1

Empirically establishing drug exposure records directly from untargeted metabolomics data

2 3

3 Haoqi Nina Zhao^{1,2,†}, Kine Eide Kvitne^{2,3,†}, Corinna Brungs^{4†}, Siddharth Mohan², Vincent Charron-4 Lamoureux^{1,2}, Wout Bittremieux^{1,2,5}, Runbang Tang², Robin Schmid^{1,2,4}, Santosh Lamichhane^{2,6}, 5 Yasin El Abiead^{1,2}, Mohammadsobhan S. Andalibi^{7,8,9}, Helena Mannochio-Russo^{1,2}, Madison 6 7 Ambre¹⁰, Nicole E. Avalon¹¹, MacKenzie Bryant¹⁰, Andrés Mauricio Caraballo-Rodríguez^{1,2}, Martin 8 Casas Maya¹⁰, Loryn Chin¹², Ronald J. Ellis^{7,8}, Donald Franklin⁸, Sagan Girod¹³, Paulo Wender P 9 Gomes^{1,2,14}, Lauren Hansen¹⁰, Robert Heaton⁸, Jennifer E. ludicello⁸, Alan K. Jarmusch^{1,2,15}, Lora Khatib⁷, Scott Letendre^{9,16}, Sarolt Magyari^{2,17}, Daniel McDonald¹⁰, Ipsita Mohanty^{1,2}, Andrés 10 Cumsille^{2,18}, David J. Moore^{8,9}, Prajit Rajkumar², Dylan H. Ross^{19,20}, Harshada Sapre², 11 Mohammad Reza Zare Shahneh²¹, Sydney P. Thomas^{1,2}, Caitlin Tribelhorn¹⁰, Helena M. Tubb¹⁰, 12 Corinn Walker¹⁰, Crystal X. Wang^{8,9}, Shipei Xing^{1,2}, Jasmine Zemlin,^{1,2,22} Simone Zuffa^{1,2}, David 13 S. Wishart^{12,23}, Rima Kaddurah-Daouk^{24,25,26}, Mingxun Wang²¹, Manuela Raffatellu^{10,22,27}, Karsten 14 Zengler^{11,10,22,28}, Tomáš Pluskal⁴, Libin Xu¹⁹, Rob Knight^{10,22,29,30,31}, Shirley M. Tsunoda², Pieter C. 15 16 Dorrestein^{1,2,22*} 17 ¹ Collaborative Mass Spectrometry Innovation Center, Skaggs School of Pharmacy and 18 19 Pharmaceutical Sciences, University of California San Diego, La Jolla, CA, USA 20 ² Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, 21 La Jolla, CA, USA ³ Department of Pharmacy, University of Oslo, Oslo, Norway 22 23 ⁴ Institute of Organic Chemistry and Biochemistry of the Czech Academy of Sciences, Prague, 24 **Czech Republic** 25 ⁵ Department of Computer Science, University of Antwerp, Antwerp, Belgium 26 ⁶ Turku Bioscience Centre, University of Turku and Åbo Akademi University, Tykistönkatu 6A, 27 20520 Turku, Finland ⁷ Department of Neurosciences, University of California San Diego, La Jolla, CA, USA 28 ⁸ Department of Psychiatry, University of California San Diego, La Jolla, CA, USA 29 ⁹ HIV Neurobehavioral Research Program, University of California San Diego, La Jolla, CA, USA 30 ¹⁰ Department of Pediatrics, University of California San Diego, La Jolla, CA, USA 31 32 ¹¹ Scripps Institution of Oceanography, University of California San Diego, La Jolla, CA, USA 33 ¹² Department of Bioengineering, University of California San Diego, La Jolla, California, USA. 34 ¹³ Department of Biological Sciences, University of Alberta, Edmonton, AB T6G 2E9, Canada ¹⁴ Faculty of Chemistry, Federal University of Pará, Belém, PA, Brazil 35 ¹⁵ Immunity, Inflammation, and Disease Laboratory, Division of Intramural Research, National 36 37 Institute of Environmental Health Sciences, National Institutes of Health, Research Triangle Park,

- 38 NC, USA
- ¹⁶ Department of Medicine, University of California San Diego, La Jolla, CA, USA.
- 40 ¹⁷ Institute of Microbiology, Eidgenössische Technische Hochschule (ETH) Zürich, Vladimir-
- 41 Prelog-Weg 4, 8093 Zürich, Switzerland
- ¹⁸ Department of Microbiology and Cell Sciences, University of Florida, Museum Drive,
 Gainesville, FL, USA
- 44 ¹⁹ Department of Medicinal Chemistry, University of Washington, Seattle, WA, USA

- 45 ²⁰Current address: Biological Sciences Division, Pacific Northwest National Laboratory, Richland,
- 46 WA, USA
- ²¹ Department of Computer Science and Engineering, University of California Riverside,
 Riverside, CA, USA
- 49 ²² Center for Microbiome Innovation, University of California San Diego, La Jolla, CA, USA.
- 50 ²³ Department of Computing Science, University of Alberta, Edmonton, AB T6G 2E8, Canada
- 51 ²⁴ Department of Psychiatry and Behavioral Sciences, Duke University, Durham, NC, 27708, USA
- 52 ²⁵ Duke Institute of Brain Sciences, Duke University, Durham, NC, USA
- 53 ²⁶ Department of Medicine, Duke University, Durham, NC, USA
- ²⁷ Chiba University, UC San Diego Center for Mucosal Immunology, Allergy, and Vaccines (CU UCSD cMAV), La Jolla, CA, USA
- Program in Materials Science and Engineering, University of California, San Diego, 9500
 Gilman Drive, La Jolla, CA 92093-0418, USA
- ²⁹ Department of Computer Science and Engineering, University of California San Diego, La Jolla,
 CA, USA
- ³⁰ Shu Chien-Gene Lay Department of Bioengineering, University of California San Diego, La
 Jolla, CA, USA
- 62 ³¹ Halıcıoğlu Data Science Institute, University of California San Diego, La Jolla, CA, USA
- [†]Haoqi Nina Zhao, Kine Eide Kvitne, and Corinna Brungs contributed equally to this work.
- 64 * Author to whom correspondence should be addressed.

65 Abstract

66 Despite extensive efforts, extracting information on medication exposure from clinical records 67 remains challenging. To complement this approach, we developed the tandem mass spectrometry (MS/MS) based GNPS Drug Library. This resource integrates MS/MS data for 68 drugs and their metabolites/analogs with controlled vocabularies on exposure sources, 69 70 pharmacologic classes, therapeutic indications, and mechanisms of action. It enables direct 71 analysis of drug exposure and metabolism from untargeted metabolomics data independent of clinical records. Our library facilitates stratification of individuals in clinical studies based 72 on the empirically detected medications, exemplified by drug-dependent microbiota-derived 73 74 *N*-acyl lipid changes in a cohort with human immunodeficiency virus. The GNPS Drug Library holds potential for broader applications in drug discovery and precision medicine. 75

76 Main Text

77 Growing evidence suggests that the chemical exposome plays a critical role in shaping 78 human health, with drugs being a significant source of chemical exposure that carries 79 profound health implications.¹ According to a recent survey by the Center for Disease Control and Prevention, nearly half (45.7%) of the U.S. population used at least one prescription drug 80 81 in the past 30 days.² Drug concentrations in human blood are on par with those of endogenous and dietary molecules,³ and have important impacts on the metabolic states 82 and microbiome composition.⁴⁻⁷ Clinical research typically relies on medical records or self-83 reporting surveys to assess drug exposure,⁸ but these methods are costly and often 84 incomplete.^{8–10} They often overlook over-the-counter medications and supplements, and fail 85 to account for patient adherence. Additionally, they miss drug usage not documented in 86 medical records, such as those purchased online,^{10,11} acquired across borders,^{12,13} or 87 consumed through secondary use of leftover drugs. Medical records are also incapable of 88 89 documenting drugs introduced into the food supply that are unknowingly consumed, such as 90 the antifungal natamycin used both to treat fungal eye infections and as a preservative for 91 dairy products. Additionally, the varying half-lives of drugs and their metabolites further 92 complicate exposure assessment, as some drugs are rapidly eliminated from the body while others can persist for months.14,15 93

94 Untargeted metabolomics offers the opportunity to complement clinical records by 95 empirically establishing the presence of drugs and their metabolites directly from biological samples. However, liquid chromatography-tandem mass spectrometry (LC-MS/MS) based 96 97 annotations, which rely on reference MS/MS library matches, are difficult to interpret. For 98 example, annotation may return complex IUPAC chemical name like а 99 "(2R,3S,4R,5R,8R,10R,11R,12S,13S,14R)-11-[(2S,3R,4S,6R)-4-(dimethylamino)-3-

100 hydroxy-6-methyloxan-2-yl]oxy-2-ethyl-3,4,10-trihydroxy-13-[(2R,4R,5S,6S)-5-hydroxy-4-

101 methoxy-4,6-dimethyloxan-2-yl]oxy-3,5,6,8,10,12,14-heptamethyl-1-oxa-6-

102 azacvclopentadecan-15-one". A text search in the right reference resource or an open web 103 search can hopefully link this IUPAC name to "azithromycin", the drug name used in clinical 104 settings. A second search of the term "azithromycin" is then required to connect the name to 105 its therapeutic role, in this case an antibiotic originally isolated from a bacterium. While this 106 example involves a simple name and a limited number of identifiers for azithromycin, other 107 compounds, like penicillin G or aspirin, have hundreds of synonyms and identifiers in 108 chemistry databases such as PubChem, making the identification process more challenging. 109 This task must be repeated for every obtained annotation, which can range from hundreds to 110 thousands in a given metabolomics experiment, to find all detected drugs in a dataset.

Even when the drugs are identified, the interpretation of their biological impacts requires extensive literature and web searches to understand the therapeutic roles of the drugs and their mechanisms of action. Public databases, such as DrugBank,^{16,17} DrugCentral,¹⁸ DailyMed,¹⁹ and KEGG DRUG,²⁰ can assist in interpretation, but the pharmacologic information is often provided as plain text or combinatorial classifications that require manual organization before downstream analysis. Although, in principle, large language models or similar text mining strategies can assist in this, the results of such models still need manual verification to confirm accuracy.^{21–23} In addition, it is not uncommon that only metabolized versions of a drug is present in the sample, leading to missed drug exposure if only the parent drug is considered.²⁴ Unfortunately, with very few exceptions, reference MS/MS libraries include only the parent drug but not the drug metabolites due to challenges in obtaining reference standards for these metabolites.

The absence of MS/MS spectra for many drug metabolites, along with other challenges mentioned above, makes it very difficult to accurately annotate all drug exposures in biological specimens. For instance, stratifying a cohort based on antibiotic exposure perhaps to better understand microbiome changes or as an exclusion criterion for clinical studies - requires identifying all antibiotics and their metabolites present in the samples. This is currently challenging due to the lack of resources that provide objective, systematic, and efficient readouts of drugs in untargeted metabolomics experiments.

130 To address this gap and to enable data science strategies on drug readouts, we 131 curated the Global Natural Product Social Molecular Networking (GNPS) Drug Library, a 132 collection of reference spectra for drugs and their metabolites/analogs (including parent ion 133 masses and MS/MS spectra) along with structured pharmacologic metadata including 134 exposure source, pharmacologic class, therapeutic indication, and mechanism of action. This 135 comprehensive resource will enable further data science analysis to empirically - and 136 retroactively - determine drug exposure using untargeted metabolomics data, complementing 137 the information available in clinical records.

The creation of this library involved three key steps: 1) collecting MS/MS spectra of drugs and drug metabolites from publicly available MS/MS reference libraries; 2) finding MS/MS spectra analogs of those drugs in publicly accessible untargeted metabolomics data to enhance coverage of the metabolized versions of drugs; and 3) linking each MS/MS spectrum of a drug to controlled-vocabulary metadata - the key component of this resource that facilitates efficient data interpretation (**Figure 1a**).

144 The reference MS/MS spectra of drugs and their known metabolites were collected 145 from two of the largest open-access mass spectral libraries, namely the GNPS Spectral 146 Library²⁵ and MSⁿLib²⁶. For all the MS/MS spectra in the GNPS and MSⁿLib, metadata enrichment was first performed against PubChem (for synonyms),²⁷ DrugCentral,¹⁸ the Broad 147 Institute Drug Repurposing Hub databases.²⁸ ChEMBL (for pharmacologic information).²⁹ and 148 149 DrugBank (for pharmacologic information and the Anatomical Therapeutic Chemical Classification code).^{16,17,30} This process utilized the available metadata in the GNPS Spectral 150 151 Library and MSⁿLib, including the chemical structures (e.g., SMILES or InChI), database identifiers (e.g., DrugBank ID or ChEMBL ID), and compound names. Based on the enriched 152 153 metadata regarding clinical phases, all MS/MS spectra of drugs and compounds in clinical 154 trials were compiled into the centralized GNPS Drug Library (see method details in 155 Supplementary Text 1), resulting in 99,122 MS/MS reference spectra for 4,723 unique 156 compounds. The compound names in the GNPS Drug Library were automatically curated 157 and set to the first synonym in PubChem. We note that the term "drug" is used here in a broad 158 sense, as the GNPS Drug Library includes not only prescribed and over-the-counter medications but also compounds currently in clinical trials, drugs that have been withdrawn,as well as substances with potential for abuse (e.g., cocaine, fentanyl).

161 Given that drug metabolites are largely overlooked in the initial search, we performed 162 a second "partial name match" to include metabolites that retain the full drug names. For 163 example, by searching for the name "venlafaxine" (a serotonin and norepinephrine reuptake 164 inhibitor used to treat depression and various anxiety disorders), we obtained reference 165 spectra for five of its metabolites, including "N-desmethylvenlafaxine", "Оdesmethylvenlafaxine", "N,O-didesmethylvenlafaxine", "N,N-didesmethylvenlafaxine", and 166 "venlafaxine N-oxide". Using this strategy combined with manual inspection of the results, we 167 captured 2,080 reference spectra for the metabolites of 110 drugs. Lastly, we added the 168 169 MS/MS spectra collected in the development of dmCCS,³¹ a collision cross section database for drugs and their metabolites where human liver microsomes and S9 fraction were used for 170 171 in vitro generation of drug metabolites. In total, 4,087 spectra for the metabolites of 470 drugs were included in the GNPS Drug Library (Figure 1a). 172





Figure 1. The GNPS Drug Library and connected pharmacologic metadata. a, The GNPS Drug Library comprises four key resources: Drug MS/MS reference spectra, drug metabolite MS/MS reference spectra, propagated drug analogs derived from public metabolomics datasets, and pharmacologic metadata connected to each reference spectrum. b, FastMASST analog search of drug spectra against public metabolomics studies yielded propagated drug-analogous MS/MS spectra, which were filtered by removing analogs for drugs with endogenous and food sources (source filter), removing mass offsets unexplained by common metabolic pathways (mass offset filter), removing 181 analogs with GNPS library matches (library match filter), removing analogs connected to multiple 182 drugs with dissimilar structures after spectra clustering (drug similarity filter), and removing analogs 183 with unrealistic drug exposure indications (dataset testing). c, Illustration of each filter employed in 184 curating FastMASST analog match results. d, Frequency of mass offsets in the propagated drug 185 analog library. The mass offsets were grouped by unit mass and stacked based on the number of 186 analog spectra. The most frequently observed mass offsets are colored while the rests are black. e, An example of structural modification sites predicted by ModiFinder.³² Purple color highlights modified 187 188 spectra and substructures, while the green color highlights unmodified ones. f, Overview of the 189 ontology-based drug metadata, highlighting common pharmaceutical classes and specific drugs in the 190 neurology/psychiatry category. Width of the bars and lines reflects the number of unique drug 191 structures in each class. g, The top 20 most detected pharmacologic classes in fecal samples from 192 the American Gut Project,³³ a cohort of the general population from the United States (US), Europe, 193 and Australia (1,993 individuals). h. Detected therapeutic drug class patterns by age and sex (1,845 194 individuals with age and sex information; age 46 ± 18 years [range 3-93], with 53% being female). 195 Detection of cardiovascular drugs increased with age, while analgesics, antihistamines, and antibiotics were detected across all ages.^{34,35} Analgesics were more frequently detected in females, consistent 196 197 with the literature.^{36,37} and drugs for erectile dysfunction were detected only in males. NSAID, non-198 steroidal anti-inflammatory drugs; ACE, angiotensin converting enzyme; SSRI, selective serotonin 199 reuptake inhibitor; PPI, proton pump inhibitor; DHFR, dihydrofolate reductase; HSV, herpes simplex 200 virus; SNRI, serotonin and norepinephrine reuptake inhibitor.

201

Despite the extensive collection effort, metabolite reference spectra were only available for 10% of the drugs included in the GNPS Drug Library. We hypothesized that unannotated drug metabolites are present in public untargeted metabolomics data. We further hypothesized that spectral alignment strategies can be used to find the modified versions of the drugs.^{38–40} In other words, public untargeted metabolomics data could be used to create a reference library of candidate drug metabolites that will facilitate the drug exposure readout in future datasets.

209 Based on MS/MS spectral alignment using two computational methods: repositoryscale molecular networking⁴¹ and fast Mass Spectrometry Search Tool (fastMASST) with 210 211 analog search,^{42,43} we retrieved all possible MS/MS spectra analogous to drugs from the 212 GNPS/MassIVE public repository (~2,700 LC-MS/MS datasets).²⁵ These spectra represent 213 drug-related molecules potentially derived from metabolism (host or microbiome), abiotic 214 processes, and adducts of drugs from MS measurements. We obtained analogous MS/MS spectra for 14.6% of the 103,209 reference spectra for drug and drug metabolites (>5.5 million 215 216 drug-analog spectral pairs).

217 In testing of the propagated analog library, we identified the need for additional filters 218 to enhance its relevance to drug exposure (Figure 1b-c). First, it is not possible to determine 219 the sources of exogenously supplied chemicals that are also produced endogenously or 220 derived from the diet. Consequently, structural analogs of drugs with endogenous or dietary 221 sources were excluded from the propagated drug analog library (e.g., analogous MS/MS 222 spectra of testosterone used to treat hypogonadism, or caffeine used as a stimulant drug, 223 were excluded). Second, propagated analogs with uncommon or unexplained mass offsets 224 (precursor mass difference between the propagated analog and the connected drug) were

excluded. Mass offsets were obtained from UNIMOD,⁴⁴ from a community-curated list of 225 explainable delta masses (Table S1), and from the Host Gut Microbiota Metabolism 226 227 Xenobiotics Database,⁴⁵ and were manually curated for those relevant to drug metabolism 228 (e.g., 14.02 Da, methylation; 176.03 Da, glucuronidation) or mass spectrometry adducts (e.g., 229 17.03 Da, ammonium adduct: see **Table S2** for the 156 mass offsets that were included). 230 MS/MS spectra of propagated analogs were excluded if the mass offsets were not in the 231 customized list, or when the mass offsets occurred fewer than ten times. Third, since drugs 232 within the same pharmacologic family often have similar structures, they can be identified as 233 analogs of each other through spectral alignments. Therefore, we excluded MS/MS spectra 234 with matches to the GNPS library from the propagated analog annotations. For example, a propagated analog of quinapril, an angiotensin converting enzyme (ACE) inhibitor, had a 235 236 spectral match to ramipril, another ACE inhibitor (Figure 1c). Excluding these analog 237 annotations ensures that they do not overwrite library matches of known drugs and 238 metabolites. Fourth, if one propagated analog spectrum is connected to multiple drugs after 239 spectral clustering, the drugs need to be structurally similar to accept the shared analog. We 240 illustrate this with the propagated analog (m/z 287.133, with a formula C₁₇H₁₉ClN₂) that is 241 connected to both hydroxyzine $[C_{21}H_{27}CIN_2O_2]$, mass offset 88.05 Da $(C_4H_8O_2)$] and chlorcyclizine [C₁₈H₂₁ClN₂, mass offset 14.02 Da (CH₂)], which share the core structure 242 243 (Figure 1c). Finally, we tested the propagated drug analog library against 12 public LC-244 MS/MS datasets to filter out analogs that have unrealistic drug exposure indications. The 245 selected datasets represent a broad range of human tissue types and biofluids, including 246 fecal (n=5), breast milk (n=2), plasma (n=3), skin (n = 1), and brain (n=1), as well as multiple 247 mouse tissues (n = 2; with metadata confirming no drugs were used). We observed analogs 248 of tocofersolan (a synthetic vitamin E derivative), iloprost (a synthetic prostacyclin mimetic), 249 desonide (a synthetic topical corticosteroid), medroxyprogesterone (a synthetic progestin), 250 and vidarabine (an adenosine analog used as an antiviral) in >50% of the human fecal samples from the American Gut Project (n = 1,993 individuals), a cohort of the general 251 population. The connected drugs for these analogs are derivatives of endogenous or food 252 derived molecules and are unlikely to be used by more than half of the population. Therefore, 253 254 these analogs cannot be confidently linked to drug exposures and were excluded. Analogs of 255 polidocanol (a synthetic long-chain fatty alcohol used as anesthetics) were observed in >70% of 2,463 human milk samples. They are likely surfactants/contaminants with the polyethylene 256 glycol structural units⁴⁶ and thus were excluded from the propagated drug analog library 257 (Figure 1c). 258

259 After all filtering steps, 3,234 clustered MS/MS spectra representing propagated 260 analogs of 577 drugs were retained in the final drug analog library. We observed that 75% of 261 the propagated analogs occurred at least once in the same data file with the corresponding 262 parent drugs (Table S3). The most common mass offsets in the drug analog library 263 correspond to a gain of 15.99 Da, which can be interpreted as the gain of an oxygen (e.g., 264 oxidative metabolism of the drug), followed by a gain or loss of 14.02 Da (CH₂, 265 (de)methylation), a gain of 1.00 Da (13C isotope), a gain of 2.00 Da (O/S/Cl/Br/2C isotopes), 266 and a gain or loss of 28.03 Da (C₂H₄, (de)ethylation; Figure 1d). Notably, it is possible that 267 such analog spectra are MS/MS of other ion forms of the parent drug, such as isotopes, 268 different adducts, in/post-source fragments, or multimers, rather than drug metabolites or 269 structural analogs. However, their indications in drug exposure remain the same and thus we 270 did not separate drug metabolites and instrument adducts in drug exposure stratification. To 271 extend structural hypotheses for the drug analogs, we employed the newly developed 272 ModiFinder,³² which leverages the shifted MS/MS fragment peaks in the MS/MS alignment to 273 predict the most likely location of the structural modifications. We were able to predict the 274 partial location of the modification for 61.5% of the analog spectra. We demonstrate examples 275 where ModiFinder predictions agree with expert manual interpretation of the MS/MS spectra 276 (Figure 1e, S1).

Connecting drug detections to their therapeutic indications typically requires expert 277 278 knowledge and/or extensive literature searches. The GNPS Drug Library addresses this 279 challenge by providing controlled-vocabulary metadata together with the specific drug 280 annotations. This allows users to annotate all drugs in an untargeted metabolomics dataset 281 and directly obtain a table with exposure sources, pharmacologic classes, therapeutic 282 indications, and mechanisms of action of the drugs, with their structures and names in a data 283 science ready format (Figure 1f, S2). Particularly, the "exposure source" information 284 categorizes the drugs in a combination of five classes, namely medical, endogenous, food, 285 personal care, and industrial sources, which was based on the source categorizations from 286 the Chemical Functional Ontology (ChemFOnt) database⁴⁷ and modified manually - by parsing of web pages and scientific literature - to increase compound coverage and improve 287 288 accuracy and consistency. This categorization allows distinguishing endogenous or food 289 sourced molecules (for the non-analogous spectra only). Examples include deoxycholic acid, 290 an endogenous molecule also used for liver disease, and lactitol, a food sweetener also used 291 as a laxative. Using the GNPS Drug Library metadata, such annotations can be separated 292 from those molecules used exclusively as drugs, which have entirely different exposure 293 implications.

294 Through structural and name matches, we extracted the pharmacologic classes of 900 drugs from the U.S. Food and Drug Administration (FDA) and the therapeutic areas, 295 296 therapeutic indications, and mechanisms of action for 3,894 drugs from the Broad Institute Drug Repurposing Hub.²⁸ However, we noticed substantial variability in the extracted 297 298 information (e.g., inconsistent therapeutic areas assigned to drugs within the same 299 pharmacological class; the sulfonamide antimicrobials sulfamethizole, sulfamethazine, and 300 sulfacetamide were categorized as infectious disease, gastroenterology, and ophthalmology, 301 respectively), or insufficient metadata for several drugs (e.g., common therapeutic indications 302 missing). Therefore, this metadata was further manually curated by expert clinical 303 pharmacologists to enhance and clean up the information retrieved from databases. This 304 manual curation increased the metadata coverage to 4,560 drugs. Drugs without associated 305 metadata are typically those that have been withdrawn from the market (e.g., indoprofen), 306 were in drug development but never marketed (e.g., tarafenacin), or are under development 307 but do not yet have regulatory approval (e.g., firsocostat).

308 In total, 735 drugs in the GNPS Drug Library (38,001 spectra) were identified with 309 endogenous or dietary sources. The final metadata of the drug library covers 27 unique 310 therapeutic areas, 571 pharmacological classes, 920 therapeutic indications, and 823 311 mechanisms of action (Figure 1f, S2). Therapeutic areas of neurology/psychiatry, infectious 312 disease, and cardiology have the highest number of included drugs (Figure 1f) and reference 313 spectra (Figure S2). We note that these incidences reflect the availability of the reference 314 spectra but not the prevalence of these drugs in the general population. Combining the exposure source and therapeutic area, we noticed that fewer drugs related to infection and 315 316 neurology/psychiatry have endogenous or food sources, while higher portions of drugs used 317 for gastroenterology (e.g., deoxycholic acid, riboflavin) and dermatology (e.g., salicylic acid, 318 nicotinamide) are endogenous and/or can come from food-derived molecules.

319 In order to assess the utility of the GNPS Drug Library for detecting drugs known to 320 be consumed, we analyzed two pharmacokinetic datasets where healthy individuals received 321 certain drugs followed by time-series sampling. In the first study, 10 participants received a single oral dose of diphenhydramine.⁴⁸ The drug was not detected in plasma and skin 322 samples before administration, but was detected in all individuals post-administration over 323 324 the course of 24 hours (Figure S3a). In plasma, detection frequencies peaked at 1-2 hours 325 (Figure S3a), aligning with the reported time to maximum concentration (~2 hours) for diphenhydramine.⁴⁹ In skin, peak detection occurred at 10-12 hours (Figure S3a), reflecting 326 327 the delayed deposition to skin compared to plasma for orally administered drugs. In the 328 second study, 14 participants received a cocktail of oral caffeine, midazolam, and 329 omeprazole.⁵⁰ These drugs were detected in plasma from 100%, 69%, and 100% of 330 participants, respectively (Figure S3b). The detection frequencies in fecal samples were 331 below 25%. The same participants began a 7-day course of oral cefprozil (administered twice 332 daily; day 2-8) the day after the cocktail drug administration. Cepprozil was detected in fecal samples with increasing frequencies from 0% at day 2 to 43% at day 9. These results 333 334 demonstrate that the GNPS Drug Library can reliably detect consumed drugs and that detection is both biofluid and time dependent. The results also emphasize the need to 335 336 establish medication exposures empirically in the context of the analyzed samples, as clinical 337 studies rarely recorded the time between drug intake and sample collection.

338 Connected with public untargeted metabolomics data, the GNPS Drug Library can 339 reveal distinct drug exposure profiles among individuals with different disease, age, and sex. 340 For different disease studies, we used the human disease ontology identifier (DOID) curated in ReDU, a controlled-vocabulary metadata for public metabolomics datasets,⁵¹ and searched 341 342 for the drugs and drug analogs using fastMASST.⁴² Samples from individuals with 343 inflammatory bowel disease, Kawasaki disease, and dental caries were characterized by high 344 detection frequencies of antibiotics (Figure S4a). Skin swabs of patients with psoriasis were 345 characterized by antifungals. Samples from people with human immunodeficiency virus (HIV) 346 showed high frequency of antivirals, and samples from individuals with Alzheimer's disease 347 were characterized by cardiology and neurology/psychiatry drugs, all consistent with the 348 expected drug usage of people with these diseases (Figure S4a).

349 To investigate drug exposures among different age and sex groups, we profiled 1,993 fecal samples from the American Gut Project,³³ with participants from the United States (US), 350 351 Europe, and Australia with age 46 ± 18 years (range 3-93; 53% female). A total of 75 different 352 drugs were detected; the most frequently detected pharmacologic classes included 353 histamine-1 receptor antagonist (allergy), angiotensin II-receptor blocker (cardiology), ACE 354 inhibitor (cardiology), beta-adrenergic receptor inhibitor (cardiology), statin (lipid-lowering), 355 non-steroidal anti-inflammatory drug (NSAID; analgesics), and selective serotonin reuptake inhibitor (SSRI; antidepressant), which matches with the most commonly prescribed drug 356 classes in these regions (Figure 1g).^{52–54} There were more drugs per individual noted in the 357 US cohort compared to the European and Australian cohorts (chi-square test; χ^2 (8, n = 1,903) 358 = 33, p = 5.3×10^{-5} , Figure S4b). When connected with age and sex information, the drug 359 360 detection agrees with the expected usage patterns of different drug classes (Figure 1h). For 361 example, cardiovascular drugs were detected more frequently with increasing age, while analgesics, antihistamines, and antibiotics were detected across all ages.^{34,35} We also 362 observed that analgesics, such as NSAIDs and paracetamol, were more frequently detected 363 364 in females (chi-square test; χ^2 (1, n = 1,958) = 15.4, p = 8.54 x 10⁻⁵), consistent with the 365 literature,^{36,37} and that drugs for erectile dysfunction were detected only in males. Overall, 366 empirical drug readout using untargeted metabolomics, facilitated by the GNPS Drug Library, 367 demonstrated good specificity among individuals with different disease, age, and sex.

368 The GNPS Drug Library can allow the discovery of previously uncharacterized drug metabolites. To illustrate this, we analyzed fecal samples from the HIV Neurobehavioral 369 Research Center (HNRC) cohort (n = 322; age 55 ± 12 years), which included both people 370 with HIV (n = 222) and people without HIV (n = 100). Among the 17,729 unique MS/MS 371 spectra obtained, 493 were annotated with the GNPS Drug Library. After removing drugs that 372 373 could be from endogenous or food sources (because we cannot assess whether they were 374 given as a medication) and grouping annotations of drugs, metabolites, or analogs, 169 375 unique drugs remained. Antiretroviral drugs (ARVs; drugs for the treatment of HIV), drugs for 376 cardiovascular disease, and drugs for anxiety and depression were the most frequently 377 detected categories (Figure 2a, S5a). Despite the high rates of viral suppression with the 378 advent of antiretroviral therapies (ART; a combination of ARVs to treat HIV), people with HIV have disproportionately high rates of depression and cardiovascular diseases, 55-58 reflected 379 380 in the observation of antidepressants and cardiovascular drugs in these samples.

381 Interestingly, 33% of the drugs were annotated together with their metabolites or 382 analogs, and the occurrences of drug metabolites/analogs aligned with those of the parent drugs (Figure 2a). For example, darunavir (an ARV) had no annotated metabolites but was 383 384 observed with 10 analogs (Figure S5b). Retention time and peak shape analysis indicated 385 that two of the darunavir analogs are in-source fragments (as judged by overlapping retention 386 times),⁵⁹ while the others remain unknown metabolites or isomers of this drug (**Figure S5c,d**). 387 For the analogs that are not in-source fragments, 63-100% (median 96%) of their occurrences 388 were together with the darunavir parent drug. The observations of darunavir analogs without 389 the parent drug are perhaps related to the timepoint of sample collection or to individuals with 390 an ultrarapid metabolizer phenotype, impairing the detection of the parent drug. Nevertheless,

this observation highlights the utility of drug metabolites and analogs to increase the sensitivity of drug exposure readouts via untargeted metabolomics. We note that the HNRC dataset was added to the GNPS/MassIVE public repository after the development of the drug analog library. Therefore, analog mining via existing public metabolomics datasets can facilitate the discovery of uncharacterized metabolites in new data.

396 To further investigate the potential metabolic sources of the observed drug analogs. 397 we cultured darunavir and 12 other drugs with a defined and complex synthetic microbial community of 111 bacterial species commonly found in the human gut.⁶⁰ Except clindamycin 398 399 (an antibiotic), all drugs observed with three or more metabolites/analogs that were present 400 in >10% samples were incubated (10 drugs in total; Table S4); omeprazole, loratadine, and 401 terbinafine were additionally included because their analogs were frequently observed in 402 samples without the respective parent drugs. Shared analogs were observed for 10 of the 13 403 drugs between the fecal samples and the microbial incubations. Among them, 404 metabolites/transformation products were observed for 4 drugs (ritonavir, atorvastatin, 405 abacavir, and omeprazole; Figure 2b, S6), while the rest of the analogs were in-source fragments based on retention time correlation analysis.⁵⁹ The ritonavir, atorvastatin, and 406 407 abacavir analogs increased in intensity with increased microbial incubation time (Figure S6a-408 d), indicating microbial metabolism as a possible source and consistent with their observation 409 in fecal samples. The omeprazole analog (m/z 330.127) appeared to be an abiotic 410 transformation product because it was already present at t=0 cultures, and its intensity 411 decreased with increased incubation time (Figure S6e-g). This is consistent with the fast 412 activation of omeprazole (m/z 346.122), a proton-pump inhibitor and a prodrug, to the reactive 413 sulphenamide product (m/z 330.127) at low pH.⁶¹ Rapid photolysis and hydrolysis of 414 omeprazole has also been reported in abiotic environments with a major deoxygenation transformation product (m/z 330.127).^{62,63} 415



416

417 Figure 2. Drug exposures in the HIV Neurobehavioral Research Center (HNRC) cohort with 418 connections to microbial metabolism and endogenous metabolites. From the HNRC cohort, 322 fecal samples were analyzed with 222 samples from people with HIV and 100 samples from people 419 420 without HIV. a, Peak area visualization of drugs detected with metabolites and analogs. Each column 421 represents one sample and each row represents one drug annotation. Drug annotations were grouped based on the parent drugs and separated by gap spaces. Drug annotations were denoted based on 422 423 their types (as drug, drug metabolites, or drug analogs) and the pharmacologic classes of the parent 424 drugs. All annotated ion/adduct forms of the parent drugs were visualized, leading to multiple rows of 425 parent annotations for some drugs. Asterisks on the drug name mark parent drug annotations 426 confirmed with commercial standards based on retention time and MS/MS spectral matches. Raw 427 peak areas were log-transformed. b, Retention time and MS/MS spectra mirror matches for drug 428 analogs observed in both the fecal samples and the drug microbial incubations. Purple traces

429 represent the fecal samples, while red traces represent the drug microbial incubation. Blue traces 430 represent mixtures of the fecal samples and the microbial incubations at 1:1 volume ratio. The atomic 431 changes of the drug analogs were based on [M+H]⁺ ion of the parent drug. c, Hierarchical clustering 432 of the samples from people with HIV (n = 222) based on detected antiretroviral drugs (ARV). Each row 433 represents one detected ARV, with peak areas summed for the drug, metabolite, and analog 434 detections followed by log-transformation (visualized with the same color scale as panel a). ARVs 435 detected in <10% of samples are not shown. Each column represents one sample, clustered into four 436 groups by hierarchical clustering with Ward's linkage and Euclidean distance. d, Sample-to-sample 437 peak areas of N-acyl lipids in people with HIV, separated by the clusters derived from ARV detections 438 shown in panel c. For each compound, the peak area in each sample was standardized to the 439 maximum value observed across all samples. A non-parametric Kruskal-Wallis test followed by 440 pairwise Wilcoxon test and Benjamini-Hochberg correction for multiple comparisons were performed. 441 P-values < 0.05 were noted in the figure. Boxplots showcase the median value, first (lower) and third 442 (upper) quartiles, and whiskers indicate the error range as 1.5 times the interguartile range.

443

444 The GNPS Drug Library can enable stratification based on drug profiles, which facilitates discovery of connections between drug exposures and endogenous metabolites. 445 446 *N*-acyl lipids are a class of signaling molecules made by host-associated microbiota⁶⁴ that play important roles in the immune system,⁶⁵ memory function,⁶⁶ and insulin regulation of the 447 human body.^{67–69} Our recent ongoing work found that the levels of histamine N-acyl lipids 448 449 differed by HIV serostatus. Specifically, we observed higher levels of histamine-C2:0, 450 histamine-C3:0, and histamine-C6:0 in people living with HIV than people without HIV.⁷⁰ To 451 investigate whether these differences were related to drug exposures, we further stratified samples in this dataset by their ARV exposure profiles. As expected, the ARV profiles clearly 452 separated based on the HIV serostatus (Figure S7a). High intensities of different ARVs were 453 454 observed in fecal samples from people with HIV, while ARVs were only occasionally observed 455 in people without HIV with low intensities. ARVs observed in people without HIV include 456 tenofovir, maraviroc, atazanavir, and raltegravir, which are commonly used for prophylaxis (Figure S7a).^{71,72} To control for the HIV serostatus and investigate the effects of ARV 457 458 exposure, we excluded samples from people without HIV and stratified the people with HIV 459 (n = 222) based on their ARV co-occurrences. Four distinct ARV exposure groups were observed based on hierarchical clustering that agreed well with the different combination 460 461 antiretroviral therapy (cART) regimens (Figure 2c). For example, Group 1 (n = 48), characterized by lamivudine, abacavir, and dolutegravir exposures, corresponded to the 462 dolutegravir/abacavir/lamivudine treatment regimen.⁷³ Group 2 (n = 58) with emtricitabine, 463 464 darunavir, ritonavir, and cobicistat exposures, agreed with the darunavir/ritonavir regimen⁷⁴ 465 and the darunavir/cobicistat/emtricitabine/tenofovir regimen.⁷⁵ Group 3 (n = 79), 466 characterized by emtricitabine and dolutegravir exposures, may be related to the 467 dolutegravir/emtricitabine/tenofovir treatment regimen (Group 2).⁷⁶ Group 4 (n = 37) were 468 without apparent ARV exposures, possibly due to poor adherence, severe comorbidities, HIV 469 elite control, or ARVs not included in the GNPS Drug Library or not amenable with LC-MS/MS 470 detections (Figure 2c). Notably, we observed that the levels of histamine-C2:0 previously associated with HIV serostatus,⁷⁰ along with the levels of eleven other *N*-acyl lipids, were 471

472 significantly different in the four ARV exposure groups (Kruskal-Wallis test, p-value < 0.05: 473 see specific p-value in Figure 2d). This suggests that exposure to different classes of ARV 474 among people with HIV, in addition to HIV serostatus itself, might contribute to changed levels 475 of these N-acyl lipids. We emphasize that these patterns could not have been revealed 476 without the empirical drug readouts from untargeted metabolomics. Clinical research records 477 may not document exposures to individual drugs and often do not provide quantitative 478 information on the exposure levels. For example, metadata for the HNRC cohort on current 479 ARV usage, which is based on self-reports, documented drug usage as "ARV-naïve" (never received ARV), "no ARV" (no current ARV use), "non-HAART" (currently using less than three 480 ARVs), and "HAART" (currently using three or more ARVs). Based on these classifications, 481 no significant differences were observed for the 52 N-acyl lipids detected in these samples 482 483 (Figure S7b). Without the empirical drug readout, enabled by the GNPS Drug Library, the 484 effects of drugs on microbial N-acyl lipid levels would be overlooked.

485 We anticipate the GNPS Drug Library to play a key role in precision medicine by 486 enhancing our understanding of the effects of drugs across a wide range of phenotypes, 487 including endogenous metabolism, gut and skin microbiome, pharmacokinetics, and drug-488 drug interactions. The empirical drug readouts from the GNPS Drug Library can enhance the 489 clinical metadata by providing sample-to-sample comparisons of the relative abundance of 490 individual drugs, which can be flexibly summarized at multiple ontology levels depending on 491 user-defined questions. The mass spectrometry community will play a key role in the 492 evolution of this resource through the continued deposition of reference libraries and 493 expansion of the public metabolomics datasets for analog searches. By harnessing the power 494 of public data and data science-ready metadata, we can unlock opportunities to deepen our 495 understanding of the intricate relationships between xenobiotic exposure and human 496 biological systems.

497 It is important to understand that the use of the GNPS Drug Library holds certain 498 limitations. The current library only supports MS/MS-based annotations to level 2/3 according 499 to the 2007 Metabolomics Standards Initiative.⁷⁷ This generally means that spectra of drug 500 isomers may be annotated as the drug. Key drugs with important clinical implications should 501 be checked for retention time matching and be quantified with analytical standards should the 502 scientific question warrant this. The GNPS Drug Library can only capture drugs that are 503 detectable in the specific biological matrix of choice (e.g., brain samples and urine will likely 504 have different drug exposure readouts) and drugs that are ionizable with the chosen mass 505 spectrometry setup. When constructing the drug analog library, we designed the filters to 506 retain analog spectra that can be as confidently linked to drug exposure as possible, at the 507 likely cost of excluding true positives. For example, metabolism pathways specific to 508 substructures infrequently captured or missing in our customized delta mass list will be 509 excluded. Metabolites shared between drugs that cannot be connected by the applied 510 structural similarity scores (the Tanimoto score)⁷⁸ will be lost. As this is an evolving resource, 511 we encourage the community to not only add to, but also report any inconsistencies in the 512 library and the metadata they may notice.

513

514 Acknowledgements:

515 This project was enabled in part by the Alzheimer's Gut Microbiome Project (AGMP) and the 516 Data Infrastructure and Molecular Atlas for AD: Connection Exposome, Gut Microbiome, and 517 Metabolome supplement funded wholly or in part by the following grants thereto: 518 1U19AG063744 and 3U19AG063744-04S1 and awarded to Dr. Kaddurah-Daouk at Duke 519 University in partnership with multiple academic institutions. As such, the investigators within 520 the AGMP and the Exposome Supplement, not listed specifically in this publication's author's 521 list, provided data along with its pre-processing and prepared it for analysis, but did not 522 participate in analysis or writing of this manuscript. A listing of AGMP Investigators can be 523 found at https://alzheimergut.org/meet-the-team/. A complete listing of ADMC investigators 524 can be found at: https://sites.duke.edu/adnimetab/team/. We also thank the support by NIH 525 for the Maternal and Pediatric Precision in Therapeutics project P50HD106463, the 526 development of tools for structure elucidation R01DK136117, the Collaborative Microbial 527 Metabolite Center U24DK133658. The HIV Neurobehavioral Research Center (HNRC) is 528 supported by Center award P30MH062512 from NIMH. This research was supported in part 529 by the Intramural Research Program of the NIH, National Institute of Environmental Health 530 Sciences (ZIC ES103363). C.B. was supported by the Czech Academy of Sciences PPLZ 531 fellowship number L200552251. V.C.L is supported by Fonds de recherche du Québec -532 Santé (FRQS) Postdoctoral fellowship (335368). N.E.A was supported in part by the National Center for Complementary and Integrative Health of the NIH under award number 533 534 F32AT011475. A.M.C.-R. and P.C.D. were supported by the Gordon and Betty Moore 535 Foundation grant GBMF12120. M.R. was supported by the NIH grant R37 Al126277. T.P. 536 was supported by the Czech Science Foundation (GA CR) grant 21-11563M and by the 537 European Union's Horizon 2020 research and innovation programme under Marie 538 Skłodowska-Curie grant agreement No. 891397.

539

540 **Disclosures**:

541 R.S.: R.S. is a co-founder of mzio GmbH.

542 D.M.: D.M. is a consultant for BiomeSense, Inc., has equity and receives income. The terms

of these arrangements have been reviewed and approved by the University of California, San
Diego in accordance with its conflict of interest policies.

545 R.K.-D.: R.K.-D. is an inventor on a series of patents on use of metabolomics for the diagnosis

546 and treatment of CNS diseases and holds equity in Metabolon Inc., Chymia LLC and 547 PsyProtix.

- 548 M.W.: M.W. is a co-founder of Ometa Labs LLC
- 549 T.P.: T.P. is a co-founder of mzio GmbH.
- 550 R.K.: R.K. is a scientific advisory board member, and consultant for BiomeSense, Inc., has
- equity and receives income. He is a scientific advisory board member and has equity in
- 552 GenCirq. He is a consultant for DayTwo, and receives income. He has equity in and acts as
- a consultant for Cybele. He is a co-founder of Biota, Inc., and has equity. He is a cofounder
- of Micronoma, and has equity and is a scientific advisory board member. The terms of these

555 arrangements have been reviewed and approved by the University of California, San Diego 556 in accordance with its conflict of interest policies.

- 557 S.M.T.: S.M.T. receives research funding from Veloxis Pharmaceuticals.
- P.C.D.: P.C.D. is a scientific advisor and holds equity in Cybele, and bileOmix, and is a
 Scientific Co-founder, and advisor and holds equity in Ometa, Arome, and Enveda with prior
- 560 approval by UC-San Diego.
- 561

562 Author contributions:

563 H.N.Z., C.B., P.C.D. conceptualized the method. H.N.Z., C.B., R.S., T.P. developed the MS/MS library for drugs. H.N.Z., W.B., R.T., R.S. developed the MS/MS library for drug 564 565 analogs. H.N.Z., S.M., H.S., P.R., curated exposure source metadata. K.E.K. curated 566 pharmacological metadata. H.N.Z., R.T., Y.E.A., H.M.-R., N.E.A., A.M.C.-R., P.W.P.G., S.M., 567 I.M., A.C., S.P.T., S.Z., curated drug name match results. H.N.Z., K.E.K., S.L., H.M.-R., L.K., 568 S.X., performed data analyses. Y.E.A., M.S.A., C.W., A.K.J., D.M., helped with data 569 interpretation. V.C.-L., H.N.Z., L.C., C.W., M.R., K.Z. performed microbial incubation 570 experiments. D.H.R., L.X. contributed MS/MS reference spectra. M.S.A., R.J.E., D.F., R.H., 571 J.I., S.L., D.J.M. developed the clinical cohort of human immunodeficiency virus (HIV) 572 infection. M.R.Z.S., M.W. performed ModiFinder analysis. S.G., D.S.W. provided support on 573 exposure source annotation. R.K.-D. supervised the consortium providing access to the 574 Alzheimer's disease cohort and acquired funding. R.K. supervised sample handling and DNA 575 data acquisition for the Alzheimer's disease, HIV and American Gut Project cohorts, and 576 acquired funding. H.N.Z., K.E.K., C.B., P.C.D. drafted the manuscript. P.C.D., S.M.T. 577 acquired funding and supervised this project. All authors reviewed and edited the manuscript.

578

579 Data availability:

580 The MGF spectral files for the GNPS Drug Library and the associated metadata of controlled 581 (.csv) can be downloaded from Zenodo archive under doi: vocabularies 10.5281/zenodo.13892289. The downloaded MGF spectral files can be added to personal 582 583 GNPS folders and used directly for library matching. Data used to validate the empirical drug 584 readouts are publicly available in GNPS/MassIVE under the accession numbers 585 MSV000085944, MSV000084008, and MSV000082493. Data used to profile drug exposures by age and sex are available at MSV000080673. Data for fecal samples from the HNRC 586 587 cohort are available at MSV000092833. Data for the drug bacterial cultures are available at 588 MSV000095331. Data for HNRC fecal samples analyzed with the bacterial cultures are 589 available at MSV000096012. Data for co-migration of the bacterial cultures and fecal samples 590 are available at MSV000096013. Due to human subject protection constraint, metadata for 591 the HNRC cohort will be provided upon request to HNRC: https://hnrp.hivresearch.ucsd.edu. 592

593 **Code availability:** The code used to guery reference spectra of drugs is available on GitHub 594 under the MSⁿ library project²⁶ (https://github.com/corinnabrungs/msn_tree_library). The 595 matches provided code used to filter the drug analog is on GitHub 596 (https://github.com/ninahaoqizhao/Manuscript_GNPS_Drug_Library). The code used for

597datasetanalysiscanbefoundonGitHub598(https://github.com/ninahaoqizhao/Manuscript_GNPS_Drug_Libraryand599https://github.com/kinekvitne/manuscript_drug_library).and

600

601 Supplementary Figure

602 Supplementary Figure S1. Additional examples of structural modification sites of the drug 603 analogs predicted by ModiFinder; Supplementary Figure S2. Overview of the ontology-based 604 drug metadata based on the numbers of reference spectra; Supplementary Figure S3. 605 Empirical drug readout in healthy individuals receiving specific drugs; Supplementary Figure 606 S4. Drug exposure profiles among cohorts with different diseases and geolocations by re-607 analyzing public metabolomics data; Supplementary Figure S5. Drug analog annotations in 608 fecal samples from people with human immunodeficiency virus; Supplementary Figure S6. 609 Drug analogs observed in human fecal samples can be produced by microbial metabolism; 610 Supplementary Figure S7. Comparison of sample clustering based on empirical drug records 611 from the GNPS Drug Library or on clinical metadata; Supplementary Figure S8. Retention 612 time and MS/MS spectra mirror matches for drugs observed in the HNRC cohort with 613 analytical standards.

614

615 Supplementary Table

Supplementary Table S1. Community-curated list of delta mass interpretation; Supplementary Table S2. Delta masses accepted in the drug analog library; Supplementary Table S3. Percentage of drug analogs demonstrating co-occurrence with the parent drugs in MASST search, separated by delta masses; Supplementary Table S4. Source of drugs used in synthetic microbial community incubation; Supplementary Table S5. Bacterial strains used in the six synthetic microbial communities; Supplementary Table S6. Composition of the BHI medium for anaerobic microbial cultures.

624 References

- Vermeulen, R., Schymanski, E. L., Barabási, A.-L. & Miller, G. W. The exposome and health:
 Where chemistry meets biology. *Science* 367, 392–396 (2020).
- 627 2. National Center for Health Statistics (U.S.), Health, United States, 2019 (U.S. Centers for
 628 Disease Control and Prevention, 2021).
- 3. Rappaport, S. M., Barupal, D. K., Wishart, D., Vineis, P. & Scalbert, A. The blood exposome
 and its role in discovering causes of disease. *Environ. Health Perspect.* **122**, 769–774 (2014).
- 4. de la Cuesta-Zuluaga, J., Boldt, L. & Maier, L. Response, resistance, and recovery of gut
 bacteria to human-targeted drug exposure. *Cell Host Microbe* 32, 786–793 (2024).
- 5. Verdegaal, A. A. & Goodman, A. L. Integrating the gut microbiome and pharmacology. *Sci. Transl. Med.* 16, eadg8357 (2024).
- 6. Vich Vila, A. et al. Impact of commonly used drugs on the composition and metabolic function
 of the gut microbiota. *Nat. Commun.* **11**, 362 (2020).
- 637 7. Maier, L. et al. Extensive impact of non-antibiotic drugs on human gut bacteria. *Nature* 555,
 638 623–628 (2018).

- 8. Sinnott, S.-J., Polinski, J. M., Byrne, S. & Gagne, J. J. Measuring drug exposure: concordance
 between defined daily dose and days' supply depended on drug class. *J. Clin. Epidemiol.* 69,
 107–113 (2016).
- 642 9. Althubaiti, A. Information bias in health research: definition, pitfalls, and adjustment methods.
 643 *J. Multidiscip. Healthc.* 9, 211–217 (2016).
- 644 10. Gong, Y. et al. Over-the-counter antibiotic sales in community and online pharmacies, China.
 645 *Bull. World Health Organ.* 98, 449–457 (2020).
- Mackey, T. K. et al. Multifactor quality and safety analysis of antimicrobial drugs sold by online
 pharmacies that do not require a prescription: Multiphase observational, content analysis, and
 product evaluation study. *JMIR Public Health Surveill.* 8, e41834 (2022).
- 12. Essigmann, H. T. et al. Epidemiology of antibiotic use and drivers of cross-border procurement
 in a Mexican American border community. *Front. Public Health* **10**, (2022).
- 13. Walmsley, B. et al. The PrEP You Want: A web-based survey of online cross-border shopping
 for HIV prophylaxis medications. *J. Med. Internet Res.* 21, e12076 (2019).
- 14. Vauquelin, G. & Charlton, S. J. Long-lasting target binding and rebinding as mechanisms to
 prolong in vivo drug action. *Br. J. Pharmacol.* **161**, 488–508 (2010).
- 15. Copeland, R. A. The drug-target residence time model: a 10-year retrospective. *Nat. Rev. Drug Discov.* 15, 87–95 (2016).
- 657 16. Wishart, D. S. et al. DrugBank: a comprehensive resource for in silico drug discovery and
 658 exploration. *Nucleic Acids Res.* 34, D668–D672 (2006).
- 17. Wishart, D. S. et al. DrugBank 5.0: a major update to the DrugBank database for 2018. *Nucleic Acids Res.* 46, D1074–D1082 (2018).
- 18. Ursu, O. et al. DrugCentral: online drug compendium. *Nucleic Acids Res.* 45, D932–D939
 (2017).
- 663 19. DailyMed. <u>https://dailymed.nlm.nih.gov/dailymed/index.cfm</u>.
- 20. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new perspectives
 on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* 45, D353–D361 (2017).
- 666 21. Wan, M. et al. TnT-LLM: Text mining at scale with large language models. Preprint at 667 https://doi.org/10.48550/arXiv.2403.12173 (2024).
- 668 22. Thirunavukarasu, A. J. et al. Large language models in medicine. *Nat. Med.* 29, 1930–1940
 669 (2023).
- 670 23. Clusmann, J. et al. The future landscape of large language models in medicine. *Commun.*671 *Med.* 3, 1–8 (2023).
- 672 24. Prakash, C., Shaffer, C. L. & Nedderman, A. Analytical strategies for identifying drug
 673 metabolites. *Mass Spectrom. Rev.* 26, 340–369 (2007).
- 674 25. Wang, M. et al. Sharing and community curation of mass spectrometry data with Global
 675 Natural Products Social Molecular Networking. *Nat. Biotechnol.* 34, 828–837 (2016).
- 676 26. Brungs, C. et al. Efficient generation of open multi-stage fragmentation mass spectral libraries.
 677 Preprint at https://doi.org/10.26434/chemrxiv-2024-l1tqh (2024).
- 678 27. Kim, S. et al. PubChem 2023 update. *Nucleic Acids Res.* **51**, D1373–D1380 (2023).
- 679 28. Corsello, S. M. et al. The Drug Repurposing Hub: a next-generation drug library and
 680 information resource. *Nat. Med.* 23, 405–408 (2017).
- 29. Zdrazil, B. et al. The ChEMBL Database in 2023: a drug discovery platform spanning multiple
 bioactivity data types and time periods. *Nucleic Acids Res.* 52, D1180–D1192 (2024).

- 683 30. Nahler, G. Anatomical therapeutic chemical classification system (ATC). in *Dictionary of* 684 *Pharmaceutical Medicine* 8–8 (Springer, Vienna, 2009).
- 31. Ross, D. H., Seguin, R. P., Krinsky, A. M. & Xu, L. High-throughput measurement and
 machine learning-based prediction of collision cross sections for drugs and drug metabolites. *J. Am. Soc. Mass Spectrom.* 33, 1061–1072 (2022).
- 688 32. Shahneh, M. R. Z. et al. ModiFinder: Tandem mass spectral alignment enables structural 689 modification site localization. *J. Am. Soc. Mass Spectrom.* (2024).
- 33. McDonald, D. et al. American Gut: an open platform for citizen science microbiome research.
 mSystems 3, e00031-18 (2018).
- 69234. Products Data Briefs Number 347 August 2019.
- 693 https://www.cdc.gov/nchs/products/databriefs/db347.htm (2019).
- 35. Helms, P. J., Ekins Daukes, S., Taylor, M. W., Simpson, C. R. & McLay, J. S. Utility of routinely
 acquired primary care data for paediatric disease epidemiology and pharmacoepidemiology. *Br. J. Clin. Pharmacol.* 59, 684–690 (2005).
- 36. Isacson, D. & Bingefors, K. Epidemiology of analgesic use: a gender perspective. *Eur. J. Anaesthesiol. Suppl.* 26, 5–15 (2002).
- 37. Anthony, M. et al. Gender and age differences in medications dispensed from a national chain
 drugstore. *J. Womens Health* **17**, 735–743 (2008).
- 38. Quinn, R. A. et al. Molecular networking as a drug discovery, drug metabolism, and precision
 medicine strategy. *Trends Pharmacol. Sci.* 38, 143–154 (2017).
- 39. Chen, L. et al. Metabolite discovery through global annotation of untargeted metabolomics
 data. *Nat. Methods* 18, 1377–1385 (2021).
- 40. Yu, N. et al. Nontarget discovery of antimicrobial transformation products in wastewater based
 on molecular networks. *Environ. Sci. Technol.* 57, 8335–8346 (2023).
- 41. Bittremieux, W. et al. Open access repository-scale propagated nearest neighbor suspect
 spectral library for untargeted metabolomics. *Nat. Commun.* 14, 8488 (2023).
- 42. Wang, M. et al. Mass spectrometry searches using MASST. *Nat. Biotechnol.* 38, 23–26
 (2020).
- 43. Mongia, M. et al. Fast mass spectrometry search and clustering of untargeted metabolomics
 data. *Nat. Biotechnol.* 1–6 (2024).
- 44. Creasy, D. M. & Cottrell, J. S. Unimod: Protein modifications for mass spectrometry. *Proteomics* 4, 1534–1536 (2004).
- 45. Kolodnitsky, A. S., Ionov, N. S., Rudik, A. V., Filimonov, D. A. & Poroikov, Vladimir. V.
 HGMMX: Host Gut Microbiota Metabolism Xenobiotics Database. *J. Chem. Inf. Model.* 63, 6463–6468 (2023).
- 46. Keller, B. O., Sui, J., Young, A. B. & Whittal, R. M. Interferences and contaminants
 encountered in modern mass spectrometry. *Anal. Chim. Acta* 627, 71–81 (2008).
- 47. Wishart, D. S. et al. ChemFOnt: the chemical functional ontology resource. *Nucleic Acids Res.*51, D1220–D1229 (2023).
- 48. Panitchpakdi, M. et al. Non-invasive skin sampling detects systemically administered drugs in
 humans. *PloS One* 17, e0271794 (2022).
- 49. Scavone, J. M., Greenblatt, D. J., Harmatz, J. S., Engelhardt, N. & Shader, R. I.
 Pharmacokinetics and pharmacodynamics of diphenhydramine 25 mg in young and elderly
 volunteers. *J. Clin. Pharmacol.* 38, 603–609 (1998).

- 50. Jarmusch, A. K. et al. Enhanced characterization of drug metabolism and the influence of the
 intestinal microbiome: A pharmacokinetic, microbiome, and untargeted metabolomics study. *Clin. Transl. Sci.* 13, 972–984 (2020).
- 51. Jarmusch, A. K. et al. ReDU: a framework to find and reanalyze public mass spectrometry
 data. *Nat. Methods* 17, 901–904 (2020).
- 52. Audi, S. et al. The 'top 100' drugs and classes in England: an updated 'starter formulary' for
 trainee prescribers. *Br. J. Clin. Pharmacol.* 84, 2562–2571 (2018).
- 53. Fuentes, A. V., Pineda, M. D. & Venkata, K. C. N. Comprehension of top 200 prescribed drugs
 in the US as a resource for pharmacy teaching, training and practice. *Pharm. J. Pharm. Educ. Pract.* 6, 43 (2018).
- 54. Wylie, C. E., Daniels, B., Brett, J., Pearson, S.-A. & Buckley, N. A. A national study on
 prescribed medicine use in Australia on a typical day. *Pharmacoepidemiol. Drug Saf.* 29,
 1046–1053 (2020).
- 55. Winston, A. & Spudich, S. Cognitive disorders in people living with HIV. *Lancet HIV* 7, e504–
 e513 (2020).
- 56. Arseniou, S., Arvaniti, A. & Samakouri, M. HIV infection and depression. *Psychiatry Clin. Neurosci.* 68, 96–109 (2014).
- 57. Donati, K. de G., Rabagliati, R., Iacoviello, L. & Cauda, R. HIV infection, HAART, and
 endothelial adhesion molecules: current perspectives. *Lancet Infect. Dis.* 4, 213–222 (2004).
- 58. Heaton, R. K. et al. Twelve-year neurocognitive decline in HIV is associated with
 comorbidities, not age: a CHARTER study. *Brain J. Neurol.* 146, 1121–1131 (2023).
- 59. Schmid, R. et al. Ion identity molecular networking for mass spectrometry-based
 metabolomics in the GNPS environment. *Nat. Commun.* 12, 3832 (2021).
- 60. Cheng, A. G. et al. Design, construction, and in vivo augmentation of a complex gut
 microbiome. *Cell* 185, 3617-3636.e19 (2022).
- 61. Olbe, L., Carlsson, E. & Lindberg, P. A proton-pump inhibitor expedition: the case histories of
 omeprazole and esomeprazole. *Nat. Rev. Drug Discov.* 2, 132–139 (2003).
- 62. Boix, C., Ibáñez, M., Sancho, J. V., Niessen, W. M. A. & Hernández, F. Investigating the
 presence of omeprazole in waters by liquid chromatography coupled to low and high
 resolution mass spectrometry: degradation experiments. *J. Mass Spectrom.* 48, 1091–1100
 (2013).
- 63. Mathew, M., Gupta, V. D. & Bailey, R. E. Stability of omeprazole solutions at various ph values
 as determined by high-performance liquid chromatography. *Drug Dev. Ind. Pharm.* 21, 965– 971 (1995).
- 64. Gentry, E. C. et al. Reverse metabolomics for the discovery of chemical structures from
 humans. *Nature* 626, 419–426 (2024).
- 65. Chang, F.-Y. et al. Gut-inhabiting Clostridia build human GPCR ligands by conjugating
 neurotransmitters with diet- and human-derived fatty acids. *Nat. Microbiol.* 6, 792–805 (2021).
- 66. Mann, A. et al. Palmitoyl serine: An endogenous neuroprotective endocannabinoid-like entity
 after traumatic brain injury. *J. Neuroimmune Pharmacol.* **10**, 356–363 (2015).
- 67. Grevengoed, T. J. et al. *N*-acyl taurines are endogenous lipid messengers that improve
 glucose homeostasis. *Proc. Natl. Acad. Sci.* **116**, 24770–24778 (2019).

- 68. Waluk, D. P., Vielfort, K., Derakhshan, S., Aro, H. & Hunt, M. C. *N*-Acyl taurines trigger insulin
 secretion by increasing calcium flux in pancreatic β-cells. *Biochem. Biophys. Res. Commun.*430, 54–59 (2013).
- 77269. Aichler, M. et al. *N*-acyl taurines and acylcarnitines cause an imbalance in insulin synthesis773and secretion provoking β cell dysfunction in type 2 diabetes. *Cell Metab.* **25**, 1334-1347.e4774(2017).
- 775 70. Mannochio-Russo, H. et al. The microbiome diversifies *N*-acyl lipid pools including short 776 chain fatty acid derived. In preparation.
- 777 71. Inciarte, A. et al. Post-exposure prophylaxis for HIV infection in sexual assault victims. *HIV* 778 *Med.* 21, 43–52 (2020).
- 779 72. Kamitani, E. et al. Growth in proportion and disparities of HIV PrEP use among key
 780 populations identified in the United States national goals: systematic review and meta781 analysis of published surveys. *J. Acquir. Immune Defic. Syndr.* 84, 379 (2020).
- 782 73. Walmsley, S. L. et al. Dolutegravir plus abacavir-lamivudine for the treatment of HIV-1
 783 infection. *N. Engl. J. Med.* 369, 1807–1818 (2013).
- 784 74. Clotet, B. et al. Efficacy and safety of darunavir-ritonavir at week 48 in treatment-experienced
 785 patients with HIV-1 infection in POWER 1 and 2: a pooled subgroup analysis of data from two
 786 randomised trials. *Lancet Lond. Engl.* 369, 1169–1178 (2007).
- 787 75. Huhn, G. D. et al. Darunavir/cobicistat/emtricitabine/tenofovir alafenamide in a rapid-initiation
 788 model of care for human immunodeficiency virus type 1 infection: Primary analysis of the
 789 DIAMOND study. *Clin. Infect. Dis.* **71**, 3110–3117 (2020).
- 790 76. Sax, P. E. et al. Abacavir–lamivudine versus tenofovir–emtricitabine for initial HIV-1 therapy.
 791 *N. Engl. J. Med.* 361, 2230–2240 (2009).
- 792 77. Sumner, L. W. et al. Proposed minimum reporting standards for chemical analysis Chemical
 793 Analysis Working Group (CAWG) Metabolomics Standards Initiative (MSI). *Metabolomics* 3,
 794 211–221 (2007).
- 78. Bajusz, D., Rácz, A. & Héberger, K. Why is Tanimoto index an appropriate choice for
 fingerprint-based similarity calculations? *J. Cheminformatics* 7, 20 (2015).
- 797 79. Bento, A. P. et al. An open source chemical structure curation pipeline using RDKit. J.
 798 Cheminformatics 12, 51 (2020).
- 80. Zuffa, S. et al. microbeMASST: a taxonomically informed mass spectrometry search tool for
 microbial metabolomics data. *Nat. Microbiol.* 9, 336–345 (2024).
- 81. Bittremieux, W. et al. Universal MS/MS visualization and retrieval with the metabolomics
 spectrum resolver web service. Preprint at https://doi.org/10.1101/2020.05.09.086066 (2020).
- 803 82. Bittremieux, W., Laukens, K., Noble, W. S. & Dorrestein, P. C. Large-scale tandem mass
 804 spectrum clustering using fast nearest neighbor searching. *Rapid Commun. Mass Spectrom.*805 e9153 (2021).
- 806 83. Schmid, R. et al. Integrative analysis of multimodal mass spectrometry data in MZmine 3. *Nat.*807 *Biotechnol.* 1–3 (2023).
- 808 84. Nothias, L.-F. et al. Feature-based molecular networking in the GNPS analysis environment.
 809 Nat. Methods 17, 905–908 (2020).
- 810 85. McDonald, D. et al. Extreme Dysbiosis of the Microbiome in Critical Illness. *mSphere* 1, e00199-16 (2016).

- 812 86. Brennan, C. et al. Clearing the plate: a strategic approach to mitigate well-to-well 813 contamination in large-scale microbiome studies. *mSystems* **0**, e00985-24 (2024).
- 814 87. Chambers, M. C. et al. A cross-platform toolkit for mass spectrometry and proteomics. *Nat.*
- 815 *Biotechnol.* **30**, 918–920 (2012).