

Deep Learning Approach to Identifying Breast Cancer Subtypes Using High-Dimensional Genomic Data

Runpu Chen[§], Le Yang[§], Steve Goodison[†], Yijun Sun^{§¶†*}

[§]Department of Computer Science and Engineering

[¶]Department of Microbiology and Immunology

[†]Department of Biostatistics

The State University of New York at Buffalo, Buffalo, NY 14203

[†]Department of Health Sciences Research

Mayo Clinic, Jacksonville, FL 32224

Abstract

Motivation: Cancer subtype classification has the potential to significantly improve disease prognosis and develop individualized patient management. Existing methods are limited by their ability to handle extremely high-dimensional data and by the influence of misleading, irrelevant factors, resulting in ambiguous and overlapping subtypes.

Results: To address the above issues, we proposed a novel approach to disentangling and eliminating irrelevant factors by leveraging the power of deep learning. Specifically, we designed a deep learning framework, referred to as DeepType, that performs joint supervised classification, unsupervised clustering and dimensionality reduction to learn cancer-relevant data representation with cluster structure. We applied DeepType to the METABRIC breast cancer dataset and compared its performance to state-of-the-art methods. DeepType significantly outperformed the existing methods, identifying more robust subtypes while using fewer genes. The new approach provides a framework for the derivation of more accurate and robust molecular cancer subtypes by using increasingly complex, multi-source data.

Availability and implementation: An open-source software package for the proposed method is freely available at www.acsu.buffalo.edu/~yijunsun/lab/DeepType.html.

Bioinformatics

submitted in January 2019, revised in August 2019

*Please address all correspondence to Dr. Yijun Sun (yijunsun@buffalo.edu).

22 1 Introduction

23 Human cancer is a heterogeneous disease initiated by random somatic mutations and driven by multiple
24 genomic alterations (Hanahan and Weinberg, 2011; Sun *et al.*, 2017). In order to move towards personal-
25 ized treatment regimes, cancers of specific tissues have been divided into subtypes based on the molecular
26 profiles of primary tumors (Sørliie *et al.*, 2001; Curtis *et al.*, 2012; Parker *et al.*, 2009). The premise is that
27 patients of the same molecular subtypes are likely to have similar disease etiology, responses to therapy,
28 and clinical outcomes. Thus, molecular subtyping can reveal information valuable for a range of cancer
29 studies from etiology and tumor biology to prognosis and personalized medicine.

30 Most early work on molecular subtyping has been performed on data obtained from breast cancer
31 tissues (Sørliie *et al.*, 2001, 2003). Typically, breast cancer is not lethal immediately, and thus there
32 is an opportunity to assist with prognostication and patient management using molecular information.
33 Molecular subtyping of breast cancer initially focused on mRNA data obtained from microarray platforms
34 and parsed molecular profiles to stratify patients according to clinical outcomes (Sørliie *et al.*, 2001).
35 Refinement of the subtype categories through validation in independent datasets identified five broad
36 subtypes, including normal-like, luminal A, luminal B, basal, and HER2+, each with distinct clinical
37 outcomes (Sørliie *et al.*, 2003; Parker *et al.*, 2009). These early studies completely altered our views of
38 breast cancer and offered a foundation for the development of therapies tailored to specific subtypes.
39 However, possibly due to the small number of tumor samples used in initial analyses and the technical
40 limitations of the methods used for gene selection and clustering analysis, several large-scale benchmark
41 studies have demonstrated that the current stratification of breast cancer is only approximate, and that
42 the high degree of ambiguity in existing subtyping systems induces uncertainty in the classification of new
43 patients (Weigelt *et al.*, 2010; Mackay *et al.*, 2011).

44 The desire for levels of accuracy that can ultimately lead to clinical utility continues to drive the
45 field to refine breast cancer subtypes (Parker *et al.*, 2009; Haibe-Kains *et al.*, 2012; Shen *et al.*, 2013; Sun
46 *et al.*, 2014, 2017) and to identify molecular subtypes in other cancers (Abeshouse *et al.*, 2015). The recent
47 establishment of international cancer genome consortia (Cancer Genome Atlas Network, 2012; Abeshouse
48 *et al.*, 2015; Curtis *et al.*, 2012) has generally overcome the sample-size issue. In this paper, we focus mainly
49 on developing methods to address the computational challenges associated with detecting cancer related
50 genes and biologically meaningful subtypes using high-dimensional genomics data. Molecular subtyping
51 can be formulated as a supervised-learning problem, that is, to use established tumor subtypes as class
52 labels to perform gene selection and construct a model for the classification of new patients. However, as
53 mentioned above, current subtyping systems provide only a rough stratification of cancer, and supervised-
54 learning based approaches may not enable us to identify novel subtypes. This is because the primary goal
55 of supervised learning is to identify genes to achieve the maximum separation of samples from different
56 subtypes, and genes that support novel subtypes can be considered irrelevant and removed. Consequently,
57 most existing methods were developed within the unsupervised-learning framework. Representative work
58 includes SparseK (Witten and Tibshirani, 2010), iCluster (Shen *et al.*, 2009, 2013) and non-negative
59 matrix factorization (Kormaksson *et al.*, 2012). A major issue with existing methods is that there is
60 no guarantee that subtypes identified through *de novo* clustering are biologically relevant. Presumably,
61 genomics data records all ongoing biological processes in a cell or tissue, where multiple factors interact
62 with each other in a complex and entangled manner. Tumor samples can be grouped based on factors that
63 are not related to the actual disease (e.g., race and eye color). A possible way to address the issue is to
64 use previously established results to guide the detection of new subtypes. However, as the name suggests,
65 *de novo* clustering completely ignores results from previous efforts. Another major limitation is that
66 for computational considerations most existing methods perform data dimensionality reduction through
67 linear transformation (e.g., feature weighting used in SparseK (Witten and Tibshirani, 2010)). Thus,
68 they cannot adequately deal with complex non-linear data and extract pertinent information to detect
69 subtypes residing in non-linear manifolds in a high-dimensional space. Finally, some existing methods do
70 not scale well to handle high-dimensional data. For example, iCluster (Shen *et al.*, 2009, 2013) involves
71 matrix inversion and thus can only process a few thousands of genes. A commonly used practice is to

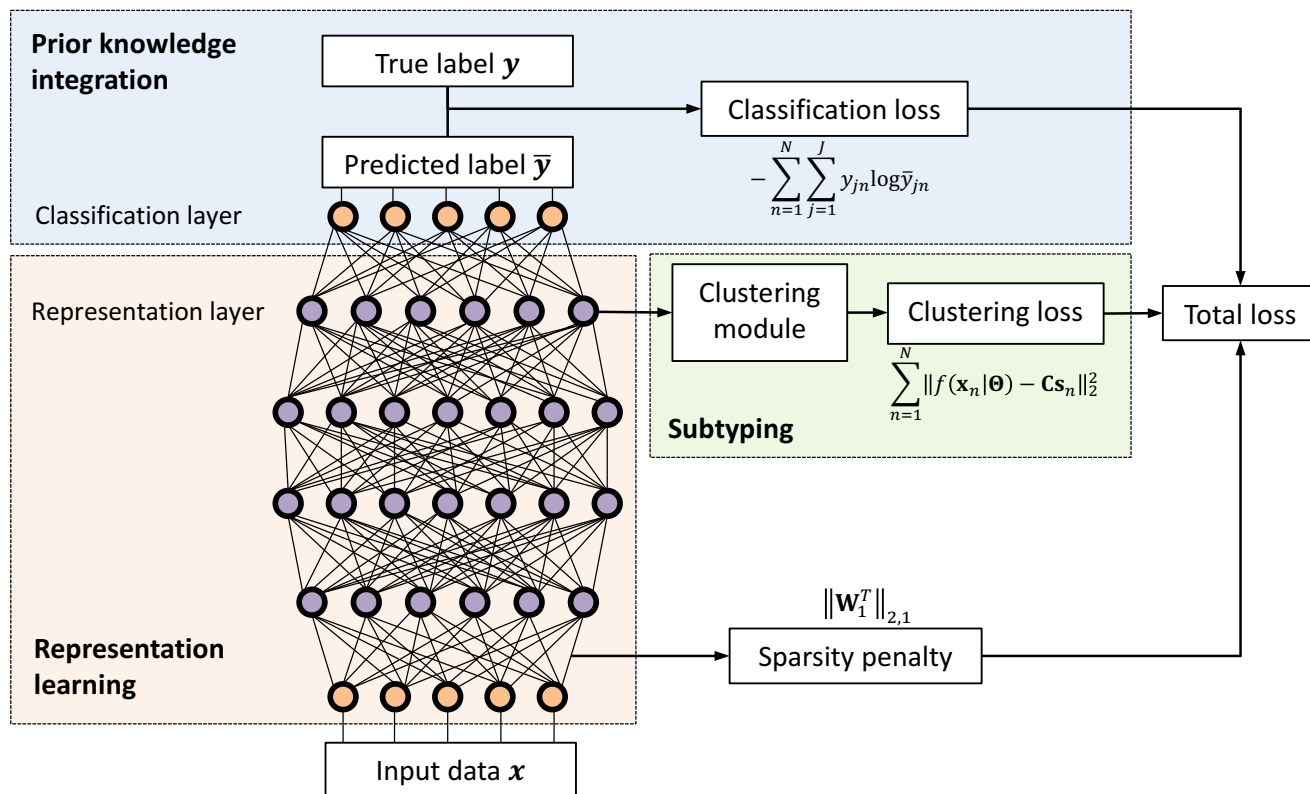


Figure 1: Overview of the proposed deep learning based method for cancer molecular subtyping. It consists of three major components: representation learning, prior knowledge integration, and subtyping. The first part maps raw genomics data onto a representation space, the second part incorporates prior biological knowledge to guide representation learning, and the third part generates subtyping results. The network parameters are learned by minimizing a unified objective function consisting of a classification loss, a clustering loss and a sparsity penalty.

72 perform preprocessing and retain only the most variant genes (Curtis *et al.*, 2012). However, there is no
 73 guarantee that low-variant genes contain no information and the cut-offs used to select variant genes were
 74 usually set somehow arbitrarily.

75 The above observations motivated us to develop a novel deep-learning based approach, referred to as
 76 DeepType, that performs cancer subtyping through joint supervised and unsupervised learning but ad-
 77 dresses their respective limitations. Due to the ability to learn good data representation, deep learning has
 78 recently achieved state-of-the-art performance in computer vision, pattern recognition and bioinformatics
 79 (LeCun *et al.*, 2015; Zheng *et al.*, 2019). For our purpose, by leveraging the power of a multi-layer neural
 80 network for representation learning, we map raw genomics data into a space where clusters can be easily
 81 detected. To ensure the biological relevance of detected clusters, we incorporate prior biological knowledge
 82 to guide representation learning. We train the neural network by minimizing a unified objective function
 83 consisting of a classification loss, a clustering loss and a sparsity penalty. The training process can be
 84 easily performed by using a mini-batch gradient descent method. Thus, our method can handle large
 85 datasets with extremely high dimensionality. Although the idea of using deep learning for clustering is
 86 not new (see, e.g., Xie *et al.* (2016)), to the best of our knowledge, this work represents the *first* attempt to
 87 use deep learning to perform joint supervised and unsupervised learning for cancer subtype classification.
 88 A large-scale experiment was performed that demonstrated that DeepType significantly outperformed the
 89 existing approaches. The new approach provides a framework for the derivation of more accurate and
 90 robust molecular cancer subtypes by using increasingly complex genomic data.

91 2 Methods

92 In this section, we present a detailed description of the proposed method for cancer subtype identification.
 93 We also propose novel procedures for optimizing the associated objective function and estimating the
 94 hyper-parameters.

95 2.1 Deep Learning for Cancer Subtype Identification

96 Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ be a cohort of tumor samples and $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_N]$ be a rough stratification of the
 97 samples (e.g., subtyping results from previous studies), where $\mathbf{x}_n \in \mathbb{R}^D$ is the n -th sample and $\mathbf{y}_n \in \mathbb{R}^J$ is
 98 the corresponding class label vector with $y_{jn} = 1$ if \mathbf{x}_n belongs to the j -th group and 0 otherwise. Our goal
 99 is to identify a small set of cancer related genes and perform clustering analysis on the detected genes to
 100 refine existing classification systems and detect novel subtypes. To this end, we utilize the representation
 101 power of a multi-layer neural network to project raw data onto a representation space where clusters
 102 can be easily detected. As discussed above, clusters identified through unsupervised learning may not be
 103 biologically relevant. To address the issue, we impose an additional constraint that the detected clusters
 104 are concordance with previous results. Specifically, we cast it as a supervised-learning problem, that is,
 105 to find a representation space where the class labels can be accurately predicted.

106 Figure 1 depicts the network structure of the proposed method. It consists of an input layer, M
 107 hidden layers, a classification layer and a clustering module. The M -th hidden layer is designated as the
 108 representation layer, the output of which is fed into the classification layer and the clustering module.
 109 Mathematically, the neural network can be described as follows:

$$\begin{aligned} \mathbf{o}_1 &= \text{sigmoid}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1), \\ \mathbf{o}_m &= \text{sigmoid}(\mathbf{W}_m \mathbf{o}_{m-1} + \mathbf{b}_m), 2 \leq m \leq M, \\ \bar{\mathbf{y}} &= \text{softmax}(\mathbf{W}_{m+1} \mathbf{o}_M + \mathbf{b}_{m+1}), \end{aligned} \quad (1)$$

110 where \mathbf{W}_m , \mathbf{b}_m , and \mathbf{o}_m are the weight matrix, bias term and output of the m -th layer, respectively, and
 111 $\bar{\mathbf{y}}$ is the output of the classification layer. For the purpose of this study, we use sigmoid and softmax as
 112 the activation functions for the hidden and classification layers, respectively. For notational convenience,
 113 let $\Theta = \{(\mathbf{W}_m, \mathbf{b}_m)\}_{m=1}^M$ and denote $f(\mathbf{x}|\Theta) : \mathbb{R}^D \rightarrow \mathbb{R}^{D_M}$ as the mapping function that projects raw
 114 data onto a representation space, where D_M is the number of the nodes in the representation layer and
 115 $D_M \ll D$.

116 We optimize network parameters Θ through joint supervised and unsupervised learning by minimizing
 117 an objective function that consists of a classification loss, a clustering loss and a regularization term. The
 118 classification loss measures the discrepancy between the predicted and given class labels. By construction,
 119 the j -th element of $\bar{\mathbf{y}}_n$ can be interpreted as the probability of \mathbf{x}_n belonging to the j -th group. Thus, we
 120 use the cross entropy to quantify the classification loss:

$$L_{\text{classification}} = - \sum_{n=1}^N \sum_{j=1}^J y_{jn} \log \bar{y}_{jn}. \quad (2)$$

121 We use the K -means method (Lloyd, 1982) to detect clusters in the representation space. The loss
 122 function optimized by K -means is given by

$$L_{\text{clustering}} = \sum_{n=1}^N \|f(\mathbf{x}_n|\Theta) - \mathbf{C}\mathbf{s}_n\|_2^2, \quad (3)$$

123 subject to $\sum_{k=1}^K s_{kn} = 1$, $s_{kn} \in \{0, 1\}$, $\forall k, \forall n$, where K is the number of clusters, \mathbf{C} is a center matrix
 124 with each column representing a cluster center, and \mathbf{s}_n is a binary vector where $s_{kn} = 1$ if \mathbf{x}_n is assigned

125 to cluster k and 0 otherwise. Finally, we impose an $\ell_{2,1}$ -norm regularization (Nie *et al.*, 2010) on the
 126 weight matrix of the first layer to control the model complexity and to select cancer related genes:

$$L_{\text{sparsity}} = \|\mathbf{W}_1^T\|_{2,1} = \sum_{j=1}^D \sqrt{\sum_{i=1}^{D_2} W_{1ij}^2}, \quad (4)$$

127 where W_{1ij} is the ij -th element of \mathbf{W}_1 and D_2 is the number of the nodes in the second layer. The
 128 $\ell_{2,1}$ -norm regularization has an effect of automatically determining the number of nodes activated in the
 129 input layer, and thus the number of genes used in downstream subtyping analysis.

130 Combining the above three losses, we obtain the following novel formulation for cancer subtype iden-
 131 tification:

$$\begin{aligned} & \min_{\{\Theta, \mathbf{S}, \mathbf{C}\}} \sum_{n=1}^N \|f(\mathbf{x}_n|\Theta) - \mathbf{C}\mathbf{s}_n\|_2^2 + \lambda \|\mathbf{W}_1^T\|_{2,1} \\ & \text{subject to } -\sum_{n=1}^N \sum_{j=1}^J y_{jn} \log \bar{y}_{jn} \leq \zeta, \sum_{k=1}^K s_{kn} = 1, s_{kn} \in \{0, 1\}, \forall k, \forall n, \end{aligned} \quad (5)$$

132 where $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_N]$ and λ is a regularization parameter that controls the sparseness of weight matrix
 133 \mathbf{W}_1 . The above formulation can be interpreted as finding a representation space to minimize the clustering
 134 loss while maintaining the classification loss smaller than a user defined upper bound ζ . For ease of
 135 optimization, we move the classification-loss constraint to the objective function and write the problem
 136 in the following equivalent form:

$$\begin{aligned} & \min_{\{\Theta, \mathbf{S}, \mathbf{C}\}} -\sum_{n=1}^N \sum_{j=1}^J y_{jn} \log \bar{y}_{jn} + \alpha \sum_{n=1}^N \|f(\mathbf{x}_n|\Theta) - \mathbf{C}\mathbf{s}_n\|_2^2 + \lambda \|\mathbf{W}_1^T\|_{2,1} \\ & \text{subject to } \sum_{k=1}^K s_{kn} = 1, s_{kn} \in \{0, 1\}, \forall k, \forall n, \end{aligned} \quad (6)$$

137 where α is a tradeoff parameter that controls the balance between the classification and clustering perfor-
 138 mance. In the following sections, we describe how to solve the above optimization problem and estimate
 139 the hyper-parameters.

140 2.2 Optimization

141 The above optimization problem contains three sets of variables, namely, network parameters Θ , assign-
 142 ment matrix \mathbf{S} , and cluster centers \mathbf{C} . It is difficult to solve the problem directly since the parameters
 143 are coupled and \mathbf{S} is a binary matrix. To address the issue, we partition the variables into two groups,
 144 i.e., Θ and (\mathbf{S}, \mathbf{C}) , and employ an alternating optimization strategy to solve the problem. Specifically, we
 145 first perform pre-training to initialize the network by ignoring the clustering module (i.e., setting $\alpha = 0$).
 146 Then, we fix Θ and transform the problem into

$$\min_{\{\mathbf{C}, \mathbf{S}\}} \sum_{n=1}^N \|f(\mathbf{x}_n|\Theta) - \mathbf{C}\mathbf{s}_n\|_2^2, \text{ subject to } \sum_{k=1}^K s_{kn} = 1, s_{kn} \in \{0, 1\}, \forall k, \forall n, \quad (7)$$

147 which can be readily solved by using the standard K -means method. Then, we fix (\mathbf{S}, \mathbf{C}) and write the
 148 problem as

$$\min_{\Theta} -\sum_{n=1}^N \sum_{j=1}^J y_{jn} \log \bar{y}_{jn} + \alpha \sum_{n=1}^N \|f(\mathbf{x}_n|\Theta) - \mathbf{C}\mathbf{s}_n\|_2^2 + \lambda \|\mathbf{W}_1^T\|_{2,1}, \quad (8)$$

149 which can be optimized through back-propagation by using the mini-batch based stochastic gradient
 150 descent method (Kingma and Ba, 2014). The above procedures iterate until convergence.

151 2.3 Parameter Estimation

152 We describe how to estimate the three hyper-parameters of the proposed method, namely regularization
153 parameter λ , tradeoff parameter α , and number of clusters K . In order to avoid a computationally
154 expensive three-dimensional grid search, we first ignore the clustering module by setting $\alpha = 0$ and
155 perform supervised learning to estimate λ . The rationale is that previous subtyping results could provide
156 us with sufficient information to determine the value of λ . Specifically, we randomly partition training
157 data into ten equally-sized sub-datasets, perform ten-fold cross-validation and estimate λ by using the
158 one-standard-error rule (Hastie *et al.*, 2009). Once we determine the value of λ , we perform K -means
159 analysis on the outputs of the representation layer and pre-estimate the number of clusters, denoted
160 as \tilde{K} , as the one that maximizes the average silhouette width (Wiwie *et al.*, 2015). Since the data
161 representation is obtained through supervised learning, which tends to group samples with the same
162 labels together, \tilde{K} is likely to be the lower bound of the true value. Let $K_i = \tilde{K} + i, 0 \leq i \leq T$. For
163 each K_i , we train a deep-learning model by using different α values and record the corresponding ten-fold
164 cross-validation classification errors. By design, α controls the tradeoff between the classification and
165 clustering performance, and the classification error increases with the increase of α . Again, by using the
166 one-standard-error rule, for each K_i , we find the largest α , denoted as α_i , that results in a classification
167 error that is within one standard deviation of the one obtained by setting $\alpha = 0$ (i.e., we require that the
168 obtained classifier does not perform significantly worse than the existing subtyping system), and record
169 the corresponding average silhouette width s_i . Once we run over all possible K_i , we obtain $T + 1$ triplets
170 $\{K_i, \alpha_i, s_i\}_{i=0}^T$. Finally, we determine the number of clusters K and the tradeoff parameter α as the pair
171 that yields the largest average silhouette width. The pseudo-code of the proposed procedure is given in
172 Algorithm S1, and the proposed procedure performed quite well in our numerical experiment (see Figure
173 S1).

174 3 Experiments

175 We conducted a large-scale experiment on breast and bladder cancers to demonstrate the effectiveness of
176 the proposed method. Due to space limit, here we report only the results of the breast cancer study and
177 present the bladder cancer results in Supplementary Data.

178 3.1 Experiment Setting

179 The breast cancer dataset was obtained from the METABRIC study (Curtis *et al.*, 2012), which contains
180 the expression profiles of 25,160 genes from 1,989 primary breast tumor samples and 144 normal breast
181 tissue samples. It is probably the largest single breast cancer dataset assayed to date. For computational
182 convenience, we retained only the top 20,000 most variant genes for the downstream analysis. For model
183 construction and performance evaluation, we randomly partitioned the data into a training and test
184 datasets, containing 80% and 20% of the samples, respectively. In this study, we used the PAM50
185 subtypes (Parker *et al.*, 2009) as class labels in the training process. We designed a four-layer neural
186 network for the joint supervised and unsupervised learning. The numbers of the nodes in the input
187 layer, the two hidden layers and the output layer were set to 20,000, 1,024, 512, and 6, respectively. We
188 employed the Adam method (Kingma and Ba, 2014) to tune the parameters of the model. The learning
189 rate was set to 1e-3, the numbers of training epochs for model initialization and the joint supervised and
190 unsupervised training were set to 300 and 1,500, respectively, and the batch size was set to 256. By
191 using the method proposed in Section 2.3, the number of clusters K , the tradeoff parameter α and the
192 regularization parameter λ were estimated to be 11, 1.2 and 0.006, respectively (see Figure S1). To ensure
193 that the constructed model did not overfit the data, we tracked the training and validation losses in the
194 training process (see Figure S2), and no sign of over-fitting was observed.

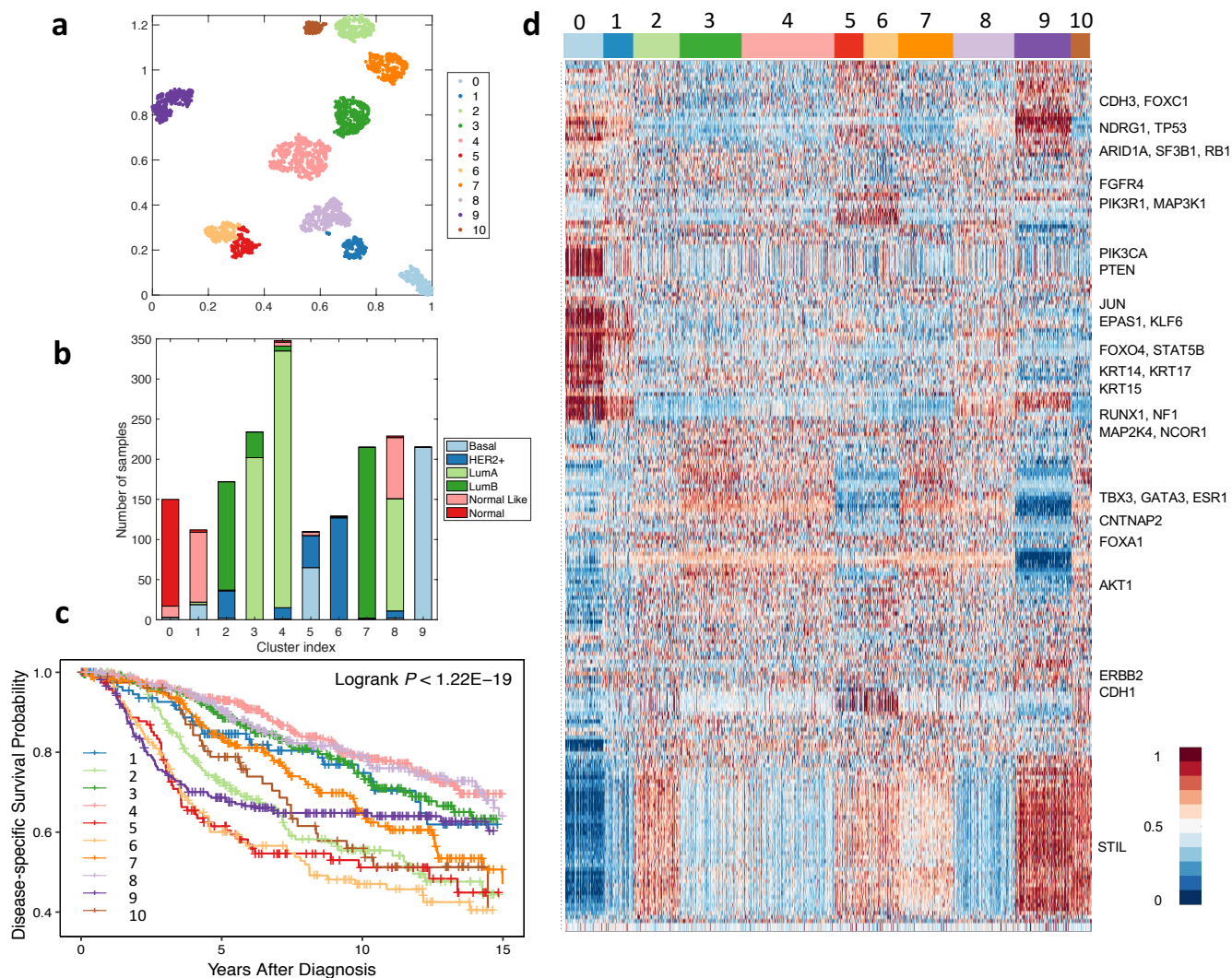


Figure 2: DeepType identified ten clinically relevant breast cancer subtypes. (a) The sample distributions of the identified clusters visualized by t-SNE. Nearly all of the normal tissue samples were grouped into a single cluster (i.e., Cluster 0), and the tumor samples were grouped into ten well-separated clusters, labelled as DeepType 1-10. (b) The PAM50 composition of the identified clusters. (c) Survival data analysis showed that the ten identified subtypes were associated with distinct clinical outcomes. (d) The heatmap of the 218 selected genes showed that the identified clusters exhibited distinct transcriptional characteristics on several gene modules. The samples were arranged by the clustering assignments, and the expression levels were linearly scaled into $[0, 1]$ across samples.

195 3.2 Clinically Relevant Subtypes Revealed by DeepType

196 By applying the proposed method to the breast cancer dataset, a total of 218 genes were selected and 11
 197 clusters were detected. To visualize the identified clusters, we applied t-SNE (van der Maaten and Hinton,
 198 2008) to the outputs of the representation layer. Figure 2(a-b) present the sample distributions of the
 199 identified clusters and their PAM50 compositions, respectively. We can see that nearly all of the normal
 200 tissue samples were grouped into a single cluster (i.e., Cluster 0), and the tumor samples were grouped
 201 into ten well-separated clusters, labelled as DeepType 1-10. To demonstrate the clinical relevance of the
 202 identified tumor subtypes, a disease-specific survival data analysis was performed. Figure 2(c) shows that
 203 the ten subtypes were associated with distinct prognostic outcomes (logrank test, p -value $< 1.22e-19$).
 204 Further internal and external validation analysis of the detected clusters is presented in Sections 3.3 and

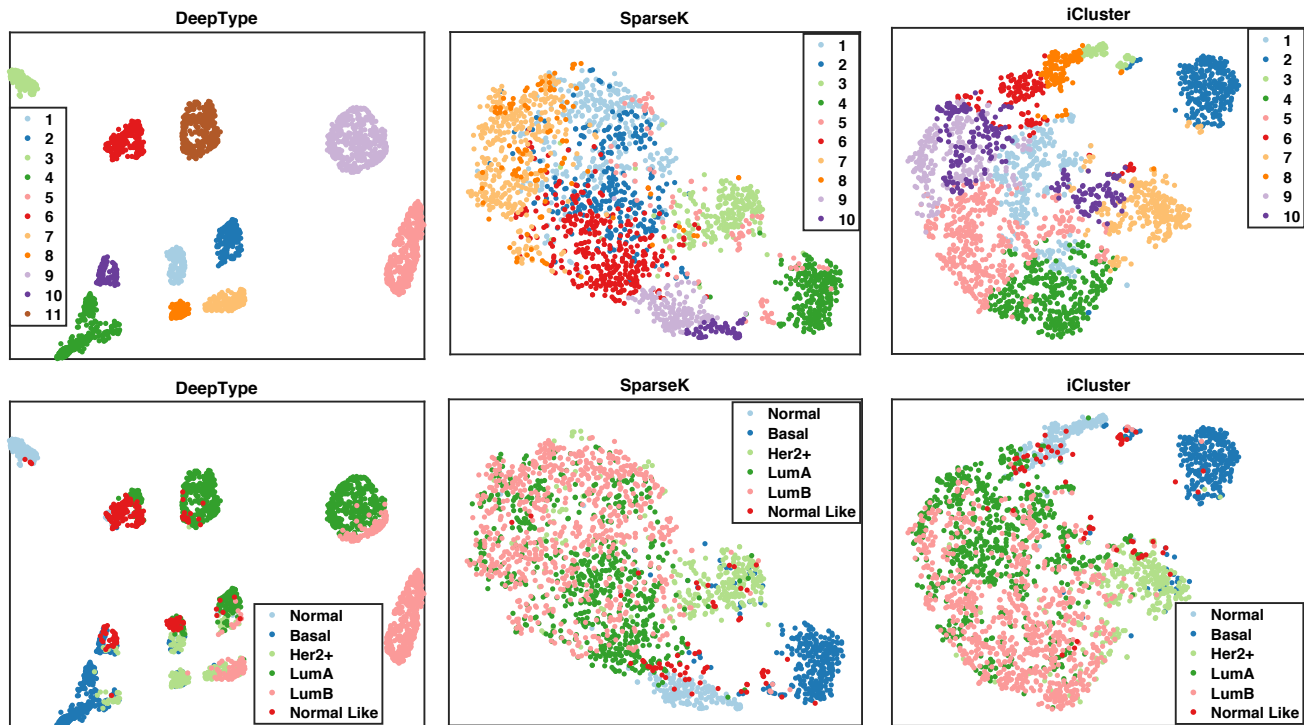


Figure 3: Visualization of the sample distributions of the clusters detected by three methods applied to data containing 10,000 most variant genes. Each sample was color-coded by its clustering assignment (top) and PAM50 label (bottom). DeepType revealed a clear eleven-cluster structure including a cluster comprising primarily normal tissue samples.

205 3.4.

206 Figure 2(d) represents the heatmap of the 218 selected genes. The descriptions of the genes are given
 207 in Table S1. The detected subtypes contain distinct transcriptional characteristics associated with several
 208 gene co-expression modules and key cancer genes. Most normal-like samples were grouped into DeepType
 209 1, and have an expression pattern similar to normal samples. The luminal A samples were divided into
 210 DeepTypes 3, 4 and 8 with low expression on the *STIL* module (key gene: *STIL*) and intermediate
 211 expression on the *GATA3* module (key genes: *TBX3*, *GATA3*, *ESR1*, *CNTNAP2* and *FOXA1*). Among
 212 the three subtypes, the expression of the *KRT* family (key genes: *KRT14*, *KRT15* and *KRT17*) were
 213 highest in DeepType 8, intermediate in DeepType 4, and lowest in DeepType 3. The luminal B samples
 214 were partitioned into DeepTypes 2, 7 and 10, with intermediate to high expression of the *GATA3* and
 215 *STIL* gene modules, and low expression of *CDH3* and *FOXC1*. Among the three subtypes, the expression
 216 of the genes in the *STIL* module was highest in DeepType 10, intermediate in DeepType 2 and lowest in
 217 DeepType 7. DeepTypes 5 and 6, which were dominated by mixed HER2+/basal and HER2+ samples,
 218 respectively, had very high expression on *ERBB2* and *CDH1* and low expression on *TBX3*, *GATA3* and
 219 *ESR1* genes. DeepType 9, composed entirely of basal samples, had low expression in the *GATA3* module
 220 and high expression in the *STIL* and *KRT* modules. The distinct expression patterns and prognostic
 221 outcomes of the detected clusters suggest that the proposed method is able to detect new breast cancer
 222 subtypes beyond the PAM50 classification, and a further analysis could reveal information of the breast
 223 cancer molecular taxonomy in a higher level of resolution.

224 3.3 Comparison Study

225 To further demonstrate the effectiveness of the proposed method, we compared it with two state-of-the-
 226 art methods, namely SparseK (Witten and Tibshirani, 2010) and iCluster (Shen *et al.*, 2009). Both

Table 1: External evaluation of subtypes identified by three methods applied to datasets with a various number of input genes. iCluster failed on datasets with 15,000 and 20,000 genes. DeepType significantly outperformed SparseK (p -value $\leq 7.7e-14$) and iCluster (p -value $\leq 1.3e-19$, Wilcoxon rank-sum test).

		Average Purity				NMI			
		5000	10000	15000	20000	5000	10000	15000	20000
PAM50	DeepType	0.86	0.80	0.85	0.87	0.62	0.56	0.62	0.61
	SparseK	0.65	0.68	0.64	0.63	0.39	0.40	0.37	0.38
	iCluster	0.43	0.65	/	/	0.08	0.34	/	/
Histological Grade	DeepType	0.67	0.67	0.66	0.67	0.12	0.12	0.12	0.13
	SparseK	0.63	0.66	0.65	0.63	0.11	0.10	0.12	0.09
	iCluster	0.55	0.63	/	/	0.04	0.11	/	/
NPI	DeepType	0.57	0.55	0.60	0.58	0.08	0.09	0.10	0.07
	SparseK	0.56	0.55	0.59	0.58	0.07	0.09	0.07	0.07
	iCluster	0.56	0.57	/	/	0.04	0.09	/	/
GGI	DeepType	0.69	0.68	0.70	0.69	0.15	0.16	0.14	0.13
	SparseK	0.69	0.70	0.70	0.69	0.12	0.14	0.13	0.12
	iCluster	0.68	0.67	/	/	0.04	0.11	/	/
Oncotype DX	DeepType	0.88	0.85	0.87	0.86	0.25	0.26	0.27	0.24
	SparseK	0.75	0.78	0.78	0.76	0.14	0.16	0.15	0.13
	iCluster	0.64	0.74	/	/	0.06	0.13	/	/

227 methods perform feature selection and clustering analysis simultaneously, and iCluster was also used in
 228 the METABRIC study (Curtis *et al.*, 2012). The source code of the two methods was downloaded from
 229 the CRAN website^{1,2}. Following Shen *et al.* (2013), we tuned the parameters of iCluster (i.e., the number
 230 of clusters K and the sparsity penalty coefficient λ) by maximizing the reproducibility index. SparseK
 231 also contains two parameters, the number of clusters K and the ℓ_1 regularization parameter λ . By using
 232 the method described in Witten and Tibshirani (2010), we first estimated the optimal λ for each K , and
 233 then determined the value of the optimal K based on gap statistic (Tibshirani *et al.*, 2001). To test the
 234 ability of the three methods to handle high-dimensional data, we generated four datasets each containing
 235 a different number of the most variant genes, ranging from 5,000, 10,000, 15,000 and 20,000. Although
 236 we herein considered only gene expression data, it is possible to perform cancer subtyping by integrating
 237 genomics data from different platforms. Therefore, the ability to handle high-dimensional data is an
 238 important consideration in algorithm development. Below, we performed a series of quantitative and
 239 qualitative analyses to compare the performance of the three methods.

240 We first visualized the sample distributions of the clusters detected by the three methods (Figure
 241 3). Since iCluster failed on the datasets with 15,000 and 20,000 genes due to the need of performing
 242 matrix inversion of high-dimensional data, we considered only the results generated by using the dataset
 243 with 10,000 genes. We can see that DeepType identified eleven well-defined clusters, nearly all normal
 244 tissue samples were grouped into a single cluster, and the clusters that composed of tumor samples were
 245 well-separated and highly concordant with the PAM50 labels. In contrast, for SparseK and iCluster, the
 246 normal tissue samples were grouped into multiple clusters, which suggests that genes unrelated to cancer
 247 were selected. Moreover, the tumor samples with different PAM50 labels overlapped considerably, and
 248 did not exhibit a clear clustering structure.

249 We then performed a series of external and internal evaluations of the clusters detected by the three
 250 methods. For external evaluation, we assessed the concordance between the identified cancer subtypes and
 251 some widely used clinical and prognostic characteristics of breast cancer, including the PAM50 subtype
 252 (Parker *et al.*, 2009), histological grade, Nottingham prognostic index (NPI) (Haybittle *et al.*, 1982),

¹<https://cran.r-project.org/web/packages/iCluster/index.html>

²<https://cran.r-project.org/web/packages/sparcl/index.html>

Table 2: Internal evaluation of subtypes identified by three methods applied to datasets with a various number of input genes. The Davies–Bouldin index is a value in $[0, \text{inf})$, and a smaller value suggests a better clustering scheme. DeepType significantly outperformed SparseK ($p\text{-value} \leq 7.8\text{e-}5$) and iCluster ($p\text{-value} \leq 7.8\text{e-}5$, Wilcoxon rank-sum test).

	Silhouette width			Davies–Bouldin index		
	DeepType	SparseK	iCluster	DeepType	SparseK	iCluster
5000	0.48	0.17	0.33	1.01	1.88	1.79
10000	0.48	0.22	0.33	0.87	1.94	1.23
15000	0.44	0.19	/	0.69	1.92	/
20000	0.63	0.15	/	0.67	2.31	/

Table 3: The numbers of genes selected by DeepType, iCluster and SparseK applied to datasets containing a various number of input genes.

# of input genes	DeepType	SparseK	iCluster
5000	182	949	521
10000	239	982	728
15000	250	918	/
20000	218	886	/

253 gene expression grade index (GGI) (Sotiriou *et al.*, 2006) and the Oncotype DX prognostic test (Sparano
 254 *et al.*, 2018) (see Table S2 for a detailed description). Specifically, we used average purity and normalized
 255 mutual information (NMI) to evaluate the extent to which the identified subtypes matched the above
 256 described characteristics. The results are reported in Table 1. Our analysis showed that the subtypes
 257 identified by DeepType were highly concordant with the clinical variables and prognostic information. In
 258 all cases, the results generated by DeepType matched the PAM50 labels to the highest degree. This is
 259 expected since the PAM50 labels were used in training DeepType. Our method also produced the highest
 260 agreement with the histological grades, NPI and GGI. Notably, when compared with Oncotype DX, the
 261 average purities and NMI scores of DeepType were much higher than the other two methods. This is
 262 highly significant since while both NPI and GGI provide some values in predicting the clinical outcomes
 263 of breast cancer patients, Oncotype DX is the only test supported by level II evidence (Sparano *et al.*,
 264 2018). We performed a Wilcoxon rank-sum test to compare the overall performance of DeepType and
 265 the two competing methods. The p -values are $7.7\text{e-}14$ (DeepType vs. SparseK) and $1.3\text{e-}19$ (DeepType
 266 vs. iCluster).

267 We next performed internal evaluation of the subtypes identified by the three methods. Internal
 268 evaluation utilizes only the intrinsic information of cluster assignments to assess the quality of obtained
 269 clusters, and compactness and separability are the two most important considerations (Halkidi *et al.*,
 270 2001). A compact and separable clustering structure means that samples in each cluster are homogeneous
 271 and different clusters are far away from each other, allowing new patients to be assigned with high certainty
 272 and low ambiguity. For the purpose of this study, we used the silhouette width (Wiwie *et al.*, 2015) and the
 273 Davies–Bouldin index (Davies and Bouldin, 1979) to quantify the cluster compactness and separability.
 274 The results are reported in Table 2. In all cases, DeepType resulted in the highest silhouette width and
 275 the lowest Davies–Bouldin index, which is consistent with the visualization result presented in Figure 3.
 276 To compare the overall performance, the Wilcoxon rank-sum test was performed. DeepType significantly
 277 outperformed SparseK ($p\text{-value} \leq 7.8\text{e-}5$) and iCluster ($p\text{-value} \leq 7.8\text{e-}5$). Our analysis suggested that
 278 our method resulted in subtypes with significantly higher cluster quality than the competing methods.

279 Finally, we compared the ability of the three methods to select relevant genes from high-dimensional
 280 data for clustering analysis. Table 3 reports the numbers of genes selected by the three methods applied to
 281 the data with a various number of input genes. Notably, while DeepType achieved the best result in terms

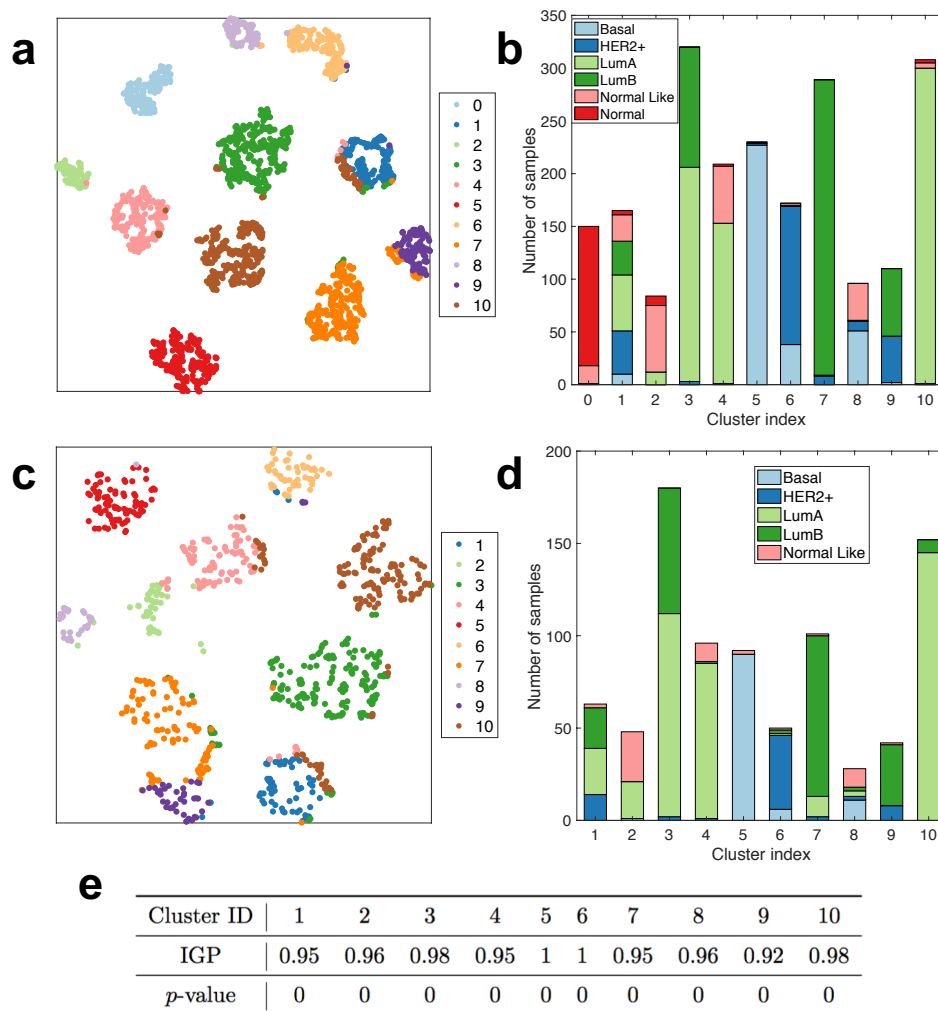


Figure 4: Results of a validation study performed on the METABRIC and SUPERTAM datasets. (a-d) The clusters detected in the METABRIC (top row) and SUPERTAM datasets (middle row) were compact and well-separated and had very similar PAM50 compositions. (e) In-group proportion (IGP) scores and *p*-values (computed based on 1,000 permutations) showed that the clusters identified in the METABRIC data were reproducible in the SUPERTAM data.

282 of both internal and external criteria, it selected the fewest genes in all cases. For clinical applications, the
 283 ability to select fewer genes can help to develop a more economic clinical assay for breast cancer subtype
 284 identification.

285 3.4 Validation Study

286 To demonstrate the generalization capability of the proposed method, we performed a validation study
 287 using the METABRIC data for training and SUPERTAM data (Haibe-Kains *et al.*, 2012) for testing. The
 288 SUPERTAM dataset contains the expression profiles of 13,092 genes from 856 breast tumor samples. Prior
 289 to the analysis, we identified 10,087 genes present in both datasets and used ComBat (Johnson *et al.*, 2007)
 290 to remove batch effects. Using the expression measures of the selected genes, we trained a deep-learning
 291 model using the METABRIC dataset and identified eleven clusters including one comprising dominantly
 292 normal samples. We then applied the constructed model to the validation dataset and classified each
 293 sample into one of the eleven clusters using the nearest shrunken centroid classifier (Tibshirani *et al.*,
 294 2002). Since the SUPERTAM data does not contain normal samples, only three samples were classified

295 into the normal cluster and thus omitted in the further analysis. Figure 4(a-d) presents the sample
296 distributions and PAM50 compositions of the identified clusters. We observed that the clusters detected
297 in the two datasets were compact and well-separated and had similar PAM50 compositions. To provide
298 a quantitative analysis of the reproducibility of the detected clusters, we employed the strategy proposed
299 in Kapp and Tibshirani (2006) and calculated the in-group proportion (IGP) score and p -value for each
300 cluster (Figure 4(e)). Our analysis showed that the identified clusters were reproducible (p -value = 0)
301 and that the proposed method generalizes well on independent datasets.

302 3.5 Robustness Analysis

303 DeepType detects disease molecular subtypes through joint supervised and unsupervised learning, where
304 the class labels from previous studies are usually error-prone. To investigate how DeepType performs in
305 the presence of label noise, we performed a robustness analysis where we corrupted the PAM50 labels
306 of a certain percentage of randomly selected samples in the METABRIC training dataset, constructed a
307 deep-learning model using the corrupted data, and applied the model to the test dataset. To assess the
308 performance of the constructed model, we computed the Rand index by comparing the cluster assignments
309 of the test samples with their PAM50 labels and those obtained by using the original training dataset
310 (i.e., no corrupted labels). To remove random variations, the experiment was repeated five times. Figure
311 S3 presents the results obtained by using the training data containing a varying percentage of corrupted
312 labels ranging from 0% to 20%. We can see that DeepType performed similarly with up to 10% label
313 errors. Considering that the PAM50 label set itself contains an unknown percentage of errors, our method
314 is very robust against label noise.

315 4 Discussion

316 In this paper, we developed a deep-learning based approach for cancer subtype identification that addresses
317 some technical limitations of existing methods. The new method performed significantly better than two
318 commonly used approaches in terms of both internal and external evaluation criteria. By leveraging the
319 power of deep learning, the new method is able to handle data with extremely high dimensionality. We
320 further demonstrated that the method generalizes well on independent datasets and is very robust against
321 label noise.

322 The proposed method has several limitations that are worthwhile to mention. Usually, training a
323 deep-learning model requires a large amount of data. The method is thus not applicable to cancers for
324 which only a small number of samples have been assayed. In this study, we applied the method to breast
325 and bladder cancers where molecular subtypes are well established and thus can be used to guide the
326 detection of new subtypes. However, for many other cancers, molecular subtypes have not yet been well
327 established. It is possible to use other clinical variables (e.g., tumor grade) to guide the identification
328 of cancer subtypes and we have showed that our approach performed well in the presence of label noise.
329 Further investigations are warranted to explore such possibilities.

330 In this paper, we presented a proof-of-concept study considering only gene expression data. Several
331 studies have recently demonstrated that combining cross-platform data could provide more information
332 for cancer subtype identification (see, e.g., Shen *et al.* (2013) and Zhang *et al.* (2012)). It is possible to
333 use deep learning to integrate genomics data from different platforms, including mRNA, copy number,
334 somatic mutation and methylation, for cancer subtyping. However, currently there are debates on how to
335 design a network to process multiple data types (Wang *et al.*, 2015). As the future work, we will perform
336 a large-scale experiment to look into this issue to identify the optimal network structure for genomics
337 data analysis. It is expected that more accurate and robust cancer subtypes would be revealed.

338 References

- 339 Abeshouse, A., Ahn, J., Akbani, R., Ally, A., Amin, S., Andry, C. D., Annala, M., Aprikian, A., Armenia,
340 J., Arora, A., *et al.* (2015). The molecular taxonomy of primary prostate cancer. *Cell*, **163**(4), 1011–
341 1025.
- 342 Cancer Genome Atlas Network (2012). Comprehensive molecular portraits of human breast tumours.
343 *Nature*, **490**(7418), 61–70.
- 344 Curtis, C., Shah, S. P., Chin, S.-F., Turashvili, G., Rueda, O. M., Dunning, M. J., Speed, D., Lynch,
345 A. G., Samarajiwa, S., Yuan, Y., *et al.* (2012). The genomic and transcriptomic architecture of 2,000
346 breast tumours reveals novel subgroups. *Nature*, **486**(7403), 346–352.
- 347 Davies, D. L. and Bouldin, D. W. (1979). A cluster separation measure. *IEEE Transactions on Pattern*
348 *Analysis and Machine Intelligence*, **1**(2), 224–227.
- 349 Haibe-Kains, B., Desmedt, C., Loi, S., Culhane, A. C., Bontempi, G., Quackenbush, J., and Sotiriou,
350 C. (2012). A three-gene model to robustly identify breast cancer molecular subtypes. *Journal of the*
351 *National Cancer Institute*, **104**(4), 311–325.
- 352 Halkidi, M., Batistakis, Y., and Vazirgiannis, M. (2001). On clustering validation techniques. *Journal of*
353 *Intelligent Information Systems*, **17**(2-3), 107–145.
- 354 Hanahan, D. and Weinberg, R. A. (2011). Hallmarks of cancer: the next generation. *Cell*, **144**(5),
355 646–674.
- 356 Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning*. Springer, New
357 York.
- 358 Haybittle, J., Blamey, R., Elston, C., Johnson, J., Doyle, P., Campbell, F., Nicholson, R., and Griffiths,
359 K. (1982). A prognostic index in primary breast cancer. *British Journal of Cancer*, **45**(3), 361–366.
- 360 Johnson, W. E., Li, C., and Rabinovic, A. (2007). Adjusting batch effects in microarray expression data
361 using empirical bayes methods. *Biostatistics*, **8**(1), 118–127.
- 362 Kapp, A. V. and Tibshirani, R. (2006). Are clusters found in one dataset present in another dataset?
363 *Biostatistics*, **8**(1), 9–31.
- 364 Kingma, D. P. and Ba, J. (2014). Adam: a method for stochastic optimization. In *International Conference*
365 *on Learning Representations*, pages 1–13.
- 366 Kormaksson, M., Booth, J. G., Figueroa, M. E., and Melnick, A. (2012). Integrative model-based cluster-
367 ing of microarray methylation and expression data. *The Annals of Applied Statistics*, **6**(3), 1327–1347.
- 368 LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *Nature*, **521**(7553), 436–444.
- 369 Lloyd, S. (1982). Least squares quantization in PCM. *IEEE Transactions on Information Theory*, **28**(2),
370 129–137.
- 371 Mackay, A., Weigelt, B., Grigoriadis, A., Kreike, B., Natrajan, R., A’hern, R., Tan, D. S., Dowsett, M.,
372 Ashworth, A., and Reis-Filho, J. S. (2011). Microarray-based class discovery for molecular classification
373 of breast cancer: analysis of interobserver agreement. *Journal of the National Cancer Institute*, **103**(8),
374 662–673.
- 375 Nie, F., Huang, H., Cai, X., and Ding, C. H. (2010). Efficient and robust feature selection via joint
376 $\ell_{2,1}$ -norms minimization. In *Advances in Neural Information Processing Systems*, pages 1813–1821.

- 377 Parker, J. S., Mullins, M., Cheang, M. C., Leung, S., Voduc, D., Vickery, T., Davies, S., Fauron, C., He,
378 X., Hu, Z., *et al.* (2009). Supervised risk predictor of breast cancer based on intrinsic subtypes. *Journal*
379 *of Clinical Oncology*, **27**(8), 1160–1167.
- 380 Shen, R., Olshen, A. B., and Ladanyi, M. (2009). Integrative clustering of multiple genomic data types
381 using a joint latent variable model with application to breast and lung cancer subtype analysis. *Bioin-*
382 *formatics*, **25**(22), 2906–2912.
- 383 Shen, R., Wang, S., and Mo, Q. (2013). Sparse integrative clustering of multiple omics data sets. *The*
384 *Annals of Applied Statistics*, **7**(1), 269 – 294.
- 385 Sørlie, T., Perou, C. M., Tibshirani, R., Aas, T., Geisler, S., Johnsen, H., Hastie, T., Eisen, M. B., Van
386 De Rijn, M., Jeffrey, S. S., *et al.* (2001). Gene expression patterns of breast carcinomas distinguish
387 tumor subclasses with clinical implications. *Proceedings of the National Academy of Sciences of the*
388 *United States of America*, **98**(19), 10869–10874.
- 389 Sørlie, T., Tibshirani, R., Parker, J., Hastie, T., Marron, J., Nobel, A., Deng, S., Johnsen, H., Pesich, R.,
390 Geisler, S., *et al.* (2003). Repeated observation of breast tumor subtypes in independent gene expression
391 data sets. *Proceedings of the National Academy of Sciences of the United States of America*, **100**(14),
392 8418–8423.
- 393 Sotiriou, C., Wirapati, P., Loi, S., Harris, A., Fox, S., Smeds, J., Nordgren, H., Farmer, P., Praz, V.,
394 Haibe-Kains, B., *et al.* (2006). Gene expression profiling in breast cancer: understanding the molecular
395 basis of histologic grade to improve prognosis. *Journal of the National Cancer Institute*, **98**(4), 262–272.
- 396 Sparano, J. A., Gray, R. J., Makower, D. F., Pritchard, K. I., Albain, K. S., Hayes, D. F., Geyer Jr, C. E.,
397 Dees, E. C., Goetz, M. P., Olson Jr, J. A., *et al.* (2018). Adjuvant chemotherapy guided by a 21-gene
398 expression assay in breast cancer. *New England Journal of Medicine*, **379**(2), 111–121.
- 399 Sun, Y., Yao, J., Nowak, N., and Goodison, S. (2014). Cancer progression modeling using static sample
400 data. *Genome Biology*, **15**(8), 440.
- 401 Sun, Y., Yao, J., Yang, L., Chen, R., Nowak, N. J., and Goodison, S. (2017). Computational approach
402 for deriving cancer progression roadmaps from static sample data. *Nucleic Acids Research*, **45**(9), e69.
- 403 Tibshirani, R., Walther, G., and Hastie, T. (2001). Estimating the number of clusters in a data set via
404 the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **63**(2),
405 411–423.
- 406 Tibshirani, R., Hastie, T., Narasimhan, B., and Chu, G. (2002). Diagnosis of multiple cancer types
407 by shrunken centroids of gene expression. *Proceedings of the National Academy of Sciences*, **99**(10),
408 6567–6572.
- 409 van der Maaten, L. and Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning*
410 *Research*, **9**, 2579–2605.
- 411 Wang, W., Arora, R., Livescu, K., and Bilmes, J. (2015). On deep multi-view representation learning. In
412 *International Conference on Machine Learning*, pages 1083–1092.
- 413 Weigelt, B., Mackay, A., A’hern, R., Natrajan, R., Tan, D. S., Dowsett, M., Ashworth, A., and Reis-Filho,
414 J. S. (2010). Breast cancer molecular profiling with single sample predictors: a retrospective analysis.
415 *The Lancet Oncology*, **11**(4), 339–349.
- 416 Witten, D. M. and Tibshirani, R. (2010). A framework for feature selection in clustering. *Journal of the*
417 *American Statistical Association*, **105**(490), 713–726.

- 418 Wiwie, C., Baumbach, J., and Röttger, R. (2015). Comparing the performance of biomedical clustering
419 methods. *Nature Methods*, **12**(11), 1033–1038.
- 420 Xie, J., Girshick, R., and Farhadi, A. (2016). Unsupervised deep embedding for clustering analysis. In
421 *International Conference on Machine Learning*, pages 478–487.
- 422 Zhang, S., Liu, C.-C., Li, W., Shen, H., Laird, P. W., and Zhou, X. J. (2012). Discovery of multi-
423 dimensional modules by integrative analysis of cancer genomic data. *Nucleic Acids Research*, **40**(19),
424 9379–9391.
- 425 Zheng, W., Yang, L., Genco, R. J., Wactawski-Wende, J., Buck, M., and Sun, Y. (2019). SENSE: Siamese
426 neural network for sequence embedding and alignment-free comparison. *Bioinformatics*, **35**(11), 1820–
427 1828.

Supplementary Data

Deep Learning Approach to Identifying Breast Cancer Subtypes Using High-Dimensional Genomic Data

Runpu Chen, Le Yang, Steve Goodison, Yijun Sun*

Algorithm 1: Hyper-parameter estimation (\mathbf{X} , \mathbf{Y} , \mathcal{A} , \mathcal{L} , T)

Input: training data \mathbf{X} , class labels \mathbf{Y} , T , $\mathcal{A} = \{a_1, \dots, a_J\}$, $\mathcal{L} = \{\lambda_1, \dots, \lambda_L\}$

Output: estimated parameters α^* , λ^* , K^*

```
1 Randomly partition  $(\mathbf{X}, \mathbf{Y})$  into ten folds  $\{(\mathbf{X}_i, \mathbf{Y}_i)\}_{i=1}^{10}$ ;  
2 Estimate  $\lambda^*$  through ten-fold cross validation;  
3 Compute average classification error  $e_0$  and one standard error  $\sigma_0$ ;  
4 Estimate  $\tilde{K}$  by maximizing average silhouette width;  
5 for  $i = 0$  to  $T$  do  
6    $K_i = \tilde{K} + i$ ;  
7   for  $j = 1$  to  $J$  do  
8      $\alpha = a_j$ ;  
9     Solve Problem (6);  
10    Compute average classification error  $e_j$ ;  
11    Compute average silhouette width  $\tilde{s}_j$ ;  
12  end  
13   $j^* = \arg \max_{1 \leq j \leq J} j$ , subject to  $e_j \leq e_0 + \sigma_0$ ;  
14  /* one-standard-error rule */  
15   $\alpha_i = a_{j^*}$ ;  
16   $s_i = \tilde{s}_{j^*}$ ;  
17 end  
18  $i^* = \arg \max_{0 \leq i \leq T} s_i$ ;  
19  $K^* = K_{i^*}$ ;  
20  $\alpha^* = \alpha_{i^*}$ 
```

*Please address all correspondence to Dr. Yijun Sun (yijunsun@buffalo.edu).

Table S1: DeepType identified 218 genes to be informative for breast cancer subtyping

See the attached Excel file

Table S2: Clinical and prognostic characteristics of breast cancer

Characteristics	Class label
PAM50 subtype	basal, HER2+, luminal A/B, normal-like
Histological grade	1, 2, 3
NPI	1, 2, 3, 4
GGI	low risk, high risk
Oncotype DX	low risk, intermediate risk, high risk

Table S3: Parameters of iCluster and SparseK used in the breast and bladder cancer experiments

	breast cancer	bladder cancer
iCluster	$\lambda = 0.01, K = 10$	$\lambda = 0.003, K = 5$
SparseK	$\lambda = 10, K = 10$	$\lambda = 10, K = 3$

Table S4: The indexes of the samples in the training and test datasets used in the breast cancer experiment (for the reproducibility purpose)

See the attached Excel file

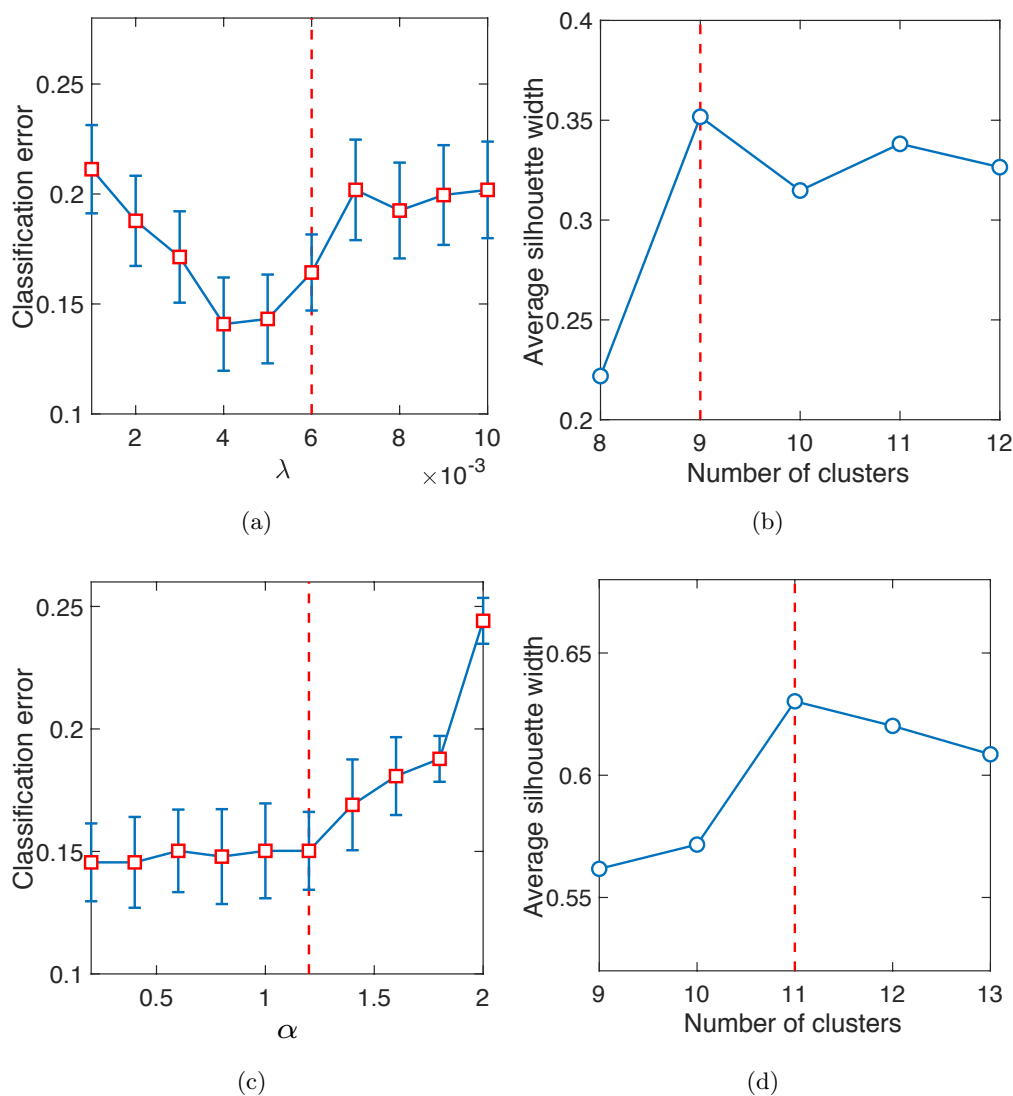


Figure S1: Hyper-parameter estimation. (a) The regularization parameter λ was estimated to be 0.006 based on the one-standard-error rule. (b) The number of clusters was pre-estimated to be 9 based on the average silhouette width. (c) We searched a range of values to estimate the number of clusters. For each $K \geq 9$, we trained a deep learning model by using different α values and estimated the optimal α by using the one-standard-error rule. The figure presents an example showing that the optimal α was estimated to be 1.2 for $K = 11$. (d) The number of clusters was finally determined to be 11 by maximizing the average silhouette width. See Section 2.3 for a detailed description of hyper-parameter estimation.

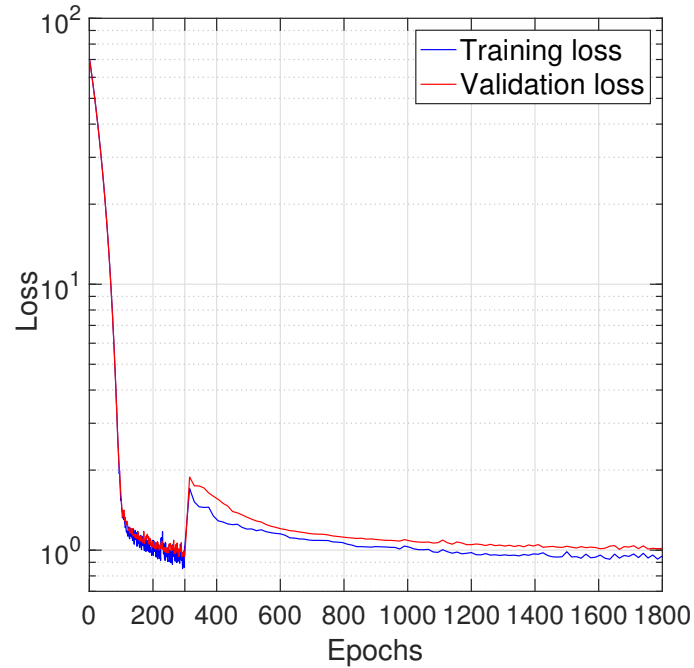


Figure S2: Curves of training and validation losses vs. epochs. No sign of over-fitting was observed. The first 300 epochs are for model initialization based on supervised learning and the remaining 1500 epochs are for model optimization based on joint supervised and unsupervised learning.

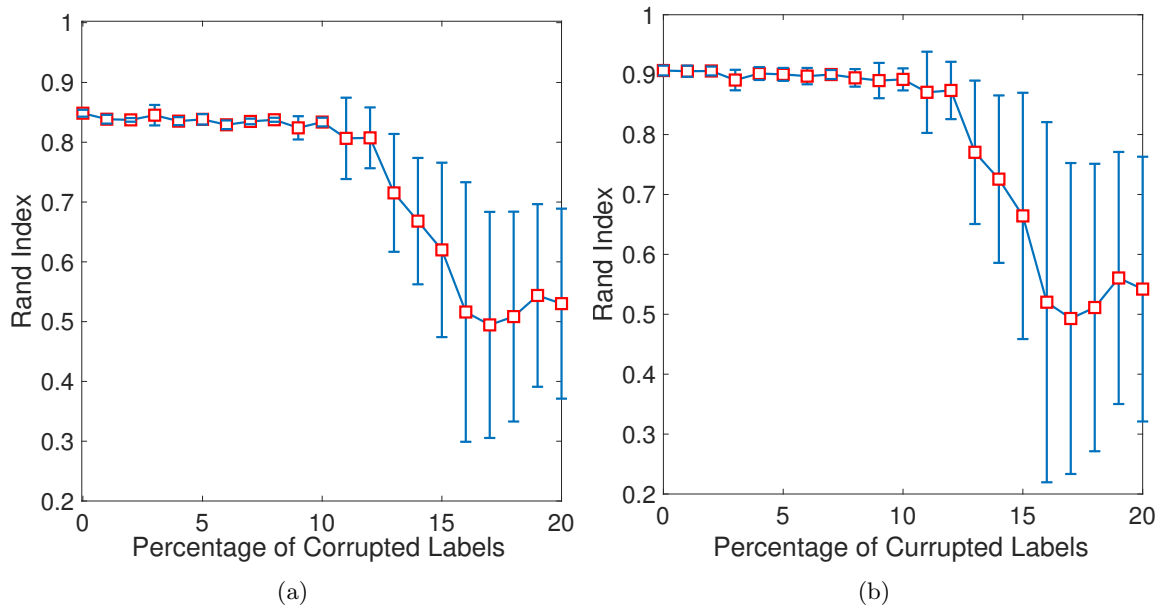


Figure S3: Rand indexes computed by comparing the cluster assignments of the test samples with their PAM50 labels (a) and those obtained by using the original training dataset (b). DeepType performed similarly when up to 10% of the class labels of the training data were corrupted.

S1 Bladder Cancer Study

S1.1 Experiment Setting

The bladder cancer dataset was obtained from the TCGA project, which contains the expression profiles of 20,241 genes from 427 bladder tumor samples. We classified each tumor sample into one of the two UNC subtypes (Damrauer *et al.*, 2014), namely luminal and basal, using the R package BLCAsubtyping (Kamoun *et al.*, 2019), and used the UNC subtypes as the class labels. For model construction and performance evaluation, we randomly partitioned the data into a training and test dataset, containing 80% and 20% of the samples, respectively. Since the sample size is small, in order to avoid possible overfitting, we used only the top 10,000 most variant genes, and designed a three-layer neural network with 10,000, 32 and 2 nodes in the input layer, the hidden layer and the output layer, respectively. The number of clusters K , the tradeoff parameter α and the regularization parameter λ were estimated to be 4, 0.005 and 0.002, respectively. Other experiment settings were similar to those used in the breast cancer study.

S1.2 Clinically Relevant Subtypes Revealed by DeepType

By applying DeepType to the bladder cancer dataset, a total of 156 genes were selected and 4 clusters were detected. The descriptions of the selected genes are given in Table S5. To visualize the detected clusters, we applied t-SNE (van der Maaten and Hinton, 2008) to the outputs of the representation layer. Figure S4(a-b) presents the sample distributions of the identified clusters and their UNC-subtype compositions, respectively. We can see that the tumor samples were grouped into four well-defined clusters, including two luminal dominated clusters (labeled as luminal 1 and 2) and two basal dominated clusters (labeled as basal 1 and 2). To demonstrate the clinical relevance of the identified tumor subtypes, a survival data analysis was performed. Figure S4(c) shows that the four subtypes are associated with distinct prognostic outcomes (logrank test, p -value < 0.0001). Further internal and external validation analysis is presented in Section S1.3.

Figure S4(d) presents the heatmap of the 156 selected genes. We can clearly see two modules, one containing key genes *MSN*, *TNC* and *MUC16*, and the other containing key genes *TOX3* and *PDX1*. The difference in the gene expressions in the two modules divided the samples into two broad categories, i.e., basal and luminal. Specifically, the basal samples have high expressions in the *MSN* module and low expressions in the *TOX3* module, and the luminal samples are on the contrary. Within each UNC subtype, luminal 1 has higher expressions in the *MSN* module than luminal 2, while basal 1 has higher expressions in the *TOX3* module than basal 2. Using t -test and the Benjamini-Hochberg procedure (Benjamini and Hochberg, 1995), we identified 99 genes differentially expressed between luminal 1 and luminal 2, and 112 genes between basal 1 and basal 2 ($FDR \leq 0.05$). Our analysis showed that DeepType is able to identify novel bladder cancer subtypes beyond the UNC subtyping system that are associated with distinct expression patterns.

S1.3 Comparison Study

For comparison, we applied SparseK and iCluster to the bladder cancer dataset. The parameters of the two methods were estimated in the same way as that used in the breast cancer study. We first visualized the sample distributions of the clusters detected by the three methods (Figure S5). As with the breast cancer study, we can see that the clusters identified by DeepType are much more compact than those detected by SparseK and iCluster. Then, we performed a series of external and internal evaluations of the quality of the identified clusters. For external evaluation, we assessed the concordance between cluster assignments and the UNC subtypes, the tumor pathological stages, and the risk of tumor recurrence computed based on a three-gene signature proposed in (Liu *et al.*, 2017) (see Table S6 for a detailed description). For internal evaluation, we used the silhouette width and Davies-Bouldin index to quantify the compactness and separateness of the obtained clusters. The results are reported in Tables S7 and S8. In terms of external criteria, our method performed significantly better than SparseK and slightly better than iCluster. In terms of internal criteria, DeepType resulted in the highest silhouette width and the lowest Davies-Bouldin index, which is consistent with the visualization result presented in Figure S5. Our analysis suggested that our method resulted in subtypes with significantly higher cluster quality than the competing methods.

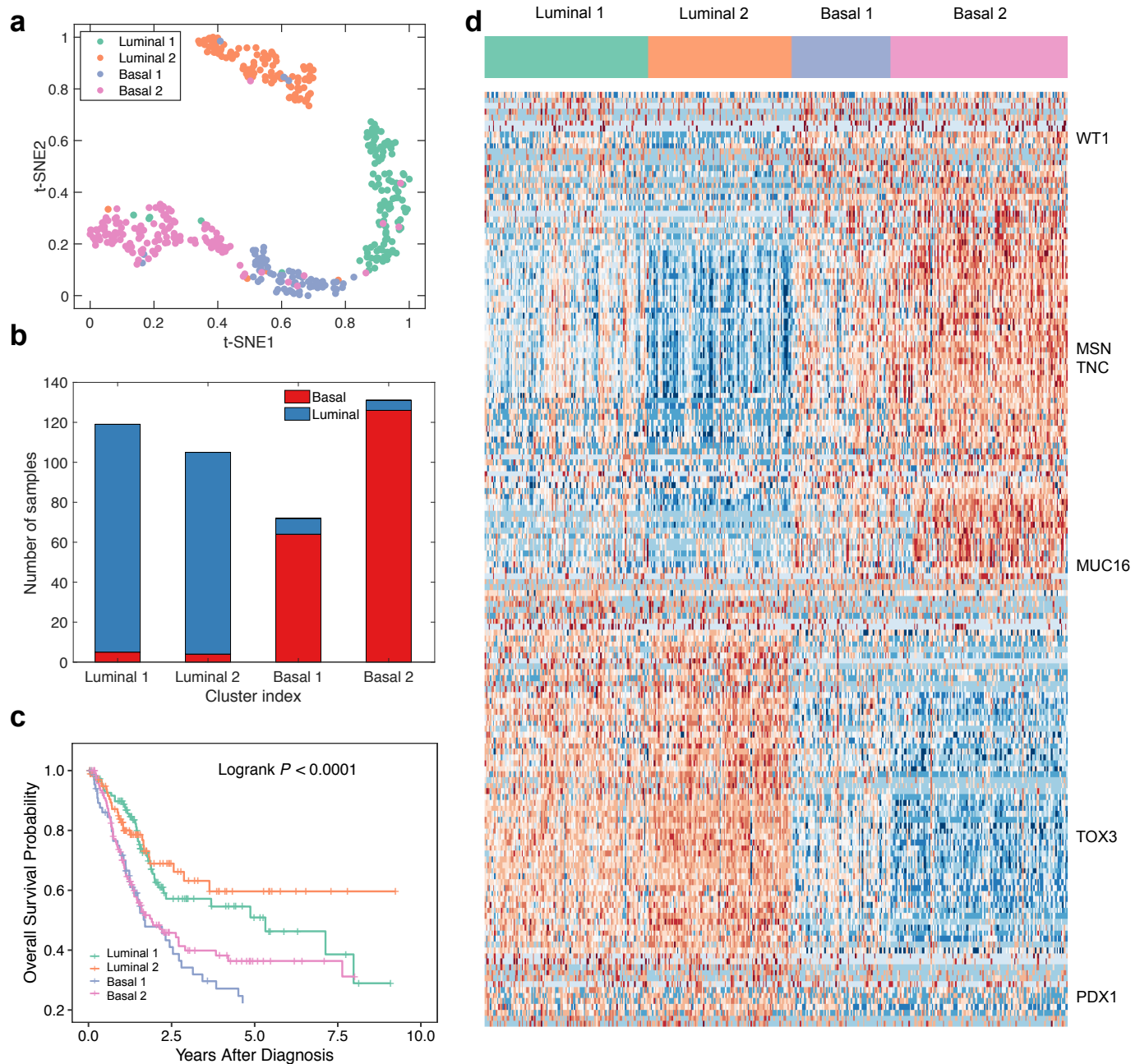


Figure S4: DeepType identified four clinically relevant bladder cancer subtypes. (a) The sample distribution of the identified clusters visualized by t-SNE. (b) The UNC-subtype composition of the four clusters. (c) The four identified subtypes are associated with distinct clinical outcomes. (d) The heatmap of the 156 selected genes showed that the identified clusters exhibited distinct transcriptional characteristics on several gene modules and key genes. The samples were arranged by the clustering assignments, and the expression data are linearly normalized into the scale of [0, 1] across samples.

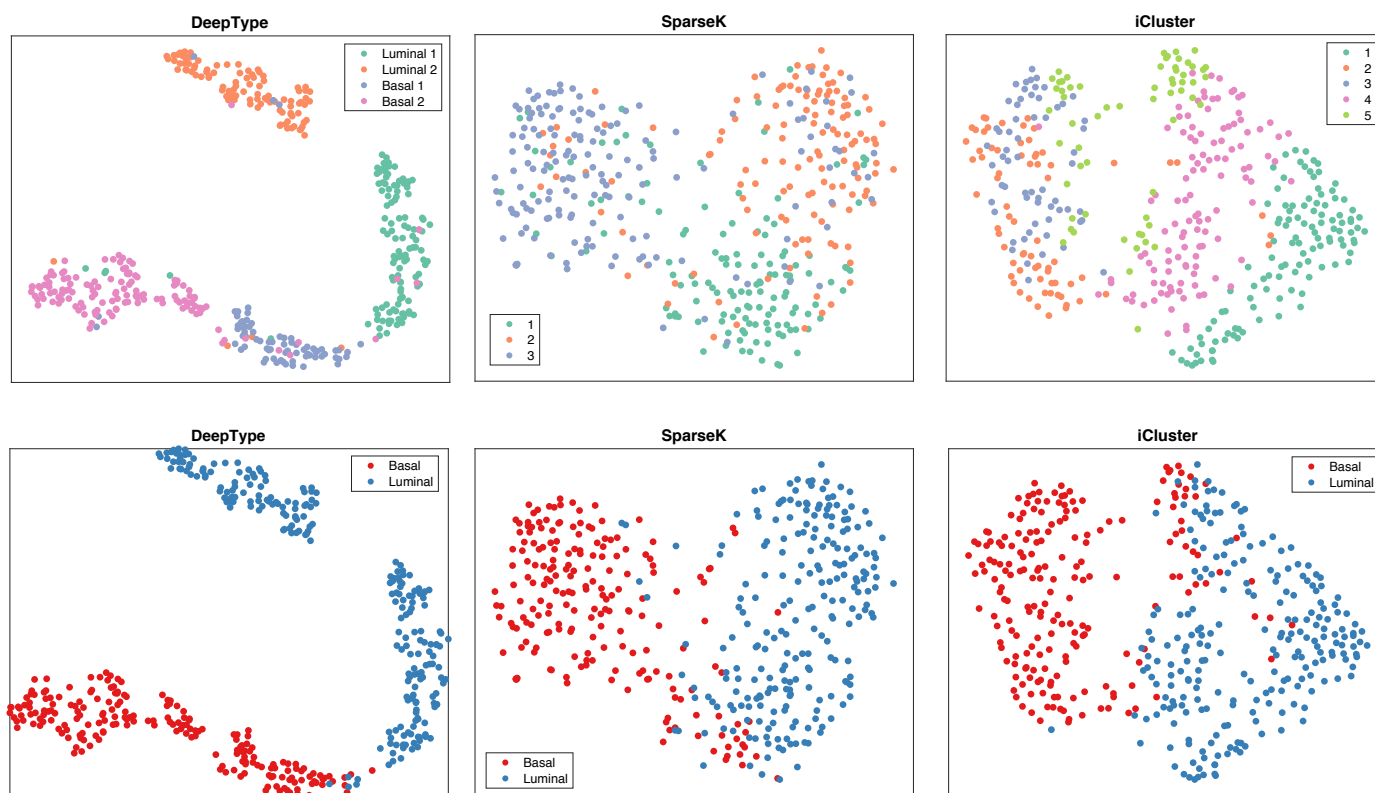


Figure S5: Visualization of the sample distributions of the clusters detected by three methods applied to the bladder cancer dataset. Each sample is color-coded by its clustering assignment (top) and UNC-subtype label (bottom). DeepType revealed a clear four-cluster structure.

Table S5: DeepType identified 156 genes to be informative for bladder cancer subtyping

See the attached Excel file

Table S6: Clinical and prognostic characteristics of bladder cancer

Characteristics	Class label
UNC subtype	basal, luminal
Pathological stage	1, 2, 3
Recurrence risk	low risk, high risk

Table S7: External evaluation of subtypes identified by three methods applied to the bladder cancer dataset. DeepType significantly outperformed SparseK ($p \leq 0.001$) and iCluster ($p \leq 0.03$, Wilcoxon rank-sum test).

	Average Purity			NMI		
	DeepType	SparseK	iCluster	DeepType	SparseK	iCluster
UNC subtypes	0.95	0.75	0.91	0.51	0.16	0.43
Pathological stage	0.45	0.42	0.44	0.04	0.02	0.04
Recurrence risk	0.78	0.66	0.75	0.20	0.08	0.20

Table S8: Internal evaluation of subtypes identified by three methods applied to the bladder cancer dataset. The Davies-Bouldin index is a value in $[0, \infty)$, and a smaller value suggests a better clustering scheme. DeepType outperforms the competing methods by a large margin.

	DeepType	SparseK	iCluster
Silhouette width	0.45	0.19	0.04
Davias-Bouldin index	1.00	2.56	3.29

References

- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, **57**(1), 289–300.
- Damrauer, J. S., Hoadley, K. A., Chism, D. D., Fan, C., Tiganelli, C. J., Wobker, S. E., Yeh, J. J., Milowsky, M. I., Iyer, G., Parker, J. S., *et al.* (2014). Intrinsic subtypes of high-grade bladder cancer reflect the hallmarks of breast cancer biology. *Proceedings of the National Academy of Sciences*, **111**(8), 3110–3115.
- Kamoun, A., de Reynies, A., Allory, Y., Sjödaahl, G., Robertson, A. G., Seiler, R., Hoadley, K. A., Al-Ahmadie, H., Choi, W., Groeneveld, C. S., *et al.* (2019). A consensus molecular classification of muscle-invasive bladder cancer. *bioRxiv*, **488460**.
- Liu, Q., Diao, R., Feng, G., Mu, X., and Li, A. (2017). Risk score based on three mrna expression predicts the survival of bladder cancer. *Oncotarget*, **8**(37), 61583.
- van der Maaten, L. and Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, **9**, 2579–2605.