

1 Computational Strategies in Nutrigenetics: 2 Constructing a Reference Dataset of Nutrition- 3 Associated Genetic Polymorphisms

4 Giovanni Maria De Filippis^a (0009-0002-8395-0724), Maria Monticelli^{a,b} (0000-0003-3136-2138), Alessandra
5 Pollice^a (0000-0003-1619-5925), Tiziana Angrisano^a (0000-0002-5940-5278), Bruno Hay Mele^{a*} (0000-0001-
6 5579-183X *), Viola Calabro^{`a} (0000-0002-6508-8889)

7 ^a Department of Biology, University of Napoli "Federico II", Complesso Universitario Monte Sant'Angelo, Via Cinthia,
8 80126 Napoli, Italy

9 ^b Institute of Biomolecular Chemistry (ICB), National Research Council (CNR), Via Campi Flegrei 34, 80078 Pozzuoli,
10 Italy

11

12 Abstract

13 **Objective:** This study aimed to build a comprehensive dataset of human genetic polymorphisms associated
14 with nutrition by integrating data from multiple sources, including the LitVar database, PubMed, and the
15 GWAS catalog. Such a resource could facilitate the exploration of genetic polymorphisms associated with
16 nutrition-related traits.

17 **Methods:** We developed a Python pipeline to streamline the integration and analysis of genetic
18 polymorphism data associated with nutrition. We employed the MeSH ontology as a framework to aggregate
19 relevant genetic data. The pipeline comprises five distinct modules that go through the following steps: data
20 extraction from LitVar and PubMed articles, generation of a joint dataset by data merging, generation of
21 comprehensive MeSH term lists, filtering of the joint dataset using the selected MeSH sets, lexical analysis
22 and augmentation of the dataset with data from of the GWAS catalog dataset.

23 **Results:** We successfully aggregated a wide range of papers and data on genetic polymorphism and nutrition-
24 related traits into a single dataset. Cross-referencing with the GWAS catalog dataset provided information
25 about possible effects or risk alleles associated with the identified genetic polymorphisms. The nutrigenetic
26 dataset we developed is a tool for nutritionists and researchers, serving as a preliminary benchmark for
27 personalized nutrition interventions based on genetic testing.

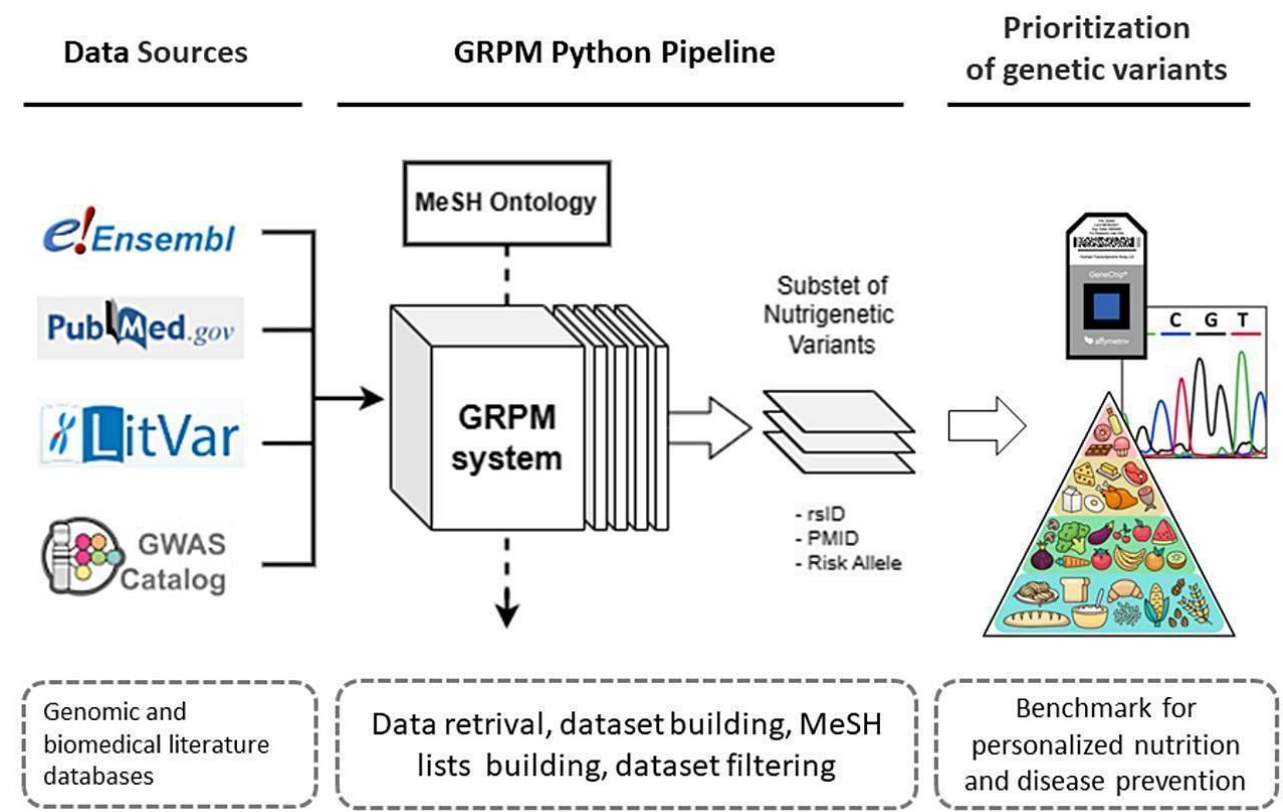
28 **Conclusion:** The pipeline presented here consolidates and organizes information on genetic polymorphisms
29 associated with nutrition, enabling comprehensive analysis and exploration of gene-diet interactions.

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

30 Overall, the method contributes to advancing personalized nutrition interventions and nutrigenomics
31 research. The flexible nature of the system allows its application to other investigations related to genetic
32 polymorphisms.

33 **Keywords:** Nutrigenetics, Genetic polymorphisms, Personalized nutrition, Gene-diet interactions, Data
34 integration, MeSH ontology

Graphical Abstract



35 1. Introduction

36 Nutrition is critical to health and disease (1). Emerging evidence suggests that genetic polymorphisms
37 significantly impact an individual's response to different nutrients and dietary patterns by affecting nutrient
38 bioavailability and metabolism (2). Moreover, it has been demonstrated that common gene variations are
39 linked to complex chronic health issues significantly affected by nutritional factors (3). Advancements in
40 genomics technologies and the subsequent availability of large-scale genetic data have fueled interest in
41 identifying and categorizing genetic polymorphisms associated with nutrition-related traits (4). The field of
42 nutrigenetics was thus born to comprehend how genetic variations influence an individual's nutritional
43 requirements, metabolism, and health outcomes (5). By considering an individual's genetic profile,
44 healthcare professionals and nutritionists can provide tailored dietary advice and interventions that optimize
45 nutrient bioavailability and promote better health outcomes in that individual (6). Nutrigenetic associations
46 imply that specific genetic polymorphisms can induce susceptibility to chronic diseases. The response to
47 specific nutrients or dietary patterns may be crucial in determining health outcomes (6).

48 Recent literature contains extensive data on nutrition-associated genetic polymorphisms (2,7). However,
49 these data are often scattered, diverse in format, and lack a standardized curation process. Such
50 complications hinder data integration, limit information extraction and synthesis, and pose a barrier to data
51 utilization in decision support systems (8). Integrating available data and overcoming the limits of self-
52 reported methods in research is crucial for accurate omics data integration, nutrigenetics, and nutrigenomics
53 research, especially in clinical settings (8). Therefore, there is a need to develop curated and consolidated
54 resources that integrate nutrition-associated genetic polymorphism data, along with omics data, to advance
55 personalized nutrition interventions and clinical decision-making.

56 Here, we built a structured dataset of human genetic polymorphisms associated with nutrition by integrating
57 data from different established sources: the LitVar database (9), which contains curated information on
58 genetic variations and their functional effects; the Pubmed-Medline database, which provides structured
59 MeSH ontology annotations; and the GWAS catalog dataset, which reports human variant-traits associations.
60 Our dataset includes data from Medline studies associated with nutrition-related genetic polymorphism.
61 Specific sets of MeSH terms related to nutrition physiology, nutrition-related diseases, prevention through
62 diet, and eating behavior were used to extract subsets of genes and their single-nucleotide polymorphisms
63 (SNPs) potentially associated with nutrition-related traits. Cross-referencing with the GWAS catalog dataset
64 (10) provided information about effect/risk alleles associated with the collected studies. The database was
65 curated to ensure data quality, consistency, and relevance to nutrition and nutrigenomics research, thus
66 providing a valuable resource to investigate the intricate interplay between genetics and nutrition.

67

68 2. Methods

69 2.1 General Presentation

70 This study uses Medical Subject Headings (MeSH) to connect genetic polymorphism data from various
71 resources. *MeSH ontology* (11) is a standardized and controlled vocabulary that offers descriptors utilized in
72 biomedicine and informatics to classify and categorize biomedical literature and data. The standardized and
73 controlled nature of MeSH terms makes them highly adaptable for broader utilization in analyzing scientific
74 indicators. MeSH terms associated with a document can be thought of as references to a collection of
75 knowledge stored as documents in a database (12).

76 To build a structured dataset that aligns genomic data and scientific papers, we connected the LitVar and
77 PubMed databases through shared MeSH terms. The LitVar database is a comprehensive and publicly
78 accessible resource that collects information on genetic variations and their associated scientific literature.
79 It aims to bridge the gap between genomic data and the relevant literature by aggregating and organizing
80 information on genetic variants from a wide range of sources (9). PubMed serves as the principal repository
81 of biomedical literature. Extracting data from PubMed is crucial for various research purposes, such as
82 literature reviews, data mining, and knowledge discovery (13).

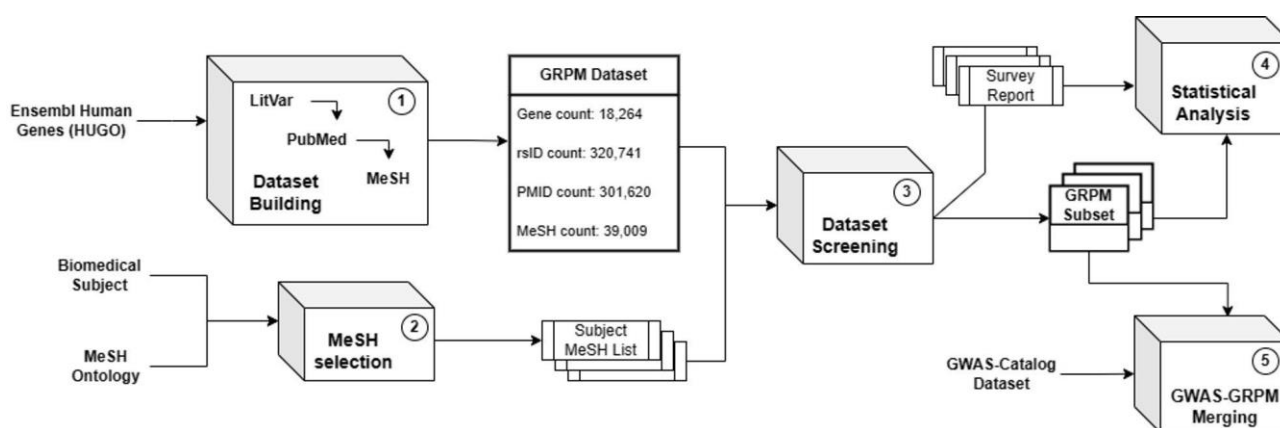
83 We developed a Python pipeline that leverages the MeSH ontology as a crucial framework to aggregate
84 genetic polymorphism data for a topic of interest effectively. We named the tool "GRPM system" (Gene-RsID-
85 PMID-MeSH) and specifically designed it to streamline the integration and analysis of genetic polymorphism
86 data associated with a given biomedical field, such as nutrition. With the increasing importance of genetic
87 factors in understanding nutrition-related traits, the GRPM system could help researchers and nutritionists
88 explore and analyze such data efficiently.

89 The system is written in Python, exploits the Jupyter Notebook format (14), and comprises five modules
90 designed for specific purposes (Figure 1). These modules support the following operations:

- 91 1. Data Retrieval and Merging: facilitates data extraction, integration, and consolidation from source
92 databases, including LitVar and PubMed, ensuring a comprehensive collection of genetic
93 polymorphisms associated with topic-related traits (GRPM dataset).
- 94 2. MeSH Term List Creation: generates coherent MeSH term lists and links them to the collection,
95 enabling efficient exploration of the GRPM dataset. This enrichment empowers users to access
96 specific genetic information relevant to the chosen topic easily.
- 97 3. Dataset Filtering with Selected MeSH Terms: enables the user to screen collected data using selected
98 MeSH terms. This way, users can refine their search and focus on specific areas of interest within the
99 GRPM dataset.

- 100 4. Statistical Analysis: assigns each gene a relative measure of interest based on the number and
101 proportion of associated findings gathered in the GRPM dataset. This calculation enables the
102 prioritization of genes related to a chosen topic, aiding further investigation or personalized nutrition
103 interventions.
- 104 5. GWAS Dataset Incorporation: integrates the GWAS catalog dataset into the collection, providing
105 insights into effect/risk alleles associated with identified genetic polymorphisms. This process
106 enriches our nutrigenetic resource with supplemental data for further analysis.

107 More detailed information about the GRPM system, including implementation details and usage instructions,
108 is available on GitHub¹.



109 **Figure 1:** A graphical overview of the GRPM system workflow showcasing the input data and interactions between the
110 five modules. A brief quantitative description of the GRPM Dataset is also shown.

111

112 2.2 GRPM Dataset Building

113 The first module uses the LitVar Application Programming Interface (API)² to retrieve all polymorphisms for
114 each human gene³ within the LitVar database alongside all associated PubMed Identifiers (PMIDs). These
115 PMIDs were subsequently employed as queries on PubMed to obtain MEDLINE data. We utilized the NBIB
116 parser⁴, a Python package designed explicitly for parsing MEDLINE-PubMed data, to streamline the data
117 collection. The collected data were ultimately consolidated into a single CSV file (from now on called “GRPM
118 ds”), serving as the primary source against which MeSH term queries can be launched to retrieve genes and
119 polymorphisms associated with specific contexts.

¹ GRPM_system (github.com): https://github.com/johndef64/GRPM_system

² LitVar API Docs: <https://www.ncbi.nlm.nih.gov/CBBresearch/Lu/Demo/LitVar/api.html>

³ The HGCN names lists of human protein-coding genes, RNA genes, and pseudogenes were retrieved from the Ensembl database through the BioMart application.

⁴ nbib · PyPI: <https://pypi.org/project/nbib/>

120 **2.3 Dataset filtering and screening**

121 The GRPM system is designed to retrieve subsets of genes and polymorphisms from GRPM ds, employing a
122 user-defined list of MeSH terms as a hook. Careful selection of the MeSHs is crucial at this stage: the list must
123 represent the chosen search field out of the total complex of terms in GRPM ds. For this purpose, we referred
124 to the subset of 21.705 MeSH terms related to LitVar publications retrieved instead of the complete MeSH
125 ontology (348.733 terms)⁵. This subset collects MeSH terms linked to papers exploring the associations
126 between genetic variants and biomedical traits.

127 The second module is built to subset this extensive MeSH collection with representatives of a particular
128 biomedical field. We utilized the natural language processing capabilities of Generative Pre-trained
129 Transformer (GPT) provided by OpenAI⁶ through its API. The procedure involved generating simple lists of
130 words through one or more biologically relevant prompts (see Supplementary Materials). These lists are used
131 to extract the real MeSH terms related to the subject from our dataset (70-400 MeSHs extracted for each
132 query). Subsequently, the extracted MeSH terms were manually screened to eliminate ambiguous and bias-
133 generating terms. This filtering process ensures that only appropriate and meaningful terms are utilized for
134 the subsequent full dataset screening.

135 The screening of the entire dataset using the selected MeSH set (our *query*) is obtained by running our third
136 Jupyter module. It takes approximately three to four hours on an average workstation using MeSH lists
137 ranging in size from 100 to 400 terms. When the procedure ends, the system generates a comprehensive
138 report and a curated "Survey Dataset" that captures the essential association data. The reports generated
139 from various surveys are subjected to individual analysis and comparative examination in the fourth module.

140 **2.4 Gene Interest Index (GI)**

141 We consider a gene "interesting" if its related SNPs are associated with a substantial number of PMIDs (*i.e.*,
142 scientific papers) that include relevant MeSH terms and if the ratio between these relevant papers and the
143 total number of papers associated with the gene is sufficiently high.

144 To assess the relevance of the gene set retrieved for the chosen topic, it is crucial to consider the MeSH set
145 employed as a single entity rather than independently, given the difference in the relative importance of
146 terms. To define a gene as "interesting" based on its associated MeSH terms from related LitVar studies, we
147 propose scaling the number of detected PMIDs (PubMed IDs) by all the PMIDs associated with that gene in

⁵ The complete Medical Subject Headings dataset can be downloaded at <https://www.nlm.nih.gov/mesh/meshhome.html>

⁶ OpenAI. (2023). ChatGPT (Mar 14 version) [Large language model]. <https://chat.openai.com/chat>

148 LitVar. This approach helps minimize selection bias caused by extensively studied genes associated with more
149 MeSH terms than others — these terms could not be directly correlated with the query topic.

150 Given the set of genes $L(i)$ retrieved with the query (j), we introduce the following indices:

- 151 1. Pg_i : The total number of PMIDs associated with gene i ;
- 152 2. $Pm_{i,j}$: The number of i -related PMIDs containing at least one MeSH from the query j ;
- 153 3. Pm_{\max} : The highest $Pm_{i,j}$ value across all the genes in L ;
- 154 4. $Pm\ score_{i,j}$: the $Pm_{i,j}$ value normalized Pm_{\max} ;
- 155 5. $Pm\ ratio_{i,j}$: the ratio of Pm to Pg . It measures the proportion of matching PMIDs to the total PMIDs
156 associated with the gene.

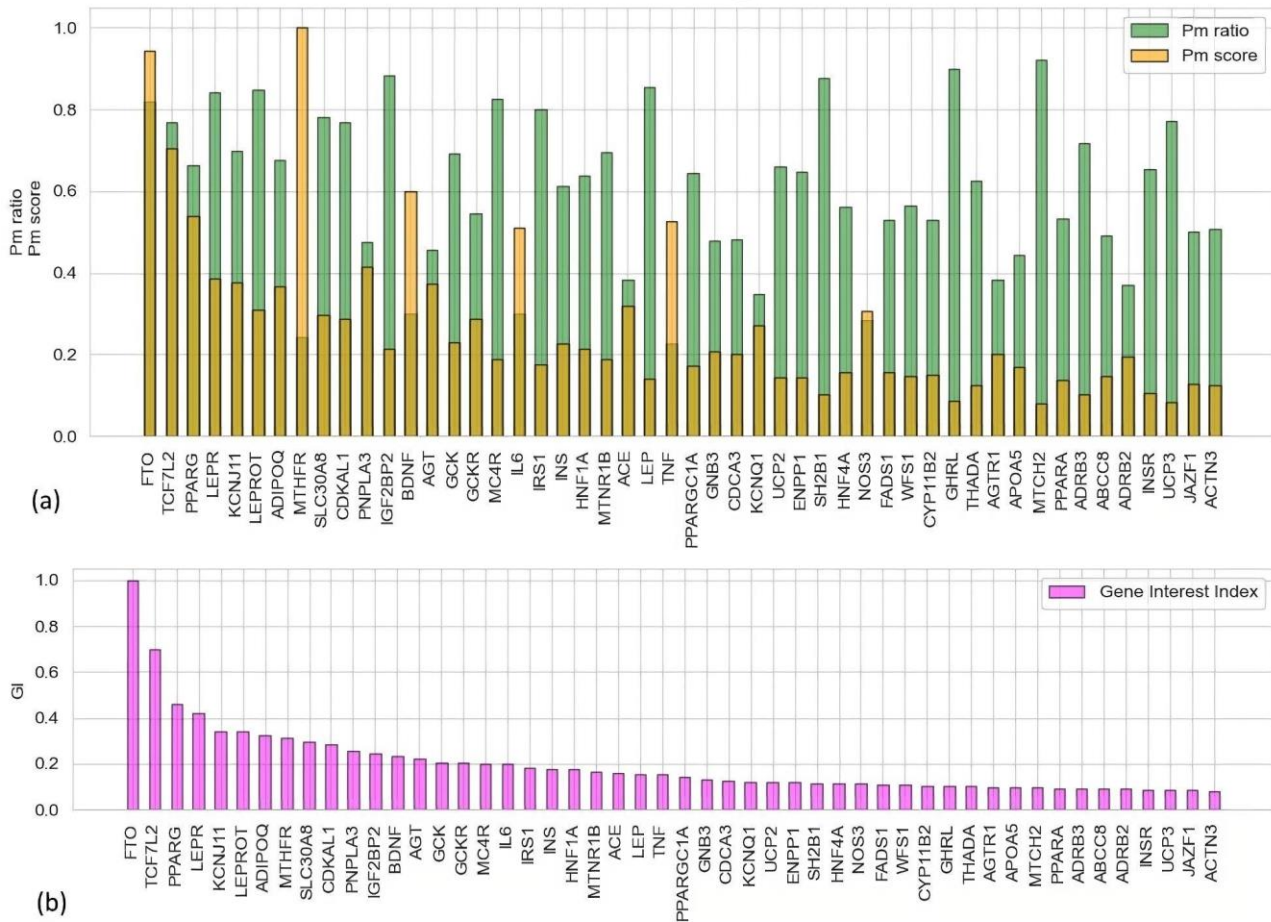
157 Based on these indices, we introduce the “Gene Interest Value” (GV), calculated as the product of “Pm score”
158 and “Pm ratio” and its normalized form, the “Gene Interest Index” (GI), which is adjusted relative to the
159 maximum value obtained in the survey. The ratio serves as a modifier in determining the level of interest for
160 each gene.

$$161 \quad GV_{ij} = Pm\ index \cdot Pm\ ratio = \frac{Pm_{ij}}{Pm_{\max}} \cdot \frac{Pm_{ij}}{Pg_i} \quad (1.1)$$

$$162 \quad GI = \frac{GV}{GV_{\max}} \quad (1.2)$$

163 By integrating the Pm score and Pm ratio, the GI method acts as a coherent measure of gene relevance.
164 Figure 2 visually represents an example of gene prioritization obtained through the Index using the “Obesity
165 and Weight Control” MeSH list as a reference. Panel (a) shows the Pm ratio (green) and Pm score (yellow). It
166 highlights the importance of considering both indexes, which produce different orders. In Panel (b), the gene
167 relevance-based sorting achieved with the GI is presented, and it is possible to appreciate the highest
168 prioritization performance versus the other two. The integrated assessment provided by the GI method
169 allows for more accurate gene prioritization, leading to a deeper understanding of gene-gene interactions
170 and potential therapeutic targets in obesity and weight control management. Another example of gene
171 prioritization through GI is presented in Supplementary Materials (Figure S1).

172



173 **Figure 2:** Visual representation of Gene Interest Index (GI) calculation through *Pm* score and *Pm* ratio, using as reference
 174 the results from "Obesity and Weight Control" MeSH query (see results section). Panel (a) shows the matching PMID
 175 ratio and overall matching PMID score. Panel (b) displays the gene relevance sort achieved with the GI.

176 2.5 GWAS Catalog data integration

177 While examining every paper to unravel the associated effect allele for each SNP ID (rsID) can be time-
 178 consuming, an initial indication of the potential effect allele is valuable for conducting preliminary studies.
 179 To address this issue, we leveraged Ensembl GWAS Catalog⁷ data (10).

180 To integrate the GWAS data with the GRPM dataset, we followed a specific workflow (fifth module). First, we
 181 retrieved the GRPM Survey data. Then, we applied a Gene Interest (GI) cut-off of 0.0125 to the GRPM Survey
 182 data to prioritize the relevant genes. Most protein-coding genes have cited works with at least one of the
 183 MeSHs of the lists used, but this does not imply that it is relevant to consider them all. By setting a GI
 184 threshold value, we aimed to prioritize the genes that most fit with our tailored MeSH terms, allowing us to
 185 focus on genes and their SNPs that demonstrated a higher degree of interest and relevance in the field of
 186 nutrigenetic dietary advice. This approach helped remove noise or irrelevant results from the search process,

⁷ GWAS Catalog (ebi.ac.uk): <https://www.ebi.ac.uk/gwas/docs/file-downloads>

187 allowing us to focus on genes more likely to provide valuable insights into gene-diet interactions and
188 personalized nutrition interventions. At this point, the filtered GRPM Survey and GWAS dataset were merged
189 based on the rsIDs. The merge was efficiently aligned with the GRPM MeSH terms through a correspondence
190 dictionary. We subsequently utilized the Natural Language Toolkit (NLTK)⁸ to tokenize the MeSH terms (and
191 all their possible synonyms) and GWAS-mapped traits to perform the alignment. Finally, we retrieved the
192 strongest SNP-risk allele for each rsID using the correspondence dictionary. This information serves as an
193 initial indication and can be beneficial for conducting further studies, whether in a clinical or in-silico setting,
194 based on the identified associations.

195 **3. Results**

196 **3.1 GRPM Dataset Statistics**

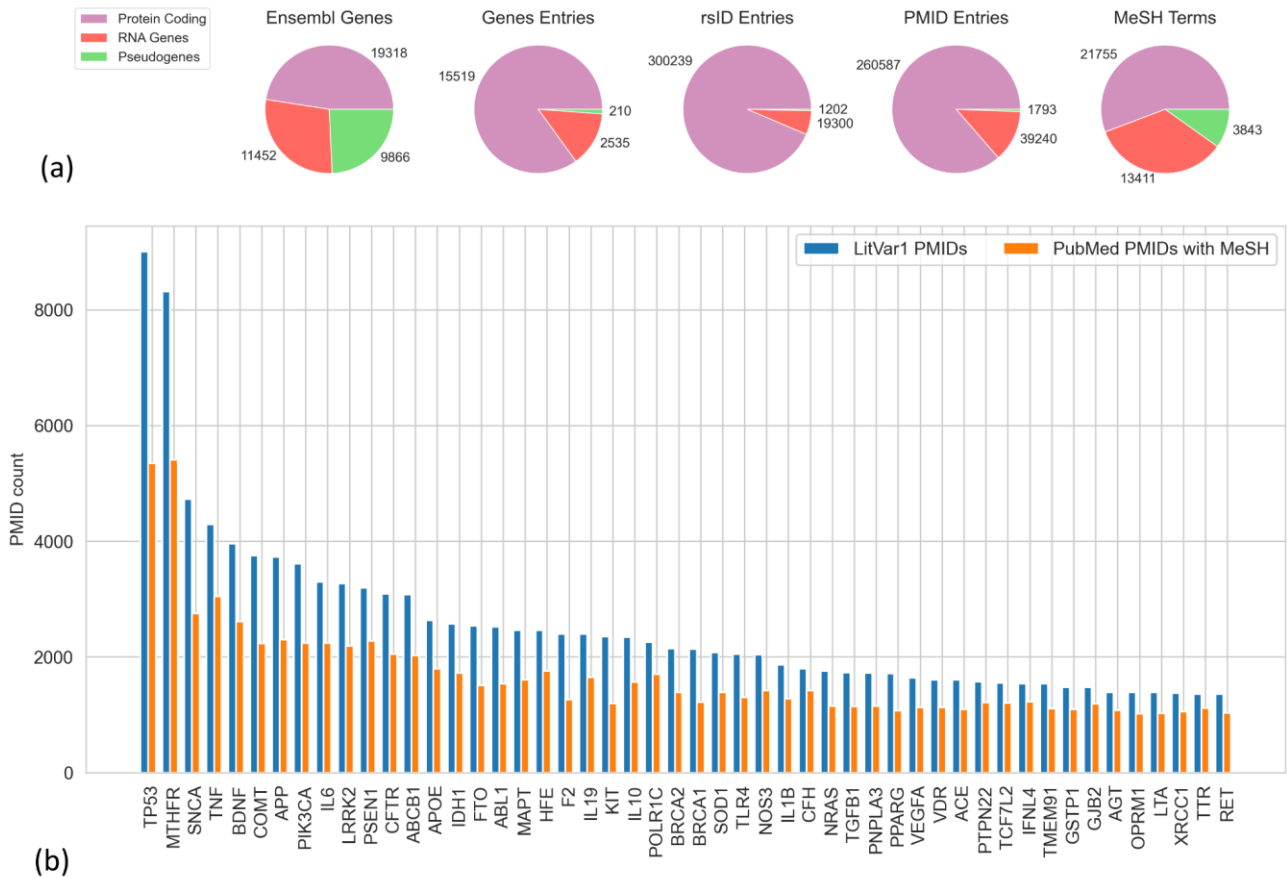
197 When utilizing the dataset, it is essential to account for each gene's relative richness in PMID. The genes that
198 garnered the most attention in research are associated with a higher number of PMIDs, resulting in more
199 MeSH annotations. Hence, to ensure data normalization, it is imperative to obtain a comprehensive
200 understanding of the most represented genes within the dataset.

201 We divided the GRPM dataset⁹ into three distinct partitions based on the gene type included: Protein Coding
202 Genes, RNA Genes, and Pseudogenes. Since some genes are absent in LitVar, the dataset covers 80% of
203 human protein-coding genes. Figure 3 (a) shows the statistics for each partition, including the number of
204 genes, rsID (reference SNP ID number) entries, PMID entries, and MeSH terms associated with each gene
205 type. These genes were considered for further analysis and integration with other data sources. It is possible
206 to appreciate how pseudogenes are associated with a negligible relative fraction of variants and publications
207 (0.003% and 0.005%, respectively), although they account for 1/3 of Ensembl Genes entries. Figure 3 (b) shows
208 the analysis of the fifty most extensively studied and represented genes within the GRPM dataset. Among
209 them, TP53 and MTHFR concentrate the highest number of PMIDs, with and without MeSHs. This deviation
210 is expected since both encode proteins extremely relevant to human metabolism and health — TP53
211 (ENSG00000141510) encodes the Cellular tumor antigen p53 (P53_HUMAN), a tumor suppressor that
212 monitors DNA integrity and initiates cellular responses to prevent tumor formation; MTHFR
213 (ENSG00000177000) encodes the methylenetetrahydrofolate reductase (MTHR_HUMAN), that takes part in
214 homocysteine metabolism and is essential for methylation reactions.

⁸ NLTK - Natural Language Toolkit: <https://www.nltk.org/>

⁹ Available on Zenodo at <https://zenodo.org/record/8205724> DOI: 10.5281/zenodo.8205724

215 A conspicuous portion of LitVar PMIDs (~77%) PMIDs extracted from LitVar is associated with MeSH terms,
 216 ensuring consistency and reliability for further investigations.



217 **Figure 3** Statistics concerning the primary GRPM dataset. In (a), the diagram showcases the distribution of genes, rsID,
 218 PMID, and MeSH terms across the three dataset partitions. (b) depicts the top fifty most represented genes in the
 219 dataset, along with their occurrence in PMID counts and PMID associated with MeSH terms.
 220

221 3.2 Nutrigenetic Dataset

222 Our study aimed to create a nutrigenetic dataset using a collection of MeSH terms related to different aspects
 223 of nutrition. We utilized specific MeSH terms that covered nutrition physiology, nutrition-related diseases,
 224 disease prevention through diet, and eating behavior. This approach allowed us to process LitVar-PubMed
 225 data in a way relevant to personalized nutritional approaches. In Table 1, we describe each selected field of
 226 interest in the context of a personalized nutritional approach.

227 **Table 1:** Categories of nutrition-related MeSH terms used to build the nutrigenetic database.

Category	Description	MeSH count
General Nutrition	A broad range of topics related to nutrition, including dietary patterns, nutrient requirements, nutritional status, and the impact of nutrition on overall health and well-being.	413
Obesity, Weight Control, and Compulsive Eating	Terms related to weight management, including obesity, weight loss strategies, and disorders such as binge eating or compulsive overeating.	243
Cardiovascular Health and Lipid Metabolism	Terms related to nutrition and cardiovascular health, including the impact of dietary factors on lipid metabolism, cholesterol levels, and the prevention of cardiovascular diseases.	319
Diabetes Mellitus Type II and Metabolic Syndrome	Terms related to type II diabetes and metabolic syndrome. Including dietary interventions, glucose metabolism, insulin resistance, and related complications.	528
Vitamin and Micronutrients Metabolism and Deficiency-Related Diseases	Terms related to the metabolism of essential vitamins and micronutrients, the impact of deficiencies on health, and the development of associated diseases.	175
Eating Behavior and Taste Sensation	Terms related to individual eating behaviors, including factors influencing food choices, taste preferences, satiety, and appetite regulation.	292
Food Intolerances	Terms related to adverse reactions to specific foods, such as lactose intolerance or gluten sensitivity. Explores the genetic and physiological factors underlying food intolerances and their impact on dietary choices.	145
Food Allergies	Examines the genetic basis of food allergies, identification of allergenic components, and strategies for managing allergic reactions through diet.	65
Diet-induced Oxidative Stress	Explores the relationship between dietary factors and oxidative stress and investigates the impact of diet on oxidative stress levels and its health implications.	77
Xenobiotics Metabolism	Focuses on the metabolism of foreign substances (xenobiotics) in the body, including drugs, environmental toxins, and dietary components.	170

228

229 Figure S2 presents an overview of the most interesting genes with their relative MeSH and PMID values on
 230 the ten nutrient lists used in our study. This figure offers a representative example of the analysis conducted
 231 and allows for a quick comparison of gene relevance across different nutrient-related traits. The Figure allows

232 for a quick comparison of gene relevance across different nutrient-related traits, showing the relative
233 richness of the associated MeSH terms and the papers associated with each gene.

234 Table 2 provides statistical data on the ten nutritional MeSH term lists employed. The table left section
235 provides information on data retrieved with every nutritional MeSH term list scanning the GRPM dataset. In
236 a large dataset like this, there can be a wide range of genes and associated MeSH terms, making it challenging
237 to identify the genes of highest interest and relevance to the specific research focus. For this reason, the right
238 section presents the same data filtered Gene Interest (GI) < 0.0125, representing the most significant
239 matches available for further investigation and nutrigenetic applications. Applying a GI cut-off was necessary
240 to select the most relevant and meaningful search results.

241 **Table 2:** Statistical data associated with the ten nutritional MeSH lists applied to the build dataset.

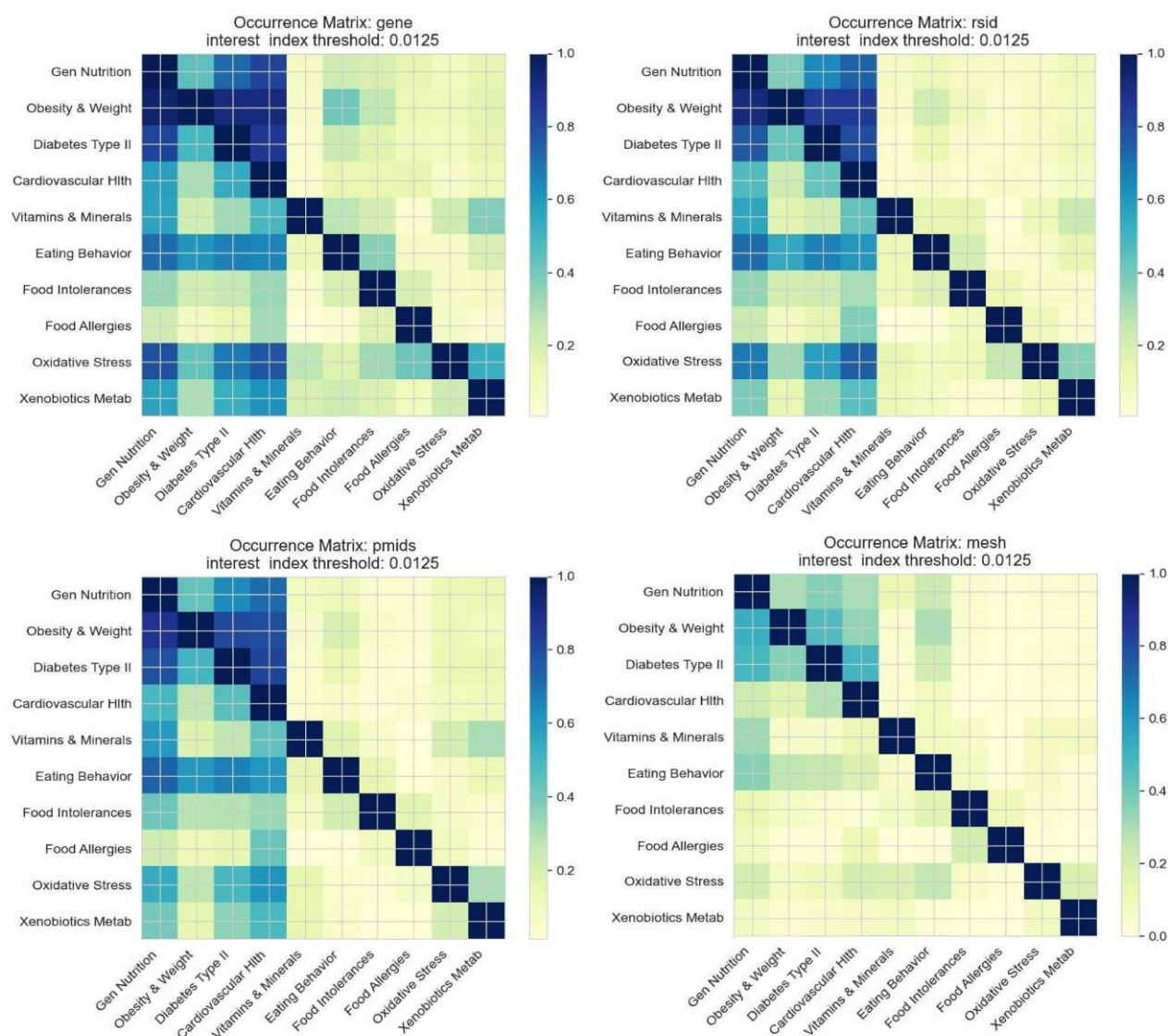
label	All MeSH matching in DB				Interesting entries based on GI threshold (0.0125)			
	#gene	#rsID	#PMID	#MeSH	#gene	#rsID	#PMID	#MeSH
General Nutrition	11,560	83,288	62,473	413	686	26,456	44,859	397
Obesity & Weight Control	9,713	53,879	35,563	243	317	10,842	22,123	230
Diabetes Type II & MS	10,717	68,844	49,896	319	603	22,270	36,198	297
Cardiovascular Health	12,368	105,598	85,065	528	975	41,931	66,113	521
Vitamins & Minerals	4,045	16,857	11,941	175	89	3,525	6,882	147
Eating Behavior	5,525	20,607	13,734	292	211	4,252	7,241	256
Food Intolerances	4,040	14,117	7,416	145	392	5,008	4,726	125
Food Allergies	4,681	16,777	11,032	65	451	6,289	7,762	64
Oxidative Stress	5,156	20,919	19,295	77	75	2,559	10,058	60
Xenobiotics Metab	7,115	35,686	27,237	170	173	7,159	14,171	151

242

243

244

245 We constructed co-occurrence matrices to explore the extent of overlap within the data obtained from the
246 ten nutritional MeSH lists. Figure 4 shows the co-occurrence patterns between genes, rsID, PMID, and MeSH
247 terms among filtered results by $GI > 0.0125$. The heatmaps¹⁰ in Figure 4 show interesting overlap patterns
248 among the different MeSH lists. For example, the results obtained from the "obesity" lists show a high overlap
249 with the "diabetes" and "cardiovascular health" lists, sharing approximately 80% of the items. The latter
250 shows a 40-60% overlap in this first group of three.



251 Figure 4: Co-occurrence patterns among genes, rsID, PMID, and MeSH terms observed in the filtered results
252 ($GI > 0.0125$) from the ten nutritional MeSH lists used for screening. The color key represents the amount of
253 overlap between datasets scaled over the total size of each dataset.

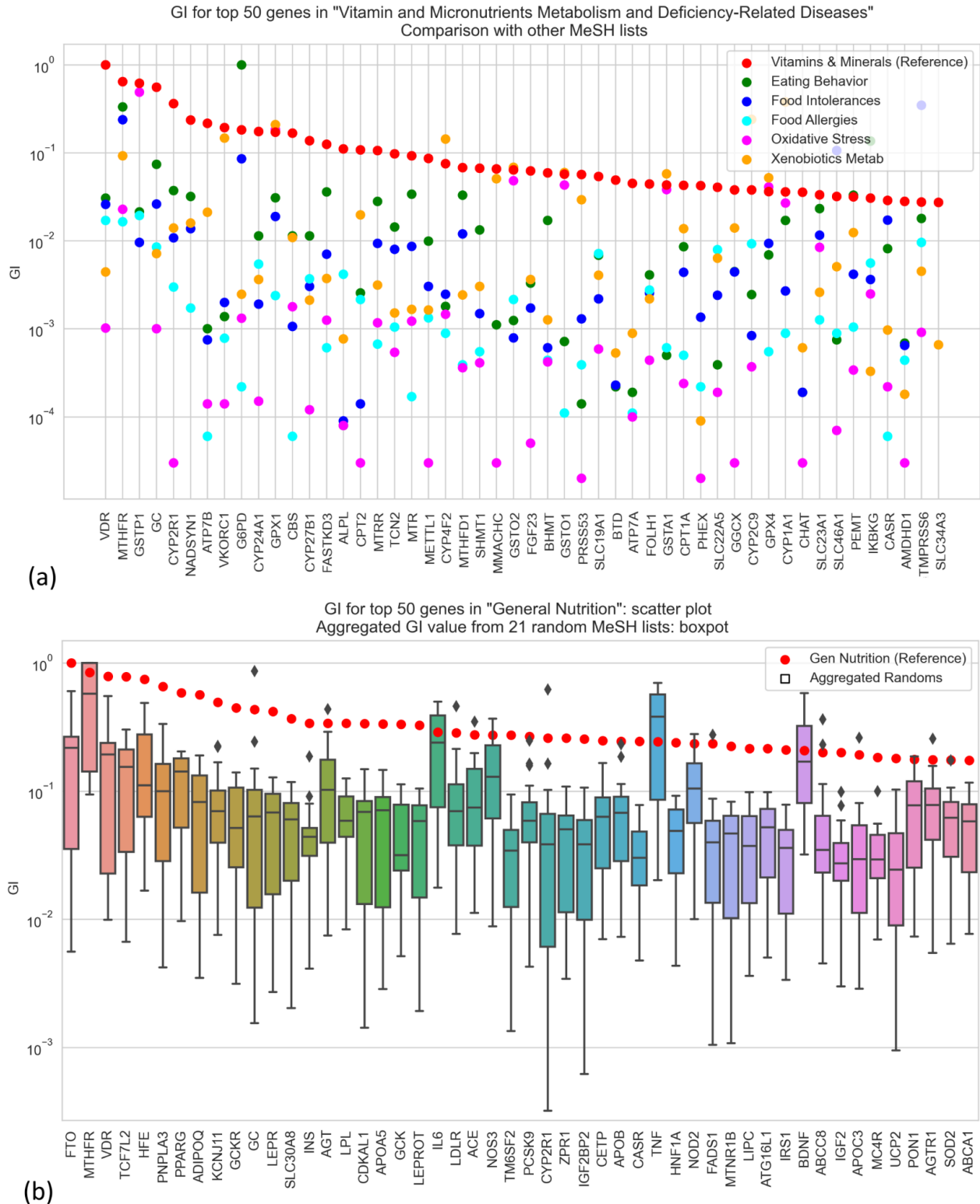
¹⁰ In the heatmap, we calculated row-wise (i.e., dataset) relative values as the ratio between number of shared entities and total number of entities of the row. This approach allows for a focused examination of co-occurrence within specific categories, providing a granular understanding of the relationships and interactions at the row level.

254 **3.3 Method Validation**

255 Figure 5 (a) shows the fifty most interesting genes in the "Vitamin and Micronutrient Metabolism" list
256 compared to the values of interest on results obtained from 5 other nutritional MeSH term lists. The genes
257 extracted using the GRPM system show specificity for the list of MeSH terms used as a reference model. This
258 behavior suggests the GRPM system can identify genes related to that particular nutritional aspect. However,
259 some of them are more interesting in other MeSH term lists than the one taken as a reference, meaning they
260 are more closely associated with other nutritional features or specific biological processes. This observation
261 suggests the complexity of gene regulation in nutrient metabolism and underscores the importance of
262 considering a broader range of nutritional MeSH terms to gain a comprehensive understanding of the
263 biological system under consideration. Supplementary Materials provide additional GI comparison results,
264 showcasing the comparison among another MeSH list utilized in our study (Figure S3).

265 To ensure the accuracy and reliability of the data collected, we compared the results obtained with
266 biologically consistent MeSH lists with those obtained with twenty random MeSH lists¹¹ of the same size.
267 Figure 5 (b) provides an example of the comparison results on the "General Nutrition" list. Another reference
268 is shown in Supplementary Materials (S4).

¹¹ Containing 450 random terms each, based on 21,705 LitVar MeSH in the GRPM dataset.



269 **Figure 5:** (a) Comparison of Gene Interest Index (GI) for the top fifty genes on the "Vitamin and Micronutrient
270 Metabolism" MeSH list with GI obtained from five different lists of nutrition MeSH terms. (b) Comparison between GI
271 obtained from the "General Nutrition" list and GI obtained using twenty randomly generated MeSH lists of the same
272 size represented as boxplots. In both graphs the y-axis is logarithmic; genes are ordered by decreasing interest relative
273 to the reference list.

274 3.4 GWAS data integration

275 Cross-referencing data between the GRPM dataset and the GWAS catalog dataset provides indicative
276 information about possible risk alleles associated with the collected studies. **Table 3** shows a sample of this
277 integration on the GRPM "General Nutrition" dataset, along with the corresponding GWAS catalog
278 information, such as Mapped Trait and Strongest Risk Allele. The preliminary results obtained through NLTK
279 show congruence between the MeSH associated with the PMID and the mapped GWAS trait.

280 The merging process on the "General Nutrition" dataset, with a GI cut-off of 0.0125, resulted in the following
281 statistics: The number of genes from the LitVar database was 365, while the number of genes with mapping
282 information was 359. There were 186 MeSH terms, 467 mapped traits, 1155 disease/traits, and 1678
283 identified strongest SNP-risk alleles.

284 **Table 3:** Sample of fifteen GRPM associations from the "General Nutrition" list merged on rsID with the GWAS Catalog
285 dataset.

GRPM Data			GWAS Data			
LITVAR GENE	LIVAR RSID	LITVAR PMID	PUBMED MeSH	MAPPED TRAIT	DISEASE/TRAIT	STRONGEST SNP-RISK ALLELE
SEM1	rs7781370	22698912	Body Mass Index	body height	Height	rs7781370-T
MTHFR	rs9651118	33213085	Hypertension	red blood cell distribution width	Red cell distribution width	rs9651118-T
CADM2	rs13078807	25893265	Pediatric Obesity	obesity	Obesity	rs13078807-G
C1QTNF6	rs229533	25751624	Diabetes Mellitus, Type 1	type 2 diabetes mellitus	Type 1 diabetes	rs229533-C
GHR	rs6184	34074802	Body Mass Index	body height	Height	rs6184-A
TMEM258	rs102275	31636271	Lipoproteins, LDL	high density lipoprotein cholesterol measurement	Fasting total cholesterol in large HDL	rs102275-C
SHROOM3	rs56281442	34502231	Diabetes Mellitus, Type 2	type 2 diabetes mellitus	Type 2 diabetes	rs56281442-G

LDLR	rs6511720	27973560	Coronary Artery Disease	coronary artery disease	Coronary artery disease	rs6511720-T
TP53INP1	rs10097617	30054458	Diabetes Mellitus, Type 2	type 2 diabetes mellitus	Type 2 diabetes	rs10097617-T
PROX1	rs2075423	32390949	Diabetes, Gestational	type 2 diabetes mellitus	Type 2 diabetes	rs2075423-G
HFE	rs1799945	11473047	Blood Glucose	systolic blood pressure	Systolic blood pressure	rs1799945-G
ESR1	rs6902771	31213659	Body Mass Index	body weight	Weight	rs6902771-T
LEP	rs17151919	33631239	Leptin	leptin measurement	circulating leptin levels	rs17151919-A
CYP2R1	rs7129781	23456391	Vitamin D3 24-Hydroxylase	vitamin D measurement	Vitamin D levels	rs7129781-C
CYP2C19	rs4494250	27618448	Hypertension	diastolic blood pressure	Diastolic blood pressure	rs4494250-A

286

287 4. Discussion

288 Understanding how genetic variations influence individual nutritional requirements, metabolism, and health
289 outcomes is crucial for developing personalized nutrition interventions (15). By considering an individual's
290 genetic profile, healthcare operators and nutritionists can provide tailored dietary advice, optimizing nutrient
291 bioavailability, and promoting better health outcomes, thus preventing chronic diseases such as obesity (16),
292 diabetes (17), or cardiovascular diseases (18).

293 Personalized approaches in healthcare and prevention are at the forefront of translational bioinformatics
294 (19). The development of computational methods and tools for consolidating and analyzing information from
295 multiple sources enhances data integration, enabling the translation of findings into personalized nutrition
296 interventions and disease prevention strategies (20,21).

297 The GRPM resource is an integrated platform for diverse data sources related to genetic polymorphisms
298 associated with nutrition. This resource enables efficient retrieval, merging, and analysis, facilitating
299 comprehensive research in nutrigenomics. Our approach leverages data mining and merging techniques to

300 identify relevant studies on nutrition-associated genetic polymorphisms based on specific MeSH term sets.
301 The pipeline allowed us to extract subsets of genes and associated SNPs linked to nutrition-related traits.
302 Machine learning algorithms could further enhance the analysis of these data to uncover intricate gene-diet
303 interactions by enabling the discovery of patterns and correlations and by producing predictive models (22).
304 Indeed, although focused primarily on genetic polymorphism and associated literature, integrating our study
305 with multi-omics data could provide further insight into the interplay between genetics and nutrition, thus
306 providing a more holistic understanding of such a complex biological field (22).

307 It is essential to highlight that this method's potential applications extend beyond the scope of our study.
308 Our approach can be employed to gather specific genetic polymorphisms associated with various health or
309 biological dysfunctions, empowering healthcare practitioners to tailor interventions based on an individual's
310 genetic profile (23).

311 As remarked by Floris and co-workers, the current approach adopted by many companies in nutrigenetic
312 counseling still relies on a limited set of genes and polymorphisms for genetic testing and counseling (24).
313 However, as sequencing costs continue to decrease and sequencing technologies become more accessible
314 than before, it is no longer justifiable to base nutrigenetic panels on a small number of genetic markers (24).
315 Using a limited set of genes and polymorphisms may overlook significant genetic variations that affect an
316 individual's response to nutrients and dietary patterns. Our study addresses the limitations of current
317 approaches in nutrigenetics by consolidating and standardizing information on genetic polymorphisms
318 associated with nutrition. Our nutrigenetic dataset offers a broader scope and coverage, improving the global
319 understanding of the interplay between genetics and nutrition-related traits.

320 Within GRPM ds, MeSH term richness across nutrigenetic categories is quite heterogeneous, confirming that
321 some categories may have a broader range of MeSH terms, covering various aspects of nutrition, while others
322 may have a narrower focus, addressing specific subtopics within nutrition.

323 Regarding our nutrigenetic dataset, particularly concerning the genes found on the MeSH lists described in
324 Table 2, it is worth mentioning that the most interesting genes identified in our study are well-known in the
325 field of nutrigenetics. These genes have been extensively studied and represented in literature. The fact that
326 these genes consistently appear in multiple MeSH lists further supports their relevance concerning nutrition-
327 related traits.

328 From a broader standpoint, the analysis revealed a high degree of overlap between genes associated with
329 specific nutritional aspects (Figure 4). For example, the "obesity" lists exhibit high overlap with the "diabetes"
330 and "cardiovascular health" lists, indicating shared genetic polymorphisms and pathways. This finding is not
331 surprising given the close relationship between obesity, diabetes, and cardiovascular health, as these

332 conditions often coexist and share common genetic and physiological factors. The high overlap suggests that
333 shared genetic polymorphisms and pathways may be involved in these conditions.

334 Conversely, the lower degree of overlap between specific lists could be attributed to the specificity of these
335 lists, focusing on more specific biological processes or conditions with distinct genetic underpinnings. This
336 behavior can be attributed to several factors. Firstly, these lists may be described by fewer MeSH terms,
337 leading to a narrower focus and less overlap with other lists. Secondly, these lists may pertain to more specific
338 biological processes or conditions with distinct genetic underpinnings than the broader conditions captured
339 by the other lists.

340 The results obtained from cross-referencing the GRPM dataset with the GWAS catalog dataset further
341 validated the findings, providing information about potential risk alleles associated with the identified genetic
342 polymorphisms. The preliminary results obtained through NLTK show congruence between the MeSH
343 associated with the PMID and the mapped GWAS trait (Table 3).

344 We thoroughly validated our method to ensure the validity and reliability of our results. Firstly, we compared
345 the most interesting genes associated with each MeSH list to the results obtained from other nutritional
346 MeSH term lists (Figure 5, A). This validation step demonstrated the effectiveness of the GRPM system in
347 identifying genes specifically related to the chosen nutritional aspect. However, we also observed that
348 specific genes showed higher interest in other MeSH term lists, highlighting the complexity of gene regulation
349 in nutrient metabolism and the need to consider multiple aspects of nutrition.

350 Furthermore, to ensure our data's accuracy and reliability, we compared the results obtained with
351 biologically consistent MeSH term lists and randomly aggregated MeSH terms (Figure 5, B). This comparison
352 was crucial in assessing the significance of our findings. Our results demonstrated that using biologically
353 consistent MeSH lists led to meaningful results. The genes and variants identified using these lists were
354 relevant to the specific nutritional context studied. On the other hand, when we adopted randomly
355 aggregated MeSH terms lacking biological cohesion, the results were non-significant and lacked meaningful
356 associations. This comparison highlights the importance of utilizing biologically relevant MeSH lists to retrieve
357 and prioritize genuinely relevant to the chosen biomedical context.

358 However, it is essential to acknowledge the limitations of this approach. One limitation of our approach is
359 the reliance on available literature and databases. The accuracy and reliability of the build resource depend
360 on the quality and completeness of the data retrieved from various sources, as well as the accuracy of the
361 data structuring and integration process. Our method relies on data from Medline studies, which may be
362 subject to publication bias (25). The data quality and consistency of retrieved data heavily depend on the
363 quality of the original studies and the curation process. Despite efforts to ensure data quality, the original
364 studies' inconsistencies, errors, and biases may still be present in the constructed dataset.

365 Moreover, our dataset is limited to the data available in the LitVar database, GWAS-Catalog, and other
366 sources used in our study (26). As a result, it may not encompass all potential genetic polymorphisms
367 associated with nutrition-related traits. Relying on available literature and data collection databases has
368 limitations (27). Despite our efforts to minimize MeSH attribution bias, the dataset could not contain all the
369 relevant literature. Inconsistencies, errors, and biases in the original studies may be transferred to the
370 constructed dataset. Finally, the dataset may cover only some populations and ethnicities, which could limit
371 its applicability to diverse populations with different genetic backgrounds (28).

372 Furthermore, it is essential to acknowledge gene-environment interactions' complex and multifactorial
373 nature, including interactions with dietary factors (29). While our dataset captures a subset of the possible
374 interactions, it may not encompass their full complexity. In interpreting the associations between genetic
375 polymorphisms and nutrition-related traits, it is crucial to consider other factors, such as environmental
376 influences, epigenetic modifications, and gene-gene interactions (30). Therefore, the complexity of gene-
377 environment interactions, including interactions with dietary factors, requires further investigation beyond
378 the scope of this research.

379 Finally, to improve the validity of the study's results, it is essential to assess the quality and scientific validity
380 of the literature sources through established criteria. Future research could follow the scientific validity
381 assessment criteria described by Grimaldi and co-workers (31) to ensure the reliability of individual sources.

382

383 **5. Conclusion**

384 Our study presents a comprehensive approach for building a nutrigenetic dataset by integrating data from
385 multiple sources through the MeSH ontology. In conclusion, the study presents a methodological framework
386 and a valuable resource for exploring genetic polymorphisms associated with nutrition-related traits. The
387 integration and analysis of genetic polymorphism data in the field of nutrition holds great promise for
388 advancing personalized nutrition interventions. This study highlights the importance of utilizing standardized
389 curation processes and comprehensive datasets in nutrigenomics research. The curated dataset presented
390 here fills a gap in the existing literature by providing a consolidated and standardized resource for
391 investigating genetic polymorphisms associated with nutrition-related traits.

392 This study exemplifies the importance of translational bioinformatics in nutrition and nutrigenomics,
393 underscoring the significance of computational approaches in consolidating and analyzing information from
394 multiple sources, thereby facilitating comprehensive research in the field of nutrigenomics.

395 Bibliography

- 396 1. Chen Y, Michalak M, Agellon LB. Importance of Nutrients and Nutrient Metabolism on Human Health.
397 Yale J Biol Med. 2018 Jun 28;91(2):95–103.
- 398 2. KIANI AK, BONETTI G, DONATO K, KAFTALLI J, HERBST KL, STUPPIA L, et al. Polymorphisms, diet and
399 nutrigenomics. J Prev Med Hyg. 2022 Oct 17;63(2 Suppl 3):E125–41.
- 400 3. Loktionov A. Common gene polymorphisms and nutrition: emerging links with pathogenesis of
401 multifactorial chronic diseases (review). J Nutr Biochem. 2003 Aug;14(8):426–51.
- 402 4. Mathers JC. Nutrigenomics in the modern era. Proc Nutr Soc. 2017 Aug;76(3):265–75.
- 403 5. Kaput J. Nutrigenomics research for personalized nutrition and medicine. Curr Opin Biotechnol. 2008
404 Apr;19(2):110–20.
- 405 6. Fenech M, El-Sohemy A, Cahill L, Ferguson LR, French TAC, Tai ES, et al. Nutrigenetics and
406 nutrigenomics: viewpoints on the current status and applications in nutrition research and practice. J
407 Nutr Nutr. 2011;4(2):69–89.
- 408 7. Comerford KB, Pasin G. Gene–Dairy Food Interactions and Health Outcomes: A Review of Nutrigenetic
409 Studies. Nutrients. 2017 Jul 6;9(7):710.
- 410 8. Singh V. Current challenges and future implications of exploiting the omics data into nutrigenetics and
411 nutrigenomics for personalized diagnosis and nutrition-based care. Nutrition. 2023 Jun 1;110:112002.
- 412 9. Allot A, Peng Y, Wei CH, Lee K, Phan L, Lu Z. LitVar: a semantic search engine for linking genomic variant
413 data in PubMed and PMC. Nucleic Acids Res. 2018 Jul 2;46(Web Server issue):W530–6.
- 414 10. MacArthur J, Bowler E, Cerezo M, Gil L, Hall P, Hastings E, et al. The new NHGRI-EBI Catalog of
415 published genome-wide association studies (GWAS Catalog). Nucleic Acids Res. 2017 Jan 4;45(Database
416 issue):D896–901.
- 417 11. Use of MeSH in Cataloging [Internet]. [cited 2023 Aug 3]. Available from: [https://wayback.archive-
418 it.org/org-350/20200228165129/https://www.nlm.nih.gov/tsd/cataloging/MeSH_CatPractices.html](https://wayback.archive-it.org/org-350/20200228165129/https://www.nlm.nih.gov/tsd/cataloging/MeSH_CatPractices.html)
- 419 12. Leydesdorff L, Comins JA, Sorensen AA, Bornmann L, Hellsten I. Cited references and Medical Subject
420 Headings (MeSH) as two different knowledge representations: clustering and mappings at the paper
421 level. Scientometrics. 2016;109(3):2077–91.
- 422 13. Khare R, Leaman R, Lu Z. Accessing Biomedical Literature in the Current Information Landscape.
423 Methods Mol Biol Clifton NJ. 2014;1159:11–31.
- 424 14. Pimentel JF, Murta L, Braganholo V, Freire J. Understanding and improving the quality and
425 reproducibility of Jupyter notebooks. Empir Softw Eng. 2021;26(4):65.
- 426 15. Reddy VS, Palika R, Ismail A, Pullakhandam R, Reddy GB. Nutrigenomics: Opportunities & challenges for
427 public health nutrition. Indian J Med Res. 2018 Nov;148(5):632–41.
- 428 16. O'Rahilly S, Farooqi IS. Human obesity as a heritable disorder of the central control of energy balance.
429 Int J Obes 2005. 2008 Dec;32 Suppl 7:S55–61.
- 430 17. Wang DD, Hu FB. Precision nutrition for prevention and management of type 2 diabetes. Lancet
431 Diabetes Endocrinol. 2018 May;6(5):416–26.

- 432 18. Peña-Romero AC, Navas-Carrillo D, Marín F, Orenes-Piñero E. The future of nutrition: Nutrigenomics
433 and nutrigenetics in obesity and cardiovascular diseases. *Crit Rev Food Sci Nutr*. 2018 Nov
434 22;58(17):3030–41.
- 435 19. Di Camillo B, Giugno R. From translational bioinformatics computational methodologies to personalized
436 medicine. *J Biomed Inform*. 2022 Sep 1;133:104170.
- 437 20. Lee KH, Kim JH. Evolution of Translational Bioinformatics: lessons learned from TBC 2016. *BMC Med*
438 *Genomics*. 2017 May 24;10(Suppl 1):32.
- 439 21. Tenenbaum JD. Translational Bioinformatics: Past, Present, and Future. *Genomics Proteomics*
440 *Bioinformatics*. 2016 Feb;14(1):31–41.
- 441 22. Zhang Y, Zhu Q, Liu H. Next generation informatics for big data in precision medicine era. *BioData Min*.
442 2015;8:34.
- 443 23. Carrasco-Ramiro F, Peiró-Pastor R, Aguado B. Human genomics projects and precision medicine. *Gene*
444 *Ther*. 2017 Sep;24(9):551–61.
- 445 24. Floris M, Cano A, Porru L, Addis R, Cambedda A, Idda ML, et al. Direct-to-Consumer Nutrigenetics
446 Testing: An Overview. *Nutrients*. 2020 Feb 21;12(2):566.
- 447 25. Carter AO, Griffin GH, Carter TP. A survey identified publication bias in the secondary literature. *J Clin*
448 *Epidemiol*. 2006 Mar;59(3):241–5.
- 449 26. Lachance J, Tishkoff SA. SNP ascertainment bias in population genetic analyses: why it is important, and
450 how to correct it. *BioEssays News Rev Mol Cell Dev Biol*. 2013 Sep;35(9):780–6.
- 451 27. Williams WG. Uses and limitations of registry and academic databases. *Semin Thorac Cardiovasc Surg*
452 *Pediatr Card Surg Annu*. 2010;13(1):66–70.
- 453 28. Nilsson PD, Newsome JM, Santos HM, Schiller MR. Prioritization of Variants for Investigation of
454 Genotype-Directed Nutrition in Human Superpopulations. *Int J Mol Sci*. 2019 Jul 18;20(14):3516.
- 455 29. Virolainen SJ, VonHandorf A, Viel KCMF, Weirauch MT, Kottyan LC. Gene-environment interactions and
456 their impact on human health. *Genes Immun*. 2023 Feb;24(1):1–11.
- 457 30. Cole BS, Hall MA, Urbanowicz RJ, Gilbert-Diamond D, Moore JH. Analysis of Gene-Gene Interactions.
458 *Curr Protoc Hum Genet*. 2017 Oct 18;95:1.14.1-1.14.10.
- 459 31. Grimaldi KA, van Ommen B, Ordovas JM, Parnell LD, Mathers JC, Bendik I, et al. Proposed guidelines to
460 evaluate scientific validity and evidence for genotype-based dietary advice. *Genes Nutr*. 2017;12:35.

461 **Declaration of Generative AI and AI-assisted technologies in the** 462 **writing process**

463 During the preparation of this work, the author(s) used the Generative Pre-trained Transformer (GPT)
464 provided by OpenAI API to facilitate the subsetting of the extensive MeSH collection with representatives of
465 a particular biomedical field. After using this service, the authors reviewed and edited the content as needed
466 and take full responsibility for the publication's content.

467

468 **Author contributions**

469 **GMDF:** Writing - original draft, Conceptualization, Data curation, Software, Visualization, Writing – review &
470 editing; **MM:** Writing - original draft, supervision, Investigation, Writing – review & editing; **AP:** Supervision,
471 Writing – review & editing, Validation; **TA:** Supervision, Writing – review & editing, Validation; **BHM:** Writing
472 - original draft, Conceptualization, Formal analysis, Validation, Writing – review & editing; **VC:** Funding
473 acquisition, Conceptualization, Supervision, Writing – review & editing.

474 **Declaration of Competing Interest**

475 The authors declare that they have no known competing financial interests or personal relationships that
476 could have appeared to influence the work reported in this paper.

477 **Acknowledgment**

478 Our research was supported by Federazione Nazionale Degli Ordini dei Biologi.