
DEEP LEARNING BASED SEMANTIC TEXTUAL SIMILARITY FOR APPLICATIONS IN MACHINE TRANSLATION DOMAIN

THARINDU RANASINGHE

A thesis submitted in partial fulfilment of the requirements of the
University of Wolverhampton for the degree of Doctor of Philosophy

2021

This work or any part thereof has not previously been presented in any form to the University or to any other body whether for the purposes of assessment, publication or for any other purpose (unless otherwise indicated). Save for any express acknowledgements, references and/or bibliographies cited in the work, I confirm that the intellectual content of the work is the result of my own efforts and of no other person.

The right of Tharindu Ranasinghe to be identified as author of this work is asserted in accordance with ss.77 and 78 of the Copyright, Designs and Patents Act 1988. At this date copyright is owned by the author.

Signature:

Date:

ABSTRACT

An abstract is a synopsis of the thesis, and it goes in the file `abstract.tex`.

ACKNOWLEDGEMENTS

Your acknowledgements should go in `ack.tex`.

We would like to acknowledge Donald Craig at Memorial University, Newfoundland who published the meta-thesis on which this template is based. You can find Donald's work on his web site, here: <http://www.cs.mun.ca/~donald/metathesis/>.

CONTENTS

Abstract	ii
Acknowledgements	iv
List of Tables	x
List of Figures	xii
List of Code Listings	xiv
I Semantic Textual Similarity	1
1 Introduction	3
1.1 What is Semantic Textual Similarity?	3
1.2 Related Work	3
1.3 Datasets	3
1.3.1 English Datasets	3
1.4 Applications	4
2 State of the Art Methods	5
2.1 Introduction	5
2.2 Related Work	5
2.3 Improving State of the Art STS Methods	5
2.3.1 Portability to Other Languages	5
2.3.2 Portability to Other Domains	5
2.4 Conclusions	5

3	Sentence Encoders	7
3.1	Introduction	7
3.2	Related Work	7
3.3	Exploring Sentence Encoders in English STS	7
3.4	Portability to Other Languages	7
3.5	Portability to Other Domains	7
3.6	Conclusions	7
4	Siamese Neural Networks	9
4.1	Introduction	9
4.2	Related Work	9
4.3	MAGRU: Improving Siamese Neural Networks	9
4.3.1	Portability to Other Languages	9
4.3.2	Portability to Other Domains	9
4.4	Conclusions	9
5	Transformers	11
5.1	Introduction	11
5.2	Related Work	11
5.3	Exploring Transformers in English STS	11
5.4	Exploring Transformers for STS in Other Languages	11
5.5	Exploring Transformers for STS in Other Domains	11
5.6	Conclusions	11

II	Applications - Translation Memories	13
1	Introduction	15
1.1	What is Translation Memory?	15
1.2	Datasets	15
1.3	Related Work	15
1.4	STS for Translation Memories	15
2	Sentence Encoders for Translation Memories	17
2.1	Introduction	17
2.2	Methodology	17
2.3	Results and Evaluation	17
III	Applications - Translation Quality Estimation	19
1	Introduction	21
1.1	What is Translation Quality Estimation?	21
1.2	Datasets	21
1.3	Related Work	21
1.4	STS for Translation Quality Estimation	21
2	TransQuest: STS Architectures for QE	23
2.1	Introduction	23
2.2	Methodology	23
2.3	Results and Evaluation	23
	Bibliography	24

LIST OF TABLES

LIST OF FIGURES

LISTINGS

Part I

Semantic Textual Similarity

CHAPTER 1

INTRODUCTION

1.1 What is Semantic Textual Similarity?

[1]

1.2 Related Work

1.3 Datasets

We experimented with several datasets throughout the experiments in the Semantic Textual Similarity Section. In order to maintain the versatility of our methods we experimented with several English datasets as well as several non English datasets and several datasets from different domains which we will introduce in this section.

1.3.1 English Datasets

1. **SICK dataset** - The SICK data contains 9927 sentence pairs with a 5,000/4,927 training/test split which were employed in the SemEval 2014 Task1: Evaluation of Compositional Distributional Semantic Models on Full Sentences through Semantic Relatedness and Textual Entailment [2]. The dataset has two types of annotations: Semantic Re-

latedness and Textual Entailment. We only use Semantic Relatedness annotations in our research. SICK was built starting from two existing datasets: the 8K ImageFlickr data set ¹ [3] and the SemEval-2012 STS MSR-Video Descriptions dataset ² [4]. The 8K ImageFlickr dataset is a dataset of images, where each image is associated with five descriptions. To derive SICK sentence pairs the organisers randomly selected 750 images and sampled two descriptions from each of them. The SemEval2012 STS MSR-Video Descriptions data set is a collection of sentence pairs sampled from the short video snippets which compose the Microsoft Research Video Description Corpus ³. A subset of 750 sentence pairs were randomly chosen from this data set to be used in SICK.

In order to generate SICK data from the 1,500 sentence pairs taken from the source data sets, a 3-step process was applied to each sentence composing the pair, namely (i) normalisation, (ii) expansion and (iii) pairing [2]. In normalisation

2. STS 2017 Dataset

1.4 Applications

¹The 8K ImageFlickr data set is available at <http://hockenmaier.cs.illinois.edu/8k-pictures.html>

²The SemEval-2012 STS MSR-Video Descriptions dataset is available at <https://www.cs.york.ac.uk/semeval-2012/task6/index.html>

³The Microsoft Research Video Description Corpus is available to download at <https://research.microsoft.com/en-us/downloads/38cf15fd-b8df-477e-a4e4-a4680caa75af/>

CHAPTER 2

STATE OF THE ART METHODS

2.1 Introduction

[5]

2.2 Related Work

2.3 Improving State of the Art STS Methods

2.3.1 Portability to Other Languages

2.3.2 Portability to Other Domains

2.4 Conclusions

CHAPTER 3

SENTENCE ENCODERS

3.1 Introduction

[6]

3.2 Related Work

3.3 Exploring Sentence Encoders in English STS

3.4 Portability to Other Languages

3.5 Portability to Other Domains

3.6 Conclusions

CHAPTER 4

SIAMESE NEURAL NETWORKS

4.1 Introduction

[7]

4.2 Related Work

4.3 MAGRU: Improving Siamese Neural Networks

4.3.1 Portability to Other Languages

4.3.2 Portability to Other Domains

4.4 Conclusions

CHAPTER 5

TRANSFORMERS

5.1 Introduction

[8]

5.2 Related Work

5.3 Exploring Transformers in English STS

5.4 Exploring Transformers for STS in Other Languages

5.5 Exploring Transformers for STS in Other Domains

5.6 Conclusions

Part II

Applications - Translation Memories

CHAPTER 1

INTRODUCTION

1.1 What is Translation Memory?

[9]

1.2 Datasets

1.3 Related Work

1.4 STS for Translation Memories

CHAPTER 2

SENTENCE ENCODERS FOR TRANSLATION MEMORIES

2.1 Introduction

[10]

2.2 Methodology

2.3 Results and Evaluation

Part III

Applications - Translation Quality Estimation

CHAPTER 1

INTRODUCTION

1.1 What is Translation Quality Estimation?

1.2 Datasets

1.3 Related Work

[11]

1.4 STS for Translation Quality Estimation

CHAPTER 2

TRANSQUEST: STS ARCHITECTURES FOR QE

2.1 Introduction

[12]

2.2 Methodology

2.3 Results and Evaluation

BIBLIOGRAPHY

- [1] Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada, August 2017. Association for Computational Linguistics.

- [2] Marco Marelli, Luisa Bentivogli, Marco Baroni, Raffaella Bernardi, Stefano Menini, and Roberto Zamparelli. SemEval-2014 task 1: Evaluation of compositional distributional semantic models on full sentences through semantic relatedness and textual entailment. In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 1–8, Dublin, Ireland, August 2014. Association for Computational Linguistics.

- [3] Cyrus Rashtchian, Peter Young, Micah Hodosh, and Julia Hockenmaier. Collecting image annotations using Amazon’s mechanical turk. In *Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and Language Data with Amazon’s Mechanical Turk*, pages 139–147, Los Angeles, June 2010. Association for Computational Linguistics.

BIBLIOGRAPHY

- [4] Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. SemEval-2012 task 6: A pilot on semantic textual similarity. In **SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 385–393, Montréal, Canada, 7-8 June 2012. Association for Computational Linguistics.
- [5] Tharindu Ranasinghe, Constantin Orasan, and Ruslan Mitkov. Enhancing unsupervised sentence similarity methods with deep contextualised word representations. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019)*, pages 994–1003, Varna, Bulgaria, September 2019. INCOMA Ltd.
- [6] Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc Barrault, and Antoine Bordes. Supervised learning of universal sentence representations from natural language inference data. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 670–680, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.
- [7] Tharindu Ranasinghe, Constantin Orasan, and Ruslan Mitkov. Semantic textual similarity with Siamese neural networks. In *Proceedings of the International Conference on Recent Advances in Natural Language Pro-*

BIBLIOGRAPHY

- cessing (RANLP 2019)*, pages 1004–1011, Varna, Bulgaria, September 2019. INCOMA Ltd.
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [9] Peter J. Arthern. Machine translation and computerized terminology systems: A translator’s viewpoint. *Translating and the Computer, Proceedings of a Seminar, London 14th November 1978. Amsterdam: North-Holland Publishing Company*, pages 77–108, 1979.
- [10] Tharindu Ranasinghe, Constantin Orasan, and Ruslan Mitkov. Intelligent translation memory matching and retrieval with sentence encoders. In *Proceedings of the 22nd Annual Conference of the European Association for Machine Translation*, pages 175–184, Lisboa, Portugal, November 2020. European Association for Machine Translation.
- [11] Fabio Kepler, Jonay Trénous, Marcos Treviso, Miguel Vera, and André F. T. Martins. OpenKiwi: An open source framework for quality estimation. In *Proceedings of the 57th Annual Meeting of the Association*

BIBLIOGRAPHY

- for Computational Linguistics: System Demonstrations*, pages 117–122, Florence, Italy, July 2019. Association for Computational Linguistics.
- [12] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online, July 2020. Association for Computational Linguistics.