

Steffen Remus, Gerold Hintz, Darina Benikova, Thomas Arnold, Judith Eckle-Kohler, Christian M. Meyer, Margot Mieskes, and Chris Biemann www.aiphes.tu-darmstadt.de

Introd	uction EmpiriST	POS Parsing Results -		
Scenario: Process German web data and computer mediated communication data (social media, CMC) Training Web: 5,109 Data: CMC: 6,034 Tokens + POS tags Task: 1. Tokenization 2. POS tagging		Tokenization:Tokenization:Genre Rec Prec F1 Rank CMC 99.30 98.62 98.62 98.62Colspan="2">Colspan="2"Tokenization:Genre Rec Prec F1 Rank CMC 99.30 98.62 98.69 99.76 2WEB 99.63 99.89 99.76 2WEB 99.61 98.99 99.76 2 <th <="" colspan="2" td=""></th>		
1. Conservativ category cha 'poster A0' => 'poster	Fokenization e splits at Unicode ange positions, e.g. = 'ILZsLuNd' r', '', 'A', '0'	1.6. one reature vector per token test tark. vypes 4 5. simple shallow features such as word position and casing. Reuse regex rulesets from tokenization vypes 4 1 6. Unicode categories of characters AKW 60 Tiger was added for training, where simple STTS tag transformation rules were applied 40000		
2. Lookahead - Compiled - Contains - Contains 	<pre>merge list: from Wikipedia abbreviations, emotioons, etc. ad-list.txt merge rules: regular expressions Merge rules (+) e Merge rules (+) e Merge rules (-) trives det critices det cr</pre>	Error Analysis Common Tokenization Errors: • Rules are underspecified • linget: • die • festerschulung • die • linget: * die • Rules are overspecified Common Tagging Errors: • Rules are overspecified Error class confusion of function word tags 22 13.6 • Mistagged MZ as NN & vice-versa 20 12.5 • Insittagged MZ as NN & vice-versa 20 12.5 • Insittagged MZ as NN & vice-versa 20 12.5 • In order panetuation tag 20 12.5 • In order panetuation		

https://github.com/AIPHES

 Freely available, Open Source, Permissive Apache V2 License References: Michael Beißwenger, Sabine Bartsch, Stefan Evert, and Kay-Michael Würzner. 2016. Empirist 2015: A shared task on the automatic linguistic annotation of computer-mediated communication, social media and web corpora. In Proceedings of the Research Timining Group "Adaptive Preparation of Infor- mation from Informational Conference of the German Society for Computational Linguistics and Language Technology (GSCL- 2015), Essen, Germany. 		Further Information			Acknowledgments	\square
	 Freely available, Open Source, Perm References: Michael Be Michael Be Beinsenger, Sabine Battsch, Stef automatic linguistic annotation of comput 10th Web as Corpus Workshop (WAC-X) Be Darina Benikova, Sef Muhie Yimam, and In International Conference of the German S Germany. 	issive Apache V2 License an Evert, and Kay-Michael Würzner. 20 er-mediated communication, social med lin, Germany. Thris Biemann. 2015. GermaNER: Free of horis Jiemann. 2015. GermaNER: Free of lociety for Computational Linguistics and i	116. Empirist 2015: A shared task on the ia and web corpora. In Proceedings of the pen German named entity recognition tool. Language Technology (GSCI- 2015), Essen,	Thi Germa Re P Hetei grant Instit	is work has been supported by the an Research Foundation as part of the esearch Training Group "Adaptive reparation of Infor- mation from rogeneous Sources" (AIPHES) under No. GRK 1994/1 and by the German ute for Educational Research (DIPF) under the KDSL program.	

Steffen Remus, Gerold Hintz, Darina Benikova, Thomas Arnold, Judith Eckle-Kohler, Christian M. Meyer, Margot Mieskes, and Chris Biemann





WAC-X @ ACL 2016, Berlin, Germany

ΤE

UNIVER

DARMST



WAC-X @ ACL 2016, Berlin, Germany









Steffen Remus, Gerold Hintz, Darina Benikova, Thomas Arnold, Judith Eckle-Kohler, Christian M. Meyer, Margot Mieskes, and Chris Biemann www.aiphes.tu-darmstadt.de





Steffen Remus, Gerold Hintz, Darina Benikova, Thomas Arnold, Judith Eckle-Kohler, Christian M. Meyer, Margot Mieskes, and Chris Biemann www.aiphes.tu-darmstadt.de





Steffen Remus, Gerold Hintz, Darina Benikova, Thomas Arnold, Judith Eckle-Kohler, Christian M. Meyer, Margot Mieskes, and Chris Biemann www.aiphes.tu-darmstadt.de

Introd	uction EmpiriST	POS Parsing Results -		
Scenario: Process German web data and computer mediated communication data (social media, CMC) Training Web: 5,109 Data: CMC: 6,034 Tokens + POS tags Task: 1. Tokenization 2. POS tagging		Tokenization:Tokenization:Genre Rec Prec F1 Rank CMC 99.30 98.62 98.62 98.62Colspan="2">Colspan="2"Tokenization:Genre Rec Prec F1 Rank CMC 99.30 98.62 98.69 99.76 2WEB 99.63 99.89 99.76 2WEB 99.61 98.99 99.76 2 <th <="" colspan="2" td=""></th>		
1. Conservativ category cha 'poster A0' => 'poster	Fokenization e splits at Unicode ange positions, e.g. = 'ILZsLuNd' r', '', 'A', '0'	1.6. one reature vector per token test tark. vypes 4 5. simple shallow features such as word position and casing. Reuse regex rulesets from tokenization vypes 4 1 6. Unicode categories of characters AKW 60 Tiger was added for training, where simple STTS tag transformation rules were applied 40000		
2. Lookahead - Compiled - Contains - Contains 	<pre>merge list: from Wikipedia abbreviations, emotioons, etc. ad-list.txt merge rules: regular expressions Merge rules (+) e Merge rules (+) e Merge rules (-) trives det critices det cr</pre>	Error Analysis Common Tokenization Errors: • Rules are underspecified • linget: • die • festerschulung • die • linget: * die • Rules are overspecified Common Tagging Errors: • Rules are overspecified Error class confusion of function word tags 22 13.6 • Mistagged MZ as NN & vice-versa 20 12.5 • Insittagged MZ as NN & vice-versa 20 12.5 • Insittagged MZ as NN & vice-versa 20 12.5 • In order panetuation tag 20 12.5 • In order panetuation		

https://github.com/AIPHES

 Freely available, Open Source, Permissive Apache V2 License References: Michael Beißwenger, Sabine Bartsch, Stefan Evert, and Kay-Michael Würzner. 2016. Empirist 2015: A shared task on the automatic linguistic annotation of computer-mediated communication, social media and web corpora. In Proceedings of the Research Timining Group "Adaptive Preparation of Infor- mation from Informational Conference of the German Society for Computational Linguistics and Language Technology (GSCL- 2015), Essen, Germany. 		Further Information			Acknowledgments	
	 Freely available, Open Source, Perm References: Michael Be Michael Be Beinsenger, Sabine Battsch, Stef automatic linguistic annotation of comput 10th Web as Corpus Workshop (WAC-X) Be Darina Benikova, Sef Muhie Yimam, and In International Conference of the German S Germany. 	issive Apache V2 License an Evert, and Kay-Michael Würzner. 20 er-mediated communication, social med lin, Germany. Thris Biemann. 2015. GermaNER: Free of horis Jiemann. 2015. GermaNER: Free of lociety for Computational Linguistics and i	116. Empirist 2015: A shared task on the ia and web corpora. In Proceedings of the pen German named entity recognition tool. Language Technology (GSCI- 2015), Essen,	Thi Germa Re P Hetei grant Instit	is work has been supported by the an Research Foundation as part of the esearch Training Group "Adaptive reparation of Infor- mation from rogeneous Sources" (AIPHES) under No. GRK 1994/1 and by the German ute for Educational Research (DIPF) under the KDSL program.	





- Michael Beißwenger, Sabine Bartsch, Stefan Evert, and Kay-Michael Winzner. 2016. Empirist 2015: A shared task on the automatic linguistic anontation of computer-mediated communication, social media and web corpora. In Proceedings of the 10th Web as Corpus Workshop (WAC-X), Berlin, Germany.
 Darina Benikova, Seid Muhie Yiman, and Chris Biemann. 2015. GermanNER: Free open German named entity recognition tool.
- Dating Setukowa Setukowa in and an and satis behavior. 2013. Gentakeka. Pree Open Gental name entry recognition tool. In International Conference of the German Society for Computational Linguistics and Language Technology (GSL-2015), Essen, Germany.

Preparation of Infor- mation from

Heterogeneous Sources" (AIPHES) under

grant No. GRK 1994/1 and by the German

Institute for Educational Research (DIPF)

under the KDSL program.

TECHNISCHE UNIVERSITÄT

DARMSTADT

AIPHES