

The DiDi Project:
**Collecting, Annotating, and Analysing South
Tyrolean CMC Data**

ird-cmc-rennes: Social Media and CMC Corpora for the eHumanities

egon w. stemle<egon.stemle@eurac.edu>



Institute for Specialised Communication
and Multilingualism
EURAC, Bozen/Bolzano, Italy

EURAC
research

October 23rd, 2015

- **Aivars Glaznieks**
Writing research, Language acquisition, Corpus linguistics, Sociolinguistics
- **Jennifer-Carmen Frey**
Technology Enhanced Learning and Media Literacy, Language Acquisition, Corpus Linguistics, Sociolinguistics
- **Myself**
Cognitive Science, Computational Linguistics, Artificial Intelligence
- **Sabrina Galasso, Nicole Stuckey**
interns (for around 3 months, each)



- 1 The Project DiDi
- 2 Data Acquisition
- 3 Data Annotation, Processing and Access
- 4 Outlook

- DiDi: Digital Natives and Digital Immigrants
Writing on social network sites: a corpus-based observation of the current language use in South Tyrol, with particular consideration of the writers' age.
- Period: 01.06.2013 - 31.07.2015

AUTONOME
PROVINZ
BOZEN
SÜDTIROL



PROVINCIA
AUTONOMA
DI BOLZANO
ALTO ADIGE

The project was financed by *the Autonomous Province of Bozen/Bolzano, South Tyrol, Department for the Promotion of Education, Universities and Research, Provincial Law 13.Dec.2006, No. 14 'Research and Innovation'*.

Research Questions

- 1 How do South Tyrolean users of SNS (with L1 German) use German in quasi-public and private communication?
- 2 How do they use other languages?
- 3 Are there differences in language use that can be explained by age? numerical age (young vs. old) versus digital age (novice vs. experienced user)?

Finding Answers

- 1 Construct a corpus of South Tyrolean language use (and sociolinguistic metadata)
- 2 Analyse the data



Spotted: Südtirol

September 2

I hon gestern 1.09 pa hockey in sterzing a gruppe gitschn gsegn de sein do trainingsloger glab i und oane fa de hot an grauen pulli unkop, i hatse gerne gsehn konnmer jemand helfn?

See Translation

Like · Comment · Share

Hannes Rainer likes this.

Top Comments -



Write a comment...



Pass sellm besser afn luki au bitte:D

See Translation

Like · Reply · 3 · September 2 at 11:20pm via mobile



Nimm wose kregn knsh Lukas Tötsch

See Translation

Like · Reply · 2 · September 2 at 11:20pm via mobile



Jo mindigshnts 3 😊

Like · Reply · 2 · September 2 at 11:17pm via mobile



Na enk foln sochn in 😊

Like · Reply · 2 · September 2 at 11:14pm via mobile



Wirtschaftskrise:D

See Translation

Like · Reply · 1 · September 2 at 11:18pm via mobile



Zu sein geburtstog werter zuastechn;)

See Translation

Data and more Data

What and more What

- SMS, e-mail, chat room, Face- book, Twitter, online forum
- meta and sociolinguistic data

Participants

Pick, Reach, Convince

- selection criteria
- print, phone, radio/TV, online
- kind of incentive

Legal and Ethical Issues

- legal situation with (private) data
- adhere to ethical standards

didi Digital Natives und Digital Immigrants

EURAC research
Itallano deutsch

Wie schreibt Südtirol auf Facebook?

Im Rahmen unseres Projektes DiDi wollen wir Forscherinnen und Forscher des **EURAC-Instituts für Fachkommunikation und Mehrsprachigkeit** die Sprachgewohnheiten der Südtirolerinnen und Südtiroler auf Facebook untersuchen.

Besonders interessiert uns dabei...

- Verwendung von unterschiedlichen Dialekten Südtirols
- Verwendung von Deutsch, Italienisch und anderen Sprachen
- Schreibstil von jüngeren und älteren Usern
- Verwendung von Smiley's, Abkürzungen, etc.
- Groß- und Kleinschreibung
- andere Stimmfärbungen

Dafür benötigen wir Südtirolerinnen und Südtiroler, die uns ihre Facebook-Texte (natürlich vollständig anonym) zur Verfügung stellen und bereit sind, ein paar Fragen zu Ihren Nutzungsgewohnheiten im Internet zu beantworten. Das Ausfüllen des Fragebogens wird max. 10 Minuten dauern.

Die Ergebnisse der Untersuchung werden ab Mai 2015 auf www.eurac.edu/didi veröffentlicht.

Weitere Informationen
Möchtest du mehr über das Projekt erfahren? Besuche unsere [Webseite!](#)

Verwendung der Daten:
Mit der Teilnahme an der Untersuchung werden einmalig deine in Facebook gespeicherten Sprachdaten (je nach Auswahl Pinwandbeiträge und/oder Inbox-Messages) des Jahres 2013 abgerufen. Dabei lesen wir nur deine eigenen Beiträge aus sowie die Anzahl der sprachlichen Verteilung der Folgekommentare. Beiträge deiner Freunde werden nicht auslesen. Die Daten werden nur für wissenschaftliche Zwecke verwendet. Es geht uns dabei ausschließlich darum, die Sprache und den Stil der Kommunikation zu untersuchen, der Inhalt der Gespräche wird hierbei außer Acht gelassen. Die Daten werden zudem größtenteils automatisch mit Hilfe spezieller computerlinguistischer Software ausgewertet. (Mehr dazu in unseren [Datenschutzbestimmungen](#).)

Anonymisierung der erhobenen Daten:
Die Datenerhebung erfolgt vollständig anonym. Allen Teilnehmenden am Forschungsprojekt wird eine ID-Nummer zugewiesen, die keinen Rückschluss auf die einzelnen Teilnehmenden erlaubt. Alle Daten werden mit Hilfe der ID-Nummer in anonymisierter Form gespeichert und für weitere Projektschritte zur Verfügung gestellt. In den Sprachdaten werden alle Hinweise auf die Verfasserin bzw. den Verfasser und andere Personen (d.h. Personen- und Ortsnamen, E-Mail-Adressen und andere Kontaktadressen) anonymisiert (d.h. durch willkürliche Namen ersetzt). Im Fragebogen wird weder der Name, das Geburtsdatum oder der Geburtsort noch der genaue Wohnort erfragt, eine eindeutige Identifizierung der teilnehmenden Person ist damit nicht möglich.

Für die Untersuchung möchte ich folgende Daten aus dem Jahr 2013 zur Verfügung stellen:

- Pinwandbeiträge
- eigene Mitteilungen aus meiner Inbox

Teilnehmen

Autonome Provinz Bozen - Südtirol

Provincia Autonoma di Bolzano Alto Adige

Das Projekt wird finanziert von der Autonomen Provinz Bozen - Südtirol, Abteilung Bildungsförderung, Universität und Forschung, Landesgesetz vom 13. Dezember 2006, Nr. 14 "Forschung und Innovation".

DiDi

Deutsch ▼

Persönliche Daten

1. Sind Sie weiblich oder männlich?

- weiblich
 männlich

2. Wann wurden Sie geboren?

Jahr:

3. Wo haben Sie ihren Lebensmittelpunkt?

- in Südtirol
 außerhalb Südtirols

4. Welche Sprache ist Ihre erste Sprache (Muttersprache)?
(Falls Sie mehrere Muttersprachen haben sollten, geben Sie bitte diese an!)

- Deutsch
 Italienisch
 Ladinisch
 andere



Zurück

Weiter

Powered by
[Opinio Survey Software](#)

DiDi

Deutsch ▼

Persönliche Daten

5. Sprechen Sie einen Südtiroler Dialekt?

- nein
 ja



Zurück

Weiter

Powered by
[Opinio Survey Software](#)

DIDI

Deutsch ▼

Persönliche Daten

6. Welchem Gebiet kann man Ihren Dialekt am ehesten zuordnen?

Orte/Täler ▼

- Orte/Täler
- Bozen
- Meran
- Eisacktal
- Etschtal
- Grödnertal
- Passeiertal
- Pustertal
- Samtal
- Überetsch - Unterland
- Ultental
- Vinschgau
- Wipptal

Zurück Weiter

Powered by
[Opinio Survey Software](#)

Solbai, Ziofungo, und *lol*...

Vielen Dank für die Teilnahme!

Aktuelle Informationen zum Projektverlauf finden sich auf www.eurac.edu/didi. Dort werden ab Mai 2015 auch die Ergebnisse des Projekts veröffentlicht.

Bitte hilf uns die Didi-App in Südtirol bekannt zu machen!

Das Gelingen unseres neuartigen Projekts die Sprache auf Facebook in Südtirol festzustellen hängt stark davon ab, dass wir die Daten die wir von dir bekommen haben auch mit denen von anderen vergleichen können. Dabei ist es wichtig, dass so viele Südtirolerinnen und Südtiroler wie möglich am Projekt teilnehmen. Du kannst uns dabei helfen und das Projekt deinen Bekannten empfehlen Z.B. indem du das App auf deiner Facebook-Pinnwand teilst!

**Südtirolerinnen und
Südtiroler haben bereits
teilgenommen!**

**Teile die Didi-App mit
deinen Freunden!**

The DiDi Facebook App

Facebook Wall Post and Comments

Egon W Stemle updated his profile picture.
April 18, 2011

- Public
- Friends of Friends
- Friends
- Friends except Acquaintances
- Only Me
- Custom
- Close Friends
- University of Trento
- See all lists...

Like Comment Share

Wieviel sind auf dich ausgesetzt? Lohnt sichs schon?
April 18, 2011 at 8:13am · Like

wie geil.
April 18, 2011 at 9:19am · Like

Egon W Stemle der preis ist heiss: ich hab' mir vor langer zeit zwei paar extra beisser in die kieferaschen gesteckt; wie sich nun herausgestellt hat, nicht allzu clever - hat einem damals aber natürlich keiner gesagt! zwei gib'ts noch zu hohlen...
April 18, 2011 at 10:04am · Like

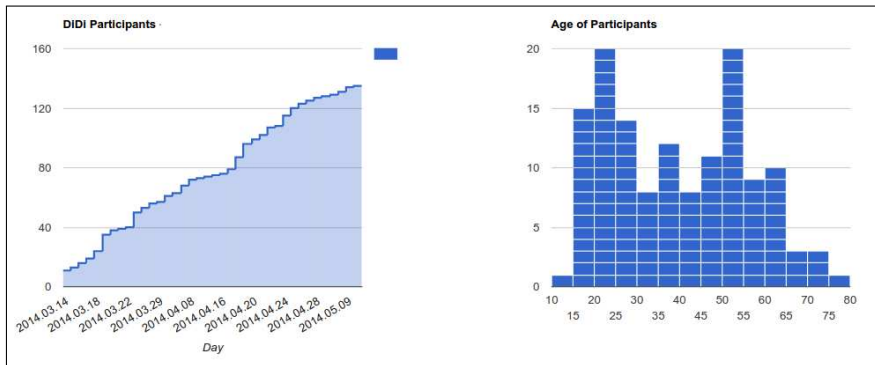
Egon W Stemle (ist ja schwer zu erkennen: also die oberen beiden sind noch da...)
April 18, 2011 at 10:06am · Like

au au au. wie das doch einen gleich verändert, nicht?
April 18, 2011 at 10:11am · Like

Write a comment...

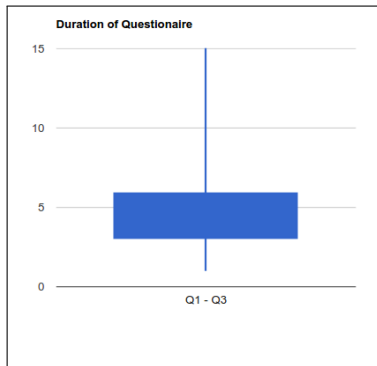
Monitoring of Data Acquisition Progress

Number of Participants and their Age



Monitoring of Data Acquisition Progress

Sanity Check: Duration of the Questionnaire



Automatic Processing

Tokenisation *ttokenizer*

POS Tagging, Lemmatisation *IMS TreeTagger*

Manual Annotation

Normalisation

WHY: corpus search, POS tagging (pretest: 0.5 -> 0.9), dialect identification

Anonymisation

indicate and classify issues in text

POS Corrections

(most frequent) OOVs, South Tyrolean lexemes

Semi-Automatic Post-Processing

Anonymisation

remove metadata (ids, attachments) and other users' data

CMC Phenomena

emoticons, @mentions, iterations, . . .

Language and Dialect ID

langid.py+manual post-processing, heuristic-driven dialect annotation

POS Corrections

CMC, anonymised tokens

It seems that perfection is reached not when there is nothing left to add, but when there is nothing left to take away.

A. de Saint Exupéry

original	anonym	token_corr	comments	norm1	norm2
<comment id=1234_4405601>					
3		3,5			
,		---			
5m		---			
		m			
sind					
die					
Betonsäulen					
Rudi	p				
</comment>					
<comment id=2345_4407568>					
...					
sel				_stir	
muasch				musst	
grod				gerade	
du					
sogn				sagen	
lssi	p				
!!!!					
</comment>					

/ 135 Displaying Results 1 - 10 of 1350 Result for: stir=/*/

4 Path: didi > 54625_10201437841060438 (normtok 3 - 17) left context: right context: 7

grid

p	0					2									
token	Südtirol	Online	(@stol_it);	"	wenn	lei	die	Hälfte	von	die	Versprechen	wird	wohr
normtok	Südtirol	Online	(@stol_it);	"	wenn	lei	die	Hälfte	von	den	Versprechen	wird	wahr
lemma	Südtirol	_unknown_	(_unknown_	_unknown_	"	wenn	lei	die	Hälfte	von	die	Versprechen	werden	wahr
pos	NE	NE	\$(NN	XY	\$(KOUS	FM	ART	NN	APPR	ART	NN	VAFIN	ADJD
stir															
cmc				at_organisation	emoticon										

norm

[0] Retweeted Stol Südtirol Online (@stol_it):

[2] " wenn lei die Hälfte von den Versprechen wird wahr , geht es uns gut In den nächsten 5 Jahren " - Messner-Windschnur #Wahlkampfblue ...

orig

[0] Retweeted Stol Südtirol Online (@stol_it):

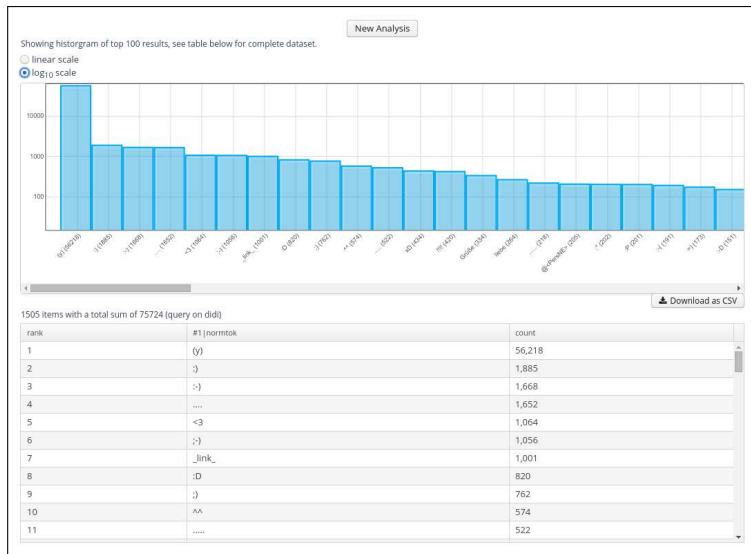
[2] " wenn lei die Hälfte von die Versprechen wird wohr , geats ins guat in die nächsten 5 johr " - Messner-Windschnur #Wahlkampfblue ...

nfo for salt:/didi/54625_10201437841060438 + x

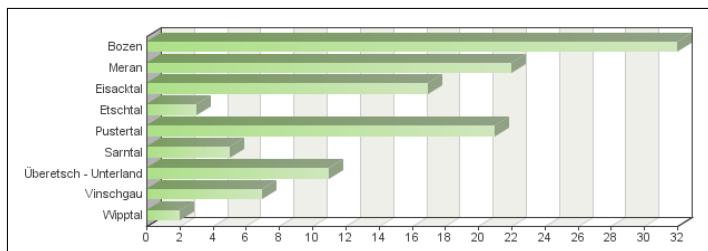
Metadata

document: 54625_10201437841060438

Name	Value
PA_Alter	36
PA_Ausbildungsabschluss	Oberschule mit Matura
PA_Beruf	ArbeitnehmerIn
PA_Dialektsprecher_ITA	no data
PA_Dialektsprecher_STIR	Bozen
PA_Geburtsjahr	1977
PA_Geschlecht	m
PA_L1_Deutsch	1
PA_L1_Italienisch	0
PA_L1_Ladinisch	0
PA_L1_andere	0
PA_Lebensmittelpunkt_STIR	1
PA_Schule	no data
RespondentId	54625
StartDate	2014.03.14 12:24
ZA_Frequenz_Blog	min. 1x pro Monat



Can we find known features of spoken South Tyrolean German in our CMC data?



- **Lexis**
nicht (not) [net] or [nit] or [it, et] (W to E)
- **Vocalism**
-er as [r], [ər] or [o] (W to E)
ei as [ua] or [oa] (W to E)
ge- as [gə], [gi] or [ga] (N to S)

- Follow-up research project
“Factors of Language Choice and Language Change”
- Use revised STTS POS tag set
- Export corpus in TEI format
- Further data analysis

Project Homepage

<http://www.eurac.edu/didi/>

ANNIS Corpus Interface

<http://commul.eurac.edu/annis/didi/>



Jennifer-Carmen Frey, Aivars Glaznieks, and Egon W. Stemle.

The DiDi Corpus of South Tyrolean CMC Data.

In *Proceedings of the 2nd Workshop on Natural Language Processing for Computer-Mediated Communication / Social Media at GSCL2015 (NLP4CMC2015)*, Essen, September 2015. German Society for Computational Linguistics & Language Technology.

<https://sites.google.com/site/nlp4cmc2015/NLP4CMC-2015.pdf>.



Jennifer-Carmen Frey, Egon W. Stemle, and Aivars Glaznieks.

Collecting language data of non-public social media profiles.

In Gertrud Faaß and Josef Ruppenhofer, editors, *Workshop Proceedings of the 12th Edition of the KONVENS Conference*, pages 11-15, Hildesheim, Germany, October 2014. Universitätsverlag Hildesheim, Germany.

<http://www.uni-hildesheim.de/konvens2014/data/konvens2014-workshop-proceedings.pdf>.



Aivars Glaznieks and Egon Stemle.

Challenges of building a CMC corpus for analyzing writer's style by age: The DiDi project.

Journal for Language Technology and Computational Linguistics (JLCL), 29(2):31-57, December 2014.

http://www.jlcl.org/2014_Heft2/2GlaznieksStemle.pdf.