

# KSBi-BIML 2026

Bioinformatics & Machine Learning(BIML)  
Workshop for Life Scientists

생명정보학 & 머신러닝 워크샵 (온라인)



## Introduction to ConnectivityMap

전민지 \_ 고려대학교



**KSBI**  
KOREAN SOCIETY FOR  
BIOINFORMATICS

| 한국생명정보학회



본 강의 자료는 한국생명정보학회가 주관하는 BIML 2026 워크샵을 목적으로 제작된 것으로 해당 목적 이외의 다른 용도로 사용할 수 없음을 분명하게 알립니다.

이를 다른 사람과 공유하거나 복제, 배포, 전송할 수 없으며 만약 이러한 사항을 위반할 경우 발생하는 **모든 법적 책임은 행위자 본인에게 있음**을 알립니다.

# KSBI-BIML 2026

## Bioinformatics & Machine Learning (BIML) Workshop for Life Scientists

한국생명정보학회가 주최하는 BIML-2026 동계 Bioinformatics & Machine Learning 교육 워크숍에 여러분을 초대합니다.

BIML 워크숍은 생명정보학 연구자들이 최신 AI바이오 분야의 인공지능 기반 분석 기술과 바이오 데이터 분석 기법을 이론과 실습을 통해 체계적으로 배울 수 있는 전문 교육 프로그램입니다. 2015년에 시작된 BIML 워크숍은 올해로 12년 차를 맞이하며, 국내 생명정보학 분야의 최초이자 최고 수준의 교육 프로그램으로 자리 잡았습니다. 이번 워크숍은 크게 인공지능바이오(AI바이오) 분야와 디지털바이오 분야, 두 분야로 구성됩니다.

AI바이오 분야에서는 생명정보 분석에 폭넓게 응용되고 있는 다양한 인공지능 기반 자료 모델링 기법을 다룰 예정입니다. 특히, 인공지능 심층학습을 활용한 단백질 구조 예측, 유전체 분석, 신약 개발에 대한 이론 및 실습 강의를 진행됩니다.

또한 디지털바이오 분야에서는 단일세포오믹스, 공간오믹스, 멀티오믹스, 메타오믹스에 대한 강의도 마련되어 있어, 연구자들의 분석 역량 강화에 실질적인 도움을 줄 것으로 기대됩니다.

또한 2024년부터 추가된 의료정보 자료 분석을 다루는 강의를 올해도 지속해서 운영하고자 합니다. 이는 최근 의료정보 자료 분석에 관한 연구 수요 증가를 반영한 것으로, 관련 연구를 수행하는 의과학자 및 의료정보 연구자들에게 유용한 지침을 제공할 것입니다.

또한, 올해도 생명정보학 기술의 다양화에 발맞춰 온라인 강좌를 대폭 확대했습니다. 올해는 무료 강좌 10개를 포함한 총 40개 이상의 강좌가 개설되며, 연구 주제에 맞는 강좌 추천과 강연료 할인 혜택도 제공합니다.

BIML-2026는 국내 주요 연구 중심 대학의 전임 교수 및 각 분야 최고 전문가들의 강의로 구성되어 있으며, 기초 이론부터 최신 연구 동향까지 아우르는 심도 있는 교육의 장이 될 것으로 확신합니다.

여러분의 많은 관심과 참여를 기대합니다!

2026년 2월

한국생명정보학회장 류 성 호

## Introduction to ConnectivityMap

ConnectivityMap 플랫폼은 3만 가지의 약물에 대한 300만 개 유전자 발현 프로파일을 포함하고 있어, 연구자들이 약물과 질병, 유전자 간의 복잡한 상호작용을 이해할 수 있게 도와주고 있다. 연구자들은 ConnectivityMap을 사용하여 기존 약물이 새로운 질병에 대해 어떤 효과를 보일 수 있는지 예측할 수 있으며, 이는 약물 개발 과정을 가속화하고 비용을 절감하는 데 큰 도움이 되고 있다.

ConnectivityMap 튜토리얼 강의를 통해 이러한 대규모 데이터셋을 탐색하고 분석하는 방법을 배우게 되며, 실제 사례 연구를 통해 ConnectivityMap이 어떻게 실제 연구에 적용될 수 있는지 배운다. 또한 이론적 지식뿐만 아니라 실제적인 적용 능력을 함양하는 것을 목표로 한다.

강의는 다음의 내용을 포함한다:

- ConnectivityMap 데이터의 이해
- ConnectivityMap 데이터의 활용
- ConnectivityMap 데이터의 응용

\* 교육생준비물: 노트북 (메모리 8GB 이상, 디스크 여유공간 30GB 이상)

\* 강의 난이도: 초급

\* 강의: 전민지 교수 (고려대학교 의과대학)

# Curriculum Vitae

Speaker Name: **Minji Jeon, Ph.D.**



## ► Personal Info

Name Minji Jeon  
Title Assistant Professor  
Affiliation Korea University

## ► Contact Information

Address 161, Jeongneung-ro, Seongbuk-gu, Seoul, 02708  
Email mjjeon@korea.ac.kr

---

## Research Interest

AI-driven drug discovery, machine learning, bioinformatics

## Educational Experience

2012 B.S. in Computer Science, Korea University, Korea  
2014 M.S. in Interdisciplinary Graduate Program in Bioinformatics, Korea University, Korea  
2018 Ph.D. in Computer Science, Korea University, Korea

## Professional Experience

2018-2019 Research Professor, Korea University, Korea  
2020-2022 Postdoctoral Fellow, Icahn School of Medicine at Mount Sinai, USA  
2022- Assistant Professor, Korea University, Korea

## Selected Publications (5 maximum)

1. Zhaoping Xiong<sup>†</sup>, **Minji Jeon**<sup>†</sup>, Robert J Allaway<sup>†</sup>, Jaewoo Kang, Donghyeon Park, Jinhyuk Lee, Hwisang Jeon, Miyoung Ko, Hualiang Jiang, Minyue Zheng, Aik Choon Tan, Xindi Guo, The Multi-Targeting Drug DREAM Challenge Community, Kristen K Dang, Alex Tropsha, Chana Hecht, Tirtha K. Das, Heather A. Carlson, Ruben Abagyan, Justin Guinney, Avner Schlessinger\*, Ross Cagan\* "Crowdsourced identification of multi-target kinase inhibitors for RET- and TAU-based disease: the Multi-Targeting Drug DREAM Challenge" *PLoS computational biology* 17.9 (2021): e1009302.
2. **Minji Jeon**<sup>†</sup>, Kathleen M. Jagodnik<sup>†</sup>, Eryk Kropiwnicki, Daniel J. Stein, Avi Ma'ayan\* "Prioritizing Pain-Associated Targets with Machine Learning" *Biochemistry* 60.18 (2021): 1430-1446.
3. **Minji Jeon**, Donghyeon Park, Jinhyuk Lee, Hwisang Jeon, Miyoung Ko, Sunkyu Kim, Yonghwa Choi, Aik-Choon Tan, Jaewoo Kang\* "ReSimNet: Drug Response Similarity Prediction using Siamese Neural Networks" *Bioinformatics* 35.24 (2019): 5249-5256.
4. Michael Patrick Menden, Dennis Wang, Yuanfang Guan, Michael Mason, Bence Szalai, Krishna C Bulusu, Thomas Yu, Jaewoo Kang, **Minji Jeon**, Russ Wolfinger, Tin Nguyen, Mikhail Zaslavskiy, AstraZeneca-Sanger Drug Combination DREAM Consortium, In Sock Jang, Zara Ghazoui, Mehmet Eren Ahsen, Robert Vogel, Elias Chaibub Neto, Thea Norman, Eric KY Tang, Mathew J Garnett, Giovanni Di Veroli, Stephen Fawell, Gustavo Stolovitzky, Justin Guinney, Jonathan R Dry, Julio Saez-Rodriguez\*, "Community assessment to advance computational prediction of cancer drug combinations in a pharmacogenomic screen" *Nature Communications*, 10.1 (2019): 2674.
5. **Minji Jeon**, Sunkyu Kim, Sungjoon Park, Heewon Lee, Jaewoo Kang\* "In silico drug combination discovery for personalized cancer therapy" *BMC systems biology*, 2018, 12.2: 16.

# KSBi-BIML 2024

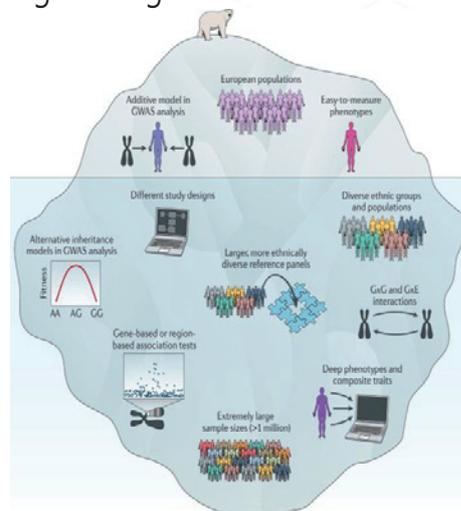
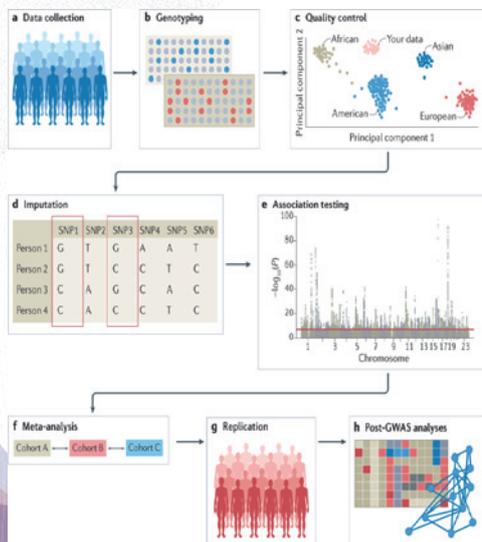
## Introduction to ConnectivityMap

고려대학교 의과대학 전민지



## Genome-based Disease Research

- GWAS (Genome Wide Association Study)
  - To identify genomic variants that are statistically associated with a risk for a disease or a particular trait
  - Limitations: association with disease is generally not sufficient to establish causality or to provide mechanistic and circuit-level biological insights



Uffelmann, Emil, et al. *Nature Reviews Methods Primers* (2021)  
Tam, Vivian, et al. *Nature Reviews Genetics* (2019)

# ConnectivityMap Concept

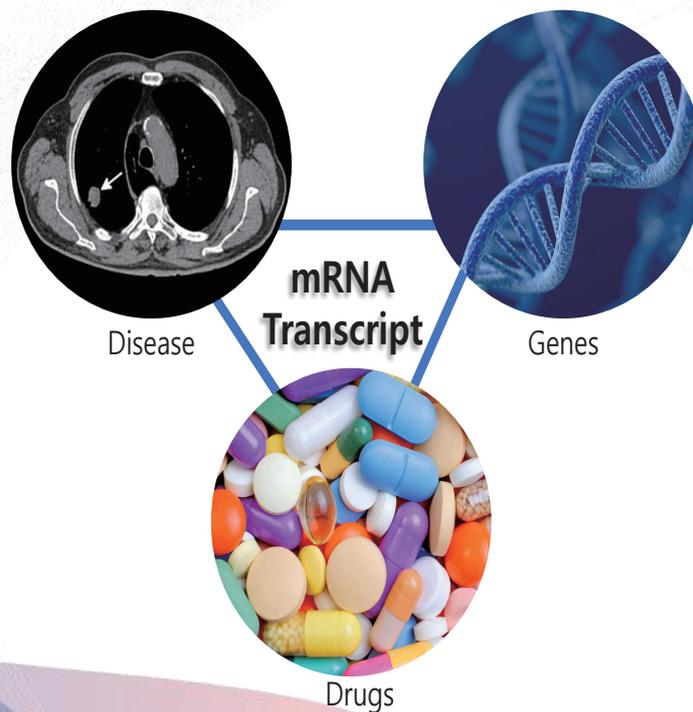
- ConnectivityMap: Linking disease, therapeutics and cell physiology

The screenshot shows the top portion of a Science journal article. The title is "The Connectivity Map: Using Gene-Expression Signatures to Connect Small Molecules, Genes, and Disease". The authors listed are Justin Lamb, Emily D. Crawford, David Peck, Joshua W. Model, Rene C. Blat, Matthew J. Wrobel, Jim Lerner, Jean Philippe Brunet, Arvind Subramanian, and Todd R. Golub. The abstract begins with: "To pursue a systematic approach to the discovery of functional connections among diseases, genetic perturbation, and drug action, we have created the first installment of a reference collection of gene-expression profiles from cultured human cells treated with bioactive small molecules, together with pattern-matching software to mine these data. We demonstrate that this 'Connectivity Map' resource can be used to find connections among small molecules sharing a mechanism of action, chemicals and physiological processes, and diseases and drugs. These results indicate the feasibility of the approach and suggest the value of a large-scale community Connectivity Map project."

Lamb, Justin, et al. *science* 2006 3

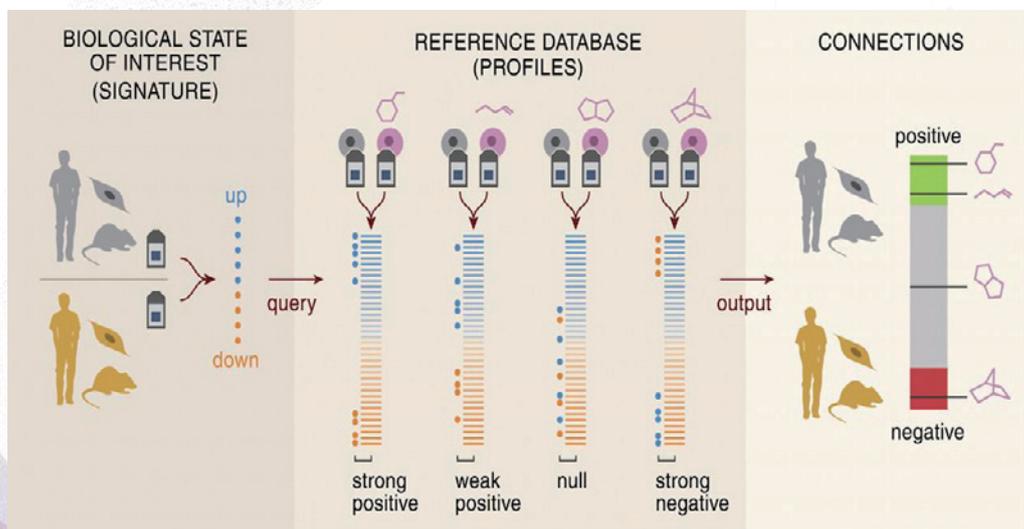
# ConnectivityMap Concept

- Linking disease, therapeutics and cell physiology



## ConnectivityMap Concept

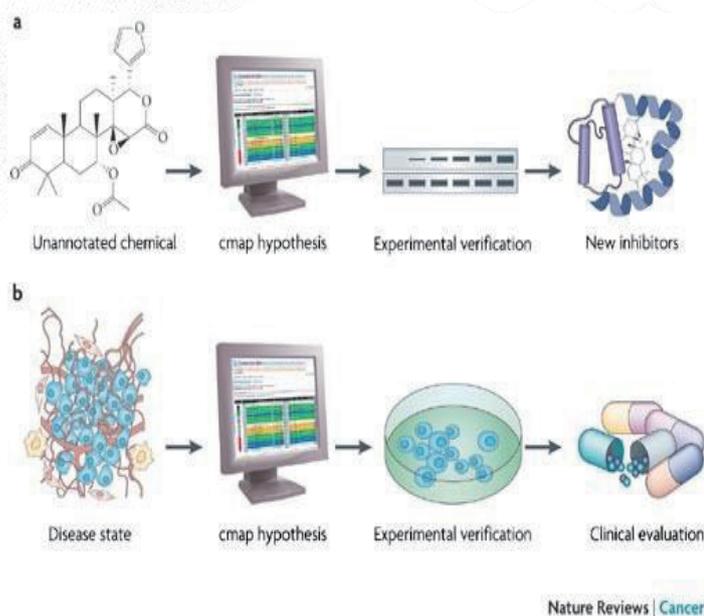
- Gene expression data could be used for the functional annotation of small molecules and genes



5

## Applications of the ConnectivityMap

- The Connectivity Map is a tool for the bench researcher



Lamb, Justin. et al. *Nature reviews cancer* (2007)

6

## Definitions

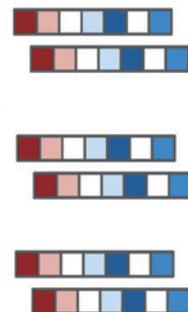
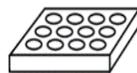
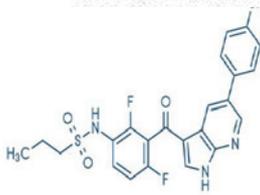
- Perturbation: an alteration of the function of a biological system, induced by external or internal mechanisms.
- Perturbagens: perturbing agents that are screened in an assay (e.g., small molecules, shRNA etc)
- Gene Signatures: differential expression of genes between two conditions, a control condition and a perturbation condition



7

## ConnectivityMap v1

- Perturb cells and measure cellular responses



**164 small molecules**  
FDA-approved drugs  
nondrug bioactive tool compounds

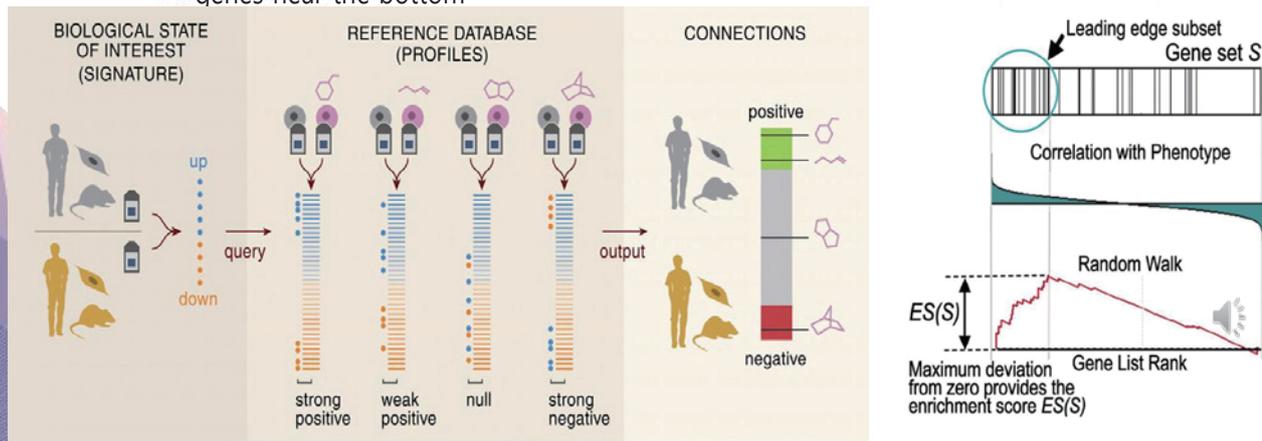
**4 cell lines**  
MCF7 (breast cancer)  
PC3 (prostate cancer)  
HL60 (leukemia)  
SKMEL5 (melanoma)

**564 microarray profiles**  
10 uM concentration  
6 or 12 hours  
with controls  
-> **453 gene signatures**

8

# ConnectivityMap v1

- Query
  - Input: gene signature (query signature)
  - Search: rank-based pattern matching strategy based on Kolmogorov-Smirnov statistics
    - The genes on the array are rank-ordered according to their differential expression relative to the control
    - The query signature is then compared to each rank-ordered list to determine whether up-regulated query genes tend to appear near the top of the list and down-regulated query genes near the bottom



9

## Results: HDAC inhibitors

- Query gene signatures
  - 13 signatures downloaded from T24 (breast carcinoma) cells treated with Trichostatin A
- Results
  - New molecules: HC toxin, valproic acid

**Fig. 2. HDAC Inhibitors.** (A) HDAC inhibitors are highly ranked with an external HDAC inhibitor signature. The "bar-view" is constructed from 453 horizontal lines, each representing an individual treatment instance, ordered by their corresponding connectivity scores with the Glaser *et al.* (24) signature (+1, top; -1, bottom). All valproic acid ( $n = 18$ ), trichostatin A ( $n = 12$ ), vorinostat ( $n = 2$ ), and HC toxin ( $n = 1$ ) instances in the data set are colored in black. Colors applied to the remaining instances reflect the sign of their scores (green, positive; gray, null; red, negative). The rank, name [instance id], concentration, cell line, and connectivity score for each of the selected HDAC inhibitor instances is shown. Unabridged results from this query are provided as Result S1. (B) Chemical structures.

rank	perturbagen	dose	cell	score
1	vorinostat [1000]	10 $\mu$ M	MCF7	1
2	trichostatin A [873]	1 $\mu$ M	MCF7	0.969
3	trichostatin A [992]	100 nM	MCF7	0.931
4	trichostatin A [1050]	100 nM	MCF7	0.929
5	vorinostat [1058]	10 $\mu$ M	MCF7	0.917
6	trichostatin A [981]	1 $\mu$ M	MCF7	0.915
7	HC toxin [909]	100 nM	MCF7	0.914
8	trichostatin A [1112]	100 nM	MCF7	0.908
9	trichostatin A [1072]	1 $\mu$ M	MCF7	0.906
10	trichostatin A [1014]	1 $\mu$ M	MCF7	0.893
11	trichostatin A [332]	100 nM	MCF7	0.882
12	trichostatin A [331]	100 nM	MCF7	0.846
13	trichostatin A [448]	100 nM	PC3	0.788
14	valproic acid [345]	10 mM	MCF7	0.743
15	valproic acid [23]	1 mM	MCF7	0.735
16	valproic acid [1047]	1 mM	MCF7	0.733
17	trichostatin A [413]	100 nM	ssMCF7	0.725
18	valproic acid [410]	10 mM	HL60	0.725
19	valproic acid [458]	1 mM	PC3	0.680
33	valproic acid [409]	1 mM	HL60	0.634
39	valproic acid [1020]	500 $\mu$ M	MCF7	0.619
52	valproic acid [348]	2 mM	MCF7	0.582
61	valproic acid [1078]	500 $\mu$ M	MCF7	0.563
71	valproic acid [629]	1 mM	SKMEL5	0.539
72	valproic acid [347]	500 $\mu$ M	MCF7	0.539
73	valproic acid [989]	1 mM	MCF7	0.538
76	valproic acid [433]	1 mM	PC3	0.528
89	trichostatin A [364]	100 nM	HL60	0.507
92	valproic acid [497]	1 mM	ssMCF7	0.501
297	valproic acid [348]	50 $\mu$ M	MCF7	0
388	valproic acid [994]	200 $\mu$ M	MCF7	0
403	valproic acid [1002]	50 $\mu$ M	MCF7	0
419	valproic acid [1060]	50 $\mu$ M	MCF7	-0.537

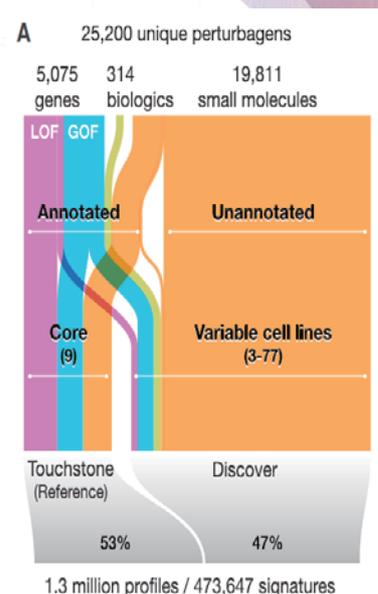
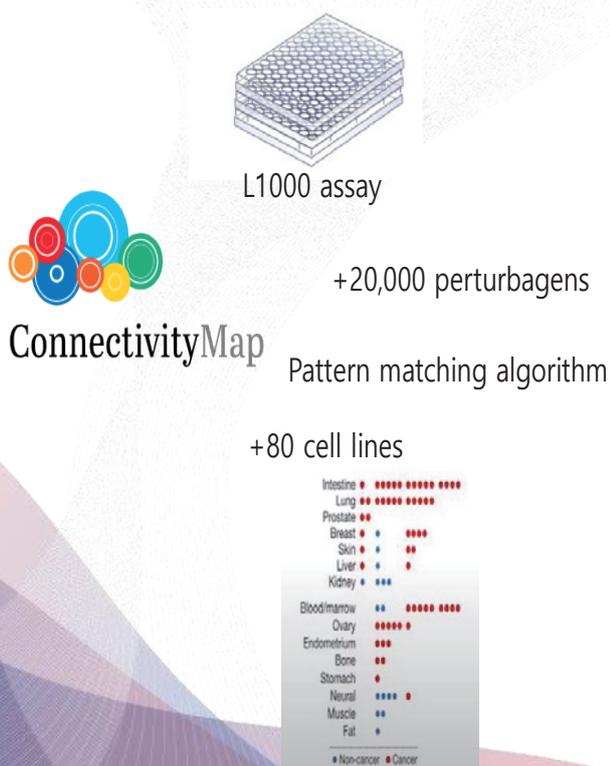
10

## Results

- anthelmintic drug parbendazole as an inducer of osteoclast differentiation (Brum et al., 2015)
- celastrol as a leptin sensitizer (Liu et al., 2015)
- compounds targeting COX2 and ADRA2A as potential diabetes treatments (Zhang et al., 2015)
- small molecules that mitigate skeletal muscular atrophy (Dyle et al., 2014) and spinal muscular atrophy (Farooq et al., 2009)
- new therapeutic hypotheses for the treatment of inflammatory bowel disease (Dudley et al., 2011) and cancer (Singh et al., 2016; Muthuswami et al., 2013; Wang et al., 2008; Schnell et al., 2015; Fortney et al., 2015; Wang et al., 2011; Churchman et al., 2015; Rosenbluth et al., 2008; Saito et al., 2009; Stockwell et al., 2012)

11

## Next-Generation ConnectivityMap (CMap v2)

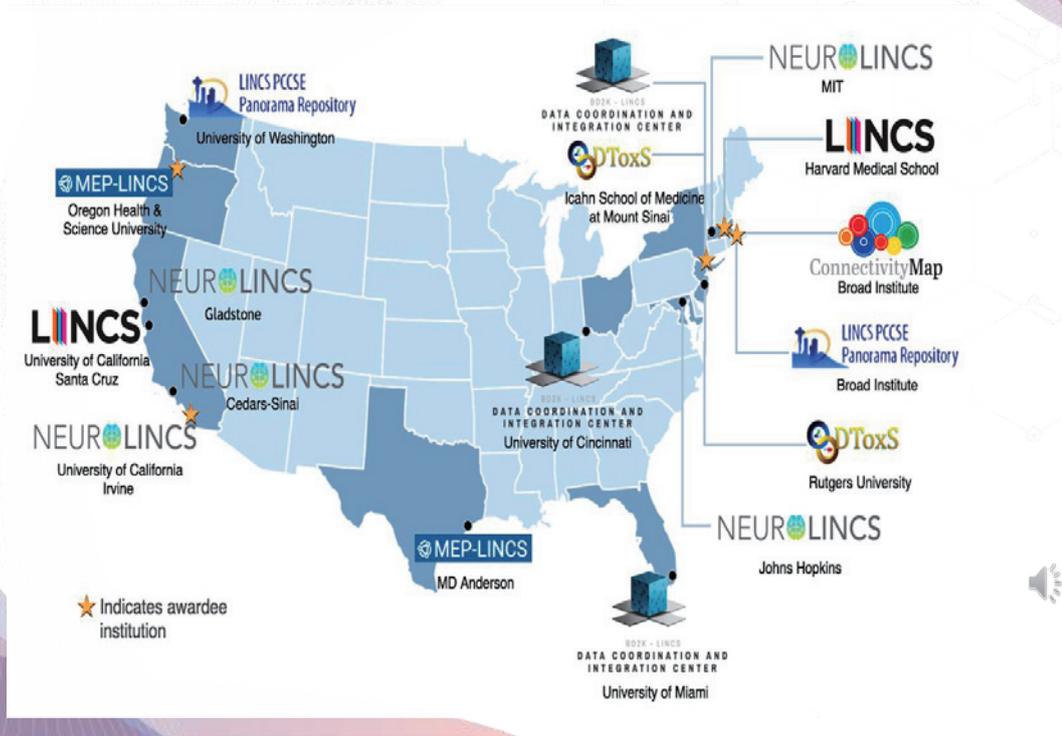


CMap version 2 with 1.3 M profiles

12

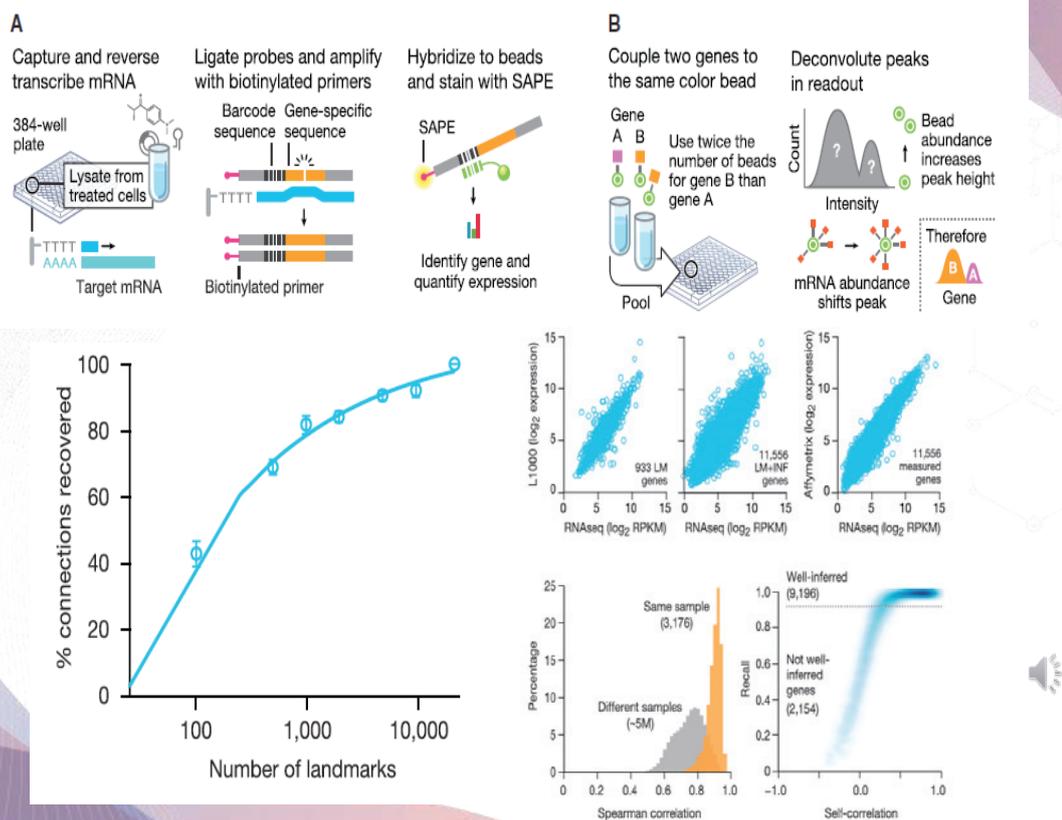
# LINCS Consortium

- The Library of Integrated Network-Based Cellular Signatures (LINCS)



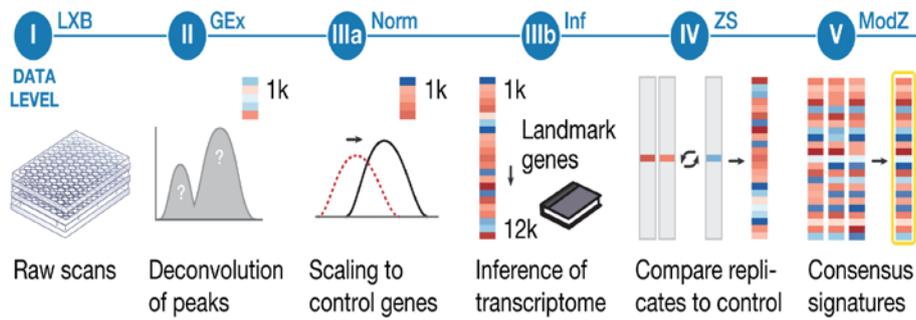
13

# L1000 Assay Enables Massive Scale Up of Data

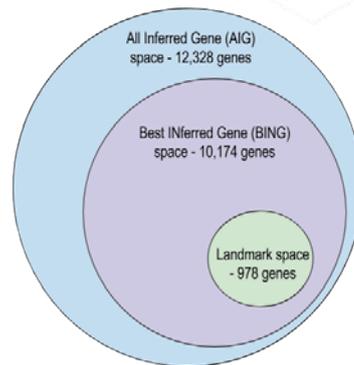


14

# L1000 Assay Enables Massive Scale Up of Data



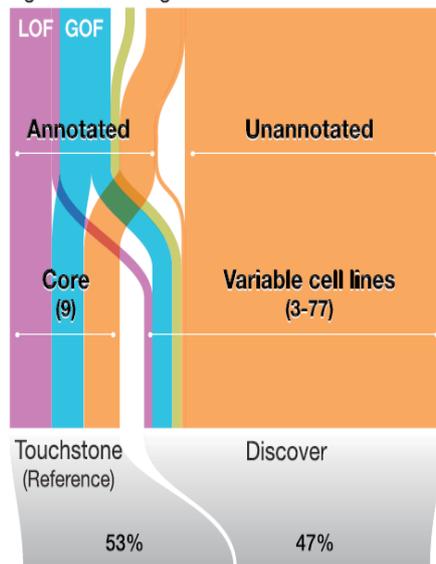
	L1000	RNA-Seq
Cost	\$	\$\$\$
Library prep?	Fast	Slow
Detection of non-abundant transcripts	No change	Needs deeper coverage
"Accuracy"	High	High
Evidence of additional benefit		?



15

# 1M Profiles of CMap v2

**A** 25,200 unique perturbagens  
 5,075 genes    314 biologics    19,811 small molecules



1.3 million profiles / 473,647 signatures

Perturbagen Type	pert_type designation in metadata files
Compound	trt_cp
Peptides and other biological agents (e.g. cytokine)	trt_lig
shRNA for loss of function (LoF) of gene	trt_sh
Consensus signature from shRNAs targeting the same gene	trt_sh.cgs
cDNA for overexpression of wild-type gene	trt_oe
cDNA for overexpression of mutated gene	trt_oe.mut
CRISPR for LLoF	trt_xpr
Controls - vehicle for compound treatment (e.g. DMSO)	ctl_vehicle
Controls - vector for genetic perturbation (e.g. empty vector, GFP)	ctl_vector
Controls - consensus signature from shRNAs that share a common seed sequence	trt_sh.css
Controls - consensus signature of vehicles	ctl_vehicle.cns
Controls - consensus signature of vectors	ctl_vector.cns
Controls - consensus signature of many untreated wells	ctl_untrt.cns
Controls - Untreated cells	ctl_untrt

16

## 1M Profiles of CMap v2

- small molecule compounds
  - ~1,300 FDA-approved drugs
  - ~5,585 bioactive tool compounds
  - 2000+ screening hits
- knocking-down genes (shRNA) or over-expressing genes
  - ~900 target/pathways of FDA-approved drugs
  - ~600 candidate disease genes
  - 500+ community nominations
- cells including primary cell lines, cancer cell lines, stem cell lines, and differentiated cell lines from different tissue types

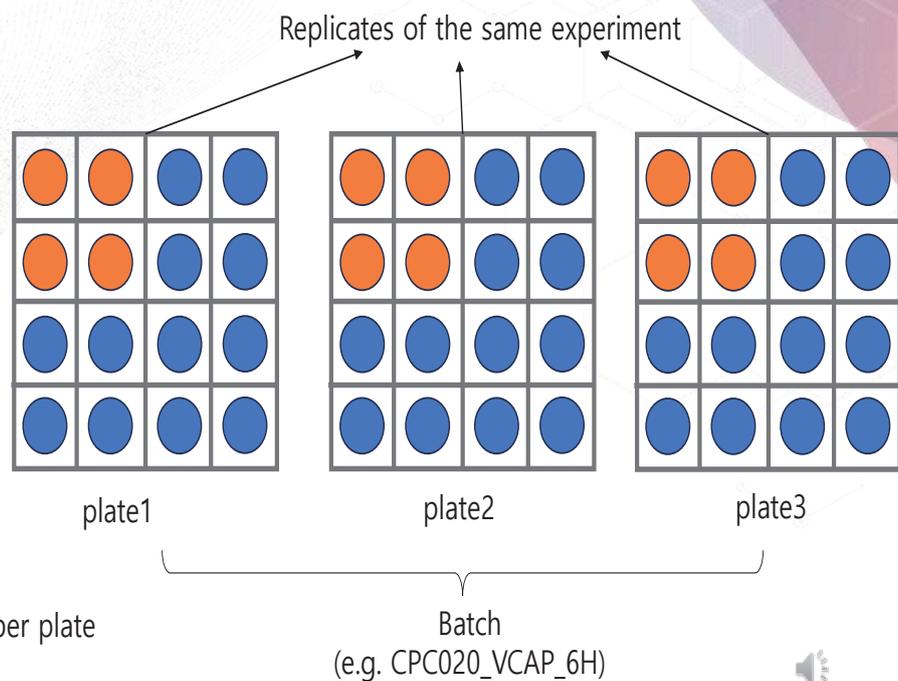


Big Data Science with the BD2K-LINCS DCIC by Avi Ma'ayan

17

## 1M Profiles of CMap v2

-  Control replicate
-  Experiment replicate



- 384 wells in each plate
- About 18 control replicates per plate
- About 2~4 plates per batch
- About 366 experiments per batch



Big Data Science with the BD2K-LINCS DCIC by Avi Ma'ayan

18

## Data Levels

Level 1 (LXB)	Raw fluorescent intensity (FI) values measured for every bead detected by Luminex scanners. Each 384-well plate generates 384 LXB files, where each file contains a fluorescent intensity value for each observed bead in the well.
Level 2 (GEX)	Gene expression levels for the 978 landmark genes, deconvoluted from the measured fluorescent intensity values.
Level 3 (NORM, INF)	NORM - Gene expression are <b>normalized</b> to invariant gene set curves and quantile normalized across each plate. INF- Additional values for <b>11,350 additional genes</b> not directly measured in the L10000 assay are inferred based on the normalized values for the 978 landmark genes.
Level 4 (ZS)	<b>Z-scores</b> for each gene based on Level 3 with respect to the entire plate population. This comparison of profiles to their appropriate population control generates a list of differentially expressed genes.
Level 5 (MODZ)	replicate-collapsed z-score vectors based on Level 4. Replicate collapse generates one differential expression vector, which we term a <b>signature</b> . Connectivity analyses are performed on signatures.

19

## ConnectivityMap Score

- Computing similarities - Weighted Connectivity Score (WTCS)

$$w_{q,r} = \begin{cases} (ES_{up} - ES_{down}) / 2, & \text{if } \text{sgn}(ES_{up}) \neq \text{sgn}(ES_{down}) \\ 0, & \text{otherwise} \end{cases}$$

Where  $ES_{up}$  is the enrichment of  $q_{up}$  in  $r$  and  $ES_{down}$  is the enrichment of  $q_{down}$  in  $r$ . WTCS ranges between -1 and 1

- Normalization of Connectivity Scores (NCS)

- Given a vector of WTCS values  $w$  resulting from a query, we normalize the values within each cell line ( $c$ ) and perturbation type ( $t$ ) to obtain normalized connectivity scores (NCS) as

$$NCS_{c,t} = \begin{cases} w_{c,t} / \mu_{c,t}^+ & \text{if } \text{sgn}(w_{c,t}) > 0 \\ w_{c,t} / \mu_{c,t}^- & \text{otherwise} \end{cases}$$

20

## ConnectivityMap Score

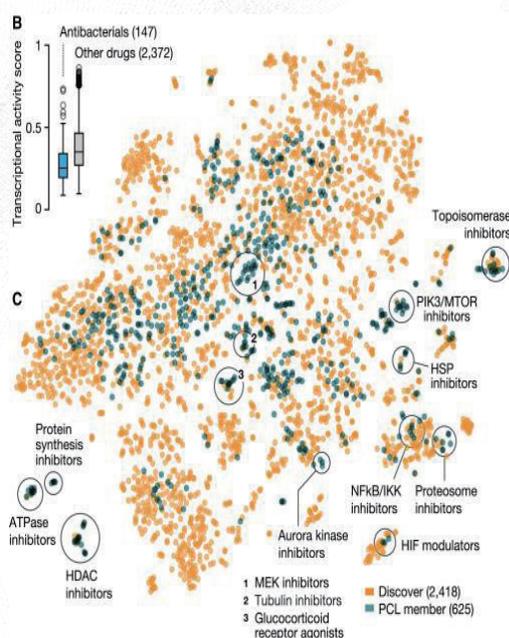
- Connectivity Map Score (Tau  $\tau$ )
  - by comparing each observed NCS value  $n_{cs_{q,r}}$  between the query  $q$  and a reference signature  $r$  to a distribution of NCS values representing the similarities between a reference compendium of queries ( $Q_{ref}$ ) and  $r$
  - Tau ( $\tau$ ) that ranges from  $-100$  to  $+100$  and represents the percentage of queries in  $Q_{ref}$  with a lower  $|NCS|$  than  $|n_{cs_{q,r}}|$

$$\tau_{q,r} = \text{sgn}(n_{cs_{q,r}}) \frac{100}{N} \sum_{i=1}^N [ |n_{cs_{i,r}}| < |n_{cs_{q,r}}| ]$$

21

## Results

- Discovery of MOA of Unannotated Small Molecules



**Figure 5. Characterizing Known and Unexpected Activities of Small Molecules**

(A) HDAC inhibitor PCL substructure. Hierarchical clustering of pairwise connectivities of the HDAC inhibitor PCL members reveals substructure within the class. Pan-HDAC inhibitors cluster together, distinct from more isoform-selective compounds.

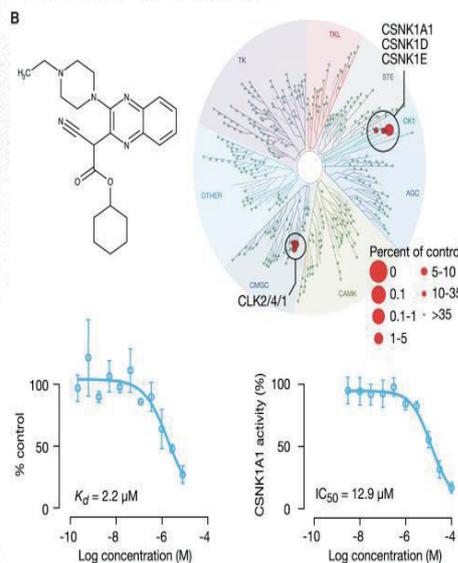
(B) Antibacterials exhibit lower transcriptional activity than other drugs. Distributions of the maximum TAS per compound for 147 antibacterials and 2,372 known drugs in CMap-Touchstone (TS). The antibacterials' TAS distribution is significantly lower ( $p < 3^{-11}$ ) than that of other drugs.

(C) Comparison of unannotated compounds with known drugs. t-SNE projection of the signatures of 2,418 unannotated but transcriptionally active compounds (orange) with PCL members (teal). Some unannotated compounds occupy regions not covered by drugs, presenting opportunities for novel chemical development.

22

# Results

- Discovery of a Selective CSNK1A1 inhibitor



23

# Results

**Table 4.** An overview of the application of CMap for a number of different diseases

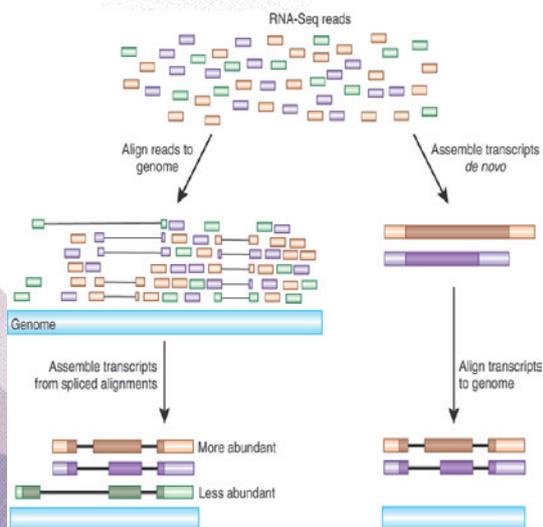
Disease	Method	Data set	Result	Drug	Reference
CNS injuries	CMap tool	Human MCF7 breast adenocarcinoma (GSE34331)	The findings show the hypothesis that inhibition of calmodulin signaling might allow neurons to alleviate substrate derived neurite growth restriction and CNS regeneration.	Calmodulin and piperazine phenothiazine (repurposed)	[54]
GBM	Pathway analysis and CMap tool	GBM data sets (GSE4290, GSE7696, GSE14805, GSE15824 and GSE16011)	Investigated antitumor drugs in GBM cell lines and identify novel drugs that can suppress GBM tumors.	Thioridazine	[55]
Gaucher disease (GD1)	Pathway analysis and CMap tool	GD1 mouse (GSE2308)	Predicted highly enriched anti-helminthic compounds for new drug action on GD1 and repurposing.	Albendazole and oxamniquine	[52]
Ovarian cancer	CMap tool	MCF7 and PC3 cell lines (GSE5258)	Found a compound as PI3K/AKT pathway inhibitor that shows the mechanism of cancer therapeutics.	Thioridazine	[56]
Stem cell leukemia (SCL)	GSEA and CMap tool	hESCs cell lines (GSE54508)	Found two HDAC inhibitors as potential inducers that can be used in treating SCL and acute megakaryoblastic leukemias.	Trichostatin A and suberoylanilide hydroxamic acid	[57]
T-cell acute lymphoblastic leukemia (T-ALL)	GSEA and CMap tool	Human and mouse T-ALL cell lines (GSE12948, GSE8416 and GSE14618)	Identified interconnecting regulatory pathways as therapeutic targets for T-ALL.	HDAC, PI3K and HSP90 inhibitors	[51]
Prostate cancer	CMap tool	Celastrol- and gedunin-treated cell lines (GSE5505 and GSE5508)	Identified target pathways of androgen receptor (AR) signaling and modulation of HSP90 MoA.	Celastrol and gedunin	[17]
Gastric cancer	Hierarchical clustering and CMap tool	Yonsei gastric cancer (GSE13861)	Predicted two possible drug candidates for gastric cancer therapy.	Vorinostat and trichostatin A	[53]
Myelomatosis	CMap tool	Human myeloma cell lines (GSE14011)	Found a drug with potential to induce suppression of cyclin D2 promoter regulation.	Pristimerin	[58]
AML	CMap tool	AML data (GSE7538)	Predicted novel treatment of human primary AML with parthenolide and transcriptional response of cells.	Celastrol	[59]

# 세포주 약물 반응 데이터의 활용



## How to get gene signatures?

- The most fundamental problem:
  - Detecting and interpreting differences in the abundance of genes or other genomic features between experimental conditions, cell types, or disease states

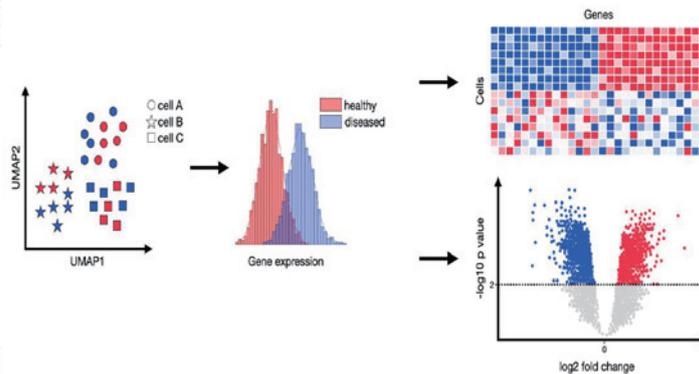


	Wild-type		Mutant	
	Mouse 1	Mouse 2	Mouse 1	Mouse 2
Gene 1	45	60	30	39
Gene 2	0	4	3	7
Gene 3	1010	800	3099	3450
...	...	...	...	...



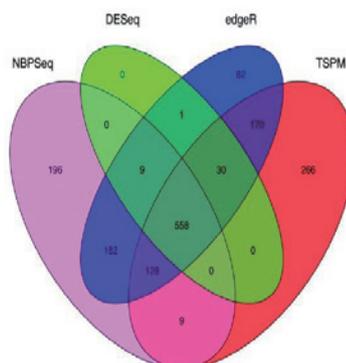
## Differential expression analysis

- Differential expression analysis tests thousands of hypotheses (one test for each gene) for gene activity changes between conditions (case and control).
- Factors affecting analysis power include limited biological replicates, non-normal distribution of read counts, and measurement uncertainty for lowly expressed genes.



## Differential expression analysis

- There are a number of software packages that have been developed for differential expression analysis of RNA-seq data
- Tools like edgeR and DESeq2 overcome these limitations using statistical models based on negative binomial distribution.
- There is no one method that performs optimally under all conditions

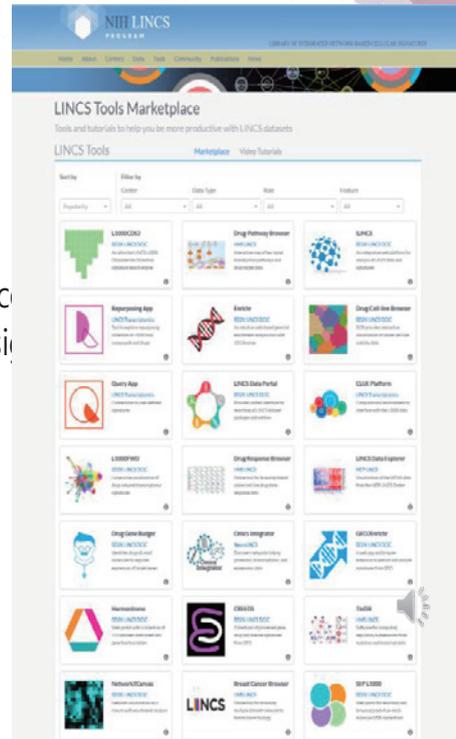


Soneson, Charlotte, and Mauro Delorenzi  
*BMC bioinformatics* (2013)

# NIH LINCS Program

- LINCS Tools
  - >50 tools developed by the consortium
  - Interactive visualizations
  - Data querying and browsing
  - Various analysis workflows
- LINCS Data
  - Transcriptomics, proteomics, epigenomics, imaging and more
  - L1000 data contain >3 million samples, >1 million signatures
  - Programmatic access via APIs
  - User interfaces for querying and viewing signatures

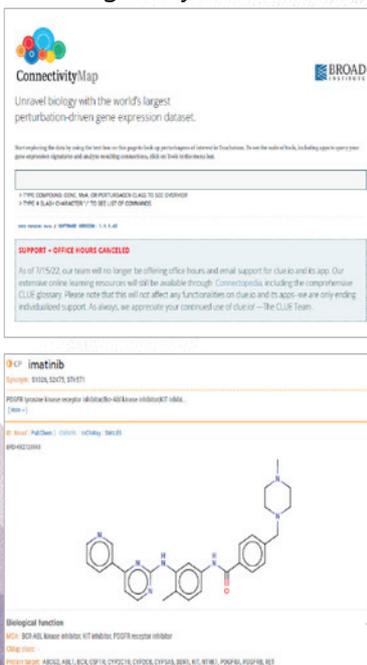
<https://lincsproject.org/LINCS/tools>



30

# ConnectivityMap

- Visit [clue.io](http://clue.io)
  - Single keyword search



**/conn "imatinib"**

Cell Lines Summary

The "score" column is the median connectivity score across the selected perturbagens and cell lines. The top and bottom scores are shown separately for perturbagen classes, compounds, genetic perturbations (overexpressions and knockdowns).

score	type	name	description
99.03	compound	imatinib	PDGFR tyrosine kinase receptor inhibitor, Src-42k kinase inhibitor, KIT inhibitor, ABL kinase inhibitor, apoptosis stimulant, breast cancer resistance protein inhibitor, colony stimulating factor receptor inhibitor
99.02	compound	suripide	cardiac amyloid inhibitor, dopamine receptor, dopamine receptor antagonist
99.02	compound	lorazepam	chloride channel agonist, GABA benzodiazepine site receptor agonist, GABA receptor agonist, potassium channel agonist
99.02	compound	pericidin	cyclooxygenase inhibitor
99.02	compound	YM-155	survivin inhibitor, XIAP expression inhibitor
99.02	compound	afatinib	antitubulin, interacts with heme and with COX3
99.02	compound	delestatinA	histamine receptor antagonist
99.02	compound	celastrol	cell wall synthesis inhibitor
99.02	compound	propranolol	β1/β2 adrenergic receptor inhibitor, PDK activator
99.02	compound	fludocortisone	glucocorticoid receptor agonist, mineralocorticoid receptor agonist
99.02	compound	ibuprofen	CXCR2 chemokine agonist, coagulation stimulant, cysteine stimulant
99.02	compound	bicalutamide	glucocorticoid receptor agonist, corticosteroid hormone receptor agonist, immunosuppressant
99.02	compound	PD-80599	MEK inhibitor, MAP kinase inhibitor
99.02	compound	cholic-acid	benzothiazole inhibitor, unidentified pharmacological activity
99.02	compound	levamisole	guanine receptor antagonist, histamine receptor antagonist, serotonin receptor antagonist
99.02	compound	desoxycurcumin	mineralocorticoid receptor agonist
99.02	compound	trypsin/histatin-5/lysine	heparin activation inhibitor
99.02	compound	PD-184352	MEK inhibitor, MAP kinase inhibitor
99.02	compound	AS-605240	Phosphatidylinositol 3-kinase (PI3K) inhibitor, PI3K inhibitor
99.02	compound	SAS-60148316	HIV integrase inhibitor

**Compound**

score	type	name	description
99.02	compound	imatinib	PDGFR tyrosine kinase receptor inhibitor, Src-42k kinase inhibitor, KIT inhibitor, ABL kinase inhibitor, apoptosis stimulant, breast cancer resistance protein inhibitor, colony stimulating factor receptor inhibitor
99.02	compound	suripide	cardiac amyloid inhibitor, dopamine receptor, dopamine receptor antagonist
99.02	compound	lorazepam	chloride channel agonist, GABA benzodiazepine site receptor agonist, GABA receptor agonist, potassium channel agonist
99.02	compound	pericidin	cyclooxygenase inhibitor
99.02	compound	YM-155	survivin inhibitor, XIAP expression inhibitor
99.02	compound	afatinib	antitubulin, interacts with heme and with COX3
99.02	compound	delestatinA	histamine receptor antagonist
99.02	compound	celastrol	cell wall synthesis inhibitor
99.02	compound	propranolol	β1/β2 adrenergic receptor inhibitor, PDK activator
99.02	compound	fludocortisone	glucocorticoid receptor agonist, mineralocorticoid receptor agonist
99.02	compound	ibuprofen	CXCR2 chemokine agonist, coagulation stimulant, cysteine stimulant
99.02	compound	bicalutamide	glucocorticoid receptor agonist, corticosteroid hormone receptor agonist, immunosuppressant
99.02	compound	PD-80599	MEK inhibitor, MAP kinase inhibitor
99.02	compound	cholic-acid	benzothiazole inhibitor, unidentified pharmacological activity
99.02	compound	levamisole	guanine receptor antagonist, histamine receptor antagonist, serotonin receptor antagonist
99.02	compound	desoxycurcumin	mineralocorticoid receptor agonist
99.02	compound	trypsin/histatin-5/lysine	heparin activation inhibitor
99.02	compound	PD-184352	MEK inhibitor, MAP kinase inhibitor
99.02	compound	AS-605240	Phosphatidylinositol 3-kinase (PI3K) inhibitor, PI3K inhibitor
99.02	compound	SAS-60148316	HIV integrase inhibitor

**Genetic**

score	type	name	description
99.03	genetic	ZNF92	Zinc finger, C2H2-type, zinc finger protein 92
99.03	genetic	RNF42	Receptor interacting protein kinase (RIPK) family, receptor-interacting serine/threonine kinase 2
99.03	genetic	RAD50	RAD50 homolog (D. cerevisiae)
99.03	genetic	FOXO4L3	forkhead box O4-like 3
99.03	genetic	STAT4	SH2 domain containing, signal transducer and activator of transcription 4
99.03	genetic	SLC3A2	SLC3 family, solute carrier family 3 (anion/cationic and neutral amino acid transport), member 2
99.03	genetic	RPN1	ribosomal protein P1
99.03	genetic	PleioD1	Pleio domain containing 1
99.03	genetic	SKP4	GPCR / Class A, RhoGuan family peptide receptors, relasin/relasin-like family peptide receptor 4
99.03	genetic	MYL6	EF-hand domain containing, myosin light chain 6B, uterine, smooth muscle and non-muscle
99.03	genetic	ATP7C1	ATPases, F-type, ATP synthase, iron transporting, mitochondrial F1 complex, gamma polypeptide 1
99.03	genetic	SOX5	SRY (sex determining region Y)-box 5, SRY (sex determining region Y)-box 5
99.03	genetic	MLK1	MLK1 (cellular), mammalian embryonic leukemia zipper kinase
99.03	genetic	STAT2	SH2 domain containing, signal transducer and activator of transcription 2, 113kDa
99.03	genetic	AGTR1	Angiotensin receptors, angiotensin II receptor, type 1
99.03	genetic	VCAN	Immunoglobulin superfamily V-set domain containing, versican
99.03	genetic	EP18	eukaryotic translation initiation factor 18
99.03	genetic	ZNF19	Zinc finger, C2H2-type, zinc finger protein 19
99.03	genetic	PCID2	proteoglycan-3, 2-oxoglutarate 5-lyase-pore 2
99.03	genetic	BRF2	BRF2, subunit of RNA-polymerase III transcription initiation factor, BRF1-like

31

# ConnectivityMap

- Search tools

**Tools**

-  **Query**  
Find perturbagens that give rise to similar (or opposing) expression signatures
-  **Touchstone**  
Explore connectivities between signatures from ~3,000 drugs and genetic loss/gain of function of ~2,000 genes that make up the CMap touchstone (reference) database
-  **Data Library**  
Explore datasets available through clue.io including L1000 cohorts and related perturbational information
-  **Repurposing**  
Explore our collection of ~5000 drugs and tool compounds to find potential drug repurposing opportunities to improve disease treatments
-  **Morpheus**  
Explore, analyze, and annotate heat maps. Choose an existing dataset or upload your own data (for example, gene expression or connectivity scores)
-  **Metadata Browser**  
CLUE Metadata browser

# ConnectivityMap

- Query signature
  - To find perturbagens that give rise to similar (or opposing) expression signatures
  - Take some time.. ~50min

**Query**

Query CMap for reference perturbagen signatures most similar (or dissimilar) to your sample.

**Note that choosing 'Latest' from the query parameters section below. will run the query against our beta dataset released on (Dec 17, 2020)**

1) Name your query  
Please note that names must contain only alphanumeric characters. Any non-alphanumeric characters will be stripped.

test

2) Query parameters  
Gene expression (L1000) Touchstone Individual query  
Latest

3) Load a collection of Ensembl Gene IDs from Litteralis for up-regulated gene sets (and optionally a collection for down-regulated gene sets). At any time you may choose an [example](#) to fill in the boxes for the individual query.

**UP-regulated genes**

Load from my files Enter 1-150 genes for optimal results.

- TARGP1
- APP
- RAPP1GAP
- SPN1
- DNALJ3
- FOSD1
- CSRP1
- MBNL2

**DOWN-regulated genes (optional)**

Load from my files Please note that 150 is a maximum.

- UP1000
- POLE2
- COX2SA
- BTK
- SFT1
- KIF14
- TSG1
- ASPB2
- HMOX1

Invalid gene: Not valid HUGO symbol or Ensembl ID, not used in query

Valid gene: Valid HUGO symbol or Ensembl ID and part of BING space, used in query

Valid but not used in query: Valid HUGO symbol or Ensembl ID not part of BING space, not used in query

More information can be found in this [Connectopedia](#) article

id	name	cell_name	pert_type	pert_dose	pert_time	moa	nsample	tas	raw_cs	for_q_nlog10
JUN	H1299	trt_oe	---	96 h	-666		3.00	0.66	0.62	15.65
BRD-K98645085	PC3	trt_cp	20 uM	24 h	-666		3.00	0.57	0.60	3.66
TG-101348	PC3	trt_cp	10 uM	24 h	JAK inhibitor FLT3 inhibitor		2.00	0.52	0.59	3.33
Isaalloid	HA1E	trt_cp	10 uM	24 h	Bacterial permeability inducer		3.00	0.61	0.58	3.25
JUN	H1299	trt_oe	---	96 h	-666		3.00	0.60	0.58	3.20
HG-9-91-01	NPC	trt_cp	3.33 uM	24 h	-666		3.00	0.53	0.58	3.19
PFI-1	NPC	trt_cp	10 uM	24 h	Bromodomain inhibitor		3.00	0.58	0.57	3.17
sunlitinib	HT29	trt_cp	10 uM	24 h	FLT3 inhibitor KIT inhibitor PDGFR inhibitor RET inh		3.00	0.68	0.57	3.13
volasertib	HCC515	trt_cp	3.33 uM	24 h	PLK inhibitor		3.00	0.55	0.57	3.09
BMS-833923	A375	trt_cp	10 uM	24 h	Smoothened receptor antagonist		3.00	0.63	0.56	3.06
GSK-1070916	A549	trt_cp	10 uM	24 h	Aurora kinase inhibitor		2.00	0.55	0.56	3.06
AKT1	H1299	trt_oe	---	96 h	-666		2.00	0.64	0.56	3.00
I-BET-762	HBL1	trt_cp	10 uM	24 h	Bromodomain inhibitor		3.00	0.63	0.55	2.99
PTTG1	HT29	trt_sh	---	96 h	-666		4.00	0.43	0.55	2.98
Isaalloid	THP1	trt_cp	10 uM	24 h	Bacterial permeability inducer		2.00	0.40	0.55	2.98
BRD-A19037878	HL60	trt_cp	0.37 uM	2 h	-666		6.00	0.85	0.55	2.96
F10	HT29	trt_xpr	-666	96 h	-666		3.00	0.16	0.55	2.94
I-BET-151	A375	trt_cp	3.33 uM	24 h	Bromodomain inhibitor		3.00	0.69	0.54	2.93
MK-2206	HCC515	trt_cp	10 uM	24 h	AKT inhibitor		3.00	0.48	0.54	2.93

4) Review and submit. Only valid genes will be used in your query.

SUBMIT

- 16 -





# L1000CDS2

## • Query and results

L1000CDS2 is being developed by the Malayan Lab at the Institute School of Medicine at Mount Sinai for the BD3K-LINCS DDC and the KDC-DD. The L1000CDS2 test is currently released at the LINCS 1,000 small molecule expression profiles generated at the Broad Institute by the Connectivity Map Group. This work is supported by NIH Grants U54-HL127624 and U54CA195071.

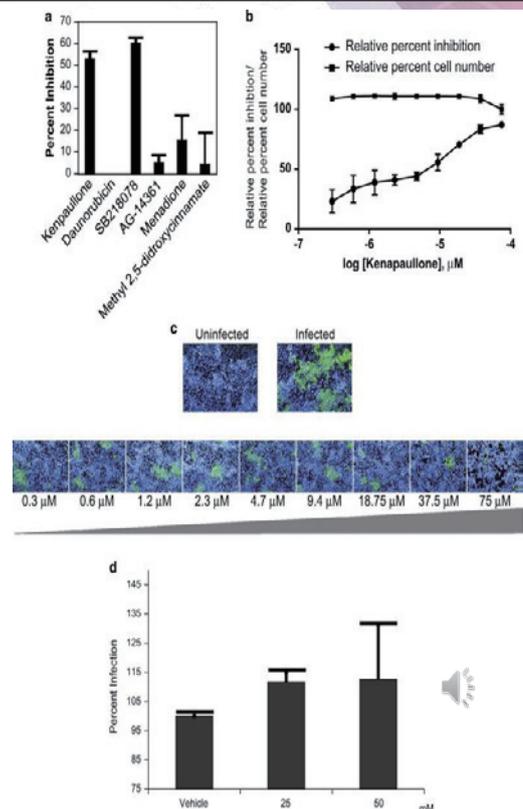
# L1000CDS2

## • Query and results

- Input: Ebola virus signatures
- Top candidate: Kenpaullone

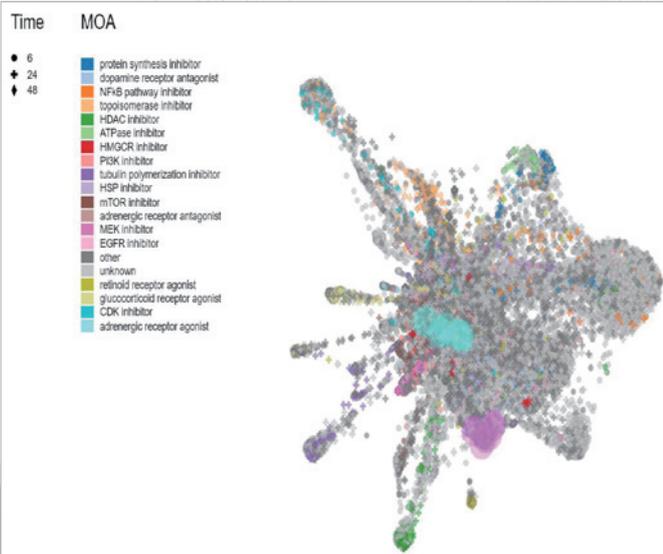
**Table 1. Top five predicted drugs at each time point**

Drug name	Cosine distance	Batch	Cell line	Time point (h)	Concentration ( $\mu$ M)
<b>30 min</b>					
Kenpaullone	1.2584	CPC002	HA1E	6	10
0800-0289	1.1978	CPC014	A549	24	10
BRD-K37312348	1.1965	CPC016	HT29	6	10
SB 225002	1.1896	CPC001	HA1E	24	10
10006350	1.1856	CPC012	MCF7	24	10
<b>60 min</b>					
Kenpaullone	1.2756	CPC002	HA1E	6	10
PD 166793	1.2619	CPC001	HA1E	24	10
BAPTA-AM	1.2564	CPC001	HA1E	24	10
BRD-U74615290	1.2396	CPC014	HT29	6	10
NGC00229596-01	1.2303	CPC008	HT29	6	10
<b>120 min</b>					
Kenpaullone	1.3489	CPC002	HA1E	6	10
NSC 23766	1.3478	CPC006	PC3	24	160
Rosiglitazone	1.3386	CPC006	HA1E	24	80
7-hydroxy-2, 3, 4, 5-tetrahydro-1H-[1]benzofuro[2, 3-c]azepin-1-one	1.3351	CPC007	A549	24	10
LY 364947	1.3183	CPC003	PC3	24	10





# L1000FWD



Shape by:

Show detailed results:

Search compounds:

Type the name of a drug:

Signature Similarity Search:

Up genes:

- ZNF238
- ACACA
- ACAT2
- ACLY
- ACSL1

Signature similarity search results

Similar signatures:  Opposite signatures:

Show: 10 entries

sig_id	drug	similarity score	p-value	q-value	Z-score	combined score
CPC006_VCAP_8HBRD-K70792160-003-02-0-24	10-DEBC	0.8338	1.73e-139	7.43e-136	-1.73	269.79
CPC008_VCAP_8HBRD-AT5301702-001-01-8-10	BRD-AT5301702	0.5463	2.50e-74	5.35e-70	-1.74	128.06
CPC001_VCAP_8HBRD-A84413429-003-01-8-10	NTNC8	0.5278	6.71e-71	9.57e-67	-1.68	117.56
CPC001_VCAP_8HBRD-K8221984-001-01-8-10	T-69478	0.5000	1.75e-62	8.24e-58	-1.61	96.58
CPC020_VCAP_8HBRD-K40645748-003-11-8-10	metformine	0.4815	7.67e-62	2.35e-58	-1.85	112.89
CPC001_VCAP_8HBRD-K82289640-003-01-8-10	lyamline	0.4815	2.43e-64	1.48e-60	-1.72	109.44
CPC007_A375_24HBRD-K33065493-001-02-8-10	BRD-K33065493	0.4815	1.02e-63	4.87e-60	-1.78	111.95
ERG005_VCAP_24HBRD-K55127134-300-06-2-20	fluphenazine	0.4815	1.26e-62	4.91e-59	-1.91	118.03
ERG005_VCAP_48HBRD-K55127134-300-06-2-20	fluphenazine	0.4815	1.25e-64	8.90e-61	-1.94	124.14
CPC008_VCAP_8HBRD-K43692206-001-01-2-10	VU-0413243.1	0.4815	3.66e-63	1.56e-59	-1.75	109.29

Showing 1 to 10 of 50 entries

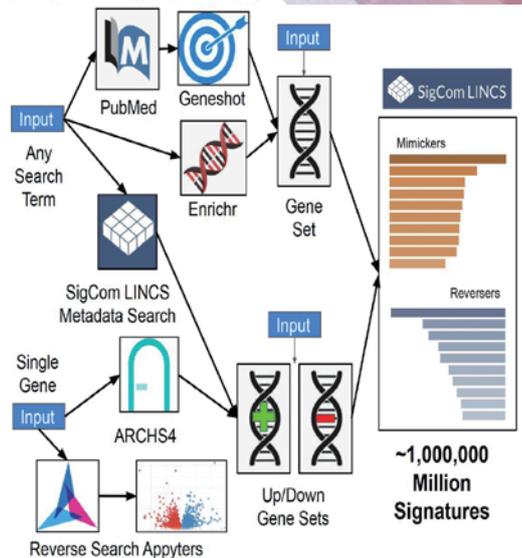
Previous 1 2 3 4 5 Next

Use Ctrl+C or Option+C to copy the link below to share.

<https://maayanlab.cloud/l1000fwd/result/64fbf8bc9967702c2b3a2>

# SigCom LINCS

- A cloud-agnostic platform for serving, storing, and que
  - >1 million signatures from
  - LINCS
  - GTEx (Genotype-Tissue Expression)
    - Signatures of aging
  - GEO (Gene Expression Omnibus)
    - RNA-seq signatures automatically extracted from GEO
  - CREEDS (Crowd Extracted Expression of Differential Signatures)
    - Participants extracted signatures from GEO



<https://maayanlab.cloud/sigcom-lincs>



# SigCom LINCS

The screenshot displays the SigCom LINCS web interface. On the left, there are input fields for 'Up Genes' and 'Down Genes', a search bar, and a 'Search gene' button. Below these are statistics for datasets (431), small molecules (33866), CRISPR knockouts (7492), data and signature generation centers (7), shRNAs (4957), gene over-expressions (4082), signatures (1113059), cell lines (353), and ligands (329). The main area shows 'Signature Similarity Search Results' with three columns: 'Automatic Human GEO RNA-seq Signatures', 'GTEx Tissue-Specific Age Group Signatures', and 'LINCS L1000 Chemical Perturbations (2021)'. A table at the bottom lists results with columns for gene name, time, and scores.

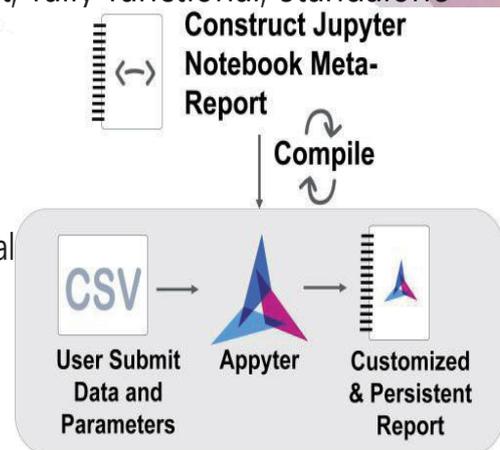
Gene	Time	Score 1	Score 2	Score 3
state... WCAP	48 h	31.21	12.87	-10.39
MCLF05...	6 h	26.34	12.34	-9.509
HEK293	24 h	25.00	8.298	-9.036
HEK293	24 h	21.69	10.21	-8.525
MCF10A	24 h	21.66	9.392	-8.486
YAPC	24 h	21.46	10.29	-8.478
VCAP	48 h	21.23	10.36	-8.433
HEC251	24 h	21.00	9.077	-8.337
HELA	24 h	20.88	9.427	-8.327

Evangelista, et al. NAR 2022

44

# Appyters

- Transform Jupyter Notebooks into lightweight, fully functional, standalone web-based bioinformatics applications
- Enable reusable workflows for
  - Customized machine learning pipelines
  - Analyzing omics data
  - Producing publishable figures and interactive visualizations

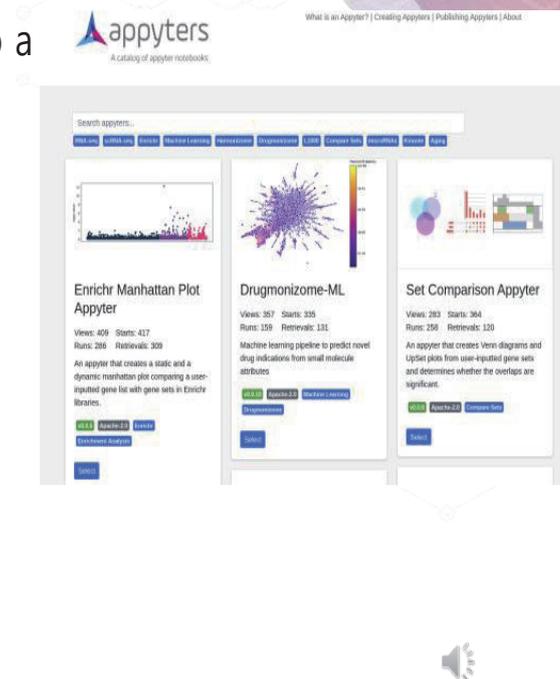


Clarke DJB et al. Patterns (2021)

45

# Appyter Catalog

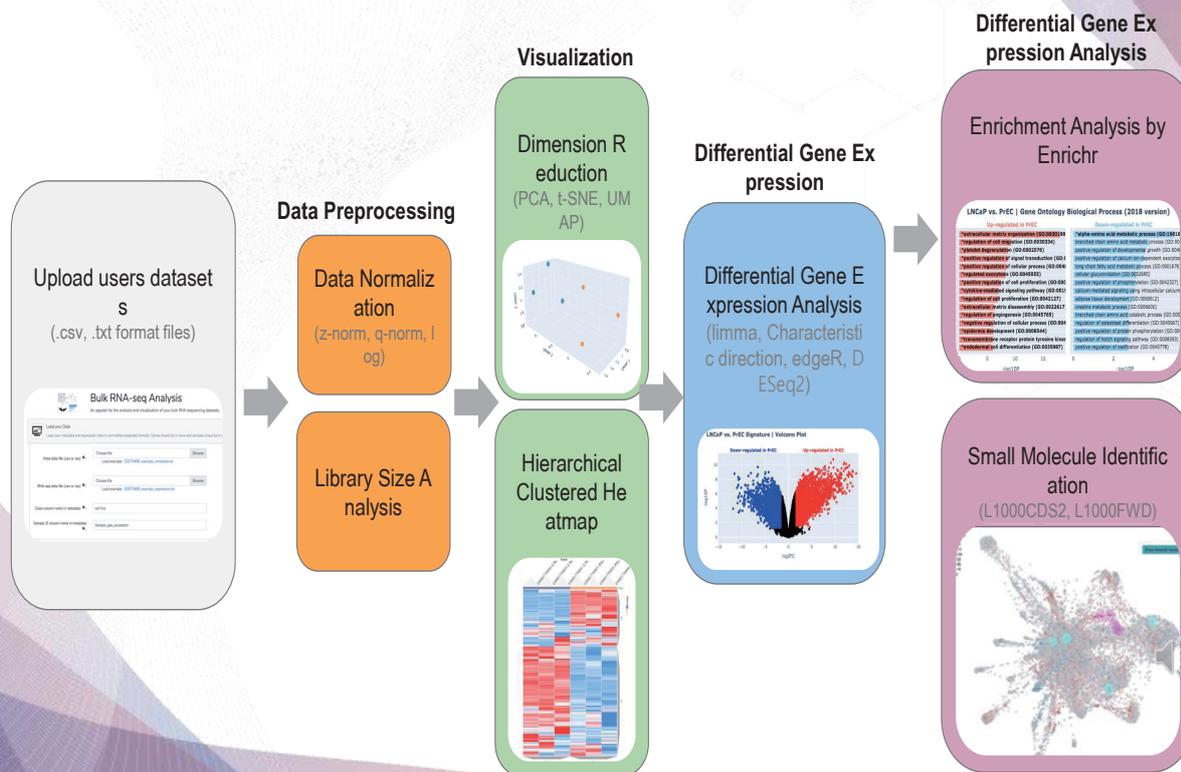
- Integrates all available Appyters (>100) into a
  - Github pull requests
  - Standardized, machine-validatable requirements
- Allows for categorization and search
- <https://appyters.maayanlab.cloud>



Clarke DJB et al. *Patterns* (2021) 46

# Bulk RNA-seq Analysis Appyter

URL: [https://appyters.maayanlab.cloud/#/Bulk\\_RNA\\_seq](https://appyters.maayanlab.cloud/#/Bulk_RNA_seq)



# Input for Bulk RNA-seq Analysis Appyter

Series GSE154613

Query DataSets for GSE154613

**Status** Public on Jul 17, 2020

**Title** Modulating the transcriptional landscape of SARS-CoV-2 as an effective method for developing antiviral compounds

**Organism** [Homo sapiens](#)

**Experiment type** Expression profiling by high throughput sequencing

**Summary** The pandemic of severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) has imposed a significant burden on the human population. To understand the virus and the disease it causes we sought to interfere with the transcriptional response of the infected host. Utilizing the expression pattern of SARS-CoV-2-infected cells, we identified a region in gene expression space that was unique to virus infection and inversely proportional to the transcriptional footprint of known compounds characterized in the Library of Integrated Network-based Cellular Signatures (LINCS). Here we demonstrate the successful identification of compounds that display efficacy in blocking SARS-CoV-2 replication based on their ability to counteract the virus-induced transcriptional landscape. These compounds were found to potently reduce viral load despite having no impact on viral entry or modulation of the host antiviral response in the absence of virus. RNA-Seq profiling implicated the induction of the cholesterol biosynthesis pathway as the underlying mechanism of inhibition and suggested that targeting this aspect of host biology may significantly reduce SARS-CoV-2 viral load.

URL: [https://appymers.maayanlab.cloud/#/Bulk\\_RNA\\_seq](https://appymers.maayanlab.cloud/#/Bulk_RNA_seq)



URL: [https://appymers.maayanlab.cloud/#/Bulk\\_RNA\\_seq](https://appymers.maayanlab.cloud/#/Bulk_RNA_seq)

# Input for Bulk RNA-seq Analysis Appyter

Supplementary file	Size	Download	File type/resource
GSE154613_RAW.tar	7.6 Mb	<a href="#">(http)(custom)</a>	TAR (of TXT)

[SRA Run Selector](#)

Raw data are available in SRA

Processed data provided as supplementary file

Custom GSE154613\_RAW.tar archive:

Supplementary file	File size
<input type="checkbox"/> GSM4675765_ACE2-A549_Amolipine_drugonly_3.counts.genes.txt.gz	116.2 Kb
<input checked="" type="checkbox"/> GSM4675766_ACE2-A549_Berbamine_COV2_1.counts.genes.txt.gz	117.7 Kb
<input checked="" type="checkbox"/> GSM4675767_ACE2-A549_Berbamine_COV2_2.counts.genes.txt.gz	117.2 Kb
<input checked="" type="checkbox"/> GSM4675768_ACE2-A549_Berbamine_COV2_3.counts.genes.txt.gz	116.8 Kb
<input type="checkbox"/> GSM4675769_ACE2-A549-Berbamine-drugonly-1.counts.genes.txt.gz	114.4 Kb
<input type="checkbox"/> GSM4675770_ACE2-A549-Berbamine-drugonly-2.counts.genes.txt.gz	116.3 Kb
<input type="checkbox"/> GSM4675771_ACE2-A549-Berbamine-drugonly-3.counts.genes.txt.gz	114.5 Kb
<input checked="" type="checkbox"/> GSM4675772_ACE2-A549-DMSO-COV2-1.counts.genes.txt.gz	116.1 Kb
<input checked="" type="checkbox"/> GSM4675773_ACE2-A549-DMSO-COV2-2.counts.genes.txt.gz	116.3 Kb
<input checked="" type="checkbox"/> GSM4675774_ACE2-A549-DMSO-COV2-3.counts.genes.txt.gz	116.5 Kb
<input type="checkbox"/> GSM4675775_ACE2-A549_Loperamide_COV2_1.counts.genes.txt.gz	117.3 Kb
<input type="checkbox"/> GSM4675776_ACE2-A549_Loperamide_COV2_2.counts.genes.txt.gz	117.2 Kb
<input type="checkbox"/> GSM4675777_ACE2-A549_Loperamide_COV2_3.counts.genes.txt.gz	117.6 Kb
<input type="checkbox"/> GSM4675778_ACE2-A549-Loperamide-drugonly-1.counts.genes.txt.gz	116.0 Kb
<input type="checkbox"/> GSM4675779_ACE2-A549-Loperamide-drugonly-2.counts.genes.txt.gz	115.6 Kb
<input type="checkbox"/> GSM4675780_ACE2-A549-Loperamide-drugonly-3.counts.genes.txt.gz	115.1 Kb
<input type="checkbox"/> GSM4675781_ACE2-A549-Mock-1.counts.genes.txt.gz	116.7 Kb
<input type="checkbox"/> GSM4675782_ACE2-A549-Mock-2.counts.genes.txt.gz	116.6 Kb
<input type="checkbox"/> GSM4675783_ACE2-A549-Mock-3.counts.genes.txt.gz	115.9 Kb

Select All

6 file(s), 700.7 Kb

	A	B	C	D	E	F	G
1							
2	DOX11L1	0	0	1	0	0	0
3	WASH7P	23	68	15	17	98	108
4	FAM138A	0	0	0	0	0	0
5	FAM138F	0	0	0	0	0	0
6	OR4F5	0	0	0	0	0	0
7	LOC729737	1	6	5	3	13	3
8	LOC1001322	0	0	0	0	0	0
9	LOC1001320	0	0	0	0	0	0
10	LOC1001333	34	79	33	29	93	74
11	OR4F29	0	0	0	0	0	0
12	OR4F16	0	0	0	0	0	0
13	OR4F3	0	0	0	0	0	0
14	LOC1002880	12	26	12	20	43	26
15	LINC00115	2	9	3	5	13	12
16	LOC543837	169	235	166	186	249	266
17	FAM41C	0	0	0	0	0	0
18	LOC1001304	27	36	26	17	39	32
19	SAMD11	156	410	139	125	525	420
20	NOX2L	925	1457	829	833	1908	1974
21	KLHL17	66	131	45	43	152	151
22	PLEKHN1	19	34	21	14	45	49
23	Clorf170	3	9	9	6	14	19

Sample_id	Class
GSM4675774	Ctrl
GSM4675767	Case
GSM4675773	Ctrl
GSM4675772	Ctrl
GSM4675766	Case
GSM4675768	Case



# Input for Bulk RNA-seq Analysis Appyter

- The bulk RNA-seq analysis Appyter starts with an expression matrix of raw read counts and metadata

**Bulk RNA-seq Analysis**  
An appyter for the analysis and visualization of your bulk RNA sequencing datasets.

**Load your Data**  
Load your metadata and expression data in comma/tab separated formats. Genes should be in rows and samples should be in columns

Meta data file (csv or tsv):    
Load example: [GSE70466\\_example\\_metadata.txt](#)

RNA-seq data file (csv or tsv):    
Load example: [GSE70466\\_example\\_expression.txt](#)

Class column name in metadata:

Sample ID column name in metadata:

URL: [https://appymaayanlab.cloud/#/Bulk\\_RNA\\_seq](https://appymaayanlab.cloud/#/Bulk_RNA_seq)

# Input for Bulk RNA-seq Analysis Appyter

**Select Normalization Methods**

Filter genes?  Yes  No

Low expression threshold:

logCPM normalization?  Yes  No

log normalization?  Yes  No

Z normalization?  Yes  No

Quantile normalization?  Yes  No

**Select Visualization Parameters**

Interactive plots?  Yes  No

Visualization Methods:

Genes for Dimension Reduction:

Gene List for Clustergrammer (Optional):

Genes for clustergrammer:

**Select Differentially Expressed Gene Analysis Parameters**

Differential expression analysis method:

Differential expression analysis plotting method:

P-value threshold:

logFC threshold:

Maximum genes for Enrich:

Enrich Libraries (upto 2):

Top ranked gene sets:

Small molecule analysis method:

Genes for L1000CDS2 or L1000FWD:

Top ranked drugs from L1000CDS2 or L1000FWD:

URL: [https://appymaayanlab.cloud/#/Bulk\\_RNA\\_seq](https://appymaayanlab.cloud/#/Bulk_RNA_seq)

# Results from Bulk RNA-seq Analysis Appyter

3D PCA plot for samples

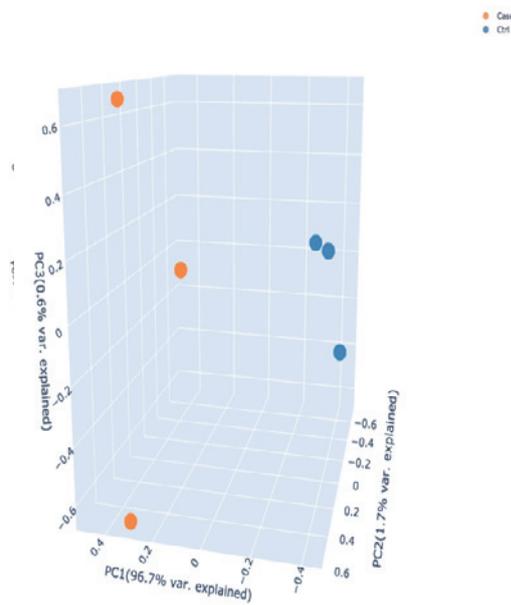


Figure 1. 3D UMAP plot for samples using 2500 genes having largest variance. The figure displays an interactive, three-dimensional scatter plot of the data. Each point represents an RNA-seq sample. Samples with similar gene expression profiles are closer in the three-dimensional space. If provided, sample groups are indicated using different colors, allowing for easier interpretation of the results.

URL: <https://appytters.maayanlab.cloud/#/Bulk RNA seq>

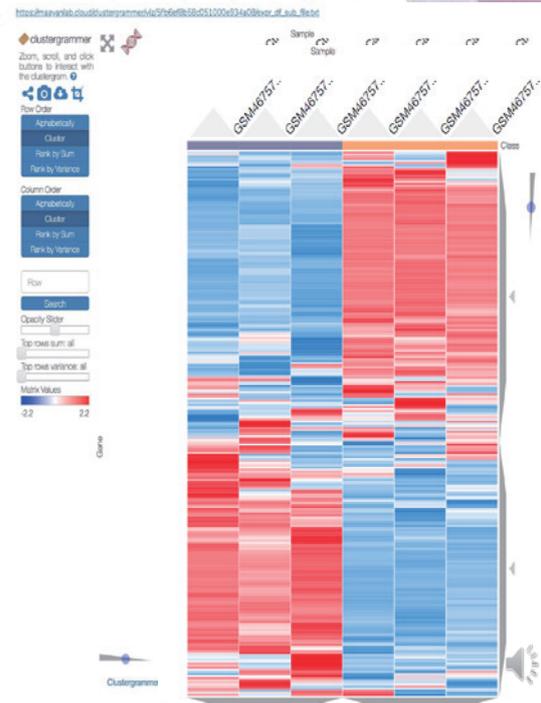


Figure 2. Clustered heatmap plot. The figure contains an interactive heatmap displaying gene expression for each sample in the RNA-seq dataset. Every row of the heatmap represents a gene, every column represents a sample, and every cell displays normalized gene expression values. The heatmap additionally features color bars beside each column which represent prior knowledge of each sample, such as the tissue of origin or experimental treatment.

# Results from Bulk RNA-seq Analysis Appyter

Ctrl vs. Case Signature | Volcano Plot

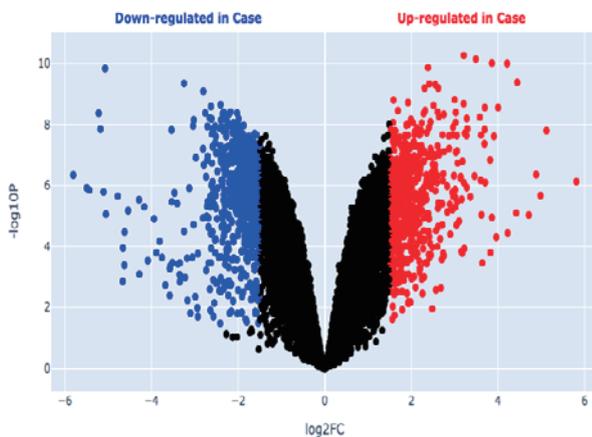
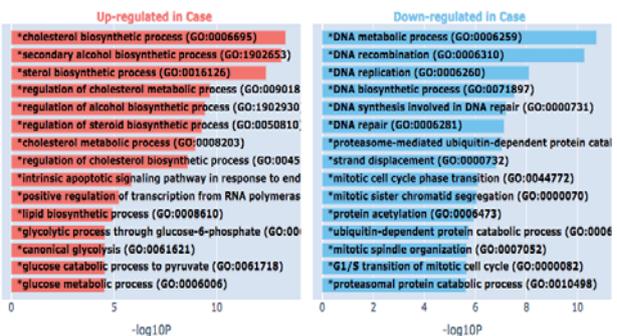
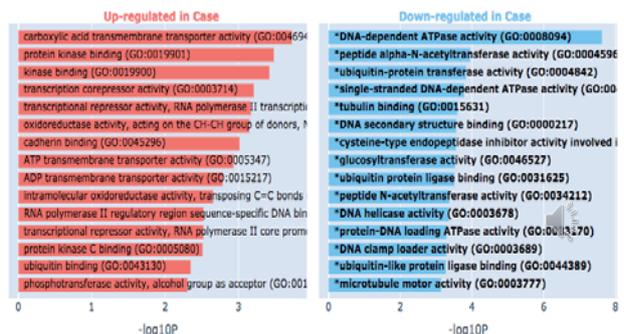


Figure 4. Volcano plot for Ctrl vs. Case. The figure contains an interactive scatter plot which displays the  $\log_2$ -fold changes and statistical significance of each gene calculated by performing a differential gene expression analysis. Genes with  $\log_2FC > 1.0$  and  $p$ -value  $< 0.05$  in red and genes with  $\log_2FC < -1.0$  and  $p$ -value  $< 0.05$  in blue. Additional information for each gene is available by hovering over it.

Ctrl vs. Case | Gene Ontology Biological Process (2018 version)

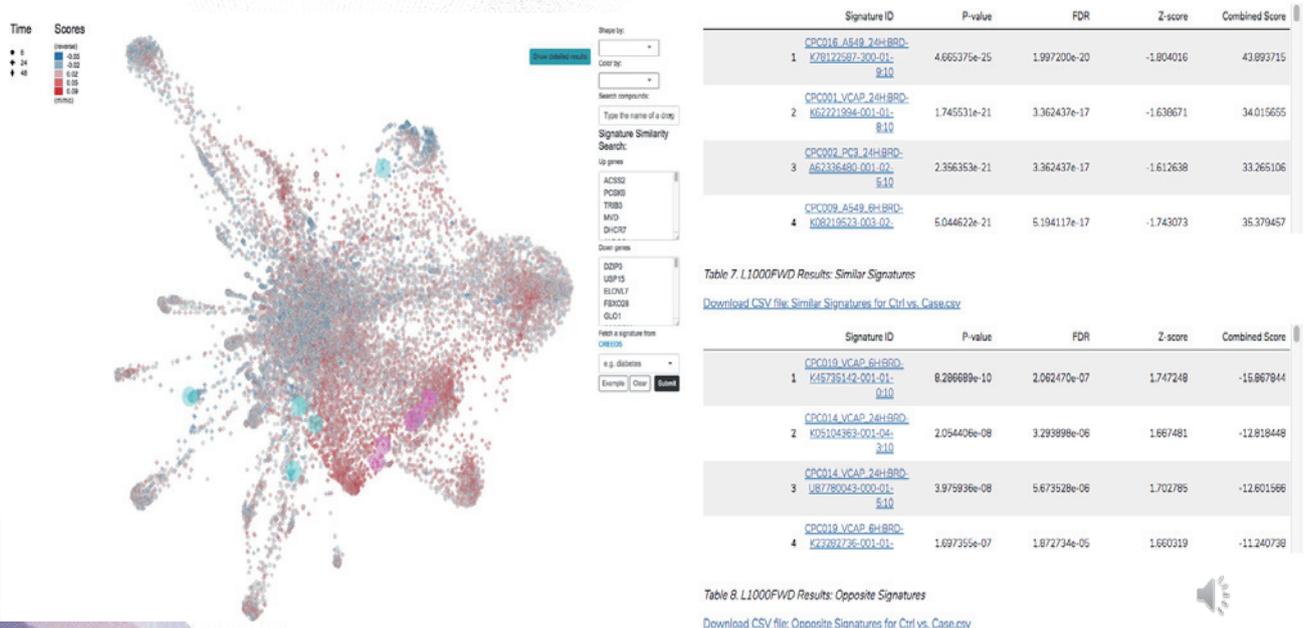


Ctrl vs. Case | Gene Ontology Molecular Function (2018 version)



URL: <https://appytters.maayanlab.cloud/#/Bulk RNA seq>

# Results from Bulk RNA-seq Analysis Appyter



URL: [https://appyters.maayanlab.cloud/#/Bulk\\_RNA\\_seq](https://appyters.maayanlab.cloud/#/Bulk_RNA_seq)

## 세포주 약물 반응 데이터의 응용



Data and text mining

## ReSimNet: drug response similarity prediction using Siamese neural networks

Minji Jeon<sup>1,†</sup>, Donghyeon Park<sup>1,†</sup>, Jinhyuk Lee<sup>1</sup>, Hwisang Jeon<sup>2</sup>,  
Miyoung Ko<sup>1</sup>, Sunkyu Kim<sup>1</sup>, Yonghwa Choi<sup>1</sup>, Aik-Choon Tan<sup>3</sup> and  
Jaewoo Kang<sup>1,2,\*</sup>

<sup>1</sup>Department of Computer Science and Engineering, <sup>2</sup>Interdisciplinary Graduate Program in Bioinformatics, Korea University, Seoul 02841, South Korea and <sup>3</sup>Division of Medical Oncology, Department of Medicine, Translational Bioinformatics and Cancer Systems Biology Laboratory, University of Colorado Anschutz Medical Campus, Aurora, CO 12801, USA

\*To whom correspondence should be addressed.

<sup>†</sup>The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

Associate Editor: Alfonso Valencia

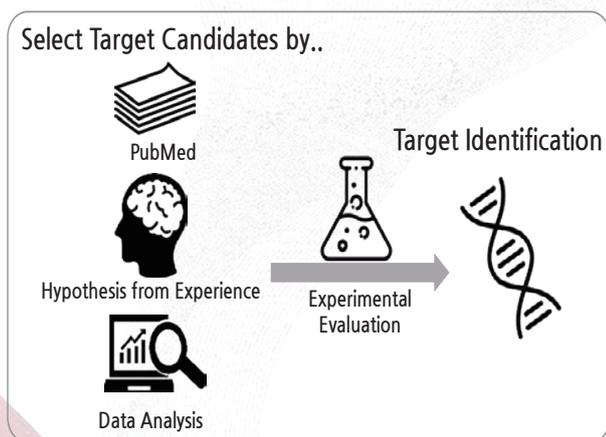
Received on July 10, 2018; revised on April 2, 2019; editorial decision on May 11, 2019; accepted on May 16, 2019



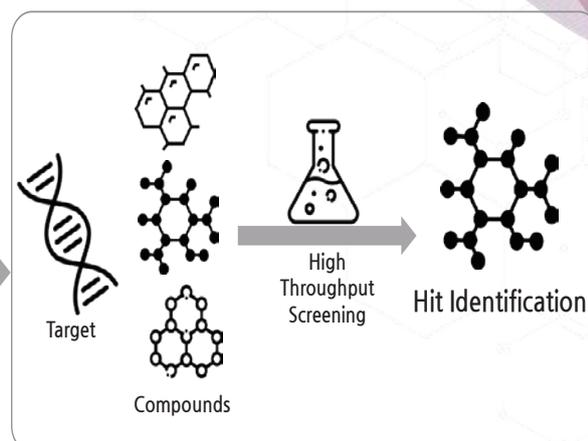
56

## Traditional Drug Discovery Pipeline: Ligand-based Drug Discovery

### Target Discovery

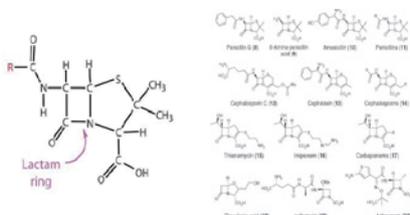


### Drug Candidate Discovery



## Introduction

- Ligand-based drug discovery
  - To find common structures among prototype drugs that bind to a desired target
  - To design structural analogs of the prototype drugs, and evaluate them



- **Target discovery** is essential
- Unable to work for **undruggable targets**
- **Limited range of drug candidates**: excludes drug candidates that may have similar effects despite binding to another protein

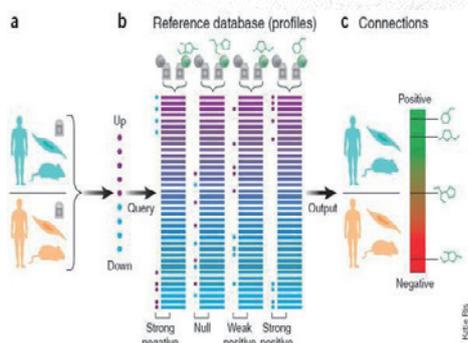
Developed a drug discovery model that can find functional analogs



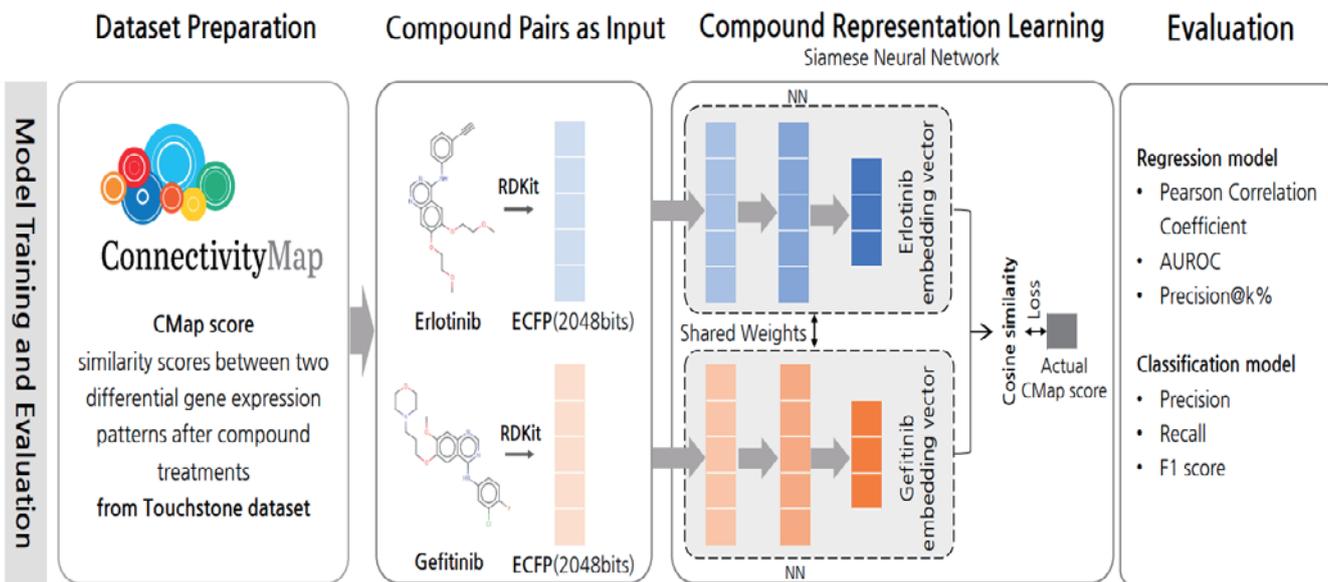
[https://application.wiley-vch.de/books/sample/3527312579\\_c01.pdf](https://application.wiley-vch.de/books/sample/3527312579_c01.pdf)

## Gene Expression-based Drug Discovery

- ConnectivityMap (CMap) by LINCS Consortium
  - Gene expression profile-based drug repurposing platform
  - Over 3 million gene expression profiles before and after drug treatments, shRNA knockdown, CRISPR knockout, etc. and 1 million gene signatures
  - User's up- and down-regulated genes in the signature DB to find chemical compounds/shRNA/CRISPR that can mimic or reverse them
  - Provide transcriptional response-based similarity scores of compound pairs (CMap Scores)
- However, as a drug repurposing platform, it is difficult to discover novel drug candidates



# Approach



Bioinformatics, 2019

# Results

Example pairs with similar structures but low actual CMap scores

Perturbagen ID: BRD-A45333338 Compound Name: periplocymarin	Perturbagen ID: BRD-K48164639 Compound Name: cholic acid	Perturbagen ID: BRD-K24846665 Compound Name: carmoxirole	Perturbagen ID: BRD-K33289131 Compound Name: CAY 10618
CMap score: -0.818 Response similarity score using MolSift: -0.136 Structure similarity score using MolVec: -0.907		CMap score: 0.234 Response similarity score using MolSift: -0.243 Structure similarity score using MolVec: -0.907	

Example pairs with dissimilar structures but high actual CMap scores

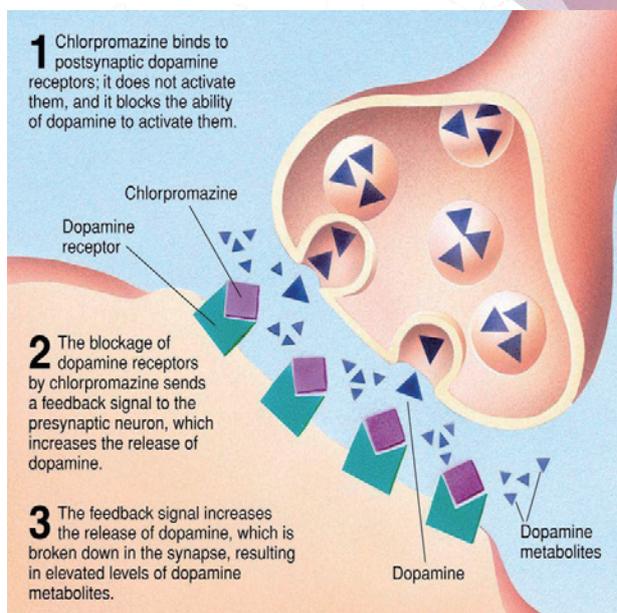
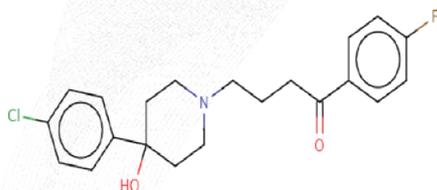
Perturbagen ID: BRD-K36877881	Perturbagen ID: BRD-K58295457				
Perturbagen ID: BRD-A41112154 CMap score: 1 Response similarity score using MolSift: 0.969 Structure similarity score using MolVec: 0.625	Perturbagen ID: BRD-K3797930 CMap score: 0.935 Response similarity score using MolSift: 0.808 Structure similarity score using MolVec: 0.382	Perturbagen ID: BRD-A87121327 CMap score: 0.989 Response similarity score using MolSift: 0.878 Structure similarity score using MolVec: 0.365	Perturbagen ID: BRD-K38963228 CMap score: 0.983 Response similarity score using MolSift: 0.907 Structure similarity score using MolVec: 0.605	Perturbagen ID: BRD-K21302415 CMap score: 0.983 Response similarity score using MolSift: 0.972 Structure similarity score using MolVec: 0.702	Perturbagen ID: BRD-K548137181 CMap score: 0.917 Response similarity score using MolSift: 0.811 Structure similarity score using MolVec: 0.293



Bioinformatics, 2019

## Drug discovery case study 1: Haloperidol

- Haloperidol
  - FDA-approved drug for schizophrenia
  - Dopamine receptor antagonist
  - Structure:



## Drug discovery case study 1: Haloperidol

**Table 2.** Top drug candidates for Haloperidol from the ZINC15 database

ZINC15 ID	ZINC15 name	Predicted similarity score by ReSimNet	Similarity score by ECFP	# of articles	Description
ZINC2516029	Chlorohaloperidol	0.995	0.884	6	Chlorohaloperidol targets the Dopamine D2 receptor <sup>a</sup>
ZINC601270	Bromperidol	0.967	0.792	66	Bromperidol is an FDA-approved drug for dementia, depression, schizophrenia, anxiety disorders and psychosomatic disorders (Yasui-Furukori <i>et al.</i> , 2002)
ZINC4214827	Amiperone	0.961	0.704	0	Amiperone targets the Dopamine D3 receptor and D3 is a potential target of Parkinson's disease and schizophrenia (Varady <i>et al.</i> , 2003)
ZINC538026	Moperone	0.955	0.792	11	Moperone is a Dopamine D2 receptor antagonist <sup>b</sup>
ZINC35851465	Cyantraniliprole	0.946	0.098	0	—
ZINC1481990	Budipine	0.942	0.172	1	Budipine is used in the treatment of Parkinson's disease (Klockgether <i>et al.</i> , 1993)
ZINC12494203	B-Hydoxychoolate	0.938	0.113	0	—
ZINC3824281	Ganaxolone	0.933	0.129	1	Ganaxolone is one of neurosteroids and is used for epilepsy (Nohria and Giller, 2007)
ZINC3812988	Butorphanol	0.933	0.200	9	Butorphanol is a neuropsychiatric agent (Iyengar <i>et al.</i> , 1987)
ZINC2041178	2,3-Dibromopropanol	0.93	0.158	0	—

<sup>a</sup><https://pubchem.ncbi.nlm.nih.gov/compound/173712#section=ChEMBL-Target-Tree>.

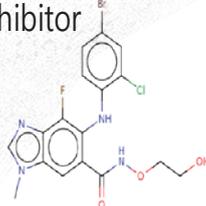
<sup>b</sup>[https://www.kegg.jp/dbget-bin/www\\_bget?D01105](https://www.kegg.jp/dbget-bin/www_bget?D01105).



# Drug discovery case study 2: Selumetinib

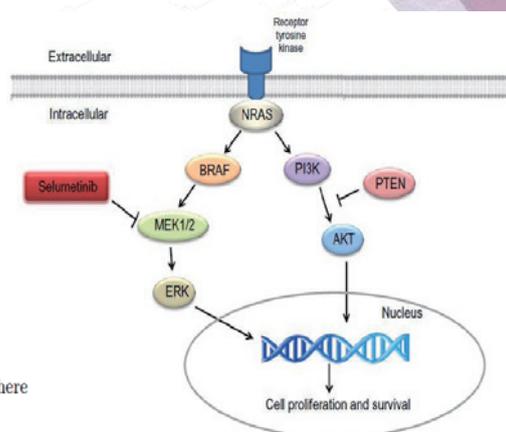
## • Selumetinib

- FDA-approved cancer drug
- MoA: MEK1/2 inhibitor
- Structure:



**Table S 11.** Top drug candidates for Selumetinib from the ZINC15 dataset. There were only two drug candidates with a score > 0.9.

ZINC15 ID	ZINC15 name	Predicted similarity score by ReSimNet	Similarity score by ECFP	# of articles	Description
ZINC38460704	binimetinib	0.993	0.841	4	Binimetinib is a MEK1/2 inhibitor [1]
ZINC43154039	buparlisib	0.937	0.084	2	Buparlisib is a PI3K inhibitor. PI3K is in a pathway parallel to the MEK pathway [2]



Bioinformatics, 2019

Jeon et al. BMC Bioinformatics (2022) 23:374  
<https://doi.org/10.1186/s12859-022-04895-5>

BMC Bioinformatics

RESEARCH

Open Access

## Transforming L1000 profiles to RNA-seq-like profiles with deep learning

Minji Jeon<sup>1,2</sup>, Zhuorui Xie<sup>1</sup>, John E. Evangelista<sup>1</sup>, Megan L. Wojciechowicz<sup>1</sup>, Daniel J. B. Clarke<sup>1</sup> and Avi Ma'ayan<sup>1\*</sup>

\*Correspondence: [avi.maayan@mssm.edu](mailto:avi.maayan@mssm.edu)

<sup>1</sup> Department of Pharmacological Sciences, Mount Sinai Center for Bioinformatics, Icahn School of Medicine at Mount Sinai, One Gustave L. Levy Place, Box 1603, New York, NY 10029, USA

<sup>2</sup> Department of Medicine, Korea University College of Medicine, Seoul, Republic of Korea

### Abstract

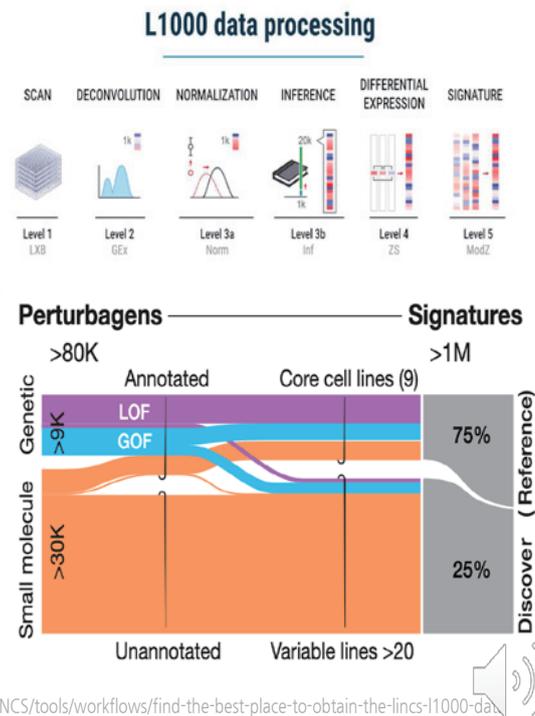
The L1000 technology, a cost-effective high-throughput transcriptomics technology, has been applied to profile a collection of human cell lines for their gene expression response to > 30,000 chemical and genetic perturbations. In total, there are currently over 3 million available L1000 profiles. Such a dataset is invaluable for the discovery of drug and target candidates and for inferring mechanisms of action for small molecules. The L1000 assay only measures the mRNA expression of 978 landmark genes while 11,350 additional genes are computationally reliably inferred. The lack of full genome coverage limits knowledge discovery for half of the human protein coding genes, and the potential for integration with other transcriptomics profiling data. Here we present a Deep Learning two-step model that transforms L1000 profiles to RNA-seq-like profiles. The input to the model are the measured 978 landmark genes while the output is a vector of 23,614 RNA-seq-like gene expression profiles. The model first transforms the landmark genes into RNA-seq-like 978 gene profiles using a modified CycleGAN model applied to unpaired data. The transformed 978 RNA-seq-like landmark genes are then extrapolated into the full genome space with a fully connected neural network model. The two-step model achieves 0.914 Pearson's correlation coefficients and 1.167 root mean square errors when tested on a published paired L1000/RNA-seq dataset produced by the LINCS and GTEx programs. The processed RNA-seq-like profiles are made available for download, signature search, and gene centric reverse search with unique case studies.

**Keywords:** L1000, RNA-seq, Gene expression translation, Generative adversarial networks



## Introduction – The LINCS L1000 Data

- L1000 assay
  - The L1000 assay only measures the mRNA expression
  - An additional 11,350 genes are computationally
- L1000 data
  - An expression profile is generated from a single point
  - Signatures (differentially expressed genes) are ca. plate



## Introduction – Limitations & Challenges

### Limitations of the L1000 Data

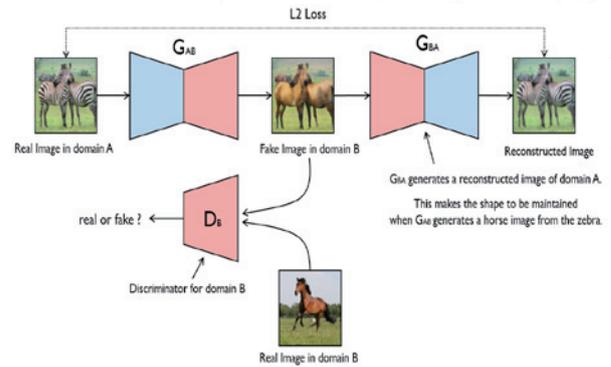
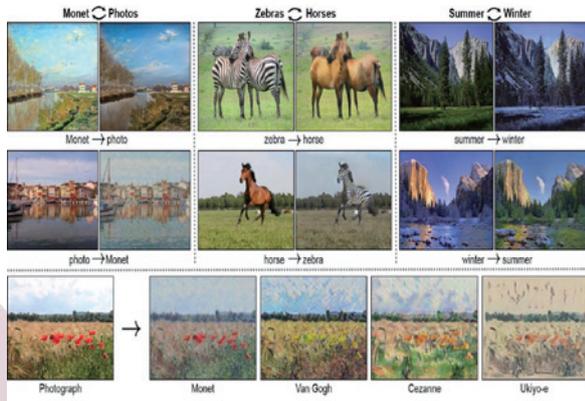
- Half of the protein-coding genes are missing from the L1000 data
  - This limits knowledge discovery about those missing genes
  - This limits the potential for integration of the CMap data with other transcriptomics profiling data

### Project Challenges

- Transform L1000 profiles to RNA-seq-like profiles at the full genome scale with Deep Learning
  - How can we train a model when there are not a lot paired L1000 and RNA-seq profiles for training?
  - How can we demonstrate that predicted RNA-seq-like profiles contain knowledge?
  - What new applications the predicted RNA-seq-like profiles can provide and how we can demonstrate these applications as use cases?
  - How the trained model compares to other published models and simpler baseline models?

## Methods – CycleGAN: Supervised Learning for Unpaired Data

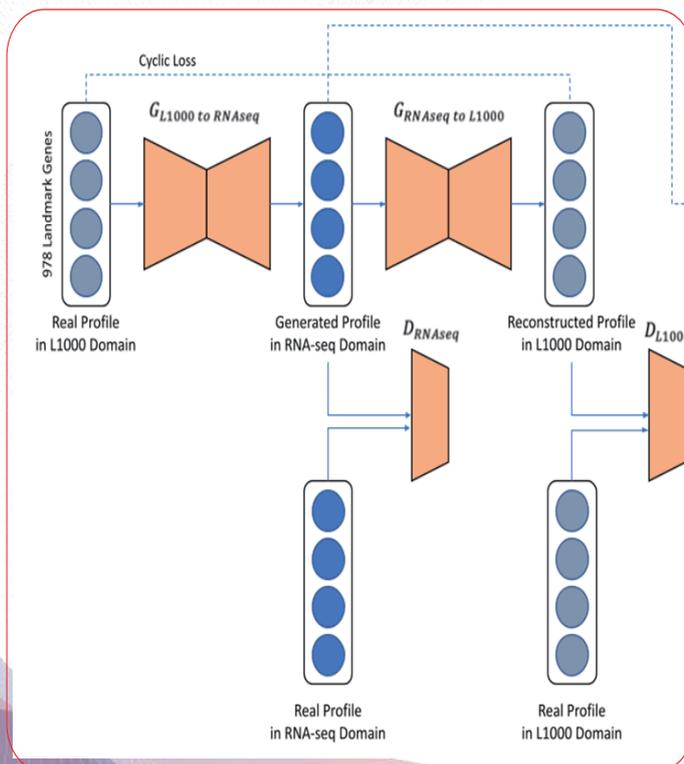
- CycleGAN (Zhu and Park et al.): An approach for learning to translate an image from a source domain X to a target domain Y in the absence of paired samples



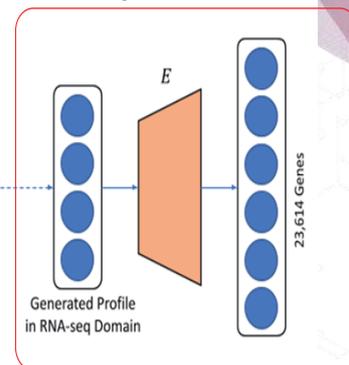
Zhu and Park et al., *ICCV 2017*

## Methods – Two-Step Model for L1000 to RNA-seq-like Profiles

Step 1



Step 2

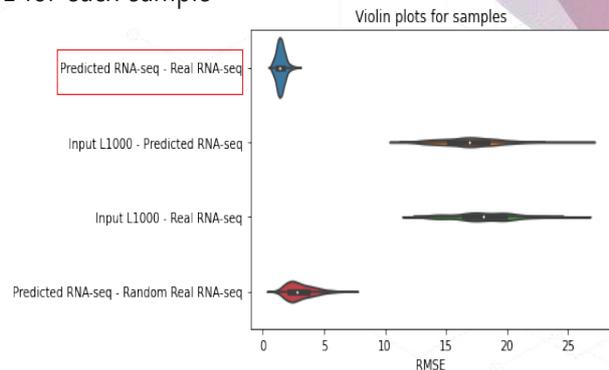
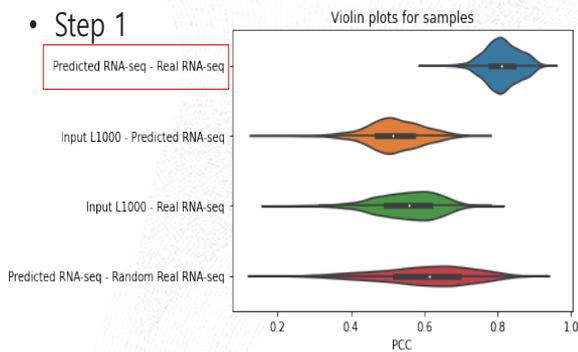


BMC Bioinformatics, 2022

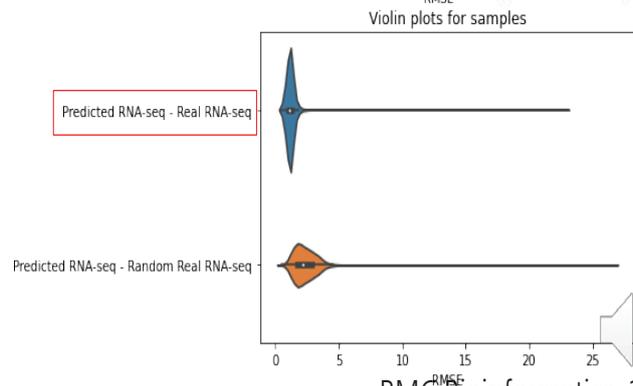
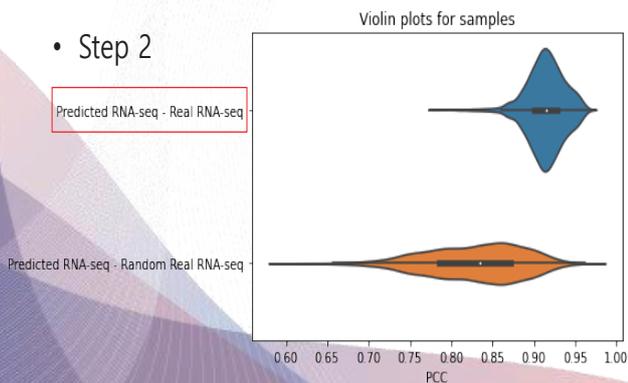
## Results - GTEx Paired Samples

- Evaluation Metrics: The average of PCC and RMSE for each sample

### Step 1



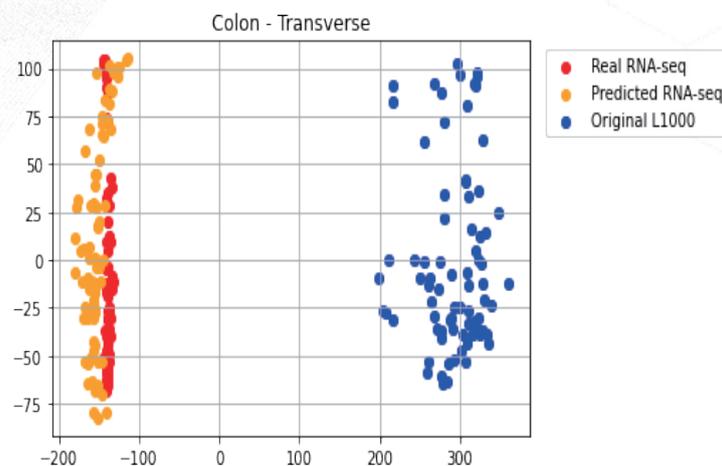
### Step 2



BMC Bioinformatics, 2022

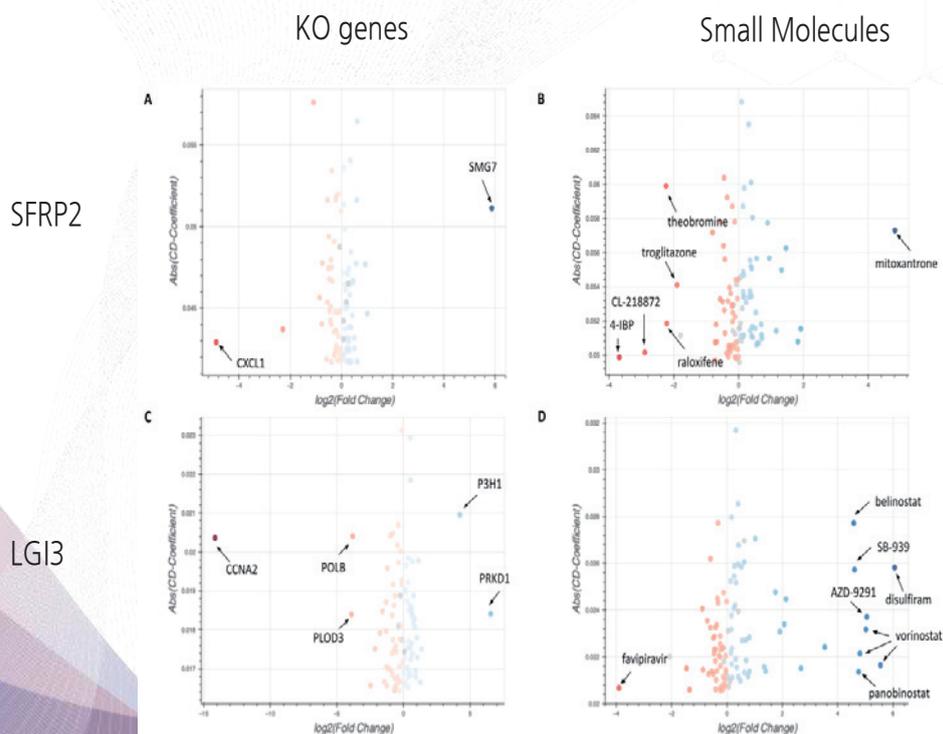
## Results - PCA of GTEx L1000-RNA-seq Paired Samples

- GTEx and LINCS profiled the same postmortem tissue samples using L1000 and RNA-seq
  - 2,929 samples from 53 tissues downloaded from GSE92743
- PCA plot of real RNA-seq profiles, predicted RNA-seq profiles, and original L1000 profiles from colon in GTEx using 11,780 common genes



BMC Bioinformatics, 2022

# Results – Gene Centric Signature Reverse Search

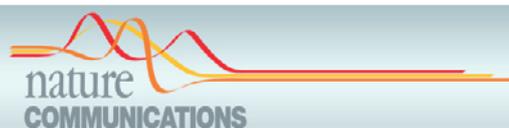


SFRP2

LGI3

[https://appytters.maayanlab.cloud/L1000\\_RNAseq\\_Gene\\_Search/](https://appytters.maayanlab.cloud/L1000_RNAseq_Gene_Search/)

BMC Bioinformatics, 2022



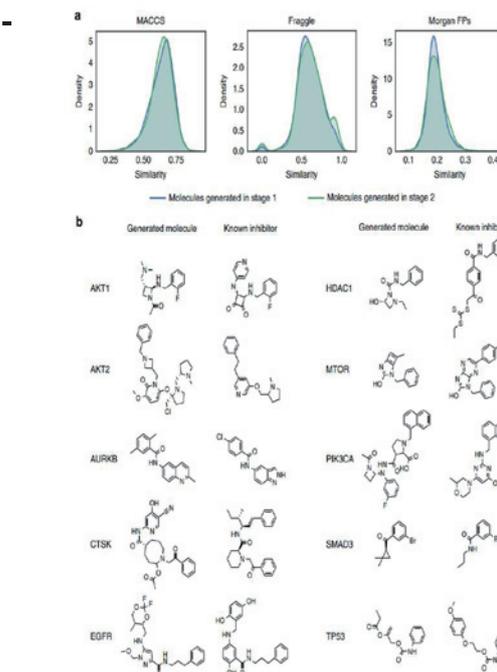
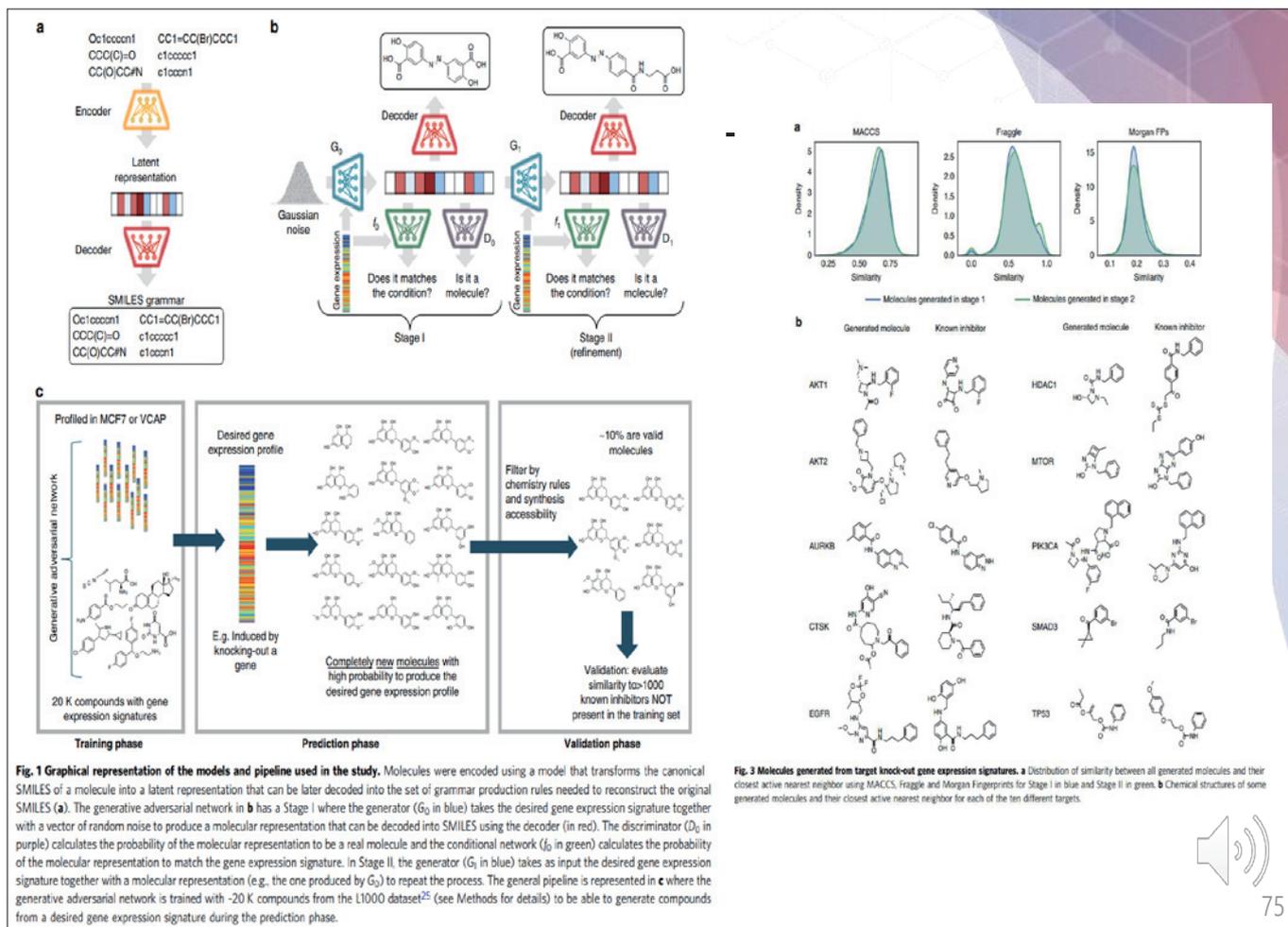
## ARTICLE

<https://doi.org/10.1038/s41467-019-13807-w> OPEN

# De novo generation of hit-like molecules from gene expression signatures using artificial intelligence

Oscar Méndez-Lucio<sup>1,2\*</sup>, Benoit Baillif<sup>1</sup>, Djork-Arné Clevert<sup>3</sup>, David Rouquié<sup>1,5\*</sup> & Joerg Wichard<sup>4,5\*</sup>

Finding new molecules with a desired biological activity is an extremely difficult task. In this context, artificial intelligence and generative models have been used for molecular de novo design and compound optimization. Herein, we report a generative model that bridges systems biology and molecular design, conditioning a generative adversarial network with transcriptomic data. By doing so, we can automatically design molecules that have a high probability to induce a desired transcriptomic profile. As long as the gene expression signature of the desired state is provided, this model is able to design active-like molecules for desired targets without any previous target annotation of the training compounds. Molecules designed by this model are more similar to active compounds than the ones identified by similarity of gene expression signatures. Overall, this method represents an alternative approach to bridge chemistry and biology in the long and difficult road of drug discovery.



75

## Predicting mechanism of action of novel compounds using compound structure and transcriptomic signature coembedding

Gwanghoon Jang<sup>1</sup>, Sungjoon Park<sup>1,\*</sup>, Sanghoon Lee<sup>1</sup>, Sunkyu Kim<sup>1</sup>, Sejeong Park<sup>1</sup> and Jaewoo Kang<sup>1,2,\*</sup>

<sup>1</sup>Department of Computer Science and Engineering, Korea University, Seoul, Republic of Korea and <sup>2</sup>Interdisciplinary Graduate Program in Bioinformatics, Korea University, Seoul, Republic of Korea

\*To whom correspondence should be addressed.

### Abstract

**Motivation:** Identifying mechanism of actions (MoA) of novel compounds is crucial in drug discovery. Careful understanding of MoA can avoid potential side effects of drug candidates. Efforts have been made to identify MoA using the transcriptomic signatures induced by compounds. However, these approaches fail to reveal MoAs in the absence of actual compound signatures.

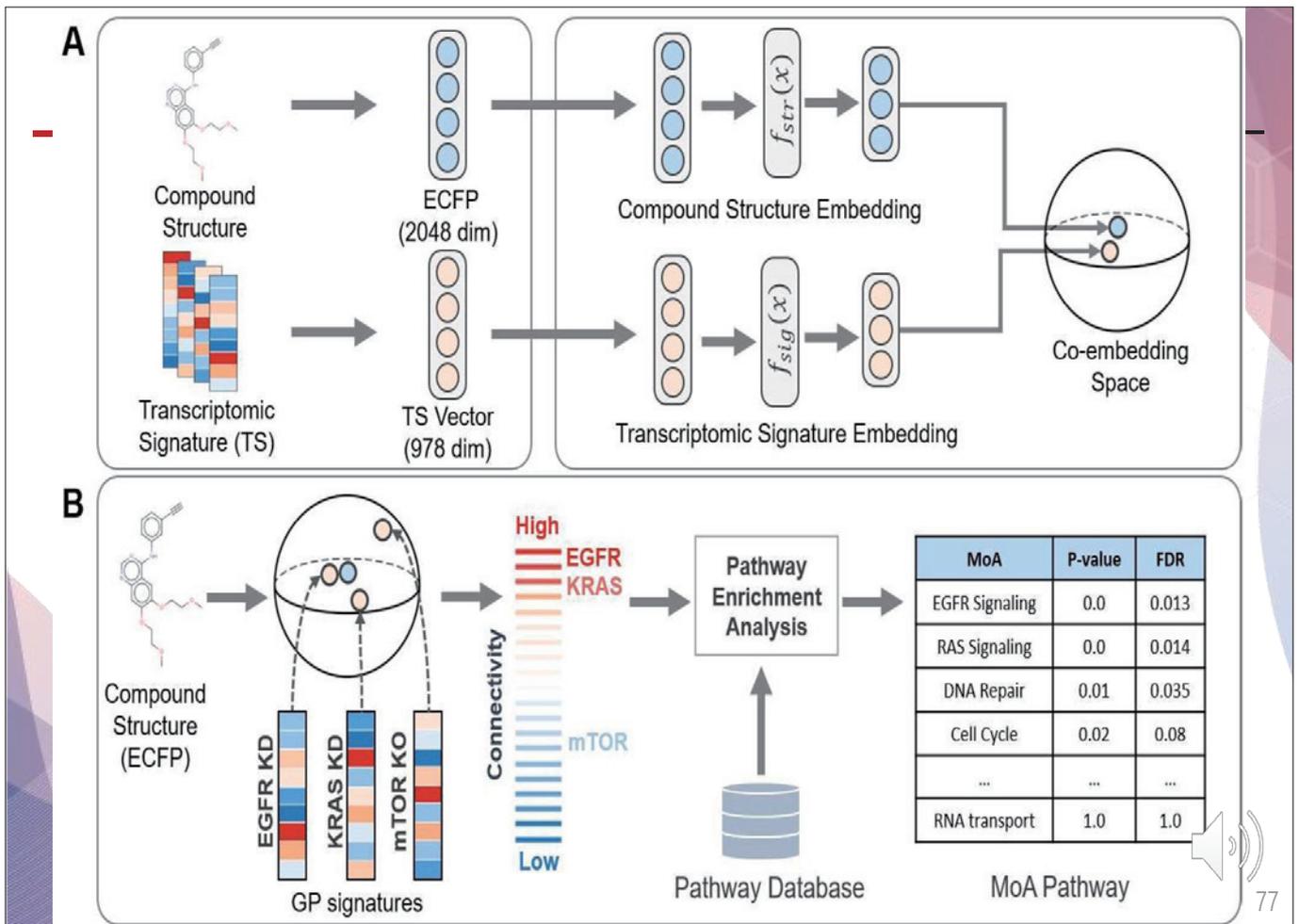
**Results:** We present MoAble, which predicts MoAs without requiring compound signatures. We train a deep learning-based coembedding model to map compound signatures and compound structure into the same embedding space. The model generates low-dimensional compound signature representation from the compound structures. To predict MoAs, pathway enrichment analysis is performed based on the connectivity between embedding vectors of compounds and those of genetic perturbation. Results show that MoAble is comparable to the methods that use actual compound signatures. We demonstrate that MoAble can be used to reveal MoAs of novel compounds without measuring compound signatures with the same prediction accuracy as that with measuring them.

**Availability and implementation:** MoAble is available at <https://github.com/dmis-lab/maoble>

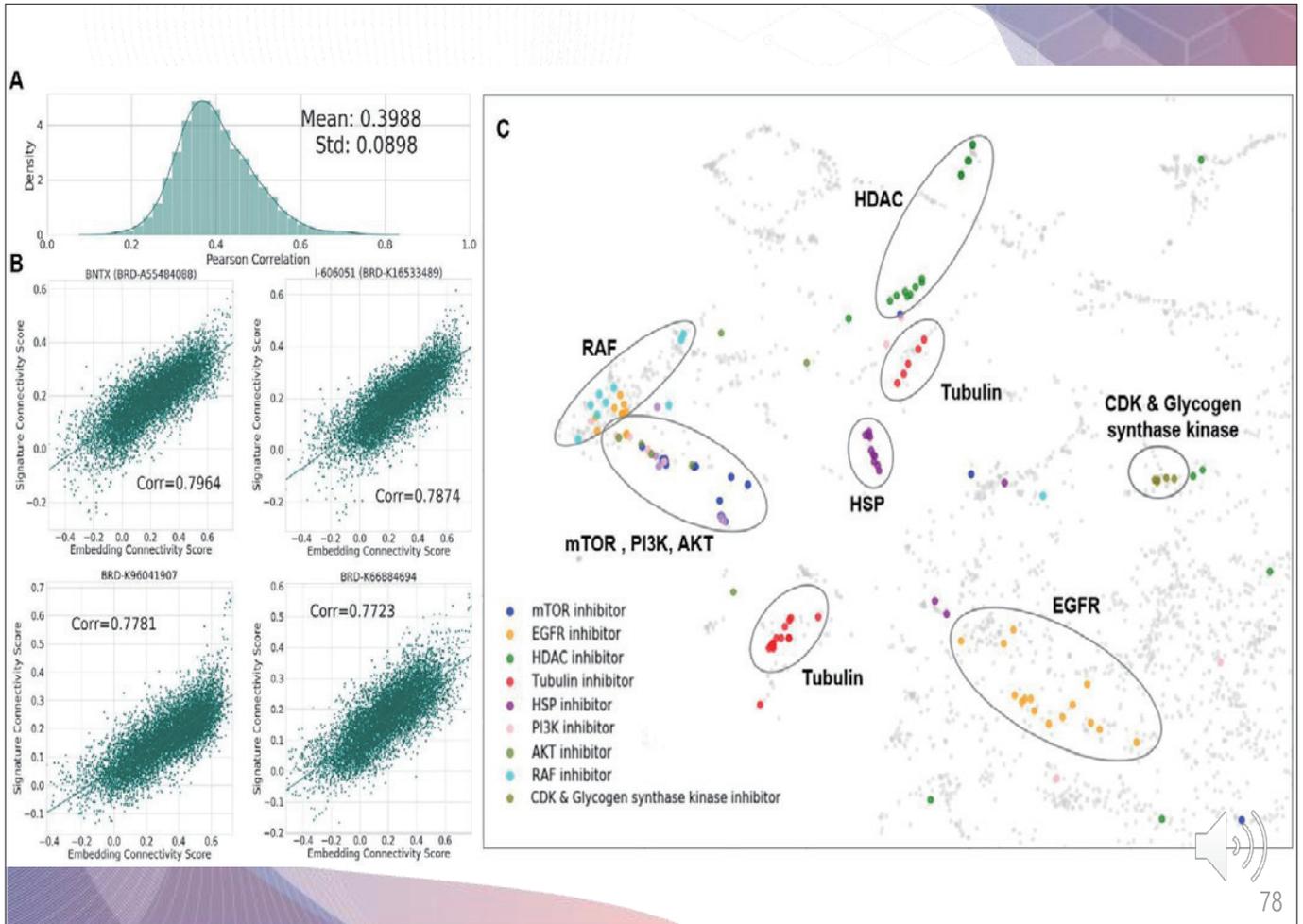
**Contact:** : sungjoonpark@korea.ac.kr or kangj@korea.ac.kr

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

76



77



78

# DeepSide: A Deep Learning Approach for Drug Side Effect Prediction

Onur Can Uner<sup>1</sup>, Halil Ibrahim Kuru<sup>1</sup>, R. Gokberk Cinbis<sup>1</sup>, Ozgur Tastan<sup>1</sup>, and A. Ercument Cicek<sup>1</sup>

**Abstract**—Drug failures due to unforeseen adverse effects at clinical trials pose health risks for the participants and lead to substantial financial losses. Side effect prediction algorithms have the potential to guide the drug design process. LINCS L1000 dataset provides a vast resource of cell line gene expression data perturbed by different drugs and creates a knowledge base for context specific features. The state-of-the-art approach that aims at using context specific information relies on only the high-quality experiments in LINCS L1000 and discards a large portion of the experiments. In this study, our goal is to boost the prediction performance by utilizing this data to its full extent. We experiment with 5 deep learning architectures. We find that a multi-modal architecture produces the best predictive performance among multi-layer perceptron-based architectures when drug chemical structure (CS), and the full set of drug perturbed gene expression profiles (GEX) are used as modalities. Overall, we observe that the CS is more informative than the GEX. A convolutional neural network-based model that uses only SMILES string representation of the drugs achieves the best results and provides 13.0% macro-AUC and 3.1% micro-AUC improvements over the state-of-the-art. We also show that the model is able to predict side effect-drug pairs that are reported in the literature but was missing in the ground truth side effect dataset. DeepSide is available at <http://github.com/OnurUner/DeepSide>.

**Index Terms**—Drug side effect prediction, deep learning, LINCS



2020 IEEE International Conference on Big Data and Smart Computing (BigComp)

## A Drug-induced Liver Injury Prediction Model using Transcriptional Response Data with Graph Neural Network

Doyeong Hwang<sup>\*†</sup>, Minji Jeon<sup>\*†</sup>, Jaewoo Kang<sup>\*†§</sup>

<sup>\*</sup>Department of Computer Science and Engineering

<sup>†</sup>Interdisciplinary Graduate Program in Bioinformatics

Korea University

Seoul, Republic of Korea

Emails: {desertbeagle, mijeon, kangj}@korea.ac.kr

<sup>†</sup>Both authors contributed equally

<sup>§</sup>Corresponding author

**Abstract**—Drug-Induced Liver Injury (DILI) is a major cause of failed drug candidates in clinical trials and withdrawal of approved drugs from the market. Therefore, machine learning-based DILI prediction can be key in increasing the success rate of drug discovery because drug candidates that are predicted to potentially induce liver injury can be rejected before clinical trials. However, existing DILI prediction models mainly focus on the chemical structures of drugs. Since we cannot determine whether a drug will cause liver injury based solely on its structure, DILI prediction based on the transcriptional effect of a drug on a cell is necessary.

Several machine learning models trained on datasets such as Liver Toxicity Knowledge Base (LTKB) [4] and Open TG-GATEs [5] were previously proposed for DILI prediction. Most of the previously proposed machine learning models are trained on drug structure information for predicting DILI [6]–[9]. However, such models used for DILI prediction do not consider genetic information or the structures and complex biological mechanisms of drugs [10]. Therefore, these models cannot predict whether a drug will cause liver injury based solely on its structure.



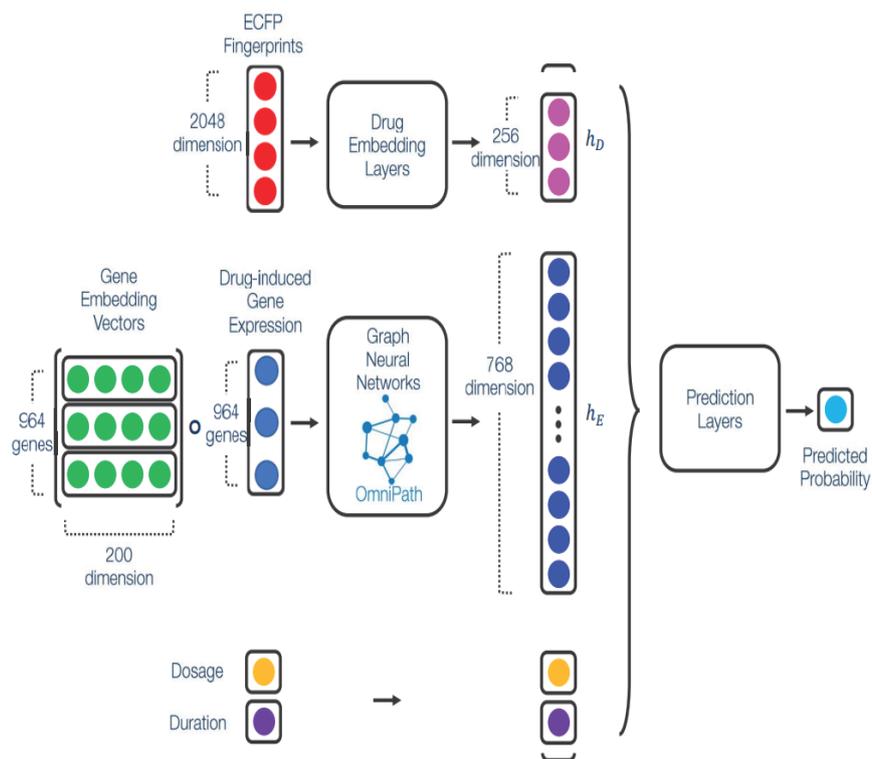


Fig. 1. The overall architecture of GLIT. The structure of a drug, the drug-induced gene expression level, and the dosage and duration of drug administration are used as inputs for GLIT. A drug embedding vector is extracted from the last layer of a neural network with three layers, and a drug-induced gene expression embedding vector is extracted from a graph neural network using a biological knowledge graph called OmniPath. The two embedding vectors are concatenated with the dosage and duration of drug administration information, and fed to the prediction layers of GLIT to predict DILI.



81

RESEARCH ARTICLE

Open Access

# DeSIDE-DDI: interpretable prediction of drug-drug interactions using drug-induced gene expressions

Eunyoung Kim and Hojung Nam\*



## Abstract

Adverse drug-drug interaction (DDI) is a major concern to polypharmacy due to its unexpected adverse side effects and must be identified at an early stage of drug discovery and development. Many computational methods have been proposed for this purpose, but most require specific types of information, or they have less concern in interpretation on underlying genes. We propose a deep learning-based framework for DDI prediction with drug-induced gene expression signatures so that the model can provide the expression level of interpretability for DDIs. The model engineers dynamic drug features using a gating mechanism that mimics the co-administration effects by imposing attention to genes. Also, each side-effect is projected into a latent space through translating embedding. As a result, the model achieved an AUC of 0.889 and an AUPR of 0.915 in unseen interaction prediction, which is competitively very accurate and outperforms other state-of-the-art methods. Furthermore, it can predict potential DDIs with new compounds not used in training. In conclusion, using drug-induced gene expression signatures followed by gating and translating embedding can increase DDI prediction accuracy while providing model interpretability. The source code is available on Git-Hub (<https://github.com/GIST-CSBL/DeSIDE-DDI>).

**Keywords:** Drug-drug interaction, Polypharmacy side effects, In silico prediction, Deep learning



82

# Predicting Cellular Responses to Novel Drug Perturbations at a Single-Cell Resolution

Leon Hetzel<sup>1,3</sup>, Simon Böhm<sup>4,3</sup>, Niki Kilbertus<sup>2,4</sup>,  
Stephan Günemann<sup>2</sup>, Mohammad Lotfollahi<sup>1,5</sup>, and Fabian Theis<sup>1,3</sup>

{leon.hetzel, simon.boehm, niki.kilbertus}@helmholtz-muenchen.de  
s.guenemann@tum.de, {mohammad.lotfollahi, fabian.theis}@helmholtz-muenchen.de

<sup>1</sup>Department of Mathematics, Technical University of Munich  
<sup>2</sup>Department of Computer Science, Technical University of Munich  
<sup>3</sup>Helmholtz Center for Computational Health, Munich  
<sup>4</sup>Helmholtz AI, Munich  
<sup>5</sup> Wellcome Sanger Institute, Cambridge

## Abstract

Single-cell transcriptomics enabled the study of cellular heterogeneity in response to perturbations at the resolution of individual cells. However, scaling high-throughput screens (HTS) to measure cellular responses for many drugs remains a challenge due to technical limitations and, more importantly, the cost of such multiplexed experiments. Thus, transferring information from routinely performed bulk RNA HTS is required to enrich single-cell data meaningfully. We introduce chemCPA, a new encoder-decoder architecture to study the perturbational effects of unseen drugs. We combine the model with an architecture surgery for transfer learning and demonstrate how training on existing bulk RNA HTS datasets can improve generalisation performance. Better generalisation reduces the need for extensive and costly screens at single-cell resolution. We envision that our proposed method will facilitate more efficient experiment designs through its ability to generate in-silico hypotheses, ultimately accelerating drug discovery.



84

(1) Encoder-Decoder: (2) Attribute embeddings: (3) Adversarial classifiers:

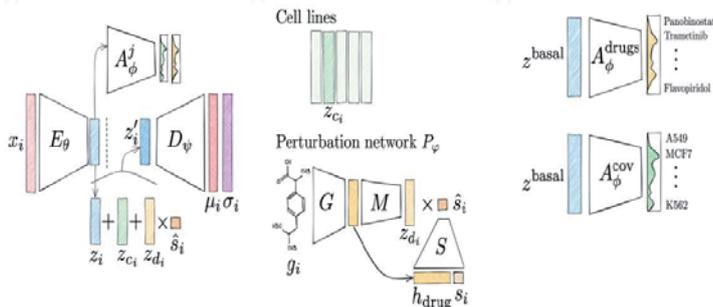


Figure 1: Architecture of chemCPA. The model consists of three parts: (1) the encoder-decoder architecture, (2) the attribute embeddings, and (3) the adversarial classifiers. The molecule encoder  $G$  can be any graph- or language-based model as long as it generates fixed-sized embeddings  $h_{\text{drugs}}$ . The MLPs  $S$  and  $M$  are trained to map the embeddings to the perturbational latent space. There,  $z_{d_i}$  is added to the basal state  $z_i$  and the covariate embedding  $z_{c_i}$ . In this work, the latter always corresponds to cell lines. The basal state  $z_i = E_\theta(x_i)$  is trained to be invariant through adversarial classifiers  $A_\phi^j$  and the decoder  $D_\psi$  gives rise to the Gaussian likelihood  $\mathcal{N}(x_i | \mu_i, \sigma_i)$ .

Table 1: Comparison of multiple models on their performance on generalisation to unseen drug-covariate combinations for dosage values of  $1 \mu\text{M}$  and  $10 \mu\text{M}$ .

Dose	Model	$\mathbb{E}[r^2]$ all	$\mathbb{E}[r^2]$ DEGs	Median $r^2$ all	Median $r^2$ DEGs
$1 \mu\text{M}$	Baseline	0.69	0.51	0.82	0.62
	scGen	0.73	0.59	0.77	0.68
	CPA	0.72	0.54	<b>0.86</b>	0.67
	chemCPA	0.74	0.60	<b>0.86</b>	0.66
	chemCPA pretrained	<b>0.77</b>	<b>0.68</b>	0.85	<b>0.76</b>
$10 \mu\text{M}$	Baseline	0.50	0.29	0.48	0.12
	scGen	0.62	0.47	0.66	0.49
	CPA	0.54	0.34	0.52	0.26
	chemCPA	0.71	0.58	0.77	0.64
	chemCPA pretrained	<b>0.76</b>	<b>0.68</b>	<b>0.82</b>	<b>0.79</b>



감사합니다  
[mjjeon@korea.ac.kr](mailto:mjjeon@korea.ac.kr)  
<https://medai.korea.ac.kr>



## 2강 실습 링크1

[https://colab.research.google.com/drive/1dryAvI-OyQ\\_XoodhcopkzAM0syb4F5AC?usp=sharing](https://colab.research.google.com/drive/1dryAvI-OyQ_XoodhcopkzAM0syb4F5AC?usp=sharing)

## 2강 실습 링크2

<https://colab.research.google.com/drive/1YBuRAZ5TwcJ4Lq5zeUOG9j50SMSibMM8?usp=sharing>

## 2강 실습 데이터 다운로드 링크

<https://drive.google.com/file/d/1Olc2FvdGIJC2G1SwHjb5DstD3ozlB11b/view?usp=shari>

# KSBI-BIML 2024

## Introduction to ConnectivityMap

고려대학교 의과대학 전민지



### DESeq2

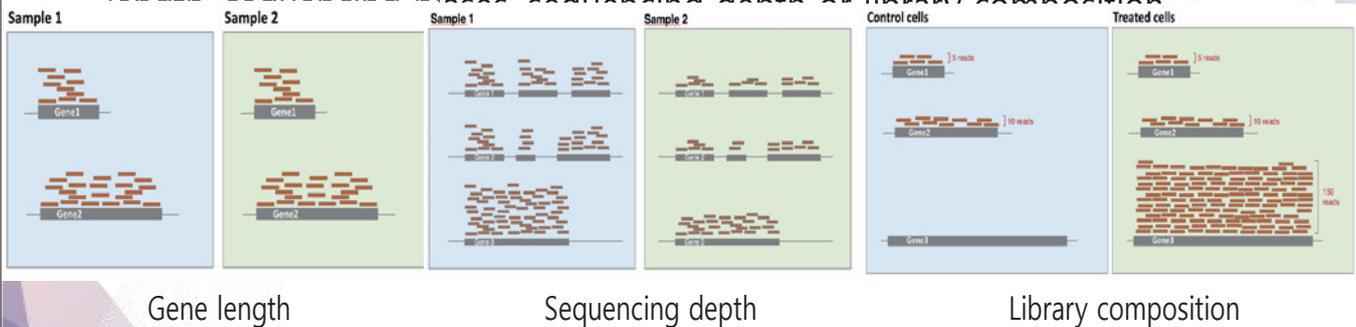
---

- Models estimate parameters (mean and dispersion) to describe count data distribution accurately.
- Dispersion parameter  $\alpha$  affects variance; edgeR and DESeq2 employ similar empirical Bayes methods to shrink  $\alpha$  towards similar-gene dispersions, enhancing differential expression test results.
- DESeq2 workflow:
  - Normalize read counts by computing size factors, addressing differences in library sizes and library compositions.
  - Calculate dispersion estimate for each gene.
  - Plot dispersion estimates of genes against mean normalized counts, and fit a line.
  - Shrunk dispersion values of each gene towards the fitted line.
  - A Generalized Linear Model accounts for confounding variables and negative binomial distribution, fitting count data.
  - For a contrast (e.g., drug-A treated vs. untreated), differential expression test assesses log fold change of normalized gene counts.
  - P-values are adjusted for multiple testing.



# DESeq2: Normalization

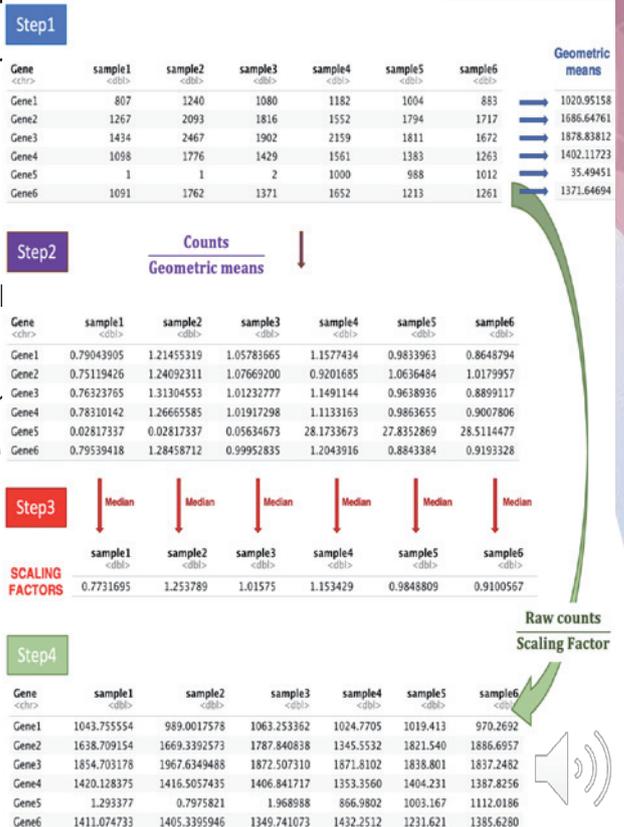
- DESeq2 expects as an input a matrix of raw counts (un-normalized counts).
- These counts are supposed to reflect gene abundance (what we are interested in)
- But they are also dependent on other less interesting factors such as gene length, sequencing biases, sequencing depth, or library composition



# DESeq2: Normalization

- Step 1: DESeq2 creates a pseudo-reference wise geometric mean (for each gene).
- Step 2: For every gene in every sample, ratios are calculated.
- Step 3: The median value of all ratios for scale factor for that sample.
- Step 4: Normalized counts can be obtained values in a given sample by that sample

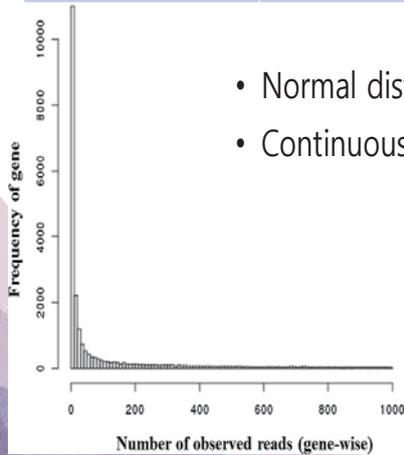
$$\hat{s}_j = \text{median}_i \frac{k_{ij}}{\left(\prod_{v=1}^m k_{iv}\right)^{1/m}}$$



## DESeq2: Count Modeling

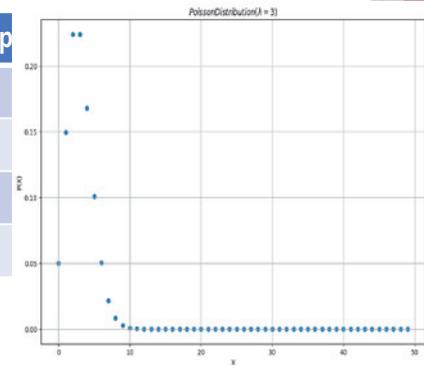
Count matrix K

	Sample1	Sample2	Samp
Gene1			
Gene2			
Gene3			



- Normal distribution (X)
- Continuous (X)

Fit Poisson distribution



- Number of cases is large
- Probability of an event happening is low
- Selecting mRNA from a large number of mRNA read
- mean = variance



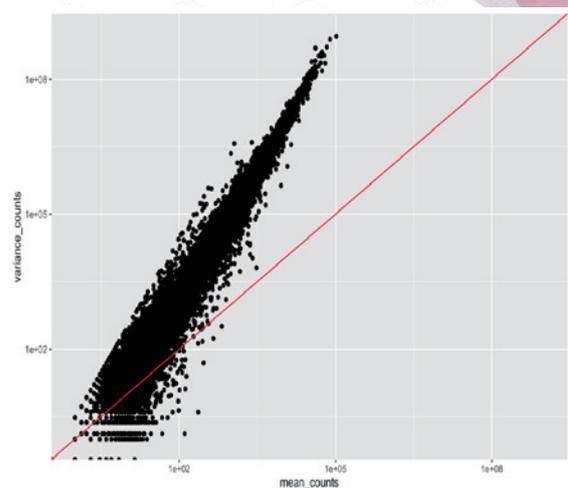
## DESeq2: Count Modeling

- But RNA-seq data does not fit Poisson distribution
  - variance  $\neq$  mean
- Negative Binomial distribution
  - variance  $>$  mean

$$K_{ij} \sim NB(\text{mean} = \mu_{ij}, \text{dispersion} = \alpha_i)$$

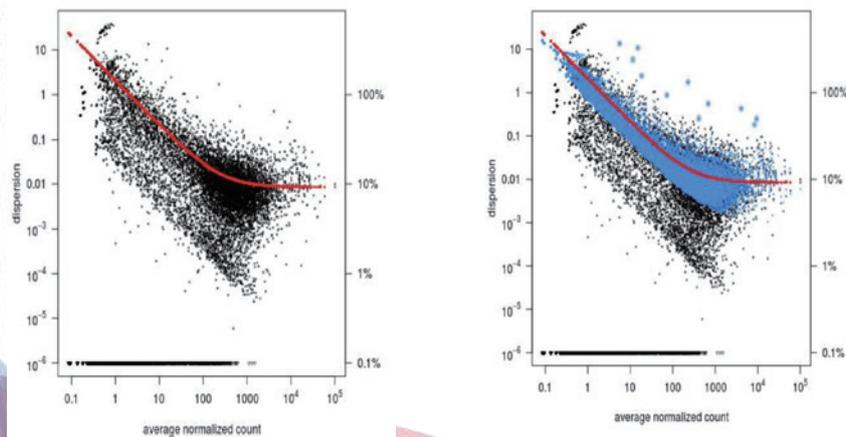
$$VAR(K_{ij}) = \mu_{ij} + \alpha_i \cdot \mu_{ij}^2$$

$K_{ij}$  = count of gene  $i$  for sample  $j$



## DESeq2: Count Modeling

- Dispersion estimation
  - Dispersion: When comparing gene expression levels between groups, it is important to know also its within-group variability
  - RNA-seq experiments typically have only few replicates
  - It is difficult to estimate within-group variability
  - Solution: pool information across genes with are expressed at similar level
    - assumes that genes with similar average expression strength have similar dispersion



## DESeq2: Generalized linear model

- Generalized linear model:

$$\log_2(q_{ij}) = \beta_0 + \beta_1 \cdot x_j + \epsilon$$

$$q_{ij} = \frac{\mu_{ij}}{\text{SizeFactor}_j}$$

$\beta_0$  is the log<sub>2</sub> expression level in the reference (control samples)

$\beta_1$  is the log<sub>2</sub>FC between treated and control cells

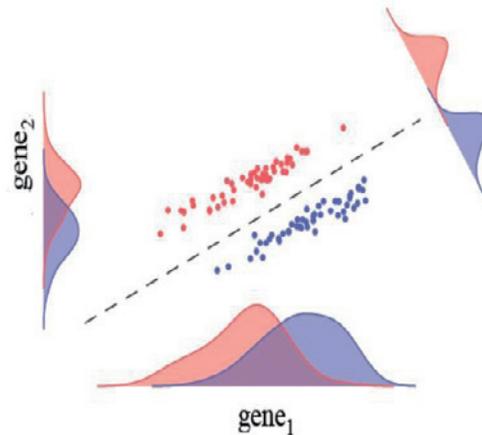
$x_j = 0$  if sample  $j$  is the control sample

$x_j = 1$  if sample  $j$  is the treated sample



## Characteristic Direction

- Characteristic Direction (Clark et al. 2014, BMC bioinformatics)
  - Genes do not function in isolation but as part of a complex network of interactions
    - This leads to significant correlations
  - Univariate approaches can miss some structure in the data
  - Multivariate approaches are sensitive to the curse of dimensionality

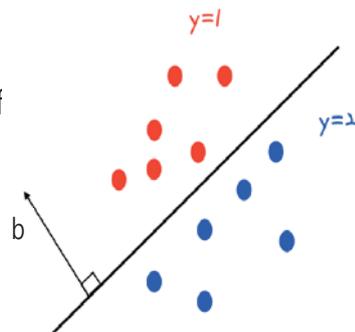


## Characteristic Direction

- Linear discriminant analysis
  - Bayes rules for the classification probability
  - The contribution of each  $b$  can be interpreted as quantifying the relative contribution of each component to the total differential expression giving the significance of the corresponding gene

Normal vector = direction of characteristics in gene expression data

$$\sum_{i=1}^p \hat{b}_i^2 \equiv 1$$



Linear classification boundary

