

1 Gene annotation bias impedes biomedical research

2 Winston A. Haynes^{1,2,3}, Aurelie Tomczak^{1,2}, and Purvesh Khatri^{1,2,*}

¹ Stanford Institute for Immunity, Transplantation, and Infection,
Stanford University, Stanford, California, USA

²Stanford Center for Biomedical Informatics Research, Department of Medicine,
Stanford University, Stanford, California, USA

³Biomedical Informatics Training Program,
Stanford University, Stanford, California, USA

*To whom correspondence should be addressed; E-mail: pkhatri@stanford.edu

3 **1 Abstract**

4 We found tremendous inequality across gene and protein annotation resources. We observe that
5 this bias leads biomedical researchers to focus on richly annotated genes instead of those with
6 the strongest molecular data. We advocate for researchers to reduce these biases by pursuing
7 data-driven hypotheses.

8 **2 Introduction**

9 After analyzing samples with a high throughput technology, the de facto first step is to perform
10 pathway or network analysis to identify biological processes that are statistically enriched in the
11 data.¹ Researchers typically form hypotheses for their follow up experiments based on the genes
12 or proteins involved in the enriched processes. Commonly used resources for identifying gene
13 functions and interactions include the Gene Ontology (GO),² Reactome,³ Comparative Toxi-

14 cogenomics Database (CTD),⁴ DrugBank,⁵ Protein Data Bank (PDB),⁶ Pubpular,⁷ and NCBI
15 GeneRIF. Since these resources are created by curation of the scientific literature, they typi-
16 cally only contain functional annotations for genes with published experimental data. Although
17 GO includes predicted functional annotations for genes, they are considered of low quality.⁸
18 Consequently, researchers select those genes or proteins for further validation that have prior
19 experimental evidence, which, in turn, leads to more functional annotations for those genes at
20 the expense of under-studied genes.

21 We hypothesized that this experimental paradigm has led to a gene-centric disease research
22 bias where hypotheses are confounded by the streetlight effect of looking for “answers where
23 the light is better rather than where the truth is more likely to lie”.⁹⁻¹¹ To test this hypothesis,
24 we examined the annotation inequality for the human genome across a number of biomedical
25 databases using gini coefficient, which is a measure of inequality such that high coefficient
26 value indicates higher inequality.¹²

27 **3 Results**

28 **3.1 Gene annotation inequality persists across databases**

29 Despite the tremendous growth of GO from 20,826 annotations in 2004 for 7,524 human
30 genes to 122,926 annotations for 16,173 genes in 2017, annotation inequality in GO has in-
31 creased from a gini coefficient of 0.34 in 2004 to 0.50 in 2017. The growth in inequality
32 over time validates that genes with existing annotations continue to receive even more anno-
33 tations. Pathway databases, including Reactome (gini=0.33)³ and the CTD (gini=0.47),⁴ have
34 a similarly high level of inequality. Indeed, every gene annotation resource we examined dis-
35 played a similarly high level of annotation inequality, including: CTD chemical-gene associa-
36 tions (gini=0.63);⁴ PDB 3D protein structures (gini=0.68);⁶ DrugBank drug-gene associations
37 (gini=0.70);⁵ GeneRIF gene publication annotations (gini=0.79); and Pubpular disease-gene

38 publication associations (gini=0.82).^{7,13} We calculated global gene annotation gini coefficient
39 of 0.63 when considering the number of annotations pooled across all these databases. When
40 comparing annotation inequality in gene resources to income inequality in the world, we ob-
41 served that the inequality index for many of the gene resources is higher than any nation in the
42 Organisation for Economic Co-operation and Development (OECD).¹⁴

43 **3.2 Annotation inequality bias affects biomedical research**

44 Next, we explored whether disease research may be affected by the inequality in gene annotation
45 databases. General concern that most published findings are false,¹⁵ many results are inflated,¹⁶
46 and research funding is being wasted^{17,18} has led to searches for solutions that will yield re-
47 producible and clinically relevant findings. Using a multi-cohort analysis framework,^{19,20} we
48 have repeatedly demonstrated that it can identify novel disease-gene relationships that lay out-
49 side the "halo of the streetlight", and which have diagnostic, prognostic, and therapeutic utility
50 across diverse diseases including cancer,²¹⁻²³ organ transplantation,¹⁹ infectious diseases,²⁴⁻²⁷
51 and autoimmunity.²⁸

52 In our manually curated meta-analyses of 104 distinct human conditions, we have inte-
53 grated transcriptome data from over 41,000 patients and 619 studies to calculate an effect size
54 for disease-gene associations.²⁰ Our meta-analyses covered diverse classes of human condi-
55 tions, such as cancer, autoimmune disease, viral infection, neurodegenerative and psychiatric
56 disorders, pregnancy, and obesity. For these conditions, we extracted all disease gene asso-
57 ciations with at least ten publications.^{7,13} Published disease-gene associations exhibited no
58 significant correlation with differential gene expression false discovery rate (FDR) rank [Figure
59 2A, Spearman's correlation= -0.005, p = 0.716]. Overall, only 19.5% of published disease-gene
60 associations were identified in gene expression meta-analyses at a FDR of 5% [Figure S1a].
61 This result is consistent with previous publications that have successfully replicated between

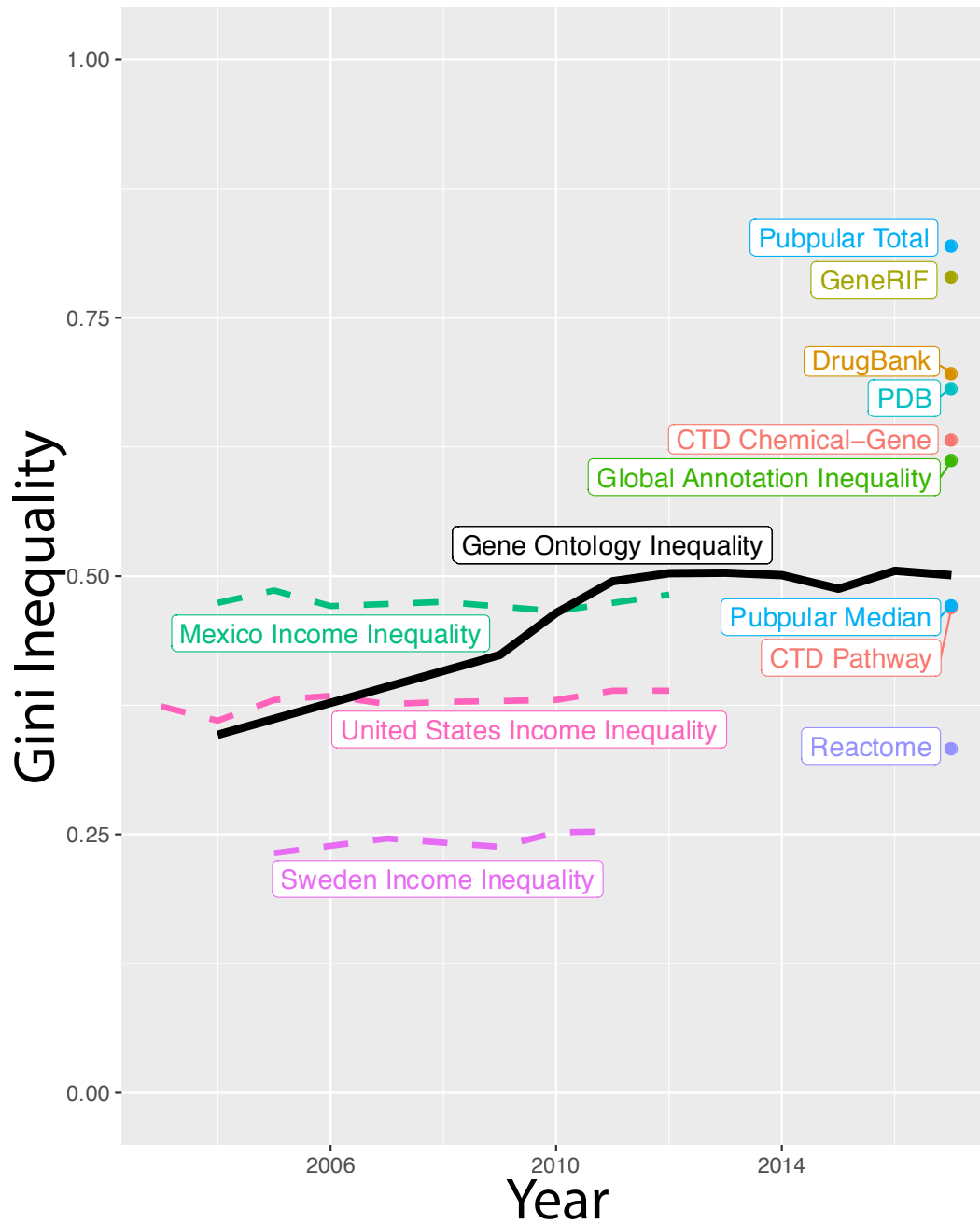


Figure 1: **Inequality in gene annotations.** We measured the gini coefficient across a variety of gene annotation resources. For comparison, we also displayed gini coefficients for income inequality across a sample of OECD nations.

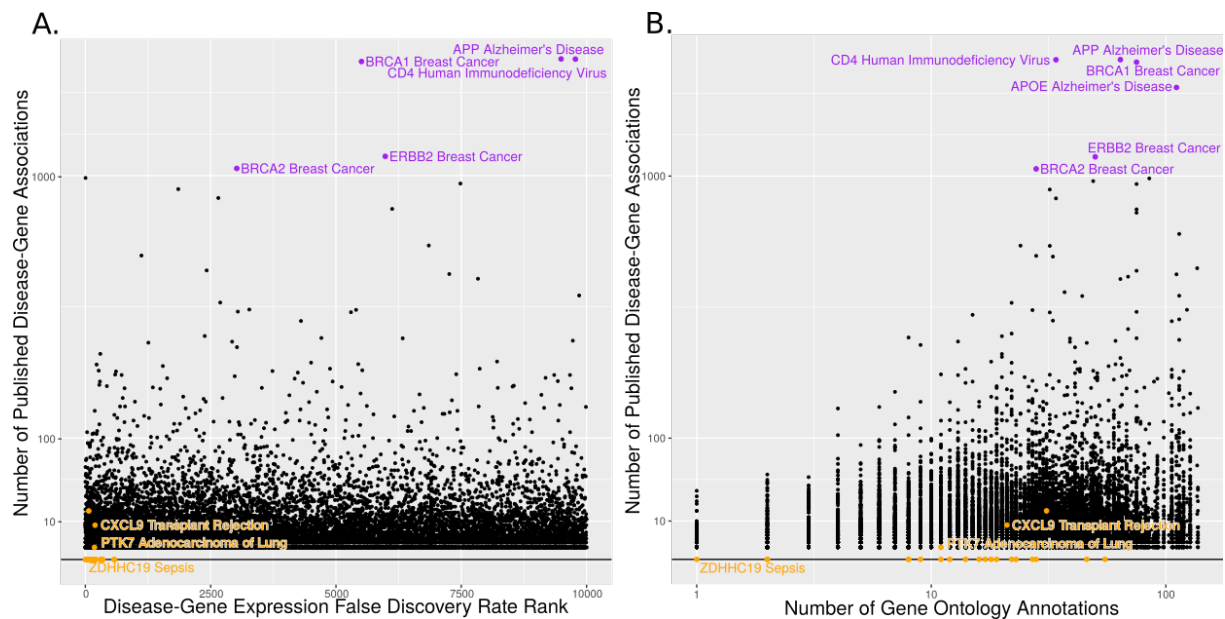


Figure 2: Published Disease-Gene Associations Not Reflected in Molecular Data. (A) The number of publications for every disease-gene pair was not significantly correlated with the gene expression meta-analysis effect size FDR rank [Spearman's correlation = -0.005, $p = 0.716$]. (B) The number of publications for every disease-gene pair correlated with the number of non-inferred from electronic annotation (non-IEA) Gene Ontology annotations [Spearman's correlation = 0.100, $p=8.7e-13$]. Orange points represent disease-gene associations published in our prior meta-analyses.^{19,23,24} Purple points have at least 1000 publications. See also Figure S1.

62 11% – 25% of research studies.^{29,30}

63 To observe whether this phenomenon was specific to gene expression, we extracted genome
64 wide significant single nucleotide polymorphisms (SNPs) from the GWAS catalog.³¹ We ob-
65 served a nominally significant correlation between the number of publications and SNP p-
66 values, indicating moderate concordance between genetic mutations and disease-gene publi-
67 cations [Figure S1b, Spearman's correlation = -0.127, $p = 0.015$].

68 Based on these results, we hypothesized that the lack of correlation with molecular evidence
69 may have been an artifact of research bias towards well-characterized genes. Therefore, we ex-
70 amined correspondence between publications about a disease-gene pair and existing knowledge
71 about that gene as indicated by the number of GO annotations. Indeed, the number of GO
72 annotations for a gene of interest was significantly correlated with the published disease-gene
73 associations [Figure 2B, Spearman's correlation = 0.100, $p=8.7e-13$], but not with gene expres-
74 sion effect size FDR rank in disease [Figure S1c, Spearman's correlation = -0.010, $p = 0.136$].²

75 Many of the highly published disease-gene associations may have been studied for reasons
76 that would not be directly reflected in gene expression analysis, including BRCA1 in breast
77 cancer and CD4 in human immunodeficiency virus. The more troubling bias occurs when
78 associations with strong molecular evidence have no publication record. Disease-gene associ-
79 ations we have reported in our published meta-analyses were typically novel findings with few
80 Gene Ontology annotations, despite having extremely low false discovery rates and high effect
81 sizes^{19,21,24} [orange points in Figure 2]. We observed similar patterns when we performed the
82 same analysis on similar publication and GWAS data from HuGE Navigator^{32,33} [Figures S1d,
83 S1e, S1f].

84 **4 Discussion**

85 Collectively, our results provide evidence of a strong research bias in literature that focuses on
86 well annotated genes instead of the genes with the most significant disease relationship in terms
87 of both gene expression and genetic variation. While focusing research on the best character-
88 ized genes may be natural because it is easy to formulate a mechanistic hypothesis of the gene's
89 function in disease, we propose that omics-era researchers should instead allow data to drive
90 their hypotheses. Our prior work shows that expanding research outside of the streetlight of
91 well characterized genes identifies novel disease-gene relationships, leads to successful repur-
92 posings of drugs, and provides clinically actionable diagnostics.^{19,21–27,34} To enable researchers
93 to pursue data-driven hypotheses, we have made our gene expression meta-analysis data pub-
94 licly available at [<http://metasignature.stanford.edu>] where it may be explored based on either
95 diseases or genes of interest. By focusing on genes with the strongest molecular evidence in-
96 stead of the most annotations, researchers will break the self-perpetuating annotation inequality
97 cycle that results in research bias.

98 **5 Materials and Methods**

99 **5.1 Gini coefficient calculation.**

100 We calculated the gini coefficients using the R package `ineq`.³⁵ We included all human genes
101 with at least one annotation in the gini calculations. We used the Entrez Gene list downloaded in
102 February 2017 of 20,698 current, protein-coding, human genes as our source of human genes.

103 We calculated the number of annotations for each human gene in the Gene Ontology.² We
104 only considered the biological process and molecular function categories and excluded terms
105 with evidence codes IEA and ND. Duplicate annotations that only differ in evidence codes
106 were counted once. We calculated the number of annotations in the January 2004 release and

107 annually for the January 2009-2017 releases.

108 We manually downloaded gene-publication data in August 2016 from Pubpular for 102 of
109 the diseases in our gene expression database.^{7,13} "Pubpular Total" refers to the inequality of
110 gene-publication data across all diseases. "Pubpular Median" refers to the median inequality of
111 gene-publication for each disease.

112 We downloaded Reactome pathway data from the complete database release 59.³ We
113 downloaded data in MySQL format and parsed pathways into UniProt identifiers using custom
114 scripts. We converted UniProt identifiers to gene names using the UniProt identifier conversion
115 tool.³⁶ We calculated the number of pathways including each gene name.

116 We downloaded the CTD⁴ data in February 2017, with the chemical-gene associations and
117 the gene-pathway associations. We calculated the number of chemical-gene and gene-pathway
118 associations for each gene name.

119 We downloaded GeneRIFs from the NCBI in February 2017. We included all human
120 GeneRIFS (Tax ID: 9606). We calculated the number of GeneRIFs for each gene.

121 We downloaded the gene names associated with protein structures from the Protein Data
122 Bank⁶ in February 2017 and calculated the number of structures per gene name.

123 We downloaded the DrugBank⁵ database version 5.0.5 and identified all drugs with known
124 activities on human genes. We calculated the number of drugs targeting each gene.

125 We downloaded OECD¹⁴ nation income inequality gini coefficients from the July 2016 data
126 release at [http://www.oecd.org/social/income-distribution-database.](http://www.oecd.org/social/income-distribution-database.htm)
127 [htm.](http://www.oecd.org/social/income-distribution-database.htm)

128 The code and data we used to run this analysis is available at [http://khatrilab.](http://khatrilab.stanford.edu/researchbias)
129 [stanford.edu/researchbias.](http://khatrilab.stanford.edu/researchbias)

130 **5.2 Gene expression data collection and meta-analysis.**

131 Gene expression meta-analysis data was compiled from the MetaSignature database.²⁰ MetaSig-
132 nature includes data from manual meta-analysis of over 41,000 samples, 619 studies, and 104
133 diseases. Briefly, relevant data were downloaded from Gene Expression Omnibus and Array-
134 Express.^{37,38} Cases and controls were manually labeled for each disease and meta-analysis was
135 performed using the MetaIntegrator package.²⁰ We used the Hedges' g summary effect size,
136 standard error, and false discovery rate which the MetaIntegrator package calculates for every
137 gene.

138 **5.3 Data collection for disease-gene publications and SNP data.**

139 We downloaded the number of publications for each disease-gene relationship from PubPular
140 and HuGE Navigator in August 2016 for as many of the 104 disease in MetaSignature as were
141 present in the databases (102 in PubPular and 81 in HuGE).^{7,13,32} PubPular gave the top 261
142 gene associations, and HuGE gave all known associations. For all correlations, we only consid-
143 ered disease-gene associations with at least 10 publications to limit false positive associations.

144 We downloaded disease-SNP relationships, including gene mappings, odds ratios, and p-
145 values, from the GWAS Catalog and HuGE Navigator for 61 and 54, respectively, of the 103
146 diseases in MetaSignature.^{31,33} From Gene Ontology, we calculated the counts of non-Inferred
147 from Electronic Annotation annotations for all the genes in the MetaSignature database.² The
148 Spearman rank correlation was used for all correlations.

149 Our plots show the top 10,000 gene associations for each disease by effect size FDR rank.
150 Correlation calculations do not include a similar limit.

151 **6 Acknowledgements**

152 We thank Paul J. Utz for feedback about the manuscript and figures and Alex Schrenchuk for
153 computer support. WAH is funded by the National Science Foundation Graduate Research
154 Fellowship under Grant No. DGE-114747. PK is funded by the the Bill and Melinda Gates
155 Foundation, and NIAID grants 1U19AI109662, U19AI057229 and U54I117925.

156 **7 Author Contributions**

157 Conceptualization, W.A.H., A.T., and P.K.; Methodology, W.A.H., A.T., and P.K.; Software,
158 W.A.H. and A.T.; Investigation, W.A.H. and A.T.; Data Curation, W.A.H. and A.T.; Writing-
159 Original Draft, W.A.H. and P.K.; Writing- Reviewing and Editing, W.A.H., A.T., and P.K.;
160 Visualization, W.H.; Funding Acquisition, P.K.

161 **References and Notes**

- 162 ¹ Khatri, P., Sirota, M. & Butte, A. J. Ten years of pathway analysis: current
163 approaches and outstanding challenges. *PLoS computational biology* **8**, e1002375
164 (2012). URL <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002375>\#pcbi-1002375-g003.
165
- 166 ² Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene Ontology
167 Consortium. *Nature genetics* **25**, 25–9 (2000). URL <http://dx.doi.org/10.1038/75556>.
168
- 169 ³ Croft, D. *et al.* The Reactome pathway knowledgebase. *Nucleic acids research* **42**, D472–
170 7 (2014). URL [http://nar.oxfordjournals.org/content/42/D1/D472.](http://nar.oxfordjournals.org/content/42/D1/D472.abstract)
171 abstract.
- 172 ⁴ Davis, A. P. *et al.* The Comparative Toxicogenomics Database’s 10th year anniversary:
173 update 2015. *Nucleic acids research* **43**, D914–20 (2015). URL <http://nar.oxfordjournals.org/content/43/D1/D914.short>.
174
- 175 ⁵ Wishart, D. S. *et al.* DrugBank: a comprehensive resource for in silico
176 drug discovery and exploration. *Nucleic acids research* **34**, D668–72 (2006).
177 URL <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1347430&tool=pmcentrez&rendertype=abstract>.
178
- 179 ⁶ Berman, H. M. *et al.* The Protein Data Bank. *Nucleic acids research* **28**, 235–
180 42 (2000). URL <http://www.ncbi.nlm.nih.gov/pubmed/10592235><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC102472>.
181

- 219 ²² Mazur, P. K. *et al.* Combined inhibition of BET family proteins and histone deacetylases as a
220 potential epigenetics-based therapy for pancreatic ductal adenocarcinoma. *Nature Medicine*
221 **21**, 1163–1171 (2015). URL [http://www.nature.com/doifinder/10.1038/](http://www.nature.com/doifinder/10.1038/nm.3952)
222 [nm.3952](http://www.nature.com/doifinder/10.1038/nm.3952).
- 223 ²³ Chen, R. *et al.* A meta-analysis of lung cancer gene expression identifies PTK7 as a survival
224 gene in lung adenocarcinoma. *Cancer Research* **74**, 2892–2902 (2014). URL [http://](http://www.ncbi.nlm.nih.gov/pubmed/24654231)
225 www.ncbi.nlm.nih.gov/pubmed/24654231.
- 226 ²⁴ Sweeney, T. E., Shidham, A., Wong, H. R. & Khatri, P. A comprehensive time-course-based
227 multicohort analysis of sepsis and sterile inflammation reveals a robust diagnostic gene set.
228 *Science Translational Medicine* **7**, 287ra71 (2015). URL [http://stm.sciencemag.](http://stm.sciencemag.org/content/7/287/287ra71.abstract)
229 [org/content/7/287/287ra71.abstract](http://stm.sciencemag.org/content/7/287/287ra71.abstract).
- 230 ²⁵ Andres-Terre, M. *et al.* Integrated, Multi-cohort Analysis Identifies Conserved Transcrip-
231 tional Signatures across Multiple Respiratory Viruses. *Immunity* **43**, 1199–1211 (2015). URL
232 <http://www.cell.com/article/S1074761315004550/fulltext>.
- 233 ²⁶ Sweeney, T. E., Braviak, L., Tato, C. M. & Khatri, P. Genome-wide expression for diagnosis
234 of pulmonary tuberculosis: a multicohort analysis. *The Lancet Respiratory Medicine* **4**, 213–
235 224 (2016).
- 236 ²⁷ Sweeney, T. E., Wong, H. R. & Khatri, P. Robust classification of bacterial and viral infections
237 via integrated host gene expression diagnostics. *Science translational medicine* **8**, 346ra91
238 (2016). URL <http://www.ncbi.nlm.nih.gov/pubmed/27384347>.
- 239 ²⁸ Lofgren, S. *et al.* Integrated, multicohort analysis of systemic sclerosis identifies robust tran-
240 scriptional signature of disease severity. *JCI Insight* **1** (2016). URL [https://insight.](https://insight.jci.org/articles/view/89073)
241 [jci.org/articles/view/89073](https://insight.jci.org/articles/view/89073).
- 242 ²⁹ Prinz, F., Schlange, T. & Asadullah, K. Believe it or not: how much can we rely on published
243 data on potential drug targets? *Nature Reviews Drug Discovery* **10**, 712–712 (2011). URL
244 <http://www.nature.com/doifinder/10.1038/nrd3439-cl>.
- 245 ³⁰ Begley, C. G. & Ellis, L. M. Drug development: Raise standards for preclinical can-
246 cer research. *Nature* **483**, 531–3 (2012). URL [http://www.nature.com/nature/](http://www.nature.com/nature/journal/v483/n7391/full/483531a.html#t1)
247 [journal/v483/n7391/full/483531a.html#t1](http://www.nature.com/nature/journal/v483/n7391/full/483531a.html#t1).
- 248 ³¹ Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait asso-
249 ciations. *Nucleic acids research* **42**, D1001–6 (2014). URL [http://www.ncbi.](http://www.ncbi.nlm.nih.gov/pubmed/24316577)
250 [nlm.nih.gov/pubmed/24316577](http://www.ncbi.nlm.nih.gov/pubmed/24316577)[http://www.pubmedcentral.nih.gov/](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3965119)
251 [articlerender.fcgi?artid=PMC3965119](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3965119).
- 252 ³² Yu, W., Clyne, M., Khoury, M. J. & Gwinn, M. Phenopedia and Genopedia: disease-centered
253 and gene-centered views of the evolving knowledge of human genetic associations. *Bioin-*
254 *formatics* **26**, 145–146 (2010). URL [http://bioinformatics.oxfordjournals.](http://bioinformatics.oxfordjournals.org/cgi/doi/10.1093/bioinformatics/btp618)
255 [org/cgi/doi/10.1093/bioinformatics/btp618](http://bioinformatics.oxfordjournals.org/cgi/doi/10.1093/bioinformatics/btp618).
- 256 ³³ Yu, W. *et al.* GWAS Integrator: a bioinformatics tool to explore human genetic associations
257 reported in published genome-wide association studies. *European Journal of Human Genet-*
258 *ics* **19**, 1095–1099 (2011). URL [http://www.nature.com/doifinder/10.1038/](http://www.nature.com/doifinder/10.1038/ejhg.2011.91)
259 [ejhg.2011.91](http://www.nature.com/doifinder/10.1038/ejhg.2011.91).

- 260 ³⁴ Li, M. D., Burns, T. C., Morgan, A. A. & Khatri, P. Integrated multi-cohort transcriptional
261 meta-analysis of neurodegenerative diseases. *Acta neuropathologica communications* **2**, 93
262 (2014). URL [http://www.pubmedcentral.nih.gov/articlerender.fcgi?](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4167139&tool=pmcentrez&rendertype=abstract)
263 [artid=4167139&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4167139&tool=pmcentrez&rendertype=abstract).
- 264 ³⁵ Zeileis, A. *ineq: Measuring Inequality, Concentration, and Poverty* (2014). URL [https:](https://cran.r-project.org/package=ineq)
265 [//cran.r-project.org/package=ineq](https://cran.r-project.org/package=ineq).
- 266 ³⁶ UniProt: the universal protein knowledgebase. *Nucleic Acids Research* **45**, D158–D169
267 (2017). URL [https://academic.oup.com/nar/article-lookup/doi/10.](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkw1099)
268 [1093/nar/gkw1099](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkw1099).
- 269 ³⁷ Brazma, A. *et al.* ArrayExpress—a public repository for microarray gene ex-
270 pression data at the EBI. *Nucleic Acids Research* **31**, 68–71 (2003). URL
271 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=165538&tool=pmcentrez&rendertype=abstract)
272 [165538&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=165538&tool=pmcentrez&rendertype=abstract).
- 273 ³⁸ Edgar, R. Gene Expression Omnibus: NCBI gene expression and hybridization array
274 data repository. *Nucleic Acids Research* **30**, 207–210 (2002). URL [http://nar.](http://nar.oxfordjournals.org/content/30/1/207.short)
275 [oxfordjournals.org/content/30/1/207.short](http://nar.oxfordjournals.org/content/30/1/207.short).

276 **8 Supplementary Materials**

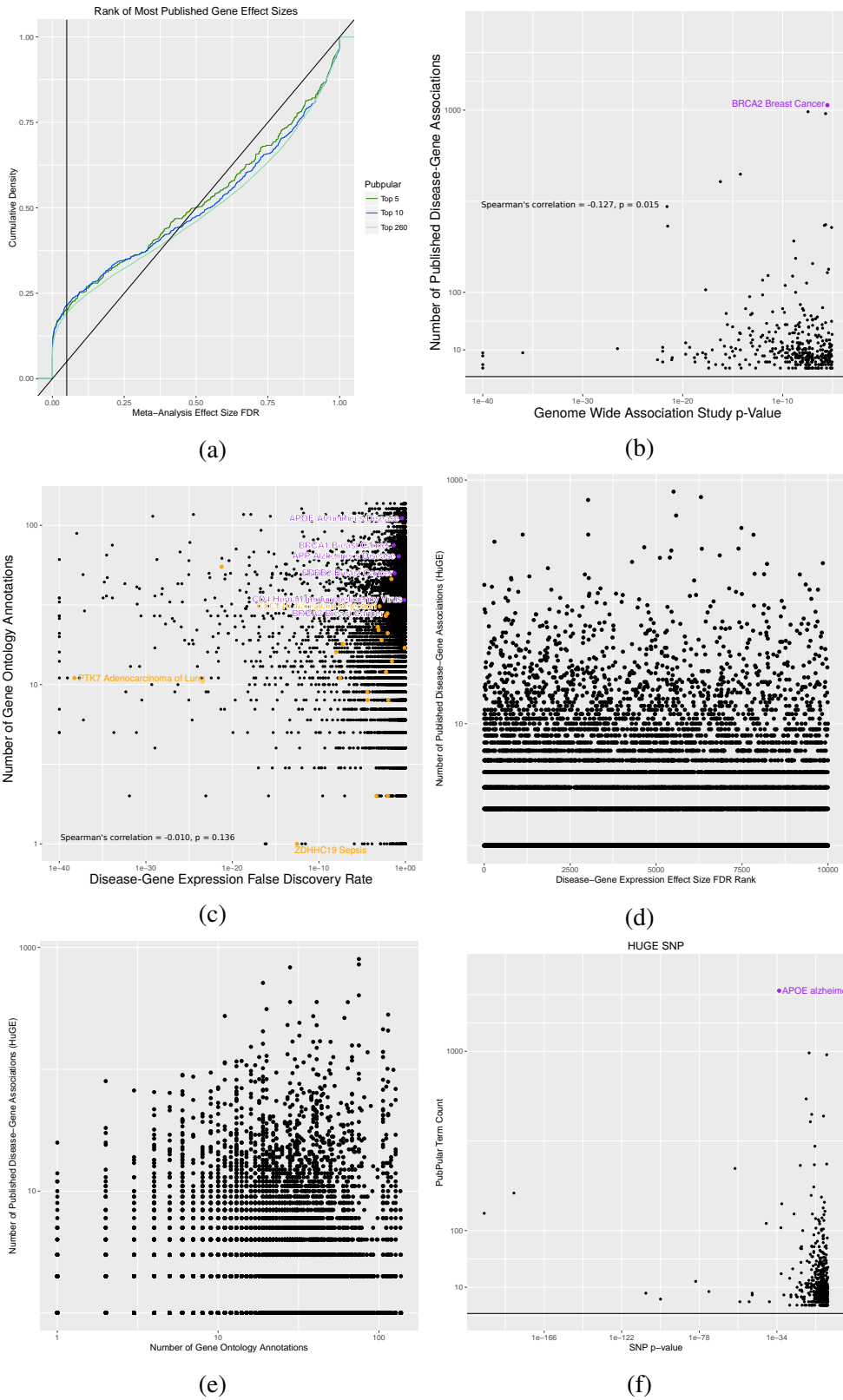


Figure S1: (caption on next page)

Figure S1: Related to Figure 2. (a) Only 20% of published disease-gene associations have gene effect size FDR of less than 5%. Cumulative distributions of the top 5, 10, and 250 disease-gene associations for each disease from PubPular database. Vertical line at a gene expression meta-analysis effect size FDR of 5%. (b) The number of publications for every disease-gene pair is nominally significantly correlated with published results from SNP GWAS from the GWAS catalog [Spearman's correlation = -0.127, $p = 0.015$]. (c) The number of gene ontology annotations for every gene is not correlated with the gene expression meta-analysis effect size false discovery rate (FDR) rank [Spearman's correlation = -0.010, $p = 0.156$]. (d) The number of publications for every disease-gene pair vs. the gene expression meta-analysis effect size FDR rank based on the HuGE Navigator data [Spearman's correlation = 0.032, $p = 0.169$]. (e) The number of publications for every disease-gene pair vs. the number of non-inferred from electronic annotation (non-IEA) Gene Ontology annotations based on the HuGE Navigator data [Spearman's correlation = 0.211, $p = 2.2e-16$]. (f) The number of publications for every disease-gene pair vs. SNP GWAS p-value based on the HuGE Navigator data [Spearman's correlation = -0.127, $p = 0.015$].

277