

## Interoperable and scalable metabolomics data analysis with microservices

Payam Emami Khoonsari<sup>1</sup>, Pablo Moreno<sup>2</sup>, Sven Bergmann<sup>3,4</sup>, Joachim Burman<sup>5</sup>, Marco Capuccini<sup>6,7</sup>, Matteo Carone<sup>7</sup>, Marta Cascante<sup>8,9</sup>, Pedro de Atauri<sup>8,9</sup>, Carles Foguet<sup>8,9</sup>, Alejandra Gonzalez-Beltran<sup>10</sup>, Thomas Hankemeier<sup>11</sup>, Kenneth Haug<sup>2</sup>, Sijin He<sup>2</sup>, Stephanie Herman<sup>1,7</sup>, David Johnson<sup>10</sup>, Namrata Kale<sup>2</sup>, Anders Larsson<sup>12,7</sup>, Steffen Neumann<sup>13,14</sup>, Kristian Peters<sup>13</sup>, Luca Pireddu<sup>15</sup>, Philippe Rocca-Serra<sup>10</sup>, Pierrick Roger<sup>16</sup>, Rico Rueedi<sup>3,4</sup>, Christoph Ruttkies<sup>13</sup>, Nouredin Sadawi<sup>17</sup>, Reza M Salek<sup>2</sup>, Susanna-Assunta Sansone<sup>10</sup>, Daniel Schober<sup>13</sup>, Vitaly Selivanov<sup>8,9</sup>, Etienne A. Thévenot<sup>16</sup>, Michael van Vliet<sup>11</sup>, Gianluigi Zanetti<sup>15</sup>, Christoph Steinbeck<sup>2,18</sup>, Kim Kultima<sup>1</sup>, Ola Spjuth<sup>7\*</sup>

1. Department of Medical Sciences, Clinical Chemistry, Uppsala University, Uppsala, Sweden
2. European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Cambridge, United Kingdom
3. Department of Computational Biology, University of Lausanne, Lausanne, Switzerland
4. Swiss Institute of Bioinformatics, Lausanne, Switzerland
5. Department of Neuroscience, Uppsala University, Uppsala, Sweden
6. Department of Information Technology, Uppsala, Sweden
7. Department of Pharmaceutical Biosciences, Uppsala University, Uppsala, Sweden
8. Department of Biochemistry and Molecular Biomedicine, Faculty of Biology, Universitat de Barcelona, Barcelona, Spain
9. Institute of Biomedicine of the Universitat de Barcelona (IBUB) and Associated Unit to CSIC, Barcelona, Spain
10. Oxford e-Research Centre, Department of Engineering Science, University of Oxford, Oxford, United Kingdom
11. Division of Analytical Biosciences, Leiden Academic Centre for Drug Research, Leiden University, Leiden, The Netherlands
12. National Bioinformatics Infrastructure Sweden, Uppsala University, Uppsala, Sweden
13. Department of Stress and Developmental Biology, Leibniz Institute of Plant Biochemistry, Halle, Germany
14. German Centre for Integrative Biodiversity Research (iDiv), Halle-Jena-Leipzig, Germany
15. CRS4: Center for Advanced Studies, Research and Development in Sardinia, Pula, Italy
16. CEA, LIST, Laboratory for Data Analysis and Systems' Intelligence, MetaboHUB, Gif-sur-Yvette, France
17. Faculty of Medicine, Department of Surgery & Cancer, Imperial College London, London, United Kingdom
18. Friedrich-Schiller-University, Jena, Germany

\* Corresponding author

## **Abstract**

Developing a robust and performant data analysis workflow that integrates all necessary components whilst still being able to scale over multiple compute nodes is a challenging task. We introduce a generic method based on the microservice architecture, where software tools are encapsulated as Docker containers that can be connected into scientific workflows and executed in parallel using the Kubernetes container orchestrator. The access point is a virtual research environment which can be launched on-demand on cloud resources and desktop computers. IT-expertise requirements on the user side are kept to a minimum, and established workflows can be re-used effortlessly by any novice user. We validate our method in the field of metabolomics on two mass spectrometry studies, one nuclear magnetic resonance spectroscopy study and one fluxomics study, showing that the method scales dynamically with increasing availability of computational resources. We achieved a complete integration of the major software suites resulting in the first turn-key workflow encompassing all steps for mass-spectrometry-based metabolomics including preprocessing, multivariate statistics, and metabolite identification. Microservices is a generic methodology that can serve any scientific discipline and opens up for new types of large-scale integrative science.

## Introduction

Metabolomics studies measure the occurrence, concentrations and changes of small molecules (metabolites) in organisms, organs, tissues, cells and cellular compartments. Metabolite abundances are assayed in the context of environmental or dietary changes, disease or other conditions<sup>1</sup>. Metabolomics experimental measurements are performed using a variety of spectroscopic methods: the two most common ones are Mass Spectrometry (MS) and Nuclear Magnetic Resonance (NMR). The use of metabolomics as a molecular phenotyping technique is growing across all biomedical domains, due to its ability to reflect the influence of external factors to which an organism is exposed, such as stress, nutrition and disease, subsumed under the term ‘exposome’<sup>2</sup>. Metabolomics data analysis has matured over the years, but is still largely developed and performed at a laboratory level with the use of conventional computing solutions and little standardisation for reproducible research. The PhenoMeNal (**Phenome and Metabolome aNalysis**) project (<http://phenomenal-h2020.eu/home/about>) was conceived to ameliorate this situation by bringing advancements in computing architecture and technology into a modern and easily deployed e-infrastructure – i.e., a computing environment combining hardware and software technology as well as required protocols and data resources – tailored specifically for efficient processing and analysis of molecular phenotype data.

Metabolomics is, as most other omics technologies, characterized by the use of high-throughput experiments that produce large amounts of data<sup>3</sup>. With increasing data size and number of samples, the analysis process becomes intractable for desktop computers due to requirements on

compute cores, memory, storage etc. As a result, large-scale computing infrastructures have become important components in scientific projects<sup>4</sup>. Moreover, making use of such complex computing resources in an analysis workflow presents its own challenges, including achieving efficient job parallelism and scheduling as well as error handling<sup>5</sup>. In addition, configuring the necessary software tools and chaining them together into a complete re-runnable analysis workflow commonly requires substantial IT-expertise, while creating portable and fault-tolerant workflows with a robust audit trail is even more difficult.

Currently, the most common large-scale computational infrastructures in science are shared High-Performance Computing (HPC) systems. Such systems are usually designed primarily to support computationally intensive batch jobs – e.g., for the simulation of physical processes – and are managed by specialized system administrators. This model leads to rigid constraints on the way these resources can be used. For instance, the installation of software must undergo approval and may be restricted, which contrasts with the needs in omics analysis where a multitude of software components of various versions – and their dependencies – are needed, and where these need to be continuously updated.

Cloud computing offers a compelling alternative to shared HPC systems, with the possibility to instantiate and configure on-demand resources such as virtual computers, networks, and storage, together with operating systems and software tools. Users only pay for the time the virtual resources are used, and when they are no longer needed they can be released and incur no further costs for usage or ownership. A few examples of cloud-based systems for metabolomics include

XCMS ONLINE<sup>6</sup>, Chorus ([chorusproject.org](http://chorusproject.org)) and The Metabolomics Workbench ([www.metabolomicsworkbench.org](http://www.metabolomicsworkbench.org)), all of which provide virtual environments that scale with computational demands. However, these applications provide limited flexibility in terms of incorporating and maintaining tools as well as constructing and using customizable workflows.

Along with infrastructure provisioning, software provisioning – i.e., installing and configuring software for users – has also advanced. Consider, for instance, containerization<sup>7</sup>, which allows entire applications with their dependencies to be packaged, shipped and run on a computer but isolated from one another in a way analogous to virtual machines, yet much more efficiently. Containers are more compact, and since they share the same operating system kernel, they are fast to start and stop and incur little overhead in execution. These traits make them an ideal solution to implement light-weight *microservices*, a software engineering methodology in which complex applications are divided into a collection of smaller, loosely coupled components that communicate over a network<sup>8</sup>. Microservices share many properties with traditional always-on web services found on the Internet, but microservices are generally smaller, portable and can be started on-demand within a separate computing environment. Another important feature of microservices is that they have a technology-agnostic communication protocol, and hence can serve as building blocks that can be combined and reused in multiple ways<sup>9</sup>.

Microservices are highly suitable to run in elastic cloud environments that can dynamically grow or shrink on demand, enabling applications to be scaled-up by simply starting multiple parallel instances of the same service. However, to achieve effective scalability a system needs to be

appropriately sectioned into microservice components and the data to be exchanged between the microservices needs to be defined for maximum efficiency— both being challenging tasks.

In this manuscript, we present a method which uses components for metabolomics data analysis encapsulated as microservices and connected into computational workflows to provide complete, ready-to-run, reproducible data analysis solutions that can be easily deployed on desktop computers as well as public and private clouds. Our approach requires virtually no involvement in the setup of computational infrastructure and no special IT skills from the user. We validate the method on four metabolomics studies and show that it enables scalable and interoperable data analysis.

## **Results**

### *Microservices*

In order to construct a microservice architecture for metabolomics we used Docker<sup>10</sup> (<https://www.docker.com/>) containers to encapsulate a large suite of software tools (See Table S1). To automate the instantiation of this cloud-portable microservice-based system and its components for metabolomics analysis, we developed a Virtual Research Environment (VRE) which uses Kubernetes (<https://kubernetes.io/>) to orchestrate containers over multiple compute nodes. Scientists can interact with the microservices programmatically via an Application Programming Interface (API) or via a web-based graphical user interface (GUI), as illustrated in Figure 1. To connect microservices into computational workflows, the two frameworks Galaxy<sup>11</sup> and Luigi (<https://github.com/spotify/luigi>) were adapted to execute jobs on Kubernetes. Galaxy

is a web-based interface for individual tools and allows users to share workflows and analysis histories. Luigi on the other hand focuses on scheduled execution, monitoring, visualization and the implicit dependency resolution of tasks<sup>12</sup>. These basic infrastructure services, together with the Jupyter notebook<sup>13</sup> interactive programming environment, are deployed as long-running services in the VRE, whereas the other analysis tools are deployed as transient compute jobs to be used on-demand. System and client applications were developed for launching the VRE on desktop computers, public and private cloud providers, automating all steps required to instantiate the virtual infrastructures. The PhenoMeNal consortium maintains a web portal (<https://portal.phenomenal-h2020.eu>) providing a GUI for launching VREs on a selection of the largest public cloud providers, including Amazon Web Services, Microsoft Azure and Google Cloud Platform, or on private OpenStack-based installations. The containers provisioned by PhenoMeNal comprise tools built as open source software that are available in a public repository such as GitHub, and are subject to continuous integration testing. The containers that satisfy testing criteria are pushed to a public container repository, and containers that are included in stable VRE releases are also pushed to Biocontainers<sup>9</sup>.

***Demonstrator 1: Scalability of microservices in a cloud environment in the analysis of a human renal proximal tubule cells dataset***

The objective of this analysis was to demonstrate the scalability of an existing workflow on a large dataset (MetaboLights<sup>14</sup> ID: MTBLS233, <http://www.ebi.ac.uk/metabolights/MTBLS233><sup>15</sup>). The experiment includes 528 mass spectrometry samples from whole cell lysates of human renal proximal tubule cells that were

pre-processed through a five-step workflow (consisting of peak picking, feature finding, linking, file filtering and exporting) using the OpenMS software<sup>16</sup> as illustrated in Figure 2. This preprocessing workflow was reimplemented using Docker containers and run using the Luigi workflow engine. Scalability of concurrent running tools (on 40 Luigi workers, each worker receives tasks from the scheduler and executes them) was measured using weak scaling efficiency (WSE), where the workload assigned to each worker stays constant and additional workers are used to solve a larger total problem. The WSE was computed to reach 88% with an execution time of approximately four hours (online methods, Figure S2), compared with the ideal case of 100% where linear scaling is achieved if the run time stays constant while the workload is increased. In addition, the final result of the workflow (online methods, Figure S3) was identical to that presented by the original MTBLS233 study (Ranninger et al.<sup>15</sup>) in negative ionization mode. However, in the positive ionization mode, one  $m/z$  feature was found in a different group ( $m/z$  range 400-1000) than it was originally reported by Ranninger et al. ( $m/z$  range 200-400).

### ***Demonstrator 2: Start-to-end LC-MS-analysis workflow on Multiple Sclerosis data***

The objective of this analysis was to demonstrate interoperability as well as to present a real-world scenario in which patients' data are processed using a microservices-based platform. We used a dataset consisting of 37 clinical cerebrospinal fluid (CSF) samples including thirteen relapsing-remitting multiple sclerosis (RRMS) patients and 14 secondary progressive multiple sclerosis (SPMS) patients as well as ten non-multiple sclerosis controls. Twenty-six quality controls (19 blank and 7 dilution series samples) were also added to the experiment. In addition,



8 pooled CSF samples containing MS/MS data were included in the experiment for improving identification (MetaboLights ID: MTBLS558, <http://www.ebi.ac.uk/metabolights/MTBLS558>). The samples were processed and analysed on the Galaxy platform<sup>11</sup> applying the Liquid Chromatography-MS (LC-MS) workflow illustrated in Figure 3, running in a PhenoMeNal VRE behind the Uppsala University Hospital firewall to be compliant with local ELSI (Ethics, Legal, Social implications) regulations. The result of multivariate analysis showed a clear difference (Figure 4A) in the metabolic constitution between the three disease groups of RRMS, SPMS and non-multiple sclerosis controls. In addition, the univariate analysis resulted in a total of three metabolites being significantly altered ( $p < 0.05$ ) between multiple sclerosis subtypes and control samples, namely alanyltryptophan and indoleacetic acid with higher and linoleoyl ethanolamide with lower abundance in both RRMS and SPMS compared to controls (Figure 4B).

### ***Demonstrator 3: 1D NMR-analysis workflow on human type 2 diabetes mellitus data***

This NMR-based metabolomics study was originally performed by Salek et al.<sup>17</sup> on urine of type 2 diabetes mellitus (T2DM) patients and controls (MetaboLights ID: MTBLS1, <http://www.ebi.ac.uk/metabolights/MTBLS1>). In total, 132 samples (48 T2DM and 84 controls) were processed using the workflow shown in Figure 5. A total of 726 metabolites were quantified and used to perform Orthogonal Projections to Latent Structures Discriminant Analysis (OPLS-DA) which resulted in a clear separation between T2DM and controls (Figure 5), reproducing previous findings<sup>17</sup>.

### ***Demonstrator 4: Start-to-end fluxomics workflow on HUVEC cells under hypoxia***

The purpose of this demonstrator was to show the integrated use of separately developed tools covering subsequent steps of the study of metabolic fluxes based on  $^{13}\text{C}$  stable isotope-resolved metabolomics (*SIRM*)<sup>18,19,20</sup>. Here we implemented the analysis of flux distributions in HUVEC cells under hypoxia (MetaboLights ID: MTBLS412, <http://www.ebi.ac.uk/metabolights/MTBLS412>), from raw mass spectra contained in netCDF files, using the workflow illustrated in Figure 6. The result was a detailed description of the magnitudes of the fluxes through the reactions accounting for glycolysis and pentose phosphate pathway.

## Discussion

Implementing the different tools and processing steps of a data analysis workflow as separate services that are made available over a network was in the spotlight in the early 2000's<sup>21</sup> as service-oriented architectures (SOA) in science. At that time, web services were commonly deployed on physical hardware and exposed and consumed publicly over the internet. However, it soon became evident that this architecture did not fulfill its promises as it did not scale well from a computational perspective. In addition, the web services were not portable and mirroring them was complicated (if at all possible). Furthermore, API changes and frequent services outage made it frustrating to connect them into functioning computational workflows. Ultimately, the ability to replicate an analysis on local and remote hardware (such as a computer cluster) was very difficult due to heterogeneity in the computing environments.

At first sight microservices might seem similar to abovementioned SOA web services, but microservices are generally executed in virtual environments (abstracting over OS and hardware architectures) in such a way that they are only instantiated and executed on-demand, and then terminated when they are no longer needed. This makes such virtual environments inherently portable and they can be launched on demand on different platforms (e.g., a laptop, a powerful physical server or an elastic cloud environment). A key aspect is that workflows are still executed identically, agnostic of the underlying hardware platform. Container-based microservices provide a wide flexibility in terms of versioning, allowing the execution of newer and older versions of each container as needed for reproducibility. Since all software dependencies are encompassed within the container, which is versioned, the risk of workflow failure due to API changes is minimized. An orchestration framework such as Kubernetes further allows for managing errors in execution and transparently handles the restarting of services. Hence, technology has caught up with service-oriented science, and microservices have taken the methodology to the next level, alleviating many of the previous problems related to scalability, portability and interoperability of software tools. This is advantageous in the context of omics analysis, which produces multidimensional data sets reaching beyond gigabytes, on into terabytes, leading to ever-increasing demand on processing performance<sup>22,23</sup>.

In Demonstrator 1, we showed that microservices enable highly efficient and scalable data analyses by executing individual modules in parallel, and that they effectively harmonize with on-demand elasticity of the cloud computing paradigm. The reached scaling efficiency of ~88% indicates remarkable performance achieved on generic cloud providers. Furthermore, although

our results in positive ionization model was slightly different to that of Ranninger et al.<sup>15</sup>, the results of our analysis were reproducible regardless of the platform used to perform the computations, indicating a level of replicability of study results and reusability of workflows in the analysis that - to the best of our knowledge - has never been reported before in metabolomics data analysis.

In addition to the fundamental demand for high performance, the increased throughput and complexity of omics experiments has led to a large number of sophisticated computational tools<sup>24</sup>, which in turn necessitates integrative workflow engines<sup>25</sup>. In order to integrate new tools in such workflow engines, compatibility of the target environment, tools and APIs needs to be considered<sup>25</sup>. Containerization facilitates this by providing a platform-independent virtual environment for developing and running the individual tools. However, the problem of compatibility between tools/APIs and data formats remains and needs to be tackled by international consortia similarly to what PhenoMeNal addresses in metabolomics by promoting and strictly adhering to FAIR Data Principles<sup>26</sup>. PhenoMeNal also overcomes the currently non-trivial task of instantiating the complete microservice environments through a web portal that allows for convenient deployment of the VRE on public cloud providers. Moreover, using the PhenoMeNal web portal, microservices and VREs can be deployed on a trusted private cloud instance or a local physical server on an internal network, such as within a hospital network, allowing for levels of isolation and avoiding transfer of data across untrusted networks which often are requirements in the analysis of sensitive data. This was highlighted in Demonstrator 2, where a complete start-to-end workflow was run on the Galaxy platform on a secure server at

Uppsala University Hospital, Sweden, leading to the identification of novel disease fingerprints in the CSF metabolome of RRMS and SPMS patients. It is worth mentioning that the selected metabolites were part of tryptophan metabolism (alanyltryptophan and indoleacetic acid) and endocannabinoids (linoleoyl ethanolamide), both of which have been previously implicated in multiple sclerosis<sup>27-32</sup>. However, since the cross-validated predictive performance ( $Q2Y = 0.286$ ) is not much higher than some of the models generated after random permutation of the response (Figure 4A), the quality of the model needs to be confirmed in a future study on an independent cohort of larger size.

In Demonstrator 3, we highlighted the fact that the microservice architecture is indeed domain-agnostic and is not limited to a particular assay technology, i.e. mass spectrometry. Using a fully automated 1D NMR workflow, we showed that the pattern of the metabolite expression is different between type 2 diabetic and healthy controls, and that a large number of metabolites contribute to such separation. The preprocessing of NMR-based experiments can be performed with minimal effort on other studies (i.e. simply by providing a MetaboLights accession number), leading to the capability to re-analyze data and compare the results with the original publication findings. Furthermore, it demonstrates the value of standardised dataset descriptions using nmrML<sup>33</sup> and ISA format<sup>34,35</sup> for representing NMR based studies, as well as the potential of the PhenoMenNal VRE to foster reproducibility.

A complete understanding of metabolic function implies a complete metabolic profile, but also knowledge of the associated distribution of metabolic fluxes in the metabolic network. In

Demonstrator 4, the microservices architecture is applied to deal with flux distributions derived from the application of stable isotope resolved metabolomics. Here we showed high rate of glycolysis in cell cultured in hypoxia which is consistent with the one expected for endothelial cells<sup>36</sup> and also further confirmation on how these cells maintain energy in low oxygen environments and without oxidative phosphorylation<sup>37,38</sup>.

While microservices are not confined to metabolomics and generally applicable to a large variety of applications, there are some important implications and limitations of the method. Firstly, tools need to be containerized in order to operate in the environment. This is however not particularly complex, and an increasing number of developers provide containerized versions of their tools on public container repositories such as Dockerhub or Biocontainers<sup>9</sup>. Secondly, uploading data to a cloud-based system can take a considerable amount of time, and having to re-do this every time a VRE is instantiated can be time-consuming. This can be alleviated by using persistent storage on a cloud resource, but the availability of such storage varies between different cloud providers. Further, the storage system can become a bottleneck when many services try to access a shared storage. We observe that using a distributed storage system with multiple storage nodes can drastically increase performance, and the PhenoMeNal VRE comes with a distributed storage system by default. When using a workflow system to orchestrate the microservices, stability and scalability are inherently dependent on the workflow system's job runner. We observed that in the Galaxy workflow engine, executing a large number of jobs resulted in the VRE becoming unresponsive whereas the Luigi engine did not have these shortcomings. Although this problem can be resolved by defining the required resources in the

Galaxy job runner for each tool, the issue of knowing how much computational resources a specific tool needs remains. This can be partially addressed by tool/workflow developers to estimate the required resources for their tools and workflows. With cloud and microservices maturing, workflow systems will need to evolve and further embrace the new possibilities of these infrastructures. Also, not all research can be easily pipelined, for example exploratory research might be better carried out in an ad-hoc manner than with workflows and the overhead this implies. A Jupyter Notebook as used in in Demonstrator 1 or embedded in Galaxy<sup>39</sup> constitute promising ways to make use of microservices for interactive analysis.

In summary, we showed that microservices allow for efficiently scaling up analyses on multiple computational nodes, enabling the processing of large data sets. By applying a number of data (mzML<sup>40</sup> , nmrML) and metadata standards (ISA serialisations for study descriptions<sup>34,35</sup>), we also demonstrated a level of interoperability which has never been achieved in the context of metabolomics, by providing completely automated start-to-end analysis workflows for mass spectrometry and NMR data. The PhenoMeNal VRE realizes the notion of “bringing compute to the data” by enabling the instantiation of complete virtual infrastructures close to large datasets that could not be uploaded over the internet, and can also be launched close to ELSI sensitive data that is not allowed to leave a secure computing environment. While the current PhenoMeNal VRE implementation uses Docker for software containers and Kubernetes for container orchestration, the microservice methodology is general and not restricted to these frameworks. In addition, we emphasise that the presented methodology goes beyond

metabolomics and can be applied to virtually any field, lowering the barriers for taking advantage of cloud infrastructures and opening up for large-scale integrative science.

## **Acknowledgements**

This research was supported by The European Commission's Horizon 2020 programme funded under grant agreement number 654241 (PhenoMeNal), The Swedish Research Council FORMAS, Uppsala Berzelii Technology Centre for Neurodiagnostics and Åke Wiberg Foundation. We kindly acknowledge contributions to cloud resources by SNIC Science Cloud, Embassy Cloud, and CityCloud. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

## **Author contributions**

KK, MAC, MC, PEK, SH contributed to Demonstrator 1. CR, KK, KP, PEK, SH, SN contributed to Demonstrator 2. KK designed the study in Demonstrator 2. JB performed collection of samples and characterization of the multiple sclerosis cohort. SH performed mass spectrometry experiment in Demonstrator 2. DS, KP, PEK, PM, RMS, contributed to Demonstrator 3. AGB, CF, DJ, MCA, MVV, PDA, PM, PRS, SAS, TH and VS contributed to Demonstrator 4. GZ, LP, PEK and PM contributed to development of Galaxy in Kubernetes. AL and MC contributed to the development of Luigi in Kubernetes. AL, MAC, MC and NS developed KubeNow. PM contributed to Galaxy-Kubernetes. EAT and PR contributed to containerizing of Workflow4Metabolomics tools. AGB, DJ, PRS and SAS contributed to ISA-API. DJ, EAT, KP, MVV, NS, OS, PEK, PM, PR, PRS, RMS, RR and SB were involved in



testing the containers and the VRE. PM, SIH and KH were involved in development and maintenance of the portal. MVV, PM and RMS contributed to the release. NK coordinated the PhenoMeNal project. CS conceived and managed PhenoMeNal project. OS conceived and coordinated the project and e-infrastructure. All authors contributed to manuscript writing.

### **Competing Financial Interests statement**

The authors declare no competing financial interest.

## References

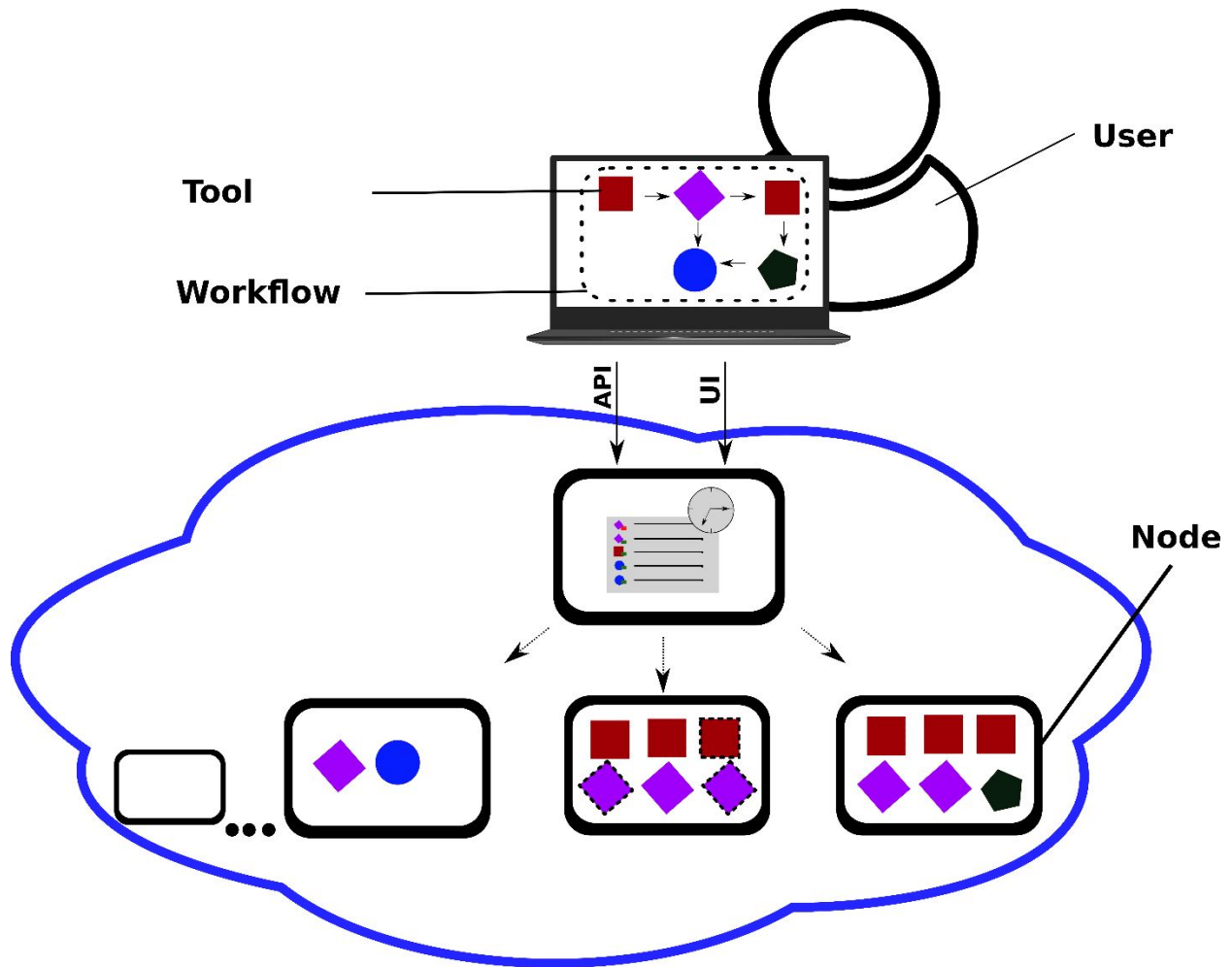
1. Nicholson, J. K. & Wilson, I. D. Opinion: understanding ‘global’ systems biology: metabonomics and the continuum of metabolism. *Nat. Rev. Drug Discov.* **2**, 668–676 (2003).
2. Manrai, A. K. *et al.* Informatics and Data Analytics to Support Exposome-Based Discovery for Public Health. *Annu. Rev. Public Health* **38**, 279–294 (2017).
3. Montenegro-Burke, J. R. *et al.* Data Streaming for Metabolomics: Accelerating Data Processing and Analysis from Days to Minutes. *Anal. Chem.* **89**, 1254–1259 (2017).
4. Liew, C. S. *et al.* Scientific Workflows. *ACM Computing Surveys* **49**, 1–39 (2016).
5. Suplatov, D., Popova, N., Zhumatiy, S., Voevodin, V. & Švedas, V. Parallel workflow manager for non-parallel bioinformatic applications to solve large-scale biological problems on a supercomputer. *J. Bioinform. Comput. Biol.* **14**, 1641008 (2016).
6. Warth, B. *et al.* Metabolizing Data in the Cloud. *Trends Biotechnol.* **35**, 481–483 (2017).
7. Silver, A. Software simplified. *Nature* **546**, 173–174 (2017).
8. Newman, S. *Building Microservices*. (‘O’Reilly Media, Inc.’, 2015).
9. da Veiga Leprevost, F. *et al.* BioContainers: an open-source and community-driven framework for software standardization. *Bioinformatics* **33**, 2580–2582 (2017).
10. Merkel, D. Docker: Lightweight Linux Containers for Consistent Development and Deployment. *Linux J.* **2014**, (2014).
11. Goecks, J., Nekrutenko, A., Taylor, J. & Galaxy Team. Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol.* **11**, R86 (2010).
12. Leipzig, J. A review of bioinformatic pipeline frameworks. *Brief. Bioinform.* **18**, 530–536 (2017).
13. Kluyver, T. *et al.* Jupyter Notebooks ? a publishing format for reproducible computational workflows. (2016).

14. Haug, K. *et al.* MetaboLights--an open-access general-purpose repository for metabolomics studies and associated meta-data. *Nucleic Acids Res.* **41**, D781–6 (2013).
15. Ranninger, C. *et al.* Improving global feature detectabilities through scan range splitting for untargeted metabolomics by high-performance liquid chromatography-Orbitrap mass spectrometry. *Anal. Chim. Acta* **930**, 13–22 (2016).
16. Sturm, M. *et al.* OpenMS - an open-source software framework for mass spectrometry. *BMC Bioinformatics* **9**, 163 (2008).
17. Salek, R. M. *et al.* A metabolomic comparison of urinary changes in type 2 diabetes in mouse, rat, and human. *Physiol. Genomics* **29**, 99–108 (2007).
18. Buescher, J. M. *et al.* A roadmap for interpreting (13)C metabolite labeling patterns from cells. *Curr. Opin. Biotechnol.* **34**, 189–201 (2015).
19. Niedenführ, S., Wiechert, W. & Nöh, K. How to measure metabolic fluxes: a taxonomic guide for 13 C fluxomics. *Curr. Opin. Biotechnol.* **34**, 82–90 (2015).
20. King, Z. A. *et al.* Escher: A Web Application for Building, Sharing, and Embedding Data-Rich Visualizations of Biological Pathways. *PLoS Comput. Biol.* **11**, e1004321 (2015).
21. Foster, I. Service-oriented science. *Science* **308**, 814–817 (2005).
22. Marx, V. Biology: The big challenges of big data. *Nature* **498**, 255–260 (2013).
23. Schadt, E. E., Linderman, M. D., Sorenson, J., Lee, L. & Nolan, G. P. Computational solutions to large-scale data management and analysis. *Nat. Rev. Genet.* **11**, 647–657 (2010).
24. Berger, B., Peng, J. & Singh, M. Computational solutions for omics data. *Nat. Rev. Genet.* **14**, 333–346 (2013).
25. Di Tommaso, P. *et al.* Nextflow enables reproducible computational workflows. *Nat. Biotechnol.* **35**, 316–319 (2017).
26. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and

- stewardship. *Sci Data* **3**, 160018 (2016).
27. Lovelace, M. D. *et al.* Current Evidence for a Role of the Kynurenine Pathway of Tryptophan Metabolism in Multiple Sclerosis. *Front. Immunol.* **7**, 246 (2016).
  28. Lim, C. K. *et al.* Kynurenine pathway metabolomics predicts and provides mechanistic insight into multiple sclerosis progression. *Sci. Rep.* **7**, 41473 (2017).
  29. Amirkhani, A. *et al.* Interferon-beta affects the tryptophan metabolism in multiple sclerosis patients. *Eur. J. Neurol.* **12**, 625–631 (2005).
  30. Centonze, D. *et al.* The endocannabinoid system is dysregulated in multiple sclerosis and in experimental autoimmune encephalomyelitis. *Brain* **130**, 2543–2553 (2007).
  31. Zamberletti, E., Rubino, T. & Parolaro, D. The Endocannabinoid System and Schizophrenia: Integration of Evidence. *Curr. Pharm. Des.* **18**, 4980–4990 (2012).
  32. Baker, D. & Pryce, G. The endocannabinoid system and multiple sclerosis. *Curr. Pharm. Des.* **14**, 2326–2336 (2008).
  33. Schober, D. *et al.* nmrML: a community supported open data standard for the description, storage, and exchange of NMR data. *Anal. Chem.* (2017). doi:10.1021/acs.analchem.7b02795
  34. Rocca-Serra, P. *et al.* Data standards can boost metabolomics research, and if there is a will, there is a way. *Metabolomics* **12**, 14 (2016).
  35. Sansone, S.-A. *et al.* Toward interoperable bioscience data. *Nat. Genet.* **44**, 121–126 (2012).
  36. Iyer, N. V. *et al.* Cellular and developmental control of O<sub>2</sub> homeostasis by hypoxia-inducible factor 1 $\alpha$ . *Genes Dev.* **12**, 149–162 (1998).
  37. Eelen, G., de Zeeuw, P., Simons, M. & Carmeliet, P. Endothelial cell metabolism in normal and diseased vasculature. *Circ. Res.* **116**, 1231–1244 (2015).
  38. Polet, F. & Feron, O. Endothelial cell metabolism and tumour angiogenesis: glucose and glutamine as essential fuels and lactate as the driving force. *J. Intern. Med.* **273**, 156–165 (2013).

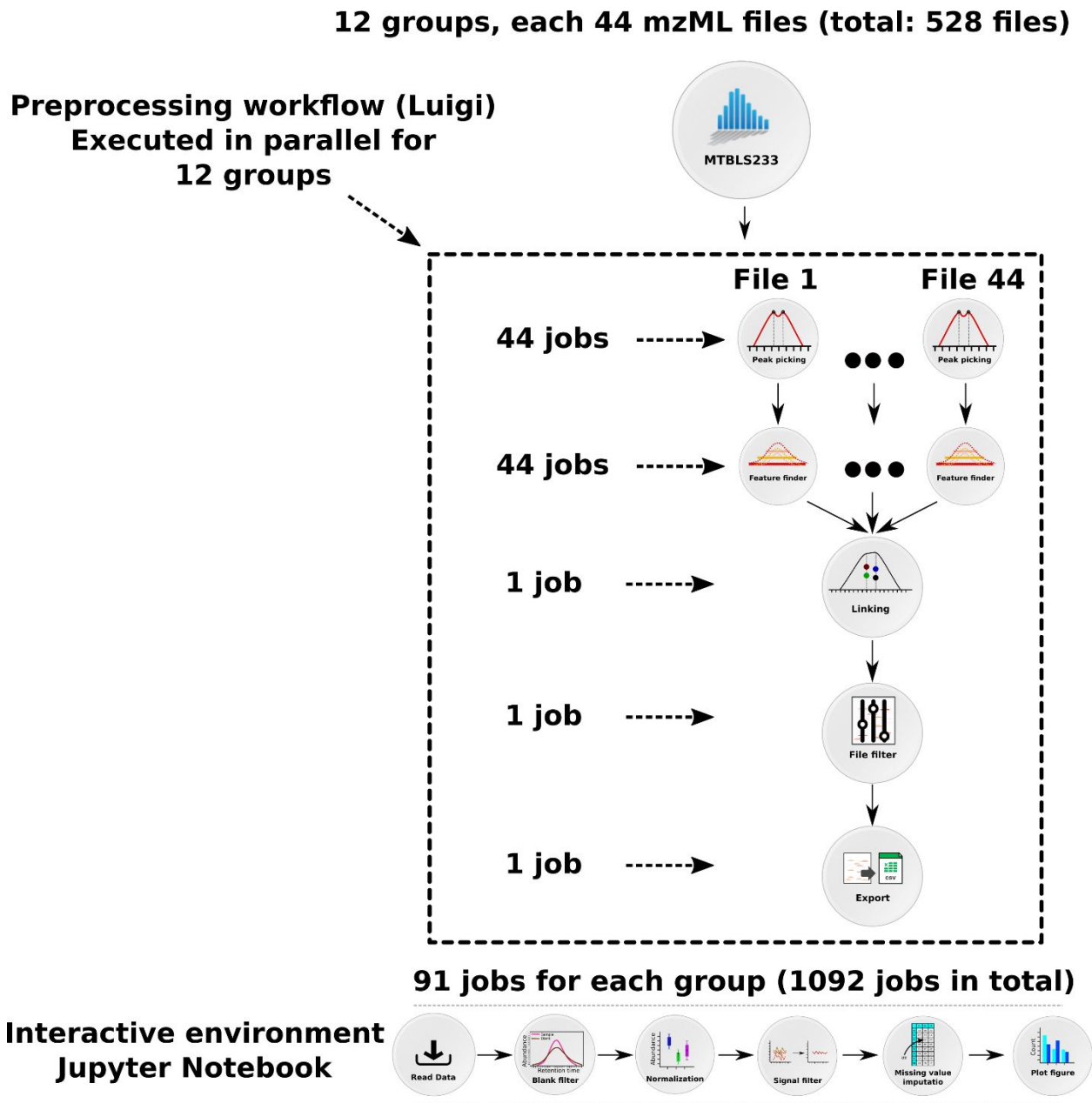
39. Grüning, B. A. *et al.* Jupyter and Galaxy: Easing entry barriers into complex data analyses for biomedical researchers. *PLoS Comput. Biol.* **13**, e1005425 (2017).
40. Martens, L. *et al.* mzML--a community standard for mass spectrometry data. *Mol. Cell. Proteomics* **10**, R110.000133 (2011).
41. Smith, C. A., Want, E. J., O'Maille, G., Abagyan, R. & Siuzdak, G. XCMS: Processing Mass Spectrometry Data for Metabolite Profiling Using Nonlinear Peak Alignment, Matching, and Identification. *Anal. Chem.* **78**, 779–787 (2006).
42. Kuhl, C., Tautenhahn, R., Böttcher, C., Larson, T. R. & Neumann, S. CAMERA: an integrated strategy for compound spectra extraction and annotation of liquid chromatography/mass spectrometry data sets. *Anal. Chem.* **84**, 283–289 (2012).
43. Giacomoni, F. *et al.* Workflow4Metabolomics: a collaborative research infrastructure for computational metabolomics. *Bioinformatics* **31**, 1493–1495 (2015).
44. Thévenot, E. A., Roux, A., Xu, Y., Ezan, E. & Junot, C. Analysis of the Human Adult Urinary Metabolome Variations with Age, Body Mass Index, and Gender by Implementing a Comprehensive Workflow for Univariate and OPLS Statistical Analyses. *J. Proteome Res.* **14**, 3322–3335 (2015).
45. Jacob, D., Deborde, C., Lefebvre, M., Maucourt, M. & Moing, A. NMRProcFlow: a graphical and interactive tool dedicated to 1D spectra processing for NMR-based metabolomics. *Metabolomics* **13**, 36 (2017).

## Main figures



**Figure 1:** Overview of the components in a microservices-based framework. Complex applications are divided into smaller, focused and well-defined (micro-) services. These services are independently deployable and can communicate with each other, which allows to couple them into data processing workflows. The user can interact with the framework programmatically via an Application Program Interface (API) or via a graphical user interface (GUI) to construct or run workflows of different services, which are executed independently. Multiple instances of services can be launched to execute tasks in parallel, which effectively can be used to scale analysis over multiple compute nodes. When run in an

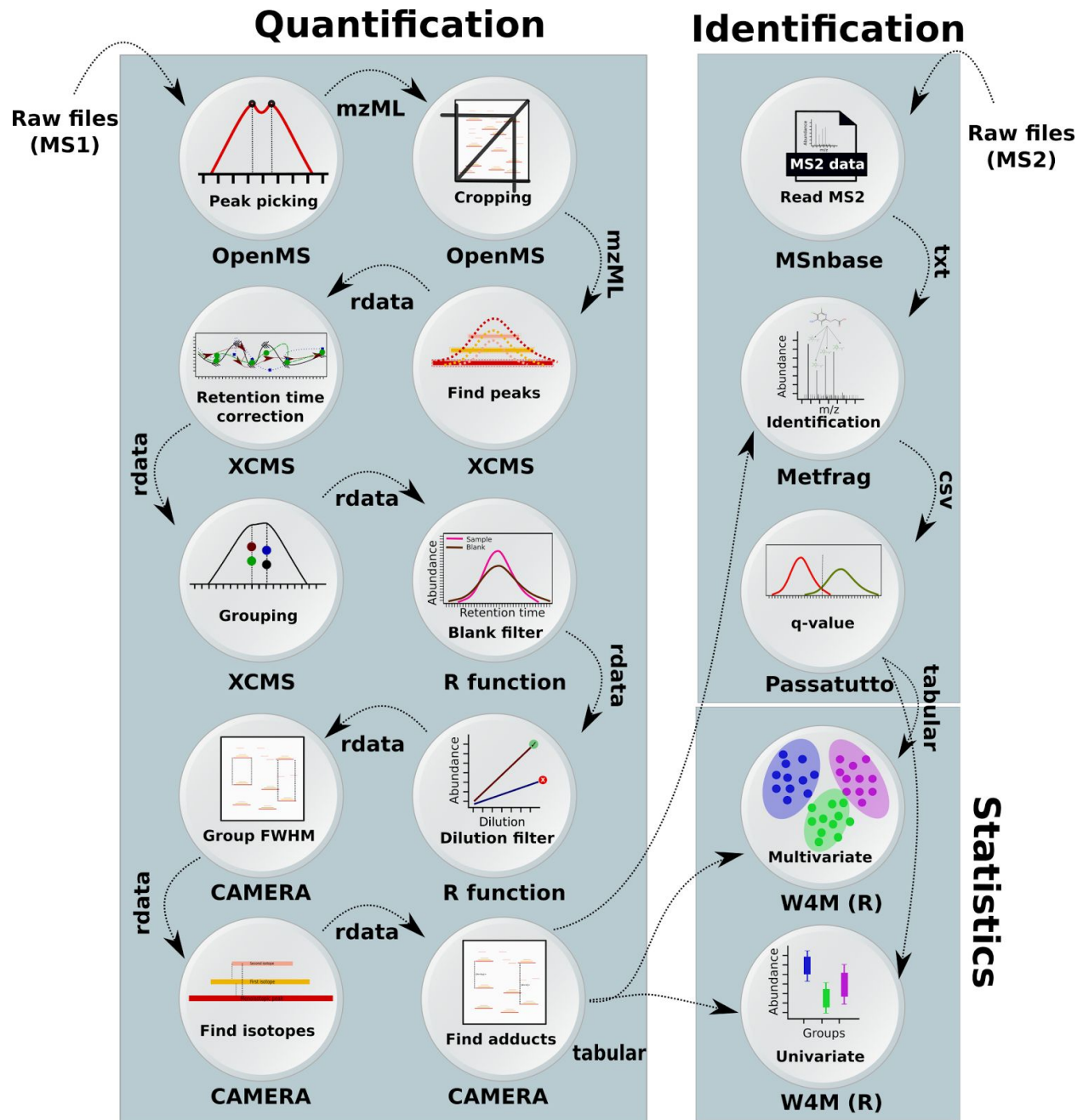
elastic cloud environment, virtual resources can be added or removed depending on the computational requirements.



**Figure 2.** Diagram of scalability-testing on the metabolomics dataset (MetaboLights ID: MTBLS233) in Demonstrator 1 to illustrate the scalability of a microservice approach. The preprocessing workflow is

composed of 5 OpenMS tasks that were run in parallel over the 12 groups in the dataset using the Luigi workflow system. The first two tasks, peak picking (528 tasks) and feature finding (528 tasks), are trivially parallelizable, hence they were run concurrently for each sample. The subsequent feature linking task needs to process all of the samples in a group at the same time, therefore 12 of these tasks were run in parallel. In order to maximize the parallelism, each feature linker container (microservice) was run on 2 CPUs. Feature linking produces a single file for each group, that can be processed independently by the last two tasks: file filter (12 tasks) and text exporter (12 tasks), resulting in total of 1092 tasks. The downstream analysis consisted of 6 tasks that were carried out in a Jupyter Notebook. Briefly, the output of preprocessing steps was imported into R and the unstable signals were filtered out. The missing values were imputed and the resulting number of features were plotted.

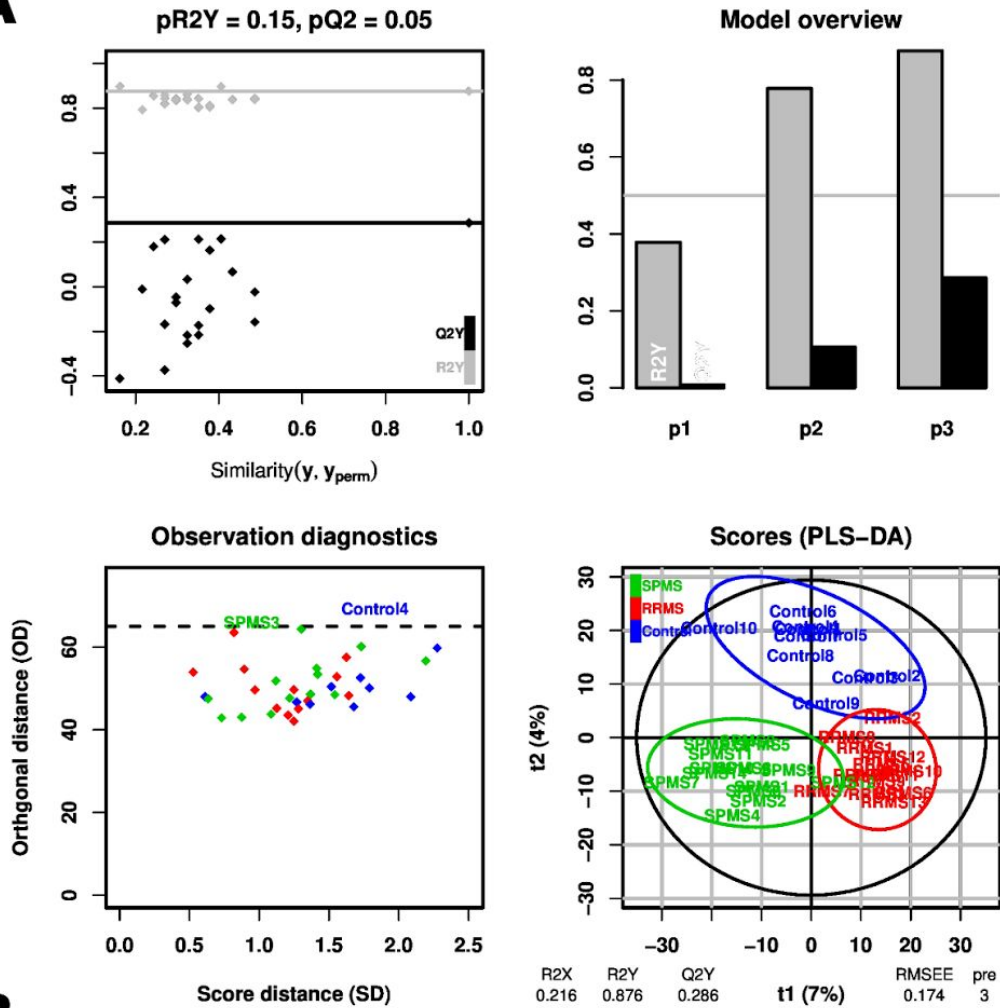




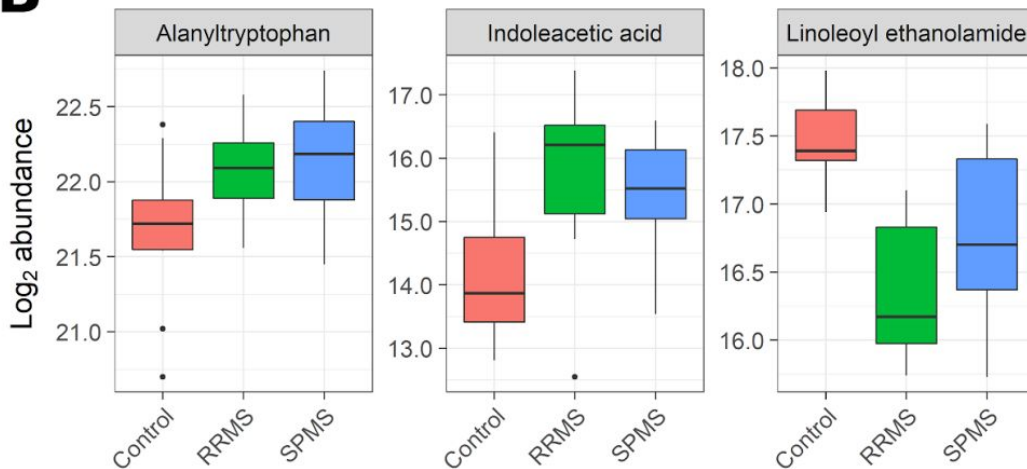
**Figure 3.** Overview of the workflow used to process multiple-sclerosis samples in Demonstrator 2, where a workflow was composed of the microservices using the Galaxy system. The data was centroided and limited to a specific mass over charge ( $m/z$ ) range using OpenMS tools. The mass traces quantification and retention time correction was done via XCMS<sup>41</sup>. Unstable signals were filtered out based on the blank and dilution series samples using an in-house function (implemented in R). Annotation

of the peaks was performed using CAMERA<sup>42</sup>. To perform the metabolite identification, the tandem spectra from the MS/MS samples in mzML format were extracted using MSnbase and passed to MetFrag. The MetFrag scores were converted to q-values using Passatutto software. The result of identification and quantification were used in “Multivariate” and “Univariate” containers from Workflow4Metabolomics<sup>43</sup> to perform Partial Least Squares Discriminant Analysis (PLS-DA)<sup>44</sup>.

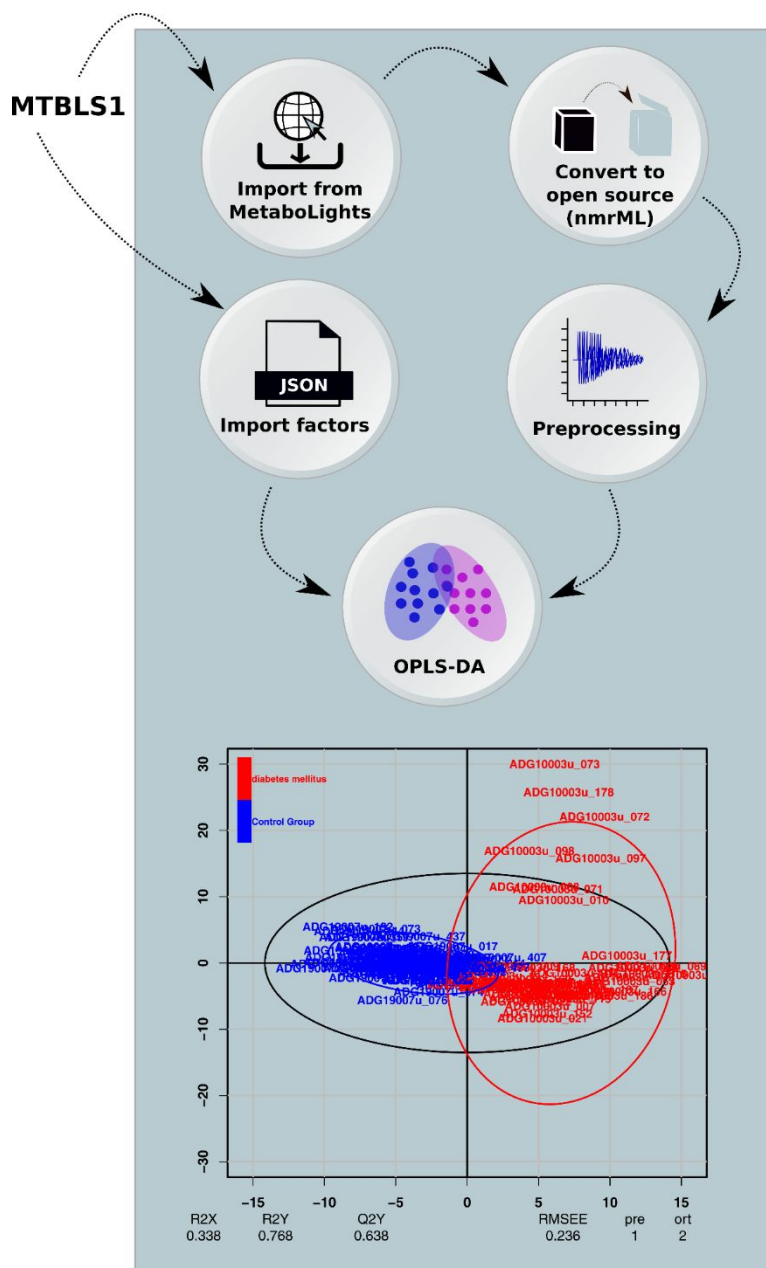
**A**



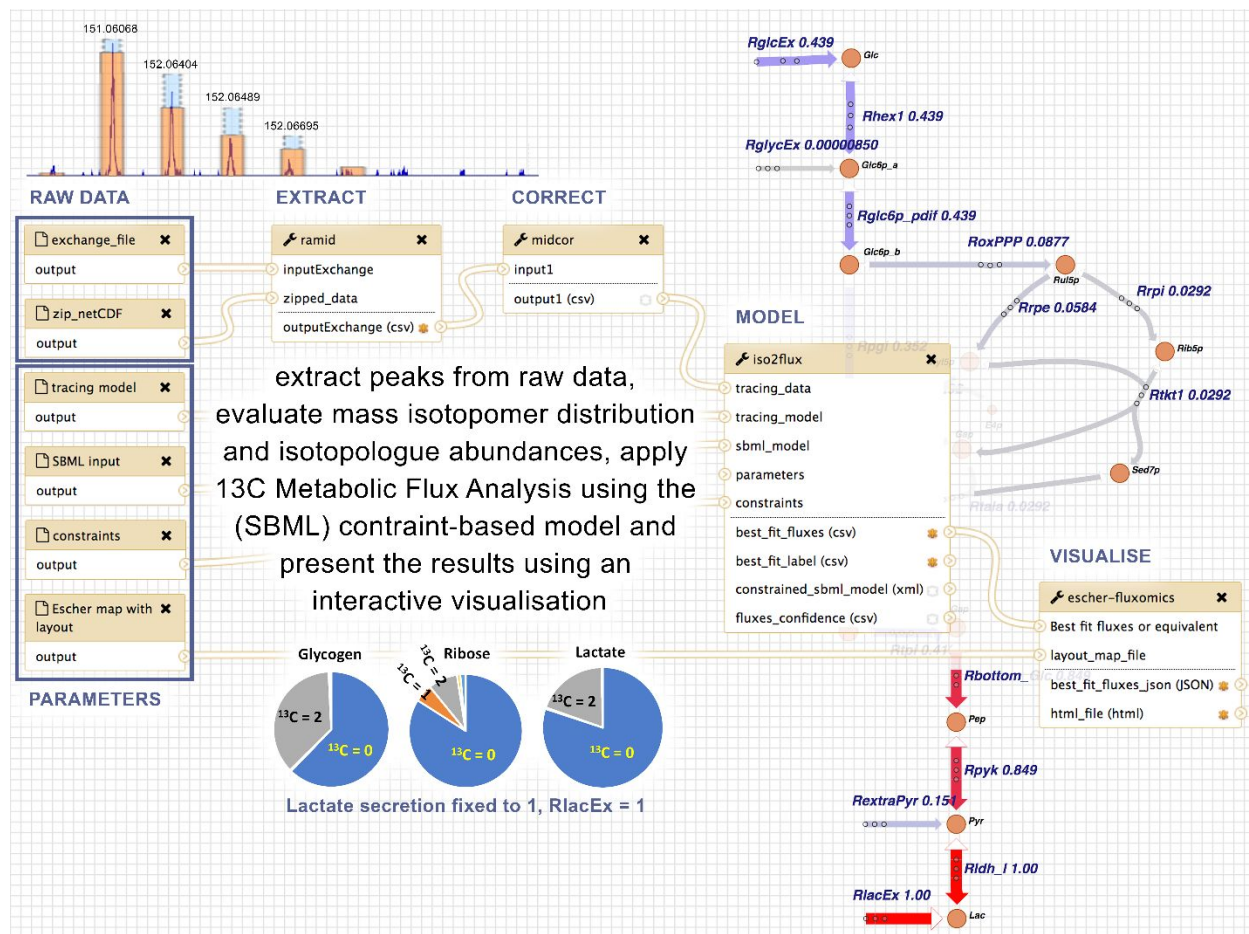
**B**



**Figure 4.** The results from analysis of multiple sclerosis data in Demonstrator 2, presenting new scientifically useful biomedical knowledge. A) The PLS-DA results suggest that the metabolite distribution in the RRMS and SPMS samples are different to controls. B) Three metabolites were identified as differentially regulated between multiple sclerosis subtypes and control samples, namely alanyltryptophan and indoleacetic acid with higher and Linoleoyl ethanolamide with lower abundance in both RRMS and SPMS compared to controls. Abbr., RRMS: relapsing-remitting multiple sclerosis, SPMS: secondary progressive multiple sclerosis.



**Figure 5.** Overview of the NMR workflow in Demonstrator 3. The raw NMR data was automatically import from Metabolights (ISA-Tab) database and converted to open source nmrML format. The preprocessing was performed using the nmr1d package part of nmrprocflow<sup>45</sup> tools. All study factors were imported from MetaboLights and were fed to the multivariate node to perform an OPLS-DA.



**Figure 6:** Overview of the workflow for fluxomics, with Ramid, Midcor, Iso2Flux and Escher-fluxomics tools supporting subsequent steps of the analysis. The example refers to HUVEC cells incubated in the presence of [1,2-<sup>13</sup>C<sub>2</sub>]glucose and label (<sup>13</sup>C) propagation to glycogen, RNA ribose and lactate measured by mass spectrometry. Ramid reads the raw netCDF files, corrects baseline and extracts the peak intensities. The resulting peak intensities are corrected (natural abundance, overlapping peaks) by Midcor, which provides isotopologue abundances. Isotopologue abundances, together with a model description (SBML model, tracing data, constraints), are used by Iso2Flux to provide flux distributions through glycolysis and

pentose-phosphate pathways, which are shown as numerical values associated to a metabolic scheme of the model by the Escher-fluxomics tool.