




FULL ARTICLE

Cloud-based archived metabolomics data: A resource for in-source fragmentation/annotation, meta-analysis and systems biology

Amelia Palermo¹ | Tao Huan^{1,2} | Duane Rinehart¹ | Markus M. Rinschen¹ |
Shuzhao Li⁵ | Valerie B. O'Donnell⁶ | Eoin Fahy³ | Jingchuan Xue¹ |
Shankar Subramaniam³ | H. Paul Benton¹ | Gary Siuzdak^{1,4} 

¹Scripps Center for Metabolomics, The Scripps Research Institute, La Jolla, California, USA

²Department of Chemistry, University of British Columbia, Vancouver, British Columbia, Canada

³Department of Bioengineering, University of California San Diego, La Jolla, California, USA

⁴Department of Chemistry, Molecular and Computational Biology, The Scripps Research Institute, La Jolla, California, USA

⁵The Jackson Laboratory for Genomic Medicine, Farmington, Connecticut, USA

⁶Systems Immunity Research Institute, Cardiff University, Cardiff, UK

Correspondence

H. Paul Benton and Gary Siuzdak, Scripps Center for Metabolomics, The Scripps Research Institute, 10550 North Torrey Pines Rd., La Jolla, CA, 92037, USA.

Email: hpbenton@scripps.edu and siuzdak@scripps.edu

Funding information

National Institutes of Health, Grant/Award Numbers: R35 GM130385, P30 MH062261, P01 DA026146, U01 CA235493; Ecosystems and Networks Integrated with Genes and Molecular Assemblies

Webpage: <https://masspec.scripps.edu/>

Abstract

Archived metabolomics data represent a broad resource for the scientific community. However, the absence of tools for the meta-analysis of heterogeneous data types makes it challenging to perform direct comparisons in a single and cohesive workflow. Here, we present a framework for the meta-analysis of metabolic pathways and interpretation with proteomic and transcriptomic data. This framework facilitates the comparison of heterogeneous types of metabolomics data from online repositories (eg, XCMS Online, Metabolomics Workbench, GNPS, and MetaboLights) representing tens of thousands of studies, as well as locally acquired data. As a proof of concept, we apply the workflow for the meta-analysis of (a) independent colon cancer studies, further interpreted with proteomics and transcriptomics data, (b) multimodal data from Alzheimer's disease and mild cognitive impairment studies, demonstrating its high-throughput capability for the systems level interpretation of metabolic pathways. Moreover, the platform has been modified for improved knowledge dissemination through a collaboration with Metabolomics Workbench and LIPID MAPS. We envision that this meta-analysis tool combined with our in-source fragmentation/annotation (ISA) technology will help overcome the primary bottleneck in analyzing diverse datasets and facilitate the full exploitation of archival metabolomics data for addressing a broad array of questions in metabolism research and systems biology.

KEYWORDS

archived data, meta-analysis, metabolic pathways, metabolomics, proteomics, systems biology, transcriptomics

1 | INTRODUCTION

Metabolites are the prime drivers of biological activity as they regulate enzyme reactions,¹ protein activation, and gene/protein expression.² Ultimately, metabolites provide an accessible functional readout for the activity of the system and in themselves modulate the phenotype.³ In line with this, the meta-analysis of untargeted high-resolution mass spectrometry (MS) metabolomic data obtained from distinct studies can be used to

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2020 The Authors. *Analytical Science Advances* published by WILEY-VCH Verlag GmbH & Co. KGaA, Weinheim.



obtain a better understanding of the altered metabolic processes and active endogenous metabolites affecting the system over a broad population of samples. This type of analysis requires the generation and/or recollection of multiple metabolomic data sets across several independent studies, to provide a more comprehensive picture than an individual study. In some cases, the data sets required for the meta-analysis have already been generated and made available on public databases. In this regard, several data storage infrastructures have been recently developed to address the raising call for metabolomics data sharing and currently encompass more than 1000 untargeted high-resolution data sets. Emerging open-access ecosystems include MetaboLights,⁴ MetabolomicsWorkbench,⁵ Metabolonote,⁶ Global Natural Products Social Molecular Networking (GNPS),⁷ and metabolomic data aggregation services, such as metabolomeXchange⁸ (<http://www.metabolomexchange.org/site/>) and Omics Discovery Index (<http://www.omicsDI.org>).⁹ In addition, the LIPID MAPS service provides a link into MetabolomicsWorkbench to support the direct deposition of lipidomics data (www.lipidmaps.org).^{10,11}

These publicly available data sets can reduce the workload for data re-collection as well as foster transparency and collaboration between researchers. However, owing to the absence of tools for their cohesive meta-analysis and to the heterogeneity of the stored data, that are often obtained by different types of MS-based metabolome profiling workflows, each study remains only partially utilized for comparative analyses.

Currently, the meta-analysis of metabolomic pathways is carried out by comparing and analyzing the results reported in published papers (eg, fold change comparison, absolute concentrations from targeted studies), thus ignoring the total content of information on metabolites contained in the raw profiling data. Moreover, the interpretation of the meta-analysis findings in the context of proteomic and transcriptomic dysregulations remains a manual task as no systems level data interpretation tool currently provides this functionality. For example, depending on data type, there are many tools for integration of multi-omics data available including correlation analysis, multivariate comparison, regression/machine learning for sample classification.

Here, we report MetaXCMS framework, to enable the meta-analysis of heterogeneous types of archived untargeted, high-resolution MS data across metabolomics, proteomics, transcriptomics, and genomics. The XCMS Online metabolomics platform^{12,13} is an environment for the direct re-analysis, in-source fragmentation/annotation^{47,48} and comparison of data from transcriptomics and metabolomics repositories and/or acquired locally, to gather insights into the dysregulated active metabolites and pathways over independent studies and populations of subjects/samples. We deploy this workflow by integrating and interpreting at systems level archival data sets from two independent colon cancer studies obtained from the XCMS Online Public repository.¹⁴ In addition, we tested this framework in the meta-analysis of archival multimodal metabolomics data acquired from plasma samples from patients with Alzheimer's disease, mild cognitive impairment, and cognitive normal patients, from the MetabolomicsWorkbench.¹⁵

2 | RESULTS

2.1 | Workflow for the meta-analysis of archival metabolomic data and systems level integration

We developed Meta XCMS framework for the meta-analysis and interpretation of archival metabolomics data by developing and combining different bioinformatic modules to be facilitated with the XCMS Online platform (Figure 1).

In the meta-analysis workflow, the raw MS data sets from individual studies can be uploaded in XCMS Online to perform data processing and analysis, including peak detection, retention time alignment, putative annotation, and statistical significance testing, to a final list of detected and dysregulated metabolic features.¹³ At this level, the user can set the processing and statistical parameters depending on the analytical platform employed for metabolome profiling and on the statistical test needed for that study. The processed jobs can then be selected and downloaded to be inputted into the Meta XCMS framework for further analysis. Multiple metabolomic analysis (analytical modalities) for a given study can be combined together for comprehensive coverage, for example, on lipid and central carbon metabolism (eg, combining data obtained by reversed phase chromatography coupled with positive electrospray ionization (ESI)-MS, and hydrophilic interaction liquid chromatography in negative ESI-MS, etc).^{16,17} Moreover, metabolomics data sets obtained from high-resolution untargeted studies archived in the Metabolomics Workbench can be directly uploaded for meta-analysis, and the user interface also supports the upload of metabolomics data sets in text/tsv format, obtained through alternative data preprocessing workflows.

The Meta XCMS framework code is based on mummichog version 2.0.7^{17,18} and leverages itself on this open source platform. Briefly, the mummichog algorithm performs a Fisher's exact test on the number of metabolites jointly dysregulated in the studies as opposed to the total metabolites in the pathway, to predict active pathways directly from putative metabolic features. To allow for a multi-file input, the code and algorithm were adjusted. Using either data from the Metabolomics Workbench or XCMS Online, data are read into the system via convenient tsv/csv peak list file formats. Data read into the system are first parsed, each feature is tagged to its input file to allow for tracing throughout the system. Several possible adducts combinations are calculated onto the feature masses. These features are queried against the pathway database, each file is processed separately and ranked on their corresponding *P*-value score to statistically eliminate false positives. Once a list of possible hits is obtained, features that match the neutral mass in the pathway, they are merged between the different files. This is done such that any compound that is seen in more than one file is merged together and the best *P*-value score is taken. This method allows for the expanded coverage and keeps the statistical

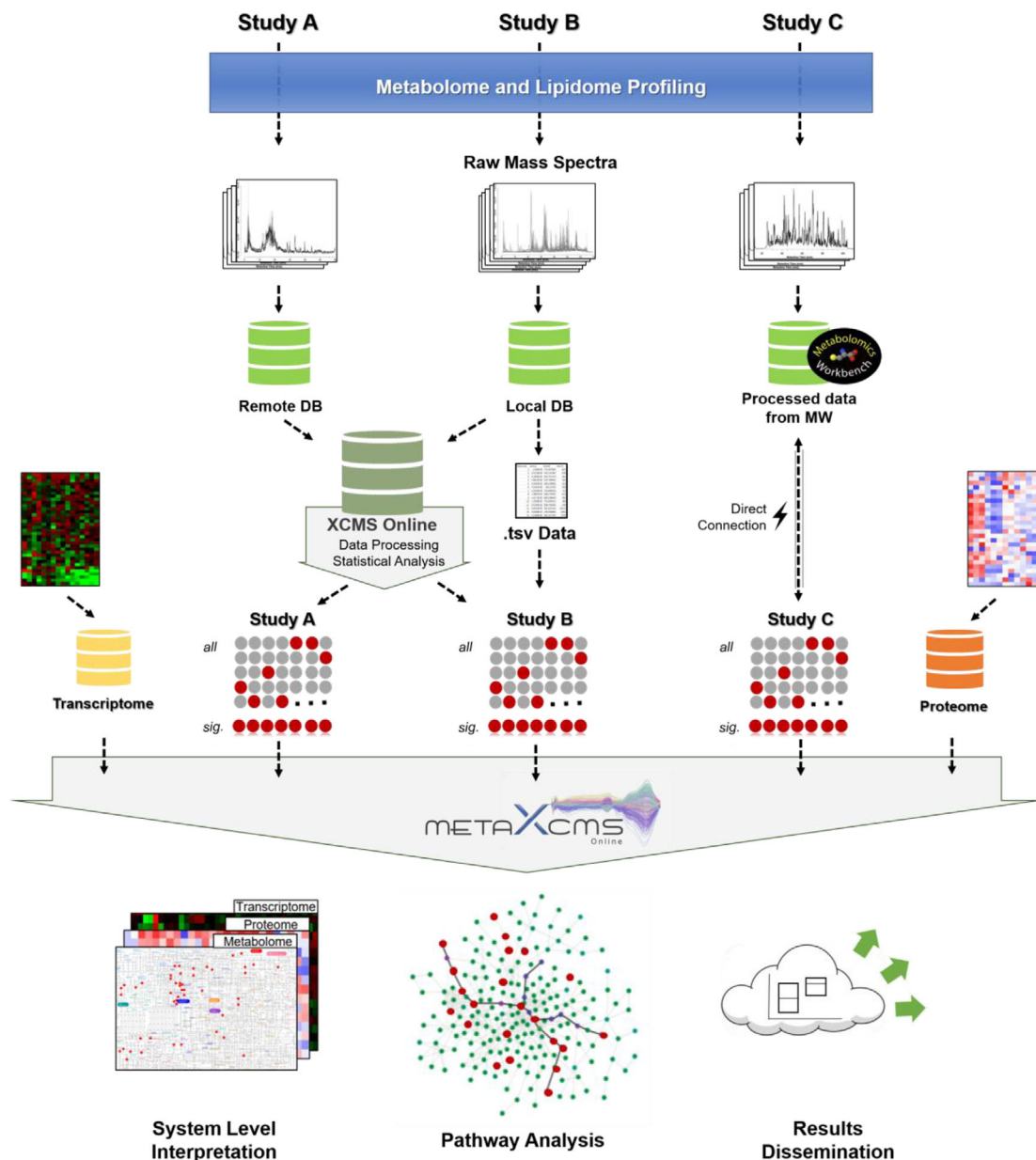


FIGURE 1 Workflow for the meta-analysis of heterogeneous metabolomics data combined with our in-source fragmentation/annotation (ISA) technology^{47,48} in XCMS Online: metabolomics data sets from public repositories are uploaded in the XCMS Online database and processed for metabolic features detection and statistical analysis. The jobs are then used as input for meta-analysis and integration with proteomics and transcriptomics data in Meta XCMS framework. Results can be shared in the XCMS Online cloud for knowledge dissemination

validation. Next, the regular *mummichog* process continues with the statistical validation of the matched pathways. Finally, the output is processed to simplify future analyses.

For each individual study included in the meta-analysis, the user can set a specific significance threshold (P -value), m/z tolerance for metabolite putative annotation, and filter the metabolic feature list according to a specified intensity threshold. We recommend to carefully choose these parameters considering the size of each individual study and the type of metabolomics platform used.

Notably, lipids comprise around a third of all metabolites, but they require distinct processing approaches for accurate annotation and pathway prediction. In particular, removing spurious MS signals is critical to improve statistical power, especially for lipids where multiple forms of the same molecule can be detected/exist. LipidFinder¹⁹ has been recently developed at LIPID MAPS¹¹ to alleviate these artifacts. Here, we suggest using LipidFinder post processing of XCMS outputs. This helps to further broaden the output of lists of putative structures and their categories for more accurate meta-analysis of lipidomics data.

As part of the Meta XCMS framework, we suggest performing multi-omic data integration by superimposing user-uploaded transcriptomic and proteomic data sets onto the dysregulated pathways. Using a list of dysregulated genes (as gene symbols or loci) and proteins (as UniProt accession

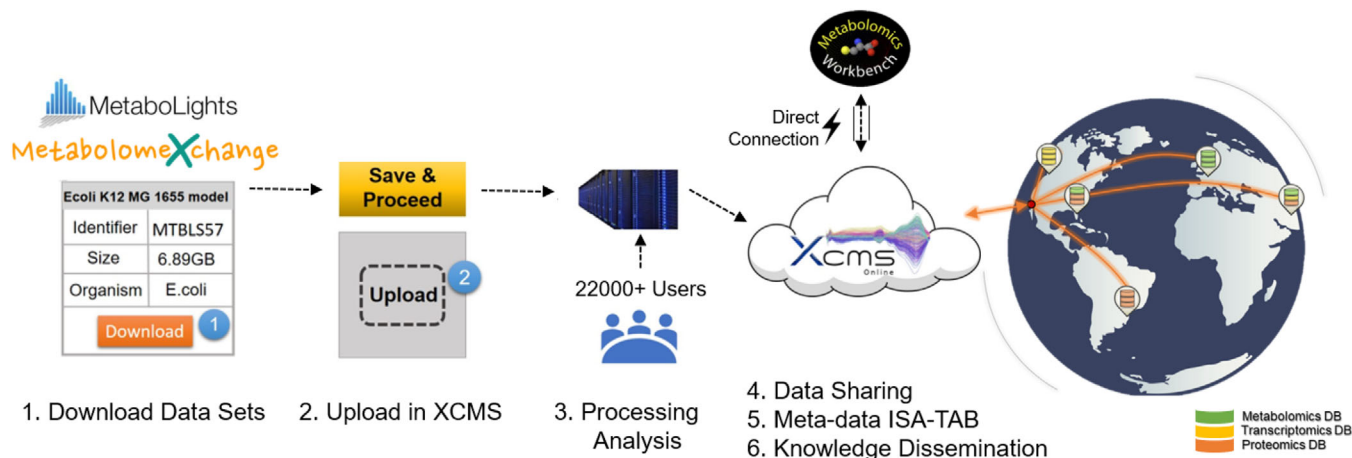


FIGURE 2 Developments of the meta-XCMS framework to enhance archival metabolomics data processing, archiving and sharing for meta-analysis and systems level interpretation

IDs or gene symbols) obtained from studies targeting a given biological question users can generate improved confidence on the pathway hits. The integration with proteomics and transcriptomics results offers the possibility to gauge a systems level mechanistic understanding of pathways dysregulation and metabolite activity in the investigated biological system.²⁰

The downloaded “metabolite results” table reports the list of all the overlapping dysregulated metabolites detected in the studies included in the meta-analysis and used for pathway prediction.

The “pathway results” table showcases the output for the metabolic pathways jointly dysregulated in the studies. For each metabolic pathway, it reports the number of overlapping dysregulated metabolites detected in different studies with respect to the total number of metabolites in the pathway. By clicking on the number of “shared metabolites” the complete list of metabolic features is shown. Entries can be further filtered based on the adduct type, study group or pairwise job group. Moreover, this table reports the overlapping dysregulated genes and proteins from the uploaded proteomics and transcriptomics data for each metabolic pathway, thus providing a rapid glance on the biological process from a system-wide perspective.

2.2 | Expanding the capability of XCMS Online for the meta-analysis of archived data

To allow the meta-analysis of metabolomics data sets obtained from disparate sources, we enabled easy parsing file options of tsv/csv in Meta-XCMS framework (Figures 1 and 2). First, the user can select studies processed in the XCMS Online private space or in the XCMS Online Public. These data are then easily downloaded to be parsed into the Meta-XCMS framework. Alternatively, the Metabolomics Workbench data can also be used. On many of the studies there are already existing outputs of identified metabolites or feature lists. In these instances where a metabolite is already identified it will be read into the system and used as a confirmed metabolite of the pathway analysis.

This strategy is aimed at fostering data dissemination and at actively promoting the full exploitation of archived metabolomics data through stimulating further meta-analysis for results validation or for generating novel hypothesis.

2.3 | Analysis of archived metabolomics data

We tested the workflow in the systems level meta-analysis of two independent colon cancer tissue metabolomics data sets by leveraging archival data from the XCMS Online Public repository, and for the meta-analysis of Alzheimer’s disease and mild cognitive impairment studies in plasma, including heterogeneous profiling data from the Metabolomics Workbench.

2.3.1 | Colon cancer

Several previous studies have pinpointed the multifaced metabolic reprogramming underlining colon cancer.^{14,24,25} Here, we performed the meta-analysis of archived untargeted metabolomics data from a study investigating the role of bacterial biofilms in colon cancer¹⁴ (Study A) and a second colon cancer study, recently performed in our laboratory (Study B) (Figure 3a). Study A involved 30 subjects diagnosed with colon cancer from stage 3 to 4 (18 females and 12 males, 61-88 years old) and was available in the XCMS Online Public repository,¹⁴ while Study B involved 19 subjects

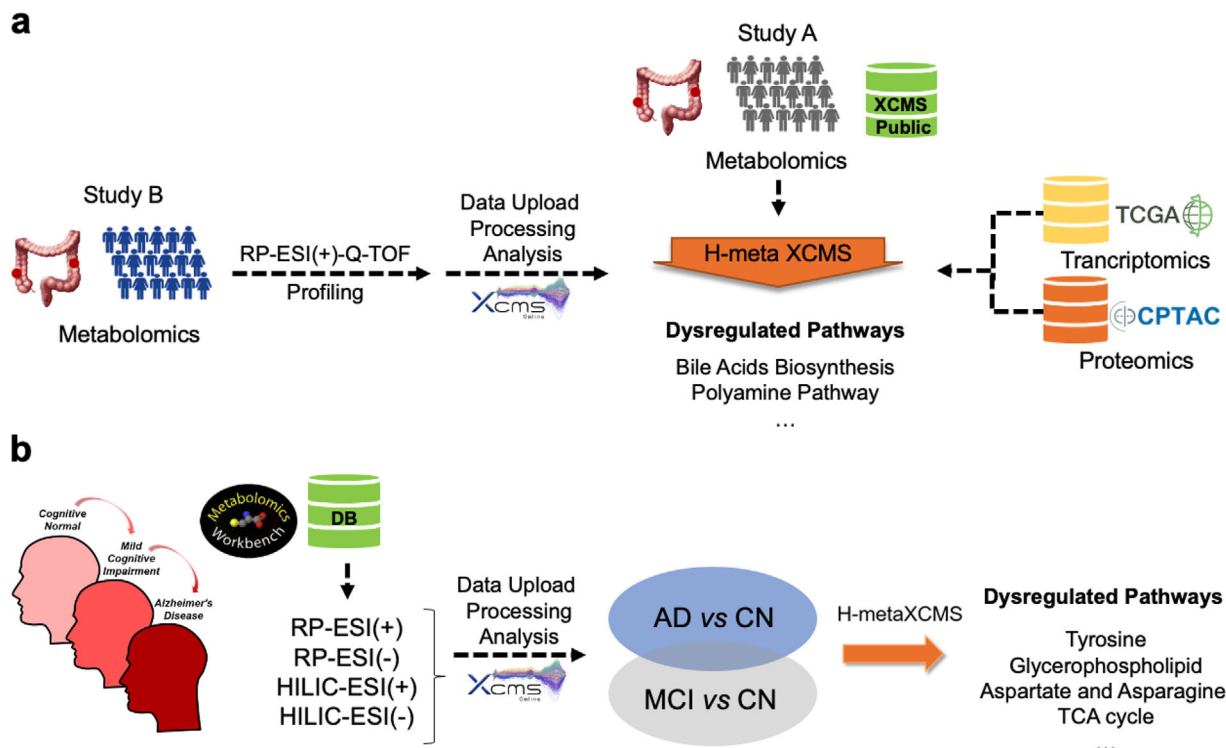


FIGURE 3 a, Meta-analysis of colon cancer metabolomics studies, and systems level interpretation with proteomics and transcriptomics data; b, cohesive re-analysis of heterogeneous archival metabolomics data from AD and MCI, compared with CN

diagnosed with colon cancer (13 females and 6 males, 62-92 years old). More details on study design and samples collection are available in Table S1. Both studies used similar platforms for metabolome profiling (reversed phase chromatography coupled with ESI(+)-quadrupole time-of-flight (Q-TOF) MS) and sample preparation protocols, therefore we expected comparable metabolome coverage and overlapping dysregulations (excluding inter-population heterogeneity).

First, raw data from study B were uploaded in XCMS Online and processed as a paired job. This job and the archival job (Study A) were then selected as input for meta-analysis in meta XCMS framework (Figure 3a). The results unveiled the presence of 30 metabolic pathways with at least 10 dysregulated putative metabolic features across study A and B (Table S2; Figure 4a), among these are glycerophospholipids metabolism, aspartate and asparagine metabolism, glycine, serine and alanine metabolism, carnitine shuttle, tyrosine metabolism, steroidal hormones, and bile acids. The dysregulation of the glycerophospholipid metabolism has been previously confirmed correlating with altered viability, proliferation, and colorectal cancer development.²⁶ The meta-analysis also highlighted the dysregulation of the aspartate and asparagine pathway that includes spermine/spermidine biosynthesis and degradation (polyamine pathway) where N1-acetylspermidine, N1-acetylspermine, spermidine, and N1,N12-diacetylspermidine, spermidine dialdehyde, spermic acid were found jointly upregulated, a finding consistent with previous work (Figure 4c).^{14,30} Of note, in the bile acids biosynthesis pathway taurine and taurochenodeoxycholate were upregulated (Figure 4b). Increased levels of conjugated bile acids have been previously reported to highly associate with colon cancer.²⁷ In particular, taurochenodeoxycholate can be hydrolyzed releasing taurine, a sulfur amino acid further transformed by the gut microbiota to form compounds with genotoxic activity (eg, H₂S), and colon tumor promoters (deoxycholic acid).^{28,29}

To Interpret this evidence in light of the variations occurring at proteomic and transcriptomic level, colon cancer data sets obtained from The Cancer Genome Atlas²⁴ and The CPTAC Proteomics Data Portal³² were uploaded and processed in meta XCMS framework. Approximately 90% of the upregulated metabolic pathways were further supported by dysregulated proteins and gene transcripts (Table S3). For instance, both polyamines and bile acids pathway dysregulations were confirmed (Figure 3a).

2.3.2 | Alzheimer's disease

Alzheimer's disease (AD) is a progressive neurodegenerative disorder of unknown etiology.³³ AD and dementia patients are usually subject to a long pre-AD period known as mild cognitive impairment (MCI). Here, we used public available metabolomics data obtained from previous longitudinal studies performed at the Mayo Clinic Study of Aging (MCSA) and Mayo Clinic Alzheimer Disease Research Center (ADRC).¹⁵ Plasma samples were from AD, MCI, and cognitive normal (CN) subjects (15 individuals/group). The metabolomics data sets and meta-data were publicly available in the

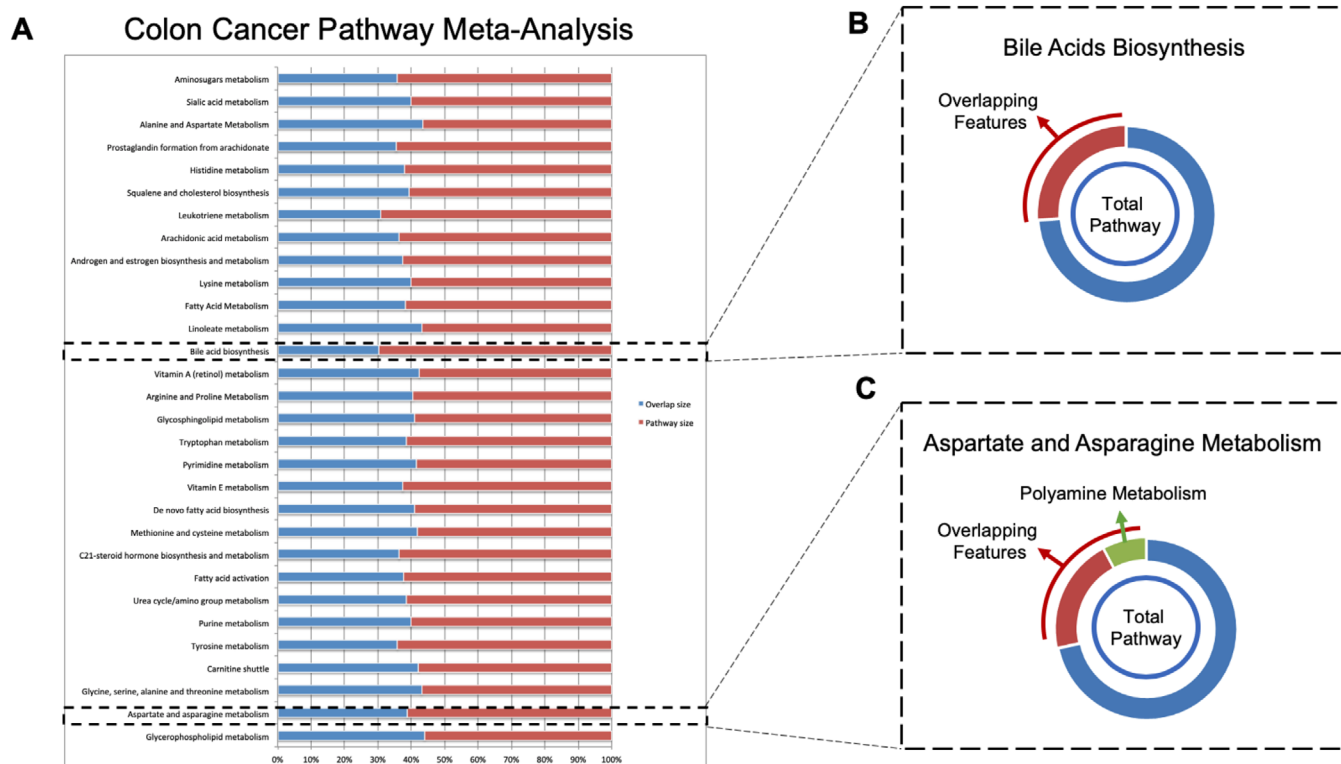


FIGURE 4 Meta-analysis across independent colon-cancer studies predicts 30 dysregulated metabolic pathways (a), including bile acids biosynthesis (b) and polyamine metabolism (c)

Metabolomics Workbench and formerly generated by LC-Q-TOF MS in four analytical modalities (hydrophilic interaction liquid chromatography [HILIC] and reversed phase liquid chromatography [RP], in both positive and negative ESI modes) for comprehensive metabolome coverage. We downloaded the raw data sets from the Metabolomics Workbench repository and uploaded them in the XCMS Online for processing and statistical analysis to extract significant metabolic variations in AD versus CN and MCI versus CN. The resulting XCMS Online jobs were then used as input for the Meta XCMS framework to detect shared metabolic changes at different disease stages.

Meta XCMS framework predicted 24 dysregulated metabolic pathways with at least 10 metabolic features dysregulated in the AD versus CN and the MCI versus CN groups, over a total of 101 paths (Table S4). We manually compared the output with the dysregulated pathways reported in the original publication.¹⁵ Our development predicted the dysregulation of tyrosine, glycerophospholipid, aspartate and asparagine, glycine, serine and alanine metabolism, urea cycle, tryptophan, and purine metabolism, together with the other pathways reported in Table S4. The original work reported 50 total dysregulated pathways, of which nine were consistently predicted across AD versus CN and MCI versus CN (Figure 5). Our approach predicted a total of 101 dysregulated pathways, of which 24 pathways were reported in the original publication, demonstrating the efficiency of the workflow in identifying jointly dysregulated metabolic pathways from heterogeneous archived metabolomic data.

3 | DISCUSSION

Archived metabolomics data are a rich source of information for second-order analysis by the scientific community. However, the heterogeneity of the data and the lack of tools for their cohesive re-analysis and interpretation hinders their full utilization. To address this, we developed a framework for archival data re-processing, analysis, integration, and interpretation at systems level. The workflow moves from the metabolomics tools available in the XCMS Online, further combining them with a bioinformatic development specifically designed for the meta-analysis of heterogeneous metabolomics data.

A key aspect of the workflow is the use of a pathway-centric approach to the meta-analysis, which allows the direct prediction of the dysregulated metabolic pathways from putative metabolic features jointly detected in the archived data/studies. This is performed through the embedment of a recently developed tool for metabolic pathway prediction from putative annotations of metabolic features extracted from different types of metabolomics data sets.^{17,18} This tool allows for higher confidence in the putative pathway enrichment results by estimating the probability of a pathway being dysregulated on the basis of the total number of dysregulated metabolites detected. It is also worth noting that, when attempting



MCI and AD Pathway Meta-Analysis

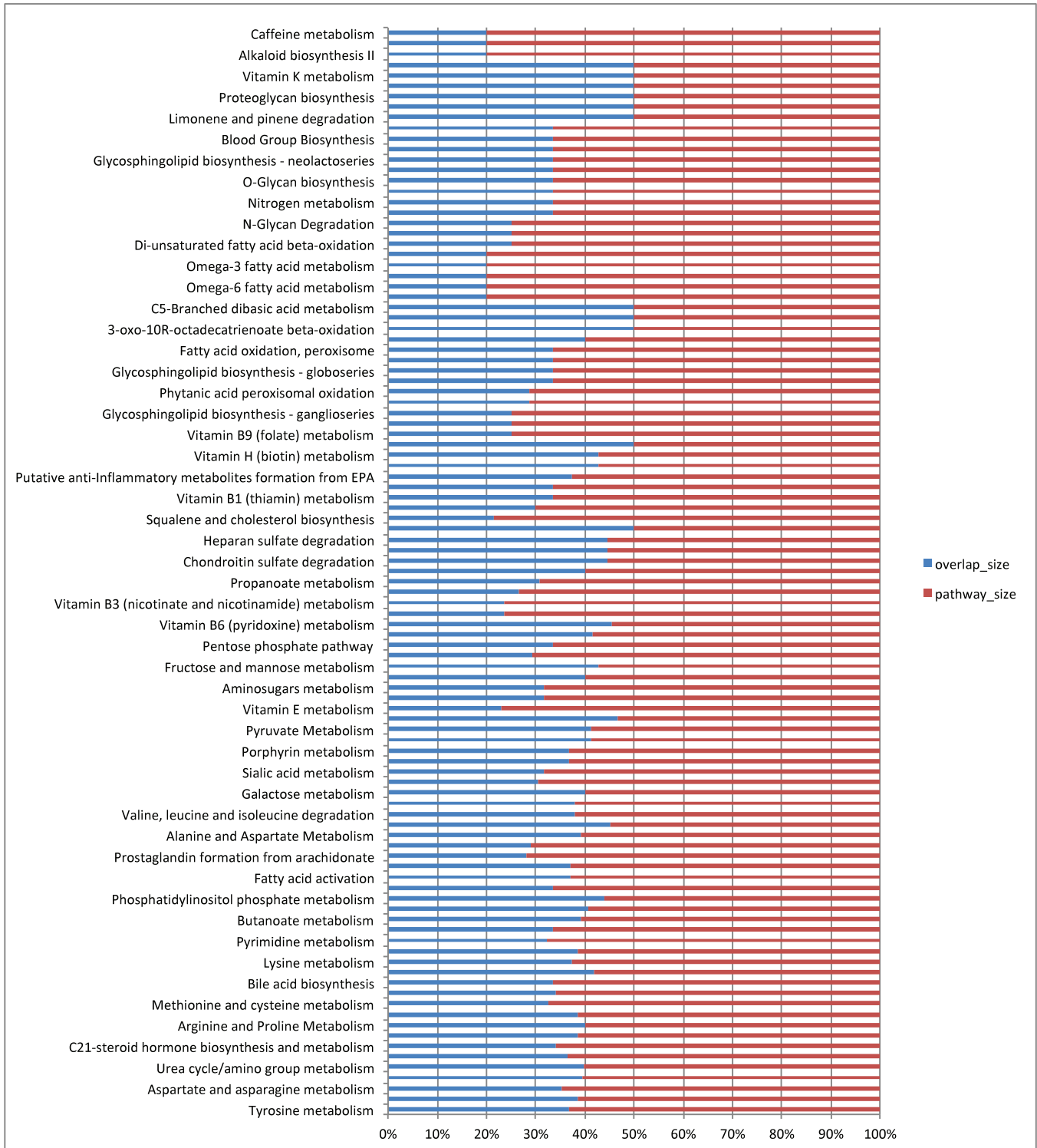


FIGURE 5 Meta-analysis across AD versus CN and MCI versus CN studies in plasma predicts 101 dysregulated metabolic pathways



to re-analyze archival metabolomics data, the physical samples are not directly accessible and often no longer available. This makes it unfeasible to perform further MS fragmentation experiments and metabolite identity confirmation. In this scenario, performing direct pathway prediction analysis represents a practical strategy to bypass this limitation and directly formulate biological hypothesis from the archived data, to be later tested through independent targeted studies or biochemical assays.

The use of a pathway-centric approach also streamlines the interpretation of the metabolomic data by superimposing the dysregulated proteins and transcripts to each metabolic pathway. This provides a rapid glance on the system in the light of other omics regulatory levels, introducing the possibility for the orthogonal confirmation of the insights extracted from the archived metabolomics data.

Despite recent efforts aimed at standardizing metabolome profiling and reporting,^{21,35,36} a widely adopted consensus on untargeted MS-based workflows in the perspective of meta-analysis is still missing.^{37,38} The metabolic profiles are indeed usually acquired by a variety of different analytical solutions,^{16,39-41} thus introducing heterogeneity in the type of available data. For example, Metabolomics Workbench⁵ currently stores >190 untargeted high-resolution MS studies, for a total of ~400 different analyses (ie, different analytical modalities including ESI positive and negative acquisition modes), while Metabolights⁴ stores ~350 among GC- and LC-MS based studies. This heterogeneity complicates the development of bioinformatics solutions for the automated meta-analysis, as each different metabolomics platform calls for specific data processing and metabolite annotation pipelines. To circumvent this limitation, we designed the workflow in a modular fashion: each study can be processed as pairwise XCMS Online job using different processing and statistical settings and the resulting jobs can be used as input for further meta-analysis. This allows the use of raw data acquired by different metabolomics platforms and modalities. For example, the user can upload both lipidomics and metabolomics data obtained by different chromatographic or ionization modes for comprehensive metabolome coverage and improved pathway prediction.¹⁷ This, together with the ability to perform multi-omics data integration, represents a fundamental advantage over our previous development for the meta-analysis of metabolomics data.^{42,43}

Besides raw data sets, the workflow also supports the direct comparison of untargeted studies already processed and available in the XCMS Online Public cloud and in the Metabolomics Workbench. Of note, several data sets currently archived in public databases are not compliant with the ISA guidelines for meta-data reporting, unearthing the need of harmonized and more pragmatic guidelines for metabolomics data sharing in public repositories.³⁷

The meta-analysis of metabolomics data sets and their interpretation at systems level has the potential of streamlining different types of study comparisons. For example, a meta-analytical approach can be used for (a) providing further validation of metabolite dysregulations in the context of independent set of samples (eg, in biomarker studies); (b) stimulating the generation of biological hypothesis from the re-analysis of archived untargeted studies; (c) streamlining the exclusion of experimental artifacts to reduce the list of dysregulated metabolites before performing time consuming structural elucidation⁴²; (d) excluding metabolic dysregulations due to physiologic heterogeneity in different populations of subjects, therefore taking a step towards the identification of therapeutic targets and biomarkers of broad applicability. In particular, in biomarker discovery the automatic integration of multiple archival metabolomics studies can be a cost-effective strategy to minimize the interstudy bias introduced by genetic and environmental factors.⁴⁴⁻⁴⁶ This strategy is not limited to archived data, since the difficulty in cross-laboratory comparison has impeded the biomedical applications of metabolomics. In an emergency situation like the current COVID-19 pandemic, the pathway-centric meta-analysis can be important for identifying scientific consensus in a timely manner.

As proof of concept, we demonstrated the utility of the workflow in two meta-analytical studies. First, we leveraged archival data sets from a previous study available in the XCMS Online Public,¹⁴ for the autonomous comparison with a colon cancer study recently performed in our laboratory. The workflow allowed a rapid glance on metabolic pathways jointly dysregulated and validation at systems level (eg, the bile acid and the polyamine pathways). In the second example, we applied the workflow for the re-analysis of archival biomarker studies obtained from the Metabolomics Workbench database. The workflow permits the streamlined and autonomous prediction of metabolic pathways changed in both AD and MCI patients in plasma (pre-AD) in agreement and beyond the results previously obtained by manual meta-analysis.¹⁵

In summary, there are many challenges in the analysis of diverse datasets including variability in experimental designs as well as information types that are largely platform dependent. However, by combining a fully automated workflow including in-source fragmentation/annotation^{47,48} with an enhanced strategy for data storage and direct connection to the Metabolomics Workbench data repository, the described approach can provide a solution for meta-analysis, with the ultimate goal of maximizing the usage and dissemination of information-rich archival metabolomics data. With the growing number of metabolomics, proteomics and genomics data generated to cover a wide range of biological questions, this workflow paves the way to unlock biological insights in the era of “big data” and “open science.”

4 | MATERIALS AND METHODS

4.1 | Meta-analysis framework

The system has been built on a local based Flask system with code based on mummichog version 2.0.7 running on python 2.7.¹⁸ Several files were altered to allow for a multi-file input and hosting on a web frontend. Metabolomics Workbench data is processed using either csv/tsv formats



directly into the software or via XCMS processing to csv/tsv files. Once the data are read each feature is processed against all possible adducts for search masses. Using the mummichog algorithm, the search masses are searched against the pathway database, each file is processed separately and ranked on their corresponding *P*-value score. Once a list of possible hits is obtained, they are merged between the different files such that any compound that is seen in more than one file is merged together and the best *P*-value score is taken. Now, the regular mummichog process continues with the statistical validation of the altered pathways. Finally, the output is processed to make further analysis simpler and result are downloaded by the user. Framework code has been made available via github at <https://github.com/hpbenton/archive-mummi>

4.2 | Metabolomics profiling data

The colon cancer study A was available in the XCMS Online Public Space as a processed job. The raw MS data from high-resolution metabolome profiling for the colon cancer study B were archived in our laboratory, and previously obtained as part of a pilot colon cancer study (Study B). Both profiling studies were performed in RP-ESI(+)-Q-TOF profiling. More details on study A and B experimental setups can be found in previous published work.^{14,17} Metabolome profiles for AD, MCI, and CN plasma samples were downloaded from the Metabolomics Workbench repository,⁵ uploaded and re-processed in the XCMS Online. These studies were performed at the Mayo Clinic Study of Aging (MCSA) and Mayo Clinic Alzheimer Disease Research Center (ADRC) and more details on study design and experimental procedures are available in previous published work.¹⁵

4.3 | Data processing and re-analysis

Raw archival data sets were uploaded as .mzXML files in the XCMS Online and processed as pairwise jobs. Before processing, the profiling data were manually examined for assessing the quality and the parameters for further processing and analysis. Colon cancer study A was already processed and available in the XCMS Online Public. The XCMS jobs were used for further meta-analysis in the Meta XCMS framework. *P*-value, intensity, and ppm error settings are reported in Supporting Information Text S1.

4.4 | Proteomics and transcriptomics data

Transcriptomics data were obtained from The Cancer Genome Atlas (TCGA)³¹ in the frame of a previous colon cancer study involving 22 subjects (22 colon cancer tissue samples vs 22 paired normal tissues). Dysregulated genes were selected based on a *P*-value cut off of 0.01 and fold change > 4. A total of 7138 dysregulated transcripts were included in the final data set for upload in the XCMS Online as gene symbols. Proteomics data were obtained by the Clinical Proteomic Tumor Analysis Consortium (CPTAC), involving 90 patients affected by colon cancer (90 colon cancer tissue samples and 90 paired normal tissues). Dysregulated proteins were filtered by *P*-value < .01 and fold change > 2, obtaining a total of 2545 dysregulated proteins uploaded in the XCMS Online as UniProt accession IDs for multiomic analysis.

ACKNOWLEDGMENTS

This research was partially funded by National Institutes of Health grants R35 GM130385, P30 MH062261, P01 DA026146 and U01 CA235493; and by Ecosystems and Networks Integrated with Genes and Molecular Assemblies (ENIGMA), a Scientific Focus Area Program at Lawrence Berkeley National Laboratory for the U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research, under contract number DE-AC02-05CH11231.

CONFLICT OF INTEREST

The authors declare no conflict of interest.

DATA AVAILABILITY

As referred to in the text, the data is available through multiple sources.

ORCID

Gary Siuzdak  <https://orcid.org/0000-0002-4749-0014>

REFERENCES

1. Hackett SR, Zanotelli VRT, Xu W, et al. Systems-level analysis of mechanisms regulating yeast metabolic flux. *Science*. 2016;354:aaf2786.
2. Johnson CH, Ivanisevic J, Siuzdak G. Metabolomics: beyond biomarkers and towards mechanisms. *Nat Rev Mol Cell Biol*. 2016;17:451-459.



3. Joyce AR, Palsson BO. The model organism as a system: integrating 'omics' data sets. *Nat Rev Mol Cell Biol.* 2006;7:198-210.
4. Haug K, Salek RM, Conesa P, et al. MetaboLights - An open-access general-purpose repository for metabolomics studies and associated meta-data. *Nucleic Acids Res.* 2013;41:D781-D786.
5. Sud M, Fahy E, Cotter D, et al. Metabolomics Workbench: an international repository for metabolomics data and metadata, metabolite standards, protocols, tutorials and training, and analysis tools. *Nucleic Acids Res.* 2016;44:D463-D470.
6. Ara T, Enomoto M, Arita M, et al. Metablonote: a Wiki-based database for managing hierarchical metadata of metabolome analyses. *Front Bioeng Biotechnol.* 2015;3:38.
7. Wang M, Carver JJ, Phelan VV, et al. Sharing and community curation of mass spectrometry data with global natural products social molecular networking. *Nat Biotechnol.* 2016;34:828-837.
8. Haug K, Salek RM, Steinbeck C. Global open data management in metabolomics. *Curr Opin Chem Biol.* 2017;36:58-63.
9. Perez-Riverol Y, Bai M, da Veiga Leprevost F, et al. Discovering and linking public omics data sets using the Omics Discovery Index. *Nat Biotechnol.* 2017;35:406-409.
10. O'Donnell VB, Dennis EA, Wakelam MJO, Subramaniam S. LIPID MAPS: serving the next generation of lipid researchers with tools, resources, data, and training. *Sci Signal.* 2019;12:eaaw2964.
11. Fahy E, Alvarez-Jarreta J, Brasher CJ, et al. LipidFinder on LIPID MAPS: peak filtering, MS searching and statistical analysis for lipidomics. *Bioinformatics.* 2019;35:685-687.
12. Tautenhahn R, Patti GJ, Rinehart D, Siuzdak G. XCMS online: a web-based platform to process untargeted metabolomic data. *Anal Chem.* 2012;84:5035-5039.
13. Forsberg EM, Huan T, Rinehart D, et al. Data processing, multi-omic pathway mapping, and metabolite activity analysis using XCMS Online. *Nat Protoc.* 2018;13:633-651.
14. Johnson CH, Dejea CM, Edler D, et al. Metabolism links bacterial biofilms and colon carcinogenesis. *Cell Metab.* 2015;21:891-897.
15. Trushina E, Dutta T, Persson XMT, Mielke MM, Petersen RC. Identification of altered metabolic pathways in plasma and CSF in mild cognitive impairment and Alzheimer's disease using metabolomics. *PLoS One.* 2013;8:e63644.
16. Ivanisevic J, Zhu Z-J, Plate L, et al. Toward omic scale metabolite profiling: a dual separation mass spectrometry approach for coverage of lipids and central carbon metabolism. *Anal Chem.* 2013;85:6876-6884.
17. Huan T, Palermo A, Ivanisevic J, et al. Autonomous multimodal metabolomics data integration for comprehensive pathway analysis and systems biology. *Anal Chem.* 2018;90:8396-8403.
18. Li S, Park Y, Duraisingham S, et al. Predicting network activity from high throughput metabolomics. *PLoS Comput Biol.* 2013;9:e1003123.
19. O'Connor A, Brasher CJ, Slatter DA, et al. LipidFinder: a computational workflow for discovery of lipids identifies eicosanoid-phosphoinositides in platelets. *JCI Insight.* 2017;2:e91634.
20. Guijas C, Montenegro-Burke JR, Warth B, Spilker M E, Siuzdak G. Metabolomics activity screening for identifying metabolites that modulate phenotype. *Nat Biotechnol.* 2018;36:316-320.
21. Rocca-Serra P, Salek RM, Arita M, et al. Data standards can boost metabolomics research, and if there is a will, there is a way. *Metabolomics.* 2016;12:14.
22. Côté R, Reisinger F, Martens L, et al. The ontology lookup service: bigger and better. *Nucleic Acids Res.* 2010;38:W155-W160.
23. Larralde M, Lawson TN, Weber RJM, et al. mzML2ISA & nmrML2ISA: generating enriched ISA-Tab metadata files from metabolomics XML data. *Bioinformatics.* 2017;33(16):2598-2600.
24. Brown DG, Rao S, Weir TL, et al. Metabolomics and metabolic pathway networks from human colorectal cancers, adjacent mucosa, and stool. *Cancer Metab.* 2016;4:11.
25. Denkert C, Budczies J, Weichert W, et al. Metabolite profiling of human colon carcinoma-deregulation of TCA cycle and amino acid turnover. *Mol Cancer.* 2008;7:72.
26. Yan G, Li L, Zhu B, Li Y. Lipidome in colorectal cancer. *Oncotarget.* 2016;7:33429-33439.
27. Degirolamo C, Modica S, Palasciano G, Moschetta A. Bile acids and colon cancer: solving the puzzle with nuclear receptors. *Trends Mol Med.* 2011;17:564-572.
28. Ridlon JM, Wolf PG, Gaskins HR. Taurocholic acid metabolism by gut microbes and colon cancer. *Gut Microbes.* 2016;7:201-215.
29. Devkota S, Wang Y, Musch MW, et al. Dietary-fat-induced taurocholic acid promotes pathobiont expansion and colitis in Il10^{-/-} mice. *Nature.* 2012;4877405:104-108.
30. Wallace HM, Caslake R. Polyamines and colon cancer. *Eur J Gastroenterol Hepatol.* 2001;13:1033-1039.
31. Tomczak K, Czerwińska P, Wiznerowicz M. The Cancer Genome Atlas (TCGA): an immeasurable source of knowledge. *Wspolczesna Onkol.* 2015;19:A68-A77.
32. Edwards NJ, Oberti M, Thangudu RR, et al. The CPTAC data portal: a resource for cancer proteomics research. *J Proteome Res.* 2015;14:2707-2713.
33. Selkoe DJ. Alzheimer's disease: genes, proteins, and therapy. *Physiol Rev.* 2001;81:741-766.
34. Caspi R, Billington R, Fulcher CA, et al. The MetaCyc database of metabolic pathways and enzymes. *Nucleic Acids Res.* 2018;44:D471-D480.
35. Fiehn O, Amberg A, Barrett D, et al. The metabolomics standards initiative (MSI). *Metabolomics.* 2007;3:175-178.
36. Peng J, Chen YT, Chen CL, Li L. Development of a universal metabolome-standard method for long-term LC-MS metabolome profiling and its application for bladder cancer urine-metabolite- biomarker discovery. *Anal Chem.* 2014;86:6540-6547.
37. Spicer RA, Salek R, Steinbeck C. Comment: a decade after the metabolomics standards initiative it's time for a revision. *Scientific Data.* 2017. <https://doi.org/10.1038/sdata.2017.138>.
38. Spicer RA, Salek R, Steinbeck C. Compliance with minimum information guidelines in public metabolomics repositories. *Sci Data.* 2017;4:170137.
39. Martin JC, Maillot M, Mazerolles G, et al. Can we trust untargeted metabolomics? Results of the metabo-ring initiative, a large-scale, multi-instrument inter-laboratory study. *Metabolomics.* 2015;11:807-821.
40. Cajka T, Smilowitz JT, Fiehn O. Validating quantitative untargeted lipidomics across nine liquid chromatography-high-resolution mass spectrometry platforms. *Anal Chem.* 2017;89:12360-12368.
41. Telu KH, Yan X, Wallace WE, Stein SE, Simón-Manso Y. Analysis of human plasma metabolites across different liquid chromatography/mass spectrometry platforms: cross-platform transferable chemical signatures. *Rapid Commun Mass Spectrom.* 2016;30:581-593.
42. Patti GJ, Tautenhahn R, Siuzdak G. Meta-analysis of untargeted metabolomic data from multiple profiling experiments. *Nat Protoc.* 2012;7:508-516.



43. Cajka T, Fiehn O. Comprehensive analysis of lipids in biological systems by liquid chromatography-mass spectrometry. *Trends Analyt Chem.* 2014;61:192-206.
44. Drucker E, Krapfenbauer K. Pitfalls and limitations in translation from biomarker discovery to clinical utility in predictive and personalised medicine. *EPMA J.* 2013;4:7.
45. Mamas M, Dunn WB, Neyses L, Goodacre R. The role of metabolites and metabolomics in clinically applicable biomarkers of disease. *Arch Toxicol.* 2011;85:5-17.
46. Xia J, Broadhurst DI, Wilson M, Wishart DS. Translational biomarker discovery in clinical metabolomics: an introductory tutorial. *Metabolomics.* 2013;9:280-299.
47. Xue J, Domingo-Almenara X, Guijas C, Palermo A, Rinschen MM, Isbell J, Benton HP, Siuzdak G. Enhanced in-source fragmentation annotation enables novel data independent acquisition and autonomous METLIN molecular identification. *Anal Chem.* 2020;92(8):6051-6059.
48. Domingo-Almenara X, Montenegro-Burke JR, Guijas C, Majumder EL-W, Benton HP, Siuzdak G. Autonomous METLIN-guided in-source fragment annotation for untargeted metabolomics. *Anal Chem.* 2019;91(5):3246-3253.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

How to cite this article: Palermo A, Huan T, Rinehart D, et al. Cloud-based archived metabolomics data: A resource for in-source fragmentation/annotation, meta-analysis and systems biology. *Anal Sci Adv.* 2020;1:70–80. <https://doi.org/10.1002/ansa.202000042>