

FWO Research Project G0E8222N

Project Name Engineering superior microbial chassis strains and monoterpene synthases for highly efficient monoterpene production - FWO Research Project G0E8222N

Grant Title G0E8222N

Principal Investigator / Researcher Kevin Verstrepen

Project Data Contact Karin Voordeckers, karin.voordeckers@kuleuven.be

Description Monoterpenes are widely used in food, pharmaceuticals and cosmetics. The monoterpene linalool for example is mainly used in flavors and fragrances, green solvents, platform chemicals and active agents and its market size is projected to surpass USD 1.9 billion by 2024. However, traditional methods of monoterpene production - plant extraction and chemical synthesis - are not sustainable and inefficient, due to cumbersome processes, high energy consumption and low yields. Large-scale and low-cost microbial production of monoterpenes is therefore crucial to establish a sustainable and economically viable supply of monoterpenes. However, despite several research efforts, it has become clear that rational engineering of the heterologous monoterpene pathways alone will not lead to economically viable yields. This is mainly due to (i) fundamental shortcomings of current production strains, which have low pathway fluxes, (ii) monoterpene cytotoxicity and (iii) suboptimal natural monoterpene synthases. This project uses the complementary expertise of both partners to comprehensively tackle these problems, using high-throughput screening methods, cutting-edge synthetic biology tools and enzyme engineering. Ultimately, we will create a portfolio of chassis strains that can serve as 'plug-and-play' strains for production of a plethora of monoterpenes. As proof-of-concept, we will develop production strains for three high-value monoterpenes: linalool, limonene and cineole.

Institution KU Leuven

1. General Information

Name applicant

Kevin Verstrepen

FWO Project Number & Title

G0E8222N - Engineering superior microbial chassis strains and monoterpene synthases for highly efficient monoterpene production

Affiliation

- KU Leuven

VIB – KU Leuven Center for Microbiology

Laboratory for Systems Biology and Laboratory of Genetics and Genomics – prof. Kevin Verstrepen

2. Data description

Will you generate/collect new data and/or make use of existing data?

- Generate new data
- Reuse existing data

Describe in detail the origin, type and format of the data (per dataset) and its (estimated) volume. This may be easiest in a table (see example) or as a data flow and per WP or objective of the project. If you reuse existing data, specify the source of these data. Distinguish data types (the kind of content) from data formats (the technical format).

The research and technical staff will use, generate, collect, process, analyze and store the data listed below, as detailed in the project description.

The following datasets will be used:

1. Experimental data

Dataset 1.1. Strains (WP1)

In-lab strain collections of *Saccharomyces cerevisiae* and *Yarrowia lipolytica* strains, present in Verstrepen lab. Data type: biological samples, yeast strains. Dataset size: 1567 *Saccharomyces*

cerevisiae isolates and 17 *Yarrowia lipolytica* isolates.

The following datasets will be generated:

1. Experimental data

Dataset 1.1. - Digital images (WP1-WP2-WP3-WP4)

Microscopy pictures, gel scans, graphs, illustrations, figures. File type: .tiff and .jpeg. Estimated file size: 5 Gb

Dataset 1.2. - Cytometry data (WP2)

Flow Cytometry and fluorescence-activated cell sorting (FACS) data (for example, lipid content of microbial cells); File type: .fsc and .csv files. Estimated file size: 1 Gb

Dataset 1.3. - Spectroscopy data (WP1-WP2-WP3-WP4)

light scattering data (for example, absorbance measurements to track microbial growth in different conditions; DPPH measurements). File type: .csv file. Estimated file size: 1 Gb

Dataset 1.4. - Omics data (WP2)

Whole genome sequencing of yeast strains. File type: .ab; .fasta and .qual. Estimated file size: 20 Gb

Dataset 1.5. Chromatography data (WP1, WP2, WP3, WP4)

HPLC, GC-MS measurements of compounds (including monoterpenes) produced by yeast strains. File size: .csv and .ms file. Estimated file size: 1 Gb

Dataset 1.6. - Plasmids (WP2)

Estimated number of plasmids created during project: 200

Dataset 1.7. - Strains (WP1, WP2, WP3, WP4)

Bacteria strains, yeast strains created during this project. For bacteria and yeast strains, this includes lab strains, natural/clinical/industrial isolates, variant libraries, optimized clones (i.e. site-specific mutants constructed through genome engineering). Estimated number of strains generated during project: 15.000. Please note that a large fraction of these stem from variant libraries, which will be stored in pools rather than as individual samples.

2. Derived and compiled data (WP1, WP2, WP3, WP4)

Dataset 2.1 - Research documentation

Research documentation generated by the research and technical staff or collected from online sources and from collaborators, including laboratory notes and protocols. File type: . doc; Estimated file size: 2 Gb

Dataset 2.2 - Manuscripts

File type: .doc and .pdf; estimated file size: 1 Gb

Datset 2.3 Algorithms and scripts

File type: .r and .py, estimated file size: 1 Gb

These datasets represent an important source of information for the laboratory of prof. Kevin Verstrepen, for scientists, journalists and higher education teachers working in the field of microbiology, but also for non-profit organizations and industries active in this broad field.

3. Legal and ethical issues

Will you use personal data? If so, shortly describe the kind of personal data you will use. Add the reference to your file in KU Leuven's Register of Data Processing for Research and Public Service Purposes (PRET application). Be aware that registering the fact that you process personal data is a legal obligation.

- No

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)

- No

Does your work possibly result in research data with potential for tech transfer and valorisation? Will IP restrictions be claimed for the data you created? If so, for what data and which restrictions will be asserted?

- Yes

Ownership of the generated data belongs to KU Leuven & VIB in accordance with the framework agreement between these institutes; copyright of the data belongs to Kevin Verstrepen. It should be noted that part of the research is carried out in collaboration with TIB, more information on the agreement between KU Leuven, VIB and TIB can be found in the answer to the next question.

We identify in an early phase the valorization potential of research lines and our group has a vast network of industrial contacts to efficiently start the route to commercialization. The lab is supported in this matter by Dr. Stijn Spaepen, IOF innovation manager responsible for research valorization. For research with valorization potential, the host lab actively protects its IP by filling patent applications with support from the IP department of VIB. Type of data with potential for tech transfer and valorization include yeast strains generated during the timeframe of this project (including phenotypic data), sequencing information generated during the timeframe of this project and (analysis) data and models hereof derived. Valorization potential includes a) licensing of (improved) strains or information on linking a specific sequence variant to a phenotype and b) creation or participation in start-up companies.

Do existing 3rd 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?

- Yes

Existing agreements between VIB and KU Leuven do not restrict publication of data. There is no IP on the generated strains that would prevent us from storing the strains, performing the anticipated experiments or publishing the results.

It should be noted that part of the research is carried out in collaboration with TIB in China. Hence, there is a collaboration agreement in place between VIB, KU Leuven and TIB (the Chinese partner in this FWO project). This agreement stipulates the following:

"Results are owned by the Party that generated them. Each Party is entitled to valorise its own Result at its own discretion. The Parties shall not exploit or valorise jointly owned Results without first putting in place a joint ownership agreement in which Parties shall specify in which Territory and/or Field each Party take the lead for commercialization and how any remuneration of the exploitation of the Jointly Owned Results shall be shared with the other Party."

For all strains used in this project, we will check, together with René Custers (Regulatory & Responsible Research Manager at VIB) whether we need to comply with the Nagoya protocol on Access to Genetic Resources and the Fair and Equitable Sharing of Benefits Arising from their Utilization. In case we need to comply, we will meet certain diligence obligations under EU Regulation (EU) No 511/2014 (apply for 'Prior Informed Consent' document and negotiate 'Mutually Agreed Terms' with Competent National Authority of the provider country).

4. Documentation and metadata

What documentation will be provided to enable reuse of the data collected/generated in this project?

Data will be generated following standardized protocols. Data will be documented by the research and technical staff at the time of data collection and analysis, by taking careful notes in their (electronic) laboratory notebook. Cryotubes of biological samples (bacterial and yeast strains) stored at -80°C will be labelled with a reference number that links to an entry in our strain database.

Digital data files will be stored on KU Leuven servers and will be accompanied with a READ ME text file that contains relevant metadata for understanding and re-use of data.

The READ ME file will be structured as follows:

PROJECT INFORMATION

P01. Project name/nickname, including name of funder, type of grant, grant number

P02. Name and contact of the Principal Investigator and scientists (the whole team) involved in the project with indication of the lead scientist

- P03. Small description of the project
P04. (Archive) Dropbox folder link associated with this project if any

GENERAL INFORMATION

- G01. Names of file(s) or dataset(s) that this README file describes
G02. Date of creation/last update of the README file.
For each update/change made to the files, an extra line is added with date and main changes made.
G03. Created by (name)
G04. Description of the dataset

FILE OVERVIEW

- F01. Software used to generate the data, including software version used
F02. Software necessary to open the file
F03. Relationship between the files
For example, for scripts used, a list of files generated by each script will be included.

METHODOLOGICAL INFORMATION

- M01. Date (beginning-end) and place of data collection
M02. Data collecting method
M03. Information about data processing methods and scripts used.
M04. Information about the instrument, calibration, settings used
M05. Information about limitations of the dataset (missing values, ...), information that ensures correct interpretation of the dataset
M06. People involved in the creation or processing of the dataset

DATA ACCESS AND SHARING

- A01. Confidentiality information/restrictions on use of data

DATA SPECIFIC INFORMATION (ABOUT THE DATA THEMSELVES)

- D01. Full names and definitions for columns and rows.
This can be provided in a csv file.
D02. Explanation of abbreviations
D03. Explanation of strain codes and sample names
Crosslink to lab strainlist if possible.

RELATIONSHIPS

- R01. Publications based on this dataset
R02. This dataset derives from... (other dataset)
R03. This dataset is related to... (documents, dataset)

Cryotubes and plates with biological strains are labeled with a reference number that links to an entry in our strain database (stored on both a professional Dropbox account as well as KU Leuven servers, see also below). All relevant information on the specific strains (strain ID, genetic information, origin of strain) is included in this database.

Will a metadata standard be used? If so, describe in detail which standard will be used. If no, state in detail which metadata will be created to make the data easy/easier to find and reuse.

- Yes

Since there is no formally acknowledged metadata standard specific to our discipline, Dublin Core Metadata will be used. Moreover, we will closely monitor MIBBI (Minimum Information for Biological and Biomedical Investigations) for metadata standards that are more specific to our data.

5. Data storage and backup during the FWO project

Where will the data be stored?

Biological samples: Cryotubes with strains will be stored in a -80° freezer present in the Verstrepen lab at KU Leuven (Gaston Geenslaan 1, Heverlee). Costs are covered by general lab expenses. Unauthorized people do not have access to strains.

Data (digital files) generated in this project will be stored in a Dropbox Business Advanced account for processing and analyses; following secure data transfer, modern data encryption standards, and encrypted block storage (256-bit AES and SSL/TLS encryption. For more details see: <https://www.dropbox.com/business/trust>

Sequencing data will be stored on an internal lab server (present in host lab) as well as on a secure Dropbox Business account for processing and analyses.

All the relevant algorithms, scripts and software code driving the project will be stored on a secure Dropbox account. Scripts used for analysis will also be stored in Jupyter notebook (jupyter.org - an open source web application to store and share scripts), in github or in the GitLab service of KU Leuven.

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), the Protein Database (for protein sequences).

How is backup of the data provided?

Strains are backed up in a compressed form (96-well plates) in the Verstrepen lab and on a 2nd location (Kasteelpark Arenberg 20, Heverlee) at KU Leuven. Data (digital files) are automatically backed up by the secure Dropbox Business Advanced account cloud backup services. Additionally, project data and sequencing data will be backed up to KU Leuven servers, with automatic back-up procedures.

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

We estimate the total data generated during this project to amount to around 30 Gb.

Dropbox Business offers unlimited storage and back-up capacity in their clouds.

There is sufficient storage and back-up capacity on all KU Leuven servers:

- the "L-drive" is an easily scalable system, built from General Parallel File System (GPFS) cluster with NetApp eseries storage systems, and a CTDB samba cluster in the front-end.
- the "J-drive" is based on a cluster of NetApp FAS8040 controllers with an Ontap 9.1P9 operating system.
- KU Leuven also offers a 2TB onedrive already included in our group's office 365 subscription

What are the expected costs for data storage and back up during the project? How will these costs be covered?

The total estimated cost of data storage during the project is 4000 €. This estimation is based on the following costs:

- Yeast/bacteria strains are easily kept alive for several weeks. This costs on average 5 euro. When no experiments are planned with a specific strain, cryopreservation will thus be used to safeguard strains, prevent genetic drift, loss of transgene and potential contaminations. - 80°C freezers are present in the lab of prof. Verstrepen and costs are included in general lab costs.
- The costs associated with a Dropbox Business account has been negotiated by the lab to 10 USD/month/user.
- The costs of digital data storage are as follows: 173,78€/TB/Year for the "L-drive" and 519€/TB/Year for the "J-drive".

Additional funding for storage has not been requested, therefore the costs will be covered by the lab's internal budget.

Data security: how will you ensure that the data are securely stored and not accessed or modified by unauthorized persons?

- For biological samples: unauthorized people do not have access to the strain collections.
- For the digital data: Access to data stored on the Dropbox Business Advanced cloud is granted based on role based access control and all access requires layers of authentication that includes strong passwords, SSH keys, 2 factor authentication, and one time passcodes. Dropbox safeguards data with document watermarking, granular content permissions and policies, document watermarking, and legal holds.
- Both the "L-drive" and "J-drive" KU Leuven 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.

6. Data preservation after the FWO project

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. In addition, in line with KU Leuven regulation, we will preserve the data for at least 5 additional years.

Where will the data be archived (= stored for the longer term)?

As a general rule, datasets will be made openly accessible, whenever possible via existing platforms that support FAIR data sharing (www.fairsharing.org), at the latest at the time of publication.

For all other datasets, long term storage will be ensured as follows:

- Biological data: yeast and bacterial strains will be stored locally in the laboratory (-80°C). Other biological and chemical samples: storage at 4°C and/or as frozen samples as appropriate
- Digital datasets: files will be stored on Dropbox account and the KU Leuven "L-drive".

What are the expected costs for data preservation during the retention period of 5 years? How will the costs be covered?

The total estimated cost of data storage during the 10 years after the end of the project is estimated to be around 8000 €. These costs will be covered from the general lab budget.

7. 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 3rd party, legal restrictions)?

- Yes. Specify:

We aim at communicating our results in top journals that require full disclosure of all included data. Biological material will be shared upon simple request following publication, unless we identify valuable IP, in which case we will first protect commercial exploitation, either through patenting or via an MTA that restricts the material from commercial use.

Which data will be made available after the end of the project?

Participants to the present project are 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, unless we identify valuable IP, in which case we will first protect commercial exploitation, either through patenting or via an MTA that restricts the material from commercial

use.

Where/how will the data be made available for reuse?

- In an Open Access repository
- Other (specify):

We aim at communicating our results in top journals that require full disclosure upon publication of all included data, either in the main text, in supplementary material or in a data repository if requested by the journal and following deposit advice given by the journal. Depending on the journal, accessibility restrictions may apply. Proper links to datasets will be provided in the corresponding publications.

As a general rule, datasets will be made openly accessible via existing platforms that support FAIR data sharing (www.fairsharing.org). Sharing policies for specific research outputs are detailed below:

- Biological data: Bacteria and yeast strains will be shared upon simple request following publication, unless we identify valuable IP. In this case, we will first protect commercial exploitation, either through patenting or via an MTA that restricts the material from commercial use.
- Datasets will be deposited in open access repositories.
- 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 deposited in lab notebook are accessible to the PI and the research staff involved in the project, and will be made available upon request.
- Manuscripts: We opt for open access publications where possible. Publications will be automatically listed in our institutional repository, Lirias 2.0, based on the authors name and ORCID ID.
- Algorithms and scripts: As soon as a manuscript is publicly available, algorithms and scripts will be deposited in a github repository.
- Nucleic acid and protein sequences: 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), the Protein Database (for protein sequences).

When will the data be made available?

- Upon publication of the research results

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.

Who will be able to access the data and under what conditions?

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 conform 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, a Creative Commons Attribution (CC-BY) or an ODC Public Domain Dedication and License, 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.

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.

What are the expected costs for data sharing? How will the costs be covered?

It is the intention to minimize data management costs by implementing standard procedures 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. A budget for publication costs has been requested in this project.

8. Responsibilities

Who will be responsible for data documentation & metadata?

Metadata will be documented by the research and technical staff involved in the project at the time of data collection and analysis, by taking careful notes in their laboratory notebook that refer to specific datasets and by ensuring detailed read me files and other documentation and metadata. Ms. Veerle Saels (senior lab technician) and Ms. Eef Lemmens (personal assistant to prof. Verstrepen) will follow this up.

Who will be responsible for data storage & back up during the project?

The research and technical staff involved in the project will ensure data storage and back up. Mr. Nico Vangoethem (IT responsible) will follow this up.

Who will be responsible for ensuring data preservation and reuse ?

The PI (Kevin Verstrepen) is responsible for data preservation and sharing, with support from the research and technical staff involved in the project. He will also be assisted by Mr. Nico Vangoethem and Ms. Eef Lemmens in these aspects.

Who bears the end responsibility for updating & implementing this DMP?

The PI (Kevin Verstrepen) bears the end responsibility of updating & implementing this DMP.