

1 Large language models identify causal genes in 2 complex trait GWAS

3
4 Authors: Suyash S. Shringarpure^{1,*}, Wei Wang¹, Sotiris Karagounis¹, Xin Wang¹, Anna C.
5 Reisetter², Adam Auton¹, Aly A. Khan^{1,3,4,*}

6
7 Affiliations:

8 ¹ 23andMe Inc., 223 N Mathilda Ave, Sunnyvale, CA 94086, USA

9 ² Therapeutics Division, 23andMe, 349 Oyster Point Blvd, South San Francisco, CA 94080, USA

10 ³ Departments of Pathology, and Family Medicine, University of Chicago, Chicago, IL 60637,
11 USA

12 ⁴ Institute for Population and Precision Health, University of Chicago, Chicago, IL 60637, USA

13

14 * Corresponding authors: sshringarpure@23andme.com; aakhan@uchicago.edu

15 Abstract

16 Identifying underlying causal genes at significant loci from genome-wide association studies
17 (GWAS) remains a challenging task. Literature evidence for disease-gene co-occurrence,
18 whether through automated approaches or human expert annotation, is one way of nominating
19 causal genes at GWAS loci. However, current automated approaches are limited in accuracy
20 and generalizability, and expert annotation is not scalable to hundreds of thousands of
21 significant findings. Here, we demonstrate that large language models (LLMs) can accurately
22 identify genes likely to be causal at loci from GWAS. By evaluating the performance of GPT-3.5
23 and GPT-4 on datasets of GWAS loci with high-confidence causal gene annotations, we show
24 that these models outperform state-of-the-art methods in identifying putative causal genes.
25 These findings highlight the potential of LLMs to augment existing approaches to causal gene
26 discovery.

27 Main

28
29 Genome-wide association studies (GWAS) have identified many regions of the genome
30 associated with complex traits, enhancing our understanding of trait biology. However, a
31 significant limitation of GWAS is the difficulty in pinpointing the underlying causal gene for a
32 given association. Approaches to causal gene identification from GWAS loci use a broad range
33 of information including functional annotation, colocalization with quantitative trait loci (QTL)
34 datasets, biological insights, and literature evidence. Literature mining for the co-occurrence of a
35 (disease, gene) pair in a publication potentially provides evidence for the causal role of the gene
36 in the disease, and may recapitulate the knowledge that an expert biologist or clinician might
37 use to identify the causal gene at a GWAS locus. However, current literature mining approaches

38 (Kafkas, Dunham, and McEntyre 2017; Tirunagari et al. 2024) for causal gene prioritization have
39 been evaluated in limited settings or through related tasks such as drug/gene entity recognition
40 and normalization, and their generalizability to all datasets is unclear.

41
42 Large language models (LLMs) are deep learning models trained on large text corpora,
43 initially to predict masked/next words from a sentence, and then subsequently trained for a large
44 number of tasks including text generation, summarization, and question-answering. Recent
45 studies have demonstrated their capability to perform biomedical tasks (Sarwal et al. 2023),
46 including summarizing gene function (Chen and Zou 2024), medical question answering
47 (Singhal et al. 2023), cell-type annotation (Hou and Ji 2024), and identifying causal genetic
48 factors from murine experimental data (Tu et al. 2023). We hypothesize that large language
49 models like GPT-3.5 (Brown et al. 2020) and GPT-4 (OpenAI et al. 2024) offer a systematic way
50 to mining literature and identifying causal genes at GWAS loci, as their training datasets include
51 scientific literature and other sources of information about genetics. This approach would enable
52 efficient and scalable annotation of likely causal genes at GWAS loci using literature evidence,
53 which is impractical through expert human annotation.

54
55 We performed a systematic evaluation of GPT-3.5 and GPT-4 for the task of causal
56 gene identification, and compared their performance to state-of-the-art computational methods
57 (Supplementary Figure 1). Four evaluation datasets containing 641 to 1692 GWAS loci, with
58 ground-truth annotations of causal genes based on different criteria, were used to test the
59 generalizability of the LLM-based approach (Supplementary Table 1). Given that LLMs are
60 trained on large undisclosed text datasets and it is hard to verify which datasets were used in
61 their training, we selected several well-studied benchmark datasets, including a newly curated
62 dataset created after the training period of GPT-3.5 and GPT-4, as well as a benchmark dataset
63 that is not available on the internet (referred to as GWAS catalog and Weeks et al. respectively,
64 Methods).

65
66 The input prompt to the LLMs contained a generic description of an expert geneticist
67 seeking to identify the causal gene, followed by the name of the GWAS phenotype and a list of
68 all genes within 500 kbp of the lead variant at the locus (Figure 2, Methods). Additionally, the
69 LLM was instructed to output the name of the causal gene, a confidence score between 0 and
70 1, and a short reason for the choice. We queried the LLMs using available APIs to use specific
71 model versions and ensure reproducibility of results (Methods).

72
73 For comparison, we evaluated other state-of-the-art approaches for predicting the causal
74 gene, including the polygenic priority score (PoPS, Weeks et al. 2023), locus-to-gene score
75 (L2G, Mountjoy et al. 2021), OpenTargets text-mining (Tirunagari et al. 2024), and the 'nearest
76 gene' method (Stacey et al. 2019), all applied on the same datasets. We evaluated all methods
77 based on agreement of predictions with the ground-truth annotations of causal genes in the
78 original dataset, using precision, recall, and F-score (the harmonic mean of precision and recall)
79 as performance measures.

80

81 We found that GPT-3.5 is competitive with existing methods, and GPT-4 outperforms
82 existing methods on all four datasets in F-score (Figure 2a and 2b, Supplementary Table 2 and
83 Supplementary Table 8, $p = 8.1e-8, 3.84e-15, 8.72e-4, 1.04e-11$ for a paired Wilcoxon signed-
84 rank test on OpenTargets, Pharmaprojects, Weeks et al., and GWAS catalog respectively, when
85 comparing GPT-4 with the best non-LLM approach). The performance of the LLM-based
86 methods correlates with a number of factors related to the GWAS locus. We observed that
87 prediction accuracy is negatively correlated with the number of genes in the locus, a relationship
88 observed for all methods (Figure 2c, Supplementary Table 2). LLM-based methods showed a
89 positive correlation between prediction accuracy and the number of publications for the causal
90 gene, and a similar correlation was observed in other methods (Figure 2d, Supplementary Table
91 2). Additionally, we also noticed a small number of obvious hallucinations, where the LLMs
92 reported a causal gene that was not included in the set of provided genes at a locus (fewer than
93 2% of all loci for GPT-3.5, fewer than 0.9% for GPT-4, Supplementary Table 3).

94
95 Next we assessed the LLM-based methods' ability to generate a confidence score
96 associated with their prediction. We observed that LLMs were well-calibrated at higher
97 confidence levels (≥ 0.8), but overly optimistic in their confidence estimates at lower
98 confidence levels (0.5-0.7), with GPT-4 showing better calibration than GPT-3.5 (Supplementary
99 Figure 2, Supplementary Table 4). We found that focusing only on the high-confidence
100 predictions from GPT-4 (confidence ≥ 0.8) allowed additional improvements to precision, with
101 improvements from 10% to 43% across the datasets (Supplementary Figure 7).

102
103 We examined the purported reasoning underlying the LLMs' confidence and predictions,
104 which might shed light on their ability to interpret complex prompts. The reasons provided by the
105 LLMs for correct predictions include phrases describing the functions of genes ("is involved in",
106 "the regulation of", "in immune response" etc.), or their association with the phenotype ("is
107 associated with", "is implicated in" etc.) among others (Supplementary Table 9). Since LLMs are
108 sensitive to their input prompts, we tested the impact of prompt structure on our results. In our
109 sensitivity analysis, we found that the LLM-based methods maintained their performance even
110 when provided with only a minimal prompt containing the output format, task instruction,
111 phenotype name, and gene names, without any additional context (Supplementary Figure 3).
112 This suggests that the LLMs internally contained most of the information needed for the causal
113 gene identification task without requiring substantial context.

114
115 To further explore the internal model representation of genetic and phenotypic
116 associations, we examined the embeddings of phenotypes and genes using the "text-
117 embedding-3-large" model (OpenAI 2024). LLMs represent words as points in a high-
118 dimensional embedding space, where similarity in these representations can capture semantic
119 relationships (Mikolov et al. 2013). In the context of causal gene identification, we hypothesized
120 that causal genes are likely to be proximal to the phenotypes they influence in the embedding
121 space. To test this hypothesis, we use pre-computed high-dimensional embeddings of LLM-
122 generated gene and phenotype descriptions at a GWAS locus. The predicted causal gene at the
123 locus is the gene which is most similar to the phenotype in the embedding space. We found that

124 using similarity of pre-computed embeddings alone achieves about 75-90% of the performance
125 of GPT-3.5 (Supplementary Figure 4).

126
127 To illustrate the role of embeddings in prediction, we present t-SNE projections for a
128 locus associated with LDL cholesterol from the Weeks et al. dataset in Figure 3a. For this locus,
129 which has 12 candidates for the causal gene, the text embeddings of the *PCSK9* gene are
130 closest to those of LDL cholesterol (as measured through cosine similarity). Consequently, the
131 LLM-based approaches correctly nominate *PCSK9* as the causal gene at the locus.
132 Supplementary Figure 5 shows a barplot of the gene-phenotype similarities for this locus. Figure
133 3b quantifies the similarity between the causal gene and phenotype in the embedding space for
134 all loci in our evaluation datasets. We observe that the causal gene is most similar to the
135 phenotype for 40-70% of all examples, depending on the evaluation set, among all genes at the
136 locus. Extending this further, we find that the causal gene is among the top 5 most similar genes
137 to the phenotype in the embedding space for 75%-93% of all examples (Supplementary Figure
138 6).

139
140 Although similarity in high-dimensional embedding space explains a large proportion of
141 the performance of LLM-based approaches, we find that LLMs improve on these phenotype
142 embeddings. For instance, the GWAS catalog dataset contains 250 loci for a phenotype related
143 to sex differences that is only described as “Multi-trait sex score” (originally defined as the sum
144 of multiple quantitative traits, weighted by their respective sex-difference effect sizes). The
145 embedding-based approach makes incorrect predictions for most of these loci (precision for
146 phenotype = 0.05) while GPT-4 achieves a precision of 0.65 for the same phenotype
147 (Supplementary Table 5). This suggests that with the additional context of the task description
148 and the gene information, the LLM is correctly able to infer that the short phenotype description
149 refers to sexual dimorphism and sex-specific traits.

150
151 In considering the broad applicability of this approach for causal gene prediction beyond
152 currently published phenotypes, we found a couple of notable failure modes when examining
153 phenotypes where GPT-4 had very low precision. First, GPT-4 had a precision of 0.08 for the
154 “Total protein” phenotype. The reasons provided by the LLM suggest an incomplete
155 understanding of the short phenotype description, interpreting it only in terms of broad protein
156 levels rather than specifically about protein levels in blood (Supplementary Table 6). This
157 indicates that more specific phenotype descriptions could improve the performance of LLMs.

158
159 In a second instance, GPT-4 had a precision of 0.0 for the phenotype “Neonatal
160 circulating Complement Component 4 (C4) protein concentration”. For this phenotype, all the
161 GWAS loci in the evaluation dataset include the *C4A* gene in the 500 kbp window around the
162 lead variant, though they have different causal genes based on coding variant signals. Since the
163 *C4A* gene is a major part of the complement system, the LLM-based approach always predicts
164 *C4A* as the causal gene. This suggests that combining LLMs with functional annotation data
165 could improve causal gene prioritization. Supplementary Table 5 shows the precision and recall
166 for all methods stratified by phenotype.

167

168 While we have demonstrated the potential of LLMs for causal gene identification, it is
169 important to consider the limitations and challenges associated with their use in our study. The
170 lack of disclosure regarding the training datasets for most LLMs makes it challenging to verify
171 whether our evaluation datasets were somehow included in their training. To mitigate this risk,
172 we introduced a study design using data curated after the LLM training period as well as a
173 benchmark dataset not available on the internet. Additionally, while LLMs return plausible
174 reasons along with their predictions, it is challenging to pinpoint the exact information the LLM
175 used to make each prediction. This is a current shortcoming, but it also points to an area of
176 active research in the field. As reasoning capabilities continue to improve in LLMs, we anticipate
177 that their ability to provide transparent and verifiable explanations will evolve rapidly.

178
179 Overall, our study demonstrates for the first time that LLMs can significantly enhance
180 causal gene identification and GWAS interpretation by systematically incorporating literature
181 evidence. LLM-based approaches, requiring only the location of the lead variant, can be applied
182 to any GWAS locus. With improved prompting, these methods could become valuable tools for
183 causal gene identification. Furthermore, combining LLM-based approaches with functional data
184 holds the potential to create even more accurate and robust methods for identifying causal
185 genes, advancing our understanding of genetic contributions to complex traits.

186
187
188

189 References

- 190 Brown, Tom B., Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
191 Arvind Neelakantan, et al. 2020. "Language Models Are Few-Shot Learners." arXiv.
192 <https://doi.org/10.48550/arXiv.2005.14165>.
- 193 Chen, Yiqun, and James Zou. 2024. "GenePT: A Simple But Effective Foundation Model for
194 Genes and Cells Built From ChatGPT." bioRxiv.
195 <https://doi.org/10.1101/2023.10.16.562533>.
- 196 Hou, Wenpin, and Zhicheng Ji. 2024. "Assessing GPT-4 for Cell Type Annotation in Single-Cell
197 RNA-Seq Analysis." *Nature Methods*, March, 1–4. <https://doi.org/10.1038/s41592-024-02235-4>.
- 199 Kafkas, Şenay, Ian Dunham, and Johanna McEntyre. 2017. "Literature Evidence in Open
200 Targets - a Target Validation Platform." *Journal of Biomedical Semantics* 8 (1): 20.
201 <https://doi.org/10.1186/s13326-017-0131-3>.
- 202 Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. "Efficient Estimation of Word
203 Representations in Vector Space." arXiv. <http://arxiv.org/abs/1301.3781>.
- 204 Minikel, Eric Vallabh, Jeffery L. Painter, Coco Chengliang Dong, and Matthew R. Nelson. 2024.
205 "Refining the Impact of Genetic Evidence on Clinical Success." *Nature*, April.
206 <https://doi.org/10.1038/s41586-024-07316-0>.
- 207 Mountjoy, Edward, Ellen M. Schmidt, Miguel Carmona, Jeremy Schwartzenruber, Gareth Peat,
208 Alfredo Miranda, Luca Fumis, et al. 2021. "An Open Approach to Systematically
209 Prioritize Causal Variants and Genes at All Published Human GWAS Trait-Associated
210 Loci." *Nature Genetics* 53 (11): 1527–33. <https://doi.org/10.1038/s41588-021-00945-5>.
- 211 OpenAI. 2024. "New Embedding Models and API Updates." *New Embedding Models and API
212 Updates* (blog). January 2024. <https://openai.com/index/new-embedding-models-and->

- 213 api-updates/
214 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia
215 Leoni Aleman, et al. 2024. "GPT-4 Technical Report." arXiv.
216 <https://doi.org/10.48550/arXiv.2303.08774>.
217 Sarwal, Varuni, Viorel Munteanu, Timur Suhodolschi, Dumitru Ciorba, Eleazar Eskin, Wei
218 Wang, and Serghei Mangul. 2023. "BioLLMBench: A Comprehensive Benchmarking of
219 Large Language Models in Bioinformatics." bioRxiv.
220 <https://doi.org/10.1101/2023.12.19.572483>.
221 Singhal, Karan, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung,
222 Nathan Scales, et al. 2023. "Large Language Models Encode Clinical Knowledge."
223 *Nature* 620 (7972): 172–80. <https://doi.org/10.1038/s41586-023-06291-2>.
224 Stacey, David, Eric B Fauman, Daniel Ziemek, Benjamin B Sun, Eric L Harshfield, Angela M
225 Wood, Adam S Butterworth, Karsten Suhre, and Dirk S Paul. 2019. "ProGeM: A
226 Framework for the Prioritization of Candidate Causal Genes at Molecular Quantitative
227 Trait Loci." *Nucleic Acids Research* 47 (1): e3–e3. <https://doi.org/10.1093/nar/gky837>.
228 Tirunagari, Santosh, Shyamasree Saha, Aravind Venkatesan, Daniel Suveges, Annalisa
229 Buniello, David Ochoa, Johanna McEntyre, Ellen McDonagh, and Melissa Harrison.
230 2024. "Lit-OTAR Framework for Extracting Biological Evidences from Literature."
231 bioRxiv. <https://doi.org/10.1101/2024.03.06.583722>.
232 Tu, Tao, Zhouqing Fang, Zhuanfen Cheng, Svetolik Spasic, Anil Palepu, Konstantina M.
233 Stankovic, Vivek Natarajan, and Gary Peltz. 2023. "Genetic Discovery Enabled by A
234 Large Language Model." bioRxiv. <https://doi.org/10.1101/2023.11.09.566468>.
235 Weeks, Elle M., Jacob C. Ulirsch, Nathan Y. Cheng, Brian L. Trippe, Rebecca S. Fine, Jenkai
236 Miao, Tejal A. Patwardhan, et al. 2023. "Leveraging Polygenic Enrichments of Gene
237 Features to Predict Genes Underlying Complex Traits and Diseases." *Nature Genetics*
238 55 (8): 1267–76. <https://doi.org/10.1038/s41588-023-01443-6>.
239

240 Methods

241

242 Evaluation datasets

243

244 We used four datasets for evaluation (1) the OpenTargets gold-standard dataset, (2) the
245 Citeline Pharmaprojects dataset of drug targets and their approved indications, (3) an evaluation
246 set created by Weeks et al. based on proximity to fine-mapped coding variant associations, and
247 (4) an evaluation set we created using associations added to the GWAS catalog from
248 manuscripts published after April 2023.

249

250 The OpenTargets gold-standard dataset is based on GWAS loci for which there is high
251 confidence (through functional criteria) in the causal gene. We downloaded the OpenTargets
252 gold-standard dataset from <https://github.com/opentargets/genetics-gold-standards/> (file:
253 https://github.com/opentargets/genetics-gold-standards/blob/master/gold_standards/processed/gwas_gold_standards.191108.tsv). We
254 subsetted the dataset to only rows that had a high-confidence annotation (highest_confidence =
255 "High"). We added gene symbol annotations to the dataset with GENCODE release 43, and
256

257 excluded rows where a gene symbol was not found. This resulted in a dataset with 851 rows. To
258 create descriptions of the phenotypes, we combined trait information from standard trait names
259 and reported trait names. We appended reported trait names to standard trait names if the
260 reported trait names are non-redundant, and we also removed irrelevant text from reported trait
261 names, such as “[EA]”, “non-cancer illness code, self-reported”, “GWAS/MetaboChip 2012” and
262 “conditional on rs7709212”.

263
264 The Pharmaprojects dataset contains drug targets and indications, as well as their drug
265 development stage. We used the Pharmaprojects dataset released by Minikel et al. (Minikel et
266 al. 2024) and created additional mappings of disease indications to EFO (Supplementary
267 Notes). We subsetted the dataset to only rows that have a mapped MeSH or EFO term and that
268 correspond to launched drugs (hcat = “Launched”). We added gene symbol annotations to the
269 dataset with GENCODE release 43, and excluded rows where a gene symbol was not found.
270 This resulted in a dataset with 1692 causal gene - phenotype pairs. Since the Pharmaprojects
271 data does not contain GWAS information, we created synthetic GWAS hits near the target gene
272 for each row to mimic real GWAS hits in terms of the distance between the GWAS hit and the
273 underlying causal gene. More details about the methods for creation of the synthetic GWAS hits
274 can be found in the Supplementary Note.

275
276 The Weeks et al. evaluation set contains non-coding credible sets that are within 500 kbp of a
277 high-confidence (posterior inclusion probability > 0.5) fine-mapped coding association in the
278 same GWAS in UK Biobank (Weeks et al. 2023). We obtained this dataset by emailing the
279 authors. Since this dataset is not available directly on the internet, it is the least likely to have
280 been directly used for training the LLMs evaluated in our study (though the underlying GWAS
281 summary statistics and publication could be part of the LLM training data). We used the gene
282 symbol annotations provided by the authors. This resulted in a dataset with 1348 causal gene -
283 phenotype pairs.

284
285 The GWAS catalog dataset contains non-coding lead variants that are within 500 kbp of a
286 coding lead variant in the same GWAS study. We downloaded the GWAS catalog lead variant
287 file version “gwas_catalog_v1.0.2-associations_e111_r2024-03-11.tsv” from
288 <https://www.ebi.ac.uk/gwas/docs/file-downloads> (download date: 03/19/2024). To avoid the
289 possibility of these results being included in LLM training datasets, we subset the data to only
290 include associations added to the GWAS catalog for manuscripts published after April 30, 2023,
291 the latest reported training cutoff date among the LLMs we evaluated in our study. This resulted
292 in a dataset with 641 causal gene - phenotype pairs.

293

294 Generating LLM input from datasets

295
296 For input to the LLMs, we converted each (phenotype, lead variant, causal gene) triplet
297 corresponding to a GWAS locus to a (phenotype, list of genes in locus) pair.

298

299 For OpenTargets, Pharmaprojects, and the GWAS catalog datasets, we identified all genes
300 within 500 kbp of the lead variant (using GENCODE release 43) based on the smallest distance
301 between the gene body and the lead variant. For the Weeks et al. dataset, the dataset provided
302 by the authors included all genes within 500 kbp of the lead variant, so we used that list of
303 genes without any additional processing.

304
305 To avoid leaking information about the lead variant location or causal gene position through the
306 ordering of gene symbols by physical position, we sorted all gene symbols lexicographically
307 before including them in the prompt to the LLMs.

308
309 As an additional sanity check to verify that the LLM relies on phenotype information and not just
310 on gene information (such as location) for predictions, we randomly sampled phenotypes at the
311 evaluation loci and provided those as input to GPT-4. As expected, we found that LLM
312 performance is considerably degraded from an average precision of 63% across all datasets to
313 an average precision of 32%, along with an increase in number of obvious hallucinations from
314 8.25 to 34 (Supplementary Table 10).

315

316 LLM prompts for identifying causal genes

317

318 To query LLMs to identify the causal genes for a (phenotype, list of genes in locus) pair, we
319 used a prompt describing the task for the LLM along with the phenotype and gene list. The
320 basic prompt we used had two components, a system prompt (for general behavior), and a user
321 prompt (pair-specific). For an example pair ("Morning person",[A,B,C,D]), the prompts are
322 described below.

323

324 System prompt:

325 You are an expert in biology and genetics.

326 Your task is to identify likely causal genes within a locus for a given GWAS phenotype based on
327 literature evidence.

328 From the list, provide the likely causal gene (matching one of the given genes), confidence (0:
329 very unsure to 1: very confident), and a brief reason (50 words or less) for your choice.

330 Return your response in JSON format, excluding the GWAS phenotype name and gene list in
331 the locus. JSON keys should be 'causal_gene','confidence','reason'.

332 Your response must start with '{' and end with '}'.

333

334 User prompt:

335 Identify the causal gene.

336 GWAS phenotype: {Morning person}

337 Genes in locus: {A},{B},{C},{D}

338

339 LLMs evaluated

340
341 For our experiments, we evaluated LLMs from OpenAI, GPT-3.5 and GPT-4. For GPT-3.5, we
342 used the model version “gpt-3.5-turbo-0125”, which has been trained on data up to September
343 2021. For GPT-4, we tested the model version “gpt-4-1106-preview”, which has been trained on
344 data up to April 2023, and model version “gpt-4-0613”, which has been trained on data up to
345 September 2021. We found that the “gpt-4-0613” model performed comparably or better than
346 the “gpt-4-1106-preview” model and had fewer obvious hallucinations (Supplementary Table 7),
347 hence we report the results for “gpt-4-0613” as GPT-4 in the main text.

348
349 To make outputs reproducible, we queried all LLMs with temperature set to 0.
350

351 Post-processing of LLM results

352
353 We consider an LLM prediction to be an obvious hallucination if the predicted gene was not in
354 the list of genes at the locus provided to the LLM. These are easily detectable, and we set such
355 predictions to NA, excluding them from downstream evaluation.
356

357 Comparison to other methods

358
359 We compared our results to several state-of-the-art methods for each dataset. First, we
360 evaluated the “nearest gene” method for all datasets. We evaluated text-mining on the
361 OpenTargets and Pharmaprojects dataset, where there was sufficient overlap between the
362 phenotypes in the dataset and phenotypes included in the text-mining scores. We also added
363 prediction based on “locus-to-gene” (L2G) scores to the OpenTargets gold-standard dataset,
364 and prediction based on polygenic priority score (PoPS) to the Weeks et al. evaluation dataset.
365 Both these methods were previously evaluated to have good performance on their respective
366 datasets.

367
368 To obtain the nearest gene prediction, we computed the distance between the genes in each
369 locus and the lead variant (using GENCODE release 43 and reference genome hg38 for the
370 OpenTargets, Pharmaprojects, and GWAS catalog datasets), and defined the nearest gene as
371 the gene with the least distance from the lead variant based on gene body. For the Weeks et al.
372 data, we directly used the nearest gene predictions provided by the authors (downloaded from
373 <https://www.finucanelab.org/data>, file:
374 [https://www.dropbox.com/sh/o6t5jprvxb8b500/AACqCux_jJbF9F56ozhzzk pia/results/UKB_AiIM
375 methods_GenePrioritization.txt.gz?dl=0](https://www.dropbox.com/sh/o6t5jprvxb8b500/AACqCux_jJbF9F56ozhzzk pia/results/UKB_AiIM_methods_GenePrioritization.txt.gz?dl=0)). In the Weeks et al. data, we found that “nearest gene”
376 nominated multiple genes at about 4% of loci (51 of 1348) due to the lead variant position being
377 within multiple gene bodies. To simplify evaluation, we randomly chose a single gene out of all
378 nominated genes as the prediction at such loci. We found that this had a minor impact on

379 precision and recall compared to that reported in the original publication (precision = 0.47 vs
380 previously reported 0.46, recall = 0.47 vs previously reported 0.48).

381
382 To get the text mining prediction, we first downloaded (gene,disease) co-occurrence information
383 and scores (Supplementary Note) from OpenTargets. We aggregated the scores of each
384 gene/disease pair by summing over all scores for the pair from all publications. We defined the
385 predicted causal gene from text mining at a locus as the gene with the largest aggregated score
386 in each locus for the GWAS phenotype.

387
388 The L2G score is a machine learning approach trained using fine-mapped genetics and
389 functional genomics data on 445 gold-standard curated GWAS loci (Mountjoy et al. 2021). The
390 predicted causal gene based on L2G score was defined as the gene with maximal L2G score in
391 each locus (indexed by the phenotype and the lead variant). We downloaded the L2G scores
392 from <https://ftp.ebi.ac.uk/pub/databases/opentargets/genetics/latest/l2g/> (download date:
393 03/26/2024).

394
395 PoPS is a similarity-based gene prioritization method that leverages gene-level summary
396 statistics and incorporates data about genes from a variety of sources to produce a phenotype-
397 gene level prioritization score. We obtained PoPS scores for the 1348 evaluation loci by
398 emailing the authors of the original publication (Weeks et al. 2023). The predicted causal gene
399 based on the PoPS score was defined as the gene with the highest PoPS score in the locus.

400

401 Evaluation of predictions

402
403 For evaluation of predictions, we computed precision (proportion of predicted causal genes that
404 were annotated as causal in the dataset), recall (proportion of annotated causal genes identified
405 among predictions), and F-score (harmonic mean of precision and recall). Some methods
406 (nearest gene, PoPS, LLM-based approaches) made a single causal gene prediction for each
407 example, while others (text mining, L2G score) only made predictions at a subset of all
408 examples. To compute these metrics, the predictions from each method were converted to
409 0/1/NA, with 1 assigned if the method made a prediction and it matched the annotated causal
410 gene, 0 assigned if the method made a prediction but it did not match the annotated causal
411 gene, and NA assigned if the method did not make a prediction for an example. For all methods
412 except text mining, we found that NAs were only a small proportion of the predictions
413 (Supplementary Table 4).

414
415 Precision = (Number of predictions scored 1) / (Number of non-NA predictions)

416
417 Recall = (Number of predictions scored 1) / (Number of examples in dataset)

418
419 F-score = $2 * \text{Recall} * \text{Precision} / (\text{Recall} + \text{Precision})$

420

421 For each metric, we computed 95% confidence intervals using the bootstrap with 1,000
422 samples.

423
424 To compare the performance of GPT-4 to the best non-LLM approach, we used a Wilcoxon
425 signed rank test with continuity correction, and we report p-values for the alternative hypothesis
426 that the GPT-4 performance is higher than the highest non-LLM approach performance. We
427 also report a p-value from a McNemar test for the same comparison with the alternative
428 hypothesis of a difference between performance for the two methods.
429

430 Factors affecting prediction accuracy

431
432 To assess the impact of locus complexity, we used the number of genes in the locus as a
433 measure of its complexity, and computed the Pearson correlation between the number of genes
434 in the locus and whether the prediction was correct.

435
436 To assess the impact of publication count, we used publication counts for causal genes to
437 evaluate correlations with prediction performance. We downloaded the “gene2pubmed” file from
438 <https://ftp.ncbi.nlm.nih.gov/gene/DATA/> (download date: January 26, 2024). We subset the file
439 to only include human genes (tax_id = 9606), and count the number of publications per gene.
440 Since many examples might share the same causal gene, for each causal gene, we computed
441 the proportion of loci predicted correctly. To account for the heavy-tailed distribution of
442 publication count by gene, we computed the Spearman correlation between the number of
443 publications for the causal gene and the proportion of correct predictions at loci with the
444 specified causal gene.

445
446 We used the bootstrap with 1,000 samples to compute 95% confidence intervals for both
447 correlations. We used the same resampled dataset for confidence interval calculations for
448 precision, recall, F-score and both correlations.

449 Calibration analysis

450
451 To assess calibration, we used the confidence scores provided by the LLM. We computed
452 precision for all predictions at a given confidence score, for any confidence score with at least 5
453 predictions. Supplementary Table 4 shows all unique confidence scores predicted by the LLMs,
454 with their counts and precision estimates. We computed a standard error for each precision by
455 assuming it to be the mean of a binomial distribution with size given by the number of
456 predictions with that score. 95% confidence intervals were calculated as precision +/- 1.96*se.
457

458 Summarizing reasons provided by the LLM for correct predictions

459

460 Across all 4 datasets, the LLMs make thousands of correct predictions. To summarize these for
461 easy inspection, we concatenated all reason strings provided by GPT-4 for predictions
462 evaluated to be correct. We then identified trigrams in the concatenated string using the ngram
463 package in R, and sorted them in decreasing order of frequency of occurrence to identify the
464 most common ones. We report the top 200 most frequent trigrams in Supplementary Table 9.
465

466 Sensitivity analysis for prompt structure

467
468 To examine how the prompt affects the prediction accuracy of the LLMs, we experimented with
469 replacing the prompt by a minimal alternative containing only the output format, task description,
470 and the locus information. We conducted these experiments with GPT-3.5 (model version “gpt-
471 3.5-turbo-0125”). The detailed prompt is described below, using the same example locus as
472 earlier.

473
474
475
476 System prompt:
477 From the list, provide the likely causal gene (matching one of the given genes), confidence (0:
478 very unsure to 1: very confident), and a brief reason (50 words or less) for your choice.
479 Return your response in JSON format, excluding the GWAS phenotype name and gene list in
480 the locus. JSON keys should be 'causal_gene','confidence','reason'.
481 Your response must start with '{' and end with '}'.

482
483 User prompt:
484 Identify the causal gene.
485 GWAS phenotype: {Morning person}
486 Genes in locus: {A},{B},{C},{D}
487

488 Embeddings for genes and phenotypes

489
490 We generated descriptions for genes and phenotypes using GPT-3.5 (model version “gpt-3.5-
491 turbo-0125”) by using the prompt below.

492
493 "You are an expert in biology and genetics.
494 Your task is to provide biologically relevant information about the query below in 300 words or
495 less.

496
497 Query: {text}."
498

499 The “text” value was generated by adding the entity type (gene/phenotype) to the beginning of
500 the entity. So the value of text for the *BRCA1* gene would be “gene *BRCA1*”, and that for the
501 breast cancer phenotype would be “phenotype breast cancer”.

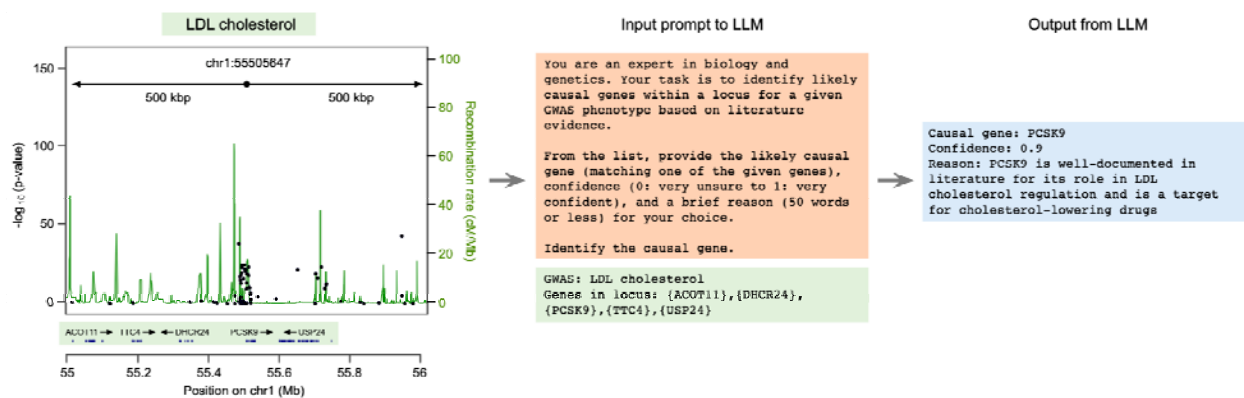
502
503 We provided the generated descriptions as input text to get embeddings using the OpenAI text
504 embedding model “text-embedding-3-large”. This produced embeddings with 3,072 dimensions.
505

506 Embedding-based causal gene prediction

507
508 Using the computed embeddings, we provided the embedding-based causal gene prediction.
509 We calculated the dot product (identical to cosine similarity since the embeddings are
510 normalized to length 1) between the phenotype embedding and the gene embedding, and we
511 defined the predicted causal gene as the gene with the largest dot product of phenotype and
512 gene embeddings among all the genes in each loci.

513
514 To create the t-SNE plot for the *PCSK9* locus, we ran t-SNE on the 13 data points (1 phenotype
515 + 12 genes in the locus) using the Rtsne package with a perplexity of 4.
516

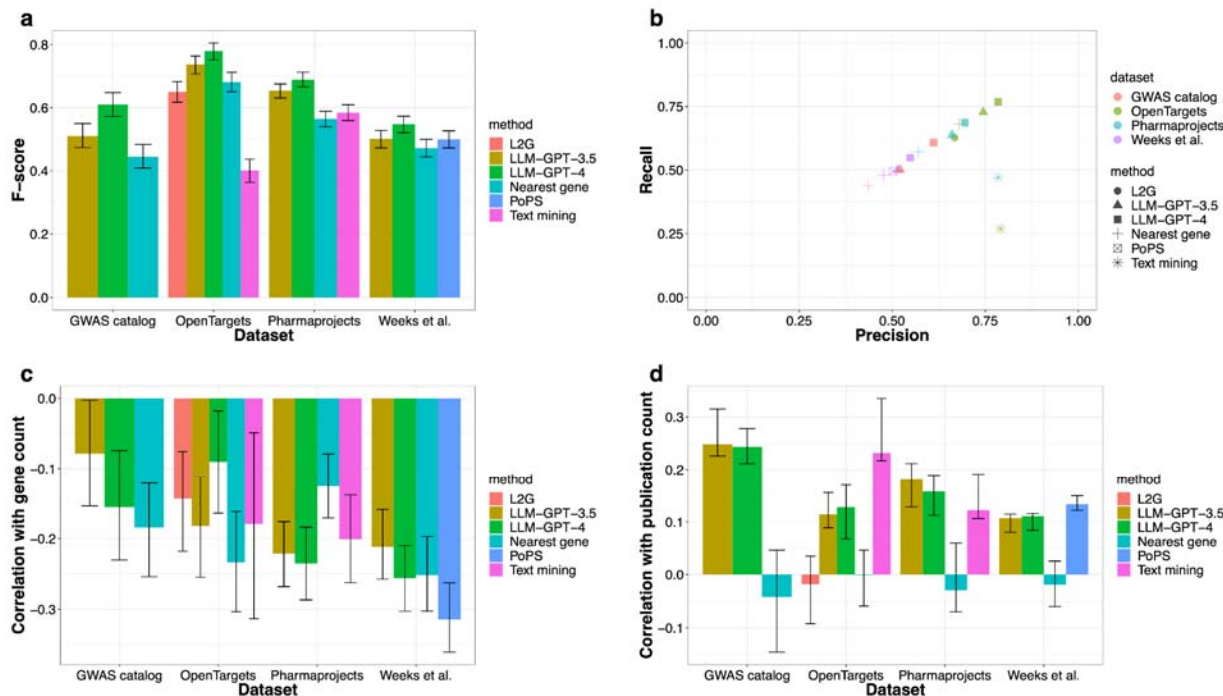
517 Figures



518
519
520 Figure 1: Schematic illustrating the use of LLMs for identifying causal genes at GWAS loci, with
521 an example of an association for LDL cholesterol near the *PCSK9* gene. For a GWAS locus, a
522 500 kbp window is extended on either side of the lead variant (indicated by the black dot), and
523 all genes within this window are considered as candidates. The LLM is then provided with the
524 phenotype name and the alphabetical list of genes and is instructed to predict the causal gene,
525 its confidence, and the reason for its choice. The part of the input prompt colored in orange is
526 the same for all loci, while the part in green is modified for each locus. For readability, output
527 formatting instructions are excluded from the figure.

528
529

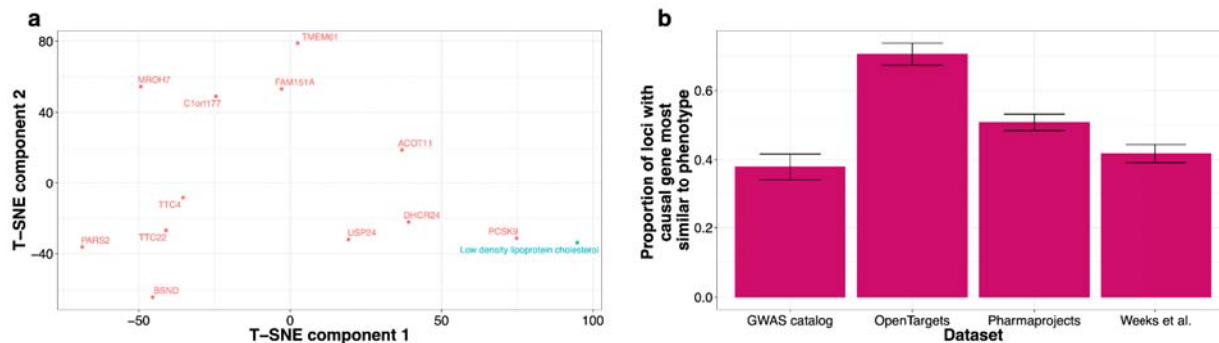
530
531
532
533
534



535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553

Figure 2: Performance of LLMs and other methods on evaluation datasets, and the factors affecting performance

- Performance of all methods on evaluation datasets as measured by F-score (a combination of precision and recall)
- Precision-recall plot showing performance of all methods on each dataset
- Impact of locus complexity on performance for all methods, measured by correlation of prediction accuracy with number of genes at locus. All methods, including LLMs, perform worse at loci with more candidate genes.
- Impact of number of publications for the causal gene on performance for all methods, measured by correlation of prediction accuracy with number of publications for causal gene. Most methods show a positive correlation of performance with how well-studied the causal gene is, except the nearest gene and L2G score methods.



554
555
556
557
558
559
560
561
562
563

Figure 3: Text embeddings of genes and phenotypes partially explain performance of LLMs

- a) T-SNE plot to visualize text embeddings of the genes and phenotype at a locus for LDL cholesterol in the Weeks et al. data. The causal gene *PCSK9* is closest to the LDL cholesterol phenotype in the text embedding space as measured by cosine similarity.
- b) Proportion of examples in the evaluation datasets where the causal gene is most similar to the phenotype in the text embedding space as measured by cosine similarity.

564 Supplementary Tables

- 565
- 566 1 - Information about datasets
- 567 2 - Performance metrics for all methods and datasets
- 568 3 - Number of hallucinations per dataset for the LLM approaches
- 569 4 - Confidence scores predicted by LLMs and precision estimates for each score value
- 570 5 - Precision and recall stratified by phenotype for all methods and datasets
- 571 6 - GPT-4 provided reasons for examples about the phenotype 'Total protein' from the Weeks et
- 572 al. dataset
- 573 7 - Comparison of two GPT-4 version in number of hallucinations and performance metrics
- 574 8 - P-values from paired Wilcoxon's signed rank test and McNemar test comparing performance
- 575 of the LLM-based approach to the best non-LLM approach.
- 576 9 - Trigrams from the reasons reported by GPT-4 for correct predictions, along with their counts,
- 577 for the 200 most frequent trigrams
- 578 10 - Performance of LLMs on datasets with scrambled phenotypes
- 579
- 580
- 581

582 Acknowledgements

583 We would like to thank the employees of 23andMe for making this work possible. We would like
584 to acknowledge Elle Weeks, Hilary Finucane and Jacob Ulirsch for sharing the evaluation
585 dataset from their publication. We would also like to acknowledge David Hinds, Steve Pitts,

586 Bertram Koelsch, Michael Holmes, Stella Aslibekyan, Nick Eriksson for comments on the
587 manuscript.

588

589 The following members of the 23andMe Research Team contributed to this study:
590 Stella Aslibekyan, Adam Auton, Elizabeth Babalola, Robert K. Bell, Jessica Bielenberg, Ninad
591 S. Chaudhary, Zayn Cochinwala, Sayantan Das, Emily DelloRusso, Payam Dibaeinia, Sarah L.
592 Elson, Nicholas Eriksson, Chris Eijsbouts, Teresa Filshtein, Pierre Fontanillas, Davide Foletti,
593 Will Freyman, Zach Fuller, Julie M. Granka, Chris German, Éadaoin Harney, Alejandro
594 Hernandez, Barry Hicks, David A. Hinds, M. Reza Jabalameli, Ethan M. Jewett, Yunxuan Jiang,
595 Sotiris Karagounis, Lucy Kaufmann, Matt Kmiecik, Katelyn Kukar, Alan Kwong, Keng-Han Lin,
596 Yanyu Liang, Bianca A. Llamas, Aly Khan, Steven J. Micheletti, Matthew H. McIntyre, Meghan
597 E. Moreno, Priyanka Nandakumar, Dominique T. Nguyen, Jared O'Connell, Steve Pitts, G.
598 David Poznik, Alexandra Reynoso, Shubham Saini, Morgan Schumacher, Leah Selcer, Anjali J.
599 Shastri, Jingchunzi Shi, Suyash Shringarpure, Keaton Stagaman, Teague Sterling, Qiaojuan
600 Jane Su, Joyce Y. Tung, Susana A. Tat, Vinh Tran, Xin Wang, Wei Wang, Catherine H.
601 Weldon, Amy L. Williams, Peter Wilton.

602

603 S.S., W.W., S.K., X.W., A.R., A.A., A.K., are employed by and hold stock or stock options in
604 23andMe, Inc.

605

606 Data availability

607 We will share all processed datasets used in our analysis, as well as the prediction results from
608 all methods on all datasets, intermediate outputs like gene and phenotype embeddings using
609 Zenodo (doi: 10.5281/zenodo.11391053).

610

611 All source data were openly available. Download links:

612 1) OpenTargets - <https://github.com/opentargets/genetics-gold-standards/>

613 2) Pharmaprojects - https://github.com/ericminikel/genetic_support

614 3) Weeks et al. - <https://www.finuanelab.org/data>

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

616 Code availability

617 We will share the scripts we used to query the LLM, as well as the scripts we use to compute
618 our evaluation metrics using Zenodo (doi: 10.5281/zenodo.11391053). All prompts we used are
619 included in the manuscript or supplementary materials.