

High resolution mass spectrometry to investigate the human exposome: Where are We?

Begoña Talavera Andújar, PhD^{*1} , Emma L. Schymanski, Dr. rer. nat.^{*1} 

¹Luxembourg Centre for Systems Biomedicine (LCSB), University of Luxembourg, Belvaux, Luxembourg

*Corresponding authors: Begoña Talavera Andújar, PhD, Luxembourg Centre for Systems Biomedicine (LCSB), University of Luxembourg, Belvaux, Luxembourg (begona.talavera@uni.lu) and Emma L. Schymanski, Dr. rer. nat., Luxembourg Centre for Systems Biomedicine (LCSB), University of Luxembourg, Belvaux, Luxembourg (emma.schymanski@uni.lu)

Abstract

Despite the significant role of the environment in health and disease, the accurate assessment of environmental exposures remains underdeveloped compared to genetic factors. To address this, the concept of the exposome was introduced in 2005 as a complement to the genome. High-resolution mass-spectrometry (HRMS) has emerged as a key technology for the comprehensive assessment of the chemical exposome. However, non-target HRMS based exposomics still faces numerous challenges, with the majority of features detected by HRMS (the “dark matter” of the chemical exposome) remaining unannotated. The lack of standardized workflows across the field often results in poorly comparable studies. Nevertheless, many positive developments have arisen in recent years, with open data revealing interesting trends in HRMS coverage. This review will examine and discuss the entire non-target HRMS exposomics workflow, from experimental design (study design and sample preparation) to the computational analysis and biological interpretation. It will also delve into key concepts, including the sometimes blurred distinction between the metabolome and the chemical exposome and the importance of exposomics within the “omics cascade”. Visualizations are used to support this discussion, including a detailed look at the chemical coverage of key categories of open exposomics resources. The review ends by exploring current challenges and strategies to advance towards harmonized exposomics studies, which are essential for greater biological insights and personalized medicine goals.

Key words: exposomics; non-target screening; high-resolution mass-spectrometry (HRMS); liquid chromatography (LC); cheminformatics.

Introduction

The phenotype of an individual arises from the interplay between genes and environment. Since only a small proportion of chronic diseases can be attributed solely to genetic factors, it is now hypothesized that the majority (70%–90%) are influenced by environmental factors, many of which remain unknown.^{1,2} Despite the significant role of the environment in health and disease, the accurate assessment of many environmental exposures remains underdeveloped compared to genetic factors.² Christopher Wild proposed the concept of the exposome in 2005 as a complement of the genome, defining it as “all life-course exposures (including risk factors), from the prenatal period onwards”.³ This concept was extended in 2014 by Miller and Jones to: “the cumulative measure of environmental influences and associated biological responses throughout the lifespan, including exposures from the environment, diet, behavior, and endogenous processes”.⁴ Recent collaborative discussions in 2024 and 2025 resulted in further revised definitions of the exposome and its research goals, to foster a common understanding in the field.^{5–7}

Figure 1 shows the three different research domains that have been described within the exposome.⁸ The *internal exposome* comprises the internal biological processes as a result of an

exposure such as oxidative stress, metabolism, and microbiome changes.^{8–10} The internal chemical exposome also includes environmentally derived chemicals plus their transformation products found in cells, tissues, organs or organisms.² The *general external exposome* includes social, economic, and environmental factors, while the *specific external exposome* encompasses the individual’s immediate local environment (such as diet, alcohol, infectious agents, pollutants).^{8–10}

Unlike the genome, which remains relatively stable over time, the exposome varies on different timescales, requiring complex study designs.² The integration of different “omics” layers, known as multi-omics analysis, can offer a more comprehensive understanding of a biological system’s state.^{11,12} Among these “omics” (see Figure 2A), metabolomics and exposomics emerge as the terminal downstream outcomes of the genome, reflecting the phenotype of a cell, tissue or organism, in response to diverse genetic or environmental influences over life.^{13,14} Metabolomics and exposomics aim to investigate small molecules, typically ranging between 50 and 1200 Da,² although this boundary can also be constrained by analytical reality such as instrument range of *eg*, $m/z = 1000$ or 2000 .¹⁵ Given the shared analytical scope, metabolomics and exposomics can be conceptually grouped as “small molecule omics.” Some of these small molecules are closely

Received: September 28, 2025; Revised: December 04, 2025; Accepted: December 10, 2025

© The Author(s) 2025. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

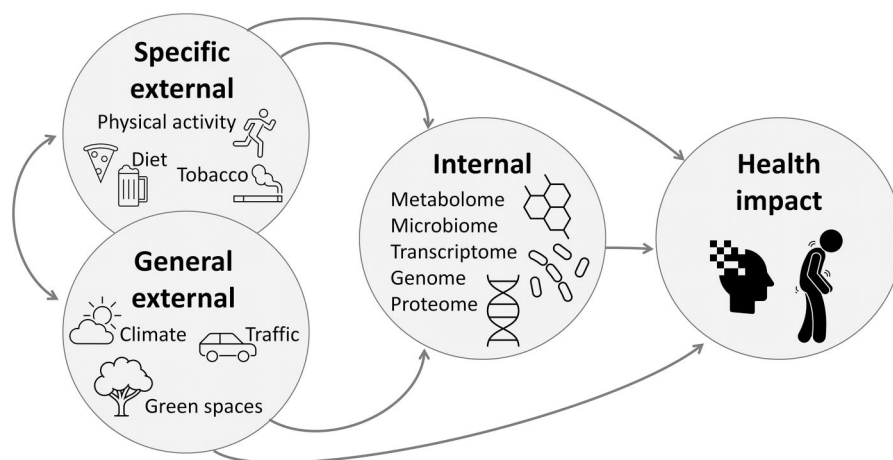


Figure 1. The exposome concept and how the specific external, general external and internal exposome contribute towards health impacts. Modified from.^{9,10}

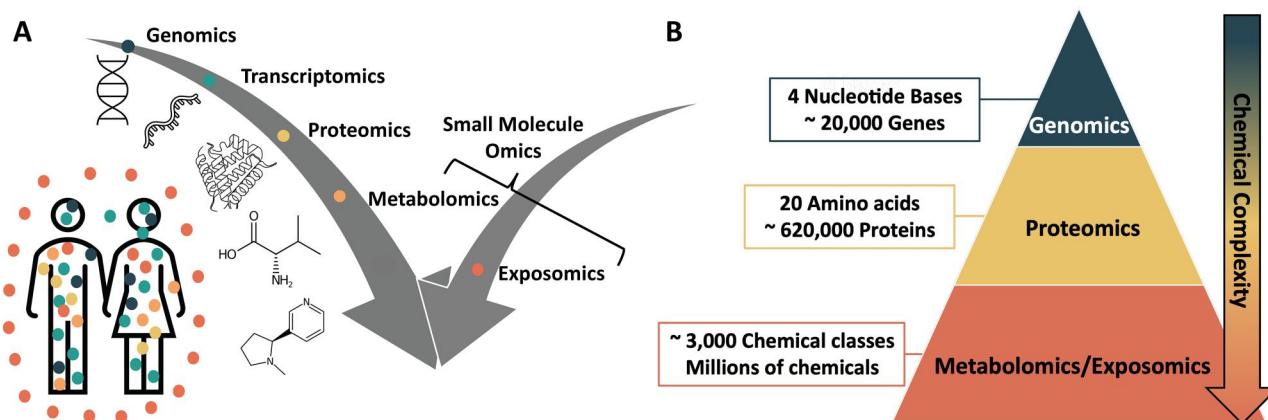


Figure 2. (A) The omics cascade from genome onwards, adapted from.^{11,12} (B) The differences in chemical complexity of the different omics, adapted from.^{16,17} Note that the colours represent the different “omes” (genes, proteins, metabolites).

related with the genome and proteome, leading some to call them the “canaries of the genome”.¹⁶ This metaphor underscores that even minor alterations, such as a single base change in a gene (especially in a metabolic enzyme), can potentially lead to a 10 000 fold-change in concentrations of specific small molecules.¹⁶

The disparity in complexity in the different “omics” layers is shown in Figure 2B, which helps explain why genomics and transcriptomics are the most mature omics, followed by proteomics, then small molecule omics.^{17,18} Current automated high-throughput techniques are able to (almost) completely cover the genome and proteome.¹⁷ Gene sequencing requires a DNA sequencer, while protein characterization can be performed on a single type of high-resolution mass spectrometer.¹⁶ In contrast, a wide range of analytical instrumentation is required to capture small molecules, explored further in section 2.4 below.

Since the border between the metabolome and exposome is not always clear-cut,^{2,19,20} a distinction can be made between environmental exposure and the resulting biological response² (see Figure 3). The metabolome refers to the complete set of small molecules (ie, metabolites) present within a biological sample such as a cell, tissue, organ or organism at a given time.¹⁷ As these metabolites arise from both endogenous and exogenous sources, the metabolome is considered a key measure for exposome research.¹⁹ In this context, the metabolome can be

considered as a subset of the exposome, often referred to as the internal chemical exposome.²⁰ Thus, while all metabolites found in a biological sample can be considered part of the exposome, not all chemicals within the exposome belong to the metabolome. In contrast, the exposome is a broader concept that encompasses not only chemicals but also all physical, biological, and psychosocial influences that may impact health, as described above and illustrated in Figure 1.

This review will specifically delve into common workflows for measuring the *chemical exposome*, with a particular focus on non-target high-resolution mass spectrometry (HRMS) coupled to liquid chromatography (LC). Although standardized workflows remain an area of active research (discussed further in^{2,20,21}) they can be divided into three main components (Figure 4), which have been used to construct the subsequent sections of this article: (1) *experimental workflow*, which encompasses experimental design, sample collection, sample preparation and data acquisition; (2) *computational workflow*, which includes data pre-processing and compound annotation; and (3) *statistical analysis and biological interpretation*.

A recent and extensive exposomics review was published in 2024 by Lai et al.²¹ however the present review has a different focus. Specifically, this manuscript emphasizes the conceptual, and often blurred, distinctions between the metabolome and the exposome, as well as the analytical and computational

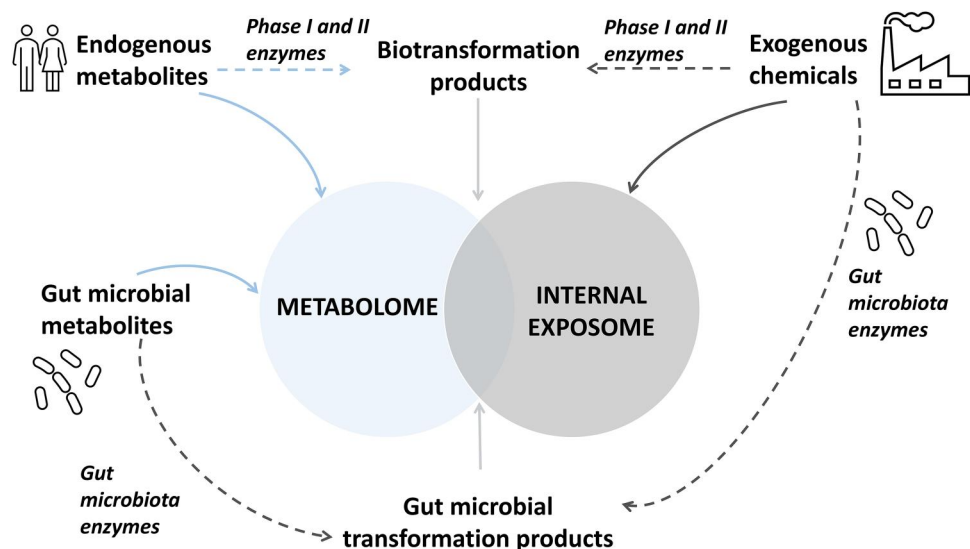


Figure 3. The chemicals part of the metabolome, exposome and the overlap. Adapted from.^{2,20} Note that gut microbiota is illustrated as a major microbial contributor to metabolic processes, however other microbiota such as saliva, nasal and skin, among others, can also contribute.

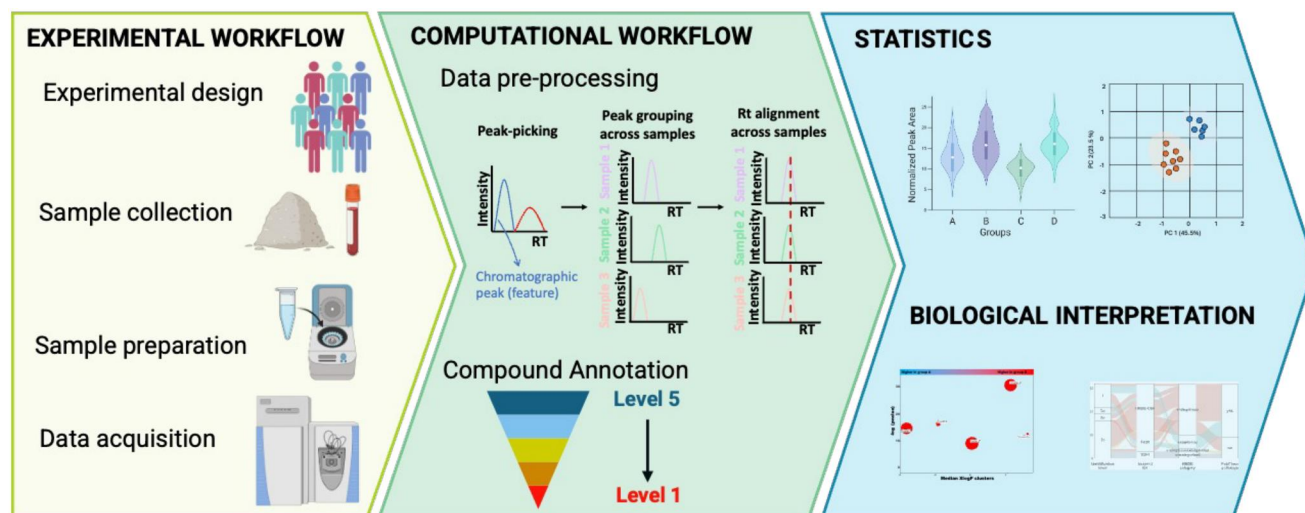


Figure 4. Common workflow steps to investigate the human chemical exposome.

workflows employed in non-target HRMS-based exposomics studies. These workflows are illustrated with visualizations to facilitate understanding, and the manuscript also discusses key challenges and strategies to advance towards harmonized studies. This review and the references provided are intended for readers entering the exposomics field, such as PhD students and early-career researchers but will be also helpful to any researcher investigating the human exposome.

Experimental workflows

Experimental design

Individuals may encounter millions of chemical exposures throughout their lifetime, where some exposures may exert life-long consequences, while childhood exposures may increase the risk of developing disease later in life.¹⁹ The design of prospective longitudinal studies (Figure 5), which collect samples at various life stages including perinatal, childhood, and adulthood, holds inherent advantages over those that gather single samples from individuals who have already developed the disease.^{19,22}

Although challenging, promising initiatives are currently underway to provide such longitudinal data, such as All of Us²³ and HELIX.²⁴ While all exposomics studies should undergo independent validation to ensure the reproducibility, this poses significant challenges due to variations in exposures across populations.²² To ensure robust scientific conclusions, several key factors must be considered, such as the number of samples, the types of samples (matrix selection) and their storage conditions.²¹

The number of samples required for a small molecule omics experiment depends on the biological variability of the studied system and the analytical variability of the technology employed.²⁵ Generally, hundreds or thousands of subjects need to be investigated to obtain robust results, although sample numbers of 3-20 per group can be suitable for generating preliminary and/or pilot data.²⁵ Such preliminary data can then be used to estimate the number of samples needed to achieve certain statistical power (eg, 0.8), using open tools such as G*Power,²⁶ MetaboAnalyst²⁷ and MultiPower,²⁸ which estimates sample sizes for multi-omics studies.

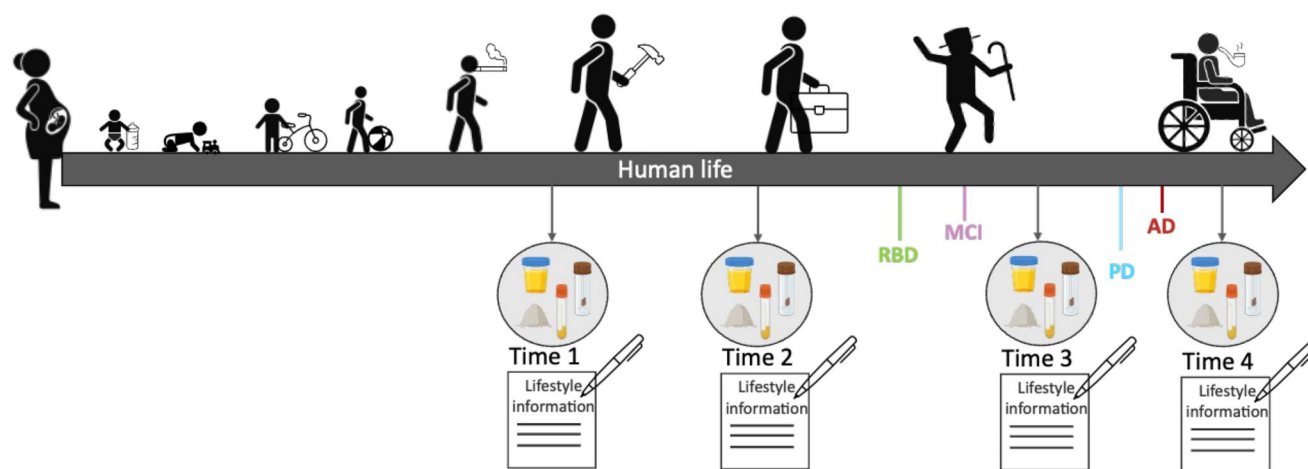


Figure 5. Human life timeline (top) and proposed longitudinal study for an exposomics study where both environmental and biological samples are collected. AD; Alzheimer's disease, MCI; Mild Cognitive Impairment, PD; Parkinson's disease, RBD; REM- sleep behavior disorder. Note that while neurodegenerative diseases were chosen for illustrative purposes, this conceptual framework can be applied to other diseases such as cancer.

Sample collection

The matrix selection is a critical aspect of the study design, affected by multiple factors including the cost of storage and collection devices, research question, technical feasibility and expected presence of specific chemicals, as some chemicals accumulate specifically in some tissues.^{21,29} To date, blood (plasma or serum) and urine are the most studied matrices.^{21,22} While biological samples can provide insights into potential biological responses (ie, endogenous compounds) to toxicants or pollutants (ie, exogenous compounds), these exogenous compounds are typically present at trace levels compared to the endogenous ones and thus stretch the dynamic range of current instrumentation,³⁰ as detailed below. Consequently, it may be advantageous to also incorporate non-target screening of environmental samples such as dust to first build a picture of potential contaminants, before analysing biological samples.

The collection and storage of many biological samples (at -80°C) for long-term longitudinal studies (Figure 5) can be costly. Dried Blood Spot (DBS) samples are compact, easy to collect and can be shipped at ambient temperature. However, the analysis of DBS is not as standardized as for plasma or serum, and it is not yet clear if there is sufficient sample material for non-target analysis.³¹ Urine is a well-studied matrix, easily accessible and less invasive than blood,²² while fecal analysis can yield valuable insights into gut microbiota metabolites, with relatively easy sample collection. Although cerebrospinal fluid (CSF) provides valuable information on brain metabolism, the collection of CSF samples is invasive, and the recruitment of healthy volunteers is more difficult, such that studies with large sample sizes involve significant commitments and dedication. The concentrations of exogenous chemicals in this matrix are typically very low compared with other biological samples such as blood. Other sample types that can be collected to answer specific questions include hair, teeth and nails, which can help to understand historical exposures, although determining the timing of exposure can be laborious.² The study of the specific external exposome can be done by collecting household dust samples with a vacuum cleaner,³² although these samples are prone to variability. Passive samplers such as silicone wristbands allow a more individual assessment of the exposure to chemicals, although the silicone used in wristbands can lead to analytical interferences if

they are not properly cleaned before use and the sample preparation may be more tedious compared to household dust.^{33,34}

Depending on the aim, two different approaches can be used when analysing the chemical exposome: (1) *target* or *targeted* studies, which focus on identifying a limited number of specific chemicals (hypothesis driven) and (2) *non-target* or *untargeted* studies, which are hypothesis generating and aim to identify as many compounds as possible.^{21,35} Typically, non-target studies are *semiquantitative*, which use peak height or area as direct read-outs or estimate the concentrations relative to quantifiable compounds such as internal standards spiked at known concentrations. In contrast, target studies frequently employ *absolute quantitation*, which involves calculating the exact concentration via calibration curves.²¹ The sections below focus on non-target HRMS studies, which can be used to generate hypotheses and insights for designing subsequent quantitative targeted studies.

Sample preparation

Sample preparation (or pretreatment) prior to instrumental analysis aims to reduce interferences, separate and concentrate analytes.³⁶ Due to the chemical diversity of the exposome, there is no universal method that captures all chemicals present in a sample²⁵ and a combination of different sample preparation approaches can cover a broader range of the chemical space.²¹ However, this is limited by increased costs in resources and time.

Sample preparation methods for non-target analysis of biological/environmental samples should be:³⁷ (1) unselective (to cover a wide range of chemicals), (2) simple and fast (to prevent chemical loss/degradation during preparation), (3) reproducible, and, for biological samples, (4) incorporate a quenching step. The rapid stopping or quenching of metabolism is an essential step to produce a stable extract that reflects the endogenous metabolite levels present in the original biological system.^{37,38} The sample preparation method is highly dependent on the matrix (eg, cells, plasma, or tissues), and analytical platform employed. For instance, although protein precipitation is the first step for plasma samples, tissues must be homogenized first.³⁹ Furthermore, while Gas Chromatography coupled with MS (GC-MS) often requires a derivatization step to make the compounds sufficiently volatile for the analysis,⁴⁰ the sample preparation for

Liquid Chromatography coupled with MS (LC-MS) analysis is simpler and typically does not require this step.³⁹

Liquid-liquid extraction (LLE), often including a protein precipitation step, and dilute and shoot (DNS) are frequently employed sample pretreatment methods in exposomics studies. Other sample preparation methods include solid-phase extraction (SPE) and dispersive solid-phase extraction, such as QuEChERS.^{21,36} While LLE and DNS offer broad analytical coverage, they are susceptible to matrix effects and interferences, which can limit reproducibility and sensitivity. In contrast, SPE and QuEChERS typically enhance sensitivity and reproducibility, albeit often at the cost of reduced chemical coverage.³⁶ Further information can be found in the NORMAN guidance for non-target screening¹⁵ and other publications on sample preparation for metabolomics^{25,37,39} and exposomics.^{21,36}

Data acquisition

While HRMS has emerged as the leading technique to investigate the chemical exposome, there is no “one size fits all” analytical method and a combination of different separation and ionization platforms is needed to capture the relevant chemical space (Figure 6).³⁵ Flow-injection is very fast, but cannot separate any isobars, while chromatographic separation introduces analytes slowly into the mass spectrometer, separating more isobars and reducing the risk of ion suppression, source fouling and coelution. A mobile phase, either liquid or gas, transports the analytes through a stationary phase-fixed system^{2,2,21} (see top part of Figure 6).

GC is frequently used for the analysis of volatile and thermally stable compounds such as fatty acids, and organic compounds (see left part of Figure 6 for examples). Due to the high temperatures, a derivatization step is often required before the analysis, which may result in compound loss^{39,41,42} and changes the resulting mass spectra. The most common ionization technique applied in GC is Electron Ionization (EI), which provides robust, and highly reproducible fragmentation patterns.³⁹

Currently, LC coupled with an electrospray ionisation (ESI) source is probably the most used HRMS-based platform for non-target studies due to the soft ionization process, high dynamic range and versatility.² LC-HRMS does not typically require derivatization and is highly applicable to the analysis of a broad range of medium to very polar compounds,^{39,41} with several examples shown in the right part of Figure 6 Using both positive (+) and negative (-) ionization modes increases the coverage. Reversed Phase (RP) columns are widely employed to separate polar and medium polar compounds, providing relatively reliable, robust and reproducible results. Hydrophilic Interaction Liquid Chromatography (HILIC) columns improve the separation of very polar compounds minimally retained by RP,^{21,39,42} but can be less reproducible. Alternative ionisation techniques such as Atmospheric Pressure Photoionization (APPI) and Atmospheric Pressure Chemical Ionization (APCI) can extend the LC range almost into the GC range (Figure 6).

Despite recent progress, current analytical platforms still lack the dynamic range to simultaneously detect trace-level exogenous compounds and high-abundance endogenous metabolites (see Figure 7 for an example based on LC-ESI-HRMS). The

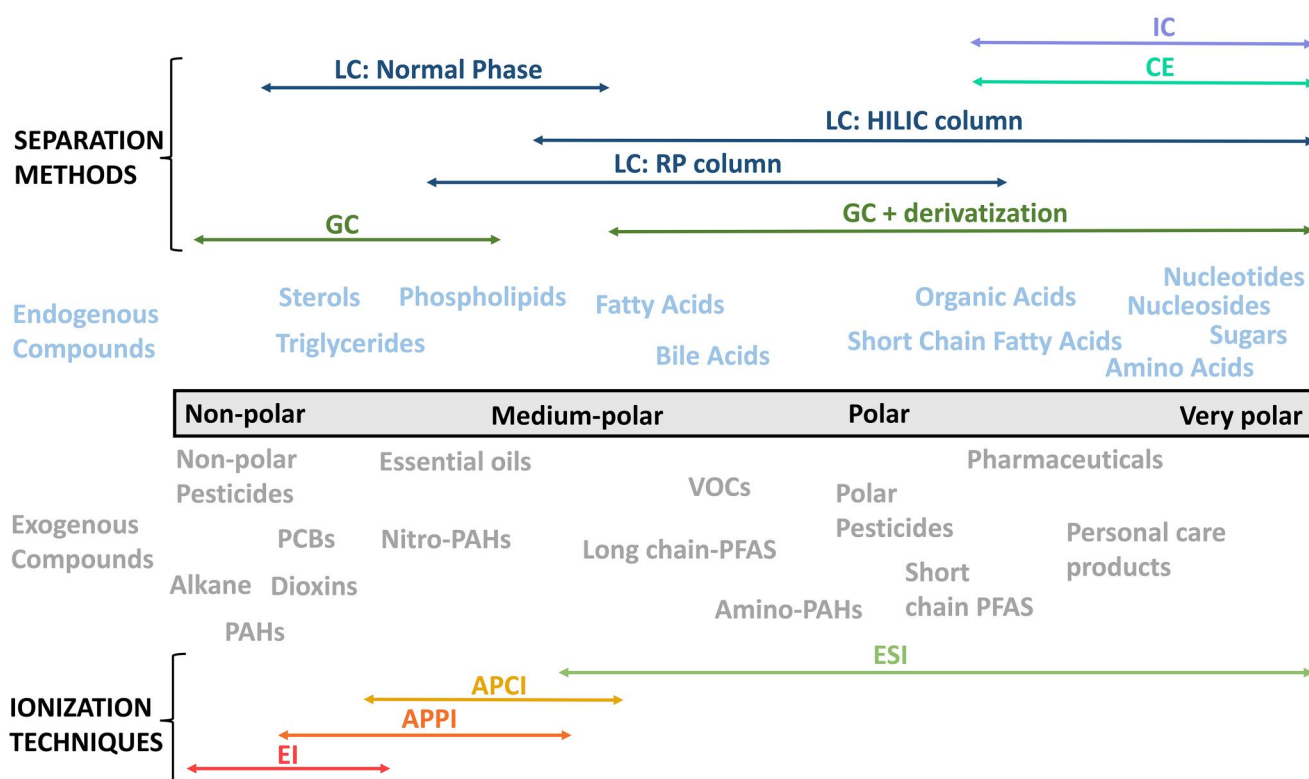


Figure 6. Separation analytical methods and their applicability range based on the polarity of the chemicals are displayed on the top part of the figure, while the different ionization techniques and their applicability range are displayed on the bottom. Examples of potentially endogenous (blue) and exogenous (gray) compounds are displayed. Adapted from Zeki et al.³⁹ and Hollender et al.¹⁵ Abbreviations: GC, gas chromatography; HILIC, Hydrophilic interaction chromatography; LC, liquid chromatography; RP, Reversed Phase; EI, electron ionization; ESI, electrospray ionization; APCI, atmospheric pressure chemical ionization; APPI, atmospheric pressure photoionization; VOCs, Volatile Organic Compounds; PFAS, Perfluoroalkyl and Polyfluoroalkyl Substances; PAHs, Polycyclic Aromatic Hydrocarbons; IC, Ion Exchange; CE, Capillary Electrophoresis.

Human metabolome and internal exposome

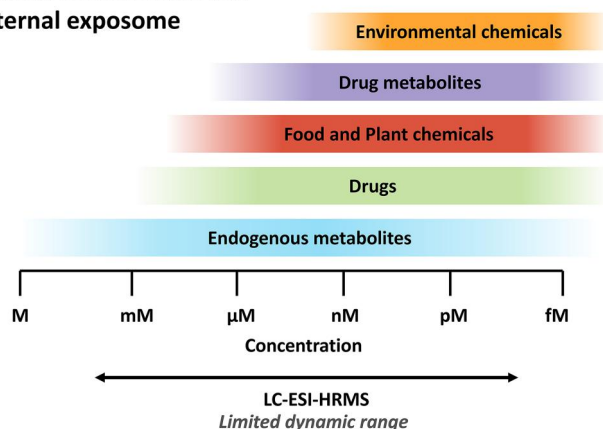


Figure 7. Concentration range of the different components of the human metabolome and internal chemical exposome as well as coverage by LC-ESI-HRMS. Adapted from.^{2,17}

integration of Ion Mobility Spectrometry (IMS) into HRMS workflows has been gaining more attention as it can improve the dynamic range and throughput of the analysis.^{2,18} Coupling IMS with HRMS offers the ability to resolve isomers or isobars that are difficult to distinguish using only HRMS, and adds Collision Cross Section (CCS) information to the retention time, MS1 and MS2 data, providing complementary information for compound identification and/or annotation.^{2,35} Depending on the instrument, including IMS also reduces the complexity in the MS1 and MS2 spectra,²¹ but comes at a cost of sensitivity and more difficult data analysis.^{2,43} Typical commercial IMS instruments do not yet offer sufficient resolution to resolve isomers within measurement and prediction errors.^{44,45}

Quadrupole Time-of-Flight (Q-TOF) and Orbitrap are the most commonly used mass analyzers in non-target studies, as they provide high resolution.^{2,39,42,46} A Q-TOF analyzer maintains the selectivity of the quadrupole and provides a mass resolution of approximately 40 000–60 000. Orbitrap analyzers reach resolutions between 250 000 to even 1 000 000 for ions with m/z below 300.⁴⁶ Triple Quadrupole (QQQ) analyzers are often used for targeted analysis to confirm potential biomarkers due to high sensitivity and selectivity,^{41,42} but are not suitable for non-target studies due to their low resolution.³⁹

Two acquisition methods are often used in non-target HRMS studies: Data-Dependent Acquisition (DDA), and Data-Independent Acquisition (DIA). DDA is typically more common, since MS2 spectra can be associated with a specific precursor,^{2,47,48} resulting in simpler processing compared with DIA.^{47,49} While low intensity features may not be selected for fragmentation, this can be alleviated with the creation of inclusion lists, supported by either open software solutions such as iterative exclusion lists from IE-Omics,⁵⁰ or instrument software such as AcquireX from ThermoFisher. DIA mode operates in a less-selective manner, as the instrument selects all the peaks within the isolation window, regardless of the peak intensity.^{2,47} Thus, the MS2 spectra are information-rich collections including fragments from low intensity ions, but they lack the precursor-fragment information⁴⁷ and require deconvolution to link the specific precursor to the MS2 spectra.^{2,48} This can be performed by open software such as MS-DIAL.⁵¹ All-Ion Fragmentation (AIF), and Sequential Window Acquisition of all Theoretical

Fragmentation Spectra (SWATH) are two commonly employed DIA methods.

In practice, method selection depends strongly on the research question, instrument availability and analytical expertise, among others. While LC-MS (eg, QQQ) remains the preferred option for the quantification of a specific and small number of compounds (target studies), LC-HRMS and GC-HRMS (eg, Orbitrap) are employed for non-target exposomics studies generating larger and more complex datasets that make the data preprocessing steps more challenging. Therefore, the analytical considerations detailed above, together with the references provided, can serve as practical guidance for selecting the most appropriate analytical platform and acquisition mode for a given a study.

Computational workflows

Non-target HRMS data pre-processing

The main objective of data pre-processing (or feature detection) is to transform the raw data files into a format that simplifies the access to the distinct characteristics of every observed feature in each sample analyzed, ie, a feature list.⁵² A “feature” does not necessarily refer to a specific chemical, but rather it typically refers to a peak (or signal) identified at a specific retention time, and m/z , containing spectral information such as MS1 and/or MS2.³⁵ Features are also called “ m/z features”, “ion features” or “ion peaks”.²¹ The resulting feature list (Figure 8, bottom right) can be then employed for various purposes, including feature prioritization, compound annotation, and statistics and typically contain the retention time, m/z , and peak area and/or intensity of each feature from each raw data file.^{21,52–54}

Data pre-processing steps include data conversion (if vendor software is not used), centroiding, filtering step for noise removal, generation of Extracted Ion Chromatograms (EICs), peak picking, peak grouping across samples, and retention time alignment (Figure 8).^{2,54} Software approaches include vendor software such as Compound Discoverer (ThermoFisher), MassHunter Profinder (Agilent), MetaboScape (Bruker) and Progenesis QI (Waters), as well as open software including MS-DIAL,⁵¹ MZmine,⁵⁵ OpenMS,⁵⁶ XCMS,⁵⁷ and patRoan.⁵⁸ Recent reviews provide a comprehensive overview of the state of the art of data pre-processing software.^{15,20,59,60} Since different instrument vendors use different formats, conversion of the raw data files into an open format such as mzML or mzXML⁵³ is required for open approaches, typically using ProteoWizard,^{51,62} although some software now embed this conversion, including MS-DIAL⁵¹ and patRoan.⁵⁸ Data acquisition can be done in either profile or centroid mode; the centroiding of profile data is an important data reduction step (see top part of Figure 8),^{15,54} often performed together with the data conversion using ProteoWizard.⁶¹ While ProteoWizard offers both vendor-specific or general algorithms, the use of vendor-specific algorithms is recommended¹⁵ The quality of the conversion using ProteoWizard hinges on the vendor type—while the ThermoFisher conversion yields high quality results, this is not the case for Waters.⁶³

Once data is centroided, a filtering step is applied to suppress or reduce the random analytical noise, which is always present in the acquired MS data.^{53,64,65} Noise is often removed via smoothing.^{64,66} Methods include linear weighted moving average,⁵¹ Savitzky-Golay smoothing,⁶⁷ moving average,⁶⁷ and binomial filter.⁶⁸ MS-DIAL uses the linear weighted moving average by default, although it supports all the aforementioned approaches.⁵¹ After the EIC generation, peak picking is performed to determine the area under the peak.^{52,53,66} This requires the

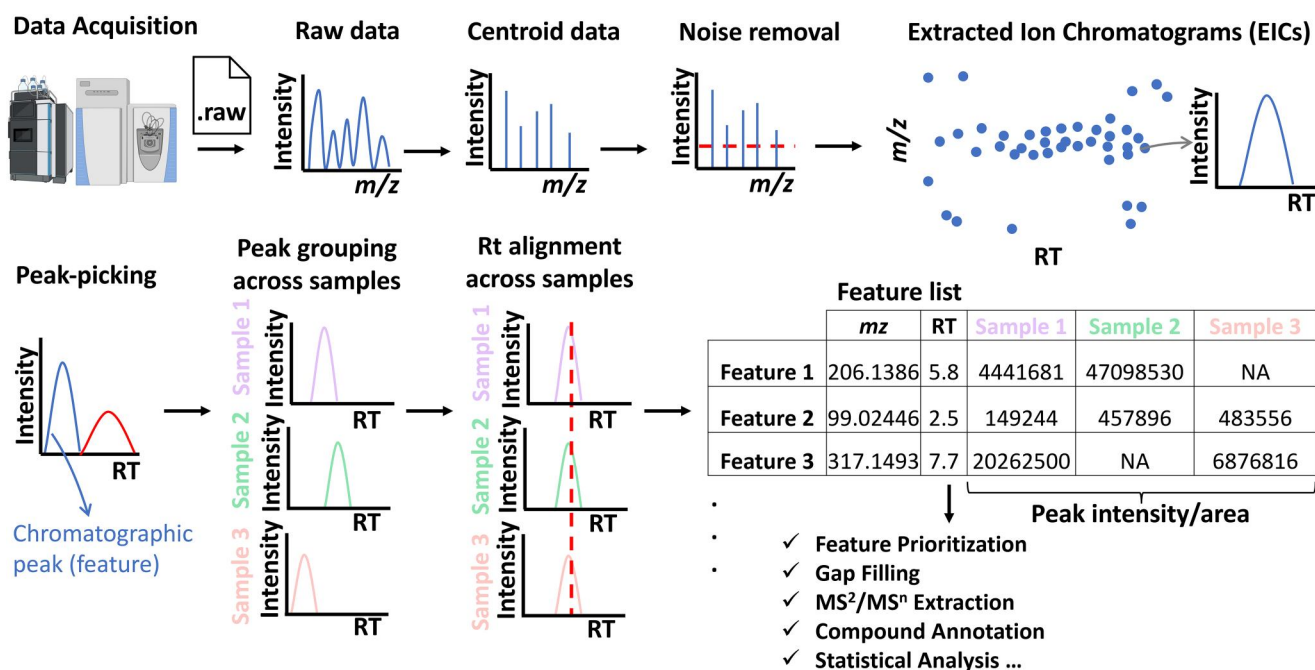


Figure 8. Data preprocessing steps for LC-HRMS raw data. First, data is centroided and noise is removed. Next, EICs are generated and a peak-picking algorithm is applied to detect true peaks. Finally, peaks are grouped across samples, and retention time alignment is performed. After that, a gap filling step can be performed to reduce the number of missing values. Adapted from.^{54,60,66} Abbreviations: RT, Retention Time; *m/z*, mass-to-charge ratio.

establishment of criteria (ie, parameters) to distinguish true peaks from noise⁶⁴ such as setting a minimum intensity threshold and an estimated chromatographic peak width, establishing a maximum *m/z* error (eg, 0.001 Da for a Q-TOF instrument and 5 ppm for an Orbitrap), and ensuring that identified peaks are present in a significant proportion of the samples.^{54,64,65} Next, the selected peaks are grouped across samples,⁵⁴ and an alignment step is needed to correct retention time differences between runs (samples).⁵³ The alignment algorithm often requires a reference sample (eg, pooled Quality Control [QC] sample) to correct the retention time differences, and this choice can have a significant impact on the results.^{53,65}

The selection of parameters for data preprocessing is a critical step in the exposomics analysis,⁵⁴ as inadequate parameter selection can lead to biased results. Since exposomics seeks to identify highly abundant endogenous compounds along with exogenous small molecules that are often present at trace levels (see Figure 7), there are no universal parameters that can effectively capture the chemical complexity of the human exposome currently.²¹ While approaches such as IPO⁶⁹ can be used to assist in parameter selection, they should be employed with a degree of caution, as they tend to discard low abundant or rare peaks, which may be trace level exogenous molecules of relevance for the exposomics question. Recently proposed quality assurance/quality control (QA/QC) guidelines support the exposomics community in their data preprocessing choices,⁵⁴ see further details in Section 5.

Compound annotation

The annotation of small molecules remains a major challenge in non-target HRMS based metabolomics and exposomics studies, as the majority of features detected (around 80%) remain unknown,⁷⁰ although experts continue to debate the fraction. High abundant peaks, with MS² available, usually represent <10% of the detected features by non-target HRMS, while the majority (60%-70%) are low abundant peaks without fragmentation

information.⁷¹ *Annotation* involves associating an identified MS feature with a specific chemical identity, while *identification* is the process of verifying that the annotated compound corresponds to the proposed chemical (ie, confirming the annotation with the reference standard).⁷² However, the lack of availability of chemical reference standards for given molecules of interest represents a major bottleneck in exposomics/metabolomics studies, complicating biological interpretations.⁷³

The large number of unannotated features in non-target studies is a common subject of debate. Recently, Giera et al.⁷⁴ suggested that most (>70%) of the unannotated features in non-target LC-HRMS experiments are in-source fragments (ISFs), concluding that the dark metabolome may be smaller than previously thought. However, the results were based on a 931K compound library that is not publicly disclosed, while the formation and recognition of ISFs is highly dependent on multiple parameters including source voltages, instrument design, matrix type, extraction methods, analyte concentration and data preprocessing workflow.⁷⁵ Moreover, many of the ISFs observed in highly concentrated synthetic chemical standard mixes may not be detected in complex biological samples. Previous studies, in contrast, have reported that ISFs entail ~2%-25% of all detected ions.⁷⁵⁻⁷⁷ While ISFs can be problematic, they can also be intentionally enhanced to improve annotation confidence⁷⁸ and modern workflows are increasingly accounting for them.⁷⁵ These different estimates highlight the current lack of consensus and harmonization in the field. While high prevalence of ISFs would imply that the dark metabolome is smaller than previously thought, assuming that most of the features identified are known, lower estimates would indicate that although ISFs remain relevant, they may represent a relatively small subset of detected features.

Non-target HRMS studies generate vast amounts of data, necessitating various computational strategies such as suspect screening and non-target screening (top part of Figure 9A). While suspect screening approaches use lists of chemicals that could

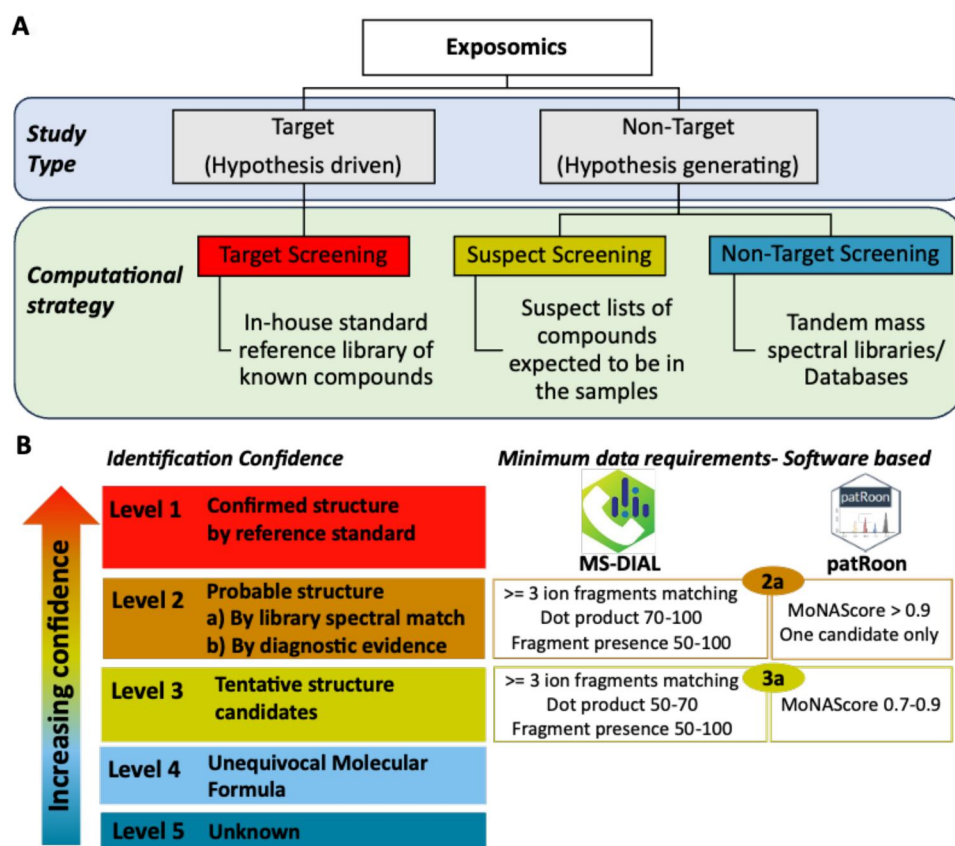


Figure 9. (A) Generic computational workflow for target and non-target exposomics studies.^{15,35,79,80} (B) Identification confidence levels by Schymanski et al.⁸¹ (left), and proposed minimum data requirements for level 2a and 3a annotations using MS-DIAL and patRoön software (right). Similar approaches could be done with other software such as MZmine⁵⁵ and LipidMatch.⁸²

potentially be present in the samples (ie, suspect lists of certain chemical classes or organisms/matrices) for more efficient discovery, non-target screening approaches aim to identify as many compounds as possible via tandem mass spectral libraries or database search (eg, via in silico fragmentation software), as shown in the bottom part of Figure 9A. A detailed glossary of terms can be found in the 2023 NORMAN guidance.¹⁵

Compound annotation is a fundamental step to convert raw HRMS data into meaningful biological information.^{35,71} Since the amount of information available for identification varies, it is essential that the confidence assignment of each feature is transparent.³⁵ Several confidence level schemes exist, including the 2007 Metabolomics Standards Initiative (MSI),⁸³ the 2014 guidelines for HRMS data with an environmental focus⁸¹ (left part of Figure 9B), 2020 guidelines for IMS⁸⁴ and 2022 guidelines for PFAS⁸⁵ and GC-HRMS.⁸⁶ The 2014 levels shown in Figure 9B range from Level 1 (confirmed structure with reference standard), to Level 5 (only a m/z is known).⁸⁷ According to metabolomics terminology conventions, only Level 1 can be considered identifications, while the rest (Level 2-5) are annotations. In downstream analyses, it is recommended to base biological interpretations primarily on Level 1 identifications, however, this is often hampered by the limited availability of chemical standards. Consequently, Level 2 annotations are frequently used as the next most reliable basis for interpretation, whereas Level 3-5 should be interpreted with caution, serving mainly for prioritization for future validation efforts.

Several criteria are used to assign the identification level of each feature, including the MS1, retention time, fragmentation

pattern (MS2), CCS (if available) and experimental data.⁷¹ However, although different classification systems have been proposed,^{38,81,83,88,89} currently there is no standardized system for compound annotation integrated in the processing workflows, making the comparison of results between studies challenging.⁷¹ This necessitates a degree of “translation” between software outputs (see eg, Figure 9B, right). While the use of False Discovery Rates (FDR), as done in other omics fields, has been proposed, estimating total FDR for compound identification in small molecule omics is still a nascent challenge in the field.^{90,91} The use of identification probability was proposed recently as an alternative to the identification levels.⁹¹ However, the probability depends on the reference library size and treats all candidates equally, such that smaller reference libraries can artificially lead to high identification probabilities, while larger libraries (such as those more applicable to exposomics) will potentially yield too many apparently equally-valid candidates with very low probabilities to support meaningful outcomes. Identification levels can be assigned automatically in some patRoön workflows,^{58,92} while the NORMAN Network also trialled an automated assignment system⁹³ and Boatman et al. recently published a checklist to facilitate automation for PFAS identification with IMS.⁹⁴

Compound annotation via tandem mass spectral libraries search

The fastest and most accurate (and thus most common) strategy for compound annotation is to compare the experimental mass spectra (MS2), with standard mass spectral libraries.^{35,47,70} Compound annotation via spectral library searching is based on

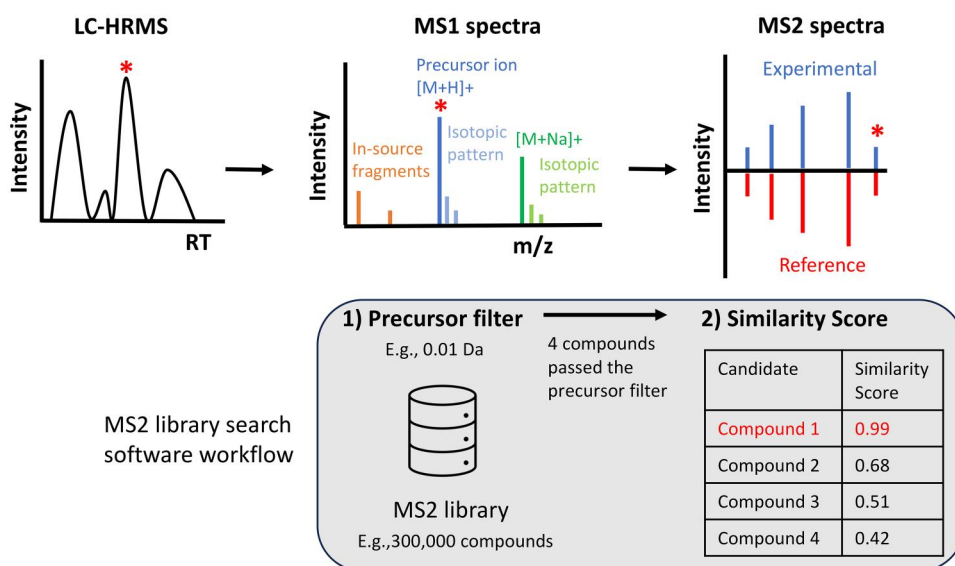


Figure 10. Exemplified workflow of MS2 library search software.⁴⁷ For a given experimental MS1, first the precursor filter is applied to remove all the candidates outside the tolerance window (eg, 0.01 Da). Subsequently, the similarity algorithm ranks the experimental MS2 spectra against the remaining library spectra candidates (four in this example) and calculates a similarity score. Note that the adducts shown ($[M+H]^+$ and $[M+Na]^+$) are illustrative examples and other adducts (eg, $[M+K]^+$, $[M+NH_4]^+$ in positive mode, $[M-H]^-$, $[M+Cl]^-$ in negative mode) can occur based on the matrix and acquisition settings, among others. Abbreviations: RT, retention time; m/z , mass-to-charge ratio.

the premise that molecules generate a reproducible “fingerprint” under specific fragmentation conditions^{72,95} (see Figure 10). Good matches between the experimental and the library MS2 spectra can lead to Level 2a annotations.⁷² If the MS2 library is created *in-house*, ie, the experimental and library spectra are acquired under the same conditions, this can lead to Level 1 annotations when the retention times also match. Commonly used spectral libraries include MassBank,⁹⁶ MassBank of North America (MoNA),⁹⁷ GNPS,⁹⁸ mzCloud,⁹⁹ METLIN,¹⁰⁰ and NIST.¹⁰¹ While GC-EI mass spectra have been standardized for over 60 years, LC-MS spectra are less standardized due to instrument variability and differences in acquisition parameters such as collision energy. As a result, MS2 libraries often contain multiple entries for each compound.^{72,91}

Tandem mass spectral libraries are typically generated by the analysis of chemical reference standards.⁷² Therefore, compound annotation by this approach is hampered by the limited availability of mass spectra data due to lack of standards, and/or lack of open data. Of the 122 million compounds in PubChem (September 2025), only 36 242, or 0.03% have open LC-MS/MS data available.¹⁰² Nonetheless, there has been substantial progress in the quality and quantity of mass spectral libraries in recent years. Automated spectral library searching and matching can be performed using various open and commercial software with different algorithms. Typically, these software tools work in a two-step procedure: (1) MS1 filter, (eg, 0.01 Da) which can remove up to 99% of false candidates and speeds up searching, and (2) similarity algorithm, which ranks the experimental MS2 spectra against the remaining library spectra and calculates a similarity score, see bottom part of Figure 9.⁴⁷ Ideally, scores should be able to distinguish true and false positive matches.⁷²

The most common spectral match score is the cosine score, which converts two MS2 spectra (observed and reference) into two equally size vectors through mass peak binning, and then calculates the dot product which ranges from 0 to 999 (or 0 to 0.999).^{47,73} A score of 999 indicates a perfect match between the two spectra, while a score of 0 indicates no match.⁴⁷ Newer

scoring approaches include the Entropy score¹⁰³ and a range of new machine learning algorithms.⁷³ The right part of Figure 9B shows how these scores can be applied to annotate compounds using different software. Several open software approaches support spectral library search, including MS-DIAL,⁵¹ MZmine,⁵⁵ openMS,⁵⁶ and XCMS,⁵⁷ while commercial examples include Progenesis QI (Waters) and MetaboScape (Bruker). The *in silico* fragmentation software MetFrag also integrates a library search and calculates an Exact Spectral Similarity score (also known as “MoNA score”).¹⁰⁴ However, the variability in output results, including similarity scores, across different software can make the identification levels obtained poorly comparable. For instance, MS-DIAL and patRoos provide similarity scores, dot product and MoNA scores, respectively, with a proposed minimum data requirements for annotation shown in Figure 9B. Importantly, even when the spectral similarity score is high (>0.9), the identity of the compound must be confirmed with a chemical reference standard to classify it as Level 1.²¹ Otherwise, the match should be considered an annotation rather than an identification (Level 2 or below), since several isomers can have very similar spectra such that only retention time or other orthogonal parameters can distinguish between them.

In silico approaches for compound annotation

In silico fragmentation software supports compound annotation of candidates beyond those in mass spectral databases. These methods typically involve matching the experimental spectra against a selection of candidates obtained from known compound databases (discussed in the next section). Approaches such as MetFrag,¹⁰⁴ Mass Frontier, MS-FINDER¹⁰⁵ and CFM-ID¹⁰⁶ fragment the candidates and match the resulting spectra with the experimental spectrum, while approaches such as CSI: FingerID¹⁰⁷ and SIRIUS use the experimental spectrum to generate fingerprints, which are matched with the fingerprints of the candidates to rank the structure candidates.^{70,108} MetFrag uses a bond dissociation approach to generate fragments for each candidate, which are compared with the experimental spectra to

determine which are the best candidates.¹⁰⁴ Mass Frontier (Thermo) uses rule-based fragmentation prediction, complementary to the bond disconnection approach. CFM-ID¹⁰⁶ is a machine learning-based approach that can predict fragments and intensities, and thus can be used to generate *in silico* libraries of the given spectrum type used during the training.⁴⁷

In silico spectral libraries can help to overcome the limited number spectra in MS2 libraries and avoid the need for “on the fly” calculations in each workflow. *In silico* spectral libraries can be generated via *eg*, quantum chemistry, machine learning, heuristic-based, and chemical reaction-based methods. Heuristic approaches are best applied to compounds with consistent fragmentation patterns such as lipids.⁴⁷ LipidBlast,¹⁰⁹ integrated within MS-DIAL⁵¹ and LipidMatch,⁸² is an *in silico* library containing more than 200 000 spectra generated using a heuristic approach. LipidMatch⁸² is a rule-based software that incorporates various libraries to facilitate the lipid annotation. CFM-ID¹⁰⁶ has been used to *eg*, generate *in silico* EI-MS and MS/MS spectra of small molecules in HMDB.¹¹⁰ In general, predictions should only be used for compounds within or close to the domain of the training set/rule sets used.

In silico approaches for compound annotation typically yield Level 3 annotations or below, but they can be upgraded to Level 2 with the support of a good tandem mass spectral library match, such as the combined approach with MoNA used in MetFrag. *In silico* annotations often serve an important early role in the elucidation process, guiding subsequent activities such as the interpretation, prioritization and even acquisition of reference standards.⁷³

Compound databases for compound annotation

Due to the extreme chemical diversity of the chemical exposome, the database selection plays a critical role in the annotation process. This aims to reduce both false positives (ie, incorrect annotations) and false negatives (absence of the correct structure in the database). While the Chemical Abstract Service (CAS)¹¹¹ database is the largest chemical registry containing over 290 million organic substances (September 2025),¹⁵ it is not freely available or compatible with open software approaches. ChemSpider¹¹² and PubChem¹¹³ contain over 128 and 122 million chemicals respectively (September 2025), making them the two largest freely available chemical databases. However, due to user quota limitations on ChemSpider, PubChem currently emerges as the most feasible large chemical database for integration into open software workflows.¹⁵

The use of smaller subsets of chemicals helps in the annotation process as these contain known molecules specific to certain domains. For example, PubChemLite for Exposomics (PCL), a subset of PubChem containing 442 379 chemicals (version 2.0.0)^{87,114,115} and the Blood Exposome Database¹¹⁶ (67 291 compounds) aid exposomics researchers in identifying relevant chemicals. Metabolite discovery in humans is facilitated by the HMDB,¹¹⁰ while KEGG¹¹⁷ and MetaCyc¹¹⁸ establish connections with proteomics and transcriptomics disciplines.⁷⁰ Lipidomics studies are supported by LIPID MAPS Structure Database (LMSD),¹¹⁹ which contains 49 790 unique lipid structures (July 2025), which is thus the largest public lipid-specific database. The Human Microbial Metabolome Database (MiMeDB) facilitates the study of small molecules produced by the human microbiome,¹²⁰ while the CompTox Chemicals Dashboard is a collection of 1 254 895 chemicals relevant for computational toxicity efforts.¹²¹ The coverage of HMDB, PCL, CompTox and the Blood Exposome database is explored in Figure 11A, which shows the

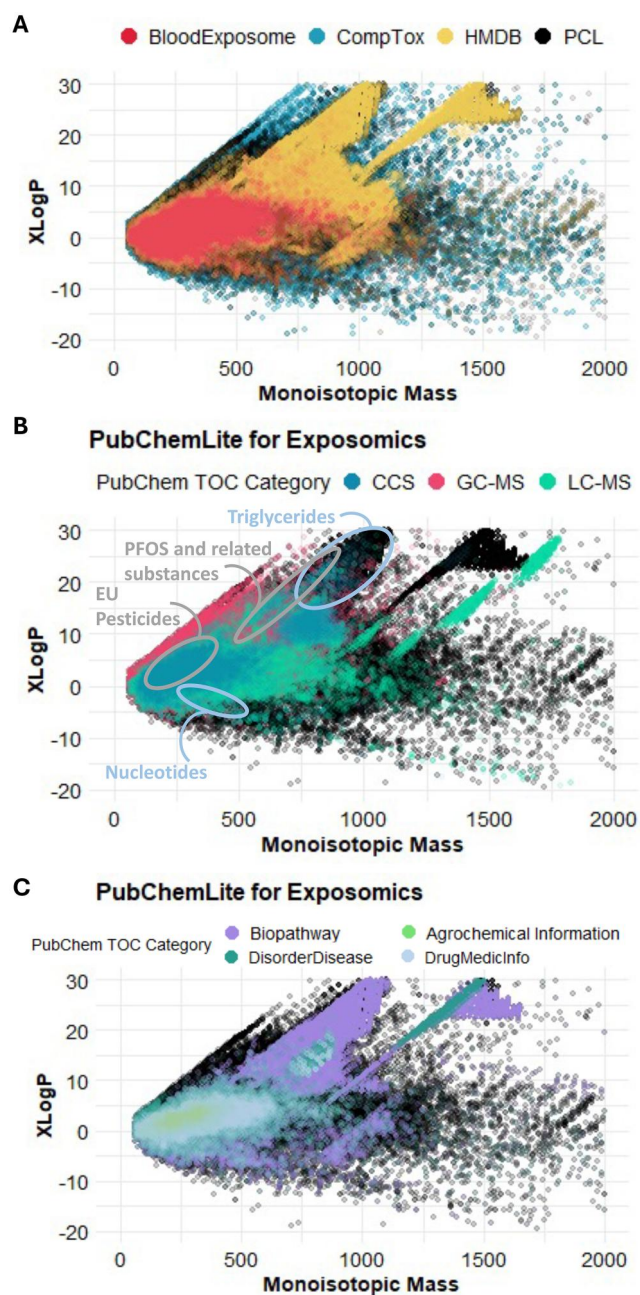


Figure 11. (A) Overlaid dot plot showing the coverage of the Blood Exposome Database, CompTox, the Human Metabolome Database (HMDB) and PubChemLite for Exposomics (PCL). (B) Overlaid dot plot of PubChemLite (PCL) displaying the coverage of the compounds with LC-MS, GC-MS as well as CCS information. (C) Overlaid plot of PCL displaying the compounds with Pathway information (Biopathway), associated disorders and diseases (DisorderDisease), agrochemical information and drug and medication information (DrugMedicInfo). R code for data visualization can be found in the GitLab repository (<https://gitlab.com/uniluxembourg/lcsb/eci/exposomics-plots>).¹²²

small, polar focus of the Blood Exposome database Versus the wider range of chemicals in HMDB (including lipid-based molecules) and the even wider range in CompTox and PCL that may not be detectable in humans, but may still influence health outcomes.

As discussed above (see section 2.4 and Figure 6), a combination of analytical techniques is necessary to capture the chemical

exposome. Figure 11B displays the chemicals within PCL that contain LC-MS and GC-MS spectral information, as well as CCS records in PubChem. This highlights the complementarity of the techniques, while GC-MS effectively identifies part of the non-polar side of the exposome, LC-MS covers a wider range of polar compounds. Notably, there are still many areas of the exposome that have no spectral information available, including very lipophilic molecules (middle and right top of the plot) and highly polar compounds (bottom right of the plot). Interestingly, CCS values are available for a good number of small molecules, providing an additional parameter to support compound annotation and identification in non-target HRMS. Several classes have been highlighted in Figure 11B, showing that both endogenous (triglycerides and nucleotides, in blue) as well as exogenous (perfluorooctanesulfonate [PFOS] and pesticides, in grey) can be captured with both LC-MS and GC-MS.

Figure 11C shows the overlaid plot of PCL along with four of the major PCL categories. Interestingly, the chemical space covered by the Biopathway category (Figure 11C) and HMDB (Figure 11A) are very similar, although both contain patches with little experimental data in PubChem (mass ~750-1750, XlogP 20-30). The Disorder and Disease category also overlaps more with the DrugMedicInfo category than perhaps expected, indicating that a large portion of this section may indicate drug availability and not additional disease insights. The Agrochemical information and DrugMedicInfo categories (Figure 11C) cluster in the left part of the plot, representing low molecular weight molecules with medium polarity, similar to the patterns observed in the Blood Exposome Database (Figure 11A).

One important fraction of the chemical exposome are transformation products, which are often overlooked.¹²³ These products, derived from biotic or abiotic reactions, can have toxicity and persistence profiles that differ significantly from their parent compounds and may even be more toxic in some cases.¹²⁴ Currently, structural elucidation of transformation products is typically performed by a combination of experimental and in silico approaches.¹²³ The “Transformations” section in PubChem, which contains documented parent-transformation product relationships from ChEMBL and the NORMAN Suspect List Exchange (NORMAN-SLE)¹²⁵ is included in PubChemLite,⁸⁷ while software such as patRoom can aid in identifying transformation products from databases or by in silico prediction.¹²⁶ Open source software such as BioTransformer¹²⁷ or ShinyTPs¹²⁴ help generate transformation product databases or suspect lists via in silico prediction and text mining, respectively.

Suspect screening approaches can be used to search for compounds of interest that are expected (“suspected”) to be in the samples, facilitated through the use of suspect lists. This can be considered a form of prioritization, as it reduces the number of compounds to be investigated and assists in discovering potentially relevant results. Although this approach was initially employed in environmental and toxicology sciences to expedite non-target screening, it has been increasingly applied in metabolomics and exposomics studies.^{2,15,21} Different platforms exist to exchange suspect lists, such as the NORMAN-SLE¹²⁵ and the CompTox Chemicals Dashboard.¹²⁸ Furthermore, specific lists of chemicals can be generated in PubChem by literature mining.¹⁵ Suspect screening can be facilitated by software such as patRoom,⁵⁸ which performs the automatic annotation of “suspects” based on pre-defined rules. The automatic assignment of identification levels is a key feature of patRoom, enhancing the reproducibility and leading to more transparent and comparable results.^{58,126} Other software options allowing suspect screening

include MZmine,⁵⁵ and Compound Discoverer (ThermoFisher). However, there is a risk in focusing too narrowly on distinct suspect classes in exposomics, depending on the study context, since recent studies have shown that significant chemicals for disease penetrance arise from many different information categories and suspect lists.³²

Statistics & biological interpretation

Non-target HRMS exposomics studies generate a large amount of data, necessitating a combination of univariate and multivariate statistical analysis to identify significant differences between groups.¹²⁹ Data pre-treatment (eg, normalization, scaling) is essential before applying statistics, as detailed in the next subsection.

Data pre-treatment

Data pre-treatment consists of transforming the HRMS feature list into a suitable state for subsequent statistical analysis.⁸³ This process aims to reduce the effects of technical and measurement errors while enhancing relevant biological variations.¹³⁰ Common pre-treatment methods include normalization, centering, scaling (eg, Pareto scaling), and transformations (eg, log and power).^{52,131} Multifunctional open tools such as MetaboAnalyst,²⁷ XCMS online,⁵⁷ and MS-DIAL⁵¹ implement some pre-treatment steps, while various statistical packages are available for data pre-treatment in C/C++, Java, R, and Python.¹³²

Normalization strategies are used to remove or correct unwanted systematic variations between samples, making them more comparable.^{53,83,130,133,134} Normalization in metabolomics and exposomics is more challenging compared to genomics and proteomics, due to the vast complexity of the chemical space.¹³⁴ For instance, there is no standard method to measure the total amount of chemicals in a sample to normalize in the way total protein amount is used in proteomics.^{66,133,134} Normalization can be performed either pre-acquisition or post-acquisition of HRMS data and is generally divided into sample-based and data-based approaches, reviewed recently elsewhere.¹³³

Univariate and multivariate statistical analysis

Univariate statistical tests aim to identify changes in individual molecules and work on the assumption of statistical independence.^{129,135,136} These approaches are commonly used to initially assess the potential relationships between exposures and disease phenotypes.²¹ Different tests are available to investigate differences across groups (or changes over time) and the choice depends on the data distribution and the experimental design (number of groups and type, ie, matched or unmatched).^{66,137} Univariate analysis can also be employed to investigate the association between specific small molecules of interest and other variables, such as known clinical parameters, environmental exposures, or microbial species, among others. For this purpose, various similarity tests, including Pearson’s correlation (parametric test) and Spearman’s correlation (non-parametric test) can be used.¹³⁸ Univariate statistics are also widely used within Exposome-Wide Association Studies (ExWAS).²¹ In an ExWAS, a large number of exposures are successively and independently tested for their association with a specific health outcome, using statistical approach analogous to Genome-wide association studies (GWAS).¹³⁹ However, although univariate statistics are widely employed to test individual chemicals, non-target exposomics studies do not generate univariate data, as chemicals are not

independent of each other, necessitating the use of multivariate statistics.¹³⁵

Multivariate statistical methods encompass both supervised and unsupervised methods.⁸³ Unsupervised methods are generally exploratory in nature, used to find patterns and generate hypotheses, while supervised methods are more confirmatory (hypothesis testing), widely used for biomarker identification, classification, and prediction.¹³² Unsupervised methods are an effective approach to explore and visualize the structure of the dataset. Principal Component Analysis (PCA)⁶⁶ stands out as one of the most widely employed methods, used to reduce the number of dimensions in the data, facilitating data exploration and visualization. It is often employed as a pre-processing step, to check the data quality, before applying a supervised method.^{52,140} Pooled QC samples that cluster tightly, ideally at the origin of the scores plot, are indicative of high-quality data¹⁴⁰ (Figure 12A).

Supervised methods are used to identify the independent variables (molecules) that best discriminate the groups under study (dependent variables). PLS-DA and orthogonal-PLS-DA (oPLS-DA) are some of the most commonly used methods. As shown in Figure 12B, they maximize the differences between groups. However, the main drawback of these supervised methods is their susceptibility to overfitting. Therefore, it is highly recommended to perform validation analysis to avoid finding false relationships and misinterpretation of the data.¹²⁹ Ideally, these classification models require splitting the dataset into training set, validation set, and test set. Thus, individual models are tested and evaluated on unique datasets and then applied to the entire dataset and/or to other datasets.^{129,135} This approach, however, is often limited when working with small sample sizes or cases such as rare diseases or mutations, which may prevent the dataset from being split effectively.

Special statistical considerations for exposomics studies

Missing data pose a particular challenge in exposomics studies, where exposures are often analysed jointly. As the number of included exposures (such as chemicals identified by HRMS) increases, the number of complete cases can decline rapidly. This issue is further exacerbated in longitudinal exposomics studies where participant numbers are typically highest at baseline, but dropouts accumulate over time, leading to progressively fewer observations and a greater proportion of missing data at later time points. Therefore, the use of imputation techniques, such as multiple imputation is recommended, as detailed by Santos et al.¹³⁹

Exposomics studies often involve analyzing samples collected at different time points and therefore processed in different analytical batches. This can introduce batch effects, which occur when quantitative measurements differ systematically between batches due to irrelevant factors.¹⁴¹ Batch effects can be caused by multiple sources at any step of the experimental workflow (Figure 4), from sample collection to sample preparation and data acquisition by the analytical platform (eg, LC-HRMS), which is often the primary source of variation.¹⁴¹ These batch effects are almost unavoidable when working with HRMS platforms. Thus, it is very important to identify the unwanted variations (eg, via PCA, as shown above in Figure 12A) and correct them. Although several strategies have been proposed to correct batch effects (eg, using ISs or pooled QC samples), it remains an active area of research with no standardized approach.^{141,142} A recent review discusses potential sources of batch effects in multiomics studies as well as different batch effect correction algorithms (BECAs).¹⁴³

Exposomics data typically contains many covariates such as confounders (eg, age, sex, medication), which must be accounted

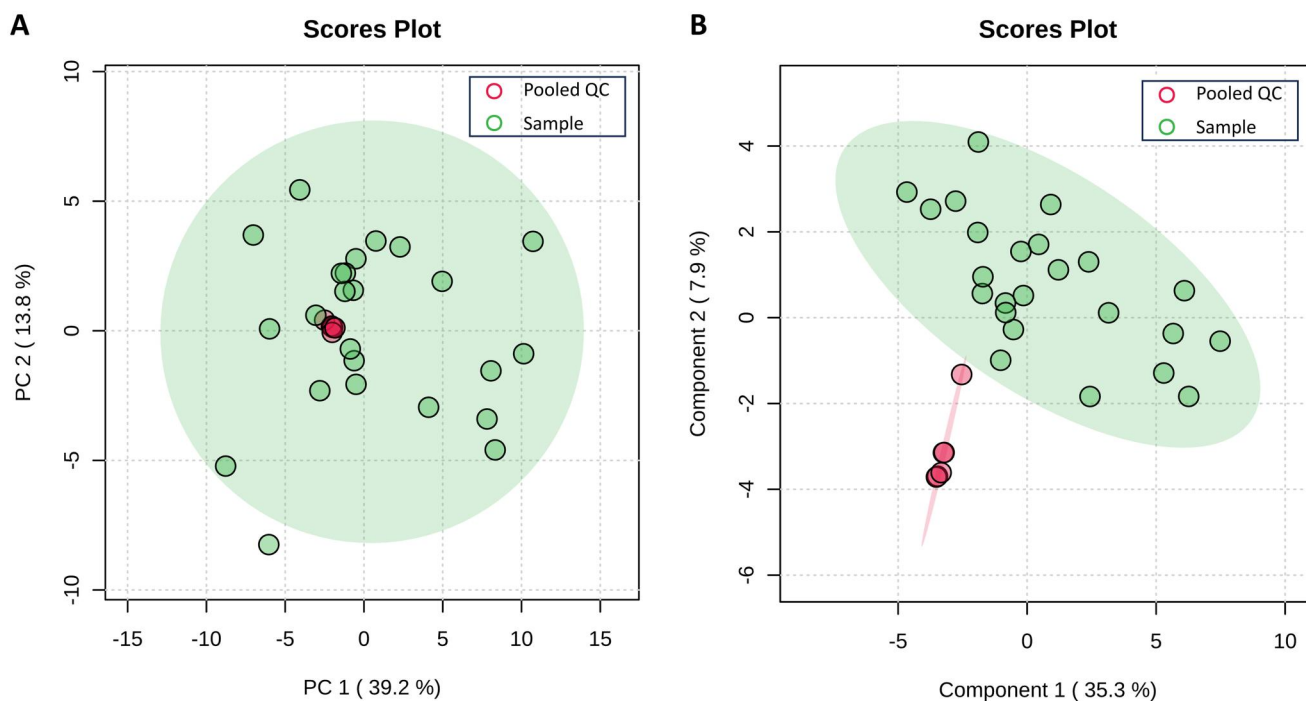


Figure 12. Multivariate statistical approaches applied in exposomics studies. (A) PCA scores plot of where all QC samples cluster tightly near the origin, which is indicative of quality data, affirming that the instrument variation was effectively corrected. PCA is widely used as a QA/QC tool to detect outliers and assess batch effects. (B) PLS-DA scores plot of the same dataset, included here to illustrate how supervised methods maximize the differences between predefined groups. Unlike PCA, PLS-DA is not used for QA/QC but is applied in downstream analyses such as classification, biomarker discovery, and hypothesis testing. Further details regarding QA/QC measures are discussed in Broadhurst et al.¹⁴⁰

in the statistical analyses. For this purpose, linear models can be employed to identify potential chemicals associated with a particular outcome (eg, Parkinson's disease or diabetes), while accounting for any number of covariates.¹⁴⁴

Network based approaches can help interpret the behaviour of chemicals or diseases that are related, and to provide insights into their mechanisms. Specifically, in exposomics, network approaches allow the identification of correlated exposures and potential biological changes associated with multiple exposures and health effects.^{139,144}

Mendelian randomization (MR) has become a valuable method in exposomics for assessing causal relationships between exposures and health outcomes. MR uses genetic variants, such as single nucleotide polymorphisms (SNPs), as instrumental variables to estimate the causal effect of a modifiable exposure on an outcome.¹⁴⁴ MR relies on the assumptions that genetic variants are strongly associated with the exposure, are not associated with confounders, and influence the outcome only through the exposure.^{21,144} Recent studies¹⁴⁵⁻¹⁴⁸ have successfully applied this approach in various exposomics contexts.

Since chemical exposures typically occur in complex mixtures rather than individually, various statistical methods can be applied to account for multiple exposures and their potential combined effects.^{21,149} These include, mixture analysis approaches, dimension reduction techniques or Bayesian model averaging. Further details can be found in Maitre et al.¹⁴⁹ In addition, machine learning approaches such as random forest, support vector machine, and gradient boosting can be used for biomarker selection, classification, and predictive modelling.^{21,144}

Biological interpretation

Biological interpretation of exposomics data is still a major challenge.¹⁵⁰ This is primarily due to the fact that the majority of the detected features by non-target HRMS remain unknown or, if annotated, little additional information is available. Another key challenge is distinguishing between exogenous and endogenous chemicals that reflect the biological response to the exposure,²⁰ where it is increasingly plausible that some detected chemicals originate from both exogenous and endogenous sources. Therefore, a multi-faceted approach is essential for interpreting small molecule omics data, encompassing different computational methods for statistical analysis, functional analysis, chemical classification, data integration, and data visualization, among others.¹³⁵

Network modelling and pathway mapping are key approaches for providing biological context to the data by enhancing the understanding of relationships between small molecules.^{151,152} MetaboAnalyst^{27,138} is a widely employed platform allowing data processing and interpretation through various functionality, including pathway, enrichment and network analysis. These approaches are valuable for data exploration and hypothesis generation, but the results require further validation. Pathway mapping is often limited by incomplete and manually curated pathway databases, leading to variability in results across different databases (eg, KEGG and MetaCyc),¹³⁵ while various metabolites can belong to different pathways. This analysis also excludes exogenous chemicals (ie, the chemical exposome). An alternative approach is ChemRICH,¹⁵⁰ which uses MeSH and Tanimoto substructure chemical similarity coefficients to cluster small molecules into non-overlapping chemical groups. In contrast with pathways analysis, ChemRICH sets have a self-contained size (based on the chemicals found in a particular study), therefore *P*-values do not rely on the size of a background

database, such as KEGG. Furthermore, the analysis can also place exposome chemicals into metabolite sets. However, results from ChemRICH cannot yet be directly integrated with genomics or proteomics results.^{135,150} Another strategy that does not require compound identification is mummichog.¹⁵³ This method uses as input peak lists, which are queried against a database to identify all the potential matches to metabolic pathways and networks.¹⁵² Linear models are useful for complex exposomics studies as they can account for covariates like age, sex, and occupation, helping to identify chemicals associated with specific metadata of interest. A recent review summarizes different computational methods including linear models with covariate adjustment, dimensionality reduction, and neural networks, among others, that support exposomics data analysis and interpretation.¹⁴⁴

Multi-omics approaches offer a more comprehensive understanding of the biological system by integrating small molecule omics with other omics data acquired from the same samples. Currently, there is a wide array of tools available for the integration of multi-omics data, such as mixOmics.¹⁵⁴ For instance, combining metabolomics with metagenomics can help elucidate the role of bacteria derived metabolites.¹⁵¹ In this context, tools like microbeMASST can help identifying the potential microbial origin of annotated chemicals by mapping known and unknown MS2 spectra to potential microbial producers.¹⁵⁵

From history to harmonization: Addressing challenges between the metabolome and exposome

Currently, the lack of harmonization not only in the annotation process but also throughout the entire non-target HRMS workflow leads to poorly comparable metabolomics and exposomics studies. As shown in Figure 13, metabolomics and exposomics are still young and rapidly growing fields of research, where workflow standardization is an ongoing task. To address this need, several initiatives have emerged separately for each discipline. In 2005, the MSI was formed,¹⁵⁶ proposing a series of minimum reporting standards for data analysis in metabolomics two years later.⁵² A network of global metabolomics repositories arose in 2015 (COordination of Standards in MetabOmicS [COSMOS]).¹⁵⁷ Since exposomics is a newer field, there is no "Exposomics Standard Initiative" yet. However, in the last few years various European, American and international initiatives have been developed, including The Human Early-Life Exposome (HELIX),²⁴ EXPOsOMICS,¹⁵⁸ The European Human Exposome Network (EHEN)¹⁵⁹ the Network for EXposomics in the United States (NEXUS),¹⁶⁰ HERCULES,¹⁶¹ Human Health Exposure Analysis Resource (HHEAR) network,¹⁶² and the International Human Exposome Network (IHEN).¹⁶³

QA/QC procedures should be implemented throughout the entire metabolomics and exposomics workflow, from sample preparation to data acquisition and data pre-processing, to ensure that the analysis is consistent, comparable, reproducible, precise and accurate. While well established QA/QC procedures in metabolomics serve as valuable benchmarks, such as the 2018 guidelines by Broadhurst et al.¹⁴⁰ (Figure 13), they may require adaptation for exposomics studies or research endeavours that integrate both metabolomics and exposomics. Notably, filtering features based on QC pooled samples may inadvertently exclude low-abundance small molecules, including exogenous compounds typically present at trace levels, which may be significant in certain groups or conditions. To address this issue, an

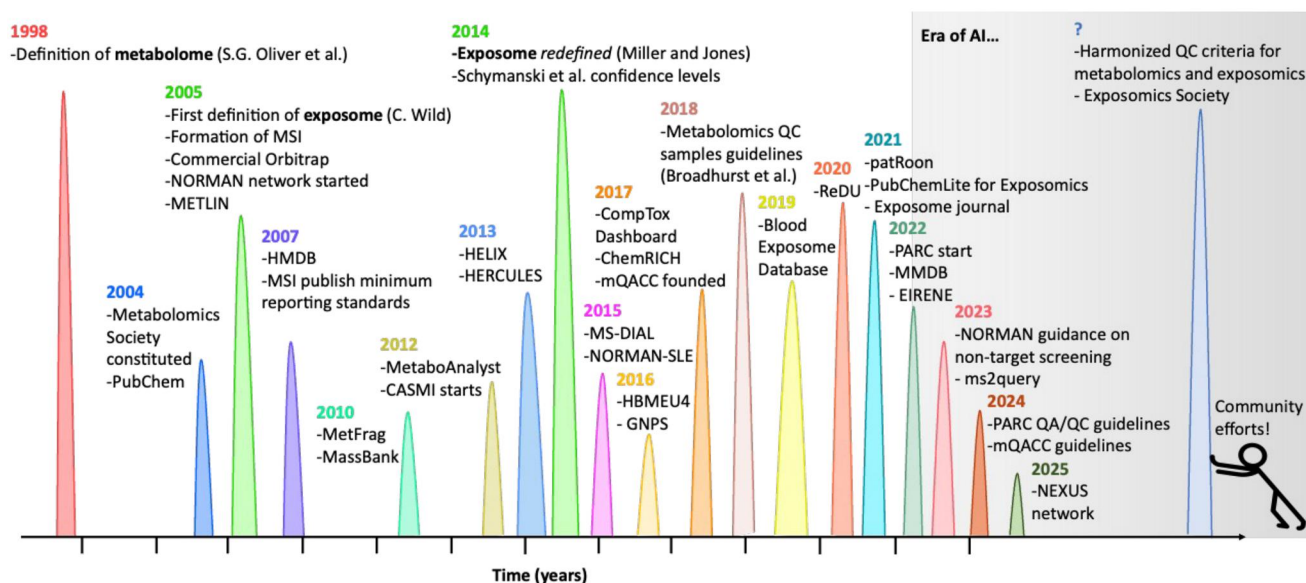


Figure 13. The metabolome and exposome timeline. Only the peak heights of the Metabolome (1998), Exposome (2014) and last peak (?) are intentionally emphasized for significance. The sizes of the other peaks were adjusted for aesthetic reasons and do not reflect their importance. A representative selection of tools is displayed here for illustrative purposes, and the authors acknowledge the presence of many other numerous tools and initiatives that have emerged in recent years. Adapted from.¹⁶⁴ Abbreviations: Artificial Intelligence; AI.

alternative approach involves preparing tailored pooled QC samples for specific study groups, such as cases and controls.¹⁶⁵ Currently, the usage of QC pooled samples varies greatly across studies accompanied by inconsistencies in the reporting of their preparation methods, as recently highlighted by Broeckling et al.¹⁶⁶ Therefore, it would be beneficial to establish guidelines for the preparation, usage, and reporting of QC samples in metabolomics and exposomics studies. A minimum set of QC measures as well as a standardized method for reporting them should be required in future studies. Table 1 summarizes a structure summary of recommended QA/QC for each stage of the non-target HRMS metabolomics and exposomics workflow, offering a practical starting point for harmonized QA/QC reporting in future exposomics studies.

The NORMAN network released their guidance on suspect and non-target screening for environmental monitoring in 2023,¹⁵ which has many parallel applications in the exposomics community. This guidance offers recommendations for all steps in a non-target screening experiment, from sample preparation to HRMS and data analysis and reporting. Additionally, the PARC initiative (environmental and biomonitoring communities) proposed their harmonized QA/QC procedures for data pre-processing of non-target and suspect screening LC-HRMS data in 2024.⁵⁴ These are very good starting points to be followed by the exposomics community, and these efforts indicate that the field is moving towards standardized workflows, which will hopefully emerge in the coming years.

Finally, artificial intelligence (AI) approaches (grey area of Figure 13), including Machine Learning (ML) approaches, are increasingly integrated into metabolomics and exposomics fields. AI can support various steps of the non-target HRMS workflow, such as feature prioritization, compound annotation, prediction modelling and pathways analysis. Although some of these approaches are still in their early stages and not yet widely and readily integrated into small molecule omics workflows, they are key starting points for supporting the new era of omics research.^{174,175} The shift from having most of the features

unannotated to actively predicting the structure and biological impact of thousands of molecules is essential for translating exposomics research findings into precision medicine.

Conclusions and future perspectives

Genes only explain a small fraction of chronic diseases, highlighting the need for a broader approach to better understand health and disease etiology.¹⁷⁶ In this context, the significant global health burden of environmental factors, such as the 9 million premature deaths attributed to pollution in 2019,¹⁷⁷ underscores the critical role of environmental factors. Consequently, exposomics has emerged as a complementary field to genomics to study environmental drivers of disease. However, translating these omics insights into clinical practice requires overcoming several challenges including the large number of unannotated features, the lack of harmonization, poor reproducibility across studies, and ethical, legal and social issues related to human exposomics data.¹⁷⁵ To address these challenges, substantial investment in exposomics research and the establishment of large-scale international consortia are urgently needed. Recent initiatives such as NEXUS are pivotal in this regard, fostering collaboration among scientists, policymakers and funders to accelerate innovation and standardize practices.

The increasing availability of open data (Figure 11) on relevant environmental chemicals and the matrices in which they are detected, is crucial to assess the suitability of our current analytical approaches to capture the chemical exposome. This data will enable a critical evaluation of whether existing methods can effectively capture relevant chemicals and downstream biological processes or whether additional analytical techniques, such as supercritical fluid chromatography¹⁷⁸ or two-dimensional chromatography, will be needed for a comprehensive characterization of the human exposome. However, while open data is important it is not always possible, particularly for human datasets, where privacy and ethical constraints may limit open

Table 1. Summary of some recommended QA/QC procedures for non-target-HRMS exposomics and metabolomics.

Workflow step	QA/QC procedures
Sample collection	<ul style="list-style-type: none"> • Develop detailed Standard Operational Protocols (SOPs) for sampling and storage. • Careful selection of appropriate sampling materials (eg, avoiding tubes with components such as phthalates and Polyethylene Glycol (PEG) that may interfere with the analysis). • Train personnel performing the sample collection.^{140,167} • Define specific acceptance/rejection criteria for samples upon arrival at the laboratory.¹⁶⁷ • Report collection method, storage temperature and number of freeze-thaw cycles, if applicable.
Sample preparation	<ul style="list-style-type: none"> • Use Internal Standards (IS). Ideally a mixture of multiple IS covering the range of the chemical space to be investigated should be added to each sample, at predetermined concentrations.^{15,140} • Use pooled QC samples, prepared by taking a small aliquot of each study sample and mixing it into homogenous pooled sample. Pooled QC samples are used to condition the analytical platform,¹⁴⁰ assess the analytical performance,² correct batch effects,^{2,140,166} support metabolite identification,¹⁶⁶ and filter low quality data,¹⁴⁰ although this may result in the loss of low abundant features such as pollutants. Thus, it is not recommended for exposomics studies, where low-abundant features may be relevant. Report how pooled samples were prepared and used. • Include blank samples to prevent false positives. Different blanks can be prepared such as extraction blanks (or process blanks) and system suitability blanks (instrument blanks). Further details can be found in the 2023 NORMAN guidance.¹⁵ • Report sample preparation protocol including the reconstitution solvent and volume, IS employed and concentrations.¹⁶⁸
Data acquisition	<ul style="list-style-type: none"> • Ensure instrument is calibrated before starting the analysis.¹⁶⁸ • Carefully plan the injection order to ensure the quality and precision of the analysis. A graphical example of a sequence is given in.¹⁴⁰ System suitability blanks (eg, MilliQ water) are injected to check that the instrument is working properly. System conditioning pooled QC samples can be injected to equilibrate the instrument before sample analysis. Injecting pooled QC samples throughout the sequence (eg, every 5 or 10 samples) helps measure the precision of the system (eg, stable retention times), correct for systematic bias, and support the later data pre-processing (eg, feature filtering).^{15,140} Randomization of all the samples across the sequence reduces systematic errors due to carryover.¹⁵ • Report the instrument configuration including the LC system, ionization source (ESI, APPI ...), and MS analyzer (eg, Orbitrap). • Report the LC-HRMS method details including mobile phase, gradient, flow rate, column temperature, sample volume injected, acquisition mode (eg, DDA), polarity (eg, positive), and <i>m/z</i> range, among others. Further details can be found in Viant et al.¹⁶⁸
Data pre-processing	<ul style="list-style-type: none"> • Following existing reporting guidelines for data pre-processing such as the METabolomics standaRds Initiative in Toxicology (MERIT)¹⁶⁸ and the Metabolomics Quality Assurance and Quality Control Consortium (mQACC).¹⁶⁹ • Share raw data in public repositories (eg, MetaboLights¹⁷⁰ and GNPS¹⁷¹) to ensure data reproducibility and reusability.⁵⁴ • Use open software computational workflows integrating data pre-processing, compound annotation, and statistical analysis, to enhance reproducibility by minimizing manual data curation.⁵⁴ • Use benchmark datasets to evaluate pre-processing algorithms.⁵⁴ Packages such as IPO⁶⁹ and Meta-Clean¹⁷² help optimize data pre-processing parameters. • Consider data pre-processing QA/QC guidelines proposed in.⁵⁴
Compound Annotation	<ul style="list-style-type: none"> • Assigning identification levels to the detected features is not trivial process and requires appropriate QA/QC procedures. Document the confidence level assigned to each feature and the annotation system employed.⁷² • Manual curation of all the annotations, coupled with the provision of evidence supporting the confidence level (eg, <i>m/z</i>, retention time, and MS2) is important for mitigating false positives and enhancing the reliability of the results.¹⁷³ • Report software employed, and parameters (eg, mass tolerance) used for compound annotation. The mass spectral library, suspect lists, and/or chemical databases utilized should be clearly stated, along with their respective versions.¹⁵ Spectra included in MS2 libraries should be curated by filtering, noise removal, and recalibration to ensure the quality of the reference spectra.⁷²
Statistical analysis	<ul style="list-style-type: none"> • Report all statistical methods transparently; the use of open source-software are preferred. • In-house generated code, such as R scripts, should be shared in public repositories (eg, GitHub) fostering transparency and reproducibility.¹⁶⁸ • Univariate statistics: report median and Relative Standard Deviation (RSD) of each feature across the pooled QC samples. • Multivariate statistics: performing and reporting PCA is recommended to confirm that the pooled QC samples cluster tightly, indicating high quality data.

access. To address these challenges, tiered or control data access, embargo periods and data use agreements can be employed.

Over the last years, the field has evolved significantly, bringing the future of exposomics and metabolomics closer to the well-established omics fields such as genomics, transcriptomics, and proteomics. The exposomics field is currently moving from proof-of-concept studies, with low sample sizes, to ExWAS, which include thousands of participants and a broad range of

environmental exposures and disease endpoints.¹⁷⁹ A key example is the HELIX study, which focuses on a cohort of more than 1000 mother-child pairs and demonstrated that early life exposures can lead to biological responses detectable through different omics layers.^{24,180} These findings highlight that exposomics can identify early life biomarkers of exposure, which can improve our understanding of health and disease status and promote public health policies. Another recent ExWAS explored

environmental exposures associated with aging in the UK Biobank.¹⁸¹ Twenty-five independent exposures were associated with mortality and proteomic aging. Notably, this study revealed a greater impact of the exposome on variation in mortality than polygenic risk scores.^{179,181} This progress is essential for a future where integrating data from different omics will allow for a more comprehensive assessment of an individual's disease risk, thereby facilitating personalized medicine. In the long term, the findings from exposomics studies will be instrumental in guiding policymakers to implement measures that protect future generations from harmful environmental exposures.

Glossary

Exposomics: Refers to the study of the exposome, ie, characterization of small molecules environmental derived and its transformation products within an entity (cell, tissue, or organism).

Feature: Peak (or signal) identified at a specific retention time, and *m/z*, containing spectral information such as MS1 and/or MS2, and intensity.

Metabolomics: Systematic and comprehensive study of low molecular weight molecules in a particular biological sample.^{39,40}

MS1: In Mass Spectrometry refers to the full scan information of the precursor ion (also known as parent ion) including the information of adducts, isotopic pattern and in-source fragments.

MS2: Also known as MS/MS refers to the fragmentation pattern of the precursor ion.

Non-target: This can refer to non-target study or non-target compound. *Non-target study* also known as untarget or untargeted, refers to discovery-based studies aiming to identify as many compounds as possible (known and unknown), whereas *non-target compound* refers to a compound for which no target or suspect identity can be assigned readily. The term not-target screening refers to the computational strategy searching for a broad range of compounds via tandem mass spectral libraries or database search.

Small molecule omics: Refers to the study of low molecule weight molecules in a particular biological or environmental sample. It encompasses metabolomics and exposomics.

Suspect: Suspects can refer to known compounds that are expected to be in the sample but with insufficient standard information to be identified. *Suspect screening* refers to a computational strategy searching for known chemicals that are expected to be in the sample.

Target: This can refer to target study or target compound. *Target study* also known as targeted, refers to validation-based studies focused on a limited number of known compounds. *Target compound* refers to known compound, preselected for the analysis, with reference standard data, including MS2 and retention time, available for the unequivocal identification.

Acknowledgments

The current and former members of the Environmental Cheminformatics (ECI) group at the LCB are acknowledged for their valuable inputs and discussion, some of which have guided and inspired the ideas discussed and figures presented in this manuscript. We wish to acknowledge that AI was used by BTA to assist with generating the code for a few visualizations (Figures 11 and 12), and for support to polish some writing and finding synonyms for contents of original text.

Author contributions

Begoña Talavera Andújar (Conceptualization [Equal], Data curation [Lead], Formal analysis [Equal], Investigation [Equal], Methodology [Equal], Software [Equal], Validation [Equal], Visualization [Lead], Writing—original draft [Lead], Writing—review & editing [Equal]), Emma Schymanski (Conceptualization [Equal], Data curation [Equal], Formal analysis [Equal], Funding acquisition [Lead], Investigation [Equal], Methodology [Equal], Project administration [Lead], Resources [Lead], Software [Equal], Supervision [Lead], Validation [Equal], Visualization [Equal], Writing—original draft [Equal], Writing—review & editing [Equal])

Funding

BTA acknowledges the support of the “Microbiomes in One Health” PhD training program, which is supported by the PRIDE doctoral research funding scheme (PRIDE/11823097) of the Luxembourg National Research Fund (FNR). ELS acknowledges funding support from the Luxembourg National Research Fund (FNR) for project A18/BM/12341006.

Conflict of interest

The authors have no conflict of interest. Emma Schymanski holds the position of Associate Editor and did not participate in the peer-review or make any editorial decisions for this manuscript.

Data availability

No new data were generated or analysed in support of this research. The code associated with Figure 11 (which retrieves relevant data from open repositories) is available in the ECI GitLab repository (<https://gitlab.com/uniluxembourg/lcsb/eci/exposomics-plots>).

Disclosure

This manuscript is based on unpublished portions of the PhD thesis of Begoña Talavera Andújar, defended on July 2024 and publicly available at the University of Luxembourg repository ORBilu (<https://orbilu.uni.lu/handle/10993/61575>).

References

1. Rappaport SM, Smith MT. Environment and disease risks. *Science*. 2010;330:460-461. <https://doi.org/10.1126/science.1192603>
2. David A, Chaker J, Price EJ, et al. Towards a comprehensive characterisation of the human internal chemical exposome: challenges and perspectives. *Environ Int*. 2021;156:106630. <https://doi.org/10.1016/j.envint.2021.106630>
3. Wild CP. Complementing the genome with an “exposome”: the outstanding challenge of environmental exposure measurement in molecular epidemiology. *Cancer Epidemiol Biomarkers Prev*. 2005;14:1847-1850. <https://doi.org/10.1158/5-9965.EPI-5-0456>
4. Miller GW, Jones DP. The nature of nurture: refining the definition of the exposome. *Toxicol Sci*. 2014;137:1-2. <https://doi.org/10.1093/toxsci/kft251>

5. Miller GW, Banbury Exposomics Consortium. Integrating exposomics into biomedicine. *Science*. 2025;388:356-358. <https://doi.org/10.1126/science.adr0544>
6. Safarlou CW, Jongsma KR, Vermeulen R. Reconceptualizing and defining exposomics within environmental health: expanding the scope of health research. *Environ Health Perspect*. 2024;132:95001. <https://doi.org/10.1289/EHP14509>
7. Vermeulen R. Human exposome research: potential, limitations and public policy implications. Think Tank, European Parliament, Accessed August 29, 2025. [https://www.europarl.europa.eu/thinktank/en/document/EPRS_STU\(2025\)765791](https://www.europarl.europa.eu/thinktank/en/document/EPRS_STU(2025)765791)
8. Wild CP. The exposome: from concept to utility. *Int J Epidemiol*. 2012;41:24-32. <https://doi.org/10.1093/ije/dyr236>
9. Vrijheid M. The exposome: a new paradigm to study the impact of environment on health. *Thorax*. 2014;69:876-878. <https://doi.org/10.1136/thoraxjnl-3-204949>
10. Guillien A, Ghosh M, Gille T, Dumas O. The exposome concept: how has it changed our understanding of environmental causes of chronic respiratory diseases? *Breathe (Sheff)*. 2023;19:230044. <https://doi.org/10.1183/20734735.4-2023>
11. La Cognata V, Morello G, Cavallaro S. Omics data and their integrative analysis to support stratified medicine in neurodegenerative diseases. *Int J Mol Sci*. 2021;22:4820. <https://doi.org/10.3390/ijms22094820>
12. Babu M, Snyder M. Multi-omics profiling for health. *Mol Cell Proteom*. 2023;22:100561. <https://doi.org/10.1016/j.mcpro.2023.100561>
13. Botas A, Campbell HM, Han X, Maletic-Savatic M. Metabolomics of neurodegenerative diseases. In: Hurley MJ (ed.) *International Review of Neurobiology*. Elsevier; 2015:53-80.
14. Roberts LD, Souza AL, Gerszten RE, Clish CB. Targeted metabolomics. *Curr Protoc Mol Biol*. 2012;98:30.2.1-30.2.24. <https://doi.org/10.1002/0471142727.mb3002s98>
15. Hollender J, Schymanski EL, Ahrens L, et al. NORMAN guidance on suspect and non-target screening in environmental monitoring. *Environ Sci Eur*. 2023;35:75. <https://doi.org/10.1186/s12302-3-00779-4>
16. Wishart DS. Metabolomics for investigating physiological and pathophysiological processes. *Physiol Rev*. 2019;99:1819-1875. <https://doi.org/10.1152/physrev.00035.2018>
17. Wishart DS. Advances in metabolite identification. *Bioanalysis*. 2011;3:1769-1782. <https://doi.org/10.4155/bio.11.155>
18. Metz TO, Baker ES, Schymanski EL, et al. Integrating ion mobility spectrometry into mass spectrometry-based exposome measurements: what can it add and how far can it go? *Bioanalysis*. 2017;9:81-98. <https://doi.org/10.4155/bio-2016-0244>
19. Walker DI, Valvi D, Rothman N, Lan Q, Miller GW, Jones DP. The metabolome: a key measure for exposome research in epidemiology. *Curr Epidemiol Rep*. 2019;6:93-103.
20. Balcells C, Xu Y, Gil-Solsona R, Maitre L, Gago-Ferrero P, Keun HC. Blurred lines: crossing the boundaries between the chemical exposome and the metabolome. *Curr Opin Chem Biol*. 2024;78:102407. <https://doi.org/10.1016/j.cbpa.2023.102407>
21. Lai Y, Koelmel JP, Walker DI, et al. High-resolution mass spectrometry for human exposomics: expanding chemical space coverage. *Environ Sci Technol*. 2024;58:12784-12822. <https://doi.org/10.1021/acs.est.4c01156>
22. Zhang P, Carlsten C, Chaleckis R, et al. Defining the scope of exposome studies and research needs from a multidisciplinary perspective. *Environ Sci Technol Lett*. 2021;8:839-852. <https://doi.org/10.1021/acs.estlett.1c00648>
23. All of Us Research Program, National Institutes of Health (NIH). In: Us Res. Program NIH. Accessed May 10, 2024, 2020. <https://allofus.nih.gov/future-health-begins-all-us>
24. Vrijheid M, Slama R, Robinson O, et al. The human early-life exposome (HELIX): project rationale and design. *Environ Health Perspect*. 2014;122:535-544. <https://doi.org/10.1289/ehp.1307204>
25. Barnes S, Benton HP, Casazza K, et al. Training in metabolomics research. I. Designing the experiment, collecting and extracting samples and generating metabolomics data. *J Mass Spectrom*. 2016;51:461-475. <https://doi.org/10.1002/jms.3782>
26. Faul F, Erdfelder E, Lang A-G, Buchner A. G*Power 3: a flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behav Res Methods*. 2007;39:175-191. <https://doi.org/10.3758/BF03193146>
27. Pang Z, Lu Y, Zhou G, et al. MetaboAnalyst 6.0: towards a unified platform for metabolomics data processing, analysis and interpretation. *Nucleic Acids Res*. 2024;52:W398-W406. <https://doi.org/10.1093/nar/gkae253>
28. Tarazona S, Balzano-Nogueira L, Gómez-Cabrero D, et al. Harmonization of quality metrics and power calculation in multi-omic studies. *Nat Commun*. 2020;11:3092. <https://doi.org/10.1038/s41467-020-16937-8>
29. Dennis KK, Marder E, Balshaw DM, et al. Biomonitoring in the era of the exposome. *Environ Health Perspect*. 2017;125:502-510. <https://doi.org/10.1289/EHP474>
30. Rappaport SM, Barupal DK, Wishart D, Vineis P, Scalbert A. The blood exposome and its role in discovering causes of disease. *Environ Health Perspect*. 2014;122:769-774. <https://doi.org/10.1289/ehp.1308015>
31. Jacobson TA, Kler JS, Bae Y, et al. A state-of-the-science review and guide for measuring environmental exposure biomarkers in dried blood spots. *J Expo Sci Environ Epidemiol*. 2023;33:505-523. <https://doi.org/10.1038/s41370-022-00460-7>
32. Talavera Andújar B, Pereira SL, Bhanu Busi S, et al. Exploring environmental modifiers of LRRK2-associated Parkinson's disease penetrance: an exposomics and metagenomics pilot study on household dust. *Environ Int*. 2024;194:109151. <https://doi.org/10.1016/j.envint.2024.109151>
33. Running LS, Kordas K, Aga DS. Use of wristbands to measure exposure to environmental pollutants in children: recent advances and future directions. *Curr Opin Environ Sci Health*. 2023;32:100450. <https://doi.org/10.1016/j.coesh.2023.100450>
34. Waclawik M, Rodzaj W, Wielgomas B. Silicone wristbands in exposure assessment: analytical considerations and comparison with other approaches. *Int J Environ Res Public Health*. 2022;19:1935. <https://doi.org/10.3390/ijerph19041935>
35. Schrimpe-Rutledge AC, Codreanu SG, Sherrod SD, McLean JA. Untargeted metabolomics strategies—challenges and emerging directions. *J Am Soc Mass Spectrom*. 2016;27:1897-1905. <https://doi.org/10.1007/s13361-6-1469-y>
36. Gu Y, Peach JT, Warth B. Sample preparation strategies for mass spectrometry analysis in human exposome research: current status and future perspectives. *TrAC Trends Anal Chem*. 2023;166:117151. <https://doi.org/10.1016/j.trac.2023.117151>
37. Vuckovic D. Current trends and challenges in sample preparation for global metabolomics using liquid chromatography-mass spectrometry. *Anal Bioanal Chem*. 2012;403:1523-1548. <https://doi.org/10.1007/s00216-2-6039-y>
38. Alosekh S, Aharoni A, Brotman Y, et al. Mass spectrometry-based metabolomics: a guide for annotation, quantification and best reporting practices. *Nat Methods*. 2021;18:747-756. <https://doi.org/10.1038/s41592-021-01197-1>
39. Zeki ÖC, Eylem CC, Reçber T, Kir S, Nemitlu E. Integration of GC-MS and LC-MS for untargeted metabolomics profiling. *J Pharm Biomed Anal*. 2020;190:113509. <https://doi.org/10.1016/j.jpba.2020.113509>

40. Fiehn O. Metabolomics by gas chromatography–mass spectrometry: combined targeted and untargeted profiling. *Curr Protoc Mol Biol*. 2016;114:30.4.1-30.4.32. <https://doi.org/10.1002/0471142727.mb3004s114>
41. Zhang A, Sun H, Wang P, Han Y, Wang X. Modern analytical techniques in metabolomics analysis. *Analyst*. 2012;137:293-300. <https://doi.org/10.1039/C1AN15605E>
42. Wang JH, Byun J, Pennathur S. Analytical approaches to metabolomics and applications to systems biology. *Semin Nephrol*. 2010;30:500-511. <https://doi.org/10.1016/j.semnephrol.2010.07.007>
43. Zhang X, Quinn K, Cruickshank-Quinn C, Reisdorph R, Reisdorph N. The application of ion mobility mass spectrometry to metabolomics. *Curr Opin Chem Biol*. 2018;42:60-66. <https://doi.org/10.1016/j.cbpa.2017.11.001>
44. Izquierdo-Sandoval D, Fabregat-Safont D, Lacalle-Bergeron L, Sancho JV, Hernández F, Portoles T. Benefits of ion mobility separation in GC-APCI-HRMS screening: from the construction of a CCS library to the application to real-world samples. *Anal Chem*. 2022;94:9040-9047. <https://doi.org/10.1021/acs.analchem.2c01118>
45. Ropartz D, Fanuel M, Ujma J, Palmer M, Giles K, Rogniaux H. Structure determination of large isomeric oligosaccharides of natural origin through multipass and multistage cyclic traveling-wave ion mobility mass spectrometry. *Anal Chem*. 2019;91:12030-12037. <https://doi.org/10.1021/acs.analchem.9b03036>
46. Collins SL, Koo I, Peters JM, Smith PB, Patterson AD. Current challenges and recent developments in mass spectrometry–based metabolomics. *Annu Rev Anal Chem (Palo Alto Calif)*. 2021;14:467-487. <https://doi.org/10.1146/annurev-anchem-091620-015205>
47. Kind T, Tsugawa H, Cajka T, et al. Identification of small molecules using accurate mass MS/MS search. *Mass Spectrom Rev*. 2018;37:513-532. <https://doi.org/10.1002/mas.21535>
48. Guo J, Huan T. Comparison of full-scan, data-dependent, and data-independent acquisition modes in liquid chromatography–mass spectrometry based untargeted metabolomics. *Anal Chem*. 2020;92:8072-8080. <https://doi.org/10.1021/acs.analchem.9b05135>
49. Yang Y, Yang L, Zheng M, Cao D, Liu G. Data acquisition methods for non-targeted screening in environmental analysis. *TrAC Trends Anal Chem*. 2023;160:116966. <https://doi.org/10.1016/j.trac.2023.116966>
50. Koelmel JP, Kroeger NM, Gill EL, et al. Expanding lipidome coverage using LC-MS/MS data-dependent acquisition with automated exclusion list generation. *J Am Soc Mass Spectrom*. 2017;28:908-917. <https://doi.org/10.1007/s13361-7-1608-0>
51. Tsugawa H, Cajka T, Kind T, et al. MS-DIAL: data-independent MS/MS deconvolution for comprehensive metabolome analysis. *Nat Methods*. 2015;12:523-526. <https://doi.org/10.1038/nmeth.3393>
52. Goodacre R, Broadhurst D, Smilde AK, et al. Proposed minimum reporting standards for data analysis in metabolomics. *Metabolomics*. 2007;3:231-241. <https://doi.org/10.1007/s11306-007-0081-3>
53. Katajamaa M, Orešič M. Data processing for mass spectrometry-based metabolomics. *J Chromatogr A*. 2007;1158:318-328. <https://doi.org/10.1016/j.chroma.2007.04.021>
54. Lennon S, Chaker J, Price EJ, et al. Harmonized quality assurance/quality control provisions to assess completeness and robustness of MS1 data preprocessing for LC-HRMS-based suspect screening and non-targeted analysis. *TrAC Trends Anal Chem*. 2024;174:117674. <https://doi.org/10.1016/j.trac.2024.117674>
55. Schmid R, Heuckeroth S, Korf A, et al. Integrative analysis of multimodal mass spectrometry data in MZmine 3. *Nat Biotechnol*. 2023;41:447-449. <https://doi.org/10.1038/s41587-023-01690-2>
56. Sturm M, Bertsch A, Gröpl C, et al. OpenMS—an open-source software framework for mass spectrometry. *BMC Bioinformatics*. 2008;9:163. <https://doi.org/10.1186/1471-2105-9-163>
57. 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. <https://doi.org/10.1021/ac300698c>
58. Helmus R, Ter Laak TL, Van Wezel AP, De Voogt P, Schymanski EL. patRoön: open source software platform for environmental mass spectrometry based non-target screening. *J Cheminform*. 2021;13:1. <https://doi.org/10.1186/s13321-020-00477-w>
59. Misra BB. New software tools, databases, and resources in metabolomics: updates from 2020. *Metabolomics*. 2021;17:49. <https://doi.org/10.1007/s11306-021-01796-1>
60. Renner G, Reuschenbach M. Critical review on data processing algorithms in non-target screening: challenges and opportunities to improve result comparability. *Anal Bioanal Chem*. 2023;415:4111-4123. <https://doi.org/10.1007/s00216-3-04776-7>
61. Chambers MC, Maclean B, Burke R, et al. A cross-platform toolkit for mass spectrometry and proteomics. *Nat Biotechnol*. 2012;30:918-920. <https://doi.org/10.1038/nbt.2377>
62. ProteoWizard: Home. Accessed April 2, 2024, <https://proteowizard.sourceforge.io/index.html>
63. Holman JD, Tabb DL, Mallick P. Employing ProteoWizard to convert raw mass spectrometry data. *Curr Protoc Bioinformatics*. 2014;46:13.24.1-13.24.9. <https://doi.org/10.1002/0471250953.bi1324s46>
64. Karaman I, Climaco Pinto R, Graça G. Metabolomics data preprocessing: from raw data to features for statistical analysis. In: *Comprehensive Analytical Chemistry*. Elsevier; 2018:197-225.
65. Boccard J, Veuthey J, Rudaz S. Knowledge discovery in metabolomics: an overview of MS data handling. *J Sep Sci*. 2010;33:290-304. <https://doi.org/10.1002/jssc.200900609>
66. Santamaria G, Pinto FR. Bioinformatic analysis of metabolomic data: from raw spectra to biological insight. *BioChem*. 2024;4:90-114. <https://doi.org/10.3390/biochem4020005>
67. Savitzky A, Golay MJE. Smoothing and differentiation of data by simplified least squares procedures. *Anal Chem*. 1964;36:1627-1639. <https://doi.org/10.1021/ac60214a047>
68. Lommen A. MetAlign: interface-driven, versatile metabolomics tool for hyphenated full-scan mass spectrometry data preprocessing. *Anal Chem*. 2009;81:3079-3086. <https://doi.org/10.1021/ac900036d>
69. Libiseller G, Dvorzak M, Kleb U, et al. IPO: a tool for automated optimization of XCMS parameters. *BMC Bioinformatics*. 2015;16:118. <https://doi.org/10.1186/s12859-015-0562-8>
70. Blaženović I, Kind T, Ji J, Fiehn O. Software tools and approaches for compound identification of LC-MS/MS data in metabolomics. *Metabolites*. 2018;8:31. <https://doi.org/10.3390/metabo8020031>
71. Chaleckis R, Meister I, Zhang P, Wheelock CE. Challenges, progress and promises of metabolite annotation for LC-MS-based metabolomics. *Curr Opin Biotechnol*. 2019;55:44-50. <https://doi.org/10.1016/j.copbio.2018.07.010>
72. Oberacher H, Sasse M, Antignac J-P, et al. A European proposal for quality control and quality assurance of tandem mass spectral libraries. *Environ Sci Eur*. 2020;32:43. <https://doi.org/10.1186/s12302-020-00314-9>
73. de Jonge NF, Mildau K, Meijer D, et al. Good practices and recommendations for using and benchmarking computational metabolomics metabolite annotation tools. *Metabolomics*. 2022;18:103. <https://doi.org/10.1007/s11306-022-01963-y>

74. Giera M, Aisporna A, Uritboonthai W, Siuzdak G. The hidden impact of in-source fragmentation in metabolic and chemical mass spectrometry data interpretation. *Nat Metab.* 2024;6:1647-1648. <https://doi.org/10.1038/s42255-024-01076-x>
75. El Abiead Y, Rutz A, Zuffa S, et al. Discovery of metabolites pre-vents amid in-source fragmentation. *Nat Metab.* 2025;7:435-437. <https://doi.org/10.1038/s42255-025-01239-4>
76. Schmid R, Petras D, Nothias L-F, et al. Ion identity molecular networking for mass spectrometry-based metabolomics in the GNPS environment. *Nat Commun.* 2021;12:3832. <https://doi.org/10.1038/s41467-021-23953-9>
77. Xu Y-F, Lu W, Rabinowitz JD. Avoiding misannotation of in-source fragmentation products as cellular metabolites in liquid chromatography-mass spectrometry-based metabolomics. *Anal Chem.* 2015;87:2273-2281. <https://doi.org/10.1021/ac504118y>
78. Xue J, Domingo-Almenara X, Guijas C, et al. Enhanced in-source fragmentation annotation enables novel data independent acquisition and autonomous METLIN molecular identification. *Anal Chem.* 2020;92:6051-6059. <https://doi.org/10.1021/acs.analchem.0c00409>
79. Krier J, Singh RR, Kondić T, et al. Discovering pesticides and their TPs in Luxembourg waters using open cheminformatics approaches. *Environ Int.* 2022;158:106885. <https://doi.org/10.1016/j.envint.2021.106885>
80. Hollender J, Bourgin M, Fenner KB, et al. Exploring the behaviour of emerging contaminants in the water cycle using the capabilities of high resolution mass spectrometry. *Chimia (Aarau).* 2014;68:793-798. <https://doi.org/10.2533/chimia.2014.793>
81. Schymanski EL, Jeon J, Gulde R, et al. Identifying small molecules via high resolution mass spectrometry: communicating confidence. *Environ Sci Technol.* 2014;48:2097-2098. <https://doi.org/10.1021/es5002105>
82. Koelmel JP, Kroeger NM, Ulmer CZ, et al. LipidMatch: an automated workflow for rule-based lipid identification using untargeted high-resolution tandem mass spectrometry data. *BMC Bioinformatics.* 2017;18:331. <https://doi.org/10.1186/s12859-017-1744-3>
83. Sumner LW, Amberg A, Barrett D, et al. Proposed minimum reporting standards for chemical analysis. *Metabolomics.* 2007;3:211-221. <https://doi.org/10.1007/s11306-7-0082-2>
84. Celma A, Sancho JV, Schymanski EL, et al. Improving target and suspect screening high-resolution mass spectrometry workflows in environmental analysis by ion mobility separation. *Environ Sci Technol.* 2020;54:15120-15131. <https://doi.org/10.1021/acs.est.0c05713>
85. Charbonnet JA, McDonough CA, Xiao F, et al. Communicating confidence of per- and polyfluoroalkyl substance identification via high-resolution mass spectrometry. *Environ Sci Technol Lett.* 2022;9:473-481. <https://doi.org/10.1021/acs.estlett.2c00206>
86. Koelmel JP, Xie H, Price EJ, et al. An actionable annotation scoring framework for gas chromatography-high-resolution mass spectrometry. *Exposome.* 2022;2:osac007. <https://doi.org/10.1093/exposome/osac007>
87. Schymanski EL, Kondić T, Neumann S, Thiessen PA, Zhang J, Bolton EE. Empowering large chemical knowledge bases for exposomics: PubChemLite meets MetFrag. *J Cheminform.* 2021;13:19. <https://doi.org/10.1186/s13321-021-00489-0>
88. Creek DJ, Dunn WB, Fiehn O, et al. Metabolite identification: are you sure? And how do your peers gauge your confidence? *Metabolomics.* 2014;10:350-353. <https://doi.org/10.1007/s11306-014-0656-8>
89. Liebisch G, Vizcaino JA, Köfeler H, et al. Shorthand notation for lipid structures derived from mass spectrometry. *J Lipid Res.* 2013;54:1523-1530. <https://doi.org/10.1194/jlr.M033506>
90. Scheubert K, Hufsky F, Petras D, et al. Significance estimation for large scale metabolomics annotations by spectral matching. *Nat Commun.* 2017;8:1494. <https://doi.org/10.1038/s41467-017-01318-5>
91. Metz TO, Chang CH, Gautam V, et al. Introducing "identification probability" for automated and transferable assessment of metabolite identification confidence in metabolomics and related studies. *Anal Chem.* 2025;97:1-11. <https://doi.org/10.1021/acs.analchem.4c04060>
92. Talavera Andújar B, Aurich D, Aho VTE, et al. Studying the Parkinson's disease metabolome and exposome in biological samples through different analytical and cheminformatics approaches: a pilot study. *Anal Bioanal Chem.* 2022;414:7399-7419. <https://doi.org/10.1007/s00216-2-04207-z>
93. Alygizakis N, Lestremau F, Gago-Ferrero P, et al. Towards a harmonized identification scoring system in LC-HRMS/MS based non-target screening (NTS) of emerging contaminants. *TrAC Trends Anal Chem.* 2023;159:116944. <https://doi.org/10.1016/j.trac.2023.116944>
94. Boatman AK, Chappel JR, Kirkwood-Donelson KI, et al. Updated guidance for communicating PFAS identification confidence with ion mobility spectrometry. *Environ Sci Technol.* 2025;59:17711-17721. <https://doi.org/10.1021/acs.est.5c01354>
95. Bittremieux W, Wang M, Dorrestein PC. The critical role that spectral libraries play in capturing the metabolomics community knowledge. *Metabolomics.* 2022;18:94. <https://doi.org/10.1007/s11306-022-01947-y>
96. Horai H, Arita M, Kanaya S, et al. MassBank: a public repository for sharing mass spectral data for life sciences. *J Mass Spectrom.* 2010;45:703-714. <https://doi.org/10.1002/jms.1777>
97. MassBank of North America. Accessed September 11, 2023. <https://mona.fiehnlab.ucdavis.edu/>
98. 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. <https://doi.org/10.1038/nbt.3597>
99. mzCloud—Advanced Mass Spectral Database. Accessed April 28, 2024. <https://www.mzcloud.org/>
100. Guijas C, Montenegro-Burke JR, Domingo-Almenara X, et al. METLIN: a technology platform for identifying knowns and unknowns. *Anal Chem.* 2018;90:3156-3164. <https://doi.org/10.1021/acs.analchem.7b04424>
101. Stein SE, Scott DR. Optimization and testing of mass spectral library search algorithms for compound identification. *J Am Soc Mass Spectrom.* 1994;5:5:859-866. [https://doi.org/10.1016/1044\(94\)87009-8](https://doi.org/10.1016/1044(94)87009-8)
102. PubChem Classification Browser. Accessed May 19, 2024. <https://pubchem.ncbi.nlm.nih.gov/classification/#hid=72>
103. Li Y, Kind T, Folz J, Vaniya A, Mehta SS, Fiehn O. Spectral entropy outperforms MS/MS dot product similarity for small-molecule compound identification. *Nat Methods.* 2021;18:1524-1531. <https://doi.org/10.1038/s41592-021-01331-z>
104. Ruttkies C, Schymanski EL, Wolf S, Hollender J, Neumann S. MetFrag relaunched: incorporating strategies beyond in silico fragmentation. *J Cheminform.* 2016;8:3. <https://doi.org/10.1186/s13321-016-0115-9>
105. Tsugawa H, Kind T, Nakabayashi R, et al. Hydrogen rearrangement rules: computational MS/MS fragmentation and structure elucidation using MS-FINDER software. *Anal Chem.* 2016;88:7946-7958. <https://doi.org/10.1021/acs.analchem.6b00770>

106. Allen F, Pon A, Wilson M, Greiner R, Wishart D. CFM-ID: a web server for annotation, spectrum prediction and metabolite identification from tandem mass spectra. *Nucleic Acids Res.* 2014;42:W94-W99. <https://doi.org/10.1093/nar/gku436>
107. Dührkop K, Shen H, Meusel M, Rousu J, Böcker S. Searching molecular structure databases with tandem mass spectra using CSI: FingerID. *Proc Natl Acad Sci U S A.* 2015;112:12580-12585. <https://doi.org/10.1073/pnas.1509788112>
108. Dührkop K, Fleischauer M, Ludwig M, et al. SIRIUS 4: a rapid tool for turning tandem mass spectra into metabolite structure information. *Nat Methods.* 2019;16:299-302. <https://doi.org/10.1038/s41592-019-0344-8>
109. Kind T, Liu K-H, Lee DY, DeFelice B, Meissen JK, Fiehn O. LipidBlast in silico tandem mass spectrometry database for lipid identification. *Nat Methods.* 2013;10:755-758. <https://doi.org/10.1038/nmeth.2551>
110. Wishart DS, Feunang YD, Marcu A, et al. HMDB 4.0: the human metabolome database for 2018. *Nucleic Acids Res.* 2018;46:D608-D617. <https://doi.org/10.1093/nar/gkx1089>
111. CAS REGISTRY, CAS. Accessed April 30, 2024. <https://www.cas.org/cas-data/cas-registry>
112. Pence HE, Williams A. ChemSpider: an online chemical information resource. *J Chem Educ.* 2010;87:1123-1124. <https://doi.org/10.1021/ed100697w>
113. Bolton EE, Wang Y, Thiessen PA, Bryant SH. Chapter 12 - PubChem: integrated platform of small molecules and biological activities. In: Wheeler RA, Spellmeyer DC, eds. *Annual Reports in Computational Chemistry*. Elsevier; 2008:217-241.
114. Bolton E, Schymanski E, Kondic T, Thiessen P, Zhang J. PubChemLite for Exposomics. 2020.
115. Elapavalore A, Ross DH, Grouès V, et al. PubChemLite Plus Collision Cross Section (CCS) values for enhanced interpretation of nontarget environmental data. *Environ Sci Technol Lett.* 2025;12:166-174. <https://doi.org/10.1021/acs.estlett.4c01003>
116. Barupal DK, Fiehn O. Generating the blood exposome database using a comprehensive text mining and database fusion approach. *Environ Health Perspect.* 2019;127:97008. <https://doi.org/10.1289/EHP4713>
117. Kanehisa M, Goto S, Hattori M, et al. From genomics to chemical genomics: new developments in KEGG. *Nucleic Acids Res.* 2006;34:D354-D357. <https://doi.org/10.1093/nar/gkj102>
118. Caspi R, Foerster H, Fulcher CA, et al. The MetaCyc Database of metabolic pathways and enzymes and the BioCyc collection of Pathway/Genome Databases. *Nucleic Acids Res.* 2008;36:D623-D631. <https://doi.org/10.1093/nar/gkm900>
119. Sud M, Fahy E, Cotter D, et al. LMSD: LIPID MAPS structure database. *Nucleic Acids Res.* 2007;35:D527-D532. <https://doi.org/10.1093/nar/gkl838>
120. Wishart DS, Oler E, Peters H, et al. MiMeDB: the human microbial metabolome database. *Nucleic Acids Res.* 2023;51:D611-D620. <https://doi.org/10.1093/nar/gkac868>
121. Williams AJ, Grulke CM, Edwards J, et al. The CompTox Chemistry Dashboard: a community data resource for environmental chemistry. *J Cheminform.* 2017;9:61. <https://doi.org/10.1186/s13321-017-0247-6>
122. Talavera Andújar B. uniluxembourg/LCSB/Environmental Cheminformatics/Exposomics plots. GitLab. In: GitLab. Accessed July 21, 2025. <https://gitlab.com/uniluxembourg/lcsb/eci/exposomics-plots>
123. Samanipour S, Barron LP, van Herwerden D, Praetorius A, Thomas KV, O'Brien JW. Exploring the chemical space of the exposome: how far have we gone? *JACS Au.* 2024;4:2412-2425. <https://doi.org/10.1021/jacsau.4c00220>
124. Palm EH, Chirsir P, Krier J, et al. ShinyTPs: curating transformation products from text mining results. *Environ Sci Technol Lett.* 2023;10:865-871. <https://doi.org/10.1021/acs.estlett.3c00537>
125. NORMAN-SLE. <https://www.norman-network.com/nds/SLE/>
126. Helmus R, Van De Velde B, Brunner AM, Ter Laak TL, Van Wezel AP, Schymanski EL. patRoom 2.0: improved non-target analysis workflows including automated transformation product screening. *Joss.* 2022;7:4029. <https://doi.org/10.21105/joss.04029>
127. Djoumbou-Feunang Y, Fiamoncini J, Gil-de-la-Fuente A, Greiner R, Manach C, Wishart DS. BioTransformer: a comprehensive computational tool for small molecule metabolism prediction and metabolite identification. *J Cheminform.* 2019;11:2. <https://doi.org/10.1186/s13321-018-0324-5>
128. CompTox Chemicals Dashboard Chemical Lists. Accessed April 30, 2024. <https://comptox.epa.gov/dashboard/chemical-lists>
129. Gertsman I, Barshop BA. Promises and pitfalls of untargeted metabolomics. *J Inherit Metab Dis.* 2018;41:355-366. <https://doi.org/10.1007/s10545-017-0130-7>
130. Karaman I. Preprocessing and pretreatment of metabolomics data for statistical analysis. In: Sussulini A, ed. *Metabolomics: From Fundamentals to Clinical Applications*. Springer International Publishing; 2017:145-161.
131. van den Berg RA, Hoefsloot HC, Westerhuis JA, Smilde AK, van der Werf MJ. Centering, scaling, and transformations: improving the biological information content of metabolomics data. *BMC Genomics.* 2006;7:142. <https://doi.org/10.1186/1471-2164-7-142>
132. Ren S, Hinzman AA, Kang EL, Szczesniak RD, Lu LJ. Computational and statistical analysis of metabolomics data. *Metabolomics.* 2015;11:1492-1513. <https://doi.org/10.1007/s11306-015-0823-6>
133. Misra BB. Data normalization strategies in metabolomics: current challenges, approaches, and tools. *Eur J Mass Spectrom (Chichester).* 2020;26:165-174. <https://doi.org/10.1177/1469066720918446>
134. Wu Y, Li L. Sample normalization methods in quantitative metabolomics. *J Chromatogr A.* 2016;1430:80-95. <https://doi.org/10.1016/j.chroma.2015.12.007>
135. Barupal DK, Fan S, Fiehn O. Integrating bioinformatics approaches for a comprehensive interpretation of metabolomics datasets. *Curr Opin Biotechnol.* 2018;54:1-9. <https://doi.org/10.1016/j.copbio.2018.01.010>
136. Chen Y, Li E-M, Xu L-Y. Guide to metabolomics analysis: a bioinformatics workflow. *Metabolites.* 2022;12:357. <https://doi.org/10.3390/metabo12040357>
137. Vinaixa M, Samino S, Saez I, Duran J, Guinovart JJ, Yanes O. A guideline to univariate statistical analysis for LC/MS-based untargeted metabolomics-derived data. *Metabolites.* 2012;2:775-795. <https://doi.org/10.3390/metabo2040775>
138. Pang Z, Zhou G, Ewald J, et al. Using MetaboAnalyst 5.0 for LC-HRMS spectra processing, multi-omics integration and covariate adjustment of global metabolomics data. *Nat Protoc.* 2022;17:1735-1761. <https://doi.org/10.1038/s41596-022-00710-w>
139. Santos S, Maitre L, Warembourg C, et al. Applying the exposome concept in birth cohort research: a review of statistical approaches. *Eur J Epidemiol.* 2020;35:193-204. <https://doi.org/10.1007/s10654-020-00625-4>
140. Broadhurst D, Goodacre R, Reinke SN, et al. Guidelines and considerations for the use of system suitability and quality control samples in mass spectrometry assays applied in untargeted clinical metabolomic studies. *Metabolomics.* 2018;14:72. <https://doi.org/10.1007/s11306-018-1367-3>

141. Han W, Li L. Evaluating and minimizing batch effects in metabolomics. *Mass Spectrom Rev.* 2022;41:421-442. <https://doi.org/10.1002/mas.21672>
142. Fan S, Kind T, Cajka T, et al. Systematic error removal using random forest for normalizing large-scale untargeted lipidomics data. *Anal Chem.* 2019;91:3590-3596. <https://doi.org/10.1021/acs.analchem.8b05592>
143. Yu Y, Mai Y, Zheng Y, Shi L. Assessing and mitigating batch effects in large-scale omics studies. *Genome Biol.* 2024;25:254. <https://doi.org/10.1186/s13059-024-03401-9>
144. Chang L, Ewald J, Hui F, Bayen S, Xia J. A data-centric perspective on exposomics data analysis. *Exposome.* 2024;4:osae005. <https://doi.org/10.1093/exposome/osae005>
145. Huang S-Y, Yang Y-X, Chen S-D, et al. Investigating causal relationships between exposome and human longevity: a Mendelian randomization analysis. *BMC Med.* 2021;19:150. <https://doi.org/10.1186/s12916-021-02030-4>
146. Domenighetti C, Sugier P-E, Ashok Kumar Sreelatha A, Comprehensive Unbiased Risk Factor Assessment for Genetics and Environment in Parkinson's Disease (Courage-PD) Consortium, et al. Dairy Intake and Parkinson's Disease: a Mendelian Randomization Study. *Mov Disord.* 2022;37:857-864. <https://doi.org/10.1002/mds.28902>
147. Li D, Zhou L, Cao Z, et al. Associations of environmental factors with neurodegeneration: an exposome-wide Mendelian randomization investigation. *Ageing Res Rev.* 2024;95:102254. <https://doi.org/10.1016/j.arr.2024.102254>
148. Zhao Y-L, Hao Y-N, Ge Y-J, et al. Variables associated with cognitive function: an exposome-wide and mendelian randomization analysis. *Alzheimers Res Ther.* 2025;17:13. <https://doi.org/10.1186/s13195-025-01670-5>
149. Maitre L, Guimbaud J-B, Warembourg C, Exposome Data Challenge Participant Consortium, et al. State-of-the-art methods for exposure-health studies: results from the exposome data challenge event. *Environ Int.* 2022;168:107422. <https://doi.org/10.1016/j.envint.2022.107422>
150. Barupal DK, Fiehn O. Chemical Similarity Enrichment Analysis (ChemRICH) as alternative to biochemical pathway mapping for metabolomic datasets. *Sci Rep.* 2017;7:14567. <https://doi.org/10.1038/s41598-017-15231-w>
151. Johnson CH, Ivanisevic J, Siuzdak G. Metabolomics: beyond biomarkers and towards mechanisms. *Nat Rev Mol Cell Biol.* 2016;17:451-459. <https://doi.org/10.1038/nrm.2016.25>
152. Xia J. Computational strategies for biological interpretation of metabolomics data. In: Sussulini A, ed. *Metabolomics: From Fundamentals to Clinical Applications.* Springer International Publishing; 2017:191-206.
153. Li S, Park Y, Duraisingham S, et al. Predicting network activity from high throughput metabolomics. *PLOS Comput Biol.* 2013;9:e1003123. <https://doi.org/10.1371/journal.pcbi.1003123>
154. Rohart F, Gautier B, Singh A, Cao K-AL. mixOmics: an R package for 'omics feature selection and multiple data integration. *PLOS Comput Biol.* 2017;13:e1005752. <https://doi.org/10.1371/journal.pcbi.1005752>
155. Zuffa S, Schmid R, Bauermeister A, et al. microbeMASST: a taxonomically informed mass spectrometry search tool for microbial metabolomics data. *Nat Microbiol.* 2024;9:336-345. <https://doi.org/10.1038/s41564-023-01575-9>
156. Fiehn O, Robertson D, Griffin J, et al. The metabolomics standards initiative (MSI). *Metabolomics.* 2007;3:175-178. <https://doi.org/10.1007/s11306-007-0070-6>
157. Salek RM, Neumann S, Schober D, et al. COordination of Standards in MetabOmicS (COSMOS): facilitating integrated metabolomics data access. *Metabolomics.* 2015;11:1587-1597. <https://doi.org/10.1007/s11306-015-0810-y>
158. Vineis P, Chadeau-Hyam M, Gmuender H, EXPOsOMICS Consortium, et al. The exposome in practice: design of the EXPOsOMICS project. *Int J Hyg Environ Health.* 2017;220:142-151. <https://doi.org/10.1016/j.ijheh.2016.08.001>
159. The European Human Exposome Network. Home—The European Human Exposome Network (EHEN). Accessed May 3, 2024. <https://www.humanexposome.eu/>
160. NEXUS. <https://www.nexus-exposomics.org/>
161. Niedzwiecki MM, Miller GW. HERCULES: an academic center to support exposome research. In: Dagnino S, Macherone A, eds. *Unraveling the Exposome: A Practical View.* Springer International Publishing, 2019:339-348.
162. Westat. Human Health Exposure Analysis Resource (HHEAR). Accessed July 15, 2025. <https://hhearprogram.org/>
163. International Human Exposome Network. IHEN—The International Human Exposome Network. In: IHEN. Accessed July 15, 2025. <https://humanexposome.net/>
164. Aharoni A, Goodacre R, Fernie AR. Plant and microbial sciences as key drivers in the development of metabolomics research. *Proc Natl Acad Sci U S A.* 2023;120:e2217383120. <https://doi.org/10.1073/pnas.2217383120>
165. Frigerio G, Moruzzi C, Mercadante R, Schymanski EL, Fustinoni S. Development and application of an LC-MS/MS untargeted exposomics method with a separated pooled quality control strategy. *Molecules.* 2022;27:2580. <https://doi.org/10.3390/molecules27082580>
166. Broeckling CD, Beger RD, Cheng LL, et al. Current practices in LC-MS untargeted metabolomics: a scoping review on the use of pooled quality control samples. *Anal Chem.* 2023;95:18645--18654. <https://doi.org/10.1021/acs.analchem.3c02924>
167. Caballero-Casero N, Belova L, Vervliet P, et al. Towards harmonised criteria in quality assurance and quality control of suspect and non-target LC-HRMS analytical workflows for screening of emerging contaminants in human biomonitoring. *TrAC Trends Anal Chem.* 2021;136:116201. <https://doi.org/10.1016/j.trac.2021.116201>
168. Viant MR, Ebbels TMD, Beger RD, et al. Use cases, best practice and reporting standards for metabolomics in regulatory toxicology. *Nat Commun.* 2019;10:3041. <https://doi.org/10.1038/s41467-019-10900-y>
169. Kirwan JA, Gika H, Beger RD, metabolomics Quality Assurance and Quality Control Consortium (mQACC), et al. Quality assurance and quality control reporting in untargeted metabolic phenotyping: mQACC recommendations for analytical quality management. *Metabolomics.* 2022;18:70. <https://doi.org/10.1007/s11306-022-01926-3>
170. Haug K, Cochrane K, Nainala VC, et al. MetaboLights: a resource evolving in response to the needs of its scientific community. *Nucleic Acids Res.* 2020;48:D440-D444. <https://doi.org/10.1093/nar/gkz1019>
171. Leao TF, Clark CM, Bauermeister A, et al. Quick-start infrastructure for untargeted metabolomics analysis in GNPS. *Nat Metab.* 2021;3:880-882. <https://doi.org/10.1038/s42255-021-00429-0>
172. Chetnik K, Petrick L, Pandey G. MetaClean: a machine learning-based classifier for reduced false positive peak detection in untargeted LC-MS metabolomics data. *Metabolomics.* 2020;16:117. <https://doi.org/10.1007/s11306-020-01738-3>

173. Mosley JD, Schock TB, Beecher CW, et al. Establishing a framework for best practices for quality assurance and quality control in untargeted metabolomics. *Metabolomics*. 2024;20:20. <https://doi.org/10.1007/s11306-023-02080-0>
174. Petrick LM, Shomron N. AI/ML-driven advances in untargeted metabolomics and exposomics for biomedical applications. *Cell Rep Phys Sci*. 2022;3:100978. <https://doi.org/10.1016/j.xcrp.2022.100978>
175. Wan M, Simonin EM, Johnson MM, et al. Exposomics: a review of methodologies, applications, and future directions in molecular medicine. *EMBO Mol Med*. 2025;17:599-608. <https://doi.org/10.1038/s44321-025-00191-w>
176. Rappaport SM. Genetic factors are not the major causes of chronic diseases. *PLoS One*. 2016;11:e0154387. <https://doi.org/10.1371/journal.pone.0154387>
177. Fuller R, Landrigan PJ, Balakrishnan K, et al. Pollution and health: a progress update. *Lancet Planet Health*. 2022;6:e535-e547. [https://doi.org/10.1016/S2542-5196\(22\)00090-0](https://doi.org/10.1016/S2542-5196(22)00090-0)
178. Tisler S, Savvidou P, Jørgensen MB, Castro M, Christensen JH. Supercritical fluid chromatography coupled to high-resolution mass spectrometry reveals persistent mobile organic compounds with unknown toxicity in wastewater effluents. *Environ Sci Technol*. 2023;57:9287-9297. <https://doi.org/10.1021/acs.est.3c00120>
179. Wild CP. The exposome at twenty: a personal account. *Exposome*. 2025;5:osaf003. <https://doi.org/10.1093/exposome/osaf003>
180. Maitre L, Bustamante M, Hernández-Ferrer C, et al. Multi-omics signatures of the human early life exposome. *Nat Commun*. 2022;13:7024. <https://doi.org/10.1038/s41467-022-34422-2>
181. Argentieri MA, Amin N, Nevado-Holgado AJ, et al. Integrating the environmental and genetic architectures of aging and mortality. *Nat Med*. 2025;31:1016-1025. <https://doi.org/10.1038/s41591-024-03483-9>