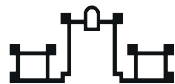


Proceedings

of the 10th International
Conference on CMC and
Social Media Corpora for
the Humanities



September 14 -15, 2023
University of Mannheim
Germany



UNIVERSITY
OF MANNHEIM
School of Humanities

IDS

LEIBNIZ-INSTITUT FÜR
DEUTSCHE SPRACHE

Funded by

DFG

Deutsche
Forschungsgemeinschaft
German Research Foundation

Louis Cotgrove
Laura Herzberg
Harald Lüngen
Ines Pisetta (eds.)

Proceedings of the 10th International Conference on CMC and Social Media Corpora for the Humanities

14–15 September 2023, University of Mannheim, Germany

Editors: Louis Cotgrove, Laura Herzberg, Harald Längen, Ines Pisetta

Published by: Leibniz-Institut für Deutsche Sprache
Mannheim, 2023

DOI: <https://doi.org/10.14618/1z5k-pb25>

ISBN: 978-3-937241-95-1

This work is licensed under a Creative Commons “Attribution 4.0. International” license.

Conference website: <https://www.uni-mannheim.de/cmc-corpora2023>

The CMC-Corpora 2023 conference is funded by *Deutsche Forschungsgemeinschaft* under the project number 524949653.

Preface

Following the successes of the ninth conference in 2022 held in the wonderful Santiago de Compostela, Spain, we are pleased to present the proceedings of the 10th edition of International Conference on CMC and Social Media Corpora for the Humanities (CMC-2023). The focal point of the conference is to investigate the collection, annotation, processing, and analysis of corpora of computer-mediated communication (CMC) and social media.

Our goal is to serve as the meeting place for a wide variety of language-oriented investigations into CMC and social media from the fields of linguistics, philology, communication sciences, media studies, and social sciences, as well as corpus and computational linguistics, language technology, textual technology, and machine learning.

This year's event is the largest so far with 45 accepted submissions: 32 papers and 13 poster presentations, each of which were reviewed by members of our ever-growing scientific committee. The contributions were presented in five sessions of two or three streams, and a single poster session. The talks in these proceedings cover a wide range of topics, including the corpora construction, digital identities, digital knowledge-building, digitally-mediated interaction, features of digitally-mediated communication, and multimodality in digital spaces.

As part of the conference, we were delighted to include two invited talks: an international keynote speech by Unn Røynealand from the University of Oslo, Norway, on the practices and perceptions of researching dialect writing in social media, and a national keynote speech by Tatjana Scheffler from the Ruhr-University of Bochum on analysing individual linguistic variability in social media and constructing corpora from this data. Additionally, participants could take part in a workshop on processing audio data for corpus linguistic analysis. This volume contains abstracts of the invited talks, short papers of oral presentations, and abstracts of posters presented at the conference.

We wish to thank all colleagues who contributed to the conference and proceedings this year for their fascinating and varied presentations, posters, and keynote talks. We would also like to thank the members of the international scientific committee for their support and help in reviewing the many submissions this year. Thanks also go to the Leibniz-Institute for the German Language and the University of Mannheim for providing administrative support and the wonderful locations for the conference this year, and a big thank you to the German Research Foundation (DFG) for their financial contribution to the conference.

We hope that the tenth edition of the conference series can build on the successes of the previous editions and we are looking forward to the next decade of CMC-Corpora conferences!

Mannheim, September 14 2023

On behalf of the organising committee

Jutta Bopp,
Louis Cotgrove,
Laura Herzberg,
Harald Lungen, and
Andreas Witt

Committees

Local Organizing Committee in Mannheim

Jutta Bopp	IDS Mannheim
Louis Cotgrove	IDS Mannheim
Laura Herzberg	University of Mannheim
Harald Lüngen	IDS Mannheim
Andreas Witt	University of Mannheim & IDS Mannheim

International Steering Committee of the Conference series

Steven Coats	University of Oulu
Julien Longhi	Cergy-Pontoise Université
Lieke Verheijen	Radboud University
Reinhild Vandekerckhove	University of Antwerp

Scientific Committee

Paul Baker	Lancaster University
Adrien Barbaresi	Berlin-Brandenburgische Akademie der Wissenschaften
Michael Beißwenger	University of Duisburg-Essen
Mario Cal Varela	Universidade de Santiago de Compostela
Steven Coats	University of Oulu
Luna DeBruyne	Ghent University
Orphée DeClercq	Ghent University
Francisco Javier Fernández Polo	University of Santiago de Compostela
Jenny Frey	EURAC Research Bolzano
Alexandra Georgakopoulou-Nunes	King's College London
Klaus Geyer	University of Southern Denmark
Aivars Glaznieks	EURAC Research Bolzano
Jan Gorisch	IDS Mannheim
Claire Hardaker	Lancaster University
Iris Hendrickx	Radboud University Nijmegen
Axel Herold	Berlin-Brandenburgische Akademie der Wissenschaften
Lisa Hilde	University of Antwerp
Mai Hodac	Université Toulouse
Wolfgang Imo	University of Hamburg
Paweł Kamocki	IDS Mannheim
Erik-Tjong Kim-Sang	Netherlands eScience Center
Alexander Koenig	CLARIN ERIC
Florian Kunneman	Vrije Universiteit Amsterdam
Marc Kupietz	IDS Mannheim
Els Lefever	Ghent University
Julien Longhi	Cergy-Pontoise Université
Maja Miličević-Petrović	University of Bologna
Nelleke Oostdijk	Radboud University

Céline Poudat
Thomas Proisl
Ines Rehbein
Sebastian Reimann
Unn Røyne land
Jan Oliver Rüdiger
Müge Satar
Tatjana Scheffler
Stefania Spina
Egon Stemle
Caroline Tagg
Simone Ueberwasser
Reinhild Vandekerckhove
Lieke Verheijen

Université Côte d'Azur
Friedrich-Alexander-Universität Erlangen-Nürnberg
University of Mannheim
Ruhr-Universität Bochum
University of Oslo
IDS Mannheim
Newcastle University
Ruhr-Universität Bochum
Università per Stranieri di Perugia
EURAC Research Bolzano
The Open University
University of Zurich
University of Antwerp
Radboud University

Contents

Keynotes

Unn Røyneland: Eye dialect in social media – practices and perceptions	1
Tatjana Scheffler: Individual linguistic variability in social media	2

Posters

Aminat Babayode, Laurens Bosman, Nicole Chan, Katharina Ehret, Ivan Fong, Noelle Harris, Alissa Hewton, Danica Reid, Maite Taboada and Rebekah Wong: Structural properties of podcasts as an emerging register of computer-mediated communication	3
Marie-Louise Bartsch and Irina Mostovaia: Ellipsis of the subject pronoun ich ('I') in German WhatsApp chats: A usage-based approach	7
Elizaveta Kibisova: Building corpora of Russian fake and genuine news for linguistic analysis	9
Aenne Knierim: Tracing Perceptions of Black History by Comparison of Two Corpora	10
Katharina Pabst, Aida Alanzi, Johanna Aminoff, Raisa Tayib and Derek Denis: Zooming in on emerging norms: Preliminary findings from a cross-linguistic investigation of videoconferencing	13
Sebastian Reimann, Lina Rodenhausen, Tatjana Scheffler and Frederik Elwert: ChrisTof: A Novel Corpus of Christian Online Forums	15
Hannah J. Seemann, Sara Shahmohammadi, Tatjana Scheffler and Manfred Stede: Building a Parallel Discourse-annotated Multimedia Corpus	17
Sarah Steinsiek: Negotiating knowledge in cooperative learning scenarios: a multimodal approach to practices of computer-mediated and face-to-face communication in the university classroom	18

Jenia Yudytska: "Linguistic features, device affordances, and contextual factors: A mixed-methods, two-corpora approach""	20
Yinglei Zang: Deontic Authority in Computer-mediated Communication Between University Teachers and Students: A Comparative Study of German and Chinese	21
Talks	
Selenia Anastasi, Tim Fischer, Florian Schneider and Chris Biemann: IDA - Incel Data Archive: a multimodal comparable corpus for exploring extremist dynamics in online interaction.	23
Tianyi Bai: The Reply Function in WhatsApp Chat Communication	29
Michael Beißwenger, Eva Gredel, Lena Rebhan and Sarah Steinsiek: Ellipsis Points in Messaging Interactions and on Wikipedia Talk Pages	33
Laura Bothe: The representation of the 'Jew' as enemy in French public Telegram channels within the identitarian-conspiratorial milieu	39
Bruno Machado Carneiro, Michele Linardi and Julien Longhi: "Studying Socially Unacceptable Discourse Classification (SUD) through different eyes: ""Are we on the same page ?""	45
Steven Coats: A Pipeline for the Large-Scale Acoustic Analysis of Streamed Content	51
Louis Cotgrove: "megageil, mega geil, and voll mega: Intensification in YouTube comments""	55
Selcen Erten: Exploring register variation in Turkish web corpus	60

Annamaria Fabian and Igor Trost,:	
"Digital Corpus Linguistic Analysis of the Language of Inclusion, Discrimination and Exclusion of people with disability in social media – in a German corpus of 214.926 Tweets on #disability and #inclusion between 2007-2023""	
.....	65
Anne Ferger, Andre Frank Krause and Karola Pitsch:	
Workflows and Methods for Creating Structured Corpora of Multimodal Interaction	
.....	73
Francisco Javier Fernández Polo:	
Balancing expert and peer-student identities in online discussion forums	
.....	78
Shir Finkelstein and Hadar Netz:	
"‘Hebrew level: Bibist.’: Online Hebrew language corrections as a tool for ‘civilized’ bashing	
.....	83
Carolina Flinz, Eva Gredel and Laura Herzberg:	
A Corpus study on the negotiation of pronominal address on talk pages of the German, French, and Italian Wikipedia	
.....	86
Florian Frenken:	
A Multivariate Register Perspective on Reddit: Exploring Lexicogrammatical Variation in Online Communities	
.....	91
Prakhar Gupta, Lliana Doudot, Romain Loup and Aris Xanthos:	
Collecting and de-identifying half a million WhatsApp messages	
.....	96
Teemu Helenius,:	
Acquiring, Analyzing, and Understanding Multimodal TikTok Short Video Data: The Case of Online Sex Worker Visibility Management	
.....	102
Sangwan Jeon:	
MigrTwit Corpora. (Im)Migration Tweets of French Politics.	
.....	108
André Frank Krause, Anne Ferger and Karola Pitsch:	
Anonymization of Persons in Videos of Authentic Social Interaction: Machine Learning Model Selection and Parameter Optimization.	
.....	112
Lothar Lemnitzer and Antonia Hamdi:	
"‘Also ehrlich’ – From adjectival use to interactive discourse marker	
.....	118

Rosa Lorés: The recontextualization of expert knowledge: intertextual patterns in digital science dissemination	124
Harald Lungen and Laura Herzberg: Studying the distribution of reply relations in Wikipedia talk pages	131
Martti Mäkinen: MMWAH! Compiling a Corpus of Multilingual / Multimodal WhatsApp Discussions by Swedish-speaking Young Adults in Finland	136
Rachel McCullough, Daniel Drylie, Mindi Barta and Daniel Smith: CoDEC-M: the multi-lingual Manosphere subcorpus of the Corpus of Digital Extremism and Conspiracies	140
Iliia Moshnikov and Eugenia Rykova: Little Big Data: Karelian Twitter Corpus	142
Anastasiia Piroh: Multimodal Intertextual Practices in Video Film Reviews	148
Ana Eugenia Sancho-Ortiz: Scientific communication on social media: Analysing Twitter for knowledge recontextualisation	154
Ulrike Schneider and Oliver Watteler: Can I Publish my Social Media Corpus? Legal Considerations for Data Publication	160
Laurel Stvan: Collecting Health Memes for a Subcorpus of Peer Health Discourse	166
Ludovic Tanguy, Céline Poudat and Lydia-Mai Ho-Dac: Specific behaviours in Wikipedia talk pages: some insights from extreme cases	171
Ralia Thoma: “Don’t be afraid of Greeklisch”: Adolescent students’ transliteration practices	176
Eva Triebel: Not an expert, but not a fan either. A corpus-based study of negative self-identification as epistemic index in web forum interaction.	182

Reinhild Vandekerckhove, Sarah Bernolet, Astrid De Wit and Tanja Mortelmans: Towards a more inclusive approach of digital literacy: social media writing at an older age	187
Jiayi Zhou: Phonetic Metaphor of Chinese Emojis: An Approach of Neologism Formation	190
Author index	194
Keyword index	195

A Corpus study on the negotiation of pronominal address on talk pages of the German, French, and Italian Wikipedia

Carolina Flinz¹, Eva Gredel², Laura Herzberg³

¹Università degli Studi di Milano, ²University of Duisburg-Essen, ³University of Mannheim
E-mail: carolina.flinz@unimi.it, eva.gredel@uni-due.de, herzberg@uni-mannheim.de

Abstract

The adequate use of social deixis is highly dependent on the situation and context and has therefore always been at the center of linguistic pragmatics. So far, principles of pronominal address have mainly been modelled with a focus on oral, co-present interaction. The use of pronominal address in computer-mediated communication with its translocal and partially anonymous contexts is still a research gap. In this context, it is particularly interesting that different digital platforms have developed specific customs or netiquettes regarding the appropriate use of address pronouns. This paper asks, from a contrastive perspective, how the appropriate use of address pronouns is negotiated on talk pages of the German, French, and Italian Wikipedia. The corpus study is based on the multilingual Wikipedia corpora of the Leibniz Institute for the German Language.

Keywords: Social Deixis, Corpus Linguistics, Wikipedia

1. Introduction¹

The appropriate use of socio-deictic signs is highly dependent on the situation and context and has always been at the center of linguistic pragmatics (cf. Nübling et al. 2017: 205). However, principles of pronominal address have so far been mainly modelled with a focus on oral interaction where speakers are co-present (cf. Kretzenbacher 2010). The use of pronominal address in computer-mediated communication (CMC) with its translocal and (partially) anonymous contexts poses special challenges for writers and has been considered in only a few initial studies (cf. Gredel 2023). This paper aims to fill this research gap by analyzing meta-discourses on pronominal address in the CMC genre of Wikipedia talk pages. With the multilingual Wikipedia corpora of the Leibniz Institute for the German Language, digital language resources are used that allow a contrastive approach to this object of investigation.

The languages German, French, and Italian², which are considered in this paper, each have a binary system of pronominal address containing an “intimate” pronoun (GER: *du*, FR: *tu*, IT: *tu*) and a more “formal” pronoun (German: *Sie*, FR: *vous*, IT: *Lei*). In oral face-to-face interaction, the selection of the appropriate pronoun in each communicative dyad is generally linked to variables such as social status, age, gender, and conversation situation of the interaction partners (cf. Nübling et al. 2017, 205). In CMC, these variables are not always apparent to writers, so the selection of appropriate pronouns must follow other principles. This corpus study focuses on this aspect through the analysis of meta-discourses.

Regarding CMC, it is interesting that different customs or netiquettes for the use of the appropriate address pronouns

have developed on various digital platforms (cf. Gredel 2023). On the multilingual Wikipedia, there are differences between the netiquettes of the considered language versions. However, there is no consensus on these netiquettes, and they are subject to controversial discussions. Based on the Wikipedia corpora of the Leibniz Institute for the German Language, this article explores whether and how writers negotiate the use of address pronouns on Wikipedia talk pages. It also analyses which aspects of the use of pronominal address are being discussed on talk pages of the German, French, and Italian Wikipedia.

2. Social Deixis

The concept of socio-deixis focused on here, which is often traced back to Fillmore (1975, 76), is characterized by Levinson as follows: “Social deixis concerns the encoding of social distinctions that are relative to participant-roles, particularly aspects of the social relationship holding between speaker and addressee(s)” (cf. Levinson 1983, 63). Central aspects of pronominal address, which in some cases pose communicative dilemmas for interaction partners, are the nature and timing of the transition from the distance form to the familiar form (cf. Mühlhäusler & Harré 1990, 142f.; cf. Simon 2003, 125). Based on corpus data, it can be shown that the unidirectionality of this transition from the distance form (in German: *Sie*) to the familiar form (in German: *du*) in CMC is not always given (cf. Gredel 2023). In specific situations or in the context of transitions from digital communication to oral interaction, the unidirectionality may be suspended in a communicative dyad, as shown in Example (1): User „Iste“ mentors new authors in Wikipedia. On his user talk page, he is asked by a new author for advice on editing Wikipedia pages and is

¹ The three authors have written the paper jointly. Carolina Flinz is responsible for the data and analyses of Italian, Eva Gredel for German, Laura Herzberg for France. Introduction (§1) and Conclusion and Outlook (§5) were written jointly.

² Social deixis can be expressed in Italian by the pronouns of second person singular *tu*, *ti* and plural *voi*, *vi*, or of third person singular *Lei*, *Le* and plural

Loro (cf. Milano 2015: 70). The ones used in contemporary Italian are *tu* and *Lei*. The use of *Loro* addressing more than one person is rare while the form *Voi* was used mainly in the past and during the fascist period. It has almost completely disappeared, even if it is being used nowadays, it's restricted to southern regional Italians and the ecclesiastic sphere (cf. Serianni 2000: 7).

addressed by this person with the formal pronoun *Sie*. In the example below, user “Iste” goes directly to the informal *du* by referring to the Wikipedia Netiquette (*WP:DU*), but marks it with an emoticon (*du ;-)* to mitigate a potential face threatening act (FTA, cf. Brown/Levinson 2007: 60-66). At the same time, he offers the other author, whom he is addressing as a newbie (*Neuling*), to return to the pronominal form of address with *Sie* if he is very uncomfortable with the informal form (*du*) of address:

- (1) Zunächst einmal ist es hier in der Wikipedia üblich, dass die Benutzer sich untereinander duzen (WP:DU); wenn du ;-) trotzdem gesiezt werden möchtest, dann sag einfach Bescheid. Grundsätzlich ist es durchaus möglich, auch als Neuling einen passablen Artikel zu schreiben, der nicht wieder gelöscht wird, wenn man sich vorher etwas eingelesen hat. (User Iste, 08.11.2011, WUD17/I65.44315)
First of all, it is usual here in Wikipedia for users to address each other by the pronoun du (WP:DU); if you ;-) would still like to be addressed by Sie, then just let me know. Basically, it is possible for a newbie to write a decent article that will not be deleted again if you have read up a bit beforehand.

In (1), the offer is made to suspend the unidirectionality of the transition from the formal to the informal pronoun in this communicative dyad – contrary to the Wikiquote.

A similar case is also found in the Italian corpus, in which the new user says that he has been used the informal form “tu” from the beginning and asks whether he operated correctly (2). The answer is as follows:

- (2) Per convenzione, wikipedia usa "di default" il tu per non distinguere gli utenti (ricordarsi chi vuole il lei e chi vuole il tu è impossibile, c'è solo da impazzirci) e anche per un discorso che, su wikipedia, non ci sono distinzioni di sorta, cioè un amministratore è importante tanto quanto un utente. Se vuoi che ti do del lei basta chiedere, ma nel lungo periodo dubito fortemente di ricordarmene (però farò uno sforzo). (User Aluong, 6 lug 2006, WUI15/A06.20577)
By convention, wikipedia uses 'by default' 'tu' so as not to distinguish between users (remembering who wants 'Lei' and who wants 'tu' is impossible, it's just going crazy) and also for a discourse that, on wikipedia, there are no distinctions whatsoever; i.e. an administrator is just as important as a user. If you want me to call you 'she', just ask, but in the long run I doubt very much that I will remember (but I will make an effort).

The preference for the informal form is due to two reasons: first, it is practical, i.e. not to confuse and not to distinguish between users, since everyone, both administrators and users are important equals.

The cases in which a user doesn't want the informal form, are usually highly marked and are often thematized metalinguistically, as demonstrated in Example (3) as well.

- (3) Bonjour. Pour commencer, je n'apprécie pas d'être tutoyée quand je ne connais pas. [...]

Bien à vous. --chansonnette [causer avec dame éliane]
 4 avril 2013 à 17:41 (CEST), WDF15/A31.85510
Hello. First of all, I don't appreciate being on first-name terms when I don't know someone. [...]
All the best.

In (3), the importance of using the appropriate socio-deictic sign is marked in a discussion page posting. User “chansonnette” starts off her posting with an explicit change of topic, also marked linguistically by “first of all” (*pour commencer*), to openly show her displeasure about the chosen form of addressing. The factor of “unfamiliarity” between her and another user as well as the associated custom of falling back to the formal *you* in similar interactions, e.g., in oral face-to-face interaction, are used as “chansonnette”'s arguments for preferring the formal *you* variant. Interestingly, she also finishes her posting by using “bien à vous” (*all the best*), a very formal form of wishing farewell, again with *vous* underlining her preference for formal addressing signs.

Previous work has also described contexts in which address pronouns are used not reciprocally or symmetrically, but asymmetrically (cf. Clyne et al. 2003). These aspects of pronominal address have predominantly been discussed for oral interaction. This paper is one of the first to use corpus linguistic methods to investigate which rules and principles can be empirically reconstructed for pronominal forms of address in CMC.

Kretzbacher (2010: 3) names a total of four methodological approaches to metalinguistic comments on the topic of social deixis: These are focus group discussions, network interviews, participant observation of uncontrolled public interactions, and observation and questioning of German-language Internet forums. From our point of view, corpus linguistic approaches are particularly suitable because in this way metalinguistic comments from hundreds or thousands of internet users can be considered, whereas in studies with the above-mentioned methodological approaches far fewer informants could be considered (in Kretzenbacher 2010 is n = 190). In the following, the data and the method will now be described in detail.

3. Method and Data

The talk pages of Wikipedia share features of CMC genres such as a dialogic structure and an informal writing style with non-standard language (cf. Storrer 2017). There are two types of Wikipedia talk pages, whose data are considered in this study based on the multilingual corpora by the Leibniz Institute for the German Language: article talk pages, where authors negotiate online encyclopedic content and user talk pages, where the contributions of individual authors are discussed. These two types of talk pages will be considered for the study. The metadata for the corpora used are as follows, cf. Table 1:

Language	Article talk pages	User talk pages
GER	373,161,686 (wdd17)	309,390,966 (wud17)
FR	138,068,162 (wdf15)	374,390,445 (wuf15)
IT	52,070,465 (wdi15)	130,067,969 (wui15)

Table 1: Size of the corpora in tokens and corpus abbreviations³ (DeReKo 2022 in COSMAS II 2022).

To be able to investigate meta-discourses and thus the negotiation of appropriate address pronouns, we use the following search strings when conducting queries in COSMAS II:

- GER: *&siezen* and *&duzen*
- FR: *vouvoyer* and *tutoyer*⁴
- IT: *dare del L/lei* and *dare del tu*⁵

Regarding the topic of social deixis and the relevant variables for choosing the appropriate address pronoun, it should be noted for the language data at hand that Wikipedia authors do not have to disclose their offline identity on Wikipedia. Against this background, guidelines have been developed for Wikipedia, which diverge depending on the language version. For example, in the German Wikipedia, it is recommended that authors generally use the informal *du* form of address (Wikipedia 2023a). In French, the usage of the formal *vous* form of address is greatly reflected upon. There are specific user boxes which can be implemented on a user page that indicate how a user wishes to be addressed (e.g., Wikipedia 2023b for *vouvoyer* boxes). Although there are also users who prefer the informal *tu*, the formal *vous* form of address

still plays a rather important role in Wikipedia user addressing. In several user surveys no consensus could be generated, so both forms of addressing continue to be used depending on a user’s preference (Wikipedia 2023c). Investigating whether, and if so, to what extent users explicitly address these conflicting priorities will shed light on the use of the appropriate address pronouns in multilingual CMC environments.⁶

In the Italian Wikipedia there are no explicit guidelines, but the preference of the informal *tu* is deductible. First, there is often the focus on “welcome and inclusion” (cf. Wikipedia 2023d) and secondly, all help and guide pages are written addressing the reader with the informal *tu*. Finally, in the box which summarizes the principles of good communication, the informal *tu* is used (cf. Wikipedia 2023e). In the talk pages (*Bar*), the preferences are addressed as well as the possibility of a survey; however, there seems to be a consensus on the informal *tu*, by leveraging on the fact that Wikipedia is a project between colleagues, and it allows everyone to feel equal, regardless of their social or cultural status (cf. Wikipedia 2023f).

The targeted corpus study focuses on this issue, which will be examined in more detail for each language in Section 4.

4. Negotiation of social deixis on Wikipedia talk pages

In the following, it will be quantitatively demonstrated to what extent corpus hits referring to a meta-discourse on social deixis can be found in the three languages under consideration. For the German language version, it can be shown that both corpora (wdd17 and wud17) contain hits for both search strings (*&siezen* and *&duzen*), cf. Table 2:

³ The corpus abbreviations read as follows, *wdd17* is the Wikipedia corpus of German (*deutsch*) article talk (*Diskussion*) pages created from a 2017 Wikipedia dump; *wud17* represents the user discussion pages.

⁴ All inflected forms were queried in a rather complex REG# (regular expression) search string:

```
#REG(^tuto(ie|nt|r(a(i(s(en)?t)s)?i?(ez|ons)ont)?s)?y(a(i(s(en)?t)?nt|s(s(e|nt)s)?i(ez|ons)))?)?â(mes|t(es)?é(es)?er|èrènt|i?(ez|ons)))$) oder
#REG(^vouvo(ie|nt|r(a(i(s(en)?t)s)?i?(ez|ons)ont)?s)?y(a(i(s(en)?t)?nt|s(s(e|nt)s)?i(ez|ons)))?)?â(mes|t(es)?é(es)?er|èrènt|i?(ez|ons)))$).
```

⁵ All inflected forms were queried in a rather complex REG# (regular expression) search string:

```
#REG(^d(â(nno)?a(i|n((d|n)o|te)|r(à|a(i|nno)|e(bbe(ro)?i|mm?o|st(e|i)te)?|ò)|t(a|e|i|o)|v(a(m|n)o|te|i)?i|o)?|e(mmo|s(s(e(ro)?i(mo)?))|st(e|i)tt(e(ro)?i))|i(a(m|n)?o|te)?e(d(e(ro)?i))|o|ò)$) /+w1 del /+w1 (tu oder lei oder Lei).
```

⁶ The multilingualism of the CMC environment *Wikipedia* with its approximately 300 language versions and the numerous interlanguage links between them has two relevant dimensions: On the one hand, there are frequent citations of other language versions or posts in other languages on the discussion pages. In addition, many authors who speak foreign languages do not only contribute to the language version of their mother tongue, but also edit in several language versions at the same time.

Language	Wikipedia name space	Corpus abbr.: search term	Occurrences	pMW ⁷	Texts
German	Talk page	wdd17: &#x73;#x69;#x65;#x70;#x65;#x72	322	0.86	208
		wdd17: &#x64;#x75;#x7a;#x65;#x6e	993	2.66	682
	User talk page	wud17: &#x73;#x69;#x65;#x70;#x65;#x72	395	1.21	290
		wud17: &#x64;#x75;#x7a;#x65;#x6e	2,052	6.29	1,659
French	Talk page	wdf15: v#x6f;#x75;#x76;#x6f;#x79;#x65;#x72	103	0.75	95
		wdf15: t#x75;#x6f;#x79;#x65;#x72	449	3.25	426
	User talk page	wuf15: v#x6f;#x75;#x76;#x6f;#x79;#x65;#x72	200	0.53	181
		wuf15: t#x75;#x6f;#x79;#x65;#x72	1,655	4.42	885
Italian	Talk page	wdi15: d#x61;#x72;#x65;#x20;#x64;#x65;#x6c;#x20;#x4c;/#x6c;#x65;#x69	29	0.56	9
		wdi15: d#x61;#x72;#x65;#x20;#x64;#x65;#x6c;#x20;#x74;#x75	61	1.17	61
	User talk page	wui15: d#x61;#x72;#x65;#x20;#x64;#x65;#x6c;#x20;#x4c;/#x6c;#x65;#x69	84	1.61	82
		wui15: d#x61;#x72;#x65;#x20;#x64;#x65;#x6c;#x20;#x74;#x75	372	7.14	308

Table 2: Results of the search queries in the Wikipedia corpora (DeReKo 2022 in COSMAS II 2022).

For the French language version, Table 2 shows that both search strings (*vouvoyer* and *tutoyer*) yield results which, however, differ in their frequency: For both Wikipedia name spaces, i. e. the article talk pages as well as the user talk pages, inflected forms of the informal address *tutoyer* are more frequently discussed than forms of the formal variant *vouvoyer*. It becomes clear that French Wikipedia authors debate the means of addressing with each other; in sum more often on their own talk pages than on the article talk pages.

The Italian language version contains hits for both search strings in both corpora, cf. Table 2. In particular *dare del tu* is more discussed than the formal form variant *dare del Lei/lei*.

5. Conclusion and outlook

In all three language versions of Wikipedia considered, there are indications that authors negotiate social deixis – and specifically the pronouns of address – in the sense of a meta-discourse. This contribution shows the extent to which there are differences between the three language versions of Wikipedia. When comparing all three languages, the frequencies of discussing socio-deictic signs meta-linguistically are significantly different between the

German, French and Italian Wikipedia language versions⁸. In both analysed Wikipedia subcorpora, i.e. the Wikipedia article talk pages on the one hand and the article talk pages on the other hand, a greater deal of discussions about addressing styles takes place on the user talk pages, with the informal *you* variant being discussed more frequently than formal *you* variant.

6. References

- Brown, Penelope/ Levinson, Stephen C. (2007): *Gesichtsbedrohende Akte*. In: Herrmann, Steffen/Krämer, Sybille/ Kuch, Hannes (eds.): *Verletzende Worte. Die Grammatik sprachlicher Missachtung*. Bielefeld: Transcript Verlag, 59–88.
- Clyne, Michael/Kretzenbacher, Heinz/Norby, Catrin/Warren, Jane: *Address in Some Western European Languages*. In: *Proceedings of the 2003 Conference of the Australian Linguistic Society*, 1–10.
- Clyne, Michael/Norrby, Catrin/Warren, Jane: *Language and human relations: styles of address in contemporary language*. Cambridge 2009.
- COSMAS I/II (2022): *Corpus Search, Management and Analysis System*, <http://www.ids-mannheim.de/cosmas2/>, ©1991-2022 Leibniz-Institut für Deutsche Sprache, Mannheim.
- DeReKo (2022): *Deutsches Referenzkorpus / Archiv der*

⁷ The abbreviation *pMW* stands for *occurrences per million words*. It is a measure of relative occurrence frequencies that are also normalized to a common base (one million current word forms). This allows for comparing frequencies in corpora of different sizes. To calculate pMW values, we need to divide the raw frequency by the total number of words in the corpus and multiply the result by one million.

⁸ This holds for testing between the three languages, with the chi-square statistic being 87.5197. The p-value is < 0.00001. The result is significant at $p < .05$ for comparing together the frequencies of the formal *you* variant as well as the informal *you* variant between the three languages, with the chi-square statistic being 61.361. The p-value is < 0.00001. The result is significant at $p < .05$, cf. <https://www.socscistatistics.com/tests/chisquare2/default2.aspx>. For each language, the differences in frequencies between the two analysed corpus types, i.e. Wikipedia article talk pages and user talk pages, are significant for the formal *you* variant in German and French, e.g. for the formal *you* variant in German, *Sie*, the difference between the name spaces is significant with the chi-square statistic being 27.5725. The p-value is < 0.00001. The result is significant at $p < .05$; French: The chi-square statistic is 7.6534. The p-value is .005667. The result is significant at $p < .05$; not for Italian: The chi-square statistic is 0.4735. The p-value is .491403. The result is not significant at $p < .05$.

- Korpora geschriebener Gegenwartssprache 2022-I (Release from 08.03.2022).
- Da Milano, Federica (2015): Italian. In: Jungbluth, Konstanze/Da Milano, Federica (eds.): *Manual of Deixis in Romance Languages*. Berlin/Boston: De Gruyter, 59–74.
- Fillmore, Charles J.: *Santa Cruz lectures on deixis* 1971. Mimeo, Bloomington 1975.
- Gredel, Eva (2023): Siezt du noch oder duzt du schon? In: Korpusstudie zum Gebrauch und zur Aushandlung sozial-deiktischer Zeichen auf digitalen Plattformen. In: Meier-Vieracker, S./ Bülow, L./ Marx, K./ Mroczynski, R. (Hg.): *Digitale Pragmatik*. Berlin/ Heidelberg: Metzler, 39–57.
- Levinson, Stephen C.: *Pragmatics*. Cambridge 1983.
- Mühlhäusler, Peter/Harré, Rom: *Pronouns and people: The linguistic construction of social and personal identity*. Oxford/Cambridge 1990.
- Nübling, Damaris/ Dammel, Antje/ Duke, Janet/ Szczepaniak, Renata: *Historische Sprachwissenschaft des Deutschen. Eine Einführung in die Prinzipien des Sprachwandels*. Tübingen 2013.
- Nübling, Damaris/ Dammel, Antje/ Duke, Janet/ Szczepaniak, Renata: *Historische Sprachwissenschaft des Deutschen. Eine Einführung in die Prinzipien des Sprachwandels*. Tübingen 2017.
- Serianni, Luca (2000): Gli allocutivi di cortesia. In: *La Crusca per voi*. N. 20 aprile 2000.
- Simon, Horst J.: Für eine grammatische Kategorie >Respekt<. *Synchronie, Diachronie und Typologie der deutschen Anredepronomina*. Tübingen 2003.
- Storrer, A. (2017). Grammatische Variation in Gespräch, Text und internetbasierter Kommunikation. In M. Konopka & A. Wöllstein (Eds.), *Grammatische Variation. Empirische Zugänge und theoretische Modellierung*, Berlin/New York: e Gruyter, pp. 105–125.
- Wikipedia 2023a: Warum sich hier alle duzen. https://de.wikipedia.org/wiki/Wikipedia:Warum_sich_hier_alle_duzen
- Wikipedia 2023b: Utilisateur vouvoient. https://fr.wikipedia.org/wiki/Cat%C3%A9gorie:Utilisateur_vouvoient
- Wikipedia 2023c: Vouvoyer ou tutoyer sur Wikipédia?. https://fr.wikipedia.org/wiki/Wikip%C3%A9dia:Sondage/Vouvoyer_ou_tutoyer_sur_Wikip%C3%A9dia%3F
- Wikipedia 2023d: Aiuto:Pagina di discussione https://it.wikipedia.org/wiki/Aiuto:Pagina_di_discussione
- Wikipedia 2023e: Wikipedia:Wikiquote <https://it.wikipedia.org/wiki/Wikipedia:Wikiquote>
- Wikipedia 2023f: Wikipedia:Bar/Discussioni https://it.wikipedia.org/wiki/Wikipedia:Bar/Discussioni/Per_favore:_Datemi_il_%22Lei%22