# FWO DMP Template - Flemish Standard Data Management Plan

## Version KU Leuven

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 | **Steven Van Belleghem 0000-0001-9399-1007** |
| Contributor name(s) (+ ORCID) & roles | **Frederik Hendrickx 0000-0002-1176-0318 co-Promotor** |
| Project number [1] & title | Unravelling structural genome variation evolution by exploring functional pan-genomics |
| Funder(s) GrantID [2] | G020025N |
| Affiliation(s) | ☒ KU Leuven<br>☐ Universiteit Antwerpen<br>☐ Universiteit Gent<br>☐ Universiteit Hasselt<br>☐ Vrije Universiteit Brussel<br>☒ Other: RBINS<br>ROR identifier KU Leuven: 05f950310 |

---

[1] "Project number" refers to the institutional project number. This question is optional. 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.

| Please provide a short project description | Structural variations (SVs) in genomes, such as insertions, deletions, inversions, duplications, and translocations, are prevalent yet often overlooked in population genetic studies. Three critical questions remain largely unaddressed: 1) "How much SVs are present within and between populations?", 2) "What is the effect of SVs on genome functioning and phenotypic development?", and 3) "What is the evolutionary significance of SVs?". Recent technological advances have revolutionized our ability to cost-effectively and confidently identify all types of SVs. Here, we propose to apply these technologies to Pogonus (Europe) and Calosoma (Galápagos) beetles, which each possess a diversity of elaborate structural rearrangements linked to habitat-specific adaptation. By collecting chromosomal genome assemblies across their geographic ranges, we aim to initially delineate the comprehensive spectrum of SVs in both species. Subsequently, we will collect gene expression, chromatin accessibility, chromatin conformation, and DNA methylation data for key developmental stages and integrate these data into a composite pan-genome assembly. This approach will allow us to assess how SVs result in changes to gene content, gene regulation and development. Finally, we will reconstruct the evolution of these SVs across the unique evolutionary history of each species. Ultimately, this project is expected to enhance our understanding of the rol and evolution of SVs in natural populations. |

## 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 [3].

| | | | | *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 |
| | | ☒ Generate new data <br> ☒ Reuse existing data | ☒ Digital <br> ☐ Physical | ☐ Audiovisual <br> ☒ Images <br> ☐ Sound <br> ☒ Numerical <br> ☒ Textual <br> ☐ Model <br> ☐ Software <br> ☐ Other: | DNA/RNA sequences, gene expression profiles, population genetics images & metrics | ☐ < 1 GB <br> ☐ < 100 GB <br> ☐ < 1 TB <br> ☐ < 5 TB <br> ☒ > 5 TB <br> ☐ NA | |
| DNA-seq short | Short-read DNA sequences of beetles | Generate new data | Digital | Observational | .fastq | <5Tb | |
| DNA-seq long | Long-read DNA sequences of beetles | Generate new data | Digital | Observational | .fastq | >5Tb | |
| RAD-tag | RAD-tag sequencing for association mapping | Generate new data | Digital | Observational | .fastq | <1Tb | |
| ATAC-seq | Chromatin | Generate new data | Digital | Observational | .fastq | <1Tb | |

---

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

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | accessibility profiles of beetles | | | | | | |
| RNA-seq | Gene expression profiles of beetles, developmental changes | Generate new data | Digital | Observational | .fastq | <1Tb | |
| Iso-seq | Long read RNA sequencing | Generate new data | Digital | Observational | .fastq | <1Tb | |
| Phenotypic measurements | Phenotypic measurements of beetles | Generate new data | Digital | Observational | .txt | <1Gb | |
| Processed sequenced data | Sequence data mapped to reference genome | Generate new data | Digital | Observational | .bam | <5Tb | |
| Genotype calls | Genotypes from sequence data relative to reference genome | Generate new data | Digital | Observational | .vcf | <5Tb | |
| Omni-C data | Chromatin contacts in the genome | Generate new data | Digital | Observational | .fastq | <1Tb | |
| DNA methylation | Methylation status obtained from Pacbio | Generate new data | Digital | Observational | .fastq | <1Tb | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | data | | | | | | |
| Pan genome | Pan genome assembly with SVs | Generate new data | Digital | Observational | .fasta .xmfa .maf | <1Tb | |
| Bioinformatic analyses pipeline | Data processing scripts | Generate new data | Digital | Observational | .py .sh .R | <1Gb | |

*GUIDANCE:*
*The data description forms the basis of your entire DMP, so make sure it is detailed and complete. It includes digital and physical data and encompasses the whole spectrum ranging from raw data to processed and analysed data including analysis scripts and code. Physical data are all materials that need proper management because they are valuable, difficult to replace and/or ethical issues are associated. Materials that are not considered data in an RDM context include your own manuscripts, theses and presentations; documentation is an integral part of your datasets and should described under documentation/metadata.*
*RDM Guidance on data*

| 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. | **Pogonus chalceus original draft genome: NCBI BioProject accession code: PRJNA381601**<br>**Genome assembly provided by Sanger institute: Tree of Life QC: Species Report - Pogonus chalceus**<br>**Calosoma resequencing data: NCBI BioProject accession code: PRJNA706924**<br>**Calosoma draft genome: NVBI: JAGJTL000000000** |
|---|---|
| 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, refer to specific datasets or data types when appropriate and provide the relevant ethical approval number. | ☐ Yes, human subject data; provide SMEC or EC approval number:<br>☐ Yes, animal data; provide ECD reference number:<br>☐ Yes, dual use; provide approval number:<br>☒ No<br>Additional information: |

| | |
|---|---|
| Will you process personal data[4]? If so, please refer to specific datasets or data types when appropriate and provide the KU Leuven or UZ Leuven privacy register number (G or S number). | ☐ Yes (provide PRET G-number or EC S-number below) <br> ☒ No <br> Additional information: |
| Does your work have potential for commercial valorization (e.g. tech transfer, for example spin-offs, commercial exploitation, …)? <br> 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)? <br> If so, please explain to what data they relate and what restrictions are in place. | ☐ Yes <br> ☒ No <br> If yes, please explain: |
| 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

---

[4] See Glossary Flemish Standard Data Management Plan

| | |
|---|---|
| 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).<br><br>*RDM guidance on documentation and metadata*. | **For genomic sequence data (including short- and long-read DNA, RNA, and ATAC), NCBI metadata standards for data submission will be used. This includes:**<br>**- Sample ID (NCBI accession and lab ID)**<br>**- Sample origin**<br>**- Date of sampling**<br>**- Location of sampling**<br>**- Tissue**<br>**- Life stage**<br>**- Sex**<br>**- Sequence type and sequence platform**<br>**Data processing pipelines stored on GitHub will be accompanied with readme.txt files.**<br>**Phenotypic measurements of beetles will be stored in excel sheets including clone ID, clone origin, phenotypic value, experimenter and date of experiment.**<br>**An electronic lab book in which all metadata is shared and curated to common standards will be kept through a lab group on Microsoft teams.**<br>**I will also keep printed version of my protocols for the fieldwork, the molecular biology experiments and bioinformatic pipelines.** |

| | |
|---|---|
| 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>**1. Raw sequencing reads for: NCBI metadata standards for data submission ('SRA_metadata.xlsx' and 'invertebrate.1.0.xlsx' file)**<br>**2. Reference genome assemblies: NCBI metadata standards for data submission**<br><br>If no, please specify (where appropriate per dataset or data type) which metadata will be created:<br><br>**3. Processed sequenced data: external harddrives (in duplicates) and KU Leuven VSC Tier-1 Data Storage for at least 10 years**<br>**4. Bioinformatic analyses pipeline: readme.txt files with descriptions of the pipelines. Stored on Github, external harddrive and Microsoft cloud service.**<br>**5. Phenotypic measurements of beetles: excel sheets including clone ID and phenotypic value. Stored on Github, external harddrive and Microsoft cloud service.** |

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

| | |
|---|---|
| Where will the data be stored?<br><br>*Consult the [interactive KU Leuven storage guide](#) to find the most suitable storage solution for your data.* | ☐ Shared network drive (J-drive)<br>☒ Personal network drive (I-drive)<br>☒ OneDrive (KU Leuven)<br>☐ Sharepoint online<br>☐ Sharepoint on-premis<br>☒ Large Volume Storage<br>☐ Digital Vault<br>☐ Other: |
| How will the data be backed up?<br><br>*WHAT STORAGE AND BACKUP PROCEDURES WILL BE IN PLACE TO PREVENT DATA LOSS?* | ☒ Standard back-up provided by KU Leuven ICTS for my storage solution<br>☒ Personal back-ups I make (specify)<br>☐ Other (specify)<br><br>**At Ku Leuven**<br>**Personal computer**<br>**External harddrives**<br>**KU Leuven VSC Tier-1 Data Storage**<br><br>**Online:**<br>**Public submission of sequence read data to NCBI (with embargo until publication)**<br>**Github** |
| 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>**We have already secured 32 TB of storage on the iRODs server in Tier 1 of KU Leuven for our lab and purchased a 5TB personal external hard drives (both currently in use).**<br><br>If no, please specify: |

| | |
|---|---|
| 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.*<br>*Guidance on security for research data* | **The KU Leuven VSC Tier-1 Data Storage has a strict permissions system (read/write/modify), which will be limited to the primary collectors of the data.**<br>**External hardddrives will be stored safely lockers in the lab space or in the office of Steven Van Belleghem.**<br>**Personal computers will be protected password protected.** |
| What are the expected costs for data storage and backup during the research project? How will these costs be covered? | **Costs for data storage are included in the funded project and allow purchasing external harddrives and KU Leuven VSC Tier-1 Data Storage until the end of the project.** |

## 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...).<br><br>*Guidance on data preservation* | ☒ All data will be preserved for 10 years according to KU Leuven RDM policy<br>☐ All data will be preserved for 25 years according to CTC recommendations for clinical trials with medicinal products for human use and for clinical experiments on humans<br>☐ Certain data cannot be kept for 10 years (explain)<br><br><br>**All raw sequencing will be publicly stored >25 years.**<br>**1. Raw sequencing reads for:**<br>**- Short-read DNA sequences (.fastQ)**<br>**- Long-read DNA sequences (.fastQ)**<br>**- ATAC-seq (.fastQ)**<br>**- RNA-seq (.fastQ)**<br>**- Hi-C (.fastQ)**<br>**2. Reference genome assemblies (.fasta)**<br><br>**Processed sequencing data will be stored for at least 10 years using the KU Leuven VSC Tier-1 Data Storage and external harddrives**<br>**3. Processed sequence data:**<br>**- read alignments to genomes (.bam)**<br>**- gene expression counts (.counts)**<br>**- open chromatin regions (.bw .bed)**<br>**- PAN genome alignments (.maf)**<br><br>**Protocols, pipelines and phenotypic measurements will be stored indefinitely on Github, onedrive, and external harddrives, and as supplementary material with published studies.**<br>**4. Bioinformatic analyses pipeline**<br>**5. Phenotypic measurements of beetles** |

| | |
|---|---|
| Where will these data be archived (stored and curated for the long-term)?<br><br>*Dedicated data repositories are often the best place to preserve your data. Data not suitable for preservation in a repository can be stored using a KU Leuven storage solution, consult the interactive KU Leuven storage guide.* | ☐ KU Leuven RDR<br>☒ Large Volume Storage (longterm for large volumes)<br>☐ Shared network drive (J-drive)<br>☐ Other (specifiy):<br><br>**1. Raw sequencing reads for: NCBI (National Center for Biotechnology Information)**<br>**2. Reference genome assemblies: NCBI (National Center for Biotechnology Information)**<br>**3. Processed sequence data: external harddrives and KU Leuven VSC Tier-1 Data Storage**<br>**4. Bioinformatic analyses pipeline: GitHub**<br>**5. Phenotypic measurements of beetles: Github and publication supplements** |
| What are the expected costs for data preservation during the expected retention period? How will these costs be covered? | **NCBI: Free**<br>**GitHub: Free (up to 2gb, which is sufficient for storing pipelines and scripts)**<br>**KU Leuven VSC Tier-1 Data Storage: ~35 euro per Tb. 32Tb is supported by the FWO Data component of the Flemish Tier-1 supercomputing platform until 2027.** |

| |
|---|
| **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, as open data<br>☐ Yes, as embargoed data (temporary restriction)<br>☒ Yes, as restricted data (upon approval, or 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. | Processed sequence data will be stored in at least two location using external harddrives and KU Leuven VSC Tier-1 Data Storage. These data can be reproduced from raw reads and the available processing pipelines (shared on GitHub). These processed data types are not typically uploaded to NCBI for storage and sharing, but after publication, these data will also be freely available upon request and shared through web transfers. |
| 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. | ☐ KU Leuven RDR<br>☒ Other data repository (specify)<br>☐ Other (specify)<br><br>**1. Raw sequencing reads for: NCBI (National Center for Biotechnology Information)**<br>**2. Reference genome assemblies: NCBI (National Center for Biotechnology Information)**<br>**3. Processed sequenced data: external harddrives and KU Leuven VSC Tier-1 Data Storage**<br>**4. Bioinformatic analyses pipeline: GitHub**<br>**5. Phenotypic measurements of beetles: Github and publication supplements** |
| When will the data be made available? | ☒ Upon publication of research results<br>☐ Specific date (specify)<br>☐ Other (specify) |
| 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>*Check the RDR guidance on licences for data and software sources code or consult the License selector tool to help you choose.* | ☒ CC-BY 4.0 (data)<br>☐ Data Transfer Agreement (restricted data)<br>☐ MIT licence (code)<br>☐ GNU GPL-3.0 (code)<br>☐ Other (specify) |

| | |
|---|---|
| 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, a PID will be added upon deposit in a data repository<br>☐ My dataset already has a PID<br>☐ No |
| What are the expected costs for data sharing? How will these costs be covered? | **NCBI: Free**<br>**GitHub: Free (up to 2gb, which is sufficient for storing pipelines and scripts)**<br>**KU Leuven VSC Tier-1 Data Storage: ~35 euro per Tb. 32Tb is supported by the FWO Data component of the Flemish Tier-1 supercomputing platform until 2027.** |

| 7. Responsibilities | |
|---|---|
| Who will manage data documentation and metadata during the research project? | **Steven Van Belleghem / Maria Madrid** |
| Who will manage data storage and backup during the research project? | **Steven Van Belleghem / Maria Madrid / Frederik Hendrickx** |
| Who will manage data preservation and sharing? | **Steven Van Belleghem / Maria Madrid / Frederik Hendrickx** |
| Who will update and implement this DMP? | **Steven Van Belleghem  / Maria Madrid / Frederik Hendrickx** |