

1 **A Machine-Curated Database of Genome-Wide**

2 **Association Studies**

3 Volodymyr Kuleshov^{1,2,*}, Jialin Ding¹, Christopher Vo¹, Braden Hancock¹, Alexander
4 Ratner¹, Yang Li³, Christopher Ré¹, Serafim Batzoglou¹, Michael Snyder²

6 **Affiliations**

7 ¹Department of Computer Science, Stanford University, Stanford, CA

8 ²Department of Genetics, Stanford University School of Medicine, Stanford, CA

9 ³Department of Medicine, University of Chicago, Chicago, IL

10

11 Correspondence should be addressed to Volodymyr Kuleshov

12 (kuleshov@cs.stanford.edu)

13

14 **Abstract**

15 Tens of thousands of genotype-phenotype associations have been discovered to date, yet
16 not all of them are easily accessible to scientists. Here, we describe GwasKB, a novel
17 machine reading system that automatically collects and synthesizes genetic associations
18 from the scientific literature into a structured database. GwasKB helps curators by
19 automatically collecting >3,000 previously documented open-access relations (with an
20 estimated recall of 60-80%) as well as >2,000 associations not present in existing human-
21 curated repositories (with an estimated precision of 82-89%). Our system represents the
22 largest fully automated GWAS curation effort, and is made possible by a novel paradigm
23 for constructing machine learning systems called data programming. Our results
24 demonstrate both the importance and the feasibility of automating the curation of
25 scientific literature.

26

27 Genome-wide association studies (GWAS) are widely used for measuring the effects of
28 genomic mutations on human traits¹. Despite revealing tens of thousands of genotype-
29 phenotype associations, not all GWAS results are available to scientists in a structured
30 form amenable to downstream analyses.

31

32 Multiple efforts are underway to catalogue published GWAS associations^{2,3}, but it is as
33 yet unclear how far we are from a complete GWAS catalogue. Currently, even the most
34 exhaustive databases vary in their scope: hundreds to thousands of variants may be
35 present in one repository, but absent in all others^{2,3}. Variants that are omitted in a
36 database are effectively lost for downstream analyses, and as more studies are published,
37 the number of these “dark variants” is expected to increase. This limits the pace of
38 scientific research and represents an inefficient use of research funding.

39

40 Here, we describe GwasKB, a machine reading system that automatically collects and
41 synthesizes thousands of genotype-phenotype associations into a structured database.
42 Our system represents the largest GWAS machine curation effort, and is made possible
43 by a novel paradigm for constructing machine learning systems called data programming.
44 When deployed on a set of 589 open-access GWAS publications, GwasKB recovers (at
45 an estimated recall of 60-80%, depending on stringency criteria) >3,000 known
46 associations that were validated in existing GWAS databases, and finds >2,000
47 associations (with an estimated precision of 82-89%) currently absent in existing
48 repositories. The number of these new variants corresponds to about 20% of all open-
49 access associations recorded in the most up-to-date human-curated database, GWAS
50 Catalog.

51

52 We make available to curators an open-source implementation of GwasKB and we also
53 provide an online tool for browsing the associations found by our system¹. We anticipate
54 that these associations will be used by scientists in future work. More generally, we
55 demonstrate that modern machine reading algorithms have matured to the point of

¹ An online interface to our machine-curated database is available at <http://gwaskb.stanford.edu/>

56 significantly improving biomedical curation efforts. Finally, our system may form the
57 basis for further efforts to curate Mendelian mutations and other data.

58

59 **Automating Biomedical Literature Curation with GwasKB**

60

61 The results of genome-wide association studies are used to estimate disease risks^{4,5}, to
62 understand the function of specific genomic regions^{6,7}, and to train predictors for the
63 effects of new mutations⁸. Overall, about 2,500-3,000 studies have been performed to
64 date^{2,3}; they have reported tens of thousands of associations that are manually collected in
65 databases like GWAS Catalog² and GWAS Central³.

66

67 However, curating the results of GWAS studies is challenging, as it requires time,
68 domain expertise, and can be prone to errors. As a result, independent human curation
69 efforts are often not consistent, and even the largest GWAS databases are incomplete.

70

71 *The GwasKB Machine Reading System*

72

73 We propose that the process of collecting and synthesizing the findings of GWAS studies
74 can be made significantly more efficient using automated machine reading technologies.

75 We demonstrate this by introducing GwasKB, an automated system that extracts
76 genotype-phenotype relations from the biomedical literature and places them in a
77 structured SQL database (Figure 1).

78

79 Specifically, GwasKB collects three main pieces of information: genetic variants (as
80 defined by their RSID), their associated phenotypes, and their p-values. We support our
81 findings with evidence from publications (identified by their Pubmed ID), which can take
82 the form of a sentence excerpt or a location in a table.

83

84 Several challenges arise when curating GWAS studies. For one, there is no universally
85 adopted threshold for the significance of genotype-phenotype associations. GwasKB
86 reports all (rsid, phenotype) associations that are significant at $p < 10^{-5}$ in at least one

87 experiment in the study (such as in one cohort or one statistical model) and it records all
88 the other p-values relevant to that association. Our threshold of $p < 10^{-5}$ is the same as the
89 one used in the GWAS Catalog.

90

91 A second difficulty arises when describing the study phenotype. Phenotypes can be very
92 general (e.g., “heart disease”) or highly specific (e.g., “high systolic blood pressure”), and
93 existing databases often differ in their level detail. GwasKB addresses this issue by
94 providing simple and detailed phenotypes, i.e. a high-level description that applies to
95 every variant in the paper (e.g. “effects of proteins on inflammation”), and, when
96 available, a detailed description for specific variants (e.g., the name of a specific protein).

97

98 Lastly, a third difficulty is posed by copyright restrictions. With GwasKB, we restrict
99 ourselves to open-access papers, which represent approximately 25% of all the studies
100 that have been published to date. All open-access publications are catalogued by the
101 PubMed Central (PMC) repository and are made publicly available in XML format.
102 GwasKB takes these XML documents as input, although any paper in HTML format may
103 be parsed by our system after minor preprocessing. In the current version, we also discard
104 any associated files that need to be processed through proprietary software. However, the
105 principles of our system extend to all kinds of studies.

106

107 *On The Design of GwasKB*

108

109 GwasKB was designed to extract three key pieces of information: genetic variants, their
110 phenotypes, and their p-values. We have structured GwasKB into a set of five
111 components that extract this information.

112

113 The first component of GwasKB parses the title and abstract of every paper to identify a
114 simple phenotype that will be associated with all its variants. The second component
115 parses the body of the paper to find tuples of RSIDs and their associated detailed
116 phenotypes. Often, the detailed phenotype is abbreviated (e.g. BMI) and a third
117 component attempts to resolve these abbreviations (e.g. output “body mass index”). A

118 fourth component extracts p-values in the form of (rsid, p-value) tuples. Finally, the fifth
119 component constructs a single structured database from all these results.

120

121 Each GwasKB component has three stages: parsing, candidate generation, and
122 classification (Figure 2). Parsing is performed with Snorkel⁹, a knowledge base
123 construction framework for documents with richly formatted data (data expressed via
124 textual, structural, tabular, and/or visual cues), such as XML documents. Content is first
125 parsed for structure---the XML tree is traversed and converted into a hierarchical data
126 model with text assigned to tables, cells, paragraphs, sentences, etc. Then each sentence
127 or cell is parsed for content using the Stanford CoreNLP pipeline¹⁰, which performs
128 sentence tokenization, part-of-speech tagging, and syntactic parsing. In candidate
129 generation, we identify in the text mentions of some target relation (e.g., p-value/rsid
130 pairs). This is done by generating a large set of substrings from the text of the paper,
131 some of which could contain our target relation. Regular expressions or dictionaries are
132 used to identify candidates that may be valid instances of the relation we are looking for
133 (erring on the side of high recall over high precision). Finally, in the classification stage,
134 we determine which of these candidates are actually correct relation mentions using a
135 machine learning classifier. We use a Naive Bayes classifier with a small number of
136 hand-crafted features (between 4 and 12) and we train the model using the recently
137 proposed data-programming paradigm¹¹

138

139 One of the most significant bottlenecks in developing machine learning-based
140 applications today is the challenge of collecting large sets of hand-labeled training data.
141 Data programming is a newly proposed paradigm for training models using higher-level,
142 less precise supervision to avoid this bottleneck. In this approach, users write a set of
143 *labeling functions*: black-box functions that label data points, and that can subsume a
144 wide variety of heuristic approaches such as *distant supervision*¹²—where an external
145 knowledge base is used to label data points—regular expression patterns, heuristic rules,
146 and more. These labeling functions are assumed to be better than random, but otherwise
147 may have arbitrary accuracies, may overlap, and may conflict. A generative model is
148 used to learn their accuracies and correlations from unlabeled data. The predictions of

149 this model can then be used for classification, or to generate labels for a second,
150 discriminative model. We refer the reader to the appendix and to the full data
151 programming paper¹¹ for more details.

152

153 *Reproducibility*

154

155 In order to make our results fully reproducible, we have released Jupyter notebooks that
156 can be used to run GwasKB, generate the database of associations and recreate most of
157 our figures and tables. The notebooks and the source code of GwasKB are freely
158 available on GitHub at github.com/kuleshov/gwasdb.

159

160 In addition, we have built an interactive website that enables users to browse associations
161 that have been extracted by GwasKB. Users can search the data by study, phenotype or
162 variant rsid. The entire dataset can also easily be exported in text or SQL format.

163

164 **Machine Reading Helps Automate GWAS Curation**

165

166 We next demonstrate how our system can significantly help humans synthesize and
167 understand findings from the biomedical literature. We deploy GwasKB on all the open-
168 access papers listed in the GWAS Catalog database (589 in total), which is the most
169 complete set of such papers that we could access. For evaluation, we also use the GWAS
170 Central database. We use $p < 10^{-5}$ in at least one cohort or study methodology as our
171 significance cutoff, and assess both the precision and the recall of our system (see Table
172 1).

173

174 *GwasKB Recovers Up To 80% of Curated Open-Access Associations*

175

176 GWAS Central and GWAS Catalog contain respectively 3008 and 4023 accessible
177 associations in our set of 589 studies. These are variants whose RSID is contained in the
178 open-access XML content made available through PubMed Central. We also define
179 mappings between GwasKB phenotypes and phenotypes from GWAS Central and

180 GWAS Catalog (see Methods). These databases often use different levels of precision to
181 describe phenotypes (e.g. “smoking behaviors” vs. “cigarette packs per day”); therefore,
182 we also specify whether our reported phenotype is exact or approximate; in the latter
183 case, it is still useful, but lacks some detail. Table 2 contains examples of relations
184 extracted by GwasKB.

185

186 Among the set of open-access papers, GwasKB recovered 2487 (82%) relations with
187 approximately correct phenotypes from GWAS Central and 3245 (81%) relations from
188 the GWAS Catalog. It also recovered 1890 (63%) relations with full accuracy from
189 GWAS Central and 2762 (69%) relations from GWAS Catalog. A number of known
190 associations were not correctly recovered because their reported phenotype was incorrect
191 (89 in GWAS Central and 147 in GWAS Catalog). In the remaining cases, we were not
192 able to report the variant itself. Overall, GwasKB recovered 81-82% of accessible
193 associations at a level of quality that will be useful in many applications.

194

195 *Machine Curation Uncovers Many Associations Not Found by Human Curators*

196

197 In total, GwasKB discovered 6422 relations within the 589 input papers, 2959 (46%) of
198 which could not be mapped to GWAS Catalog or GWAS Central. Notably, many of these
199 appeared to be valid.

200

201 We investigated this further by first manually inspecting a random subset of 100 novel
202 relations (with independent validation from two independent annotators). We found that
203 82 relations fully met the specifications of our system, 11 were incorrect, and 7 were
204 originally identified by a different study (and referenced as background material). Most
205 of the errors can be attributed to incorrect phenotypes. Of the 82 relations matching
206 system specifications, 60 appeared to satisfy the same criteria as GWAS Central or
207 GWAS Catalog relations from the same paper, while 22 were not significant at 10^{-5} in all
208 cohorts. The latter may have been omitted by human curators for this reason.

209

210 **Novel Variants Found by GwasKB Are Correlated with Genomic Function**

211

212 *Linkage Disequilibrium Between Variants from GwasKB and from Existing Databases*

213

214 To validate the novel variants found by our system, we conducted a series of analyses
215 aimed at characterizing the variants' function. First, we reasoned that detected variants
216 may be in linkage disequilibrium (LD) with known variants (because they originate from
217 the same LD block), or among themselves, thereby inflating our number of truly novel
218 associations.

219

220 We estimated LD from the Thousand Genomes dataset (Supplementary Methods); Figure
221 3 shows the histogram of r^2 distances between each novel variant, and its closest variant
222 in the GWAS Catalog. The distribution of r^2 scores is highly multimodal, with large
223 peaks at $r^2=1$, and many more at $r^2=0$.

224

225 Using a threshold of $r^2 > 0.5$, we filtered our set of new [pmid, rsid, phen, pvalue]
226 associations from 3170 to 1494 by removing variants in LD with known manually
227 curated variants; of the 1676 variants that we eliminated, 765 were not in the 1000
228 Genomes database or their closest previously known variant was not in the database; the
229 remaining 911 SNPs were in LD with known variants. We further reduced this set to
230 1304 associations by eliminating novel variants that were in LD with each other. Thus,
231 although many variants are in LD with known variants, over 40% of our discovered
232 variants do not appear to be linked to variants previously identified in GWAS databases.

233

234 We argue that it is preferable to curate both novel and known variants, since we do not
235 know which mutation in an LD block is truly causal and the r^2 cutoff for defining LD
236 blocks is somewhat arbitrary and may vary. We think that filtering should be performed
237 by the user, depending on their goal; this is also the approach taken by the GWAS
238 Central repository. Moreover, if the authors of a GWAS study report multiple variants in
239 LD, we believe that it is better to report their findings as they are, rather than introducing
240 additional bias through our own filtering.

241

242 *Comparison to Alternative Approaches for Estimating Variant Significance*

243

244 Our second analysis focuses on the biological function of the novel variants. We focus on
245 two large classes of phenotypes: neurodegenerative diseases (ND; 27 traits, including
246 Autism, Alzheimer's, Parkinson's, etc.) and autoimmune disorders (AI; 23 traits,
247 including Diabetes, Arthritis, Lupus, etc.); for the analyses below, we consider the subset
248 of variants that are not in LD with any variant in the GWAS Catalog or GWAS Central
249 (283 ND SNPs and 155 AI SNPs).

250

251 We also collected two sets of genes that were found to be highly expressed in brain cells
252 as well as in blood cells; specifically, we reasoned that SNPs associated to
253 neuropsychiatric and autoimmune diseases should be more highly enriched near genes
254 expressed in brain and immune cells, respectively. Indeed, we found that variants
255 associated with ND diseases (32 ND SNPs in total) occurred significantly more often
256 within 200Kbp of genes with preferential brain expression, while variants associated with
257 AU traits (15 variants in total) were found much more frequently in near genes with
258 preferential blood expression ($p < 0.05$; see Supplementary Material).

259

260 We should note however that the vast majority of ND and AU variants were found far
261 from coding regions. To test whether this set of SNPs also make biological sense, we
262 used GREAT¹³, a tool which annotates the function of variants in intergenic areas of the
263 genome. In particular, GREAT links intergenic regions with Disease Ontology (DO)
264 terms, and outputs terms that are significantly enriched for a particular set of variants.
265 When we applied GREAT to ND SNPs, we found a strong enrichment in regions known
266 to play a role in ND-related phenotypes, such as cognitive disease ($p < 10^{-32}$), dementia (p
267 $< 10^{-23}$), and neurodegenerative disease ($p < 10^{-23}$). Similarly, AI variants were
268 significantly associated with AI-related terms, the most significant of which were disease
269 by infectious agent ($p < 10^{-27}$), viral infectious disease ($p < 10^{-19}$), and autoimmune
270 disease ($p < 10^{-17}$). In fact, the top 20 DO terms for either set of variants were all
271 exclusively associated with the correct family of phenotypes (Supplementary Tables 1,2).
272 Hence, our predicted variants were highly consistent with these external annotations.

273

274 *Examining the Effect Sizes of Novel GwasKB Variants*

275

276 Finally, we analyzed the magnitude with which novel variants affect their predicted
277 phenotypes and other, related traits. Specifically, we used freely available GWAS
278 summary statistics from the LD Hub project¹⁴ to assess the distribution of SNP effect
279 sizes across novel variants and compared them to those of random SNPs. We focused on
280 the 11 most frequent traits in our dataset for which summary statistics were available; for
281 each trait, we identified an LD Hub study that provides effect sizes (in the form of beta
282 coefficients or log odds ratios) for that trait. Figure 4 compares the distribution of effect
283 sizes of the novel variants identified by GwasKB to the distribution of effects sizes for all
284 SNPs, again restricting to variants that show no LD with other variants in GWAS
285 databases. Whereas the distribution of random SNPs is centered around zero, as one
286 would expect, novel SNP effect sizes appear to follow a different distribution
287 (Kolmogorov-Smirnov Test; see Figure 4 and Supplementary Figures 1,2) and tend to
288 have significantly higher magnitudes than expected.

289

290 We also examined the effects of GwasKB variants on phenotypes which are known to be
291 related to their primary, predicted trait. For each pair of diseases, we took the set of
292 variants that GwasKB found to be associated with the first disease, and computed their
293 average absolute effect size using summary statistics from the second disease; in several
294 cases, variants that we determined to be associated with one trait (e.g. Obesity) also had
295 large effect sizes on related traits (e.g. BMI).

296

297 Specifically, we used a permutation test to compute the probability of observing the
298 absolute average effect size among novel variants within a random set of SNPs; Figure 5
299 shows the resulting matrix of p-values (we only include traits for which we computed at
300 least one small p-value). In particular, we found that three traits (Obesity, BMI, Type 2
301 Diabetes) shared variants with high effect sizes. These three phenotypes are known to be
302 highly correlated.

303

304 Interestingly, we also observed an unusual correlation between LDL Cholesterol levels
305 and Alzheimer's disease. To investigate this further, we repeated the same analysis using
306 variants that have been confirmed by the GWAS Catalog (Supplementary Figures 3-5).
307 The resulting matrix resembles closely that of novel variants and also shows correlation
308 between Alzheimer's and LDL cholesterol. We also found a novel variant (rs6857) that
309 was previously found to be associated with Alzheimer's¹⁵; our system also correctly
310 determined that is associated with LDL¹⁶ at $p < 10^{-7}$; this association is notably missing
311 from current manually-curated databases.

312

313 **Discussion**

314

315 *The importance of curation.* If GWAS associations are not recorded in a database, they
316 are effectively missing for many practical purposes, e.g. for training machine learning
317 systems (to predict SNP function). GWAS studies are also costly (often involving
318 genotyping tens of thousands of subjects), and it thus a waste of research funding to not
319 fully record their results.

320

321 An alternative to curation to ask authors to directly report their findings online. This is
322 already possible within GWAS Central, although in practice not all authors do this, and
323 hence the database is far from complete. In addition, past studies still need to be curated.
324 An ideal solution appears to involve a combination of authors, machines, and curators.

325

326 *Hand-curation is a difficult task.* Why do manual curation efforts miss certain
327 associations? Curating papers is often a tedious task involving browsing through highly
328 technical material in search of short snippets of text. Humans are generally not well-
329 suited to this kind of work: they may accidentally skip table rows, or become tired and
330 skip a paragraph. Curation also requires understanding advanced technical concepts such
331 as linkage disequilibrium or multiple hypothesis testing. This makes the task unsuitable
332 for crowdsourcing approaches.

333

334 *Machines may outperform humans.* Computers, on the other hand, don't suffer from the
335 aforementioned limitations: they excel at repetitive work and only need to be
336 programmed by experts once. Crucially, even though machines make errors, these errors
337 are systematic, not random: one may follow an iterative process of fixing these errors and
338 redeploying the system, until a sufficient level of accuracy is reached. Redeploying our
339 system takes on the order of hours, while asking humans to return and correct their errors
340 would take at least months.

341

342 Of course, humans also have many advantages over machines. Indeed, the sets of
343 GwasKB and human-curated associations were quite distinct, with thousands of relations
344 present in one set, but not the other. The most accurate and complete GWAS database is
345 in fact a combination of both sources. In the future, we see curation as a collaboration
346 between humans and machines.

347

348 *Biomedical information extraction.* Extracting structured relations from unstructured text
349 is subject of the field of information extraction¹⁷ (IE). Information extraction is widely
350 used in diverse domains such as news¹⁸, finance¹⁹, geology²⁰, and in the biomedical
351 domain. In the biomedical setting, IE systems have been used to parse electronic medical
352 records²¹, identify drug-drug interactions²², and associate genotypes with drug response²³.
353 A considerable amount of effort has gone into uncovering gene/disease associations from
354 biomedical literature²⁴. Our approach, however, takes a different approach, as it attempts
355 to identify the effects of individual mutations. Recently, Jain et al. applied information
356 extraction to the GWAS domain²⁵; their work focused on creating extractors for two
357 specific relations: paper phenotypes and subject ethnicities; these extractors achieved an
358 87% precision-at-2 and a 83% F1-score on the two tasks, respectively. In contrast, our
359 works introduces an end-to-end system that extracts full (phenotype, rsid, pvalue)
360 relations comparable to ones found in hand-curated databases.

361

362 *Applications beyond GWAS studies.* Dozens of literature curation efforts are currently
363 underway in cancer genomics, pharmacogenomics, and many other fields. Our findings
364 hint at the possibility of using machine curation there as well.

365

366 The GWAS domain is in many ways easier than others since variants have standardized
367 identifiers and a lot of information is structured in tables. Nonetheless, it allows us to
368 demonstrate the importance of machine curation and to develop a core system that can be
369 generalized to other domains. Within the GWAS setting, our system can be further
370 improved by extracting additional relations (e.g. risk alleles, odds ratios).

371

372 **Conclusion**

373

374 In summary, we have introduced in this work a new machine reading system for
375 extracting structured databases from publications describing genome-wide association
376 studies, and we have used it to both recover many known relations, as well as a number
377 of associations that were not present in any existing repository.

378

379 Our results demonstrate how machine reading algorithms may help human curators
380 synthesize the large amount of knowledge contained in the biomedical literature. This
381 knowledge can be made widely accessible using new systems that combine the efforts of
382 both human and machines, thus accelerating the pace of discovery in science.

383

384

385 **Online Methods**

386

387 *Detailed Description of GwasKB*

388

389 GwasKB is implemented in Python on top of the Snorkel information extraction
390 framework¹¹. Snorkel provides utilities for parsing XML documents and training machine
391 learning classifiers. GwasKB extends the parsers/classifiers in Snorkel and applies them
392 to the GWAS extraction task. Below, we provide additional details on the various
393 components of GwasKB

394

395 *Identifying simple phenotypes.* We parse paper titles and abstracts and generate
396 candidates from the EFO, Snomed and Mesh ontologies. We use 11 labeling functions
397 (LFs), which include the following: is the mention in the title; is the mention less than 5
398 characters; does the mention contain nouns; is the mention in the first half of the
399 sentence, etc. We include the full list of labeling functions in our open-source GitHub
400 repository. The high-level phenotype is the set of three highest scoring mentions
401 exceeding a user-specified score threshold or the single highest mention if none exceeds
402 the threshold; this enables us to handle multiple valid phenotypes.

403

404 *Identifying precise phenotypes.* We only parse tables and generate candidates from cells
405 whose header contains the words "phenotype", "trait", or "outcome". Candidate p-values
406 are generated by matching a regular expression; candidate relations consist of
407 horizontally aligned phenotype and p-value candidates. We use three labeling functions:
408 is the candidate mostly a number; is the header of the cell (indicating it's in a phenotype
409 column) very long; does the mention contain words referring to an rsid. The module is
410 described in more detail on GitHub.

411

412 *Resolving Acronyms.* We resolve acronyms by looking at the entire paper, including
413 tables and the main natural language text in the body of the paper. We extract candidates
414 from aligned pairs table cells, where one row is labeled "phenotype", "trait", or
415 "description", while the other is labeled "abbreviation", "acronym", or "phenotype". We

416 generate candidates from the main text using a regular expression. Our labeling functions
417 include the following: is the candidate all in caps; does the candidate match to the
418 Snomed dictionary; does the acronym candidate consist of the letters of each word of the
419 phenotype candidate; is one a prefix of the other; etc. The module for resolving
420 abbreviations is also described on GitHub.

421

422 *Identifying p-values.* We again generate candidates from tables; SNP candidates are
423 generated using a regular expression; p-value candidates are ones that match one of three
424 regular expressions (see GitHub); candidate relations consist of horizontally aligned SNP
425 and p-value candidates (with at most one rsid per row). These candidates were accurate
426 and we report them all.

427

428 *Mapping Phenotypes Across Databases*

429

430 In order to compare against GWAS Central and GWAS Catalog, we define mappings
431 between GwasKB phenotypes and ones used in these two repositories. These mappings
432 are tables with about 800 entries each that also indicate whether the mapping is fully or
433 partially correct (e.g. "smoking behaviors" vs "packs per day"). We define the latter as
434 conceptually containing the precise label while also being not so broad as to be useless.
435 See also our earlier discussion on high- and low-level phenotypes. To confirm the
436 validity of our mappings, we asked an independent annotator to label 100 random table
437 entries; their concordance with our labels was 95%. These mappings are available in our
438 GitHub repository.

439

440 *Understanding The Errors Of GwasKB Components*

441

442 *Simple phenotype extraction.* Errors at this stage mostly occur when the true phenotypes
443 are not found in our candidate dictionaries (e.g. "genome-wide association study in
444 bipolar patients"; we can only generate the candidate "bipolar disorder"). The second
445 major source of error are phenotypes mentioned only in passing (e.g. "high body fat is a

446 risk for diabetes" when diabetes is not the phenotype whose association is being
447 reported).

448

449 To estimate the precision of this module, we first restrict ourselves to (paper, rsid,
450 phenotype) relations produced by GwasKB that are also confirmed by an existing
451 database, in the sense that the variant specified by the rsid occurs in *some* relation
452 associated with the paper (but not necessarily one with the same phenotype). Then, we
453 look at the fraction of these relations whose phenotype is also correct (at the approximate
454 level). This gives precisions of 97% in the GWAS Catalog and 96% in GWAS Central.

455

456 *Detailed phenotype extraction.* Most errors occur because we do not correctly resolve
457 acronyms or because low-level phenotypes are not in tables (but rather only in text).
458 Acronyms are not resolved most often because the shortened symbol is not clearly related
459 to the full expression (e.g. CYS5 for Cysteine proteinase inhibitor 5 precursor), and they
460 are presented in tables with confusing formatting. We estimate precision in the same way
461 as for simple phenotypes, but this time, we require that phenotype agree fully. Precision
462 was 73% in GWAS Central, the database with the most precise phenotypes. In GWAS
463 Central, it was 82%.

464

465 *p-values.* To evaluate p-value extraction accuracy, we labeled by hand 100 random
466 relations and found that our rule-based extraction procedure had a precision of 98%.
467 Errors occurred when p-values referred to other entities in the row, such as haplotypes.
468 Note also that oftentimes, variants and their p-values are only provided in text but not in
469 tables. This was the primary reason why we failed to report the rsid's of 584 (15%)
470 GWAS Catalog and 432 (14%) GWAS Central associations.

471

472 *Error Analysis Over 100 New Relations*

473

474 Of the incorrect relations, 7 were due to incorrect phenotype labels (but the underlying
475 SNP was significant) and 4 were due to table parsing errors (the p-value was extracted
476 incorrectly). Of the 22 variants that were not significant in all cohorts, 18 could be

477 identified as such via extracted tags and relations. We also determined that 60 relations
478 were correct because they were either described as “significant” in the paper text (in
479 addition to having $p < 10^{-5}$ in all cohorts) or they had essentially the same or higher level
480 of significance as SNPs that were included in GWAS Catalog or GWAS Central.

481

482 To confirm the accuracy of our analysis, we asked two independent annotators with
483 expertise in genomics to label a random subset of 50 associations out of the ones
484 analyzed above. For each annotator, respectively 47 and 48 out of 50 labels were
485 consistent with ours. We publish our 100 samples and their annotation on GitHub; for
486 each example, we add a justification for our label.

487

488 *Estimating the Precision of GwasKB*

489

490 We estimate our overall precision at 92%: we consider the 3463 relations confirmed by
491 existing databases as correct, and estimate the error rate on the other relations to be 18%
492 (incorrect and repeat relations).

493

- 494 1. Bush, W.S. & Moore, J.H. Chapter 11: Genome-Wide Association Studies. *PLoS*
495 *Comput. Biol.* **8**, 1-11 (2012).
- 496 2. Welter, D. et al. The NHGRI GWAS Catalog, a curated resource of SNP-trait
497 associations. *Nucleic Acids Res.* **42**, D1001-D1006 (2014).
- 498 3. Beck, T., Hastings, R.K., Gollapudi, S., Free, R.C. & Brookes, A.J. GWAS Central: a
499 comprehensive resource for the comparison and interrogation of genome-wide
500 association studies. *European Journal of Human Genetics* **22**, 949-952 (2013).
- 501 4. Cariaso, M. & Lennon, G. SNPedia: a wiki supporting personal genome
502 annotation, interpretation and analysis. *Nucleic Acids Res.* **40**, D1308-D1312
503 (2012).
- 504 5. *Promethease*.
- 505 6. de Leeuw, C.A., Mooij, J.M., Heskes, T. & Posthuma, D. MAGMA: Generalized
506 Gene-Set Analysis of GWAS Data. *PLoS Comput. Biol.* **11**, e1004219+ (2015).
- 507 7. Weng, L. et al. SNP-based pathway enrichment analysis for genome-wide
508 association studies. *BMC Bioinformatics* **12**, 99+ (2011).
- 509 8. Zhou, J. & Troyanskaya, O. Predicting effects of noncoding variants with deep
510 learning-based sequence model. *Nat Meth* **12**, 931-934 (2015).
- 511 9. Ratner, A.J., Bach, S.H., Ehrenberg, H.R. & Ré, C. Snorkel: Fast training set
512 generation for information extraction. *Proceedings of the 2017 ACM International*
513 *Conference on Management of Data* 1683-1686 (2017).
- 514 10. Manning, C.D. et al. The Stanford CoreNLP Natural Language Processing Toolkit.
515 *Association for Computational Linguistics (ACL) System Demonstrations* 55-60
516 (2014).
- 517 11. Ratner, A.J., De Sa, C.M., Wu, S., Selsam, D. & Ré, C. Data programming: Creating
518 large training sets, quickly. *Advances in Neural Information Processing Systems*
519 3567-3575 (2016).
- 520 12. Mintz, M., Bills, S., Snow, R. & Jurafsky, D. Distant supervision for relation
521 extraction without labeled data. *Proceedings of the Joint Conference of the 47th*
522 *Annual Meeting of the ACL and the 4th International Joint Conference on Natural*
523 *Language Processing of the AFNLP: Volume 2-Volume 2* 1003-1011 (2009).
- 524 13. McLean, C.Y. et al. GREAT improves functional interpretation of cis-regulatory
525 regions. *Nat Biotechnol* **28**, 495-501 (2010).
- 526 14. Zheng, J. et al. LD Hub: a centralized database and web interface to perform LD
527 score regression that maximizes the potential of summary level GWAS data for
528 SNP heritability and genetic correlation analysis. *Bioinformatics* **33**, 272-279
529 (2017).
- 530 15. Yu, C. et al. Comprehensive analysis of APOE and selected proximate markers for
531 late-onset Alzheimer's disease: patterns of linkage disequilibrium and
532 disease/marker association. *Genomics* **89**, 655-665 (2007).
- 533 16. Deshmukh, H.A. et al. Genome-wide association study of genetic determinants of
534 LDL-c response to atorvastatin therapy: importance of Lp (a). *Journal of lipid*
535 *research* **53**, 1000-1011 (2012).
- 536 17. Moens, M. *Information Extraction: Algorithms and Prospects in a Retrieval*
537 *Context*. (Springer Netherlands: 2009).

- 538 18. Tumarkin, R. & Whitelaw, R.F. News or Noise? Internet Postings and Stock Prices.
539 *Financial Analysts Journal* **57**, 41-51 (2001).
- 540 19. Das, S. & Chen, M. Yahoo! for Amazon: Extracting Market Sentiment from Stock
541 Message Boards. *Proceedings of the Asia Pacific Finance Association Annual*
542 *Conference (APFA)* (2001).
- 543 20. Zhang, C. et al. GeoDeepDive: Statistical Inference Using Familiar Data-processing
544 Languages. *Proceedings of the 2013 ACM SIGMOD International Conference on*
545 *Management of Data* 993-996 (2013).doi:10.1145/2463676.2463680
- 546 21. Zhou, X., Han, H., Chankai, I., Prestrud, A. & Brooks, A. Approaches to Text Mining
547 for Clinical Medical Records. *Proceedings of the 2006 ACM Symposium on Applied*
548 *Computing* 235-239 (2006).doi:10.1145/1141277.1141330
- 549 22. Percha, B., Garten, Y. & Altman, R.B. Discovery and explanation of drug-drug
550 interactions via text mining. *Pacific Symposium on Biocomputing. Pacific*
551 *Symposium on Biocomputing* 410-421 (2012).
- 552 23. Rinaldi, F., Schneider, G. & Clematide, S. Relation Mining Experiments in the
553 Pharmacogenomics Domain. *J. of Biomedical Informatics* **45**, 851-861 (2012).
- 554 24. Pletscher-Frankild, S., Pallega, A., Tsafou, K., Binder, J.X. & Jensen, L.J. DISEASES:
555 Text mining and data integration of disease-gene associations. *bioRxiv* 008425+
556 (2014).doi:10.1101/008425
- 557 25. Jain, S. et al. Weakly supervised learning of biomedical information extraction
558 from curated data. *BMC Bioinformatics* **17**, S1 (2016).
- 559
560
561

Biomedical Publication

ARTICLES							
nature genetics							
Genome-wide association study of blood pressure and hypertension							
Table 1 Genome-wide association results for SBP-associated SNPs with P							
CHARGE meta-analysis, SBP							
SNP identifier	Chr	Position	Gene	MAF	Beta	s.e.	P
rs2681492	12	88537220	ATP2B1	0.20	-1.26	0.19	3.0E-11
rs2681472	12	88533090	ATP2B1	0.18	-1.29	0.19	3.5E-11
rs11105354	12	88550654	ATP2B1	0.18	-1.30	0.20	3.7E-11

Here we report results of a genome-wide association study of systolic (SBP) blood pressure



Structured Database

Variant	rs2681492
Simple phenotype	Hypertension Blood pressure
Detailed phenotype	Systolic
p-value	3.0e-11
Source	PMID: 19430479, Tbl. 1

562

563

564

565

566

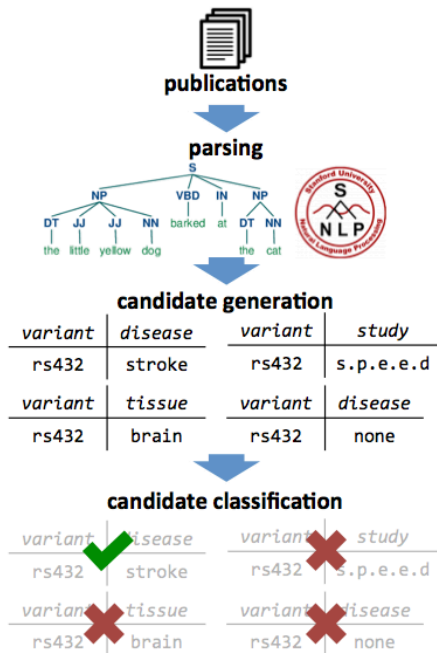
567

568

569

Figure 1: The GwasKB machine reading system. GwasKB takes as input a set of biomedical publications retrieved from PubMed Central (left) and automatically creates a structured database of GWAS associations described in these publications (right). For each association, the system identifies a genetic variant (purple), a high-level phenotype (pertaining to all variants in the publication), a detailed low-level phenotype (specific to individual variants, if available; red), and a p-value (orange). Acronyms are also resolved (red).

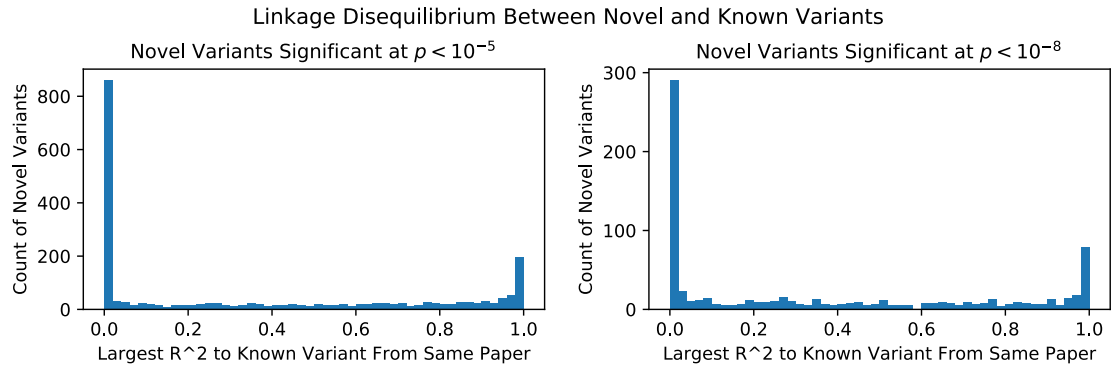
570



571

572 **Figure 2:** General structure of a GwasKB module. The system contains separate modules
 573 for extracting variants, phenotypes, p-values, and for resolving acronyms. Each module
 574 consists of three stages. At the parsing stage, we process papers using the Stanford
 575 CoreNLP pipeline, performing full syntactic parsing. Next, given a target relation (e.g.,
 576 variant-phenotype), we generate a large set of candidates, some of which could be correct
 577 instances of the target object on relation. Then, at the classification stage, we determine
 578 which candidates are correct using a machine learning classifier.

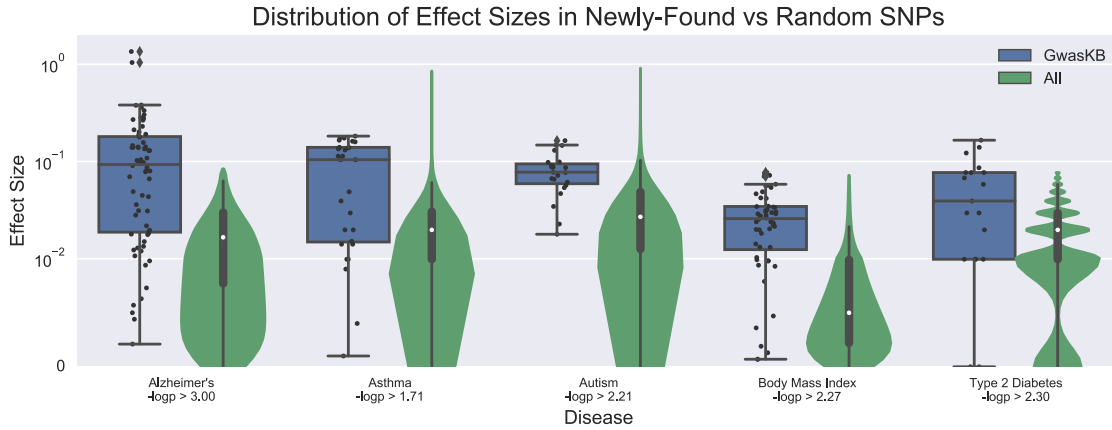
579



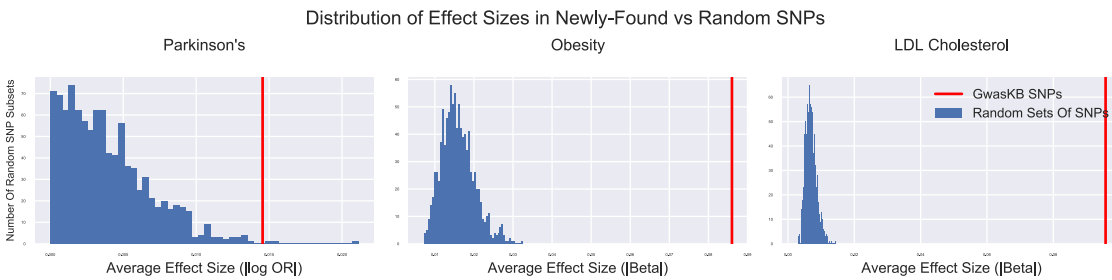
580

581 **Figure 3:** Linkage disequilibrium between GwasKB variants not present in existing
 582 human curated databases and variants from the GWAS Catalog. We use the 1000
 583 Genomes dataset to estimate the r^2 metric between pairs of variants, and report distances
 584 from each GwasKB variant to the most correlated GWAS Catalog SNP reported in the
 585 same paper. The distribution of r^2 scores is highly multimodal; many GwasKB variants
 586 are uncorrelated ($r^2=0$) with GWAS Catalog SNPs.

587



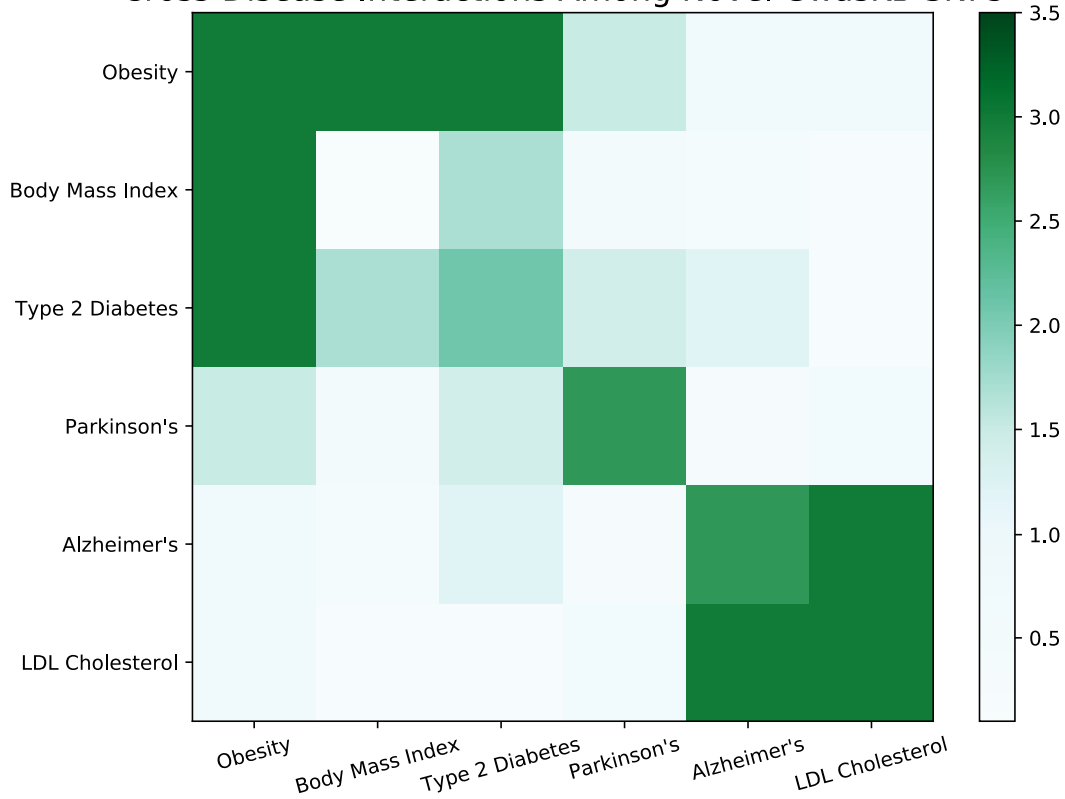
588



589

590 **Figure 4:** Visualizing the effect sizes of variants identified by GwasKB. *Top:* We
 591 compare the distribution of effect sizes (absolute values of beta coefficients or log odds
 592 ratios; data from LD Hub) of variants identified by GwasKB (blue) to that of all variants
 593 (green) for multiple traits. Blue variant effect sizes cluster away from zero and follow a
 594 different distribution (Kolmogorov-Smirnov test). *Bottom:* We subsample 1000 random
 595 sets of variants with the same number of elements as the set of GwasKB SNPs for a given
 596 disease; the average effect size of GwasKB variants (red) is higher than that of the
 597 random subsets (blue). In all settings, we only look at novel GwasKB variants not present
 598 in existing human-curated repositories.

Cross-Disease Interactions Among Novel GwasKB SNPs



599

600 **Figure 5:** Visualizing the effects of variants identified by GwasKB for pairs of related
601 phenotypes. For each pair of phenotypes, we compute the average absolute effect size of
602 GwasKB SNPs from the first phenotype (left) using summary statistics from the second
603 phenotype (right; summary statistics were obtained from the LD hub). The heat map
604 displays the log-probability of observing an equal or greater effect size by sampling
605 random variants (we thus compute p-values using a one-sided permutation test). Variants
606 predicted by GwasKB to be associated with Obesity, BMI, or Type 2 Diabetes also have
607 significant effects sizes for other, related diseases within this trio. In this analysis, we
608 only look at novel GwasKB variants not present in existing human-curated repositories.

609

610

Database	Statistics over open-access papers		
	Papers	Associations	Unique Associations
GWAS Catalog	589	8,384	>2,026
GWAS Central	516	5,914	>364
GwasKB (ours)	589	6,231	>2,777

611

612 **Table 1:** Numbers of associations contained in different GWAS databases. Unique
613 associations are contained in one database and in none of the others. Human curated
614 databases (GWAS Catalog and GWAS Central) significantly differ in their scope. Our
615 machine-curated repository (GwasKB) automatically recovers a large fraction of known
616 results and also finds a comparable number of unique associations.

617

	Source	Simple phenotype	Precise phenotype	p-value
Study	<i>Genome-wide pharmacogenomic study of metabolic side effects to antipsychotic drugs.</i>			
rs17661538	GwasKB	Antipsychotic drugs / Metabolic side effects	Clozapine - Triglycerides	1.00E-06
	GwasCat	Clozapine-induced change in triglycerides		1.00E-06
Study	<i>Genome-wide meta-analysis identifies seven loci associated with platelet aggregation in response to agonists.</i>			
rs12566888	GwasKB	Platelet aggregation	-	5.00E-19
	GwasCat	Platelet aggregation, epinephrine		5.00E-19
Study	<i>A genome-wide association study of the Protein C anticoagulant pathway.</i>			
rs13130255	GwasKB	Protein C	funcPS	3.00E-06
	GwasCat	Anticoagulant levels (funcPS)		3.00E-06
Study	<i>Genome-wide association study of CSF levels of 59 Alzheimer's disease candidate proteins: significant associations with proteins involved in amyloid processing and inflammation.</i>			
rs948399	GwasKB	Proteins Involved / Inflammation / Alzheimer's Disease	metalloproteinase-3	1.00E-07

618

619 **Table 2:** Examples of associations identified by GwasKB. Associations can be classified
620 as correct (rs17661538), partially correct (rs12566888; the precise phenotype is missing)
621 and incorrect (rs13130255). We also compare these associations to their corresponding
622 entries in the GWAS Catalog. The last entry (rs948399) is an example of a previously
623 undocumented association discovered by our system.

624

625 **Contributions**

626

627 V.K. conceived the study. B.H., V.K., and A.R. developed modules for the Snorkel
628 system. V.K. developed the GwasKB system. V.K., J.D., and C.V. performed
629 computational analysis. J.D. developed the web interface. V.K. and Y.L. wrote the paper.
630 Y.L., C.R., S.B., and M.S. supervised the study.

631

632 **Competing interests**

633

634 None declared.

635