

RESEARCH ARTICLE

Predicting bioprocess targets of chemical compounds through integration of chemical-genetic and genetic interactions

Scott W. Simpkins¹, Justin Nelson¹, Raamesh Deshpande², Sheena C. Li³, Jeff S. Piotrowski^{3*}, Erin H. Wilson², Abraham A. Gebre⁴, Hamid Safizadeh^{2,5}, Reika Okamoto³, Mami Yoshimura³, Michael Costanzo⁶, Yoko Yashiroda³, Yoshikazu Ohya⁴, Hiroyuki Osada³, Minoru Yoshida³, Charles Boone^{3,6*}, Chad L. Myers^{1,2*}

1 University of Minnesota-Twin Cities, Bioinformatics and Computational Biology Graduate Program, Minneapolis, Minnesota, United States of America, **2** University of Minnesota-Twin Cities, Department of Computer Science and Engineering, Minneapolis, Minnesota, United States of America, **3** RIKEN Center for Sustainable Resource Science, Wako, Saitama, Japan, **4** University of Tokyo, Department of Integrated Biosciences, Graduate School of Frontier Sciences, Kashiwa, Chiba, Japan, **5** University of Minnesota, Department of Electrical and Computer Engineering, Minneapolis, Minnesota, United States of America, **6** University of Toronto, Donnelly Centre, Toronto, Ontario, Canada

* Current Address: Yumanity Therapeutics, Cambridge, MA, United States of America
* charlie.boone@utoronto.ca (CB); chadm@umn.edu (CLM)



OPEN ACCESS

Citation: Simpkins SW, Nelson J, Deshpande R, Li SC, Piotrowski JS, Wilson EH, et al. (2018) Predicting bioprocess targets of chemical compounds through integration of chemical-genetic and genetic interactions. *PLoS Comput Biol* 14(10): e1006532. <https://doi.org/10.1371/journal.pcbi.1006532>

Editor: Jean-Philippe Mr. Vert, Google Inc, FRANCE

Received: March 22, 2018

Accepted: September 26, 2018

Published: October 30, 2018

Copyright: © 2018 Simpkins et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: The complete set of inputs and results for the primary analysis performed in this manuscript are available from the Dryad Digital Repository at: [doi:10.5061/dryad.nr2cf12](https://doi.org/10.5061/dryad.nr2cf12).

Funding: This work was partially supported by the National Institutes of Health (<https://www.nih.gov/>) (R01HG005084, R01GM104975) and the National Science Foundation (<https://www.nsf.gov/>) (DBI 0953881). SWS is supported by an NSF Graduate Research Fellowship (00039202), an NIH

Abstract

Chemical-genetic interactions—observed when the treatment of mutant cells with chemical compounds reveals unexpected phenotypes—contain rich functional information linking compounds to their cellular modes of action. To systematically identify these interactions, an array of mutants is challenged with a compound and monitored for fitness defects, generating a chemical-genetic interaction profile that provides a quantitative, unbiased description of the cellular function(s) perturbed by the compound. Genetic interactions, obtained from genome-wide double-mutant screens, provide a key for interpreting the functional information contained in chemical-genetic interaction profiles. Despite the utility of this approach, integrative analyses of genetic and chemical-genetic interaction networks have not been systematically evaluated. We developed a method, called CG-TARGET (Chemical Genetic Translation via A Reference Genetic nETwork), that integrates large-scale chemical-genetic interaction screening data with a genetic interaction network to predict the biological processes perturbed by compounds. In a recent publication, we applied CG-TARGET to a screen of nearly 14,000 chemical compounds in *Saccharomyces cerevisiae*, integrating this dataset with the global *S. cerevisiae* genetic interaction network to prioritize over 1500 compounds with high-confidence biological process predictions for further study. We present here a formal description and rigorous benchmarking of the CG-TARGET method, showing that, compared to alternative enrichment-based approaches, it achieves similar or better accuracy while substantially improving the ability to control the false discovery rate of biological process predictions. Additional investigation of the compatibility of chemical-genetic and genetic interaction profiles revealed that one-third of observed chemical-genetic interactions contributed to the highest-confidence biological process predictions and that negative chemical-genetic interactions overwhelmingly formed the basis of these predictions. We

Biotechnology training grant (T32GM008347), and a one-year fellowship from the University of Minnesota Bioinformatics and Computational Biology Graduate Program (<https://r.umn.edu/academics-research/graduate-programs/bicb>). SCL and JSP are supported by a RIKEN (<http://www.riken.jp/en/>) Foreign Postdoctoral Research Fellowship. SCL is supported by a RIKEN CSRS (<http://www.csrs.riken.jp/en/>) Research Topics for Cooperative Projects Award (201601100228), and a RIKEN FY2017 Incentive Research Projects Grant. YO is supported through Grants-in-Aid for Scientific Research (15H04402) from the Ministry of Education, Culture, Sports, Science and Technology, Japan (www.mext.go.jp/en/). CB and YO are supported by JSPS KAKENHI grant number 15H04483 (<http://www.jsps.go.jp/english/>). CB and YY are supported by a JSPS Grant-in-Aid for Scientific Research on Innovative Areas (17H06411). CLM and CB are fellows in the Canadian Institute for Advanced Research (CIFAR, <https://www.cifar.ca/>) Genetic Networks Program. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: I have read the journal's policy and the authors of this manuscript have the following potential competing interests: Authors Simpkins, Nelson, and Myers have licensed the CG-TARGET software for commercial use (z.umn.edu/cgtarget).

also present experimental validations of CG-TARGET-predicted tubulin polymerization and cell cycle progression inhibitors. Our approach successfully demonstrates the use of genetic interaction networks in the high-throughput functional annotation of compounds to biological processes.

Author summary

Understanding how chemical compounds affect biological systems is of paramount importance as pharmaceutical companies strive to develop life-saving medicines, governments seek to regulate the safety of consumer products and agrichemicals, and basic scientists continue to study the fundamental inner workings of biological organisms. One powerful approach to characterize the effects of chemical compounds in living cells is chemical-genetic interaction screening. Using this approach, a collection of cells—each with a different defined genetic perturbation—is tested for sensitivity or resistance to the presence of a compound, resulting in a quantitative profile describing the functional effects of that compound on the cells. The work presented here describes our efforts to integrate compounds' chemical-genetic interaction profiles with reference genetic interaction profiles containing information on gene function to predict the cellular processes perturbed by the compounds. We focused on specifically developing a method that could scale to perform these functional predictions for large collections of thousands of screened compounds and robustly control the false discovery rate. With chemical-genetic and genetic interaction screens now underway in multiple species including human cells, the method described here can be generally applied to enable the characterization of compounds' effects across the tree of life.

Introduction

The discovery of chemical compounds with desirable and interesting biological activity advances our understanding of how compounds and biological systems interact. Chemical-genetic interaction profiling enables this discovery by measuring the response of defined gene mutants to chemical compounds [1–8]. Specifically, a chemical-genetic interaction profile refers to the set of gene mutations that confer sensitivity (a negative chemical-genetic interaction) or resistance (a positive interaction) to a compound and provides functional insights into the compound's mode(s) of action. Recent advances in DNA sequencing technology have enabled dramatic increases in the throughput of chemical-genetic interaction screens (into the range of thousands of compounds) via multiplexed analysis of pooled mutant libraries [6,7,9].

Similarly, genetic interactions identify pairs of gene mutations whose combined phenotypes are more or less severe than expected given the phenotypes of the individual mutants. In *S. cerevisiae*, the vast majority of all possible gene double-mutant pairs have been constructed and scored for fitness-based genetic interactions, yielding a global compendium of genome-wide genetic interaction profiles that quantitatively describe each gene's function. Similarity between two genes' genetic interaction profiles implies that these genes perform similar functions, enabling the functional annotation of uncharacterized genes and the construction of a global hierarchy of cellular function [5,10].

The global genetic interaction network in *S. cerevisiae* provides a resource for interpreting chemical-genetic interaction profiles across a broad range of cellular function, as the chemical-

genetic interaction profile of a compound should resemble the genetic interaction profile of its cellular target or target processes [2,5]. Importantly, this approach to interpretation does not depend on reference chemical-genetic interaction profiles and thus enables the discovery of compounds with novel modes of action. Previous small and large-scale chemical-genetic interaction studies have employed various computational methods to provide more informative clustering of the resulting interaction matrices [3,11] and even predict perturbed protein complexes [12] or direct protein targets [13]. However, the integration of chemical-genetic and genetic interaction profiles has only been performed in the context of relatively small studies [2,5].

Here, we present the use of genetic interaction profiles to systematically interpret chemical-genetic interaction profiles on a large scale. To this end, we developed a computational method, called CG-TARGET (Chemical Genetic Translation via A Reference Genetic nET-work), that integrates chemical-genetic and genetic interaction profiles to predict the biological processes perturbed by compounds. In a recent publication [14], we applied this method to a chemical-genetic interaction screen of nearly 14,000 compounds in *S. cerevisiae* [14], using profiles from the global yeast genetic interaction network [5,10] to interpret the chemical-genetic interaction profiles. Here, we show that CG-TARGET recapitulates known information for well-characterized compounds and showed a marked improvement in false discovery rate control compared to alternative, enrichment-based approaches. Additionally, we experimentally validated two different mode-of-action predictions, one in an *in vitro* system using mammalian proteins, confirming both the accuracy of the predictions and the potential to translate them across species. CG-TARGET is available, free for non-commercial use, at <https://github.com/csbio/CG-TARGET>.

Results

Overview of datasets used in this study

We obtained chemical-genetic interaction profiles from a recent large-scale chemical-genetic interaction screen in *S. cerevisiae* [14]. Profiles were obtained in two batches, labeled “RIKEN” and “NCI/NIH/GSK” to reflect the compound libraries screened—for RIKEN, the RIKEN Natural Product Depository [15], and for NCI/NIH/GSK, plated libraries from the NCI Open Chemical Repository, the NIH Clinical Collection, and the GlaxoSmithKline Published Kinase Inhibitor Set [16]. The RIKEN compounds were primarily natural products and derivatives—mostly uncharacterized—but also contained ~200 approved drugs and chemical probes from which we selected a well-characterized subset for benchmarking. The NCI/NIH/GSK compounds were more characterized, having been tested against the NCI-60 cancer cell line panel (NCI collections), tested in clinical trials (NIH Clinical Collection) or designed to inhibit human kinases (GSK)—but their specific modes of action remained primarily uncharacterized. The final datasets consisted of interaction scores for 8418 RIKEN compounds and 3565 NCI/NIH/GSK compounds (with 5724 and 2128 negative control conditions, respectively) screened against a diagnostic set of ~300 haploid gene deletion mutants selected to optimally capture the information in the complete *S. cerevisiae* non-essential deletion collection [14,17]. Each profile contained z-scores that reflected the deviation of each strain’s observed abundance from expected abundance in the presence of a compound.

Genetic interaction profiles were obtained from a recently assembled, genome-wide compendium of genetic interaction profiles in *S. cerevisiae* [5]. These profiles were generated through the systematic analysis of double mutant fitness and consist of epsilon scores that reflect the deviation of each double mutant’s observed fitness from that expected given the single mutant fitness values, assuming a multiplicative null model [18]. Profiles were filtered to

the ~35% with the highest signal, and we mapped these 1505 high-signal “query” genes to Gene Ontology biological process terms [19,20] to define the bioprocess targets of compounds. (see [Materials and Methods](#)).

Predicting perturbed bioprocesses from chemical-genetic interaction profiles

We developed CG-TARGET (Chemical Genetic Translation via A Reference Genetic nET-work) to predict the biological processes perturbed by compounds in our recently-generated dataset of ~12,000 chemical-genetic interaction profiles (Fig 1). CG-TARGET requires three input datasets: 1) chemical-genetic interaction profiles; 2) genetic interaction profiles; and 3) a mapping from the query genes in the genetic interaction profiles to gene sets representing coherent biological processes (referred to as “bioprocesses”). Predicting the bioprocesses perturbed by a particular compound involves four distinct steps. First, a control set of resampled chemical-genetic interaction profiles is generated, each of which consists of one randomly-sampled interaction score per gene mutant across all compound treatment profiles in the chemical-genetic interaction dataset; these profiles thus provide a means to account for variance in each mutant strain observed upon treatment with bioactive compound but not upon treatment with experimental controls (DMSO with no active compound). Second, “gene-target” prediction scores between each compound and query gene are generated by computing an inner product between all chemical-genetic interaction profiles (comprising compound treatment, experimental control, and random profiles) and all L_2 -normalized query genetic interaction profiles; normalizing only the genetic interaction profiles results in gene-target scores that should be more robust to noise in the chemical-genetic data [21] and reflect the overall strength of each chemical-genetic profile as well as its similarity to gene mutants’ profiles. Third, these “gene-target” prediction scores are aggregated into bioprocess predictions; a z-score and empirical p-value for each compound-bioprocess prediction are obtained by mapping the gene-target prediction scores to the genes in the bioprocess of interest and comparing these scores to those from shuffled gene-target prediction scores and to distributions of the scores derived from experimental control and resampled profiles. Finally, the false discovery rates for these predictions are estimated by comparing, across a range of significance

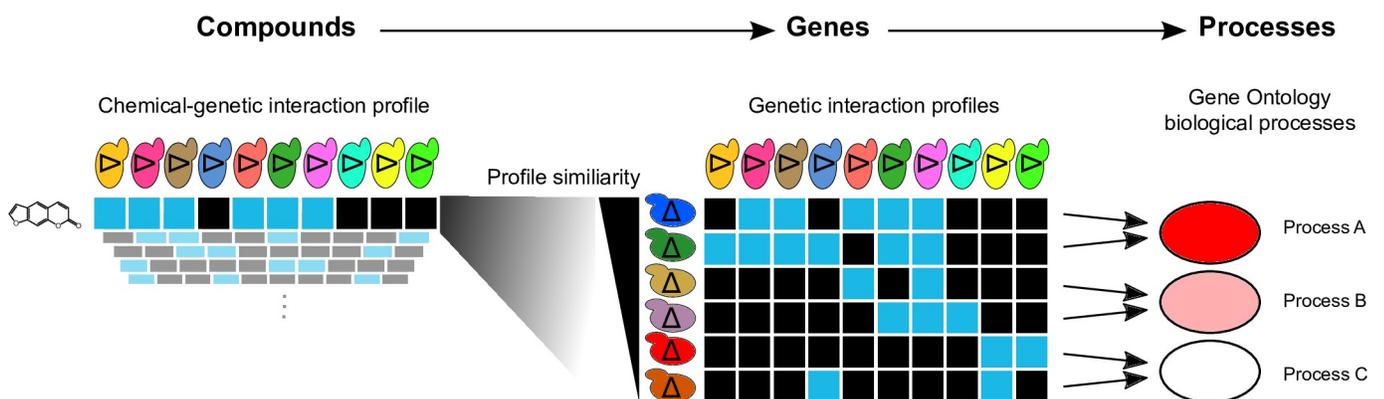


Fig 1. Overview of the integration of chemical-genetic and genetic interaction networks for bioprocess target prediction using CG-TARGET. Chemical-genetic interaction profiles, obtained by measuring the sensitivity or resistance of a library of gene mutants to a chemical compound, are compared against genetic interaction profiles consisting of double mutant interaction scores. The resulting similarities are aggregated at the level of biological processes to predict the bioprocess(es) perturbed by the compound. Better agreement between chemical-genetic and genetic interaction profiles leads to stronger bioprocess predictions. Each blue box represents a negative chemical-genetic (i.e. sensitivity) or genetic interaction, while each black box represents the absence of an interaction. Stronger bioprocess predictions are depicted with a darker red.

<https://doi.org/10.1371/journal.pcbi.1006532.g001>

thresholds, the frequency at which experimental control and randomly resampled profiles predict bioprocesses versus that of compound treatment profiles (see [Materials and Methods](#)). A schematic representation of the method is provided as [S1 Fig](#).

Application to and evaluation on large-scale chemical-genetic interaction data

To provide a baseline for benchmarking the performance of CG-TARGET on these large screens, we implemented two simple, enrichment-based approaches for predicting bioprocess-level targets. The “direct enrichment” approach tested for enrichment of GO biological processes among each compound’s 20 strongest negative chemical-genetic interactors, providing a comparison to methods that do not incorporate genetic interaction profiles. The “gene-target enrichment” approach tested for the enrichment of GO biological processes among the top-*n* gene-target prediction scores for each compound, enabling a comparison of CG-TARGET’s z-score-based approach to enrichment on the gene-target scores. For the comparisons to gene-target enrichment below, we selected *n* = 20 as it showed the best overall performance across a range of values of *n* ([S2 Fig](#)).

We applied CG-TARGET to the RIKEN and NCI/NIH/GSK chemical-genetic interaction screens, identifying 848 out of 8418 compounds (10%) from the RIKEN screen and 705 of 3565 compounds (20%) from the NCI/NIH/GSK screen with at least one prediction that achieved false discovery rates of 25 and 27%, respectively (referred to as “high-confidence” compounds and predictions) ([Table 1](#), [Fig 2](#)). Measured using the RIKEN dataset, this rate of discovery at $FDR \leq 25\%$ was over 4-fold higher in terms of number of discovered compounds than that of direct enrichment (190 compounds) and over 100-fold higher than that of gene-target enrichment (7 compounds, [Fig 3A](#)). In all cases, the false discovery rates derived from resampled profiles were more conservative than those derived from experimental controls, suggesting that some sources of variance in each gene mutant’s interaction scores arose only upon treatment with compound and therefore could not be corrected using only solvent controls.

In addition to assessing false discovery rate control relative to baseline methods, we also assessed prediction accuracy. We performed the first of these comparisons against the direct enrichment predictions by asking if the top prediction for each of 35 well-characterized compounds matched what was known about that compound. For direct enrichment, the top prediction for 11 of these 35 compounds matched its known mode of action, with only 6 of these compounds passing the $FDR \leq 25\%$ criteria that would enable their discovery in a large-scale screen ([S1 Table](#)). In contrast, CG-TARGET matched 17 of these compounds to their known mode of action, with 16 passing the $FDR \leq 25\%$ discovery threshold.

We then compared CG-TARGET to gene-target enrichment using two measures of accuracy. The first accuracy-based evaluation was performed on genetic interaction profiles with

Table 1. The number of compounds discovered at selected false discovery rates upon application of CG-TARGET to data from two large-scale chemical-genetic interaction screens. The “RIKEN” screen consisted of 8418 total compounds from the RIKEN Natural Product Depository, and the “NCI/NIH/GSK” screen consisted of 3565 compounds across 6 chemical compound collections from the National Cancer Institute, National Institutes of Health, and GlaxoSmithKline.

Dataset	RIKEN		NCI/NIH/GSK	
	p-value	number of compounds	p-value	number of compounds
FDR cutoff				
0.00	$< 2 \times 10^{-5}$	434	$< 2 \times 10^{-5}$	352
0.05	2×10^{-5}	505	4×10^{-5}	405
0.10	8×10^{-5}	598	1.6×10^{-4}	494
0.25*	2.8×10^{-4}	848	4.7×10^{-4}	705

*This cutoff is 0.27 for the NCI/NIH/GSK dataset

<https://doi.org/10.1371/journal.pcbi.1006532.t001>

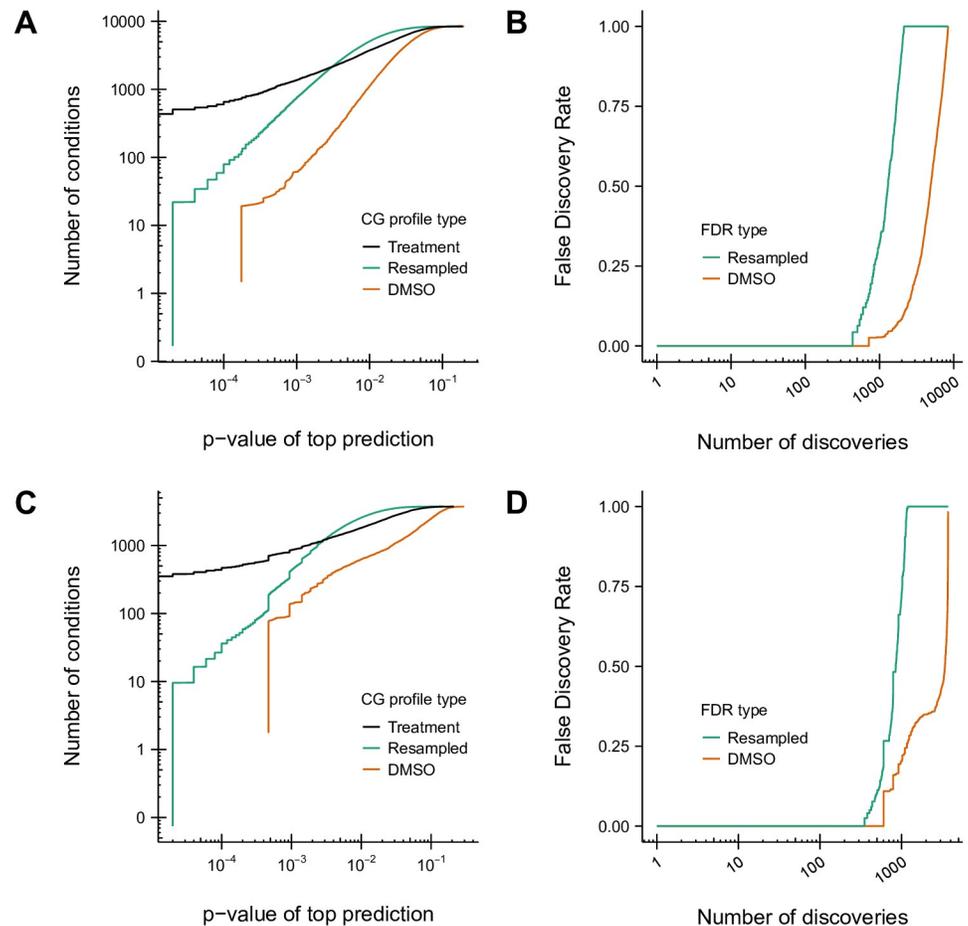


Fig 2. Rate of compound discovery and control of the false discovery rate for the prediction of bioprocesses from chemical-genetic interaction profiles. Perturbed bioprocesses were predicted using CG-TARGET for compounds, experimental controls (DMSO), and resampled chemical-genetic interaction profiles from the RIKEN and NCI/NIH/GSK datasets. (A) The number of compounds, experimental controls, and randomly resampled chemical-genetic interaction profiles discovered with at least one bioprocess prediction passing the given significance thresholds, for the RIKEN dataset. (B) DMSO and resampled profile-derived estimates of the false discovery rate of biological process predictions, for the RIKEN dataset, given the number of discovered compounds. Values were calculated from (A). (C-D) Same as (A-B), respectively, but for the NCI/NIH/GSK dataset.

<https://doi.org/10.1371/journal.pcbi.1006532.g002>

added noise, which provided a means to both simulate chemical-genetic interaction profiles and annotate them with gold-standard GO biological process annotations for evaluation. For the second accuracy-based evaluation, we assigned each of the aforementioned well-characterized compounds to a “gold standard” bioprocess term and evaluated the ranks of each compound’s gold-standard bioprocess within its list of bioprocess predictions. We note that neither of these methods were particularly suitable for comparing CG-TARGET to direct enrichment, as 1) the assumption of alignment between chemical-genetic and genetic interaction profiles was implicit in the generation of the simulated profiles and 2) we anticipated that spurious rank differences would result from differences in the size (~300 genes for direct versus ~1500 genes for CG-TARGET) and composition (about half of the former in the latter) of the two gene universes that defined the bioprocess term sets.

CG-TARGET performed comparably to the best-performing gene-target enrichment method using our measures of accuracy. This is first shown in the evaluation of these methods’ respective abilities to predict a gold-standard annotated bioprocess as the top prediction for

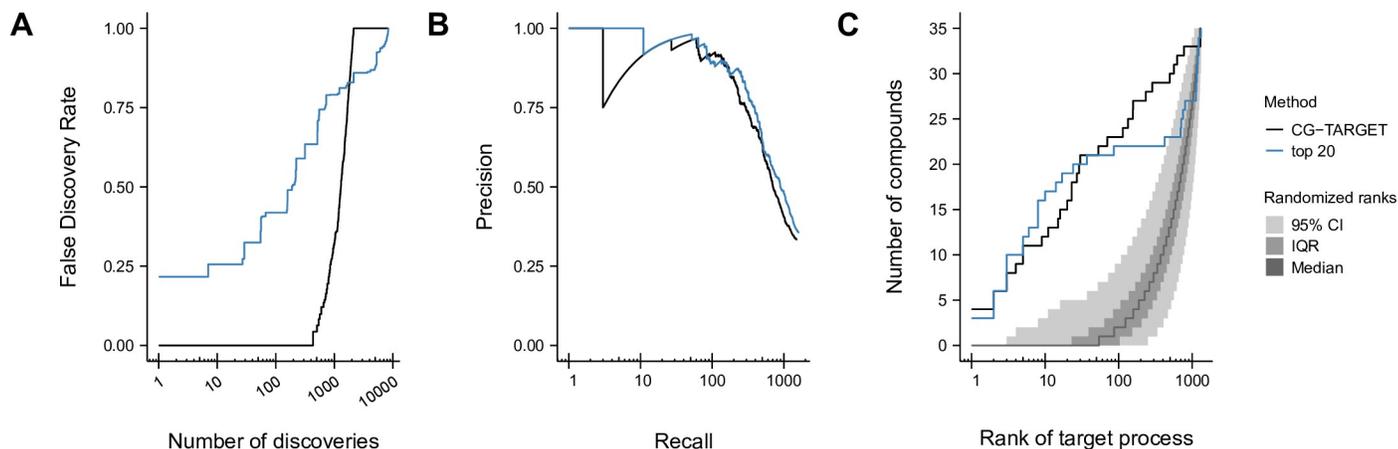


Fig 3. Comparison of CG-TARGET performance versus gene-target enrichment. Perturbed bioprocesses were predicted using both CG-TARGET and a method that calculated enrichment on the set of each compound's 20 most similar genetic interaction profiles ("top 20"). (A) Bioprocess prediction false discovery rate estimates derived from resampled chemical-genetic interaction profiles, performed on compounds from the RIKEN dataset. (B) Precision-recall analysis of the ability to recapitulate gold-standard annotations within the set of top bioprocess predictions for ~4500 simulated compounds. Each simulated compound was designed to target one query gene in the genetic interaction network and thus inherited gold-standard biological process annotations from its target gene. (C) For each of 35 well-characterized compounds in the RIKEN dataset with literature-derived, gold-standard biological process annotations, we determined the rank of its gold-standard bioprocess within its list of predictions. The number of compounds for which a given rank (or better) was achieved is plotted. The grey ribbons represent the median, interquartile range (25th to 75th percentiles), and 95% confidence interval of 10,000 rank permutations.

<https://doi.org/10.1371/journal.pcbi.1006532.g003>

each simulated chemical-genetic interaction profile. Specifically, CG-TARGET performed nearly as well as the top-20 gene-target enrichment method across both low and high recall values (Fig 3B). Both methods captured a gold-standard annotation as the top predicted bioprocess for approximately 34% of the simulated compounds (33.4% and 35.6% for CG-TARGET and top-20 gene-target enrichment, respectively), which represented more than a 22-fold enrichment over the background expectation of 1.5% (the average number of gold-standard bioprocess annotations per simulated compound divided by the number of bioprocesses).

For the 35 gold-standard compound-bioprocess pairs, we observed that both CG-TARGET and gene-target enrichment captured the gold-standard bioprocess for 6 and 21 (out of 35) compounds above ranks of 2 and 40 (out of 1329), respectively, with slightly decreased performance for CG-TARGET between these rank thresholds (Fig 3C, Table 2). The significance of these rank values was evaluated by randomizing the order of each compound's bioprocess predictions 10,000 times and recalculating the ranks. Both methods achieved similar results in this respect, with CG-TARGET and gene-target enrichment respectively identifying 22 and 21 gold-standard compounds with significantly better ranks than the random expectation. The two methods also performed similarly when comparing the "effective rank" of each compound's gold-standard bioprocess, with CG-TARGET and gene-target enrichment respectively identifying 20 and 22 compounds for which the gold-standard or a closely-related bioprocess achieved a rank of 5 or better. Despite the similar performance in rank space, however, none of the 21 significantly-ranked predictions made by gene-target enrichment achieved $FDR \leq 25\%$, compared to 16 out of 22 for CG-TARGET (Table 2).

Characterizing performance with respect to individual bioprocess terms

In addition to benchmarking CG-TARGET's ability to prioritize gold-standard annotated bioprocesses for specific compounds, we also benchmarked its ability to prioritize compounds that perturb specific bioprocesses. Specifically, each GO term was evaluated based on the ranks of the predictions for the simulated chemical-genetic interaction profiles derived from genes

Table 2. Evaluation of predictions made by CG-TARGET, and comparison to a baseline enrichment approach, for literature-derived, gold-standard compound-process annotations. The target bioprocess rank was determined by its position in the list of all bioprocess predictions for each gold-standard compound, with the significance computed empirically by shuffling the bioprocesses and re-computing the rank (bold p-values indicate significance, $p < 0.05$). Asterisks indicate cases in which the false discovery rate of the gold-standard compound-process prediction was less than 25%. The “top-20 enrichment” approach was selected as a baseline for comparison. The “effective rank” of a compound-bioprocess prediction represents the top rank within the compound’s list of predictions among bioprocesses that are similar to the original bioprocess.

Compound	GO ID	GO term	CG-TARGET			top-20 enrichment		
			Target process rank	Rank significance	Effective rank	Target process rank	Rank significance	Effective rank
5-Fluorocytosine	GO:0032774	RNA biosynthetic process	27	0.0208	2	3	0.0027	1
Aclacinomycin A	GO:0071103	DNA conformation change	1	*0.0009	1	86	0.0643	2
Acriflavine	GO:0006259	DNA metabolic process	30	*0.0238	1	5	0.0042	1
Benomyl	GO:0007017	microtubule-based process	2	*0.0015	2	8	0.0056	2
Blasticidin S	GO:0006412	translation	772	0.5842	57	1311	0.9883	247
Bortezomib	GO:0030163	protein catabolic process	3	0.0026	1	8	0.0084	1
Brefeldin A	GO:0006888	ER to Golgi vesicle-mediated transport	565	0.4207	32	1172	0.8818	169
Caffeine	GO:0031929	TOR signaling cascade	1	*0.0007	1	1	0.0007	1
Calcofluor White	GO:0071554	cell wall organization or biogenesis	624	0.4675	90	1127	0.8526	176
Camptothecin	GO:0071103	DNA conformation change	16	*0.0114	4	6	0.0040	1
Cisplatin	GO:0006260	DNA replication	134	0.1018	23	10	0.0071	1
Daunorubicin	GO:0006260	DNA replication	70	0.0530	21	1210	0.9092	178
FK228	GO:0006325	chromatin organization	23	*0.0169	2	17	0.0131	2
Fluconazole	GO:0008202	steroid metabolic process	114	0.0870	12	708	0.5333	187
Furazolidone	GO:0006260	DNA replication	20	*0.0148	4	5	0.0034	1
Gramicidin S	GO:0071554	cell wall organization or biogenesis	286	0.2186	39	1151	0.8705	173
Griseofulvin	GO:0007017	microtubule-based process	1291	0.9718	227	750	0.5673	216
Haloperidol	GO:0008202	steroid metabolic process	5	*0.0035	2	37	0.0279	6
Hedamycin	GO:0006281	DNA repair	4	*0.0029	1	3	0.0022	1
Hydroxyurea	GO:0006260	DNA replication	29	0.0239	6	1236	0.9269	1
Itraconazole	GO:0008202	steroid metabolic process	234	0.1786	29	696	0.5239	193
Latrunculin B	GO:0007010	cytoskeleton organization	11	*0.0083	1	8	0.0068	2
Micafungin	GO:0071554	cell wall organization or biogenesis	495	0.3718	47	1134	0.8577	150
Mitomycin	GO:0006260	DNA replication	15	0.0104	4	2	0.0014	1
MMS	GO:0006281	DNA repair	3	*0.0022	1	3	0.0022	1
Mycophenolic acid	GO:0006259	DNA metabolic process	1	*0.0006	1	3	0.0025	1
Nigericin	GO:0048193	Golgi vesicle transport	157	0.1158	13	1	0.0007	1
Nocodazole	GO:0007017	microtubule-based process	2	*0.0015	2	14	0.0100	3
Oligomycin A	GO:0009268	response to pH	9	0.0075	2	2	0.0012	1
Podophyllotoxin	GO:0007017	microtubule-based process	53	0.0411	6	800	0.6038	157
Polyoxin D	GO:0071554	cell wall organization or biogenesis	1302	0.9788	225	1168	0.8828	173
Rapamycin	GO:0031929	TOR signaling cascade	156	0.1140	8	422	0.3117	9
Trichostatin A	GO:0006325	chromatin organization	23	*0.0169	3	24	0.0173	1
Tunicamycin	GO:0070085	glycosylation	1	*0.0005	1	1	0.0005	1
Tyrocidine B	GO:0071554	cell wall organization or biogenesis	5	*0.0040	1	2	0.0019	1
		Num with significant rank		22			21	
		Num with significant rank and FDR < 25%		16			0	

<https://doi.org/10.1371/journal.pcbi.1006532.t002>

annotated to that GO term. The 100 best-performing terms represented a diversity of bioprocesses related to the proteasome, glycolipid metabolism, DNA replication and repair, replication and division checkpoints, RNA splicing, microtubules, Golgi and vesicle transport, and chromatin state (S3 Fig). In contrast, the 100 worst-performing terms were bioprocesses primarily related to carbohydrate, nucleotide, and coenzyme/cofactor metabolism, as well as the mitochondria, transmembrane transport, and protein synthesis and localization (S4 Fig). The best-performing terms were also significantly smaller than the worst-performing ones (8 and 35 genes on average, respectively; rank-sum p-value $< 2.2 \times 10^{-16}$), which, given the fact that we would expect the power to increase with gene set size assuming the corresponding set was still functionally coherent, suggests that our method identifies functionally specific signal. Interestingly, the relatively poor performance of many metabolism-related bioprocess terms may result from the fact that the chemical-genetic and genetic interaction screens were both performed in relatively rich medium, precluding analysis of condition-specific phenotypes for genes only required for growth in minimal medium. While the set of best-performing terms did include a diverse range of bioprocesses, the possibility of “blind spots” should always be considered when interpreting the predictions made by CG-TARGET, as they may lead to false negative results that either exclude interesting compounds (e.g. those whose primary modes of action affect carbohydrate metabolism) or mask potential side effects of compounds whose primary modes of action are more easily observed by this method.

Application of CG-TARGET to protein complexes refines functional specificity of mode-of-action predictions

The prediction of perturbed protein complexes offers the opportunity to enhance the specificity of GO biological process predictions (especially for overly-general bioprocess terms) and investigate functional space not accessible by bioprocess annotations. As such, we investigated the potential to expand the use of CG-TARGET to the prediction of perturbed protein complexes. When CG-TARGET was applied to predict protein complex targets for the RIKEN screen data, 714 compounds were identified with at least one high-confidence ($FDR \leq 25\%$) complex prediction, 604 of which also occurred in our original set of RIKEN compounds with high-confidence bioprocess predictions. Similar, but not completely overlapping, sets of genes (Jaccard index > 0.2) contributed to the top 5 of both bioprocess and protein complex predictions for more than one third of these compounds (219; 36%); this suggested that the two standards possessed both shared and complementary functional information that could be used to improve predictions.

We observed that protein complex predictions narrowed down less-specific bioprocess terms and enabled predictions in places where bioprocess annotations were sparser. To assess the ability to refine bioprocess prediction specificity, we mapped each protein complex to the childless bioprocess terms that completely encompassed them and looked for substantial improvements in prediction strength from the bioprocess to its protein complex “child.” We observed several instances in which bioprocess predictions with $FDR > 25\%$ (not high confidence) could be converted to high-confidence predictions by refining the bioprocess term to a constituent protein complex. For example, we saw substantial gains for the following bioprocess-to-complex combinations (sizes in parentheses): “mRNA polyadenylation” (bioprocess, not high confidence; size 8) to “mRNA cleavage factor matrix” (complex, high confidence; size 4); “cytoplasmic translation” (51) to “cytoplasmic ribosomal large subunit” (24); “vacuolar acidification” (14) to “H⁺-transporting ATPase, Golgi/vacuolar” (5); and “regulation of fungal-type cell wall organization” (8) to “PKC pathway” (4) (S2 Table). Importantly, 27 of the 110 compounds with high-confidence protein complex but not bioprocess predictions achieved their high-confidence status purely based on protein complex predictions that enhanced the

specificity of a non-high-confidence, overlapping bioprocess prediction. Additionally, a separate set of 22 out of 110 compounds achieved high-confidence status based solely on predictions to protein complexes that did not strongly overlap with any bioprocesses (Jaccard < 0.2), demonstrating that the current set of protein complex annotations enabled predictions in functional space that was not well captured by a GO biological process term.

Predicting perturbed protein complexes also provided the opportunity to compare our method's performance against that of a previous, protein complex-based method called PCBA (Protein Complex-based Bayesian factor Analysis) [12]. PCBA was designed to infer the compound-induced activities of protein complexes (and thus predict compound mode of action) by linking them to observed mutant fitnesses via genetic and physical interactions. The authors highlighted six compounds in their study, five of which also possessed a high-confidence (FDR \leq 25%) CG-TARGET-based protein complex prediction. For the PCBA-based mode-of-action predictions, only two of the six compounds (benomyl and nocodazole) could be matched to their known modes of action based on protein complex activity scores alone—the remainder required additional interpretation based on the mutants that were linked to the perturbed complexes through physical or genetic interactions (S3 Table). In contrast, CG-TARGET directly generated protein complex predictions related to the known modes of action for four of the five compounds with high-confidence predictions, using only the diagnostic set of ~300 mutants (PCBA used ~3000-mutant whole-genome profiles). While the two studies used different sources of chemical-genetic profiles and protein complex annotations (which precluded more rigorous comparisons), these limited examples suggest that CG-TARGET performs at least comparably to PCBA and possibly better when focusing just on the protein complex scores. In addition, CG-TARGET can utilize arbitrary gene sets (including highly-overlapping GO biological process terms), while factor analysis-based methods such as PCBA are generally restricted to non-overlapping gene sets due to identifiability issues [12].

Assessing the compatibility of chemical-genetic and genetic interaction profiles

Our evaluations of CG-TARGET support the premise of the method that genetic interaction profiles can be used as a tool to interpret chemical-genetic interaction profiles. However, we sought to better understand the extent to which these two types of profiles actually agree with one another, and if their systematic differences could shed light on the limits of the core assumption behind our method (i.e. that chemicals mimic the interaction profiles of their genetic targets). To investigate the compatibility of chemical-genetic and genetic interaction profiles, we quantified the contribution of individual gene mutants in the chemical-genetic interaction profiles to the prediction of individual bioprocesses. For a single compound and predicted bioprocess, these “importance scores” were obtained by 1) computing a mean genetic interaction profile across all L_2 -normalized query genetic interaction profiles that possessed an inner product of 2 or higher with the chemical-genetic interaction profile and mapped to the predicted bioprocess, and 2) computing the Hadamard product (elementwise multiplication) between this mean genetic interaction profile and the compound's chemical-genetic interaction profile. Each score could have been positive, indicating agreement in the sign of chemical-genetic and genetic interactions for a gene mutant, or negative, indicating that the interactions did not agree for that gene mutant. As such, the importance scores summarized the concordance between chemical-genetic and genetic interaction profiles, conditioned on an individual compound and a perturbed bioprocess of interest.

We use the prediction of NPD4142, a compound from the RIKEN Natural Product Depository, to the “mRNA transport” bioprocess to illustrate how the overlap between chemical-

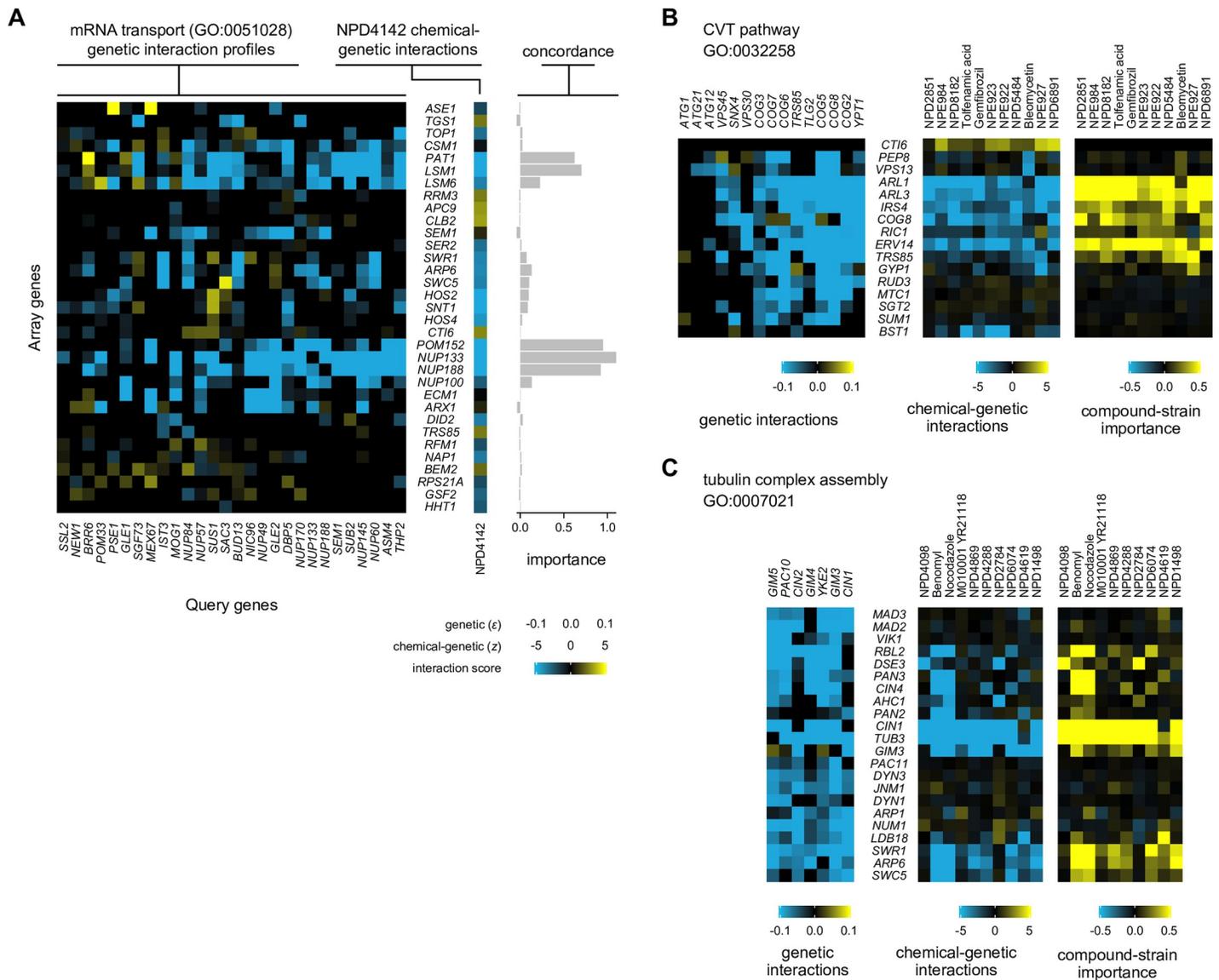


Fig 4. Detailed analysis of the contribution of individual gene mutants to biological process predictions. Each panel shows, for a bioprocess and either a compound (A) or a set of compounds (B-C) predicted to perturb that bioprocess, the subset of the respective chemical-genetic and L_2 -normalized genetic interaction profiles with signal. The importance profiles are the row-wise mean of the Hadamard product (elementwise multiplication) of each chemical-genetic interaction profile and the genetic interaction profiles for query genes with which it possessed an inner product of 2 or higher that are annotated to the GO term; they reflect the strength of each strain's contribution to the bioprocess prediction. For all panels, a query gene from the genetic interaction network was selected if it contributed to the importance score calculation for any selected compound; query genes were ordered from left to right in ascending order of their inner products (or their average, for B-C) with the selected chemical-genetic interaction profile(s). Each strain (row) was included if it passed at least one of three criteria: 1) the magnitude of its mean genetic interaction score across the selected query genes exceeded 0.04; 2) the magnitude of its chemical-genetic interaction score (for B-C, the mean of such scores) exceeded 2.5; or 3) its importance score exceeded 0.1 (for B-C, the mean of such scores). (A) Schematic showing the prediction of the "mRNA transport" bioprocess (GO:0051028) for chemical compound NPD4142. (B) Schematic showing the prediction of "CVT pathway" (FDR < 1%) for compounds whose top prediction was to that term. (C) Schematic showing the prediction of "tubulin complex assembly" (FDR < 1%).

<https://doi.org/10.1371/journal.pcbi.1006532.g004>

genetic and genetic interactions led to bioprocess predictions (Fig 4A). A qualitative examination revealed that, indeed, NPD4142 possessed a pattern of chemical-genetic interactions similar to the genetic interactions for the query genes annotated to mRNA transport. More quantitatively and as expected, we observed that the contribution of each gene mutant to a bioprocess prediction depended on the strength of its chemical-genetic interaction with

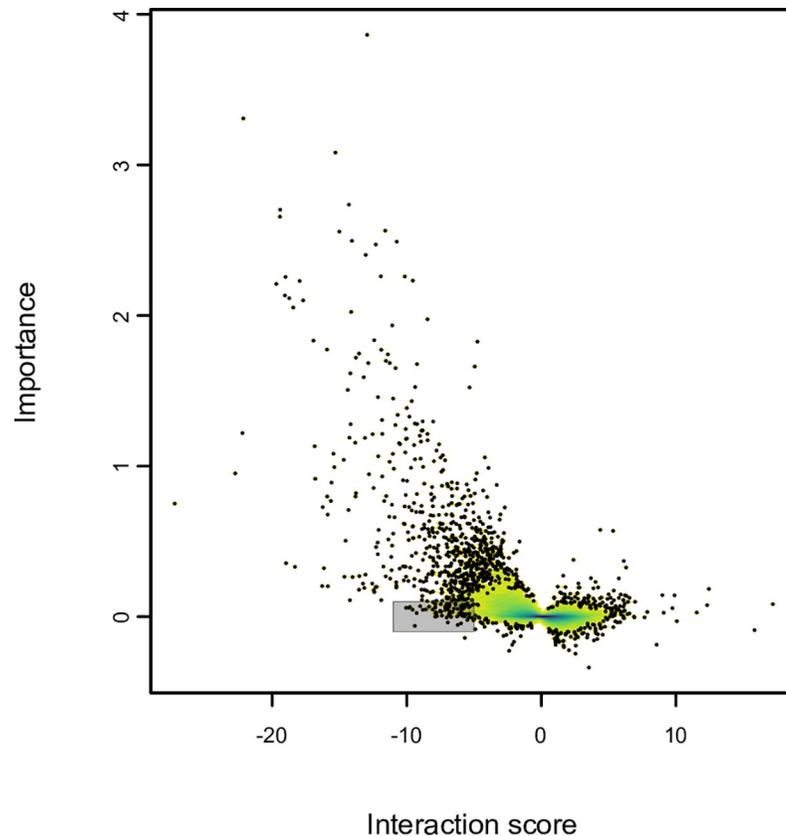


Fig 5. Global visualization of the contribution of chemical-genetic interactions to CG-TARGET bioprocess predictions. Chemical-genetic interaction profiles and their corresponding importance score profiles (see Fig 4 legend) were gathered for each of 130 diverse compounds from the high confidence set ($FDR \leq 25\%$) and their associated top bioprocess predictions. Importance is plotted as a function of chemical-genetic interaction score. One thousand points from the regions of lowest density (white) are plotted, with only density plotted in the remaining higher-density regions. Density increases in order of white, yellow, green, and violet. The shaded region highlights strains with strong negative (≤ -5) chemical-genetic interactions and no contribution (± 0.1) to a compound's top bioprocess prediction.

<https://doi.org/10.1371/journal.pcbi.1006532.g005>

NPD4142 and the number and intensity of its genetic interactions with the mRNA transport query genes. Chemical-genetic interactions with mutants of *POM152*, *NUP133*, and *NUP188*, which encode components of the nuclear pore that facilitate import and export of molecules such as mRNA, were the most important, followed by interactions with mutants in the Lsm1-7-Pat1 complex, which is involved in the degradation of cytoplasmic mRNA.

Using this approach to assess the importance of individual mutants in the chemical-genetic profile, we globally analyzed the contribution of chemical-genetic interactions to each compound's top bioprocess prediction (Fig 5). We performed this analysis twice: first, on all HCS compounds, and second, on a diverse subset of 130 compounds to correct for potential functional biases in the full set [14]. We present here the results from the 130-compound subset, although the results for the full set were qualitatively similar. For each compound, an average of 42% of its chemical-genetic interactions contributed to its top bioprocess prediction (chemical-genetic interaction cutoff ± 2.5 , importance score cutoff $+0.1$)—a fraction that increased substantially (to 78%) when limiting the analysis to each compound's strong interactions that contributed strongly (chemical-genetic interaction cutoff ± 5 , importance score cutoff $+0.5$).

Overall, we observed that more than one-third of chemical-genetic interactions (1112 / 3129) contributed to a top bioprocess prediction (chemical-genetic interaction cutoff ± 2.5 ; importance score cutoff +0.1). Strikingly, negative chemical-genetic interactions much more frequently contributed to a bioprocess prediction: approximately one-half (1071 / 2112) of negative chemical-genetic interactions contributed as compared to only ~4% (41 / 1017) of positive chemical-genetic interactions at the same cutoff. Furthermore, we observed differences in how the signs within chemical-genetic and mean genetic interaction profiles could disagree with each other despite the global profile similarity that led to bioprocess prediction, with positive chemical-genetic interactions contributing negatively to bioprocess predictions (importance score cutoff < -0.1) over 10 times more frequently than negative interactions (1.9% vs. 0.14%). This trend of negative chemical-genetic interactions supporting strong bioprocess predictions was even more pronounced when restricting this analysis to strong interactions (chemical-genetic interaction cutoff ± 5 ; importance score cutoff +0.5), where negative interactions comprised essentially the entire set of contributing chemical-genetic interactions (219 / 220, 99.5%). These observations were also supported by analyses in which we predicted perturbed bioprocesses using only negative or positive chemical-genetic interactions, finding that negative chemical-genetic interactions were the primary drivers of bioprocess predictions and overwhelmingly responsible for their accuracy [14]. We conclude that negative interactions in chemical-genetic interaction profiles contain the large majority of the functional information necessary to predict modes of action.

Negative chemical-genetic interactions also contained information reflecting general effects of chemical perturbations. Specifically, we identified nine mutant strains that exhibited strong negative chemical genetic interactions (z-score < -5) yet were enriched for a lack of contribution (importance score < 0.1) to bioprocess predictions (hypergeometric test, Benjamini-Hochberg FDR ≤ 0.05 ; shaded region of Fig 5). Manual inspection of these mutants revealed connections to the high osmolarity glycol (HOG) pathway, cell polarity (cytoskeletal actin polarization, kinetochore and chromosome segregation), and other stress response mechanisms (S4 Table). As the HOG pathway is important for the cellular response to high osmolarity and other stresses [22–24], and repolarization of the cytoskeleton is required for cells to adapt and continue dividing after stress [25,26], we hypothesize that many of these overrepresented mutants interact negatively with compounds due to an impaired ability to respond to external stress. This chemical perturbation-specific information may complement or even completely obscure the chemical-genetic signature of a compound's primary mode of action, potentially complicating the interpretation of chemical-genetic interaction profiles using a genetic interaction network.

We compared the concordance of chemical-genetic and genetic interaction profiles across multiple compounds predicted to the same bioprocess, revealing that some bioprocesses were predicted by homogenous sets of chemical-genetic interaction profiles while others were much more heterogeneous despite their predicted targeting of the same bioprocess. For example, predictions made to the “CVT pathway” (FDR $< 1\%$) depended almost entirely on a suite of strong negative chemical-genetic interactions with *ARL1*, *ARL3*, and *ERV13*, with contributions from *IRS4* and *COG8* (Fig 4B). This uniformity in the prediction of a bioprocess is contrasted by the diversity of profiles captured within “tubulin complex assembly” predictions (Fig 4C). Compounds with top predictions to this term could potentially be partitioned into three classes, divided according to strong contributions from: 1) *CIN1/TUB3*, *PAN3/CIN4*, and the SWR1 complex (known tubulin polymerization inhibitors Benomyl and Nocodazole); 2) *CIN1/TUB3* and *DSE2* (NPD4098 and NPD2784); or 3) only *CIN1/TUB3* (all remaining compounds except NPD4619). Interestingly, the structures of the compounds in each of the former two groups are distinct from those in the other groups, suggesting that the observed

diversity in these compounds' functional profiles is mechanistically derived from their structures.

Experimental validation of compound-bioprocess predictions

Phenotypic analysis of cell cycle progression. The genes and pathways that govern the cell cycle are highly conserved throughout eukaryotes, enabling researchers to infer from yeast how cells in higher organisms integrate internal and external signals to decide when to divide [27]. As such, compounds that inhibit the progression of the cell cycle in yeast may enable a better understanding of the eukaryotic cell cycle or even form the basis for new therapeutic approaches for cancer, in which the cell division cycle is dysregulated [28,29]. We observed that compounds from the RIKEN Natural Product Depository were enriched for predictions to cell cycle-related bioprocesses [14], especially to the “mitotic spindle assembly checkpoint” that occurs at the beginning of M phase. After manual inspection of these compounds' chemical-genetic interaction profiles, we selected 17 to test if our predictions validated experimentally. Specifically, we looked for increases in the percentage of cells in the G2 phase of the cell cycle (via fluorescence-activated cell sorting) and two budding phenotypes (bud size and % cells with large buds) for yeast treated with compound, together indicative of arrest at the G2/M checkpoint of the cell cycle (Fig 6A–6C). Indeed, 6 of the 17 selected compounds induced increases in any and all phenotypes, while 0 out of 10 bioactive control compounds (with high-confidence predictions to bioprocesses not related to cell cycle signaling and progression) induced increases in any of these phenotypes ($p < 0.05$, one-sided Fisher exact test). As compounds can activate the G2/M checkpoint in multiple ways (e.g. induction of DNA damage, inhibition of chromosome segregation), the set of compounds with spindle assembly checkpoint predictions can serve as a resource for studying the diversity of mechanisms by which cell cycle progression is arrested at this checkpoint and which of these may have therapeutic potential. In addition to our study of G2/M checkpoint-activating compounds, we also selected two compounds with high-confidence predictions to the term “cell-cycle phase” (mutually exclusive with mitotic spindle assembly checkpoint), one of which (NPD7834) was observed to arrest cells in G1 phase (Fig 6A–6C).

Inhibition of tubulin polymerization. Compounds that disrupt microtubules are useful for studying cell organization and division and remain promising candidates as antitumor agents [30–32]. As such, we focused on all compounds with the strongest predictions to “tubulin complex assembly” (FDR < 1%) and tested them for activity in an *in vitro*, mammalian (porcine) tubulin polymerization assay (Fig 7A). Like the previous validation experiment, a negative control set of compounds was selected at random to contain high-confidence compounds (bioprocess predictions with FDR \leq 25%) whose predictions were not related to microtubule assembly or related bioprocesses. We observed that the novel compound NPD2784 strongly inhibited tubulin polymerization, nearly as well as the drug nocodazole and more strongly than the microtubule probe benomyl. In addition, the entire set of compounds predicted to perturb tubulin complex assembly showed significantly increased inhibition of tubulin polymerization when compared to the negative control compounds ($p < 0.006$, Wilcoxon rank-sum test). Strikingly, all newly-annotated compounds were structurally novel, with a maximum structural similarity of 0.25 (computed using Braun-Blanquet similarity on all-shortest-path fingerprints of length 8) to six compounds representative of major classes of microtubule-perturbing agents (Fig 7B) [33]. Thus, we would not have identified these compounds based on structural similarity to well-characterized compounds. However, among the compounds selected for validation (known and newly-annotated microtubule-perturbing agents), we did observe that

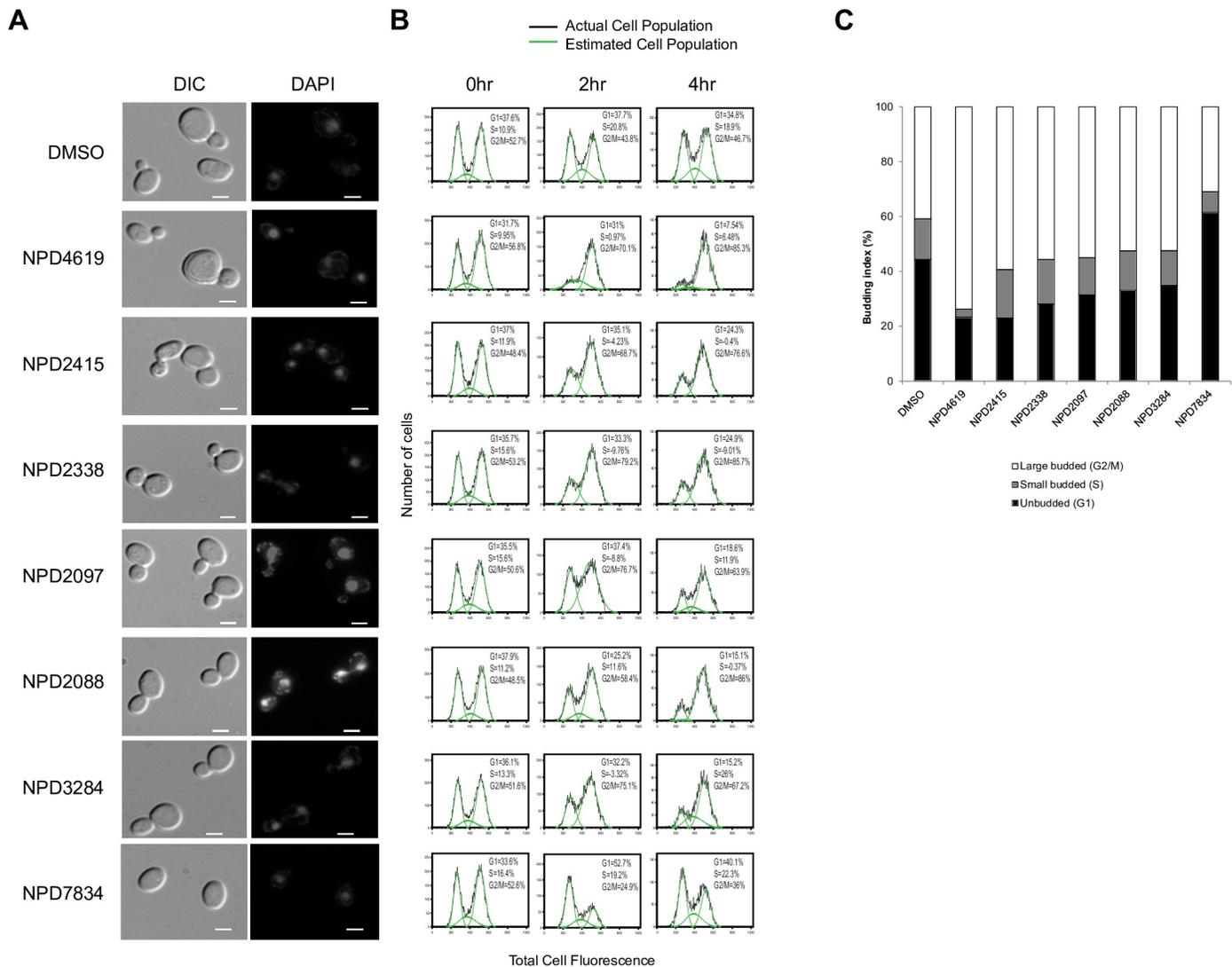


Fig 6. *In vivo* experimental validation of cell cycle-related biological process predictions. Phenotypic validation of cell cycle-related predictions, performed on drug-hypersensitive yeast treated with solvent control (DMSO) or compounds predicted to perturb the cell cycle. Out of 17 compounds predicted to arrest cells in G2/M phase, data are shown for the 6 that exhibited increases in the relevant phenotypes in any and all assays. Data for NPD7834 are also shown. (A) Differential interference contrast microscopy (DIC) and fluorescence upon DAPI staining showing bud size and DNA localization, respectively, after compound treatment. The scale bar represents a distance of 5 μ m. (B) FACS analysis of cell populations in different cell cycle phases at 0, 2, and 4 hours after compound treatment. The green curve overlay represents the estimated cell population in G1, S and G2/M phases. (C) Budding index percentages induced by treatment with compound or solvent control.

<https://doi.org/10.1371/journal.pcbi.1006532.g006>

structural similarity was predictive of the top 20% of chemical-genetic profile similarities (AUPR = 0.43 vs. 0.2 for a random classifier). This suggests that slight differences in function are influenced by structure and further exploration of compounds with similar structures may yield even more tubulin polymerization inhibitors. With this experimental validation, we have demonstrated the ability of CG-TARGET, and a genetic interaction network in general, to capture a shared mode of action across diverse compounds that can be biochemically-validated. Furthermore, we note that this validation was achieved with a mammalian tubulin assay, demonstrating the power of yeast chemical genomics coupled with CG-TARGET to predict modes of action that translate broadly to other species, including mammalian systems.

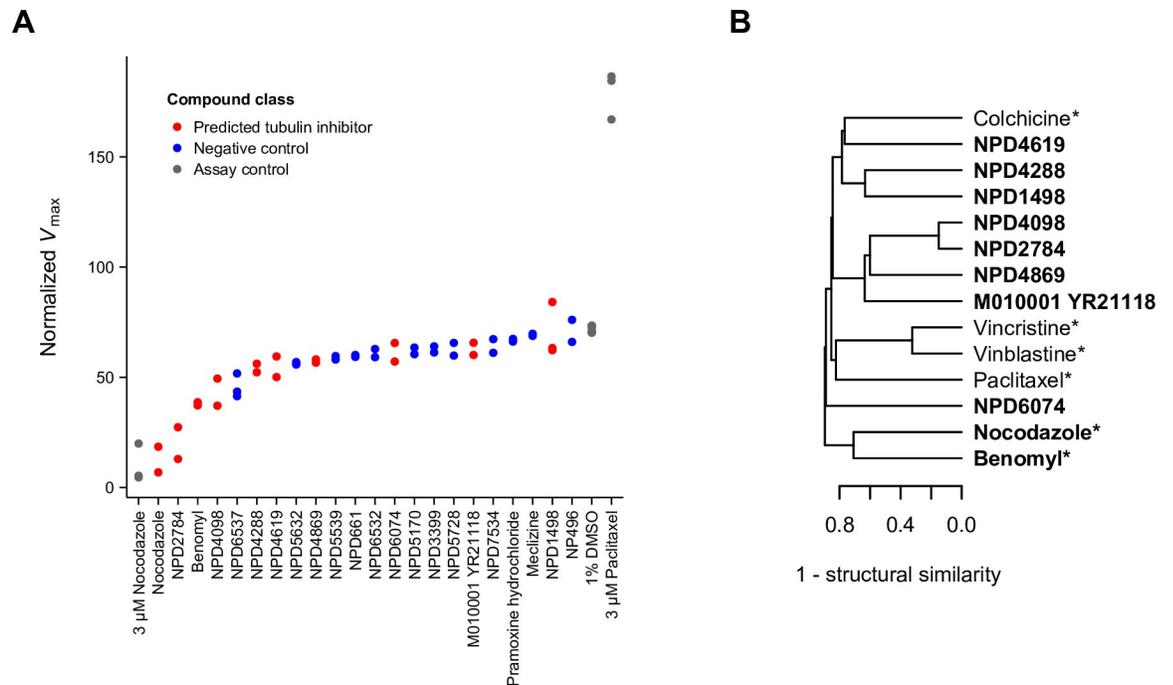


Fig 7. *In vitro* experimental validation of “tubulin complex assembly” biological process predictions. (A) *In vitro* inhibition of tubulin polymerization by compounds predicted to perturb “tubulin complex assembly” (FDR < 1%; red) compared to randomly-selected negative control compounds with high-confidence predictions to bioprocesses not related to chromosome segregation, kinetochore, spindle assembly, and microtubules (blue). V_{max} values reflecting the maximum rate of tubulin polymerization for each compound from independent replicate experiments are plotted. Assay positive and negative control compounds are colored grey. (B) Structural similarity-based hierarchical clustering of compounds tested in (A). Single linkage was used in combination with (1 - structural similarity) as the distance metric; as such, the structural similarity of the two most similar compounds at each junction can be inferred directly from the dendrogram. Compounds predicted to perturb “tubulin complex assembly” (FDR < 1%) are in bold, and known microtubule-perturbing agents are marked with an asterisk. Structural similarity was calculated as the Braun-Blanquet similarity coefficient on all-shortest-path chemical fingerprints of length 8 (see [Materials and Methods](#)).

<https://doi.org/10.1371/journal.pcbi.1006532.g007>

Discussion

The scaling of chemical-genetic interaction screens from tens or hundreds of compounds to tens of thousands of compounds has provided the opportunity, and the necessity, to develop better methods for interpreting the interaction profiles and prioritizing high-confidence compounds. We developed a method, CG-TARGET, to address this need and applied it in a recent study to predict perturbed biological processes for 1522 out of nearly 14,000 compounds screened for chemical-genetic interactions [14]. Our rigorous benchmarking of CG-TARGET showed that, in terms of accuracy, it outperformed direct enrichment on chemical-genetic interactions, and in terms of false discovery rate control, it outperformed both enrichment-based alternatives (direct enrichment and gene-target enrichment) by identifying at least 4-fold more compounds at FDR \leq 25%. Multiple experimental validations have further supported the accuracy of the method and its usefulness for functionally annotating previously uncharacterized compounds, with validations of predicted tubulin polymerization and mitotic checkpoint inhibitors presented here. The companion paper describes additional experimental validations, including one performed on 67 compounds based on linking bioprocess predictions to the stage of induced arrest in an orthogonal cell cycle assay [14].

This study is, to our knowledge, the first systematic evaluation of the ability of genetic interaction profiles to interpret chemical-genetic interaction profiles at a large scale. The results of this study are encouraging, as a genome-wide compendium of genetic interaction profiles

provides a much more comprehensive and unbiased resource for profile interpretation than a limited set of gold standard compounds. Aggregating the compound-gene similarities into compound-bioprocess predictions not only provided for increased statistical confidence but also allowed for direct functional annotation of compounds without direct protein targets (e.g. DNA-damaging or membrane-disrupting agents). Interestingly, enrichment on compound-gene similarities performed similarly to CG-TARGET in ranking bioprocess predictions for individual compounds but performed much worse on the task of prioritizing these predictions across compounds. CG-TARGET likely excelled here because it accounts both for the chemical-genetic profile strength in compound-gene similarity calculations and for the effects of general signals that arise upon treatment with bioactive compound. These general signals could be amplified through their similarity to a large cluster of profiles in the genetic interaction network and were the specific motivation for incorporating resampled profiles into the prediction scheme.

Genetic interaction-based interpretation of chemical-genetic interaction profiles has revealed broad insights into chemical function and provided interesting directions for further exploration, but some questions remain to be addressed about the limits of the technique. In the companion paper, we used the results from CG-TARGET to characterize the distribution of predicted perturbed functions for entire chemical libraries, revealing a general depletion of compound action in the nucleus and an enrichment of activity near the cell wall and membrane [14]. Additionally, we investigated the hypothesis that the profile of a compound with multiple independent modes of action would resemble a combination of distinct genetic interaction profiles, which led us to a compound whose independent predictions to cell wall and DNA perturbation were both validated (the top 20 dual-process predictions are included as Supplementary Table 2 in [14]). Indeed, we observed broad compatibility between chemical-genetic and genetic interaction profiles, the overwhelming basis of which was contributed by negative chemical-genetic interactions. However, we observed exceptions to this compatibility for genes to which perturbations may reduce the ability of cells to deal with external stress. In general, the fact that chemicals may induce stresses that cannot be recapitulated with genetic perturbations represents a potential blind spot in our approach, but one that could possibly be remedied by including specific stress conditions in the compendium of profiles used for interpretation. We do note, however, that every observed chemical-genetic or genetic interaction essentially represents an increased or decreased ability to deal with a particular stress, and many of our predictions are successful because the stresses induced by genetic and chemical perturbations overlap.

While we demonstrated here the ability to predict perturbed bioprocesses for compounds and prioritize the highest-confidence predictions, many further steps are required to identify lead compounds and ultimately develop molecular probes or pharmaceutical agents. Perturbing a biological process does not necessarily require perturbing a specific protein target, and as such, further refinements to our methods are needed to identify specific molecular targets (i.e. proteins) and prioritize the compounds most likely to perturb a small number of defined targets in the cell. We envision the use of multiple functional standards with CG-TARGET, such as biological processes and protein complexes as demonstrated here, to improve our ability to predict compound mode of action at different levels of resolution and predict the compounds that exert specific versus general effects in the cell. Different modes of chemical-genetic interaction screening can provide support in this endeavor, as heterozygous diploid mutant strains, gene overexpression strains, and/or spontaneous compound-resistant mutants can provide evidence for the direct, essential cellular target(s) of a compound [1,7]. Regardless of the limitations in predicting precise molecular targets, information about the bioprocesses perturbed by an entire library would be useful in selecting the compounds most amenable to activity

optimization and off-target effect minimization in the development of a pharmaceutical agent or molecular probe.

The approach described here can be translated to work in other species for which obtaining functional information on compounds would be useful. For example, genome-wide deletion collections have been developed for *Escherichia coli* [34] and *Schizosaccharomyces pombe* [35] and used to perform chemical-genetic interaction screens [36,37] as well as genetic interaction mapping [38–41]. Such efforts are even underway in human cell lines, enabled by genome-wide CRISPR screens [42–47]. Furthermore, future efforts to interpret chemical-genetic interaction profiles in a new species need not wait for the completion of a comprehensive, all-by-all genetic interaction network as exists in *S. cerevisiae*, as our work highlights the ability of a diagnostic set of gene mutants to capture functional information and predict perturbed biological processes. From the discovery of urgently-needed antibacterial or antifungal agents, to the treatment of orphan diseases or a better understanding of drug and chemical toxicity, the combination of chemical-genetic and genetic interactions in a high-throughput format, with appropriate analysis tools, offers a means to achieve these goals via the discovery of new compounds with previously uncharacterized modes of action.

Materials and Methods

Datasets

Chemical-genetic interaction data. Chemical-genetic interaction profiles were obtained from a recent study [14], in which nearly 14,000 compounds were screened for chemical-genetic interactions across ~300 haploid yeast gene deletion strains. The chemical-genetic interaction profiles consisted of two sub-datasets: 1) the “RIKEN” dataset, containing chemical-genetic interaction profiles spanning 289 deletion strains for 8418 compounds from the RIKEN Natural Product Depository [15] and 5724 negative experimental controls (solvent control, DMSO); and 2) the “NCI/NIH/GSK” dataset, containing chemical-genetic interactions spanning 282 deletion strains for 3565 compounds from the NCI Open Chemical Repository, the NIH Clinical Collection, and the GSK kinase inhibitor collection [16], as well as 2128 negative experimental control profiles. The solvent control profiles consisted of biological and technical replicate profiles.

Genetic interaction data. The genetic interaction dataset was obtained from a recently assembled *S. cerevisiae* genetic interaction map [5,10]; it was filtered to contain quantitative fitness observations for double mutants obtained upon crossing 1505 high-signal query gene mutants into an array of 3827 array gene mutants. The procedure for selecting the 1505 high-signal query genes out of the larger pool of 4382 is described in [14]. Briefly, each query profile was required to possess at least 40 significant genetic interactions, a sum of cosine similarity scores with all other query profiles greater than 2, and a sum of inner products with all other query profiles greater than 2. The final genetic interaction dataset used in this study was filtered to contain only array strains present in the chemical-genetic interaction datasets.

GO Biological Processes and protein complexes. A subset of terms from the “biological process” ontology within the Gene Ontology annotations [20] were used as the bioprocesses. Query genes from the *S. cerevisiae* genetic interaction dataset were mapped to biological process terms using annotations from the *Saccharomyces cerevisiae* Genome Database [19]. Both Gene Ontology and *S. cerevisiae* annotations were downloaded on September 12, 2013 from their respective databases via Bioconductor in R [48]. Terms were propagated using “is_a” relationships, such that each gene was also annotated to all parents of its direct biological process annotations. The final set of bioprocesses consisted of the terms with 4–200 gene annotations from the set of 1505 high-signal query genes in the genetic interaction dataset. For

benchmarking against the “direct enrichment” baseline method, the set of bioprocesses also consisted of terms with 4–200 gene annotations but mapped from the ~300 diagnostic deletion mutants present in the chemical-genetic interaction profiles.

Protein complex annotations were obtained from [10]. Complexes with 3 or more genes annotated to them were used as the input biological processes for CG-TARGET-based protein complex predictions.

Gold-standard compound-process annotations. Biological processes were assigned to 35 primarily antifungal compounds with chemical-genetic interaction profiles in the RIKEN dataset, based on known information about their modes of action. Bioprocess terms were selected to be specific to the compounds’ modes of action where applicable.

Predicting perturbed bioprocesses from chemical-genetic interaction profiles

Our method to predict biological processes perturbed by compounds is briefly summarized in the recent study that contains its original application to a large-scale chemical-genetic interaction dataset, generating the bioprocess predictions that are subjected to further rigorous benchmarking in this manuscript [14]. The method is more formally described here. [S1 Fig](#) and [S5 Table](#) respectively provide a schematic representation and reference for variables and symbols.

At a high-level, CG-TARGET predicts the bioprocesses perturbed by compounds in three major steps (after generating a set of randomly resampled profiles to use as a control). First, chemical-genetic interaction profiles are compared to genetic interaction profiles to generate compound-gene similarity scores. Second, these similarity scores are aggregated into compound-bioprocess scores, which are compared against score distributions derived from negative experimental control profiles, randomly resampled profiles, and randomization of the gene labels on the compound-gene scores. Finally, false discovery rate estimates are computed by comparing the rates, across a range of p-value thresholds, at which discoveries are made for negative control and randomly resampled profiles versus the discovery rate for compound-derived profiles.

Notation. We first clarify here a few uses of mathematical notation that simplify the explanation of the methods. First, the i^{th} row and column vectors of a matrix A are denoted as A_{i^*} and $A_{^*,j}$, respectively. Second, the Iverson bracket is used to convert logical propositions into values of 1 or 0, depending on if the logical proposition is true or false, respectively. This is used to simplify expressions for counting the number of elements in a vector that meet given criteria. Specifically, for a logical proposition L , the definition of the Iverson bracket is:

$$[L] = \begin{cases} 1 & \text{if } L \text{ is true} \\ 0 & \text{if } L \text{ is false} \end{cases} \quad (1)$$

The following section introduces different types of chemical-genetic interaction profiles α , β , and γ , which respectively reference treatment, negative control, and randomly resampled profiles (or scores that derive from these profiles). Instead of individually specifying which of these types are involved in each equation, we use the symbols a and b to respectively denote that a particular variable is actually multiple variables representing all profiles (C_a expands to C_{α} , C_{β} , and C_{γ}) or just the control profiles (C_b expands to C_{β} and C_{γ}). Additionally, the symbol c represents statistics derived from both types of control profiles and an additional set of statistics, denoted as δ , derived from the shuffling of gene labels (c expands to β , γ , and δ).

Data representation and overview of procedure. CG-TARGET requires chemical-genetic interaction profiles, genetic interaction profiles, and a mapping from genes to biological processes, all of which will be represented as matrices here (illustrated in [S1 Fig](#), along with

example matrix dimensions and a graphical description of the bioprocess prediction procedure). For chemical-genetic interaction matrices, let us consider an $n_m \times n_\alpha$ matrix of compound treatment profiles C_α , an $n_m \times n_\beta$ matrix of negative experimental control profiles C_β , and an $n_m \times n_\gamma$ matrix of resampled profiles C_γ , where n_m is the number of mutant strains in each chemical-genetic interaction profile, n_α is the number of profiles derived from treatment with compound, n_β is the number of profiles derived from negative experimental controls, and n_γ is the number of chemical-genetic interaction profiles resampled from C_α . The matrix G of genetic interaction profiles is $n_q \times n_q$ and the binary matrix B of gene to bioprocess mappings is $n_q \times n_p$, where n_m is the number of mutant strains in the chemical-genetic interaction and genetic interaction profiles, n_q is the number of genetic interaction profiles, and n_p is the number of bioprocesses in B annotated from the n_q genetic interaction profiles in G .

To predict perturbed biological processes, chemical-genetic interaction matrices for each profile type $a \in \{\alpha, \beta, \gamma\}$ are first converted to matrices of compound-gene similarity scores and then to matrices containing the sums of these compound-gene similarity scores for each compound-process pair. Three different z-score/p-value matrix pairs are then computed for each profile type a , two of which are derived from the control chemical-genetic interaction profile types $b \in \{\beta, \gamma\}$ (“control-derived” z-scores/p-values) and one of which is derived by randomizing the scores within each compound’s vector of compound-gene similarity scores (“within-compound” z-scores/p-values, denoted as δ). The z-score and p-value matrices across all scoring approaches $c \in \{\beta, \gamma, \delta\}$ are then combined into a final z-score/p-value matrix pair for each profile type a . The false discovery rate is estimated by comparing the rate of prediction for the treatment profiles α against that of the control profiles $b \in \{\beta, \gamma\}$ across a range of p-value thresholds. For the comparison of CG-TARGET to an enrichment-based approach, one enrichment factor/p-value matrix pair replaces the final z-score/p-value matrix pair for each profile type a , with the same false discovery rate calculations occurring afterward.

Resampled chemical-genetic interaction profiles. We construct a matrix C_γ wherein each compound-mutant interaction was drawn randomly with replacement from that mutant’s set of interaction scores across treatment (not negative control) conditions. Where $\text{rand}(x)$ is a function to randomly sample one value from x , and $\{1..n_x\}$ is the set of integers between 1 and n_x inclusive, C_γ is denoted by:

$$(C_\gamma)_{ij} = (C_\alpha)_{i, \text{rand}(\{1..n_\alpha\})}. \tag{2}$$

For this study, C_γ consisted of 50,000 resampled profiles (S1 Fig).

Mapping the similarity between chemical-genetic and genetic interaction profiles onto biological processes. An L_2 column-normalized genetic interaction matrix G' is constructed from the genetic interaction matrix G by:

$$G'_{ij} = \frac{G_{ij}}{\|G_{*j}\|_2}. \tag{3}$$

Compound-gene similarity scores are then computed as the inner product between each chemical-genetic interaction profile and L_2 -normalized genetic interaction profile:

$$S_a = (C_a)^T G'. \tag{4}$$

Compound-process scores are computed as the inner product between each compound’s vector of compound-gene similarity scores and each process’ vector of binary gene annotations. Each compound-process score is thus the sum of a compound’s gene similarity scores

within each process, which is denoted by:

$$X_a = S_a B. \tag{5}$$

Computing statistics on biological process predictions with CG-TARGET. For each compound-process score, we compute a z-score and empirical p-value based on the distribution of that process' scores across the two types of control profiles ("control-derived") and also upon shuffling the gene labels of the compound-gene scores and recomputing compound-process scores ("within-profile"). The two control-derived z-scores require vectors containing the mean and standard deviation of each process' scores across the control profiles, as denoted by:

$$\begin{aligned} (u_b)_j &= \frac{1}{n_b} \sum_{i=1}^{n_b} (X_b)_{ij} \\ (v_b)_j &= \sqrt{\frac{1}{n_b - 1} \sum_{i=1}^{n_b} ((X_b)_{ij} - (u_b)_j)^2} \end{aligned} \tag{6}$$

The resulting control-derived z-score matrices are computed as:

$$(Z^*_{(a,b)})_{ij} = \frac{(X_a)_{ij} - (u_b)_j}{(v_b)_j}. \tag{7}$$

The p-value that accompanies each control-derived compound-process z-score is computed by counting the number of times the compound-process score is less than or equal to the control-derived scores for that process, as denoted by:

$$(P_{Z^*_{(a,b)}})_{ij} = \frac{1}{n_b} \sum_{k=1}^{n_b} [(X_a)_{ij} \leq (X_b)_{kj}]. \tag{8}$$

Each within-profile compound-process z-score compares the mean of the compound's gene similarity scores within the process to the mean and standard deviation the compound's entire set of gene similarity scores. These compound-wise means and standard deviations are denoted as the following w_a and y_a vectors, respectively:

$$\begin{aligned} (w_a)_i &= \frac{1}{n_q} \sum_{j=1}^{n_q} (S_a)_{ij} \\ (y_a)_i &= \sqrt{\frac{1}{n_q - 1} \sum_{j=1}^{n_q} ((S_a)_{ij} - (w_a)_i)^2} \end{aligned} \tag{9}$$

The within-profile compound-process z-scores are computed as follows, where d is a vector containing the sizes of each process term:

$$\begin{aligned} d_j &= \sum_{i=1}^{n_q} B_{ij} \\ (Z^*_{(a,\delta)})_{ij} &= \frac{(X_a)_{ij}/d_j - (w_a)_i}{(y_a)_i/\sqrt{d_j}} \end{aligned} \tag{10}$$

The p-value that accompanies each within-profile compound-process z-score is computed by counting the number of times that the compound-process score is less than or equal to compound-process scores in a distribution that results from recomputing these scores after

randomly permuting the compound's gene similarity scores. Where ${}^k S_a$ represents the k^{th} row-wise permutation (out of n_l total permutations) of the compound-gene similarity score matrix S_a , the within-profile compound-process p-value matrix is denoted by:

$${}^k X_a = {}^k S_a B$$

$$(P_{Z^*(a,\delta)})_{ij} = \frac{1}{n_l} \sum_{k=1}^{n_l} [(X_a)_{ij} \leq ({}^k X_a)_{ij}] \quad (11)$$

Ultimately, the different p-values and z-scores for each compound-process pair are combined into one p-value and z-score for that pair. These scores are combined such that the largest (least significant) p-value is chosen along with its associated z-score. If multiple p-values tie for the largest value, then the one with the smallest associated z-score is chosen. As such, the resulting combination of p-value and z-score represents the most conservative estimate of the strength and significance of the prediction from compound to perturbed biological process.

To combine the p-values and z-scores, a matrix $Psources_a$ is first created to determine, for each compound-process pair, which p-value and z-score matrices will contribute the final p-value and z-score. For each z-score/p-value scoring approach in c , each entry of this matrix is denoted by:

$$f_p(\epsilon) = (P_{Z^*(a,\epsilon)})_{ij}$$

$$f_z(\epsilon) = (Z_{(a,\epsilon)}^*)_{ij} \quad (12)$$

$$(Psources_a)_{ij} = \operatorname{argmin}_{c'} f_z(c') \text{ where } c' \in \operatorname{argmax}_{c \in \{\beta,\gamma,\delta\}} f_p(c).$$

The resulting final p-value and z-score matrices for each profile type $a \in (\alpha, \beta, \gamma)$ are then:

$$(Z_{(a)})_{ij} = (Z_{(a,(Psources_a)_{ij})}^*)_{ij}$$

$$(P_{Z(a)})_{ij} = (P_{Z^*(a,(Psources_a)_{ij})})_{ij} \quad (13)$$

Computing biological process enrichments. Two enrichment-based methods for predicting biological processes perturbed by compounds were also implemented to provide appropriate baselines for assessing the performance of CG-TARGET. The “direct enrichment” method computed, for each compound, biological process enrichment on the 20 mutants with the strongest negative chemical-genetic interactions. The “gene-target enrichment” method computed, for each compound, biological process enrichment within the genes that contributed the top n compound-gene similarity scores for each compound. For either of these approaches, two sets of matrices are computed, $E_{(a,n)}$ and $P_{E(a,n)}$, which respectively contain the enrichment factor and hypergeometric p-value for each compound and biological process pair. For gene-target enrichment, we computed enrichments for $n \in \{10, 20, 50, 100, 200, 300, 400, 600, 800\}$.

First, a binary matrix $X_{(a,k)}^{\text{top}}$ is derived from the matrix of compound-gene similarity scores X_a , such that in each row, the positions corresponding to the top n scores are set to 1 and the remaining positions are set to 0. This is denoted as:

$$(X_{(a,n)}^{\text{top}})_{ij} = [(X_a)_{ij} \geq (\operatorname{sortDesc}((X_a)_{i,*}))_n] \quad (14)$$

where $\operatorname{sortDesc}(x)$ is a function that returns the values in a vector x sorted in descending order. The final enrichment factor and p-value matrices are then computed as:

$$(E_{(a,n)})_{ij} = \frac{((X_{(a,n)}^{\text{top}})_{i,*} B_{*j}) n_q}{(\sum B_{*j}) n} \quad (15)$$

$$(P_{E(a,n)})_{ij} = 1 - \operatorname{hygeCDF}(n_q, \sum B_{*j}, n, ((X_{(a,n)}^{\text{top}})_{i,*} B_{*j}) - 1)$$

where $B_{\cdot,j}$ is a binary vector of gene annotations for the j^{th} bioprocess and $\text{hygeCDF}(N, K, n, k)$ is the cumulative hypergeometric distribution given a population size of N with K success states and n draws with k observed successes.

Estimating the false discovery rate. The false discovery rates of the compound-process predictions are estimated by comparing, using the entire range of observed p-values as thresholds, the number of compounds with at least one bioprocess prediction against the number of experimental controls and resampled profiles with at least one bioprocess prediction at each threshold. We compute a false discovery rate matrix FDR_b for the treatment profiles α against each control profile type $b \in \{\beta, \gamma\}$. This FDR_b matrix is individually computed for the CG-TARGET-based compound-process predictions as well as for the enrichment-based compound-process predictions (using the p-value matrices $P_{Z(a)}$ and $P_{E(a,n)}$); for simplicity, we do not change the notation of FDR_b to reflect if the false discovery rate values were computed on the output from CG-TARGET or our baseline enrichment-based approaches.

The first step in computing the false discovery rate is obtaining a vector $ptop_a$ that contains the smallest process prediction p-value for each compound. Additionally, the union of all observed p-values p_{all} defines the universe of p-values for which corresponding false discovery rates will be computed. Given p-value matrices P_a ($P_{Z(a)}$ or $P_{E(a,n)}$ for one value of n) and a function $\text{sortAsc}()$ that returns the input values sorted in ascending order, the vectors $ptop_a$ and p_{all} are given by:

$$\begin{aligned} (ptop_a)_i &= \min((P_a)_{i,*}) \\ p_{all} &= \text{sortAsc}\left(\bigcup_{i,j,a \in \{\alpha, \beta, \gamma\}} (P_a)_{i,j}\right). \end{aligned} \tag{16}$$

We then compute a mapping from each observed p-value to its corresponding false discovery rate, with mappings generated with respect to each control profile type $b \in \{\beta, \gamma\}$. First, a vector of false discovery rates r^*_b is computed, each value corresponding to a p-value threshold in p_{all} , by dividing the fraction of treatment profiles with one or more bioprocess predictions that pass the threshold by the fraction of control profiles that also pass the threshold. As the p-values in the vector p_{all} are monotonically increasing, it is desirable for the false discovery rate to increase monotonically with the p-value. However, it is possible for the false discovery rate to decrease as p-value increases (if the fraction of treatment profiles passing the threshold increases faster than the fraction of control profiles passing the threshold), and thus we adjust each false discovery rate value in the vector r^*_b to be the minimum of its current value or any value at a larger index to generate a new vector r_b (similar to the Benjamini-Hochberg procedure [49]). The final p-value to false discovery rate mappings can be written as a function of the p-value p , with the procedure to generate these mappings given by:

$$\begin{aligned} (r^*_b)_i &= \frac{\frac{1}{n_b} \sum_{j=1}^{n_b} [(ptop_b)_j \leq (p_{all})_i]}{\frac{1}{n_x} \sum_{j=1}^{n_x} [(ptop_x)_j \leq (p_{all})_i]} \\ r_b &= \text{rev}(\text{cumMin}(\text{rev}(r^*_b))) \\ f_{FDR(b)}(p) &= (r_b)_{\{i: (r_b)_i=p\}} \end{aligned} \tag{17}$$

Given this mapping of p-value to false discovery rate, the resulting matrices of false discovery rates with respect to control profile types $b \in \{\beta, \gamma\}$ are given by:

$$(FDR_b)_{i,j} = f_{FDR(b)}((P_a)_{i,j}). \tag{18}$$

Computational evaluation of bioprocess predictions

Performance on simulated chemical-genetic interaction profiles. We generated a set of simulated chemical-genetic interaction profiles derived from genetic interaction profiles [14]. Each simulated chemical-genetic interaction profile was a query genetic interaction profile augmented with noise sampled from a Gaussian distribution with a mean of 0 and a variance for each array gene twice that of the same array gene in the genetic interaction dataset. Three simulated profiles were generated based on each query gene, resulting in 4515 total profiles. Because each simulated chemical-genetic interaction profile was derived from a query genetic interaction profile, it inherited the gold-standard bioprocess annotations from its parent genetic interaction profile in subsequent benchmarking efforts.

We then used CG-TARGET and each top-*n* enrichment method to predict perturbed bioprocesses for this set of 4515 simulated chemicals x 289 deletion mutants. For each simulated chemical, its top bioprocess prediction was compared to the set of inherited gold-standard bioprocess annotations, counting as a true positive if the top prediction matched an existing annotation and a false positive if it did not. Precision-recall curves were then generated by sorting the list of each simulated chemical's top bioprocess predictions (p-value ascending, z-score or enrichment factor descending) and computing the precision (true positives / (true positives + false positives)) and recall (true positives) at each point in this list.

Performance on gold-standard compound-bioprocess annotations. The predicted perturbed bioprocesses for each of the gold-standard compounds were sorted, first in ascending order by their p-value and then descending order by their z-score (for CG-TARGET) or enrichment factor (top-*n* enrichment), and the rank of each compound's gold-standard bioprocess annotation was recorded. To assess the significance of each rank, each pair of p-value and z-score was randomly assigned to a new bioprocess (without replacement), the lists re-ordered, and the ranks of each compound's target bioprocess re-computed. The empirical p-value for each gold-standard compound-process pair was computed as the number of times the rank from the shuffled bioprocesses achieved the same or better rank as the observed rank, divided by the number of randomizations. These randomizations were also used as a baseline against which to compare the number of compounds (out of 35) that achieved a given rank, as seen in Fig 3 and S1 Fig; the displayed ribbons were generated by calculating, for each rank, the relevant percentiles on the distribution of compounds with randomized predictions that achieved that rank. The "effective rank" of a compound's gold-standard bioprocess annotation was determined as the minimum rank of any bioprocess term with which it possessed sufficient gene annotation similarity (overlap index ≥ 0.4 , where the overlap index of two sets is defined as the size of the intersection divided by the size of the smaller set).

Characterizing performance with respect to individual bioprocess terms. For each propagated GO biological process term used for bioprocess prediction, we gathered all predictions to that term across the 4515 simulated chemical-genetic interaction profiles and sorted the predictions in ascending order by p-value and then in descending order by z-score. The area under the precision-recall curve (AUPR) was calculated across this sorted list of simulated compounds, with a true positive defined as the occurrence of a simulated compound that was annotated to the bioprocess (via the simulated compound's parent gene). To obtain the final evaluation statistic for each GO term, this AUPR was divided by the AUPR of a random classifier, which is equal to the number of true positives divided by the total number of simulated compounds.

Assessing the compatibility of chemical-genetic and genetic interaction profiles

Analysis of bioprocess prediction drivers in chemical-genetic interaction data. Given a compound and a predicted bioprocess, a profile of “importance scores” describes the contribution of each gene mutant to that compound’s bioprocess prediction. To obtain this score, a mean genetic interaction profile was first computed across all L_2 -normalized genetic interaction profiles annotated to the biological process for which the inner product with the compound’s chemical-genetic interaction profile was 2 or greater. The importance score profile was then obtained by taking the Hadamard product (elementwise multiplication) between this mean genetic interaction profile and the compound’s chemical-genetic interaction profile.

Overrepresentation analyses of gene mutants with strong chemical-genetic and/or genetic interactions. After restricting the data to the top biological process prediction for each compound, gene mutants that possessed strong, negative chemical-genetic interaction scores (z-score < -5) were assessed for overrepresentation with respect to the number of times they did not contribute (importance score within ± 0.1) to a compound’s top bioprocess prediction. Specifically, the number of times each strain occurred inside and outside the region described above (grey box in Fig 5) was compared to the number of times all strains occurred inside and outside the region using a hypergeometric test, using all strains with interaction z-scores < -5 as the background set. Details on the genes overrepresented in this region are given in S4 Table.

Experimental validation of compound-bioprocess predictions

Phenotypic analysis of cell cycle progression. To examine the effect of compounds on arresting cells in G2/M phase, we looked for differences in budding index and cell DNA content between compounds predicted to perturb the cell cycle versus negative control compounds. Seventeen compounds with high-confidence predictions to the bioprocess term “mitotic spindle assembly checkpoint” and strong negative chemical-genetic interactions with *PAT1* and *LSM6* (a common signature for compounds with this bioprocess prediction) were selected for validation. Additionally, ten bioactive (growth inhibition 50–80% compared to DMSO control) compounds with high confidence predictions (false discovery rate $\leq 25\%$) to bioprocess terms not related to cell cycle signaling and progression were selected as negative controls. Two compounds predicted to perturb “cell cycle phase” were also tested in these experiments. All compounds were tested at a concentration of 10 $\mu\text{g/mL}$, which was also the concentration used for chemical genomic screening [14].

To quantify budding index, logarithmically-growing *pdr1 Δ pdr3 Δ snq2 Δ* cells were transferred to fresh galactose-containing medium (YPGal) containing compounds and incubated at 25 °C for 4 hours. The budding status of at least 200 cells was visually determined under the microscope. The percentage of the budded cells in no compound or compound-treated samples was counted.

For flow cytometry analysis, log phase *pdr1 Δ pdr3 Δ snq2 Δ* cells were grown in YPGal media in the presence or absence of a compound for 4 hours; they were then fixed in 70% ethanol for 1 hour at 25 °C. Cells were collected by centrifugation, washed, and resuspended in buffer containing RNase A (0.25 mg/mL in 50 mM Tris, pH 7.5) for 1.5 hours. Cells were further incubated in 20 μl of 20 mg/ml proteinase K at 50 °C for 1 hour. Samples were then stained with propidium iodide, briefly sonicated, and measured using FACSCalibur ver 2.0 (Becton Dickinson, CA, USA).

The proportions of predicted active compounds and negative controls with positive phenotypic results were compared using the `prop.test` function in R to assess significance.

Inhibition of tubulin polymerization. *In vitro* tubulin polymerization assays using a fluorescent-based porcine tubulin polymerization assay (Cytoskeleton, BK011P) were performed following manufacturer specifications. Compounds were tested at a concentration of 10 $\mu\text{g/ml}$ (with the exception of assay controls), which was identical to the concentration used for chemical genomic screening. All ten compounds predicted to perturb “tubulin complex assembly” with the minimum estimated false discovery rate ($\text{FDR} < 1\%$) were selected for testing. Twelve compounds with predictions of false discovery rate $\leq 25\%$ to any bioprocess except those related to chromosome segregation, kinetochore, spindle assembly, and microtubules were randomly selected as negative controls.

The degree of tubulin polymerization inhibition was summarized in a single V_{max} statistic for each compound treatment replicate. The V_{max} for each compound’s fluorescence time-course was calculated as the maximum change in fluorescence between consecutive time points, which were measured at 1-minute intervals. Three batches of experiments were performed in total (resulting in $N \geq 2$ for each compound), and we normalized the V_{max} values in each batch by subtracting the difference between that batch’s mean DMSO (solvent control) V_{max} and the overall mean DMSO V_{max} . To determine if the tubulin-predicted compounds inhibited polymerization to a significantly greater degree than the controls, we calculated the mean of the normalized V_{max} values for each compound and performed a one-sided Wilcoxon rank-sum to test for a difference in the ranks of these values between the two classes of compounds.

Chemical structure similarities between each pair of compounds selected for tubulin polymerization validation were obtained by first computing an all-shortest-paths fingerprint with path length 8 for each compound [50]. Similarities were computed on the fingerprints using the Braun-Blanquet similarity coefficient, which is defined as the size of the intersection divided by the size of the larger set. In a recent study, this combination of structure descriptor and similarity coefficient performed well when evaluated globally on our entire chemical-genetic interaction dataset [51]. Chemical structures are available from the MOSAIC database [52].

Supporting information

S1 Fig. Schematic representation of CG-TARGET bioprocess prediction procedure Further details on the presented procedures, including equations, are given in “Predicting the biological processes perturbed by compounds” in Materials and Methods. (PDF)

S2 Fig. Performance comparison of CG-TARGET versus baseline enrichment approaches Perturbed biological processes were predicted using both CG-TARGET and methods that calculated enrichment on the set of each compound’s n most similar genetic interaction profiles (“top n ,” $n \in \{10, 20, 50, 100, 200, 300, 400, 600, 800\}$). (A) Bioprocess prediction false discovery rate estimates derived from resampled chemical-genetic interaction profiles, performed on compounds from the RIKEN dataset. (B) Precision-recall analysis of the ability to recapitulate gold-standard annotations within the set of top bioprocess predictions for ~ 4500 simulated compounds. Each simulated compound was designed to target one query gene in the genetic interaction network and thus inherited gold-standard bioprocess annotations from its target gene. (C) For each of 35 well-characterized compounds in the RIKEN dataset with literature-derived, gold-standard bioprocess annotations, we determined the rank of its gold-standard bioprocess within its list of predictions. The number of compounds for which a given rank (or better) was achieved is plotted. The grey ribbons represent the median, interquartile range

(25th to 75th percentiles), and 95% confidence interval of 10,000 rank permutations. (PDF)

S3 Fig. Induced GO hierarchy of the 100 best-performing GO biological process terms, evaluated using simulated chemical-genetic interaction profiles Each term was evaluated using precision-recall statistics (area under the precision-recall curve divided by the area under a curve produced by a random classifier) to analyze its ability to rank simulated chemical-genetic interaction profiles from which it was annotated as a gold-standard bioprocess. Green nodes represent the 100 best-performing GO biological process terms, yellow nodes represent terms for which predictions were made but did not rank among the top 100, and white nodes represent terms in the Biological Process ontology that were not selected for bioprocess prediction. Hovering the mouse over each node reveals its GO ID and name. (HTML)

S4 Fig. Induced GO hierarchy of the 100 worst-performing GO biological process terms, evaluated using simulated chemical-genetic interaction profiles Same as [S3 Fig](#), but for the 100-worst performing GO biological process terms. (HTML)

S1 Table. Comparison of CG-TARGET GO biological process mode-of-action predictions to direct GO enrichment on chemical-genetic interaction profiles Each row shows the top prediction for one of 35 well-characterized compounds, with predictions generated by either enrichment on the top 20 negative chemical-genetic interaction scores (“direct enrichment”) or using CG-TARGET. Gold-standard bioprocess annotations for the compounds, with literature support, were used to qualitatively determine if each compound’s top bioprocess prediction matched what was known about that compound. For direct enrichment, the association p-value was derived from the hypergeometric CDF and the Benjamini-Hochberg FDR was computed for each compound’s set of enrichments. All false discovery rates were generated by comparing the rate of resampled profile predictions to the rate of treatment profile predictions across the range of observed p-values. Driver genes are the members of a bioprocess that led to its prediction. (XLSX)

S2 Table. Using protein complexes to refine CG-TARGET GO biological process mode-of-action predictions Compounds, GO biological processes, and protein complexes are shown if the mode-of-action prediction to the protein complex was stronger than that to the associated GO biological process (comparison first based on p-value, then on z-score in the case of a tie). Protein complexes were limited to those of size 4 or greater whose gene annotations were a subset of those for the corresponding GO biological process term. The final column indicates compounds that did not achieve a false discovery rate of 25% or less for any GO biological process mode-of-action predictions but did for at least one protein complex prediction (with “HCS” denoting “high confidence set”). (XLSX)

S3 Table. Comparison of CG-TARGET protein complex predictions to Protein Complex-based, Bayesian factor Analysis (PCBA) Mode-of-action predictions were highlighted for six compounds in the PCBA study [12], all of which were also included in this study. For the CG-TARGET-based predictions, only the top protein complex prediction for each compound was retained. For the PCBA-based predictions here, the highlighted modes of action were based on 1) protein complexes with predicted altered activity in the presence of compound and 2) gene ontology enrichments performed directly on the strains (filtered by their

contributions to the inference of protein complex activity).
(XLSX)

S4 Table. Overrepresentation analysis of mutant strains with strong negative chemical-genetic interactions and no contribution to top bioprocess predictions Overrepresentation within the shaded region of Fig 5 was evaluated using a hypergeometric test to compare the occurrence of one strain versus all strains inside and outside of the region, with the background containing only strains that possessed strong (z -score < -5) negative chemical-genetic interactions. The compounds and top bioprocess predictions associated with each strain's occurrences in the region are given, as well as the appropriate background list of strains and information on the gene deleted in each strain.
(XLSX)

S5 Table. Reference for variables and symbols used to describe the CG-TARGET method in Materials and Methods Any instance of the symbols a , b , or c should be expanded into individual expressions for each of the members of those respective sets (i.e. C_a becomes C_{α} , C_{β} , and C_{γ}). Where two of these symbols appear in a variable's subscript, that variable exists for all pairwise combinations of those set members.
(DOCX)

S1 Data. Descriptions of supporting datasets uploaded to the Dryad digital repository All supporting datasets are available from the Dryad digital repository at the following link: <https://dx.doi.org/10.5061/dryad.nr2cf12>. They represent all raw data and final results for the application of CG-TARGET to the RIKEN chemical-genetic interaction screen.
(DOCX)

Acknowledgments

SWS would like to thank Henry Neil Ward for his proofreading of the manuscript. Computing resources and data storage services were partially provided by the Minnesota Supercomputing Institute and the UMN Office of Information Technology, respectively. Software licensing services were provided by the UMN Office for Technology Commercialization.

Author Contributions

Conceptualization: Scott W. Simpkins, Justin Nelson, Raamesh Deshpande, Sheena C. Li, Jeff S. Piotrowski, Hiroyuki Osada, Minoru Yoshida, Charles Boone, Chad L. Myers.

Data curation: Scott W. Simpkins, Raamesh Deshpande.

Formal analysis: Scott W. Simpkins.

Funding acquisition: Michael Costanzo, Yoko Yashiroda, Yoshikazu Ohya, Hiroyuki Osada, Minoru Yoshida, Charles Boone, Chad L. Myers.

Investigation: Scott W. Simpkins, Justin Nelson, Raamesh Deshpande, Sheena C. Li, Jeff S. Piotrowski, Erin H. Wilson, Abraham A. Gebre, Hamid Safizadeh, Reika Okamoto, Mami Yoshimura, Yoko Yashiroda.

Methodology: Scott W. Simpkins.

Project administration: Yoko Yashiroda, Hiroyuki Osada, Minoru Yoshida, Charles Boone, Chad L. Myers.

Resources: Michael Costanzo, Yoshikazu Ohya, Charles Boone.

Software: Scott W. Simpkins, Raamesh Deshpande.

Supervision: Michael Costanzo, Yoko Yashiroda, Yoshikazu Ohya, Minoru Yoshida, Charles Boone, Chad L. Myers.

Validation: Sheena C. Li, Jeff S. Piotrowski, Abraham A. Gebre, Reika Okamoto, Mami Yoshimura, Yoko Yashiroda, Yoshikazu Ohya.

Visualization: Scott W. Simpkins, Justin Nelson, Raamesh Deshpande.

Writing – original draft: Scott W. Simpkins.

Writing – review & editing: Scott W. Simpkins, Justin Nelson, Raamesh Deshpande, Sheena C. Li, Jeff S. Piotrowski, Abraham A. Gebre, Hamid Safizadeh, Michael Costanzo, Yoshikazu Ohya, Hiroyuki Osada, Minoru Yoshida, Charles Boone, Chad L. Myers.

References

- Giaever G, Shoemaker DD, Jones TW, Liang H, Winzeler EA, Astromoff A, et al. Genomic profiling of drug sensitivities via induced haploinsufficiency. *Nat Genet.* 1999 Mar; 21(3):278–83. <https://doi.org/10.1038/6791> PMID: 10080179
- Parsons AB, Brost RL, Ding H, Li Z, Zhang C, Sheikh B, et al. Integration of chemical-genetic and genetic interaction data links bioactive compounds to cellular target pathways. *Nat Biotechnol.* 2004 Jan; 22(1):62–9. <https://doi.org/10.1038/nbt919> PMID: 14661025
- Parsons AB, Lopez A, Givoni IE, Williams DE, Gray CA, Porter J, et al. Exploring the mode-of-action of bioactive compounds by chemical-genetic profiling in yeast. *Cell.* 2006 Aug 11; 126(3):611–25. <https://doi.org/10.1016/j.cell.2006.06.040> PMID: 16901791
- Hillenmeyer ME, Fung E, Wildenhain J, Pierce SE, Hoon S, Lee W, et al. The chemical genomic portrait of yeast: uncovering a phenotype for all genes. *Science.* 2008 Apr 18; 320(5874):362–5. <https://doi.org/10.1126/science.1150021> PMID: 18420932
- Costanzo M, Baryshnikova A, Bellay J, Kim Y, Spear ED, Sevier CS, et al. The genetic landscape of a cell. *Science.* 2010 Jan 22; 327(5964):425–31. <https://doi.org/10.1126/science.1180823> PMID: 20093466
- Hoepfner D, Helliwell SB, Sadlish H, Schuierer S, Filipuzzi I, Brachat S, et al. High-resolution chemical dissection of a model eukaryote reveals targets, pathways and gene functions. *Microbiol Res.* 2014 Mar; 169(2–3):107–20. <https://doi.org/10.1016/j.micres.2013.11.004> PMID: 24360837
- Lee AY, St Onge RP, Proctor MJ, Wallace IM, Nile AH, Spagnuolo PA, et al. Mapping the cellular response to small molecules using chemogenomic fitness signatures. *Science.* 2014 Apr 11; 344(6180):208–11. <https://doi.org/10.1126/science.1250217> PMID: 24723613
- Wildenhain J, Spitzer M, Dolma S, Jarvik N, White R, Roy M, et al. Prediction of Synergism from Chemical-Genetic Interactions by Machine Learning. *Cell Syst.* 2015 Dec; 1(6):383–95. <https://doi.org/10.1016/j.cels.2015.12.003> PMID: 27136353
- Smith AM, Heisler LE, Mellor J, Kaper F, Thompson MJ, Chee M, et al. Quantitative phenotyping via deep barcode sequencing. *Genome Res.* 2009 Oct; 19(10):1836–42. <https://doi.org/10.1101/gr.093955.109> PMID: 19622793
- Costanzo M, VanderSluis B, Koch EN, Baryshnikova A, Pons C, Tan G, et al. A global genetic interaction network maps a wiring diagram of cellular function. *Science.* 2016 Sep 23; 353(6306).
- Flaherty P, Giaever G, Kumm J, Jordan MI, Arkin AP. A latent variable model for chemogenomic profiling. *Bioinforma Oxf Engl.* 2005 Aug 1; 21(15):3286–93.
- Han S, Kim D. Inference of protein complex activities from chemical-genetic profile and its applications: predicting drug-target pathways. *PLoS Comput Biol.* 2008 Aug 29; 4(8):e1000162. <https://doi.org/10.1371/journal.pcbi.1000162> PMID: 18769708
- Hillenmeyer ME, Ericson E, Davis RW, Nislow C, Koller D, Giaever G. Systematic analysis of genome-wide fitness data in yeast reveals novel gene function and drug action. *Genome Biol.* 2010; 11(3):R30. <https://doi.org/10.1186/gb-2010-11-3-r30> PMID: 20226027
- Piotrowski JS, Li SC, Deshpande R, Simpkins SW, Nelson J, Yashiroda Y, et al. Functional annotation of chemical libraries across diverse biological processes. *Nat Chem Biol.* 2017 Sep; 13(9):982–93. <https://doi.org/10.1038/nchembio.2436> PMID: 28759014

15. Kato N, Takahashi S, Nogawa T, Saito T, Osada H. Construction of a microbial natural product library for chemical biology studies. *Curr Opin Chem Biol.* 2012 Apr; 16(1–2):101–8. <https://doi.org/10.1016/j.cbpa.2012.02.016> PMID: 22406171
16. Drewry DH, Willson TM, Zuercher WJ. Seeding collaborations to advance kinase science with the GSK Published Kinase Inhibitor Set (PKIS). *Curr Top Med Chem.* 2014; 14(3):340–2. <https://doi.org/10.2174/1568026613666131127160819> PMID: 24283969
17. Deshpande R, Nelson J, Simpkins SW, Costanzo M, Piotrowski JS, Li SC, et al. Efficient strategies for screening large-scale genetic interaction networks. *bioRxiv [Internet].* 2017 Jul 5; Available from: <http://biorxiv.org/content/early/2017/07/05/159632.abstract>
18. Baryshnikova A, Costanzo M, Kim Y, Ding H, Koh J, Toufighi K, et al. Quantitative analysis of fitness and genetic interactions in yeast on a genome scale. *Nat Methods.* 2010 Dec; 7(12):1017–24. <https://doi.org/10.1038/nmeth.1534> PMID: 21076421
19. Cherry JM, Hong EL, Amundsen C, Balakrishnan R, Binkley G, Chan ET, et al. Saccharomyces Genome Database: the genomics resource of budding yeast. *Nucleic Acids Res.* 2012 Jan; 40(Database issue):D700–705. <https://doi.org/10.1093/nar/gkr1029> PMID: 22110037
20. Gene Ontology Consortium. Gene Ontology Consortium: going forward. *Nucleic Acids Res.* 2015 Jan; 43(Database issue):D1049–1056. <https://doi.org/10.1093/nar/gku1179> PMID: 25428369
21. Deshpande R, Vandersluis B, Myers CL. Comparison of profile similarity measures for genetic interaction networks. *PloS One.* 2013; 8(7):e68664. <https://doi.org/10.1371/journal.pone.0068664> PMID: 23874711
22. Brewster JL, Gustin MC. Hog1: 20 years of discovery and impact. *Sci Signal.* 2014 Sep 16; 7(343):re7. <https://doi.org/10.1126/scisignal.2005458> PMID: 25227612
23. Marques JM, Rodrigues RJ, de Magalhães-Sant'ana AC, Gonçalves T. Saccharomyces cerevisiae Hog1 protein phosphorylation upon exposure to bacterial endotoxin. *J Biol Chem.* 2006 Aug 25; 281(34):24687–94. <https://doi.org/10.1074/jbc.M603753200> PMID: 16790423
24. Lawrence CL, Botting CH, Antrobus R, Coote PJ. Evidence of a New Role for the High-Osmolarity Glycerol Mitogen-Activated Protein Kinase Pathway in Yeast: Regulating Adaptation to Citric Acid Stress. *Mol Cell Biol.* 2004 Apr 15; 24(8):3307–23. <https://doi.org/10.1128/MCB.24.8.3307-3323.2004> PMID: 15060153
25. Lillie SH, Brown SS. Immunofluorescence localization of the unconventional myosin, Myo2p, and the putative kinesin-related protein, Smy1p, to the same regions of polarized growth in Saccharomyces cerevisiae. *J Cell Biol.* 1994 May; 125(4):825–42. PMID: 8188749
26. Chowdhury S, Smith KW, Gustin MC. Osmotic stress and the yeast cytoskeleton: phenotype-specific suppression of an actin mutation. *J Cell Biol.* 1992 Aug; 118(3):561–71. PMID: 1639843
27. Nurse P. Universal control mechanism regulating onset of M-phase. *Nature.* 1990 Apr 5; 344(6266):503–8. <https://doi.org/10.1038/344503a0> PMID: 2138713
28. Collins K, Jacks T, Pavletich NP. The cell cycle and cancer. *Proc Natl Acad Sci U S A.* 1997 Apr 1; 94(7):2776–8. PMID: 9096291
29. Visconti R, Della Monica R, Grieco D. Cell cycle checkpoint in cancer: a therapeutically targetable double-edged sword. *J Exp Clin Cancer Res CR.* 2016 Sep 27; 35(1):153. <https://doi.org/10.1186/s13046-016-0433-9> PMID: 27670139
30. Denning DP, Hirose T. Anti-tubulins DEpendably induce apoptosis. *Nat Cell Biol.* 2014 Aug; 16(8):741–3. <https://doi.org/10.1038/ncb3012> PMID: 25082198
31. Jackson JR, Patrick DR, Dar MM, Huang PS. Targeted anti-mitotic therapies: can we improve on tubulin agents? *Nat Rev Cancer.* 2007 Feb; 7(2):107–17. <https://doi.org/10.1038/nrc2049> PMID: 17251917
32. La Regina G, Bai R, Coluccia A, Famigliani V, Pelliccia S, Passacantilli S, et al. New pyrrole derivatives with potent tubulin polymerization inhibiting activity as anticancer agents including hedgehog-dependent cancer. *J Med Chem.* 2014 Aug 14; 57(15):6531–52. <https://doi.org/10.1021/jm500561a> PMID: 25025991
33. Lu Y, Chen J, Xiao M, Li W, Miller DD. An overview of tubulin inhibitors that interact with the colchicine binding site. *Pharm Res.* 2012 Nov; 29(11):2943–71. <https://doi.org/10.1007/s11095-012-0828-z> PMID: 22814904
34. Baba T, Ara T, Hasegawa M, Takai Y, Okumura Y, Baba M, et al. Construction of Escherichia coli K-12 in-frame, single-gene knockout mutants: the Keio collection. *Mol Syst Biol.* 2006; 2:2006.0008.
35. Kim D-U, Hayles J, Kim D, Wood V, Park H-O, Won M, et al. Analysis of a genome-wide set of gene deletions in the fission yeast Schizosaccharomyces pombe. *Nat Biotechnol.* 2010 Jun; 28(6):617–23. <https://doi.org/10.1038/nbt.1628> PMID: 20473289

36. Kapitzky L, Beltrao P, Berens TJ, Gassner N, Zhou C, Wüster A, et al. Cross-species chemogenomic profiling reveals evolutionarily conserved drug mode of action. *Mol Syst Biol* [Internet]. 2010 Dec 21 [cited 2017 Feb 28]; 6. Available from: <http://msb.embopress.org/cgi/doi/10.1038/msb.2010.107>
37. French S, Mangat C, Bharat A, Côté J-P, Mori H, Brown ED. A robust platform for chemical genomics in bacterial systems. *Mol Biol Cell*. 2016 Mar 15; 27(6):1015–25. <https://doi.org/10.1091/mbc.E15-08-0573> PMID: 26792836
38. Babu M, Arnold R, Bundalovic-Torma C, Gagarinova A, Wong KS, Kumar A, et al. Quantitative genome-wide genetic interaction screens reveal global epistatic relationships of protein complexes in *Escherichia coli*. *PLoS Genet*. 2014 Feb; 10(2):e1004120. <https://doi.org/10.1371/journal.pgen.1004120> PMID: 24586182
39. Roguev A, Bandyopadhyay S, Zofall M, Zhang K, Fischer T, Collins SR, et al. Conservation and rewiring of functional modules revealed by an epistasis map in fission yeast. *Science*. 2008 Oct 17; 322(5900):405–10. <https://doi.org/10.1126/science.1162609> PMID: 18818364
40. Frost A, Elgort MG, Brandman O, Ives C, Collins SR, Miller-Vedam L, et al. Functional repurposing revealed by comparing *S. pombe* and *S. cerevisiae* genetic interactions. *Cell*. 2012 Jun 8; 149(6):1339–52. <https://doi.org/10.1016/j.cell.2012.04.028> PMID: 22682253
41. Ryan CJ, Roguev A, Patrick K, Xu J, Jahari H, Tong Z, et al. Hierarchical modularity and the evolution of genetic interactomes across species. *Mol Cell*. 2012 Jun 8; 46(5):691–704. <https://doi.org/10.1016/j.molcel.2012.05.028> PMID: 22681890
42. Estoppey D, Hewett JW, Guy CT, Harrington E, Thomas JR, Schirle M, et al. Identification of a novel NAMPT inhibitor by CRISPR/Cas9 chemogenomic profiling in mammalian cells. *Sci Rep*. 2017 Feb 16; 7:42728. <https://doi.org/10.1038/srep42728> PMID: 28205648
43. Deans RM, Morgens DW, Ökesli A, Pillay S, Horlbeck MA, Kampmann M, et al. Parallel shRNA and CRISPR-Cas9 screens enable antiviral drug target identification. *Nat Chem Biol*. 2016 May; 12(5):361–6. <https://doi.org/10.1038/nchembio.2050> PMID: 27018887
44. Blomen VA, Májek P, Jae LT, Bigenzahn JW, Nieuwenhuis J, Staring J, et al. Gene essentiality and synthetic lethality in haploid human cells. *Science*. 2015 Nov 27; 350(6264):1092–6. <https://doi.org/10.1126/science.aac7557> PMID: 26472760
45. Wang T, Birsoy K, Hughes NW, Krupczak KM, Post Y, Wei JJ, et al. Identification and characterization of essential genes in the human genome. *Science*. 2015 Nov 27; 350(6264):1096–101. <https://doi.org/10.1126/science.aac7041> PMID: 26472758
46. Hart T, Chandrashekar M, Aregger M, Steinhart Z, Brown KR, MacLeod G, et al. High-Resolution CRISPR Screens Reveal Fitness Genes and Genotype-Specific Cancer Liabilities. *Cell*. 2015 Dec 3; 163(6):1515–26. <https://doi.org/10.1016/j.cell.2015.11.015> PMID: 26627737
47. Horlbeck MA, Xu A, Wang M, Bennett NK, Park CY, Bogdanoff D, et al. Mapping the Genetic Landscape of Human Cells. *Cell*. 2018 Aug 9; 174(4):953–967.e22. <https://doi.org/10.1016/j.cell.2018.06.010> PMID: 30033366
48. Huber W, Carey VJ, Gentleman R, Anders S, Carlson M, Carvalho BS, et al. Orchestrating high-throughput genomic analysis with Bioconductor. *Nat Methods*. 2015 Jan 29; 12(2):115–21. <https://doi.org/10.1038/nmeth.3252> PMID: 25633503
49. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J R Stat Soc Ser B Methodol*. 1995; 57(1):289–300.
50. Hinselmann G, Rosenbaum L, Jahn A, Fechner N, Zell A. jCompoundMapper: An open source Java library and command-line tool for chemical fingerprints. *J Cheminformatics*. 2011 Jan 10; 3(1):3.
51. Safizadeh H, Simpkins SW, Nelson J, Myers CL. Improving prediction of compound function from chemical structure using chemical-genetic networks. *bioRxiv* [Internet]. 2017 Mar 1; Available from: <http://biorxiv.org/content/early/2017/03/01/112698.abstract>
52. Nelson J, Simpkins SW, Safizadeh H, Li SC, Piotrowski JS, Hirano H, et al. MOSAIC: a chemical-genetic interaction data repository and web resource for exploring chemical modes of action. *Bioinformatics*. 2017; <https://doi.org/10.1093/bioinformatics/btx732> PMID: 29206899