ACKNOWLEDGEMENTS

Writing this section has inspired me since the very beginning. Perhaps because I was wondering how my Ph.D. journey would shape this text at the end. And here is how it goes.

This thesis is the result of my past four years work as a Ph.D. researcher. Although I am sure this work is not flawless, I must say that finishing up my Ph.D. during Covid-19 gave me immense satisfaction. Starting my Ph.D. on an exceptionally sunny day in April 2018 in Galway, I had a sweet episode of my life full of new ideas in research, rewarding experiences in engineering, papers, and conferences that were typical of Ph.D. life at the time. To my surprise, I even rarely struggled with the type of challenges that some of my fellows complained about, such as making progress, carrying out experiments, catching deadlines to submit papers, or having a life while doing a Ph.D. However, the impact of the physical restrictions and psychological burden that Covid-19 imposed on everyone, including myself, were undeniable.

2020 was a particularly harsh year for everyone but even worse for international postgraduate students living in Ireland. Doing a Ph.D. during Covid-19 did not only require taking care of regular activities but also learning to cope in the new mainly virtual world while living in a shared house and dealing with the emerging mental health issues far away from the loved ones. Now that I look back, I feel that I should be gleeful that those daunting days are gone, despite the long-lasting effects.

Regardless of the problems, I enjoyed every single day of this journey and will always remember it with sweet memories. Living in Ireland was not without challenge, but at least, it taught me to be more grateful for having things that we usually take for granted: the generous sun, a more predictable weather, good food, fresh artisan bread, decent accommodation, museums, and a more affordable place to live. Nevertheless, I will miss this place for the friendly Irish people, the grasslands that make running more joyful, the biodiversity and the wildlife, particularly those robins singing at night!

This work could not have been accomplished without the help of many people. First and foremost, I would like to thank my esteemed supervisor, Dr. John M^cCrae, for his invaluable supervision, wise pragmatism, support, tutelage, and kindness during the course of my Ph.D. degree. My gratitude extends to the ELEXIS project for the funding opportunity to undertake my studies at the School of Computer Science at NUI Galway. My sincere thanks also go to Dr. Mathieu d'Aquin and Dr. Paul Buitelaar who played an important role as my graduate research committee members, and to the examination committee members, Dr. Marie-Claude L'Homme and Dr. James McDermott.

This journey could not be this memorable without my friends at Insight: Tobias Daudert, Joana Barros, Omnia Zayed, Adrian Ó Dubhghaill, Oksana Dereza, Sarah Carter, Niki Pavlopoulou, and Brendan Smith. Many thanks to the administrative staff at Insight and ELEXIS for their support: Hilda Fitzpatrick, Christiane Leahy-Coen, and Anna Woldrich. My Ph.D. experience would have certainly been different without the compassionate friends who helped me in many ways, especially by hosting me with open arms when I had difficulty finding accommodation: Housam Ziad and his lovely family, Aftab Alam and Dr. Theodorus Fransen. Likewise, thanks to Jalal Sajadi and Daban Q. Jaff for their true friendship over the years. Special thanks to Dr. Sanni Nimb who hosted me during my visit to the Society for Danish Language and Literature and the Centre for Language Technology in Copenhagen in 2019 which was an enriching experience in many ways. I would also like to warmly thank Dr. Mathieu Constant for hosting me in Nancy in 2021. Throughout my visit, I had a great time sharing my office with Pauline Gillet and Charlène Weyh.

The accomplishment of my Ph.D. studies is important to me in two other ways too. First, it marks the end of my formal education which has been continuously going on since I remember. I partially owe this to France which provided me with free education and unique, eye-opening, rich, and unforgettable experiences. Also, I am indebted to Dr. Kyumars S. Esmaili who encouraged me during my bachelor's to follow my passion for languages, linguistics, and computer science by studying natural language processing. Second, it ends as the fourth decade of my life begins, so do my plans to contribute to society and have a more significant role "to make the world a better place", something that I whimsically told my parents when I left them and I am sometimes contemplative if ever it will be eventual! I will stay optimistic.

This brings me to the most important people in my life. There is no word to properly express my gratitude to my family back in Kurdistan, particularly to my parents who always believed in me, endorsed me in my choices, and even sacrificed themselves to make sure that their children have a fair chance to achieve their goals, more than their own generation. I understand that pursuing my studies imposed thousands of kilometers of distance between us and years of untogetherness that could not have passed this way without their support and unconditional love. *Spas û xoşewîstî!* I am also deeply grateful to my family in Greece who have always supported me, have given new meaning to my life with their kind hearts, positive vibes and love. *Άπειρα ευχαριστώ!*

Finally, I would like to thank the love of my life, Ioanna, from the bottom of my heart for being my best friend and companion and for her patience and encouragement over the years. It is impossible to imagine this journey getting to an end without her unending support and love. To her, I would say: *Je t'aime*.