# FWO DMP Template - Flemish Standard Data Management Plan

Project supervisors (from application round 2018 onwards) and fellows (from application round 2020 onwards) will, upon being awarded their project or fellowship, be invited to develop their answers to the data management related questions into a DMP. The FWO expects a **completed DMP no later than 6 months after the official start date** of the project or fellowship. The DMP should not be submitted to FWO but to the research co-ordination office of the host institute; FWO may request the DMP in a random check.

At the end of the project, the **final version of the DMP** has to be added to the final report of the project; this should be submitted to FWO by the supervisor-spokesperson through FWO's e-portal. This DMP may of course have been updated since its first version. The DMP is an element in the final evaluation of the project by the relevant expert panel. Both the DMP submitted within the first 6 months after the start date and the final DMP may use this template.

The DMP template used by the Research Foundation Flanders (FWO) corresponds with the Flemish Standard Data Management Plan. This Flemish Standard DMP was developed by the Flemish Research Data Network (FRDN) Task Force DMP which comprises representatives of all Flemish funders and research institutions. This is a standardized DMP template based on the previous FWO template that contains the core requirements for data management planning. To increase understanding and facilitate completion of the DMP, a standardized **glossary** of definitions and abbreviations is available via the following link.

| 1. General Project Information | |
|---|---|
| Name Grant Holder & ORCID | **Freek Van de Velde (0000-0003-3050-2207), supervisor** |
| Contributor name(s) (+ ORCID) & roles | **Dirk Speelman (0000-0003-1561-1851), co-supervisor**<br>**Anthe Sevenants (0000-0002-5055-770X), researcher** |
| Project number[1] & title | 3H220497, Connecting morphosyntax and lexical semantics with Elastic Net regression |
| Funder(s) GrantID[2] | ZKE1867 - G059922N |
| Affiliation(s) | ☑ KU Leuven<br>☐ Universiteit Antwerpen<br>☐ Universiteit Gent<br>☐ Universiteit Hasselt<br>☐ Vrije Universiteit Brussel<br>☐ Other:<br>Provide ROR[3] identifier when possible: |

---

[1] "Project number" refers to the institutional project number. This question is optional since not every institution has an internal project number different from the GrantID. Applicants can only provide one project number.

[2] Funder(s) GrantID refers to the number of the DMP at the funder(s), here one can specify multiple GrantIDs if multiple funding sources were used.

[3] Research Organization Registry Community. https://ror.org/

| | |
|---|---|
| Please provide a short project description | This project proposes to use regularization methods from machine learning, more specifically Elastic Net regression (and its siblings Ridge and Lasso), to look into lexical semantic effects in morphosyntactic alternances. These regularization techniques apply shrinkage to the coefficients and can thus be used for variable selection, especially when the number of predictors is very large. In variationist studies, this is often the case if one wishes to enter lexemes associated with a construction into a regression model to predict constructional variants. We combine the Elastic Net regulator with k-fold cross-validation - a standard procedure - to avoid overfitting. Our approach mitigates the various drawbacks present in alternative approaches that are currently used in variationist linguistics, like random factors in mixed models and collostructional analysis. We look at ten multifactorially driven alternances from Dutch. The project offers a transparent pipeline that can easily be extrapolated to other case studies, and to other languages. |

| | | | | **2. Research Data Summary** | | | |
|---|---|---|---|---|---|---|---|

List and describe all datasets or research materials that you plan to generate/collect or reuse during your research project. For each dataset or data type (observational, experimental etc.), provide a short name & description (sufficient for yourself to know what data it is about), indicate whether the data are newly generated/collected or reused, digital or physical, also indicate the type of the data (the kind of content), its technical format (file extension), and an estimate of the upper limit of the volume of the data[4].

| | | | | *ONLY FOR DIGITAL DATA* | *ONLY FOR DIGITAL DATA* | *ONLY FOR DIGITAL DATA* | *ONLY FOR PHYSICAL DATA* |
|---|---|---|---|---|---|---|---|
| Dataset Name | Description | New or Reused | Digital or Physical | Digital Data Type | Digital Data Format | Digital Data Volume (MB, GB, TB) | Physical Volume |
| Corpora and datasets | We will both use existing text corpora, viz. SoNaR (Oostdijk N., M. Reynaert, V. Hoste & I. Schuurman (2013). The Construction of a 500-Million-Word Reference Corpus of Contemporary Written Dutch. In: P. Spyns & J.Odijk (red.), *Essential Speech and Language* | ☒ Generate new data<br>☒ Reuse existing data | ☒ Digital<br>☐ Physical | ☒ Observational<br>☐ Experimental<br>☐ Compiled/ aggregated data<br>☐ Simulation data<br>☒ Software<br>☐ Other<br>☐ NA | ☐ .por<br>☐ .xml<br>☐ .tab<br>☒ .csv<br>☐ .pdf<br>☒ .txt<br>☐ .rtf<br>☐ .dwg<br>☐ .tab<br>☐ .gml<br>☐ other:<br>☐ NA | ☐ < 100 MB<br>☐ < 1 GB<br>☒ < 100 GB<br>☐ < 1 TB<br>☐ < 5 TB<br>☐ < 10 TB<br>☐ < 50 TB<br>☐ > 50 TB<br>☐ NA | |

---

[4] Add rows for each dataset you want to describe.

| | | *Technology for Dutch. Theory and Applications of Natural Language Processing.* Heidelberg: Springer, 219-247.<br><br>and Dutch C-CLAMP (Piersoul, Jozefien, Robbert De Troij & Freek Van de Velde. 2021. '150 years of written Dutch: the construction of the Dutch Corpus of Contemporary and Late Modern Periodicals'. *Nederlandse* | | | | | |
|---|---|---|---|---|---|---|---|

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | *Taalkunde* 26(3): 339-362.<br><br>We will reuse datasets from different Dutch alternation studies, carried out as part of PhD projects or as part of MA theses or as part of individual research within the research unit (QLVL). Exactly which datasets are suitable for the research project at hand will be established. | | | | | | |
| Feature matrices | The matrices of features that are deduced from the corpus | New data | Digital | Compiled/ aggregated data | .npy (NumPy binary format) | < 100 MB | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | data. | | | | | | |
| Jupyter Notebooks | The Python files used to create the feature matrices. | New data | Digital | Software | .ipynb (IPython Notebook) | < 100 MB | |
| R scripts | The R scripts used to analyse the data statistically. | New data | Digital | Software | .R (R script source file) | < 100 MB | |

| | |
|---|---|
| If you reuse existing data, please specify the source, preferably by using a persistent identifier (e.g. DOI, Handle, URL etc.) per dataset or data type. | Can only be decided after inspection of the datasets (see above) |

---

[5] These data are generated by combining multiple existing datasets.

| | |
|---|---|
| Are there any ethical issues concerning the creation and/or use of the data (e.g. experiments on humans or animals, dual use)? If so, please describe these issues further and refer to specific datasets or data types when appropriate. | ☐ Yes, human subject data<br>☐ Yes, animal data<br>☐ Yes, dual use<br>☒ No<br>If yes, please describe: |
| Will you process personal data[6]? If so, briefly describe the kind of personal data you will use. Please refer to specific datasets or data types when appropriate. If available, add the reference to your file in your host institution's privacy register. | ☐ Yes<br>☒ No<br>If yes:<br><br>- Short description of the kind of personal data that will be used:<br>- Privacy Registry Reference: |
| Does your work have potential for commercial valorization (e.g. tech transfer, for example spin-offs, commercial exploitation, …)? If so, please comment per dataset or data type where appropriate. | ☐ Yes<br>☒ No<br>If yes, please comment: |
| Do existing 3rd party agreements restrict exploitation or dissemination of the data you (re)use (e.g. Material/Data transfer agreements, research collaboration agreements)? If so, please explain to what data they relate and what restrictions are in place. | ☐ Yes<br>☒ No<br>If yes, please explain: |

---

[6] See Glossary Flemish Standard Data Management Plan

| Are there any other legal issues, such as intellectual property rights and ownership, to be managed related to the data you (re)use?<br>If so, please explain to what data they relate and which restrictions will be asserted. | ☐ Yes<br>☒ No<br>If yes, please explain: |
|---|---|

| 3. Documentation and Metadata | |
|---|---|
| Clearly describe what approach will be followed to capture the accompanying information necessary to keep **data understandable and usable**, for yourself and others, now and in the future (e.g. in terms of documentation levels and types required, procedures used, Electronic Lab Notebooks, README.txt files, Codebook.tsv etc. where this information is recorded). | • **All source code scripts will be commented as much as possible to explain the inner workings of the scripts written for the research workflow.**<br>• **If the scripts are hosted on a collaboration platform such as GitHub, the repository's README file will explain the nature of the scripts and how to use them.** |
| Will a metadata standard be used to make it easier to **find and reuse the data**?<br><br>If so, please specify which metadata standard will be used. If not, please specify which metadata will be created to make the data easier to find and reuse.<br><br>*REPOSITORIES COULD ASK TO DELIVER METADATA IN A CERTAIN FORMAT, WITH SPECIFIED ONTOLOGIES AND VOCABULARIES, I.E. STANDARD LISTS WITH UNIQUE IDENTIFIERS.* | ☐ Yes<br>☒ No<br>If yes, please specify (where appropriate per dataset or data type) which metadata standard will be used:<br><br>If no, please specify (where appropriate per dataset or data type) which metadata will be created:<br>The GitHub README's contents should make the data discoverable with search engines. |

| 4. Data Storage & Back-up during the Research Project | |
|---|---|
| Where will the data be stored? | • GitHub (researcher's personal account and/or P.I. personal account)<br>• The researcher's laptop<br>• The researcher's personal home server<br>• The QLVL server<br>• OneDrive |
| How will the data be backed up?<br><br>*WHAT STORAGE AND BACKUP PROCEDURES WILL BE IN PLACE TO PREVENT DATA LOSS? DESCRIBE THE LOCATIONS, STORAGE MEDIA AND PROCEDURES THAT WILL BE USED FOR STORING AND BACKING UP DIGITAL AND NON-DIGITAL DATA DURING RESEARCH.[7]*<br><br>*REFER TO INSTITUTION-SPECIFIC POLICIES REGARDING BACKUP PROCEDURES WHEN APPROPRIATE.* | The GitHub repositories provide *version control*, which means that anyone can look at the incremental history of how the data were created. In addition, the GitHub platform will ensure that the data will be available far into the future (they act as a de facto cloud host).<br>We also back up all data by the home institution (KU Leuven) subscription to OneDrive |
| Is there currently sufficient storage & backup capacity during the project? If yes, specify concisely. If no or insufficient storage or backup capacities are available, then explain how this will be taken care of. | ☒ Yes<br>☐ No<br>If yes, please specify concisely: KU Leuven ensures a personal data limit exceeding the current projected size of all data<br><br>If no, please specify: |

---

[7] Source: Ghent University Generic DMP Evaluation Rubric: https://osf.io/2z5g3/

| How will you ensure that the data are securely stored and not accessed or modified by unauthorized persons?<br><br>*CLEARLY DESCRIBE THE MEASURES (IN TERMS OF PHYSICAL SECURITY, NETWORK SECURITY, AND SECURITY OF COMPUTER SYSTEMS AND FILES) THAT WILL BE TAKEN TO ENSURE THAT STORED AND TRANSFERRED DATA ARE SAFE. [7]* | • The researcher's GitHub account uses Multi-Factor Authentication (MFA), which means attackers cannot easily break into the account using only a password. In addition, the password for the GitHub account is behind a password manager (which also uses MFA).<br>• The researcher's laptop is secured using a strong password.<br>• The researcher's personal server is not directly connected to the internet.<br>• The institution's OneDrive is also protected |
|---|---|
| What are the expected costs for data storage and backup during the research project? How will these costs be covered? | There are no expected costs for data storage. GitHub data storage is free. Should more data storage be needed beyond the GitHub provided space, the (free) KU Leuven offerings or CLARIN will be used. |

| **5. Data Preservation after the end of the Research Project** ||
|---|---|
| Which data will be retained for at least five years (or longer, in agreement with other retention policies that are applicable) after the end of the project? In case some data cannot be preserved, clearly state the reasons for this (e.g. legal or contractual restrictions, storage/budget issues, institutional policies...). | All data should be able to be preserved for five years or longer. |
| Where will these data be archived (stored and curated for the long-term)? | All data should remain available on GitHub. In addition, the data will be stored on the university's central servers (with automatic back-up procedures) for at least 10 years, conform the KU Leuven RDM policy. |

| What are the expected costs for data preservation during the expected retention period? How will these costs be covered? | The expected costs will be determined once more data has been generated for the project. We do not expect the costs to exceed 1000 euros, based on experience with other research projects. The bench fee should be able to cover these costs. |
| --- | --- |

| 6. Data Sharing and Reuse | |
|---|---|
| Will the data (or part of the data) be made available for reuse after/during the project? Please explain per dataset or data type which data will be made available.<br><br>*Note that 'available' does not necessarily mean that the data set becomes openly available, conditions for access and use may apply. Availability in this question thus entails both open & restricted access. For more information: https://wiki.surfnet.nl/display/standards/info-eu-repo/#infoeurepo-AccessRights* | ☒ Yes, in an Open Access repository<br>☐ Yes, in a restricted access repository (after approval, institutional access only, …)<br>☐ No (closed access)<br>☐ Other, please specify: |
| If access is restricted, please specify who will be able to access the data and under what conditions. | |
| Are there any factors that restrict or prevent the sharing of (some of) the data (e.g. as defined in an agreement with a 3rd party, legal restrictions)? Please explain per dataset or data type where appropriate. | ☐ Yes, privacy aspects<br>☐ Yes, intellectual property rights<br>☐ Yes, ethical aspects<br>☐ Yes, aspects of dual use<br>☐ Yes, other<br>☒ No<br><br>If yes, please specify: |
| Where will the data be made available?<br>If already known, please provide a repository per dataset or data type. | To be decided later |

| | |
|---|---|
| When will the data be made available?<br><br>*THIS COULD BE A SPECIFIC DATE (DD/MM/YYYY) OR AN INDICATION SUCH AS 'UPON PUBLICATION OF RESEARCH RESULTS'.* | To be decided later |
| Which data usage licenses are you going to provide? If none, please explain why.<br><br>*A DATA USAGE LICENSE INDICATES WHETHER THE DATA CAN BE REUSED OR NOT AND UNDER WHAT CONDITIONS. IF NO LICENCE IS GRANTED, THE DATA ARE IN A GREY ZONE AND CANNOT BE LEGALLY REUSED. DO NOTE THAT YOU MAY ONLY RELEASE DATA UNDER A LICENCE CHOSEN BY YOURSELF IF IT DOES NOT ALREADY FALL UNDER ANOTHER LICENCE THAT MIGHT PROHIBIT THAT.*<br><br>*EXAMPLE ANSWER: E.G. "DATA FROM THE PROJECT THAT CAN BE SHARED WILL BE MADE AVAILABLE UNDER A CREATIVE COMMONS ATTRIBUTION LICENSE (CC-BY 4.0), SO THAT USERS HAVE TO GIVE CREDIT TO THE ORIGINAL DATA CREATORS." [8]* | To be decided later |
| Do you intend to add a PID/DOI/accession number to your dataset(s)? If already available, please provide it here.<br><br>*INDICATE WHETHER YOU INTEND TO ADD A PERSISTENT AND UNIQUE IDENTIFIER IN ORDER TO IDENTIFY AND RETRIEVE THE DATA.* | ☐ Yes<br>☐ No<br>If yes:<br><br>To be decided later |
| What are the expected costs for data sharing? How will these costs be covered? | To be decided later |

---

[8] Source: Ghent University Generic DMP Evaluation Rubric: https://osf.io/2z5g3/

| 7. Responsibilities | |
|---|---|
| Who will manage data documentation and metadata during the research project? | Anthe Sevenants |
| Who will manage data storage and backup during the research project? | Anthe Sevenants |
| Who will manage data preservation and sharing? | Anthe Sevenants |
| Who will update and implement this DMP? | Anthe Sevenants |