



HAL
open science

A community resource for paired genomic and metabolomic data mining

Michelle A. Schorn, Stefan Verhoeven, Lars Ridder, Florian Huber, Deepa D. Acharya, Alexander A. Aksenov, Gajender Aleti, Jamshid Amiri Moghaddam, Allegra T. Aron, Saefuddin Aziz, et al.

► To cite this version:

Michelle A. Schorn, Stefan Verhoeven, Lars Ridder, Florian Huber, Deepa D. Acharya, et al.. A community resource for paired genomic and metabolomic data mining. *Nature Chemical Biology*, 2021, 10.1038/s41589-020-00724-z . pasteur-03153704

HAL Id: pasteur-03153704

<https://pasteur.hal.science/pasteur-03153704>

Submitted on 26 Feb 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

OPEN

A community resource for paired genomic and metabolomic data mining

Genomics and metabolomics are widely used to explore specialized metabolite diversity. The Paired Omics Data Platform is a community initiative to systematically document links between metabolome and (meta)genome data, aiding identification of natural product biosynthetic origins and metabolite structures.

Michelle A. Schorn, Stefan Verhoeven, Lars Ridder, Florian Huber, Deepa D. Acharya, Alexander A. Aksenov, Gajender Aleti, Jamshid Amiri Moghaddam, Allegra T. Aron, Saefuddin Aziz, Anelize Bauermeister, Katherine D. Bauman, Martin Baunach, Christine Beemelmans, J. Michael Beman, María Victoria Berlanga-Clavero, Alex A. Blacutt, Helge B. Bode, Anne Boullie, Asker Brejnrod, Tim S. Bugni, Alexandra Calteau, Liu Cao, Víctor J. Carrión, Raquel Castelo-Branco, Shaurya Chanana, Alexander B. Chase, Marc G. Chevrette, Leticia V. Costa-Lotufó, Jason M. Crawford, Cameron R. Currie, Bart Cuypers, Tam Dang, Tristan de Rond, Alyssa M. Demko, Elke Dittmann, Chao Du, Christopher Drozd, Jean-Claude Dujardin, Rachel J. Dutton, Anna Edlund, David P. Fewer, Neha Garg, Julia M. Gauglitz, Emily C. Gentry, Lena Gerwick, Evgenia Glukhov, Harald Gross, Muriel Gugger, Dulce G. Guillén Matus, Eric J. N. Helfrich, Benjamin-Florian Hempel, Jae-Seoun Hur, Marianna Iorio, Paul R. Jensen, Kyo Bin Kang, Leonard Kaysser, Neil L. Kelleher, Chung Sub Kim, Ki Hyun Kim, Irina Koester, Gabriele M. König, Tiago Leao, Seoung Rak Lee, Yi-Yuan Lee, Xuanji Li, Jessica C. Little, Katherine N. Maloney, Daniel Männle, Christian Martin H., Andrew C. McAvoy, William W. Metcalf, Hosein Mohimani, Carlos Molina-Santiago, Bradley S. Moore, Michael W. Mullowney, Mitchell Muskat, Louis-Félix Nothias, Ellis C. O'Neill, Elizabeth I. Parkinson, Daniel Petras, Jörn Piel, Emily C. Pierce, Karine Pires, Raphael Reher, Diego Romero, M. Caroline Roper, Michael Rust, Hamada Saad, Carmen Saenz, Laura M. Sanchez, Søren Johannes Sørensen, Margherita Sosio, Roderich D. Süssmuth, Douglas Sweeney, Kapil Tahlan, Regan J. Thomson, Nicholas J. Tobias, Amaro E. Trindade-Silva, Gilles P. van Wezel, Mingxun Wang, Kelly C. Weldon, Fan Zhang, Nadine Ziemert, Katherine R. Duncan, Max Crüsemann, Simon Rogers, Pieter C. Dorrestein, Marnix H. Medema and Justin J. J. van der Hooft

Interactions between bacteria, fungi, plants, and animals, as well as their environments are often facilitated through specialized metabolites, also known as natural products. These specialized metabolites are molecules naturally produced by organisms that are not strictly required for survival but may confer an advantage to the producing organism, such as the inhibition of nearby species competing for nutritional resources. The chemical structures and functions, as well as the biosynthetic origins of such metabolites, are largely hidden, especially in complex environments. To understand and harness these chemical interactions, it is crucial to study their genetic and structural bases. However, the confident recognition, dereplication, and prioritization of specialized metabolites in complex mixtures remains very challenging. While individual efforts to interpret the chemical and genetic languages have been largely successful

in connecting genes and molecules^{1,2}, large-scale correlations leveraging complementary chemical and genomic data have yet to be realized.

The research community has generated a wealth of genomic and metabolomic data, which has been deposited in dedicated repositories, and tools for mining these data separately are being developed rapidly. Platforms such as the antibiotics and Secondary Metabolite Analysis Shell (antiSMASH)³ and PRediction Informatics for Secondary Metabolomes (PRISM)⁴ use genomic information to annotate biosynthetic gene clusters (BGCs), a set of genes that encode the producing framework for metabolites of diverse chemical classes, such as polyketides, peptides and terpenoids. The antiSMASH database and the Joint Genome Institute's (JGI's) Integrated Microbial Genomes and Microbiomes (IMG/M)/Atlas of Biosynthetic Gene Clusters (IMG/ABC) database⁵ contain tens

of thousands of BGCs identified in publically available genomes, while the Minimum Information about a Biosynthetic Gene cluster (MIBiG)⁶ database connects over 2,000 BGCs to the specialized metabolites for which they encode the biosynthetic pathways. On the metabolomics side, mass spectrometry (MS) has become the most commonly used technique for performing high-throughput measurements³. Data repositories and analysis platforms such as the Global Natural Product Social Molecular Networking-Mass Spectrometry Interactive Virtual Environment (GNPS-MassIVE)⁷, MetaboLights⁸, and the Metabolomics Workbench⁹ facilitate the sharing, processing, and analysis of MS data. These platforms, along with spectral libraries², such as the GNPS spectral library, METLIN, MassBank, and the commercially available NIST library, provide resources for reference mass spectra of a wide range of chemical structures, thereby aiding

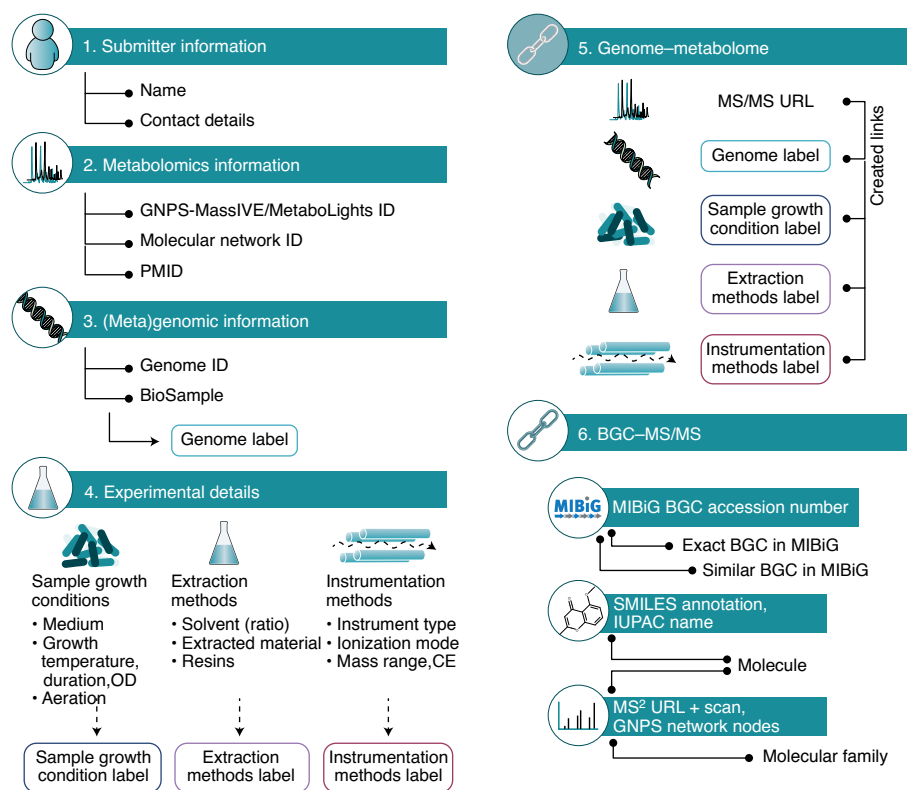


Fig. 1 | Overview of the Paired Omics Data Platform. The PoDP links genomic and metabolomic data deposited in public repositories, accompanied by minimal metadata. The platform documents basic submitter information as well as accession numbers of the metabolome and genome data. Standardized logging of key experimental details enables users to better search and compare datasets, and user-defined labels of genome and experimental details (sections 3 and 4) enable straightforward submission of multiple links. The core of the platform consists of links between the genomic and metabolomic datasets and links between BGCs and MS/MS spectra, which facilitate data integration. PMID, PubMed ID; OD, optical density; CE, collision energy; SMILES, simplified molecular-input line-entry; IUPAC, International Union of Pure and Applied Chemistry; MS², MS/MS fragmentation.

metabolite annotation. Together, these resources provide the basis for sharing and reusing genomic and metabolomic data and structural annotations and have spurred the development of numerous algorithms for mining these information-dense data.

Several studies and tools have started to explore the combination of genomic and metabolomic data to enhance metabolite annotation, dereplication, and prioritization workflows. While MS-based metabolomics provides increasing amounts of information related to the metabolite structures present in complex mixtures, it faces inherent limitations with respect to structural identification. To address this, several tools, such as GNPS-based molecular networking⁷ and mass spectrometry to latent dirichlet allocation (MS2LDA) substructure discovery¹⁰, have been proposed that computationally exploit tandem mass spectrometry (MS/MS) fragmentation spectra to map relationships

between metabolites in networks and identify (shared) substructures, thereby facilitating metabolite annotation. Genomics has also been used to provide complementary structural information through the biotransformations encoded in biosynthetic machinery¹, as well as a way to link specialized metabolites to their producers via BGCs that are mined from genome sequences from known organisms. Integrative strategies have been described for bacterial¹¹, fungal¹², and plant¹³ specialized metabolites. A series of tools and approaches, mostly targeting biosynthetically modular natural products such as peptides and glycosides, have been introduced over the last decade to integrate genome and metabolome data, such as peptidogenomics¹¹, MetaMiner¹⁴, GRAPE-GARLIC¹⁵ and metabologenomics¹⁶. These tools show the potential of combined omics approaches to accelerate natural product discovery.

It has become standard procedure to deposit genomic information to public databases, such as the National Center for Biotechnology Information's (NCBI's) GenBank¹⁷ or JGI's IMG/M⁵, and it is becoming increasingly common to submit mass spectrometry data to repositories such as GNPS-MassIVE⁷, MetaboLights⁸ or Metabolomics workbench⁹. However, there is currently no straightforward way to connect different types of omics data that are derived from the same biological source. It often takes extensive literature review to determine which omics data belong to the same species, organism, or sample, and therefore constitute 'paired' datasets, making reuse of these data challenging and time consuming. Additionally, there is no straightforward way to obtain consistent metadata for such links. To facilitate large-scale, effective integration of these data, it is vital to have a community-driven online resource that stores annotated links between paired datasets. Here, we refer to paired data as genomic data (specifically a genome or metagenome assembly) and metabolomic data (specifically MS/MS data) that originate from the same source. So far, no such platform supporting natural product discovery has been available.

The value of integrating different data types and organizing sample metadata is increasingly recognized by the scientific community. For example, the BioStudies¹⁸ and BioSample¹⁹ databases facilitate the capture and organization of various omics data types and sample information. In particular, the BioStudies database supports linkage between genomics and metabolomics studies; however, links between genome-mining resources, such as MIBiG, and natural product metabolomics platforms, such as GNPS-MassIVE, are currently not documented in this database.

Here we introduce the Paired Omics Data Platform (PoDP) to streamline access to paired omics data so that both humans and computers can access and read paired datasets and can also record and exploit validated links between BGCs and metabolites (<https://pairedomicsdata.bioinformatics.nl/>). In addition to linking these omics data types, the platform stores essential metadata (i.e., growth media, extraction solvent, and ionization mode) using existing ontology where available, thus facilitating reuse of for-the-user relevant sections of paired data. This platform will boost the successful integration of unsupervised data-mining strategies to fine-tune the structural annotation of modular natural product classes and include yet-unknown classes of natural products. This will aid in structural and functional

Box 1 | Preliminary submissions to the PoDP

Early contributors to the PoDP seeded the platform with 70 projects from over 45 labs in 10 countries. These 70 projects encompass:



- 1,268 genomes
- 1,306 metagenomes
- 42 metagenome-assembled genomes



- 4,853 paired (meta)genomes and metabolomes
- 114 validated links between BGCs and MS/MS spectra



- 155 sample growth conditions
- 100 extraction methods
- 75 instrumentation methods

Because all data linked in the platform must already be in a public database, many early contributors made data public that had previously not been public. Some early contributors went a step further and actually acquired genomic or metabolomic data to complement already public data and make paired datasets. Submitting genomic and metabolomic data to the PoDP will increase visibility of those data and allow researchers to adhere to FAIR data principles.

annotations of natural products and the genes responsible for their production, and we anticipate that this will help uncover the potential producers of molecules in nature. Finally, registering these links in a standardized way gives the community an invaluable resource of Findable, Accessible, Interoperable, and Reusable (FAIR)²⁰ data.

Standards for paired data

The aim of the PoDP is to connect public metabolomics datasets to their genomic origins. The PoDP does not store any metabolomics or genomics datasets, but captures metadata defining pairs of omics datasets in existing public databases and platforms already validated and utilized by the genomics and metabolomics communities. The PoDP consists of a six-section form for easy and quick input of data (Fig. 1). The metadata is organized in projects that can consist of multiple related experiments, identified by their MassIVE accession or MetaboLights study identifier. The (meta)genomes(s) used in these experiments can all be added to the same project via a public database identifier (e.g., a NCBI GenBank accession number or JGI Genome ID), with the user creating easy-to-recall genome labels for each (meta)genome. Minimal metadata with information about sample preparation and data collection are recorded in a modular way, allowing for multiple experimental

set-ups within one project. Furthermore, through BioSample accession IDs, metadata stored elsewhere can be linked to (meta) genome(s) as well. User-specified metadata labels are also used for easy recall in the linking step, in which a URL for a specific set of MS spectra is linked with the genome label and metadata labels to create a genome–metabolome link. To create a BGC–MS/MS link, a MIBiG identifier for the same or similar BGC can be linked with a MS/MS URL and scan number of a single measured molecule or molecular network nodes (representing unique measured molecules) in a molecular family (a group of structurally related molecules identified by similar fragmentation patterns). This approach thus stimulates the submission of validated gene clusters to the MIBiG repository in order to make a BGC–MS/MS link in the PoDP.

By obtaining iterative feedback from a group of early users from various research groups, we narrowed down the required metadata in the PoDP to the minimum information needed to make meaningful links between genomic and metabolomic data relevant to the community. Capturing the full range of relevant variables in any given experiment in a standardized and machine-readable format would lead to a very complex and tedious data entry process. Therefore, a balance was struck between flexible and user-friendly data

entry, maintaining machine readability for future large-scale analyses. By standardizing and connecting to ontologies only the most relevant information that could substantially affect the metabolites produced, extracted, and detected by MS, we arrived at a set of minimal metadata required for submission.

To enable machine readability of the data, ontologies are used to standardize response options wherever possible. This ensures that a global community can use the same term for a given piece of metadata and use these ontologies to make accurate and meaningful selections of data to analyze. For example, researchers can reliably select and obtain only datasets that use tryptic soy broth for culture or only metagenomic datasets derived from aquatic invertebrates, or just the fraction of paired datasets in which the MS data was obtained in positive ionization mode. For metadata categories with numerous options, in which all possibilities cannot be captured by standard ontologies, an “Other” category is provided for further explanation. Free text entered in the “Other” boxes is inherently not machine-readable but gives an option for customization by the user and can help to keep important but non-standardized records of the paired data. Furthermore, all fields including the “Other” boxes can be searched to find projects containing specific data.

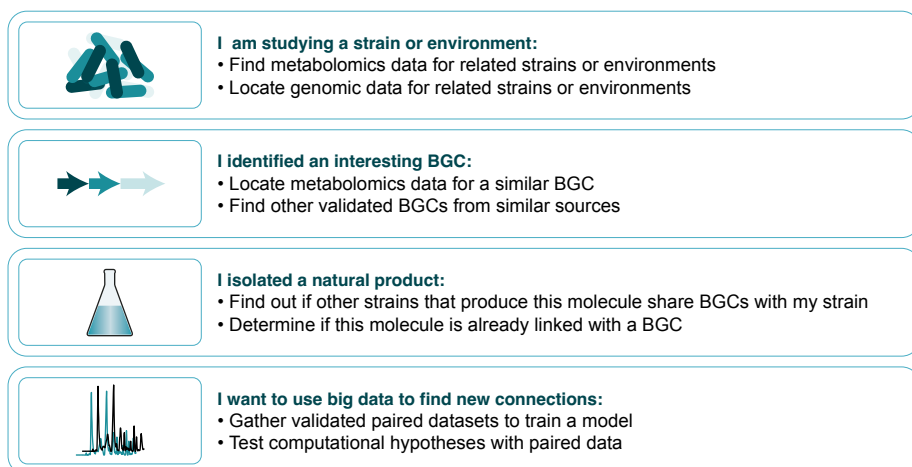


Fig. 2 | Example use cases of the Paired Omics Data Platform. Users may approach the PoDP using genomic or metabolomic data (or using metadata) and exploit the links provided to generate new hypotheses about their primary data. Specifically, genomic data may enable new hypotheses about the structures or biosynthetic pathways for an identified molecule or mass feature, while metabolomic data may provide new hypotheses regarding the product(s) of a BGC. Integrative computational approaches allow scaling these analyses to systematic and comprehensive efforts.

Preliminary dataset statistics

An initial call to deposit paired datasets in the PoDP was met with enthusiasm from the research community. Over 45 laboratories from 10 countries have contributed 70 paired datasets. Those 70 projects (Box 1) contain 4,853 MS samples associated with sequenced source material. Of the more than 2,600 different genomic sources deposited, 1,306 are metagenomes, 1,268 are genomes, and 42 are metagenome-assembled genomes. The integrative collection of over 4,800 genome–metabolome links is accompanied by metadata: 155 sample preparation methods, 100 extraction methods, and 75 instrumentation methods. Furthermore, 114 links between BGCs and their associated MS/MS spectra are registered in the platform. These community-curated data are regularly archived to a Zenodo dataset and made available for download in JSON format.

The PoDP encourages adherence to FAIR principles²⁰, requiring data to already be deposited in databases and made publicly available before being entered in the PoDP. Presence of a project in the PoDP will increase the findability of those data, results, and publications, while allowing researchers to perform new analyses on existing publicly available data without the need to generate new data. As part of this community effort, a number of projects deposited in the PoDP made their data publicly available to allow submission into the platform; thus far, over 680 metabolomics samples and over 70

genomic sources, including five BGCs newly uploaded to MIBiG, were made public. For example, the PoDP stimulated the upload of metabolomics data to MassIVE for a collection of 120 sequenced *Streptomyces* strains for which genomics data was previously published²¹. In another example, 20 metagenomes from marine sediments were made public for the platform. Additionally, some datasets were acquired and made publicly available expressly for deposition into the PoDP. In one case, a research group with 44 already sequenced cyanobacterial strains²² was inspired to acquire metabolomics data for each strain so that the paired data could be uploaded to the PoDP.

To better view the data encompassed by the PoDP, users can search for projects under the “List” tab, using keywords to find studies of interest. For example, to find paired data resulting from a *Streptomyces* or *Salinispora* species, searching for the genera (“*Streptomyces* | *Salinispora*”) will result in the projects (currently 18) that measured *Streptomyces* or *Salinispora* strains. Likewise, to compare projects that used methanol to extract cell pellets, searching “methanol + cells” retrieves projects that used methanol to extract cell pellets. To obtain more detail on the metadata contained in each project, users can navigate to the project page by clicking on the project identifier. There, users can find details about the genome or metagenome when clicking on the label, which will then provide a link to the publically available genomic data.

Likewise, the publically available MS data can be downloaded directly from the link provided. Clicking on the Sample Growth, Extraction, and Instrumentation Methods labels will display the corresponding metadata.

Applications of the platform

The PoDP can be used in both basic and advanced ways. In a basic way, researchers from across disciplines can apply linked data for numerous applications (Fig. 2). With linked data, we refer to a BGC that can be experimentally linked to a MS/MS spectrum or a molecular family. For example, a natural product chemist who isolates a molecule from a cyanobacterium can use the PoDP to find mass spectra from genetically similar cyanobacteria for comparative metabolomics analyses. A biologist who has identified a BGC of interest and has MS data for the producing strain can download data for the products of similar BGCs and their products to determine whether the BGC is novel and/or to guide molecule isolation. Scientists from all fields can find reliable paired data for use in their own research while also contributing their data for future community use. The importance of consistent metadata cannot be underestimated, and we welcome the development of curated resources such as the Natural Products Atlas²³ that aim to create coherent records for microbial natural products. Combined with the PoDP, this gives researchers complementary resources to mine for natural product structures, their producers, and available omics data.

Furthermore, more advanced applications are possible utilizing large-scale computational approaches (Fig. 2). Several algorithmic strategies to link genomics and metabolomics data to chart specialized metabolic diversity have been suggested, including correlation- and feature-based matching². Both types of linking benefit from systematically curated datasets of related organisms with BGCs and metabolites occurring in various samples. With the PoDP in place now, these strategies can be used more effectively to select appropriate datasets to start mining for novel links. Moreover, algorithms to score and rank links between BGCs and metabolites are easier to develop and benchmark: for example, a new set of scores was recently proposed using a number of PoDP datasets with validated BGC–metabolite links to demonstrate the effect of the novel scoring system within the newly introduced NPLinker framework²⁴.

Moving forward with FAIR data

The amount of preliminary data deposited and the enthusiasm from the community for the PoDP reaffirm the need for such a repository of paired public datasets. Feedback from early users also indicated an eagerness to include additional kinds of data in the future. Presently, the PoDP is expressly for linking MS/MS data and whole-genome or metagenome data. Potentially, the PoDP could be developed to include other types of spectral data, like full scan (MS¹) metabolomics mass spectrometry data and NMR, as well as proteomics data. Additionally, different kinds of genomic data could be facilitated, including 16S rRNA or other amplicon sequences, transcriptomic data, and genetic manipulation or heterologous expression data. Such additions will further fuel integrated omics analysis tools and approaches, a field that has gained much traction recently²⁵.

The PoDP requires researchers to deposit their data in public databases, stimulating the upload of data by early users, which is exemplified by more than 1,800 GNPS-MassIVE and MetaboLights submissions just prior to submitting these data in the PoDP. As a FAIR data platform, the PoDP not only facilitates reuse of data, but also promotes the work of researchers who submit their data to the PoDP, through increased publication visibility. Future efforts to (re)use these data by connecting to other platforms and programs for analyzing paired data, such as NPLinker²⁴, will further the field of natural product prediction and discovery.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

Each project can be downloaded from the website individually as a JSON file. The (meta)genome and metabolome datasets can be found in their public repositories. All PoDP projects are archived monthly to Zenodo at <https://doi.org/10.5281/zenodo.3736430>.

Code availability

The software is licensed under the Apache 2.0 open source license and the source code can be found on GitHub (<https://github.com/iomega/paired-data-form>), which includes the dependencies of the software. Each software release is archived to Zenodo at <https://doi.org/10.5281/>

[zenodo.2656630](https://doi.org/10.5281/zenodo.2656630). A full description of how the platform was built can be found on <https://pairedomicsdata.bioinformatics.nl/methods>. □

Michelle A. Schorn^{1,2,7,3}, Stefan Verhoeven^{1,3,7,3}, Lars Ridder³, Florian Huber^{1,3}, Deepa D. Acharya⁴, Alexander A. Aksenov^{1,5}, Gajender Aleti⁶, Jamshid Amiri Moghaddam^{1,7}, Allegra T. Aron⁵, Saefuddin Aziz^{8,9}, Anelize Bauermeister^{5,10}, Katherine D. Bauman¹¹, Martin Baunach¹², Christine Beemlanna⁷, J. Michael Beman^{13,14}, Maria Victoria Berlanga-Clavero^{1,15}, Alex A. Blacutt¹⁶, Helge B. Bode^{17,18,19,20}, Anne Boullie²¹, Asker Brejnrod⁵, Tim S. Bugni^{1,22}, Alexandra Calteau²³, Liu Cao^{1,24}, Víctor J. Carrión^{25,26}, Raquel Castelo-Branco^{27,28,29}, Shaurya Chanana^{1,4}, Alexander B. Chase^{1,11}, Marc G. Chevrette⁴, Leticia V. Costa-Lotufo^{1,10}, Jason M. Crawford^{1,30,31,32}, Cameron R. Currie^{1,33,34}, Bart Cuypers^{35,36}, Tam Dang³⁷, Tristan de Rond¹¹, Alyssa M. Demko^{1,11}, Elke Dittmann^{1,12}, Chao Du^{1,25}, Christopher Drozd¹⁶, Jean-Claude Dujardin³⁶, Rachel J. Dutton^{38,39}, Anna Edlund^{40,41}, David P. Fewer^{1,29}, Neha Garg⁴², Julia M. Gauglitz⁵, Emily C. Gentry⁵, Lena Gerwick¹¹, Evgenia Glukhov¹¹, Harald Gross⁸, Muriel Gugger²¹, Dulce G. Guillén Matus^{1,11}, Eric J. N. Helfrich^{1,17,19,43,44}, Benjamin-Florian Hempel^{37,45}, Jae-Seoun Hur⁴⁶, Marianna Iorio⁴⁷, Paul R. Jensen¹¹, Kyo Bin Kang^{1,48}, Leonard Kayser^{8,49}, Neil L. Kelleher^{1,50}, Chung Sub Kim^{1,30,31,51}, Ki Hyun Kim^{1,51}, Irina Koester⁵², Gabriele M. König⁵³, Tiago Leao^{5,11}, Seoung Rak Lee^{51,54}, Yi-Yuan Lee^{1,24}, Xuanji Li⁵⁵, Jessica C. Little⁵⁶, Katherine N. Maloney^{1,57}, Daniel Männle^{1,8,49,58}, Christian Martin H.⁵⁹, Andrew C. McAvoy^{1,42}, William W. Metcalf^{1,60}, Hosein Mohimani²⁴, Carlos Molina-Santiago¹⁵, Bradley S. Moore^{1,11,39,61}, Michael W. Mullaney^{1,50}, Mitchell Muskat^{1,11}, Louis-Félix Nothias^{1,62}, Ellis C. O'Neill⁶², Elizabeth I. Parkinson^{1,63}, Daniel Petras^{1,5,52}, Jörn Piel^{1,43}, Emily C. Pierce^{1,38}, Karine Pires^{1,64}, Raphael Reher^{1,11}, Diego Romero^{1,15}, M. Caroline Roper¹⁶, Michael Rust^{1,43}, Hamada Saad^{8,65}, Carmen Saenz^{1,66}, Laura M. Sanchez^{1,56}, Søren Johannes Sørensen^{1,55}, Margherita Sosio⁴⁷, Roderich D. Süssmuth³⁷, Douglas Sweeney¹¹, Kapil Tahlan⁶⁷, Regan J. Thomson^{1,50}, Nicholas J. Tobias^{19,68}, Amaro E. Trindade-Silva⁶⁹, Gilles P. van Wezel^{1,25}, Mingxun Wang⁵, Kelly C. Weldon^{5,39}, Fan Zhang²², Nadine Ziemert^{49,58}, Katherine R. Duncan^{1,70},

Max Crüsemann^{1,53}, Simon Rogers^{1,71}, Pieter C. Dorrestein^{1,5,11,39,72}, Marnix H. Medema^{1,2} and Justin J. J. van der Hooft^{1,2}

¹Laboratory of Microbiology, Department of Agricultural and Food Sciences, Wageningen University, Wageningen, the Netherlands. ²Bioinformatics Group, Wageningen University, Wageningen, the Netherlands. ³Netherlands eScience Center, Amsterdam, the Netherlands. ⁴Wisconsin Institute for Discovery and Department of Plant Pathology, University of Wisconsin-Madison, Madison, WI, USA. ⁵Collaborative Mass Spectrometry Innovation Center, Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, La Jolla, CA, USA. ⁶Department of Psychiatry, University of California San Diego, San Diego, CA, USA. ⁷Leibniz Institute for Natural Product Research and Infection Biology e.V. Hans-Knöll-Institute (HKI), Jena, Germany. ⁸Pharmaceutical Biology Department, Pharmaceutical Institute, Eberhard Karls University Tübingen, Tübingen, Germany. ⁹Microbiology Department, Biology Faculty, Jenderal Soedirman University, Purwokerto, Indonesia. ¹⁰Instituto de Ciências Biomédicas, Universidade de São Paulo, São Paulo, Brazil. ¹¹Center for Marine Biotechnology and Biomedicine, Scripps Institution of Oceanography, University of California San Diego, La Jolla, CA, USA. ¹²University of Potsdam, Institute of Biochemistry and Biology, Potsdam-Golm, Germany. ¹³Department of Life and Environmental Sciences, University of California Merced, Merced, CA, USA. ¹⁴Sierra Nevada Research Institute, University of California Merced, Merced, CA, USA. ¹⁵Instituto de Hortofruticultura Subtropical y Mediterránea "La Mayora", Universidad de Málaga-Consejo Superior de Investigaciones Científicas, Departamento de Microbiología, Universidad de Málaga, Málaga, Spain. ¹⁶Department of Microbiology and Plant Pathology, University of California Riverside, Riverside, CA, USA. ¹⁷Molecular Biotechnology, Department of Biosciences, Goethe University Frankfurt, Frankfurt am Main, Germany. ¹⁸Buchmann Institute for Molecular Life Sciences, Goethe University Frankfurt, Frankfurt am Main, Germany. ¹⁹Senckenberg Gesellschaft für Naturforschung, Frankfurt am Main, Germany. ²⁰Max-Planck-Institute for Terrestrial Microbiology, Department of Natural Products in Organismic Interactions, Marburg, Germany. ²¹Institut Pasteur, Collection of Cyanobacteria, Paris, France. ²²Pharmaceutical Sciences Division, University of Wisconsin-Madison, Madison, WI, USA. ²³Laboratoire d'Analyses Bioinformatiques pour la Génomique et le Métabolisme, Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry, Université Paris-Saclay, Evry, France. ²⁴Computational Biology Department, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA, USA. ²⁵Microbial Biotechnology, Institute of Biology, Leiden University, Leiden, the Netherlands. ²⁶Department of Microbial Ecology, Netherlands Institute of Ecology,

Wageningen, the Netherlands. ²⁷Interdisciplinary Centre of Marine and Environmental Research, University of Porto, Porto, Portugal. ²⁸Faculty of Sciences, University of Porto, Porto, Portugal. ²⁹Department of Microbiology, University of Helsinki, Helsinki, Finland. ³⁰Department of Chemistry, Yale University, New Haven, CT, USA. ³¹Chemical Biology Institute, Yale University, West Haven, CT, USA. ³²Department of Microbial Pathogenesis, Yale School of Medicine, New Haven, CT, USA. ³³Department of Bacteriology, University of Wisconsin-Madison, Madison, WI, USA. ³⁴Department of Energy Great Lakes Bioenergy Research Center, Wisconsin Energy Institute, University of Wisconsin-Madison, Madison, WI, USA. ³⁵Adrem Data Lab, Department of Computer Science, University of Antwerp, Antwerp, Belgium. ³⁶Molecular Parasitology Unit, Department of Biomedical Sciences, Institute of Tropical Medicine, Antwerp, Belgium. ³⁷Technische Universität Berlin, Institut für Chemie, Berlin, Germany. ³⁸Division of Biological Sciences, University of California San Diego, La Jolla, CA, USA. ³⁹Center for Microbiome Innovation, University of California San Diego, La Jolla, CA, USA. ⁴⁰J. Craig Venter Institute, Genomic Medicine Group, La Jolla, CA, USA. ⁴¹Department of Pediatrics, School of Medicine, University of California San Diego, La Jolla, CA, USA. ⁴²School of Chemistry and Biochemistry, Georgia Institute of Technology, Atlanta, GA, USA. ⁴³Institute of Microbiology, Eidgenössische Technische Hochschule (ETH) Zürich, Zürich, Switzerland. ⁴⁴Department of Biological Chemistry and Molecular Pharmacology, Harvard Medical School, Harvard University, Boston, MA, USA. ⁴⁵Charité, University Medicine Berlin, Berlin-Brandenburg Center for Regenerative Therapy (BCRT), Campus Virchow Klinikum, Berlin, Germany. ⁴⁶Korean Lichen Research Institute, Sunchon National University, Sunchon, Republic of Korea. ⁴⁷Naicons Srl, Milano, Italy. ⁴⁸College of Pharmacy, Sookmyung Women's University, Seoul, Korea. ⁴⁹German Centre for Infection Research (DZIF), Tübingen, Germany. ⁵⁰Department of Chemistry, Northwestern University, Evanston, IL, USA. ⁵¹School of Pharmacy, Sungkyunkwan University, Suwon, Republic of Korea. ⁵²Scripps Institution of Oceanography, University of California San Diego, La Jolla, CA, USA. ⁵³Institute for Pharmaceutical Biology, University of Bonn, Bonn, Germany. ⁵⁴Department of Chemistry, Princeton University, Princeton, NJ, USA. ⁵⁵Section of Microbiology, University of Copenhagen, Copenhagen, Denmark. ⁵⁶Department of Pharmaceutical Sciences, University of Illinois at Chicago, Chicago, IL, USA. ⁵⁷Department of Chemistry, Point Loma Nazarene University, San Diego, CA, USA. ⁵⁸Interfaculty Institute for Microbiology and Infection Medicine Tübingen, Microbiology and Biotechnology, University of Tübingen, Tübingen, Germany. ⁵⁹Centro de Biodiversidad y Descubrimiento de Drogas, Instituto de Investigaciones Científicas y Servicios de Alta

Tecnología, Panama, Republic of Panama. ⁶⁰Carl R. Woese Institute for Genomic Biology and Department of Microbiology, University of Illinois at Urbana-Champaign, Urbana, IL, USA. ⁶¹Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, La Jolla, CA, USA. ⁶²School of Chemistry, University of Nottingham, Nottingham, UK. ⁶³Department of Chemistry and Department of Medicinal Chemistry and Molecular Pharmacology, Purdue University, West Lafayette, IN, USA. ⁶⁴Instituto Federal de Santa Catarina, Florianópolis, Santa Catarina, Brazil. ⁶⁵Phytochemistry and Plant Systematics Department, Division of Pharmaceutical Industries, National Research Centre, Cairo, Egypt. ⁶⁶The Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark. ⁶⁷Department of Biology, Memorial University of Newfoundland, St. John's, Canada. ⁶⁸LOEWE-Centre for Translational Biodiversity Genomics, Frankfurt am Main, Germany. ⁶⁹Departamento de Fisiologia e Farmacologia, Faculdade de Medicina, Universidade Federal do Ceará, Fortaleza, Ceará, Brazil. ⁷⁰University of Strathclyde, Strathclyde Institute of Pharmacy and Biomedical Sciences, Glasgow, UK. ⁷¹School of Computing Science, University of Glasgow, Glasgow, UK. ⁷²Department of Pharmacology and Pediatrics, University of California San Diego, La Jolla, CA, USA. ⁷³These authors contributed equally: Michelle A. Schorn, Stefan Verhoeven. ✉e-mail: pdorrstein@health.ucsd.edu; marnix.medema@wur.nl; justin.vanderhooft@wur.nl

Published online: 15 February 2021

<https://doi.org/10.1038/s41589-020-00724-z>

References

- Tietz, J. I. & Mitchell, D. A. *Curr. Top. Med. Chem.* **16**, 1645–1694 (2016).
- van der Hooft, J. J. J. et al. *Chem. Soc. Rev.* **49**, 3297–3314 (2020).
- Blin, K. et al. *Nucleic Acids Res.* **47**, W81–W87 (2019).
- Skinnider, M. A., Merwin, N. J., Johnston, C. W. & Magarvey, N. A. *Nucleic Acids Res.* **45**, W49–W54 (2017).
- Palaniappan, K. et al. *Nucleic Acids Res.* **48**, D422–D430 (2020).
- Kautsar, S. A. et al. *Nucleic Acids Res.* **48**, D454–D458 (2020).
- Wang, M. et al. *Nat. Biotechnol.* **34**, 828–837 (2016).
- Haug, K. et al. *Nucleic Acids Res.* **48**, D440–D444 (2020).
- Sud, M. et al. *Nucleic Acids Res.* **44**, D463–D470 (2016).
- van der Hooft, J. J. J., Wandý, J., Barrett, M. P., Burgess, K. E. V. & Rogers, S. *Proc. Natl. Acad. Sci. USA* **113**, 13738–13743 (2016).
- Kersten, R. D. et al. *Nat. Chem. Biol.* **7**, 794–802 (2011).
- Hautbergue, T., Jamin, E. L., Debrauwer, L., Puel, O. & Oswald, I. P. *Nat. Prod. Rep.* **35**, 147–173 (2018).
- Jeon, J. E. et al. *Cell* **180**, 176–187.e19 (2020).
- Cao, L. et al. *Cell Syst.* **9**, 600–608.e4 (2019).
- Dejong, C. A. et al. *Nat. Chem. Biol.* **12**, 1007–1014 (2016).
- Goering, A. W. et al. *ACS Cent. Sci.* **2**, 99–108 (2016).
- Benson, D. A. et al. *Nucleic Acids Res.* **41**, D36–D42 (2013).
- Sarkans, U. et al. *Nucleic Acids Res.* **46**, D1266–D1270 (2018).
- Barrett, T. et al. *Nucleic Acids Res.* **40**, D57–D63 (2012).
- Wilkinson, M. D. et al. *Sci. Data* **3**, 160018 (2016).
- Chevrete, M. G. et al. *Nat. Commun.* **10**, 516 (2019).
- Shih, P. M. et al. *Proc. Natl. Acad. Sci. USA* **110**, 1053–1058 (2013).
- van Santen, J. A. et al. *ACS Cent. Sci.* **5**, 1824–1833 (2019).
- Eldjárn, G. H. et al. Ranking microbial metabolomic and genomic links using correlation-based and feature-based scoring functions. Preprint at [bioRxiv](https://doi.org/10.1101/2020.06.12.148205) <https://doi.org/10.1101/2020.06.12.148205> (2020).

25. Misra, B. B., Langefeld, C., Olivier, M. & Cox, L. A. *J. Mol. Endocrinol.* **62**, R21–R45 (2019).

Acknowledgements

The research reported in this publication was supported by an ASDI eScience Grant (ASDI.2017.030) from the Netherlands eScience Center (to J.J.v.d.H. and M.H.M.), a National Institutes of Health (NIH) Genome to Natural Products Network supplementary award (no. U01GM110706 to M.H.M.), a Wageningen Graduate School Postdoc Talent Program fellowship (to M.A.S.), a Marie Skłodowska-Curie Individual Fellowship from the European Union (MSCA-IF-EF-ST-897121 to M.A.S.), the National Science Foundation (NSF) (1817955 to L.M.S. and 1817887 to R.J.D.), a Fundação para a Ciência e Tecnologia (FCT) fellowship (SFRH/BD/136367/2018 to R.C.B.), the National Cancer Institute of the NIH (award no. F32CA221327 to M.W.M.), the University of California, San Diego, Scripps Institution of Oceanography, and two grants from the NIH (Awards GM118815 and 107550 to L.G.), and the National Center for Complementary and Integrative Health of the NIH (award no. R01AT009143 to R.J.T. and N.L.K.).

Author contributions

J.J.v.d.H., M.H.M., and P.C.D. conceived the concept and managed the project. S.V., M.A.S., and J.J.v.d.H. wrote code and developed the platform. M.H.M., L.R., F.H., and J.J.v.d.H. supervised the platform building. All other authors contributed data to the platform, tested it and provided suggestions on how to improve the platform. M.A.S., S.V., P.C.D., M.H.M., and J.J.v.d.H. wrote the manuscript, and all authors contributed to editing the manuscript.

Competing interests

M.H.M. is a co-founder of Design Pharmaceuticals and a member of the scientific advisory board of Hexagon Bio. P.C.D. is a member of the scientific advisory boards of Sirenas and Cybele. N.L.K., W.W.M., and R.J.T. are on the board of directors of MicroMGx. M.W. is a founder of Ometa Labs LLC. A.A.A. is a consultant for Ometa Labs, Clarity Genomics and co-founder of Arome Sciences Inc. William Gerwick, spouse of L.G., has an equity interest in Sirenas Marine Discovery, Inc. and NMRFinder, companies that may potentially benefit from the research results, and also serves on the companies' respective scientific advisory boards. The terms of this last arrangement have been reviewed and approved by the University of California, San Diego (USA), in accordance with its conflict of interest policies.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41589-020-00724-z>.



Open Access This article is licensed under a Creative Commons Attribution 4.0

International License, which permits use,

sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

The Paired Omics Data Platform is running at <https://pairedomicsdata.bioinformatics.nl/>.

The software is licensed under the Apache 2.0 open source licence and the source code can be found on GitHub (<https://github.com/iomega/paired-data-form>), which includes the dependencies of the software. Each software release is archived to Zenodo with <https://doi.org/10.5281/zenodo.2656630> as DOI. Installation guide available at <https://github.com/iomega/paired-data-form/blob/master/README.md>. A demo dataset can be loaded from the project submission form.

The current Paired Omics Data Platform was built using the following software: React (v16.13.1). The platform runs using Docker Compose (v1.25.4) with containers for the web application, web service and redis queue. The web service has an OpenAPI (v3.0.3) specification (<https://www.openapis.org/>) which can be used to submit and retrieve projects in a programmatic manner. Additional information on the genomes entered into the platform was retrieved from GenBank using the public genome identifiers in a project. The Paired Omics Data Platform offers textual searches using elastic search (v7.6.2).

Data analysis

Access to the data and data analysis is done through the same software as data collection, so <https://github.com/iomega/paired-data-form> and <https://doi.org/10.5281/zenodo.2656630>

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The Paired Omics Data Platform is running at <https://pairedomicsdata.bioinformatics.nl/>.

Each project can be downloaded from the website individually as a JSON file. The publicly available (meta)genome and metabolome datasets can be found in their public repositories (NCBI GenBank, RefSeq, JGI IMG, ENA, MGnify, MassIVE, MetaboLights). All PoDP projects are archived monthly to Zenodo with <https://doi.org/10.5281/zenodo.3736430> as DOI.

There are no restrictions to access any of the data and we put in a large effort to make the Paired Omics Data Platform adhere to the FAIR data principles.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

- Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

- Sample size
- Data exclusions
- Replication
- Randomization
- Blinding

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

- | n/a | Included in the study |
|-------------------------------------|--|
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Antibodies |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Eukaryotic cell lines |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Palaeontology |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Animals and other organisms |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Human research participants |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Clinical data |

Methods

- | n/a | Included in the study |
|-------------------------------------|---|
| <input checked="" type="checkbox"/> | <input type="checkbox"/> ChIP-seq |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Flow cytometry |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> MRI-based neuroimaging |