

Supplementary Online Materials

Rui-Sheng Wang, Joseph Loscalzo*

Department of Medicine, Brigham and Women's Hospital & Harvard Medical School, Boston,
MA

*Corresponding Author: E-mail: jloscalzo@partners.org

Compiling a human protein interactome

The interactome we used in this study is from [1] and also enhanced by incorporating additional data sources. Specifically, the protein-protein interactions are derived from several high-throughput yeast-two-hybrid studies [2-4] and also combined with binary interactions from IntAct and MINT databases [5, 6] and literature-curated interactions obtained by low throughput experiments reported in the IntAct, MINT, HPRD, BioGRID databases [5-8], as well as the CCSB Human Interactome (HI-2012, http://interactome.dfci.harvard.edu/H_sapiens/). CORUM, literature curated protein interactions (LCI) from the CCSB, and experimentally determined human protein complexes are also included in the set of protein-protein interactions [9, 10]. Protein-DNA regulatory interactions are taken from the TRANSFAC database [11], and kinase-substrate interactions are obtained from the PhosphositePlus database [12]. Metabolic enzyme-coupled interactions (two enzymes that share adjacent reactions) are derived from the KEGG and BiGG databases as compiled previously [13]. In addition, protein interactions from 3D structural prediction and signaling interactions are also included in the construction of the interactome [14, 15].

Network analyses and implementation

In this study, most of the network analyses were performed using Python, with the assistance of a Python package, NetworkX [16]. It contains many built-in network analysis algorithms, such as shortest path algorithms, subgraph induction, and random graph generators, etc. We readily used these algorithms for examining the proximity between drug targets and MI disease proteins. Specifically, we created an empty graph under the Python environment after we imported NetworkX:

```
>>> import networkx as nx
>>> G=nx.Graph()
```

We then used functions `G.add_node()` and `G.add_edge()` to add nodes (proteins) and edges (interactions) into the empty graph we created. Finally, the proximity between drug targets and MI disease proteins was examined by using the function `nx.shortest_path()` which returns all the shortest paths between a source node (drug target) and a target node (disease protein). By counting the number of pairs of drug targets and disease proteins that have shortest path lengths of 1 and 2, we obtained the number of pairs of drug targets and disease proteins that have interactions or have common neighbors. The average shortest path length can also be calculated from the output of `nx.shortest_path()`.

In addition, we used a null model to assess the significance of emergent properties. The null model keeps the human interactome unchanged and randomly selects 1,000 pairs of random protein sets of the same size as MI-related drug targets and MI disease proteins respectively:

```
>>> import random as rd
>>> r=rd.randint(0,N-1)
```

where N is the number of proteins in the human interactome. The topological properties of random protein sets are then compared with those of the sets of real drug targets and disease proteins. Specifically, we calculated proximity measures (i.e., the number of interactions, the number of protein pairs that have common neighbors, and the average shortest path length) between each random drug target set and random disease protein set. The 1000 proximity values from random protein sets form a null normal distribution after we fit the histograms using the function `normfit(data)` in Matlab (Mathworks, Inc). The significance of emergent properties of observations (i.e. P -value) was obtained by comparing them with null models:

```
>>> [muhat,sigmahat] = normfit(data)
>>> p=1-normcdf(x,muhat,sigmahat)
```

where “data” store the 1000 proximity values from random protein sets, and “x” is the observed topological parameter.

To find modules densely connecting MI-related drug targets and MI disease proteins, we constructed a bipartite network using the interactions between them, and used the Louvain method to maximize the modularity function Q [17, 18] defined for characterizing the modularity of complex networks:

$$Q = \frac{1}{2m} \sum_{i,j} (A_{i,j} - \frac{k_i k_j}{m}) \delta(c_i, c_j),$$

where m is the total number of edges; the $A_{i,j}$ are the adjacency matrix elements; k_i and k_j are the degrees of node i and node j , respectively; c_i and c_j are the module indices of node i and node j , respectively; and δ , the delta function, is equal to 1 if $c_i = c_j$, and is otherwise equal to 0. The Louvain method was implemented by the Python package NetworkX:

```
>>> import community
>>> import networkx as nx
>>> C_best = community.best_partition(G)
```

where “G” is the graph for the bipartite network. This method works as hierarchical clustering and returns network partitions at the different levels by using the function:

```
>>> dendo = community.generate_dendrogram(G)
>>> C_i=community.partition_at_level(dendo, i)
>>> M=community.modularity(C_best, G)
```

which allows us to examine the modularity value of partitions at different level and choose a best partition.

Statistical tests and tools

When we assessed the biological relevance of DTD modules, we created a control for significance, i.e., we randomly selected an interaction set with the same number of interactions as the module and calculated the enrichment of drug pairs that have similar side effects or therapeutic effects in the random modules. In addition, GO-based functional similarity of pairs of MI-related drug targets and MI disease proteins was quantified by GS2 (GO-based similarity of gene sets) developed in [19]:

```
>>> from pyGS2 import get_go_graph
>>> tree = get_go_graph(open('go_daily-termdb.obo-xml'))
>>> s=tree.GS2([gene1, gene2])
```

where the GO annotation file `go_daily-termdb.obo-xml` was download from the Gene Ontology database [20]. Unless otherwise specified, when we assess the significance of emergent properties of observations by comparing them with null models (random controls), all

P-values are obtained by fitting the histograms into normal distributions using the ‘normfit’ command in Matlab (Mathworks, Inc):

```
>>> [muhat,sigmahat] = normfit(data)
>>> p=1-normcdf(x,muhat,sigmahat)
```

where “data” store the values from random controls, and “x” is the observed topological parameter.

The proximity between MI-related drug targets and MI disease proteins using the interactions from STRING v10

In this study, we chose to use the comprehensive human interactome we compiled from different databases. This interactome consists of diverse types of physical molecular interactions and has recently been shown by us to have great potential in deciphering disease-disease [1]. One of the main reasons that we used this specific interactome is because the majority of this consolidated interactome is derived from unbiased experimental detection of physical protein-protein interactions. There are many other protein interaction databases with higher coverage, such as STRING [21] and HAPPI [22]. However, these databases contain a large number of predicted, rather than experimentally ascertained, interactions, many of which are of uncertain statistical confidence. While experimentally derived protein-protein interactions can be associated with significant false positive ratios, predicted protein interactions can be even less reliable. Moreover, only a subset of the interactions in STRING represents physical interactions; it contains largely functional interactions predicted from gene expression correlations and other datasets. Since our study focuses on identifying potential pathways underlying drug actions, predicted functional, as compared to physical, interactions are sub-optimal for this purpose.

Nevertheless, although STRING contains many predicted physical and functional interactions, it can still be used to assess the proximity of MI-related drug targets to MI disease proteins. We, therefore, downloaded human protein functional links (protein-protein interactions) from STRING v10. Each protein-protein interaction has a confidence score in range [100,1000]. There are 8,548,002 protein interactions for *Homo sapiens*. We used 900 as the confidence threshold and obtained a set of 205,450 protein-protein interactions after we mapped the protein aliases to HGNC gene names and removed redundancy in the data. We then examined the closeness relationships between MI-related drug targets and MI disease genes using the interactions in this subset of STRING v10 and found that the conclusions remain the same as those obtained using our consolidated human interactome. Specifically, we identified 1,605 interactions between MI-related drug targets and MI disease proteins, which is significantly greater than the number of interactions between two random sets of the same size ($P < 1.0 \times 10^{-16}$), as shown in Figure S6 (A). Figure S6 (B) shows that there are significantly more pairs of MI-related drug targets and MI disease proteins with common neighbors than protein pairs from two random sets ($P < 1.0 \times 10^{-16}$). The average shortest path length between MI-related drug targets and MI disease proteins is 4.20, significantly smaller than that between two random sets of the same size [$P < 1.0 \times 10^{-16}$, Figure S6 (C)]. We also assessed the closeness relationship between control drug targets and MI disease proteins and the closeness relationships between MI-related drug targets and control disease proteins using the interactions from STRING v10. The results, shown in Figure S7 (A-C), indicate that MI-related drug targets are significantly closer to MI disease proteins than control drug targets. Figure S8 (A-C) shows that MI disease proteins are significantly closer to MI-related drug targets than control disease proteins. This validation using another interaction

database further confirms the closeness relationships between MI-related drug targets and MI disease genes at a systems level.

Impact of module detection methods on the results

A network module is conceptually defined as a group of nodes in the network that are more densely interconnected than to the rest of the network. There is not a strict mathematical definition for network modules. It is not uncommon for two different module-finding techniques to give different results, as each method has its own design principles based on network topology; however, we expect the results should be largely concordant if the general principles upon which the analysis is performed hold. Whether drug targets/disease proteins are included or excluded from the modules depends on their connections to other drug targets/disease proteins. The majority of drug targets and disease proteins would be included in the modules; only those with peripheral connections may be excluded if we use a different module-finding method.

The differences in results from two module-finding techniques arise from two sources: method design principles, and the intrinsic properties of the biological network under consideration. A number of studies have shown that biological networks are not simply modular; rather, they display strong multi-scale modularity or hierarchical modularity [23, 24]. Different network partitions may have the same modularity values, which defines the problem of multi-solution limitation common to all module-based network analyses. A module-detection method incorporating biological knowledge may, therefore, be useful in reducing the impact of the multi-solution problem.

The method we used in this study is the Louvain method for community detection. It is a widely used greedy optimization method that maximizes the modularity function Q [17, 18]. To

examine the impact of module-finding techniques, we used simulated annealing to maximize another modularity measure, modularity density, D [25, 26]. The results, summarized in Table S3, give us more small modules: 20 modules with more than five proteins, 12 of them are 100% contained within our large DTD modules, confirming the multiple-scale modularity mentioned above. The majority of other modules are more than 85% contained within our DTD modules, indicating the robustness of our DTD modules. Large modules provide a more complete overall view of the pathways, and small modules give better resolution, but may be incomplete and, thereby, lose some information.

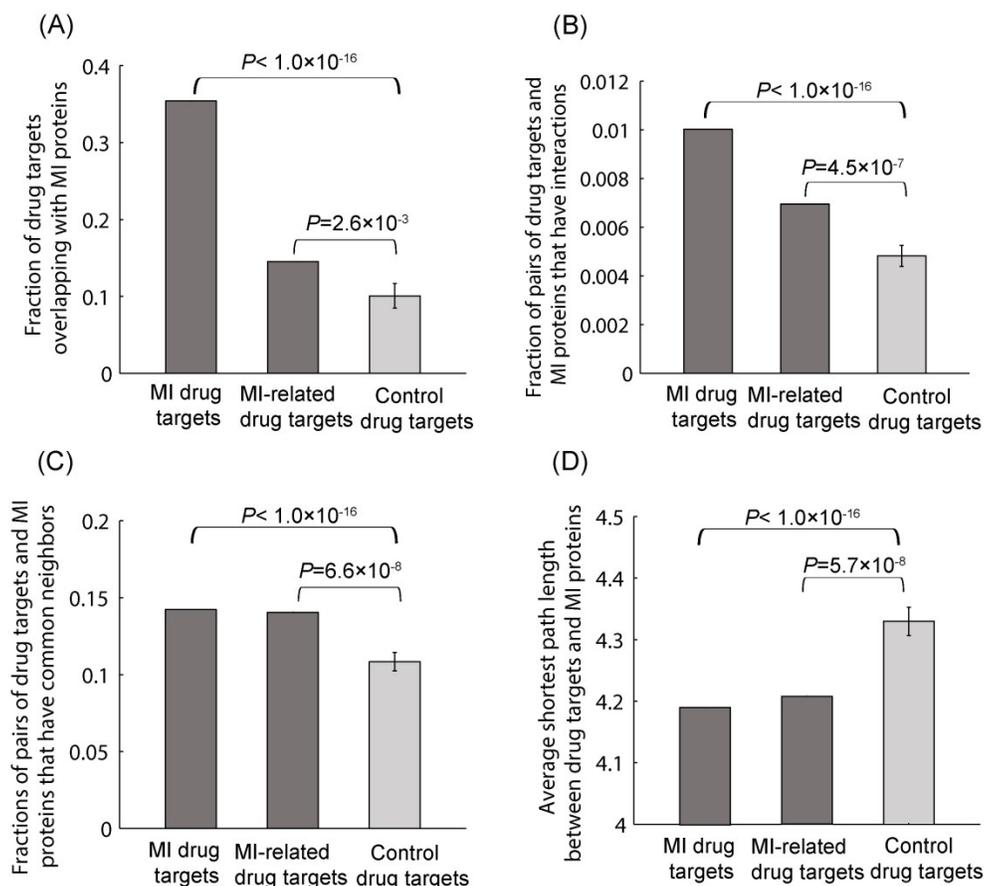


Figure S1. MI-related drug targets are closer to MI disease proteins than control drug targets in the interactome even after removing MI drug targets. (A) Compared to control drug targets, MI-related drug targets have greater overlap with MI disease proteins. (B) Compared to control drug targets, MI-related drug targets have more interactions with MI disease proteins. (C) Compared to control drug targets, there are more pairs of MI-related drug targets and MI disease proteins that have common neighbors in the interactome. (D) Compared to control drug targets, MI-related drug targets have a smaller average shortest path length with MI disease proteins.

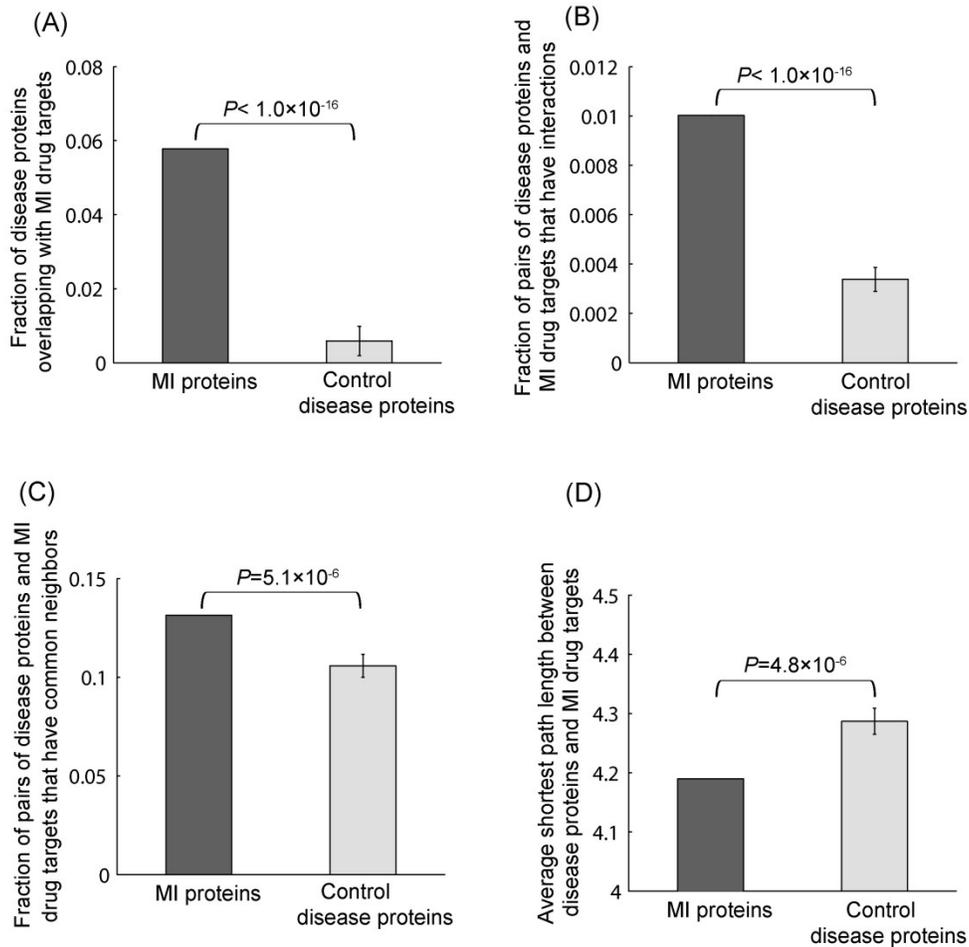


Figure S2. MI disease proteins are closer to MI drug targets than control disease proteins.

(A) Compared to control disease proteins, MI disease proteins have greater overlap with MI drug

targets. (B) Compared to control disease proteins, MI disease proteins have more interactions with MI drug targets. (C) Compared to control disease proteins, there are more pairs of MI disease proteins and MI drug targets that have common neighbors in the interactome. (D) Compared to control disease proteins, MI disease proteins have a smaller average shortest path length with MI drug targets in the interactome.

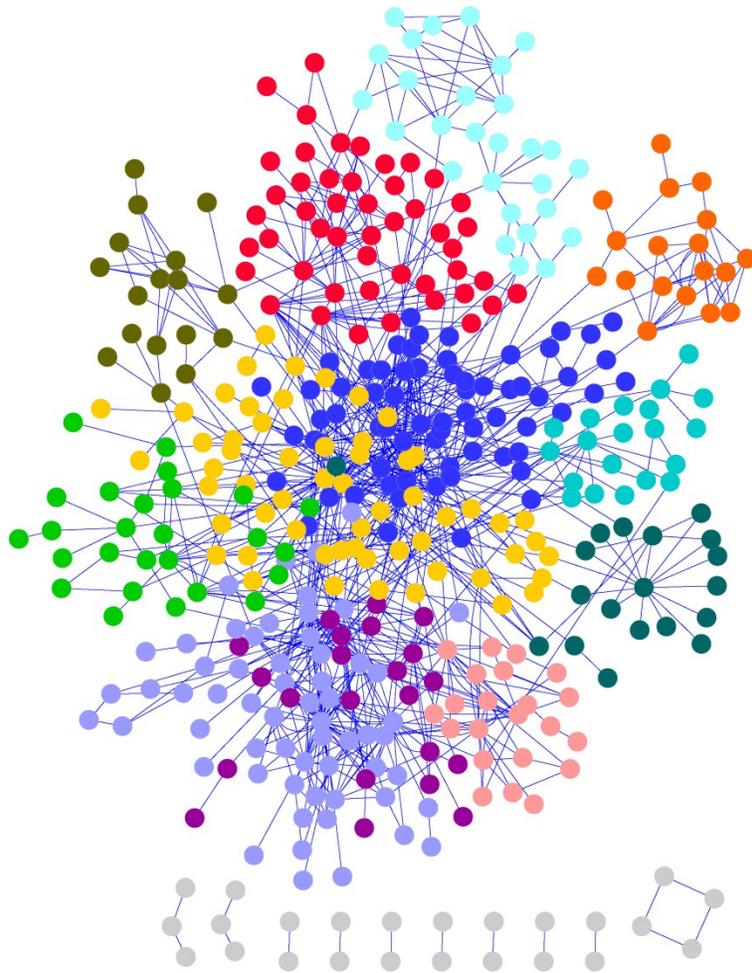
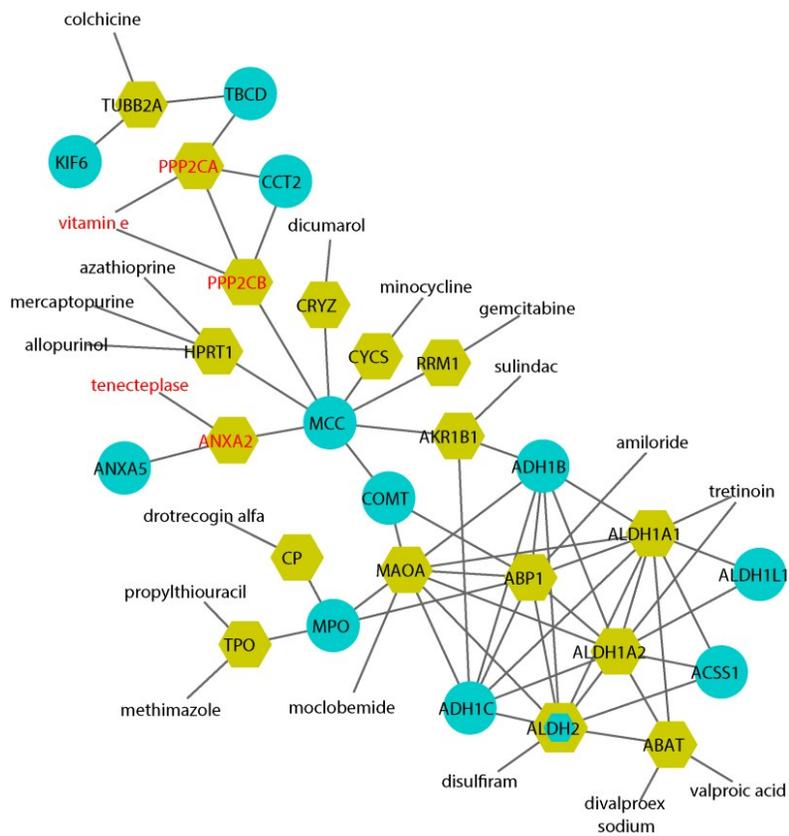
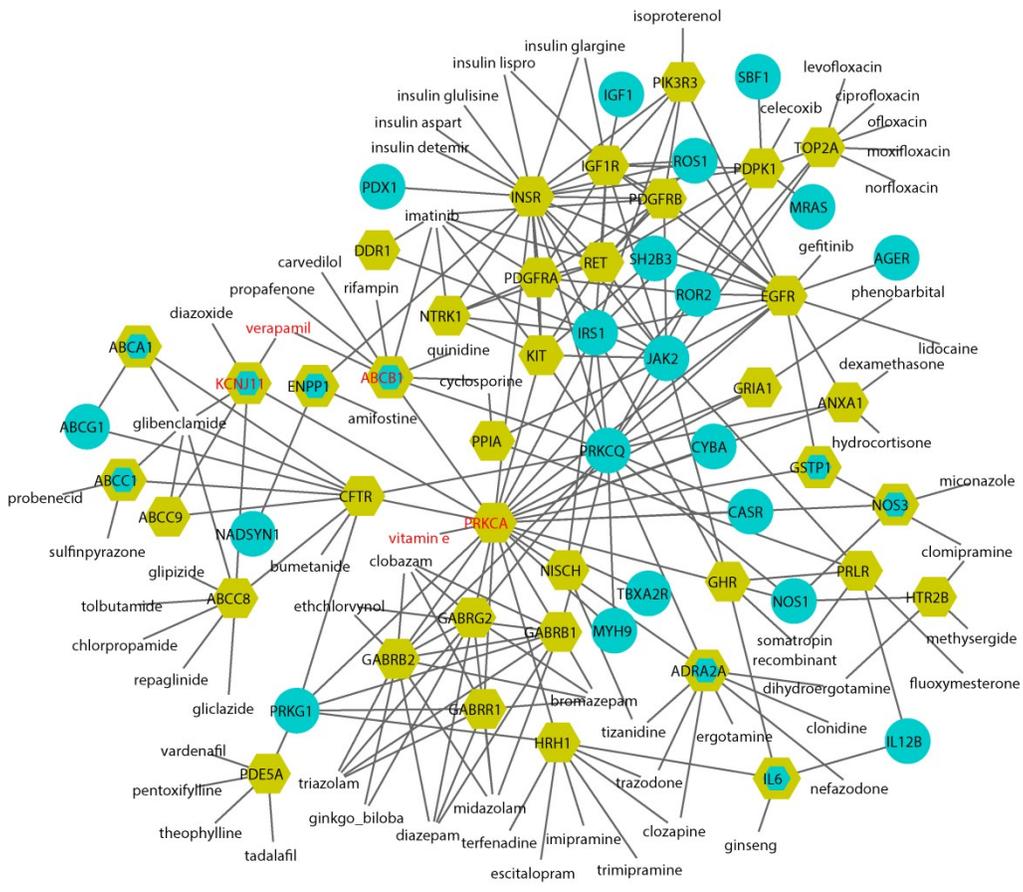


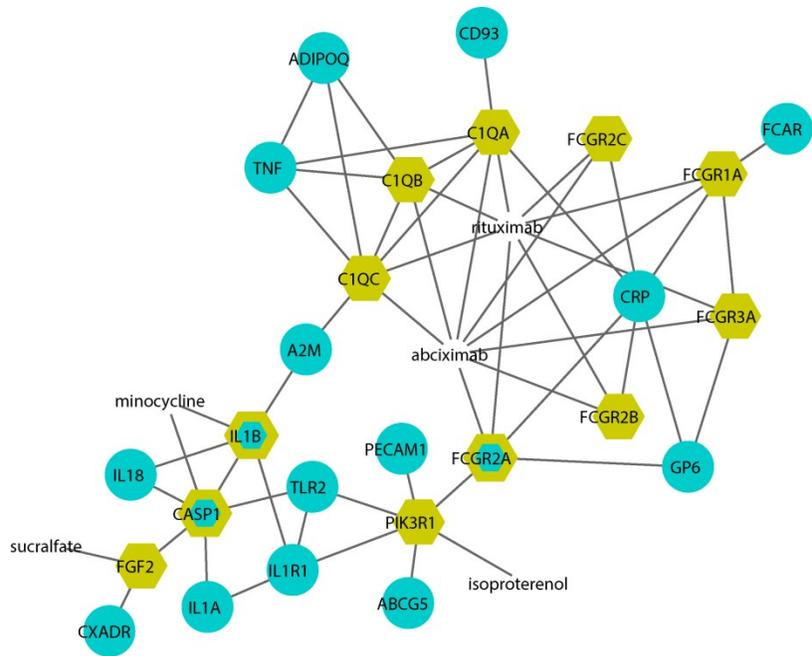
Figure S3. The modular bipartite network of MI-related drug targets and MI disease proteins. Nodes of the same color define a module. The nodes in gray are isolated (orphan) nodes.



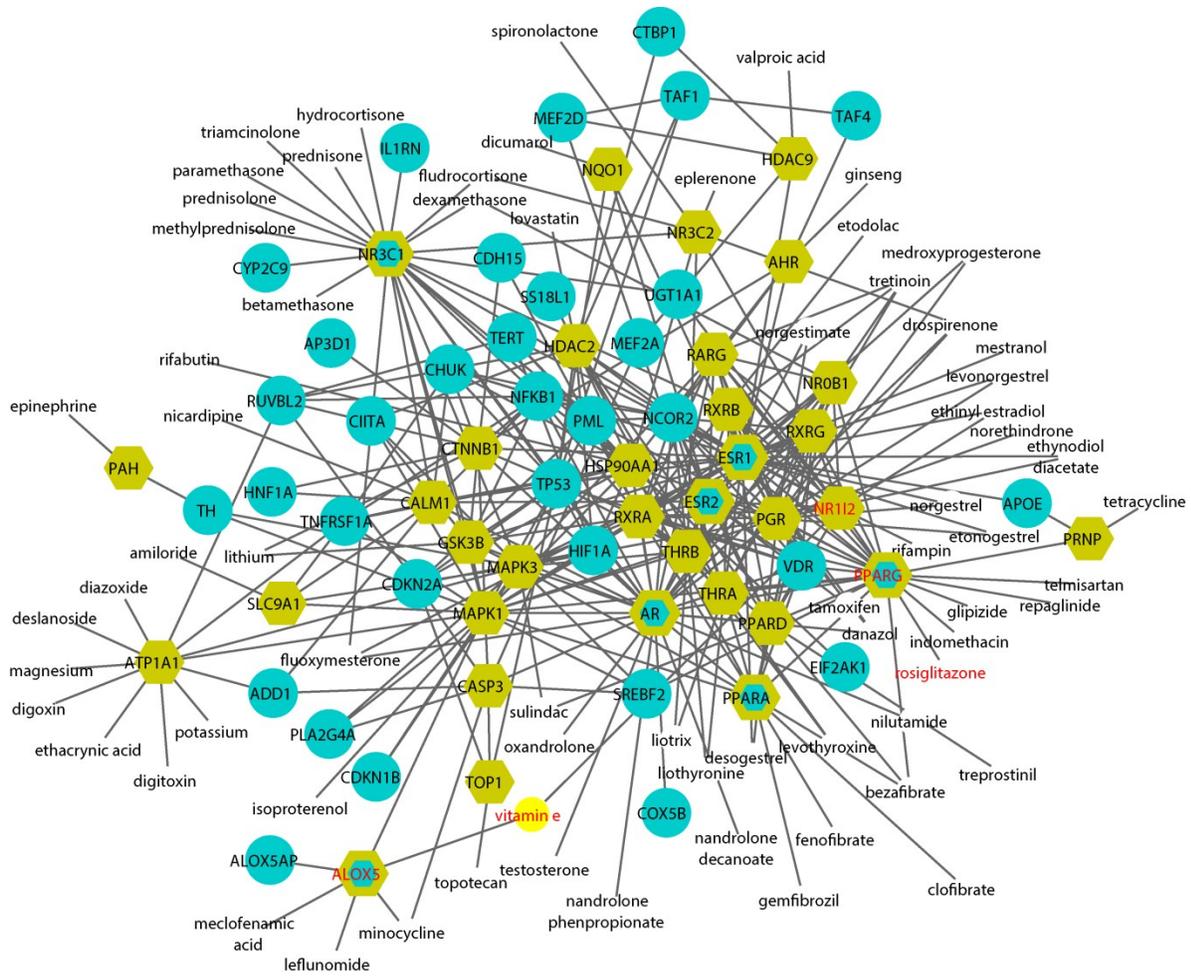
Module 1



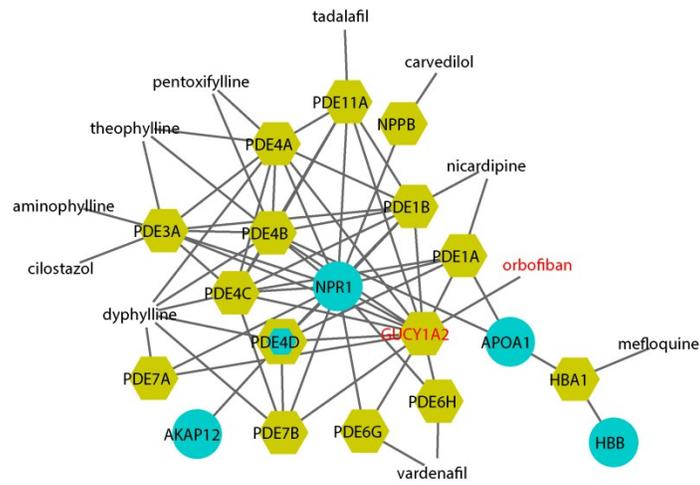
Module 2



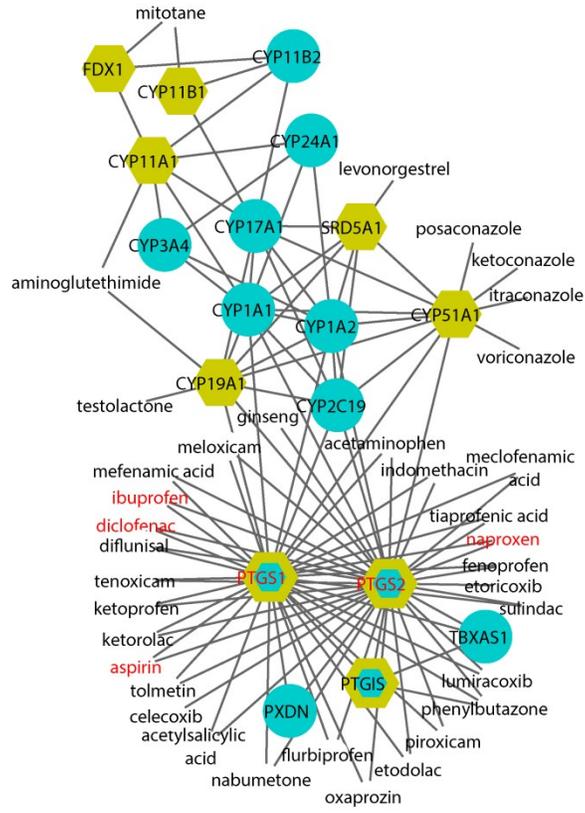
Module 4



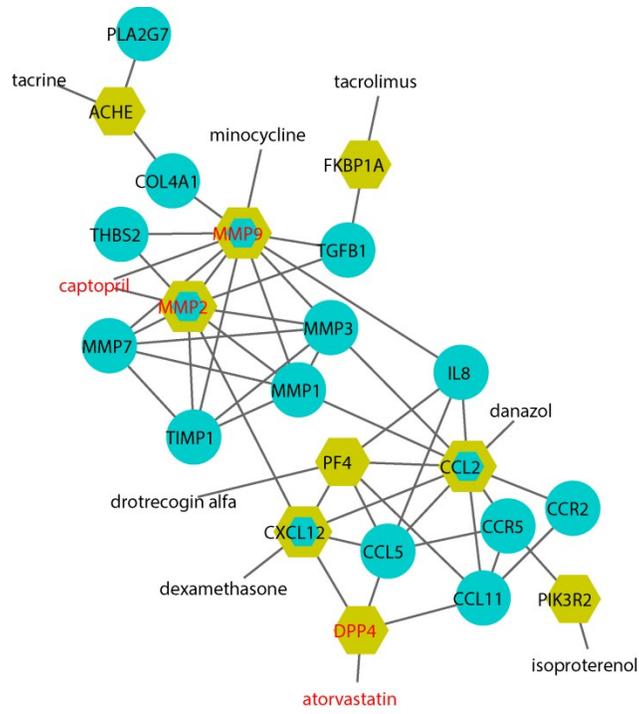
Module 5



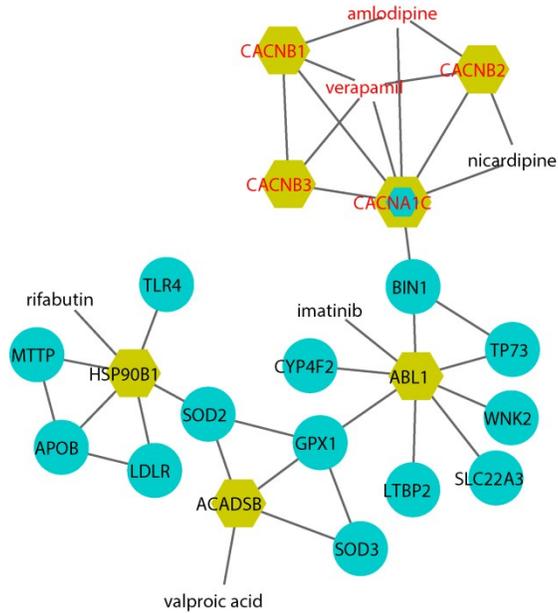
Module 6



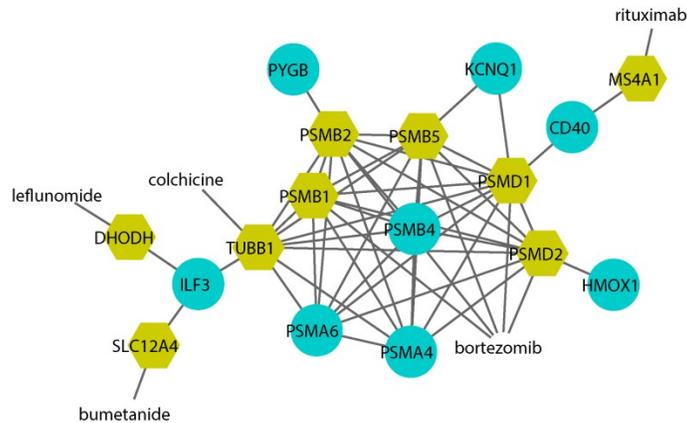
Module 7



Module 8



Module 11



Module 12

Figure S4. The DTD modules. The bold text represents MI drugs or MI drug targets. The blue nodes represent MI disease proteins, and the yellow nodes denote MI-related drug targets. The nodes with both colors are both drug targets and MI disease proteins. The nodes with only labels (without node shapes) are drugs. MI drugs and MI drug targets are denoted in red.

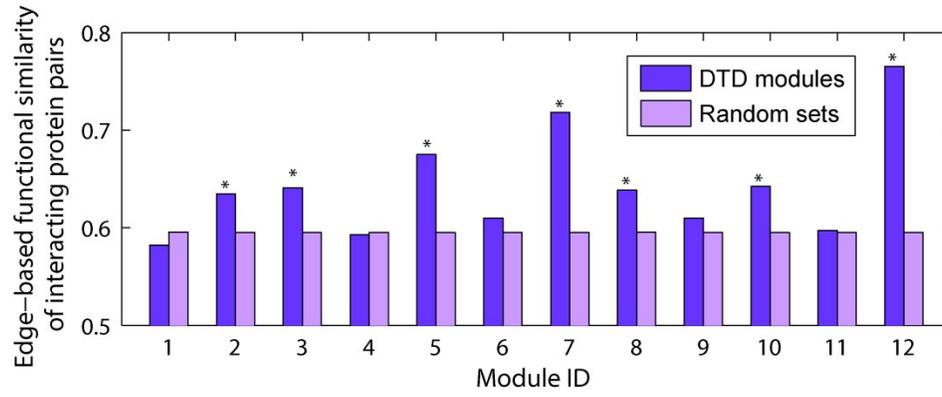


Figure S5. Functional similarity of interacting protein pairs in the DTD modules, $*P < 0.05$.

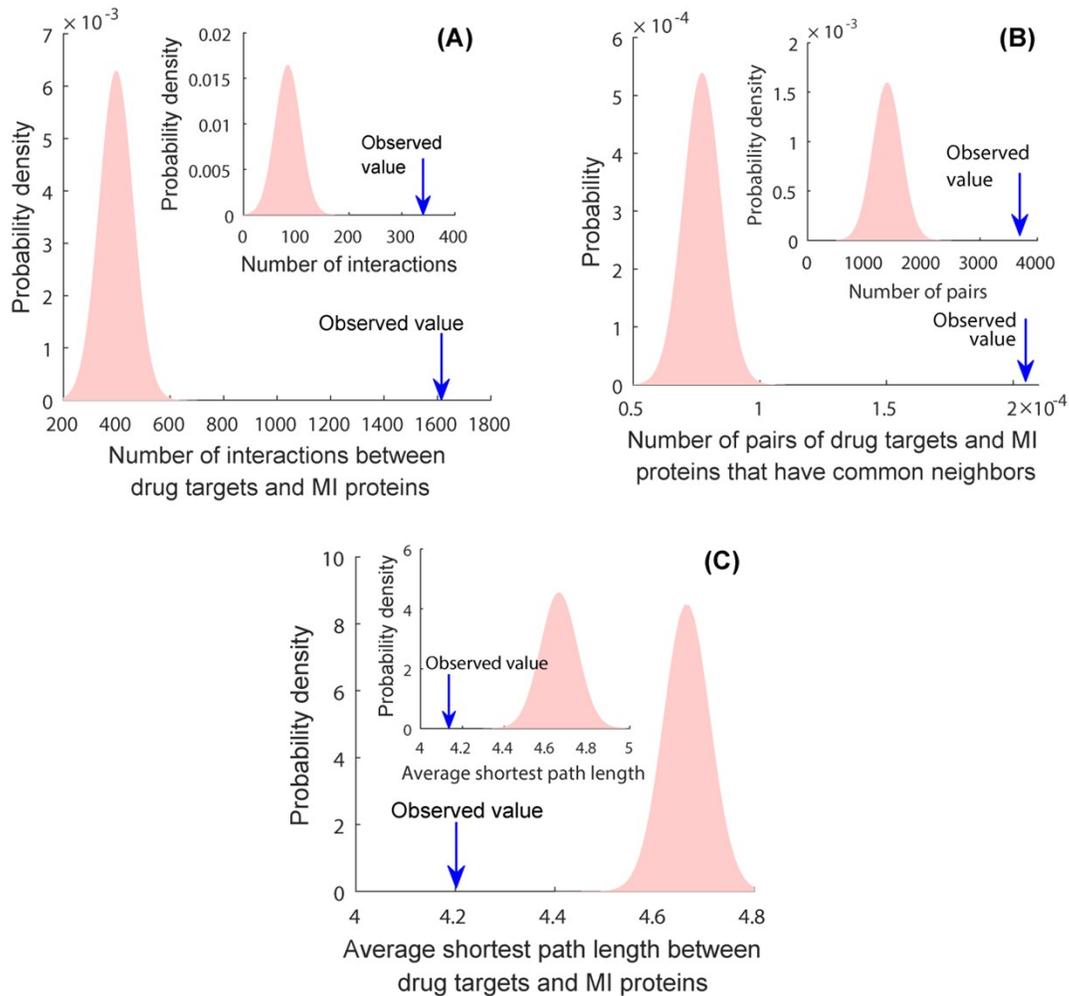


Figure S6. The closeness relationships between MI(-related) drug targets and MI disease proteins using the interactions from STRING v10. (A) MI-related drug targets (MI drug target, inset) and MI disease proteins have significantly more interactions than expected by chance. **(B)** There are significantly more pairs of MI-related drug targets (MI drug target, inset) and MI disease proteins with common neighbors than expected by chance. **(C)** The average shortest path length between MI-related drug targets (MI drug target, inset) and MI disease proteins is significantly smaller than that between two random gene sets.

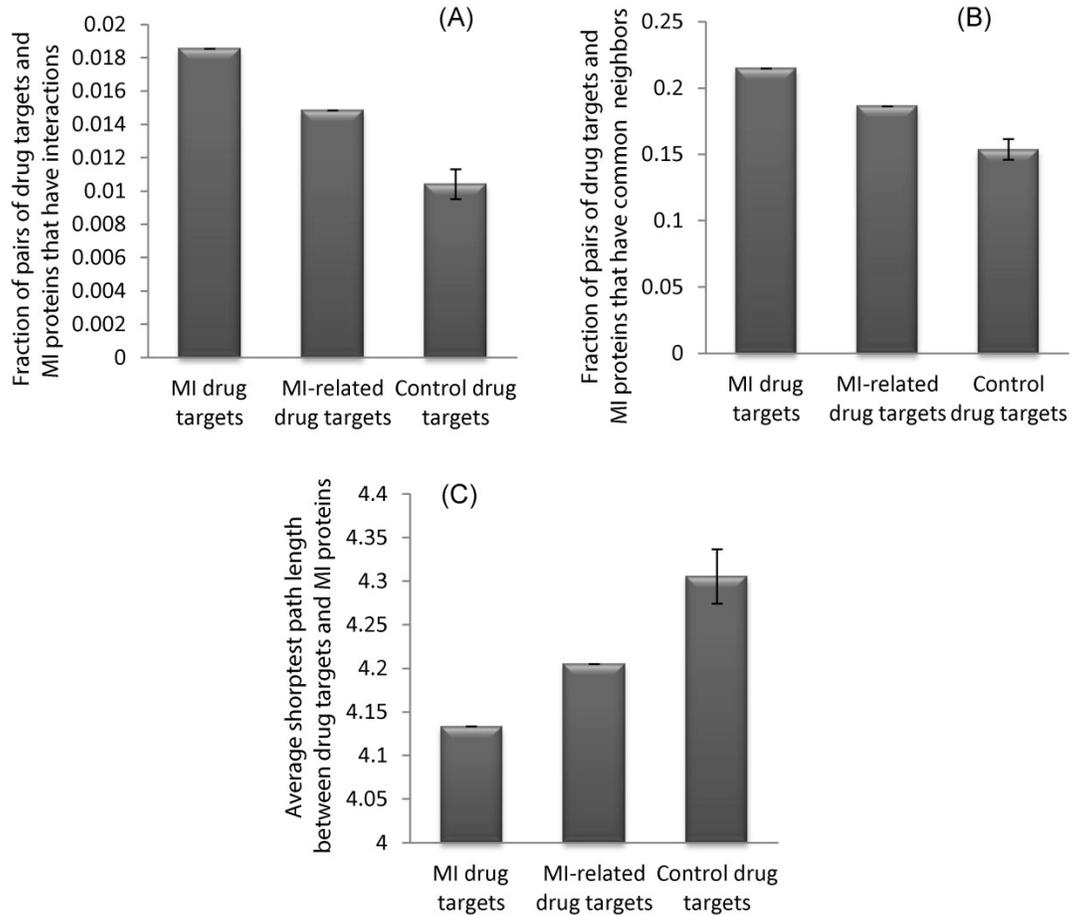


Figure S7. The proximity between control drug targets and MI disease proteins using the interactions from STRING v10. (A) Compared to control drug targets, MI(-related) drug targets have more interactions with MI disease proteins ($P < 1.0 \times 10^{-16}$ for MI drug targets, and $P = 3.8 \times 10^{-7}$ for MI-related drug targets). (B) Compared to control drug targets, there are more pairs of MI(-related) drug targets and MI disease proteins that have common neighbors ($P = 1.1 \times 10^{-15}$ for MI drug targets, and $P = 1.2 \times 10^{-5}$ for MI-related drug targets). (C) Compared to control drug targets, MI(-related) drug targets have a smaller average shortest path length with MI disease proteins ($P = 1.6 \times 10^{-8}$ for MI drug targets, and $P = 6.2 \times 10^{-4}$ for MI-related drug targets).

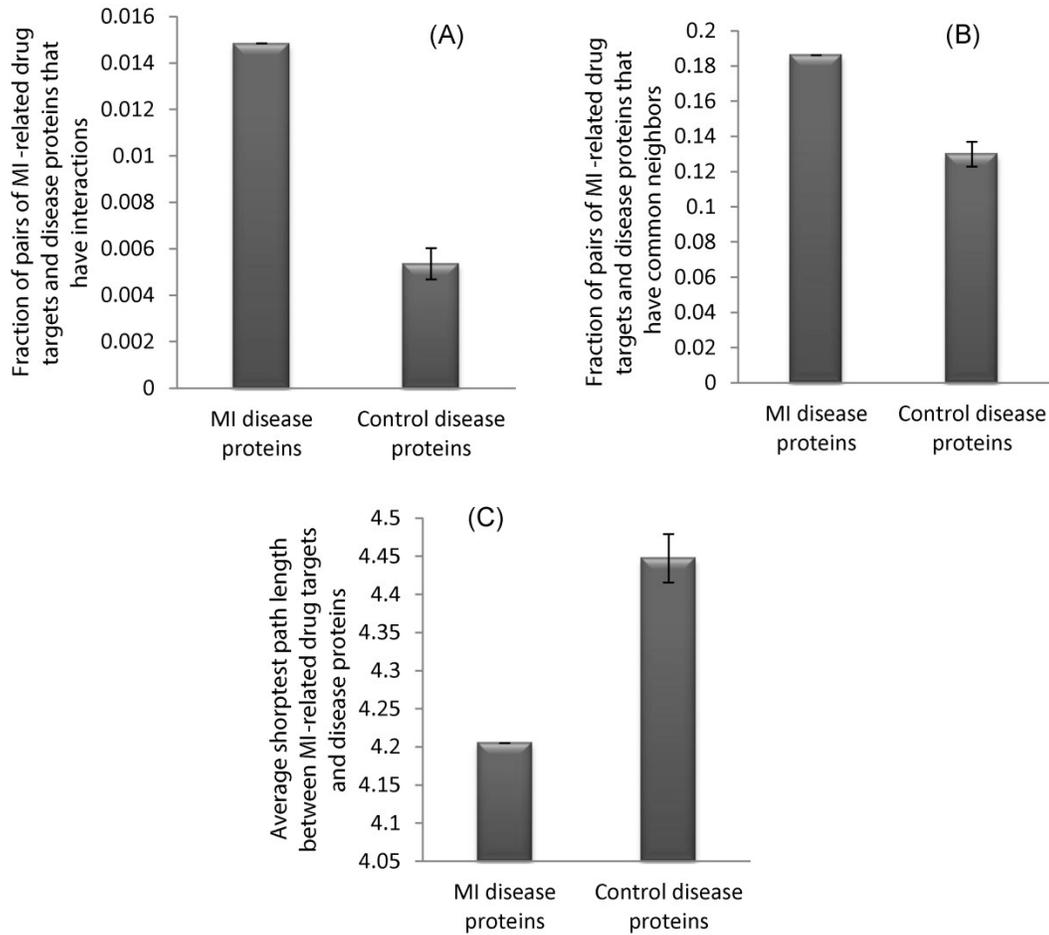


Figure S8. The proximity between MI-related drug targets and control disease proteins using the interactions from STRING v10. (A) Compared to control disease proteins, MI disease proteins have more interactions with MI-related drug targets ($P < 1.0 \times 10^{-16}$). (B) Compared to control disease proteins, there are more pairs of MI disease proteins and MI-related drug targets that have common neighbors ($P = 1.1 \times 10^{-15}$). (C) Compared to control disease genes, MI disease proteins have a smaller average shortest path length with MI-related drug targets ($P = 1.1 \times 10^{-14}$).

Table S1. Contingency table for the enrichment of module non-MI drugs with cardiovascular effects, P -value = 1.9×10^{-3} (Chi-squared test). Non-module drugs have targets in the interactome.

	Cardiovascular effects	No cardiovascular effects	Total
Module drugs	133	91	224
Non-module drugs	13	28	41
Total	146	119	265

Table S2. Contingency table for the enrichment of cardiovascular-associated proteins in drug targets, P -value = 4.5×10^{-3} (Chi-squared test).

	Cardiovascular proteins	No cardiovascular proteins	Total
Module targets	146	45	191
Non-module targets	63	42	105
Total	209	87	296

Table S3. Comparison of DTD modules with modules derived by optimization of modularity density D .

D module IDs	Module size	Overlap with DTD modules	D module IDs	Module size	Overlap with DTD modules
D_1	29	Module 2 (86.2%)	D_{11}	15	Module 3(100%)
D_2	11	Module 3(90.9%)	D_{12}	36	Module 9 (41.7%) Module 10 (55.6%)
D_3	17	Module 8 (52.9%) Module 9 (29.4%)	D_{13}	13	Module 6 (100%)
D_4	73	Module 2 (23.3%) Module 3 (19.2%) Module 5 (30.1%)	D_{14}	7	Module 7 (100%)
D_5	26	Module 9 (100%)	D_{15}	9	Module 2 (88.9%)
D_6	37	Module 5 (97.3%)	D_{16}	6	Module 1 (100%)
D_7	7	Module 9 (100%)	D_{17}	11	Module 7 (100%)
D_8	8	Module 11 (87.5%)	D_{18}	10	Module 12 (100%)
D_9	11	Module 2 (100%)	D_{19}	7	Module 11 (100%)
D_{10}	7	Module 4 (100%)	D_{20}	8	Module 8 (100%)

References

1. Menche, J., et al., *Uncovering disease-disease relationships through the incomplete interactome*. Science, 2015. **347**(6224): p. 1257601.
2. Stelzl, U., et al., *A human protein-protein interaction network: a resource for annotating the proteome*. Cell, 2005. **122**(6): p. 957-68.
3. Rual, J.F., et al., *Towards a proteome-scale map of the human protein-protein interaction network*. Nature, 2005. **437**(7062): p. 1173-8.
4. Venkatesan, K., et al., *An empirical framework for binary interactome mapping*. Nat Methods, 2009. **6**(1): p. 83-90.
5. Aranda, B., et al., *The IntAct molecular interaction database in 2010*. Nucleic Acids Res, 2010. **38**(Database issue): p. D525-31.
6. Licata, L., et al., *MINT, the molecular interaction database: 2012 update*. Nucleic Acids Res, 2012. **40**(Database issue): p. D857-61.
7. Keshava Prasad, T.S., et al., *Human Protein Reference Database--2009 update*. Nucleic Acids Res, 2009. **37**(Database issue): p. D767-72.
8. Chatr-Aryamontri, A., et al., *The BioGRID interaction database: 2013 update*. Nucleic Acids Res, 2013. **41**(Database issue): p. D816-23.
9. Ruepp, A., et al., *CORUM: the comprehensive resource of mammalian protein complexes--2009*. Nucleic Acids Res, 2010. **38**(Database issue): p. D497-501.

10. Havugimana, P.C., et al., *A census of human soluble protein complexes*. Cell, 2012. **150**(5): p. 1068-81.
11. Matys, V., et al., *TRANSFAC: transcriptional regulation, from patterns to profiles*. Nucleic Acids Res, 2003. **31**(1): p. 374-8.
12. Hornbeck, P.V., et al., *PhosphoSitePlus: a comprehensive resource for investigating the structure and function of experimentally determined post-translational modifications in man and mouse*. Nucleic Acids Res, 2012. **40**(Database issue): p. D261-70.
13. Lee, D.S., et al., *The implications of human metabolic network topology for disease comorbidity*. Proc Natl Acad Sci U S A, 2008. **105**(29): p. 9880-5.
14. Zhang, Q.C., et al., *Structure-based prediction of protein-protein interactions on a genome-wide scale*. Nature, 2012. **490**(7421): p. 556-60.
15. Vinayagam, A., et al., *A directed protein interaction network for investigating intracellular signal transduction*. Sci Signal, 2011. **4**(189): p. rs8.
16. Harenberg, S., et al., *Community detection in large-scale networks: a survey and empirical evaluation*. Wiley Interdisciplinary Reviews: Computational Statistics, 2014. **6**: p. 426-439.
17. Newman, M.E., *Modularity and community structure in networks*. Proc Natl Acad Sci U S A, 2006. **103**(23): p. 8577-82.
18. Blondel, V.D., et al., *Fast unfolding of communities in large networks*. Journal of Statistical Mechanics: Theory and Experiment 2008. **10**: p. P10008.
19. Ruths, T., D. Ruths, and L. Nakhleh, *GS2: an efficiently computable measure of GO-based similarity of gene sets*. Bioinformatics, 2009. **25**(9): p. 1178-84.
20. Ashburner, M., et al., *Gene ontology: tool for the unification of biology. The Gene Ontology Consortium*. Nat Genet, 2000. **25**(1): p. 25-9.
21. Szklarczyk, D., et al., *STRING v10: protein-protein interaction networks, integrated over the tree of life*. Nucleic Acids Res, 2015. **43**(Database issue): p. D447-52.
22. Chen, J.Y., S. Mamidipalli, and T. Huan, *HAPPI: an online database of comprehensive human annotated and predicted protein interactions*. BMC Genomics, 2009. **10 Suppl 1**: p. S16.
23. Clauset, A., C. Moore, and M.E.J. Newman, *Hierarchical structure and the prediction of missing links in networks*. Nature, 2008. **453**(7191): p. 98-101.
24. Meunier, D., R. Lambiotte, and E.T. Bullmore, *Modular and hierarchically modular organization of brain networks*. Front Neurosci, 2010. **4**: p. 200.
25. Guimera, R. and L.A. Nunes Amaral, *Functional cartography of complex metabolic networks*. Nature, 2005. **433**(7028): p. 895-900.
26. Li, Z., et al., *Quantitative function for community detection*. Phys Rev E Stat Nonlin Soft Matter Phys, 2008. **77**(3 Pt 2): p. 036109.