

Manifold Fitting Reveals Metabolomic Heterogeneity and Disease Associations in UK Biobank Populations

Bingjie Li^{a,1}, Jiaji Su^{a,1}, Runyu Lin^{a,1}, Shing-Tung Yau^{b,2}, and Zhigang Yao^{a,2}

This manuscript was compiled on April 18, 2025

Nuclear magnetic resonance (NMR)-based metabolic biomarkers provide comprehensive insights into human metabolism; however, extracting biologically meaningful patterns from such high-dimensional data remains a significant challenge. In this study, we propose a manifold-fitting-based framework to analyze metabolic heterogeneity within the UK Biobank population, utilizing measurements of 251 NMR biomarkers from 212,853 participants. Initially, our method clusters these biomarkers into seven distinct metabolic categories that reflect the modular organization of human metabolism. Subsequent manifold fitting to each category unveils underlying low-dimensional structures, elucidating fundamental variations from basic energy metabolism to hormone-mediated regulation. Importantly, three of these manifolds clearly stratify the population, identifying subgroups with distinct metabolic profiles and associated disease risks. These subgroups exhibit consistent links with specific diseases, including severe metabolic dysregulation and its complications, as well as cardiovascular and autoimmune conditions, highlighting the intricate relationship between metabolic states and disease susceptibility. Supported by strong correlations with demographic factors, clinical measurements, and lifestyle variables, these findings validate the biological relevance of the identified manifolds. By utilizing a geometrically informed approach to dissect metabolic heterogeneity, our framework enhances the accuracy of population stratification and deepens our understanding of metabolic health, potentially guiding personalized interventions and preventive healthcare strategies.

metabolic manifolds | geometric decomposition | manifold fitting | disease risk prediction | blood metabolomics | population heterogeneity

Nuclear magnetic resonance (NMR)-based metabolomics is transforming our understanding of human metabolic health by enabling the high-throughput, simultaneous quantification of a broad spectrum of circulating metabolites—including lipids, amino acids, and glycolysis-related compounds—at a population scale. This approach provides a holistic and cost-effective snapshot of systemic metabolism, reflecting both genetic and environmental influences (1). Such NMR-derived metabolic signatures have been shown to correlate with a wide array of clinical outcomes. For instance, specific plasma metabolite profiles have been linked to early atherosclerotic changes in subclinical cardiovascular disease, thereby supporting earlier interventions and improved risk stratification (2). In addition, metabolic biomarkers identified through plasma profiling have enhanced the prediction of diabetic complications such as diabetic retinopathy, informing targeted screening protocols and individualized patient management (3).

The UK Biobank has integrated extensive phenotypic and genetic data on over half a million participants, and ongoing initiatives have expanded these resources to include NMR-based metabolic profiling (4). With these data in hand, researchers have the opportunity to move beyond studying individual biomarkers toward a deeper understanding of metabolic heterogeneity within large, diverse populations. A more comprehensive perspective on biomarkers enables the characterization of distinct metabolic categories and the exploration of how genetic, lifestyle, and environmental factors shape metabolic profiles. By harnessing these large-scale metabolomic datasets, it is possible to uncover potential new disease mechanisms, refine disease subtyping, and improve the precision of risk prediction models. In this way, NMR metabolomics can guide clinical decision-making and public health strategies—identifying individuals at elevated risk for specific diseases long before clinical symptoms manifest, and pointing toward targeted lifestyle or therapeutic interventions.

Significance Statement

This study is the first to utilize a manifold-fitting framework within NMR-based metabolomics to explore metabolic heterogeneity in the UK Biobank population. Our method clusters 251 metabolic biomarkers into seven distinct categories that reflect the modular organization of human metabolism. Applying manifold fitting reveals low-dimensional structures in each category, capturing crucial metabolic variations associated with diverse disease risks. Notably, fitted manifolds in three categories distinctly stratify the population, each identifying two subgroups with unique metabolic profiles linked to a broad spectrum of diseases, from metabolic complications to cardiovascular and autoimmune disorders. This nuanced stratification enhances our understanding of the interactions between metabolism and disease, potentially guiding personalized health interventions and advancing preventive medicine strategies.

Author affiliations: ^aDepartment of Statistics and Data Science, National University of Singapore, Singapore 117546, Singapore; ^bYau Mathematical Sciences Center, Jinghai, Tsinghua University, Beijing 100084, China

Z.Y. designed research; B.L., J.S., R.L. and Z.Y. performed research; B.L., J.S., and R.L. analyzed data; Z.Y. contributed research methodology; S.-T.Y. supervised the project; and B.L., J.S., R.L., S.-T.Y. and Z.Y. wrote the paper.

The authors declare no competing interest.

¹B.L., J.S., and R.L. contributed equally to this work.

²To whom correspondence should be addressed. Email: styau@tsinghua.edu.cn or zhigang.yao@nus.edu.sg.

125 Emerging research leveraging NMR-derived metabolic
126 biomarkers has taken two primary analytical directions. On
127 the one hand, population-level studies have been focused on
128 building predictive and stratification models—often using
129 supervised statistical and machine learning approaches—
130 to estimate disease risks or classify individuals based on
131 known outcomes (5). On the other hand, unsupervised
132 techniques, such as clustering or factorization methods,
133 have been employed to dissect metabolic heterogeneity by
134 identifying latent subgroups with distinct biochemical profiles,
135 independent of predefined clinical labels (6). Although these
136 methods have provided valuable insights, several limitations
137 remain. Many commonly used analytic techniques rely on
138 linear assumptions or simple dimension reduction tools (e.g.,
139 principal component analysis), which may fail to capture
140 the complex, nonlinear relationships that characterize high-
141 dimensional metabolic data. Addressing these shortcomings
142 will require adopting more sophisticated models capable of
143 navigating nonlinear solution spaces, such as kernel-based
144 methods (7), deep learning architectures (8), or advanced
145 dimension reduction techniques (9–11), thereby enabling a
146 more nuanced understanding of metabolic heterogeneity and
147 its links to health and disease.

148 Despite the high-dimensional nature of NMR-based
149 metabolomic measurements—which can encompass hundreds
150 of metabolites—the underlying biochemical pathways governing
151 these metabolites are constrained by the organism’s
152 metabolic pathways and regulatory networks (12–14). As a result,
153 the observed metabolic variation likely resides on a lower-
154 dimensional manifold embedded within the high-dimensional
155 space (15, 16), making metabolomic data an ideal candidate
156 for manifold fitting approaches (17–19). Indeed, manifold
157 fitting, an invention that estimates the underlying manifold,
158 has already demonstrated its capacity to capture the intrinsic
159 geometry of complex, high-dimensional datasets, as evidenced
160 by its successful application to modeling nonlinear data
161 structures (20) and enhancing single-cell RNA sequencing
162 analysis through improved clustering and visualization (21).

163 A key advantage of manifold fitting for metabolomic
164 data analysis lies in its ability to reconstruct a smooth
165 manifold directly in the ambient measurement space, thereby
166 retaining all metabolically relevant information while filtering
167 out measurement noise. The complex, nonlinear relationships
168 between metabolites—shaped by substrate-product
169 transformations, regulatory feedback loops, and pathway
170 cross-talk—are likely to form a smooth manifold that could
171 capture the key underlying degrees of freedom in cellular
172 metabolism (22). Unlike traditional dimension reduction
173 techniques, which often rely on transformations that risk
174 losing information, modern manifold fitting methods (19–
175 21) operate directly in the original feature space. Their
176 flexible neighborhood definitions enable them to adapt to
177 diverse metabolite distributions, faithfully representing the
178 complexity and nuance of metabolic networks.

179 The aim of this study is to develop a comprehensive
180 analytical method to elucidate population-level metabolic
181 heterogeneity from multiple perspectives. Our approach
182 builds on the manifold-fitting framework, capitalizing on
183 the intrinsic properties of metabolic biomarkers. Given that
184 metabolic biomarkers naturally form distinct categories—each
185 characterized by coordinated patterns of variation reflecting

187 specific biological processes and associated disease risks—
188 we first implement an unsupervised clustering approach
189 to divide the metabolites into seven categories based on
190 their biological relevance and coordination patterns. This
191 categorization enables a modular understanding of human
192 metabolism, from basic energy metabolism to complex
193 lipoprotein regulation. Applying manifold fitting to these
194 categories reveals underlying low-dimensional structures that
195 remained consistent across the population. Remarkably,
196 we identify distinct population stratification patterns that
197 demonstrate strong associations with disease risks and health
198 outcomes. Through comprehensive analysis of these metabolically
199 defined subgroups, we uncover not only their unique
200 disease susceptibility profiles but also differential responses
201 to lifestyle factors. This integrative approach allows us to
202 bridge the gap between metabolic patterns and actionable
203 health management strategies, particularly for high-risk
204 populations. Specifically, by examining how lifestyle factors
205 modify disease risks in metabolically vulnerable subgroups,
206 we provide insights into targeted intervention strategies. To
207 our knowledge, this work represents the first application of
208 manifold fitting to large-scale metabolomic data, setting the
209 stage for more precise population stratification while enabling
210 a deeper understanding of the interplay between metabolic
211 profiles, lifestyle factors, and disease risks that can guide
212 personalized prevention efforts.

Results

213 **A brief overview of our framework.** We investigate metabolic
214 heterogeneity in a large-scale population cohort ($n = 212,853$)
215 from the UK Biobank with 251 NMR-measured metabolic
216 biomarkers and examine the associations of these biomarkers
217 with lifestyle factors and clinical outcomes. Our analytical
218 framework comprises four sequential phases: metabolic
219 biomarker clustering, manifold fitting for each biomarker
220 category, heterogeneity visualization, and characterization
221 of metabolically distinct subgroups in relation to health
222 outcomes and lifestyle factors.

223 Following data acquisition, an unsupervised clustering
224 approach is implemented for the 251 metabolic biomarkers
225 (**Fig. 1A**). This clustering strategy is based on the
226 biological premise that metabolic pathways exhibit varying
227 degrees of interconnectedness, with some pathways exhibiting
228 strong regulatory coupling while others operate independently
229 (23). Hence, identifying categories of highly interconnected
230 metabolites while maintaining inter-category independence
231 will enhance the detection of underlying manifold structures.
232 The number of metabolite categories is optimized through
233 silhouette coefficient maximization, yielding seven distinct
234 metabolic categories, which are denoted as **C1–C7** for
235 simplicity (*Materials and Methods*). The composition of these
236 categories are detailed in subsequent analyses. This modular
237 decomposition of the metabolome enables multi-dimensional
238 characterization of population heterogeneity across distinct
239 metabolic pathways, facilitating the identification of new
240 disease associations through pathway-specific analyses.

241 Subsequently, manifold fitting is performed on each
242 metabolic category to explore their underlying low-
243 dimensional structures (**Fig. 1B**), resulting in seven distinct
244 manifolds, which are denoted as **M1–M7** correspondingly
245 (*Materials and Methods*). Our objective is to identify principal
246

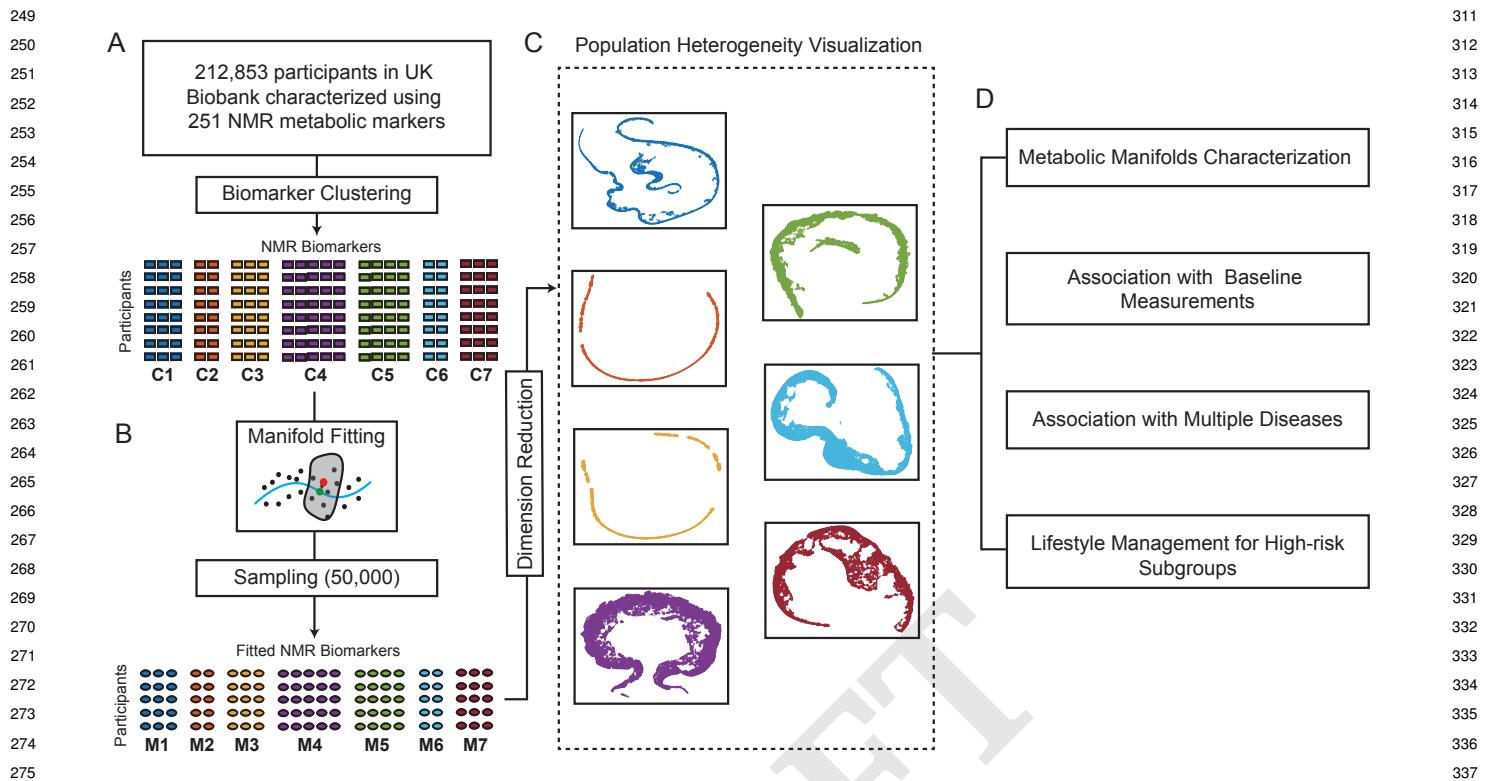


Fig. 1. Manifold-fitting-based framework for metabolic profiling and population heterogeneity in UK Biobank. (A) A total of 212,853 UK Biobank participants are characterized using 251 metabolic biomarkers. These biomarkers are divided into seven categories (**C1–C7**) based on their population-level associations. (B) Manifold fitting is applied to identify the underlying intrinsic structures. Seven manifolds (**M1–M7**) are extracted from each of the seven categories. (C) Population heterogeneity is visualized using UMAP dimension reduction. The visualization reveals distinct patterns across all seven manifolds. (D) Population health profiles are analyzed based on the heterogeneity discovered from the manifolds. These analyses encompass metabolic manifolds characterization, association with baseline measurements, association with multiple diseases, and lifestyle management for high-risk subgroups.

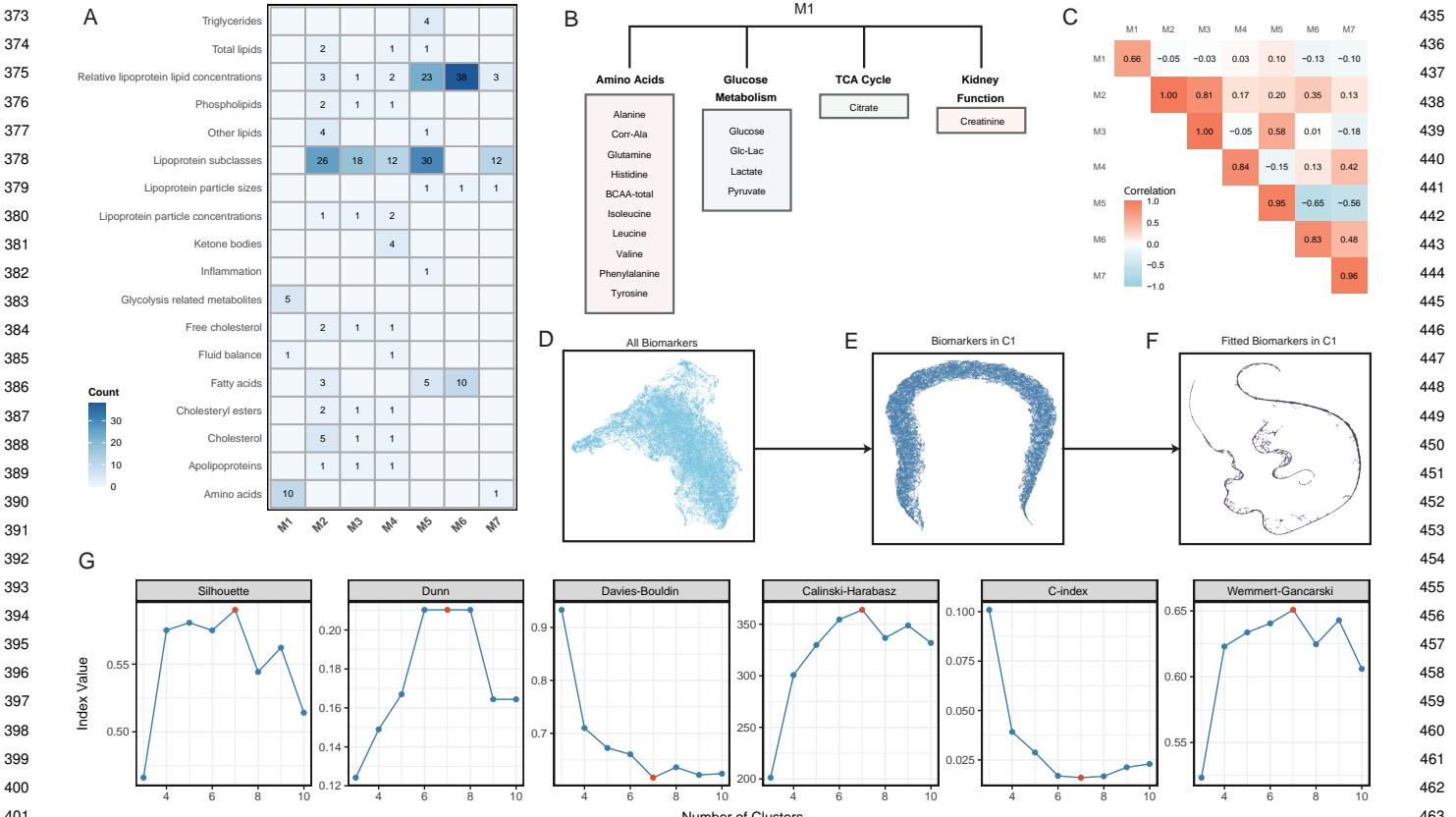
metabolic variations within the population while mitigating confounding variability that could be classified as noise. From a biological perspective, this noise stems from multiple sources: not only inherent technical measurement errors in the detection methodology and fluctuations in time-sensitive biomarkers induced by transient lifestyle factors during the measurement process (24), but also intrinsic variability from diverse complex biological processes within organisms, such as stochastic fluctuations in cellular signaling pathways, randomness in gene expression, molecular-level stochastic events, and inter-individual biological differences (25). Despite these multi-source noises, the manifold fitting procedure is able to minimize their impact on population heterogeneity characterization, leading to data with more pronounced structural features and well-defined distributions. The manifold fitting procedure reduces the impact of noise on population heterogeneity characterization, leading to more pronounced data structures and well-defined distributions.

After applying manifold fitting, we randomly select a subset of 50,000 participants for further dimension reduction and visualization using uniform manifold approximation and projection (UMAP, (11)). This process generates two-dimensional embeddings, resulting in seven distinct projections that capture population heterogeneity across various metabolic domains (Fig. 1C). Four of these projections (**M1, M2, M3, M5**) exhibit remarkable topological discontinuities, manifesting as well-defined, discrete population substructures.

Subsequent density-based clustering analysis of the reduced-dimensional representations reveals robust population stratification into binary subgroups within **M1, M2**, and **M5**. The remaining metabolic manifolds demonstrate quasi-arc trajectories, suggesting continuous phenotypic variation along metabolic axes. The projection results are more informative compared to dimension reduction result of neural network based methods (see for *SI Appendix Fig. S9*).

The population heterogeneity can be applied to a variety of downstream health analyses (Fig. 1D), such as characterization of distinct metabolic manifold structures to establish metabolic reference states, investigation of associations with baseline clinical measurements to validate biological relevance, comprehensive analysis of relationships with multiple disease states to identify potential metabolic risk factors, and development of targeted lifestyle management strategies for identified high-risk subgroups (*Materials and Methods*). This systematic analytical framework enables a multi-perspective understanding of how metabolic heterogeneity intersects with health outcomes and lifestyle factors, potentially informing personalized intervention strategies.

Characterization of metabolic manifolds. The seven distinct metabolic manifolds we have identified contain different sets of metabolites. Fig. 2A illustrates the distribution of metabolic biomarkers across these manifolds. Notably, **M6** exhibits the highest density of relative lipoprotein lipid concentrations



403 **Fig. 2.** Characterization of metabolic manifolds and their low-dimensional structure. (A) Heatmap showing the distribution of metabolic biomarkers across seven manifolds (M1–M7). (B) Detailed composition of M1, illustrating four distinct metabolic modules: amino acids, glucose metabolism, TCA cycle, and kidney function markers. (C) Mean Pearson correlation coefficients between biomarkers across different manifolds, demonstrating strong intra-manifold correlations and weak inter-manifold correlations. (D) Initial high-dimensional structure of all 251 biomarkers without clear pattern. (E) Emergence of low-dimensional characteristics in the first metabolite category (C1) after biomarker clustering. (F) Final manifold structure revealing two distinct subgroups after manifold fitting. (G) The optimal number of clusters identified by six clustering validation metrics. Higher values indicate better clustering for Silhouette score, Dunn index, Calinski-Harabasz index, and Wermert-Gancarski index, while lower values are optimal for Davies-Bouldin index and C-index. All metrics consistently indicate seven clusters as optimal (red dots).

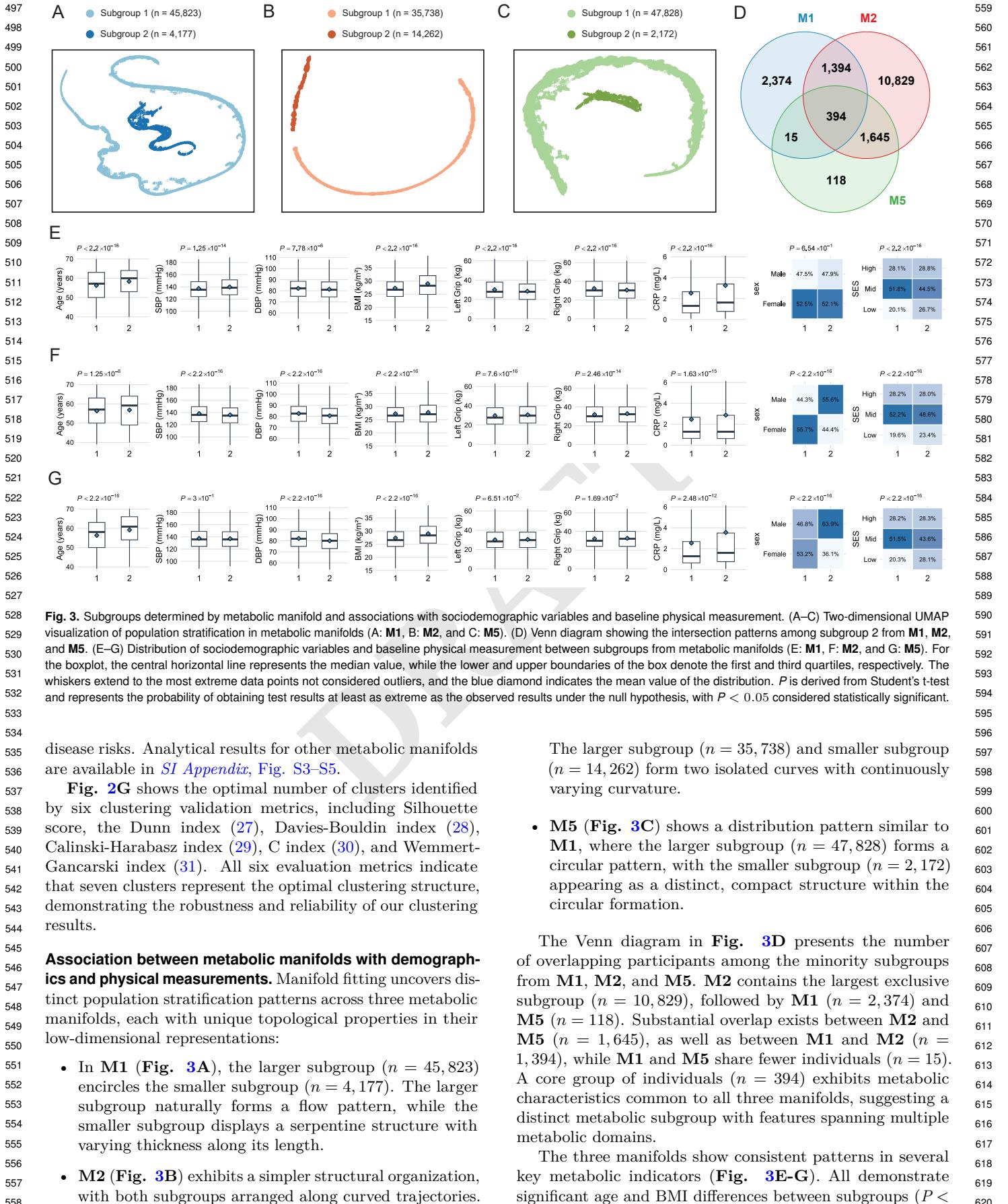
410 (contain 38 metabolic biomarkers, $n_b = 38$), while M1 shows significant enrichment in amino acids ($n_b = 10$) and glycolysis-related metabolites ($n_b = 5$). Lipoprotein subclasses demonstrate heterogeneous distribution patterns, with substantial representation in M2 ($n_b = 26$), M3 ($n_b = 18$), M5 ($n_b = 30$), and M7 ($n_b = 12$). Here, we use M1 as an example to demonstrate the biological rationale behind our metabolic marker categorization. The comprehensive description of metabolite compositions for all other manifolds can be found in *SI Appendix, Table S2–S3*.

411 **Fig. 2B** delineates the biological composition of M1, which comprises four distinct metabolic modules: amino acids, glucose metabolism, tricarboxylic acid (TCA) cycle, and kidney function biomarkers. The amino acid module encompasses ten key metabolites, including branched-chain amino acids, aromatic amino acids, and other essential and non-essential amino acids. The glucose metabolism module contains fundamental glycolytic intermediates and glucose-lactate ratio. The TCA cycle is represented by citrate, while kidney function is monitored through creatinine levels. M1 captures a fundamental metabolic network centered on cellular energy metabolism and protein homeostasis, which reveals the intricate interplay between amino acid metabolism, glucose utilization, and energy production through the

412 TCA cycle. The clustering of these metabolites suggests coordinated regulation of energy substrate utilization (26).

413 The biological validity of our biomarker clustering can be further substantiated through correlation analysis. **Fig. 2C** presents the mean Pearson correlation coefficients between biomarkers assigned to different manifolds. Global analysis reveals strong intra-manifold correlations while maintaining relatively weak inter-manifold correlations, supporting the modular organization of the metabolome. For instance, biomarkers within M1 exhibit a high average correlation coefficient (0.66), whereas correlations between M1 and other manifolds remain consistently low (less than 0.13).

414 The manifold fitting procedure substantially enhances the 415 interpretability of the data. Initially, the 251 biomarkers 416 do not show any obvious low-dimensional structure, as 417 shown in the UMAP visualization **Fig. 2D**. However, after 418 clustering the biomarkers, the first metabolite category (C1) 419 begins to exhibit emerging low-dimensional characteristics, 420 as depicted in the UMAP plot **Fig. 2E**, despite residual 421 variance in other directions. Following manifold fitting, the 422 data show pronounced directionality and segregate into two 423 distinct, disconnected subgroups, as illustrated in **Fig. 2F**. 424 Subsequent analyses have demonstrated that these subgroups 425 exhibit significant differences in metabolic profiles, and 426



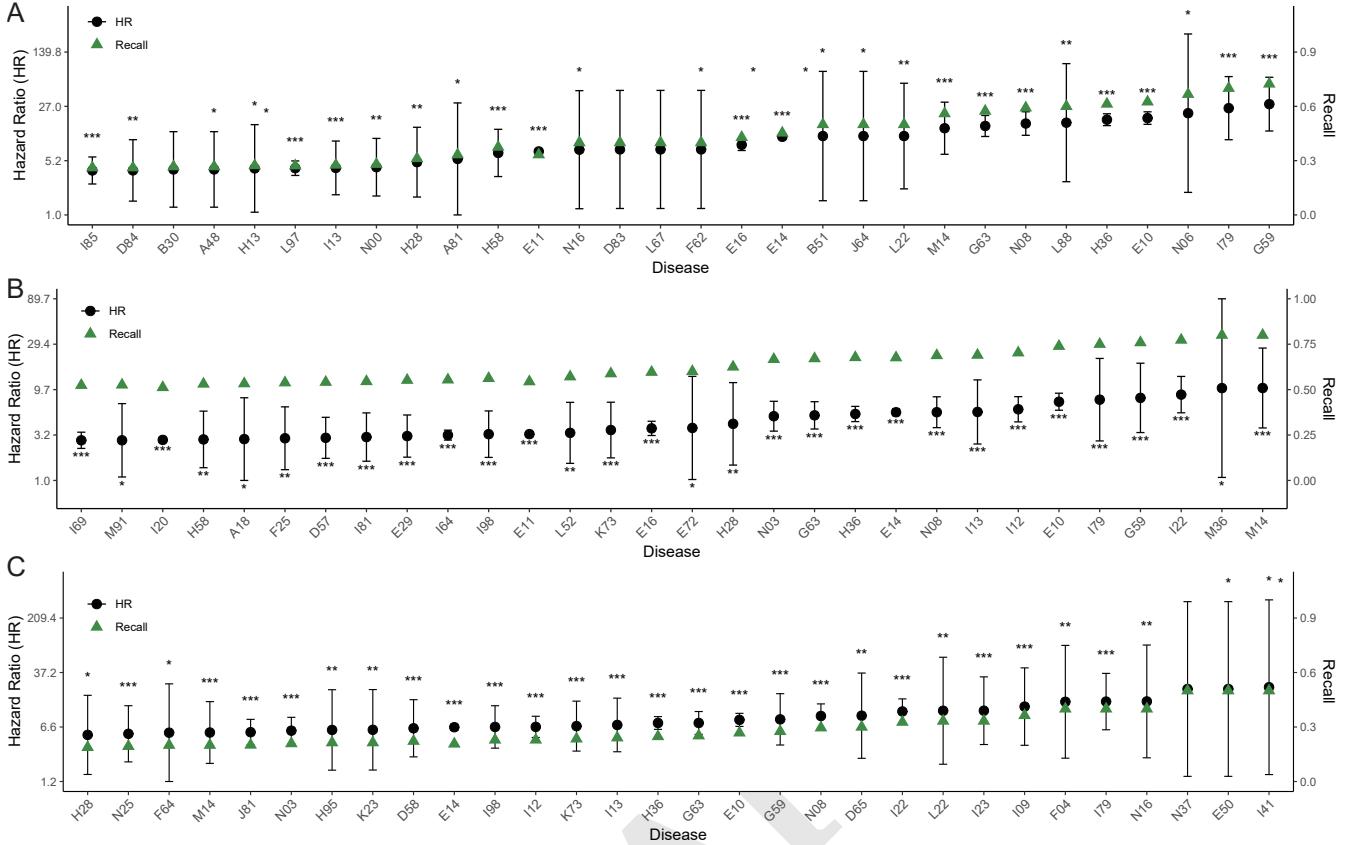


Fig. 4. Comparison of the hazard ratios (HRs) of diseases and their corresponding recall in three identified high-risk subgroups relative to their respective reference subgroups from metabolic manifolds (A: M1, B: M2, and C: M5). HRs represent the risk of disease incidence from baseline assessment (2007–2010) until right censoring. Each panel displays diseases sorted by descending HRs, with black circles and error bars representing the estimation of HRs and their 95% confidence intervals (CIs), respectively, and green triangles indicating the recall (percentage of all disease cases captured) of the high-risk subgroup among all patients with each disease. The number of asterisks (*) represents the level of statistical significance: one asterisk (*) indicates $P < 0.05$, two asterisks (**) indicate $P < 0.01$, and three asterisks (****) indicate $P < 0.001$.

2×10^{-16}), with Subgroup 2 consistently showing higher values. Blood pressure and grip strength measurements show systematic variations across manifolds, though with differing magnitudes. CRP levels are consistently elevated in Subgroup 2, suggesting a common inflammatory component. All manifolds display clear socioeconomic status (SES) stratification, indicating a robust link between metabolic profiles and social determinants of health.

Each manifold captures distinct aspects of metabolic heterogeneity. **M1** shows balanced gender distribution (47.5% vs 47.6% male) despite clear clinical differences (BMI: 26.64 vs 28.15 kg/m²; CRP: 1.28 vs 1.63 mg/L). In contrast, **M2** exhibits the most pronounced gender disparity (46.8% vs 63.9% male) but more modest clinical differences (BMI: 26.65 vs 27.04 kg/m²). **M5** has the strongest coupling between clinical and demographic differences, with the largest age gap (58 vs 61 years) and substantial BMI differences (26.68 vs 28.35 kg/m²), alongside clear demographic stratification.

Association between metabolic manifolds with multiple diseases. Among the 923 diseases analyzed through Cox proportional hazards regression models, the three metabolic manifolds demonstrate significant positive associations ($HR > 1$, $P < 0.05$) with 269 (29.1%), 287 (31.1%), and 298 (32.3%) diseases in **M1**, **M2**, and **M5**, respectively. For each manifold,

we conducted separate survival analyses using membership in the high-risk cluster as the predictor variable and time-to-disease onset as the outcome. Protective effects ($HR < 1$, $P < 0.05$) are relatively rare, observed in only 7 (0.8%), 23 (2.5%), and 18 (1.9%) diseases. Despite representing small subgroups (8.4%, 28.5%, and 4.3% of the total population), these manifolds show high capture rates for severe conditions, with **M2** identifying the most diseases with high recall rates (81 diseases with recall > 50%) compared to **M1** (15 diseases) and **M5** (6 diseases), suggesting its broader coverage of disease spectrum. Detailed results including all hazard ratios, confidence intervals, and P-values are available in *SI Appendix, Table S4–S6*.

M1, **M2**, and **M5**, each representing distinct population stratification patterns in their low-dimensional representations, demonstrate significant associations with both type 1 (E10; HR: 7.18–18.93) and type 2 diabetes (E11; HR: 3.24–6.89), accompanied by their microvascular complications including retinopathy (H36; HR: 5.31–18.02), nephropathy (N08; HR: 5.55–16.04), and neuropathy (G63; HR: 5.14–28.87). Here, hazard ratios (HRs) represent the risk of disease incidence from baseline assessment (2007–2010) until right censoring at the earliest of: date of disease diagnosis, date of death, date of loss to follow-up, or the last date of data collection. These associations are particularly robust in

745 conditions with large sample sizes, such as type 2 diabetes
746 (E11; $n = 4,519$) and retinal disorders (H36; $n = 504$),
747 as evidenced by their narrow CIs and highly significant P
748 ($P < 0.001$). Correspondence and information between ICD-
749 10 codes and disease names is available in *SI Appendix, Table*
750 **S1**.

751 Despite these commonalities, each manifold captures
752 distinct disease spectrums and risk gradients. **M1** (**Fig.**
753 **4A**), with its highest hazard ratios ($HR > 15$) for multiple
754 conditions, predominantly highlights severe metabolic dys-
755 regulation and its complications, particularly in neurological
756 (mononeuropathy (G59): $HR = 28.87$, 95% CI: 12.79–65.17,
757 recall = 72.4%) and vascular disorders (arterial diseases (I79):
758 $HR = 25.56$, 95% CI: 9.82–66.51, recall = 70.0%). In contrast,
759 **M2** (**Fig. 4B**) exhibits a more moderate but broader risk
760 profile (most $HR < 6$), distinctively emphasizing cardiovascular
761 and autoimmune conditions, with notably elevated risks
762 for arthropathies (M14; $HR = 10.05$, 95% CI: 3.77–26.77,
763 recall = 80.0%) and subsequent myocardial infarction (I22;
764 $HR = 8.55$, 95% CI: 5.47–13.34, recall = 77.3%). **M5** (**Fig.**
765 **4C**) demonstrates an intermediate pattern with a clear risk
766 stratification, showing significant but generally lower hazard
767 ratios compared to **M1**, while maintaining a comprehensive
768 coverage of metabolic complications. Notably, the high recall
769 rates (ranging from 60–80% for severe conditions) across
770 all manifolds indicate their effectiveness in identifying the
771 majority of high-risk patients despite representing relatively
772 small subgroups. Even though the remaining five manifolds
773 lack two-subgroup structures, they still exhibit population-
774 specific patterns in disease onset risk, as demonstrated in *SI*
775 **Appendix, Fig. S1–S2**.

776 These distinct risk patterns have important clinical impli-
777 cations for personalized prevention strategies. **M1** identifies
778 a subgroup requiring aggressive complication prevention and
779 intensive monitoring, particularly for neurological and vas-
780 cular complications. **M2** suggests a need for comprehensive
781 cardiovascular protection and autoimmune surveillance in its
782 high-risk subgroup. **M5** indicates a requirement for stratified
783 prevention strategies across multiple systems, with particular
784 attention to the progression of metabolic syndrome. These
785 findings highlight the potential of metabolic manifolds as
786 complementary biomarkers for patient stratification, suggest-
787 ing that their combined use might enable more precise risk
788 prediction and personalized intervention strategies.

789 **Associations between lifestyle and disease in high-risk**
790 **metabolic subgroups.** To translate our metabolic insights
791 into actionable prevention strategies, we investigate the
792 relationship between lifestyle factors and disease outcomes
793 in high-risk populations identified through manifold analysis.
794 Combining the high-risk subgroups of **M1**, **M2**, and **M5**, we
795 have identified a comprehensive population that demonstrates
796 increased susceptibility across multiple disease domains.

797 Our analysis focuses on three key lifestyle factors, including
798 sleep patterns, physical activity and smoking status, and their
799 associations with four major diseases: diabetes mellitus (E11,
800 **Fig. 5A**), chronic ischemic heart disease (I25, **Fig. 5B**),
801 chronic obstructive pulmonary disease (J44, **Fig. 5C**), and
802 chronic kidney disease (N18, **Fig. 5D**). The results reveal
803 consistent patterns of lifestyle-associated risk across these
804 conditions, suggesting potential targets for intervention in
805 metabolically vulnerable populations.

806 Sleep patterns emerge as a significant modifier of disease
807 risk. Individuals with unhealthy sleep patterns show con-
808 sistently higher disease incidence rates compared to those
809 with healthy sleep patterns, with particularly pronounced
810 differences in diabetes (14.48% vs 10.77%) and chronic kidney
811 disease (9.36% vs 7.10%). This pattern suggests that sleep
812 quality may be a crucial mediator of metabolic health,
813 potentially through its effects on energy homeostasis and
814 inflammatory pathways.

815 Physical activity levels demonstrate similarly striking
816 associations with disease outcomes. The physically inactive
817 group shows elevated risk across all studied conditions, with
818 the most notable differences observed in diabetes (14.09%
819 vs 8.90%) and ischemic heart disease (11.53% vs 9.72%).
820 These findings underscore the importance of regular physical
821 activity as a protective factor against metabolic dysfunction,
822 even in populations with underlying metabolic vulnerability.

823 Smoking status shows the most dramatic impact on disease
824 risk, particularly for respiratory conditions. The contrast
825 is most striking for chronic obstructive pulmonary disease,
826 where smokers show more than four-fold higher incidence
827 rates compared to non-smokers (9.01% vs 2.03%). However,
828 the impact of smoking extends beyond respiratory health,
829 with elevated risks observed across all studied conditions,
830 including diabetes (14.21% vs 9.65%) and chronic kidney
831 disease (9.48% vs 6.14%).

Materials and Methods

832 **UK Biobank cohort and data access.** The UK Biobank is a compre-
833 hensive biomedical database (*SI Appendix, A*) that provides global
834 access to data from approximately half a million participants
835 aged between 40 and 69 years at baseline (32). This resource
836 encompasses a wide array of health-related information, including
837 questionnaire data on socio-economic status, lifestyle factors, and
838 cognitive assessments, as well as measurements of heart and lung
839 function, body size, and composition. Additionally, a variety of
840 biochemical and imaging data are available.

841 The UK Biobank blood samples are collected at baseline in 22
842 assessment centres across the UK from 2007 to 2010. Protocols
843 for handling and storing these samples are detailed in (33).
844 Nightingale Health Plc. engages in the biomarker profiling of
845 baseline plasma samples for the entire cohort. This profiling
846 employs the Nightingale Health NMR biomarker platform (*SI*
847 **Appendix, B**), the specifics of which are well-documented in (34, 35).
848 The main procedural steps in the biomarker analysis and quality
849 controlling are outlined in (4). This NMR metabolic biomarker
850 dataset is made available to the research community via the UK
851 Biobank as of March 2021. For this study, data are requested under
852 UK Biobank project 146760 in March 2024. After the exclusion
853 of incomplete data, a total of 251 NMR metabolic biomarkers
854 from 212,853 participants are utilized. As the primary target of
855 analysis, disease outcome is defined based on the first occurrence
856 of 3-character ICD-10 codes using the hospital inpatient records
857 from the UK Biobank. Additional lifestyle-related indicators, such
858 as sleep and exercise levels, are also included in the analysis.

859 **Refining data intrinsic structures.** The proposed methodological
860 framework employs all 251 NMR metabolic biomarkers to improve
861 our understanding of their relationships and properties in a
862 biological context, avoiding the selection of a subset. To cope
863 with the high dimension and complex covariance structure among
864 the biomarkers, we first partition them into several low-dimensional
865 subspaces. Subsequently, we apply manifold fitting to each
866 subspace to elucidate their intrinsic low-dimensional structures,
867 facilitating further analysis. This process is illustrated in **Fig.**
868 **1A–C**.

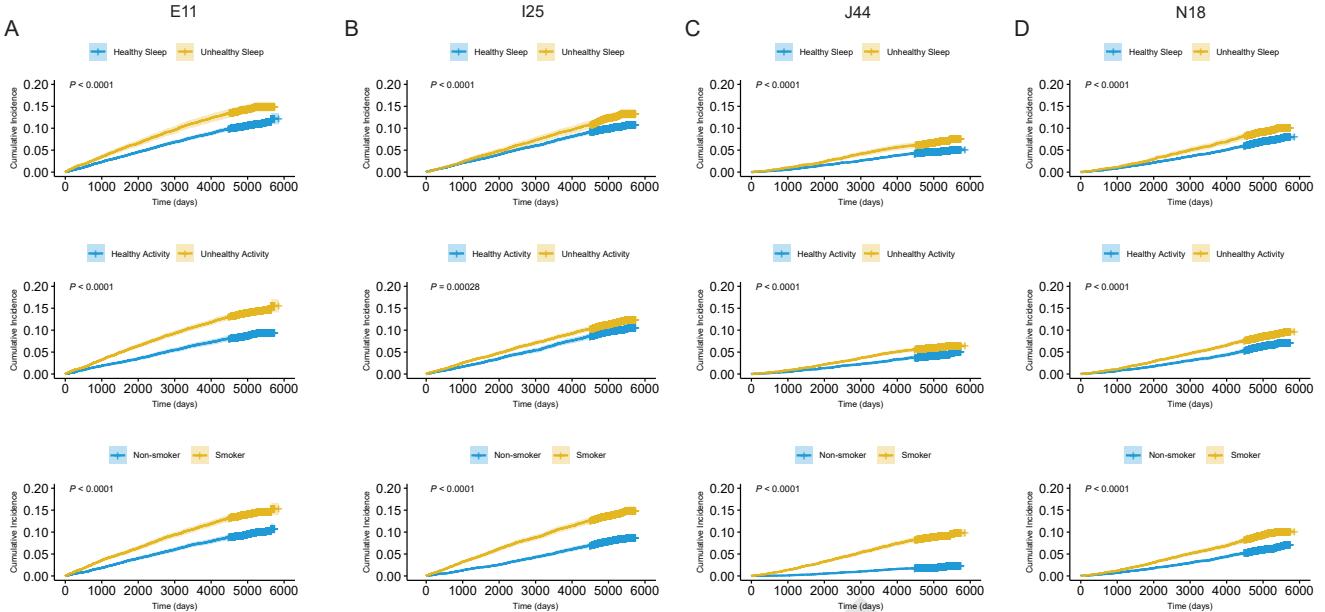


Fig. 5. Kaplan-Meier survival curves showing cumulative disease incidence stratified by sleep patterns, physical activity, and smoking status in metabolically vulnerable populations. Disease-specific incidence rates for (A) diabetes mellitus (E11), (B) chronic ischemic heart disease (I25), (C) chronic obstructive pulmonary disease (J44), and (D) chronic kidney disease (N18). Cross-mark represent right censoring marks indicating participants who were lost to follow-up, died from other causes, or reached the end of the study period without experiencing the event of interest.

Partitioning metabolic biomarkers into seven categories. The 251 metabolic biomarkers under study are involved in diverse metabolic pathways, with some directly measured and others calculated from these direct measurements. This variability contributes to a complex covariance structure, necessitating the partitioning of the biomarker space into multiple subspaces, or, in other words, cluster the biomarkers into multiple categories. Each biomarker is initially treated as a 212,853-dimensional vector, with pairwise distances calculated using a correlation metric. Dimension reduction and two-dimensional UMAP visualization reveal nonlinear dependencies and clustering features among the biomarkers (SI Appendix, Fig. S6). Following this, we recalculate a Manhattan distance matrix from the two-dimensional UMAP projections for hierarchical clustering using the average linkage method. We validate these clusters by calculating silhouette scores for cluster sizes ranging from 2 to 50, identifying the optimal size based on the highest average silhouette score. This process delineates seven distinct biomarker categories (**C1–C7**), each characterized by a low-dimensional nonlinear correlation structure or modeled around a latent manifold, exhibiting minimal inter-group correlations.

Fitting latent manifold in each subspace. The covariance matrices and two-dimensional visualizations for most metabolic biomarker categories exhibit strong nonlinear correlations within their respective subspaces. These dependencies imply that the relationships among biomarkers can be effectively refined and more accurately described using manifold fitting. This method leverages the structural features of the samples to enhance characterizations of principal variations by aligning the data towards a low-dimensional latent manifold. This alignment serves to filter out noise and irrelevant information, ensuring that the refined data remains consistent within the original space. The effectiveness of manifold fitting has been validated across various fields, and for more technical and theoretical details, refer to (19).

Take category **C1** as an instance, we represent the data in **C1** as an N by D_1 matrix $(x_{ij})_{N \times D_1}$, where N is the sample size, D_1 stands for the number of biomarkers in **C1**, and x_{ij} denotes the value of the j th biomarker in **C1** for the i th investigated individual. For each individual, let $x_i = (x_{i1}, \dots, x_{iD_1})^\top$ be a D_1 -dimensional vector in \mathbb{R}^{D_1} . The manifold fitting algorithm includes two primary steps: direction estimation and contraction estimation. Consider a sample point $z \in \{x_i : i = 1, \dots, N\}$, the direction from z to the

latent manifold is estimated using a reference point defined as:

$$\mu_z = \frac{\sum_{i=1}^N x_i \mathbf{1}_{\|x_i - z\| \leq r_1}}{\sum_{i=1}^N \mathbf{1}_{\|x_i - z\| \leq r_1}},$$

where $\|\cdot\|$ stands for the Euclidean norm and r_1 is a neighborhood radius. We then define a hyper-cylinder elongated along the vector $\mu_z - z$ by construct a projection matrix onto $\mu_z - z$ as:

$$\Pi_z = \frac{(\mu_z - z)(\mu_z - z)^\top}{\|\mu_z - z\|^2},$$

with which we decompose the vector $x_i - z$, for each x_i , into two directions:

$$u_i = \Pi_z(x_i - z), \quad v_i = x_i - z - u_i,$$

yielding the fitted version of z :

$$z' = \frac{\sum_{i=1}^N x_i \mathbf{1}_{\|u_i\| \leq r_2, \|v_i\| \leq r_1}}{\sum_{i=1}^N \mathbf{1}_{\|u_i\| \leq r_2, \|v_i\| \leq r_1}}, \quad [1]$$

where $r_2 \gg r_1$ is another radius along the direction of $\mu_z - z$. As proposed by (19), with appropriately chosen parameters r_1 and r_2 , z' is significantly closer to the latent manifold compared to z . This approach allows us to filter out less important information and improve subsequent analyses without prior knowledge of the dimension of the latent manifold.

Furthermore, two-dimensional UMAP visualizations of some categories of biomarkers represent uniformly distributed within a circular area, possibly due to the large variability in biomarker distribution ranges, possibly skewed by dominant components. To reduce these effects, we implement a rank-based transformation. Continuing to use **C1** as an example, the empirical cumulative distribution function for the j th biomarker is defined as:

$$\hat{F}_j(t) = \frac{1}{N} \sum_{i=1}^N \mathbf{1}_{x_{ij} \leq t},$$

where $\mathbf{1}_{x_{ij} \leq t}$ is the indicator function of $x_{ij} \leq t$. Subsequently, the raw observation x_{ij} is transformed to $y_{ij} \in [0, 1]$ by taking

$$y_{ij} = \hat{F}_j(x_{ij}), \quad \text{for } i = 1, \dots, N, \quad j = 1, \dots, D_1. \quad [2]$$

The visualization of transformed data reveals stronger nonlinear dependencies in biomarker groups that are not characterized by the initial two-dimensional UMAP visualizations of the raw data. The manifold fitting is then performed on the dataset $\{y_{ij}\}$ and the resulting $\{y'_{ij}\}$ are transformed back to the measurement space using the quantile function

$$x'_{ij} = \widehat{F}_j^{-1}(y'_{ij}) =: \inf\{t \in \mathbb{R} : \widehat{F}_j(t) \geq y'_{ij}\}. \quad [3]$$

By applying the manifold fitting and transformation procedure to categories **C1–C7**, we generate the refined datasets on metabolic manifolds **M1–M7**, where the intrinsic geometric structures of the data are effectively refined. This entire process is encapsulated in the pseudo-code provided in Algorithm 1. Based on the initial visualizations, we determine that categories **C1** and **C4** do not require transformation. According to the recommendation from (19), the radii are set to be $r_1 = 5\sigma/\log_{10}(N)$ and $r_2 = 10\sigma\sqrt{\log_e(1/\sigma)/\log_{10}(N)}$, with a parameter σ . For the two categories without transformation, we fix σ as 1.5, while for the remaining transformed categories, σ is adjusted to $0.8\sqrt{D_j/D_1}$ for the j th category, where D_j denotes the number of biomarkers in that category.

Algorithm 1 Manifold fitting for each biomarker category

Input: Data set of interest (one of **C1–C7**): $\{x_i\}_{i=1}^N$; neighborhood radii r_1 and r_2 .

- 1: Draw the two-dimensional visualization of $\{x_i\}_{i=1}^N$;
- 2: Determine if transformation is needed
- 3: **if** so **then**
- 4: Calculate $z_{ij} \leftarrow \widehat{F}_j(x_{ij})$ with Eq. (2);
- 5: **else**
- 6: Copy $z_{ij} \leftarrow x_{ij}$;
- 7: Update z'_i from z_i with Eq. (1), for $i = 1, \dots, N$;
- 8: **if** so **then**
- 9: Calculate $x'_{ij} \leftarrow \widehat{F}_j^{-1}(z'_{ij})$ with Eq. (3);
- 10: **else**
- 11: Copy $x'_{ij} \leftarrow z'_{ij}$;

Return: Fitted set $\{x'_i\}_{i=1}^N$ (the corresponding one of **M1–M7**).

Heterogeneity analysis and phenotype associations. Following manifold fitting, we project each manifold (**M1–M7**) into two dimensions using UMAP. Due to the computational complexity, we randomly sample 50,000 participants to generate UMAP visualizations and conduct subsequent analyses in the main text; UMAP projections for the complete population are provided in *SI Appendix Fig. S7–S8*. Four of these projections (**M1**, **M2**, **M3**, **M5**) exhibit marked topological discontinuities, manifesting as well-defined, discrete population substructures. We then apply density-based spatial clustering of applications with noise (DBSCAN, (36)) to these reduced-dimensional representations, which reveals robust stratification into binary subgroups within manifolds **M1**, **M2**, and **M5**. For these three categories, we examine associations with demographic and clinical variables, including sex, age, body mass index (BMI), systolic and diastolic blood pressure, bilateral hand grip strength, and sleep quality metrics. For continuous variables, we assess between-group differences using two-tailed t-tests and visualize distributions through boxplots; for categorical variables, we compare frequency distributions using chi-squared tests.

Following population stratification, we conduct survival analysis on diseases recorded in the UK Biobank using Cox proportional hazards model. For each disease, we calculate the follow-up time as the interval between the baseline assessment date and either the date of the first disease diagnosis or the end of the follow-up period for censored cases. Disease status is determined using the International Classification of Diseases 10th revision (ICD-10) codes from UK Biobank health records. We construct Cox models incorporating subgroup assignments as predictors, with

subgroup 2 designated as the high-risk group. For each disease, we compute hazard ratios with 95% CIs and corresponding P . To assess the practical utility of our stratification approach, we calculate the recall rate for each disease, which is defined as the proportion of total disease cases captured within subgroup 2. We specifically focus on diseases showing both statistical significance ($p < 0.05$) and clinical relevance ($HR > 1$), ranking them by HR magnitude to identify conditions most strongly associated with metabolic vulnerability. This comprehensive analysis enables us to identify both the statistical strength (through HRs and P) and practical significance (through recall rates) of the metabolic subgroup stratification in disease risk prediction.

For the high-risk subgroups identified through manifold fitting and clustering, we further investigate the relationship between lifestyle factors (*SI Appendix, D*) and disease outcomes (*SI Appendix, C*) through survival analysis. We first combine the high-risk subgroups from manifolds **M1**, **M2**, and **M5** to create a comprehensive high-risk subgroup. For each disease of interest, we exclude individuals who had the disease at baseline and calculate the follow-up time from the baseline assessment date to either the date of disease diagnosis or the end of follow-up period.

We stratify this high-risk population based on three lifestyle factors: sleep patterns (healthy/unhealthy sleep), physical activity levels (healthy/unhealthy activity), and smoking status (smoker/non-smoker). For each lifestyle factor, we generate Kaplan-Meier curves to visualize cumulative disease incidence over time, with differences between groups assessed using log-rank tests. We also calculate the crude incidence rates for each lifestyle category by dividing the number of incident cases by the total number of individuals in that category. This analysis enables us to quantify the potential impact of modifiable lifestyle factors on disease risk within metabolically vulnerable populations.

Discussion

Our study presents a new analytical framework for understanding metabolic heterogeneity in large populations through manifold fitting, offering several key advances in both methodological approaches and biological insights. By applying this framework to the UK Biobank cohort, we have demonstrated how complex metabolic relationships can be effectively captured and interpreted through low-dimensional manifolds, revealing distinct population substructures with meaningful clinical correlates.

Our approach offers several advantages over traditional metabolomic analysis methods. First, by operating directly in the original feature space, manifold fitting preserves the interpretability of metabolic measurements while effectively reducing noise and capturing underlying biological structure. Second, the identification of discrete subgroups within continuous metabolic variations provides a natural framework for patient stratification that could inform precision medicine approaches. Third, the strong associations between manifold-based subgroups and disease outcomes suggest that these metabolic patterns may serve as early indicators of disease risk, potentially enabling more targeted prevention strategies.

The clinical implications of our findings are particularly relevant for precision medicine. The distinct disease associations observed across different manifolds suggest that metabolic risk factors may cluster in patterns that are not captured by traditional clinical measurements. For instance, the identification of subgroups with elevated risk for specific disease types (metabolic complications, cardiovascular diseases, or autoimmune conditions) suggests the potential for more nuanced approaches to patient risk stratification and preventive care. Although some subgroups show substantial

overlap across manifolds, closer inspection reveals highly asymmetric population compositions, with certain manifolds capturing more specialized high-risk subpopulations. This hierarchical organization highlights how different models may uncover varying levels of metabolic vulnerability, from broader risk profiles to more severely affected individuals with amplified disease burdens.

Looking forward, our framework opens several promising avenues for future research. First, genetic analysis of the identified metabolic subgroups could provide crucial insights into the hereditary components of metabolic heterogeneity. By conducting genome-wide association studies within each manifold-defined subgroup, we could identify genetic variants that contribute to specific metabolic patterns (37). This could be particularly informative for understanding the genetic architecture of complex metabolic traits and their relationship to disease risk (38).

Further exploration could focus on the longitudinal stability of these metabolic manifolds and their potential as predictive biomarkers. Time-series analysis of metabolic profiles could reveal how individuals transition between different metabolic states and whether these transitions correlate with disease onset or progression. This could lead to the development of early warning systems for metabolic dysfunction and more precise timing of preventive interventions (39).

Data, Materials, and Software Availability. The Nightingale Health NMR biomarker data has been publicly available to the UK Biobank resource since Spring 2021 and can be accessed through the UK Biobank portal (<https://biobank.ndph.ox.ac.uk/showcase/label.cgi?id=220>). Approved researchers may access these data in accordance with the UK Biobank data-access protocol. Detailed information on the data access process is available at <https://www.ukbiobank.ac.uk/enable-your-research/apply-for-access>.

The code of our analysis method, implemented in R and MATLAB, includes a demonstration of the pipeline, all necessary intermediate results, their corresponding final results, and all evaluation functions used in this study. This implementation is available at <https://github.com/zhihang-yao/MF-Metabolomic-Heterogeneity> (40). Additional results are provided in the SI Appendix.

ACKNOWLEDGMENTS. Z.Y. has been supported by the Singapore Ministry of Education Tier 2 grant (A-0008520-00-00 and A-8001562-00-00) and the Tier 1 grant (A8000987-00-00 and A-8002931-00-00) at the National University of Singapore; B.L. and J.S. are postdoctoral researchers supported by grant A-8001562-00-00. R.L. is a doctoral student supported by a Research Scholarship at the National University of Singapore.

1. P Soininen, et al., High-throughput serum nmr metabolomics for cost-effective holistic studies on systemic metabolism. *Analyst* **134**, 1781–1785 (2009).
2. PA Chevill, et al., Plasma metabolomic profiling in subclinical atherosclerosis: the diabetes heart study. *Cardiovasc. Diabetol.* **20**, 1–12 (2021).
3. S Yang, et al., Plasma metabolomics identifies key metabolites and improves prediction of diabetic retinopathy: Development and validation across multinational cohorts. *Ophthalmology* **131**, 1436–1446 (2024).
4. H Julkunen, et al., Atlas of plasma nmr biomarkers for health and disease in 118,461 individuals from the UK biobank. *Nat. Commun.* **14**, 604 (2023).
5. T Buergerl, et al., Metabolomic profiles predict individual multidisease outcomes. *Nat. Medicine* **28**, 2309–2320 (2022).
6. W Zhang, et al., Classification of osteoarthritis phenotypes by metabolomics analysis. *BMJ open* **4**, e006286 (2014).
7. C Brouard, et al., Fast metabolite identification with input output kernel regression. *Bioinformatics* **32**, i28–i36 (2016).
8. P Sen, et al., Deep learning meets metabolomics: a methodological perspective. *Briefings Bioinforma.* **22**, 1531–1542 (2021).
9. M Belkin, P Niyogi, Using manifold structure for partially labeled classification in *Advances in Neural Information Processing Systems*, eds. S Becker, S Thrun, K Obermayer. (MIT Press), Vol. 15, pp. 1–8 (2002).
10. L Van der Maaten, G Hinton, Visualizing data using t-SNE. *J. Mach. Learn. Res.* **9** (2008).
11. L McInnes, J Healy, J Melville, UMAP: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426* (2018) (Accessed on 20 May 2024).
12. ML Mo, BØ Palsson, Understanding human metabolic physiology: a genome-to-systems approach. *Trends Biotechnol.* **27**, 37–44 (2009).
13. K Suhré, C Gieger, Genetic variation in metabolic phenotypes: study designs and applications. *Nat. Rev. Genet.* **13**, 759–769 (2012).
14. JK Nicholson, ID Wilson, Understanding 'global' systems biology: Metabonomics and the continuum of metabolism. *Nat. Rev. Drug Discov.* **2**, 668–676 (2003).
15. MD Ritchie, et al., Multifactor-dimensionality reduction reveals high-order interactions among estrogen-metabolism genes in sporadic breast cancer. *The Am. J. Hum. Genet.* **69**, 138–147 (2001).
16. KH Liland, Multivariate methods in metabolomics—from pre-processing to dimension reduction and statistical analysis. *TrAC Trends Anal. Chem.* **30**, 827–841 (2011).
17. C Fefferman, S Ivanov, Y Kurylev, M Lassas, H Narayanan, Fitting a putative manifold to noisy data in *Conference on Learning Theory*, eds. B Sébastien, P Vianney, R Philippe. (PMLR), pp. 688–720 (2018).
18. C Fefferman, S Ivanov, M Lassas, H Narayanan, Fitting a manifold of large reach to noisy data. *J. Topol. Analysis* **0**, 1–82 (2023).
19. Z Yao, J Su, B Li, S-T Yau, Manifold fitting. *arXiv preprint arXiv:2304.07680* (2023) (Accessed on 10 June 2024).
20. Z Yao, J Su, S-T Yau, Manifold fitting with cyclegan. *Proc. Natl. Acad. Sci.* **121**, e2311436121 (2024).
21. Z Yao, B Li, Y Lu, S-T Yau, Single-cell analysis via manifold fitting: A framework for rna clustering and beyond. *Proc. Natl. Acad. Sci.* **121**, e2400002121 (2024).
22. R Steuer, et al., From structure to dynamics of metabolic pathways: application to the plant mitochondrial tca cycle. *Bioinformatics* **23**, 1378–1385 (2007).
23. R Steuer, J Kurths, O Feihl, W Weckwerth, Observing and interpreting correlations in metabolomic networks. *Bioinformatics* **19**, 1019–1026 (2003).
24. A Izadpanah, et al., A short-term diet and exercise intervention ameliorates inflammation and markers of metabolic health in overweight/obese children. *Am. J. Physiol. Metab.* **303**, E542–E550 (2012).
25. M Kaern, TC Elston, WJ Blake, JJ Collins, Stochasticity in gene expression: from theories to phenotypes. *Nat. Rev. Genet.* **6**, 451–464 (2005).
26. CB Newgard, Metabolomics and metabolic diseases: where do we stand? *Cell Metab.* **25**, 43–56 (2017).
27. JC Dunn, Well-separated clusters and optimal fuzzy partitions. *J. Cybern.* **4**, 95–104 (1974).
28. DL Davies, DW Boulton, A cluster separation measure. *IEEE Transactions on Pattern Analysis Mach. Intell.* pp. 224–227 (1979).
29. T Caliński, J Harabasz, A dendrite method for cluster analysis. *Commun. Stat. Methods* **3**, 1–27 (1974).
30. FE Harrell Jr, KL Lee, DB Mark, Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Stat. Medicine* **15**, 361–387 (1996).
31. B Desgranges, MB Desgranges, Package 'clustercrit' (2018).
32. C Sudlow, et al., UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLOS Med.* **12**, e1001779 (2015).
33. P Elliott, TC Peakman, The UK biobank sample handling and storage protocol for the collection, processing and archiving of human blood and urine. *Int. J. Epidemiol.* **37**, 234–244 (2008).
34. P Würtz, et al., Quantitative serum nuclear magnetic resonance metabolomics in large-scale epidemiology: a primer on-omic technologies. *Am. J. Epidemiol.* **186**, 1084–1096 (2017).
35. P Soininen, AJ Kangas, P Würtz, T Suna, M Ala-Korpela, Quantitative serum nuclear magnetic resonance metabolomics in cardiovascular epidemiology and genetics. *Circ. Cardiovasc. Genet.* **8**, 192–206 (2015).
36. M Ester, HP Kriegel, J Sander, X Xu, A density-based algorithm for discovering clusters in large spatial databases with noise in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, KDD'96*, eds. E Simoudis, J Han, U Fayyad. (AAAI Press), p. 226–231 (1996).
37. E Uffelmann, et al., Genome-wide association studies. *Nat. Rev. Methods Primers* **1**, 59 (2021).
38. H Zhong, X Yang, LM Kaplan, C Molony, EE Schadt, Integrating pathway analysis and genetics of gene expression for genome-wide association studies. *The Am. J. Hum. Genet.* **86**, 581–591 (2010).
39. VP Mäkinen, et al., Longitudinal metabolomics of increasing body-mass index and waist-hip ratio reveals two dynamic patterns of obesity pandemic. *Int. J. Obes.* **47**, 453–462 (2023).
40. B Li, J Su, R Lin, S-T Yau, Z Yao, Github - zhihang-yao/mf-Metabolomic-Heterogeneity (2024) <https://github.com/zhihang-yao/MF-Metabolomic-Heterogeneity>, Deposited on 29 December 2024.