

# NOGEA: Network-Oriented Gene Entropy Approach for Dissecting Disease Comorbidity and Drug Repositioning

Zihu Guo<sup>1,#,a</sup>, Yingxue Fu<sup>1,#,b</sup>, Chao Huang<sup>1,#,c</sup>, Chunli Zheng<sup>2,#,d</sup>, Ziyin Wu<sup>1,#,e</sup>, Xuotong Chen<sup>1,f</sup>, Shuo Gao<sup>1,g</sup>, Yaohua Ma<sup>2,h</sup>, Mohamed Shahen<sup>3,i</sup>, Yan Li<sup>4,j</sup>, Pengfei Tu<sup>5,k</sup>, Jingbo Zhu<sup>6,l</sup>, Zhenzhong Wang<sup>7,m</sup>, Wei Xiao<sup>7,\*,n</sup>, Yonghua Wang<sup>1,2,\*,o</sup>

<sup>1</sup> College of Life Science, Northwest A & F University, Yangling, Shaanxi 712100, China.

<sup>2</sup> College of Life Science, Northwest University, Xi'an, Shaanxi 710069, China.

<sup>3</sup> Zoology Department, Faculty of Science, Tanta University, Tanta 31527, Egypt.

<sup>4</sup> Key Laboratory of Industrial Ecology and Environmental Engineering (MOE), Department of Materials Sciences and Chemical Engineering, Dalian University of Technology, Dalian, Liaoning 116024, China.

<sup>5</sup> State Key Laboratory of Natural and Biomimetic Drugs, School of Pharmaceutical Sciences, Peking University, Beijing 100191, China.

<sup>6</sup> School of Food Science and Technology, Dalian Polytechnic University, Dalian, Liaoning 116034, China.

<sup>7</sup> State Key Laboratory of New-tech for Chinese Medicine Pharmaceutical Process, Lianyungang, Jiangsu 222001, China.

# Equal contribution.

\* Corresponding author(s).

E-mail: xw\_kanion@163.com (Xiao W), yh\_wang@nwafu.edu.cn (Wang Y).

**Running title:** Guo Z et al / Network-Oriented Gene Entropy Approach

## 28    **Abstract**

29    Rapid development of high-throughput technologies has permitted the identification  
 30    of an increasing number of disease-associated genes (DAGs), which are important for  
 31    understanding disease initiation and developing precision therapeutics. However,  
 32    DAGs often contain large amounts of redundant or false positive information, leading  
 33    to difficulties in quantifying and prioritizing potential relationships between these  
 34    DAGs and human diseases. In this study, a network-oriented gene entropy approach  
 35    (NOGEA) is proposed for accurately inferring master genes that contribute to specific  
 36    diseases by quantitatively calculating their perturbation abilities on directed disease-  
 37    specific gene networks. In addition, we confirmed that the master genes identified by  
 38    NOGEA have a high reliability for predicting disease-specific initiation events and  
 39    progression risk. Master genes may also be used to extract the underlying information  
 40    of different diseases, thus revealing mechanisms of disease comorbidity. More  
 41    importantly, approved therapeutic targets are topologically localized in a small  
 42    neighborhood of master genes on the interactome network, which provides a new way  
 43    for predicting new drug-disease associations. Through this method, 11 old drugs were  
 44    newly identified and predicted to be effective for treating pancreatic cancer and then  
 45    validated by *in vitro* experiments. Collectively, the NOGEA was useful for  
 46    identifying master genes that control disease initiation and co-occurrence, thus  
 47    providing a valuable strategy for drug efficacy screening and repositioning. NOGEA  
 48    codes are publicly available at <https://github.com/guozihuaa/NOGEA>.

49  
 50    **KEYWORDS:** Systems pharmacology; Gene entropy; Disease gene network;  
 51    Disease comorbidity; Drug repositioning

## 52 Introduction

53 The onset and progression of most complex diseases often involves the dysfunction of  
 54 thousands of genes as well as certain altered interactions among them. High-  
 55 throughput technologies such as gene expression profiling and whole genome  
 56 sequencing have permitted the identification of an increasing number of disease  
 57 associated genes (DAGs) [1], which may provide valuable insight into mechanisms of  
 58 disease initiation and progression. However, as the existing DAGs are usually derived  
 59 from multiple sources, they often contain large amounts of redundant or false positive  
 60 information [2] due to collection bias and noise, such that causal relationships among  
 61 these genes in most cases remain elusive. Therefore, identifying master genes that  
 62 control disease state transitions from large numbers of DAGs plays a critical role in  
 63 understanding disease initiation mechanisms. In addition, complex diseases show  
 64 considerable comorbidity [3]. The master gene defects in one disease may initiate  
 65 cascades of interactions that lead to the co-occurrence of multiple diseases in a given  
 66 patient. Pharmacological targeting of the DAG module on the human interactome has  
 67 proven to be a valuable strategy for drug efficacy screening [4]. At present, it is  
 68 unclear whether the identification of master genes will further facilitate the network-  
 69 based drug repositioning.

70 Recent trends in omics technologies and complex biological networks have led to  
 71 a proliferation of attempts to find the master genes for different diseases. For  
 72 example, genome-wide association studies (GWAS) have emerged as a powerful tool  
 73 for detecting sequence variation associated with many human traits and diseases [5].  
 74 Due to the low-frequency of many mutations, GWAS usually require large cohort  
 75 sizes to attain sufficient statistical power. More importantly, GWAS identify only the  
 76 genetic risk factors associated with disease, rather than the master genes of the disease  
 77 phenotypes because patient genomes contain a certain proportion of “passenger  
 78 mutations” [6] and the initiation of many diseases is often triggered by the interplay  
 79 between genetic and non-genetic factors. Transcriptome analysis is considered to be  
 80 an effective complement of GWAS for its ability to capture non-genetic perturbations  
 81 to the organism. Yet variations in mRNA expression are sometimes caused by  
 82 aberrant protein activity of upstream regulators such as transcription factors, making it  
 83 difficult to directly identify the master gene set using transcriptome profiling [7].

84 Recently, gene co-expression-based approaches have been proposed to construct  
 85 context-specific regulatory networks [8] and a local network entropy measure has  
 86 been developed based on co-expression networks for identifying master genes [9].

87 While these approaches provide new ways to find master genes, building a highly  
88 confident co-expression regulatory network often requires large sample sizes, which  
89 are usually not available for relatively rare diseases. To overcome this limitation,  
90 protein-protein interaction (PPI) network-based approaches have been developed to  
91 infer master genes that are important for disease-related biological processes, such as  
92 predicting therapeutic targets [10] or driver genes [11]. Some topological parameters  
93 such as the degree and betweenness centrality of the nodes are usually used as  
94 important measures to screen master genes [12]. However, current approaches are  
95 based mainly on the constant global undirected interactome, ignoring the fact that  
96 disease initiation and therapeutics are frequently context-dependent, depending on  
97 specific tissues or pathological microenvironment [13]. Therefore, some genes that  
98 exhibit important topological properties on the interaction network, such as the hub  
99 genes [14], will be automatically selected as key regulators for disease state initiation  
100 and maintenance , leading to a possible increase in false positive master genes.  
101 Conversely, some classes of genes presenting as upstream regulators of a signaling  
102 cascade, such as the G protein-coupled receptors [15], may be identified as  
103 dispensable genes due to their relatively low degrees on the interactome, thus  
104 decreasing the sensitivity for distinguishing core ones from the giant pool of DAGs.

105 In this study, we have developed a network-oriented gene entropy approach to  
106 quantify the perturbation or regulatory ability of each DAG in distinct disease  
107 contexts by assembling and interrogating disease-specific regulatory networks. Master  
108 genes for each disease, whose altered expression was sufficient for disease state  
109 transitions, were identified as those genes that exhibited high entropy values by our *in*  
110 *silico* method, and were further adopted to investigate comorbidity and causal  
111 relationships among different diseases. We further confirmed that existing effective  
112 drugs are most likely to target the local module of master genes on the interactome.  
113 Using these methods, we have identified 11 old drugs as potent anticancer agents for  
114 pancreatic cancer treatment.

115

116

## 117 **Results and Discussion**

### 118 **Computation of gene entropy in disease networks**

119 To identify master genes in distinct disease contexts, a network-oriented gene entropy  
120 approach (NOGEA) was developed (Figure 1A and 1B). Briefly, Shannon entropy  
121 theory was applied to quantify the amount of disorder within intracellular signals in

each disease specific context, which was subsequently factorized as the summation of contribution for each DAG. First, directed disease specific gene networks for 293 diseases were constructed to reflect the distinct disease contexts by mapping all DAGs (Table S1) to a previously established directed PPI network (Table S2) [20]. A directed network visualizes the hierarchy of intracellular signal transduction between the interacting proteins, and hence clearly reflects the importance of each DAG in a certain physiological and pathological context. The regulation likelihood between each pair of DAGs was then calculated based on the directed distance on the PPI network to generate a probability-based signaling flux matrix (Figure 1A). Finally, the perturbation ability of each DAG in a disease-specific context was calculated by the network-oriented gene entropy metric (Methods, Figure 1B). The distribution of entropy values for all DAGs is illustrated as a histogram in Figure S1, and the perturbation ability of each DAG was then ranked based on their entropy values (Table S1).

To efficiently explore the biological features of each entropy distribution, all DAGs were classified as “Master”, “Interim” or “Redundant” genes which represent high, medium and low entropy genes, respectively. We created an entropy value curve for each disease and then identified two inflection points as thresholds to separate the low, medium and high entropy genes, respectively (Methods). We then merged the master genes of all diseases into a whole master gene set. Interim and redundant genes from different diseases were treated in the same way to obtain the whole interim and redundant gene sets, respectively. As a result, 798 master, 1,962 interim, and 1,387 redundant genes were obtained (Figure 1C, Table S3).

In order to verify whether the master genes play a key role in disease initiation and development, enrichment analyses were performed using several well-established gene clusters (Table S4). We observed that there was an overrepresentation ( $z$ -score=22.61) of disease-causing mutation-associated proteins among all master genes, which was higher than the enrichment score of both interim and redundant genes (Figure 1D). The essential genes were demonstrated to play critical roles in human diseases [28], and the master genes were enriched in essential genes, whose  $z$ -score was two times larger than the enrichment score of the redundant genes (Figure 1D). More importantly, we found that master genes were highly enriched in cancer-associated genes; whereas, redundant genes showed less enrichment (Figure 1D). Further KEGG analysis of the master genes showed that these genes were mainly enriched in pathways with close relationships with cancer initiation and progression

(Figure S2). For example, PI3K-AKT signaling pathway (has:04151), which is commonly perturbed in cancers, were found among the top five enriched pathways ( $P < 10e-30$ ). In a recent study, genes on the interactome were classified into different node types, in which “indispensable” nodes were found to be key players in mediating the transition of disease states. As shown in Figure S3A, we found that master genes were highly enriched in “indispensable” genes, but redundant genes were enriched among the “dispensable” genes. Consistent with these observations, the master genes were highly enriched in “critical” genes that acted as driver nodes in all control configurations (Figure S3B) [26]. Further dissection of all different functional classes within signaling proteins revealed that the master genes were most likely enriched in kinases and membrane receptors (Figure 1E). In summary, the results indicated that the master genes are preferred key regulators in disease initiation and development, reflecting the reliability of the NOGEA method.

Traditional network topology parameters, such as the connective degree and betweenness centrality, are commonly used as baseline methods for characterizing the importance of nodes in biological networks [29]. To validate the effectiveness of NOGEA, we compared it with four baseline methods (the connective degree, connective in-degree, connective out-degree and betweenness centrality-based methods) and four newly proposed methods (Katz [30], Catapult [30], HANRD [31] and GPS [32]), all of which are network-based methods for prioritizing disease genes. We first compared the AUROCs between different methods (Methods) and found that NOGEA significantly outperformed both the baseline methods and the newly proposed methods (Figure 1F). We further evaluated AUPRC, area under the precision-recall curve, for each method. NOGEA consistently surpassed all other methods, overmatching the second-best method by ~10% (Figure 1F).

Correlations between gene entropy values and four network topology parameters were assessed using Pearson's correlation coefficients (PCC). For most diseases, we observed that the PCCs between gene entropy values and network topology parameters were relatively small ( $<0.25$ , Figure S4A). Nonetheless, significant correlation values were observed between the in-degree connective ( $R^2=0.051$ ,  $P<1.0e-15$ , Figure 1G), out-degree connective ( $R^2=0.274$ ,  $P<1.0e-15$ , Figure 1H), degree connective (sum of in and out-degree,  $R^2=0.155$ ,  $P<1.0e-15$ , Figure 1I) and betweenness centrality ( $R^2=0.031$ ,  $P<1.0e-15$ , Figure 1J) for genes in the primary directed PPI network versus gene entropy values. Fisher's exact test was then applied to further determine whether gene entropy is associated with traditional network

topology parameters. Specifically, we constructed a contingency table to classify the disease genes into different bins based on their entropy values and network parameter values (Figure 1K). We found that gene entropy was significantly associated with traditional network topology parameters, including connective degree ( $P < 0.01$ ), connective in-degree ( $P < 0.01$ ), connective out-degree ( $P < 0.01$ ) and betweenness centrality ( $P < 0.01$ ). All these results demonstrate that master genes prefer to possess high topology parameter values, indicating relative consistency between gene entropy and the four network topology parameters.

To investigate variation of the regulatory role of a specific gene in different diseases, we calculated the divergence-degree of gene entropy across diseases using the coefficient of variation (CV) (Table S1, Figure S4B). The results show that up to 60% of the genes have a high CV ( $>15\%$ ), indicating the distinct roles these genes play in different disease contexts. We then examined the entropy value variation of the shared genes in different diseases, and observed that these genes usually exhibit similar entropy values in distinct diseases within the same disease category. For example, corticotropin-releasing hormone receptor 1 (CRHR1) is related to eight mental health-associated diseases with different entropy rank scores (rank  $>0.80$ ), including anxiety and depressive disorders (Table S1), which is consistent with its major role in mental disorders [33]. We also observed a low entropy rank score for CRHR1 in pulmonary disease (rank = 0.55), indicating variation in its regulatory role in distinct disease contexts. Further, we found that ~15% genes have approximately equal rank scores in their associated diseases. For instance, interleukin 4 receptor (IL4R) and phosphatidylinositol-4,5-bisphosphate 3-kinase catalytic subunit alpha (PIK3CA) had high rank scores in their associated diseases (Table S1), especially for neoplasms, suggesting crucial roles for these genes in these diseases. In summary, NOGEA provided a new way to explore the regulatory role of each DAG in distinct disease contexts.

## **NOGEA for exploring disease comorbidity**

Exploration of the underlying mechanisms of comorbidity, which refers to the coexistence of multiple diseases or disorders, is difficult due to complex interactions among environmental, lifestyle and treatment-related factors [34]. In addition, disease comorbidity includes not only the co-occurrence of multiple diseases, but also the potential cause-and-effect relationships among these diseases. Thus, uncovering the diseases' co-occurrence and causal relationships along with underlying mechanisms is



of great significance for their prevention and treatment. Using experiment-based approaches or mathematical models, previous studies explored the molecular features of disease comorbidity for several diseases, including from gastritis to gastric cancer [35] and from diabetes to cancer [36]. However, existing experiment-based methods to explore the underlying mechanisms for co-occurrence and causal relationships remain costly, labour-intensive, or focused on a small fraction of molecular features. Comparatively, mathematical models provide novel ways to reveal disease comorbidity using multi-omics data; however, these models are difficult to apply in other diseases, due to the lack of multi-scale information for these diseases.

The results discussed above demonstrate that NOGEA-inferred master genes are closely associated with disease onset and development, prompting us to investigate whether the network entropy-based approach would be capable of uncovering the molecular basis of disease co-occurrence. Therefore, we constructed a new master gene disease network (M-GDN), where edge would link two different diseases if they shared at least one master gene (Table S5). For comparison, we constructed five other disease networks: the redundant gene-based disease network (R-GDN), the interim gene-based disease network (I-GDN), the all genes-based disease network (A-GDN), the traditional hereditary disease network (THDN) and the random disease genes network (RGN).

To test whether the M-GDN would provide an accurate picture of disease comorbidity, we evaluated the Tanimoto similarity between these networks and the human disease comorbidity network (HDCN), which was extracted from the Medicare Claims Database and constructed in a recent study [3]. The M-GDN showed the highest similarity with the HDCN (higher than that of R-GDN and THDN) at a significantly higher level than expected based on the random values (Figure 2A), which indicates that those genes most associated with disease comorbidity tended to be master genes with high entropy rather than arbitrary disease genes. In contrast to previous THDN models, M-GDN considers genetic factors as well as genes that respond to environmental, lifestyle, and/or treatment-related factors, thus providing a more comprehensive solution for exploring the comorbidity of disease. Furthermore, in view of the impact of cellular network interactions on disease comorbidity, we extended our result to a PPI-based M-GDN (Table S6), where two diseases were linked if the master gene of one disease directly interacted with genes of the other disease in the PPI network. Consistent with the above results, the PPI-based M-GDN demonstrated the best predictive ability in identifying disease comorbidity. We then



observed that the inferred underlying molecular mechanisms of disease comorbidity are in accordance with current pathobiological knowledge (Figure 2B). For example, M-GDN confirmed the conclusion that AKT1 mutations lead to schizophrenia and type 2 diabetes mellitus (with rank scores of 0.96 and 0.94 in schizophrenia and type 2 DM, respectively) [37]. We also observed in the M-GDN that ADRB2 mutations may lead to asthma and obesity (with rank scores of 0.95 and 0.97 in asthma and obesity, respectively), which is consistent with a previous study [38]. These results suggest that M-GDN helps bridge the gap between bench-based biological discovery and bedside clinical solutions, and thus may provide new insights into the mechanisms of disease comorbidity.

Recent reports in the literature suggest that mutations in the IRS1 gene are closely related to the comorbidity of type 2 DM and obesity [39]. The M-GDN revealed that, in addition to IRS1, PTGS2 also plays a crucial role in the co-morbidities of these diseases. It is well known that PTGS2 influences the inflammatory response, which is closely connected with the comorbidity of type 2 DM and obesity [40]. Another example is the comorbidity of leukemia and cardiomyopathy, whose underlying mechanisms remain unclear. Interestingly, FAS is involved in the regulation of cell apoptosis, which affects left ventricular function [41] while PRKCA enhances cell resistance [42] and regulates cardiac contractility and an increased risk for heart failure. More importantly, the FAS-PRKCA interaction has been identified as the top connected cross-talk PPI by *in situ* proximity ligation assays [43]. These results demonstrate that the interaction between FAS and PRKCA may account for the comorbidity of leukemia and cardiomyopathy.

Next, we investigated the molecular basis of disease causal relationships from the perspective of directed biological networks. As an illustration, we constructed a directed comorbidity network (Table S7, Figure 2C) centered on Parkinson's disease. We observed high co-occurrence risk between Parkinson's and other diseases including Alzheimer's disease. Recent research suggests that these diseases are related to the accumulation of common proteins in the brain, such as alpha-synuclein protein [44]. Using alcoholism and Parkinson's disease as an example, we observed a significant directed interaction from alcoholism to Parkinson's disease ( $P < 0.01$ ), but not vice versa. This result is consistent with recent clinical studies, which suggest that alcoholism may be an inducer of Parkinson's disease [45]. A subsequent network analysis further discovered that the aberration of alcoholism master genes may lead to the modification of most Parkinson's disease's master genes (Figure 2D). Collectively,

NOGEA is potentially useful for investigating mechanisms underlying disease comorbidity as well as their causal relationships.

# **NOGEA can infer drug-disease associations**

Recently, several state-of-the-art network-based methods were proposed to investigate the relationships between drugs and diseases, such as the network proximity approach and network inference algorithm [4, 46]. In this study, we assessed relationships between DAGs and drug targets based on the gene network entropy to evaluate the effects of drugs on each disease. For each drug-disease relationship, we calculated the drug disturbance entropy (DDE) parameter, which represents potential therapeutic effects of the drug (Methods, Table S8-S10). To further investigate DDE's effectiveness, we evaluated the correlation between the DDE value and the hits by known drug-disease interactions (DDIs), and found the occurrence number of known DDIs increased with increasing DDE values (Figure 3A). Consistent with previous research [4], a highly significant correlation occurred between DDE values and the enrichment of known drug-disease interactions ( $R^2=0.75$ ,  $P=2.2e-16$ ) (Figure 3B), indicating a high likelihood that a drug will successfully treat a disease if the drug is capable of strongly perturbing the local module of master genes on the interactome.

To validate the utility of DDE for distinguishing known drug-disease pairs from the unknown drug-disease pairs, we compared the AUC of ROC curves for different drug-disease prediction methods (Methods). To obtain a robust AUC estimation, the drug-disease set was split into a training set and a testing set according to a given fraction coefficient for developing and validating the model, respectively. We compared the DDE's performance with several other state-of-the-art methods [4, 46], including the network inference algorithm (NIA), network proximity approach (NPA), network kernel approach (NKA), network shortest approach (NSA), network center approach (NCA), and network separation approach (NSEA). As illustrated in Fig. 3C, DDE exhibited the best performance (average AUROC=70%) in discriminating known and unknown drug-disease pairs, significantly outperforming the other approaches. Interestingly, we noticed that the NIA, which appeared to be the second-best method (average AUROC=68%), was also able to construct a directed disease-specific gene network and identify master genes before predicting the drug-disease associations. A compressive comparison between the two methods demonstrated their connection and difference (as seen in Supplementary Note 2,

Figure S5, Table S11-S12). Collectively, these results suggest that DDE is effective for predicting drug-disease associations.

Pancreatic cancer is a refractory malignant carcinoma of the digestive tract with a 5-year survival rate of ~4% [47] that modestly responds to very few existing chemotherapy treatment options. Revisiting the complex interaction pattern between drug targets and pancreatic cancer genes in a systemic manner is essential for developing more effective therapeutic regimens. Therefore, we used pancreatic cancer as an example to explore the utility of NOGEA for drug-disease association inference. By measuring the entropy of each pancreatic cancer gene in the pancreatic cancer specific network (Figure 3D, Figure 3E), we found that those genes with high entropy such as MET, KDR, and EGFR may play more important roles than the lower entropy genes for pancreatic cancer treatment. As reported in previous studies [48], EGFR-mediated signaling is involved in the tumorigenesis of pancreatic cancer, and the preclinical data support EGFR inhibition as a potential treatment strategy for pancreatic cancer. In addition, c-Met protein, which is coded by the MET gene, is a marker of pancreatic cancer stem cells and thus a therapeutic target [49]. KDR (VEGFR-2) is known to be crucial for embryonic vasculature development by modulating endothelial cell proliferation and migration [50]. Moreover, the CD44 gene is a potentially interesting prognostic marker and therapeutic target in pancreatic cancer [51].

To investigate differences in the targeting patterns between effective drugs and other less-effective drugs from a network-based perspective, we constructed a gene entropy map for pancreatic cancer. We first calculated the linkage strength between drug targets and pancreatic cancer genes for two FDA-approved drugs: Axitinib and Erythromycin (Figure 3D). Axitinib binds to FLT4, FLT1 and KDR, which was identified as a pancreatic cancer master gene by NOGEA. The DDE of Axitinib to pancreatic cancer is 37.6, suggesting that targets of Axitinib are more closely related to pancreatic cancer genes than expected by chance. Conversely, the DDE of Erythromycin (whose efficacy remains unknown) to pancreatic cancer is 1.1. Even though this drug inhibits ABCB1, ALB and KCNH2, the disease proteins and drug targets are not closer than expected by randomly selecting protein sets. However, some drugs that do not directly inhibit the pancreatic cancer master genes may still have the potential to be effective drugs. For example, Sirolimus, which is currently in phase II clinical trials, targets three proteins (FKBP1A, FGF2 and MTOR) but no known pancreatic cancer genes. Nevertheless, Sirolimus has a high DDE value of

12.1 due to the relatively strong perturbation of high entropy genes such as CD44 and EGFR via FGF2 (Figure 3E). Drugs such as Pravastatin (DDE=-0.7) are predicted to be ineffective pancreatic cancer drugs due to their weak perturbation of nearly all pancreatic cancer genes (Figure 3E). Collectively, these results suggest that NOGEA may be capable of identifying the core genes among many DAGs that provide the basis for rational drug discovery.

### **Pancreatic cancer drug screening**

Due to the encouraging performance of the drug disturbance entropy metric for accurately inferring drug-disease associations, we screened potentially effective drugs for pancreatic cancer treatment. We first calculated and prioritized DDE values for all FDA-approved drugs (Table S13-S14). From top 10% of these drugs, we selected 19 molecules that were not known to be associated with pancreatic cancer for further experimental validation. The half-maximal inhibitory concentration ( $IC_{50}$ ) of a molecule, an important metric to measure its response to certain cancer cell lines, has been widely applied in the screening of potential anti-proliferative agents in preclinical cancer pharmacogenomics. The BxPC3 human pancreatic cancer cell line, which has been frequently used in the study of pancreatic cancer and screening of chemo preventive agents [52], was used in our *in vitro* study to evaluate its response to the candidate drugs. We identified 11 candidate drugs that inhibit BxPC3 cell lines in a dose dependent manner and exhibit low  $IC_{50}$  values ( $<100 \mu\text{M/L}$ , Figure S6, Figure 4A-4C), demonstrating their efficacies for inhibiting pancreatic cancer cell proliferation and potential for pancreatic cancer therapy *in vivo*. One drug for example, Vinorelbine, is a drug that has already been approved for non-small-cell lung cancer treatment [53]. In our study, Vinorelbine exhibited a low  $IC_{50}$  value of  $1.55 \text{ nM/L}$  (Figure 4A). Conversely, some non-classical anticancer drugs also displayed acceptable suppressive effects on BxPC3. Additional drugs, including Saquinavir, which is mainly used with other medications for HIV/AIDS treatment or prevention [54], and Celecoxib, a drug mainly used for treatment of pain and inflammation in adults [55], showed  $IC_{50}$  values of  $22.63 \mu\text{M/L}$  (Figure 4B) and  $45.36 \mu\text{M/L}$  (Figure 4C), respectively. These results indicate that our model has the capacity to predict proper drug candidates for disease therapy.

Transcriptional expression analysis was conducted to validate our hypothesis that efficient drugs tend to perturb the master genes directly or through their targets. We first identified 1,335 differentially expressed genes (referred to as SAQDEGs) after

Saquinavir treatment (Figure S7A, Table S15). The pancreatic cancer master genes (n=849) that were most likely to be perturbed by Saquinavir were named SAQPEGs and further incorporated with their corresponding neighbor genes on the interactome (Table S15). Finally, a hypergeometric test was used to assess the overlap between SAQDEGs and SAQPEGs. These analyses revealed that the differentially expressed genes were significantly enriched for SAQPEGs (Figure 4D,  $P < 0.01$ ). Results for Celecoxib were similar to those for Saquinavir (Figure S7B, Figure 4E), suggesting a close relationship between genes perturbed by the efficient drugs and the local module of master genes.

Finally, to demonstrate the reliability of the DDE approach for extensive screening of pancreatic cancer candidate drugs, we conducted a literature mining analysis to evaluate the association between the candidate drugs (top 10%) and pancreatic cancer based on our previous reports [56] (Methods). We observed that 8 of the top 10 candidate drugs were anticancer agents that showed significant literature mining correlation scores with pancreatic cancer ( $P < 0.01$ , Table S14). In addition, most anticancer candidate drugs (~85%) were significantly associated with pancreatic cancer (Figure 4F, Table S14), suggesting the sensitivity of this model. Interestingly, an analysis of the categories of these candidate drugs revealed that the largest proportion, 44/224 (19.6%), were assigned to Central Nervous System Agents (CNSA). For example, Celecoxib, which was sensitive to the BxPC3 cell lines as mentioned above (Figure 4C), also acts as a CNSA. In general, these results indicate that DDE provides a rational strategy for drug repurposing due to its capacity to quantify drug targeting tendencies on the interactome.

## Materials and methods

### Data set collection

The DAGs for all diseases were obtained from four publicly available databases including KEGG Disease [16], Comparative Toxicogenomics Database [17], Therapeutic Target Database [18] and PharmGKB [19]. All disease names and their corresponding IDs were standardized by mapping to Medical Subject Headings ontology (MeSH; [www.nlm.nih.gov/mesh/](http://www.nlm.nih.gov/mesh/)) and official gene symbols for these DAGs were retrieved from GeneCards (<http://www.genecards.org/>). We then conducted a disease filtering process to ensure disease specificity. We first removed diseases with

levels  $< 2$  on the MeSH tree structures, such as “Nervous System Diseases” and “Cardiovascular Diseases”, as these disease types are too broad. Tanimoto similarity (ratio between the number of shared DAGs and the number of joined DAGs) was then computed for each disease pair and used to remove diseases showing high similarity ( $>0.50$ ) with its descendant disease. The weighted directed PPI network was constructed using data from a previous study [20], which consisted of 13,684 weighted interactions among 6082 proteins. The DAGs were then mapped to corresponding proteins in the PPI network, and those diseases with at least 20 DAGs in the human interactome were retained, for they are likely to induce a module on the network. As a result, we obtained 11,414 disease-gene associations between 274 diseases and 2848 protein-coding genes. For each disease, we manually extracted drug-disease associations from the drug indication information in DrugBank [21]. In addition, we obtained drug-target interactions for all FDA-approved drugs from DrugBank. To construct a disease comorbidity network, we retrieved disease pairs with comorbidity relationships from a recent study [3] of 665 diseases and their corresponding genes extracted from Online Mendelian Inheritance in Man (OMIM) [22].

## The disease-specific network-oriented gene entropy approach (NOGEA)

**Construction of a flux matrix based on the expectation of the Bernoulli distribution.** To construct the directed disease-specific gene networks, DAGs were mapped to the directed PPI network. For any given disease  $D$ , whose  $m$  associated genes can be mapped to the directed PPI network, an initial DAG vector  $V^{(D)} = \{V_1^{(D)}, \dots, V_i^{(D)}, \dots, V_m^{(D)}\}$  was generated to represent the disease, where  $V_i^{(D)}$  is the  $i$ -th DAG. The directed shortest path between two DAGs of disease  $D$  was calculated using the “igraph” package [23] based on the R 3.32 environment (r-project.org). For a given DAG pair  $V_i^{(D)}$  and  $V_j^{(D)}$ ,  $I_{(i,j)}$  is a random variable that obeys the Bernoulli distribution and represents the interaction or information transfer between node pair  $V_i^{(D)}$  to  $V_j^{(D)}$ . The distribution function of  $I_{(i,j)}$  is defined as

$$p(I_{(i,j)} = a; d_{(i,j)}, \omega) = (e^{-\omega * d_{(i,j)}})^a (1 - e^{-\omega * d_{(i,j)}})^{1-a} \quad (1)$$

where  $a = 1$  or  $0$ , indicating whether signal transduction exists between node pair  $V_i^{(D)}$  and  $V_j^{(D)}$ , and  $\omega$  is a scale parameter to adjust the likelihood for different distances. In addition,  $d_{(i,j)}$  is the directed distance between the given node pair  $V_i^{(D)}$



469 and  $V_j^{(D)}$ . It is the number of edges in a directed shortest path connecting them, and  
 470 was calculated using the “igraph” package based on Dijkstra's algorithm, reflecting  
 471 the possibility of the pairwise regulatory relationship from  $V_i^{(D)}$  to  $V_j^{(D)}$ . The details  
 472 for determining the optimal scale parameter are presented in Supplementary Note 1.  
 473 Therefore, the space of "possible" values assumed by  $I(i, j)$  is  $\{0, 1\}$ , and if  $a = 1$ ,  
 474  $p(a; d_{(i,j)}, \omega)$  represents the likelihood that there is a signaling flux between the node  
 475 pair. In the field of network communication, it is widely accepted that the success rate  
 476 of signal propagation decays exponentially with increasing distance [24]. In addition,  
 477 previous studies have demonstrated that exponential decay is a popular kernel to  
 478 characterize the network influence between two nodes [25]. Previously, we used the  
 479 exponential component to evaluate the association between two nodes in protein-  
 480 protein networks [26]. Thus, we believe that the success probability of the signal  
 481 transduction between two proteins decays exponentially with the increase of their  
 482 distance and the exponential component  $e^{-\omega * d_{(i,j)}}$  is useful for representing the  
 483 success probability. In this way, the stochastic information flux matrix for a given  
 484 disease is obtained by a simplified formula Eq. (2)

$$485 \quad P(I; d, \omega) = \{p(I_{(i,j)} = 1; d_{(i,j)}, \omega)\}_{(m \times m)} = \{e^{-\omega * d_{(i,j)}}\}_{(m \times m)} \quad (2)$$

486 And,  $p(I_{(i,j)} = 1; d_{(i,j)}, \omega)$  is equal to the expectation of  $I_{(i,j)}$ , where

$$487 \quad E(p(I_{(i,j)}; d_{(i,j)}, \omega)) = e^{-\omega * d_{(i,j)}} \quad (3)$$

488 The expectation was subsequently used to estimate the distribution of signaling  
 489 fluxing. For a given disease D with  $m$  associated genes, the biological signaling may  
 490 flux between any node pair (DAG)  $V_i^{(D)}$  and  $V_j^{(D)}$ . We then assumed that the edge (or  
 491 the node pair) through which the signals fluxes is a random variable  $F$ , and its event  
 492 space is

$$493 \quad \{f_{(i,j)} | 1 \leq i \leq m, 1 \leq j \leq m, i \neq j\} = \{f_{(1,2)}, \dots, f_{(i,j)}, \dots, f_{(m,m-1)}\} \quad (4)$$

494 where  $f_{(i,j)}$  represents signals that may be transferred from DAG  $V_i^{(D)}$  to  $V_j^{(D)}$ .

495 **Normalization of the fluxing matrix.** The probability distribution of signal  
 496 fluxing was estimated from

$$497 \quad p(F = f_{(i,j)}) = \frac{1}{Z} * E(p(I_{(i,j)}; d_{(i,j)}, \omega)) = \frac{1}{Z} * e^{-\omega * d_{(i,j)}} \quad (5)$$

498 where  $Z$  is the normalization constant or partition function, and

$$499 \quad Z = \sum_{i=1}^m \sum_{j=1, j \neq i}^m e^{-\omega * d_{(i,j)}} \quad (6)$$

500 to ensure that the sum of the probability is 1.



501

502 **Definition and calculation of disease gene entropy.** Based on the probability  
503 distribution of signal fluxing, we calculated the entropy for a given disease  $S^{(D)}$  in  
504 terms of the weighted Shannon entropy formula, which can be interpreted as the  
505 degree of disorder or complexity for the disease specific context,

$$506 \quad S^{(D)} = - \frac{\sum_{i=1}^m \sum_{j=1, j \neq i}^m p(f_{(i,j)}) * k_j^{out} \log p(f_{(i,j)})}{(m-1) \sum_{j=1}^m k_j^{out}} \quad (7)$$

507 where  $k_j^{out}$  is the out-degree of node  $V_j^{(D)}$  in the directed PPI network, which was  
508 calculated using the “igraph” package. Interestingly, we found that the disease  
509 entropy  $S^{(D)}$  can be factorized as shown in Eq. (8),

$$510 \quad S^{(D)} = \sum_{i=1}^m S_i^{(D)} \quad (8)$$

511 where  $S_i^{(D)}$  is the gene entropy of gene  $V_i^{(D)}$ , which is obtained by

$$512 \quad S_i^{(D)} = - \frac{\sum_{j=1, j \neq i}^m p(f_{(i,j)}) * k_j^{out} \log p(f_{(i,j)})}{(m-1) \sum_{j=1}^m k_j^{out}} \quad (9)$$

513 Therefore,  $S_i^{(D)}$  is a sub-entropy of disease entropy  $S^{(D)}$ , and is considered as the  
514 “disorder contribution” to a disease specific context.

515 **Gene entropy value normalization.** Through the above procedure, a gene  
516 entropy map was established for 293 diseases. For any given disease  $D$ , the gene  
517 entropy z-scores were calculated, making the gene entropy values of different diseases  
518 comparable,

$$519 \quad ZS_i^{(D)} = \frac{S_i^{(D)} - \mu(S_i^{(D)})}{\delta(S_i^{(D)})} \quad (10)$$

520 where  $\mu(S_i^{(D)})$  and  $\delta(S_i^{(D)})$  are the estimation of the expectation and standard  
521 deviations of  $S_i^{(D)}$  for disease  $D$ . In addition, to assess the disturbance capability of a  
522 gene in a disease-specific network in a more intuitive manner, we calculated the rank  
523 score for all DAGs according to their entropy values, which range from 0 to 1 and  
524 reflect their likelihood as master genes.

525 **Rank score calculation of gene entropy.** The gene entropy values for disease  $D$   
526 were sorted in an ascending order, and a rank list was generated:

$$527 \quad RL^{(D)} = \{rl(S_1^{(D)}), \dots, rl(S_i^{(D)}), \dots, rl(S_m^{(D)})\} \quad (11)$$

528 where the  $rl(S_i^{(D)})$  is the rank value of  $S_i^{(D)}$ . Note that those genes that possess equal  
529 entropy values have the same rank values. For example, if there are  $k$  genes

530  $\{V_{i+1}^{(D)}, \dots, V_{i+k}^{(D)}\}$  possessing equal entropy values  $\{S_{i+1}^{(D)}, \dots, S_{i+k}^{(D)}\}$ , their rank values  
531 were determined by equation (12):

$$532 \quad rl(S_{i+1}^{(D)}) = \dots = rl(S_{i+k}^{(D)}) = \frac{\sum_{j=1}^k po(S_{i+j}^{(D)})}{k} \quad (12)$$

533 where  $po(S_{i+j}^{(D)})$  is the position of  $S_{i+j}^{(D)}$  in the ascending entropy value list. Based on  
534 the rank list, rank score vector  $RS^{(D)}$  was generated by Eq. (13):

$$535 \quad RS^{(D)} = \left\{ \frac{rl(S_i^{(D)}) - \min(RL(S^{(D)}))}{\max(RL^{(D)}) - \min(RL^{(D)})} \right\}_{(1 \times m)} \quad (13)$$

536 where  $\max(RL^{(D)})$  and  $\min(RL^{(D)})$  are the maximum and minimum of  $RL^{(D)}$ ,  
537 respectively.

538

539 **Disease-gene classification based on the gene entropy value.** To  
540 comprehensively explore the biological meaning of the entropy, we divided all DAGs  
541 into three groups based on their entropy values using an adaptive approach. Briefly,  
542 we created an entropy value curve for each disease, and identified two inflection  
543 points in the curve as thresholds. Specifically, for each disease D, we ranked each  
544 gene entropy value ( $S_i^{(D)}$ ) in ascending order. Then we mapped each entropy value  
545 onto a two-dimensional coordinate system such that the lowest entropy value ( $S_1^{(D)}$ )  
546 became coordinate  $(1, S_1^{(D)})$ , the second lowest value became  $(2, S_2^{(D)})$ , and so on, until  
547 the maximum entropy value ( $S_{max}^{(D)}$ ) was reached. Two inflection points, individually  
548 defined as the threshold points of most rapid increase from the low to the medium and  
549 from the medium to the high entropy values, were identified in the entropy value  
550 curve from the interval of 10th to 50th percentile and 51st to 90th percentile,  
551 respectively, of all entropy values. The entropy value corresponding to this threshold  
552 was used as an adaptive disease-specific classification threshold. Master genes of all  
553 diseases were then merged and adopted as the whole master gene set to explore their  
554 common biological meanings. Interim and redundant genes from different diseases  
555 were treated in the same way to obtain the whole interim and redundant gene sets,  
556 respectively. Therefore, some genes may belong to all three gene sets (master, interim  
557 and redundant), because they play different roles in distinct disease contexts.

558

559

560 **Disease comorbidity relationship evaluation**

A real human disease comorbidity network (HDCN) was constructed in which nodes represented diseases and edges represented the reported comorbidity relationships, respectively. We then built five different types of inferred disease comorbidity networks to compare with the HDCN. First, a master gene disease network (M-GDN) was constructed, where edges linked two different diseases only if they shared at least one high entropy gene. We then constructed the redundant gene disease network (R-GDN), the interim gene disease network (I-GDN), the whole genes-based disease network (A-GDN) and the traditional hereditary disease network (THDN), respectively. A Tanimoto coefficient was used to evaluate the similarity between different networks as shown in Eq. (15),

$$T(A, B) = \frac{|E(A) \cap E(B)|}{|E(A)| + |E(B)| - |E(A) \cap E(B)|} \quad (15)$$

where  $A$  and  $B$  are different networks,  $E(\cdot)$  represents the edge set of a given network and  $|E(\cdot)|$  is the number of edges in the net. To assess the significance of the similarity of different networks, the random disease genes network was randomly generated 1,000 times and compared with the HDCN using equation (15). In the random disease genes network, each disease involves a random sampling gene set of the same size as the disease in A-GDN.

Previous research has demonstrated that cellular interaction links result in statistically significant comorbidity patterns [3]. Therefore, we believe that the directed interaction strength from the DAGs of one disease to another in the directed cellular network can reflect the causal relationship between the two diseases. To evaluate whether a causal relationship exists between two diseases, we estimated the significance of the interaction strength between the DAGs of the disease pairs using the Monte Carlo method. We first defined a raw causal relationship score (RCRS) for two given diseases: D1 and D2,

$$RCRS(D1 \rightarrow D2) = \sum_{i \in D1, j \in D2} p(I_{(i,j)}; d_{(i,j)}) * \varphi(p(I_{(i,j)}; d_{(i,j)})) \quad (16)$$

where  $p(I_{(i,j)}; d_{(i,j)})$  was calculated by equation (1),  $d_{(i,j)}$  is the directed distance between master gene pair  $V_i^{(D1)}$  and  $V_j^{(D2)}$ , and  $\varphi(p(I_{(i,j)}; d_{(i,j)}))$  is an indicator function. In addition,  $\varphi(p)$  was calculated as

$$\varphi(p) = \begin{cases} 1, & p \geq p_{cut} \\ 0, & p < p_{cut} \end{cases} \quad (17)$$

where  $p_{cut}$  is a threshold, below which the probability was discarded and considered not contributive to the overall interaction and  $p_{cut}$  was determined according to a

previous study [27]. We then used a normalized causal relationship score (NCRS) to quantify the risk that disease  $D1$  will induce disease  $D2$ . The  $NCRS$  is defined in Eq. (18)

$$NCRS(D1 \rightarrow D2) = \frac{RCRS(D1 \rightarrow D2) - \mu(RCRS(D1 \rightarrow D2))}{\delta(RCRS(D1 \rightarrow D2))} \quad (18)$$

where  $\mu(RCRS(D1 \rightarrow D2))$  and  $\delta(RCRS(D1 \rightarrow D2))$  are the estimation of the expectation and standard deviations of  $RCRS$  under the same condition, respectively. Then, Monte Carlo simulation was performed 1,000 times to estimate the  $\mu(RCRS(D1 \rightarrow D2))$  and  $\delta(RCRS(D1 \rightarrow D2))$  by randomly sampling the same number of genes as  $D1$  and  $D2$ . In each simulation, the values, the average and standard deviations of  $RCRS$  were calculated. To assess whether the causal relationship from disease  $D1$  to  $D2$  was significant, the P-value of  $RCRS(D1 \rightarrow D2)$  was further calculated as shown in Eq. (19):

$$p(RCRS(D1 \rightarrow D2)) = \frac{n_{RCRS(random) > RCRS(D1 \rightarrow D2)} + 1}{N_{total} + 1} \quad (19)$$

where  $N_{total}$  is the total number of simulations, and  $n_{RCRS(random) > RCRS(D1 \rightarrow D2)}$  is the number of random  $RCRS$  values that are larger than  $RCRS(D1 \rightarrow D2)$ . The  $RCRS$  value for the significance of P-values was set to 0.01. Finally, for a disease pair  $D1$  and  $D2$ , if both  $RCRS(D1 \rightarrow D2)$  and  $RCRS(D2 \rightarrow D1)$  were significant ( $P < 0.01$ ), the two diseases were considered to be co-occurrent; whereas, if only one was significant ( $P < 0.01$ ), we determined that a causal relationship exists between the two diseases.

# Drug disturbance entropy (DDE)

To quantify the effects of a drug on each disease based on the gene network entropy, we applied an ensemble approach, referred to as drug disturbance entropy (DDE), to evaluate the relationship between drug targets and disease proteins (encoded by disease genes) on the interactome. We first evaluated the linkage strength between each DAG and drug target on the interactome, which was then transformed to a probability. The perturbation value for each target and DAG was defined as the product of the strength probability and the DAG entropy,

$$T_{(t,i)} = p(I_{(t,i)} = 1; d_{(t,i)}) * S_i \quad (20)$$

where  $p(I_{(t,i)} = 1; d_{(t,i)})$  represents the strength probability between drug target  $t$  and DAG  $V_i^{(D)}$ ,  $S_i$  is the entropy value of DAG  $V_i^{(D)}$ , and  $d_{(t,i)}$  is the distance between

target  $t$  and DAG  $V_i^{(D)}$ . The raw disturbance entropy, which represents an estimate of a drug's therapeutic effects through distinct targets, was defined as

$$ET(T, V^{(D)}) = \sum_{t \in T, i \in G} T_{(t,i)} * \varphi(T_{(t,i)}) \quad (21)$$

where  $T_{(t,i)}$  is the perturbation entropy between target  $t$  and DAG  $V_i^{(D)}$ , and  $\varphi(T_{(t,i)})$  is an indicator function as shown in Eq. (22)

$$\varphi(T_{(t,i)}) = \begin{cases} 1, & T_{(t,i)} \geq T_{cut} \\ 0, & T_{(t,i)} < T_{cut} \end{cases} \quad (22)$$

where  $T_{cut}$  is a cut-off threshold of the disturbance entropy. The threshold of the perturbation value was determined by extensive sampling, and relationships with a perturbation value below this threshold were discarded. The remaining values were summed as the raw DDE of the drug to the disease. The advantage of this procedure is that weak relationships are eliminated, which greatly reduces noise and improves the robustness of the measure. By sampling across the range of  $T_{cut}$  choices, the threshold that led to the highest ROC AUC was chosen. We obtained the proper  $T_{cut}$  as  $0.89 * \max(T_{(t,i)})$  by evaluating the performance of predictions of drug-disease associations. Detailed information for determining  $T_{cut}$  is depicted in Supplementary Note 1.

To avoid possible high DDE that may be caused by a large number of drug targets and DAGs, we converted raw DDE to a size-bias-free value using the mean and standard deviation of raw DDE modeled from sets of random molecules, so that the potential therapeutic effects between distinct drugs and diseases could be evaluated under the same metric. The raw drug disturbance entropy was transformed to a size-bias-free score under formula (23)

$$ET^*(T, V^{(D)}) = \frac{ET(T, V^{(D)}) - \mu(ET(T, V^{(D)}))}{\delta(ET(T, V^{(D)}))} \quad (23)$$

where  $T$  and  $V^{(D)}$  are the drug target set and the disease-associated gene set respectively;  $\mu(ET(T, V^{(D)}))$  and  $\delta(ET(T, V^{(D)}))$  are the estimation of the expectation and standard deviations of DDE under this condition, respectively.

The estimation procedure of  $\mu(ET(T, V^{(D)}))$  and  $\delta(ET(T, V^{(D)}))$  are as follows: For each pair of  $(T, V^{(D)})$ , we constructed 1,000 random set pairs with  $|T|$  targets and  $|V^{(D)}|$  DAGs, preserving the degree distribution of the randomized targets and disease proteins. To avoid repeatedly choosing the same nodes during the degree-preserving random selection, we used a binning approach as described in a previous report [4].

657

658

659

## 660 **Conclusion**

661 Disease phenotypes typically result from interactions among multiple complex  
 662 environmental and genetic factors. The occurrence, development and treatment of a  
 663 disease usually involves hundreds of genes [29]. Presently, we proposed a network-  
 664 oriented gene entropy approach (NOGEA) for accurately inferring master genes that  
 665 contribute to specific diseases by quantitatively calculating their perturbation abilities  
 666 on directed disease-specific gene networks. Our results confirm that that master genes  
 667 are enriched in gene sets that account for disease onset and development. This may  
 668 imply that at a molecular level, those master genes with high entropy values are the  
 669 underlying start-points of the disease state, impacting those redundant genes with low  
 670 entropy through a directed disease-specific gene network. Interestingly, the  
 671 comorbidity prediction model built using the master genes showed the best agreement  
 672 with the independent clinical data set compared to the model established using the  
 673 whole disease gene set. This indicates that our method may decrease the influence of  
 674 noise and improve the efficiency for extracting more important genes from massive  
 675 genomic data sets. Finally, through this method, 11 old drugs were newly identified  
 676 and predicted to be effective for treating pancreatic cancer and then validated by *in*  
 677 *vitro* experiments. However, it remains challenging to simulate the complex contents  
 678 of the tumor microenvironment *in vitro*, making it difficult to comprehensively  
 679 evaluate drug response using IC<sub>50</sub>. Therefore, despite our encouraging results, future  
 680 work focusing on *in vivo* validation before clinical use is needed.

681 Although the identified master genes may be important for elucidating  
 682 mechanisms of disease progression and drug screening, we acknowledge that it is  
 683 difficult to directly evaluate the accuracy of NOGEA for identifying master genes at  
 684 this stage due to the lack of ‘gold standard’ reference data sets. Nevertheless, the  
 685 availability of more personal genome data in the future will allow for construction of  
 686 patient-specific networks, NOGEA will provide new opportunities to identify patient-  
 687 specific master genes and promote the development of personalized medicine.  
 688 Emerging deep learning methods may become powerful techniques for exploring  
 689 poly-pharmacy side effects [57] and discovering disease gene associations [58] from  
 690 massive data sets [59]. Because gene entropy values can be used as novel disease  
 691 feature data, we expect that integrating deep learning with NOGEA will significantly

improve the accuracy for determining disease-drug or disease-disease associations.  
Extending the systematic approach presented here from signal drugs to multiple drugs  
may pave the way toward a better understanding of drug combinations.

## **Authors' contributions**

YHW and WX formulated the idea of the paper and supervised the research. ZHG and CH performed the research and drew the figures. XTC and SG collected data. ZYW and XTC performed laboratory experiments. ZHG, YXF, CH and CLZ wrote the paper. YXF, MS, PFT, YHM, JBZ, YL and ZHW revised the paper. All authors reviewed the manuscript.

## **Competing interests**

The authors have declared no competing interests.

## **Acknowledgements**

The research was supported by the National Natural Science Foundation of China (NO. U1603285) and the National Science and Technology Major Project of China (grant number 2019ZX09201004-001).

We thank TopEdit ([www.topeditsci.com](http://www.topeditsci.com)) for its linguistic assistance during the preparation of this manuscript.



# References

- [1] Jimenez-Sanchez G, Childs B, Valle D. Human disease genes. *Nature* 2001;409:853-855.
- [2] Debret G, Jung C, Hugot JP, Pascoe L, Victor JM, Lesne A. Genetic susceptibility to a complex disease: the key role of functional redundancy. *Hist Phil Life Sci* 2011;33:497-514.
- [3] Park J, Lee DS, Christakis NA, Barabási AL. The impact of cellular networks on disease comorbidity. *Mol Syst Biol* 2009;5:262.
- [4] Guney E, Menche J, Vidal M, Barabasi A. Network-based in silico drug efficacy screening. *Nat Commun* 2016;7:10331.
- [5] Todorovic M, Newman JR, Shan J, Bentley S, Wood SA, Silburn PA, et al. Comprehensive assessment of genetic sequence variants in the antioxidant 'master regulator' nrf2 in idiopathic Parkinson's disease. *PLoS One* 2015;10:e0128030.
- [6] Thomas K. High-throughput gene expression and mutation profiling: current methods and future perspectives. *Breast Care* 2013;8:401-406.
- [7] Alvarez MJ, Shen Y, Giorgi FM, Lachmann A, Ding BB, Ye BH, et al. Functional characterization of somatic mutations in cancer using network-based inference of protein activity. *Nat Genet* 2016;48:838.
- [8] Walsh LA, Alvarez MJ, Sabio EY, Reyngold M, Makarov V, Mukherjee S, et al. An integrated systems biology approach identifies TRIM25 as a key determinant of breast cancer metastasis. *Cell Rep* 2017;20:1623-1640.
- [9] West J, Bianconi G, Severini S, Teschendorff AE. Differential network entropy reveals cancer system hallmarks. *Sci Rep* 2012;2:802.
- [10] Reilly MT, Cunningham KA, Natarajan A. Protein-protein interactions as therapeutic targets in neuropsychopharmacology. *Neuropsychopharmacology* 2013;34:247-248.
- [11] Porta-Pardo E, Garcia-Alonso L, Hrabe T, Dopazo J, Godzik A. A Pan-cancer catalogue of cancer driver protein interaction interfaces. *PLoS Comput Biol* 2015;11:e1004518.
- [12] Vidal M, Cusick ME, Barabási A-L. Interactome networks and human disease. *Cell* 2011;144:986-998.
- [13] Greene CS, Arjun K, Wong AK, Emanuela R, Zelaya RA, Himmelstein DS, et al. Understanding multicellular function and disease with human tissue-specific networks. *Nat Genet* 2015;47:569-576.
- [14] Allen JD, Yang X, Min C, Luc G, Guanghua X. Comparing statistical methods for constructing large scale gene networks. *PLoS One* 2012;7:e29348.
- [15] Enrique R. Mitogenic signaling pathways induced by G protein-coupled receptors. *J Cell Physiol* 2010;213:589-602.
- [16] Kanehisa M, Goto S, Sato Y, Kawashima M, Furumichi M, Tanabe M. Data, information, knowledge and principle: back to metabolism in KEGG. *Nucleic Acids Res* 2014;42:199-205.
- [17] Davis AP, Grondin CJ, Johnson RJ, Sciaky D, King BL, McMorran R, et al. The comparative toxicogenomics database: update 2017. *Nucleic Acids Res* 2017;45:D972-D978.
- [18] Zhu F, Shi Z, Qin C, Tao L, Liu X, Xu F, et al. Therapeutic target database update 2012: a resource for facilitating target-oriented drug discovery. *Nucleic Acids Res* 2012;40:D1128-D1136.

763 [19] Whirlcarrillo M, Mcdonagh EM, Hebert JM, Gong L, Sangkuhl K, Thorn CF, et  
764 al. Pharmacogenomics knowledge for personalized medicine. *Clin Pharmacol Ther*  
765 2012;92:414-417.

766 [20] Vinayagam A, Stelzl U, Foulle R, Plassmann S, Zenkner M, Timm J, et al. A  
767 directed protein interaction network for investigating intracellular signal transduction.  
768 *Sci Signal* 2011;4:rs8.

769 [21] Wishart DS, Feunang YD, Guo AC, Lo EJ, Marcu A, Grant JR, et al. DrugBank  
770 5.0: a major update to the DrugBank database for 2018. *Nucleic Acids Res* 2017;46.

771 [22] McKusick, Victor A. Mendelian Inheritance in Man and its online version,  
772 OMIM. *Am J Hum Genet* 2007;80:588-604.

773 [23] Csardi G, Nepusz T. The igraph software package for complex network research.  
774 *Int J* 2006;complex systems.

775 [24] Takeoka M, Guha S, Wilde MM. Fundamental rate-loss tradeoff for optical  
776 quantum key distribution. *Nat Commun* 2015;5:5235.

777 [25] Cohen E, Dellling D, Pajor T, Werneck RF. Distance-based influence in  
778 networks: computation and maximization. *Comput Sci* 2014.

779 [26] Li P, Huang C, Fu Y, Wang J, Wu Z, Ru J, et al. Large-scale exploration and  
780 analysis of drug combinations. *Bioinformatics* 2015;31:2007.

781 [27] Keiser MJ, Roth BL, Armbruster BN, Ernsberger P, Irwin JJ, Shoichet BK.  
782 Relating protein pharmacology by ligand chemistry. *Nat Biotechnol* 2007;25:197-206.

783 [28] Dickerson JE, Zhu A, Robertson DL, Hentges KE. Defining the Role of Essential  
784 Genes in Human Disease. *PLoS One* 2011;6:e27368.

785 [29] Albert-László B, Natali G, Joseph L. Network medicine: a network-based  
786 approach to human disease. *Nat Rev Genet* 2011;12:56-68.

787 [30] Singh-Blom UM, Natarajan N, Tewari A, Woods JO, Dhillon IS, Marcotte EM.  
788 Prediction and validation of gene-disease associations using methods inspired by  
789 social network analyses. *PLoS One* 2013;8.

790 [31] Rao A, Vg S, Joseph T, Kotte S, Sivadasan N, Srinivasan R. Phenotype-driven  
791 gene prioritization for rare diseases using graph convolution on heterogeneous  
792 networks. *BMC Med Genomics* 2018;11:57.

793 [32] Copps KD, White MF. Regulation of insulin sensitivity by serine/threonine  
794 phosphorylation of insulin receptor substrate proteins IRS1 and IRS2. *Diabetologia*  
795 2012;55:2565-2582.

796 [33] Rogers J, Raveendran M, Fawcett GL, Fox AS, Shelton SE, Oler JA, et al.  
797 CRHR1 genotypes, neural circuits and the diathesis for anxiety and depression. *Mol*  
798 *Psychiatry* 2013;18:700-707.

799 [34] Lee D-S, Park J, Kay KA, Christakis NA, Oltvai ZN, Barabási A-L. The  
800 implications of human metabolic network topology for disease comorbidity. *P Natl*  
801 *Acad Sci USA* 2008;105:9880-9885.

802 [35] Guo Y, Nie Q, MacLean AL, Li Y, Lei J, Li S. Multiscale modeling of  
803 inflammation-induced tumorigenesis reveals competing oncogenic and oncoprotective  
804 roles for inflammation. *Cancer Res* 2017;77:6429-6441.

805 [36] Wu D, Hu D, Chen H, Shi G, Fetahu IS, Wu F, et al. Glucose-regulated  
806 phosphorylation of TET2 by AMPK reveals a pathway linking diabetes to cancer.  
807 *Nature* 2018;559:637-641.

808 [37] Liu Y, Li Z, Zhang M, Deng Y, Yi Z, Shi T. Exploring the pathogenetic  
809 association between schizophrenia and type 2 diabetes mellitus diseases based on  
810 pathway analysis. *BMC Med Genomics* 2013;6 Suppl 1:S17.

811 [38] Danielewicz H. What the genetic background of individuals with asthma and  
812 obesity can reveal: ss  $\beta$ 2-adrenergic receptor gene polymorphism important? *Pediatr*  
813 *Allergy Immunol Pulmonol* 2014;27:104.

814 [39] Schmitz-Peiffer C, Whitehead J. IRS-1 regulation in health and disease. *IUBMB*  
815 *Life* 2003;55:367-374.

816 [40] Bastard J, Maachi M, Lagathu C, Kim MJ, Caron M, Vidal H, et al. Recent  
817 advances in the relationship between obesity, inflammation, and insulin resistance.  
818 *Eur Cytokine Netw* 2006;17:4-12.

819 [41] Sheppard R, Bedi M, Kubota T, Semigran MJ, Dec W, Holubkov R, et al.  
820 Myocardial expression of fas and recovery of left ventricular function in patients with  
821 recent-onset cardiomyopathy. *J Am Coll Cardiol* 2005;46:1036-1042.

822 [42] Ruvoilo PP, Deng X, Carr BK, May WS. A functional role for mitochondrial  
823 protein kinase  $\text{Ca}$  in Bcl2 phosphorylation and suppression of apoptosis. *J Biol Chem*  
824 1998;273:25436-25442.

825 [43] Chen T, Lin K, Chen C, Lee S, Lee P, Liu Y, et al. Using an in situ proximity  
826 ligation assay to systematically profile endogenous protein-protein interactions in a  
827 pathway network. *J Proteome Res* 2014;13:5339-5346.

828 [44] Guo JL, Covell DJ, Daniels JP, Iba M, Stieber A, Zhang B, et al. Distinct  $\alpha$ -  
829 synuclein strains differentially promote tau inclusions in neurons. *Cell* 2013;154:103-  
830 117.

831 [45] Bettiol SS, Rose TC, Hughes CJ, Smith LA. Alcohol consumption and  
832 parkinson's disease risk: a review of recent findings. *J Parkinsons Dis* 2015;5:425-  
833 442.

834 [46] Zickenrott S, Angarica VE, Upadhyaya BB, del Sol A. Prediction of disease-  
835 gene-drug relationships following a differential network analysis. *Cell Death Dis*  
836 2016;7:e2040-e2040.

837 [47] Sivakumar S, De SI, Chlon L, Markowitz F. Master regulators of oncogenic  
838 KRAS response in pancreatic cancer: an integrative network biology analysis. *PLoS*  
839 *Med* 2017;14:e1002223.

840 [48] Kelley RK, Ko AH. Erlotinib in the treatment of advanced pancreatic cancer.  
841 *Biologics: Targets & Therapy* 2008;2:83-95.

842 [49] Li C, Wu JJ, Hynes M, Dosch J, Sarkar B, Welling TH, et al. c-Met is a marker  
843 of pancreatic cancer stem cells and therapeutic target. *Gastroenterology*  
844 2011;141:2218.

845 [50] Korc M. Pathways for aberrant angiogenesis in pancreatic cancer. *Mol Cancer*  
846 2003;2:8-8.

847 [51] Li X, Zhang X, Zheng L, Guo W. Expression of CD44 in pancreatic cancer and  
848 its significance. *Int J Clin Exp Pathol* 2015;8:6724-6731.

849 [52] De Soto JA, Mullins R. The use of PARP inhibitors as single agents and as  
850 chemosensitizers in sporadic pancreatic cancer. *J Clin Oncol Off J Am Soc Clin*  
851 *Oncol* 2011;29:e13542.

852 [53] Listed N. Effects of vinorelbine on quality of life and survival of elderly patients  
853 with advanced non-small-cell lung cancer. The Elderly Lung Cancer Vinorelbine  
854 Italian Study Group. *J Natl Cancer Inst* 1999;91:66-72.

855 [54] Merry C, Barry MG, Mulcahy F, Tjia JF, Halifax KL, Heavey J, et al. Ritonavir  
856 pharmacokinetics alone and in combination with saquinavir in HIV-infected patients.  
857 *AIDS* 1998;12:325-327.

858 [55] Tindall E. Celecoxib for the treatment of pain and inflammation: the preclinical  
859 and clinical results. *J Am Osteopath Assoc* 1999;99:S13-17.

860 [56] Lounkine E, Keiser MJ, Whitebread S, Mikhailov D, Hamon J, Jenkins J, et al.  
861 Large scale prediction and testing of drug activity on side-effect targets. *Nature*  
862 2012;486:361.

863 [57] Li Y, Huang C, Ding L, Li Z, Gao X. Deep learning in bioinformatics:  
864 introduction, application, and perspective in the big data era. *Methods* 2019.

865 [58] Zitnik M, Agrawal M, Leskovec J. Modeling polypharmacy side effects with  
866 graph convolutional networks. *Bioinformatics* 2018;34:i457-i466.  
867 [59] Li Y, Kuwahara H, Yang P, Song L, Gao X. PGCN: Disease gene prioritization  
868 by disease and gene embedding through graph convolutional neural networks.  
869 *bioRxiv* 2019:532226.

870

871

## 872 **Figure legends**

### 873 **Figure 1 Computation and characterization of network-oriented gene entropy in** 874 **disease-specific networks**

875 A. Construction of directed disease-specific gene networks by mapping disease  
876 genes to the directed PPI network and normalizing the interaction strength. **B.**  
877 Calculation of the perturbation ability (gene entropy) of each gene. **C.** The Venn plot  
878 of the disease gene from different classes; Master: the master genes, Interim: the  
879 interim genes, Redundant: the redundant genes. **D.** Enrichment result (z-score) of  
880 master, interim and redundant genes in the context of OMIM, cancer and essential  
881 genes. **E.** Enrichment result (z-score) of master, interim and redundant entropy genes  
882 in the context of kinase, membrane receptor (MR), transcription factor (TF). **F.**  
883 Comparison of NOGEA performance with other methods for disease gene  
884 prioritization using AUROC and AUPRC. **G.** DAG entropy values versus their in-  
885 degree in the primary directed PPI network. **H.** DAG entropy values versus their out-  
886 degree in the primary directed PPI network. **I.** DAG entropy values versus their  
887 betweenness in the primary directed PPI network. **J.** DAG entropy values versus their  
888 degree (sum of in- and out-degree) in the primary directed PPI network. **K.**  
889 Assessment of the association between gene entropy and four commonly used  
890 network topology parameters.

891

### 892 **Figure 2 Exploration of disease comorbidity using network entropy**

893 A. Distribution of Tanimoto similarities between HDCN and other disease-disease  
894 networks (M-GDN, I-DGN, R-DGN, A-DGN, THDN and RGN). **B.** The inferred  
895 molecular basis of disease comorbidity relationships. Brown and blue nodes represent  
896 master genes inferred by NOGEA; green nodes represent diseases. **C.** The  
897 comorbidity of Parkinson's disease. In this figure, the width of the edge represents the  
898 likelihood of disease comorbidity, arrows represent the inferred causative disease-  
899 disease associations, and the color of the nodes depicts the disease category from  
900 MESH. **D.** The molecular basis of the comorbidity between Parkinson's disease and

alcoholism. The nodes represent the master genes of the disease and the directed links describe the direction from the directed PPI network.

### **Figure 3 Drug-disease association inference based on the disease gene entropy**

A. The hits number by known DDIs in each ranked drug-disease pair bin. B. The correlation between average DDE score in each bin and the hits enrichment fold for known DDIs. C. AUROC for drug-disease predictions using different methods. D. The interaction between drug targets and pancreatic cancer genes. The width of the links, the shade of the pancreatic cancer genes nodes, and the size of the node describe the interaction strength, entropy value, and degree of each node in the human interactome, respectively. E. The entropy value rank plot of pancreatic cancer genes (right); the heat map describes the shortest distance between the drug targets and pancreatic cancer genes of four drugs (left).

### **Figure 4 Screening of potential efficient drugs for pancreatic cancer treatment**

A-C. Cell inhibition rate curves against BxPC3 for Vinorelbine, Saquinavir and Celecoxib, respectively. D. The number and significance of overlapped genes between differentially expressed genes and the inferred effect genes after Saquinavir treatment. E. The number and significance of overlapped genes between the differentially expressed genes and the inferred effect genes after Celecoxib treatment. F. The overlapped drug number between each category and the top 10% of efficient drugs. Red bar: number of literature mining significant drugs; AIA: Anti-Inflammatory Agents, AIANS: Anti-Inflammatory Agents (Non-Steroidal); ANA: Antineoplastic Agents; ANIA: Antineoplastic and Immunomodulating Agents; ARA: Antirheumatic Agents; CVA: Cardiovascular Agents; CNSA: Central Nervous System Agents; HTA: Hypotensive Agents; PNSA: Peripheral Nervous System Agents; SSA: Sensory System Agents.

## **Supplementary material**

### **Figure S1 Distribution of gene entropy values for all DAGs**

Histogram plots showing the distribution of gene entropy values for all DAGs before (left) and after (right) normalization. The x-axis shows the range of gene entropy values, and the y-axis shows the count of genes possessing different entropy values.

### **Figure S2 KEGG pathway enrichment results**

936 X-axis: the top 20 significantly enriched 'KEGG pathway terms' of the master genes;  
937 y-axis: significance of the enrichment  $[-\log(P\text{-value})]$ .

938

### 939 **Figure S3 The disease-gene enrichment analysis for different classifications**

940 Enrichment results (z-score) of master, interim and redundant genes in the context of  
941 gene sets for critical (**A**), redundant (**A**), indispensable (**B**) and dispensable (**B**) genes.

942

### 943 **Figure S4 The property of the disease-gene entropy concept**

944 A. The correlation between entropy value and topology property for each disease. In  
945 this figure, each point represents a disease. The coordinate of each point represents  
946 the Pearson's correlation coefficient (PCC) for the gene entropy values versus the  
947 in-degree (x-axis) and the out-degree (y-axis) of the disease-associated genes  
948 (DAGs). The size and the color represent PCC for the gene entropy values versus  
949 degree (sum of in- and out-degree) and betweenness, respectively. **B.** The  
950 distribution and cumulative probability of the coefficient of variation for the  
951 DAGs among different disease contexts.

952

### 953 **Figure S5 Rank scores for the top 20% of high entropy genes for three diseases**

954 Bar plots show the rank scores of the top 20% of high entropy genes for systemic  
955 lupus (CD4 cells) (top), systemic lupus (B cells) (middle) and rheumatoid arthritis (B  
956 cells) (bottom). Red bars represent the rank scores of the core genes retrieved from  
957 NIA.

958

### 959 **Figure S6 The dose-response curve of the BxPC3 cell of 8 drugs**

960 **A-H.** The dose-response curve of BxPC3 cells for 8 drugs that have not been  
961 associated with pancreatic cancer. X-axis: the concentration of each drug; y-axis: the  
962 percent inhibition rate of the BxPC3 cells.

963

### 964 **Figure S7 The heat map of microarray experiment results**

965 A. Differentially expressed genes between the Saquinavir (saq1, saq2) treated BxPC3  
966 cell group and the control group (con1, con2). Color represent the relative  
967 expression of the differentially expressed genes. **B.** Differentially expressed genes  
968 between the Celecoxib (cel1, cel2) treated BxPC3 cell group and the control group  
969 (con1, con2).

970



971 **Figure S8 Estimation of the scale parameter  $\omega$**

972 Selected parameters ( $\omega=1.1$ ) that showed the highest mean AUROC and were thus  
973 used for further analysis.

974

975 **Figure S9 Characterization of gene entropy features with different scale**  
976 **parameters  $\omega$**

977 A. Normalized probability of different distances with scale parameter  $\omega$  ranging from  
978 0 to 4. **B.** Coronary disease gene entropy values with different scale parameters,  
979  $\omega=0$  (top) and  $\omega=10$  (bottom). **C.** Coronary disease gene entropy values with scale  
980 parameter  $\omega$  ranging from 0 to 10.

981

982 **Figure S10 Performance of the drug-disease relationship predictions using**  
983 **different scale parameters**

984 The box plot shows the AUROC for drug-disease predictions using different scale  
985 parameters. To account for the heterogeneous degree distribution of the directed  
986 interactome, we preserved the degree of randomized targets and disease genes.

987

988 **Table S1 Full list of disease-gene associations used in this study**

989 Entropy value: the entropy value calculated using NOGEA in a specific disease; rank  
990 score: the rank score for each gene entropy in a specific disease. This table also  
991 includes topology parameters of the DAGs in the directed global PPI network, i.e., the  
992 undirected degree, the in-degree, the out-degree and the betweenness centrality. In  
993 addition, this list includes the mean and standard deviations of the entropy among  
994 different diseases for a disease gene, the number of the gene-associated diseases and  
995 the coefficient of variation of the disease gene among different diseases. The evidence  
996 for the disease-gene associations was retrieved from CTD, TTD and PharmGKB.

997

998 **Table S2 List of the directed protein-protein interactions**

999 The list was obtained from a recent study as described in the paper, and each row  
1000 presents a directed edge.

1001

1002 **Table S3 Classification of the disease-associated genes**

1003 This list includes all the disease-gene relations used in this study. Genes of each  
1004 disease were assigned to master, interim and redundant groups according to their  
1005 entropy values.



1006

1007 **Table S4 Gene sets used for enrichment**

1008 This table lists all 8 different gene sets used for enrichment analysis, which contains  
1009 1707 OMIM genes, 2186 predicted cancer genes, 1750 essential genes, 1551  
1010 transcription factors, 366 kinases, 249 membrane receptors, 1336 druggable genes and  
1011 982 FDA targets, respectively. All gene sets were obtained from a recent study  
1012 (PMCID: PMC4983807).

1013

1014 **Table S5 Inferred comorbidity relationships of disease pairs from the shared**  
1015 **genes**

1016 This table lists all inferred comorbidity relationships involving master genes. As  
1017 described in the paper, if two diseases shared a master gene, they were considered to  
1018 be co-morbid diseases. Shared master genes are also listed.

1019

1020 **Table S6 Inferred comorbidity relationships of disease pairs from the interacting**  
1021 **gene pairs**

1022 This table lists all inferred comorbidity relationships involving master genes. As  
1023 described in the paper, if master genes of two diseases directly interact with each  
1024 other on the interactome, they were treated as co-morbid diseases. Interacting master  
1025 gene pairs are also listed.

1026

1027 **Table S7 Inferred causal or co-occurrence relationships between Parkinson's**  
1028 **and other diseases.**

1029 Results of the inferred relationships correspond with Figure 2C. This table lists all  
1030 inferred causal or cooccurrence relationships between Parkinson's disease and other  
1031 diseases. The validated relationships are marked as "YES". The "positive sim" is the  
1032 likelihood from "V1" to "V2" and the "negative sim" is the likelihood from "V2" to  
1033 "V1".

1034

1035 **Table S8 Information for all FDA-approved drugs that were used in the present**  
1036 **study**

1037 This table lists all FDA approved drugs that were used in the present work and their  
1038 corresponding IDs in other databases.

1039

1040 **Table S9 List of drug-target relationships used in the present study**

1041 This table lists all FDA drug-target relationships used in this study.

1042

1043 **Table S10 Drug-disease information**

1044 This table includes FDA drug indications, drug names and corresponding MESH IDs  
1045 inferred from the indication information.

1046

1047 **Table S11 Gene rank list for three diseases**

1048 This table lists the gene rank scores and core genes for systemic lupus (CD4 cells),  
1049 systemic lupus (B cells) and rheumatoid arthritis (B cells).

1050

1051 **Table S12 Drug disturbance entropy (DDE) for each FDA-approved drug**  
1052 **associated with three diseases**

1053 This table lists the value of DDE calculated using NOGEA for each FDA-approved  
1054 drug associated with the systemic lupus (CD4 cells), systemic lupus (B cells) and  
1055 rheumatoid arthritis (B cells).

1056

1057 **Table S13 FDA-approved drugs and their categories**

1058 This table lists all present FDA approved drugs and their corresponding categories  
1059 retrieved from the DrugBank database.

1060

1061 **Table S14 The DDE for each FDA-approved drug associated with pancreatic**  
1062 **cancer and the literature mining results**

1063 This table lists all the DDE scores calculated using NOGEA. The result of literature  
1064 mining contains the number of articles derived by searching each drug name,  
1065 “pancreatic cancer” as well as both search terms, respectively. The P-values were  
1066 assessed using the hypergeometric test.

1067

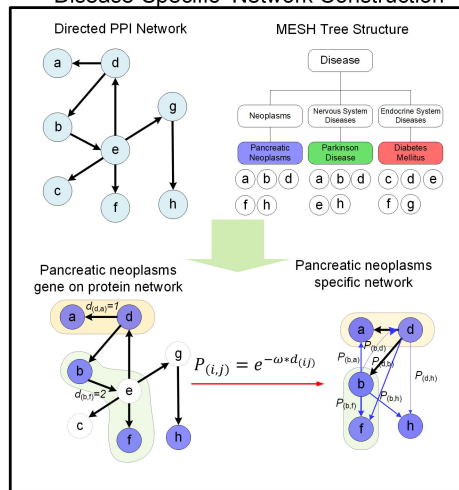
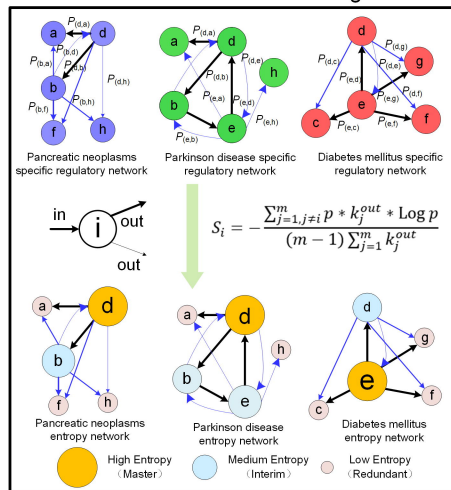
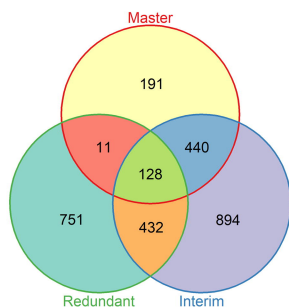
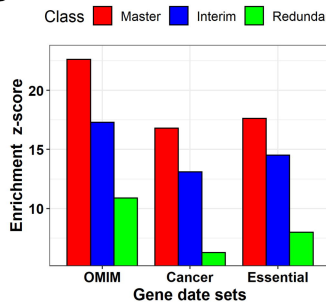
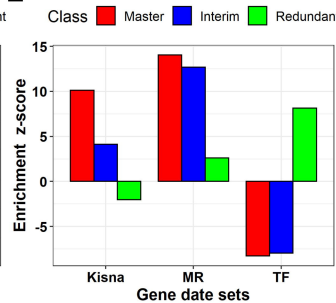
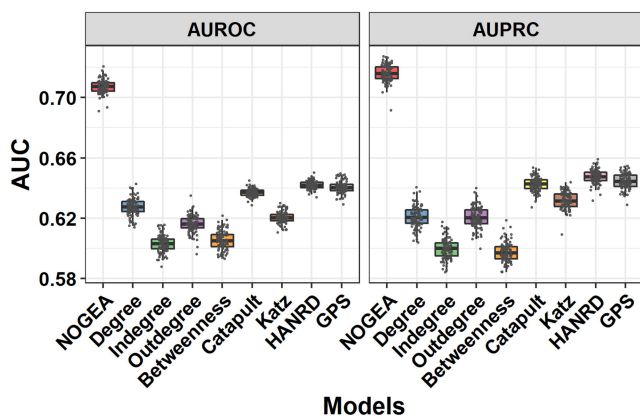
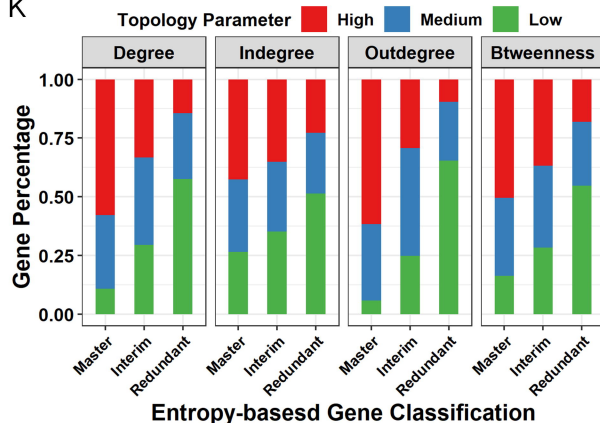
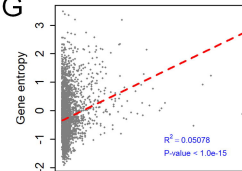
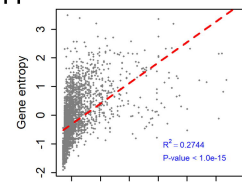
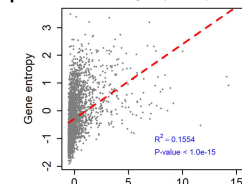
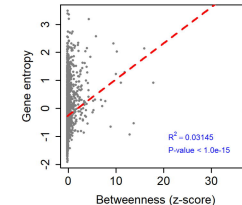
1068 **Table S15 Differentially expressed genes and the predicted effected genes after**  
1069 **treatment with Saquinavir and Celecoxib.**

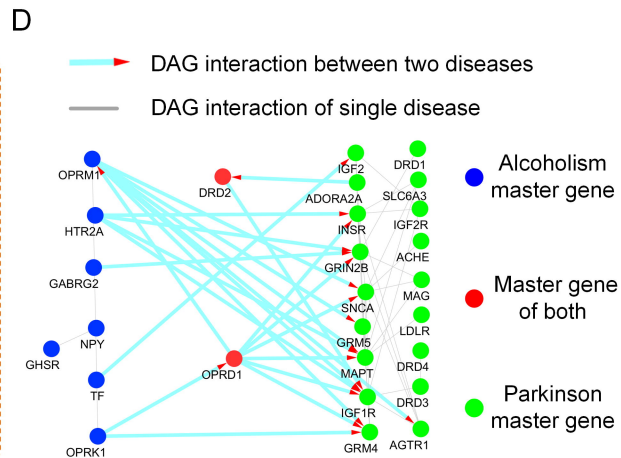
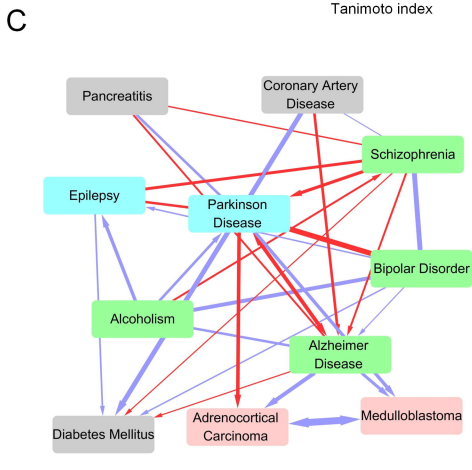
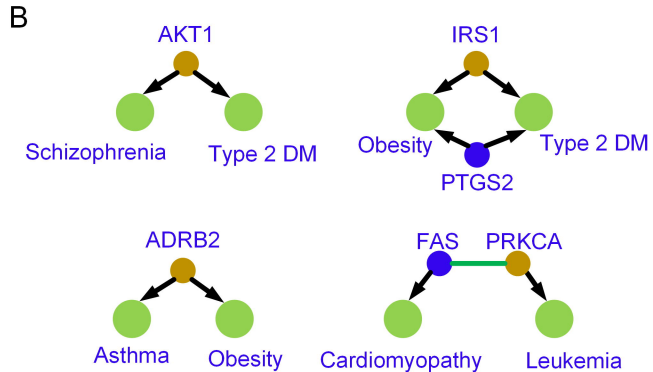
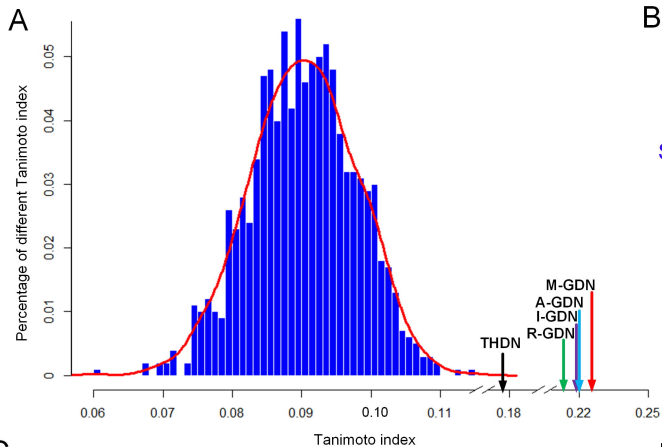
1070 CELDEG: the differentially expressed gene after treatment with Celecoxib. CELPEG:  
1071 the predicted effected gene after treatment with Celecoxib. SAQDEG: the  
1072 differentially expressed gene after treat with Saquinavir. SAQPEG: the predicted  
1073 effected gene after treatment with Saquinavir.

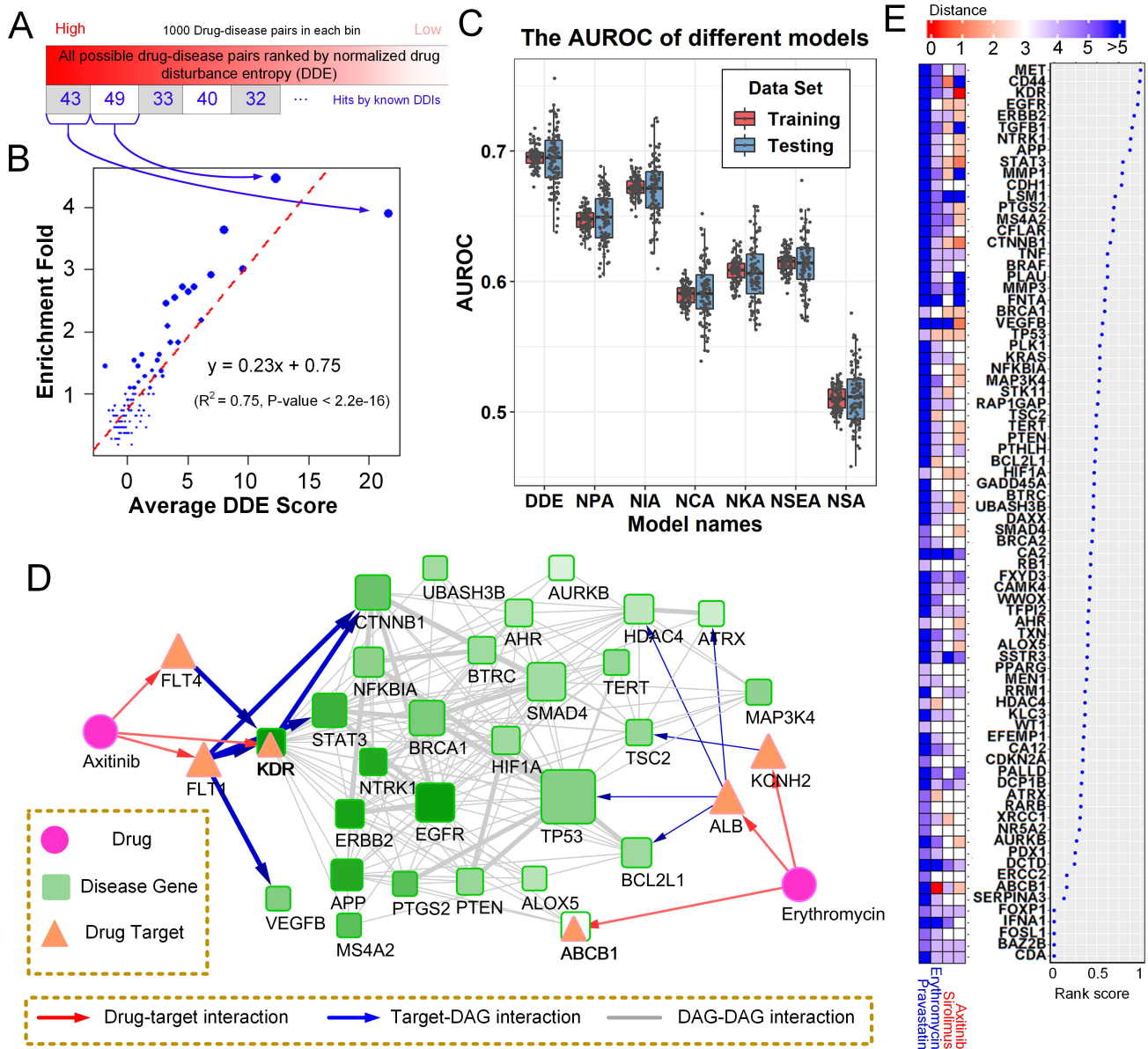
1074

1075 **Table S16 Release versions of the database used in this study.**

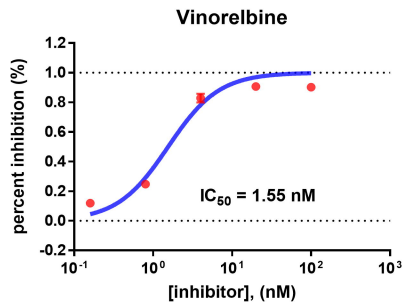
1076 This table lists all the databases and corresponding versions that were used in this  
1077 study.

**A** Disease-Specific Network Construction**B** Disease Gene Prioritizing**C****D****E****F****The AUROC and AUPRC of Different Models****K****G****H****I****J**

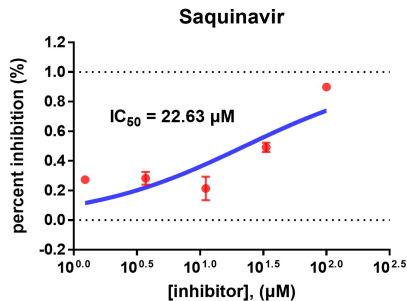




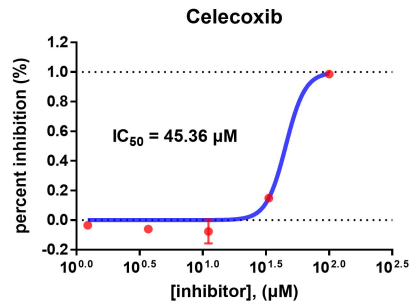
A



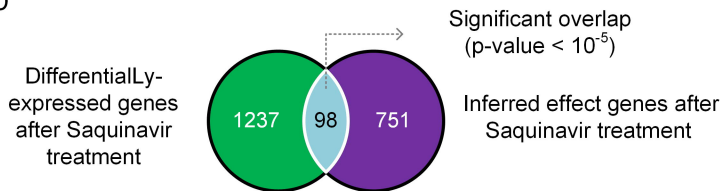
B



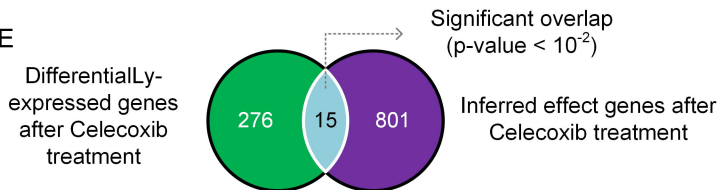
C



D



E



F

