

Identification of disease treatment mechanisms through the multiscale interactome

Camilo Ruiz,^{1,2} Marinka Zitnik,³ Jure Leskovec^{1,4,*}

¹ Computer Science Department, Stanford University, Stanford, CA 94305, USA

² Bioengineering Department, Stanford University, Stanford, CA 94305, USA

³ Biomedical Informatics Department, Harvard University, Boston, MA 02115, USA

⁴ Chan Zuckerberg Biohub, San Francisco, CA 94158, USA

*Corresponding author. Email: jure@cs.stanford.edu

1 **Most diseases disrupt multiple proteins, and drugs treat such diseases by restoring the func-**
2 **tions of the disrupted proteins. How drugs restore these functions, however, is often unknown**
3 **as a drug's therapeutic effects are not limited only to the proteins that the drug directly tar-**
4 **gets. Here, we develop the multiscale interactome, a powerful approach to explain disease**
5 **treatment. We integrate disease-perturbed proteins, drug targets, and biological functions**
6 **into a multiscale interactome network, which contains 478,728 interactions between 1,661**
7 **drugs, 840 diseases, 17,660 human proteins, and 9,798 biological functions. We find that**
8 **a drug's effectiveness can often be attributed to targeting proteins that are distinct from**
9 **disease-associated proteins but that affect the same biological functions. We develop a ran-**
10 **dom walk-based method that captures how drug effects propagate through a hierarchy of**
11 **biological functions and are coordinated by the protein-protein interaction network in which**
12 **drugs act. On three key pharmacological tasks, we find that the multiscale interactome pre-**
13 **dicts what drugs will treat a given disease more effectively than prior approaches, identifies**
14 **proteins and biological functions related to treatment, and predicts genes that interfere with**
15 **treatment to alter drug efficacy and cause serious adverse reactions. Our results indicate that**
16 **physical interactions between proteins alone are unable to explain the therapeutic effects of**
17 **drugs as many drugs treat diseases by affecting the same biological functions disrupted by**
18 **the disease rather than directly targeting disease proteins or their regulators. We provide**
19 **a general framework for identifying proteins and biological functions relevant in treatment,**
20 **even when drugs seem unrelated to the diseases they are recommended for.**

21 Complex diseases, like cancer, disrupt dozens of proteins that interact in underlying bio-
22 logical networks [1–4]. Treating such diseases requires practical means to control the networks
23 that underlie the disease [5–7]. By targeting even a single protein, a drug can affect hundreds of
24 proteins in the underlying biological network. To achieve this effect, the drug relies on physical
25 interactions between proteins. The drug binds a target protein, which physically interacts with
26 dozens of other proteins, which in turn interact with dozens more, eventually reaching the proteins
27 disrupted by the disease [8–10]. Networks capture such interactions and are a powerful paradigm
28 to investigate the intricate effects of disease treatments and how these treatments translate into
29 therapeutic benefits, revealing insights into drug efficacy [10–15], side effects [16], and effective
30 combinatorial therapies for treating the most dreadful diseases, including cancers and infectious
31 diseases [17–19].

32 However, existing systematic approaches assume that, for a drug to treat a disease, the pro-
33 teins targeted by the drug need to be *close* to or even need to *coincide* with the disease-perturbed
34 proteins [10–14] (Figure 1). As such, current approaches fail to capture biological functions,
35 through which target proteins can restore the functions of disease-perturbed proteins and thus treat
36 a disease [20–25] (Supplementary Fig. 3). Moreover, current systematic approaches are “black-
37 boxes:” they predict treatment relationships but provide little biological insight into how treatment
38 occurs. This suggests an opportunity for a systematic, explanatory approach. Indeed for particu-
39 lar drugs and diseases, custom networks have demonstrated that incorporating specific biological
40 functions can help explain treatment [26–29].

41 Here we present the multiscale interactome, a powerful approach to explain disease treat-
42 ment. We integrate disease-perturbed proteins, drug targets and biological functions in a mul-
43 tiscale interactome network. The multiscale interactome uses the physical interaction network
44 between 17,660 human proteins, which we augment with 9,798 biological functions, in order to
45 fully capture the fundamental biological principles of effective treatments across 1,661 drugs and
46 840 diseases.

47 To identify how a drug treats a disease, our approach uses biased random walks which model
48 how drug effects spread through a hierarchy of biological functions and are coordinated by the
49 protein-protein interaction network in which drugs act. In the multiscale interactome, drugs treat
50 diseases by propagating their effects through a network of physical interactions between proteins
51 and a hierarchy of biological functions. For each drug and disease, we learn a diffusion profile,
52 which identifies the key proteins and biological functions involved in a given treatment. By com-

53 paring drug and disease diffusion profiles, the multiscale interactome provides an interpretable
54 basis to identify the proteins and biological functions that explain successful treatments.

55 We demonstrate the power of the multiscale interactome on three key tasks in pharmacology.
56 First, we find the multiscale interactome predicts which drugs can treat a given disease more accu-
57 rately than existing methods that rely on physical interactions between proteins (i.e. a molecular-
58 scale interactome). This finding indicates that our approach accurately captures the biological
59 functions through which target proteins affect the functions of disease-perturbed proteins, even
60 when drugs are distant to diseases they are recommended for. The multiscale interactome also
61 improves prediction on entire drug classes, such as hormones, that rely on biological functions and
62 thus cannot be accurately represented by approaches which only consider physical interactions be-
63 tween proteins. Second, we find that the multiscale interactome is a “white-box” method with the
64 ability to identify proteins and biological functions relevant in treatment. Finally, we find that the
65 multiscale interactome predicts what genes alter drug efficacy or cause serious adverse reactions
66 for a given treatment and identifies biological functions that help explain how these genes interfere
67 with treatment.

68 Our results indicate that the failure of existing approaches is not due to algorithmic limita-
69 tions but is instead fundamental. We find that a drug can treat a disease by influencing the behaviors
70 of proteins that are *distant* from the drug’s direct targets in the protein-protein interaction network.
71 We find evidence that as long as those proteins affect the same biological functions disrupted by
72 the disease proteins, the treatment can be successful. Thus, physical interactions between proteins
73 alone are unable to explain the therapeutic effects of drugs, and functional information provides an
74 important component for modeling treatment mechanisms. We provide a general framework for
75 identifying proteins and biological functions relevant in treatment, even when drugs seem unrelated
76 to the diseases they are recommended for.

77 **Results**

78 **The multiscale interactome represents the effects of drugs and diseases on proteins and bio-**
79 **logical functions.** The multiscale interactome models drug treatment by integrating both physical
80 interactions between proteins and a multiscale hierarchy of biological functions. Crucially, many
81 treatments depend on biological functions (Supplementary Fig. 3) [20–24]. Existing systematic
82 network approaches, however, primarily model physical interactions between proteins [10–14],
83 and thus cannot accurately model such treatments (Figure 1a, Supplementary Fig. 1).

84 Our multiscale interactome captures the fact that drugs and diseases exert their effects through
85 both proteins and biological functions (Figure 1b). In particular, the multiscale interactome is
86 a network in which 1,661 drugs interact with the human proteins they primarily target (8,568
87 edges) [30,31] and 840 diseases interact with the human proteins they disrupt through genomic al-
88 terations, altered expression, or post-translational modification (25,212 edges) [32]. Subsequently,
89 these protein-level effects propagate in two ways. First, 17,660 proteins physically interact with
90 other proteins according to regulatory, metabolic, kinase-substrate, signaling, and binding rela-
91 tionships (387,626 edges) [33–39]. Second, these proteins alter 9,798 biological functions accord-
92 ing to a rich hierarchy ranging from specific processes (i.e. embryonic heart tube elongation) to
93 broad processes (i.e. heart development). Biological functions can describe processes involving
94 molecules (i.e. DNA demethylation), cells (i.e. the mitotic cell cycle), tissues (i.e. muscle at-
95 rophy), organ systems (i.e. activation of the innate immune response), and the whole organism
96 (i.e. anatomical structure development) (34,777 edges between proteins and biological functions,
97 22,545 edges between biological functions; Gene Ontology) [40,41]. By modeling the effect of
98 drugs and diseases on both proteins and biological functions, our multiscale interactome can model
99 the range of drug treatments that rely on both [20–24].

100 Overall, our multiscale interactome provides a large, systematic dataset to study drug-disease
101 treatments. Nearly 6,000 approved treatments (i.e., drug-disease pairs) spanning almost every
102 category of human anatomy are compiled [31,42,43], exceeding the largest prior network-based
103 study by 10X [13] (Anatomical Therapeutic Classification; Supplementary Fig. 4).

104 **Propagation of the effects of drugs and diseases through the multiscale interactome.** To learn
105 how the effects of drugs and diseases propagate through proteins and biological functions, we
106 harnessed network diffusion profiles (Figure 1c). A network diffusion profile propagates the effects
107 of a drug or disease across the multiscale interactome, revealing the most affected proteins and

108 biological functions. The diffusion profile is computed by biased random walks that start at the
109 drug or disease node. At every step, the walker can restart its walk or jump to an adjacent node
110 based on optimized edge weights. The diffusion profile $\mathbf{r} \in \mathbb{R}^{|V|}$ measures how often each node
111 in the multiscale interactome is visited, thus encoding the effect of the drug or disease on every
112 protein and biological function.

113 Diffusion profiles contribute three methodological advances. First, diffusion profiles provide
114 a general framework to adaptively integrate physical interactions between proteins and a hierarchy
115 of biological functions. When continuing its walk, the random walker jumps between proteins
116 and biological functions at different hierarchical levels based on optimized edge weights. These
117 edge weights encode the relative importance of different types of nodes: w_{drug} , w_{disease} , w_{protein} ,
118 $w_{\text{biological function}}$, $w_{\text{higher-level biological function}}$, $w_{\text{lower-level biological function}}$. These weights are hyperparame-
119 ters which we optimize when predicting the drugs that treat a given disease (Methods). For drug
120 and disease treatments, these optimized edge weights encode the knowledge that proteins and bi-
121 ological functions at different hierarchical levels have different importance in the effects of drugs
122 and diseases [20, 21]. By adaptively integrating both proteins and biological functions in a hierar-
123 chy, therefore, diffusion profiles model effects that rely on both.

124 Second, diffusion profiles provide a mathematical formalization of the principles governing
125 how drug and disease effects propagate in a biological network. Drugs and diseases are known to
126 generate their effects by disrupting or binding to proteins which recursively affect other proteins
127 and biological functions. The effect propagates via two principles [8, 9]. First, proteins and bio-
128 logical functions closer to the drug or disease are affected more strongly. Similarly in diffusion
129 profiles, proteins and biological functions closer to the drug or disease are visited more often since
130 the random walker is more likely to visit them after a restart. Second, the net effect of the drug
131 or disease on any given node depends on the net effect on each neighbor. Similarly in diffusion
132 profiles, a random walker can arrive at a given node from any neighbor.

133 Finally, comparing diffusion profiles provides a rich, interpretable basis to predict pharma-
134 cological properties. Traditional random walk approaches predict properties by measuring the
135 proximity of drug and disease nodes [9]. By contrast, we compare drug and disease diffusion
136 profiles to compare their effects on proteins and biological functions, a richer comparison. Our ap-
137 proach is thus consistent with recent machine learning advances which harness diffusion profiles
138 to represent nodes [44].

139 **The multiscale interactome accurately predicts which drugs treat a disease.** By comparing
140 the similarity of drug and disease diffusion profiles, the multiscale interactome predicts what drugs
141 treat a given disease up to 40% more effectively than molecular-scale interactome approaches
142 (AUROC 0.705 vs. 0.620, +13.7%; Average Precision 0.091 vs. 0.065, +40.0%; Recall@50 0.347
143 vs. 0.264, +31.4%) (Figure 2a, b, Methods). Note that drug-disease treatment relationships are
144 never directly encoded into our network. Instead, the multiscale interactome learns to effectively
145 predict drug-disease treatment relationships it has never posted previously seen.

146 Moreover, the multiscale interactome accurately models classes of drugs that rely on biolog-
147 ical functions and which molecular-scale interactome approaches thus cannot model effectively.
148 Indeed, the top overall performing drug classes (i.e., sex hormones, modulators of the genital
149 system; Supplementary Fig. 6) and the top drug classes for which the multiscale interactome out-
150 performs the molecular-scale interactome (i.e., pituitary, hypothalamic hormones and analogues;
151 Figure 2c, Supplementary Fig. 7) harness biological functions that describe processes across the
152 body. For example, Vasopressin, a pituitary hormone, treats urinary disorders by binding receptors
153 which trigger smooth muscle contraction in the gastrointestinal tract, free water reabsorption in
154 the kidneys, and contraction in the vascular bed [30, 45, 46]. Treatment by Vasopressin, and by
155 pituitary and hypothalamic hormones more broadly, relies on biological functions that describe
156 processes across the body and that are modeled by the multiscale interactome.

157 **The multiscale interactome identifies proteins and biological functions relevant in complex**
158 **treatments.** Existing interactome approaches to systematically study treatment are “black boxes:”
159 they predict what drug treats a disease but cannot explain how the drug treats the disease through
160 specific proteins and biological functions [10–15] (Figure 2d). By contrast, drug and disease dif-
161 fusion profiles identify proteins and biological functions relevant to treatment (Figure 2e, Sup-
162 plementary Note 3). For a given drug and disease, we identify proteins and biological functions
163 relevant to treatment by inducing a subgraph on the k most frequently visited nodes in the drug and
164 disease diffusion profiles which correspond to the proteins and biological functions most affected
165 by the drug and disease.

166 Gene expression signatures validate the biological relevance of diffusion profiles (Figure
167 2f). We find that drugs with more similar diffusion profiles have more similar gene expression
168 signatures (Spearman $\rho = 0.392$, $p = 5.8 \times 10^{-7}$, $n = 152$) [47, 48], indicating that diffusion
169 profiles reflect the effects of drugs on proteins and biological functions.

170 Furthermore, case studies validate the proteins and biological functions that diffusion pro-
171 files identify as relevant to treatment. Consider the treatment of Hyperlipoproteinemia Type III
172 by Rosuvastatin (i.e., Crestor). In Hyperlipoproteinemia Type III, defects in apolipoprotein E
173 (APOE) [49–51] and apolipoprotein A-V (APOA5) [52,53] lead to excess blood cholesterol, even-
174 tually leading to the onset of severe arteriosclerosis [50]. Rosuvastatin is known to treat Hyper-
175 lipoproteinemia Type III by inhibiting HMG-CoA reductase (HMGCR) and thereby diminishing
176 cholesterol production [54,55]. Crucially, diffusion profiles identify proteins and biological func-
177 tions that recapitulate these key steps (Figure 2g). Notably, there is no direct path of proteins
178 between Hyperlipoproteinemia and Rosuvastatin. Instead, treatment operates through biological
179 functions (i.e., cholesterol biosynthesis and its regulation). Consistently, the multiscale interac-
180 tome identifies Rosuvastatin as a treatment for Hyperlipoproteinemia far more effectively than a
181 molecular-scale interactome approach, ranking Rosuvastatin in the top 4.33% of all drugs rather
182 than the top 72.7%. The multiscale interactome explains treatments that rely on biological func-
183 tions, a feat which molecular-scale interactome approaches cannot accomplish.

184 Similarly, consider the treatment of Cryopyrin-Associated Periodic Syndromes (CAPS) by
185 Anakinra. In Cryopyrin-Associated Periodic Syndromes, mutations in NLRP3 and MME lead to
186 immune-mediated inflammation through the Interleukin-1 beta signaling pathway [56]. Anakinra
187 treats Cryopyrin-Associated Syndromes by binding IL1R1, a receptor which mediates regulation of
188 the Interleukin-1 beta signaling pathway and thus prevents excessive inflammation [30,57]. Again,
189 diffusion profiles identify proteins and biological functions that recapitulate these key steps (Fig-
190 ure 2h). Crucially, diffusion profiles identify the regulation of inflammation and immune system
191 signaling, complex biological functions which are not modelled by molecular-scale interactome
192 approaches. Again, the multiscale interactome identifies Anakinra as a treatment for CAPS far
193 more effectively than a molecular-scale interactome approach, ranking Anakinra in the top 10.9%
194 of all drugs rather than the top 71.8%.

195 **The multiscale interactome identifies genes that alter patient-specific drug efficacy and cause**
196 **adverse reactions.** A key goal of precision medicine is to understand how changes in genes
197 alter patient-specific drug efficacy and cause adverse reactions [58] (Figure 3a). For particu-
198 lar treatments, detailed mechanistic models have been developed which can predict and explain
199 drug resistance among genes already identified as relevant to treatment [26–29]. More systemati-
200 cally, however, current tools of precision medicine struggle to predict the genes that interfere with

201 patient-specific treatment [59] and explain how such genes interfere with treatment [60].

202 We find that genetic variants that alter drug efficacy and cause serious adverse reactions occur
203 in genes that are highly visited in the corresponding drug and disease diffusion profiles (Figure
204 3b). We define the treatment importance of a gene according to the visitation frequency of the
205 corresponding protein in the drug and disease diffusion profiles (Methods). Genes that alter drug
206 efficacy and cause adverse reactions exhibit substantially higher treatment importance scores than
207 other genes (median network importance = 0.912 vs. 0.513; $p = 2.95 \times 10^{-107}$, Mood's median
208 test), indicating that these treatment altering genes occur at highly visited nodes. We thus provide
209 evidence that the topological position of a gene influences its ability to alter drug efficacy or cause
210 serious adverse reactions.

211 We find that the network importance of a gene in the drug and disease diffusion profiles pre-
212 dictates whether that gene alters drug efficacy and causes adverse reactions for that particular treat-
213 ment (AUROC = 0.79, Average Precision = 0.82) (Figure 3c). Importantly, the knowledge that a
214 gene alters a given treatment is never directly encoded into our network. Instead, diffusion profiles
215 predict treatment altering relationships that the multiscale interactome has never previously seen.
216 Our diffusion profiles thereby provide a systematic approach to identify genes with the potential
217 to alter treatment. Our finding is complementary to high-resolution, temporal approaches such as
218 discrete dynamic models which model drug resistance and adverse reactions by first curating genes
219 and pathways deemed relevant to a particular treatment [26–29]. Diffusion profiles may help pro-
220 vide candidate genes and pathways for inclusion in these detailed approaches, including genes not
221 previously expected to be relevant. New treatment altering genes, if validated experimentally and
222 clinically, could ultimately affect patient stratification in clinical trials and personalized therapeutic
223 selection [61].

224 Finally, we find that when a gene in a diseased patient alters the efficacy of one indicated drug
225 but not another, that gene primarily targets the genes important to treatment for the resistant drug
226 (Figure 3d, e). Overall, 71.0% of the genes known to alter the efficacy of one indicated drug but not
227 another exhibit higher network importance in the altered treatments than in the unaltered treatment.
228 We thus provide a network formalism explaining how changes to genes can alter efficacy and cause
229 adverse reactions in only some drugs indicated to treat a disease.

230 Consider Benazepril and Diltiazem, two drugs indicated to treat Hypertensive Disease (Fig-
231 ure 3f). A mutation in the AGT gene alters the efficacy of Benazepril but not Diltiazem [62–64].
232 Indeed, our approach gives higher treatment importance to AGT in treatment by Benazepril than in

233 treatment by Diltiazem, ranking AGT as the 45th most important gene for Benazepril treatment but
234 only the 418th most important gene for Diltiazem treatment. Moreover, our approach explains why
235 AGT alters the efficacy of Benazepril but not Diltiazem (Figure 3f). Diltiazem primarily operates
236 at a molecular-scale, inhibiting various calcium receptors (CACNA1S, CACNA1C, CACNA2D1,
237 CACNG1) which trigger relaxation of the smooth muscle lining blood vessels and thus lower blood
238 pressure [30,65–67]. By contrast, Benazepril operates at a systems-scale: Benazepril binds to ACE
239 which affects the renin-angiotensin system, a systems-level biological function that controls blood
240 pressure through hormones [30,68,69]. Crucially, AGT or Angiotensinogen, is a key component of
241 the renin-angiotensin system [69–71]. Therefore, AGT affects the key biological function used by
242 Benazepril to treat Hypertensive Disease. By contrast, AGT plays no role in the calcium receptor
243 driven pathways used by Diltiazem. Thus when a gene alters the efficacy of a drug, the multiscale
244 interactome can identify biological functions that may help explain the alteration in treatment.

245 Discussion

246 The multiscale interactome provides a general approach to systematically understand how drugs
247 treat diseases. By integrating physical interactions and biological functions, the multiscale interac-
248 tome improves prediction of what drugs will treat a disease by up to 40% over physical interactome
249 approaches [10, 13]. Moreover, the multiscale interactome systematically identifies proteins and
250 biological functions relevant to treatment. By contrast, existing systematic network approaches are
251 “black-boxes” which make predictions without providing mechanistic insight. Finally, the mul-
252 tiscale interactome predicts what genes alter drug efficacy or cause severe adverse reactions for
253 drug treatments and identifies biological functions that may explain how these genes interfere with
254 treatment.

255 The multiscale interactome demonstrates that integrating biological functions into the inter-
256 actome improves the systematic modeling of drug-disease treatment. Historically, systematic ap-
257 proaches to study treatment via the interactome have primarily focused on physical interactions be-
258 tween proteins [8–10, 13]. Here, we find that integrating biological functions into a physical inter-
259 actome improves the systematic modeling of nearly 6,000 treatments. We find drugs and drug cate-
260 gories which depend on biological functions for treatment. More broadly, incorporating biological
261 functions may improve systematic approaches that currently use physical interactions to study dis-
262 ease pathogenesis [72–75], disease comorbidities [6], and drug combinations [22–24]. Harnessing
263 the multiscale interactome in these settings may thus help answer key pharmacological questions.
264 Moreover, the multiscale interactome can be readily expanded to add additional node types rele-
265 vant to the problem at hand (i.e. microRNAs to study cancer initiation and progression [76]). Our
266 finding is consistent with systematic studies which demonstrate, in other contexts, that networks
267 involving functional information can strengthen prediction of cellular growth [25, 77], identifica-
268 tion of gene function [78–80], inference of drug targets [81], and general discovery of relationships
269 between biological entities [82, 83].

270 Moreover, we find that diffusion profiles incorporating both proteins and biological functions
271 provide predictive power and interpretability in modeling drug-disease treatments. Diffusion pro-
272 files predict what drugs treat a given disease and identify proteins and biological functions relevant
273 to treatment. In other pharmacological contexts, diffusion profiles incorporating proteins and bi-
274 ological functions may thus improve systematic approaches which currently employ proximity or
275 other non-interpretable methods [6, 16, 17, 33]. In studying the efficacy of drug combinations [17],
276 diffusion profiles may identify synergistic effects on key biological functions. In studying the

277 adverse reactions of drug combinations [16], diffusion profiles may identify biological functions
278 which help explain polypharmacy side effects. In disease comorbidities [6, 33], diffusion profiles
279 may predict new comorbidities and identify biological functions which help explain the develop-
280 ment of the comorbidity.

281 Finally, our study shows that both physical interactions and biological functions can propa-
282 gate the effects of drugs and diseases. We find that many drugs neither directly target the proteins
283 associated with the disease they treat nor target proximal proteins. Instead, these drugs affect the
284 same biological functions disrupted by the disease. This view expands upon the current view of
285 indirect effects embraced in other biological phenomena. In the omnigenic model of complex
286 disease [84, 85], for example, hundreds of genetic variants affect a complex phenotype through
287 indirect effects that propagate through a regulatory network of physical interactions. Our results
288 suggest that the multiscale interactome, incorporating both physical interactions and biological
289 functions, may help propagate indirect effects in complex disease. Altogether, the multiscale in-
290 teractome provides a general computational paradigm for network medicine.

291 **Data availability.** All data used in the paper, including the multiscale interactome, approved
292 drug-disease treatments, drug and disease classifications, gene expression signatures, and pharma-
293 cogenomic relationships is available at github.com/snap-stanford/multiscale-interactome.

294 **Code availability.** Python implementation of our methodology is available at [github.com/snap-](https://github.com/snap-stanford/multiscale-interactome)
295 [stanford/multiscale-interactome](https://github.com/snap-stanford/multiscale-interactome). The code is written in Python. Please read the README for
296 information on downloading and running the code.

297 **Author contributions.** C.R., M.Z., and J.L. designed research; C.R., M.Z., and J.L. performed
298 research; C.R., M.Z., and J.L. analyzed data; and C.R., M.Z., and J.L. wrote the paper.

299 **Corresponding author.** Correspondence should be addressed to J.L. (jure@cs.stanford.edu).

300 **Competing interests.** The authors declare no competing interests.

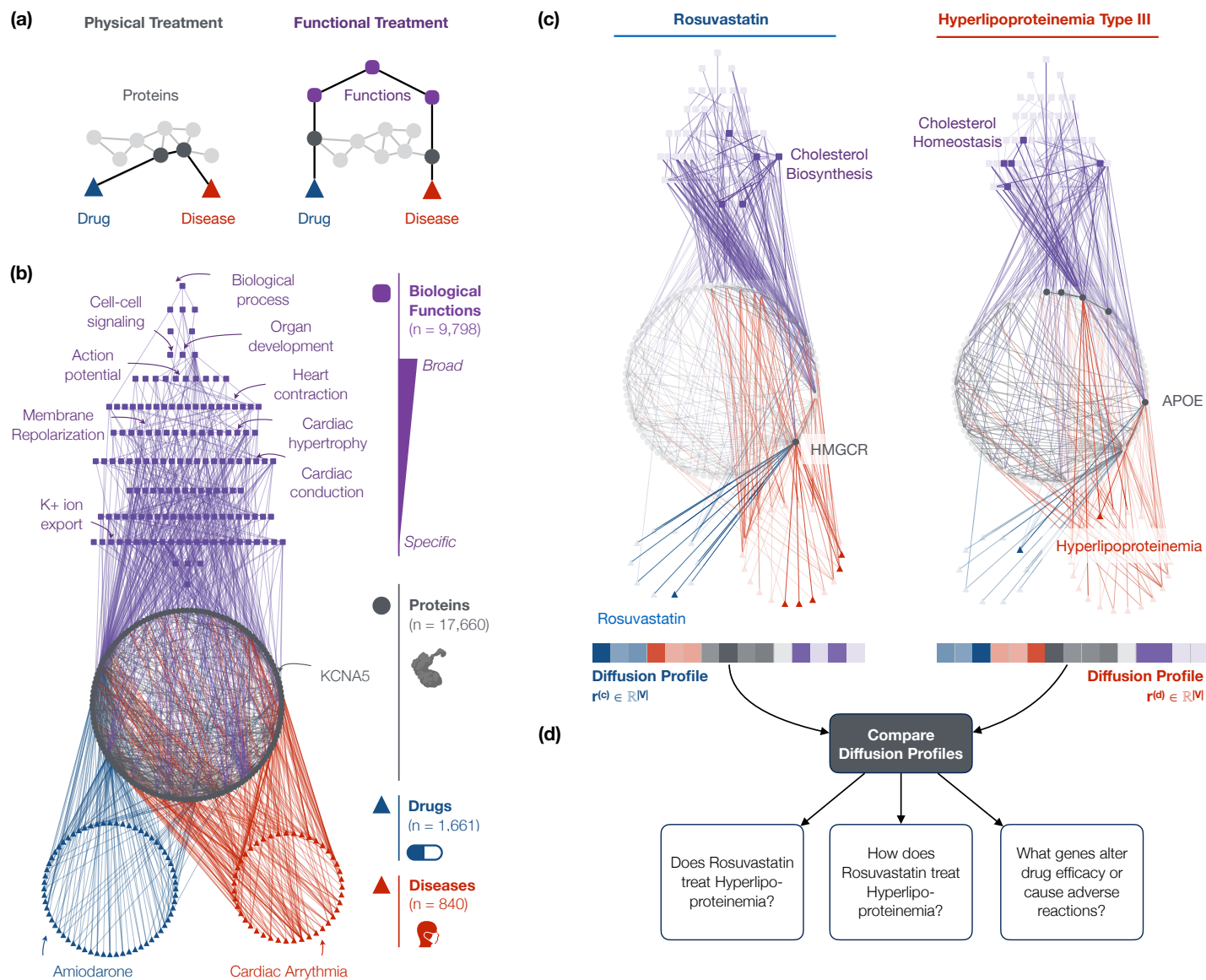


Figure 1: The multiscale interactome models drug treatment through both proteins and biological functions. (a) Existing systematic network approaches assume that drugs treat diseases by targeting proteins that are proximal to disease proteins in a network of physical interactions [10–14]. However, drugs can also treat diseases by targeting distant proteins that affect the same biological functions (Supplementary Fig. 3) [20–25]. (b) The multiscale interactome models drug-disease treatment by integrating both proteins and a hierarchy of biological functions (Supplementary Fig. 1). (c) The diffusion profile of a drug or disease captures its effect on every protein and biological function. The diffusion profile propagates the effect of the drug or disease via random walks which adaptively explore proteins and biological functions based on optimized edge weights. Ultimately, the visitation frequency of a node corresponds to the drug or disease’s propagated effect on that node (Methods). (d) By comparing the diffusion profiles of a drug and disease, we compare their effects on both proteins and biological functions. Thereby, we predict whether the drug treats the disease (Figure 2a-c), identify proteins and biological functions related to treatment (Figure 2d-h), and identify which genes alter drug efficacy or cause dangerous adverse reactions (Figure 3). For example, Hyperlipoproteinemia Type III’s diffusion profile reveals how defects in APOE affect cholesterol homeostasis, a hallmark of the excess blood cholesterol found in patients [49–53]. The diffusion profile of Rosuvastatin, a treatment for Hyperlipoproteinemia Type III, reveals how binding of HMG-CoA Reductase (HMGCR) reduces the production of excess cholesterol [54, 55]. By comparing these diffusion profiles, we thus predict that Rosuvastatin treats Hyperlipoproteinemia Type III, identify the HMGCR and APOE-driven cholesterol metabolic functions relevant to treatment, and predict that mutations in APOE and HMGCR may interfere with treatment and thus alter drug efficacy or cause dangerous adverse reactions.

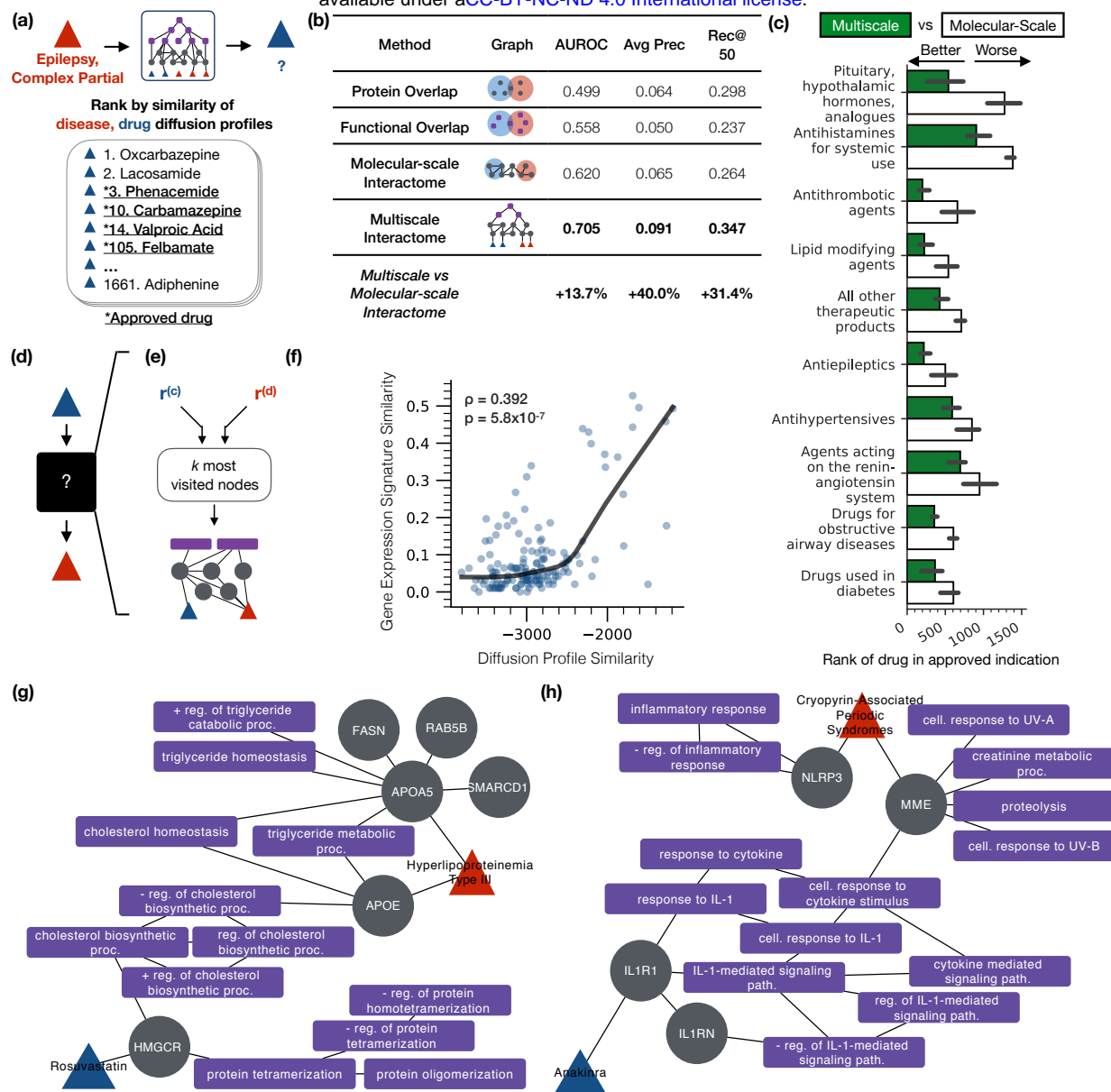


Figure 2: The multiscale interactome accurately predicts what drugs treat a disease and systematically identifies proteins and biological functions related to treatment. (a) To predict whether a drug treats a disease, we compare the drug and disease diffusion profiles according to a correlation distance. (b) By incorporating both proteins and biological functions, the multiscale interactome improves predictions of what drug will treat a given disease by up to 40% over molecular-scale interactome approaches [13]. Reported values are averaged across five-fold cross validation (Methods). (c) The multiscale interactome outperforms the molecular-scale interactome most greatly on drug classes that are known to harness biological functions which describe processes across the body (i.e., pituitary, hypothalamic hormones and analogues; median and 95% CI shown). (d) Existing interactome approaches are “black boxes”: they predict what drug treats a disease but do not explain how the drug treats the disease through specific biological functions [10–15]. (e) By contrast, the diffusion profiles of a drug and disease reveal the proteins and biological functions relevant to treatment. For each drug and disease pair, we induce a subgraph on the k most frequently visited nodes in the drug and disease diffusion profiles to explain treatment. (f) Drugs with more similar diffusion profiles have more similar gene expression signatures (Spearman $\rho = 0.392$, $p = 5.8 \times 10^{-7}$, $n = 152$), suggesting that drug diffusion profiles capture their biological effects. (g) The multiscale interactome explains treatments that molecular-scale interactome approaches cannot faithfully represent. Rosuvastatin treats Hyperlipoproteinemia Type III by binding to HMG CoA reductase (HMGCR) which drives a series of cholesterol biosynthetic functions affected by Hyperlipoproteinemia Type III [49–55]. (h) Anakinra treats Cryopyrin-Associated Periodic Syndromes by binding to IL1R1 which regulates immune-mediated inflammation through the Interleukin-1 beta signaling pathway [30, 57]. Inflammation is a hallmark of Cryopyrin-Associated Periodic Syndromes [56].

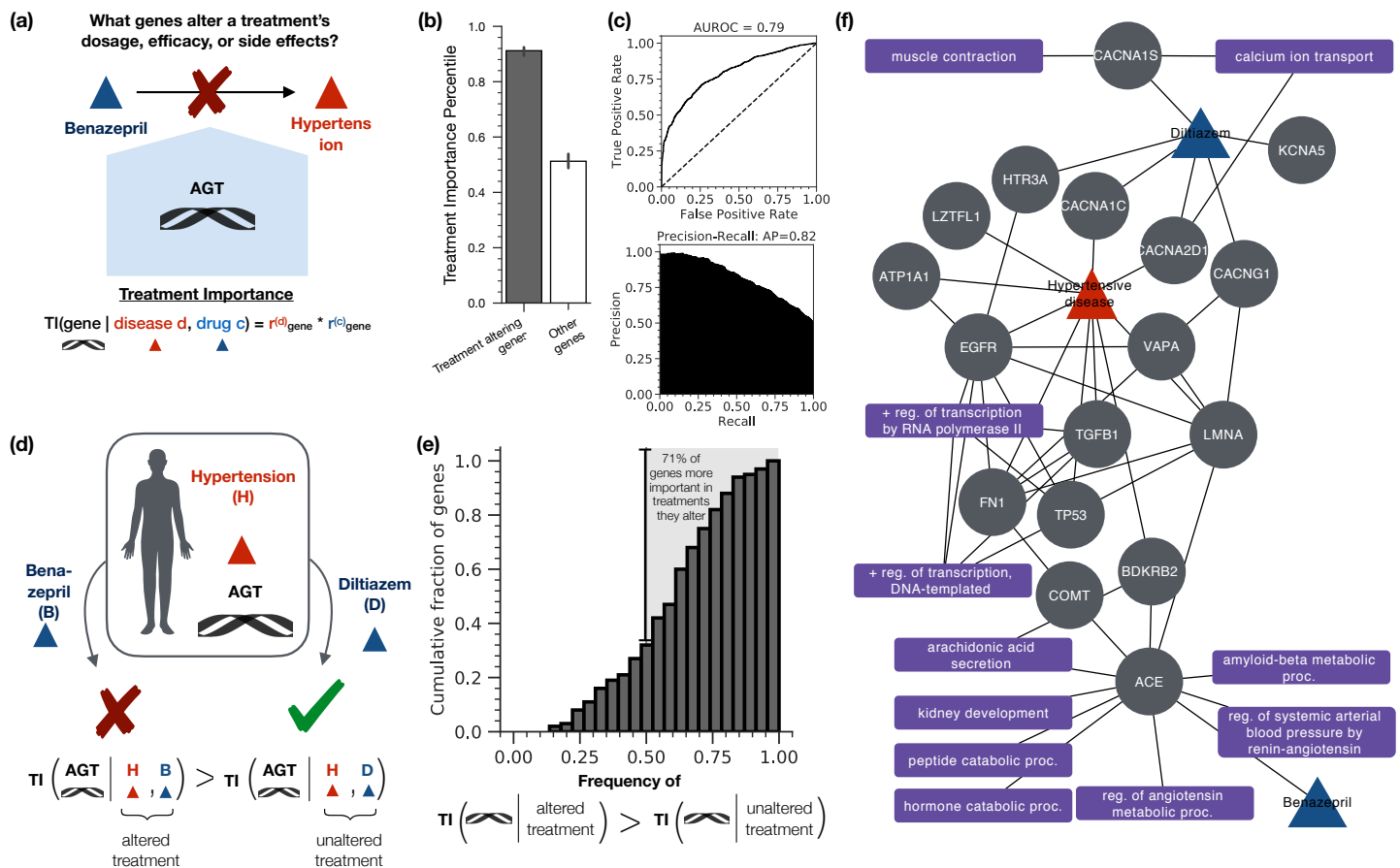


Figure 3: Diffusion profiles identify which genes alter drug efficacy and cause serious adverse reactions and identify biological functions that help explain the alteration in treatment. (a) Genes alter drug efficacy and cause serious adverse reactions in a range of treatments [61]. A pressing need exists to systematically identify genes that alter drug efficacy and cause serious adverse reactions for a given treatment and explain how these genes interfere with treatment [59]. (b) Genetic variants alter drug efficacy and cause serious adverse reactions by targeting genes of high network importance in treatment (median network importance of treatment altering genes = 0.912 vs. 0.513 $p = 2.95 \times 10^{-107}$; Mood's median test; median and 95% CI shown). We define the network treatment importance of a gene according to its visitation frequency in the drug and disease diffusion profiles (Methods). (c) The treatment importance of a gene in the drug and disease diffusion profiles predicts whether that gene alters drug efficacy and causes serious adverse reactions for that particular treatment (AUROC = 0.79, Average Precision = 0.82). (d) Genes uniquely alter efficacy in one indicated drug but not another by primarily targeting the genes and biological functions used in treatment by the affected drug. In patients with Hypertensive Disease, a mutation in AGT alters the efficacy of Benazepril but not Diltiazem. Indeed, AGT exhibits a higher network importance in Benazepril treatment than in Diltiazem treatment, ranked as the 45th most important gene rather than the 418th most important gene. (e) Overall, 71.0% of genes known to alter efficacy in one indicated drug but not another exhibit higher network importance in treatment by the affected drug. (f) Diffusion profiles can identify biological functions that may help explain alterations in treatment. Shown are the proteins and biological functions identified as relevant to the treatment of Hypertensive Disease by Benazepril and Diltiazem. AGT, which uniquely alters the efficacy of Benazepril, is a key regulator of the renin-angiotensin system, a biological function harnessed by Benazepril in treatment but not by Diltiazem [69–71].

301 **Methods**

302 **The multiscale interactome.** The multiscale interactome captures how drugs use both a net-
303 work of physical interactions and a rich hierarchy of biological functions to treat diseases. In
304 the multiscale interactome, 1,661 drugs connect to the proteins they target (8,568 edges) [30, 31].
305 840 diseases connect to the proteins they disrupt through genomic alterations, altered expression,
306 or post-translational modification (25,212 edges) [32]. 17,660 proteins connect to other proteins
307 based on physical interactions such as regulatory, metabolic, kinase-substrate, signaling, or bind-
308 ing relationships (387,626 edges) [33–39]. Proteins connect to the 9,798 biological functions they
309 affect (22,545 edges) [40, 41]. Finally, biological functions connect to each other in a rich hierar-
310 chy ranging from specific processes (i.e. embryonic heart tube elongation) to broad processes (i.e.
311 heart development) (22,545 edges) [40, 41]. Biological functions can describe processes involving
312 molecules (i.e. DNA demethylation), cells (i.e. the mitotic cell cycle), tissues (i.e. muscle atro-
313 phy), organ systems (i.e. activation of the innate immune response), and the whole organism (i.e.
314 anatomical structure development).

315 We visualize a representative subset of the multiscale interactome using Cytoscape [86] (Fig-
316 ure 1b).

317 **Drug–protein interactions.** We map drugs to their protein targets using DrugBank [30] and the
318 Drug Repurposing Hub [31]. For DrugBank, we map the Uniprot Protein IDs to Entrez IDs using
319 HUGO [87]. For the Drug Repurposing Hub, we map drugs to their DrugBank IDs using the drug
320 names and DrugBank’s “drugbank_approved_target_uniprot_links.csv” file. We map protein targets
321 to Entrez IDs using HUGO [87]. We filter drug-target relationships to only include proteins that are
322 represented in the network of physical interactions between proteins (see Methods: Protein–protein
323 interactions). All drug-target interactions are provided in Supplementary Data 1.

324 **Disease–protein interactions.** We map diseases to genes they affect through genomic alterations,
325 altered expression, or post-translational modification by using DisGeNet [32]. To ensure high-
326 quality disease-gene associations, we only consider the “curated” set of disease-gene associations
327 provided by DisGeNet which draws from expert-curated repositories: UniProt, the Comparative
328 Toxicogenomics Database, Orphanet, the Clinical Genome Resource (ClinGen), Genomics Eng-
329 land PanelApp, the Cancer Genome Interpreter (CGI), and the Psychiatric Disorders Gene Asso-
330 ciation Network (PsyGeNET). We exclude all disease-gene associations that are inferred, based
331 on orthology relationships from animal models, or based on computational-mining of the litera-

332 ture. Ultimately, diseases are associated with genes they affect via genomic alteration, alteration
333 of expression, or post-translational modification according to the DisGeNet relationship ontology.
334 To avoid circularity in the analysis, we remove disease-gene associations marked as therapeutic.
335 Finally, we filter disease-gene relationships to only consider genes whose protein products were
336 present in the network of physical interactions between proteins (see Methods: Protein-protein
337 interactions). All disease-protein interactions are provided in Supplementary Data 2.

338 **Protein-protein interactions.** We generate a network of 387,626 physical interactions between
339 17,660 proteins by compiling seven major databases. Across all databases, we only consider hu-
340 man proteins and their interactions; only allow protein-protein interactions with direct experi-
341 mental evidence; and only allow *physical* interactions between proteins, filtering out genetic and
342 indirect interactions between proteins such as those identified via synthetic lethality experiments.
343 All protein-protein interactions are provided in Supplementary Data 3.

344 1. *The Biological General Repository for Interaction Datasets [34]* (BioGRID; 309,187 in-
345 teractions between 16,352 proteins). BioGRID manually curates both physical and genetic
346 interactions between proteins from 71,713 high- and low-throughput publications. We map
347 BioGRID proteins to Entrez IDs by using HUGO [87]. We only include protein-protein
348 interactions from BioGRID that result from experiments indicating a *physical* interaction
349 between the proteins, as described by BioGRID [34], and ignore protein-protein interactions
350 indicating a *genetic* interaction between the proteins. We use the “BIOGRID-ORGANISM-
351 Homo_sapiens-3.5.178.tab” file.

352 2. *The Database of Interacting Proteins [36]* (DIP; 4,235 interactions between 2,751 proteins).
353 DIP only considers physical protein-protein interactions with experimental evidence and
354 curates these from the literature. We map the UniProt ID of each protein to its Entrez ID
355 by using HUGO [87]. We allow all experimental methods from DIP since they all capture
356 physical interactions [36].

357 3. *The Human Reference Protein Interactome Mapping Project.* We integrate four protein-
358 protein interaction networks from the Human Reference Protein Interactome Mapping
359 Project that were generated through high-throughput yeast two hybrid assays (HI-I-05 [39]:
360 2,611 interactions between 1,522 proteins; HI-II-14 [35] 13,426 interactions between 4,228
361 proteins; Venkatesan-09 [37]: 233 interactions between 229 proteins; Yu-11 [38] 1,126 in-

362 teractions between 1,126 proteins). Since protein-protein interactions in all four networks
363 result from a yeast two-hybrid system, all protein-protein interactions are physical and ex-
364 perimentally verified. We thus include all protein-protein interactions across these networks.
365 Proteins are already provided with their Entrez ID so no mapping is required.

366 4. *Menche-2015* [33] (138,425 interactions between 13,393 proteins). Finally, we integrate
367 the physical protein-protein interaction network compiled by Menche et al. [33]. Menche
368 et al. compiles different types of physical protein-protein interactions from a range of
369 sources. In all cases, protein-protein interactions result from direct experimental evidence.
370 Menche et al. compiles regulatory interactions from the TRANSFAC database; binary inter-
371 actions from a series of high-throughput yeast-two-hybrid datasets as well as the IntAct and
372 MINT databases; literature curated interactions from IntAct, MINT, BioGRID, and HPRD;
373 metabolic-enzyme coupled interactions from KEGG and BIGG; protein complex interac-
374 tions from CORUM; kinase-substrate interactions from PhosphositePlus; and signaling in-
375 teractions from [88]. All proteins are provided in Entrez format and thus do not require
376 further mapping.

377 **Protein – biological function interactions.** We map proteins to the biological functions they
378 affect by using the human version of the Gene Ontology [40, 41] (7,993 proteins; 6,387 biologi-
379 cal functions; 34,777 edges). We only allow experimentally verified associations between genes
380 and biological functions according to the following IDs: EXP – inferred from experiment, IDA
381 – inferred from direct assay, IMP – inferred from mutant phenotype, IGI – inferred from genetic
382 interaction, HTP – high throughput experiment, HDA – high throughput direct assay, HMP – high
383 throughput mutant phenotype, and HGI – high throughput genetic interaction. We exclude any
384 protein–biological function relationships that are inferred from physical interactions to avoid re-
385 dundancy with the physical network of interacting proteins. We also exclude protein–biological
386 function relationships inferred from gene expression patterns since the Gene Ontology states that
387 such interactions are challenging to map to specific proteins [40, 41]. To prevent circularity, we
388 further ignore all associations based on phylogenetically inferred annotations or various compu-
389 tational analyses (sequence or structural similarity, sequence orthology, sequence alignment, se-
390 quence modeling, genomic context, reviewed computational analysis). Finally, we ignore associ-
391 ations based on author statements, curator inference, electronic annotations (i.e. automated anno-
392 tations), and those for which no biological data was available. Some biological functions in the

393 Gene Ontology have multiple synonymous IDs. For each biological function, we use the “master
394 IDs” provided by GOATOOLS [89]. All protein – biological function interactions are provided in
395 Supplementary Data 4.

396 **Biological function – biological function interactions.** We construct a hierarchy of biological
397 functions by using the Gene Ontology’s Biological Processes [40, 41]. The Gene Ontology rep-
398 represents a curated hierarchy of biological functions, where highly specific biological functions are
399 children of more general biological functions according to numerous relationship types. For ex-
400 ample, “negative regulation of response to interferon-gamma” $\xrightarrow{\text{is a}}$ “negative regulation of innate
401 immune response” $\xrightarrow{\text{is a}}$ “negative regulation of immune response” $\xrightarrow{\text{negatively regulates}}$ “immune re-
402 sponse.” We allow relationships between biological functions of the following types: regulates,
403 positively regulates, negatively regulates, part of, and is a. In order to allow the model to focus on
404 the biological functions most relevant to treatment, we only consider biological functions which
405 are associated with at least one drug target or one disease protein, either directly or implicitly
406 through their children. All biological function – biological function interactions are provided in
407 Supplementary Data 5.

408 **Constructing dataset of approved drug-disease treatments.** We construct a dataset of 5,926
409 unique, approved drug-disease pairs, exceeding the largest prior network-based study by 10X [13].
410 We source approved drug-disease pairs from the Drug Repurposing Database [42] ($n_{\text{pairs}} = 2,538$;
411 $n_{\text{drugs}} = 996$, $n_{\text{diseases}} = 463$), the Drug Repurposing Hub [31] ($n_{\text{pairs}} = 1,449$; $n_{\text{drugs}} =$
412 908 , $n_{\text{diseases}} = 265$), and the Drug Indication Database [43] ($n_{\text{pairs}} = 3,304$; $n_{\text{drugs}} = 1,147$,
413 $n_{\text{diseases}} = 615$). In all cases, we filter drug-disease pairs to ensure that only FDA-approved
414 treatment relationships are included.

415 We extract approved drug-disease pairs from each database as follows. In all cases, drugs are
416 mapped to DrugBank IDs [30] and diseases are mapped to unique identifiers from the National Li-
417 brary of Medicine [90] (NLM UMLS CUIDs: NLM Unified Medical Language System Controlled
418 Unique Identifier):

419 1. *The Drug Repurposing Database* is a gold-standard database of drug-disease pairs extracted
420 from drug labels and the American Association of Clinical Trials Database [42]. Drugs
421 and diseases in the Drug Repurposing Database are provided with DrugBank IDs and NLM
422 UMLS CUIDs so no additional mapping is required. We extract only the drug and disease
423 pairs designated as “Approved” treatment relationships.

424 2. *The Broad Institute's Drug Repurposing Hub* is a hand-curated collection of drug-disease
425 pairs compiled from drug labels, DrugBank, the NCATS NCGC Pharmaceutical Collection
426 (NPC), Thomson Reuters Integrity, Thomson Reuters Cortellis, Citeline Pharmaprojects,
427 the FDA Orange Book, ClinicalTrials.gov, and PubMed [31]. We map drugs to DrugBank
428 IDs by comparing their provided names and PubChem IDs to DrugBank's external links
429 mapping [30]. We map diseases to UMLS CUIDs by using the UMLS Metathesaurus's
430 REST API [90]. Finally, we only include drug-disease pairs with a "Launched" clinical
431 phase attribute, indicating FDA approval.

432 3. *The Drug Indication Database* provides drug-indications relationships from DailyMed,
433 DrugBank, the Pharmacological Actions sections of the Medical Subject Headings, the Na-
434 tional Drug File Reference Terminology, the Physicians' Desk Reference, the Chemical En-
435 tities of Biological Interest (ChEBI), the Comparative Toxicogenomics Database, the Ther-
436 apeutic Claims section of the USP Dictionary of United States Adopted Names and Inter-
437 national Drug Names, and the World Health Organization Anatomic-Therapeutic-Chemical
438 classification) [43]. The Drug Indication Database captures both diseases and non-disease
439 medical conditions (i.e. pregnancy) for which a drug is used. Additionally, the Drug In-
440 dication Database captures both treatment relationships between drugs and indications as
441 well as prevention, management, and diagnostic relationships. We filter the Drug Indication
442 Database to only include *approved* treatment relationships between drugs and *diseases*.

443 We map drugs to DrugBank IDs by using the provided CAS and ChEBI IDs as well as Drug-
444 Bank's external links mapping [30]. Indications are already provided with UMLS CUIDs.

445 We filter indications to only include diseases in two ways. First, we only consider indica-
446 tions with a UMLS semantic type of "B2.2.1.2.1 Disease or Syndrome", "B2.2.1.2 Patho-
447 logic Function", or "B2.2.1.2.1.2 Neoplastic Process." Second, we only consider indications
448 present in DisGeNet, a database mapping diseases to their associated genes [32].

449 To ensure that drug-disease relationships specifically represent treatment relationships, we
450 filter drug-disease pairs based on the "indication subtype." We remove drug-indication pairs
451 where the indication subtype described is not treatment (i.e. preventative/prophylaxis, di-
452 agnosis, adjunct, palliative, reduction, causes/inducing/associated, and mechanism). We ad-
453 ditionally remove all drug indication pairs from the Comparative Toxicogenomics Database
454 (CTD). The goal of CTD is to provide broad chemical-disease associations published in the

455 literature [91]. Concurrently, CTD does not subset these chemical-disease associations into
456 drug-disease relationships that represent FDA-approved treatments.

457 Finally, we remove overly broad diseases from the Drug Indication Database. We remove
458 disease categories (i.e. diseases with “Diseases” in their name such as “Cardiovascular Dis-
459 eases” and “Metabolic Diseases”). We also remove diseases with more than 130 approved
460 drugs (i.e. Disorder of Eye – 290 approved drugs).

461 After compiling approved drug-disease treatment pairs, we remove treatments for which
462 drugs rely on binding to non-human proteins (i.e. viral or bacterial proteins) to induce their effect.
463 The multiscale interactome only models human proteins and biological functions. The multiscale
464 interactome is thus not designed to model treatments which rely on binding to viral or bacterial
465 proteins. To remove such treatments, we map all disease UMLS CUIDs to their corresponding Dis-
466 ease Ontology ID [92]. We then remove diseases corresponding to the “disease by infectious agent
467 category” of the Disease Ontology. The Disease Ontology does not map many UMLS CUIDs to
468 corresponding Disease Ontology IDs. We thus manually curate the final list of diseases to remove
469 additional infectious diseases: malaria, bacterial septicemia, fungal infection, coccidiosis, gon-
470 orrhea, gastrointestinal roundworms, shingles, lice, gastrointestinal parasites, tapeworm, syphilis,
471 genital herpes, lungworms, fungicide, fungal keratosis, yeast infection, laryngitis, enterocolitis,
472 protozoan infection, African trypanosomiasis, sepsis, Chagas disease, mites, bacterial vaginosis,
473 scabies, pinworm, equine protozoal myeloencephalitis (EPM), microsporidiosis, and ringworm.

474 Finally, we filter approved drug-disease treatment pairs to only include drugs with at least one
475 known target in DrugBank [30] or the Drug Repurposing Hub [31] and diseases with at least one as-
476 sociated gene in the curated version of DisGeNet [32] as these are the only drugs and diseases that
477 the multiscale interactome represents (see Methods: Drug–protein interactions, Disease–protein
478 interactions).

479 Ultimately, we achieve a dataset of 5,926 approved drug-disease pairs, exceeding the largest
480 prior network-based study by 10X [13]. All approved drug-disease pairs are provided in Supple-
481 mentary Data 6.

482 **Learning drug and disease diffusion profiles.** We propagate the effects of each drug and disease
483 across the multiscale interactome by using network diffusion profiles. A drug or disease diffusion
484 profile learns the proteins and biological functions most affected by each drug or disease. Each
485 drug or disease diffusion profile is computed through biased random walks that start at the drug or

486 disease node. At every step, the random walker can restart its walk or jump to an adjacent node
487 based on optimized edge weights. After many walks, the diffusion profile measures how often
488 every node was visited, thus representing the effect of the drug or disease on that node.

489 By using optimized edge weights, diffusion profiles learn to adaptively inte-
490 grate proteins and biological functions. Diffusion profiles rely on a set of scalar
491 weights which encode the relative importance of different types of nodes: $W =$
492 $\{w_{\text{drug}}, w_{\text{disease}}, w_{\text{protein}}, w_{\text{biological function}}, w_{\text{higher-level biological function}}, w_{\text{lower-level biological function}}\}$. These
493 weights are hyperparameters which we optimize when predicting the drugs that treat a given
494 disease (see Methods: Model selection and optimization of scalar weights). When a random
495 walker continues its walk, it picks the next node to jump to based on the relative values of
496 these weights. For example, if a random walker is at a protein and has both protein and
497 biological function neighbors, it is $\frac{w_{\text{protein}}}{w_{\text{biological function}}}$ times more likely to jump to the protein neigh-
498 bors than the biological function neighbors. Notice that proteins connect to drugs, diseases,
499 proteins, and biological functions, making $\{w_{\text{drug}}, w_{\text{disease}}, w_{\text{protein}}, w_{\text{biological function}}\}$ the relevant
500 weights for a random walker currently at a protein. By contrast, biological functions connect
501 to proteins, higher-level biological functions, and lower-level biological functions, making
502 $\{w_{\text{protein}}, w_{\text{higher-level biological function}}, w_{\text{lower-level biological function}}\}$ the relevant weights for a random walker
503 at a biological function. By providing separate weights for higher- and lower-level biological
504 functions, the random walker learns to explore different levels of the hierarchy of biological
505 functions and integrate them appropriately.

506 Diffusion profiles represent a general methodology to propagate signals through a hetero-
507 geneous biological network. By carefully defining edge weights and the nodes that the random
508 walker restarts to, diffusion profiles can be used in a wide range of biological tasks. Here, we de-
509 fine edge weights for drug, disease, protein, and biological function node types, yet more or fewer
510 weights can be used based on the problem of interest. Similarly, here, the random walker jumps
511 to the initial drug or disease node after a restart, but in reality, it can restart to any node or any set
512 of nodes. The edge weights and restart nodes thus make diffusion profiles a flexible approach to
513 propagate signals across a heterogeneous biological network, with applicability to a wide range of
514 problems in systems biology and pharmacology.

515 **Computing drug and disease diffusion profiles through power iteration.** Mathematically, we
516 compute diffusion profiles through a matrix formulation with power iteration [93–95]. The diffu-

517 sion profile computation takes as input:

- 518 1. $G = (V, E)$ the unweighted, undirected multiscale interactome with V nodes and E edges.
- 519 2. $W = \{w_{\text{drug}}, w_{\text{disease}}, w_{\text{protein}}, w_{\text{biological function}}, w_{\text{higher-level biological function}}, w_{\text{lower-level biological function}}\}$
520 the set of scalar weights which encode the relative likelihood of the walker jumping from
521 one node type to another when continuing its walk.
- 522 3. α which represents the probability of the walker continuing its walk at a given step rather
523 than restarting.
- 524 4. $\mathbf{s} \in \mathbb{R}^{|V|}$ a restart vector which sets the probability the walker will jump to each node after a
525 restart; here, \mathbf{s} is a one-hot vector encoding the drug or disease of interest.
- 526 5. ϵ the tolerance allowed for convergence of the power iteration computation.

527 The diffusion profile computation outputs $\mathbf{r} \in \mathbb{R}^{|V|}$, a drug- or disease-diffusion profile which
528 measures the frequency with which the random walker visits each node. Note that $\sum_i \mathbf{r}_i = 1$.

529 Before computing the diffusion profile of a drug or disease of interest, we prepro-
530 cess the multiscale interactome in order to only allow biologically meaningful walks. Dif-
531 fusion profiles are designed to capture how a drug or disease of interest propagates its ef-
532 fect by recursively affecting proteins and biological functions. Notice that drugs and dis-
533 eases do not propagate their effect by using other drugs and diseases as intermediates.
534 Therefore, we disallow paths that have drugs and diseases as intermediate nodes. To ac-
535 complish this mathematically, we convert $G = (V, E)$ to a directed graph G' where all
536 previously undirected edges are replaced by edges in both directions (i.e. edges now
537 include drug \leftrightarrow protein, disease \leftrightarrow protein, protein \leftrightarrow protein, protein \leftrightarrow biological function, and
538 lower-level biological function \leftrightarrow higher-level biological function). We then make the drug or dis-
539 ease of interest a source node (i.e. no in-edges) and all other drugs and diseases sink nodes (i.e. no
540 out-edges). In G' , a random walker starts at the drug or disease of interest and recursively walks
541 to proteins and biological functions. If the walker reaches any other drug or disease node, it must
542 restart its walk.

543 Next, we encode G' and the set of scalar weights W into a biased transition matrix $\mathbf{M} \in$
544 $\mathbb{R}^{|V| \times |V|}$. Each entry \mathbf{M}_{ij} denotes the probability $p_{i \rightarrow j}$ a random walker jumps from node i to node
545 j when continuing its walk. Consider a random walker at node i jumping to neighbor j of type t .
546 Let T be the set of all node types adjacent to node i . We compute $p_{i \rightarrow j}$ in two steps.

1. First, we compute the probability of the random walker jumping to a node of type t rather than a node of a different type. w_t is the weight of node type t as specified in W :

$$p_t = \frac{w_t}{\sum_{t' \in T} w_{t'}}.$$

2. Second, we compute the probability that the random walker jumps to node j rather than to another adjacent node of type t . Let n_t be the number of adjacent nodes of type t :

$$\mathbf{M}_{ij} = p_{i \rightarrow j} = \frac{p_t}{n_t}.$$

After constructing \mathbf{M} , we finally compute the diffusion profile through power iteration as shown in Algorithm 1. The key equation is:

$$\mathbf{r}^{(k+1)} = \overbrace{(1 - \alpha)\mathbf{s}}^{\text{Restart walk}} + \alpha \left(\underbrace{\mathbf{r}^{(k)}\mathbf{M}}_{\text{from node with out-edges}} + \underbrace{\mathbf{s} \sum_{j \in J} \mathbf{r}_j^{(k)}}_{\text{from node without out-edges}} \right).$$

547 At each step, the random walker can restart its walk at the drug or disease node according to
 548 $(1 - \alpha)\mathbf{s}$ or continue its walk. If the random walker continues its walk from a node with out-edges,
 549 then it jumps to an adjacent node according to $\alpha(\mathbf{r}^{(k)}\mathbf{M})$. If the random walker continues its walk
 550 from a node without out-edges (i.e. a sink node), then it restarts its walk according to $\alpha(\mathbf{s} \sum_{j \in J} \mathbf{r}_j^{(k)})$,
 551 where J is the set of sink nodes in the graph. At every iteration, $\sum_i \mathbf{r}_i = 1$.

552 Code for the power iteration implementation is available at [github.com/snap-](https://github.com/snapstanford/multiscale-interactome)
 553 [stanford/multiscale-interactome](https://github.com/snapstanford/multiscale-interactome). We use a tolerance of $\epsilon = 1 \times 10^{-6}$.

Algorithm 1 Diffusion profiles through power iteration

```

% Initialize diffusion profile
 $\mathbf{r}_i^{(0)} = \frac{1}{|V|} \forall i$ 
% While not converged
while  $\|\mathbf{r}^{(k+1)} - \mathbf{r}^{(k)}\|_1 > \epsilon$  do
    % Start new walk at drug or disease node or continue walk.
     $\mathbf{r}^{(k+1)} = (1 - \alpha)\mathbf{s} + \alpha(\mathbf{r}^{(k)}\mathbf{M} + \mathbf{s} \sum_{j \in J} \mathbf{r}_j^{(k)})$ 
end while

```

554 **Predicting what drugs will treat a given disease with diffusion profiles.** For a drug to treat a
555 disease, it must affect proteins and biological functions similar to those disrupted by the disease.
556 The diffusion profiles of the drug $\mathbf{r}^{(c)}$ and the disease $\mathbf{r}^{(d)}$ encode the effect of the drug and the
557 disease on proteins and biological functions. Therefore, comparing $\mathbf{r}^{(c)}$ and $\mathbf{r}^{(d)}$ allows us to predict
558 what drugs treat a given disease.

559 For each drug and each disease, we compute the diffusion profile as described above. For
560 each disease, we then rank-order the drugs most likely to treat the disease based on the similarity
561 of the drug and disease diffusion profiles $\text{SIM}(\mathbf{r}^{(c)}, \mathbf{r}^{(d)})$ and a series of baseline methods.

562 We test five metrics of vector similarity:

- 563 1. L2 norm: $\sqrt{\sum_i |\mathbf{r}_i^{(c)} - \mathbf{r}_i^{(d)}|^2}$,
- 564 2. L1 norm: $\sum_i |\mathbf{r}_i^{(c)} - \mathbf{r}_i^{(d)}|$,
- 565 3. Canberra distance: $\sum_i \frac{|\mathbf{r}_i^{(c)} - \mathbf{r}_i^{(d)}|}{|\mathbf{r}_i^{(c)}| + |\mathbf{r}_i^{(d)}|}$,
- 566 4. Cosine similarity: $\frac{\mathbf{r}^{(c)} \cdot \mathbf{r}^{(d)}}{\|\mathbf{r}^{(c)}\|_2 \|\mathbf{r}^{(d)}\|_2}$,
- 567 5. Correlation distance: $1 - \frac{(\mathbf{r}^{(c)} - \bar{\mathbf{r}}^{(c)}) \cdot (\mathbf{r}^{(d)} - \bar{\mathbf{r}}^{(d)})}{\|(\mathbf{r}^{(c)} - \bar{\mathbf{r}}^{(c)})\|_2 \|(\mathbf{r}^{(d)} - \bar{\mathbf{r}}^{(d)})\|_2}$.

568 We additionally test two proximity metrics. In particular, we consider the visitation fre-
569 quency of the drug node i in the disease diffusion profile as: $\mathbf{r}_i^{(d)}$. We also consider the visitation
570 frequency of the drug node i in the disease diffusion profile multiplied by the visitation frequency
571 of the disease node j in the drug diffusion profile: $\mathbf{r}_i^{(d)} * \mathbf{r}_j^{(c)}$.

572 **Baseline metrics to predict what drugs will treat a disease.** To predict what drugs will treat a
573 given disease, we consider baselines that measure (1) the overlap between drug targets and disease

574 proteins, (2) the overlap between the functions of drug targets and disease proteins, and (3) the
575 state-of-the-art proximity metric on a molecular-scale interactome. First, we compute the “protein
576 overlap” baseline which we define as the Jaccard Similarity between the set of drug targets T and
577 the set of disease proteins S : $\frac{|T \cap S|}{|T \cup S|}$. Second, we compute the “functional overlap” baseline which
578 we define as SimIC which measures the semantic similarity between the GO terms U associated
579 with the drug targets and the GO terms V associated with the disease proteins [96]. We tested 17
580 functional overlap baselines, of which this was the best performing (Methods: Baseline metrics
581 of functional overlap between drug targets and disease proteins) (Supplementary Fig. 5). Third,
582 we compute the state-of-the-art proximity metric on a molecular-scale interactome which is the
583 closest distance metric in [10, 13]. Let T be the set of drug targets, S be the set of disease proteins,
584 and $l(s, t)$ be the shortest path length between nodes s and t . The state-of-the-art proximity metric
585 first computes the “closest” distance $d(S, T) = \frac{1}{|T|} \sum_{t \in T} \min_{s \in S} l(s, t)$ between S and T . Next,
586 this distance is compared to a reference distance distribution which measures $d(S, T)$ when S and
587 T are randomly permuted to sets of proteins that match the size and degrees of the original disease
588 proteins and drug targets in the network. Finally, the state-of-the-art proximity metric is computed
589 by taking a z-score of $d(S, T)$ with respect to the reference distribution: $z(S, T) = \frac{d(S, T) - \mu_{d(S, T)}}{\sigma_{d(S, T)}}$.

590 **Baseline metrics of functional overlap between drug targets and disease proteins.** We tested
591 17 baseline methods that predict what drugs treat a disease by considering the biological functions
592 affected by drug targets and disease proteins (Supplementary Fig. 5).

593 First, we tested baseline methods that compare the functional overlap between drug targets
594 and disease proteins. Let U and V be the sets of Gene Ontology (GO) terms associated with drug
595 targets and disease proteins respectively. Let U' and V' be the multisets of GO terms associated
596 with drug targets and disease proteins respectively. Let U'' and V'' be the sets of GO terms enriched
597 among drug targets and disease proteins according to Gene Set Enrichment Analysis (GSEA) re-
598 spectively [89, 97]. Note that in the multisets U' and V' , U'_i and V'_i correspond to the number of
599 occurrences of the i^{th} element in the multiset.

600 We measure the following baselines:

- 601 • The Jaccard Similarity or Intersection between the set of GO terms associated with the drug
602 targets and the set of GO terms associated with the disease proteins: $\frac{|U \cap V|}{|U \cup V|}$ or $|U \cap V|$
- 603 • The Jaccard Similarity or Intersection between the multiset of GO terms associated with the
604 drug targets and the multiset of GO terms associated with the disease proteins: $\frac{\sum_i \min(U'_i, V'_i)}{\sum_i \max(U'_i, V'_i)}$

605 or $\sum_i \min(U'_i, V'_i)$

- 606 • The Jaccard Similarity or Intersection between the set of GO terms enriched among drug
607 targets and the set of GO terms enriched among disease proteins according to Gene Set
608 Enrichment Analysis [89, 97]: $\frac{|U'' \cap V''|}{|U'' \cup V''|}$ or $|U'' \cap V''|$
- 609 • The Z-scored Jaccard Similarity or Intersection between the set of GO terms associated with
610 the drug targets and the set of GO terms associated with the disease proteins: $z\left(\frac{|U \cap V|}{|U \cup V|}\right)$ or
611 $z(|U \cap V|)$
- 612 • The Z-scored Jaccard Similarity or Intersection between the multisets of GO terms asso-
613 ciated with the drug targets and the set of GO terms associated with the disease proteins:
614 $z\left(\frac{\sum_i \min(U'_i, V'_i)}{\sum_i \max(U'_i, V'_i)}\right)$ or $z\left(\sum_i \min(U'_i, V'_i)\right)$

615 We compute reference distributions for z-scored metrics by following the approach in [10,
616 13]. Specifically, we randomly permute the set of disease proteins S and the set of drug targets T
617 to sets of proteins that match the size and degrees of the original disease proteins and drug targets
618 in the network. We then generate the GO sets and multisets that correspond to the permuted S and
619 T , compute the relevant baseline metric, and repeat this for random permutations of S and T to
620 generate a reference distribution. Finally, we compute a z-score by comparing the baseline metric
621 for the true S and T to the reference distribution.

622 Second, we tested baseline methods that calculate the semantic similarity between the GO
623 terms associated with the drug targets and those associated with the disease proteins [98]. Consider
624 U and V , the sets of GO terms directly associated with drug targets and disease proteins respec-
625 tively. Semantic similarity methods first define a similarity $\text{sim}(u, v)$ between a GO term directly
626 associated with drug targets u and a GO term directly associated with disease proteins v . The
627 similarity of the sets U and V are subsequently calculated by aggregating across the similarities of
628 pairwise GO terms u and v .

629 We used the following semantic similarity metrics as as they are among the most common
630 and best-performing metrics in a variety of settings [98].

- 631 • The Resnik Similarity [99, 100] between u and v measures the information content of
632 the most informative common ancestor between u and v . $\text{sim}(u, v) = \text{Resnik}(u, v) =$
633 $\text{IC}[\text{MICA}(u, v)]$

634 – Let $p(u)$ be the fraction of proteins in the multiscale interactome that are associated
635 with a GO term u or its descendants. The information content IC of term u is defined as
636 $IC(u) = -\log[p(u)]$. The Maximum Informative Common Ancestor (MICA) between
637 two GO terms u and v is defined as $MICA(u, v) = \operatorname{argmax}_{x \in \text{ancestors}(u, v)} IC(x)$.

638 • simIC [96] integrates both the information content of GO terms and the structural in-
639 formation of the GO hierarchy to determine the similarity between GO terms u and v :

640
$$\text{sim}(u, v) = \text{simIC}(u, v) = \frac{2 \log[p(MICA(u, v))]}{\log[p(u)] + \log[p(v)]} \left(1 - \frac{1}{1 + IC[MICA(u, v)]}\right)$$

641 • simGIC [101] which considers the information content of all common ancestors of the GO
642 terms directly associated with the drug targets U and the GO terms directly associated with
643 the disease proteins V . $\text{sim}(u, v) = \text{simGIC}(U, V) = \frac{\sum_{x \in A(U) \cap A(V)} IC(x)}{\sum_{x \in A(U) \cup A(V)} IC(x)}$.

644 – Here, $A(X)$ is the set of terms within X and all their ancestors in the GO hierarchy.

645 We aggregated the Resnik Similarity and simIC across U and V by using the average, maxi-
646 mum, and best match average approaches.

647 • Average: $\frac{1}{|U||V|} \sum_{u \in U} \sum_{v \in V} \text{sim}(u, v)$

648 • Max: $\max_{u, v \in U \times V} \text{sim}(u, v)$

649 • Best Match Average [102]: $\frac{1}{|U|+|V|} \left[\sum_{u \in U} \max_{v \in V} \text{sim}(u, v) + \sum_{v \in V} \max_{u \in U} \text{sim}(u, v) \right]$

650 **Evaluating predictions of what drugs will treat a disease.** We evaluate how effectively a model
651 ranks the drugs that will treat a disease by using AUROC, Average Precision, and Recall@50.
652 For each disease, a model produces a ranked list of drugs. We identify the drugs approved to
653 treat the disease and, consistent with prior literature, assume that other drugs cannot treat the
654 disease [11–14]. For each disease, we then compute the model’s AUROC, Average Precision, and
655 Recall@50 values based on the ranked list of drugs. We report the model’s performance across
656 diseases by reporting the median of the AUROC, the mean of the Average Precision, and the mean
657 of the Recall@50 values across diseases.

658 To ensure robust results, we perform five-fold cross validation. We split the drugs into five
659 folds and create training and held-out sets of the drugs and their corresponding indications. We
660 compute the above evaluation metrics separately on the training and held-out sets. Ultimately, we
661 report all performance metrics on the held-out set, averaged across folds (Figure 2b).

662 **Model selection and optimization of scalar weights.** The diffusion profiles of
663 each drug and disease depend on the scalar weights used to compute them $W =$
664 $\{w_{\text{drug}}, w_{\text{disease}}, w_{\text{protein}}, w_{\text{biological function}}, w_{\text{higher-level biological function}}, w_{\text{lower-level biological function}}\}$ and the
665 probability α of continuing a walk. Similarly, how effectively diffusion profiles predict what
666 drugs treat a given disease depends on the similarity metric used to compare drug and disease
667 diffusion profiles. We optimize the prediction model across the scalar weights W , the probability
668 of continuing a walk α , and the comparison metrics by performing a sweep and selecting the
669 model with the highest median AUROC on the training set, averaged across folds.

670 After initial coarse explorations for each hyperparameter, we sweep across 486 combina-
671 tions of hyperparameters sampled linearly within the following ranges $w_{\text{drug}} \in [3, 9]$, $w_{\text{disease}} \in$
672 $[3, 9]$, $w_{\text{protein}} \in [3, 9]$, $w_{\text{higher-level biological function}} \in [1.5, 4.5]$, $w_{\text{lower-level biological function}} \in [1.5, 4.5]$, $\alpha \in$
673 $[0.85, 0.9]$ and set $w_{\text{biological function}} = w_{\text{higher-level biological function}} + w_{\text{lower-level biological function}}$. We also
674 sweep across the seven comparison metrics described above. We repeat this procedure for both
675 the multiscale interactome and the molecular-scale interactome to identify the best diffusion-
676 based model for both. The optimal weights for the molecular-scale interactome are $w_{\text{drug}} =$
677 4.88 , $w_{\text{disease}} = 6.83$, $w_{\text{protein}} = 3.21$ with $\alpha = 0.854$ and use the L1 norm to compare $\mathbf{r}^{(c)}$ and $\mathbf{r}^{(d)}$
678 (Figure 2c, Supplementary Note 1). The optimal weights for the multiscale interactome are $w_{\text{drug}} =$
679 3.21 , $w_{\text{disease}} = 3.54$, $w_{\text{protein}} = 4.40$, $w_{\text{higher-level biological function}} = 2.10$, $w_{\text{lower-level biological function}} =$
680 4.49 , $w_{\text{biological function}} = 6.58$ with $\alpha = 0.860$ and use the correlation distance to compare $\mathbf{r}^{(c)}$
681 and $\mathbf{r}^{(d)}$ (Figure 2b, c). We utilize these optimal weights for the multiscale interactome for all
682 subsequent sections. Optimized diffusion profiles are provided in Supplementary Data 10.

683 Additional information on selecting the edge weight ranges is provided as Supplementary
684 Note 2.

685 **Evaluating predictions of what drugs will treat a disease by drug category.** We analyze the
686 multiscale interactome’s predictive performance across drug categories by using the Anatomical
687 Therapeutic Chemical Classification (ATC) [103]. We map all drugs to their ATC class by using
688 DrugBank’s XML database “full_database.xml” [30]. We use the second level of the ATC classi-
689 fication and only consider categories with at least 20 drugs. For the drugs in each ATC Level II
690 category, we compute the rank of the drugs for the diseases they are approved to treat. We conduct
691 this analysis twice, first to understand the overall performance of the best multiscale interactome
692 model (Supplementary Fig. 6) and second to understand the differential performance of the best

693 multiscale interactome model compared to the best molecular-scale interactome model using dif-
694 fusion profiles (Figure 2c; Supplementary Fig. 7). The ATC classification for the drugs in our
695 study is provided in Supplementary Data 7.

696 **Diffusion profiles identify proteins and biological functions related to treatment.** For a given
697 drug-disease pair, diffusion profiles identify the proteins and biological functions related to treat-
698 ment. For each drug-disease pair, we select the top k proteins and biological functions in the drug
699 diffusion profile and in the disease diffusion profile. To explain the relevance of these proteins and
700 biological functions to treatment, we induce a subgraph on these nodes and remove any isolated
701 components. We set $k = 10$ for the case studies in Figures 2g, 2h, and 3f. We focus on these
702 nodes since the nodes ranked most highly in the diffusion profiles have the highest propagated
703 effect and are thus considered the most relevant to treatment. Additionally, these top nodes also
704 capture a substantial fraction of the overall visitation frequency in the diffusion profile (i.e. about
705 50% for Figures 2g, 2h). We additionally include the rankings of the top 20 proteins and biological
706 functions for each case study as Supplementary Fig. 16-18.

707 **Validation of diffusion profiles through gene expression signatures.** To validate drug diffusion
708 profiles, we compare drug diffusion profiles to the drug gene expression signatures present in the
709 Broad Connectivity Map [47, 48] (Figure 2f).

710 We map drugs in the Broad Connectivity Map to DrugBank IDs using PubChem IDs, drug
711 names, and the DrugBank “approved_drug_links.csv” and “drugbank_vocabulary.csv” files [30].

712 Drugs in the Broad Connectivity Map have multiple gene expression signatures based on the
713 cell line, the drug dose, and the time of exposure. However, drugs only have a single diffusion
714 profile. We thus only consider drugs where activity is consistent across cell lines and select a
715 single representative gene expression signature for each drug. To accomplish this, we follow Broad
716 Connectivity Map guidelines [47, 48] as described next. For drugs:

- 717 1. We only consider drugs with similar signatures across cell lines (an inter-cell connectivity
718 score ≥ 0.4) and with activity across many cell lines (an aggregated transcriptional activity
719 score ≥ 0.3).
- 720 2. We only consider drugs that are members of the “touchstone” dataset: the drugs that are
721 the most well-annotated and systematically profiled across the Broad’s core cell lines at
722 standardized conditions. The Broad Connectivity Map specifically recommends the “touch-
723 stone” dataset as a reference.

724 For gene expression signatures, we utilize the Level 5 Replicate Consensus Sig-
725 natures provided by the Broad Connectivity Map. Each gene expression signature
726 captures the z-scored change in expression of each gene across replicate experiments
727 (“GSE92742_Broad_LINCS_Level5_COMPZ.MODZ_n473647x12328.gctx”). For these gene ex-
728 pression signatures:

- 729 1. We only consider genes whose expression is measured directly rather than inferred (i.e.
730 “landmark” genes).
- 731 2. We only consider signatures that are highly reproducible and distinct ($\text{distil_cc_q75} \geq 0.2$
732 and $\text{pct_self_rank_q25} \leq 0.1$).
- 733 3. We require that each signature be an “exemplar” signature for the drug as indicated by the
734 Broad Connectivity Map (i.e. a highly reproducible, representative signature).
- 735 4. We require that each signature be sufficiently active (i.e. have a transcriptional activity score
736 ≥ 0.35) and result from at least 3 replicates ($\text{distil_n_sample_thresh} \geq 3$).
- 737 5. In cases where multiple signatures meet these criteria for a given drug, we select the signature
738 with the highest transcriptional activity score.

739 The gene expression signatures we ultimately use for each drug are provided in Supplemen-
740 tary Data 8.

Finally, we compare the similarity of drugs based on their diffusion profiles and their gene expression signatures. We compare the similarity of drug diffusion profiles by the Canberra distance, multiplied by -1 so higher values indicate higher similarity. We compare the similarity of drug gene expression signatures based on the overlap in the 25 most upregulated genes U and 25 most downregulated genes D :

$$\frac{1}{2} \left[\frac{|U_{\text{drug1}} \cap U_{\text{drug2}}|}{|U_{\text{drug1}} \cup U_{\text{drug2}}|} + \frac{|D_{\text{drug1}} \cap D_{\text{drug2}}|}{|D_{\text{drug1}} \cup D_{\text{drug2}}|} \right].$$

741 We use rank transformed gene expression signatures and diffusion profiles. We only allow the
742 comparison of gene expression signatures that are in the same cell, with the same dose, and at the
743 same exposure time. Ultimately, we measure the Spearman Correlation between the similarity of
744 the drugs as described by the drug diffusion profiles and the similarity of the drugs as described
745 the gene expression signatures.

746 **Compiling genetic variants that alter treatment.** We compile genetic variants that alter treat-
747 ment by using the Pharmacogenomics Knowledgebase (PharmGKB) [64]. PharmGKB is a gold-
748 standard database mapping the effect of genetic variants on treatments. PharmGKB is manually
749 curated from a range of sources, including the published literature, the Allele Frequency Database,
750 the Anatomical Therapeutic Chemical Classification, ChEBI, ClinicalTrials.gov, dbSNP, Drug-
751 Bank, the European Medicines Agency, Ensembl, FDA Drug Labels at DailyMed, GeneCard,
752 HC-SC, HGNC, HMDB, HumanCyc Gene, LS-SNP, MedDRA, MeSH, NCBI Gene, NDF-RT,
753 PMDA, PubChem Compound, RxNorm, SnoMed Clinical Terminology, and UniProt KB.

754 We use PharmGKB’s “Clinical Annotations” which detail how variants at the gene level al-
755 ter treatments. PharmGKB’s “clinical_ann_metadata.tsv” file provides triplets of drugs, diseases,
756 and genetic variants known to alter treatment. Treatment alteration occurs when a genetic vari-
757 ant alters the efficacy, dosage, metabolism, or pharmacokinetics of treatment or otherwise causes
758 toxicity or an adverse drug reaction. We map genes to their Entrez ID using HUGO, drugs to
759 their DrugBank ID using PharmGKB’s “drugs.tsv” and “chemicals.tsv” files, and diseases to their
760 UMLS CUIDs by using PharmGKB’s “phenotypes.tsv” file. To ensure consistency with the ap-
761 proved drug-disease pairs we previously compiled, we only consider (drug, disease, gene) triplets
762 in which the drug and disease are part of an FDA-approved treatment. Ultimately, we obtain 1,223
763 drug-disease-gene triplets with 201 drugs, 94 diseases, and 455 genes. All drug-disease-gene
764 triplets are provided in Supplementary Data 9.

Computing treatment importance of a gene based on diffusion profiles. We define the treat-
ment importance (TI) of gene i as the product of the visitation frequency of the corresponding
protein in the drug and disease diffusion profiles. For a treatment composed of drug compound c
and disease d , the treatment importance of gene i is:

$$\text{TI}(i|c, d) = \mathbf{r}_i^{(c)} * \mathbf{r}_i^{(d)}.$$

765 We define the treatment importance percentile as the percentile rank of $\text{TI}(i|c, d)$ compared
766 to all other genes for the same drug and disease. Intuitively, gene i is important to a treatment if
767 the corresponding protein is frequently visited in both the drug and disease diffusion profiles.

768 **Comparing treatment importance of treatment altering genetic mutations vs other genetic**
769 **mutations.** We compare the treatment importance of genes known to alter a treatment with the

770 treatment importance of other genes (Figure 3b). In particular, we compare the set of (drug, disease,
771 gene) triplets where the gene is known to alter the drug-disease treatment with an equivalently sized
772 set of (drug, disease, gene) triplets where the gene is not known to alter treatment. We construct
773 the latter set by sampling drugs, diseases, and genes uniformly at random that are not known to
774 alter treatment from PharmGKB [64]. The drugs and diseases in all triplets correspond to approved
775 drug-disease pairs. Thereby, we construct a distribution of the treatment importance for “treatment
776 altering genes” and a distribution of the treatment importance for “other genes” (Figure 3b).

777 **Predicting genes that alter a treatment based on treatment importance.** We evaluate the abil-
778 ity of treatment importance to predict the genes that will alter a given treatment (Figure 3c). For
779 each (drug, disease, gene) triplet, we use the treatment importance of the gene $\text{TI}(i|c, d)$ to predict
780 whether the gene alters treatment or not for that drug-disease pair (i.e. binary classification). We
781 use the set of positive and negative (drug, disease, gene) triplets constructed previously (see Meth-
782 ods: Comparing treatment importance of treatment altering genetic mutations vs other genetic
783 mutations). We assess performance using AUROC and Average Precision (Figure 3c).

784 **Comparing treatment importance of genes that alter one drug indicated to treat a disease but**
785 **not another.** We analyze how often a gene has a higher treatment importance in the treatments it
786 alters than in those it does not alter (Figure 3e).

Formally, let i be a gene. Consider a triplet $(d, c_{\text{altered}}, c_{\text{unaltered}})$ of a disease d , a drug c_{altered}
approved to treat the disease whose treatment is altered due to a mutation in i , and a drug $c_{\text{unaltered}}$
approved to treat the disease whose treatment is not altered due to a mutation in i . Let n_{triplets} be
the total number of such triplets for gene i . For each gene i , we measure the fraction f of triplets
 $(d, c_{\text{altered}}, c_{\text{unaltered}})$ for which the treatment importance of i is higher in the (c_{altered}, d) treatment than
in the $(c_{\text{unaltered}}, d)$ treatment, as shown below. We only consider genes for which $n_{\text{triplets}} \geq 100$.

$$f[\text{TI}(i|c_{\text{altered}}, d) > \text{TI}(i|c_{\text{unaltered}}, d)] = \frac{\sum_{\forall(d, c_{\text{altered}}, c_{\text{unaltered}})} \mathbb{1}\{\text{TI}(i|c_{\text{altered}}, d) > \text{TI}(i|c_{\text{unaltered}}, d)\}}{n_{\text{triplets}}}.$$

787 **Analyzing whether distant proteins can have common biological functions.** We analyzed
788 whether two proteins can be more distant than expected by random chance in a physical protein-
789 protein interaction (PPI) network yet affect the same function (Supplementary Fig. 2). To run this
790 analysis, we first compute the set of all protein pairs that are both present in the protein-protein
791 interaction network described previously (Methods: Protein-protein interactions) and are also as-

792 sociated with a common biological function. We only consider direct associations of proteins to
793 biological functions (i.e. we do not propagate associations up the GO hierarchy) in order to ensure
794 that shared biological functions are specific and not generic (i.e. shared associations with the GO
795 term 'Biological Process').

796 For each protein pair with a common biological function, we then:

- 797 1. Compute the shortest path distance in the PPI network between these two proteins.
- 798 2. Construct a reference distribution of shortest paths for these two protein pairs by following
799 the approach in [10, 13]. Specifically, we randomly sample other proteins in the network
800 with similar degree to the original proteins and measure the shortest path distance. These
801 randomly sampled proteins do *not* necessarily share a common biological function.
- 802 3. Using the true shortest path distance between the proteins and the random reference distribu-
803 tion, we compute a z-score. The z-score captures whether the proteins with a shared function
804 are closer or further away than expected by random chance in the PPI network.

805 **Construction of alternative multiscale interactomes that explicitly represent cells, tissues,**
806 **and organs.** We constructed three alternative multiscale interactomes which explicitly represent
807 cells, tissues, and organs. In these alternative multiscale interactomes, the nodes and edges in the
808 original multiscale interactome are all present. Additionally, (1) human cells, tissues, and organs
809 are added as additional nodes; (2) edges between these cell, tissue, and organ nodes are added
810 according to relationships defined in established anatomical ontologies; and (3) edges between GO
811 biological function nodes and cell, tissue, and organ nodes are added according to relationships
812 provided in Gene Ontology Plus (GO Plus) [104]. GO Plus maintains a curated set of relationships
813 between the biological functions in GO and the cell, tissue, and organ nodes present in two key
814 anatomical ontologies: Uberon and the Cell Ontology. We thus constructed three alternative mul-
815 tiscale interactomes incorporating human subsets of Uberon, the Cell Ontology, and both Uberon
816 and the Cell Ontology.

- 817 1. *Multiscale Interactome + Uberon*: Uberon is an ontology covering anatomical struc-
818 tures in animals [105, 106]. Uberon nodes include tissues (i.e. cardiac muscle tis-
819 sue UBERON:0001133), organs (i.e. heart UBERON:0000948), and organ systems (i.e.
820 cardiovascular system UBERON:0004535). We utilized GO Plus (i.e. “go-plus.owl”)

821 to link GO biological function nodes present in our original network to Uberon nodes
822 present in a human-specific subset of Uberon (i.e. “subsets/human-view.obo”). Edges be-
823 tween Uberon nodes, which encode anatomical relationships, were also added according to
824 “subsets/human-view.obo”.

825 2. *Multiscale Interactome + Cell Ontology*: The Cell Ontology is an ontology for the represen-
826 tation of in vivo cell types [107, 108]. Nodes consist primarily of cell types and their hierar-
827 chical relationships (i.e. epithelial cell CL:0000066, epithelial cell of pancreas CL:0000083,
828 pancreatic A cell CL:0000171). We utilized a human-specific subset of the Cell Ontology
829 previously prepared by the Human Cell Atlas Ontology [109]. We utilized GO Plus to link
830 GO biological function nodes in our original network to Cell Ontology terms and the Cell
831 Ontology (i.e. “cl-basic.obo”) to link Cell Ontology terms with one another.

832 3. *Multiscale Interactome + Uberon + Cell Ontology*: The “Multiscale Interactome + Uberon
833 + Cell Ontology” network contains all nodes and edges present in our original network as
834 well as nodes and edges added via GO Plus, Uberon, and Cell Ontology as described above.

835 **Prediction of what drugs treat a given disease in alternative multiscale interactomes.** We
836 evaluate the ability of diffusion profiles to predict what drugs treat a given disease in the alternative
837 multiscale interactomes (see Methods: Construction of alternative multiscale interactomes that
838 explicitly represent cells, tissues, and organs). Given the presence of new node types, we modify
839 the edge weight hyperparameters used in the calculation of diffusion profiles. We then sweep
840 over the full set of edge weight hyperparameters according to the broad hyperparameter sweep
841 described in Supplementary Note 2, in which we sample 586 combinations of hyperparameters
842 sampled linearly in the range [1, 100]. The new sets of edge weight hyperparameters and their
843 optimal values are present below:

844 1. *Multiscale Interactome + Uberon*: The optimal weights for “Multiscale Interactome +
845 Uberon” are $w_{\text{drug}} = 55.2$, $w_{\text{disease}} = 27.3$, $w_{\text{protein}} = 76.8$, $w_{\text{biological function}} = 66.1$, $w_{\text{uberon}} =$
846 82.2 , $w_{\text{higher-level biological function or uberon}} = 67.1$, $w_{\text{lower-level biological function or uberon}} = 45.7$ with $\alpha =$
847 0.76 and use the correlation distance to compare $\mathbf{r}^{(c)}$ and $\mathbf{r}^{(d)}$.

848 2. *Multiscale Interactome + Cell Ontology*: The optimal weights for “Multiscale In-
849 teractome + Cell Ontology” are $w_{\text{drug}} = 39.0$, $w_{\text{disease}} = 17.1$, $w_{\text{protein}} =$

850 72.4, $w_{\text{biological function}} = 60.0$, $w_{\text{cell ontology}} = 23.1$, $w_{\text{higher-level biological function or cell ontology}} =$
851 25.7, $w_{\text{lower-level biological function or cell ontology}} = 22.8$ with $\alpha = 0.83$ and use the correlation dis-
852 tance to compare $\mathbf{r}^{(c)}$ and $\mathbf{r}^{(d)}$.

853 3. *Multiscale Interactome + Uberon + Cell Ontology*: The optimal weights
854 for “Multiscale Interactome + Uberon + Cell Ontology” are $w_{\text{drug}} =$
855 60.2, $w_{\text{disease}} = 12.8$, $w_{\text{protein}} = 42.3$, $w_{\text{biological function}} = 78.4$, $w_{\text{uberon}} =$
856 70.0, $w_{\text{cell ontology}} = 91.7$, $w_{\text{higher-level biological function or uberon or cell ontology}} =$
857 26.7, $w_{\text{lower-level biological function or uberon or cell ontology}} = 76.1$ with $\alpha = 0.82$ and use the cor-
858 relation distance to compare $\mathbf{r}^{(c)}$ and $\mathbf{r}^{(d)}$.

References

- 859
860
861 1. Huttlin, E. L. *et al.* Architecture of the human interactome defines protein communities and
862 disease networks. *Nature* **545**, 505–509 (2017).
- 863 2. Creixell, P. *et al.* Pathway and network analysis of cancer genomes. *Nature Methods* **12**, 615
864 (2015).
- 865 3. Parikshak, N. N., Gandal, M. J. & Geschwind, D. H. Systems biology and gene networks in
866 neurodevelopmental and neurodegenerative disorders. *Nature Reviews Genetics* **16**, 441–458
867 (2015).
- 868 4. Leiserson, M. D. *et al.* Pan-cancer network analysis identifies combinations of rare somatic
869 mutations across pathways and protein complexes. *Nature Genetics* **47**, 106–114 (2015).
- 870 5. Nikolsky, Y., Nikolskaya, T. & Bugrim, A. Biological networks and analysis of experimental
871 data in drug discovery. *Drug Discovery Today* **10**, 653–662 (2005).
- 872 6. Hu, J. X., Thomas, C. E. & Brunak, S. Network biology concepts in complex disease comor-
873 bidities. *Nature Reviews Genetics* **17**, 615–629 (2016).
- 874 7. Hormozdiari, F., Penn, O., Borenstein, E. & Eichler, E. E. The discovery of integrated gene
875 networks for autism and related disorders. *Genome Research* **25**, 142–154 (2015).
- 876 8. Barabási, A.-L., Gulbahce, N. & Loscalzo, J. Network medicine: a network-based approach
877 to human disease. *Nature Reviews Genetics* **12**, 56–68 (2011).
- 878 9. Cowen, L., Ideker, T., Raphael, B. J. & Sharan, R. Network propagation: a universal amplifier
879 of genetic associations. *Nature Reviews Genetics* **18**, 551–562 (2017).
- 880 10. Cheng, F. *et al.* Network-based approach to prediction and population-based validation of in
881 silico drug repurposing. *Nature Communications* **9**, 2691 (2018).
- 882 11. Pushpakom, S. *et al.* Drug repurposing: progress, challenges and recommendations. *Nature*
883 *Reviews Drug Discovery* **18**, 41–58 (2019).
- 884 12. Lotfi Shahreza, M., Ghadiri, N., Mousavi, S. R., Varshosaz, J. & Green, J. R. A review of
885 network-based approaches to drug repositioning. *Briefings in Bioinformatics* **19**, 878–892
886 (2018).
- 887 13. Guney, E., Menche, J., Vidal, M. & Barabasi, A.-L. Network-based in silico drug efficacy
888 screening. *Nature Communications* **7**, 10331 (2016).
- 889 14. Wang, W., Yang, S., Zhang, X. & Li, J. Drug repositioning by integrating target information
890 through a heterogeneous network model. *Bioinformatics* **30**, 2923–2930 (2014).
- 891 15. Luo, Y. *et al.* A network integration approach for drug-target interaction prediction and
892 computational drug repositioning from heterogeneous information. *Nature Communications*
893 **8**, 573 (2017).
- 894 16. Zitnik, M., Agrawal, M. & Leskovec, J. Modeling polypharmacy side effects with graph
895 convolutional networks. *Bioinformatics* **34**, i457–i466 (2018).
- 896 17. Cheng, F., Kovacs, I. A. & Barabasi, A.-L. Network-based prediction of drug combinations.
897 *Nature Communications* **10**, 1197 (2019).

- 898 18. Hu, Y. *et al.* Optimal control nodes in disease-perturbed networks as targets for combination
899 therapy. *Nature Communications* **10**, 2180 (2019).
- 900 19. Firestone, A. J. & Settleman, J. A three-drug combination to treat BRAF-mutant cancers.
901 *Nature Medicine* **23**, 913–914 (2017).
- 902 20. Zhao, S. & Iyengar, R. Systems pharmacology: network analysis to identify multiscale
903 mechanisms of drug action. *Annual Review of Pharmacology and Toxicology* **52**, 505–521
904 (2012).
- 905 21. Walpole, J., Papin, J. A. & Peirce, S. M. Multiscale computational models of complex
906 biological systems. *Annual Review of Biomedical Engineering* **15**, 137–154 (2013).
- 907 22. van Hasselt, J. C. & Iyengar, R. Systems pharmacology: defining the interactions of drug
908 combinations. *Annual Review of Pharmacology and Toxicology* **59**, 21–40 (2019).
- 909 23. Han, K. *et al.* Synergistic drug combinations for cancer identified in a CRISPR screen for
910 pairwise genetic interactions. *Nature Biotechnology* **35**, 463–474 (2017).
- 911 24. Jia, J. *et al.* Mechanisms of drug combinations: interaction and network perspectives. *Nature*
912 *Reviews Drug Discovery* **8**, 111–128 (2009).
- 913 25. Yu, M. K. *et al.* Translation of genotype to phenotype by a hierarchy of cell subsystems. *Cell*
914 *Systems* **2**, 77–88 (2016).
- 915 26. Zañudo, J. G. T., Scaltriti, M. & Albert, R. A network modeling approach to elucidate drug
916 resistance mechanisms and predict combinatorial drug treatments in breast cancer. *Cancer*
917 *Convergence* **1**, 5 (2017).
- 918 27. Zañudo, J. G., Steinway, S. N. & Albert, R. Discrete dynamic network modeling of oncogenic
919 signaling: Mechanistic insights for personalized treatment of cancer. *Current Opinion in*
920 *Systems Biology* **9**, 1–10 (2018).
- 921 28. Trachana, K. *et al.* Taking systems medicine to heart. *Circulation Research* **122**, 1276–1289
922 (2018).
- 923 29. Montagud, A. *et al.* Conceptual and computational framework for logical modelling of bio-
924 logical networks deregulated in diseases. *Briefings in Bioinformatics* **20**, 1238–1249 (2019).
- 925 30. Wishart, D. S. *et al.* DrugBank 5.0: a major update to the DrugBank database for 2018.
926 *Nucleic Acids Research* **46**, D1074–D1082 (2017).
- 927 31. Corsello, S. M. *et al.* The Drug Repurposing Hub: a next-generation drug library and infor-
928 mation resource. *Nature Medicine* **23**, 405–408 (2017).
- 929 32. Piñero, J. *et al.* DisGeNET: a comprehensive platform integrating information on human
930 disease-associated genes and variants. *Nucleic Acids Research* **45**, D833–D839 (2016).
- 931 33. Menche, J. *et al.* Uncovering disease-disease relationships through the incomplete interac-
932 tome. *Science* **347**, 1257601 (2015).
- 933 34. Oughtred, R. *et al.* The BioGRID interaction database: 2019 update. *Nucleic Acids Research*
934 **47**, D529–D541 (2019).
- 935 35. Rolland, T. *et al.* A proteome-scale map of the human interactome network. *Cell* **159**, 1212–
936 1226 (2014).

- 937 36. Salwinski, L. *et al.* The Database of Interacting Proteins: 2004 update. *Nucleic Acids Re-*
938 *search* **32**, D449–D451 (2004).
- 939 37. Venkatesan, K. *et al.* An empirical framework for binary interactome mapping. *Nature*
940 *Methods* **6**, 83 (2009).
- 941 38. Yu, H. *et al.* Next-generation sequencing to generate interactome datasets. *Nature Methods*
942 **8**, 478 (2011).
- 943 39. Rual, J.-F. *et al.* Towards a proteome-scale map of the human protein–protein interaction
944 network. *Nature* **437**, 1173–1178 (2005).
- 945 40. Consortium, G. O. The Gene Ontology resource: 20 years and still GOing strong. *Nucleic*
946 *Acids Research* **47**, D330–D338 (2018).
- 947 41. Ashburner, M. *et al.* Gene Ontology: tool for the unification of biology. *Nature Genetics* **25**,
948 25–29 (2000).
- 949 42. Brown, A. S. & Patel, C. J. A standard database for drug repositioning. *Scientific Data* **4**,
950 170029 (2017).
- 951 43. Sharp, M. E. Toward a comprehensive drug ontology: extraction of drug-indication relations
952 from diverse information sources. *Journal of Biomedical Semantics* **8**, 2 (2017).
- 953 44. Donnat, C., Zitnik, M., Hallac, D. & Leskovec, J. Learning structural node embeddings via
954 diffusion wavelets. In *Proceedings of the 24th ACM SIGKDD International Conference on*
955 *Knowledge Discovery & Data Mining*, 1320–1329 (2018).
- 956 45. Nielsen, S. *et al.* Vasopressin increases water permeability of kidney collecting duct by
957 inducing translocation of aquaporin-CD water channels to plasma membrane. *Proceedings*
958 *of the National Academy of Sciences* **92**, 1013–1017 (1995).
- 959 46. Holmes, C. L., Landry, D. W. & Granton, J. T. Science review: vasopressin and the cardio-
960 vascular system part 1–receptor physiology. *Critical Care* **7**, 427–434 (2003).
- 961 47. Subramanian, A. *et al.* A next generation connectivity map: L1000 platform and the first
962 1,000,000 profiles. *Cell* **171**, 1437–1452 (2017).
- 963 48. Lamb, J. *et al.* The Connectivity Map: using gene-expression signatures to connect small
964 molecules, genes, and disease. *Science* **313**, 1929–1935 (2006).
- 965 49. Utermann, G., Jaeschke, M. & Menzel, J. Familial hyperlipoproteinemia type III: Deficiency
966 of a specific apolipoprotein (APO E-III) in the very-low-density lipoproteins. *FEBS Letters*
967 **56**, 352–355 (1975).
- 968 50. Utermann, G. *et al.* Polymorphism of apolipoprotein E: Genetics of hyperlipoproteinemia
969 type III. *Clinical Genetics* **15**, 37–62 (1979).
- 970 51. Ghiselli, G., Schaefer, E. J., Gascon, P. & Breser, H. Type III hyperlipoproteinemia associated
971 with apolipoprotein E deficiency. *Science* **214**, 1239–1241 (1981).
- 972 52. Wang, J. *et al.* APOA5 genetic variants are markers for classic hyperlipoproteinemia pheno-
973 types and hypertriglyceridemia. *Nature Clinical Practice Cardiovascular Medicine* **5**, 730–
974 737 (2008).

- 975 53. Evans, D., Seedorf, U. & Beil, F. Polymorphisms in the apolipoprotein a5 (APOA5) gene
976 and type III hyperlipidemia. *Clinical Genetics* **68**, 369–372 (2005).
- 977 54. Moghadasian, M. H. Clinical pharmacology of 3-hydroxy-3-methylglutaryl coenzyme a-
978 reductase inhibitors. *Life Sciences* **65**, 1329–1337 (1999).
- 979 55. Holdgate, G., Ward, W. & McTaggart, F. Molecular mechanism for inhibition of 3-hydroxy-
980 3-methylglutaryl CoA (HMG-CoA) reductase by rosuvastatin. *Biochemical Society Trans-
981 actions* **31**, 528–531 (2003).
- 982 56. Shinkai, K., McCalmont, T. & Leslie, K. Cryopyrin-associated periodic syndromes and
983 autoinflammation. *Clinical and Experimental Dermatology: Clinical Dermatology* **33**, 1–9
984 (2008).
- 985 57. Kone-Paut, I. & Galeotti, C. Anakinra for cryopyrin-associated periodic syndrome. *Expert
986 Review of Clinical Immunology* **10**, 7–18 (2014).
- 987 58. Ashley, E. A. Towards precision medicine. *Nature Reviews Genetics* **17**, 507–522 (2016).
- 988 59. Goldstein, D. B., Tate, S. K. & Sisodiya, S. M. Pharmacogenetics goes genomic. *Nature
989 Reviews Genetics* **4**, 937–947 (2003).
- 990 60. Hansen, N. T., Brunak, S. & Altman, R. Generating genome-scale candidate gene lists for
991 pharmacogenomics. *Clinical Pharmacology & Therapeutics* **86**, 183–189 (2009).
- 992 61. Karczewski, K. J., Daneshjou, R. & Altman, R. B. Chapter 7: Pharmacogenomics. *PLOS
993 Computational Biology* **8**, e1002817 (2012).
- 994 62. Su, X. *et al.* Association between angiotensinogen, angiotensin II receptor genes, and blood
995 pressure response to an angiotensin-converting enzyme inhibitor. *Circulation* **115**, 725–732
996 (2007).
- 997 63. Yu, H. *et al.* A core promoter variant of angiotensinogen gene and interindividual variation
998 in response to angiotensin-converting enzyme inhibitors. *Journal of the Renin-Angiotensin-
999 Aldosterone System* **15**, 540–546 (2014).
- 1000 64. Barbarino, J. M., Whirl-Carrillo, M., Altman, R. B. & Klein, T. E. PharmGKB: a world-
1001 wide resource for pharmacogenomic information. *Wiley Interdisciplinary Reviews: Systems
1002 Biology and Medicine* **10**, e1417 (2018).
- 1003 65. Nayler, W. G. & Dillon, J. Calcium antagonists and their mode of action: an historical
1004 overview. *British Journal of Clinical Pharmacology* **21**, 97S–107S (1986).
- 1005 66. Sutton, M. S. J. & Morad, M. Mechanisms of action of diltiazem in isolated human atrial and
1006 ventricular myocardium. *Journal of Molecular and Cellular Cardiology* **19**, 497–508 (1987).
- 1007 67. O'Connor, S. E., Grosset, A. & Janiak, P. The pharmacological basis and pathophysiological
1008 significance of the heart rate-lowering property of diltiazem. *Fundamental & Clinical
1009 Pharmacology* **13**, 145–153 (1999).
- 1010 68. Balfour, J. A. & Goa, K. L. Benazepril. *Drugs* **42**, 511–539 (1991).
- 1011 69. Lavoie, J. L. & Sigmund, C. D. Minireview: overview of the renin-angiotensin system—an
1012 endocrine and paracrine system. *Endocrinology* **144**, 2179–2183 (2003).

- 1013 70. Caulfield, M. *et al.* Linkage of the angiotensinogen gene to essential hypertension. *New*
1014 *England Journal of Medicine* **330**, 1629–1633 (1994).
- 1015 71. Jeunemaitre, X. *et al.* Molecular basis of human hypertension: role of angiotensinogen. *Cell*
1016 **71**, 169–180 (1992).
- 1017 72. Sanchez-Vega, F. *et al.* Oncogenic signaling pathways in The Cancer Genome Atlas. *Cell*
1018 **173**, 321–337 (2018).
- 1019 73. Jones, D. Pathways to cancer therapy. *Nature Reviews Drug Discovery* **7**, 875–876 (2008).
- 1020 74. Jones, S. *et al.* Core signaling pathways in human pancreatic cancers revealed by global
1021 genomic analyses. *Science* **321**, 1801–1806 (2008).
- 1022 75. Parsons, D. W. *et al.* An integrated genomic analysis of human glioblastoma multiforme.
1023 *Science* **321**, 1807–1812 (2008).
- 1024 76. Di Leva, G., Garofalo, M. & Croce, C. M. MicroRNAs in cancer. *Annual Review of Pathol-*
1025 *ogy: Mechanisms of Disease* **9**, 287–314 (2014).
- 1026 77. Ma, J. *et al.* Using deep learning to model the hierarchical structure and function of a cell.
1027 *Nature Methods* **15**, 290 (2018).
- 1028 78. Cho, H., Berger, B. & Peng, J. Compact integration of multi-network topology for functional
1029 analysis of genes. *Cell Systems* **3**, 540–548 (2016).
- 1030 79. Wang, S., Cho, H., Zhai, C., Berger, B. & Peng, J. Exploiting ontology graph for predicting
1031 sparsely annotated gene function. *Bioinformatics* **31**, i357–i364 (2015).
- 1032 80. Mi, H., Muruganujan, A., Casagrande, J. T. & Thomas, P. D. Large-scale gene function
1033 analysis with the PANTHER classification system. *Nature Protocols* **8**, 1551 (2013).
- 1034 81. Yamanishi, Y., Kotera, M., Kanehisa, M. & Goto, S. Drug-target interaction prediction from
1035 chemical, genomic and pharmacological data in an integrated framework. *Bioinformatics* **26**,
1036 i246–i254 (2010).
- 1037 82. Balaji, S., McClendon, C., Chowdhary, R., Liu, J. S. & Zhang, J. IMID: integrated molecular
1038 interaction database. *Bioinformatics* **28**, 747–749 (2012).
- 1039 83. Bell, L., Chowdhary, R., Liu, J. S., Niu, X. & Zhang, J. Integrated bio-entity network: a
1040 system for biological knowledge discovery. *PLOS One* **6**, e21474 (2011).
- 1041 84. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An expanded view of complex traits: from polygenic
1042 to omnigenic. *Cell* **169**, 1177–1186 (2017).
- 1043 85. Liu, X., Li, Y. I. & Pritchard, J. K. Trans effects on gene expression can drive omnigenic
1044 inheritance. *Cell* **177**, 1022–1034 (2019).
- 1045 86. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular
1046 interaction networks. *Genome Research* **13**, 2498–2504 (2003).
- 1047 87. Yates, B. *et al.* Genenames.org: the HGNC and VGNC resources in 2017. *Nucleic Acids*
1048 *Research* **45**, D619–D625 (2016).
- 1049 88. Vinayagam, A. *et al.* A directed protein interaction network for investigating intracellular
1050 signal transduction. *Science Signaling* **4**, rs8–rs8 (2011).

- 1051 89. Klopfenstein, D. V. *et al.* GOATOOLS: A python library for gene ontology analyses. *Scientific Reports* **8**, 1–17 (2018).
1052
- 1053 90. Bodenreider, O. The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic Acids Research* **32**, D267–D270 (2004).
1054
- 1055 91. Davis, A. P. *et al.* The Comparative Toxicogenomics Database: update 2019. *Nucleic Acids Research* **47**, D948–D954 (2019).
1056
- 1057 92. Schriml, L. M. *et al.* Human Disease Ontology 2018 update: classification, content and workflow expansion. *Nucleic Acids Research* **47**, D955–D962 (2019).
1058
- 1059 93. Langville, A. N. & Meyer, C. D. A survey of eigenvector methods for web information retrieval. *SIAM Review* **47**, 135–161 (2005).
1060
- 1061 94. Page, L., Brin, S., Motwani, R. & Winograd, T. The PageRank citation ranking: Bringing order to the web. Tech. Rep., Stanford InfoLab (1999).
1062
- 1063 95. Hagberg, A., Swart, P. & S Chult, D. Exploring network structure, dynamics, and function using NetworkX. Tech. Rep., Los Alamos National Lab.(LANL), Los Alamos, NM (United States) (2008).
1064
1065
- 1066 96. Li, B., Luo, F., Wang, J. Z., Feltus, F. A. & Zhou, J. Effectively integrating information content and structural relationship to improve the GO-based similarity measure between proteins. In Arabnia, H. R. *et al.* (eds.) *International Conference on Bioinformatics & Computational Biology, BIOCOMP 2010, July 12-15, 2010, Las Vegas Nevada, USA, 2 Volumes*, 166–172 (CSREA Press, 2010).
1067
1068
1069
1070
- 1071 97. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences* **102**, 15545–15550 (2005).
1072
1073
- 1074 98. Pesquita, C. Semantic similarity in the Gene Ontology. In *The Gene Ontology Handbook*, 161–173 (Humana Press, New York, NY, 2017).
1075
- 1076 99. Lord, P. W., Stevens, R. D., Brass, A. & Goble, C. A. Investigating semantic similarity measures across the Gene Ontology: the relationship between sequence and annotation. *Bioinformatics* **19**, 1275–1283 (2003).
1077
1078
- 1079 100. Resnik, P. Semantic similarity in a taxonomy: An information-based measure and its application to problems of ambiguity in natural language. *Journal of Artificial Intelligence Research* **11**, 95–130 (1999).
1080
1081
- 1082 101. Pesquita, C. *et al.* Metrics for GO based protein semantic similarity: a systematic evaluation. In *BMC Bioinformatics*, vol. 9, S4 (Springer, 2008).
1083
- 1084 102. Azuaje, F., Wang, H. & Bodenreider, O. Ontology-driven similarity approaches to supporting gene functional assessment. In *Proceedings of the ISMB'2005 SIG Meeting on Bio-ontologies*, 9–10 (2005).
1085
1086
- 1087 103. Organization, W. H. *et al.* The Anatomical Therapeutic Chemical Classification System with defined daily doses-ATC/DDD (2009).
1088
- 1089 104. Consortium, G. O. Gene Ontology Consortium: going forward. *Nucleic Acids Research* **43**, D1049–D1056 (2015).
1090

- 1091 105. Mungall, C. J., Torniai, C., Gkoutos, G. V., Lewis, S. E. & Haendel, M. A. Uberon, an
1092 integrative multi-species anatomy ontology. *Genome Biology* **13**, R5 (2012).
- 1093 106. Haendel, M. A. *et al.* Unification of multi-species vertebrate anatomy ontologies for com-
1094 parative biology in Uberon. *Journal of Biomedical Semantics* **5**, 21 (2014).
- 1095 107. Bard, J., Rhee, S. Y. & Ashburner, M. An ontology for cell types. *Genome Biology* **6**, R21
1096 (2005).
- 1097 108. Diehl, A. D. *et al.* The Cell Ontology 2016: enhanced content, modularization, and ontology
1098 interoperability. *Journal of Biomedical Semantics* **7**, 1–10 (2016).
- 1099 109. Welter, D., Jupp, S. & Osumi-Sutherland, D. Human Cell Atlas Ontology. In *Proceedings of*
1100 *the 9th International Conference on Biological Ontology (ICBO)* (2018).