

1 **Supplementary Information:**

2  
3 **Proteome-scale detection of small-molecule protein interactions using data from gene**  
4 **expression perturbations**

5  
6  
7 Nicolas A. Pabon<sup>1,†</sup>, Yan Xia<sup>2,†</sup>, Sam Estabrooks<sup>3</sup>, Zhaofeng Ye<sup>4</sup>, Amanda K. Herbrand<sup>5</sup>,  
8 Evelyn Süß<sup>5</sup>, Ricardo M. Biondi<sup>5</sup>, Victoria A. Assimon<sup>6</sup>, Jason E. Gestwicki<sup>6</sup>, Jeffrey L.  
9 Brodsky<sup>3</sup>, Carlos J. Camacho<sup>1,\*</sup>, and Ziv Bar-Joseph<sup>2,\*</sup>

10  
11  
12 <sup>1</sup> Department of Computational and Systems Biology, University of Pittsburgh , Pittsburgh, PA 15213

13 <sup>2</sup> Machine Learning Department, School of Computer Science, Carnegie Mellon University, 15213

14 <sup>3</sup> Department of Biological Sciences, University of Pittsburgh, Pittsburgh, PA 15213

15 <sup>4</sup> School of Medicine, Tsinghua University, Beijing, China 100084

16 <sup>5</sup> Department of Internal Medicine I, Universitätsklinikum Frankfurt, 60590 Frankfurt, Germany

17 <sup>6</sup> Department of Pharmaceutical Chemistry, University of California, San Francisco, San Francisco, CA  
18 94158

19  
20  
21 <sup>†</sup> These two authors contributed equally

22 <sup>\*</sup> To whom correspondence should be addressed

23  
24 **Availability**

25 Supplementary Methods, Results, Data and Matlab code are available at the supporting website  
26 <http://sb.cs.cmu.edu/Target2/>.

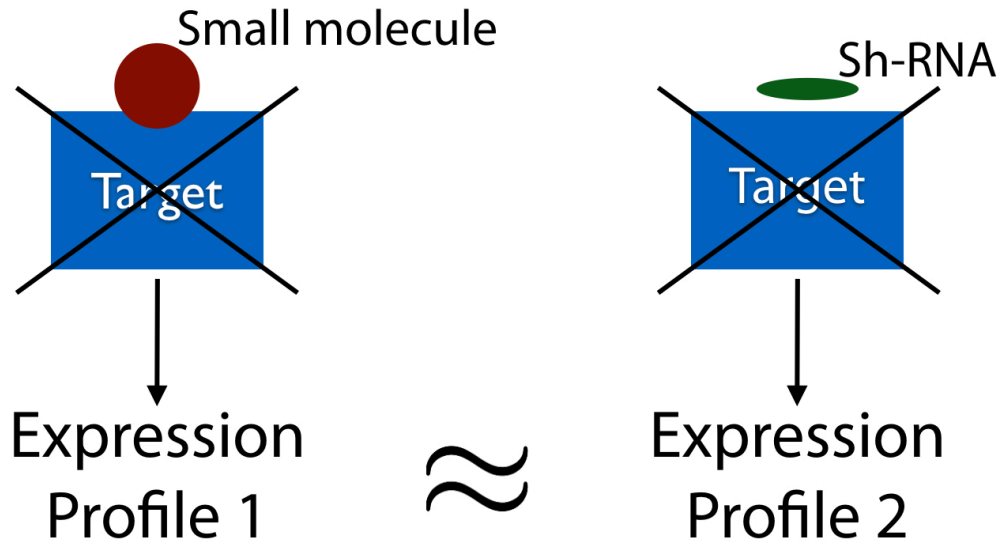
27  
28 **Contact**

29 Ziv Bar-Joseph: [zivbj@cs.cmu.edu](mailto:zivbj@cs.cmu.edu)

30 Carlos J. Camacho: [ccamacho@pitt.edu](mailto:ccamacho@pitt.edu)

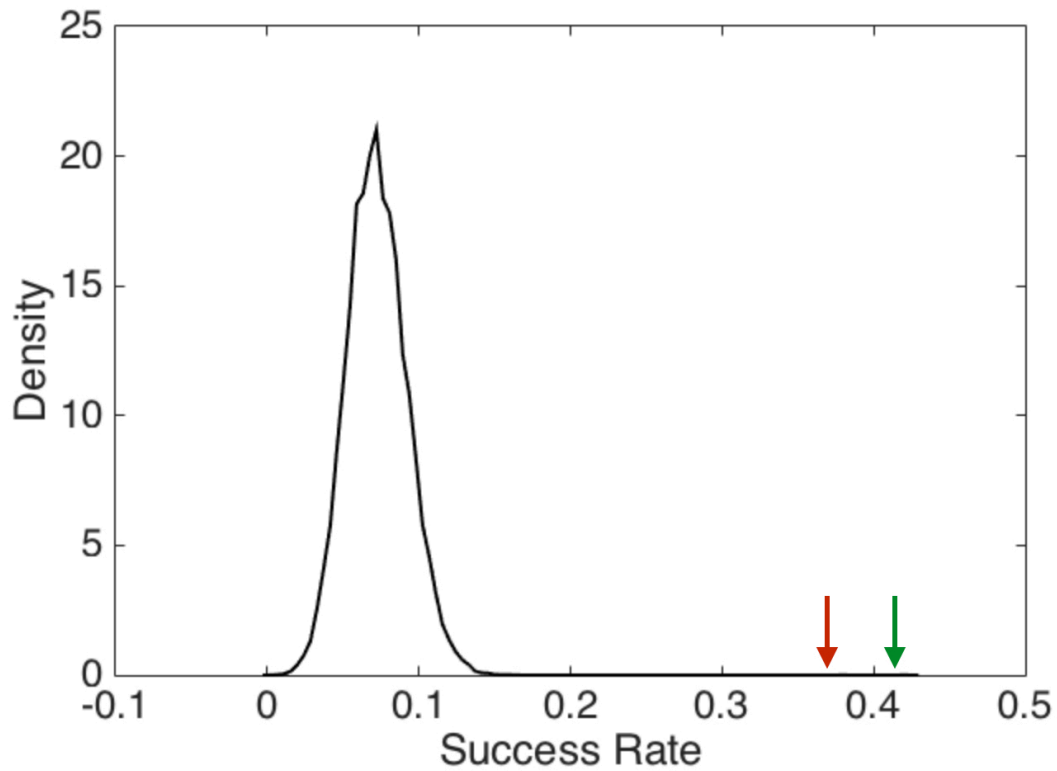
31  
32  
33  
34  
35  
36

1 **Supplementary Figure 1. Hypothesis of our genomic target prediction**  
2 **analysis: the effect of inhibiting a target by a small molecule is similar to**  
3 **that of knocking down the same target with sh-RNA.** These effects can be  
4 manifested by the resulting expression profiles as well as the activation/inhibition  
5 of specific pathways. As we discuss in detail in the Online Methods, this  
6 assumption leads to the construction of several different types of features that we  
7 use to evaluate the similarity of the effects of the two treatments within the LINCS  
8 dataset.  
9



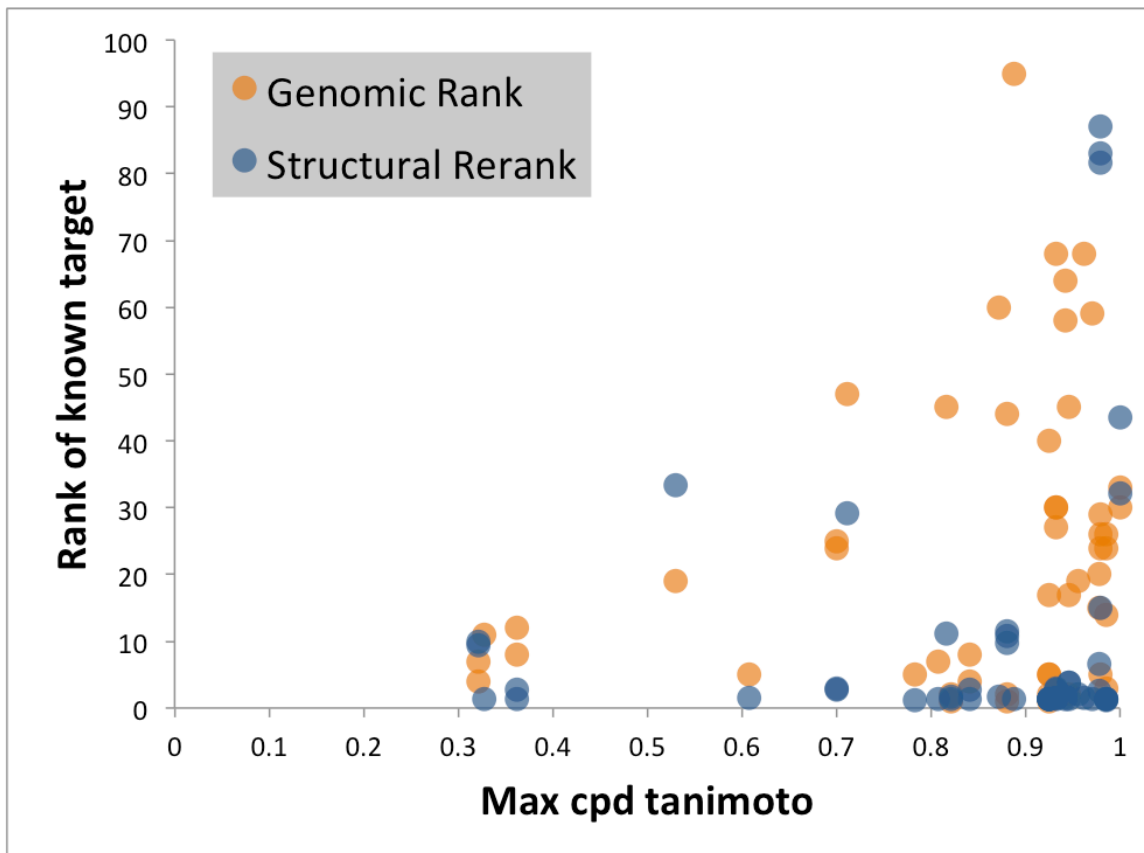
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

1 **Supplementary Figure 2. Comparing the random forest approaches with a**  
2 **random classifier** for predicting known targets of the 152 drugs in the validation  
3 set. The red arrow indicates the success rate of on-the-fly random forest and the  
4 green arrow represents the two-level random forest.  
5



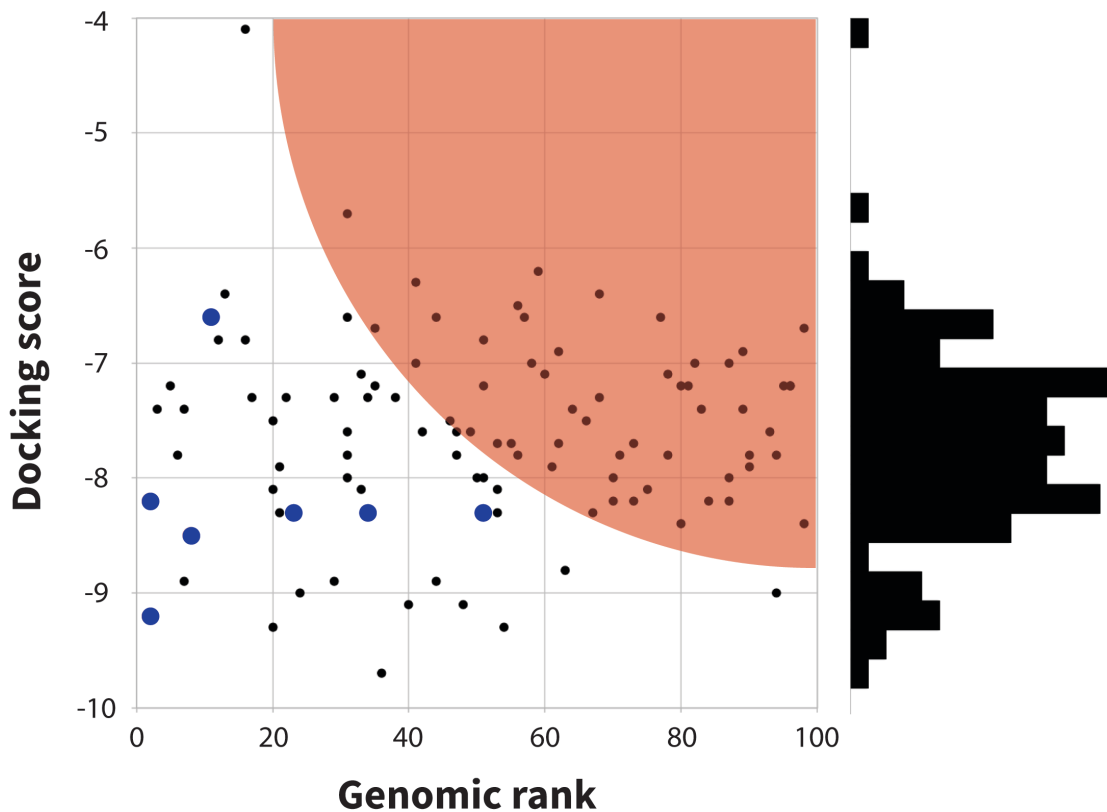
6  
7  
8  
9  
10

1 **Supplementary Figure 3. Correlation of target prediction accuracy and**  
2 **“structural uniqueness” of the query compound with respect to the training**  
3 **compounds.** Each point in the plot represents one of the 53 compounds in our  
4 enrichment analysis. The structural uniqueness of a compound (x-axis) is defined  
5 as its maximum Tanimoto distance to any of the training compounds. The  
6 predicted ranking of the known target for each compound is shown on the y-axis.  
7 Orange and blue points represent the ranking pre- and post- structural filtering.  
8



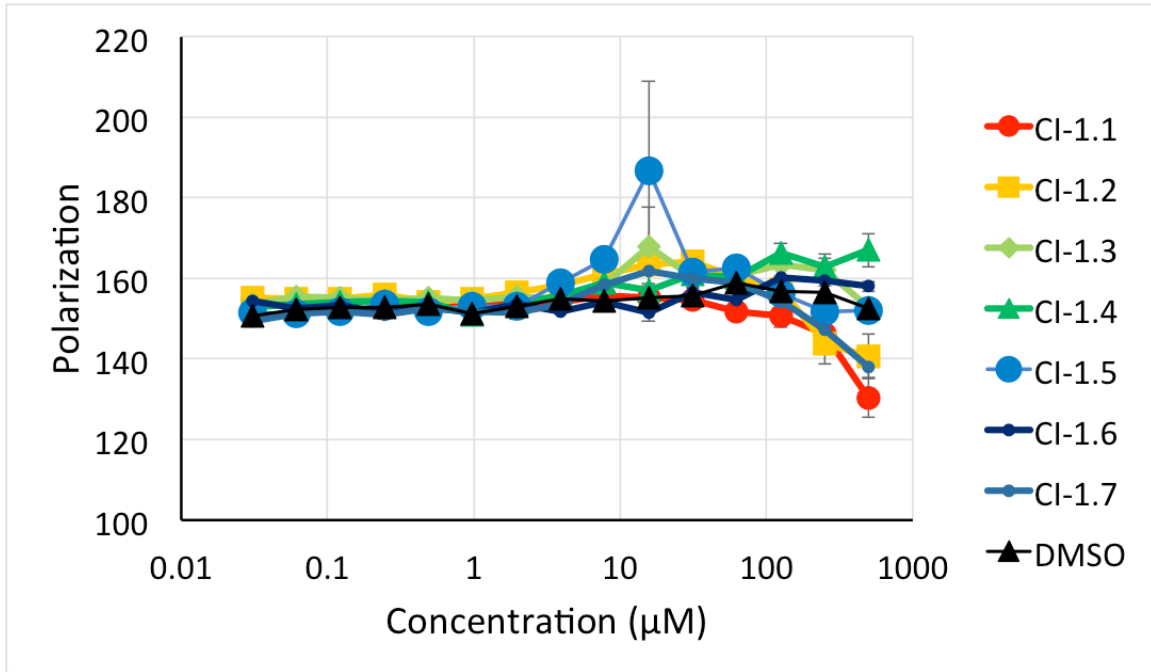
9  
10  
11

1 **Supplementary Figure 4. Result of the target-centric screen against CHIP.**  
2 The plot on the left shows the 104 compounds predicted by random forest to bind  
3 CHIP, plotted according to the rank of CHIP in their predicted targets list (x -  
4 axis), vs. their CHIP docking score (y - axis). The shaded red area of the plot  
5 represents compounds that were filtered out of analysis due to low rank/score.  
6 The blue dots represent the compounds that were purchased for experimental  
7 validation. The histogram on the right shows the distribution of compounds by  
8 docking score.  
9



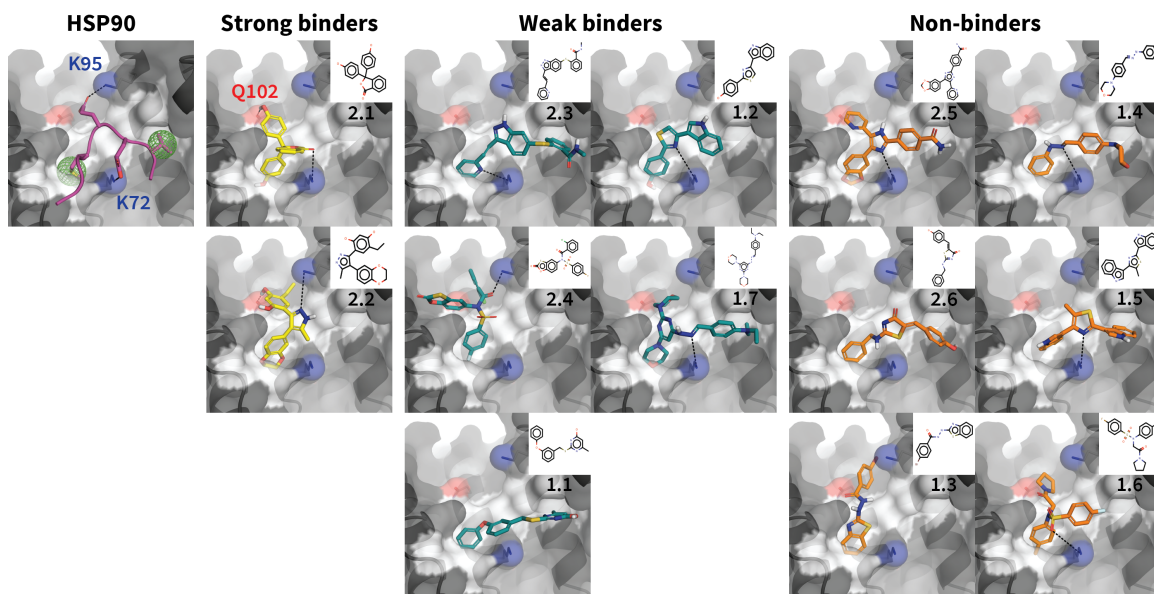
10  
11

1 **Supplementary Figure 5.** Disruption of CHIP binding to chaperone peptide  
2 measured by fluorescence polarization. Results are the average and standard  
3 error of the mean of two experiments each performed in triplicate.  
4



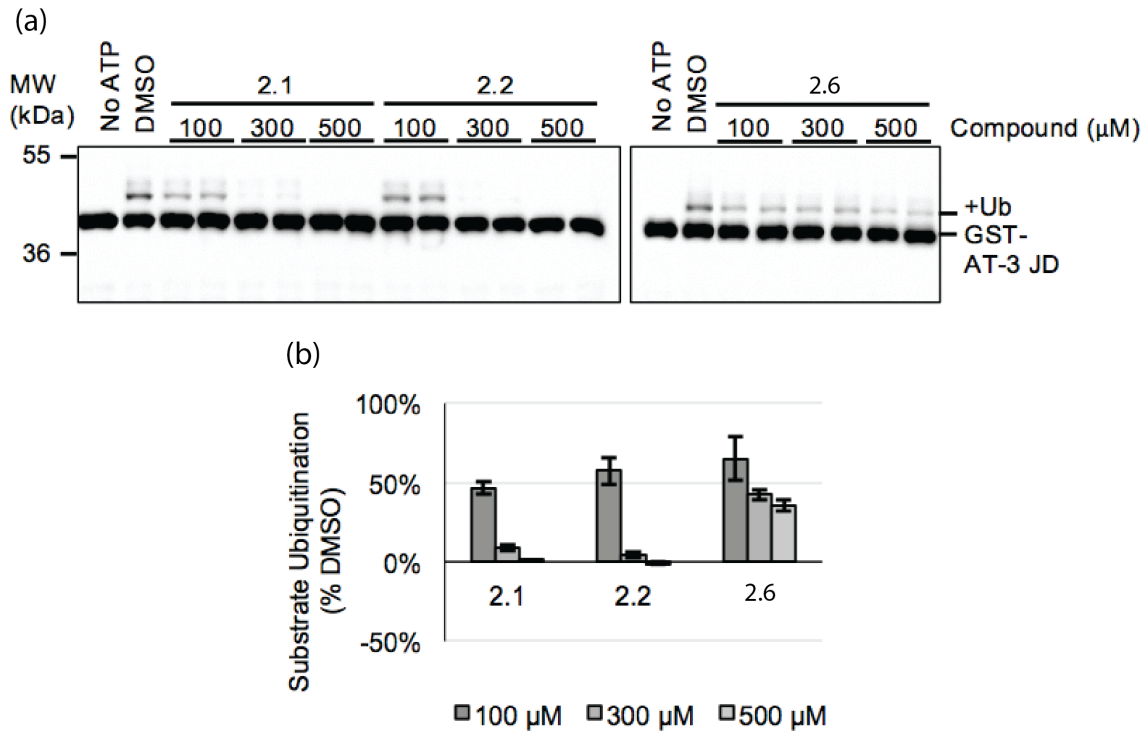
5  
6  
7

1 **Supplementary Figure 6. Comparison of virtual screens against CHIP.**  
2 **HSP90** shows structure of the CHIP (grey) - HSP90 (magenta) interface (PDB  
3 ID: 2C2L<sup>1</sup>), indicating the hydrophobic (green spheres) and polar contact (blue  
4 surface / dashed lines) pharmacophores used to screen the ZINC database.  
5 **Strong binders** show predicted binding modes for compounds 2.1 and 2.2 from  
6 the LINCS screen, which showed the strongest FP signal and robust inhibition of  
7 CHIP ligases activity. Interestingly, 2.1 and 2.2 are the only predicted hits to  
8 make a novel hydrogen bond to CHIP residue Q102, a contact whose importance  
9 is not obvious from the cocrystal structure show predicted binding modes. **Weak**  
10 **binders** show predicted binding modes for compounds 2.3 and 2.4 from the  
11 LINCS screen, and compounds 1.1, 1.2, and 1.7 from the ZINC screen, which  
12 showed modest FP signal. **Non-binders** show predicted binding modes for non-  
13 binding LINCS compounds 2.5 and 2.6, and non-binding ZINC compounds 1.3 –  
14 1.6.  
15



16  
17  
18

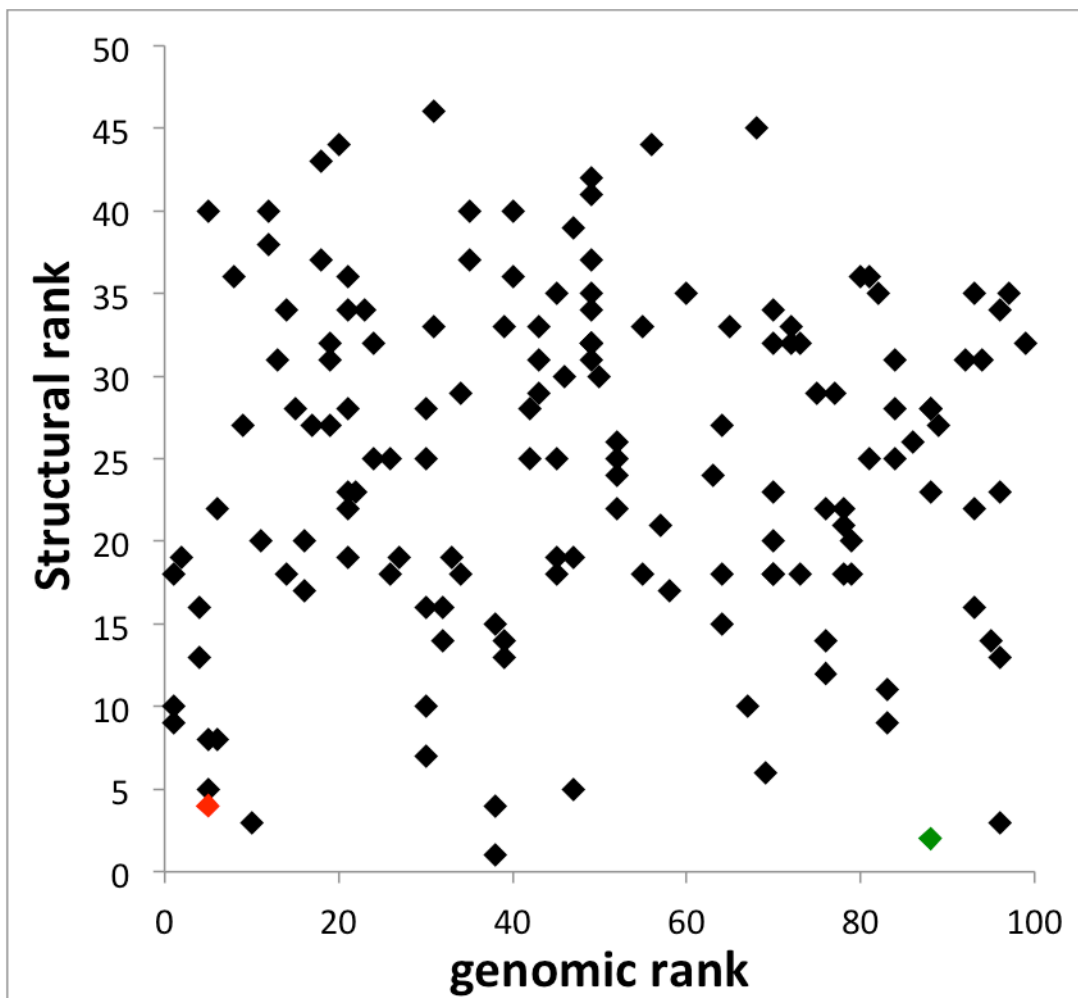
1 **Supplementary Figure 7. Predicted CHIP Inhibitors Prevent Ubiquitination of an**  
 2 **Alternate Substrate.** (A) Anti-GST western blot showing AT-3 JD substrate  
 3 ubiquitination by CHIP in reactions treated with compounds. (B) Quantification of all  
 4 reactions as in A treated with up to 500  $\mu$ M compound 2.1, 2.2, or 2.6, normalized to  
 5 ubiquitination by a DMSO treated control (all compounds: N=4).  
 6



7  
8  
9  
10

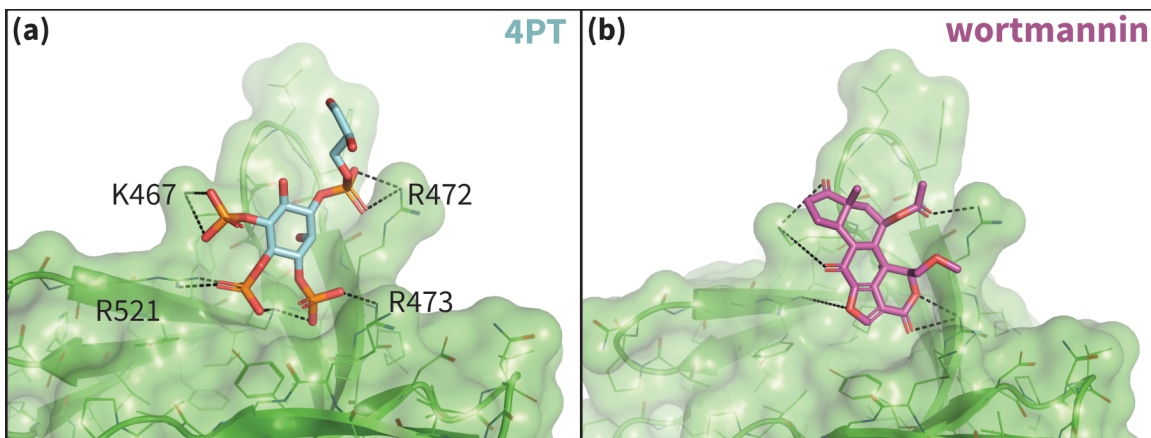


1 **Supplementary Figure 8.** Pipeline-predicted targets for the drug wortmannin.  
2 Each point on the plot represents one structural model of one potential target  
3 predicted in the top 100 for wortmannin by the random forest regressor. The  
4 random forest ranking for each target (x axis) is plotted against the docking score  
5 ranking (y axis). The red dot indicates the ranking of the known target PIK3CA.  
6 The green dot indicates the ranking for the previously unknown target, PDPK1.  
7



8  
9  
10

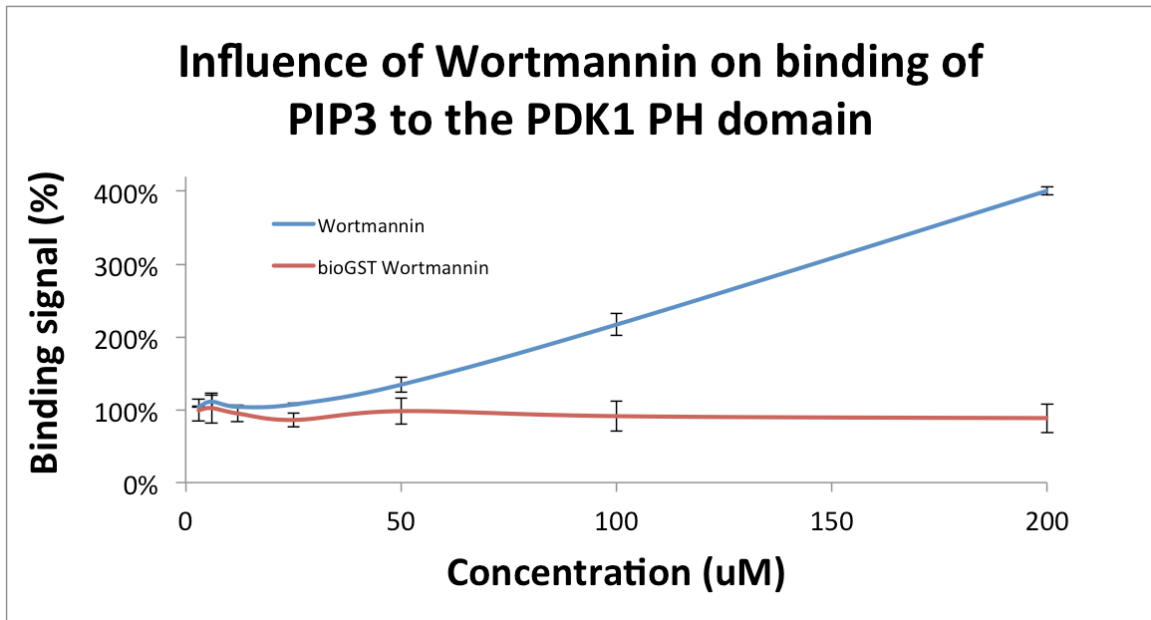
1 **Supplementary Figure 9.** Docking model of wortmannin bound to the PH  
2 domain of PDK1. (a) Cocrystal structure (PDB ID: 1W1G<sup>2</sup>) of the PH domain of  
3 PDK1 bound to the 4PT ligand which mimics the head group of it's natural ligand  
4 PIP3. Dashed lines indicate key polar interactions. (b) Docking model of  
5 wortmannin bound to the PDPK1 PH domain, which captures many of the same  
6 polar interactions seen in the cocystal.  
7



8  
9  
10

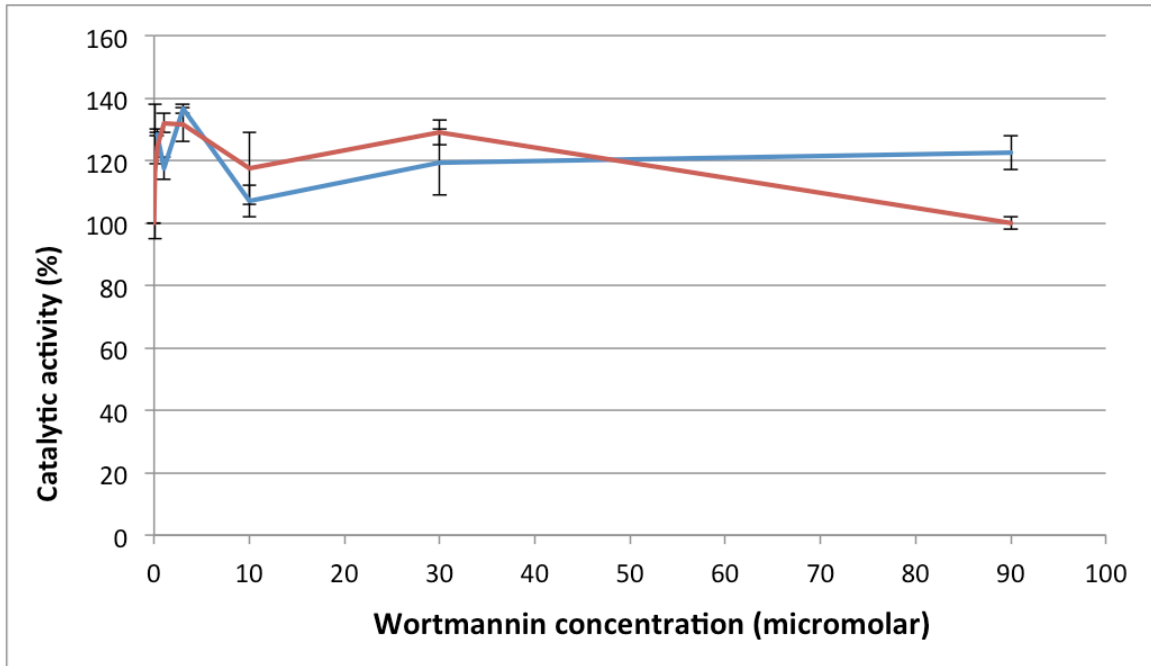
1  
2  
3  
4

**Supplementary Figure 10.** Alphascreen PDK1-PIP3 interaction-displacement assay results for increasing concentrations of wortmannin.



5  
6  
7  
8

1 **Supplementary Figure 11.** Effect of wortmannin on the in-vitro phosphorylation  
2 of the substrate T308tide by the isolated catalytic domain of PDK1. The two lines  
3 are from two replicates of the activity assay, with error bars representing the  
4 standard error on the mean from two parallel runs for each replicate.  
5



6  
7  
8  
9  
10  
11

1 **Supplementary Table 1. Results of testing our random forest classifier on**  
2 **the 123 FDA approved drugs profiled in 4-6 LINCS cell lines**, after having  
3 trained our model on the 29 FDA approved drugs profiled in all 7 LINCS cell  
4 lines. The rank of the highest-ranking known target for each compound is listed  
5 next to their LINCS ID. We achieve top-100 predictions for 32 drugs, a 26%  
6 success rate.  
7  
8 **\*\* Table provided as a separate tabular data file *supp\_table\_1.txt***  
9  
10

1 **Supplementary Table 2. Structural enrichment of random forest predictions**  
2 **for validation hits and comparison with existing methods.** Our 63 `hits' are  
3 listed with their LINCS ID and the number of top-100 predicted targets that had  
4 structures available in the PDB. The ranking of the known targets are shown after  
5 our genomic random forest target prediction (GEN), and after our structural re-  
6 ranking (STR), along with the percentile rankings produced by alternative target  
7 prediction methods HTDocking (HTD) and PharmMapper (PHM). STR, HTD, and  
8 PHM values of 100 indicate that the structure of the known target either is not  
9 known or was not included in the set of potential targets used by the method.

10

11

12 **\*\* Table provided as a separate tabular data file *supp\_table\_2.txt***

13

14

15

16

17

18

1 **Supplementary Table 3. Predicted CHIP-targeting compounds out of 104**  
2 **candidate molecules.** 'CHIP RANK' indicates the ranking of CHIP in the  
3 random-forest predicted list of potential targets for each compound. 'CPD RANK'  
4 indicates the structure-based ranking of the compound after docking of all 104  
5 candidate compounds to the HSP90 binding site on the CHIP-TPR domain.  
6

Cpd #	NAME	ID	CHIP RANK	CPD RANK
2.1	phenolphthalein	BRD_K19227686	2	22
2.2	HSP90_inhibitor	BRD_K65503129	2	4
2.3	axitinib	BRD_K29905972	8	13
2.4	BRD_K59556282	BRD_K59556282	11	92
2.5	SB_431542	BRD_K67298865	34	17
2.6	MW_STK33_2B	BRD_K78930611	51	16

7  
8  
9  
10

1 **Supplementary Table 4.** Symbols and notations

<b>Symbol</b>	<b>Meaning</b>
$d$	Index for a drug
$c$	Index for a cell line
$g$	Index for a gene
$N_D$	Total number of genes
$N_C$	Total number of cell lines
$C_d$	The set of cell line indeces for drug $d$
$P_d$	The set of protein target indeces for drug $d$
$G_c$	The set of knockdown gene indeces for cell line $c$
$T_d$	The intersection of knockdown gene indeces $G_c$ for all cell lines in $C_d$
$N_{dc}$	Number of experiments for applying drug $d$ to cell line $c$
$N_{gc}$	Number of experiments for knocking down gene $g$ in cell line $c$
$N_g$	Neighbors, or protein-protein interaction partners, of gene $g$
$\Delta$	Drug-response data
$\Gamma$	Gene-knockdown data
$\Psi$	Control data
$\Omega$	Full feature data
$X_d$	Training data derived from drug $d$
$y_d$	Training label derived from drug $d$
$V_d$	Negative (non-target) genes for drug $d$

2  
3  
4  
5



1 **Supplementary Table 5.** Summary of constructed feature sets. Note that  
 2 different feature sets can have different dimensions (some contain values for  
 3 each of the cell lines, etc...). The exact dimension and content of each feature  
 4 set is discussed in the text.  
 5

<b>Feature Name</b>	<b>Symbol</b>	<b>Meaning</b>
Correlation	$f_{cor}$	Correlation between a drug treatment experiment and a gene knockdown experiment
Cell Selection	$f_{CS}$	Correlation between a drug treatment experiment and the control experiment for the cell line
PPI Correlation	$f_{PC}$	Fraction of the known binding partners of a gene in the top $X$ correlated knockdown experiments
PPI Expression	$f_{PE}$	The average or the max (absolute value) expression for the known binding partners of a gene

6

1 **Supplementary Table 6.** The number of drugs profiled for different number of  
2 cell lines included in the validation dataset. While several drugs were profiled in  
3 at least four of these cell lines (152), only 29 were profiled in all seven cell lines.  
4

# Cells	7	6	5	4	Total <sup>5</sup>
# Drugs	29	30	42	51	152

1 **Supplementary Table 7.** Seven cell lines were included in the validation dataset.  
2 The number of drugs, knockdown genes, and control experiment are shown. For  
3 a given cell line, we only include drugs that have their target knockdown  
4 experiments available in that cell line.  
5

<b>Cell Line</b>	<b>Drugs</b>	<b>Knockdowns</b>	<b>Controls</b>
A549	188	11947	52
MCF7	180	12031	54
VCAP	175	13225	56
HA1E	172	11968	53
A375	143	11696	58
HCC515	129	7828	52
HT19	96	10185	52

1 **Supplementary Table 8. The cellular localization of successful and**  
 2 **unsuccessful drug targets enriched by gene ontology.**  
 3

	<b>Cellular Component</b>	<b>p-value</b> <sup>4</sup> <sub>5</sub>
<b>Successful Targets</b>	proteasome core complex	7.81E-37
	proteasome core	1.10E-28
	proteasome alpha-subunit	5.68E-18
	cytosol	7.53E-12
	protein complex	1.88E-11
<b>Failed Targets</b>	transmembrane transporter complex	7.77E-15
	sodium-exchanging ATPase complex	4.42E-14
	cation-transporting ATPase complex	8.74E-13
	plasma membrane part	2.19E-11
	chloride channel complex	2.33E-09

## 1 Supplementary Results

### 3 **Gene ontology analysis of protein targets**

4 While the success rate of our Random Forest genomic analysis is promising, there are  
5 still several drugs for which we fail to correctly identify the target. We attempted to  
6 determine if the genomic data we used is more appropriate to specific drug / protein  
7 characteristics. By characterizing the set of drugs and / or proteins for which we expect  
8 the method to be more accurate we improve the ability of experimentalists to use our  
9 methods when studying one of these molecules.

10  
11 We divided the 152 drugs in our training data into “successful” predictions (the 63 drugs  
12 for which the correct target was ranked in the top 100), and “unsuccessful” predictions.  
13 We also divided the known targets into those that were correctly predicted and those  
14 that were not. We considered several different ways to characterize small molecules  
15 including molecular weight, solubility, and hydrophobicity, but none of these seemed to  
16 significantly correlate with our “successful” and “unsuccessful” classifications. Next, we  
17 used gene ontology (Online Methods) to test for enrichment of “successful” and  
18 “unsuccessful” targets. Interestingly, we found that “successful” targets were  
19 significantly associated with intracellular categories, while the “unsuccessful” targets  
20 were mostly associated with transmembrane and extracellular categories  
21 (Supplementary Table 8).

22  
23 Based on this result we further incorporated cellular component as a feature in our two-  
24 level random forest. We encode this feature by assigning 1 to the intracellular genes and  
25 -1 to the extracellular ones (see Online Methods for detail). We ran the two-level  
26 random forest with this additional feature included and demonstrated that the cellular  
27 component increases the number of top 100 genes to 66 and top 50 genes to 55.

## 1 Supplementary Methods

### 3 **Extracting experiments from LINCS**

4 After determining the subsets of small molecules and cell lines, we obtained the  
5 associated experiment identifiers known as “distil IDs” from LINCS meta-  
6 information. We included only the reproducible distil IDs known as “Gold” IDs.  
7 We then extracted the corresponding signature values from LINCS using the  
8 L1000 Analysis Tools (l1ktools)<sup>1</sup>. We only extracted the signature values of the  
9 978 “landmark” genes because their expression was directly measured, whereas  
10 the values of other genes were imputed from the data of these landmark genes.

### 12 Drug response experiments

13 There exist multiple experiments (distil IDs) corresponding to a combination of  
14 drug  $d$  and cell line  $c$  (applying drug  $d$  to cell line  $c$ ). Denote the  $N_{dc}$  as the  
15 number of experiments for the combination  $d,c$ . We extracted a matrix of  
16 signature values of size  $978 \times N_{dc}$  (number of landmark genes  $\times$  number of  
17 experiments) per combination. We next took the median of signature values  
18 across different experiments, and obtained a  $978 \times 1$  signature vector per  
19 combination. The overall drug-response data  $\Delta$ , therefore, is implemented as a  
20 MATLAB structure with  $D = 152$  entries, each containing the following fields.

```
21  
22 |         name:  $PertID_d$  (string)  
23 |         cells:  $Cells_{C_d}$  ( $|C_d| \times 1$  string array)  
24 |         signature:  $\Delta_{d..}$  ( $978 \times |C_d|$ )  
25 |
```

26 | where  $PertID_d$  is the unique internal identifier of a small molecule  $d$  in LINCS.  
27 |  $\Delta_{d..}$  contains the expression values of drug  $d$  across  $C_d$  different cell lines. The  
28 |  $Cells_{C_d}$  field contains cell line names corresponding to the column of  $\Delta_{d..}$ .

### 30 Gene knockdown (KD) experiments

31 We follow a similar protocol to extract the signature values of gene knockdown  
32 experiments. Denote  $N_{gc}$  as the number of experiments for the combination of  
33 gene  $g$  and cell line  $c$  (knocking down gene  $g$  in cell line  $c$ ). Then, for each  
34 combination of  $g$  and  $c$  we extracted signature values of size  $978 \times N_{gc}$ . After  
35 taking the medians across different experiments, we obtain a  $978 \times 1$  vector per  
36 combination. The overall gene knockdown data  $\Gamma$  has  $C = 7$  entries and each  
37 entry contains the following fields:

```
38  
39 |         name:  $Cells_c$  (string)  
40 |         genes:  $Symbols_{G_c}$  ( $|G_c| \times 1$  string array)  
41 |         signature:  $\Gamma_{c..}$  ( $978 \times |G_c|$ )  
42 |
```

43 | where  $Cells_c$  is the name of the cell line indexed by  $c$ .  $\Gamma_{c..}$  contains the signature  
44 | values of the knockdown of genes in cell line  $c$ . The  $Symbols_{G_c}$  field is a subset of

---

<sup>1</sup> <https://github.com/cmap/l1ktools>

1 | gene symbols corresponding to the column identifiers of  $\Gamma_{c..}$  under the HGNC  
2 | naming scheme.

3

#### 4 | Control experiments

5 | We also extracted the signatures of control experiments. The signature values for  
6 | each cell line were extracted and we obtained a  $978 \times 1$  vector after taking the  
7 | medians. We denote the overall control experiment data as  $\Psi$ .  $\Psi$  is of size  $978 \times$   
8 |  $C$  and implemented with the following format:

9

10 |           **name:**  $Cells_c$  (string)

11 |           **control:**  $\Psi_c$  ( $978 \times 1$ )

12

13 | where  $\Psi_c$  is the signature column vector for a cell line  $c$ .

14

15

#### 16 | **Extracting and integrating features from different data sources**

17

##### 18 | Correlation feature

19 | The correlation feature, denoted as  $f_{cor}$ , is constructed as follows:

20

21 | - For each drug  $d$  in  $\Delta$  ( $\Delta_{d..}$ ):

22

23 |     - Denote  $T_d$  as the intersection of gene symbol indices for cells in  $C_d$ :

24

$$T_d = \bigcup_{c \in C_d} G_c$$

25

26 | - Obtain the knockdown signature values of  $T_d$  from  $\Gamma$ . Denote this data  
27 | matrix as  $\Gamma_{C_d \cdot T_d}$ , which is of size  $|C_d| \times 978 \times |T_d|$ , where for each cell line in  
28 |  $C_d$  there is a signature matrix of size  $978 \times |T_d|$ .

29

30

31 | - Compute the Pearson's correlation between  $\Delta_{d..}$  ( $978 \times |C_d|$ ) and  $\Gamma_{C_d \cdot T_d}$  ( $|C_d|$   
32 |  $\times 978 \times |T_d|$ ). Specifically, for each cell line  $c \in C_d$ , we compute the  
33 | correlation between  $\Delta_{d \cdot c}$  and  $\Gamma_{c \cdot T_d}$ , and obtain a correlation vector of size  $|T_d|$ .  
34 | This is the correlation between the responses of the cells to the drug  
35 | treatment and their response to the gene KD. Each entry in this vector is the  
36 | correlation of 978 landmark genes of the drug  $d$  in one cell line ( $\Delta_{d \cdot c}$ ) and a  
37 | knockdown of gene  $g$  in the same cell line ( $\Gamma_{c \cdot g}$ ). In other words, if we collect  
38 | these correlation vectors for all cell lines in  $C_d$  and denote the overall  
39 | correlation feature as  $f_{cor}$ :

40

$$f_{cor}(d, g, c) = \text{corr}(\Delta_{d \cdot c}, \Gamma_{c \cdot g}) \quad \forall g \in T_d$$

41

42

1 | The correlation feature for one drug  $d$ ,  $f_{cor}(d, \cdot)$ , has a dimension of  $|T_d| \times$   
2 |  $|C_d|$ .

#### 4 | Cell selection feature

5 | The cell selection feature, denoted as  $f_{CS}$ , is computed as follows:

6 |

7 | - For each drug  $d$  in  $\Delta$  ( $\Delta_{d\cdot}$ ):

8 |

9 |     - For each cell line  $c$  in  $C_d$ :

10 |

11 |         - Compute the correlation between  $\Delta_{d\cdot c}$  and  $\Psi_c$

12 |

$$f_{CS}(d, c) = corr(\Delta_{d\cdot c}, \Psi_c)$$

13 |

14 |  $f_{CS}(d, \cdot)$  produces a  $|C_d| \times 1$  vector, and each entry corresponds to the correlation  
15 | between the drug-response and control experiments for one cell line in  $C_d$ . This  
16 | feature is used to determine the relevance of the drug to the cell type being  
17 | studied.

18 |

#### 19 | PPI correlation score:

20 | The PPI correlation Score, denoted as  $f_{PC}$  is constructed as follows:

21 |

22 | - For each drug  $d$  in  $\Delta$  ( $\Delta_{d\cdot}$ ):

23 |

24 |     - Obtain  $T_d$ , as defined above.

25 |

26 |     - For each cell line  $c$  in  $C_d$ :

27 |

28 |         - Sort  $T_d$  in descending order using the correlation values  $f_{cor}(d, \cdot, c)$

29 |

30 |         - Denote the sorted gene symbol indices for cell line  $c$  as  $\sigma_c(T_d)$

31 |

32 |     - For each knockdown gene  $g$  in  $T_d$ :

33 |

34 |         - Obtain the set of neighbor gene symbol indices from the PPI  
35 | adjacency list, and denote it as  $N_g$ .

36 |

37 |         - Compute  $f_{PC}$  as:

38 |

$$f_{PC}(d, g, c) = \frac{|N_g \cap \sigma_c(T_d)_{1:100}|}{|N_g \cap \sigma_c(T_d)| + 50}$$

39 |

40 |  $f_{PC}(d, g, c)$  has the same dimension as  $f_{cor}$  ( $|T_d| \times |C_d|$ ). It reflects the fraction of  
41 | gene  $g$ 's binding partners that are more correlated with drug  $d$  in the context of  
42 | cell line  $c$ . We use 50 as the pseudo-count to penalize hub proteins, which have  
43 | substantially more neighbors than others.



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43

## PPI expression score

We compute two types of PPI expression scores, denoted as  $f_{PE_{max}}$  and  $f_{PE_{avg}}$ , as follows:

- For each drug  $d$  in  $\Delta$  ( $\Delta_{d..}$ ):
  - For each knockdown gene  $g$  in  $T_d$ :
    - Obtain  $N_g$ , as above (the list of neighbors, or interaction partners, of  $g$ )
    - For each cell line  $c$  in  $C_d$ :
      - Find the set of signature values for the neighbors of  $g$ ,  $\Delta_{d,N_g,c}$  (size  $|N_g| \times 1$ )
      - Compute the two PPI expression scores as:

$$f_{PE_{max}}(d, g, c) = \max(\Delta_{d,N_g,c})$$

$$f_{PE_{avg}}(d, g, c) = \text{avg}(\Delta_{d,N_g,c})$$

## Feature data structure

We combined the features for all drugs in a MATLAB structure  $\Omega$ .  $\Omega$  has  $D$  entries, and each entry  $\Omega^{(d)}$  has the following fields:

- name:**  $PertID_d$  (string)
- targets:**  $P_d$  (protein targets for  $d$ )
- cells:**  $Cells_{C_d}$  ( $|C_d| \times 1$  string array)
- genes:**  $T_d$  (common genes across  $G_c$ )
- correlation:**  $f_{cor}(d, \cdot)$  ( $|T_d| \times |C_d|$ )
- PPI correlation:**  $f_{PC}(d, \cdot)$  ( $|T_d| \times |C_d|$ )
- max PPI expression:**  $f_{PE_{max}}(d, \cdot)$  ( $|T_d| \times |C_d|$ )
- avg PPI expression:**  $f_{PE_{avg}}(d, \cdot)$  ( $|T_d| \times |C_d|$ )
- cell selection:**  $f_{CS}(d, \cdot)$  ( $|C_d| \times 1$ )

There are a total of  $D = 152$  drugs in  $\Omega$ , and the number of drugs with different values of  $|C_d|$  are summarized in Supplementary Table 6.

## **Subcellular Localization Assignment**

We obtained the cellular localization of genes from the Gene Ontology Consortium. The GO database provides web services to query genes in terms of their associated biological processes, cellular components and molecular

1 functions in a species-independent manner<sup>2</sup>. We further assign the locations as  
 2 either “intracellular” (inside of cell) or “extracellular” (outside of cell). The detailed  
 3 assignments are shown in Supplementary Table 8.

## 4 5 6 **Classification procedure**

### 7 8 Criterion of successful classification

9 Due to the intrinsic noise from the data, we define a successful classification for a  
 10 drug if any of its correct targets is enriched into the top  $K$  ranked genes, where  $K$   
 11 can be either 50 or 100.

### 12 13 Analysis of feature importance

14 The evaluation of single features was performed using the drugs that have been  
 15 applied on all seven cell lines. There are 29 of these drugs from  $\Omega$ . We sort  
 16 (descendingly) the common genes  $T_d$  for a drug  $d$  and cell line  $c$  using an  
 17 individual feature  $f(d, \cdot, c)$ , where  $f$  is either  $f_{cor}$  or  $f_{PC}$ . Denote  $\sigma_d(g, c)$  as the  
 18 ranking of a gene  $g \in T_d$  in the context of cell line  $c$ . Then, we define the overall  
 19 ranking of a gene,  $\sigma_d(g)$ , to be the best ranking across all seven cell lines:  
 20  $\sigma_d(g) = \min(\sigma_d(g, c))$  for  $c \in C_d$ .

### 21 22 Constructing training dataset

23 Next, we wish to learn and evaluate classifiers that predict drug targets using all  
 24 features from the feature dataset  $\Omega$ . We first construct a training data set (design  
 25 matrix  $X$  and its associated labels  $y$ ) from the feature dataset  $\Omega$ .

26  
27 For each drug  $d$  in  $\Omega$ , we select the rows corresponding to the targets in  $P_d$  from  
 28 the other feature matrices and concatenate them into a row vector. The same cell  
 29 selection vector is appended to every row of targets. These rows are assigned  
 30 with a positive label 1. We then randomly sampled 100 non-target genes  
 31 (denoted as  $v_d$ ) and construct the row vectors the same way as the target genes,  
 32 and these rows are assigned with a negative label 0. In other words, the training  
 33 matrix and label vector constructed from a drug  $d$  are of the following format:

$$34 \quad X_d = \begin{bmatrix} f_{cor}(d, P_{d1, \cdot}) & f_{PC}(d, P_{d1, \cdot}) & f_{PE_{max}}(d, P_{d1, \cdot}) & f_{PE_{avg}}(d, P_{d1, \cdot}) & f_{CS}(d, \cdot) \\ f_{cor}(d, P_{d2, \cdot}) & f_{PC}(d, P_{d2, \cdot}) & f_{PE_{max}}(d, P_{d2, \cdot}) & f_{PE_{avg}}(d, P_{d2, \cdot}) & f_{CS}(d, \cdot) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ f_{cor}(d, P_{dm, \cdot}) & f_{PC}(d, P_{dm, \cdot}) & f_{PE_{max}}(d, P_{dm, \cdot}) & f_{PE_{avg}}(d, P_{dm, \cdot}) & f_{CS}(d, \cdot) \\ f_{cor}(d, v_{d1, \cdot}) & f_{PC}(d, v_{d1, \cdot}) & f_{PE_{max}}(d, v_{d1, \cdot}) & f_{PE_{avg}}(d, v_{d1, \cdot}) & f_{CS}(d, \cdot) \\ f_{cor}(d, v_{d2, \cdot}) & f_{PC}(d, v_{d2, \cdot}) & f_{PE_{max}}(d, v_{d2, \cdot}) & f_{PE_{avg}}(d, v_{d2, \cdot}) & f_{CS}(d, \cdot) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ f_{cor}(d, v_{d100, \cdot}) & f_{PC}(d, v_{d100, \cdot}) & f_{PE_{max}}(d, v_{d100, \cdot}) & f_{PE_{avg}}(d, v_{d100, \cdot}) & f_{CS}(d, \cdot) \end{bmatrix}; y_d = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

35  
36

<sup>2</sup> <http://geneontology.org/page/go-enrichment-analysis>

1 | where  $m = |P_d|$ , the total number of targets for drug  $d$ . Therefore, the training  
2 | matrix  $X_d$  for drug  $d$  is of size  $(m + 100) \times 5|C_d|$ , and label vector  $y_d$  has length  
3 |  $(m + 100)$ .

4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15

1 References

2

3 1 Zhang, M. *et al.* Chaperoned ubiquitylation--crystal structures of the CHIP U  
4 box E3 ubiquitin ligase and a CHIP-Ubc13-Uev1a complex. *Mol. Cell* **20**, 525-  
5 538, doi:10.1016/j.molcel.2005.09.023 (2005).

6 2 Komander, D. *et al.* Structural insights into the regulation of PDK1 by  
7 phosphoinositides and inositol phosphates. *EMBO J.* **23**, 3918-3928,  
8 doi:10.1038/sj.emboj.7600379 (2004).

9