

1 **TITLE PAGE**

2

3 **Full Title**

4 Systematic interrogation of diverse Omic data reveals interpretable, robust, and generalizable
5 transcriptomic features of clinically successful therapeutic targets

6 **Short Title**

7 Omic features of successful therapeutic targets

8 **Authors**

9 Andrew D. Rouillard¹, Mark R. Hurle¹, Pankaj Agarwal^{1*}

10 **Affiliations**

11 ¹Computational Biology, GSK, King of Prussia, PA, USA

12 **Corresponding Author**

13 *pankaj.agarwal@gsk.com (PA)

14

15 **Author Contributions**

16 Andrew D. Rouillard: data curation, formal analysis, methodology, software, visualization,
17 writing (original draft preparation), writing (review and editing)

18

19 Mark R. Hurle: conceptualization, data curation, methodology, writing (review and editing)

20

21 Pankaj Agarwal: conceptualization, methodology, supervision, writing (review and editing)

22

23 **ABSTRACT**

24

25 Target selection is the first and pivotal step in drug discovery. An incorrect choice may not
26 manifest itself for many years after hundreds of millions of research dollars have been spent. We
27 collected a set of 332 targets that succeeded or failed in phase III clinical trials, and explored
28 whether Omic features describing the target genes could predict clinical success. We obtained
29 features from the recently published comprehensive resource: Harmonizome. Nineteen features
30 appeared to be significantly correlated with phase III clinical trial outcomes, but only 4 passed
31 validation schemes that used bootstrapping or modified permutation tests to assess feature
32 robustness and generalizability while accounting for target class selection bias. We also used
33 classifiers to perform multivariate feature selection and found that classifiers with a single
34 feature performed as well in cross-validation as classifiers with more features (AUROC=0.57
35 and AUPR=0.81). The two predominantly selected features were mean mRNA expression across
36 tissues and standard deviation of expression across tissues, where successful targets tended to
37 have lower mean expression and higher expression variance than failed targets. This finding
38 supports the conventional wisdom that it is favorable for a target to be present in the tissue(s)
39 affected by a disease and absent from other tissues. Overall, our results suggest that it is feasible
40 to construct a model integrating interpretable target features to inform target selection. We
41 anticipate deeper insights and better models in the future, as researchers can reuse the data we
42 have provided to improve methods for handling sample biases and learn more informative
43 features. Code, documentation, and data for this study have been deposited on GitHub at
44 <https://github.com/arouillard/omic-features-successful-targets>.

45

46 **AUTHOR SUMMARY**

47

48 Drug discovery often begins with a hypothesis that changing the abundance or activity of a
49 target—a biological molecule, usually a protein—will cure a disease or ameliorate its symptoms.
50 Whether a target hypothesis translates into a successful therapy depends in part on the
51 characteristics of the target, but it is not completely understood which target characteristics are
52 important for success. We sought to answer this question with a supervised machine learning
53 approach. We obtained outcomes of target hypotheses tested in clinical trials, scoring targets as
54 successful or failed, and then obtained thousands of features (i.e. properties or characteristics) of
55 targets from dozens of biological datasets. We statistically tested which features differed
56 between successful and failed targets, and built a computational model that used these features to
57 predict success or failure of targets in clinical trials. We found that successful targets tended to
58 have more variable mRNA abundance from tissue to tissue and lower average abundance across
59 tissues than failed targets. Thus, it is probably favorable for a target to be present in the tissue(s)
60 affected by a disease and absent from other tissues. Our work demonstrates the feasibility of
61 predicting clinical trial outcomes from target features.

62

63 **INTRODUCTION**

64

65 More than half of drug candidates that advance beyond phase I clinical trials fail due to lack of
66 efficacy (1, 2). One possible explanation for these failures is sub-optimal target selection (3).

67 Many factors must be considered when selecting a target for drug discovery (4, 5). Intrinsic
68 factors include the likelihood of the target to be tractable (can the target's activity be altered by a
69 compound, antibody, or other drug modality?), safe (will altering the target's activity cause
70 serious adverse events?), and efficacious (will altering the target's activity provide significant
71 benefit to patients?). Extrinsic factors include the availability of investigational reagents and
72 disease models for preclinical target validation, whether biomarkers are known for measuring
73 target engagement or therapeutic effect, the duration and complexity of clinical trials required to
74 prove safety and efficacy, and the unmet need of patients with diseases that might be treated by
75 modulating the target.

76

77 Over the past decade, technologies have matured enabling high-throughput genome-,
78 transcriptome-, and proteome-wide profiling of cells and tissues in normal, disease, and
79 experimentally perturbed states. In parallel, researchers have made substantial progress curating
80 or text-mining biomedical literature to extract and organize information about genes and
81 proteins, such as molecular functions and signaling pathways, into structured datasets. Taken
82 together, both efforts have given rise to a vast amount of primary, curated, and text-mined data
83 about genes and proteins, which are stored in online repositories and amenable to computational
84 analysis (6, 7).

85

86 To improve the success rate of drug discovery projects, researchers have investigated whether
87 any features of genes or proteins are useful for target selection. These computational studies can
88 be categorized according to whether the researchers were trying to predict tractability (8, 9),
89 safety (10-13), efficacy (no publications to our knowledge), or overall success (alternatively
90 termed “drug target likeness”) (8, 13-26). Closely related efforts include disease gene prediction,
91 where the goal is to predict genes mechanistically involved in a given disease (27-32), and
92 disease target prediction, where the goal is to predict genes that would make successful drug
93 targets for a given disease (33-35).

94

95 To our knowledge, we report the first screen for features of genes or proteins that distinguish
96 targets of approved drugs from targets of drug candidates that failed in clinical trials. In contrast,
97 related prior studies have searched for features that distinguish targets of approved drugs from
98 the rest of the genome (or a representative subset) (13, 15-25). Using the remainder of the
99 genome for comparison has been useful for finding features enriched among successful targets,
100 but it is uncertain whether these features are specific to successful targets or are enriched among
101 targets of failed drug candidates as well. Our study aims to fill this knowledge gap by directly
102 testing for features that separate targets by clinical outcome, expanding the scope of prior studies
103 that have investigated how genetic disease associations (36) and publication trends (37) of
104 targets correlate with clinical outcome.

105

106 Our work has five additional innovative characteristics. First, we included only targets of drugs
107 that are presumed to be selective (no documented polypharmacology) to reduce ambiguity in
108 assigning clinical trial outcomes to targets. Second, we included only phase III failures to enrich

109 for target efficacy failures, as opposed to safety and target engagement failures, which are more
110 common in phase I and phase II (2). Third, we excluded targets of assets only indicated for
111 cancer, as studies have observed that features of successful targets for cancer differ from features
112 of successful targets for other indications (22, 23), moreover, cancer trials fail more frequently
113 than trials for other indications (2). Fourth, we interrogated a diverse and comprehensive set of
114 features, over 150,000 features from 67 datasets covering 16 feature types, whereas prior studies
115 have examined only features derived from protein sequence (16-18, 24, 25), protein-protein
116 interactions (13, 15, 18-23), Gene Ontology terms (13, 15, 16), and gene expression profiles (15,
117 19, 21, 25). Fifth, because targets of drugs and drug candidates do not constitute a random
118 sample of the genome, we implemented a suite of tests to assess the robustness and
119 generalizability of features identified as significantly separating successes from failures in the
120 biased sample.

121
122 A handful of the initial 150,000+ features passed our tests for robustness and generalizability to
123 new targets or target classes. Interestingly, these features were predominantly derived from gene
124 expression datasets. *Notably, two significant features were discovered repeatedly in multiple*
125 *datasets: successful targets tended to have lower mean mRNA expression across tissues and*
126 *higher expression variance than failed targets.* We also trained a classifier to predict phase III
127 success probabilities for untested targets (no phase III clinical trial outcomes reported for drug
128 candidates that selectively modulate these targets). We identified 943 targets with sufficiently
129 unfavorable expression characteristics to be predicted twice as likely to fail in phase III clinical
130 trials as past phase III targets. Furthermore, we identified 2,700,856 target pairs predicted with
131 99% consistency to have a 2-fold difference in success probability. Such pairwise comparisons

132 may be useful for prioritizing short lists of targets under consideration for a therapeutic program.

133 We conclude this paper with a discussion of the biases and limitations faced when attempting to

134 analyze, model, or interpret data on clinical trial outcomes.

135

136 **RESULTS**

137

138 **Examples of successful and failed targets obtained from phase III clinical trial reports**

139

140 We extracted phase III clinical trial outcomes reported in Pharmaprojects (38) for drug
141 candidates reported to be selective (single documented target) and tested as treatments for non-
142 cancer diseases. We grouped the outcomes by target, scored targets with at least one approved
143 drug as successful ($N_S=259$), and scored targets with no approved drugs and at least one
144 documented phase III failure as failed ($N_F=72$) (Supplementary Table S1). The target success
145 rate (77%) appears to be inflated relative to typically reported phase III success rates (58%) (2)
146 because we scored targets by their best outcome across multiple trials.

147

148 **Comprehensive and diverse collection of target features obtained from the Harmonizome**

149

150 We obtained target features from the Harmonizome (39), a recently published collection of
151 features of genes and proteins extracted from over 100 Omics datasets. We limited our analysis
152 to 67 datasets that are in the public domain or GSK had independently licensed (Table 1). Each
153 dataset in the Harmonizome is organized into a matrix with genes labeling the rows and features
154 such as diseases, phenotypes, tissues, and pathways labeling the columns. We included the mean
155 and standard deviation calculated along the rows of each dataset as additional target features.
156 These summary statistics provide potentially useful and interpretable information about targets,
157 such as how many pathway associations a target has or how variable a target's expression is
158 across tissues.

159

160

Table 1. Datasets tested for features significantly separating successful targets from failed targets.

Dataset	Feature Type	Total Genes	Covered Samples	Total Features	Covered Features	Reduced Features
Roadmap Epigenomics Cell and Tissue DNA Methylation Profiles	cell or tissue DNA methylation	13835	227	26	26	4
Allen Brain Atlas Adult Human Brain Tissue Gene Expression Profiles	cell or tissue expression	17979	287	416	416	2
Allen Brain Atlas Adult Mouse Brain Tissue Gene Expression Profiles	cell or tissue expression	14248	287	2234	2234	2
BioGPS Human Cell Type and Tissue Gene Expression Profiles	cell or tissue expression	16383	320	86	86	2
BioGPS Mouse Cell Type and Tissue Gene Expression Profiles	cell or tissue expression	15443	313	76	76	2
GTEx Tissue Gene Expression Profiles	cell or tissue expression	26005	328	31	31	2
GTEx Tissue Sample Gene Expression Profiles	cell or tissue expression	19250	301	2920	2920	2
HPA Cell Line Gene Expression Profiles	cell or tissue expression	15868	259	45	45	1
HPA Tissue Gene Expression Profiles	cell or tissue expression	17496	314	33	33	2
HPA Tissue Protein Expression Profiles	cell or tissue expression	15788	266	46	46	11
HPA Tissue Sample Gene Expression Profiles	cell or tissue expression	16742	300	123	123	2
HPM Cell Type and Tissue Protein Expression Profiles	cell or tissue expression	7274	94	6	6	2
ProteomicsDB Cell Type and Tissue Protein Expression Profiles	cell or tissue expression	2776	28	55	55	5
Roadmap Epigenomics Cell and Tissue Gene Expression Profiles	cell or tissue expression	12824	164	59	59	6
TISSUES Curated Tissue Protein Expression Evidence Scores	cell or tissue expression	16216	317	645	245	106
TISSUES Experimental Tissue Protein Expression Evidence Scores	cell or tissue expression	17922	316	245	244	44
TISSUES Text-mining Tissue Protein Expression Evidence Scores	cell or tissue expression	16184	330	4189	2974	2118
ENCODE Histone Modification Site Profiles	cell or tissue histone modification sites	22382	330	437	432	91
Roadmap Epigenomics Histone Modification Site Profiles	cell or tissue histone modification sites	21032	313	385	295	282
ENCODE Transcription Factor Binding Site Profiles	cell or tissue transcription factor binding sites	22845	330	1681	1591	723
JASPAR Predicted Transcription Factor Targets	cell or tissue transcription factor binding sites	21547	330	113	80	77
COMPARTMENTS Curated Protein Localization Evidence Scores	cellular compartment associations	16738	330	1465	228	105
COMPARTMENTS Experimental Protein Localization Evidence Scores	cellular compartment associations	6495	73	61	37	10
COMPARTMENTS Text-mining Protein Localization Evidence Scores	cellular compartment associations	14375	330	2083	877	545
GO Cellular Component Annotations	cellular compartment associations	16757	328	1549	208	124
LOCATE Curated Protein Localization Annotations	cellular compartment associations	9639	269	80	50	20
LOCATE Predicted Protein Localization Annotations	cellular compartment associations	19747	325	26	23	10
CTD Gene-Chemical Interactions	chemical interactions	11125	321	9518	2222	2042
Guide to Pharmacology Chemical Ligands of Receptors	chemical interactions	899	209	4896	189	52
Kinativ Kinase Inhibitor Bioactivity Profiles	chemical interactions	232	9	28	28	25
KinomeScan Kinase Inhibitor Targets	chemical interactions	287	10	75	75	72

CMPA Signatures of Differentially Expressed Genes for Small Molecules	chemical perturbation differentially expressed genes	12148	300	6102	5066	5065
ClinVar SNP-Phenotype Associations	disease or phenotype associations	2458	143	3293	3	2
CTD Gene-Disease Associations	disease or phenotype associations	21582	331	6327	2926	2116
dbGAP Gene-Trait Associations	disease or phenotype associations	5668	147	512	51	49
DISEASES Curated Gene-Disease Assocation Evidence Scores	disease or phenotype associations	2252	115	772	94	49
DISEASES Experimental Gene-Disease Assocation Evidence Scores	disease or phenotype associations	4055	131	352	106	43
DISEASES Text-mining Gene-Disease Assocation Evidence Scores	disease or phenotype associations	15309	330	4630	2559	1850
GAD Gene-Disease Associations	disease or phenotype associations	10705	318	12780	1189	980
GAD High Level Gene-Disease Associations	disease or phenotype associations	8016	314	20	19	16
GWAS Catalog Gene-Disease Associations	disease or phenotype associations	4356	127	1009	30	28
GWASdb SNP-Disease Associations	disease or phenotype associations	11805	253	587	252	126
GWASdb SNP-Phenotype Associations	disease or phenotype associations	12488	261	824	397	150
HPO Gene-Disease Associations	disease or phenotype associations	3158	171	6844	1187	667
HuGE Navigator Gene-Phenotype Associations	disease or phenotype associations	12055	322	2755	1241	1153
MPO Gene-Phenotype Associations	disease or phenotype associations	7798	299	8581	2434	1444
OMIM Gene-Disease Associations	disease or phenotype associations	4553	209	6177	5	4
GeneSigDB Published Gene Signatures	gene signatures or modules	19723	331	3517	1363	1313
MSigDB Cancer Gene Co-expression Modules	gene signatures or modules	4869	135	358	135	95
MiRTarBase microRNA Targets	microRNA targets	12086	218	598	93	91
TargetScan Predicted Conserved microRNA Targets	microRNA targets	14923	283	1539	1020	791
TargetScan Predicted Nonconserved microRNA Targets	microRNA targets	18210	324	1541	1534	1236
GO Biological Process Annotations	pathway, function, or process associations	15717	328	13214	2436	1215
GO Molecular Function Annotations	pathway, function, or process associations	15777	327	4164	367	204
HumanCyc Pathways	pathway, function, or process associations	932	41	288	11	8
KEGG Pathways	pathway, function, or process associations	7016	298	303	185	179
PANTHER Pathways	pathway, function, or process associations	1962	138	147	40	39
Reactome Pathways	pathway, function, or process associations	9005	309	1814	289	159
Wikipathways Pathways	pathway, function, or process associations	4958	263	301	140	137
DEPOD Substrates of Phosphatases	phosphatase interactions	293	19	114	13	9
NURSA Protein Complexes	protein complex associations	9785	141	1798	1182	1181
InterPro Predicted Protein Domain Annotations	protein domain associations	18002	329	11017	119	63
BioGRID Protein-Protein Interactions	protein interactions	15270	306	15272	1191	1163
DIP Protein-Protein Interactions	protein interactions	2709	140	2711	32	24
Guide to Pharmacology Protein Ligands of Receptors	protein interactions	187	46	213	5	4

IntAct Biomolecular Interactions	protein interactions	12303	269	12305	422	417
GTEx eQTL	SNP eQTL targets	7898	107	7817	2	1
TOTALS	NA	NA	NA	174228	44092	28562

161

162 The datasets contained a total of 174,228 features covering 16 feature types (Table 1). We
163 restricted our analysis to 44,092 features that had at least three non-zero values for targets
164 assigned a phase III outcome. Many datasets had strong correlations among their features. To
165 reduce feature redundancy and avoid excessive multiple hypothesis testing while maintaining
166 interpretability of features, we replaced each group of highly correlated features with the group
167 mean feature and assigned it a representative label (Fig 1, Supplementary Table S2). The number
168 of features shrunk to 28,562 after reducing redundancy.

169

170 **Fig 1. Feature Selection Pipeline.** Each dataset took the form of a matrix with genes labeling the rows and features
171 labeling the columns. We appended the mean and standard deviation computed across all features as two additional
172 features. **Step 1:** We filtered the columns to eliminate redundant features, replacing each group of correlated
173 features with the group average feature, where a group was defined as features with squared pair-wise correlation
174 coefficient $r^2 \geq 0.5$. If the dataset mean feature was included in a group of correlated features, we replaced the group
175 with the dataset mean. **Step 2:** We filtered the rows for targets with clinical trial outcomes of interest: targets of
176 selective drugs approved for non-cancer indications (successes) and targets of selective drug candidates that failed in
177 phase III clinical trials for non-cancer indications (failures). **Step 3:** We tested the significance of each feature as an
178 indicator of success or failure using permutation tests to quantify the significance of the difference between the
179 means of the successful and failed targets. We corrected for multiple hypothesis testing using the Benjamini-
180 Yekutieli method to control the false discovery rate at 0.05 within each dataset. **Step 4:** We “stressed” the
181 significant features with additional tests to assess their robustness and generalizability. For example, we used
182 bootstrapping to estimate probabilities that the significance findings will replicate on similar sets of targets.

183

184 **Target features tested for correlation with phase III outcome**

185

186 We performed permutation tests (40, 41) on the remaining 28,562 target features to find features
187 with a significant difference between the successful and failed targets, and we corrected p-values
188 for multiple hypothesis testing using the Benjamini-Yekutieli method (42) (Fig 1, Supplementary
189 Table S2). We used permutation testing to apply the same significance testing method to all
190 features, since they had heterogeneous data distributions. We detected 19 features correlated with
191 clinical outcome at a within-dataset false discovery rate of 0.05 (Table 2). The significant
192 features were derived from 7 datasets, of which 6 datasets were gene expression atlases: Allen
193 Brain Atlas adult human brain tissues (43, 44), Allen Brain Atlas adult mouse brain tissues (43,
194 45), BioGPS human cell types and tissues (46-48), BioGPS mouse cell types and tissues (46-48),
195 Genotype-Tissue Expression Project (GTEx) human tissues (49, 50), and Human Protein Atlas
196 (HPA) human tissues (51). The remaining dataset, TISSUES (52), was an integration of
197 experimental gene and protein tissue expression evidence from multiple sources. Two
198 correlations were significant in multiple datasets: successful targets tended to have lower mean
199 expression across tissues and higher expression variance than failed targets.

200

201 **Table 2. Features significantly correlated with phase III outcome.**

Dataset	Feature	Corr Pval	Correlation Sign	Correlated Target Classes (and sign)	Repl Prob (Bootstrap)	Repl Prob (Class Holdout Bootstrap)	Repl Prob (Within Class Permutation Bootstrap)
BioGPS Human Cell Type and Tissue Gene Expression Profiles	[mean]	0.001	-1	GPCRs (-1)	0.89	0.98	0.83
BioGPS Human Cell Type and Tissue Gene Expression Profiles	stdv	0.010	-1	GPCRs (-1), Integrins (+1)	0.69	0.56	0.32
BioGPS Mouse Cell Type and Tissue Gene Expression Profiles	[mean]	0.042	-1	GPCRs (-1)	0.55	0.74	0.56
Allen Brain Atlas Adult Human Brain Tissue Gene Expression Profiles	[mean]	0.006	-1	GPCRs (-1)	0.78	0.80	0.78
Allen Brain Atlas Adult Mouse Brain Tissue Gene Expression Profiles	r3 roof plate	0.002	-1	None	0.88	1.00	0.89
Allen Brain Atlas Adult Mouse Brain Tissue Gene Expression Profiles	[mean]	0.007	-1	None	0.76	1.00	0.79
GTEx Tissue Gene Expression Profiles	[mean]	0.014	-1	GPCRs (-1)	0.65	0.60	0.76
GTEx Tissue Gene Expression Profiles	stdv	0.014	+1	GPCRs (+1)	0.69	0.94	0.76

HPA Tissue Gene Expression Profiles	[mean]	0.004	-1	GPCRs (-1)	0.80	0.90	0.85
HPA Tissue Gene Expression Profiles	stdv	0.004	+1	None	0.81	1.00	0.81
TISSUES Experimental Tissue Protein Expression Evidence Scores	bone marrow	0.001	-1	GPCRs (-1)	0.92	0.96	<i>0.66</i>
TISSUES Experimental Tissue Protein Expression Evidence Scores	[hematopoietic cells]	0.001	-1	GPCRs (-1), Integrins (+1)	0.93	1.00	<i>0.72</i>
TISSUES Experimental Tissue Protein Expression Evidence Scores	[mean]	0.001	-1	GPCRs (-1)	0.85	0.99	<i>0.76</i>
TISSUES Experimental Tissue Protein Expression Evidence Scores	[epithalamus and pineal gland]	0.012	-1	None	<i>0.73</i>	0.97	<i>0.49</i>
TISSUES Experimental Tissue Protein Expression Evidence Scores	erythroid cell	0.015	-1	None	<i>0.68</i>	0.94	<i>0.45</i>
TISSUES Experimental Tissue Protein Expression Evidence Scores	[t-lymphocyte]	0.017	-1	None	<i>0.65</i>	0.95	<i>0.65</i>
TISSUES Experimental Tissue Protein Expression Evidence Scores	[miscellaneous tissues]	0.017	-1	GPCRs (-1)	<i>0.64</i>	<i>0.64</i>	<i>0.63</i>
TISSUES Experimental Tissue Protein Expression Evidence Scores	[thymus and thorax]	0.017	-1	Integrins (+1)	<i>0.60</i>	<i>0.37</i>	<i>0.44</i>
TISSUES Experimental Tissue Protein Expression Evidence Scores	adrenal cortex	0.043	-1	None	<i>0.44</i>	<i>0.62</i>	<i>0.45</i>
Footnotes							
Abbreviations: Corr Pval = p-value corrected for multiple hypothesis testing, Repl Prob = replication probability.							
[Square brackets] denote groups of features.							
[miscellaneous tissues] is a heterogeneous group of digestive, respiratory, urogenital, reproductive, nervous, cardiovascular, and hematopoietic system tissues.							
White background indicates features that passed all tests for robustness and generalizability.							
Gray background indicates features that failed at least one test for robustness or generalizability. <i>Strikethrough italics</i> indicates the failed test(s).							

202

203 **Significant features tested for robustness to sample variation and generalization across**
204 **target classes**

205

206 Because targets of drugs and drug candidates do not constitute a random sample of the genome,
207 features that separate successful targets from failed targets in our sample may perform poorly as
208 genome-wide predictors of success versus failure. We performed three analyses to address this
209 issue (Fig 1).

210

211 *Robustness to sample variation*

212

213 We used bootstrapping (53, 54) (sampling with replacement from the original set of examples to
214 construct sets of examples equal in size to the original set) to investigate how robust our
215 significance findings were to variation in the success and failure examples. For each dataset that
216 yielded significant features in our primary analysis, we repeated the analysis on 1000 bootstrap
217 samples and quantified the replication probability (55) of each feature as the fraction of
218 bootstraps yielding a significant correlation with phase III outcome at a within-dataset false
219 discovery rate of 0.05. Twelve features had less than 80% probability (considered a strong
220 replication probability in (55)) that their correlation with clinical outcome will generalize to new
221 examples (Table 2).

222

223 *Robustness to target class variation*

224

225 We tested if any of the significance findings depended upon the presence of targets from a single
226 target class in our sample. We obtained target class labels (i.e. gene family labels) from the
227 HUGO Gene Nomenclature Committee (56), tested if any target classes were significantly
228 correlated with phase III outcome, and then tested if these classes were correlated with any
229 features. The GPCR and integrin classes were correlated with phase III outcome as well as
230 several features (Table 2). This raised the possibility that instead of these features being genome-
231 wide indicators of clinical outcome, they were simply reflecting the fact that many GPCRs have
232 succeeded (62/70, $p < 0.05$) or that integrins have failed (3/3, $p < 0.01$). To test this possibility, we
233 repeated the bootstrapping procedure described above to obtain replication probabilities, except
234 excluded GPCRs and integrins from being drawn in the bootstrap samples. Six features had less

235 than 80% probability that their correlation with clinical outcome will generalize to new target
236 classes (Table 2).

237

238 *Generalization across target classes*

239

240 In the preceding analysis, we checked one target class at a time for its impact on our significance
241 findings. To broadly test whether features generalize across target classes, we repeated the
242 permutation testing described in our initial analysis, but only shuffled the success/failure labels
243 within target classes, inspired by the work of Epstein et al. (57) on correcting for confounders in
244 permutation testing. By generating a null distribution with preserved ratio of successes to failures
245 within each target class, features must correlate with clinical outcome within multiple classes to
246 be significant, while features that discriminate between classes will not be significant. We
247 repeated the modified permutation tests on 1000 bootstrap samples to obtain replication
248 probabilities. We rejected fifteen features that had less than 80% probability that their correlation
249 with clinical outcome generalizes across target classes (Table 2). This set of fifteen features
250 included all features with less than 80% replication probability in either of the previous two tests.
251 The remaining robust and generalizable features were: 1) mean mRNA expression across tissues
252 (HPA and BioGPS human tissue expression datasets), 2) standard deviation of expression across
253 tissues (HPA human tissue expression dataset), and 3) expression in r3 roof plate (Allen Brain
254 Atlas adult mouse brain tissue expression dataset). The r3 roof plate expression profile was
255 correlated with mean expression across tissues in the Allen Brain Atlas dataset ($r^2=0.47$), falling
256 just below the $r^2=0.5$ cut-off that would have grouped r3 roof plate with the mean expression
257 profile during dimensionality reduction.

258

259 **Classifier-based assessment of feature usefulness and interpretability**

260

261 Statistical significance did not guarantee the remaining features would be useful in practice for
262 discriminating between successes and failures. To test their utility, we trained a classifier to
263 predict target success or failure, using cross-validation to select a model type (Random Forest or
264 logistic regression) and a subset of features useful for prediction. Because we used all targets
265 with phase III outcomes for the feature selection procedure described above, simply using the
266 final set of features to train a classifier on the same data would yield overly optimistic
267 performance, even with cross-validation. Therefore, we implemented a nested cross-validation
268 routine to perform both feature selection and model selection (58).

269

270 *Cross-validation routine*

271

272 The outer loop of the cross-validation routine had five steps (Fig 2): 1) separation of targets with
273 phase III outcomes into training and testing sets, 2) univariate feature selection using the training
274 set, 3) aggregation of features from different datasets into a single feature matrix, 4) classifier-
275 based feature selection and model selection using the training set, and 5) evaluation of the
276 classifier on the test set. Step 4 used an inner loop with 5-fold cross-validation repeated 20 times
277 to estimate the performance of different classifier types (Random Forest or logistic regression)
278 and feature subsets (created by incremental feature elimination). The simplest classifier (least
279 number of features, with logistic regression considered simpler than Random Forest) with cross-
280 validation values for area under the receiver operating characteristic curve (AUROC) and area

281 under the precision-recall curve (AUPR) within 95% of maximum was selected. The outer loop
282 used 5-fold cross-validation repeated 200 times, which provided 1000 train-test cycles for
283 estimating the generalization performance of the classifier and characterizing the consistency of
284 the selected features and model type.

285

286 **Fig 2. Modeling Pipeline.** We trained a classifier to predict phase III clinical trial outcomes, using 5-fold cross-
287 validation repeated 200 times to assess the stability of the classifier and estimate its generalization performance. For
288 each fold of cross-validation, modeling began with the non-redundant features for each dataset. **Step 1:** We split the
289 targets with phase III outcomes into training and testing sets. **Step 2:** We performed univariate feature selection
290 using permutation tests to quantify the significance of the difference between the means of the successful and failed
291 targets in the training examples. We controlled for target class as a confounding factor by only shuffling outcomes
292 within target classes. We accepted features with adjusted p-values less than 0.05 after correcting for multiple
293 hypothesis testing using the Benjamini-Yekutieli method. **Step 3:** We aggregated significant features from all
294 datasets into a single feature matrix. **Step 4:** We performed incremental feature elimination with an inner 5-fold
295 cross-validation loop repeated 20 times to select the type of classifier (Random Forest or logistic regression) and
296 smallest subset of features that had cross-validation area under the receiver operating characteristic curve (AUROC)
297 and area under the precision-recall curve (AUPR) values within 95% of maximum. **Step 5:** We refit the selected
298 model using all the training examples and evaluated its performance on the test examples.

299

300 *Classifier consistency*

301

302 Simple models were consistently selected for the classifier (Table 3, Supplementary Table S3).
303 In 1000 train-test cycles, a logistic regression model with one feature was selected most the time
304 (66%), followed in frequency by a logistic regression model with two features (8%), a Random
305 Forest model with two features (8%), and a logistic regression model with three features (6%).
306 Other combinations of model type (logistic regression or Random Forest) and number of features

307 (ranging from 1 to 8) appeared 11% of the time (each 4% or less). For one of the train-test cycles
308 (0.1%), no significant features were found in the univariate feature selection step, resulting in a
309 null model. Note that the logistic regression models were selected primarily because we imposed
310 a preference for simple and interpretable models, not because they performed better than
311 Random Forest models. The Random Forest model tended to perform as well as the logistic
312 regression model on the inner cross-validation loop, with AUROC = 0.62 ± 0.06 for Random
313 Forest and 0.63 ± 0.05 for logistic regression (Supplementary Table S4).

314

315 **Table 3. Distribution of train-test cycles by classifier type and number of selected features.**

Selected Features		Selected Model Type		Total
		Logistic Regression	Random Forest	
Selected Features	1	662	5	667
	2	82	84	166
	3	57	41	98
	4	22	2	24
	5	24	1	25
	6	11	0	11
	7	6	0	6
	8	2	0	2
	Total	866	133	999*

Footnotes

* 1 train-test cycle yielded no significant features for modeling

316

317 Gene expression features were consistently selected for the classifier (Table 4, Supplementary
318 Table S3). Mean mRNA expression across tissues and standard deviation of expression across
319 tissues had frequencies of 69% and 59%, respectively. More precisely, 36% of the models used
320 mean mRNA expression across tissues as the only feature, 31% used standard deviation of
321 expression as the only feature, and 12% used mean and standard deviation as the only two
322 features. Other expression features appeared in 21% of the models. These expression features

323 tended to be correlated with mean expression across tissues (median $r^2=0.49$). Disease
324 association features appeared in 0.4% of the models.

325

326 **Table 4. Number of train-test cycles in which feature was selected for the classifier.**

Feature Type	Feature	Count
cell or tissue expression	mean across tissues	685
cell or tissue expression	standard deviation across tissues	585
cell or tissue expression	other	214
disease or phenotype associations	mean across diseases	2
disease or phenotype associations	other	2
pathway, function, or process associations	any	1

327

328 *Classifier performance*

329

330 The classifier consistently had better than random performance in cross-validation (Fig 3, Table
331 5, Supplementary Table S5). The 2.5th, 50th, and 97.5th percentiles for AUROC were 0.51, 0.57,
332 and 0.61. For comparison, a random ordering of targets would yield an AUROC of 0.50. The
333 receiver operating characteristic curve showed that there was no single cut-off that would
334 provide satisfactory discrimination between successes and failures (Fig 3A). For an alternative
335 view, we used kernel density estimation (59) to fit distributions of the probability of success
336 predicted by the classifier for the successful, failed, and unlabeled targets (Fig 3B,
337 Supplementary Table S1). The distributions for successes and failures largely overlapped, except
338 in the tails.

339

340 **Fig 3. Classifier Performance.** (A) Receiver operating characteristic (ROC) curve. The solid black line indicates
341 the median performance across 200 repetitions of 5-fold cross-validation and the gray area indicates the range of the
342 2.5 and 97.5 percentiles. The dotted black line indicates the performance of random rankings. (B) Distributions of
343 the probability of success predicted by the classifier for the successful, failed, and unlabeled targets. (C) Precision-

344 recall curve for success predictions. **(D)** Precision-recall curve for failure predictions. **(E)** Pairwise target
345 comparisons. For each pair of targets, we computed the fraction of repetitions of cross-validation in which Target B
346 had a higher predicted probability of success greater than Target A. The heatmap illustrates this fraction, thresholded
347 at 0.95 or 0.99, plotted as a function of the median predicted probabilities of success of two targets. The upper left
348 region is where the classifier is 95% (above solid black line) or 99% (above dotted blue line) consistent in predicting
349 greater probability of success of Target B than Target A. **(F)** Relationship between features and phase III outcomes.
350 Heat map showing the projection of the predicted success probabilities onto the two dominant features selected for
351 the classifier: mean expression across tissues and standard deviation of expression across tissues. Red, white, and
352 blue background colors correspond to 1, 0.5, and 0 success probabilities. Red plusses and blue crosses mark the
353 locations of the success and failure examples. It appears the model has learned that failures tend to have high mean
354 expression and low standard deviation of expression across tissues, while successes tend to have low mean
355 expression and high standard deviation of expression. The success and failure examples are not well separated,
356 indicating that we did not discover enough features to fully explain why targets succeed or fail in phase III clinical
357 trials.

358

359 **Table 5. Classifier performance statistics.**

Statistic	2.5 Percentile	Median	97.5 Percentile
True Positives (TP)	91	220	243
False Positives (FP)	16	52	65
True Negatives (TN)	5	16	52
False Negatives (FN)	1	24	154
True Positive Rate (TPR)	0.370	0.903	0.995
False Positive Rate (FPR)	0.232	0.762	0.928
False Negative Rate (FNR)	0.005	0.096	0.630
True Negative Rate (TNR)	0.072	0.237	0.768
Misclassification Rate (MCR)	0.206	0.241	0.542
Accuracy (ACC)	0.458	0.759	0.794
False Discovery Rate (FDR)	0.149	0.194	0.213
Positive Predictive Value (PPV)	0.787	0.806	0.851
False Omission Rate (FOMR)	0.233	0.583	0.741
Negative Predictive Value (NPV)	0.259	0.417	0.767
Area Under Receiver Operating Characteristic Curve (AUROC)	0.512	0.574	0.615
Area Under Precision-Recall Curve (AUPR)	0.777	0.811	0.836

Positive Likelihood Ratio (PLR)	1.058	1.184	1.619
Negative Likelihood Ratio (NLR)	0.086	0.402	0.819
Diagnostic Odds Ratio (DOR)	1.748	3.066	13.344
Risk Ratio (RR)	1.143	1.387	3.447
Matthews Correlation Coefficient (MCC)	0.100	0.178	0.251

360

361 We attempted to identify subsets of targets with high positive predictive value (PPV) or high
362 negative predictive value (NPV). The median PPV rose as high as 0.99, but uncertainty in the
363 PPV was so large that we could not be confident in identifying any subset of targets with a
364 predicted success rate better than the historical 0.77 (Fig 3C). The median NPV rose to 0.40,
365 roughly twice the historical failure rate of 0.23. Furthermore, at 0.40 median NPV, 99% of the
366 cross-validation repetitions had an NPV greater than the historical failure rate (Fig 3D). Using
367 this cut-off, we identified 943 unlabeled targets expected to be twice as likely to fail in phase III
368 clinical trials as past phase III targets.

369

370 We reasoned that a more practical use of the classifier would be to make pair-wise comparisons
371 among a short list of targets already under consideration for a therapeutic program. To assess the
372 utility of the classifier for this purpose, for every pair of targets T_A and T_B , we computed the
373 fraction of cross-validation runs in which the classifier predicted greater probability of success
374 for T_B than T_A . We identified 67,270,678 target pairs (39%) with at least a 0.1 difference in
375 median success probability where the classifier was 95% consistent in predicting greater
376 probability of success for T_B than T_A . The classifier was 99% consistent for 41528043 target
377 pairs (24%). Requiring at least a 2-fold difference in median success probability between T_B and
378 T_A reduced these counts to 2,730,437 target pairs (1.6%) at 95% consistency and 2,700,856
379 target pairs (1.6%) at 99% consistency. We visualized these results by plotting the 95% and 99%
380 consistency fraction thresholds smoothly interpolated as a function of the median predicted

381 probabilities of success of T_A and T_B (Fig 3E). For a median probability of success of T_A around
382 0.2, T_B must have a median probability of success of 0.5 or greater at the 99% threshold. For
383 lower T_A success probabilities, the T_B success probability must be even higher because there is
384 greater uncertainty about the low T_A probabilities. For higher T_A success probabilities, the T_B
385 success probability at the 99% threshold increases steadily until a T_A success probability of about
386 0.6, where the T_B success probability reaches 1. For T_A success probabilities above 0.6, no
387 targets are predicted to have greater probability of success with 99% consistency.

388

389 *Feature interpretation*

390

391 To interpret the relationship inferred by the classifier between target features and outcomes, we
392 created a heatmap of the probability of success predicted by the classifier projected onto the two
393 features predominantly selected for the model: mean expression and standard deviation of
394 expression across tissues (Fig 3F). The probability of success was high in the subspace with low
395 mean expression and high standard deviation of expression, and transitioned to low probability in
396 the subspace with high mean expression and low standard deviation of expression. This trend
397 appeared to be consistent with the distribution of the success and failure examples in the space.

398

399 **DISCUSSION**

400

401 **Gene expression predicts phase III outcome**

402

403 We searched over 150,000 target features from 67 datasets covering 16 feature types for
404 predictors of target success or failure in phase III clinical trials (Table 1, Fig 1). We found
405 several features significantly correlated with phase III outcome, robust to re-sampling, and
406 generalizable across target classes (Table 2). To assess the usefulness of such features, we
407 implemented a nested cross-validation routine to select features, train a classifier to predict the
408 probability a target will succeed in phase III clinical trials, and estimate the stability and
409 generalization performance of the model (Figs 2 and 3, Tables 3, 4, and 5). Ultimately, we found
410 two features useful for predicting success or failure of targets in phase III clinical trials.
411 Successful targets tended to have low mean mRNA expression across tissues and high standard
412 deviation of mRNA expression across tissues (Fig 3F). These features were significant in
413 multiple gene expression datasets, which increased our confidence that their relationship to phase
414 III outcome was real, at least for the targets in our sample, which included only targets of
415 selective drugs indicated for non-cancer diseases.

416

417 One interpretation of why the gene expression features were predictive of phase III outcome is
418 that they are informative of the specificity of a target's expression across tissues. A target with
419 tissue specific expression would have a high standard deviation relative to its mean expression
420 level. Tissue specific expression has been proposed by us and others as a favorable target
421 characteristic in the past (4, 14, 60-62), but the hypothesis had not been evaluated empirically

422 using examples of targets that have succeeded or failed in clinical trials. For a given disease, if a
423 target is expressed primarily in the disease tissue, it is considered more likely that a drug will be
424 able to exert a therapeutic effect on the disease tissue while avoiding adverse effects on other
425 tissues. Additionally, specific expression of a target in the tissue affected by a disease could be
426 an indicator that dysfunction of the target truly causes the disease.

427

428 The distribution of the success and failure examples in feature space (Fig 3F) partially supports
429 the hypothesis that tissue specific expression is a favorable target feature. Successes were
430 enriched among targets with low mean expression and high standard deviation of expression
431 (tissue specific expression), and failures were enriched among targets with high mean expression
432 and low standard deviation of expression (constitutive expression). However, it does not hold in
433 general that, at any given mean expression level, targets with high standard deviation of
434 expression tend to be more successful than targets with low standard deviation of expression.
435 Nevertheless, our results encourage further investigation of the relationship between tissue
436 specific expression and clinical trial outcomes. Deeper insight may be gleaned from analysis of
437 gene expression features explicitly designed to quantify specificity of a target's expression in the
438 tissue(s) affected by the disease treated in each clinical trial.

439

440 **Caveats and limitations**

441

442 Latent factors (variables unaccounted for in this analysis) could confound relationships between
443 target features and phase III outcomes. For example, diseases pursued vary from target to target,
444 and a target's expression across tissues may be irrelevant for diseases where drugs can be

445 delivered locally or for Mendelian loss-of-function diseases where treatment requires systemic
446 replacement of a missing or defective protein. Also, clinical trial failure rates vary across disease
447 classes (2). Although we excluded targets of cancer therapeutics from our analysis, we otherwise
448 did not control for disease class as a confounding explanatory factor. Modalities (e.g. small
449 molecule, antibody, antisense oligonucleotide, gene therapy, or protein replacement) and
450 directions (e.g. activation or inhibition) of target modulation also vary from target to target and
451 could be confounding explanatory factors or alter the dependency between target features and
452 outcomes.

453

454 The potential issues described above are symptoms of the fact that our analysis (and any analysis
455 of clinical trial outcomes) attempts to draw conclusions from a small (roughly 300 targets) and
456 biased sample (63, 64). Latent factors such as target classes, disease classes, modalities, and
457 directions of target modulation are not uniformly represented in the sample, yet correlations
458 between target features and clinical trial outcomes likely depend on these factors. Unfortunately,
459 attempts to stratify, match, or otherwise control for these factors are limited by the sample size.
460 (The number of combinations of target class, disease class, modality, and direction of modulation
461 exceeds the sample size.) We employed several tests to build confidence that our findings
462 generalize across target classes, but did not address other latent factors. Consequently, we cannot
463 be sure that conclusions drawn from this study apply equally to targets modulated in any
464 direction, by any means, to treat any disease. For specific cases, expert knowledge and common
465 sense should be relied upon to determine whether conclusions from this study (or similar studies)
466 are relevant.

467

468 Another limitation is selection bias (63, 64). Targets of drugs are not randomly selected from the
469 genome and cannot be considered representative of the population of all possible targets.
470 Likewise, diseases treated by drugs are not randomly chosen; therefore, phase III clinical trial
471 outcomes for each target cannot be considered representative of the population of all possible
472 outcomes. Although we implemented tests to build confidence that our findings can generalize to
473 new targets and new target classes, ultimately, no matter how we dissect the sample, a degree of
474 uncertainty will always remain about the relevance of any findings for new targets that lack a
475 representative counterpart in the sample.

476
477 Additionally, data processing and modeling decisions have introduced bias into the analysis. For
478 example, we scored each target as successful or failed by its best outcome in all applicable
479 (selective drug, non-cancer indication) phase III clinical trials. This approach ignores nuances. A
480 target that succeeded in one trial and failed in all others is treated as equally successful as a target
481 that succeeded in all trials. Also, the outcome of a target tested in a single trial is treated as
482 equally certain as the outcome of a target tested in multiple trials. Representing target outcomes
483 as success rates or probabilities may provide better signal for discovering features predictive of
484 outcomes.

485
486 Another decision was to use datasets of features as we found them, rather than trying to reason
487 about useful features that could be derived from the original data. Because of the breadth of data
488 we interrogated, the effort and expertise necessary to hand engineer features equally well across
489 all datasets exceeded our resources. Others have had success hand engineering features for
490 similar applications in the past, particularly with respect to computing topological properties of

491 targets in protein-protein interaction networks (18, 20, 21). This analysis could benefit from such
492 efforts, potentially changing a dataset or feature type from yielding no target features correlated
493 with phase III outcomes to yielding one or several useful features (22). On a related point,
494 because we placed a priority on discovering interpretable features, we performed dimensionality
495 reduction by averaging groups of highly correlated features and concatenating their (usually
496 semantically related) labels. Dimensionality reduction by principal components analysis (65) or
497 by training a deep auto-encoder (66) could yield more useful features, albeit at the expense of
498 interpretability.

499
500 We cannot stress enough the importance of taking care not to draw broad conclusions from our
501 study, particularly with respect to the apparent dearth of features predictive of target success or
502 failure. We examined only a specific slice of clinical trial outcomes (phase III trials of selective
503 drugs indicated for non-cancer diseases) summarized in a particular way (net outcome per target,
504 as opposed to outcome per target-indication pair). Failure of a feature to be significant in our
505 analysis should not be taken to mean it has no bearing on target selection. For example, prior
506 studies have quantitatively shown that genetic evidence of disease association(s) is a favorable
507 target characteristic (3, 36), but we did not find a significant correlation between genetic
508 evidence and target success in phase III clinical trials. Our finding is consistent with the work of
509 Nelson et al. (36), who investigated the correlation between genetic evidence and drug
510 development outcomes at all phases and found a significant correlation overall and at all phases
511 of development except phase III. As a way of checking our work, we applied our methods to test
512 for features that differ between targets of approved drugs and the remainder of the druggable
513 genome (instead of targets of phase III failures), and we recovered the finding of Nelson et al.

514 that targets of approved drugs have significantly more genetic evidence than the remainder of the
515 druggable genome (Supplementary Table S6). This example serves as a reminder to be cognizant
516 of the domain of applicability of research findings. Though we believe we have performed a
517 rigorous and useful analysis, we have shed light on only a small piece of a large and complex
518 puzzle.

519
520 Advances in machine learning enable and embolden us to create potentially powerful predictive
521 models for target selection. However, as described in the limitations, scarce training data are
522 available, the data are far from ideal, and we must be cautious about building models with biased
523 data and interpreting their predictions. For example, many features that appeared to be
524 significantly correlated with phase III clinical trial outcomes in our primary analysis did not hold
525 up when we accounted for target class selection bias. This study highlights the need for both
526 domain knowledge and modeling expertise to tackle such challenging problems.

527

528 **Conclusion**

529
530 Our analysis revealed several features that significantly separated targets of approved drugs from
531 targets of drug candidates that failed in phase III clinical trials. This suggested that it is feasible
532 to construct a model integrating multiple interpretable target features derived from Omics
533 datasets to inform target selection. Only features derived from tissue expression datasets were
534 promising predictors of success versus failure in phase III, specifically, mean mRNA expression
535 and standard deviation of expression across tissues. Although these features were significant at a
536 false discovery rate cut-off of 0.05, their effect sizes were too small to be useful for classification

537 of the majority of untested targets, however, even a two-fold improvement in target quality can
538 dramatically increase R&D productivity (67). We identified 943 targets predicted to be twice as
539 likely to fail in phase III clinical trials as past phase III targets, and, therefore, should be flagged
540 as having unfavorable expression characteristics. We also identified 2,700,856 target pairs
541 predicted with 99% consistency to have a 2-fold difference in success probability, which could
542 be useful for prioritizing short lists of targets with attractive disease relevance.

543

544 It should be noted that our analysis was not designed or powered to show that specific datasets or
545 data types have no bearing on target selection. There are many reasons why a dataset may not
546 have yielded any significant features in our analysis. In particular, data processing and filtering
547 choices could determine whether or not a dataset or data type has predictive value. Also, latent
548 factors, such as target classes, disease classes, modalities, and directions of target modulation,
549 could confound or alter the dependency between target features and clinical trial outcomes.
550 Finally, although we implemented tests to ensure robustness and generalizability of the target
551 features significantly correlated with phase III outcomes, selection bias in the sample of targets
552 available for analysis is a non-negligible limitation of this study and others of its kind.
553 Nevertheless, we are encouraged by our results and anticipate deeper insights and better models
554 in the future, as researchers improve methods for handling sample biases and learn more
555 informative features.

556

557 **METHODS**

558

559 **Data**

560

561 *Clinical Outcomes*

562

563 We extracted data from Citeline's Pharmaprojects database (38) (downloaded May 27, 2016),
564 reformatting available XML data into a single tab-delimited form having one row for each asset
565 (i.e. drug or drug candidate)/company combination. For each asset, known targets, identified
566 with EntrezGene (68) IDs and symbols, and indications are reported. We obtained 107,120 asset-
567 indication pairs and 37,211 asset-target pairs, correcting a single outdated EntrezGene ID, for
568 SCN2A, which we updated from 6325 to 6326.

569

570 An overall pipeline status of each asset (e.g. "Launched", "Discontinued", "No Development
571 Reported") is reported in a single field ("Status"), and detailed information for each indication
572 being pursued is dispersed throughout several other fields (e.g., "Key Event Detail",
573 "Overview", etc.). While many assets have been tried against a single indication, and thus the
574 status of the asset-indication pair is certain, the majority (N=61,107) of asset-indication pairs are
575 for assets with multiple indications. For those pairs, we used a combination of string searching of
576 these fields and manual review of the results to determine the likely pipeline location and status
577 of each indication. For example, we excluded efforts where a trial of an asset was reported as
578 planned, but no further information was available. Asset-indication pairs were thus assigned a
579 status of Successful ("Launched", "Registered", or "Pre-registration"), Failed ("Discontinued",

580 “No Development Reported”, “Withdrawn”, or “Suspended”), or In Progress, consisting of
581 9,337, 72,269 and 25,159 pairs, respectively. We then used the pipeline location to assign each
582 asset-indication pair to one of 10 outcomes: Succeeded, In Progress-Preclinical, In Progress-
583 Phase I, In Progress-Phase II, In Progress-Phase III, Failed-Preclinical, Failed-Phase I, Failed-
584 Phase II, Failed-Phase III, and Failed-Withdrawn. We discarded indications which were
585 diagnostic in nature or unspecified, mapping the remainder to Medical Subject Headings (MeSH)
586 (69). We also observed that only 24% of the failures reported in Pharmaprojects are clinical
587 failures, suggesting a clinical success rate of nearly 35%, much higher than typically cited (67).

588

589 We joined the list of asset-indication-outcome triples with the list of asset-target pairs to produce
590 a list of asset-target-indication-outcome quadruples. We then filtered the list to remove: 1) assets
591 with more than one target, 2) non-human targets, 3) cancer indications (indications mapped to
592 MeSH tree C04), and 4) outcomes labeled as In Progress at any stage or Failed prior to Phase III.
593 We scored the remaining targets (N=331) as Succeeded (N=259), if the target had at least one
594 successful asset remaining in the list, or Failed (N=72), otherwise.

595

596 *Target Features*

597

598 We obtained target features from the Harmonizome (39), a recently published collection of
599 features of genes and proteins extracted from over 100 Omics datasets. We downloaded (on June
600 30, 2016) a subset of Harmonizome datasets that were in the public domain or GSK had
601 independently licensed (Table 1). Each dataset was structured as a matrix with genes labeling the
602 rows and features such as diseases, phenotypes, tissues, and pathways labeling the columns.

603 Genes were identified with EntrezGene IDs and symbols, enabling facile integration with the
604 clinical outcome data from Pharmaprojects. Some datasets were available on the Harmonizome
605 as a “cleaned” version and a “standardized” version. In all instances, we used the cleaned
606 version, which preserved the original data values (e.g. gene expression values), as opposed to the
607 standardized version, in which the original data values were transformed into scores indicating
608 relative strengths of gene-feature associations intended to be comparable across datasets. The
609 data matrices were quantitative and filled-in (e.g. gene expression measured by microarray),
610 quantitative and sparse (e.g. protein expression measured by immunohistochemistry), or
611 categorical (i.e. binary) and sparse (e.g. pathway associations curated by experts). We
612 standardized quantitative, filled-in features by subtracting the mean and then dividing by the
613 standard deviation. We scaled quantitative, sparse features by dividing by the mean. We included
614 the mean and standard deviation calculated along the rows of each dataset as additional target
615 features. We excluded features that had fewer than three non-zero values for the targets with
616 phase III clinical trial outcomes. The remaining features, upon which our study was based, have
617 been deposited at <https://github.com/arouillard/omic-features-successful-targets>.

618

619 **Dimensionality Reduction**

620

621 Our goals in performing dimensionality reduction were to identify groups of highly correlated
622 features, avoid excessive multiple hypothesis testing, and maintain interpretability of features.
623 For each dataset, we computed pair-wise feature correlations (r) using the Spearman correlation
624 coefficient (70-72) for quantitative, filled-in datasets, and the cosine coefficient (71, 72) for
625 sparse or categorical datasets. We thresholded the correlation matrix at $r^2=0.5$ (for the Spearman

626 correlation coefficient, this corresponds to one feature explaining 50% of the variance of another
627 feature, and for the cosine coefficient, this corresponds to one feature being aligned within 45
628 degrees of another feature) and ordered the features by decreasing number of correlated features.
629 We created a group for the first feature and its correlated features. If the dataset mean was
630 included in the group, we replaced the group of features with the dataset mean. Otherwise, we
631 replaced the group of features with the group mean and assigned it the label of the first feature
632 (to indicate that the feature represents the average of features correlated with the first feature),
633 while also retaining a list of the labels of all features included in the group. We continued
634 through the list of features, repeating the grouping process as described for the first feature,
635 except excluding features already assigned to a group from being assigned to a second group.

636

637 **Feature Selection**

638

639 We performed permutation tests (40, 41) to find features with a significant difference between
640 successful and failed targets. We used permutation testing in order to apply the same significance
641 testing method to all features. The features in our collection had heterogeneous shapes of their
642 distributions and varying degrees of sparsity, and therefore no single parametric test would be
643 appropriate for all features. Furthermore, individual features frequently violated assumptions
644 required for parametric tests, such as normality for the t-test (for continuous-valued features) or
645 having at least five observations in each entry of the contingency table for the Chi-squared test
646 (for categorical features). For each feature, we performed 10^5 success/failure label permutations
647 to obtain a null distribution for the difference between the means of successful and failed targets,
648 and then calculated an empirical two-tailed p-value as the fraction of permutations that yielded a

649 difference between means at least as extreme as the actual observed difference. We used the
650 Benjamini-Yekutieli method (42) to correct for multiple hypothesis testing within each dataset
651 and accepted features with corrected p-values less than 0.05 as significantly correlated with
652 phase III clinical trial outcomes, thus controlling the false discovery rate at 0.05 within each
653 dataset.

654

655 **Feature Robustness and Generalizability**

656

657 *Robustness to sample variation*

658

659 We used bootstrapping (53, 54) to investigate how robust our significance findings were to
660 variation in the success and failure examples. We created a bootstrap sample by sampling with
661 replacement from the original set of examples to construct an equal sized set of examples. For
662 each dataset that yielded significant features in our primary analysis, we repeated the analysis on
663 the bootstrap sample and recorded whether the features were still significant at the
664 aforementioned 0.05 false discovery rate cut-off. We performed this procedure on 1000 bootstrap
665 samples and quantified the replication probability (55) of each feature as the fraction of
666 bootstraps showing a significant correlation between the feature and phase III clinical trial
667 outcomes. We accepted features with replication probabilities greater than 0.8 (55) as robust to
668 sample variation.

669

670 *Robustness to target class variation*

671

672 We tested if any of the significance findings depended upon the presence of targets from a single
673 target class in our sample. We obtained target class labels (i.e. gene family labels) from the
674 HUGO Gene Nomenclature Committee (56) (downloaded April 19, 2016) and created binary
675 features indicating target class membership. Using the same permutation testing and multiple
676 hypothesis testing correction methods described above for feature selection, we tested if any
677 target classes were significantly correlated with phase III clinical trial outcomes. Then, we tested
678 if the significant target classes were correlated with any significant features. Such features might
679 be correlated with clinical outcome only because they are surrogate indicators for particular
680 target classes that have been historically very successful or unsuccessful, as opposed to the
681 features being predictors of clinical outcome irrespective of target class. To test this possibility,
682 we performed a bootstrapping procedure as described above, except did not allow examples from
683 target classes correlated with clinical outcome to be drawn when re-sampling. Thus, the modified
684 bootstrapping procedure provided replication probabilities conditioned upon missing information
685 about target classes correlated with clinical outcome. We accepted features with replication
686 probabilities greater than 0.8 as robust to target class variation.

687

688 *Generalization across target classes*

689

690 We implemented a modified permutation test, inspired by the approach of Epstein et al. (57) to
691 correct for confounders in permutation testing, to select features correlated with phase III clinical
692 trial outcomes while controlling for target class as a confounding explanatory factor. In the
693 modified permutation test, success/failure labels were shuffled only within target classes, so the
694 sets of null examples had the same ratios of successes to failures within target classes as in the

695 set of observed examples. Consequently, features had to correlate with clinical outcome within
696 multiple classes to be significant, while features that discriminated between classes would not be
697 significant. We performed bootstrapping as described previously to obtain replication
698 probabilities for the significant features, in this case conditioned upon including target class as an
699 explanatory factor. We accepted features with replication probabilities greater than 0.8 as
700 generalizable across target classes represented in the sample.

701

702 **Clinical Outcome Classifier**

703

704 We trained a classifier to predict target success or failure in phase III clinical trials, using a
705 procedure like the above for initial feature selection, then using cross-validation to select a model
706 type (Random Forest or logistic regression) and subset of features useful for prediction. We used
707 an outer cross-validation loop with 5-folds repeated 200 times, yielding a total of 1000 train-test
708 cycles, to estimate the generalization performance and stability of the feature selection and
709 model selection procedure (58). Each train-test cycle had five steps: 1) splitting examples into
710 training and testing sets, 2) univariate feature selection on the training data, 3) aggregation of
711 significant features from different datasets into a single feature matrix, 4) model selection and
712 model-based (multivariate) feature selection on the training data, and 5) evaluation of the
713 classifier on the test data.

714

715 *Step 2: Univariate feature selection*

716

717 Beginning with the non-redundant features obtained from dimensionality reduction, we
718 performed modified permutation tests to find features with a significant difference between
719 successful and failed targets in the training examples. As described above, for the modified
720 permutation test, success/failure labels were shuffled only within target classes. This was done to
721 control for target class as a confounding factor that might explain correlations between phase III
722 outcomes and features. For each feature, we performed 10^4 success/failure label permutations
723 and calculated an empirical two-tailed p-value. We corrected for multiple hypothesis testing
724 within each dataset and accepted features with corrected p-values less than 0.05.

725

726 *Step 3. Feature aggregation*

727

728 Significant features from different datasets, each having different target coverage, had to be
729 aggregated into a single feature matrix prior to training a classifier. When features from many
730 datasets were aggregated, we found that the set of targets with no missing data across all features
731 could become very small. To mitigate this, we excluded features from non-human datasets and
732 small datasets (fewer than 2,000 genes). We also excluded features from the Allen Brain Atlas
733 human brain expression atlas, unless there were no other significant features, because we noticed
734 it had poor coverage of targets with phase III outcomes (287) compared to other expression
735 atlases, such as BioGPS (320), GTEx (328), and HPA (314), which almost always yielded
736 alternative significant expression-based features. After aggregating features into a single matrix,
737 we min-max scaled the features so that features from different datasets would have the same
738 range of values (from 0 to 1).

739

740 To reduce redundancy in the aggregated feature matrix, we grouped features as described for the
741 primary analysis. We used the cosine coefficient to compute pair-wise feature correlations
742 because some features were sparse. Instead of replacing groups of correlated features with the
743 group mean, we selected the feature in each group that was best correlated with phase III
744 outcomes, because we preferred not to create features derived from multiple datasets.

745

746 *Step 4. Model selection and model-based feature selection*

747

748 We hypothesized that a Random Forest classifier (73) would be a reasonable model choice
749 because the Random Forest model does not make any assumptions about the distributions of the
750 features and can seamlessly handle a mixture of quantitative, categorical, filled-in, or sparse
751 features. Furthermore, we expected each train-test cycle to yield only a handful of significant
752 features. Consequently, we would have 10- to 100-fold more training examples than features and
753 could potentially afford to explore non-linear feature combinations. We also trained logistic
754 regression classifiers and used an inner cross-validation loop (described below) to choose
755 between Random Forest and logistic regression for each train-test cycle of the outer cross-
756 validation loop. We used the implementations of the Random Forest and logistic regression
757 classifiers available in the Scikit-learn machine learning package for Python. To correct for
758 unequal class sizes during training, the loss functions of these models weighted the training
759 examples inversely proportional to the size of each example's class.

760

761 We performed incremental feature elimination with an inner cross-validation loop to 1) choose
762 the type of classifier (Random Forest or logistic regression) and 2) choose the smallest subset of

763 features needed to maximize the performance of the classifier. First, we trained Random Forest
764 and logistic regression models using the significant features aggregated in Step 2, performing 5-
765 fold cross-validation repeated 20 times to obtain averages for the area under the receiver
766 operating characteristic curve (AUROC) and area under the precision recall curve (AUPR). We
767 also obtained average feature importance scores from the Random Forest model. Next, we
768 eliminated the feature with lowest importance score and trained the models using the reduced
769 feature set, performing another round of 5-fold cross-validation repeated 20 times to obtain
770 AUROC, AUPR, and feature importance scores. We continued eliminating features then
771 obtaining cross-validation performance statistics and feature importance scores until no features
772 remained. Then, we found all models with performance within 95% of the maximum AUROC
773 and AUPR. If any logistic regression models satisfied this criterion, we selected the qualifying
774 logistic regression model with fewest features. Otherwise, we selected the qualifying Random
775 Forest model with fewest features.

776

777 *Step 5. Classifier evaluation*

778

779 For each train-test cycle, after selecting a set of features and type of model (Random Forest or
780 logistic regression) in Step 4, we re-fit the selected model to the training data and predicted
781 success probabilities for targets in the test set as well as unlabeled targets. For each round of 5-
782 fold cross-validation, we computed the classifier's receiver operating characteristic curve,
783 precision-recall curve, and performance summary statistics, including the true positive rate, false
784 positive rate, positive predictive value, negative predictive value, and Matthews correlation
785 coefficient.

786

787 We computed distributions of the log odds ratios predicted by the classifier (log of the ratio of
788 the predicted probability of success over the probability of failure) for the successful, failed, and
789 untested (unlabeled) targets, aggregating predicted probabilities from the 200 repetitions of 5-
790 fold cross-validation. Histograms of the log odds ratios for the three groups of targets were
791 roughly bell-shaped, so we fit the distributions using kernel density estimation (59) with a
792 Gaussian kernel and applied Silverman's rule for the bandwidth. We transformed the fitted
793 distributions from a function of log odds ratio to a function of probability of success using the
794 rule $\text{pdf}(x) = \text{pdf}(y) * |dy/dx|$.

795

796 We created a heatmap of the probability of success predicted by the classifier projected onto the
797 two dominant features in the model: mean mRNA expression across human tissues and standard
798 deviation of mRNA expression across human tissues. We examined the heatmap to interpret the
799 classifier's decision function and assess its plausibility.

800

801 To more concretely assess the usefulness of the classifier, we found the probability cut-off
802 corresponding to the maximum median positive predictive value and determined the number of
803 unlabeled targets predicted to succeed at that cut-off. Likewise, we found the probability cut-off
804 corresponding to the maximum median negative predictive value and determined the number of
805 unlabeled targets predicted to fail at that cut-off. We also created a heatmap illustrating the
806 separation needed between the median predicted success probabilities of two targets in order to
807 be confident that one target is more likely to succeed than the other. This heatmap was created by

808 calculating the fraction of times Target B had greater probability of success than Target A across
809 the 200 repetitions of 5-fold cross-validation, for all pairs of targets.

810

811 **Implementation**

812

813 Computational analyses were written in Python 3.4.5 and have the following package
814 dependencies: Fastcluster 1.1.20, Matplotlib 1.5.1, Numpy 1.11.3, Requests 2.13.0, Scikit-learn
815 0.18.1, Scipy 0.18.1, and Statsmodels 0.6.1. Code, documentation, and data have been deposited
816 on GitHub at <https://github.com/arouillard/omic-features-successful-targets>.

817

818

819 **ACKNOWLEDGMENTS**

820

821 Many thanks to Dr. Subhas Chakravorty for assisting with access to and processing of the
822 Pharmaprojects data and to Dr. David Cooper for helpful direction regarding nested cross-
823 validation.

824

825

826 **REFERENCES**

827

- 828 1. Arrowsmith J, Miller P. Trial watch: phase II and phase III attrition rates 2011-2012. *Nat Rev Drug Discov.* 2013;12(8):569.
- 829 2. Harrison RK. Phase II and phase III failures: 2013-2015. *Nat Rev Drug Discov.* 2016;15(12):817-8.
- 830 3. Cook D, Brown D, Alexander R, March R, Morgan P, Satterthwaite G, et al. Lessons learned from the fate of AstraZeneca's drug pipeline: a five-dimensional framework. *Nat Rev Drug Discov.* 2014;13(6):419-31.
- 831 4. Gashaw I, Ellinghaus P, Sommer A, Asadullah K. What makes a good drug target? *Drug Discov Today.* 2011;16(23-24):1037-43.
- 832 5. Bunnage ME, Gilbert AM, Jones LH, Hett EC. Know your target, know your molecule. *Nat Chem Biol.* 2015;11(6):368-72.
- 833 6. Rouillard AD, Wang Z, Ma'ayan A. Abstraction for data integration: Fusing mammalian molecular, cellular and phenotype big datasets for better knowledge extraction. *Comput Biol Chem.* 2015;59 Pt B:123-38.
- 834 7. Rigden DJ, Fernandez-Suarez XM, Galperin MY. The 2016 database issue of Nucleic Acids Research and an updated molecular biology database collection. *Nucleic Acids Res.* 2016;44(D1):D1-6.
- 835 8. Abi Hussein H, Geneix C, Petitjean M, Borrel A, Flatters D, Camproux AC. Global vision of druggability issues: applications and perspectives. *Drug Discov Today.* 2017;22(2):404-15.

- 848 9. Fauman EB, Rai BK, Huang ES. Structure-based druggability assessment--identifying
849 suitable targets for small molecule therapeutics. *Curr Opin Chem Biol.* 2011;15(4):463-8.
- 850 10. Perez-Lopez AR, Szalay KZ, Turei D, Modos D, Lenti K, Korcsmaros T, et al. Targets of
851 drugs are generally, and targets of drugs having side effects are specifically good spreaders of
852 human interactome perturbations. *Sci Rep.* 2015;5:10182.
- 853 11. Iwata H, Mizutani S, Tabei Y, Kotera M, Goto S, Yamanishi Y. Inferring protein
854 domains associated with drug side effects based on drug-target interaction network. *BMC Syst
855 Biol.* 2013;7(Suppl 6):S18.
- 856 12. Wang X, Thijssen B, Yu H. Target essentiality and centrality characterize drug side
857 effects. *PLoS Comput Biol.* 2013;9(7):e1003119.
- 858 13. Kotlyar M, Fortney K, Jurisica I. Network-based characterization of drug-regulated
859 genes, drug targets, and toxicity. *Methods.* 2012;57(4):499-507.
- 860 14. Kandoi G, Acencio ML, Lemke N. Prediction of Druggable Proteins Using Machine
861 Learning and Systems Biology: A Mini-Review. *Front Physiol.* 2015;6:366.
- 862 15. Costa PR, Acencio ML, Lemke N. A machine learning approach for genome-wide
863 prediction of morbid and druggable human genes based on systems-level data. *BMC Genomics.*
864 2010;11(Suppl 5):S9.
- 865 16. Bakheet TM, Doig AJ. Properties and identification of human protein drug targets.
866 *Bioinformatics.* 2009;25(4):451-7.
- 867 17. Li Q, Lai L. Prediction of potential drug targets based on simple sequence properties.
868 *BMC Bioinformatics.* 2007;8:353.

- 869 18. Li ZC, Zhong WQ, Liu ZQ, Huang MH, Xie Y, Dai Z, et al. Large-scale identification of
870 potential drug targets based on the topological features of human protein-protein interaction
871 network. *Anal Chim Acta*. 2015;871:18-27.
- 872 19. Jeon J, Nim S, Teyra J, Datti A, Wrana JL, Sidhu SS, et al. A systematic approach to
873 identify novel cancer drug targets using machine learning, inhibitor design and high-throughput
874 screening. *Genome Med*. 2014;6(7):57.
- 875 20. Zhu M, Gao L, Li X, Liu Z, Xu C, Yan Y, et al. The analysis of the drug-targets based on
876 the topological properties in the human protein-protein interaction network. *J Drug Target*.
877 2009;17(7):524-32.
- 878 21. Yao L, Rzhetsky A. Quantitative systems-level determinants of human genes targeted by
879 successful drugs. *Genome Res*. 2008;18(2):206-13.
- 880 22. Mora A, Donaldson IM. Effects of protein interaction data integration, representation and
881 reliability on the use of network properties for drug target prediction. *BMC Bioinformatics*.
882 2012;12(13):294.
- 883 23. Mitsopoulos C, Schierz AC, Workman P, Al-Lazikani B. Distinctive Behaviors of
884 Druggable Proteins in Cellular Networks. *PLoS Comput Biol*. 2015;11(12):e1004597.
- 885 24. Xu H, Xu H, Lin M, Wang W, Li Z, Huang J, et al. Learning the drug target-likeness of a
886 protein. *Proteomics*. 2007;7(23):4255-63.
- 887 25. Bull SC, Doig AJ. Properties of protein drug target classes. *PLoS One*.
888 2015;10(3):e0117955.
- 889 26. Li S, Yu X, Zou C, Gong J, Liu X, Li H. Are Topological Properties of Drug Targets
890 Based on Protein-Protein Interaction Network Ready to Predict Potential Drug Targets? *Comb
891 Chem High Throughput Screen*. 2016;19(2):109-20.

- 892 27. Ghiassian SD, Menche J, Barabasi AL. A DIseASE MOdule Detection (DIAMOnD)
893 algorithm derived from a systematic analysis of connectivity patterns of disease proteins in the
894 human interactome. *PLoS Comput Biol.* 2015;11(4):e1004120.
- 895 28. Yang P, Li X, Chua HN, Kwoh CK, Ng SK. Ensemble positive unlabeled learning for
896 disease gene identification. *PLoS One.* 2014;9(5):e97079.
- 897 29. Carson MB, Lu H. Network-based prediction and knowledge mining of disease genes.
898 *BMC Med Genomics.* 2015;8(Suppl 2):S9.
- 899 30. Zhu C, Wu C, Aronow BJ, Jegga AG. Computational approaches for human disease gene
900 prediction and ranking. *Adv Exp Med Biol.* 2014;799:69-84.
- 901 31. Piro RM, Di Cunto F. Computational approaches to disease-gene prediction: rationale,
902 classification and successes. *FEBS J.* 2012;279(5):678-96.
- 903 32. Moreau Y, Tranchevent LC. Computational tools for prioritizing candidate genes:
904 boosting disease gene discovery. *Nat Rev Genet.* 2012;13(8):523-36.
- 905 33. Emig D, Ivliev A, Pustovalova O, Lancashire L, Bureeva S, Nikolsky Y, et al. Drug
906 target prediction and repositioning using an integrated network-based approach. *PLoS One.*
907 2013;8(4):e60618.
- 908 34. Sun J, Zhu K, Zheng W, Xu H. A comparative study of disease genes and drug targets in
909 the human protein interactome. *BMC Bioinformatics.* 2015;16(Suppl 5):S1.
- 910 35. Suthram S, Dudley JT, Chiang AP, Chen R, Hastie TJ, Butte AJ. Network-based
911 elucidation of human disease similarities reveals common functional modules enriched for
912 pluripotent drug targets. *PLoS Comput Biol.* 2010;6(2):e1000662.
- 913 36. Nelson MR, Tipney H, Painter JL, Shen J, Nicoletti P, Shen Y, et al. The support of
914 human genetic evidence for approved drug indications. *Nat Genet.* 2015;47(8):856-60.

- 915 37. Heinemann F, Huber T, Meisel C, Bundschus M, Leser U. Reflection of successful
916 anticancer drug development processes in the literature. *Drug Discov Today*. 2016;21(11):1740-
917 4.
- 918 38. Pharmaprojects [Internet]. 2017. Available from:
919 <https://pharmaintelligence.informa.com/products-and-services/data-and-analysis/pharmaprojects>.
- 920 39. Rouillard AD, Gundersen GW, Fernandez NF, Wang Z, Monteiro CD, McDermott MG,
921 et al. The harmonizome: a collection of processed datasets gathered to serve and mine
922 knowledge about genes and proteins. *Database (Oxford)*. 2016;2016.
- 923 40. Ernst MD. Permutation Methods: A Basis for Exact Inference. *Statistical Science*.
924 2004;19(4):676-85.
- 925 41. Phipson B, Smyth GK. Permutation P-values should never be zero: calculating exact P-
926 values when permutations are randomly drawn. *Stat Appl Genet Mol Biol*. 2010;9:Article39.
- 927 42. Benjamini Y, Yekutieli D. The control of the false discovery rate in multiple testing
928 under dependency. *Annals of Statistics*. 2001;29(4):1165-88.
- 929 43. Sunkin SM, Ng L, Lau C, Dolbeare T, Gilbert TL, Thompson CL, et al. Allen Brain
930 Atlas: an integrated spatio-temporal portal for exploring the central nervous system. *Nucleic
931 acids research*. 2013;41(Database issue):D996-D1008.
- 932 44. Hawrylycz MJ, Lein ES, Gulyas B, Bongianni AL, Shen EH, Ng L, Miller JA, et al. An
933 anatomically comprehensive atlas of the adult human brain transcriptome. *Nature*.
934 2012;489(7416):391-9.
- 935 45. Lein ES, Hawrylycz MJ, Ao N, Ayres M, Bensinger A, Bernard A, et al. Genome-wide
936 atlas of gene expression in the adult mouse brain. *Nature*. 2007;445(7124):168-76.

- 937 46. Wu C, MacLeod I, Su AI. BioGPS and MyGene. info: organizing online, gene-centric
938 information. *Nucleic acids research*. 2012;gks1114.
- 939 47. Su AI, Cooke MP, Ching KA, Hakak Y, Walker JR, Wiltshire T, et al. Large-scale
940 analysis of the human and mouse transcriptomes. *Proceedings of the National Academy of
941 Sciences of the United States of America*. 2002;99(7):4465-70.
- 942 48. Su AI, Wiltshire T, Batalov S, Lapp H, Ching KA, Block D, et al. A gene atlas of the
943 mouse and human protein-encoding transcriptomes. *Proceedings of the National Academy of
944 Sciences of the United States of America*. 2004;101(16):6062-7.
- 945 49. Consortium G. The Genotype-Tissue Expression (GTEx) project. *Nature genetics*.
946 2013;45(6):580-5.
- 947 50. Consortium G. Human genomics. The Genotype-Tissue Expression (GTEx) pilot
948 analysis: multitissue gene regulation in humans. *Science*. 2015;348(6235):648-60.
- 949 51. Uhlen M, Fagerberg L, Hallstrom BM, Lindskog C, Oksvold P, Mardinoglu A, et al.
950 Proteomics. *Tissue-based map of the human proteome*. *Science*. 2015;347(6220):1260419.
- 951 52. Santos A, Tsafou K, Stolte C, Pletscher-Frankild S, O'Donoghue SI, Jensen LJ.
952 Comprehensive comparison of large-scale tissue expression datasets. *PeerJ*. 2015;3:e1054.
- 953 53. Efron B, Tibshirani R. *An introduction to the bootstrap*. New York: Chapman and Hall;
954 1991.
- 955 54. Calmettes G, Drummond GB, Vowler SL. Making do with what we have: use your
956 bootstraps. *J Physiol*. 2012;590(15):3403-6.
- 957 55. Jaffe AE, Storey JD, Ji H, Leek JT. Gene set bagging for estimating the probability a
958 statistically significant result will replicate. *BMC Bioinformatics*. 2013;14:360.

- 959 56. Gray KA, Yates B, Seal RL, Wright MW, Bruford EA. Genenames.org: the HGNC
960 resources in 2015. *Nucleic Acids Res.* 2015;43(Database issue):D1079-85.
- 961 57. Epstein MP, Duncan R, Jiang Y, Conneely KN, Allen AS, Satten GA. A permutation
962 procedure to correct for confounders in case-control studies, including tests of rare variation. *Am
963 J Hum Genet.* 2012;91(2):215-23.
- 964 58. Varma S, Simon R. Bias in error estimation when using cross-validation for model
965 selection. *BMC Bioinformatics.* 2006;7:91.
- 966 59. Scott DW. *Multivariate Density Estimation: Theory, Practice, and Visualization*: John
967 Wiley and Sons, Inc.; 1992.
- 968 60. Kumar V, Sanseau P, Simola DF, Hurle MR, Agarwal P. Systematic Analysis of Drug
969 Targets Confirms Expression in Disease-Relevant Tissues. *Sci Rep.* 2016;6:36205.
- 970 61. Lage K, Hansen NT, Karlberg EO, Eklund AC, Roque FS, Donahoe PK, et al. A large-
971 scale analysis of tissue-specific pathology and gene expression of human disease genes and
972 complexes. *Proc Natl Acad Sci U S A.* 2008;105(52):20870-5.
- 973 62. Magger O, Waldman YY, Ruppin E, Sharan R. Enhancing the prioritization of disease-
974 causing genes through tissue specific protein interaction networks. *PLoS Comput Biol.*
975 2012;8(9):e1002690.
- 976 63. Grimes DA, Schulz KF. Bias and causal associations in observational research. *The
977 Lancet.* 2002;359(9302):248-52.
- 978 64. Tripepi G, Jager KJ, Dekker FW, Zoccali C. Selection bias and information bias in
979 clinical research. *Nephron Clin Pract.* 2010;115(2):c94-9.
- 980 65. Groth D, Hartmann S, Klie S, Selbig J. Principal components analysis. *Methods Mol
981 Biol.* 2013;930:527-47.

- 982 66. Hinton GE, Salakhutdinov RR. Reducing the dimensionality of data with neural
983 networks. *Science*. 2006;313(5786):504-7.
- 984 67. Hurle MR, Nelson MR, Agarwal P, Cardon LR. Trial watch: Impact of genetically
985 supported target selection on R&D productivity. *Nature reviews Drug discovery*.
986 2016;15(9):596-7.
- 987 68. Brown GR, Hem V, Katz KS, Ovetsky M, Wallin C, Ermolaeva O, et al. Gene: a gene-
988 centered information resource at NCBI. *Nucleic Acids Res*. 2015;43(Database issue):D36-42.
- 989 69. Coletti MH, Bleich HL. Medical subject headings used to search the biomedical
990 literature. *J Am Med Inform Assoc*. 2001;8(4):317-23.
- 991 70. Spearman C. The Proof and Measurement of Association between Two Things.
992 *American Journal of Psychology*. 1904;15(1):72-101.
- 993 71. Frades I, Matthiesen R. Overview on techniques in cluster analysis. *Methods Mol Biol*.
994 2010;593:81-107.
- 995 72. Deshpande R, Vandersluis B, Myers CL. Comparison of profile similarity measures for
996 genetic interaction networks. *PLoS One*. 2013;8(7):e68664.
- 997 73. Breiman L. Random Forests. *Machine Learning*. 2001;45:5-32.
- 998
- 999

1000 **SUPPORTING INFORMATION**

1001

1002 **S1. Supplementary Table S1.** List of targets with their phase III outcome labels and predicted
1003 success probabilities for 200 cross-validation repetitions.

1004

1005 **S2. Supplementary Table S2.** List of non-redundant features with their similar features and p-
1006 values from the basic permutation test.

1007

1008 **S3. Supplementary Table S3.** List of classifier attributes (selected features, selected model type,
1009 and test performance) for 1000 train-test cycles.

1010

1011 **S4. Supplementary Table S4.** Comparison of inner cross-validation loop AUROC and AUPR
1012 values between Random Forest and logistic regression models for 1000 train-test cycles.

1013

1014 **S5. Supplementary Table S5.** List of classifier test performance statistics for 200 cross-
1015 validation repetitions.

1016

1017 **S6. Supplementary Table S6.** Cases illustrating how the significance of genetic evidence (and
1018 likely other types of evidence) as a predictor of target success depends on which targets are
1019 compared.

1020





