



Sarah Schulz

Senior Software Development Engineer, NLP

Born on: 20th August 1987

Place of Birth: Göppingen

Nationality: German

Experience

- 10/2018 – **Senior AI Engineer**, *Amboss GmbH*, Berlin, Medical Knowledge Modelling.
01/2020
- 10/2018 – **Senior Development Engineer, NLP**, *Ada Health GmbH*, Berlin, Medical Knowledge Acquisition.
01/2020
- 01/2018 – **Postdoctoral Researcher**, *University of Stuttgart*, Stuttgart, Institute for Natural Language Processing (IMS).
09/2018
- 05/2015 **PhD student**, *University of Stuttgart*, Stuttgart, Institute for Natural Language Processing (IMS).
–12/2017
- 2013–2015 **PhD student**, *Ghent University*, Ghent, Applied Linguistics – Language and Translation Technology Team (LT3).

Education

- 2010–2013 **Master of Arts**, *Eberhard-Karls-Universität*, Tübingen, 1,3.
International Studies in Computational Linguistics
- 2007–2010 **Bachelor of Arts**, *Friedrich-Alexander-Universität*, Erlangen-Nürnberg, 1,7.
Studies in Theatre and Media Science and German
- 2004–2007 **Abitur**, *Fritz-Ruoff-Schule*, Nürtingen, 1,4.
- 1998–2004 **Mittlere Reife**, *Realschule*, Neuffen, 1,8.

PhD thesis

- title **The Taming of the Shrew - Non-Standard Text Processing in the Digital Humanities**
- supervisors Prof. Dr. Jonas Kuhn
- description Digital Humanities as a research field that offers a variety of texts that are different from the standard on a lexical and syntactic level, challenges NLP to develop more flexible ways for automatic processing.

Master thesis

title **On Tokenization** - *Automatic extraction of semantic multiword units for meaning-based tasks*

supervisors Prof. Dr. Walt Detmar Meurers & PD. Dr. Frank Richter

description An approach is presented to automatically find meaning units in a text to enable semantic tokenization. It uses statistical methods borrowed from multiword expression extraction and combines it with a new approach based on unity checking of meaning of phrases.

Bachelor thesis

title **Zwischen Latein und Volkssprache** - *Der Deutsche Ptolemäus auf dem Weg zu einer vielschichtigen digitalen Edition durch die Implementierung von XLink*

supervisors Prof. Dr. Mechthild Habermann & Prof. Dr. Günther Görz

description In this thesis, the digital edition of the Early New High German text “Der deutsche Ptolemäus” was extended by extra-information about the connection between the German and Latin vocabulary. This information was included in a glossary that was implemented with the help of Xlink.

Professional experience

since **AI Engineer**, Amboss GmbH, Berlin.

03/2020 Information architecture of medical knowledge
Task:

- conceptualization of medical knowledge
- integration of external resources into an existing knowledge space

10/2018 – **Senior Software Development Engineer NLP**, Ada Health GmbH, Berlin.

01/2020 Knowledge Acquisition from unstructured text
Task:

- processing of biomedical text
- annotation of data
- machine learning

01/2018– **Postdoctoral Researcher**, University of Stuttgart - Institute for Natural Language

09/2018 *Processing*, Stuttgart.
Postdoctoral Researcher - project: Center for Reflected Text Analytics (CRETA)
Task:

- Digital Humanities
- interdisciplinary work
- machine learning

05/2015– **PhD student**, University of Stuttgart - Institute for Natural Language Processing ,

12/2017 Stuttgart.
PhD student - project: Center for Reflected Text Analytics (CRETA)
Task:

- processing of mainly historical text in a Digital Humanities context
- interdisciplinary work
- machine learning

12/2014– **Intern**, The Swiss AI Lab - IDSIA (Istituto Dalle Molle di Studi sull'Intelligenza

01/2015 *Artificiale*), Lugano.
Research Internship Neural Networks and Language Processing

- 04/2013– **PhD student**, *Ghent University - Language and Translation Technology Team*,
04/2015 Ghent.
PhD student - project: Automatic Monitoring for Cyberspace Applications
Task:
 - automatic text normalization of user-generated content
 - machine translation
 - machine learning
- 03/2012– **intern**, *European Academy of Bolzano - Institute for Specialised Communication and*
08/2012 *Multilingualism*, Bolzano.
Intern at the project 'Korpus Südtirol'
Tasks:
 - acquisition and digitalization of data for Korpus Südtirol
 - actualization and modification of Wortwarte developed by Lothar Lemnitzer and adaption for South Tyrolean German
 - compilation of a web-based corpus for South Tyrolean German
- 2010 –2012 **student assistant**, *Department of Linguistics*, Tübingen.
Proceeding with the project of the German wordnet 'GermaNet'
- 2009–2010 **student assistant**, *Department of Artificial Intelligence*, Erlangen.
Assistent at a course about XML and digital editions
- 09/2009– **intern**, *Department of Artificial Intelligence*, Erlangen.
10/2009 Assistance at different projects
Tasks:
 - digital data acquisition and evaluation
 - preparation of a first test run of a navigation and recommendation system for mobile phones
- 2008–2009 **student assistant**, *Department of German Linguistics*, Erlangen.
Investigation of language change by reference to editions of 'Melusine'

Main research interests

- Digital Health
- Explainability of AI
- Machine learning
- Digital Humanities
- Non-Standard Text

Languages

German	fluent	<i>native language</i>
English	fluent	<i>9 years at school, master's program in English</i>
Dutch	fluent	<i>private lessons, 2 years in Belgium</i>
French	very basics	<i>3 years at school</i>
Latin	very basics	<i>2 semesters at university</i>
Italian	very basics	<i>1 semester at university, 6 months in South Tyrol</i>
Portuguese	very basics	<i>1 semester at community college</i>

Computer skills

Programming	<ul style="list-style-type: none"> ○ Python ○ Java ○ Bashscript ○ Gnu R ○ SQL 	<i>excellent</i> <i>basic</i> <i>basics</i> <i>basics</i> <i>basics</i>	Mark-up	<ul style="list-style-type: none"> ○ LaTeX ○ XML ○ HTML 	<i>excellent</i> <i>good</i> <i>basics</i>
Python libraries	<ul style="list-style-type: none"> ○ sklearn ○ nltk ○ pandas ○ keras 				
Agile: Scrum	<ul style="list-style-type: none"> ○ Git ○ Jira ○ Confluence ○ Notion 				
Office	<ul style="list-style-type: none"> ○ Linux ○ OpenOffice 				

Memberships, Reviewing and Participation

- Scientific Committee LREC 2020
- Participant at Dagstuhl Seminar 19452: Machine Learning Meets Visualization to Make Artificial Intelligence Interpretable
- Participant at Dagstuhl Seminar 17192: Human-Like Neural-Symbolic Computing
- Reviewer for Digital Humanities im deutschsprachigen Raum (DHd) 2018 in Köln
- Programme Committee 12th Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities (LaTeCH) 2018
- Scientific Committee LREC 2018
- Reviewer for Workshop CEx@AI*IA 2017 in Bari
- Reviewer for Digital Humanities (DH) 2017 in Montreal
- Programme Committee 11th Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities (LaTeCH) 2017
- Reviewer for Digital Humanities im deutschsprachigen Raum (DHd) 2017 in Bern
- Scientific Committee LREC 2016
- Digital Humanities (DH) 2016 in Krakow
- Programme Committee 10th Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities (LaTeCH) 2016
- Reviewer for Journal of Artificial Intelligence Research (JAIR)

Organizing/Managing

- Dagstuhl Seminar 19172: Computational Creativity meets Digital Humanities
- Co-Manager of Interdisciplinary College 2019 (13th to 20th of March 2020, Günne)
- Co-Manager of Interdisciplinary College 2019 (12th to 19th of March 2019, Günne)
- Co-Manager of Interdisciplinary College 2018 (9th to 16th of March 2018, Günne)
- Co-Manager of Interdisciplinary College 2017 (10th to 17th of March 2017, Günne)

Interests

sports pilates, hiking, mountain biking, running, climbing

Publications

- [1] Orphée De Clercq, Sarah Schulz, Bart Desmet, Els Lefever, and Véronique Hoste. Normalization of Dutch User-Generated Content. In *Proceedings of the 9th International Conference on Recent Advances in Natural Language Processing*, Hissar, Bulgaria, 2013.
- [2] Orphée De Clercq, Schulz Schulz, Bart Desmet, and Véronique Hoste. Towards Shared Datasets for Normalization Research. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, Reykjavik, Iceland, May 2014. European Language Resources Association (ELRA).
- [3] Marisa Delz, Benjamin Layer, Sarah Schulz, and Johannes Wahle. Overgeneralization of verbs — The change of the German verb system. In *Proceedings of the 9th International Conference on the Evolution of Language*, Evolang IX, pages 96–103, Kyoto, Japan, 3 2012.
- [4] Bart Desmet, Orphée De Clercq, Marjan Van de Kauter, Sarah Schulz, Cynthia Van Hee, and Veronique Hoste. *Taaltechnologie 2.0: sentimentanalyse en normalisatie*, pages 157–161. Beschouwingen uit een talenhuis : opstellen over onderwijs en onderzoek in de vakgroep Vertalen, Tolken en Communicatie aangeboden aan Rita Godyns. Academia Press, 2014.
- [5] D. Doran, S.C. Schulz, and T. R. Besold. What does explainable ai really mean? a new conceptualization of perspectives. *CEUR Workshop Proceedings*, 2071, March 2018. Copyright © 2018 for this paper by its authors. Copying permitted for private and academic purposes. Proceedings of the First International Workshop on Comprehensibility and Explanation in AI and ML 2017 co-located with 16th International Conference of the Italian Association for Artificial Intelligence (AI*IA 2017) Bari, Italy, November 16th and 17th, 2017.
- [6] Janis Pagel, Nils Reiter, Ina Rösiger, and Sarah Schulz. A Unified Annotation Workflow for Diverse Goals. In Sandra Kübler and Heike Zinsmeister, editors, *Proceedings of the Workshop on Annotation in Digital Humanities, co-located with ESSLLI 2018*, August 2018.
- [7] Nils Reiter, Sarah Schulz, Gerhard Kremer, Roman Klinger, Gabriel Viehhauser, and Jonas Kuhn. Teaching Computational Aspects in the Digital Humanities Program at University of Stuttgart – Intentions and Experiences. In *Proceedings of the Workshop on Teaching NLP for Digital Humanities (Teach4DH 2017) co-located with GSCL 2017*, pages 43–48, Berlin, Germany, September 2017.
- [8] Ina Roesiger, Sarah Schulz, and Nils Reiter. Towards coreference for literary text: Analyzing domain-specific phenomena. In *Proceedings of the Second Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*, pages 129–138, Santa Fe, New Mexico, August 2018. Association for Computational Linguistics.

- [9] Sarah Schulz. Named-Entity Recognition for User-Generated Content. In *Proceedings of European Summer School in Logic Language and Computation 2014 Student Session*. Springer, 2014.
- [10] Sarah Schulz. *The Taming of the Shrew - non-standard text processing in the Digital Humanities*. PhD thesis, University of Stuttgart, 2018.
- [11] Sarah Schulz and Mareike Keller. Code-switching ubiquie est - language identification and part-of-speech tagging for historical mixed text. In *Proceedings of the 10th SIGHUM Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities*, pages 43–51, Berlin, Germany, August 2016. Association for Computational Linguistics.
- [12] Sarah Schulz and Nora Ketschik. From 0 to 10 million annotated words: part-of-speech tagging for Middle High German. *Language Resources and Evaluation*, 53(4):837–863, Dec 2019.
- [13] Sarah Schulz and Jonas Kuhn. Learning from within? comparing PoS tagging approaches for historical text. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 4316–4322, Portorož, Slovenia, May 2016. European Language Resources Association (ELRA).
- [14] Sarah Schulz and Jonas Kuhn. Multi-modular domain-tailored OCR post-correction. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2716–2726, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.
- [15] Sarah Schulz, Verena Lyding, and Lionel Nicolas. Compiling a diverse web corpus for South Tyrolean German - STirWaC. In *Proceedings of the 8th Web as Corpus Workshop*, pages 37–45, Lancaster, UK, 2013.
- [16] Sarah Schulz, Guy De Pauw, Orphée De Clercq, Bart Desmet, Véronique Hoste, Walter Daelemans, and Lieve Macken. Multimodular text normalization of dutch user-generated content. *ACM Trans. Intell. Syst. Technol.*, 7(4), July 2016.