

Data Management Plan

FIRM-CENTRIC TECHNOLOGY ROADMAPS (FTRMS) AND TECHNOLOGY VALUATION – Bruno Cassiman (G010324N)

1. General Project Information	
Name Grant Holder & ORCID	Bruno Cassiman (0000-0001-9602-6755)
Project number & title	G010324N FIRM-CENTRIC TECHNOLOGY ROADMAPS (FTRMS) AND TECHNOLOGY VALUATION
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: ROR identifier KU Leuven: 05f950310
Please provide a short project description	<p>The project proposes to develop Firm-centric Technology Roadmaps (FTRMs) and to improve technology valuation models based on them. Anticipating technology evolutions has become a critical business skill. Evaluating a potential technology opportunity in relation to the organization's own technological capabilities is critical for success. We address three gaps in the literature. First, we develop a methodology for quantifying Technology Roadmaps (TRMs) and white spaces therein based on patent and publication information. Second, we link these TRMs to critical firm-specific scientific and technological capabilities. Third, we develop a method to estimate the economic value of technologies based on these FTRMs. Altogether, this will allow answering important economic and business strategy questions related to technology development and entry-timing decisions of the organization. To this end, we leverage recent Natural Language Processing and Machine Learning methodologies based on patent and publication text.</p>

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.

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
Patentsview data	Descriptive and full text data from USPTO patents.	<input type="checkbox"/> Generate new data <input checked="" type="checkbox"/> Reuse existing data	<input checked="" type="checkbox"/> Digital <input type="checkbox"/> Physical	<input type="checkbox"/> Audiovisual <input type="checkbox"/> Images <input type="checkbox"/> Sound <input checked="" type="checkbox"/> Numerical <input checked="" type="checkbox"/> Textual <input type="checkbox"/> Model <input type="checkbox"/> Software <input type="checkbox"/> Other:	.csv Python file	<input type="checkbox"/> < 1 GB <input checked="" type="checkbox"/> < 100 GB <input type="checkbox"/> < 1 TB <input type="checkbox"/> < 5 TB <input type="checkbox"/> > 5 TB <input type="checkbox"/> NA	/
PATSTAT	More comprehensive descriptives and full text data on patents (beyond USPTO)	Reuse existing data	Digital	Numerical and textual	.csv Python file	<1 TB	/
KPSS data	Results of commercial value of patents from the findings of Kogan et al. (2017)	Reuse existing data	Digital	Numerical	.dta Python file	<100 GB	/
DISCERN data	database of publicly listed	Reuse existing data	Digital	Numerical	.dta	<100 GB	/

	U.S. headquartered firms matched to assignees of patents from the United States Patent and Trademark Office (USPTO) and scientific publications from the Web of Science for the period 1980-2015.						
Orbis	We use financial and firm information from the Orbis historical database. (which covers a large chunk of non listed firms)	Reuse existing data	Digital	Numerical	.csv Python file	<100 GB	/
Compustat	We use financial and firm information from compustat, provided by WRDS	Reuse existing data	Digital	Numerical	.csv Python file	<1 GB	/

AHG data	We use data from Arts, Hou and Gomez (2020) on novelty of patents	Reuse existing data	Digital	Numerical and Textual	.csv	<100 GB	/
OpenAlex	We use data on Scientific publications from OpenAlex. OpenAlex is a bibliographic catalogue of scientific papers, authors and institutions accessible in open access mode	Reuse of existing data	Digital	Numerical and Textual	JSON (database)	<100 GB	/
Revenue and Patents data IMEC	We use data on contract revenue related to patents of IMEC	Reuse of existing data	Digital	Numerical	.csv Python file	<100 GB	/

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 be described under documentation/metadata.

[RDM Guidance on data](#)

<p>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.</p>	<p>Patentsview (csv): https://patentsview.org/download/data-download-tables PATSTAT: https://inspire.wipo.int/patstat-online KPSS (dta): https://github.com/KPSS2017/Technological-Innovation-Resource-Allocation-and-Growth-Replication-Kit ORBIS: https://bib.kuleuven.be/ebib/collectie/data/databanken/orbis_global DISCERN: https://zenodo.org/records/4320782 AHG: https://zenodo.org/records/3515985 OpenAlex: https://openalex.org/</p>
<p>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.</p>	<p><input type="checkbox"/> Yes, human subject data; provide SMEC or EC approval number: <input type="checkbox"/> Yes, animal data; provide ECD reference number: <input type="checkbox"/> Yes, dual use; provide approval number: <input checked="" type="checkbox"/> No Additional information: The data we use is either publicly available or licensable for any party (give a fee applies).</p>
<p>Will you process personal data¹? 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).</p>	<p><input type="checkbox"/> Yes (provide PRET G-number or EC S-number below) <input checked="" type="checkbox"/> No Additional information: We will strictly only use firm and patent information. Individual inventors will not be analyzed in this project.</p>
<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 If yes, please comment: We only use (semi)publicly available data which depending on the results could be used in the deployment of a tool to predict the value for patents and related whitespaces in technology and technology portfolios. To this end we will use the public (and licensed) data to train our models but relate the predictions to each firms proprietary data in the final version of the spin-off. (Which is beyond the scope of this project)</p>

¹ See Glossary Flemish Standard Data Management Plan

<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)?</p> <p>If so, please explain to what data they relate and what restrictions are in place.</p>	<p><input checked="" type="checkbox"/> Yes <input type="checkbox"/> No</p> <p>If yes, please explain: The Orbis and WRDS data has to be licensed by the KUL, currently they have a good relationships and they are licenses which the KUL is eager to extend. However, if at any point Orbis restricts access the KUL will have to destroy their copies and SQL application on Orbis. The access towards WRDS could also be restricted. However, in both those cases I would have to secure my own access, potentially through the use of a Bench Fee.</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?</p> <p>If so, please explain to what data they relate and which restrictions will be asserted.</p>	<p><input checked="" type="checkbox"/> Yes <input type="checkbox"/> No</p> <p>If yes, please explain: Since we are all in order with our licensed data (Orbis/WRDS) being licensed correctly and all other data being publicly available there is no issue with regards to the usage of this data for this research project. However, in the future possible commercial exploitation some new licenses will probably have to be negotiated in order to comply with the regulations from the data providers. We will also have access to proprietary data from IMEC.</p>

3. Documentation and Metadata

<p>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).</p> <p><i>RDM guidance on documentation and metadata.</i></p>	<p>I will provide a file that shows which data was collected from where and how it relates to each project. Furthermore, in each of the coding files the sources are mentioned and I will provide comments in each step in order to further enhance the reusability of the code as well as help with versioning. When publishing the code there will be README.txt files provided for each of the important steps to enhance replication. Replication will be made as easy as reading in the files sequentially to get the results (given that the data paths are adjusted of course).</p>
<p>Will a metadata standard be used to make it easier to find and reuse the data?</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: A README.txt file will provide the necessary context for the data as well as specify the 2 unique identifying variables which we will use. Firstly, we have Patstat's patent_id (As well as patentsview patent_id). Secondly firms are identified with the unique GVKEY provided by WRDS computstat. These unique identifiers should be enough to link all of our data together. Further data clarifications on potential variables we will develop will also be included in the README.txt file.</p>

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

<p>Where will the data be stored?</p> <p><i>Consult the interactive KU Leuven storage guide to find the most suitable storage solution for your data.</i></p>	<p> <input type="checkbox"/> Shared network drive (J-drive) <input type="checkbox"/> Personal network drive (I-drive) <input checked="" type="checkbox"/> OneDrive (KU Leuven) <input type="checkbox"/> Sharepoint online <input type="checkbox"/> Sharepoint on-premis <input checked="" type="checkbox"/> Large Volume Storage <input type="checkbox"/> Digital Vault <input checked="" type="checkbox"/> Other: MSI server (Patstat data) </p>
<p>How will the data be backed up?</p> <p><i>WHAT STORAGE AND BACKUP PROCEDURES WILL BE IN PLACE TO PREVENT DATA LOSS?</i></p>	<p> <input checked="" type="checkbox"/> Standard back-up provided by KU Leuven ICTS for my storage solution <input checked="" type="checkbox"/> Personal back-ups I make (the source data will be backed up on the KU Leuven onedrive but results and preprocessed data is also backed up on an external SSD drive) <input type="checkbox"/> Other (specify) </p>
<p>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.</p>	<p> <input checked="" type="checkbox"/> Yes <input type="checkbox"/> No </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.</i></p> <p>Guidance on security for research data</p>	<p>The data used will be kept securely on the MSI server as well as personal network drives from the KUL. These are regularly updated and maintained by the FEB ICT personnel as such they are some of the default options at the KUL. When collaborating with others I will use the Shared network drive :J, these can only be accessed via VPN connection to the KU Leuven network and follows the same maintenance and update scrutiny from the ICT personnel.</p>
<p>What are the expected costs for data storage and backup during the research project? How will these costs be covered?</p>	<p>If any further data storage is required we will tap into the bench fee provided by FWO for the yearly workings of this project. As such there will always be budget in the case that storage costs would suddenly increase.</p>

5. Data Preservation after the end of the Research Project	
<p>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> <p>Guidance on data preservation</p>	<p><input checked="" type="checkbox"/> All data will be preserved for 10 years according to KU Leuven RDM policy</p> <p><input type="checkbox"/> 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</p> <p><input type="checkbox"/> Certain data cannot be kept for 10 years (explain)</p>

<p>Where will these data be archived (stored and curated for the long-term)?</p> <p><i>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.</i></p>	<p><input checked="" type="checkbox"/> KU Leuven RDR</p> <p><input type="checkbox"/> Large Volume Storage (longterm for large volumes)</p> <p><input checked="" type="checkbox"/> Shared network drive (J-drive)</p> <p><input type="checkbox"/> Other (specify):</p>
<p>What are the expected costs for data preservation during the expected retention period? How will these costs be covered?</p>	<p>The expected costs would be the standard costs related to storing projects on KU Leuven RDR, this would mean that no extra costs would be imposed on this project.</p>

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 & RESTRICTED ACCESS. FOR MORE INFORMATION: https://wiki.surfnet.nl/display/STANDARDS/INFO-EU-REPO/#INFOEUREPO-ACCESSRIGHTS</i></p>	<p><input type="checkbox"/> Yes, as open data</p> <p><input type="checkbox"/> Yes, as embargoed data (temporary restriction)</p> <p><input checked="" type="checkbox"/> Yes, as restricted data (upon approval, or 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>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 <input checked="" type="checkbox"/> Yes, intellectual property rights <input type="checkbox"/> Yes, ethical aspects <input type="checkbox"/> Yes, aspects of dual use <input type="checkbox"/> Yes, other <input type="checkbox"/> No </p> <p>If yes, please specify: As part of the licensing agreement we are not allowed to share the Orbis data to third parties, the results of our analysis can be published freely. The data on contract revenue of IMEC is confidential and shared with the project, but not available for sharing.</p>
<p>Where will the data be made available? If already known, please provide a repository per dataset or data type.</p>	<p> <input checked="" type="checkbox"/> KU Leuven RDR <input type="checkbox"/> Other data repository (specify) <input type="checkbox"/> Other (specify) </p>
<p>When will the data be made available?</p>	<p> <input checked="" type="checkbox"/> Upon publication of research results <input type="checkbox"/> Specific date (specify) <input type="checkbox"/> Other (specify) </p>

<p>Which data usage licenses are you going to provide? If none, please explain why.</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>Check the RDR guidance on licences for data and software sources code or consult the License selector tool to help you choose.</p>	<input type="checkbox"/> CC-BY 4.0 (data) <input type="checkbox"/> Data Transfer Agreement (restricted data) <input type="checkbox"/> MIT licence (code) <input type="checkbox"/> GNU GPL-3.0 (code) <input type="checkbox"/> Other (specify)
<p>Do you intend to add a PID/DOI/accession number to your dataset(s)? If already available, please provide it here.</p> <p><i>INDICATE WHETHER YOU INTEND TO ADD A PERSISTENT AND UNIQUE IDENTIFIER IN ORDER TO IDENTIFY AND RETRIEVE THE DATA.</i></p>	<input checked="" type="checkbox"/> Yes, a PID will be added upon deposit in a data repository <input type="checkbox"/> My dataset already has a PID <input type="checkbox"/> No
<p>What are the expected costs for data sharing? How will these costs be covered?</p>	

7. Responsibilities	
Who will manage data documentation and metadata during the research project?	Bruno Cassiman
Who will manage data storage and backup during the research project?	Bruno Cassiman
Who will manage data preservation and sharing?	Bruno Cassiman
Who will update and implement this DMP?	Bruno Cassiman

