



# DATA MANAGEMENT PLAN

**SFB 1287 / 2021 / Phase1 / Project A03**



**SFB  
1287**

**Limits of Variability in Language**  
Cognitive, Grammatical, and Social Aspects

# TABLE OF CONTENTS

<b>General Information .....</b>	<b>3</b>
(1) Scholarly / Scientific Requirements .....	3
(2) Dataset.....	3
(3) Re-using .....	5
<b>Technical Information .....</b>	<b>6</b>
(4) Dataset Storage.....	6
(5) File Formats .....	7
<b>Publication .....</b>	<b>8</b>
(6) Publishing / Sharing.....	8
(7) Quality Measures.....	9
<b>Legal and Ethics .....</b>	<b>10</b>
(8) Personal Data .....	10
(9) Anonymisation .....	10
(10) Property Rights .....	11
<b>Storage and long-term preservation .....</b>	<b>12</b>
(11) Why?.....	12
(12) Archiving Reasons.....	12
(13) Digital Archive.....	13
<b>References.....</b>	<b>14</b>

# GENERAL INFORMATION

## (1) Are there requirements regarding the data management from your scholarly/scientific community?

❖ yes

(a) If yes, what are the requirements?

❖ DFG Guidelines on the Handling of Research Data

## (2) What kind of dataset is it?

---

*The following questions collect information on the data that is produced or used in the project. They also help to estimate the value of the data in terms of potential re-use and long-term preservation. In the case of personal data, the principle of data minimization (Art. 5 EU General Data Protection Regulation) allows the collection of personal data only when there are no other reasonable means to clarify the research question (re-use of existing data would be such a reasonable means). Also, there shall be no more information collected than necessary. The information regarding the data collected, produced, or used in the project is gathered along the datasets. The definition of these datasets is an important conceptual decision that has to be done individually and carefully for each project.*

---



<b>Name of Experiment / Acronym / Number:</b>	Anaphoric Distance: Story Continuation Experiment
<b>PI or responsible person (head of the study):</b>	Manfred Stede
<b>Other persons involved:</b>	Berfin Aktas
<b>Subject area:</b>	discourse, coreference
<b>Method / Type of data:</b>	Narrative continuations in text and audio format
<b>Anonymizable data:</b>	yes
<b>Participants:</b>	native speakers of English (recruited via crowd sourcing)
<b>Short description (of the study):</b>	The impact of anaphoric distance on referential choice is tested across oral and written media. Native speakers of English are recruited via a crowdsourcing company. Participants continued to short stories either by providing oral or written responses.
<b>Other comments:</b>	
<b>Time for data collection (approximate):</b>	October-November 2020
<b>Time for data analysis (approximate):</b>	May-July 2021
<b>Related publications:</b>	<p>Das, D., Scheffler, T., Bourgonje, P. &amp; Stede, M. 2018 Constructing a Lexicon of English Discourse Connectives. K. Komtani, D. Litman, K. Yu, A. Papangelis, L. Cavedon, &amp; M. Nakano (eds.), Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue, pp. 360-365. *</p> <p>Stede, M., Scheffler, T., &amp; Mendes, A. 2019 Connective-Lex: A Web-Based Multilingual Lexical Resource for Connectives. Discourse [Online], 24. doi:10.4000/discours.10098</p>



**Keywords (used in publications):**

**Link to preregistration:**

**Funding reference:**

Funded by the Deutsche  
Forschungsgemeinschaft (DFG,  
German Research Foundation)  
– Project number 317633480  
– SFB 1287

### **(3) Which individuals, groups or institutions could be interested in re-using this dataset? What consequences does the reuse potential have for the provision of the data later?**

---

*It is important to specify whether the data will be permitted for reuse. But legal impediments, such as privacy, and copyright must be taken into account.*

---

The data collected via the crowdsourcing experiment that might be of interest to researchers at or outside of University of Potsdam, for example to:

- ❖ Perform further experiments on the impact of anaphoric distance on referential choice across different production media
- ❖ Analyze the differences in oral vs written narratives

## TECHNICAL INFORMATION

### (4) Where is the dataset stored during the project?

---

*Please delete all project-files from source that is not part of the University of Potsdam.*

---

<b>for raw data:</b>	SFB1287 - File Server Box.UP - Cloud (University of Potsdam) Computer in the laboratory
<b>for analysis data:</b>	Researcher´s Computer Box.UP - Cloud (University of Potsdam)
<b>for further documentation, related code, or software:</b>	SFB1287 - File Server Researcher´s Computer Box.UP - Cloud (University of Potsdam)

(a) If data is stored on lab or personal computers, please describe the backup strategy.

- ❖ Sharing data through box UP
- ❖ Regular backups are taken to an external hard drive.
- ❖ The university cloud and SFB-File Server is backed up regularly.

## (5) Which file formats are used?

---

*When choosing a data format, one should consider the consequences for collaborative use, long-term preservation as well as reuse. It is advisable to use formats that are standardised, open, non-proprietary, and well-established in the respective scholarly community. A table with recommended file formats can be found in found in Kristin Briney, *Data Management for Researchers*, Pelargic, 2015, pages 133-134.*

---

- ❖ Participants' text responses: .xlsx
- ❖ Participants' audio responses: .webm, .ogg, .wav
- ❖ Audio Stimuli: .mp3

# PUBLICATION

## (6) Will this dataset be published or shared?

---

*anonymizable data*

- ID will be removed (and or code-list will be destroyed) [legally correct: code list will be destroyed as soon as possible without jeopardizing experiment; exception: follow-up study planned, if so, talk to UP data protection officer (Dr. David Kneis; <mailto:datenschutz@uni-potsdam.de>) on how to do this correctly]
- Publication of anonymized data and code on OSF or RADAR (University of Potsdam) (or as required by the Journal)

---

*non-anonymizable data*

- on RADAR (University of Potsdam) but not accessible from the outside world

❖ yes

### (a) If yes, under which terms of use or license will the dataset be published or shared?

---

*The options refer to the licenses of the Creative Commons family. If data is anonymised / pseudonymized, it's probably not legally required, but might be good in terms of research ethics to adjust consent forms / subject information sheets.*

---



**Principal investigator of the study assures that the consent form / subject information sheets support publishing of the data.**

**for data:**

CC 0 (Public Domain) (recommended)



for scripts: -

for software: -

(b) If yes, when will the data be published?

---

*Recommended procedure: Upload data and obtain digital identifier (e.g., DOI, OSF link) when submitting the first paper; thus, you can cite the data in the paper. If necessary, restrict public access (embargo) until last paper published (max. 2 years).*

---

❖ when the paper is published

(c) If no, please explain why not. Please differentiate between legal and contractual reasons and voluntary restrictions.

## **(7) Which measures of quality assurance are taken for this dataset?**

❖ The participants qualified via a qualification experiment where their responses are evaluated manually. Engagement of the participants are checked via asking the result of mathematical expressions and/or asking basic questions about a text. All the audio responses are listened and low-effort/incoherent responses are eliminated.

## LEGAL AND ETHICS

### (8) Does this dataset contain personal data?

---

*The EU General Data Protection Regulation (GDPR) defines in Art. 4 personal data as "any information relating to an identified or identifiable natural person". An identifiable natural person is "one who can be identified, directly or indirectly, in particular by reference to an identifier such as a name, an identification number, location data, an online identifier or to one or more factors specific to the physical, physiological, genetic, mental, economic, cultural or social identity of that natural person".*

*For advice and/or trainings on how to comply with privacy regulations - including proper anonymization & pseudonymization - you can always contact UP's privacy officer, Dr David Kneis, at [datenschutz@uni-potsdam.de](mailto:datenschutz@uni-potsdam.de). As the privacy officer, his perspective is focussed more on the legal side of things than the research ethics or technical aspects.*

---

❖ yes

### (9) Are the data anonymised?

---

*Anonymised data: ID will be removed (and or code-list will be destroyed) [legally correct: code list will be destroyed as soon as possible without jeopardizing experiment; exception: follow-up study planned, if so, definitely talk to UP data protection officer on how to do this correctly]*

---

**for raw data:** no

**for analysis data:** no

**for published data:** yes

## (10) Does the project use and/or produce data that is protected by intellectual or industrial property rights?

---

Measurement data has no intellectual property, so usually, the answer here will be „no“. Data or software can be subject to intellectual or industrial property rights. Applicable laws differ broadly even within EU. According to the German copyright law (UrhG) works of literature, scholarship and the arts that can be regarded as a “personal intellectual creation” are protected by copyright. Mere data, e.g., measured data or survey data, and metadata (except in some cases descriptive metadata) are not protected by copyright. § 2 of the UrhG lists the following kinds of protected works (list is not concluded):

- 
- linguistic works such as written works, speeches, and computer programs
  - works or the fine arts including works of the applied arts as well as sketches of such works
  - works of photography
  - descriptions and illustrations of scholarly or technical nature such as drawings, plans, maps, sketches, tables, and three-dimensional representations
- 

According to § 3, copyright is also applicable to translations and other modifications or adaptations of work if they are individual intellectual creations of the editor. Finally, according to § 4 copyright also extends to collected editions and database works. Collected editions are: “collections of work, data or other independent elements that are individual intellectual creations based on the selection and arrangement of the elements”.

Database works are defined as “collected editions, the elements of which are arranged in a systematic or methodical way and can be accessed individually by electronic means or in other ways”.

---

❖ yes

(a) If yes, please explain which!

- ❖ licensed corpora (e.g., we bought the LDC license for the CallHome corpus)

# STORAGE AND LONG-TERM PRESERVATION

## **(11) Does this dataset have to be preserved for long-term?**

---

*The DFG expects primary data that is the basis of a publication to be stored in the researcher's own institution or an appropriate nationwide infrastructure long-term (for at least 10 years).*

---

- ❖ At least 10 years after the end of the first funding period of the SFB1287.

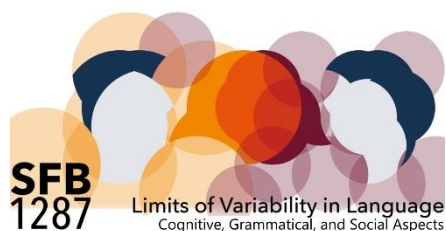
## **(12) What are the reasons this dataset has to be preserved for the long-term?**

- ❖ Used in a publication / proof of good scientific practice
- ❖ Re-use (if anonymizable data) in subsequent projects or by others
- ❖ By DFG requirements

**(13) Where will the data (including metadata, documentation, and relevant code) be stored or archived after the end of the project?**

❖ SFB1287 - File Server

## REFERENCES



**SFB 1287**

<https://www.sfb1287.uni-potsdam.de>



**University of Potsdam**

<https://www.uni-potsdam.de>

Funded by

**DFG** Deutsche  
Forschungsgemeinschaft  
German Research Foundation

**Deutsche Forschungsgemeinschaft e.V.**

<https://www.dfg.de>