

# 1 Mining influential genes based on deep learning

2 **Lingpeng Kong<sup>1</sup>, Yuanyuan Chen<sup>2</sup>, Cong Pian<sup>2</sup>, Mingmin Xu<sup>1</sup>, Zutan Li<sup>1</sup>, Jingya Fang<sup>1</sup>, Liangyun**  
3 **Zhang<sup>2\*</sup>**

4 1 College of Agriculture, Nanjing Agricultural University, Jiangsu, Nanjing 210095, China

5 2 Department of Mathematics, College of Science, Nanjing Agricultural University, Nanjing 210095,  
6 China

7 \* To whom correspondence should be addressed: zlyun@njau.edu.cn

## 8 Abstract

9 **Background:** Currently, large-scale gene expression profiling has been successfully applied to the  
10 discovery of functional connections among diseases, genetic perturbation, and drug action. To address  
11 the cost of an ever-expanding gene expression profile, a new, low-cost, high-throughput reduced  
12 representation expression profiling method called L1000 was proposed, with which one million profiles  
13 were produced. Although a set of ~1,000 carefully chosen landmark genes that can capture ~80% of  
14 information from the whole genome has been identified for use in L1000, the robustness of using these  
15 landmark genes to infer target genes is not satisfactory. Therefore, more efficient computational methods  
16 are still needed to deep mine the influential genes in the genome.

17 **Results:** Here, we propose a computational framework based on deep learning to mine a subset of genes  
18 that can cover more genomic information. Specifically, an AutoEncoder framework is first constructed  
19 to learn the non-linear relationship between genes, and then DeepLIFT is applied to calculate gene  
20 importance scores. Using this data-driven approach, we have re-obtained a landmark gene set. The result  
21 shows that our landmark genes can predict target genes more accurately and robustly than that of L1000  
22 based on two metrics (mean absolute error (MAE) and Pearson correlation coefficient (PCC)). This  
23 reveals that the landmark genes detected by our method contain more genomic information.

24 **Conclusions:** We believe that our proposed framework is very suitable for the analysis of biological big  
25 data to reveal the mysteries of life. Furthermore, the landmark genes inferred from this study can be used

26 for the explosive amplification of gene expression profiles to facilitate research into functional  
27 connections.

28 **Keywords:** landmark genes, deep learning, AutoEncoder, DeepLIFT

29

## 30 **Background**

31 One of the fundamental challenges that has emerged throughout biomedicine is the need to establish  
32 relationships between disease, physiological processes and the role of small molecule therapies. To  
33 address this problem, a genomic signature is required that should have sufficiently high complexity to  
34 provide a rich description for all biological states, including those that are physiological, related to  
35 disease, or induced with a chemical, and that should be generated in a low-cost and high-throughput way.  
36 Gene expression profiling has been widely applied in medicine and biology to elucidate the response  
37 mechanism of cells to diseases, genetic interference and drug therapy[1, 2]; using this technique, the  
38 Connectivity Map (CMap) project has been proposed a systematic approach to discover functional  
39 connections among diseases, genetic perturbation, and drug action. Meanwhile, this study also suggested  
40 the value of a large-scale community CMap project[3].

41 Higher requirements have been put forward for the scale of the CMap project, and a diversity of chemical  
42 perturbations, genetic perturbations, and cell types await to be characterized. Unfortunately, although the  
43 price of commercial gene expression microarrays has been decreasing steadily, the high cost of profiling  
44 thousands of samples makes this prospect difficult. Therefore, how to reduce the cost of acquiring gene  
45 expression profiles is the first problem to be solved.

46 Previous studies have shown that although there are a large number of genes in the genome, most of their  
47 expression patterns are highly correlated[4, 5]. Cluster analysis of single-cell RNA-Seq indicated that  
48 genes from the same cluster showed similar expression patterns under different conditions[6]. Given such  
49 high similarity, researchers from the Library of Integrated Network-Based Cellular Signatures (LINCS)  
50 program hypothesized that it is possible to capture any cellular state at a low cost by measuring a reduced

51 representation of the transcriptome[7]. Using Affymetrix HG-U133A microarray data from the Gene  
52 Expression Omnibus (GEO)[8], these researchers applied an iterative peel-off procedure based cluster  
53 analysis to identify the subset of universally informative transcripts termed ‘landmark genes’. According  
54 to the LINCS analysis, a set of ~1,000 genes was finally identified as landmark genes, which was  
55 sufficient to recover 82% of the information in the full transcriptome. Then, the expression profile of the  
56 target genes was inferred by a linear regression algorithm, which was subsequently improved several  
57 times to improve the reliability of prediction[9, 10]. Finally, based on the ~1,000 landmark genes, a new,  
58 low-cost, high-throughput reduced representation expression profiling method called L1000 was  
59 proposed, with which one million profiles were reported for the first time.

60 Cluster analysis mostly measures the similarity between variables by linear distance, such as Euclidean  
61 distance. As nonlinear regulatory relationships between genes are very common in biology[11], it is  
62 difficult for the ~1,000 landmark genes inferred by cluster analysis to fully represent genomic  
63 information. Therefore, a new computational method with the capacity to capture the non-linear  
64 relationships of genes is needed to re-mine the influential genes that cover more information about the  
65 genome.

66 Deep learning, a non-linear network structure using multi-layer non-linear functions, has recently  
67 emerged based on big data, and academic interest has increased rapidly since the early 2000s[12].  
68 Furthermore, the recent success of deep learning in diverse fields such as image and speech  
69 recognition[13, 14], natural language processing[15, 16], and bioinformatics[17, 18] suggests its ability  
70 to learn hierarchical nonlinear patterns on large data sets. Deep learning can be divided into supervised  
71 learning and unsupervised learning. The former mainly includes deep neural network (DNN),  
72 convolutional neural network (CNN) and recurrent neural network (RNN) and is mainly used for  
73 classification tasks such as transcription factor binding site prediction[19], promoter prediction[20] and  
74 predicting the effects of noncoding variants[21]. The most representative of the latter is AutoEncoder,  
75 which is commonly used for dimension reduction[22] to analyse high-dimensional gene expression  
76 data[23, 24] and to integrate heterogeneous data[25, 26]. As a non-linear feature extraction method,  
77 AutoEncoder is capable of learning more useful features than linear feature extraction methods, such as  
78 principal component analysis (PCA).

79 Despite deep neural networks become increasingly popular, there is still a "black box" nature that hinders  
80 their application when interpretability is paramount. Understanding how an input feature affects a  
81 particular input can lead to new scientific discoveries. Therefore, multiple studies have been conducted  
82 to explain this "black box"[27-29]. Similarly, DeepLIFT is an efficient and effective method for  
83 computing importance scores in a neural network by comparing the activation of each neuron to a  
84 reference activation[30]. This method has been successfully applied to visualize splice site-related motifs  
85 from a trained CNN model[31].

86 Here, we present a deep learning framework to mine a gene set that can cover more genomic information.  
87 Specifically, we first constructed an AutoEncoder framework using ~130,000 gene expression profiles  
88 from the GEO Affymetrix microarray platform for training to learn the complex regulatory relationships  
89 across genes. Using this model, ~22,000 dimensional expression data were reduced to only 100.  
90 Clustering analysis of lung cancer showed that these 100 dimensional features well represent the  
91 biological information of gene expression data. Then, DeepLIFT was applied to measure the impact of  
92 each input layer neuron on the bottleneck layer neurons by providing an importance score. Using this  
93 data-driven approach, we obtained a list of genes that were sorted based on the importance score. By  
94 extracting genes from top to bottom, a new landmark gene set with the same number of genes as the  
95 original set from L1000 was finally identified. To compare the two landmark gene sets, we next used D-  
96 GEX[10] as a prediction model to infer the expression profiles of the target genes (besides the landmark  
97 genes) based on the landmark genes. The result shows that our landmark gene set can predict target  
98 genes more accurately and reliably than that of L1000 by comparing two performance metrics, MAE and  
99 PCC. Therefore, the landmark genes inferred by our method truly contain more information about the  
100 genome and are more suitable for expanding the scale of the CMap project.

## 101 **Results and discussion**

### 102 **Performance evaluation of the AutoEncoder model**

103 After training the AutoEncoder model with GEO-based training samples (99909), we use reserved test  
104 samples (11100) to evaluate its predictive power in both gene and sample dimensions. In terms of genes,  
105 we use MAE and PCC to measure the prediction error and similarity of each gene. As shown in Figure

106 1A, the average MAE and PCC of all genes are 0.2222 and 0.7627, respectively, and the permutation test  
107 shows that there is a significant high similarity between the predicted value and the real value of almost  
108 all genes (21696/22268). In terms of samples, we collect 237 lung cancer samples from the GEO database  
109 as new test samples, including 49 normal samples, 58 lung adenocarcinoma (ADC) samples and 130  
110 lung squamous cell carcinoma (SCC) samples. Then, we take the expression profiles of these samples as  
111 the input of the trained AutoEncoder and use the output of the bottleneck layer to cluster the samples.  
112 Figure 1B shows that the low dimensional space mapped by the trained AutoEncoder well retains the  
113 biological information of the samples. All of these results show that our trained AutoEncoder can learn  
114 the non-linear relationships between genes well.

115 **Figure 1.** Performance evaluation of the AutoEncoder model in both gene (A) and sample dimensions (B). (A) The  
116 density plots of the predictive error (MAE) and the similarity (PCC) of all genes. (B) The circular diagram of  
117 clustering for three types of samples, including normal (Normal), lung adenocarcinoma (ADC) and lung squamous  
118 cell carcinoma (SCC).

119

## 120 **Comparison of the landmark genes**

121 First, we analyse the degree of overlap between our landmark genes (called D1000) and the landmark  
122 genes from L1000 (called L1000) and find that only 129 genes are shared. In addition, to evaluate the  
123 performance of the landmark genes inferred by our method, we use them as input to infer the expression  
124 profile of the target genes using a deep learn-based method, D-GEX. Then, we also use the MAE and  
125 PCC of each common target gene (9163) to compare D1000 with L1000. We define MAE and PCC of  
126 the target genes inferred from L1000 and D1000 as  $MAE_{L1000}$ ,  $MAE_{D1000}$ ,  $PCC_{L1000}$ , and  $PCC_{D1000}$ ,  
127 respectively. As shown in Figure 2A and 2B, compared with  $MAE_{L1000}$  with a value of 0.1129-1.0524,  
128 the  $MAE_{D1000}$  value range is 0.0994-0.6681, and the paired t-test shows that  $MAE_{D1000}$  is significantly

129 lower than  $MAE_{L1000}$  ( $p < 0.01$ ). Similarly, as shown in Figure 2C and 2D, compared with  $PCC_{L1000}$   
130 with a value of 0.0006-0.9875, the  $PCC_{D1000}$  value range is 0.4764-0.9905, and the paired t-test shows  
131 that  $PCC_{D1000}$  is significantly higher than  $PCC_{L1000}$  ( $p < 0.01$ ). Furthermore, all  $PCC_{D1000}$  pass the  
132 permutation test, but 44 target genes fail in  $PCC_{L1000}$ . These results show that the new landmark genes  
133 inferred from our method can predict target genes more accurately and robustly than the old landmark  
134 genes.

135 **Figure 2.** The density plot (A, C) and scatter plot (B, D) are used for comparison of the landmark genes inferred  
136 from our method (labelled as “D1000”) and that of L1000 (labelled as “L1000”) in terms of MAE (A, B) and PCC  
137 (C, D). In B and D, each dot represents a predicted target gene, and the red dot indicates that D1000 is better than  
138 L1000.

139

#### 140 **Cross-platform generalization analysis of the landmark genes**

141 RNA-Seq is another high-throughput sequencing platform that has gradually become the standard for  
142 gene expression profiling. Next, to explore the ability to use landmark genes inferred from the  
143 microarray-based GEO dataset to infer target genes from the RNA-Seq-based expression profiling, we  
144 download a RNA-Seq-based gene expression profiling containing 2,921 samples from GTEx database,  
145 and the predicted target genes are analysed. The results indicate that the average MAE and PCC of all  
146 target genes are 0.4590 and 0.7790 (Figure 3), respectively, and that 92.51% of the target genes pass the  
147 permutation test, which shows that the landmark genes have excellent cross-platform generalization.

148 **Figure 4.** Cross-platform generalization analysis of the landmark genes inferred from our method.

149

#### 150 **Functional analysis of the landmark genes**

151 Finally, to analyse whether the landmark genes suggested by our data-driven approach based on the  
152 analysis of 129,158 samples are enriched in particular known biological categories, we study their  
153 molecular functions from the perspective of Gene Ontology (GO). Given that the landmark genes cover  
154 most information about the genome, we infer that the landmark genes, when considered as a set, are  
155 dominated by either very few functions or many functions.

156 To test this inference, we use the R Bioconductor package clusterProfiler (v3.10.1) to apply  
157 hypergeometric statistics between the 943 landmark genes and a database of 1,645 gene sets that come  
158 from molecular function terms compiled in Gene Ontology. As shown in Figure 4, we observe only 34  
159 functional categories, most of which tend to be basic and generic, such as "DNA binding transcription  
160 factor binding", "GDP binding", "enzyme inhibitor activity" and "protease binding", and contain only a  
161 small fraction of the landmark genes (e.g., "cell adhesion molecule binding" contains 61 of 943  
162 landmarks). The results show that no particular functional category dominates the landmark genes.

163

164 **Figure 4.** Enriched GO molecular functions term by using the landmark genes as a set.

165

## 166 **Conclusion**

167 The central dogma of molecular biology states that the flow of genetic information is "DNA to RNA to  
168 protein". Current biological studies, such as genomic studies including variable splicing and single  
169 nucleotide polymorphisms, and epigenomic studies including methylation and histone modification, are  
170 all ultimately concerned with the regulation of gene expression. Therefore, gene expression patterns can  
171 reflect almost every aspect of life activities and can be used as genomic signatures to discover the  
172 functional connections among diseases, genetic perturbation, and drug action.

173 In this study, we proposed a deep learning-based method to detect influential genes in the genome to  
174 obtain large-scale expression profiles at lower costs. In a nutshell, this is a question of feature selection.  
175 The computing framework we designed combines AutoEncoder and DeepLIFT to assess the impact of  
176 each gene in the genome, which was carried out using a data-driven approach in an unbiased manner  
177 rather than selecting transcripts based on prior biological knowledge. AutoEncoder is a nonlinear feature  
178 extraction method, which is transformed into a feature selection method by the application of DeepLIFT.  
179 The results show that using our landmark gene set can predict target genes more accurately and robustly  
180 than the gene set inferred from cluster analysis and reflects the advantages of deep learning in nonlinear  
181 computation.

182 We believe that our proposed framework is very suitable for the analysis of biological big data to reveal  
183 the mysteries of life. Furthermore, the landmark genes inferred from this study can be used for the  
184 explosive amplification of gene expression profiles to facilitate research into functional connections.

## 185 **Methods**

186 In this study, our goal is to extract ~1,000 influential genes from ~22,000 genes, which is a feature  
187 selection problem. Although many feature selection methods such as subset selection[32] and random  
188 forest[33], which are usually used in classification tasks, can effectively filter out redundant features,  
189 they cannot effectively capture the nonlinear relationship between features. In view of the above  
190 problems, we designed a computational framework as follows.

## 191 **Data sources**

192 In Table 1, three publicly available datasets are used for our analysis: the microarray-based GEO dataset,  
193 the RNA-Seq-based GTEx dataset and the lung cancer subtype dataset. The first two were downloaded  
194 from [https://cbcl.ics.uci.edu/public\\_data/D-GEX/](https://cbcl.ics.uci.edu/public_data/D-GEX/); the latter, from the GEO database.

195 First, the microarray-based GEO dataset is used to train AutoEncoder. This dataset contains 129,158 gene  
196 expression profiles, each of which contains 22,268 probes corresponding to 978 landmark genes and  
197 21,290 target genes. The original expression data are quantile normalized to a range of values between 4  
198 and 15 to remove technical variation[34]. Considering that a dataset containing a large number of  
199 redundant samples with high similarity corresponds to low statistical representativeness[35], the k-means

200 clustering program is used to remove duplicated profiles. Finally, the remaining 111,009 samples are  
201 randomly divided into ~90% (99909) for training and ~10% (11100) for testing.

202 Next, cross-platform performance can be evaluated based on the RNA-Seq dataset from GTEx, which  
203 contains 2,921 gene expression profiles of various tissue samples produced on RNA-Seq platform in the  
204 format of reads per kilobase per million (RPKM). We refer to the pre-processing protocol used in D-  
205 GEX for cross-platform data matching and joint quantile normalization. The 22,268 probes are finally  
206 matched to 10463 genes based on Gencode V12 annotations, including 943 landmark genes and 9520  
207 target genes.

208 Finally, the lung cancer subtype dataset is used to verify whether AutoEncoder can effectively learn  
209 biological information. This dataset contains 237 gene expression profiles from the GSE4573 and  
210 GSE10072 microarray datasets, including 49 normal samples, 58 lung adenocarcinoma (ADC) samples  
211 and 130 lung squamous cell carcinoma (SCC) samples.

212 **Table 1.** Three expression datasets from the GEO and GTEx databases.

Dataset	Sample size	Platform	Database
1	129,158	Microarray	GEO
2	2,921	RNA-Seq	GTEx
3	237	Microarray	GEO

213

#### 214 **A brief summary of the computational framework**

215 Our computational framework mainly consists of two parts, AutoEncoder-based and DeepLIFT-based  
216 (Figure 5). In the AutoEncoder-based part, we use ~130,000 gene expression profiles to train an  
217 AutoEncoder that is composed of two steps, encoder and decoder. However, AutoEncoder is a feature  
218 extraction method that transforms data from the original, high-dimensional space to a relatively low-  
219 dimensional space. In other words, new features are generally different from original features. Here, the  
220 encoder compresses the 22268 dimensional samples to 100 dimensions. In the DeepLIFT-based part, we  
221 use DeepLIFT to compute the importance scores of each input layer neuron on the bottleneck layer

222 neurons. Then, we rank the genes based on the average importance scores, and the new landmark genes  
223 (see Additional file 1) can be identified by selecting the top 943 genes (the same number as the L1000).

224 **Figure 5.** The workflow for mining influential genes based deep learning. (A) The architecture and parameter  
225 settings of AutoEncoder. (B) Application of DeepLIFT to compute the importance scores in the Encoder network  
226 and use of D-GEX as a baseline method to predict target genes for performance evaluation.

227

## 228 **AutoEncoder framework**

229 AutoEncoder is a multi-task unsupervised feed-forward neural network with multiple stacked hidden  
230 layers, which is composed of two parts, an encoder and a decoder (Figure 1A). Considering a dataset  $X$   
231 with  $m$  samples and  $n$  features, the encoder  $E|_{X \rightarrow Y}$  aims to map the original data  $X$  to the reduced  
232 representation  $Y$  through the bottleneck layer, and the purpose of the decoder  $D|_{Y \rightarrow X}$  is tuned to  
233 reconstruct the original data  $X$  from the low-dimensional representation  $Y$  by minimizing the difference  
234 between  $X$  and  $\hat{X}$ .

235 Specifically, we use the Python Keras library to implement an AutoEncoder with three hidden layers of  
236 500, 100, and 500 nodes. For a given layer  $l$ , we use sigmoid as the activation function.

$$237 \quad o = f_l(x) = \text{sigmoid}(W_l x + b_l)$$

238 Where  $x$  is an input vector of size  $d$ ,  $W_l$  is the weight matrix of size  $p \times d$ , and  $b_l$  is an intercept  
239 vector of size  $p$ . Given a set of gene expression profiles  $S = (s_1, \dots, s_m)$  containing  $m$  samples, where  
240  $s_m = (g_{m1}, \dots, g_{mn})$  denotes each gene expression profile containing  $n$  genes, the input vector  $s_m$  is  
241 reconstructed to  $\hat{s}_m$  through a series of matrix transformations of multiple network layers. Training an  
242 AutoEncoder involves finding parameters  $\theta = (W, b)$  minimizing a specific loss function. Here, we use  
243 Mean Absolute Error (MAE) as the loss function.

$$244 \quad MAE(g, \hat{g}) = \frac{1}{n} \sum_{i=1}^n |g_i - \hat{g}_i|$$

245 To control overfitting, we add an L2 regularization penalty  $\alpha=1e-6$  on the weight vector. Thus, the loss  
 246 function above becomes:

$$247 \quad loss(g, \hat{g}) = \frac{1}{n} \sum_{i=1}^n |g_i - \hat{g}_i| + \sum_{j=1}^k \alpha \|W_j\|_2^2$$

248 Finally, AutoEncoder is trained using the Adam[36] optimization algorithm with 100 epochs and 10%  
 249 dropout.

### 250 Evaluation metrics

251 Given a test set  $T = (s_1, \dots, s_m)$  containing  $m$  samples, we use two different metrics for the evaluation  
 252 of predicted expression. For each gene  $g_j$ , the definition of MAE is:

$$253 \quad MAE_j = \frac{1}{m} \sum_{i=1}^m |g_{ij} - \hat{g}_{ij}|$$

254 The following equation shows the definition of PCC:

$$255 \quad PCC_j = \rho(g_j, \hat{g}_j) = \frac{\sum_{i=1}^m (g_{ij} - \mu_j)(\hat{g}_{ij} - \hat{\mu}_j)}{\sqrt{\sum_{i=1}^m (g_{ij} - \mu_j)^2} \sqrt{\sum_{i=1}^m (\hat{g}_{ij} - \hat{\mu}_j)^2}}$$

256 where  $PCC_j$  indicates the Pearson correlation coefficient for the  $j$ -th predicted gene and  $\mu_j, \hat{\mu}_j$  are the  
 257 mean of  $g_j, \hat{g}_j$  respectively.

258 The Pearson correlation coefficient, an absolute measure of similarity between genes, does not in itself  
 259 reflect how uncommon that similarity is. Hence, we apply a permutation test to aid in the interpretation  
 260 of similarity. Briefly, in addition to computing the  $PCC_j$  between  $g_j$  and  $\hat{g}_j$ , we also compute the  
 261  $PCC_{null}$  between the  $\hat{g}_j$  and any gene other than  $g_j$  as a reference distribution of similarity values. After

262 that, we compare  $PCC_j$  to  $PCC_{null}$ , and if the fraction of  $PCC_{null}$  that is higher than  $PCC_j$  is lower  
263 than 0.01,  $g_j$  and  $\hat{g}_j$  are considered to be significantly correlated.

264

#### 265 **Additional files**

266 **Additional file 1:** A file in TXT format including a list of genes sorted based on contribution scores.

#### 267 **Abbreviations**

268 CMap: Connectivity Map; MAE: mean absolute error; PCC: Pearson correlation coefficient; GEO: Gene  
269 Expression Omnibus; ADC: lung adenocarcinoma; SCC: lung squamous cell carcinoma.

#### 270 **Declaration**

#### 271 **Ethics approval and consent to participate**

272 Not applicable.

#### 273 **Consent for publication**

274 Not applicable.

#### 275 **Availability of data and materials**

276 Three publicly available datasets are used for our analysis: the microarray-based GEO dataset, the RNA-  
277 Seq-based GTEx dataset and the lung cancer subtype dataset. The first two were downloaded from  
278 [https://cbcl.ics.uci.edu/public\\_data/D-GEX/](https://cbcl.ics.uci.edu/public_data/D-GEX/); the latter, from the GEO database with GSE4573 and  
279 GSE10072.

#### 280 **Competing interests**

281 The authors declare that they have no competing interests.

#### 282 **Funding**

283 This work is supported by the Natural Science Foundation of China (11571173).

284 **Authors' contributions**

285 LPK: Conceived and designed the experiments; LPK and ZTL: Performed the experiments; YYC and  
286 JYF: Analyzed the data; PC and MMX: Contributed reagents/materials/analysis tools; LPK and LYZ:  
287 Wrote the paper. All authors read and approved the final manuscript.

288 **Acknowledgments**

289 This work was supported by the high-performance computing platform of Bioinformatics Center,  
290 Nanjing Agricultural University.

291 **References**

- 292 1. Darmanis S, Sloan SA, Zhang Y, Enge M, Caneda C, Shuer LM, Hayden Gephart MG, Barres  
293 BA, Quake SR: **A survey of human brain transcriptome diversity at the single cell level.**  
294 *Proc Natl Acad Sci U S A* 2015, **112**(23):7285-7290.
- 295 2. Calon A, Lonardo E, Berenguer-Llergo A, Espinet E, Hernando-Momblona X, Iglesias M,  
296 Sevillano M, Palomo-Ponce S, Tauriello DV, Byrom D *et al*: **Stromal gene expression defines**  
297 **poor-prognosis subtypes in colorectal cancer.** *Nat Genet* 2015, **47**(4):320-329.
- 298 3. Lamb J, Crawford ED, Peck D, Modell JW, Blat IC, Wrobel MJ, Lerner J, Brunet JP,  
299 Subramanian A, Ross KN *et al*: **The Connectivity Map: using gene-expression signatures to**  
300 **connect small molecules, genes, and disease.** *Science (New York, NY)* 2006, **313**(5795):1929-  
301 1935.
- 302 4. Ntranos V, Kamath GM, Zhang JM, Pachter L, Tse DN: **Fast and accurate single-cell RNA-**  
303 **seq analysis by clustering of transcript-compatibility counts.** *Genome Biol* 2016, **17**(1):112.
- 304 5. Heimberg G, Bhatnagar R, El-Samad H, Thomson M: **Low Dimensionality in Gene**  
305 **Expression Data Enables the Accurate Extraction of Transcriptional Programs from**  
306 **Shallow Sequencing.** *Cell Syst* 2016, **2**(4):239-250.
- 307 6. Shah S, Lubeck E, Zhou W, Cai L: **In Situ Transcription Profiling of Single Cells Reveals**  
308 **Spatial Organization of Cells in the Mouse Hippocampus.** *Neuron* 2016, **92**(2):342-357.
- 309 7. Subramanian A, Narayan R, Corsello SM, Peck DD, Natoli TE, Lu X, Gould J, Davis JF, Tubelli  
310 AA, Asiedu JK *et al*: **A Next Generation Connectivity Map: L1000 Platform and the First**  
311 **1,000,000 Profiles.** *Cell* 2017, **171**(6):1437-1452 e1417.
- 312 8. Edgar R, Domrachev M, Lash AE: **Gene Expression Omnibus: NCBI gene expression and**  
313 **hybridization array data repository.** *Nucleic Acids Res* 2002, **30**(1):207-210.
- 314 9. Wang X, Ghasedi Dizaji K, Huang H: **Conditional generative adversarial network for gene**  
315 **expression inference.** *Bioinformatics* 2018, **34**(17):i603-i611.
- 316 10. **Gene expression inference with deep learning.** *Bioinformatics* 2016.
- 317 11. Brunel H, Gallardo-Chacon JJ, Buil A, Vallverdu M, Soria JM, Caminal P, Perera A: **MISS: a**  
318 **non-linear methodology based on mutual information for genetic association studies in**  
319 **both population and sib-pairs analysis.** *Bioinformatics* 2010, **26**(15):1811-1818.
- 320 12. Min S, Lee B, Yoon S: **Deep learning in bioinformatics.** *Brief Bioinform* 2017, **18**(5):851-869.
- 321 13. Krizhevsky A, Sutskever I, Hinton GE: **ImageNet classification with deep convolutional**

- 322 **neural networks.** *Communications of the ACM* 2017, **60**(6):84-90.
- 323 14. Chorowski J, Bahdanau D, Serdyuk D, Cho K, Bengio Y: **Attention-Based Models for Speech**  
324 **Recognition.** *Adv Neur In* 2015, **28**.
- 325 15. Li JW, Luong MT, Jurafsky D: **A Hierarchical Neural Autoencoder for Paragraphs and**  
326 **Documents.** *Proceedings of the 53rd Annual Meeting of the Association for Computational*  
327 *Linguistics and the 7th International Joint Conference on Natural Language Processing, Vol 1*  
328 2015, **1**:1106-1115.
- 329 16. Collobert R, Weston J, Bottou L, Karlen M, Kavukcuoglu K, Kuksa P: **Natural Language**  
330 **Processing (Almost) from Scratch.** *J Mach Learn Res* 2011, **12**:2493-2537.
- 331 17. Kelley DR, Snoek J, Rinn JL: **Basset: learning the regulatory code of the accessible genome**  
332 **with deep convolutional neural networks.** *Genome Res* 2016, **26**(7):990-999.
- 333 18. Kalkatawi M, Magana-Mora A, Jankovic B, Bajic VB: **DeepGSR: an optimized deep-learning**  
334 **structure for the recognition of genomic signals and regions.** *Bioinformatics* 2019,  
335 **35**(7):1125-1132.
- 336 19. Zhou J, Lu Q, Gui L, Xu R, Long Y, Wang H: **MTTFsite: cross-cell type TF binding site**  
337 **prediction by using multi-task learning.** *Bioinformatics* 2019, **35**(24):5067-5077.
- 338 20. Umarov R, Kuwahara H, Li Y, Gao X, Solovyev V: **Promoter analysis and prediction in the**  
339 **human genome using sequence-based deep learning models.** *Bioinformatics* 2019,  
340 **35**(16):2730-2737.
- 341 21. Zhou J, Troyanskaya OG: **Predicting effects of noncoding variants with deep learning-based**  
342 **sequence model.** *Nat Methods* 2015, **12**(10):931-934.
- 343 22. Gligorijevic V, Barot M, Bonneau R: **deepNF: deep network fusion for protein function**  
344 **prediction.** *Bioinformatics* 2018, **34**(22):3873-3881.
- 345 23. Chen L, Cai C, Chen V, Lu X: **Learning a hierarchical representation of the yeast**  
346 **transcriptomic machinery using an autoencoder model.** *BMC Bioinformatics* 2016, **17**  
347 **Suppl 1**:9.
- 348 24. Khalili M, Alavi Majd H, Khodakarim S, Ahadi B, Hamidpour M: **PREDICTION OF THE**  
349 **THROMBOEMBOLIC SYNDROME: AN APPLICATION OF ARTIFICIAL NEURAL**  
350 **NETWORKS IN GENE EXPRESSION DATA ANALYSIS.** *ARCHIVES OF ADVANCES IN*  
351 *BIOSCIENCES (JOURNAL OF PARAMEDICAL SCIENCES)* 2016, **7**(2):15-22.
- 352 25. Miotto R, Li L, Kidd BA, Dudley JT: **Deep Patient: An Unsupervised Representation to**  
353 **Predict the Future of Patients from the Electronic Health Records.** *Sci Rep* 2016, **6**:26094.
- 354 26. Chen Q, Song X, Yamada H, Shibasaki R: **Learning Deep Representation from Big and**  
355 **Heterogeneous Data for Traffic Accident Inference;** 2016.
- 356 27. Zeiler M, Fergus R: **Visualizing and Understanding Convolutional Neural Networks,** vol.  
357 8689; 2013.
- 358 28. Springenberg J, Dosovitskiy A, Brox T, Riedmiller M: **Striving for Simplicity: The All**  
359 **Convolutional Net.** 2014.
- 360 29. Simonyan K, Vedaldi A, Zisserman A: **Deep Inside Convolutional Networks: Visualising**  
361 **Image Classification Models and Saliency Maps.** *preprint* 2013.
- 362 30. Shrikumar A, Greenside P, Kundaje A: **Learning Important Features Through Propagating**  
363 **Activation Differences.** 2017.
- 364 31. Zuallaert J, Godin F, Kim M, Soete A, Saeys Y, De Neve W: **SpliceRover: interpretable**  
365 **convolutional neural networks for improved splice site prediction.** *Bioinformatics* 2018,

- 366           **34(24):4180-4188.**
- 367    32.    John GH, Kohavi R, Pflieger K: **Irrelevant Features and the Subset Selection Problem.** In:  
368           *Machine Learning Proceedings 1994.* Edited by Cohen WW, Hirsh H. San Francisco (CA):  
369           Morgan Kaufmann; 1994: 121-129.
- 370    33.    Liaw A, Wiener M: **Classification and Regression by RandomForest.** *Forest* 2001, **23.**
- 371    34.    Bolstad BM, Irizarry RA, Astrand M, Speed TP: **A comparison of normalization methods for**  
372           **high density oligonucleotide array data based on variance and bias.** *Bioinformatics* 2003,  
373           **19(2):185-193.**
- 374    35.    Chen W, Lv H, Nie F, Lin H: **i6mA-Pred: identifying DNA N6-methyladenine sites in the**  
375           **rice genome.** *Bioinformatics* 2019, **35(16):2796-2800.**
- 376    36.    Kingma D, Ba J: **Adam: A Method for Stochastic Optimization.** *International Conference on*  
377           *Learning Representations* 2014.
- 378