

Supplementary Materials:  
**A novel network regularized matrix  
decomposition method to detect  
mutated cancer genes in tumour  
samples with inter-patient heterogeneity**

Jianing Xi<sup>1,+</sup>, Ao Li<sup>1,2,+,\*</sup> and Minghui Wang<sup>1,2</sup>

<sup>1</sup>School of Information Science and Technology, University of Science and Technology of China, Hefei AH230027, China and

<sup>2</sup>Centers for Biomedical Engineering, University of Science and Technology of China, Hefei AH230027, China.

\*E-mail: aoli@ustc.edu.cn

<sup>+</sup>These authors contributed equally to this work

## Contents

<b>1</b>	<b>Supplementary Material</b>	<b>2</b>
1.1	Configuration settings of mCGfinder, HotNet2 and ReMIC . . .	2
1.2	Ranking of the genes detected by mCGfinder, HotNet2 and ReMIC	2
1.3	Normalized Adjacency and Laplacian matrix normalization . . .	2
1.4	Proof: $(\ \mathbf{s}_r\ _2^2 \mathbf{I}_p + \lambda_L \mathbf{L})$ is an invertible matrix . . . . .	3
1.5	Significance test through a semiexact estimation . . . . .	3
1.6	Input data: TCGA somatic mutation data . . . . .	4
<b>2</b>	<b>Supplementary Figures and Captions</b>	<b>6</b>
<b>3</b>	<b>Supplementary Tables and Captions</b>	<b>13</b>

# 1 Supplementary Material

## 1.1 Configuration settings of mCGfinder, HotNet2 and ReMIC

In comparison analysis, the detection results of mCGfinder are selected by default threshold 0.05. The tuning parameter used to balance the fitness of the model and the smoothness of the scores of connected genes is set to 0.1.

In HotNet2, the parameter of insulated heat-diffusion  $\beta$  is set to 0.45 for the iRefIndex network [1] as suggested [2]. The permuted networks of the gene interaction network are generated by HotNet2. The scores of the genes are set to their mutation frequencies, and the numbers of delta permutations and significance permutations are set to 100 as default settings [2].

In ReMIC, the number of iterations to generate the bag of random mutation score is set to 20000, which satisfy the condition that it is larger than number of genes (totally 12129 genes in iRefIndex network). The number of permutations is set to 10000. The pseudocount is set to true. The diffusion strength  $\beta$  is set to 0.01, 0.02 and 0.03 respectively as suggested in ReMIC [3] (scale parameter of 0 is excluded since its result is equivalent to gene mutation frequencies). Since the performance of ReMIC with  $\beta = 0.03$  is relative better than the performance with  $\beta = 0.01$  or 0.02, we use the results of ReMIC with  $\beta = 0.03$  in the comparison study.

## 1.2 Ranking of the genes detected by mCGfinder, HotNet2 and ReMIC

The genes detected by mCGfinder are ranked by negative logarithm of q-values. The q-values of genes are calculated by the significance test and false discovery rate control [4] in mCGfinder.

The genes detected by HotNet2 are ranked by the values of the minimum edge weight parameter  $\delta$ , which have been used to calculate the true positive rates and false positive rates in previous study [2].

The genes detected by ReMIC are sorted by the negative logarithm of p-values, which are calculated by permutation test in ReMIC [3].

## 1.3 Normalized Adjacency and Laplacian matrix normalization

The gene interaction network in iRefindex [1] is an undirected, unweighted graph  $G = (\mathbf{V}, \mathbf{E})$  without graph loops  $(i, i)$  or multiple edges from one node to another, where  $\mathbf{V}$  is the vertex set,  $p = |\mathbf{V}|$ , and  $\mathbf{E}$  is the edge set. Then the symmetric normalized adjacency matrix of the graph  $G$  is an  $p \times p$  symmetric matrix with one row and column for each node defined by [5]:

$$\mathbf{A} = \mathbf{D}^{-1/2} \mathbf{A}^{(adj)} \mathbf{D}^{-1/2},$$

where  $\mathbf{D} = \text{Diag}(d_1, \dots, d_i, \dots, d_p)$  for  $d_i$  the degree of node  $i$  in the graph  $G$  and  $\mathbf{A}^{(adj)}$  is the original adjacency matrix of graph  $G$ . Therefore, the diagonal

elements  $A_{ij}$  of  $\mathbf{A}$  are equal the elements  $A_{ij}^{(adj)}$  of  $\mathbf{A}^{adj}$  divided by the square root of the product of  $d_i$  and  $d_j$ , i.e.

$$A_{ij} = A_{ij}^{adj} / \sqrt{d_i d_j}. \quad (1)$$

The normalized Laplacian Matrix is

$$\mathbf{L} = \mathbf{I} - \mathbf{A} = \mathbf{D}^{-1/2}(\mathbf{D} - \mathbf{A}_{adj})\mathbf{D}^{-1/2} = \mathbf{D}^{-1/2}\mathbf{L}_G\mathbf{D}^{-1/2},$$

for  $\mathbf{L}_G$  the un-normalized Laplacian. Therefore, the diagonal elements  $L_{ij}$  of  $\mathbf{L}$  are equal the degree of vertex  $i$  and off-diagonal elements  $L_{ij}$  are -1 if vertex  $i$  is adjacent to  $j$  and 0 otherwise [5], i.e.

$$L_{ij} = \begin{cases} 1 & \text{if } i = j \text{ and } d_j \neq 0, \\ -\frac{1}{\sqrt{d_i d_j}} & \text{if } i \text{ and } j \text{ are adjacent,} \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

#### 1.4 Proof: $(\|\mathbf{s}_r\|_2^2 \mathbf{I}_p + \lambda_L \mathbf{L})$ is an invertible matrix

*Proof.* Note that graph Laplacian matrix  $\mathbf{L}$  ( $p \times p$ ) is positive semidefinite. Thus, through eigendecomposition, it can be factorized as  $\mathbf{L} = \mathbf{P}^T \mathbf{\Lambda} \mathbf{P}$ , where  $\mathbf{P}$  is an orthogonal matrix and  $\mathbf{\Lambda}$  is a diagonal matrix whose diagonal entries are the eigenvalues of  $\mathbf{L}$ . Due to positive semidefinite, all diagonal entries of diagonal matrix  $\mathbf{\Lambda}$  are nonnegative.

Because of the matrix orthogonality  $\mathbf{P}^T \mathbf{P} = \mathbf{I}_p$ , the matrix  $(\|\mathbf{s}_r\|_2^2 \mathbf{I}_p + \lambda_L \mathbf{L})$  can be factorized as

$$(\|\mathbf{s}_r\|_2^2 \mathbf{I}_p + \lambda_L \mathbf{L}) = \mathbf{P}^T (\|\mathbf{s}_r\|_2^2 \mathbf{I}_p + \lambda_L \mathbf{\Lambda}) \mathbf{P},$$

where  $(\|\mathbf{s}_r\|_2^2 \mathbf{I}_p + \lambda_L \mathbf{\Lambda})$  is also a diagonal matrix. Note that  $\|\mathbf{s}_r\|_2^2$  is always positive,  $\lambda_L$  is a nonnegative tuning parameter and all diagonal entries of matrix  $\mathbf{\Lambda}$  are nonnegative. Consequently, all diagonal entries of  $(\|\mathbf{s}_r\|_2^2 \mathbf{I}_p + \lambda_L \mathbf{\Lambda})$  are positive, suggesting that it is a positive definite matrix. Therefore, the investigated matrix  $(\|\mathbf{s}_r\|_2^2 \mathbf{I}_p + \lambda_L \mathbf{L})$  is invertible.

#### 1.5 Significance test through a semiexact estimation

In brief, we define  $\mathbf{X}_{\text{net}} := \left[ (\|\mathbf{s}_r\|_2^2 \mathbf{I}_p + \lambda_L \mathbf{L})^{-1} \mathbf{X}^T \right]^T = \mathbf{X} (\|\mathbf{s}_r\|_2^2 \mathbf{I}_p + \lambda_L \mathbf{L})^{-1}$  as the network influenced matrix. For the  $r$ -th component, the coefficients of gene score vector  $\mathbf{g}_r$  can be calculated by the summation of the entries of a subset of rows of the network influenced matrix  $\mathbf{X}_{\text{net}}$ , where the rows are indicated by the sample indicator vector  $\mathbf{s}_r$  of the investigated component. To assess which of these mutated genes are statistically significant in a subset of samples, we follow the procedure of previous studies [6, 7] and identify the genes of which the scores can disprove the null hypothesis that their values of the gene score vector coefficients are only contributed by background mutations alone. Since the random background mutations could occur anywhere in the genome,

we model the null distribution by recalculating the gene score vectors across all combinations of permutations of the network influenced matrix  $X_{\text{net}}$  within rows (samples) indicated by  $\mathbf{s}_r$  of the  $r$ -th component [6, 7]. Under the null hypothesis, the arrangement of entries in  $X_{\text{net}}$  is independent between the indicated samples (rows). Accordingly, by permuting the entries in the rows indicated by  $\mathbf{s}_r$  of the matrix  $X_{\text{net}}$ , we can generate a conservative, high estimate of the null distribution which contains information from both the somatic mutations and the network context.

Since large numbers of permutations is usually time consuming, we follow the procedure proposed in previous approaches [6, 7] by using a semi-exact estimate of this null distribution instead of simulating the null distribution by performing each of these permutations in turn. The distribution of the sum of across the indicated rows equals the convolution of the distributions of entries in all the indicated rows. For the investigated component, we approximate these distributions by generating histograms  $\mathbf{h}_i^r(x_{\text{net}})$  for the  $i$ -th indicated row of the network influenced matrix  $X_{\text{net}}$ , where the number of bins is set to  $10^5$ . The final distribution for coefficient values of the gene score vector is calculated by  $\mathbf{H}^r = \mathbf{h}_1^r \otimes \mathbf{h}_2^r \dots \otimes \mathbf{h}_{l_r}^r$ , where  $1, 2, \dots, l_r$  is the indices of rows indicated by the  $r$ -th component (totally  $l_r$  samples included in the investigated component). By comparing the coefficients of the estimated vector  $\mathbf{g}_r$  to the distribution above, we can assign the p-value for each investigated gene by the sum of the tail of the null distribution estimated above.

## 1.6 Input data: TCGA somatic mutation data

In this study, we use TCGA somatic mutation data to evaluate the performances of mCGfinder. To prevent mutagens or carcinogens involved in cancer treatment which could cloud the origin of the cancer, TCGA have strict sample criteria in acquiring tissue samples, such as “sample from primary tumor was necessary” and “neoadjuvant treatment was not allowable”:

<https://cancergenome.nih.gov/cancersselected/biospeccriteria>

Therefore, to the best of our knowledge, the underlying genomics of primary, untreated tumor samples in TCGA is not affected by chemotherapy.

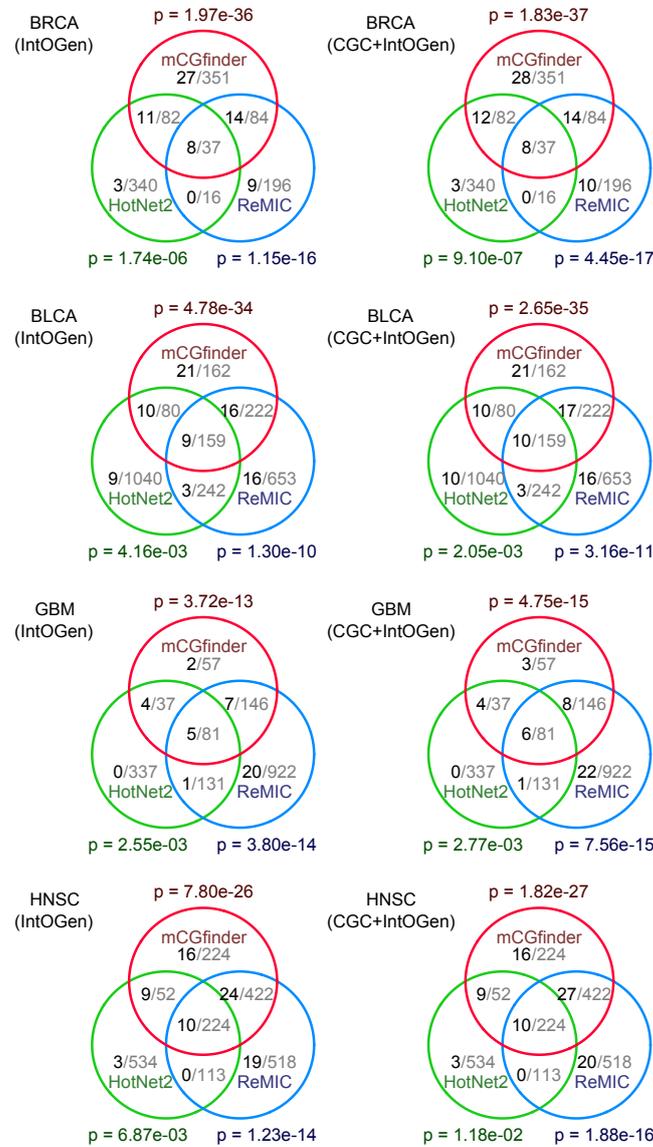
In consistent with HotNet2 and ReMIC, somatic mutation data are required as the input of mCGfinder. Therefore, somatic mutations from raw data files should be filtered to remove polymorphisms as described in previous study [8]. More detailed information of the datasets of the investigated cancers are provided in Supplementary Table S5.

## References

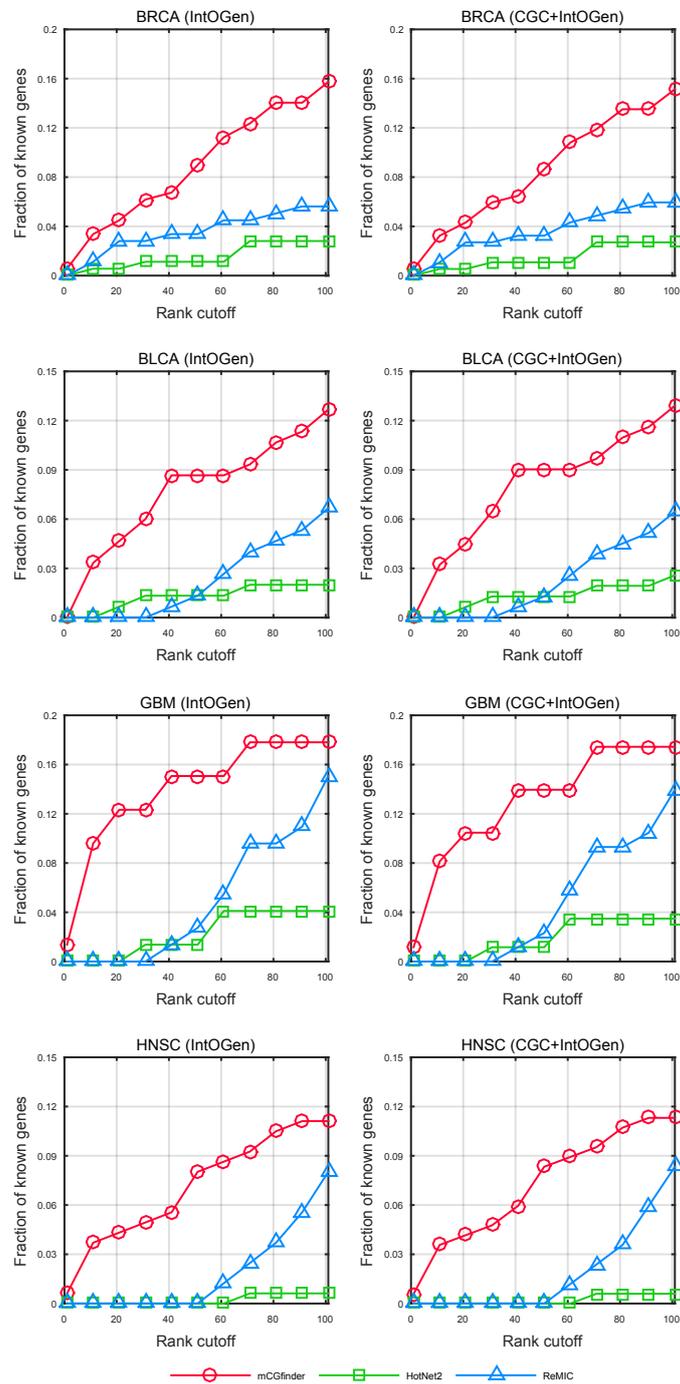
- [1] Sabry Razick, George Magklaras, and Ian M Donaldson. irefindex: a consolidated protein interaction database with provenance. *BMC bioinformatics*, 9(1):1, 2008.
- [2] Mark D Leiserson, Fabio Vandin, Hsin-Ta Wu, Jason R Dobson, and Benjamin R Raphael. Pan-cancer identification of mutated pathways and protein complexes. *Cancer Research*, 74(19 Supplement):5324–5324, 2014.

- [3] Sepideh Babaei, Marc Hulsman, Marcel Reinders, and Jeroen de Ridder. Detecting recurrent gene mutation in interaction network context using multi-scale graph diffusion. *BMC bioinformatics*, 14(1):1, 2013.
- [4] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the royal statistical society. Series B (Methodological)*, pages 289–300, 1995.
- [5] Fan RK Chung. *Spectral graph theory*, volume 92. American Mathematical Soc., 1997.
- [6] Jianing Xi and Ao Li. Discovering recurrent copy number aberrations in complex patterns via non-negative sparse singular value decomposition. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 13(4):656–668, 2016.
- [7] Rameen Beroukhi, Gad Getz, Leia Nghiemphu, Jordi Barretina, Teli Hsueh, David Linhart, Igor Vivanco, Jeffrey C Lee, Julie H Huang, Sethu Alexander, et al. Assessing the significance of chromosomal aberrations in cancer: methodology and application to glioma. *Proceedings of the National Academy of Sciences*, 104(50):20007–20012, 2007.
- [8] Cancer Genome Atlas Network et al. Comprehensive molecular portraits of human breast tumours. *Nature*, 490(7418):61–70, 2012.
- [9] P Andrew Futreal, Lachlan Coin, Mhairi Marshall, Thomas Down, Timothy Hubbard, Richard Wooster, Nazneen Rahman, and Michael R Stratton. A census of human cancer genes. *Nature Reviews Cancer*, 4(3):177–183, 2004.
- [10] Abel Gonzalez-Perez, Christian Perez-Llamas, Jordi Deu-Pons, David Tamborero, Michael P Schroeder, Alba Jene-Sanz, Alberto Santos, and Nuria Lopez-Bigas. Intogen-mutations identifies cancer drivers across tumor types. *Nature methods*, 10(11):1081–1082, 2013.
- [11] Jingchun Zhu, J Zachary Sanborn, Stephen Benz, Christopher Szeto, Fan Hsu, Robert M Kuhn, Donna Karolchik, John Archie, Marc E Lenburg, Laura J Esserman, et al. The ucsc cancer genomics browser. *Nature methods*, 6(4):239–240, 2009.

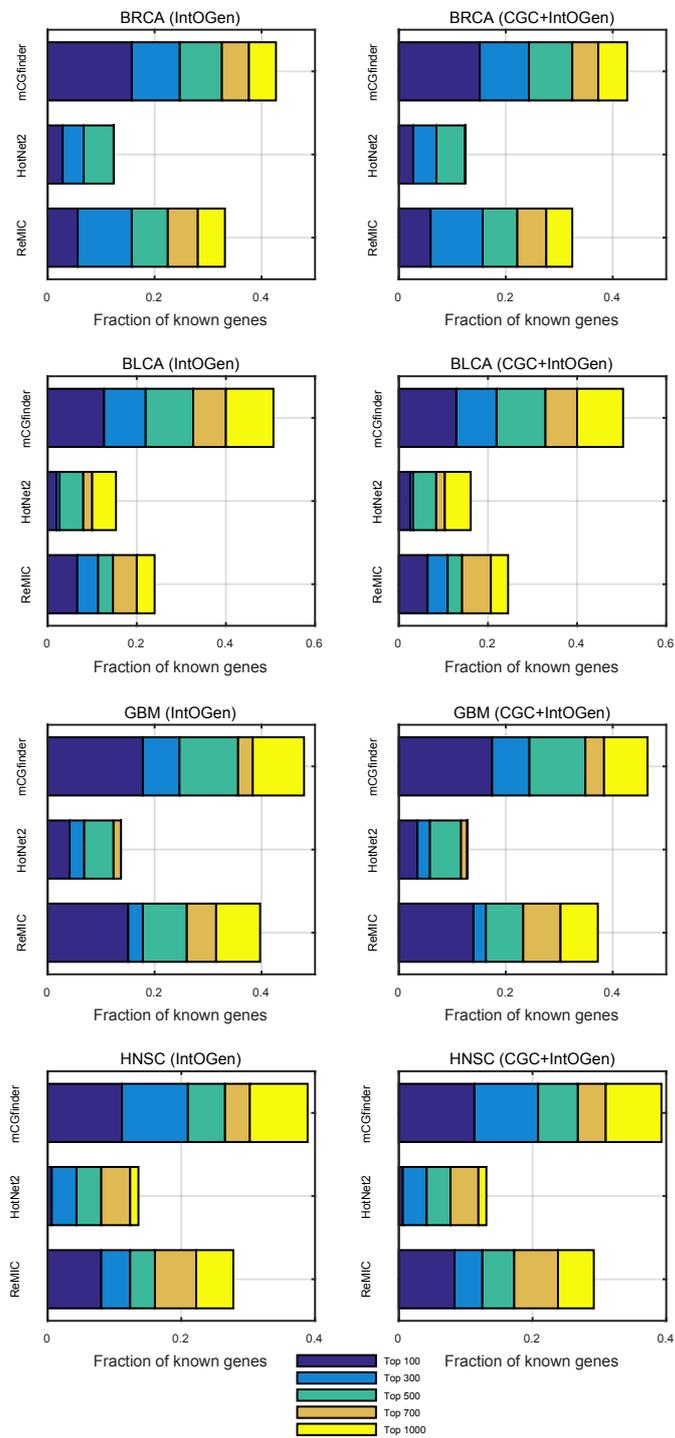
## 2 Supplementary Figures and Captions



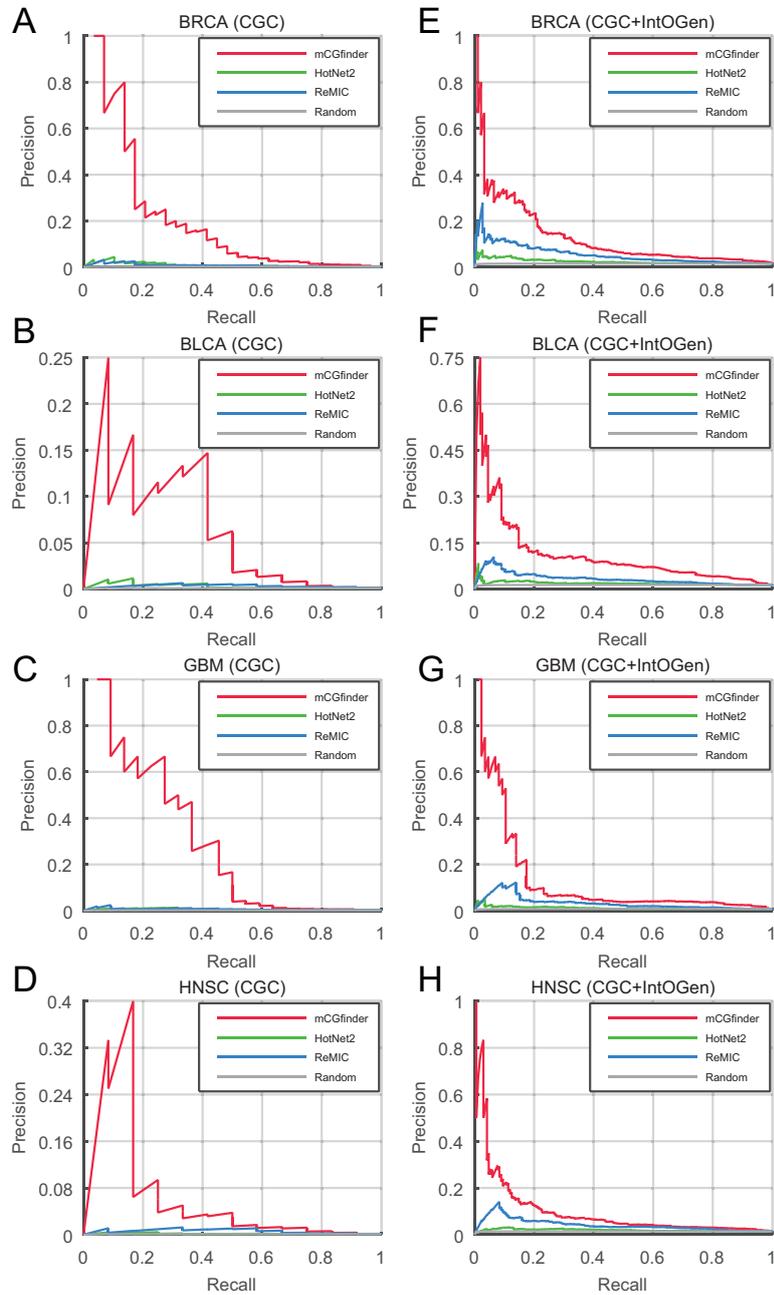
**Supplementary Figure S1.** Venn diagrams of intersections between the genes detected by mCGfinder (red circle), HotNet2 (green circle) and ReMIC (blue circle) on TCGA somatic mutation datasets of BRCA (first row), BLCA (second row), GBM (third row) and HNSC (fourth row). In each region, The gray and black numbers in each region of the Venn diagrams indicate the number of detected genes and the number of genes also reported in IntOGen gene lists (first column) and the combined lists of the two databases (second column) respectively. The p-values next to the circles of methods are the enrichment significance of the detection results for the validation known cancer gene lists.



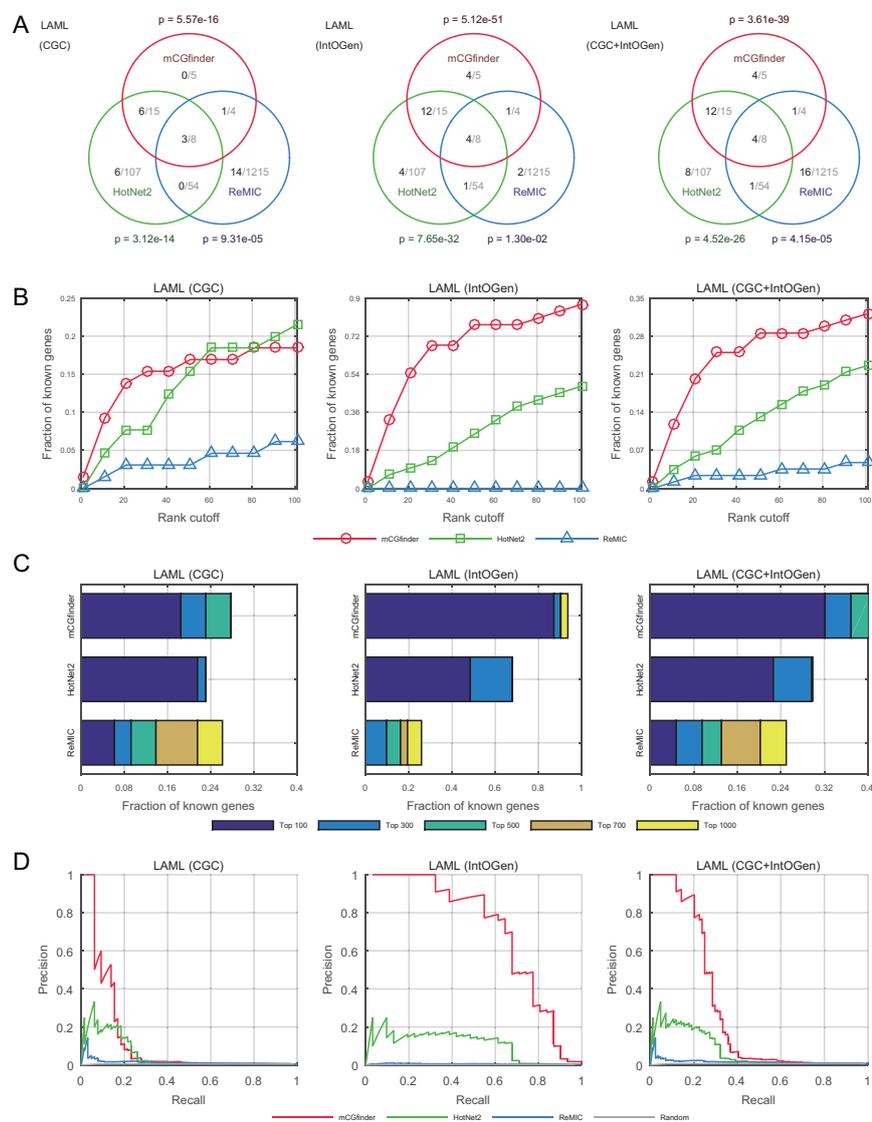
**Supplementary Figure S2.** Rank cutoff curves of top 100 candidates in mCGfinder (red line with circle markers), HotNet2 (green line with square markers) and ReMIC (blue line with triangle markers) results, describing the relation between various cutoffs and the fraction of known IntOGen cancer genes (first column) and the combined genes from the two databases (second column) that are ranked above this cutoff in BRCA (first row), BLCA (second row), GBM (third row) and HNSC (fourth row).



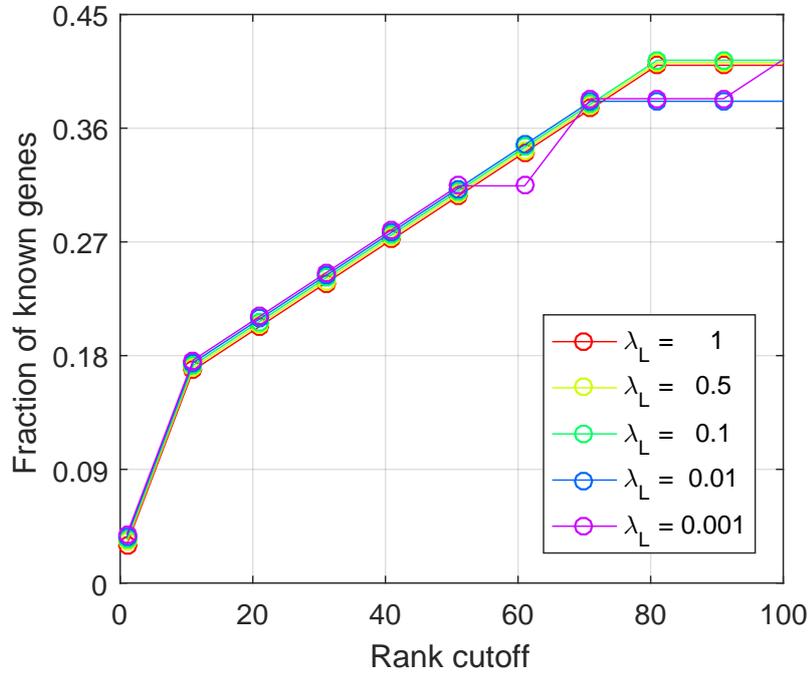
**Supplementary Figure S3.** Cumulative fractions of known cancer genes reported in IntOGen (first column) and the combined genes lists from the two databases (second column) within the top 100, 300, 500, 700 and 1000 genes on BRCA (first row), BLCA (second row), GBM (third row) and HNSC (fourth row).



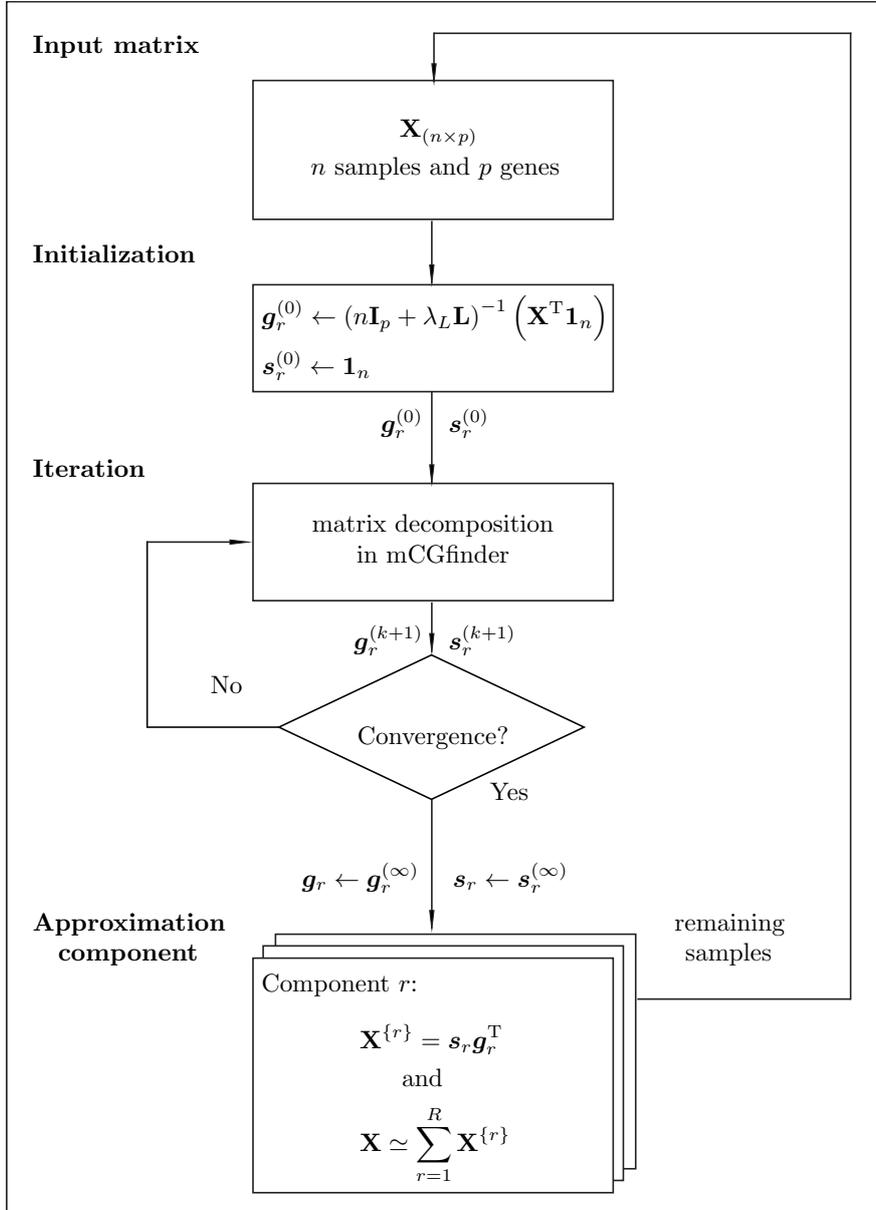
**Supplementary Figure S4.** Precision-recall curves for the three methods on BRCA (A and E), BLCA (B and F), GBM (C and G) and HNSC (D and H) data, validated by CGC (A-D) and the combined genes lists from both CGC and IntOGen (E-H) respectively. where each point indicates the precision and recall at a different rank in the prediction. Red, green and blue lines represent the curves of mCGfinder, HotNet2 and ReMIC results.



**Supplementary Figure S5.** The detection results on TCGA LAML data given by mCGfinder (red), HotNet2 (green) and ReMIC (blue), validated by CGC (left column), IntOGen (middle column) and the combined genes lists from the two databases (right column). (A) Venn diagrams of intersections between the genes detected by the three methods. (B) Rank cutoff curves of top 100 candidates detected by the three methods. (C) Cumulative fractions of known cancer genes within the top 100, 300, 500, 700 and 1000 genes. (D) Precision-recall curves for the three methods.



**Supplementary Figure S6.** The detection results of mCGfinder with different tuning parameter values on BRCA data. The results show that the performance of mCGfinder is not sensitive to the selection of the tuning parameter when the value varies from 1 to 0.001.



**Supplementary Figure S7.** Flowchart of the component-by-component decomposition strategy in mCGfinder model. After initialization, mCGfinder estimates the sample indicator vector  $\mathbf{s}$  and the gene score vector  $\mathbf{g}$  through the iterative estimation procedure until convergence. Then the component can be obtained through the outer product of the two vectors  $\mathbf{s}\mathbf{g}^T$ . We repeat this procedures aforementioned on the remaining samples, until all samples are assigned.

### 3 Supplementary Tables and Captions

**Supplementary Table S1.** The full lists of Cancer Gene Census (CGC) annotated cancer genes [9] detected by mCGfinder on BRCA (A), BLCA (B), GBM (C), HNSC (D) and LAML (E) data, sorted by their rank in mCGfinder result.

(A) CGC genes in BRCA results

Gene Symbol	Rank	q-value	also detected by
PIK3CA	1	$\leq 1.00e-159$	HotNet2/ReMIC
TP53	2	$\leq 1.00e-159$	HotNet2/ReMIC
GATA3	4	1.96e-99	HotNet2
MAP3K1	5	2.19e-44	
CDH1	9	1.18e-28	HotNet2
NCOR1	21	1.03e-11	
MAP2K4	29	2.87e-10	HotNet2
CTCF	32	2.04e-08	
AKT1	45	8.21e-07	
RB1	53	1.64e-06	ReMIC
ARID1A	69	1.44e-05	HotNet2
FOXA1	73	4.48e-05	HotNet2
TBX3	104	2.57e-04	
ARID1B	154	2.67e-03	HotNet2
BRCA2	240	1.10e-02	
ERBB2	340	1.82e-02	ReMIC
CASP8	392	4.03e-02	
MAP3K13	456	4.03e-02	

(B) CGC genes in BLCA results

Gene Symbol	Rank	q-value	also detected by
KDM6A	4	6.80e-41	HotNet2
STAG2	12	3.94e-17	HotNet2/ReMIC
LRP1B	26	1.76e-11	ReMIC
FGFR3	30	2.01e-10	HotNet2
ERBB3	34	2.07e-09	ReMIC
TSC1	96	9.24e-06	mCGfinder
NOTCH2	340	6.68e-03	HotNet2/ReMIC
NOTCH1	537	2.40e-02	ReMIC

(C) CGC genes in GBM results

Gene Symbol	Rank	q-value	also detected by
PTEN	1	1.99e-159	HotNet2
TP53	2	7.24e-141	HotNet2
EGFR	4	1.42e-109	HotNet2/ReMIC
PIK3CA	6	3.70e-38	HotNet2/ReMIC

*continued on next page*

continued from previous page

PIK3R1	8	4.14e-33	HotNet2/ReMIC
NF1	9	1.80e-31	ReMIC
ATRX	14	2.12e-17	ReMIC
IDH1	17	3.59e-13	HotNet2
PDGFRA	32	4.51e-08	ReMIC
STAG2	33	4.51e-08	ReMIC
LZTR1	66	7.43e-05	
ROS1	283	2.36e-02	HotNet2/ReMIC

(D) CGC genes in HNSC results

Gene Symbol	Rank	q-value	also detected by
FAT1	3	6.39e-114	ReMIC
NOTCH1	5	1.54e-78	HotNet2/ReMIC
FAT4	32	3.46e-16	ReMIC
NFE2L2	80	2.45e-09	HotNet2/ReMIC
TGFBR2	144	8.53e-07	
CTCF	160	4.80e-06	HotNet2
ERBB3	421	2.33e-03	ReMIC
BCORL1	615	1.07e-02	ReMIC
MTOR	753	2.95e-02	ReMIC

(E) CGC genes in LAML results

GeneSymbol	Rank	q-value	also detected by
NPM1	1	5.18e-165	mCGfinder/HotNet2
DNMT3A	2	1.03e-93	mCGfinder/HotNet2/ReMIC
FLT3	3	4.11e-87	mCGfinder/ReMIC
RUNX1	4	2.22e-45	mCGfinder/HotNet2
PTPN11	9	5.07e-10	mCGfinder/HotNet2/ReMIC
NRAS	10	1.13e-09	mCGfinder/HotNet2
CEBPA	15	3.79e-06	mCGfinder/HotNet2
KIT	16	4.17e-06	mCGfinder/HotNet2
RAD21	17	7.77e-05	mCGfinder/HotNet2
KRAS	23	1.45e-03	mCGfinder/HotNet2/ReMIC

**Supplementary Table S2.** The full lists of Integrative Onco Genomics (IntOGen) annotated cancer genes [10] detected by mCGfinder on BRCA (A), BLCA (B), GBM (C), HNSC (D) and LAML (E) data, sorted by their rank in mCGfinder result.

(A) IntOGen genes in BRCA results			
Gene Symbol	Rank	q-value	also detected by
PIK3CA	1	<i>leq</i> 1.00e-159	HotNet2/ReMIC
TP53	2	<i>leq</i> 1.00e-159	HotNet2/ReMIC
GATA3	4	1.96e-99	HotNet2
MAP3K1	5	2.19e-44	
MLL3	8	2.45e-34	HotNet2/ReMIC
CDH1	9	1.18e-28	HotNet2
MACF1	20	1.03e-11	HotNet2/ReMIC
NCOR1	21	1.03e-11	
PTEN	27	1.34e-10	HotNet2/ReMIC
CBFB	28	2.87e-10	HotNet2
MAP2K4	29	2.87e-10	HotNet2
CTCF	32	2.04e-08	
RUNX1	44	2.82e-07	
AKT1	45	8.21e-07	
ASPM	46	1.64e-06	HotNet2
NF1	51	1.64e-06	ReMIC
PIK3R1	52	1.64e-06	ReMIC
RB1	53	1.64e-06	ReMIC
ATM	58	4.76e-06	
AKAP9	60	1.37e-05	
ARID1A	69	1.44e-05	HotNet2
TBL1XR1	71	2.58e-05	
FOXA1	73	4.48e-05	HotNet2
ANK3	76	8.87e-05	HotNet2
ASH1L	77	8.87e-05	
MYH14	92	8.87e-05	HotNet2
SETDB1	95	8.87e-05	
SVEP1	96	8.87e-05	ReMIC
SF3B1	102	1.07e-04	
TBX3	104	2.57e-04	
CCAR1	113	5.16e-04	HotNet2
MLL2	128	5.16e-04	HotNet2/ReMIC
RBM5	137	5.16e-04	
RPGR	139	5.16e-04	HotNet2
BRCA1	157	2.67e-03	ReMIC
CAD	159	2.67e-03	
CHD4	162	2.67e-03	
MGA	186	2.67e-03	
AHNAK	224	1.10e-02	ReMIC
BRCA2	240	1.10e-02	
EGFR	262	1.10e-02	ReMIC
MTOR	285	1.10e-02	

continued on next page

continued from previous page

MYH11	288	1.10e-02	ReMIC
NOTCH2	298	1.10e-02	ReMIC
SETD2	318	1.10e-02	HotNet2/ReMIC
STAG2	320	1.10e-02	HotNet2/ReMIC
TAF1	322	1.10e-02	
HCFC1	337	1.16e-02	
ERBB2	340	1.82e-02	ReMIC
PIK3CB	346	2.45e-02	
KRAS	354	3.19e-02	
MYB	355	3.19e-02	
ZFP36L1	370	3.63e-02	
ATR	384	4.03e-02	ReMIC
CASP8	392	4.03e-02	
MLL	462	4.03e-02	
MLLT4	464	4.03e-02	ReMIC
MYH9	466	4.03e-02	
SMARCA4	513	4.03e-02	ReMIC
FN1	553	4.83e-02	ReMIC

(B) IntOGen genes in BLCA results

Gene Symbol	Rank	q-value	also detected by
TP53	2	5.81e-116	
ARID1A	3	7.86e-44	HotNet2/ReMIC
KDM6A	4	6.80e-41	HotNet2
RB1	7	3.34e-26	
ELF3	11	3.94e-17	HotNet2
STAG2	12	3.94e-17	HotNet2/ReMIC
EP300	15	5.55e-16	HotNet2
CDKN1A	29	2.01e-10	
FGFR3	30	2.01e-10	HotNet2
ERBB3	34	2.07e-09	ReMIC
ERCC2	35	2.07e-09	
FAT1	36	2.07e-09	ReMIC
AHNAK	41	2.10e-08	ReMIC
NCOR2	64	1.41e-06	ReMIC
ARHGAP35	73	9.24e-06	HotNet2/ReMIC
FBXW7	78	9.24e-06	HotNet2/ReMIC
HSP90AA1	83	9.24e-06	
TRIO	94	9.24e-06	ReMIC
TSC1	96	9.24e-06	
ANK3	102	5.70e-05	HotNet2/ReMIC
CHEK2	105	5.70e-05	
NFE2L2	115	5.70e-05	HotNet2/ReMIC
TP53BP1	176	3.13e-04	ReMIC
TXNIP	178	3.13e-04	HotNet2
APC	184	5.49e-04	ReMIC
ATR	186	5.49e-04	ReMIC

continued on next page

continued from previous page

CDKN2A	204	1.47e-03	ReMIC
MGA	236	1.47e-03	ReMIC
MYH10	242	1.47e-03	
RHOA	257	1.47e-03	ReMIC
SETD2	261	1.47e-03	HotNet2/ReMIC
SF3B1	262	1.47e-03	
CAD	300	6.68e-03	ReMIC
CNOT1	304	6.68e-03	HotNet2
HLA-A	326	6.68e-03	HotNet2
NUP98	341	6.68e-03	
SMC1A	357	6.68e-03	
CLTC	384	1.64e-02	
FN1	391	1.64e-02	ReMIC
CLSPN	423	1.77e-02	
HSP90AB1	429	1.87e-02	
ARID1B	430	1.91e-02	HotNet2/ReMIC
ACTB	438	2.40e-02	
AFF4	441	2.40e-02	HotNet2
CDK12	461	2.40e-02	HotNet2/ReMIC
CHD3	464	2.40e-02	
CHD9	465	2.40e-02	ReMIC
CLASP2	466	2.40e-02	HotNet2
EIF2AK3	481	2.40e-02	HotNet2
MAP3K1	520	2.40e-02	
MAP3K4	521	2.40e-02	ReMIC
MECOM	524	2.40e-02	
MLH1	527	2.40e-02	
NAP1L1	532	2.40e-02	
NOTCH1	537	2.40e-02	ReMIC
PTEN	556	2.40e-02	

(C) IntOGen genes in GBM results

Gene Symbol	Rank	q-value	also detected by
PTEN	1	1.99e-159	HotNet2
TP53	2	7.24e-141	HotNet2
EGFR	4	1.42e-109	HotNet2/ReMIC
PIK3CA	6	3.70e-38	HotNet2/ReMIC
PIK3R1	8	4.14e-33	HotNet2/ReMIC
NF1	9	1.80e-31	ReMIC
RB1	11	2.84e-28	
ATRX	14	2.12e-17	ReMIC
IDH1	17	3.59e-13	HotNet2
STAG2	33	4.51e-08	ReMIC
CHD8	36	5.77e-07	HotNet2/ReMIC
KDR	64	7.43e-05	HotNet2/ReMIC
RPL5	68	7.43e-05	
PTPN11	158	4.59e-03	ReMIC

continued on next page

continued from previous page

AKAP9	184	2.36e-02	ReMIC
BRAF	194	2.36e-02	ReMIC
BRCA1	195	2.36e-02	ReMIC
CLOCK	202	2.36e-02	HotNet2

(D) IntOGen genes in HNSC results

Gene Symbol	Rank	q-value	also detected by
TP53	1	$p=1.00e-159$	HotNet2
FAT1	3	6.39e-114	ReMIC
CDKN2A	4	9.29e-109	HotNet2/ReMIC
NOTCH1	5	1.54e-78	HotNet2/ReMIC
PIK3CA	6	1.54e-78	HotNet2/ReMIC
CASP8	11	6.94e-37	
NSD1	12	2.24e-31	ReMIC
LAMA2	23	2.34e-21	HotNet2/ReMIC
FBXW7	36	3.44e-15	HotNet2/ReMIC
EP300	42	3.51e-14	
HRAS	43	3.51e-14	
MACF1	44	3.51e-14	HotNet2/ReMIC
NOTCH2	48	4.09e-12	HotNet2/ReMIC
HLA-A	58	3.66e-11	HotNet2
ATR	67	3.06e-10	ReMIC
KALRN	78	2.45e-09	ReMIC
NFE2L2	80	2.45e-09	HotNet2/ReMIC
EPHA2	86	1.91e-08	HotNet2
EGFR	103	1.32e-07	ReMIC
CYLD	116	2.04e-07	
MYH9	135	8.53e-07	ReMIC
TGFBR2	144	8.53e-07	
ARID2	157	4.80e-06	HotNet2/ReMIC
CTCF	160	4.80e-06	HotNet2
ATRX	184	9.47e-06	ReMIC
APC	204	2.63e-05	ReMIC
ATM	205	2.63e-05	ReMIC
HLA-B	214	2.63e-05	HotNet2
RASA1	223	2.63e-05	ReMIC
SMARCA4	225	2.63e-05	ReMIC
SPTAN1	227	2.63e-05	ReMIC
NCOR1	247	8.22e-05	
BAZ2B	266	1.29e-04	ReMIC
KDM6A	282	1.29e-04	HotNet2
TRIO	301	1.29e-04	ReMIC
FN1	324	4.63e-04	ReMIC
CHD9	347	5.76e-04	ReMIC
ARID1B	403	2.33e-03	ReMIC
CUL3	415	2.33e-03	
MEF2C	438	2.33e-03	

continued on next page

continued from previous page

PBRM1	451	2.33e-03	HotNet2
RAC1	461	2.33e-03	ReMIC
TAOK2	470	2.33e-03	
APAF1	523	8.96e-03	
PCDH18	586	8.96e-03	ReMIC
NF1	608	9.44e-03	ReMIC
CIITA	677	2.01e-02	
ARHGAP35	688	2.95e-02	HotNet2/ReMIC
BRCA1	694	2.95e-02	ReMIC
DICER1	704	2.95e-02	HotNet2
DNMT3A	707	2.95e-02	ReMIC
MTOR	753	2.95e-02	ReMIC
PABPC3	765	2.95e-02	ReMIC
WHSC1	810	2.95e-02	
B2M	826	3.58e-02	
ARFGEF2	836	4.07e-02	HotNet2
BRWD1	839	4.07e-02	
CUL1	843	4.07e-02	
HSPA8	853	4.07e-02	

(E) IntOGen genes in LAML results

Gene Symbol	Rank	q-value	also detected by
NPM1	1	5.18e-165	mCGfinder/HotNet2
DNMT3A	2	1.03e-93	mCGfinder/HotNet2/ReMIC
FLT3	3	4.11e-87	mCGfinder/ReMIC
RUNX1	4	2.22e-45	mCGfinder/HotNet2
IDH2	5	1.48e-35	mCGfinder
IDH1	6	1.99e-26	mCGfinder
TP53	7	4.08e-25	mCGfinder/HotNet2
TET2	8	8.60e-11	mCGfinder/HotNet2
PTPN11	9	5.07e-10	mCGfinder/HotNet2/ReMIC
NRAS	10	1.13e-09	mCGfinder/HotNet2
ASXL1	12	1.47e-07	mCGfinder
WT1	13	1.55e-07	mCGfinder/HotNet2
CEBPA	15	3.79e-06	mCGfinder/HotNet2
KIT	16	4.17e-06	mCGfinder/HotNet2
RAD21	17	7.77e-05	mCGfinder/HotNet2
STAG2	18	7.82e-05	mCGfinder/HotNet2/ReMIC
U2AF1	19	6.90e-04	mCGfinder/HotNet2
KRAS	23	1.45e-03	mCGfinder/HotNet2/ReMIC
PHF6	24	2.41e-03	mCGfinder
SUZ12	26	3.38e-03	mCGfinder/HotNet2
PRPF8	30	1.82e-02	mCGfinder/HotNet2



**Supplementary Table S4.** AUC scores of PR-curves of the detection results for IntOGen genes in BRCA, BLCA, GBM, HNSC and LAML.

Database	IntOGen				
Method	BRCA	BLCA	GBM	HNSC	LAML
mCGfinder	13.0%	10.5%	13.2%	9.6%	66.7%
HotNet2	2.4%	1.8%	1.1%	1.9%	9.4%
ReMIC	5.3%	3.2%	3.0%	3.8%	0.6%
Random	1.5%	1.3%	0.6%	1.3%	0.2%
Database	Union(CGC, IntOGen)				
Method	BRCA	BLCA	GBM	HNSC	LAML
mCGfinder	13.0%	10.9%	12.8%	9.9%	26.5%
HotNet2	2.5%	1.9%	1.2%	1.9%	6.7%
ReMIC	5.4%	3.3%	3.2%	4.1%	1.6%
Random	1.5%	1.3%	0.7%	1.4%	0.7%

**Supplementary Table S5.** The detailed information of TCGA somatic mutation datasets of BRCA, BLCA, GBM, HNSC and LAML respectively. The datasets are downloaded from the UCSC Cancer Genomics Browser [11]: <https://genome-cancer.soe.ucsc.edu/proj/site/hgHeatmap/>

Title	TCGA bladder urothelial carcinoma (BLCA) gene-level nonsilent somatic mutation (broad automated)
Dataset	BLCA gene-level mutation (broad automated)
Dataset ID	TCGA.BLCA_mutation_broad_gene
Domain	TCGA
Origin	Bladder
Disease	bladder urothelial carcinoma
Sample Type	tumor
Data Type	somatic mutation
Clinical Cohort	TCGA Bladder Cancer
N	238
Version	2015-02-24
Title	TCGA breast invasive carcinoma (BRCA) gene-level nonsilent somatic mutation (wustl)
Dataset	BRCA gene-level mutation (wustl)
Dataset ID	TCGA.BRCA_mutation_wustl_gene
Domain	TCGA
Origin	Breast
Disease	breast invasive carcinoma
Sample Type	tumor
Data Type	somatic mutation
Clinical Cohort	TCGA Breast Cancer
N	776
Version	2015-02-24
Title	TCGA glioblastoma multiforme (GBM) gene-level nonsilent somatic mutation (broad)
Dataset	GBM gene-level mutation (broad)
Dataset ID	TCGA.GBM_mutation_broad_gene
Domain	TCGA
Origin	Brain
Disease	glioblastoma multiforme
Sample Type	tumor
Data Type	somatic mutation
Clinical Cohort	TCGA Glioblastoma
N	291
Version	2015-02-24

*continued on next page*

continued from previous page

Title	TCGA head & neck squamous cell carcinoma (HNSC) gene-level nonsilent somatic mutation (broad automated)
Dataset	HNSC gene-level mutation (broad automated)
Dataset ID	TCGA_HNSC.mutation_broad_gene
Domain	TCGA
Origin	Head and Neck region
Disease	head & neck squamous cell carcinoma
Sample Type	tumor
Data Type	somatic mutation
Clinical Cohort	TCGA Head and Neck Cancer
N	509
Version	2015-02-24
Title	TCGA acute myeloid leukemia (LAML) gene-level nonsilent somatic mutation (wustl)
Dataset	LAML gene-level mutation (wustl)
Dataset ID	CGA_LAML.mutation_wustl_gene
Domain	TCGA
Origin	White blood cell
Disease	acute myeloid leukemia
Sample Type	tumor
Data Type	somatic mutation
Clinical Cohort	TCGA Acute Myeloid Leukemia
N	197
Version	2015-02-24