

## 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	<b>Pieter Vanmechelen</b>	<b>0000-0002-3733-2880</b>
Contributor name(s) (+ ORCID) & roles	<b>Giovanni Samaey</b> <b>Geert Lombaert</b>	<b>Promotor</b> <b>Co-promotor</b> <b>0000-0001-8433-4523</b> <b>0000-0002-9273-3038</b>
Project number <sup>1</sup> & title	3E200531 Multilevel Markov Chain Monte Carlo methods for Bayesian full-field data assimilation	
Funder(s) GrantID <sup>2</sup>	<b>1SD1823N</b>	
Affiliation(s)	<input checked="" type="checkbox"/> KU Leuven <input type="checkbox"/> Universiteit Antwerpen <input type="checkbox"/> Universiteit Gent <input type="checkbox"/> Universiteit Hasselt <input type="checkbox"/> Vrije Universiteit Brussel <input type="checkbox"/> Other: Provide ROR <sup>3</sup> identifier when possible:	
Please provide a short project description	Markov chain Monte Carlo (MCMC) methods are an indispensable tool in sampling probability distributions for which a direct sampling is impossible. In Bayesian data assimilation, these methods are used to sample the posterior probability distribution for a parameter that needs to be estimated based on an accept-reject procedure, given a prior distribution for the parameter and a model that allows computing the likelihood of the data for a given value of the parameter. This project deals with the development and analysis of MCMC methods when the available data has a full-field nature. Markov chain Monte Carlo methods face two problems when confronted with full field data. First, computing the likelihood can be a very expensive step in MCMC sampling, especially when the corresponding model contains multiple time scales and is high-dimensional. Second, the acceptance rates of standard MCMC procedures drop as the data dimensionality increases. In this project, we will develop and analyse multilevel Markov chain Monte Carlo methods that increase the computational efficiency of Bayesian posterior sampling by making use of a hierarchy of levels at different resolution, and by explicitly exploiting spatial correlation that is present in full-field data. While being application-independent in nature, the developed methods will be applied to problems in damage assessment in civil engineering.	

<sup>1</sup> "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.

<sup>2</sup> 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.

<sup>3</sup> Research Organization Registry Community. <https://ror.org/>

## 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<sup>4</sup>.

Dataset Name	Description	New or Reused	Digital or Physical	ONLY FOR DIGITAL DATA	ONLY FOR DIGITAL DATA	ONLY FOR DIGITAL DATA	ONLY FOR PHYSICAL DATA
				Digital Data Type	Digital Data Format	Digital Data Volume (MB, GB, TB)	Physical Volume
		<input type="checkbox"/> Generate new data <input type="checkbox"/> Reuse existing data	<input type="checkbox"/> Digital <input type="checkbox"/> Physical	<input type="checkbox"/> Observational <input type="checkbox"/> Experimental <input type="checkbox"/> Compiled/aggregated data <input type="checkbox"/> Simulation data <input type="checkbox"/> Software <input type="checkbox"/> Other <input type="checkbox"/> NA	<input type="checkbox"/> .por <input type="checkbox"/> .xml <input type="checkbox"/> .tab <input type="checkbox"/> .csv <input type="checkbox"/> .pdf <input type="checkbox"/> .txt <input type="checkbox"/> .rtf <input type="checkbox"/> .dwg <input type="checkbox"/> .tab <input type="checkbox"/> .gml <input type="checkbox"/> other: <input type="checkbox"/> NA	<input type="checkbox"/> < 100 MB <input type="checkbox"/> < 1 GB <input type="checkbox"/> < 100 GB <input type="checkbox"/> < 1 TB <input type="checkbox"/> < 5 TB <input type="checkbox"/> < 10 TB <input type="checkbox"/> < 50 TB <input type="checkbox"/> > 50 TB <input type="checkbox"/> NA	
Simulation code	MATLAB scripts	New data	Digital	Software	.m	< 1 GB	/
Numerical experiments	Results produced by running MATLAB scripts	New data	Digital	Experimental data	.mat	< 100 GB	/

<sup>4</sup> Add rows for each dataset you want to describe.

Manuscripts	Paper and thesis text files	New data	Digital	Text	.tex and .pdf	< 1 GB	/
Papers	Literature study	Existing data	Digital	Text	.pdf	< 100 GB	/

  

**GUIDANCE:**

DATA CAN BE DIGITAL OR PHYSICAL (FOR EXAMPLE BIOBANK, BIOLOGICAL SAMPLES, ...). DATA TYPE: DATA ARE OFTEN GROUPED BY TYPE (OBSERVATIONAL, EXPERIMENTAL ETC.), FORMAT AND/OR COLLECTION/GENERATION METHOD.

EXAMPLES OF DATA TYPES: OBSERVATIONAL (E.G. SURVEY RESULTS, SENSOR READINGS, SENSORY OBSERVATIONS); EXPERIMENTAL (E.G. MICROSCOPY, SPECTROSCOPY, CHROMATOGRAMS, GENE SEQUENCES); COMPILED/AGGREGATED DATA<sup>5</sup> (E.G. TEXT & DATA MINING, DERIVED VARIABLES, 3D MODELLING); SIMULATION DATA (E.G. CLIMATE MODELS); SOFTWARE, ETC.

EXAMPLES OF DATA FORMATS: TABULAR DATA (.POR, .SPSS, STRUCTURED TEXT OR MARK-UP FILE XML, .TAB, .CSV), TEXTUAL DATA (.RTF, .XML, .TXT), GEOSPATIAL DATA (.DWG, .GML, ..), IMAGE DATA, AUDIO DATA, VIDEO DATA, DOCUMENTATION & COMPUTATIONAL SCRIPT.

DIGITAL DATA VOLUME: PLEASE ESTIMATE THE UPPER LIMIT OF THE VOLUME OF THE DATA PER DATASET OR DATA TYPE.

PHYSICAL VOLUME: PLEASE ESTIMATE THE PHYSICAL VOLUME OF THE RESEARCH MATERIALS (FOR EXAMPLE THE NUMBER OF RELEVANT BIOLOGICAL SAMPLES THAT NEED TO BE STORED AND PRESERVED DURING THE PROJECT AND/OR AFTER).

  

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.	The only existing data reused is the collection of papers used in the literature study. Sources are various journals in which these papers were published.
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.	<input type="checkbox"/> Yes, human subject data <input type="checkbox"/> Yes, animal data <input type="checkbox"/> Yes, dual use <input checked="" type="checkbox"/> No If yes, please describe:

<sup>5</sup> These data are generated by combining multiple existing datasets.

<p>Will you process personal data<sup>6</sup>? 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.</p>	<p><input type="checkbox"/> Yes  <input checked="" type="checkbox"/> No</p> <p>If yes:</p> <ul style="list-style-type: none"> <li>- Short description of the kind of personal data that will be used:</li> <li>- Privacy Registry Reference:</li> </ul>
<p>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.</p>	<p><input checked="" type="checkbox"/> Yes  <input type="checkbox"/> No</p> <p>If yes, please comment: The source code used to generate results will be made publicly available, hence can be used by others as example for commercial implementation of our methods. However, the code is not directly usable commercially as it is designed with academic purposes in mind.</p>
<p>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.</p>	<p><input type="checkbox"/> Yes  <input checked="" type="checkbox"/> No</p> <p>If yes, please explain:</p>
<p>Are there any other legal issues, such as intellectual property rights and ownership, to be managed related to the data you (re)use?  If so, please explain to what data they relate and which restrictions will be asserted.</p>	<p><input type="checkbox"/> Yes  <input checked="" type="checkbox"/> No</p> <p>If yes, please explain:</p>

---

<sup>6</sup> See Glossary Flemish Standard Data Management Plan

### 3. Documentation and Metadata

<p>Clearly describe what approach will be followed to capture the accompanying information necessary to keep <b>data understandable and usable</b>, 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).</p>	<p><b>Code is documented and has an accompanying README. The code is in version control, where the git history allows to trace back changes. Computational experiments are explained carefully in papers and are tested for reproducibility. Additionally, published code is accompanied with a dockerfile such that a virtual machine can be set up with the right environment to run the published codes. Using these codes and dockerfile based virtual machines, each of the published computational experiments can be reproduced.</b></p>
<p>Will a metadata standard be used to make it easier to <b>find and reuse the data</b>?</p> <p>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.</p> <p><i>REPOSITORIES COULD ASK TO DELIVER METADATA IN A CERTAIN FORMAT, WITH SPECIFIED ONTOLOGIES AND VOCABULARIES, I.E. STANDARD LISTS WITH UNIQUE IDENTIFIERS.</i></p>	<p><input checked="" type="checkbox"/> Yes  <input type="checkbox"/> No</p> <p>If yes, please specify (where appropriate per dataset or data type) which metadata standard will be used:  For computational experiments we encode metadata about the running environment for the code in dockerfiles.</p> <p>If no, please specify (where appropriate per dataset or data type) which metadata will be created:  Metadata about how to reproduce results will be documented in papers and README files.</p> <p><b>Note:</b> in the future, we plan to automate metadata generation in our research group using iRODS. The pilot for this change in data management flow still has to be set up.</p>

### 4. Data Storage & Back-up during the Research Project

<p>Where will the data be stored?</p>	<p>Data is stored locally on a departmentally issued laptop and depending on data type, back-ups are stored:</p> <ul style="list-style-type: none"> <li>• For version-controlled files, the KU Leuven Gitlab service is used</li> <li>• Literature is maintained in a reference manager (Zotero)</li> <li>• Other files are backed up on the KU Leuven OneDrive or the NextCloud server of the Department of Computer Science.</li> </ul> <p>After the research:</p> <ul style="list-style-type: none"> <li>• Where possible, all relevant data/code will be made publicly available with the corresponding papers</li> <li>• Manuscripts and all corresponding data are stored with read-only access at the Department of Computer Science.</li> </ul>
<p>How will the data be backed up?</p> <p><i>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.<sup>7</sup></i></p> <p><i>REFER TO INSTITUTION-SPECIFIC POLICIES REGARDING BACKUP PROCEDURES WHEN APPROPRIATE.</i></p>	<p>All files are synced with departmental file servers to ensure daily backup. The sole exception is the papers of the literature study which are stored in Zotero. This all is in accordance with the wider NUMA DMP policy, used by almost all researchers in our section.</p>

<sup>7</sup> Source: Ghent University Generic DMP Evaluation Rubric: <https://osf.io/2z5g3/>

<p>Is there currently sufficient storage &amp; 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.</p>	<p><input checked="" type="checkbox"/> Yes  <input type="checkbox"/> No</p> <p>If yes, please specify concisely: All data is digital and currently backed up with plenty of storage space left. If in future this space is insufficient, more can be requested at the departmental level.</p> <p>If no, please specify:</p>
<p>How will you ensure that the data are securely stored and not accessed or modified by unauthorized persons?</p> <p><i>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. <sup>7</sup></i></p>	<p>This issue is covered by the NUMA data management plan.</p>
<p>What are the expected costs for data storage and backup during the research project? How will these costs be covered?</p>	<p>Costs are fully covered by the NUMA research section.</p>

## 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...).	<p><b>Where possible, all relevant data/code will be made publicly available with corresponding papers.</b></p> <p><b>Manuscripts and all corresponding data are stored with read-only access at the Department of Computer Science.</b></p>
Where will these data be archived (stored and curated for the long-term)?	<b>Department of Computer Science</b>
What are the expected costs for data preservation during the expected retention period? How will these costs be covered?	<b>Costs will be fully covered by the NUMA research section.</b>

## 6. Data Sharing and Reuse

<p>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.</p> <p><i>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 &amp; RESTRICTED ACCESS. FOR MORE INFORMATION: <a href="https://wiki.surfnet.nl/display/STANDARDS/INFO-EU-REPO/#INFOEU-REPO-ACCESSRIGHTS">HTTPS://WIKI.SURFNET.NL/DISPLAY/STANDARDS/INFO-EU-REPO/#INFOEU-REPO-ACCESSRIGHTS</a></i></p>	<p><input checked="" type="checkbox"/> Yes, in an Open Access repository</p> <p><input checked="" type="checkbox"/> Yes, in a restricted access repository (after approval, institutional access only, ...)</p> <p><input type="checkbox"/> No (closed access)</p> <p><input type="checkbox"/> Other, please specify:</p>
<p>If access is restricted, please specify who will be able to access the data and under what conditions.</p>	<p>The promotor of this project, Giovanni Samaey, has full access to all repositories and can regulate access to the data.</p>
<p>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.</p>	<p><input type="checkbox"/> Yes, privacy aspects</p> <p><input type="checkbox"/> Yes, intellectual property rights</p> <p><input type="checkbox"/> Yes, ethical aspects</p> <p><input type="checkbox"/> Yes, aspects of dual use</p> <p><input type="checkbox"/> Yes, other</p> <p><input checked="" type="checkbox"/> No</p> <p>If yes, please specify:</p>
<p>Where will the data be made available? If already known, please provide a repository per dataset or data type.</p>	<p><b>KU Leuven gitlab servers</b></p>

<p><b>When will the data be made available?</b></p> <p><i>THIS COULD BE A SPECIFIC DATE (DD/MM/YYYY) OR AN INDICATION SUCH AS 'UPON PUBLICATION OF RESEARCH RESULTS'.</i></p>	<p><b>Open access repositories will be created upon publication of the corresponding paper. Restricted access repositories can be made available upon request.</b></p>
<p><b>Which data usage licenses are you going to provide? If none, please explain why.</b></p> <p><i>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.</i></p> <p><i>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." <sup>8</sup></i></p>	<p><b>This is still under consideration.</b></p>
<p><b>Do you intend to add a PID/DOI/accession number to your dataset(s)? If already available, please provide it here.</b></p> <p><i>INDICATE WHETHER YOU INTEND TO ADD A PERSISTENT AND UNIQUE IDENTIFIER IN ORDER TO IDENTIFY AND RETRIEVE THE DATA.</i></p>	<p><input checked="" type="checkbox"/> Yes  <input type="checkbox"/> No</p> <p>If yes: We will use a KU Leuven system that is currently in development.</p>

<sup>8</sup> Source: Ghent University Generic DMP Evaluation Rubric: <https://osf.io/2z5g3/>

What are the expected costs for data sharing? How will these costs be covered?	<b>Costs are not precisely known, but expected to be modest.</b>
---	--

7. Responsibilities	
Who will manage data documentation and metadata during the research project?	<b>Pieter Vanmechelen</b>
Who will manage data storage and backup during the research project?	<b>Pieter Vanmechelen</b>
Who will manage data preservation and sharing?	<b>Giovanni Samaey</b>
Who will update and implement this DMP?	<b>Pieter Vanmechelen</b>