



University of Dundee

A reference map of potential determinants for the human serum metabolome

Bar, Noam; Korem, Tal; Weissbrod, Omer; Zeevi, David; Rothschild, Daphna; Leviatan, Sigal

Published in:
Nature

DOI:
[10.1038/s41586-020-2896-2](https://doi.org/10.1038/s41586-020-2896-2)

Publication date:
2020

Document Version
Peer reviewed version

[Link to publication in Discovery Research Portal](#)

Citation for published version (APA):

Bar, N., Korem, T., Weissbrod, O., Zeevi, D., Rothschild, D., Leviatan, S., Kosower, N., Lotan-Pompan, M., Weinberger, A., Le Roy, C. I., Menni, C., Visconti, A., Falchi, M., Spector, T. D., The IMI DIRECT consortium, Adamski, J., Franks, P. W., Pedersen, O., & Segal, E. (2020). A reference map of potential determinants for the human serum metabolome. *Nature*, *588*, 135-140. <https://doi.org/10.1038/s41586-020-2896-2>

General rights

Copyright and moral rights for the publications made accessible in Discovery Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

A reference map of potential determinants for the human serum metabolome

Noam Bar^{1,2,+}, Tal Korem^{1,2,3,4,5+}, Omer Weissbrod^{1,2,6}, David Zeevi^{1,2,7}, Daphna Rothschild^{1,2}, Sigal Leviatan^{1,2}, Noa Kosower^{1,2}, Maya Lotan-Pompan^{1,2}, Adina Weinberger^{1,2}, Caroline I Le Roy⁸, Cristina Menni⁸, Alessia Visconti⁸, Mario Falchi⁸, Tim D Spector⁸, The IMI DIRECT consortium, Jerzy Adamski^{9,10,11}, Paul W Franks^{12,13}, Oluf Pedersen¹⁴, Eran Segal^{1,2,*}

Author affiliations

¹Department of Computer Science and Applied Mathematics, Weizmann Institute of Science, Rehovot 7610001, Israel

²Department of Molecular Cell Biology, Weizmann Institute of Science, Rehovot 7610001, Israel

³Department of Systems Biology, Columbia University Irving Medical Center, New York, NY 10032, USA

⁴Department of Obstetrics and Gynecology, Columbia University Irving Medical Center, New York, NY 10032, USA

⁵CIFAR Azrieli Global Scholars program, CIFAR, Toronto, Canada

⁶Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA 02115, USA

⁷Center for Studies in Physics and Biology, The Rockefeller University, New York, NY 10065, USA

⁸Department for Twin Research & Genetic Epidemiology, King's College London, London, UK

⁹Research Unit Molecular Endocrinology and Metabolism, Genome Analysis Center, Helmholtz Zentrum München, German Research Center for Environmental Health, Ingolstädter Landstraße 1, 85764 Neuherberg, Germany

¹⁰Lehrstuhl für Experimentelle Genetik, Technische Universität München, 85350 Freising-Weihenstephan, Germany

¹¹Department of Biochemistry, Yong Loo Lin School of Medicine, National University of Singapore, 8 Medical Drive, Singapore 117597, Singapore

¹²Lund University Diabetes Center, Department of Clinical Sciences, Lund University, Malmö, Sweden

¹³Harvard Chan School of Public Health, Boston, MA, USA

¹⁴The Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and Medical Sciences, University of Copenhagen, 2200 Copenhagen, Denmark

A list of authors and their affiliations appears at the end of the paper.

+These authors contributed equally to this work.

*to whom correspondence should be addressed: eran.segal@weizmann.ac.il

Summary

The serum metabolome contains a plethora of biomarkers and causative agents of various diseases, some endogenously produced and some uptaken from the environment¹. The origin of specific compounds is known, including metabolites that are highly heritable^{2,3} or influenced by the gut microbiome⁴; lifestyle choices such as smoking⁵; and dietary consumption⁶. However, we still have a poor understanding of the key determinants of most metabolites. Here, we measured the levels of 1251 metabolites in serum samples from a unique and deeply phenotyped healthy human cohort of 491 individuals. We applied machine learning algorithms to predict metabolite levels in held-out subjects based on host genetics, gut microbiome, clinical parameters, diet, lifestyle and anthropometric measurements. Notably, we obtained statistically significant predictions for over 76% of the profiled metabolites. Diet and microbiome had the strongest predictive power, and each explained hundreds of metabolites, with over 50% of the variance explained in some metabolites. We further validated microbiome related predictions by showing a high replication rate in two geographically independent cohorts^{7,8} that were not available to us when building the algorithms. Using feature attribution analysis⁹, we uncover specific dietary and bacterial interactions. We demonstrate that some of these interactions may be causal, as some metabolites we predicted to be positively associated with bread increased following a randomized clinical trial of bread intervention. Overall, our results unravel potential determinants of over 800 metabolites, paving the way towards mechanistic understanding of alterations in metabolites under different conditions and to designing interventions for manipulating circulating metabolite levels.

Results

We used mass spectrometry to profile serum samples from 491 healthy individuals for whom we previously collected extensive clinical, lifestyle, dietary, genetics and gut microbiome data¹⁰ (**Extended Data Table 1**; Methods). Our untargeted metabolomics measured the levels of 1251 metabolites, covering a wide range of biochemicals including lipids, amino acids, xenobiotics, carbohydrates, peptides, nucleotides and approximately 30% unidentified compounds (**Extended Data Fig. 1a, Supplementary Table 1**; Methods). Of note, to further interpret unidentified metabolites and aid in biomarker discovery, we designed models that accurately predict their candidate biological pathway (**Supplementary Note 1, Extended Data Fig. 2, Supplementary Table 2-5**). Most measured metabolites were prevalent across the cohort, including 498 metabolites detected in all samples, and 1104 metabolites detected in over 50% of the samples (**Extended Data Fig. 1b**). Following quality control (Methods), 475 individuals with high quality data were included in subsequent analyses.

To validate the accuracy of our metabolomic measurements, we compared the levels of creatinine and cholesterol to measurements obtained using standardized lab tests (Methods) performed independently on a second blood sample taken on the same visit, and found good agreement ($R=0.87$, creatinine; $R=0.79$, cholesterol, **Extended Data Fig. 1c,d**). We further found that samples taken one week apart for 20 participants were significantly correlated (median Spearman $\rho=0.68$, $\text{std}=0.06$), in contrast to samples from different participants that showed no correlation (median Spearman $\rho=0.05$, $\text{std}=0.12$; Methods; **Extended Data Fig. 1e**). These results validate the reproducibility and accuracy of our data, are consistent with previous work showing long-term stability in the human metabolome¹¹, and confirm that this metabolic profile is a unique, person-specific signature.

Robust predictions of serum metabolites

We trained gradient boosted decision trees¹² (GBDT) algorithms that predict metabolite levels in held-out subjects (Methods; **Supplementary Note 2**). GBDT systematically outperformed linear models (Lasso; Methods), with a median and maximum Explained Variance (EV) gain of 8.3 and 43.2%, respectively, for prediction with diet data, and 4.6 and 14.9% for microbiome data (**Extended Data Fig. 3**). Notably, our predictions for over 76% of the metabolite groups tested were statistically significant with at least one feature group, following multiple hypothesis correction (Methods). The largest number of metabolites (335) were significantly explained by diet-related features, and 182 by the microbiome (**Fig. 1a,b**). Our models explained over 10% of the variance for 543 metabolite groups (**Fig. 1d**), with a median EV of 10.2% (range 0-73.5%; **Supplementary Table 6**), and over 50% of the variance explained for 17 metabolites.

We next checked, for each feature group, whether any type of metabolites was enriched with superior predictions (**Fig. 1c**; Methods). We found that clinical data better predicted metabolites classified as blood lipids, amino acids and peptides as opposed to xenobiotics and unidentified compounds, on which it performed worse than on other metabolites. In contrast, microbiome data better explained levels of xenobiotics ($p<10^{-4}$) and unidentified compounds ($p<0.001$), highlighting its potential for explaining the origin of the large number of unidentified compounds. We further found that predictions based on clinical data were significantly correlated with those of diet (Spearman $\rho=0.30$, $p<10^{-20}$), and had a weaker correlation with predictions made with the microbiome ($R=0.21$, $p<10^{-11}$). Predictions based on microbiome data had the highest correlation to predictions based on diet ($R=0.44$, $p<10^{-20}$). Finally, we found that metabolites associated with genetics could not be predicted by other feature groups, and there was a weak correlation between the prediction accuracy of a model containing all other features (“full model”, Methods) and the heritability of metabolites ($R=0.09$, $p<0.005$). Altogether, each feature group was especially informative with respect to a different set of metabolites (**Extended data Fig. 4,5a**).

To estimate the relative predictive power of each feature group across all metabolites, we built models predicting the principal metabolomic components (**Extended data Fig. 5b**). Diet had the largest predictive power, inferring 48.9% of the variance explained by a model containing all features (Methods), while lifestyle factors explained only 1.9% of that EV (**Fig. 1e**). Notably, microbiome data had 30.8% of the full model predictive power. As a large portion of it did not overlap with the predictions of other data, these results highlight its importance in predicting and potentially determining serum metabolites levels.

Replication in external cohorts

To test the robustness and reproducibility of our gut microbiome-based models, we validated their accuracy in two geographically independent cohorts (Methods): 1,004 samples of healthy senior British participants (TwinsUK⁷) and 245 samples of northern Europeans with T2D (IMI-DIRECT⁸; **Extended Data Table 1**). Validation data were not available to us while developing the prediction models, which were trained only on samples from the Israeli cohort. We obtained predictions for metabolites that had statistically significant predictions (FDR<0.1) with $R^2 > 5\%$ in the Israeli cohort (107 metabolites in TwinsUK, 50 in IMI-DIRECT), using only microbiome data from the validation cohorts. Notably, 95 of 107 and 28 of 50 predictions replicated (FDR<0.1) in the healthy TwinsUK and T2D IMI-DIRECT cohorts, respectively, including the top 60 and 28 of the top 50 predictions (**Fig. 2; Supplementary Table 7,8**). We note that most replicated associations are accompanied by a reduction in effect size, which is expected, particularly due to study specific biases. These results indicate that our models unravel robust associations between serum metabolites and the gut microbiome, despite the vast differences between both the populations and the protocols and staff used to assemble these cohorts. Finally, most significant associations between metabolite levels and body mass index (BMI) also replicated in the TwinsUK cohort¹³ with high accuracy (Pearson $R=0.85$, $p < 10^{-10}$, **Extended Data Fig. 5c; Supplementary Table 9**).

Diet and microbiome models are independent

As the diet modulates the gut microbiome¹⁴, we compared the EV of metabolites obtained by models based on either. Although some metabolites, mostly related to coffee consumption, were significantly predicted by both the diet and the microbiome, many were not (**Supplementary Table 10**). Furthermore, adding microbiome data to a diet-based prediction model improved its accuracy in 66% of cases (median and max gain of 2.1 and 62.6% respectively; **Supplementary Table 11**), while adding permuted data reduced performance in 82% of cases (median and max gain of -1.7, and 7.4% respectively; **Extended Data Fig. 5d-f**). Finally, 34 metabolites were significantly predicted only by the microbiome. Altogether, these results suggest that the gut microbiome may be modulating the production of many circulating metabolites, independent of diet.

We next used feature attribution analysis (SHAP⁹; Methods) to interpret these models, infer the drivers of each prediction, and examine interactions between different predictive factors (**Supplementary Note 3, Extended Data Fig. 6**). We found dozens of diet and bacterial features that were strongly predictive of blood metabolites in our models (**Fig. 3a; Extended Data Fig. 7a-g**). Notably, the reported consumption of coffee (both long- and short-term; Methods) had higher importance compared to other dietary features for a large number of xenobiotics and unidentified compounds. This included metabolites from the xanthine metabolism pathway such as paraxanthine (Diet prediction Pearson $R=0.64$, $p < 10^{-20}$) and caffeine ($R=0.68$, $p < 10^{-20}$), as previously reported¹⁵. These metabolites were also significantly predicted using microbiome data, with a *Clostridiceae* species being the main predictor. Another strong feature was long-term fish consumption which accurately predicted the levels of several blood lipids such as 3-carboxy-4-methyl-5-propyl-2-furanpropionic acid (Diet $R=0.71$, $p < 10^{-20}$), a uremic toxin that accumulates in the serum of chronic kidney disease (CKD) patients¹⁶ and was also suggested to prevent and reverses steatosis¹⁷. X-16124 (Microbiome $R=0.77$, $p < 10^{-20}$) and X-11850 ($R=0.7$, $p < 10^{-20}$), are two unidentified metabolites which were accurately predicted by

microbiome data, and specifically by bacteria from the *Eggerthellaceae* family and *Clostridium* genus, respectively. Microbiome data was also highly predictive of the uremic toxins phenylacetylglutamine ($R=0.63$, $p<10^{-20}$) and indoxyl sulfate, ($R=0.37$, $p<10^{-20}$) previously reported in association with cardiovascular disease¹⁸ and CKD¹⁹; these predictions were driven by a *Lachnospiraceae* species.

To assess if a few important taxa are sufficient for accurate prediction, we defined the “main predictor” of each metabolite as the taxa with the maximal mean absolute SHAP value. 19 bacterial taxa were the main predictors for the top 50 microbiome-predicted metabolites (Prediction $R>0.4$; **Supplementary Table 12**). One *Clostridiceae* species was the main predictor of 22 of these, which are also strongly associated with coffee consumption in diet-based models. *Clostridium* sp. CAG:138 was the main predictor of 5 metabolites, including phenylacetylcarnitine ($R=0.47$, $p<10^{-20}$) and p-cresol-glucuronide ($R=0.64$, $p<10^{-20}$) as previously reported²⁰. Other taxa, however, were the “main predictor” of only 1-2 top metabolites, demonstrating that many different bacteria are required to accurately predict the levels of different metabolites. Among the main bacterial predictors of the top 100 metabolites, 89 belonged to Firmicutes, highlighting their strong predictive power. Of note, although Bacteroidetes is the second most abundant phylum in our cohort (**Extended Data Fig. 8a**), none of its species were among these main predictors.

To check whether “main predictors” are sufficient for accurate prediction, we compared, for each metabolite, the accuracy of a full microbiome model to the accuracy of a model based only on its main predictor (**Fig. 3b**). We found that a “main predictor”-based model could only explain a median of 36% of the full model EV. Cinnamoylglycine, for example, is significantly predicted using microbiome data ($R=0.49$, $p<10^{-20}$), yet its main-predictor-based model fails to provide a significant prediction. In contrast, some metabolites are exclusively predicted by a single bacterial species, such as the unidentified metabolite X-16124, for which a model based on an *Eggerthellaceae* species explained 93% of the variance of a full model. Indeed, in 95% of the individuals where this bacteria was detectable, X-16124 was also detectable in serum, compared to only 23% of individuals for which this bacteria was not detected (Mann-Whitney U , $p<10^{-20}$; **Extended Data Fig. 8b**).

Novel genetic-metabolomics associations

Multiple genome wide association studies found that human genetics influence serum metabolites^{2,3,21–24}. The median serum metabolite ACE-heritability, using the traditional twin model, was estimated to be 25%, while the median narrow-sense heritability, based only on discovered genetic loci, was estimated to be 2.1%². As we measured multiple molecules which were not yet identified in these studies, we searched for associations between their levels and single nucleotide polymorphisms (SNPs; **Supplementary Note 4**). Notably, we found 68 statistically significant associations ($p<5\times 10^{-11}$ for all), of which, to the best of our knowledge (Methods), 22 have not been previously reported (**Supplementary Table 13**). These include ethylmalonate, a branched fatty acid which was reported in association with anorexia nervosa²⁵, and was associated with rs2066938, that explained 50% of its variance. This SNP is a 3' UTR variant of the gene UNC119B, which we also found to be associated with butyrylcarnitine, replicating previous reports². Other examples include 2'-O-methyluridine and 2'-O-methylcytidine, both nucleotides involved in pyrimidine metabolism, which we found to associate with a missense variant in the PHYHD1 gene, and were previously reported to be negatively correlated with PHYHD1 expression²⁶. We further found that X-21441, which we predicted as an androgenic steroid (**Supplementary Note 1**), was associated with rs8187710, a missense variant in the ABCC2 gene, explaining 11% of its variance (**Extended Data Fig. 9**). rs8187710 was previously demonstrated to be associated with nonalcoholic fatty liver disease (NAFLD)²⁷. Interestingly, X-21441 was also negatively correlated with age in our cohort (Pearson $R=-0.3$, $p<10^{-7}$), independent of the genotype (**Extended Data Fig. 9c**), suggesting that X-21441 might be an independent metabolic risk factor, mediating the genetic susceptibility of NAFLD and chronological age, a known risk factor for NAFLD²⁸.

Proof-of-concept clinical validation

As a proof-of-concept analysis examining whether some of the feature-metabolite interactions we uncovered may be causal. We used our diet-based models to select the top 5% of metabolites that were either positively or negatively associated with normal consumption of white or whole-wheat bread (**Fig. 4a,b**; Methods). We then analyzed the serum metabolome from the beginning and end of a previously conducted week-long intervention²⁹, in which two randomized groups of ten healthy individuals increased their consumption of either whole-grain sourdough bread or industrial white bread, respectively (**Fig. 4a**; Methods). Notably, we found that metabolites that were positively associated with the consumption of whole-wheat bread in our discovery cohort increased significantly more following the sourdough bread intervention (median fold-change (FC) 1.62) than metabolites that were negatively associated with it (median FC 0.66; Mann-Whitney U , $p < 10^{-10}$; **Fig. 4c**). Moreover, we found no statistically significant differences when comparing the mean FC of these metabolites under the white bread intervention ($p > 0.1$; **Fig. 4c**).

Some metabolites whose levels increased following the sourdough bread intervention were previously linked to consumption of whole-grain wheat flour. A notable example is betaine, an amino acid which has been shown to improve vascular risk factors³⁰ and is also highly abundant in wheat bran and germ³¹. We found that the mean FC in betaine levels in the sourdough bread group was 6.16, as opposed to 0.82 in the white bread group (Mann-Whitney U , $p < 0.004$; **Fig. 4d**). Another example is cytosine, for which the mean FC was far greater in the sourdough bread compared to the white bread group, 78.5 vs. 0.53, respectively ($p < 0.002$; **Fig. 4e**). To the best of our knowledge, unlike betaine, cytosine levels were not previously linked to bread consumption.

A similar analysis using metabolites that were associated with white bread consumption in our cohort could not find significant changes in these metabolites following intervention, potentially due to high baseline white wheat consumption in the typical diet of the study population. Overall, these results suggest that some of the associations that we found here might be causal.

Discussion

Although our cohort is not the largest in which serum metabolomics were measured, it is, to our knowledge, the only one in which these measurements were coupled with such a diverse array of potential determinants. Still, it has several limitations. First, while drug intake was shown to have a large effect on the serum metabolome profile³², our cohort was healthy, with limited drug intake. We are therefore likely underestimating its potential impact on blood metabolites. Second, replication of results are still required for predictions by most factors other than the microbiome. Third, due to the lack of reliable annotations, we have not associated metabolites with specific enzymes; this could be addressed in subsequent experimental studies by focusing on strongly predictive taxa. Finally, since this study is mainly based on observational data, interpretation of interactions should be made with caution, and the associations cannot be considered as causal.

Taken together, our results unravel a comprehensive list of potential determinants for circulating blood metabolites. Many of the associations and interactions detected here replicated previously reported findings, supporting the validity of our results. The vast majority of them, however, are novel, making them a useful resource for future studies, either for improving molecular understanding of health and disease, or for forming the basis of interventional studies aimed at altering the levels of blood metabolites.

Main Text References

1. Psychogios, N. *et al.* The human serum metabolome. *PLoS ONE* **6**, e16957 (2011).
2. Shin, S.-Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat. Genet.* **46**, 543–550 (2014).
3. Long, T. *et al.* Whole-genome sequencing identifies common-to-rare variants associated with human blood metabolites. *Nat. Genet.* **49**, 568–578 (2017).
4. Wikoff, W. R. *et al.* Metabolomics analysis reveals large effects of gut microflora on mammalian blood metabolites. *Proc Natl Acad Sci USA* **106**, 3698–3703 (2009).
5. Xu, T. *et al.* Effects of smoking and smoking cessation on human serum metabolite profile: results from the KORA cohort study. *BMC Med.* **11**, 60 (2013).
6. Playdon, M. C. *et al.* Comparing metabolite profiles of habitual diet in serum and urine. *Am. J. Clin. Nutr.* **104**, 776–789 (2016).
7. Moayyeri, A., Hammond, C. J., Hart, D. J. & Spector, T. D. The UK Adult Twin Registry (TwinsUK Resource). *Twin Res. Hum. Genet.* **16**, 144–149 (2013).
8. Koivula, R. W. *et al.* Discovery of biomarkers for glycaemic deterioration before and after the onset of type 2 diabetes: rationale and design of the epidemiological studies within the IMI DIRECT Consortium. *Diabetologia* **57**, 1132–1142 (2014).
9. Lundberg, S. M., Erion, G. G. & Lee, S.-I. Consistent Individualized Feature Attribution for Tree Ensembles. *arXiv* (2018).
10. Zeevi, D. *et al.* Personalized nutrition by prediction of glycemic responses. *Cell* **163**, 1079–1094 (2015).
11. Yousri, N. A. *et al.* Long term conservation of human metabolic phenotypes and link to heritability. *Metabolomics* **10**, 1005–1017 (2014).
12. Ke, G. *et al.* LightGBM: A Highly Efficient Gradient Boosting Decision Tree. <https://papers.nips.cc/paper/6907-lightgbm-a-highly-efficient-gradient-boosting-decision-tree.pdf> (2017).
13. Cirulli, E. T. *et al.* Profound Perturbation of the Metabolome in Obesity Is Associated with Health Risk. *Cell Metab.* **29**, 488-500.e2 (2019).
14. David, L. A. *et al.* Diet rapidly and reproducibly alters the human gut microbiome. *Nature* **505**, 559–563 (2014).
15. Azam, S., Hadi, N., Khan, N. U. & Hadi, S. M. Antioxidant and prooxidant properties of caffeine, theobromine and xanthine. *Med. Sci. Monit.* **9**, BR325-30 (2003).
16. Tsutsumi, Y. *et al.* Renal disposition of a furan dicarboxylic acid and other uremic toxins in the rat. *J. Pharmacol. Exp. Ther.* **303**, 880–887 (2002).
17. Prentice, K. J. *et al.* CMPF, a Metabolite Formed Upon Prescription Omega-3-Acid Ethyl Ester Supplementation, Prevents and Reverses Steatosis. *EBioMedicine* **27**, 200–213 (2018).
18. Nemet, I. *et al.* A Cardiovascular Disease-Linked Gut Microbial Metabolite Acts via Adrenergic Receptors. *Cell* **180**, 862-877.e22 (2020).
19. Hung, S.-C., Kuo, K.-L., Wu, C.-C. & Tarng, D.-C. Indoxyl sulfate: A novel cardiovascular risk factor in chronic kidney disease. *J. Am. Heart Assoc.* **6**, (2017).
20. Evenepoel, P., Meijers, B. K. I., Bammens, B. R. M. & Verbeke, K. Uremic toxins originating from colonic microbial metabolism. *Kidney Int. Suppl.* S12-9 (2009). doi:10.1038/ki.2009.402
21. Yousri, N. A. *et al.* Whole-exome sequencing identifies common and rare variant metabolic QTLs in a Middle Eastern population. *Nat. Commun.* **9**, 333 (2018).
22. Suhre, K. *et al.* Human metabolic individuality in biomedical and pharmaceutical research. *Nature* **477**, 54–60 (2011).
23. Gieger, C. *et al.* Genetics meets metabolomics: a genome-wide association study of metabolite profiles in human serum. *PLoS Genet.* **4**, e1000282 (2008).
24. Kettunen, J. *et al.* Genome-wide association study identifies multiple loci influencing human serum

- metabolite levels. *Nat. Genet.* **44**, 269–276 (2012).
25. Capo-chichi, C. D. *et al.* Riboflavin and riboflavin-derived cofactors in adolescent girls with anorexia nervosa. *Am. J. Clin. Nutr.* **69**, 672–678 (1999).
 26. Darst, B. F., Lu, Q., Johnson, S. C. & Engelman, C. D. Integrated analysis of genomics, longitudinal metabolomics, and Alzheimer's risk factors among 1,111 cohort participants. *Genet. Epidemiol.* **43**, 657–674 (2019).
 27. Sookoian, S., Castaño, G., Gianotti, T. F., Gemma, C. & Pirola, C. J. Polymorphisms of MRP2 (ABCC2) are associated with susceptibility to nonalcoholic fatty liver disease. *J. Nutr. Biochem.* **20**, 765–770 (2009).
 28. Hamaguchi, M. *et al.* Aging is a risk factor of nonalcoholic fatty liver disease in premenopausal women. *World J. Gastroenterol.* **18**, 237–243 (2012).
 29. Korem, T. *et al.* Bread Affects Clinical Parameters and Induces Gut Microbiome-Associated Personal Glycemic Responses. *Cell Metab.* **25**, 1243-1253.e5 (2017).
 30. Olthof, M. R., van Vliet, T., Boelsma, E. & Verhoef, P. Low dose betaine supplementation leads to immediate and long term lowering of plasma homocysteine in healthy men and women. *J. Nutr.* **133**, 4135–4138 (2003).
 31. Craig, S. A. S. Betaine in human nutrition. *Am. J. Clin. Nutr.* **80**, 539–549 (2004).
 32. Liu, J. *et al.* Integration of epidemiologic, pharmacologic, genetic and gut microbiome data in a drug-metabolite atlas. *Nat. Med.* **26**, 110–117 (2020).

Acknowledgments

We thank past and present members of the Segal group for useful discussions. N.B. received a PhD scholarship for Data Science by the Israeli Council for Higher Education (CHE) via the Weizmann Data Science Research Center and is supported by a research grant from Madame Olga Klein – Astrachan. T.K. is a CIFAR Azrieli Global Scholar in the Humans & the Microbiome Program. E.S. is supported by the Crown Human Genome Center, by D. L. Schwarz, J.N. Halpern and L. Steinberg, and by grants funded by the European Research Council and the Israel Science Foundation. The work leading to this publication has received support from the Innovative Medicines Initiative Joint Undertaking under grant agreement n°115317 (DIRECT), resources of which are composed of financial contribution from the European Union's Seventh Framework Programme (FP7/2007-2013) and EFPIA companies' in kind contribution. We thank Dr. Avirup Dutta for introducing us to the DIRECT consortium dataset.

Author Contributions

N.B. and T.K. conceived the project, designed and conducted all analyses, interpreted the results, wrote the manuscript and are listed in random order. O.W. and D.Z. designed statistical analyses. D.R. and S.L. conducted microbiome analysis. N.K. coordinated and designed data collection. M.L.-P. and A.W. developed protocols, performed microbiome sequencing, and processed serum samples. A.W. designed the project and oversaw sample collection and processing. C.L.R., C.M., A.V. M.F. and T.D.S. performed the replication analysis on the TwinsUK cohort. J.A, P.W.F and O.P performed the replication analysis on the IMI-DIRECT cohort. E.S. conceived and directed the project and analyses, designed the analyses, interpreted the results and wrote the manuscript.

Competing interests

The authors declare no competing interests.

Main Figure Legends

Figure 1 | Diet, gut microbiome, genetics and clinical data predict the levels of most serum metabolites. Figure refers to 5-fold cross validation predictions of metabolite levels based on separate models for each feature group. **(a)** Box and swarm plots (centre, median; box, interquartile range (IQR); whiskers, $1.5 \times \text{IQR}$) showing the EV (R^2) of the top 50 significantly predicted metabolites (when available) of each feature group (group names below panel c). **(b)** Heatmap with color gradient from left to right corresponding to the 95% confidence interval (CI) for EV, for each metabolite (y-axis) by every feature group (x-axis). Only metabolites with significant predictions ($\text{FDR} < 10\%$) are shown, and their number per group is shown on top. P-values and CIs were estimated using bootstrapping (Methods). **(c)** Enrichment of metabolite types in the predictions by each feature group (two-sided Mann-Whitney U test; Methods). Only significant enrichments are shown ($p < 0.05$ after 10% FDR correction). Exact p-values are written in each cell. **(d)** A histogram of the number of metabolites (y-axis) with any value of EV (x-axis) obtained using the full model. Inset shows the EV range of 0.3-0.8. **(e)** The fraction of total EV (x-axis) of each feature group (y-axis) compared to the total EV of a model with all feature groups excluding genetics (full model). Total EV is the sum of the EV of the first 15 metabolite principal components (PCs) weighted by the EV of each PC (Methods).

Figure 2 | Validation of microbiome-based predictions of metabolites in two independent cohorts. (a, c) R^2 of predicted metabolites in our cohort (x-axis), versus the rate of replicated associations in the replication cohorts (a, TwinsUK; b, IMI-DIRECT) computed as the fraction of significant replications out of all predictions with equal or higher prediction R^2 in our cohort (left y-axis; blue; $\text{FDR} < 10\%$), and the cumulative number of metabolites (right y-axis; red). **(b, d)** Spearman correlation between true and predicted levels of metabolites in our cohort (x-axis) versus the same correlations in the replication cohorts (y-axis; b, TwinsUK; d, IMI-DIRECT). Metabolites are colored by the replication success (replicated - blue, not replicated - red; $\text{FDR} < 10\%$).

Figure 3 | Diet and gut microbiome data independently explain a wide range of biochemicals. (a) Subset of heatmap showing the directional mean absolute SHAP values (Methods) of various features (x-axis) computed from 5-fold cross validation models that predict metabolite levels (y-axis) using two separate models, one based on diet and another on gut microbiome data. Positive (negative) SHAP values indicate that higher (lower) feature values lead, on average, to higher predicted values. Shown are the top 100 predicted metabolites using diet and gut microbiome, and the top 30 features by maximum mean absolute SHAP value across all metabolites. See extended heatmap in **Extended Data Fig. 7g**. **(b)** The EV of every metabolite from microbiome-based prediction models (x-axis) compared to using only the top predictor of that metabolite, selected as the feature with the largest mean absolute SHAP value (y-axis). Dashed red palette lines mark different y:x ratios. PAGln, phenylacetylglutamine.

Figure 4 | Metabolites explained by bread increase following an intervention that increases bread consumption. (a) Measuring metabolites and routine white- and whole-wheat consumption. We analyzed samples from the first week of a randomized controlled trial²⁹, in which 10 participants increased consumption of whole-grain sourdough bread and 10 others increased consumption of industrial white bread. **(b)** Histogram of directional mean absolute SHAP values of whole-wheat bread consumption for metabolites computed based on held-out samples from our cohort. The top 5% ($n=59$; blue) positively associated metabolites and the top 5% ($n=59$; red) negatively associated metabolites are used for further analysis. **(c)** Box plots (centre, median; box, IQR; whiskers, $1.5 \times \text{IQR}$) showing the mean FC of the top 5% positively (blue) and negatively (red) associated metabolites, separated by intervention group. They show a significantly higher mean FC for the top 5% positively- vs. negatively- associated metabolites under sourdough bread intervention (two-sided

Mann-Whitney U , $p=5 \cdot 10^{-11}$). **(d-e)** The FC (y-axis) of both betaine (d; two-sided Mann-Whitney U , $p=0.0036$) and cytosine (e; $p=0.0014$) were higher in the sourdough bread group compared to the white bread group.

Methods

Description of cohorts

We analyzed banked samples from two previously collected cohorts^{10,29}, for a total of 491 Israeli individuals. Studies were approved by Tel Aviv Sourasky Medical Center Institutional Review Board (IRB), approval numbers TLV-0658-12, TLV-0050-13 and TLV-0522-10; Kfar Shaul Hospital IRB, approval number 0-73. All participants signed written informed consent forms. Full study designs, including inclusion and exclusion criteria were described elsewhere^{10,29}. In brief, participants in both studies were healthy individuals aged between 18 and 70. The participants answered detailed medical, lifestyle and nutritional questionnaires, provided stool and serum samples for metagenomic sequencing and metabolomics, were genotyped, underwent a comprehensive blood test, and for a period of at least one week, recorded all of their daily activities and nutritional intake in real-time using their smartphones with a specialized app provided to them²⁹. Both blood and stool samples were not taken under strict fasting conditions. 16 samples of participants for which microbiome data was not available to us were excluded from all analyses. Meetings in which participants provided blood samples took place in two different centers, Weizmann (45% of participants) and Tel-Aviv (55% of participants). All meetings in Weizmann took place within the first half of the day, while most meetings in Tel-Aviv took place during the second half of the day (82% of the participants).

Feature groups

The “diet” feature group includes both answers for a detailed food frequency questionnaire (FFQ) aimed at capturing long term dietary habits, and the daily mean consumption of different food types, computed over a week based on real-time logging. In both cases we kept only items which were reported to be consumed at least once by at least 1% of our participants, resulting in 670 different food types from logging, and 141 different items from the FFQ.

The “macronutrients” feature group includes the daily mean consumption of macronutrients (lipids, proteins, carbohydrates), calories and water, calculated from real-time logging.

The “anthropometrics” feature group includes weight, BMI, waist and hips circumference, and waist to hips ratio (WHR).

The “cardiometabolic” feature group includes systolic and diastolic blood pressure, heart rate in beats per minute and a glycemic status as previously described³³.

The “drugs” feature group includes 30 binary features representing the intake of 20 common medications as reported in questionnaires, in addition to 10 medication groups as previously described³³. We included only drugs reported to be used by at least 1% of our participants.

The “clinical data” feature group includes the age and sex of the participants, and the following feature groups described above: anthropometrics, cardiometabolic, and drugs.

The “lifestyle” feature group includes smoking status (current, past), stress levels obtained from questionnaires, and the daily mean sleeping time, exercise time and midday sleep time based on real time logging.

The “time of day” feature is a binary feature indicating whether the sample was taken during the first half of the day.

The “seasonal effects” feature is the month in which the sample was taken. In some analyses we also grouped months by season (Winter: December - February; Spring: March - May; Summer: June - August; Fall: September - November).

The “microbiome” feature group includes bacterial relative abundance calculated both by considering coverage (see below), and by MetaPhlAn2³⁴, as well as the first 10 principal components computed over the

log transformed relative abundance of a bacterial gene catalog³⁵ as previously described^{33,36}. Preprocessing steps are described below.

We further defined a full model that included all of the above.

Metabolomics profiling and preprocessing

Metabolite concentrations were measured in serum samples by Metabolon, Inc., Durham, North Carolina, USA, by using an untargeted LC/MS platform as previously described^{2,37,38}. A total of 540 serum samples were profiled, 19 of which were control samples (technical replicate) pooled from several individuals. The other 521 serum samples belonged to 491 participants.

We removed from further analysis 27 metabolites with less than 10 measurements across our cohort, and 54 metabolites that we found to have significantly different distributions in samples collected in two different recruitment centers (Mann-Whitney U $p < 0.05/1251$; Bonferroni corrected; **Supplementary Table 14**). For the remaining 1170 metabolites, we performed robust standardization (subtracting the median and dividing by the standard deviation) over the log (base 10) transformed levels, followed by clipping outlier samples which were farther than 5 standard deviations. We next used two separate normalization schemes, one for single metabolites, which we subsequently used in the feature attribution analysis, and the second for metabolite groups, which we used for global and enrichment analyses.

For single metabolites, we regressed metabolite levels against storage times (only for metabolites present in at least 50 samples), and finally, imputed missing values as the minimum value per metabolite. For the second scheme, metabolites were grouped by correlation with a Spearman rho threshold of 0.85. This is done in order to handle possible bias resulting from uncertainty of metabolite assignments and a high rate of highly correlated mass spectrometry peaks, and resulted in 1067 metabolite groups, 982 of which are singletons. The value of the metabolite group was set to the mean. The category of each metabolite group was assigned based on majority vote, where unidentified compounds were excluded from the vote unless all metabolites in the group were unidentified.

Microbiome preprocessing

Sample collection, DNA extraction, and sequencing of the samples in this study was previously described^{10,29,33}. Briefly, we used only samples which were collected using swabs, filtered metagenomic reads containing Illumina adapters, filtered low-quality reads and trimmed low-quality read edges. We detected host DNA by mapping with GEM³⁹ to the human genome (hg19) with inclusive parameters, and removed human reads. We subsampled all samples to have 10 million reads.

Bacterial relative abundance estimation was performed by mapping bacterial reads to species-level genome bins (SGB) representative genomes (**Supplementary Table 15**)⁴⁰. We selected all SGB representatives from groups with at least 5 genomes, and for these representatives genomes kept only unique regions as a reference data set. Mapping was performed using bowtie2⁴¹ and abundance was estimated by calculating the mean coverage of unique genomic regions across the 50 percent most densely covered areas as previously described^{36,42}. Feature names include the lowest taxonomy level identified.

Comparing metabolomics to lab tests

We compared the levels of both creatinine and cholesterol which we previously obtained via standard lab tests¹⁰ with their metabolomic levels. Since the lab tests were performed by two different labs, we centered the tests by reducing from the value of each sample the mean of all tests taken in the lab in which it was performed. We then performed a standardization of the resulting measurements. The metabolomic profiling and the lab tests were performed on two samples taken at the same blood draw.

Correlation of metabolic profiles within and between individuals

We compared the Spearman correlations between standardized metabolomic profiles of the same participant taken one week apart ($n=20$) to correlations between standardized metabolomic profiles of different individuals ($n=475$). Each pair of samples taken from the same participant was run in the same metabolomic batch. In the group of different individuals, only pairs of individuals from the same batch were included (resulting in a total of 3835 such pairs), and were further stratified by sex.

Predictive models of metabolite groups

We used gradient boosting decision trees from the LightGBM (version 2.1.2) package¹², in order to predict the levels of 1067 metabolite groups based on 7 feature groups in held-out subjects. In order to estimate the EV of each metabolite group we ran a 5-fold cross validation (CV) model using each feature group as input, and evaluated the results using the coefficient of determination (R^2). For all prediction results except those based on human genetics (Methods) we computed 95% confidence intervals and p-values via 1000 iterations of bootstrapping⁴³. In each bootstrap iteration, we performed a random 5-fold cross validation, where in each fold we randomly sampled (with replacement) a group of subjects from the training set to have the same size as the current training set. We next used this set in order to train our model and evaluated the model's performance on the set of subjects in the remaining fold. Finally, we computed the coefficient of determination between the measured values of the metabolite and the concatenation of the CV's predicted values as obtained from the bootstrapping iteration. We applied the Fisher transformation to the estimations of the explained variance we got from bootstrapping in order to induce normality⁴⁴, and then computed a standard error, and estimated the p-values via the normal CDF using the Wald test⁴⁵, such that our null hypothesis is that the explained variance should distribute normally with zero mean. Confidence intervals were computed empirically from the bootstrapping results. We corrected p-values of predictions for multiple hypotheses using the Benjamini-Hochberg procedure over all feature groups (10% FDR). In all CV and bootstrapping runs we used a fixed and predetermined set of hyperparameters: For the microbiome and diet feature groups: *learning_rate*=0.005, *max_depth*=default, *feature_fraction*=0.2, *num_leaves*=default, *min_data_in_leaf*=15, *metric*=L2, *early_stopping_rounds*=None, *n_estimators*=2000, *bagging_fraction*=0.8, *bagging_freq*=1; for other feature groups: *learning_rate*=0.01, *max_depth*=5, *feature_fraction*=0.8, *num_leaves*=25, *min_data_in_leaf*=15, *metric*=L2, *early_stopping_rounds*=None, *n_estimators*=200, *bagging_fraction*=0.9, *bagging_freq*=5.

Human genetics based prediction models

To obtain the predictions based on human genetics, we used a similar 5-fold CV scheme, in which in every fold we calculated the associations between SNPs and metabolite levels within the training fold, and then trained a model only on the top 10 SNPs which reached genome-wide significance (Bonferroni adjusted). For folds where no SNP reached the significance level, we assigned every sample in the test fold with the mean metabolite level of the training fold. Due to high complexity and running time issues, p-values and confidence intervals were not computed based on bootstrapping, rather we estimated the p-values of the Pearson correlation between the true and predicted metabolite levels. Metabolites for which the R^2 was negative were assigned a p-value of 1.

Testing for SNP associations with metabolites

Genotype processing and imputation of 413 individuals were described previously³³. We performed genome wide associations for single metabolites ($n=1170$) and calculated the p-value and the estimated effect sizes using plink (v1.07). When declaring a genome-wide significance for the SNP-metabolite associations we used

a conservative Bonferroni adjustment procedure to control for the false discovery rate due to the large number of SNPs tested ($p < (5 \times 10^{-8}) / 1170$). We performed all genome wide associations using imputed genotypes.

For named molecules, their chemical identification, super and sub pathways are presented as provided by Metabolon. For unidentified molecules, super and sub pathways are estimated based on our biological pathway classifier. We did our best to scan the available literature for known associations between genetic loci and metabolites before reporting an association as novel. The main resources included the GWAS Catalog⁴⁶ and the GWAS server^{2,22}.

Pathway category enrichment analysis

For each pathway category we used a Mann-Whitney *U* test comparing the prediction accuracy of metabolites from that category compared to prediction accuracy of metabolites from other categories. Direction of enrichment was determined by the sign of the Mann-Whitney *U* test statistic. We considered only metabolite groups for which at least one feature group had a significant prediction (after correcting for multiple hypotheses), resulting in 982 metabolite groups.

Validation of metabolite predictions based on microbiome

We validated the robustness of the associations between the gut microbiome composition and the levels of circulating metabolites in two independent cohorts in which we had access to both metagenomics sequencing. Serum metabolomics in these cohorts were performed using the same Metabolon platform that we used for the discovery cohort. The first validation cohort included 1,004 samples of healthy participants from the TwinsUK cohort⁷, for which there was an average of 0.9 ± 1.3 years gap between the collection of faecal and blood samples. The second validation cohort included 245 samples of participants of European ancestry with type 2 diabetes (T2D) from the IMI-DIRECT consortium⁸. Data from both these validation cohorts were not available to us while developing the prediction models. The metagenomics samples from both cohorts went through the exact same analysis pipeline as our discovery cohort to extract the bacterial features which our prediction models were based on. We then applied our models on these data to obtain the metabolite predictions for both cohorts. Only metabolites which were significantly predicted based only on microbiome data with $R^2 > 5\%$ ($FDR < 0.1$) in our discovery cohort were considered for further analysis (107 metabolites out of 678 in TwinsUK, 50 metabolites out of 261 in IMI-DIRECT). Within every validation cohort, we performed robust standardization (subtracting the median and dividing by the standard deviation) over the log (base 10) transformed levels, followed by clipping outlier samples which were farther than 5 standard deviations, and finally, imputed missing values as the minimum value per metabolite. The analysis of these geographically distinct cohorts holds multiple potential sources of noise, including different methods, centers and staff involved in assembling these cohorts, as well as different cohort demographics, clinical manifestations, different genetic background and dietary and lifestyle preferences. Therefore, we defined a successful replication as one which restores the original ranking of the participants as dictated by the true levels of the metabolite in hand. Hence, for every metabolite, in each validation cohort, we computed the Spearman correlation between its true levels and its predicted levels. A replication was considered significant if the FDR adjusted p-value of the Spearman correlation was lower than 0.1 and the correlation coefficient was strictly positive.

Feature attribution analysis

In order to explain the output of our machine learning models and find specific associations between features and metabolite levels we used SHAP (SHapley Additive exPlanations)⁴⁷, a recently introduced framework for interpreting predictions, which assigns each feature an importance value for a particular prediction. Briefly, for a specific prediction, a feature's SHAP value is defined as the change in the expected value of the model's

output when this feature is observed vs when it is missing. It is computed using a sum that represents the impact of each feature being added to the model averaged over all possible orderings of features being introduced. Shapley values based analysis in gut microbiome data was recently demonstrated to be useful, as it allowed for the estimation of complex contributions of gut microbiome taxa to functional shifts, while maintaining global community composition properties⁴⁸.

Individual SHAP values were computed for held-out subjects in 5-fold CV using the module TreeExplainer (version 0.24.0)^{9,49}, based on models trained only on features from the respective feature group. Before training, we standardized the levels of target metabolites, so that SHAP values from different models would be comparable (they are measured in the same units as the target). In each CV fold we ran a random hyperparameter search consistent of 10 iterations using the module RandomizedSearchCV from sklearn (version 0.20.4), and chose the best model for predicting the held out subjects and computing SHAP values. In all feature attribution analyses we used the ungrouped list of 1170 metabolites.

For every feature, we computed the mean absolute SHAP value across all instances in a specific model, reflecting the mean impact of each feature on the predictions and serving as a feature importance measure. We further used these values to compute directional mean absolute SHAP values, by multiplying them with the sign of the Spearman correlation between the population feature and the target. Here, positive values indicate that higher feature values lead, on average, to higher predicted values, while negative values indicate that lower feature values lead, on average, to higher predicted values.

When performing feature attribution analysis with gut microbiome data as input, we only included the relative abundance of SGB representative genomes as features, taking only features which were present in over 5% of the samples, resulting with 753 bacterial taxa. When using diet as input, we only considered features which were present in at least 5% of the samples, resulting with 398 food types from logging and items from the FFQ.

Comparing gradient boosting decision trees with a linear model

We compared the EV of every single metabolite obtained for a GBDT and a Lasso regression model. The EV of all models were calculated in 5-fold CV, where in each fold we ran a hyperparameter search consistent of 10 iterations as described above. We used LightGBM as the GBDT model, and Lasso regression (sklearn, version 0.20.4) as the linear model, since its regularization scheme is better suited for a large number of features, as in the case of diet and gut microbiome composition. Since GBDT handles missing values well, we first imputed all missing values as the median of each feature to assure a fair comparison. When applying the models on the microbiome data, we used \log_{10} transformed values.

Estimating relative predictive power of feature groups

In order to estimate the relative predictive power of different feature groups we first applied a principal component analysis over the metabolite groups data to get the first 400 PCs which constitute >99% of the total variance in the data (**Extended Data Fig. 5b**). We then used 5-fold CV prediction models as described above to predict the PCs based on the different feature groups independently. As baseline, we used the full model, which consists of all features combined to predict the levels of the PCs, and estimated the overall fraction of variance explained by: $\frac{\sum_{i=1}^{nPC} EV_i \times PC_i}{\sum_{i=1}^{nPC} PC_i}$, where EV_i is the fraction of EV that the model recovers for PC i . PC_i is the fraction of variance that PC i explains out of the overall variation in the data. nPC is the number of the first PCs, those which capture the most variation. For the features we have collected, we defined this sum obtained for the full model as the total explainable variance in circulating blood metabolites. Next, for every feature group we computed a similar expression and calculated the relative predictive power by dividing this expression by that of the full model. The estimates we present are for $nPC = 15$, as the overall EV of the

full model that we estimated using the first 15 PCs constitutes over 97% of the overall EV of the full model based on all 400 PCs.

Biological sub pathway prediction

We used gradient boosting decision trees from the LightGBM (version 2.1.2) package¹², in order to build a multiclass classifier to predict the biological sub pathways of metabolites as annotated by Metabolon. When developing the classifier, we only considered named metabolites from biological sub pathways which include over 10 metabolites each in our data, resulting with 28 sub pathways covering a total of 572 named molecules (sub pathway size range 11-44). The rationale behind this is that we tried to find the balance between covering as many metabolites and types of metabolites possible while keeping the number of classes reasonable.

We trained our model in a leave-out-out CV scheme, where in every training fold we used 20% of the training samples as internal validation to perform an early stopping of 50 rounds. We then obtained a soft max of size 28 per metabolite, representing the probabilities of every metabolite being labeled as one of the 28 sub pathways. For the prediction of the unidentified molecules, we used a model trained once using all 572 metabolites. The features used for the training of the model included the normalized levels of metabolites across our main discovery cohort, the mean raw count of the metabolite and the fraction of missing values across the discovery cohort. In addition, to capture the associations between metabolites and their predictive features, we included the directional mean absolute SHAP values for every pair of metabolite-feature computed from the “full model” as described above. The final vectors of probabilities were determined as an ensemble of three models, the first, trained only on the SHAP values, the second, trained only on metabolite levels, means and fraction of missing values, and the third, trained on all combined. Finally, the mean of these three models was computed.

When evaluating the performance of our classifier on the named labeled molecules, we concatenated all vectors of probabilities resulting from the leave-one-out procedure. For every sub pathway we computed a classification report including the classification precision ($TP / (TP + FP)$), recall ($TP / (TP + FN)$) and f1-score ($2 * (precision * recall) / (precision + recall)$), to account for the imbalanced class sizes. The overall accuracy was computed as the fraction of metabolites with correctly assigned labels out of all metabolites from all sub pathways which were included in the training phase. In all runs we used a fixed and predetermined set of hyperparameters (*objective=multiclass, num_leaves=25, max_depth=4, learning_rate=0.005, bagging_fraction=0.8, feature_fraction=0.8, bagging_freq=1, bagging_seed=2018, class_weight=balanced, n_estimators=2000, early_stopping_rounds=50*); TP, True Positive; FP, False Positive; FN, False Negative.

Characterization of unidentified metabolites by Metabolon

Characterization of unidentified metabolites was done as previously described²¹. Briefly, identification of tentative structural features for unidentified biochemicals incorporates a detailed analysis of mass spec data, i.e., gathering information such as the accurate monoisotopic mass, the elution time and fragmentation pattern of the primary ion, and correlation to other molecules. The accurate monoisotopic mass is used to identify a likely structural formula for the unidentified biochemical, which is then used to search against chemical structure databases. When a candidate structure fits the accurate monoisotopic mass and fragmentation data, an authentic standard is commercially purchased or synthesized (when possible). Confirmation of a proposed structure is based on a match to three primary criteria, including co-elution with the unidentified molecule of interest, and a high degree match to both the accurate monoisotopic mass and fragmentation pattern.

Interaction networks

We used a graphical layout in order to visualize the associations of features with the levels of metabolites. The nodes are either metabolites or features, and the edges are the directional mean absolute SHAP values

computed from models trained only on features from the respective feature group as described above. All networks were constructed using Cytoscape⁵⁰. The threshold for presenting SHAP values as edges was determined as 0.12, keeping the network sparse enough for convenience of visualization.

Analysis of bread intervention

In order to find the associations between metabolite levels and the consumption of both types of bread in the study cohort we computed the directional mean absolute SHAP values of the reported consumption of both white and whole-wheat bread for all metabolites. The SHAP values were computed in cross validation from models based only on the reported consumption of each type of bread. We ranked the metabolites according to their directional mean absolute SHAP value for each type of bread and used the top 5% positively and negatively driven metabolites for further analysis. The prediction models were constructed using 458 samples of distinct individuals, a subset of our cohort from which we excluded all samples of individuals which participated in the intervention study.

For each metabolite in every individual, we computed the FC of metabolite levels between the samples taken at the end of the first week of intervention and the start of that week. Prior to computing FC we imputed missing values with the minimum per metabolite and standardized their log (base 10) transformed levels. Furthermore, for each intervention group, we computed the mean FC of every metabolite based on the 10 samples from that group. We then compared the mean FC of the top 5% positively and negatively driven metabolites mentioned above within each intervention group by performing a rank sum test (two-sided Mann-Whitney U) over the mean FC.

For comparing the FC of betaine and cytosine between the two intervention groups, we used a two-sided Mann-Whitney U test.

LMM-based estimates of the explained variance of metabolites using gut microbiome

For the in-sample estimation of EV for metabolites based on gut microbiome we used a linear mixed model framework that we had recently developed³³. Briefly, we used GCTA⁵¹, a tool used in statistical genetics for the estimating of SNP-based genetic kinship. Instead of a matrix of host SNPs, as is commonly used in GCTA, we used a kinship matrix computed over the presence-absence of microbial species which were also used as features in the out-of-sample prediction models. We added the storage time as a covariate to the model. P-values were computed using RL-SKAT⁵².

Statistical analysis

For all statistical analysis and prediction models we used Python 2.7.8 with packages: pandas 0.23.4, numpy 1.14.2, scikit-learn 0.20.4, scipy 1.1.0, shap 0.24.0, LightGBM 2.1.2.

Data Availability

The raw metagenomic sequencing data is available from the European Nucleotide Archive (ENA; <https://www.ebi.ac.uk/ena>): PRJEB11532, PRJEB17643, and for the TwinsUK: PRJEB32731. The raw metabolomics data and the phenotypic data is available from the European Genome-phenome Archive (EGA; <https://ega-archive.org/>): EGAS00001004512. Known links between genetic loci and serum metabolites were taken from the GWAS Catalog⁴⁶ (<https://www.ebi.ac.uk/gwas/>) and the GWAS server^{2,22} (<http://metabolomics.helmholtz-muenchen.de/gwas/>).

Code Availability

Analysis source code is available at <https://github.com/noambar/SerumMetabolomePredictions>.

Methods References

33. Rothschild, D. *et al.* Environment dominates over host genetics in shaping human gut microbiota. *Nature* **555**, 210–215 (2018).
34. Segata, N. *et al.* Metagenomic microbial community profiling using unique clade-specific marker genes. *Nat. Methods* **9**, 811–814 (2012).
35. Li, J. *et al.* An integrated catalog of reference genes in the human gut microbiome. *Nat. Biotechnol.* **32**, 834–841 (2014).
36. Zeevi, D. *et al.* Structural variation in the gut microbiome associates with host health. *Nature* (2019).
37. Bridgewater BR, E. A. High Resolution Mass Spectrometry Improves Data Quantity and Quality as Compared to Unit Mass Resolution Mass Spectrometry in High-Throughput Profiling Metabolomics. *Metabolomics* **04**, (2014).
38. Zierer, J. *et al.* The fecal metabolome as a functional readout of the gut microbiome. *Nat. Genet.* **50**, 790–795 (2018).
39. Marco-Sola, S., Sammeth, M., Guigó, R. & Ribeca, P. The GEM mapper: fast, accurate and versatile alignment by filtration. *Nat. Methods* **9**, 1185–1188 (2012).
40. Pasolli, E. *et al.* Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* **176**, 649-662.e20 (2019).
41. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
42. Korem, T. *et al.* Growth dynamics of gut microbiota in health and disease inferred from single metagenomic samples. *Science* **349**, 1101–1106 (2015).
43. Efron, B. & Tibshirani, R. J. *An Introduction to the Bootstrap*. (Chapman and Hall, 1994). doi:10.1007/978-1-4899-4541-9
44. Fisher, R. A. Frequency Distribution of the Values of the Correlation Coefficient in Samples from an Indefinitely Large Population. *Biometrika* **10**, 507 (1915).
45. Wald, A. Sequential tests of statistical hypotheses. *Ann. Math. Statist.* **16**, 117–186 (1945).
46. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
47. Lundberg, S. & Lee, S.-I. A Unified Approach to Interpreting Model Predictions. *arXiv* (2017).
48. Manor, O. & Borenstein, E. Systematic characterization and analysis of the taxonomic drivers of functional shifts in the human microbiome. *Cell Host Microbe* **21**, 254–267 (2017).
49. GitHub - slundberg/shap: A unified approach to explain the output of any machine learning model. at <<https://github.com/slundberg/shap>>
50. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
51. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
52. Schweiger, R. *et al.* RL-SKAT: An Exact and Efficient Score Test for Heritability and Set Tests. *Genetics* **207**, 1275–1283 (2017).

Extended Data Legends

Extended Data Table 1 | Basic characteristics and demographics of our main and validation cohorts.

Extended Data Figure 1 | Accurate and reproducible untargeted serum metabolomics. (a) Breakdown of the 1251 measured metabolites by type. (b) Number of samples (y-axis) in which each metabolite (x-axis) was identified, sorted by prevalence. (c-d) Mass-spectrometry measurements (y-axis) versus standardized lab tests results (x-axis; Methods) for creatinine (a; Pearson $R=0.87$, $p<10^{-20}$) and cholesterol (b; $R=0.79$, $p<10^{-20}$). (e) Spearman correlations (y-axis; centre, median; box, IQR; whiskers, $1.5\times IQR$) between standardized metabolomic profiles (Methods) of different individuals ($n=475$; median Spearman ρ 0.05, standard deviation [std]=0.12) stratified by sex, and between standardized metabolomic profiles of the same participant ($n=20$; median Spearman ρ 0.68, $std=0.06$) taken one week apart. C&V, cofactors and vitamins; a.u., arbitrary units.

Extended Data Figure 2 | Biological sub pathway prediction of unidentified molecules. Figure panels refer to the results of a leave-one-out cross validation prediction model of metabolites' sub pathways based on their normalized levels, raw mean, percentage of missing values, and SHAP values (Methods). Results shown are for a model trained using only sub pathways which include over 10 molecules in our data (28 sub pathways, 572 named metabolites). (a) The overall accuracy of the sub pathway classifier (y-axis) when a success is considered as having the true label in one of the top k predictions (x-axis). (b) The log loss of the classifier (y-axis) computed over the resulting soft max (raw probabilities; blue) and a dichotomous matrix where for every metabolite we only keep the top predicted sub pathway as 1 and zero-out all other predictions (red). (c) The overall accuracy of the model (left y-axis; blue) and the corresponding fraction of metabolites (right y-axis; red) when considering only metabolites for which the classifier predicted a maximal probability above some threshold (x-axis). (d) A confusion matrix showing the predicted sub pathways (x-axis), determined as the label with the highest probability per metabolite, versus the true annotated sub pathways (y-axis). Each cell in the matrix counts the number of metabolites of a certain true sub pathway (y-axis) which were assigned with some predicted sub pathway (x-axis) by our model. The rightmost column is the sum of every row and represents the number of metabolites annotated for every sub pathway. The matrix is ordered by the higher order biological pathway (super pathway). Cell colors are log scaled. (e) Classification results summarizing the f1-score, precision and recall per sub pathway. Rows correspond to the sub pathway annotation in panel d. (f) For every sub pathway (y-axis) shown are the fraction of metabolites truly annotated as such (black), predicted as such by the classifier (blue; out of the named molecules in the support of the model), and the fraction of unidentified molecules predicted as such (out of all unidentified molecules). M., Metabolism; Xeno., Xenobiotics; Ptds, Peptides; AAs, Amino Acids.

Extended Data Figure 3 | Comparative analysis of linear versus nonlinear models and in-sample versus out-of-sample predictions. (a) Metabolite prediction R^2 of GBDT vs Lasso regression models using diet data. Shown are only metabolites for which both models achieved significant predictions with R^2 above 0.05. (b) Histogram of the differences between the R^2 of GBDT compared to Lasso regression using the diet data. (c) The levels of the metabolite hydroxy-CMPF* (y-axis; centre, median; box, IQR; whiskers, $1.5\times IQR$) vs the monthly consumption of cooked, baked or grilled fish as reported in a food frequency questionnaire. The comparison of Spearman and Pearson correlation coefficients suggests that the relationship between the metabolite and the numerical values of the question are monotonic yet non-linear, which explains why GBDT performs better in predicting the levels of hydroxy-CMPF* from diet data. The x-axis is not in scale. (d, e) Same as a-b for microbiome. (f) Estimations of gut microbiome explainability (b^2) of metabolite levels obtained via applying a linear mixed model on the bacterial species composition as previously described (y-axis) versus the explained variance (R^2) of metabolites from out-of-sample prediction models based on the same gut

microbiome data. Shown are only metabolites with significant b^2 estimates (5% FDR). **(g)** Histogram of the differences between the b^2 estimates and the R^2 of out-of-sample prediction using the gut microbiome data. GBDT, Gradient Boosting Decision Trees; a.u., arbitrary units.

Extended Data Figure 4 | Comparison of explained variance of metabolites for every pair of feature groups. Every panel shows a dot plot of the explained variance of the metabolite groups (y-axis) from models based on every pair of feature groups (x-axis). Panels on the diagonal shows the marginal distribution of explained variance of metabolite groups for a certain feature group.

Extended Data Figure 5 | Comparative analysis of different feature groups. **(a)** Spearman correlations computed between the EV of metabolites for every pair of feature groups. **(b)** The proportion of variance explained by each of the first 400 principal components (left y-axis; black) and their cumulative EV (right y-axis; blue). **(c)** R^2 multiplied by the sign of the Pearson correlation coefficient (x-axis) between metabolite levels and BMI in our study, versus the mean R^2 multiplied by the sign of the Pearson correlation coefficient (y-axis) of BMI associated metabolites recently reported by a different group¹³. Shown are 36 (out of 49) BMI associated metabolites that were also measured in this cohort. P-value for the Pearson correlation, $p=7 \cdot 10^{-11}$. Line and shaded coloring represent the fitting of a linear model and the 95% confidence interval. **(d)** The EV of every metabolite from prediction models based on the gut microbiome (x-axis) versus diet (y-axis). Dashed red line is $y=x$. **(e)** Same for prediction models based on both gut microbiome and diet (x-axis) compared to using only diet (y-axis). **(f)** Same for prediction models based on diet and permuted gut microbiome (x-axis) compared to using only diet (y-axis).

Extended Data Figure 6 | Networks of interactions between phenotypes explain diverse metabolites. Interactions between features from different feature groups predictive of similar metabolites are presented in a graphical layout, in which nodes are either metabolites or features, and edges are the directional mean absolute SHAP values (Methods) computed from models trained only on features from the respective feature group. Circular nodes - metabolites; predictive feature nodes - squares; both colored by relevant categories. Shown are only edges with a mean absolute SHAP value greater than 0.12. **(a)** Network of associations for the following feature groups: macronutrients, diet, microbiome, lifestyle, drugs and seasonal effects. **(b)** A large group of metabolites whose predictions are mainly driven by the reported consumption of coffee and the relative abundance of a bacteria from the Clostridiales order. **(c)** Metabolites explained by seasonal fruit consumption. **(d)** Selected examples of interactions between metabolites and features in predictive models.

Extended Data Figure 7 | Specific dietary features and bacterial taxa underlie the accurate prediction of circulating metabolites. **(a-f)** Predicted (y-axis) vs measured (x-axis) levels (arbitrary units) of X-16124 (a; Pearson $R=0.77$, $p<10^{-20}$), phenylacetylglutamine (b; $R=0.63$, $p<10^{-20}$), p-cresol-glucuronide (c; $R=0.64$, $p<10^{-20}$), caffeine (d; $R=0.68$, $p<10^{-20}$), hydroxy-CMPF (e; $R=0.72$, $p<10^{-20}$) and stachydrine (f; $R=0.5$, $p<10^{-20}$). Predictions of a-c are based only on microbiome data, and colored by the relative abundance of the bacterial taxa having the highest mean absolute SHAP value for each metabolite. Predictions of d-f are based only on diet data, and colored by the reported consumption of the dietary item having the highest mean absolute SHAP value for each metabolite. P-values for prediction were estimated via bootstrapping. **(g)** Heatmap showing the directional mean absolute SHAP values (Methods) of various features (x-axis) computed from 5-fold cross validation models that predict metabolite levels (y-axis) using two separate models, one based on diet and another on gut microbiome data. Positive (negative) SHAP values indicate that higher (lower) feature values lead, on average, to higher (higher) predicted values. Shown are the top 150 predicted metabolites using diet and gut microbiome, and the top 40 features by maximum mean absolute SHAP value across all metabolites. C&V, Cofactors and vitamins; AAs, Amino Acids.

Extended Data Figure 8 | Distribution of phyla and a taxa from the *Eggerthellaceae* family. (a) Stacked bar plots per sample (x-axis) showing the relative abundance of bacterial phyla (y-axis). Samples are sorted by the relative abundance of the most abundant phylum, *Firmicutes*. *Bacteroidetes* is the second most abundant phylum in our cohort. Relative abundance of a phylum is computed as the sum over relative abundances of all bacterial features belonging to that phylum. **(b)** The levels of the unidentified compound X-16124 in individuals for which the bacterial taxa from the *Eggerthellaceae* family was detectable in stool versus individuals for which it was not ($p < 10^{-20}$, two-sided Mann-Whitney *U*).

Extended Data Figure 9 | The unidentified molecule X-21441 associates with rs8187710 independent of age. (a) A table showing the coefficients, standard errors and p-values resulted from a multiple linear regression model with levels of the unidentified molecule X-21441 as the dependent variable, the allele dosage of rs8187710 (0-2) and age (years) as the independent variables: $y_{X-21441} = constant + \beta_1 \circ rs8187710 + \beta_2 \circ Age$. **(b)** The levels of X-21441 (y-axis) versus the genotype of the participants (x-axis). Number of participants with each genotype is indicated below the tick labels. The explained variance of X-21441 by rs8187710 as estimated using plink (Methods) is indicated on the upper right corner of the panel. **(c)** The levels of X-21441 (y-axis; centre, median; box, IQR; whiskers, 1.5×IQR) versus the age of the participants (x-axis) colored by genotype of participants. Line and shaded coloring represent the fitting of a linear model and the 95% confidence interval. SE, Standard Error; a.u., arbitrary units.

Additional DIRECT consortium members

Henrik Vestergaard^{14,15}, Manimozhiyan Arumugam¹⁴, Torben Hansen¹⁴, Kristine Allin¹⁴, Tue Hansen¹⁴, Mungwan Hong¹⁶, Jochen Schwenk¹⁶, Ragna Haussler¹⁶, Matilda Dale¹⁶, Toni Giorgino¹⁷, Marianne Rodriguez¹⁸, Mandy Perry¹⁹, Rachel Nice¹⁹, Timothy McDonald^{19,20}, Andrew Hattersley²⁰, Angus Jones²⁰, Ulrike Graefemody²¹, Patrick Baum²², Rolf Grempler²², Cecilia Engel Thomas^{23,24,25}, Federico De Masi^{24,25}, Caroline Anna Brorsson^{24,25}, Gianluca Mazzoni^{24,25}, Rosa Allesøe^{24,25}, Simon Rasmussen^{24,25}, Valborg Gudmundsdóttir^{24,25}, Agnes Martine Nielsen^{24,25}, Karina Banasik^{24,25}, Konstantinos Tsirigos^{24,25}, Birgitte Nilsson^{24,25}, Helle Pedersen^{24,25}, Søren Brunak^{24,25}, Tugce Karaderi^{24,25}, Agnete Lundgaard^{24,25}, Joachim Johansen^{24,25}, Ramneek Gupta²⁵, Peter Wad Sackett²⁵, Joachim Tillner²⁶, Thorsten Lehr²⁷, Nina Scherer²⁷, Christiane Dings²⁷, Iryna Sihinevich²⁷, Heather Loftus²⁸, Louise Cabrelli²⁸, Donna McEvoy²⁹, Andrea Mari³⁰, Roberto Bizzotto³⁰, Andrea Tura³⁰, Leen 't Hart^{31,32,33}, Koen Dekkers³², Nienke van Leeuwen³², Slieker Roderick^{32,33}, Femke Rutters³³, Joline Beulens³³, Giel Nijpels³³, Anitra Koopman³³, Sabine van Oort³³, Lenka Groeneveld³³, Leif Groop³⁴, Petra Elders³⁵, Ana Viñuela³⁶, Anna Ramisch³⁶, Emmanouil (Manolis) Dermitzakis³⁶, Beate Ehrhardt³⁷, Christopher Jennison³⁷, Philippe Froguel^{38,39}, Mickael Canouil³⁹, Amelie Bonneford³⁹, Ian McVittie⁴⁰, Dianne Wake⁴⁰, Francesca Frau⁴¹, Hans-Henrik Staerfeldt⁴², Peter Sackett⁴², Kofi Adragni⁴³, Melissa Thomas⁴³, Han Wu⁴³, Imre Pavo⁴⁴, Birgit Steckel-Hamann⁴⁴, Henrik Thomsen⁴⁵, Giuseppe (Nick) Giordano⁴⁶, Hugo Fitipaldi⁴⁶, Martin Ridderstråle⁴⁶, Azra Kurbasic⁴⁶, Naeimeh Atabaki Pasdar⁴⁶, Hugo Pomares-Millan⁴⁶, Pascal Mutie⁴⁶, Robert Koivula^{46,47}, Nicky McRobert⁴⁷, Mark McCarthy^{47,48,49}, Agata Wesolowska-Andersen⁴⁷, Anubha Mahajan⁴⁹, Moustafa Abdalla⁴⁹, Juan Fernandez⁴⁹, Reinhard Holl⁵⁰, Alison Heggie⁵¹, Harshal Deshmukh⁵¹, Anita Hennige⁵², Susaana Bianzano⁵², Barbara Thorand^{53,54}, Sapna Sharma^{54,55}, Harald Grallert^{54,55}, Jonathan Adam⁵⁵, Martina Troll⁵⁵, Andreas Fritsche⁵⁶, Anita Hill⁵⁷, Claire Thorne⁵⁷, Michelle Hudson⁵⁷, Teemu Kuulasmaa⁵⁸, Jagadish Vangipurapu⁵⁸, Markku Laakso⁵⁸, Henna Cederberg⁵⁸, Tarja Kokkola⁵⁸, Yunlong Jiao⁵⁹, Stephen Gough⁵⁹, Neil Robertson⁵⁹, Helene Verkindt⁶⁰, Violeta Raverdi⁶⁰, Robert Caiazzo⁶⁰, Francois Pattou⁶⁰, Margaret White⁶¹, Louise Donnelly⁶¹, Andrew Brown⁶¹, Colin Palmer⁶¹, David Davtian⁶¹, Adem Dawed⁶¹, Ian Forgie⁶¹, Ewan Pearson⁶¹, Hartmut Ruetten⁶², Petra Musholt⁶², Jimmy Bell⁶³, Louise Thomas⁶³, Brandon Witcher⁶³, Mark Haid⁶⁴, Claudia Nicolay⁶⁵, Miranda Mourby⁶⁶, Jane Kaye^{66,67}, Nisha Shah⁶⁶, Harriet Teare⁶⁶, Gary Frost⁶⁸, Bernd Jablonka⁶⁹, Mathias Uhlen⁷⁰, Rebeca Eriksen⁷¹, Josef Vogt⁷², Avirup Dutta⁷², Anna Jonsson⁷², Line Engelbrechtsen⁷², Annemette Forman⁷², Nadja Sondertoft⁷², Nathalie de Preville⁷³, Tania Baltauss⁷³, Mark Walker⁷⁴, Johann Gassenhuber⁷⁵, Maria Klintenber⁷⁶, Margit Bergstrom⁷⁶, Jorge Ferrer⁷⁷

¹⁵Bornholms Hospital, Rønne, Denmark

¹⁶Affinity Proteomics, Science for Life Laboratory, School of Engineering Sciences in Chemistry, Biotechnology and Health, KTH - Royal Institute of Technology, Solna, Sweden

¹⁷Biophysics Institute (IBF-CNR), National Research Council of Italy, Milan, Italy, and Department of Biosciences, University of Milan, Milan, Italy

¹⁸Biotech & Biomarkers Research Department, Institut de Recherches Internationales Servier, Croissy sur Seine, France

¹⁹Blood Sciences, Royal Devon and Exeter NHS Foundation Trust, Exeter, United Kingdom

²⁰Institute of Clinical and Biological Sciences, University of Exeter Medical School, Exeter, United Kingdom

²¹Boehringer Ingelheim International GmbH, Therapeutic Area CNS, Retinopathies and Emerging Areas, Ingelheim am Rhein, Germany

²²Boehringer Ingelheim International GmbH, Translational Medicine & Clinical Pharmacology, Biberach an der Riss, Germany

²³Affinity Proteomics, Science for Life Laboratory, School of Biotechnology, KTH - Royal Institute of Technology, Box 1031, SE-171 21 Solna, Sweden

- ²⁴Disease Systems Biology Program, Novo Nordisk Foundation Center for Protein Research, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark
- ²⁵Section for Bioinformatics, Department of Health Technology, Technical University of Denmark, Kgs. Lyngby, Denmark
- ²⁶Clinical Operations, Sanofi-Aventis Deutschland GmbH, Frankfurt, Germany
- ²⁷Clinical Pharmacy, Saarland University, Saarbrücken, Germany
- ²⁸Clinical Research Centre, Ninewells Hospital and Medical School, University of Dundee, Dundee, Scotland, United Kingdom
- ²⁹Clinical Research Facility, Royal Victoria Infirmary, Newcastle upon Tyne, United Kingdom
- ³⁰CNR Institute of Neuroscience, Padova, Italy
- ³¹Department of Biomedical Data Sciences, Molecular Epidemiology section, Leiden University Medical Center, Leiden, The Netherlands
- ³²Department of Cell and Chemical Biology, Leiden University Medical Center, Leiden, The Netherlands
- ³³Department of Epidemiology & Biostatistics, Amsterdam UMC- location VUmc, Amsterdam Public Health Research Institute, Amsterdam, The Netherlands
- ³⁴Department of Clinical Sciences, Diabetes & Endocrinology Unit, Lund University, Skåne University Hospital Malmö, CRC, 91-12, 205 02, Malmö, Sweden
- ³⁵Department of General Practice, Amsterdam UMC- location VUmc, Amsterdam Public Health research institute, Amsterdam, The Netherlands
- ³⁶Department of Genetic Medicine and Development, University of Geneva Medical School, Geneva, Switzerland
- ³⁷Department of Mathematical Sciences, University of Bath, Bath, United Kingdom
- ³⁸Department of Metabolism, Digestion and Reproduction, Imperial College London, London, United Kingdom
- ³⁹Université de Lille, INSERM UMR 1283, CNRS UMR 8199, Institut Pasteur de Lille, EGID, Lille, France
- ⁴⁰Diabetes Research Network, Royal Victoria Infirmary, Newcastle upon Tyne, United Kingdom
- ⁴¹Digital and Data Sciences, Sanofi-Aventis Deutschland GmbH, Frankfurt, Germany
- ⁴²Center for Biological Sequence Analysis, Department of Systems Biology, Technical University of Denmark, Kongens Lyngby, Denmark
- ⁴³Eli Lilly and Company, Indianapolis, Indiana, USA
- ⁴⁴Eli Lilly Regional Operations GmbH, Vienna, Austria
- ⁴⁵Faculty of Medical and Health Sciences, University of Copenhagen, Copenhagen
- ⁴⁶Genetic and Molecular Epidemiology Unit, Lund University Diabetes Centre, Department of Clinical Sciences, CRC, Lund University, SUS, Malmö, Sweden
- ⁴⁷Oxford Centre for Diabetes, Endocrinology and Metabolism, Radcliffe Department of Medicine, University of Oxford, Oxford, United Kingdom
- ⁴⁸Oxford NIHR Biomedical Research Centre, Oxford University Hospitals NHS Foundation Trust, John Radcliffe Hospital, Oxford, United Kingdom
- ⁴⁹Wellcome Centre for Human Genetics, University of Oxford, Oxford, United Kingdom
- ⁵⁰Institute for Epidemiology and Medical Biometry, ZIBMT, University of Ulm, Ulm, Germany
- ⁵¹Institute of Cellular Medicine, Newcastle University, Newcastle upon Tyne, United Kingdom
- ⁵²Boehringer Ingelheim International GmbH, Medicine Cardiometabolism and Respiratory, Ingelheim am Rhein, Germany
- ⁵³Institute of Epidemiology II, Research Unit of Diabetes Epidemiology, Helmholtz Zentrum München, Neuherberg, Germany
- ⁵⁴German Center for Diabetes Research (DZD), Neuherberg, Germany
- ⁵⁵Research Unit of Molecular Epidemiology, Institute of Epidemiology, Helmholtz Zentrum München, Neuherberg, Germany

- ⁵⁶Medizinische Universitätsklinik Tübingen, Eberhard Karls Universität Tübingen, Tübingen, Germany
- ⁵⁷NIHR Exeter Clinical Research Facility, University of Exeter Medical School, Exeter, United Kingdom
- ⁵⁸Internal Medicine, Institute of Clinical Medicine, University of Eastern Finland, Kuopio, Finland
- ⁵⁹Oxford Centre for Diabetes, Endocrinology and Metabolism, University of Oxford, Oxford, UK, OX3 7LJ
- ⁶⁰Inserm, Univ Lille, CHU Lille, Lille Pasteur Institute, EGID, Lille, France
- ⁶¹Population Health & Genomics, School of Medicine, University of Dundee, Dundee, United Kingdom
- ⁶²R&D Global Development, Translational Medicine & Clinical Pharmacology (TMCP), Sanofi-Aventis Deutschland GmbH, Frankfurt, Germany
- ⁶³Research Centre for Optimal Health, Department of Life Sciences, University of Westminster, London, United Kingdom
- ⁶⁴Research Unit of Molecular Endocrinology and Metabolism, Helmholtz Zentrum München, Neuherberg, Germany
- ⁶⁵Lilly Deutschland GmbH, Bad Homburg, Germany
- ⁶⁶Centre for Health, Law and Emerging Technologies (HeLEX), Faculty of Law, University of Oxford, Oxford, United Kingdom
- ⁶⁷Technologies (HeLEX), Melbourne Law School, University of Melbourne, Carlton, Victoria, Australia
- ⁶⁸Section for Nutrition Research, Faculty of Medicine, Hammersmith Campus, Imperial College London, London, United Kingdom
- ⁶⁹Strategy & Innovation, Sanofi-Aventis Deutschland GmbH, Frankfurt, Germany
- ⁷⁰Systems Biology, Science for Life Laboratory, School of Engineering Sciences in Chemistry, Biotechnology and Health, KTH - Royal Institute of Technology, Solna, Sweden
- ⁷¹Section for Nutrition Research, Division of Digestive Diseases, Department of Metabolism, Digestion and Reproduction, Faculty of Medicine, Imperial College London, UK
- ⁷²The Novo Nordisk Center for Basic Metabolic Research, Section of Metabolic Genetics, Faculty of Health and Medical Science, University of Copenhagen, Copenhagen, Denmark
- ⁷³Translational & Clinical Research, Metabolism Innovation Pole, Institut de Recherches Internationales Servier, Suresnes Cedex, France
- ⁷⁴Translational and Clinical Research Institute, Faculty of Medical Sciences, Newcastle University, Newcastle, United Kingdom
- ⁷⁵Diabetes Division, Sanofi-Aventis Deutschland GmbH, Frankfurt, Germany
- ⁷⁶VO Endokrinologi, Enheten för diabetesstudier Lasarettsgatan 15, Skånes Universitetssjukhus i Lund, Sweden
- ⁷⁷Institut d'Investigacions Biomediques August Pi i Sunye, Centre Esther Koplowitz, Barcelona, Spain