

ONCSC 7700 – 013

Special Topic: Cancer Data Science Journal Club

Starting:

Sept 4, 2025

9:30am – 10:30am

Class: Bi-weekly

Course Director: Aik Choon Tan, Ph.D.

**Course Teaching Assistant:
Min Hu**

<http://tanlab.org/teaching/ONCSC7700/>

Class: Thursday 9:30am - 10:30am
Venue: HCI Research South Conference Room 4C

COURSE TOPICS

1. TOPIC 1 - GENE SET ANALYSIS

Classic & Concept: (AC Tan - 9/4/2025)

- GSEA Paper: Subramanian et al. (2005). Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. PNAS. [\[PDF\]](#).
- Mootha et al. (2003). PGC-1-responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes. Nat. Genetics. [\[PDF\]](#)
- Liberzon et al. (2015). The Molecular Signature Database (MSigDB) hallmark gene set collection. Cell Systems. 1(6): 417-425. [\[PDF\]](#)

Variations & Expansions:

- Irizarry et al. (2009). Gene set enrichment analysis made simple. Stat Methods Med Res. 18(6): 565-575. [\[PDF\]](#)
- Hanzelmann et al. (2013). GSVA: gene set variation analysis for microarray and RNA-seq data. BMC Bioinformatics. 14: Article number: 7. [\[PDF\]](#)
- Powers et al. (2018). GSEA-InContext: identifying novel and common patterns in expression experiments. Bioinformatics. 34(13): i555-i564. [\[PDF\]](#)
- Cousins et al. (2023). Gene set proximity analysis: expanding gene set enrichment analysis through learned geometric embeddings, with drug-repurposing applications in COVID-19. Bioinformatics. 39(1): btac735. [\[PDF\]](#)

Applications:

- Ma et al. (2020). Integrative differential expression and gene set enrichment analysis using summary statistics for scRNA-seq studies. Nature Communications. 11: Article number: 1585. [\[PDF\]](#)
- Franchini et al. (2023). Single-cell gene set enrichment analysis and transfer learning for functional annotation of scRNA-seq data. NAR Genomics and Bioinformatics, 5(1): lqad024. [\[PDF\]](#)
- Fan et al. (2024). irGSEA: the integration of single-cell rank-based gene set enrichment analysis. Briefings in Bioinformatics. 25(4): bbae243. [\[PDF\]](#)

AI:

- Joachimiak et al. (2024). Gene Set Summarization using Large Language Models. ArXiv:2305.13338v3. [\[PDF\]](#)
- Wang et al. (2025). GeneAgent: self-verification language agent for gene-set analysis using domain databases. Nature Methods. 22:1677-1685. [\[PDF\]](#)

Presentation Schedule

Date	Presenter	Module
09/04/2025	AC Tan	Gene Set Analysis
09/18/2025	Jenny Ge	Gene Set Analysis
10/02/2025	Min Hu	Gene Set Analysis
10/16/2025	Hui Huang	Gene Set Analysis
10/31/2025	Yemi Imodoye	TBD
11/13/2025	Paulina Jaimes	TBD
11/27/2025	Wontak Kim	TBD
12/11/2025	Swagata Maity	TBD
01/08/2026	Xiandong Peng	TBD

PGC-1 α -responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes

Vamsi K Mootha^{1,2,3,10}, Cecilia M Lindgren^{1,4,10}, Karl-Fredrik Eriksson⁴, Aravind Subramanian¹, Smita Sihag¹, Joseph Lehar¹, Pere Puigserver⁵, Emma Carlsson⁴, Martin Ridderstråle⁴, Esa Laurila⁴, Nicholas Houstis¹, Mark J Daly¹, Nick Patterson¹, Jill P Mesirov¹, Todd R Golub^{1,5}, Pablo Tamayo¹, Bruce Spiegelman⁵, Eric S Lander^{1,6}, Joel N Hirschhorn^{1,7,8}, David Altshuler^{1,2,7,9,11} & Leif C Groop^{4,11}

DNA microarrays can be used to identify gene expression changes characteristic of human disease. This is challenging, however, when relevant differences are subtle at the level of individual genes. We introduce an analytical strategy, Gene Set Enrichment Analysis, designed to detect modest but coordinate changes in the expression of groups of functionally related genes. Using this approach, we identify a set of genes involved in oxidative phosphorylation whose expression is coordinately decreased in human diabetic muscle. Expression of these genes is high at sites of insulin-mediated glucose disposal, activated by PGC-1 α and correlated with total-body aerobic capacity. Our results associate this gene set with clinically important variation in human metabolism and illustrate the value of pathway relationships in the analysis of genomic profiling experiments.

RESULTS

We used DNA microarrays to profile expression of over 22,000 genes in skeletal muscle biopsy samples from 43 age-matched males (**Table 1**), 17 with normal glucose tolerance (NGT), 8 with impaired glucose tolerance (IGT) and 18 with DM2. We obtained samples at the time of diagnosis (before treatment with hypoglycemic medication) and under the controlled conditions of a hyperinsulinemic euglycemic clamp. When assessed with either of two different analytical techniques^{3,6} that take into account the multiple comparisons implicit in microarray analysis, no single gene had a significant difference in expression between the diagnostic categories (data not shown). This result is consistent with smaller studies^{7,8} that did not identify any individual gene whose expression difference was significant when corrected for the large number of hypotheses tested^{9,10}.

Gene Set Enrichment Analysis

To test for sets of related genes that might be systematically altered in diabetic muscle, we devised a simple approach called Gene Set Enrichment Analysis (GSEA), which we introduce here (**Fig. 1**) and describe in more detail elsewhere (A.S. *et al.*, manuscript in preparation). The method combines information from the members of previously defined sets of genes (for example, biological pathways) to increase signal relative to noise and improve statistical power.

Outline

- Methods for Identifying Gene Set Analysis
 - Motivation
 - Statistics
 - Gene Set Enrichment Analysis (GSEA)

Limitations of Candidate Gene Analysis

- Functional genomics technologies such as expression profiling using microarrays provide a global approach to understanding cellular processes in different biological phenotypes.
- Candidate genes analyses
 - Gene lists
 - Number of genes range from hundred to thousands
 - Sifting through gene list is a daunting task to group these genes into functional groups (*ad hoc* analysis)
 - Bias and require expert knowledge

Motivation

- Genes must act in concert to drive various cellular processes.
- Gene expression alterations might be revealed at the level of biological pathways or co-regulated gene sets (functional groups).
- Gene set analysis –
 - more objective and robust.
 - able to discover sets of coordinated differentially expressed genes among pathway members and their association to a specific biological phenotype.
 - provide new insights linking biological phenotypes to their underlying molecular mechanisms.
 - suggesting new hypotheses about pathway membership and connectivity.

How to find Sets of Co-ordinately Differentially Expressed Genes

- High-throughput “omics” data (e.g. microarray gene expression, RNA-seq etc) full with noise.
- Real biological signals might be subtle.
- (*Assumption*) Genes with similar function or participate in a biological process have similar expression patterns.
- *Goal*: Find these sets of genes from high-throughput “omics” data

Gene Set Analysis

- Samples with known biological phenotypes (e.g. *class labels*)
- High-throughput measurements of data points (e.g. *gene expressions*)
- Set of genes involved in biological processes or cellular functions or pathways (e.g. *gene sets*)
- Compare the *gene expressions* of various *class labels* to find differentially expressed *gene sets*.

Gene Set Enrichment Analysis (GSEA)

- *Goal:* to detect modest but coordinated expression changes in pre-specified sets of related genes (gene sets).
- Gene set can be all the genes involved in
 - specific pathway (obtained from Pathway databases such as KEGG, BioCARTA, REACTOME etc)
 - specific gene ontology class (obtained from Gene Ontology category)
 - specific chromosome locations
 - specific transcriptional regulated targets (e.g. transcription factor targets, miRNA targets)
 - specific gene signatures (obtained from published papers or your own experiments)
 - Specific drug targets (obtained from experiments of drug-gene interactions)

Original Publication of GSEA

PGC-1 α -responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes

Vamsi K Mootha^{1,2,3,10}, Cecilia M Lindgren^{1,4,10}, Karl-Fredrik Eriksson⁴, Aravind Subramanian¹, Smita Sihag¹, Joseph Lehar¹, Pere Puigserver⁵, Emma Carlsson⁴, Martin Ridderstråle⁴, Esa Laurila⁴, Nicholas Houstis¹, Mark J Daly¹, Nick Patterson¹, Jill P Mesirov¹, Todd R Golub^{1,5}, Pablo Tamayo¹, Bruce Spiegelman⁵, Eric S Lander^{1,6}, Joel N Hirschhorn^{1,7,8}, David Altshuler^{1,2,7,9,11} & Leif C Groop^{4,11}

DNA microarrays can be used to identify gene expression changes characteristic of human disease. This is challenging, however, when relevant differences are subtle at the level of individual genes. We introduce an analytical strategy, Gene Set Enrichment Analysis, designed to detect modest but coordinate changes in the expression of groups of functionally related genes. Using this approach, we identify a set of genes involved in oxidative phosphorylation whose expression is coordinately decreased in human diabetic muscle. Expression of these genes is high at sites of insulin-mediated glucose disposal, activated by PGC-1 α and correlated with total-body aerobic capacity. Our results associate this gene set with clinically important variation in human metabolism and illustrate the value of pathway relationships in the analysis of genomic profiling experiments.

NATURE GENETICS VOLUME 34 | NUMBER 3 | JULY 2003 267-273

[Citations: >11311]

Original Publication of GSEA

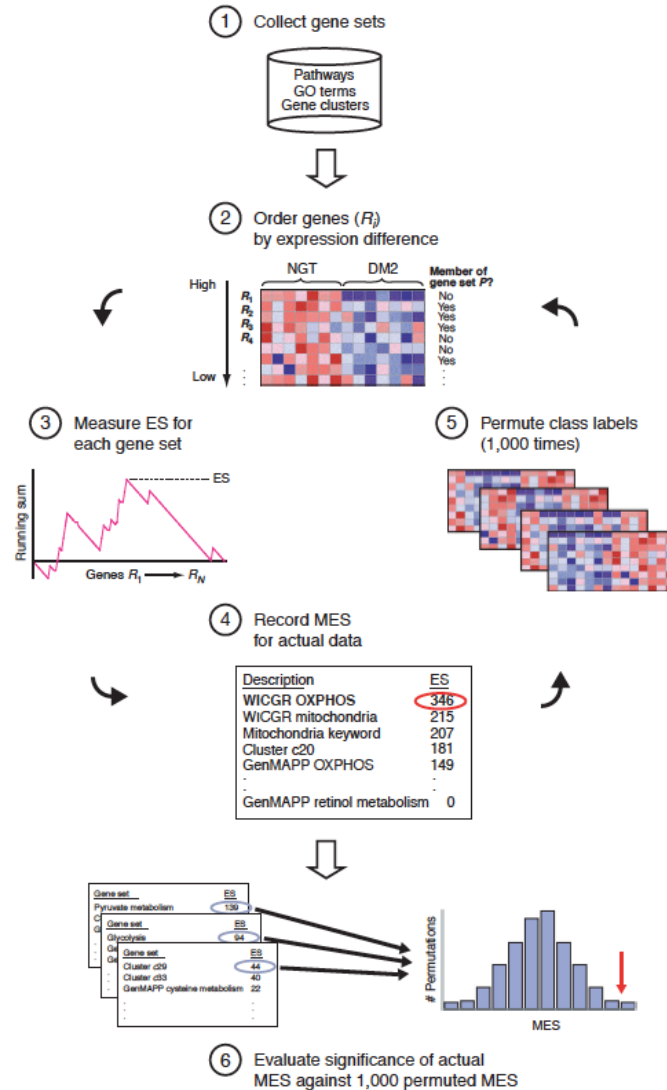


Figure 1 Schematic overview of GSEA. The goal of GSEA is to determine whether any *a priori* defined gene sets (step 1) are enriched at the top of a list of genes ordered on the basis of expression difference between two classes (for example, highly expressed in individuals with NGT versus those with DM2). Genes R_1, \dots, R_N are ordered on the basis of expression difference (step 2) using an appropriate difference measure (for example, SNR). To determine whether the members of a gene set S are enriched at the top of this list (step 3), a Kolmogorov-Smirnov (K-S) running sum statistic is computed: beginning with the top-ranking gene, the running sum increases when a gene annotated to be a member of gene set S is encountered and decreases otherwise. The ES for a single gene set is defined as the greatest positive deviation of the running sum across all N genes. When many members of S appear at the top of the list, ES is high. The ES is computed for every gene set using actual data, and the MES achieved is recorded (step 4). To determine whether one or more of the gene sets are enriched in one diagnostic class relative to the other (step 5), the entire procedure (steps 2–4) is repeated 1,000 times, using permuted diagnostic assignments and building a histogram of the maximum ES achieved by any pathway in a given permutation. The MES achieved using the actual data is then compared to this histogram (step 6, red arrow), providing us with a global P value for assessing whether any gene set is associated with the diagnostic categorization.

Original Publication of GSEA

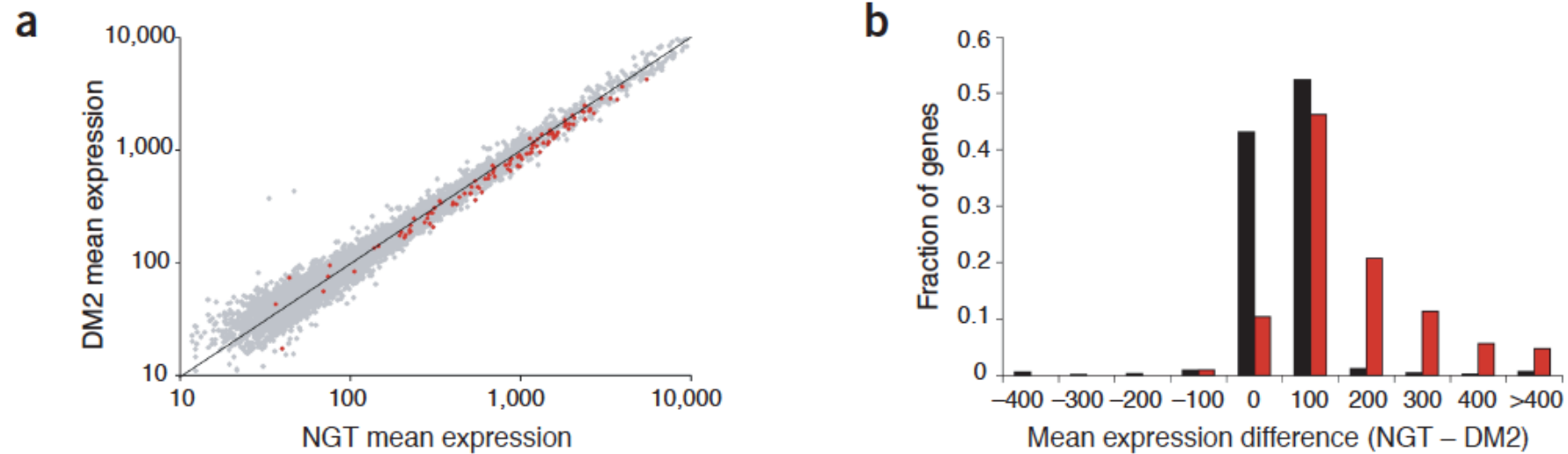


Figure 2 OXPHOS gene expression is reduced in diabetic muscle. (a) The mean expression of all genes (gray) and of OXPHOS genes (red) is plotted for individuals with DM2 versus those with NGT. (b) Histogram of mean gene expression level differences between individuals with NGT and DM2, using the data from a, for all genes (black) and for OXPHOS genes (red).

Original Publication of GSEA

peroxisome proliferator-activated receptor γ coactivator 1 α (PGC-1 α , encoded by PPARGC1)

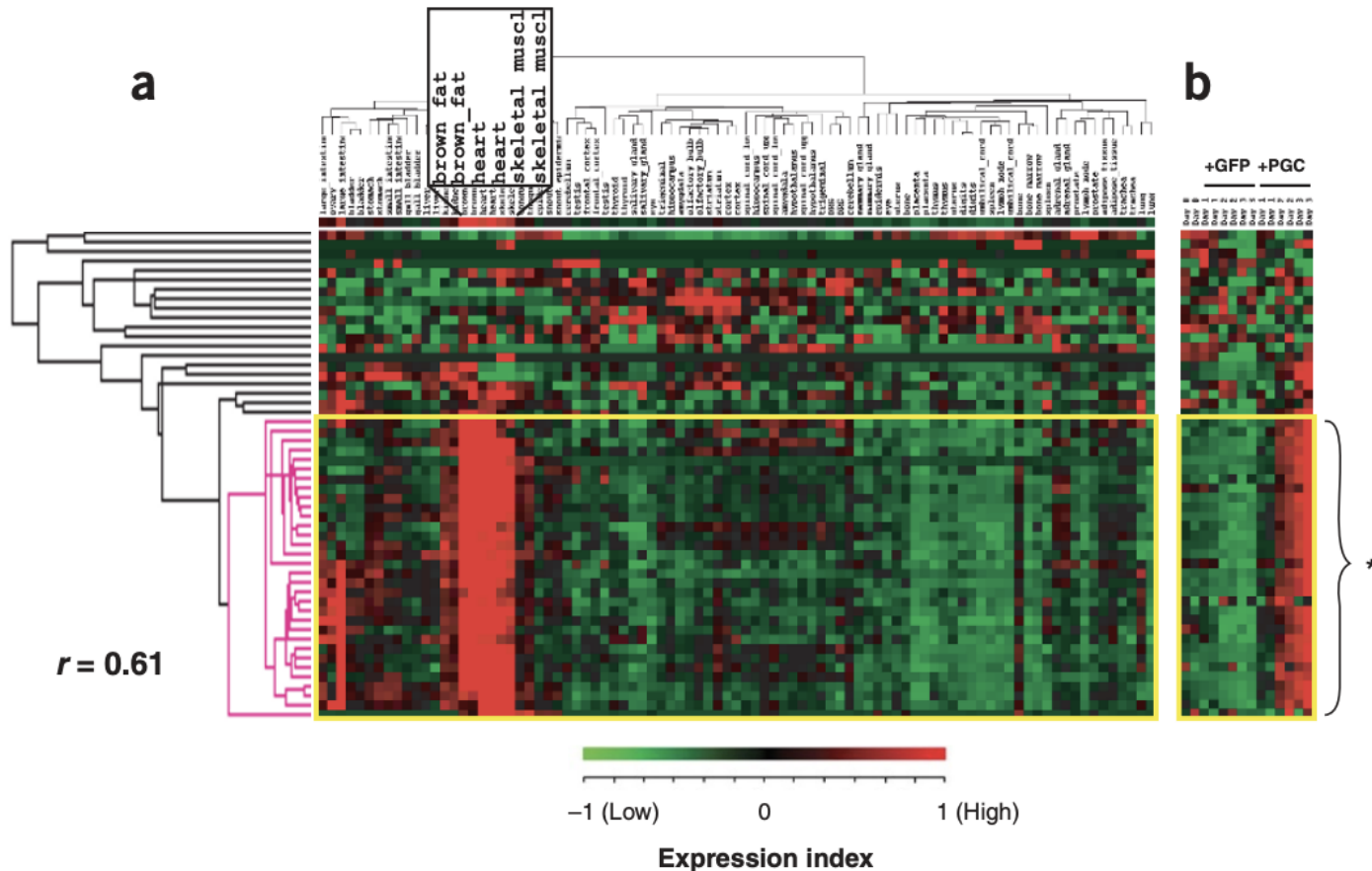


Figure 3 OXPHOS-CR represents a coregulated subset of OXPHOS genes responsive to the transcriptional coactivator PGC-1 α . (a) Normalized expression profile of 52 mouse homologs of the human OXPHOS genes across the mouse expression atlas¹². These 52 genes were hierarchically clustered³². The pink tree on the left corresponds to a subcluster with a correlation coefficient of 0.65. We call the human homologs of these mouse genes the OXPHOS-CR set. The human homologs of this tightly coregulated cluster, marked with an asterisk and delimited with a yellow box, are *ATP5J*, *ATP5L*, *ATP5O*, *COX5B*, *COX6A2*, *COX7A1*, *COX7B*, *COX7C*, *CYC1*, *CYCS*, *GRIM19*, *HSPC051*, *NDUFA2*, *NDUFA5*, *NDUFA7*, *NDUFA8*, *NDUFB3*, *NDUFB5*, *NDUFB6*, *NDUFC1*, *NDUFS2*, *NDUFS3*, *NDUFS5*, *SDHA*, *SDHB*, *UQCRB* and *UQCRC1*. (b) Normalized expression profile of OXPHOS mouse homologs in a mouse skeletal muscle cell line during a 3-d time course in response to PGC-1 α . The expression profile includes infection with control vectors (expressing GFP) or with vectors expressing PGC-1 α before infection (d 0) and 1, 2 and 3 d after adenoviral infection, all done in duplicate.

Original Publication of GSEA

a

Predictor(s)	R^2_{adj}	P value
Diabetes status	0.28	0.0006
¹ OXPHOS-CR	0.22	0.0012
Diabetes status, OXPHOS-CR	0.33	0.0004
² <i>UQCRB</i>	0.31	<0.0001
Diabetes status, <i>UQCRB</i>	0.38	0.0001

¹ Addition of OXPHOS-CR improves the model with $P = 0.05$.

² Addition of *UQCRB* improves the model with $P = 0.03$.

b

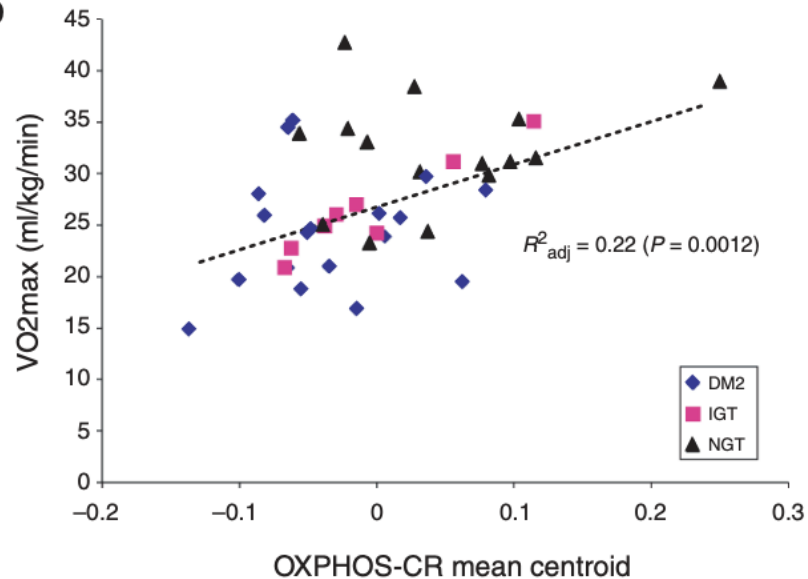


Figure 4 OXPHOS-CR predicts total-body aerobic capacity (VO2max).

(a) Linear regression was used to model VO2max with diabetes status, the mean centroid of OXPHOS-CR gene expression, expression of *UQCRB* or in combination as explanatory (predictor) variables. The explanatory power and significance of the model are shown in the table. (b) Linear regression of VO2max against the mean centroid of OXPHOS-CR gene expression.

GSEA paper (PNAS 2005)

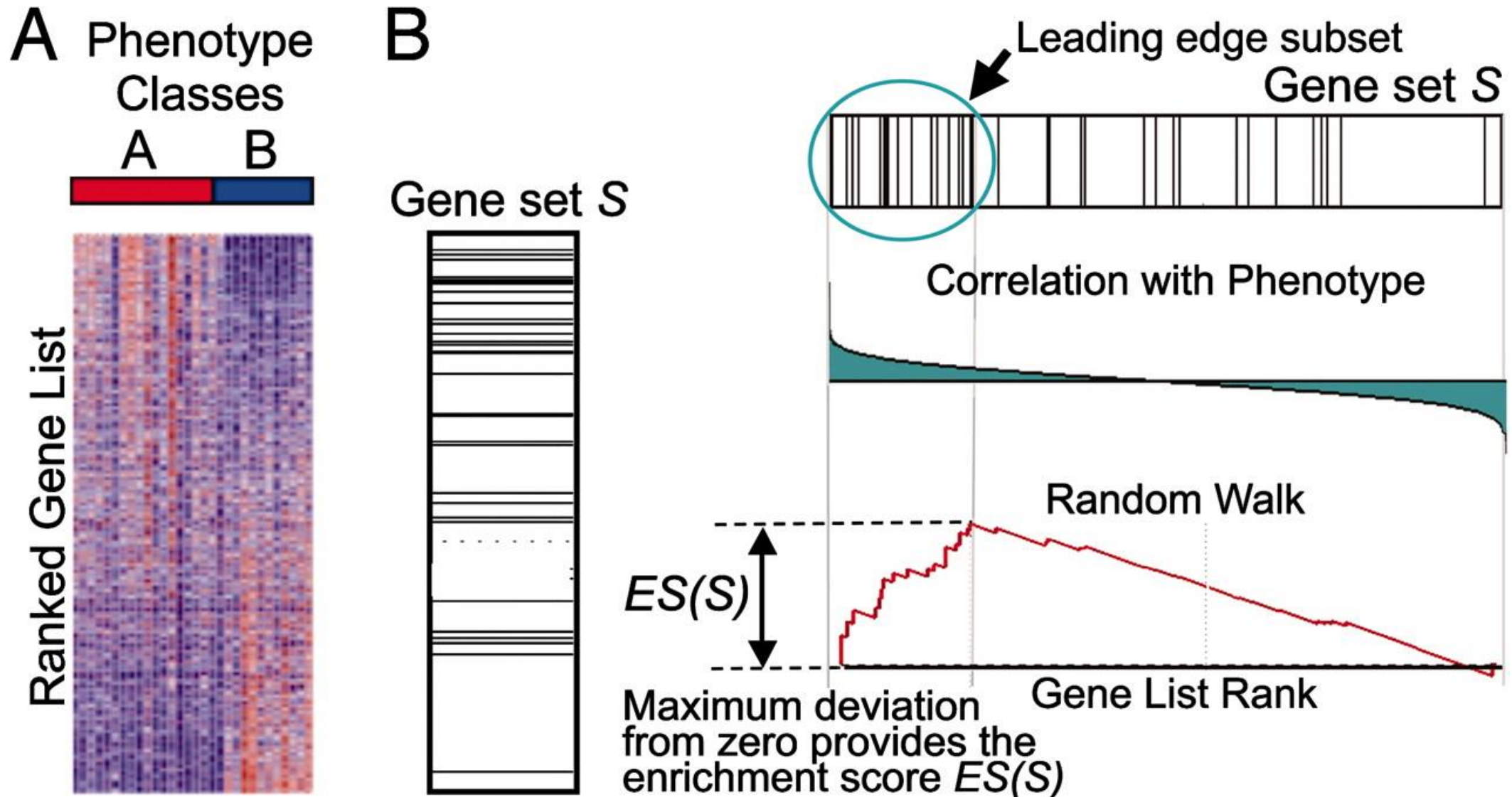
Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles

Aravind Subramanian^{a,b}, Pablo Tamayo^{a,b}, Vamsi K. Mootha^{a,c}, Sayan Mukherjee^d, Benjamin L. Ebert^{a,e}, Michael A. Gillette^{a,f}, Amanda Paulovich^g, Scott L. Pomeroy^h, Todd R. Golub^{a,e}, Eric S. Lander^{a,c,i,j,k}, and Jill P. Mesirov^{a,k}

PNAS October 25, 2005 vol. 102 no. 43 15545–15550

[Citations: >53471]

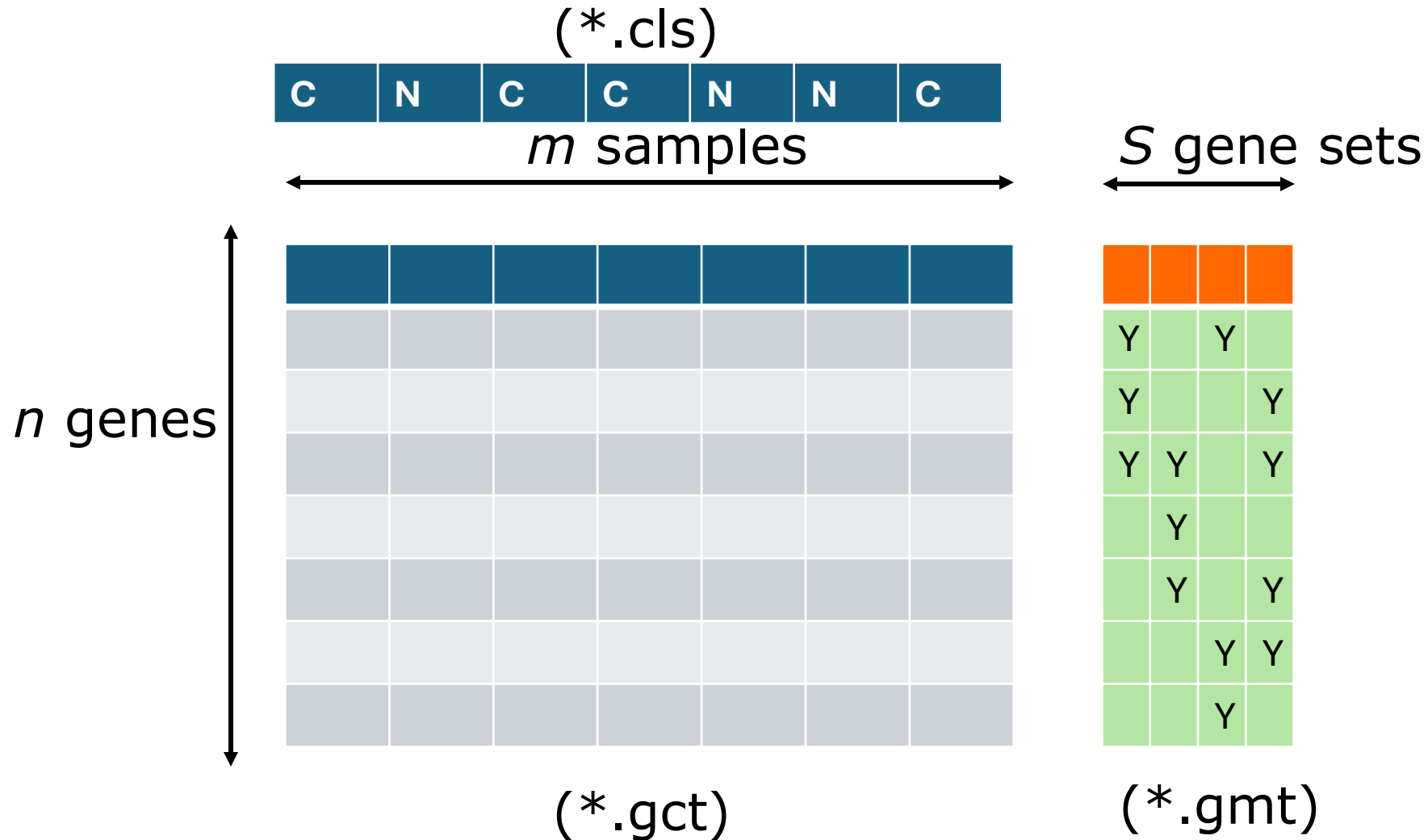
GSEA paper (PNAS 2005)



GSEA Algorithm: Step 1

- Collect gene sets (*.gmt) from databases
- Compile gene expression data (*.gct)
- Define class labels for your samples (*.cls)
- Microarray chip definition file (*.chip) [*Not applicable to RNA-seq*]

Required Files for GSEA

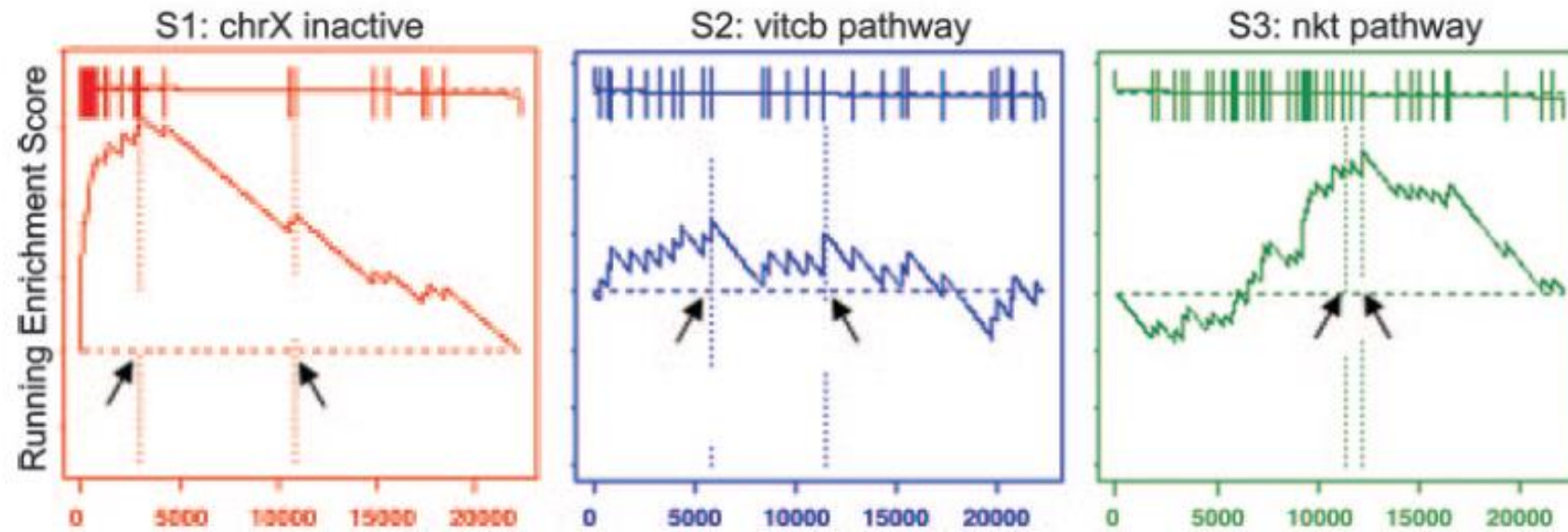


GSEA Algorithm: Step 2

- Rank genes based on their expression differences between the two phenotypes (in GSEA, the measurement is Signal-to-noise ratio)
- Compute Enrichment Score (ES)
 - Compute cumulative sum over ranked genes:
 - Increase sum when gene in set, decrease it otherwise.
 - Magnitude of increment depends on correlation of gene with phenotype.
 - Maximum deviation from zero = enrichment score

GSEA Algorithm: Step 2

Enrichment Plot



Results on SETPATHWAY (Enrichment Plot)

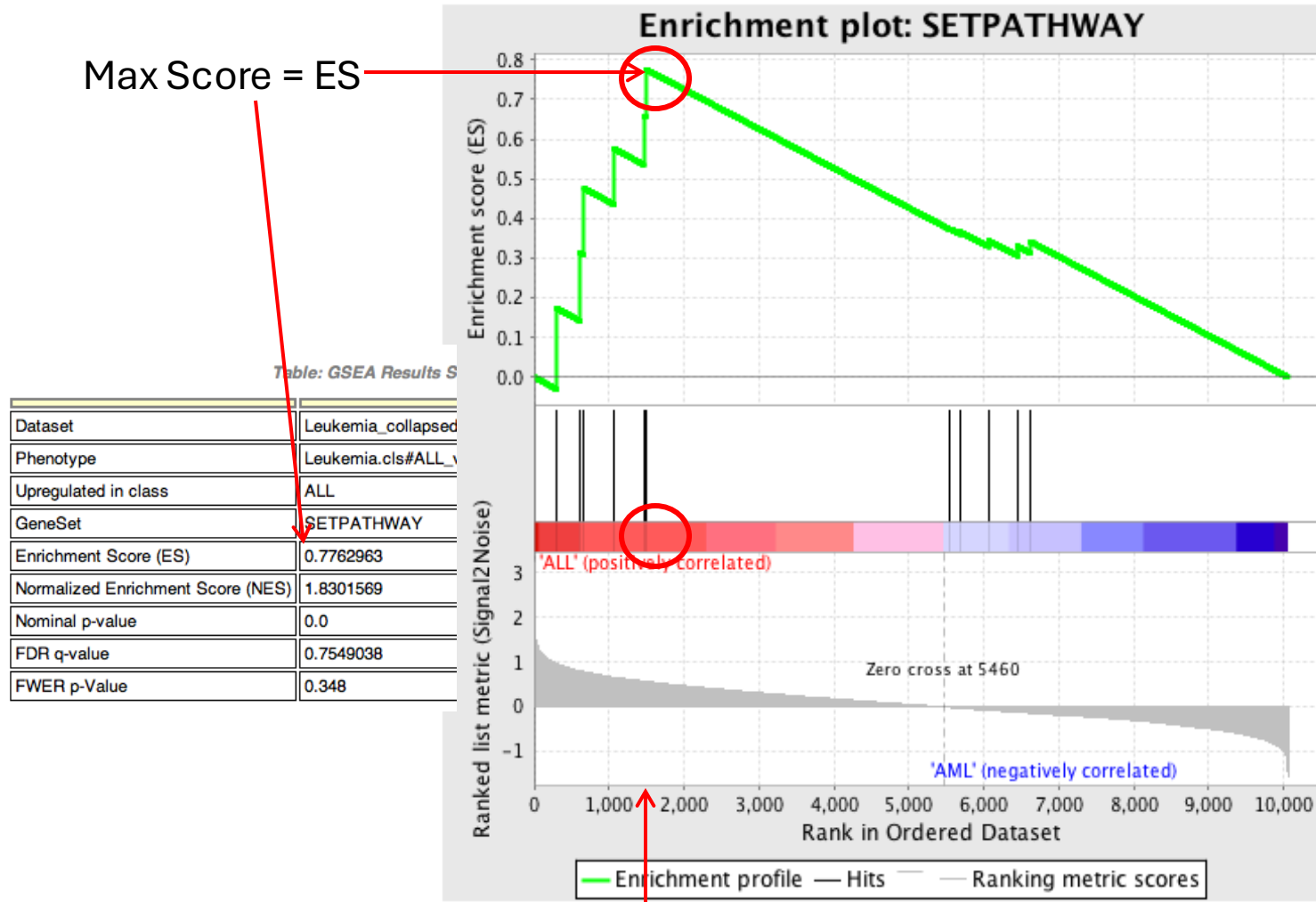


Fig 1: Enrichment plot: SETPATHWAY
Profile of the Running ES Score & Positions of GeneSet Members on the Rank Ordered List

GSEA Algorithm: Step 3

- Compute the Significance of the Gene Sets by Permutation Test
- Permutation Test (n times)
 - Permute phenotype labels
 - Permute gene sets
- For each permutation, compute ES score
- Compare ES score for actual data to distribution of ES scores from permuted data

GSEA Algorithm: Step 4

- Adjustment for multiple hypothesis testing:
 - Normalize the ES accounting for size of each gene set, yielding normalized enrichment score (NES)
 - Control proportion of false positives by calculating FDR corresponding to each NES, by comparing tails of the observed and null distributions for the NES.

GSEA Leading Edge Analysis

- Genes might be involved in different pathways/gene sets.
- Selected core genes might be identified in top enriched gene sets.

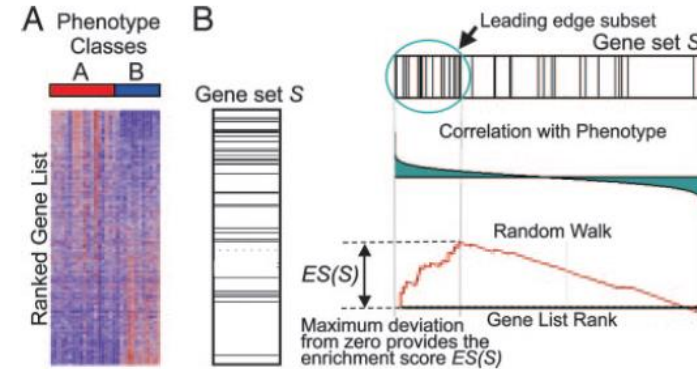
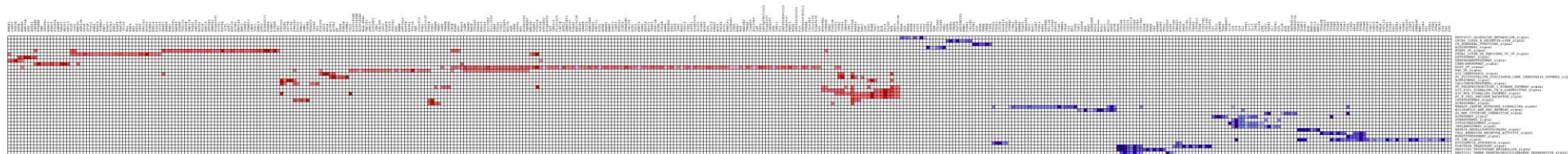
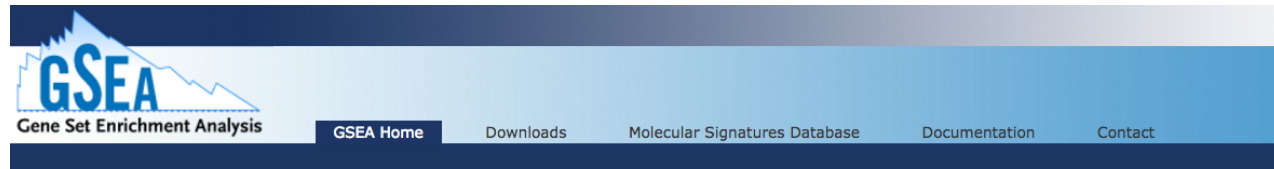


Fig. 1. A GSEA overview illustrating the method. (A) An expression data set sorted by correlation with phenotype, the corresponding heat map, and the "gene tags," i.e., location of genes from a set S within the sorted list. (B) Plot of the running sum for S in the data set, including the location of the maximum enrichment score (ES) and the leading-edge subset.



GSEA Homepage

<http://www.broadinstitute.org/gsea/index.jsp>



Overview

Gene Set Enrichment Analysis (GSEA) is a computational method that determines whether an a priori defined set of genes shows statistically significant, concordant differences between two biological states (e.g. phenotypes).

From this web site, you can:

- **Download** the GSEA software and additional resources to analyze, annotate and interpret enrichment results.
- **Explore the Molecular Signatures Database (MSigDB)**, a collection of annotated gene sets for use with GSEA software.
- **View documentation** describing GSEA and MSigDB.

What's New

19-Sep-2016: The Documentation section of our website is temporarily offline. We are working to resolve the issue and get it back as soon as possible.

15-Aug-2016: The first beta of the next major GSEA Desktop release is available, with SVG plots, Cytoscape 3.3+ support and [much more](#).

17-Jun-2016: You can now follow [@GSEA_MSigDB](#) on Twitter!

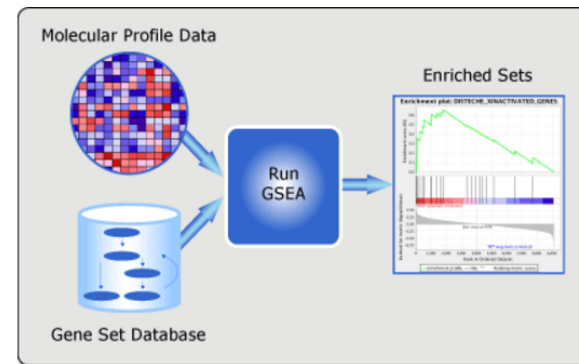
29-Feb-2016: The Sunday 28-Feb-2016 maintenance is complete on the GSEA/MSigDB website. Thanks for your patience!

13-Jan-2016: Version 5.1 of the Molecular Signatures Database (MSigDB) is now available. It includes the addition of 2,962 gene sets to the C7 collection of immunologic signatures, as well as a number of updates and corrections. See the [Release Notes](#) for details.

23-Dec-2015: Our [paper](#) describing the generation of the Hallmarks collection and examples of its use for GSEA was published in Cell Systems.

10-Dec-2015: We have confirmed that GSEA v2.2.0 and newer are compatible with Java 8 and produce equivalent results. Its use is highly recommended.

[Follow @GSEA_MSigDB](#)



Registration

Please [register](#) to download the GSEA software and view the MSigDB gene sets. After registering, you can log in at any time using your email address. Registration is free. Its only purpose is to help us track usage for reports to our funding agencies.

Contributors

GSEA and MSigDB are maintained by the [GSEA team](#) with the support of our MSigDB Scientific Advisory Board. Our thanks to our many contributors. Funded by: [National Cancer Institute](#), [National Institutes of Health](#), [National Institute of General Medical Sciences](#).



Citing GSEA



To cite your use of the GSEA software, please reference Subramanian, Tamayo, et al. (2005, *PNAS* 102, 15545-15550) and Mootha, Lindgren, et al. (2003, *Nat Genet* 34, 267-273).

GSEA Implementation

<http://www.broadinstitute.org/gsea/index.jsp>

Software

There are several options for GSEA software. All options implement exactly the same algorithm. Usage recommendations and installation instructions are listed below. Current Java implementations of GSEA require Java 7 or 8. Java 8 is recommended.

javaGSEA Desktop Application	<ul style="list-style-type: none">▶ Easy-to-use graphical user interface▶ Runs on any desktop computer (Windows, Mac OS X, Linux etc.) that supports Java 7 or 8. Java 8 is recommended.▶ Produces richly annotated reports of enrichment results▶ Integrated gene sets browser to view gene set annotations, search for gene sets and map gene sets between platforms	Launch with 1GB (for 32 or 64-bit Java)  memory: 
javaGSEA Java Jar file	<ul style="list-style-type: none">▶ Command line or offline usage. See our User Guide for details.▶ Runs on any platform that supports Java 7 or 8. Java 8 is recommended.▶ We recommend using the 'Launch' buttons above instead of this mode for most users	download gsea2-2.2.2.jar
BETA javaGSEA Java Jar file	<ul style="list-style-type: none">▶ The first Beta version of the next major GSEA Desktop release, with SVG plots, Cytoscape 3.3+ support for the Enrichment Map, and much more.▶ Our tests show this Beta version produces equivalent results, but use the Production version if you have concerns. At a minimum, verification with the Production version before publication is strongly recommended.▶ Please contact us with bugs or other feedback. We will aim to address problems as soon as possible in future Beta releases.▶ Runs only on the command line. See our User Guide for details.▶ Runs on any platform that supports Java 7 or 8. Java 8 is recommended.	BETA download gsea2-3.0_beta_1.jar
R-GSEA R Script	<ul style="list-style-type: none">▶ Usage from within the R programming environment▶ Easily inspect, learn and tweak the algorithm▶ Incorporate GSEA into your own data analysis pipeline▶ Programmatically call the open source GSEA R API▶ Note that this script has not been updated since 2005 and may not work as-is with modern R distributions.▶ Click here to learn more about the R-GSEA script	download GSEA-P-R.1.0.zip
GenePattern GSEA Module	<ul style="list-style-type: none">▶ Use GSEA from within GenePattern▶ Use GSEA in concert with a large suite of other analytics found in GenePattern (a powerful and flexible analysis platform developed at the Broad Institute)	GenePattern site

For details on the GSEA algorithm and software refer to the [Documentation](#).
For details on the latest release refer to the [Release Notes](#).

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>



MSigDB
Molecular Signatures
Database

Molecular Signatures Database v5.1

Overview

The Molecular Signatures Database (MSigDB) is a collection of annotated gene sets for use with GSEA software. From this web site, you can

- ▶ **Search** for gene sets by keyword.
- ▶ **Browse** gene sets by name or collection.
- ▶ **Examine** a gene set and its annotations. See, for example, the [ANGIOGENESIS](#) gene set page.
- ▶ **Download** gene sets.
- ▶ **Investigate** gene sets:
 - ▶ **Compute overlaps** between your gene set and gene sets in MSigDB.
 - ▶ **Categorize** members of a gene set by gene families.
 - ▶ **View the expression profile** of a gene set in any of the three provided public expression compendia.

Registration

Please [register](#) to download the GSEA software and view the MSigDB gene sets. After registering, you can log in at any time using your email address. Registration is free. Its only purpose is to help us track usage for reports to our funding agencies.

Current Version

MSigDB database v5.1 updated January 2016. [Release notes](#).
GSEA/MSigDB web site v5.0 released March 2015

Contributors

The MSigDB is maintained by the [GSEA team](#) with the support of our [MSigDB Scientific Advisory Board](#). We also welcome and appreciate contributions to this shared resource and encourage users to submit their gene sets to genesets@broadinstitute.org. Our thanks to our many [contributors](#).

Funded by: [National Cancer Institute](#), [National Institutes of Health](#), [National Institute of General Medical Sciences](#).



Collections

The MSigDB gene sets are divided into 8 major collections:

H **hallmark gene sets** are coherently expressed signatures derived by aggregating many MSigDB gene sets to represent well-defined biological states or processes.

C1 **positional gene sets** for each human chromosome and cytogenetic band.

C2 **curated gene sets** from online pathway databases, publications in PubMed, and knowledge of domain experts.

C3 **motif gene sets** based on conserved cis-regulatory motifs from a comparative analysis of the human, mouse, rat, and dog genomes.

C4 **computational gene sets** defined by mining large collections of cancer-oriented microarray data.

C5 **GO gene sets** consist of genes annotated by the same GO terms.

C6 **oncogenic signatures** defined directly from microarray gene expression data from cancer gene perturbations.

C7 **immunologic signatures** defined directly from microarray gene expression data from immunologic studies.

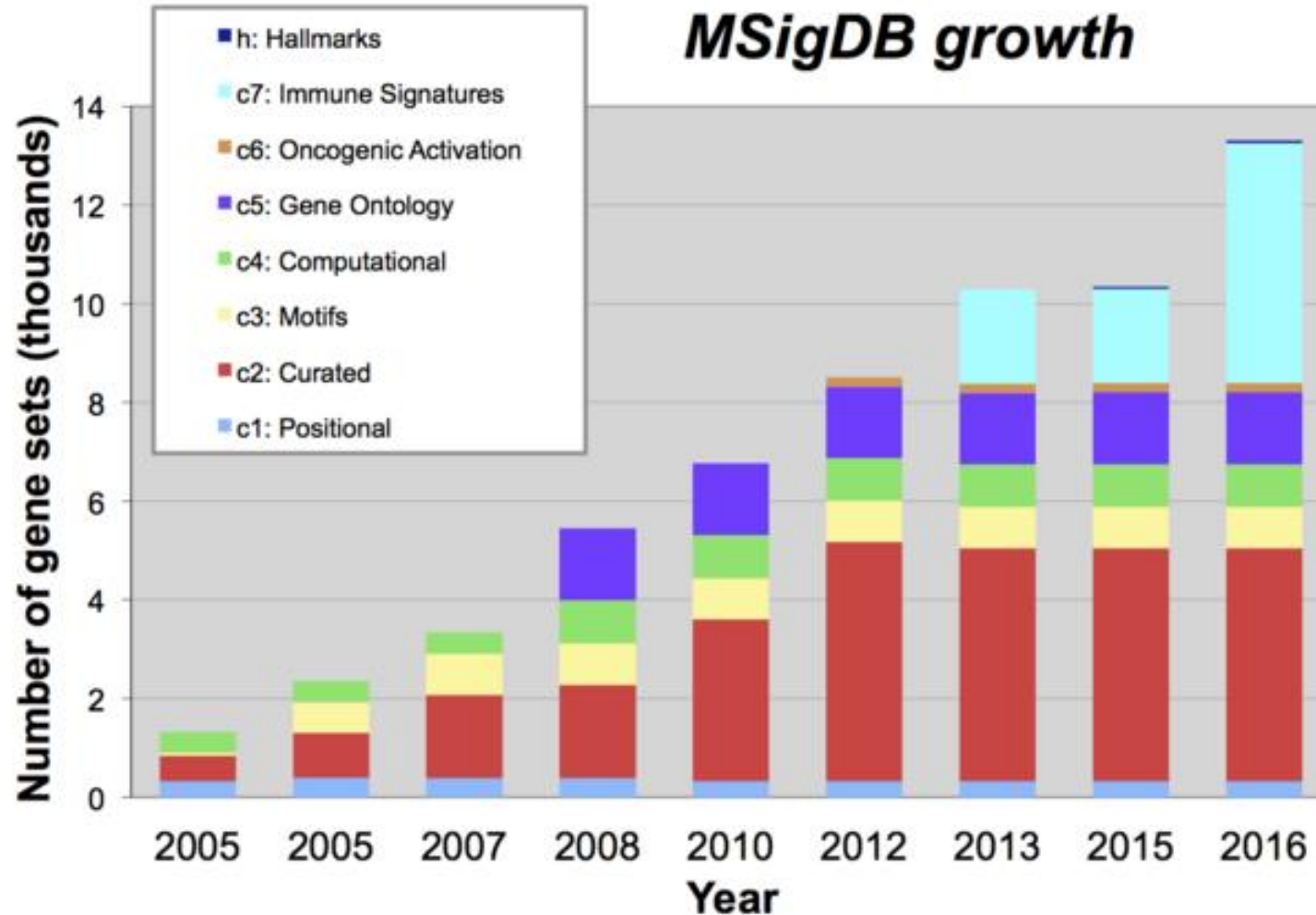
Citing the MSigDB

To cite your use of the Molecular Signatures Database (MSigDB), please reference Subramanian, Tamayo, et al. (2005, *PNAS* 102, 15545-15550) and also the source for the gene set as listed on the gene set page.

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Currently contains 13,311 gene sets organized into 8 categories.



Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Browse Gene Sets





















Gene set name:

By first letter:

1 2 3 4 5 6 7 8 9 0 **A** B C D E F G H I J K L M N O P Q R S T U V W X Y Z

By collection:

[\[about the MSigDB collections\]](#)

- ▶ **H** (hallmark gene sets, 50 gene sets) 
- ▶ **C1** (positional gene sets, 326 gene sets) 
 - ▶ by chromosome: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 X Y
- ▶ **C2** (curated gene sets, 4726 gene sets) 
 - ▶ **CGP** (chemical and genetic perturbations, 3396 gene sets) 
 - ▶ **CP** (Canonical pathways, 1330 gene sets) 
 - ▶ **CP:BIOCARTA** (BioCarta gene sets, 217 gene sets) 
 - ▶ **CP:KEGG** (KEGG gene sets, 186 gene sets) 
 - ▶ **CP:REACTOME** (Reactome gene sets, 674 gene sets) 
- ▶ **C3** (motif gene sets, 836 gene sets) 
 - ▶ **MIR** (microRNA targets, 221 gene sets) 
 - ▶ **TFT** (transcription factor targets, 615 gene sets) 
- ▶ **C4** (computational gene sets, 858 gene sets) 
 - ▶ **CGN** (cancer gene neighborhoods, 427 gene sets) 
 - ▶ **CM** (cancer modules, 431 gene sets) 
- ▶ **C5** (GO gene sets, 1454 gene sets) 
 - ▶ **BP** (GO biological process, 825 gene sets) 
 - ▶ **CC** (GO cellular component, 233 gene sets) 
 - ▶ **MF** (GO molecular function, 396 gene sets) 
- ▶ **C6** (oncogenic signatures, 189 gene sets) 
- ▶ **C7** (immunologic signatures, 4872 gene sets) 

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Click on a gene set name to view its gene set page.

[HALLMARK_ADIPOGENESIS](#)
[HALLMARK_ALLOGRAFT_REJECTION](#)
[HALLMARK_ANDROGEN_RESPONSE](#)
[HALLMARK_ANGIOGENESIS](#)
[HALLMARK_APICAL_JUNCTION](#)
[HALLMARK_APICAL_SURFACE](#)
[HALLMARK_APOPTOSIS](#)
[HALLMARK_BILE_ACID_METABOLISM](#)
[HALLMARK_CHOLESTEROL_HOMEOSTASIS](#)
[HALLMARK_COAGULATION](#)
[HALLMARK_COMPLEMENT](#)
[HALLMARK_DNA_REPAIR](#)
[HALLMARK_E2F_TARGETS](#)
[HALLMARK_EPITHELIAL_MESENCHYMAL_TRANSITION](#)
[HALLMARK_ESTROGEN_RESPONSE_EARLY](#)
[HALLMARK_ESTROGEN_RESPONSE_LATE](#)
[HALLMARK_FATTY_ACID_METABOLISM](#)

[HALLMARK_G2M_CHECKPOINT](#)
[HALLMARK_GLYCOLYSIS](#)
[HALLMARK_HEDGEHOG_SIGNALING](#)
[HALLMARK_HEME_METABOLISM](#)
[HALLMARK_HYPOXIA](#)
[HALLMARK_IL2_STAT5_SIGNALING](#)
[HALLMARK_IL6_JAK_STAT3_SIGNALING](#)
[HALLMARK_INFLAMMATORY_RESPONSE](#)
[HALLMARK_INTERFERON_ALPHA_RESPONSE](#)
[HALLMARK_INTERFERON_GAMMA_RESPONSE](#)
[HALLMARK_KRAS_SIGNALING_DN](#)
[HALLMARK_KRAS_SIGNALING_UP](#)
[HALLMARK_MITOTIC_SPINDLE](#)
[HALLMARK_MTORC1_SIGNALING](#)
[HALLMARK_MYC_TARGETS_V1](#)
[HALLMARK_MYC_TARGETS_V2](#)
[HALLMARK_MYOGENESIS](#)

[HALLMARK_NOTCH_SIGNALING](#)
[HALLMARK_OXIDATIVE_PHOSPHORYLATION](#)
[HALLMARK_P53_PATHWAY](#)
[HALLMARK_PANCREAS_BETA_CELLS](#)
[HALLMARK_PEROXISOME](#)
[HALLMARK_PI3K_AKT_MTOR_SIGNALING](#)
[HALLMARK_PROTEIN_SECRETION](#)
[HALLMARK_REACTIVE_OXYGEN_SPECIES_PATHWAY](#)
[HALLMARK_SPERMATOGENESIS](#)
[HALLMARK_TGF_BETA_SIGNALING](#)
[HALLMARK_TNFA_SIGNALING_VIA_NFKB](#)
[HALLMARK_UNFOLDED_PROTEIN_RESPONSE](#)
[HALLMARK_UV_RESPONSE_DN](#)
[HALLMARK_UV_RESPONSE_UP](#)
[HALLMARK_WNT_BETA_CATENIN_SIGNALING](#)
[HALLMARK_XENOBIOTIC_METABOLISM](#)

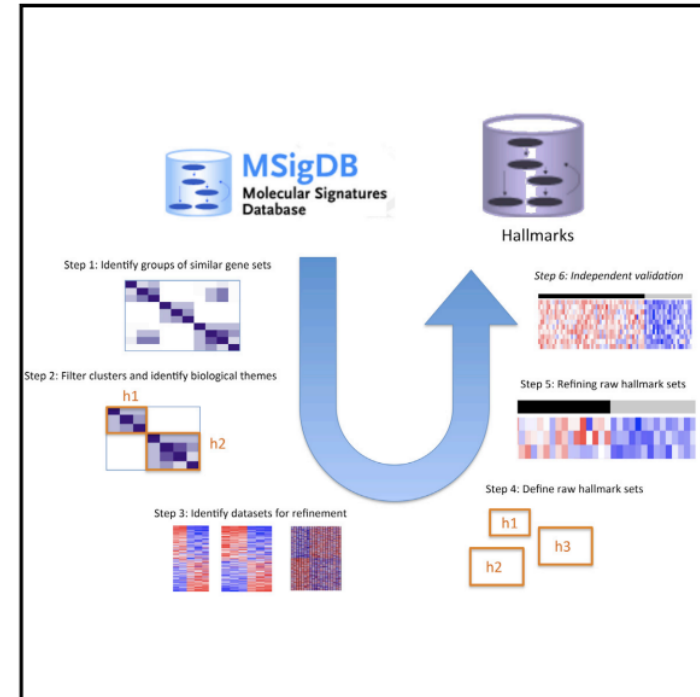
MSigDB Hallmark Gene Set

Cell Systems

Report

The Molecular Signatures Database Hallmark Gene Set Collection

Graphical Abstract



Authors

Arthur Liberzon, Chet Birger,
Helga Thorvaldsdóttir,
Mahmoud Ghandi, Jill P. Mesirov,
Pablo Tamayo

Correspondence

jmesirov@ucsd.edu (J.P.M.),
ptamayo@ucsd.edu (P.T.)

In Brief

Through extensive automated and manual curation, Liberzon et al. provide a refined and concise collection of “hallmark” gene sets from the Molecular Signatures Database for gene set enrichment analysis.

[Citations: >11800]

Highlights

- We generate 50 “hallmark” gene sets from the Molecular Signature Database (MSigDB)
- This required a hybrid approach combining computation with manual expert curation
- The hallmarks reduce redundancy and produce more robust enrichment analysis results
- We plan to move forward with a program to enhance and expand the hallmarks collection

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Gene Set: HALLMARK_KRAS_SIGNALING_UP

Standard name	HALLMARK_KRAS_SIGNALING_UP
Systematic name	M5953
Brief description	Genes up-regulated by KRAS activation.
Full description or abstract	
Collection	H: hallmark gene sets
Source publication	
Exact source	
Related gene sets	(show 14 founder gene sets for this hallmark gene set)
External links	
Organism	Homo sapiens
Contributed by	Arthur Liberzon (Broad Institute)
Source platform	HUMAN_GENE_SYMBOL
Dataset references	(show 5 hallmark refinement datasets) (show 1 hallmark validation datasets)
Download gene set	format: grp text gmt gmx xml
Compute overlaps ?	(show collections to investigate for overlap with this gene set)
Compendia expression profiles ?	Human tissue compendium (Novartis) Global Cancer Map (Broad Institute) NCI-60 cell lines (National Cancer Institute)
Advanced query	Further investigate these 200 genes
Gene families ?	Categorize these 200 genes by gene family
Show members	(show 200 members mapped to 200 genes)
Version history	5.0: First introduced

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Show members


(hide 200 members mapped to 200 genes)

Original Member	Entrez Gene Id	Gene Symbol	Gene Description
ABCB1	5243	ABCB1	ATP-binding cassette, sub-family B (MD...
ACE	1636	ACE	angiotensin I converting enzyme (pepti...
ADAM17	6868	ADAM17	ADAM metallopeptidase domain 17
ADAM8	101	ADAM8	ADAM metallopeptidase domain 8
ADAMDEC1	27299	ADAMDEC1	ADAM-like, decysin 1
AKAP12	9590	AKAP12	A kinase (PRKA) anchor protein 12
AKT2	208	AKT2	v-akt murine thymoma viral oncogene ho...
ALDH1A2	8854	ALDH1A2	aldehyde dehydrogenase 1 family, membe...
ALDH1A3	220	ALDH1A3	aldehyde dehydrogenase 1 family, membe...
AMMECR1	9949	AMMECR1	Alport syndrome, mental retardation, m...
ANGPTL4	51129	ANGPTL4	angiopoietin-like 4
ANKH	56172	ANKH	ankylosis, progressive homolog (mouse)
ANO1	55107	ANO1	anoctamin 1, calcium activated chlorid...
ANXA10	11199	ANXA10	annexin A10
APOD	347	APOD	apolipoprotein D
ARG1	383	ARG1	arginase, liver
ATG10	83734	ATG10	ATG10 autophagy related 10 homolog (S....
AVL9	23080	AVL9	AVL9 homolog (S. cerevisiae)
BIRC3	330	BIRC3	baculoviral IAP repeat containing 3
BMP2	650	BMP2	bone morphogenetic protein 2
BPGM	669	BPGM	2,3-bisphosphoglycerate mutase
BTBD3	22903	BTBD3	BTB (POZ) domain containing 3
BTC	685	BTC	betacellulin

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Search Gene Sets

Search by keyword, collection, organism, or contributor. 

See the [Browse Gene Sets page](#) for an alphabetical list of gene sets and collections, or to search by gene set name.

Keywords:

*(supports boolean operators AND and OR,
and
wildcard searches with *)*

Search Filters:

collection

- all collections
- H: hallmark gene sets
- C1: positional gene sets
- C2: curated gene sets
- CGP: chemical and genetic perturbations
- CP: Canonical pathways
- CP:BiOCARTA: BioCarta gene sets
- CP:KEGG: KEGG gene sets
- CP:REACTOME: Reactome gene sets
- C3: motif gene sets

organism

- all organisms
- Danio rerio
- Homo sapiens
- Macaca mulatta
- Mus musculus
- Rattus norvegicus

contributor

- all contributors
- Aristoteles University of Thessaloniki
- BioCarta
- Broad Institute
- Columbia University
- Dana-Farber Cancer Institute
- Giannina Gaslini Institute
- GO
- Johns Hopkins University School of Medicine
- KEGG

control-click to select multiple lines

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

click on rows to select gene sets, click a gene set name to view the gene set page

☐ select all 305 ☒ 0 gene sets selected Select An Action...

<< < **1** 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 > >> 10

name	# genes	description	collections	organism	contributor
AACTTT_UNKNOWN	1890	Genes with promoter regions [-2kb,2kb] around transcription start site containing motif AACTTT. Motif does not match any known transcription factor	C3 TFT	Homo sapiens	Broad Institute
ABE_VEGFA_TARGETS_30MIN	29	Genes up-regulated in HUVEC cells (endothelium) at 30 min after VEGFA [GeneID=7422] stimulation.	C2 CGP	Homo sapiens	University of Washington
ACCATTT,MIR-522	160	Targets of MicroRNA ACCATTT,MIR-522	C3 MIR	Homo sapiens	Broad Institute
ACTGAAA,MIR-30A-3P,MIR-30E-3P	201	Targets of MicroRNA ACTGAAA,MIR-30A-3P,MIR-30E-3P	C3 MIR	Homo sapiens	Broad Institute
AGCATTA,MIR-155	134	Targets of MicroRNA AGCATTA,MIR-155	C3 MIR	Homo sapiens	Broad Institute
ATAGGAA,MIR-202	102	Targets of MicroRNA ATAGGAA,MIR-202	C3 MIR	Homo sapiens	Broad Institute
ATATGCA,MIR-448	212	Targets of MicroRNA ATATGCA,MIR-448	C3 MIR	Homo sapiens	Broad Institute
ATGCAGT,MIR-217	115	Targets of MicroRNA ATGCAGT,MIR-217	C3 MIR	Homo sapiens	Broad Institute
BENPORATH_CYCLING_GENES	648	Genes showing cell-cycle stage-specific expression [PMID=12058064].	C2 CGP	Homo sapiens	Broad Institute
BIOCARTA_TEL_PATHWAY	18	Telomeres, Telomerase, Cellular Aging, and Immortality	C2 CP CP:BIOCARTA	Homo sapiens	BioCarta

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Investigate Gene Sets

Gain further insight into the biology behind a gene set by using the following tools:

- ▶ **compute overlaps** with other gene sets in MSigDB ([more...](#))
- ▶ **display the gene set expression profile** based on a selected compendium of expression data ([more...](#))
- ▶ **categorize** members of the gene set by gene families ([more...](#))

Gene Identifiers

Compute Overlaps

- ☐ **H:** hallmark gene sets [?](#)
- ☐ **C1:** positional gene sets [?](#)
- ☐ **C2:** curated gene sets [?](#)
 - ☐ **CGP:** chemical and genetic perturbations [?](#)
 - ☐ **CP:** Canonical pathways [?](#)
 - ☐ **CP:BIOCARTA:** BioCarta gene sets [?](#)
 - ☐ **CP:KEGG:** KEGG gene sets [?](#)
 - ☐ **CP:REACTOME:** Reactome gene sets [?](#)
- ☐ **C3:** motif gene sets [?](#)
 - ☐ **MIR:** microRNA targets [?](#)
 - ☐ **TFT:** transcription factor targets [?](#)
- ☐ **C4:** computational gene sets [?](#)
 - ☐ **CGN:** cancer gene neighborhoods [?](#)
 - ☐ **CM:** cancer modules [?](#)
- ☐ **C5:** GO gene sets [?](#)
 - ☐ **BP:** GO biological process [?](#)
 - ☐ **CC:** GO cellular component [?](#)
 - ☐ **MF:** GO molecular function [?](#)
- ☐ **C6:** oncogenic signatures [?](#)
- ☐ **C7:** immunologic signatures [?](#)

show top 10 genesets

with FDR q-value below

[compute overlaps](#)

Compendia expression profiles

- ☒ Human tissue compendium (Novartis)
- ☐ Global Cancer Map (Broad Institute)
- ☐ NCI-60 cell lines (National Cancer Institute)

[display expression profile](#)

Gene families

[show gene families](#)

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Investigate Gene Sets

Gain further insight into the biology behind a gene set by using the following tools:

- **compute overlaps** with other gene sets in MSigDB ([more...](#))
- **display the gene set expression profile** based on a selected compendium of expression data ([more...](#))
- **categorize** members of the gene set by gene families ([more...](#))

Gene Identifiers

KRAS
BRAF
NRAS
MAP2K1
MAP2K2
MAPK1
CCND1
PIK3CA
AKT1
PTEN
MTOR

Compute Overlaps

- ☒ **H:** hallmark gene sets [?](#)
- ☐ **C1:** positional gene sets [?](#)
- ☒ **C2:** curated gene sets [?](#)
 - ☒ **CGP:** chemical and genetic perturbations [?](#)
 - ☒ **CP:** Canonical pathways [?](#)
 - ☒ **CP:BIOCARTA:** BioCarta gene sets [?](#)
 - ☒ **CP:KEGG:** KEGG gene sets [?](#)
 - ☒ **CP:REACTOME:** Reactome gene sets [?](#)
- ☐ **C3:** motif gene sets [?](#)
 - ☐ **MIR:** microRNA targets [?](#)
 - ☐ **TFT:** transcription factor targets [?](#)
- ☐ **C4:** computational gene sets [?](#)
 - ☐ **CGN:** cancer gene neighborhoods [?](#)
 - ☐ **CM:** cancer modules [?](#)
- ☐ **C5:** GO gene sets [?](#)
 - ☐ **BP:** GO biological process [?](#)
 - ☐ **CC:** GO cellular component [?](#)
 - ☐ **MF:** GO molecular function [?](#)
- ☐ **C6:** oncogenic signatures [?](#)
- ☐ **C7:** immunologic signatures [?](#)

show genesets

with FDR q-value below

[compute overlaps](#)

Compendia expression profiles

- ☒ Human tissue compendium (Novartis)
- ☐ Global Cancer Map (Broad Institute)
- ☐ NCI-60 cell lines (National Cancer Institute)

[display expression profile](#)

Gene families

[show gene families](#)

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>











Compute Overlaps for Selected Genes

Collections	# Overlaps Shown	# Gene Sets in Collections	# Genes in Comparison (n)	# Genes in Universe (N)
C2, H	10	4776	11	45956

Click the gene set name to see the gene set page. Click the number of genes [in brackets] to download the list of genes.

Color bar shading from light green to black, where lighter colors indicate more significant FDR q-values (< 0.05) and black indicates less significant FDR q-values (≥ 0.05).

Save to: [Excel](#) | [GenomeSpace](#)

Gene Set Name [# Genes (K)]	Description	# Genes in Overlap (k)	k/K	p-value ?	FDR q-value ?
KEGG_GLIOMA [65]	Glioma	11		1.85×10^{-32}	8.85×10^{-29}
KEGG_PROSTATE_CANCER [89]	Prostate cancer	11		7.56×10^{-31}	1.8×10^{-27}
KEGG_ENDOMETRIAL_CANCER [52]	Endometrial cancer	10		1.5×10^{-29}	2.39×10^{-26}
KEGG_ACUTE_MYELOID_LEUKEMIA [60]	Acute myeloid leukemia	10		7.16×10^{-29}	8.55×10^{-26}
KEGG_MELANOMA [71]	Melanoma	10		4.39×10^{-28}	4.19×10^{-25}
REACTOME_SIGNALING_BY_FGFR [112]	Genes involved in Signaling by FGFR	10		5.37×10^{-26}	4.28×10^{-23}
KEGG_NON_SMALL_CELL_LUNG_CANCER [54]	Non-small cell lung cancer	9		1.16×10^{-25}	7.91×10^{-23}
REACTOME_SIGNALING_BY_FGFR_IN_DISEASE [127]	Genes involved in Signaling by FGFR in disease	10		1.98×10^{-25}	1.18×10^{-22}
REACTOME_NGF_SIGNALLING_VIA_TRKA_FROM_OM_THE_PLASMA_MEMBRANE [137]	Genes involved in NGF signalling via TRKA from the plasma membrane	10		4.35×10^{-25}	2.31×10^{-22}
KEGG_PATHWAYS_IN_CANCER [328]	Pathways in cancer	11		2.07×10^{-24}	9.19×10^{-22}

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Gene/geneset overlap matrix

Entrez Gene Id	Gene Symbol	KEGG_GLIOMA	KEGG_PROSTATE_CANCER	KEGG_ENDOMETRIAL_CANCER	KEGG_ACUTE_MYELOID_LEUKEMIA	KEGG_MELANOMA	REACTOME_SIGNALING_BY_FGFR	KEGG_NON_SMALL_CELL_LUNG_CANCER	REACTOME_SIGNALING_BY_FGFR_IN_DISEASE	REACTOME_NGF_SIGNALING_VIA_TRKA_FROM_THE_PLASMA_MEMBRANE	KEGG_PATHWAYS_IN_CANCER	Entrez	Source	Gene Description
5594	MAPK1											5594	S	mitogen-activated protein kinase 1
5604	MAP2K1											5604	S	mitogen-activated protein kinase kinase 1
3845	KRAS											3845	S	v-Ki-ras2 Kirsten rat sarcoma viral oncogene homolog
4893	NRAS											4893	S	neuroblastoma RAS viral (v-ras) oncogene homolog
5605	MAP2K2											5605	S	mitogen-activated protein kinase kinase 2
207	AKT1											207	S	v-akt murine thymoma viral oncogene homolog 1
5290	PIK3CA											5290	S	phosphoinositide-3-kinase, catalytic, alpha polypeptide
673	BRAF											673	S	v-raf murine sarcoma viral oncogene homolog B1
595	CCND1											595	S	cyclin D1
5728	PTEN											5728	S	phosphatase and tensin homolog
2475	MTOR											2475	S	mechanistic target of rapamycin (serine/threonine kinase)

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

Gene Identifiers

KRAS
BRAF
NRAS
MAP2K1
MAP2K2
MAPK1
CCND1
PIK3CA
AKT1
PTEN
MTOR

Compute Overlaps

- ☐ **H**: hallmark gene sets
- ☐ **C1**: positional gene sets
- ☐ **C2**: curated gene sets
 - ☐ **CGP**: chemical and genetic perturbations
 - ☐ **CP**: Canonical pathways
 - ☐ **CP:BIOCARTA**: BioCarta gene sets
 - ☐ **CP:KEGG**: KEGG gene sets
 - ☐ **CP:REACTOME**: Reactome gene sets
- ☐ **C3**: motif gene sets
 - ☐ **MIR**: microRNA targets
 - ☐ **TFT**: transcription factor targets
- ☐ **C4**: computational gene sets
 - ☐ **CGN**: cancer gene neighborhoods
 - ☐ **CM**: cancer modules
- ☐ **C5**: GO gene sets
 - ☐ **BP**: GO biological process
 - ☐ **CC**: GO cellular component
 - ☐ **MF**: GO molecular function
- ☐ **C6**: oncogenic signatures
- ☐ **C7**: immunologic signatures

show top 10 genesets

with FDR q-value below

compute overlaps

Compendia expression profiles

- ☒ Human tissue compendium (Novartis)
- ☐ Global Cancer Map (Broad Institute)
- ☐ NCI-60 cell lines (National Cancer Institute)

display expression profile

Gene families

show gene families

<http://software.broadinstitute.org/gsea/msigdb/>

- ### Display Expression for Selected Genes

Figure 1. Heatmap of the expression of 100 genes in 100 cell lines. The genes are listed on the left, and the cell lines are listed on the right. The color scale ranges from -2 (blue) to 2 (red). The heatmap shows that the expression of these genes is highly variable across the cell lines, with some genes being highly expressed in many cell lines and others being highly expressed in only a few. The cell lines are grouped into three main clusters: 1) Cell lines derived from human cancer (e.g., A549, H1299, H1975, H2009, H2173, H226, H23, H246, H292, H460, H520, H596, H660, H661, H727, H753, H773, H827, H837, H858, H925, H959, H973, H975, H976, H977, H978, H979, H980, H981, H982, H983, H984, H985, H986, H987, H988, H989, H990, H991, H992, H993, H994, H995, H996, H997, H998, H999, H1000), 2) Cell lines derived from non-human sources (e.g., A549, H1299, H1975, H2009, H2173, H226, H23, H246, H292, H460, H520, H596, H660, H661, H727, H753, H773, H827, H837, H858, H925, H959, H973, H975, H976, H977, H978, H979, H980, H981, H982, H983, H984, H985, H986, H987, H988, H989, H990, H991, H992, H993, H994, H995, H996, H997, H998, H999, H1000), and 3) Cell lines derived from normal human cells (e.g., A549, H1299, H1975, H2009, H2173, H226, H23, H246, H292, H460, H520, H596, H660, H661, H727, H753, H773, H827, H837, H858, H925, H959, H973, H975, H976, H977, H978, H979, H980, H981, H982, H983, H984, H985, H986, H987, H988, H989, H990, H991, H992, H993, H994, H995, H996, H997, H998, H999, H1000).

kinases oncogenes translocated genes tumor suppressors

Molecular Signatures Database (MSigDB)

<http://software.broadinstitute.org/gsea/msigdb/>

View Gene Families for Selected Genes

The following table provides a functional overview of the MSigDB gene sets by categorizing their genes into a small number of carefully chosen "gene families". To categorize the genes in a gene set, use the gene set page or the Investigate Gene Sets page.

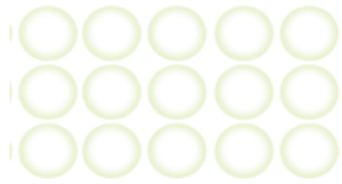
Click on a gene family or gene family intersection to retrieve annotations for those genes.

	cytokines and growth factors	transcription factors	homeodomain proteins	cell differentiation markers	protein kinases	translocated cancer genes	oncogenes	tumor suppressors
tumor suppressors	0	0	0	0	0	0	0	1
oncogenes	0	0	0	0	2	2	6	
translocated cancer genes	0	0	0	0	1	2		
protein kinases	0	0	0	0	6			
cell differentiation markers	0	0	0	0				
homeodomain proteins	0	0	0					
transcription factors	0	0						
cytokines and growth factors	0							

Members of these "gene families" share a common feature such as homology or biochemical activity. They do not necessarily have common origins. For the source of each "gene family" definition, [click here](#).


Curated Gene Signatures (GeneSigDB)

<http://genesigdb.org/genesigdb/>



GeneSigDB

Curated Gene Signatures



[Home](#) [Browse](#) [Analyze My Genes](#) [Download](#) [Support](#) [Contact Us](#)

Publication Search ?

Search the full text of articles to retrieve a list of publications and the gene signatures they describe. Enter one or more search terms, such as author name, article title, journal name, or keywords.

[Search Publications](#) (e.g.: basal breast cancer)

OR

Gene Search ?

Search gene annotations to retrieve genes listed in GeneSigDB gene signatures.

[Search Genes](#) (e.g.: BRCA*, BRCA1)

[Citations: >150]

The **Gene Signature DataBase** is a searchable database of fully traceable, standardized, annotated gene signatures which have been manually curated from publications that are indexed in [PubMed](#). Enter a search term above to get started.

News

September, 2011: GeneSigDB Data and Website Update

We continue to expand. So far we have read and processed almost 3,000 publications to extract 3,515 genes signatures from 1,604 publications. See [GeneSigDB Release 4 release notes](#)

We have a new tag cloud [Browse](#) feature to enable easy browsing of GeneSigDB.

Additional [download](#) formats. Download GeneSigDB as an R/Bioconductor data file, gmt or compressed flat file formats.

GeneSigDB Data Release 4

Gene Signatures: 3515
Published Articles: 1604
Genes (Human): 20,523
Tissues and Diseases: More than 50
Species: 3

Systems biology

DSigDB: drug signatures database for gene set analysis

Minjae Yoo^{1,†}, Jimin Shin^{1,†}, Jihye Kim¹, Karen A. Ryall¹, Kyubum Lee², Sunwon Lee², Minji Jeon², Jaewoo Kang² and Aik Choon Tan^{1,2,*}

¹Department of Medicine, Translational Bioinformatics and Cancer Systems Biology Laboratory, Division of Medical Oncology, University of Colorado Anschutz Medical Campus, Aurora, CO 80045, USA and ²Department of Computer Science and Engineering, Korea University, Seoul 136-713, South Korea

*To whom correspondence should be addressed.

[†]The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

Associate Editor: Jonathan Wren

Received on February 17, 2015; revised on April 30, 2015; accepted on May 13, 2015

Abstract

Summary: We report the creation of Drug Signatures Database (DSigDB), a new gene set resource that relates drugs/compounds and their target genes, for gene set enrichment analysis (GSEA). DSigDB currently holds 22 527 gene sets, consists of 17 389 unique compounds covering 19 531 genes. We also developed an online DSigDB resource that allows users to search, view and download drugs/compounds and gene sets. DSigDB gene sets provide seamless integration to GSEA software for linking gene expressions with drugs/compounds for drug repurposing and translational research.

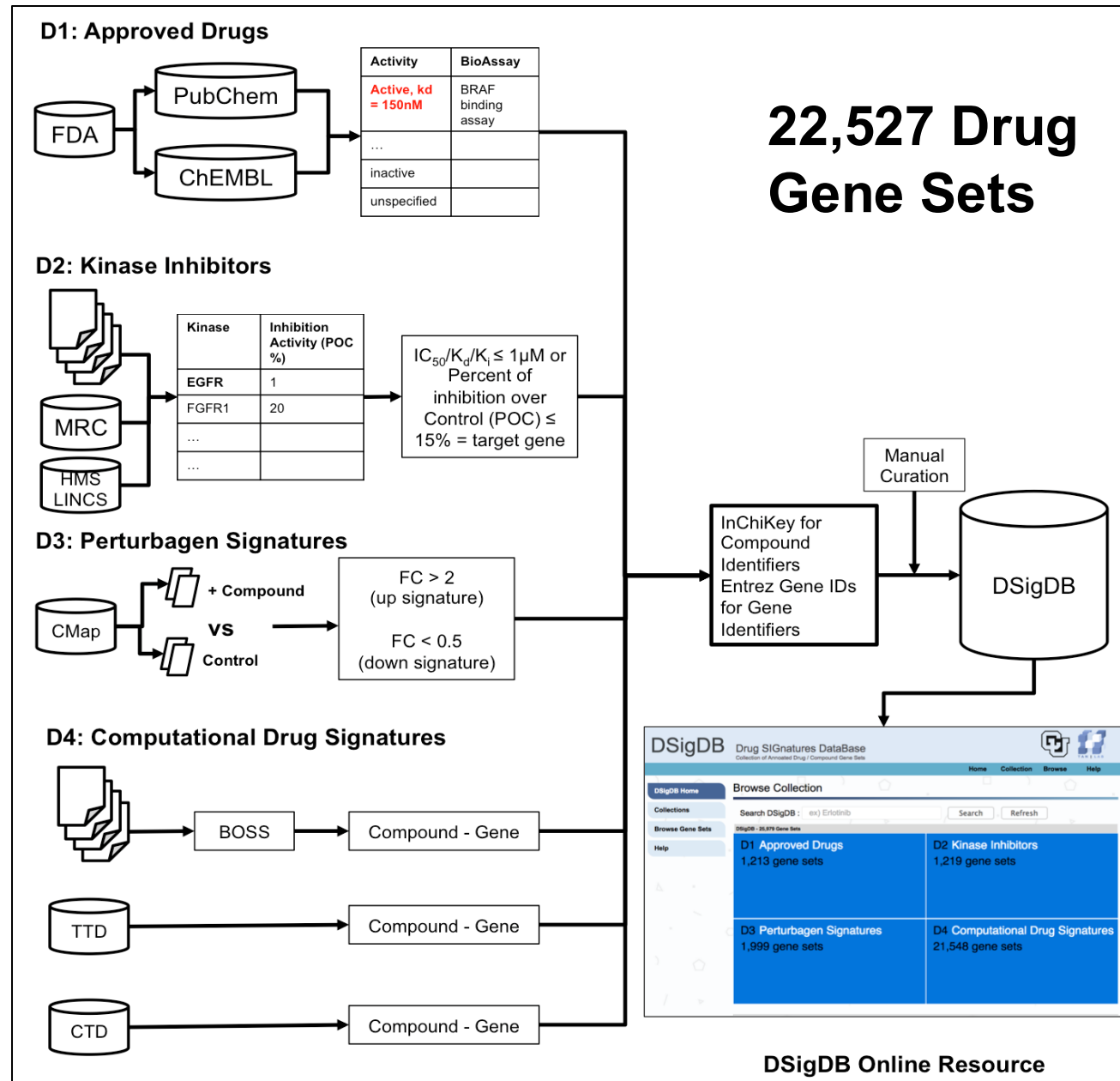
Availability and implementation: DSigDB is freely available for non-commercial use at <http://tanlab.ucdenver.edu/DSigDB>.

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

Contact: aikchoon.tan@ucdenver.edu

[Citations: >570]

DSigDB Workflow



DSigDB

The screenshot displays the DSigDB (Drug SIGNatures DataBase) website. The header includes the site name, a description, and navigation links. A sidebar on the left provides quick access to various sections. The main content area, titled 'Browse Collection', features a search bar with the example text 'ex) Erlotinib' and buttons for 'Search' and 'Refresh'. Below the search bar, a table lists four categories of gene sets, each with a count.

DSigDB - 22,527 Gene Sets	
D1 Approved Drugs 1,202 gene sets	D2 Kinase Inhibitors 1,220 gene sets
D3 Perturbagen Signatures 1,998 gene sets	D4 Computational Drug Signatures 18,107 gene sets

<http://tanlab.ucdenver.edu/DSigDB/>

DSigDB

The screenshot shows the DSigDB (Drug SIGNatures DataBase) website. The header includes the site name, a navigation bar with links to Home, Search Gene, Collection, Browse, Download, and Help, and logos for the University of Texas at Austin and the Tan Lab. A left sidebar contains links to DSigDB Home, Search Gene, Collections, Browse Gene Sets, Download, and Help. The main content area, titled 'Browse Collection', features a search bar with the text 'ex) Erlotinib' and buttons for 'Search' and 'Refresh'. Below the search bar, a table displays four major collections: D1 Approved Drugs (1,202 gene sets), D2 Kinase Inhibitors (1,220 gene sets), D3 Perturbagen Signatures (1,998 gene sets), and D4 Computational Drug Signatures (18,107 gene sets). A search results section at the bottom provides instructions on how to use the drug names in the table.

DSigDB Drug SIGNatures DataBase
Collection of Annotated Drug / Compound Gene Sets

Home Search Gene Collection Browse Download Help

DSigDB Home

Search Gene

Collections

Browse Gene Sets

Download

Help

Browse Collection

Search DSigDB : Search Refresh

DSigDB - 22,527 Gene Sets

D1 Approved Drugs 1,202 gene sets	D2 Kinase Inhibitors 1,220 gene sets
D3 Perturbagen Signatures 1,998 gene sets	D4 Computational Drug Signatures 18,107 gene sets

Search Result

Drug Name - Click on a drug name to view its gene set page.

Menu (points to top navigation bar)

Search Box (points to search input field)

Menu (points to left sidebar)

Zoomable table to view the different sub-collections with the four major collections. (points to the collection table)

Results Display (points to search result instructions)

Drug Search

Browse Collection

Search DSigDB :

DSigDB - 22,527 Gene Sets

D1 Approved Drugs
1,202 gene sets

D2 Kinase Inhibitors
1,220 gene sets

D3 Perturbagen Signatures
1,998 gene sets

D4 Computational Drug Signatures
18,107 gene sets

Search Result

Drug Name - Click on a drug name to view its gene set page.

Collection	Source	Representative Name	Synonym
D1	D1	Erlotinib Hydrochloride	Erlotinib Hydrochloride
D2	FDA	Erlotinib	Erlotinib
	Kinome Scan	Erlotinib	Erlotinib
	RBC	Erlotinib	Erlotinib
D4	BOSS	Erlotinib	Erlotinib
	CTD	Erlotinib	Erlotinib
	TTD	Erlotinib	Erlotinib
Unique Gent Set for "Erlotinib"		gmt	text

Gene Search

- DSigDB Home
- Search Gene**
- Collections
- Browse Gene Sets
- Download
- Help

Search Gene

Search Gene (19,531) :

Show entries

Gene

Source

Chemical Name

Search your gene name

Gene

Source

Chemical Name

Page 1 of 1 (Total 1 Data Sets)

Previous

1

Next

Gene Search Result

Search Gene		
Search Gene (19,531) : <input type="text" value="EGFR"/>		
<input type="button" value="Search"/>		
Show <input type="text" value="15"/> entries		
Gene	Source	Chemical Name
EGFR	D1	chlorpromazine
EGFR	D1	afatinib
EGFR	D1	thioridazine
EGFR	D1	vandetanib
EGFR	D1	baciguent
EGFR	D1	levodopa
EGFR	D1	hexachlorophene
EGFR	D1	zafirlukast
EGFR	D1	erlotinib hydrochloride
EGFR	D1	miconazole
EGFR	D1	tamoxifen
EGFR	D1	crystal violet
EGFR	D1	methyldopa
EGFR	D1	dobutamine
EGFR	D1	crizotinib
Gene	Source	Chemical Name
Page 1 of 42 (Total 616 Data Sets)		
Previous <input type="button" value="1"/> <input type="button" value="2"/> <input type="button" value="3"/> <input type="button" value="4"/> <input type="button" value="5"/> ... <input type="button" value="42"/> Next		

Browsing Collection

(A)

DSigDB - 22,527 Gene Sets / D2		
FDA 28 gene sets	HMS LINCS 90 gene sets	
MRC 157 gene sets	Roche 570 gene sets	
GSK 204 gene sets	Kinome Scan 72 gene sets	RBC 99 gene sets

(B)

DSigDB - 22,527 Gene Sets / D2		
FDA 28 gene sets	HMS LINCS 90 gene sets	
MRC 157 gene sets	Roche 570 gene sets	
GSK 204 gene sets	Kinome Scan 72 gene sets	RBC 99 gene sets

Browsing Collection

DSigDB - 22,527 Gene Sets / D2

FDA 28 gene sets	HMS LINCS 90 gene sets		
MRC 157 gene sets	Roche 570 gene sets		
	Kinome Scan 72 gene sets	RBC 99 gene sets	
GSK 204 gene sets			

Search Result : FDA

Drug Name - Click on a drug name to view its gene set page.

Afatinib	Axitinib	Bosutinib	Cabozantinib
Ceritinib	Crizotinib	Dabrafenib	Dasatinib
Erlotinib	Gefitinib	Ibrutinib	Imatinib
Lapatinib	Lenvatinib	Nilotinib	Nintedanib
Palbociclib	Pazopanib	Ponatinib	Regorafenib
Ruzolitinib	Sirolimus	Sorafenib	Sunitinib
Tofacitinib	Trametinib	Vandetanib	Vemurafenib

Compound Webpage

Gene Set: D2 : FDA - Gefitinib

CollectionD2 : FDA

Chemical NameGefitinib

FDA	NPC	WHO	Indian	Australia	China	Traditional Herbal	Clinical Trail
Approved	Not	Not	Approved	Not	Not	Not	Not

Molecular Weight446.902 g/mol

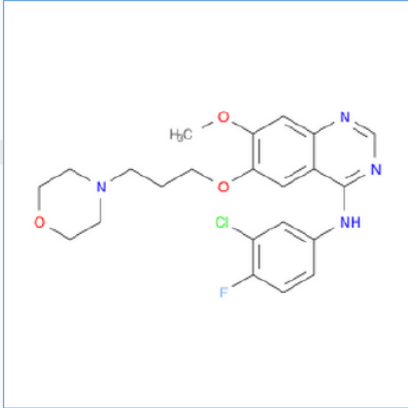
Hydrogen Bond Donor Count1

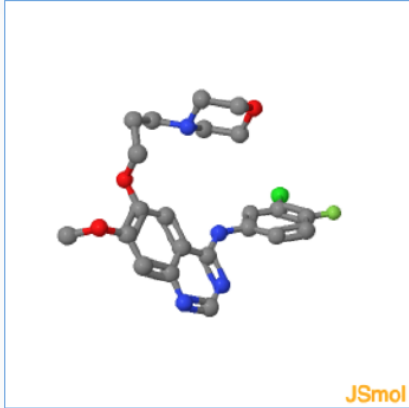
Hydrogen Bond Acceptor Count6

cLogP4.2865

Lipinski RuleTrue

Structure





InChI

InChIKey

Links

CAS Num : 184475-35-2

Gene (40 / 41)
⊞ More

Value Type	Value↑	Concentration	Gene	PMID / Source
Kd	0.520	nM	EGFR(del_L747-T751,Sins)	22037378
Kd	0.540	nM	EGFR(del_E746-A750)	22037378
Kd	0.570	nM	EGFR(del_L747-E749,A750P)	22037378
Kd	0.570	nM	EGFR(del_L747-S752,P753S)	22037378

Download gene sets

[gmt](#), [text](#), [Detailed text](#)

Compound Webpage

Gene Set: D2 : FDA - Gefitinib

Drug Information

Clinical Development Stage of the Drug

FDA	NPC	WHO	Indian	Australia	China	Traditional Herbal	Clinical Trail
Approved	Not	Not	Approved	Not	Not	Not	Not

Molecular Information of the Compound

Molecular Weight: 446.902 g/mol

Hydrogen Bond Donor Count: 1

Hydrogen Bond Acceptor Count: 6

cLogP: 4.2865

Lipinski Rule: True



2D and 3D Chemical Structure

Compound unique identifiers

InChI: 1S/C22H24ClFN4O3/c1-29-20-13-19-16(12-21(20)31-8-2-5-28-6-9-30-10-7-28)22(26-14-25-19)27-15-3-4-18(24)17(23)11-15/h3-4,11-14H,2,5-10H2,1H3,(H,25,26,27)

InChIKey: XGALLCVXEZPNRQ-UHFFFAOYSA-N

External Links

Links:  

CAS Num : 184475-35-2

Gene membership and quantitative inhibition data

Value Type	Value†	Concentration	Gene	PMID / Source
Kd	0.520	nM	EGFR(del_L747-T751,Sins)	22037378
Kd	0.540	nM	EGFR(del_E746-A750)	22037378
Kd	0.570	nM	EGFR(del_L747-E749,A750P)	22037378
Kd	0.570	nM	EGFR(del_L747-S752,P753S)	22037378

Download gene sets: [gmt](#), [text](#), [Detailed text](#)

Download Gene Set

Link to PubMed/Source

Collections

Collection	Description	Unique Number of Genes	Number of Gene Sets	Download
DSigDB	All Gene Sets.	19,531	22,527	GMT File
D1 : FDA Approved (browse 1,202 gene sets)	FDA Approved Drug Gene Sets.	1,288	1,202	GMT File
D2 : Kinase Inhibitors	Kinase Inhibitors Gene Sets based on in vitro kinase profiling assays.	407	1,220	GMT File
FDA (browse 28 gene sets)	FDA Approved Kinase Inhibitors.	341	28	GMT File
HMS LINC (browse 90 gene sets)	Kinase inhibition assays extracted from HMS LINC database.	381	90	GMT File
MRC (browse 157 gene sets)	Kinase inhibition assays extracted from MRC Kinome Inhibition database.	137	157	GMT File
GSK (browse 204 gene sets)	GSK Published Kinase Inhibitor Set (PKIS), kinase inhibitors used as chemical probes.	116	204	GMT File
Roche (browse 570 gene sets)	Kinase Inhibitors profiled by Roche.	153	570	GMT File
RBC (browse 99 gene sets)	Kinase Inhibitors profiled by Reaction Biology Corporation.	246	99	GMT File
KinomeScan (browse 72 gene sets)	Kinase Inhibitors profiled by DiscoveryRx using KinomeScan technology.	374	72	GMT File
D3 : Perturbagen Signatures (browse 1,998 gene sets)	7,064 gene expression profiles from three cancer cell lines perturbed by 1,309 compounds from CMap (build 02).	11,137	1,998	GMT File
CMap (browse 1,998 gene sets)	7,064 gene expression profiles from three cancer cell lines perturbed by 1,309 compounds from CMap (build 02).	11,137	1,998	GMT File
D4 : Computational Drug Signatures	Drug signatures extracted from literatures using a mixture of manual curation and by automatic computational approaches.	18,854	18,107	GMT File
BOSS (browse 2,114 gene sets)	Text mining approach of drug-gene targets using Biomedical Object Search System (BOSS).	3,354	2,114	GMT File
CTD (browse 5,163 gene sets)	Curation of targets from Comparative Toxicogenomics Database (CTD).	18,700	5,163	GMT File
TTD (browse 10,830 gene sets)	Manual curation of targets from the Therapeutics Targets Database (TTD).	1,389	10,830	GMT File

Take home message

- Genes don't act alone to drive biological processes
- Gene set analysis such as GSEA can identify set of coordinately and subtle expressed genes participated in a functional group compared to Candidate Gene Analysis
- *Biology trumps statistics* – if you can validate the gene sets

Download Collections

Download

DSigDB provides several options for downloading the data.

Current Release

The current data release of DSigDB is Release 1 (released May 2015).

DSigDB Release 1 (Updated)

- DSigDBv1.0.gmt
- DSigDBv1.0.txt
- DSigDBv1.0 Detailed.txt

Compound : Gefitinib

EPHA6
STK10
MKNK1
EGFR
RIPK2
MAP2K5
HIPK4
ABL1
FLT3
CSNK1E
GAK
LYN
IRAK1

Drug	Gene	Type	Source	
Gefitinib	EGFR	Kd=40.0 (nM)	FDA	
Gefitinib	EGFR	Kd=0.54 (nM)	FDA	
Gefitinib	EGFR	Kd=0.98 (nM)	FDA	
Gefitinib	ABL1	Kd=460.0 (nM)	FDA	
Gefitinib	CDK7	Kd=610.0 (nM)	FDA	
Gefitinib	EGFR	Kd=140.0 (nM)	FDA	
Gefitinib	ABL1	Kd=680.0 (nM)	FDA	
Gefitinib	ABL1	Kd=360.0 (nM)	FDA	
Gefitinib	LCK	Kd=630.0 (nM)	FDA	
Gefitinib	ABL1	Kd=480.0 (nM)	FDA	
Gefitinib	MKNK1	Kd=290.0 (nM)	FDA	
Gefitinib	SBK1	Kd=560.0 (nM)	FDA	
Gefitinib	SLK	Kd=920.0 (nM)	FDA	
Gefitinib	EGFR	Kd=1.1 (nM)	FDA	
Gefitinib	ABL1	Kd=230.0 (nM)	FDA	
Gefitinib	IRAK4	Kd=540.0 (nM)	FDA	
Gefitinib	ERBB3	Kd=790.0 (nM)	FDA	
Gefitinib	GAK	Kd=13.0 (nM)	FDA	
Gefitinib	ABL1	Kd=780.0 (nM)	FDA	
Gefitinib	LYN	Kd=990.0 (nM)	FDA	
Gefitinib	IRAK1	Kd=69.0 (nM)	FDA	
Gefitinib	CHEK2	Kd=800.0 (nM)	FDA	

Use Case Example: EGFRwt NSCLC

Cell line	Histology	EGFR	KRAS	Gefitinib IC50 ($\mu\text{mol/L}$)
Sensitive				
H358	BAC	Wild-type	Mutant	0.18
H322	BAC	Wild-type	Wild-type	0.25
Calu-3	Adenocarcinoma	Wild-type	Wild-type	0.3
H1334	Large	Wild-type	Wild-type	0.3
H1648	Adenocarcinoma	Wild-type	Wild-type	0.38
HCC78	Adenocarcinoma	Wild-type	Wild-type	0.4
H2126	Large	Wild-type	Wild-type	1
HCC193	Adenocarcinoma	Wild-type	Wild-type	1.5
HCC95	Adenocarcinoma	Wild-type	Wild-type	1.9
Resistant				
H125	Adenosquamous	Wild-type	Wild-type	4.8
HCC44	Adenocarcinoma	Wild-type	Mutant	7.9
H1703	Squamous	Wild-type	Wild-type	8
HCC15	Squamous	Wild-type	Wild-type	9.4
A549	Adenocarcinoma	Wild-type	Wild-type	9.6
H157	Squamous	Wild-type	Mutant	12.8
H460	Large	Wild-type	Mutant	12.9
H520	Squamous	Wild-type	Wild-type	13.6
H1299	Large	Wild-type	Wild-type	14.7

(Adapted from Coldren et al MCR 2006)

Enriched Gene Sets

Enriched in Sensitive Group (p < 0.05)

GENE SET NAME	GENE SET SIZE	Normalized Enrichment Score	Nominal p-val	Intended targets	Inhibiting EGFR/ERBB2/ERBB3 based on Kinase Inhibition Assays		
					EGFR	ERBB2	ERBB3
CI-1033_KINOME SCAN	28	1.78	0.0000	EGFR/ERBB2	Yes	Yes	Yes
AZD-9291_LINCS	43	1.63	0.0125	EGFR	Yes	Yes	No
ZM-447439_LINCS	41	1.55	0.0285	AURKA	Yes	Yes	Yes
AZD-2171_KINOME SCAN	42	1.52	0.0271	VEGFR2/PDGFRα/PDGFRβ	Yes	No	Yes
SB-203580_KINOME SCAN	18	1.52	0.0496	p38-alpha	Yes	No	No
WH-4-023_LINCS	124	1.48	0.0101	LCK	Yes	Yes	Yes
PP-242_KINOME SCAN	111	1.48	0.0116	MTOR/PIK3CA	Yes	Yes	No
CABOZANTINIB_FDA	45	1.47	0.0313	VEGFR2,MET	No	No	No
VANDETANIB_FDA	51	1.47	0.0355	RET/VEGFR2/EGFR	Yes	No	Yes
HG-9-91-01_LINCS	137	1.46	0.0179	SIK1	Yes	Yes	Yes
AZ-628_LINCS	51	1.46	0.0265	BRAF	Yes	No	No
VANDETANIB_KINOME SCAN	51	1.45	0.0448	RET/VEGFR2/EGFR	Yes	No	Yes
EXEL-2880/GSK-1363089_KINOME SCAN	131	1.45	0.0147	MET/AXL/VEGFR2	Yes	No	Yes
PD-173955_KINOME SCAN	105	1.40	0.0305	ABL1/SRC	Yes	No	Yes
BOSUTINIB_LINCS	69	1.40	0.0490	ABL1/SRC	Yes	Yes	Yes
R406_KINOME SCAN	183	1.39	0.0123	SYK,FLT3	Yes	No	No

Enriched in Resistant Group (p < 0.05)

GENE SET NAME	GENE SET SIZE	Normalized Enrichment Score	Nominal p-val	Intended targets	Inhibiting EGFR/ERBB2/ERBB3 based on Kinase Inhibition Assays		
					EGFR	ERBB2	ERBB3
KINOME_858_ROCHE	17	-1.81	0.0041	NA	No	No	No
KINOME_1901_ROCHE	17	-1.67	0.0149	NA	No	No	No
CHEMBL2062936_ROCHE	16	-1.66	0.0042	NA	No	No	No
KINOME_1242_ROCHE	16	-1.62	0.0198	NA	No	No	No
KINOME_1221_ROCHE	19	-1.62	0.0194	NA	No	No	No
KINOME_866_ROCHE	15	-1.56	0.0395	NA	No	No	No
RO-3306_MRC	29	-1.45	0.0455	CDK1	No	No	No

RO-3306 Sensitivity (From GDSC)

