

## University of Groningen

### One Model to Rule them All

Bjerva, Johannes

**IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.**

*Document Version*

Publisher's PDF, also known as Version of record

*Publication date:*  
2017

[Link to publication in University of Groningen/UMCG research database](#)

*Citation for published version (APA):*

Bjerva, J. (2017). *One Model to Rule them All: multitask and Multilingual Modelling for Lexical Analysis*. [Thesis fully internal (DIV), University of Groningen]. Rijksuniversiteit Groningen.

**Copyright**

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

**Take-down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

*Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.*

# **One Model to Rule them all**

**Multitask and Multilingual Modelling for Lexical Analysis**

Johannes Bjerva



university of  
 groningen



The work in this thesis has been carried out under the auspices of the Center for Language and Cognition Groningen (CLCG) of the Faculty of Arts of the University of Groningen.



Groningen Dissertations in Linguistics 164

ISSN: 0928-0030

ISBN: 978-94-034-0224-6 (printed version)

ISBN: 978-94-034-0223-9 (electronic version)

© 2017, Johannes Bjerva

Document prepared with  $\text{\LaTeX}2_{\epsilon}$  and typeset by pdf $\text{\LaTeX}$   
(Droid Serif and Lato fonts)

Cover art: *Cortical Columns*. © 2014, Greg Dunn

21K, 18K, and 12K gold, ink, and dye on aluminized panel.

Printed by Off Page ([www.offpage.nl](http://www.offpage.nl)) on G-print 115g paper.



rijksuniversiteit  
 groningen

# **One Model to Rule them all**

## **Multitask and Multilingual Modelling for Lexical Analysis**

### **Proefschrift**

ter verkrijging van de graad van doctor aan de  
Rijksuniversiteit Groningen  
op gezag van de  
rector magnificus prof. dr. E. Sterken  
en volgens besluit van het College voor Promoties.

De openbare verdediging zal plaatsvinden op  
donderdag 7 december 2017 om 14.30 uur

door

**Johannes Bjerva**

geboren op 21 maart 1990  
te Oslo, Noorwegen

**Promotor**

Prof. dr. ing. J. Bos

**Copromotor**

Dr. B. Plank

**Beoordelingscommissie**

Prof. dr. A. Søgaaard

Prof. dr. J. Tiedemann

Prof. dr. L. R. B. Schomaker

# Acknowledgements

This has been a bumpy ride, to say the least, and having reached the end of this four-year journey, I owe a debt of gratitude to all of my friends, family, and colleagues. Your support and guidance has certainly helped smooth out most of the ups and downs I've experienced.

First of all, I would like to thank my PhD supervisors. Johan, the freedom you allowed me during the past four years has been one of the things I've appreciated the most in the whole experience. This meant that I could pursue the track of research I felt was the most interesting, without which writing a whole book would have been much more arduous – thank you! Barbara, thank you for agreeing to join as co-supervisor so late in my project, and for putting in so much time during the last couple of months of my PhD. The thesis would likely have looked quite different if you hadn't started in Groningen when you did. I owe you a huge thanks, especially for the final weeks – reading and commenting on the whole thesis in less than 24 hours during your vacation. Just, wow!

Next, I would like to thank Anders Søggaard, Jörg Tiedemann, and Lambert Schomaker for agreeing to form my thesis assessment committee. I feel honoured that you took the time to read this thesis, and that you deemed it scientifically sound. I also want to thank everyone who I have collaborated with throughout these years, both in Groningen and in Stockholm. Thanks Calle for agreeing to work with a sign language novice such as myself. Raf, it was an enlightening

experience to work with you and see the world of ‘real’ humanities research. Robert, we should definitely continue with our one-week shared task submissions (with various degrees of success).

Most of the last few years were spent at the computational linguistics group in Groningen. I would especially like to thank all of my fellow PhD students throughout the years. Thanks Kilian, Noortje, Valerio, and Harm for welcoming me with open arms when I joined the group. Special thanks to Kilian for being so helpful with answering all of my billions and billions of questions involved in finishing this thesis. Also, a special thanks to both Noortje and Harm for the times we shared when I had just moved here (especially that first New Year’s eve!). Hessel, thanks for all the help with administrative matters while I was abroad, especially for sending a gazillion of travel declarations for me. Rik and Anna, it was great getting to know you both better during the last few months – hopefully you find the bookmark to be sufficiently sloth-y. Dieke and Rob, thanks for being such great laid back drinking buddies and travel companions. To all of you, and Pauline – I hope we will continue the tradition of going all out whenever I come to visit Groningen! A big thanks to the rest of the computational linguistics group, especially Gertjan, Malvina, Gosse, John, Leonie, and Martijn. Also thanks to all other PhD students who started with me: Luis, Simon, Raf, Ruben, Aynur, Jonne, and everyone else whose name I’ve failed to mention. A special thanks to Ben and Kaja - I hope we keep up our board-game centred visits to one another, no matter where we happen to pursue our careers.

I spent most of my final year in Stockholm, and it was great being in that relaxed atmosphere during one of the more intense periods of the past few years. Most of all, I am sincerely grateful to Calle, to Johan (and Klara, Iris, and Vive), and to Bruno. Your support is truly invaluable, and by my lights, there is not much more to say than *<dank>*:(since a pal’s always needed). I’d also like to thank Johan and

Calle especially for agreeing to observing weird Dutch traditions by being my paranymphs. Josefina, Elísabet, and David, thank you all for being there for me during the past year. Thanks to the computational linguistics group for hosting me during this period, especially Mats, Robert, Kina, and Gintarė. Finally, thanks to everyone else at the Department of Linguistics at Stockholm University.

Having moved to Copenhagen this autumn has been an extremely pleasant experience. I would like to thank the entire CoAStal NLP group for simply being the coolest research group there is. Especially, I'd like to thank Isabelle for making the whole process of moving to Denmark so easy. Thanks to Anders, Maria, Mareike, Joachim, Dirk, and Ana for being so welcoming. I'd also like to thank everyone else at the image section at DIKU, especially the running buddies at the department, and most especially Kristoffer and Niels.

Finally, I would like to thank my family. This bumpy ride would have been challenging to get through without their support through everything. Kiitos, mamma! Takk/tack Aksel, Paulina, Julian, och Lucas! Takk Olav, og takk Amanda!

Let's see where the journey goes next!

Copenhagen, November 2017



# Contents

<b>Contents</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Chapter guide . . . . .	4
1.2 Publications . . . . .	6
<b>I Background</b>	<b>11</b>
<b>2 An Introduction to Neural Networks</b>	<b>13</b>
2.1 Introduction . . . . .	14
2.2 Representation of NNs, terminology, and notation . . .	15
2.3 Feed-forward Neural Networks . . . . .	17
2.3.1 Feature representations . . . . .	22
2.3.2 Activation Functions . . . . .	23
2.3.3 Learning . . . . .	24
2.4 Recurrent Neural Networks . . . . .	32
2.4.1 Long Short-Term Memory . . . . .	36
2.4.2 Common use-cases of RNNs in NLP . . . . .	39
2.5 Convolutional Neural Networks . . . . .	44
2.5.1 Local receptive fields . . . . .	45
2.5.2 Weight sharing . . . . .	46
2.5.3 Pooling . . . . .	48
2.6 Residual Networks . . . . .	50
2.7 Neural Networks and the Human Brain . . . . .	51
2.8 Summary . . . . .	53
<b>3 Multitask Learning and Multilingual Learning</b>	<b>55</b>
3.1 Multitask Learning . . . . .	56

---

3.1.1	Non-neural Multitask Learning . . . . .	57
3.1.2	Neural Multitask Learning . . . . .	58
3.1.3	Effectivity of Multitask Learning . . . . .	62
3.1.4	When MTL fails . . . . .	63
3.2	Multilingual Learning . . . . .	63
3.2.1	Human Annotation . . . . .	65
3.2.2	Annotation Projection . . . . .	65
3.2.3	Model Transfer . . . . .	66
3.2.4	Model Transfer with Multilingual Input Representations . . . . .	67
3.2.5	Continuous Space Word Representations . . . . .	69
3.3	Outlook . . . . .	76
<b>II</b>	<b>Multitask Learning</b>	<b>79</b>
<b>4</b>	<b>Multitask Semantic Tagging</b>	<b>81</b>
4.1	Introduction . . . . .	82
4.2	Semantic Tagging . . . . .	83
4.3	Method . . . . .	87
4.3.1	Inception model . . . . .	88
4.3.2	Deep Residual Networks . . . . .	90
4.3.3	Modelling character information and residual bypass . . . . .	90
4.3.4	System description . . . . .	92
4.4	Evaluation . . . . .	95
4.4.1	Experiments on semantic tagging . . . . .	96
4.4.2	Experiments on Part-of-Speech tagging . . . . .	96
4.4.3	The Inception architecture . . . . .	96
4.4.4	Effect of pre-trained embeddings . . . . .	98
4.5	Discussion . . . . .	98
4.5.1	Performance on semantic tagging . . . . .	98
4.5.2	Performance on Part-of-Speech tagging . . . . .	99
4.5.3	Inception . . . . .	100

4.5.4	Residual bypass . . . . .	100
4.5.5	Pre-trained embeddings . . . . .	101
4.6	Conclusions . . . . .	101
<b>5</b>	<b>Information-theoretic Perspectives on Multitask Learning</b>	<b>103</b>
5.1	Introduction . . . . .	104
5.2	Information-theoretic Measures . . . . .	105
5.2.1	Entropy . . . . .	105
5.2.2	Conditional Entropy . . . . .	105
5.2.3	Mutual Information . . . . .	107
5.2.4	Information Theory and MTL in NLP . . . . .	107
5.3	Data . . . . .	108
5.3.1	Morphosyntactic Tasks . . . . .	109
5.3.2	Semantic Tasks . . . . .	111
5.4	Method . . . . .	113
5.4.1	Architecture and Hyperparameters . . . . .	113
5.4.2	Experimental Overview . . . . .	114
5.4.3	Replicability and Reproducibility . . . . .	114
5.5	Results and Analysis . . . . .	115
5.5.1	Morphosyntactic Tasks . . . . .	115
5.5.2	Language-dependent results . . . . .	118
5.5.3	Semantic Tasks . . . . .	118
5.6	Conclusions . . . . .	118
<b>III</b>	<b>Multilingual Learning</b>	<b>123</b>
<b>6</b>	<b>Multilingual Semantic Textual Similarity</b>	<b>125</b>
6.1	Introduction . . . . .	126
6.2	Cross-lingual Semantic Textual Similarity . . . . .	128
6.3	Method . . . . .	131
6.3.1	Multilingual word representations . . . . .	131
6.3.2	System architecture . . . . .	132
6.3.3	Data for Semantic Textual Similarity . . . . .	135

---

6.4	Experiments and Results . . . . .	135
6.4.1	Comparison with Monolingual Representations .	136
6.4.2	Single-source training . . . . .	137
6.4.3	Multi-source training . . . . .	138
6.4.4	Results on SemEval-2017 . . . . .	140
6.4.5	Results on SemEval-2016 . . . . .	141
6.5	Conclusions . . . . .	142
<b>7</b>	<b>Comparing Multilinguality and Monolinguality</b>	<b>143</b>
7.1	Introduction . . . . .	144
7.2	Semantic Tagging . . . . .	146
7.2.1	Background . . . . .	146
7.2.2	Data . . . . .	146
7.2.3	Method . . . . .	147
7.2.4	Experiments and Analysis . . . . .	148
7.2.5	Summary of Results on Semantic Tagging . . . .	154
7.3	Tagging Tasks in the Universal Dependencies . . . . .	154
7.3.1	Data . . . . .	154
7.3.2	Method . . . . .	155
7.3.3	Results and Analysis . . . . .	157
7.3.4	Summary of Results on the Universal Depend- encies . . . . .	158
7.4	Morphological Inflection . . . . .	160
7.4.1	Method . . . . .	160
7.4.2	Results and Analysis . . . . .	162
7.4.3	Summary of Results on Morphological Inflection	164
7.5	Estimating Language Similarities . . . . .	164
7.5.1	Data-driven Similarity . . . . .	165
7.5.2	Lexical Similarity . . . . .	166
7.5.3	Results and Analysis . . . . .	167
7.5.4	When is Multilinguality Useful? . . . . .	170
7.6	Conclusions . . . . .	171

<b>IV Combining Multitask and Multilingual Learning</b>	<b>173</b>
<b>8 Multitask Multilingual Learning</b>	<b>175</b>
8.1 Combining Multitask Learning and Multilinguality . . .	176
8.2 Data . . . . .	179
8.2.1 Labelled data . . . . .	179
8.2.2 Unlabelled data . . . . .	180
8.3 Method . . . . .	181
8.3.1 Architecture . . . . .	181
8.3.2 Hyperparameters . . . . .	181
8.4 Experiments and Analysis . . . . .	182
8.5 Discussion . . . . .	193
8.6 Conclusion . . . . .	195
<b>V Conclusions</b>	<b>197</b>
<b>9 Conclusions</b>	<b>199</b>
9.1 Part II - Multitask Learning . . . . .	199
9.2 Part III - Multilingual Learning . . . . .	201
9.3 Part IV - Combining Multitask Learning and Multilin- guality . . . . .	202
9.4 Final words . . . . .	203
<b>Appendices</b>	<b>205</b>
<b>A Correlation figures for all languages in Chapter 5</b>	<b>207</b>
<b>B Bibliographical abbreviations</b>	<b>211</b>
<b>Bibliography</b>	<b>213</b>
<b>Summary</b>	<b>239</b>
<b>Samenvatting</b>	<b>243</b>