

Studienarbeit

Syxtus Gaal

Automatisierte Erstellung neuer Sprachkorpora

Ein Beispiel anhand des Lëtzebuergeschen



Bachelor + Master
Publishing

Syxtus Gaal

Automatisierte Erstellung neuer Sprachkorpora: Ein Beispiel anhand des Lëtzebuergeschen

Originaltitel der Studienarbeit: Automatische phonetische Annotation - ein HMM-basierter Aligner für das Lëtzebuergesche

ISBN: 978-3-86341-642-3

Herstellung Bachelor + Master Publishing, ein Imprint der Diplomica® Verlag GmbH, Hamburg, 2012

Zugl. Universität Stuttgart, Stuttgart, Deutschland, Studienarbeit, April 2008

Dieses Werk ist urheberrechtlich geschützt. Die dadurch begründeten Rechte, insbesondere die der Übersetzung, des Nachdrucks, des Vortrags, der Entnahme von Abbildungen und Tabellen, der Funksendung, der Mikroverfilmung oder der Vervielfältigung auf anderen Wegen und der Speicherung in Datenverarbeitungsanlagen, bleiben, auch bei nur auszugsweiser Verwertung, vorbehalten. Eine Vervielfältigung dieses Werkes oder von Teilen dieses Werkes ist auch im Einzelfall nur in den Grenzen der gesetzlichen Bestimmungen des Urheberrechtsgesetzes der Bundesrepublik Deutschland in der jeweils geltenden Fassung zulässig. Sie ist grundsätzlich vergütungspflichtig. Zuwiderhandlungen unterliegen den Strafbestimmungen des Urheberrechtes.

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Die Informationen in diesem Werk wurden mit Sorgfalt erarbeitet. Dennoch können Fehler nicht vollständig ausgeschlossen werden, und die Diplomarbeiten Agentur, die Autoren oder Übersetzer übernehmen keine juristische Verantwortung oder irgendeine Haftung für evtl. verbliebene fehlerhafte Angaben und deren Folgen.

© Bachelor + Master Publishing, ein Imprint der Diplomica® Verlag GmbH, Hamburg, 2012

<http://www.diplom.de>, Hamburg 2012
Printed in Germany

Hiermit erkläre ich, dass ich die vorliegende Arbeit selbständig verfasst habe und dabei keine andere als die angegebene Literatur verwendet habe.

Alle Zitate und sinngemäßen Entlehnungen sind als solche unter genauer Angabe der Quelle gekennzeichnet.

Danksagung

Es sind die Ideen mehrerer Personen, die diese Arbeit mitgeprägt haben. Peter Gilles führte mich in die lëtzebuergesche Phonetik ein, betreute fachlich diese Arbeit sowie stellte eine Sammlung wertvoller Sprachaufnahmen bereit. Stefan Rapp und Antje Schweizer stellten das Programm zur Verfügung, auf dem diese Arbeit aufbaut und opfertens hilfsberet ihre Zeit, um mich darin einzuarbeiten. Wolfgang Wokurek hat einen sinnvollen Skopus dieses Projekts definiert. Die Korrekturen von Nils Herda trugen zu einer klaren Strukturierung, sowie einer hohen Qualität des Textes bei. Okko Buss ist Autor des Transkriptionsprogramms, das ich während der Korpuserstellung verwendet habe. Wojciech Przystas machte mich bereits beim ersten Entwurf auf potentielle Fehlerquellen aufmerksam und erteilte wertvolle L^AT_EX-Tipps. Diese Arbeit verwendet z.T. seine L^AT_EX-Codefragmente.

Bei allen diesen Personen möchte ich mich zutiefst bedanken.

Inhaltsverzeichnis

1	Einleitung	5
1.1	Motivation	5
1.2	Aufgabenstellung	5
2	Lëtzebuergesch	6
2.1	Geschichte Luxemburgs	6
2.2	Luxemburgisch	7
2.3	Dialekte	8
2.4	Koiné	8
2.5	Phonetik der luxemburgischen Koiné	8
2.6	Phonetik des Deutschen	9
2.7	Gegenüberstellung des Deutschen und des Lëtzebuergesch	11
3	Phonetische Alignierung als Erkennungsproblem	12
3.1	Toolgestützte Spracherkennung mit dem Aligner	12
3.2	Funktionsweise	12
3.3	Phoneminventar und Aussprachelexikon	13
3.4	Abbildung der lëtzebuergesch Phoneme auf deutsche Sprachlaute	14
4	Aufbereitung der Sprachdaten	14
4.1	Vorhandene Sprachdaten	14
4.2	Aufteilung der Datensätze	16
4.3	Vorbereitung der Daten - Gilles-Korpus	16
4.3.1	Orthographische Transkription	17
4.3.2	Erstellung eines Aussprachelexikons	18
4.3.3	Generierung der phonetischen Transkriptionen	19
4.3.4	Korpusstruktur	21
4.4	Vorbereitung der Daten - das 6000-Wieder-Korpus	22
4.4.1	Bereinigung des Datensatzes	22
4.4.2	Aufteilung für Test und Training	23
4.5	Zusammenfassung	23