

**DATA MANAGEMENT PLAN (DMP)**

**FWO PROJECT: 11N3122N**

**Thi Xuan Ai Pham, KU Leuven**

**Laboratory for Epigenetic Reprogramming**

**Stem Cell Institute, KU Leuven**

**Principle investigator: Vincent Pasque**

**NOVEL STEM CELL-BASED HUMAN EMBRYO MODELS AND SINGLE-CELL TECHNOLOGIES TO STUDY EARLY EMBRYONIC DEVELOPMENT**

1. Data description	
Will you generate/collect new data and/or make use of existing data?	X Generate new data X Reuse existing data

## DMP for 1170722N

<p>Describe the origin, type and format of the data (per dataset) and its (estimated) volume</p>	<p>I will generate, collect, process, analyze and store the data listed below, as detailed in the project description. (not relevant to this section)</p> <p>The following datasets will be generated:</p> <p><b>1. General</b> The major aim of the project is the analysis of single-cell epigenomic and single-cell transcriptomic data. This data will be acquired on in vitro models, including cells obtained by reprogramming somatic cells into induced pluripotent stem cells or cells derived from induced pluripotent stem cells and/or embryonic stem cell lines.</p> <p><b>1. Experimental data</b>  <b>Dataset 1.1. – Digital images</b> Microscopy pictures, gel scans, graphs, illustrations, figures.  <b>Dataset 1.2. – Cytometry data</b> Flow Cytometry and fluorescence-activated cell sorting (FACS) data  <b>Dataset 1.3. – Omics data</b> Genomics, transcriptomics, epigenomics data.</p> <p><b>2. Derived and compiled data</b>  <b>Dataset 2.1 – Research documentation</b> Research documentation generated by myself and collaborators or collected from online sources and from collaborators, including ethical approval documents, laboratory notes, protocols.  <b>Dataset 2.2 – Manuscripts</b> The results will be published as BioRxiv preprints and articles in peer reviewed journals.  <b>Dataset 2.3 – Algorithms and scripts</b> Algorithms and scripts to investigate the single-cell data and to integrate different omics layers (transcriptome/epigenome) will be designed.  <b>Dataset 2.4 – Software</b> The combined set of scripts, algorithms, visualization tools &amp; computer programs.</p> <p><b>3. Canonical data</b>  <b>Dataset 3.1 – Nucleic acid sequences</b></p> <p>These datasets are the backbone of my PhD project, but after publishing will also represent an important source of information for others studying early human embryonic development.</p> <p>Data will be stored in the following formats:</p>
--	--

## DMP for 1170722N

	<ul style="list-style-type: none"><li>- Text files: Plain text data (Unicode, .txt), MS Word (.doc/.docx), Adobe Portable Document Format (.pdf), LaTeX (.tex) format;</li><li>- Quantitative tabular data: comma-separated value files (.csv), tab-delimited file (.tsv), delimited text (.txt), MS Excel (.xls/.xlsx);</li><li>- Digital images in raster formats: uncompressed TIFF (.tif/.tiff), JPEG (.jpg), JPEG 2000 (.jp2), Adobe Portable Document Format (.pdf), bitmap (.bmp), .gif;</li><li>- Digital images in vector formats: scalable vector graphics (.svg), encapsulated postscript (.eps), Scalable Vector Graphics (.svg), Adobe Illustrator (.ai);</li><li>- Flow cytometry data: Flow Cytometry Standard (.fcs);</li><li>- Nucleotide sequences: raw sequence data trace (.ab1), text-based format (.fasta(.gz)/.fa(.gz)) and accompanying QUAL file (.qual), Genbank format (.gb/.gbk);</li><li>- Next generation sequencing raw data: binary base call format (.bcl)</li><li>- Sequence alignment data: (.sam), .bam</li><li>- Coverage data: .bed, .bg, .bedGraph, .bw, .bigwig</li><li>- Structural variations data: .vcf(.gz), .bcf</li><li>- Read/UMI count data: .tsv(.gz), Matrix Market format (.mtx), .loom, .rds(.gz)</li><li>- Nucleic acid samples resulting from (single-cell) nucleic acid amplification, or sequence library preparations will be stored in labeled tubes or SBS plates in -20C freezers. We have electronic laboratory databases that will keep the physical storage address of these samples.</li></ul> <p>Raw as well as processed data will be submitted to a public repository in the aforementioned described standard formats, to enable sharing and long-term validity of the data.</p>
--	--

2. Ethical and legal issues	
Will you use personal data? If so, shortly describe the kind of personal data you will use AND add the reference to your file in your host institution's privacy register.	<input checked="" type="checkbox"/> No <input type="checkbox"/> Yes
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, add the reference to the formal approval by the relevant ethical review committee(s).	<input type="checkbox"/> No <input checked="" type="checkbox"/> Yes <ul style="list-style-type: none"> <li>- Ethical approval S64962: Characterization and functional assessment of the regulatory landscape of embryonic lineages using blastoids as a model.</li> <li>- Ethical approval (S66375) to work on human embryos with Dr. Laurent David (in Nantes)</li> <li>- Ethical approval to work on human embryos with Dr. Hilde Van Devele (in Brussels) (S66184): "Chromatin regulators perturbations in human embryos (CREPE). Human embryo inhibitor (PRC2) treatment"</li> </ul>

## DMP for 1170722N

<p>Does your work possibly result in research data with potential for tech transfer and valorization?</p> <p>Will IP restrictions be claimed for the data you created? If so, for what data and which restrictions will be asserted?</p>	<p><input type="checkbox"/> No <input checked="" type="checkbox"/> Yes</p> <p>I do not exclude that the proposed work could result in research data with potential for tech transfer and valorization. The KU Leuven has a policy to actively monitor research data for such potential. If there is substantial potential, the invention will be thoroughly assessed, and in a number of cases the invention will be IP protected (mostly patent protection or copyright protection). As such the IP protection does not withhold the research data from being made public. In the case a decision is taken to file a patent application it will be planned so that publications need not be delayed.</p>
<p>Do existing 3<sup>rd</sup> party agreements restrict dissemination or exploitation of the data you (re)use? If so, to what data do they relate and what restrictions are in place?</p>	<p><input checked="" type="checkbox"/> No <input type="checkbox"/> Yes</p> <p>No third-party agreement restricts dissemination or exploitation of the data from this project.</p>

### 3. Documentation and metadata

<p>What documentation will be provided to enable understanding and reuse of the data collected/generated in this project?</p>	<p>Documentation will consist of notes in the electronic laboratory notebook (E-notebook) that refer to specific datasets. Those notes will describe the biological samples used, experimental setup and protocols used, sequences generated, the links to the specific computer location as well as the names of the respective datasets.</p> <p>I will also maintain a metadata sheet with the connection between lab samples and files on our data storage, so that data files, lab samples, and experimental notes remain properly linked.</p> <p>Algorithms, scripts and software usage will be documented, e.g. using Jupyter Notebooks. When scripts, algorithms and software tools are finalized, they will be additionally described in manuscripts and on GitHub (see <a href="https://www.github.com/pasquelab">www.github.com/pasquelab</a> for our previous datasets and scripts).</p>
---	---

## DMP for 1170722N

<p>Will a metadata standard be used? If so, describe in detail which standard will be used. If not, state in detail which metadata will be created to make the data easy/easier to find and reuse.</p>	<p><input type="checkbox"/> No  <input checked="" type="checkbox"/> Yes</p> <p>Sequencing data types require particular metadata, such as data submitted to EGA, GEO, SRA, ArrayExpress, or ENA. Local data that is not (yet) submitted to these resources will be based on generalized metadata schema such as Dublin Core or DataCite, including the following elements:</p> <ul style="list-style-type: none"> <li>• Title: free text</li> <li>• Creator: Last name, first name, organization</li> <li>• Date and time reference</li> <li>• Subject: Choice of keywords and classifications</li> <li>• Description: Text explaining the content of the data set and other contextual information needed for the correct interpretation of the data, the software(s) (including version number) used to produce and to read the data, the purpose of the experiment, etc.</li> <li>• Format: Details of the file format,</li> <li>• Resource Type: data set, image, audio, etc.</li> <li>• Identifier: DOI (when applicable)</li> <li>• Access rights: closed access, embargoed access, restricted access, open access.</li> </ul> <p>When depositing data in a local or public repository, the final dataset will be accompanied by this information under the form of a README.txt document. This file will be located in the top level directory of the dataset and will also list the contents of the other files and outline the file-naming convention used (see section 7 below). This will allow the data to be understood by other members of the laboratory and add contextual value to the dataset for future reuse.</p>
--	---

4. Data storage & backup during the FWO project	
Where will the data be stored?	<p>Digital files will be stored on KU Leuven servers.</p> <ul style="list-style-type: none"> <li>- Omics data: omics data generated during the project will either be stored on KU Leuven servers or on the Flemish Supercomputer Centre (VSC), initially in the staging + archive area and later only in the archive area (archive is mirrored).</li> <li>- Vectors: As a general rule at least two independently obtained clones will be preserved for each vector, both under the form of purified DNA (in -20°C freezer) and as a bacteria glycerol stock (-80°C). All published vectors and the associated sequences will be sent to the non-profit plasmid repository Addgene, which will take care of vector storage and shipping upon request.</li> <li>- Algorithms, scripts and software: All the relevant algorithms, scripts and software code driving the project will be stored in private online git repositories of the PI (e.g., <a href="https://github.com/pasquelab">https://github.com/pasquelab</a>). As soon as the manuscript is publicly available, the repository will be changed to a public repository.</li> <li>- Nucleic acid and protein sequences: All nucleic acid and protein sequences generated during the project will be stored on KU Leuven servers. Upon publication, all sequences supporting a manuscript will be made publicly available via repositories such as the GenBank database or the European Nucleotide Archive (nucleotide sequences from primers / new genes / new genomes), NCBI Gene Expression Omnibus (microarray data / RNA-seq data / CHIPseq data), the Protein Database (for protein sequences), the EBI European Genome-phenome Archive (EGA) for personally identifiable (epi)genome and transcriptome sequences.</li> <li>- Other data files (Digital images, cytometry data, etc..) will be stored on local KU Leuven servers and PI computers.</li> </ul>
How will the data be backed up?	<p>KU Leuven drives are backed-up according to the following scheme:</p> <ul style="list-style-type: none"> <li>- data stored on the “L-drive” is backed up daily using snapshot technology, where all incremental changes in respect of the previous version are kept online; the last 14 backups are kept.</li> <li>- data stored on the “J-drive” is backed up hourly, daily (every day at midnight) and weekly (at midnight between Saturday and Sunday); in each case the last 6 backups are kept.</li> <li>- All omics data stored on the Flemish Supercomputer Centre (VSC) will be transferred on a monthly basis to the archive area which is mirrored.</li> </ul>

## DMP for 1170722N

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.	<input type="checkbox"/> No <input checked="" type="checkbox"/> Yes <p>There is sufficient storage and back-up capacity on all KU Leuven servers:</p> <ul style="list-style-type: none"> <li>- the “L-drive” is an easily scalable system, built from General Parallel File System (GPFS) cluster with NetApp series storage systems, and a CTDB samba cluster in the front-end.</li> <li>- the “J-drive” is based on a cluster of NetApp FAS8040 controllers with an Ontap 9.1P9 operating system.</li> <li>- the Staging and Archive on VSC are also sufficiently scalable (petabyte scale)</li> </ul>
What are the expected costs for data storage and backup during the project? How will these costs be covered?	<p>The total estimated cost of data storage during the project is ~9,000 EUR. This estimation is based on the following costs:</p> <ul style="list-style-type: none"> <li>- The costs of digital data storage are as follows: 128,39€/TB/Year for the “L-drive” and 519EUR/TB/Year for the “J-drive”.</li> <li>- The cost of VSC archive is 70 EUR/TB/Year, and staging 130EUR/TB/Year.</li> <li>- We expect costs to drop slightly during the coming four years. Additional budget for compute and data storage is budgeted for in ongoing projects, and will be costed in complementary project applications.</li> </ul>
Data security: how will you ensure that the data are securely stored and not accessed or modified by unauthorized persons?	<ul style="list-style-type: none"> <li>- The “L-drive” and “J-drive” servers are accessible only by laboratory members and are mirrored in the second ICTS datacenter for business continuity and disaster recovery so that a copy of the data can be recovered within an hour.</li> <li>- The VSC storage is only accessible to VSC accounts, and specifically our volume is only accessible to group members.</li> <li>- No personal data will be stored on the VSC nor local servers.</li> </ul>



<b>5. Data preservation after the end of the FWO project</b> KU Leuven expects that data generated during the project are retained for a period of minimally 5 years after the end of the project, in as far as legal and contractual agreements allow.	
Which data will be retained for the expected 5 year period after the end of the project? In case only a selection of the data can/will be preserved, clearly state the reasons for this (legal or contractual restrictions, physical preservation issues, ...).	The minimum preservation term of 5 years after the end of the project will be applied to all datasets.
Where will these data be archived (= stored for the long term)?	<p>As a general rule, datasets will be made openly accessible, whenever possible via existing platforms that support FAIR data sharing (<a href="http://www.fairsharing.org">www.fairsharing.org</a>), at the latest at the time of publication or preprint deposition.</p> <p>For all other datasets, long term storage will be ensured as follows:</p> <ul style="list-style-type: none"> <li>- Large sequencing/omics data: will be stored on VSC archive</li> <li>- Small digital files: files will be stored on the "L-drive".</li> <li>- Developed algorithms and software will be stored on VSC archive and/or L-drive, as well on public repositories such as Github.com.</li> <li>- Third-party software and algorithms that are used are referenced by their version numbers (e.g., in our Jupyter notebooks), and are installed as modules on the VSC and/or containers (Docker, Singularity) on the VSC, to ensure reproducibility.</li> </ul>

## DMP for 1170722N

What are the expected costs for data preservation during these 5 years? How will the costs be covered?	<p>The total estimated cost of data storage during the 5 years after the end of the is 2,500 EUR. This estimation is based on 70EUR/Tb/year. The storage after the project is much smaller because during the project a large working space is needed, and post-publication data will be made accessible via open access platforms.</p> <p>The costs for this data preservation will be upfront paid for from this FWO project</p>
--	--

### 6. Data sharing and reuse

Are there any factors restricting or preventing the sharing of (some of) the data (e.g. as defined in an agreement with a 3 <sup>rd</sup> party, legal restrictions)?	<p><input checked="" type="checkbox"/> No</p> <p><input type="checkbox"/> Yes</p>
Which data will be made available after the end of the project?	<p>The PI in the present project is committed to publish research results to communicate them to peers and to a wide audience. All research outputs supporting publications will be made openly accessible. Depending on their nature, some data may be made available prior to publication, either on an individual basis to interested researchers and/or potential new collaborators, or publicly via repositories (e.g. negative data).</p>

## DMP for 1170722N

Where/how will the data be made available for reuse?	<p>X In an Open Access repository</p> <p>X In a restricted access repository</p> <p><input type="checkbox"/> Upon request by mail</p> <p>X Other (specify): open-access publications in peer-reviewed journals, including supplemental information</p> <p>As a general rule, datasets will be made openly accessible via existing platforms that support FAIR data sharing (<a href="http://www.fairsharing.org">www.fairsharing.org</a>). Sharing policies for specific research outputs are detailed below:</p> <ul style="list-style-type: none"> <li>- Omics datasets will be deposited in open access repositories such as the EMBL-EBI platform for genomics and epigenomics data, or the NCBI Gene Expression Omnibus (GEO), the EBI ArrayExpress databases for functional genomics data or the EBI European Genome-phenome Archive (EGA) for personally identifiable genetic and phenotypic data. Data at EGA will be collected from individuals whose consent agreements authorise data release only for specific research use to bona fide researchers.</li> <li>- Other digital datasets that support publications (including image, cytometry data, and simulation data) will be made publicly available via an open research data platform such as Mendeley Data or Zenodo.</li> <li>- Research documentation: All protocols used to generate published data will be described in the corresponding manuscript(s), and the related documentation will be included as supplementary information. These data and all other documents (daily logs, raw data) deposited in the E-Notebook are accessible to the PI and the research staff, and will be made available upon request.</li> <li>- Manuscripts: All scientific publications will be shared openly. Manuscripts submitted for publication will be deposited in a pre-print server such as bioRxiv, arXiv, Nature Precedings or ASAPbio). At the time of publication, research results will be summarized on the PI's website (<a href="http://www.kuleuven.be/pasquelab">http://www.kuleuven.be/pasquelab</a>) and post-print pdf versions of publications will be made available there if allowed by copyright agreements, possibly after an embargo as determined by the publisher. Before the end of the embargo or in cases where sharing the post-print is not allowed due to copyright agreements, a pre-print version of the manuscript will be made available. Publications will also be automatically added to our institutional repository, Lirias 2.0, based on the authors name and ORCID ID (the metadata will be added, not the full manuscripts).</li> <li>- Algorithms, scripts and software: All the relevant algorithms, scripts and software code driving the project will be made available to public repositories such as <a href="http://www.github.com/pasquelab">www.github.com/pasquelab</a></li> <li>- Nucleic acid and protein sequences: All nucleic acid and protein sequences generated during the project will made publicly available</li> </ul>
--	--

## DMP for 1170722N

	<p>via repositories such as the GenBank database or the European Nucleotide Archive (nucleotide sequences from primers / new genes / new genomes), NCBI Gene Expression Omnibus (microarray data / RNA-seq data / CHIPseq data), the Protein Database (for protein sequences), the EBI European Genome-phenome Archive (EGA) for personally identifiable (epi)genome and transcriptome sequences.</p> <ul style="list-style-type: none"><li>- Data that do not support publication will be either deposited in an open access repository or made available upon request by email.</li></ul>
--	---

## DMP for 1170722N

When will the data be made available?	<ul style="list-style-type: none"><li><input type="checkbox"/> Immediately after the end of the project</li><li>X Upon publication of the research results</li><li><input type="checkbox"/> After an embargo period.</li></ul> <p>As a general rule all research outputs will be made openly accessible at the latest at the time of publication. No embargo will be foreseen unless imposed e.g. by pending publications, potential IP requirements – note that patent application filing will be planned so that publications need not be delayed - or ongoing projects requiring confidential data. In those cases, datasets will be made publicly available as soon as the embargo date is reached.</p>
---------------------------------------	---

## DMP for 1170722N

Who will be able to access the data and under what conditions?	<p>Whenever possible, datasets and the appropriate metadata will be made publicly available through repositories that support FAIR data sharing. As detailed above, metadata will contain sufficient information to support data interpretation and reuse, and will be conform to community norms. These repositories clearly describe their conditions of use (typically under a Creative Commons CC0 1.0 Universal (CC0 1.0) Public Domain Dedication or an ODC Public Domain Dedication and Licence, with a material transfer agreement when applicable). Interested parties will thereby be allowed to access data directly, and they will give credit to the authors for the data used by citing the corresponding DOI. A CC-BY license will be opted for when possible. For data shared directly by the PI, a material transfer agreement (and a non-disclosure agreement if applicable) will be concluded with the beneficiaries in order to clearly describe the types of reuse that are permitted.</p> <p>For KU Leuven data submitted to the EBI European Genome-phenome Archive (EGA), which operates under controlled access, the data access/submission requests will be received by the operational committee and processed in consultation with the PIs produced data. The advisory committee will provide general guidance in terms of policies and will be referred to in handling controversial cases.</p>
What are the expected costs for data sharing? How will these costs be covered?	<p>It is the intention to minimize data management costs by implementing standard procedures (internally generated operational procedures (<a href="https://en.wikipedia.org/wiki/Standard_operating_procedure">https://en.wikipedia.org/wiki/Standard_operating_procedure</a>)) e.g. for metadata collection and file storage and organization from the start of the project, and by using free-to-use data repositories and dissemination facilities whenever possible. Data management costs will be covered by the laboratory budget.</p>

### 7. Responsibilities

Who will be responsible for the data documentation & metadata?	<p>Metadata will be documented by the research and technical staff at the time of data collection and analysis, by taking careful notes in the electronic laboratory notebook (E-notebook) that refer to specific datasets.</p>
--	---

## DMP for 1170722N

Who will be responsible for data storage & back up during the project?	I will ensure data storage and back up, with support from ICTS, HPC, and gbiomed-IT staff.
Who will be responsible for ensuring data preservation and sharing?	The PI is responsible for data preservation and sharing, with support from ICTS, HPC, and gbiomed-IT staff.
Who bears the end responsibility for updating & implementing this DMP?	The PI is ultimately responsible for all data management during and after data collection, including implementing and updating the DMP.