

This is the peer reviewed version of the following article: Vandemoortele B., Vanessa, V. **Molecular systems biology approaches to investigate mechanisms of gut–brain communication in neurological diseases.** *Eur J Neurol.* 2023; 30: 3622-3632., which has been published in final form at <https://doi.org/10.1111/ene.15819>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions. This article may not be enhanced, enriched or otherwise transformed into a derivative work, without express permission from Wiley or by statutory rights under applicable legislation. Copyright notices must not be removed, obscured or modified. The article must be linked to Wiley's version of record on Wiley Online Library and any embedding, framing or otherwise making available the article or pages thereof by third parties from platforms, services and websites other than Wiley Online Library must be prohibited.

Molecular systems biology approaches to investigate mechanisms of gut-brain communication in neurological diseases

Vandemoortele Boris^{1,2,3} & Vermeirssen Vanessa^{1,2,3,*}

¹ Lab for Computational Biology, Integromics and Gene Regulation (CBIGR), Cancer Research Institute Ghent (CRIG), Ghent, Belgium

² Department of Biomedical Molecular Biology, Ghent University, Ghent, Belgium

³ Department of Biomolecular Medicine, Ghent University, Ghent, Belgium

*Corresponding author: vanessa.vermeirssen@ugent.be

Total word count: 6916 (4998 main text)

Short running title: Systems biology for the gut-brain axis

1 **Abstract (250)**

2

3 While the incidence of neurological disease is increasing worldwide, treatment remains mostly
4 limited to symptom management. The gut-brain axis, which encompasses the communication
5 routes between microbiota, gut, and brain, has emerged as a crucial area of investigation for
6 identifying new preventive and therapeutic targets in neurological disease. Due to the inter-
7 organ, systemic nature of the gut-brain axis, together with the multitude of biomolecules and
8 microbial species involved, molecular systems biology approaches are required to accurately
9 investigate the mechanisms of gut-brain communication. High-throughput omics profiling,
10 together with computational methodologies such as dimensionality reduction or clustering,
11 machine learning, network inference and genome-scale metabolic models, allow to discover
12 novel biomarkers and elucidate mechanistic insights. In this review, we introduce the general
13 concepts of experimental and computational methodologies for gut-brain axis research and
14 discuss their applications, mainly in human cohorts. We further highlight important aspects
15 concerning rational study design, sampling procedures and data modalities relevant for gut-
16 brain communication, strengths and limitations of methodological approaches and some future
17 perspectives. In conclusion, we review how multi-omics analysis, together with advanced data
18 mining, are essential to functionally characterize the gut-brain axis in neurological disease and
19 finally put forward novel preventive or therapeutic strategies.

20

1 **Introduction**

2

3 Despite decades of research, treatment of neurological diseases is limited to symptom
4 management, and most drugs achieve only moderate efficacy. Hence, there is an urgent unmet
5 medical need for novel, cost-efficient disease-modifying treatments, of which patients are most
6 likely to benefit if administered early in disease progression. In this respect, the gut-brain axis
7 (GBA) is an exciting research area that opens up possibilities for preventive and therapeutic
8 strategies. The GBA refers to the communication routes and interrelationship between
9 microbiota and inflammation in the gut, and neuroinflammation and neurological disease in the
10 brain¹. Due to the inter-organ, systemic nature of the GBA, together with the diversity of
11 molecular features and microbial species involved, systems biology approaches are required
12 to capture underlying molecular mechanisms.

13

14 Technological advancements in high-throughput molecular profiling enable the cost-efficient,
15 high-throughput analysis of multiple biomolecules and molecular interactions in parallel, such
16 as genome, epigenome, transcriptome, proteome, metabolome and interactome. In addition,
17 latest developments in computational methodologies allow performing multi-omics data
18 integration to capture a multi-modal view and resolve the flow of signaling information across
19 multiple regulatory layers. In this way, information is obtained that exceeds that of the sum of
20 the individual omics, such that GBA molecular systems biology emerges as an exciting new
21 framework to study complex neurological diseases. However, data integration and
22 interpretation have emerged as new bottlenecks. Integrating multi-omics datasets is
23 challenging due to heterogeneity in terms of size, format, dimensionality, noisiness and
24 information content. Generally, we distinguish between multi-stage and multi-dimensional
25 integration, where the former merges different modalities in subsequent steps and the latter
26 combines them at once. The most ideal scenario is to preserve specific properties of each data
27 modality and to integrate different omics data and environmental context simultaneously to
28 identify coordinated behavior between the different levels. Broadly, multi-omics data
29 integration can be based on pairwise statistical association, clustering or dimensionality
30 reduction, network inference, machine learning and composite methodologies. Pairwise
31 statistical association focusses on the interaction between pairs of omics data; clustering or
32 dimensionality reduction transforms the data into a common space of lower dimensions to find
33 patterns in the data; network inference maps the data on to graphs representing interactions
34 between biomolecules; and machine learning predicts or classifies through a model that is
35 optimized iteratively using input data. In addition, prior information in the form of expert
36 biological knowledge can be taken into account, such as in genome-scale metabolic models
37 (GEMs), which are of particular interest in the context of GBA and use gene-protein reaction

1 rules to link genes encoding enzymes to curated metabolic pathways. These GEMs can be
2 constructed for both the host and the gut microbiome, enabling to model host-microbe and
3 microbe-microbe interactions².

4
5 Multi-omics integration serves several goals in GBA research: understanding the molecular
6 mechanisms at play, classifying disease versus normal, subtyping within a disease, predicting
7 biomarkers for diagnosis and prognosis, identifying causal effects and detecting molecular
8 drivers. Some methodologies are more suited for a given purpose than others: while clustering
9 or dimensionality reduction and machine learning lean more towards biomarker discovery,
10 patient classification and subtyping, statistical association, network inference and GEMs allow
11 to elucidate molecular mechanisms. Therefore, in the choice of methodology, the biological
12 question has to be kept in mind. In this work we review how omics and multi-omics analysis,
13 together with advanced bioinformatics and machine learning, are essential to functionally
14 characterize the GBA and put forward novel preventive or therapeutic strategies for
15 neurological disorders. First, we address specific experimental design challenges in the
16 context of GBA studies with a focus on human cohorts, and summarize the characteristics of
17 various data modalities. Next, we introduce statistical concepts behind computational
18 methodologies, and show their application in recently published studies focusing on
19 neurological disease. In addition, we discuss the strengths and weaknesses of the investigated
20 approaches, as well as potential challenges and future directions of molecular systems biology
21 in GBA research.

22 23 **Experimental study design along the gut-brain axis**

24 Different routes have been identified by which crosstalk between gut and brain, and between
25 the immune system and the nervous system, occurs. Gut microbiota and their metabolites can
26 directly activate the vagus nerve, which connects to the central nervous system. They can also
27 influence the generation, maturation and function of immune cells, which can migrate to the
28 brain¹. In addition, microbial products can induce the release of neurotransmitters and peptide
29 hormones from enteroendocrine cells. Moreover, several microbes and microbial metabolites
30 can pass through the intestinal barrier, enter systemic circulation, cross the brain barriers and
31 act as neuroimmunomodulatory signals in the brain^{3,4}. This especially occurs upon dysbiosis
32 and inflammation when permeability of gut and brain barriers is increased. Both cell-
33 autonomous and circulating metabolites serve as signals that transmit information about
34 environmental changes to individual cells to induce appropriate adjustments in gene
35 regulation. This implicates that constructing a complete and accurate GBA model will require
36 multi-omics molecular data sampled from multiple tissues and liquids such as blood, stool and
37 cerebrospinal fluid (CSF)^{4,5} (**Figure 1A**). Ideally, one has access to both gut and brain tissue

1 biopsies, a challenging objective in human cohorts, hence researchers often turn to animal
2 models. In addition, perturbation experiments, which are better suited to demonstrate
3 causality, are more feasible in animal models. The initial design and the types of data
4 generated in a given study thus determine which analyses can be performed, and which
5 biological questions can be answered. Next to the multi-omics data to profile, an additional
6 consideration regarding experimental design is where, in which individuals and when to sample
7 these data.

8

9 Extensive and heterogenous patient cohorts better represent the overall population, but will
10 result in a higher within-group variability. Moreover, multi-omics data, including microbial
11 diversity, are influenced by a myriad of confounding factors besides disease status, e.g. age,
12 sex, geographical location, diet, past and current drug treatments, and even physical exercise,
13 often leading to contradictory findings across studies^{6,7}. A rigorous analysis of and correction
14 for potential confounding factors is thus essential, and findings should not readily be
15 extrapolated across heterogenous populations. Next, also sampling location and time points
16 must be considered. Longitudinal sampling over a given time period is recommended and
17 allows to monitor the dynamics of disease progression, biomarker presence and the effect of
18 environmental factors. Molecular snapshots taken at a single time point lack these features,
19 but are much easier to obtain, especially if sampling requires invasive procedures or essential
20 tissues. Whole blood and serum, next to CSF, represent an established GBA communication
21 route, making it an interesting resource to identify potential messenger-molecules⁴. Stool
22 samples on the other hand can be considered a functional readout of the gut microbiota⁴.
23 Circulating immune cells can be isolated from blood and CSF, after which they can be
24 phenotypically profiled using multi-omics⁸. Finally, tissue biopsies can be analyzed by high-
25 throughput molecular profiling in different anatomical locations, i.e. gut and brain. Preferably,
26 all omics are measured on the same subjects, resulting in coupled data that is well suited for
27 integration⁵. All samples ideally originate from the same biological material i.e. from the same
28 tissue or liquid biopsy at the same time, in order to avoid batch effects between the different
29 omics data. In a so-called split sample study design, samples taken from the same biological
30 material are divided for different omics analyses. In a replicate-matched study design, samples
31 from different biological replicates within the same experiment are used for different omics
32 analyses e.g. mutually exclusive omics analyses in which sample preparation for one omics
33 impedes profiling of a second omics in the same sample. In the GBA context, often a source-
34 matched study is conducted, where different samples of the same individual are chosen for
35 different analyses e.g. transcriptomics on gut tissue and metabolomics on plasma⁵. Whether
36 these samples are best taken at a single time point depends on the biological question at hand,
37 as different omics modalities are subject to different time scales of change. For example,

1 changing metabolite levels might result from fast post-translational feedback mechanisms
2 within metabolic pathways, while transcriptomic changes resulting from these altered
3 metabolite levels are likely observable only after a given period of time.

5 **Different data modalities synergistically define the neurological disease phenotype**

6 Uncovering the unknown function of a single gene is a monumental task, but gives only limited
7 insight. Genes do not act in isolation but are embedded in highly complex biological systems.
8 Transcription is facilitated by regulatory factors such as transcription factors, chromatin-
9 modifying enzymes and nucleosome remodeling complexes. Furthermore, metabolites, as
10 nutrients, products of the host proteome or derived from the gut microbiota, have broader roles
11 in cellular signaling than simply being sources of fuel and building blocks. Metabolites are
12 known to influence gene regulation as ligands for signaling or regulatory factors, and as
13 substrates or cofactors of DNA or histone modifying enzymes and chromatin remodelers⁹
14 **(Figure 1B)**. Especially short-chain fatty acids (SCFAs) modulate brain homeostasis and
15 neuroinflammation by affecting microglia activation, astrocyte and oligodendrocyte function
16 and Treg expansion through chromatin remodeling, histone deacetylase inhibition and ligand
17 binding to G protein-coupled receptors¹⁰. As an example, butyrate acts as an endogenous
18 inhibitor of histone deacetylases¹¹; and butyrate or butyrate-producing bacteria have also been
19 reported to be differentially abundant in numerous neurological diseases such as multiple
20 sclerosis (MS)¹², major depressive disorder (MDD)^{6,7}, Alzheimer's disease (AD)¹³ and
21 Parkinson's disease (PD)¹⁴. Interestingly, these SCFAs are not produced by the host but by
22 specific microbial species, thus linking the gut microbiome to human gene regulation and
23 neurological disease.

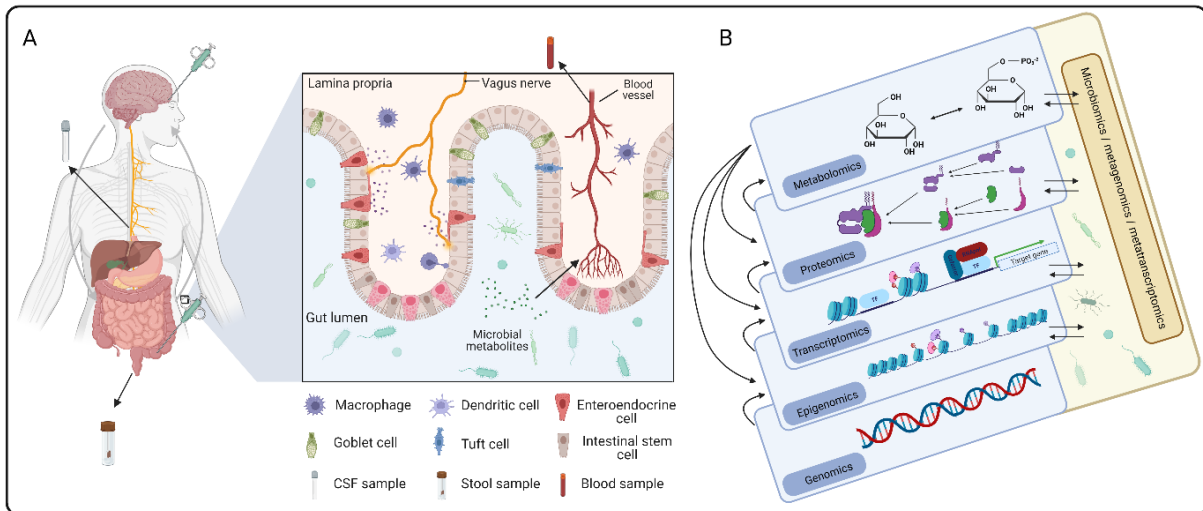
24
25 Specifically in the GBA context, **microbiomics/metagenomics and metatranscriptomics**
26 data are an essential resource. The analysis of 16S ribosomal RNA , referred to as
27 microbiomics, characterizes the presence of micro-organisms up to genus-level accuracy.
28 However, microbiomics provides insufficient resolution, as distinct species within the
29 *Bacteroides* genus differentially impact depression-like behavior¹⁵. Whole metagenome
30 shotgun sequencing, although more expensive, provides up to strain-level sensitivity as well
31 as insights into the functional potential of the identified micro-organisms¹⁶. Metatranscriptomics
32 in its turn reveals which microbial genes are actively being transcribed in a given context or
33 individual. **Metabolomics** data represent the phenotype of an entire biological regulatory
34 cascade due to its implicit integration of genomics, epigenomics, transcriptomics, proteomics
35 and even metagenomics data. Although often profiled in easily accessible bodily fluids such
36 as serum or stool samples, untargeted metabolomics data are extremely complex to interpret
37 and require highly-skilled scientists to annotate mass spectrometry peaks to biological

1 metabolites¹⁷. Targeted metabolomics on the other hand are more easily interpretable, but are
2 limited to known metabolites, and thus represent only a fraction of the true metabolic diversity.
3 In addition to metabolomics and microbiome data which can be easily obtained from liquid
4 biopsies, transcriptomes and epigenomes can be profiled from tissue biopsies and circulating
5 immune cells. These data can reveal potential gene regulatory mechanisms through
6 association of gene expression with specific transcription factor binding sites within regions of
7 open chromatin. If proteomics data are added, which are also a complex data type, the effect
8 of post-translational modifications on protein activity and half-life can be incorporated.

9

10 Molecular data resulting from transcriptomics, proteomics, metabolomics, epigenomics and
11 microbiomics/metagenomics assays are highly dimensional, implying the amount of molecular
12 features i.e. genes or metabolites measured greatly outnumbers the number of observations.
13 This is referred to as the curse of dimensionality¹⁸. Additionally, data encounter sparsity at
14 several levels. First, there is sparsity at the sample or patient level, as not all omics are profiled
15 across all samples. Second, there is sparsity within each omics dataset, as not all molecular
16 features are measured in all samples due to technical limitations. Furthermore, different omics
17 suffer from different degrees of sparsity, with mass spectrometry-based techniques often
18 resulting in more sparse data matrices. Data of different modalities result in discrete or
19 continuous, numerical or categorical variables with different ranges that cannot be directly
20 compared such that data first need to be transformed. In addition, more omics data are
21 collected at single cell level instead of at tissue level (i.e. bulk), allowing to dissect complex
22 tissues into specific cell states. However, these data present with high heterogeneity,
23 dimensionality, overdispersion and sparsity¹⁹. Hence, omics data preprocessing should be
24 performed with great care. This includes quality control by removal of batch effects resulting
25 from technical artefacts and removal of low signal features, normalization across individuals,
26 scaling of the different data modalities, and selection of highly variable features within each
27 dataset.

28



1
 2 **Figure 1: The gut-brain axis encompasses systems-level, inter-organ regulatory processes. A:** The gut-
 3 brain axis connects different organs and tissues through three main communication routes. First, the vagus nerve
 4 connects the brain to the enteric nervous system, which in its turn is connected to the central nervous system.
 5 Second, epithelial and immune cells are in direct contact with microbial cells and molecules, after which immune
 6 cells can travel to anatomically distant organs. Third, microbial metabolites can enter the bloodstream and reach
 7 the blood brain barrier through systemic circulation. Ideally, omics data are sampled from gut and brain tissue, stool,
 8 blood and cerebrospinal fluid. **B:** Biological systems are regulated through feedback mechanisms between multiple
 9 molecular levels, and need to be profiled by different omics techniques. Genetic information is encoded in the DNA,
 10 which is tightly packaged into chromatin. Epigenetic modification of histones such as acetylation and methylation
 11 can result in open chromatin, which can be actively transcribed by RNA polymerase in combination with
 12 transcriptional regulators. mRNA is then translated into protein, which often function in multi-protein complexes.
 13 These proteins facilitate cellular metabolism, and metabolites in turn regulate the activity of proteins such as
 14 transcription factors and histone modifiers. Created with BioRender.com

1 **Computational strategies for biomarker discovery and patient stratification**

2
3 Neuroinflammation is well known to present with gut inflammation and dysbiosis, and
4 increasing evidence indicates gut dysfunction may precede neurological symptoms even by
5 decades²⁰. This opens up new possibilities towards early screening . Screening methods find
6 their origin in the differential abundance of specific molecular markers in a patient compared
7 to control group or in healthy versus inflamed tissue. Based on the data modalities, the number
8 of individuals enrolled and the eventual goal of the study, three classes of methods can be
9 used to classify patients and identify descriptive molecular compounds (**Figure 2A, Table 1**).

10 The first class of methods, already widely applied to GBA, aims to identify **differential**
11 **abundances** of specific individual molecular features based on a single omics. These
12 statistical tests assess for each individual feature whether it is present in a higher or lower
13 abundancy in one group compared to the other, and whether the observed difference reflects
14 biological signal or random noise. The microbiota can be additionally characterized by two
15 distinct diversity measures: alpha- and beta-diversity. Alpha-diversity is a measure for the
16 richness of a microbial community, i.e. the number of different genera/species/strains, and is
17 represented by indices such as the Shannon or Fisher index⁷. Beta-diversity reflects
18 differences in microbial abundance of specific genera/species/strains between groups. In a
19 mouse model of autism spectrum disorder, differences in stool beta-diversity were reported
20 together with a >40-fold increase of the metabolite 4-ethylphenylsulfate (4EPS) in serum
21 metabolomics compared to a healthy control group, an effect that could be restored by treating
22 mice with *Bacteriodes Fragilis*²¹. A follow-up study reported that gut-derived 4EPS could enter
23 systemic circulation and subsequently the brain, where it altered oligodendrocyte maturation,
24 neuronal myelination and increased anxiety-like behavior³. Similar studies using
25 metabolomics, microbiomics/metagenomics and/or transcriptomics have been performed in
26 extensive patient cohorts for other neurological disorders, for example depletion of the
27 *Coprococcus* and *Dialister* taxa has been reported across multiple MDD patient cohorts⁶.
28 However, methods designed to identify differentially abundant molecular features suffer from
29 the dimensionality curse, i.e. the fact that each dataset contains many more features than
30 observations. Also, they fail to identify interactions between features and cannot integrate
31 multiple data modalities.

32 A second class of methods is **dimensionality reduction and clustering**. During Principle
33 Component Analysis (PCA), the most intuitive approach, each data point in a high-dimensional
34 space is mapped onto a novel set of axes, or principal components (PC), which are chosen
35 through rotation of the original axes to capture maximal variance in the original data. Since the
36 first PC captures the highest proportion of the total variance, a number of PCs can be chosen

1 to represent the data in lower dimensional space. This implicates that data points in the
2 reduced dimension matrix are actually linear combinations of the original data points, so that
3 they lose their one-on-one relationship with individual molecular features. Other methods such
4 as Principal Coordinate analysis (PCoA), UMAP or t-SNE have been specifically optimized for
5 high-dimensional and sparse biological data and can be used to create intuitive visualizations
6 of underlying data clusters, trajectories and patterns of interest¹⁸. PCoA is commonly used to
7 visualize microbiome beta-diversity due to its intrinsic visualization of dissimilarity between
8 samples and robust performance on highly sparse data. Using a human discovery cohort,
9 microbiomics and PCoA, Liang and colleagues found significant microbiome differences in
10 individuals characterized by mild or no cognitive impairment, which were mainly caused by
11 species belonging to the *Firmicutes*, *Bacteroidetes* and *Proteobacteria* phyla²². Dimensionality
12 reduction techniques can further be used to project distinct data modalities into a common low-
13 dimensional space, facilitating data integration across omics, not only capturing signals shared
14 by all omics data, but also those emerging from the complementarity of the various omics²³.
15 Multi-omics factor analysis (MOFA) infers a set of (hidden) factors that capture both technical
16 and biological sources of variability across multi-omics data, allowing the identification of
17 sample subgroups through clustering, and identification of highly informative features across
18 multiple omics at once²⁴. Clark et al. profiled metabolome, proteome, lipidome, one-carbon
19 metabolism and inflammatory markers in the CSF of a cohort comprising adults with normal
20 cognition, mild cognitive impairment and mild dementia. MOFA identified 5 hidden factors
21 within the multi-omics datasets, to which protein 14-3-3 zeta/delta, clusterin, interleukin-15 and
22 transgelin-2 contributed substantially. Addition of these four MOFA-selected features to a
23 reference classification model for AD pathology resulted in an increased sensitivity and
24 specificity²⁵.

25 **Machine learning** comprises methods such as logistic regression, support vector machines,
26 Bayesian models, random forests (RFs) and boosting, which identify the most informative
27 features in a given dataset by assigning them a higher weight in the final model. Samples in
28 the original dataset are typically divided into a well-balanced training and testing set, after
29 which a model can be iteratively trained until it achieves both good prediction specificity and
30 sensitivity. Logistic regression is designed for binary classification tasks, and model
31 parameters can be interpreted as feature importance. Levi et al. applied logistic regression to
32 serum metabolomics data and achieved near-perfect separation of MS patients and healthy
33 controls¹². RFs are ensembles of decision trees, which are flowchart-like decision structures
34 that aim to maximize the differences between groups at each additional split. A major
35 advantage of RFs is their high interpretability through feature prioritization methods, which
36 estimate the effect of removing or replacing a given molecular feature during the classification

1 procedure²⁶. Removing or replacing highly-informative features will results in a substantial drop
2 in classification performance, thus identifying the degree to which specific features are
3 characteristic to the groups under study. Finally, gradient boosting decision tree models are a
4 variation to RFs, in which new trees are added to the ensemble of classifiers additively instead
5 of randomly. These have classified dysbiotic and non-dysbiotic microbiomes using predicted
6 metabolite secretion fluxes in a gut inflammation cohort, after which chorismate, D-ribose, L-
7 lactate and phenol were identified as the most informative metabolites²⁷. Levi and colleagues
8 further used gradient boosting to predict metabolite levels using only microbiome data, and
9 found 26 metabolites to be associated with the microbiome. One of these, indolepropionate,
10 was present in lower concentrations in the serum of MS patients compared to controls,
11 although there was no significant difference in the abundance of indolepropionate-producing
12 microbiota. However, MS patients' microbiomes had a lower abundance of indolelactate-
13 producing species, an intermediate of tryptophan to indolepropionate catabolism¹².

14

15 **Computational strategies to elucidate molecular mechanisms of GBA communication** 16 **in neurological disease**

17

18 Differential molecular abundances, hidden data patterns or features informative for machine
19 learning in themselves lack biological interpretability and do not result in disease pathology
20 insights. In order to achieve a mechanistic understanding of regulatory processes and to
21 identify causal effects or molecular drivers, the biological relationships between different omics
22 data types needs to be considered (**Figure 2B, Table 1**).

23

24 **Pairwise integration** of biological data can identify interactions between different omics types
25 and thus across regulatory layers in the cell and the organism. Two broad categories can be
26 distinguished, namely genetics of intermediate trait analysis, and correlation analysis between
27 two modalities such as the microbiome and metabolome. The analysis of expression
28 quantitative trait loci (eQTL) and DNA variants links genetic variations to transcriptomic
29 alterations by assessing statistical associations between genomic polymorphisms and the
30 expression levels of, often nearby, genes²⁸. However, the analysis of metabolic and microbial
31 trait loci is likely more informative in the GBA context, although this has not yet been reported.
32 Correlation analyses on the other hand can reveal functional interactions between microbiota,
33 metabolites and genes. Within the IRONMET MDD patient cohort, lower dietary and circulating
34 proline levels were associated to lower Patient Health Questionnaire-9 (PHQ9) scores.
35 Circulating proline levels could further be positively associated to species from the
36 *Parabacteroides* and *Prevotella* genera, and negatively with *Actinobacteria* and SCFA-
37 producing species. Remarkably, patients with high dietary proline but low circulating proline

1 levels had a microbial signature associated with low PHQ9 scores, suggesting systems-level
2 interactions between the microbiome, metabolome and MDD symptoms⁷. Of particular interest
3 in the GBA context is the analysis of correlations between microbial abundances and
4 circulating metabolites. Such analyses revealed three ‘metabolite type-bacterial taxa’
5 correlated pairs in a model of MDD, next to associations between SCFAs and differentially
6 abundant microbial genera²⁹. Using longitudinal multi-omics data in an MS patient cohort,
7 Cantoni et al. found significant correlations between the gut microbiome and host blood
8 immune profiles in healthy controls but not in patients, suggesting disruptions of immune-
9 microbiome homeostatic interactions in MS³⁰. In an AD patient cohort, correlations between
10 differentially abundant microbial genera and CSF biomarkers have been reported, such that
11 easily accessible microbiome data might provide an alternative for the invasive lumbar
12 puncture in the future³¹. Finally, Liang et al. found significant correlations between serum
13 metabolites and metagenomic pathways enriched in individuals with impaired cognition²²,
14 which were identified using a LASSO logistic regression model²². However, although statistical
15 associations reveal insightful interactions between the microbiome, metabolome, gene
16 expression and phenotypic traits, they cannot identify causative features. Additional
17 (perturbation) experiments, longitudinal data or more advanced bioinformatics frameworks,
18 often including expert biological knowledge as prior information, are needed to identify
19 causality and true biological mechanisms of action.

20

21 **Integrated regulatory networks** provide an intuitive method to study molecular interactions
22 across omics types. These networks are constructed of nodes, which represent omics
23 features, and edges between the nodes reflecting correlations, regulatory or functional
24 interactions. Uncovering these edges in a given biological context comes down to unraveling
25 statistical dependencies between molecular features, using methods such as Bayesian
26 statistics, regression, mutual information and correlation on transcriptome data and requires a
27 large number of bulk observations. Single cell omics datasets on the other hand inherently
28 contain the required statistical variability between features such that patient-specific regulatory
29 networks can be constructed from a single sample, and regulatory programs can be inferred
30 in a cell-type specific manner. The inclusion of multi-omics in the network inference process,
31 such as regulator binding information or protein-protein interactions, results in more accurate
32 biological networks^{32,33}. Although the search for the best method is still subject of research, we
33 and others have shown that different methods add complementary information to the inference
34 of robust regulatory networks, which advocates for the construction of ensemble networks³⁴.
35 The popular methodology Weighted Gene Coexpression Network Analysis (WGCNA)
36 constructs coexpression modules from single omics, after which associations between
37 coexpression modules, other omics or phenotypic traits can be assessed³⁵. This has allowed

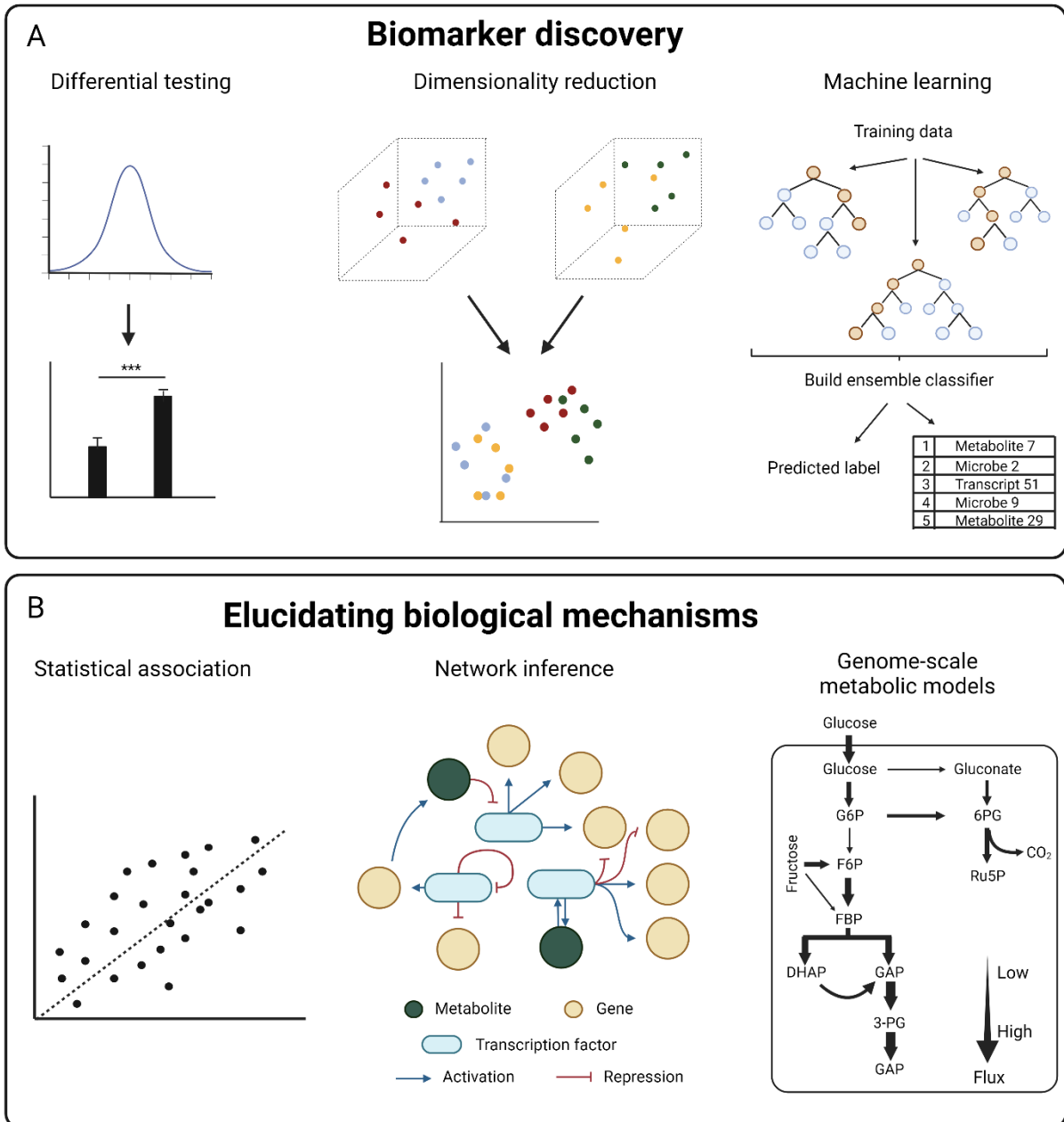
1 clustering of both serum and CSF metabolites into coexpression modules that could be
2 associated to MS severity⁸ or neurotoxicity³⁶. Similarly, coexpression modules have been
3 constructed from blood metabolomics in an AD patient cohort, revealing significant correlations
4 between amino acids, short-, medium- and long-chain acylcarnitines on the one hand, and AD
5 severity and cognitive traits on the other hand. Further integration with transcriptome data
6 highlighted *cpt1a*, which encodes a protein involved in the rate-limiting step of acylcarnitine
7 transport into mitochondria, and might thus account for the observed accumulation of
8 medium/long-chain acylcarnitines during AD progression³⁷. Badam et al. compared several
9 module detection tools such as WGCNA and clique-based methods to propose a framework
10 for multi-omics and genetic risk factor integration in MS³⁸. Zheng and colleagues made use of
11 non-human primates displaying depressive-like behavior and WGCNA to construct separate
12 metagenomics and metabolomics coexpression modules, which were correlated to
13 depressive-like behavior, revealing a functional interaction between the gut microbiome, lipid
14 metabolism and depressive-like behavior³⁹. Interestingly, WGCNA can also be applied on
15 different tissues and/or liquids, making it appealing in the GBA context. However, similar to
16 studies that assess correlations between features in omics datasets, WGCNA represents
17 merely associations. Directed regulatory networks, in which upstream regulators and
18 downstream targets are characterized, add a more causal interpretation of the data. The
19 integrative multi-omics module network inference algorithm Lemon-Tree builds coexpression
20 modules across samples using a model-based Gibbs sampler, after which potential regulators
21 are assigned through an ensemble of decision trees⁴⁰. Interestingly, expression data need not
22 come from a single tissue nor need potential regulators come from the same data modality, as
23 they can be any omics profiled on the same samples^{40,41}. This is particularly useful in the GBA
24 context, as it allows to model gene expression in both gut and brain, and assign microbial
25 signatures or circulating metabolites as potential regulators. Regulatory network inference on
26 neuroinflammation cohorts can thus theoretically be exploited to predict causal relationships
27 between the gut microbiota, circulating metabolites and gene expression patterns in the brain.
28 Overall, a critical aspect is to go beyond statistical associations and identify direct causal
29 relationships. Indeed, as insightful as associations between microbial abundances, metabolite
30 levels and transcriptional programs are, these fail to provide information on the directionality
31 of the interaction and cannot identify driving mechanisms. Transcription dynamics information
32 which is inherently present in single cell- or bulk time series transcriptome data, enables causal
33 network inference, as does the inclusion of other omics data, such as chromatin accessibility
34 data or regulator binding information.

35

36 Finally, expert-curated biological databases allow to exploit biological knowledge to achieve
37 context-specific causality. Several databases provide curated knowledge on transcriptional

1 regulatory interactions (DoRothEA, OmniPath), interactions between metabolites and
2 transcriptional regulators (STITCHdb, ASD), kinase-substrate interactions and metabolic
3 interactions (Recon3D), or interactions between genes and metabolic pathways (KEGG), and
4 these can be used as prior knowledge in integrated regulatory networks or pathway-based
5 models³². **GEMs**, such as the human Recon3D, use gene-protein-reaction associations based
6 on known genomic information and metabolic pathways to construct a stoichiometry matrix of
7 metabolites and reactions, creating a mathematical representation of a given metabolic
8 network. Through the addition of constraints such as nutrient input, upper- and lower-bound
9 reaction fluxes, and additional omics data, mainly transcriptomics, an optimization approach
10 can be applied to infer context-specific metabolic pathway activity⁴². Especially exciting in the
11 GBA context, GEMs can also be constructed for the microbiota, allowing to model microbe-
12 microbe and host-microbe metabolic interactions^{43,44}. The inclusion of personalized
13 microbiome, transcriptome and metabolome data allows contextualizing GEMs to patient-
14 specific models, in which the effects of dysbiosis on for example secreted and circulating
15 metabolites can be studied²⁷. These models can be further extended towards whole-body
16 metabolic reconstructions including tissue-specific GEMs and transport routes between
17 anatomically distant organ systems⁴². Using microbiome data and resulting personalized
18 microbial community models, Baldini and colleagues identified PD associated changes in the
19 predicted secretion potential of nine microbial metabolites including GABA, an effect mainly
20 explained by the *Akkermansia*, *Acidaminococcus* and *Roseburia* genera⁴⁵. Furthermore,
21 personalized microbial community GEMs have been used to describe increased circulating
22 homoserine levels and altered sulfur metabolism in PD patients. These changes could largely
23 be explained by differences in *Akkermansia muciniphila* abundance, as this species accounted
24 for over 50% of the variance in predicted secretion potential of methionine and hydrogen
25 sulfide¹⁴.

26



1
 2 **Figure 2: Overview of computational multi-omics analysis methods for biomarker discovery or the**
 3 **elucidation of biological mechanisms. A:** Differential testing identifies features that are differentially abundant in
 4 one group compared to another, and assesses whether this is the effect of biological signal or random noise.
 5 Dimensionality reduction methods allow to project different data modalities in a common latent space, and thus
 6 analyze them together. Machine learning allows to pinpoint biomarkers and classify observations into (sub)groups.
 7 A major advantage is the interpretability through feature prioritization methods. **B:** Pairwise statistical association
 8 models interactions between different omics but cannot identify causality. Network inference tools uncover
 9 regulatory interactions in a data-driven manner. Genome-scale metabolic models allow to study flux through
 10 metabolic networks and model host-microbe interactions. Created with BioRender.com

11

1 **Conclusion & future perspectives**

2
3 Due to technological advancements in high-throughput biology and computer science, GBA is
4 becoming an exciting field of study in neurological disease. Gut-brain communication involves
5 multiple organs, and matched sampling of the appropriate tissues and liquids at suitable time
6 points in extensive patient cohorts is recommended. Experimental design must be carefully
7 considered before and during each study to adequately deal with sample and omics
8 heterogeneity and confounding factors. Ideally, all omics are longitudinally profiled in all
9 individuals, and metadata must be collected and shared conform ethical standards. The study
10 design determines the downstream analysis and the conclusions that can ultimately be drawn.
11 Different data modalities complement one another in the description of the complete
12 neurological disease phenotype and allow for an enhanced biomarker discovery and the
13 elucidation of mechanistic insights. Since the generation of multi-omics data from patient
14 cohorts is a costly process often requiring valuable biological material, data should be shared
15 and made publicly available according to FAIR principles (Findable, Accessible, Interoperable
16 and Reusable)⁴⁶. This should allow for a more robust biomarker identification, as current
17 studies often suffer from limited sample sizes and incomplete control for confounding factors^{6,7}.
18 Advancements in state-of-the art omics profiling techniques at single molecule, single cell and
19 spatial level hold great promise, as these enable to study lowly-abundant molecular features,
20 rare cell types and spatial heterogeneity within tissues. Furthermore, improved untargeted
21 metabolomics analysis techniques will results in many more metabolites being profiled, as
22 today only a fraction of the total metabolomic diversity can accurately be identified. In addition,
23 stable isotope tracing in vivo with ¹³C-labeled nutrients or metabolites⁴⁷ is a powerful
24 methodology to demonstrate the causality and route of gut-brain communication in
25 neurological disease.

26
27 Novel joint dimensionality reduction and machine learning are increasingly being applied in
28 GBA to discover novel biomarkers. With the increase in multi-omics data for GBA, deep
29 learning approaches bear great potential to detect novel biomarkers. Integrated network
30 inference methods like WGCNA and Lemon-Tree are especially interesting in the GBA context,
31 since they enable multi-omics integration across tissues revealing mechanistic insights. Efforts
32 such as the Virtual Metabolic Human⁴⁹ database provide an invaluable resource for data
33 integration, but knowledge-based models by themselves are restricted to curated biological
34 information, and are therefore limited in uncovering novel biology. Ideally data integration is a
35 combination of supervised and unsupervised learning, such as the combination of
36 unsupervised network inference and supervised GEMs, thus contextualizing and extending
37 expert biological knowledge in a data-driven manner⁵⁰. This should result in more accurate

1 GBA models, allowing wet-lab researchers to prioritize hypotheses and most efficiently make
2 use of available resources.

3
4 Overall, integration of different data modalities, especially when profiled in distinct organ
5 systems, is still in its infancy and often lacks robustness. Currently, there is no optimal tool that
6 is broadly applicable to different types of research questions, and general guidance in the field
7 is lacking. Furthermore, there is an urgent need for the development of novel computational
8 approaches tailored towards the study design and data complexity inherent to GBA research.
9 Today, molecular systems biology and omics integration have already revealed significant
10 insights regarding gut-brain communication for numerous neurological diseases, as reviewed
11 here. The GBA can certainly be considered a basis for treating neurological diseases, and
12 likely holds great potential towards the development of disease-modifying therapeutics and
13 personalized medicine.

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

1 **Acknowledgements**

2 None.

3

4 **Conflict of interest**

5 Disclosure: None.

6

7 **Funding**

8 This work was supported by a grant from the Ghent University Special Research Fund
9 BOF/STA/201909/030 'Multi-omics data integration to elucidate the causes of complex
10 diseases'.

11 **Author contributions**

12 Author contributions are reported according to the CRediT taxonomy (BV = Boris
13 Vandemoortele, VV = Vanessa Vermeirssen). Conceptualization: BV, VV; Supervision: VV;
14 Visualization: BV; Writing - original draft: BV, VV; Writing – review and editing: BV, VV.

15

16 **Data availability statement**

17 Not applicable.

18

19 **ORCID**

20 Boris Vandemoortele <https://orcid.org/0000-0002-4352-0765>

21 Vanessa Vermeirssen <https://orcid.org/0000-0002-1975-0712>

22

23

24

25

26

27

28

29

30

31

1 **References**

- 2 1. Agirman, G., Yu, K. B. & Hsiao, E. Y. Signaling inflammation across the gut-brain axis.
3 *Science* **374**, 1087–1092 (2021).
- 4 2. Richelle, A. *et al.* Model-based assessment of mammalian cell metabolic functionalities
5 using omics data. *Cell Rep. Methods* **1**, 100040 (2021).
- 6 3. Needham, B. D. *et al.* A gut-derived metabolite alters brain activity and anxiety behaviour
7 in mice. *Nature* 1–7 (2022) doi:10.1038/s41586-022-04396-8.
- 8 4. Lai, Y. *et al.* High-coverage metabolomics uncovers microbiota-driven biochemical
9 landscape of interorgan transport and gut-brain communication in mice. *Nat. Commun.* **12**,
10 6000 (2021).
- 11 5. Cavill, R., Jennen, D., Kleinjans, J. & Briedé, J. J. Transcriptomic and metabolomic data
12 integration. *Brief. Bioinform.* **17**, 891–901 (2016).
- 13 6. Valles-Colomer, M. *et al.* The neuroactive potential of the human gut microbiota in quality
14 of life and depression. *Nat. Microbiol.* **4**, 623–632 (2019).
- 15 7. Mayneris-Perxachs, J. *et al.* Microbiota alterations in proline metabolism impact
16 depression. *Cell Metab.* **34**, 681-701.e10 (2022).
- 17 8. Fitzgerald, K. C. *et al.* Multi-omic evaluation of metabolic alterations in multiple sclerosis
18 identifies shifts in aromatic amino acid metabolism. *Cell Rep. Med.* **2**, 100424 (2021).
- 19 9. Schwartzman, J. M., Thompson, C. B. & Finley, L. W. S. Metabolic regulation of chromatin
20 modifications and gene expression. *J. Cell Biol.* **217**, 2247–2259 (2018).
- 21 10. Dalile, B., Van Oudenhove, L., Vervliet, B. & Verbeke, K. The role of short-chain fatty acids
22 in microbiota–gut–brain communication. *Nat. Rev. Gastroenterol. Hepatol.* **16**, 461–478
23 (2019).
- 24 11. Davie, J. R. Inhibition of histone deacetylase activity by butyrate. *J. Nutr.* **133**, 2485S-
25 2493S (2003).
- 26 12. Levi, I. *et al.* Potential role of indolelactate and butyrate in multiple sclerosis revealed by
27 integrated microbiome-metabolome analysis. *Cell Rep. Med.* **2**, 100246 (2021).

- 1 13. Liu, J. *et al.* Neuroprotective Effects of Clostridium butyricum against Vascular Dementia
2 in Mice via Metabolic Butyrate. *BioMed Res. Int.* **2015**, 412946 (2015).
- 3 14. Hertel, J. *et al.* Integrated Analyses of Microbiome and Longitudinal Metabolome Data
4 Reveal Microbial-Host Interactions on Sulfur Metabolism in Parkinson's Disease. *Cell Rep.*
5 **29**, 1767-1777.e8 (2019).
- 6 15. Zhang, Y. *et al.* Bacteroides species differentially modulate depression-like behavior via
7 gut-brain metabolic signaling. *Brain. Behav. Immun.* **102**, 11–22 (2022).
- 8 16. Pérez-Cobas, A. E., Gomez-Valero, L. & Buchrieser, C. Metagenomic approaches in
9 microbial ecology: an update on whole-genome and marker gene sequencing analyses.
10 *Microb. Genomics* **6**, mgen000409 (2020).
- 11 17. De Paepe, E. *et al.* A validated multi-matrix platform for metabolomic fingerprinting of
12 human urine, feces and plasma using ultra-high performance liquid-chromatography
13 coupled to hybrid orbitrap high-resolution mass spectrometry. *Anal. Chim. Acta* **1033**, 108–
14 118 (2018).
- 15 18. Malepathirana, T., Senanayake, D., Vidanaarachchi, R., Gautam, V. & Halgamuge, S.
16 Dimensionality reduction for visualizing high-dimensional biological data. *Biosystems* **220**,
17 104749 (2022).
- 18 19. Luecken, M. D. & Theis, F. J. Current best practices in single-cell RNA-seq analysis: a
19 tutorial. *Mol. Syst. Biol.* **15**, e8746 (2019).
- 20 20. Pellegrini, C. *et al.* Prodromal Intestinal Events in Alzheimer's Disease (AD): Colonic
21 Dysmotility and Inflammation Are Associated with Enteric AD-Related Protein Deposition.
22 *Int. J. Mol. Sci.* **21**, 3523 (2020).
- 23 21. Hsiao, E. Y. *et al.* The microbiota modulates gut physiology and behavioral abnormalities
24 associated with autism. *Cell* **155**, 1451–1463 (2013).
- 25 22. Liang, X. *et al.* Gut microbiome, cognitive function and brain structure: a multi-omics
26 integration analysis. *Transl. Neurodegener.* **11**, 49 (2022).
- 27 23. Cantini, L. *et al.* Benchmarking joint multi-omics dimensionality reduction approaches for
28 the study of cancer. *Nat. Commun.* **12**, 124 (2021).

- 1 24. Argelaguet, R. *et al.* Multi-Omics Factor Analysis—a framework for unsupervised
2 integration of multi-omics data sets. *Mol. Syst. Biol.* **14**, e8124 (2018).
- 3 25. Clark, C., Dayon, L., Masoodi, M., Bowman, G. L. & Popp, J. An integrative multi-omics
4 approach reveals new central nervous system pathway alterations in Alzheimer’s disease.
5 *Alzheimers Res. Ther.* **13**, 71 (2021).
- 6 26. Fabris, F., Doherty, A., Palmer, D., de Magalhães, J. P. & Freitas, A. A. A new approach
7 for interpreting Random Forest models and its application to the biology of ageing.
8 *Bioinformatics* **34**, 2449–2456 (2018).
- 9 27. Heinken, A., Hertel, J. & Thiele, I. Metabolic modelling reveals broad changes in gut
10 microbial metabolism in inflammatory bowel disease patients with dysbiosis. *Npj Syst. Biol.*
11 *Appl.* **7**, 1–11 (2021).
- 12 28. Patel, D. *et al.* Cell-type-specific expression quantitative trait loci associated with Alzheimer
13 disease in blood and brain tissue. *Transl. Psychiatry* **11**, 1–17 (2021).
- 14 29. Tian, T. *et al.* Multi-omics data reveals the disturbance of glycerophospholipid metabolism
15 caused by disordered gut microbiota in depressed mice. *J. Adv. Res.* **39**, 135–145 (2021).
- 16 30. Cantoni, C. *et al.* Alterations of host-gut microbiome interactions in multiple sclerosis.
17 *eBioMedicine* **76**, (2022).
- 18 31. Vogt, N. M. *et al.* Gut microbiome alterations in Alzheimer’s disease. *Sci. Rep.* **7**, 13537
19 (2017).
- 20 32. Dugourd. Causal integration of multi-omics data with prior knowledge to generate
21 mechanistic hypotheses. *Mol. Syst. Biol.* **17**, e9730 (2021).
- 22 33. Loers, J. U. & Vermeirssen, V. SUBATOMIC: a SUBgraph BAsed mulTi-OMIcs clustering
23 framework to analyze integrated multi-edge networks. *BMC Bioinformatics* **23**, 363 (2022).
- 24 34. Vermeirssen, V., De Clercq, I., Van Parys, T., Van Breusegem, F. & Van de Peer, Y.
25 Arabidopsis Ensemble Reverse-Engineered Gene Regulatory Network Discloses
26 Interconnected Transcription Factors in Oxidative Stress[W]. *Plant Cell* **26**, 4656–4679
27 (2014).

- 1 35. Langfelder, P. & Horvath, S. Eigengene networks for studying the relationships between
2 co-expression modules. *BMC Syst. Biol.* **1**, 54 (2007).
- 3 36. Ntranos, A. *et al.* Bacterial neurotoxic metabolites in multiple sclerosis cerebrospinal fluid
4 and plasma. *Brain* awab320 (2021) doi:10.1093/brain/awab320.
- 5 37. Horgusluoglu, E. *et al.* Integrative metabolomics-genomics approach reveals key
6 metabolic pathways and regulators of Alzheimer's disease. *Alzheimers Dement.* **n/a**,
7 (2021).
- 8 38. Badam, T. V. S. *et al.* A validated generally applicable approach using the systematic
9 assessment of disease modules by GWAS reveals a multi-omic module strongly
10 associated with risk factors in multiple sclerosis. *BMC Genomics* **22**, 631 (2021).
- 11 39. Zheng, P. *et al.* The gut microbiome modulates gut–brain axis glycerophospholipid
12 metabolism in a region-specific manner in a nonhuman primate model of depression. *Mol.*
13 *Psychiatry* **26**, 2380–2392 (2021).
- 14 40. Bonnet, E., Calzone, L. & Michoel, T. Integrative Multi-omics Module Network Inference
15 with Lemon-Tree. *PLOS Comput. Biol.* **11**, e1003983 (2015).
- 16 41. Erola, P., Björkegren, J. L. M. & Michoel, T. Model-based clustering of multi-tissue gene
17 expression data. *Bioinformatics* **36**, 1807–1813 (2020).
- 18 42. Thiele, I. Personalized whole-body models integrate metabolism, physiology, and the gut
19 microbiome. *Mol. Syst. Biol.* **16**, e8982 (2020).
- 20 43. Magnúsdóttir, S. & Thiele, I. Modeling metabolism of the human gut microbiome. *Curr.*
21 *Opin. Biotechnol.* **51**, 90–96 (2018).
- 22 44. Baldini, F. *et al.* The Microbiome Modeling Toolbox: from microbial interactions to
23 personalized microbial communities. *Bioinformatics* **35**, 2332–2334 (2019).
- 24 45. Baldini, F. *et al.* Parkinson's disease-associated alterations of the gut microbiome predict
25 disease-relevant changes in metabolic functions. *BMC Biol.* **18**, 62 (2020).
- 26 46. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and
27 stewardship. *Sci. Data* **3**, 160018 (2016).

- 1 47. Lund, P. J. *et al.* Stable isotope tracing in vivo reveals a metabolic bridge linking the
2 microbiota to host histone acetylation. *Cell Rep.* **41**, (2022).
- 3 48. Li, Y., Ge, X., Peng, F., Li, W. & Li, J. J. Exaggerated false positives by popular differential
4 expression methods when analyzing human population samples. *Genome Biol.* **23**, 79
5 (2022).
- 6 49. Noronha, A. *et al.* The Virtual Metabolic Human database: integrating human and gut
7 microbiome metabolism with nutrition and disease. *Nucleic Acids Res.* **47**, D614–D624
8 (2019).
- 9 50. Chung, C. H., Lin, D.-W., Eames, A. & Chandrasekaran, S. Next-Generation Genome-
10 Scale Metabolic Modeling through Integration of Regulatory Mechanisms. *Metabolites* **11**,
11 606 (2021).
- 12